

# **BGP, Not As Easy As 1-2-3**

Ashley Flavel

*Thesis submitted for the degree of*

*Doctor of Philosophy*

*in*

*Applied Mathematics*

*at*

*The University of Adelaide*

Faculty of Engineering, Computer and Mathematical Sciences



October 6, 2009

# Signed Statement

This work 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 has been made in the text.

I consent to this copy of my thesis, when deposited in the University Library, being available for loan and photocopying, subject to the provisions of the Copyright Act 1968.

SIGNED: ..... DATE: .....

# Acknowledgements

I would firstly like to acknowledge the support of my primary supervisor, Associate Professor Matthew Roughan. It was he who first inspired me to pursue a researching career. In addition, his financial and social assistance helped gain me an internship at AT&T Research that reinvigorated my desire to work on real-world problems with real-world solutions.

I would also like to acknowledge the support of my co-supervisor, Professor Nigel Bean. His ability to interpret my sometimes incoherent description of a problem into one that was simple and easy-to-understand was a major factor in the success of this thesis.

Dr. Olaf Maennel, as fellow co-supervisor, provided an alternative perspective to inter-domain routing to my primary supervisor A/Prof Roughan. His pragmatic approach to inter-domain routing and in-depth knowledge of its quirks and corner cases was an excellent source to verify the sanity of my ideas. All three of my supervisors provided contrasting perspectives which although at times was frustrating allowed the development of ideas that were not only theoretically sound but pragmatic. I would like to thank them all for their friendship and support over the last few years.

Dr. Aman Shaikh of AT&T Research Labs was a major contributor to the ideas in this thesis. The four months I spent at AT&T Research and the subsequent year-long collaboration resulted in the best work of my thesis. Aman also assisted with my transport to and from the laboratories in which time many of the most important research breakthroughs came (as well as many discussions on cricket). AT&T also deserve acknowledgement for providing me access to their

commercially sensitive network data.

This thesis was reviewed by three of my most admired networking researchers:- Dr. Timothy Griffin of The University of Cambridge, Professor Jennifer Rexford of Princeton University and Dr. Steve Uhlig of Technische Universität Berlin. I am sincerely privileged and greatly appreciate their time and useful comments.

For financial assistance I am very grateful to the Australian Research Council's Communications Research Network (ACoRN). ACoRN provided me with substantial financial assistance for my research visit to AT&T. I specifically would like to thank ACoRN's Adelaide University representative Belinda Chiera. In addition, I would like to acknowledge the financial support of the Australian Research Council through grant DP0557066.

Dr. Flo Rice, although not related to my technical research, was a pivotal figure in its development. She provided me with accommodation for the time I spent at AT&T and became a good friend.

The staff and students at the Teletraffic Research Centre at the University of Adelaide provided me with an excellent research atmosphere. I would specifically like to thank Dr. Jeremy McMahon for frequent useful discussions.

I would like to thank Maxine and Mark Wong See for proof reading several sections of this thesis. A special acknowledgement is reserved for Carrie Kelly who devoted a significant amount of time and energy to proof read the entire thesis.

Lastly, I would like to thank my family for their support throughout my time as a student. Without them my chance to reach this point would not be possible. They have enabled me to reach a level of education that will give me incredible future opportunities.

# Dedication

I dedicate this thesis to my family. My desire to succeed comes from them, and without their love and support this thesis would not have been possible.

