



**CHARACTERISATION OF THE  
CAPSULAR POLYSACCHARIDE  
BIOSYNTHESIS LOCI OF  
*STREPTOCOCCUS PNEUMONIAE*  
SEROGROUP 19**

by

**Judy Kay Morona B. Sc. (Adelaide)**

A thesis submitted for the degree of Doctor of Philosophy

Department of Microbiology and Immunology  
University of Adelaide

November, 1998

# ABSTRACT

In this thesis, the genetic loci encoding capsular polysaccharide synthesis (*cps*) have been characterised for all members of *Streptococcus pneumoniae* serogroup (19F, 19A, 19B and 19C). In each serotype, the *cps* locus is located in the *S. pneumoniae* chromosome between *dexB* and *aliA* and appears to be arranged as a single transcriptional unit. The arrangement of the genes within the *cps19* loci is highly conserved with 13 genes (*cps19A-H*, and *K-O*) common to all four serogroup 19 members. These genes encode functions required for the synthesis of the shared trisaccharide component of the group 19 capsular polysaccharide (CPS) repeat unit structures. Furthermore, the genetic differences between the group 19 *cps* loci identified are consistent with the differences in the CPS structures of individual serotypes. Functions have been assigned to nearly all the *cps19* gene products, based on either gene complementation or similarity to other proteins with known functions. This has enabled biosynthetic pathways for production of all four group 19 CPSs to be proposed.

Nearly all of the common genes from types 19F, 19B and 19C are >95% identical to each other. However, closely related homologues of *cps19fI* and *J*, which encode the type 19F polysaccharide polymerase and repeat unit transporter, respectively, are not found in the type 19B and 19C *cps* loci. In type 19B and 19C this region of the *cps* locus (between *cps19bH* and *cps19bK*) contains five genes which encode a unique polysaccharide polymerase and repeat unit transporter, as well as two additional putative glycosyl transferases and a protein which may be involved in synthesis of an activated ribose precursor. Transformation studies indicated that these five genes encode all of the functions required to convert a type 19F pneumococcus to type 19B. The type 19C *cps*

locus differs from the 19B *cps* locus only in the insertion of a glucosyl transferase gene (*cps19cS*) between *cps19cK* and *cps19cL*. Transformation studies have shown that the presence of this gene accounts for the additional glucose side chain in the otherwise identical repeat unit structures. The type 19C *cps* locus contains 19 genes, and at 21 kb it is the largest pneumococcal capsule gene cluster characterised to date.

Although the *cps19a* and *cps19f* loci are identical in the number and arrangement of the genes present, the similarity between individual genes varies from 70% to 99% identity (for both the nucleotide and the deduced amino acid sequences). This sequence divergence is surprising given that the only difference between their CPS repeat units is the glycosidic linkage which joins the repeat units together ( $\alpha(1\rightarrow2)$  for 19F and  $\alpha(1\rightarrow3)$  for 19A). Theoretically, only a difference in the *cps19aI* gene, which presumably encodes the polysaccharide polymerase responsible for this linkage, is required to change a type 19F pneumococcus into type 19A. Indeed, this was demonstrated by a transformation event in which the region of the *cps19a* locus encoding Cps19aH and Cps19aI replaced the homologous portion of the *cps19f* locus was sufficient to convert CPS type from 19F to 19A. Given that Cps19fH and Cps19aH are >95% identical, it seems probable that Cps19aI (79% identity) is solely responsible for the observed alteration in CPS type.

The serotype specificity of the *cps19f* genes was investigated by Southern hybridisation analysis of chromosomal DNA from other *S. pneumoniae* serotypes. Large variations in the hybridisation patterns were obtained with the different gene-specific probes. Probes specific for sequences flanking *cps19f* hybridised with all the serotypes tested. However, within the *cps* loci, only *cps19fA* and *cps19fB* were common to all serotypes. Based on the Southern hybridisation analysis a protocol for PCR amplification of *cps* loci was developed and used to amplify the *cps* regions from a variety of pneumococcal serotypes. Direct sequencing of the 5' end of the PCR products was undertaken and identified two classes of *cpsC* gene. Southern hybridisation studies with

*cps19aC*- and *D*-specific gene probes demonstrated that homologues of the first four genes in the *cps* locus, *cpsA-D*, are present in all serotypes and that all the *cps* loci tested evolved from one of two clonal origins which contained either class I or class II *cpsC* and *D* genes. The *cpsE* gene, which encodes a glucose-1-phosphate transferase, is also conserved (in the two distinct classes) in the *cps* loci of all serotypes tested which contain glucose in their CPS, except type 3.

The sequence analysis of the various *cps* loci presented in this thesis provides further evidence that in nature frequent recombination occurs between different *cps* loci resulting in either complete exchange of the *cps* locus or exchange of only part of the *cps* locus, and could potentially result in the expression of new capsular serotypes.

This thesis contains no material which has been accepted for the award of any other degree or diploma in any University or other tertiary institution and to the best of my knowledge and belief, contains no material previously published or written by another person, except where due reference is made in the text.

I give consent to this copy of my thesis, when deposited in the University library being available for loan and photocopying.

Judy Kay Morona

# ACKNOWLEDGMENTS

I would like to thank my supervisors Assoc. Prof. James Paton and Prof. Paul Manning for the opportunity to undertake my postgraduate studies. In particular, thank you James for your encouragement and your confidence in me, and for your help with this thesis. This project has been both challenging and rewarding.

I would also like to thank Gianni Pozzi, Francesco Iannelli, Rubens López, Ernesto García and Alexander Tomasz for kindly providing copies of their manuscripts prior to publication. This has enabled a much more complete comparison of the various pneumococcal *cps* loci. Also thank you to Mike Gratten, Denise Murphy and Andrew Lawrence for serotyping various pneumococcal isolates and transformants.

To my colleagues, past and present, especially Anne Berry, Rob Lock, David Ogunniyi, David Miller, Ursula Talbot, Angelo Guidolin, Adrienne Paton, Rebekah Miller, Chris Vindurampulle, Matthew Woodrow and Stuart McKessar, thank you for the camaraderie. Those elusive results always seemed more attainable after a celebratory glass of champagne. To David M., thank you for your help and enthusiasm over the last few months. To Anne, I especially thank you for all of your help and friendship over the years.

To Renato, my husband, a very special thank you, all of your enthusiastic discussions have been invaluable. Thank you for your support and your help with the housework. I couldn't have done it without you.

To my daughter Stephanie, thank you for your patience and love, I know I have locked myself into the study far too often recently.

# ABBREVIATIONS

Abbreviations which are acceptable to the Journal of Bacteriology are used in this thesis without definition in the text. Additional abbreviations are defined when first used in the text and are listed below.

Amp	ampicillin
AP	alkaline phosphatase
BA	blood agar
Cml	chloramphenicol
CPS	capsular polysaccharide
<i>cps</i>	capsule biosynthesis genes
C-terminus	carboxy terminus
DIG	digoxigenin
DOC	deoxycholate
ECA	enterobacterial common antigen
EPS	exopolysaccharide
Ery	erythromycin
ExoIII	exonuclease III
Gal	galactose
Glc	glucose
GlcA	glucuronic acid
GlcNAc	<i>N</i> -acetyl glucosamine
Hib vaccine	<i>Haemophilus influenzae</i> type b vaccine
HPLC	high-performance liquid chromatography
InPCR	inverse PCR
Kan	kanamycin
LB	Luria-Bertani broth
LB+KA	Luria-Bertani broth containing ampicillin and kanamycin
LOS	lipooligosaccharide
LPS	lipopolysaccharide
LR-PCR	long-range PCR
ManNAc	<i>N</i> -acetyl mannosamine
N-terminus	amino terminus
ORF	open reading frame
PBP	penicillin binding proteins
Rha	rhamnose
Rib	ribose
Rif	rifampicin
SC	free secretory component
Str	streptomycin
TEMED	<i>N,N,N',N'</i> -tetramethyl-ethylene-diamine
Tet	tetracycline
THY	Todd-Hewitt broth supplemented with 0.5% (wt/v) yeast extract
und-P	undecaprenol-phosphate

# TABLE OF CONTENTS

<b>1. INTRODUCTION</b>	<b>1</b>
1.1 <i>Streptococcus pneumoniae</i> , a Distinguished History	1
1.1.1 Development of the Gram stain	1
1.1.2 Humoral immunity and phagocytosis	2
1.1.3 The concept of serotyping	2
1.1.4 The observation of antibiotic resistance in bacteria	2
1.1.5 Recognition of the first non-protein antigen	3
1.1.6 Bacterial transformation	3
1.1.7 Discovery of DNA as the genetic material	4
1.2 Pneumococcal Disease	5
1.2.1 Colonisation of the nasopharynx	5
1.2.2 Risk factors for pneumococcal disease	6
1.2.3 Incidence of pneumococcal disease	7
1.2.4 Pathogenesis of pneumococcal infections	8
1.2.5 Pneumococcal virulence factors	8
1.2.5.1 Capsule	8
1.2.5.2 Cell wall	9
1.2.5.3 Autolysin	9
1.2.5.4 Pneumolysin	9
1.2.5.5 Neuraminidase	10
1.2.5.6 Hyaluronidase	11
1.2.5.7 Pneumococcal surface protein A (PspA)	11
1.2.5.8 PsaA	11
1.2.5.9 Choline binding protein A (CbpA)	12
1.3 The Cell Surface of the Pneumococcus	12
1.3.1 The cell wall	13
1.3.2 Teichoic acid	14
1.3.3 The plasma membrane	15
1.3.4 Proteins associated with the cell surface	16
1.4 The Pneumococcal Capsule	16
1.4.1 Capsule serotyping	17
1.4.2 Chemical structure of the capsular polysaccharide repeat units	18
1.4.3 The chemical structure of group 19 CPS	20
1.4.4 The relationship between serotype and virulence	21

1.5 Distribution of Capsule Types	21
1.5.1 Capsule types associated with paediatric pneumococcal disease	22
1.5.2 Capsule types associated with invasive pneumococcal disease in adults	24
1.5.3 Geographic distribution of invasive pneumococci	24
1.5.4 Changes in serotype distribution over time	25
1.6 Antimicrobial Resistance in the Pneumococcus	26
1.6.1 Geographic distribution of antimicrobial resistance	26
1.6.2 The spread of antibiotic resistance	27
1.6.3 Antibiotic resistance among different pneumococcal serotypes	27
1.7 The Pneumococcal Polysaccharide Vaccine	28
1.7.1 Antibody response	28
1.7.2 Conjugate polysaccharide-protein vaccines	29
1.8 Polysaccharide Biosynthesis	30
1.8.1 Some common features in polysaccharide biosynthesis	31
1.8.2 Polysaccharide biosynthesis via lipid-linked repeat unit intermediates	32
1.8.3 Polysaccharide biosynthesis involving ABC transporters	34
1.8.4 Polysaccharide synthesis via processive glycosyl transferases	34
1.8.5 Polysaccharide synthesis not involving lipid-linked intermediates	35
1.8.5 Capsule gene clusters	36
1.9 The Capsule Locus of <i>S. pneumoniae</i>	37
1.9.1 Pneumococcal capsular transformation	37
1.9.2 Localisation of the type 19F capsule locus	38
1.9.3 Partial cloning and sequencing of the type 19F capsule locus	39
1.9.4 Characterisation of <i>cps19fA-G</i>	40
1.9.4.1 Cps19fA	40
1.9.4.2 Cps19fB	41
1.9.4.3 Cps19fC	41
1.9.4.4 Cps19fD	41
1.9.4.5 Cps19fE	42
1.9.4.6 Cps19fF	43
1.9.4.7 Cps19fG	43
1.10 Aims of this Thesis	43
<b>2. MATERIALS AND METHODS</b>	<b>45</b>
2.1 Bacterial strains and cloning vectors	45
2.1.1 Bacterial strains	45
2.1.2 Cloning vectors	46
2.1.3 Growth media	47

2.2	Chemicals and reagents	48
2.2.1	General chemicals	48
2.2.2	Antibiotics	48
2.2.3	Enzymes	48
2.2.4	Oligodeoxynucleotides	49
2.3	Serotyping of pneumococcal strains	51
2.4	Bacterial transformation	51
2.4.1	Preparation of competent <i>E. coli</i> strains	51
2.4.1.1	RbCl <sub>2</sub> method	51
2.4.1.2	Electroporation method	52
2.4.1.3	CaCl <sub>2</sub> method	52
2.4.2	Transformation of competent <i>E. coli</i> strains	53
2.4.3	Electroporation of <i>E. coli</i> strains	53
2.4.4	Transformation of pneumococcal strains	54
2.4.5	Preparation of competence factor for pneumococcal transformation	54
2.5	DNA extraction procedures	55
2.5.1	Plasmid isolation	55
2.5.2	Plasmid isolation for pneumococcal transformation	56
2.5.3	Preparation of pneumococcal chromosomal DNA	57
2.6	Analysis and manipulation of DNA	58
2.6.1	Restriction endonuclease digestion of DNA	58
2.6.2	Calculation of sizes of restriction fragments	58
2.6.3	Analytical and preparative separation of restriction fragments	59
2.6.4	Isolation of restriction fragments from agarose gels	59
2.6.5	<i>In vitro</i> cloning	59
2.7	Southern hybridisation	60
2.7.1	Labelling of DNA probes	60
2.7.2	DNA hybridisation by Southern blotting	61
2.7.2.1	Transfer of DNA to nylon membrane	61
2.7.2.2	High stringency hybridisation	61
2.7.2.3	Low stringency hybridisation	62
2.7.2.4	Immunological detection	62
2.8	DNA sequencing and analysis	62
2.8.1	Nested deletions	62
2.8.2	DNA sequencing using dye-labelled primers	64
2.8.3	Sequencing with dye-labelled terminators	65
2.8.4	DNA sequencing	65
2.8.5	DNA sequence analysis	65

2.9 PCR Amplification	66
2.9.1 PCR amplification using Taq polymerase	66
2.9.2 Long range PCR	66
2.9.3 Inverse PCR	67
2.9.4 Purification of PCR products	67
2.10 Insertion-duplication mutagenesis	68
2.11 Plasmid insertion/rescue	69
2.12 Protein and LPS analysis	70
2.12.1 T7 expression	70
2.12.2 LPS preparation	70
2.12.3 SDS-PAGE	71
2.12.4 Western blotting	71
<b>3. CHARACTERISATION OF THE COMPLETE TYPE 19F CAPSULE LOCUS</b>	<b>73</b>
3.1 Introduction	73
3.2 Results	73
3.2.1 Isolation, cloning and DNA sequencing of the type 19F <i>cps</i> locus	73
3.2.2 Analysis of the <i>cps19f</i> DNA sequence	76
3.2.3 Characterisation of <i>cps19fG-O</i>	77
3.2.3.1 <i>cps19fG</i>	77
3.2.3.2 <i>cps19fH</i>	79
3.2.3.3 <i>cps19fI</i>	80
3.2.3.4 <i>cps19fJ</i>	83
3.2.3.5 <i>cps19fK</i>	85
3.2.3.6 <i>cps19fLMNO</i>	87
3.2.4 Insertion-duplication mutagenesis of <i>cps19fG-O</i> genes	96
3.2.5 T7 expression of several Cps19f proteins	98
3.2.6 Complementation of an <i>E. coli rffE</i> mutant with Cps19fK	100
3.2.7 Complementation of <i>S. flexneri rfbBDAC</i> with <i>cps19fLMNO</i>	100
3.2.8 Putative biosynthetic pathway for <i>S. pneumoniae</i> type 19F CPS	102
3.3 Conclusions	104
<b>4. ANALYSIS OF CAPSULE LOCI FROM VARIOUS PNEUMOCOCCAL SEROTYPES</b>	<b>106</b>
4.1 Introduction	106

4.2 Results	107
4.2.1 Serotype specificity of the <i>cps19f</i> genes and flanking DNA sequences	107
4.2.1.1 Presence of <i>cps19f</i> homologues in members of serogroup 19	107
4.2.1.2 Presence of <i>cps19f</i> homologues in other serotypes	109
4.2.1.3 Analysis of the intergenic regions	110
4.2.2 Amplification of capsule loci by LR-PCR	111
4.2.3 Southern hybridisation analysis of LR-PCR products	113
4.2.4 Partial DNA sequencing of the LR-PCR products	116
4.2.5 Phylogenetic analysis of <i>cpsC</i> sequences	120
4.3 Conclusions	120
<b>5. CHARACTERISATION OF THE <i>S. PNEUMONIAE</i> TYPE 19B CAPSULE LOCUS</b>	<b>124</b>
5.1 Introduction	124
5.2 Results	125
5.2.1 Isolation of the type 19B specific <i>cps</i> genes	125
5.2.2 Characterisation of the <i>cps19b</i> genes	128
5.2.2.1 <i>cps19bP</i>	129
5.2.2.2 <i>cps19bI</i>	129
5.2.2.3 <i>cps19bQ</i>	129
5.2.2.4 <i>cps19bR</i>	132
5.2.2.5 <i>cps19bJ</i>	133
5.2.3 Serotype specificity of the <i>cps19b</i> genes	134
5.2.4 Capsule type switching by transformation	135
5.3 Conclusions	136
<b>6. MOLECULAR BASIS FOR PNEUMOCOCCUS TYPE 19A</b>	<b>139</b>
6.1 Introduction	139
6.2 Results	140
6.2.1 PCR amplification and sequencing of the type 19A locus	140
6.2.2 Analysis of the <i>cps19a</i> locus	141
6.2.3 Comparison of the <i>cps19a</i> locus from <i>S. pneumoniae</i> strains 19A1 and 19A2	148
6.2.4 Serotype specificity of the <i>cps19a</i> genes	149
6.2.5 Capsule transformation from 19F to 19A	153
6.2.6 Sequence variation in the 5' intergenic region of serogroup 19	155
6.3 Conclusions	159

<b>7. CHARACTERISATION OF THE <i>S. PNEUMONIAE</i> TYPE 19C SPECIFIC CPS REGION</b>	<b>163</b>
7.1 Introduction	163
7.2 Results	164
7.2.1 Isolation of the type 19C specific <i>cps</i> region	164
7.2.2 Characterisation of Cps19cS	168
7.2.3 Serotype specificity of <i>cps19cS</i>	171
7.2.4 Transformation of <i>S. pneumoniae</i> type 19B to type 19C	171
7.2.5 Characterisation of IS19C, located in the 5' intergenic region	172
7.3 Conclusions	174
<b>8. DISCUSSION</b>	<b>176</b>
8.1 Introduction	176
8.2 Analysis of the <i>S. pneumoniae</i> group 19 <i>cps</i> loci	176
8.2.1 Comparison of the group 19 <i>cps</i> loci	176
8.2.2 Distribution of the <i>cps19</i> genes in other pneumococcal serotypes	181
8.2.3 G+C content	181
8.2.4 Transcription and translation of the <i>cps19</i> loci	182
8.3 The <i>cps</i> loci from other <i>S. pneumoniae</i> serotypes	184
8.3.1 The <i>S. pneumoniae</i> type 1 <i>cps</i> locus	184
8.3.2 The <i>S. pneumoniae</i> type 2 <i>cps</i> locus	185
8.3.3 The <i>S. pneumoniae</i> type 3 <i>cps</i> locus	187
8.3.4 The <i>S. pneumoniae</i> type 4 <i>cps</i> locus	187
8.3.5 The <i>S. pneumoniae</i> type 14 <i>cps</i> locus	188
8.3.6 The <i>S. pneumoniae</i> type 23F <i>cps</i> locus	189
8.3.7 The <i>S. pneumoniae</i> type 33F <i>cps</i> locus	190
8.4 Analysis of the <i>cps</i> loci from <i>S. pneumoniae</i>	191
8.4.1 Organisation of pneumococcal <i>cps</i> loci	191
8.4.2 Comparison of the <i>cps</i> genes and their protein products	194
8.4.3 Capsular transformation <i>in vivo</i>	198
8.4.4 The presence of IS elements in the 5' and 3' intergenic regions	200
8.4.5 Transcription of the pneumococcal <i>cps</i> locus	200
8.4.6 Regulation of CPS production in <i>S. pneumoniae</i>	202
8.4.7 The function of CpsC and CpsD	205
8.4.8 Relationship of the <i>S. pneumoniae</i> <i>cps</i> locus to other Gram-positive <i>cps</i> loci	206
8.5 Future studies	209

<b>REFERENCES</b>	<b>210</b>
<b>APPENDIX I</b>	
The nucleotide and deduced amino acid sequence of the distal portion of the <i>cps19f</i> locus.	228
<b>APPENDIX II</b>	
The nucleotide and deduced amino acid sequence of the serotype specific region of the <i>cps19b</i> locus.	232
<b>APPENDIX III</b>	
The nucleotide and deduced amino acid sequence of the complete <i>cps19a</i> locus.	235
<b>APPENDIX IV</b>	
The nucleotide and amino acid sequence of the <i>cps19a</i> locus from <i>cps19aJ-aliA</i> from <i>S. pneumoniae</i> strain 19A2.	240
<b>APPENDIX V</b>	
The nucleotide sequence of Rx1-19F, Rx1-19A.1, Rx1-19A.2, Rx1-19A.3 and Rx1-19A in the regions where recombination between the <i>cps19f</i> and <i>cps19a</i> loci has occurred.	243
<b>APPENDIX VI</b>	
The nucleotide and deduced amino acid sequence of the central (type-specific) portion of <i>cps19c</i> locus.	250
<b>APPENDIX VII</b>	
Publications	252

# Chapter 1



## INTRODUCTION

### 1.1 *Streptococcus pneumoniae*, a Distinguished History

*Streptococcus pneumoniae*, or the pneumococcus, was first isolated in 1880 by Sternberg in the United States and Pasteur in France. Both independently isolated “roughly lancet-shaped pairs of coccoid bacteria” from the blood of rabbits injected with human saliva. Early intensive scientific investigations of the pneumococcus yielded many important discoveries concerning both cellular and molecular biology (Austrian, 1981a; Watson *et al.*, 1993).

Some of these fundamental achievements, which can be attributed to studies on the pneumococcus, are summarised below.

#### 1.1.1 Development of the Gram stain

Gram (1884) examined sections of lung tissue from people who had died of pneumonia and discovered that these sections contained many coccoid bacteria that retained an aniline-gentian violet stain, “the cocci of croupous pneumonia”. Although the significance was not realised at the time, he also noted the presence of an encapsulated coccoid bacterium that did not retain the stain (later identified as *Klebsiella pneumoniae*).

Today, most clinically important bacteria are recognised as either Gram-negative or Gram-positive.

### **1.1.2 Humoral immunity and phagocytosis**

Klempner and Klempner (1891) showed that rabbits injected with heat-killed pneumococci were immune to reinfection with the same strain but not to reinfection with different clinical isolates. They also demonstrated that rabbits were protected against primary infection by infusion of serum from an immunised rabbit. Issaëff (1893) demonstrated that the protective serum was not bactericidal, rather it promoted the ingestion of pneumococci by the phagocytic cells of the immune system. Neufeld and Rimpau (1904) demonstrated that this phenomenon results from the exposure of the bacteria, rather than the phagocytes, to the serum. These observations described opsonisation, the interaction of antibodies and complement with the bacterial surface.

### **1.1.3 The concept of serotyping**

Neufeld (1902) demonstrated both macroscopic agglutination and microscopic swelling of the pneumococcal capsule after addition of specific antiserum to a cell suspension. Although the capsule swelling or “quellung” reaction became the accepted technique for routinely serotyping pneumococcal isolates, it was not used for this purpose until much later (Armstrong, 1931).

### **1.1.4 The observation of antibiotic resistance in bacteria**

The pneumococcus was one of the first bacteria to be recognised as being able to acquire resistance to antibiotics. Morgenroth and Kaufman (1912) isolated pneumococci resistant to optichin, a derivative of quinine, from mice which had been treated with this

antibiotic. Moore and Chesney (1917) also described optichin resistant pneumococci isolated from humans. However, these findings pre-dated the development and introduction of safe, effective antibiotic treatments and were largely forgotten. There were reports of antibiotic resistant pneumococcal isolates in the 1940s (Frisch *et al.*, 1943; Eriksen, 1945), however, the significance of clinical resistance to antibiotics was not recognised until twenty years later when a penicillin resistant pneumococcus was isolated in Australia from the sputum of a patient (Hansman and Bullen, 1967). Antibiotic resistant pneumococci are becoming more prevalent and the significance of this will be discussed in section 1.6.

### **1.1.5 Recognition of the first non-protein antigen**

The discovery of a “soluble substance of the Pneumococcus” in the blood and urine from pneumonia patients (Dochez and Avery, 1917) lead to the eventual recognition of the first non-protein antigen, the pneumococcal capsular polysaccharide (CPS). Heidelberger and Avery (1923) identified this soluble substance as a complex carbohydrate or polysaccharide which was responsible for serological reactivity. However, it was widely believed that only proteins could be immunogenic leading to the assumption that the immunogenicity was due to a contaminating protein, even though it was resistant to the effects of boiling or trypsin digestion. Heidelberger (1927) later concluded that this highly reactive substance was distributed over the surface of the bacteria and was responsible for type specificity. Heidelberger and colleagues then showed that this capsule was immunogenic by immunising mice with the capsular material to protect them from subsequent pneumococcal challenge.

### **1.1.6 Bacterial transformation**

The phenomenon of transformation of pneumococci from one serotype to another

was first described by Griffith (1928). His studies were designed to determine the requirement for reversion of unencapsulated or rough pneumococcal variants to the encapsulated or smooth form. When mice were injected with live rough pneumococci derived from one serotype and heat-killed cells from another, a proportion of the mice succumbed to an infection caused by pneumococci of the same serotype as that of the heat-killed inoculate. It wasn't long before capsular transformation was demonstrated *in vitro* by Dawson and Sia (1931), and repeated by Alloway (1932; 1933) using pneumococcal cell extracts.

### **1.1.7 Discovery of DNA as the genetic material**

The hunt for this “transforming principle” led to one of the most significant discoveries in biological science this century. In 1944, Avery, McLeod and McCarty published a paper which showed conclusively that DNA was the determinant of capsular transformation, and therefore, was the carrier of genetic information (Avery *et al.*, 1944). A paper published two years later (McCarty and Avery, 1946) silenced the sceptics by showing that deoxyribonuclease destroyed the biological activity of the “transforming principle”. In this paper, McCarty wrote with great insight: “It remains one of the challenging problems for future research to determine what sort of configurational or structural differences can be demonstrated between desoxyribonucleates of separate specificities.” McCarty suggested that “the nucleic acid of the pneumococcus is concerned with innumerable other functions of the bacterial cell” thus “discovering the chemical basis of biological specificity of desoxyribonucleic acids becomes extremely complex, since a given preparation will represent a mixture of a large number of entities of diverse specificity.” Hotchkiss (1951) demonstrated this “diverse specificity” of DNA by introducing a trait other than capsule type into a pneumococcus. He induced penicillin

resistance in a penicillin-sensitive pneumococcus by transfer of DNA from a penicillin-resistant strain.

## 1.2 Pneumococcal Disease

More than a century after its initial isolation *S. pneumoniae* is still an important cause of life-threatening, invasive diseases such as pneumonia, bacteraemia and meningitis, with high morbidity and mortality throughout the world. *S. pneumoniae* is also a leading cause of less serious but highly prevalent infections such as otitis media and sinusitis.

An overview of the mechanisms involved in pneumococcal disease and the pathogenesis of *S. pneumoniae* is presented below. Several comprehensive reviews on this topic have been published recently (Musher, 1992; Paton *et al.*, 1993; Tuomanen *et al.*, 1995; AlonsoDeVelasco *et al.*, 1995; Watson and Musher, 1996).

### 1.2.1 Colonisation of the nasopharynx

*S. pneumoniae* is carried, asymptotically, in the upper respiratory tract by many healthy individuals. The transmission of pneumococci from a carrier to another person is dependent on the frequency and intimacy of their contact (Riley and Douglas, 1981). Thus among adults, carrier rates are highest for those living in crowded conditions such as in barracks and dormitories. This person to person spread often occurs concurrently with viral infection of the upper respiratory tract. Carriage rates are highest among preschool children and tend to decrease with increasing age (AlonsoDeVelasco *et al.*, 1995). Day-care centre attendance is associated with increased carriage rates of pneumococci and thus the risk of invasive pneumococcal disease for such children is increased (Takala *et al.*,

1995; Cherian *et al.*, 1994). Virtually everyone is colonised by pneumococci at some stage. The average duration of carriage is about six weeks in an adult, but can be greater than one year in some people (Musher, 1992). There are two consequences of carriage, either seroconversion and subsequent elimination of the pneumococcus, or invasion of the organism leading to pneumococcal disease. Most infections do not result from prolonged carriage, but probably occur within the first week of colonisation (Musher, 1992).

### **1.2.2 Risk factors for pneumococcal disease**

Pneumococcal disease occurs most frequently in the extreme ages of life. High rates of invasive pneumococcal disease occur in children under two. The rate then decreases with increasing age with teenage children and young adults having the lowest rates of disease. The rate then increases with middle age, again reaching a high level in the elderly (Musher, 1992).

People with functional or anatomic asplenia are highly susceptible to pneumococcal bacteraemia because of reduced capacity to clear encapsulated bacteria from the blood. Increased risk of pneumococcal disease is also associated with certain medical conditions that result in either impaired pulmonary clearance mechanisms (such as cigarette smoking, emphysema, chronic bronchitis, chronic pulmonary disease and viral respiratory infections such as influenza), or reduced immune responses (such as leukaemia, multiple myeloma, lymphoma, HIV infection, Hodgkin's disease, organ or bone marrow transplantation, prolonged use of systemic corticosteroids and chronic renal failure). Other medical conditions that increase susceptibility to pneumococcal disease include chronic cardiovascular diseases, chronic liver diseases, diabetes mellitus (which is often associated with either cardiovascular or renal dysfunction), malnutrition, and alcoholism. People with chronic cerebrospinal fluid leakage resulting from congenital lesions, skull fractures or

neurosurgical procedures are also at risk of recurrent pneumococcal meningitis (Centers for Disease Control and Prevention, 1997).

### 1.2.3 Incidence of pneumococcal disease

*S. pneumoniae* is the most common cause of acute otitis media in children aged less than five years in the United States (Bluestone *et al.*, 1992). In the United States more than half of all children suffer acute otitis media during their first year of life, and nearly half have had at least three attacks before their third birthday (Teele *et al.*, 1989). Although acute otitis media does not usually progress to invasive disease, it has a considerable impact upon health care costs.

*S. pneumoniae* is the commonest cause of community-acquired bacterial pneumonia. It is thought that at least 500,000 cases of pneumococcal pneumonia occur annually in the United States. The precise figure is difficult to ascertain because in many cases of pneumonia, the aetiological agent remains unidentified. In the United States, *S. pneumoniae* has been estimated to be responsible for 25-35% of all cases of bacterial pneumonia requiring hospitalisation (Fang, 1990), and concurrent bacteraemia occurs in 10-25% of adult patients. *S. pneumoniae* is also a leading cause of bacterial meningitis in the United States (Wenger *et al.*, 1990).

Invasive pneumococcal disease causes 40,000 deaths a year in the United States alone (Centers for Disease Control and Prevention, 1997). Despite appropriate antimicrobial therapy and intensive medical care, the overall case fatality rates are about 15-20% for bacteraemia and meningitis. This rate increases to 30-40% among elderly patients (Hook *et al.*, 1983; Wenger *et al.*, 1990). As the prevalence of multiply antibiotic resistant pneumococci increases, pneumococcal disease will become more difficult to manage, potentially increasing the mortality rate further.

### 1.2.4 Pathogenesis of pneumococcal infections

The first stage of pneumococcal infection is nasopharyngeal colonisation. The virulence of the colonising strain and the immune status of the host are important determining factors in the onset of pneumococcal disease. However, the mechanisms which enable the pneumococcus to infect the lung or to migrate directly into the bloodstream are poorly understood (AlonsoDeVelasco *et al.*, 1995). Certain risk factors, such as the presence of a respiratory viral infection, greatly enhance the spread of pneumococci to the lungs. Epithelial damage caused by the virus also aids the spread of the pneumococcus into the bloodstream. From the blood, the pneumococcus may spread to the meninges. The pneumococcus may also enter the meninges directly from the nasopharynx either if the dura mater has been compromised or as a complication of sinusitis (Musher, 1992).

### 1.2.5. Pneumococcal virulence factors

Many pneumococcal components are known to contribute to the pathogenicity of this organism. Various virulence factors and their role in pneumococcal disease are described below.

#### 1.2.5.1 Capsule

The polysaccharide capsule protects the pneumococcus from phagocytosis and is recognised as the major virulence factor of *S. pneumoniae*. All clinical isolates of *S. pneumoniae* are smooth. Although rough strains can be maintained *in vitro*, such strains are almost completely avirulent. Avery and Dubos (1931) demonstrated that enzymatic depolymerisation of the capsule increased the LD<sub>50</sub> of a type 3 pneumococcus by 10<sup>6</sup> fold. More recently a similar effect on the virulence of a type 3 pneumococcus was achieved by

transposon mutagenesis of a gene essential for capsule production (Watson and Musher, 1990).

#### **1.2.5.2 Cell Wall**

Both cell wall peptidoglycan and teichoic acid induce inflammatory responses in the host. Teichoic acid activates the alternate pathway of complement which enhances vascular permeability, and recruits and activates leucocytes (AlonsoDeVelasco *et al.*, 1995). The cell wall components also induce production of cytokines such as interleukin-1, and stimulate production of platelet-activating factor (Tuomanen *et al.*, 1995). The acute inflammatory response generated is thought to cause the tissue damage responsible for the high morbidity and mortality associated with pneumococcal disease (Musher, 1992). In fact, the tissue damage caused by pneumococcal disease can be mimicked by injecting purified pneumococcal cell wall components into animals (Tuomanen *et al.*, 1985).

#### **1.2.5.3 Autolysin**

Autolysin, a N-acetylmuramic acid-L-alanine amidase, cleaves the covalent bond between the glycan chains and the peptide side chains in the cell wall. Autolysin mutants tend to grow in short chains rather than discreet diplococci suggesting that autolysin may be involved in cell separation (Paton *et al.*, 1993). The lytic tendency of pneumococci in stationary phase culture, or in the presence of deoxycholate (DOC), is caused by the activation of autolysin. Autolysin is thought to contribute to pneumococcal pathogenesis because of its ability to lyse the cell. Unrestrained growth of the pneumococcus at the site of infection leads to autolysis and the release of cell wall fragments and virulence proteins such as pneumolysin resulting in an inflammatory response and progression of pneumococcal disease (AlonsoDeVelasco *et al.*, 1995).

#### **1.2.5.4 Pneumolysin**

Pneumolysin is a thiol-activated cytolysin which is produced by nearly all clinical

*S. pneumoniae* isolates. Pneumolysin is a member of the thiol-activated cytolysin family of toxins produced by several Gram-positive genera. However, pneumolysin differs from the other members of this family in that it is not secreted, but remains cytoplasmic until it is released by autolysis. Pneumolysin is a bifunctional toxin, and in addition to its cytotoxic properties it is capable of directly activating the classical complement pathway by binding to the Fc region of human IgG (Paton *et al.*, 1984; Mitchell *et al.*, 1991). It has been shown to have toxic effects on many cell types, which undoubtedly contributes to the pathogenesis of pneumococcal disease (for reviews see Paton *et al.*, 1993; Paton, 1996). These properties include inhibition of the bactericidal activity of leucocytes (Paton and Ferrante, 1983), blockade of proliferative responses and Ig production by lymphocytes (Ferrante *et al.*, 1984), reduction of ciliary beating of human respiratory epithelium (Steinfort *et al.*, 1989), and direct cytotoxicity for respiratory endothelial and epithelial cells (Rubins *et al.*, 1992; 1993). Inactivation of the pneumolysin gene in a *S. pneumoniae* type 2 strain increased the LD<sub>50</sub> (in mice) approximately 100-fold (Berry *et al.*, 1989) and immunisation with purified pneumolysin increases the survival time of mice challenged intranasally with virulent pneumococci (Paton *et al.*, 1983; 1993).

#### 1.2.5.5 Neuraminidase

All clinical *S. pneumoniae* isolates produce one or more neuraminidases. To date, two different enzymes have been identified in the pneumococcus (Lock *et al.*, 1988a; Camara *et al.*, 1991; Berry *et al.*, 1996). These enzymes cleave the terminal sialic acid residues from glycolipids, glycoproteins, and oligosaccharides on cell surfaces. This has the potential to cause great damage to the host and, indeed, the purified protein has been shown to be toxic for mice (Lock *et al.*, 1988b). Another possible role in pathogenesis for neuraminidase could be the unmasking of cell-surface receptors for pneumococcal adhesins (Paton *et al.*, 1993).

### 1.2.5.6 Hyaluronidase

Most clinical isolates produce this enzyme which is associated with the cell surface. It degrades hyaluronic acid, which is found associated with mammalian connective tissue. Thus, hyaluronidase may contribute to pathogenesis by degradation of the connective tissue, allowing greater access to host tissues, thereby facilitating invasion (Paton *et al.*, 1993; Berry *et al.*, 1994).

### 1.2.5.7 Pneumococcal surface protein A (PspA)

PspA and autolysin, are members of a family of pneumococcal choline binding proteins which contain C-terminal repeated choline-binding domains (García *et al.*, 1986; Yother and Briles, 1992). Although its precise function is unknown, PspA is a highly variable protective antigen, both in its molecular weight and its antigenicity (Crain *et al.*, 1990; Waltman *et al.*, 1990). PspA-negative mutant pneumococci are more readily cleared from the blood of mice than wild type strains (McDaniel *et al.*, 1987). The relative importance of PspA to pathogenicity of pneumococci appears to vary from strain to strain. More dramatic reductions in virulence were observed when PspA genes in strains of serotypes 3 and 5 were inactivated compared to a type 2 strain (Briles *et al.*, 1988).

### 1.2.5.8 PsaA

PsaA is a lipoprotein which was initially thought to be an adhesin based on sequence homology with putative lipoprotein adhesins of *S. sanguis* and *S. parasanguis* (Sampson *et al.*, 1994). Pneumococcal *psaA*-negative mutants, unable to express PsaA, are virtually avirulent for mice (Berry and Paton, 1996), and immunisation with the purified protein protects mice from challenge with virulent pneumococci (Talkington *et al.*, 1996). The *psaA* gene is part of an operon which was recently demonstrated to encode a manganese transporter (Dintilhac *et al.*, 1997). Thus, its role in pathogenesis could be explained either by a requirement for manganese for regulation of expression of other

virulence factors, or by growth retardation due to an inability to scavenge this metal *in vivo* (Paton, 1998).

#### 1.2.5.9 Choline binding protein A (CbpA)

CbpA is also a pneumococcal choline binding protein (Rosenow *et al.*, 1997). *S. pneumoniae* *cbpA*-negative mutants exhibited reduced adherence to cytokine-activated cells and endothelial cells *in vitro*, and failed to bind to immobilised sialic acid or lacto-N-neotetraose which are known pneumococcal ligands on these cells, suggesting that CbpA may be an adhesin (Rosenow *et al.*, 1997). Hammerschmidt *et al* (1997) independently described a protein, designated SpsA, which binds specifically to human secretory IgA and the free secretory component (SC). Interaction between the pneumococcal cell surface and the SC may interfere with the activity of secretory IgA, and may also directly facilitate adherence, thereby promoting colonisation of the nasopharyngeal mucosa. The demonstrated degree of heterogeneity, in the amino (N-) terminal region, of SpsA among different pneumococcal strains (Hammerschmidt *et al.*, 1997) and the high degree of sequence similarity between CbpA and SpsA indicate that they are the same protein (Paton, 1998). The adhesive properties of CbpA are consistent with the observation that carriage rates of pneumococcal CbpA mutants were reduced 100-fold in an animal model, whereas, there was no reduction of virulence in an intraperitoneal model of sepsis (Rosenow *et al.*, 1997).

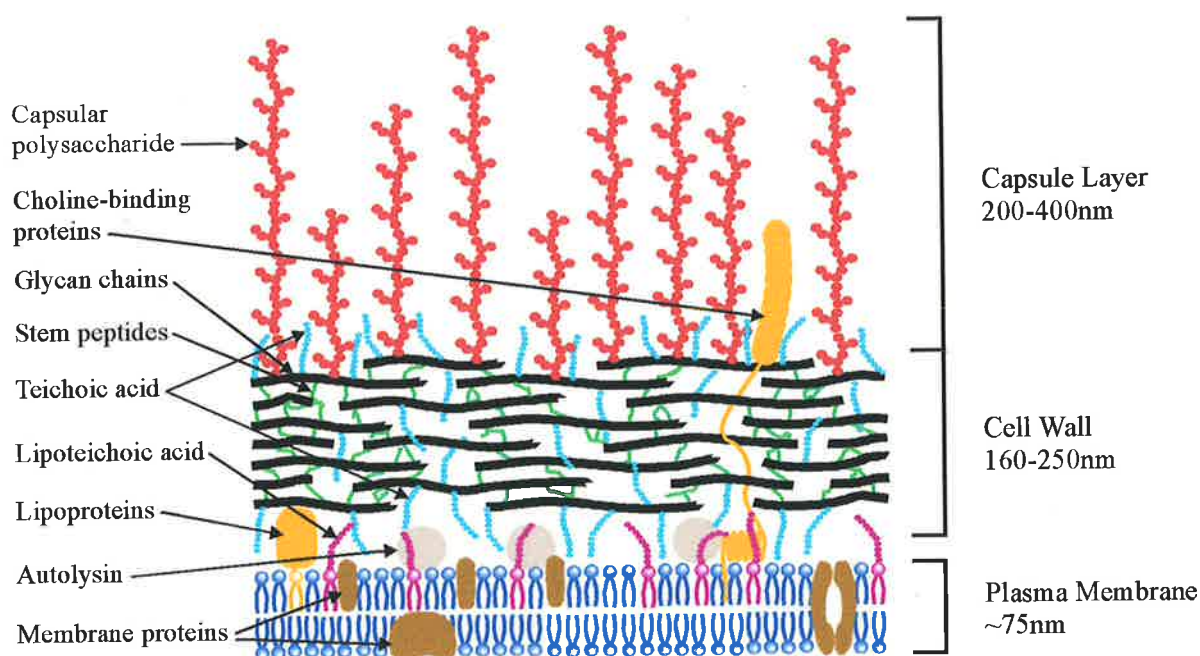
### 1.3 The Cell Surface of the Pneumococcus

A schematic representation of the structure of the pneumococcal cell surface is shown in **Fig. 1.1**. It consists of three layers, the innermost layer is the plasma membrane

which is followed by the peptidoglycan cell wall layer. The outermost layer consists of the polysaccharide capsule which is covalently linked to the peptidoglycan in the cell wall. Properties of each component of the cell surface are described below.

### 1.3.1 The cell wall

The pneumococcal cell wall is typical of a Gram-positive bacterium. It consists of several layers of peptidoglycan, comprising glycan chains of alternating *N*-acetylglucosamine (GlcNAc) and *N*-acetylmuramic acid cross-linked with species-specific peptide bridges. These stem peptides have alanine as the first residue, linked to *N*-acetylmuramic acid. High-performance liquid chromatography (HPLC) has revealed the complexity of the peptidoglycan supramolecular structure for both Gram-positive and Gram-negative bacteria. HPLC analysis has identified 18 different muropeptides of highly conserved molar ratios in *S. pneumoniae* (Garcia-Bustos *et al.*, 1987). This profile is

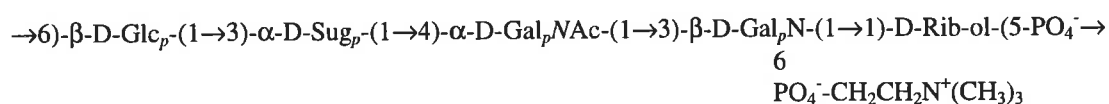


**Fig. 1.1. The pneumococcal cell surface.** Diagrammatic representation of the pneumococcal cell surface showing possible arrangements of constituents. The thickness of each layer is as described by Tomasz (1981). Representation of cell wall, teichoic acid and capsular polysaccharide adapted from Hammond *et al.* (1984), Fig 2.11.

found in all pneumococci, regardless of serotype, geographic origin or isolation date, with the exception of penicillin-resistant strains which have different mucopeptide profiles (Garcia-Bustos and Tomasz, 1990). Resistant pneumococci have alterations in the penicillin binding proteins (PBPs) which result in reduced affinity for the drug. PBPs are enzymes (transpeptidases or carboxypeptidases) that are involved in the assembly of bacterial cell walls. The interaction between PBPs and penicillin is due to the structural similarity between the  $\beta$ -lactam ring of penicillin and the carboxy (C) -terminal D-alanyl-D-alanine residues in the cell wall stem peptides. Thus, it is highly likely the altered PBPs have altered affinities for their natural substrates resulting in an altered peptidoglycan structure in penicillin resistant pneumococci (Severin and Tomasz, 1996).

### 1.3.2 Teichoic acid

The teichoic acid of the pneumococcus, the C-polysaccharide, was first described by Tillett and Francis (1930), although its chemical structure (Fig. 1.2) was not determined until fifty years later (Jennings *et al.*, 1980).



**Fig. 1.2. Structure of the C-polysaccharide.** Glc<sub>p</sub>, glucose; Sug<sub>p</sub>, 2-acetamido-4-amino-2,4,6, trideoxy-D-galactose; Gal<sub>p</sub>NAC, *N*-acetyl galactosamine; Gal<sub>p</sub>N, galactosamine; Rib-ol, ribitol; PO<sub>4</sub><sup>-</sup>-CH<sub>2</sub>CH<sub>2</sub>N<sup>+</sup>(CH<sub>3</sub>)<sub>3</sub>, phosphorylcholine.

The teichoic acid is covalently linked to the pneumococcal peptidoglycan via *N*-acetylmuramic acid (Tomasz, 1981). The phosphorylcholine in the teichoic acid is a recognition site for autolysin, an *N*-acetylmuramic acid-L-alanine amidase, which cleaves

the bond between alanine and *N*-acetylmuramic acid.

Teichoic acid is uniformly distributed on both the inner and the outer surfaces of the cell wall, and possibly within it as well. The thickness of the teichoic acid layer varies between different pneumococcal strains (Sørensen *et al.*, 1988). The level of teichoic acid substitution depends on the degree of cross-linking within the cell wall peptidoglycan (Fischer and Tomasz, 1985).

Interestingly, the individual components of the teichoic acid can all be found as a constituent in at least one pneumococcal CPS type. Thus, it is a possibility that some of the genes involved in teichoic acid and CPS synthesis may be shared. It is also possible that the capsule (*cps*) genes evolved by duplication of the teichoic acid locus in an ancestral *S. pneumoniae* strain.

### 1.3.3 The plasma membrane

The plasma membrane consists of a lipid bilayer closely associated with the inner layer of the cell wall. Complex foldings of the plasma membrane, or mesosomal structures, have been observed in the pneumococcus (Tomasz, 1981). However the function of these mesosomal structures is unknown.

The Forssman (F) antigen is a lipid-linked teichoic acid, identical in structure to C-polysaccharide, and is found in the outer leaflet of the plasma membrane. Approximately 20% of the teichoic acid content of the pneumococcus is lipid-linked. The pneumococcal F antigen is a highly specific inhibitor of autolysin, suggesting that it may play an important role in regulating cellular autolysis. When the association between autolysin and the F antigen is disrupted, either by cessation of cell wall biosynthesis or treatment with detergents such as DOC, autolysin becomes free to interact with the choline in teichoic acid. This interaction enables autolysin to cleave the covalent bond between the glycan

chains and the stem peptides, resulting in destruction of the cell wall and cell lysis (Giudicelli and Tomasz, 1984).

### 1.3.4 Proteins associated with the cell surface

*S. pneumoniae* is predicted to contain a large number of surface exposed proteins but to date only a small number have been investigated. PspA, PspB and CbpA have been shown to be exposed on the cell surface of pneumococci. PspA and CbpA are choline-binding proteins which are structurally similar, with almost identical C-termini (Yother and Briles, 1992; Rosenow *et al.*, 1997; Hammerschmidt *et al.*, 1997). They are anchored to the cell surface by interaction with the F antigen, which is associated with the plasma membrane, via the repeated choline binding motifs in the C-termini. This was demonstrated when a *S. pneumoniae* mutant lacking phosphocholine in its F antigen no longer retained PspA on the cell surface but released it into the culture medium (Yother *et al.*, 1998). Several other membrane associated proteins have also been identified in the pneumococcus. These include the PBPs which are integral membrane proteins involved with peptidoglycan synthesis, and autolysin, a choline-binding protein, which also interacts with the choline moieties in the F antigen. Lipoproteins, such as PsaA (Sampson, *et al.*, 1994), are also found associated with the pneumococcal plasma membrane.

## 1.4 The Pneumococcal Capsule

The outermost layer of the pneumococcus consists of a polysaccharide capsule, which is covalently linked to the peptidoglycan in the cell wall (Sørensen *et al.*, 1990).

The nature of this covalent linkage is unknown. The polysaccharide capsule is the thickest layer, completely masking other structures during exponential growth.

#### 1.4.1. Capsule serotyping

Serological analysis of *S. pneumoniae* isolates over the past four decades has resulted in the recognition of 90 structurally distinct capsule types (Henrichsen, 1995). Originally two typing systems were employed; the American system, which simply assigned serotype numbers in the order identified, and the Danish system which groups antigenically similar serotypes. The Danish system is now most widely used and will be used throughout this thesis.

Since the 1930s, serotyping of pneumococcal isolates has been routinely performed using the quellung (or capsular swelling) reaction with type-specific antisera. When the specific antibody binds to the capsule of a pneumococcus, it results in a change in its refractive index; the cells appear swollen and are clearly visible under bright field or phase-contrast microscopy.

A list of type designations and their antigenic formulas are shown in **Table 1.1**. Antigenic formulas are determined by the pattern of antigenic reactivity of factor sera. Factor sera are developed by cross-absorption of type-specific antisera (Henrichsen, 1995). Antigen 'a' is a factor characteristic of the types of a group or an individual type. The letters b, c, d, etc. indicate additional partial antigens which are characteristic of only some types within a group. For example the two factor sera 6b and 6c are required to distinguish types 6A and 6B and are made by reciprocal absorption of type-specific antisera. The factor sera often cross-react with types from a different group, for example, factor serum 7h reacts with serotypes 7B, 7C, 19B, 19C, 24F, 24B and 40. The antigenic cross-reactivity between serotypes is illustrated in **Table 1.1**.

**Table 1.1. Antigenic Formulas for the 90 serotypes of pneumococci**

Group	Type	Antigenic Formula	Group	Type	Antigenic Formula	Group	Type	Antigenic Formula
1		1a	15	15F	15a, 15b, 15c, 15f	28	28A	28a, 28c, 23d
2		2a		15A	15a, 15c, 15d, 15g	29		29a, 29b, 13b
3		3a		15B	15a, 15b, 15d, 15e, 15h	31		31a, 20b
4		4a		15C	15a, 15d, 15e	32	32F	32a, 27b
5		5a	16	16F	16a, 16b, 11d		32A	32a, 32b, 27b
6	6A	6a, 6b		16A	16a, 16c	33	33F	33a, 33b, 33d
	6B	6a, 6c	17	17F	17a, 17b		33A	33a, 33b, 33d, 20b
7	7F	7a, 7b		17A	17a, 17c		33B	33a, 33c, 33d, 33f
	7A	7a, 7b, 7c	18	18F	18a, 18b, 18c, 18f		33C	33a, 33c, 33e
	7B	7a, 7d, 7e, 7h		18A	18a, 18b, 18d		33D	33a, 33c, 33d, 33f, 6a
	7C	7a, 7d, 7f, 7g, 7h		18B	18a, 18b, 18e, 18g	34		34a, 34b
8		8a		18C	18a, 18b, 18c, 18e	35	35F	35a, 35b, 34b
9	9A	9a, 9c, 9d	19	19F	19a, 19b, 19d		35A	35a, 35c, 20b
	9L	9a, 9b, 9c, 9f		19A	19a, 19c, 19d		35B	35a, 35c, 29b
	9N	9a, 9b, 9e		19B	19a, 19c, 19e, 7h		35C	35a, 35c, 20b, 42a
	9V	9a, 9c, 9d, 9g		19C	19a, 19c, 19f, 7h	36		36a, 9e
10	10F	10a, 10b	20		20a, 20b, 7g	37		37a
	10A	10a, 10c, 10d	21		21a	38		38a, 25b
	10B	10a, 10b, 10c, 10d, 10e	22	22F	22a, 22b	39		39a, 10d
	10C	10a, 10b, 10c, 10f		22A	22a, 22c	40		40a, 7g, 7h
11	11F	11a, 11b, 11e, 11g	23	23F	23a, 23b, 18b	41	41F	41a, 41b
	11A	11a, 11c, 11d, 11e		23A	23a, 23c, 15a		41A	41a
	11B	11a, 11b, 11f, 11g		23B	23a, 23b, 23d	42		42a, 20b, 35c
	11C	11a, 11b, 11c, 11d, 11f	24	24F	24a, 24b, 24d, 7h	43		43a, 43b
	11D	11a, 11b, 11c, 11e		24A	24a, 24c, 24d	44		44a, 44b, 12b, 12d
12	12F	12a, 12b, 12d		24B	24a, 24b, 24e, 7h	45		45a
	12A	12a, 12c, 12d	25	25F	25a, 25b	46		46a, 12c, 44b
	12B	12a, 12b, 12c, 12e		25A	25a, 25c, 38a	47	47F	47a, 35a, 35b
13		13a, 13b	27		27a, 27b	47	47A	47a, 43b
14		14a	28	28F	28a, 28b, 16b, 23d	48		48a

Adapted from Henrichsen (1995), Table 6.

## 1.4.2 Chemical structure of the capsular polysaccharide repeat units

The structure of the type 3 repeat unit was determined in 1941 (Reeves and Goebel, 1941), followed by type 8 in 1957 (Jones and Perry, 1957). Since then the repeat unit structures of more than half of the capsule types have been determined, mostly in the 1980s (reviewed by van Dam *et al.*, 1990). The chemical structures of the capsular repeat units vary greatly, ranging from linear polymers of one or two saccharides to complex polysaccharides with additional side chains.

Most pneumococcal capsules are negatively charged and possess acidic components such as: D-glucuronic acid (GlcA) (including types 1, 2, 3, 5, 8, 9A, 9N and 9V), phosphodiester bonds (including types 6A, 6B, 11A, 15F, 19F, 19A, and 23F), and pyruvate (type 4). To illustrate the diversity, the CPS repeat unit structures for some of the clinically most significant serotypes are shown in Table 1.2 (Lee, 1987).

**Table 1.2 Chemical structures of capsular polysaccharide repeat units**

Type	Chemical Structure	Reference
1	$\rightarrow 3\text{-}\alpha\text{-D-Sug}_p\text{-}(1\rightarrow 4)\text{-}\alpha\text{-D-Gal}_p\text{A}\text{-}(1\rightarrow 3)\text{-}\alpha\text{-D-Gal}_p\text{A}\text{-}(1\rightarrow$ (+0.3 OAc)	Lindberg <i>et al.</i> , 1980
2	$\rightarrow 4\text{-}\beta\text{-D-Glc}_p\text{-}(1\rightarrow 3)\text{-}\alpha\text{-L-Rha}_p\text{-}(1\rightarrow 3)\text{-}\alpha\text{-L-Rha}_p\text{-}(1\rightarrow 3)\text{-}\beta\text{-L-Rha}_p\text{-}(1\rightarrow$ 2 ↑ 1 $\alpha\text{-D-Glc}_p\text{A}\text{-}(1\rightarrow 6)\text{-}\beta\text{-D-Glc}_p$	Jansson <i>et al.</i> , 1988
3	$\rightarrow 3\text{-}\beta\text{-D-Glc}_p\text{A}\text{-}(1\rightarrow 4)\text{-}\beta\text{-D-Glc}_p\text{-}(1\rightarrow$	Reeves and Goebel, 1941
4	$\rightarrow 3\text{-}\beta\text{-D-Man}_p\text{NAc}\text{-}(1\rightarrow 3)\text{-}\alpha\text{-L-Fuc}_p\text{NAc}\text{-}(1\rightarrow 3)\text{-}\alpha\text{-D-Gal}_p\text{NAc}\text{-}(1\rightarrow 4)\text{-}\alpha\text{-D-Gal}_p\text{-}(1\rightarrow$ 2 3 X H <sub>3</sub> C COOH	Jones and Currie, 1988
5	$\rightarrow 4\text{-}\beta\text{-D-Glc}_p\text{-}(1\rightarrow 4)\text{-}\alpha\text{-L-Fuc}_p\text{NAc}\text{-}(1\rightarrow 3)\text{-}\beta\text{-D-Hex}_p\text{-}(1\rightarrow$ 3 ↑ 1 $\alpha\text{-L-Pne}_p\text{NAc}\text{-}(1\rightarrow 3)\text{-}\beta\text{-D-Glc}_p\text{A}$	Jansson <i>et al.</i> , 1985
6A	$\rightarrow 2\text{-}\alpha\text{-D-Gal}_p\text{-}(1\rightarrow 3)\text{-}\alpha\text{-D-Glc}_p\text{-}(1\rightarrow 3)\text{-}\alpha\text{-L-Rha}_p\text{-}(1\rightarrow 3)\text{-D-Rib-ol}\text{-}(5\text{-PO}_4^- \rightarrow$	Rebers and Heidelberg, 1961
6B	$\rightarrow 2\text{-}\alpha\text{-D-Gal}_p\text{-}(1\rightarrow 3)\text{-}\alpha\text{-D-Glc}_p\text{-}(1\rightarrow 3)\text{-}\alpha\text{-L-Rha}_p\text{-}(1\rightarrow 4)\text{-D-Rib-ol}\text{-}(5\text{-PO}_4^- \rightarrow$	Kenne <i>et al.</i> , 1979
7F	$\rightarrow 6\text{-}\alpha\text{-D-Gal}_p\text{-}(1\rightarrow 3)\text{-}\beta\text{-L-Rha}_p\text{-}(1\rightarrow 4)\text{-}\beta\text{-D-Glc}_p\text{-}(1\rightarrow 3)\text{-}\beta\text{-D-Gal}_p\text{NAc}\text{-}(1\rightarrow$ 2 2 4 ↑ OAc ↑ 1 1 $\beta\text{-D-Gal}_p$ $\alpha\text{-D-Glc}_p\text{NAc}\text{-}(1\rightarrow 2)\text{-}\alpha\text{-L-Rha}_p$	Moreau <i>et al.</i> , 1988
8	$\rightarrow 4\text{-}\beta\text{-D-Glc}_p\text{A}\text{-}(1\rightarrow 4)\text{-}\beta\text{-D-Glc}_p\text{-}(1\rightarrow 4)\text{-}\alpha\text{-D-Glc}_p\text{-}(1\rightarrow 4)\text{-}\alpha\text{-D-Gal}_p\text{-}(1\rightarrow$	Jones and Perry, 1957
9N	$\rightarrow 4\text{-}\alpha\text{-D-Glc}_p\text{A}\text{-}(1\rightarrow 3)\text{-}\alpha\text{-D-Glc}_p\text{-}(1\rightarrow 3)\text{-}\beta\text{-D-Man}_p\text{NAc}\text{-}(1\rightarrow 4)\text{-}\beta\text{-D-Glc}_p\text{-}(1\rightarrow 4)\text{-}\alpha\text{-D-Glc}_p\text{NAc}\text{-}(1\rightarrow$	Jones <i>et al.</i> , 1985
9V	$\rightarrow 4\text{-}\alpha\text{-D-Glc}_p\text{A}\text{-}(1\rightarrow 3)\text{-}\alpha\text{-D-Gal}_p\text{-}(1\rightarrow 3)\text{-}\beta\text{-D-Man}_p\text{NAc}\text{-}(1\rightarrow 4)\text{-}\beta\text{-D-Glc}_p\text{-}(1\rightarrow 4)\text{-}\alpha\text{-D-Glc}_p\text{-}(1\rightarrow$ +OAc +OAc	Perry <i>et al.</i> , 1981
14	$\rightarrow 6\text{-}\beta\text{-D-Glc}_p\text{NAc}\text{-}(1\rightarrow 3)\text{-}\beta\text{-D-Gal}_p\text{-}(1\rightarrow 4)\text{-}\beta\text{-D-Glc}_p\text{-}(1\rightarrow$ 4 ↑ 1 $\beta\text{-D-Gal}_p$	Lindberg <i>et al.</i> , 1977
18C	$\rightarrow 4\text{-}\beta\text{-D-Glc}_p\text{-}(1\rightarrow 4)\text{-}\beta\text{-D-Gal}_p\text{-}(1\rightarrow 4)\text{-}\alpha\text{-D-Glc}_p\text{-}(1\rightarrow 3)\text{-}\alpha\text{-L-Rha}_p\text{-}(1\rightarrow$ 2 3 ↑ ↑ 1 PO <sub>4</sub> <sup>-</sup> -1-Glyc-ol AcO3- $\alpha\text{-D-Glc}_p$	Lugowski and Jennings, 1984
23F	$\rightarrow 4\text{-}\beta\text{-D-Glc}_p\text{-}(1\rightarrow 4)\text{-}\beta\text{-D-Gal}_p\text{-}(1\rightarrow 4)\text{-}\beta\text{-L-Rha}_p\text{-}(1\rightarrow$ 2 3 ↑ ↑ 1 PO <sub>4</sub> <sup>-</sup> -2-Glyc-ol $\alpha\text{-L-Rha}_p$	Richards and Perry, 1988

Abbreviations: Sug<sub>p</sub>, 2-acetamido-4-amino-2,4,6, trideoxy-D-galactose; Hex<sub>p</sub>, 2-acetamido-2,6 dideoxy-D-xylo-hexos-4-ulose; Pn<sub>p</sub>NAc, 2-acetamido-2,6 dideoxy-L-talose; Glc<sub>p</sub>, glucose; Gal<sub>p</sub>, galactose; Rha<sub>p</sub>, rhamnose; Glc<sub>p</sub>A, glucuronic acid; Gal<sub>p</sub>A, galacturonic acid; Man<sub>p</sub>NAc, N-acetyl mannosamine; Fuc<sub>p</sub>NAc, N-acetyl fucosamine; Gal<sub>p</sub>NAc, N-acetyl galactosamine; Glc<sub>p</sub>NAc, N-acetyl glucosamine; ribitol; Glyc-ol, glycerol; OAc, O-acetyl; PO<sub>4</sub><sup>-</sup>, phosphate.

### 1.4.3 The chemical structure of group 19 CPS

Of particular relevance to this thesis is the chemical structure of the CPS from group 19 pneumococci. Group 19 consists of four immunologically cross-reactive *S. pneumoniae* types (19F, 19A, 19B and 19C). The CPS repeat unit structures for the members of group 19 are shown in **Table 1.3**. Two different structures have been proposed for type 19A. The expression of these two different repeat units has been reported to be dependent on culture conditions (Lee *et al.*, 1987).

Group 19 pneumococci are a clinically significant cause of pneumococcal disease, although prevalence of the component serotypes varies. In one study (Robbins *et al.*, 1983), group 19 pneumococci accounted for 7% of all isolates from cases of invasive disease. Of these, 65% were caused by 19F, 34% by 19A, 1% by 19B; 19C was a very rare cause of disease in this study.

**Table 1.3 Chemical structures of the capsular polysaccharide repeat units for group 19**

Type	Chemical Structure	Reference
19F	$\rightarrow 4\text{-}\beta\text{-D-Man}_p\text{NAc-(1}\rightarrow 4\text{)-}\alpha\text{-D-Glc}_p\text{-(1}\rightarrow 2\text{)-}\alpha\text{-L-Rha}_p\text{-(1-PO}_4^- \rightarrow$	Ohno <i>et al.</i> , 1980
19A	$\rightarrow 4\text{-}\beta\text{-D-Man}_p\text{NAc-(1}\rightarrow 4\text{)-}\alpha\text{-D-Glc}_p\text{-(1}\rightarrow 3\text{)-}\alpha\text{-L-Rha}_p\text{-(1-PO}_4^- \rightarrow$	Katzenellenbogen and Jennings, 1983
	$\rightarrow 4\text{-}\beta\text{-D-Man}_p\text{NAc-(1}\rightarrow 4\text{)-}\alpha\text{-D-Glc}_p\text{-(1}\rightarrow 2\text{)-}\alpha\text{-L-Rha}_p\text{-(1-PO}_4^- \rightarrow$	Lee and Fraser, 1980
	$\beta\text{-D-Glc}_p\text{NAc-(1}\rightarrow 3\text{)-}\beta\text{-D-Gal}_p\text{-(1-PO}_4^- \quad \text{PO}_4^- \rightarrow 1\text{)-}\alpha\text{-L-Fuc}_p$	
19B	$\rightarrow 4\text{-}\beta\text{-D-Man}_p\text{NAc-(1}\rightarrow 4\text{)-}\beta\text{-D-Glc}_p\text{-(1}\rightarrow 4\text{)-}\beta\text{-D-Man}_p\text{NAc-(1}\rightarrow 4\text{)-}\alpha\text{-L-Rha}_p\text{-(1-PO}_4^- \rightarrow$	Beynon <i>et al.</i> , 1991
	$\begin{array}{c} 3 \\ \uparrow \\ 1 \\ \beta\text{-D-Rib}_p\text{-(1}\rightarrow 4\text{)-}\alpha\text{-L-Rha}_p \end{array}$	
19C	$\beta\text{-D-Glc}_p$	Beynon <i>et al.</i> , 1991
	$\begin{array}{c} 1 \\ \downarrow \\ 6 \\ \rightarrow 4\text{-}\beta\text{-D-Man}_p\text{NAc-(1}\rightarrow 4\text{)-}\beta\text{-D-Glc}_p\text{-(1}\rightarrow 4\text{)-}\beta\text{-D-Man}_p\text{NAc-(1}\rightarrow 4\text{)-}\alpha\text{-L-Rha}_p\text{-(1-PO}_4^- \rightarrow \\ 3 \\ \uparrow \\ 1 \\ \beta\text{-D-Rib}_p\text{-(1}\rightarrow 4\text{)-}\alpha\text{-L-Rha}_p \end{array}$	

Abbreviations: Glc<sub>p</sub>, glucose; Rha<sub>p</sub>, rhamnose; Man<sub>p</sub>NAc, N-acetyl mannosamine; Gal<sub>p</sub>, galactose; Fuc<sub>p</sub>NAc, N-acetyl fucosamine; Glc<sub>p</sub>NAc, N-acetyl glucosamine; Rib<sub>p</sub>, ribose; PO<sub>4</sub><sup>-</sup>, phosphate.

#### 1.4.4 The relationship between serotype and virulence

The importance of the capsule in pneumococcal infection is due to its antiphagocytic properties (van Dam *et al.*, 1990). Immunity to pneumococcal infection is serotype dependent and the presence of specific antibody to the capsule results in opsonisation and rapid clearance of the invading pneumococci.

Not all pneumococcal serotypes are equally invasive, and both the capsule type and the amount of CPS expressed affect the relative virulence of pneumococci (Austrian, 1981a). MacLeod and Krauss (1950) showed that the quantity of CPS produced by each of several mutants influenced their virulence in mice. Recently, the CPS serotype expressed by otherwise isogenic pneumococcal strains, generated by *in vitro* or *in vivo* transformation, was shown to directly affect their virulence for mice (Kelly *et al.*, 1994; Nesin *et al.*, 1998). However, the genetic background of these strains was also important, and virulence appeared to be dependent on a combination of CPS and other genetic factors (Kelly *et al.*, 1994). The difference in virulence between serotypes is probably due in part to differences in the ability of the CPS to prevent activation of the alternative complement pathway, to resist phagocytosis and to induce an antibody response (AlonsoDeVelasco *et al.*, 1995).

### 1.5 Distribution of Capsule Types

In 1979 a worldwide pneumococcal typing surveillance study was initiated by the World Health Organisation. Since then, more than 25,000 strains have been typed in Denmark alone. Even though a large amount of information about serotype distribution is available, analysis of this data is complicated by the fact that serotype prevalence varies

with site of infection, age, time and geographical location (Nielsen and Henrichsen, 1992).

### 1.5.1 Capsule types associated with paediatric pneumococcal disease

Many epidemiological studies investigating the distribution of pneumococcal serotypes causing paediatric invasive disease have been reported in the literature. The results from many of these studies are summarised in **Table 1.4**. Serotypes/groups 6, 14,

**Table 1.4 Incidence of invasive disease in children**

Study	Age	No. of Isolates	Most common serotypes/groups (% of total)														
			1	2	3	4	5	6	7	9	11	12	14	18	19	23	
Europe <sup>a</sup> , 1982-1987	0-14	2,390	7	0.8	1.9	3.4	3.2	15.6	6.9	4.1	0.8	1.7	13.3	8.2	12	6.3	
Spain <sup>b</sup> , 1979-1993	0-5	167	3	0.6		2.4	7.2	17.4	3	4.8		0	2.6	4.8	17.4	15	
Belgium <sup>b</sup> , 1991	0-5	77	5.2	0		2.6	1.3	7.8	6.5	9		0	28.6	9.1	15.6	5.2	
Denmark <sup>b</sup> , 1991-1992	0-14	164	6.7	0		1.8	1.8	20.1	6.7	6.7		0.6	13.4	16.5	10.4	2.4	
Finland <sup>b</sup> , 1985-1989	0-2	235	0	0		6	0	18.4	8.5	6		0.3	21.5	7.6	17.4	7.3	
Boston, Mass. <sup>c</sup> , 1957-1969	0-15	71	<2	0	0	8.5	2.8	11.3	<2	2.8	0	2.8	29.6	12.7	12.7	5.6	
Chicago, Illinois <sup>d</sup> , 1967-1976	0-15	293	4.4	0	3.8	5.8	<2	16	4.1	7.9	0	<2	17.1	8.9	10.6	9.6	
Birmingham, Alabama <sup>e</sup> , 1975-1978	0-15	114	3.5	0	0	7.1	0.8	20.2	2.6	4.4	0	0	15.8	13.1	9.6	12.3	
USA <sup>f</sup> , 1978-1994	0-6	3570	1.1	0	0.8	6.7	0.2	17.2	1.6	7.2		0.6	28	8.2	14	6.9	
Australia <sup>g</sup> , 1970-1979	0-18	219	4.6	0		9.1	0.5	10.5	2.3	6.4		1.4	22.4	11.9	10	7.8	
Australian Aborigines <sup>h</sup> , 1989-1994	0-14	114	5.3	0	0	6.1	0	13.1	4.4	11.4	0	0	8.8	9.6	7	5.3	8.8
Israel <sup>i</sup> , 1988-1990	0-12	213	13.4	0		2.7	16.5	10.7	5.4	10.7		3.1	14.7	6.7	9.4	3.6	
South Africa <sup>b</sup> , 1987-1991	0-5	1138	10.5	1.1		3.5	1.7	28	2.8	1.8		1.1	19.2	3.9	13	3.2	
Mexico <sup>b</sup> , 1992-1993	0-5	120	0.8	2.5		0.8	0.8	17.5	0.8	6.7		0.8	9.2	2.5	14.2	20.8	
Brazil <sup>i</sup> , 1977-1988	0-2	308	6.4	1.9		2.6	10.3	18.1	1.9	5.1		1.6	10.7	10.5	8	5.4	
Uruguay <sup>k</sup> , 1987-1989	0-5	48	6.3	0		2.1	14.6	4.2	4.2	10.5		0	39.6	0	0	2.1	
Gambia <sup>b</sup> , 1991	0-5	59	10.2	0		0	8.5	15.3	0	5.1		8.5	32.2	1.7	5.1	0	
Egypt <sup>lb</sup> , 1977-1978, 1992	0-6	86	31.4	2.3		2.3	4.7	8.1	3.5	9.3		5.8	7	3.5	2.3	1.2	
Pakistan <sup>m</sup> , 1986-1990	0-5	168	0.6	0		0	4.2	10.7	0.6	6.5		0	1.2	2.4	38.7	4.2	
Papua New Guinea <sup>no</sup> , 1980-1987	0-5	151	2.6	7.3		0.7	10.6	6.6	13.2	2		2.6	6.6	0.7	7.3	7.3	
Rwanda <sup>p</sup> , 1984-1990	0-15	130	22.3	0.7		0	14.6	10	0.8	0.8		0.8	14.6	3.1	3.8	0.8	

a. Nielsen and Henrichsen (1992); b. Snaidack *et al.* (1995); c. Burke *et al.* (1971); d. Jacobs *et al.* (1979); e. Gray *et al.* (1979); f. Butler *et al.* (1995); g. Hansman (1983); h. Gratten *et al.* (1996); i. Dagan *et al.* (1992); j. Tauney *et al.* (1990); k. Mogdasy *et al.* (1992); l. Guirguis *et al.* (1990); m. Mastro *et al.* (1991); n. Gratten and Montgomery (1991); o. Gratten *et al.* (1985); p. Bogaerts *et al.* (1993). The six most common serotypes in each study are shaded.

18, 19 and 23 are consistently among the most common disease causing isolates in children, although their rank order differs significantly. Within serogroups 6, 18, 19 and 23, serotypes 6B, 18C, 19F and 23F are the most common isolates (Nielsen and Henrichsen, 1992; Butler *et al.*, 1995). Interestingly, these serotypes are also the most poorly immunogenic in children (Douglas *et al.*, 1983). Other serotypes/groups, such as type 1, 4, 5, 7F, 9 and 12F, are a considerable cause of disease in only some of the studies. For example, type 1 is a common cause of disease in many developing countries and in parts of Europe but not in the USA, whereas type 4 is most commonly isolated in the USA. In the USA, serotypes 6B, 14, 18C, 19F and 23F along with types 3 and 19A are the most common pneumococcal isolates causing otitis media as shown in **Table 1.5** (Butler *et al.*, 1995). Generally the most common disease causing serotypes are also those most frequently carried in the nasopharynx (**Table 1.6**). However, there are some exceptions, 11A, for example, is a rare cause of pneumococcal disease in children but it is frequently carried. In Papua New Guinea, groups 15 and 33 (not shown in **Table 1.6**) are also

**Table 1.5 Incidence of otitis media in children**

Study	Age	No. of Isolates	Most common serotypes/groups (% of total)													
			1	2	3	4	5	6	7	9	11	12	14	18	19	23
USA <sup>a</sup> , 1978-1994	0-6	249			9	4		14		3			15	2	22	13
Birmingham, Alabama <sup>b</sup> 1975-1978	0-15	396	1.8	0.3	10.6	3.5	0.3	9.6	2.5	4.3	1.8	0	9.6	4.5	28.2	11.9

a. Butler *et al.* (1995); b. Gray *et al.* (1979). The six most common serotypes in each study are shaded.

**Table 1.6 Carriage of pneumococcal strains**

Study	Age	No. of Isolates	Most common serotypes/groups (% of total)													
			1	2	3	4	5	6	7	9	11	12	14	18	19	23
Birmingham, Alabama <sup>a</sup> 1975-1978	0-15	245	2	0	14.3	5.7	0	22.9	1.2	3.7	8.2	0.4	8.2	10.2	20.4	16.3
Birmingham, Alabama <sup>b</sup> 1974-1975	0-2	573	0.2	0	2.8	3.3	0	27	1.2	4.5	4.2	0	7	2.6	17.4	15.4
Papua New Guinea <sup>c</sup> 1985-1987	0-7	1449	0.1	0.1	1.6	1.9	0.3	26.6	1.7	3.2	1	0.8	8.1	1.5	23.3	14.1

a. Gray *et al.* (1979); b. Gray, *et al.* (1980); c. Montgomery *et al.* (1990). The six most common serotypes in each study are shaded.

commonly isolated from the nasopharynx with carriage rates of 5.6% and 4.8% respectively (Montgomery *et al.*, 1990), but these serotypes rarely cause disease.

### 1.5.2 Capsule types associated with invasive pneumococcal disease in adults

Although there have been many epidemiological studies investigating the distribution of pneumococcal serotypes causing disease, few report serotype prevalence only in adults. The serotypes responsible for most adult pneumococcal disease are quite different to those that cause paediatric disease. Serotypes 1, 3, 4, 7F, 8 and 14 were associated with most disease in adults in the studies summarised in **Table 1.7**.

**Table 1.7 Incidence of invasive disease in adults**

Study Country, Year	Age	No. of Isolates	Most common serotypes (% of total)														
			1	2	3	4	5	6	7	8	9	11	12	14	18	19	23
Europe and Israel <sup>a</sup> , 1982-1987	>14	6,376	8.9	0.4	9.1	6.6	1.5	6.3	7.8	5.1	6.8	2.2	3.2	7.9	2.2	6.6	4.4
Australian Aborigines <sup>b</sup> , 1989-1994	>14	107	15.9	0	12.1	7.5	0	0.9	13.1	6.5	3.8	0.9	9.3	1.9	0.9	0.9	0.9

a. Nielsen and Henrichsen (1992); b. Gratten *et al.* (1996). Shaded boxes represent the six most common serotypes in each study.

### 1.5.3 Geographic distribution of invasive pneumococci

The serotypes that most commonly cause invasive pneumococcal disease in different countries are summarised in **Table 1.8**. Serotypes/groups 6, 14 and 19, are consistently important in all geographic areas. Others occur more frequently in some countries. Types 1 and 5 are a major cause of disease in developing countries, such as Africa, China and Israel, whereas group 18 is only important in developed countries (Snaideck *et al.*, 1995). The serotype distribution among indigenous and disadvantaged minorities in developed countries reflects that of the country to which they belong.

However, they often live in poor conditions and their disease rates are far greater than the rest of the population (McIntyre, 1997).

**Table 1.8 Incidence of invasive disease (all ages)**

Study Country, Year	No. of Isolates	Most common serotypes (% of total)														
		1	2	3	4	5	6	7	8	9	11	12	14	18	19	23
Japan <sup>a</sup> , 1984	430	1.2	0	12.7	2.4	1.2	11.8	1.4	0.3	5.6	4.1	1.4	4.9	2.5	13	6.8
Taiwan <sup>a</sup> , 1984	170	2.9	0	5.6	3.5	0.6	9.5	1.8	1.8	4.1	2.9	0	10	1.2	10.5	7.1
China <sup>a</sup> , 1984	448	9.8	7.8	5.6	1.1	14.1	11.6	2.5	2	1.6	0.9	2.7	7.6	1.8	9.2	7.4
Europe and Israel <sup>b</sup> , 1982-1987	10,298	8	0.6	6.7	5.3	2.3	9.1	7.6	3.7	5.9	1.8	2.7	8.7	3.8	8.2	5.2
Australian Aborigines <sup>c</sup> , 1989-1994	221	11.3	0	6.4	7.4	0	9.4	5.9	4.4	8.4	0.5	8.9	6.4	4.4	3.4	5.4
South Africa <sup>d</sup> , 1979-1986	4766	24.2	4	5.9	4.4	3.7	18.5	5	6.3	2.7	0.7	2.2	9.3	4.3	10.2	1.6
USA <sup>e</sup> , 1970-1983	4676	5	0.2	7.1	10.6	1.6	6.4	5.9	9.8	5.4	1	4	9.2	4.4	5.5	4.7
CDC <sup>e</sup> , 1978-1983	1900	3.2	0	4.5	9.4	0.8	11.7	3.5	3.5	7.8	0.7	2.2	14.7	7.5	11	7.1

a. Lee (1987); b. Nielsen and Henrichsen (1992); c. Gratten *et al.* (1996); d. Klugman and Koornhof (1988); e. Robbins *et al.* (1983). Shaded boxes represent the six most common serotypes in each study.

### 1.5.4 Changes in serotype distribution over time

Surveys of the serotype distribution in the 1930s showed that most of the pneumococcal disease recorded was caused by a small number of types, principally 1, 2, 3, 5, 7 and 8 (Finland and Barnes, 1977). The apparent small number of serotypes causing disease may have been partly due to the limitations in the number of serotypes that could be positively identified at the time. Types 1, 2 and 3 were the most commonly isolated types in Boston City Hospital before 1950 (Finland and Barnes, 1977), but collectively accounted for less than 5% of all isolates in 1979-1982 (Barry *et al.*, 1984). The occurrence of type 5 was also declining. Recent investigations in the developing countries have shown a type distribution similar to that in the industrialised countries before 1950 (Klugman and Koornhof, 1988; Snaideck *et al.*, 1995; Scott *et al.*, 1996; Sankilampi, 1997).

## 1.6 Antimicrobial Resistance in the Pneumococcus

The emergence of clinically significant resistance to antibiotics was first noted in 1967 when a penicillin resistant pneumococcus was isolated, in Australia, from the sputum of a patient (Hansman and Bullen, 1967). Subsequently, penicillin resistant pneumococci were isolated in New Guinea, and the prevalence of penicillin resistant isolates from this region increased dramatically from 12% in 1970 (Hansman *et al.*, 1974) to 33% in 1980 (Gratten *et al.*, 1980). Multiply resistant pneumococci (resistant to three or more classes of antibiotics) were first isolated in South Africa (Jacobs *et al.*, 1978); these strains showed a range of resistance patterns to several antibiotics including penicillin, tetracycline and chloramphenicol. Today, antibiotic resistant pneumococci are common throughout the world and are rapidly increasing in prevalence (Appelbaum, 1992). This is becoming a major clinical problem as management of multiply resistant pneumococcal infections become more complicated and requires the use of more expensive alternative antimicrobial agents. The impact of increased antimicrobial resistance on the mortality rate has not yet been clearly defined (Centers for Disease Control and Prevention, 1997).

### 1.6.1 Geographical distribution of antimicrobial resistance

Antibiotic resistance in pneumococci is increasing all over the world. The pattern of antibiotic resistance is uneven both between countries and within countries. South Africa, Spain and Eastern Europe report rates of penicillin resistance of up to 50% of all isolates, compared with less than 10% in most of Western Europe and USA (Schreiber and Jacobs, 1995). In Australia, resistance to antibiotics is also increasing. In 1989 less than 2% of all isolates were penicillin resistant, but by 1994 that had increased to 5.8% (Collignon and Bell, 1996). Interestingly, 8% of all isolates were found to be resistant to at

least three antibiotics, including penicillin, erythromycin, tetracycline, and/or chloramphenicol. In some Aboriginal communities, levels of antibiotic resistant pneumococci are much higher. One study, in Darwin, showed penicillin resistance to be as high as 30% (Skull *et al.*, 1996). Recent data from the Women's and Children's Hospital (WCH) in Adelaide shows that 40% of *S. pneumoniae* isolates from WCH patients are resistant to penicillin, compared with only 15% two years ago (J. Bell, WCH, personal communication).

### **1.6.2 The spread of antibiotic resistance**

There is a link between antibiotic usage and the increasing prevalence of antibiotic resistance in pneumococci. One study from Iceland (Aarson *et al.*, 1996) confirms this linkage among children attending day-care. In this study, the overall pneumococcal carriage rate was 52.7%, with 9.7% of isolates being penicillin resistant. Only 1.8% of children who had no courses of antibiotics in the previous year carried penicillin resistant strains compared with 14% who had antibiotics, and 61% of those who had received antibiotics in the previous two to seven weeks.

### **1.6.3 Antibiotic resistance among different pneumococcal serotypes**

Antibiotic resistance is most prevalent in *S. pneumoniae* serotypes/groups 6, 9, 14, 19 and 23, which are commonly associated with carriage in children (Table 1.4, Schreiber and Jacobs, 1995). These serotypes accounted for more than 80% of the penicillin resistant infections among Aboriginal children in Darwin (Skull *et al.*, 1996).

## 1.7 The Pneumococcal Polysaccharide Vaccine

The first commercially available pneumococcal vaccine was licensed in the USA in 1977. This vaccine consisted of 50 µg of purified capsular polysaccharide of 14 of the most common disease related serotypes (Hilleman *et al.*, 1981). These 14 serotypes (1, 2, 3, 4, 6A, 7F, 8, 9N, 12F, 14, 18C, 19F, 23F and 25F) accounted for 68% of all cases of pneumococcal disease in the USA. In 1983, the 14-valent vaccine was replaced by a 23-valent formulation containing 25 µg of the purified capsular polysaccharide from serotypes 1, 2, 3, 4, 5, 6B, 7F, 8, 9N, 9V, 10A, 11A, 12F, 14, 15B, 17F, 18C, 19F, 19A, 20, 22F, 23F and 33F, which are responsible for more than 85% of all cases of pneumococcal disease in North America and Europe (Robbins *et al.*, 1983). The halving of the amount of each capsular polysaccharide did not affect the antibody response to the vaccine.

### 1.7.1 Antibody response

Pneumococcal CPS induces type-specific antibodies that enhance opsonisation, phagocytosis and killing of pneumococci. After vaccination, a type-specific antibody response develops in >80% of healthy young adults (Musher *et al.*, 1990) and this may persist for up to 10 years (Mufson *et al.*, 1987). However, immune responses are usually not consistent among all 23 serotypes in the vaccine and the level of antibody response required for protection against pneumococcal disease has not been clearly defined.

The immune response to the pneumococcal vaccine is dependent on the age, health and immunocompetency of the individual. The antibody response to the vaccine is generally good in healthy young adults, but it is often diminished or absent in young children aged <2 years and immunocompromised individuals (Centers for Disease Control

and Prevention, 1997). Bacterial polysaccharides induce antibodies primarily by T-cell independent mechanisms and are therefore poorly immunogenic in children aged <2 years whose immune systems are immature. The immune response to some common paediatric pneumococcal serotypes is also suboptimal in children aged 2-5 years (Douglas *et al.*, 1983) and does not reach "adult levels" until age 8 to 10 years (Paton *et al.*, 1986).

### 1.7.2 Conjugate polysaccharide-protein vaccines

Conjugation to a protein carrier dramatically increases the immunogenicity of polysaccharide antigens in young children, converting the polysaccharide from a T-cell independent to a T-cell dependent antigen, capable of eliciting immunological memory. The success of this approach has been demonstrated with *Haemophilus influenzae* type b (Hib). Since the licensing of the first conjugate Hib vaccine in 1987, for use in children at least 18 months of age, the incidence of Hib meningitis has significantly decreased. An American study (Adams *et al.*, 1993) shows a sharp decrease of 56% in the incidence of Hib meningitis between 1989 and 1991. Interestingly, there was also a substantial decrease in the incidence of Hib meningitis in infants aged less than one year during this same period even though the vaccine was not licensed for use in this age group until 1990. Hib conjugate vaccination has been shown to lower the rate of nasopharyngeal carriage (Takala *et al.*, 1991), thus, vaccination of older children can reduce the probability of spread of Hib to younger unvaccinated children.

Development of a pneumococcal conjugate vaccine is more complicated than for Hib because of the numerous pneumococcal serotypes that cause disease. The number of pneumococcal serotypes that could be included in a conjugate vaccine formulation is probably limited. The pneumococcal serotypes that cause most paediatric pneumococcal disease are also those that are commonly antibiotic resistant, thus initial vaccination against

these serotypes would provide the greatest benefit to the community. Several pneumococcal conjugate formulations are currently being evaluated.

An Israeli study, involving vaccination of children 12 to 18 months with a conjugate vaccine consisting of polysaccharides of types 4, 6B, 9V, 14, 18C, 19F and 23F conjugated to the outer membrane complex of *Neisseria meningitidis* serogroup B, showed a significant reduction in the nasopharyngeal carriage rates one year after vaccination (Dagan *et al.*, 1996). The carriage rates of pneumococcal types included in the vaccine and antibiotic resistant pneumococci were significantly lower than in the control group which had received the current 23-valent vaccine. However, as expected, the carriage rate of pneumococcal serotypes not included in the conjugate vaccine was unaffected.

Another study, investigating the immunogenicity in Gambian infants of a pentavalent vaccine consisting of types 6B, 14, 18C, 19F and 23F conjugated to the diphtheria toxin mutant protein CRM<sub>197</sub>, also showed a significant reduction in the nasopharyngeal carriage of pneumococci belonging to the vaccine serotypes two years after vaccination (Leach *et al.*, 1996; Obaro *et al.*, 1996). However, in this study, the carriage rate of non-vaccine serotypes was significantly higher than in the unvaccinated control group. Thus, widespread use of conjugate pneumococcal vaccines may alter the pneumococcal serotype distribution associated with colonisation, and potentially also the types associated with invasive disease.

## 1.8 Polysaccharide Biosynthesis

Bacteria produce a large variety of polysaccharides associated with the outer surface of the cell. These include capsular polysaccharide (CPS), lipopolysaccharide O-antigen

(LPS), lipooligosaccharide (LOS), exopolysaccharide (EPS), enterobacterial common antigen (ECA), and cell wall polysaccharides such as peptidoglycan and teichoic acid. These polysaccharides are diverse in their structure and are often species and type-specific antigens.

### 1.8.1 Some common features in polysaccharide biosynthesis

Most polysaccharides are synthesised from activated sugar intermediates. Monosaccharides can be activated by the addition of nucleoside diphosphates by nucleotidyltransferases. Specific activated nucleotide sugars can also be synthesised by enzymatic modification of the monosaccharide part of a nucleotide sugar to produce derivatives such as deoxy sugars, aminodeoxy sugars and glycuronic acids (Shibaev, 1986). However, the biosynthesis of some polysaccharides does not involve nucleotide sugars, such as glucan biosynthesis in *Streptococcus mutans* (as discussed in section 1.8.5).

Biosynthesis of many polysaccharides involves a reversible first step, that is transfer of sugar 1-phosphate to the lipid carrier undecaprenol phosphate (und-P). Only a limited number of sugars are known to be transferred to und-P. RfbP from *Salmonella enterica* serovar typhimurium (Jiang *et al.*, 1991) has been shown to transfer galactose (Gal)-1-phosphate to und-P, and several functional homologues have been identified in other polysaccharide systems. Some of the identified homologues have been shown to transfer a different sugar, glucose (Glc)-1-phosphate, to a lipid carrier, possibly und-P (Kolkman *et al.*, 1996). Another enzyme known to initiate polysaccharide synthesis in Gram-negative bacteria is Rfe which transfers GlcNAc-1-phosphate to und-P (Meier-Dieter *et al.*, 1990).

Polymerisation of the polysaccharide usually occurs via one of two different mechanisms. These two main mechanisms differ in the direction of chain elongation,

which occurs at either the reducing or the non-reducing end of the polysaccharide chain (Whitfield, 1995). They are described in greater detail below. However, at this stage, the precise mechanisms involved in subsequent transport across the outer membrane (in Gram-negative bacteria) and covalent attachment (if required) of the polysaccharide to the outer surface are poorly understood.

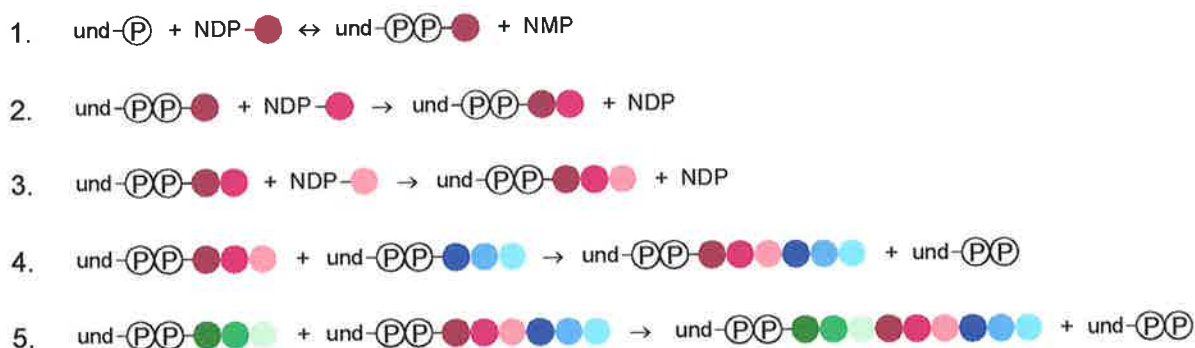
### **1.8.2 Polysaccharide biosynthesis via lipid-linked repeat unit intermediates**

Biosynthesis of many polysaccharides, including: LPS O-antigen in *S. enterica*, *Escherichia coli* and *Shigella flexneri* (Whitfield, 1995), peptidoglycan (Fuchs-Cleveland and Gilvarg, 1976), the EPS xanthan gum in *Xanthomonas campestris* (Ielpi *et al.*, 1993), EPS from *Lactococcus lactis* (van Kranenburg *et al.*, 1997), succinoglycan from *Rhizobium meliloti* (Reuber and Walker, 1993) and CPS from *Aerobacter aerogenes* (Troy *et al.*, 1971), involves assembly of an oligosaccharide repeat unit intermediate on und-P (**Fig. 1.3A**). Biosynthesis of the repeat units occurs through a series of interdependent and sequential reactions, each mediated by a specific glycosyl transferase.

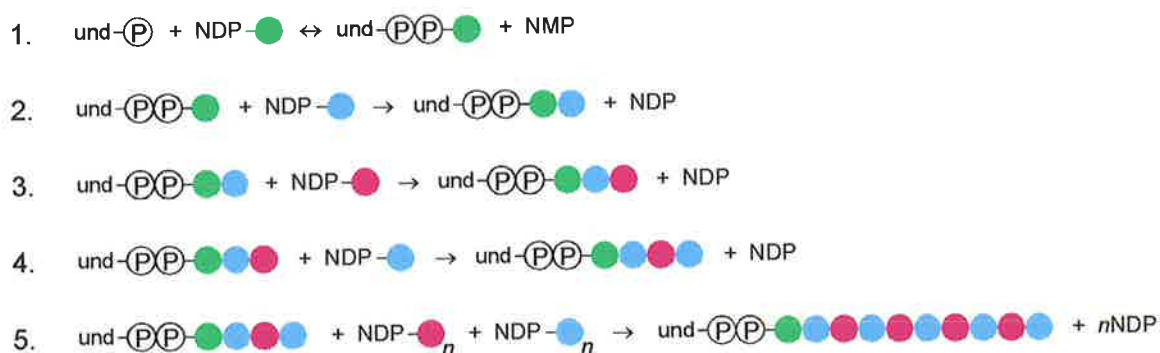
The initial reaction involves the transfer of sugar 1-phosphate from the nucleotide-diphosphate precursor to und-P (step 1). The subsequent transferase reactions involve the formation of glycosidic linkages to the und-P-P-sugar intermediate by transfer of the sugar moiety only, releasing the nucleotide-phosphate (steps 2 & 3). This method allows the synthesis of complex repeat unit structures with complete fidelity, due to the substrate specificity of each of the glycosyl transferases. The und-P-P-oligosaccharide repeat units are synthesised on the cytoplasmic side of the plasma membrane and are translocated to the outer surface by the repeat unit transporter, an integral membrane protein. The polysaccharide polymerase then links the repeat units to form a polysaccharide chain. The polysaccharide polymerase transfers the growing polysaccharide chain from one und-P

carrier to the non-reducing end of a newly synthesised und-P-P-oligosaccharide repeat unit (steps 4 & 5). Growth of the polymer therefore occurs at the reducing terminus (the end closest to und-P) in a block-wise fashion (Whitfield, 1995).

### A.



### B.



**Figure 1.3. Polymerisation of polysaccharides.** Two distinct polymerisation models have been identified which can be distinguished by the direction of chain elongation. **A.** Polysaccharide biosynthesis via lipid-linked repeat unit intermediates. Polymerisation of a polysaccharide with a trisaccharide repeat unit in a block-wise fashion where the growth of the polymer occurs at the reducing terminus or the end closest to und-P. Coloured circles represent sugars and groups of three similarly coloured circles represent the trisaccharide repeat units. **B.** Polysaccharide biosynthesis involving ABC transporters. Polymerisation of a polysaccharide with a disaccharide repeat unit is illustrated. Polymerisation is initiated by the addition of a sugar-1-phosphate, that is not necessarily present in the repeat unit structure of the polysaccharide, and continues by sequential addition of sugars at the non-reducing terminus or the end furthest from und-P. Different coloured circles represent different sugars. (Adapted from Whitfield, 1995)

### 1.8.3 Polysaccharide biosynthesis involving ABC transporters

Polymerisation at the non-reducing end (Fig. 1.3B) involves the sequential transfer of sugars to the non-reducing end of the und-P polymer. Polysaccharides synthesised by this method tend to have simple repeat unit structures of one or two sugars. Some polysaccharides synthesised in this way include: CPS from *E. coli* K1, K4, K5, K92 and *Haemophilus influenzae* type b (Roberts, 1996), and *N. meningitidis* group B (Frosch *et al.*, 1989), O-polysaccharide from *K. pneumoniae* serotype O1 (Bronner *et al.*, 1994), teichoic acid from *Bacillus subtilis* (Lazarevic and Karamata, 1995), amylovoran from *Erwinia amylovora* (Bugert and Geider, 1995), and Nod factors from *R. meliloti* (Vázquez *et al.*, 1993). The pathway is initiated by the transfer of a sugar-1-phosphate residue to und-P (step 1). This sugar is not necessarily present in the repeat unit structure of the polysaccharide. The sugars are then transferred directly from NDP-sugar precursors to the growing polysaccharide chain by specific glycosyl transferases (steps 2 to 5). Growth of the polymer occurs at the end furthest from the und-P carrier (Whitfield, 1995).

Transport of the completed polysaccharide chain across the cytoplasmic membrane into the periplasm is facilitated by ABC transporters. ABC transporters are ubiquitous and are involved in a diverse range of import and export systems. Those involved in polysaccharide export form a closely related sub-class within this large group of proteins (Reizer *et al.*, 1992). They are a two component system, comprising an integral membrane protein and a hydrophilic ATP-binding protein. The transport complex consists of two of each component (Reizer *et al.*, 1992).

