# History-Based Batch Job Scheduling on a Network of Interactively Used Workstations

**Inauguraldissertation**

zur
Erlangung der Würde eines Doktors der Philosophie
vorgelegt der
Philosophisch-Naturwissenschaftlichen Fakultät
der Universität Basel

von
Andreas Wespi
aus Schüpfheim, LU

Basel, 1999

Genehmigt von der Philosophisch-Naturwissenschaftlichen Fakultät auf Antrag der Herren

Prof. Dr. Helmar Burkhart

Prof. Dr. Clemens H. Cap

Basel, den 17. November 1998

Prof. Dr. Stefan M. Schmid, Dekan

Research Reports in Computer Science

Band 7

**Andreas Wespi**

# History-Based Batch Job Scheduling on a Network of Interactively Used Workstations

Printed in Germany.

To Sarah, our beloved little baby,
who stayed with us for only a short time.

# Acknowledgements

I am deeply indebted to my advisor Prof. Dr. Helmar Burkhart for his guidance, understanding, and support during the evolution of this thesis. I am also grateful to Prof. Dr. Clemens H. Cap, who kindly agreed to be my co-advisor, for his interest in my work and the animated discussions we had.

Many thanks to the IBM Zurich Research Laboratory, which gave me the opportunity to work on this thesis. My sincere thanks go particularly to Dr. Ernst Rothauser for his continuous help, guidance, encouragement, and the many hours of lively and fruitful discussions. I also thank my colleagues of the Information Systems group for their technical advice and understanding, especially when my data collection processes filled up too much disk space. I am indebted to Lilli Pavka for proofreading this text.

I am grateful to all those many people who contributed to this thesis with advice and expertise.

I thank my parents, who always supported me in good times and bad, and who taught me to love to learn. Most of all I thank my wife Rahel, who supported me with her understanding even though I often tried her patience working late in the evenings and on weekends.

# Abstract

Networks of workstations are the computing environment at many sites nowadays. Workstations are used for daily interactive work but also, owing to their excellent price/performance ratio, as dedicated computing servers. As it is known that many workstations sit idle in a network of interactively used workstations, approaches have been developed to make use of the idle resources by employing them for computation-intensive tasks.

Distributed Batch Job Scheduling Systems (DBJSSs) have been built that allow batch jobs to be run on interactively used workstations in a way that the impact on interactive users is kept to a minimum. Whereas batch job scheduling in a dedicated cluster of computing servers is well understood, scheduling in a network of interactively used workstations, due to the unpredictability of the available resources, is quite different and has only marginally been investigated in the past. As a consequence, either simple scheduling strategies are implemented in current DBJSSs or the system administrators are given the means to implement their own scheduling strategy without any further assistance.

I have analyzed the workstation utilization data of the distributed computing environment of the IBM Zurich Research Laboratory for more than one year. Based on a new approach, I have evaluated the impact of different definitions of the workstation idle state on the total amount of resources available for batch job processing. Furthermore, I could show that many workstations have a regular usage profile.

This observation has led to the development of novel, history-based batch job scheduling algorithms for networks of interactively used workstations. To quantify the performance benefits of the newly developed algorithms, the SIByL simulation environment has been built to analyze different batch job scheduling approaches under the same environment

conditions.

Simulation results show that history-based batch job scheduling algorithms outperform the traditional scheduling approaches implemented in many DBJSSs. The traditional approaches are based only on the state of the workstations at the time when the scheduling decision takes place, whereas history-based algorithms also consider the workstation usage in the past. For moderately loaded systems and environments where job migration is of limited use, the history-based algorithms result in mean job turnaround times that are about 30% shorter than the times obtained with the traditional scheduling approaches.

History-based batch job scheduling results in an increased amount of data to be processed to find a good job placement. I show how current DBJSSs can be easily enhanced so that they support history-based batch job scheduling without generating a lot of additional network load or additional computing load on the client workstations. Increased computing load is placed only on the server side. Furthermore, the SIByL simulation tool can be integrated in the DBJSS server such that, based on prior workstation usage data and job execution statistics, the best scheduling algorithm is selected adaptively out of a set of predefined algorithms.

# Contents

# List of Figures

# List of Tables