

INTRODUCING ENHANCED FULLY-ADAPTIVE ROUTING  
DECISIONS WITHIN TORUS-MESH AND HYPERCUBE  
INTERCONNECT NETWORKS

by

CHRISTOPHER L LYDICK

B.S., Kansas State University, 2006

---

A THESIS

submitted in partial fulfillment of the  
requirements for the degree

MASTER OF SCIENCE

Department of Electrical and Computer Engineering  
College of Engineering

KANSAS STATE UNIVERSITY  
Manhattan, Kansas  
2008

Approved by:

Major Professor  
Don Gruenbacher

# Copyright

Christopher L Lydick

2008

# Abstract

The method for communicating within an interconnection network, or fabric of connections between nodes, can be as diverse as are the applications which utilize them. Because of dynamic traffic loads on these interconnection networks, fully-adaptive routing algorithms have been shown to exploit locality while balancing loads and softening the effects of hotspots. One issue which has been overlooked is the impact of data traveling along the periphery of a selected minimal routable quadrant (MRQ) within these fully-adaptive algorithms. As data aligns with the destination in the x, y, and z dimensions for instance, the data then traverses the periphery of an MRQ. For each dimension that this occurs, the data is given one less choice for routing around hotspots which could appear later along the path. By weighting the decision of selecting a next-hop by avoiding the periphery of the selected MRQ, the data then has more options for avoiding hotspots. One hybridized routing algorithm which borrows heavily from CQR (an efficient and stable fully-adaptive algorithm), is introduced within this work. Enhanced CQR with Periphery Avoidance, attempts to weight the routing decision for a next hop using both output queues and the proximity to the periphery of the MRQ. This fully-adaptive algorithm is tested using simulations and a laboratory research cluster using a USB interconnect in the hypercube topology. It is also compared against other static, oblivious, and adaptive algorithms. Thor's Tack Hammer, the Kansas State University research cluster, is also benchmarked and discussed as an inexpensive and dependable parallel system.

# Table of Contents

<b>Table of Contents</b>	<b>iv</b>
<b>List of Figures</b>	<b>viii</b>
<b>List of Tables</b>	<b>xii</b>
<b>Acknowledgements</b>	<b>xiii</b>
<b>Dedication</b>	<b>xiv</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Interconnection Networks . . . . .	1
1.2 High Performance Computing Clusters . . . . .	2
1.3 Routing within Interconnection Networks . . . . .	5
1.4 Contributions . . . . .	6
<b>2 Previous Work</b>	<b>8</b>
2.1 Undeliverable Data . . . . .	8
2.1.1 Deadlock . . . . .	8
2.1.2 Livelock . . . . .	9
2.1.3 More on Undeliverable Data . . . . .	10
2.2 Static Routing Algorithms . . . . .	11
2.2.1 Dimension Ordered Routing, DOR . . . . .	12
2.2.2 Direction Ordered Routing, DIR . . . . .	12
2.3 Adaptive Routing Algorithms . . . . .	13
2.3.1 Chaotic . . . . .	13
2.3.2 Minimally Adaptive . . . . .	14
2.3.3 Fully Adaptive: GOAL, GAL, CQR . . . . .	15
2.4 Oblivious Routing Algorithms . . . . .	17
2.4.1 Valiant’s Algorithm . . . . .	17
2.4.2 Minimal Oblivious . . . . .	18
2.5 Traffic Patterns . . . . .	18
2.6 Latency and Throughput Analyses . . . . .	19
<b>3 Enhancing Current Adaptive Routing Algorithms</b>	<b>21</b>
3.1 Issues with Previous Algorithms . . . . .	21
3.2 Introducing Enhanced CQR with Periphery Avoidance . . . . .	22
3.2.1 Path Diversity and Periphery Avoidance . . . . .	22



3.2.2	Periphery Avoidance Decision Function . . . . .	24
<b>4</b>	<b>The Research Cluster: Thor's Tack Hammer</b>	<b>28</b>
4.1	Topology . . . . .	28
4.1.1	Subnetting and Addressing Scheme . . . . .	30
4.2	Node and Cluster Specifications . . . . .	35
4.3	Node and Cluster Performance . . . . .	39
4.3.1	USB Interconnect Throughput . . . . .	40
4.3.2	Linpack Benchmark . . . . .	40
4.3.3	Power Consumption . . . . .	42
<b>5</b>	<b>Laboratory Results</b>	<b>44</b>
5.1	Methodology for Tests . . . . .	44
5.1.1	A Packet's Progression through the Interconnection Network . . . . .	46
5.2	Laboratory Results . . . . .	48
5.2.1	Dimension Ordered Routing (DOR) Results . . . . .	51
5.2.2	Direction Ordered Routing (DIR) Results . . . . .	52
5.2.3	Minimal Oblivious Routing Results . . . . .	52
5.2.4	Minimal Adaptive Routing Results . . . . .	53
5.2.5	CQR Routing Results . . . . .	53
5.2.6	Enhanced CQR Routing Results . . . . .	54
5.2.7	Comparison of All Results . . . . .	57
<b>6</b>	<b>Simulation Results</b>	<b>59</b>
6.1	Methodology for Tests . . . . .	59
6.2	Results from the 3-ary 3-cube Simulations . . . . .	61
6.3	Results from the 4-ary 3-cube Simulations . . . . .	64
6.4	Results from the 5-ary 3-cube Simulations . . . . .	67
6.4.1	Comparison of All Results . . . . .	70
<b>7</b>	<b>Conclusions</b>	<b>71</b>
7.1	Future Work . . . . .	73
	<b>References</b>	<b>76</b>
<b>A</b>	<b>Appendix A: Linux BASH Scripts</b>	<b>77</b>
A.1	USB Interconnect Scripts . . . . .	77
A.1.1	file: interconnectUp . . . . .	77
A.1.2	file: usbUp . . . . .	79
A.1.3	file: usbPing . . . . .	82
A.2	Data Analysis Scripts . . . . .	84
A.2.1	file: interconnectTime . . . . .	84
A.2.2	file: interconnectRouters . . . . .	85
A.2.3	file: nodeSetup . . . . .	86

A.2.4	file: interconnectThroughput . . . . .	87
<b>B</b>	<b>Appendix B: C Programs</b>	<b>91</b>
B.1	file: server.c . . . . .	91
B.2	file: loser.c . . . . .	111
B.3	file: injector.c . . . . .	114
B.4	file: router.h . . . . .	116
<b>C</b>	<b>Appendix C: Matlab Simulation Scripts</b>	<b>117</b>
C.1	User Interface Scripts . . . . .	117
C.1.1	file: do_load.m . . . . .	117
C.1.2	file: view_load.m . . . . .	123
C.2	Routing Scripts . . . . .	125
C.2.1	file: dim2i.m . . . . .	125
C.2.2	file: i2dim.m . . . . .	126
C.2.3	file: dimord_route.m . . . . .	127
C.2.4	file: dirord_route.m . . . . .	129
C.2.5	file: minobliv_route.m . . . . .	132
C.2.6	file: do_minadaptive.m . . . . .	133
C.2.7	file: do_cqr1.m . . . . .	137
C.2.8	file: do_cqr2.m . . . . .	141
C.2.9	file: cqr_addPacket.m . . . . .	145
C.2.10	file: cqr_chooseQuadrant.m . . . . .	146
C.2.11	file: cqr_getRandomValue.m . . . . .	155
C.2.12	file: cqr_transformTrafficMatrix.m . . . . .	156
<b>D</b>	<b>Appendix D: Full Laboratory Results</b>	<b>157</b>
D.1	Matlab Plot Scripts . . . . .	157
D.1.1	file: convertLoad.m . . . . .	157
D.1.2	file: view_load.m . . . . .	159
D.1.3	file: i2dim.m . . . . .	159
D.1.4	file: dim2i.m . . . . .	159
D.2	Static Algorithms . . . . .	159
D.2.1	Dimension Ordered Routing Results . . . . .	159
D.2.2	Direction Ordered Routing Results . . . . .	163
D.3	Oblivious Algorithms . . . . .	167
D.3.1	Minimal Oblivious Routing Results . . . . .	167
D.4	Adaptive Algorithms . . . . .	171
D.4.1	Minimal Adaptive Routing Results . . . . .	171
D.4.2	CQR Routing Results . . . . .	175
D.4.3	Enhanced CQR Routing Results . . . . .	179

<b>E</b>	<b>Appendix E: Full Matlab Simulation Results</b>	<b>183</b>
E.1	Results from the 3-ary 3-cube Topology . . . . .	183
E.1.1	Static Algorithms . . . . .	183
E.1.2	Oblivious Algorithms . . . . .	190
E.1.3	Adaptive Algorithms . . . . .	194
E.2	Results from the 4-ary 3-cube Topology . . . . .	206
E.2.1	Static Algorithms . . . . .	206
E.2.2	Oblivious Algorithms . . . . .	213
E.2.3	Adaptive Algorithms . . . . .	217
E.3	Results from the 5-ary 3-cube Topology . . . . .	229

# List of Figures

1.1	A Torus . . . . .	3
1.2	4-ary 2-cube Slice of a Tori-Mesh Network . . . . .	4
2.1	Deadlock Example . . . . .	9
2.2	Livelock Example . . . . .	10
2.3	Minimal Adaptive Quadrants . . . . .	14
2.4	GAL Injection Queues . . . . .	16
3.1	The Ratio of Periphery Links to that of Non-Periphery Links . . . . .	25
4.1	Thor's Tack Hammer Node Connectivity . . . . .	29
4.2	An Addressing and Subnetting Example . . . . .	34
4.3	Thor's Tack Hammer Network Topology . . . . .	35
4.4	Single Node: The VIA VT310-DP . . . . .	36
4.5	Adaptec USB 2.0 Card used on Thor's Tack Hammer . . . . .	37
4.6	USB 2.0 NetLink Cables used on Thor's Tack Hammer . . . . .	38
5.1	An Example of Excellent Laboratory Results . . . . .	49
5.2	An Example of OK Laboratory Results . . . . .	50
5.3	An Example of Poor Laboratory Results . . . . .	51
5.4	CQR Routing (FL) TTH Results . . . . .	54
5.5	CQR Routing (FL-HS) - TTH Results . . . . .	55
5.6	Enhanced CQR Routing (FL) TTH Results . . . . .	56
5.7	Enhanced CQR Routing (FL-HS) TTH Results . . . . .	56
6.1	CQR Routing (FL) 3-ary 3-cube Simulation Results . . . . .	62
6.2	CQR Routing (FL-HS) 3-ary 3-cube Simulation Results . . . . .	62
6.3	Enhanced CQR Routing (FL) 3-ary 3-cube Simulation Results . . . . .	63
6.4	Enhanced CQR Routing (FL-HS) 3-ary 3-cube Simulation Results . . . . .	63
6.5	CQR Routing (FL) 4-ary 3-cube Simulation Results . . . . .	65
6.6	CQR Routing (FL-HS) 4-ary 3-cube Simulation Results . . . . .	65
6.7	ECQR Routing (FL) 4-ary 3-cube Simulation Results . . . . .	66
6.8	ECQR Routing (FL-HS) 4-ary 3-cube Simulation Results . . . . .	66
6.9	CQR Routing (FL) 5-ary 3-cube Simulation Results . . . . .	68
6.10	CQR Routing (FL-HS) 5-ary 3-cube Simulation Results . . . . .	68
6.11	ECQR Routing (FL) 5-ary 3-cube Simulation Results . . . . .	69
6.12	ECQR Routing (FL-HS) 5-ary 3-cube Simulation Results . . . . .	69
D.1	Dimension Ordered Routing (NN) TTH Results . . . . .	159

D.2	Dimension Ordered Routing (UR) TTH Results	160
D.3	Dimension Ordered Routing (BC) TTH Results	160
D.4	Dimension Ordered Routing (TP) TTH Results	161
D.5	Dimension Ordered Routing (TOR) TTH Results	161
D.6	Dimension Ordered Routing (FL) TTH Results	162
D.7	Direction Ordered Routing (NN) TTH Results	163
D.8	Direction Ordered Routing (UR) TTH Results	164
D.9	Direction Ordered Routing (BC) TTH Results	164
D.10	Direction Ordered Routing (TP) TTH Results	165
D.11	Direction Ordered Routing (TOR) TTH Results	165
D.12	Direction Ordered Routing (FL) TTH Results	166
D.13	Minimal Oblivious Routing (NN) TTH Results	167
D.14	Minimal Oblivious Routing (UR) TTH Results	168
D.15	Minimal Oblivious Routing (BC) TTH Results	168
D.16	Minimal Oblivious Routing - Transpose Traffic - TTH Results	169
D.17	Minimal Oblivious Routing (TOR) TTH Results	169
D.18	Minimal Oblivious Routing (FL) TTH Results	170
D.19	Minimal Adaptive Routing (NN) TTH Results	171
D.20	Minimal Adaptive Routing (UR) TTH Results	172
D.21	Minimal Adaptive Routing (BC) TTH Results	172
D.22	Minimal Adaptive Routing (TP) TTH Results	173
D.23	Minimal Adaptive Routing (TOR) TTH Results	173
D.24	Minimal Adaptive Routing (FL) TTH Results	174
D.25	Minimal Adaptive Routing - (FL-HS) - TTH Results	174
D.26	CQR Routing (NN) TTH Results	175
D.27	CQR Routing (UR) TTH Results	176
D.28	CQR Routing (BC) TTH Results	176
D.29	CQR Routing (TP) TTH Results	177
D.30	CQR Routing (TOR) TTH Results	177
D.31	CQR Routing (FL) TTH Results	178
D.32	CQR Routing (FL-HS) - TTH Results	178
D.33	Enhanced CQR Routing (NN) TTH Results	179
D.34	Enhanced CQR Routing (UR) TTH Results	180
D.35	Enhanced CQR Routing (BC) TTH Results	180
D.36	Enhanced CQR Routing (TP) TTH Results	181
D.37	Enhanced CQR Routing (TOR) TTH Results	181
D.38	Enhanced CQR Routing (FL) TTH Results	182
D.39	Enhanced CQR Routing (FL-HS) TTH Results	182
E.1	Dimension Ordered Routing (NN) 3-ary 3-cube Simulation Results	184
E.2	Dimension Ordered Routing (UR) 3-ary 3-cube Simulation Results	184
E.3	Dimension Ordered Routing (BC) 3-ary 3-cube Simulation Results	185
E.4	Dimension Ordered Routing (TP) 3-ary 3-cube Simulation Results	185

E.5	Dimension Ordered Routing (TOR) 3-ary 3-cube Simulation Results . . . . .	186
E.6	Dimension Ordered Routing (FL) 3-ary 3-cube Simulation Results . . . . .	186
E.7	Direction Ordered Routing (NN) 3-ary 3-cube Simulation Results . . . . .	187
E.8	Direction Ordered Routing (UR) 3-ary 3-cube Simulation Results . . . . .	187
E.9	Direction Ordered Routing (BC) 3-ary 3-cube Simulation Results . . . . .	188
E.10	Direction Ordered Routing (TP) 3-ary 3-cube Simulation Results . . . . .	188
E.11	Direction Ordered Routing (TOR) 3-ary 3-cube Simulation Results . . . . .	189
E.12	Direction Ordered Routing (FL) 3-ary 3-cube Simulation Results . . . . .	189
E.13	Minimal Oblivious Routing (NN) 3-ary 3-cube Simulation Results . . . . .	190
E.14	Minimal Oblivious Routing (UR) 3-ary 3-cube Simulation Results . . . . .	191
E.15	Minimal Oblivious Routing (BC) 3-ary 3-cube Simulation Results . . . . .	191
E.16	Minimal Oblivious Routing (TP) 3-ary 3-cube Simulation Results . . . . .	192
E.17	Minimal Oblivious Routing (TOR) 3-ary 3-cube Simulation Results . . . . .	192
E.18	Minimal Oblivious Routing (FL) 3-ary 3-cube Simulation Results . . . . .	193
E.19	Minimal Adaptive Routing (NN) 3-ary 3-cube Simulation Results . . . . .	194
E.20	Minimal Adaptive Routing (UR) 3-ary 3-cube Simulation Results . . . . .	195
E.21	Minimal Adaptive Routing (BC) 3-ary 3-cube Simulation Results . . . . .	195
E.22	Minimal Adaptive Routing (TP) 3-ary 3-cube Simulation Results . . . . .	196
E.23	Minimal Adaptive Routing (TOR) 3-ary 3-cube Simulation Results . . . . .	196
E.24	Minimal Adaptive Routing (FL) 3-ary 3-cube Simulation Results . . . . .	197
E.25	Minimal Adaptive Routing (FL-HS) 3-ary 3-cube Simulation Results . . . . .	197
E.26	CQR Routing (NN) 3-ary 3-cube Simulation Results . . . . .	198
E.27	CQR Routing (UR) 3-ary 3-cube Simulation Results . . . . .	198
E.28	CQR Routing (BC) 3-ary 3-cube Simulation Results . . . . .	199
E.29	CQR Routing (TP) 3-ary 3-cube Simulation Results . . . . .	199
E.30	CQR Routing (TOR) 3-ary 3-cube Simulation Results . . . . .	200
E.31	CQR Routing (FL) 3-ary 3-cube Simulation Results . . . . .	200
E.32	CQR Routing (FL-HS) 3-ary 3-cube Simulation Results . . . . .	201
E.33	Enhanced CQR Routing (NN) 3-ary 3-cube Simulation Results . . . . .	202
E.34	Enhanced CQR Routing (UR) 3-ary 3-cube Simulation Results . . . . .	202
E.35	Enhanced CQR Routing (BC) 3-ary 3-cube Simulation Results . . . . .	203
E.36	Enhanced CQR Routing (TP) 3-ary 3-cube Simulation Results . . . . .	203
E.37	Enhanced CQR Routing (TOR) 3-ary 3-cube Simulation Results . . . . .	204
E.38	Enhanced CQR Routing (FL) 3-ary 3-cube Simulation Results . . . . .	204
E.39	Enhanced CQR Routing (FL-HS) 3-ary 3-cube Simulation Results . . . . .	205
E.40	Dimension Ordered Routing (NN) 4-ary 3-cube Simulation Results . . . . .	206
E.41	Dimension Ordered Routing (UR) 4-ary 3-cube Simulation Results . . . . .	207
E.42	Dimension Ordered Routing (BC) 4-ary 3-cube Simulation Results . . . . .	207
E.43	Dimension Ordered Routing (TP) 4-ary 3-cube Simulation Results . . . . .	208
E.44	Dimension Ordered Routing (TOR) 4-ary 3-cube Simulation Results . . . . .	208
E.45	Dimension Ordered Routing (FL) 4-ary 3-cube Simulation Results . . . . .	209
E.46	Direction Ordered Routing (NN) 4-ary 3-cube Simulation Results . . . . .	210
E.47	Direction Ordered Routing (UR) 4-ary 3-cube Simulation Results . . . . .	210

E.48	Direction Ordered Routing (BC) 4-ary 3-cube Simulation Results . . . . .	211
E.49	Direction Ordered Routing (TP) 4-ary 3-cube Simulation Results . . . . .	211
E.50	Direction Ordered Routing (TOR) 4-ary 3-cube Simulation Results . . . . .	212
E.51	Direction Ordered Routing (FL) 4-ary 3-cube Simulation Results . . . . .	212
E.52	Minimal Oblivious Routing (NN) 4-ary 3-cube Simulation Results . . . . .	213
E.53	Minimal Oblivious Routing (UR) 4-ary 3-cube Simulation Results . . . . .	214
E.54	Minimal Oblivious Routing (BC) 4-ary 3-cube Simulation Results . . . . .	214
E.55	Minimal Oblivious Routing (TP) 4-ary 3-cube Simulation Results . . . . .	215
E.56	Minimal Oblivious Routing (TOR) 4-ary 3-cube Simulation Results . . . . .	215
E.57	Minimal Oblivious Routing (FL) 4-ary 3-cube Simulation Results . . . . .	216
E.58	Minimal Adaptive Routing (NN) 4-ary 3-cube Simulation Results . . . . .	217
E.59	Minimal Adaptive Routing (UR) 4-ary 3-cube Simulation Results . . . . .	218
E.60	Minimal Adaptive Routing (BC) 4-ary 3-cube Simulation Results . . . . .	218
E.61	Minimal Adaptive Routing (TP) 4-ary 3-cube Simulation Results . . . . .	219
E.62	Minimal Adaptive Routing (TOR) 4-ary 3-cube Simulation Results . . . . .	219
E.63	Minimal Adaptive Routing (FL) 4-ary 3-cube Simulation Results . . . . .	220
E.64	Minimal Adaptive Routing (FL-HS) 4-ary 3-cube Simulation Results . . . . .	220
E.65	CQR Routing (NN) 4-ary 3-cube Simulation Results . . . . .	221
E.66	CQR Routing (UR) 4-ary 3-cube Simulation Results . . . . .	221
E.67	CQR Routing (BC) 4-ary 3-cube Simulation Results . . . . .	222
E.68	CQR Routing (TP) 4-ary 3-cube Simulation Results . . . . .	222
E.69	CQR Routing (TOR) 4-ary 3-cube Simulation Results . . . . .	223
E.70	CQR Routing (FL) 4-ary 3-cube Simulation Results . . . . .	223
E.71	CQR Routing (FL-HS) 4-ary 3-cube Simulation Results . . . . .	224
E.72	Enhanced CQR Routing (NN) 4-ary 3-cube Simulation Results . . . . .	225
E.73	Enhanced CQR Routing (UR) 4-ary 3-cube Simulation Results . . . . .	225
E.74	Enhanced CQR Routing (BC) 4-ary 3-cube Simulation Results . . . . .	226
E.75	Enhanced CQR Routing (TP) 4-ary 3-cube Simulation Results . . . . .	226
E.76	Enhanced CQR Routing (TOR) 4-ary 3-cube Simulation Results . . . . .	227
E.77	Enhanced CQR Routing (FL) 4-ary 3-cube Simulation Results . . . . .	227
E.78	Enhanced CQR Routing (FL-HS) 4-ary 3-cube Simulation Results . . . . .	228
E.79	Minimal Adaptive Routing (FL) 5-ary 3-cube Simulation Results . . . . .	229
E.80	Minimal Adaptive Routing (FL-HS) 5-ary 3-cube Simulation Results . . . . .	230
E.81	CQR Routing (FL) 5-ary 3-cube Simulation Results . . . . .	230
E.82	CQR Routing (FL-HS) 5-ary 3-cube Simulation Results . . . . .	231
E.83	Enhanced CQR Routing (FL) 5-ary 3-cube Simulation Results . . . . .	231
E.84	Enhanced CQR Routing (FL-HS) 5-ary 3-cube Simulation Results . . . . .	232

# List of Tables

2.1	Traffic Patterns Used to Test Routing Algorithms . . . . .	18
2.2	Minimal Adaptive and Chaos Routing Algorithms Compared by Throughput	19
2.3	Minimal Adaptive and Chaos Routing Algorithms Compared by Latency . .	19
2.4	Adaptive Routing Algorithms Compared by Throughput and Latency . . . .	19
4.1	Linpack Test Results of Various Systems . . . . .	41
4.2	Linpack/HPL Test Results of Thor's Tack Hammer . . . . .	42
4.3	Power Consumption Results of Thor's Tack Hammer . . . . .	42
5.1	Thor's Tack Hammer Results - Dimension Ordered Routing . . . . .	52
5.2	Thor's Tack Hammer Results - Direction Ordered Routing . . . . .	52
5.3	Thor's Tack Hammer Results - Minimal Oblivious . . . . .	53
5.4	Thor's Tack Hammer Results - Minimal Adaptive . . . . .	53
5.5	Thor's Tack Hammer Results - CQR Routing . . . . .	54
5.6	Thor's Tack Hammer Results - Enhanced CQR . . . . .	55
5.7	Comparison of All Laboratory Results . . . . .	58
6.1	Simulation Results - 3-ary 3-cube - Average Utilizations . . . . .	61
6.2	Simulation Results - 3-ary 3-cube - Standard Deviations . . . . .	61
6.3	Simulation Results - 4-ary 3-cube - Average Utilizations . . . . .	64
6.4	Simulation Results - 4-ary 3-cube - Standard Deviations . . . . .	64
6.5	Simulation Results - 5-ary 3-cube - Average Utilizations . . . . .	67
6.6	Simulation Results - 5-ary 3-cube - Standard Deviations . . . . .	67



# Acknowledgments

The success of this project could not have happened without the excellent supervision and support of Dr. Don Gruenbacher, who not only served as my major professor and advisor, but also encouraged me and inspired me to strive towards my goals. His selfless attitude and dedication towards academics is definitely appreciated.

I would also like to take this opportunity to thank my other advising professors, Dr. Caterina Scoglio and Dr. Sanjoy Das. Both of these individuals have given excellent guidance and support throughout my bachelor's and master's degrees.

# Dedication

To my family.

Most notably, my sisters, my closest friends, and my Grammy. You all are the light in my life, my motivation, and the unending source of support that I've come to greatly appreciate. I love you all so very much.

# Chapter 1

## Introduction

Computational systems in the world today are being limited by a different factor than they once were. In the past, things such as gate switching speed, memory size, and thermal properties have been (and will always be) issues when aiming towards faster and more efficient computers. But, aside from those things, communication (both between systems and within them) has become a tremendous issue. The study of interconnection networks, or the fabric that connects various nodes within a system, aims at finding answers to these complex communication questions [8]. Then within these interconnection networks, one must examine the efficiency behind the actual data exchange, or routing. Once the routing has been covered, flow control must then be integrated to work well with the routing algorithms on a particular interconnection network [16].

The rest of the chapter will outline the specifics of interconnection networks, and will dive into a specific use of these networks within high performance computing clusters. Furthermore, the routing (which is a core discussion of this research work) is briefly introduced as are some of the initial results stemming from the evaluation of the routing algorithms.

### 1.1 Interconnection Networks

The realm of interconnection networks ranges tremendously, as do their applications. These interconnection networks can be viewed both macroscopically and microscopically. For instance, by zooming out and looking at the macro-level, the entire Internet can be viewed

as an interconnection network with end-users and servers as the nodes, and the links, routers, and switches as the fabric connecting them. Different graphs or topologies of interconnection networks serve different purposes, and are very closely tied to their specific application. But within the different topologies there remains one quality which is central to all: balance the load as efficiently as possible, thus taking advantage of the diverse number of links, and increasing the overall ability for nodes to communicate.

At a more microscopic level, one can see that communication fabrics play a very large role in the efficiency of other parallel systems. For instance, the latest-and-greatest processors available through the popular manufacturers are called multi-core processors. These processors have multiple processing "cores" which attempt to compute or process data in parallel using a common communication fabric or interconnection network.

Though these two examples vary greatly in both the size of the communication fabric and the number of nodes within the network, both are solid examples of current applications of interconnection networks.

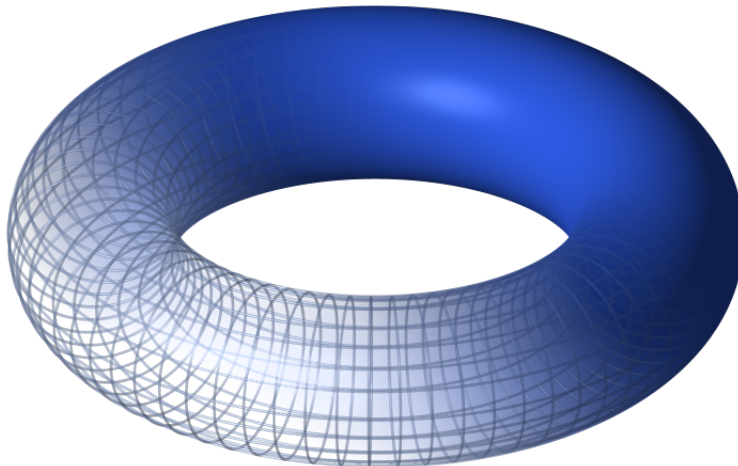
## 1.2 High Performance Computing Clusters

High Performance Computers (HPC), commonly referred to as "supercomputers" are generally used to process large simulations or tackle problems that were once considered impossible or improbable to solve in a non-parallel environment. HPCs are traditionally only found within companies or laboratories with large budgets because of their tremendous building and maintenance costs. The smaller clusters traditionally are thought of as computing clusters, as the HPC term tends to refer to the largest of these parallel systems. The largest (and fastest) 500 HPCs are tracked and cataloged within the Top500.Org website.

The applications on HPCs can and do vary widely. For example, corporations who manufacture automobiles could use an HPC or smaller cluster to effectively simulate prototype crash tests without having to physically build a prototype (which could only be good for one crash test). Instead, with a similar initial investment, these automotive corporations can

crash these simulated prototypes as many times as they desire. They can also collect massive amounts of data from these simulations, helping the company's efficiency and vehicle safety.

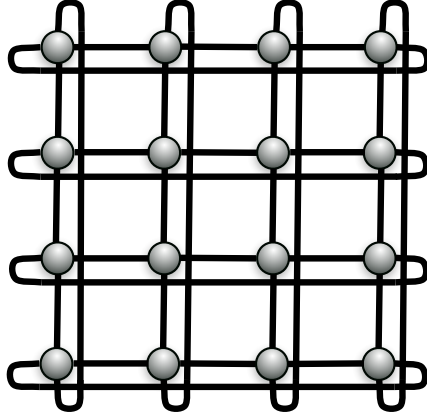
Within HPCs, and more specifically, massively parallel processing computers (MPP), there exist many different ways to connect the nodes: butterfly networks, torus networks, non-blocking and blocking networks, among others. As intuition would state, the diverse number of interconnection possibilities matches the diverse applications for these HPCs. The work focused here centers on the network topology and interconnection fabric that is most popular within the HPC systems currently used by the US Department of Energy (US DOE) [4] [10]. This is a torus-mesh configuration.



**Figure 1.1:** *Here is a torus shape, showing lines where the individual slices may reside. (Graphic public domain, <http://en.wikipedia.org/wiki/Image:Torus.png>.)*

The best way to explain this particular topology is to first consider a torus. A graphic of this shape is shown in Figure 1.1, but in layman's terms, its shape is often described as a product of two circles. Those familiar with tasty pastries like to compare it to a doughnut.

This network is composed of slices, connected in what one would consider in an  $X$ - $Y$  plane. For the example in Figure 1.2, there are four nodes both in the  $X$  and  $Y$  directions.



**Figure 1.2:** *This figure shows a 4-ary 2-cube slice of a tori-mesh network. These slices connect to other slices in an up/down (or +/- Z) direction to complete the torus shape.*

These nodes are also connected to other nodes in a  $Z$  direction. These slices continue to connect to other slices in the positive and negative  $Z$  directions and finally loop back around to the first, creating the torus shape. The defining characteristic that separates Tori-Mesh networks from that of hypercube networks is the irregular ratio between the number of slices in the  $Z$  direction and the number of nodes in the  $X$  and  $Y$  directions. If all three values are equal then that topology is a hypercube.

The Torus-Mesh interconnection network has many characteristics that make it attractive to the HPC community, and US DOE in general. One such characteristic is the high path diversity. Because of its configuration, there are generally many equivalent shortest paths between a single source and destination [8]. This improves a particular job or simulation's resiliency in being completed even when failures or hotspots occur. Resiliency is a quality which is of high importance within the US DOE on these large HPC systems [6].

Another quality that makes this particular network appealing to HPCs is the individual node's ability to quickly communicate with its neighbors with very little delay [8]. Most of the time when running simulations, nodes need only to communicate with their direct neighbors.

Lastly, and quite possibly most important, because of the tremendous investment in

creating these systems, tori-mesh networks are generally easy to expand.

A very critical part of creating an efficient HPC is assigning efficient routing policies. This next section discusses just that.

### 1.3 Routing within Interconnection Networks

There are quite a few categories that routing algorithms in general fit into. Routing algorithms can be source or hop-based in that the decisions that data traverses can be determined at the source or intermittently through the network. Routing algorithms may also be dynamic or static in that the paths do or don't change throughout time. All routing algorithms within this research are hop-based and includes both static and dynamic implementations.

Routing within a Torus-Mesh network can be viewed within three separate categories. The first two are the dynamic algorithms, which causes the individual nodes (which traditionally also do the routing for other nodes' data in transit) to make decisions on-the-fly as to a packet's next hop. These decisions can be done *obliviously* (oblivious routing), or *adaptively* by examining various indicators of network availability (adaptive routing).

The last option for routing is referred to as a static approach. In this particular option, the data from source to destination travels a single path and usually cannot deviate from that path.

There are a couple of reasons why routing on this network (or any network for that matter) is so critical. While a network's topology is an important factor in the efficiency of communications, the routing mechanism is the limiting reagent in this mix: an efficient routing policy can make for a very efficient HPC, while a poor policy can make it very inefficient. Alongside with efficient communications, there is a need to evenly balance load between links and nodes within a network [1]. Again, a good routing policy will achieve this.

There are two large HPCs that will be frequently referenced within this work. The first, Sandia National Laboratory's Red Storm/Thor's Hammer MPP currently has 26,569

processing nodes, each with a 2.4 GHz dual-core processor and all nodes connected in a tori-mesh interconnect network. This HPC is ranked as the 6th highest performance computer in the world based upon the Linpack Benchmark (as of Nov 2007). This benchmark is a globally recognized standard for cataloging HPCs and is discussed later in this thesis.

Another well-known parallel computer within the HPC community is Lawrence Livermore National Laboratory's BlueGene/L MPP. This parallel system currently has 212,992 nodes, each with a PowerPC 440 700MHz processor. This system was ranked as the highest performing parallel computer in the world (as of Nov 2007).

Both rankings were taken from the independent group called Top500.org, which catalogs and tests the high-performance computers who wish to be ranked.

The Red Storm supercomputer at Sandia National Laboratories, as well as Blue Gene/L at Lawrence Livermore National Laboratories can implement any routing algorithm appropriate for their topologies, but research has pointed towards the stability of static routing, specifically an algorithm called Dimension Ordered Routing (DOR) [1] [2] [4]. Because of the complexity of these systems, along with the laboratories' experience with dynamic routing algorithms, stability tends to be one of the deciding factors when considering which algorithms to use.

Flow-control is one of the last pieces of the puzzle in terms of creating an overall efficient HPC. Unfortunately, the scope of this research work does not extend into that area.

## 1.4 Contributions

The work contained within this section briefly discusses and introduces the specific work of this research thesis and the contributions that accompany them.

Adaptive routing algorithms give a large amount of flexibility in two key areas of efficient HPC usage: proper load balancing, and the effective usage of a network's interconnect (thusly exploiting locality). These routing algorithms, which are introduced and discussed later in this thesis, have shown excellent characteristics for effective operation on hypercube



and torus-mesh networks (among others).

An HPC system's overall reliability and availability is paramount in the efficient operation of large HPC systems [6]. Hotspots, which are naturally-occurring and unavoidable issues, occur when higher-than-average traffic demands are aimed at a single node. Multiple hotspots can occur within a network at any given time, and adaptive algorithms are the only algorithms which can adapt and respond to those traffic problems [5]. Therefore, this work aimed at increasing a routing algorithm's ability to make better decisions and to more effectively route around hotspots.

Along with introducing an enhanced decision function within these routing algorithms, a research cluster and simulation environment were created to verify those ideas. Each are introduced and explained in depth within this work.

An itemized list of contributions for this work include:

- Implement a simulation-based environment within Matlab to compare traffic and delay results for each algorithm given various hypercube topologies;
- Build a research cluster: a scaled version of the Red Storm/Thor's Hammer HPC;
- Implement a torus-mesh network on this research cluster using a USB interconnect;
- Implement various static, oblivious, and adaptive algorithms using a set of socket-layer C programs, validating and verifying previous work, and demonstrating this work's implementation differences;
- Introduce Enhanced CQR with Periphery Avoidance, which makes intermediate routing decisions by routing productively but avoiding the perimeter when possible;
- Discuss the impact of this new enhancement and demonstrate evidence of its potential.

# Chapter 2

## Previous Work

This chapter will attempt to provide a well-rounded literature review of the focused area of this thesis. Next is a discussion on undeliverable data – a critical issue with routing algorithms. Finally there is a discussion of both static and dynamic routing algorithms as well as evidence of their reported performance.

### 2.1 Undeliverable Data

With the high path diversity and the large number of routing decisions within tori-mesh networks, these issues of livelock and deadlock are of great concern and must be considered within each specific routing algorithm. These next two sections briefly introduce the problem of undeliverable data.

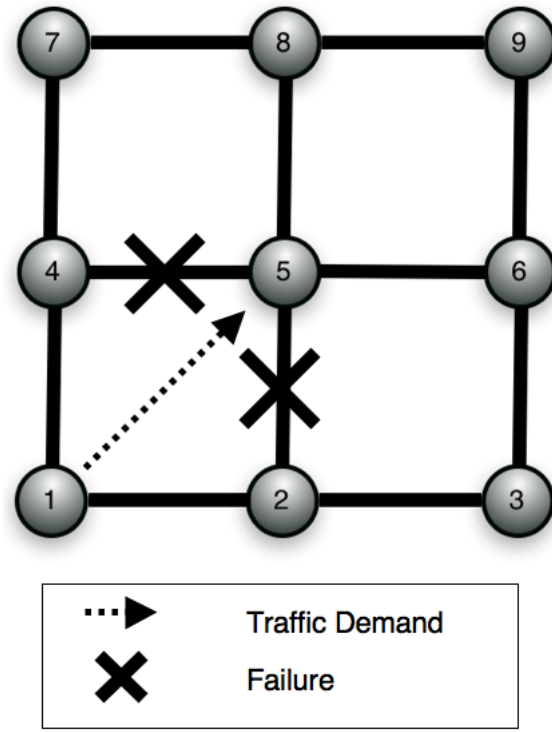
#### 2.1.1 Deadlock

Deadlock, as it is defined by [8], occurs when a set of agents holding resources are waiting on another set of resources such that a cycle of waiting agents is formed. Most networks are designed to avoid deadlock, but it is also possible to recover from deadlock by detecting and breaking cyclic wait-for relationships.

Deadlock is more of an issue on routing algorithms that do not allow for misrouting of data. When this data reaches a node or link it cannot traverse through, it stops its progress, and is usually dropped [9]. Because of finite buffers, and random failures, deadlock

is usually an issue to be considered within routing algorithms. This particular issue can also be subdued by implementing a few policies. As outlined within [9], deadlock prevention, avoidance, and recovery techniques can all reduce the problem of deadlock.

Figure 2.1 shows a deadlock example.



**Figure 2.1:** *This example shows how deadlock occurs, using the same traffic demand as before but this data is unable to reach the destination because of its inability to misroute. This data will end up being dropped at node 4 or node 2 after being unable to reach the destination.*

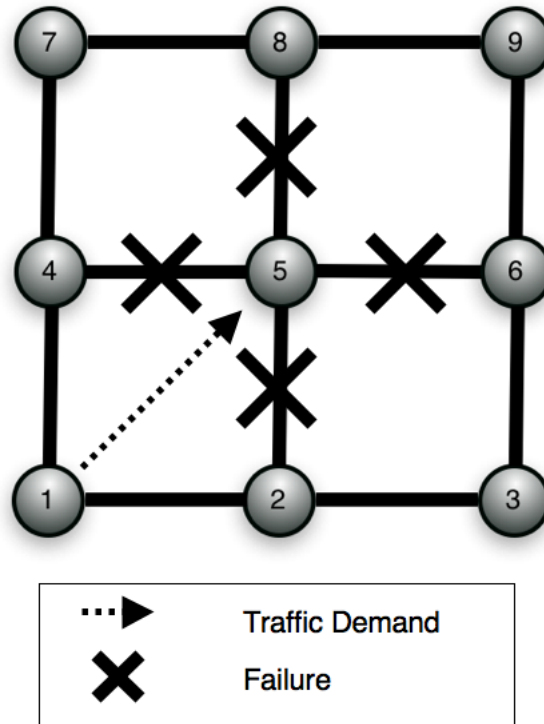
### 2.1.2 Livelock

Livelock, as it is defined within [8], occurs when a packet is not able to make progress in the network and is never delivered to its destination. Unlike Deadlock, though, a livelock packet continues to move through the network.

Livelock tends to occur most frequently when data is allowed to be misrouted, or routed non-minimally [9]. More specifically, livelock becomes a reality when data is forever in

transit towards its destination, and continues to move within the network, but does not reach the destination within finite time. Usually this issue comes about because of over-utilized channels or link failures. There are two schemes towards tackling the issue of livelock within routing algorithms. By using only minimal paths, and having great restrictions on non-minimal paths, livelock can be less of an issue [9].

Figure 2.2 shows a livelock example.



**Figure 2.2:** A simple livelock example showing a traffic demand from node 1 to node 5. The traffic is unable to reach the destination, regardless of its ability to misroute up to node 8, or node 6. This data will end up being misrouted indefinitely around the destination.

### 2.1.3 More on Undeliverable Data

Deadlock and livelock are two issues playing tug-of-war with one another. By reducing an algorithm's likelihood of falling into deadlock, it increases its likelihood of being vulnerable to livelock, and vice versa. These inversely-related issues are of great importance and will

be part of the discussion of each algorithm listed within this chapter.

## 2.2 Static Routing Algorithms

As it was stated before there is some appeal to static routing algorithms within these large parallel systems. Because of their simplicity, these algorithms tend to be very easy to implement within hardware [8]. As the name would imply, the path that data takes from a source to destination will not change through time. The difference between the algorithms depends on how these static routes are chosen. But once chosen, the path from source to destination cannot change.

While static, or source-based routing may sound simple, it is really only efficient when the administrator or operator knows a general idea of traffic patterns prior to any operations. Under those ideal conditions, source-based routing can very effectively balance traffic throughout the network. In terms of simulations, it would not be too terribly difficult to determine these traffic patterns prior to starting the simulations - thus, these routing techniques seem to be a pretty good choice on these HPCs.

As always, there are never any free lunches. Assuming that the operator has correctly predicted a traffic pattern and has generated all routing tables for a given set of nodes which will be used for a simulation, all notion of stability goes out the window when one of two things occur: (1) a node loses its ability to communicate with a set of nodes, or (2) a link becomes over- utilized, rendering it unusable. These two topics form the foundation of this research which are referred to as nodal failures and hot-spots, respectively.

Both of the following static approaches were initially developed because of their ease and simplicity in terms of hardware implementation [9]. It's no wonder why they were initially attractive to large HPC systems.

### 2.2.1 Dimension Ordered Routing, DOR

The most published application of static or source-based routing algorithms within these DOE HPC systems is what is called Dimension Ordered Routing (DOR) [1] [8] [9]. At the source, the node pre-computes a set of preferred directions to route the packet within. For instance, if a source and destination pair exist within the same  $Z$  dimension, only differing within the  $X$  and  $Y$  dimensions, it would be preferred not to change coordinates within the  $Z$  dimension [8] [9].

Once the preferred directions have been calculated, this information is placed within the header of the data and it is directed into the corresponding  $X$  direction first. Once the packet has successfully reached the same  $X$  coordinate as the destination, it then routes in the  $Y$  direction (if it needs to). Following the  $Y$  direction, intermediate nodes then continue to route in the  $Z$  direction (again, if necessary) until it finally reaches the destination.

So, in the case above, if the pre-computed preferred directions were  $[+X, -Y, 0]$ , this packet would fully route in the positive  $X$  direction and then the negative  $Y$  direction last. In terms of all possible directions, this is the order that data is routed within the network:  $[+X, -X, +Y, -Y, +Z, -Z]$ .

Dimension Ordered Routing is proven to be free from deadlocks [7] by balancing the load on virtual channels. The issue shown in [7] demonstrates that by making this modification the structure of the network and the diversity of usable channels changes. This has shown to create non-uniform patterns. These virtual channels are only necessary on torus or wrap-around networks, so with regular mesh-type networks, the full benefits of DOR can be witnessed.

### 2.2.2 Direction Ordered Routing, DIR

Direction Ordered Routing (DIR) is very similar to that of Dimension Ordered Routing. The difference in these two resides in the order in which the routing occurs once the preferred directions have been calculated. These calculations are done just as they are with DOR,

but instead of fully routing within the  $X$  direction before progressing to the  $Y$  direction and so on, DIR routes positively within all directions first, then routes negatively within the directions last [8] [9]. The order that data is routed within the network is as follows:  $[+X, +Y, +Z, -X, -Y, -Z]$

The same issues with regard to deadlock apply to that of DIR as they do with DOR.

## 2.3 Adaptive Routing Algorithms

Adaptive routing algorithms seem to be where the most recent research has been heading in terms of these HPC systems. Though there is some stigma with regard to their overall stability within complex networks, it is the work of this research to try to show highlights of using these routing algorithms. The balance between intuition of network conditions and making routing decisions based on that intuition is definitely taken at the expense of simplicity.

### 2.3.1 Chaotic

Chaotic routing uses what is called deflection routing [8], randomly granting contending data access to channels and deflecting (or misrouting) the data that is denied access to the requested channel. These particular deflections may or may not be productive (meaning an extra hop may be incurred by deflection).

Deadlock is not as much of an issue for this algorithm, as the data can be easily misrouted [5] [8]. Livelock can be handled within this routing scheme by introducing time stamps and battle scars [5]. Time stamps indicate when the packets were introduced into the network, and when contending packets request the same channel, the router can grant the packet with the oldest packet. Battle scar implementations require packets to contain information as to how many times they have been misrouted. Again, the router can grant packets with larger battle scars access to channels over those with smaller or zero battle scars [5].

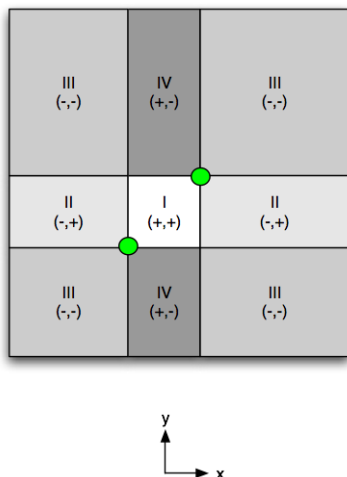
### 2.3.2 Minimally Adaptive

The minimal adaptive routing algorithm first starts the data progression by designating the minimal quadrant, which contains all possible shortest paths from source to destination. This quadrant is of size  $\prod^k \Delta k$ , where  $k$  is the number of dimensions which source and destination differ by greater than zero. Any node within this quadrant is considered to be a possible intermediate node along some of many possible shortest paths.

Figure 2.3.2 shows the an example on a  $n = 2$  torus network. Quadrant I is the most minimal because it contains all minimal paths. Once this minimal quadrant has been properly identified, all routing occurs progressively within it.

For instance, if  $s = (2, 3)$  and  $d = (5, 6)$ , a packet must route minimally within Quadrant I by either routing positively in the  $X$  direction or positively in the  $Y$  direction. The decision as to whether to pick the  $X$  or  $Y$  direction depends on local characteristics such as output queue length. Once a next hop has been determined, the data is then routed to the next hop, and the process repeats until it reaches the destination.

**Figure 2.3:** *This figure shows the possible quadrants and the distances data must traverse while within them. Quadrant I is the most minimal, if  $\Delta x$  and  $\Delta y$  are equal, then quadrants II and IV are the same. Quadrant III is the least minimal in this example. (Figure adapted from [12].)*





Once this quadrant has been determined, routing is done within it by choosing the next progressive hop by analyzing some local measurements or congestion indicators - typically output queue length [8].

Typically this particular algorithm fares well towards local load balancing, but has no insight of congestion further down the path, limiting its ability to achieve global load balancing. By using virtual-channels, deadlock can be avoided [3].

### 2.3.3 Fully Adaptive: GOAL, GAL, CQR

The essence of the best algorithms is epitomized within the characteristics of fully adaptive routing algorithms for these networks [13]. By exploiting locality and the network's wrap-around feature, an efficient routing algorithm should be able to have the flexibility to misroute data to balance the traffic load. Fully adaptive algorithms do just this.

The work that was done by Singh, et. al, [12] [13] [14], has collectively progressed into what seems like a solid attempt at finding a stable, reliable, and fully adaptive routing algorithm on tori-mesh networks for a diverse set of traffic patterns.

The first of the three algorithms, Globally Oblivious Adaptive Locally (GOAL), works very similarly to that of Minimal Adaptive with a few changes. First, it finds all possible routable quadrants and scales the probability of selecting a particular quadrant based on how minimum it is with respect to a particular source/destination pair. Once those probabilities are assigned and a particular quadrant has been selected, routing is done minimally within it [12]. This algorithm uses a virtual-channels implementation (called \*-channels) to ensure deadlock and cycle-free operation [12].

One problem with this algorithm, which was a similar problem with Minimal Adaptive, was its inability to recognize global congestion. The answer to this comes through in the next algorithm that was developed by this same group, Globally Adaptive Load-Balanced (GAL).

This algorithm is implemented by keeping a globally visible set of input queues for each

node, called injection queues. Each node has as many injection queues as it has inputs, and when a packet begins its commute from source to destination, the injection queue on the destination node relative to the link the packet will be received on is increased. As with GOAL, GAL first determines the routable quadrant, then selects this quadrant based on distance (as did GOAL), but it also considers the injection queues. This way, nodes can keep a general idea of how much total traffic is being sent in any particular direction, and can attempt to balance that out globally [14].

**Figure 2.4:** *A graphical representation of how GAL Injection queues appear globally. The value increments as packets are sent toward a particular interface, and they decrease as packets reach a destination through a particular interface.*

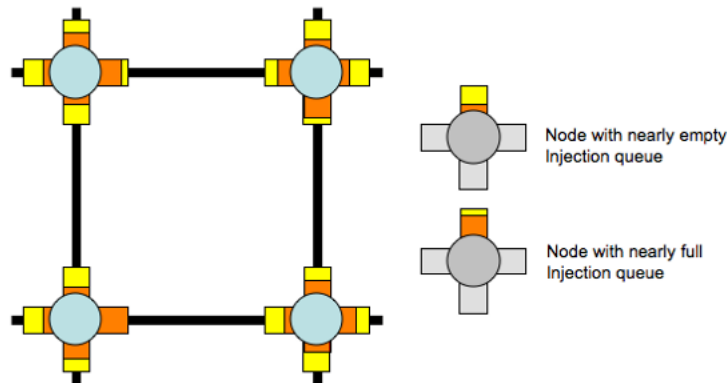


Figure 2.4 shows graphically how these injection queues may look.

Another versatile option for this algorithm to reduce its complexity is by subnetting the injection queues. By grouping multiple interfaces together with a single queue, thus reducing the number of queues to store globally, the authors showed that this technique had little impact on the overall performance of this algorithm [14]. This algorithm is also said to be deadlock and cycle-free [14].

Finally, a better meld of the global and local qualities that GAL and GOAL provided was integrated into a simple-to-implement and efficient protocol called Adaptive Channel Queue

Routing (CQR) [13]. This routing algorithm uses local information in the form of output queues (as opposed to the GAL algorithm which implemented input injection queues which were globally accessible). CQR uses these output queues as estimators for global congestion, which was proven within the work of [13]. The work showed that even in the absence of global information, local information can provide good estimations of global congestion.

Also, relieving the algorithm from accessing and changing global information made it much simpler to implement (not to mention lower its overhead). It was also shown to quickly adjust to changing traffic patterns and the changeover from minimal routing to non-minimal routing[13]. CQR was also shown to be deadlock and cycle-free [13].

## 2.4 Oblivious Routing Algorithms

This final section of routing algorithms within these HPC networks is referred to as Oblivious Routing. These routing algorithms tend to hybridize aspects of the adaptive and static algorithms. For instance, oblivious routing techniques tend to keep hardware implementations simpler than adaptive algorithms (a static characteristic), while the paths that data takes within the network can change throughout time (which is an adaptive characteristic).

The next subsections introduce two popular oblivious algorithms, Valliant's Algorithm and Minimal Oblivious.

### 2.4.1 Valiant's Algorithm

Valiant's Algorithm came about in an attempt to balance loads within any particular topology [8]. The algorithm works by first selecting at random an intermediate node,  $x$ , and once successfully routing to that intermediate node, then routing from  $x$  to the destination. As it's stated, any arbitrary routing algorithm can be used to get data from source to the intermediate node and then from the intermediate node to the destination [8].

But with the addition of better network load balance, locality is clearly taken out of the picture [12]. And as it was stated before, the DOE HPCs have a desire to maintain locality

Name	Description
Nearest Neighbor (NN)	Nodes send data to only their neighbors with equal probability.
Uniform Random (UR)	Nodes send data to random destinations.
Bit Compliment (BC)	$(x, y, z)$ sends data to $(k - x - 1, k - y - 1, k - z - 1)$ .
Transpose (TP)	$(x, y, z)$ sends to all possible permutations of $x, y$ , and $z$ .
Tornado (TOR)	$(x, y, z)$ sends to $(x + \frac{k}{2} - 1, y, z)$ .
Flood (FL)	$(x, y, z)$ sends to all other nodes within the network.
HotSpot (HS)	$n$ number of nodes have $k$ times more traffic demand than others.

**Table 2.1:** *These are the traffic patterns outlined in [12] [13] [14], which have been used to validate and verify routing algorithms on hypercube and torus mesh networks.*

because of the nature of the applications. This makes Valiant’s Algorithm a possibly good application in other networks, but not within Redstorm, BlueGene/L, or the like.

Deadlock is said to be avoided using two subnetworks for the two steps, and using two virtual channels within each step [15].

## 2.4.2 Minimal Oblivious

The Minimal Oblivious algorithm works very closely to that of the Minimal Adaptive algorithm listed above, in that the traffic will traverse one of many possible minimum shortest paths. The algorithm first calculates its minimal quadrant, and then selects a random intermediate node somewhere within the minimum quadrant. Once that intermediate node is selected, the data is routed first to the intermediate node and secondly from the intermediate node to the destination, all done randomly, minimally, and oblivious to any network indicators or local conditions [8].

## 2.5 Traffic Patterns

In order to analyze these above algorithms, there are a few traffic patterns that are used to test and validate an algorithm’s ability to route packets across a given topology effectively. These traffic patterns usually serve as a worst-case scenario for traffic demands, but they are real-world possibilities for demands none-the-less.

Traffic Pattern	Minimal Oblivious	Chaos Routing
UR	0.6	1.0
BC	0.3	0.5
TP	0.5	0.5
HS	0.6	0.9

**Table 2.2:** This table shows results obtained from [5], comparing the Chaos and Minimal Adaptive routing algorithms with a 100% load applied the simulated hypercube network. The results show the fraction of data successfully transmitted given the traffic pattern at 100% load.

Traffic Pattern	Minimal Oblivious	Chaos Routing
UR	600 (cycles)	550 (cycles)
BC	2200	1000
TP	690	900
HS	600	500

**Table 2.3:** This table shows the results obtained from the same conditions as in Table 2.2. This table shows the latency (in terms of extra cycles) seen when observing the throughput in the previous table.

## 2.6 Latency and Throughput Analyses

This section attempts to show previous work, comparing the algorithms of interest in terms of throughput and latency. All algorithms were compared using similar approaches, by analyzing how well the algorithms reacted given specified traffic demands.

Tables 2.2 and 2.3 show useful data taken from [5], comparing the latency and throughput

Traffic Pattern	Minimal Adaptive	GOAL	GAL	CQR
NN/UR	1.0	0.75	1.0	1.0
TR/BC/TOR	0.33	0.53	0.53	0.53
Average	0.63	0.67	0.73	0.7
HS	0.46	0.48	0.49	0.49
Low Load Latency	4.45	6.17	4.45	4.45

**Table 2.4:** This table compares the adaptive routing algorithms in terms of throughput (shown in fractions given 100% loads), as well as latency during low loads. This information was obtained from [13] [14] [12].

of Minimal Oblivious and Chaos routing algorithms. The results show Chaotic routing meeting or exceeding the performance of Minimal Oblivious.

Table 2.4 compares the adaptive algorithms in terms of throughput. All three of the Stanford group's algorithms show consistent progression over that of Minimal Adaptive. CQR ends up being the most resilient algorithm in terms of these chosen traffic demands when concerning throughput and latency.

# Chapter 3

## Enhancing Current Adaptive Routing Algorithms

Because adaptive algorithms are able to exploit locality as well as misrouting data to balance the load, they are very attractive algorithms for implementation. Their stability is necessary for traffic changes and CQR, Minimal Adaptive, and other fully adaptive algorithms alike are able to transition between those traffic differences well [13].

### 3.1 Issues with Previous Algorithms

As it was stated before, adaptive algorithms have two qualities which make them attractive. They have the ability to take advantage of locality while capitalizing on path diversity by balancing load. As for CQR, once a quadrant has been selected at the source, the data in transit is not allowed to deviate outside of that quadrant and must make productive moves towards the destination [13].

A productive move is defined as intermediately choosing a next hop so that the distance between the data in transit and the intended destination decreases [13]. If the distance stays the same or increases, it cannot be considered a productive hop.

Adaptive algorithms such as CQR could more effectively route around hotspots or nodal failures if they encountered these network abnormalities away from the edges of the selected quadrants. In fact for every "edge" of the quadrant the data hits, it reduces the possible

decisions it can make further down the path.

Take for instance a simple example where a selected quadrant has  $\Delta x = 4$  and  $\Delta y = 2$ . If the first two hops within that quadrant are in the  $y$ -direction, then the final four hops must be in the  $x$ -direction. If the data encounters a hotspot along the  $x$ -direction, it has no choice but to continue routing through the congestion. The next section introduces the concept further, called Periphery Avoidance.

## 3.2 Introducing Enhanced CQR with Periphery Avoidance

Along with the flexibility that adaptive routing algorithms offer, there are some drawbacks. Minimal Adaptive, for instance, offers great flexibility in routing decisions near the center of the selected quadrant, while offering less flexibility near the periphery of the quadrants. This is because once data has reached the periphery of its intended routing quadrant (it can also be thought of as aligning in a similar dimension with the destination), it reduces the ability to route around congestion and hotspots. By offering the traffic a path away from the periphery of the selected quadrant, we increase data's ability to adaptively route around hot-spots.

