

Scheduling Algorithms for Input-Queued Cell Switches

by

Nicholas William McKeown

B.Eng (University of Leeds) 1986

M.S. (University of California at Berkeley) 1992

A thesis submitted in partial satisfaction of the
requirements for the degree of

Doctor of Philosophy

in

Engineering-Electrical Engineering
and Computer Sciences

in the

GRADUATE DIVISION

of the

UNIVERSITY of CALIFORNIA at BERKELEY

Committee in charge:

Professor Jean Walrand, Chair

Professor Pravin P. Varaiya

Professor Ronald W. Wolff

1995

This thesis of Nicholas William McKeown is approved:

Chair

Date

Date

Date

University of California at Berkeley

1995

Scheduling Algorithms for Input-Queued Cell Switches

© 1995

by

Nicholas William McKeown

Abstract

Scheduling Algorithms for Input-Queued Cell Switches

by

Nicholas William McKeown

Doctor of Philosophy in Engineering-Electrical Engineering and Computer Sciences

University of California at Berkeley

Professor Jean Walrand, Chair

The algorithms described in this thesis are designed to schedule cells in a very high-speed, parallel, input-queued crossbar switch. We present several novel scheduling algorithms that we have devised, each aims to match the set of inputs of an input-queued switch to the set of outputs more efficiently, fairly and quickly than existing techniques.

In Chapter 2 we present the simplest and fastest of these algorithms: SLIP — a parallel algorithm that uses rotating priority (“round-robin”) arbitration. SLIP is simple: it is readily implemented in hardware and can operate at high speed. SLIP has high performance: for uniform i.i.d. Bernoulli arrivals, SLIP is stable for any admissible load, because the arbiters tend to *desynchronize*. We present analytical results to model this behavior. However, SLIP is not always stable and is not always monotonic: adding more traffic can actually make the algorithm operate more efficiently. We present an approximate analytical model of this behavior. SLIP prevents starvation: all contending inputs are eventually served. We present simulation results, indicating SLIP’s performance. We argue that SLIP can be readily implemented for a 32x32 switch on a single chip.

In Chapter 3 we present *i*-SLIP, an iterative algorithm that improves upon SLIP by converging on a maximal size match. The performance of *i*-SLIP improves with up to $\log_2 N$ iterations. We show that although it has a longer running time than SLIP, an *i*-SLIP scheduler is little more complex to implement.

In Chapter 4 we describe maximum or maximal *weight* matching algorithms based on the occupancy of queues, or waiting times of cells. These algorithms are stable over a wider range of traffic loads. We describe two algorithms, *longest queue first* (LQF) and *oldest cell first* (OCF) and consider their performance. We prove that LQF, although too complex to implement in hardware, is stable under all admissible i.i.d. offered loads. We consider two implementable, iterative algorithms *i*-LQF and *i*-OCF which converge on a maximal weight matching. Finally, we present two interesting implementations of the Gale-Shapley algorithm, designed to solve the *stable marriage problem*.

To my parents,

and

My Wife,

My Love,

My Le.

Table of Contents

Acknowledgements	viii
-------------------------------	-------------

CHAPTER 1

Introduction	1
1 Problem Statement	1
2 Motivation.....	4
2.1 Datapath for an Input-Queued Switch	4
2.2 Controlling the Datapath.....	5
3 Background.....	7
3.1 Input vs. Output Queueing.....	7
3.2 Overcoming Head-of-Line Blocking	8
3.3 Previous Scheduling Work	9
3.3.1 Maximum Size Matching.....	9
3.3.2 Neural Network Algorithms	10
3.3.3 Scheduling into the Future.....	11
3.3.4 Parallel Iterative Matching.....	12
3.4 Simple Comparison of Previous Techniques.....	14
4 Outline of Thesis.....	16

CHAPTER 2

The SLIP Algorithm

with a Single Iteration	18
1 Introduction.....	18
2 Basic Round-Robin Matching Algorithm.....	20
2.1 Performance of RRM for Bernoulli Arrivals.....	20
3 The SLIP Algorithm	23
4 Simulated Performance of SLIP	24
4.1 Bernoulli Traffic	24
4.2 “Bursty” Traffic	26
4.3 As a Function of Switch Size.....	28
4.4 Burst Reduction	30
5 Analysis of SLIP Performance.....	32

5.1	Convergence to Time-Division Multiplexing Under Heavy Load	32
5.2	Desynchronization of Arbiters	32
5.3	Stability of SLIP	35
5.3.1	Drift Analysis of a 2x2 SLIP Switch: First Approximation	38
5.3.2	Drift Analysis of a 2x2 SLIP Switch: Second Approximation ...	41
5.4	Approximate Delay Model for 2x2 SLIP Switch	41
6	Variations on SLIP	44
6.1	Prioritized SLIP	44
6.2	Threshold SLIP	45
6.3	Weighted SLIP	46
6.4	Least Recently Used	46
7	Implementing SLIP	49
7.1	Prioritized SLIP	51