### 1.8.4 Polysaccharide synthesis via processive glycosyl transferases

LPS O:54 from *S. enterica* serovar borreze (Keenleyside and Whitfield, 1996), cellulose from *Acetobacter xylinum* and *Agrobacterium tumifaciens* (Saxena *et al.*, 1994),

hyaluronic acid from *Streptococcus pyogenes* (DeAngelis *et al.*, 1993), type 3 CPS from *S. pneumoniae* (Arrecubieta *et al.*, 1995; Dillard *et al.*, 1995), alginate from *Pseudomonas aeruginosa* (Maharaj *et al.*, 1993) and chitin from *Saccharomyces cerevisiae* and *Candida albicans* (Cabib *et al.*, 1983) are all polysaccharides with simple mono- or di-saccharide repeat unit structures. Biosynthesis of these polysaccharides requires a processive glycosyl transferase or synthase. These glycosyl transferases share some common structural features including four transmembrane domains and a large hydrophilic domain with two distinct catalytic sites capable of forming two different glycosidic linkages (Keenleyside and Whitfield, 1996). However, there may be some differences in the mechanism of polysaccharide biosynthesis between different synthases. Biosynthesis of hyaluronic acid in *S. pyogenes* and type 3 CPS from *S. pneumoniae* have been shown to require only the synthase itself for production of CPS (Arrecubieta *et al.*, 1996; DeAngelis and Weigel, 1994). Whereas, studies on the biosynthesis of O:54 LPS from *S. enterica* serovar borreze have shown that the initial steps of polysaccharide synthesis resemble those for ABC-transporter dependent pathways (Keenleyside and Whitfield, 1996). Biosynthesis of O:54 LPS is dependent on the presence of a functional *rfe* gene, thus the initial step in polysaccharide synthesis involves the transfer of sugar 1-phosphate to und-P. The polysaccharide chain is then polymerised by the synthase. No dedicated export systems have been identified for these polysaccharides. These proteins all have four C-terminal transmembrane domains (Cabib *et al.*, 1983) which could conceivably form a pore in the cytoplasmic membrane, enabling extrusion of the polysaccharide chain as it is synthesised.

### **1.8.5 Polysaccharide synthesis not involving lipid-linked intermediates**

Some EPSs, such as glucan (a Glc polymer) and fructan (a fructose polymer) from *S. mutans* and levan (a homopolymer of fructose) from *E. amylovora* are synthesised

extracellularly and do not involve activated sugar precursors or lipid-linked intermediates (Gross *et al.*, 1992; Mukasa, 1986).

In *S. mutans*, glucan is synthesised from sucrose by glucosyl transferases (GTFs) which are exported from the cell via a signal peptide and convert sucrose to glucan and fructose. They have a catalytic site thought to be involved in sucrose hydrolysis (Mooser *et al.*, 1991) and a series of C-terminal direct repeats thought to be involved in glucan binding (Mooser and Wong, 1988). These repeats are similar to those identified in the C-termini of other Gram-positive ligand-binding proteins, such as those found in pneumococcal choline binding proteins (Wren, 1991).

### 1.8.6 Capsule gene clusters

Capsule gene clusters have been cloned and sequenced from a variety of Gram-negative and Gram-positive bacteria. These include *E. coli* K1 (Silver *et al.*, 1984), K4 (Drake *et al.*, 1990), K5 (Petit *et al.*, 1995) K7, K12, K92 (Roberts *et al.*, 1986), K10, and K54 (Pearce and Roberts, 1995), *Haemophilus influenzae* type b (Kroll *et al.*, 1989), *N. meningitidis* group B (Frosch *et al.*, 1989), *Salmonella typhi* Vi antigen (Hashimoto *et al.*, 1993), *E. amylovora* (Bugert and Geider, 1995), *Erwinia stewartii* (Dolph *et al.*, 1988), *Pseudomonas solanacearum* (Huang and Schell, 1995), *K. pneumoniae* (Arakawa *et al.*, 1995), and *Staphylococcus aureus* types 1 (Lin *et al.*, 1994), 5 (Sau *et al.*, 1997), and 8 (Sau and Lee, 1996). The genes involved in capsule biosynthesis for all of these organisms are clustered together on the chromosome. The organisation of these gene clusters share some similarities. Generally, the genes which are common to all serotypes are located near the 5' and 3' ends of the locus, with the serotype specific genes in the centre. This arrangement may have contributed to the high degree of capsule diversity, with recombination events between the conserved regions resulting in acquisition

of new genes and the expression of a new capsule serotype.

## 1.9 The Capsule Locus of *S. pneumoniae*

As outlined above, capsule production requires a complex pathway including transport into the cell and/or synthesis of the component monosaccharides, activation of each to a nucleotide precursor, co-ordinated transfer of each sugar, in sequence, to the repeating oligosaccharide and subsequent polymerisation, export and attachment to the cell surface. Thus, the loci encoding capsule production would be expected to consist of a large number of genes. However, when the work for this thesis was initiated, very little was known about these loci. Only a single gene from the type 3 capsule locus (Arrecubieta *et al.*, 1994) and the first seven genes of the type 19F locus (Guidolin *et al.*, 1994) had been cloned and sequenced. Since then the entire sequences for the capsule loci of several different pneumococci have been completed. These will be discussed in detail in chapter 8.

### 1.9.1 Pneumococcal capsular transformation

Classical genetic studies carried out by Austrian *et al.*, (1959) demonstrated that the *S. pneumoniae* genes required for biosynthesis and expression of capsular polysaccharide are closely linked on the pneumococcal chromosome and are transferred as a cassette during transformation. In the vast majority of cases, incorporation of the donor *cps* locus into the recipient cell chromosome was accompanied by loss of the original *cps* locus, consistent with genetic recombination between homologous sequences flanking the serotype-specific *cps* genes.

Transformation of a non-encapsulated derivative of a type 3 strain with DNA from

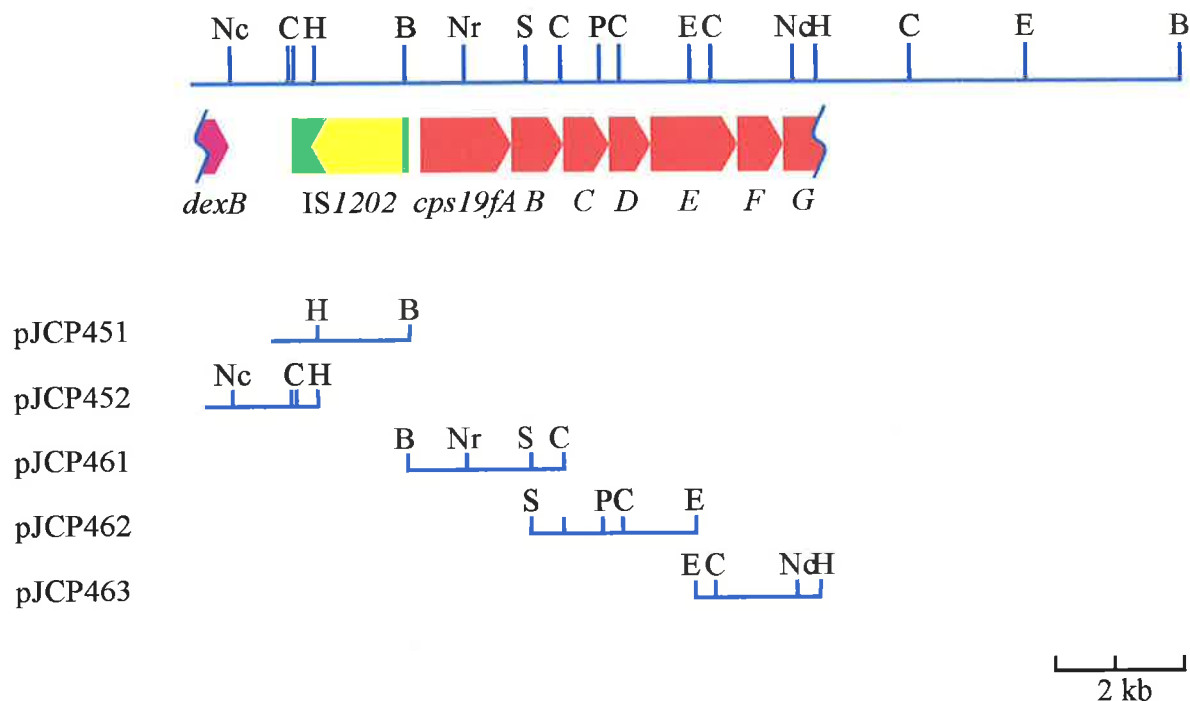
a type 1 strain occasionally resulted in encapsulated transformants expressing both type 1 and type 3 CPS. Such binary encapsulated strains contained the defective type 3 *cps* locus as well as the functional type 1 locus (Austrian *et al.*, 1959). Binary encapsulation could also be achieved by transformation of the type 3 mutant with DNA from other serotypes whose CPS contained GlcA, a component of type 3 CPS. The recipient strains used in these studies had mutations in the gene which encoded the enzyme required for synthesis of UDP-GlcA. Thus, in the binary transformants, the UDP-GlcA required for synthesis of type 3 CPS could be supplied by the enzyme encoded by the second *cps* locus (Bernheimer *et al.*, 1967).

Binary transformants always expressed abundant amounts of the type 3 CPS but only small amounts of the second CPS and represented atypical recombination of the second *cps* locus into the type 3 chromosome (Bernheimer *et al.*, 1967). However, unrelated pneumococci could only be transformed to the binary phenotype by using DNA from binary strains, never by using a mixture of the DNA from the two component serotypes. This suggests a degree of linkage between the two *cps* loci in the binary strains, and possibly a tandem arrangement in some cases (Bernheimer *et al.*, 1967). Binary transformants were only ever achieved using type 3 pneumococci, thus the type 3 *cps* locus and/or its flanking sequences may differ in some way to that of other pneumococci.

### 1.9.2 Localisation of the type 19F capsule locus

The fortuitous discovery of the insertion sequence *IS1202* closely linked to the type 19F capsule locus (Morona *et al.*, 1994a) enabled the isolation and characterisation of part of this locus. As shown in **Fig. 1.4**, the capsule locus is located downstream from *dexB*, a gene which is part of the multiple sugar metabolism locus and encodes a glucan-1,6- $\alpha$ -glucosidase. *IS1202* is located in the 1 kb intergenic region and is only 87 bp from the start

of the first gene in the capsule locus.



**Fig. 1.4** Physical map of the *S. pneumoniae* chromosome in the vicinity of the 19F *cps* locus (Guidolin *et al.*, 1994). Arrows represent potential ORFs and the box represents *IS1202*. Gene designations are indicated below the map; genes *cps19fB-G* are abbreviated to *B-G*, respectively. Restriction sites are as follows: B, *Bam*HI; C, *Cla*I; E, *Eco*RI; H, *Hind*III; Nc, *Nco*I; Nr, *Nru*I; P, *Pst*I; S, *Sph*I. Various subclones of this region, designated pJCP451-2, pJCP461-3, are as shown.

### 1.9.3 Partial cloning and sequencing of the 19F capsule locus

The first part of the type 19F capsule locus was cloned and sequenced (Guidolin *et al.*, 1994). Examination of the sequence data revealed the presence of six complete open reading frames (ORFs) and one incomplete ORF, as shown in **Fig. 1.4**. An almost perfect consensus promoter sequence was detected 30 bp upstream of the first ORF in this operon. Each of the ORFs is preceded by a ribosome binding site and they are all closely coupled, being separated by 1-15 nucleotides.

## 1.9.4 Characterisation of *cps19fA-G*

The genes *cps19fA-G* are part of an operon essential for type 19F capsule biosynthesis as disruption of these ORFs resulted in the loss of type 19F capsule production (Guidolin *et al.*, 1994). The predicted functions of the seven *cps19f* ORFs, based on database homology searches (summarised in **Table 1.9** and discussed in detail below), are also consistent with the involvement of these proteins in capsule biosynthesis. Thus the genes *cps19fA-G* form the first part of the type 19F capsule locus.

**Table 1.9** Predicted functions of *S. pneumoniae cps19f* ORFs

ORF	Predicted function	Similar proteins (% identity)
Cps19fA	regulation?	<i>B. subtilis</i> LytR <sup>a</sup> (27.6%)
Cps19fB	unknown	<i>S. agalactiae</i> CpsA <sup>b</sup> (63.8%)
Cps19fC	chain length regulation /export	<i>S. agalactiae</i> CpsB <sup>b</sup> (46%) <i>R. meliloti</i> ExoP <sup>c</sup> (22.5% to N-terminus)
Cps19fD	chain length regulation /export	<i>S. agalactiae</i> CpsC <sup>b</sup> (59.2%) <i>R. meliloti</i> ExoP <sup>c</sup> (29.9% to C-terminus)
Cps19fE	glucose-1-phosphate transferase	<i>S. agalactiae</i> CpsD <sup>b</sup> (45.6%) <i>S. enterica</i> serovar typhimurium RfbP <sup>d</sup> (31.4%) <i>X. campestris</i> Xps2a <sup>e</sup> (24.1%)
Cps19fF	<i>N</i> -acetyl mannosamine transferase	<i>B. subtilis</i> TagA <sup>f</sup> (31.2%) <i>E. coli</i> K12 RffM <sup>g</sup> (33%) <i>S. enterica</i> serovar typhimurium RffM <sup>h</sup> (32.5%)
Cps19fG <sup>j</sup>	unknown	<i>Haemophilus influenzae</i> LicD <sup>i</sup> (38%)

a. Lazarevic *et al.* (1992); b. Rubens *et al.* (1993); c. Glucksmann *et al.* (1993); d. Jiang *et al.* (1991); e. Ielpi *et al.* (1993); f. Mauel *et al.* (1991); g. Daniels *et al.* (1992); h. Lu and Abdelal (1993); i. Weiser *et al.* (1989); j. Cps19fG ORF is truncated.

### 1.9.4.1 Cps19fA

Cps19fA exhibits significant similarity (27.6% identity) to the entire *B. subtilis* *lytR* gene product (Lazarevic *et al.*, 1992). *LytR* is a basic protein that is thought to act as a transcriptional attenuator of *lytABC* (autolysin) expression and is thought to be membrane bound via a N-terminal anchoring domain. Cps19fA is also basic and contains three

hydrophobic segments near its N-terminus which may be membrane anchoring domains. Interestingly, the predicted LytR protein lacks this hydrophobic region of the sequence, but has a single hydrophobic region near its N-terminus. Thus Cps19fA may play a role as a regulator of capsule gene expression although there is no direct evidence for this.

#### **1.9.4.2 Cps19fB**

The *cps19fB* gene product has 63.8% identity with the product of the *Streptococcus agalactiae* (group B streptococcus) *cpsA* gene. The *S. agalactiae* *cpsA* gene has been shown to be involved in the production of type III capsular polysaccharide in group B streptococci (Rubens *et al.*, 1993). No other significant sequence similarities were detected, and the functions of either of these gene products is not known.

#### **1.9.4.3 Cps19fC**

The *cps19fC* gene product has 46% identity with the *S. agalactiae* CpsB protein (Rubens *et al.*, 1993), and 22.5% identity with the N-terminal portion of the *R. meliloti* ExoP protein (Glucksmann *et al.*, 1993). Cps19fC also has some similarity to proteins involved in chain length regulation of polysaccharides, including a conserved motif (amino acids 168-187 in Cps19fC) found near the C-terminus in this type of protein (Kolkman *et al.*, 1997b; Becker *et al.*, 1995). Cps19fC probably regulates the chain length of type 19F polysaccharide.

#### **1.9.4.4 Cps19fD**

The *cps19fD* gene product has significant similarity (59.2% identity) to the *S. agalactiae* *cpsC* gene product (Rubens *et al.*, 1993). In addition, the Cps19fD protein exhibits similarity (29.9% identity) to the C-terminal portion of the *R. meliloti* ExoP protein (Glucksmann *et al.*, 1993). Cps19fC and Cps19fD are homologous to first and last thirds of ExoP, respectively. This may indicate that ExoP has multiple functional domains, which are encoded by separate genes in streptococci. Cps19fD and the C-terminal domain

of ExoP contain the ATP-binding motifs GXGKTT and YIIVD (amino acids 46-51 and 148-152 in Cps19fD) which have been described previously (Guidolin *et al.*, 1994; Kolkman *et al.*, 1997b). *R. meliloti* ExoP-truncated mutants (missing the C-terminal domain) produce less succinoglycan with a reduced molecular weight (Becker *et al.*, 1995). This suggests that Cps19fD may function, in conjunction with Cps19fC, to regulate the export and chain length of type 19F CPS.

#### 1.9.4.5 Cps19fE

The *cps19fE* gene product is relatively hydrophobic and exhibits significant similarity to known glycosyl transferases. The galactosyl transferase encoded by the *S. agalactiae cpsD* gene (Rubens *et al.*, 1993) shows 45.6% identity to the *cps19fE* gene product. Cps19fE also exhibits similarity to RfbP (31.4%) from *S. enterica* serovar typhimurium strain LT2 (Jiang *et al.*, 1991), an und-P Gal-1-phosphate transferase which catalyses the initial step in O-antigen biosynthesis, and the *X. campestris xps2a (gumD)* gene product (24.1%), a Glc-1-phosphate transferase which catalyses the first step in xanthan gum synthesis (Ielpi *et al.*, 1993). Examination of the hydrophobicity plot of Cps19fE revealed that the N-terminal portion of the protein has 3 hydrophobic segments which may be potential membrane spanning domains; this suggests that Cps19fE is anchored to the bacterial membrane. Interestingly, although the N-terminal portions of Cps19fE, RfbP, and Xps2a have little amino acid sequence similarity, the hydropathy plots for these regions are very similar. Studies with RfbP (Wang *et al.*, 1996) have demonstrated that the hydrophilic C-terminal region is required for the transferase activity of these proteins. Kolkman *et al.* (1996) demonstrated that Cps14E, which is almost identical (95.8%) to Cps19fE, is a glucosyl transferase. Thus Cps19fE is a glucosyl-transferase which adds Glc, the first sugar in the polysaccharide repeat unit, to the lipid carrier.

#### 1.9.4.6 Cps19fF

The *cps19fF* gene product has similarity (31-33% identity) to three other proteins, namely the *B. subtilis tagA* gene product, an *N*-acetyl mannosamine (ManNAc) transferase which is needed for cell wall teichoic acid biosynthesis (Mauel *et al.*, 1991), and the *E. coli* K12 and *S. enterica* serovar typhimurium strain LT2 *rffM* gene product, which is a putative *N*-acetyl-D-mannosaminuronic acid transferase, involved in the synthesis of ECA (Daniels *et al.*, 1992; Lu and Abdelal, 1993). Since the type 19F CPS contains ManNAc, Cps19fF is probably the transferase which catalyses the addition of this sugar in the synthesis of the type 19F polysaccharide.

#### 1.9.4.7 Cps19fG

The *cps19fG* gene (incomplete) encodes a truncated protein which exhibits 38% identity with the N-terminal end of the LicD protein of *Haemophilus influenzae* (Weiser *et al.*, 1989). However, the precise function of LicD is unknown.

## 1.10 Aims of this Thesis

As mentioned previously (section 1.9), when this study commenced, the only capsule biosynthesis genes from *S. pneumoniae* that had been analysed were a single gene from the type 3 locus (Arrecubieta *et al.*, 1994) and what appeared to be the first six genes of the type 19F locus (Guidolin *et al.*, 1994). Knowledge of the factors determining or regulating capsular polysaccharide biosynthesis is likely to provide new insights into the mechanisms of pathogenesis of pneumococcal disease. Moreover, comparison of the genetic structure of the loci encoding synthesis of immunologically-related capsular polysaccharides may provide information on the mechanism whereby *S. pneumoniae*

acquired the capacity to synthesise 90 distinct capsular serotypes. Accordingly, the work described in this thesis was aimed at complete characterisation of the *cps* loci for all four members of *S. pneumoniae* serogroup 19 (serotypes 19F, 19A, 19B and 19C).

## Chapter 2

# MATERIALS AND METHODS

## 2.1 Bacterial strains and cloning vectors

### 2.1.1 Bacterial strains

The *E. coli* strains used in this study are described in **Table 2.1**. Bacteriophage N4 was obtained from Prof. L.B. Rothman-Denes and Dr. D.R. Kiino (Kiino *et al.*, 1993).

**Table 2.1** *E. coli* K12 strains

Strain	Relevant genotype	Source / Reference
DH1	F <sup>-</sup> , <i>gyrA96</i> , <i>recA1</i> , <i>relA2</i> , <i>endA1</i> , <i>thi1</i> , <i>hsdR1</i> , <i>supE44</i> , $\lambda$ <sup>-</sup>	Hanahan (1983)
DH5	F <sup>-</sup> , <i>deoR</i> , <i>recA1</i> , <i>endA1</i> , <i>hsdR17</i> ( $r_k^-$ , $m_k^-$ ), <i>supE44</i> , <i>thi1</i> , <i>gyrA96</i> , <i>relA1</i> , $\lambda$ <sup>-</sup>	Bethesda Research Laboratories, Gaithersburg, MD
DH5 $\alpha$	F <sup>-</sup> , $\phi$ 80 <i>dlacZ</i> $\Delta$ M15, $\Delta$ ( <i>lacZYA-argF</i> ), U169, <i>deoR</i> , <i>recA1</i> , <i>endA1</i> , <i>hsdR17</i> ( $r_k^-$ , $m_k^-$ ), <i>supE44</i> , <i>thi1</i> , <i>gyrA96</i> , <i>relA1</i> , $\lambda$ <sup>-</sup>	Bethesda Research Laboratories, Gaithersburg, MD
MC4100	F <sup>-</sup> , <i>araD139</i> , $\Delta$ ( <i>argF-lac</i> )U169, <i>rpsL150</i> , <i>rclA1</i> , <i>flbB5301</i> , <i>deoC1</i> , <i>ptsF25</i> , <i>rbsR</i> , <i>thiA1</i>	<i>E. coli</i> Genetic Stock Center, Yale University, CT (Silhavy <i>et al.</i> , 1984)
KI8828	MC4100 <i>nfrC2</i>	Prof. L.B. Rothman-Denes and Dr. D.R. Kiino (Kiino <i>et al.</i> , 1993)
SØ874	F <sup>-</sup> , <i>lacZ2286</i> , <i>trp-49</i> , $\Delta$ ( <i>sbcB-rfb</i> )86, <i>upp-12</i> , <i>relA1</i> , <i>rspL150</i> , $\lambda$ <sup>-</sup>	Neuhard and Thomassen (1976)

*S. pneumoniae* strains and clinical isolates used in this study are described in **Table 2.2**, all other clinical isolates were from the Department of Microbiology, Women's and Children's Hospital, Adelaide, South Australia.

**Table 2.2** *S. pneumoniae* strains

Strain	Serotype	Source / Reference
Rx1	rough	Shoemaker and Guild (1974)
Rx1-19F	19F	Morona <i>et al.</i> (1994a)
SSZ-19F	19F	Dr. Chi-Jen Lee, Center for Biologics, FDA, Bethesda, Md., USA
1777/39 (19A1)	19A	Dr. Jorgen Henrichsen, Statens Seruminstitut, Copenhagen, Denmark
19A2	19A	Dr. Chi-Jen Lee, Center for Biologics, FDA, Bethesda, Md., USA
clinical isolates	19A	Mike Gratten, ARI Research and Reference Unit, Centre for Public Health Sciences, Queensland Health, Brisbane, Australia
19B	19B	Dr. Chi-Jen Lee, Center for Biologics, FDA, Bethesda, Md., USA
19C	19C	Dr. Chi-Jen Lee, Center for Biologics, FDA, Bethesda, Md., USA

FDA, Food and Drug Administration; ARI, Acute Respiratory Infections.

### 2.1.2 Cloning vectors

The bacterial cloning vectors used in this study are listed in **Table 2.3**.

**Table 2.3** Bacterial plasmids and cloning vectors

Plasmid/Vector	Antibiotic Resistance	Reference
pBluescript SK+	Amp	Stratagene, La Jolla, CA.
pBluescript KS+	Amp	Stratagene, La Jolla, CA.
pGEM-7Zf(+)	Amp	Promega Corporation, Madison, WI.
pK184	Kan	Jobling and Holmes (1990)
pK194	Kan	Jobling and Holmes (1990)
pVA891	Ery, Cml	Macrina <i>et al.</i> (1983)
pGP1-2	Kan	Tabor and Richardson (1985)

Amp, ampicillin; Kan, kanamycin; Ery, erythromycin; Cml, chloramphenicol.

### 2.1.3 Growth media

*E. coli* strains were grown in Luria-Bertani (LB) broth with or without 1.5% (w/v) Bacto-agar (Difco Laboratories, Detroit, MI, USA). LB broth was prepared by dissolving 10 g tryptone (Difco), 5 g yeast extract (Difco), and 5 g NaCl in deionised water to a final volume of 1 litre, adjusted to pH 7.5 with NaOH, and sterilised by autoclaving at 121°C for 15 min. LB agar was prepared similarly, but with the addition of 15 g agar before autoclaving. Where appropriate, chloramphenicol (Cml), ampicillin (Amp), kanamycin (Kan), or erythromycin (Ery) was added to the growth medium at concentrations of 25, 50, 25 and 10 µg/ml, respectively.

*E. coli* strains were preserved as follows: a 1.5 ml aliquot from a fresh overnight culture was mixed with 375 µl 80% (v/v) glycerol and stored at -80°C.

Pneumococci were routinely grown in either Todd-Hewitt broth supplemented with 0.5% (w/v) yeast extract (THY), or on blood agar (BA). THY was prepared by dissolving 36.4 g of Todd-Hewitt broth (Oxoid Limited, Basingstoke, England) and 5 g yeast extract (Difco) in deionised water to a final volume of 1 litre, and sterilised by autoclaving. BA was prepared by dissolving 40 g blood agar base no. 2 (Oxoid) in deionised water to a final volume of 1 litre, sterilised by autoclaving and cooled to 55°C before the addition of 25 ml of defibrinated horse blood (Amadeus International, Adelaide, Australia). Where appropriate, Ery was added to media at a concentration of 0.2 µg/ml. When pneumococci were grown to stationary phase in THY, 1% (w/v) choline (Aldrich Chemical Company, NSW, Australia) was added to the medium.

Pneumococcal strains were grown in serum broth for 4 h at 37°C and then stored at -80°C. Serum broth was prepared by dissolving 10 g Bacto peptone (Difco), 5 g NaCl and 8 g Lab Lemco powder (Oxoid) into 1 litre of deionised water, and the pH was adjusted to 7.4. After sterilisation by autoclaving, 10% (v/v) donor horse serum (CSL, Victoria,

Australia) was added, and the broth was filter sterilised.

## **2.2 Chemicals and reagents**

### **2.2.1 General chemicals**

Most chemicals used were AnalaR grade and were purchased from Ajax Chemicals, (NSW, Australia). Tris was purchased from Progen Industries (Queensland, Australia). Maleic acid, SDS, DTT and N,N,N',N'-tetramethyl-ethylene-diamine (TEMED) were purchased from Sigma Chemical Company (St. Louis, MD., USA). Acrylamide and ammonium persulphate were purchased from Bio-Rad Laboratories (Hercules, CA, USA). The deoxyribonucleoside triphosphates (dATP, dCTP, dGTP and dTTP), X-gal, IPTG, and herring sperm DNA were purchased from Boehringer Mannheim (Mannheim, Germany). DOC was purchased from BDH Biochemicals (Poole, England).

### **2.2.2 Antibiotics**

Amp was purchased from CSL (Victoria, Australia); Kan sulphate was purchased from Progen Industries; Cml and Ery were purchased from Boehringer Mannheim; rifampicin (Rif), tetracycline (Tet) and streptomycin sulphate (Str) were purchased from Sigma.

### **2.2.3 Enzymes**

RNase A, lysozyme, pronase and proteinase K were purchased from Boehringer-Mannheim. All restriction endonucleases were purchased from either Boehringer

Mannheim, or Progen and used according to the manufacturers' recommendations. Other DNA modifying enzymes, T4 DNA ligase, Klenow fragment of DNA polymerase I and Taq polymerase were purchased from Boehringer Mannheim. Sequencing kits using either dye-labelled primer or dye-labelled terminators were purchased from Applied Biosystems.

## 2.2.4 Oligodeoxynucleotides

The oligodeoxynucleotides used are listed in **Table 2.4** and were purchased from Life Technologies (Gaithersburg, MD, USA).

**Table 2.4 Oligodeoxynucleotides**

Name	Sequence (5'-3')	Homologous to:
DEXB	ataaagcTTCCATGGGATGCTTTCTGT	<i>dexB</i> (35-16 nt from 3' end)
CPSA2	ataggatCCCTAGCAAGGCAACTAGTA	< <i>cps19fA</i> (nt 335-354)
CPS5'	TGATGTTCAAGGTATAGGTGTTAATCA	before <i>cps19fA</i> (nt 146-169)
CPSX1	atactcgAGATGGAGCACCAAATCATTC	<i>cps19fG</i> (nt 6,245-6,265)
J5	atgggatCCACTACACTTGATGGGGGT	<i>cps19fF</i> (nt 5,224-5,243)
J7	TTATTATATTTTgcATgCCTTTAATGGC	<i>cps19fI</i> (nt 7,772-7,799)
J8	TAAAACTAGACaAGCTyTAGCAAAA	< <i>cps19fE</i> (nt 4,584-4,608)
J9	TTCCTTCTaGATTTGTAAAGATATT	<i>cps19fJ</i> (nt 9,502-9,526)
J10	CTAAGAACAAGcTTGTATATTTCCCT	< <i>cps19fI</i> (nt 7,919-7,943)
J11	GATggATCCCATCAATCCAACCCAAGT	< <i>cps19fJ</i> (nt 9,108-9,129)
J12	CTCGCGCTGCAgCAAAACAAC	<i>cps19fL</i> (nt 11,597-11,614)
J14	CAACTTTgAAITCACCAGTTATTTTTTG	<i>cps19fH</i> (nt 7,515-7,542)
J15	GACTCAgGTAcCAAATTTAAAAGAAGAA	<i>cps19fI</i> (nt 8,881-8,907)
J16	AAATGGATccAAAGGTAGAAATGTAATC	< <i>cps19fJ-K</i> (nt 10,379-10,406)
J17	CACGGgGATcCTACTACGACATAT	<i>cps19fK</i> (nt 10,697-10,720)
J18	AATGATCTGcATgCGTTCTGTATCTC	< <i>cps19fK</i> (nt 11,166-11,191)
J19	AATAGTCAGTCgaCCTACTATTGGAC	<i>cps19fO</i> (nt 15,340-15,364)
J20	AAGTAgGAtCCTGATATTTCCCAGC	< <i>cps19fG</i> (nt 6,096-6,120)

(continued overleaf)

Table 2.4 Oligodeoxynucleotides (continued)

Name	Sequence (5'-3')	Homologous to:
J21	TGCTGGA <u>Agc</u> TTGAAATTGGTTGGG	<i>cps19fK</i> (nt 11,341-11,365)
J22	AATTG <u>Gaa</u> TTCTTTTATAGATTTAACACAAG	< <i>cps19fH</i> (nt 6,743-6,772)
J24	TTGAGAA <u>AGc</u> TTTTGAAGAGGTGAAA	<i>cps19fO</i> (nt 14,260-14,285)
J25	GTACG <u>ga</u> TCCAAGTCGGACGACC	< <i>cps19fO</i> (nt 14,682-14,704)
J26	TAGTGAGAA <u>Tc</u> TCTATCCTATCTTCTA	<i>aliA</i> (nt 16,127-16,154)
J27	GGCT <u>g</u> GatCCAAACCATCTGGACTT	< <i>cps19fL</i> (nt 11,787-11,811)
J28	GAAAAA <u>a</u> GCTTCCACTTGGATTTC	<i>cps19fM</i> (nt 12,541-12,566)
J29	GCTTGGATC <u>c</u> GCGTAGTTCACAAA	< <i>cps19fM</i> (nt 12,880-12,903)
J30	TGACTT <u>Aa</u> GCTtGAAACCGCGCGA	<i>cps19fN</i> (nt 13,136-13,159)
J31	GCCT <u>g</u> GAiCCAGCTTCAAAGTTG	< <i>cps19fN</i> (nt 13,946-13,946)
J32	CGTCTTTT <u>ac</u> TAGTTGTTGGATATCAAT	< <i>aliA</i> (nt 17,853-17,880)
J33	ACAAGA <u>A</u> TCgATTACGTATTTAGTTGGT	<i>aliA</i> (nt 17,339-17,366)
J34	CTATTTGGATCCGTAGCTTTGGCA	< <i>aliA</i> (nt 17,057-17,080)
J36	CAATAATGTCACGCCCGCAAGGGCAAGT	< <i>aliA</i> (nt 16,463-16,490)
J37	TGGTGGAA <u>AGC</u> ttAGAAAGAAGCTG	<i>cps19fN</i> (nt 14,002-14,026)
J38	TCTAGCG <u>Ga</u> TCcCAGTTTACCAAGC	<IG3' region (nt 15,568-15,592)
J39	TAGTTCATGTAGTTGCAAGTGACATGCACAA	<i>cps19fB</i> (nt 2,190-2,220)
J42	AAGTC <u>a</u> TcGATGAAACTATTTCTTGTG	<i>cps19bH</i> (nt 1,499-1,525)
J43	TATATGATTGT <u>A</u> TcGATTATCATGTGGC	<i>cps19bP</i> (nt 2,693-2,720)
J44	TAATCAGCTGAT <u>c</u> gATCAGCTCCGCTC	< <i>cps19bP</i> (nt 3,395-3,421)
J45	AATGG <u>A</u> tcCATTTTTATCGCTTCTGGAC	< <i>cps19bK</i> (nt 9,273-9,302)
J47	CCTATACCATTTG <u>a</u> ATTCAATATAATTCTAG	< <i>cps19bJ</i> (nt 7,830-7,857)
J49	AATGTG <u>g</u> AtCCTACATATGTATTAACAG	<i>cps19bR</i> (nt 6,886-6,913)
J70	GATGGT <u>gga</u> TCCTGTTTTAGATTTATTTGG	<i>cps19aK</i> (nt 12,114-12,143)
J72	TATCGTCA <u>t</u> CGATAAAATTCTTCACCG	< <i>cps19aL</i> (nt 13,515-13,540)
J87	CTCATACT <u>a</u> GTCCACTGTTGGTG	< <i>cps19aE</i> (nt 6,458-6,480)
J88	CTCAGG <u>at</u> CCCGCCTGTACCC	<i>cps19aL</i> (nt 13,266-13,286)
J92	TTGGATAAGCTtGAAAAAATCGGATT	<i>cps19aB</i> (nt 3,441-3,467)
J93	TACTGAA <u>a</u> GCTTTGCTAGTTTTTCAC	< <i>cps19aE</i> (nt 5,330-5,355)
J94	AACACGATAGAAATCGATGTATTTTC	<i>cps19aC</i> (nt 3,932-3,956)
J95	TTTTG <u>g</u> AtCCTGCACGCGCAAAAG	< <i>cps19aD</i> (nt 4,797-4,820)

< denotes primer is complementary to the sequence of the indicated gene. Lower case letters indicate base substitutions to insert restriction sites. Cleavage sites for restriction endonucleases *EcoRI* (GAATTC), *BamHI* (GGATCC), *HindIII* (AAGCTT), *ClaI* (ATCGAT), *XbaI* (TCTAGA), *PstI* (CTGCAG), *SphI* (GCATGC), *SalI* (GTCGAC), *KpnI* (GGTACC), *SpeI* (ACTAGT) and *XhoI* (CTCGAG) are underlined. The sequences to which the primers are homologous are available as follows: *dexB-cps19fA* and *cps19fA-G* are available under GenBank accession nos U04047 and U09239, respectively; *cps19fG-O*, *cps19b* and *cps19a* are as described in Appendices I, II and III, respectively, nt denotes nucleotide position.

## 2.3 Serotyping of pneumococcal strains

Production of capsule by pneumococci was assessed by quellung reaction, using diagnostic pneumococcal typing sera produced by Statens Seruminstitut, Copenhagen, Denmark. Serum broth was inoculated with cells from an overnight BA plate and incubated at 37°C for 4 h. Then 10 µl serum broth culture was mixed with 10 µl of the typing sera and microscopic examination for capsular swelling confirmed the production of type-specific capsule. This was performed by staff from the Dept. of Microbiology and Infectious Diseases, Women's and Children's Hospital, Adelaide.

Pneumococci belonging to serogroup 19 were serotyped (19F, 19A, 19B or 19C) by quellung reaction, using factor specific antisera obtained from Statens Seruminstitut, Copenhagen, Denmark, and was performed by Mr. M. Gratten and Ms. Denise Murphy, Public Health Bacteriology Laboratory, Centre for Public Health Sciences, Queensland Health.

## 2.4 Bacterial transformation

### 2.4.1 Preparation of competent *E. coli* strains

#### 2.4.1.1 RbCl<sub>2</sub> method

*E. coli* K12 DH5α was made competent as follows. A 100 ml LB broth was inoculated with 1 ml of an overnight culture of DH5α and grown to very early log phase ( $A_{600}$  0.04). The cells were then centrifuged in a Hettich Universal/K2S centrifuge at 4,000 x g for 15 min at 4°C. The pelleted cells were resuspended in 40 ml of TfbI (30 mM

KOAc, 100 mM RbCl<sub>2</sub>, 10 mM CaCl<sub>2</sub>, 50 mM MnCl<sub>2</sub>, adjusted to pH 5.8 with 0.2 M acetic acid and filter sterilised) and placed on ice for 5 min. The cells were again centrifuged at 4,000 × g for 15 min at 4°C. The cells were then resuspended in 2 ml of TfbII (10 mM MOPS, 75 mM CaCl<sub>2</sub>, 10 mM RbCl<sub>2</sub>, 15% (v/v) glycerol, adjusted to pH 6.5 with KOH and filter sterilised) and placed on ice for 15 min. The cells were dispensed into sterile 1.5 ml microcentrifuge tubes in 50 µl volumes, quickly frozen on ethanol/dry ice for 5 min, and transferred to -70°C. The bacterial cells remained viable with no apparent loss of competence for at least 6 months. The cells were thawed at room temperature before being placed on ice and transformed as described in section 2.4.2.

#### **2.4.1.2 Electroporation method**

DH5(pGP1.2) was made competent for transformation by electroporation, as follows. Briefly, 100 ml fresh LB broth was inoculated with 10 ml of an overnight broth culture. The culture was grown at 30°C with vigorous shaking until a density of A<sub>600</sub> = 0.5 to 0.8 was achieved. Before harvesting, the bacterial cells were chilled on ice for 15 min and then centrifuged at 4,000 × g for 15 min at 4°C. The supernatant was drained away as much as possible and the pellet was resuspended in 100 ml ice-cold water and centrifuged as before. The last step was repeated, and the cells were resuspended in 50 ml ice cold water. The resultant pellet was then resuspended in 20 ml of ice-cold 10% (v/v) glycerol, centrifuged as before and resuspended in a final volume of 2 ml of ice-cold 10% (v/v) glycerol. Cells were quickly frozen in ethanol/dry ice in 100 µl aliquots in 1.5 ml microfuge tubes, and stored at -70°C. The cells were stable for about six months under these conditions.

#### **2.4.1.3 CaCl<sub>2</sub> method**

All other *E. coli* strains were made competent for transformation with plasmid DNA according to the method described by Brown *et al.* (1979). An overnight culture (in

LB broth) was diluted 1:20 into fresh LB broth and incubated with aeration for 2 h until the culture reached an  $A_{600}$  of 0.04 (very early log phase). The bacterial cells were pelleted by centrifugation at  $4,000 \times g$  for 15 min at  $4^{\circ}\text{C}$  and resuspended in 0.5 volumes of cold 100 mM  $\text{MgCl}_2$ , centrifuged again and resuspended in a 0.1 volumes of cold 100 mM  $\text{CaCl}_2$ . The bacterial cells were allowed to stand for 1 to 2 h on ice before use.

#### **2.4.2 Transformation of competent *E. coli* strains**

Competent cells (50  $\mu\text{l}$ ) were mixed with 15  $\mu\text{l}$  DNA and kept on ice for 30 min. The cell/DNA mixture was then heat-shocked at  $42^{\circ}\text{C}$  for 3 min, and then again placed on ice for a further 10 min. The bacterial cells were incubated at  $37^{\circ}\text{C}$  for 1 h after the addition of 0.5 ml LB. The bacterial cell suspension was then plated onto selection plates.

#### **2.4.3 Electroporation of *E. coli* strains**

Transformation of electroporation competent cells was performed as described in the Electrocell Manipulator<sup>®</sup> 600 (ECM<sup>®</sup> 600) Operation Manual (BTX, San Diego, CA., USA). Frozen bacteria were gently thawed at room temperature and placed on ice before the addition of 1 to 2  $\mu\text{g}$  of DNA (dissolved in Hi-pure water), the cells were mixed thoroughly and allowed to stand for 1 min. The cell suspension was transferred to a chilled sterile BTX disposable electroporation cuvette (P/N 620, 2 mm gap) and this was then placed in the safety chamber of the ECM<sup>®</sup> 600 which was set at 2.50 kV/resistance high voltage and 129  $\Omega$  resistance. The cell suspension was then pulsed for 5-6 msec at a charging voltage of 2.45 kV. After electroporation the cell suspension was mixed with 960  $\mu\text{l}$  LB, incubated at  $30^{\circ}\text{C}$  for 1 h without shaking, and then plated onto selective medium.

#### 2.4.4 Transformation of pneumococcal strains

*S. pneumoniae* strains were grown at 37°C, in 10 ml THY broth, to a density of  $3 \times 10^8$  cells/ml ( $A_{600}$  of 0.27). The culture was diluted 100-fold into competence medium containing 10% (v/v) glycerol and stored at -70°C in 500 µl aliquots. After thawing a 500 µl aliquot of cells, 500 µl of competence factor was added and incubated at 37°C for 20 min. The cells were then mixed with 100 µl of DNA prepared as described in section 2.5.2 or 2.5.3 and incubated at 37°C for 2 h. Transformants were selected by plating the cells on blood agar plates containing 0.2 µg/ml erythromycin and incubating at 37°C for up to 48 h.

#### 2.4.5 Preparation of competence factor for pneumococcal transformation

Pneumococcal strains produce one or more competence-inducing pheromones which activate competence and the uptake of exogenous DNA (Pozzi *et al.*, 1996). Production of these pheromones occurs spontaneously at a particular point in the growth cycle and when harvested can be used to induce competence in other pneumococcal strains. The highly transformable *S. pneumoniae* strain Rx1 was used to prepare competence factor, as described previously by Yother *et al.* (1986).

A 10 ml culture of Rx1 in THY was grown, at 37°C, to a density of  $3 \times 10^8$  cells ( $A_{600}$  of 0.27). A 5 ml aliquot, with 10% (v/v) glycerol, was stored at -70°C, for 16 h. The 5 ml aliquot was thawed at 37°C and diluted into 500 ml competence medium (THY; 0.2% (w/v) bovine serum albumin; 0.01%  $\text{CaCl}_2$ ; adjust to pH 7.8). The 500 ml culture was incubated at 37°C for 30 min and then 18 ml samples were taken at 10 min intervals. Each 18 ml sample was mixed with 2.6 ml 80% (v/v) glycerol and stored at -70°C immediately; a 100 µl aliquot of each time point was stored separately in a 1.5 ml microfuge tube for further analysis.

To assess competence, the 100  $\mu$ l samples were incubated at 37°C for 10 min, then 10  $\mu$ l chromosomal DNA (from an Ery resistant pneumococcus) was added and the samples were incubated for a further 30 min at 37°C. After incubation, serial 10-fold dilution to  $10^{-6}$  were carried out in a microtitre tray and plated onto BA plates containing 0.2  $\mu$ g/ml Ery, and incubated at 37°C for 16 h. The time points with the highest number of Ery-resistant transformants and thus the highest transformation frequencies were prepared for use as competence factor. The frozen aliquots were thawed and the bacteria pelleted by centrifugation at  $4,000 \times g$  for 15 min at 4°C. The supernatant (or competence factor) was carefully removed and filter sterilised through a 0.45  $\mu$ m filter, dispensed into 500  $\mu$ l aliquots and stored at -70°C.

## 2.5 DNA extraction procedures

### 2.5.1 Plasmid isolation

Plasmid isolation procedures were based on the alkaline lysis method (Morelle, 1989). A 1.5 ml aliquot of an *E. coli* overnight culture was transferred to a microcentrifuge tube, pelleted by centrifugation at 13,000 rpm for 2 min in a microfuge (Biofuge 13, Heraeus Instruments, Germany), and resuspended in 100  $\mu$ l of a solution of 50 mM Glc, 25 mM Tris-HCl (pH 8.0), 10 mM EDTA, and lysozyme at a final concentration of 4  $\mu$ g per ml. After 5 min incubation at room temperature, 200  $\mu$ l of a solution of 1% (w/v) SDS and 200 mM NaOH was added. The tubes were then mixed by inverting several times and incubated on ice for 5 min. After the addition of 400  $\mu$ l 7.5 M ammonium acetate, pH 4.8, the tubes were again mixed by inversion, incubated for 10 min on ice, and cellular debris

was then removed by centrifugation at 13,000 rpm for 10 min. Supernatants were then transferred to a clean microfuge tube and the plasmid DNA precipitated by the addition of 500  $\mu$ l of ice-cold isopropanol. The tubes were then mixed by inverting several times and left standing at room temperature for 15 min. The plasmid DNA was pelleted by centrifugation at 13,000 rpm for 15 min, the pellet was washed with 70% (v/v) ethanol, and dried under vacuum. The pellet was resuspended in 20-50  $\mu$ l of sterile water (Baxter, NSW, Australia) and stored at 4°C.

Plasmid DNA to be used in sequencing reactions (Applied Biosystems) was further purified as follows. The plasmid DNA prepared as above was mixed with 1  $\mu$ l RNase A (10 mg/ml in H<sub>2</sub>O) and incubated at 37°C for 15-30 min. The DNA was then mixed with an equal volume of 7.5 M ammonium acetate and 2 volumes of isopropanol. The sample was mixed by vortex, left standing at room temperature for 15 min, and then centrifuged for 15 min at 13,000 rpm. The supernatant was discarded, the DNA pellet was washed with 75% ethanol, 25% 50 mM Na acetate and microfuged for 10 min. The pellet was then washed in 100% ethanol and microfuged for 5 min. The pellet was dried under vacuum and resuspended in the original volume of TE (10 mM Tris, 1 mM EDTA, pH 8.0). Alternatively, they were prepared using the plasmid mini-kit (Qiagen), following the procedure outlined in the Qiagen plasmid handbook.

### **2.5.2 Plasmid isolation for pneumococcal transformation**

A 100 ml overnight culture was centrifuged at 4,000  $\times$  g for 15 min at 4°C and the pellet was resuspended in 4 ml 50 mM Glc, 25 mM Tris (pH 8.0), 10 mM EDTA (pH 8.0), 2 mg/ml lysozyme and incubated on ice for 10 min, then 8 ml 0.2 M NaOH, 1% (w/v) SDS was added, mixed by inversion and placed on ice for 10 min. Next, 6 ml 1.5 M Na acetate (pH 4.8) was added, mixed by inversion and placed on ice for a further 10 min. The debris

was pelleted by centrifugation at  $4,000 \times g$  for 15 min at  $4^{\circ}\text{C}$  and the supernatant was carefully decanted into a clean tube. The supernatant was mixed with 6 ml of ice-cold isopropanol and placed on ice for 15 min. The DNA was pelleted by at  $4,000 \times g$  for 20 min at  $4^{\circ}\text{C}$ , resuspended in 2 ml of 0.1 M Na acetate, 0.05 M Tris (pH 8.0) and transferred to microfuge tubes (500  $\mu\text{l}$  per tube) and re-precipitated by the addition of 1 ml 100% alcohol. The tubes were placed on ice for 15 min and the DNA was pelleted by centrifugation at  $4,000 \times g$  for 15 min at  $4^{\circ}\text{C}$ . The DNA was resuspended in a total volume of 400  $\mu\text{l}$  TE and used to transform *S. pneumoniae*.

### 2.5.3 Preparation of pneumococcal chromosomal DNA

*S. pneumoniae* chromosomal DNA was prepared using the Wizard<sup>®</sup> Genomic DNA Purification Kit (Promega Corporation). The protocol was modified to optimise recovery of chromosomal DNA from pneumococci as follows.

The growth from a BA plate incubated at  $37^{\circ}\text{C}$  overnight was resuspended in 10 ml THY containing 1% (w/v) choline and grown at  $37^{\circ}\text{C}$  for 3-4 h. The culture was pelleted by centrifugation at  $4,000 \times g$  for 20 min at  $4^{\circ}\text{C}$ , resuspended in 200  $\mu\text{l}$  50 mM EDTA (pH 8.0) plus 0.1% (w/v) DOC and the suspension was divided between two microfuge tubes before incubation at  $37^{\circ}\text{C}$  for 10 min. After 10 min, the suspension was translucent, indicating that most of the bacteria had lysed, then 200  $\mu\text{l}$  50 mM EDTA (pH 8.0) and 300  $\mu\text{l}$  Nuclei Lysis Solution (from the kit) was added and gently mixed by pipette. The lysates were then incubated at  $80^{\circ}\text{C}$  for 10 min to ensure all the bacteria had lysed. The lysates were cooled to room temperature, 1.5  $\mu\text{l}$  of a 10 mg/ml RNase A solution was added, mixed by inversion and incubated at  $37^{\circ}\text{C}$  for 15-60 min. The lysates were then cooled to room temperature and 100  $\mu\text{l}$  Protein Precipitation Solution (from the kit) was added and

mixed by inversion. The lysates were placed on ice for 5 min and then centrifuged at 13,000 rpm for 10 min. The supernatant was transferred to a clean microfuge tube and 600  $\mu$ l isopropanol was added and mixed by inversion. The precipitated chromosomal DNA appeared as white fluffy strands and was pelleted by centrifugation at 13,000 rpm for 10 min. The supernatant was carefully removed and the pellet was washed with 600  $\mu$ l 70% (v/v) ethanol. The DNA was again pelleted by centrifugation at 13,000 rpm for 5 min and the supernatant was drained from the tube, removing as much as possible. The pellet was resuspended in 50-100  $\mu$ l TE and allowed to resuspend at 4°C overnight.

## **2.6 Analysis and manipulation of DNA**

### **2.6.1 Restriction endonuclease digestion of DNA**

DNA samples were digested using the restriction enzyme buffers recommended by the manufacturers. About 100 to 500 ng of plasmid DNA, or purified restriction fragments, were incubated with 2 units of each restriction enzyme in a final volume of 10  $\mu$ l at 37°C for 2 h, with the exception of *Sma*I digests, which were carried out at 30°C. The reactions were terminated by heating at 65°C for 10 min. Prior to electrophoresis on an agarose gel, a one-tenth volume of tracking dye consisting of 25% (w/v) Ficoll (type 400), 0.25% (w/v), bromophenol blue, 0.25% (w/v) xylene cyanol, was added. For the digestion of chromosomal DNA, 4  $\mu$ g of DNA was incubated with 4 units of each enzyme in a final volume of 20  $\mu$ l at 37°C for 16 h.

### **2.6.2 Calculation of sizes of restriction fragments**

The sizes of restriction enzyme fragments were calculated by comparing their

relative mobility with that of *EcoRI*-digested *B. subtilis* bacteriophage SPP1 DNA, and/or with that of *HindIII*-digested  $\lambda$  phage DNA obtained from Bresatec (Adelaide, South Australia). The calculated molecular sizes of the *EcoRI*-digested SPP1 DNA fragments were: 8.51 kb, 7.35 kb, 6.11 kb, 4.84 kb, 3.59 kb, 2.81 kb, 1.95 kb, 1.86 kb, 1.51 kb, 1.39 kb, 1.16 kb, 0.98 kb, 0.72 kb, 0.48 kb, and 0.36 kb (Ratcliffe et al., 1979). The calculated molecular sizes of the *HindIII*-digested phage  $\lambda$  DNA fragments were: 23.13 kb, 9.42 kb, 6.56 kb, 4.36 kb, 2.32 kb, 2.03 kb, 0.56 kb and 0.13 kb (Sanger et al., 1983).

### **2.6.3 Analytical and preparative separation of restriction fragments**

Electrophoresis of digested DNA was carried out at room temperature on horizontal, 0.8% -1.5% (w/v) agarose (DNA grade agarose, Progen Industries) gels with a Tris-borate-EDTA (TBE; 67 mM Tris, 22 mM boric acid and 2 mM EDTA at pH 8.0) buffer system, containing 0.5  $\mu$ g/ml ethidium bromide, as described by Maniatis et al. (1982). DNA bands were visualised by trans-illumination with UV light and photographed using either a Polaroid MP4 camera with type 667 positive film or Mitsubishi thermal paper (K65HM) on a Mitsubishi video copy processor fitted to a Tracktel GDS-2 camera.

### **2.6.4 Isolation of restriction fragments from agarose gels**

DNA restriction fragments were excised from agarose gels and purified as detailed in the Qiaex handbook using the Qiaex DNA Gel Extraction Kit (Qiagen).

### **2.6.5 *In vitro* cloning**

The DNA fragment to be subcloned (about 200 ng) was cleaved in either single or double restriction enzyme digests. This was combined with 20 ng of similarly cleaved

vector DNA, then ligated with 1 unit of T4 DNA ligase in a volume of 15  $\mu$ l in a buffer containing a final concentration of 20 mM Tris-HCl (pH 7.5), 10 mM MgCl<sub>2</sub>, 10 mM DTT, and 600 nM ATP at 4°C for 16 h. The ligated DNA was then used directly for transformation of *E. coli* strains. Wherever possible, transformants were screened for insertional inactivation of any interrupted genetic marker prior to plasmid DNA isolation. If the interrupted marker conferred antibiotic resistance, the transformants were screened for sensitivity to the antibiotic. Where the DNA insert was cloned into a multicloning site near the start of the  $\beta$ -galactosidase gene, transformants were selected using indicator plates containing IPTG and X-Gal. The transformants containing only vector DNA formed blue colonies whereas those with additional DNA inserted in the vector formed white colonies in these plates.

## 2.7 Southern hybridisation

### 2.7.1 Labelling of DNA probes

DNA probes were labelled with digoxigenin (DIG) using the DNA Labelling Kit purchased from Boehringer Mannheim. The DNA to be labelled was added to a microfuge tube and denatured by heating to 95°C for 10 min and then rapidly cooled on ice for 1 min. Then 2  $\mu$ l hexanucleotide mixture, 2  $\mu$ l dNTP mixture (containing DIG-11-dUTP) and 1  $\mu$ l Klenow enzyme was added to the DNA. Water was added to give a final volume of 20  $\mu$ l. The mix was incubated at 37°C overnight. The labelled DNA mixture was heated at 95°C for 10 min and cooled rapidly on ice before being added to the hybridisation mix.

## 2.7.2 DNA hybridisation by Southern blotting

### 2.7.2.1 Transfer of DNA to nylon membrane

Restriction fragments from single and/or double digests of DNA were separated electrophoretically on a 1% (w/v) agarose gel. DIG labelled lambda DNA, restricted with *HindIII*, was used as a DNA molecular size marker. Unidirectional transfer of DNA from agarose gels to Hybond-N<sup>+</sup> nylon membrane (Amersham, England) was performed as described by Southern (1975) and modified by Maniatis *et al.* (1982). The DNA was fixed to the nylon membrane by placing the membrane on a piece of Whatman filter paper soaked in 0.4 M NaOH for 5-30 min. The membrane was then washed in 5x SSC (150 mM NaCl, 15 mM Na<sub>3</sub> citrate, pH 7.0 for 1x SSC), and then in deionised water before prehybridisation.

### 2.7.2.2 High stringency hybridisation

The membrane was prehybridised for 1-6 h at 42°C in a prehybridisation solution consisting of 50% (v/v) formamide, 5x SSC (pH 7.0), 2% (w/v) blocking reagent (Boehringer Mannheim) and 0.01% (w/v) N-lauroylsarcosine. The membrane was then hybridised overnight at 42°C in hybridisation solution consisting of 50% (v/v) formamide, 5x SSC (pH 7.0), 5x Denhardt's reagent, 1% (w/v) SDS, 100 µg/ml heparin and 100 µg/ml of single stranded herring sperm DNA (Sigma) (Maniatis *et al.*, 1982) to which the denatured DIG labelled DNA probe had been added. After hybridisation, high stringency washing of the membrane was carried out as follows. The membrane was washed twice for 5 min with 2x SSC at 65°C, and then twice for 5 min with a solution of 0.2x SSC and 0.1% (w/v) SDS at 65°C. These conditions were able to detect DNA sequences with ≥90% homology to the DIG labelled DNA probe, based on the actual DNA sequences detected.

### **2.7.2.3 Low stringency hybridisation**

Low stringency hybridisation was essentially the same as high stringency hybridisation, except it was performed at room temperature. After hybridisation, the membrane was washed twice for 15 min with 2× SSC at 50°C. These conditions were able to detect DNA sequences with ≥70% homology to the DIG labelled DNA probe, based on the actual DNA sequences detected.

### **2.7.2.4 Immunological detection**

DIG labelled probe which annealed to the DNA on the membrane was detected as follows. After stringency washing of the membranes, they were washed briefly in DIG Buffer 1 (100 mM maleic acid, 150 mM NaCl, pH 7.5) and then incubated for 30 min at room temperature in 100 ml of DIG Buffer 2 (1% blocking reagent in DIG Buffer 1). The membrane was then incubated with 6 µl anti-DIG Fab fragment, alkaline phosphatase (AP) conjugate in 30 ml DIG Buffer 2 for 1 h. The membrane was then washed twice in 100 ml DIG Buffer 1 and then briefly in DIG Buffer 3 (100 mM Tris, 100 mM NaCl, 50 mM MgCl<sub>2</sub>, pH 9.5). The DIG label was detected by incubating the membrane in the dark with 10 ml DIG Buffer 3 containing 45 µl 4-nitro blue tetrazolium (NBT; 75 mg/ml in 70% (v/v) dimethylformamide) and 35 µl 5-bromo-4-chloro-3-indolylphosphate, toluidinium salt (X-phosphate; 50 mg/ml in 100% dimethylformamide).

## **2.8 DNA sequencing and analysis**

### **2.8.1 Nested deletions**

In order to obtain sequence of both strands of various clones containing

pneumococcal DNA, a series of nested deletion derivatives were constructed by the method of Henikoff (1984) using an 'Erase-a-base' kit (Promega, Madison, WI) according to the manufacturer's instructions.

Approximately 5 µg of DNA was digested, in a total volume of 20 µl, with two restriction enzymes. One which generated exonuclease I resistant 3' ends (such as *Bst*XI, *Kpn*I and *Sph*I) and one which generated exonuclease I sensitive ends. This ensured that exonuclease digestion would proceed in only one direction. An aliquot of 1 µl was analysed by agarose gel electrophoresis to ensure that digestion was complete and that the plasmid had been linearised.

Prior to digestion with exonuclease III (Exo III), the restricted DNA was diluted to 60 µl in Exo III 1x buffer (66 mM Tris-HCl (pH 8.0), 0.66 mM MgCl<sub>2</sub>) and incubated at 37°C. After the addition of 400 units of Exo III, 2.5 µl aliquots were removed at 30 s intervals and added to 7.5 µl of S1 nuclease mix (40.5 mM potassium acetate, pH 4.6, 338 mM NaCl, 1.35 mM ZnSO<sub>4</sub>, 6.75% (v/v) glycerol containing 2.25 units of S1 nuclease) and placed on ice. After all the samples had been taken, the tubes were incubated at room temperature for 30 min. After the addition of 1 µl of S1 stop buffer (300 mM Tris, 50 mM EDTA) the samples were incubated at 65°C for 10 min to inactivate the S1 nuclease. A 2 µl aliquot from each of the time points was analysed by agarose gel electrophoresis to determine the extent of Exo III digestion. The samples which had been appropriately digested were processed further. The samples were incubated at 37°C for 1 min prior to the addition of 1 µl of Klenow mix (30 µl of 20 mM Tris (pH 8.0), 100 mM MgCl<sub>2</sub> with 5 units of Klenow DNA polymerase). The samples were incubated at 37°C for 3 min before the addition of 1 µl of dNTP mix (0.125 mM of each of dATP, dCTP, dGTP and dTTP) followed by a further 5 min incubation at 37°C. The samples were then transferred to room

temperature and mixed with 40  $\mu$ l of ligase mix (50 mM Tris, pH 7.6, 10 mM MgCl<sub>2</sub>, 1 mM ATP, 5% (w/v) PEG 8000, 1 mM DTT, with 5 units/ml T4 DNA ligase). After 1 h the samples were placed at 4°C overnight.

The samples were then transformed (as described in section 2.4.2) into competent *E. coli* DH5 $\alpha$ , prepared using the RbCl<sub>2</sub> method (section 2.4.1.1), and plated onto selection plates. Plasmid DNA from the transformants was characterised by restriction enzyme analysis. Plasmids which contained appropriate deletions were subjected to DNA sequencing.

## 2.8.2 DNA Sequencing using dye-labelled primers

Sequencing reactions were carried out using 1  $\mu$ g of double stranded plasmid DNA with the Applied Biosystems Prism dye-primer cycle sequencing ready reaction kit as detailed below.

Reagent	A	C	G	T
Ready Reaction premix	4	4	8	8
DNA template ( $\mu$ l)	1	1	2	2
Total volume ( $\mu$ l)	5	5	10	10

Samples were placed in a Hybaid® Touchdown® Thermal Cycler. The program used consisted of: rapid thermal ramp to 95°C for 30 sec, rapid thermal ramp to 55°C for 30 sec, rapid thermal ramp to 70°C for 60 sec and was repeated for 15 cycles.

After the cycles were completed, the four separate reactions for each DNA sample were mixed together with 80  $\mu$ l of 95% (v/v) ethanol, mixed thoroughly, and kept on ice for 15 min to precipitate the DNA. The DNA was pelleted at 15,000 rpm for 30 min in a

microcentrifuge, washed in 250  $\mu$ l of 70% (v/v) ethanol, centrifuged for another 10 min, dried *in vacuo* for 3 min and stored at -20°C.

### 2.8.3 Sequencing with dye-labelled terminators

Double-stranded DNA template (1  $\mu$ g in a volume of up to 10  $\mu$ l), 8  $\mu$ l of Applied Biosystems Prism dye-terminator premix and 3.2 pM of primer were added and the volume adjusted to 20  $\mu$ l with deionised H<sub>2</sub>O. Sequencing reactions were carried out using 1  $\mu$ g of double-stranded plasmid DNA in a Hybaid® Touchdown® Thermal Cycler as follows; 96°C for 30 sec, 50°C for 30 sec, 60°C for 4 min repeated for 25 cycles.

After the cycles were completed, each reaction mixture was precipitated with 50  $\mu$ l of 95% (v/v) ethanol and 2  $\mu$ l of 3 M sodium acetate (pH 4.8) and placed on ice for 10 min. The DNA was pelleted at 15,000 rpm for 30 min in a microfuge at 4°C, and the pellet was washed in 250  $\mu$ l of ice-cold 70% (v/v) ethanol. Samples were dried *in vacuo* and stored at -20°C.

### 2.8.4 DNA sequencing

The sequencing reactions were then sent to the DNA sequencing service provided by the Molecular Pathology Unit, IMVS, Adelaide, Australia. DNA sequence data were obtained using an Applied Biosystems model 373A automated DNA sequencer.

### 2.8.5 DNA sequence analysis

The sequence was analysed using DNASIS and PROSIS Version 7.0 software (Hitachi Software Engineering, San Bruno, CA.). The programs BLASTX (Altschul *et al.*, 1990) and gapped BLAST (Altschul *et al.*, 1997) was used to translate DNA sequences and

conduct homology searches of the protein databases available at the National Center for Biotechnology Information, Bethesda, Md. Amino acid sequence alignments were performed with the program CLUSTAL (Higgins and Sharp, 1988) and CLUSTAL W (Thompson *et al.*, 1994). The program PROFILEGRAPH (Hofmann and Stöffel, 1989) was used to align hydropathy plots generated according to the method of Kyte and Doolittle (1982). Phylogenetic analysis was carried out using the program MEGA (Kumar *et al.*, 1994).

## **2.9 PCR Amplification**

### **2.9.1 PCR amplification using Taq polymerase**

PCR amplification was performed in a Hybaid® Touchdown® Thermal Cycler and the 50 µl reaction volume contained PCR buffer (1.5 mM MgCl<sub>2</sub>, 10 mM Tris (pH 8.4), 50 mM KCl, and 100 µg per ml of gelatine), 1.5 U of Taq polymerase, 1 µM of each primer, 100 ng of DNA template, and 200 mM of each of the four dNTPs. The program consisted of 25 cycles, and each cycle consisted of: denaturation at 95°C for 30 sec, annealing at 50°C for 30 sec, and elongation at 72°C for 1-4 min, 1 min for each kb of expected product. A 5 µl sample of the PCR product was analysed by agarose gel electrophoresis.

### **2.9.2 Long range PCR**

The Expand™ Long Template PCR System (Boehringer Mannheim, Germany) was used for long range PCR (LR-PCR), according to the manufacturer's instructions, in a Hybaid® Touchdown® Thermal Cycler. The LR-PCR program used consisted of 25 cycles,

after an initial incubation at 92°C for 2 min to allow complete denaturation of the DNA template. The first 10 cycles consisted of denaturation at 92°C for 30 sec, annealing initially at 55°C for 1 min, decrementing by 1°C per cycle to 45°C at cycle 10, and elongation at 68°C for 5-20 min. Cycles 11-25 consisted of denaturation at 92°C for 30 sec, annealing at 55°C for 1 min at cycle 11, decrementing to 45°C at cycle 25, and elongation at 68°C for 5-20 min at cycle 11, increasing by 1 min per cycle to cycle 25. The length of the extension time varied with the expected size of the PCR product, approximately 1 min for every kb to be amplified, up to a maximum of 20 min.

### **2.9.3 Inverse PCR**

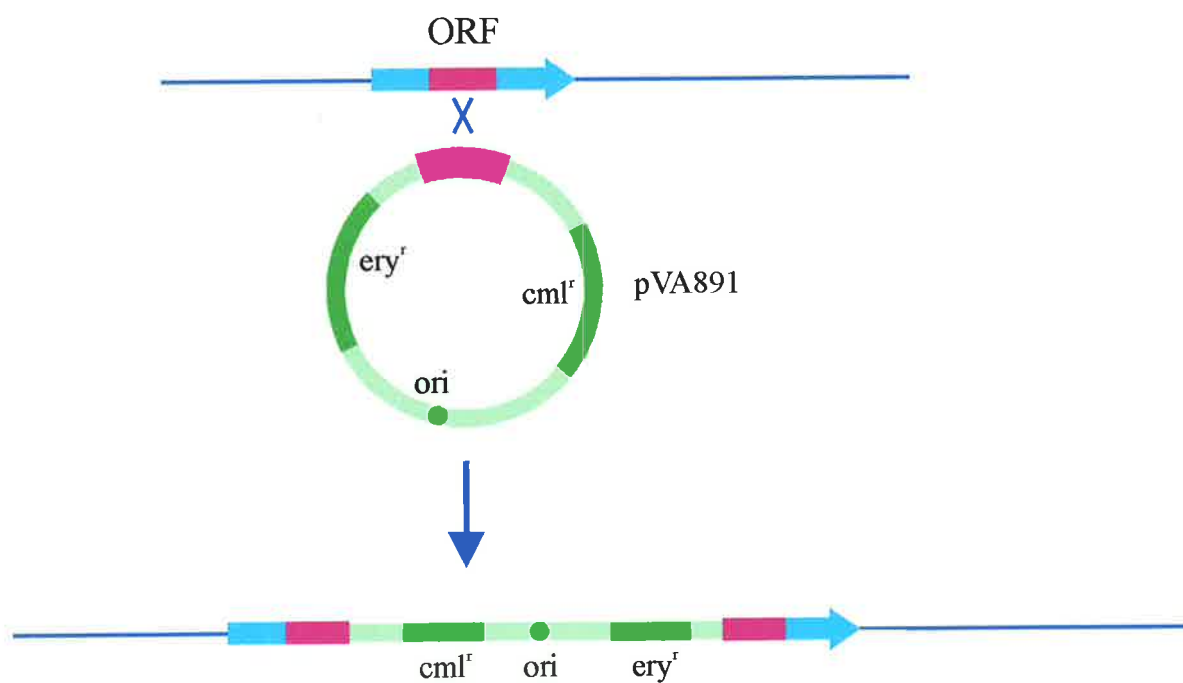
The method used for inverse PCR (InPCR) was that described by Ochman *et al.* (1988). Briefly, chromosomal DNA (5 µg) was digested to completion with an appropriate restriction enzyme. The cleaved DNA was then self-religated at 16°C for 16 h at a concentration of 10 µg/ml. Amplification was performed using 1 µg of ligated DNA and 50 pM of each appropriate primer using either Taq polymerase or LR-PCR as described above.

### **2.9.4 Purification of PCR products**

PCR products were purified for further analysis using a PCR Clean Up Kit (QIAGEN), according to the manufacturer's instructions. If the PCR product required concentration, several individual PCR reactions were pooled and purified and eluted into a final volume of 50 µl.

## 2.10 Insertion-duplication mutagenesis

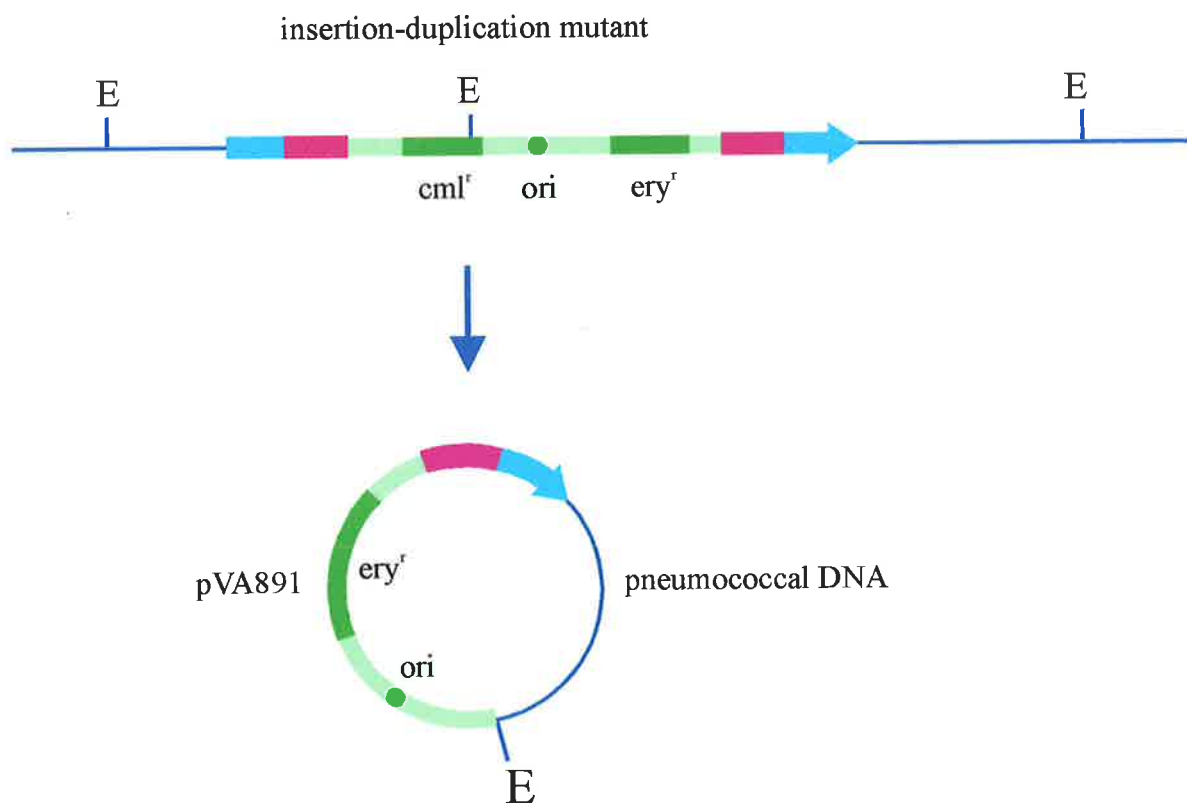
This method was used to mutagenise pneumococci by insertion of pVA891 (Table 2.3) into pneumococcal ORFs. This required the cloning an appropriate internal fragment of the pneumococcal ORF into pVA891, which encodes Cml and Ery resistance, but cannot replicate in *S. pneumoniae* (Macrina *et al.*, 1983), and transforming this plasmid construct into Rx1-19F as described in section 2.4.4. As the pVA891 replicon cannot function in pneumococci, Ery-resistant transformants are the result of a homologous recombination event directed by the cloned fragment of pneumococcal DNA that leads to the integration of the pVA891 plasmid into the host chromosome, and consequent disruption of the gene of interest. Furthermore, the cloned segment of pneumococcal DNA is duplicated and flanks the integrated copy of pVA891 as illustrated in Fig. 2.1.



**Fig. 2.1. Insertion-duplication mutagenesis.** Recombination between the pneumococcal chromosome and the homologous DNA fragment (shown in pink) inserted into pVA891 (shown in green) results in the insertion of pVA891 into the chromosome and interruption of the ORF. The relative positions of the origin of replication and the antibiotic resistance markers for Cml and Ery are shown in dark green.

## 2.11 Plasmid insertion/rescue

This procedure was used to isolate pneumococcal DNA which flanks the pVA891 insertion site of an insertion-duplication mutant. Plasmid sequences along with flanking DNA were recovered from the mutants by digestion of chromosomal DNA (5  $\mu$ g) with an appropriate restriction enzyme, self-religation of the cleaved DNA (at a concentration of 10  $\mu$ g/ml) as shown in Fig. 2.2, followed by transformation into *E. coli* DH5 $\alpha$  (section 2.4.2) and selection for either Cml or Ery resistance. The pneumococcal DNA was subsequently subcloned into pGEM-7Zf(+), pBluescript SK+, or pBluescript KS+ for further analysis.



**Fig. 2.2. Plasmid insertion/rescue.** Chromosomal DNA from a pneumococcal insertion-duplication mutant is digested with an appropriate restriction endonuclease (E) and self-ligated to recircularise the pVA891 replicon (shown in green) containing additional pneumococcal DNA from the adjoining region. The relative positions of the origin of replication and the antibiotic resistance markers for Cml and Ery are shown in dark green.

## 2.12 Protein and LPS analysis

### 2.12.1 T7 expression

DH5(pGP1-2) derivatives were grown in 10 ml LB+KA (LB broth containing Amp and Kan) at 30°C for 16 h. The bacteria were then washed in 10 ml LB+KA and resuspended in 10 ml LB+KA. One ml was diluted to 30 ml in LB+KA and grown at 30°C to an OD of 0.3. The culture was then divided in two. Half was heat shocked at 42°C for 30 min to induce T7 expression and then grown at 37°C for 2 h after the addition of 200 µg/ml Rif. The other 15 ml was uninduced and grown at 30°C for the duration of the experiment. The bacteria were pelleted by centrifugation at 4,000 × g for 15 min at 4°C and resuspended in 50 µl water and 350 µl loading buffer (consisting of 62.5 mM Tris-HCl (pH 6.8) 2% (w/v) SDS, 10% (v/v) glycerol, 4% (v/v) β-mercaptoethanol, and 0.1% (w/v) bromophenol blue), heated at 100°C for 5 min, and stored at -20°C before SDS-PAGE analysis.

### 2.12.2 LPS preparation

One ml of an overnight culture was pelleted at 15,000 rpm for 2 min in a microfuge, the pellet was resuspended in 50 µl loading buffer and heated at 100°C for 5 min. Then 10 µl of proteinase K solution (2.5 mg/ml proteinase K in loading buffer) was added and incubated at 56°C for 4 h. The samples were stored at -20°C before SDS-PAGE analysis.

### 2.12.3 SDS-PAGE

SDS-PAGE was performed as described by Laemmli (1970) using Protean II gel apparatus (Bio-Rad Laboratories) and a Model 2000/200 power supply (Bio-Rad Laboratories). The 10 or 15% separating polyacrylamide gels consisted of 10 or 15% acrylamide, 142 mM Tris (pH 8.8), 0.05% SDS, 0.04% ammonium persulphate and 0.014% TEMED. The stacking gel consisted of 3% acrylamide 62 mM Tris (pH 6.8), 0.05% SDS, 0.04% ammonium persulphate and 0.011% TEMED. The polyacrylamide gels were 15 cm long, 14 cm wide, and 1.5 mm thick. The prepared samples were loaded using a loading syringe (Hamilton Company, Reno, Nevada, USA) prior to electrophoresis in electrode buffer (25 mM Tris, 25 mM glycine, 0.1% SDS, pH 8.3) at 350 V for 2-3 h. Proteins were stained overnight at room temperature in 0.1% (w/v) Coomassie Brilliant Blue R250, 10% (v/v) acetic acid, 25% (v/v) methanol, with gentle agitation. The gel was destained by gentle agitation in several changes of 10% (v/v) acetic acid, 10% (v/v) isopropanol at 65°C. The molecular weight markers (Bio-Rad) were myosin (200 kDa),  $\beta$ -galactosidase (116.25 kDa), phosphorylase b (97.4 kDa), serum albumin (66.2 kDa), ovalbumin (45 kDa), carbonic anhydrase (31 kDa), trypsin inhibitor (21.5 kDa), lysozyme (14.4 kDa) and aprotinin (6.5 kDa).

### 2.12.4 Western blotting

Samples were electrophoretically transferred from SDS-PAGE gels onto nitrocellulose filters as described by Towbin *et al.* (1979), in transblot buffer (25 mM Tris-HCl, 200 mM glycine, 5% (v/v) methanol, pH 8.3) at 350 mA for 2 h in a BioRad Transblot apparatus fitted with a water-cooled coil.

The nitrocellulose filter was then agitated for 20 min in 100 ml of TTBS (20 mM Tris-HCl, 150 mM NaCl, .05% (v/v) Tween-20, pH 7.4) containing 5% (w/v) skim milk.

This solution was then discarded and replaced with 50 ml TTBS containing 0.02% (w/v) skim milk and 10  $\mu$ l of anti-*S. flexneri* type 4 serum (Murex, UK), the filter was then incubated at room temperature, with agitation, overnight. The filter was then washed for 10 min in three changes of 100 ml TTBS. The filter was then incubated at room temperature with agitation for 60 min in 50 ml TTBS containing 25  $\mu$ l goat anti-rabbit IgG-AP conjugate. The filter was then washed for 5 min in four changes of 100 ml TTBS and once in DIG Buffer 3. The filter was then developed as described in section 2.7.2.4 except that colour development occurred much faster and it was not performed in the dark.

## Chapter 3

# CHARACTERISATION OF THE COMPLETE TYPE 19F CAPSULE LOCUS

### 3.1 Introduction

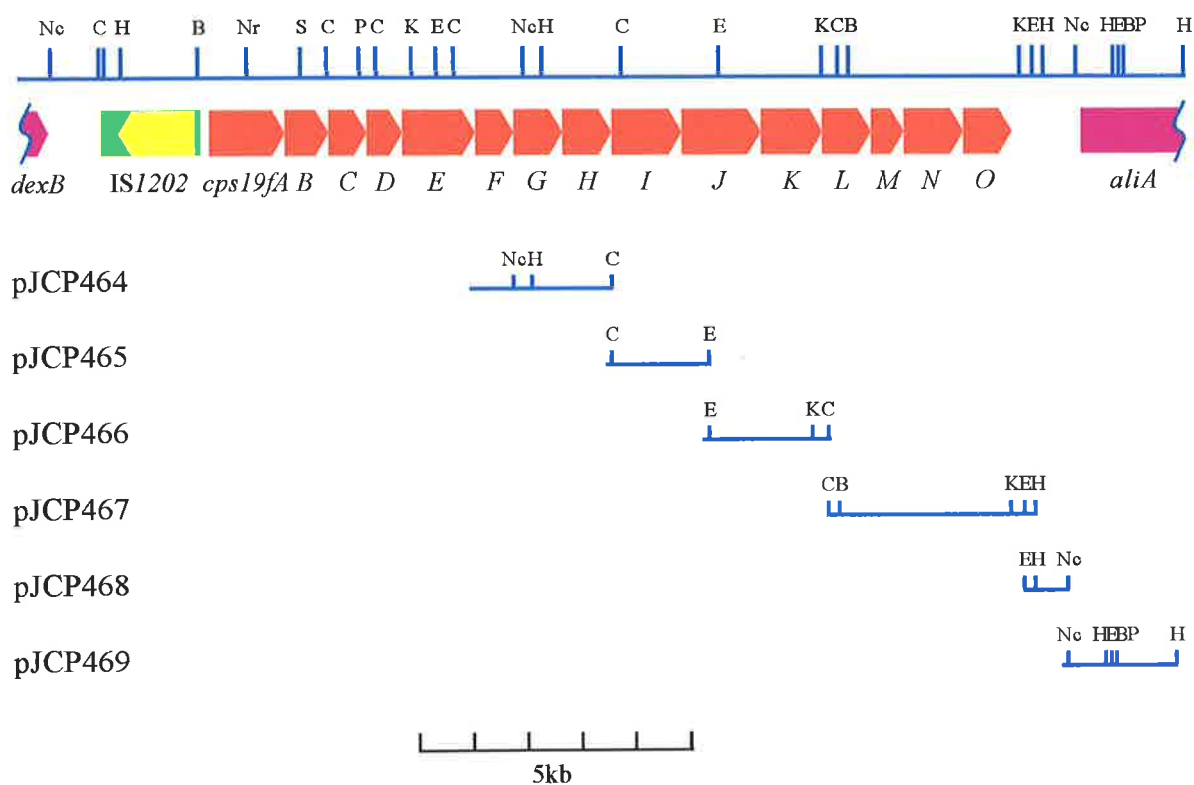
The cloning and sequencing of the first six genes of the type 19F *cps* locus has been reported previously (section 1.9.2; Guidolin *et al.*, 1994). At that time, no complete pneumococcal *cps* locus had been described for any serotype. Accordingly, the work described in this chapter was aimed at isolating, sequencing and where possible, functionally characterising the remainder of the pneumococcal type 19F *cps* locus.

### 3.2 Results

#### 3.2.1 Isolation, cloning and DNA sequencing of the type 19F *cps* locus

The strategy used to obtain the remainder of the type 19F locus was to sequentially isolate DNA fragments from *S. pneumoniae* Rx1-19F (Morona *et al.*, 1994a) using either plasmid insertion/rescue or InPCR, and then clone them into pBluescript SK+ or pGEM-7Zf(+) for further analysis. The resultant recombinant plasmids pJCP464, pJCP465,

pJCP466, pJCP467, pJCP468 and pJCP469 are shown in **Fig. 3.1** and described in detail below.



**Fig. 3.1. Physical map of the chromosome of *S. pneumoniae* Rx1-19F in the vicinity of the *cps19f* locus.** Boxed arrows represent potential ORFs and the closed box represents the insertion element *IS1202*. Gene designations are indicated below the map; *cps19fB-O* are abbreviated to *B-O*, respectively. Restriction sites are as follows; B, *Bam*HI; C, *Cla*I; E, *Eco*RI; H, *Hind*III; K, *Kpn*I; Ne, *Nco*I; Nr, *Nru*I; P, *Pst*I; S, *Sph*I. The regions of DNA subcloned into various recombinant plasmids are shown below the map.

The plasmid pJCP464 (**Fig. 3.1**) was constructed using plasmid insertion/rescue (described in section 2.11). The pVA891 replicon was excised, along with flanking pneumococcal DNA, from a derivative of *S. pneumoniae* Rx1-19F in which *cps19fF* had been interrupted by insertion of pVA891 (Guidolin *et al.*, 1994). The *Cla*I restricted DNA was recircularised and transformed into *E. coli* DH5 $\alpha$ . The 2.65 kb pneumococcal DNA insert was excised from the rescued pVA891 derivative and subcloned into pBluescript

SK+. This plasmid was used to generate a series of nested deletion derivatives (section 2.8.1) which were then sequenced (section 2.8.2).

The available sequence data was then used to design primers J7 and J8 (**Table 2.4**) which were used for InPCR amplification (section 2.9.3) of *EcoRI* restricted and circularised chromosomal DNA to obtain a 2.0 kb segment of downstream DNA, overlapping the insert of pJCP464 by 100 bp. This PCR product was cloned into pGEM-7Zf(+) to generate pJCP465. Sequencing of the insert of this plasmid enabled design of primers J9 and J10 (**Table 2.4**) which were then used for InPCR amplification of *ClaI* restricted and circularised chromosomal DNA to obtain a further 2.2 kb of flanking DNA (again overlapping by 100 bp). This PCR product was also cloned into pGEM-7Zf(+) to generate pJCP466 (**Fig. 3.1**) and sequenced.

The next plasmid, pJCP467 (**Fig. 3.1**), was obtained by rescue of the pVA891 replicon plus flanking DNA from a derivative of *S. pneumoniae* Rx1-19F in which *cpsI9fJ* had been interrupted (described in section 3.2.4). The *HindIII*-restricted DNA was recircularised and transformed into *E. coli* DH5 $\alpha$ , followed by subcloning of the 3.85-kb *ClaI-HindIII* pneumococcal DNA insert into pGEM-7Zf(+). Nested deletion derivatives of pJCP467 were sequenced. Sequence data across the *ClaI* site was obtained by dye terminator sequencing (section 2.8.3) of the rescued pVA891 based plasmid with the primer J12 (**Table 2.4**).

Two further rounds of InPCR were then used to construct pJCP468 and pJCP469. The PCR product obtained from *NcoI* restricted and circularised chromosomal DNA with primers J19 and J20 (**Table 2.4**) was cloned in pGEM-7Zf(+) to generate pJCP468 (**Fig. 3.1**) and then sequenced. Primers J26 and J27 (**Table 2.4**) were designed for inverse PCR amplification of *ClaI* restricted and circularised chromosomal DNA. The 4-kb PCR product was cloned into pBluescript SK+ to generate pJCP469 (**Fig. 3.1**). However only

the first 2.1 kb (to the *Hind*III site) was sequenced using the M13 forward primer and primers J32, J33 and J34 (Table 2.4) as well as sequence of a subclone containing the internal 1.2-kb *Hind*III fragment using the M13 forward and reverse primers. Only the part of pJCP469 which was sequenced is shown in Fig. 3.1.

### 3.2.2 Analysis of the *cps19f* DNA sequence

Both strands of the pneumococcal DNA inserts of each of the above plasmids (or nested derivatives thereof) were subjected to sequence analysis in order to compile the complete sequence of the remainder of *cps19f*, as shown in Appendix I. Examination of the compiled sequence revealed the presence of a further nine potential open reading frames (ORFs), which have been designated *cps19fG* to *cps19fO*. Each ORF is preceded by a ribosome binding site and the majority are very closely linked. The only potentially significant intergenic gaps occur between *cps19fJ* and *cps19fK* (63 nucleotides) and between *cps19fN* and *cps19fO* (65 nucleotides). However, potential stemmed-loop structures were not found immediately downstream of *cps19fJ* or *cps19fN* and no obvious promoter sequences were seen immediately upstream of *cps19fK* or *cps19fO*. The first six genes of the *cps19f* locus are also closely linked (intergenic distances 1-15 nucleotides) (Guidolin *et al.*, 1994) and so co-transcription of the entire locus remains a possibility.

A large region (1,458 nucleotides) downstream of *cps19fO* did not appear to contain any significant ORFs on either DNA strand, but the region from nucleotides 15,390-15,630 contained numerous stemmed-loop structures reminiscent of transcription terminators. There was, however, an additional ORF commencing at nucleotide 16,446, which was preceded by a ribosome binding site and by -10 and -35 promoter sequences. Comparison of sequence data for this gene with those deposited on GenBank indicated 97.3% DNA identity with that reported for the pneumococcal *aliA* gene, which encodes an

oligopeptide binding protein (Alloing *et al.*, 1994) and so is unrelated to CPS biosynthesis. This information, combined with earlier results (section 1.9) show that in *S. pneumoniae* type 19F, the *cps* locus is located between *dexB* and *aliA* in the chromosome, as shown in Fig. 3.1.

### 3.2.3 Characterisation of *cps19fG-O*

The locations and several properties of each of the ORFs designated *cps19fG-O* are summarised in Table 3.1. Significant similarities with other known proteins, revealed by comparison with sequence databases, are described below.

**Table 3.1. Summary of ORFs *cps19fG-O*.**

ORF	Location in sequence	Predicted MW	No. amino acids	Hydrophobicity index <sup>a</sup>	Predicted pI	%G+C content <sup>b</sup>
<i>cps19fG</i>	5,883-6,693	31,647	269	-0.39	8.43	36.3
<i>cps19fH</i>	6,694-7,572	34,474	292	-0.54	7.80	30.3
<i>cps19fI</i>	7,573-8,910	51,734	445	0.68	9.59	29.7
<i>cps19fJ</i>	8,933-10,354	55,055	473	0.81	9.83	29.7
<i>cps19fK</i>	10,418-11,506	40,950	362	-0.30	5.48	35.2
<i>cps19fL</i>	11,545-12,414	32,215	289	-0.21	4.69	42.3
<i>cps19fM</i>	12,415-13,011	22,379	198	-0.40	5.05	41.5
<i>cps19fN</i>	13,021-14,070	39,053	349	-0.45	5.16	42.1
<i>cps19fO</i>	14,136-14,986	32,330	283	-0.50	4.71	41.5

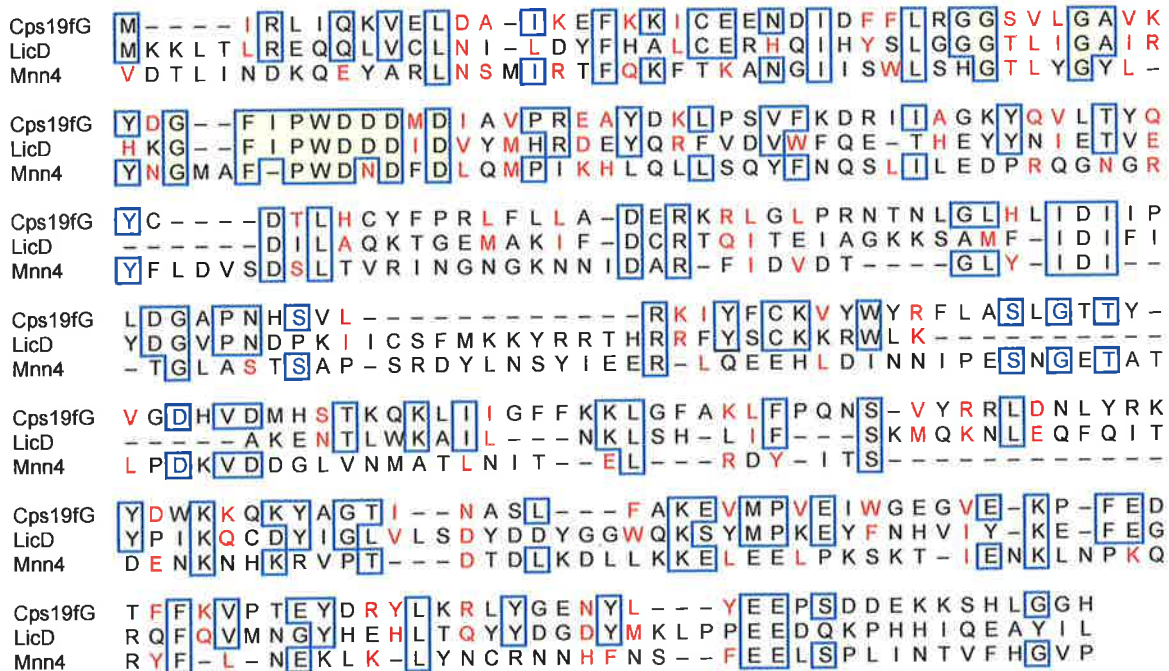
a. According to Kyte and Doolittle (1982), as implemented in PROSIS.

b. Percent guanine plus cytosine (G+C) of coding region.

#### 3.2.3.1 *cps19fG*

The *cps19fG* gene encodes a putative 31.6 kDa protein which exhibits 25.9% identity with the LicD protein of *H. influenzae*, which is encoded by the *licD* gene of its LPS locus (Weiser *et al.*, 1989). However, the precise function of LicD is unknown.

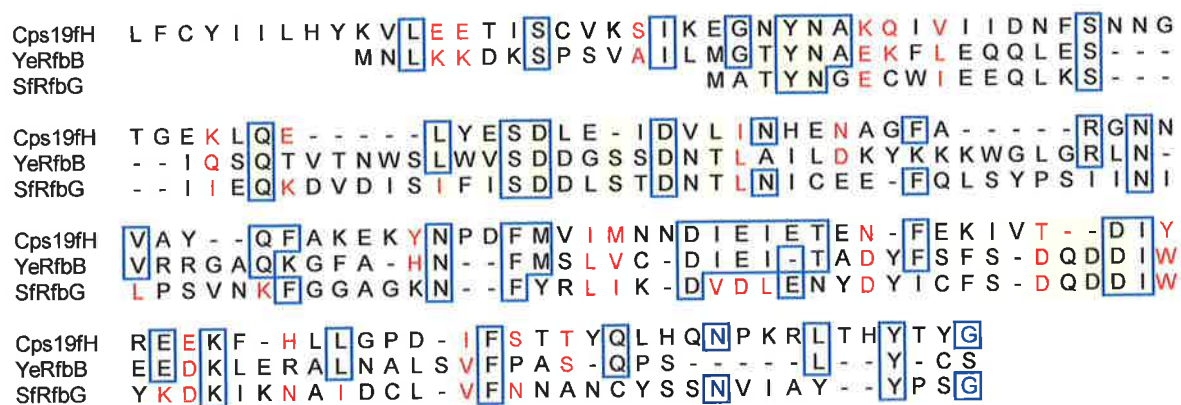
Cps19fG also has a low level of similarity to the central portion of the Mnn4 gene product which is required for the addition of mannosylphosphate to *N*-linked oligosaccharides in *S. cerevisiae* (Odani *et al.*, 1996). The alignment between these three proteins is shown in **Fig. 3.2**; the similarity is greatest at the N-terminal end of Cps19fG where a conserved motif has been identified. This motif includes aspartate (D) residues, which are predicted to be present in the catalytic sites of glycosyl transferases (Saxena *et al.*, 1995). The structure of the *H. influenzae* LPS is highly variable, but Phillips *et al.* (1993) suggest it may be phosphorylated. The type 19F CPS and the cell wall mannans in *S. cerevisiae* are both phosphorylated. Thus it is possible that Cps19fG, LicD and Mnn4 are all involved in the phosphorylation of polysaccharides.



**Fig. 3.2. Alignment of Cps19fG.** Alignment of Cps19fG with *H. influenzae* LicD (LicD) (Weiser *et al.*, 1989) and amino acids 465-730 of *S. cerevisiae* Mnn4 (Mnn4) (Odani *et al.*, 1996) using the default settings of the program CLUSTAL (Higgins and Sharp, 1988) and enhanced by manual adjustment. Residues identical to those in Cps19fG are boxed; similar residues are shown in red; - indicates absence of a residue. The yellow shaded regions correspond to an apparently conserved motif in these proteins.

### 3.2.3.2 *cps19fH*

The *cps19fH* gene encodes a putative 34.5 kDa protein, assuming that translation is initiated at the TTG codon, as shown in **Appendix I**. The nearest in-frame ATG codon is a further 90 codons downstream (the expected translation product would be approximately 20 kDa) and is not preceded by a ribosome binding site. Cps19fH has a limited degree of identity with rhamnosyl transferases from *Yersinia enterocolitica* and *S. flexneri*, as shown in the alignment of these proteins in **Fig. 3.3**. Interestingly, the more conserved regions correspond to a motif previously identified in a number of rhamnosyl and other 6-deoxyhexosyl transferases (Morona *et al.*, 1995). This suggests that Cps19fH is likely to be the rhamnosyl transferase involved in incorporation of rhamnose (Rha) into the type 19F CPS.



**Fig. 3.3. Alignment of Cps19fH.** Alignment of Cps19fH, from amino acid position 1 to 143, with *Y. enterocolitica* RfbB (YeRfbB) (Zhang *et al.*, 1993) and *S. flexneri* RfbG (SfRfbG) (Morona *et al.*, 1995) using the default settings of the program CLUSTAL (Higgins and Sharp, 1988) and enhanced by manual adjustment. Residues identical to those in Cps19fH are boxed; similar residues are shown in red; - indicates absence of a residue. The yellow shaded regions correspond to those found to be most conserved amongst a variety of Gram-negative rhamnosyl and 6-deoxy-hexosyl transferases (Morona *et al.*, 1995).

### 3.2.3.3 *cps19fl*

The *cps19fl* gene encodes a putative 51.7 kDa protein, which has similarity to Rfc proteins (O-antigen polymerases) from a variety of Gram-negative bacteria, as shown in **Table 3.2**. Although the overall similarity between Cps19fl and the various Rfc proteins is low (14.0 - 20.1% identity) it is as strong as the degree of identity within the Gram-negative species (15.5-19.8%). The relationship between Cps19fl and Rfc proteins is even more apparent when similar as well as identical amino acids are considered, as exemplified by the alignment with *E. coli* K12 Rfc (Stevenson *et al.*, 1994) shown in **Fig. 3.4**. The hydropathy plots for the various proteins (**Fig. 3.5**) also illustrate the marked similarity, each having 10-12 hydrophobic, potentially membrane-spanning, domains. It therefore seems probable that Cps19fl is the polysaccharide polymerase.