Analytically speaking, a majority of the paths within a quadrant fall along the periphery of the network, but many choices do exist within the interior of a quadrant. The next section discusses path diversity as it relates to this overall concept of periphery avoidance and introduces a decision function in which queue sizes are weighted with the distance from the periphery of a quadrant to choose a next-hop.

### 3.2.1 Path Diversity and Periphery Avoidance

The work from [8] outlines mathematical expressions for calculating the number of possible distinct routes within a torus-mesh or hypercube network. Below is the expression from [8] for a 2-Dimensional network:



$$|R_{sd}| = \binom{\Delta x + \Delta y}{\Delta x} \quad (3.1)$$

And for a 3-Dimensional network, the expression for finding the number of possible distinct routes between source and destination is:

$$|R_{sd}| = \binom{\Delta x + \Delta y + \Delta z}{\Delta x} \cdot \binom{\Delta y + \Delta z}{\Delta y} \quad (3.2)$$

Given the possible routes that exist between source and destination within a quadrant for the values for  $\Delta x$ ,  $\Delta y$ , and  $\Delta z$ , a subset of these routes use the interior of the network while another disjoint subset use one or more parts of the periphery. The next equation was created to demonstrate that relationship:

$$|R_{sd}| = |R_{sd\bar{p}}| + |R_{sdp}| \quad (3.3)$$

The expression for  $|R_{sdp}|$  describes the routes which use the periphery during one or more hops, while  $|R_{sd\bar{p}}|$  describes the routes which route only within the interior of the quadrant.

It was easier to think in terms of counting the paths which use the interior of the quadrant than developing an expression for counting the number of paths which use the periphery of the quadrant. For a 2-Dimensional network, that expression is:

$$|R_{sd\bar{p}}| = \binom{(\Delta x - 1) + (\Delta y - 1)}{(\Delta x - 1)} + 2 \cdot n \quad (3.4)$$

And for a 3-Dimensional network the expression is:

$$|R_{sd\bar{p}}| = \binom{(\Delta x - 1) + (\Delta y - 1) + (\Delta z - 1)}{(\Delta x - 1)} \cdot \binom{(\Delta y - 1) + (\Delta z - 1)}{(\Delta y - 1)} + 2 \cdot n \quad (3.5)$$

An example 2-D network with a quadrant of size  $\Delta x = 4$  and  $\Delta y = 3$  demonstrates the relationship between these values. Therefore, given those values and using Equations 3.3 and 3.4, the counts for the various routes are:

$$\begin{aligned} |R_{sd}| &= \binom{4+3}{4} = 35 \\ |R_{sd\bar{p}}| &= \binom{3+2}{3} + 2 \cdot 2 = 14 \\ &\vdots \\ |R_{sdp}| &= 21 \end{aligned}$$

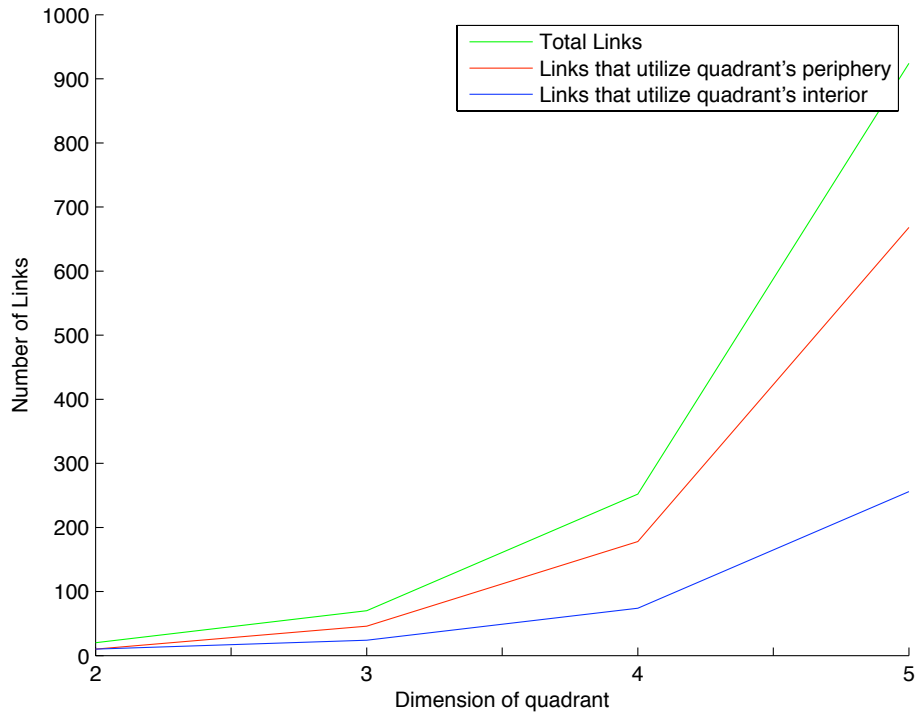
This example depicts the general trend of the relationship between  $|R_{sd}|$ ,  $|R_{sd\bar{p}}|$ , and  $|R_{sdp}|$ . Because of the high number of links which utilize areas of the periphery, and given that hot-spots are common occurrences avoiding these links whenever possible should definitely improve the data's ability to route around localized hotspots when encountering them.

Figure 3.1 shows the relationship between the number of links which utilize the periphery and the number of links that do not.

### 3.2.2 Periphery Avoidance Decision Function

Traditional CQR uses two decision phases when sending data through the network. The first phase is done at the source, when the quadrant to route within is selected. Once this decision has been made, the data carries with it a vector of size  $n$  (e.g.  $v = \{x_1, x_2, \dots, x_n\}$ ), referred to as the minimal direction vector. The values within this vector indicate whether routing is allowed within a given dimension. If a particular dimension carries a value of  $+1$ , then the data may route in the positive direction in that dimension. Similarly, if the value is  $-1$ , the data may route in the negative direction of that dimension. Finally, if the value is  $0$ , routing is not allowed in that dimension.

For this discussion, the minimal direction vector must not only contain the allowed directions to route within that dimension, but must also contain the  $\Delta$  values for that dimension. This enables the decision function, which is explained below, to weight the links near the center of the quadrant higher than those which reside on the periphery. As the



**Figure 3.1:** *This figure demonstrates the ratio between the links which utilize the periphery of a selected quadrant to that of the links which utilize the interior of the quadrant. As the quadrant size increases, the ratio converges to 4:1.*

data progresses through the network and decides a next hop, this vector must be modified to reflect the actual  $\Delta$ 's at each intermediate hop. This enables the decision function to constantly evaluate and aim towards keeping the data in transit away from the periphery of the network.

The Enhanced CQR decision function as it relates to the concept of Periphery Avoidance uses the same first phase as traditional CQR, but uses a slightly different method for choosing intermediate nodes within a network on the second phase. Before considering this enhanced version, first consider the old decision function for CQR:

$$Next\_Hop = \min(Q(x_1), Q(x_2), \dots, Q(x_n)) \forall x_i \in n \quad (3.6)$$

Where  $Q(x_i)$  is:

$$Q(x_i) := k \text{ s.t. } \begin{cases} k = \text{Output queue value at direction } x_i \text{ for } \Delta x_i > 0 \\ k = +\infty \text{ for } \Delta x_i = 0 \end{cases} \quad (3.7)$$

The function  $Q(x_i)$  returns  $+\infty$  when  $x_i$  is not a productive direction, otherwise, it returns the value of the output queue in the dimension  $x_i$ . Therefore, a selection is made to ensure the next hop is productive, and relative to the shortest output queues.

To create a decision function which incorporates both output queue length and periphery avoidance, the following equation was developed:

$$Next\_Hop = \min(Q_p(x_1), Q_p(x_2), \dots, Q_p(x_n)) \forall x_i \in n \quad (3.8)$$

Where  $Q_p(x_i)$  is:

$$Q_p(x_i) := k \text{ s.t. } \begin{cases} k = (Q(x_i) + 1) \cdot (1 - \frac{\Delta_i}{\Delta_{total}}) \text{ for } \Delta x_i > 0 \\ k = +\infty \text{ for } \Delta x_i = 0 \end{cases} \quad (3.9)$$

The value for  $\Delta_{total}$  is the sum of all the absolute values in each dimension. One is added to the queue lengths to ensure that zero-length queues do not compromise the function's ability to avoid the periphery.

First, consider a selected quadrant with direction vector  $v = \{\Delta x = 4, \Delta y = 10, \Delta z = 0\}$ . This quadrant yields a total of 1001 unique minimum paths from source to destination as per Equation 3.1. When considering the number of paths which avoids the periphery, that value is found to be 220 as per Equation 3.4. Analyzing the decision of the next hop in terms of traditional CQR, the  $x$  and  $y$  are then selected based on the smallest output queue length.

But, this is where Periphery Avoidance can take advantage of having such a difference in the  $\Delta$ 's. Of course, output queues need to be considered, as they are an indicator of

network congestion, but routing in the  $y$  direction should be preferred as it avoids the periphery better than choosing the  $x$  direction for this example.

Therefore, for the enhanced CQR the values obtained from the decision function are:

$$Q_p(x) = (Q(x) + 1) \cdot 0.714$$

$$Q_p(y) = (Q(y) + 1) \cdot 0.285$$

Here, it's demonstrated that if queue sizes are equivalent, the  $y$  direction will be selected as the minimum value, and the next hop.

Considering another example which depicts how the decision function factors on deciding a next hop based on a quadrant with  $v = \{\Delta x = 4, \Delta y = 4, \Delta z = 0\}$ , with equivalent output queues (as with the previous example), the number of paths between source and destination are 70, with 20 paths within the interior of the quadrant.

Therefore the decision function yields:

$$Q_p(x) = (Q(x) + 1) \cdot 0.500$$

$$Q_p(y) = (Q(y) + 1) \cdot 0.500$$

This causes the output queue values to be the sole deciding factor in this scenario.

These examples demonstrate the basic functionality of this decision function and its ability to weight both output queue lengths and routing towards the interior of the quadrant – both of which are important in load balancing and avoiding hotspots within the network.

# Chapter 4

## The Research Cluster: Thor's Tack Hammer

### 4.1 Topology

This section describes the topology of the KSU research cluster, Thor's Tack Hammer. The topology of this cluster is one of the characteristics of this parallel system that helps explain its overall performance, and operability.

In an attempt to keep the individual nodes diskless (or operate without the use of individual hard drives), it was necessary to have a single central server connected to each node, providing PXE-booting and NFS file sharing. This server, named Sandlab, provided those services to the nodes of Thor's Tack Hammer, as well as providing the cluster with a firewall from the Internet.

Along with this management network exists the USB interconnection network. This network uses USB patch cables to connect the nodes in the hypercube or torus-mesh topology. During execution of code, the particular traffic demands and routing decisions are done by routing packets within the USB interconnect network, and cannot use the management network. Both networks use separate subnets, the management network using 192.168.0.0/24 and the USB interconnect using 10.0.0.0/16.

Figure 4.1.1 shows a graphical representation of the network topology of Thor's Tack Hammer.

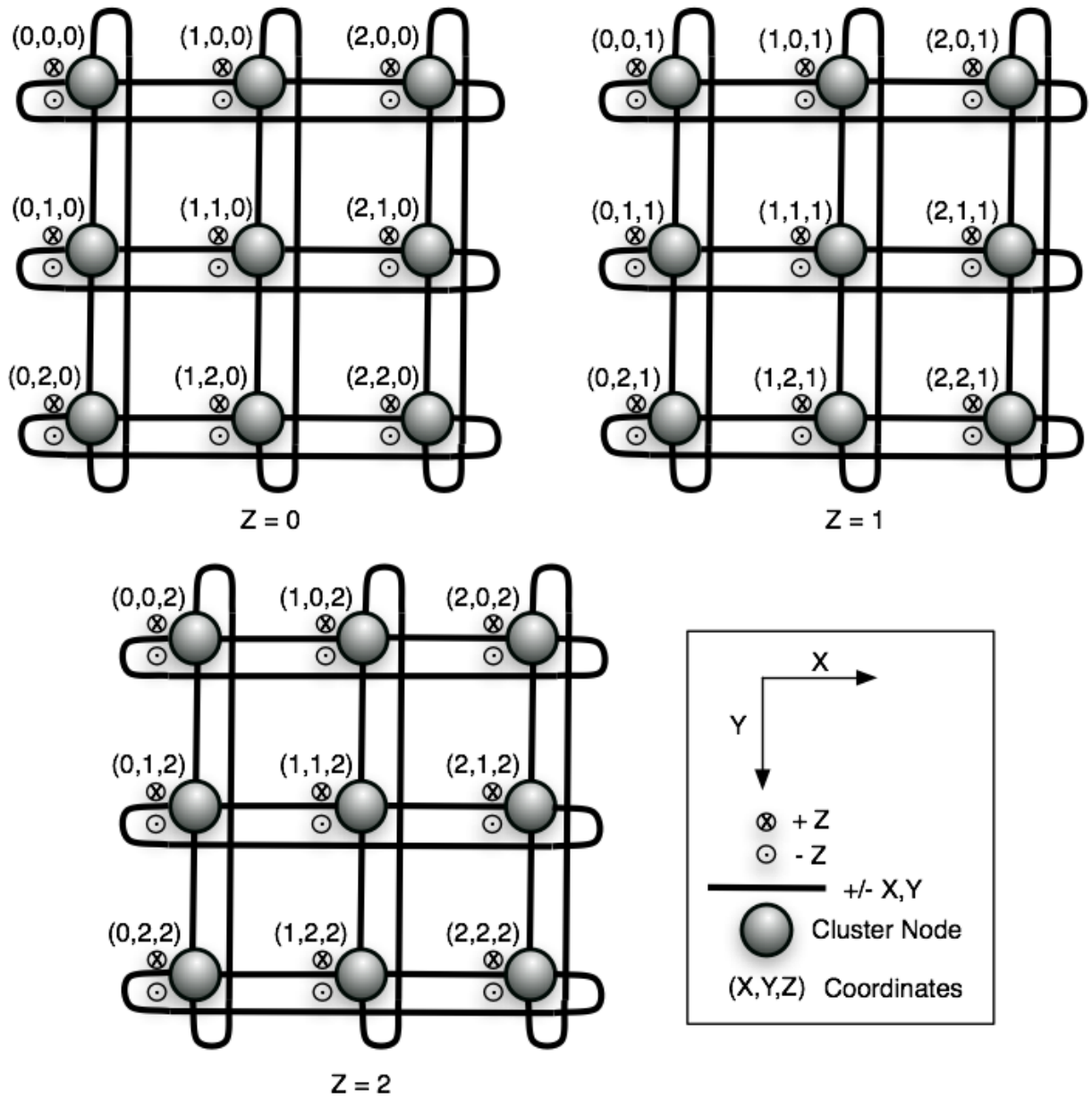


Figure 4.1: This figure shows the connectivity between nodes in Thor's Tack Hammer.

### 4.1.1 Subnetting and Addressing Scheme

Because of the version of Linux that was used, as well as the hardware chosen for the interconnection, each USB interface was able to be used like a traditional ethernet interface. Once the Linux kernel was brought up, the USB interfaces could communicate through the USB patch cables by IP.

This was one of the reasons for choosing these particular patch cables for the interconnect. They were fairly inexpensive (approx. \$5-10US per cable), and Linux had drivers for easily transmitting data across the cables.

Subnetting between the nodes was implemented to make communication between a node and its neighbors easy. A subnetting scheme was initially developed by a colleague at Sandia National Laboratories, and is explained below. By creating exclusive subnets for each USB link, it could be asserted that data would be transmitted along the USB interconnect and not the management network. This same scheme was used throughout the research as the standard for connectivity between nodes.

There are first two things to consider when understanding this subnetting scheme. There are slightly different algorithms for calculating a node's local interface addresses and the neighbors' interface addresses. The first part of this discussion explains the calculation of a node's local address (given in x,y,z form - see Figure 4.1).

The first two octets of the addresses can be arbitrarily chosen (x.x.C.D). For Thor's Tack Hammer, '10' and '0' were chosen as the first and second octets respectively. The last two octets are partitioned into 4-bit values, the x-value, the y-value, the z-value, and the direction. Therefore, the third octet contains the x-value and y-value (4 bits + 4 bits), and the fourth octet contains the z-value and the direction (again, 4 bits + 4 bits). Appendix A gives the BASH Linux scripts that were developed, and the code for generating these values is available for reference there. The pseudo-code for calculating a node's the third octet for a positive-x interface is as follows:



$$C = (x\text{-value} \ll 4) + y\text{-value};$$

*(alternatively)*

$$C = (x\text{-value} * 16) + y\text{-value};$$

Whereas, calculating an address for the third octet in a negative-x direction is as follows:

$$C = ((x\text{-value}-1) \ll 4) + y\text{-value};$$

*(alternatively)*

$$C = ((x\text{-value}-1) * 16) + y\text{-value};$$

And calculating the third octet for a negative-y interface:

$$C = (x\text{-value} \ll 4) + (y\text{-value}-1);$$

*(alternatively)*

$$C = (x\text{-value} * 16) + (y\text{-value}-1);$$

Finally for the third octet, the wrap-around must be accounted for since the interconnect is that of a hypercube/torus-mesh. Therefore, a test is done to see if the x-value or y-value is 0. If that is the case, and it is a negative move, the following calculations are done for the negative-x interface:

$$C = ((k - 1) \ll 4) + y\text{-value};$$

*(alternatively)*

$$C = ((k - 1) * 16) + y\text{-value};$$

And similarly for the negative-y interface:

$$C = (x\text{-value} \ll 4) + (k - 1);$$

*(alternatively)*

$$C = (x\text{-value} * 16) + (k - 1);$$

The calculations for the fourth octet is similar, but static values for the direction must first be considered. They are as follows:

Positive	X:	dir=1	(0x0001)
Negative	X:	dir=2	(0x0010)
Positive	Y:	dir=5	(0x0101)
Negative	Y:	dir=6	(0x0110)
Positive	Z:	dir=9	(0x1001)
Negative	Z:	dir=10	(0x1010)

These values were chosen to allow for 4 IP addresses per subnet. The values not shown 0 (0x0000), 3 (0x0011), etc., represent the network and broadcast addresses per subnet, which are illegal addresses for the USB interfaces.

Having the appropriate direction values allows for the calculation of the fourth octet. Again, positive and negative directions will alter the way the addresses are calculated (this time with z-value as they were before with x-value and y-value). For positive-z interface the following calculation finds the appropriate octet:

$$D = (z\text{-value} \ll 4) + \text{dir};$$

And now for the negative-z interface:

$$D = ((z\text{-value}-1) \ll 4) + \text{dir};$$

Finally, when the z-value is equal to 0 and a negative-z direction is necessary the calculation must account for the network's wrap-around:

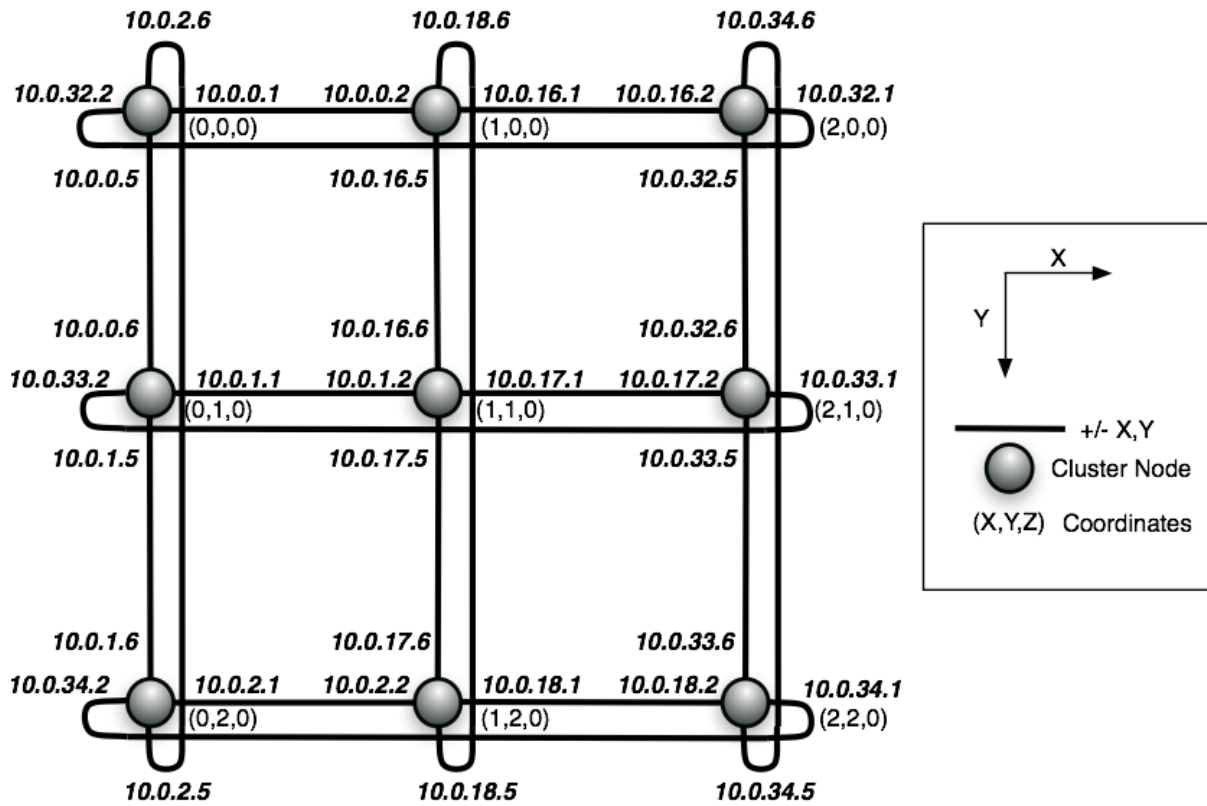
$$D = ((k - 1) \ll 4) + \text{dir};$$

In order to allow 4 addresses per subnet as explained above, the interfaces were assigned their 10.0.C.D/30 address and subnet. Calculating the neighbors is done similarly, except the direction values are toggled. Therefore the direction table shown above would change to:

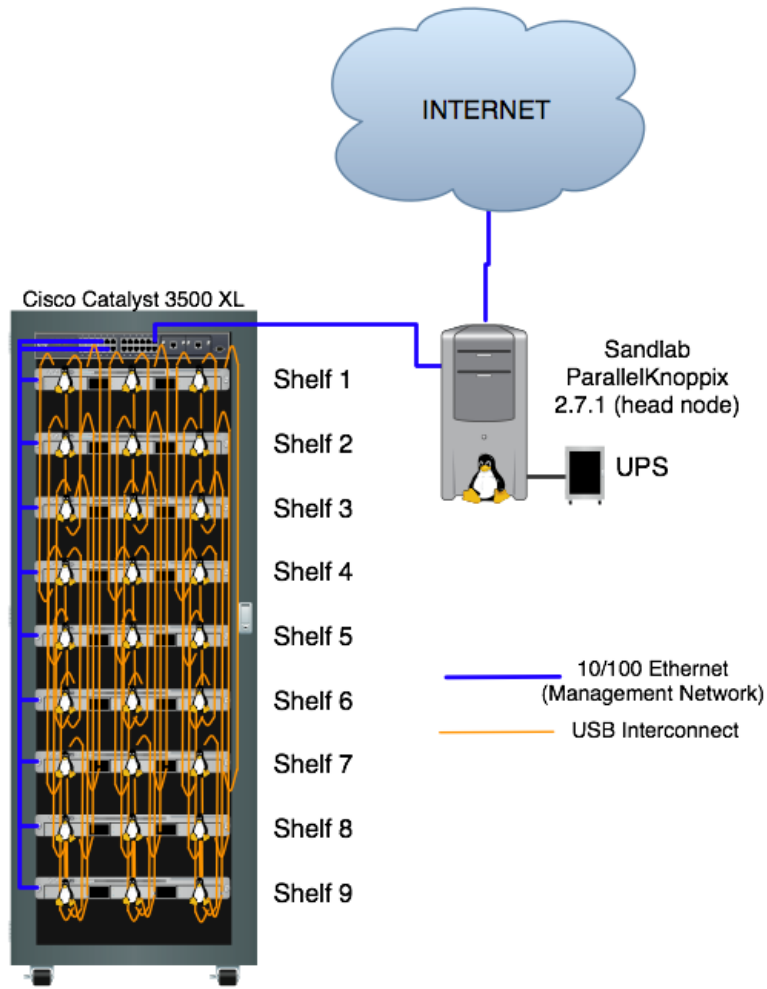
Positive	X:	dir=2	(0x0010)
Negative	X:	dir=1	(0x0001)
Positive	Y:	dir=6	(0x0110)
Negative	Y:	dir=5	(0x0101)
Positive	Z:	dir=10	(0x1010)
Negative	Z:	dir=9	(0x1001)

All other calculations for C and D are done exactly the same as outlined before.

A figure has also been provided giving an example of the results of these calculations. See Figure 4.1.1. This figure shows an entire Z-plane of assigned address, but only for the X and Y interfaces. The Z interfaces were not shown, though they are calculated just as described previously. This example should demonstrate the subnetting and addressing scheme implemented on Thor's Tack Hammer.



**Figure 4.2:** This figure demonstrates the addressing and subnetting scheme used on Thor's Tack Hammer's USB interconnect network. Though the Z interfaces are not shown, they are calculated just as described previously.

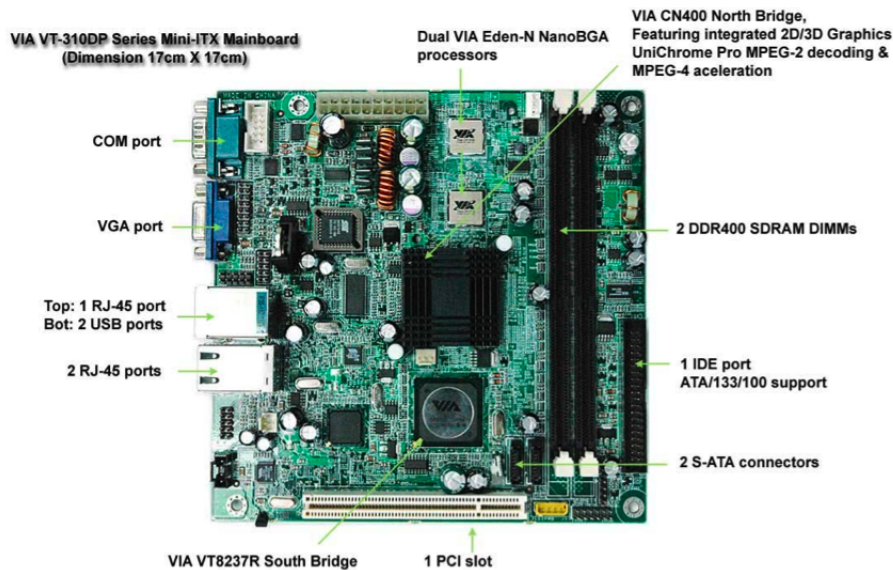


**Figure 4.3:** This figure shows the network topology of Thor's Tack Hammer. The management network (shown in blue) is in a star-type topology and the USB interconnect network (shown in orange) is in a hypercube-type topology.

## 4.2 Node and Cluster Specifications

Thor's Tack Hammer is composed of 27 nodes, each node being diskless and headless (without individual hard drives and monitors). Each node is a VIA-VT310-DP Mini-ITX (17 cm x 17 cm) Motherboard, complete with the following features:

- ◇ Dual 1GHz Via Eden<sup>TM</sup>-N Nano BGA processors
- ◇ 1GB Kingston KVR400X64C3AK2 DDR RAM
- ◇ 133 MHz front side bus
- ◇ Onboard Intel i82551QM 10/100Mbs Ethernet Adapter (used as Management Interface)
- ◇ Onboard VIA VT6103L 10/100Mbs Ethernet Adapter (unused)
- ◇ Onboard VIA VT6122 10/100/1000Mbs Ethernet Adapter (unused)
- ◇ VIA CN400 North Bridge
- ◇ VT8237R South Bridge
- ◇ Two onboard USB ports (2.0)
- ◇ NEC PCI USB card (2.0)
- ◇ Custom Award BIOS
- ◇ 27.7 Watts Total Power Consumption



**Figure 4.4:** This figure shows the VIA VT310-DP, the motherboard used for each node of Thor's Tack Hammer. (Image borrowed from VIA's site: <http://tinyurl.com/2t3lkg>.)

Each node resides with two other nodes (three total) per shelf within the rack, and all three nodes share a common power supply. To enable all nodes to come up simultaneously, wake-on-lan (WOL) was used to switch on the nodes from their off state. Simply turning off the power supply once the nodes have halted was the process for powering down the nodes.

A customized BIOS was needed to allow nodes to PXE-Boot and wake-on-lan on the same interface. This interface ended up being the Intel i82551QM, as opposed to the VIA



**Figure 4.5:** *This image is a photograph of the Adaptec USB 2.0 AUA-5100 PCI card used to increase the number of USB interfaces on each node. This card used an NEC D720101GJ Chipset. (Image borrowed from <http://tinyurl.com/37m8js>.)*

VT6122 (Gigabit Ethernet) or the VIA VT6103L interfaces. The engineers at VIA helped develop this customized BIOS for the particular needs of this project.

Coupled with the VIA VT310-DP motherboard was a USB 2.0 PCI card to expand the number of USB slots for the cluster's interconnect network. Using throughput tests, it was found that the NEC USB PCI controller worked better than the others tested. This controller did a good job of balancing the USB bandwidth equally between the connected USB patch cables.

Finally, there were many preliminary operating systems which have been used. Back when this project started in August of 2005, OpenBSD 3.8 and 3.9 were used both on Sandlab (as the PXE-server and NFS server), as well as the diskless nodes. A tutorial for OpenBSD diskless compilation can be found at:

<http://www.openbsdsupport.org/diskless.pdf>

Soon after, discussion started with the OpenBSD developers towards developing a stable EHCI Ethernet interface via the USB patch cables, as the current version at that date caused kernel panics when attempting to send the maximum amount of data over the USB patch



**Figure 4.6:** *This image is a photograph of the USB 2.0 NetLink cables used as the interconnect on Thor’s Tack Hammer. The links used an ALi M5632 Chipset. (Image borrowed from <http://tinyurl.com/37zk7w>.)*

cable.

Other operating systems which were considered were PelicanHPC, and CentOS (using OSCAR), but both were discarded after issues with the USB Ethernet drivers and the diskless architecture of Thor’s Tack Hammer. Finally, a slightly remastered version of ParallelKnoppix (PK) 2.7.1 was used on the cluster.

The additional packages and tools needed for this particular research was:

1. NTP v.4.2.4p4 (to ensure all nodes have nearly similar clocks)
2. SchedUtils v.1.5.0-1 (CPU affinity)
3. BASH scripts (to automate the cluster’s initialization - see Appendix 1)
4. C programs (to generate/route/account for traffic over the interconnect)

The management interface of each network (which was used to PXE-Boot, WOL, and maintain NFS mounts for Linux on each node) was implemented over standard Category-5 Twisted-Pair Ethernet cables. Each cluster node connected to a Cisco Catalyst 3500 series XL switch. The head node also connected to this switch.



## 4.3 Node and Cluster Performance

There are many possible ways of measuring performance on computational systems and networks. This section is dedicated towards outlining various ways in which Thor's Tack Hammer was benchmarked and ranks against other parallel systems and sequential systems.

Before delving into how Thor's Tack Hammer (or any parallel system) compares to sequential systems, it is important to first consider Amdahl's Law and its application towards understanding parallelization. Amdahl's Law has been used to find the overall expected improvement (or speedup) that can be realized through multiple processors as opposed to a single processor. This law depends greatly on the ratio of which instructions can actually be partitioned onto other processors (with zero communication between them) and computed independently.

Amdahl's Law:  $S_{Total} = \frac{1}{r_s + (\frac{r_p}{n})}$   
 $r_s$ : The fraction of instructions that must be run sequentially.  
 $r_p$ : The fraction of instructions that can be run in parallel.  
 $n$ : The number of processors running the parallel instructions.  
 $S_{Total}$ : Total resulting speed-up.

For instance, if one particular task can have 18% of the instructions run in parallel, leaving 82% needing to run concurrently, by executing them across multiple processors, the maximum speedup possible for one processor makes the equation become:

$$1 = \frac{1}{0.82 + (\frac{0.18}{1})}$$

Whereas running the same code on a parallel system with 27 processors, such as Thor's Tack Hammer, we see a maximum possible speedup of:

$$1.209 = \frac{1}{0.82 + (\frac{0.18}{27})}$$

And finally, running that code on BlueGene/L with 212,922 processors, we see a maximum speedup of:

$$1.219 = \frac{1}{0.82 + (\frac{0.18}{212,922})}$$

These results clearly depict the necessity of having a good ratio of parallel and sequential instructions in order to truly capitalize on the number of processors in an HPC. With this low ratio of parallel instructions to the sequential instructions, the difference in maximum speedup between Thor's Tack Hammer and BlueGene/L is marginal.

Along with Amdahl's Law, and its impact on the performance of Thor's Tack Hammer, there are other benchmarks available to demonstrate overall performance. These next two sections outline two other benchmarks, throughput between nodes on the USB interconnect, and the Linpack benchmark.

### **4.3.1 USB Interconnect Throughput**

Network throughput between systems can be easily measured by using tools such as Netperf and Iperf. These tools attempt to measure maximum throughput by sending as much data as possible from a source to a destination, given a static amount of time.

In order to measure the throughput between nodes using the USB interconnect, Iperf version 2.0.3 was used. There were two different tests used, first measuring how much data was able to be sent over a single USB link, and the second measuring how much data was able to be sent from a single node to all of its 6 neighbors simultaneously.

The first test attempted to measure the maximum throughput on a single USB link. That maximum throughput was found to be around 76Mbits/sec. The second test disclosed evidence that each node uses two separate USB EHCI (USB 2.0) controllers. The two controllers, a NEC USB controller (via the PCI card, 4 USB ports used), and the VIA USB controller (via the onboard USB ports, 2 USB ports used). Because of the two separate controllers, sending data over all six interfaces totalled at 230Mbits/sec to 265Mbits/sec. This also indicates that the patch cable runs slower than the controller.

### **4.3.2 Linpack Benchmark**

The Linpack benchmark, widely used and accepted as a standard measure of compute power, operates by solving a random dense system of linear equations. The benchmark is the same

System	Number of Processors	Peak GFlops/sec	Top500 Rank	HPL
LLNL-BlueGene/L (US)	212,922	596,378	1 <sup>st</sup>	Y
SNL-Red Storm (US)	26,569	127,531	6 <sup>th</sup>	Y
Thor's Tack Hammer	27	2.299	—	Y
Thor's Tack Hammer	1	0.190	—	N
iBook 1.42Ghz PPC G4	1	0.051	—	N

**Table 4.1:** *Linpack test results of various machines. These results were obtained from the November 2007 Top500.org website. The left most column indicates whether the Linpack or the Linpack/HPL version was used.*

benchmark used by the Top500.org site to rank the top 500 fastest computers in the world. It uses the BLAS library (Basic Linear Algebra Subprograms) to solve the equations using Gaussian elimination with partial pivoting. The benchmark measures how many millions of floating point operations per second were observed during the computation.

In order to facilitate a proper parallel Linpack benchmark, the Linpack/HPL version of the test tool was downloaded and executed run over MPI (Message Passing Interface).

Various parameters are used to "tune" the benchmark. These parameters vary the ways in which the problem set is partitioned among the nodes within a cluster or larger HPC. The three parameters used to benchmark this cluster are  $N$ , which specifies the number of problems to be run;  $NB$ , which specifies the block size of the problem set; and finally the  $P$  and  $Q$  parameters, which partition the entire problem set between the nodes.  $P$  and  $Q$  are multiplied together, and must be less than or equal to the number of nodes in the cluster.

Documentation that came with the downloaded Linpack/HPC tool explained that by trying all possible combinations of  $P$  and  $Q$ , the user should find one value that is the best result.

Table 4.1 shows the Linpack test results for various well-known computer systems and HPC systems. Table 4.2 show the varied Linpack/HPL test results for Thor's Tack Hammer during tuning.

Number of Processors	$N$	$NB$	$P$	$Q$	Time(sec)	Peak GFlops/sec
1	500	100	1	1	438.44	0.190
3	5000	100	1	3	179.95	0.463
27	5000	100	1	27	53.95	1.533
27	5000	100	27	1	76.89	1.099
27	5000	100	9	3	41.10	2.058
27	500	100	3	9	35.78	2.299

**Table 4.2:** *Varied Linpack/HPL tests for Thor’s Tack Hammer after adjusting the tuning parameters.*

Watts	Context During Observation
12	All nodes off; one power supply on; Cisco switch off.
68	All nodes off; all power supplies off; Cisco switch on.
187	All nodes off; all power supplies on; Cisco switch on.
270	3/27 nodes on/idle; all power supplies on; Cisco switch on.
365	6/27 nodes on/idle; all power supplies on; Cisco switch on.
930	27/27 nodes on/idle; all power supplies on; Cisco switch on.
1060	27/27 nodes on/HPL; all power supplies on; Cisco switch on.

**Table 4.3:** *This table shows the results of power consumption on Thor’s Tack Hammer. These results were observed using a Kill-A-Watt P4400 device.*

### 4.3.3 Power Consumption

This last section demonstrates the power consumption of Thor’s Tack Hammer in terms of electricity used. These results were obtained by observing the amount of watts that were being used at various points during the startup and execution of the cluster. The tool used to analyze the watts used was a P3 Kill-A-Watt P4400. The following table shows the various levels of power used and the context in which that observation took place.

Table 4.3 makes a few things clear about the consumption of power for this particular cluster. First, it is interesting to see that because of the low consumption of power per node, this entire cluster uses less electricity than that of a typical hair drier or microwave oven.

It’s also worth analyzing further the power consumption per node both with idle and at high utilization, considering all factors included within the network. These factors include power consumption by the switch for the management network and residual power

consumption by the power supplies.

Therefore we get values of 34.44 watts per node for idle processing, and 39.26 watts per node for high processor utilization for this particular cluster network. These results should indicate that operating a cluster such as this can be very energy efficient. Also, because of these results, it's also clear that not all parallel systems require complex cooling mechanisms. This cluster maintains a steady temperature only by the fans which circulate air over the processors' heat sinks.

# Chapter 5

## Laboratory Results

This next chapter explains in depth the methodology for obtaining the results comparing the various algorithms on Thor's Tack Hammer as well as presenting and interpreting those results.

### 5.1 Methodology for Tests

The method for testing these algorithms within Thor's Tack Hammer was implemented using C code and BASH scripts. The code and scripts that were written were fully commented and are outlined in the appendices.

It was first considered to use routing tables to enable inter-node routing, but upon further consideration, would be difficult if not unreasonable to implement. This is because dynamic implementations would require constant and costly calculations and changes to the routing tables for every packet received.

Along with dynamically changing the routing tables, current transport-layer network daemons which were compatible with dynamic routing algorithms were considered. These included GNU Zebra ([www.zebra.org](http://www.zebra.org)) and Quagga ([www.quagga.net](http://www.quagga.net)). These applications have been under active development for dynamic routing algorithms such as BGP, OSPF, and RIP. This was an option that was considered, but was not implemented.

Therefore, an application-layer socket program was developed to appropriately receive data, calculate the next hop, and forward the data accordingly. A few assumptions were

made to verify that this application would in fact be suitable for testing and evaluating the algorithms in question:

**Assumption 1:** *The use of TCP was necessary to guarantee correct delivery of data between nodes.* This was an assumption made after preliminary evidence that without the flow-control of TCP, UDP could not guarantee correct data delivery. It was also assumed that a layer-4 protocol such as UDP would not be a good implementation choice on a larger supercomputer for this same reason.

**Assumption 2:** *Time stamps were necessary to compare both the algorithms and the C routing programs.* The time stamps provided crucial information about any particular data's experience traversing through the network. This included actual time spent during the calculations for a data segment's next hop, as well as time spent in transit between nodes. NTP was used to allow a fine resolution of time, and to ensure all nodes had a value nearly exact to that of the other nodes. Because of drifting the estimated error between nodes ranged from 10s to 100s of micro seconds. This is supported in [11].

**Assumption 3:** *By comparing each routing algorithm's time stamps of service delay, it is possible to compare their complexities.* By using the best-performing algorithm as the baseline, the other algorithms can be compared accordingly. This of course represents a software-based calculation and should not be considered an accurate assessment of a hardware-based implementation.

**Assumption 4:** *A fork subroutine was necessary to lower the impact of blocking calls onto the processing of other data.* This enabled nodes to continue to process data in parallel over all six USB interfaces simultaneously, and limited the blocking calls to only affect the data associated with that child process.

**Assumption 5:** *Because access to a node's actual interface output queues was not available at the application level, a shared memory implementation using semaphores would be sufficient.* Output queues were implemented using semaphores and shared memory between the child processes. This was necessary for the fully adaptive algorithms.

**Assumption 6:** *This C-program implementation aligns with the queueing models of the previous work.* Even though the C-programs implement a queue using forks, and a common area of shared memory protected with semaphores, each interface has a single processor. Also, a constant value was used to decrease the queues when necessary. Therefore this C-program implementation is equivalent to the model proposed in the previous work, which modeled the queues as M/D/1 queues.

These assumptions help justify the way in which validation of the algorithms was performed. They also help explain the reasoning behind the programs' implementations.

### **5.1.1 A Packet's Progression through the Interconnection Network**

It is important to consider how data is not only created, but transmitted, routed, and cataloged throughout the interconnection network using the C programs. This section explains just that.

*Step 1: Data Injection.* Data is created in chunks of variable size, all "dummy" random data. This data includes a header of fixed size used to contain source, destination, time stamps, and other vital information for proper routing. This all occurs through a program run on the head-node (which is not part of the USB interconnect network) using `inject.c`, the Injector program. This Injector program reads a file which lists all source/destination traffic demands and paired with the specific routing algorithm, creates the data and sends the data via TCP through the management network directly to the source node. The assumption that all data is injected into the network at nearly the same time can be validated by examining time stamps.

*Step 2: Preliminary Routing.* The source node then receives the data from the management network and processes the header. This file is the Server file, or `server.c`. It calculates the next appropriate hop (given the routing algorithm and the destination), and routes accordingly through the USB interconnect network. A time stamp is added at the source node to setup an initial time.



*Step 3: Intermediate Routing.* The data continues to be received through the USB network and intermediate hops process the data's header and send it onto the next hop using the `server.c` program. The time that the data spends while being processed (time between receipt of data and transmission to next hop) at each hop is added to the previous service value(s), and the number of hops the packet has traversed is incremented. If this hop count exceeds a predetermined value, it is removed from the interconnect network and processed like it had reached the destination successfully, but flagged as a dropped packet.

*Step 4: Destination Routing.* Once the data has been received at the destination node, it forwards the header information to another program listening on the loopback interface, `loserver.c`. This program parses the header once more, and writes valid information to a file which will be parsed by the BASH scripts. A final time stamp is generated at the destination node. CPU affinity allows for this process to run on a completely separate processor than the `server.c` program.

The `server.c` program is set to listen on any interface given a specified port, and can only transmit data onto an outbound USB interconnect interface. This enables it to receive data from the management port (during Data Injection), and guarantees it will only transmit the data from that node onto the USB interconnect network. This program runs on all nodes within the interconnect network.

The `loserver.c` program is set to listen only to the loopback interface on a different port, and receives the header and writes data to files. Some of the data obtained from the header includes average queue times, number of hops the data traversed, and actual time when data was created and removed from the interconnect network. This program runs on all nodes within the interconnect network.

The `injector.c` program is used to inject data into Thor's Tack Hammer. It is run on only the head-node (Sandlab). It accepts arguments of an input file and the algorithm to route on the USB interconnect network. The input file designates in a matrix form the demand (or number of data segments) to send given every source/destination pair.

## 5.2 Laboratory Results

This section presents both validation of previous routing algorithms as well as the data gathered from the execution of Enhanced CQR with Periphery Avoidance on Thor's Tack Hammer. All of the assumptions and implementations described in the previous section apply to this data.

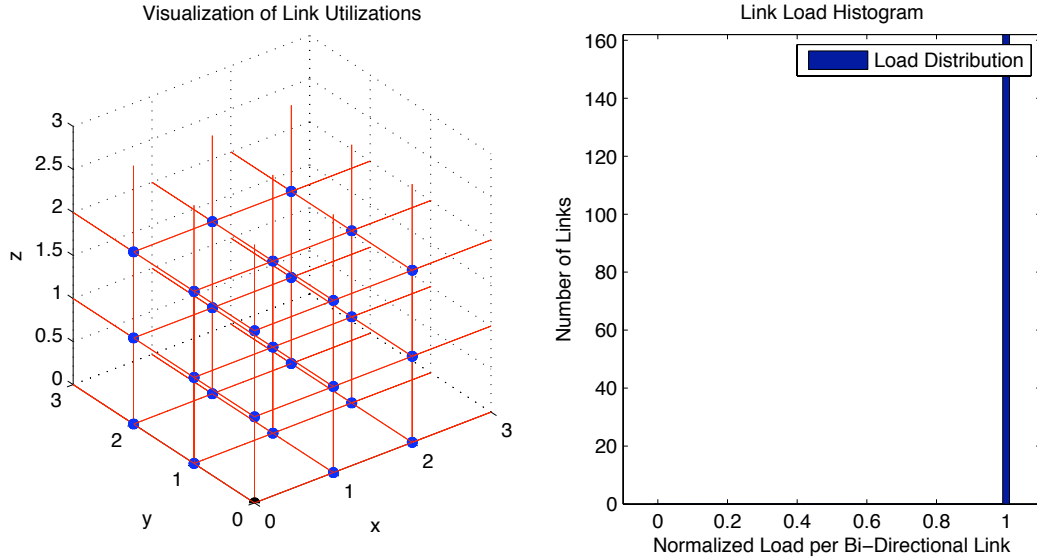
The result from the laboratory analysis are presented below by graphically showing their link-load and a histogram of all links normalized. In order to normalize these links, the maximum amount of data transmitted across a single bi-directional link was found and then all others were normalized to that value. Therefore, the highest-utilized link is given a value of 1.0 and all subsequent lower-utilized links have lower values. Those links with zero utilization then have a value of 0 within the histogram.

All algorithms were tested using the Table 2.1 before, and by using the BASH scripts in Appendix A to parse the results. These results were then placed through Matlab scripts listed in Appendix C to graphically view loads and link-utilization distributions.

Static and oblivious algorithms were not tested against the presence of hotspots, as they are unable to make decisions because of them [5]. The adaptive algorithms, however, were compared both with and without the presence of hotspots. See Table 2.1 for an explanation of traffic patterns, and sections (missing ref) for in-depth explanations of each algorithm.

There are some basic indicators of whether or not an algorithm outperforms another. First, a tightly distributed and highly utilized network indicates good load balancing. This may seem counter-intuitive but because the links are all normalized, it is desired. By having all links (or most links) at 100% utilization, that indicates that each link experienced the same traffic load as every other link. This means that for the histograms presented in Appendix E (which demonstrate the distribution of link utilizations), a large grouping near 1.0 or at 1.0 is highly desired. The tight distribution indicates a high number of links at or near the same value.

Another indicator which proves helpful in comparing the algorithms is the standard

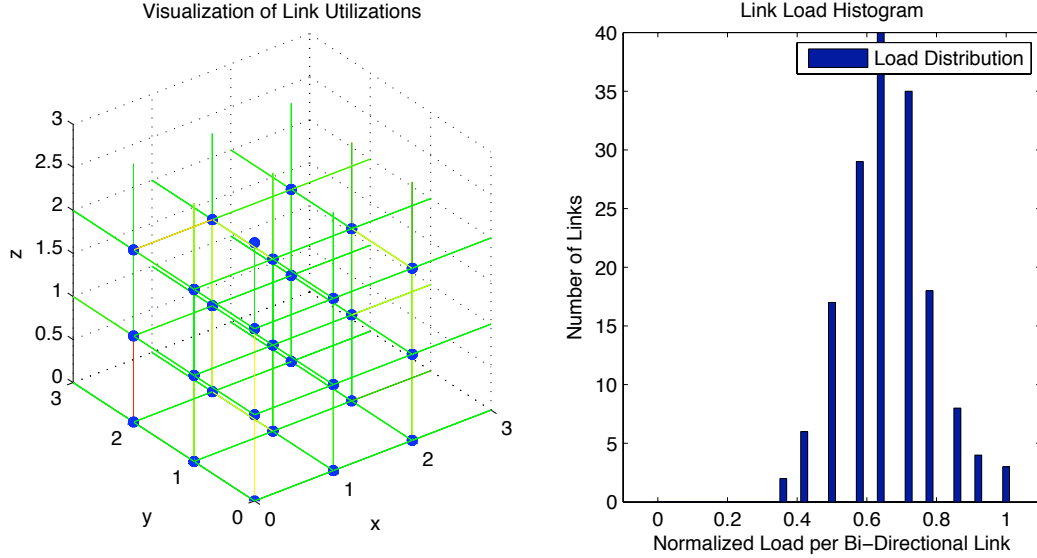


**Figure 5.1:** *This figure illustrates what are considered excellent laboratory results. Because all links experienced 1.0 utilization after being normalized, they all experienced the exact same traffic load - indicating very efficient load-balancing.*

deviations of the link distributions. Those values are presented in the next few tables. The reason for including these values is to demonstrate how "tight" the link distributions were. The optimum value we search for with standard deviations would be 0.0, indicating all links were at a singular spot. Because of the normalization techniques used, if the standard deviation was 0.0, the utilizations would all be 1.0. This situation is shown in Figure 5.2.

Lastly, by comparing the service times for each algorithm, certain assumptions can be made. In order to demonstrate these scenarios and give examples which indicate good, average, and poor results, Figures 5.2 5.2 and ?? are included below. The first, Figure 5.2 demonstrates a highly utilized network, a zero standard deviation, and illustrates the scenario for absolute optimum load-balancing. Figure 5.2 illustrates an example which is considered good, because of the assumed normal distribution with a fairly low standard deviation. Finally, Figure ?? illustrates bad results, in that many links were at 0.0 utilization (indicating many links were unused during that execution).

On another note, bi-modal distributions could indicate various things, but most certainly,



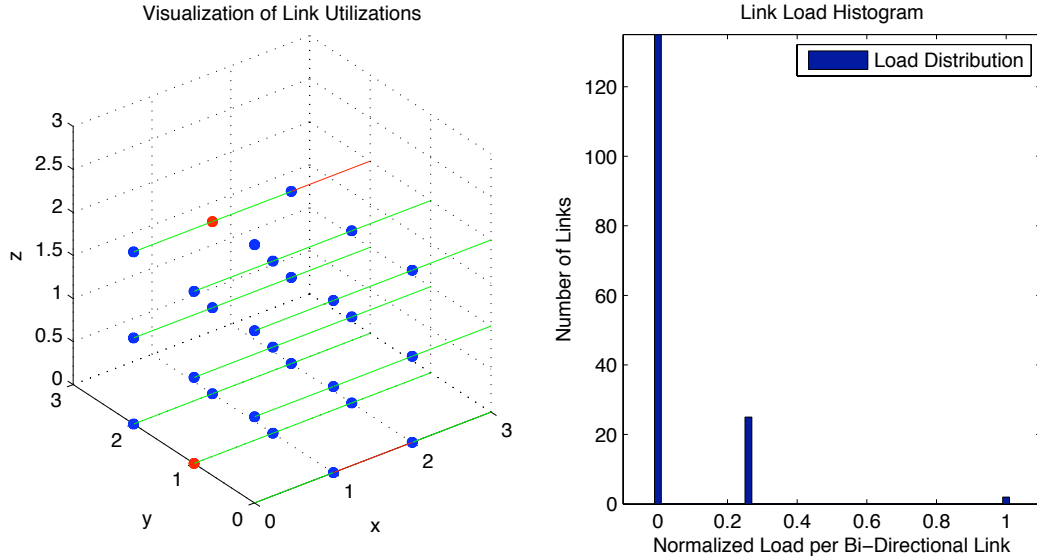
**Figure 5.2:** *This figure illustrates what are considered OK laboratory results. Because all links experienced traffic loads (none were at 0.0 utilization) and the standard deviation of the link distribution was fairly low, these results were better than poor, but worse than good or excellent.*

indicate faulty decisions given a traffic demand. In certain cases, this bi-modal characteristic is unavoidable (as it is with minimal adaptive given a tornado traffic pattern, for instance).

As it was stated in the assumptions above, this service time can be an indication as to the computational complexity and overhead associated with the implementation of a particular algorithm - both through software and more or less through hardware.

When considering the graph on the left of each result figure (both laboratory and simulation), green lines represent lowly utilized links (which range from 0 :  $(1 - std\_dev)$ , in a normal distribution that is 0-66% utilization), yellow lines represent medium utilized links (which range from  $(1 - std\_dev)$  :  $(1 - \frac{std\_dev}{2})$ , in a normal distribution, ranges from 66-83%), and red lines representing highly utilized links (which range from  $(1 - \frac{std\_dev}{2})$  : 1, in a normal distribution, ranges from 84-100%).

The exact same traffic demands were implemented for these tests as were implemented for the 3-ary 3-cube simulation presented in the previous chapter. Each individual traffic



**Figure 5.3:** *This figure illustrates what are considered poor laboratory results. There were a large number of links which experienced zero traffic (indicating poor load-balancing techniques).*

demand consists of 800KB of data and any TCP overhead associated with its transmission from source to destination. Because of 1500 Byte MTUs on the interconnect, fragmentation and reassembly procedures were necessary. The application layer header was included within the 800KB data. All results presented were the average values after executing the algorithm and traffic demand twice.

### 5.2.1 Dimension Ordered Routing (DOR) Results

Appendix D.2.1 shows the graphs associated with Dimension Ordered Routing when varying the traffic demands. The benefits of using DOR are visible in what are deemed easy or benign traffic patterns such as NN, UR, and FLOOD. The hard traffic patterns, such as TOR, BC, and TP, did not yield good load-balancing results.

Table 5.1 shows more results obtained during these tests.

Traffic Pattern	NN	UR	BC	TP	TOR	FL
Avg Service Time per Node (ms):	16.6	13.5	14.9	15.0	16.4	13.9
Avg Utilization per Link (%):	64.6	30.6	22.9	26.1	13.5	64.1
Link Utilization Std Dev:	13.8	19.2	33.0	26.9	29.7	7.5

**Table 5.1:** *This table shows the results obtained from the execution of the various traffic patterns while using Dimension Ordered Routing on Thor’s Tack Hammer*

Traffic Pattern	NN	UR	BC	TP	TOR	FL
Avg Service Time per Node (ms):	18.2	13.4	14.6	13.9	15.6	13.8
Avg Utilization per Link (%):	65.5	37.6	8.8	25.6	10.6	73.0
Link Utilization Std Dev:	14.6	23.7	18.2	25.6	23.7	8.6

**Table 5.2:** *This table shows the results obtained from the execution of the various traffic patterns while using Direction Ordered Routing on Thor’s Tack Hammer*

### 5.2.2 Direction Ordered Routing (DIR) Results

Appendix D.2.2 shows the graphs associated with Direction Ordered Routing when varying the traffic demands. The results of DIR are nearly identical to DOR, and use the same benign and hard traffic patterns as the previous section discusses.

Table 5.2 shows more results obtained during these tests. As the results show, the service times are lower than that of the adaptive and oblivious algorithms, but the utilizations are lower than the others.

### 5.2.3 Minimal Oblivious Routing Results

See Appendix D.3.1 for the graphs associated with Minimal Oblivious Routing when varying the traffic demands.

Table 5.3 shows more results obtained during these tests.

Traffic Pattern	NN	UR	BC	TP	TOR	FL
Avg Service Time per Node (ms):	19.0	14.0	13.7	14.5	11.6	15.4
Avg Utilization per Link (%):	70.6	43.5	18.3	26.3	11.6	63.3
Link Utilization Std Dev:	18.3	23.9	26.8	27.4	26.0	12.9

**Table 5.3:** *This table shows the results obtained from the execution of the various traffic patterns while using Minimal Oblivious Routing on Thor’s Tack Hammer.*

Traffic Pattern	NN	UR	BC	TP	TOR	FL	FL-HS
Avg Service Time per Node (ms):	21.8	20.6	30.0	28.6	31.3	19.6	23.7
Avg Utilization per Link (%):	60.4	38.9	11.9	23.6	12.4	53.7	33.5
Link Utilization Std Dev:	15.1	20.5	23.5	24.5	27.3	12.6	13.3

**Table 5.4:** *This table shows the results obtained from the execution of the various traffic patterns while using Minimal Adaptive Routing on Thor’s Tack Hammer.*

### 5.2.4 Minimal Adaptive Routing Results

See Appendix D.4.1 for the graphs associated with Minimal Adaptive Routing when varying the traffic demands. Because hotspots can affect the outcome of the results, the addition of hotspots were included with the FLOOD traffic pattern

Table 5.4 shows more results obtained during these tests. The results shown here indicate that the implementation of an adaptive algorithm increases the service time, but they also show that the utilization increased because of the ability to make better routing decisions and distribute the load better. This data aligns with that of the previous work in [13].

