# SCHEDULING IN METACOMPUTING SYSTEMS

By

Heath A. James, B.Sc.(Ma. & Comp. Sc.)(Hons)

July 1999

A THESIS SUBMITTED FOR THE DEGREE OF

DOCTOR OF PHILOSOPHY

IN THE DEPARTMENT OF COMPUTER SCIENCE

UNIVERSITY OF ADELAIDE

# Contents

# Abstract

The problem of scheduling in a distributed heterogeneous environment is that of deciding where to place programs, and when to start execution of the programs. The general problem of scheduling is investigated, with focus on jobs consisting of both independent and dependent programs. We consider program execution within the context of metacomputing environments, and the benefits of being able to make predictions on the performance of programs. Using the constraint of restricted placement of programs, we present a scheduling system that produces heuristically good execution schedules in the absence of complete global system state information. The scheduling system is reliant on a processor-independent global naming strategy and a single-assignment restriction for data.

Cluster computing is an abstraction that treats a collection of interconnected parallel or distributed computers as a single resource. This abstraction is commonly used to refer to the scope of resource managers, most often in the context of queueing systems. To express greater complexity in a cluster computing environment, and the programs that are run on the environment, the term *metacomputing* [74] is now being widely adopted. This may be defined as a collection of possibly heterogeneous computational nodes which have the appearance of a single, virtual computer for the purposes of resource management and remote execution.

We review current technologies for cluster computing and metacomputing, with focus on their resource management and scheduling capabilities. We critically review these technologies against our Distributed Information Systems Control World (DISCWorld).

We develop novel mechanisms that enable and support scheduling and program placement in the DISCWorld prototype. We also discuss a mechanism by which a high-level job's internal structure can be represented, and processing requests controlled. We formulate this using extended markup language (XML).

ix

To enable processing requests, which consist of a directed acyclic graph of programs, we develop a mechanism for their composition and scheduling placement. This rich data pointer and a complimentary futures mechanism are implemented as part of the DISCWorld remote access mechanism (DRAM). They also form the basis for a model of wide-area computation independent of DISCWorld.

We have implemented geospatial imagery applications using both simple RMI and common gateway interface, and the novel mechanisms developed in this thesis. After analysing and measuring our system, we find significant performance and scalability benefits.