**Table 3.2. Similarity of Cps19fl to other proteins.**

	% Identity <sup>a</sup>						
	Cps19fl	K12Rfc	SdRfc	StRfc	SfRfc	M67Rfc	M40Rfc
Cps19fl <sup>b</sup>	100	16.5 [394]	19.6 [382]	15.6 [417]	14 [350]	18.3 [394]	20.1 [318]
K12Rfc <sup>c</sup>		100	17.5 [382]	19.1 [397]	19.4 [309]	16.5 [358]	17 [335]
SdRfc <sup>d</sup>			100	15.5 [381]	16.5 [315]	18.5 [211]	16.5 [363]
StRfc <sup>e</sup>				100	18.8 [308]	19.7 [269]	17.5 [341]
SfRfc <sup>f</sup>					100	17.5 [331]	19.8 [349]
M67Rfc <sup>g</sup>						100	17.5 [360]
M40Rfc <sup>h</sup>							100

a. Percentage of identical amino acids determined with FASTA as implemented in PROSIS. Numbers in parentheses indicate the number of amino acids over which the % identity occurs. b. *S. pneumoniae* Cps19fl. c. *E. coli* K12 Rfc (Stevenson *et al.*, 1994). d. *Shigella dysenteriae* type 1 Rfc (Klena and Schnaitman, 1993). e. *S. enterica* serovar typhimurium Rfc (Collins and Hackett, 1991). f. *S. flexneri* Rfc (Morona *et al.*, 1994b). g. *S. enterica* serovar muenchen strain M67 Rfc (Brown *et al.*, 1992). h. *S. enterica* serovar montevideo strain M40 Rfc (Lee *et al.*, 1992).

```

Cps19fl  M S Y L F L L C L T L F L L T I F Y F F A F I Q D L I A P P V V M S V M F L I S S V
K12Rfc   M I Y L - V I - - S V F L I T A F I C L Y L K K D I F Y P A V C V N I I F A L V L L

Cps19fl  F A L V N S K N W N I E Y S G I A Y I L I I S G I I F S I P L M A L - K S P N F N
K12Rfc   G Y E I T S D I Y A F Q L N D A T L I F L L C N V L T F T L S C L - L T E S V L D L

Cps19fl  T E V K I A D R L I D I Q F W K I A - L T I I I D L F - I L Y L Y R K E I Y N L V L
K12Rfc   N I R K V N N A I Y S I P S K K V H N V G L L V I S F S M I Y I C M R - - - - - L

Cps19fl  S N G Y T G S N I Q W F F R N A T S Y E G E L T V R T F I R V L I R V I D V S A Y I
K12Rfc   S N Y Q F G T S L L S Y M N L I R D A D V E D T S R N F S A Y M Q P I I - L T T - -

Cps19fl  F G Y T F I N N F L I Y R H K R P K D I L L L V P L L I F I S K T L I S G G R Q D I
K12Rfc   F A L - F I W S K K F T N T K V S K T F T L L V - F I V F I F A I I L N T G K Q I V

Cps19fl  I K I L I A Y V I M M Y I Q Q K R K V G W N R V I S H K Y I H L G F V G L I A G I P
K12Rfc   F M V I I S Y - - - A F I - - - - - V G V N R V - - - K H - - - - Y V Y L I T A V G

Cps19fl  A F Y - - Y S L F L A G R S T T R T L F E S V S T Y L G G S I Q H F N Q - Y I E N P
K12Rfc   V L F S L Y M L F L R G L P G G M A Y Y - - L S M Y L V S P I I A F Q E F Y F Q Q V

Cps19fl  L D - - P G E V F - G S E T L V P I L N I L G E M G L V N Y R S T I H L E F R T L G
K12Rfc   S N S A S S H V F W F F E R L M G L L T G G V S M S L - - - - - H K E F V W V G

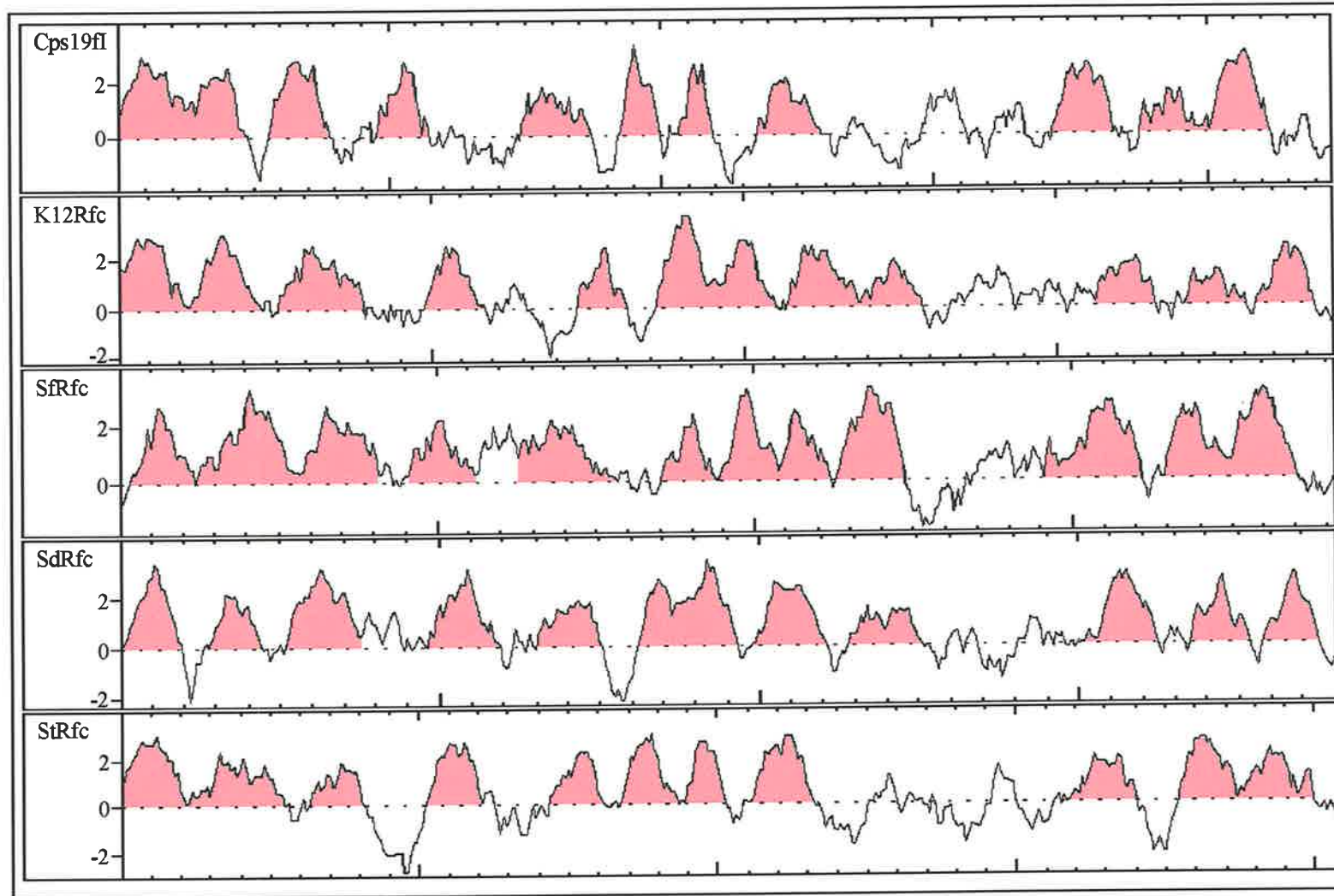
Cps19fl  V T V G N V Y T F F R R P L H D F G L V G M Y V F V F A V G A F F A I Y Y L V L R K
K12Rfc   L P T - N V Y T A F S D - - - - - - - - - Y V Y I S A - - - - - E L S Y L M M V I

Cps19fl  K Q V G F N L D I H T I I Y S Y V F Y W I F L S S I E Q Y S F T M I S L Y T L V F I
K12Rfc   - H G C I S G V L W R L S R N Y I S V K I F - - - - - Y S Y - F I Y T F S F I F Y

Cps19fl  V L V Y F M A - I F Y W C T D F K R G K L I F K I S D S S I K L K E E
K12Rfc   H E S - F M T N I S S W - I Q I T L C I V F S Q F L K A Q K I K - -

```

**Fig. 3.4. Alignment of Cps19fl.** Alignment of Cps19fl with *E. coli* K12 Rfc (K12Rfc) (Stevenson *et al.*, 1994), using the default settings of the program CLUSTAL (Higgins and Sharp, 1988). Identical residues are boxed; similar residues are shown in red; - indicates absence of a residue.



**Fig. 3.5. Hydropobicity analysis of Cps19fI.** Hydropathy plots of Cps19fI, *E. coli* K12 Rfc (K12Rfc) (Stevenson *et al.*, 1994), *S. flexneri* Rfc (SfRfc) (Morona *et al.*, 1994b), *S. dysenteriae* Rfc (SdRfc) (Klena and Schnaitman, 1993) and *S. enterica* serovar typhimurium Rfc (StRfc) (Collins and Hackett, 1991) were generated by the method of Kyte and Doolittle (1982) and aligned using PROFILEGRAPH (Hofmann and Stöffel, 1989). Positive numbers on the Y-axis indicate hydrophobic regions and putative membrane-spanning domains are shaded. The position of every 10th amino acid is marked on each X-axis.

### 3.2.3.4 *cps19fJ*

The *cps19fJ* gene encodes a putative 55.1 kDa protein with similarity to RfbX proteins of *E. coli*, *Shigella* sp. and *Yersinia* sp., as well as to the CapF protein of *S. aureus*, as shown in **Table 3.3**. The RfbX proteins are known to be involved in transport of O-antigen repeat units across the membrane (Liu *et al.*, 1996; Macpherson *et al.*, 1995). Again, the overall similarity between Cps19fJ and the various Gram-negative proteins is low (15.7-19.4% identity), but the degree of identity within the Gram-negative species is similar (16.4-31.4%). Moreover, the hydropathy plots for the Cps19fJ and the various RfbX-related proteins are very similar (**Fig. 3.6**). These are all integral membrane proteins with 10-12 hydrophobic, potentially membrane-spanning, domains. Thus, Cps19fJ is probably the polysaccharide repeat unit transporter.

**Table 3.3. Similarity of Cps19fJ to other proteins.**

	% Identity <sup>a</sup>						
	Cps19fJ <sup>b</sup>	SfRfbX <sup>c</sup>	K12RfbX <sup>d</sup>	SdRfbX <sup>e</sup>	YeTrgA <sup>f</sup>	YpRfbX <sup>g</sup>	SaCapF <sup>h</sup>
Cps19fJ	100	19.2	18.1	19.4	19.2	15.7	16.4
		[421]	[414]	[402]	[416]	[396]	[397]
SfRfbX		100	19.7	22.5	19.8	19.9	19
			[396]	[409]	[424]	[403]	[410]
K12RfbX			100	31.4	28.3	17.4	21.5
				[401]	[406]	[373]	[395]
SdRfbX				100	28.4	18.6	21.1
					[401]	[269]	[393]
YeTrsA					100	20.8	20.4
						[394]	[401]
YpRfbX						100	16.1
							[367]
SaCapF							100

a. Percentage of identical amino acids determined with FASTA as implemented in PROSIS. Numbers in parentheses indicate the number of amino acids over which the % identity occurs.

b. *S. pneumoniae* Cps19fJ.

c. *S. flexneri* RfbX (Macpherson *et al.*, 1995).

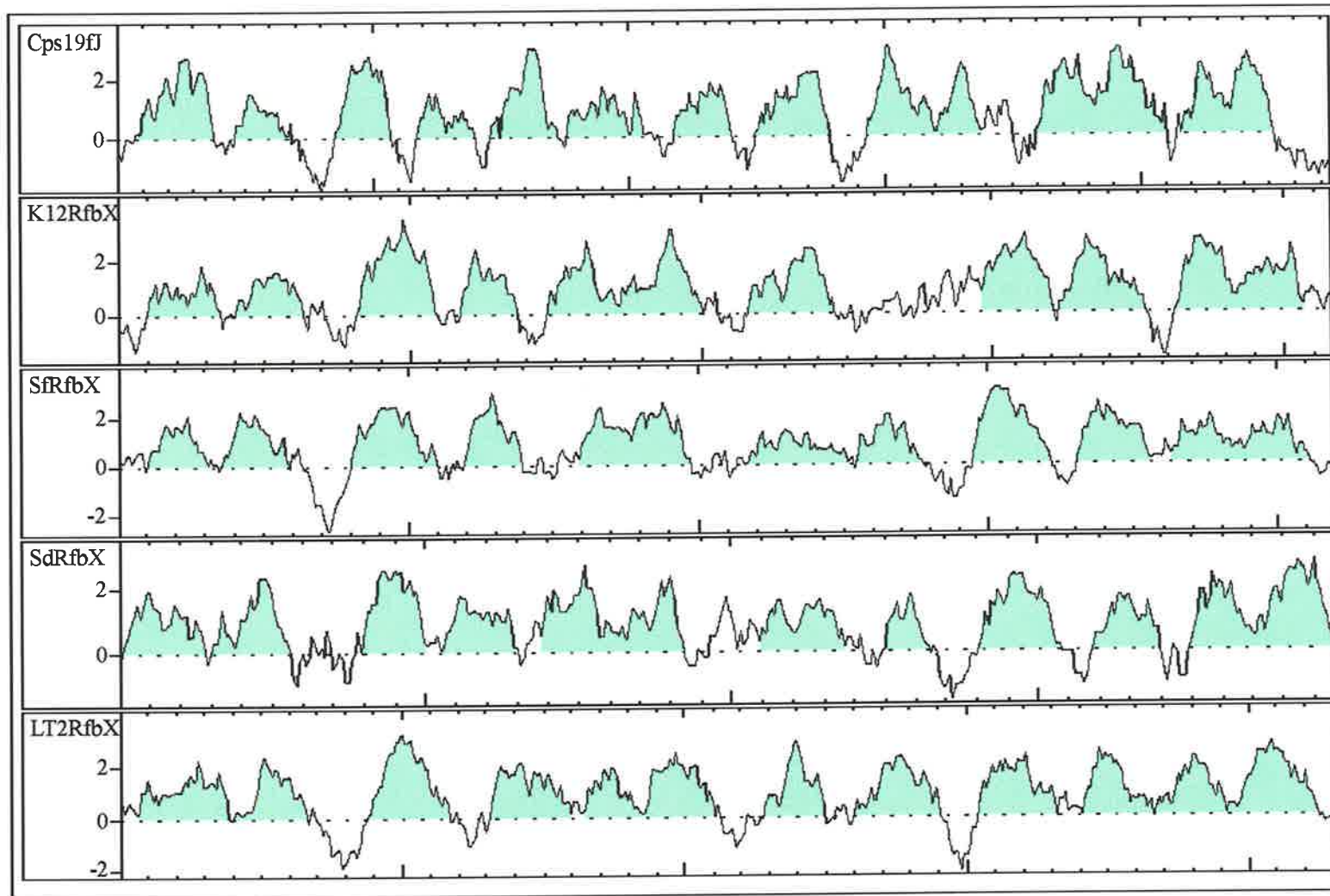
d. *E. coli* K12 RfbX (Stevenson *et al.*, 1994).

e. *S. dysenteriae* RfbX (Klena and Schnaitman, 1993).

f. *Y. enterocolitica* TrsA (Skurnik *et al.*, 1995).

g. *Y. pseudotuberculosis* RfbX (Kessler *et al.*, 1993).

h. *S. aureus* CapF (Lin *et al.*, 1994).



**Fig. 3.6. Hydropobicity analysis of Cps19fJ.** Hydropathy plots of Cps19fJ, *E. coli* K12 RfbX (K12RfbX) (Stevenson *et al*, 1994), *S. flexneri* RfbX (SfrfbX) (Macpherson *et al*, 1995), *S. dysenteriae* RfbX (SdrfbX) (Klena and Schnaitman, 1993) and *S. enterica* serovar typhimurium LT2 RfbX (LT2RfbX) (Jiang *et al*, 1991) were generated by the method of Kyte and Doolittle (1982) and aligned using PROFILEGRAPH (Hofmann and Stöffel, 1989). Positive numbers on the Y-axis indicate hydrophobic regions and putative membrane-spanning domains are shaded. The position of every 10th amino acid is marked on each X-axis.

### 3.2.3.5 *cps19fK*

The *cps19fK* gene encodes a putative 40.9 kDa protein. Cps19fK has a high degree of similarity (49% to 63% identity,) with a family of proteins including the *E. coli rffE* gene product (Table 3.4). The alignment of Cps19fK with these proteins is shown in Fig. 3.7. RffE is a UDP-GlcNAc-2-epimerase, and functions in the synthesis of UDP-ManNAc, a precursor required for the synthesis of UDP-N-acetyl mannosaminuronic acid in ECA biosynthesis (Meier-Dieter *et al.*, 1990). ManNAc is a component of the type 19F CPS and thus *cps19fK* is presumed to encode the enzyme needed to synthesise UDP-ManNAc in type 19F pneumococci.

**Table 3.4. Similarity of Cps19fK to other proteins.**

	%Identity <sup>a</sup>					
	Cps19fK <sup>b</sup>	BsYvyH <sup>c</sup>	K12RffE <sup>d</sup>	PsEpsC <sup>e</sup>	054RfbC <sup>f</sup>	Syn0624 <sup>g</sup>
Cps19fK	100	63.2	50.8	50.3	49.7	49
		[359]	[368]	[372]	[374]	[357]
BsYvyH		100	54.9	52	52.9	47
			[368]	[373]	[367]	[362]
K12RffE			100	65	54.8	47.3
				[369]	[367]	[364]
PsEpsC				100	52.6	45.8
					[367]	[369]
054RfbC					100	44.4
						[365]
Syn0624						100

a. Percentage of identical amino acids determined with FASTA as implemented in PROSIS. Numbers in parentheses indicate the number of amino acids over which the % identity occurs.

b. *S. pneumoniae* Cps19fK

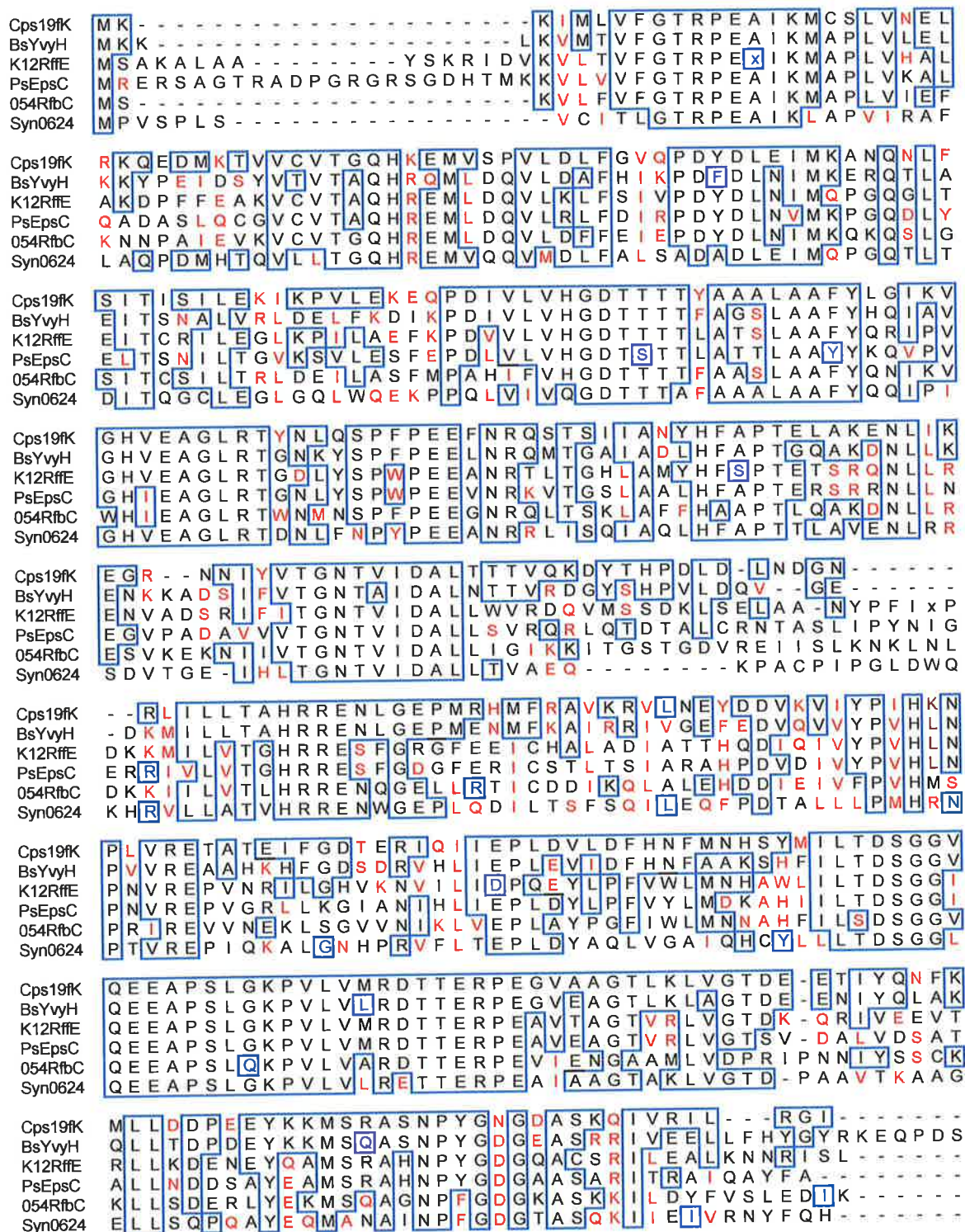
c. *B. subtilis* YvyH (Soldo *et al.*, 1993)

d. *E. coli* K12 RffE (Daniels *et al.*, 1992)

e. *P. solanacearum* EpsC (Huang and Schell, 1995)

f. *S. enterica* O54 RfbC (Keenleyside and Whitfield, 1995)

g. *Synechocystis* sp. ORF 0624 (Kaneko *et al.*, 1995)



**Fig. 3.7. Alignment of Cps19fK.** Alignment of Cps19fK with *B. subtilis* YvyH (BsYvyH) (Soldo *et al.*, 1993), *E. coli* K12 RffE (K12RffE) (Daniels *et al.*, 1992), *P. solanacearum* EpsC (PsEpsC) (Huang and Schell, 1995), *S. enterica* O54 RfbC (O54RfbC) (Keenleyside and Whitfield, 1995), and *Synechocystis* sp. ORF 0624 (Syn0624) (Kaneko *et al.*, 1995), as determined using the default settings of the program CLUSTAL (Higgins and Sharp, 1988). Residues identical to those in Cps19fK are boxed; similar residues are shown in red; - indicates absence of a residue.

### 3.2.3.6 *cps19fLMNO*

The *cps19fL*, *cps19fM*, *cps19fN* and *cps19fO* genes encode proteins of 32.2, 22.4, 39.1 and 32.3 kDa, respectively. These four genes have extensive similarity to a portion of the *S. flexneri rfb* gene cluster (*rfbBDAC*) which encodes the enzymes involved with dTDP-Rha biosynthesis (Macpherson *et al.*, 1994).

Cps19fL has a high degree of similarity (69.1-70.6% identity) to Glc-1-phosphate thymidyl transferases (RfbA) from *E. coli*, *S. flexneri* and *S. enterica* (Table 3.5). RfbA catalyses the first step in the dTDP-Rha biosynthesis pathway, suggesting *cps19fL* may also encode this enzyme. The alignment between Cps19fL and the RfbA proteins is shown in Fig. 3.8. Interestingly, there is also significant nucleotide sequence similarity (67-70% identity) between the *cps19fL* and *rfbA* genes.

Cps19fM has 24.3% to 34.6% identity to RfbC from various Gram-negative bacteria and StrM from *Streptomyces griseus* (Table 3.6, Fig. 3.9). Thus, *cps19fM* probably encodes dTDP-4-keto-6-deoxyglucose-3,5-epimerase, the third enzyme in the dTDP-Rha biosynthesis pathway.

Cps19fN has similarity to a large family of dehydratases, in particular to RfbB from various Gram-negative bacteria and StrE from *S. griseus* (Table 3.7, Fig. 3.10). An important feature of these closely related proteins is the presence of a NAD-binding domain (shaded residues in Fig. 3.10; Macpherson *et al.*, 1994). Thus, Cps19fN is probably a dTDP-Glc-4,6-dehydratase, the second enzyme in the dTDP-Rha biosynthesis pathway.

Cps19fO exhibits 27% to 39.7% identity to RfbD from various Gram-negative bacteria and StrL from *S. griseus* (Table 3.8, Fig. 3.11). These proteins also have a NAD-binding domain (shaded residues in Fig. 3.11). Thus Cps19fO is probably dTDP-L-Rha synthase, which catalyses the last step in dTDP-Rha biosynthesis.

Table 3.5. Similarity of Cps19fL to other proteins.

	% Identity <sup>a</sup>															
	Sf RfbA	K12 RfbA	O7 RfbA	M32 RfbA	LT2 RfbA	Nm RfbA	Ng RfbA	Ye RfbA	Xc RfbA	Ec RfbA	Bs SpsI	Stf TylA	Sg StrD	Sp DnrL	Sn SnoD	Sv GraD
Cps19fL <sup>b</sup>	68.8 [285]	69.4 [288]	70.6 [286]	70.5 [285]	69.1 [285]	66.9 [285]	67.6 [287]	66.1 [289]	64.1 [284]	64.3 [286]	45.3 [243]	54.4 [285]	35.4 [237]	33.6 [259]	33.3 [249]	36.2 [232]
SfRfbA <sup>c</sup>	100	89.3 [291]	93.2 [293]	91.1 [292]	89.3 [291]	63.2 [285]	63.5 [285]	64 [286]	65.9 [290]	62.4 [287]	43.1 [211]	50.7 [286]	33.9 [242]	37.9 [198]	33.5 [218]	34.9 [215]
K12RfbA <sup>d</sup>		100	94.2 [292]	92.4 [291]	92.1 [291]	65.1 [284]	65.1 [284]	65.6 [288]	68.3 [290]	64 [286]	43 [237]	51.7 [288]	34.6 [243]	31.4 [280]	33.7 [243]	36 [228]
O7RfbA <sup>e</sup>			100	95.6 [293]	93.5 [292]	65 [286]	65 [286]	65.9 [287]	67.7 [291]	63.9 [288]	42.9 [238]	52.4 [288]	35.7 [244]	32.8 [268]	33.6 [244]	35.8 [229]
M32RfbA <sup>f</sup>				100	95.2 [291]	66 [285]	66 [285]	64.7 [286]	67.2 [290]	63.1 [287]	43 [237]	51.6 [287]	35.4 [243]	33 [267]	32 [266]	36 [228]
LT2RfbA <sup>g</sup>					100	64.2 [285]	64.2 [285]	64.3 [286]	67.2 [290]	61.7 [287]	43.5 [237]	50.9 [287]	35 [243]	32.7 [266]	31.5 [267]	36 [228]
NmRfbA <sup>h</sup>						100	91.3 [288]	69.1 [285]	63.4 [284]	69.9 [286]	44.4 [239]	50.5 [285]	35.9 [237]	35.3 [232]	34.9 [232]	34.1 [232]
NgRfbA <sup>i</sup>							100	69.5 [285]	63.4 [284]	70.6 [286]	42.7 [239]	48.8 [285]	35.9 [237]	36.2 [232]	34.5 [232]	33.6 [232]
YeRfbA <sup>j</sup>								100	61.2 [286]	72.1 [287]	42.6 [244]	50 [286]	35.2 [244]	35.5 [234]	32.6 [233]	33.2 [232]
XcRfbA <sup>k</sup>									100	61.9 [286]	43.7 [245]	50.5 [289]	38.1 [247]	37.9 [232]	35.1 [248]	33.3 [243]
EcRffM <sup>l</sup>										100	42.2 [244]	51.2 [285]	35.4 [243]	36.1 [233]	32.1 [243]	33.3 [228]
BsSpsI <sup>m</sup>											100	39.5 [248]	28.9 [239]	30.8 [247]	32.5 [240]	31.3 [243]
StfTylA <sup>n</sup>												100	37.7 [236]	37.7 [236]	28.7 [268]	34.2 [231]
SgStrD <sup>o</sup>													100	59.2 [355]	54.1 [355]	51 [355]
SpDnrL <sup>p</sup>														100	49 [355]	49.2 [355]
SnSnoD <sup>q</sup>															100	57.2 [355]
SvGraD <sup>r</sup>																100

a. Percentage of identical amino acids determined with FASTA as implemented in PROSIS. Numbers in parentheses indicate the number of amino acids over which the % identity occurs.

b. *S. pneumoniae* Cps19fL

c. *S. flexneri* RfbA (Macpherson *et al.*, 1994)

d. *E. coli* K12 RfbA (Stevenson *et al.*, 1994)

e. *E. coli* O7 RfbA (Marolda and Valvano, 1995)

f. *S. enterica* E1 RfbA (Wang *et al.*, 1992)

g. *S. enterica* serovar typhimurium strain LT2 RfbA (Jiang *et al.*, 1991),

h. *N. meningitidis* RfbA (Hammerschmidt *et al.*, 1994)

i. *Neisseria gonorrhoeae* RfbA (Robertson *et al.*, 1994)

j. *Y. enterocolitica* RfbA (Zhang *et al.*, 1993)

k. *X. campestris* RfbA (Koplin *et al.*, 1993)

l. *E. coli* RffM (Daniels *et al.*, 1992)

m. *B. subtilis* SpsI P39629 (GenBank accession no. Z99123)

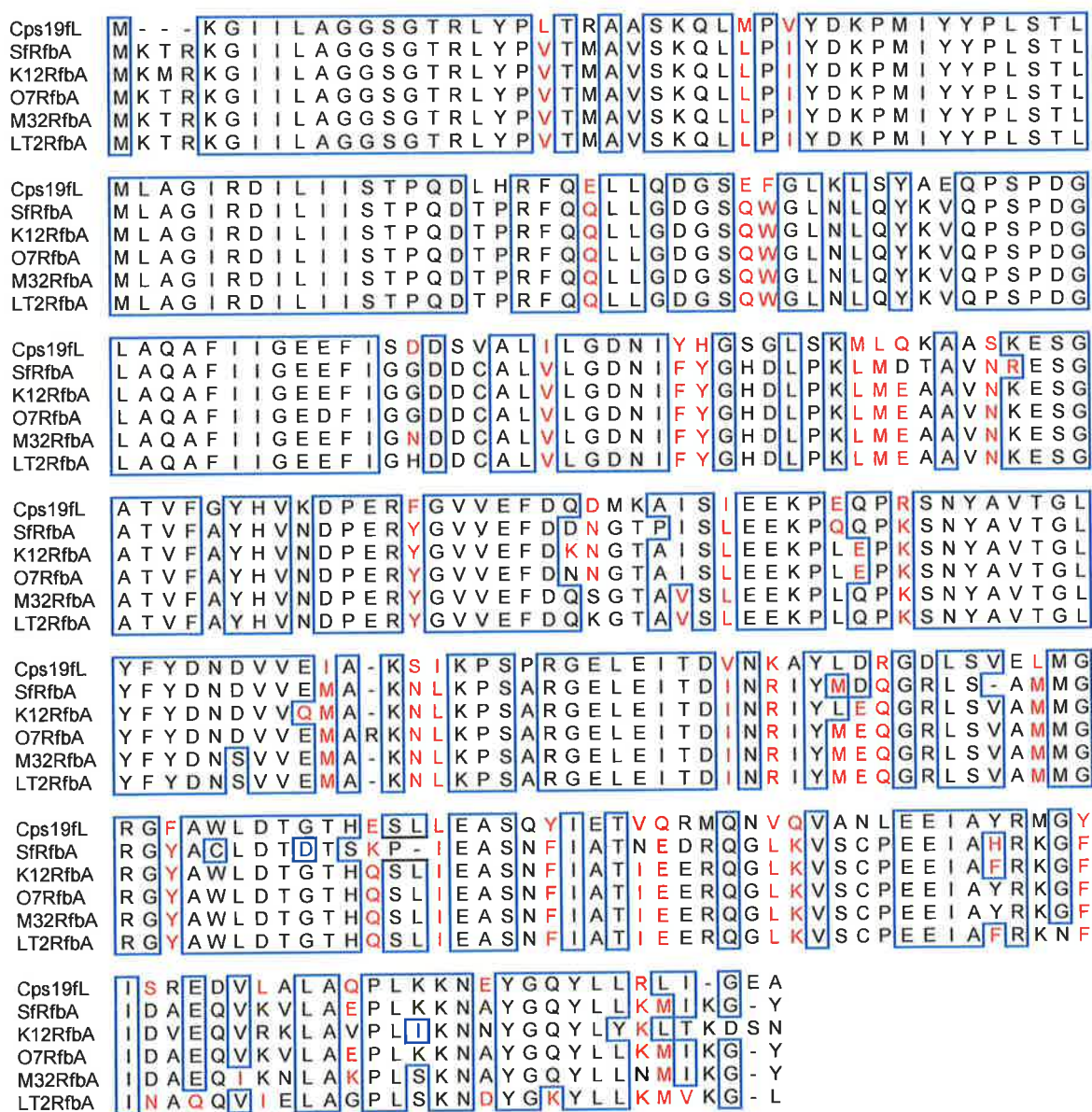
n. *Streptomyces fradiae* TylA (Merson-Davies and Cundliffe, 1994)

o. *S. griseus* StrD (Distler *et al.*, 1987)

p. *Streptomyces peucetius* DnrL (Gallo *et al.*, 1996)

q. *Streptomyces nodosus* SnoD (GenBank accession no. A25110)

r. *Streptomyces violaceoruber* GraD (Bechthold *et al.*, 1995)



**Fig. 3.8. Alignment of Cps19fL.** Alignment of the amino acid sequences of Cps19fL, *S. flexneri* RfbA (SfRfbA) (Macpherson *et al.*, 1994), *E. coli* K12 RfbA (K12RfbA) (Stevenson *et al.*, 1994), *E. coli* O7 RfbA (O7RfbA) (Marolda and Valvano, 1995), *S. enterica* E1 RfbA (M32RfbA) (Wang *et al.*, 1992) and *S. enterica* serovar typhimurium LT2 RfbA (LT2RfbA) (Jiang *et al.*, 1991), as determined using the default settings of the program CLUSTAL (Higgins and Sharp, 1988). Residues identical to those in Cps19fL are boxed; similar residues are shown in red; - indicates absence of a residue.

**Table 3.6. Similarity of Cps19fM to other proteins.**

	% Identity <sup>a</sup>								
	XcRfbC	YeRfbF	M32RfbC	LT2RfbC	NmRfbC	SfRfbC	07RfbC	K12RfbC	SgStrM
Cps19fM <sup>t</sup>	34.6 [179]	32.3 [161]	31.6 [177]	30.9 [178]	31.5 [178]	31.8 [157]	24.3 [177]	29.3 [184]	29.1 [141]
XcRfbC <sup>c</sup>	100	52.7 [150]	56.7 [157]	58 [157]	46.5 [187]	51.7 [178]	50.9 [171]	52.8 [159]	33.3 [156]
YeRfbF <sup>d</sup>		100	52.9 [166]	55.4 [175]	52 [173]	54 [174]	48 [177]	49.7 [181]	39.4 [175]
M32RfbC <sup>e</sup>			100	70.4 [189]	57.7 [175]	83.3 [168]	71 [176]	66.3 [169]	36.3 [168]
LT2RfbC <sup>f</sup>				100	58.7 [179]	71.9 [178]	65.7 [175]	66.1 [177]	35.3 [167]
NmRfbC <sup>g</sup>					100	59.1 [181]	51.2 [172]	55.6 [169]	37.9 [161]
SfRfbC <sup>h</sup>						100	73.6 [178]	66.9 [169]	35.5 [172]
07RfbC <sup>i</sup>							100	58.6 [181]	37.9 [174]
K12RfbC <sup>j</sup>								100	34.8 [178]

a. Percentage of identical amino acids determined with FASTA as implemented in PROSIS. Numbers in parentheses indicate the number of amino acids over which the % identity occurs.

b. *S. pneumoniae* Cps19fM

c. *X. campestris* RfbC (Koplin *et al.*, 1993)

d. *Y. enterocolitica* RfbF (Zhang *et al.*, 1993)

e. *S. enterica* E1 RfbC (Wang *et al.*, 1992)

f. *S. enterica* serovar typhimurium strain LT2 RfbC (Jiang *et al.*, 1991)

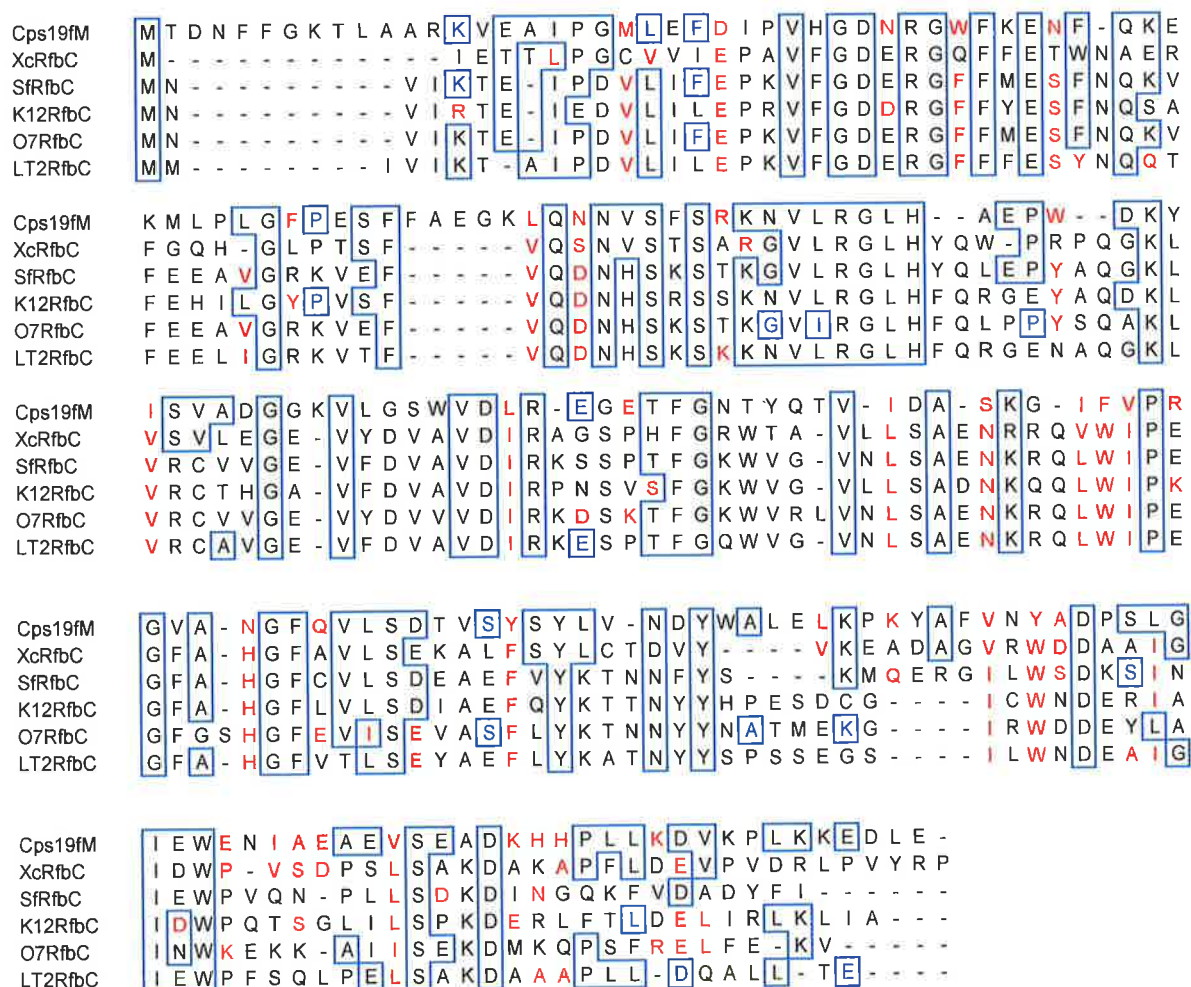
g. *N. meningitidis* RfbC (Hammerschmidt *et al.*, 1994)

h. *S. flexneri* RfbC (Macpherson *et al.*, 1994)

i. *E. coli* O7 RfbC (Marolda and Valvano, 1995)

j. *E. coli* K12 RfbC (Stevenson *et al.*, 1994)

k. *S. griseus* StrM (Distler *et al.*, 1987)



**Fig. 3.9. Alignment of Cps19fM.** Alignment of the amino acid sequences of Cps19fM, *X. campestris* RfbC (XcRfbC) (Koplin *et al.*, 1993), *S. flexneri* RfbC (SfRfbC) (Macpherson *et al.*, 1994), *E. coli* K12 RfbC (K12RfbC) (Stevenson *et al.*, 1994), *E. coli* O7 RfbC (O7RfbC) (Marolda and Valvano, 1995) and *S. enterica* serovar typhimurium LT2 RfbC (LT2RfbC) (Jiang *et al.*, 1991), as determined using the default settings of the program CLUSTAL (Higgins and Sharp, 1988). Residues identical to those in Cps19fM are boxed; similar residues are shown in red; - indicates absence of a residue.

**Table 3.7. Similarity of Cps19fN to other proteins.**

	% Identity <sup>a</sup>														
	Nm RfbB	Sy RfbB	Xc RfbB	Se Gdh	Bs SpsJ	Sv GraE	Ec RffG	LT2 RfbB	Sf RfbB	Stf TylB	K12 RfbB	Sg StrE	Ng RfbB	O7 RfbB	
Cps19fN <sup>b</sup>	46.3 [350]	42 [338]	44.1 [333]	44.7 [338]	45.8 [325]	43.6 [330]	44.8 [355]	45.5 [325]	44 [343]	43.2 [333]	45.8 [343]	42.4 [330]	43.2 [340]	46.1 [343]	
NmRfbB <sup>c</sup>	100	57.1 [343]	60.5 [349]	46.1 [343]	48.3 [333]	42.2 [339]	75.5 [351]	71.1 [352]	69.3 [352]	45.2 [343]	71.9 [352]	45.8 [345]	91.9 [345]	71.3 [352]	
SyRfbB <sup>d</sup>		100	62.1 [338]	44.9 [339]	47.9 [334]	47 [338]	58.7 [339]	54.3 [341]	54 [341]	47.8 [339]	55.1 [341]	46.5 [340]	55.4 [341]	54.8 [341]	
XcRfbB <sup>e</sup>			100	46.9 [331]	49.5 [331]	45.8 [336]	58.9 [348]	56.7 [349]	56.6 [328]	47.2 [335]	56.7 [349]	45.3 [333]	57.9 [335]	57.9 [349]	
SeGdh <sup>f</sup>				100	49.7 [314]	57.2 [318]	44.3 [348]	41 [284]	35.7 [350]	66.1 [327]	44.7 [284]	61.9 [323]	44 [341]	44.4 [284]	
BsSpsJ <sup>g</sup>					100	45.9 [314]	47.6 [334]	47 [270]	45.2 [330]	46.5 [314]	44.6 [336]	44.5 [317]	47.9 [333]	45.9 [270]	
SvGraE <sup>h</sup>						100	45 [342]	27.8 [224]	43.5 [283]	57.9 [318]	41.9 [344]	60.9 [317]	43.8 [336]	43.8 [283]	
EcRffG <sup>i</sup>							100	72.3 [354]	70.6 [354]	42.2 [348]	72.3 [354]	43.1 [341]	73.2 [329]	71.8 [354]	
LT2RfbB <sup>j</sup>								100	85 [361]	42 [281]	87.8 [361]	43.5 [271]	69.1 [346]	87.5 [361]	
SfRfbB <sup>k</sup>									100	37.7 [332]	90.3 [361]	46.5 [269]	66.8 [346]	90.6 [361]	
StfTylB <sup>l</sup>										100	39.2 [302]	67.9 [320]	44.9 [341]	37.3 [332]	
K12RfbB <sup>m</sup>											100	33.1 [311]	68.8 [346]	91.1 [361]	
SgStrE <sup>n</sup>												100	45.3 [340]	34.1 [311]	
NgRfbB <sup>o</sup>													100	68.8 [346]	
O7RfbB <sup>p</sup>														100	

a. Percentage of identical amino acids determined with FASTA as implemented in PROSIS. Numbers in parentheses indicate the number of amino acids over which the % identity occurs.

b. *S. pneumoniae* Cps19fN

c. *N. meningitidis* RfbB (Hammerschmidt *et al.*, 1994)

d. *Synechocystis* sp. RfbB (GenBank accession no. D64003)

e. *X. campestris* RfbB (Koplin *et al.*, 1993)

f. *Saccharopolyspora erythraea* Gdh (Linton *et al.*, 1995)

g. *B. subtilis* SpsJ (GenBank accession no. Z99123)

h. *S. violaceoruber* GraE (Bechthold *et al.*, 1995)

i. *E. coli* RffG (Daniels *et al.*, 1992)

j. *S. enterica* serovar typhimurium strain LT2 RfbB (Jiang *et al.*, 1991)

k. *S. flexneri* RfbB (Macpherson *et al.*, 1994)

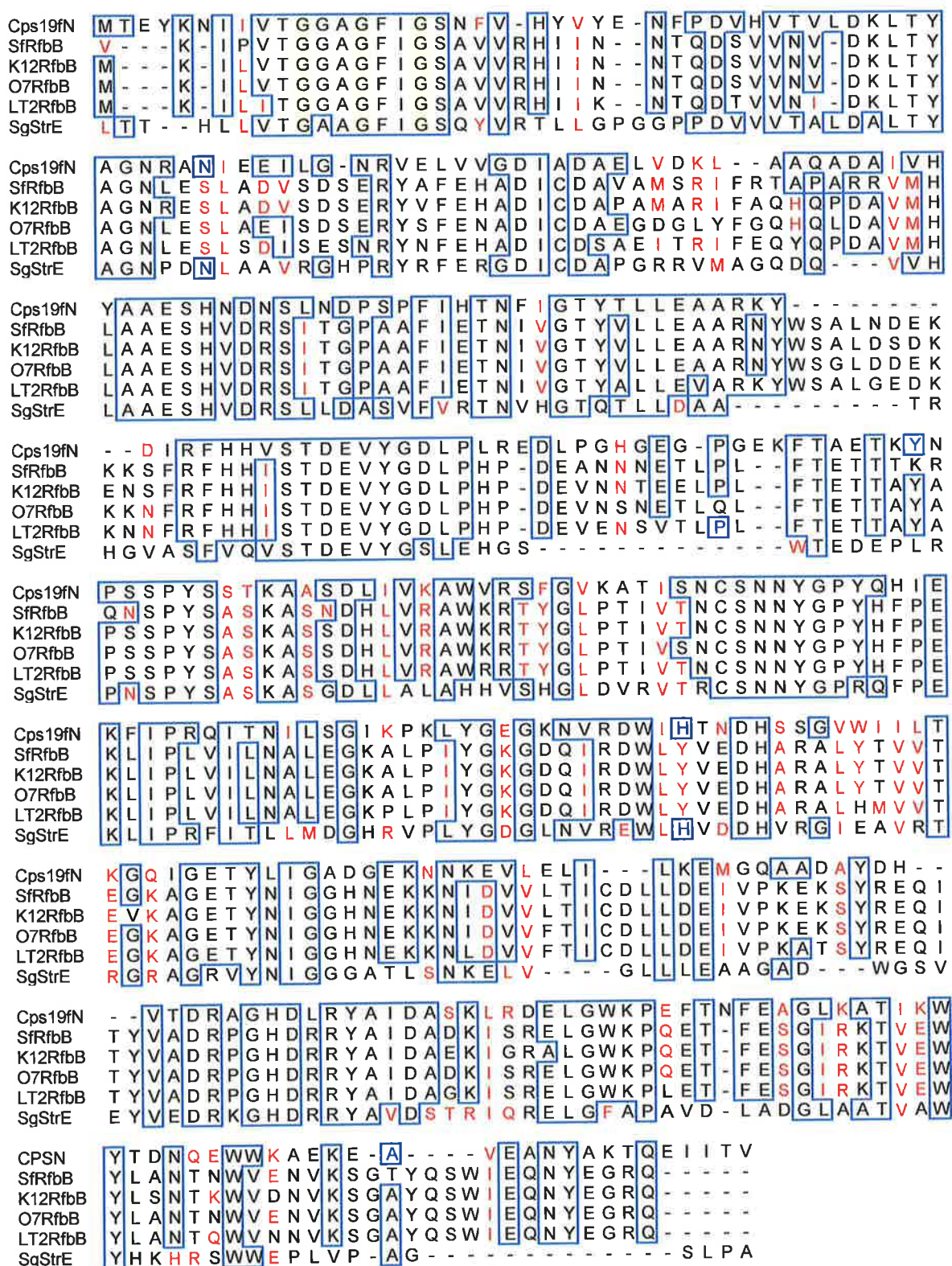
l. *S. fradiae* TylB (Merson-Davies and Cundliffe, 1994)

m. *E. coli* K12 RfbB (Stevenson *et al.*, 1994)

n. *S. griseus* StrE (Distler *et al.*, 1987)

o. *N. gonorrhoeae* RfbB (Robertson *et al.*, 1994)

p. *E. coli* O7 RfbB (Marolda and Valvano, 1995)



**Fig. 3.10. Alignment of Cps19fN.** Alignment of the amino acid sequences of Cps19fN, *S. flexneri* RfbB (SfRfbB) (Macpherson *et al.*, 1994), *E. coli* K12 RfbB (K12RfbB) (Stevenson *et al.*, 1994), *E. coli* O7 RfbB (O7RfbB) (Marolda and Valvano, 1995), *S. enterica* serovar typhimurium LT2 RfbB (LT2RfbB) (Jiang *et al.*, 1991) and *S. griseus* StrE (SgStrE) (Distler *et al.*, 1992), as determined using the default settings of the program CLUSTAL (Higgins and Sharp, 1988). Residues identical to those in Cps19fN are boxed; similar residues are shown in red; - indicates absence of a residue. The yellow shaded region corresponds to a NAD-binding motif (Macpherson *et al.*, 1994).

**Table 3.8. Similarity of Cps19fO to other proteins.**

	% Identity <sup>a</sup>								
	SgStrL	LT2RfbD	K12RfbD	BsSpsK	O7RfbD	XcRfbD	YeRfbG	SvGraD	SfRfbD
Cps19fO <sup>b</sup>	38.6 [285]	35 [286]	32.4 [296]	39.7 [257]	31.9 [298]	27.4 [252]	27 [285]	37.4 [297]	30.6 [297]
SgStrL <sup>c</sup>	100	35.4 [291]	34.4 [291]	39.9 [258]	33.2 [292]	32 [297]	32.8 [195]	49.7 [292]	32.5 [292]
LT2RfbD <sup>d</sup>		100	84.6 [299]	35.7 [244]	84.4 [301]	29.1 [299]	30.7 [231]	37.9 [290]	80.7 [300]
K12RfbD <sup>e</sup>			100	35.2 [273]	91 [301]	39.5 [299]	31.9 [229]	35.7 [294]	87.7 [300]
BsSpsK <sup>f</sup>				100	441 [276]	21.3 [227]	29.5 [251]	39.2 [278]	33.5 [236]
O7RfbD <sup>g</sup>					100	37.9 [301]	30.2 [232]	34.8 [296]	87.7 [302]
XcRfbD <sup>h</sup>						100	22.8 [263]	36.4 [305]	36.2 [301]
YeRfbG <sup>i</sup>							100	29.6 [206]	30.4 [230]
SvGraD								100	33 [294]
SfRfbD <sup>k</sup>									100

a. Percentage of identical amino acids determined with FASTA as implemented in PROSIS. Numbers in parentheses indicate the number of amino acids over which the % identity occurs.

b. *S. pneumoniae* Cps19fO

c. *S. griseus* StrL (Distler *et al.*, 1987)

d. *S. enterica* serovar typhimurium LT2 RfbD (Jiang *et al.*, 1991)

e. *E. coli* K12 RfbD (Stevenson *et al.*, 1994)

f. *B. subtilis* SpsK (GenBank accession no. Z99123)

g. *E. coli* O7 RfbD (Marolda and Valvano, 1995)

h. *X. campestris* RfbD (Koplin *et al.*, 1993)

i. *Y. enterocolitica* RfbG (Zhang *et al.*, 1993)

j. *S. violaceoruber* GraD (Bechthold *et al.*, 1995)

k. *S. flexneri* RfbD (Macpherson *et al.*, 1994)



### 3.2.4 Insertion-duplication mutagenesis of *cps19fG-O* genes

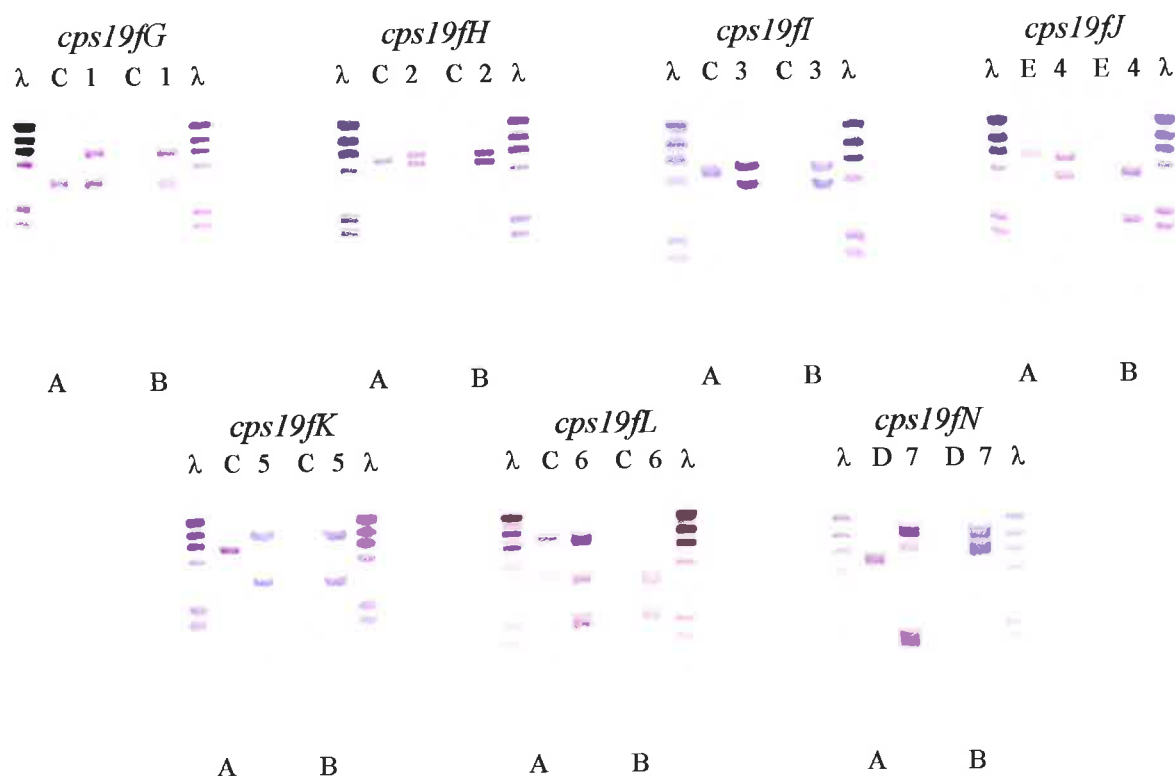
In order to confirm the involvement of the ORFs designated *cps19fG* to *cps19fO* in type 19F CPS production, the chromosomal copies of the respective genes in the encapsulated strain Rx1-19F were individually disrupted by insertion-duplication mutagenesis (see section 2.10). To achieve this, a small internal segment of each ORF was obtained either by band isolation of an appropriate restriction fragment or by PCR amplification, and then cloned into the plasmid pVA891 (Table 3.9). Recombinant plasmids were then transformed into Rx1-19F (section 2.4.4). As described in section 2.10, Ery-resistant transformants are the result of a homologous recombination event directed by the cloned fragment of pneumococcal DNA that leads to the integration of the pVA891 plasmid into the host chromosome, such that the cloned fragment is duplicated, and consequent disruption of the gene of interest, as illustrated in Fig. 2.1.

Ery-resistant *S. pneumoniae* Rx1-19F transformants were obtained for all except the *cps19fM*- and *cps19fO*-containing pVA891 derivatives. As these two genes are part of the same biochemical pathway as *cps19fL* and *cps19fN*, no further attempts were made to disrupt these genes. Correct disruption of the respective chromosomal ORF was confirmed

**Table 3.9. Summary of *cps19f* gene fragments inserted into pVA891**

ORF	Nucleotides	Fragment cloned into pVA891
<i>cps19fG</i>	6,317-6,603	<i>Hind</i> III- <i>Sau</i> 3A fragment from pJCP464
<i>cps19fH</i>	7,054-7,247	<i>Eco</i> RV- <i>Sau</i> 3A fragment from pJCP464
<i>cps19fI</i>	8,179-8,440	<i>Hae</i> III- <i>Sau</i> 3A fragment from pJCP465
<i>cps19fJ</i>	9,651-10,165	<i>Eco</i> RI- <i>Eco</i> RV fragment from pJCP466
<i>cps19fK</i>	10,708-11,177	PCR product using primers J17 and J18 (Table 2.4)
<i>cps19fL</i>	11,600-11,958	<i>Sma</i> I- <i>Bam</i> HI fragment from pJCP467
<i>cps19fM</i>	12,549-12,893	PCR product using primers J28 and J29 (Table 2.4)
<i>cps19fN</i>	13,148-13,961	PCR product using primers J30 and J31 (Table 2.4)
<i>cps19fO</i>	14,270-14,679	PCR product using primers J24 and J25 (Table 2.4)

by Southern hybridisation analysis of each insertion-duplication mutant generated (Fig 3.12). CPS gene-specific probes hybridised with two separate DNA fragments in the respective mutants but with only one fragment in the Rx1-19F control, indicating that the gene is disrupted in the mutants. As expected, the pVA891-specific probe hybridised only with the mutants. For Rx1-19F-J, Rx1-19F-L and Rx1-19F-N, the hybridisation pattern is different for the two probes because the enzyme used to restrict the chromosomal DNA was also used to clone the appropriate fragment into pVA891. In such cases, the interrupted *cps* gene and pVA891-related sequences will be on distinct chromosomal restriction fragments. The *cps19fG*, *cps19fH*, *cps19fI*, *cps19fJ*, *cps19fK*, *cps19fL* and *cps19fN* mutants all exhibited a rough phenotype and did not produce a type 19F capsule, as judged by quellung reaction, confirming that all are part of the *cps19f* locus. On the



**Fig. 3.12. Southern hybridisation analysis of insertion-duplication mutants.** Southern blots of *Cla*I-restricted chromosomal DNA from Rx1-19F (C), Rx1-19F-G (1), Rx1-19F-H (2), Rx1-19F-I (3), Rx1-19F-K (5), and Rx1-19F-L (6); *Eco*RI-restricted chromosomal DNA from Rx1-19F (E) and Rx1-19F-J (4); *Bam*HI-restricted chromosomal DNA from Rx1-19F (D) and Rx1-19F-N (7) were probed with the gene specific probes indicated above each blot (A) and with a pVA891 specific probe (B). The sizes of the DIG-labelled lambda ( $\lambda$ ) markers are as follows: 23 kb, 9.4 kb, 6.6 kb, 4.4 kb, 2.3 kb and 2.0 kb.

other hand, insertion of pVA891 into the apparent non-coding region downstream of *cps19fO* (nucleotides 15,149-15,444 were used as the target) did not interfere with type 19F capsule expression in *S. pneumoniae* Rx1-19F.

### 3.2.5 T7 expression of several Cps19f proteins

Expression of *cps19fH*, *cps19fI*, *cps19fJ*, *cps19fK*, *cps19fL*, *cps19fM* *cps19fN* and *cps19fO* cloned in *E. coli* DH5 (pGP1-2), in which the T7 RNA polymerase gene is under the control of lambda P<sub>L</sub>/cI<sub>857</sub> (Tabor and Richardson, 1985), was investigated. Expression of these genes from the vector T7 promoter was performed as described in section 2.12.1 and analysed on SDS-PAGE gels stained with Coomassie Brilliant Blue (section 2.12.3).

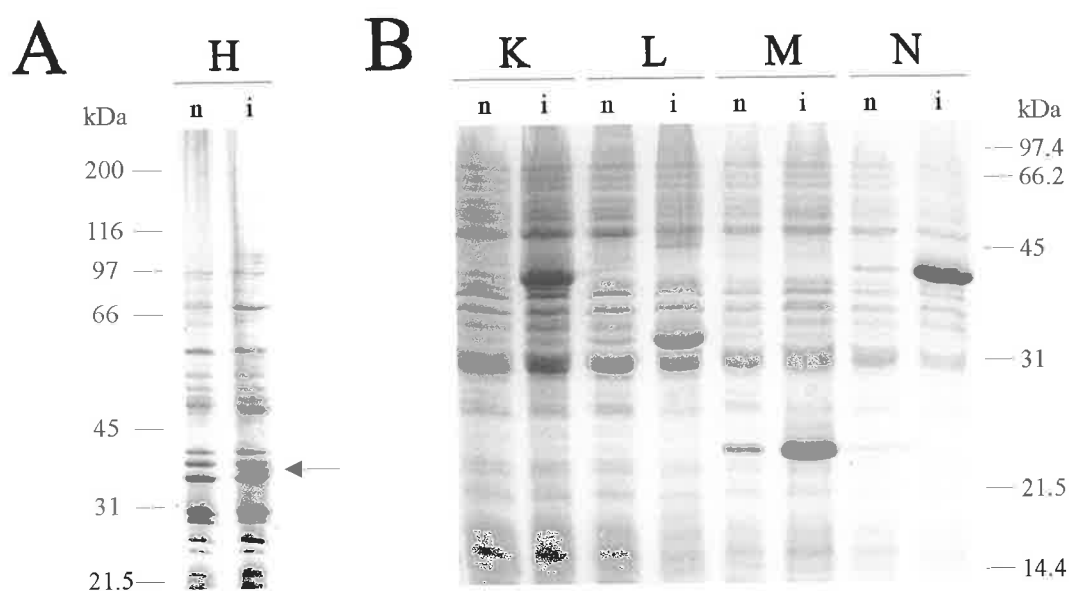
To examine the expression of *cps19fH* in *E. coli*, one of the nested deletion derivatives of pJCP464, containing *cps19fH* (nucleotides 6,578-7,846) downstream from the vector T7 promoter, was transformed into *E. coli* DH5(pGP1-2). Induction of cultures at 42°C for 30 min followed by 37°C for 2 hours (section 2.12.1) resulted in the expression of low levels of a 34-35 kDa protein in cells containing the pJCP464 derivative, which was not seen in uninduced cultures (see arrow, **Fig. 3.13A**). Both the size of the translation product and the low level of expression of *cps19fH* is consistent with initiation at the TTG codon as indicated in **Appendix I**.

The *cps19fI* and *cps19fJ* genes were PCR amplified using primer pairs J11/J14 and J15/J16, respectively (primers are described in **Table 2.4**) and cloned into pGEM-7Zf(+) such that the ORFs were downstream of the T7 promoter. These plasmids were transformed into *E. coli* DH5(pGP1-2). However, induction of cultures at 42°C did not result in increased production of any polypeptides with sizes corresponding to these ORFs, as judged by SDS-PAGE analysis (result not shown). However, this was not unexpected given previous observations that very hydrophobic proteins such as Cps19fI and Cps19fJ

are poorly expressed in this system (Morona *et al.*, 1994b).

To determine whether *cps19fK* was expressed in *E. coli*, a deletion derivative of pJCP466 (designated pJCP470), which placed nucleotides 10,304-11,720 downstream from the vector T7 promoter, was transformed into *E. coli* DH5(pGP1-2). Induction of the culture resulted in the expression of a polypeptide of approximately 41 kDa, which is consistent with the size of Cps19fK predicted from the DNA sequence (Fig. 3.13B).

Plasmids for T7 expression of *cps19fL*, *cps19fM*, *cps19fN* and *cps19fO* were constructed by cloning an appropriate PCR product into pBluescript SK+ or KS+ as shown in Table 3.10 and transforming these constructs into *E. coli* DH5(pGP1-2). Induction of the genes under T7 control (section 2.12.1) resulted in the expression of polypeptides of approximately 32, 22 and 39 kDa for *cps19fL*, *M* and *N*, respectively, as predicted from the DNA sequence (Fig. 3.13B). However, no protein product with a size corresponding to that predicted for Cps19fO was detected when cells containing the *cps19fO* construct were induced (result not shown) and this was not investigated further.



**Fig. 3.13. T7-expression of *cps19f* genes in *E. coli*.** *E. coli* DH5(pGP1-2) was transformed with various plasmids containing individual *cps19fH* (H), *cps19fK* (K), *cps19fL* (L), *cps19fM* (M) or *cps19fN* (N) ORFs. Lysates of cultures incubated at 42°C for 2 h to induce expression of T7 RNA polymerase (i), and non-induced cultures (n) were separated by SDS-PAGE and stained with Coomassie Brilliant Blue. The mobilities of molecular size markers are shown separately for panels A and B.

**Table 3.10. Plasmid constructs for T7 expression of *cps19fL-O***

Gene	Primers <sup>a</sup>	Restriction sites	<i>cps19f</i> sequence <sup>b</sup>	Vector
<i>cps19fL</i>	J21-J29	<i>Hind</i> III- <i>Eco</i> RV	11,347-12,487	pBluescript SK+
<i>cps19fM</i>	J12-J31	<i>Bam</i> HI- <i>Hind</i> II	11,968-13,134	pBluescript KS+
<i>cps19fN</i>	J25-J28	<i>Hind</i> III- <i>Xba</i> I	12,547-14,129	pBluescript SK+
<i>cps19fO</i>	J37-J38	<i>Hind</i> III- <i>Bam</i> HI	14,009-15,562	pBluescript SK+

a. Primers sequences are as described in **Table 2.4**

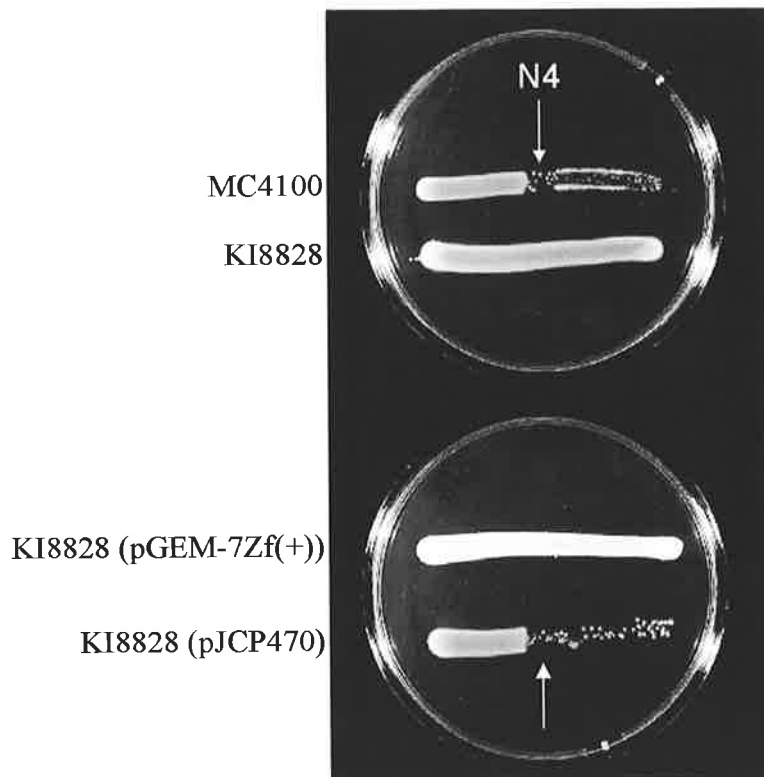
b. The corresponding nucleotide numbers of the *cps19f* sequence (**Appendix I**)

### 3.2.6 Complementation of an *E. coli rffE* mutant with *cps19fK*

Whilst the precise mechanism is not understood, *E. coli rffE* (*nfrC*) mutants are known to be resistant to infection with bacteriophage N4 (Kiino *et al.*, 1993). In order to determine whether Cps19fK and *E. coli* RffE are functional homologues, the *rffE*-negative mutant *E. coli* KI8828 (Kiino *et al.*, 1993) was transformed with pJCP470 (which contains only the complete *cps19fK* ORF), or with pGEM-7Zf(+). The various strains, including the *rffE*-positive wild type parent *E. coli* MC4100 (Silhavy *et al.*, 1984), were then examined for susceptibility to infection with N4 phage (**Fig. 3.14**). Transformation with pJCP470, but not pGEM-7Zf(+), clearly conferred susceptibility to N4 on *E. coli* KI8828. Thus, Cps19fK is a functional RffE homologue and is therefore likely to be a UDP-GlcNAc-2-epimerase.

### 3.2.7 Complementation of *S. flexneri rfbBDAC* with *cps19fLMNO*

As mentioned in section 3.2.3.6, database searches indicated that the *cps19fL-O* region had significant similarity to a portion of the *S. flexneri rfb* region (*rfbBDAC*) responsible for biosynthesis dTDP-Rha. To confirm functional homology, a portion of the *S. pneumoniae* Rx1-19F chromosome containing *cps19fL-O* (equivalent to nucleotides 11,351-15,449 in **Appendix I**) was PCR amplified using primers J21 and J36 (**Table 2.4**) and cloned into the *Hind*III and *Eco*RI sites of pK194. The recombinant plasmid



**Fig. 3.14. N4 bacteriophage susceptibility.** A suspension of bacteriophage N4 ( $10^9$  pfu/ml) was streaked vertically on both plates as indicated by the arrows. The indicated *E. coli* strain was then streaked from left to right across the phage streak and incubated for 18 h at 37°C.

(designated pJCP471), or pK194, was then transformed into *E. coli* SØ874 containing pPM2716. The latter plasmid, from which *rfbBDAC* has been deleted, is a derivative of pPM2213 which contains the complete *S. flexneri rfb* region and directs the expression of *S. flexneri* serotype 4 O-antigen in *E. coli* K12 (Macpherson *et al.*, 1994). Lysates of *E. coli* SØ874 containing pPM2213, pPM2716, pPM2716 + pJCP470, or pPM2716 + pK194, were then subjected to Western blot analysis (section 2.12.4) using a rabbit antiserum raised against *S. flexneri* serotype 4 O-antigen (Fig. 3.15). Immunoreactive O-antigen can be seen in both the pPM2213 and the pPM2716 + pJCP470 tracks, indicating that *cps19FLMNO* can complement the *S. flexneri* serotype 4 *rfbBDAC* deletion in *E. coli*. Thus, Cps19fL is a Glc-1-phosphate thymidyl transferase, Cps19fM is a dTDP-4-keto-6-deoxyglucose-3,5- epimerase, Cps19fN is a dTDP-Glc-4,6-dehydratase and Cps19fO is a dTDP-L-Rha synthase.

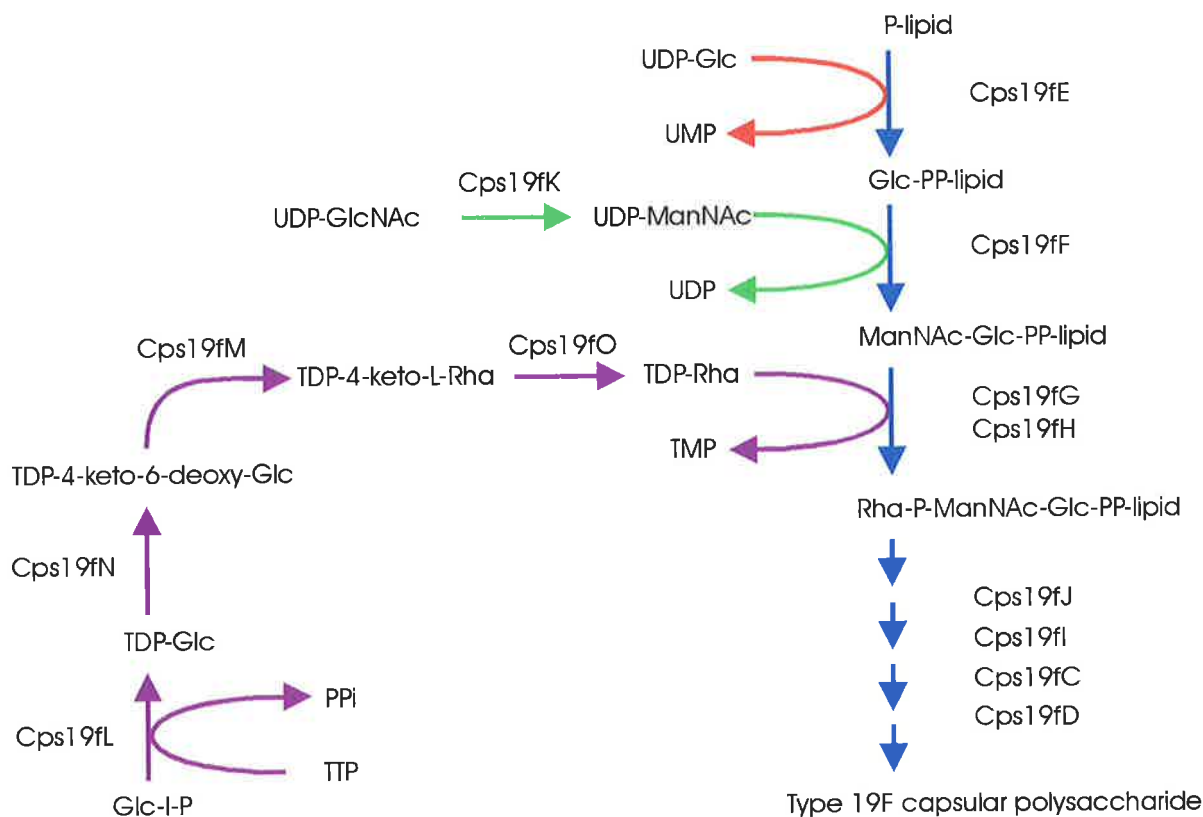
1 2 3 4



**Fig. 3.15. Complementation of *S. flexneri* *rfbBDAC* by *cps19fLMNO*.** *E. coli* lysates (prepared as described in section 2.12.2) were separated by SDS-PAGE, electroblotted onto nitrocellulose, probed with anti-*S. flexneri* type 4 serum, and developed as described in sections 2.12.3 and 2.12.4. Lanes: 1, *E. coli* SØ874 containing pPM2213; 2, *E. coli* SØ874 containing pPM2716; 3, *E. coli* SØ874 containing pPM2716 + pJCP470; 4, *E. coli* SØ874 containing pPM2716 + pK194.

### 3.2.8 Putative biosynthetic pathway for *S. pneumoniae* type 19F CPS

Functions have been previously proposed for the majority of the first six *cps19f* gene products on the basis of database comparisons (Guidolin *et al.*, 1994). Cps19fE is believed to catalyse the first step in CPS biosynthesis, namely transfer of Glc-1-phosphate to a lipid carrier (presumably und-P). Cps19fF is a putative ManNAc transferase, while Cps19fC and Cps19fD are believed to be involved with chain length regulation and export of CPS. Collectively, this information and the findings of the current study indicate that the trisaccharide biological repeat unit of type 19F CPS is:  $\rightarrow 2$ )- $\alpha$ -L-Rha<sub>p</sub>-(1-PO<sub>4</sub><sup>-</sup> $\rightarrow 4$ )- $\beta$ -D-Man<sub>p</sub>NAc-(1 $\rightarrow 4$ )- $\alpha$ -D-Glc<sub>p</sub>-(1 $\rightarrow$ ), i.e. the first sugar in the repeat unit is Glc, not Rha as shown in previously published structures (Lee and Fraser, 1980; Ohno *et al.*, 1980). A putative biosynthetic pathway for type 19F CPS is shown in **Fig. 3.16**.



**Fig. 3.16. Putative biosynthetic pathway for *S. pneumoniae* type 19F CPS.**

Biosynthesis of type 19F CPS probably occurs via a mechanism analogous to that proposed for Rol/Cld- and Rfc-dependent O-antigen assembly in *Salmonella enterica* serogroups B and E (section 1.8.2; Whitfield, 1995). The initial step, catalysed by Cps19fE, involves transfer of Glc-1-phosphate from UDP-Glc to a lipid carrier on the cytoplasmic face of the cell membrane. Cps19fE has several large hydrophobic domains in its N-terminal portion, which would anchor it to the membrane and facilitate interaction with the lipid carrier (Guidolin *et al.*, 1994). Cps19fF, G and H then catalyse the sequential transfer of the other component monosaccharide precursors (synthesised in the cytoplasm by Cps19fK, L, M, N, and O) to form the trisaccharide repeat unit. These lipid-linked repeat units are then translocated from the cytoplasmic to the extracellular side of the cell membrane and polymerised in a blockwise fashion, extending the polysaccharide at the reducing terminus. These two steps are catalysed by Cps19fJ and I, respectively, both

of which are integral membrane proteins. The process of translocation and polymerisation may be closely linked and the two proteins possibly form a complex in the membrane.

In *Salmonella* and *Shigella* O-antigen assembly, Rol/Cld has been proposed to regulate chain length by modulating interaction between the lipid-linked nascent O-antigen and either the polymerase Rfc, or RfaL, a ligase responsible for transfer of O-antigen to the lipid A core molecule (Morona *et al.*, 1995; Whitfield, 1997). Cps19fC and D may perform a similar function in *S. pneumoniae* type 19F. Pneumococcal CPS is believed to be covalently linked to the cell wall peptidoglycan (Sørensen *et al.*, 1990), but the precise nature of this linkage and the enzyme responsible are unknown.

### 3.3 Conclusions

Sequential rounds of InPCR and plasmid insertion/rescue were used to isolate the region of the *S. pneumoniae* type 19F chromosome responsible for CPS biosynthesis. The data presented in this chapter, combined with that which was described previously (Morona *et al.*, 1994a; Guidolin *et al.*, 1994) indicates that the *cps19f* locus consists of 15 genes, which are tightly clustered on the chromosome. The *cps19f* locus is flanked by *dexB* and IS1202 at the 5' end and by *aliA* at the 3' end. Arrecubieta *et al* (1995) have also shown that the *cps* locus of *S. pneumoniae* type 3 is flanked at the 5' end by *dexB*. Dillard *et al* (1995) have reported that the sequences flanking the 3' end of the *cps* locus of a different type 3 *S. pneumoniae* strain were homologous to *plpA*. Examination of the published nucleotide sequences of *plpA* (Pearce *et al.*, 1994) and *aliA* (Alloing *et al.*, 1994) indicates that these sequences describe the same gene and so the chromosomal location reported for the *cps3* locus is identical to that for *cps19f*. However, the *aliA* gene is not functional in *S.*

*pneumoniae* type 3, as it contains a large deletion at the 5' end of the gene (Dillard *et al.*, 1995).

Clues as to the likely function of the *cps19f* gene products have been provided by comparisons with known proteins whose sequences have been deposited with databases (as described above). Moreover, for *cps19fK*, *cps19fL*, *cps19fM*, *cps19fN* and *cps19fO*, the function of gene products has been confirmed by complementation of mutations in *E. coli*. This information has been used to propose a biosynthetic pathway for type 19F CPS. However, experimental confirmation of the function of the remaining proteins encoded by *cps19f* will require characterisation of the phenotypic impact of mutagenesis of the respective ORFs. Interpretation of phenotypic data obtained with the insertion-duplication mutants generated in the present study and previously (Guidolin *et al.*, 1994) is complicated by the likelihood of polar effects. Hence, further phenotypic analysis of the function of individual *cps19f* genes awaits the construction of in-frame deletion mutants, in which transcription of the remainder of the type 19F *cps* locus would not be expected to be affected.

## Chapter 4

# ANALYSIS OF CAPSULE LOCI FROM VARIOUS PNEUMOCOCCAL SEROTYPES

## 4.1 Introduction

Analysis of the *cps19f* locus identified genes encoding a range of functions, including those involved in biosynthesis of activated sugar precursors, various glycosyl transferases, a polysaccharide polymerase, and proteins possibly involved in export functions. Some of these specific functions might be expected to be required by all pneumococci, regardless of capsular serotype; others might be expected to be present in various subsets of capsular types (depending on which sugars are present in their repeat units), and some genes might be expected to be serotype-specific. To examine the specificity of individual *cps19f* genes, their distribution amongst diverse *S. pneumoniae* serotypes was examined by Southern hybridisation analysis.

This provided information on the organisation of pneumococcal CPS loci which was used to develop a LR-PCR protocol for isolation of large fragments of the *cps* loci from other *S. pneumoniae* serotypes.

## 4.2 Results

### 4.2.1 Serotype specificity of the *cps19f* genes and flanking DNA sequences

Southern hybridisation analysis (described in section 2.7.2) to determine the serotype specificity of the individual *cps19f* genes was undertaken. DNA fragments corresponding to each of the *cps19f* genes and the flanking regions were isolated from appropriate clones either by restriction enzyme digestion and subsequent purification after agarose gel electrophoresis, or by PCR amplification using specific primers. They were then labelled with DIG (section 2.7.1) and used to probe *Cla*I-restricted chromosomal DNA from representative pneumococci belonging to serotypes 2, 3, 4, 6A, 6B, 7F, 8, 9N, 9V, 12, 14, 16, 17, 18C, 22, 23F, 24, and the other members of serogroup 19 (19A, 19B and 19C). Hybridisation and washing were performed under high stringency conditions (as described in section 2.7.2.2). The results of these 76 Southern hybridisation reactions are summarised in **Table 4.1**.

Large variations in the hybridisation patterns were obtained with the different gene-specific probes. Probes specific for sequences flanking *cps19f* (*dexB*, the 5' intergenic region, the 3' intergenic region and *aliA*) hybridised with all serotypes tested (see section 4.2.1.3). However, of the genes within the *cps* locus, high stringency homologues of only *cps19fA* and *cps19fB* were present in all serotypes tested.

#### 4.2.1.1 Presence of *cps19f* homologues in other members of serogroup 19

Serogroup 19 consists of four immunologically cross-reactive serotypes 19F, 19A, 19B and 19C. The results in **Table 4.1** suggest that the *cps* loci of types 19B and 19C are very closely related to 19F. They hybridised with all the *cps19f* probes except *cps19fI* and *cps19fJ*. Hybridisation with probes specific for the two genes flanking *cps19fI* and *cps19fJ* (*cps19fH* and *cps19fK*) was weaker suggesting a lesser degree of similarity in this region.

**Table 4.1. Hybridization of type 19F *cps* genes and neighbouring sequences with chromosomal DNA from other pneumococcal serotypes.**

Type/Group	DIG labelled DNA probes <sup>a</sup>																			Presence of sugar in capsule <sup>b</sup>					
	<i>dexB</i>	IG 5'	<i>cps19fA-O</i>																IG 3'	<i>aliA</i>	Glu	ManNAc	Rha		
			A	B	C	D	E	F	G	H	I	J	K	L	M	N	O								
19F	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	*	*	*	
19A	+	+	+	+	-	-	-	-	±	±	-	±	-	-	+	+	+	+	+	+	+	+	*	*	*
19B	+	+	+	+	+	+	+	+	+	±	-	-	±	+	+	+	+	+	+	+	+	+	*	*	*
19C	+	+	+	+	+	+	+	+	+	±	-	-	±	+	+	+	+	+	+	+	+	+	*	*	*
2	+	+	+	+	-	-	-	-	-	-	-	-	-	±	+	+	+	+	+	+	+	+	*	*	*
3	+	+	+	+	+	+	-	-	-	-	-	-	-	-	-	-	-	-	±	+	+	*	*	*	
4	+	+	+	+	+	+	-	-	-	-	-	-	-	-	-	-	-	-	+	+	+	*	*	*	
6A	+	+	+	+	-	-	-	-	-	-	-	-	-	-	+	+	+	+	+	+	+	+	*	*	*
6B	+	+	+	+	-	-	-	-	-	-	-	-	-	-	+	+	+	+	+	+	+	+	*	*	*
7F	+	+	+	+	+	+	+	+	-	-	-	-	-	-	+	+	+	+	+	+	+	+	*	*	*
8	+	+	+	+	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	+	+	+	*	*	*
9N	+	+	+	+	+	+	+	+	+	-	-	-	-	-	-	-	-	-	-	+	+	+	*	*	*
9V	+	+	+	+	-	-	-	-	+	-	-	-	-	-	-	-	-	-	-	+	+	+	*	*	*
12	+	+	+	+	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	+	+	+	*	*	*
14	+	+	+	+	+	+	+	+	+	-	-	-	-	-	-	-	-	-	-	+	+	+	*	*	*
16	+	+	+	+	+	+	+	+	+	-	-	-	-	-	+	+	+	+	+	+	+	+	*	*	*
17	+	+	+	+	-	-	-	-	-	-	-	-	-	-	+	+	+	+	+	+	+	+	*	*	*
18C	+	+	+	+	+	+	+	+	+	-	-	-	-	-	+	+	+	+	+	+	+	+	*	*	*
22	+	+	+	+	-	-	-	-	-	-	-	-	-	-	+	+	+	+	+	+	+	+	*	*	*
23F	+	+	+	+	-	-	-	-	-	-	-	-	-	-	+	+	+	+	+	+	+	+	*	*	*
24	+	+	+	+	+	+	+	+	+	-	-	-	-	-	+	+	+	+	+	+	+	+	*	*	*

a. The following DNA fragments were labelled with DIG and used as probes at high stringency: 0.5 kb *Nco* I-*Eco* RI restriction fragment for *dexB* and 0.9 kb *Cla* I-*Nco* I restriction fragment for the 5' intergenic region (IG 5') from previously published sequence (Morona *et al.*, 1994); nucleotides 336-1,468, 1,571-2,380, 2,380-2,998, 3,126-3,739, 3,682-4,979, 5,225-5,725, 6,015-6,630, 6,674-7,731, 7,789-8,965, 9,013-10,278, 10,530-11,539, 11,539-12,493, 12,456-13,139, 13,139-14,134 and 14,134-14,955 for *cpsA-O* genes, respectively; nucleotides 14,955-15,449 for the 3' intergenic region (IG 3') and a 1.2 kb *Hin* dIII restriction fragment for *aliA*. Strong, weak and no hybridization are indicated by +, ± and - respectively.

b. The presence of the sugars glucose (Glu), *N*-acetyl mannosamine (ManNAc) and rhamnose (Rha) in the capsule of each pneumococcal type is indicated by an asterisk.

However, type 19A did not appear to have sequences capable of high-stringency hybridisation with the probes specific for *cps19fC*, *cps19fD*, *cps19fE*, *cps19fF*, *cps19fI*, *cps19fK*, and *cps19fL*, and hybridised only weakly with probes specific for *cps19fG*, *cps19fH* and *cps19fJ*. The apparent dissimilarity of the type 19A and type 19F *cps* loci is intriguing given the structural similarity between the two polysaccharides as shown in **Table 1.3**. To further examine the relationship between *cps19f* and *cps19a*, low stringency Southern hybridisation (section 2.7.2.3) of type 19A chromosomal DNA with the *cps19fK* and *cps19fL* gene probes was undertaken. Under these conditions, both these probes hybridised with type 19A indicating that 19A contained homologues of both *cps19fK* and *cps19fL*.

#### 4.2.1.2 Presence of *cps19f* homologues in other serotypes

Outside of serogroup 19, types 7F, 16, 18C and 24 were the most similar to type 19F, hybridising at high stringency to 9 of the 15 *cps19f* gene probes. These four serotypes all had sequences closely related to *cps19fA-E* and *cps19fL-O*, as shown in **Table 4.1**. The least similar were types 8 and 12, which hybridised only with *cps19fA* and *cps19fB*. Sequences closely related to *cps19fL*, *cps19fM*, *cps19fN* and *cps19fO*, the genes required for dTDP-L-Rha biosynthesis (section 3.2.3.6), were detected in all serotypes tested whose CPS contains L-Rha.

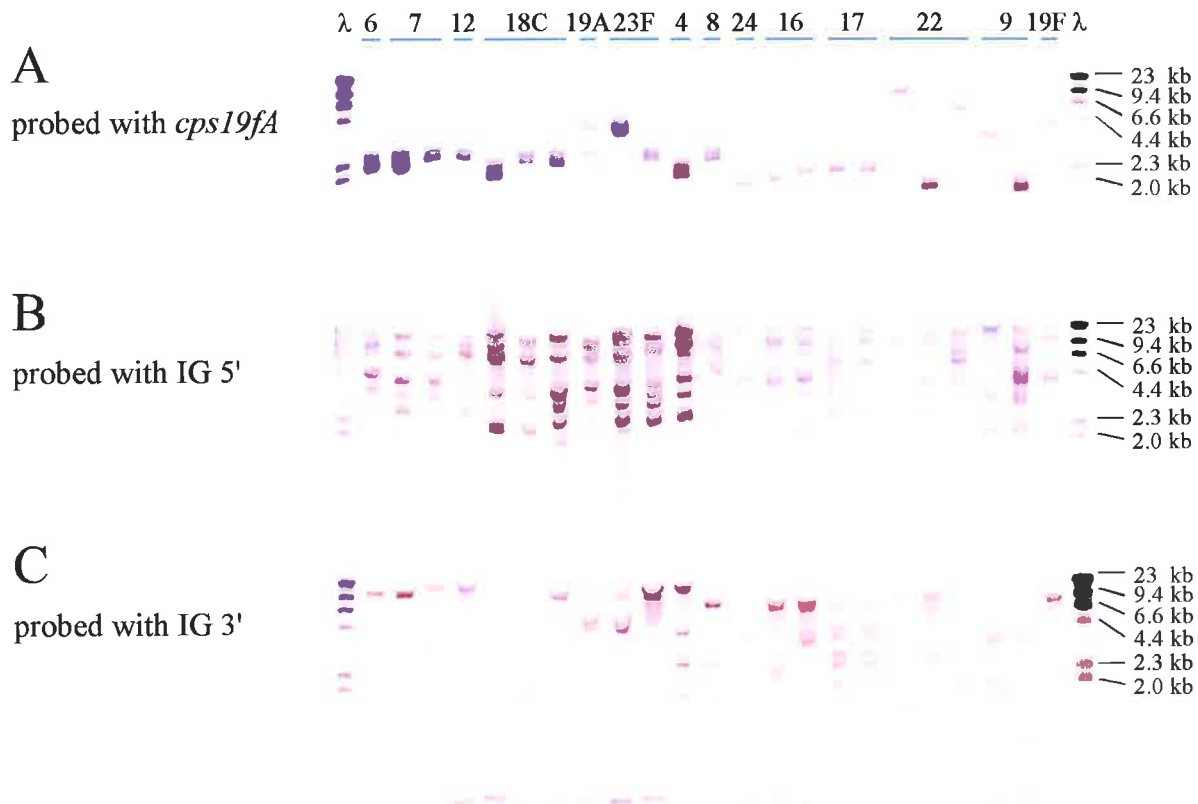
The CPSs of four of the serotypes that were tested contain ManNAc (types 4, 9N, 9V and 12). However, only the two group 9 serotypes contained a gene hybridising to *cps19fF*, which encodes the putative UDP-ManNAc transferase (section 1.9.3). This is consistent with the fact that only group 9 CPS has the same  $\beta$ -D-Man<sub>p</sub>NAc-(1→4)- $\alpha$ -D-Glc<sub>p</sub> linkage seen in group 19 CPS. It was surprising, however, that none of these serotypes contained high or low stringency homologues to *cps19fK*, which encodes a UDP-GlcNAc-2-epimerase (section 3.2.3.5). This enzyme activity is likely to be essential for

synthesis of UDP-ManNAc, and is therefore presumably encoded by a divergent gene in the other three serotypes. Interestingly, the sequence of the genome of *S. pneumoniae* type 4 became available during the course of this project (available from the TIGR Microbial Database [[http://www.tigr.org/pub/data/s\\_pneumoniae/](http://www.tigr.org/pub/data/s_pneumoniae/)]). Examination of the *cps4* locus revealed the presence of a gene (designated *cps4L*, see section 8.3.4) which has 65.9% identity to *cps19fK* and presumably encodes the UDP-GlcNAc-2-epimerase required for type 4 CPS biosynthesis.

Of the various serotypes tested, the CPS of all but type 4 contains Glc (**Table 4.1**). However, types 2, 3, 6A, 6B, 8, 9V, 12, 17, 19A, 22 and 23F also lack sequences which hybridise to *cps19fE*, encoding the type 19F putative Glc-1-phosphate transferase (section 1.9.3). It is possible that Glc may not be the start of the biological repeat unit in some or all of these serotypes. However, glucosyl transferase activity which adds Glc-1-phosphate to a lipid carrier has recently been observed in serotypes 2, 6A, 9V, 12F, 17F, 22F and 19A (Kolkman *et al.*, 1997a), suggesting the presence of a divergent or unrelated *cpsE* gene in these serotypes. As expected, both serotypes 3 and 4 lacked Glc-1-phosphate transferase activity (Kolkman *et al.*, 1997a). Type 4 CPS lacks Glc. The type 3 CPS does contain this sugar, but the capsule is synthesised by a single processive transferase (Dillard *et al.*, 1995) and so does not require a *cps19fE* homologue, as will be discussed in detail in chapter 8.

#### 4.2.1.3 Analysis of the intergenic regions

Southern hybridisation analysis with probes specific for *dexB*, *cps19fA-O* and *aliA* were very specific with the individual probes hybridising with one, or at the most two, DNA fragments for each serotype. This is illustrated in **Fig. 4.1A** in which restricted chromosomal DNA from various *S. pneumoniae* serotypes was probed with *cps19fA*. However, the probes specific for either the 5' or the 3' intergenic regions hybridised with multiple DNA fragments (**Fig. 4.1B** and **C**). Database searches with these intergenic



**Fig. 4.1 Southern hybridisation of pneumococcal chromosomal DNA.** Southern blots of *Cla*I-restricted chromosomal DNA of the indicated *S. pneumoniae* serotypes were probed with DIG-labelled probes specific for *cps19fA* (A), IG 5' (B) and IG 3' (C). The molecular size standards are shown on the right-hand side of the figure and correspond to DIG-labelled *Hind*III-digested λ phage DNA.

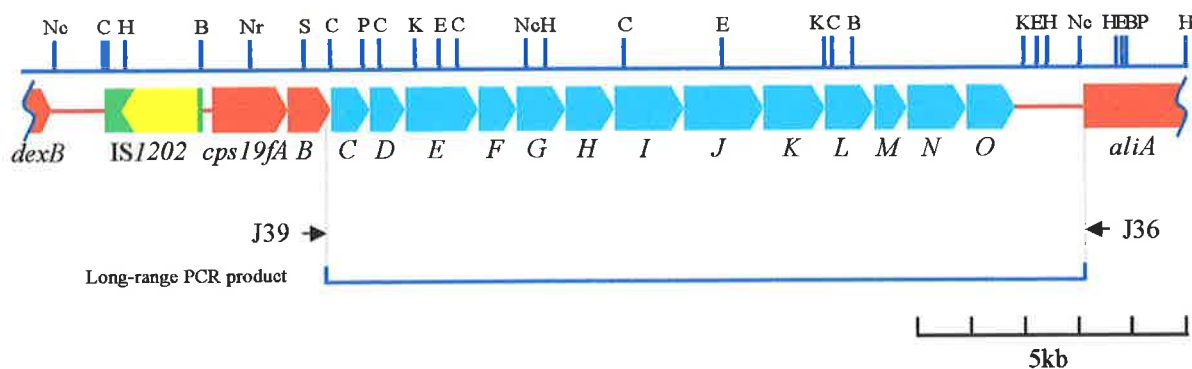
regions against the TIGR *S. pneumoniae* type 4 genomic DNA sequence showed that sections of these regions were, in fact, repeated. A region of 500 nucleotides in the 5' intergenic region, which lies 300 nucleotides downstream from the 3' end of *dexB*, is repeated six times and an adjoining region of 100 nucleotides is repeated at least 20 times. Large sections of the 3' intergenic region are also repeated at least six times in the partially sequenced *S. pneumoniae* type 4 chromosome.

#### 4.2.2 Amplification of capsule loci by LR-PCR

Southern hybridisation analysis indicated that high stringency homologues of *dexB*,

*cps19fA*, *cps19fB* and *aliA* are present in all the pneumococcal serotypes investigated. This suggested that, at least in theory, complete capsule loci of these serotypes could be amplified using LR-PCR with primers based on sequences in *dexB* and *aliA*. However, the distance between the 3' end of *dexB* and the 5' end of *aliA* in type 19F is large (19.1 kb), and the *cps19fA-O* operon itself is 14.8 kb in size (Fig. 4.2). Moreover, the expected size for many *cps* loci would be even larger than that for *cps19f* as their CPS repeat units contain up to 7 sugars, compared with 3 for type 19F (van Dam *et al.*, 1990), and so would require extra genes to encode additional biosynthetic enzymes. Thus the *dexB-aliA* region of these types would be expected to exceed 20-25 kb, and the efficiency of LR-PCR amplification might be low. To increase the likelihood of successfully amplifying LR-PCR products from larger *cps* loci the 5' primer (J39, described in Table 2.4) was based on the sequence at the 3' end of *cps19fB* reducing the size of the potential LR-PCR products by up to 5 kb (the distance between the 3' end of *dexB* and the 3' end of *cps19fB*, Fig. 4.2). Such smaller LR-PCR products are sufficient to characterise other *cps* loci as they would contain all of the type-specific portions of the loci, missing only the highly conserved *cpsA* and *B* genes. The other primer (J36, described in Table 2.4) was based on sequence at the 5' end of *aliA* (Fig. 4.2) as there are no genes which are conserved in all serotypes at the 3' end of the *cps* locus. The primers were not based on the apparently conserved 3' intergenic region because of the theoretically greater possibility of minor sequence divergence in non-coding DNA. LR-PCR amplification, using these primers was undertaken as described in section 2.9.2.