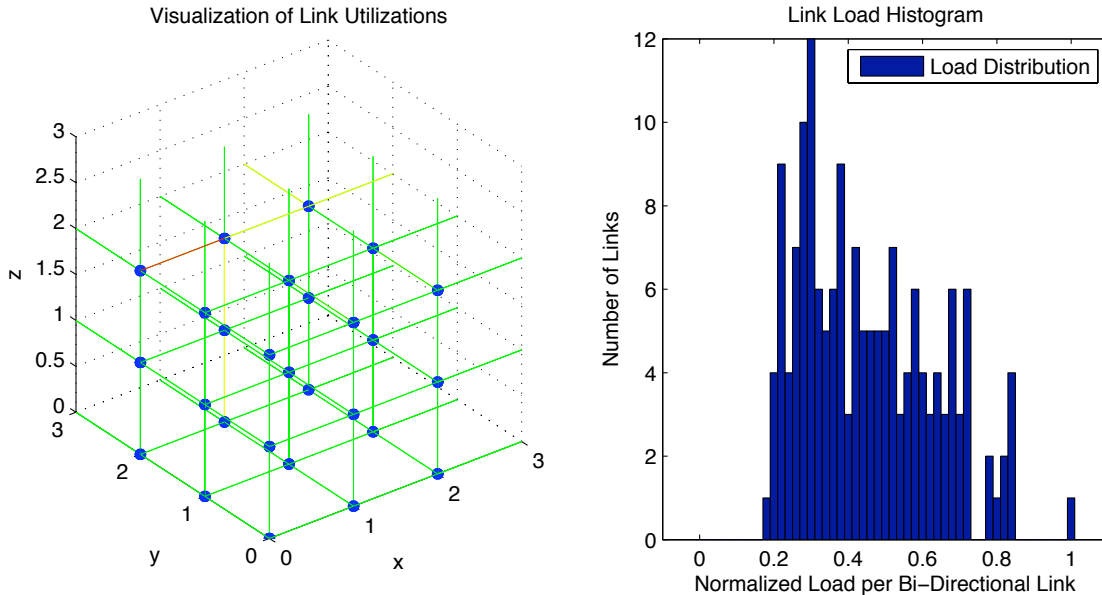
### 5.2.5 CQR Routing Results

See Appendix D.4.2 for the graphs associated with CQR Routing when varying the traffic demands. Because hotspots can affect the outcome of the results, the addition of hotspots were included with the FLOOD traffic pattern

Table 5.5 shows more results obtained during these tests. This algorithm showed a slight improvement over the utilizations from Minimal Adaptive. CQR also had similar service times, as should be expected. This data also aligned with the previous work of [13].

Traffic Pattern	NN	UR	BC	TP	TOR	FL	FL-HS
Avg Service Time per Node (ms):	24.3	20.5	26.5	22.4	32.4	23.0	22.5
Avg Utilization per Link (%):	46.8	38.8	12.1	29.9	14.3	45.4	33.1
Link Utilization Std Dev:	25.7	17.9	17.1	22.5	26.0	17.9	16.2

**Table 5.5:** This table shows the results obtained from the execution of the various traffic patterns while using CQR Routing on Thor’s Tack Hammer.



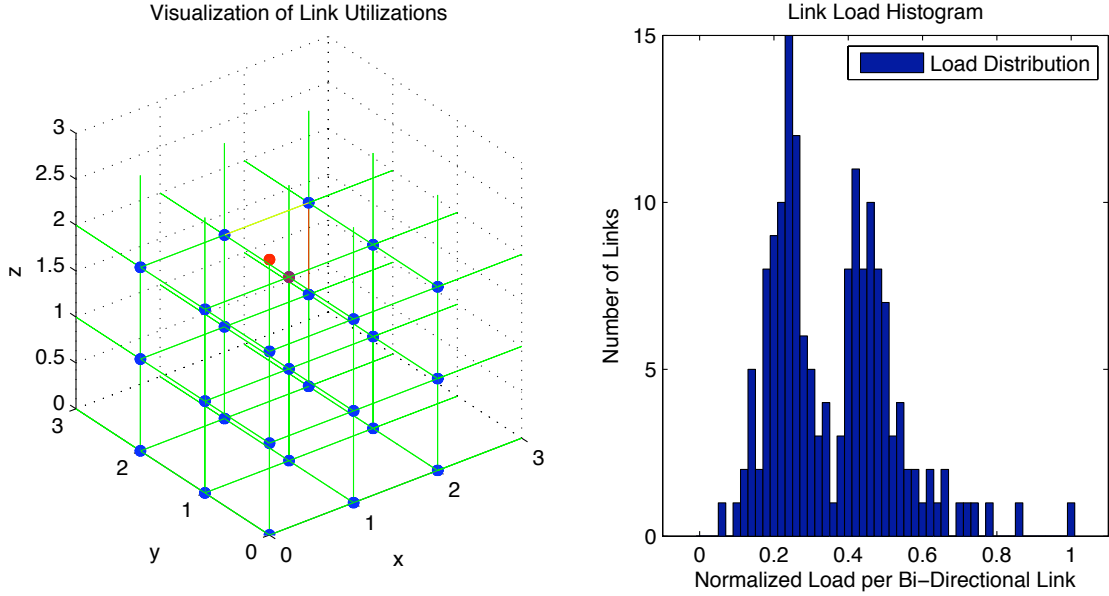
**Figure 5.4:** Results of CQR Routing using a flood traffic pattern (FL). Two individual traffic demands were assigned to each node from each node.

Also, it is worth noting that both algorithms, Minimal Adaptive and CQR, showed a heavy decrease in the average utilization when encountering hotspots. Both responded similarly to the hotspot traffic, dropping the average nearly 20%.

## 5.2.6 Enhanced CQR Routing Results

See Appendix D.4.3 for the graphs associated with Enhanced CQR Routing when varying the traffic demands. Because hotspots can affect the outcome of the results, the addition of hotspots were included with the FLOOD traffic pattern



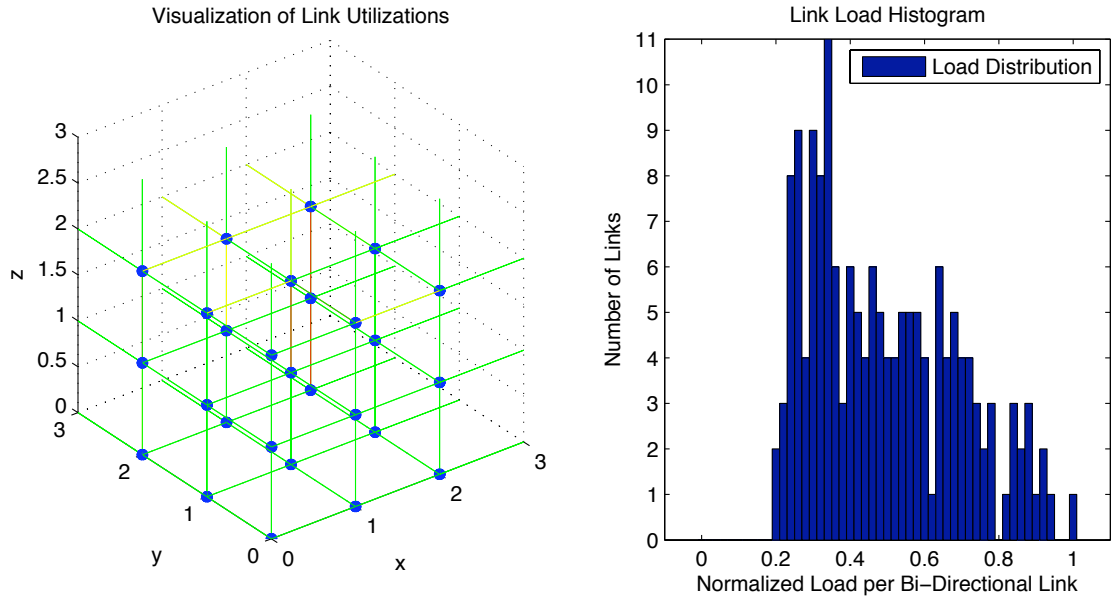


**Figure 5.5:** Results of CQR Routing using a flood traffic pattern (FL) with Hot-spots. Two individual traffic demands were assigned to each node from each node. Five percent of the nodes were made hot-spots, and they are marked as red in the figure.

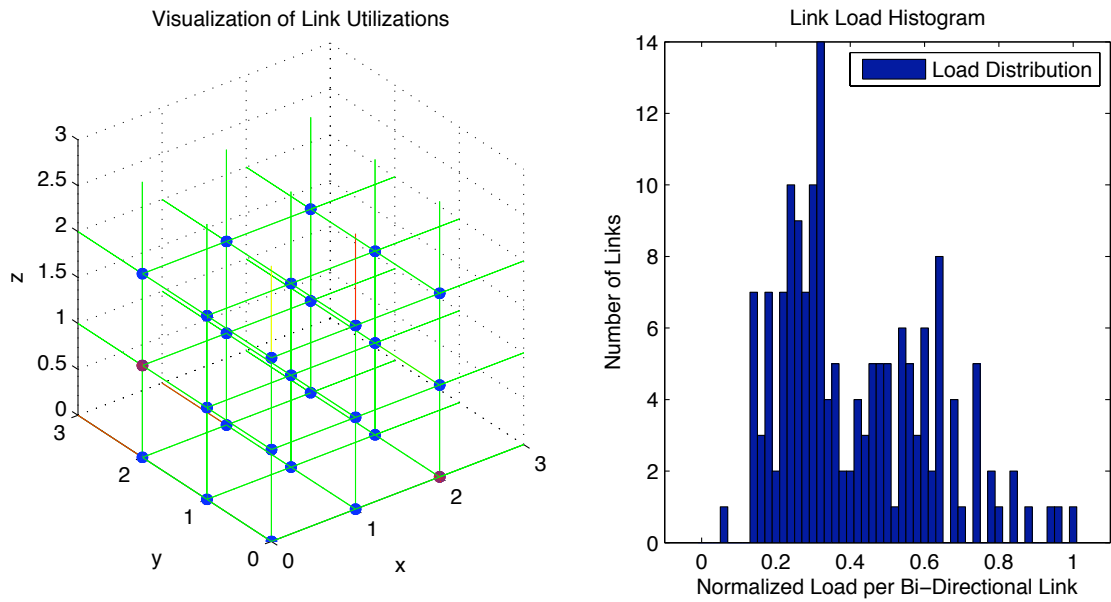
Traffic Pattern	NN	UR	BC	TP	TOR	FL	FL-HS
Avg Service Time per Node (ms):	23.7	22.6	27.7	21.1	35.1	21.7	16.5
Avg Utilization per Link (%):	45.1	43.4	14.4	32.6	16.4	49.2	41.4
Link Utilization Std Dev:	23.2	19.0	21.3	22.6	27.4	19.8	20.0

**Table 5.6:** This table shows the results obtained from the execution of the various traffic patterns while using Enhanced CQR Routing on Thor’s Tack Hammer.

Table ?? shows more results obtained during these tests. These results indicate that this enhanced version of CQR responds very similarly to that of traditional CQR in each tested traffic pattern. As for the FLOOD traffic pattern with hotspots, Enhanced CQR with Periphery Avoidance shows great potential, as it exceeds traditional CQR’s results with a 10% higher average utilization. These results do show that in the presence of hotspots, Enhanced CQR with its modified decision function, gives better results when encountering hotspots within a network.



**Figure 5.6:** Results of Enhanced CQR Routing using a flood traffic pattern (FL). Two individual traffic demands were assigned to each node from each node.



**Figure 5.7:** Results of Enhanced CQR Routing using a flood traffic pattern (FL) with Hot-spots. Two individual traffic demands were assigned to each node from each node. Five percent of the nodes were made hot-spots, and they are marked as red in the figure.

## 5.2.7 Comparison of All Results

In order to easily compare all of the results presented above, an elementary metric was developed. Two of the values obtained from results have an inverse relationship based on the following:

1. Good simulation results have high average utilizations. Because these values are all normalized, high utilizations indicate good load-balancing occurred. On the other hand, low utilizations indicate heavy-tails on the link utilization distribution, meaning a few links were very highly utilized compared to the links with average utilizations.
2. Good simulation results have low service times. As algorithms become more complex, their service times will increase. This is because it becomes more complex to calculate a next hop along the path. Therefore as the service times increase, the algorithm becomes more costly to implement.

The metric that was developed to compare the results uses this inverted relationship between link utilizations and service times to give a very elementary indication of overall performance. This metric does not weight the utilizations and service times, but that would not be difficult to include in future work.

$$m = \frac{Utilization}{Service\_Time} \quad (5.1)$$

Where

*Utilization* is the normalized utilization between 0.0 and 1.0,

*Service\_Time* is the service time measured in milliseconds.

As Table 5.7 shows, DOR and DIR algorithms performed very similarly, which was expected. For all algorithms, the hard traffic patterns gave much lower values, while the benign pat-

Routing Algorithm	NN	UR	BC	TP	TOR	FL	FL-HS
DOR	3.89	2.27	1.54	1.74	0.82	4.61	–
DIR	3.59	2.80	0.60	1.84	0.68	5.29	–
MO	3.72	3.11	1.34	1.81	0.10	4.11	–
MA	2.77	1.88	0.84	0.82	0.40	2.74	1.41
CQR	1.93	1.89	0.46	1.33	0.44	1.97	1.47
ECQR	1.90	1.92	0.52	1.54	0.47	2.27	2.50

**Table 5.7:** *By using the metric explained within Equation 5.1, it was possible to do an elementary comparison between all laboratory results. These results should only be used to demonstrate some level of increased performance when the routing algorithms gave higher values.*

terns gave higher values. CQR and ECQR also performed very similarly, except in the presence of hotspots, which ECQR showed increased performance over CQR.

# Chapter 6

## Simulation Results

This next chapter explains in depth the methodology for obtaining the results comparing the various algorithms by using Matlab scripts.

### 6.1 Methodology for Tests

The method for creating and testing the various routing algorithms on different topologies was made possible by creating a simulation environment within Matlab. The beginnings of this environment was initially developed by Dr. Don Gruenbacher, but was heavily modified to meet the needs of this research work.

By utilizing the Matlab functionality of matrices, it was possible to easily implement a simulation which would test an algorithm's ability to balance load, and minimize delay between source and destination given a specific topology.

The process by which this simulation environment operates is outlined in a numerical form below:

1. The user inputs simulation parameters such as the specific traffic demand to simulate, the routing algorithm to implement,  $k$  and  $n$  values for the topology, and probabilities for hotspots and/or nodal failures.
2. Given those values, the script first generates a traffic-demand matrix, which is  $k^n \times k^n$  in size, using rows as sources and columns as destinations. Hotspots are included

within the calculation of this matrix.

3. Depending on the routing algorithm that was selected, the traffic matrix is then passed to its corresponding function, which returns a load matrix and delay matrix. The load matrix is also of size  $k^n \times k^n$ , but only has values within  $(i, j)$  where  $i$  and  $j$  are neighbors and sent or received data during the simulation. The delay matrix is of size  $1 \times k^n$ , and includes the value for each node regarding how many packets it had to route non-minimally, enqueue beyond another time step, or drop (depending on the routing algorithm) – each of these possibilities represent some form of delay.
4. For the static algorithms, paths are calculated and simply added to the global traffic matrix. That matrix is then normalized by finding the highest utilized link and dividing the rest by that value. The same is done for oblivious algorithms. For adaptive algorithms, a different approach was necessary. Each simulation for the adaptive algorithms included individual time-steps, progressing every data segment one single hop, and then doing so until each segment successfully reaches its destination. This particular methodology gives way for the data to make adaptive routing decisions as it progresses through the network. Simply calculating a full path and adding it to the resulting traffic matrix would not provide an accurate model for this type of algorithm. The traffic matrix for adaptive algorithms are also normalized as they were for static algorithms.
5. Once the traffic matrix for a simulation has been calculated, it is then processed through another script which displays the link utilizations in a graphical form, and specific to the topology. Coupled with that display is a histogram depicting the distribution of the link utilizations.

Given the simulation environment described above and the code given in Appendix C, three topologies were fully examined. First, a 3-ary 3-cube was examined, which would coincide with the results from Thor’s Tack Hammer (also a 3-ary 3-cube topology). Sec-

Routing Algorithm	NN	UR	BC	TP	TOR	FL	FL-HS
DOR	100%	34.1	33.0	26.7	16.7	100	–
DIR	100%	33.2	8.3	33.3	16.7	100	–
MO	100%	31.7	11.1	26.7	16.7	60.0	–
MA	100%	49.9	16.7	44.4	16.7	75.0	45.8
CQR	39.4%	46.8	19.4	41.8	16.7	65.9	35.2
ECQR	39.4%	52.6	18.4	46.6	16.7	66.2	43.6

**Table 6.1:** *This table demonstrates the average utilizations for each routing algorithm given the 3-ary 3-cube topology of the simulation environment.*

Routing Algorithm	NN	UR	BC	TP	TOR	FL	FL-HS
DOR	0.0	21.2	47.0	27.8	37.7	0.0	–
DIR	0.0	20.2	16.7	32.0	37.8	0.0	–
MO	0.0	22.9	20.4	27.7	37.4	14.8	–
MA	0.0	17.6	25.4	26.5	37.3	9.11	13.8
CQR	30.4	17.7	23.8	16.5	37.4	12.5	12.0
ECQR	30.3	17.2	20.1	20.4	37.4	13.8	13.5

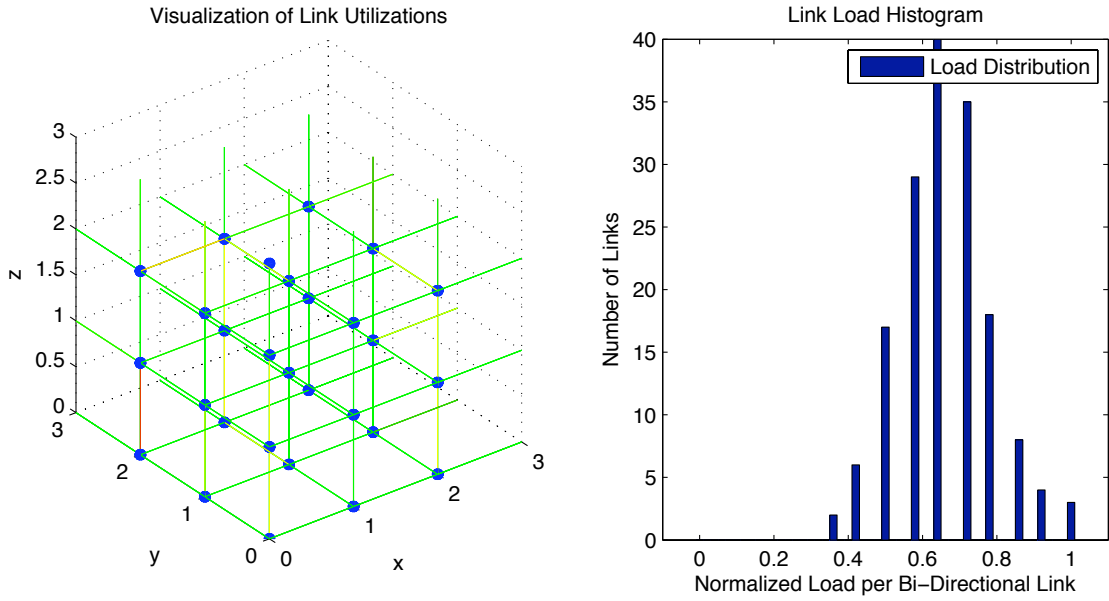
**Table 6.2:** *This table demonstrates the standard deviations for each routing algorithm given the 3-ary 3-cube topology of the simulation environment.*

only, a larger network of 4-ary 3-cube was simulated, while lastly a 5-ary 3-cube was also simulated. Those results are provided in the next few sections.

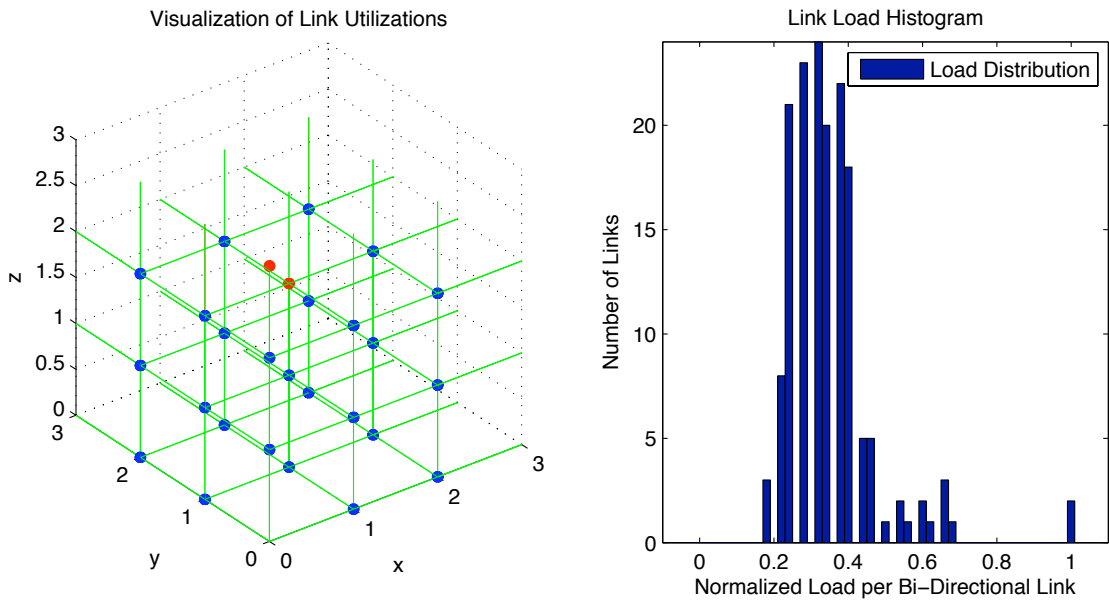
## 6.2 Results from the 3-ary 3-cube Simulations

Tables 6.1 and 6.2 show all results obtained from the simulations in a 3-ary 3-cube hypercube topology. All tests used the exact traffic demands that were used within the 3-ary 3-cube laboratory results of the previous chapter, also shown in Appendix D.

The first table shows the average utilizations, and the second table shows the standard deviation for that specific simulation. The graphs of these results are available in Appendix E.1.

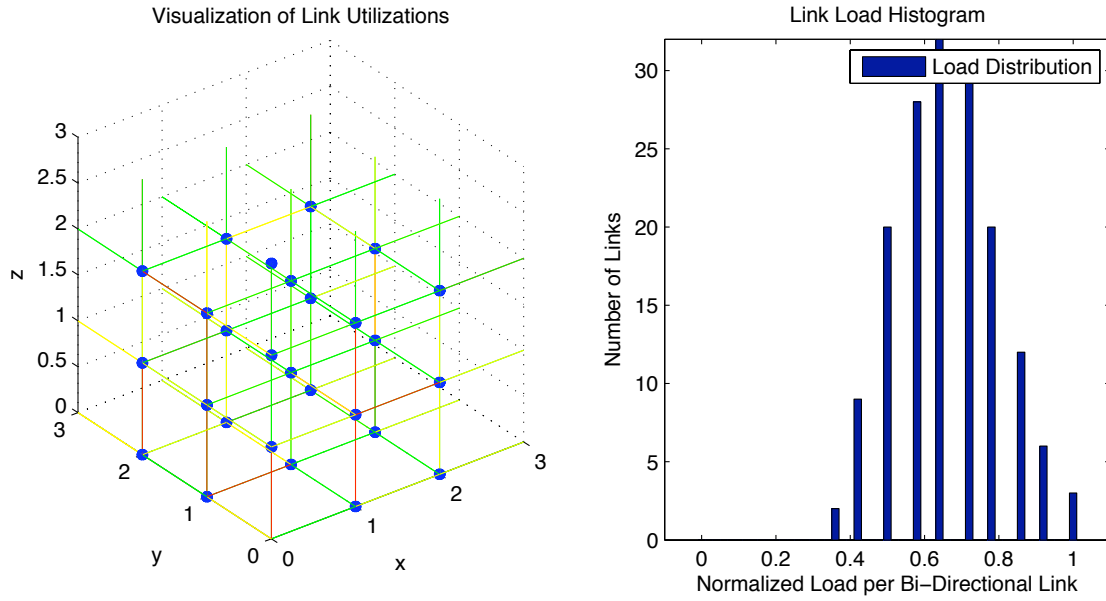


**Figure 6.1:** Results of CQR Routing using the flood traffic demand (FL). Two individual traffic demands were assigned to each node from each node.

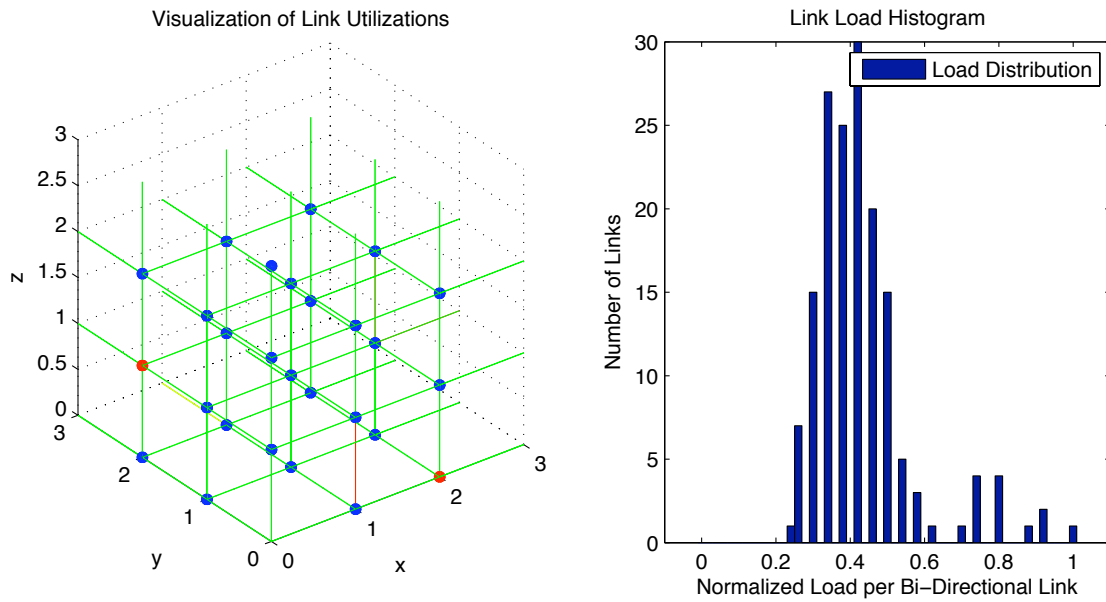


**Figure 6.2:** Results of CQR Routing using the flood traffic demand (FL) with Hot-spots. Two individual traffic demands were assigned to each node from each node. Five percent of the nodes were made hot-spots, and they are marked as red in the figure.





**Figure 6.3:** Results of Enhanced CQR Routing using the flood traffic demand (FL). Two individual traffic demands were assigned to each node from each node.



**Figure 6.4:** Results of Enhanced CQR Routing using the flood traffic demand (FL) with Hot-spots. Two individual traffic demands were assigned to each node from each node. Five percent of the nodes were made hot-spots, and they are marked as red in the figure.

Routing Algorithm	NN	UR	BC	TP	TOR	FL	FL-HS
DOR	100%	30.1	50.0	25.0	16.7	80.0	–
DIR	100%	32.8	12.5	25.0	16.7	80.0	–
MO	100%	31.0	16.7	22.2	16.7	59.0	–
MA	100%	58.4	25.0	44.4	16.7	85.3	50.3
CQR	34.2%	56.6	25.0	47.6	16.7	75.6	45.9
ECQR	34.4%	51.2	35.0	46.0	16.7	69.5	41.1

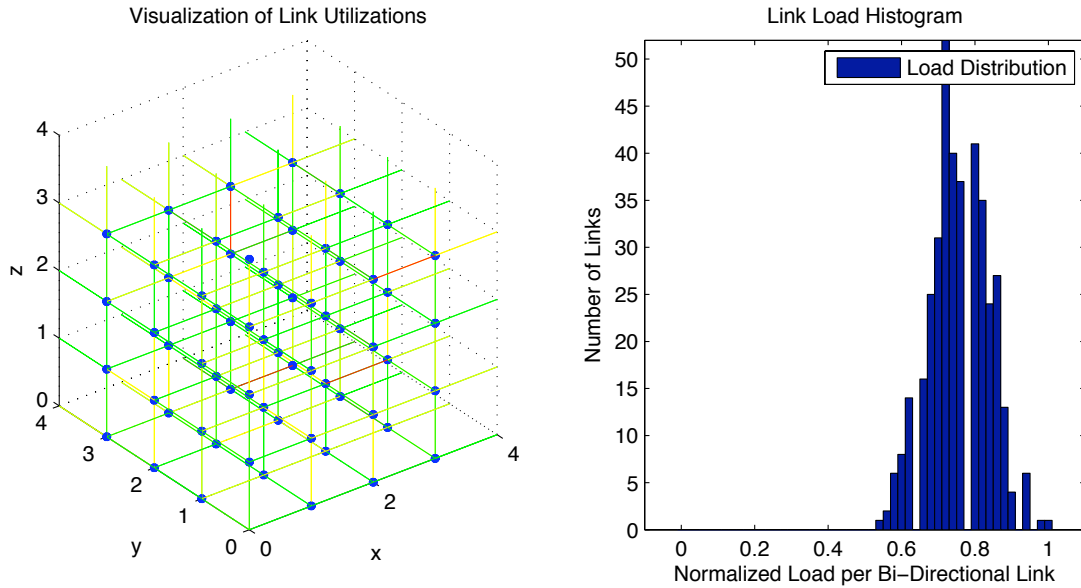
**Table 6.3:** *This table demonstrates the average utilizations for each routing algorithm given the 4-ary 3-cube topology of the simulation environment.*

Routing Algorithm	NN	UR	BC	TP	TOR	FL	FL-HS
DOR	0.0	17.5	50.0	25.4	37.3	6.4	–
DIR	0.0	16.9	23.4	21.2	37.3	7.8	–
MO	0.0	18.5	26.7	18.5	37.3	16.9	–
MA	0.0	13.4	29.5	18.3	37.3	5.4	10.4
CQR	21.5	12.8	23.5	16.5	37.3	8.01	10.2
ECQR	31.5	14.8	23.5	17.8	37.3	15.1	12.2

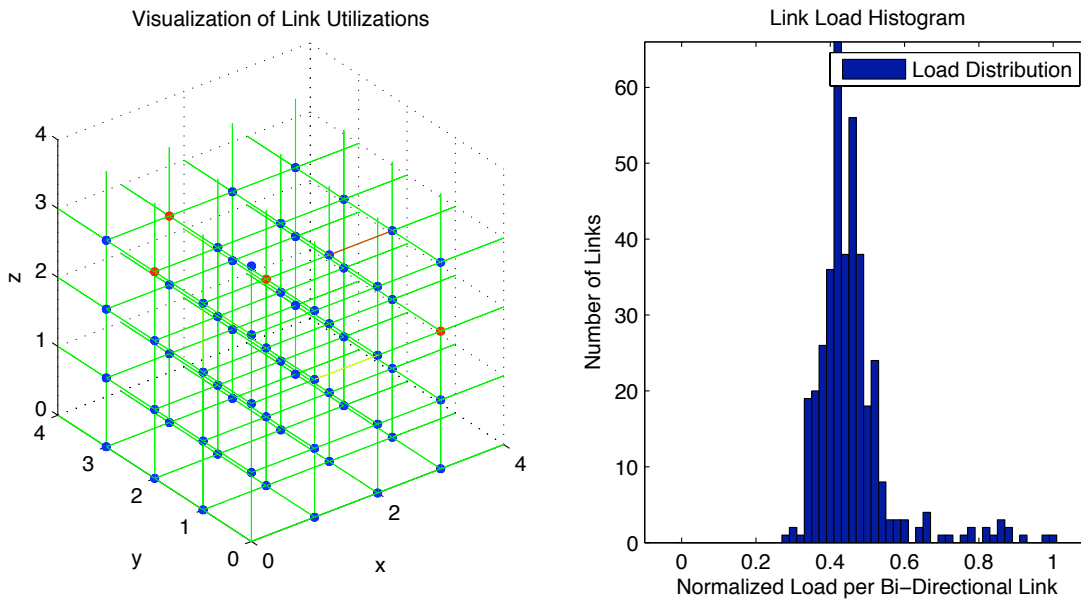
**Table 6.4:** *This table demonstrates the standard deviations for each routing algorithm given the 4-ary 3-cube topology of the simulation environment.*

### 6.3 Results from the 4-ary 3-cube Simulations

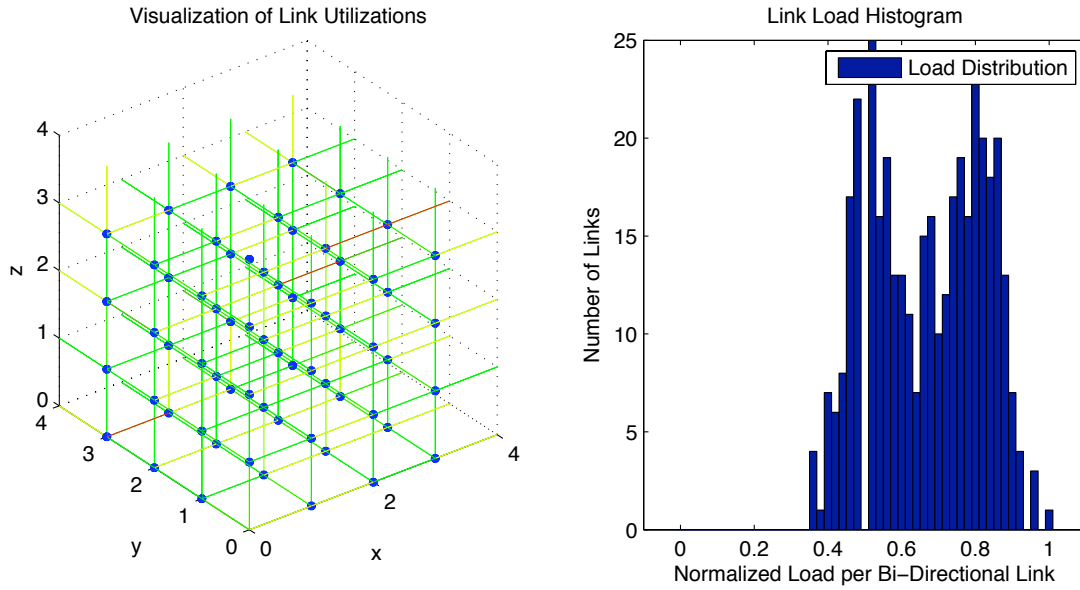
Tables 6.3 and 6.4 show all results obtained from the simulations in a 4-ary 3-cube hypercube topology. The first table shows the average utilizations, and the second table shows the standard deviation for that specific simulation. The graphs of these results are available in Appendix E.2.



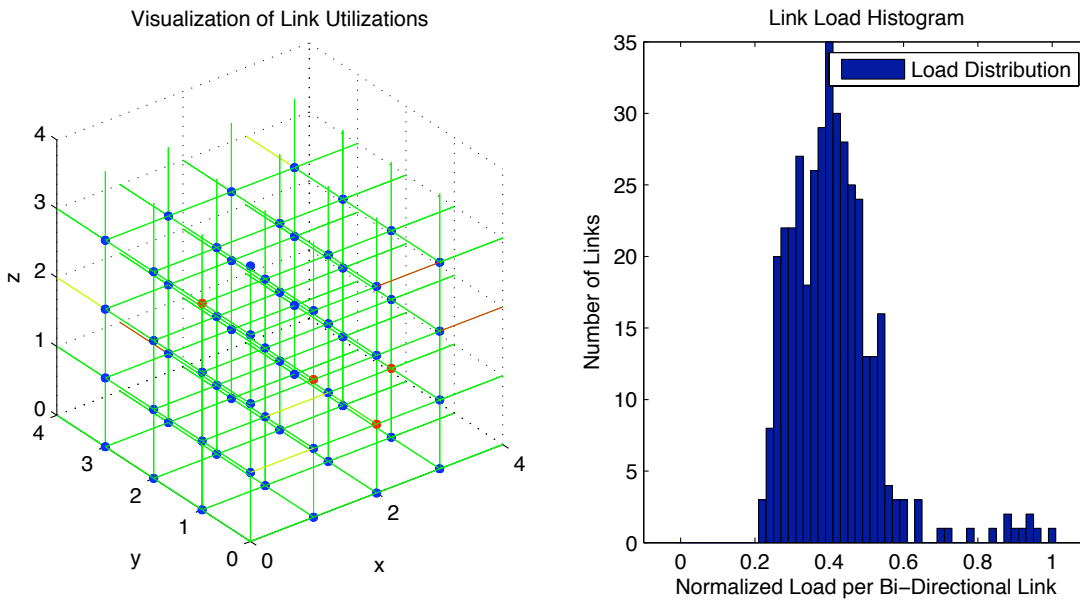
**Figure 6.5:** Results of CQR Routing using the flood traffic demand (FL). Two individual traffic demands were assigned to each node from each node.



**Figure 6.6:** Results of CQR Routing using the flood traffic demand (FL) with Hot-spots. Two individual traffic demands were assigned to each node from each node. Five percent of the nodes were made hot-spots, and they are marked as red in the figure.



**Figure 6.7:** Results of Enhanced CQR Routing using the flood traffic demand (FL). Two individual traffic demands were assigned to each node from each node.



**Figure 6.8:** Results of Enhanced CQR Routing using the flood traffic demand (FL) with Hot-spots. Two individual traffic demands were assigned to each node from each node. Five percent of the nodes were made hot-spots, and they are marked as red in the figure.

Routing Algorithm	FL	FL-HS
MA	87.2%	53.7
CQR	86.0%	51.0
ECQR	81.8%	49.1

**Table 6.5:** *This table demonstrates the average utilizations for each routing algorithm given the 5-ary 3-cube topology of the simulation environment.*

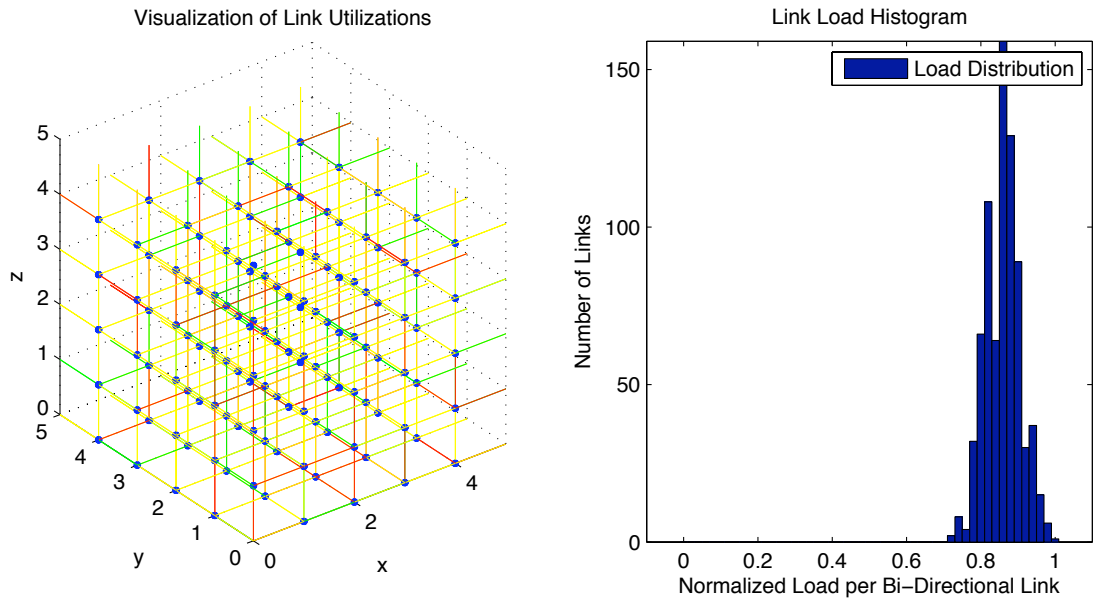
Routing Algorithm	FL	FL-HS
MA	4.5	9.0
CQR	4.6	9.3
ECQR	5.2	9.3

**Table 6.6:** *This table demonstrates the standard deviations for each routing algorithm given the 5-ary 3-cube topology of the simulation environment.*

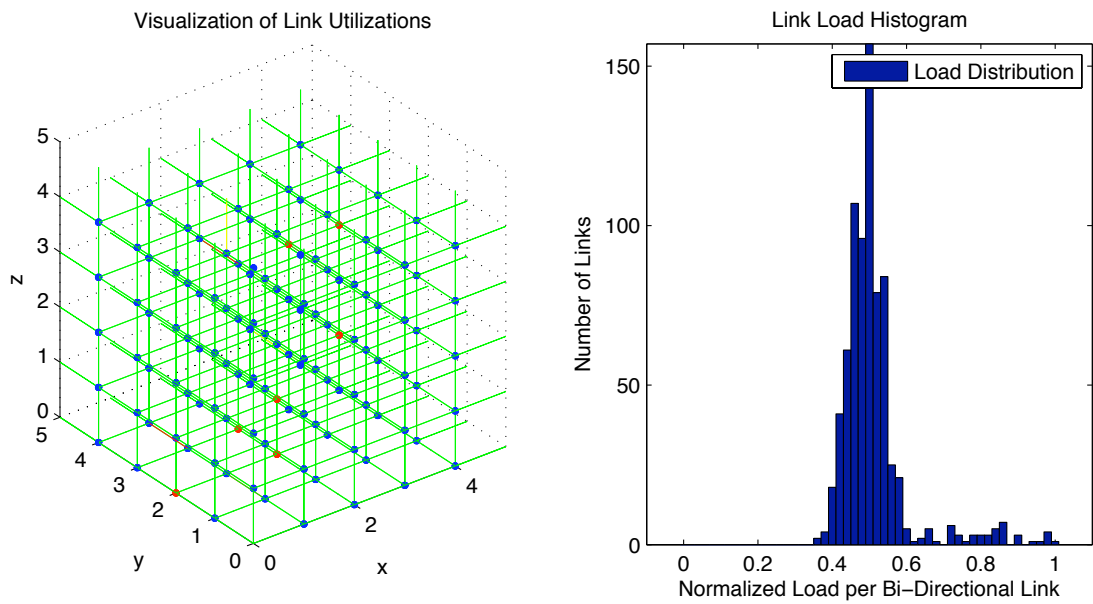
## 6.4 Results from the 5-ary 3-cube Simulations

Tables 6.5 and 6.6 show all results obtained from the simulations in a 5-ary 3-cube hypercube topology. The first table shows the average utilizations, and the second table shows the standard deviation for that specific simulation. The graphs of these results are available in Appendix E.3.

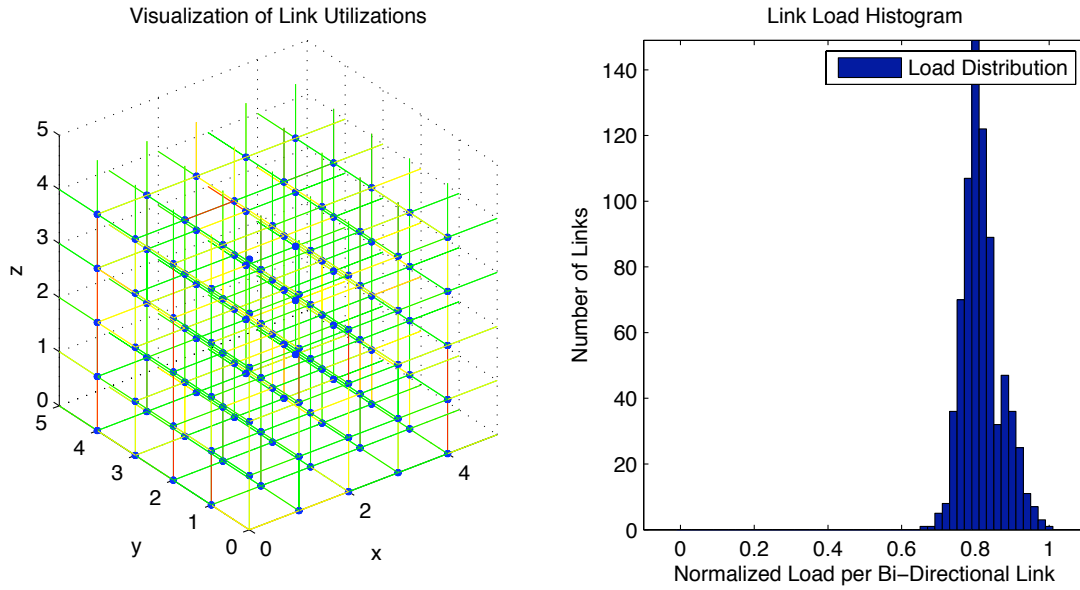
These tables are much smaller than the previous two topologies, solely because of the large execution times involved in the simulations of this network size. Because of this, only the adaptive algorithms are shown, along with the FL and FL-HS traffic patterns. All others have been excluded.



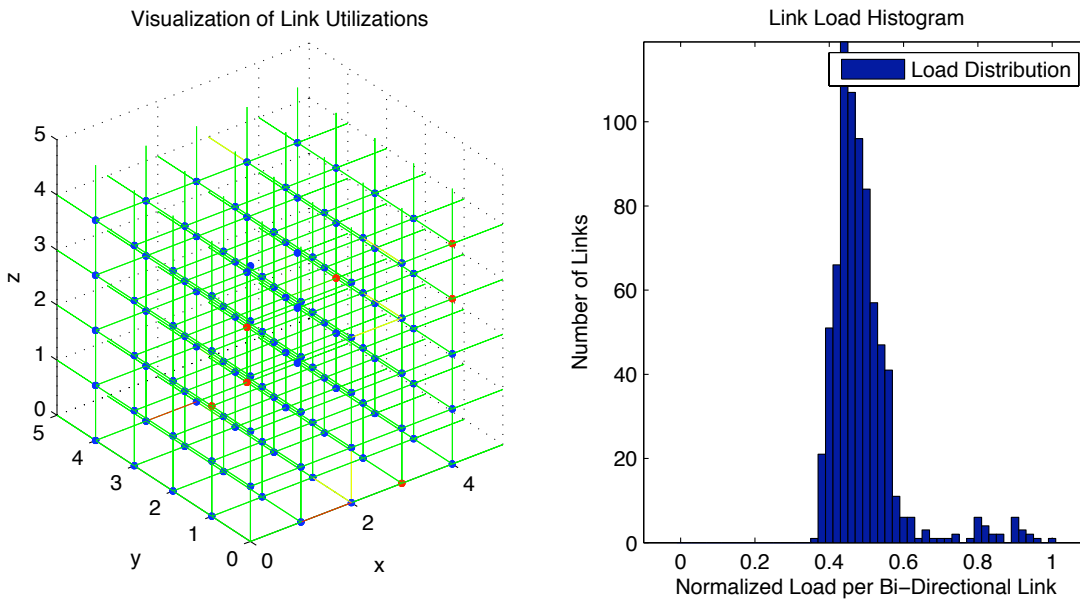
**Figure 6.9:** Results of CQR Routing using the flood traffic demand (FL). Two individual traffic demands were assigned to each node from each node.



**Figure 6.10:** Results of CQR Routing using the flood traffic demand (FL) with Hot-spots. Two individual traffic demands were assigned to each node from each node. Five percent of the nodes were made hot-spots, and they are marked as red in the figure.



**Figure 6.11:** Results of Enhanced CQR Routing using the flood traffic demand (FL). Two individual traffic demands were assigned to each node from each node.



**Figure 6.12:** Results of Enhanced CQR Routing using the flood traffic demand (FL) with Hot-spots. Two individual traffic demands were assigned to each node from each node. Five percent of the nodes were made hot-spots, and they are marked as red in the figure.

### 6.4.1 Comparison of All Results

Taking into consideration the laboratory results obtained within Chapter 5, the results obtained from these simulations should indicate a more-or-less best-case representation of the algorithms. Within the simulation environment the various OSI layers have been ignored, while they definitely had impact on the laboratory experiments. For instance, the layer 3 functionality such as window-sizing and segmentation were not considered within these simulations.

It is visible in the results that there is balancing game that occurs between ensuring locality and efficient load-balancing between the fully-adaptive algorithms. Because of Minimal Adaptive's inability to route non-minimally, its Nearest Neighbor (NN) results showed 100% utilization, while CQR and ECQR have to make decisions as to when to route non-minimally in order to balance the load more efficiently. The values which cause the change-over between minimal/non-minimal between CQR/ECQR should be further analyzed, as the efficiency could possibly increase (thus, increasing the average utilization).

It's also worth pointing out some scaling issues which may be occurring. The results have indicated that with the 3-ary 3-cube topology, ECQR increased performance exceeded that of CQR. As the network became larger with the 4-ary 3-cube and 5-ary 3-cube topologies, this difference shifted. This could indicate that the enhancements made within this work do not scale (or aren't shown to scale here), or that there exists a certain scenario or specific ratio of hotspots to non-hotspots, which depict this increased performance.

The bit complement (BC) traffic pattern has an interesting affect on the Direction Routing Algorithm (DIR), in that its level of performance is lower than that of Dimension Ordered Routing (DOR). These oddities are confirmed on the various topologies, as well as through the laboratory experiments.

Finally, the tornado (TOR) traffic pattern has been shown (both here, and in previous work) to be considerably difficult to respond to. Neither of these algorithms outperformed another when implementing this traffic pattern, as was expected.



# Chapter 7

## Conclusions

As interconnect networks become more prevalent within electrical and computer engineering, research such as this will continue to help drive towards more efficient routing implementations. Routing decisions, which is the main topic of discussion here, is only one contribution towards higher efficient and more productive data transfer between systems. Whether these systems are cores of a multi-core processor, or nodes within a HPC system, the method for data exchange between the individual systems is essential to increasing performance.

This research thesis contributed to this efficiency by:

- Implementing a simulation-based environment within Matlab to compare the various algorithms' implementations within hypercube topologies;
- Building a research cluster: a scaled version of the Red Storm/Thor's Hammer HPC;
- Implementing a torus-mesh network on this research cluster using a USB interconnect network;
- Implementing various static, oblivious, and adaptive routing algorithms using a set of socket-layer C programs, validating and verifying previous work, and demonstrating this work's implementation differences;
- Introducing Enhanced CQR with Periphery Avoidance, which makes intermediate routing decisions by routing productively but avoiding the perimeter when possible;

- Discussing the impact of this new enhancement and demonstrate evidence of its potential.

The Matlab simulation environment which was developed for this research work, and came from a basic implementation courtesy of Dr. Don Gruenbacher, was thoroughly extended from that early state to fully implement a broad set of routing algorithms and network traffic demands. All extensions which were necessary to this work include:

- Adding a user-friendly GUI, enabling user-input values to quickly expedite the simulation;
- Extending the routing capabilities from only static and oblivious routing to adaptive and fully-adaptive routing simulation capabilities;
- Adding capabilities for generating hotspots and nodal failures during simulations; and
- Adding analysis metrics such as usage, and delay for further analysis.

Singularly using only the Matlab scripts or the laboratory results from Thor's Tack Hammer may not have been sufficient for a high level of confidence in reporting the results of this thesis. Therefore both were invaluable in the analysis of both previous and new routing algorithms.

The research cluster was built from scratch, all commercially available components, and was shown to be an efficient way to implement a small parallel system. The power consumed during execution also suggested that this particular cluster could be very attractive to those who wish to dip their toes into the pool of high performance computing clusters, while not spending too much in implementation and maintenance. Also because of its low power consumption, no extra means for temperature control were necessary (as they would be in most cases).

The results which came from the Matlab simulations and the laboratory results indicated that the concept of periphery avoidance within routable quadrants was one which enabled

better hotspot avoidance. This was clear for the 3-ary 3-cube topology, but was unclear on larger topologies within the simulation environment. The advantages may not be as necessary for larger quadrants, as the number of minimal paths between source and destination increase as the quadrant sizes increase – therefore, for larger networks, Periphery Avoidance did not have a positive impact on CQR routing.

Further work should be done to analyze the complete impact of including these new enhancements for intermediate routing decisions while data progresses through a routable quadrant. This work has shown that enhancing CQR to include a function of periphery avoidance helped to avoid hotspots within the network.

## 7.1 Future Work

There are many possible recommendations for future work as it relates to this particular work, but some very interesting ideas include:

- *Analyzing the effects of layer-4 implementations.* Though this work used TCP exclusively within the laboratory tests, it would be worth analyzing the true impact of another layer-4 implementation.
- *Implementing a readily-available dynamic routing implementation such as Zebra and Quagga.* By analyzing the routing algorithms on another routing daemon, further validation could be achieved.
- *Implementing the dynamic algorithm by using Linux IP Routing Tables.* Though assumptions were made indicating that such an implementation would be computationally taxing, it would be worth analyzing further and seeing the direct impact of such an implementation.
- *Apply minimization or optimal control techniques.* By doing this, and applying weights towards the variables of periphery avoidance and output queues, an optimum decision

function could be found to help further increase Enhanced CQR's ability to route around hotspots while effectively balancing the network load.

- *Analyzing other interconnect implementations such as ethernet, firewire, or SATA.* All had been discussed as possibilities, but the implementation of either could make a research cluster such as Thor's Tack Hammer more appealing to a variety of users.
- *Hybridizing this research with that of overlay networks.* By using these new fully-adaptive routing decisions, routing within overlay networks could benefit from this research.
- *Including the evaluation of flow control operations in the presence of hotspots and difficult traffic patterns.* Though this research direction was not included within this work, the implementation of good flow control mechanisms can definitely improve the efficiency of a parallel computing cluster. Further analysis combining the enhancements to CQR with good flow-control practices could further demonstrate the benefits of implementing adaptive routing algorithms.
- *Analyzing the effects of neighboring hotspots verses non-neighboring hotspots.* Though this was not analyzed within this work, analysis of hotspots of neighboring nodes could impact the results of any hotspot analysis.
- *Further analyze the scalability of Enhanced CQR with Periphery Avoidance.* As it was indicated within the previous chapter, larger networks did not indicate the benefits of using ECQR as they did with the 3-ary 3-cube network (both simulation and laboratory results). Also, including more emphasis in analytically determining at what percentage of hot-spots ECQR's benefits were clearly seen has yet to be found. This could also turn into another minimization or optimal control problem.

# Bibliography

- [1] N. Adiga, M. Blumrich, D. Chen, P. Coteus, and A. Gara. Blue gene/l torus interconnection network. *IBM J. Res. and Dev*, Jan 2005.
- [2] G. Almasi, C. Archer, J. Castanos, and J. Gunnels. Design and implementation of message-passing services for the blue gene/l supercomputer. *IBM Journal of Research and Development*, Jan 2005.
- [3] P. Berman, L. Gravano, G. Pifarre, and J. Sanz. Adaptive deadlock-and livelock-free routing with all minimal paths in torus networks. *Proceedings of the fourth annual ACM symposium on Parallel Algorithms and Architectures*, Jan 1992.
- [4] M. Blumrich, D. Chen, P. Coteus, A. Gara, and M. Giampapa. Design and analysis of the bluegene/l torus interconnection network. *IBM Research Report RC23025 (W0312-022) December*, Jan 2003.
- [5] K. Bolding, M. Fulgham, and L. Snyder. The case for chaotic adaptive routing. *Computers*, Jan 1997.
- [6] Camp, B. Cole, E. DeBenedictis, R. Leland, J. Tomkins, and et. al. Architectural specification for massively parallel computers: an experience and measurement-based approach. *Concurrency and Computation Practice and Experience*, Jan 2005.
- [7] W. Dally, H. Aoki, and C. MIT. Deadlock-free adaptive routing in multicomputer networks using virtual channels. *Parallel and Distributed Systems*, Jan 1993.
- [8] W. J. Dally. *Principles and Practices of Interconnection Networks*. Morgan Kaufmann, Jan 2004.

- [9] J. Duato, S. Yalamanchili, and L. Ni. *Interconnection Networks*. Morgan Kaufmann, Jan 2003.
- [10] A. Hoisie, G. Johnson, D. Kerbyson, M. Lang, and S. Pakin. A performance comparison through benchmarking and modeling of three leading supercomputers: blue gene/l, red storm, and purple. *Proceedings of the 2006 ACM/IEEE conference on Supercomputing*, Jan 2006.
- [11] D. Mills. A brief history of ntp time: confessions of an internet timekeeper. *ACM Computer Communications Review*, Jan 2003.
- [12] A. Singh, W. Dally, A. Gupta, and B. Towles. Goal: a load-balanced adaptive routing algorithm for torus networks. *Computer Architecture*, Jan 2003.
- [13] A. Singh, W. Dally, A. Gupta, and B. Towles. Adaptive channel queue routing on k-ary n-cubes. *Proceedings of the sixteenth annual ACM symposium on Parallel Algorithms and Architectures*, Jan 2004.
- [14] A. Singh, W. Dally, B. Towles, and A. Gupta. Globally adaptive load-balanced routing on tori. *Computer Architecture Letters*, Jan 2004.
- [15] L. Valiant. A scheme for fast parallel communication. *SIAM Journal on Computing*, Jan 1982.
- [16] J. Vetter, S. Alam, T. D. Jr, and M. Fahey. Early evaluation of the cray xt3 at ornl. *Proceedings of the 47th Cray User Group Conference*, Jan 2005.

# Appendix A

## Appendix A: Linux BASH Scripts

### A.1 USB Interconnect Scripts

This first section lists the BASH scripts used to correctly initialize the USB interconnection network on Thor's Tack Hammer. At the head of each script contains its functionality.

The file triplets.txt, which was used within these scripts but was not included, simply listed the hardware address for each node, and the corresponding x, y, and z values. It also included the order for which each interface was brought up within Linux. It was found that this was the same order for each node.