# Contents

<b>Abstract</b>	<b>xix</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Preliminary Background . . . . .	8
1.2 Thesis Roadmap . . . . .	9
1.3 Statement of Research Contributions . . . . .	10
1.3.1 Publications Arising From This Thesis . . . . .	12
<b>2 Background</b>	<b>15</b>
2.1 General Routing Protocols . . . . .	17
2.1.1 Link-State Protocols . . . . .	17
2.1.2 Distance-Vector Protocols . . . . .	18
2.2 Routing in the Internet . . . . .	18
2.3 Border Gateway Protocol . . . . .	20
2.3.1 BGP Decision Process . . . . .	22
2.3.2 BGP Operation . . . . .	25
2.3.3 Internal vs External BGP . . . . .	26
2.4 BGP, Not as Easy as 1-2-3? . . . . .	30
2.4.1 The Stable Paths Problem . . . . .	31
2.4.2 MED Oscillation . . . . .	34
2.4.3 iBGP Oscillation . . . . .	37
2.4.4 BGP in the Wild . . . . .	39
2.4.5 BGP Data . . . . .	40

<b>3</b>	<b>Where's Waldo? Practical Searches for Stability in iBGP</b>	<b>43</b>
3.1	Introduction . . . . .	43
3.2	Related Work . . . . .	46
3.3	Background . . . . .	47
3.3.1	iBGP Recap . . . . .	47
3.3.2	Best Path Selection . . . . .	48
3.3.3	Interior Gateway Protocol . . . . .	49
3.3.4	Physical Graph . . . . .	50
3.3.5	Signaling Graph . . . . .	50
3.3.6	Egress Instance . . . . .	50
3.4	Stability . . . . .	51
3.4.1	Complexity of Determining Signaling Correctness . . . . .	52
3.5	Router Reliance Graph . . . . .	52
3.5.1	Reliance Rules for a Route Reflector Topology . . . . .	53
3.5.2	Co-reliance Groups . . . . .	54
3.6	Where Can An Oscillation Occur? . . . . .	57
3.7	Algebraic Description of Co-reliance Groups . . . . .	60
3.7.1	Reducing the Size of Co-reliance Groups . . . . .	62
3.7.2	Oscillation Detection . . . . .	64
3.7.3	Oscillation Classes . . . . .	64
3.7.4	Reliances between Co-reliance Groups . . . . .	70
3.8	Oldest-Route Tie-breaker . . . . .	74
3.9	Prioritizing Egress Instances . . . . .	80
3.9.1	Proving the Stability of an Egress Instance . . . . .	81
3.9.2	Proving the Stability of a Configuration . . . . .	82
3.9.3	Checking the Stability of the Current Network . . . . .	83
3.9.4	Checking the Stability of the Current Network with Limited Measurement Infrastructure . . . . .	85
3.9.5	Practical Implementation . . . . .	86
3.9.6	Online Tool . . . . .	88

3.10	Preventing BGP Oscillation . . . . .	94
3.11	Three-Or-More-Level Route-Reflector Hierarchies . . . . .	94
3.11.1	Greater than Three-Level Hierarchies . . . . .	104
3.12	Discussion . . . . .	106
<b>4</b>	<b>Humpty Dumpty: Putting iBGP Back Together Again</b>	<b>109</b>
4.1	Introduction . . . . .	110
4.2	Related Work . . . . .	111
4.3	Two-Level Route-Reflector Reliance Graph . . . . .	114
4.4	General Route-Reflector Reliance Graph . . . . .	116
4.4.1	Notation Recap . . . . .	117
4.4.2	Reliance Rules for Route Reflection . . . . .	119
4.5	Finding the Actual Solution . . . . .	124
4.5.1	Ordering of Routers Within a Co-reliance Group . . . . .	125
4.5.2	Breaking Ties . . . . .	130
4.5.3	Dynamic IGP . . . . .	130
4.6	Evaluation . . . . .	131
4.7	Generalized Topologies . . . . .	133
4.7.1	Route-Reflection with MED . . . . .	133
4.7.2	Full Mesh . . . . .	136
4.7.3	Confederations . . . . .	137
4.8	Discussion . . . . .	139
<b>5</b>	<b>Peer Dragnet: Analysis of BGP Peering Policies</b>	<b>141</b>
5.1	Introduction . . . . .	141
5.2	Background . . . . .	145
5.3	Related Work . . . . .	147
5.4	Data Collection . . . . .	148
5.4.1	BGP Routes . . . . .	149
5.4.2	IGP Distance Information . . . . .	149
5.4.3	iBGP Topology Information . . . . .	150