All pneumococcal serotypes tested in Table 4.1, except type 3 (the 5' end of the *aliA* gene has been deleted in all type 3 strains tested [Caimano *et al.*, 1998]), and a type 20 isolate, were used as templates for LR-PCR. The PCR reactions were analysed by electrophoresis on a 1% agarose gel. PCR products were obtained from at least one pneumococcal isolate of types 2, 4, 6A, 6B, 8, 9N, 14, 18C, 19F, 19A, 19B, 20 and 23F but



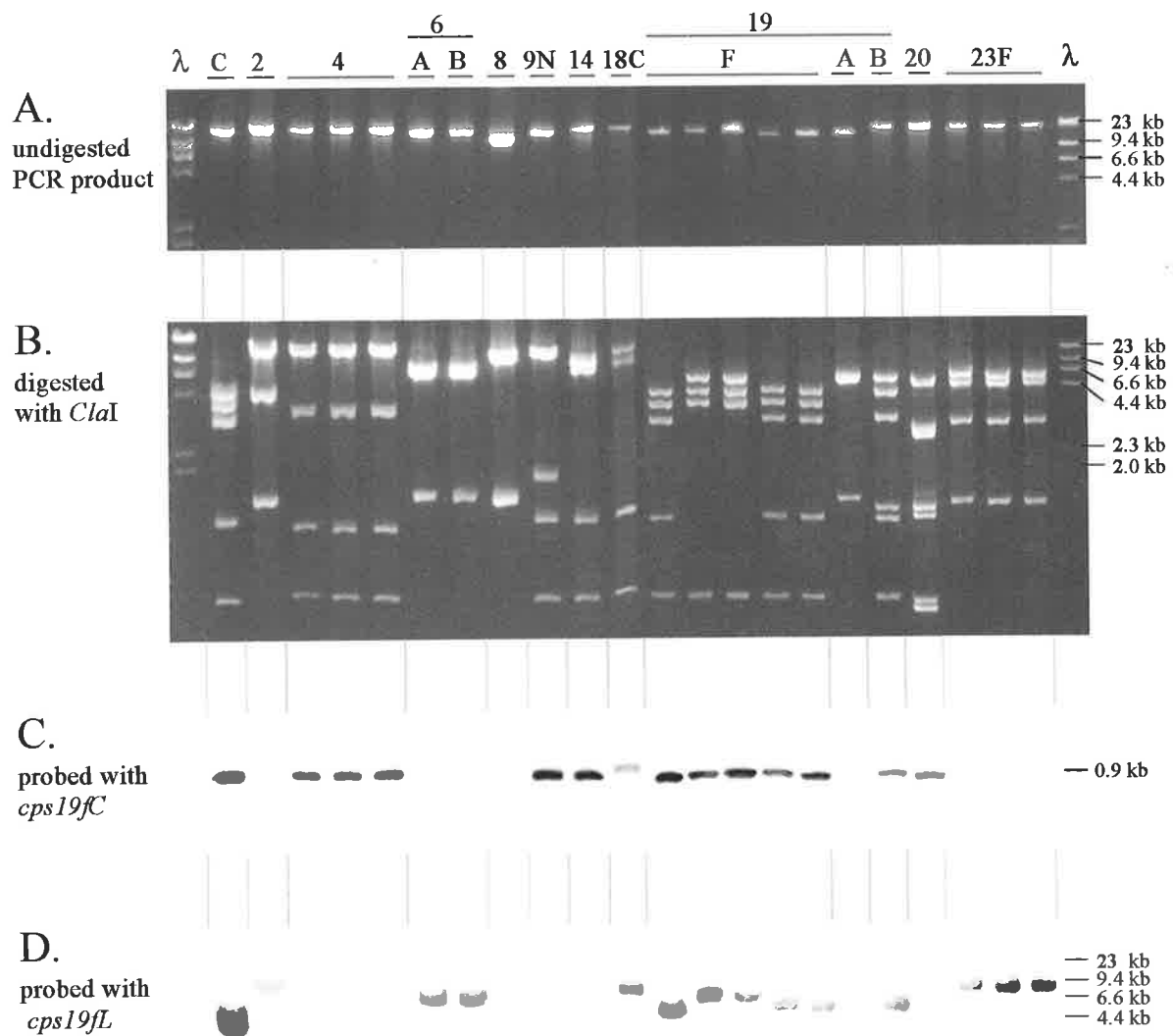
**Fig. 4.2.** The *dexB-aliA* region of the type 19F chromosome. The conserved regions including those within the *cps19f* locus are indicated in red. The position of the two primers used for long-range PCR, J39 and J36 (Table 2.4), are indicated with arrows.

not from types 7F, 9V, 12, 16, 17, 19C, 22 and 24. Analysis of the DNA fragments reveals that the LR-PCR products ranged in size from 15-20 kb as shown in Fig. 4.3A. The PCR products were also digested with *ClaI* and electrophoresed on a 1% agarose gel (Fig. 4.3B). Identical restriction patterns were obtained for three different isolates of serotypes 4 and 23F. However, a restriction site polymorphism was observed in two of the five LR-PCR products from different type 19F strains (Fig. 4.3B).

### 4.2.3 Southern hybridisation analysis of LR-PCR products

In order to confirm that they contained *cps*-related sequences, the LR-PCR products from the various *S. pneumoniae* serotypes were restricted with *ClaI*, and subjected to Southern hybridisation analysis using probes specific for *cps19fC* and *cps19fL* (Fig. 4.3C and D).

The *cps19fC* probe hybridised at high stringency with a 0.9-kb DNA fragment in types 4, 9N, 14, 18C, 19F, 19B and 20. Both the hybridisation pattern and the size of the DNA fragment which hybridised with the *cps19fC* probe are consistent with the Southern hybridisation data obtained when probing *ClaI*-restricted chromosomal DNA with the *cps19fC* probe (Table 4.2). Chromosomal hybridisation data were not available for the



**Fig. 4.3. LR-PCR products.** PCR product, undigested (A), or digested with endonuclease *ClaI* (B), was electrophoresed on a 1% agarose gel in the presence of ethidium bromide. *ClaI* restricted PCR product was subjected to Southern hybridisation analysis using DIG-labelled probes specific for *cps19fC* (C) or *cps19fL* (D). The molecular size standards are shown on the right-hand side of the figure and correspond to *HindIII*-digested  $\lambda$  phage DNA.

type 20 isolate, from which the LR-PCR product was amplified, as it had not been included in the previous Southern hybridisation study (Table 4.1).

The *cps19fL* probe hybridised with DNA fragments ranging in size from 6 to 10 kb in the *ClaI*-restricted PCR products from types 2, 6A, 6B, 18C, 19F, 19B and 23F. Hybridisation was consistent with that obtained from Southern hybridisation with *ClaI*-restricted chromosomal DNA from these isolates, although, the sizes of the restriction fragments differ (Table 4.2). The size of this *ClaI* fragment is affected because there is no *ClaI* site between *cps19fL* and the end of the PCR product in type 19F.

As expected, neither the *cps19fC* nor the *cps19fL* probes hybridised at high stringency with the LR-PCR products obtained from types 8 and 19A. Previous Southern hybridisation data (Table 4.1) indicated that *S. pneumoniae* types 8 and 19A do not contain DNA sequences highly homologous to *cps19fC* or *cps19fL*.

**Table 4.2. Sizes of *ClaI*-restricted fragments which hybridise to *cps19fC* and *cps19fL* probes.**

Serotype	Sizes <sup>a</sup> of <i>ClaI</i> -restricted fragments which hybridise to:			
	<i>cps19fC</i> <sup>b</sup>		<i>cps19fL</i> <sup>b</sup>	
	LR-PCR product	Chrom. DNA	LR-PCR product	Chrom. DNA
2	-	-	10 kb	15 kb
4	0.9 kb	0.9 kb	-	-
6A	-	-	8 kb	10 kb
6B	-	-	8 kb	10 kb
8	-	-	-	-
9N	0.9 kb	0.9 kb	-	-
14	0.9 kb	0.9 kb	-	-
18C	0.9 kb	0.9 kb	10 kb	15 kb
19F	0.9 kb	0.9 kb	6-9 kb	10-12 kb
19A	-	-	-	-
19B	0.9 kb	0.9 kb	6 kb	10 kb
20	0.9 kb	NA <sup>c</sup>	-	NA
23F	-	-	10 kb	10 kb

a. Size is approximate.

b. DIG labelled probes specific for *cps19fC* and *cps19fL* as described in Table 4.1.

c. Not analysed.

#### 4.2.4 Partial DNA sequencing of the LR-PCR products

As additional confirmation that the LR-PCR products contain *cps* related genes, they were subjected to one round of sequence analysis using dye-terminator sequencing reactions with the J39 primer (located at the 5' end of the PCR product). This provided sequence data for the 3' end of *cpsB* and the 5' portion of *cpsC* for all the LR-PCR products except type 18C. No sequence data were obtained using this latter template, presumably due to the low yield of the PCR product obtained. Analysis of the various sequences and those published for types 1, 3 and 14 (Muñoz *et al.*, 1997; Arrecubieta *et al.*, 1995; Kolkman *et al.*, 1996) showed that there are two distinct *cpsC* genes in these loci, designated class I and class II. Types 1, 3, 9N, 14, 19F, 19B and 20 have class I *cpsC* genes which exhibit >95% identity to *cps19fC*, whereas types 2, 6A, 6B, 8, 19A and 23F have class II *cpsC* genes which exhibit 72-74% identity to *cps19fC*, but >95% identity to each other (Fig. 4.4). Table 4.3 shows the sequence similarity between the various *cpsC* genes. The sequences obtained from the LR-PCR products also included the last 75 nucleotides of *cpsB*; this region can also be separated into the same two classes as described above (Fig. 4.4).

An interesting exception is found in type 4, the *cps4C* gene of which is a hybrid consisting of a class II 5' region and a class I 3' region, with a distinct cross-over point between nucleotides 345 and 354 of the type 4 sequence (Fig. 4.4). Comparison of the TIGR type 4 sequence data (section 8.3.4) with the 19F *cps* sequence showed another point of divergence within the *cpsB* gene. The comparison of the sequences of *cps19fB* and *C*, and *cps4B* and *C*, is shown in Fig. 4.5. A region of 852 nucleotides, including most of *cpsB* and part of *cpsC*, shares approximately 74% similarity between *cps19f* and *cps4*, whereas flanking regions exhibit greater than 95% identity. This may have arisen as a consequence of recombination between a class I *cps* locus and a DNA fragment

```

...cpsB →                                     *                                     cpsC →
cps19f AAATATGGAGCGAAAAAGCAAAGAACTTTTGTAGATAATCCAGAAAATATAATGGATCAATTAATTTAGGAGAAAATGAAGGAACAAAAC
cap3   .....A.....G.....
cap1   .....G.....
cps9n  .....G.....T.....
cps19b .....G.....
cps20  .....G.....T.....G.....
cps14  (T).....C.....G.....
cps4   C.....AAGCG..G..TC.G.....A...C...TC.A...G...C...TG...A...T
cps8   C.....AAGCG..G..TC.G.....A...C...TC...G...C...TG...A...
cps2   C.....AAGCG..G..TC.G.....A...C...TC...G...C...TG...A...
cps23f C.....AAGCG..G..TC.G.....A...C...TC...G...C...TG...A...
cps6b  C.....AAGCG..G..TC.G.....A...C...TC...G...C...TG...A...
cps6a  C.....AAGCG..G..TC.G.....A...C...TC...G...C...TG...A...
cps19a C.....AAGCG..G..TC.G.....A...C...TC...G...C...TG...A...

cps19f TTTGGAAATCGATGTATTGCAACTATTACAGAGCTTTATGGAAGAAAGTTGGTCATTTTATTAGTGGCAATTATAACTTCTTCAGTTGCTTTGCCTAC
cap3   .....
cap1   .....
cps9n  .....G.....
cps19b .....T.....
cps20  .....A.....
cps14  .....C.....A.....
cps4   GA.A.....T...T...G.T.A.AGC..G...C.C.C.AA.G...A...C.G.G..AGG.G.G.GG...A.T
cps8   GA.A.....T...T...T.A.A.C.G...C.C.C.AA.G...C.G.G..AGG.G.G.GG...A.T
cps2   GA.A.....T...T...T.A.A.C.G...C.C.C.AA.G...C.G.G..AGG.G.G.GG...A.T
cps23f GA.A.....T...T...T.A.A.C.G...C.C.C.AA.G...A...C.G.G..AGG.A.G.GG...A.T
cps6b  GA.A.....T...T...T.A.A.C.G...C.C.C.AA.G...C.G.G..AAG.G.G.GG...A.T
cps6a  GA.A.....T...T...T.A.A.C.G...C.C.C.AA.G...C.G.G..AAG.GTG.GG...A.T
cps19a GA.A.....T...T...C.T.A.A.C.G...CAC..C.AA.A...C.G.G..AGGG.G.GA...A.T

cps19f AGTACTTTGTATCAAACTGAGTTTACTAGTATGACTCGGATTTATGTAGTTAACCGTGATCAGGAGAGAAGTCTGGTTTAAACCAATCAAGACTTGC
cap3   .....C.....A.....
cap1   .....C.....A.....C.....
cps9n  .....C.....A...A.....
cps19b .....C.....A...A.....
cps20  .....C.....T.....G.....
cps14  .....C.....
cps4   C.....A..G.T.G.A.A.A..G...CC..G.A...C...G..T.CA..A...C..C.G..G.G.A...G.T
cps8   C.....A..G.T.G.A.A.A..G...CC..G.A...C...G..T.CA..A...C..G..GC.G.A...G.T
cps2   C.....A..G.T.G.A.A.A..G...CC..G.A...C...G..T.CA..A...C...G..GC.G.A...G.T
cps23f C.....A..G.T.G.A.A.A..G...CC..G.A...C...G..T.CA..A...C..C.G..G.G.A...G.T
cps6b  C.....A..G.T.G.A.A.A..G...CC..G.A...C...G..T.CA..A...C..C.G..GC.G.A...G.T
cps6a  C.....A..G.T.G.A.A.A..G...CC..G.A...C...G..T.CA..A...C..C.G..GC.G.A...G.T
cps19a C.....A..G.T.G.A.A.A..A..C.CC..G.T...C...C...A...A...T..C.G..AC.G..G...G.

↓                                     ↓
cps19f AGGCAGGATCATCTTGGTTAAAGACTATCGTGAAATTATCCTATCGCAGGATGTTTGGAGGAAGTTGTTTCTGATTGAACTAGATTTGACGCCAA
cap3   .....A.....A.....
cap1   .....G..A.....G.....A.....T.....
cps9n  .....A.....A.....
cps19b .....A.T.ATC...A...C...G...T...A...GA.GA...G...C...
cps20  .....T...CG...T...
cps14  .....A.....T.....
cps4   A.T.ATC...A...C...G...T...CA...AA...A.CGA.AA...GT.G..CA..C.AG...
cps8   T.ATC...A...C...G...T...CA...AA...A.CGA.AA...GT.G..CA..C.AG...
cps2   T.ATC...A...C...G...T...CA...AA...A.CGA.AA...GT.G..CA..C.AG...
cps23f A.T.ATC...A...C...G...T...CA...AA...A.CGA.AA...GT.G..A..C.AG...
cps6b  A.T.ATC...A...C...G...T...CA...AA...A.CGA.AA...GT.G..A..C.AG...
cps6a  A.T.ATC...A...C...G...T...CA...AA...A.CGA.AA...GT.G..A..C.AG...
cps19a A.T.ATC...A...C..C...T...A...AA.G..A.CGA.AA...GT.G..A..C.AG...

cps19f AGATTTGGCTAATAAAATTAAGTAACAGTACCAGTTGATACCCGATTTGTCTCTGTTTCAGTTAGTGTGATCGAGTTCTCTGAAGAGGCAAGCCGTATCGCT
cap3   .....G.....A.....
cap1   .....G.....A.....
cps9n  .....G.....A.....
cps19b .....G.....A.....G.....
cps20  .....G.....
cps14  .....G.....A.....
cps4   G.....G.....
cps8   ACG..A...GC...G..C...G..T...AC..C..T...C...AA..C..T..C.AG...AA.CAG..A..G..A..C..T
cps2   ACG..A...GC...G..C...G..T...AC..C..T...C...AA..C..T..C.AG...AA.CAG..A..G..A..C..T
cps23f ACG..A..C.GC...G..C...GG.T.G...C..C..T...C...AA..C..T..C.AG...AA.CAG..A..G..A..C..T
cps6b  ACG..A..C.GC...G..C...GG.T...C..C..T...C...AA..C..T..C.AG...AA.CAG..A..G..A..C..T
cps6a  ACG..A..C.GC...G..C...GG.T...C..C..T...C...AA..C..T..C.AG...AA.CAG..A..G..A..C..T
cps19a ACG..AA...GC...G..GC...G..T...CC..C..T...C...AA..C..T..C.AG...AA.CAA..A..G..A..C..T..C..T

```

**Fig. 4.4. Comparison of class I and class II *cps* sequences.** The 500 nucleotides of sequence shown (100 nucleotides per line) corresponds to nucleotides 2,273-2,772 of the *cps19f* sequence. The following sequences are available under the GenBank accession numbers as indicated: *cps19f*, U09239; *cap3*, Z47210; *cap1*, Z83335; *cps14*, X85785; *cps23f*, AF030373; *cps19a*, AF094575. The *cps4* sequence is available from the TIGR Microbial database. The stop codon of the *cpsB* gene is indicated with an asterisk. The start codon of *cpsC* is underlined. "(T)" denotes an extra nucleotide, and "-" denotes the absence of a nucleotide in the *cps14* DNA sequence. The arrows indicate the region where the cross-over between class I and class II sequences has occurred in *cps4*.

(approximately 852 nucleotides long) from a class II *cps* locus, resulting in a mosaic *cpsB-C* region. Analysis of the available type 23F sequence data (Coffey *et al.*, 1998a) indicated that the *cps23f* locus (which has a class II *cpsC* gene) also diverges from *cps19f* (which has a class I *cpsC* gene) within the *cpsB* gene, but 98 nucleotides further downstream from the point of divergence for *cps4*. This suggests that the point of sequence divergence from class I to class II within *cpsB* may vary between different serotypes.

**Table 4.3. Similarity between *cpsC* sequences from various pneumococcal serotypes.**

	% Identity <sup>a</sup>													
	<i>cps19f</i>	<i>cps3</i>	<i>cap1</i>	<i>cps9n</i>	<i>cps19b</i>	<i>cps20</i>	<i>cps14</i>	<i>cps4</i>	<i>cps8</i>	<i>cps2</i>	<i>cps23f</i>	<i>cps6b</i>	<i>cps6a</i>	<i>cps19a</i>
<i>cps19f</i>	100	98.5	97.9	97.9	96.7	96.5	97.9	82.4	73.9	73.7	73.3	73.1	72.9	72.2
<i>cps3</i>		100	98.6	99	97.9	96.9	98.3	83	74.3	74.1	73.7	73.5	73.3	72.7
<i>cap1</i>			100	98.1	97.1	95.9	97.3	82.6	73.7	73.5	73.7	73.5	73.3	72.7
<i>cps9n</i>				100	97.7	96.7	97.7	82.4	74.1	73.9	73.5	73.3	73.1	72.5
<i>cps19b</i>					100	95.6	96.5	84.5	75.4	75.2	75.2	75	74.9	73.9
<i>cps20</i>						100	95.8	81.2	73.7	73.5	73.1	72.9	72.7	72.5
<i>cps14</i>							100	82.6	73.9	73.7	73.4	73.2	73	72.4
<i>cps4</i>								100	88.6	88.8	89	88.8	88.6	85.3
<i>cps8</i>									100	99.8	97.7	98.3	98.1	94
<i>cps2</i>										100	97.9	98.5	98.3	94.2
<i>cps23f</i>											100	98.6	98.5	93.6
<i>cps6b</i>												100	99.8	94.2
<i>cps6a</i>													100	94
<i>cps19a</i>														100

a. Percentage of DNA identity over 517 nucleotides, equivalent to nucleotides 2,273-2,789 of the *cps19f* sequence (GenBank accession no U09239), as determined in DNASIS.

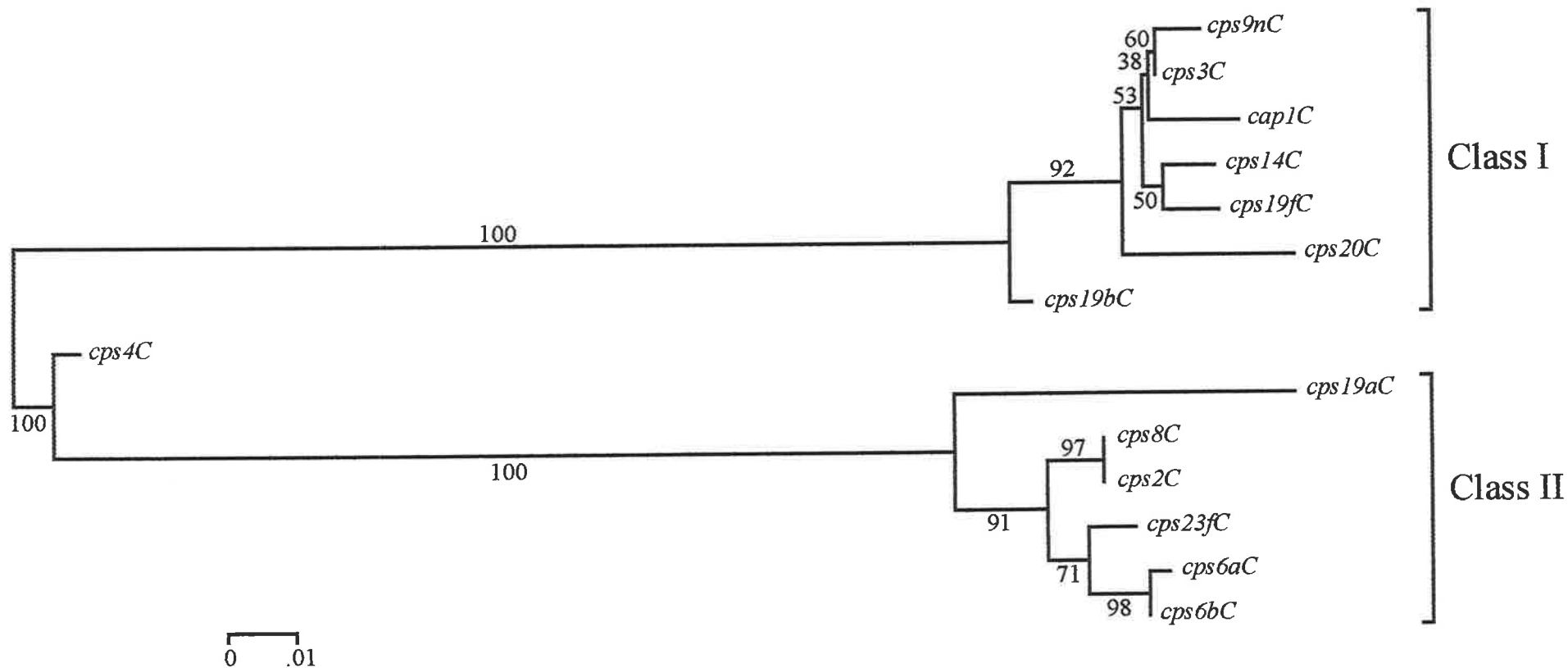


### 4.2.5 Phylogenetic analysis of *cpsC* sequences

To further investigate the differences between the two classes of *cpsC* sequences, their phylogenetic relationship was investigated. An alignment of the partial *cpsC* sequences was generated using CLUSTAL W (Thompson *et al.*, 1994; data not shown) and this alignment was used to generate a phylogenetic tree using the Neighbour-Joining method, and the distance measure of Tamura and Nei (1993), as implemented in the program MEGA (Kumar *et al.*, 1994). The tree in **Fig. 4.6** shows two highly significant clusters of *cpsC* sequences, and confirms the observations initially made on the basis of sequence homology that the *cpsC* genes are divided into two classes. The *cps4C* sequence forms a third cluster; as described above, this gene is a hybrid of the two *cpsC* classes and has a recombination cross-over point near nucleotide 345 (as shown in **Fig. 4.4**) within the *cpsC* gene. The *cps19bC* gene is also separated from the other class I *cpsC* sequences; *cps19bC* also appears to have a mosaic structure with a small region of class II sequence (nucleotides 409-444 in **Fig. 4.4**) which is presumably the result of a recombination event.

## 4.3 Conclusions

The Southern hybridisation data described in this chapter indicate that the DNA flanking the *cps* locus is highly conserved among pneumococci of diverse serotypes. Using LR-PCR with primers specific for common regions, it was possible to amplify the major portion of the capsule loci from several different pneumococcal serotypes. Southern hybridisation analysis confirmed that the large PCR products obtained did indeed contain *cps*-related DNA. Moreover, direct sequencing of the PCR products identified two



**Fig. 4.6. Phylogenetic tree of *cpsC* sequences.** The *cpsC* gene sequences were aligned in CLUSTAL W (Thompson *et al.*, 1994), and the phylogenetic tree generated using MEGA (Kumar *et al.*, 1994), as described in the text. The numbers associated with the branches are Bootstrapping Confidence Limits, resulting from 500 replications, as defined in MEGA. The scale represents the number of nucleotide substitutions per site.

apparent classes of the *cpsC* gene. This was confirmed by phylogenetic analysis of the sequence data. The presence of the *cpsC* gene in all *cps* loci examined is consistent with the important proposed role of CpsC as a chain length regulator in pneumococcal CPS production (section 1.9.3). At this stage, it is not possible to determine whether the differences between class I and class II *cpsC* genes is functionally significant. Translation of the genes indicates a similar degree of amino acid sequence divergence between class I and class II CpsC proteins (approximately 70% identity). Interestingly, even small differences between the functionally homologous Rol (Wzz) proteins of *Shigella* species have previously been shown to impact on the modal chain length of the LPS O-antigen (Klee *et al.*, 1997).

The type 4 and 19B *cpsC* sequences both show evidence of recombination within the *cps* loci. Two recent studies have demonstrated that natural recombination events involving exchange of entire *cps* loci (or major portions thereof) have resulted in switching of capsule type (e.g. from 23F to 19F) by multiply drug-resistant pneumococcal clones on numerous occasions (Coffey *et al.*, 1998a; Nesin *et al.*, 1998). The current study indicates that recombination events involving small fragments within pneumococcal *cps* loci may also be common in nature, and may represent a mechanism whereby additional serotype diversity is generated.

A comparison between the *cps19f* and *cps4* sequences revealed that the distinction between class I and class II sequences can be extended to include the 3' region of *cpsB*, where the J39 primer sequence is located (Fig. 4.5). The 3' region of *cps19fB* exhibits only 75% similarity with the 3' region of *cps4B*. It is possible that the *cps* loci from some serotypes may be even more variable in this region, such that the J39 primer is non-functional, and could explain why a LR-PCR product was not obtained in some of the serotypes tested. This suggests that the efficiency of the LR-PCR could be significantly

improved by using a primer sequence from the 5' region of *cpsI9fB*, which appears to be more conserved, in lieu of J39 (Fig. 4.5).

Although LR-PCR successfully obtained PCR products from 13 serotypes, it was unsuccessful from 8 others. This may be simply due to the size of *cpsB-aliA* regions from these serotypes. Indeed, the CPS of the serotypes for which LR-PCR products were not obtained all have large repeat units containing 5 to 7 sugars. No attempts were made to modify the LR-PCR amplification conditions used, but it is likely that such modifications could improve the amplification efficiency of the larger *cps* loci. The possibility of gene rearrangements within the *cps* loci of these strains also cannot be excluded.

The LR-PCR products that have so far been obtained provide a ready source of DNA to study the as yet uncharacterised *cps* genes required for CPS production from serotypes 6A, 6B, 8, 9N, 18C and 20. Additionally, with improvements, the LR-PCR protocol has the potential to amplify the *cps* loci of other serotypes. This would extend its utility for the epidemiological and clonal analysis of clinical isolates.

## Chapter 5

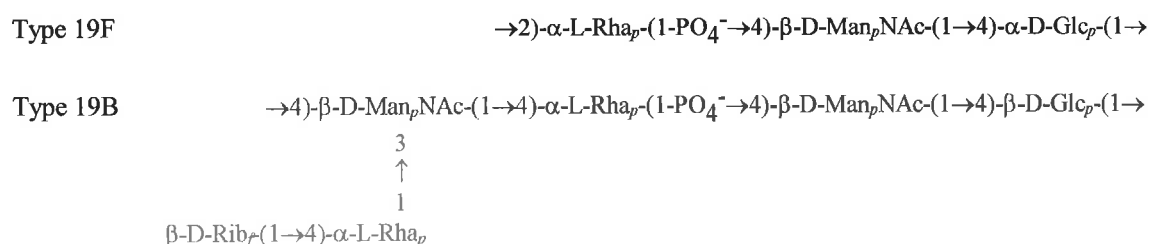
# CHARACTERISATION OF THE *S. PNEUMONIAE* TYPE 19B CAPSULE LOCUS

### 5.1 Introduction

The structure of the type 19B CPS is more complex than that of type 19F. The oligosaccharide repeat unit of type 19B CPS has a fourth sugar (ManNAc) in its backbone, with a disaccharide side-chain, as shown in **Fig. 5.1**. Thus the type 19B capsule locus (*cps19b*) would be predicted to contain extra genes required for biosynthesis of the more complicated repeat unit; one would also predict that a distinct polysaccharide repeat unit transporter and polysaccharide polymerase would be required.

In chapter 4, Southern hybridisation was used to demonstrate that all but 2 of the 15 genes in the *S. pneumoniae* type 19F capsule locus hybridised to the DNA of a type 19B *S. pneumoniae* strain. These 2 genes, *cps19fI* and *cps19fJ*, encode the polysaccharide polymerase and polysaccharide repeat unit transporter, respectively, and are located together near the middle of the *cps19f* locus. DNA from the 2 flanking genes, *cps19fH* and *cps19fK*, hybridised only weakly to 19B DNA, whereas all other *cps19f* genes hybridised strongly at high stringency. These data suggested that the type 19B capsule locus does contain different polysaccharide repeat unit transporter and polysaccharide

polymerase genes. It was not obvious, however, where the postulated additional genes required for type 19B CPS biosynthesis were located and hence characterisation of the *cps19b* locus was undertaken.



**Fig. 5.1. Biological repeat units of pneumococcal type 19F and type 19B capsular polysaccharide.** The order of the sugars in the repeat units have been altered compared to the published chemical structures for 19F (Ohno *et al.*, 1980) and 19B (Beynon *et al.*, 1991) as it has been determined that glucose is the first sugar in the biological repeat unit. D-Glu<sub>p</sub>, glucose; D-Man<sub>p</sub>NAc, *N*-acetyl mannosamine; L-Rha<sub>p</sub>, rhamnose; D-Rib<sub>p</sub>, ribose; PO<sub>4</sub><sup>-</sup>, phosphate;.

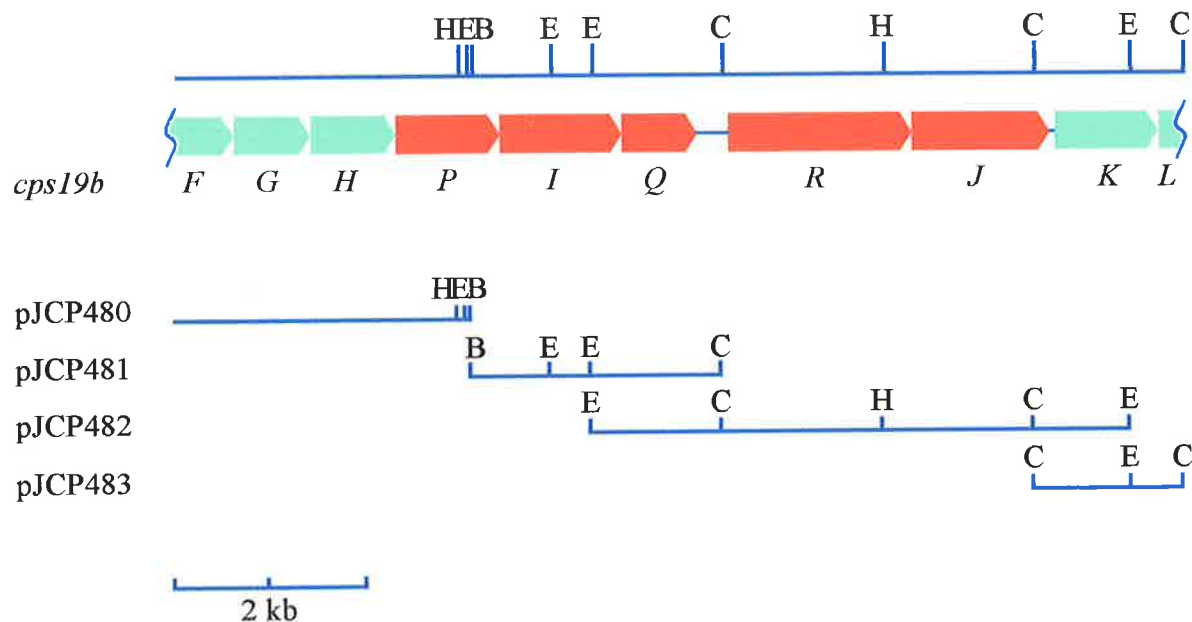
## 5.2 Results

### 5.2.1 Isolation of the type 19B-specific *cps* genes

In order to isolate the type-specific portion of the *cps19b* locus, primers J5 which is homologous to *cps19fF* and J27 which is complementary to *cps19fL* (Table 2.4) were used to amplify the intervening region of the 19B chromosome using LR-PCR with an extension time of 12 min per cycle (see section 2.9.2). The LR-PCR product (from the 3' end of *cps19bB* to *aliA*) which was described in chapter 4 was not used because it contains a large amount of DNA (*cps19bC-cps19bG* and *cps19bL-aliA*) which was predicted (on the basis of hybridisation data) to be greater than 95% identical to *cps19f*. The J5/J27-directed PCR product was approximately 10.5 kb in size. This is 4 kb larger than the equivalent region of *cps19f*, indicating that extra ORFs may be present in this region of the *cps19b* locus. A

map of the 10.5-kb PCR product was generated, using the restriction enzymes *Bam*HI, *Cla*I, *Hind*III and *Eco*RI, and appropriate fragments of the PCR product were then cloned into pBluescript KS+, generating four recombinant plasmids with inserts as shown in **Fig. 5.2.**

## 5.2.



**Fig. 5.2. Physical map of part of the *S. pneumoniae* type 19B capsule locus (*cps19b*).** Green arrows represent ORFs also present in *cps19f* and red arrows represent type 19B-specific ORFs. Gene designations are indicated below the map; *cps19bF-L* are abbreviated to *F-L*, respectively. Restriction sites are as follows; B, *Bam*HI; C, *Cla*I; E, *Eco*RI; H, *Hind*III. The regions of DNA subcloned into various recombinant plasmids are shown below the map.

Both strands of the pneumococcal DNA inserts of each of the above plasmids (or nested deletion derivatives thereof) were subjected to sequence analysis (section 2.8) in order to compile the sequence of this portion of the *cps19b* locus, as shown in **Appendix II.** The sequence across the *Bam*HI site at the junction of the inserts of pJCP480 and pJCP481 was obtained by dye terminator sequencing (section 2.8.3) of the PCR product with the primer J44 (**Table 2.4**). Examination of the compiled sequence revealed, as



ORFs, between *cps19bH* and *cps19bK*, which have been designated *cps19bP*, *I*, *Q*, *R* and *J*. Each ORF is preceded by a ribosome binding site and the majority are very closely linked. The only potentially significant intergenic gap of 204 nucleotides occurs between *cps19bQ* and *cps19bR*. However, no potential stemmed-loop structures or obvious promoter sequences were found in this region. As predicted, the sequence at the 3' end of the 10.5 kb *cps19b* PCR product again shows similarity to the *cps19f* sequence, starting in the vicinity of nucleotide 9,233 (**Fig. 5.3B**); this is immediately before the start of the *cps19bK* gene, which has 93% identity to *cps19fK*.

### 5.2.2 Characterisation of the *cps19b* genes

The locations and several properties of each of the type 19B-specific ORFs, *cps19bP*, *I*, *Q*, *R* and *J*, are summarised in **Table 5.1**. Significant similarities with other known proteins, revealed by comparison with sequence databases, are described below.

**Table 5.1. Summary of ORFs *cps19bP-J*.**

ORF	Location in sequence <sup>a</sup>	Predicted MW	No. amino acids	Hydrophobicity index <sup>b</sup>	Predicted pI	% G+C content <sup>c</sup>
<i>cps19bP</i>	2,350-3,432	43,334	361	-0.41	8.69	29.5
<i>cps19bI</i>	3,451-4,695	48,667	414	0.77	9.73	27.2
<i>cps19bQ</i>	4,703-5,605	34,876	300	-0.26	8.20	29.7
<i>cps19bR</i>	5,809-7,746	76,348	645	-0.30	8.65	27.2
<i>cps19bJ</i>	7,736-9,181	53,851	481	0.91	9.89	29.5

a. Nucleotide numbers correspond to *cps19b* sequence as shown in Appendix II.

b. According to Kyte and Doolittle (1982), as implemented in PROSIS.

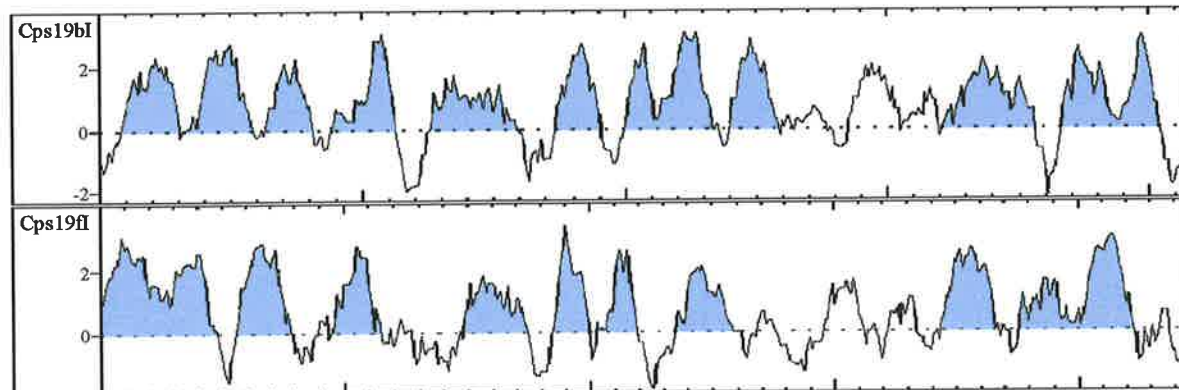
c. Percent guanine plus cytosine (G+C) of coding region.

### 5.2.2.1 *cps19bP*

The *cps19bP* gene encodes a putative 43.3 kDa protein. Database searches with Cps19bP did not reveal any significant similarity to any other proteins.

### 5.2.2.2 *cps19bI*

The *cps19bI* gene encodes a putative 48.7 kDa protein. Database searches with Cps19bI also failed to identify any significant similarity to any other proteins. However, Cps19bI is a very hydrophobic protein and the hydropathy plot (Fig. 5.4) illustrates the marked similarity between Cps19bI and Cps19fI, each having 10-12 hydrophobic, potentially membrane-spanning domains. This is a typical hydropathy profile for Rfc-like proteins and suggests that Cps19bI, like Cps19fI, may be a polysaccharide polymerase.



**Fig. 5.4. Hydropathy plots of Cps19bI and Cps19fI.** The hydropathy plots were generated by the method of Kyte and Doolittle (1982) and aligned using PROFILEGRAPH (Hofmann and Stöffel, 1989). Positive numbers on the Y-axis indicate hydrophobic regions. The position of every 10th amino acid is marked on each X-axis.

### 5.2.2.3 *cps19bQ*

The *cps19bQ* gene encodes a putative 34.9 kDa protein, which has similarity with rhamnosyl transferases from *Leptospira interrogans*, *Shigella dysenteriae*, and *S. flexneri*,

and a 6-deoxyaltrosyl transferase from *Y. enterocolitica* as shown in **Table 5.2**. The alignment of Cps19bQ with the other proteins shows several regions of similarity, including a motif previously identified in rhamnosyl and other 6-deoxyhexosyl transferases (section 3.2.3.2; Morona *et al.*, 1995), as shown in **Fig. 5.5**. This motif contains conserved aspartate residues and is reminiscent of the catalytic sites identified in RfbA<sub>O:54</sub> of *S. enterica* serovar borreze (Keenleyside and Whitfield, 1996). A second conserved motif was also identified (**Fig. 5.5**, yellow shaded amino acids), suggesting that these proteins form a closely related sub-group of this type of transferase. Cps19bQ is the second putative rhamnosyl transferase in the *cps19b* locus. Cps19bH, which has 91.1% identity with Cps19fH, is also proposed to be a rhamnosyl transferase, adding Rha to the repeat unit backbone (section 3.2.3.2). Thus Cps19bQ is predicted to add Rha (the first sugar of the disaccharide side-chain [see **Fig. 5.1**]) to the distal ManNAc.

**Table 5.2. Similarity of Cps19bQ to other proteins.**

	% Identity <sup>a</sup>				
	Cps19bQ <sup>b</sup>	SdRfbQ <sup>c</sup>	LiRfbF <sup>d</sup>	YeRfbC <sup>e</sup>	SfRfbF <sup>f</sup>
Cps19bQ	100	24	22.9	22.1	20.2
		[267]	[245]	[213]	[257]
SdRfbQ		100	31.3	30	30.2
			[297]	[300]	[298]
LiRfbF			100	28.9	30.6
				[308]	[284]
YeRfbC				100	29.5
					[298]
SfRfbF					100

a. Percentage of identical amino acids determined with FASTA as implemented in PROSIS. Numbers in parentheses indicate the number of amino acids over which the % identity occurs.

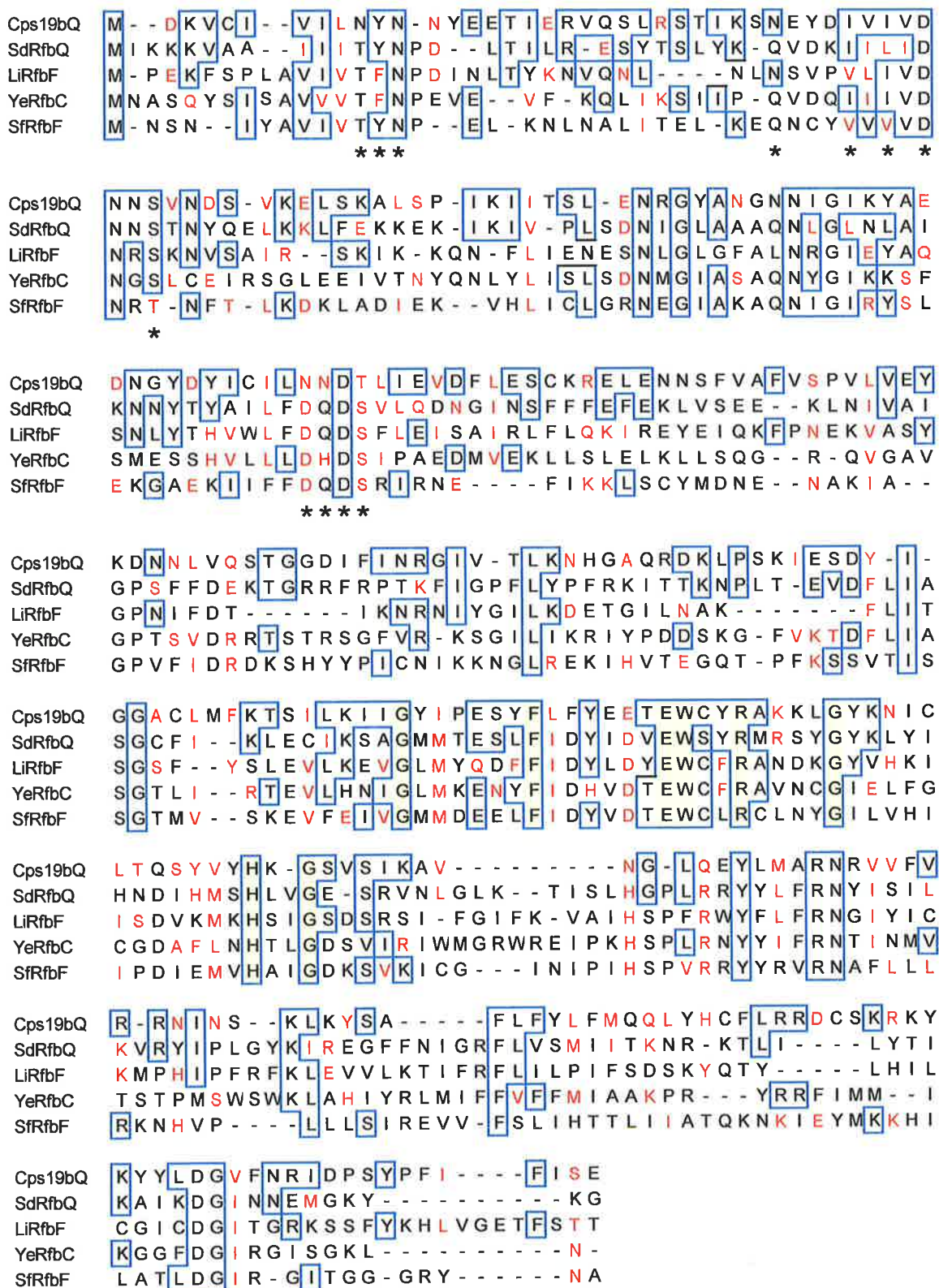
b. *S. pneumoniae* Cps19bQ.

c. *S. dysenteriae* RfbQ (Klena and Schnaitman, 1993).

d. *L. interrogans* RfbF (Mitchison *et al.*, 1997).

e. *Y. enterocolitica* RfbC (Zhang *et al.*, 1993).

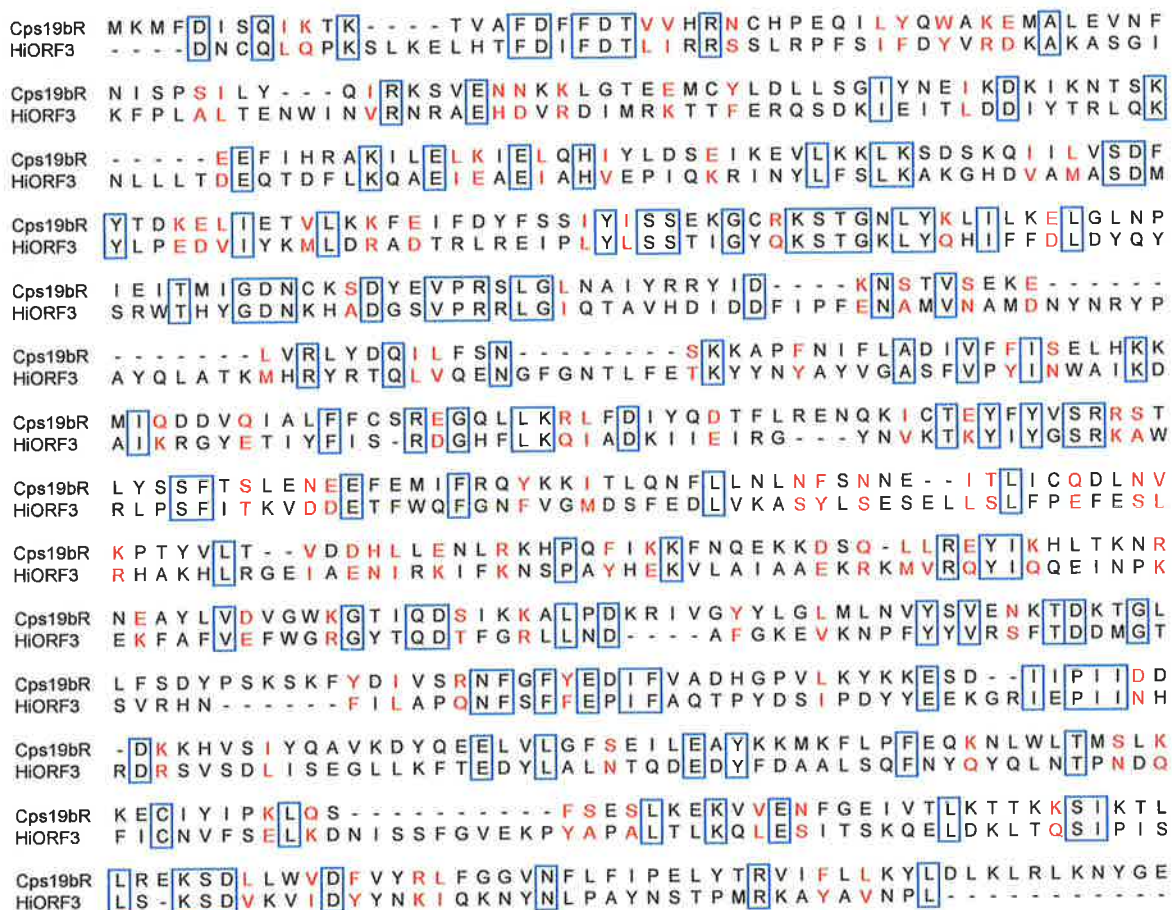
f. *S. flexneri* RfbF (Macpherson *et al.*, 1994).



**Fig. 5.5. Alignment of Cps19bQ with other proteins.** Alignment of Cps19bQ with *S. dysenteriae* RfbQ, (SdRfbQ) (Klena and Schnaitman, 1993), *L. interrogans* RfbF, (LiRfbF) (Mitchison *et al.*, 1997), *Y. enterocolitica* RfbC, (YeRfbC) (Zhang *et al.*, 1993) and *S. flexneri* RfbF, (SfRfbF) (Morona *et al.*, 1995) using the default settings of the program CLUSTAL (Higgins and Sharp, 1988) and enhanced by manual adjustment. Residues identical to Cps19bQ are boxed; similar residues are shown in red; - indicates absence of a residue. The asterisks correspond to a motif found in a variety of rhamnosyl and 6-deoxy-hexosyl transferases (section 3.2.3.2; Morona *et al.*, 1995). The yellow shaded region corresponds to an additional motif shared by this sub-group of transferases.

### 5.2.2.4 *cps19bR*

The *cps19bR* gene encodes a putative 76.3 kDa protein. Cps19bR has similarity (21.4% identity, 40.3% similarity) to the central portion (amino acids 366-1071) of hypothetical protein 3 from the capsule locus of *H. influenzae* type b (Van Eldere *et al.*, 1995) as shown in **Fig. 5.6**. Although the function of hypothetical protein 3 is not known, it is located in the serotype-specific Region II of the capsule locus. The only sugar in the *H. influenzae* type b capsule also found in the pneumococcal type 19B CPS is ribose (Rib). Thus it is possible that Cps19bR and protein 3 could function in the synthesis of an activated Rib precursor. However, the nature of the activated Rib precursor that is used in CPS biosynthesis in bacteria is not known (Van Eldere *et al.*, 1995), and the possibility that both proteins are ribosyl transferases cannot be excluded.



**Fig. 5.6.** Alignment of Cps19bR with *H. influenzae* ORF 3 (HiORF3). Alignment was performed using the default settings of the program CLUSTAL (Higgins and Sharp, 1988). Residues identical to Cps19bR are boxed; similar residues are shown in red; - indicates absence of a residue.

### 5.2.2.5 *cps19bJ*

The *cps19bJ* gene encodes a putative 53.9 kDa protein with low level similarity to RfbX proteins from *E. coli*, *S. dysenteriae*, *Y. enterocolitica*, to the CapF protein of *S. aureus* and to Cps19fJ, as shown in **Table 5.3**. The RfbX proteins are known to be involved in export of O-antigen repeat units (Liu *et al.*, 1996; Macpherson *et al.*, 1995). The hydropathy plots for RfbX-like proteins are all very similar, with 10-12 hydrophobic, membrane-spanning domains, and the similarity between Cps19bJ and Cps19fJ is shown in **Fig. 5.7**. Thus, Cps19bJ is likely to be the polysaccharide repeat unit transporter.

**Table 5.3. Similarity of Cps19bJ to other proteins.**

	% Identity <sup>a</sup>						
	Cps19bJ <sup>b</sup>	K12RfbX <sup>c</sup>	YeTrsA <sup>d</sup>	SdRfbX <sup>e</sup>	SaCapF <sup>f</sup>	YeRfbX <sup>g</sup>	Cps19fJ <sup>h</sup>
Cps19bJ	100	23.3	21.6	22.3	21	19	18.2
		[404]	[402]	[394]	[395]	[420]	[406]
K12RfbX		100	28.3	31.4	21.5	16.8	18.1
			[406]	[401]	[395]	[386]	[414]
YeTrsA			100	28.4	20.4	17.5	19.2
				[401]	[401]	[406]	[416]
SdRfbX				100	21.1	16.7	19.4
					[393]	[377]	[402]
SaCapF					100	20	16.4
						[404]	[397]
YeRfbX						100	18.1
							[425]
Cps19fJ							100

a. Percentage of identical amino acids determined with FASTA as implemented in PROSIS. Numbers in parentheses indicate the number of amino acids over which the % identity occurs.

b. *S. pneumoniae* Cps19bJ.

c. *E. coli* K12 RfbX (Stevenson *et al.*, 1994)

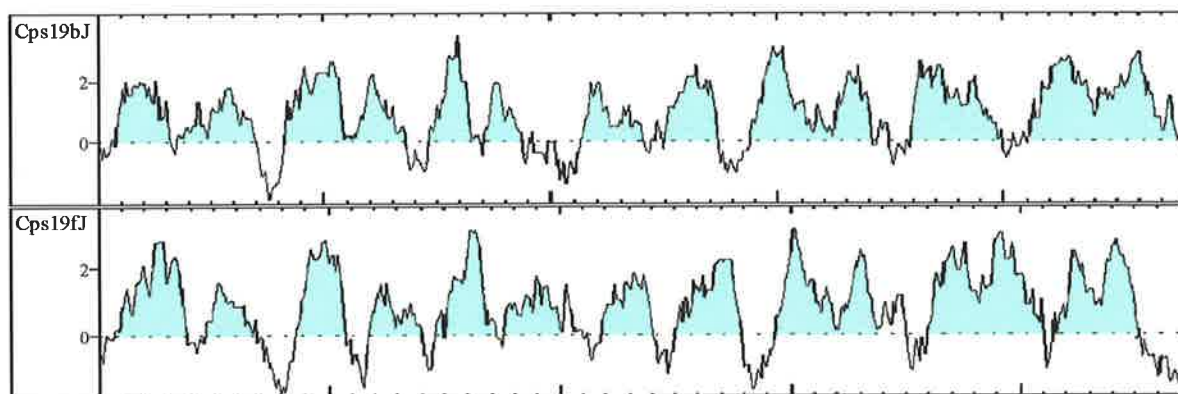
d. *Y. enterocolitica* TrsA (Skurnik *et al.*, 1995)

e. *S. dysenteriae* RfbX (Klena and Schnaitman, 1993).

f. *S. aureus* CapF (Lin *et al.*, 1994).

g. *Y. enterocolitica* RfbX (Zhang *et al.*, 1997).

h. *S. pneumoniae* Cps19fJ (section 3.2.3.4)



**Fig. 5.7. Hydropathy plots of Cps19bJ and Cps19fJ.** Hydropathy plots were generated by the method of Kyte and Doolittle (1982) and aligned using PROFILEGRAPH (Hofmann and Stöffel, 1989). Positive numbers on the Y-axis indicate hydrophobic regions. The position of every 10th amino acid is marked on each X-axis.

### 5.2.3 Serotype specificity of the *cps19b* genes

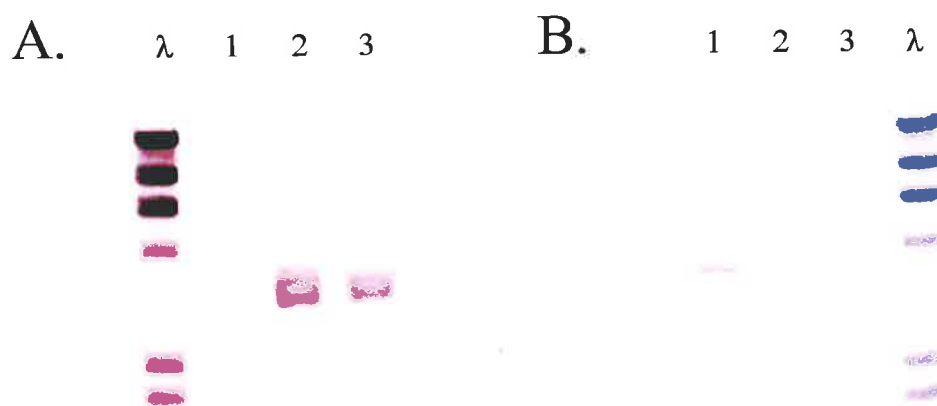
To examine the relationship between *cps19b* and encapsulation loci of other *S. pneumoniae* serotypes, DNA fragments corresponding to the individual *cps19b* genes (described in Table 5.4) were labelled with DIG and used to probe (at high stringency) Southern blots of restricted chromosomal DNA from representative pneumococci belonging to the following types/groups: 2, 3, 4, 6, 7F, 7B, 8, 9, 12, 14, 16, 17, 18, 19F, 19A, 19C, 22, 23 and 24. None of the probes hybridised to DNA from any serotype tested, except to the closely related type 19C (Table 1.3), which has high stringency homologues for all 5 genes (result not shown).

**Table 5.4 DNA fragments used as gene-specific probes**

Probe	Fragment used as probe:
<i>cps19bP</i>	<i>SpeI-BamHI</i> fragment from pJCP480 (nucleotides 2,305-3,188 of <i>cps19b</i> sequence)
<i>cps19bI</i>	<i>PvuII-HindII</i> fragment from pJCP481 (nucleotides 3,412-4,545 of <i>cps19b</i> sequence)
<i>cps19bQ</i>	<i>HindII-ClaI</i> fragment from pJCP481 (nucleotides 4,540-5,718 of <i>cps19b</i> sequence)
<i>cps19bR</i>	<i>ClaI-NsiI</i> fragment from pJCP482 (nucleotides 5,713-7,766 of <i>cps19b</i> sequence)
<i>cps19bJ</i>	<i>NsiI-PvuII</i> fragment from pJCP482 (nucleotides 7,760-9,046 of <i>cps19b</i> sequence)

### 5.2.4 Capsule type switching by transformation

To determine if the genes sequenced were sufficient for type 19B CPS biosynthesis, the 10.5-kb PCR product was transformed into an unencapsulated, Ery-resistant, derivative of Rx1-19F, designated Rx1-19F-I, in which the *cps19fl* gene had been disrupted by insertion-duplication mutagenesis using pVA891, as described in section 3.2.4. Several smooth transformants were checked for Ery sensitivity, indicating loss of the pVA891 sequence. Southern hybridisation was used to confirm the absence of both pVA891 and the *cps19fl* gene, and the presence of each of the *cps19bP*, *I*, *Q*, *R* and *J*, genes in one transformant, designated Rx1-19B. The hybridisation results, using probes specific for *cps19fl* and *cps19bl*, are shown in Fig. 5.8; other results not shown. The production of a type 19B capsule by Rx1-19B was then confirmed by quellung reaction. This shows that it is possible to alter capsule production from type 19F to type 19B by replacing part of the *cps* locus, and that the region of *cps19b* described in this study determines the 19B serotype.



**Fig. 5.8. Southern hybridisation of Rx1-19B.** *Cla*I-restricted chromosomal DNA of *S. pneumoniae* strains Rx1-19F (1), Rx1-19B (2), and 19B (3) was probed with DIG labelled probes specific for *cps19bl* (A) and *cps19fl* (B). The sizes of the DIG-labelled lambda ( $\lambda$ ) markers are as follows: 23 kb, 9.4 kb, 6.6 kb, 4.4 kb, 2.3 kb and 2.0 kb.

### 5.3 Conclusions

When the presence of the 15 *cps19f* genes in other pneumococcal serotypes was examined, the hybridisation patterns showed blocks of *cps19f* genes which hybridised to at least one other serotype flanking blocks of genes which did not (**Table 3.1**). This suggests that different serotypes may have evolved as a consequence of the replacement of one cluster of genes within the capsule locus with an alternative gene cluster. The work described in this chapter demonstrates that the five additional type-specific genes in the *cps19b* locus are indeed grouped together. Transformation of Rx1-19F-I with the 10.5-kb *cps19b* PCR product yielded a transformant, Rx1-19B, which expresses type 19B capsule. Thus, this cluster of five genes (*cps19bPIQRJ*) is sufficient to encode all of the additional and/or distinct functions required for production of type 19B rather than type 19F CPS.

The chemical structure of the *S. pneumoniae* type 19B capsule is considerably more complex than that for type 19F. The type 19B backbone has an additional ManNAc which also carries a (1→3) linked  $\beta$ -D-Rib<sub>f</sub>-(1→4)- $\alpha$ -L-Rha<sub>p</sub> side-chain (**Fig. 5.1**). Therefore, biosynthesis of type 19B CPS would be predicted to require several additional and/or different enzymes; 3 transferases for the addition of the 3 extra sugars, at least one enzyme for the synthesis of the activated Rib precursor, a distinct polysaccharide repeat unit transporter and a distinct polysaccharide polymerase.

Analysis of the predicted protein products from the type 19B ORFs identified candidates for several of these enzymes. These are the rhamnosyl transferase (Cps19bQ), needed for the addition of Rha to the distal ManNAc, the polysaccharide repeat unit transporter (Cps19bJ), and the polysaccharide polymerase (Cps19bI). Cps19bR is the most likely candidate for the enzyme required for the synthesis of the activated Rib precursor. However, the two transferases needed for the addition of the distal ManNAc and Rib

remain unidentified. It is possible that Cps19bP functions as one of these transferases. However, there are no other ORFs in the type-specific region of *cps19b* which could encode the other transferase.

One possible explanation for the absence of an ORF encoding a third transferase is as follows. In type 19B CPS both ManNAc sugars are (1→4) linked. Cps19bF, which is almost identical to Cps19fF, a putative ManNAc transferase (section 1.9.3), could be responsible for the addition of both ManNAc residues. The ability of a single transferase to transfer the same sugar to what appears to be different acceptors has been proposed previously. For example, RfbG in *Shigella flexneri* is thought to add Rha, via the same linkage, at two separate positions within the O-antigen repeat unit (Morona *et al.*, 1995). Interestingly, type 19F and type 19A, have closely related *cps19J* (repeat unit transporter) genes (section 4.2.1) and their capsules are similar, as both have a trisaccharide backbone containing only one ManNAc (**Table 1.3**). On the other hand, both type 19B and type 19C, which also have highly homologous *cps19J* genes, have a tetrasaccharide backbone containing two ManNAc residues (**Table 1.3**). Given that Cps19fJ and Cps19bJ are quite distinct, the difference in the CPS backbone might be explained by the specificity of these polysaccharide repeat unit transporters. Cps19fJ may have absolute specificity for the trisaccharide repeat unit and may thus prevent Cps19fF from adding an additional ManNAc sugar to the backbone. On the other hand, Cps19bJ may be specific for the tetrasaccharide backbone and may only transport the repeat unit after the addition of both the distal ManNAc by Cps19bF and the disaccharide side-chain by Cps19bQ and Cps19bP. This hypothesis could be investigated by the construction of an Rx1-19F in-frame *cps19fJ* deletion mutant. Biochemical analysis of this strain should determine the presence of either a trisaccharide or a tetrasaccharide repeat unit linked to the lipid carrier. The presence of a tetrasaccharide would support the above hypothesis.

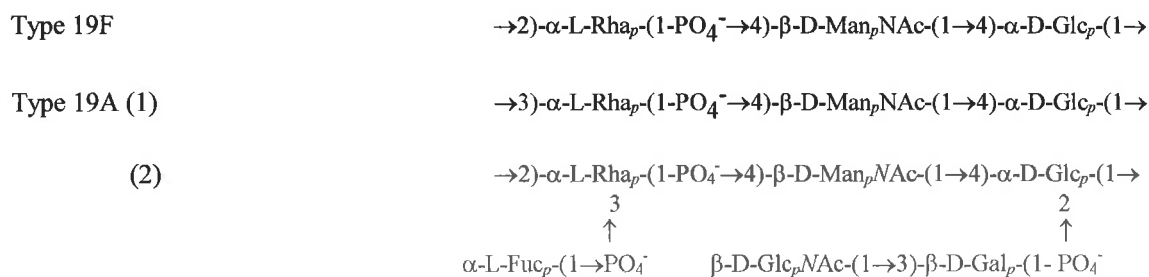
Alternatively, Cps19bR could be the ribosyl transferase, while Cps19bP could be the second ManNAc transferase. However, the activated Rib precursor is unlikely to be ubiquitous to the pneumococcus, as Rib is not a component of any essential cell polysaccharide, such as teichoic acid. Cps19bR is a relatively large protein and could be bi-functional, capable of both synthesising the activated Rib precursor and transferring it to the repeat unit. Construction and characterisation of defined in-frame deletion mutants in all of the type 19B-specific genes described in this chapter are required to determine their precise functions.

## Chapter 6

# ANALYSIS OF THE *S. PNEUMONIAE* TYPE 19A *CPS* LOCUS

## 6.1 Introduction

Analysis of purified type 19A CPS has yielded two distinct putative structures (**Fig. 6.1**). One is the same as type 19F except for a 1→3 linkage (rather than 1→2) between Glc and Rha (Katzenellenbogen and Jennings, 1983). This difference would necessitate an alteration only in the specificity of the polysaccharide polymerase (Cps19fI). The alternative structure involves the same trisaccharide backbone as type 19F, but with additional  $\beta$ -D-Glc<sub>p</sub>NAc-(1→3)- $\beta$ -D-Gal<sub>p</sub>-(1-PO<sub>4</sub><sup>-</sup>→2) and  $\alpha$ -L-Fuc<sub>p</sub>-(1-PO<sub>4</sub><sup>-</sup>→3) side chains attached to the Glc and Rha, respectively (Lee and Fraser, 1980). This would necessitate a number of additional enzyme activities not found in the *cps* locus of type 19F strains. Interestingly, individual type 19A strains were subsequently reported to be capable of producing either structural type, depending on the growth conditions (Lee *et al.*, 1987). Sequence analysis of the type 19A *cps* locus was undertaken in an attempt to provide a molecular explanation for these findings.



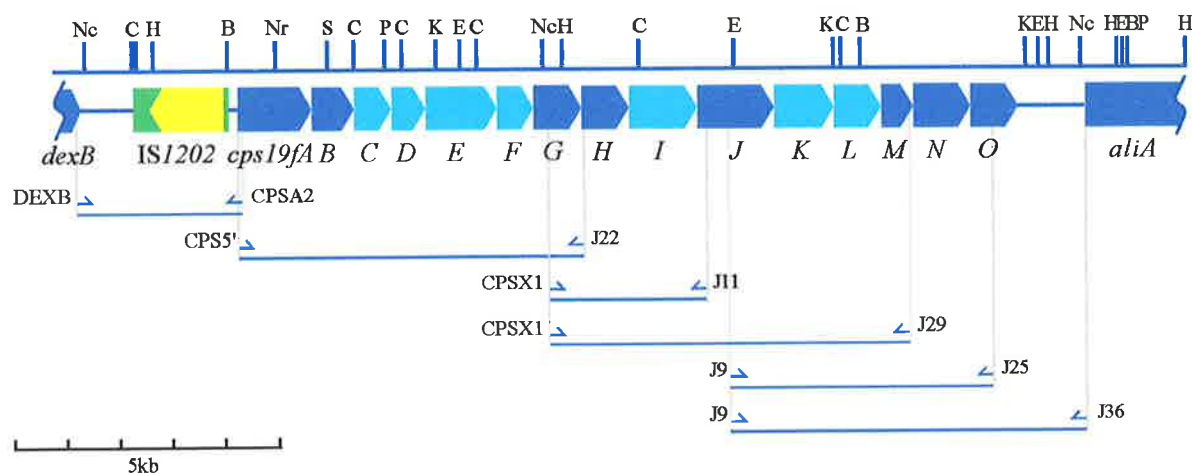
**Fig. 6.1. Biological repeat units of pneumococcal type 19F and type 19A capsular polysaccharide.** The order of the sugars in the repeat units have been altered compared to the published chemical structures for 19F (Ohno *et al.*, 1980) and 19A (Lee and Fraser, 1980; Lee *et al.*, 1987) reflecting the fact that glucose is the first sugar in the biological repeat unit. D-Glu<sub>p</sub>, glucose; D-Man<sub>p</sub>NAc, *N*-acetyl mannosamine; L-Rha<sub>p</sub>, rhamnose; D-Gal<sub>p</sub>, galactose; D-Glc<sub>p</sub>NAc, *N*-acetyl glucosamine; L-Fuc<sub>p</sub>, fucose; PO<sub>4</sub><sup>-</sup>, phosphate.

## 6.2 Results

### 6.2.1 PCR amplification and sequencing of the type 19A locus

Southern hybridisation data (section 4.2.1) suggested that the *cps* loci of type 19F and 19A are significantly different. Only 7 of the 15 *cps19f* gene-specific probes (*cps19fA*, *B*, *G*, *H*, *M*, *N* and *O*) hybridised at high stringency to 19A chromosomal DNA. It was assumed that the arrangement of these conserved genes within the two loci would be similar. Thus, a series of over-lapping DNA fragments containing type 19A-specific genes flanked by conserved sequences were generated by LR-PCR (section 2.9.2) using primers based on the *cps19f* sequence. The primers used are described in **Table 2.4** and a map of the PCR products spanning the entire *cps19a* locus is shown in **Fig. 6.2**. DNA from two different type 19A clinical isolates (19A1, obtained from J. Henrichsen, and 19A2, obtained from C. J. Lee [**Table 2.2**]) was used as template. Interestingly, the PCR products amplified from regions between *cps19A* and *cps19J* (using primer pairs CPS5'/J22 and CPSXI/J11 as described in **Table 2.4**) were identical in size for both type 19A isolates and type 19F, but the PCR products obtained from the 5' intergenic regions (using primers

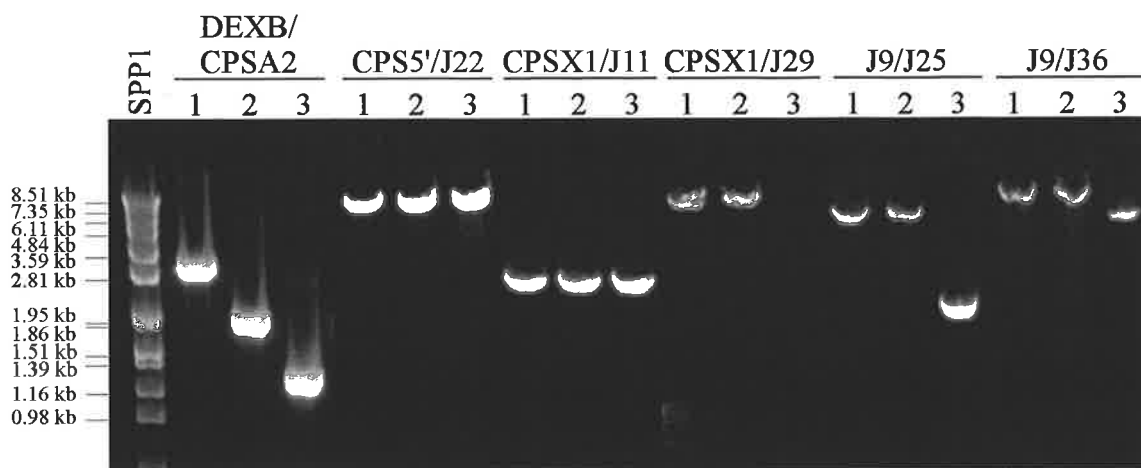
DEXB and CPSA2, as described in **Table 2.4**) varied between all three strains (**Fig. 6.3**). The 19A2 isolate also differed in the 3' region of the *cps19a* locus; the PCR products obtained from the *cps19J* to *cps19O* region (using primer pairs CPSXI/J29, J9/J25 and J9/J36, as described in **Table 2.4**) were either smaller or absent from 19A2 (**Fig. 6.3**). This suggests that part of this region of the *cps19a* locus of this strain may have been deleted. The PCR products obtained from this region of the *cps19a* locus (using primers J88 and J36, described in **Table 2.4**) from six Australian type 19A isolates were either identical in size or larger (probably indicating the presence of an IS element in the 3' intergenic region) than that from 19A1 and 19F (**Fig. 6.4**), suggesting that this part of the *cps19a* locus in 19A2 is atypical.



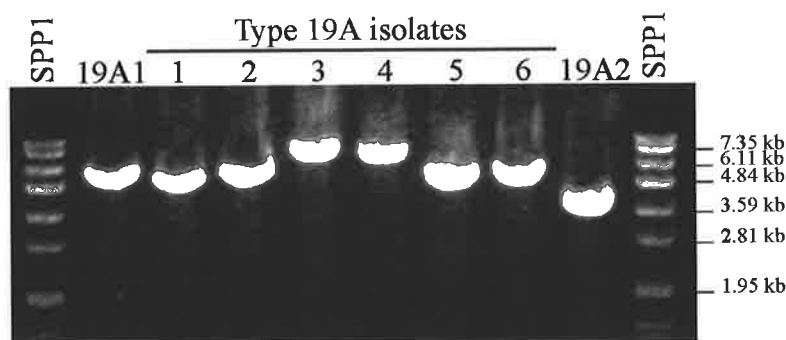
**Fig. 6.2.** Schematic representation of the PCR products amplified using *cps19f* primers. The *cps19f* genes which hybridised (at high stringency) to type 19A chromosomal DNA are shown in dark blue and the remainder of the locus is shown in pale blue. The primers used are described in **Table 2.4**.

## 6.2.2 Analysis of the *cps19a* locus

The sequence of the PCR products from 19A1 were determined using specifically designed primers, as described in section 2.8. Analysis of the compiled sequence (**Appendix III**) revealed that the *cps19f* and *cps19a* loci are very closely related. The



**Fig. 6.3.** Comparison of the PCR products obtained from *S. pneumoniae* Rx1-19F, 19A1 and 19A2. The PCR products from Rx1-19F (1), 19A1 (2) and 19A2 (3) were electrophoresed in a 0.8% agarose gel and stained with ethidium bromide. *Eco*RI-digested SPP1 bacteriophage DNA was used as a size marker (section 2.6.2); the approximate sizes of the PCR products are indicated on the left side of the figure.



**Fig. 6.4.** Comparison of the PCR products obtained from the 3' region of the *cps19a* locus from various type 19A isolates. The PCR products were electrophoresed in a 0.8% agarose gel and stained with ethidium bromide. *Eco*RI-digested SPP1 bacteriophage DNA was used as a size marker (section 2.6.2) and the approximate sizes of the fragments are indicated on the right side of the figure.

*cps19a* locus has the same number of ORFs organised in identical order to those in *cps19f* with homologies to the *cps19f* genes ranging from 70.1% to 99.4% identity. The sizes, G+C content and % identity of the *cps19a* and *cps19f* protein products are shown in **Table 6.1**.

Notwithstanding the overall similarity between the *cps19a* and *cps19f* loci, several interesting differences between the two loci were noted. Whereas the start codon for both *cps19aG* and *cps19aH* is TTG, only *cps19fH* has a TTG start codon in the *cps19f* locus.

The intergenic gaps between the *cps19a* genes and the *cps19f* genes are all similar, except between *cps19aK* and *cps19aL* which is much larger (152 nucleotides) compared to that between *cps19fK* and *cps19fL* (38 nucleotides). The largest variation between the *cps19a* and *cps19f* loci occurs in the 5' intergenic region. This region in the 19A1 strain has several deletions compared to the same region in type 19F, but the 3' intergenic regions of types 19F and 19A1 are almost identical (96.7% identity). The differences in the 5' intergenic region will be discussed in section 6.2.5.

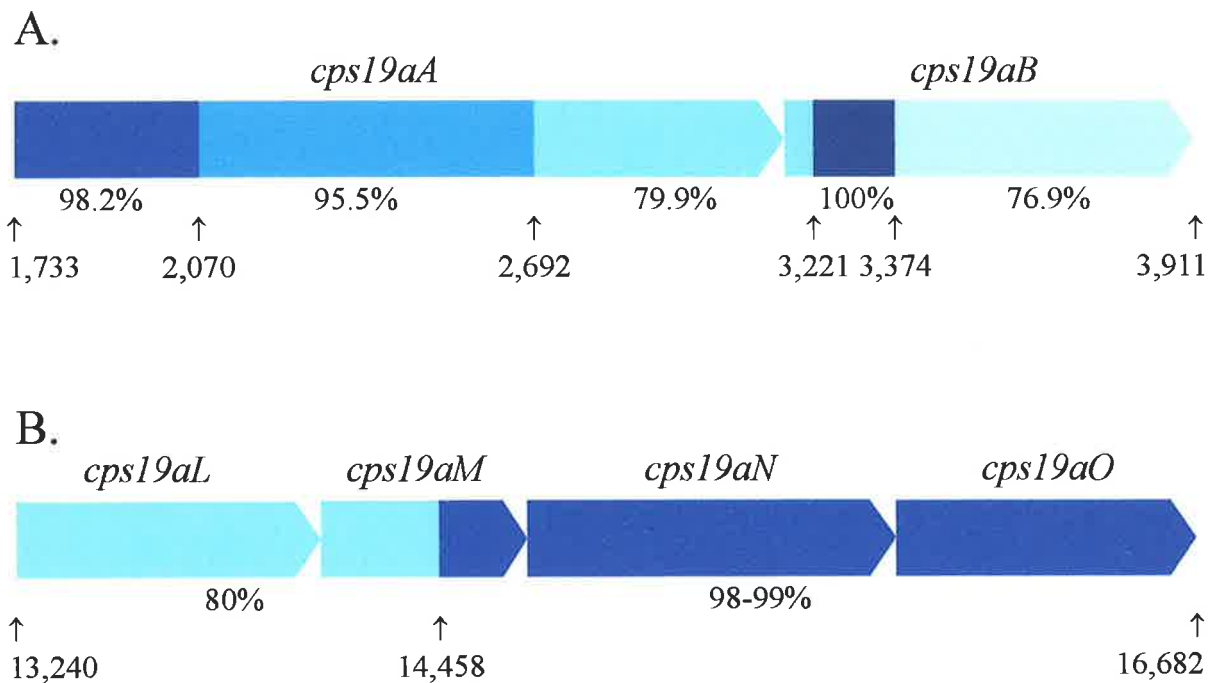
**Table 6.1. Comparison of the *cps19a* and *cps19f* ORFs.**

<i>cps19a</i> ORFs	Predicted size			<i>cps19f</i> ORFs	Predicted size			% Identity	
	Da	no. aa	%G+C		Da	no. aa	%G+C	DNA	aa
Cps19aA	53,576	481	39.5	Cps19fA	53,572	481	38.1	90.5	92.3
Cps19aB	28,138	243	41.3	Cps19fB	28,352	243	38	82	85.2
Cps19aC	25,473	230	42.1	Cps19fC	25,497	230	38.2	70.1	71.7
Cps19aD	25,155	229	41	Cps19fD	24,947	227	34.5	73	80.2
Cps19aE	51,971	453	37.7	Cps19fE	52,595	455	33.2	71.2	70.5
Cps19aF	28,273	247	34.1	Cps19fF	28,155	247	33.6	78.9	82.9
Cps19aG	31,195	266	37.2	Cps19fG	31,647	269	36.3	90.9	93.6
Cps19aH	34,455	292	32.2	Cps19fH	34,474	292	30.3	90.8	95.2
Cps19aI	51,604	444	32.7	Cps19fI	51,734	445	29.7	78.5	80.7
Cps19aJ	54,650	474	33	Cps19fJ	55,055	473	29.7	82.3	83.3
Cps19aK	40,749	362	36.9	Cps19fK	40,950	362	35.2	85.2	92.8
Cps19aL	32,242	289	43.3	Cps19fL	32,215	289	42.3	79.9	92.4
Cps19aM	22,408	198	41.2	Cps19fM	22,379	198	41.5	87.6	94.4
Cps19aN	39,086	349	42.4	Cps19fN	39,053	349	42.1	98.2	99.1
Cps19aO	32,330	283	41.3	Cps19fO	32,330	283	41.5	99.4	99.3

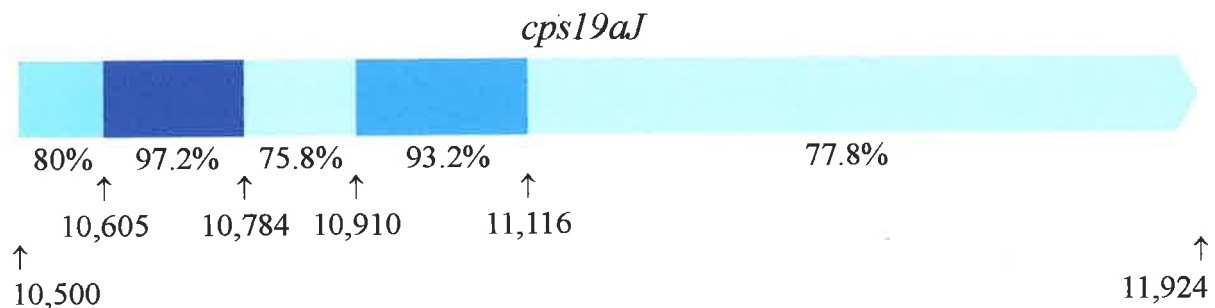
The *cps19a* and *cps19f* sequences were examined to identify potential cross-over points where recombination between the *cps* loci may have occurred. A cross-over point was identified within the *cps19M* gene; the first 348 nucleotides of *cps19aM* have 80.3%

identity to *cps19fM* whereas the remainder of *cps19aM* is 98% identical to *cps19fM*, as shown in **Fig. 6.5B**. However, there is no clearly identifiable cross-over point at the 5' end of the loci; the *cps19aAB* genes present a mosaic pattern with small regions of varying degrees of identity to the *cps19fAB* genes, ranging from 76.6% to 100% as shown in **Fig. 6.5A**. This suggests that the *cps19a* locus and the type 19A serotype may be the result of several recombination events between the ancestral *cps* locus (possibly type 19F) and exogenous DNA. Some of these recombination events may have involved small DNA fragments that did not affect the serotype, while others resulted in the exchange of larger regions of the capsule locus, which may have altered the expressed serotype. The *cps19aB* gene is almost identical (except for the first 42 nucleotides) to the *cps4B* gene which was described in chapter 4, and shows the same point of sequence divergence from the *cps19fB* gene. A small region of *cps19aB* (nucleotides 3,221-3,374) has 100% identity to *cps19fB*. This region presumably accounts for the high stringency hybridisation of the *cps19aB* DNA to a *cps19fB* probe (**Table 4.1**) as there is only 76.7% identity between the remainder of the *cps19aB* and *cps19fB* genes. The highly conserved region may either encode a functionally important domain in the *cps19B* gene product or may simply be the result of a recombination event.

The overall identity between *cps19fJ* and *cps19aJ* is only 82%, which is insufficient for the *cps19fJ* probe to hybridise to the *cps19aJ* gene under high stringency conditions. However, on closer examination of the sequences, two small regions (nucleotides 10,605-10,784 and 10,910-11,116) at the 5' end of *cps19aJ* have >90% DNA sequence identity (97.6% and 93.2%, respectively) to *cps19fJ* (**Fig. 6.6**), which presumably accounts for the Southern hybridisation data obtained previously (**Table 4.1**). Similarly, Cps19aJ and Cps19fJ have only 83.3% amino acid identity, but when their hydropathy profiles are compared they are almost indistinguishable (**Fig. 6.7**). This suggests that the divergence



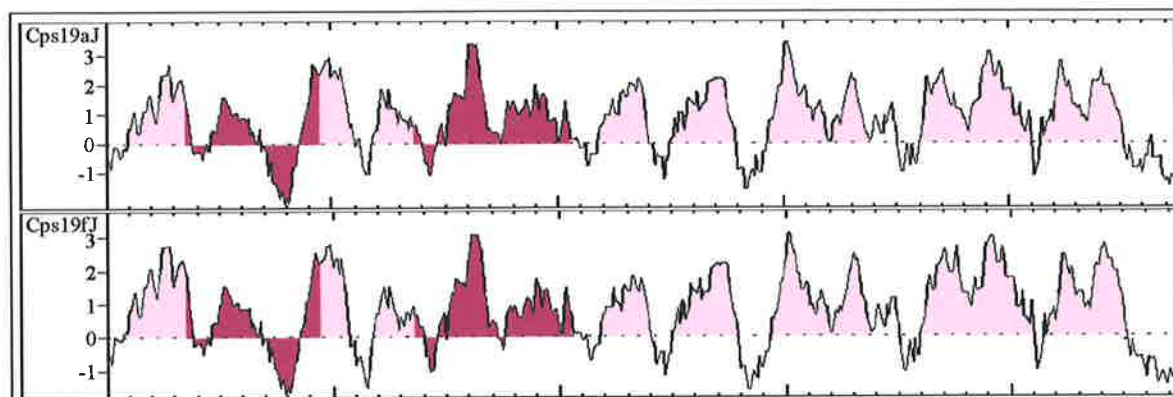
**Fig. 6.5. Diagrammatic representation of the similarity of the *cps19a* locus to *cps19f* in the regions where recombination between the loci may have occurred. A. The *cps19aAB* region. There are several possible recombination points in this region of the locus. B. The *cps19aL-O* region. There is a single cross-over point within *cps19aM*. Increasing similarity is represented by progressively darker shades of blue and the % identity is shown under the individual coloured regions. The arrows indicate the points of divergence and the number below the arrow corresponds to the nucleotide number in the *cps19a* sequence as shown in Appendix III.**



**Fig. 6.6. Diagrammatic representation of the similarity between *cps19aJ* and *cps19fJ* sequences. Increasing similarity of *cps19aJ* to *cps19fJ* is represented by progressively darker shades of blue and the % identity is shown under the individual coloured regions. The arrows indicate the points of divergence and the number below the arrow corresponds to the nucleotide number of the *cps19a* sequence as shown in Appendix III.**

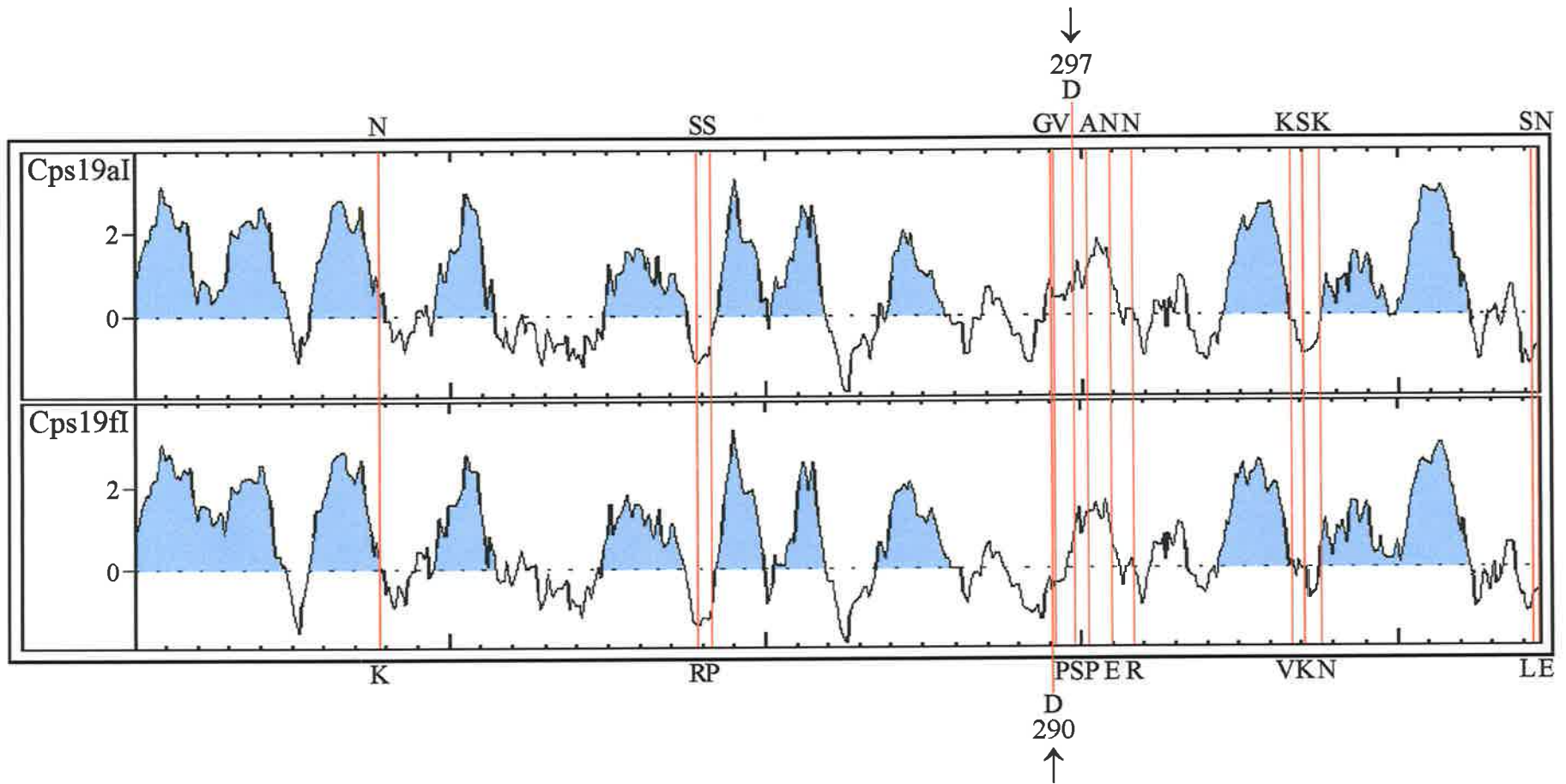
between these two sequences has resulted in conservative amino acid substitutions, which do not impact upon the function of Cps19aJ or Cps19fJ, both of which transport the same

trisaccharide repeat unit across the membrane. It is also tempting to speculate that the highly conserved regions (with >90% identity) may be important for the function of these proteins. The first conserved region incorporates the second membrane-spanning domain and the following hydrophilic domain; and the second conserved region includes a small hydrophilic domain as well as the fifth and sixth membrane-spanning domains (Fig. 6.7).



**Fig. 6.7.** Hydropathy profiles of Cps19aJ and Cps19fJ showing the conserved regions. The hydrophobic domains are shown in pink except in the highly conserved regions which are shown in maroon. Hydropathy plots were generated by the method of Kyte and Doolittle (1982) and aligned using PROFILEGRAPH (Hofmann and Stöffel, 1989). Positive numbers on the Y-axis indicate hydrophobic regions. The position of every 10th amino acid is marked on each X-axis.

The putative polysaccharide polymerases, Cps19aI and Cps19fI, are predicted to form different glycosidic linkages in type 19A ( $\alpha 1 \rightarrow 3$ ) and type 19F ( $\alpha 1 \rightarrow 2$ ) CPS, respectively. These two proteins are 80.7% identical and their amino acid sequences were examined to identify any potentially significant differences between them. Their hydropathy profiles are almost identical as shown in Fig. 6.8. A cluster of non-conservative amino acid substitutions were located in the region between amino acids 290 and 320 of Cps19fI and Cps19aI (Fig. 6.8). No such clustering of non-conservative amino acid substitutions was observed when comparing either Cps19fH/Cps19aH or Cps19fJ/Cps19aJ (results not shown). This region is predicted to be on the outer surface of the cytoplasmic membrane, based on the topology of the O-antigen polymerase (Rfc) from



**Fig. 6.8. Comparison of the hydropobicity profiles of Cps19aI and Cps19fI.** The hydropathy profiles of Cps19aI and Cps19fI were generated by the method of Kyte and Doolittle (1982) and aligned using PROFILEGRAPH (Hofmann and Stöffel, 1989). Positive numbers on the Y-axis indicate hydrophobic regions and putative membrane-spanning domains are shaded. The position of every 10th amino acid is marked on each X-axis. The red lines indicate the position of non-conserved amino acid changes with the amino acid indicated in single letter code above the line for Cps19aI and below the line for Cps19fI. The position of the Aspartate (D) residue is indicated by arrows.

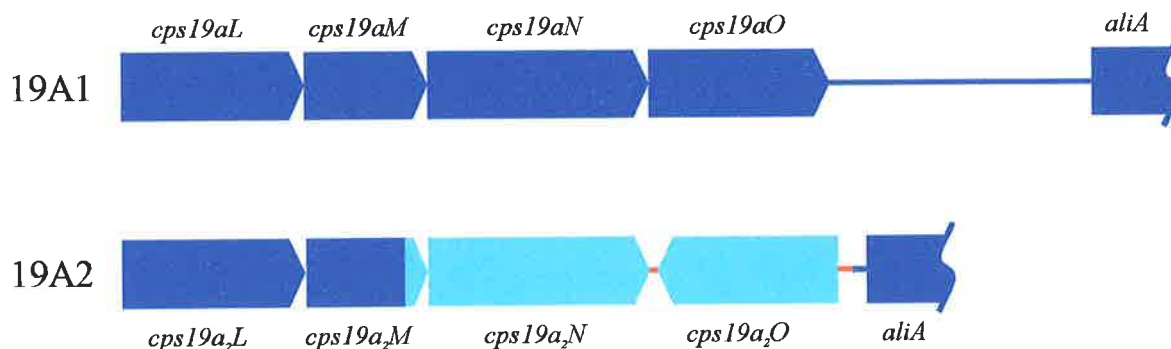
*S. flexneri* (Daniels *et al.*, 1998). A comparison of the hydropathy profiles of Cps19fI and *S. flexneri* Rfc was shown previously (Fig. 3.5). The CPS repeat units are predicted to be transported across the cytoplasmic membrane prior to polymerisation (section 1.8.2). Thus, the external location of the non-conserved regions in Cps19aI and Cps19fI is consistent with that of the putative catalytic site in these proteins. Of particular interest is the single aspartate (D) residue present in this region of both proteins, although its position is not conserved. Aspartate residues are predicted to be required for the formation of glycosidic linkages, and have been identified in the catalytic sites of  $\beta$ -glycosyl transferases (Saxena *et al.*, 1995). As shown in Fig. 6.8, the position of the aspartate residue varies by seven amino acids between the two proteins. This is equivalent to two turns of an  $\alpha$ -helix and would position the two aspartate residues on the same side of an  $\alpha$ -helix in both proteins, but in a different spacial position, perhaps facilitating formation of the altered glycosidic linkage formed by these two proteins.

### 6.2.3 Comparison of the *cps19a* loci from *S. pneumoniae* strains 19A1 and 19A2

The 19A2 PCR product obtained using primers J9 and J36, which amplified the 3' region of the *cps19a* locus encoding the dTDP-Rha biosynthesis genes (*cps19L-O*), was smaller in size to that from both Rx1-19F and 19A1 (Fig. 6.3). To identify the deletion present in 19A2, this PCR product was sequenced and this region of the locus (designated *cps19a<sub>2</sub>*) is shown in Appendix IV.