#### A.1.1 file: interconnectUp

```
1 #!/bin/bash
2
3 # ++++++ #
4 # File: interconnectUp
5 # Author: Chris Lydick
6 # Date: Feb 18, 2008
7 #
8 # This script is run on the head-node and brings up
9 # all nodes' usb interfaces and then pings them. All
10 # information is routed to a log.txt file for analysis.
11 #
12 # ++++++ #
13
14
15 # These variables allow the script to color-code the success/failures.
16 NORMAL="^[[0;39m"
17 RED="^[[1;31m"
18 GREEN="^[[1;32m"
19
20 # Print the date at the top of the log.txt file.
21 echo `date` >> log.txt
22
23 # Copy the triplets.txt file if it is not in the default location.
24 if [ -e /pkhome/pkhome/triplets.txt ]; then
25 echo "/pkhome/pkhome/triplets.txt was found."
26 else
27 cp -f triplets.txt /pkhome/pkhome/triplets.txt
28 echo "/pkhome/pkhome/triplets.txt was not found."
29 echo "Copied successfully."
```

```

30 fi
31
32 # For all nodes in the LAM, bring up the interfaces. If
33 # an '.oops.usbup' file was touched, we see that as a flag indicating an error.
34 for i in `cat /pkhome/pkhome/tmp/bhosts|head -n 27`
35 do
36 echo -n "Bringing up USB interfaces for $i..."
37 echo "Bringing up USB interfaces for $i..." >> log.txt
38 ssh knoppix@$i "sudo /pkhome/pkhome/usbUp" >> log.txt
39 if [ -e /pkhome/pkhome/tmp/.oops.usbup ]; then
40 echo "$RED failed. $NORMAL"
41 echo "failed." >> log.txt
42 rm -f /pkhome/pkhome/tmp/.oops.usbup
43 else
44 echo "$GREEN success. $NORMAL";
45 echo "success." >> log.txt
46 fi
47 done
48
49 # Wait 10 seconds for the routing tables to fully initialize before pingng
50 # the neighbors.
51 sleep 10
52
53 # Now, ping the nodes... output all information to the logs.
54 for i in `cat /pkhome/pkhome/tmp/bhosts|head -n 27`
55 do
56 echo "Pinging hosts from $i."
57 echo "Pinging hosts from $i." >> log.txt
58 ssh knoppix@$i "sudo /pkhome/pkhome/usbPing" >> log.txt
59 done
60
61 # User interface.
62 echo ""
63 echo "You may want to verify all interfaces were brought up without"
64 echo "error by viewing the file log.txt."
65
66 # This makes partitioning the logs simple.
67 echo "===== " >> log.txt
68 echo "" >> log.txt
69 echo "" >> log.txt

```



## A.1.2 file: usbUp

```
1 #/bin/bash
2
3 # ++++++ #
4 # File: usbUp
5 # Author: Chris Lydick
6 # Date: Feb 18, 2008
7 #
8 # This script brings up the USB Ethernet Interfaces and
9 # automatically assigns the addresses/subnets.
10 #
11 # ++++++ #
12
13 k=3
14 n=3
15 PRIMARYOCTET="10.0"
16 USBSUBNET="255.255.255.252"
17
18 getTriplet()
19 {
20
21 # First find our management hardware address
22 for i in `ifconfig eth2 |grep HWaddr`
23 do
24 if [ `echo $i |head -c 5` == "00:E0" ]; then
25 hwaddr=$i;
26 fi
27 done
28
29 # Given the hardware address, lookup our triplet information in a
30 # global file, "triplets.txt"
31 triplet=`cat /home/knoppix/triplets.txt |grep $hwaddr |tail -c 6`
32
33 # get x,y,z
34 z=`echo $triplet |head -c 1`
35 x=`echo $triplet |tail -c 2`
36 y=`echo $triplet |tail -c 4 |head -c 1`
37 echo "triplet: $x,$y,$z"
38
39 # This next area allows for the triplets.txt file to contain
40 # information as to the specific order for usb0-5 as they're
41 # brought up. If all nodes are correctly booted, they should
42 # have a consistent order.
43 order=`cat /home/knoppix/triplets.txt |grep $hwaddr |head -c 11 |tr ' .' ' '`
44 first=`echo ${order:0:1}`;
45 second=`echo ${order:2:1}`;
46 third=`echo ${order:4:1}`;
47 fourth=`echo ${order:6:1}`;
48 fifth=`echo ${order:8:1}`;
49 sixth=`echo ${order:10:1}`;
50 }
51
52 # This function calculates both local and neighbor addresses
53 # for the x-directions (+/-).
54 getXs()
55 {
56 dir_mod_pos=1;
57 dir_mod_neg=2;
58 POSCOCTET=$((($x<<4)+$y);
59 if [ $x -eq 0 ]; then
60 NEGCOCCTET=$((($k-1)<<4)+$y);
61 else
62 NEGCOCCTET=$((($x-1)<<4)+$y);
63 fi
```

```

64 POSDOCTET=$((z<<4)+dir_mod_pos];
65 NEGDOCTET=$((z<<4)+dir_mod_neg];
66 POSXIP="$PRIMARYOCTET.$POSCOCTET.$POSDOCTET";
67 NEGXIP="$PRIMARYOCTET.$NEGCOCTET.$NEGDOCTET";
68 NEXTPOSX="$PRIMARYOCTET.$POSCOCTET.${POSDOCTET+1}";
69 NEXTNEGX="$PRIMARYOCTET.$NEGCOCTET.${NEGDOCTET-1}";
70 echo "The Positive X interface is: $POSXIP";
71 echo "The Negative X interface is: $NEGXIP";
72 }
73
74 # This function calculates both local and neighbor addresses
75 # for the y-directions (+/-).
76 getYs()
77 {
78     dir_mod_pos=5;
79     dir_mod_neg=6;
80     POSCOCTET=$((x<<4)+y];
81     if [ $y -eq 0 ]; then
82     NEGCOCTET=$((x<<4)+(k-1)];
83     else
84     NEGCOCTET=$((x<<4)+(y-1)];
85     fi
86     POSDOCTET=$((z<<4)+dir_mod_pos];
87     NEGDOCTET=$((z<<4)+dir_mod_neg];
88     POSYIP="$PRIMARYOCTET.$POSCOCTET.$POSDOCTET";
89     NEGYIP="$PRIMARYOCTET.$NEGCOCTET.$NEGDOCTET";
90     NEXTPOSY="$PRIMARYOCTET.$POSCOCTET.${POSDOCTET+1}";
91     NEXTNEGY="$PRIMARYOCTET.$NEGCOCTET.${NEGDOCTET-1}";
92     echo "The Positive Y interface is: $POSYIP";
93     echo "The Negative Y interface is: $NEGYIP";
94 }
95
96 # This function calculates both local and neighbor addresses
97 # for the z-directions (+/-).
98 getZs()
99 {
100     dir_mod_pos=9;
101     dir_mod_neg=10;
102     COCTET=$((x<<4)+y];
103     POSDOCTET=$((z<<4)+dir_mod_pos];
104     if [ $z -eq 0 ]; then
105     NEGDOCTET=$((k-1)<<4)+dir_mod_neg];
106     else
107     NEGDOCTET=$((z-1)<<4)+dir_mod_neg];
108     fi
109     POSZIP="$PRIMARYOCTET.$COCTET.$POSDOCTET";
110     NEGZIP="$PRIMARYOCTET.$COCTET.$NEGDOCTET";
111     NEXTPOSZ="$PRIMARYOCTET.$COCTET.${POSDOCTET+1}";
112     NEXTNEGZ="$PRIMARYOCTET.$COCTET.${NEGDOCTET-1}";
113     echo "The Positive Z interface is: $POSZIP";
114     echo "The Negative Z interface is: $NEGZIP";
115 }
116
117 getTriplet;
118 getXs;
119 getYs;
120 getZs;
121
122 # List all usb interfaces in $usbs, and count
123 for i in `ifconfig -a |grep usb`
124 do
125     if [ `echo $i |head -c 3` == "usb" ]; then
126     count=`echo ${count+1}`;
127     usbs="$usbs $i";
128     fi

```

```

129 done
130
131 # Bring up usb interfaces
132 count_x=0;
133 if [ $count == 6 ]; then
134 #echo $usbs
135 for j in $usbs
136 do
137 echo -n "Bringing up interface $j...";
138 sudo ifconfig $j up;
139 case "$count_x" in
140 "$first" ) sudo ifconfig $j inet $POSYIP netmask $USBSUBNET;;
141 "$second" ) sudo ifconfig $j inet $NEGYIP netmask $USBSUBNET;;
142 "$third" ) sudo ifconfig $j inet $POXIP netmask $USBSUBNET;;
143 "$fourth" ) sudo ifconfig $j inet $NEGXIP netmask $USBSUBNET;;
144 "$fifth" ) sudo ifconfig $j inet $NEGZIP netmask $USBSUBNET;;
145 "$sixth" ) sudo ifconfig $j inet $POSZIP netmask $USBSUBNET;;
146 * ) echo "Found an extra case in the loop.";;
147 esac
148 [ $? -eq 0 ] && echo "success." || touch /home/knoppix/tmp/.oops.usbup
149 count_x='echo ${count_x+1}';
150 done
151 # This occurs when there aren't exactly 6 usb interfaces. Consider rebooting node.
152 # A file is touched which can be seen by the head-node, indicating an error.
153 else
154 echo "A problem occurred. Not all interfaces appear to be working."
155 echo "Only $count interfaces were found."
156 echo "Exiting..."
157 touch /home/knoppix/tmp/.oops.usbup
158 fi

```

### A.1.3 file: usbPing

```
1 #/bin/bash
2
3 # ++++++ #
4 # File: usbPing
5 # Author: Chris Lydick
6 # Date: Feb 18, 2008
7 #
8 # This script allows a node to ping all of its closest
9 # neighbors (within 1 hop).
10 #
11 # ++++++ #
12
13 k=3
14 n=3
15
16 PRIMARYOCTET="10.0"
17 USBSUBNET="255.255.255.252"
18
19 getTriplet()
20 {
21 # First find our management hardware address
22 for i in `ifconfig eth2 |grep HWaddr`
23 do
24 if [ `echo $i |head -c 5` == "00:E0" ]; then
25 hwaddr=$i;
26 fi
27 done
28
29 # Given the hardware address, lookup our triplet information in a
30 # global file, "triplets.txt"
31 triplet=`cat /home/knoppix/triplets.txt |grep $hwaddr |tail -c 6`
32
33 # get x,y,z
34 z=`echo $triplet |head -c 1`
35 x=`echo $triplet |tail -c 2`
36 y=`echo $triplet |tail -c 4 |head -c 1`
37 echo "triplet: $x,$y,$z"
38 }
39
40 # This function calculates both local and neighbor addresses
41 # for the x-directions (+/-).
42 getXs()
43 {
44 dir_mod_pos=1;
45 dir_mod_neg=2;
46 POSCOCTET=$((($x<<4)+$y);
47 if [ $x -eq 0 ]; then
48 NEGCOCTET=$((($k-1)<<4)+$y);
49 else
50 NEGCOCTET=$((($x-1)<<4)+$y);
51 fi
52 POSDOCTET=$((($z<<4)+$dir_mod_pos);
53 NEGDOCTET=$((($z<<4)+$dir_mod_neg);
54 POSXIP="$PRIMARYOCTET.$POSCOCTET.$POSDOCTET";
55 NEGXIP="$PRIMARYOCTET.$NEGCOCTET.$NEGDOCTET";
56 NEXTPOSX="$PRIMARYOCTET.$POSCOCTET.$[$POSDOCTET+1]";
57 NEXTNEGX="$PRIMARYOCTET.$NEGCOCTET.$[$NEGDOCTET-1]";
58 }
59
60 # This function calculates both local and neighbor addresses
61 # for the y-directions (+/-).
62 getYs()
63 {
```

```

64  dir_mod_pos=5;
65  dir_mod_neg=6;
66  POSCOCTET=$((x<<4)+$y);
67  if [ $y -eq 0 ]; then
68  NEGCOCTET=$((x<<4)+($k-1));
69  else
70  NEGCOCTET=$((x<<4)+($y-1));
71  fi
72  POSDOCTET=$((z<<4)+$dir_mod_pos];
73  NEGDOCTET=$((z<<4)+$dir_mod_neg];
74  POSYIP="$PRIMARYOCTET.$POSCOCTET.$POSDOCTET";
75  NEGYP="$PRIMARYOCTET.$NEGCOCTET.$NEGDOCTET";
76  NEXTPOSY="$PRIMARYOCTET.$POSCOCTET.${POSDOCTET+1}";
77  NEXTNEGY="$PRIMARYOCTET.$NEGCOCTET.${NEGDOCTET-1}";
78  }
79
80  # This function calculates both local and neighbor addresses
81  # for the z-directions (+/-).
82  getZs()
83  {
84  dir_mod_pos=9;
85  dir_mod_neg=10;
86  COCTET=$((x<<4)+$y);
87  POSDOCTET=$((z<<4)+$dir_mod_pos];
88  if [ $z -eq 0 ]; then
89  NEGDOCTET=$((($k-1)<<4)+$dir_mod_neg];
90  else
91  NEGDOCTET=$((($z-1)<<4)+$dir_mod_neg];
92  fi
93  POSZIP="$PRIMARYOCTET.$COCTET.$POSDOCTET";
94  NEGZIP="$PRIMARYOCTET.$COCTET.$NEGDOCTET";
95  NEXTPOSZ="$PRIMARYOCTET.$COCTET.${POSDOCTET+1}";
96  NEXTNEGZ="$PRIMARYOCTET.$COCTET.${NEGDOCTET-1}";
97  }
98
99  getTriplet;
100 getXs;
101 getYs;
102 getZs;
103
104 # Ping all neighbors - we filter out all lines except where we successfully
105 # or unsuccessfully transmitted one ICMP packet.
106 echo -n "+x $NEXTPOSX: "
107 ping -q -c 1 $NEXTPOSX |grep transmitted
108 echo -n "-x $NEXTNEGX: "
109 ping -q -c 1 $NEXTNEGX |grep transmitted
110 echo -n "+y $NEXTPOSY: "
111 ping -q -c 1 $NEXTPOSY |grep transmitted
112 echo -n "-y $NEXTNEGY: "
113 ping -q -c 1 $NEXTNEGY |grep transmitted
114 echo -n "+z $NEXTPOSZ: "
115 ping -q -c 1 $NEXTPOSZ |grep transmitted
116 echo -n "-z $NEXTNEGZ: "
117 ping -q -c 1 $NEXTNEGZ |grep transmitted

```

## A.2 Data Analysis Scripts

The following scripts enabled efficient analysis of the data that was to be sent or had been sent across the USB Interconnect. As with before, the head of each script contains information describing its functionality.

### A.2.1 file: interconnectTime

```
1 #!/bin/bash
2
3 # ++++++ #
4 # File: interconnectTime
5 # Author: Chris Lydick
6 # Date: Mar 3, 2008
7 #
8 # This script reports the estimated and maximum error
9 # of time from the headnode (using NTP).
10 #
11 # ++++++ #
12
13
14 for i in `cat /pkhome/pkhome/tmp/bhosts | head -n 27`
15 do
16 echo $i
17 ssh $i "ntptime | grep error"
18 done
```

## A.2.2 file: interconnectRouters

```
1 #!/bin/bash
2
3 # ++++++ #
4 # File: interconnectRouters
5 # Author: Chris Lydick
6 # Date: Mar 3, 2008
7 #
8 # This script compiles the most recent version of server.c,
9 #   loser.c, and injector.c, and executes the script nodeSetup
10 #   locally. It also grabs usb byte counts (line 27).
11 #
12 # ++++++ #
13
14 # This script first kills all loser and server processes
15 #   and then restarts the server, sending output to /dev/null.
16 cd /mnt/sda2/usb/InterconnectScripts/
17
18 gcc -o /pkhome/pkhome/loser server/loser.c
19 gcc -o /pkhome/pkhome/server server/server.c -lm
20 gcc -o server/injector server/injector.c
21
22 j=26;
23 for i in `cat /pkhome/pkhome/tmp/bhosts |head -n 27`
24 do
25 echo "$i, addr:$j"
26 ssh $i "/pkhome/pkhome/nodeSetup $j"
27 ssh $i "cat /proc/net/dev |grep usb" > server/logs/$i-before
28 j=$((j-1));
29 done
30
```

## A.2.3 file: nodeSetup

```
1 #!/bin/bash
2
3 # ++++++ #
4 # File: nodeSetup
5 # Author: Chris Lydick
6 # Date: Mar 3, 2008
7 #
8 # This script first kills all previous server/loserver
9 # processes and then re-initializes them with the newly
10 # compiled versions. It then sets the CPU affinity for
11 # "server" exclusively to the first CPU and "loserver"
12 # exclusively to the second CPU. All output is redirected
13 # to /dev/null, and servers are placed in background.
14 #
15 # ++++++ #
16
17 sudo pkill server
18 /pkhomes/pkhomes/loserver $1 1>/dev/null &
19 /pkhomes/pkhomes/server $1 1>/dev/null &
20 taskset -p 01 'pgrep -x server' 1>/dev/null
21 taskset -p 02 'pgrep -x loserver' 1>/dev/null
22
23 exit
24
```



## A.2.4 file: interconnectThroughput

```
1 #!/bin/bash
2
3 # ++++++ #
4 # File: interconnectThroughput
5 # Author: Chris Lydick
6 # Date: Mar 3, 2008
7 #
8 # This script parses all data that was sent during a run
9 # and returns statistics given their timestamps.
10 #
11 # ++++++ #
12
13 # This check verifies we've not run it since resetting the
14 # USB byte counts.
15 if [ -e server/logs/192.168.0.200-before ]; then
16 echo -n ""
17 else
18 echo "No before files found. You must re-run interconnectRouters to regenerate these files."
19 exit;
20 fi
21
22 echo -n "Getting latest USB byte counts..."
23
24 # For each node, get the latest USB byte counts.
25 for i in `cat /pkhome/pkhome/tmp/bhosts |head -n 27`
26 do
27 ssh $i "cat /proc/net/dev |grep usb" > server/logs/$i-after
28 done
29 echo "done."
30
31 allnums="0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26"
32 nums="1 2 3 4 5 6"
33 diffs=""
34
35 # For each node...
36 for i in `cat /pkhome/pkhome/tmp/bhosts |head -n 27`
37 do
38 # And for each of the 6 interfaces...
39 for j in $nums
40 do
41 # Get the jth byte count at the start.
42 var1=`cat server/logs/$i-before |head -n $j |tail -n 1|tr ":" " " `
43 set -- $var1
44 shift; shift; shift; shift; shift; shift; shift; shift; shift; shift;
45 trans_before=$1;
46
47 # Get the jth byte count at the end of run.
48 var2=`cat server/logs/$i-after |head -n $j |tail -n 1|tr ":" " " `
49 set -- $var2
50 shift; shift; shift; shift; shift; shift; shift; shift; shift; shift;
51 trans_after=$1;
52
53 # Only analyze the transmitted data. To do both would be redundant.
54 trans=$((trans_after-trans_before));
55 diffs="$diffs $trans"
56 done
57
58 done
59
60 # For ease, the differences between before-bytes and after-bytes can be copied/pasted
61 # to Matlab.
62 echo "The following can be pasted to MatLab for a histogram of the link utilizations."
63 echo ""
```

```

64 echo ""
65 #diffs='echo $diffs |tr " " "\n"'
66 echo "x = [$diffs]"
67 echo "x = x./(max(x))"
68 echo "x = [x 0]"
69 echo "hist(x,100)";
70 echo ""
71 echo ""
72
73 # Remove these logs, force the user to re-run interconnectRouters script.
74 rm -fr server/logs/*
75
76 # This loop was necessary for large numbers of packets. 'cat' has a finite space
77 # when it comes to data that is passed to it.
78 cd /pckhome/pckhome/router/routerlogs/
79 for i in $allnums
80 do
81 hop_1="$hop_1 'echo "" |cat \grep -d recurse -l "hops: 1" $i/ \grep trans'"
82 hop_2="$hop_2 'echo "" |cat \grep -d recurse -l "hops: 2" $i/ \grep trans'"
83 hop_3="$hop_3 'echo "" |cat \grep -d recurse -l "hops: 3" $i/ \grep trans'"
84 avg_queue="$avg_queue 'echo "" |cat \find $i/ |grep SUCCESS\grep avg_queue'"
85 done
86
87 # First find all data dropped by the C programs.
88 alldrops='find . |grep DROP'
89 set -- $alldrops
90 echo "actual drops: $#";
91 rm -f 'find . |grep DROP' || echo -n "";
92
93 # Find all unknown data dropped by the C programs.
94 allunknown='find . |grep UNKNOWN'
95 set -- $allunknown
96 echo "actual unknown: $#";
97 rm -f 'find . |grep UNKNOWN' || echo -n "";
98
99 # Parse the average queueing times.
100 set -- $avg_queue
101 total_drops=0;
102 queue=0;
103 total_packets=0;
104 i=0;
105 j=0;
106
107 while [ $# -gt 0 ]
108 do
109 shift
110 i=${i+1};
111 queue=${queue+$i};
112 shift
113 done
114 if [ $i -gt 0 ]; then
115 queue=${queue/$i};
116 echo "avg queue time: $queue"
117 fi
118
119 # Parse the headers from data which traversed one hop.
120 set -- $hop_1
121 hops=0;
122 i=0;
123 j=0;
124 while [ $# -gt 0 ]
125 do
126 shift
127 # If the time was greater than 1.5 seconds, we
128 # assume there was some kind of timeout. This

```

```

129 # data will be discarded, counted as a drop.
130 if [ $1 -lt 1500000 ]; then
131 i=${i+1};
132 hops=${hops+$1};
133 else
134 j=${j+1};
135 fi
136 shift
137 done
138 total_packets=${total_packets+$j+$i};
139 if [ $i -gt 0 ]; then
140 hops=${hops/$i};
141 echo "avg 1-hop time: $hops"
142 fi
143 total_drops=${total_drops+$j};
144 all_avg=${i*hops};
145
146 # Parse the headers from data which traversed two hops.
147 set -- $hop_2
148 hops=0;
149 i=0;
150 j=0;
151 while [ $# -gt 0 ]
152 do
153 shift
154 if [ $1 -lt 1500000 ]; then
155 i=${i+1};
156 hops=${hops+$1};
157 else
158 j=${j+1};
159 fi
160 shift
161 done
162 total_packets=${total_packets+$j+$i};
163 if [ $i -gt 0 ]; then
164 hops=${hops/$i};
165 echo "avg 2-hop time: $hops"
166 fi
167 total_drops=${total_drops+$j};
168 all_avg=${all_avg + ($hops*$i)};
169
170 # And finally, for three hops.
171 set -- $hop_3
172 hops=0;
173 i=0;
174 j=0;
175 while [ $# -gt 0 ]
176 do
177 shift
178 if [ $1 -lt 1500000 ]; then
179 i=${i+1};
180 hops=${hops+$1};
181 else
182 j=${j+1};
183 fi
184 shift
185 done
186 total_packets=${total_packets+$j+$i};
187 if [ $i -gt 0 ]; then
188 hops=${hops/$i};
189 echo "avg 3-hop time: $hops"
190 fi
191 total_drops=${total_drops+$j};
192 all_avg=${(all_avg + ($hops*$i))/(total_packets-total_drops)};
193

```

```
194 # Echo all global information
195 echo "avg trans time: $all_avg";
196 echo "total drops: $total_drops"
197 echo "total packets: $[${total_packets}/2]"
198
199 # Remove all data packets, we're done.
200 cd /pkhme/pkhme/router/routerlogs/
201 rm -f `find . |grep SUCCESS`
202
```

# Appendix B

## Appendix B: C Programs

### B.1 file: server.c

```
1 /* ##### *
2 * Filename: router.c *
3 * Author: Chris Lydick *
4 * Date: Mar 1, 2008 *
5 * Usage: ./server [address] *
6 * Notes: Portions borrowed from http://tinyurl.com/2w36o4 *
7 * and http://tinyurl.com/6bu7s *
8 * This is a program which was created to route data over *
9 * a torus-mesh/hypercube network using: *
10 * 1. Dimension Ordered Routing *
11 * 2. Direction Ordered Routing *
12 * 3. Minimal Adaptive Routing *
13 * 4. Minimal Oblivious Routing *
14 * 5. CQR Routing *
15 * 6. Enhanced CQR *
16 * ##### */
17
18 #include "router.h"
19
20 /* ===== */
21 /* ===== FUNCT DECLARATION ===== */
22 /* ===== */
23 void initNeighborhood(void);
24 void child_handler(int s);
25 void myperror(char *x);
26 int makeDecision(char *hdr, char *ra, char *dst);
27 void send_packet(char* packet, char* host, unsigned long long timestamp);
28 int dimord(int s, int d);
29 void sem_lock(int id);
30 void sem_unlock(int id);
31 int get_outputQueue(int sem_id, struct queue *q);
32 int inc_outputQueue(int sem_id, struct queue *q);
33 int dec_outputQueue(int sem_id, struct queue *q);
34 void adjust_queues(int sem_id, struct queue *q);
35 int min_choose(int x, int y, int z);
36 int adaptive_choose(int x, int y, int z);
37 int adaptive_periphery_choose(int x, int y, int z);
38 void make_vector(char v[], char dir[], int x_val, int y_val, int z_val);
39 int min(int x, int y);
40 double min_d(double x, double y);
41
42 /* ===== */
```

```

43 /* ===== GLOBAL VARIABLES ===== */
44 /* ===== */
45 int myx, myy, myz, k, n;
46 char *addr, next_posx[12], next_posy[12], next_posz[12], next_negx[12], next_negy[12], next_negz[12];
47 struct queue *q_posx, *q_posy, *q_posz, *q_negx, *q_negy, *q_negz;
48 int sem_posx, sem_posy, sem_posz, sem_negx, sem_negy, sem_negz;
49
50 /* ===== */
51 /* ===== Main Function ===== */
52 /* ===== */
53 // This is the main fuction which is executed.
54 // It's purpose is to:
55 // (1): Accept TCP Packets of fixed size on port 3490.
56 // (2): Analyzes the packet...
57 // (a): if: Hop Count is higher than MAXHOP: drop, route to localhost.
58 // (b): else if: packet is destined for this host: route to localhost
59 // (c): else: make routing decision, and send packet onto next hop.
60 int main(int argc, char* argv[])
61 {
62 int sockfd, new_fd, numbytes, hop_nu, remSocket, yes;
63 int shmid;
64 struct sockaddr_in my_addr, their_addr, next_addr;
65 union semun sem1, sem2, sem3, sem4, sem5, sem6;
66 struct hostent *remHost;
67 struct timeval t;
68 struct sigaction sa;
69 char* shmem;
70 char newbuf[MAXDATASIZE], buf[MAXDATASIZE], header[HEADERSIZE+1];
71 char h_src[3], h_dest[3], h_ra[2], h_hops[3], h_queue[7];
72 char h_queue_t[7], d_time[7], time_ms_t[17], time_ms[10];
73 unsigned long long t1, t2;
74 socklen_t sin_size;
75
76
77 k=3; n=3; yes=1;
78 if (argc > 1)
79 addr = argv[1];
80 else myperror("usage: ./server [node number]");
81
82 if ((sockfd = socket(AF_INET, SOCK_STREAM, 0)) == -1)
83 myperror("socket");
84 if (setsockopt(sockfd, SOL_SOCKET, SO_REUSEADDR, &yes, sizeof(int)) == -1)
85 perror("setsockopt");
86
87 // Create 6 semaphores, all with one semaphore each, readable by the owner.
88 // these correspond to the semaphores for each output queue on each USB interface
89 sem_posx = semget(SID_POSX, 1, IPC_CREAT | 0600);
90 sem_posy = semget(SID_POSY, 1, IPC_CREAT | 0600);
91 sem_posz = semget(SID_POSZ, 1, IPC_CREAT | 0600);
92 sem_negx = semget(SID_NEGX, 1, IPC_CREAT | 0600);
93 sem_negy = semget(SID_NEGY, 1, IPC_CREAT | 0600);
94 sem_negz = semget(SID_NEGZ, 1, IPC_CREAT | 0600);
95
96 // Now, initialize the semaphores to 1.
97 sem1.val = 1;
98 if (semctl(sem_posx, 0, SETVAL, sem1) == -1)
99 myperror("sem1");
100 sem2.val = 1;
101 if (semctl(sem_posy, 0, SETVAL, sem2) == -1)
102 myperror("sem2");
103 sem3.val = 1;
104 if (semctl(sem_posz, 0, SETVAL, sem3) == -1)
105 myperror("sem3");
106 sem4.val = 1;
107 if (semctl(sem_negx, 0, SETVAL, sem4) == -1)

```

```

108 myperror("sem4");
109 sem5.val = 1;
110 if (semctl(sem_negy, 0, SETVAL, sem5) == -1)
111 myperror("sem5");
112 sem6.val = 1;
113 if (semctl(sem_negz, 0, SETVAL, sem6) == -1)
114 myperror("sem6");
115
116 // Create a shared memory block of appropriate size.
117 shmId = shmget(IPC_PRIVATE, sizeof(struct queue)*6, IPC_CREAT | IPC_EXCL | 0600);
118 // Grab the address where it exists.
119 shmem = shmat(shmId, NULL, 0);
120
121 // Adjust the queue pointers to the appropriate positions within
122 // the shared memory.
123 q_posx = (struct queue*) shmem;
124 q_posy = (struct queue*) shmem+(sizeof(struct queue)*1);
125 q_posz = (struct queue*) shmem+(sizeof(struct queue)*2);
126 q_negx = (struct queue*) shmem+(sizeof(struct queue)*3);
127 q_negy = (struct queue*) shmem+(sizeof(struct queue)*4);
128 q_negz = (struct queue*) shmem+(sizeof(struct queue)*5);
129
130 // Set all the values of each queue to 0.
131 q_posx[0].val = 0;
132 q_negx[0].val = 0;
133 q_posy[0].val = 0;
134 q_negy[0].val = 0;
135 q_posz[0].val = 0;
136 q_negz[0].val = 0;
137
138 // Grab the time, and set all lastUpdated to the value of time now.
139 gettimeofday(&t, NULL);
140 t1 = (unsigned long long)t.tv_sec * 1000000 + (unsigned long long)t.tv_usec;
141 q_posx[0].lastUpdated = t1;
142 q_posy[0].lastUpdated = t1;
143 q_posz[0].lastUpdated = t1;
144 q_negx[0].lastUpdated = t1;
145 q_negy[0].lastUpdated = t1;
146 q_negz[0].lastUpdated = t1;
147
148 // Setup my server listen port.
149 my_addr.sin_family = AF_INET;
150 my_addr.sin_port = htons(ROUTEPORT);
151 my_addr.sin_addr.s_addr = INADDR_ANY;
152 memset(my_addr.sin_zero, '\0', sizeof my_addr.sin_zero);
153
154 // Bind a socket to my listen port.
155 if (bind(sockfd, (struct sockaddr *)&my_addr, sizeof my_addr) == -1)
156 myperror("bind");
157 // Allow a backlog of connections to BACKLOG
158 if (listen(sockfd, BACKLOG) == -1)
159 myperror("listen");
160 // Enable child_handler to reap all dead children lingering.
161 sa.sa_handler = child_handler;
162 sigemptyset(&sa.sa_mask);
163 sa.sa_flags = SA_RESTART;
164 if (sigaction(SIGCHLD, &sa, NULL) == -1)
165 myperror("sigaction");
166 // Initialize our neighborhood - calculate all IP addresses of interfaces and neighbors
167 initNeighborhood();
168 // main while loop.
169 while(1) {
170 sin_size = sizeof their_addr;
171 // Accept a new connection, immediately fork, and let child do work.
172 if ((new_fd = accept(sockfd, (struct sockaddr *)&their_addr, &sin_size)) == -1) {

```

```

173 perror("accept");
174 continue; }
175 if (!fork()) {
176 // child doesn't need the listener socket.
177 close(sockfd);
178 // receive the packet from the sender.
179 if ((numbytes=recv(new_fd, buf, MAXDATASIZE-1, 0)) == -1)
180 perror("recv");
181 // get time the packet was received
182 gettimeofday(&t, NULL);
183 t1 = (unsigned long long)t.tv_sec * 1000000 + (unsigned long long)t.tv_usec;
184 memcpy(newbuf, buf, MAXDATASIZE);
185 strncpy(&header[0], &newbuf[0], HEADERSIZE); header[HEADERSIZE+1] = '\0';
186 strncpy(&h_dest[0], &newbuf[2], 2); h_dest[2] = '\0';
187 strncpy(&h_ra[0], &newbuf[4], 1); h_ra[1] = '\0';
188 strncpy(&h_hops[0], &newbuf[5], 2); h_hops[2] = '\0';
189 strncpy(&h_queue[0], &newbuf[16], 6); h_queue[6] = '\0';
190 strncpy(&time_ms[0], &newbuf[7], 9); time_ms[9] = '\0';
191 // if the header does not contain a time value for beginning,
192 // insert the timestamp now.
193 if (atoi(time_ms) == 0) {
194 sprintf(time_ms_t, "%llu", t1);
195 strncpy(&time_ms[0], &time_ms_t[7], 9); time_ms[9] = '\0';
196 strncpy(&newbuf[7], &time_ms[0], 9);
197 }
198 // if this packet is destined for itself, or has exceeded the max
199 // number of hops, route the header to localhost & remove from network.
200 if ((atoi(h_dest) == atoi(addr)) || (atoi(h_hops) >= MAXHOPS)) {
201 remHost=gethostbyname("127.0.0.1");
202 remSocket=socket(AF_INET, SOCK_STREAM, 0);
203 next_addr.sin_family = AF_INET;
204 next_addr.sin_port = htons(LOPORT);
205 next_addr.sin_addr = *((struct in_addr *)remHost->h_addr);
206 gettimeofday(&t, NULL);
207 t2 = (unsigned long long)t.tv_sec * 1000000 + (unsigned long long)t.tv_usec;
208 if ((t2-t1) >= 10000)
209 sprintf(d_time, "%llu", t2-t1);
210 else
211 sprintf(d_time, "0%llu", t2-t1);
212 sprintf(h_queue_t, "%d", atoi(d_time) + atoi(h_queue));
213 strncpy (&header[16], &h_queue_t[0], 6);
214 if (connect(remSocket, (struct sockaddr *)&next_addr, sizeof(next_addr)) < 0)
215 perror("connecting to localhost");
216 if (send(remSocket, header, HEADERSIZE-1, 0) < 0)
217 perror("sending to localhost");
218 close(remSocket);
219
220 }
221 else {
222 // Fork off children to run the adjust_queues in parallel. That way if one queue
223 // is slow to return the semaphore, others aren't blocked.
224 if (atoi(h_ra) > 4) {
225 if (!fork()){
226 if (!fork()){
227 if (!fork()){
228 if (!fork()){
229 if (!fork()){
230 if (!fork()){
231 adjust_queues(sem_posx, q_posx);
232 exit(0);}
233 adjust_queues(sem_negx, q_negx);
234 exit(0);}
235 adjust_queues(sem_posy, q_posy);
236 exit(0);}
237 adjust_queues(sem_negy, q_negy);

```



```

238  exit(0);}
239  adjust_queues(sem_posz, q_posz);
240  exit(0);}
241  adjust_queues(sem_negz, q_negz);
242  exit(0);}
243  }
244  // add one to the hop.
245  newbuf[6] = newbuf[6]+1;
246  switch(makeDecision(newbuf,h_ra,h_dest)){
247  case GO_POSX:
248  // send the packet to the next hop in the +x direction
249  send_packet(newbuf, next_posx, t1);
250  // If a dynamic algorithm is being used, increment queue.
251  if (atoi(h_ra) > 3) {
252  inc_outputQueue(sem_posx, q_posx);
253  printf("incrementing posx\n");}
254  break;
255  case GO_NEGX:
256  send_packet(newbuf, next_negx, t1);
257  if (atoi(h_ra) > 3) {
258  inc_outputQueue(sem_negx, q_negx);
259  printf("incrementing negx\n");}
260  break;
261  case GO_POSY:
262  send_packet(newbuf, next_posy, t1);
263  if (atoi(h_ra) > 3) {
264  inc_outputQueue(sem_posy, q_posy);
265  printf("incrementing posy\n");}
266  break;
267  case GO_NEGY:
268  send_packet(newbuf, next_negy, t1);
269  if (atoi(h_ra) > 3) {
270  inc_outputQueue(sem_negy, q_negy);
271  printf("incrementing negy\n");}
272  break;
273  case GO_POSZ:
274  send_packet(newbuf, next_posz, t1);
275  if (atoi(h_ra) > 3) {
276  inc_outputQueue(sem_posz, q_posz);
277  printf("incrementing posz\n");}
278  break;
279  case GO_NEGZ:
280  send_packet(newbuf, next_negz, t1);
281  if (atoi(h_ra) > 3) {
282  inc_outputQueue(sem_negz, q_negz);
283  printf("incrementing negz\n");}
284  break;
285  default:
286  printf("Chose the default area...\n");
287  remHost=gethostbyname("127.0.0.1");
288  remSocket=socket(AF_INET, SOCK_STREAM, 0);
289  next_addr.sin_family = AF_INET;
290  next_addr.sin_port = htons(LOPORT);
291  next_addr.sin_addr = *((struct in_addr *)remHost->h_addr);
292  gettimeofday(&t, NULL);
293  t2 = (unsigned long long)t.tv_sec * 1000000 + (unsigned long long)t.tv_usec;
294  if ((t2-t1) >= 10000)
295  sprintf(d_time, "%llu", t2-t1);
296  else
297  sprintf(d_time,"0%llu", t2-t1);
298  sprintf(h_queue_t, "%d", atoi(d_time) + atoi(h_queue));
299  strncpy (&header[16],&h_queue_t[0], 6);
300  if (connect(remSocket,(struct sockaddr *)&next_addr, sizeof(next_addr)) < 0)
301  perror("connecting to localhost");
302  if (send(remSocket,header,HEADERSIZE-1,0) < 0)

```

```

303 perror("sending to localhost");
304 close(remSocket);
305 break;
306
307
308 }
309
310 }
311 // Child is now done, close socket and exit.
312 close(new_fd);
313 exit(0);
314 }
315 // Parent doesn't need this.
316 close(new_fd);
317 }
318 return 0;
319 }
320
321 /* ===== */
322 /* ===== funct: makeDecision() ===== */
323 /* ===== */
324 // This function given the header, routing algorithm
325 // and the destination, calculates the next hop.
326 // It also takes into account any output queues if
327 // a dynamic algorithm is used.
328 int makeDecision(char *hdr, char *ra, char *dest)
329 {
330 int x,y,z;
331 int db_addr;
332 int dir_x, dir_y, dir_z, algorithm;
333 int ret_val = -1;
334 char h_dir[4], h_v[4], h_int[4], h_x[2], h_y[2], h_z[2];
335 int sign_x, sign_y, sign_z;
336
337 algorithm = atoi(ra);
338 db_addr = (int)atoi(dest);
339 z = (int) floor(db_addr / (k*k));
340 y = (int) floor(db_addr / k) % k;
341 x = (int) db_addr % k;
342
343 switch (algorithm) {
344
345 case 0: // Not Used.
346 break;
347 case 1: // Dimension Ordered Routing
348 dir_x = dimord(myx,x);
349 dir_y = dimord(myy,y);
350 dir_z = dimord(myz,z);
351 if ((dir_x=dimord(myx,x))!=0)
352 if (dir_x == -1)
353 ret_val = GO_NEGX;
354 else
355 ret_val = GO_POSX;
356 else if ((dir_y=dimord(myy,y))!=0)
357 if (dir_y == -1)
358 ret_val = GO_NEGY;
359 else
360 ret_val = GO_POSY;
361 else if ((dir_z=dimord(myz,z))!=0)
362 if (dir_z == -1)
363 ret_val = GO_NEGZ;
364 else
365 ret_val = GO_POSZ;
366
367 break;

```

```

368 case 2: // Direction Ordered Routing
369 dir_x = dimord(myx,x);
370 dir_y = dimord(myy,y);
371 dir_z = dimord(myz,z);
372 if (dir_x == 1)
373 ret_val = GO_POSX;
374 else if (dir_y == 1)
375 ret_val = GO_POSY;
376 else if (dir_z == 1)
377 ret_val = GO_POSZ;
378 else if (dir_x == -1)
379 ret_val = GO_NEGX;
380 else if (dir_y == -1)
381 ret_val = GO_NEGY;
382 else if (dir_z == -1)
383 ret_val = GO_NEGZ;
384
385 break;
386 case 3: // Minimal Oblivious
387 dir_x = dimord(myx,x);
388 dir_y = dimord(myy,y);
389 dir_z = dimord(myz,z);
390 ret_val = min_choose(dir_x, dir_y, dir_z);
391 break;
392 case 4: // Minimal Adaptive
393 dir_x = dimord(myx,x);
394 dir_y = dimord(myy,y);
395 dir_z = dimord(myz,z);
396 ret_val = adaptive_choose(dir_x, dir_y, dir_z);
397 break;
398 case 5: // CQR
399 dir_x = dimord(myx,x);
400 dir_y = dimord(myy,y);
401 dir_z = dimord(myz,z);
402 strncpy(&h_v[0],&hdr[22],3); h_v[3]='\0';
403 strncpy(&h_dir[0],&hdr[25],3); h_dir[3]='\0';
404 if (atoi(h_v) == 0) { // First hop, find v.
405 make_vector(h_v, h_dir, dir_x, dir_y, dir_z);
406 strncpy(&hdr[25],&h_dir[0],3);
407 strncpy(&hdr[22],&h_v[0],3);
408 }
409 h_x[0] = h_v[0]; h_x[1] = '\0';
410 h_y[0] = h_v[1]; h_y[1] = '\0';
411 h_z[0] = h_v[2]; h_z[1] = '\0';
412 if (h_dir[0] == '0') sign_x = 1;
413 else sign_x = -1;
414 if (h_dir[1] == '0') sign_y = 1;
415 else sign_y = -1;
416 if (h_dir[2] == '0') sign_z = 1;
417 else sign_z = -1;
418 ret_val = adaptive_choose(sign_x*atoi(h_x), sign_y*atoi(h_y), sign_z*atoi(h_z));
419 switch (ret_val) {
420 case GO_POSX: h_v[0] = h_v[0] - 1; break;
421 case GO_NEGX: h_v[0] = h_v[0] + 1; break;
422 case GO_POSY: h_v[1] = h_v[1] - 1; break;
423 case GO_NEGY: h_v[1] = h_v[1] + 1; break;
424 case GO_POSZ: h_v[2] = h_v[2] - 1; break;
425 case GO_NEGZ: h_v[2] = h_v[2] + 1; break;
426 }
427 strncpy(&hdr[22],&h_v[0],3);
428 break;
429 case 6: // CQR - Periphery Avoidance
430 dir_x = dimord(myx,x);
431 dir_y = dimord(myy,y);
432 dir_z = dimord(myz,z);

```

```

433 strncpy(&h_v[0],&hdr[22],3); h_v[3]='\0';
434 strncpy(&h_dir[0],&hdr[25],3); h_dir[3]='\0';
435 if (atoi(h_v) == 0) { // First hop, find v.
436 make_vector(h_v, h_dir, dir_x, dir_y, dir_z);
437 strncpy(&hdr[25],&h_dir[0],3);
438 strncpy(&hdr[22],&h_v[0],3);
439 }
440 h_x[0] = h_v[0]; h_x[1] = '\0';
441 h_y[0] = h_v[1]; h_y[1] = '\0';
442 h_z[0] = h_v[2]; h_z[1] = '\0';
443 if (h_dir[0] == '0') sign_x = 1;
444 else sign_x = -1;
445 if (h_dir[1] == '0') sign_y = 1;
446 else sign_y = -1;
447 if (h_dir[2] == '0') sign_z = 1;
448 else sign_z = -1;
449 ret_val = adaptive_periphery_choose(sign_x*atoi(h_x), sign_y*atoi(h_y), sign_z*atoi(h_z));
450 switch (ret_val) {
451 case GO_POSX: h_v[0] = h_v[0] - 1; break;
452 case GO_NEGX: h_v[0] = h_v[0] + 1; break;
453 case GO_POSY: h_v[1] = h_v[1] - 1; break;
454 case GO_NEGY: h_v[1] = h_v[1] + 1; break;
455 case GO_POSZ: h_v[2] = h_v[2] - 1; break;
456 case GO_NEGZ: h_v[2] = h_v[2] + 1; break;
457 }
458 strncpy(&hdr[22],&h_v[0],3);
459 break;
460 case 7: // VGD-CQR
461 break;
462 default:
463 break;
464 }
465 return ret_val;
466 }
467
468 /* ===== */
469 /* ===== funct: initNeighborhood() ===== */
470 /* ===== */
471 // This function initializes the IP addresses of
472 // all the USB interfaces on this node, as well as
473 // calculating the IP addresses of each neighbor.
474 void initNeighborhood(void)
475 {
476
477 char c_octet_pos[3], c_octet_neg[3];
478 char d_octet_pos[3], d_octet_neg[3];
479 int x,y,z;
480 double db_addr, res;
481 int dir_mod_x = 1;
482 int dir_mod_y = 5;
483 int dir_mod_z = 9;
484
485 db_addr = (double)atoi(addr);
486 z = (int) floor(db_addr / (k*k));
487 y = (int) floor(db_addr / k) % k;
488 x = (int) db_addr % k;
489 myx = x;
490 myy = y;
491 myz = z;
492
493 //Xs.
494 sprintf(c_octet_pos,"%d",((x*16)+y));
495 sprintf(d_octet_pos,"%d",((z*16)+dir_mod_x)+1);
496 sprintf(d_octet_neg,"%d",((z*16)+dir_mod_x));
497 if (!x)

```

```

498 printf(c_octet_neg, "%d", ((k-1)*16)+y));
499 else
500 printf(c_octet_neg, "%d", ((x-1)*16)+y));
501 printf(next_posx, "10.0.%s.%s", c_octet_pos, d_octet_pos);
502 printf(next_negx, "10.0.%s.%s", c_octet_neg, d_octet_neg);
503
504 //Ys.
505 printf(c_octet_pos, "%d", ((x*16)+y));
506 printf(d_octet_pos, "%d", ((z*16)+dir_mod_y)+1);
507 printf(d_octet_neg, "%d", ((z*16)+dir_mod_y));
508 if (!y)
509 printf(c_octet_neg, "%d", ((x*16)+(k-1)));
510 else
511 printf(c_octet_neg, "%d", ((x*16)+(y-1)));
512 printf(next_posy, "10.0.%s.%s", c_octet_pos, d_octet_pos);
513 printf(next_negy, "10.0.%s.%s", c_octet_neg, d_octet_neg);
514
515 //Zs.
516 printf(c_octet_pos, "%d", ((x*16)+y));
517 printf(c_octet_neg, "%d", ((x*16)+y));
518 printf(d_octet_pos, "%d", ((z*16)+dir_mod_z)+1);
519 if (!z)
520 printf(d_octet_neg, "%d", ((k-1)*16)+dir_mod_z);
521 else
522 printf(d_octet_neg, "%d", ((z-1)*16)+dir_mod_z);
523 printf(next_posz, "10.0.%s.%s", c_octet_pos, d_octet_pos);
524 printf(next_negz, "10.0.%s.%s", c_octet_neg, d_octet_neg);
525 }
526
527 /* ===== */
528 /* ===== funct: send_packet() ===== */
529 /* ===== */
530 // This function sends a packet to the correct destination.
531 // Much of this code was originally repeated throughout this
532 // file, all converged here.
533 void send_packet(char* packet, char* host, unsigned long long timestamp)
534 {
535     int remSocket;
536     struct timeval t;
537     struct hostent *remHost;
538     struct sockaddr_in next_addr;
539     unsigned long long t2;
540     char d_time[7], h_queue_t[7], h_queue[7];
541
542     // copy the previous queue value from the packet
543     strncpy(&h_queue[0], &packet[16], 6); h_queue[6] = '\0';
544     remHost = gethostbyname(host);
545     // setup the outgoing socket
546     remSocket = socket(AF_INET, SOCK_STREAM, 0);
547     next_addr.sin_family = AF_INET;
548     next_addr.sin_port = htons(ROUTEPORT);
549     next_addr.sin_addr = *((struct in_addr *)remHost->h_addr);
550     // connect to the next hop, TCP handshake occurs
551     if (connect(remSocket, (struct sockaddr *)&next_addr,
552         sizeof(next_addr)) < 0)
553         myperror("connecting to next host");
554     // grab time & calculate difference, add to previous queue times
555     gettimeofday(&t, NULL);
556     t2 = ((unsigned long long)t.tv_sec * 1000000 +
557         (unsigned long long)t.tv_usec) - timestamp;
558     if (t2 >= 10000)
559         printf(d_time, "%llu", t2);
560     else
561         printf(d_time, "0%llu", t2);
562     printf(h_queue_t, "%d", atoi(d_time) + atoi(h_queue));

```

```

563 strncpy (&packet[16],&h_queue_t[0], 6);
564 // send data, close socket.
565 if(send(remSocket,packet,MAXDATASIZE-1,0) < 0)
566 myperror("sending to next host");
567 close(remSocket);
568 }
569
570 /* ===== */
571 /* ===== funct: child_handler() ===== */
572 /* ===== */
573 // This function reaps all dead child processes.
574 void child_handler(int s)
575 {
576 while(waitpid(-1, NULL, WNOHANG) > 0);
577 }
578
579 /* ===== */
580 /* ===== funct: myperror() ===== */
581 /* ===== */
582 // This function was created because these two
583 // lines were used frequently.
584 void myperror(char *x)
585 {
586 perror(x);
587 exit(1);
588 }
589
590 /* ===== */
591 /* ===== funct: dimord() ===== */
592 /* ===== */
593 // This function returns the difference between
594 // si and di in a particular dimension.
595 int dimord(int s, int d)
596 {
597 int temp;
598 int ret_val;
599
600 temp = (d-s) % k;
601 if (temp < -1)
602 temp = temp + k;
603 else if (temp > 1)
604 temp = temp - k;
605
606 if (temp < 0)
607 ret_val = -1;
608 else if (temp > 0)
609 ret_val = 1;
610 if (s == d)
611 ret_val = 0;
612
613 return ret_val;
614 }
615
616 /* ===== */
617 /* ===== funct: get_outputQueue() ===== */
618 /* ===== */
619 // This function returns the value of the output
620 // queue after obtaining the semaphore.
621 int get_outputQueue(int sem_id, struct queue *q)
622 {
623 int ret_val = -1;
624 // Receive the semaphore
625 sem_lock(sem_id);
626 // Enter the critical section
627 ret_val = q[0].val;

```

```

628 printf("queue val: %d\n", ret_val);
629 // Return the semaphore
630 sem_unlock(sem_id);
631
632 return ret_val;
633 }
634
635 /* ===== */
636 /* ===== funct: inc_outputQueue() ===== */
637 /* ===== */
638 // This function increments an output queue once
639 // it successfully receives the semaphore for that
640 // queue. It returns -1 if the queue is full.
641 int inc_outputQueue(int sem_id, struct queue *q)
642 {
643 int ret_val = -1;
644 // Receive the semaphore
645 sem_lock(sem_id);
646 // Enter the critical section
647 if (q[0].val < MAX_QUEUE)
648 {
649 q[0].val = q[0].val + 1;
650 ret_val = 1;
651 }
652 else printf("queue exceeded MAX\n");
653 // Return the semaphore
654 sem_unlock(sem_id);
655
656 return ret_val;
657 }
658
659 /* ===== */
660 /* ===== funct: adjust_queues() ===== */
661 /* ===== */
662 // This function adjusts the output queues by decrementing
663 // the queues based on their last updated time value. If
664 // it is greater than ADJUSTMENT, it is decreased that
665 // number of times. It is assumed that packets depart
666 // from the queues at a rate of one per ADJUSTMENT microseconds.
667 void adjust_queues(int sem_id, struct queue *q)
668 {
669 struct timeval t;
670 unsigned long long time;
671
672 gettimeofday(&t, NULL);
673 // Receive the semaphore
674 sem_lock(sem_id);
675 // Enter the critical section
676 time = (unsigned long long)t.tv_sec * 1000000 + (unsigned long long)t.tv_usec;
677 // remove as many packets from the output queues as we're expecting packets to leave.
678 while ((q[0].lastUpdated + ((unsigned long long)ADJUSTMENT)) <= time)
679 {
680 if (q[0].val > 0) q[0].val = q[0].val - 1;
681 q[0].lastUpdated = q[0].lastUpdated + ((unsigned long long)ADJUSTMENT);
682 }
683 // Return semaphore
684 sem_unlock(sem_id);
685
686 }
687
688 /* ===== */
689 /* ===== funct: sem_lock() ===== */
690 /* ===== */
691 // This function locks a given semaphore. It blocks
692 // until successfully obtained.