CHAPTER 3

The SLIP Algorithm

with Multiple Iterations53

1	Introduction.....	53
2	The Iterative SLIP Matching Algorithm.....	54
2.1	Description.....	54
2.2	Updating Pointers.....	55
2.3	Properties	56
3	Simulated Performance of Iterative SLIP.....	57
3.1	How Many Iterations?.....	57
3.2	Bernoulli Traffic	58
3.3	Bursty Traffic.....	63
3.4	As a Function of Switch Size.....	65
4	Variations of Iterative SLIP	66
4.1	Iterative SLIP with LRU Accept Arbiters	66
4.2	Separate Pointers for each Iteration	69
5	Implementing Iterative SLIP.....	69

CHAPTER 4

Weighted Matching

Algorithms	72
1 Introduction.....	72
2 Maximum Weight Matching.....	73
2.1 Starvation with LQF	74
2.2 Performance of LQF and OCF Algorithms	75
2.2.1 Uniform Workload.....	75
2.2.2 Non-Uniform Workload.....	75
2.2.3 Stability of 2x2 Switch	77
2.2.4 Stability of NxN Switch.....	80
3 Iterative Maximal Weight Matching Algorithms.....	81
3.1 i-LQF.....	81
3.2 Properties	82
3.3 i-OCF	82
3.4 Properties	83
3.5 Performance of i-LQF and i-OCF.....	83
3.5.1 Uniform Workload.....	83
3.5.2 Nonuniform Workload.....	83
3.5.3 Stability for 2x2 Switch.....	83
3.6 Implementation of i-LQF and i-OCF.....	84
4 Stable Marriages	87
4.1 The Stable Marriage Problem	88
4.2 The Gale-Shapley Algorithm.....	88
4.3 Analogy to Switch Scheduling.....	89
4.3.1 The GS-LQF Algorithm.....	90
4.3.2 The GS-OCF Algorithm	90
4.4 Performance of GS-LQF and GS-OCF Algorithms.....	91
4.5 Implementation of GS-LQF and GS-OCF.....	91
4.6 Egalitarian Stable Marriage	92
References	94

APPENDIX 1

Arbiter Synchronization for Single-Iteration SLIP97

APPENDIX 2

Stability of Single-Iteration

SLIP Algorithm101

- 1 Single-Step Drift Analysis of 2x2 Switch with 1 Queue101
 - 1.1 First Approximation.....101
 - 1.2 Second Approximation103
- 2 Matrix Geometric Solution for 2x2 Switch with 1 Queue.....106

APPENDIX 3

Stability of 2x2

Switch109

- 1 Stability of 2x2 Switch with 3 Active Flows (Theorem 4.2).....109
 - 1.1 Definitions.....109
 - 1.2 Problem statement.....110
 - 1.3 Solution.....110
 - 1.4 Stable Algorithms111
- 2 Relative Queue Sizes (Theorem 4.3)111

APPENDIX 4

Stability of NxN Switch

with i.i.d. Arrivals115

- 1 Definitions.....115
- 2 Main Theorem.....116
- 3 Proof.....116

Acknowledgements

For their continued guidance, support and encouragement throughout my time at Berkeley, I would like to thank my adviser Jean Walrand and Pravin Varaiya. I greatly appreciate the freedom and collegial respect you have given me and your other students.

Numerous discussions with Richard Edell lead to the design of the datapath that provided the main motivation for this thesis. Richard, you are a truly gifted engineer and it has been a pleasure to be your colleague.

I am grateful to Professor Tom Anderson for discussions about the iterative properties of the SLIP algorithm, and to Professor Venkat Anantharam for suggesting the proof in Appendix 4. I also wish to acknowledge the helpful feedback and suggestions of Chuck Thacker (DEC SRC), the inventor of parallel iterative matching. I thank Dana Randall for introducing me to the stable marriage problem, and Matthew J. Salzman (CMU) for kindly donating his code to implement the maximum match algorithms used in many of my simulations.

I am extremely grateful to John Limb, for whom I worked at Hewlett-Packard Labs in Bristol, England. John, you have been a constant source of inspiration to me; and without your encouragement I would not have gone back to school to pursue my Ph.D.

I would like to give thanks to the numerous other people at Hewlett-Packard Labs who supported me over the years, in particular John Taylor, Daniel Pitt, Steve Wright and Gwenda Ward.

Last, and definitely most, I want to thank my family. Words cannot express my thanks to my wife and parents for all your love and encouragement. I dedicate this thesis to you.