Analysis of the sequence identified a gene rearrangement in the 3' region of the *cps19a<sub>2</sub>* locus, as well as deletion of 1.4 kb of DNA between the end of *cps19aO* and the start of *aliA*, as shown in Fig. 6.9. The first 3,347 nucleotides of the *cps19a<sub>2</sub>* sequence have 99.8% identity to *cps19a*, followed by 1,185 nucleotides with 80% identity to *cps19a*

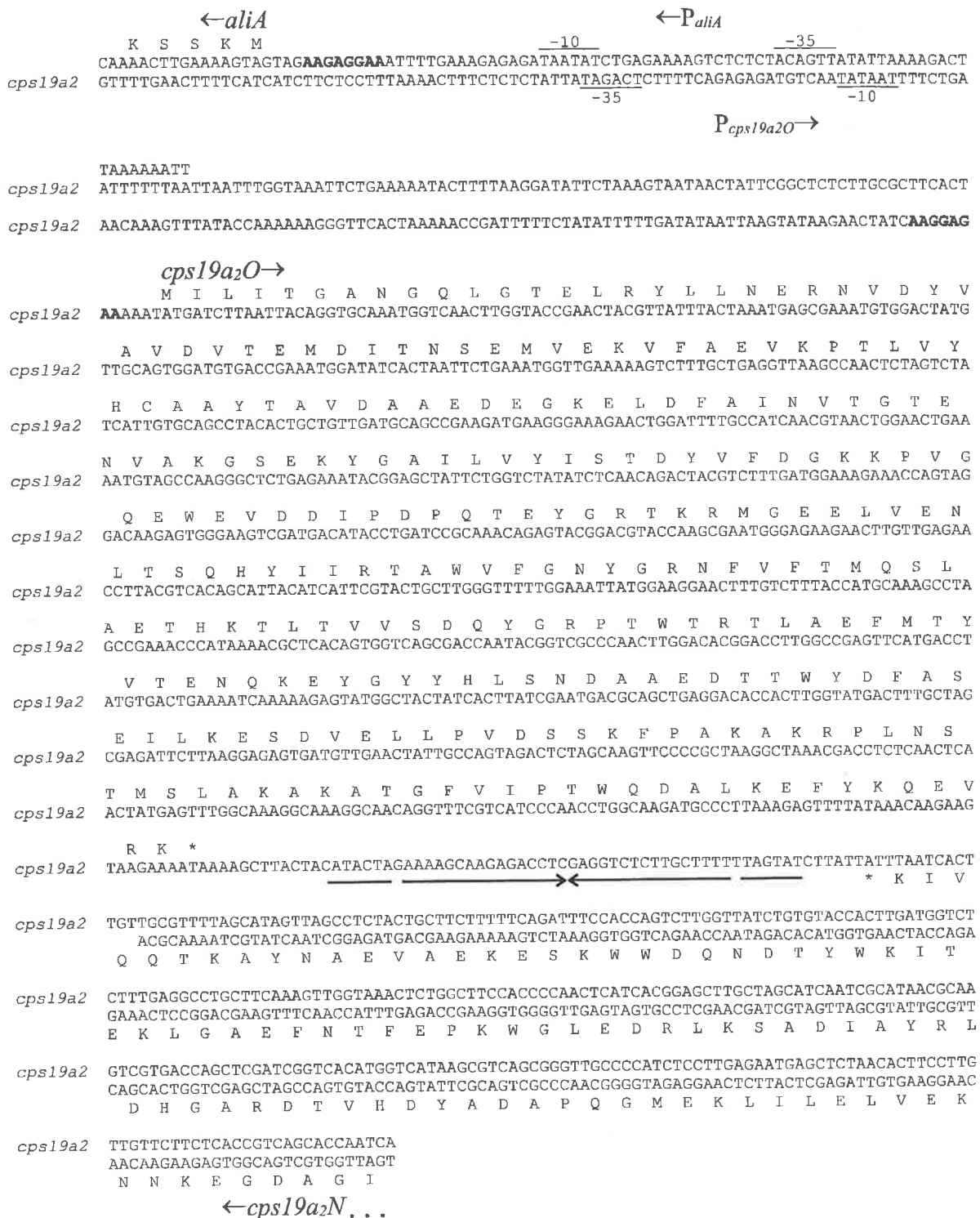
and 84% identity to *cps19f*. The remainder of the sequence then diverges until the final 94 nucleotides which are 90% identical to the same region in *cps19a*. The conserved regions contain the genes *cps19a<sub>2</sub>JKLMN*, with an approximate recombination point 120 nucleotides from the end of *cps19a<sub>2</sub>M*. The next 1.1 kb of DNA contains an inverted copy of *cps19a<sub>2</sub>O* (Fig. 6.9) with 76.4% identity to *cps19aO* and *cps19fO*. A potential promoter was identified upstream of *cps19a<sub>2</sub>O* in the same region (but on the opposite DNA strand) as that for *aliA* (Fig 6.10). There are 61 nucleotides between the stop codons of *cps19a<sub>2</sub>N* and *cps19a<sub>2</sub>O* and a stemmed-loop structure which could be a transcription terminator ( $\Delta G = -30.5$  kcal/mole) was identified in this region (Fig 6.10).



**Fig. 6.9.** Diagrammatic representation of the *cps19aL-aliA* region of *cps19a* and *cps19a<sub>2</sub>*. The dark blue region of *cps19a<sub>2</sub>* is >90% identical to the same region in *cps19a* and the pale blue region exhibits 75-80% identity to the equivalent regions in *cps19a*. The red line indicates small non-coding regions in *cps19a<sub>2</sub>* with no similarity to *cps19a*.

#### 6.2.4 Serotype specificity of the *cps19a* genes

The relationship between the *cps19aC*, *D*, and *E* genes and *cps* loci of other *S. pneumoniae* serotypes was examined by Southern hybridisation. The distribution of *cps19aA* and *B* genes was not investigated because previous Southern hybridisation data



**Fig. 6.10.** The sequence of *cps19a2*, from the 3' end of *cps19a2N* to the start of *aliA*. The sequence shown is the complementary strand of nucleotides 4,360-5,832 of the *cps19a2* sequence (Appendix IV). The double stranded sequence at the 5' and 3' ends represents the regions with homology to *cps19a*. The amino acids are shown above or below the first letter of the codon and the direction of translation is as indicated. The stemmed-loop structure is indicated by arrows below the sequence. The putative promoter regions (-10 and -35 sequences) for *aliA* and *cps19a2O* are also shown.

(section 4.2.1) indicated that homologues of the closely related *cps19fA* and *B* genes were present in all the serotypes tested. DNA fragments corresponding to the individual genes were isolated by PCR amplification using specific primers and restriction enzyme digestion (when necessary) and subsequent purification after agarose gel electrophoresis (**Table 6.2**). They were then labelled with DIG (section 2.7.1) and used to probe *ClaI*-restricted chromosomal DNA from representative pneumococci belonging to serotypes 2, 3, 4, 6A, 6B, 7F, 8, 9N, 9V, 12, 14, 16, 17, 18C, 19F, 19B, 19C, 20, 22, 23F, and 24. Hybridisation and washing conditions were at high stringency (as described in section 2.7.2.2).

**Table 6.2 DNA fragments used as gene-specific probes**

Probe	Fragment used as probe:
<i>cps19aC</i>	400 bp <i>PstI</i> fragment from J92-J95 PCR product (nt 3,932-4,356 of <i>cps19a</i> sequence)
<i>cps19aD</i>	500 bp <i>PstI</i> fragment from J93-J94 PCR product (nt 4,356-4,820 of <i>cps19a</i> sequence)
<i>cps19aE</i>	900 bp <i>ClaI</i> fragment from J87-J94 PCR product (nt 5,594-6,480 of <i>cps19a</i> sequence)
<i>cps19aK</i>	1.4 kb J70-J72 PCR product (nt 12,114-13,540 of <i>cps19a</i> sequence)

All primers used are described in **Table 2.4**. nt, nucleotide position.

The results for the hybridisation of both 19A and 19F genes are summarised in **Table 6.3**. The most remarkable feature of this table is that all the serotypes tested contained high stringency homologues of either *cps19fC-E* or *cps19aC-E*, except types 3 and 4 which do not have a high stringency homologue of either *cps19fE* or *cps19aE* (the gene which encodes the glucosyl transferase which adds Glc-1-phosphate to the lipid carrier). The absence of a *cpsE* homologue in types 3 and 4 is not surprising because the type 4 CPS does not contain Glc and the type 3 CPS is synthesised via a processive transferase as described in sections 1.8.4 and 8.4.2. Type 4 also contains a hybrid *cpsC* gene as previously described in chapter 4. The Southern hybridisation data presented here

**Table 6.3. Hybridization of *cps19fC-E* and *cps19aC-E* genes with other pneumococcal serotypes.**

Serotype	DIG labelled DNA probes					
	<i>cps19fC</i>	<i>cps19fD</i>	<i>cps19fE</i>	<i>cps19aC</i>	<i>cps19aD</i>	<i>cps19aE</i>
2	-	-	-	+	+	+
3	+	+	-	-	-	-
4	+	+	-	+	-	-
6A	-	-	-	+	+	+
6B	-	-	-	+	+	+
7F	+	+	+	-	-	-
8	-	-	-	+	+	+
9N	+	+	+	-	-	-
9V	-	-	-	+	+	+
12	-	-	-	+	+	+
14	+	+	+	-	-	-
16	+	+	+	-	-	-
17	-	-	-	+	+	+
18C	+	+	+	-	-	-
19F	+	+	+	-	-	-
19A	-	-	-	+	+	+
19B	+	+	+	-	-	-
19C	+	+	+	-	-	-
20	+	+	+	-	-	-
22	-	-	-	+	+	+
23F	-	-	-	+	+	+
24	+	+	+	-	-	-

Blue shaded '+' denotes hybridization with the gene-specific probes; - denotes no hybridization. The yellow shading represents no hybridization to either *cps19f* or *cps19a* probes.

suggests that the distinction between class I and class II genes identified for the *cpsC* gene in chapter 4, extends beyond *cpsC* to include *cpsD* and *cpsE*.

The presence of *cps19aK* homologues in serotypes 4, 9N, 9V and 12 was also examined, using a DIG-labelled *cps19aK* probe (described in Table 6.2) and low stringency conditions (as described in section 2.7.2.3). These serotypes are expected to contain a functional homologue to *cps19K* in their *cps* loci because ManNAc is a constituent sugar in their CPS. Only type 19A chromosomal DNA had hybridised with the *cps19fK* gene probe with low stringency conditions (section 4.2.1). When *cps19aK* was used as a probe, type 12 and the other members of serogroup 19 (19F, 19B and 19C) hybridised, suggesting that the gene(s) present in the serotypes 4, 9N and 9V are not closely

related to either *cps19aK* or *cps19fK*. Indeed, examination of the *cps4* sequence (section 9.3.4) reveals that *cps4L*, which probably encodes the UDP-GlcNAc-2-epimerase required for type 4 CPS biosynthesis, has 67.8% identity to the *cps19aK* gene and 65.9% identity to *cps19fK*, compared with 85.2% identity between *cps19aK* and *cps19fK*. This suggests that the UDP-GlcNAc-2-epimerase gene may have been derived from distinct sources in different serotypes.

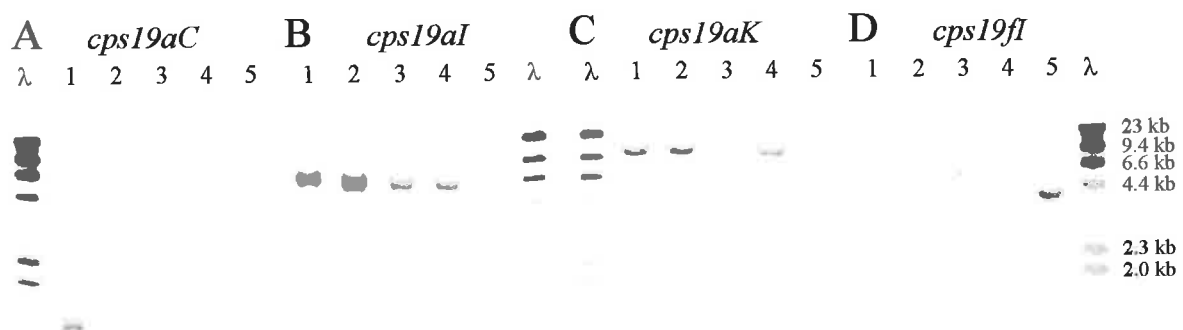
The serotype specificities of *cps19aG*, *H*, *I* and *J* were not reinvestigated as the homologues in *cps19f* were all serotype 19F and/or serogroup 19 specific.

### 6.2.5 Capsule transformation from type 19F to 19A

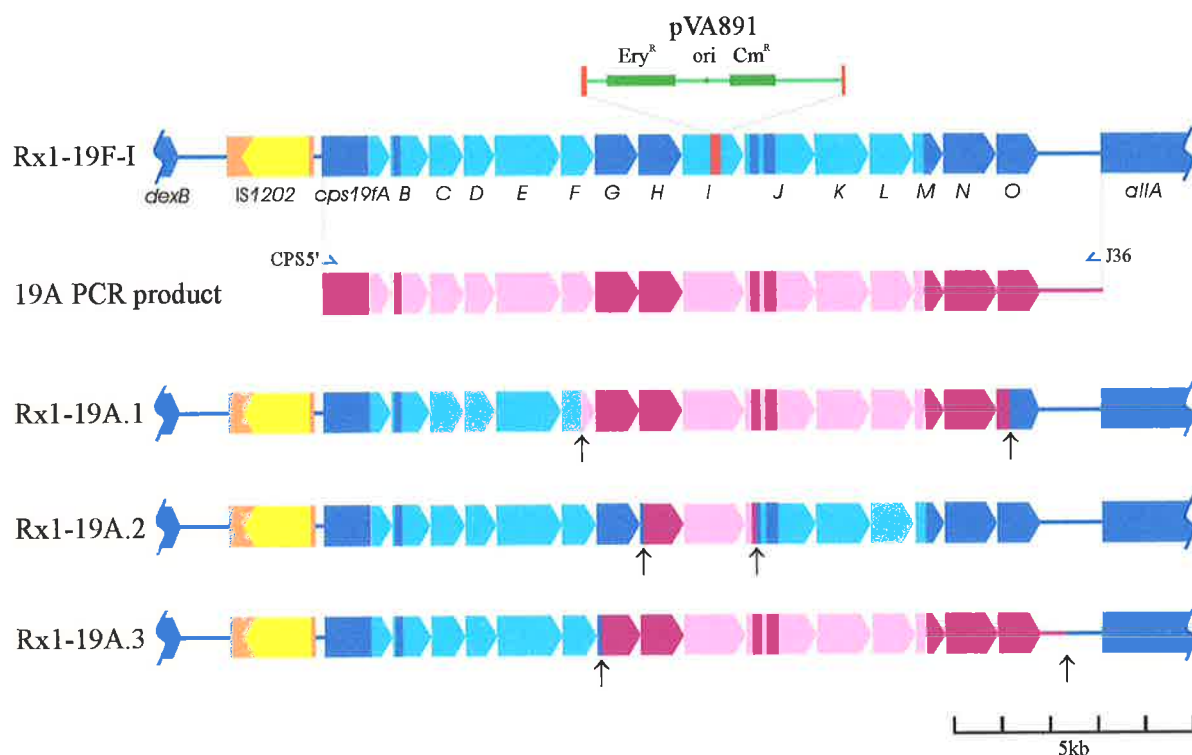
To confirm that the *cps19a* locus was sufficient for type 19A CPS biosynthesis, a 16.5 kb PCR product from the 5' end of *cps19aA* to the 5' end of *aliA*, was amplified using the primers CPS5' and J36 (**Table 2.4**). This was used to transform Rx1-19F-I, an unencapsulated, Ery-resistant derivative of Rx1-19F in which the *cps19fI* gene had been disrupted by insertion-duplication mutagenesis using pVA891, as described in section 3.2.4. Several smooth transformants were checked for Ery sensitivity, indicating loss of the pVA891 sequence. Southern hybridisation analysis was used to confirm the absence of both pVA891 and the *cps19fI* gene, and to determine if the *cps19aC*, *I* and *K* genes were present in three individual transformants (**Fig. 6.11**).

The production of a type 19A capsule by these three smooth transformants, designated Rx1-19A.1-3, was then confirmed by quellung reaction (section 2.3). The cross-over points between the *cps19f* locus and the type 19A PCR product were then identified by sequencing the regions where recombination was predicted to have occurred, based on the Southern hybridisation data in **Fig. 6.11**. The actual sequence data are shown in **Appendix V** and a diagram indicating the recombination points is shown in **Fig. 6.12**.

Two transformants (Rx1-19A.1 and 3) were similar, resulting in the exchange of a large region of the *cps19f* locus, from *cps19fG* to *cps19fN* (including *cps19fO* in Rx1-19A.3). On the other hand, Rx1-19A.2 is derived from exchange of a much smaller region of the *cps19f* locus, involving only *cps19H* and *cps19I* (Fig. 6.12). The *cps19aH* gene has 90.8% nucleotide identity to *cps19fH*, and the encoded highly conserved putative rhamnosyl transferases (95.2% amino acid identity) are predicted to be functionally identical in both type 19F and 19A CPS biosynthesis. The *cps19aI* and *cps19fI* genes are less conserved, with only 78.5% nucleotide identity, and the encoded putative polysaccharide polymerases (80.7% amino acid identity) are predicted to form different glycosidic linkages. Thus, this data shows that it is possible to alter capsule production from type 19F to type 19A by replacing only two genes in the *cps19f* locus, and that the presence of the *cps19aI* gene probably determines the 19A serotype.



**Fig. 6.11. Southern hybridisation analysis of the three Rx1-19A transformants.** 19A1 (1), Rx1-19A.1 (2), Rx1-19A.2 (3), Rx1-19A.3 (4) and Rx1-19F (5) were probed with DIG-labelled probes specific for *cps19aC* (A), *cps19aI* (B), *cps19aK* (C) and *cps19fI* (D).



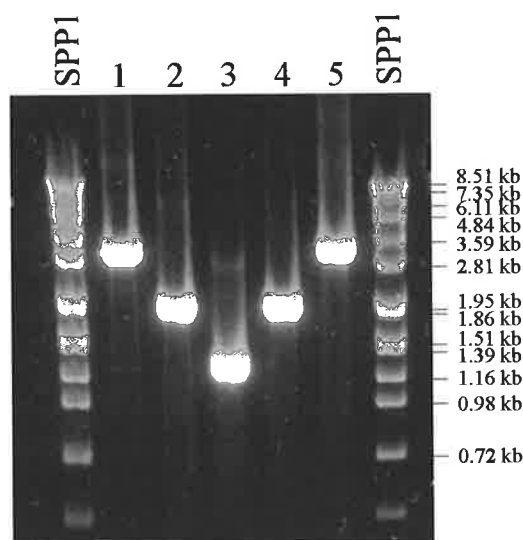
**Fig. 6.12. Diagrammatic representation of the cross-over points in the Rx1-19A transformants.** The *cps19f* locus is shown in blue and the *cps19a* locus is shown in pink. The dark colours represent regions with >90% DNA sequence identity and the light colours indicate 70-80% identity. The arrows indicate the points from which the sequence becomes *cps19a*-specific.

### 6.2.6 Sequence variation in the 5' intergenic region of serogroup 19

The 5' intergenic region of both type 19A isolates (19A1 and 19A2), as well as that for type 19F, 19B and 19C isolates were PCR amplified using the DEXB and CPSA2 primers (Table 2.4). The sizes of the PCR products varied between strains (Fig. 6.13), but those from 19A1 and 19B appeared to be the same size. The PCR products from 19A2 strain and the type 19B and 19C isolates were sequenced using specific primers and the 5' intergenic regions were compared as shown in Fig. 6.14.

Interestingly, the 5' intergenic regions of 19A1, 19B and 19C are all almost identical with three conserved deletions compared to 19F (Fig. 6.14). These three deletions remove all but 150 nucleotides of the 1-kb intergenic region between *dexB* and IS1202, as well as the 3' end of IS1202 (up to the stop codon of the putative transposase).

Another mutation at nucleotide position 2,151 (as depicted in **Fig. 6.14**) introduces a stop codon which interrupts the putative transposase in 19A1, 19B and 19C. The larger size of the PCR product obtained from 19C is due to the presence of an additional IS element,



**Fig. 6.13. PCR products of the 5' intergenic region.** The PCR products obtained, using DEXB/CPSA2 primers, from Rx1-19F (1), 19A1 (2), 19A2 (3), 19B (4) and 19C (5) were electrophoresed in a 0.8% agarose gel and stained with ethidium bromide. *Eco*RI-digested SPP1 bacteriophage DNA was used as a size marker (section 2.6.2) and the approximate sizes of the DNA fragments are indicated on the right side of the figure.

which will be described in chapter 7. This IS element inserted into the inverted repeat of *IS1202*, adjacent to the *cps19c* locus (as shown in **Fig. 6.14**).

Analysis of the derived sequences indicated that the 5' intergenic region of 19A2 is almost identical to that of 19F, except that it does not contain a copy of *IS1202* in the 5' intergenic region, although Southern hybridisation data has previously shown that this type 19A strain does contain a copy of *IS1202* in its chromosome (Morona *et al.*, 1994a). When PCR products from the 5' intergenic region from the six Australian type 19A isolates were examined by electrophoresis, they were all the same size as that from 19A2 (**Fig. 6.15**).

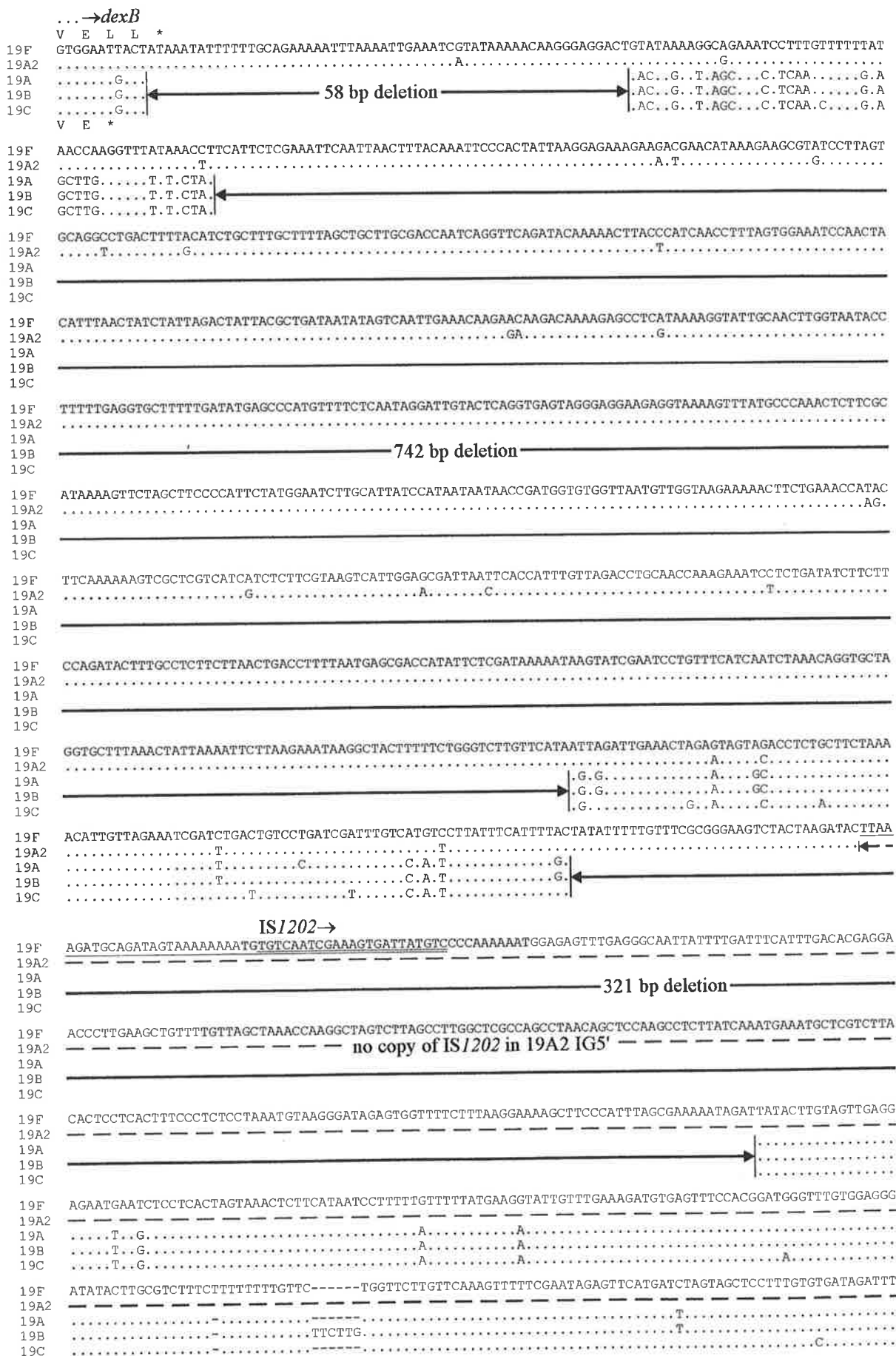
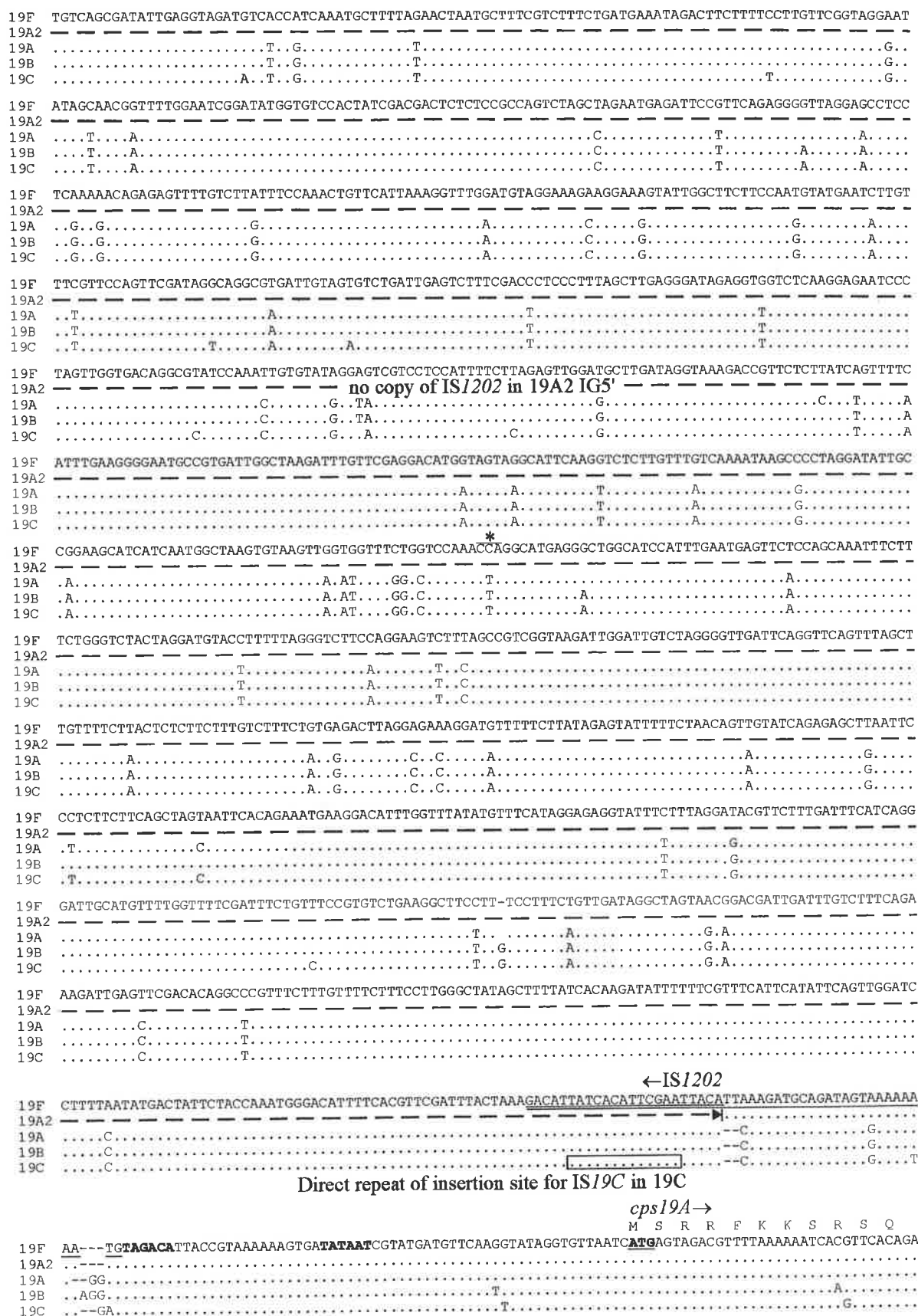
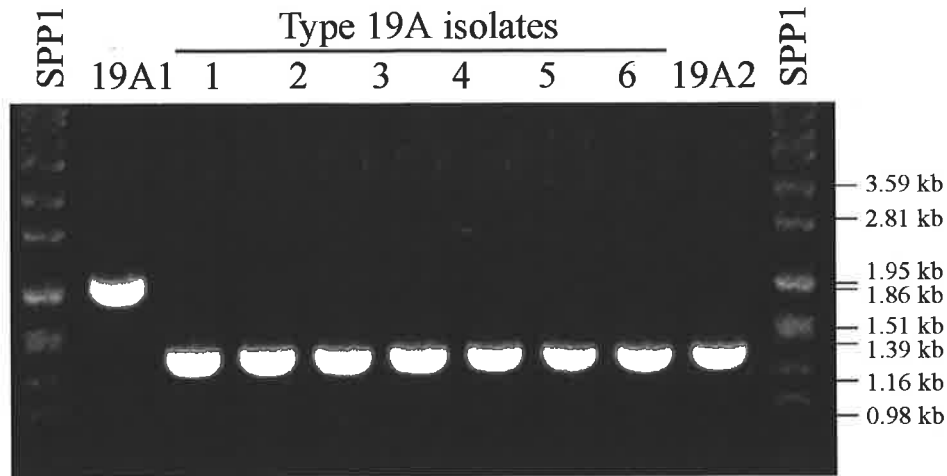


Fig. 6.14.

Continued overleaf.



**Fig. 6.14.** Comparison of the 5' intergenic sequences of *S. pneumoniae* strains Rx1-19F, 19A1, 19A2, 19B and 19C. The complete 5' intergenic sequences are shown (100 nucleotides per line) and are available under accession numbers U09239, AF094575, AF105112, AF105114, and AF105115, respectively. The stop codon in the truncated IS1202 transposase is indicated with an underlined asterisk. The -35 and -10 promoter regions are shown in bold and the start codon of *cpsA* is shown in bold and underlined. The amino acid translation of the 3' end of *dexB* and the 5' end of *cps19A* are as indicated. The direct repeat and inverted repeats associated with IS1202 are underlined and double underlined, respectively.



**Fig. 6.15. Comparison of PCR products from the 5' intergenic regions from several Australian type 19A isolates.** The PCR products were electrophoresed in a 0.8% agarose gel and stained with ethidium bromide. *EcoRI*-digested SPP1 bacteriophage DNA was used as a size marker (section 2.6.2) and the approximate sizes of the DNA fragments are indicated on the right side of the figure.

### 6.3 Conclusions

The *cps* locus from *S. pneumoniae* type 19A is similar to *cps19f*, in that it has the same number of genes arranged in the same order. However, large regions of the two loci share only 70-80% nucleotide sequence identity, suggesting that they either diverged a long time ago or that portions of the loci have separate origins. Some regions within the *cps19a* locus do have greater than 90% identity to *cps19f*, which may be a consequence of either recombination between the two loci, or perhaps due to a requirement for a higher degree of conservation in regions encoding functionally important domains. Transformation studies in section 6.2.4 have demonstrated that with the exception of *cps19aI*, all *cps19a* genes are functionally homologous to their *cps19f* counterparts. Genetic exchange of a DNA fragment containing just two genes, the *cps19aI* gene (which encodes the polysaccharide polymerase) and the highly conserved *cps19aH* gene (which encodes the rhamnosyl transferase), was sufficient to convert a strain containing an interrupted *cps19f* locus to type

19A. This is consistent with the fact that according to the structure proposed by Lee and Fraser (1980) type 19A CPS differs from type 19F only by the type of glycosidic linkage between identical trisaccharide repeat units. The studies in this chapter have shown that the genes which are present in the *cps19a* locus are sufficient for type 19A CPS biosynthesis and hence, the biosynthetic pathway for type 19A CPS is essentially identical to that proposed for type 19F CPS in **Fig. 3.16**.

Of particular interest, are the two closely related polysaccharide polymerase genes (*cps19aI* and *cps19fI*). The proteins encoded by these genes have almost identical hydrophobicity profiles. A cluster of non-conservative amino acid changes between amino acids 290-320 has been identified in a hydrophilic region of the proteins and is a likely candidate for the putative catalytic domain. These two proteins provide a perfect model for further characterisation of this catalytic domain. Exchange of this small non-conserved region between the two proteins, in which the position of the potentially important aspartate residues differ by seven amino acids, could potentially alter the glycosidic linkage formed between the repeat units and thus the CPS serotype expressed. Also site-directed mutagenesis, altering the aspartate residues, might be predicted to abrogate the function of the polymerase. It would then be possible to investigate restoring the function of the polymerase, with altered specificity, by altering the position of the aspartate residue.

No additional genes, which might be involved in type 19A CPS biosynthesis, were identified either in or adjacent to the *cps19a* locus. Thus, the extra genes required to synthesise the side-chains proposed in the alternative type 19A CPS structure proposed by Lee *et al.* (1987), as shown in **Fig. 6.1**, must be located elsewhere on the *S. pneumoniae* chromosome. It is not known if these extra putative genes are present in all pneumococci or are specific to type 19A strains. The three Rx1-19A transformants described in section 6.2.4 synthesise type 19A CPS as judged by quellung reaction. However, it is not known if

they are capable of synthesising the proposed side-chains. If they are, these extra genes must be present in Rx1 (a rough derivative of a type 2 strain) and are therefore likely to be ubiquitous to all pneumococci. Chemical analysis of the CPS of these transformants grown under conditions reported to result in the biosynthesis of the alternative type 19A CPS (Lee *et al.*, 1987) is required to resolve this issue.

In chapter 4, it was noted that the *cpsC* gene could be divided into two distinct classes (I and II), depending on the degree of similarity to *cps19fC*. Whereas the *cps19fC* gene belonged to class I, the *cps19aC* gene belonged to class II. Southern hybridisation data were obtained using *cps19aC-E* probes and chromosomal DNA from the same pneumococcal serotypes analysed with *cps19fC-E* probes in **Table 4.1**. The results indicated that all three genes were present together, as either class I (*cps19fC-E* homologues) or class II (*cps19aC-E* homologues) in most serotypes tested. The *S. pneumoniae cps* loci, can be divided into class I and class II, on the basis of their *cpsC-E* genes. Using these criteria, the *cps* loci from pneumococcal types 1, 3, 4, 14, 19F, 19B, and 19C are class I and types 2, 19A, 23F, 33F are class II. The presence of *cpsB* homologues in all loci was demonstrated previously (**Table 4.1**), and the point at which the class I and class II sequences diverged was shown to occur within this gene in the serotypes tested (section 4.2.4). However, the precise point of divergence can vary, and type 19A and 19F *cps* loci appear to diverge within the *cpsA* gene (**Fig. 6.5**).

The intergenic regions flanking the *cps* loci from all six Australian *S. pneumoniae* type 19A isolates tested appear to be almost identical to type 19F except that *IS1202* was absent. A type 19F strain which lacks *IS1202* has been previously reported (Morona *et al.*, 1994a). However, two of the type 19A isolates did appear to contain extra DNA at the 3' end of the locus, which has not been investigated further, and may indicate the presence of yet another IS element in the 3' intergenic region. The common occurrence of IS elements

in the intergenic regions flanking the *cps* loci in different *S. pneumoniae* strains is discussed in section 8.4.4. The similarity in the 5' intergenic region between Rx1-19F, 19A2 and the six Australian type 19A isolates is suggestive of a common ancestry for type 19F and 19A. The unexplained rearrangement at the 3' end of the *cps* locus of 19A2, where the *cps19a<sub>2</sub>O* gene has been inverted, is probably a subsequent event. The *S. pneumoniae* type 19B and 19C strains analysed had an identical series of deletions in the 5' intergenic region, suggesting that they share the same clonal origin. The additional IS element found in this region in 19C could have integrated either before or after the acquisition of the type 19C-specific gene(s). Strain 19A1, which has an identical 5' intergenic region to 19B, may have originated from a recombination event between an ancestral type 19B strain and the *cps19a* locus.

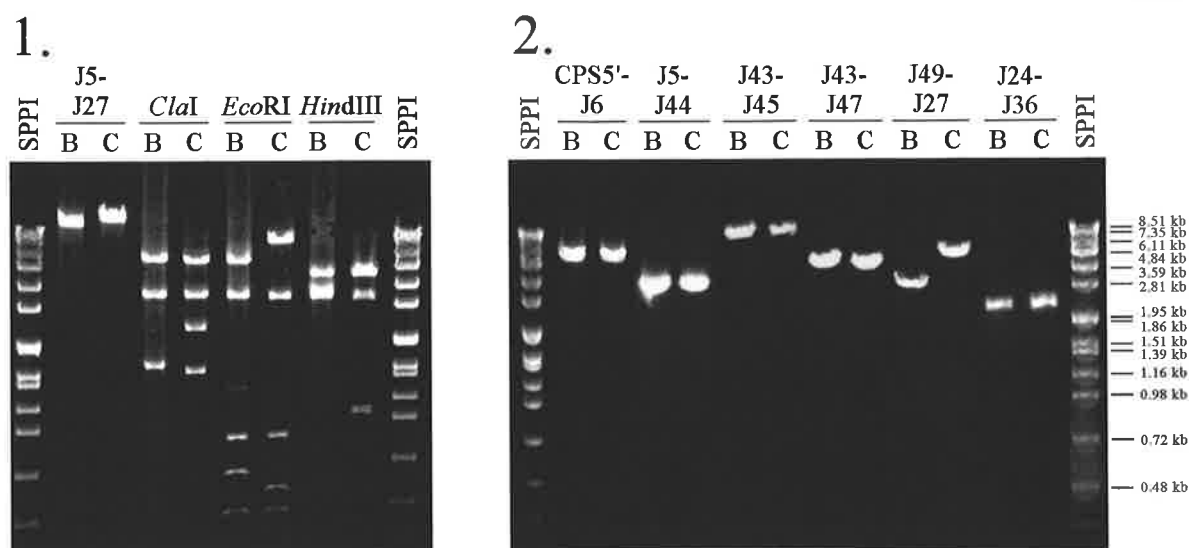
## Chapter 7

# CHARACTERISATION OF THE *S. PNEUMONIAE* TYPE 19C SPECIFIC CPS REGION

## 7.1 Introduction

The structure of the type 19C CPS is similar to that of type 19B. The only difference is that the type 19C repeat unit has an additional Glc side-chain, as shown in **Fig. 7.1**. Thus the type 19C capsule locus (*cps19c*) would be predicted to contain an extra gene required for the addition of this side-chain. Southern hybridisation analysis with *cps19f* gene specific probes has indicated that the *cps19b* and *cps19c* loci are almost identical with all but two of the genes hybridising to type 19B and 19C DNA (section 4.2.1). These two genes, *cps19fI* and *cps19fJ*, are located together near the middle of the *cps19f* locus and are replaced by the five genes *cps19bP*, *I*, *Q*, *R* and *J* in the *cps19b* locus (section 5.2.1). Homologues of these 19B-specific genes were also shown to be present in type 19C by hybridisation analysis (section 5.2.3). The extra gene required for type 19C biosynthesis may also be located in this region of the *cps19c* locus. Accordingly, the work described in this chapter was aimed at locating and characterising any additional gene(s) required for synthesis of type 19C CPS.





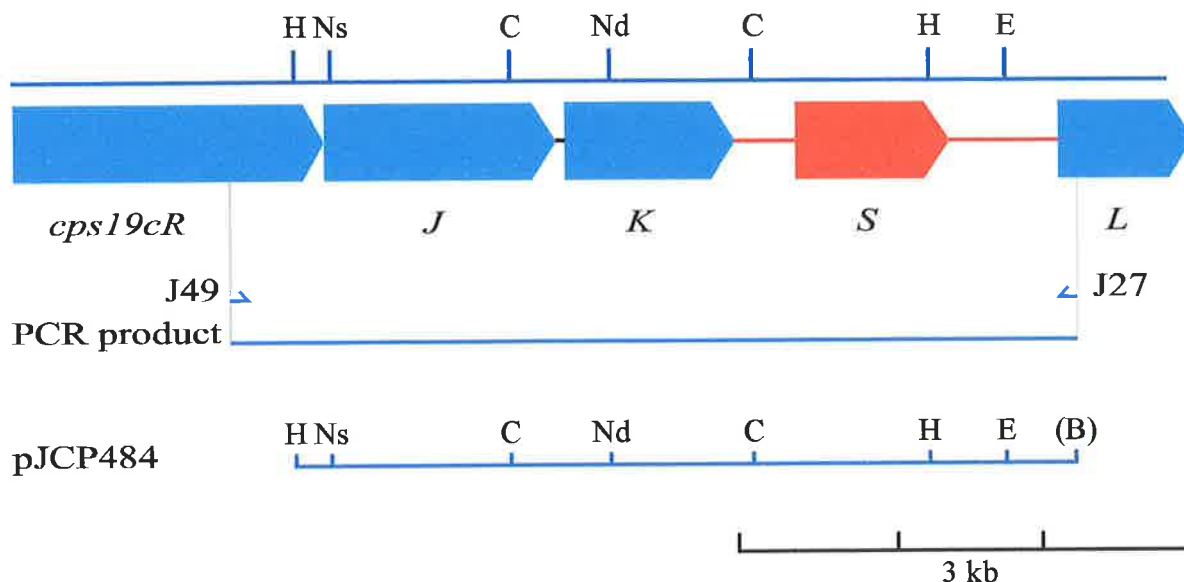
**Fig. 7.2. Comparison of PCR products obtained from *S. pneumoniae* types 19B and 19C.** PCR products obtained from 19B (B) and 19C (C) using various primers were electrophoresed in a 0.8% agarose gel and stained with ethidium bromide. Panel 1 shows the PCR products obtained using primers J5 and J27 either undigested or restricted with *ClaI*, *EcoRI* or *HindIII* prior to electrophoresis. Panel 2 shows undigested PCR products obtained using other primers as indicated. *EcoRI*-digested SPP1 bacteriophage DNA was used as a size marker (section 2.6.2); the approximate sizes of the PCR products are indicated on the left side of the figure.

**Table 7.1. PCR products obtained from *cps19b* and *cps19c*.**

PCR product <sup>a</sup>	Section of locus amplified	Size of PCR product <sup>b</sup>	
		Type 19B	Type 19C
CPS5'-J6	<i>cps19A-cps19F</i>	5.5 kb	5.5 kb
J5-J44	<i>cps19F-cps19I</i>	3.5 kb	3.5 kb
J43-J45	<i>cps19P-cps19K</i>	6.5 kb	6.5 kb
J43-J47	<i>cps19P-cps19J</i>	5 kb	5 kb
J49-J27	<i>cps19R-cps19L</i>	4 kb	6 kb
J24-J36	<i>cps19O-aliA</i>	2.5 kb	2.5 kb

a. PCR primers are described in Table 2.4. b. sizes are approximate.

between *cps19H* and *cps19K* was the same size in both types 19B and 19C, the additional gene must be located between *cps19cK* and *cps19cL* in the *cps19c* locus. The J49-J27 PCR product was purified (section 2.9.4) for further analysis. A map of the 5.3 kb PCR product was obtained, using the restriction enzymes *BamHI*, *ClaI*, *HindIII*, *NsiI*, *NdeI* and



**Fig. 7.3. Physical map of part of the *cps19c* locus.** Boxed arrows represent potential ORFs. Gene designations are indicated below the map; *cps19cB-S* are abbreviated to *B-S*, respectively. Restriction sites are as follows; B, *Bam*HI; C, *Cla*I; E, *Eco*RI; H, *Hind*III; Nd, *Nde*I; Ns, *Nsi*I. The region of DNA subcloned into pBluescript KS+ is shown below the map. The *Bam*HI restriction site is bracketed because it is generated by the J27 primer and not present in the chromosome of *S. pneumoniae* type 19C.

*Eco*RI, and it was sub-cloned into pBluescript KS+, generating pJCP484, as shown in **Fig. 7.3.**

Both strands of the pneumococcal DNA insert, and nested derivatives thereof, were subjected to sequence analysis in order to compile the sequence of the portion of the *cps19c* locus, as shown in **Appendix VI**. Examination of the compiled sequence revealed, as expected, that the first 2.9 kb of sequence at the 5' end has a high degree of similarity to the *cps19b* sequence. This region contains the homologues of *cps19bR*, *cps19bJ* and *cps19bK* (*cps19cR*, *cps19cJ* and *cps19cK*) which exhibited 98.5%, 99.7% and 94.9% identity, respectively. The sequence then diverges (at nucleotide 2,906 of the *cps19c* sequence) just prior to the end of *cps19cK*; the sequence of nucleotides 2,954-3,155 exhibits 74.8% identity to the 5' region of *cps19bL* but does not encode an ORF, and the sequence then diverges as shown in **Fig. 7.4A**. An additional potential ORF, designated *cps19cS* (shown in **Fig. 7.3**), is located between *cps19cK* and *cps19cL* and has a TTG start codon, which is preceded by a ribosome binding site, as shown in **Appendix VI**. The

A.

```

...cps19cK→
M S Q A S N P Y G K G D A S K Q I V H I L S G I *
cps19c AATGAGTCAAGCTAGTAATCCTTATGGAAAAGGTGATGCTAGTAAACAGATTGTTTCATATTTAAGCGGAATTTAAGCGAGGCCAA
cps19b AATGAGTCAAGCTAGTAATCCTTATGGAAATGGTATGCTAGTAAACAGATTGTTTCGATTTTACGTGGAATTTGAGTGTGTTTAG
M S Q A S N P Y G N G D A S K Q I V R I L R G I *
...cps19bK→
cps19c ATAAAGTAATAAAAAACACTA---TCTATAAAAGGTATTGATCTTGTAGTTGATTGCGGAACATGTTTCATATCCTTTGACTCGAG
cps19b ATAAAGTAATAACAGAAAGGTACCTACTATGAAAGGTATTATTCTAGCAGGTGGTTCGGGGACTCGCTTGTATCCTTTGACTCGCG
M K G I I L A G G S G T R L Y P L T R A
cps19bL→
cps19c CTACATAAAAATAACTTGTGCTGATTTATGATAGATCGATAATTTACTA-CTACTTTGGACATTGATGTTAGCAGTTATTAGGGAT
cps19b CTGCATCAAACAACCTGATGCCGGTTTATGATAAACCCATGATTTACTACCCACTTTCAACATTGATGTTGGCTGGAATTAGGGAT
A S K Q L M P V Y D K P M I Y Y P L S T L M L A G I R D
GATTTTATCATTACTTCAGTAAGTTAAATCG
cps19c GTTTTGATTAACTCAACTTTTCAGGATTGCGCTTGCTTCTAGATTTCTTCT-GATTTTATCATTACTTCAGTAAGTTAAATCG
cps19b ATTTTGATTATTTCCACTCCACAGGATTACATCGATTCCAAGAGCTTCTTCAAGACGGATCTGAGTTGGCTCAAACCTTCTT
I L I I S T P Q D L H R F Q E L L Q D G S E F G L K L S Y

```

B.

```

cps19c TAATCTCTTCCAGTATTTGTAGTAGAATTAAGTGTCTTGGATTAAAATAAAGAAC--ACAACACTTTTTATTCAGTGTGTGTA
cps19b AATGAGTCAAGCTAGTAATCCTTATGGAAT-GGTGATGCTAGTAAACAGATTGTTTCGATTTTACGTGGAATTTGAG-TGTGTT
M S Q A S N P Y G N G D A S K Q I V R I L R G I *
...cps19bK→
cps19cL→
M K G I I L A G G S G T R L Y P L T
cps19c TGGGTGAA--ACGAAAGGAACG-ATTGACTTATGAAAGGTATTATCTCGCGGGTGGTTCGGGGACACGTTTATATCCTTTGACT
cps19b TAGATAAAGTAATAACAGAAAGGTACCTAC-TATGAAAGGTATTATCTAGCAGGTGGTTCGGGGACTCGCTTGTATCCTTTGACT
M K G I I L A G G S G T R L Y P L T
cps19bL→
R A A S K Q L M P V Y D K P M I Y Y P L S T L M L A G I R
cps19c CGAGCTGCATCAAAGCAACTGATGCCGGTTTATGATAAACCGATGATTACTACCCACTTCAACTTTGATGTTGGCTGGGATTAG
cps19b CGCGCTGCATCAAACAACCTGATGCCGGTTTATGATAAACCCATGATTACTACCCACTTCAACATTGATGTTGGCTGGAATTAG
R A A S K Q L M P V Y D K P M I Y Y P L S T L M L A G I R
D I L I I S T P Q D L P R F K E L L Q D G S E F G I Q L S
cps19c GGATATTTTGATTATCTCAACTCCTCAAGATTTGCCTCGTTTTAAGAGCTCCTTCAAGATGGCTGAGTTGGGATTCAATTGT
cps19b GGATATTTTGATTATTTCCACTCCACAGGATTTACATCGATTCCAAGAGCTTCTTCAAGACGGATCTGAGTTGGGTCAAACTTT
D I L I I S T P Q D L H R F Q E L L Q D G S E F G L K L S

```

**Fig. 7.4. Homology between *cps19b* and *cps19c*.** Location of the approximate points of sequence divergence, at the end of *cps19K* (A), and immediately upstream of *cps19L* (B). Nucleotides 2,850-3,188 and 4,905-5,243 of the *cps19c* sequence are represented in panels A and B, respectively. The *cps19b* sequence (nucleotides 10,255-10,599 of Appendix II) is the same in both panels. The amino acid translation for the depicted ORFs are shown above the sequence for *cps19c* and below the sequence for *cps19b*. The stop codons are underlined and the ATG start codons are double underlined. The repeated sequence in *cps19c* is shown in red. Arrows indicate the approximate points of divergence. Identical amino acids are shown in bold.

closest potential ATG start codon is located 138 nucleotides downstream but it is not preceded by a ribosome binding site. As predicted, the 3' end of the *cps19c* sequence again shows similarity to the *cps19b* sequence, starting from nucleotide 5,017 (**Fig. 7.4B**); this is immediately before the start of the *cps19cL* gene, which has 90.6% identity to *cps19bL*. Thus, the first 200 nucleotides of *cps19cL* are duplicated in the *cps19c* locus, with 77.8% identity between the copies present immediately downstream of *cps19cK* and at the 5' end of *cps19cL*, as shown in **Fig. 7.4**. There are potentially significant intergenic gaps immediately before and after the *cps19cS* gene of 370 and 633 nucleotides, respectively. However, no potential stemmed-loop structures or obvious promoter sequences were found in these intergenic regions.

## 7.2.2 Characterisation of Cps19cS

The type 19C-specific ORF *cps19cS* is located between *cps19cK* and *cps19cL* in the *cps19c* locus (nucleotides 3,276-4,385) and encodes a putative 43.2 kDa protein containing 343 amino acids. This hydrophilic protein has a hydrophobicity index (according to Kyte and Doolittle [1982]) of -0.23 and a predicted pI of 5.18. The region from the 3' end of *cps19cK* to the 5' end of *cps19cL* has a G+C content of 30.4% increasing slightly to 31.4% for the *cps19cS* coding region. This is lower than the %G+C of the two flanking genes *cps19cK* (35.3%) and *cps19cL* (42.6%).

Database searches with Cps19cS found significant similarity to the C-termini of various glycosyl transferases as shown in **Table 7.2**. The alignment of the conserved C-terminal regions of these proteins is shown in **Fig. 7.5**. Interestingly, one of these glycosyl transferases, CpoA, is possibly involved in teichoic acid biosynthesis in *S. pneumoniae* (Grebe *et al.*, 1997). Cps19cS exhibits 21% identity along its entire length to WaaG (**Table 7.2**), an  $\alpha(1\rightarrow3)$  glucosyl transferase involved in LPS core biosynthesis in *E. coli*

and *S. enterica* serovar typhimurium (Heinrichs *et al.*, 1998). Thus, Cps19cS could function as the glucosyl-transferase required for the addition of the  $\beta(1\rightarrow6)$  linked Glc side-chain to the backbone in type 19C CPS biosynthesis.

**Table 7.2. Similarity of Cps19cS to other proteins.**

	% Identity <sup>a</sup>								
	Mj RfbU	Aa MtfC	Af Gal	Sa CapM	Ye TrsD	Vc RfbV	Ea AmsD	St WaaG	Sp CpoA
Cps19cS <sup>b</sup>	26 [127]	28.5 [158]	22.9 [223]	32.3 [93]	23 [222]	26.8 [142]	22.2 [221]	21 [300]	28.7 [94]
MjRfbU <sup>c</sup>	100	30.3 [195]	34.2 [190]	25.4 [244]	26.3 [179]	25.4 [213]	26.9 [119]	30.4 [161]	24.3 [173]
AaMtfC <sup>d</sup>		100	24.8 [311]	23.7 [359]	21.6 [361]	34.6 [373]	23 [183]	21.9 [233]	19.5 [246]
AfGal <sup>e</sup>			100	25.1 [175]	20.9 [344]	24.5 [372]	25.6 [203]	27.9 [183]	21.5 [251]
SaCapM <sup>f</sup>				100	25.6 [164]	20.9 [349]	17.7 [220]	13.9 [230]	15.2 [230]
YeTrsD <sup>g</sup>					100	30.3 [165]	33.1 [362]	20.7 [227]	25.4 [114]
VcRfbV <sup>h</sup>						100	26.3 [152]	22.8 [136]	22.4 [250]
EaAmsD <sup>i</sup>							100	19.8 [177]	25.2 [111]
StWaaG <sup>j</sup>								100	21.6 [180]
SpCpoA <sup>k</sup>									100

a. Percentage of identical amino acids determined with FASTA as implemented in PROSIS. Numbers in parentheses indicate the number of amino acids over which the % identity occurs.

b. *S. pneumoniae* Cps19cS

c. *Methanococcus jannaschii* RfbU (GenBank accession no. F64500)

d. *Aquifex aeolicus* MtfC (GenBank accession no. AE000693)

e. *Archaeoglobus fulgidus* galactosyl transferase (GenBank accession no. AE000983)

f. *S. aureus* CapM (Lin *et al.*, 1994)

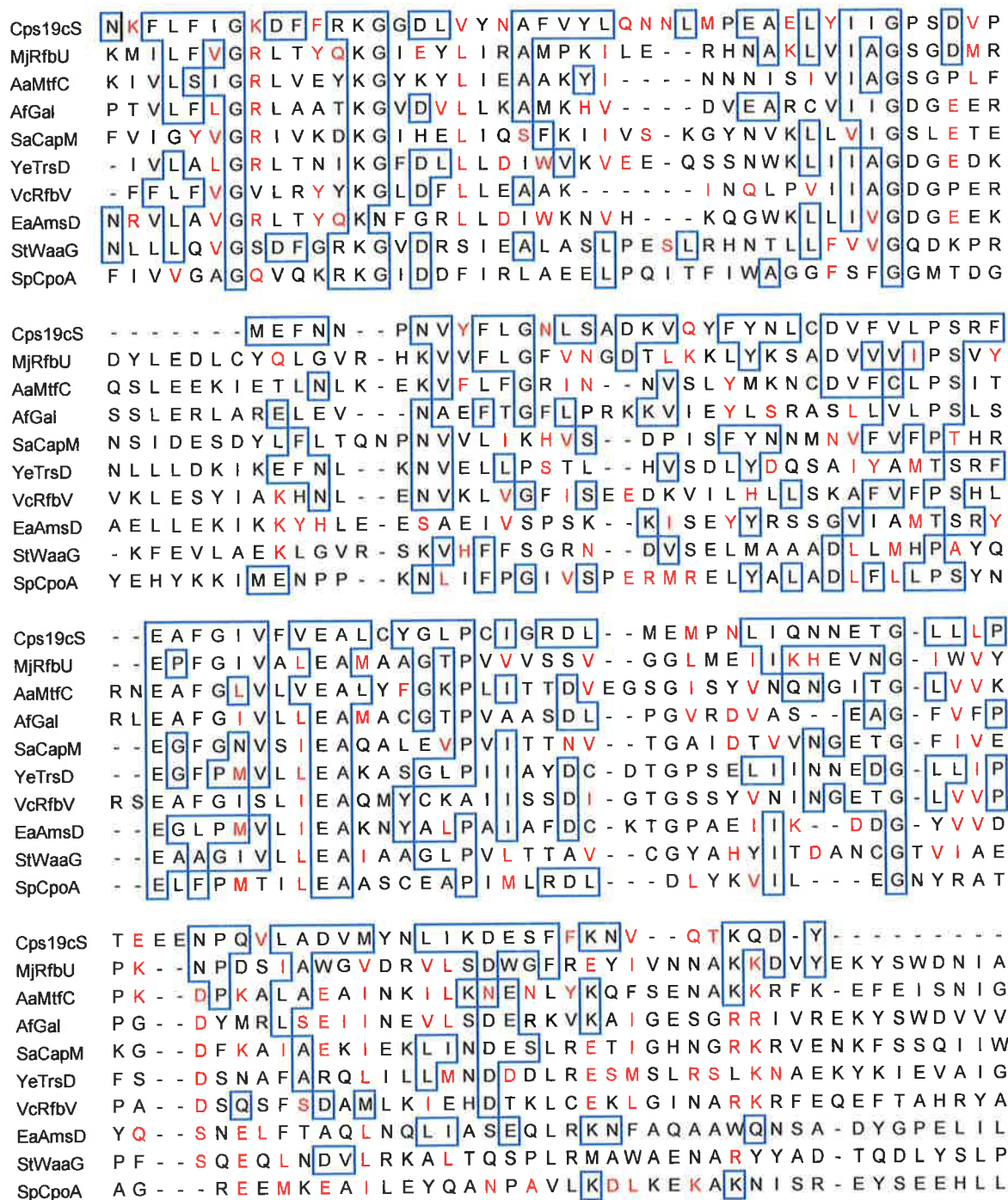
g. *Y. enterocolitica* TrsD (Skurnik *et al.*, 1995)

h. *Vibrio cholerae* RfbV (Fallarino *et al.*, 1997)

i. *E. amylovora* AmsD (Bugert and Geider, 1995)

j. *S. enterica* serovar typhimurium WaaG (Heinrichs *et al.*, 1998)

k. *S. pneumoniae* CpoA (Grebe *et al.*, 1997)



**Fig. 7.5. Alignment of Cps19cS.** Alignment of Cps19cS with *M. jannaschii* RfbU (MjRfbU) (GenBank accession no. F64500), *A. aeolicus* MtfC (AaMtfC) (GenBank accession no. AE000693), *A. fulgidus* Galtf (AfGal) (GenBank accession no. AE000983), *S. aureus* CapM (SaCapM) (Lin *et al.*, 1994), *Y. enterocolitica* TrsD (YeTrsD) (Skurnik *et al.*, 1995), *V. cholerae* RfbV (VcRfbV) (Fallarino *et al.*, 1997), *E. amylovora* AmsD (EaAmsD) (Bugert and Geider, 1995), and *S. enterica* serovar typhimurium WaaG (StWaaG) (Heinrichs *et al.*, 1998), as determined using the default settings of the program CLUSTAL. Residues identical to Cps19cS are boxed; similar residues are shown in red; - indicates absence of a residue.

### 7.2.3 Serotype specificity of *cps19cS*

To examine the relationship between *cps19cS* and encapsulation loci of other *S. pneumoniae* serotypes, a *SacI-HindIII* DNA fragment from a nested deletion derivative of pJCP484 corresponding to nucleotides 3,160-4,269 of the *cps19c* sequence was labelled with DIG and used to probe (at high stringency) Southern blots of restricted chromosomal DNA from representative pneumococci belonging to the following types/groups: 2, 3, 4, 6, 7F, 7B, 8, 9N, 9V, 12, 14, 16, 17, 18, 19F, 19A, 19B, 20, 22, 23F and 24. None of these serotypes had a high stringency homologue to *cps19cS* (result not shown). However, this is not surprising when the structures for their CPS are examined, because none contain a Glc side-chain with a  $\beta(1\rightarrow6)$  linkage (van Dam *et al*, 1990).

### 7.2.4 Transformation of *S. pneumoniae* type 19B to type 19C

In section 5.2.4, capsule production was altered from type 19F to type 19B by replacing *cps19fIJ* with the central region of *cps19b*, which contains the *cps19bPIQRJ* genes, and determines the 19B serotype. A similar approach was taken to determine whether *cps19cS* is indeed the gene responsible for the additional Glc side-chain which distinguishes type 19C CPS. A large PCR product of the *cps19c* region between *cps19cF* and *aliA* was amplified using primers J5 and J36 (Table 2.4) and transformed into Rx1-19F-I (an unencapsulated, Ery-resistant derivative of Rx1-19F, in which the *cps19fI* gene had been disrupted by insertion-duplication mutagenesis using pVA891, as described in section 3.2.4). The resultant transformant, expressing type 19C CPS, would be predicted to contain the *cps19cPIQRJ* genes required for both type 19B and 19C CPS biosynthesis as well as *cps19cS*. The *cps19cK* gene, which is located between *cps19cJ* and *cps19cL*, would also replace the almost identical *cps19fK* gene (94.9% identity). However, the

encoded UDP-GlcNAc-2-epimerase, while essential for CPS biosynthesis in all group 19 members, is not serotype determining.

A smooth transformant was checked for Ery sensitivity, indicating loss of the pVA891 sequence. Southern hybridisation was used to confirm the absence of both pVA891 and the *cps19fI* gene, and the presence of each of the *cps19cP*, *J*, and *S* genes (Fig. 7.6). The presence of *cps19cP* and *J* genes in the Rx1-19C transformant were confirmed using the *cps19bP* and *cps19bJ* gene probes because these *cps19b* genes are >95% identical to those in *cps19c*. The production of a type 19C capsule by the transformant, designated Rx1-19C, was then confirmed by quellung reaction. This shows that it is possible to alter capsule production from type 19F to type 19C by replacing *cps19fIJ* with the *cps19cPIQRJ* genes (required for both type 19B and 19C CPS biosynthesis) and the *cps19cS* gene, which determines the 19C serotype.



**Fig. 7.6. Southern hybridisation of Rx1-19F, 19B, Rx1-19C and 19C.** Rx1-19F (1), 19B (2), Rx1-19C (3), and 19C (4) were probed with DIG-labelled probes specific for *cps19bP* (A), *cps19bJ* (B), *cps19cS* (C) and *cps19fI* (D).

### 7.2.5 Characterisation of IS19C, located in the 5' intergenic region

Sequence analysis of the 5' intergenic region between *dexB* and the *cps19c* locus (section 6.2.5), indicated the presence of an additional insertion sequence designated

*IS19C*. The insertion point was within the inverted repeat of *IS1202*, just upstream of the *cps19c* locus, as shown in **Fig. 6.13**. This 1.2-kb IS element is flanked at both ends by a 13-bp direct repeat, followed by 14 bp of unique DNA and then a 14-bp inverted repeat as shown in **Fig. 7.7**. The ORF which encodes the putative transposase in *IS19C* lies in the same orientation as that for *IS1202* and opposite to the *cps19c* genes. This putative transposase has 67.5% identity to the transposase encoded by *IS1239* from *S. pyogenes*, but at the DNA sequence level these IS elements exhibit negligible similarity. The transposases from these two IS elements also share amino acid similarity to other transposases found in several different bacterial species including *IS30* from *E. coli*, as shown in **Table 7.3**.

```

1  TTCGAATGTGATATGAAAGTTCATAATGAAGTTAGCCACCTTACCTTAGTCAAGAATTAGATGTTTCACTATGTTTGAGTAAGTTGATGATTTTCATTG
                                     -35
101  ATAACAGGTTTGAACCTGTAGGCTAGGTGGCCAAGGCTAATCATAGCCTTGGTTTFAGCTGAAAAACAGGTTCAAGGGTTCCTGTTGTCAAATGAAATG
                                     -10                                     M Q E H Y T P K G K
201  TGATTTAAGGTATAAGAAAACACCTCTGTGCTATACTTGTGTTCCACCACAACACAAGGAAAGGCACAGAGATGCAAGAACATTATACCCCAAAAGGGA
                                     H L T I D N R R L I E R W K N E N K S N R E I A G L L G K A P Q T
301  AACATTTGACAATAGATAACCGTCGCTTGATTTGAGCGGTGGAAGAATGAAAATAAGTCCAATCGTGAATTCAGGCTTGTAGGAAAGGCGCCTCAAC
                                     I H N E V K R G T T L Q Q V R K G L Y K K V Y S A D Y A Q T V Y Q
401  GATTCATAATGAAGTCAAAAGAGGTACAACCTTACAACAAGTGAGAAAAGGGCTATACAAAAGGCTATTCTGCGGATTACGCACAACTGTTTACCAA
                                     F N R K R S V K K L I L T K E I R E K I L H Y H K Q K F S P E M M V
501  TTCAATCGAAAACGGTCGGTGAAAAGTTAATTTAACAAAGGAAATCAGAGAGAAGATCTTACACTATCATAAGCAAAAATTTCCCGCTGAAATGATGG
                                     N K K Q V K V G I S T I Y Y W F H N G H L R L T K A D M L Y P R K
601  TTAACAAGAAGCAAGTGAAGTTGGTATTTCAACCATCTACTACTGGTTTCATAATGGTCATTTAAGATTAACGAAGGCCGACATGCTTTATCCAGAAA
                                     R K G V K K Q A S P N F K P A G K S I E E R P D V I N L R L E N G
701  AAGGAAAGGTGTCAGAAGCAAGCTAGTCCGAACCTTAAGCCGGCAGGTAATCTATCGAAGAACGTCCTGACGTTATTAATCTTCGCTTGAAAATGGT
                                     H Y E I D T V L L T K I K N Y C L L V L T D R R S R H Q I I R L I P
801  CATTATGAAATTGATACCGTCTACTGACTAAGATAAAAAAATATTGCTGTTAGTCTTAACCGACCGGCGGAGCAGACACCAAAATTAAGGTTAATTC
                                     N K T A E S V N Q A L T L L L G E H R I L S I T A D N G S E F K R
901  CAAATAAAACTGCTGAATCTGTCAATCAGGCGCTTACGTTACTATTAGGGGAGCATCGTATTCTGTCCATTACTGCAGATAAATGGTTCGGAGTTCAAACG
                                     L S E V F P E E H I Y Y A H A Y S S W E R G S N E N H N R L I R R
1001  ATTGTCTGAGGTATTTCTTGAGGAACATATCTACTACGCACATGCTTACTCTTATGGGAGAGAGTTCAAATGAAATCATAATCGATTAATTCGGAGA
                                     W L P E G T K K T T P K E V A F I E N W I N N Y P K K C L D Y K S P
1101  TGTTTACCTGAAGGAACCAAGAAAACGACTCCGAAAGAAGTAGCTTTTATCGAAAATTTGGATTAACAACCTACCCTAAAAAATGCTTGGACTACAAGTCGC
                                     N E F L L G G *
1201  CAAATGAATTTCTTTGGTGGCTAACTTCAACTTGAATTTGGGTTTCGAATGTGATA

```

**Fig. 7.7. *IS19C*.** The sequence of *IS19C* as shown above corresponds to nucleotides 1,644-2,875 of the sequence available under GenBank accession no. AF105115. The translation of the putative transposase is shown in single letter code above the first nucleotide of each codon. The direct repeats are underlined and the inverted repeats are doubly underlined and are both shown in bold. A potential -35 and -10 promoter sequence for the putative *IS19C* transposase is underlined with a broken line. The ribosome binding site is underlined, and the amino acid translation is represented by single letter code above the first nucleotide of each codon; the stop codon is indicated by an asterisk.

**Table 7.3. Similarity of the putative transposase encoded by IS19C to other transposases.**

	% Identity <sup>a</sup>								
	IS19c	IS1239	IS1161	IS4351	IS1470	IS1070	IS1086	IS30	IS1394
IS19c <sup>b</sup>	100	67.5	36.2	33	31.5	31.1	28.8	28	26.9
		[314]	[232]	[318]	[305]	[254]	[326]	[321]	[308]
IS1239 <sup>c</sup>		100	29	31.5	31.8	30.7	28.8	26.2	29
			[314]	[321]	[333]	[241]	[326]	[321]	[328]
IS1161 <sup>d</sup>			100	27.7	36.3	34.5	26.5	28.7	28.1
				[328]	[344]	[203]	[336]	[328]	[320]
IS4351 <sup>e</sup>				100	41.2	30.7	30.5	31.6	33
					[170]	[313]	[325]	[326]	[324]
IS1470 <sup>f</sup>					100	31.5	26.8	30.8	29.2
						[302]	[340]	[334]	[332]
IS1070 <sup>g</sup>						100	30.8	28.6	33.3
							[208]	[280]	[204]
IS1086 <sup>h</sup>							100	41.1	58.3
								[326]	[326]
IS30 <sup>i</sup>								100	42.2
									[329]
IS1394 <sup>j</sup>									100

a. Percentage of identical amino acids determined with FASTA as implemented in PROSIS. Numbers in parentheses indicate the number of amino acids over which the % identity occurs.

b. *S. pneumoniae* IS19C putative transposase

c. *S. pyogenes* IS1239 putative transposase (Kapur *et al.*, 1994)

d. *Streptococcus salivarius* IS1161 putative transposase (Giffard *et al.*, 1993)

e. *Bacteroides fragilis* IS4351 putative transposase (Smith, 1987)

f. *Clostridium perfringus* IS1470 putative transposase (Brynestad *et al.*, 1997)

g. *Leuconostoc lactis* IS1070 putative transposase (Vaughan and de Vos, 1995)

h. *Ralstonia eutropha* IS1086 putative transposase (Dong *et al.*, 1992)

i. *E. coli* IS30 transposase (Dalrymple *et al.*, 1984)

j. *Pseudomonas alcaligenes* IS1394 putative transposase (GenBank accession no. U37284)

## 7.3 Conclusions

The *cps19c* locus is almost identical to the *cps19b* locus except that an extra gene (*cps19cS*) has inserted between *cps19cK* and *cps19cL*. This gene is most likely to encode the glucosyl transferase required for the addition of the Glc side-chain in the type 19C repeat unit. Interestingly, all three putative transferases involved in the addition of side-chains to type 19B and/or 19C CPS, Cps19cS, Cps19P and Cps19Q, appear to be cytoplasmic enzymes, as they lack both a leader sequence for export to the cell surface and

a hydrophobic transmembrane sequence which could anchor them to the cell membrane. Thus, the Rha-Rib disaccharide side-chain present in both type 19B and 19C CPS, and the Glc side-chain specific to type 19C CPS are most probably added to the repeat units in the cytoplasm, before translocation to the outer surface by Cps19J and subsequent polymerisation by Cps19I. It is interesting that the Glc side-chain does not appear to interfere with the function of either the repeat unit transporter (Cps19bJ and Cps19cJ) or the polysaccharide polymerase (Cps19bI and Cps19cI), as the proteins encoded by *cps19b* and *cps19c* are almost identical (99.7% and >95%, respectively) and are able to function in the biosynthesis of both the type 19B and type 19C CPS.

## Chapter 8

# DISCUSSION

## 8.1 Introduction

*S. pneumoniae* group 19 is the first group for which the *cps* loci from all the members (19F, 19A, 19B and 19C) have been completely characterised. Functions have been assigned to the majority of the *cps19* gene products, based on either gene complementation or similarity to other proteins with known functions. The ability of PCR products containing either complete or partial *cps* loci to transform pneumococci from one serotype to another demonstrated that the *cps19* loci contain sufficient genetic information for expression of type-specific CPSs.

Comparison of the *cps* loci from all four members of group 19 and to the *cps* loci from other pneumococcal serotypes will be discussed below.

## 8.2 Analysis of the *S. pneumoniae* group 19 *cps* loci

### 8.2.1 Comparison of the group 19 *cps* loci

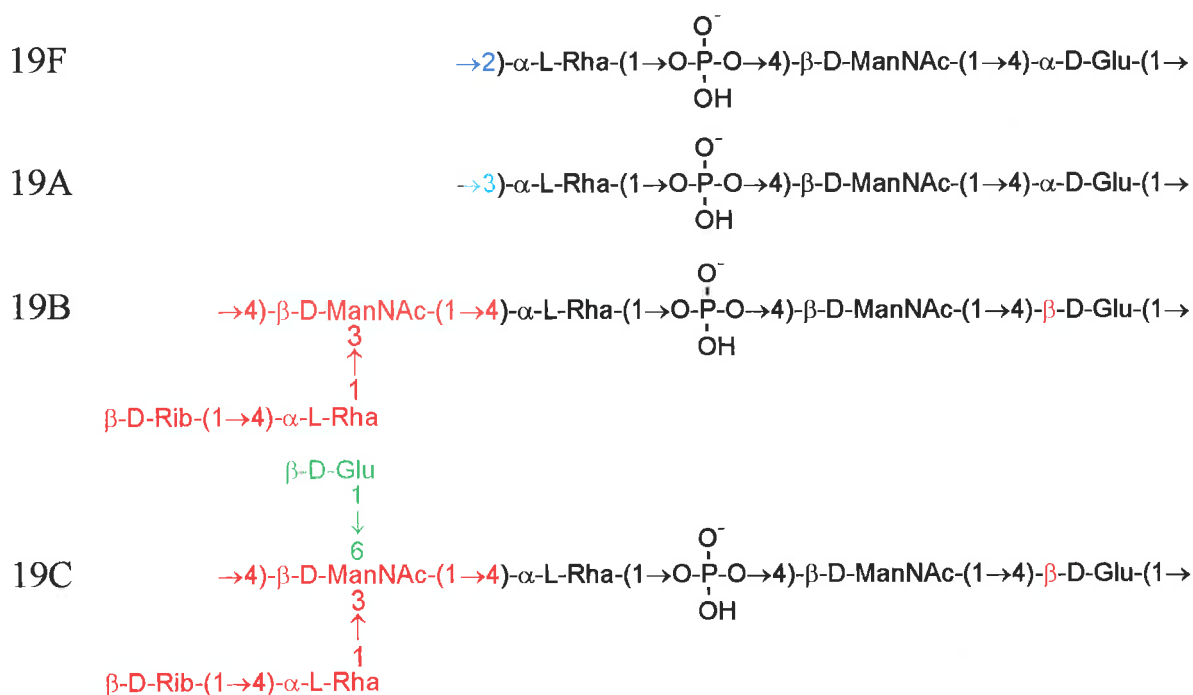
The structural similarities between the CPS repeat units from all four members of

serogroup 19 (19F, 19A, 19B and 19C) are reflected in the highly conserved arrangement of their *cps* loci, with 13 genes (*cps19A-H*, and *K-O*) common to all four serogroup members, as shown in **Fig. 8.1**. Nearly all of the common genes from types 19F, 19B and 19C are >95% identical to each other, whereas those from the type 19A *cps* locus are more divergent.

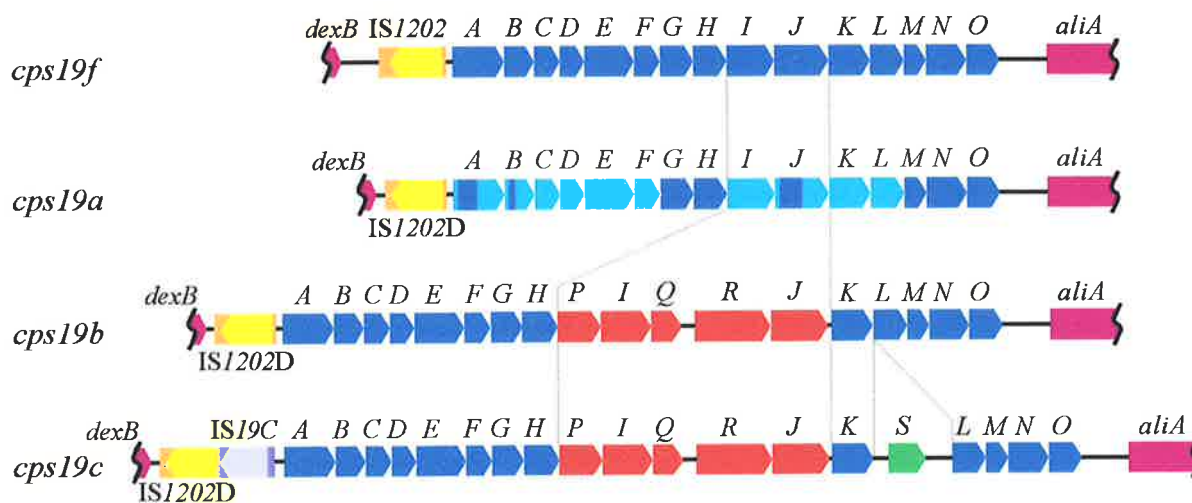
Although the *cps19a* and *cps19f* loci are identical in the number and arrangement of the genes present, the similarity between individual genes varies from 70% to 99% identity (for both the nucleotide and the deduced amino acid sequences). This sequence divergence is surprising given that the only difference between their CPS repeat units is the glycosidic linkage which joins the repeat units together ( $\alpha(1\rightarrow2)$  for 19F and  $\alpha(1\rightarrow3)$  for 19A) (**Fig. 8.1A**). Hypothetically, only a difference in the *cps19aI* gene, which presumably encodes the polysaccharide polymerase responsible for this linkage, is required to change a type 19F pneumococcus into type 19A. The results of the present study are consistent with this hypothesis. A transformation event in which the region of the *cps19a* locus encoding most of Cps19aH, all of Cps19aI and the first 76 amino acids of Cps19aJ replaced the homologous portion of the *cps19f* locus was sufficient to convert CPS type from 19F to 19A (section 6.2.4). Given that Cps19fH and Cps19aH are >95% identical, it seems probable that Cps19aI is solely responsible for the observed alteration in CPS type.

The degree of divergence between the type 19F and 19A loci suggests that they may have diverged early in their evolutionary past, or that their component genes originated from different sources. The latter alternative is supported by the fact that the G+C content of the *cps19a* locus is higher than that for *cps19f* (section 8.2.2). In chapters 4 and 6, sequence variations between the *cpsCDE* genes among different pneumococcal serotypes identified two distinct classes of *cps* loci. Type 19A was recognised as a class II locus whereas the other members of group 19 belonged to class I. Thus, the marked difference

A.



B.



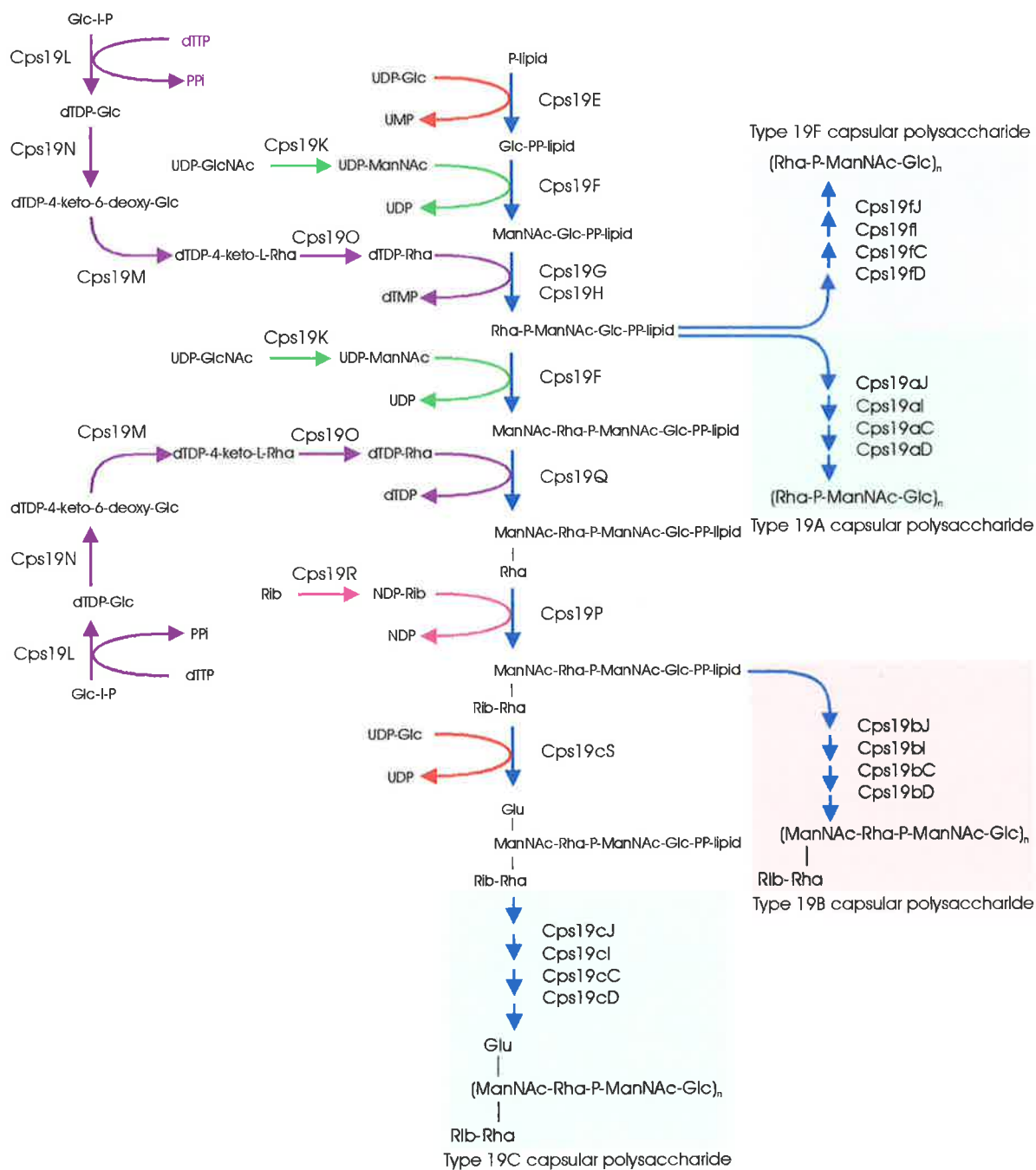
**Fig. 8.1 Comparison of the CPS structures and the *cps* loci of *S. pneumoniae* group 19.** **A. CPS repeat unit structures.** The part of the CPS repeat unit structure shown in black is common to all group 19 members. The differences in their CPS structures are highlighted in different colours: the different repeat unit linkages are shown in dark blue for 19F and light blue for 19A, the additional sugars common to 19B and 19C are shown in red, and the additional glucose side-chain in 19C is shown in green. **B. The *cps* loci.** The genes of the *cps19f* locus are shown in dark blue, and the *cps19f* homologues from other group 19 members with >90% identity are also shown in dark blue. The *cps19a* genes which have 70-90% identity to the *cps19f* homologue are shown in light blue. The 19B and 19C specific genes are shown in red and the additional 19C specific gene is shown in green.

between the genes at the 5' end of the type 19A *cps* locus and homologues in other members of group 19 is almost certainly a consequence of recombination between portions of class I and class II loci, rather than divergence over a long period of time.

Closely related homologues of *cps19fl* and *J*, which encode the type 19F polysaccharide polymerase and repeat unit transporter, respectively, are not found in the type 19B *cps* locus. In type 19B this region of the *cps* locus (between *cps19bH* and *cps19bK*) contains five genes (**Fig. 8.1B**), encoding an unrelated polymerase and repeat unit transporter, as well as two additional putative glycosyl transferases and a protein which may be involved in synthesis of an activated Rib precursor (section 5.2.2.4). Transformation studies indicated that these five genes encode all of the functions required to convert a type 19F pneumococcus to type 19B (section 5.2.4).

The type 19C locus is virtually identical to the type 19B locus, with the *cps19c* genes exhibiting >95% identity to all 18 of the *cps19b* genes, differing only in the insertion of a glucosyl transferase gene (*cps19cS*) between *cps19cK* and *cps19cL*. Transformation studies have shown that the presence of this gene accounts for the additional Glc side chain in the otherwise identical repeat unit structures (**Fig. 8.1A**). The type 19C *cps* locus contains 19 genes, and at 21 kb, it is the largest pneumococcal capsule gene cluster characterised to date.

The 13 genes common to all members of group 19 encode functions required for the synthesis of the shared trisaccharide component of the group 19 CPS structures. Furthermore, the genetic differences between the group 19 *cps* loci identified are also consistent with the differences in the CPS structures of individual serotypes. This information has been used to propose biosynthetic pathways for each of the serotypes as shown in **Fig. 8.2**. Only the functions of Cps19A and Cps19B are still unknown, and hence, are not included in **Fig. 8.2**.



**Fig. 8.2.** The putative biosynthetic pathways for 19F, 19A, 19B and 19C.

### 8.2.2 Distribution of the *cps19* genes in other pneumococcal serotypes

When the distribution of the *cps19* genes in other serotypes was examined, homologues of the first four genes, *cpsA-D*, were identified in the *cps* loci of all serotypes tested to date, and furthermore, the *cpsC* and *cpsD* genes existed as two distinct classes, designated class I and class II. The *cpsE* gene, encoding the glucosyl transferase which adds Glc-1-phosphate to the lipid carrier, is also conserved as two distinct alleles in the class I and class II *cps* loci of all serotypes tested which contain Glc in their CPS (except type 3). Thus, pneumococcal *cps* loci can be segregated into class I and class II loci which may have evolved from two distinct clonal ancestors, based on the degree of nucleotide and amino acid similarity between them. Whereas, *cpsC-E* genes in the same class share >95% identity at both the DNA and amino acid level, they share only 70-80% identity to genes belonging to the other class. The precise point at which these sequences diverge into the two distinct classes appears to vary between serotypes but probably occurs within *cpsB* in most serotypes.

Homologues of the four genes involved in dTDP-Rha biosynthesis (*cps19L-O*) were also present in all serotypes tested that contained Rha in their CPS. A homologue to *cps19fF*, a putative ManNAc transferase responsible for linking ManNAc via a  $\beta(1\rightarrow4)$  linkage to Glc in the type 19 CPS, was also present in types 9N and 9V. This was not surprising as their CPS repeat units also contain this same glycosidic linkage. The remainder of the *cps19* genes were all specific to one or more members of group 19.

### 8.2.3 G+C content

The G+C content of the individual genes in the *cps19* locus is quite variable, ranging from 27-43%, with genes of similar G+C content clustered together in a mosaic structure. The G+C content of some of the *cps* genes are atypical for pneumococci and this

suggests that *S. pneumoniae* may have acquired these genes from a different bacterial source. The G+C content of the individual genes is shown in **Table 8.1** and is represented diagrammatically in **Fig 8.3**. Interestingly, the G+C content of most of the genes in the *cps19a* locus are consistently higher (up to 6.5%) than that for the respective *cps19f* gene, suggesting a distinct ancestral origin for these two loci.

Also of interest is the low G+C content of the type 19B-specific gene cluster (*cps19bPIQRJ*), which ranges from 27.2-29.7% for the individual genes (**Table 8.1**). This is comparable to the G+C content for the *cps19fIJ* gene cluster of 29.7%, but differs considerably from the remainder of either of the loci, which ranges from 30.3-42.3% G+C in the *cps19f* locus (**Table 8.1**). Conversely, the four type 19F genes encoding dTDP-Rha biosynthesis (*cps19fLMNO*) have G+C contents ranging from 41.5-42.3%, which is significantly higher than that of flanking sequences and for pneumococcal genes in general, the average G+C content of which have been estimated at 35.1% (Bridge and Sneath, 1983). This is consistent with acquisition of these different gene clusters from distinct sources.

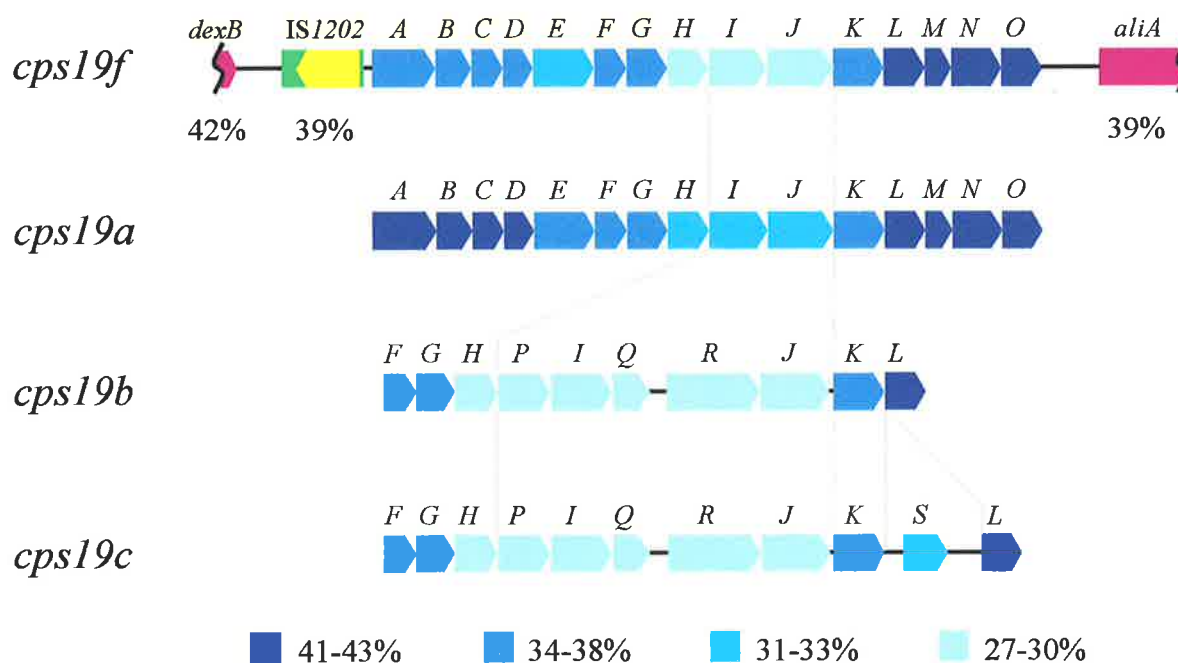
#### **8.2.4 Transcription and translation of the *cps19* loci**

All the *cps19* loci examined contained an identical consensus  $\sigma^{70}$  promoter (TAGACA-17 bp-TATAAT) just upstream of the first ORF, as described for *cps19f* (Guidolin et al., 1994). All *cps19* ORFs identified were preceded by a ribosome binding site and are closely linked with most intergenic gaps being less than 100 nucleotides. Even the four intergenic gaps larger than 100 nucleotides did not appear to contain any sequences which resemble either transcription terminators or promoter sequences. The only stemmed-loop structures likely to function as transcription terminators are found downstream of the *cps19* loci in the 3' intergenic region (section 3.2.1). Thus the entire

**Table 8.1. The G+C content of the *cps19* genes.**

Gene	%G+C			
	19F	19A	19B	19C
<i>cps19A</i>	38.1	39.5	NS	NS
<i>cps19B</i>	38	41.3	NS	NS
<i>cps19C</i>	38.2	42.1	NS	NS
<i>cps19D</i>	34.5	41	NS	NS
<i>cps19E</i>	33.2	37.7	NS	NS
<i>cps19F</i>	33.6	34.1	34.2*	NS
<i>cps19G</i>	36.3	37.2	36.6	NS
<i>cps19H</i>	30.3	32.2	30.2	NS
<i>cps19I</i>	29.7	32.7	27.2	NS
<i>cps19J</i>	29.7	33	29.5	29.9
<i>cps19K</i>	35.2	36.9	35.2	35.3
<i>cps19L</i>	42.3	43.3	41.9*	42.6*
<i>cps19M</i>	41.5	41.2	NS	NS
<i>cps19N</i>	42.1	42.4	NS	NS
<i>cps19O</i>	41.5	41.3	NS	NS
<i>cps19P</i>	-	-	29.5	NS
<i>cps19Q</i>	-	-	29.7	NS
<i>cps19R</i>	-	-	27.2	26.4*
<i>cps19cS</i>	-	-	-	31.4

NS, denotes not sequenced; "-", denotes absence of gene in that serotype; "\*" denotes only part of the gene was sequenced.



**Fig. 8.3. Diagrammatic representation of the G+C content of the group 19 *cps* loci.** Individual *cps* genes are shaded in accordance with their approximate G+C content. For flanking genes, the % G+C is shown on the map. Portions of the *cps19b* and *cps19c* loci are not shown, as they are >95% identical to *cps19f*.

loci are probably co-transcribed in each of the four serotypes.

Although most of the *cps19* genes used an ATG initiation codon, a few of the genes were predicted to use a TTG initiation codon. For one of these (*cps19fH*) use of a TTG initiation codon was confirmed by T7 expression studies in *E. coli* (section 3.2.5). The *cps19H* homologues in the *cps* loci of the other members of group 19 also have TTG initiation codons, as does *cps19aG*.

### 8.3 The *cps* loci from other *S. pneumoniae* serotypes

Since the publication of the first six genes of the *cps19f* locus (Guidolin *et al.*, 1994) sequence data for *cps* loci from *S. pneumoniae* serotypes 1, 3, 14 and 23F have been published (Arrecubieta *et al.*, 1995; Dillard *et al.*, 1995; Kolkman *et al.*, 1997b; Muñoz *et al.*, 1997; Ramirez and Tomasz, 1998). The sequence of the type 4 *cps* locus has also become available by analysis of the partial genome sequence for *S. pneumoniae* type 4 released by the TIGR Microbial Database ([http://www.tigr.org/pub/data/s\\_pneumoniae/](http://www.tigr.org/pub/data/s_pneumoniae/)). The *cps* loci from a different *S. pneumoniae* type 23F strain and from types 2 and 33F also have recently been determined and are either submitted for publication or are currently in press (Iannelli *et al.*, 1998; Llull *et al.*, 1998; Morona *et al.*, 1998). An overview of the genes present in these loci is provided below to enable a comparison with those from the group 19 *cps* loci.

#### 8.3.1 The *S. pneumoniae* type 1 *cap* locus

The *S. pneumoniae* type 1 *cps* locus (Muñoz *et al.*, 1997) has 11 ORFs, designated *cap1A-K*, and the predicted functions of their protein products are listed in **Table 8.2**. The

*cps* locus is flanked on both sides by IS1167 and redundant copies of the genes required for dTDP-Rha biosynthesis (*orf1-4*) were found downstream of the 3' copy of IS1167. The functions of the protein products encoded by both *cap1H* and *cap1I* have been deduced by examination of their hydropathy profiles (data not shown).

**Table 8.2. Predicted functions of *S. pneumoniae cap1*<sup>a</sup> ORFs**

ORF	Predicted or known function	Similar proteins (% identity)
<i>cap1A</i>	regulation?	<i>S. pneumoniae</i> Cps19fA (>98%)
<i>cap1B</i>	unknown	<i>S. pneumoniae</i> Cps19fB (>98%)
<i>cap1C</i>	chain length regulation/export	<i>S. pneumoniae</i> Cps19fC (>98%)
<i>cap1D</i>	chain length regulation/export	<i>S. pneumoniae</i> Cps19fD (>98%)
<i>cap1E</i>	galacturonic acid transferase	<i>Proteus mirabilis</i> GpsF <sup>b</sup> (33.5%) <i>S. dysenteriae</i> Rfp <sup>c</sup> (29.1%)
<i>cap1F</i>	acetyl transferase	<i>Rhizobium leguminosarium</i> NodL <sup>d</sup> (30.3%) <i>E. coli</i> K12 LacA <sup>e</sup> (29.7%)
<i>cap1G</i>	glycosyl transferase	<i>S. aureus</i> CapJ <sup>f</sup> (30.1%) <i>E. coli</i> K12 WcaL <sup>g</sup> (19.5%)
<i>cap1H</i>	polysaccharide polymerase	-
<i>cap1I</i>	repeat unit transporter	<i>M. jannaschii</i> MJ1068 <sup>h</sup> (22.7%) <i>Shigella sonnei</i> form 1 ORF3 <sup>i</sup> (22.3%)
<i>cap1J</i>	galacturonic acid epimerase	<i>V. cholerae</i> Orf9 <sup>j</sup> (53.2%) <i>E. coli</i> K12 GalE <sup>k</sup> (27.5%)
<i>cap1K</i>	UDP-glucose dehydrogenase	<i>S. pneumoniae</i> Cap3A <sup>l</sup> (61.7%)
<i>orf1</i>	glucose-1-phosphate thymidyl transferase	<i>S. pneumoniae</i> Cps19fL (96.1%)
<i>orf2</i>	dTDP-4-keto-6-deoxyglucose-3,5-epimerase	<i>S. pneumoniae</i> Cps19fM (98%)
<i>orf3</i>	dTDP-glucose-4,6-dehydratase	<i>S. pneumoniae</i> Cps19fN (98.8%)
<i>orf4</i>	dTDP-L-rhamnose synthase	<i>S. pneumoniae</i> Cps19fO (97.5%)

The sequences can be found under GenBank accession nos: a, Z83335; b, L36873; c, S27671; d, X17557; e, X51872; f, U10927; g, X56793; h, U67549; i, U34305; j, U47057; k, X06226; l, Z47210.

### 8.3.2 The *S. pneumoniae* type 2 *cps* locus

The *S. pneumoniae* type 2 *cps* locus has recently been sequenced by Iannelli *et al.* (1998), and consists of 17 ORFs, designated *cps2A-O*, *P* and *T*, and the predicted functions

of their protein products are listed in **Table 8.3**. Interestingly, the *cps2F* gene product has significant similarity to the experimentally determined galactosyl transferases, Cps14J and Cps14I, from *S. pneumoniae* type 14, as well as the putative rhamnosyl transferase RgpBc from *S. mutans*. As type 2 CPS does not contain Gal, Cps2F is probably one of the rhamnosyl transferases required for type 2 CPS biosynthesis.

**Table 8.3 Predicted functions of *S. pneumoniae* *cps2*<sup>a</sup> ORFs**

ORF	Predicted function	Similar proteins (% identity)
<i>cps2A</i>	regulation?	<i>S. pneumoniae</i> Cps19fA (96%)
<i>cps2B</i>	unknown	<i>S. pneumoniae</i> Cps19fB (84%)
<i>cps2C</i>	chain length regulation/export	<i>S. pneumoniae</i> Cps19fC (67%)
<i>cps2D</i>	chain length regulation/export	<i>S. pneumoniae</i> Cps19fD (78%)
<i>cps2E</i>	glucosyl-1-phosphate transferase	<i>S. pneumoniae</i> Cps23fE <sup>b</sup> (95%) <i>S. pneumoniae</i> Cps14E <sup>c</sup> (60%) <i>S. pneumoniae</i> Cps19fE (60%)
<i>cps2T</i>	rhamnosyl transferase	<i>S. pneumoniae</i> Cps23fT <sup>b</sup> (81%)
<i>cps2F</i>	rhamnosyl transferase	<i>S. mutans</i> RgpBc <sup>d</sup> (33%) <i>S. pneumoniae</i> Cps23fU <sup>b</sup> (23%) <i>S. pneumoniae</i> Cps14J <sup>c</sup> (19%) <i>S. pneumoniae</i> Cps14I <sup>c</sup> (18%)
<i>cps2G</i>	glycosyl transferase	<i>S. pneumoniae</i> RPN00103 <sup>e</sup> (35%) <i>S. pneumoniae</i> RPN00642 <sup>e</sup> (20%) <i>N. meningitidis</i> IcsA <sup>f</sup> (19%)
<i>cps2H</i>	polysaccharide polymerase	<i>S. pneumoniae</i> Cps14H <sup>c</sup> (21%) <i>S. pneumoniae</i> Cps19fI (20%) <i>S. flexneri</i> Rfc <sup>g</sup> (22%)
<i>cps2I</i>	glycosyl transferase	<i>N. meningitidis</i> RfaK <sup>h</sup> (25%) <i>N. meningitidis</i> IcsA <sup>f</sup> (24%)
<i>cps2J</i>	repeat unit transporter	<i>S. pneumoniae</i> Cps23fJ <sup>b</sup> (36%)
<i>cps2K</i>	UDP-glucose dehydrogenase	<i>S. pneumoniae</i> Cap1K <sup>i</sup> (89%) <i>S. pneumoniae</i> Cps3D <sup>j</sup> (Cap3A <sup>k</sup> ) (74%)
<i>cps2P</i>	UDP-galactopyranose mutase	<i>E. coli</i> Glf <sup>l</sup> (67%) <i>Mycobacterium tuberculosis</i> Glf <sup>m</sup> (58%)
<i>cps2L</i>	glucose-1-phosphate thymidyl transferase	<i>S. pneumoniae</i> Cps19fL (91%)
<i>cps2M</i>	dTDP-4-keto-6-deoxyglucose-3,5-epimerase	<i>S. pneumoniae</i> Cps19fM (96%)
<i>cps2N</i>	dTDP-glucose-4,6-dehydratase	<i>S. pneumoniae</i> Cps19fN (99%)
<i>cps2O</i>	dTDP-L-rhamnose synthase	<i>S. pneumoniae</i> Cps19fO (99%)

The sequences can be found under GenBank accession nos: a, AF026471; b, AF030373; c, X85787; d, AB010970; e, GenBank unfinished genomes (TIGR Genomic Database); f, U39810; g, X71970; h, U58765; i, Z83335; j, U15171; k, Z47210; l, U09876 ; m, U96128.