```

```

693 void sem_lock(int sem_set_id)
694 {
695     struct sembuf sem_op;
696     sem_op.sem_num = 0;
697     sem_op.sem_op = -1;
698     sem_op.sem_flg = 0;
699     semop(sem_set_id, &sem_op, 1);
700 }
701
702 /* ===== */
703 /* ===== funct: sem_unlock() ===== */
704 /* ===== */
705 // This function returns a semaphore for use by
706 // another process.
707 void sem_unlock(int sem_set_id)
708 {
709     struct sembuf sem_op;
710     sem_op.sem_num = 0;
711     sem_op.sem_op = 1;
712     sem_op.sem_flg = 0;
713     semop(sem_set_id, &sem_op, 1);
714 }
715
716 /* ===== */
717 /* ===== funct: min_choose() ===== */
718 /* ===== */
719 // This function randomly chooses an order based on
720 // all permutations of k and then returns which
721 // direction to randomly route within given the
722 // possible values passed. Eg. if a packet can
723 // minimally route +x or -y, this function picks
724 // between the two choices.
725 int min_choose(int x, int y, int z)
726 {
727     int t1;
728     int ret_val = 0;
729
730     if ((t1=rand())<0.166)
731     {
732         if (x != 0)
733             if (x > 0) ret_val = GO_POSX;
734             else ret_val = GO_NEGX;
735         else if (y != 0)
736             if (y > 0) ret_val = GO_POSY;
737             else ret_val = GO_NEGY;
738         else if (z != 0)
739             if (z > 0) ret_val = GO_POSZ;
740             else ret_val = GO_NEGZ;
741     }
742     else if (t1 < 0.333)
743     {
744
745         if (y != 0)
746             if (y > 0) ret_val = GO_POSY;
747             else ret_val = GO_NEGY;
748         else if (z != 0)
749             if (z > 0) ret_val = GO_POSZ;
750             else ret_val = GO_NEGZ;
751         else if (x != 0)
752             if (x > 0) ret_val = GO_POSX;
753             else ret_val = GO_NEGX;
754     }
755     else if (t1 < 0.5 )
756     {
757         if (z != 0)

```



```

758 if (z > 0) ret_val = GO_POSZ;
759 else ret_val = GO_NEGZ;
760 else if (x != 0)
761 if (x > 0) ret_val = GO_POSX;
762 else ret_val = GO_NEGX;
763 else if (y != 0)
764 if (y > 0) ret_val = GO_POSY;
765 else ret_val = GO_NEGY;
766 }
767 else if (t1 < 0.666)
768 {
769 if (x != 0)
770 if (x > 0) ret_val = GO_POSX;
771 else ret_val = GO_NEGX;
772 else if (z != 0)
773 if (z > 0) ret_val = GO_POSZ;
774 else ret_val = GO_NEGZ;
775 else if (y != 0)
776 if (y > 0) ret_val = GO_POSY;
777 else ret_val = GO_NEGY;
778 }
779 else if (t1 < 0.866)
780 {
781
782 if (y != 0)
783 if (y > 0) ret_val = GO_POSY;
784 else ret_val = GO_NEGY;
785 else if (x != 0)
786 if (x > 0) ret_val = GO_POSX;
787 else ret_val = GO_NEGX;
788 else if (z != 0)
789 if (z > 0) ret_val = GO_POSZ;
790 else ret_val = GO_NEGZ;
791
792 }
793 else
794 {
795 if (z != 0)
796 if (z > 0) ret_val = GO_POSZ;
797 else ret_val = GO_NEGZ;
798 else if (y != 0)
799 if (y > 0) ret_val = GO_POSY;
800 else ret_val = GO_NEGY;
801 else if (x != 0)
802 if (x > 0) ret_val = GO_POSX;
803 else ret_val = GO_NEGX;
804
805 }
806 return ret_val;
807
808 }
809
810
811
812 /* ===== */
813 /* ===== funct: adaptive_choose() ===== */
814 /* ===== */
815 // This is the function which calculates the next
816 // adaptive decision based on the possible directions
817 // and the output queues.
818 int adaptive_choose(int x, int y, int z)
819 {
820 int t, ret_val;
821 int x_queue = LARGENU;
822 int y_queue = LARGENU;

```

```

823 int z_queue = LARGENU;
824
825 ret_val = -1;
826 if (x != 0)
827 if (x > 0)
828 x_queue = get_outputQueue(sem_posx, q_posx);
829 else
830 x_queue = get_outputQueue(sem_negx, q_negx);
831 if (y != 0)
832 if (y > 0)
833 y_queue = get_outputQueue(sem_posy, q_posy);
834 else
835 y_queue = get_outputQueue(sem_negy, q_negy);
836 if (z != 0)
837 if (z > 0)
838 z_queue = get_outputQueue(sem_posz, q_posz);
839 else
840 z_queue = get_outputQueue(sem_negz, q_negz);
841
842
843 t = min(min(x_queue, y_queue), z_queue);
844
845 if (t == x_queue)
846 if (x > 0) ret_val = GO_POSX;
847 else ret_val = GO_NEGX;
848 else if (t == y_queue)
849 if (y > 0) ret_val = GO_POSY;
850 else ret_val = GO_NEGY;
851 else if (t == z_queue)
852 if (z > 0) ret_val = GO_POSZ;
853 else ret_val = GO_NEGZ;
854
855 return ret_val;
856 }
857
858 /* ===== */
859 /* ===== funct: adaptive_periphery_choose() ===== */
860 /* ===== */
861 // This function does the same as adaptive_choose, but
862 // includes the periphery calculation in the determination
863 // of the next hop.
864 int adaptive_periphery_choose(int x, int y, int z)
865 {
866
867 int ret_val;
868 double t, delta_total;
869 double x_queue = (double)LARGENU;
870 double y_queue = (double)LARGENU;
871 double z_queue = (double)LARGENU;
872
873
874 delta_total = (double)(abs(x) + abs(y) + abs(z));
875 ret_val = -1;
876 if (x != 0)
877 if (x > 0)
878 x_queue = (double)(1 + get_outputQueue(sem_posx, q_posx));
879 else
880 x_queue = (double)(1 + get_outputQueue(sem_negx, q_negx));
881 if (y != 0)
882 if (y > 0)
883 y_queue = (double)(1 + get_outputQueue(sem_posy, q_posy));
884 else
885 y_queue = (double)(1 + get_outputQueue(sem_negy, q_negy));
886 if (z != 0)
887 if (z > 0)

```

```

888 z_queue = (double)(1 + get_outputQueue(sem_posz, q_posz));
889 else
890 z_queue = (double)(1 + get_outputQueue(sem_negz, q_negz));
891
892 x_queue = x_queue * (1.0 - ((double)abs(x)/delta_total));
893 y_queue = y_queue * (1.0 - ((double)abs(y)/delta_total));
894 z_queue = z_queue * (1.0 - ((double)abs(z)/delta_total));
895
896 t = min_d(min_d(x_queue, y_queue), z_queue);
897
898 if (t == x_queue)
899 if (x > 0) ret_val = GO_POSX;
900 else ret_val = GO_NEGX;
901 else if (t == y_queue)
902 if (y > 0) ret_val = GO_POSY;
903 else ret_val = GO_NEGY;
904 else if (t == z_queue)
905 if (z > 0) ret_val = GO_POSZ;
906 else ret_val = GO_NEGZ;
907
908 return ret_val;
909
910 }
911
912
913 /* ===== */
914 /* ===== funct: make_vector() ===== */
915 /* ===== */
916 // This function creates the v and dir vectors for CQR and ECQR.
917 void make_vector(char v[], char dir[], int x_val, int y_val, int z_val){
918
919 int t_q, x_val_t, y_val_t, z_val_t;
920 double temp_min_val, min_val;
921 char v_1[4];
922 int nu_queues = 3;
923
924 // for +++
925 x_val_t = x_val;
926 y_val_t = y_val;
927 z_val_t = z_val;
928 t_q = 0;
929 if (x_val_t > 0)
930 t_q = get_outputQueue(sem_posx, q_posx);
931 else if (x_val_t < 0)
932 t_q = get_outputQueue(sem_negx, q_negx);
933 else
934 nu_queues = nu_queues - 1;
935 if (y_val_t > 0)
936 t_q = t_q + get_outputQueue(sem_posy, q_posy);
937 else if (y_val_t < 0)
938 t_q = t_q + get_outputQueue(sem_negy, q_negy);
939 else
940 nu_queues = nu_queues - 1;
941 if (z_val_t > 0)
942 t_q = t_q + get_outputQueue(sem_posz, q_posz);
943 else if (z_val_t < 0)
944 t_q = t_q + get_outputQueue(sem_negz, q_negz);
945 else
946 nu_queues = nu_queues - 1;
947 min_val=((double)abs(x_val_t)+(double)abs(y_val_t)+(double)abs(z_val_t))*((double)t_q/(double)nu_queues);
948 sprintf(v_1, "%d%d%d", abs(x_val_t), abs(y_val_t), abs(z_val_t));
949 sprintf(dir, "000");
950
951 // for +-
952 x_val_t = x_val;

```

```

953 y_val_t = y_val;
954 if (z_val > 0) z_val_t = z_val-k;
955 else if (z_val < 0) z_val_t = z_val+k;
956 else z_val_t = 0;
957 nu_queues = 3;
958 t_q = 0;
959 if (x_val_t > 0)
960 t_q = get_outputQueue(sem_posx, q_posx);
961 else if (x_val_t < 0)
962 t_q = get_outputQueue(sem_negx, q_negx);
963 else
964 nu_queues = nu_queues - 1;
965 if (y_val_t > 0)
966 t_q = t_q + get_outputQueue(sem_posy, q_posy);
967 else if (y_val_t < 0)
968 t_q = t_q + get_outputQueue(sem_negy, q_negy);
969 else
970 nu_queues = nu_queues - 1;
971 if (z_val_t > 0)
972 t_q = t_q + get_outputQueue(sem_posz, q_posz);
973 else if (z_val_t < 0)
974 t_q = t_q + get_outputQueue(sem_negz, q_negz);
975 else
976 nu_queues = nu_queues - 1;
977 temp_min_val = ((double)abs(x_val_t)+(double)abs(y_val_t)+(double)abs(z_val_t))*((double)t_q/(double)nu_queues);
978 if (temp_min_val < min_val){
979 min_val = temp_min_val;
980 sprintf(v_1, "%d%d%d", abs(x_val_t), abs(y_val_t), abs(z_val_t));
981 if (x_val_t >= 0) dir[0] = '0'; else dir[0] = '1';
982 if (y_val_t >= 0) dir[1] = '0'; else dir[1] = '1';
983 if (z_val_t >= 0) dir[2] = '0'; else dir[2] = '1';
984 }
985
986 // +-+
987 x_val_t = x_val;
988 if (y_val > 0) y_val_t = y_val-k;
989 else if (y_val < 0) y_val_t = y_val+k;
990 else y_val_t = 0;
991 z_val_t = z_val;
992 nu_queues = 3;
993 t_q = 0;
994 if (x_val_t > 0)
995 t_q = get_outputQueue(sem_posx, q_posx);
996 else if (x_val_t < 0)
997 t_q = get_outputQueue(sem_negx, q_negx);
998 else
999 nu_queues = nu_queues - 1;
1000 if (y_val_t > 0)
1001 t_q = t_q + get_outputQueue(sem_posy, q_posy);
1002 else if (y_val_t < 0)
1003 t_q = t_q + get_outputQueue(sem_negy, q_negy);
1004 else
1005 nu_queues = nu_queues - 1;
1006 if (z_val_t > 0)
1007 t_q = t_q + get_outputQueue(sem_posz, q_posz);
1008 else if (z_val_t < 0)
1009 t_q = t_q + get_outputQueue(sem_negz, q_negz);
1010 else
1011 nu_queues = nu_queues - 1;
1012 temp_min_val = ((double)abs(x_val_t)+(double)abs(y_val_t)+(double)abs(z_val_t))*((double)t_q/(double)nu_queues);
1013 if (temp_min_val < min_val){
1014 min_val = temp_min_val;
1015 sprintf(v_1, "%d%d%d", abs(x_val_t), abs(y_val_t), abs(z_val_t));
1016 if (x_val_t >= 0) dir[0] = '0'; else dir[0] = '1';
1017 if (y_val_t >= 0) dir[1] = '0'; else dir[1] = '1';

```

```

1018 if (z_val_t >= 0) dir[2] = '0'; else dir[2] = '1';
1019 }
1020
1021
1022 // +--
1023 x_val_t = x_val;
1024 if (y_val > 0) y_val_t = y_val-k;
1025 else if (y_val < 0) y_val_t = y_val+k;
1026 else y_val_t = 0;
1027 if (z_val > 0) z_val_t = z_val-k;
1028 else if (z_val < 0) z_val_t = z_val+k;
1029 else z_val_t = 0;
1030 nu_queues = 3;
1031 t_q = 0;
1032 if (x_val_t > 0)
1033 t_q = get_outputQueue(sem_posx, q_posx);
1034 else if (x_val_t < 0)
1035 t_q = get_outputQueue(sem_negx, q_negx);
1036 else
1037 nu_queues = nu_queues - 1;
1038 if (y_val_t > 0)
1039 t_q = t_q + get_outputQueue(sem_posy, q_posy);
1040 else if (y_val_t < 0)
1041 t_q = t_q + get_outputQueue(sem_negy, q_negy);
1042 else
1043 nu_queues = nu_queues - 1;
1044 if (z_val_t > 0)
1045 t_q = t_q + get_outputQueue(sem_posz, q_posz);
1046 else if (z_val_t < 0)
1047 t_q = t_q + get_outputQueue(sem_negz, q_negz);
1048 else
1049 nu_queues = nu_queues - 1;
1050 temp_min_val = ((double)abs(x_val_t)+(double)abs(y_val_t)+(double)abs(z_val_t))*((double)t_q/(double)nu_queues);
1051 if (temp_min_val < min_val){
1052 min_val = temp_min_val;
1053 sprintf(v_1, "%d%d%d", abs(x_val_t), abs(y_val_t), abs(z_val_t));
1054 if (x_val_t >= 0) dir[0] = '0'; else dir[0] = '1';
1055 if (y_val_t >= 0) dir[1] = '0'; else dir[1] = '1';
1056 if (z_val_t >= 0) dir[2] = '0'; else dir[2] = '1';
1057 }
1058
1059
1060 // --+
1061 if (x_val > 0) x_val_t = x_val-k;
1062 else if (x_val < 0) x_val_t = x_val+k;
1063 else x_val_t = 0;
1064 y_val_t = y_val;
1065 z_val_t = z_val;
1066 nu_queues = 3;
1067 t_q = 0;
1068 if (x_val_t > 0)
1069 t_q = get_outputQueue(sem_posx, q_posx);
1070 else if (x_val_t < 0)
1071 t_q = get_outputQueue(sem_negx, q_negx);
1072 else
1073 nu_queues = nu_queues - 1;
1074 if (y_val_t > 0)
1075 t_q = t_q + get_outputQueue(sem_posy, q_posy);
1076 else if (y_val_t < 0)
1077 t_q = t_q + get_outputQueue(sem_negy, q_negy);
1078 else
1079 nu_queues = nu_queues - 1;
1080 if (z_val_t > 0)
1081 t_q = t_q + get_outputQueue(sem_posz, q_posz);
1082 else if (z_val_t < 0)

```

```

1083 t_q = t_q + get_outputQueue(sem_negz, q_negz);
1084 else
1085 nu_queues = nu_queues - 1;
1086 temp_min_val = ((double)abs(x_val_t)+(double)abs(y_val_t)+(double)abs(z_val_t))*((double)t_q/(double)nu_queues);
1087 if (temp_min_val < min_val){
1088 min_val = temp_min_val;
1089 sprintf(v_1, "%d%d%d", abs(x_val_t), abs(y_val_t), abs(z_val_t));
1090 if (x_val_t >= 0) dir[0] = '0'; else dir[0] = '1';
1091 if (y_val_t >= 0) dir[1] = '0'; else dir[1] = '1';
1092 if (z_val_t >= 0) dir[2] = '0'; else dir[2] = '1';
1093 }
1094
1095
1096 // --+
1097 if (x_val > 0) x_val_t = x_val-k;
1098 else if (x_val < 0) x_val_t = x_val+k;
1099 else x_val_t = 0;
1100 y_val_t = y_val;
1101 if (z_val > 0) z_val_t = z_val-k;
1102 else if (z_val < 0) z_val_t = z_val+k;
1103 else z_val_t = 0;
1104 nu_queues = 3;
1105 t_q = 0;
1106 if (x_val_t > 0)
1107 t_q = get_outputQueue(sem_posx, q_posx);
1108 else if (x_val_t < 0)
1109 t_q = get_outputQueue(sem_negx, q_negx);
1110 else
1111 nu_queues = nu_queues - 1;
1112 if (y_val_t > 0)
1113 t_q = t_q + get_outputQueue(sem_posy, q_posy);
1114 else if (y_val_t < 0)
1115 t_q = t_q + get_outputQueue(sem_negy, q_negy);
1116 else
1117 nu_queues = nu_queues - 1;
1118 if (z_val_t > 0)
1119 t_q = t_q + get_outputQueue(sem_posz, q_posz);
1120 else if (z_val_t < 0)
1121 t_q = t_q + get_outputQueue(sem_negz, q_negz);
1122 else
1123 nu_queues = nu_queues - 1;
1124 temp_min_val = ((double)abs(x_val_t)+(double)abs(y_val_t)+(double)abs(z_val_t))*((double)t_q/(double)nu_queues);
1125 if (temp_min_val < min_val){
1126 min_val = temp_min_val;
1127 sprintf(v_1, "%d%d%d", abs(x_val_t), abs(y_val_t), abs(z_val_t));
1128 if (x_val_t >= 0) dir[0] = '0'; else dir[0] = '1';
1129 if (y_val_t >= 0) dir[1] = '0'; else dir[1] = '1';
1130 if (z_val_t >= 0) dir[2] = '0'; else dir[2] = '1';
1131 }
1132
1133
1134 //--+
1135 if (x_val > 0) x_val_t = x_val-k;
1136 else if (x_val < 0) x_val_t = x_val+k;
1137 else x_val_t = 0;
1138 if (y_val > 0) y_val_t = y_val-k;
1139 else if (y_val < 0) y_val_t = y_val+k;
1140 else y_val_t = 0;
1141 z_val_t = z_val;
1142 nu_queues = 3;
1143 t_q = 0;
1144 if (x_val_t > 0)
1145 t_q = get_outputQueue(sem_posx, q_posx);
1146 else if (x_val_t < 0)
1147 t_q = get_outputQueue(sem_negx, q_negx);

```

```

1148 else
1149 nu_queues = nu_queues - 1;
1150 if (y_val_t > 0)
1151 t_q = t_q + get_outputQueue(sem_posy, q_posy);
1152 else if (y_val_t < 0)
1153 t_q = t_q + get_outputQueue(sem_negy, q_negy);
1154 else
1155 nu_queues = nu_queues - 1;
1156 if (z_val_t > 0)
1157 t_q = t_q + get_outputQueue(sem_posz, q_posz);
1158 else if (z_val_t < 0)
1159 t_q = t_q + get_outputQueue(sem_negz, q_negz);
1160 else
1161 nu_queues = nu_queues - 1;
1162 temp_min_val = ((double)abs(x_val_t)+(double)abs(y_val_t)+(double)abs(z_val_t))*((double)t_q/(double)nu_queues);
1163 if (temp_min_val < min_val){
1164 min_val = temp_min_val;
1165 sprintf(v_1, "%d%d%d", abs(x_val_t), abs(y_val_t), abs(z_val_t));
1166 if (x_val_t >= 0) dir[0] = '0'; else dir[0] = '1';
1167 if (y_val_t >= 0) dir[1] = '0'; else dir[1] = '1';
1168 if (z_val_t >= 0) dir[2] = '0'; else dir[2] = '1';
1169 }
1170
1171
1172 // ---
1173 if (x_val > 0) x_val_t = x_val-k;
1174 else if (x_val < 0) x_val_t = x_val+k;
1175 else x_val_t = 0;
1176 if (y_val > 0) y_val_t = y_val-k;
1177 else if (y_val < 0) y_val_t = y_val+k;
1178 else y_val_t = 0;
1179 if (z_val > 0) z_val_t = z_val-k;
1180 else if (z_val < 0) z_val_t = z_val+k;
1181 else z_val_t = 0;
1182 nu_queues = 3;
1183 t_q = 0;
1184 if (x_val_t > 0)
1185 t_q = get_outputQueue(sem_posx, q_posx);
1186 else if (x_val_t < 0)
1187 t_q = get_outputQueue(sem_negx, q_negx);
1188 else
1189 nu_queues = nu_queues - 1;
1190 if (y_val_t > 0)
1191 t_q = t_q + get_outputQueue(sem_posy, q_posy);
1192 else if (y_val_t < 0)
1193 t_q = t_q + get_outputQueue(sem_negy, q_negy);
1194 else
1195 nu_queues = nu_queues - 1;
1196 if (z_val_t > 0)
1197 t_q = t_q + get_outputQueue(sem_posz, q_posz);
1198 else if (z_val_t < 0)
1199 t_q = t_q + get_outputQueue(sem_negz, q_negz);
1200 else
1201 nu_queues = nu_queues - 1;
1202 temp_min_val = ((double)abs(x_val_t)+(double)abs(y_val_t)+(double)abs(z_val_t))*((double)t_q/(double)nu_queues);
1203 if (temp_min_val < min_val){
1204 min_val = temp_min_val;
1205 sprintf(v_1, "%d%d%d", abs(x_val_t), abs(y_val_t), abs(z_val_t));
1206 if (x_val_t >= 0) dir[0] = '0'; else dir[0] = '1';
1207 if (y_val_t >= 0) dir[1] = '0'; else dir[1] = '1';
1208 if (z_val_t >= 0) dir[2] = '0'; else dir[2] = '1';
1209 }
1210 v[0] = v_1[0];
1211 v[1] = v_1[1];
1212 v[2] = v_1[2];

```

```

1213
1214 }
1215
1216 /* ===== */
1217 /* =====   funct: min_d()   ===== */
1218 /* ===== */
1219 // This function returns the minimum of the two values passed
1220 double min_d(double x, double y)
1221 {
1222 return (double)((x*(x<y)) + (y*(x>=y)));
1223 }
1224
1225 /* ===== */
1226 /* =====   funct: min()   ===== */
1227 /* ===== */
1228 // This function returns the minimum of the two values passed
1229 int min(int x, int y)
1230 {
1231 return ((x*(x<y)) + (y*(x>=y)));
1232 }
1233
1234
1235

```



## B.2 file: loser.c

```
1
2 /* ##### *
3  * Filename: loser.c      *
4  * Author:   Chris Lydick *
5  * Date:    Mar 1, 2008 *
6  * Usage:   ./loser [address] *
7  * Notes:   Portions borrowed from http://tinyurl.com/2w36o4 *
8  * *
9  * ##### */
10
11 #include "router.h"
12
13 /* ===== */
14 /* =====  FUNCT DECLARATION  ===== */
15 /* ===== */
16 void child_handler(int s);
17 void myperror(char *x);
18
19
20 /* ===== */
21 /* =====  Main Function  ===== */
22 /* ===== */
23 int main(int argc, char* argv[])
24 {
25     int sockfd, new_fd, numbytes, hop_nu, remSocket, yes;
26     struct sockaddr_in my_addr, their_addr, next_addr;
27     struct hostent *localhost;
28     socklen_t sin_size;
29     struct sigaction sa;
30     char newheader[HEADERSIZE+1], header[HEADERSIZE+1], h_src[3], h_dest[3], h_ra[2], h_hops[3], *addr, fullpath[80];
31     char time_ms[10], time_ms_t[17], h_queue[7];
32     char h_time[10];
33     FILE *fp;
34     struct timeval tv;
35     time_t curtime;
36     int p_time, f_time, delta;
37     unsigned long long time;
38
39     yes=1;
40     if (argc > 1)
41         addr = argv[1];
42     else myperror("usage: ./loser [node number]");
43
44     // Create server listen socket
45     if ((sockfd = socket(AF_INET, SOCK_STREAM, 0)) == -1)
46         myperror("socket");
47     // Set option on socket to reuse the address.
48     if (setsockopt(sockfd, SOL_SOCKET, SO_REUSEADDR, &yes, sizeof(int)) == -1)
49         perror("setsockopt");
50
51     // Set parameters for binding the socket to an address.
52     my_addr.sin_family = AF_INET;
53     my_addr.sin_port = htons(LOPORT);
54     my_addr.sin_addr.s_addr = inet_addr("127.0.0.1");
55     memset(my_addr.sin_zero, '\0', sizeof my_addr.sin_zero);
56
57     // Bind the socket and address.
58     if (bind(sockfd, (struct sockaddr *)&my_addr, sizeof my_addr) == -1)
59         myperror("bind");
60
61     // Set the socket to listen for incoming connections.
62     if (listen(sockfd, BACKLOG) == -1)
63         myperror("listen");
```

```

64
65 // This reaps all dead child processes.
66 sa.sa_handler = child_handler;
67 sigemptyset(&sa.sa_mask);
68 sa.sa_flags = SA_RESTART;
69 if (sigaction(SIGCHLD, &sa, NULL) == -1)
70 myperror("sigaction");
71
72 // Main accept() loop
73 while(1) {
74 sin_size = sizeof their_addr;
75 if ((new_fd = accept(sockfd, (struct sockaddr *)&their_addr, &sin_size)) == -1) {
76 perror("accept");
77 continue; }
78 printf("server: got connection from %s\n",inet_ntoa(their_addr.sin_addr));
79 // Fork off a child to do the work, only child will continue here.
80 if (!fork()) {
81 // Child process doesn't need the listener
82 close(sockfd);
83 if ((numbytes=recv(new_fd, header, HEADERSIZE-1, 0)) == -1)
84 perror("recv");
85 close(new_fd);
86 memcpy(newheader, header, HEADERSIZE); newheader[HEADERSIZE-1] = '\0';
87 strncpy(&h_src[0], &newheader[0],2); h_src[2] = '\0';
88 strncpy(&h_dest[0],&newheader[2],2); h_dest[2] = '\0';
89 strncpy(&h_ra[0],&newheader[4],1); h_ra[1] = '\0';
90 strncpy(&h_hops[0],&newheader[5],2); h_hops[2] = '\0';
91 strncpy(&h_time[0],&newheader[7],9); h_time[9] = '\0';
92 strncpy(&h_queue[0],&newheader[16],6); h_queue[6] = '\0';
93 gettimeofday(&tv, NULL);
94 time = (unsigned long long)tv.tv_sec * 1000000 + (unsigned long long)tv.tv_usec;
95 sprintf(time_ms_t,"%llu", time);
96 //sprintf(time_ms_t, "%ld.%ld", tv.tv_sec, tv.tv_usec);
97 strncpy(&time_ms[0],&time_ms_t[7],9); time_ms[9] = '\0';
98 //printf("time: %s\n", time_ms);
99 f_time = atoi(time_ms);
100 p_time = atoi(h_time);
101 delta = f_time - p_time;
102 //printf("f_time: %4.4f\n", f_time);
103 //printf("p_time: %4.4f\n", p_time);
104 //printf("time difference: %4.4f\n", f_time-p_time);
105 if (atoi(h_dest) == atoi(addr)){
106 sprintf(fullpath, "%s%s/SUCCESS_%s",FILE_PATH,addr,time_ms);
107 fp = fopen(fullpath, "w");
108 fprintf(fp,"from: %d\n", atoi(h_src));
109 fprintf(fp,"to: %d\n", atoi(h_dest));
110 fprintf(fp,"algorithm: %d\n", atoi(h_ra));
111 fprintf(fp,"hops: %d\n",atoi(h_hops));
112 fprintf(fp,"total_time: %d\n", delta);
113 fprintf(fp,"queue_time: %d\n", atoi(h_queue));
114 fprintf(fp,"trans_time: %d\n", delta-atoi(h_queue));
115 fprintf(fp,"avg_trans/hop: %d\n", (delta-atoi(h_queue))/atoi(h_hops));
116 fprintf(fp,"avg_queue/hop: %d\n", atoi(h_queue)/atoi(h_hops));
117 fprintf(fp,"start: %d\n", p_time);
118 fprintf(fp,"end: %d\n", f_time);
119 fprintf(fp,"header: %s\n", newheader);
120 fclose(fp);
121 }
122 else if (atoi(h_hops) >= MAXHOPS){
123 sprintf(fullpath, "%s%s/DROP_%s",FILE_PATH,addr,time_ms);
124 fp = fopen(fullpath, "w");
125 fprintf(fp,"from: %d\n", atoi(h_src));
126 fprintf(fp,"to: %d\n", atoi(h_dest));
127 fprintf(fp,"algorithm: %d\n", atoi(h_ra));
128 fprintf(fp,"hops: %d\n",atoi(h_hops));

```

```

129 fprintf(fp,"total_time: %d\n", delta);
130 fprintf(fp,"queue_time: %d\n", atoi(h_queue));
131 fprintf(fp,"trans_time: %d\n", delta-atoi(h_queue));
132 fprintf(fp,"start: %d\n", p_time);
133 fprintf(fp,"end: %d\n", f_time);
134 fprintf(fp,"header: %s\n", newheader);
135 fclose(fp);
136
137 }
138 else {
139 //perror("unrecognized packet");
140 sprintf(fullpath, "%s%s/UNKNOWN_%s",FILE_PATH,addr,time_ms);
141 fp = fopen(fullpath, "w");
142 fprintf(fp,"from: %d\n", atoi(h_src));
143 fprintf(fp,"to: %d\n", atoi(h_dest));
144 fprintf(fp,"algorithm: %d\n", atoi(h_ra));
145 fprintf(fp,"hops: %d\n",atoi(h_hops));
146 fprintf(fp,"total_time: %d\n", delta);
147 fprintf(fp,"queue_time: %d\n", atoi(h_queue));
148 fprintf(fp,"trans_time: %d\n", delta-atoi(h_queue));
149 fprintf(fp,"start: %d\n", p_time);
150 fprintf(fp,"end: %d\n", f_time);
151 fprintf(fp,"header: %s\n", newheader);
152 fclose(fp);
153 }
154 //printf("Closing this connection.");
155 exit(0);
156 }
157 close(new_fd); // parent doesn't need this
158 }
159
160 return 0;
161 }
162
163
164 void child_handler(int s)
165 {
166 while(waitpid(-1, NULL, WNOHANG) > 0);
167 }
168
169 void myperror(char *x)
170 {
171 perror(x);
172 exit(1);
173 }
174
175

```

## B.3 file: injector.c

```
1 /* ##### *
2 * Filename: injector.c *
3 * Author: Chris Lydick *
4 * Date: Mar 1, 2008 *
5 * Usage: ./injector [traffic demand file] [routing algorithm]*
6 * Notes: Portions borrowed from http://tinyurl.com/2w36o4 *
7 * *
8 * ##### */
9 #include "router.h"
10
11 /* ===== */
12 /* ===== Main Function ===== */
13 /* ===== */
14 int main(int argc, char *argv[])
15 {
16     int i, j, k, i_i, suffix, sockfd, numbytes;
17     char ra[2];
18     char buf[MAXDATASIZE], addr[14], line[81], time_ms_t[17], time_ms[10], src[3], dest[3];
19     char no_pkts[3];
20     struct hostent *he;
21     struct sockaddr_in their_addr;
22     time_t curtime;
23     struct timeval tv;
24     unsigned long long time;
25     FILE *fp;
26
27
28     if (argc != 3) {
29         fprintf(stderr, "usage: injector [traf demand file] [routing algorithm]\n");
30         exit(1);
31     }
32
33     sprintf(ra, argv[2], 1);
34
35     if ((fp=fopen(argv[1], "r")) == NULL)
36         fprintf(stderr, "file not found.\n");
37
38     // For each row in the text file, send packets with 0s in header.
39     // Fork off children to do this.
40     for (i=0; i<27; i++)
41     {
42         suffix = 200 - i;
43         fgets(line, 82, fp); line[80] = '\0';
44         if (!fork()){
45             i_i = i;
46             // Clear the header - but src is same for entire row.
47             strncpy(&buf[0], "0000000000000000000000000000000000", 40);
48             if (i_i < 10)
49                 sprintf(src, "%d", i_i);
50             else
51                 sprintf(src, "%d", i_i);
52             strncpy(&buf[0], &src[0], 2);
53             buf[4] = ra[0];
54             for (j=0; j<27; j++)
55             {
56                 if (i == j) j++;
57                 if (j == 27) break;
58                 if (j < 10)
59                     sprintf(dest, "%d", j);
60                 else
61                     sprintf(dest, "%d", j);
62                 strncpy(&buf[2], &dest[0], 2);
63                 strncpy(&no_pkts[0], &line[j*3], 2); no_pkts[2] = '\0';
```

```

64 for (k=0; k<atoi(no_pkts); k++)
65 {
66 sprintf(addr,"192.168.0.%d", suffix); addr[13] = '\0';
67 if ((he=gethostbyname(addr)) == NULL)
68 perror("gethostbyname");
69 if ((sockfd=socket(AF_INET, SOCK_STREAM, 0)) == -1)
70 perror("socket");
71 their_addr.sin_family = AF_INET;
72 their_addr.sin_port = htons(ROUTEPORT);
73 their_addr.sin_addr = *((struct in_addr *)he->h_addr);
74 memset(their_addr.sin_zero, '\0', sizeof their_addr.sin_zero);
75 if (connect(sockfd, (struct sockaddr *)&their_addr, sizeof their_addr) == -1) {
76     perror("connect"); exit(1);}
77 if (send(sockfd, buf, MAXDATASIZE-1, 0) == -1)
78     perror("send");
79 close(sockfd);
80 printf("sending a packet: %s -> %s\n", src, dest);
81
82 }
83
84 }
85     exit(0);
86 }
87     else {
88     }
89 }
90 return 0;
91 }

```

## B.4 file: router.h

```
1
2 #include <stdio.h>
3 #include <math.h>
4 #include <stdlib.h>
5 #include <unistd.h>
6 #include <errno.h>
7 #include <string.h>
8 #include <sys/types.h>
9 #include <sys/socket.h>
10 #include <netinet/in.h>
11 #include <arpa/inet.h>
12 #include <sys/wait.h>
13 #include <signal.h>
14 #include <netdb.h>
15 #include <sys/time.h>
16 #include <time.h>
17 #include <sys/shm.h>
18 #include <sys/sem.h>
19 #include <sys/ipc.h>
20
21 #define ROUTEPORT 3490
22 #define LOPORT 3491
23 #define MAXDATASIZE 800000
24 #define HEADERSIZE 50
25 #define MAXHOPS 7
26 #define BACKLOG 10
27 #define GO_POSX 1
28 #define GO_POSY 2
29 #define GO_POSZ 3
30 #define GO_NEGX 4
31 #define GO_NEGY 5
32 #define GO_NEGZ 6
33 #define FILE_PATH "/pckhome/pckhome/router/routerlogs/"
34 #define SID_POSX 45
35 #define SID_POSY 46
36 #define SID_POSZ 47
37 #define SID_NEGX 48
38 #define SID_NEGY 49
39 #define SID_NEGZ 50
40 #define MAX_QUEUE 200
41 #define ADJUSTMENT 2000000
42 #define LARGENU 1000000
43 struct queue
44 {
45     int val;
46     unsigned long long lastUpdated;
47 };
48 union semun
49 {
50     int val;
51     struct semid_ds *buf;
52     unsigned short *array;
53     struct seminfo *__buf;
54 };
```

# Appendix C

## Appendix C: Matlab Simulation Scripts

### C.1 User Interface Scripts

#### C.1.1 file: do\_load.m

```
1 function do_load();
2
3 pattern = input('\n\nEnter a traffic pattern to use... \n0:Flood Traffic (all src/dest pairs)\n1:Tornado Traffic\n
4 2:Transpose Traffic\n3:Nearest Neighbor\n4:Uniform Random Traffic\n5:Bit-Compliment Traffic\n: ');
5 option = input('\n\nEnter which routing algorithm to use... \n1:Dimension Ordered\n2:Direction Ordered\n3:Minimal
6 Oblivious\n4:Minimal Adaptive\n5:Chaos\n6:CQR\n7:Enhanced CQR - Periphery Avoidance\n: ');
7 k = input('\n\nEnter k...\n: ');
8 n = input('\n\nEnter n...\n: ');
9 runs = input('\n\nEnter the number of runs to simulate...\n: ');
10 output_text_file = 1;
11 if (option > 3)
12     failure = input('\n\nEnter the probability of nodal failures in decimal form...\n: ');
13     hs = input('\n\nEnter the percent of nodes which should be hotspots in decimal form...\n: ');
14 else
15     failure = 0.00;
16     hs = 0;
17 end
18
19 if (k^n > 50)
20     plot = input('\n\nThis seems to be a large set. Would you like to plot these? \n
21 **This operation could take a while to process!**\n0:No\n1:Yes\n: ');
22 else
23     plot = 1;
24 end
25
26 failed_nodes_queue = 0;
27 len = k^n;
28 ch_load = zeros(len, len);
29 traffic = zeros(len, len);
30 if(rem(k,2) == 1)
31     scale = 3*n*(k/4 - 1/(4*k)); % k even
32 else
33     scale = 6*k*n/4; % k odd
34 end
35
36 x = traffic/(scale);
37 traffic = x - diag(diag(x));
```

```

38
39 % All src/dest pairs
40 if(pattern == 0)
41     for i=1:k^n
42         for j=1:k^n
43             traffic(i,j) = runs;
44         end
45     end
46
47 % Tornado Traffic
48 elseif(pattern == 1)
49     for j=1:runs
50         for src=1:len
51             [x,y,z] = i2dim(src-1, k, n);
52
53             dx = rem(ceil(k/2)-1,k) + x;
54             if dx >= k
55                 dx = dx - k;
56             end
57             dy = y;
58             dz = z;
59             dest = dim2i(dx, dy, dz, k, n)+1;
60             traffic(src, dest) = traffic(src, dest) + 1;
61         end
62     end
63
64 % Transpose Traffic
65 elseif(pattern == 2)
66     for j=1:runs
67         for src=1:len
68             [x,y,z] = i2dim(src-1, k, n);
69             if(n > 2)
70                 dest = dim2i(z, x, y, k, n)+1;
71                 traffic(src, dest) = traffic(src,dest) + 1;
72                 dest = dim2i(z,y,x,k,n)+1;
73                 traffic(src, dest) = traffic(src,dest) + 1;
74                 dest = dim2i(x,z,y,k,n)+1;
75                 traffic(src, dest) = traffic(src,dest) + 1;
76                 dest = dim2i(y,x,z,k,n)+1;
77                 traffic(src, dest) = traffic(src,dest) + 1;
78                 dest = dim2i(y,z,x,k,n)+1;
79                 traffic(src, dest) = traffic(src,dest) + 1;
80             else
81                 z = 0;
82                 dest = dim2i(y, x, z, k, n)+1;
83                 traffic(src, dest) = traffic(src,dest) + 1;
84             end
85         end
86     end
87
88
89 % Nearest Neighbor Traffic
90 elseif(pattern == 3)
91     for j=1:runs
92         for src=1:len
93             [x,y,z] = i2dim(src-1, k, n);
94             if(n > 2)
95                 dest = dim2i(mod(x+1, k), y, z, k, n)+1;
96                 traffic(src,dest) = traffic(src,dest) + 1;
97                 dest = dim2i(mod(x-1, k), y, z, k, n)+1;
98                 traffic(src,dest) = traffic(src,dest) + 1;
99                 dest = dim2i(x, mod(y+1, k), z, k, n)+1;
100                traffic(src,dest) = traffic(src,dest) + 1;
101                dest = dim2i(x, mod(y-1, k), z, k, n)+1;
102                traffic(src,dest) = traffic(src,dest) + 1;

```



```

103         dest = dim2i(x, y, mod(z+1, k), k, n)+1;
104         traffic(src,dest) = traffic(src,dest) + 1;
105         dest = dim2i(x, y, mod(z-1, k), k, n)+1;
106         traffic(src,dest) = traffic(src,dest) + 1;
107
108     else
109         dest = dim2i(mod(x+1, k), y, z, k, n)+1;
110         traffic(src,dest) = traffic(src,dest) + 1;
111         dest = dim2i(mod(x-1, k), y, z, k, n)+1;
112         traffic(src,dest) = traffic(src,dest) + 1;
113         dest = dim2i(x, mod(y+1, k), z, k, n)+1;
114         traffic(src,dest) = traffic(src,dest) + 1;
115         dest = dim2i(x, mod(y-1, k), z, k, n)+1;
116         traffic(src,dest) = traffic(src,dest) + 1;
117     end
118
119 end
120 end
121
122 % Uniform Random Traffic
123 elseif(pattern == 4)
124     for j=1:runs
125         for src=1:len
126             dest = ceil(rand*len);
127             traffic(src,dest) = traffic(src,dest) + ceil(rand()*runs);
128         end
129     end
130
131 % Bit Compliment Traffic
132 elseif(pattern == 5)
133     for j=1:runs
134         for src=1:len
135             [x,y,z] = i2dim(src-1, k, n);
136             if n < 3
137                 dest = dim2i((k-x-1), (k-y-1), z, k, n)+1;
138             else
139                 dest = dim2i((k-x-1), (k-y-1), (k-z-1), k, n)+1;
140             end
141             traffic(src,dest) = traffic(src,dest) + 1;
142         end
143     end
144 end
145
146 if (hs > 0)
147     nums=(1:1:(k^n));
148     for (j=1:ceil(((k^n)*hs)))
149         rand_value = ceil(rand()*length(find(nums>=1)));
150         [rows,cols,vals] = find(nums>=1);
151         x_ = rows(rand_value);
152         y_ = cols(rand_value);
153         h_s(j) = nums(x_,y_);
154         nums(x_,y_) = 0;
155         traffic(:,h_s(j)) = traffic(:,h_s(j))*4;
156         [strng] = sprintf('Hotspot Generated at node %d',h_s(j));
157         disp(strng);
158     end
159 end
160
161 if (output_text_file == 1)
162     file_1 = fopen('/Users/lydick/Desktop/file.txt','w');
163     for (i=1:k^n)
164         fprintf(file_1,'%0d',traffic(i,1));
165         for (j=2:k^n)
166             if (j == i)
167                 fprintf(file_1,' 00');

```

```

168         else
169             fprintf(file_1,' 0%d',traffic(i,j));
170         end
171     end
172     fprintf(file_1,'\n');
173 end
174 fclose(file_1);
175 %traffic
176 end
177
178 % For each traffic dest/source pair
179 % These are delay independent/load independent visualizations.
180 %     Mostly this is the visual demand... not the traffic.
181
182 if (option <= 3)
183     for src=1:len
184         for dest=1:len
185             if(traffic(src,dest) > 0)
186                 if (option == 1)
187                     ch_load = ch_load + traffic(src,dest)*dimord_route(src-1, dest-1, k, n);
188
189                 elseif (option == 2)
190                     ch_load = ch_load + traffic(src,dest)*dirord_route(src-1, dest-1, k, n);
191
192                 elseif (option == 3)
193                     ch_load = ch_load + traffic(src,dest)*minobliv_route(src-1, dest-1, k, n);
194
195                 end
196             end
197         end
198     end
199 end
200
201
202 elseif (option == 7)
203     [ch_load, delay] = do_cqr2(traffic, k, n, failure);
204     avg_delay = mean(delay);
205     max_delay = max(delay);
206     min_delay = min(delay);
207     [strng] = sprintf('Mean Delay: \t%0.3g\nMax Delay: \t%d\nMin Delay: \t%d', avg_delay, max_delay, min_delay);
208     disp(strng);
209
210 elseif (option == 6)
211     [ch_load, delay] = do_cqr1(traffic, k, n, failure);
212     avg_delay = mean(delay);
213     max_delay = max(delay);
214     min_delay = min(delay);
215     [strng] = sprintf('Mean Delay: \t%0.3g\nMax Delay: \t%d\nMin Delay: \t%d', avg_delay, max_delay, min_delay);
216     disp(strng);
217
218 elseif (option == 5)
219     [ch_load, delay] = do_chaos(traffic, k, n, failure);
220     avg_delay = mean(delay);
221     max_delay = max(delay);
222     min_delay = min(delay);
223     [strng] = sprintf('Mean Delay: \t%0.3g\nMax Delay: \t%d\nMin Delay: \t%d', avg_delay, max_delay, min_delay);
224     disp(strng);
225
226 elseif (option == 4)
227     [ch_load, delay, failed_nodes_queue] = do_minadaptive(traffic, k, n, failure);
228     avg_delay = mean(delay);
229     max_delay = max(delay);
230     min_delay = min(delay);
231     [strng] = sprintf('Mean Delay: \t%0.3g\nMax Delay: \t%d\nMin Delay: \t%d', avg_delay, max_delay, min_delay);
232     disp(strng);

```

```

233
234 end
235
236 subplot(1,2,1);
237 max_ch_load = max(max(ch_load));
238 ch_load = ch_load / max_ch_load;
239 std_dev = mean(std(ch_load));
240 if (plot == 1)
241     view_load(ch_load, k, n, 1.0-(std_dev), 1.0-(std_dev*0.5));
242 end
243
244
245 count = 1;
246 while (count <= (k^n) && plot == 1 && n > 2)
247     [val_x,val_y,val_z] = i2dim(count,k,n);
248     if (find(failed_nodes_queue == count) >= 1)
249         [val1_x,val1_y,val1_z] = sphere(30);
250         val1_x = val1_x.*0.05; val1_y = val1_y.*0.05; val1_z = val1_z.*0.05;
251         val1_x = val1_x + val_x; val1_y = val1_y + val_y; val1_z = val1_z + val_z;
252         plot3(val1_x,val1_y,val1_z,'k');
253
254     elseif ((hs > 0) && (nums(count) == 0))
255         [val1_x,val1_y,val1_z] = sphere(30);
256         val1_x = val1_x.*0.05; val1_y = val1_y.*0.05; val1_z = val1_z.*0.05;
257         val1_x = val1_x + val_x; val1_y = val1_y + val_y; val1_z = val1_z + val_z;
258         plot3(val1_x,val1_y,val1_z,'r');
259
260     else
261         [val1_x,val1_y,val1_z] = sphere(30);
262         val1_x = val1_x.*0.05; val1_y = val1_y.*0.05; val1_z = val1_z.*0.05;
263         val1_x = val1_x + val_x; val1_y = val1_y + val_y; val1_z = val1_z + val_z;
264         plot3(val1_x,val1_y,val1_z,'b');
265     end
266     count = count + 1;
267 end
268
269 link_load = zeros(k^n*6,1);
270 if (n == 3)
271     for i=1:k^n
272         [x,y,z] = i2dim(i-1,k,n);
273         link_load(6*(i-1)+1) = ch_load(i,dim2i(mod(x+1,k),y,z,k,n)+1);
274         link_load(6*(i-1)+2) = ch_load(i,dim2i(mod(x-1,k),y,z,k,n)+1);
275         link_load(6*(i-1)+3) = ch_load(i,dim2i(x,mod(y+1,k),z,k,n)+1);
276         link_load(6*(i-1)+4) = ch_load(i,dim2i(x,mod(y-1,k),z,k,n)+1);
277         link_load(6*(i-1)+5) = ch_load(i,dim2i(x,y,mod(z+1,k),k,n)+1);
278         link_load(6*(i-1)+6) = ch_load(i,dim2i(x,y,mod(z-1,k),k,n)+1);
279     end
280 else
281     [r,c,v] = find(ch_load);
282     link_load = v;
283 end
284
285 hold off;
286 t = [0:1/50:1];
287 subplot(1,2,2);
288 [t1,t2] = hist(link_load,t);
289 hist(link_load,t);
290 axis([-0.1 1.1 0 max(t1)])
291 ylabel('Number of Links');
292 xlabel('Normalized Load per Bi-Directional Link');
293 title('Link Load Histogram');
294 legend('Load Distribution');
295 avg_thpt = mean(link_load);
296 std_thpt = std(link_load);
297 [strng] = sprintf('Avg Link Load: \t%2.2f%\nStd Dev: \t%2.2f', avg_thpt*100,std_thpt*100);

```

```
298 disp(strng);
299 error_status = 0;
300 set(gcf, 'Position', [50,500,650,300]);
```

## C.1.2 file: view\_load.m

```
1 function view_load(route, k, n, mid, high)
2 %
3 % Use: view_load(route, k, n, mid, high)
4 %     load = 2-D matrix containing load info between each pair of
5 %     nodes
6 %     k , n = size of torus
7 %     mid = threshold for coloring loads < mid to green color
8 %     high = thrshold for coloring loads > high to red color
9 %     all other loads will be colored yellow
10 %
11
12 len = k^n;
13 i = 1;
14 for r=1:len
15     for c = 1:len
16         if(route(r,c) > 0)
17             [sx,sy,sz] = i2dim(r-1, k, n);
18             [dx,dy,dz] = i2dim(c-1, k, n);
19             if(abs(sx-dx) <= 1)
20                 x(i,:) = [sx dx];
21             else
22                 x(i,:) = [k-1, k];
23             end
24
25             if(abs(sy-dy) <= 1)
26                 y(i,:) = [sy dy];
27             else
28                 y(i,:) = [k-1, k];
29             end
30
31             if(n > 2)
32                 if(abs(sz-dz) <= 1)
33                     z(i,:) = [sz dz];
34                 else
35                     z(i,:) = [k-1, k];
36                 end
37             else
38                 z(i,:) = [0 0];
39             end
40             if(route(r,c) < mid)
41                 color = 'g';
42             elseif(route(r,c) > high)
43                 color = 'r';
44             else
45                 color = 'y';
46             end
47
48             if (n > 2)
49                 plot3(x', y', z', color);
50             else
51                 plot(x',y', color);
52             end
53         end
54
55         hold on;
56         %i = i + 1;
57     end
58 end
59 end
60
61 grid on;
62 if (n > 2)
63     axis ([0 k 0 k 0 k], 'square');
```

```
64     xlabel('z');
65 else
66     axis ([0 k 0 k], 'square');
67 end
68
69 xlabel('x');
70 ylabel('y');
71 title('Visualization of Link Utilizations');
72
```

## C.2 Routing Scripts

### C.2.1 file: dim2i.m

```
1 function i = dim2i(x,y,z,k,n)
2 %
3 %
4 % Use: function i = dim2i(x, y, z, k, n)
5 % i = index position
6 % k = size of k for k-ary n-cube
7 % n = size of n for k-ary n-cube
8 %
9
10 if(n > 2)
11     i = x + k*y + k*k*z;
12 else
13     i = x + k*y;
14 end
```

## C.2.2 file: i2dim.m

```
1 function [x,varargout] = i2dim(i, k, n)
2 %
3 %
4 % Use: function [x,y,z] = index2dim(i, k, n)
5 % i = index position
6 % k = size of k for k-ary n-cube
7 % n = size of n for k-ary n-cube
8 %
9 nout = max(nargout,1)-1;
10 varargout = cell(nout, 1);
11 x = rem(i,k);
12 if(n > 1)
13     varargout(1) = num2cell(rem(floor(i/k),k));
14 end
15
16 if(n > 2)
17     varargout(2) = num2cell(floor(i/(k*k)));
18 end
```



### C.2.3 file: dimord\_route.m

```
1 function route = dimord_route(src, dest, k, n)
2 %
3 % Use: route = dimord_route(src, dest, k, n)
4 %
5
6 [sx, sy, sz] = i2dim(src, k, n);
7 [dx, dy, dz] = i2dim(dest, k, n);
8
9
10 len = k^n;
11 route = zeros(len);
12
13 previ = src;
14
15 % Route in the x-dimension
16
17 currentx = sx;
18 dirx = dimord(sx, dx, k);
19 if dirx == -2
20     dirx = sign(rand(1)-.5);
21 end
22
23 while currentx ~= dx
24     currentx = currentx + dirx;
25
26     if(currentx < 0)
27         currentx = currentx + k;
28     end
29     if(currentx >= k)
30         currentx = currentx - k;
31     end
32
33     currenti = dim2i(currentx,sy,sz,k,n);
34     route(previ+1, currenti+1) = 1;
35     previ = currenti;
36 end
37
38 % Route in the y-dimension
39 if (n > 1)
40     currenty = sy;
41     diry = dimord(sy, dy, k);
42     if diry == -2
43         diry = sign(rand(1)-.5);
44     end
45
46 while currenty ~= dy
47     currenty = currenty + diry;
48     if(currenty < 0)
49         currenty = currenty + k;
50     end
51
52     if(currenty >= k)
53         currenty = currenty - k;
54     end
55
56     currenti = dim2i(dx,currenty,sz,k,n);
57     route(previ+1, currenti+1) = 1;
58     previ = currenti;
59 end
60
61 end % y-dimension
62
63 % Route in the z-dimension
```

```
64 if (n > 2)
65 currentz = sz;
66 dirz = dimord(sz, dz, k);
67 if dirz == -2
68     dirz = sign(rand(1)-.5);
69 end
70
71 while currentz ~= dz
72     currentz = currentz + dirz;
73     if(currentz < 0)
74         currentz = currentz + k;
75     end
76     if(currentz >= k)
77         currentz = currentz - k;
78     end
79
80     currenti = dim2i(dx, dy, currentz,k,n);
81     route(previ+1, currenti+1) = 1;
82     previ = currenti;
83 end
84
85 end
```

## C.2.4 file: dirord\_route.m

```
1 function route = dirord_route(src, dest, k, n)
2 %
3 % Use: route = dirord_route(src, dest, k, n)
4 %       performs direction order routing
5 %       (+x, +y, +z, -x, -y, -z)
6 %
7
8 [sx, sy, sz] = i2dim(src, k, n);
9 [dx, dy, dz] = i2dim(dest, k, n);
10
11 len = k^n;
12 route = zeros(len);
13
14 previ = src;
15
16 % Get directions in each dimension
17 currentx = sx;
18 dirx = dimord(sx, dx, k);
19 if dirx == -2
20     dirx = sign(rand(1)-.5);
21 end
22
23 if (n > 1)
24     currenty = sy;
25     diry = dimord(sy, dy, k);
26     if diry == -2
27         diry = sign(rand(1)-.5);
28     end
29 end
30
31 if (n > 2)
32     currentz = sz;
33     dirz = dimord(sz, dz, k);
34     if dirz == -2
35         dirz = sign(rand(1)-.5);
36     end
37 end
38
39
40 % Route in the +x-dimension
41 if(dirx > 0)
42     while currentx ~= dx
43         currentx = currentx + dirx;
44
45         if(currentx < 0)
46             currentx = currentx + k;
47         end
48         if(currentx >= k)
49             currentx = currentx - k;
50         end
51
52         currenti = dim2i(currentx,sy,sz,k,n);
53         route(previ+1, currenti+1) = 1;
54         previ = currenti;
55     end
56 end
57
58 % Route in the +y-dimension
59 if(n > 1)
60     if(diry > 0)
61         while currenty ~= dy
62             currenty = currenty + diry;
63             if(currenty < 0)
```

```

64         currenty = currenty + k;
65     end
66
67     if(currenty >= k)
68         currenty = currenty - k;
69     end
70
71     currenti = dim2i(currentx,currenty,sz,k,n);
72     route(previ+1, currenti+1) = 1;
73     previ = currenti;
74 end
75     end
76 end
77
78
79 % Route in the +z-dimension
80 if(n > 2)
81     if(dirz > 0)
82         while currentz ~= dz
83             currentz = currentz + dirz;
84             if(currentz < 0)
85                 currentz = currentz + k;
86             end
87             if(currentz >= k)
88                 currentz = currentz - k;
89             end
90
91             currenti = dim2i(currentx, currenty, currentz,k,n);
92             route(previ+1, currenti+1) = 1;
93             previ = currenti;
94         end
95     end
96 end
97
98 % Route in the -x-dimension
99 if(dirx < 0)
100     while currentx ~= dx
101         currentx = currentx + dirx;
102
103         if(currentx < 0)
104             currentx = currentx + k;
105         end
106         if(currentx >= k)
107             currentx = currentx - k;
108         end
109
110         currenti = dim2i(currentx, currenty, currentz,k,n);
111         route(previ+1, currenti+1) = 1;
112         previ = currenti;
113     end
114 end
115
116 % Route in the -y-dimension
117 if(n > 1)
118     if(diry < 0)
119         while currenty ~= dy
120             currenty = currenty + diry;
121             if(currenty < 0)
122                 currenty = currenty + k;
123             end
124
125             if(currenty >= k)
126                 currenty = currenty - k;
127             end
128

```

```

129         currenti = dim2i(currentx, currenty, currentz,k,n);
130         route(previ+1, currenti+1) = 1;
131         previ = currenti;
132     end
133 end
134 end
135
136 % Route in the -z-dimension
137 if(n > 2)
138     if(dirz < 0)
139         while currentz ~= dz
140             currentz = currentz + dirz;
141             if(currentz < 0)
142                 currentz = currentz + k;
143             end
144             if(currentz >= k)
145                 currentz = currentz - k;
146             end
147
148             currenti = dim2i(currentx, currenty, currentz,k,n);
149             route(previ+1, currenti+1) = 1;
150             previ = currenti;
151         end
152     end
153 end

```

## C.2.5 file: minobliv\_route.m

```
1 function route = minobliv_route(src, dest, k, n)
2 %
3 % Use: route = minobliv_route(src, dest, k, n)
4 %         performs minimum oblivious routing
5 %
6 %
7
8 len = k^n;
9 route = zeros(len);
10
11 previ = src;
12
13 % Get coordinates of source & destination
14 [sx, sy, sz] = i2dim(src, k, n);
15 [dx, dy, dz] = i2dim(dest, k, n);
16
17 [dirx, deltax] = dimord(sx, dx, k);
18
19 if(n > 1)
20     [diry, deltay] = dimord(sy, dy, k);
21 end
22 if(n > 2)
23     [dirz, deltaz] = dimord(sz, dz, k);
24 end
25
26 % Calculate an intermediate node
27 % First find the bounding box
28 ix = round(deltax*rand(1) + sx);
29
30 if(ix < 0)
31     ix = ix + k;
32 end
33 if(ix >= k)
34     ix = ix - k;
35 end
36
37 if(n > 1)
38     iy = round(deltay*rand(1) + sy);
39     if(iy < 0)
40         iy = iy + k;
41     end
42     if(iy >= k)
43         iy = iy - k;
44     end
45
46 end
47
48 if(n > 2)
49     iz = round(deltaz*rand(1) + sz);
50     if(iz < 0)
51         iz = iz + k;
52     end
53     if(iz >= k)
54         iz = iz - k;
55     end
56 end
57
58 inter_node = dim2i(ix,iy,iz,k,n);
59 route = dimord_rand_route(src, inter_node, k, n);
60 route = route + dimord_rand_route(inter_node, dest, k, n);
61
```

## C.2.6 file: do\_minadaptive.m

```
1 function [ch_load, delay, failed_nodes_queue] = do_minadaptive(traffic, k, n, loss);
2
3 % Use: [ch_load, delay, failed_nodes_queue] = do_chaos(traffic, k, n, loss);
4
5 size = k^n;
6
7 threshold = 1;
8 flit_size = 0.1;
9 total_deflections = 0;
10 dropped = 0;
11
12 % start the queue with a value which should never occur. Pi.
13 failed_nodes_queue = pi;
14
15 total_original_traffic = 0;
16
17 ch_load = zeros(size,size);
18 delay = [1:size];
19 for (i=0:size-1)
20     delay(i+1) = 0;
21 end
22
23 failure = delay;
24
25 total_original_traffic = sum(sum(traffic>0));
26
27 for (i=0:size-1)
28     if (rand() <= loss)
29         failure(i+1) = 1;
30         [str,err] = sprintf('Nodal Failure Simulated at node: %d',i);
31         disp(str);
32         failed_nodes_queue = push(failed_nodes_queue, i);
33         dropped = dropped + sum(traffic(i+1,:)>0);
34         for (j=1:size)
35             traffic(i+1,j) = 0;
36         end
37         dropped = dropped + sum(traffic(:,i+1)>0);
38         for (j=1:size)
39             traffic(j,i+1) = 0;
40         end
41     end
42 end
43 end
44
45 total_traffic = 0;
46
47 for (i=1:size)
48     for (j=1:size)
49         total_traffic = total_traffic + traffic(i,j);
50     end
51 end
52
53
54 % While there is traffic within the temp_traffic matrix...
55
56 while (length(find(traffic>=1))~=0)
57
58     temp_ch_load = zeros(size,size);
59     temp_traffic = traffic;
60     [strng,err] = sprintf('...traffic demand of %d packets', length(find(traffic)));
61     disp(strng);
62
63     while (length(find(temp_traffic>=1))~=0)
```

```

64
65 % for one node/src pair, randomly select src/dest pair
66 rand_value = ceil(rand()*length(find(temp_traffic>=1)));
67 [rows,cols,vals] = find(temp_traffic>=1);
68 src = rows(rand_value);
69 dest = cols(rand_value);
70
71
72 % Get coordinates for those values
73 [sx, sy, sz] = i2dim(src-1, k, n);
74 [dx, dy, dz] = i2dim(dest-1, k, n);
75
76 [dirx, deltax] = dimord(sx, dx, k);
77
78 if(n > 1)
79     [diry, deltay] = dimord(sy, dy, k);
80     dirz = 0;
81
82 end
83 if(n > 2)
84     [dirz, deltaz] = dimord(sz, dz, k);
85
86 end
87
88 % these print the directions, but before selecting which direction
89 % to take, we need to consider that that path may be not as
90 % described... I don't think I'm doing this correctly... arg!!
91
92 % Next hop coordinates first equal the source location.
93 nx = sx;
94 ny = sy;
95 nz = sz;
96 test_x = -1;
97 test_y = -1;
98 test_z = -1;
99 val_test_x = Inf;
100 val_test_y = Inf;
101 val_test_z = Inf;
102 flag_misroute = 0;
103 route_direction = 'x';
104
105 % Test locations, needing to see which temp_traffic values are
106 % smallest
107
108 if (dirx ~= 0)
109     test_x = dim2i(mod(nx+dirx,k),ny,nz,k,n)+1;
110     val_test_x = temp_ch_load(src,test_x);
111     if (failure(test_x) == 1)
112         val_test_x = Inf;
113     end
114 end
115
116 if (diry ~= 0)
117     test_y = dim2i(nx,mod(ny+diry,k),nz,k,n)+1;
118     val_test_y = temp_ch_load(src,test_y);
119     if (failure(test_y) == 1)
120         val_test_y = Inf;
121     end
122 end
123
124 if (dirz ~= 0)
125     test_z = dim2i(nx,ny,mod(nz+dirz,k),k,n)+1;
126     val_test_z = temp_ch_load(src,test_z);
127     if (failure(test_z) == 1)
128         val_test_z = Inf;