5.4.4	Aggregate Traffic Data . . . . .	150
5.5	Analysis of Peering Policies . . . . .	150
5.5.1	Policy Implementation Techniques . . . . .	151
5.5.2	How Peering Links are Used . . . . .	154
5.6	Impact on Routing and Traffic . . . . .	157
5.6.1	Dealing with Routing Dynamics . . . . .	158
5.6.2	Routing Impact . . . . .	159
5.6.3	Traffic Impact . . . . .	167
5.7	Dynamics of Peering Policies . . . . .	172
5.7.1	Policy Change Detection Algorithm . . . . .	173
5.8	Operational Peer Dragnet . . . . .	177
5.9	Non-Canonical Policy Mitigation . . . . .	184
5.9.1	AS-wide BGP Route Controller . . . . .	184
5.9.2	Import Policies . . . . .	186
5.9.3	Distributed Knowledge . . . . .	186
5.10	Discussion . . . . .	190
<b>6</b>	<b>CleanBGP: Verifying the Consistency of BGP Data</b>	<b>193</b>
6.1	Introduction . . . . .	193
6.2	Data Consistency . . . . .	196
6.3	Measurement Artifacts . . . . .	197
6.3.1	Session Failures and Resets . . . . .	197
6.3.2	Incomplete Tables . . . . .	198
6.3.3	Missing Updates . . . . .	198
6.3.4	Update Ordering . . . . .	198
6.3.5	Non-atomic Table Dumps . . . . .	199
6.3.6	Other Artifacts . . . . .	199
6.4	Characterization of Artifacts . . . . .	202
6.4.1	Table Comparison . . . . .	202
6.4.2	Oldest Prefix . . . . .	204

6.4.3	State Information . . . . .	205
6.4.4	Downtime . . . . .	205
6.4.5	Session Re-establishment . . . . .	206
6.4.6	Detecting Measurement Artifacts . . . . .	207
6.5	Extended Measurement Artifacts . . . . .	208
6.6	Cleaning Data . . . . .	210
6.6.1	Session Failures/Re-establishments . . . . .	211
6.6.2	Incomplete Tables . . . . .	212
6.6.3	Missing Updates . . . . .	212
6.6.4	Update Ordering . . . . .	213
6.7	Default Parameter Selection . . . . .	213
6.7.1	Sliding Window Length . . . . .	213
6.7.2	Re-establishment Phase Thresholds . . . . .	213
6.7.3	Downtime Threshold and Bin Length . . . . .	214
6.7.4	Suspicious Bin Thresholds . . . . .	215
6.8	Automated Parameter Selection . . . . .	216
6.8.1	Sliding Window Thresholds . . . . .	216
6.8.2	Suspicious Bin Thresholds . . . . .	220
6.8.3	Discussion . . . . .	225
6.9	Results . . . . .	228
6.10	Discussion . . . . .	232
<b>7</b>	<b>Conclusion</b>	<b>233</b>
	<b>Acronyms</b>	<b>235</b>
	<b>Bibliography</b>	<b>237</b>

# List of Figures

2.2.1	Routing protocol domains . . . . .	19
2.3.1	Internal router structure . . . . .	21
2.3.2	BGP decision process . . . . .	23
2.3.3	Example route-reflector topology . . . . .	27
2.3.4	Route-reflection obscures route availability . . . . .	29
2.4.1	Griffin <i>et al.</i> 's Good Gadget . . . . .	31
2.4.2	Non-stable gadgets defined by Griffin <i>et al.</i> . . . . .	31
2.4.3	MED Oscillation . . . . .	35
2.4.4	iBGP persistent oscillation . . . . .	37
3.3.1	Edge types in the iBGP signaling graph . . . . .	51
3.5.1	An example egress instance . . . . .	56
3.7.1	Stable solutions for a two node co-reliance group . . . . .	62
3.7.2	Co-reliance group reduction . . . . .	63
3.7.3	Oscillation classes Venn diagram . . . . .	65
3.7.4	Example co-reliance groups for each oscillation class. . . . .	66
3.7.5	The state machine for the four-node 'Good' state machine in Figure 3.7.4(a) . . . . .	67
3.7.6	The state machine for the three-node 'Bad' co-reliance group in Figure 3.7.4(b) . . . . .	68
3.7.7	The state machine for the five-node 'Naughty' co-reliance group in Figure 3.7.4(c) . . . . .	69
3.7.8	The state machine for the four-node 'Asymptotically Good' co-reliance group in Figure 3.7.4(d) . . . . .	71