### 8.3.3 The *S. pneumoniae* type 3 *cps* locus

The sequence of the *S. pneumoniae* type 3 *cps3/cap3* locus contains six ORFs, only four of which are transcribed. These four genes, have been designated *cps3D*, *S*, *U* and *M* by Dillard *et al.* (1995) and *cap3B*, *A*, *C* and *D* by Arrecubieta *et al.* (1995). The predicted functions of their protein products are listed in **Table 8.4**.

**Table 8.4** Predicted functions of *S. pneumoniae* *cps3<sup>a</sup>/cap3<sup>b</sup>* ORFs

ORF	Predicted or known function	Similar proteins (% identity)
<i>cps3D/cap3A</i>	UDP-glucose dehydrogenase	<i>S. pyogenes</i> HasB <sup>c</sup> (57%) <i>E. coli</i> Ugd <sup>d</sup> (55.3%)
<i>cps3S/cap3B</i>	type 3 synthase	<i>S. pyogenes</i> HasA <sup>e</sup> (25.8%) <i>R. meliloti</i> NodC <sup>f</sup> (21%) <i>S. cerviciae</i> Csd2 <sup>g</sup> (24.7%)
<i>cps3U/cap3C</i>	glucose-1-phosphate uridylyltransferase	<i>S. pyogenes</i> HasC <sup>h</sup> (78%) <i>B. subtilis</i> GtaB <sup>i</sup> (55%)
<i>cps3M/orf3</i>	phosphoglucomutase	<i>B. subtilis</i> YhxB <sup>j</sup> (38.6%) <i>E. coli</i> PgmU <sup>k</sup> (22.6%)

The sequences can be found under GenBank accession nos: a, U66846, U15171, and U66845; b, Z47210; c, L08444; d, U90519; e, L21187; f, X01649; g, M73697; h, Q54713; i, Z99122; j, Y14079; k, AE000172.

### 8.3.4 The *S. pneumoniae* type 4 *cps* locus

The sequence of the *S. pneumoniae* type 4 *cps* locus (located on contig no. 4108 of the partial *S. pneumoniae* genome sequence available from the TIGR Microbial Database) was examined and the *cps4* locus, located between *dexB* and *aliA* in the *S. pneumoniae* chromosome, was found to have 15 ORFs, designated *cps4A-O*. The predicted functions of the protein products encoded by *cps4A-O*, based on homologies to other proteins in the GenBank database, are listed in **Table 8.5**, and are sufficient for biosynthesis of type 4 CPS repeat unit which is shown in **Table 1.2**.

**Table 8.5 Predicted functions of *S. pneumoniae* cps4 ORFs**

ORF	Predicted function	Similar proteins (% identity)
<i>cps4A</i>	regulation?	<i>S. pneumoniae</i> Cps19fA (94.2%)
<i>cps4B</i>	unknown	<i>S. pneumoniae</i> Cps19fB (86.8%)
<i>cps4C</i>	chain length regulation/export	<i>S. pneumoniae</i> Cps19fC (86.5%)
<i>cps4D</i>	chain length regulation/export	<i>S. pneumoniae</i> Cps19fD (94.3%)
<i>cps4E</i>	galactosyl-1-phosphate transferase	<i>S. aureus</i> Cap8M <sup>a</sup> (48.6%) <i>S. aureus</i> Cap5M <sup>b</sup> (49.2%) <i>S. enterica</i> serovar typhimurium RfbP <sup>c</sup> (36.5%)
<i>cps4F</i>	N-acetyl fucosamine transferase	<i>S. aureus</i> Cap8L <sup>a</sup> (23.9%) <i>S. aureus</i> Cap5L <sup>b</sup> (23.9%) <i>E. coli</i> K12 WcaI <sup>d</sup> (22.6%)
<i>cps4G</i>	N-acetyl mannosamine transferase	<i>S. aureus</i> Cap8H <sup>a</sup> (32%)
<i>cps4H</i>	N-acetyl galactosamine transferase	<i>Streptococcus thermophilus</i> EpsF <sup>e</sup> (25.8%) <i>H. influenzae</i> lsg Orf3 <sup>f</sup> (19.8%)
<i>cps4I</i>	polysaccharide polymerase	-
<i>cps4J</i>	pyruvyl transferase	<i>E. coli</i> K12 WcaK <sup>d</sup> (17.8%)
<i>cps4K</i>	repeat unit transporter	<i>E. coli</i> K12 RfbX <sup>g</sup> (20.9%) <i>Y. enterocolitica</i> TrsA <sup>h</sup> (19.1%) <i>S. enterica</i> serovar typhimurium RfbX <sup>i</sup> (20.1%)
<i>cps4L</i>	UDP-N-acetyl glucosamine-2-epimerase	<i>B. subtilis</i> YvyH <sup>j</sup> (63.2%) <i>S. aureus</i> Cap5P <sup>b</sup> (57.2%) <i>S. aureus</i> Cap8P <sup>a</sup> (57.2%) <i>E. coli</i> K12 RffE <sup>k</sup> (53%)
<i>cps4M</i>	UDP-N-acetyl galactosamine dehydratase (UDP-FucNAc synthesis)	<i>S. aureus</i> Cap8E <sup>a</sup> (61.3%) <i>S. aureus</i> Cap5E <sup>b</sup> (61%) <i>Y. enterocolitica</i> TrsG <sup>l</sup> (33.4%)
<i>cps4N</i>	UDP-N-acetyl fucosamine synthase (UDP-FucNAc synthesis)	<i>S. aureus</i> Cap8F <sup>a</sup> (45.6%) <i>S. aureus</i> Cap5F <sup>b</sup> (46.2%) <i>S. enterica</i> serovar typhimurium RfbE <sup>c</sup> (25%)
<i>cps4O</i>	UDP-N-acetyl glucosamine-4-epimerase (UDP-FucNAc synthesis)	<i>S. aureus</i> Cap8G <sup>a</sup> (58.9%) <i>S. aureus</i> Cap5G <sup>b</sup> (58.9%) <i>B. subtilis</i> YvyH <sup>j</sup> (25.3%)

The sequences can be found under GenBank accession nos: a, U73374; b, U81973; c, X56793; d, U38473; e, U40830; f, M94855; g, AF013583; h, Z47767; i, X60665; j, Z99122, k, L18799; l, Z47767.

### 8.3.5 The *S. pneumoniae* type 14 cps locus

The *S. pneumoniae* type 14 cps locus (Kolkman *et al.*, 1997b) consists of 12 ORFs, designated *cps14A-L* and the predicted functions of their protein products are listed in **Table 8.6**. The functions of the glycosyl transferases have been experimentally determined.

**Table 8.6. Predicted functions of *S. pneumoniae* *cps14*<sup>a</sup> ORFs**

ORF	Predicted or known function	Similar proteins (% identity)
<i>cps14A</i>	regulation?	<i>S. pneumoniae</i> Cps19fA (96.9%)
<i>cps14B</i>	unknown	<i>S. pneumoniae</i> Cps19fB (96.3%)
<i>cps14C</i>	chain length regulation/export	<i>S. pneumoniae</i> Cps19fC (92.6%)
<i>cps14D</i>	chain length regulation/export	<i>S. pneumoniae</i> Cps19fD (95.2%)
<i>cps14E</i>	glucosyl-1-phosphate transferase	<i>S. pneumoniae</i> Cps19fE (95.8%) <i>S. enterica</i> serovar typhimurium RfbP <sup>b</sup> (31.3%)
<i>cps14F</i>	transferase co-factor	<i>Sphingomonas</i> S88 SpsK <sup>c</sup> (33.8% to N-terminal)
<i>cps14G</i>	$\beta$ -1,4-galactosyl transferase	<i>Sphingomonas</i> S88 SpsK <sup>c</sup> (27.7% to C-terminal)
<i>cps14H</i>	polysaccharide polymerase	<i>S. flexneri</i> Rfc <sup>d</sup> (25.6%) <i>S. enterica</i> serovar typhimurium Rfc <sup>b</sup> (21%)
<i>cps14I</i>	$\beta$ -1,3-N-acetylglucosaminyl transferase	<i>R. meliloti</i> ExoO <sup>e</sup> (26.6%) <i>S. thermophilus</i> EpsI <sup>f</sup> (26.1%)
<i>cps14J</i>	$\beta$ -1,4-galactosyl transferase	<i>S. thermophilus</i> EpsI <sup>f</sup> (43.9%) <i>Y. enterocolitica</i> TrsB <sup>g</sup> (26.9%)
<i>cps14K</i>	unknown	<i>H. influenzae</i> type b Orf4 <sup>h</sup> (28.3%)
<i>cps14L</i>	repeat unit transporter	<i>S. dysenteriae</i> RfbX <sup>i</sup> (26.1%) <i>E. coli</i> K12 RfbX <sup>j</sup> (33%) <i>Y. enterocolitica</i> O:3 Wzx <sup>g</sup> (23.5%) <i>S. pneumoniae</i> Cap1I <sup>k</sup> (20.7%)
<i>orfX</i>	glycerol-phosphate transferase	<i>B. subtilis</i> TagF <sup>l</sup> (26.6%)

The sequences can be found under GenBank accession nos: a, X85785; b, X56793; c, U51197; d, X71970; e, Z22636; f, U40830; g, Z47767; h, X78559; i, L07293; j, AF013583; k, Z83335; l, Z99122.

### 8.3.6 The *S. pneumoniae* type 23F *cps* locus

The *S. pneumoniae* type 23F *cps* locus has recently been determined from two different isolates (Morona *et al.*, 1998; Ramirez and Tomasz, 1998). The two loci are almost identical but the gene designations differ. Whereas Ramirez and Tomasz (1998) have designated the 18 ORFs *cps23fA-R*, Morona *et al.* (1998) have designated the genes *cps23fA-E, I, J, L-O* and *T-Z*. The latter gene designations are used below because they take into account the previously published sequence of the dTDP-Rha biosynthesis genes (*cps23fL-O*) (Coffey *et al.*, 1998a). The predicted functions of the protein products encoded by the *cps23f* genes are listed in **Table 8.7**.

**Table 8.7. Predicted functions of *S. pneumoniae* cps23f<sup>a</sup> ORFs**

ORF	Predicted function	Similar proteins (% identity)
<i>cps23fA</i>	regulation?	<i>S. pneumoniae</i> Cps19fA (96%)
<i>cps23fB</i>	unknown	<i>S. pneumoniae</i> Cps19fB (88.5%)
<i>cps23fC</i>	chain length regulation/export	<i>S. pneumoniae</i> Cps19fC (69.6%)
<i>cps23fD</i>	chain length regulation/export	<i>S. pneumoniae</i> Cps19fD (79.3%)
<i>cps23fE</i>	glucosyl-1-phosphate transferase	<i>S. pneumoniae</i> Cps19fE (71.2%)
<i>cps23fT</i>	rhamnosyl transferase	<i>A. aeolicus</i> MtfB <sup>b</sup> (18.9%) <i>Synechocystis</i> sp. RfbW <sup>c</sup> (22.3%) <i>A. fulgidus</i> WbaZ-1 <sup>d</sup> (25.5%)
<i>cps23fI</i>	polysaccharide polymerase	-
<i>cps23fU</i>	galactosyl transferase	<i>S. pneumoniae</i> Cps14J <sup>e</sup> (33.5%) <i>S. thermophilus</i> EpsI <sup>f</sup> (36%) <i>Y. enterocolitica</i> TrsB <sup>g</sup> (29.7%)
<i>cps23fV</i>	rhamnosyl transferase	-
<i>cps23fJ</i>	repeat unit transporter	<i>S. thermophilus</i> EpsM <sup>f</sup> (32.1%) <i>S. pneumoniae</i> Cps14L <sup>e</sup> (16.7%) <i>Y. enterocolitica</i> TrsA <sup>g</sup> (17.9%)
<i>cps23fW</i>	glycerol-2-phosphate transferase	<i>B. subtilis</i> TagF <sup>h</sup> (27.8%) <i>S. pneumoniae</i> type 14 OrfX <sup>c</sup> (37%)
<i>cps23fX</i>	glyceraldehyde-2-phosphate dehydrogenase (CDP-2-glycerol biosynthesis)	<i>A. fulgidus</i> GldA <sup>i</sup> (32.3%) <i>M. jannaschii</i> GldA <sup>j</sup> (34.4%) <i>B. subtilis</i> AraM <sup>k</sup> (28.6%)
<i>cps23fY</i>	glycerol-2-phosphate cytidyltransferase (CDP-2-glycerol biosynthesis)	<i>A. fulgidus</i> RfbF <sup>l</sup> (27.5%)
<i>cps23fZ</i>	glyceraldehyde-2-phosphotransferase (CDP-2-glycerol biosynthesis)	<i>E. coli</i> NagD <sup>m</sup> (32%) <i>B. subtilis</i> AraL <sup>k</sup> (26.3%)
<i>cps23fL</i>	glucose-1-phosphate thymidyl transferase	<i>S. pneumoniae</i> Cps19fL (95.5%)
<i>cps23fM</i>	dTDP-4-keto-6-deoxyglucose-3,5-epimerase	<i>S. pneumoniae</i> Cps19fM (99%)
<i>cps23fN</i>	dTDP-glucose-4,6-dehydratase	<i>S. pneumoniae</i> Cps19fN (99.7%)
<i>cps23fO</i>	dTDP-L-rhamnose synthase	<i>S. pneumoniae</i> Cps19fO (99.3%)

The sequences can be found under GenBank accession nos: a, AF030373 and AF057294; b, .AE000693; c, D64002; d, AE001104; e, X85785; f, U40830; g, Z47767; h, Z99122; i, AE000988; j, V67518; k, Z99118; l, AE001025; m, AF052007.

### 8.3.7 The *S. pneumoniae* type 33F cps locus

The *S. pneumoniae* type 33f cps locus has been sequenced by Llull *et al.* (1998) and consists of 15 ORFs, designated *cap33fA-O*. The *cap33fA-O* genes and the predicted functions of their protein products are listed in **Table 8.8**. The authors did not assign functions to Cap33fF, G, I, K, M and O, however on closer examination of the similarities

to other proteins, putative functions could be assigned to all of these proteins, except Cap33fI.

**Table 8.8. Predicted functions of *S. pneumoniae* cps33f<sup>a</sup> ORFs**

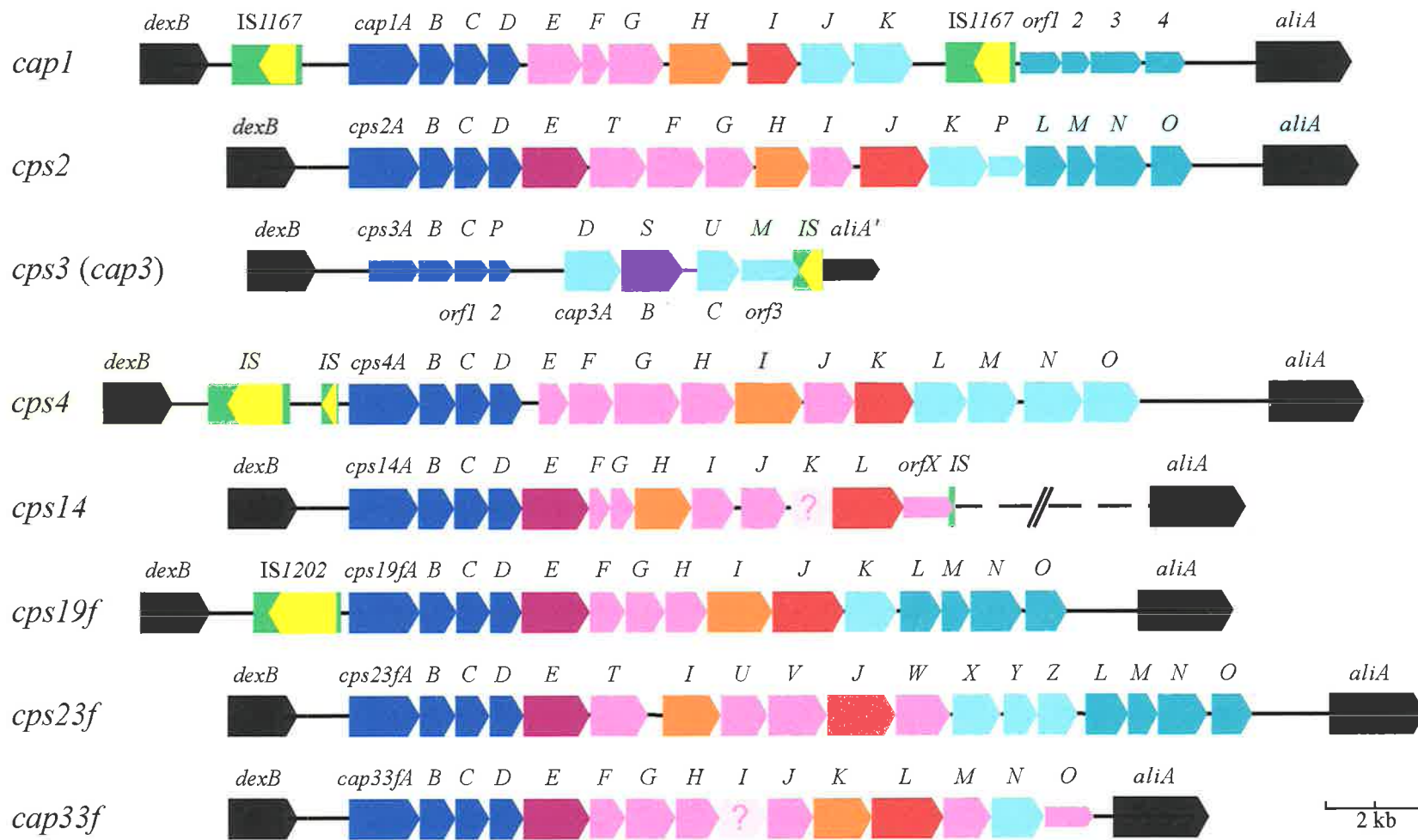
ORF	Predicted function	Similar proteins (% identity)
<i>cap33fA</i>	regulation?	<i>S. pneumoniae</i> Cps19fA (94.2%)
<i>cap33fB</i>	unknown	<i>S. pneumoniae</i> Cps19fB (85.6%)
<i>cap33fC</i>	chain length regulation/export	<i>S. pneumoniae</i> Cps19fC (71.3%)
<i>cap33fD</i>	chain length regulation/export	<i>S. pneumoniae</i> Cps19fD (78.4%)
<i>cap33fE</i>	glucosyl-1-phosphate transferase	<i>S. pneumoniae</i> Cps14E <sup>b</sup> (70.5%) <i>S. pneumoniae</i> Cps19fE (71%)
<i>cap33fF</i>	glycosyl transferase	<i>K. pneumoniae</i> RfbC <sup>c</sup> (34%)
<i>cap33fG</i>	glycosyl transferase	<i>K. pneumoniae</i> Orf7 <sup>d</sup> (28.7%)
<i>cap33fH</i>	galactosyl transferase	<i>S. pneumoniae</i> Cps14I <sup>b</sup> (26.5%) <i>S. pneumoniae</i> Cps14J <sup>b</sup> (24.9%)
<i>cap33fI</i>	unknown	<i>S. pneumoniae</i> Cps14K <sup>b</sup> (22.4%)
<i>cap33fJ</i>	galactosyl transferase	<i>S. pneumoniae</i> Cps14J <sup>b</sup> (31.2%)
<i>cap33fK</i>	polysaccharide polymerase	-
<i>cap33fL</i>	repeat unit transporter	<i>S. pneumoniae</i> Cps14L <sup>b</sup> (34.5%)
<i>cap33fM</i>	acetyl transferase	<i>E. coli</i> YinH <sup>e</sup> (23.7%)
<i>cap33fN</i>	UDP-galactopyranose mutase	<i>E. coli</i> Glf <sup>f</sup> (61.6%)
( <i>cap33fO</i> )	acetyl transferase	<i>B. subtilis</i> YkrP <sup>g</sup> (28.4%)

The sequences can be found under GenBank accession nos: a, AJ006986; b, X85787; c, L41518; d, D21242; e, P37669; f, U09876, g, Z99111.

## 8.4 Analysis of the *cps* loci from *S. pneumoniae*

### 8.4.1 Organisation of pneumococcal *cps* loci

All of the *S. pneumoniae* *cps* loci which have been sequenced to date are located on the *S. pneumoniae* chromosome between *dexB* and *aliA*, and with the exception of type 3, show a remarkably similar arrangement of the *cps* genes (Fig. 8.4). For most of these, functions of the gene products have been proposed on the basis of amino acid sequence



**Fig 8.4. The *cps* loci from *S. pneumoniae* serotypes 1, 2, 3, 4, 14, 19F, 23F and 33F.** The large boxed arrows represent ORFs involved in CPS biosynthesis and the narrow boxed arrows represent cryptic genes not involved with CPS biosynthesis in that serotype. The ORFs common to all serotypes are shown in dark blue. The specific transferases are shown in pink and the Glc-1-phosphate transferase which adds Glc to the lipid carrier is shown in dark pink. The type 3 synthase is shown in purple, the repeat unit transporters are shown in red and the polysaccharide polymerases are shown in orange. The ORFs involved in nucleotide-sugar biosynthesis are shown in pale blue and those involved in the biosynthesis of dTDP-Rha are shown a shade darker. The green boxes represent IS elements and the yellow boxed arrows represent their ORFs. The dashed line represents the unsequenced region between *dexB* and *aliA* (black boxed arrows). The two ORFs with question marks are of unknown function.

similarities with known proteins. However, Kolkman *et al.* (1996; 1997) have biochemically characterised the four glycosyl transferases encoded by genes in the type 14 *cps* locus. Also, a gene in the type 1 locus (*cap1K*) has been proven to encode a UDP-Glc dehydrogenase by complementation analysis (Muñoz *et al.*, 1997).

The first four genes (*cpsA-D*) are conserved at the 5' end of all loci except type 3, although type 3 retains extensive DNA similarity to these genes. However, there are no intact type 3 ORFs corresponding to *cpsA* and *B* due to several frame-shift mutations and deletions in this region. Two ORFs corresponding to *cpsC* and *D* are present, but *cps3P/orf2* (the *cpsD* homologue) is truncated compared to *cpsD* in other serotypes, and no promoter sequence for these ORFs has been identified in type 3 (García and López, 1997). This is consistent with the distinct mechanism of biosynthesis of type 3 CPS which is described below.

The central portion contains genes which encode the specific glycosyl transferases, the polysaccharide polymerase and the repeat unit transporter required for synthesis of the individual CPSs, and the 3' end contains the genes which encode the biosynthetic enzymes required to synthesise the individual nucleotide sugars. However, this arrangement is not absolute, as the type 14 *cps* locus does not appear to contain any genes which encode biosynthetic enzymes (the sequence downstream of *orfX* has not been determined). However, this is not surprising as type 14 CPS contains only Glc, GlcNAc and Gal, the nucleotide precursors of which are ubiquitous within the cell. Exceptions have also been identified within the group 19 *cps* loci. Types 19B and 19C contain a gene which encodes an enzyme putatively involved in the biosynthesis of an activated Rib precursor (*cps19R*) in their central region, and 19C also contains an additional putative glycosyl transferase gene (*cps19cS*) which has inserted between two genes which encode biosynthetic enzymes (**Fig. 8.1**).

Interestingly, the position of the genes which encode enzymes involved in the synthesis of dTDP-Rha is conserved, at the 3' end of the locus, in all the serotypes examined which contain Rha in their CPS. Type 1 CPS does not contain Rha, but cryptic copies of these genes (designated *orf1-4*) are present in the same position relative to *aliA* in the *cap1* locus of all type 1 isolates examined to date (Muñoz et al., 1997). These cryptic genes are separated from the *cap1* locus by a copy of *IS1167* and having no obvious promoter, are presumably not transcribed. Also *orf4* has a frame-shift mutation which would truncate the protein product. This suggests that the type 1 serotype may have arisen from a recombination event resulting in the insertion of type 1 specific DNA into the *cps* locus of a Rha-containing serotype.

#### 8.4.2 Comparison of the *cps* genes and their protein products

As already described above, the type 3 locus appears to be atypical. The type 3 CPS consists of a simple disaccharide repeat unit (**Table. 1.2**) and the *cps3/cap3* locus contains only three intact genes which are transcribed as a single unit (Arrecubieta *et al.*, 1995; Dillard *et al.*, 1995). The first gene (*cps3D* or *cap3A*) encodes a UDP-Glc dehydrogenase required for the synthesis of UDP-GlcA (Arrecubieta *et al.*, 1994). The second gene (*cps3S* or *cap3B*) encodes the type 3 synthase, a processive  $\beta$ -glycosyl transferase which links the alternating Glc and GlcA moieties via distinct glycosidic bonds (Arrecubieta *et al.*, 1995; Dillard *et al.*, 1995). There is a significant degree of amino acid sequence similarity between Cps3S/Cap3B and other bacterial polysaccharide synthases, including HasA which synthesises the hyaluronic acid capsule of group A streptococci. These synthases have a common predicted architecture and are capable of forming the two different glycosidic linkages, and extruding the growing polysaccharide chain, as it is synthesised, through a pore or channel formed in the membrane (section 1.8.4; Keenleyside and

Whitfield, 1996). Interestingly, transformation with plasmids carrying *cps3S/cap3B* alone is sufficient to direct synthesis of type 3 CPS in *E. coli* or in *S. pneumoniae* serotypes 1, 2, 5 or 8, all of which contain UDP-GlcA. In smooth heterologous *S. pneumoniae* hosts, expression of *cps3S/cap3B* resulted in binary encapsulated strains producing type 3 as well as the original CPS type (Arrecubieta *et al.* 1996). In the *E. coli* transformants, a significant proportion of the type 3 CPS appeared in the periplasm, indicating that additional type-specific genes are not required for transport of CPS across the cell membrane (Arrecubieta *et al.*, 1996). The final complete gene in the *cps3/cap3* locus (*cps3U* or *cap3C*) encodes a Glc-1-phosphate uridylyltransferase (or UDP-Glc pyrophosphorylase), which is capable of complementing an *E. coli galU* mutant, thereby confirming its function (Arrecubieta *et al.*, 1995). The product of this enzyme (UDP-Glc) is present in all pneumococci, regardless of serotype, and the pneumococcal *galU* gene (a functional homologue of *cps3U/cap3C*) has recently been identified and shown to be located elsewhere in the chromosome (Mollerach *et al.*, 1998). This is consistent with the finding that insertion-duplication mutagenesis of *cps3U/cap3C* does not abrogate type 3 CPS biosynthesis (Arrecubieta *et al.*, 1995; Dillard *et al.*, 1995). On the other hand, insertion-duplication mutagenesis of the *galU* gene, which is 77% identical to *cap3C*, resulted in loss of CPS expression in *S. pneumoniae* type 1 and type 3 (Mollerach *et al.*, 1998). Thus it appears that *cps3U/cap3C* is not capable of replacing the *galU* gene in *S. pneumoniae* either due to poor expression of the type 3 enzyme or poor affinity for the substrate (Mollerach *et al.*, 1998). The *galU* gene is present in all pneumococci tested, regardless of serotype and predicted to be absolutely essential for CPS production as most if not all serotypes contain Glc, and/or other sugars, such as Gal, GlcA, and GalA, which require UDP-Glc for biosynthesis of their activated precursors (Mills and Smith, 1965). No other *S. pneumoniae cps* locus examined to date contains a *galU* homologue. An

additional ORF, designated *cps3M*, is located 3' to *cps3U/cap3C* and has similarity to genes encoding phosphoglucomutases (Dillard *et al.*, 1995). However, it is truncated at its C-terminus and may not encode a functional enzyme; insertion-duplication mutagenesis of *cps3M* also has no impact on encapsulation (Caimano *et al.*, 1998; Dillard *et al.*, 1995).

The other *S. pneumoniae* serotypes have considerably more complex CPSs, consisting of 3-7 sugars, and as a consequence their *cps* loci are also more complex, containing 11 (*cap1*) to 19 (*cps19c*) genes. These loci all contain specific glycosyl transferases, a putative repeat unit transporter and a putative polysaccharide polymerase suggesting a common mode of CPS biosynthesis involving the sequential assembly of the CPS repeat units immobilised on a lipid carrier, followed by transport across the membrane and polymerisation prior to attachment to the cell surface. Thus, CPS biosynthesis in most *S. pneumoniae* serotypes proceeds as described in section 1.8.2 and resembles that of O-antigen biosynthesis in Gram-negative bacteria (Whitfield, 1995).

The serotype specificity of pneumococcal CPSs is determined by both the individual sugars present in the CPS and the nature of the glycosidic bonds joining them. This, in turn, is determined by both the substrate specificity and the nature of the glycosidic linkage formed by the glycosyl transferases encoded by the individual *cps* genes. Thus, each pneumococcal *cps* locus would be expected to include genes which encode mostly unique, potentially serotype determining, glycosyl transferases. For example, serotypes 4, 14, 23F and 33F all contain functionally different galactosyl transferases, which form different glycosyl linkages to different acceptor sugars, in accordance with their distinct CPS structures. Although the amino acid sequences share some similarities, there is no significant DNA similarity between these galactosyl transferases, suggesting that, as expected, they were probably acquired from diverse origins and are unlikely to have arisen from gene duplication.

However, the *cps* loci from all the serotypes examined, except types 1, 3 and 4, contain a homologue of *cps19fE* (70-99% identity). The closely related *cps14E* gene has been shown experimentally to encode a glucosyl transferase which adds Glc-1-phosphate to the lipid carrier in the cell membrane (Kolkman *et al.*, 1996), a common first step in CPS biosynthesis. The absence of a *cps19fE* homologue in types 1 and 4 is consistent with the absence of Glc in their CPS. The only other transferase genes which have significant DNA sequence identity (71.6%) are *cps2T* and *cps23fT* which encode putative  $\beta(1\rightarrow4)$  rhamnosyl transferases. Both these genes are located in the same position within their respective *cps* loci, just downstream of *cpsE*, raising the possibility that type 2 and type 23F strains may have had a common ancestor. Southern hybridisation data also indicates that a high stringency homologue of *cps19fF*, which encodes a putative ManNAc transferase, is present in type 9N and 9V strains. This transferase is predicted to link ManNAc to Glc via a  $\beta(1\rightarrow4)$  linkage in group 9 CPS production. Group 9 CPS also contains these two sugars joined by the same linkage, thus the presence of a *cps19fF* homologue in group 9 strains is not surprising.

The DNA similarity between the genes involved in nucleotide sugar biosynthesis in different serotypes is much greater than that observed for transferases. The *cap1J*, *cap3A* and *cps2K* genes all encode UDP-Glc dehydrogenase, which catalyses synthesis of UDP-GlcA; *cap1J* and *cps2K* share 75% identity, and the type 3 gene *cap3A* is slightly more divergent with only 60% identity to both *cap1J* and *cps2K*. UDP-GlcNAc-2-epimerase, which catalyses synthesis of UDP-ManNAc, is encoded by both *cps4L* and *cps19fK* which share 70% identity. Both *cps2P* and *cap33fN* (83.6% identity) are predicted to encode proteins similar to UDP-galactopyranose mutase which converts UDP-galactopyranose to UDP-galactofuranose. However, the *cps2P* gene is truncated at the 5' end compared to *cap33fN* and as UDP-galactofuranose is not required for type 2 CPS biosynthesis, *cps2P*

may be a cryptic gene within the *cps2* locus. The four dTDP-Rha biosynthesis genes are highly conserved in all serotypes containing Rha in their CPS, and have 90-99% nucleotide sequence identity. The high degree of sequence identity and the conserved position of these genes at the 3' end of the respective loci suggests the possibility of a common ancestor for these serotypes, and for type 1 which contains cryptic copies of these genes.

Both type 14 and type 33F *cps* loci contain copies of genes which encode redundant transferases at the 3' end of the loci. The redundant acetyl transferase encoded by *cap33fO* is truncated due to a frame-shift mutation in type 33F but is intact in the type 33A *cps* locus (Lull *et al.*, 1998). The *orfX* gene at the 3' end of the *cps14* locus is also truncated (Kolkman *et al.*, 1997b), but a glycerol-phosphate transferase, the predicted product of this gene, would be required for the biosynthesis of the immuno-cross reactive type 15A CPS (van Dam *et al.*, 1990). These redundant genes may be the remnants of a recombinational event within an ancestral *cps* locus, which resulted in the generation of a new serotype for which these distal genes were no longer required.

### 8.4.3 Capsular transformation *in vivo*

There is growing evidence that the phenomenon of capsular transformation first observed by Griffith (section 1.1.6) is a common phenomenon *in vivo*. Application of modern molecular typing techniques has resulted in the detection of otherwise genetically indistinguishable pneumococci expressing different capsular types. This has been particularly evident in clonal groups which are resistant to multiple antibiotics. Indeed, derivatives of a highly successful multiply-resistant type 23F clone (which originated in Spain) expressing types 3, 9N, 14, 19A and 19F capsules have been isolated (Barnes *et al.*, 1995; Coffey *et al.*, 1991; Coffey *et al.*, 1998a; Nesin *et al.*, 1998).

Completion of the *cps19f* sequence enabled further characterisation of a parental multiresistant type 23F strain and eight otherwise genetically indistinguishable type 19F clinical isolates. This involved PCR amplification, and subsequent sequencing, of the regions from *dexB-cpsB* and from *cpsL-aliA* for all of these strains (Coffey *et al.*, 1998b). Examination of polymorphisms in the conserved regions of the two *cps* loci indicated that in each case, the 5' recombination occurred upstream of *dexB*. In six of the eight type 19F strains, the 3' cross-over point was downstream of *aliA*. However, in the other two, a recombination cross-over point between the introduced type 19F sequences and the type 23F chromosome was identified; this was in *cpsM* in one strain and *cpsN* in the other (Coffey *et al.*, 1998b). Thus, capsule switching involves exchange of very large DNA fragments, ranging from at least 15 kb to over 22.5 kb. The existence of multiple cross-over points as well as additional minor polymorphisms within the type 19F-derived *cps* genes also indicated that the eight multiply resistant type 19F strains that were studied arose as a consequence of a minimum of four independent transformation events involving different type 19F donors. It therefore appears that these capsule switching events may be relatively common among pneumococci in nature (Coffey *et al.*, 1998b).

Multiple serotypes of *S. pneumoniae* are frequently carried concurrently in the human nasopharynx (Austrian, 1981b) providing ample opportunity for exchange of DNA between types. In addition to enhancing the spread of drug resistance amongst diverse capsular types, these exchanges may also provide a mechanism for evasion of serotype-specific host immune defences, such as those resulting from immunisation with pneumococcal CPS-protein conjugate vaccines which provide cover against a limited range of serotypes.

#### 8.4.4 The presence of IS elements in the 5' and 3' intergenic regions

Complete or partial IS elements have been located adjacent to 9 of the 11 *cps* loci sequenced to date, and in type 1 the locus is flanked at both the 5' and 3' ends by IS1167. The first IS element identified just upstream of the *S. pneumoniae cps* loci was IS1202 which is associated with all four members of group 19 (Morona *et al.*, 1994a). Since then several more IS elements have been identified in both the 5' and 3' intergenic regions of *S. pneumoniae cps* loci. In fact, 5' and 3' intergenic regions themselves contain sequences with some similarity to IS elements, and are repeated many times in the chromosome, although the complete IS elements are no longer present.

The *cps*-flanking regions appear to be common targets for IS elements, and this has led to the suggestion that they may play a role in horizontal transfer of *cps* genes (Kolkman *et al.*, 1998; Muñoz *et al.*, 1997). There are precedents for this in other bacteria; the *H. influenzae* type b capsule genes, for example, are located on a 17-kb compound transposon (Kroll *et al.*, 1991). However, the identification of cross-over points within the *cps* loci of the type 19F derivatives of the multiresistant type 23F *S. pneumoniae* clone referred to above confirm that at least in these cases, capsular exchange occurred as a consequence of homologous recombination rather than by transposition. Nevertheless, Muñoz *et al.* (1997) have demonstrated that IS1167 sequences flanking part of the *cap1* locus cloned in a plasmid could direct ectopic integration of these genes into copies of IS1167 located elsewhere in the pneumococcal chromosome, resulting in genetically binary strains.

#### 8.4.5 Transcription of the pneumococcal *cps* locus

As mentioned previously, the only promoter sequence identified in the *S. pneumoniae cps19f* locus is immediately upstream of *cps19fA* (Guidolin *et al.*, 1994), and the only stemmed-loop structures likely to function as transcription terminators are found

downstream of the *cpsI9f* locus in the 3' intergenic region (section 3.2.2). The consensus  $\sigma^{70}$  promoter sequence (TAGACA-17 bp-TATAAT) is also present in the same position in all of the other *cps* loci which have been sequenced, except type 3, which has a deletion of 280 nucleotides at the 5' end of the *cpsA*-related region resulting in the loss of the promoter sequence. The three *cps3* genes are transcribed as a single operon from a promoter immediately upstream of *cps3D* (*cap3A*) (Arrecubieta *et al.*, 1995). Llull *et al.* (1998) have shown that the promoter sequences identified for the *cps* loci in types 1, 3, 14 and 33F are all functional in *S. pneumoniae*. The transcriptional start site has also been determined by primer extension analysis of the type 1 locus (Muñoz *et al.*, 1997).

It is possible, however, that the level of transcription of *cps* loci varies from strain to strain. Although the -10 and -35 sequences themselves are highly conserved, Llull *et al.* (1998) have noted minor variations in flanking sequences. In type 37, for example, the -10 and -35 sequences are separated by 16 nucleotides rather than 17, and this correlated with a markedly lower promoter strength. A 4-nucleotide deletion immediately 5' to the -35 sequence in type 33F, relative to type 1, also correlated with a slight reduction in promoter strength (Llull *et al.*, 1998). The extent to which such differences impact on the level of encapsulation is uncertain, but it is curious in the light of these findings that type 37, which had the weakest *cps* promoter, produces one of the thickest capsules of all *S. pneumoniae* serotypes (Austrian, 1981b).

Many polysaccharide loci from Gram-negative bacteria (including: *E. coli kps*, *E. coli* K12 *rfa*, *S. flexneri rfb*, *S. enterica rfb*, *Vibrio cholerae rfb*, and *Y. enterocolitica rfb*) have a large leader sequence which is involved in transcriptional regulation of the locus (Marolda and Valvano, 1998). This region includes the conserved JUMPstart sequence (Hobbs and Reeves, 1994) which is part of two *cis*-acting sequences known as *ops* (for operon polarity suppressor) and are located upstream of the coding regions of all RfaH-

regulated operons (Marolda and Valvano, 1998). These leader sequences form a series of stem-loop structures; one of these contains the JUMPstart and *ops* sequences and interacts with RfaH, a transcription anti-terminator, and the RNA polymerase. In the absence of RfaH or with mutations affecting this stemmed-loop, premature termination of the transcript occurs when the RNA polymerase reaches the other stemmed-loop structures in this region. In the presence of RfaH, the JUMPstart stemmed-loop would serve to bring RfaH and possibly other factors together with the RNA polymerase complex, preventing the formation of the other stemmed-loops and hence, premature termination (Marolda and Valvano, 1998). However, in *S. pneumoniae cps* loci this leader sequence is entirely absent, with only 20 bp separating the transcriptional start site and the ATG start codon of *cpsA*. This is too short for formation of any stemmed-loop structures. Thus, transcription of the *cps* genes and hence CPS production in *S. pneumoniae* is not regulated in the same manner as polysaccharide production in Gram-negative bacteria.

#### **8.4.6 Regulation of CPS production in *S. pneumoniae***

Colonisation of the nasopharyngeal mucosa is an essential first step in the pathogenesis of pneumococcal disease, and is presumed to involve interaction between pneumococcal adhesins and specific receptors on host epithelial cells. In a proportion of cases, asymptomatic carriage progresses to invasive disease, and although the events involved are poorly understood, it is clearly a watershed in the bacteria-host relationship. In recent years evidence has emerged that this transition involves a major switch in expression of important virulence determinants, as the pneumococcus adapts to the altered microenvironment (Tuomanen and Masure, 1997). Maximal expression of capsule is clearly essential for systemic virulence, but the degree of exposure of other important pneumococcal surface structures, such as the adhesins, may also be influenced by capsular

thickness. Non-encapsulated pneumococci exhibit higher adherence to human respiratory epithelial (A549) cells *in vitro* than otherwise isogenic derivatives expressing either type 3 or type 19F capsules (Talbot *et al.*, 1996). Thus, the very feature (encapsulation) which is absolutely essential for systemic virulence of *S. pneumoniae* could be disadvantageous during the colonisation phase. Pneumococci have recently been shown to undergo a bidirectional phase variation between two distinct colonial morphologies, described as “opaque” and “transparent”. The transparent phenotype exhibits increased *in vitro* adherence to buccal epithelial cells and cytokine-activated A549 cells relative to opaque variants of the same strain, as well as an enhanced capacity to colonise the nasopharynx of infant rats (Cundell *et al.*, 1995; Weiser *et al.*, 1994). On the other hand, the opaque form is associated with massively increased virulence in animal models of systemic disease, and this correlates with increased production of CPS relative to cell wall teichoic acid compared with the transparent phenotype (Kim and Weiser, 1998). Phase variation also correlated with alteration in levels of several surface proteins, but the molecular mechanism involved is yet to be elucidated.

Clearly, the capacity to regulate CPS production is important for the survival of the pneumococcus in different host environments. The presence of a 115-bp repeated element has been noted upstream of the *cps* promoter in *S. pneumoniae* types 1, 3, 14, and 19F. This element appears to be specific to pneumococci, and copies have been found in the vicinity of other genes believed to be associated with virulence (Kolkman *et al.*, 1997b; Muñoz *et al.*, 1997). These sequences have no obvious function, although Kolkman *et al.* (1997) have suggested that they might have a regulatory function for coordinately controlled expression of virulence-related genes. However, examination of sequence data for *S. pneumoniae* type 4 released by TIGR indicates that there are at least 40 copies of sequences with > 80% identity to the element in the genome. In type 19F, the copy near

the *cps* locus is actually separated from the *cps* promoter by *IS1202*, and so it is difficult to imagine its involvement in a global regulatory mechanism. Moreover, insertion of the mutagenesis vector pVA891 into *IS1202*, which places the 115-bp element more than 8 kb upstream of the promoter, is known not to affect type 19F CPS production (Morona *et al.*, 1994a). As previously mentioned, in type 19F, *IS1202* is inserted immediately 5' to the -35 sequence of the *cps* promoter and so the above result, in all probability, precludes the involvement of any upstream region in transcriptional regulation of *cps19f*.

At present the mechanism of regulation of CPS production in *S. pneumoniae* is not understood. The only product of the *cps* loci likely to be involved is CpsA, which resembles a *B. subtilis* transcriptional attenuator. Recently, a derivative of *S. pneumoniae* type 19F with an in-frame deletion mutation in *cps19fA* was constructed. This mutant still produces type 19F CPS, as judged by quellung reaction. Also, colonies growing on BA appeared smooth, but they were much smaller than the parental type 19F strain, in spite of the lack of any obvious difference in growth rate (J. Morona, unpublished data). Additional phenotypic characterisation of this mutant, including transcriptional studies of *cps19f* and quantitation of total CPS production are in progress.

It is also possible that genes located elsewhere on the chromosome may influence CPS production in *S. pneumoniae*. Watson *et al.* (1995) suggested that a region 3' to the pneumococcal *lytA* (autolysin) gene might be involved in regulation of capsular expression, as this appeared to be the site of Tn916 insertion in a type 3 mutant which had lost the capacity to produce CPS. However, García *et al.* (1996) subsequently demonstrated that deletion of this region of the chromosome had no impact on encapsulation. An additional possibility is that regulation of CPS production is indirect, and mediated by availability of precursors or co-factors. For example, the mechanism of synthesis of the pneumococcal cell wall teichoic acid (see Fig. 1.2) is likely to be similar to CPS. Accordingly, if teichoic

acid synthesis is subject to direct regulation, this may influence the availability of common precursors (e.g. UDP-Glc or the lipid carrier) for CPS biosynthesis. Such a phenomenon would be consistent with the observation that phase variation in *S. pneumoniae* has opposite effects on the total amounts of CPS and teichoic acid (Kim and Weiser, 1998). Saturation of a common lipid carrier pool with the disaccharide precursor for peptidoglycan synthesis has previously been proposed to explain blockade of EPS production in *Streptococcus thermophilus* with a mutation in *pbp2b*, one of its penicillin-binding protein genes (Stingele and Mollet, 1996). Similarly, competition for the supply of GDP-mannose between the alginate and LPS biosynthetic pathways is thought to explain why mucoid *P. aeruginosa* strains lack the LPS O-antigen (May and Chakrabarty, 1994).

#### 8.4.7 The function of CpsC and CpsD

The *cpsC* and *D* gene products are predicted to function in CPS chain length regulation and export in an analogous fashion to ExoP from *R. meliloti* (Paulsen *et al.*, 1997). CpsC has two putative membrane-spanning hydrophobic domains which are found in all chain length regulators (Rol- and Cld-like proteins) and its topology is predicted to be similar these proteins. Thus, the N- and C-termini are predicted to be located in the cytoplasm, while the central portion is predicted to be exposed on the external side of the cell membrane. CpsD is presumed to be a cytoplasmic protein, and like the C-terminal portion of ExoP, contains a putative ATP binding domain (Guidolin *et al.*, 1994). The C-terminal "CpsD" domain of ExoP is required for synthesis of high molecular weight EPS (Becker *et al.*, 1995). This observation suggests that CpsD may interact with CpsC in the pneumococcus. Recently, derivatives of *S. pneumoniae* type 19F with in-frame deletion mutations in *cps19fC* and *cps19fD* were constructed. These mutants no longer produced

type 19F CPS, as judged by quellung reaction (J. Morona, unpublished data), indicating that both CpsC and CpsD are important for CPS production in *S. pneumoniae*.

Two distinct classes of *cpsC* and *cpsD* genes (class I and class II) have been identified in pneumococcal *cps* loci. However, the exact mechanism by which CpsC and CpsD (and their homologues) regulate polysaccharide chain length is still unknown. The effect of the two different classes on CPS chain length, and hence, capsule thickness await analysis of defined strains in which the class I and class II genes have been interchanged.

Genes which encode functional homologues of both CpsC and D are found in many polysaccharide biosynthesis loci from a large variety of Gram-negative and Gram-positive bacteria. The CpsC homologues (Rol- and Cld-like proteins) are present in all polysaccharide loci where biosynthesis occurs via lipid-linked repeat unit intermediates. However, CpsD homologues are not always present; they are not found in loci involved in O-antigen biosynthesis, but they are usually present in loci involved in CPS and EPS biosynthesis. This may have implications for the function of CpsC homologues in O-antigen and CPS biosynthesis. The average number of repeat units is much larger for CPS than for O-antigen, and CpsD may be required to enable CpsC to function such that a long polysaccharide chain can be synthesised (R. Morona, unpublished results). Gram-negative CPS and EPS biosynthesis loci generally contain a gene encoding one large protein (such as ExoP in *R. meliloti*, AmsA in *E. amylovora* and Orf6 in *K. pneumoniae* O1:K2) with separate domains homologous to either CpsC or CpsD (Whitfield *et al.*, 1997). However, in Gram-positive CPS and EPS biosynthesis loci, two separate genes encode CpsC and CpsD homologues, as discussed below.

#### **8.4.8 Relationship of the *S. pneumoniae cps* locus to other Gram-positive *cps* loci**

To date, CpsA and CpsB homologues have only been found in ORFs cloned from

other Gram-positive bacteria, indicating that their function is probably Gram-positive specific. The only protein of known function similar to either of these proteins is LytR (a transcriptional attenuator of the *B. subtilis* autolytic enzymes) which has 28% identity to Cps19fA. However, the function of CpsA in pneumococcal CPS production has yet to be determined, as does that for CpsB. Homologues to CpsC and CpsD are more wide-spread, found in both Gram-positive and Gram-negative polysaccharide biosynthesis loci, as described above.

The presence of homologues of the first four *S. pneumoniae* *cps* genes, *cpsA-D*, was examined in other Gram-positive polysaccharide loci. Interestingly, homologues to CpsA were only found in other *Streptococcus* species, with the first four genes arranged in the same order in all of the loci examined (Table 8.9). CPS and EPS biosynthesis loci from other Gram-positive bacteria contained homologues to CpsB, C and D, also in the 5' region of the loci, but the gene order differed to that in the *Streptococcus* species. In *S. aureus* types 1, 5 and 8 capsules, the *L. lactis* EPS and a polysaccharide from *B. subtilis*, the homologues are ordered *cpsC*, *D* and then *B*. However, the *cpsC* homologue is not necessarily the first gene in the locus; in both the *eps* locus from *L. lactis* and the *ywq* locus from *B. subtilis*, there are two genes preceding the *cpsC* homologue (unrelated to the pneumococcal *cpsA* or *B* genes) one of which encodes a protein (EpsR) with similarity to DNA-binding regulatory proteins (van Kranenburg *et al.*, 1997). Interestingly, an ORF designated *orfY*, located after the transcription terminator at the 3' end of the *eps* locus from *Lactococcus lactis*, and transcribed in the opposite direction to the *eps* genes, has significant similarity to LytR (23.6% identity), but not to CpsA. Any role this gene may have in EPS production in *L. lactis* has yet to be determined (van Kranenburg *et al.*, 1997). In group B streptococci, an additional gene located upstream of the *cpsA* homologue, but transcribed in the opposite direction, encodes a protein (CpsY) with similarity to the LysR

family of DNA-binding transcriptional regulators (Koskiniemi *et al.*, 1998). However, no such gene is present adjacent to any of the pneumococcal *cps* loci examined to date. Thus, it appears that distinct transcriptional regulatory mechanisms may exist for polysaccharide biosynthesis loci in different Gram-positive species.

**Table 8.9. CpsA-D homologues present in polysaccharide biosynthesis loci of Gram-positive bacteria**

Homologues to:	CpsA	CpsB	CpsC	CpsD
<i>Streptococcus pneumoniae</i> type 19F	Cps19fA	Cps19fB	Cps19fC	Cps19fD
<i>Streptococcus pneumoniae</i> type 19A	Cps19aA	Cps19aB	Cps19aC	Cps19aD
<i>Streptococcus pneumoniae</i> type 19B	Cps19bA	Cps19bB	Cps19bC	Cps19bD
<i>Streptococcus pneumoniae</i> type 19C	Cps19cA	Cps19cB	Cps19cC	Cps19cD
<i>Streptococcus pneumoniae</i> type 1 <sup>a</sup>	Cap1A	Cap1B	Cap1C	Cap1D
<i>Streptococcus pneumoniae</i> type 2 <sup>b</sup>	Cps2A	Cps2B	Cps2C	Cps2D
<i>Streptococcus pneumoniae</i> type 4 <sup>c</sup>	Cps4A	Cps4B	Cps4C	Cps4D
<i>Streptococcus pneumoniae</i> type 14 <sup>d</sup>	Cps14A	Cps14B	Cps14C	Cps14D
<i>Streptococcus pneumoniae</i> type 23F <sup>e</sup>	Cps23fA	Cps23fB	Cps23fC	Cps23fD
<i>Streptococcus pneumoniae</i> type 33F <sup>f</sup>	Cap33fA	Cap33fB	Cap33fC	Cap33fD
<i>Streptococcus thermophilus</i> <sup>g</sup>	EpsA	EpsB	EpsC	EpsD
<i>Streptococcus salivarius</i> <sup>h</sup>	CpsA	CpsB	CpsC	CpsD
<i>Streptococcus agalactiae</i> <sup>i</sup>	CpsX	CpsA	CpsB	CpsC
<i>Staphylococcus aureus</i> type 1 <sup>j</sup>		Cap1C	Cap1A	Cap1B
<i>Staphylococcus aureus</i> type 5 <sup>k</sup>		Cap5C	Cap5A	Cap5B
<i>Staphylococcus aureus</i> type 8 <sup>l</sup>		Cap8C	Cap8A	Cap8B
<i>Lactococcus lactis</i> <sup>m</sup>		EpsC	EpsA	EpsB
<i>Bacillus subtilis</i> <sup>n</sup>		YwqE	YwqC	YwqD

a, Muñoz *et al.* (1997); b, Iannelli *et al.* (1998); c, TIGR Microbial Database; d, Kolkman *et al.* (1997); e, Morona *et al.* (1998); f, Llull *et al.* (1998); g, Stinglele *et al.* (1996); h, Griffin *et al.* (1996); i, Koskiniemi *et al.* (1998) and Rubens *et al.* (1993); j, Wen *et al.* (1994); k, Sau *et al.* (1997); l, Sau and Lee (1996); m, van Kranenburg *et al.* (1997); n, Presecan *et al.* (1997).

The polysaccharide loci examined above are all presumed to share a similar mode of biosynthesis, involving the assembly of repeat units on lipid carriers which are exported prior to polymerisation. The *S. pneumoniae* type 3 and the Group A Streptococcus

hyaluronic acid CPS loci were not included in **Table 8.9** because the mode of CPS biosynthesis differs, involving a synthase which exports the growing polysaccharide chain, as described in sections 1.8.4 and 8.4.2.

The similarities between the CPS and EPS loci listed in **Table 8.9** extend over the entire locus. They all appear to be arranged as tightly linked, single transcriptional units with the common genes located at the 5' end. The serotype-specific genes encoding specific transferases, a repeat unit transporter, and a polysaccharide polymerase are located in the central portion. The genes encoding enzymes involved in biosynthesis of nucleotide sugars are generally located at the 3' end of the loci. This arrangement has some similarity to CPS loci in Gram-negative bacteria which are divided into three separate regions with the serotype-specific region in the centre (Roberts, 1996; section 1.8.6).

## 8.5 Future studies

Much remains to be learnt about the precise molecular events involved in the mechanisms of CPS biosynthesis, and about how CPS production in pneumococci is regulated. Detailed biochemical and mutational analyses are also required to confirm proposed functions of specific genes and in particular to elucidate the precise functions of the four genes at 5' end of the *cps* loci. These presumably encode important common steps in polysaccharide biosynthesis in pneumococci, as well as in certain other Gram-positive genera. As a first step, in-frame deletion mutants of *cps19fA*, *cps19fC* and *cps19fD* have already been constructed in *S. pneumoniae* strain Rx1-19F. Given the importance of capsules to the virulence of *S. pneumoniae* and several other Gram-positive pathogens, such conserved components of the CPS biosynthesis machinery may prove to be useful targets for novel antimicrobial strategies.

# REFERENCES

- Aarson, VA., Kristinsson, KG., Sigurdsson, JA., Stefansdottir, G., Molstad, S. and Gudmundsson, S. 1996. Do antimicrobials increase the carriage rate of penicillin resistant pneumococci in children? Cross sectional prevalence study. *Br. Med. J.* 313:387-391.
- Adams, WG., Deaver, KA., Cochi, SL., Plikaytis, BD., Zell, ER., Broome, CV. and Wenger, JD.; for the *Haemophilus influenzae* Study Group. 1993. Decline of childhood *Haemophilus influenzae* Type b (Hib) disease in the Hib vaccine era. *JAMA.* 269:221-226.
- Alloing, G., de Philip, P. and Claverys, J-P. 1994. Three highly homologous membrane-bound lipoproteins participate in oligopeptide transport by the Ami system of the Gram-positive *Streptococcus pneumoniae*. *J. Mol. Biol.* 241:44-58.
- Alloway, H. 1932. The transformation in vitro of R pneumococci into S forms of different specific types by use of filtered pneumococcus extracts. *J. Exp. Med.* 55:91-99.
- Alloway, H. 1933. Further observations on the use of pneumococcus extracts in effecting transformation of type in vitro. *J. Exp. Med.* 57:265-278.
- AlonsoDeVelasco, E., Verheul, AFM., Verhoef, J. and Snippe, H. 1995. *Streptococcus pneumoniae*: Virulence factors, pathogenesis, and vaccines. *Microbiol. Rev.* 59:591-603.
- Altschul, SF., Gish, W., Miller, W., Myers, EW. and Lipman, DJ. 1990. Basic local alignment search tool. *J. Mol. Biol.* 215:403-410.
- Altschul, SF., Madden, TL., Schaffer, AA., Zhang, J., Zhang, Z., Miller, W. and Lipman, DJ. 1997. Gapped BLAST and PSI-BLAST: a new generation of protein database search programs. *Nucleic Acids Res.* 25:3389-3402.
- Appelbaum, PC. 1992. Antimicrobial resistance in *Streptococcus pneumoniae*: An overview. *Clin. Infect. Dis.* 15:77-83.
- Arakawa, Y., Wacharotayankun, R., Nagatsuka, T., Ito, H., Kato, N. and Ohta, M. 1995. Genomic organisation of the *Klebsiella pneumoniae cps* region responsible for serotype K2 capsular polysaccharide synthesis in the virulent strain Chedid. *J. Bacteriol.* 177:1788-1796.
- Armstrong, RR. 1931. A swift and simple method for detecting pneumococcal "type". *Br. Med. J.* 1:214-215.
- Arrecubieta, C., García, E. and López, R. 1995. Sequence and transcriptional analysis of a DNA region involved in the production of capsular polysaccharide in *Streptococcus pneumoniae* type 3. *Gene* 167:1-7.
- Arrecubieta, C., López, R. and García, E. 1994. Molecular characterisation of cap3A, a gene from the operon required for the synthesis of the capsule of *Streptococcus pneumoniae* type 3: sequencing of mutations responsible for the unencapsulated phenotype and localisation of the capsular cluster on the pneumococcal chromosome. *J. Bacteriol.* 176:6375-6383.
- Arrecubieta, C., López, R. and García, E. 1996. Type 3-specific synthase of *Streptococcus pneumoniae* (Cap3B) directs type 3 polysaccharide biosynthesis in *Escherichia coli* and in pneumococcal strains of different serotypes. *J. Exp. Med.* 184:449-455.
- Austrian, R. 1981a. Pneumococcus: The first one hundred years. *Rev. Infect. Dis.* 3:183-189.

- Austrian, R. 1981b.** Some observations on the pneumococcus and on the current status of pneumococcal disease and its prevention. *Rev. Infect. Dis.* 3(Suppl.):S1-S17.
- Austrian, R., Bernheimer, HP., Smith, EEB. and Mills, GT. 1959.** Simultaneous production of two capsular polysaccharides by pneumococcus. II. The genetic and biochemical bases of binary capsulation. *J. Exp. Med.* 110:585-602.
- Avery, OT. and Dubos, R. 1931.** The protective action of a specific enzyme against type 3 pneumococcus infection in mice. *J. Exp. Med.* 54:73-89.
- Avery, OT., MacLeod, CM. and McCarty, M. 1944.** Studies on the chemical nature of the substance inducing transformation of pneumococcal types: Induction of transformation by a desoxyribonucleic acid fraction isolated from pneumococcus type III. *J. Exp. Med.* 79:137-157.
- Barnes, DM., Whittier, S., Gilligan, PH., Soares, S., Tomasz, A. and Henderson, FW. 1995.** Transmission of multidrug-resistant serotype 23F *Streptococcus pneumoniae* in group day care: evidence suggesting capsular transformation of the resistant strain *in vivo*. *J. Infect. Dis.* 171:890-896.
- Barry, MA., Craven, DE. and Finland, M. 1984.** Serotypes of *Streptococcus pneumoniae* isolated from blood cultures at Boston City Hospital between 1979 and 1982. *J. Infect. Dis.* 149:449-452.
- Bechthold, A., Sohng, JK., Smith, TM., Chu, X. and Floss, HG. 1995.** Identification of *Streptomyces violaceoruber* Tu22 genes involved in the biosynthesis of granaticin. *Mol. Gen. Genet.* 248:610-620.
- Becker, A., Niehaus, K. and Pühler, A. 1995.** Low-molecular-weight succinoglycan is predominantly produced by *Rhizobium meliloti* strains carrying a mutated ExoP protein characterised by a periplasmic N-terminal domain and a missing C-terminal domain. *Mol. Microbiol.* 16:191-203.
- Bernheimer, HP., Wermundsen, IE. and Austrian, R. 1967.** Qualitative differences in the behaviour of pneumococcal deoxyribonucleic acids transforming to the same capsular type. *J. Bacteriol.* 93:320-333.
- Berry, AM., Lock, RA., Thomas, SM., Rajan, DP., Hansman, D. and Paton, JC. 1994.** Cloning and nucleotide sequence of the *Streptococcus pneumoniae* hyaluronidase gene, and purification of the enzyme from recombinant *Escherichia coli*. *Infect. Immun.* 62:1101-1108.
- Berry, AM., Lock, RA. and Paton, JC. 1996.** Cloning and characterisation of *nanB*, a second *Streptococcus pneumoniae* neuraminidase gene, and purification of the NanB enzyme from recombinant *Escherichia coli*. *J. Bacteriol.* 16:4854-4860.
- Berry, AM. and Paton, JC. 1996.** Sequence heterogeneity of PsaA, a 37-kilodalton putative adhesin essential for virulence of *Streptococcus pneumoniae*. *Infect. Immun.* 64:5255-5262.
- Berry, AM., Paton, JC., Glare, EM. Hansman, D. and Catchside, DEA. 1988.** Cloning and expression of the pneumococcal neuraminidase gene in *Escherichia coli*. *Gene.* 71:299-305.
- Berry, AM., Yother, J., Briles, DE., Hansman, D. and Paton, JC. 1989.** Reduced virulence of a defined pneumolysin-negative mutant of *Streptococcus pneumoniae*. *Infect. Immun.* 57:2324-2330.
- Beynon, LM., Richards, JC., Perry, MB. and Kniskern. 1991.** Antigenic and structural relationships within group 19 *Streptococcus pneumoniae*: chemical characterisation of the specific capsular polysaccharides of type 19B and 19C. *Can. J. Chem.* 70:131-137.
- Bluestone, CD., Stephenson, JS. and Martin, LM. 1992.** Ten-year review of otitis media pathogens. *Pediatr. Infect. Dis. J.* 11:S7-S11.
- Bogaerts, J., Lepage, P., Taelman, H., Rouvroy, D., Batungwanayo, J., Kestelyn, P., Hitimana, DG., Van de Perre, P., Vandepitte, J., Verbist, L. and Verhaegen, J. 1993.** Antimicrobial susceptibility and serotype distribution of *Streptococcus pneumoniae* from Rwanda, 1984-1990. *J. Infect.* 27:157-168.

- Briles, DE., Yother, J. and McDaniel, LS. 1988.** Role of pneumococcal surface protein A in the virulence of *Streptococcus pneumoniae*. *Rev. Infect. Dis.* 3:S82-S88.
- Bridge, PD. and Sneath, PHA. 1983.** Numerical taxonomy of *Streptococcus*. *J. Gen. Microbiol.* 129:565-597.
- Bronner, D., Clarke, BR. and Whitfield, C. 1994.** Identification of an ATP-binding cassette transport system required for translocation of lipopolysaccharide O-antigen side-chains across the cytoplasmic membrane of *Klebsiella pneumoniae* serotype O1. *Mol. Microbiol.* 14:505-519.
- Brown, MCM., Weston, A., Saunders, JR. and Humphreys, GO. 1979.** Transformation of *E. coli* C600 by plasmid DNA at different phases of growth. *FEMS Microbiol. Lett.* 5:219-222.
- Brown, PK., Romana, LK. and Reeves, PR. 1992.** Molecular analysis of the *rfb* gene cluster of *Salmonella* serovar muenchen (strain M67): the genetic basis of the polymorphism between groups C2 and B. *Mol. Microbiol.* 6:1385-1394.
- Brynstad, S., Synstad, B. and Granum, PE. 1997.** The *Clostridium perfringens* enterotoxin gene is on a transposable element in type A human food poisoning strains. *Microbiol.* 143:2109-2115.
- Bugert, P. and Geider, K. 1995.** Molecular analysis of the *ams* operon required for exopolysaccharide synthesis of *Erwinia amylovora*. *Mol. Microbiol.* 15:917-933.
- Burke, JP., Klein, JO., Gezon, HM. and Finland, M. 1971.** Pneumococcal bacteraemia. *Amer. J. Dis. Child.* 121:353-359.
- Butler, JC., Breiman, RF., Lipman, HB., Hofmann, J. and Facklam, RR. 1995.** Serotype distribution of *Streptococcus pneumoniae* infections among preschool children in the United States, 1978-1994: Implications for development of a conjugate vaccine. *J. Infect. Dis.* 171:885-889.
- Cabib, E., Bowers, B. and Roberts, RL. 1983.** Vectorial synthesis of a polysaccharide by isolated plasma membranes. *Proc. Natl. Acad. Sci. USA.* 80:3318-3321.
- Caimano, MJ., Hardy, GG. and Yother, J. 1998.** Capsule genetics in *Streptococcus pneumoniae* and a possible role for transposition in the generation of the type 3 locus. *Microb. Drug. Resist.* 4:11-23.
- Camara, M., Mitchell, TJ, Andrew, PW. and Boulnois, GJ. 1991.** *Streptococcus pneumoniae* produces at least two distinct enzymes with neuraminidase activity: cloning and expression of a second neuraminidase gene in *Escherichia coli*. *Infect. Immun.* 59:2856-2858.
- Centers for Disease Control and Prevention. 1997.** Prevention of pneumococcal disease: Recommendations of the Advisory Committee on Immunisation Practices (ACIP). *MMWR* 46(No. RR-8):1-24.
- Cherian, T., Steinhoff, MC., Harrison, LH., Rohn, D., McDougal, L. and Dick, J. 1994.** A cluster of invasive pneumococcal diseases in young children in child care. *JAMA.* 271:695-698.
- Coffey, TJ., Dowson, CG., Daniels, M., Zhou, J., Martin, C., Spratt, BG. and Musser, JM. 1991.** Horizontal gene transfer of multiple penicillin-binding protein genes and capsular biosynthetic genes in natural populations of *Streptococcus pneumoniae*. *Mol. Microbiol.* 5:2255-2260.
- Coffey, TJ., Enright, MC., Daniels, M., Morona, JK., Morona, R., Hryniewicz, W., Paton, JC. and Spratt, BG. 1998a.** Recombinational exchanges at the capsular polysaccharide biosynthetic locus lead to frequent serotype changes among natural isolates of *Streptococcus pneumoniae*. *Mol. Microbiol.* 27:73-83.
- Coffey, TJ., Enright, MC., Daniels, M., Wilkinson, P., Berron, S., Fenoll, A. and Spratt, BG. 1998b.** Serotype 19A variants of the Spanish serotype 23F multiresistant clone of *Streptococcus pneumoniae*. *Microb. Drug Resist.* 4:51-55.
- Collignon, PJ. and Bell, JM. 1996.** Drug resistant *Streptococcus pneumoniae*: the beginning of the end for many antibiotics? *Med. J. Aust.* 164:64-67.

- Collins, LV. and Hackett, J. 1991.** Molecular cloning, characterisation, and nucleotide sequence of the *rfc* gene, which encodes an O-antigen polymerase of *Salmonella typhimurium*. *J. Bacteriol.* 173:2521-2529.
- Crain, MJ., Waltman, WD. II, Turner, JS., Yother, J., Talkington, DF., McDaniel, LS., Gray, BM. and Briles, DE. 1990.** Pneumococcal surface protein A (PspA) is serologically highly variable and is expressed by all clinically important capsular serotypes of *Streptococcus pneumoniae*. *Infect. Immun.* 58:3293-3299.
- Cundell, DR., Weiser, JN., Shen, J., Young, A. and Tuomanen, EI. 1995.** Relationship between colonial morphology and adherence of *Streptococcus pneumoniae*. *Infect. Immun.* 63:757-761.
- Dagan, R., Engelhard, D., Piccard, E., and the Israeli Paediatric Bacteraemia and Meningitis Group. 1992.** Epidemiology of invasive childhood pneumococcal infections in Israel. *JAMA.* 268:3328-3332.
- Dagan, R., Melamed, R., Muallem, M., Piglansky, L., Greenberg, D., Abramson, O., Mendelman, PM., Bohidar, N. and Yagupsky, P. 1996.** Reduction of nasopharyngeal carriage of pneumococci during the second year of life by a heptavalent conjugate pneumococcal vaccine. *J. Infect. Dis.* 174:1271-1278.
- Dalrymple, B., Caspers, P. and Arber, W. 1984.** Nucleotide sequence of the prokaryotic mobile genetic element IS30. *EMBO J.* 3:2145-2149.
- Daniels, C., Vindurampulle, C. and Morona, R. 1998.** Overexpression and topology of the *Shigella flexneri* O-antigen polymerase (Rfc/Wzy). *Mol. Microbiol.* 28:1211-1222.
- Daniels, DL., Plunkett III, G. Burland, VD. and Blattner, FR. 1992.** Analysis of the *Escherichia coli* genome: DNA sequence of the region from 84.5 to 86.5 minutes. *Science* 257:771-778.
- Dawson, MH. and Sia, RHP. 1931.** In vitro transformation of pneumococcal types: 1. A technique for inducing transformation of pneumococcal types in vitro. *J. Exp. Med.* 54:681-699.
- DeAngelis, PL., Papaconstantinou, J. and Weigel, PH. 1993.** Molecular cloning, identification, and sequence of the hyaluronan synthase gene from Group A *Streptococcus pyogenes*. *J. Biol. Chem.* 268:19181-19184.
- DeAngelis, PL. and Weigel, PH. 1994.** Immunochemical confirmation of the primary structure of streptococcal hyaluronan synthase and synthesis of high molecular weight product by the recombinant enzyme. *Biochem.* 33:9033-9039.
- Dillard, JP., Vandersea, MW. and Yother, J. 1995.** Characterisation of the cassette containing genes for type 3 capsular polysaccharide biosynthesis in *Streptococcus pneumoniae*. *J. Exp. Med.* 181:973-983.
- Dintilhac, A., Alloing, G., Granadel, C. and Claverys, JP. 1997.** Competence and virulence of *Streptococcus pneumoniae*: Adc and PsaA mutants exhibit a requirement for Zn and Mn resulting from inactivation of putative ABC metal permeases. *Mol. Microbiol.* 25:727-739.
- Distler, J., Mansouri, K., Mayer, G., Stockman, M. and Piepersberg, W. 1992.** Streptomycin biosynthesis and its regulation in *Streptomyces*. *Gene.* 115:105-111.
- Dochez, AR. and Avery, OT. 1917.** Soluble substance of pneumococcus origin in the blood and urine during lobar pneumonia. *Proc. Soc. Exp. Biol. Med.* 14:126-127.
- Dolph, PJ., Majerczak, DR. and Coplin, DL. 1988.** Characterisation of a gene cluster for exopolysaccharide biosynthesis and virulence in *Erwinia stewartii*. *J. Bacteriol.* 170:865-871.
- Dong, Q., Sadouk, A., van der Lelie, D., Taghavi, S., Ferhat, A., Nuyten, JM., Borremans, B., Mergeay, M. and Toussaint, A. 1992.** Cloning and sequencing of IS1086, an *Alcaligenes eutrophus* insertion element related to IS30 and IS4351. *J. Bacteriol.* 174:8133-8138.

- Douglas, RM., Paton, JC., Duncan, SJ. and Hansman, DJ. 1983. Antibody response to pneumococcal vaccination in children younger than five years of age. *J. Infect. Dis.* 148:131-137.
- Drake, CR., Roberts, IS., Jann, B., Jann, K. and Boulnois, G. 1990. Molecular cloning and expression of the genes encoding the *Escherichia coli* K4 capsular polysaccharide, a fructose-substituted chondroitin. *FEMS Microbiol. Lett.* 66:227-30.
- Eriksen, KR. 1945. Studies on induced resistance to penicillin in a pneumococcus type 1. *Acta Pathol. Microbiol. Scand.* 22:398-405.
- Fallarino, A., Mavrangelos, C., Strocher, UH. and Manning, PA. 1997. Identification of additional genes required for O-antigen biosynthesis in *Vibrio cholerae* O1. *J. Bacteriol.* 179:2147-2153.
- Fang, GD., Fine, M., Orloff, J., Arisumi, D., Yu, VL., Kapoor, W., Grayston, JT., Wang, SP., Kohler, R., Muder, RR., Yee, YC., Rihs, JD. and Vickers RM. 1990. New and emerging aetiologies for community-acquired pneumonia with implications for therapy: A prospective multicenter study of 359 cases. *Medicine.* 69:307-316.
- Ferrante, A., Rowan-Kelly, B. and Paton, JC. 1984. Inhibition of *in vitro* human lymphocyte response by the pneumococcal toxin pneumolysin. *Infect. Immun.* 46:585-589.
- Finland, M. and Barnes, MW. 1977. Changes in occurrence of capsular serotypes of *Streptococcus pneumoniae* at Boston City Hospital during selected years between 1935 and 1974. *J. Clin. Microbiol.* 5:154-166.
- Fischer, H. and Tomasz, A. 1985. Peptidoglycan cross-linking and teichoic acid attachment in *Streptococcus pneumoniae*. *J. Bacteriol.* 163:46-54.
- Frisch, AW., Price, AE. and Meyers, GB. 1943. Type VIII pneumococcus: Development of sulphadiazine resistance transmission by cross infection and persistence in carriers. *Ann. Intern. Med.* 18:271-278.
- Frosch, M. and Müller, A. 1993. Phospholipid substitution of capsular polysaccharides and mechanisms of capsule formation in *Neisseria meningitidis*. *Mol. Microbiol.* 8:483-493.
- Frosch, M., Weisgerber, C. and Meyer, T. 1989. Molecular characterisation and expression in *Escherichia coli* of the gene complex encoding the polysaccharide capsule of *Neisseria meningitidis* group B. *Proc. Natl. Acad. Sci. USA.* 86:1669-1673.
- Fuchs-Cleveland, E. and Gilvarg, CC. 1976. Oligomeric intermediates in peptidoglycan biosynthesis in *Bacillus megaterium*. *Proc. Natl. Acad. Sci. USA.* 73:4200-4204.
- Gallo, MA., Ward, J. and Hutchinson, CR. 1996. The *dnrM* gene in *Streptomyces peucetius* contains a naturally occurring frameshift mutation that is suppressed by another locus outside of the daunorubicin-production gene cluster. *Microbiol.* 142:269-275.
- García, E., Arrecubieta, C. and López, R. 1996. The *lytA* gene and the region located downstream of this gene are not involved in the formation of the type 3 capsular polysaccharide of *Streptococcus pneumoniae*. *Curr. Microbiol.* 33:133-135.
- García, P., García, JL., García, E. and López, R. 1986. Nucleotide sequence and expression of the pneumococcal autolysin gene from its own promoter in *Escherichia coli*. *Gene.* 43:265-272.
- García, E. and López, R. 1997. Molecular biology of the capsular genes of *Streptococcus pneumoniae*. *FEMS Microbiol. Lett.* 149:1-10.
- Garcia-Bustos, JF., Chait, BT. and Tomasz, A. 1987. Structure of the peptide network of pneumococcal peptidoglycan. *J. Biol. Chem.* 262:15400-15405.
- Garcia-Bustos, JF. Tomasz, A. 1990. A biological price of antibiotic resistance: Major changes in the peptidoglycan structure of penicillin-resistant pneumococci. *Proc. Natl. Acad. Sci. USA.* 87:5414-5419.

- Giffard, PM., Rathsam, C., Kwan, E., Kwan, DW., Bunny, KL., Koo, SP. and Jaques, NA. 1993.** The *fff* gene encoding the cell-bound fructosyl transferase of *Streptococcus salivarius* ATCC 25975 is preceded by an insertion sequence and followed by *FUR1* and *clpP* homologues. *J. Gen. Microbiol.* 139:913-920.
- Giudicelli, S. and Tomasz, A. 1984.** Attachment of pneumococcal autolysin to wall teichoic acids, an essential step in enzymatic wall degradation. *J. Bacteriol.* 158:1188-1190.
- Glucksmann, MA., Reuber, TL. and Walker, GC. 1993.** Genes needed for the modification, polymerisation, export, and processing of succinoglycan by *Rhizobium meliloti*: a model for succinoglycan biosynthesis. *J. Bacteriol.* 175:7045-7055.
- Gram, C. 1884.** Über die isolierte färbung der schizomyceten in schnitt- und trockenpräparaten. *Fortschr. Med.* 2:185-189.
- Gratten, M., Barker, J., Shann, F., Gerega, G., Montgomery, J., Kajoi, M. and Lupiwa, T. 1985.** The aetiology of purulent meningitis in highland children: a bacteriological study. *PNG. Med. J.* 28:233-240.
- Gratten, M. and Montgomery J. 1991.** The bacteriology of acute pneumonia and meningitis in children in Papua New Guinea: assumptions, facts and technical strategies. *PNG. Med. J.* 34:185-198.
- Gratten, M., Naraqi, S. and Hansman, D. 1980.** High prevalence of penicillin-insensitive pneumococci in Port Moresby, Papua New Guinea. *Lancet.* 2:192-195.
- Gratten, M., Torzillo, P., Morey, F., Dixon, J., Erlich, J., Harrer, J. and Henrichsen, J. 1996.** Distribution of capsular types and antibiotic susceptibility of invasive *Streptococcus pneumoniae* isolated from Aborigines in Central Australia. *J. Clin. Microbiol.* 34:338-341.
- Gray, BM., Converse III, GM. and Dillon Jr., HC. 1979.** Serotypes of *Streptococcus pneumoniae* causing disease. *J. Infect. Dis.* 140:979-983.
- Gray, BM., Converse III, GM. and Dillon Jr., HC. 1980.** Epidemiologic studies of *Streptococcus pneumoniae* in infants: Acquisition, carriage, and infection during the first 24 months of life. *J. Infect. Dis.* 142:923-933.
- Grebe, T., Paik, J. and Hakenbeck, R. 1997.** A novel resistance mechanism against beta-lactams in *Streptococcus pneumoniae* involves CpoA, a putative glycosyl transferase. *J. Bacteriol.* 179:3342-3349.
- Griffin, AM., Morris, VJ. and Gasson, MJ. 1996.** The *cpsABCDE* genes involved in polysaccharide production in *Streptococcus salivarius* ssp. thermophilus strain NCBF 2393. *Gene.* 183:23-27.
- Griffith, F. 1928.** The significance of pneumococcal types. *J. Hyg.* 27:113-159.
- Gross, M., Geier, G., Rudolph, K. and Geider, K. 1992.** Levan and levansucrase synthesised by the fireblight pathogen *Erwinia amylovora*. *Physiol. Mol. Plant Pathol.* 40:371-381.
- Guidolin, A., Morona, JK., Morona, R., Hansman, D. and Paton, JC. 1994.** Nucleotide sequence analysis of genes essential for capsular polysaccharide biosynthesis in *Streptococcus pneumoniae* type 19F. *Infect. Immun.* 62:5384-5396.
- Guirguis, NI., Helmy, MF., Mohamed, MR. and Ali, RH. 1990.** A suggested vaccine formulation for the control of pneumococcal meningitis in Egypt. *J. Egypt Public Health Assoc.* 65:291-303.
- Hammerschmidt, S., Birkholz, C., Zahringer, U., Robertson, BD., van Putten, J., Ebeling, O. and Frosch, M. 1994.** Contribution of genes from the capsule gene complex (*cps*) to lipooligosaccharide biosynthesis and serum resistance in *Neisseria meningitidis*. *Mol. Microbiol.* 11:885-896.
- Hammerschmidt, S., Taley, SR., Brandtzaeg, P. and Chatwal, S. 1997.** SpsA, a novel pneumococcal surface protein with specific binding to secretory Immunoglobulin A and secretory component. *Mol. Microbiol.* 25:1113-1124.

- Hammond, SM., Lambert, PA and Rycroft, AN. 1984.** The bacterial cell surface. Croon Helm, London.
- Hanahan, D. 1983.** Studies on transformation of *Escherichia coli* with plasmids. *J. Mol. Biol.* 166:557-580.
- Hansman, D. 1983.** Serotypes in pneumococcal disease: a ten year study in Australia 1970 through 1979. *Aust. NZ. J. Med.* 13:359-64.
- Hansman, D. and Bullen, MM. 1967.** A resistant pneumococcus [letter]. *Lancet.* 2:264-265.
- Hansman, D., Devitt, L., Miles, H. and Riley, I. 1974.** Pneumococci relatively insensitive to penicillin in Australia and New Guinea. *Med. J. Aust.* 2:353-356.
- Hashimoto, Y., Li, N., Yokoyama, H. and Ezaki, T. 1993.** Complete nucleotide sequencing and characterisation of ViaB region encoding Vi antigen in *Salmonella typhi*. *J. Bacteriol.* 175:4456-4465.
- Heidelberger, M. 1927.** Immunologically specific polysaccharides. *Chem. Rev.* 3:403-423.
- Heidelberger, M., and Avery, OT. 1923.** The soluble specific substance of pneumococcus. *J. Exp. Med.* 38:73-79.
- Heinrichs, DE., Monteiro, MA., Perry, MB. and Whitfield, C. 1998.** The assembly system for the lipopolysaccharide R2 core-type of *Escherichia coli* is a hybrid of those found in *Escherichia coli* K-12 and *Salmonella enterica*. Structure and function of the R2 WaaK and WaaL homologues. *J. Biol. Chem.* 23:8849-8859.
- Henikoff, S. 1984.** Unidirectional digestion with exonuclease III creates targeted breakpoints for DNA sequencing. *Gene.* 28:351-359.
- Henrichsen, J. 1995.** Six newly recognised types of *Streptococcus pneumoniae*. *J. Clin. Microbiol.* 33:2759-2762.
- Higgins, DG. and Sharp, PM. 1988.** CLUSTAL: a package for performing multiple sequence alignments on a microcomputer. *Gene.* 73:237-244.
- Hilleman, MR., Carlson Jr., AM., McLean, AA., Vella, PP., Weibel RE. and Woodhour, AF. 1981.** *Streptococcus pneumoniae* polysaccharide vaccine: age, and dose responses, safety, persistence of antibody, revaccination, and simultaneous administration of pneumococcal and influenza vaccines. *Rev. Infect. Dis.* 3(Suppl):S31-S42.
- Hobbs, M. and Reeves, PR. 1994.** The JUMPstart sequence: a 39 bp element common to several polysaccharide gene clusters. *Mol. Microbiol.* 12:855-856.
- Hofmann, K. and Stöffel, W. 1989.** PROFILEGRAPH: an interactive graphical tool for protein sequence analysis. *Comput. Appl. Biosci.* 5:151-153.
- Hook, EW., Horton, CA. and Scharberg, DR. 1983.** Failure of intensive care unit support to influence mortality from pneumococcal bacteraemia. *JAMA.* 249:1055-1057.
- Hotchkiss, RD. 1951.** Transfer of penicillin resistance in pneumococci by the desoxynucleate derived from resistant cultures. *Cold Spring Harb. Symp. Quant. Biol.* 16:457-461.
- Huang, J. and Schell, M. 1995.** Molecular characterisation of the *eps* gene cluster of *Pseudomonas solanacearum* and its transcriptional regulation at a single promoter. *Mol. Microbiol.* 16:977-989.
- Iannelli, F., Pearce, BJ. and Pozzi, G. 1998.** The type 2 capsule locus of *Streptococcus pneumoniae*. *J. Bacteriol.* (in press).
- Ielpi, L., Couso, RO. and Dankert, MA. 1993.** Sequential assembly and polymerisation of the polypreo-linked pentasaccharide repeating unit of the xanthan polysaccharide in *Xanthomonas campestris*. *J. Bacteriol.* 175:2490-2500.

- Issaëff, B. 1893.** Contribution à l'étude de l'immunité acquise contre le pneumocoque. Ann. de l'Institut Pasteur. 7:260-279.
- Jacobs, MR., Koornhof, HJ., Robins-Browne, RM. Stevenson, CM., Vermaak, ZA., Freiman, I., Miller, GB., Witcomb, MA., Isaacson, M., Ward, JI. and Austrian, R. 1978.** Emergence of multiply resistant pneumococci. New Engl. J. Med. 299:734-740.
- Jacobs, NM., Lerdkachornsuk, S. and Metzger, W. 1979.** 1. Pneumococcal bacteraemia in infants and children: a ten-year experience at the Cook County Hospital with special reference to the pneumococcal serotypes isolated. Pediatrics. 64:296-300.
- Jansson, PE., Lindberg, B., Andersson, M., Lindquist, U. and Henrichsen, J. 1988.** Structural studies of the capsular polysaccharide from *Streptococcus pneumoniae* type 2, a reinvestigation. Carbohyd. Res. 182:111-117.
- Jansson, PE., Lindberg, B. and Lindquist, U. 1985.** Structural studies of the capsular polysaccharide from *Streptococcus pneumoniae* type 5. Carbohyd. Res. 140:101-110.
- Jennings, HJ., Lugowski, C. and Young, NM. 1980.** Structure of the complex polysaccharide C-substance from *Streptococcus pneumoniae* type 1. Biochemistry. 19:4712-4719.
- Jiang, X-M., Neal, B., Santiago, F., Lee, SJ., Romana, LK. and Reeves, PR. 1991.** Structure and sequence of the *rfb* (O antigen) gene cluster of *Salmonella* serovar typhimurium (strain LT2). Mol. Microbiol. 5:695-713.
- Jobling, MG. and Holmes, RK. 1990.** Construction of vectors with the p15a replicon, kanamycin resistance, inducible *lacZ* $\alpha$  and pUC18 or pUC19 multiple cloning sites. Nucleic Acids Res. 18:5315-5316.
- Jones, C. and Currie, F. 1988.** The pneumococcal polysaccharide S4: A structural re-assessment. Carbohyd. Res. 184:279-284.
- Jones, C., Mulloy, B., Wilson, A., Dell, A. and Oates, JE. 1985.** Structure of the capsular polysaccharide from *Streptococcus pneumoniae* type 9. J. Chem. Soc. Perkin I:1665-1673.
- Jones, JNK. and Perry, MB. 1957.** The structure of the type VIII pneumococcus specific polysaccharide. J. Amer. Chem. Soc. 79:2787-2793.
- Kaneko, T., Tanaka, A., Sato, S., Kotani, H., Sazuka, T., Miyajima, N., Sugiura, M. and Tabata, S. 1995.** Sequence analysis of the genome of the unicellular cyanobacterium *Synechocystis* sp. strain PCC6803. I. Sequence features in the 1 Mb region from map positions 64% to 92% of the genome. DNA Res. 2:153-166.
- Kapur, V., Reda, KB., Li, LL., Rich, RR. and Musser JM. 1994.** Characterisation and distribution of insertion sequence IS1239 in *Streptococcus pyogenes*. Gene 150:135-140.
- Katzenellenbogen, E. and Jennings. HJ. 1983.** Structural determination of the capsular polysaccharide of *Streptococcus pneumoniae* type 19A (57). Carbohyd. Res. 124:235-245.
- Keenleyside, WJ. and Whitfield, C. 1995.** Lateral transfer of *rfb* genes: a mobilizable ColE1-type plasmid carries the *rfb*O:54 (O:54 antigen biosynthesis) gene cluster from *Salmonella enterica* serovar Borreze. J. Bacteriol. 177:5247-5253.
- Keenleyside, WJ. and Whitfield, C. 1996.** A novel pathway for O-polysaccharide biosynthesis in *Salmonella enterica* serovar Borreze. J. Biol. Chem. 271:28581-28592.
- Kelly, T., Dillard, JP. and Yother, J. 1994.** Effect of genetic switching of capsular type on virulence of *Streptococcus pneumoniae*. Infect. Immun. 62:1813-1819.
- Kenne, L., Lindberg, B. and Madden, JK. 1979.** Structural studies of the capsular antigen from *Streptococcus pneumoniae* type 26. Carbohyd. Res. 73:175-182.

- Kessler, AC., Haase, A. and Reeves, PR. 1993.** Molecular analysis of the 3,6-dideoxyhexose pathway genes of *Yersinia pseudotuberculosis* serogroup IIA. *J. Bacteriol.* 175:1412-1422.
- Kiino, DR., Licudine, R., Wilt, K., Yang, DHC. and Rothman-Denes, LB. 1993.** A cytoplasmic protein, NfrC, is required for bacteriophage N4 adsorption. *J. Bacteriol.* 175:7074-7080.
- Kim JO. and Weiser, JN. 1998.** Association of intrastrain phase variation in quantity of capsular polysaccharide and teichoic acid with the virulence of *Streptococcus pneumoniae*. *J. Infect. Dis.* 177:368-377.
- Klee, SR., Tzschaschel, BD., Timmis, KN. and Guzman, CA. 1997.** Influence of different *rol* gene products on the chain length of *Shigella dysenteriae* type 1 lipopolysaccharide O antigen expressed by *Shigella flexneri* carrier strains. *J. Bacteriol.* 179:2421-2425.
- Klemperer, G and Klemperer, F. 1891.** Versuche über immunisierung und heilung bei der pneumokokkeninfektion. *Berliner Klinische Wochenschrift.* 28:833-835, 869-875.
- Klena, JD. and Schnaitman, CA. 1993.** Function of the *rfb* gene cluster and the *rfe* gene in the synthesis of O-antigen by *Shigella dysenteriae* 1. *Mol. Microbiol.* 9:393-402.
- Klugman, KP. and Koornhof, HJ. 1988.** Drug resistance patterns and serogroups of serotypes of pneumococcal isolates from cerebrospinal fluid of blood, 1979-1986. *J. Infect. Dis.* 158:956-964.
- Klugman, KP. 1996.** Serotypes and clones of antibiotic-resistant pneumococci. Abstract from *Streptococcus pneumoniae: Molecular biology and mechanisms of disease.* Oeiras, Portugal, September, 1996.
- Kolkman, MAB., Morrison, DA., van der Zeijst, BAM. and Nuijten, PJM. 1996.** The capsule polysaccharide synthesis locus of *Streptococcus pneumoniae* serotype 14: identification of the glycosyl transferase gene *cps14E*. *J. Bacteriol.* 178:3736-3741.
- Kolkman, MAB., van der Zeijst, BAM. and Nuijten, PJM. 1997a.** Functional analysis of glycosyl transferases encoded by the capsular polysaccharide biosynthesis locus of *Streptococcus pneumoniae* serotype 14. *J. Biol. Chem.* 272:19502-19508.
- Kolkman, MAB., van der Zeijst, BAM. and Nuijten, PJM. 1998.** Diversity of capsular polysaccharide synthesis gene clusters in *Streptococcus pneumoniae*. *J. Biochem. (Tokyo)* 123:937-945.
- Kolkman, MAB., Wakarchuk, W., Nuijten, PJM. and van der Zeijst, BAM. 1997b.** Capsular polysaccharide synthesis in *Streptococcus pneumoniae* serotype 14: molecular analysis of the complete *cps* locus and identification of genes encoding glycosyl transferases required for the biosynthesis of the tetrasaccharide subunit. *Mol. Microbiol.* 26:197-208.
- Koplin, R., Wang, G., Hotte, B., Priefer, UB. and Puhler, A. 1993.** A 3.9-kb DNA region of *Xanthomonas campestris* pv. *campestris* that is necessary for lipopolysaccharide production encodes a set of enzymes involved in the synthesis of dTDP-rhamnose. *J. Bacteriol.* 175:7786-7792.
- Koskiniemi, S., Sellin, M. and Norgren, M. 1998.** Identification of two genes, *cpsX* and *cpsY*, with putative regulatory function on capsule expression in group B streptococci. *FEMS Immun. Med. Microbiol.* 21:159-168.
- Kroll, JS., Loynds, BM. and Moxon, ER. 1991.** The *Haemophilus influenzae* capsulation gene cluster: a compound transposon. *Mol. Microbiol.* 5:1549-1560.
- Kroll, S., Zamze, S., Loynds, B. and Moxon, E. 1989.** Common organisation of chromosomal loci for production of different capsular polysaccharides in *Haemophilus influenzae*. *J. Bacteriol.* 171:3343-3347.
- Kumar, S., Tamura, K. and Nei, M. 1994.** MEGA: Molecular Evolutionary Genetics Analysis software for microcomputers. *Comput. Appl. Biosci.* 10:189-191.

- Kyte, J. and Doolittle, RF. 1982.** A simple method for displaying the hydrophobic character of a protein. *J. Mol. Biol.* 157:105-132.
- Laemmli, UK. 1970.** Cleavage of structural proteins during the assembly of the head of bacteriophage T4. *Nature.* 227:680-685.
- Lazarevic, V. and Karamata, D. 1995.** The *tagGH* operon of *Bacillus subtilis* 168 encodes a two-component ABC transporter involved in the metabolism of two wall teichoic acids. *Mol. Microbiol.* 16:345-355.
- Lazarevic, V., Margot, P., Soldo, B. and Karamata, D. 1992.** Sequencing and analysis of the *Bacillus subtilis* *lytRABC* divergon: a regulatory unit encompassing the structural genes of the *N*-acetylmuramoyl-L-alanine amidase and its modifier. *J. Gen. Microbiol.* 138: 1949-1961.
- Leach, A., Ceesay, SJ., Banya, WAS. and Greenwood, BM. 1996.** Pilot trial of a pentavalent pneumococcal polysaccharide/protein conjugate vaccine in Gambian infants. *Pediatr. Infect. Dis. J.* 15:333-339.
- Lee, C-J. 1987.** Bacterial capsular polysaccharide biochemistry, immunity and vaccine. *Mol. Immun.* 24:1005-1019.
- Lee, C-J. and Fraser, BA. 1980.** The structures of the cross-reactive types 19 (19F) and 57 (19A) pneumococcal polysaccharides. *J. Biol. Chem.* 255:6847-6853.
- Lee, C-J., Fraser, BA., Boykins, RA. and Li JP. 1987.** Effect of culture conditions on the structure of *Streptococcus pneumoniae* type 19A (57) capsular polysaccharide. *Infect. Immun.* 55:1819-1823.
- Lee, SJ., Romana, LK. and Reeves, PR. 1992.** Cloning and structure of group C1 O antigen (*rfb* gene cluster) from *Salmonella enterica* serovar montevideo. *J. Gen. Microbiol.* 138:305-312.
- Lin, WS., Cunneen, T. and Lee, CY. 1994.** Sequence analysis and molecular characterisation of genes required for the biosynthesis of type 1 capsular polysaccharide in *Staphylococcus aureus*. *J. Bacteriol.* 176:7005-7016.
- Lindberg, B., Lönngrén, J. and Powell, DA. 1977.** Structural studies on the specific type 14 pneumococcal polysaccharide. *Carbohydr. Res.* 58:177-186.
- Lindberg, B., Lindqvist, B., Lönngrén, J. and Powell, DA. 1980.** Structural studies of the capsular polysaccharide from *Streptococcus pneumoniae* type 1. *Carbohydr. Res.* 78:111-117.
- Linton, KJ., Jarvis, BW. and Hutchinson, CR. 1995.** Cloning of the genes encoding thymidine diphosphoglucose 4,6-dehydratase and thymidine diphospho-4-keto-6-deoxyglucose 3,5-epimerase from the erythromycin-producing *Saccharopolyspora erythraea*. *Gene.* 153:33-40.
- Liu, D., Cole, RA. and Reeves, PR. 1996.** An O-antigen processing function for Wzx (RfbX): a promising candidate for O-unit flippase. *J. Bacteriol.* 178:2102-2107.
- Llull, D., López, R., García, E. and Muñoz, R. 1998.** Molecular structure of the gene cluster responsible for the synthesis of the polysaccharide capsule of *Streptococcus pneumoniae* type 33F. *Biochim. Biophys. Acta.* (in press).
- Lock, RA., Paton, JC. and Hansman, D. 1988a.** Purification and immunological characterisation of neuraminidase produced by *Streptococcus pneumoniae*. *Microb. Pathogen.* 4:33-43.
- Lock, RA., Paton, JC. and Hansman, D. 1988b.** Comparative efficacy of pneumococcal neuraminidase and pneumolysin as immunogens protective against *Streptococcus pneumoniae* infection. *Microb. Pathogen.* 5:461-467.
- Lu, C-D. and Abdelal, AT. 1993.** The *Salmonella typhimurium* uracil-sensitive mutation *use* is in *argU* and encodes a minor arginine tRNA. *J. Bacteriol.* 175:3897-3899.