```



```

129         end
130     end
131
132     flag_no_progress = 0;
133     min_val = min(val_test_y, min(val_test_z, val_test_x));
134     if (val_test_z == Inf && val_test_y == Inf && val_test_x == Inf && src ~= dest)
135         flag_no_progress = 1;
136     end
137
138     % If the min val will go over the threshold, we must misroute.
139     % First recalculate the minimum of the opposite direction.
140     if (min_val + flit_size > threshold && flag_no_progress == 0)
141         flag_misroute = 1;
142         total_deflections = total_deflections + 1;
143         traffic(src,dest) = traffic(src,dest) + 1;
144     elseif (flag_no_progress == 1)
145         dropped = dropped + 1;
146     else
147
148         % If all are the same, this randomizes the direction to route.
149         switch(ceil(rand()*3))
150             case (1)
151                 switch(min_val)
152                     case (val_test_x)
153                         route_direction = 'x';
154                         next_hop = test_x;
155                     case (val_test_y)
156                         route_direction = 'y';
157                         next_hop = test_y;
158                     case (val_test_z)
159                         route_direction = 'z';
160                         next_hop = test_z;
161                 end
162             case (2)
163                 switch(min_val)
164                     case (val_test_z)
165                         route_direction = 'z';
166                         next_hop = test_z;
167                     case (val_test_x)
168                         route_direction = 'x';
169                         next_hop = test_x;
170                     case (val_test_y)
171                         route_direction = 'y';
172                         next_hop = test_y;
173                 end
174             case (3)
175                 switch(min_val)
176                     case (val_test_y)
177                         route_direction = 'y';
178                         next_hop = test_y;
179                     case (val_test_z)
180                         route_direction = 'z';
181                         next_hop = test_z;
182                     case (val_test_x)
183                         route_direction = 'x';
184                         next_hop = test_x;
185                 end
186         end
187     end
188
189 end
190
191 if (flag_misroute == 1)
192     delay(src) = delay(src) + 1;
193 end

```

```

194
195
196     temp_traffic(src,dest) = temp_traffic(src,dest) - 1;
197     traffic(src,dest) = traffic(src,dest) - 1;
198     if (src ~= next_hop && flag_misroute == 0 && flag_no_progress == 0)
199         traffic(next_hop,dest) = traffic(next_hop,dest) + 1;
200         temp_ch_load(src,next_hop) = temp_ch_load(src,next_hop) + flit_size;
201     end
202
203
204
205     end
206
207 % This gets us to every src/dest going one hop.
208 ch_load = ch_load + temp_ch_load;
209 for (i=1:size)
210     traffic(i,i) = 0;
211 end
212
213 end
214
215 [strng, err] = sprintf('Total deflections: %d', total_deflections);
216 disp(strng);
217 [strng, err] = sprintf('Total number of dropped packets from failed nodes: %d', dropped);
218 disp(strng);
219 [strng, err] = sprintf('Total traffic demand: %d', total_original_traffic);
220 disp(strng);
221 [strng, err] = sprintf('Percent of dropped traffic: %0.03g', dropped/total_original_traffic*100);
222 disp(strng);
223
224
225
226

```

## C.2.7 file: do\_cqr1.m

```
1 function [ch_load, delay] = do_cqr1(traffic, k, n, loss)
2
3 % Use: [ch_load, delay] = do_cqr1(traffic, k, n, loss);
4 size = k^n;
5 threshold = 100;
6 flit_size = 0.1;
7 total_deflections = 0;
8 dropped = 0;
9 time = 0;
10 ch_load = zeros(size,size);
11 delay = zeros(size,1);
12 failure = zeros(size,1);
13
14 % Calculate and simulate nodal failures.
15 for i=0:size-1
16     if (rand() <= loss)
17         failure(i+1) = 1;
18         [str,err] = sprintf('Nodal Failure Simulated at node: %d',i);
19         disp(str);
20         disp(err);
21         for j=1:size
22             traffic(i+1,j) = 0;
23             traffic(j,i+1) = 0;
24         end
25     end
26 end
27
28 % get rid of any traffic src=dest
29 for i=1:k^n
30     traffic(i,i) = 0;
31 end
32
33 % Find total traffic demand
34 total_traffic = sum(sum(traffic));
35
36 temp_traffic = traffic; % preserve the original traffic matrix.
37 packet_matrix = cqr_transformTrafficMatrix(temp_traffic); % convert to the packet matrix
38 extra_temp_traffic = temp_traffic;
39 packet_queues = zeros(k^n,6);
40
41 while(length(find(temp_traffic>=1))~=0)
42
43     %Find a random packet, remove it, change the parameters, and put it
44     %back into the matrix.
45     [src,dest] = cqr_getRandomValue(temp_traffic);
46     [xdir,ydir,zdir,packet_queues] = cqr_chooseQuadrant(src,dest,n,k,packet_queues);
47     [dx,dy,dz] = i2dim(dest-1,k,n);
48     [sx,sy,sz] = i2dim(src-1,k,n);
49     [d_x,deltax] = dimord(sx,dx,k);
50     [d_y,deltay] = dimord(sy,dy,k);
51     [d_z,deltaz] = dimord(sz,dz,k);
52     time = time + abs(deltax) + abs(deltay) + abs(deltaz);
53     if (sign(d_x) ~= sign(xdir)) delay(src) = delay(src) + (k-(abs(deltax)+abs(deltax))); end
54     if (sign(d_y) ~= sign(ydir)) delay(src) = delay(src) + (k-(abs(deltay)+abs(deltay))); end
55     if (sign(d_z) ~= sign(zdir)) delay(src) = delay(src) + (k-(abs(deltaz)+abs(deltaz))); end
56     [pk1,packet_matrix,extra_temp_traffic] = cqr_removePacket(packet_matrix,src,dest,extra_temp_traffic);
57     pk1.initialized = 1;
58     pk1.x_dir = xdir;
59     pk1.y_dir = ydir;
60     pk1.z_dir = zdir;
61     [packet_matrix,extra_temp_traffic] = cqr_addPacket(pk1,packet_matrix,src,dest,extra_temp_traffic);
62     temp_traffic(src,dest) = temp_traffic(src,dest) - 1;
63 end
```

```

64
65 % we should now have fully transformed and calculated packet/traffic
66 % matrices
67
68 while (length(find(traffic>=1))~=0)
69
70     temp_ch_load = zeros(size,size);
71     temp_traffic = traffic;
72     temp_packet_matrix = packet_matrix;
73
74     while (length(find(temp_traffic>=1))~=0)
75
76         % for one node/src pair, randomly select src/dest pair
77         [src,dest] = cqr_getRandomValue(temp_traffic);
78
79         % Get coordinates for those values
80         [pkt,temp_packet_matrix,temp_traffic] = cqr_removePacket(temp_packet_matrix,src,dest,temp_traffic);
81         [pkt_global,packet_matrix,traffic] = cqr_removePacket(packet_matrix,src,dest,traffic);
82
83         dirx = pkt.x_dir;
84         diry = pkt.y_dir;
85         dirz = pkt.z_dir;
86
87         [sx, sy, sz] = i2dim(src-1, k, n);
88
89         % Next hop coordinates first equal the source location.
90         nx = sx;
91         ny = sy;
92         nz = sz;
93         test_x = -1;
94         test_y = -1;
95         test_z = -1;
96         val_test_x = Inf;
97         val_test_y = Inf;
98         val_test_z = Inf;
99         flag_misroute = 0;
100        flag_no_progress = 0;
101
102        % Test locations, needing to see which temp_traffic values are
103        % smallest
104        if (dirx ~= 0)
105            test_x = dim2i(mod(nx+dirx,k),ny,nz,k,n)+1;
106            val_test_x = temp_ch_load(src,test_x);
107            if (failure(test_x) == 1)
108                val_test_x = Inf;
109            end
110        end
111
112        if (diry ~= 0)
113            test_y = dim2i(nx,mod(ny+diry,k),nz,k,n)+1;
114            val_test_y = temp_ch_load(src,test_y);
115            if (failure(test_y) == 1)
116                val_test_y = Inf;
117            end
118        end
119
120        if (dirz ~= 0)
121            test_z = dim2i(nx,ny,mod(nz+dirz,k),k,n)+1;
122            val_test_z = temp_ch_load(src,test_z);
123            if (failure(test_z) == 1)
124                val_test_z = Inf;
125            end
126        end
127
128        min_val = min(val_test_y, min(val_test_z, val_test_x));

```

```

129
130     if (val_test_z == Inf && val_test_y == Inf && val_test_x == Inf)
131         if (src ~= dest)
132             flag_no_progress = 1;
133             [strng] = sprintf('%d-->%d :: xdir:%d ydir:%d zdir:%d TIME:%d',src,dest,dirx,diry,dirz,time);
134             disp(strng);
135         else
136             %packet reached its destination... DROP IT!
137         end
138     end
139
140     % If the min val will go over the threshold, we must misroute.
141     % First recalculate the minimum of the opposite direction.
142     if (min_val + flit_size > threshold && flag_no_progress == 0)
143         flag_misroute = 1;
144         total_deflections = total_deflections + 1;
145         delay(src) = delay(src) + 1;
146         % we queue back only globally - not locally. There is nothing
147         % we can do until the next iteration with the threshold being
148         % surpassed.
149         [packet_matrix,traffic] = cqr_addPacket(pkt_global,packet_matrix,src,dest,traffic);
150
151     elseif (flag_no_progress == 1 && src ~= dest)
152         dropped = dropped + 1;
153     elseif (src == dest)
154         disp('packet reached destination');
155     else
156
157         % If all are the same, this randomizes the direction to route.
158         switch(ceil(rand()*3))
159             case (1)
160                 switch(min_val)
161                     case (val_test_x)
162                         next_hop = test_x;
163                     case (val_test_y)
164                         next_hop = test_y;
165                     case (val_test_z)
166                         next_hop = test_z;
167                 end
168             case (2)
169                 switch(min_val)
170                     case (val_test_z)
171                         next_hop = test_z;
172                     case (val_test_x)
173                         next_hop = test_x;
174                     case (val_test_y)
175                         next_hop = test_y;
176                 end
177             case (3)
178                 switch(min_val)
179                     case (val_test_y)
180                         next_hop = test_y;
181                     case (val_test_z)
182                         next_hop = test_z;
183                     case (val_test_x)
184                         next_hop = test_x;
185                 end
186         end
187
188         % we now have our next_hop
189         [nx, ny, nz] = i2dim(next_hop-1, k, n);
190         [dx, dy, dz] = i2dim(dest-1, k, n);
191         d_x = dimord(nx,dx,k);
192         d_y = dimord(ny,dy,k);

```

```

194     d_z = dimord(nz,dz,k);
195     pkt_global.x_dir = abs(pkt_global.x_dir) * d_x;
196     pkt_global.y_dir = abs(pkt_global.y_dir) * d_y;
197     pkt_global.z_dir = abs(pkt_global.z_dir) * d_z;
198
199     z_dir_temp = pkt_global.z_dir;
200     if (n < 3)
201         z_dir_temp = 0;
202     end
203
204     if (pkt_global.x_dir == 0 && pkt_global.y_dir == 0 && z_dir_temp == 0)
205     else
206         [packet_matrix, traffic] = cqr_addPacket(pkt_global,packet_matrix,next_hop,dest,traffic);
207     end
208     temp_ch_load(src,next_hop) = temp_ch_load(src,next_hop) + flit_size;
209 end
210 end
211
212 % This gets us to every src/dest going one hop.
213 ch_load = ch_load + temp_ch_load;
214 for i=1:size
215     traffic(i,i) = 0;
216 end
217
218 end
219
220 [strng, err] = sprintf('Total deflections: %d', total_deflections);
221 disp(strng);
222 [strng, err] = sprintf('Total number of dropped packets from failed nodes: %d', dropped);
223 disp(strng);
224 [strng, err] = sprintf('Total traffic demand: %d', total_traffic);
225 disp(strng);
226 [strng, err] = sprintf('Percent of dropped traffic: %0.03g', dropped/total_traffic*100);
227 disp(strng);
228 [strng, err] = sprintf('Total number of hops: %d', time);
229 disp(strng);
230

```

## C.2.8 file: do\_cqr2.m

```
1 function [ch_load, delay] = do_cqr2(traffic, k, n, loss)
2 % Use: [ch_load, delay] = do_cqr2(traffic, k, n, loss);
3
4 size                = k^n;
5 threshold           = 100;
6 flit_size          = 0.1;
7 total_deflections  = 0;
8 dropped            = 0;
9 time = 0;
10 ch_load            = zeros(size,size);
11 delay              = zeros(size,1);
12 failure            = zeros(size,1);
13
14 % Calculate and simulate nodal failures.
15 for i=0:size-1
16     if (rand() <= loss)
17         failure(i+1) = 1;
18         [str,err] = sprintf('Nodal Failure Simulated at node: %d',i);
19         disp(str);
20         disp(err);
21         for j=1:size
22             traffic(i+1,j) = 0;
23             traffic(j,i+1) = 0;
24         end
25     end
26 end
27
28 % get rid of any traffic src=dest
29 for i=1:k^n
30     traffic(i,i) = 0;
31 end
32
33 % Find total traffic demand
34 total_traffic = sum(sum(traffic));
35
36 temp_traffic = traffic; % preserve the original traffic matrix.
37 packet_matrix = cqr_transformTrafficMatrix(temp_traffic); % convert to the packet matrix
38 extra_temp_traffic = temp_traffic;
39 packet_queues = zeros(k^n,6);
40
41 while(length(find(temp_traffic>=1))~=0)
42
43     %Find a random packet, remove it, change the parameters, and put it
44     %back into the matrix.
45     [src,dest] = cqr_getRandomValue(temp_traffic);
46     [xdir,ydir,zdir,packet_queues] = cqr_chooseQuadrant_v(src,dest,n,k,packet_queues);
47     [pk1,packet_matrix,extra_temp_traffic] = cqr_removePacket(packet_matrix,src,dest,extra_temp_traffic);
48     [dx,dy,dz] = i2dim(dest-1,k,n);
49     [sx,sy,sz] = i2dim(src-1,k,n);
50     [d_x,deltax] = dimord(sx,dx,k);
51     [d_y,deltay] = dimord(sy,dy,k);
52     [d_z,deltaz] = dimord(sz,dz,k);
53     if (sign(d_x) ~= sign(xdir)) delay(src) = delay(src) + (k-(abs(deltax)+abs(deltax))); end
54     if (sign(d_y) ~= sign(ydir)) delay(src) = delay(src) + (k-(abs(deltay)+abs(deltay))); end
55     if (sign(d_z) ~= sign(zdir)) delay(src) = delay(src) + (k-(abs(deltaz)+abs(deltaz))); end
56     pk1.initialized = 1;
57     pk1.x_dir = xdir;
58     pk1.y_dir = ydir;
59     pk1.z_dir = zdir;
60     [packet_matrix,extra_temp_traffic] = cqr_addPacket(pk1,packet_matrix,src,dest,extra_temp_traffic);
61     temp_traffic(src,dest) = temp_traffic(src,dest) - 1;
62 end
63
```

```

64 % we should now have fully transformed and calculated packet/traffic
65 % matrices
66
67
68 while (length(find(traffic>=1))~=0)
69
70     temp_ch_load = zeros(size,size);
71     temp_traffic = traffic;
72     temp_packet_matrix = packet_matrix;
73
74     while (length(find(temp_traffic>=1))~=0)
75         time = time + 1;
76         % for one node/src pair, randomly select src/dest pair
77         [src,dest] = cqr_getRandomValue(temp_traffic);
78
79         % Get coordinates for those values
80         [pkt,temp_packet_matrix,temp_traffic] = cqr_removePacket(temp_packet_matrix,src,dest,temp_traffic);
81         [pkt_global,packet_matrix,traffic] = cqr_removePacket(packet_matrix,src,dest,traffic);
82
83         dirx = pkt.x_dir;
84         diry = pkt.y_dir;
85         dirz = pkt.z_dir;
86
87         [sx, sy, sz] = i2dim(src-1, k, n);
88
89         % Next hop coordinates first equal the source location.
90         nx = sx;
91         ny = sy;
92         nz = sz;
93         test_x = -1;
94         test_y = -1;
95         test_z = -1;
96         val_test_x = Inf;
97         val_test_y = Inf;
98         val_test_z = Inf;
99         flag_misroute = 0;
100        flag_no_progress = 0;
101
102
103        % Test locations, needing to see which temp_traffic values are
104        % smallest
105        delta_total = abs(dirx) + abs(diry) + abs(dirz);
106
107        if (dirx ~= 0)
108            test_x = dim2i(mod(nx+sign(dirx),k),ny,nz,k,n)+1;
109            val_test_x = (1 + temp_ch_load(src,test_x)) * (1 - (abs(dirx) / delta_total));
110            if (failure(test_x) == 1)
111                val_test_x = Inf;
112            end
113        end
114
115        if (diry ~= 0)
116            test_y = dim2i(nx,mod(ny+sign(diry),k),nz,k,n)+1;
117            val_test_y = (1 + temp_ch_load(src,test_y)) * (1 - (abs(diry) / delta_total));
118            if (failure(test_y) == 1)
119                val_test_y = Inf;
120            end
121        end
122
123        if (dirz ~= 0)
124            test_z = dim2i(nx,ny,mod(nz+sign(dirz),k),k,n)+1;
125            val_test_z = (1 + temp_ch_load(src,test_z)) * (1 - (abs(dirz) / delta_total));
126            if (failure(test_z) == 1)
127                val_test_z = Inf;
128            end

```



```

129     end
130
131     min_val = min(val_test_y, min(val_test_z, val_test_x));
132
133     if (val_test_z == Inf && val_test_y == Inf && val_test_x == Inf)
134         if (src ~= dest)
135             flag_no_progress = 1;
136             [strng] = sprintf('packet cannot move(!) %d->%d :: xdir:%d ydir:%d zdir:%d TIME:%d',src,dest,dirx,diry,dirz,time);
137             disp(strng);
138         else
139             %packet reached its destination... DROP IT!
140         end
141     end
142
143     % If the min val will go over the threshold, we must misroute.
144     % First recalculate the minimum of the opposite direction.
145     if (min_val + flit_size > threshold && flag_no_progress == 0)
146         flag_misroute = 1;
147         total_deflections = total_deflections + 1;
148         delay(src) = delay(src) + 1;
149         % we queue back only globally - not locally. There is nothing
150         % we can do until the next iteration with the threshold being
151         % surpassed.
152         [packet_matrix,traffic] = cqr_addPacket(pkt_global,packet_matrix,src,dest,traffic);
153
154     elseif (flag_no_progress == 1 && src ~= dest)
155         dropped = dropped + 1;
156     elseif (src == dest)
157         disp('packet reached destination');
158     else
159
160         d_x1 = pkt_global.x_dir;
161         d_y1 = pkt_global.y_dir;
162         d_z1 = pkt_global.z_dir;
163         % If all are the same, this randomizes the direction to route.
164         switch(ceil(rand()*3))
165             case (1)
166                 switch(min_val)
167                     case (val_test_x)
168                         next_hop = test_x;
169                         d_x1 = (sign(pkt_global.x_dir) * -1) + pkt_global.x_dir;
170                     case (val_test_y)
171                         next_hop = test_y;
172                         d_y1 = (sign(pkt_global.y_dir) * -1) + pkt_global.y_dir;
173                     case (val_test_z)
174                         next_hop = test_z;
175                         d_z1 = (sign(pkt_global.z_dir) * -1) + pkt_global.z_dir;
176                 end
177             case (2)
178                 switch(min_val)
179                     case (val_test_z)
180                         next_hop = test_z;
181                         d_z1 = (sign(pkt_global.z_dir) * -1) + pkt_global.z_dir;
182                     case (val_test_x)
183                         next_hop = test_x;
184                         d_x1 = (sign(pkt_global.x_dir) * -1) + pkt_global.x_dir;
185                     case (val_test_y)
186                         next_hop = test_y;
187                         d_y1 = (sign(pkt_global.y_dir) * -1) + pkt_global.y_dir;
188                 end
189             case (3)
190                 switch(min_val)
191                     case (val_test_y)
192                         next_hop = test_y;
193                         d_y1 = (sign(pkt_global.y_dir) * -1) + pkt_global.y_dir;

```

```

194         case (val_test_z)
195             next_hop = test_z;
196             d_z1 = (sign(pkt_global.z_dir) * -1) + pkt_global.z_dir;
197         case (val_test_x)
198             next_hop = test_x;
199             d_x1 = (sign(pkt_global.x_dir) * -1) + pkt_global.x_dir;
200         end
201     end
202 end
203
204 % we now have our next_hop
205 [nx, ny, nz] = i2dim(next_hop-1, k, n);
206 [dx, dy, dz] = i2dim(dest-1, k, n);
207 pkt_global.x_dir = d_x1;
208 pkt_global.y_dir = d_y1;
209 pkt_global.z_dir = d_z1;
210 z_dir_temp = pkt_global.z_dir;
211 if (n < 3)
212     z_dir_temp = 0;
213 end
214
215 if (pkt_global.x_dir == 0 && pkt_global.y_dir == 0 && z_dir_temp == 0)
216     %do nothing, drop it, it's successfully reached the dest.
217 else
218     [packet_matrix, traffic] = cqr_addPacket(pkt_global,packet_matrix,next_hop,dest,traffic);
219 end
220 temp_ch_load(src,next_hop) = temp_ch_load(src,next_hop) + flit_size;
221
222 end
223 end
224
225 % This gets us to every src/dest going one hop.
226 ch_load = ch_load + temp_ch_load;
227 for i=1:size
228     traffic(i,i) = 0;
229 end
230
231 end
232
233 [strng, err] = sprintf('Total deflections: %d', total_deflections);
234 disp(strng);
235 [strng, err] = sprintf('Total number of dropped packets from failed nodes: %d', dropped);
236 disp(strng);
237 [strng, err] = sprintf('Total traffic demand: %d', total_traffic);
238 disp(strng);
239 [strng, err] = sprintf('Percent of dropped traffic: %0.03g', dropped/total_traffic*100);
240 disp(strng);
241

```

## C.2.9 file: cqr\_addPacket.m

```
1 function [packet_matrix, traffic_matrix] = cqr_addPacket(packet, packet_matrix_in, src, dest, traffic_matrix_in)
2
3 traffic_matrix_in(src,dest) = traffic_matrix_in(src,dest) + 1;
4 val = traffic_matrix_in(src,dest);
5
6 packet_matrix_in(src,dest,val).valid = packet.valid;
7 packet_matrix_in(src,dest,val).initialized = packet.initialized;
8 packet_matrix_in(src,dest,val).x_dir = packet.x_dir;
9 packet_matrix_in(src,dest,val).y_dir = packet.y_dir;
10 packet_matrix_in(src,dest,val).z_dir = packet.z_dir;
11
12 packet_matrix = packet_matrix_in;
13 traffic_matrix = traffic_matrix_in;
```

## C.2.10 file: cqr\_chooseQuadrant.m

```
1 function [dir_x,dir_y,dir_z,newqueues] = cqr_chooseQuadrant(src,dest,n,k,queues)
2 % Use: [x_dir,y_dir,z_dir,newqueues] = cqr_chooseQuadrant(src,dest,n,k,queues)
3
4     queue_x_pos = queues(src,1);
5     queue_x_neg = queues(src,2);
6     queue_y_pos = queues(src,3);
7     queue_y_neg = queues(src,4);
8     queue_z_pos = queues(src,5);
9     queue_z_neg = queues(src,6);
10
11     [sx, sy, sz] = i2dim(src-1, k, n);
12     [dx, dy, dz] = i2dim(dest-1, k, n);
13
14
15
16
17     if (n == 3)
18         % for quadrant I: (+,+,+)
19         [dir_q1x,q1_deltax] = dimord(sx, dx, k);
20         [dir_q1y,q1_deltay] = dimord(sy, dy, k);
21         [dir_q1z,q1_deltaz] = dimord(sz, dz, k);
22         nu_of_queues = 3;
23         if (sign(dir_q1x) < 0)
24             t_q = queue_x_neg;
25         elseif (sign(dir_q1x) > 0)
26             t_q = queue_x_pos;
27         else
28             t_q = 0;
29             nu_of_queues = nu_of_queues - 1;
30         end
31
32         if (sign(dir_q1y) < 0)
33             t_q = t_q + queue_y_neg;
34         elseif (sign(dir_q1y) > 0)
35             t_q = t_q + queue_y_pos;
36         else
37             nu_of_queues = nu_of_queues - 1;
38         end
39
40         if (sign(dir_q1z) < 0)
41             t_q = t_q + queue_z_neg;
42         elseif (sign(dir_q1z) > 0)
43             t_q = t_q + queue_z_pos;
44         else
45             nu_of_queues = nu_of_queues - 1;
46         end
47         t_q = t_q + 1;
48         if (nu_of_queues ~= 0) q1_val = (abs(q1_deltax) + abs(q1_deltay) + abs(q1_deltaz)) * (t_q/nu_of_queues);
49         else q1_val = Inf;
50         end
51
52         min_val = q1_val;
53         dir_x = sign(q1_deltax);
54         dir_y = sign(q1_deltay);
55         dir_z = sign(q1_deltaz);
56
57
58         % for quadrant II: (+ + -)
59         [dir_q2x,q2_deltax] = dimord(sx, dx, k);
60         [dir_q2y,q2_deltay] = dimord(sy, dy, k);
61         [dir_q2z,q2_deltaz] = dimord(sz, dz, k);
62         if (q2_deltaz > 0) q2_deltaz = q2_deltaz-k;
63         elseif (q2_deltaz < 0) q2_deltaz = q2_deltaz+k;
```

```

64     else q2_deltaz = 0;
65     end
66     nu_of_queues = 3;
67     if (sign(dir_q2x) < 0)
68         t_q = queue_x_neg;
69     elseif (sign(dir_q2x) > 0)
70         t_q = queue_x_pos;
71     else
72         t_q = 0;
73         nu_of_queues = nu_of_queues - 1;
74     end
75
76     if (sign(dir_q2y) < 0)
77         t_q = t_q + queue_y_neg;
78     elseif (sign(dir_q2y) > 0)
79         t_q = t_q + queue_y_pos;
80     else
81         nu_of_queues = nu_of_queues - 1;
82     end
83
84     if (sign(dir_q2z) < 0)
85         t_q = t_q + queue_z_neg;
86     elseif (sign(dir_q2z) > 0)
87         t_q = t_q + queue_z_pos;
88     else
89         nu_of_queues = nu_of_queues - 1;
90     end
91     t_q = t_q + 1;
92     if (nu_of_queues ~= 0) q2_val = (abs(q2_deltax) + abs(q2_deltay) + abs(q2_deltaz)) * (t_q/nu_of_queues);
93     else q2_val = Inf;
94     end
95
96     if (q2_val < min_val)
97         min_val = q2_val;
98         dir_x = sign(q2_deltax);
99         dir_y = sign(q2_deltay);
100        dir_z = sign(q2_deltaz);
101    end
102
103
104    % for quadrant III: (+ - +)
105    [dir_q3x,q3_deltax] = dimord(sx, dx, k);
106    [dir_q3y,q3_deltay] = dimord(sy, dy, k);
107    [dir_q3z,q3_deltaz] = dimord(sz, dz, k);
108    if (q3_deltay > 0) q3_deltay = q3_deltay-k;
109    elseif (q3_deltay < 0) q3_deltay = q3_deltay+k;
110    else q3_deltay = 0;
111    end
112    nu_of_queues = 3;
113    if (sign(dir_q3x) < 0)
114        t_q = queue_x_neg;
115    elseif (sign(dir_q3x) > 0)
116        t_q = queue_x_pos;
117    else
118        t_q = 0;
119        nu_of_queues = nu_of_queues - 1;
120    end
121
122    if (sign(dir_q3y) < 0)
123        t_q = t_q + queue_y_neg;
124    elseif (sign(dir_q3y) > 0)
125        t_q = t_q + queue_y_pos;
126    else
127        nu_of_queues = nu_of_queues - 1;
128    end

```

```

129
130     if (sign(dir_q3z) < 0)
131         t_q = t_q + queue_z_neg;
132     elseif (sign(dir_q3x) > 0)
133         t_q = t_q + queue_z_pos;
134     else
135         nu_of_queues = nu_of_queues - 1;
136     end
137     t_q = t_q + 1;
138     if (nu_of_queues ~= 0) q3_val = (abs(q3_deltax) + abs(q3_deltay) + abs(q3_deltaz)) * (t_q/nu_of_queues);
139     else q3_val = Inf;
140     end
141
142     if (q3_val < min_val)
143         min_val = q3_val;
144         dir_x = sign(q3_deltax);
145         dir_y = sign(q3_deltay);
146         dir_z = sign(q3_deltaz);
147     end
148
149
150     % for quadrant IV: (+ - -)
151     [dir_q4x,q4_deltax] = dimord(sx, dx, k);
152     [dir_q4y,q4_deltay] = dimord(sy, dy, k);
153     [dir_q4z,q4_deltaz] = dimord(sz, dz, k);
154     if (q4_deltay > 0) q4_deltay = q4_deltay-k;
155     elseif (q4_deltay < 0) q4_deltay = q4_deltay+k;
156     else q4_deltay = 0;
157     end
158     if (q4_deltaz > 0) q4_deltaz = q4_deltaz-k;
159     elseif (q4_deltaz < 0) q4_deltaz = q4_deltaz+k;
160     else q4_deltaz = 0;
161     end
162     nu_of_queues = 3;
163     if (sign(dir_q4x) < 0)
164         t_q = queue_x_neg;
165     elseif (sign(dir_q4x) > 0)
166         t_q = queue_x_pos;
167     else
168         t_q = 0;
169     nu_of_queues = nu_of_queues - 1;
170     end
171
172     if (sign(dir_q4y) < 0)
173         t_q = t_q + queue_y_neg;
174     elseif (sign(dir_q4y) > 0)
175         t_q = t_q + queue_y_pos;
176     else
177         nu_of_queues = nu_of_queues - 1;
178     end
179
180     if (sign(dir_q4z) < 0)
181         t_q = t_q + queue_z_neg;
182     elseif (sign(dir_q4x) > 0)
183         t_q = t_q + queue_z_pos;
184     else
185         nu_of_queues = nu_of_queues - 1;
186     end
187     t_q = t_q + 1;
188     if (nu_of_queues ~= 0) q4_val = (abs(q4_deltax) + abs(q4_deltay) + abs(q4_deltaz)) * (t_q/nu_of_queues);
189     else q4_val = Inf;
190     end
191
192     if (q4_val < min_val)
193         min_val = q4_val;

```

```

194         dir_x = sign(q4_deltax);
195         dir_y = sign(q4_deltay);
196         dir_z = sign(q4_deltaz);
197     end
198
199
200     % for quadrant V: (- + +)
201     [dir_q5x,q5_deltax] = dimord(sx, dx, k);
202     [dir_q5y,q5_deltay] = dimord(sy, dy, k);
203     [dir_q5z,q5_deltaz] = dimord(sz, dz, k);
204     if (q5_deltax > 0) q5_deltax = q5_deltax-k;
205     elseif (q5_deltax < 0) q5_deltax = q5_deltax+k;
206     else q5_deltax = 0;
207     end
208     nu_of_queues = 3;
209     if (sign(dir_q5x) < 0)
210         t_q = queue_x_neg;
211     elseif (sign(dir_q5x) > 0)
212         t_q = queue_x_pos;
213     else
214         t_q = 0;
215         nu_of_queues = nu_of_queues - 1;
216     end
217
218     if (sign(dir_q5y) < 0)
219         t_q = t_q + queue_y_neg;
220     elseif (sign(dir_q5y) > 0)
221         t_q = t_q + queue_y_pos;
222     else
223         nu_of_queues = nu_of_queues - 1;
224     end
225
226     if (sign(dir_q5z) < 0)
227         t_q = t_q + queue_z_neg;
228     elseif (sign(dir_q5z) > 0)
229         t_q = t_q + queue_z_pos;
230     else
231         nu_of_queues = nu_of_queues - 1;
232     end
233     t_q = t_q + 1;
234     if (nu_of_queues ~= 0) q5_val = (abs(q5_deltax) + abs(q5_deltay) + abs(q5_deltaz)) * (t_q/nu_of_queues);
235     else q5_val = Inf;
236     end
237
238     if (q5_val < min_val)
239         min_val = q5_val;
240         dir_x = sign(q5_deltax);
241         dir_y = sign(q5_deltay);
242         dir_z = sign(q5_deltaz);
243     end
244
245
246     % for quadrant VI: (- + -)
247     [dir_q6x,q6_deltax] = dimord(sx, dx, k);
248     [dir_q6y,q6_deltay] = dimord(sy, dy, k);
249     [dir_q6z,q6_deltaz] = dimord(sz, dz, k);
250     if (q6_deltax > 0) q6_deltax = q6_deltax-k;
251     elseif (q6_deltax < 0) q6_deltax = q6_deltax+k;
252     else q6_deltax = 0;
253     end
254     if (q6_deltaz > 0) q6_deltaz = q6_deltaz-k;
255     elseif (q6_deltaz < 0) q6_deltaz = q6_deltaz+k;
256     else q6_deltaz = 0;
257     end
258     nu_of_queues = 3;

```

```

259     if (sign(dir_q6x) < 0)
260         t_q = queue_x_neg;
261     elseif (sign(dir_q6x) > 0)
262         t_q = queue_x_pos;
263     else
264         t_q = 0;
265         nu_of_queues = nu_of_queues - 1;
266     end
267
268     if (sign(dir_q6y) < 0)
269         t_q = t_q + queue_y_neg;
270     elseif (sign(dir_q6y) > 0)
271         t_q = t_q + queue_y_pos;
272     else
273         nu_of_queues = nu_of_queues - 1;
274     end
275
276     if (sign(dir_q6z) < 0)
277         t_q = t_q + queue_z_neg;
278     elseif (sign(dir_q6x) > 0)
279         t_q = t_q + queue_z_pos;
280     else
281         nu_of_queues = nu_of_queues - 1;
282     end
283     t_q = t_q + 1;
284     if (nu_of_queues ~= 0) q6_val = (abs(q6_deltax) + abs(q6_deltay) + abs(q6_deltaz)) * (t_q/nu_of_queues);
285     else q6_val = Inf;
286     end
287
288     if (q6_val < min_val)
289         min_val = q6_val;
290         dir_x = sign(q6_deltax);
291         dir_y = sign(q6_deltay);
292         dir_z = sign(q6_deltaz);
293     end
294
295
296     % for quadrant VII: (- - +)
297     [dir_q7x,q7_deltax] = dimord(sx, dx, k);
298     [dir_q7y,q7_deltay] = dimord(sy, dy, k);
299     [dir_q7z,q7_deltaz] = dimord(sz, dz, k);
300     if (q7_deltax > 0) q7_deltax = q7_deltax-k;
301     elseif (q7_deltax < 0) q7_deltax = q7_deltax+k;
302     else q7_deltax = 0;
303     end
304     if (q7_deltay > 0) q7_deltay = q7_deltay-k;
305     elseif (q7_deltay < 0) q7_deltay = q7_deltay+k;
306     else q7_deltay = 0;
307     end
308     nu_of_queues = 3;
309     if (sign(dir_q7x) < 0)
310         t_q = queue_x_neg;
311     elseif (sign(dir_q7x) > 0)
312         t_q = queue_x_pos;
313     else
314         t_q = 0;
315         nu_of_queues = nu_of_queues - 1;
316     end
317
318     if (sign(dir_q7y) < 0)
319         t_q = t_q + queue_y_neg;
320     elseif (sign(dir_q7y) > 0)
321         t_q = t_q + queue_y_pos;
322     else
323         nu_of_queues = nu_of_queues - 1;

```



```

324     end
325
326     if (sign(dir_q7z) < 0)
327         t_q = t_q + queue_z_neg;
328     elseif (sign(dir_q7x) > 0)
329         t_q = t_q + queue_z_pos;
330     else
331         nu_of_queues = nu_of_queues - 1;
332     end
333     t_q = t_q + 1;
334     if (nu_of_queues ~= 0) q7_val = (abs(q7_deltax) + abs(q7_deltay) + abs(q7_deltaz)) * (t_q/nu_of_queues);
335     else q7_val = Inf;
336     end
337
338     if (q7_val < min_val)
339         min_val = q7_val;
340         dir_x = sign(q7_deltax);
341         dir_y = sign(q7_deltay);
342         dir_z = sign(q7_deltaz);
343     end
344
345
346     % for quadrant VIII: (- - -)
347     [dir_q8x,q8_deltax] = dimord(sx, dx, k);
348     [dir_q8y,q8_deltay] = dimord(sy, dy, k);
349     [dir_q8z,q8_deltaz] = dimord(sz, dz, k);
350     if (q8_deltax > 0) q8_deltax = q8_deltax-k;
351     elseif (q8_deltax < 0) q8_deltax = q8_deltax+k;
352     else q8_deltax = 0;
353     end
354     if (q8_deltay > 0) q8_deltay = q8_deltay-k;
355     elseif (q8_deltay < 0) q8_deltay = q8_deltay+k;
356     else q8_deltay = 0;
357     end
358     if (q8_deltaz > 0) q8_deltaz = q8_deltaz-k;
359     elseif (q8_deltaz < 0) q8_deltaz = q8_deltaz+k;
360     else q8_deltaz = 0;
361     end
362     nu_of_queues = 3;
363     if (sign(dir_q8x) > 0)
364         t_q = queue_x_neg;
365     elseif (sign(dir_q8x) < 0)
366         t_q = queue_x_pos;
367     else
368         t_q = 0;
369     end
370     nu_of_queues = nu_of_queues - 1;
371     end
372
373     if (sign(dir_q8y) > 0)
374         t_q = t_q + queue_y_neg;
375     elseif (sign(dir_q8y) < 0)
376         t_q = t_q + queue_y_pos;
377     else
378         nu_of_queues = nu_of_queues - 1;
379     end
380
381     if (sign(dir_q8z) > 0)
382         t_q = t_q + queue_z_neg;
383     elseif (sign(dir_q8x) < 0)
384         t_q = t_q + queue_z_pos;
385     else
386         nu_of_queues = nu_of_queues - 1;
387     end
388     t_q = t_q + 1;
389     if (nu_of_queues ~= 0) q8_val = (abs(q8_deltax) + abs(q8_deltay) + abs(q8_deltaz)) * (t_q/nu_of_queues);

```

```

389     else q8_val = Inf;
390     end
391
392     if (q8_val < min_val)
393         min_val = q8_val;
394         dir_x = sign(q8_deltax);
395         dir_y = sign(q8_deltay);
396         dir_z = sign(q8_deltaz);
397     end
398
399
400     else
401
402         % for quadrant I: (+ +)
403         [dir_q1x,q1_deltax] = dimord(sx, dx, k);
404         [dir_q1y,q1_deltay] = dimord(sy, dy, k);
405         nu_of_queues = 2;
406         if (sign(dir_q1x) < 0)
407             t_q = queue_x_neg;
408         elseif (sign(dir_q1x) > 0)
409             t_q = queue_x_pos;
410         else
411             t_q = 0;
412             nu_of_queues = nu_of_queues - 1;
413         end
414
415         if (sign(dir_q1y) < 0)
416             t_q = t_q + queue_y_neg;
417         elseif (sign(dir_q1y) > 0)
418             t_q = t_q + queue_y_pos;
419         else
420             nu_of_queues = nu_of_queues - 1;
421         end
422
423         t_q = t_q + 1;
424         if (nu_of_queues ~= 0) q1_val = (q1_deltax + q1_deltay) * (t_q/nu_of_queues);
425         else q1_val = Inf;
426         end
427
428         min_val = q1_val;
429         dir_x = sign(q1_deltax);
430         dir_y = sign(q1_deltay);
431
432
433
434         % for quadrant II: (+ -)
435         [dir_q2x,q2_deltax] = dimord(sx, dx, k);
436         [dir_q2y,q2_deltay] = dimord(sy, dy, k);
437         if (q2_deltay > 0) q2_deltay = q2_deltay - k;
438         elseif (q2_deltay < 0) q2_deltay = q2_deltay + k;
439         else q2_deltay = 0;
440         end
441         nu_of_queues = 2;
442         if (sign(dir_q2x) < 0)
443             t_q = queue_x_neg;
444         elseif (sign(dir_q2x) > 0)
445             t_q = queue_x_pos;
446         else
447             t_q = 0;
448             nu_of_queues = nu_of_queues - 1;
449         end
450
451         if (sign(dir_q2y) > 0)
452             t_q = t_q + queue_y_neg;
453         elseif (sign(dir_q2y) < 0)

```

```

454         t_q = t_q + queue_y_pos;
455     else
456         nu_of_queues = nu_of_queues - 1;
457     end
458
459     t_q = t_q + 1;
460     if (nu_of_queues ~= 0) q2_val = (abs(q2_deltax) + abs(q2_deltay)) * (t_q/nu_of_queues);
461     else q2_val = Inf;
462     end
463
464     if (q2_val < min_val)
465         min_val = q2_val;
466         dir_x = sign(q2_deltax);
467         dir_y = sign(q2_deltay);
468     end
469
470
471     % for quadrant III: (- +)
472     [dir_q3x,q3_deltax] = dimord(sx, dx, k);
473     [dir_q3y,q3_deltay] = dimord(sy, dy, k);
474     if (q3_deltax > 0) q3_deltax = q3_deltax - k;
475     elseif (q3_deltax < 0) q3_deltax = q3_deltax + k;
476     else q3_deltax = 0;
477     end
478     nu_of_queues = 2;
479     if (sign(dir_q3x) > 0)
480         t_q = queue_x_neg;
481     elseif (sign(dir_q3x) < 0)
482         t_q = queue_x_pos;
483     else
484         t_q = 0;
485         nu_of_queues = nu_of_queues - 1;
486     end
487
488     if (sign(dir_q3y) < 0)
489         t_q = t_q + queue_y_neg;
490     elseif (sign(dir_q3y) > 0)
491         t_q = t_q + queue_y_pos;
492     else
493         nu_of_queues = nu_of_queues - 1;
494     end
495
496     t_q = t_q + 1;
497     if (nu_of_queues ~= 0) q3_val = (abs(q3_deltax) + abs(q3_deltay)) * (t_q/nu_of_queues);
498     else q3_val = Inf;
499     end
500
501     if (q3_val < min_val)
502         min_val = q3_val;
503         dir_x = sign(q3_deltax);
504         dir_y = sign(q3_deltay);
505     end
506
507
508     % for quadrant IV: (- -)
509     [dir_q4x,q4_deltax] = dimord(sx, dx, k);
510     [dir_q4y,q4_deltay] = dimord(sy, dy, k);
511     if (q4_deltay > 0) q4_deltay = q4_deltay - k;
512     elseif (q4_deltay < 0) q4_deltay = q4_deltay + k;
513     else q4_deltay = 0;
514     end
515     if (q4_deltax > 0) q4_deltax = q4_deltax - k;
516     elseif (q4_deltax < 0) q4_deltax = q4_deltax + k;
517     else q4_deltax = 0;
518     end

```

```

519     nu_of_queues = 2;
520     if (sign(dir_q4x) > 0)
521         t_q = queue_x_neg;
522     elseif (sign(dir_q4x) < 0)
523         t_q = queue_x_pos;
524     else
525         t_q = 0;
526         nu_of_queues = nu_of_queues - 1;
527     end
528
529     if (sign(dir_q4y) > 0)
530         t_q = t_q + queue_y_neg;
531     elseif (sign(dir_q4y) < 0)
532         t_q = t_q + queue_y_pos;
533     else
534         nu_of_queues = nu_of_queues - 1;
535     end
536
537     t_q = t_q + 1;
538     if (nu_of_queues ~= 0) q4_val = (abs(q4_deltax) + abs(q4_deltay)) * (t_q/nu_of_queues);
539     else q4_val = Inf;
540     end
541
542     if (q4_val < min_val)
543         min_val = q4_val;
544         dir_x = sign(q4_deltax);
545         dir_y = sign(q4_deltay);
546     end
547
548     dir_z = 0;
549
550 end
551
552 newqueues = queues;
553
554 if (sign(dir_x) > 0)
555     newqueues(src,1) = queue_x_pos + 1;
556 elseif (sign(dir_x) < 0)
557     newqueues(src,2) = queue_x_neg + 1;
558 end
559
560 if (sign(dir_y) > 0)
561     newqueues(src,3) = queue_y_pos + 1;
562 elseif (sign(dir_y) < 0)
563     newqueues(src,4) = queue_y_neg + 1;
564 end
565
566 if (sign(dir_z) > 0)
567     newqueues(src,5) = queue_z_pos + 1;
568 elseif (sign(dir_z) < 0)
569     newqueues(src,6) = queue_z_neg + 1;
570 end
571

```

## C.2.11 file: cqr\_getRandomValue.m

```
1 function [src_,dest_,v] = cqr_getRandomValue(traff)
2 % use: [src,dest,v] = cqr_getRandomValue(traffic_matrix);
3
4     rand_value = ceil(rand()*length(find(traff>=1)));
5     [rows,cols,vals] = find(traff>=1);
6     src_ = rows(rand_value);
7     dest_ = cols(rand_value);
8     v = ceil(vals(rand_value)*rand());
```

## C.2.12 file: cqr\_transformTrafficMatrix.m

```
1 function [packet_matrix_] = cqr_transformTrafficMatrix (traf_initial)
2     %use: [packet_matrix] = cqr_transformTrafficMatrix (traffic_matrix);
3
4     traf = traf_initial;
5     size = length(traf);
6     for i=1:size
7         for j=1:size
8             for k = 1:max(max(traf))
9                 packet_matrix_(i,j,k).initialized = 0;
10                packet_matrix_(i,j,k).valid = 0;
11                packet_matrix_(i,j,k).x_dir = 0;
12                packet_matrix_(i,j,k).y_dir = 0;
13                packet_matrix_(i,j,k).z_dir = 0;
14            end
15        end
16    end
17
18    while (length(find(traf>=1))~=0)
19        [src,dest,v] = cqr_getRandomValue(traf);
20        for i=1:size
21            if (packet_matrix_(src,dest,i).valid == 0)
22                value = i;
23                break;
24            end
25        end
26        packet_matrix_(src,dest,value).valid = 1;
27        traf(src,dest) = traf(src,dest) - 1;
28    end
```

# Appendix D

## Appendix D: Full Laboratory Results

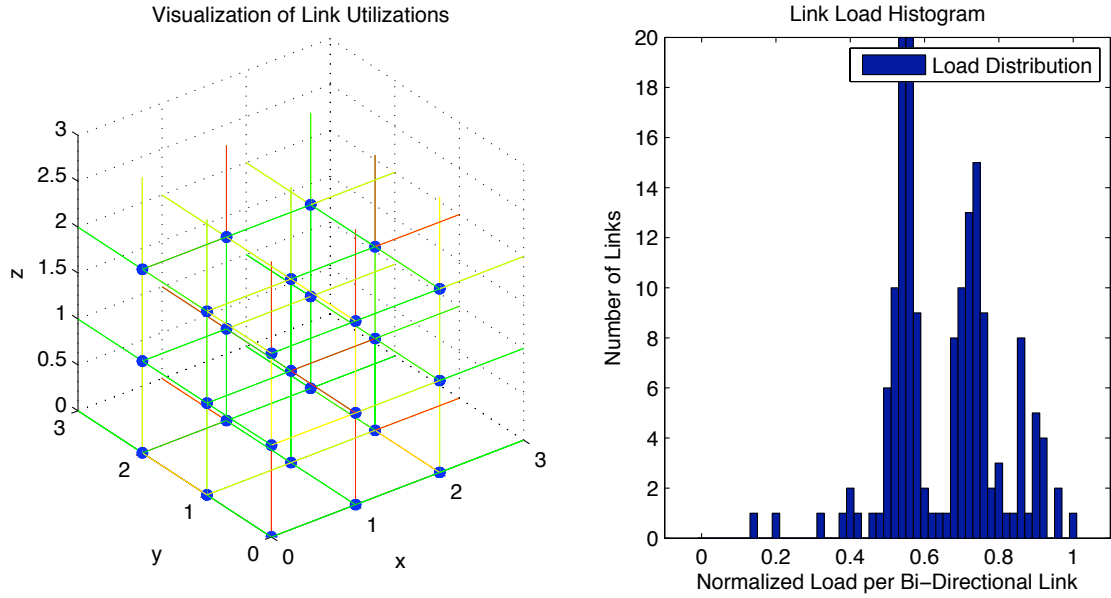
### D.1 Matlab Plot Scripts

#### D.1.1 file: convertLoad.m

```
1 function convertLoad(k,n,matrix);
2 % usage: convertLoad(k,n,load_matrix)
3 %     where, load_matrix is the load distribution taken directly
4 %     from the 'interconnectThroughput' script.
5
6 load_matrix = zeros(k^n);
7
8 for (i=1:(k^n))
9     node = 28-i;
10    negy = matrix(node,1);
11    posy = matrix(node,2);
12    posx = matrix(node,3);
13    negx = matrix(node,4);
14    negz = matrix(node,5);
15    posz = matrix(node,6);
16    [x,y,z] = i2dim((i-1),k,n);
17    node_posx = dim2i(mod(x+1,k),y,z,k,n) + 1;
18    node_negx = dim2i(mod(x-1,k),y,z,k,n) + 1;
19    node_posy = dim2i(x,mod(y+1,k),z,k,n) + 1;
20    node_negy = dim2i(x,mod(y-1,k),z,k,n) + 1;
21    node_posz = dim2i(x,y,mod(z+1,k),k,n) + 1;
22    node_negz = dim2i(x,y,mod(z-1,k),k,n) + 1;
23    load_matrix(i,node_posx) = posx;
24    load_matrix(i,node_negx) = negx;
25    load_matrix(i,node_posy) = posy;
26    load_matrix(i,node_negy) = negy;
27    load_matrix(i,node_posz) = posz;
28    load_matrix(i,node_negz) = negz;
29 end
30
31 max_ch_load = max(max(load_matrix));
32 load_matrix = load_matrix / max_ch_load;
33 std_dev = mean(std(load_matrix));
34 subplot(1,2,1);
35 view_load(load_matrix, k, n, 1.0-(std_dev), 1.0-(std_dev*0.5));
36
37 count = 0;
38
39 while (count < (k^n))
40     [val_x,val_y,val_z] = i2dim(count,k,n);
```

```
41     [val1_x,val1_y,val1_z] = sphere(30);
42     val1_x = val1_x.*0.05; val1_y = val1_y.*0.05; val1_z = val1_z.*0.05;
43     val1_x = val1_x + val_x; val1_y = val1_y + val_y; val1_z = val1_z + val_z;
44     plot3(val1_x,val1_y,val1_z,'b');
45     count = count + 1;
46 end
47 set(gcf,'Position', [50,500,650,300])
```





**Figure D.1:** Results of Dimension Ordered Routing (DOR) using the nearest neighbor traffic demand (NN). Each node had two individual traffic demands to each of its six neighbors.

### D.1.2 file: view\_load.m

See Appendix C.1.2.

### D.1.3 file: i2dim.m

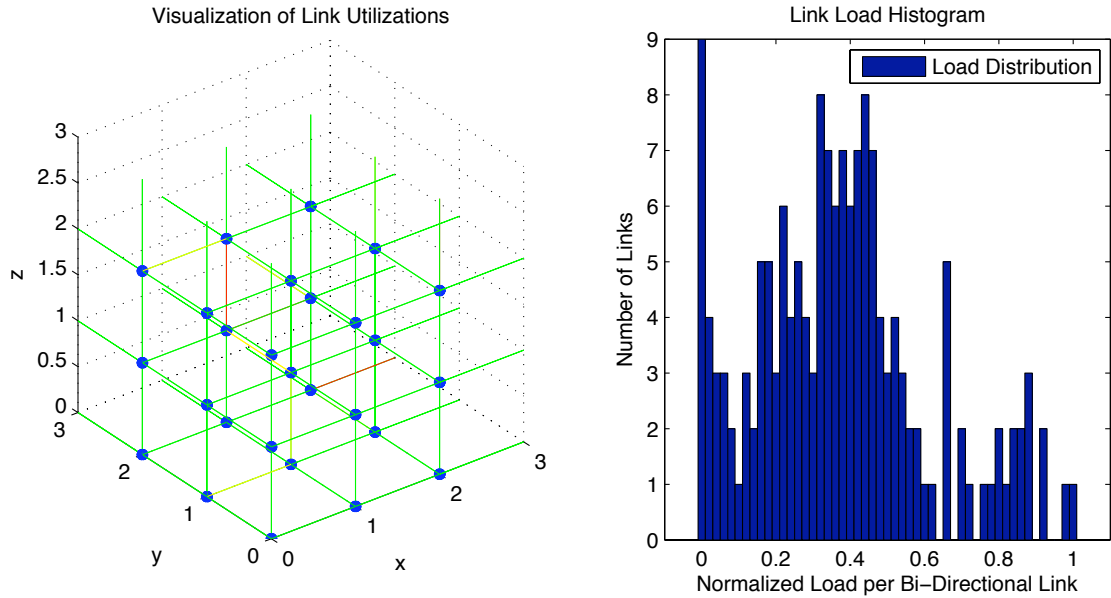
See Appendix C.2.2.

### D.1.4 file: dim2i.m

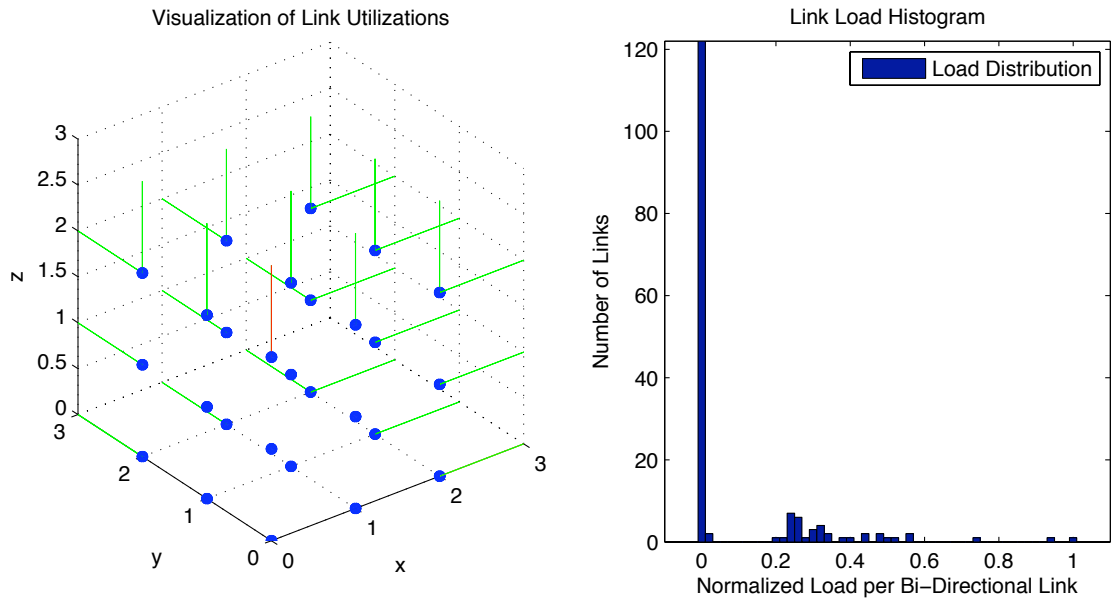
See Appendix C.2.1.

## D.2 Static Algorithms

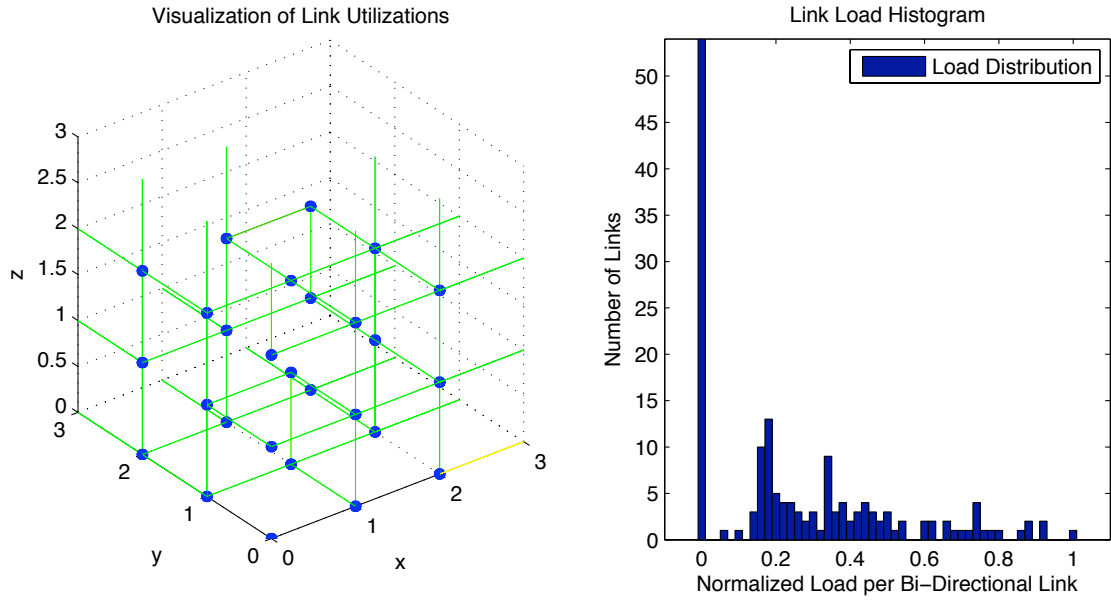
### D.2.1 Dimension Ordered Routing Results



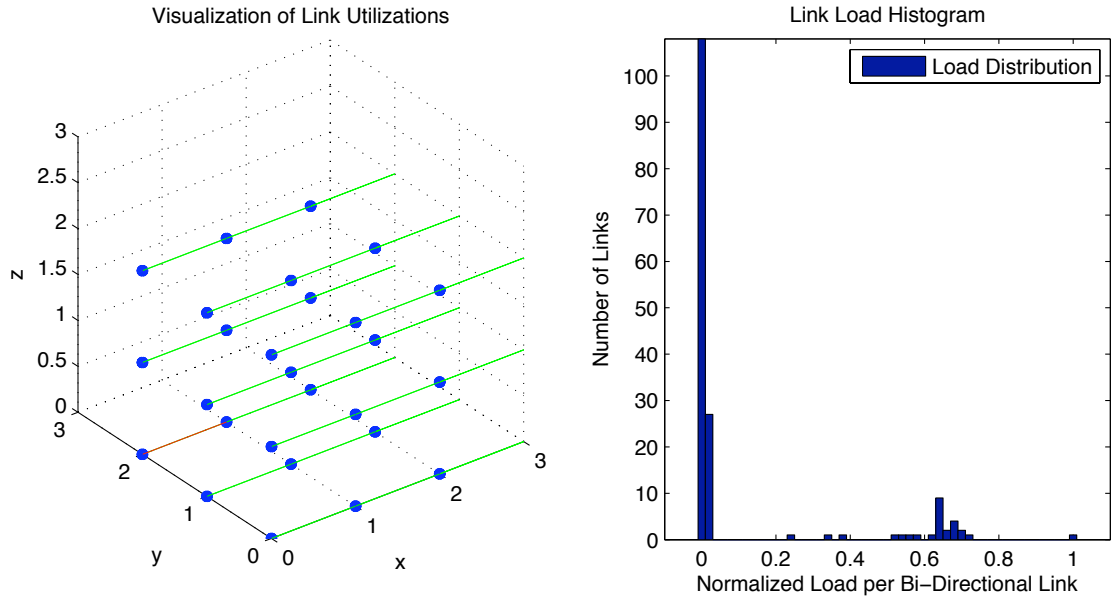
**Figure D.2:** Results of Dimension Ordered Routing (DOR) using a uniform random traffic pattern (UR). For each node, a value of 0-9 nodes was assigned as having 0-9 individual traffic demands.