3.7.9	The state machine of the four-node ‘Asymptotically Good’ co-reliance group in Figure 3.7.4(d) with inbound $i$ at node 3 . . . . .	73
3.8.1	The state machine of the single cycle three-node co-reliance group with all ‘weak’ reliances . . . . .	76
3.9.1	Prioritization of the egress instances currently used in the AS. . .	84
3.9.2	Prioritization of egress instances consistent with available measurement data. . . . .	87
3.9.3	Prioritization of egress instances for an online tool. . . . .	90
3.9.4	An example of the prioritization of egress instances . . . . .	91
3.9.5	Equivalent example to Figure 3.9.4 with a shorter sliding window	92
3.11.1	Three-level route-reflector hierarchy . . . . .	97
3.11.2	Oscillation in a three-level route-reflector hierarchy . . . . .	97
3.11.3	Oscillation in three-level route-reflector hierarchy (bottom level full-mesh). . . . .	98
3.11.4	Oscillation between levels of route-reflector topology. . . . .	99
3.11.5	An example three-level topology with three child preference paths	101
3.11.6	An example from a Tier-2 AS of a route-reflector preferring a downstream egress learned from a non-downstream router . . . .	102
3.11.7	Four-level route-reflector hierarchy . . . . .	105
3.11.8	Modified four-level route-reflector hierarchy . . . . .	105
4.2.1	Stable egress instances violating Griffin and Wilfong’s condition .	113
4.3.1	Reliances and co-reliance groups for examples in Figure 4.2.1 . . .	116
4.4.1	An example three-level route-reflector topology . . . . .	118
4.5.1	Router comparison subroutine for a non-singleton co-reliance group	127
4.5.2	Function and variable definitions used in the <code>compare_routers</code> and the network solver algorithm. . . . .	128
4.5.3	Network solver algorithm . . . . .	129
4.5.4	Subroutine <code>igp_change</code> for determining the reliance graphs requiring recalculation when an IGP distance changes. . . . .	131
4.7.1	An example topology where the MED attribute is respected . . . .	135

4.7.2	Reliance graph for a full-mesh topology . . . . .	137
4.7.3	Full-mesh topology with the MED attribute respected. . . . .	138
4.7.4	An example confederation of sub-ASes and the corresponding reliance graph. . . . .	139
5.2.1	The impact of non-canonical peering policy . . . . .	144
5.5.1	Plot of the proportion of peers implementing a non-canonical peering policy . . . . .	151
5.5.2	The techniques used by a subset of peers to de-preference routes .	152
5.5.3	Example of peers displaying different behavior modes . . . . .	155
5.6.1	The impact of non-canonical peering policy . . . . .	161
5.6.2	A CDF for the impact of two peers' non-canonical policy on all routers . . . . .	162
5.6.3	A CCDF showing the proportion of decisions affected by non- canonical policies of peers . . . . .	163
5.6.4	The possible impact of a peers policy in the absence of routes from other ASes . . . . .	164
5.6.5	The impact of the non-canonical peering policy of the two peers' from Figure 5.6.2 when routes from other ASes are unavailable . .	165
5.6.6	Example of "when good routes go bad" phenomenon . . . . .	166
5.6.7	Traffic Impact: Finding the ingress router of a flow . . . . .	169
5.6.8	Traffic Impact: Finding the egress link under a canonical policy .	170
5.6.9	The shift of traffic that would occur for various ingress PoPs if a peer were to use a canonical peering policy. . . . .	172
5.7.1	Policy changes for one peer during interval September 1, 2007 - January 14, 2008 . . . . .	173
5.8.1	Summary table of peers . . . . .	179
5.8.2	The canonical peering policy of <i>Kangaroo Corp.</i> . . . . .	180
5.8.3	Legend for a peer's de-preferencing techniques. . . . .	181
5.8.4	The non-canonical peering policy of <i>Emu Inc.</i> . . . . .	181
5.8.5	The non-canonical peering policy of <i>Platypus Tech</i> . . . . .	182