- Lugowski, C. and Jennings, HJ. 1984.** Structural determination of the capsular polysaccharide of *Streptococcus pneumoniae* type 18C (56). *Carbohydr. Res.* 131:119-129.
- MacLeod, CM. and Krauss, MR. 1950.** Relation of virulence of pneumococcal strains to the quantity of capsular polysaccharide formed *in vitro*. *J. Exp. Med.* 92:1-9.
- McCarty, M. and Avery, OT. 1946.** Studies on the chemical nature of the substance inducing transformation of pneumococcal types: II. Effect of desoxyribonuclease on the biological activity of the transforming substance. *J., Exp. Med.* 83:89-96.
- McDaniel, LS., Scott, G., Kearney, JF. and Briles, DE. 1984.** Monoclonal antibodies against protease-sensitive pneumococcal antigens can protect mice from fatal infection with *Streptococcus pneumoniae*. *J. Exp. Med.* 160:386-397.
- McDaniel, LS., Yother, J., Vijayakamur, M., McGarry, L., Guild, WR. and Briles, DE. 1987.** Use of insertional inactivation to facilitate studies of biological properties of pneumococcal surface protein A (PspA). *J. Exp. Med.* 165:381-394.
- McIntyre, P. 1997.** Epidemiology and prevention of pneumococcal disease. *Comm. Dis. Intell.* 21:41-46.
- Macpherson, DF., Manning, PA. and Morona, R. 1994.** Characterisation of the dTDP-rhamnose biosynthetic genes encoded in the *rfb* locus of *Shigella flexneri*. *Mol. Microbiol.* 11:281-292.
- Macpherson, DF., Manning, PA. and Morona, R. 1995.** Genetic analysis of the *rfbX* gene of *Shigella flexneri*. *Gene.* 155:9-17.
- Macrina, FL., Evans, RP., Tobian, JA., Hartley, DL., Clewell, DB. and Jones, KR. 1983.** Novel shuttle plasmid vehicles for *Escherichia-Streptococcus* transgeneric cloning. *Gene.* 25:145-150.
- Maharaj, R., May, TB., Wang, SK. and Chakrabarty, AM. 1993.** Sequence of the *alg8* and *alg44* genes involved in the synthesis of alginate by *Pseudomonas aeruginosa*. *Gene.* 136:267-269.
- Maniatis, T., Fritsch, EF. and Sambrook, J. 1982.** *Molecular cloning: a laboratory manual.* Cold Spring Harbor Laboratory, Cold Spring Harbor, NY.
- Marolda, CL. and Valvano, MA. 1995.** Genetic analysis of the dTDP-rhamnose biosynthesis region of the *Escherichia coli* VW187 (O7:K1) *rfb* gene cluster: identification of functional homologues of *rfbB* and *rfbA* in the *rff* cluster and correct location of the *rffE* gene. *J. Bacteriol.* 177:5539-5546.
- Marolda, CL. and Valvano, MA. 1998.** Promoter region of the *Escherichia coli* O7-specific lipopolysaccharide gene cluster: structural and functional characterisation of an upstream untranslated mRNA sequence. *J. Bacteriol.* 180:3070-3079.
- Mastro, TD., Ghafoor, A., Nomani, NK., Ishaq, Z., Anwar, F., Granoff, DM., Spika, JS., Thornsberry, C. and Facklam RR. 1991.** Antimicrobial resistance of pneumococci in children with acute lower respiratory tract infection in Pakistan. *Lancet.* 1:156-159
- Mauel, C., Young, M. and Karamata, D. 1991.** Genes concerned with synthesis of poly(glycerol phosphate), the essential teichoic acid in *Bacillus subtilis* strain 168, are organised in two divergent transcriptional units. *J. Gen. Microbiol.* 137:929-941.
- May, TB. and Chakrabarty, AM. 1994.** *Pseudomonas aeruginosa*: genes and enzymes of alginate synthesis. *Trends Microbiol.* 2:151-157.
- Meier-Dieter, U., Starman, R., Barr, K., Mayer, H. and Rick, PD. 1990.** Biosynthesis of enterobacterial common antigen in *Escherichia coli*. Biochemical characterisation of Tn10 insertion mutants defective in enterobacterial common antigen synthesis. *J. Biol. Chem.* 265:13490-13497.
- Merson-Davies, LA. and Cundliffe, E. 1994.** Analysis of five tylosine biosynthetic genes from the *tylIBA* region of the *Streptomyces fradiae* genome. *Mol. Microbiol.* 13:349-355.

- Mills, GT. and Smith, EEB. 1965.** Biosynthesis of capsular polysaccharides in the pneumococcus. *Bull. Soc. Chim. Biol.* 47:1751-1765.
- Mitchell, TJ., Andrew, PW., Saunders, FK., Smith, AN. and Boulnois, GJ. 1991.** Complement activation and antibody binding by pneumolysin via a region of the toxin homologous to a human acute phase protein. *Mol. Microbiol.* 5:1883-1888.
- Mitchison, MD. Bulach, M., Vinh, T., Rajakumar, K., Faine, S. and Adler, B. 1997.** Identification and characterisation of the dTDP-rhamnose biosynthesis and transfer genes of the lipopolysaccharide-related *rfb* locus in *Leptospira interrogans* serovar Copenhageni. *J. Bacteriol.* 179:1262-1267.
- Mogdasy, MC., Camou, T., Fajardo, C. and Hortal, M. 1992.** Colonising and invasive strains of *Streptococcus pneumoniae* in Uruguayan children: Type distribution and patterns of antibiotic resistance. *Pediatr. Infect. Dis. J.* 11:648-652.
- Mollerach, M., López, R. and García, E. 1998.** Characterisation of the *galU* gene of *Streptococcus pneumoniae* encoding a UDP-glucose pyrophosphorylase: a gene essential for capsular polysaccharide biosynthesis. *J. Exp. Med.* (in press).
- Montgomery, JM., Lehman, D., Smith, T., Michael, A., Joseph, B., Lupiwa, T., Coakley, C., Spooner, V., Best, B., Riley, ID. and Alpers, MP. 1990.** Bacterial colonisation of the upper respiratory tract and its association with acute lower respiratory tract infections in Highland children of Papua New Guinea. *Rev. Infect. Dis.* 12:S1006-S1016.
- Moore, HF. and Chesney, AM. 1917.** A study of ethylhydrocuprein (optichin) in the treatment of acute lobar pneumonia. *Ach. Intern. Med.* 19:611-682.
- Mooser, G., Hefta, SA., Paxton, RJ., Shively, JE. and Lee, TD. 1991.** Isolation and sequence of an active-site peptide containing a catalytic aspartic acid from two *Streptococcus sorbinus* alpha-glucosyl transferases. *J. Biol. Chem.* 266:8916-8922.
- Mooser, G. and Wong, C. 1988.** Isolation of a glucan binding domain of glucosyl transferase (1,6- $\alpha$ -glucan synthase) from *Streptococcus sorbinus*. *Infect. Immun.* 56:880-884.
- Moreau, M., Richards, JC., Perry, MB. and Kniskern, PJ. 1988.** Application of high resolution n. m. r. spectroscopy to the elucidation of the structure of the specific capsular polysaccharide of *Streptococcus pneumoniae* type 7F. *Carbohydr. Res.* 182:79-99.
- Morelle, G. 1989.** A plasmid extraction procedure on a miniprep scale. *Focus.* 11.1:7-8.
- Morgenroth, J. and Kaufman, M. 1912.** Arzneifestigkeit bei bakterien (pneumokokken). *Zeitschrift für Immunitätsforschung und Experimentelle Therapie.* 15:610-624.
- Morona, JK., Guidolin, A., Morona, R., Hansman, D. and Paton, JC. 1994a.** Isolation, Characterisation, and Nucleotide Sequence of IS1202, an Insertion Sequence of *Streptococcus pneumoniae*. *J. Bacteriol.* 176:4437-4443.
- Morona, JK., Miller, DC., Coffey, TJ., Vindurampulle, CJ., Spratt, BG., Morona, R. and Paton, JC. 1998.** Molecular and genetic characterisation of the capsule biosynthesis locus of *Streptococcus pneumoniae* type 23F. *Microbiol.* (Manuscript submitted).
- Morona, R., Macpherson, DF., Van den Bosch, L., Carlin, NIA. and Manning, PA. 1995.** Lipopolysaccharide with an altered O-antigen produced in *Escherichia coli* K-12 harbouring mutated, cloned *Shigella flexneri rfb* genes. *Mol. Microbiol.* 18:209-223.
- Morona, R., Mavris, M., Fallarino, A. and Manning, PA. 1994b.** Characterisation of the *rfc* region of *Shigella flexneri*. *J. Bacteriol.* 176: 733-747.
- Mufson, MA., Krause, HE., Schiffman, G. and Hughey, DF. 1987.** Pneumococcal antibody levels one decade after immunisation of healthy adults. *Amer. J. Med. Sci.* 293:279-284.

- Mukasa, H. 1986.** Properties of *Streptococcus mutans* glucosyl transferases, p 121-132. In Molecular microbiology and immunobiology of *Streptococcus mutans*. Hanada, S., Michalek, SM., Kiyono, H., Menaker, L. and McGhee, JR. Elsevier Science Publishers, New York.
- Muñoz, R., Mollerach, M., López, R. and García, E. 1997.** Molecular organisation of the genes required for the synthesis of type 1 capsular polysaccharide of *Streptococcus pneumoniae*: formation of binary encapsulated pneumococci and identification of cryptic dTDP-rhamnose biosynthesis genes. Mol. Microbiol. 25:79-92.
- Musher, DM. 1992.** Infections caused by *Streptococcus pneumoniae*: Clinical spectrum, pathogenesis, Immunity, and treatment. Clin. Infect. Rev. 14:801-809.
- Musher, DM., Luchi, MJ., Watson, DA., Hamilton, R. and Baughn, RE. 1990.** Pneumococcal polysaccharide vaccine in young adults and older bronchitics: determination of IgG responses by ELISA and the effect of adsorption of serum with non-type-specific cell wall polysaccharide. J. Infect. Dis. 161:728-735.
- Nesin M., Ramirez, M. and Tomasz, A. 1998.** Capsular transformation of a multidrug-resistant *Streptococcus pneumoniae* in vivo. J. Infect. Dis. 177:707-713.
- Neufeld, F. 1902.** Über die agglutination der pneumokokken und über die theorien der agglutination. Zeitschrift für Hygiene und Infektionskrankheiten (Leipzig). 40:54-72.
- Neufeld, F. and Rimpau, W. 1904.** Über die antikorper des streptokokken- und pneumokokken-immunserums. Deutsche Med. Wochenschrift. 30:1458-1460.
- Neuhard, J. and Thomassen, E. 1976.** Altered deoxyribonucleotide pools in P2 educants of *Escherichia coli* K-12 due to deletion of the *dcd* gene. J. Bacteriol. 126:999-1001.
- Nielsen, SV. and Henrichsen, J. 1992.** Capsular types of *Streptococcus pneumoniae* isolated from blood and CSF during 1982-1987. Clin. Infect. Dis. 15:794-798.
- Obaro, SK., Adegbola, RA., Banya, WAS. and Greenwood, BM. 1996.** Carriage of pneumococci after pneumococcal vaccination. Lancet. 348:271-272.
- Ochman, H., Gerber, AS. and Hartl, DL. 1988.** Genetic applications of an inverse polymerase chain reaction. Genetics. 120:621-623.
- Odani, T., Shimma, Y., Tanaka, A. and Jigami, Y. 1996.** Cloning and analysis of the *MNN4* gene required for phosphorylation of *N*-linked oligosaccharides in *Saccharomyces cerevisiae*. Glycobiol. 6:805-810.
- Ohno, N., Yadomae, T. and Miyazaki, T. 1980.** The structure of the type specific polysaccharide of pneumococcus type XIX. Carbohydr. Res. 80:297-304.
- Paton, JC. 1996.** The contribution of pneumolysin to the pathogenicity of *Streptococcus pneumoniae*. Trends Microbiol. 4:103-106.
- Paton, JC. 1998.** Novel pneumococcal surface proteins: role in virulence and vaccine potential. Trends Microbiol. 6:85-87.
- Paton, JC., Andrew, PW., Boulnois, GJ. and Mitchell, TJ. 1993.** Molecular analysis of the pathogenicity of *Streptococcus pneumoniae*: The role of pneumococcal proteins. Ann. Rev. Microbiol. 47:89-115.
- Paton, JC. and Ferrante, A. 1983.** Inhibition of human polymorphonuclear leucocyte respiratory burst, migration and bactericidal activity by the pneumococcal toxin, pneumolysin. Infect. Immun. 41:1212-1216.
- Paton, JC., Lock, RA. and Hansman, D. 1983.** Effect of immunisation with pneumolysin on survival time of mice challenged with *Streptococcus pneumoniae*. Infect. Immun. 40:584-552.
- Paton, JC., Rowan-Kelly, B. and Ferrante, A. 1984.** Activation of human complement by the pneumococcal toxin, pneumolysin. Infect. Immun. 43:1085-1087.

- Paton, JC., Toogood, IR., Cockington, RA. and Hansman, D. 1986.** Antibody response to pneumococcal vaccine in children aged 5 to 15 years. *Amer. J. Dis. Child.* 140:135-138.
- Paulsen, IT., Beness, AM. and Saier, MH. Jr. 1997.** Computer-based analyses of the protein constituents of transport systems catalysing export of complex carbohydrates in bacteria. *Microbiol.* 143:2685-2699.
- Pearce, BJ., Naughton, AM. and Masure, HR. 1994.** Peptide permeases modulate transformation in *Streptococcus pneumoniae*. *Mol. Microbiol.* 12:881-892.
- Pearce, R. and Roberts, IS. 1995.** Cloning and analysis of the gene clusters for the production of the *Escherichia coli* K10 and K54 antigens: identification of a new group of *serA*-linked capsule gene clusters. *J. Bacteriol.* 177:3992-3997.
- Perry, MB., Daoust, V. and Carlo, DJ. 1981.** The specific capsular polysaccharide of *Streptococcus pneumoniae* type 9V. *Can. J. Biochem.* 59:524-533.
- Petit, C., Rigg, GP., Pazzani, C., Smith, A., Sieberth, V., Stevens, M., Boulnois, G., Jann, K. and Roberts, IS. 1995.** Region 2 of the *Escherichia coli* K5 capsule gene cluster encoding proteins for the biosynthesis of the K5 polysaccharide. *Mol. Microbiol.* 17:611-620.
- Phillips, NJ., Apicella, MA., Griffiss, JM. and Gibson, BW. 1993.** Structural studies of the lipooligosaccharides from *Haemophilus influenzae* type b strain A2. *Biochem.* 32:2003-2012.
- Pozzi, G., Masala, L., Iannelli, F., Manganeli, R., Håvarstein, LS., Piccoli, L., Simon, D. and Morrison, DA. 1996.** Competence for genetic transformation in encapsulated strains of *Streptococcus pneumoniae*: two allelic variants of the peptide pheromone. *J. Bacteriol.* 178:6087-6090.
- Presecan, E., Moszer, I., Boursier, L., Cruz Ramos, HC., de la Fuente, V., Hullo, MF., Lelong, C., Schleich, S., Sekowska, A., Song, BH., Villani, G., Kunst, F., Danchin, A. and Glaser, P. 1997.** The *Bacillus subtilis* genome from *gerBC* (311 degrees) to *licR* (334 degrees). *Microbiol.* 143:3313-3328.
- Ramirez, M. and Tomasz, A. 1998.** Molecular characterisation of the complete 23F capsular polysaccharide locus of *Streptococcus pneumoniae*. *J. Bacteriol.* 180:5273-5278.
- Ratcliffe, SW., Luh, J., Ganesan, AT., Behrens, B., Thompson, R., Montenegro, MA., Morelli, G. and Trautner, TA. 1979.** The genome of *Bacillus subtilis* phage SPP1: the arrangement of restriction endonuclease generated fragments. *Mol. Gen. Genet.* 168:165-172.
- Rebers, PA. and Heidelberger, M. 1961.** The specific polysaccharide of type VI pneumococcus. 2. The repeating unit. *J. Amer. Chem. Soc.* 83:3056-3059.
- Reeves, RE. and Goebel, WF. 1941.** Chemoimmunological studies on the soluble specific substance of pneumococcus. V. The structure of the type III polysaccharide. *J. Biol. Chem.* 139:511-519.
- Reizer, J., Reizer, A. and Saier, MH. Jr. 1992.** A new sub-family of bacterial ABC-type transport systems catalysing export of drugs and carbohydrates. *Prot. Sci.* 1:1326-1332.
- Reuber, TL. and Walker, GC. 1993.** Biosynthesis of succinoglycan, a symbiotically important exopolysaccharide of *Rhizobium meliloti*. *Cell* 74:269-280.
- Richards, JC. and Perry, MB. 1988.** Structure of the specific capsular polysaccharide of *Streptococcus pneumoniae* type 23F (American type 23). *Biochem. Cell Biol.* 66:758-771.
- Riley, ID. and Douglas, RM. 1981.** An epidemiologic approach to pneumococcal disease. *Rev. Infect. Dis.* 3:233-245.
- Robbins, JB., Austrian, R., Lee, C-J., Rastogi, SC., Schiffman, G., Henrichsen, J., Mäkelä, PH., Broome, CV., Facklam, RR., Tiesjema, RH. and Parke Jr., JC. 1983.** Considerations for formulating the second-generation pneumococcal capsular polysaccharide vaccine with emphasis on the cross-reactive types within groups. *J. Infect. Dis.* 148:1136-1159.

- Roberts, IS. 1996.** The biochemistry and genetics of capsular polysaccharide production in bacteria. *Ann. Rev. Microbiol.* 50:285-315.
- Roberts, IS., Mountford, R., High, N., Bitter-Suermann, D., Jann, K., Timmis, K. and Boulnois, G. 1986.** Molecular cloning and analysis of the genes for the production of the K5, K7, K12 and K92 capsular polysaccharides in *Escherichia coli*. *J. Bacteriol.* 168:1228-1233.
- Robertson, BD., Frosch, M. and van Putten, JP. 1994.** The identification of cryptic rhamnose biosynthesis genes in *Neisseria gonorrhoeae* and their relationship to lipopolysaccharide biosynthesis. *J. Bacteriol.* 176:6915-6920.
- Rosenow, C., Ryan, P., Weiser, JN., Johnson, S., Fontan, P., Ortqvist, A. and Masure, HR. 1997.** Contribution of novel choline-binding proteins to adherence, colonisation and immunogenicity of *Streptococcus pneumoniae*. *Mol. Microbiol.* 25:819-829.
- Rubens, CE., Heggen, LM., Haft, RF. and Wessels, MR. 1993.** Identification of *cpsD*, a gene essential for type III capsule expression in group B streptococci. *Mol. Microbiol.* 8: 843-855.
- Rubins, JB., Duane, PG., Charboneau, D. and Janoff, EN. 1992.** Toxicity of pneumolysin to pulmonary endothelial cells in vitro. *Infect. Immun.* 60:1740-1746.
- Rubins, JB., Duane, PG., Clawson, D., Charboneau, D., Young, J. and Niewoehner, DE. 1993.** Toxicity of pneumolysin to pulmonary alveolar epithelial cells. *Infect. Immun.* 61:1352-1358.
- Sampson, JS., O'Connor, SP., Stinson, AR., Tharpe, JA. and Russell, H. 1994.** Cloning and nucleotide sequence of *psaA*, the *Streptococcus pneumoniae* gene encoding a 37-kilodalton protein homologous to previously reported *Streptococcus* sp. adhesins. *Infect. Immun.* 62:319-324.
- Sanger, F., Coulson, AR., Hong, GF., Hill, DF. and Petersen, GB. 1983.** Nucleotide sequence of bacteriophage lambda DNA. *J. Mol. Biol.* 162:729-773.
- Sankilampi, U. 1997.** *Streptococcus pneumoniae* and the elderly: an epidemiological and serological study focusing on the potential of the polysaccharide vaccine. PhD. Thesis. University of Oulu, Finland.
- Sau, S. Bhasin, N., Wann, ER., Lee, JC., Foster, TJ. and Lee, CY. 1997.** The *Staphylococcus aureus* allelic genetic loci for serotype 5 and 8 capsule expression contain the type-specific genes flanked by common genes. *Microbiol.* 143:2395-2405.
- Sau, S. and Lee, CY. 1996.** Cloning of type 8 capsule genes and analysis of gene clusters for the production of different capsular polysaccharides in *Staphylococcus aureus*. *J. Bacteriol.* 178:2118-2126.
- Saxena, IM., Brown, RM. Jr., Fevre, M., Geremia, RA. and Henrissat, B. 1995.** Multidomain architecture of beta-glycosyl transferases: implications for mechanism of action. *J. Bacteriol.* 177:1419-1424.
- Saxena, IM., Kudlicka, K., Okuda, K. and Brown, OM. Jr. 1994.** Characterisation of genes in the cellulose-synthesising operon (*acs* operon) of *Acetobacter xylinum*: implications for cellulose crystallisation. *J. Bacteriol.* 176:5735-5752.
- Scott, JAG., Hall, AJ., Dagan, R., Dixon, JMS., Eykyn, SJ., Fenoll, A., Hortal, M., Jetté, LP., Jorgensen, JH., Lamothe, F., Latorre, C., Macfarlane, JT., Shlaes, DM., Smart, LE. and Taunay, A. 1996.** Serogroup-specific epidemiology of *Streptococcus pneumoniae*: associations with age, sex, and geography in 7,000 episodes of invasive disease. *Clin. Infect. Dis.* 22:973-981.
- Severin, A. and Tomasz, A. 1996.** Naturally occurring peptidoglycan variants of *Streptococcus pneumoniae*. *J. Bacteriol.* 178:168-174.
- Schreiber, JR. and Jacobs, MR. 1995.** Antibiotic-resistant pneumococci. *Pediatr. Clin. N. Amer.* 42:519-537.
- Shibaev, VN. 1986.** Biosynthesis of bacterial polysaccharide chains composed of repeating units. *Adv. Carbohyd. Chem. Biochem.* 44:277-339.

- Shoemaker, NB. and Guild, WR. 1974.** Destruction of low efficiency markers is a slow process occurring at a heteroduplex stage of transformation. *Mol. Gen. Genet.* 128:283-290.
- Silhavy, T.J., Berman, ML. and Enquist, LW. 1984.** Experiments with gene fusions. Cold Spring Harbor Laboratory, Cold Spring Harbor, NY.
- Silver, RP., Aaronson, W. and Vann, WF. 1984.** Genetic and molecular analysis of *Escherichia coli* K1 antigen genes. *J. Bacteriol.* 157:568-575.
- Skull, S., Leach, J. and Currie, BJ. 1996.** *Streptococcus pneumoniae* carriage and penicillin/ceftriaxone resistance in hospitalised children in Darwin. *Aust. NZ. J. Med.* 26:391-395.
- Skurnik, M., Venho, R., Toivanen, P. and al-Hendy, A. 1995.** A novel locus of *Yersinia enterocolitica* serotype O:3 involved in lipopolysaccharide outer core biosynthesis. *Mol. Microbiol.* 17:575-594.
- Smith, CJ. 1987.** Nucleotide sequence analysis of Tn4551: use of *ermFS* operon fusions to detect promoter activity in *Bacteroides fragilis*. *J. Bacteriol.* 169:4589-4596.
- Snaideck, DH., Schwartz, B., Lipman, H., Bogaerts, J., Butler, J.C., Dagan, R., Echaniz-Aviles, G., LLOYD-Evans, N., Fenoll, A., Girgis, NI., Henrichsen, J., Klugman, K., Lehmann, D., Takala, AK., Vandepitte, J., Gove, S. and Breiman, RF. 1995.** Potential interventions of childhood pneumonia: Geographic and temporal differences in serotype and serogroup distribution of sterile site pneumococcal isolates from children - implications for vaccine strategies. *Pediatr. Infect. Dis. J.* 14:503-510.
- Soldo, B., Lazarevic, V., Margot, P. and Karamata, D. 1993.** Sequencing and analysis of the divergon comprising *gtaB*, the structural gene of UDP-glucose pyrophosphorylase of *Bacillus subtilis* 168. *J. Gen. Microbiol.* 139:3185-3195.
- Sorensen, UBS., Blum, J., Birch-Andersen, A. and Henrichsen, J. 1988.** Ultrastructural localisation of capsules, cell wall polysaccharide, cell wall proteins, and F antigen in pneumococci. *Infect. Immun.* 43:876-878.
- Sorensen, UBS., Henrichsen, J., Chen, H-C. and Szu, SC. 1990.** Covalent linkage between the capsular polysaccharide and the cell wall peptidoglycan of *Streptococcus pneumoniae* revealed by immunochemical methods. *Microb. Pathog.* 8:325-334.
- Southern, E. 1975.** Detection of specific sequences among DNA fragments separated by gel electrophoresis. *J. Mol. Biol.* 98:503-517.
- Steinfors, C., Wilson, R., Mitchell, T., Feldman, C., Rutman, A., Todd, H., Sykes, D., Walker, J., Saunders, K., Andrew, PW., Boulnois, GJ. and Cole, PJ. 1989.** Effect of *Streptococcus pneumoniae* on human respiratory epithelium in vitro. *Infect. Immun.* 57:2006-2013.
- Stevenson, G., Neal, B., Liu, D., Hobbs, M., Packer, NH., Batley, M., Redmond, JW., Lindquist, L. and Reeves, P. 1994.** Structure of the O antigen of *Escherichia coli* and the sequence of its *rfb* gene cluster. *J. Bacteriol.* 176:4144-4156.
- Stingle, F. and Mollet, B. 1996.** Disruption of the gene encoding penicillin-binding protein 2b (*pbp2b*) causes altered cell morphology and cease in exopolysaccharide production in *Streptococcus thermophilus* Sfi6. *Mol. Microbiol.* 22:357-366.
- Tabor, S. and Richardson, CC. 1985.** A bacteriophage T7 RNA polymerase/ promoter system for controlled exclusive expression of specific genes. *Proc. Natl. Acad. Sci. USA.* 82:1074-1078.
- Takala, AK., Eskola, J., Leinonen, M., Kayhty, H., Nissinen, A., Pekkanen, E. and Makela, PH. 1991.** Reduction in oropharyngeal carriage of *Haemophilus influenzae* type b (Hib) in children immunised with an Hib conjugate vaccine. *J. Infect. Dis.* 164:982-986.
- Takala, AK., Jero, J., Kela, E., Rönnerberg, P-R, Koskeniemi, E. and Eskola, J. 1995.** Risk factors for primary pneumococcal disease among children in Finland. *JAMA* 273:859-864.

- Talbot, U., Paton, AW. and Paton, JC. 1996.** Uptake of *Streptococcus pneumoniae* by respiratory epithelial cells. *Infect. Immun.* 64: 3772-3777.
- Talkington, DF., Brown, BG., Tharpe, JA., Koenig, A. and Russell, H. 1996.** Protection of mice against fatal pneumococcal challenge by immunisation with pneumococcal surface adhesin A (PsaA). *Microb. Pathogen.* 21:17-22.
- Tamara, K. and Nei. M. 1993.** Estimation of the number of nucleotide substitutions in the control region of mitochondrial DNA in humans and chimpanzees. *Mol. Biol. Evol.* 10:512-526.
- Taunay, AE., Austrian, R., Landgraf, IM., Vieira, MFP. and Melles, CEA. 1990.** Sorotipos de *Streptococcus pneumoniae* isolados de liquido cefalorraquidiano no periodo de 1977-1988 na cidade de Sao Paulo, Brasil. *Rev. Inst. Med. Trop. Sao Paulo.* 32:11-15.
- Teele, DW., Klein, JO., Rosner, B. and the Greater Boston Otitis Media Study Group. 1989.** Epidemiology of otitis media during the first seven years of life in children in Greater Boston: A prospective, cohort study. *J. Infect. Dis.* 160:83-94.
- Thompson, JD., Higgins, DG. and Gibson, TJ. 1994.** CLUSTAL W: improving the sensitivity of progressive multiple sequence alignment through sequence weighting, position specific gap penalties and weight matrix choice. *Nucleic Acids Res.* 22:4673-4680.
- Tillet, WS. and Francis, T. Jr. 1930.** Serological reaction in pneumonia with a non-protein somatic fraction of pneumococcus. *J. Exp. Med.* 52:561-571.
- Tomasz, A. 1981.** Surface components of *Streptococcus pneumoniae*. *Rev. Infect. Dis.* 3:190-211.
- Towbin, H., Staehelin, T. and Gordon, J. 1979.** Electrophoretic transfer of proteins from polyacrylamide gels to nitrocellulose sheets : procedure and some applications. *Proc. Natl. Acad. Sci. USA.* 76:4350-4354.
- Troy, FA., Frerman, FE. and Heath, EC. 1971.** The biosynthesis of capsular polysaccharide in *Aerobacter aerogenes*. *J. Biol. Chem.* 246:118-.
- Tuomanen, EI., Austrian, R. and Masure, HR. 1995.** Pathogenesis of pneumococcal infection. *New Engl. J. Med.* 332:1280-1284.
- Tuomanen, EI., Liu, H., Hengstler, B., Zak, O. and Tomasz, A. 1985.** The induction of meningeal inflammation by components of the pneumococcal cell wall. *J. Infect. Dis.* 151:859-868.
- Tuomanen, EI. and Masure, HR. 1997.** Molecular and cellular biology of pneumococcal infection. *Microb. Drug Resist.* 3:297-308.
- van Dam, JEG., Fleer, A. and Snippe, H. 1990.** Immunogenicity and immunochemistry of *Streptococcus pneumoniae* capsular polysaccharide. *Antonie van Leeuwenhoek.* 58:1-47.
- Van Eldere, J., Brophy, L., Loynds, B., Celis, P., Hancock, I., Carman, S., Kroll, JS. and Moxon. ER. 1995.** Region II of the *Haemophilus influenzae* Type b capsulation locus is involved in serotype-specific polysaccharide synthesis. *Mol. Microbiol.* 15:107-118.
- van Kranenburg, R., Marugg, JD. van Swam, II. Willem, NJ. and de Vos, WM. 1997.** Molecular characterisation of the plasmid-encoded *eps* gene cluster essential for exopolysaccharide biosynthesis in *Lactococcus lactis*. *Mol. Microbiol.* 24:387-397.
- Vaughan, EE. and de Vos, WM. 1995.** Identification and characterisation of the insertion element IS1070 from *Leuconostoc lactis* NZ6009. *Gene.* 155:95-100.
- Vázquez, M., Santana, O. and Quinto, C. 1993.** The NodI and NodJ proteins from *Rhizobium* and *Bradyrhizobium* strains are similar to capsular polysaccharide secretion proteins from Gram-negative bacteria. *Mol. Microbiol.* 8:369-377.

- Waltman, WD. II, McDaniel, LS., Gray, BM. and Briles, DE. 1990. Variation in the molecular weight of PspA (Pneumococcal surface protein A) among *Streptococcus pneumoniae*. Microb. Pathogen. 8:61-69.
- Wang, L., Lui, D. and Reeves, PR. 1996. C-terminal half of *Salmonella enterica* WbaP (RfbP) is the galactosyl-1-phosphate transferase domain catalysing the first step of O-antigen synthesis. J. Bacteriol. 178:2598-2604.
- Wang, L., Romana, LK. and Reeves, PR. 1992. Molecular analysis of a *Salmonella enterica* Group E1 *rfb* gene cluster: O Antigen and the genetic basis of the major polymorphism. Genetics. 130:429-443.
- Watson, DA., Kapur, V., Musher, DM., Jacobson, JW. and Musser, JM. 1995. Identification, cloning and sequencing of DNA essential for encapsulation of *Streptococcus pneumoniae*. Curr. Microbiol. 31:251-259.
- Watson, DA. and Musher, DM. 1990. Interruption of capsule production in *Streptococcus pneumoniae* serotype 3 by insertion of transposon Tn916. Infect. Immun. 58:3135-3138.
- Watson, DA. and Musher, DM. 1996. The pneumococcus: Sugar-coated killer. Infect. Med. 13:373-381.
- Watson, DA., Musher, DM., Jacobson, JW. and Verhoef, J. 1993. A brief history of the pneumococcus in biomedical research: A panoply of scientific discovery. Clin. Infect. Dis. 17:913-924.
- Wenger, JD., Hightower, AW., Facklam, RR., Gaventa, S., Broome, CV. and the Bacterial Meningitis Study Group. 1990. Bacterial meningitis in the United States, 1986: Report of a multistate surveillance study. J. Infect. Dis. 162:1316-1323.
- Weiser JN., Austrian, R., Sreenivasan, PK. and Masure, HR. 1994. Phase variation in pneumococcal opacity: relationship between colonial morphology and nasopharyngeal colonisation. Infect. Immun. 62:2582-2589.
- Weiser, JN., Love, JM. and Moxon, ER. 1989. The molecular mechanism of phase variation of *H. influenzae* lipopolysaccharide. Cell 59: 657-665.
- Whitfield, C. 1995. Biosynthesis of lipopolysaccharide O antigens. Trends Microbiol. 3:178-185.
- Whitfield, C., Amor, PA. and Koplín, R. 1997. Modulation of the surface architecture of Gram-negative bacteria by the action of surface polymer: lipid A-core ligase and by determinants of polymer chain length. Mol. Microbiol. 23:629-638.
- Wren, BW., Russell, RR. and Tabaqchali, S. 1991. Antigenic cross-reactivity and functional inhibition by antibodies to *Clostridium difficile* toxin A, *Streptococcus mutans* glucan-binding protein, and a synthetic peptide. Infect. Immun. 59:3151-3155.
- Yother, J. and Briles, DE. 1992. Structural properties and evolutionary relationships of PspA, a surface protein of *Streptococcus pneumoniae*, as revealed by sequence analysis. J. Bacteriol. 174:601-609.
- Yother, J., Leopold, K., White, J. and Fischer, W. 1998. Generation and properties of a *Streptococcus pneumoniae* mutant which does not require choline or analogues for growth. J. Bacteriol. 180:2093-2101.
- Yother, J., McDaniel, LS. and Briles, DE. 1986. Transformation of encapsulated *Streptococcus pneumoniae*. J. Bacteriol. 168:1463-1465.
- Zhang, L., al-Hendy, A., Toivanen, P. and Skurnik, M. 1993. Genetic organisation and sequence of the *rfb* gene cluster of *Yersinia enterocolitica* serotype O:3: similarities to the dTDP-L-rhamnose biosynthesis pathway of *Salmonella* and to the bacterial polysaccharide transport systems. Mol. Microbiol. 9:309-321.

# APPENDIX I

## The nucleotide and deduced amino acid sequence of the distal portion of the *cps19f* locus.

The nucleotide sequence is numbered in accordance with GenBank accession number U09239 and is shown from nucleotide 5,851, within the distal portion of *cps19f* to the *HindIII* site within *aliA*. The amino acid translation for each ORF is represented by single letter code above the first nucleotide of each codon. Possible ribosome binding sites are underlined.

```
...cps19f→                cps19fG→
K R G V K Y *                M I R L I Q K V E L D A I K E F K K I C E E N
5851 AAAAGAGGTTGATAAATTTGAACAGACAGAGATGATTCGCTTAATTCAAAAGTTGAATTAGATGCTATAAAAGAGTTAAAAAATCTGTGAAGAGAA
  D I D F F L R G G S V L G A V K Y D G F I P W D D D M D I A V P R
5951 TGATATAGATTTTTCCCTCCGGTGGTAGTGTACTTGGTGCAGTCAAATACGACGGCTTTATCCATGGGATGATGATATGGATATCGCTGTCCTCGT
  E A Y D K L P S V F K D R I I A G K Y Q V L T Y Q Y C D T L H C Y F
6051 GAAGCATACGACAAACTTCCAAGTGTTCCTCAAAGATAGAATTATCGCTGGGAAATATCAGGTTCTTACTTATCAATACGTGATACGTTGCATTCGTACT
  P R L F L L A D E R K R L G L P R N T N L G L H L I D I I P L D G
6151 TTCTCGACTATTCCTTTTAGCAGATGAAAGAAAACGTTTGGGCTTGCCACGAAATACCAATCTAGGATTCGATTGATTGATATCATTCCCTTAGATGG
  A P N H S V L R K I Y F C K V Y W Y R F L A S L G T T Y V G D H V
6251 AGCACCAATCATTCCGTTTTAAGAAAAGATTTACTTTTGTAAAGTATACTGGTATCGTTTTTTAGCAAGCTTAGGAACAACCTTATGTTGGCGACCATGTG
  D M H S T K Q K L I I G F F K K L G F A K L F P Q N S V Y R R L D N
6331 GATATGCAATCCACTAAGCAAAAATAATTATGGTTTCTTTAAAAAAGTAGGATTTGCAAAAATTTCCCTCAAAAATTCGTATACAGACGCTTGGATA
  L Y R K Y D W K K Q K Y A G T I N A S L F A K E V M P V E I W G E
6451 ATCTCTATAGAAAATGATGATGGAAAAGCAGAAGTATGCGGGGACTATCAATGCTTCTTTATTTGCTAAAAGAGTTATGCCAGTAGAGATTTGGGGAGA
  G V E T F F E D T F F K V P T E Y D R Y L K R L Y G E N Y L Y E E
6551 AGGAGTAGAGAAGCCTTTTGAGGATACCTTCTTTAAAGTTCCAACGGAGTATGATCGCTACCTGAAAAGACTTTACGGAGAAAATATCTTTACGAAGAG
                cps19fH→
P S D D E A E K K S H L G G H * M F C Y I I L H Y K V L E E T I S C V
6651 CCTAGTGTATGATAAAGAAATCGCATTTAGGAGGACACTAATTTGTTTTGTTATATATTTTGCATTACAAAGTCTTAGAAGAAAATTTTCTTGTGTT
  K S I K E G N Y N A K Q I V I I D N F S N N G T G E K L Q E L Y E S
6751 AAATCTATAAAGAAGGCAATTATAATGCAAAGCAAATCGTTATATGATAATTTCTCTAATAATGGTACTGGTGAAAACACAAAGAGCTTTATGAAT
  D L E I D V L I N H E N A G F A R G N N V A Y Q F A K E K Y N P D
6851 CAGATTTAGAAATGATGTTTTGATTAACCATGAAAATGCTGGTTTTGCTCGTGAAAATAATGTGGCTTATCAATTTGCTAAGGAAAAGTATAACCCCGA
  F M V I M N N D I E I E T E N F E K I V T D I Y R E E K F H L L G
6951 TTTCTATGGTTATCATGAATAACGATATTGAGATAGAAAACGAAAATTTTGA AAAAATTTGTCAGATATCTATCGTGAGGAAAAATTCCTATTGCTCGGG
  P D I F S T T Y Q L H Q N P K R L T H Y T Y G E V K A L N E K F K K
7051 CCAGATCTTCTCGACTACTTACCAACTTCACCAAAACCCAAAACGTTGACACATTTATACTTATGGAGAAAGTTAAAGCTCTAAAATGAAAAATTTAAAA
  G S Q V S L A L K I K C W L K A S K V L R T A I Y Q N R R K K G S
7151 AAGGGAGCCAAGTTAGTCTAGCTTTAAAAATCAAATGTTGGTTGAAAAGCTAGTAAAGTTCTTCGAAACAGCAATCTATCAAAAATAGACGTA AAAAAGGATC
  V D Y R K Q V E N P I L H G S F I V Y S R D F I E K E E Y A F N P
7251 AGTAGACTATAGAAAACAGGTAGAAAACCCAATCTTCCTGGTTCTTTTATTGTATATTCGAGAGATTTTATCGAAAAGAGGAGTATGCTTTTAAACCTT
  N T F F F Y Y E T E I L D Y E A E L K G Y K R I Y T P K I R V L H H Q
7351 AACACCTTCTTTTACTATGAAACAGAGATATTAGATTGAAGCTGAAATTTAAAAGGATACAAAGAAATTTATACACCTAAAATAGAGTTTTCACCATC
  N V A T N Q V Y T N L L E K T L F S N K C N F K S T S Y F L K L M
7451 AAAATGTTGCAACTAATCAAGTTTACACGAACTTGTTAGAAAACCTGTTTTTCAAATAAATGCAACTTTAAATCCACCGATTATTTTTGAAGTTGAT
                cps19fl→
K E N E D V * M S Y L F L L C L T L F L L T I F Y F F A F I Q D L
7551 GAAAGAAAACGAGGATGTTTTAAATGAGTTATTTTACTTTGCTTACATTATCTTATTGACTATATTCTATTCTTTGCTTTTATTCAAGATTA
```

I A P P V V M S V M F L I S S V F A L V N S K N W N I E Y S G I A Y  
 7651 ATTGCTCCTCCAGTAGTATGTCTGTAATGTTTCTAATAGTTCAGTATTTGCAGTGGTTAAATCAAAAACTGGAATATGAAATATAGTGAATAGCCT  
 I L I I S G I I I F S I P L M A A L K S P N F N T E V K I A D R L I  
 7751 ATATTCATAAATAGTGGTATATTATATTTTCGATTCCTTTAAATGGCAATAAAATCACCTAATTTAATACAGGTTAAGATGCTGATCGATTAAT  
 D I Q F W K I A L T I I I D L F I L Y L Y R K E I Y N L V L S N G  
 7851 TGATATCAAAATTTGGAAAATGCTCTAACTATTATAATGATCTCTTTATTTTGTATCTTTACAGGAAGGAAATATACAACCTTGTCTTAGTAAATGGA  
 Y T G S N I Q W F F R N A T S Y E G E L T V R T F I R V L I R V I D  
 7951 TATACGGGGTCAAAATTTTCAGTGGTTTTTAGAAAATGCAACGAGTATGAAGGTGAATGACAGTGCAGACTTTTATTCGAGTTCTCATTCTGTTATTG  
 V S A Y I F G Y T F I N N F L I Y R H K R P K D I L L L V P L L I  
 8051 ACGTATCTGCTTATATTTTGGATATACTTTTATTAATAATTTCTTATCTATCGCCATAAACGCCCTAAAGACATATTACTTTTAGTACCTTTATTAAT  
 F I S G K T L I S G R Q D I I K I L I A Y I V I M M Y I Q Q K R K V  
 8151 ATTTATTTCAAAAACTTTAATAATCAGGAGCGCCGGCAAGATATTAAATTTCTGATTGCCTATGTAATCATGATGTATATCCAACAAAAACGGAAAGTT  
 G W N R V I S H K A Y I H L G F V G L I A G I P A F Y Y S L F L A G R  
 8251 GGATGGAATAGAGTATATCTCATAAATATATCACCTTGGATTTGTTGGTTTAAATAGCAGGTATCCAGCATTTTACTACTCTTTTGGTTTTCGCGGTC  
 S T T R T L F E S V S T Y L G G S I Q H F N Q Y I E N P L D P G E  
 8351 GTTCAACGACTAGGACGCTATTTGAGAGTGTTCGACCTATCTAGGAGGCTCAATTCAGCATTTTAAATCAGTATATTGAAAATCCATTAGATCCGTGGTGA  
 V F G S E T L V P I L N I L G E M G L V N Y R S T I H L E F R T L  
 8451 AGTTTTTGGCAGTGAACATTTGGTGCCTATATTAATATATAGGGGAAATGGGCCAGTAAATATCGTAGTACAATTCATTTAGAATTTCCGGACACTA  
 G V T V G N V Y T F F R R P L H D F G L V G F M Y V F V F A V G A F F  
 8551 GGAGTACTGTAGAAAATGTTTATACCTTTTTTGAAGACCCCTGCATGATTTTGGTCTAGTTGGTATGTATGTATTGCTTTGCTGTAGGTGCTTTTT  
 A I Y Y L V L R K K Q V G F N L D I H T I I Y S Y V F Y W I F L S  
 8651 TTGCTATTTATTTAGTCTGAGAAAACAGGTGGTTTTTAAATTTGGATATTCATACCATTATTTATCTTATGCTTTTATTTGGATTTTATTTATC  
 S I E Q Y S F T M I S L Y T L V F I V L V Y F M A I F Y W C T D F  
 8751 ATCAATCGAGCAATCTCGTTCACAAATGATTAGTCTATATACACTTGTATTTATTGTTGGTTTACTTTTATGGCTATCTTTTACTGGGTACAGATTTT

*cps19fJ*→

K R G K L I F K I S D S S I K L K E E \* M N T K I K  
 8851 AAAAGAGGAAAAGTGAATTTTAAAAATTTCTGACTCAAGTATCAAAATAAAAGAGAATAACAGAATGTATAGGAGAGGGTAGAATAACTAAAAATAAA  
 N I I T S F S Y V I S S N L L I V L T S S S L V V L I V P K I M G V T  
 8951 AATAATAAASACTAGTTTTCTTATGTTATTTCTTCAAATCTGCATAGTTTAACTCATCACTAGTTGTTTGTATGTTTCTAAAAATAATGGGGGTAA  
 E Y S Y W I Y I F Y L T Y I G F F H L G W L I D G I W L K Y G G L  
 9091 CTGAGTACAGTTACTGGCACTTTATATTTTTATCTGACCTATATCGGTTTTTCCACTTGGGTTGGATTGATGGGATTTATCTCAAATATGGTGGCTT  
 E Y T N L D R K Q F Y S Q M I L F S S F L M L I S L V L F T L N  
 9151 AGAATATACAAATTTAGATAGAAAACAGTTTTATCTCAGATGATTTCTATTTCTTAAATGCTAATCTCGCTGATTTATTTACTTTGAACTTA  
 I T V R D E N A R Y I Y N M A I I S M T V T N L R T L V Y I L Q M  
 9251 ATAACGTAAAGGATGAAAACGCAAGATATATTTATAATATGGCTATCATCAGCATGACAGTCAAAAACCTAAGAACACTCTATGTTTATATCTTCGAGA  
 T N R L K D S S V I L I S D R V L Y V L L F M F I V F G W H E Y  
 9351 TGACAAATCGCTTGAAGGATGTTTCAGTCAATCTAATAGTGCCTTTTATATGTAACCTTTTATTCATGTTTATTTGATTTGGATGGCATGAGTA  
 K V M I W A D I L G R T F S L C M L S F W I C K D I V F Q P L S K F  
 9451 CAAGGTCATGATTTGGGCTGATATTTCTAGGTCGAACATTTCTCTCATGCTTTCCTCTGGATTGTAAAGATATTTGGTTTCAGCCCTTTGTCAAAATTT  
 I L D F K E S L D N I R V G I N L M L S N I A S S M I I G I V R M G  
 9551 ATCTTAGATTTCAAGGATCCCTTGATAAATATCCGCTGTGGAAATCAATTAATGCTATCTAACATTCGCGAGTAGCATGATTTAGGCATTTGTTCTGATGG  
 I Q W N W N I E T F G K V S L T L S I S N L L M T F I N A I G L V  
 9651 GAATTCAAATGGAATTTGAAACATTCGGGAAAAGTGCATTAACCTTTGAGTATATCTAATTTTATTAATGACTTTTATTAATGCAATTTGGATTTAGT  
 I F P L I K R T K T E N P K I Y S N L R N A L M L V M F A I L T  
 9751 TATCTTTCTTTGATAAAGCGAACAAGACTGAGAAATTAACCTAAAATTTATCTAATTTAAGAAAATGCTTTGATGTTGGTTATGTTTCGCAATCTTGCTC  
 F Y Y P L K F I L D I W L P A Y K D A L V F M A L I F P M S V Y E G  
 9851 TTCTACTACCTTTAAAATTTATTTCTTGATATTTGGCTTCTGTATAAGGATGCCTTAGTTTTCATGGCCCTAATTTTCTTATGTCAGTTTATGAAG  
 K M A L V I N T Y L K A M R M E K D I L K I N A L V S I V V  
 9951 GGAAAATGGCTTTGGTGATAAATACATATTTAAAAGCAATGAGAAATGAAAAAGACATTTCTAAAATTAATGCTTTGGTTATGTTAAGTAGTATAGTAGT  
 T L V T T L L L N N L G L T V V S I V I L L A L R S I I A E L I L  
 10051 GACATTAGTACTACTACTAAATAATTTGGGGCTGACAGTTGTATCTATAGTTATTTACTTGGCTTAAAGAAGTATAATAGCTGAATTAATTTTA  
 S K K L K I S V K Q D I A L E L L E T M T P I I F I S S S W I I A  
 10151 TCCAAAAAAGTGAAGATATCAGTCAAGCAAGACATTTAGAGTTACTTATGACGATTAATTTATTTCTTCAAGTTGGTATCTCTCTATTTGGATTG  
 V I I Y L L A Y T L Y L Y L K H K D I R M Y I E Y F K N H K K I S  
 10251 CAGTAATAATTTATTTATTTGGGTATACTTTATATTTGTATTTAAAGCACAAAGATATCAGAATGTATATAGAAATCTTTAAAATCATAAAAAATATC

*cps19fK*→

\* M K K I M L V F G T R  
 10351 ATAAAAATATATATCAATGAAATGGTAGATTACATTTCTACCTTTTTATCCATTTAGGAGGAAACGATGAAAAAGATAATGCTAGTTTTCGGTACACAGT  
 P E A I K M C S L V N E L R K Q E D M K T V V C V T G Q H K E M V S  
 10451 CCAGAAGCAATAAAAAATGTTTCATTTAGTCAATGAGTTGAGAAAACAGGAGATATGAAAACAGTTGTTTGTGTAAGTGGTCAACACAAGGAGATGGTTA  
 P V L D L F G V Q P D Y D L E I M K A N Q N L F S I T I S I L E K  
 10551 GCCCTGTTTGGATTTATTTGGAGTTCAACGACTATGATTTAGAATATGAAGGCTAATCAAACTTGTCTCTATAACAATAAGTATTTAGAAAA  
 I K P V L E K E Q P D I V L V H G D T T T T Y A A A L A A F Y L G  
 10651 AATAAAACCTGTGTAGAAAAGAACAACAGATATTTGTTTGGTTCACGGTGATACTACTACGACATATGCAGCAGCTTTAGCGGCATTTTATCTAGGA  
 I K V G H V E A G L R T Y N L Q S P P P E E F N R Q S T S I I A N Y  
 10751 ATTAAGTTGGTCAATTTGAAGCTGATACGAACCTGCAAGCTTCCAGAGAGTTCACAGACAATCGACTTCTATATTGCAAAAT  
 H F A P T E L A K E N L I K E G R N N I Y V T G N T V I D A L T T  
 10851 ATCATTTTGTCTTACAGAATTAGCTAAAAGAAAATCTAATAAAAAGAGGGTGAAGAATAATATCTATGTTACTGGAATACAGTGATGATGCACTTACAAC  
 T V Q K D Y T H P D L D L N D G N R L I L L T A H R R E N N L G E P  
 10951 TACAGTACAAAAGGATTAACACACCCGATCTAGATTTAAACGATGGAATTCGCCCTCATCTTATGACTGCTCATAGACGCGAAAATCTCGGAGAACCT  
 M R F M R F R A V K R V L N E Y D D V K V I Y P I H K N P L V R E T A  
 11051 ATGAGACATATGTTTAGAGCTGTTAAACGAGTTTTAAATGAATATGACGATGTTAAAGTAATTTATCCAATTCATAAAAATCCATTTGGTACGGGAAACAG  
 T E I F G D T E R I Q I I E P L D V L D F H N F M N H S Y M I L T  
 11151 CTACAGAAAATTTGGAGATACAGAACCTTATTCAGATCTGAACTTTAGATGTTCTGATTTTCAACCTTTATGAATCATAGTTATATGATTTAATC  
 D S G G V Q E E A P S L G K P V L V M R D T T E R P E G V A A G T  
 11251 TGACTCAGGAGGGTCCAAGAAGAGGACCTTCGTTAGGAAAAGCTGTATGGTTCATGCGAGATACGACAGAAAACCTGAAAGGAGTAGCTGCTGGAACG  
 L K L V G T D E E T I Y Q N F K M L L D D P E E Y K K M S R A S N P  
 11351 TTTGAAATTTGGTGGACTGATGAGGAGACTATTTATCAAACTTTAAGATGCTTTTAGACGATCCCGAAGAAATATAAAAAATGAGTCCGAGTAGTAATC

*cps19fL*→

Y G N G D A S K O I V R I L R G I \* M K  
 11451 CTTATGGAATGTTGATGCTAGTAAACAGATGTTTCGAATTTTACGTGGAATTTGAGTGTGTTTCAGATAAAGTAAATATAGAAAGTACCTTACTATGAAA  
 G I I L A G S G T R L Y P L T R A A S K Q L M P V Y D K P M I Y Y  
 11551 GGTATTTACTAGCAGTGGTTTCGGGACAGCTTTATATCTTTGACTCGCGTGCATCAAAACAACTTATGCCGGTTTATGATAAACCGATGATTTACT



16051 CTCAAAGTTTTGAAGGAGCTAAAGCAAGAGCTATTATTATGAGCTTATTGGAAAACAGCTAAACGTCATCAACTAAATAGTGAGAAATATCTATCCTATCT  
 16151 TCTAGAATGTCTTCCAAACGAGGAAACTCTCGTAAACAAAGAGGTTTTAGAGGCCTATTTACCATGGACTAAAGTTGTACAAGAAAAGTGCAAATAAGAA  
 16251 ATCTCCAGATTAGGAACATCCGTGAGTTCACATAATCTGGAGATTTTTCAATAGACCTCGTTATTGGGCGGTTACGATATTCATATTTTTTGCAAAGATG  
*aliA*→  
M K  
 16351 TTGTTTGAAAAATAATTTCAAAAAATCTGAAAAATCTGTTGACATCTTTCTGAAAAGAGTCTATAATGGAGAGAAAAGTTTTAAAGGAGAAAAATGATGAA  
 S S K L L A L A G V T L L A A T T L A A C S G S G S S T K G E K T  
 16451 AAGTTCAAAACTACTTGCCTTGGCGGCGTGACATTATGGCGGCGACTACTTTAGCTGCATGCTCTGGATCAGGTTCAAGCACTAAAGGTGAGAAGACA  
 L S Y I Y E T D P D N L N Y L T T A K A A T A N I T S D V V D G L L  
 16551 TTATCATAcATTTATGAGACAGACCCTGATAACCTCAACTATTTGACAACCTGCTAAGGCTGGACAGCAAATATTACCAGTACGCTGGTTGATGGGTTGC  
 E N D R Y G N F V P S M T E D W S V S K D G L T Y T Y T I R K D A  
 16651 TAGAAAATGATCGCTACGGGAACTTGTGCCGCTCTATGACTGAGGATGGTCTGTATCCAAGGATGGGTTGACTTACACTTATACTATCCGTAAGGATGC  
 K W Y T S V G V E Y A A V K A Q D F V A G L K Y A A D K K S D A L  
 16751 TAAATGGTATACTTCTGTGGGTGTAGAATACGGCGCAGTCAAGGCTCAGGACTTTGTGGCAGGCTTAAGTATGCTGCTGATAAAAAATCAGATGCTCTT  
 Y L V Q E S I K G L D G Y V K G E I K D F S Q V G I K A L D E Q T V  
 16851 TACCTTGTCAAGAATCAATCAAGGATGGATGGCTATGTCAAAGGGGAAATCAAAGATTTCTCTCAAGTAGGAATTAAGGCTCTGGATGAACAGACAG  
 Q Y T L N K P E S F  
 16951 TTCAGTACACTTTGAACAAACCAGAAAGCTT

# APPENDIX II

## The nucleotide and deduced amino acid sequence of the serotype specific region of the *cps19b* locus.

The amino acid translation for each ORF is represented by single letter code above the first nucleotide of each codon. Possible ribosome binding sites are underlined. The GenBank accession no. is AF004325.

...*cps19bF*→  
P L H L M G V N A D K I N Q C H T D E K I K K I V N E S G I I N A D  
1 C C A C T A C A C T T G A T G G G G G T G A A T G C T G A T A A A A T T A A T C A G T G C C A T A C A G A T G A G A A A A T C A A A A A A T C G T T A A T G A G T C A G G A A T C A T T A A T G C G G  
G A S V V L A S K F L G T P V P E R V A G I D L M Q C L L E L S N  
101 A T G G A G C A T C A G T T G T T C T T G C A A G T A A G T T T T A G G A A C G C C T G T T C C T G A A C G A G T A G C G G G T A T T G A T T G A T G C A A T G T C T T T A G A G T T G T C A A A  
K K G Y S V Y F F G A K E E V L Q D M L K V F K R D Y P N L I V I  
201 T A A A A A G G A T A T T C A G T T T A C T T T T T G G A G C A A A A G A A G A G T T T T G C A A G A T A T G C T C A A A G T A T T T A A G A G A G A T T A C C A A A T T G A T A G T T A T T  
G H R N G Y F S E E D E Q A I Q E D I R E K N P D F V F I G I T S P  
301 G G A C A C A G A A A T G G C T A T T T T C T G A A G A G G A T G A A C A A G C T A T T C A A G A A G A T A T T C G T G A A A A G A A C C C T G A T T T T G T G T T A T T G G A A T T A C G T C T C  
K K E Y I I Q K F M D S G V N S V F M G V G G S F D V L S G H I Q  
401 C T A A A A A G A A T A T A T T A T T C A A A A A T T A T G G A T A G T G G C G T C A A T T C G G T A T T A T G G G A G T T G G C G G T A G T T T T G A T G T C T T G C T G G T C A T A T C C A  
R A P L W M Q K S N L E W L F R V A N E P K R L F K R Y F V G N I  
501 A C G A G C A C C T C T A T G G A T G C A A A A G T C A A A T T A G A G T G G T T A T T C C G T G T A G C T A A T G A G C C T A A A C G T C T C T T A A A C G T A T T T T G T A G G G A A T A T T

*cps19bG*→  
M I R L I Q K V E L D A I K  
601 T C A T T C A T A G G A A A A G T T T T A A A A G C A A A A A G A G G T G T A A A A T A T T G A A C C A G A C A G A G A T G A T T C G C T T A A T T C A A A A A G T T G A A T T A G A T G C T A T A A A  
E F K K I C E E N D I D F F L R G G S V L G A V K Y D G F I P W D  
701 A G A G T T T A A A A A A T C T G T G A A G A G A A T G A T A T A G A T T T T T C C T C C G C G G T G G T A G T G T A C T T G G T G C A G T C A A A T A C G A C G C G T T A T T C C A T G G G A T  
D D M D I A V P R E A Y D K L P S V F K D R I I A G K Y Q V L T Y Q  
801 G A T G A T A T G G A T A T C G C T G T C C C T C G T G A A G C A T A C G A C A A A C T T C C A A G T G T T T T C A A A G A T A G A A T A T C G C T G G G A A A T A T C A G G T T C T A C T A T A C  
Y C D T L H G Y F P R L F L L E D E R K R L G L P R N T N L G L H  
901 A A T A C T G T G A T A C G T T G C A T G G C T A C T T T C C T C G A C T A T T T C T T T A G A A G A T G A A A G A A A C G T T T G G G C T T G C C A C G A A T A C C A A T C T G G G A T T G C A  
L I D I I P L D G A P N H S V L R K I Y F G K V Y W Y R F L A S L  
1001 T T T G A T T G A T A T C A T T C C T T T A G A T G G A G C A C C A A A T C A T T C G G T T T T A A G A A A G A T T A C T T T G G T A A A G T A T A C T G G T A T C G T T T T T A G C A A G T T T A  
G T T Y V G D H V D M H S T K Q K L I I G F F K K L G F A K L F P Q  
1101 G G A C A A C T T A T G T T G C G A C C A T G T G G A T A T G C A T T C C A C T A A G C A A A A C T A A T T A T T G G T T C T T T A A A A A C T A G G A T T T G C A A A C A T T T T C C C C  
N S V Y R R L D N L Y K K Y D W K K Q K Y A G T I N A S L F A K E  
1201 A A A A T T C T G T A T A C A G A C G C T T G G A T A A T C T C T A T A A A A G T A T G T T G G A A A A G C A G A A G T A T G C T G G G A C T A T C A A T G C T T C T T A T T T G C C A A A G A  
V M P V E I W G E G V E K P F E D T F F K V P T G Y D R Y L K R L  
1301 A G T T A T G C C A G T A G A G A T T T G G G G A G A G G A T A G A G A G C C T T T T G A G G A T A C C T T C T T T A A A G T T C C A C G G G T A T G A T C G C T A C C T G A A A G A C T T

*cps19bH*→  
L F C Y I I L H Y K  
1401 T A C G G A G A A A A C T A T C T T C A C G A A G C C G A G T G A T G A T G A A A A G A A A T C G C A T T T A G G G G G A T A A A A A T T T G T T T G T A T A T A T T T T A C A T T A C A A A  
V L D E T I S C V K S I K E G N S N E K Q I V I I D N F S N N G T G  
1501 G T C T T A G A T G A A A C T A T T T C T T G T T A A A T C T A T A A A A G A G C A A T T C C A A T G A A A A G C A A A T C G T T A T T A T T G A T A A T T T T C A A A A T A A T G G T A C A G  
E K L Q E L Y E S D L E I D V L I N H E N A G F A R G N N V A Y Q  
1601 G T G A A A A C T A C A A G A G C T T T A T G A A T C A G A T T A G A G A T T G A T G A T T A A C C A C G A A A A T G C T G G A T T T G C T G A G G A A A T A A T G T A G C T T A T C A  
F A K E K Y N P D F M V I M N N D I E I E T E D F E K I V T D I Y  
1701 A T T T G C T A A G G A A A G T A T A C C C C G A T T T C A T G G T T A T C A T G A A T A A C G A T A T T G A G A T A G A A A C A G A A G A T T T T G A A A A A A T C G T G A C A G A T A T C T A T  
H K E K F H L L G P D I F S T T Y Q L H Q N P K R L T H Y T Y E E V  
1801 C A C A A A G A A A A T T T C A T T T G C T A G G G C C A G A C A T T T T T C G A C A C A T T T C A G C T A C A T C A A A A T C C T A A A C G C C T G A C A C A T T A T A C T T A T G A A G A G G  
V A L N E K F K R G S Q L S L T L K I K C W L K A S K V L R T A I  
1901 T A G T G G C T C T C A A T G A A A A T T T A A A A G A G G A G C C A A T T A G T C T A A C T T T A A A A A T C A A A T G T T G G T T G A A G C T A G C A A A G T T T C G A A C A G C A A T  
Y Q N R R K K G S V D H R K Q V E N P I L H G S F I V Y S R D F I  
2001 C T A T C A A A A T A G A C C T A A A A A G G A T C A G T A G A C C A T A G A A A A C A G G T A G A A A A C C C G A T T C T C A T G G T T C C T T T A T T G T A T A T T C G A G A G A T T T T A T T  
E K E Y A F N L N T L F Y Y E T E I L D Y E A E L K G Y K R I Y T  
2101 G A A A A G A G G A G A T A T G C T T T T A A C C T T A A T A C C C T T T T A C T A T A G A A C A G A G A T A T T A G A T T A T G A A G C T G A G T T A A A A G G A T A C A A G A G A A T T A T A  
P K I K V L H H Q N V A A C N Q V Y T N L V E K T L F S N K C N F E  
2201 C A C C G A A G A T T A A G G T T T A C A C C A A T C A A A A T G T G G C A C A A A T C A G G T T T A T A C A A A C T T A G T A G A A A A A C T T T G T T T T C A A A C A A A T G C A A C T T T G A

*cps19bP*→

S T S Y F L T L M E K N E \* M M K K V L Y V T N V D W N W I K Q
2301 ATCCACTAGTATTATTTTAACTTTAATGGAGAAAAATGAATAAAAAATGATGAAAAAGTTTATATGTGACAAATGTTGATGGAAATGGATAAAACAA
R P Q F I A E N L S N F Y E M L V L Y R Y W Y N R K G L T E D R N T
2401 CGTCCCAATTTATTCGAGAAAACTTATCTAATTTTATGAGATGCTAGTCTGTGATTCGATTTGGTACAATGAAAAAGGATTTGACTGAGGATGAAAAATA
N I T N I S R I Y A L P F V N R S P K L K Q L N D K I V A W N I R
2501 CTAACATTAACAATATATCAGTATTTATGCCCCCTTTGTAAATAGTACCCGAACTAAAGCAACTGAATGATAAAAATTTGCTTGGAAATATTCG
R K V K A F P E Y V Y L T N P M Q F A S L V D N S E Q T K I I Y D
2601 AAGAAAAAGTTAAGGCTTTTAAACCGGAATATGTATATTTGACAAATCCAATGCAGTTTGCATCTCTTGTAGACAATTCAGAACAAAAATAATATATGAT
C M D Y H V A F I E N R E E R Q R L K D L E E K L V N R A N I L I L V
2701 TGTATGGATTATCATGTGGCTTTTATAGAAAAATAGAGAAGACGTCAGCGATTAAAGGATTAGAAGAAAACTGGTCAATAGAGCTAATTTAATCTTGG
S S E K L R E N I I S D Y N L E E Q V V N K I V V V R N G Y N G K I
2801 TTTCCGAGTGAAGAACTAAGAGAAAAATTTTATGATTAATACTTTGGAAGAGCAGTTAATAAAAATAGTGGTTGTTAGAAATGGTTAATGGTAAAAAT
L S I P T R H K K N N Q K L V L A Y V G T I S H W F D F D I I L R
2901 CTTAAGCATCTTACTGCGATAAAAAGAATAACCAAAGCTTGTACTTGCATATGTTGGAACATCAGTCATTTGGTTGATATCATTATTACGA
S L K D F D N I E Y N L I G P I S K A D I P E H D R I H Y L G S V P
3001 AGTTTAAAAAGTTTGGATAAATATGAATTAATTTGATGGTCCGATTAAGCAAGCTGATATTCTGAACTGATAGAAATCATTATTAGGAAGTGTAC
H E K I Y Q Y I E N A D V L I M P F Q I N D I V E A V D P V K L Y
3101 CACACGAGAAGATTATCAGTATATTGAGAATGCAGATGTTCTGATTATGCCGTTTCAAATTAATGATATTGTTGAAGCAGTGGATCCAGTCAAGTGTGA
E Y I N F K K N I L T V C Y K E I L R F E P F V Y M Y N Y L D Y
3201 TGAATATATTAATTTAAAAAAAATATACCTTACGGTATGTTATAAGGAAATCTGAGATTGAAACCAATTTGATACATGTTCAAAATTTATTAGATTAC
Q M N L L Q L I E N N N L K Y D S I A R E D F L K S N A W E K R A E
3301 CAAATGAATTTGTTGCAATGATTGAAAAATAAATTTGAAATATGACAGTATAGCTAGAGAAGATTTTTTGAAAAGTAATGCTTGGGAAAAAGACGGC

*cps19bI*→

L I H Q L I N Q L \* M K R Q K F E F I E I L Y Y F T V
3401 AGCTGATTCATCAGCTGATTAACCAATTTGTAATTTTATGAGAGAGATGAATGAAGAGACAGAAGTTTGAATTTATAGAGATTTCTATACTATTTTACAGT
M L S V G M F L M F T L S L Y W H R N L L T I L S I A L S F L M L
3501 GATGTTATCAGTGGGAATGTTCTTATGTTTACCCTCAGCTTATATTGGCAGAAAATTTAACTATTTTATCTATGCTCTCCTCATTTTTAAATGCTT
P I L I V N A K R I S K S A F I Y G T F L S I C I I Y E I L R A K T
3601 CCTACTGATTTGTAAGCTAAAAGAATTTCTAAATCTGCTTTTATCTATGTTGACTTTTTTATCTATATGATATTAAGAGCTAAAA
L Y N Y S V S N I F L A S R Q Y I W I F L F F V L I Y L F K N K Q
3701 CACTGTATAATTACAGTGTGAGTAATATTTTGGCCCTAGACAATATATGATTTTCTATTTTGTATTGATTTACCTTTTAAAAACAAACA
E N M R K I L D N T L N I F M F S L G I R A F T W F L Y T L F Q V
3801 AGAAAAATGAGAAAAATTTAGATAATACACTCAATATTTTATGTTTCTCTCGGAATTAGAGCAATTTACTTGGTTTTTATATACGTTATTTCAAGTT
E L F P S I L R E F G L W Y R N R L F S V R I D G T P I I G L L
3901 GAATTTATTTCCATCTATTTTAAAGAGAATTCGGAGATTTGTTGATCCAAATGAATTTTTCAGTACGAATAGATGGAACACCATAATTTATAATAGTGTGT
I S T F Y F K F G N R K Y F Y I F L I L M Y I T F V N Q T R V
4001 TAATTTCCACTTTTTCTATTTTAAATTTGGAATAGGAATCTTTTATTTATTTGTTCTTGATTTAATGATATTAACATTTGTAATTCAGACAAGAGT
L L V S V L I S I F L M F V F S R R T S R L L T S L S F V T I I
4101 GCTACTGGTTCTGTTTGAATTTCTTATGTTTCTTATGTTTCTAGAACTTCTAGATTGCTCACCTTTAAGTTTGTACTTAATATA
A F V Y G G G L D Y I K A Y L N I D A G T F D L G L G F R Y W E L K
4201 GCATTCGTTTATGGAGGTGATTTGATTTAATTAAGCATATTTAAATATAGATGCTGGAACATTTGATCTGGGATTTAGATTTAGATTTGGGAGTTGA
Y Y L G L F C T L G V G I L T S S N I N S F I L A G P S
4301 AATACATCTGGGTTTATGGCTATTGATTTCTGTACCCCTGGTGGTAAATCAACATCAAGTAATATAAATAGTAATTTTATTTGGCTGGGCCAAG
A V K M Y L D D L G F L E L F V Q F G V Y A A I F M Y I F Y K L
4401 TGCTGTTAAGATGATCTAGATGATTTAGGTTTTTATAGATTTATGTTCAATTTGGTGTAGCCGCAATTTTATGATGTTTATATTTTATAAATTA
I N L I L R M S N D K Y R V D R A F P I A L L T N L I I T S I S L N
4501 ATCAATTTAATTCGAAGATGCAAAATGATAAATATAGATGAGTGGACCCGCTTTTTTATGCACTTAACTAAATTTAATAATTTCAATTTCTTTAA
I F G A Q R S F S L A I V L A L I F Y Y D Y R L K N D V E N \*
4601 ATATCTTTGGAGCCGAGAGAAGTTTCTCATTGGCAATGTTCTTGCATTAATATCTACTATGACTATAGGCTGAAAAATGACGTAGAAAAATAGGTGTA

*cps19bQ*→

M D K V C I V I L N Y N N Y E E T I E R V Q S L R S T I K S N E Y
4701 TAATGGATAAAGTATGATAGTATTCTCAATTAATAAATTTGAAGAAACGATTTAGCGGTGACAAAGTTTGAAGTACTATAAATCAAAATGAGTA
D I V I V D N N S V K E L S P I K I T S L E N R
4801 TGACATCGTATTTGATAGATAAATAATTCGGTGAATGATAGTCTCAAAGAGTTATCCAAGCGGTTATCCCTATTAAGATTTATCTAGTTTGAAGAAATAGA
G Y A N G N I G I K Y A E D N G Y D Y I C I L N N D T L I E V D F
4901 GGATATCCGAATGGAACAAATTTGGAATAAAGTATGAGTGGACCAAGTATGATTTGATATTAATAAACAAGTACATTAATTTGAGTTGATTTG
L E S C K R E L E N N S F V A F V S P V L V E Y K D N N L V Q S T
5001 TTTTGAATCGTGTAAACGAGAAGTAAATAAATTTCTTCTGTTGTTAGTCCAGTGTAGTTGAATATAAAGATAAATACTGGTACAATCTAC
G G D I F I N R G I V T L K N H G A Q R D K L P S K I E S D Y I G
5101 AGTGTGATATTTTATTAATAGGGAAATTTGAATTTTAAATAATCATGTTGCTCAGAGAGACAACTTCTTCTAAATCGAAAGTACATTAATTTGGG
G A C L M F K I I G Y I P E S Y F L F Y E E T E W C Y R A
5201 GGAGCATGTTGATGTTCAAAACCTCTATCTGAAAAATTTGGATATATACCTGAAAGTACTTTCTATTTTATGAGAAACTGAATGGTGTATAGGG
K K L G Y K N I C L T Q S Y V Y H K G S V S I K A V N G L Q E Y L
5301 CTAATAAATTTAGGCTATAAAGATATGCTTACTCAAAGTATGTTGATCATAAAGGTTTGGCTCTCTATAAAAGCGGCAATGACTTCAAGAAATTTT
M A R N R V V F V R R N I N S K L K Y S A F L F Y L F M Q Q L Y H
5401 AATGGCAAGAAATAGAGTTGATTTGTTCTGAGAAAATATAAATAGCAAACTAAAGTATTTCTGCTTTTTTGTCTTATTTATTCACCAACCTTATCAT
C F L R R D C S K R K Y L L D G V F N R I D P S Y P F I F I S E
5501 TGCTTCTTGGCAAGGGATTTGTTCTAAAGAAAGTATAAATATTTATAGTGGTATTTAATAGAATGATCCATCCATCCATTTATTTTCATAAGCG
\*
5601 AATAAGTTACTACTTATAACTGTAGATGCACTAAGATAAATATAGTAGACTGAATCTAAAAATAGTACGAAATAATGCTAAAAACATTTATAGAAATTAAT
5701 TTTACTTCCATAATCGATTTGTTCTTATCTTATTTCAATCTGCTATAGATTAATAAAGTGAAGAAAGTGGTCATTTATAGGAACATAAATTTACATAGTTA

*cps19bR*→

M K M F D I S Q I K T K T V A F D F F D T V V H R N C H P E Q
5801 GGTGAAACATGAAAATGTTGATATATCTCAAATAAAAAACAAAAACAGTTGCATTTGATTTTTTGTACTGTTGTGCATAGGAATGTCATCTCTGAGCA
I L Y Q W A K E M A L E V N F N I S P S I L Y Q I R K S V E N N K
5901 AATTTTGTATCAGTGGCTAAGGAAATGGCTTTGGAGGTGAATTTCAATATATCTCCCTCAATATATATCAAAATGAAAAAGGATTCGAAAAACAATAAA
K L G T E E M C Y L D L L S G I Y N E I K D K I K N T S K E E F I H
6001 AAGTTAGGCACTGAAGAAATGTTTATCTAGACCTTTTGTCTGGAATATATAATGAAATTAAGATAAAGATAAATAATACATCGAAAGAGGATTTATTC
R A K I E L E Q H I Y L D S E I K E V L R K K L K S D S K Q
6101 ATAGAGCTAAAATCTAGATGAAAATGAAATACAACATATTTATTTGGATTCGGAATTAAGAAAGTTTAAAAAATTTGAAAAGTGAATTTCAAAAACA
I I L V S D F Y T D K E L I E T V L K K F E I F D Y F S S I Y I S
6201 GATTATCTTAGTTTCTGATTTTATACTGATAAAGAAATTAATAGAACTGATTTAAAAAAGTTTGAATTTTTGATTTTCTCTCTATCTATATTTTCG
S E K G C R K S T G N L Y K L I L K E L G L N P I E I T M I G D N C
6301 AGTGAGAAAGTTTGGCTAAATCCACAGGAAATTTATATAAGTTAATTTTAAAGAAATTAGGCTTAAATCCATAGAAATTCACATGATAGGAGATAAAT
K S D Y E V P R S L G L N A I Y R R Y I D K N S T V S E K E L V R
6401 GTAATCTGATTATGAAGTCCCGCTTCTCTAGGATTAATGCTATTTATAGACGATATATAGATAAATAATCAACAGTCTCAGAGAAAGGATTTGTAAG
L Y D Q I L F S N S K K A P N I F L A D I V F P I S E L V F H K K M
6501 ACTCTATGATCAAAATTTTCTAACAGTAAAAAGCACCCTTTAATATTTTATTTAGCTGACATCGTATTTTATTTTGGAAATTTGCATAAAAAAGAT
I Q D D V I Q I A L F C S R E G Q L L K R L F D I Y Q D T F L R E N
6601 ATTCAGGATGATGTTCAATAGCACTCTTTTTTGTCTCAAGAGAAGGACGTTATTGAAAAGACTTTTTGATATATATCAAGATACTTTTTTGAGAGAAA
Q K I C T E Y F V S R R S T L Y S S F T S L E N E E F E M I F R
6701 ATCAAAAGATTCTACAGAATTTTATGTTCTAGACGATCGCAATTTATTTCTTCTTTACTTTTATAGAAATGAAAGATTTGAGATGATCTTTTCG
Q Y K K I T L Q N F L L N L N F S N N E I T L I C Q D L N V K P T
6801 TCAATATAAAAAATTCGTTACAAAATTTCTTATTAACCTTGAATTTTCTAATAATGAAATCAATGATTTGTCAGATTTAATGTGAAACCTACA

Y V L T V D D H L L E N L R K H P Q F I K K F N Q E K K D S Q L L R  
6901 TATGATTAACAGTAGACGATCATTTATAGAGAACCTAAGAAAACCCCTCAGTTTATTAAGAAGTTCAATCAAGAGAAAAAGATGCCAATACTAC  
E Y I K H L T K N R N E A Y L V D V G W K G T I Q D K K A L A  
7001 GTGAGTATATTAAGCATTAAACAAAAACCGAAATGAGGCATATTTAGTTGATGTAGGTTGGAAGGAACGATACAGATAGTATCAAAAAGGCTCTTCC  
D K R I V G Y Y L G L M L N V Y S V E N K T D K T G L L F S D Y P  
7101 AGACAAGAGAAATAGTAGATTTCTAGGATTTGATGCTCAATGTTTATTCAGTAGAAAAATAAACCGGATAAAAAGTGGTTTATTTTCTGATATCCA  
S K S K F Y D I V S R N F G F Y E D I F V A D H G P V L K Y K K E S  
7201 AGTAAATCGAAATTTATGATATTTGAGTAGAAATTTGGTTTATGAAGATATTTTGTAGCAGATCATGGTCCAGTTTGAAGATATAAAAAAGAGA  
D I I P I I D D D K K H V S I Y Q A V K D Y Q E E L V L G F S E I  
7301 GTGATATTTCCCATTTATGACGATGATAAGAAACATGTGAGTATTTATCAGGCAGTTAAAGATTATCAAGAAGAGTTAGTATTAGGATTTTCAGAAAT  
L E A Y K K M K F L P F E Q K N L W L T M S L K K E C I Y I P K L  
7401 TTTGGAAGCTTATAAGAAATGAAGTTTCTCCCTTTGAAACAAAAAATCTGTGGTAAACGATGTGCGTAAAAAAGAAATGATTTATATACCTPAAGTTA  
Q S F S E S L K E K V E N F G E I V T L K T T K K S I K T L R E  
7501 CAATCATTTCTGAATCTTAAAAAGAAAAAGTTGTTGAAAAATTTGGTGTAGATAGTAACTCTTAAGACTACGAAAAAATCTATAAAAAAGTTATTAAGAG  
K S D L L W V D F V Y R L F G G V N F L F I P E L Y T R V I F L L  
7601 AGAAAAGCGATTTACTTTGGGTTGATTTGTTTATAGATTTGTTGGTGTCAATTTCTATTTATCCAGAAATATATACAAGAGTATATTTTTTATT  
K Y L D L K L R L K N Y G E \*

*cps19bJ*→

M G N K S I K L N A L L N I V L T L S N I I  
7701 GAAATATTTAGATTGAAATTTGAGGTTGAAAAATTTGGGAAATAATCCATAAAGTTGAATGCATTATTAATATTGCTCAGCCTATCAAAATATCAT  
F P L I T F P Y I S R I L N P N G I G L T S F F S S I G N Y G I L  
7801 TTTCCCATTAATCACTTTTCTTATATCTAGAATATGAATCCAAATGGTATAGGTTTAACTTCATTTTGTAGTCAATAGGGAATATGGTATTTTA  
L A S L G I S T Y G I K A V A S V E R D D R D K L S K V V I Q E L M I I  
7901 CTTGCTCTCGGAAATTTCAACTTATGGTATCAAGCAGTAGCAAGTGTAGAGATGATAGAGATAAGTTGTCAAAAAGTAGTACAGGAGTTAATGATTA  
N V A M S I I T T A I L L F M T I F I F I T Q L N R E F S L F L I T C  
8001 TAAACGTTGCTATGCTATAATAACAACGCAATACTATTTATGACTATATTTATAACACAATGAATAGAGAATTTTCACTCCTATTGATCACATG  
G T I L A M Q I L M P I L L I S G F S N I T G N Q I L I P M N R E  
8101 TGGGACTATTTTACTTCTCCTTTCCCTTAAATGGTTGATAGTGGAAATGAAGAATAATACGTATATTACTACTAGGTCAGTAGTGTAAAAATTTCTA  
S L I L I F L L V K R P E D Y I V F A S I S L F S S L S S N I L N L  
8201 TCATTAATTTGATTTTCTACTTGTGAAAAGGCCAGAGGATTTATTTGTTTGTCTAGTATTTTCTTCTCTAAGTTCAAAATCTTAAATC  
W H S R H F I N I K L Y K N L Q F K Y H F K P M W Y L F A S L L A  
8301 TATGGCATGCCGACATTTTCAATATAATAATATAAAAAATTTACAATTTAAATATCATTTTAAACCAATGTGGTATTTATTTGCTCATCTGTC  
N I Y T V M L D T V M L F G I N G N E A V G Y Y S V A S K V K W I  
8401 AGTAAATATTTATACTAATTTAGATACAGTGTGCTCGGTTTATTAATGGTAAATGAGGCTGTGGATACTATTCTGTGGCATCAAGGTTAAGTGGATT  
L L S L I T S I S A V L L L R L S F Y I S K N D T S N F I K M L K E  
8501 TTTGCTCTTATACATCTATTGTCAGTTTGTCTACTGAGACTTTTCAATTTTATAGTAAAAATGACACCTCGAATTTTATAAAAAATGTTAAAGG  
S S A V I F F I A I P L M V F F I V E A K D S I L L L G G S Q Y L  
8601 AGTCACTGCGGTTATATTTTATGCGATTCCATTGATGGTATCTTTATTTGTAGAGGCGAAAGATAGTATCTTATTACTAGGAGGAGTCAAGTATCT  
P A T L A M Q I L M P I L L I S G F S N I T G N Q I L I P M N R E  
7701 TCCYGGACTTTAGCGATGCAAACTACTTATGCCAATTTTACTTATTTCTGGTTTCTCGAATATTACAGGAAATCAAAATTTGATTCCAATGAATAGAGAA  
Y F M V I A V T I G A V I N L L L M P K F G I G A S V A T  
8801 AAATATTTTATGTTGCGAGTAAACGATGGTGTGTGATTAATCTTATTTTGAATCTACTGTTAATGCCTAAGTTTGAATATTATGGTGTCTCTGTGCGAA  
L F A E L S Q M T V Q L H F S K E Y L V S N I S I K S L V N V I I  
8901 CTCCTTTTGGGAAATTTGCGCAGATGACGGTACAAATACATTTTCAAAAGAAATATTAGTATCAAAATATATCGATAAAGAGTTTGGTAAATGTTGATAAT  
A T V V S I I P L I I L N Q L I T I T I P F Y S L M L A G F A F F  
9001 TGCAACAGTTGTTTCTATAATACCCTAATCTATTTGAATCAGTGTATAACGATACTATACCATTTTATCTCTAATGCTAGCAGGTTTGTCTTCTTT  
S L Y L V I L L L K E E V T I Q L F S L A K K K \*  
9101 TCATATATTTAGTAATTTCTGCTTTTATTAAGGAGGAAGTACGATTCATTTTCTCTTCTTGCAGAAAGAAATGAAATGGTTCAGAAATGAAAT

*cps19bK*→

M K I M L V F G T R P E A I K M C P L  
9201 GTATAACAATAAATAATTTAATTTAATTTAGAGGAGAAATCATGAAGATAATGCTAGTTTTGGTACACGCTCCAGAAGCGATAAAAAATGTGCCATT  
V N E L K K Q A D M E T V V C V T G Q H K E M V S P V L E L F G V  
9301 AGTGAATGAGTTGAAAAACAGGCAGATATGAAAAACAGTTTGTGTAACGTTGTTGTAACAGGATGTTAGCCCTGTTTGGAAATTTGTTGGAGTT  
Q P D Y D L E I M K A N Q T L F S I T T S I L E K I K P V L E E E Q  
9401 CAACAGATATGATTTAGAAATTTAGAAAGCTAACTCAACCTTTGTTTCTATAACAACAAGTATTTAGAAAAATTTAAACGTTTGTAGAAAGAAC  
P D I V L V H G D T T T Y A A A L A A F Y L G I K V G H V E A G  
9501 AACCAGATATTTGCTAGTTACGGTACACTACTACGACATATGACAGCAGTTTGGCAGCATTATCTAGGAATTAAGGTTGGTCAATGTTGAAGCTGG  
L R T Y N L Q S P F P E E F N R Q S T S I I A N Y H F A P T E L A  
9601 GTTACGAACTTACAACCTGCAAGTCCCTTCCAGAAGAAATTAATAGACAATCGACTTCTATTATTGCAAAATTTATCAATTTGCCCCACAGAATTAGT  
K E N L T K E G R N N V Y V T G N T V I D A L K T T V Q K K D Y T H P  
9701 AAAGAAAATCTAACAAAAGAGGTTAGAAACAATGTTTATGTGACAGGTAATACGGTAAATTTGATGACACTTAAACACTACAGTACAAAAGGATTTATACACCC  
D L D L N A D N R L I L L T A H R R E N L G E P M K H M F R A V K  
9801 CTGATTTAGATTTAAACGCTGATAATCGTCTCATTTACTGCTCATAGACGCTGAAATCTCGGAGAACCTATGAACACATGTTTAGAGCTGTTAA  
R I L N E Y D D V K V I Y P I H K N P L V R E T A A E I F G D I E  
9901 ACGAATTTTAAATGAATATGACGATGTTAAGGTAATTTTCAATTCATTAAGAAATCTTTGGTTCGTGAAACAGCTGCGGAAATTTTGGAGACATAGAA  
R I Q I E P L D V L D F H N F M N N S Y M I L T D S G V Q E E A  
10001 CGAATTCAGATTTGAACTTTAGATGTTCTGGATTTTCATAACTTTATGAATAATAGTTACATGATTTTAACTGACGAGGAGGTTGAGGAAAG  
P S L G K P V L V M R D T T E R P E G V A A G T L K L V G T D E E  
10101 CGCCTTCTGTTAGGAAAGCCTGATTTGGTATGCGAGATACGACAGAAAGACCTGAAAGAGTAGCTGCTGGAACGTTGAAATTTGGTGGGACTGATGAGGA  
T I Y Q N F K M L L D D S E E Y K K M S Q A S N P Y G N G D A S K  
10201 GACTATTTATCAAAATTTAAGATGCTTTTAGACGATTTCCGAGAAATATAAAAAATGAGTCAAGCTAGTAACTTATGGAATGGTATGCTAGTAAA

*cps19bL*→

Q I V R I L R G I \* M K G I I L A G G S G  
10301 CAGATTTGTCGAATTTTACGTGGAATTTGAGTGTGTTAGATAAAGTAAACAGAAAAGGTACCCTACTATGAAAGGATTTATTTCTAGCAGGTTGGTTCCGG  
T R L Y P L T R A A S K Q L M P V Y D K P M I Y Y P L S T L L A  
10401 GACTCGCTTGATCTCTTTGACTCGCGCTGCATCAAAACAATGATGCGGTTTATGATAAACCCATGATTTACTACCCACTTTCAACATTTGATGTTGGCT  
G I R D I L I I S T P Q D L H R  
10501 GGAATTAGGATATTTTGATTTTCCACTCCACAGGATTTACATCGAT

# APPENDIX III

## The nucleotide and deduced amino acid sequence of the complete *cps19a* locus.

The nucleotide sequence is numbered in accordance with GenBank accession number AF094575 and is shown from the 3' end of *dexB* to the 5' region of *aliA*. The amino acid translation for each ORF is represented by single letter code above the first nucleotide of each codon. The putative -35 and -10 promoter sequences are double underlined and possible ribosome binding sites are singly underlined. The sequence which corresponds to the partial copy of *IS1202* is shown in bold and the inverted repeat is underlined with a dashed line. The three arrows in the 5' intergenic region indicate the point of three separate deletions when compared to the *cps19f* sequence.

```

→dexB
V E *
1  GTGGAATGACTGACTAGAAATGAGCAAACCTCAAGTTTTTGAAGCTTGAGGTTTTTACTATAGTGGATTGAAACTAGAAATAGTGCACCTCTGCTTCTAAA
    ↑                                     ↑
101 ACATTGTTAGAAATCGATTGACTGTCCCGATCGAATTGTCCTATTCTTATTTTCATTTTGCCTATACTTGTAGTTGAGGGAATTAAGCTCCTCAGTAGTA
    ↑
201  AACTCTTCATAATCCTTTTTATTTTTATGAAGATATGTTTTGAAAGATGTGAGTTCCACGGATGGGTTTGTGGAGGGATATACTTGCCTCTTCTTTTT
301  TTGTTCTGGTCTTCTGTTCAAAGTTTTTCGAAATAGAGTTTCATGATTTAGTAGCTCCTTTGTTGATAGATTTTGTGAGGATATTGAGGTAGATGTCCTCCG
401  TCAAATGCTTTTATAACTAATGCTTTCGTCTTCTGATGAAATAGACTTCTTTTCCCTTGTTCGGTAGGGATATAGTAAACGATTTTGGAAATCGGATATGGT
501  GTCCACTATCGACGACTCTCCGCCAGTCTAGCCAGAATGAGATTCCTTTTCAGAGGGGTTAGGAACCTCCTCGAAGACAGAGATTTTGTCTTGTTC
601  AAACGTTCATTAAGGTTTGAATGTAGGAAAGCAGGAAGGTATTGGCTTCTTCCAAGGTATGAATATTGTTTTGTTCCAGTTCGATAGCCAGGCAGAGAT
701  TGTAGTGTCTGATTGAGTCTTTCGACTCTCCCTTTAGCTTGAAGGATAGAGGTTGCTCCTCAAGGAGAAATCCCTAGTTGGTGAACAGGCATCCAAACTGTG
801  TATGGGTATCGTCCCTCCATTTTCTTAGAGTTGGAGGCTTGATAGGTAAGACCGTCTCTCATCTGTTTTAATTTGAAGGGGAATGCCGTGATGGCTAA
901  GATTTGTTTCGAGGACATGATAGTAAGCATTCAAGTCTCTTGTTTTCAAATAAAGCCCTAGGATATTGCCAGAAGCATCATCAATGGCTAAGTGTAAAG
1001 TTAGATGTTTGGGCTCCAAACTAGGCATGAGGGCTGGCATCCATTTGAAATGAGTTCACCAGCAAATTTCTTCTGGGCTACTAGGATGACTTTTGTAG
1101 GGTCTTCAAAGGAAGTTTCAGCCGTCGGTAAGATTGGATTGCTTAGGGGTTGATTCAGGTTTCAGTTAGCTTGTCTTAACTCTCTCTTTGCTTTCT
1201 ATGGGACTTAGCCGACAGGATATTTTCTTATAGATATTTTTCTAACAGTAGTATCAGAGAGCTGAATTCCTTCTTTCAGCTAGCAATTCACAGAA
1301 TGAAGGACATTTGGTTTATATGTTTCATAGGAGAGGTATTTTTTATGGAGACGTTCTTTGATTTTCATCAGGATTCGATGTTTGGTTTTTCGATTTCTGT
1401 TTCCTGTCTGAAGGCTCTTTTCTTTCAGTTCATAGGCTAGTAGCAGACGATTTGATTTGCTTTTCAGAAAGATTGAGCTCGACACAGGCTCGTTTCTT
1501 TGTTTTCTTTCTTGGCTTATAGCTTTTATCACAAGATATTTTTCTGTTTCATTCATATTTCAGTTGGATCCTTTTCATATGACTATTCACCAAAATGGGA
1601 CATTTTCAGTTCGATTTACTAAAGACATTTACATTCGAATTCACAAGAATGCAGATAGTGAAAAAAGGTGTAGACATTACCGTAAAAAAGTGATATA

                                cps19aA →
                                M S R R F F K K S R S Q K V K R S V N I V L L T
1701 ATCGTATGATGTTCAAGGTATAGGTTTAAATCATGAGTAGACGTTTTAAAAAATCACGTTCCACAGAAAGTGAAGCGAAGTGTAAATATCGTTTTGCTGAC
    I Y L L L V C F L L F L I F T Y N I L A F R Y L N L V V P A L V L
1801 TATTTATTTATTTAGTTTGTTTTTTATTTGTTCTTAATCTTTACGTACAATATCCTTGGCTTTTAGATATCTTAACCTAGTGGTACCTGCGTTAGTCCATA
    L V A L V G L L L I I Y K K A E K F T I F L L V F S I L V S S V S L
1901 CTAGTTGCCCTGGTAGGGCTACTCTTGTATTATCTATAAAAAAGCTGAAAGTNTTACTATTTTCTGTTGGTGTCTCTATCTCTGACGCTCTGTGTCCG
    F A V Q Q F V G L T N R L N A T S N Y S E Y S I S V A V L A D S E
2001 TCTTTGAGTACAGCAGTTTGTGGGACTGACCAATCGTTTAAATGCGACTTCTAATTAAGTCTCAGAAATATTCGATCAGTGTGCTGTTTTAGCAGATAGTAGA

```



cps19aF →

Y I D G W T I W K D I E I L L K T V K V V L M K D G A K \* M N E
6601 TATATTGATGGTGGACAATCGGAAAGATATTGAAATTTTATGAAAGACAGTTAAAGTTGATGATGAAGGATGGAGCGAAATAGTTCATATGAATGA
R I Q I L G I T I D P L T M K E T V D A V E Q Y V L K H P L H L

cps19aG →

K V L R A K R G Y K L \* L I Q K V E L D A I K E F Q K
7401 AAAGTTTTAAGAGCGAAAAGGGGATATAAACTTTGAACCCACACAGAAATAGTTAAATTTGATTCAGAAAAGTTGAAATGGATGCATATAAAAGAGTTTCAAAA
I C K E N N I D F F L R G G S V L G A V K Y D G F I P W D D D M D

cps19aH →

Y L H E E P S D D E K K S H L G G Q \* L F C Y I I L H Y K V L E E
8201 TATCTTCACGAAGCCTAGTGTATGATAAAGAAATCGCATTTAGGAGGACAATAATTTGTTTTGTTATATATTATTCATACAAAGTCTTAGAGGAA
T I S C V K S I K E G N S N T K Q I V I I D N F S N N G T G E K L Q

cps19aI →

F L K L M K E N E G V \* M T Y L F L L C L T L F L L T F F Y F F A
9101 TTTTTGAAGTTGATGAAAGAAAACGAGGGTGTGTAATGACTTATTTATTTTACTCGCTGACCTTATTCCTTATTAACCTTCTTATTTTGGCT
F N Q D L I A P P V V M S V M F L L I S S V F A L I N V Q N W N I E Y

cps19aJ →

N S K I K N I L T N F S Y V I S S N L L T V L T S S L V V L I F P
10501 TGAATAGCAAAATAAAAATATACTAACTAATTTCTCTTATGTCATTTCTTCAAACTTTTGGACAGCTTGGACCTCTTCTGGTGTCTTAAATTTCC
K L M G V T E Y S Y T L Y I F Y L T Y I G F F H L G D G I Y