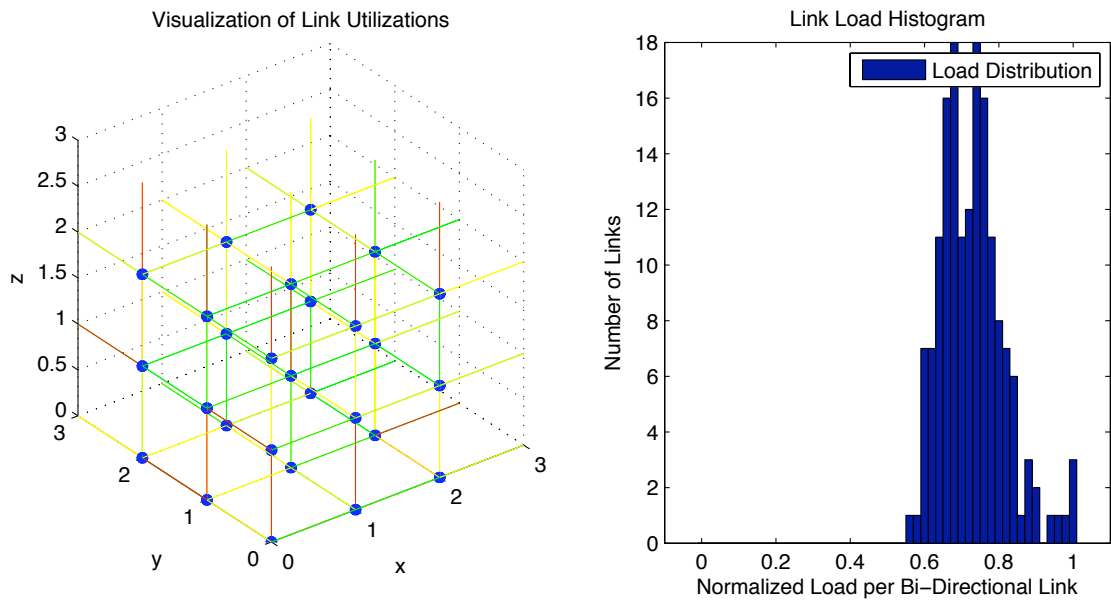
**Figure D.3:** Results of Dimension Ordered Routing (DOR) using a bit complement traffic pattern (BC). Two individual traffic demands were assigned for possible values described in Table 2.1.



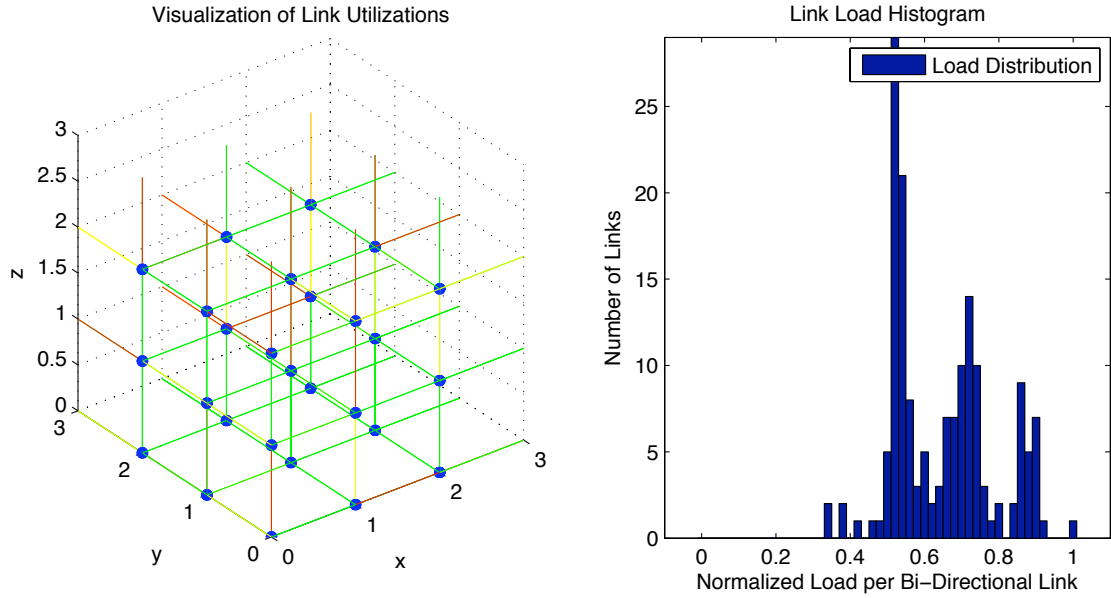
**Figure D.4:** Results of Dimension Ordered Routing (DOR) using a transpose traffic pattern (TP). Two individual traffic demands were assigned for each permutation of all  $x, y, z$  values.



**Figure D.5:** Results of Dimension Ordered Routing (DOR) using a tornado traffic pattern (TOR). Three individual traffic demands were assigned according to Table 2.1.

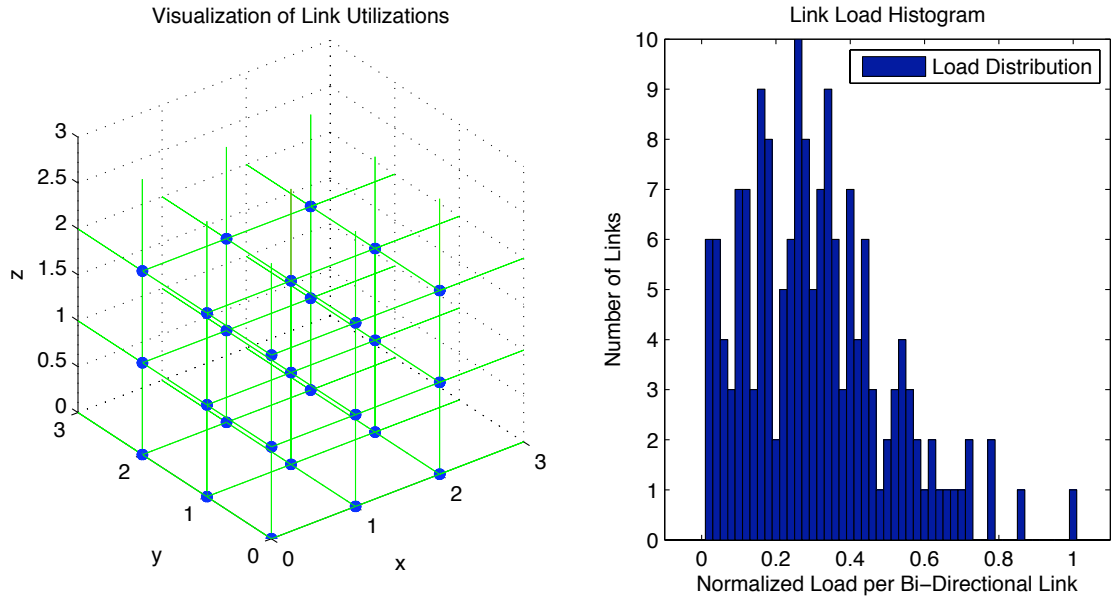


**Figure D.6:** Results of Dimension Ordered Routing (DOR) using a flood traffic pattern (FL). Two individual traffic demands were assigned to each node from each node.

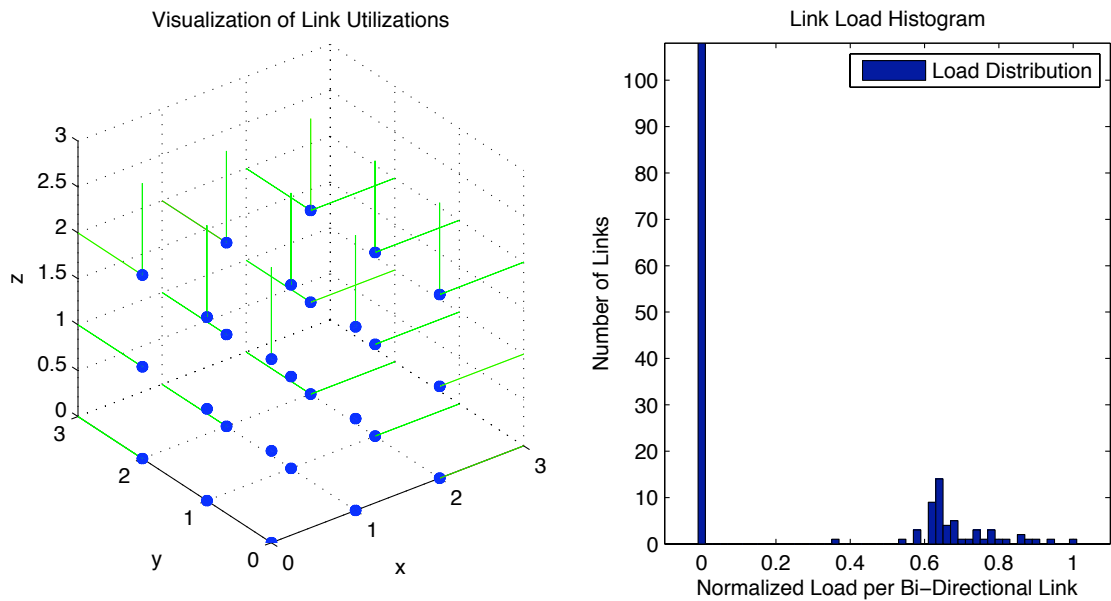


**Figure D.7:** Results of Direction Ordered Routing (DIR) using the nearest neighbor traffic demand (NN). Each node had two individual traffic demands to each of its six neighbors.

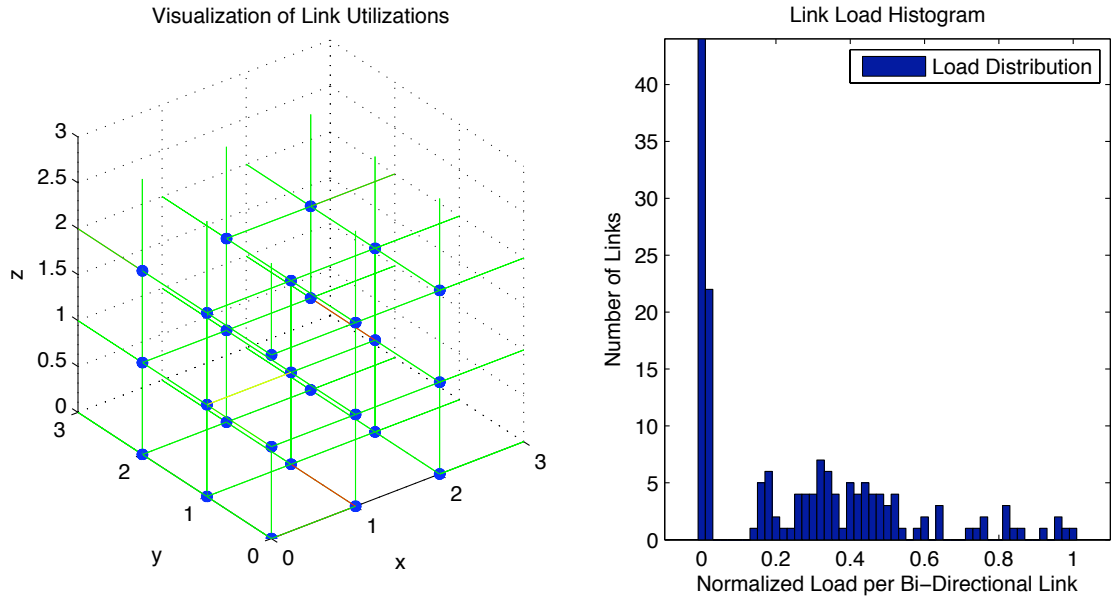
## D.2.2 Direction Ordered Routing Results



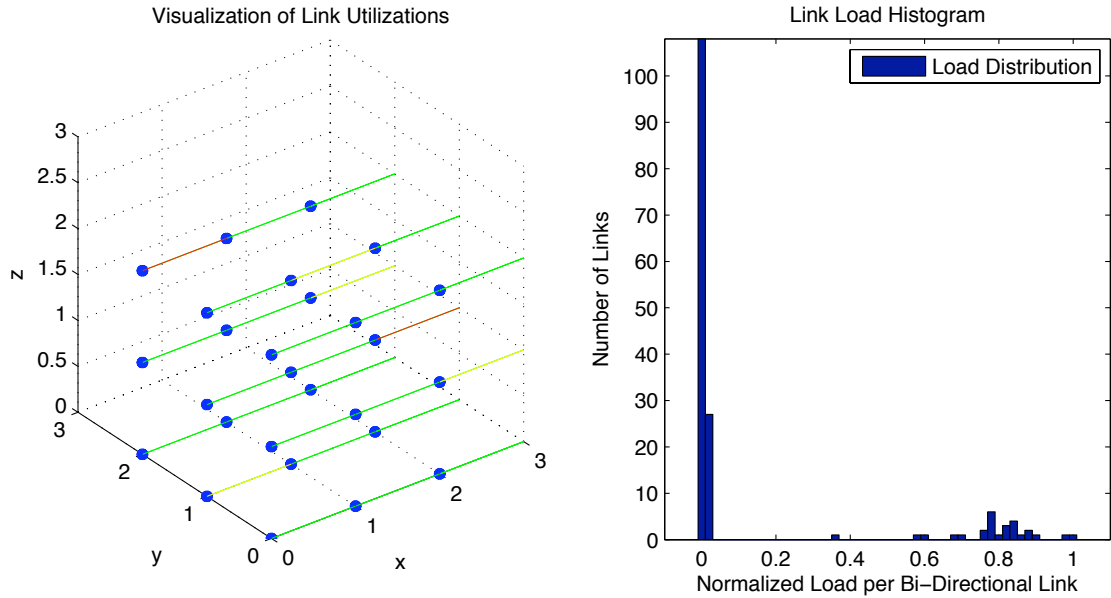
**Figure D.8:** Results of Direction Ordered Routing (DIR) using a uniform random traffic pattern (UR). For each node, a value of 0-9 nodes was assigned as having 0-9 individual traffic demands.



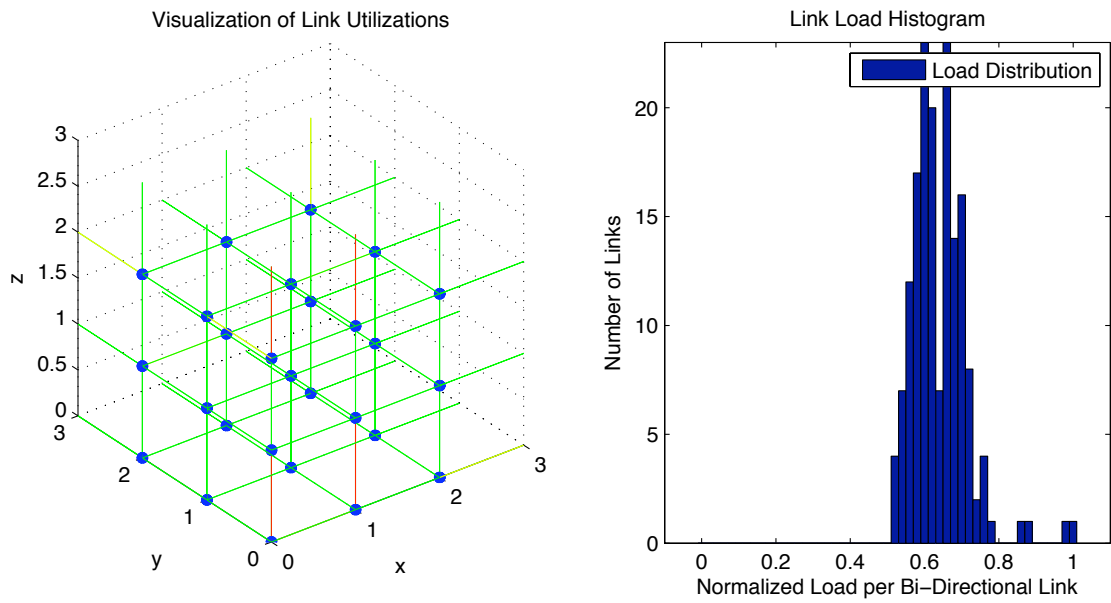
**Figure D.9:** Results of Direction Ordered Routing (DIR) using a bit complement traffic pattern (BC). Two individual traffic demands were assigned for possible values described in Table 2.1.



**Figure D.10:** Results of Direction Ordered Routing (DIR) using a transpose traffic pattern (TP). Two individual traffic demands were assigned for each permutation of all  $x, y, z$  values.

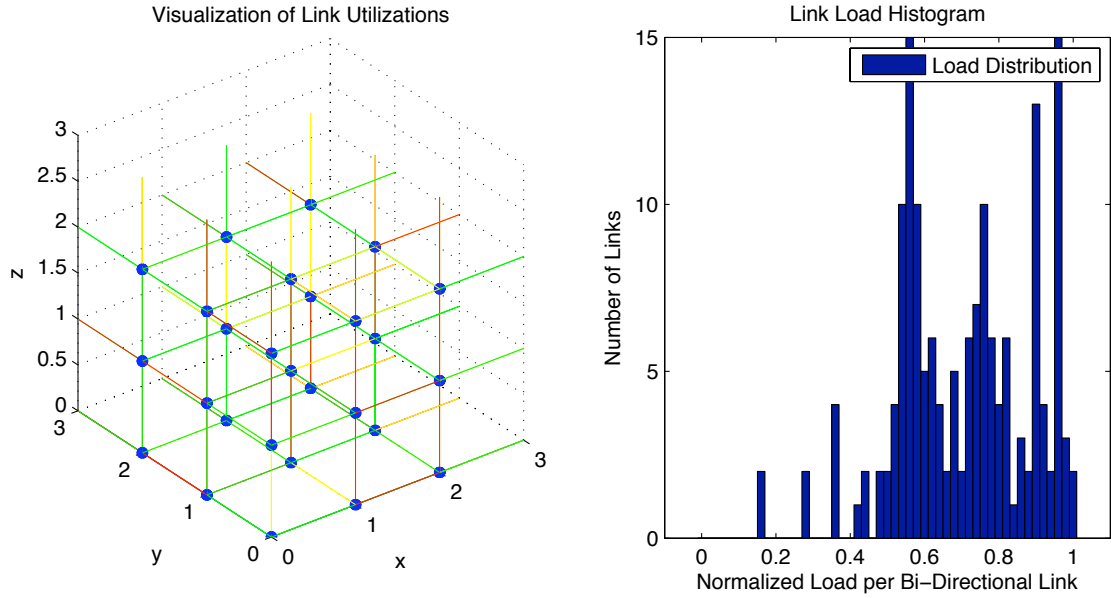


**Figure D.11:** Results of Direction Ordered Routing (DIR) using a tornado traffic pattern (TOR). Three individual traffic demands were assigned according to Table 2.1.



**Figure D.12:** Results of Direction Ordered Routing (DIR) using a flood traffic pattern (FL). Two individual traffic demands were assigned to each node from each node.

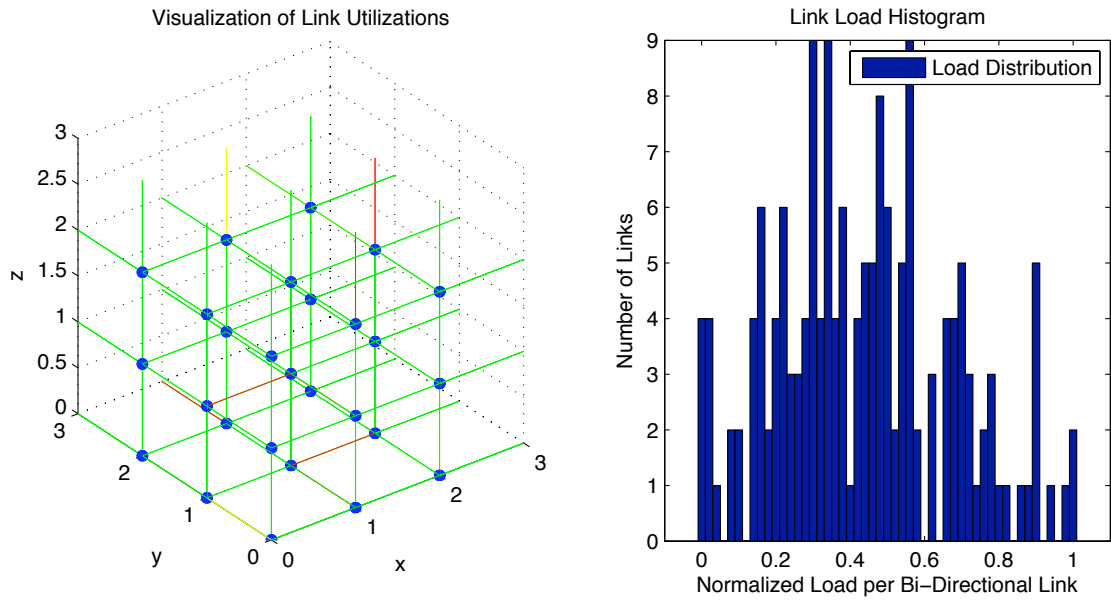




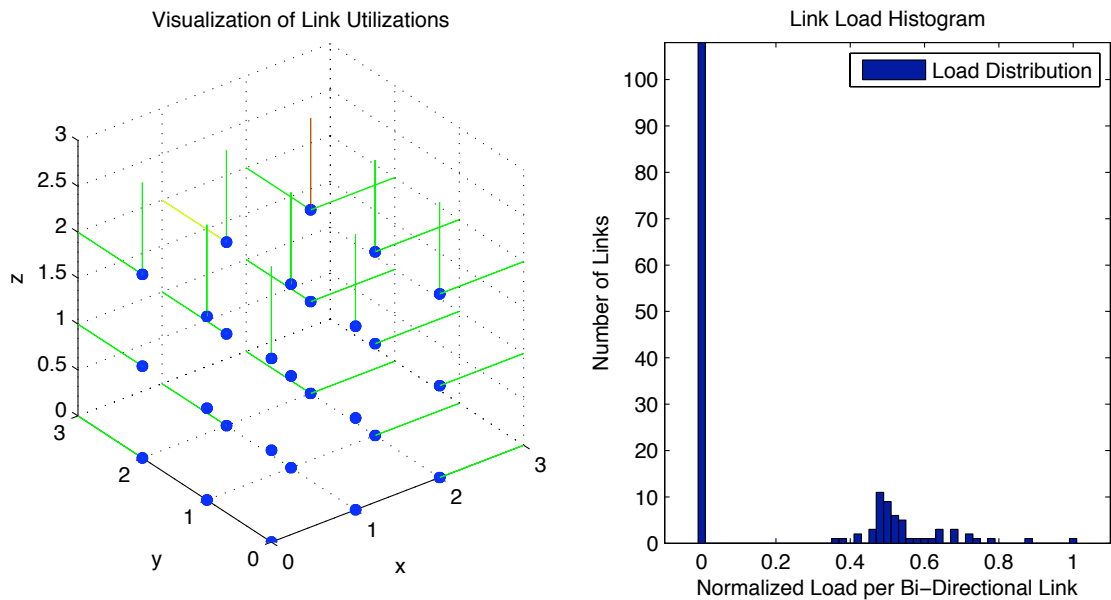
**Figure D.13:** Results of Minimal Oblivious Routing using the nearest neighbor traffic demand (NN). Each node had two individual traffic demands to each of its six neighbors.

## D.3 Oblivious Algorithms

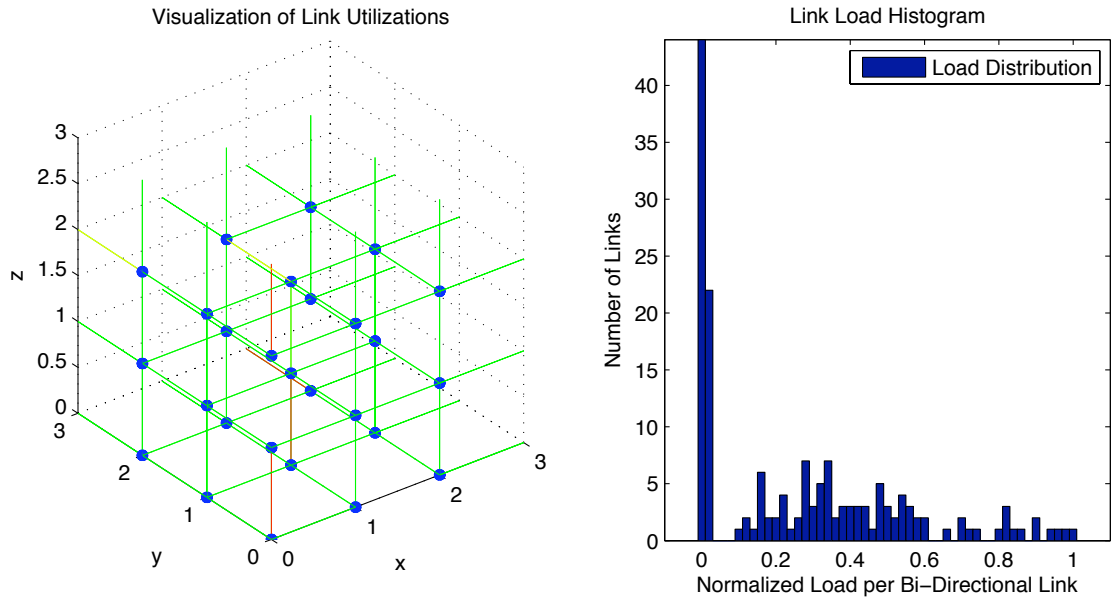
### D.3.1 Minimal Oblivious Routing Results



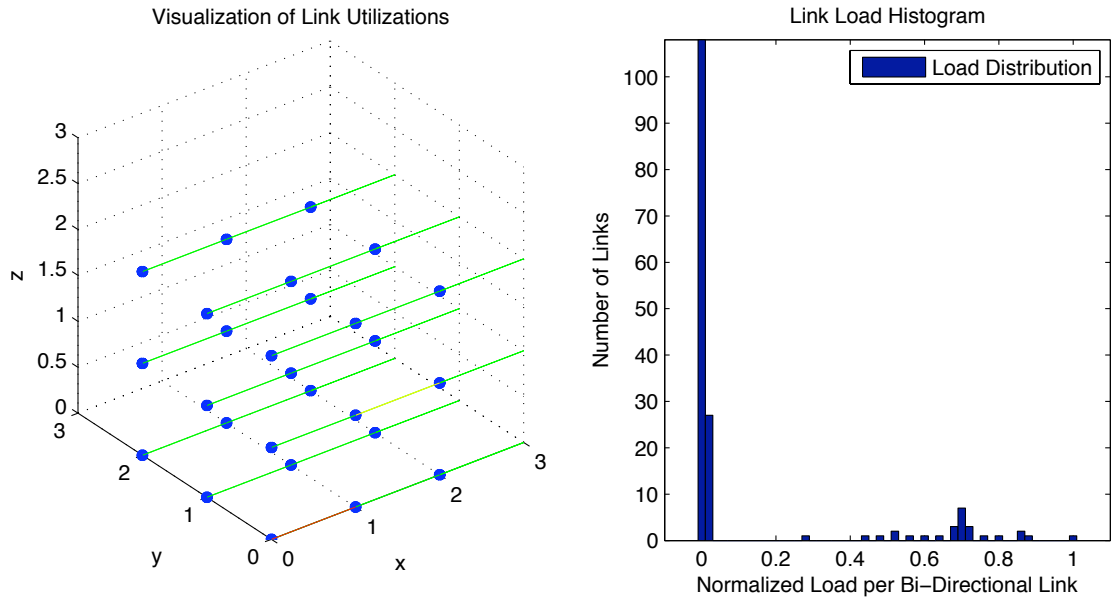
**Figure D.14:** Results of Minimal Oblivious Routing using a uniform random traffic pattern (UR). For each node, a value of 0-9 nodes was assigned as having 0-9 individual traffic demands.



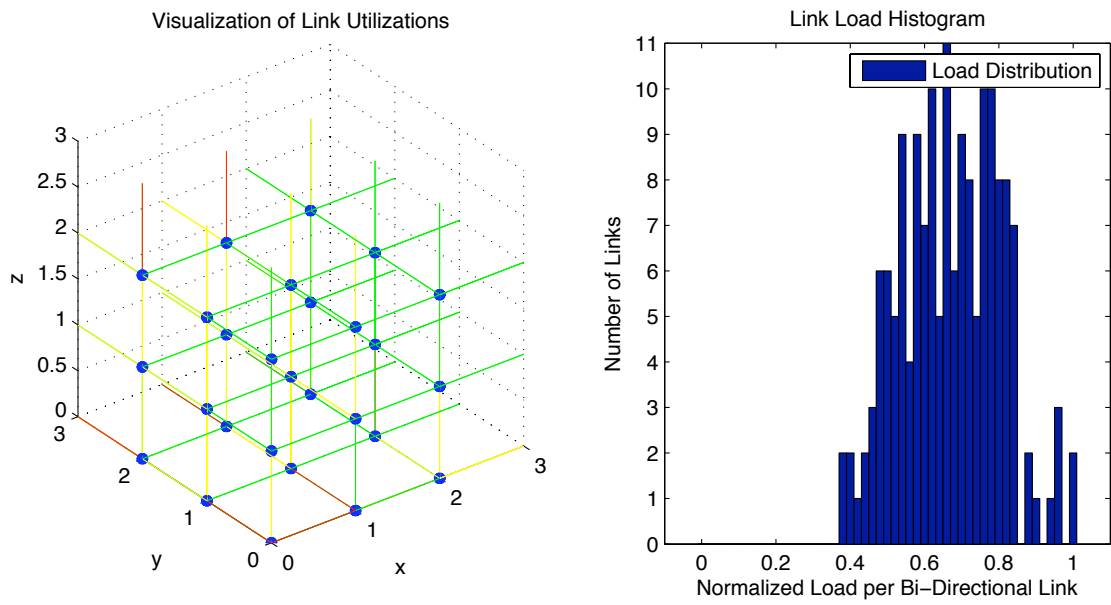
**Figure D.15:** Results of Minimal Oblivious Routing using a bit complement traffic pattern (BC). Two individual traffic demands were assigned for possible values described in Table 2.1.



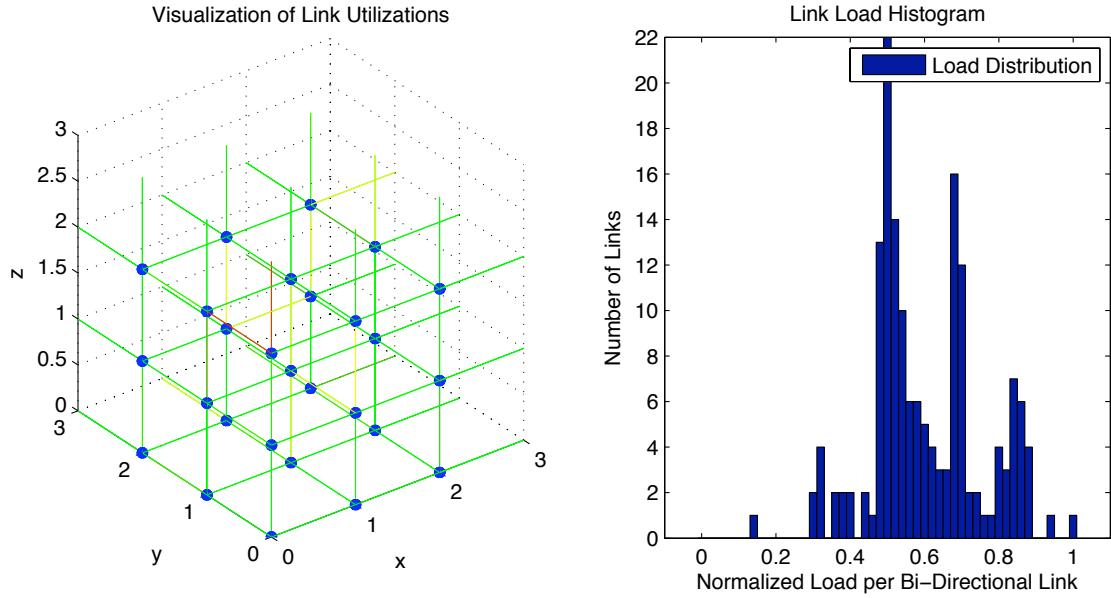
**Figure D.16:** Results of Minimal Oblivious Routing using a transpose traffic pattern (TP). Two individual traffic demands were assigned for each permutation of all  $x, y, z$  values.



**Figure D.17:** Results of Minimal Oblivious Routing using a tornado traffic pattern (TOR). Three individual traffic demands were assigned according to Table 2.1.



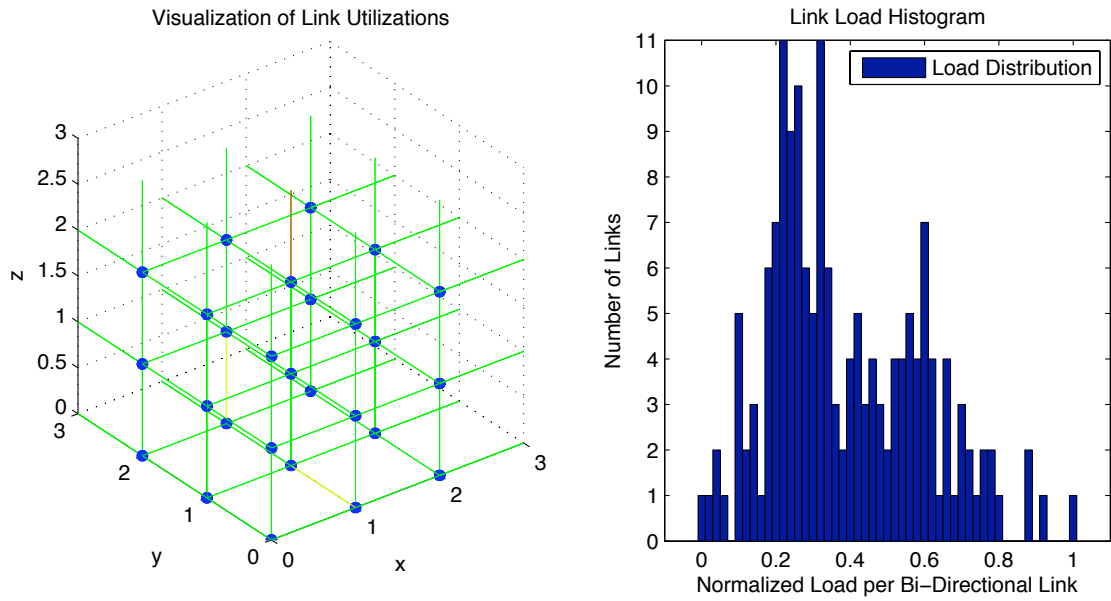
**Figure D.18:** Results of Minimal Oblivious Routing using a flood traffic pattern (FL). Two individual traffic demands were assigned to each node from each node.



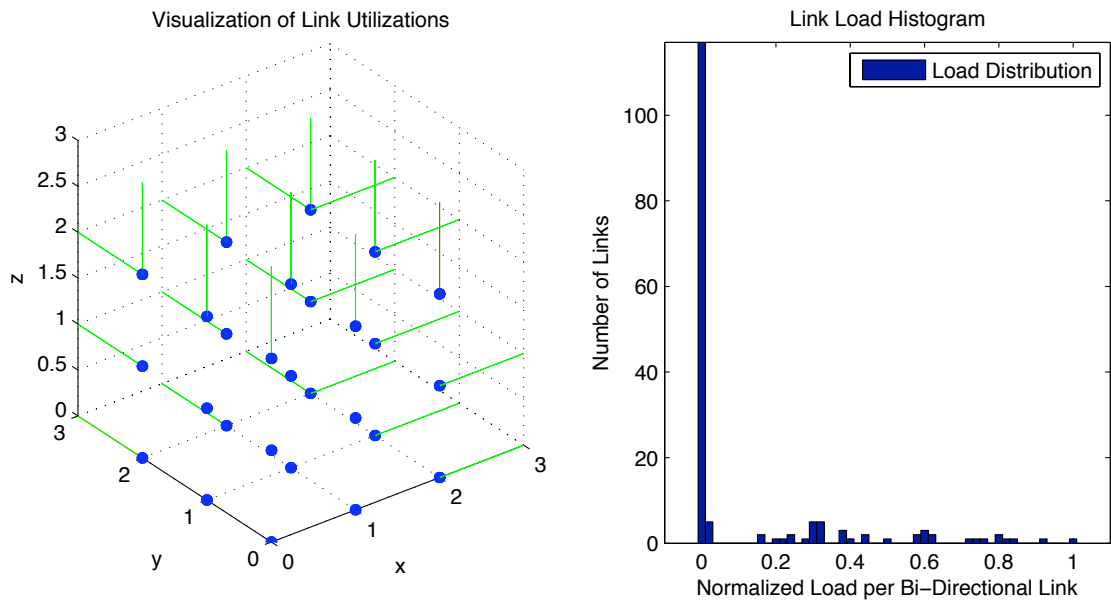
**Figure D.19:** Results of Minimal Adaptive Routing using the nearest neighbor traffic demand (NN). Each node had two individual traffic demands to each of its six neighbors.

## D.4 Adaptive Algorithms

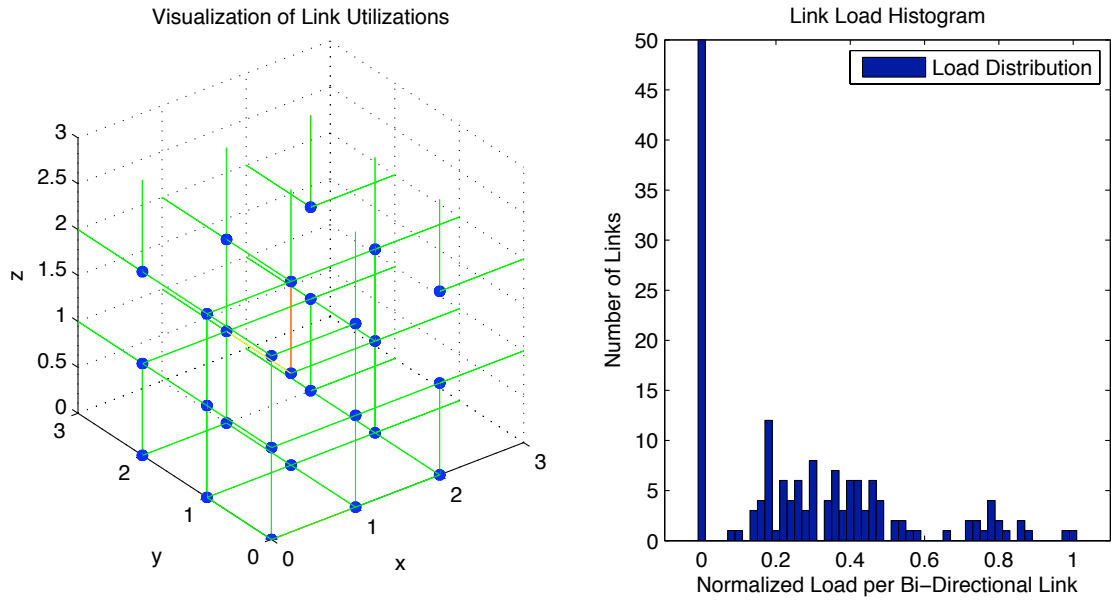
### D.4.1 Minimal Adaptive Routing Results



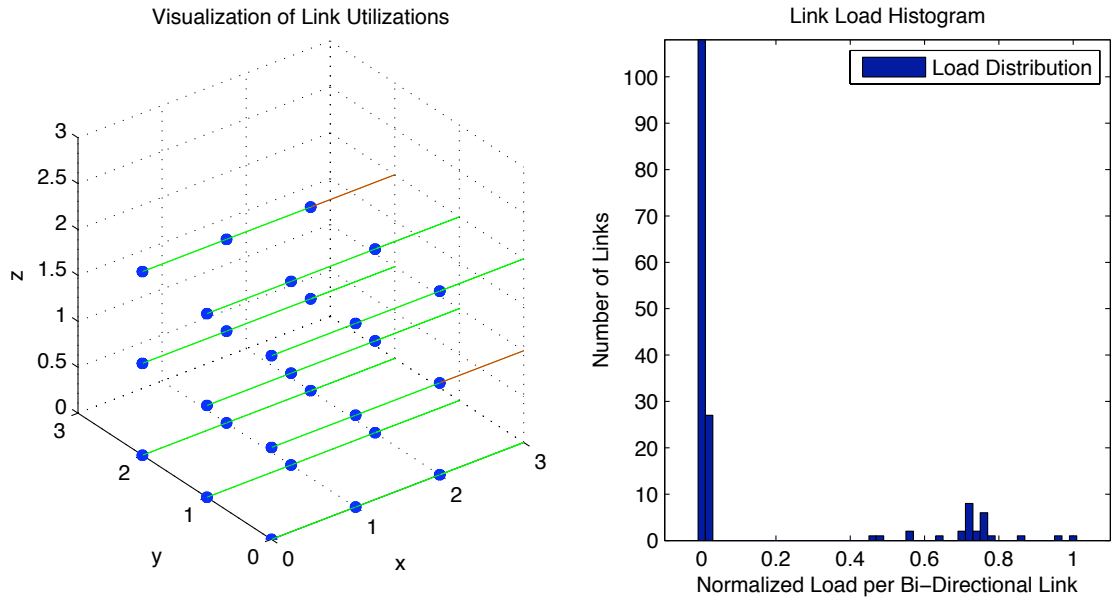
**Figure D.20:** Results of Minimal Adaptive Routing using a uniform random traffic pattern (UR). For each node, a value of 0-9 nodes was assigned as having 0-9 individual traffic demands.



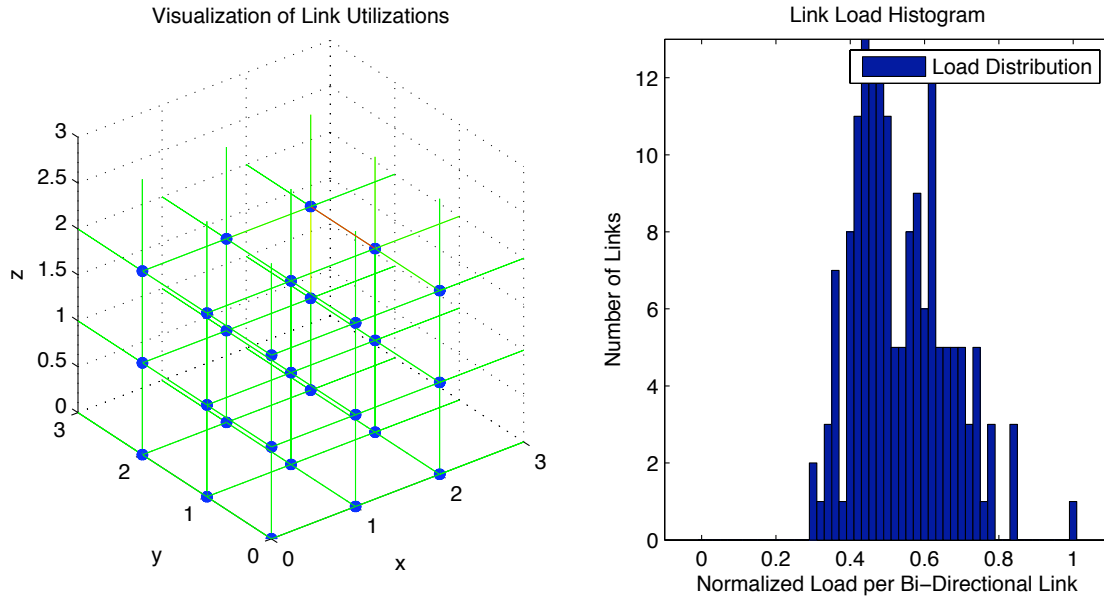
**Figure D.21:** Results of Minimal Adaptive Routing using a bit complement traffic pattern (BC). Two individual traffic demands were assigned for possible values described in Table 2.1.



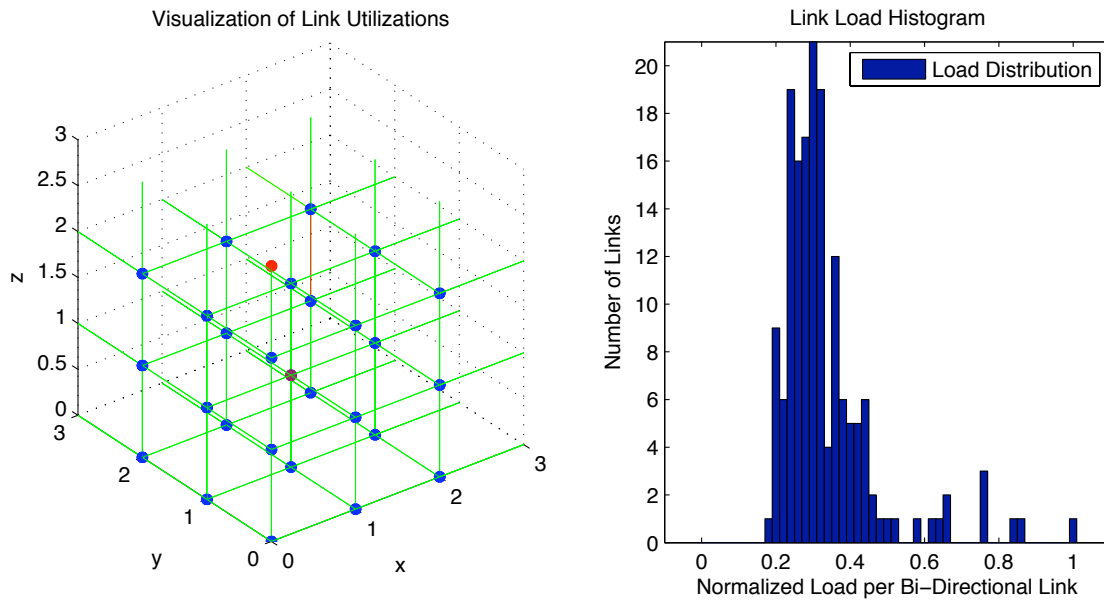
**Figure D.22:** Results of Minimal Adaptive Routing using a transpose traffic pattern (TP). Two individual traffic demands were assigned for each permutation of all  $x, y, z$  values.



**Figure D.23:** Results of Minimal Adaptive Routing using a tornado traffic pattern (TOR). Three individual traffic demands were assigned according to Table 2.1.

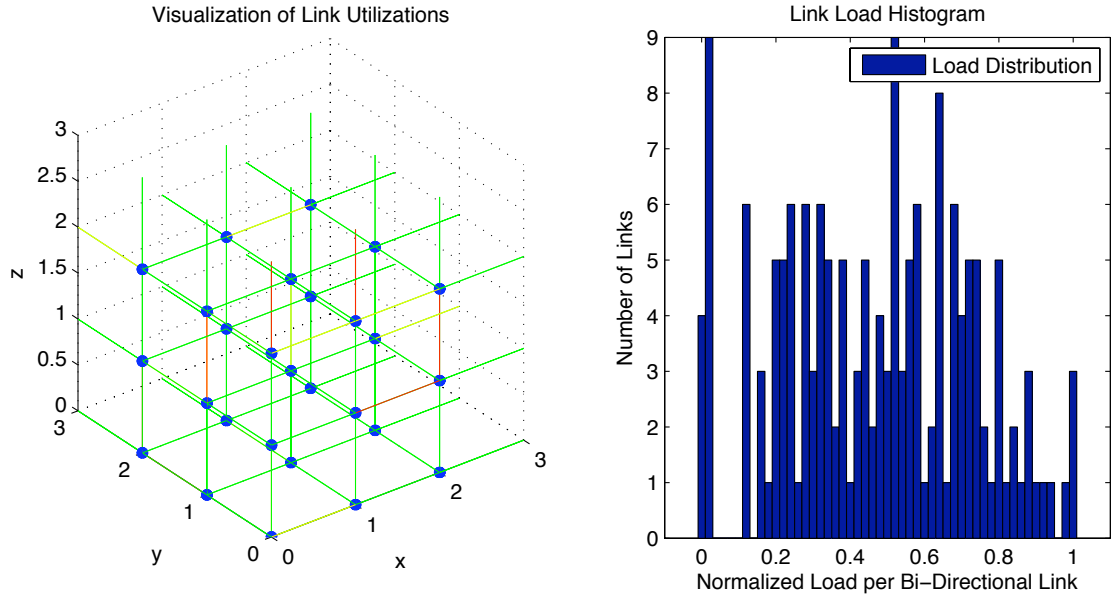


**Figure D.24:** Results of Minimal Adaptive Routing using a flood traffic pattern (FL). Two individual traffic demands were assigned to each node from each node.



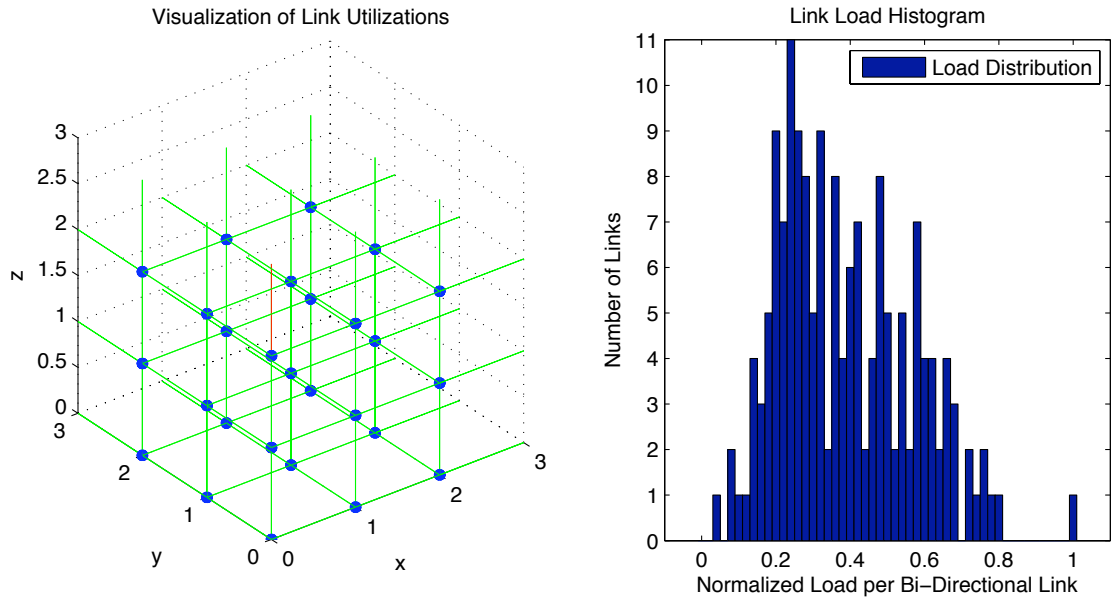
**Figure D.25:** Results of Minimal Adaptive Routing using a flood traffic pattern (FL) with Hot-spots. Two individual traffic demands were assigned to each node from each node. Five percent of the nodes were made hot-spots, and they are marked as red in the figure.



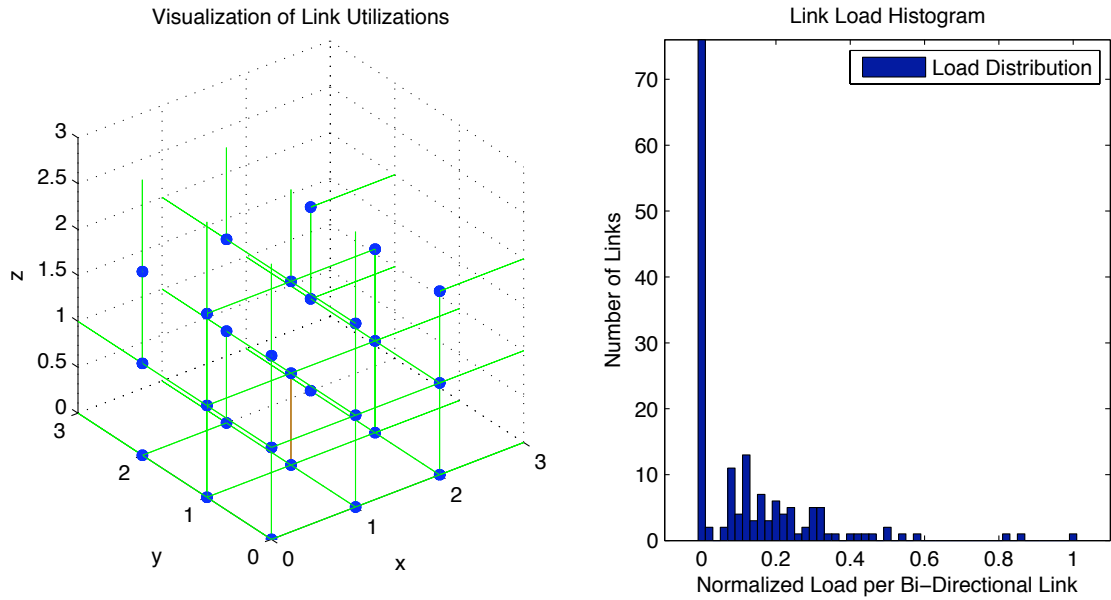


**Figure D.26:** Results of CQR Routing using the nearest neighbor traffic demand (NN). Each node had two individual traffic demands to each of its six neighbors.

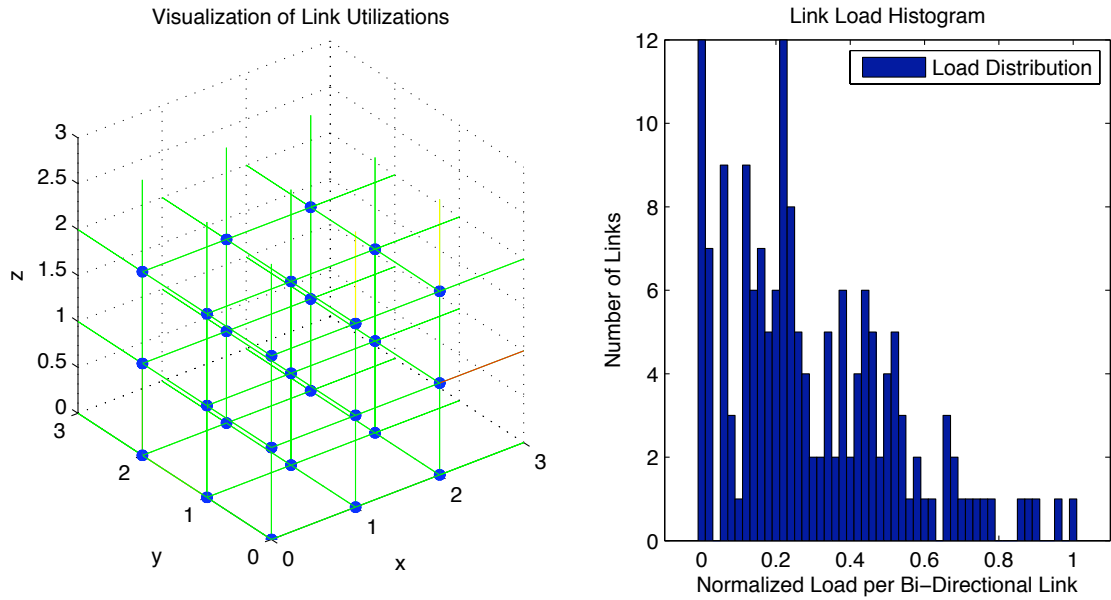
#### D.4.2 CQR Routing Results



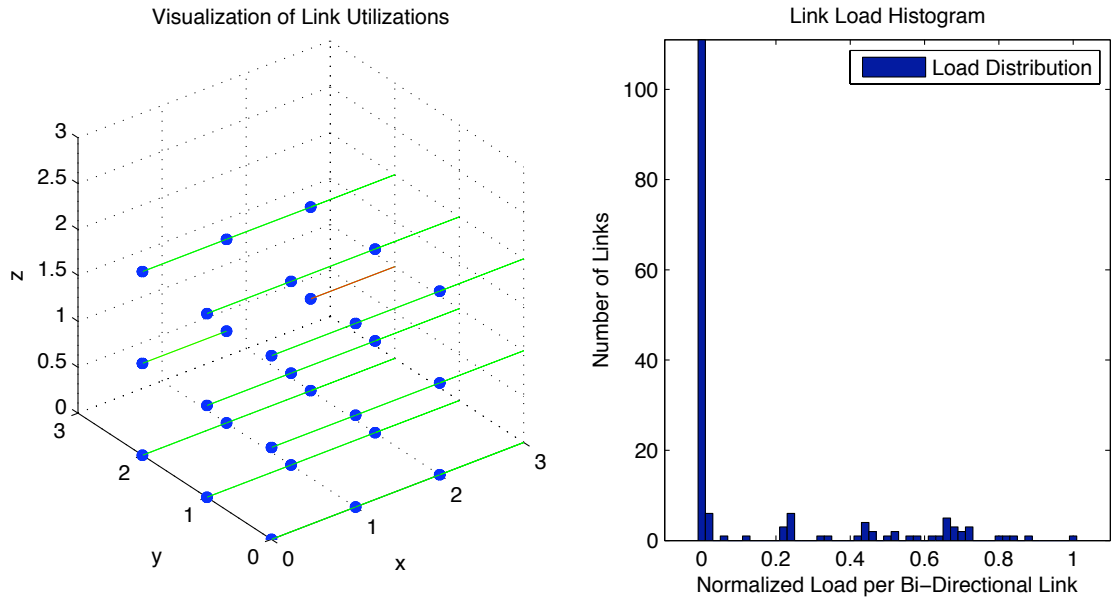
**Figure D.27:** Results of CQR Routing using a uniform random traffic pattern (UR). For each node, a value of 0-9 nodes was assigned as having 0-9 individual traffic demands.