5.8.6	The non-canonical peering policy of <i>Dingo Net</i> . . . . .	183
5.9.1	Decentralized mitigation scheme . . . . .	187
6.4.1	Consistency-check example . . . . .	203
6.4.2	Oldest prefix characteristic . . . . .	205
6.5.1	Finding the interval of extended measurement artifacts . . . . .	208
6.6.1	Detected failures in inter-table interval and the time we infer the missing withdrawal occurred. . . . .	212
6.8.1	A cartoon illustration of a monitoring BGP session to determine sliding window thresholds . . . . .	219
6.8.2	A cartoon illustration of the multi-variate threshold obtained us- ing LDA on the data points from Figure 6.8.1. . . . .	220
6.8.3	Cartoon illustration of independent thresholds obtained using LDA on the data points from Figure 6.8.1. . . . .	221
6.8.4	Example of sliding window parameter selection . . . . .	221
6.8.5	Anomalous data-points . . . . .	222
6.8.6	Cartoon illustration of monitoring BGP session to determine bin parameters . . . . .	224
6.8.7	Cartoon illustration of LDA producing an undesirable class sep- aration. . . . .	225
6.8.8	Cartoon illustration showing a desirable class separation. . . . .	226
6.8.9	Cartoon illustration using LDA to tune thresholds independently. . . . .	226
6.8.10	Example of bin parameter selection . . . . .	227

# List of Tables

2.0.1	Example Forwarding Table . . . . .	16
2.4.1	Step-by-step route selections for Bad Gadget . . . . .	33
2.4.2	Step-by-step route selections for Naughty Gadget . . . . .	33
2.4.3	Best route selection at routers 0 and 1 from Figure 2.4.3 . . . . .	36
2.4.4	Step-by-step route selection for Figure 2.4.4 . . . . .	38
3.7.1	Properties of oscillation classes. . . . .	65
3.8.1	Table showing the result of $\oplus$ for weak and strong reliances. . . . .	75
3.11.1	The egress selected by routers 0, 1 and 2 in Figure 3.11.2 . . . . .	96
3.11.2	The egress selected by routers 0, 1 and 2 in Figure 3.11.3. . . . .	98
3.11.3	The egress selected by routers 0 – 5 in Figure 3.11.4 . . . . .	100
4.4.1	Downstream egress sets for routers in example topology of Figure 4.4.1. . . . .	120
4.4.2	Reliances for example topology of Figure 4.4.1 . . . . .	122
5.5.1	Summary of peer behavior modes . . . . .	157
5.7.1	Number of snapshots before the policy change detection algo- rithm identifies a policy change . . . . .	176
6.3.1	Data characteristics of main measurement artifacts . . . . .	201
6.7.1	Default parameter settings. . . . .	214
6.9.1	Summary of consistency-check failures . . . . .	231
6.9.2	Session failure characteristics . . . . .	231

# Abstract

The Internet is literally an “Inter-Network”, that is, a network of networks. Networks can be entities including Internet Service Providers (ISPs), universities and commercial enterprises. Every network or Autonomous System (AS) has individual requirements, restrictions and capabilities to transit data traffic. No central controlling body determines how ASes connect — instead contractual agreements are established between AS pairs to govern their relationship. It is not feasible for all ASes to be physically connected to all others. Consequently, some ASes provide transit between other ASes. Such a service usually results in remuneration from one or both ASes.