10901 V Y V L Q M T N R L K D S S I I L L S D R V L Y V L L L F L F I V
CGTGTATGTTGACAGTACAAACCGCTTGAAGGATGTTCCATCATTTCTACGTGATCGCGT...
11001 F G W H E Y K V M I L A D I L G R S F S L I L S F W I C K D I V F Q
TTTGGTTGGCATGAATACAAGGTTATGATTTTGGCAGATATTCCTAGTGCATATTTCTCTCATACTTCTCTCGATTGTAAGAATATGTTGTTTC
11101 P L S K F I F N I K E S F D N I R V G I N I L M L S N I A S S L I I
AGCCTTTGTCAAAATTCATCTTTAATATAAAAAGAGCTTTTGATAAATCCGAGTTGGTATCAATTTAATGTATCCAAACATCGCAAGTAGTTGATTAAT
11201 G I V R L G I Q W N W N I E T F G K V S L T L S V S N L M L T F I
CGGTATTTCTGCTAGGAATCCAATGGAACCTGGAATATCGAAACATTCGGGAAAGTATCTCTGACCTTGAGTGGTTCCAACTATGATGATTTCAIT
11301 N A I G L V V F P L L R R T K A E N L P K I Y S N L R N V L M L I M
AATGCCATTGGACTAGTCGTCTTTCCGCTATTAAGACGACAAAAGCAGAAAACCTACCTAAGATTTATCTTAATTTAAGAAACGCTTTTGTATGCTTATCA
11401 F A I L L I Y Y P L K I V L D L W L P A Y Q D A L I F M T L I F P
TGTTCCGATTTTGTCTATTACTATCTTTAAAAATGTATTAGACCTCTGGTTGCCAGCCTATCAAGATGCCTGATTTTCATGACCCCTTATTTCC
11501 M S V Y E G K M A L V I N T Y L K A L R M E R D I L R V N A L V M
TATGTCAGTCTATGAAGGAAAATGGCATTTGGTCAATAACTTACTTAAAGGCATTAAGAATGGAAAAGAGATACCTTAGAGTCAATGACCTTAGTAATG
11601 L I S M G V T L V T T Y L L N S L E L T V V S I V V L L A L R S I I
TTGATCAGTATGGGATCTTACAACATCTGTTAAATAGTTTGGAGCTGACTGTTGTATCGATAGTTGTTTGTAGCTTTGAGAAGTATCA
11701 A E L I L S K K L D V S V K K D I V L E F L L T L V F I S S S W Y
TAGCCGAATTAATCTGCTAAAACCTGATTTTCAGTAAAGAGGATATGTAATTTAGAAATTTCTTTGACGCTGTCTTTATTTCTCAAGTTGGTA
11801 L P I G L A V I V Y T I A Y G L Y L K H E D I K Y L A Y F H K
TTTGCCGATTGGCTGGCAGTATTTGCTACACAATAGCCTACGCTTTATATCTCTATTGAAACACGAAGATATCAAAACCTATTTAGCTTACTTTAAA

cps19aK ->

11901 A S K K K T S N \* M K K I
GCTAGTAAAAACATCAAAATAAAAATATATAATGATTAAGTGGTAGATTTCTATTTCTACCCTTTTAGATATTCGGGAGGTAATGATGAAAAAAT
12001 M L V F G T R P E A I K M C P L V N E L K K H E D M E T I V C V T G
ATGTTAGTATTCGGTACACGCTCGGGAAGCCATCAAAATGTGTCATTTAGTCAATGAGTTGAAAAACACGAAGATATGGAACAATTTGTTGTTGTTACTG
12101 Q H K E M V S P V L D L F G V V P D Y D L E I M K A N Q D S G E Y K K
GACAACACAAAGAGATGTTAGTCTGTTTATGATTTATTTGGTGTGTACAGATTTATGATTTAGAAATTTAAGAGCTAACCAACCTTGTCTCTAT
12201 T T S I L E K I K P V L A K E Q P D I V L V H G D T T T Y A A A
CACAACTAGTACTTTGGAAAAGATAAAAACAGTTLTLAGAGAAGAAACACAGATATTTGCTAGTTCCAGGTACAGCAACTTACGAGCAGCC
12301 L A A F Y L G I K V G H V E A G L R T Y N L Q S P F P E E F N R Q S
TTGGCAGCATTCTATTTGGGAATTAAGTAGGACATGTTGAAGCTGTTGCGAACGTFACAATTTACAAGTCCATTTCTGAAGYATTAACAGGCAAT
12401 T S I A T Y H F A P T E L A K E N L L K E G R E N V Y V T G N T
CGACATCAATTCGAACCTTACCAATTTGCTCCAAGTGGCTAAAGAAAATCTCTTAAAAGAAAGGTAGAGAGAATGTTTATGTTGGAATAAC
12501 V I D A L T T V Q E D Y T H T H L D L N A N N R L I L L T A H R
TGTCATTGATGCTCTTACAACTACTGTTCAAGAGGATATACACACTCATTAGATTTAAACGCTAACAACTCGTCTATCTTATTGACTGCTCATAGA
12601 R E N L G E P M R H M F R A V K R V L N E Y E D V K V I Y P I H K N
CGCGAAAATCTCGGCAACCGATGAGACATATGTTTAGAGCAGTTAAACGAGTATTGAATGAATACGAAGATGTTAAAAGTCAATTTCAATTCATAAGA
12701 P L V R E T A A E I F G D T E R I Q I I E P L D V L D F H N F M N
ATCCTTTGTTACGTGAACAGCTGACAGATCTTTGGAGATACAGAAGGATTCAAAATTTAAGACCTTTGGATGTTCTGATTTCCACAACCTCATGAA
12801 Q S Y M I L T D S G G T Q E E A P S S L G K P V L V M R D T T E R P
CCAAAGTTATATGATTTTAAACAGATTTGCGGAGTTCCAGGAAGAAGCCTCTTCTTAGGAAAACCTGTATTGGTTATGCGCGATACTACTGAAAGACCT
12901 E G V A A G T L K L V Q D E E T I Y Q N F K L L L D D S G E Y K K
GAAGGATGAGCGCAGAACCTTGAATTTAGTAGGAATGACGAAGAACAATTTATCAAACTTCAAACTACTCTTAGATGATTTGGAAGAATACAAA
13001 M S Q A S N P Y G N G D A S Q I V Q I L R G I \*
AAATGATCAGCTAGCAATCCATACGCTAATGGTGTAGTGTCAACAGATTTGTTAGATTTTGGCTGGGATTTAAGAAGATTTCTTAAAAGAGTCTAA
13101 GAGGAATCCACTCCACAAACATAAAAACCTTATGCTATCATTTATGATTTAGATAGACTCATAACTTAATATTTTTATCCTAAGAAAATTCGTTTATAT

cps19aL ->

13201 M K G I I L A G G S G T R L Y P L T R A A
CAATATGAATGAAAAACAACGAAGAAGGTTATTTTCATTAAGAAGTATTTCTTGCAGGCGGCTCAGGTACCCGCTGTACCCACTTACTCGGGCTG
13301 S K Q L M P V Y D K P M I Y Y P L S T L M L A G I K D I L I I S T
CGTCAAAACAGCTGATCCGGTTTATGATAAACCATATGATTTATTTCCGTTGTCGACATTAATGTTGGCTGGAATTAAGATATTTTGATTTCTCAAC
13401 P Q D L P R F K D L D L D G S E F G I K L S Y A E Q P S P D G L A
TCCTCAAGATTTGCCCCGTTTAAAGGATTTGCTCTGGATGTTCCGAATTTGGATCAAGCTTTCCATGCGGAACAACCTAGTCCCGATGGCATTTGCT
13501 Q A F L I G E F I G D S V A L I T L G D N I Y H G P L S T M L Q
CAGGCTTTCTTATCGGTGAAGAATTTATCGGTGACGATAGCGTTGCTTGTATTTAGGCGACAATATCTACCATGGTCTGGTCTGAGCACAATGCTTC
13601 K A A K K E K G A T V F G Y Q V K D P E R F G V E F D T D M N A
AAAAAGCAGCCAAAGAAAGAGAAAGGTCGACTGTTTGGCTACCAAGTGAAGGATCCAGAGCCTTTTGGTGTGGTTGATATACAGACATGAATGC
13701 I S I E E K P E Y P R S N Y A V T G L Y F Y D N D V V E I A K Q I
TATTTCTATAGAAGAAAACCGGAGTATCCCTCCAACCTGACAGTACAGGACTGATTTCTATGATAATGATGTTGGAGATTTGCTAAACAGATC
13801 K P S A R G E L E I T D V N K A Y L N R G D L S V E L M G R G F A W
AAACCTAGTCTCGTGGCAGTTAGAAAATACAGACCTTAAACAGGCTTACCTAAATCTGGTGAACCTTTCTGTTGAGCTGATGGGGCGTGGTTTGGCT
13901 L D T G T H E S L L E A S Q Y I E T V Q R M Q N V L M L E E I
GGTTGGATACGGAAACCATGAAAGCTTGTAGAAAGCTTCTCAGTATATTGAAACGGTTACGCGTATGCAGAAGCTTCAAGTGGCAATAGAAAGAAAT
14001 S Y R M G Y I S R E D V L E L A Q L L K K N E Y G Q Y L L R L I G
TTCTATCGTATGGGTACATTTAGTCTGTAAGCTGCTGGAAATTTGGCTCAGCTTCTTAAAGAAAATGAATACGGCAATATTTGCTCCGTTTATGGA

cps19aM ->

14101 E A \* M S D N F F G K T L V V R K I D A I P G L L E F D I P V H G D
GAAGCATAGATGCAGATAATTTTGGAAAGACACTTGTGTACCGCAAGATTTGATGCTATACAGGACTGCTAGAGTTTGGATTTCCCGTTCCATGGAG
14201 N R G W F K E N F Q K E K M L P L G F P E S K F F A A G K L Q N N V
ACAATCGTGGTTGGTTTAAAGAAAATTTCCAGAAGAAAAGATGCTACCGCTTGGTTTTCCTGAAAGCTTCTTTGCTGACGGAAAACCTGCAAAATAACGT
14301 S F S R K N V L R G L H A E P W D K A Y I S V A D D G K V L G S W V
CAGCTTTTCTCGCAAAAATGTTCTTCGAGGACTCCATGCTGAACCTGGGCAAGTATATCTCTGTTGACAGCAGTGGGAAGTTTATGATTTGGGTA
14401 D L R E G E T F G N V Y Q T E I D A S K G I F V P R G V A N G F Q V
GATTTGCGTGAAGCGAGCCTTTGGAATGTTTACCAGACAGATTTGATGCAAGCAAGGAAATCTTTGTTCTCGAGGCTGAGCTAATGCTTCCAGT
14501 L V S D T V S Y S Y L V N D Y W A L E L K P K Y A F V N Y A D P S L
TTCTATCAGATACAGTGTATATAGCTATCTGGTCAATGATTTACTGGGCGCTTGAACCTCAAAACCAAGTATGCCTTTGTAAGACTACCGCTGATCCAGCCT
14601 G I E W E N I A E A E V S E A D K N H P L L K D V K P L K K E D L
TGGTATTGAATGGGAAAATATTGCAGAAGCAGAGGTTTTCAGAAGCAGATAAAAATCATCCACTACTTAAGGATGTAAAAACCTTTGAAAAAGAAGATTTG

cps19aN ->

14701 E \* M T E Y K N I I V T G G A G G F I G S N F V H Y V E N F P
GAATAAGCAAAGAATGACTGAATACAAAATATATACGTGACGGTGGAGCTGGCTTTATCGTTCTAACCTTTGCTCATATGTTTACGAGAACTTTC
14801 D V H V T V L D K L T V L A G N R A N I E I L G N R V E L V G D
CAGATGTTCAATGTCAGTCTAGATAAGTTGACTTATGCTGGAACCCGCGCAATATTGAGGAAAATTTAGGTAATCGGTTGAGTTAGTTGTTGGTGA
14901 I A D A E L V D K L A Q A D A I V H Y A A E S H N D P
CATTGCTGATGCGGAGTTGGTAGACAAGTTGGCTGCTCAAGCAGATGCTATCGTTCAATATGACCGGAAAGCCCAATGATAATTCGCTCAATGATCCA
15001 S P F I H T N F I G T Y T L L E A A R K Y D I R F H V S T D E V Y
TCGCCATTTATCATACTAATCTTGAACCTATACCTTTTGAAGCTCTGTAAGTACCTGCTGATGATATTCGCTTCCACCATGATGACAGATGAAAGTTT
15101 G D L P L R E D L P G H G E G P A E K F T A E T K Y N P S S P Y S
ATGGGGATCTCCCTTTACCGCAAGATTTGCCAGGTCATGGAGAAGGGCCGGTGCAGAAAATTTACGGCTGAAACCAAGTACAATCCAAAGCTGCGCTTACTC
15201 S T K A A S D L I V K A W R S F G V K A T I S N C S N N Y G P Y
ATCAACCAAGGCGCTCAGATTTGATGTCAAAGCTGGGTCGCTTTTGGAGTCAAGGCAACGATTTCAACCTGTTCAAAATAACTACGGTCCCTTAT
15301 Q H I E K I F P R Q I T N I L S G I K P K L Y G E G K N V R D W I H
CAACATATCGAAAAATTCATCCACGCTCAGATTACTAACCTCCTAAGTGAATTAAGCCAAAACCTTTACGGTGAAGGTAAGAATGTTGCTGACTGGATTC

T N D H S S G V W T I L T K G Q I G E T Y L I G A D G E K N N K E  
 15401 ATACCAATGACCATTCTTCAGGAGTTGGACAATCTTGACAAAAGGGCAAAATCGGTGAAACCTACTTGATTGGGGCTGATGGTGAGAAGAACAAATGAGGA  
 V L E L I L K E M G Q A T D A Y D H V T D R A G H D L R Y A I D A  
 15501 AGTTTGGAACTTATCCTTAAGGAAATGGGACAAAGCTACGGATGCCTATGATCATGTGACTGACCGTGCAGGACATGACCTTCGCTATGCGATTGATGCC  
 S K L R E E L G W K P E F T N F E A G L K A T I K W Y T D N Q E W W  
 15601 AGCAAGCTCCGTGAGGAGTTGGGGTGGAAACCTGAATTTACCAACTTTGAAGCTGGGCTCAAGGCAACAATCAAGTGGTATACAGATAACCAAGAAATGGT  
 K A E K E A V E A N Y A K T Q E I I T V \*  
 15701 GGAAGCAGAAAAAGAGCTGTTGAGCCAAATTAAGCTCAGGAGATTATTACAGTATAAAAAAGCAGGAAATAGCTGCTTTTTATTGCTATATTG

*cps19a0* →

M I L I T G A N G Q L G T E L R Y L L D E R N E  
 15801 GGAAGAGTTACATATTAGAAAGGTCTAGAGATGATTTTAAATTACAGGGCAAAATGGCCAATTAGGAACGGAACTTCGCTATTATTGGATGGAACGTAATG  
 E Y V A V D V A E M D I T D A E M V E K V F E E V K P T L V Y H C  
 15901 AAGAATACGTGGCAGTAGATGTGGCTGAGATGGACATTACCGATGCAGAAAATGGTTGAGAAAAGTTTTGAAGAGGTGAAACCGACTTTAGTCTACCACGTG  
 A A Y T A V D A A E D E G R E L D F A I N V T G T K N V A K A S E  
 16001 TGCAGCCTACCCGCTGTGATGCAGCAGAGGATGAAGGAAGAGATTGGACTTCGCCATCAATGTGACGGGGACAAAAAATGTGCAAAAAGCATCTGAA  
 K H G A T L V Y I S T D Y V F D G K K P V G Q E W E V D D R P D P Q  
 16101 AAGCATGGTGAACCTCTAGTTTATATTCTACGGACTATGCTTTGATGGTAAGAAACCGATTGGACAAGAGTGGGAAAGTTGATGACCGACAGATCCCAC  
 T E Y G R T K R M G E E L V E K H V S N F Y I I R T A W V F G N Y  
 16201 AGACAGAAATATGGCCGTAAAGCGTATGGGGGAAGATTTAGTTGAGAAGCATGTGCTAATTTCTATATTATCCGTACTGCTGGGTATTGGAAATTA  
 G K N F V F T M Q N L A K T H K T L T V V N D Q Y G R P T W T R T  
 16301 TGGCAAAAACCTTCGTTTTTACCATGCAAAATCTGCGAAAACCTATAAGACTTTAAACAGTTGTAATGACCAGTACGGTTCGTCAGCTGGACTGTACC  
 L A E F M T Y L A E N R K E F G Y Y H L S N D A T E D T T W Y D F A  
 16401 TTGGCTGAGTTCATGACCTAGCTGAAAATCGTAAGGAATTTGGTTATTATCAATTTGTCAAATGATGCGACAGAAGACACAACATGGTATGATTTG  
 V E I L K D T D V E V K P V D S S Q F P A K A K R P L N S T M S L  
 16501 CAGTTGAAATTTGAAAAGATACAGATGTCGAAGTCAAGCCAGTAGATTCCAGTCAATTTCCAGCCAAAGCTTAAACGTCGCCATAACTCAACGATGAGCCT  
 A K A K A T G F V I P T W Q D A L Q E F Y K Q E V R \*  
 16601 GGCCAAAGCCAAAGCTACTGGATTGTTTATCCCAACTGGCAAGATGCATTGCAAGAAATTTTACAACAAGAAGTGAGATAAGTAGAATGATCTTCT  
 16701 AGTCTAATAAAGAGGCAGATAAATGAACCTCCAAGGAGCTTAAGATGTACGATTATCTTGTGTTGGTGTGCTCTCTTTGGCCCATAGCTTTGGCTCAG  
 16801 CTTCTATTATCGCTCACACCATCCATCAGAAAGTTTAACTGAAAGGTACCCAATTTATCGCCAAGAAGAAGATTGGGCTAGGATGGGTTTACCAATCACACG  
 16901 TAAGGAAATCTCTAATGGCATATCAAGGCAAGTCAATACTATTTAGAGTCCCTTTATAACCTTTTACGAGAAAAGTTGTTAGAACAACCTCTCTCAT  
 17001 GCGGATGAAACCTCTTATCGGGTGTAGAGAGTGTAGCCATCTGACCTACTATTGGACCTTTTGTCTGGGAAAGCTGAGAAATCAAGCAATCACCGTGT  
 17101 ACCATCATGATCAGCGTCGAGTGGTTTAGTAGTACAAGAATTCCTAGGAGATTATCTGGCTATGTGCATTTGTGATATGTTGGCCGACTAAGCTTAGGAC  
 17201 TTTAGTCCCTTAGTCTCGCTATGCGATAGCAGTCCAAAGTTTAGGAGCAAGGCGAGCTAAGCTTTGGTAACTACGAAACCGCTAGAAGCTTATCGTCAA  
 17301 CTGAAAAGAAGCTGAACCTGTTGGATGTTGGGCGCATGTGAAGAAGGAAGTTTTTTGAAGCGCCCCAAGCAAGCGGATAAATCATCGTTAGGAGCTAAA  
 17401 GGTTTAGCTTATGTTGATCAGTTATTTGCCTTGGAAAAGAGACTGGGAGGCTTTGCTAGCTGATGAACGACTACAGAAACGTCAGAAGAGCTCCAACCC  
 17501 TAATGGAAGATTTCTTTGCTTGGTCCGCGCTCAGTCAGTTTATCGGGTCAAAAACCTAGGAAGGGCAATTAATAACAGCCCTCAAGTATAAAGAAACCTT  
 17601 TAAGACCATTTTAAAAGAGCGGACATCTGGTCTTTTCCAATAATCTAGCTGAACGCGCCATTAATCAATGTTGTTATGGGACGGAGTAAAAGAGTCCAGTGG  
 17701 ACTCTTTTAGCCTAAGCTCAGTTTAAAAAAGTGAAGGTTGTTATTTCTCAAAATTTTGAAGGAGCTAAAAGCAAGAGCTATTATTATGAGCTTATTGGAA  
 17801 ACAGCTAAACGTCATCAACTAAATAGCGAGAAATATCTATCCTATCTTCTAGAATGTCTTCCAACGAGGAACTCTCGTAAACAAGAGGTTTTAGAGG  
 17901 CTTATTTACCAATGTAATAAGTTGTACAAGAAAGTGAATAAAGAAATCTCCAGATTAGGAACATATATGAGTTCTCTAGTCTGGAGATTTTTCAATA  
 18001 TACTTCGTTATTGGACGGTTACGATATTCATATTTTTTGAAGAAGTGTGTTGAAAATAAATTTTCAAAAATTTCTGAAAATTTCTGTTGACAACTTTCTG

*aliA* →

M K S S R L F A L A G V T L L A A T T L  
 18101 AAAAGAGTCTATAATGGAGAGAAAGTTTTAAAGGAGAAAATGATGAAAAGTTCAAGACTATTGGCCCTGGCGGGCGTGACATTATTGGCGGGGACTACTTT  
 A A C S G S G S S T K G E K T F S Y I Y E T D P D N L N Y L T T A  
 18201 TAGCTGCATGCTCTGGATCAGGTTCAAGCACTAAAGGTGAGAAAACATTTCTCATACATTTATGAGACAGACCCCTGATAACCTCAACTATTGACAACCTGC  
 K A A T A N I T S N V V D G L L E N D R Y G N F V P S M A E D W S  
 18301 TAAGGCTGCGACAGCAAAATATTACCAGTAACGTGGTTGATGGTTTGTCTAGAAAATGATCGCTACGGGAACTTTGTGCCCTCTATGGCTGAGGATTTGGTCT  
 V S K D G L T Y T Y T I R K D A K W Y T S E G E E Y A A V K A Q D F  
 18401 GTATCCAAGGATGGATTGACTTACACTTATCTATCCGTAAGGATGCAAAATGGTATACTTCTGAAGGTGAAGAATACGCGGAGTCAAAAGCTCAAGACT  
 V T G L K Y A A D K K S D A L Y L V Q E S I K G V D A Y V K G E I  
 18501 TTGTAACAGGACTAAAATATGCTGCTGATAAAAAATCAGATGCTCTTACCTTGTTCAGAAATCAATCAAAGGGGTGGATGCCTATGTAAGGGGAAAT  
 K D F S Q V G I K A L D D Q T V Q Y T L N K P E S F W N S K T T M  
 18601 CAAAGATTTCTCAGGTAGGAATTAAGGCTTTGGATGATCAGACAGTTCAGTACACTTTGAACAAACAGAAAGTTTTTGGAACTCAAAAACAACCACTG  
 G V L A P V N E E L L N S K G D D S  
 18701 GGTGTGCTTGGCCAGTTAATGAAGAGTTGTTGAACCTAAAAGGGGATGATTTCT

# APPENDIX IV

## The nucleotide and amino acid sequence of the *cps19a* locus from *cps19aJ-aliA* from *S. pneumoniae* strain 19A2.

The nucleotide sequence is numbered in accordance with GenBank accession number AF105113 and starts from the centre of *cps19aJ* to the beginning of *aliA*. This sequence has been referred to as *cps19a<sub>2</sub>* throughout this thesis to distinguish it from the *cps19a* locus from *S. pneumoniae* strain 19A1. The amino acid translation for each ORF is represented by single letter code above the first nucleotide of each codon, except when translation is in the complimentary DNA strand, when it is shown below the nucleotide sequence. Possible ribosome binding sites are underlined. The promoters for *aliA* and *cps19a<sub>2</sub>O* are indicated and the stem-loop termination sequence is indicated with arrows.

```
...cps19a2J→
  S K F I F N I K E S F D N I R V G I N L M L S N I A S S L I I G I
1  TGTCAAATTCATCTTTAATATAAAAGAGTCTTTTGATAATATCCGAGTTGGTATCAATTTAATGTTATCCAACATCGCAAGTAGTTGATTATCGGTAT
  V R L G I Q W N W N I E T F G K V S L T L S V S N L L M T F I N A
101 TGTTCGTCTAGGAATCCAATGGAACGGAAATATCGAAACATTCGGAAAGTATCTCTGACCTTGAGTGTTCCTCAATCTATTGATGACTTTCATTAATGCC
  I G L V V F P L L R R T K A E N L P K I Y S N L R N V L M L I M F A
201 ATTGGACTAGTCGTCTTTCCGCTAATTAAGACGAACAAAAGCAGAAAACCTACCTAAGATTTATTCTAATTTAAGAAACGTTTGTATGCTTATCATGTTGC
  I L L I Y Y P L K I V L D L W L P A Y Q D A L I F M T L I F P M S
301 CGATTTTGCTCATTTACTATCCTTTAAAAATGTTATTAGACCTCTGGTTGCCAGCCTATCAAGATGCCTTGATTTTCATGACCCTTATTTCCCTATGTC
  V Y E G K M A L V I N T Y L K A L R M E R D I L R V N A L V M L I
401 AGTCTATGAAGGGAAAATGGCATTTGGTCAATTAATACTTACTTAAAGGCATTAAGAATGGAAAGAGATATCCTTAGAGTCAATGCCTTAGTAATGTTGATC
  S M G V T L V T T Y L L N S L E L T V V S I V V L L A L R S I I A E
501 AGTATGGGAGTGACCCTGGTTACAACATACCTGTTAAATAGTTTGGAGCTGACTGTTGATCGATAGTTGTTTGGCTAGCTTTGAGAAGTATCATAGCCG
  L I L S K K L D V S V K K D I V L E F L L T L V F I S S S W Y L P
601 AATTAATTCGTCTAAAAAAGTGGATGTTTCGGTTAAGAAGGATATTGTTATTAGAATTTCTTTTGACGCTTGCTTTATTCTTCAAGTTGGTATTGTC
  I G L A V I V Y T I A Y G L Y L Y L K H E D I K T Y L A Y F K A S
701 GATTTGGCTGGCAGTAATTGCTACACAATAGCCTACGGTTTATATCTCTATTTTGAAACACGAAGATATCAAAACCTATTTAGCTTACTTTAAAGCTAGT

                                     cps19a2K→
K K T S N *
801 AAAAAAACATCAAATTAATAATGATTAAGTGGTAGATTCTATTTCTACCGTTTTAGATATTCGGGAGGTAATGATGAAAAAATATGTTA
  V F G T R P E A I K M C P L V N E L K K H E D M E T I V C V T G Q H
901 GTATTCGGTACACGTCGGAAGCCATCAAAATGTGTCCATTAGTCAATGAGTTGAAAAAACACGAAGATATGGAACAATTTGTGTGTACTGGACAAC
```

K E M V S P V L D L F G V V P D Y D L E I M K A N Q T L F S I T T  
 1001 ACAAAAGAGATGGTTAGTCTGTTTATAGATTTATTTGGTGTGTACCAGATTATGATTAGAAAATTGAAGGCTAACCAAACCTTGTCTCTATCAACA  
 S I L E K I K P V L E K E Q P D I V L I H G D T T T T Y A A A L A  
 1101 TAGTATCTTGGAAAAGATAAAACAGTTTATAGAGAAGGAACAACAGATATTGTCCTAATTCAGGGTGACACTACGACAACCTTATGCAGCAGCCTTGGCA  
 A F Y L G I K V G H V E A G L R T Y N L Q S P F P E E F N R Q S T S  
 1201 GCATTCTATTTGGGAATTAAGTAGGACATGTGAAGCTGGTTGGCGAAGCTACAATTTACAAAGTCCATTCTCGAAGAATTTAACAGGCAATCGACAT  
 I I A T Y H F A P T E L A K E N L L K E G R E N V Y V T G N T V I  
 1301 CAATCATTGCAACTTACCATTGCTCCAACCTGAGTTGGCTAAAGAAAATCTCTTAAAGAAAGGTAGAGAGAATGTTTATGTGACTGGAAATACCTGTCAAT  
 D A L T T T V Q E D Y T H T H L D L N A N N R L I L L T A H R R E  
 1401 TGATGCTCTTACAACCTACTGTTCAAGAGGATTATACACACACTCATTAGATTTAAACGCTAACAAATCGTCTCATCTTATTGACTGCTCATAGACCGGAA  
 N L G E P M R H M F R A V K R V L N E Y E D V K V I Y P I H K N P L  
 1501 AATCTCGGCAACCGATGAGACATATGTTTAGAGCAGTTAAACGAGTATTGAATGAATACGAAGATGTTAAAGTCATTATCCAATTCATAAGAATCCCTT  
 V R E T A A E I F G D T E R I Q I I E P L D V L D F H N F M N Q S  
 1601 TGGTACGTGAACAGCTGCAGAGATCTTTGGAGATACAGAACCGATTCAAAATTTAGAACCTTTGGATGTTCTTGAATTCACCAACTTCATGAACCAAG  
 Y M I L T D S G G V Q E E A P S L G K P V L V M R D T T E R P E G  
 1701 TTATATGATTTTAAACAGATTCTGGCGGAGTTCAGGAAGAAGCCTTCTTTAGGAAAACCTGTATTGGTTATGGCGGATACTACTGAAAGACCTGAAGGA  
 V A A G T L K L V G T D E E T I Y Q N F K L L L D D S G E Y K K M S  
 1801 GTAGCGGCAAGAACCTTGAATTAGTAGGAACGACGAAGAACAATTTATCAAACTTCAAACTACTTCTAGATGATTCTGGAGAATACAAAAAATGA  
 Q A S N P Y G N G D A S Q Q I V Q I L R G I \*  
 1901 GTCAGGCTAGCAATCCATACGGTAATGGTATGCTAGTCAACAGATGTTTCCAGATTTTCCGTGGGATTTAAGAAGATTTCTTAAAGAGTCTAAGAGGAA  
 2001 TCCACTCCACAAACATAAACTCTTATGCTATCATTATGATTGAGATAGACTCATAACTTAATATTTTATCCTAAGAAAATTCGTTTCATATTCAATAT

*cps19a<sub>2</sub>L*→

M K G I I L A G G S G T R L Y P L T R A A S K  
 2101 GAATGAAAAACAACGAAGAAAGTTATTTTATTGAAAGGTATTTCTTGCAGGGGGCTCAGGTACCCGCTGTACCCTTACTCGGGCTGCGTCAA  
 Q L M P V Y D K P M I Y Y P L S T L M L A G I K D I L I I S T P Q  
 2201 AACAGCTGATGCCGGTTTATGATAAACCTATGATTTATTATCCGTTGTCGACATTAATGTTGGCTGGAATTAAGATATTTTGATTATCTCAACTCCTCAA  
 D L P R F K D L L L D G S E F G I K L S Y A E Q P S P D G L A Q A  
 2301 AGATTTGCCCCGTTTAAAGGACTTGCCTTGGATGGTCCGAATTTGGGATCAAGCTTTCCTATGCGGAACAACCTAGTCCCGATGGACTTGTCTCAGGCT  
 F L I G E E F I G D D S V A L I L G D N I Y H G P G L S T M L Q K A  
 2401 TTTCTTATCGGTGAAGAATTTATCGGTGACGATAGCGTTGCTTGAATTTTAGGGACAAATATCTACCATGGTCCCTGGTCTGAGCAAAATGCTTCAAAAAG  
 A K K E K G A T V F G Y Q V K D P E R F G V V E F D T D M N A I S  
 2501 CAGCCAAGAAAAGAGAAAGGTGCGACTGTTTTTGGCTACCAAGTGAAGGATCCAGAGCGTTTTTGGTGTGGTTGAGTTGATACAGACATGAATGCTATTTCT  
 I E E K P E Y P R S N Y A V T G L Y F Y D N D V V E I A K Q I K P  
 2601 TATAGAAGAAAAGCCGGAGTATCCTCGCTCCAACATGACAGTGACAGGACTGTATTTCTATGATAATGATGTTGTTGGAGATGCTAAACAGATCAAACCT  
 S A R G E L E I T D V N K A Y L N R G D L S V E L M G R G F A W L D  
 2701 AGTGTCTGTTGGCGAGTTAGAAAATTACAGACGTTAAACAAGGCTTACCTAAATCGTGGTGACCTTTCTGTTGAGCTGATGGGGCTGGTTTTGCTCGTTGG  
 T G T H E S L L E A S Q Y I E T V Q R M Q N V Q V A N L E E I S Y  
 2801 ATACGGGAACCCATGAAAGCTTGCTAGAAGCTTCTCAGTATATTGAACCGTTCAGCAGAACGTTCAAGTGGCAAATCTAGAAGAAAATTTCCATA  
 R M G Y I S R E D V L E L A Q P L K K N E Y G R Y L L R L I G E A  
 2901 TCGTATGGGCTACATTAGTCTGGAAGACGTGCTGGAATTTGGCTCAGCCTCTTAAGAAAATGAATACGGACGATATTTGCTCCGTTGATTGGAGAAGCA

*cps19a<sub>2</sub>M*→

\* M S D N F F G K T L V V R K I D A I P G L L E F D I P V H G D N R  
 3001 TAGATTCAGATAAATTTTTTGGAAAGACACTTGTGGTACGCAAGATTGATGCTATACCAGGACTGCTAGAGTTTGACATTTCCCGTTCATGGAGACAATC  
 G W F K E N F Q K E K M L P L G F P E S F F A A G K L Q N N V S F  
 3101 GTGGTTGGTTTAAAGAAAATTTCCAGAAGGAAAAGATGCTACCCTTGGTTTTCTGAAAGCTTCTTTGCTGCAGGGAACGCAAAAATAACGTCAGCTT  
 S R K N V L R G L H A E P W D K Y I S V A D D G K V L G S W V D L  
 3201 TTTCTCGCAAAAATGTTCTTCGAGGACTCCATGCTGAACCTTGGGACAAGTATATCTCTGTTGCAGACGATGGGAAGTTTTAGGATCTTGGTAGATTTG  
 R E G E T F G N V Y Q T E I D A S K G I F V P R G V A N G F Q V L S  
 3301 CGTGAAGGGCAGACCTTTGGAATGTTTACCAGACAGAGATTGATGCTAGCAAGGAAATTTTGTCCCTCGTGGCGTAGCCAACGGTTTTCAAGTCTCT  
 D T V S Y S Y L V N D Y W A L E L K P K Y A F V N Y A D P A L G I  
 3401 CAGATACAGTTTCCATAGCTATCTGGTCAATGACTACTGGGCTCTGAACTCAAACCCAAATATGCTTTGTTAACTATGCTGACCCAGCATTGGGAAT  
 E W E N L P E A E V S E A D K H H P L L R D V K P L T K D E L \*  
 3501 TGAGTGGGAAAACCTACCGAAGCTGAGGTTTCAGAGGCAGACAAACACCATCTCTATTAAGGGATGTCAAACCACTTACGAAAAGACGAGTTGTAAGAAAG

*cps19a<sub>2</sub>N*→

M T E Y K K I I V T G G A G F I G S N F V H Y V Y N N F P D V  
 3601 GAAACATTTAGCTGAATACAAAAAATATCGTGACAGGTGGAGCTGGTTTTATCGGTTCTAACTTTGTCCACTATGTTTACAATAACTTTCCAGATGT  
 H V T V L D K L T Y A G N R A N I E E I L G D R V E L V V G D I A  
 3701 CCATGTGACAGTCTGGACAAGCTGACTTATGACAGTAACTGTCGAATATCGAGGAAATTTTAGGCGACCGTGTGAGTTGGTTGTTGGAGATATTGCT  
 D A A L V D K L A A E A D A I V H Y A A E S H N D N S L N D P S P F  
 3801 GATGCAGCCTTGGTAGACAAGTTGGCGGCTGAAGCGGATGCTATCGTTCACTATGCGGCAGAAAGCCACAATGATAACTCGCTCAATGACCCGAGTCCAT  
 I H T N F I G T Y T L L E A A R K Y D I R F H H V S T D E V Y G D  
 3901 TTATCCACCAACTTTATCGGAACTTACACACTTTTAGAAGCGGCTCGTAAATACGACATTCGTTTCCACCATGATATGACTGATGAAGTCTATGGTGA  
 L P L R E D L P G H G E G L G E K F T A E T K Y N P S S P Y S S T  
 4001 CCTGCCTCTCGGTGAAGATTGCGCAGGTCAATGGGAAGGCTAGGTGAGAAATTTACCCTGAAACCAAGTACAATCCAAGCTCGCCTTACTCATCAACC

4101 K A A S D L I V K A W V R S F G V K A T I S N C S N N Y G P Y Q H I  
 AAGGCTGCTTCAGACTTGATCGTTAAAGCTTGGGTGCGCTCATTTGGTGTAAAGCAACGATTTCTAACTGTTCAAACAACATATGGTCCATACCAGCATA  
 4201 E K F I P R Q I T N I L S G I K P K L Y G E G K N V R D W I H T N  
 TCGAGAAGTTTCATTCGCGCCAGATTACGAATATCTTGAGCGGTATCAAGCCAAAACCTTACGGAGAAGGTAAGAAATGTGCGTGACTGGATTACACCAA  
 4301 D H S S G V W T I L T K G Q I G E T Y L I G A D G E K N N K E V L  
 TGACCATTCATCAGGCGTTTGGACGATTCTGACCAAGGGTCAAAATCGGTGAAAACCTTACTTGATTGGTGTGACGGTGAGAAGAACAACAAGGAAGTGTTA  
 4401 E L I L K E M G Q P A D A Y D H V T D R A G H D L R Y A I D A S K L  
 GAGCTCATTCTCAAGGAGATGGGGCAACCCGCTGACGCTTATGACCATGTGACCGATCGAGCTGGTCACGACTTGCCTTATGCGATTGATGCTAGCAAGC  
 4501 R D E L G W K P E F T N F E A G L K E T I K W Y T D N Q D W W K S  
 TCCGTGATGAGTTGGGTGGAAGCCAGAGTTTACCAACTTTGAAGCAGGCCCAAAGAGACCATCAAGTGGTACACAGATAACCAAGACTGGTGGAATC  
 4601 E K E A V E A N Y A K T Q Q V I K \*  
 TGAAAAAGAAGCAGTAGAGGCTAACTATGCTAAAAACGCAACAAGTGATTAAATAAAGATACTAAAAAGCAAGAGACCTCGAGGTCTCTTGCTTTTCT

4701 AGTATGTAGTAAGCTTTTATTTTCTTACTTCTTGTTTATAAACTCTTTAAGGCATCTTGCCAGGTTGGGATGACGAAACCTGTTGCCTTTGCCTTTGC  
 \* K R V E Q K Y F E K L A D Q W T P I V F G T A K A K A  
 4801 CAAACTCATAGTTGAGTTGAGAGGTCGTTTAGCCTTAGCGGGAACTTGCTAGAGTCTACTGGCAATAGTTCAACATCACTCTCCTTAAGAATCTCGCTA  
 L S M T S N L P R K A K A P F K S S D V P L L E V D S E K L I E S  
 4901 GCAAAGTCATACCAAGTGGTCTCCTCAGCTGCGTCATTCGATAAGTGATAGTAGCCATACTCTTTTGTATTTTTCAGTCACATAGGTCATGAACTCGGCCA  
 A F D Y W T T D E A A D N S L H Y Y G Y E K Q N E T V Y T M F E A  
 5001 AGGTCCGTGTCCTCAAGTTGGGCGACCGTATTGGTTCGCTGACCACTGTGAGCGTTTTATGGGTTTCGGCTAGGCTTTGTCATGGTAAAGCAAAGTTCCCTTC  
 L T R T W T P R G Y Q D S V V T L T K H T E A L S Q M T F V F N R G  
 5101 ATAATTTCCAAAAACCAAGCAGTACGAATGATGTAATGCTGTGACGTAAGGTTCTCAACAAGTTCTTCTCCATTCGCTTGGTACGTCCTGACTCTGTT  
 Y N G F V W A T R I I Y H Q S T L N E V L E E G M R K T R G Y E T  
 5201 TCGGATCAGGTATGTCATCGACTTCCCACTCTTGCTTACTGTTTCTTCCATCAAAGACGTAGTCTGTTGAGATATAGACCAGAATAGTCCGTTATT  
 Q P D P I D D V E W E Q G V P K K G D F V Y D T S I Y V L I A G Y  
 5301 TCTCAGAGCCCTTGCTACATTTTTCAGTTCAGTTACGTTGATGGCAAATCCAGTCTTTCCCTTCATCTTCGGCTGCATCAACAGCAGTGTAGGCTGC  
 K E S G K A V N E T G T V N I A F D L E K G E D E A A D V A T Y A A  
 5401 ACAATGATAGACTAGAGTTGGCTTAACTCAGCAAAGACTTTTCAACCATTTTCAAGATTTAGTGATATCCATTTTCGGTTCACATCCACTGCAACATAGTCC  
 C H Y V L T P K V E A F V K E V M E S N T I D M E T V D V A V Y D  
 5501 ACATTTGCTCATTTAGTAAATAACGTAGTTTCGGTACCAAGTTGACCATTTGCACCTGTAATTAAGATCATATTTTCTCCTTGATAGTTCTTATACTTA  
 V N R E N L L Y R L E T G L Q G N A G T I L I M

←*cps19a2O*

5601 ATTATATCAAAAATATAGAAAAATCGGTTTTAGTGAACCCTTTTTGGTATAAACTTTGTTAGTGAAGCGCAAGAGAGCCGAATAGTTATTACTTTAGA  
 5701 ATATCCTTAAAAGTATTTTTCAGAATTTACCAAATTAATTAATAAAATTCAGAAAATTTATA<sup>-35</sup>TTGACATCTCTCTGAAAAGAGTCTATA<sup>-10</sup>ATAGAGAGAAAGT<sup>-10</sup>  
<sup>-10</sup>

*P<sub>aliA</sub>*→

←*P<sub>cps19a2O</sub>*

*aliA*→

M K S S K  
 5801 TTTAAAGGAGAAGATGATGAAAAGTTCAAAAC

# APPENDIX V

The nucleotide sequence of Rx1-19F, Rx1-19A.1, Rx1-19A.2, Rx1-19A.3 and Rx1-19A in the regions where recombination between the *cps19f* and *cps19a* loci has occurred.

## A.

The nucleotide sequence of *cps19F* and *cps19G* for Rx1-19F, Rx1-19A.1, Rx1-19A.3 and Rx1-19A indicating the regions where recombination at the 5' end of the *cps* locus has occurred in Rx1-19A.1, Rx1-19A.3. The amino acid translation for Rx1-19F is represented by single letter code above the first nucleotide of each codon. The vertical arrows indicate the first nucleotide which is type 19A-specific in each transformant.

*cps19f*→

	M R D R I Q L L G V T I D L L T M N E T I D S V E Q Y V L E K R
Rx1-19F	ATGAGGGATAGAATCCAACCTTTTAGGTGTAACAATTGATTTGCTTACGATGAATGAAACGATAGATAGTGTAGAACAATATGTATTAGAAAAAAG
Rx1-19A.1	.....
Rx1-19A.3	.....
Rx1-19A	...AT..G.....T...A.A.....AA.T..T..A...CCAT.A.....A.....A.TG...GC...T.....G.....T...A.G...GCA
	P L H L M G V N A D K I N Q C H T D E K I K K I V N E S G I I N
Rx1-19F	ACCACTACACTTGATGGGGTGAATGCTGATAAAAATTAATCAGTGCATACAGATGAGAAAATCAAAAAATCGTTAATGAGTCAGGAATCATT
Rx1-19A.1	.....
Rx1-19A.3	.....
Rx1-19A	C..TT.G.....A..T..C..G.....CT..GA.....A.....A.....
	A D G A S V V L A S K F L G T P V P E R V A G I D L M Q C L L
Rx1-19F	ATGCGGATGGAGCATCAGTTGTCTTGCAAGTAAGTTTTTAGGAACGCCCTGTTCCCTGAACGAGTAGCGGGTATTGATTGATGCAATGCTCTTTTA
Rx1-19A.1	.....
Rx1-19A.3	.....
Rx1-19A	...C..C..T..T...A..T.A..G.....AC.A.....G.....T..T..A..A..CT.....A..CA.T.AC..
	E L S N K K G Y S V Y F F G A K E E V L Q D M L K V F K R D Y P
Rx1-19F	GAGTTGTCAAAATAAAAAAGGATATTCAGTTTACTTTTTTGGAGCAAAAGAAGAAGTTTTGCAAGATATGCTCAAAGTATTTAAGAGAGATTATCC
Rx1-19A.1	.....
Rx1-19A.3	.....
Rx1-19A	.....T..G.....C..G..C.....C.....A.....T.A.....T..G..A..A.....
	N L I V I G H R N G Y F S E E D E Q A I Q E D I R E K N P D F V
Rx1-19F	AAATTGATAGTTATTGGACACAGAAATGGCTATTTTCTGAAGAGGATGAAACAAGCTATTCAGAAGATATTCGTGAAAAGAACCCTGATTTTG
Rx1-19A.1	.....
Rx1-19A.3	.....
Rx1-19A	..GC.C.A.AT..G...CT.T.....A.....GCCT..A.....A..CA.....G.....AAA..G..A..A..A.....

## Rx1-19A.1



F I G I T S P K K E Y I I Q K F M D S G V N S V F M G V G S  
 Rx1-19F TGTTTATTGGAAATTACGTCCTCTAAAAAGAAATATATTATTTCAAAAATTTATGGATAGTGGCGTCAATTCCGGTATTTATGGGAGTTGGCGGTAGT  
 Rx1-19A.1 .....T.A.  
 Rx1-19A.3 .....T.A.  
 Rx1-19A .....G.A.....T.....G.....C.TC.C.....A.AA.T.G.T.G.....G.....T.A.

F D V L S G H I Q R A P L W M Q K S N L E W L F R V A N E P K R  
 Rx1-19F TTTGATGTCCTGTCTGGTCATATCCAACGAGCACCTCTATGGATGCAAAAGTCAAATTTAGAGTGGTTATCCGTGTAGCTAATGAGCCTAAACG  
 Rx1-19A.1 .....AC.A.A.G.....A.....T.AT.....TG.TC.CC.....G.T.....G.A.....  
 Rx1-19A.3 .....AC.A.A.G.....A.....T.AT.....TG.TC.CC.....G.T.....G.A.....  
 Rx1-19A .....AC.A.A.G.....A.....T.AT.....TG.TC.CC.....G.T.....G.A.....

*cps19fG*→

L F K R Y F V G N I S F I G K V L K A K R G V K Y \* M I  
 Rx1-19F TCTCTTTAAACGTTATTTGTAGGGAATATTTTCATTTCATAGGAAAAGTTTAAAGCAAAAAGAGGTGTAATAATTTGAACCAGACAGAGATGAT  
 Rx1-19A.1 .....GC.....G.....T.A.....G.....G.....G.....ATAT...CT.....C.....A.AG.  
 Rx1-19A.3 .....GC.....G.....T.A.....G.....G.....G.....ATAT...CT.....C.....A.AG.  
 Rx1-19A .....GC.....G.....T.A.....G.....G.....G.....ATAT...CT.....C.....A.AG.

## Rx1-19A.3



R L I Q K V E L D A I K E F K K I C E E N D I D F F L R G G S V  
 Rx1-19F TCGCTTAATTCAAAAGTTGAATTAGATGCTATAAAAGAGTTTAAAAAATCTGTGAAGAGAATGATATAGATTTTCCCTCCGCGGTGGTAGTG  
 Rx1-19A.1 .AAA.G.....G.....G.....C.....T.A.....CA.....C.....T.T.G.G.....  
 Rx1-19A.3 .....AAA.G.....G.....G.....C.....T.A.....CA.....C.....T.T.G.G.....  
 Rx1-19A .....AAA.G.....G.....G.....C.....T.A.....CA.....C.....T.T.G.G.....

L G A V K Y D G F I P W D D D M D I A V P R E A Y D K L P S V  
 Rx1-19F TACTTGGTGCAGTCAAAATACGACGGCTTTATTCCATGGGATGATGATATGGATATCGCTGTCCCTCGTGAAGCATACGACAAAAGTTCCAAGTGT  
 Rx1-19A.1 .....G.....T.T.....C.....T.....G.....GC.....T.....G.A.C  
 Rx1-19A.3 .....G.....T.T.....C.....T.....G.....GC.....T.....G.A.C  
 Rx1-19A .....G.....T.T.....C.....T.....G.....GC.....T.....G.A.C

F K D R I I A G K Y Q V L T Y Q Y C D T L H C Y F P R L F L L A  
 Rx1-19F TTCAAAGATAGAATTATCGCTGGGAAATATCAGGTTCTTACTTATCAACTGTGATACGTTGCATTGCTACTTCCCTCGACTATTCCTTTTAGC  
 Rx1-19A.1 .T.G.....A.C.....G.....T.C.....T.A.....CT.....A  
 Rx1-19A.3 .T.G.....A.C.....G.....T.C.....T.A.....CT.....A  
 Rx1-19A .T.G.....A.C.....G.....T.C.....T.A.....CT.....A

D E R K R L G L P R N T N L G L H L I D I I P L D G A P N H S V  
 Rx1-19F AGATGAAAGAAAACGTTTGGGCTTGCCACGAAATACCAATCTAGGATTTGCATTTGATTGATATCATTCCTTTAGATGGAGCCAAATCATTCCG  
 Rx1-19A.1 .....A.....T.....G.....T.C.....A.....AT  
 Rx1-19A.3 .....A.....T.....G.....T.C.....A.....AT  
 Rx1-19A .....A.....T.....G.....T.C.....A.....AT

L R K I Y F C K V Y W Y R F L A S L G T T Y V G D H V D M H S  
 Rx1-19F TTTTAAGAAAGATTTACTTTTGTAAAGTATACTGGTATCGTTTATTTAGCAAGCTTAGGAACAACTTATGTTGGCGACCATGTGGATATGCATTCC  
 Rx1-19A.1 .....C.....CG.....T.....T.C.....A.....  
 Rx1-19A.3 .....C.....CG.....T.....T.C.....A.....  
 Rx1-19A .....C.....CG.....T.....T.C.....A.....

T K Q K L I I G F F K K L G F A K L F P Q N S V Y R R L D N L Y  
 Rx1-19F ACTAAGCAAAAATAATTTATGGTTTCTTTAAAAAAGTAGGATTTGCAAAAATTTCCCTCAAAATTTCTGTATACAGACGCTTGATAATCTCTA  
 Rx1-19A.1 G.....C.....G.....C.....A.G.....  
 Rx1-19A.3 G.....C.....G.....C.....A.G.....  
 Rx1-19A G.....C.....G.....C.....A.G.....

R K Y D W K K Q K Y A G T I N A S L F A K E V M P V E I W G E G  
 Rx1-19F TAGAAAGTATGATTTGGAAAAGCAGAAGTATCGGGGACTATCAATGCTTCTTTATTTGCTAAAGAAGTTATGCCAGTAGAGATTTGGGGAGAAG  
 Rx1-19A.1 .A.....T.....A.T.....A.C.G.....C.....  
 Rx1-19A.3 .A.....T.....A.T.....A.C.G.....C.....  
 Rx1-19A .A.....T.....A.T.....A.C.G.....C.....

V E K P F E D T F F K V P T E Y D R Y L K R L Y G E N Y L Y E  
 Rx1-19F GAGTAGAGAAGCCTTTTGAGGATACCTTCTTTAAAGTTCCAACGGAGTATGATCGCTACCTGAAAAGACTTTACGGAGAAAATCTCTTTACGAA  
 Rx1-19A.1 .....C.....C.G.....T.....C.....  
 Rx1-19A.3 .....C.....C.G.....T.....C.....  
 Rx1-19A .....C.....C.G.....T.....C.....

E P S D D E K K S H L G G H \*  
 Rx1-19F GAGCCTAGTGTATGATGAAAAGAAATCGCATTTAGGAGGACACTAA  
 Rx1-19A.1 .....A.....  
 Rx1-19A.3 .....A.....  
 Rx1-19A .....A.....



```

Rx1-19F      GGCGCATAGCTTTGGCTCAGTTTCTATATATCGCTCACACCATCCATCAGAAGTTTAACTGTAAGGTACCCAATTATCGCCAAGAAGAAGATTGGG
Rx1-19A.1
Rx1-19A.3    .....C.....
Rx1-19A      .....C.....

Rx1-19F      CTAGGATGGGTTTACCAATCACACGTAAGGAAATCTCTAATTTGGCATATCAAGGCAAGTCAACTACTATTTAGAGTCCCTTTATAACCTTTTACGA
Rx1-19A.1
Rx1-19A.3
Rx1-19A

Rx1-19F      GAAAAGTTGTTAGAACAACTCTTCTTCATGCGGATGAAACCTCTTATCGGGTCTTAGAAAAGTGATAGTCAGCTGACCTACTATTGGACCTTTTT
Rx1-19A.1
Rx1-19A.3    .....GC...G...C..T...
Rx1-19A      .....GC...G...C..T...

Rx1-19F      GTCTGGGAAAGCTGAGAATCAAGTAATCACGCTTTACCACCATGATCAGTGTCCGAGTGGTTCGGTAGTGAAGAATTCCTAGGAGATTATTCTG
Rx1-19A.1
Rx1-19A.3    .....C.....G...T...C.....TA...A...
Rx1-19A      .....C.....G...T...C.....TA...A...
    
```

Rx1-19A.3



```

Rx1-19F      GCTATGTGCATTTGTGATATGTTGCGGCAGTAACTTAGGACTTTAGTCTCTAGTCTGTCTATGCGATAGCAGTCCAAGTTTtaggagcaaggcg
Rx1-19A.1
Rx1-19A.3    .....C.....
Rx1-19A      .....C.....

Rx1-19F      ACGCTAAGCTTGGTAAACTGCGAACCGCTAGAAGCTTATCGTCAACTGGAA-GAAGCTGAAGTGTGGATGTTGGGCGCATGTGA-GAAGGAAA
Rx1-19A.1
Rx1-19A.3    .....A.....A.....A.....G
Rx1-19A      .....A.....A.....A.....G

Rx1-19F      TTTTGTGAAGTGCACCCCAAGCAAGCAGATAAATCATCCTTAGGAGCTAAAAGTTTtagcttattgtgatcagttattttccttggaaagagactg
Rx1-19A.1
Rx1-19A.3    .....G.....G.....
Rx1-19A      .....G.....G.....

Rx1-19F      GGAGGCTTTGC-AGCTGATGAACGACTACAGAAACGTCAAGAACATCTCCAACCCCTAAATGGAAGACTTCTTTGCTTAGTGCCGTCGTCAGTCAG
Rx1-19A.1
Rx1-19A.3    .....T.....G.....G.....
Rx1-19A      .....T.....G.....G.....

Rx1-19F      TTTTATCGGGTTCAAAAC TAGGAAGGCAATGAATACAGCCTCAAGTATGAAGAAACCTTTAAGACCATTTTAAAAGACGGACATCTGGTCTT
Rx1-19A.1
Rx1-19A.3    .....A.....
Rx1-19A      .....A.....

Rx1-19F      TCCAATAATCTAGCTGAACGCGCCATTAATCATTTGGTTATGGGACGGAGTAAAAGAGTCCAGTGGACTCTTTTtagcctaagctaaatTTTAAAA
Rx1-19A.1
Rx1-19A.3    .....C.G...A...
Rx1-19A      .....C.G...A...

Rx1-19F      AGCGAGGGTGGTTATTTTCTCAAAGTTTGAAGGAGCTAAAGCAAGAGCTATTATTATGAGCTTATTGGAACAGCTAAACGTCATCAACTAAAT
Rx1-19A.1
Rx1-19A.3    ..T.....A.....
Rx1-19A      ..T.....A.....

Rx1-19F      AGTGAGAAATATCTATCCTATCTTCTAGAATGTCTTCCAACGAGGAAACTCTCGTAAACAAAGAGGTTTtagaggcctatttaccatggactaa
Rx1-19A.1
Rx1-19A.3    ..C.....T.....
Rx1-19A      ..C.....T.....

Rx1-19F      AGTTGTACAAGAAAAGTGCAAATAAGAAATCTCCAGATTAGGAACTATCCGTGAGTTCACATACTGGAGATTTTCAATAGACCTCGTTATTGG
Rx1-19A.1
Rx1-19A.3    .....ATA...T...G.....T..T...
Rx1-19A      .....ATA...T...G.....T..T...

Rx1-19F      GCGGTTACGATATTCATATTTTGTCAAAGATGTTGTTGAAAAATAATTTTCAAAAATTTGAAAATTCGTTGACATCTTCTGAAAAGAGTC
Rx1-19A.1
Rx1-19A.3    .....A.....A...
Rx1-19A      .....A.....A...
    
```

*aliA*→

```

M K S S K L L A L A G V T L L A A T T L A A
Rx1-19F      TATAATGGAGAGAAAAGTTTAAAGGAGAAAAATGATGAAAAGTTCAAAAAC TACTTGCCTTGCGGGCGTGACATTATTGGCGGCGACTACTTTAGC
Rx1-19A.1
Rx1-19A.3    .....G...T.....
Rx1-19A      .....G...T.....
    
```

## C.

The nucleotide sequence of *cps19G*, *cps19H*, *cps19I* and *cps19J* for Rx1-19F, Rx1-19A.2 and Rx1-19A. The amino acid translation for Rx1-19F is represented by single letter code above the first nucleotide of each codon. The region where recombination has occurred is indicated by vertical arrows marking the first and last nucleotides which are type 19A-specific.

*cps19fG*→

```

M I R L I Q K V E L D A I K E F K K I C E E N D I D F F L R G G
Rx1-19F ATGATTTCGCTTAATTCAAAAAGTTGAATTAGATGCTATAAAAAGAGTTTAAAAAATCTGTGAAGAGAATGATATAGATTTTTCCTCCGCGGTGG
Rx1-19A.2 .....
Rx1-19A ..AG..AAA..G.....G.....G.....C.....T...A.....CA.....C.....T..T..G..

S V L G A V K Y D G F I P W D D D M D I A V P R E A Y D K L P S
Rx1-19F TAGTGACTTGGTGCAGTCAAAATACGACGGCTTATTCATGGGATGATGATATGGATATCGCTGCCCTCGTGAAGCATACGACAAAATTCCAA
Rx1-19A.2 .....
Rx1-19A G.....G.....T..T.....C.....T...G.....GC.....T.....G

V F K D R I I A G K Y Q V L T Y Q Y C D T L H C Y F P R L F L
Rx1-19F GTGTTTTCAAGATAGAATTATCGCTGGAAATATCAGGTTCTTACTTATCAACTGTGATACGTTGCATGCTACTTTCCTCGACTATTCCTT
Rx1-19A.2 .....
Rx1-19A ..A.C..T..G.....A..C.....G.....T..C.....T..A.....CT.....

L A D E R K R L G L P R N T N L G L H L I D I I P L D G A P N H
Rx1-19F TTAGCAGATGAAAGAAAACGTTTGGGCTTGCCACGAAATACCAATCTAGGATTCGATTTGATTTGATATCATTCCTTTAGATGGAGCACCAAATCA
Rx1-19A.2 .....
Rx1-19A .....A.....A.....T...T.....G.....T..C.....

S V L R K I Y F C K V Y W Y R F L A S L G T T Y V G D H V D M H
Rx1-19F TTCGGTTTTAAGAAAAGATTTACTTTTGTAAAGTATACGTTATCGTTTTTTAGCAAGCTTAGGAACAACCTTATGTTGGCGACCAATGTTGGATATGC
Rx1-19A.2 .....
Rx1-19A .....AT.....C.....CG.....T.....T..C.....A.....

S T K Q K L I I G F F K K L G F A K L F P Q N S V Y R R L D N
Rx1-19F ATTCCACTAAGCAAAAACCTAATTTATGGTTTCTTTAAAAAACTAGGATTTGCAAAACTATTTCCCTCAAATTTCTGTATACAGACGCTTGGATTAAT
Rx1-19A.2 .....
Rx1-19A .....G.....C.....G.....C.....A.....G.....

L Y R K Y D W K K Q K Y A G T I N A S L F A K E V M P V E I W G
Rx1-19F CTCTATAGAAAGTATGATTGAAAAGCAGAAGTATCGGGGATATCAATGCTTCTTTATTTGCTAAAGAAGTTATGCCAGTAGAGATTGGGG
Rx1-19A.2 .....
Rx1-19A .....A.....T...A..T...A...C.G...C.....CA.....

E G V E K P F E D T F F K V P T E Y D R Y L K R L Y G E N Y L Y
Rx1-19F AGAAGGAGTAGAGAAGCCTTTTGAGGATACCTTCTTTAAAGTTCCAACGGAGTATGATCGCTACCTGAAAAGACTTTACGGAGAAAACATATCTTT
Rx1-19A.2 .....
Rx1-19A .....G.....C.....C.....T.....C.....

E E P S D D E K K S H L G G H * L F C Y I I L H Y K V L E E T
Rx1-19F ACGAAGAGCCTAGTGATGATGATAAAGAAATCGCATTTAGGAGGACATAATTTGTTTTGTTATATATTTTGCATFACAAAGCTTAPAGAAGAAAC
Rx1-19A.2 .....
Rx1-19A .....A.....G.....

↓

I S C V K S I K E G N Y N A K Q I V I I D N F S N N G T G E K L
Rx1-19F TATTTCTGTGTAAATCTATAAAAAGAGCAATTATAATGCAAGCAATCGTTATATTGATAAATTTCTCTAATAATGTTACTGGTGAAGAAAC
Rx1-19A.2 .....
Rx1-19A .....A.....C...A.....G.....T.....C.....G.....

Q E L Y E S D L E I D V L I N H E N A G F A R G N N V A Y Q F
Rx1-19F TACAAGAGCTTTATGAATCAGATTTAGAAATTTGATGTTTATGATTAACCATGAAAATGCTGGTTTTGCTCGTGGAAATAATGTTGGCTTATCAATTT
Rx1-19A.2 .....
Rx1-19A .....G.....C...G.....C.....C.....A.....T.....C.....A.....G...

A K E K Y N P D F M V I M N N D I E I E T E N F E K I V T D I Y
Rx1-19F GCTAAGGAAAAGTATAACCCCGATTTTCATGGTTATCATGAATAACGATATTGAGATAGAAAACAGAAAATTTGAAAAAATTTGACAGATATATCA
Rx1-19A.2 .....
Rx1-19A .....C..A..T...T.....T.....T.....G.G.....C.....

R E E K F H L L G P D I F S T T Y Q L H Q N P K R L T H Y T Y G
Rx1-19F TCCTGAGGAAAATTCATTGCTCGGGCCAGATATCTCTCCTACTTACCAACTTACCACAAAACCCAAACGGTTGACACATATATACTTATG
Rx1-19A.2 .....
Rx1-19A .....GA.....T..A..A.....T...G..G..T..G.....G.....

```

E V K A L N E K F K K G S Q V S L A L K I K C W L K A S K V L  
 Rx1-19F GAGAAGTTAAAGCTCTAAATGAAAAATTTAAAAAGGGAGCCAAAGTTAGTCTAGCTTTAAAAATCAAATGTTGGTTGAAGCTAGTAAAGTTCTT  
 Rx1-19A.2 A . . . G . . . G . . . C . . . . . G . . . . . A . . . . . T . . . . . GT . . . . .  
 Rx1-19A A . . . G . . . C . . . . . G . . . . . A . . . . . T . . . . . GT . . . . .

R T A I Y Q N R R K K G S V D Y R K Q V E N P I L H G S F I V Y  
 Rx1-19F CGAACACGAATCTATCAAATAGACGTAAAAAGGATCAGTAGACTATAGAAAAAGGATAGAAAAACCAATCTTCATGGTCTTTTATTGTATA  
 Rx1-19A.2 . . G . . . . . G . . . . . GAA . . . . . G . . . . .  
 Rx1-19A . . G . . . . . G . . . . . GAA . . . . . G . . . . .

S R D F I E K E E Y A F N P N T F F Y Y E T E I L D Y E A E L K  
 Rx1-19F TTCGAGAGATTTTATCGAAAAAGAGGAGTATGCTTTTAAACCCTAACACCTTCTTTTACTATGAAACAGAGATATAGATTATGAAGCTGAATTAA  
 Rx1-19A.2 . . T . . . . . A . . . . . T . . . . . G . . . . . T . . . . . C . . . . . T . . . . . G . . . . .  
 Rx1-19A . . T . . . . . A . . . . . T . . . . . G . . . . . T . . . . . C . . . . . T . . . . . G . . . . .

G Y K R I Y T P K I R V L H H Q N V A T N Q V Y T N L L E K T  
 Rx1-19F AAGGATACAAGAGAATTTATACACCTAAATTTAGAGTTTTCACCAATCAAATGTTGCAACTAATCAAGTTTACACGAACCTGTTAGAAAAAAC  
 Rx1-19A.2 . . . . . T . . . . . G . . . . . G . . . . . AG . . . . . C . . . . . C . . . . . T . . . . . A . . . . . AG . . . . .  
 Rx1-19A . . . . . T . . . . . G . . . . . G . . . . . AG . . . . . C . . . . . C . . . . . T . . . . . A . . . . . AG . . . . .

*cps19fl*→

L F S N K C N F K S T S Y F L K L M K E N E D V \* M S Y L F L L  
 Rx1-19F TTGTTTCAAATAAATGCAACTTTAAATCCACCAGTTATTTTGAAGTTGATGAAAGAAAACAGGATGTTTAAATGAGTTATTATTATTTTACT  
 Rx1-19A.2 . . . . . C . . . . . G . . . . . T . . . . . . . . . . G . . . . . C . . . . .  
 Rx1-19A . . . . . C . . . . . G . . . . . T . . . . . . . . . . G . . . . . C . . . . .

C L T L F L L T I F Y F F A F I Q D L I A P P V V M S V M F L I  
 Rx1-19F TTGCCTACATTATTCTTATGACTATATCTATTCTTCTTTTATCAAGATTTAATGCTCCTCAGTAGTTATGCTCGTAATGTTTCTAA  
 Rx1-19A.2 C . . . . . G . . . . . C . . . . . A . . . . . T . . . . . C . . . . . T . . . . . C . . . . . G . . . . . A . . . . . G . . . . . T . . . . . C . . . . .  
 Rx1-19A C . . . . . G . . . . . C . . . . . A . . . . . T . . . . . C . . . . . T . . . . . C . . . . . G . . . . . A . . . . . G . . . . . T . . . . . C . . . . .

S S V F A L V N S K N W N I E Y S G I A Y I L I I S G I I I F  
 Rx1-19F TTAGTTTCAAGTATTTGACCTGTTAAATCAAATAAATGGAATATGGAATAGCTATATTTCTATAATAGTGGTATTATTATTATTT  
 Rx1-19A.2 . . . . . T . . . . . C . . . . . TA . . . . . GTGC . . . . . G . . . . . TT . . . . . T . . . . . C . . . . . T . . . . . G . . . . . T . . . . .  
 Rx1-19A . . . . . T . . . . . C . . . . . TA . . . . . GTGC . . . . . G . . . . . TT . . . . . T . . . . . C . . . . . T . . . . . G . . . . . T . . . . .

S I P L M A L K S P N F N T E V K I A D R L I D I Q F W K I A L  
 Rx1-19F TCGATTCTTTAATGCGATTAAATCACCTAATTTAATACTAGGTTAAGATGCTGATCGATTAAATGATATCAATTTGGAAAAATGCTCT  
 Rx1-19A.2 . . . . . A . . . . . G . . . . . C . . . . . CT . . . . . T . . . . . T . . . . . G . . . . . C . . . . . A . . . . . CA . . . . . A . . . . . AG . . . . . GA . . . . . G . . . . . C . . . . .  
 Rx1-19A . . . . . A . . . . . G . . . . . C . . . . . CT . . . . . T . . . . . T . . . . . G . . . . . C . . . . . A . . . . . CA . . . . . A . . . . . AG . . . . . GA . . . . . G . . . . . C . . . . .

T I I I D L F I L Y L Y R K E I Y N L V L S N G Y T G S N I Q W  
 Rx1-19F AACTATTATAATGATCTCTTTATTTTGTATCTTTACAGGAAGGAAATATACAACCTTGTCTTAGTAATGGATATACGGGGTCAAATATTTCAGT  
 Rx1-19A.2 T . . . . . C . . . . . A . . . . . C . . . . . A . . . . . T . . . . . GA . . . . . G . . . . . TC . . . . . T . . . . . CA . . . . . CC . . . . . T . . . . .  
 Rx1-19A T . . . . . G . . . . . C . . . . . A . . . . . C . . . . . T . . . . . GA . . . . . G . . . . . TC . . . . . T . . . . . CA . . . . . CC . . . . . T . . . . .

F F R N A T S Y E G E L T V R T F I R V L I R V I D V S A Y I  
 Rx1-19F GGTTTTTGAAGATGCAACGAGTTATGAAAGTGAATGACAGTGCAGACTTTTATTCGAGTTCTCATTTCGTTTATTGACGTTATCTGCTTATTT  
 Rx1-19A.2 . . . . . C . . . . . T . . . . . C . . . . . GC . . . . . A . . . . . CG . . . . . G . . . . . C . . . . . A . . . . .  
 Rx1-19A . . . . . C . . . . . T . . . . . C . . . . . GC . . . . . A . . . . . CG . . . . . G . . . . . C . . . . . A . . . . .

F G Y T F I N N F L I Y R H K R P K D I L L L V P L L I F I S K  
 Rx1-19F TTGGATATACTTTTATAATAATTTCTTATCTATGCCATAACGCCCTAAAGACATATCTTTTAGTACCTTTATTAATTTATTTTCAA  
 Rx1-19A.2 . . . . . CT . . . . . C . . . . . T . . . . . A . . . . . T . . . . . TT . . . . . C . . . . . G . . . . . T . . . . . T . . . . .  
 Rx1-19A . . . . . CT . . . . . C . . . . . T . . . . . A . . . . . T . . . . . TT . . . . . C . . . . . G . . . . . T . . . . . T . . . . .

T L I S G G R Q D I I K I L I A Y V I M M Y I Q Q K R K V G W N  
 Rx1-19F AACTTTAATATCAGGAGCCGGCAAGATATTTAAATTTCTGATGCGCTATGTAATCATGATGATATCCAACAAAACCGAAAGTTGGATGGA  
 Rx1-19A.2 . . . . . C . . . . . T . . . . . T . . . . . G . . . . . TA . . . . . ATTG . . . . . A . . . . . T . . . . . A . . . . . TG . . . . . A . . . . . GCC . . . . . T . . . . .  
 Rx1-19A . . . . . C . . . . . T . . . . . T . . . . . G . . . . . TA . . . . . ATTG . . . . . A . . . . . T . . . . . A . . . . . TG . . . . . A . . . . . GCC . . . . . T . . . . .

R V I S H K Y I H L G F V G L I A G I P A F Y Y S L F L A G R  
 Rx1-19F ATAGAGTCATATCTCATAAATATATTCACCTTGATTTGTTGGTTAATAGCAGGATTTCCAGCATTTTACTACTTTTGTTTTGGCCGGTGGT  
 Rx1-19A.2 . . . . . AG . . . . . C . . . . . C . . . . . GAGA . . . . . T . . . . . A . . . . . C . . . . . G . . . . . T . . . . . G . . . . . TA . . . . . C . . . . . T . . . . .  
 Rx1-19A . . . . . AG . . . . . C . . . . . C . . . . . GAGA . . . . . T . . . . . A . . . . . C . . . . . G . . . . . T . . . . . G . . . . . TA . . . . . C . . . . . T . . . . .

S T T R T L F E S V S T Y L G G S I Q H F N Q Y I E N P L D P G  
 Rx1-19F TCAACGACTAGGACGCTATTGAGAGTGTTCGACCTATCTAGGAGGCTCAATTCAGCATTTTAAATCAGTATATTGAAAATCCATTAGATCCTGG  
 Rx1-19A.2 . . . . . T . . . . . A . . . . . TG . . . . . A . . . . . A . . . . . T . . . . . T . . . . . G . . . . . C . . . . . TA . . . . . T . . . . .  
 Rx1-19A . . . . . T . . . . . A . . . . . TG . . . . . A . . . . . A . . . . . T . . . . . T . . . . . G . . . . . C . . . . . TA . . . . . T . . . . .

E V F G S E T L V P I L N I L G E M G L V N Y R S T I H L E F R  
 Rx1-19F TGAAGTTTTGGCAGTGAAACATTGGTGCCPATATTAATATATTAGGGGAAATGGGCCTAGTTAATATCGTAGTACAATTCATTAGAAATTC  
 Rx1-19A.2 . . . . . G . . . . . GA . . . . . GT . . . . . T . . . . . AG . . . . . TA . . . . . G . . . . . T . . . . . G . . . . . TA . . . . . TC . . . . . T . . . . .  
 Rx1-19A . . . . . G . . . . . GA . . . . . GT . . . . . T . . . . . AG . . . . . TA . . . . . G . . . . . T . . . . . G . . . . . TA . . . . . TC . . . . . T . . . . .

T L G V T V G N V Y T F F R R P L H D F G L V G M Y V F V F A  
 Rx1-19F GGACACTAGGAGTTACTGTAGGAAATGTTTATACTTTTGTAGAAAGCCCTGCGATGATTTGGTCTAGTTGGTATGATGATTTGCTTTGCT  
 Rx1-19A.2 . . . . . CAGT . . . . . GA . . . . . A . . . . . T . . . . . C . . . . . C . . . . . G . . . . . C . . . . . GT . . . . . C . . . . . A . . . . .  
 Rx1-19A . . . . . CAGT . . . . . GA . . . . . A . . . . . T . . . . . C . . . . . C . . . . . G . . . . . C . . . . . GT . . . . . C . . . . . A . . . . .

V G A F F A I Y Y L V L R K K Q V G F N L D I H T I I Y S Y V F  
 Rx1-19F GTAGGTGCTTTTTGCTATTATTTATTTAGTCTGAGAAAGAAACAGGTTGGTTTTAATTTGGATATTCATACCATTTATTCTTTATGCTT  
 Rx1-19A.2 . . . . . C . . . . . TG . . . . . T . . . . . GAAAT . . . . . AA . . . . . GC . . . . . GT . . . . . C . . . . . A . . . . . AA . . . . . C . . . . .  
 Rx1-19A . . . . . C . . . . . TG . . . . . T . . . . . GAAAT . . . . . AA . . . . . GC . . . . . GT . . . . . C . . . . . A . . . . . AA . . . . . C . . . . .

Y W I F L S S I E Q Y S F T M I S L Y T L V F I V L V Y F M A I  
 Rx1-19F TTATTGGATTTTTTATCATCAATCGAGCAATCTCGTTCAATGATTAGTCTATATACACTTGTATTTATTGCTTGGTTTACTTTATGCTTA  
 Rx1-19A.2 . . . . . G . . . . . C . . . . . T . . . . . T . . . . . CA . . . . . C . . . . . T . . . . . G . . . . . T . . . . . A . . . . .  
 Rx1-19A . . . . . G . . . . . C . . . . . T . . . . . T . . . . . CA . . . . . C . . . . . T . . . . . G . . . . . T . . . . . A . . . . .

F Y W C T D F K R G K L I F K I S D S S I K L K E E \*  
 Rx1-19F TC'TTTTACTGGTACAGATTTTAAAGAGGAAAAGTATTTTAAATTTCTGACTCAAGTATCAAATTAAGAAAGAAATAACAGAATGTATAG  
 Rx1-19A.2 . . . . . T . . . . . AA . . . . . TTG . . . . . CC . . . . . C . . . . . AG . . . . . AA . . . . . C . . . . . C . . . . . CGG . . . . .  
 Rx1-19A . . . . . T . . . . . AA . . . . . TTG . . . . . CC . . . . . C . . . . . AG . . . . . AA . . . . . C . . . . . C . . . . . CGG . . . . .

*cps19fj*→

M N T K I K N I I T S F S Y V I S S N L L I V L T S S L  
 Rx1-19F -GAG-AGGGTAGATGAATACTAAAATAAAAATAAATACTAGTTTTCTTATGTTATTTCTTCAAATCTGCTCATAGTTTAAACCTCATCACT  
 Rx1-19A.2 A..G.TA..T.....GC.....C.....A...C.....C.....TT.G.C...C.G...T..CT.  
 Rx1-19A A..G.TA..T.....GC.....C.....A...C.....C.....TT.G.C...C.G...T..CT.

V V L I V P K I M G V T E Y S Y W Q L Y I F Y L T Y I G F F H L  
 Rx1-19F AGFTGTTTGGATTGTTCTAAAATAATGGGGTAACTGAGTACAGTTACTGGCACTTTATATTTTTATCTGACCTATATCGGTTTTTCCACT  
 Rx1-19A.2 G.....C..A...T.C..C...T.....A.....  
 Rx1-19A G.....C..A...T.C..C...T.....A.....



G W I D G I Y L K Y G G L E Y T N L D R K Q F Y S Q M I L F S  
 Rx1-19F TGGTGTGGATTGATGGGATTTATCTCAAATATGGTGGCTTAGAATATACAAATTTAGATAGAAAACAGTTTTTATCTCAGATGATTTCTATTTTCT  
 Rx1-19A.2 .....T.....AT.....  
 Rx1-19A .....T.....AT.....

S F L M L I S L V L F T L N L I T V R D E N A R Y I Y N M A I I  
 Rx1-19F AGTTTCTTAATGCTAATCTCGCTGGTATTATTTACTTTGAACCTAATAACTGTAAGGGATGAAAACGCAAGATATATTTATAATATGGCTATCAT  
 Rx1-19A.2 .....C.....A.T.....CT..C...G...GG...A..TG..T..G.....CA...TCG..T...C.....A...T...  
 Rx1-19A .....C.....A.T.....CT..C...G...GG...A..TG..T..G.....CA...TCG..T...C.....A...T...

S M T V T N L R T L Y V Y I L Q M T N R L K D S S V I L I S D R  
 Rx1-19F CAGCATGACAGTCACAACTTAAGAACACTCTATGTTTATATCTTGCAGATGACAAATCGCTTGAAGGATAGTTTCAGTCATTCTAATTAGTGATC  
 Rx1-19A.2 .....A..T...T.A.....T.....G.....C.G...G.T.....C.....CA...C.C.....  
 Rx1-19A .....A..T...T.A.....T.....G.....C.G...G.T.....C.....CA...C.C.....

V L Y V L L L F M F I V F G W H E Y K V M I W A D I L G R T F  
 Rx1-19F GCGTTTATATGCTACTCTTTTATTCATGTTTATTTGGATGGCATGAGTACAAGGTATGATTGGGCTGATATTCTAGGTCGAACATTT  
 Rx1-19A.2 .....G.....C.....T.....A.....T.....T...A.....T.....  
 Rx1-19A .....G.....C.....T.....A.....T.....T...A.....T.....

S L M L S F W I C K D I V F Q P L S K F I L D F K E S L D N I R  
 Rx1-19F TCTCTCATGCTTCTCTCGATTTGTAAGATATGTTGTTTCAGCCTTTGTCAAATTTATCTTAGATTTCAAGGAGTCCCTTGATAATATCCG  
 Rx1-19A.2 .....A.....C.....TA..A.A.A...TT.....  
 Rx1-19A .....A.....C.....TA..A.A.A...TT.....

V G I N L M L S N I A S S M I I G I V R M G I Q W N W N I E T F  
 Rx1-19F TGTTGGAATCAATTTAATGCTATCTAACATTTGCGAGTAGCATGATTATAGGCATTTGTTGATGGGAATTCATGGAATTTGGAATATGAAACAT  
 Rx1-19A.2 A.....T.....T.....C.....C.A.....TT.....C..T.....C.A.....C.....C.....C.....  
 Rx1-19A A.....T.....T.....C.....C.A.....TT.....C..T.....C.A.....C.....C.....C.....

G K V S L T L S I S N L L M T F I N A I G L V I F P L I K R T  
 Rx1-19F TCGGAAAAGTGTCAATTAACCTTTGAGTATATCTAATTTATTAATGACTTTTATTAATGCCATTTGATTTAGTTATCTTTCTTTGATAAAGCGAACA  
 Rx1-19A.2 .....A..TC.G..C.....G.T..C...C...G.....C.....C.....CG.....GC..AT..GA.....  
 Rx1-19A .....A..TC.G..C.....G.T..C...C...G.....C.....C.....CG.....GC..AT..GA.....

K T E N L P K I Y S N L R N A L M L V M F A I L L F Y Y P L K F  
 Rx1-19F AAGACTGAGAATTTACCTAAAATTTATTTCAATTTAAGAAATGCTTTGATGTTGGTTATGTTGCGCAATCTTGCTCTTCTACTACTCTTTAAAAT  
 Rx1-19A.2 .....AG.A..A..C.....G.....C.T.....C.TA.C.....G..T.....A.T.....A.....  
 Rx1-19A .....AG.A..A..C.....G.....C.T.....C.TA.C.....G..T.....A.T.....A.....

I L D I W L P A Y K D A L V F M A L I F P M S V Y E G K M A L V  
 Rx1-19F TATTTCTGATATTGGCTTCCCTGCTTATAAGGATGCCCTAGTTTTCATGGCCCTAATTTTCTATGTCAGTTTATGAAGGGAAAATGGCTTTGG  
 Rx1-19A.2 .....G..AT..A..CC.C...T.G..A..C...C.A.....GA.....A...T.....C.....C.....A.....  
 Rx1-19A .....G..AT..A..CC.C...T.G..A..C...C.A.....GA.....A...T.....C.....C.....A.....

I N T Y L K A M R M E K D I L K I N A L V M L T S I V V T L V  
 Rx1-19F TGATAAATACATATTTAAAGCAATGAGAATGGAAAAGACATTTCTCAAATTAATGCTTTGGTTATGTTAACTAGTATAGTAGTACATTAGTG  
 Rx1-19A.2 .....C..T.....T..C.....G...T.A.....G...T..C..T.G.G.C.....C..A..A...G.TC...G.G...CC.G..T  
 Rx1-19A .....C..T.....T..C.....G...T.A.....G...T..C..T.G.G.C.....C..A..A...G.TC...G.G...CC.G..T

T F L L L N N L G L T V V S I V I L L A L R S I I A E L I L S K  
 Rx1-19F ACTACTCTACTACTAAATAATTTGGGGCTGACAGTTGTATCTATAGTTATTTTACTTGCTTTAAGAAGTATAATAGCTGAATTAATTTTATCCAA  
 Rx1-19A.2 .....A..ATAC..GT.....G...A.....T.....G.....G...G..A...G.....C.....C.....C.G..T..  
 Rx1-19A .....A..ATAC..GT.....G...A.....T.....G.....G...G..A...G.....C.....C.....C.G..T..

K L K I S V K Q D I A L E L L M T I I F I S S S W Y L S I W I A  
 Rx1-19F AAAACTGAAGATATCAGTCAAGCAAGACATGCTTTAGAGTTACTTATGACGATTTATATTTTCTTCAAGTGGTATCTCTATTTGGATTG  
 Rx1-19A.2 .....G..TG.T..G..T...A.G..T...TA...A..T...T...C..G.C.....T.GC.G...G..C.G..  
 Rx1-19A .....G..TG.T..G..T...A.G..T...TA...A..T...T...C..G.C.....T.GC.G...G..C.G..

V I I Y L L A Y T L Y L Y L K H K D I R M Y I E Y F K N H K K  
 Rx1-19F CAGTAATAATTTATTTATTTGGCGTATACCTTATATTTGATTTAAGACACAAAGATATCAGAATGTATAGAATACTTTAAAATCATAAAAA  
 Rx1-19A.2 .....TG.C..CAC.A.A..C..CGG.....C.C.....G..A..G.....A..CC...T...CTT.....GC.AG.....  
 Rx1-19A .....TG.C..CAC.A.A..C..CGG.....C.C.....G..A..G.....A..CC...T...CTT.....GC.AG.....

I S \*  
 Rx1-19F ATATCATAAAAA  
 Rx1-19A.2 .....  
 Rx1-19A .....A..TT..

# APPENDIX VI

## The nucleotide and deduced amino acid sequence of the central (type-specific) portion of *cps19c* locus.

The nucleotide sequence is numbered in accordance with GenBank accession number AF105116 and starts from the distal portion of *cps19cR* to the beginning of *cps19cL*. The amino acid translation for each ORF is represented by single letter code above the first nucleotide of each codon. Possible ribosome binding sites are underlined.

... *cps19cR*→  
A Y K K M E F L P F E Q K N L W L M M S L K K E C I Y I P K L Q S  
1 AAGCTTATAAGAAAATGGAGTTTCTCCCTTTGAACAAAAAAATCTGTGGCTAATGATGTCGCTAAAAAAGAATGTATTTATACCTAAGTTACAATC  
F S E S L K E K V V E N L G E I V T L K T T K K S I K T L L R K N  
101 ATTTTCTGAATCTTTAAAAGAAAAGTGTGAAAAATTTGGGTGAGATAGTAACCTCTTAAGACTACGAAAAATCTATAAAAAACGTTATTAAGAAAAGAT  
S D L L W V D F V Y R L F G G V N F L F I P E L Y T R V I F L L K Y  
201 AGCGATTTACTTTGGGTTGATTCGTTTATAGATTGTTGGTGGTCAATTTCCCTATTTATCCAGAATTATATACAAGAGTTATATTTTATTGAAAT  
L D L K L R L K N Y G E \*

*cps19cJ*→  
M G N K S I K L N A L L N I V L T L S N I I F P  
301 ATTTAGATTTGAAATGAGGTTGAAAAATATGGGGAATAAATCCATAAAGTTGAATGCATTATTAATATGTCCTGACGCTATCAATATCATTTTCC  
L I T F P Y I S R I L N P N G I G L T S F F S S I G N Y G I L L A  
401 CATTAATCACTTTCCCTTATATATCTAGAATATTGAATCCAAATGGTATAGGTTAACTTCAATTTTTAGTTCAATAGGAATTAATGGTATTTACTTGC  
S L G I S T Y G I K A V A S V R D D R D K L S K V V Q E L M I I N  
501 TTCTCTGGGAATTTCAACTTACGGTATCAAAGCAGTAGCAAGTGTAGAGATGATAGAGATAAGTTGTCAAAGTAGTACAGGAGTAAATGATTATAAAC  
V A M S I I T T A I L L F M T I F I T Q L N R E F S L L L I T C G G T  
601 GTTGCTATGCTATAATAACAACCTGCAATACTATTATTTAGACTATATTTATAACACAATGAATAGAGAATTTTCACTCCATTTGATCACATGTGGGA  
I L S S P F A L N W L Y S G M E E Y T Y I T T R S V V F T I L S L  
701 CTATTTTATCTTCCCTTTCCGCTTAAATTTGGTTGTATAGTGGAAATGGAAGAATATACGTATATTACTACTAGGTCAGTGTGTTTACAATTTCTATCATT  
I L I F L L V K R P E D Y I V F A S I S L F S S L S S N I L N L W  
801 AATATTGATTTTCTACTTGTGAAAAGGCCAGAGGATATATTTGCTTAGTATTTCATTTGTTTTCTTCTCAAGTTCAAATATCTTAAATCTATGG  
H S R H F I N I K L Y K N L Q F K Y H F K P M W Y L F A S L L A V N  
901 CATAGCCGACATTTCAATTAATTAATTAATAAAAAATTTACAATTTAAATATCATTTTAAACCAATGTGGTATTTATTTGCTCATTTACTTGCAGTAA  
I Y T N L D T V M L G F I N G N E A V G Y Y S V A S K V K W I L L  
1001 ATATTTATACTAATTTTAGATACAGTGTGCTCGGTTTTATTAATGGTAATGAGGCTGTGGGATACTATTCTGTGGCATCAAAGGTTAAGTGGATTTTGGT  
S L I T S I S A V L L L P R L S F Y I S K N D T S N F I K M L K E S  
1101 TTCTCTTATACATCTATTAGTGCAGTTTGGCTACCGAGACTTTCATTTTATATTGTAATAAATGACACCTCGAATTTTAAAAATGTTAAAGGAGTCA  
S A V I F F I A I P L M V F F I V E A K D S I L L L G G S Q Y L P A  
1201 TCTGCGGTTTATATTTTATTTGCGATTTCCATTTGATGGTATTCTTTATTTGTAGAGGCCAAAGATAGTATCTTATTACTAGGAGGAAGTCAATCTTCTCT  
T L A M Q I L M P I L L I S G F S N I T G N Q I L I P M N R E K Y  
1301 CGACTTTAGCGATGCAAATACTTATGCCAATTTTACTTATTTCTGGTTTCTCGAATATTACAGGAAATCAAATATTGATTTCCAAATGAATAGAGAAAAA  
F M V A V T I G A V I N L I L N L L L M P K F G I I G A S V A T L  
1401 TTTTATGGTTGACGTAACGATTTGGTCTGTGATTAATCTTATTTTGAATCTACTGTTAATGCCTAAGTTTGAATATTGTTGCTTCTGTCGCAACTCTT  
F A E L M E Q M T V Q L H F S K E Y L V S N I S I K S L V N V I I A T  
1501 TTTGGGAATTTGTCGAGATGACCGGTACAATTACATTTTCAAAGAATATTTAGTATCAAATATATCGATAAAGAGTTGGTTAATGTGATAATTGCAA  
V V S T I P L I I L N Q L I T I T I P F Y S L M L A G F A F F S L  
1601 CAGTTGTTTCTACAATACCCTAATCATTTTGAATCAGCTGATACGATAACGATAACCATTTTATTTCTTAATGCTAGCAGGTTTGGCTTTCTTTTCATT  
Y L V I L L L K E E V T I Q L F S L L A K K K \*  
1701 ATATTTAGTAATTTCTGCTTTTATTAAGGAGGAAGTGACGATTCATTTATTTTCTTCTTTCGCAAGAAGAAGTAAATTTGGTTAGAAATTTGAATGTATA

*cps19cK*→  
M K I M L V F G T R P E A I K M C P L V N  
1801 AACAAATAAAGAAATTTAATTTATTTAGGAGGAAATCATGAAGATAATGCTAGTTTGGTACAGTCCAGAACGCAATAAAAATGTCTCCATTTAGTGA  
E L K K Q A D M E T V V C V T G Q H K E M V S P V L E L F G V Q P  
1901 ATGAGTTGAAAAACAGGCAGATATGGAAACAGTTGTTTGTGTAACCTGGTCAACACAAGGAGATGGTTAGCCCTGTTTGGAAATTTGTTGGAGTTCAACC  
D Y D L E I M K A N Q T L F S I T T S I L E K I K P V L E E E Q P  
2001 AGACTATGATTTAGAAATAATGAAAGCTAATCAAACCTGTTCTCTATAACAACAGTATTTTAGAAAAAATTAACCTGTTTTAGAAGAAGAACCAACC

D I V L V H G D T T T T Y A A A L A A F Y L G I K V G H V E A G L R  
2101 GATATTGTTCTAGTTCATGGTGACACTACTACAACATATGCAGCAGCTTTGGCAGCATTATCTTGGAAATTAAGTTGGCCATGTTGAAGCTGGGTTAC  
T Y N L Q S P F P E E F N R Q S T S I I A N Y H F A P T E L A K E  
2201 GAACTTACAACCTGCAAGTCCCTTCCCAGAAGAATTTAATAGACAATCGACTTCTATATATGCAAAATATCATTTTGGCTCCAACCTGAATTTGGCTAAAGA  
N L L K E G R E N V Y V T G N T V I D A L T T T V Q K D Y T H P D  
2301 AAATCTCTTAAAGAGGGTAGAGAGAATGTTTATGTAAGTAAATACAGTTATTTGATGCATTTACAACCTACAGTACAAAAGGATATACACACCCCTGAT  
L D L N V N T R L I L L T A H R R E N L G E P M K H M F R A V K R V  
2401 TTAGATTAAACGTTAAACTCGCTTATCTACTGACTGCATAGACGTAAGAAATCTCGGAGAACCTATGAAACACATGTTTAGAGCTGTTAAACGAG  
L N E Y D D V K V I Y P I H K N P L V R E T A T E I F G D T E R I  
2501 TCTTAAATGAATGACGATGTTAAGGTAATTTTCCAATTCATAAGAATCCCTTGGTGCCTGAAACAGCTACAGAAATTTTGGAGACACAGAAGCTAT  
Q I I E P L D V L D F H N F M N H S Y M I L T D S G G V Q E E A P  
2601 TCAGATTATTGAACCTTTGGATGTTCTGATTTTCATAATTCATGAATCATAGTTACATGATTTTAACTGATTCAGGAGGAGTTTTCAGGAAGAGGCTCCT  
S L G K P V L V M R D T T E R P E G V A A G T L K L V G T D E E T I  
2701 TCTTTAGAAAACCTGTATTTGGTCTGCGAGATACGACAGAAAGACCTGAAGGAGTAGCTGCCGGAACCTTGAATTTGGTTGGAACCTGATGAGGAGACTA  
Y Q N F K M L L D D S E E Y K K M S Q A S N P Y G K G G D A S K Q I  
2801 TTTATCAAACTTTAAGATGCTTTTAGACGATTCGGAAGAATATAAAAAATGAGTCAAGCTAGTAACTCTTATGGAAAAGGTGATGCTAATCAACAGAT  
V H I L S G I \*  
2901 TGTTCAATTTTAAAGCGAATTTAAGCGAGGCCAAATAAAGTAATAAAAAACACTATCTTATAAAAGGATTGATCTTGTAGTTGATTCGGGAACATGTT  
3001 CATATCCTTTGACTCGAGCTACATAAAAAATAACTTTGTCTGATTTATGATAGATCGATAATTTACTACTACTTTGGACATGATGTTAGCAGTTATTAGG  
3101 GATGTTTTGATTAACCTCAACTTTTCAGGATTCGCCTTGCTTCTAGATTTTCTTCTGATTTTATCATTACTTCAGTAAGTTAAATCGTCTATTTACTAG

*cps19cS*→

L K I V I P R I I  
3201 AATTAACACTTTTAAAAATCCATTGTTAGTCTCATGTTTTAGATATAGGTAACATAAAATTTAGGAGGTTGTTAGTTTGAAGATTGTAATTCGAAGAATTA  
H N K E Q L T W D W S G T I T N I K K F L G K Y E I V E E Q N I F  
3301 TTCATAACAAGAACAACCTGACCTGGGATGGTCCGGGACAATACTAATAAAAAATTTTAGGAAATACGAGATTTGTTGAGCAACAGAATATTTT  
Y T F R M N V H K V L V R L G I K K S D M S M T Y I K Y A E N Q V  
3401 CTATACTTTTGAATGAATGTCACAAAGTCTTGTTCGTTTAGTATTAATAAAATCTGATATGAGCATGACGTATATAAAATATGCTGAAAACTCAAGTT  
H L S P E D V C L T F D E F P L S F P D N P V Y I Y L Q D L N L H Y L  
3501 CATCTATCGCCAGAGGATGTTTGTCTCACGTTTGTGATGATTTCCTTTATCTTCTCTGATAATCCAGTTTATATCTATCAAGACTTAAATCTTCATTATT  
I E S S Q N N S Q S F K Y S G F Q N V P A D I L D R R M R K Q E I  
3601 TGATAGAGAGTTCTCAAAACAATAGTCAATCGTTCAAATATAGTGGTTTTCAAAACGTCCTCCGGCTGATATTCAGATAGACGAATGAGAAAAACAGGAAAT  
F Y N Q A T G I F T M S K W F S D Y L I A Q Q G L P V E K V H Y V  
3701 ATTTTATAATCAAGCTACTGGAATATTTACTATGAGTAAATGGTTTTTCAGATTACTTTGATAGCTCAACAAGGACTTCCAGTTGAAAAGTTTCATTATGTG  
G A G T N M N N L F L D H S H K E R N K F L F I G K D F F R K G G D  
3801 GGGCAGGAACAATAATGAATAATCTTTCTGACCACTCTCATAAGGAACGTAATAAGTTTTTATTTATTGGTAAAGATTTTTTCTGTAAGGAGGAG  
L V Y N A F V Y L Q N N L M P E A E L Y I I G P S D V P M E F N N  
3901 ATCTTGTATTAATGCTTTTGTCTATTTGCAAAATAATCTCATGCCGAGGAGAGTTGATACATTATAGTCTTCCAGATGTTCCGATGGAATTTAACAA  
P N V Y F L G N L S A D K V Q Y F Y N L C D V F V L P S R F E A F  
4001 TCCGAATGTTTATTTTATTTAGGTAATCTATCAGCCGATAAGGTGCAATATTTTATAATCTTTGTGATGATTTGTTTACCTTCCCGATTTGAGGCAATTT  
G I V F V E A L C Y G L P C I G R D L M E M P N L I Q N N E T G L L  
4101 GGAATGTTATTTGTTGAAGCTTTGCTATGTTTACCATGATCGGTGATTTAATGGAATGCAAAACCTAATTCAAAATAATGAAACTGGATTAT  
L P T E E E N P Q V L A D V M Y N L I K D E S F F K N V Q T K Q D  
4201 TATTACCTACTGAAGAGGAAAATCCACAGGTTTTAGCTGACGTAATGTATAATTTGATAAAAGATGAAAGCTTTTTTAAAAATGTTTCAGACTAAACAAGA  
Y Y K A E Y S W D T V A K R M I S I M K Q D M N N N L \*  
4301 TTATTTAATAAGCAGAAATATTCGTGGGACACAGTTGCCAAAAGAAATGATTTCAATATGAAGCAAGATATGAACAACATCTATAAGTGAATAAATAAAG  
4401 ATATTTAGTCTAACCGTCTATTCTACTAACAGTAATACTGAAAGAAATTTAAAATACAAGTTATACGTATAGTGTCTTCCCTAATTTCTGACATGTT  
4501 GAAAAGATTTTAAAACTAATTTAAAAGATTAAGGAAGTAAATTCAGCTTATTGAATAAAATTAATTTAGTCTAGTTATTTCTTGTCTGATCTTAT  
4601 TCAATTTTCCGAACTTTGTGAGTTGTTGCAACTACCAGAATACACAAAGAACAGAAATCTCTTAGCATCATGATGAAACAATAATCTGCATAAATTT  
4701 CAGATATAACTAGTGGTGCAGAAAGTGGATCTTTGAAAATAAGAAATCTTTGTAGAGCTATTTCACTAGCTGGTGAAGAATTTTTTCAATTTTTCGTA  
4801 CAGATTACAATAATGTTATGTTGGGTTTTAAATTTCTTGTGTTGGTGTGTTGATGAAATAGAGATGCACGTTATCGAAATTCAGATAGTGGTATATCGCT  
4901 TGCTTAATCTCTTCCAGTATTTGTAGTAGAATTAAGTGTCTTGTATTAATAAAGAACAACACTTTTTTATTCAGTGTGTTGATGTTGGTGAACCGA

*cps19cL*→

M K G I I L A G G S G T R L Y P L T R A A S K Q L M P V  
5001 AAGGAACGATGTTACTTTATGAAAGGATTATTCTCGCGGGTGGTTCGGGGACAGTTTTATATCCTTTGACTCGAGCTGCATCAAAGCAACTGATGCCGGT  
Y D K P M I Y Y P L S T L M L A G I R D I L I I S T P Q D L P R F  
5101 TTATGATAAACCGATGATTTACTACCCACTTTCAACTTTGATGTTGGCTGGGATTTAGGATATTTGATTATCTCAACTCCTCAAGATTTGCCCTCGTTTT  
K E L L Q D G S E F G I Q L S Y A E Q P S P D G L D  
5201 AAAGAGCTCCTCAAGATGGCTCTGAGTTGGGATTCAAATGTTCTTATGCAGAGCAACCAAGTCCAGATGTTTGGATCC

# APPENDIX VII

## Publications

**Morona, JK., Morona, R. and Paton, JC. 1997.** Characterization of the locus encoding the *Streptococcus pneumoniae* type 19F capsular polysaccharide biosynthetic pathway. *Mol. Microbiol.* 23:751-763.

**Morona, JK., Morona, R. and Paton, JC. 1997.** Molecular and genetic characterization of the capsule biosynthesis locus of *Streptococcus pneumoniae* type 19B. *J. Bacteriol.* 179:4953-4958.

**Morona, JK., Morona, R. and Paton, JC. 1999.** Analysis of capsule loci from various *Streptococcus pneumoniae* serotypes using long-range PCR identifies two classes of *cpsC*. *Microbiol.* (manuscript submitted).

**Morona, JK., Morona, R. and Paton, JC. 1999.** Genetics of capsular polysaccharide biosynthesis in *Streptococcus pneumoniae* types belonging to serogroup 19. (manuscript in preparation).

**Paton, JC. and Morona, JK. 1999.** *Streptococcus pneumoniae* capsular polysaccharide. In: Fischetti, V., Novick, R., Ferretti, J., Portnoy, D. and Rood, J. (eds.). Gram-positive pathogens. ASM Press, Washington DC, USA. (in press).