Unlike centrally administered networks where all nodes in the network make generic, predictable decisions, each AS has the ability to select its best route based on its own proprietary commercial agreements. Such agreements are converted to a technical policy implemented in the Border Gateway Protocol (BGP). The ability to implement policies ensures the commercial viability of the Internet, but also makes the prediction of routes difficult and even more worrisome, conflicting policies can cause undesirable BGP states where no single AS has sufficient knowledge to understand what is happening [43].

Designing new clean-slate routing protocols is one approach to improving the predictability and reliability of the Internet. However, due to the Internet’s distributed political and administrative control, significant collaboration is required to implement a new routing protocol — especially when no new protocol currently proposed has sufficiently superior flexibility, scalability or robustness. The difficulty in implementing new and improved protocols is evident in the deploy-

ment of IPv6 [23]. Although the IPv6 standard has been defined for over a decade and offers a larger address space, better security and embedded quality of service in comparison to traditional IPv4, its deployment is limited to 1200 of over 30000 ASes in the Internet [66]. Hence, it is crucial practical solutions to current problems are evolved in addition to developing clean-slate techniques. Consequently, our approach is pragmatic — designing tangible solutions to practical problems that can be implemented immediately.

In this thesis we examine and combine eBGP, iBGP, OSPF, Netflow and router configuration data to discover important aspects of routing. It is this investigation that instigated the development of a model of iBGP. iBGP is the version of BGP implemented *within* ASes to propagate routes between internal routers. It exists on a logical topology, however it interacts with the physical topology. It is this interaction which can cause persistent oscillation [49] — a system state where routers alter their decision ad infinitum. Detecting configurations which can cause this oscillation is NP-hard [49]. However, our model of iBGP introduced in this thesis benefits from the ‘designed’ structure of the iBGP topology to restrict the search space dramatically to one that is computationally feasible.

iBGP data — which is collected to analyze router decisions — is often only collected on a subset of routers due to its massive storage requirements. In addition there is a substantial amount of correlation between router decisions. Our model of iBGP discovers the dependencies between router decisions and can consequently predict the decisions of those routers for which no measurements are available. It does not rely on any assumption of operator configuration, and subsequently is able to be applied in any network scenario — not just the one originally configured. It is this feature, together with the model’s ability to use any available measurement data that makes our technique ideal for network measurement and management applications. We found our model is efficient and accurate on the network of a large Tier-2 AS, where all but seven of over 12.7 million decisions were consistent with observed data. Further we were able to predict the decision of 85% of routers where observed data was unavailable.

During our analysis, we also identified several minor configuration errors on operational routers when we predicted the “correct” outcome.

The internal state of a network can be influenced by neighboring ASes. Peering agreements are closely guarded due to their commercial sensitivity. They are implemented in BGP in the form of policies and are difficult to infer with publicly available data sources. We examined the peering policies of over 100 ASes from the perspective of a large Tier-2 AS, finding 22% differ from the canonical peering policy outlined in many peering agreements. When a policy differs from the canonical peering policy, it may result in sub-optimal routing within the Tier-2 AS. We used our model of iBGP to firstly predict the decisions of all routers under the current peering policy, before determining the changes that would have occurred under a canonical peering policy. This analysis not only provided a metric for the routing impact of a peers’ non-canonical policy, but subsequently used in combination with traffic data allowed us to determine the influence of the peer on traffic flows. Our techniques described allow an AS to fully quantify the impact of a non-canonical peering policy and adapt business arrangements appropriately.

Throughout our analysis of BGP data, we noticed several inconsistencies in the data. Although the results in the above work were insensitive to such inconsistencies, other applications requiring accurate, fine time-scale analysis of the routing state are much more sensitive. Consequently, we undertake a self-consistency check on the BGP data and examine the possible causes of such inconsistencies. We also present a mechanism to ‘clean’ the data to minimize the effects of any inconsistency.