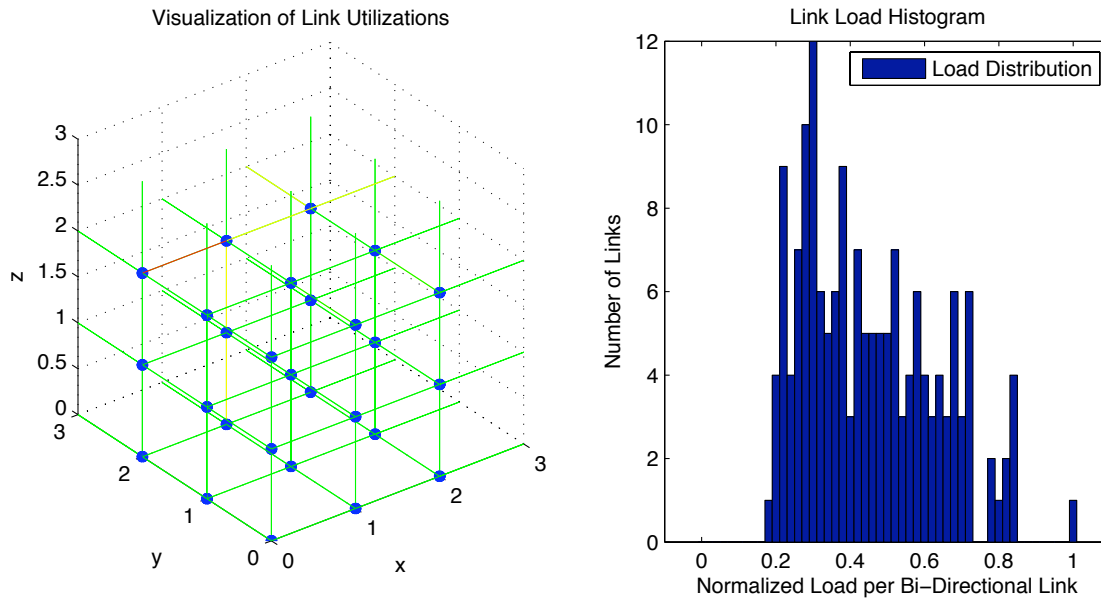
**Figure D.28:** Results of CQR Routing using a bit complement traffic pattern (BC). Two individual traffic demands were assigned for possible values described in Table 2.1.



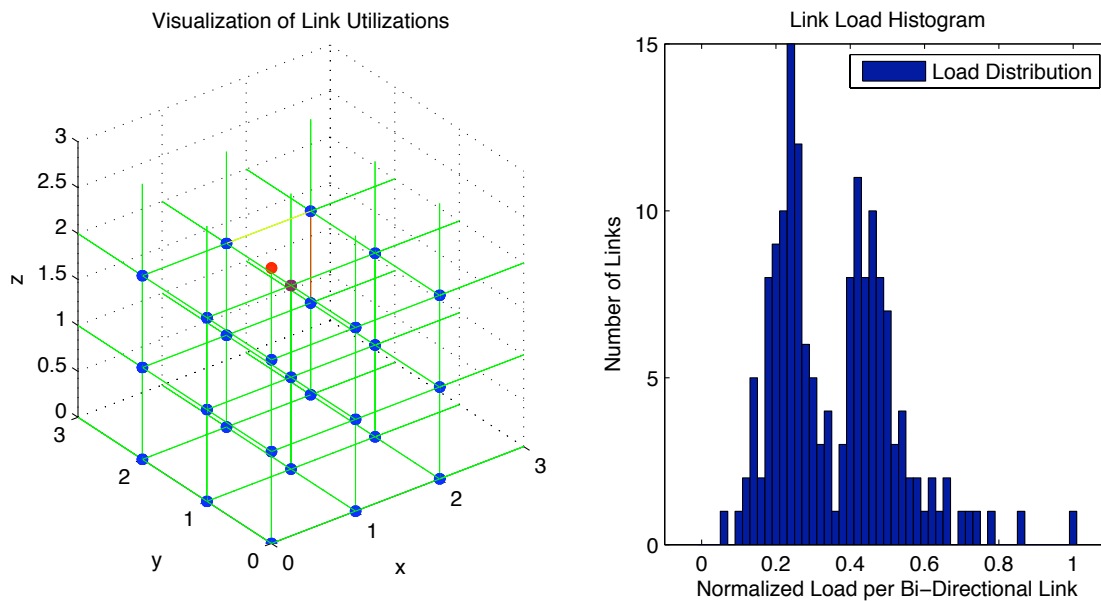
**Figure D.29:** Results of CQR Routing using a transpose traffic pattern (TP). Two individual traffic demands were assigned for each permutation of all  $x, y, z$  values.



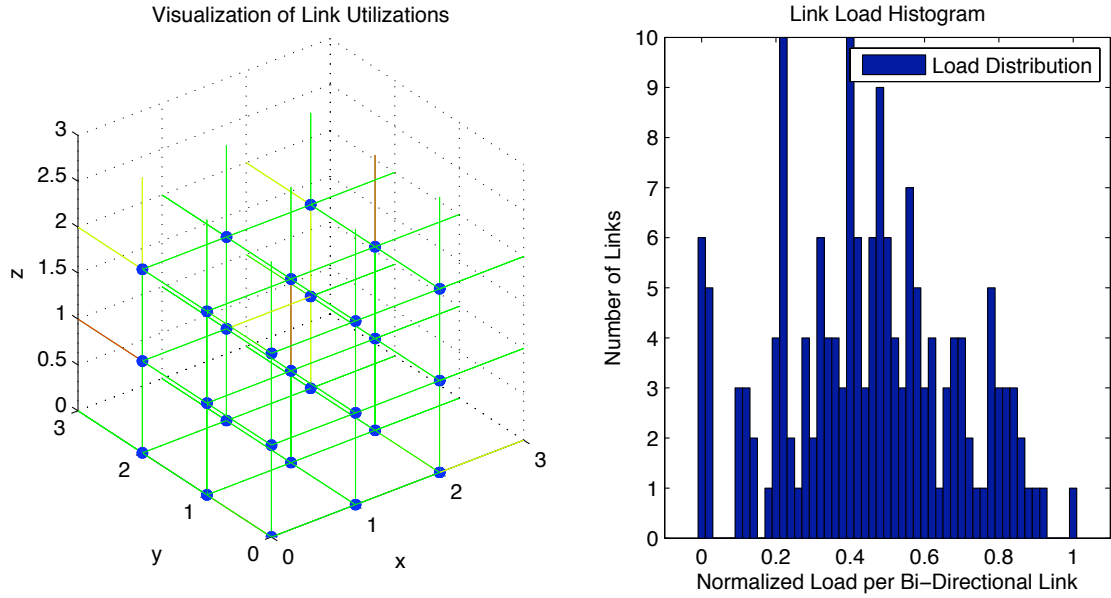
**Figure D.30:** Results of CQR Routing using a tornado traffic pattern (TOR). Three individual traffic demands were assigned according to Table 2.1.



**Figure D.31:** Results of CQR Routing using a flood traffic pattern (FL). Two individual traffic demands were assigned to each node from each node.

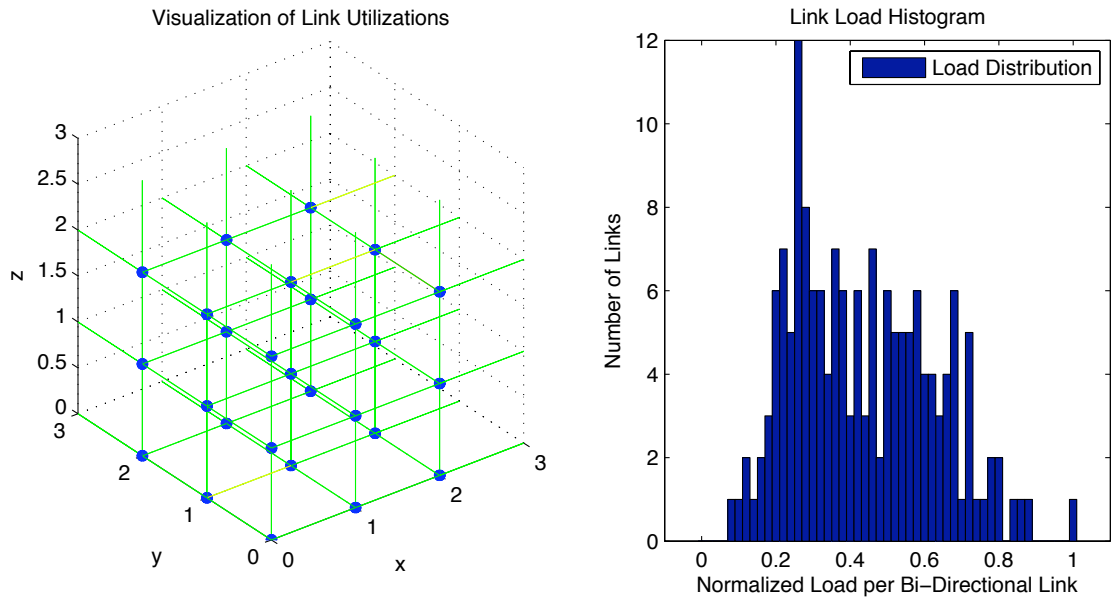


**Figure D.32:** Results of CQR Routing using a flood traffic pattern (FL) with Hot-spots. Two individual traffic demands were assigned to each node from each node. Five percent of the nodes were made hot-spots, and they are marked as red in the figure.

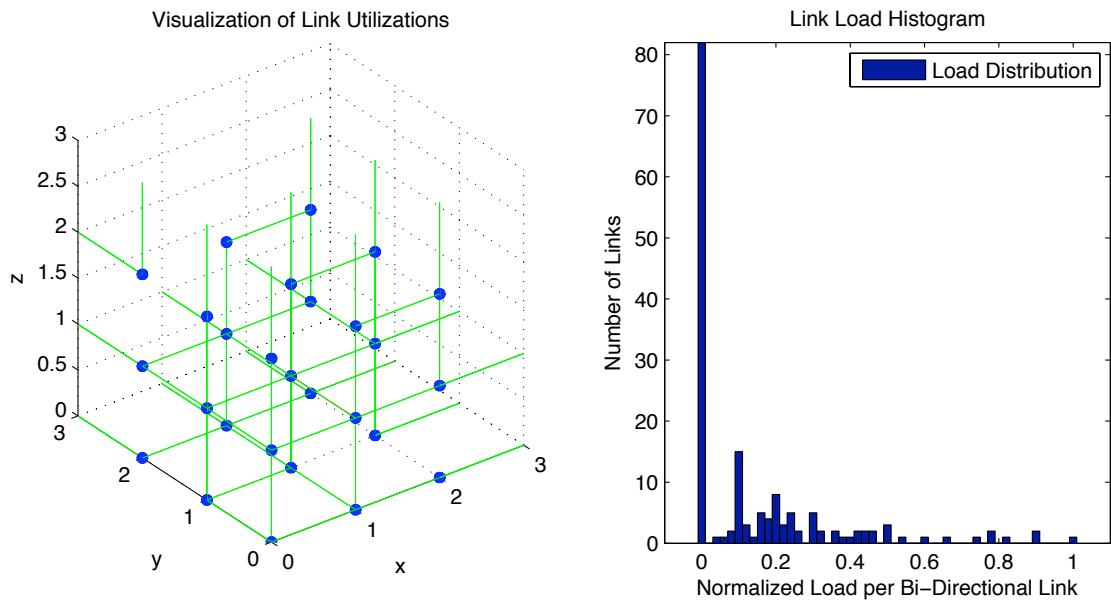


**Figure D.33:** Results of Enhanced CQR Routing using the nearest neighbor traffic demand (NN). Each node had two individual traffic demands to each of its six neighbors.

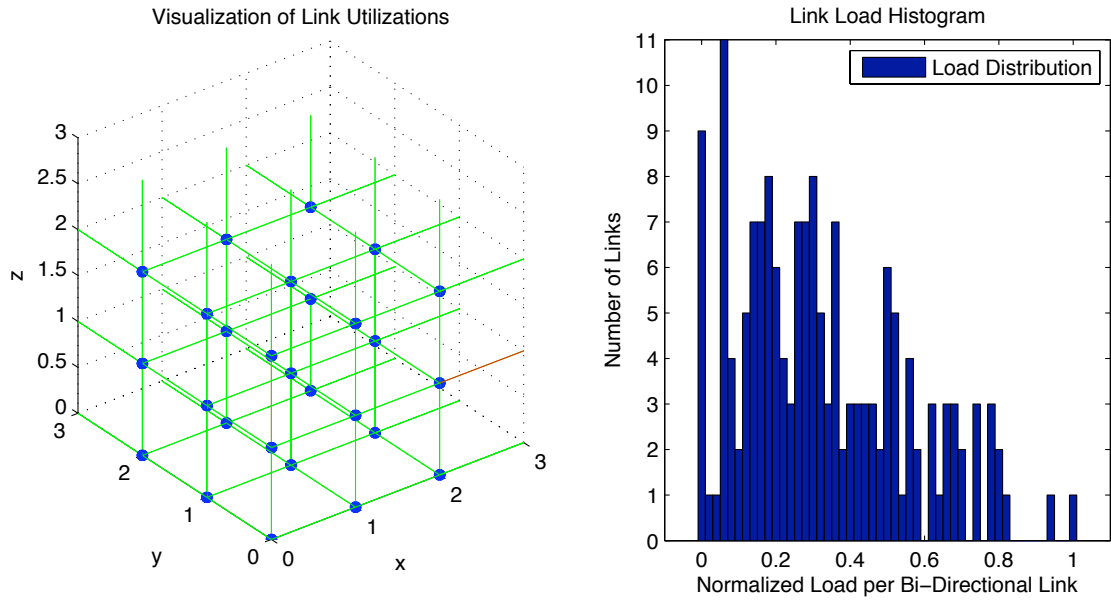
### D.4.3 Enhanced CQR Routing Results



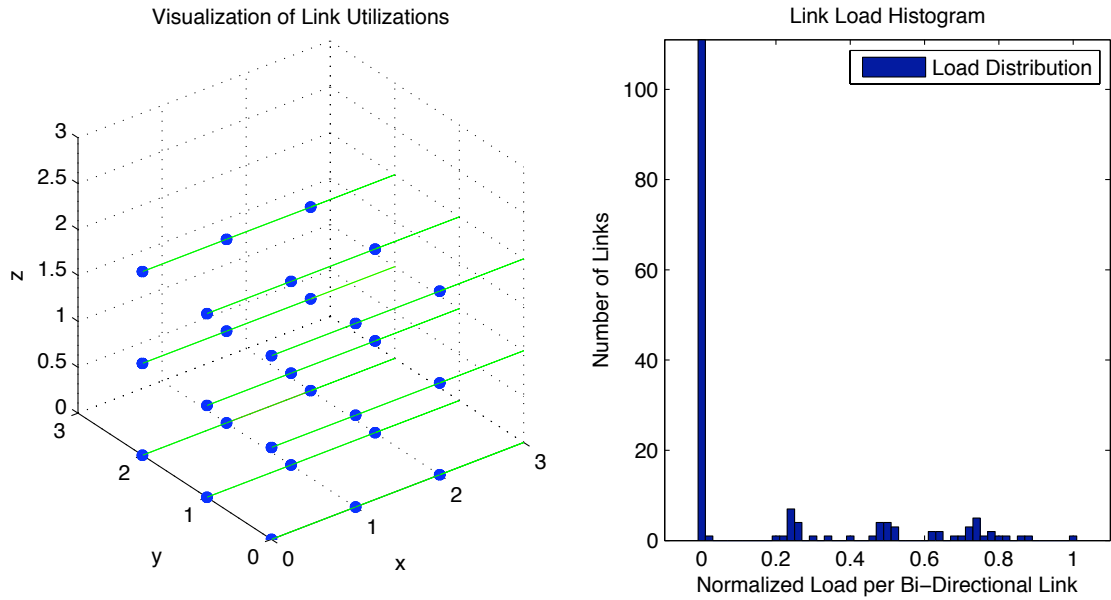
**Figure D.34:** Results of Enhanced CQR Routing using a uniform random traffic pattern (UR). For each node, a value of 0-9 nodes was assigned as having 0-9 individual traffic demands.



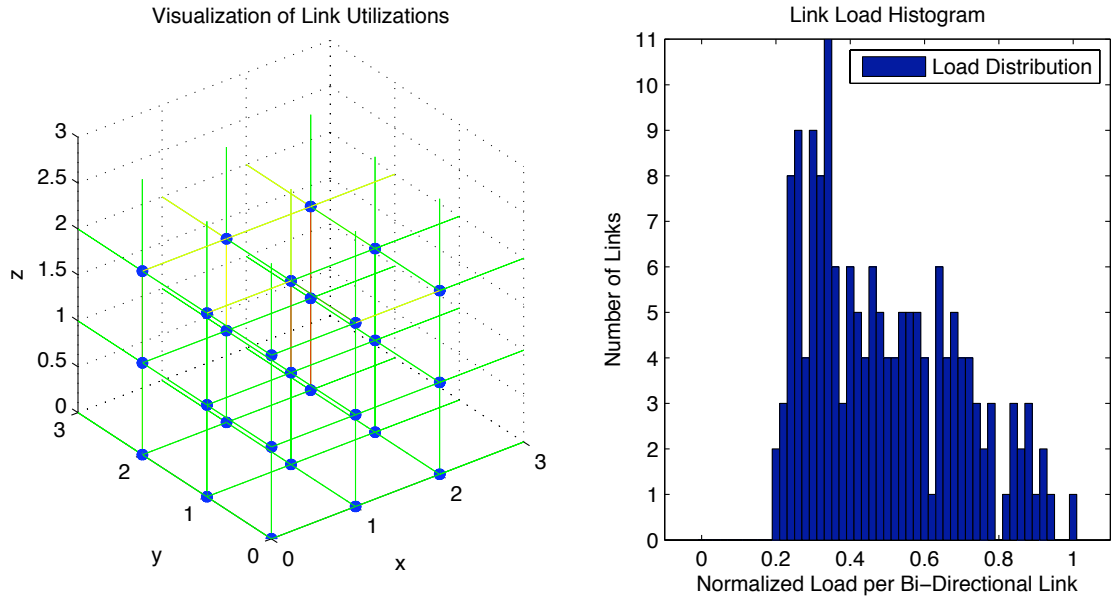
**Figure D.35:** Results of Enhanced CQR Routing using a bit complement traffic pattern (BC). Two individual traffic demands were assigned for possible values described in Table 2.1.



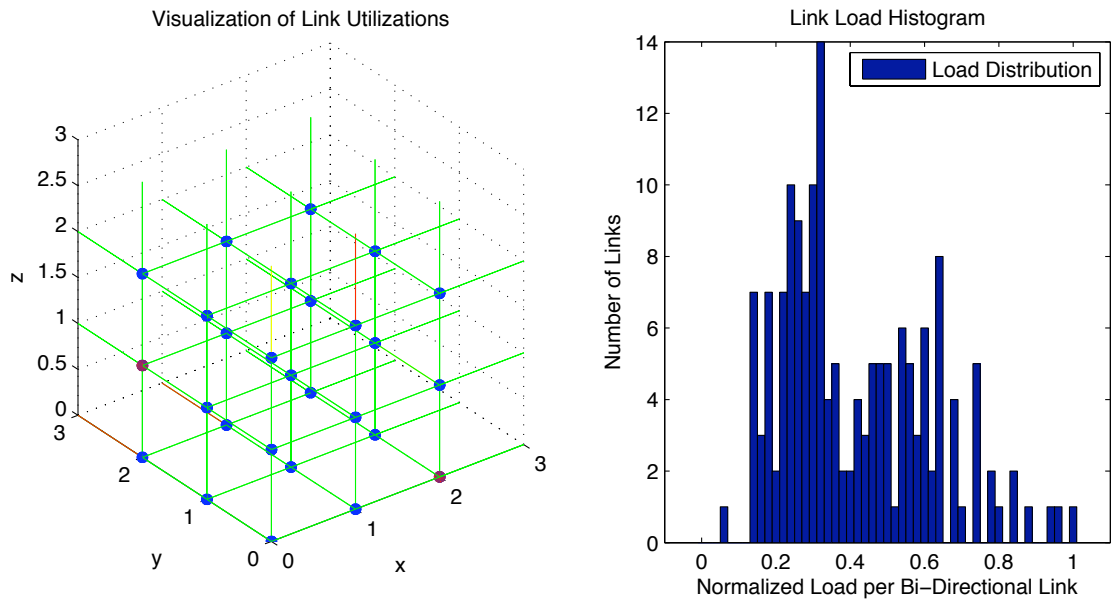
**Figure D.36:** Results of Enhanced CQR Routing using a transpose traffic pattern (TP). Two individual traffic demands were assigned for each permutation of all  $x, y, z$  values.



**Figure D.37:** Results of Enhanced CQR Routing using a tornado traffic pattern (TOR). Three individual traffic demands were assigned according to Table 2.1.



**Figure D.38:** Results of Enhanced CQR Routing using a flood traffic pattern (FL). Two individual traffic demands were assigned to each node from each node.



**Figure D.39:** Results of Enhanced CQR Routing using a flood traffic pattern (FL) with Hot-spots. Two individual traffic demands were assigned to each node from each node. Five percent of the nodes were made hot-spots, and they are marked as red in the figure.

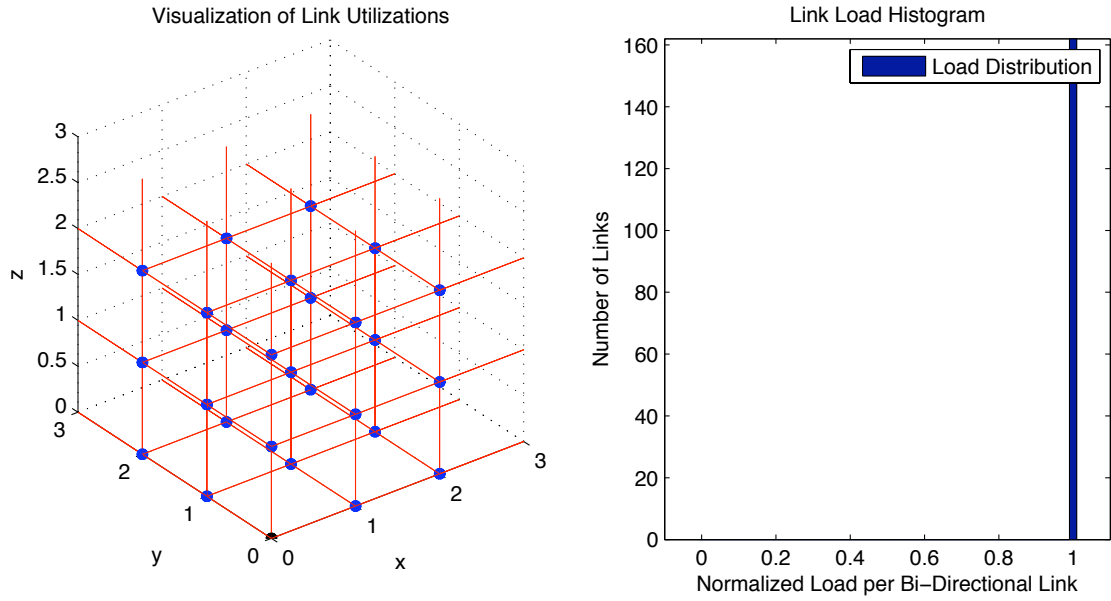


# Appendix E

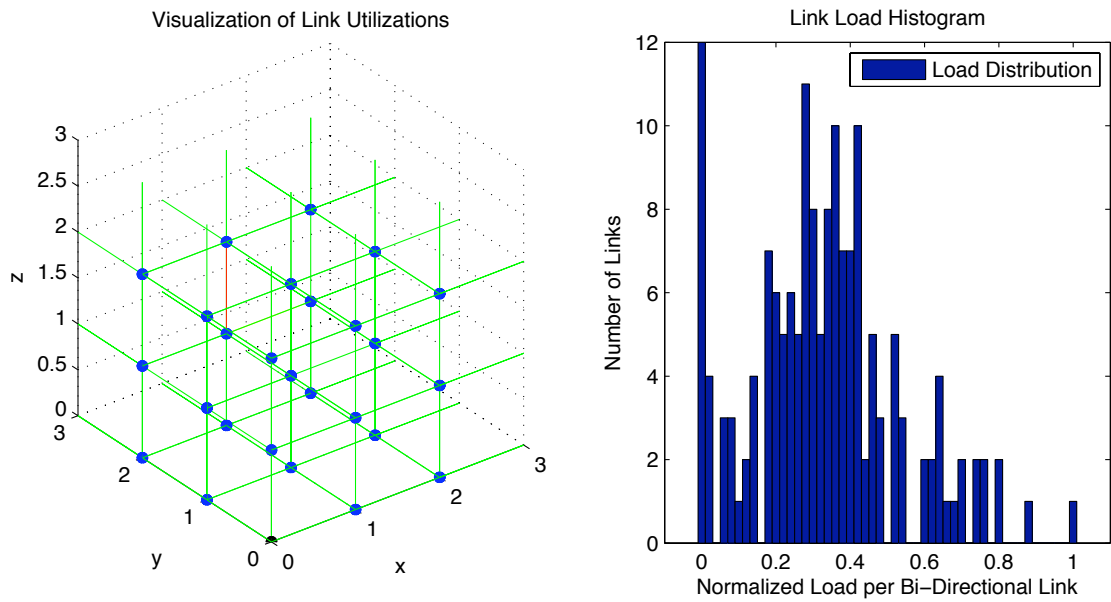
## Appendix E: Full Matlab Simulation Results

### E.1 Results from the 3-ary 3-cube Topology

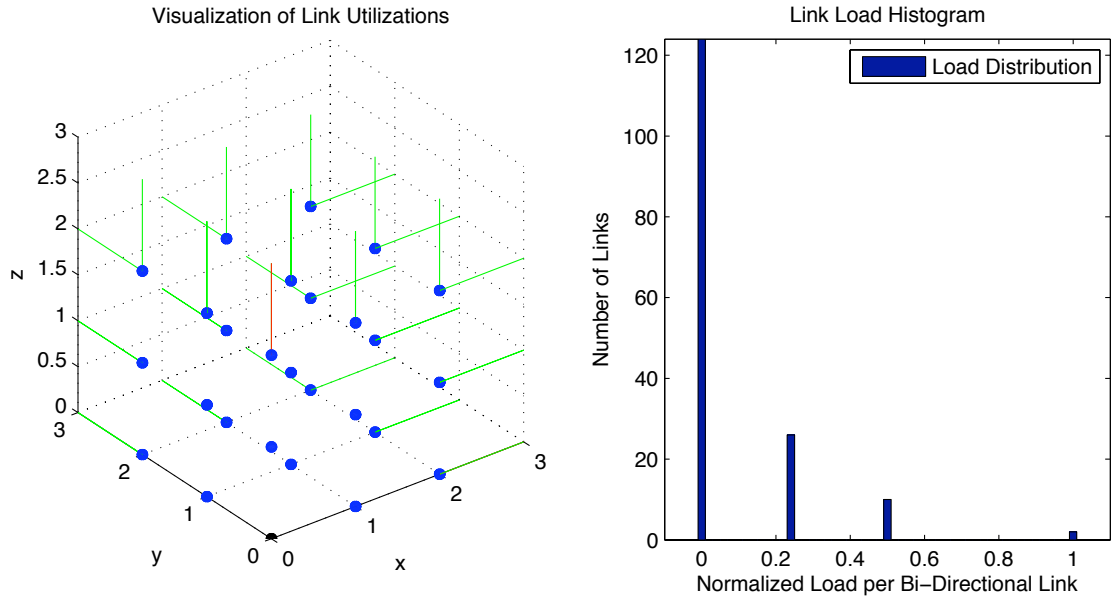
#### E.1.1 Static Algorithms



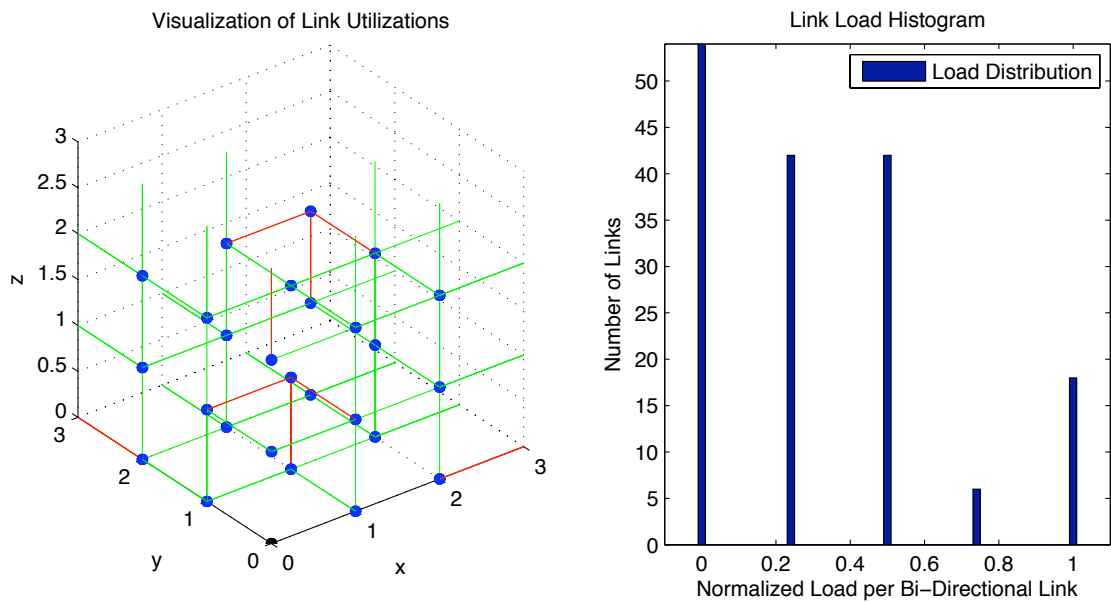
**Figure E.1:** Results of Dimension Ordered Routing (DOR) using the nearest neighbor traffic demand (NN). Each node had two individual traffic demands to each of its six neighbors.



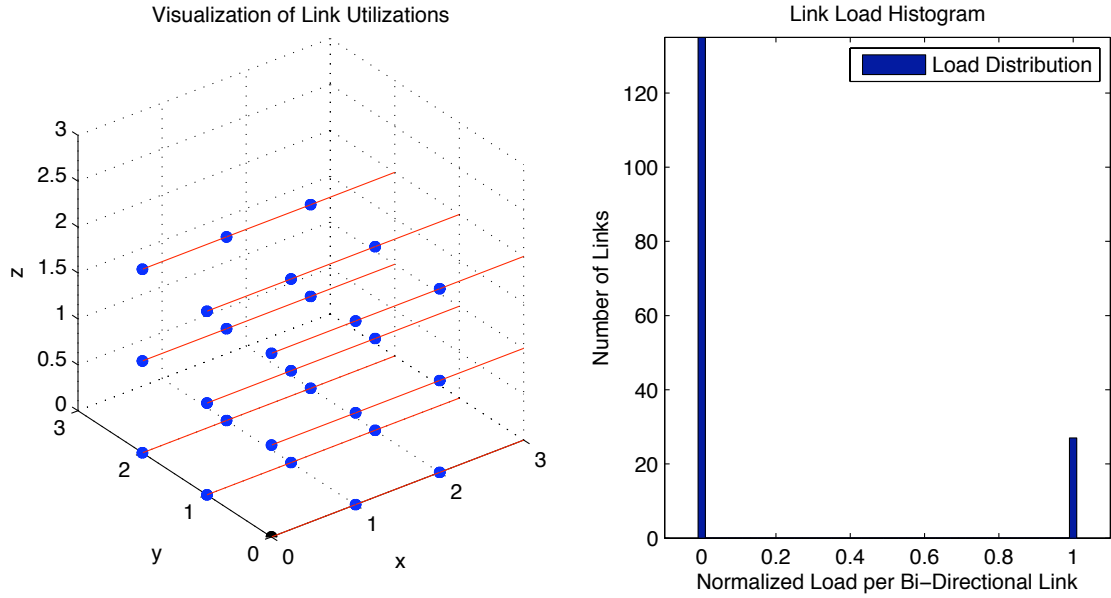
**Figure E.2:** Results of Dimension Ordered Routing (DOR) using the uniform random traffic demand (UR). For each node, a value of 0-9 was assigned as having 0-9 individual traffic demands.



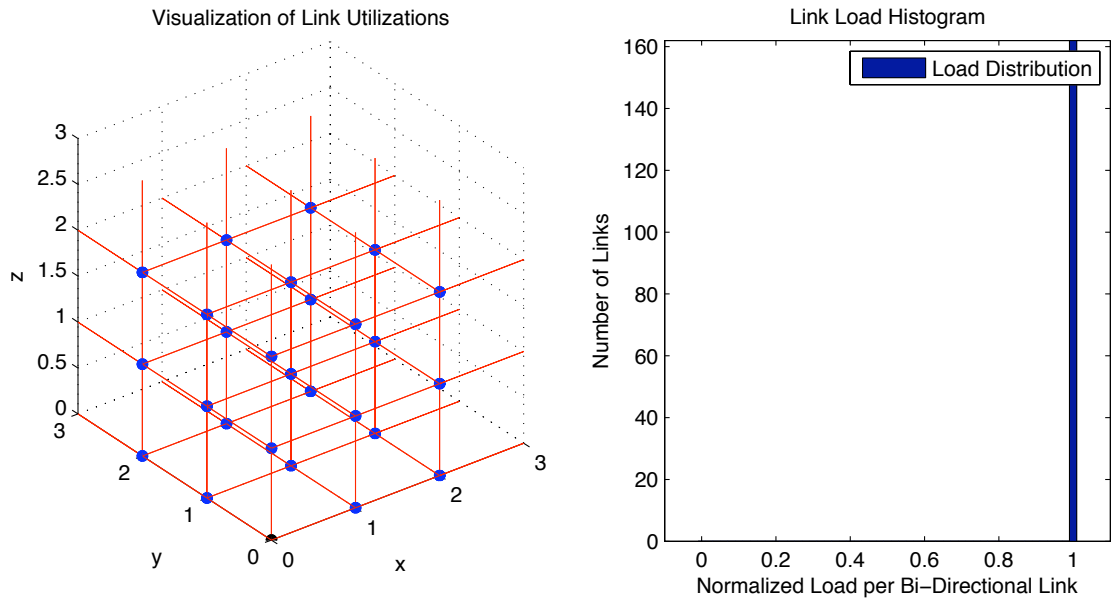
**Figure E.3:** Results of Dimension Ordered Routing (DOR) using the bit complement traffic demand (BC). Two individual traffic demands were assigned for possible values described in Table 2.1.



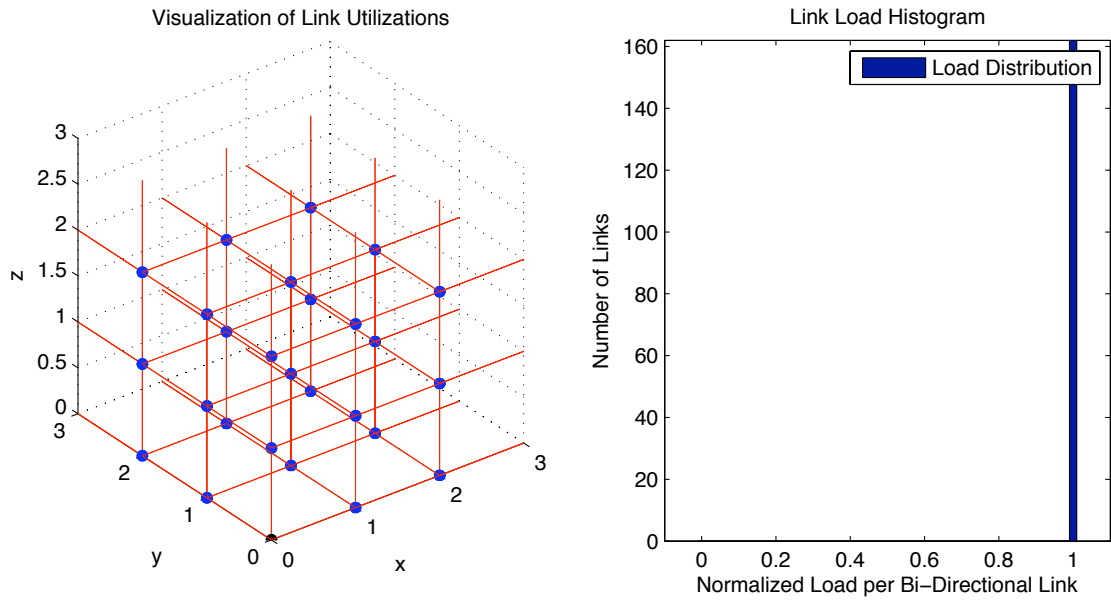
**Figure E.4:** Results of Dimension Ordered Routing (DOR) using the transpose traffic demand (TP). Two individual traffic demands were assigned for each permutation of all  $x, y, z$  values.



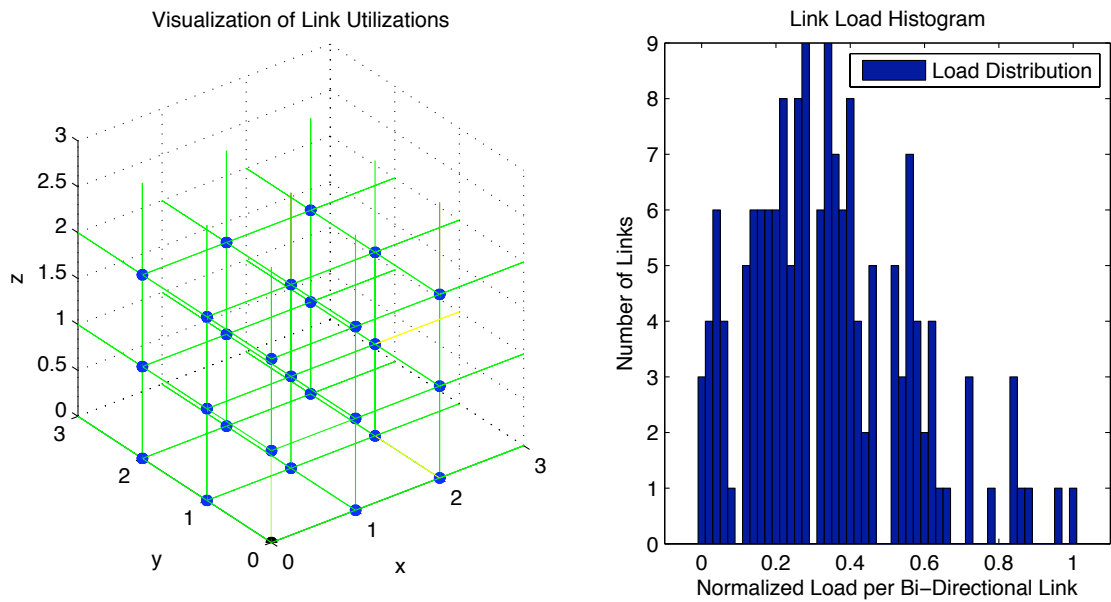
**Figure E.5:** Results of Dimension Ordered Routing (DOR) using the tornado traffic demand (TOR). Three individual traffic demands were assigned for possible values described in Table 2.1.



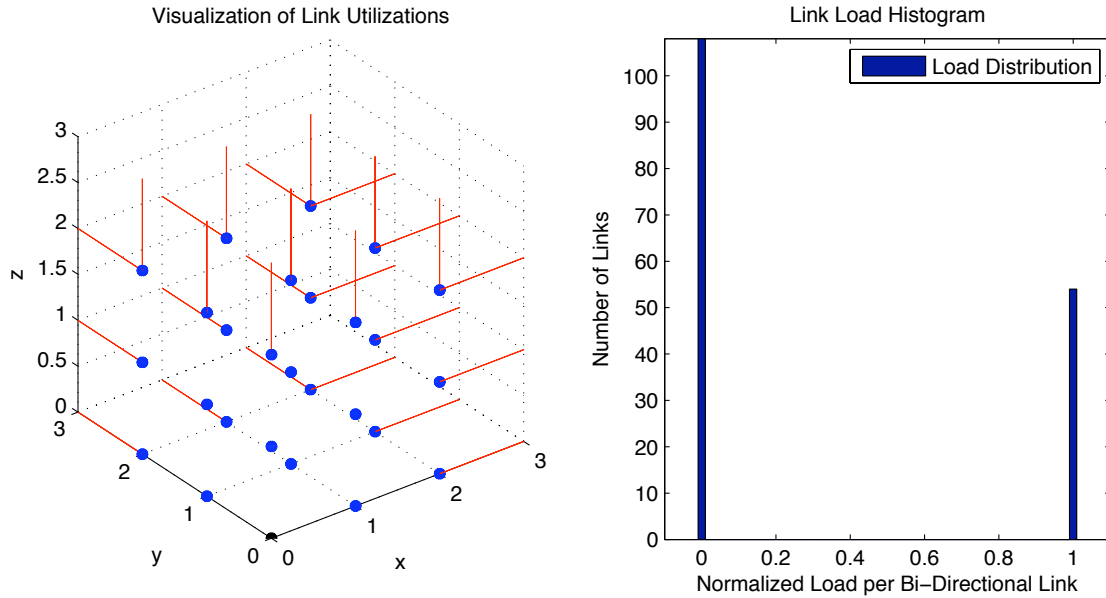
**Figure E.6:** Results of Dimension Ordered Routing (DOR) using the flood traffic demand (FL). Two individual traffic demands were assigned to each node from each node.



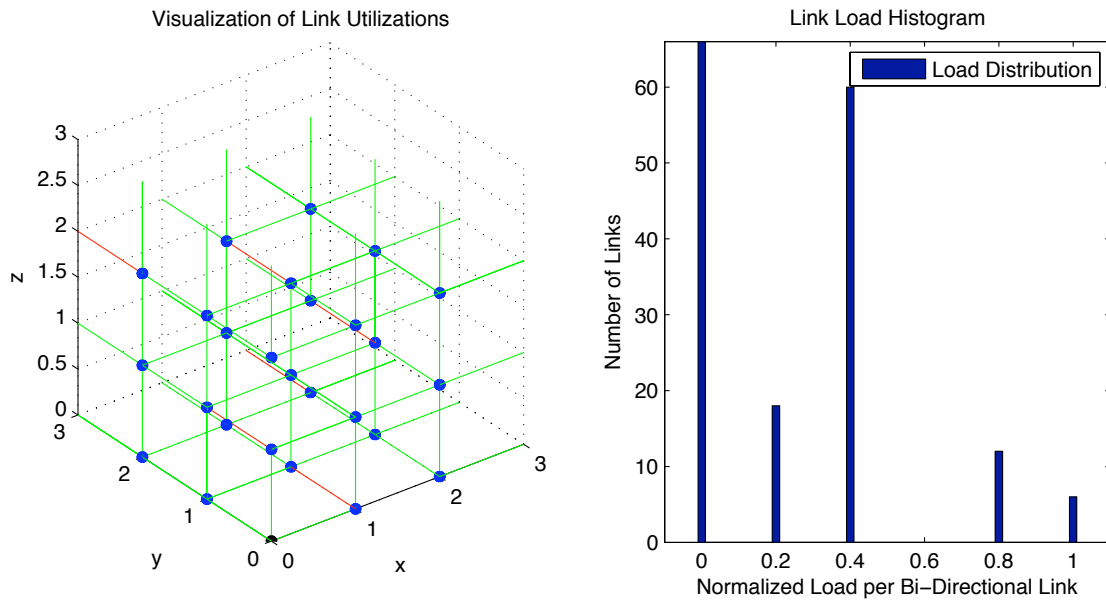
**Figure E.7:** Results of Direction Ordered Routing (DIR) using the nearest neighbor traffic demand (NN). Each node had two individual traffic demands to each of its six neighbors.



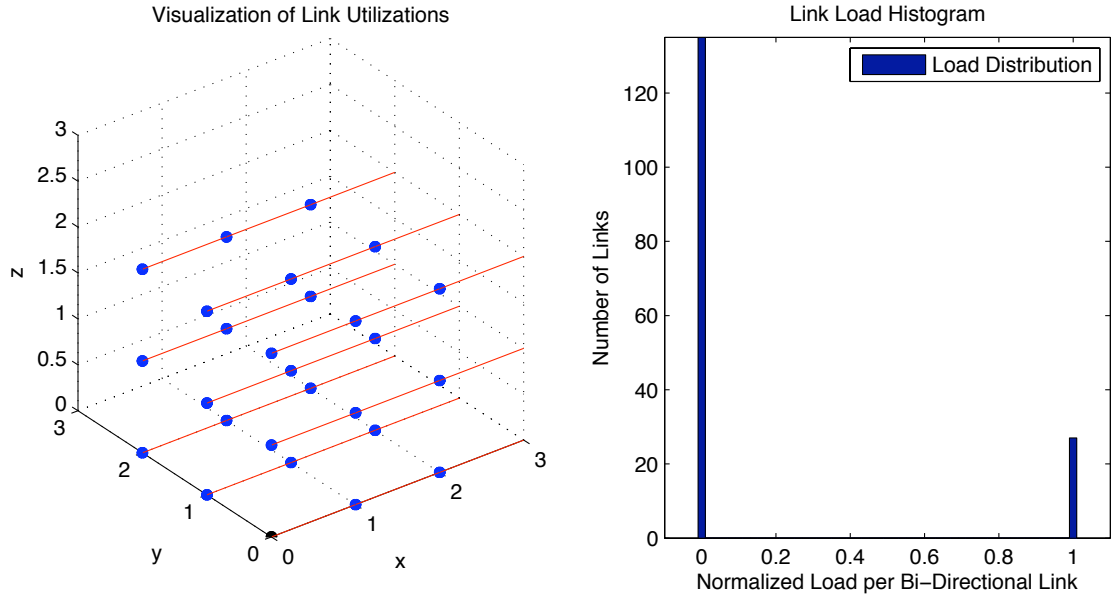
**Figure E.8:** Results of Direction Ordered Routing (DIR) using the uniform random traffic demand (UR). For each node, a value of 0-9 was assigned as having 0-9 individual traffic demands.



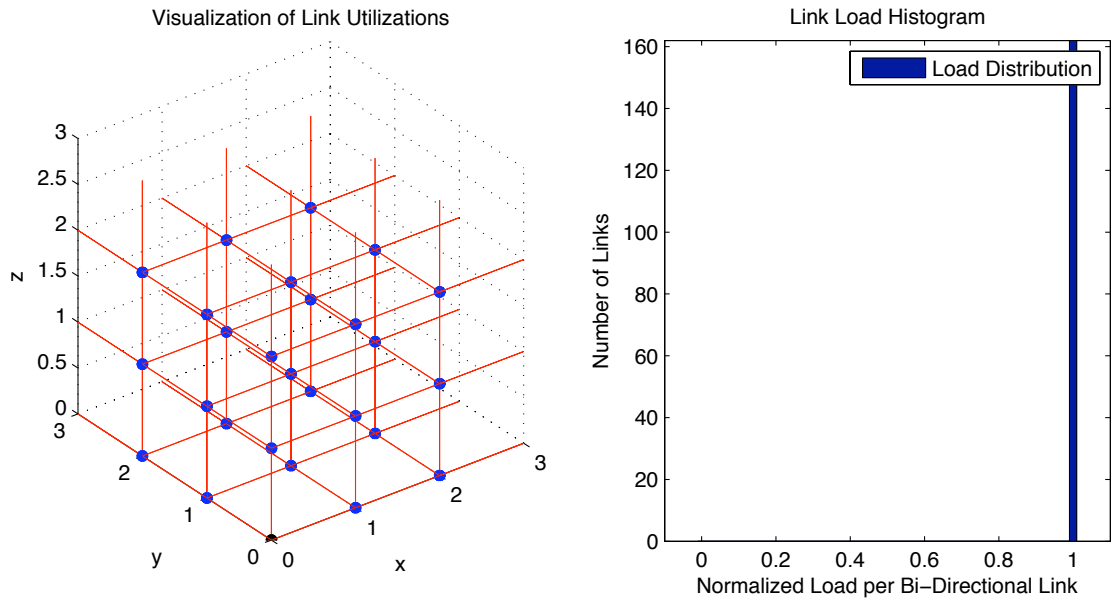
**Figure E.9:** Results of Direction Ordered Routing (DIR) using the bit complement traffic demand (BC). Two individual traffic demands were assigned for possible values described in Table 2.1.



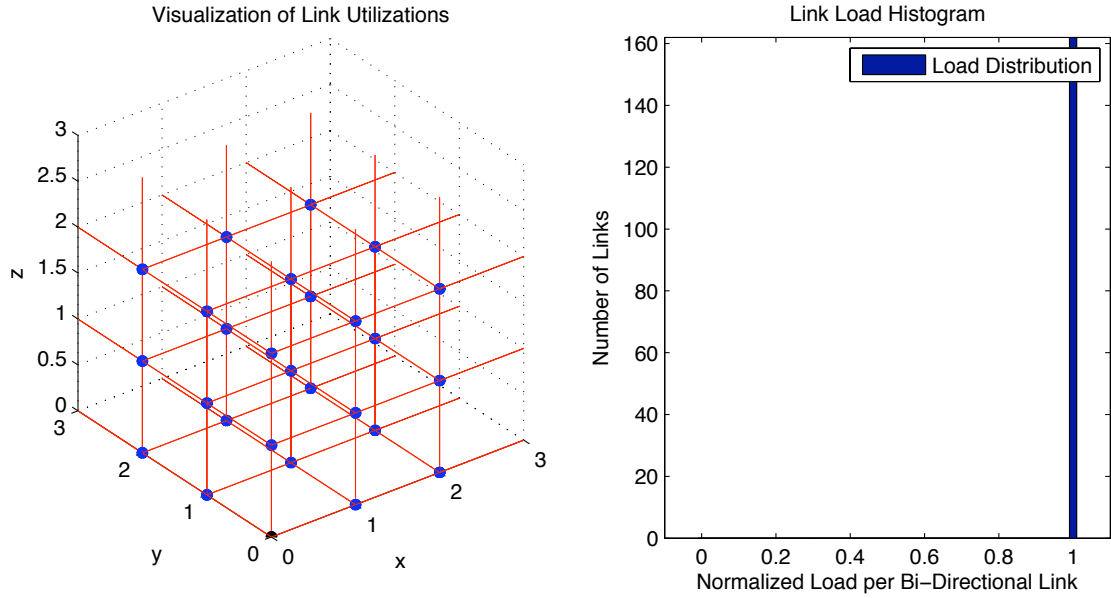
**Figure E.10:** Results of Direction Ordered Routing (DIR) using the transpose traffic demand (TP). Two individual traffic demands were assigned for each permutation of all  $x, y, z$  values.



**Figure E.11:** Results of Direction Ordered Routing (DIR) using the tornado traffic demand (TOR). Three individual traffic demands were assigned for possible values described in Table 2.1.



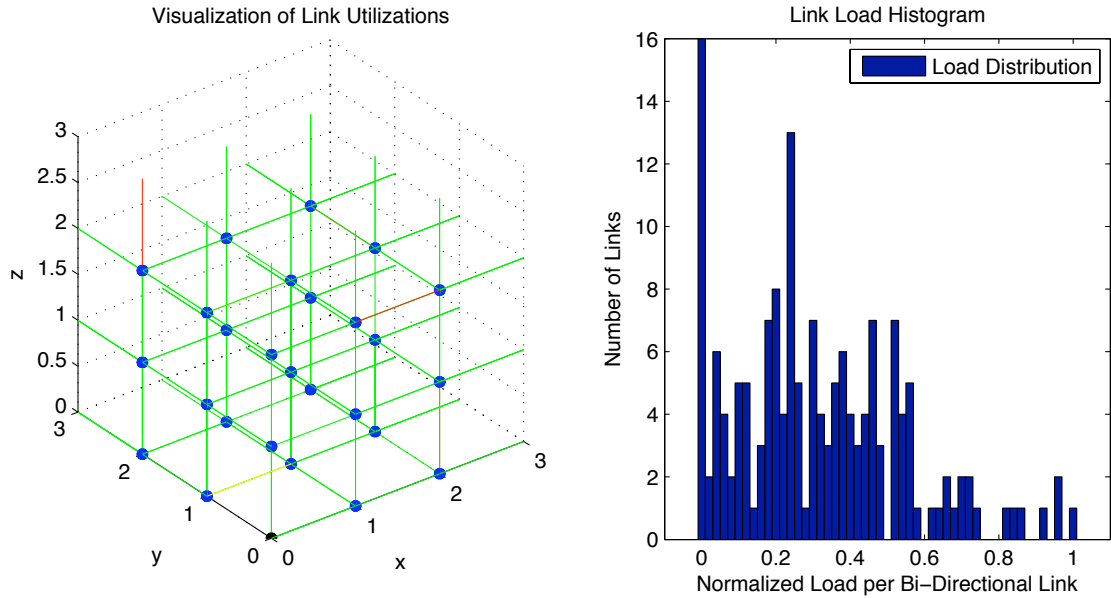
**Figure E.12:** Results of Direction Ordered Routing (DIR) using the flood traffic demand (FL). Two individual traffic demands were assigned to each node from each node.



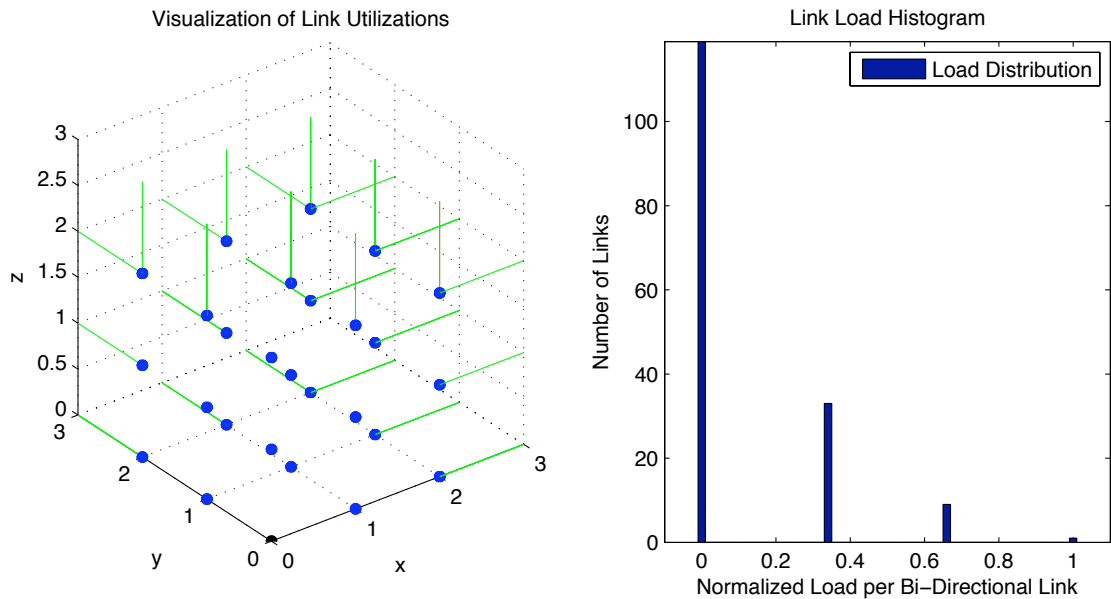
**Figure E.13:** Results of Minimal Oblivious Routing (MO) using the nearest neighbor traffic demand (NN). Each node had two individual traffic demands to each of its six neighbors.

## E.1.2 Oblivious Algorithms

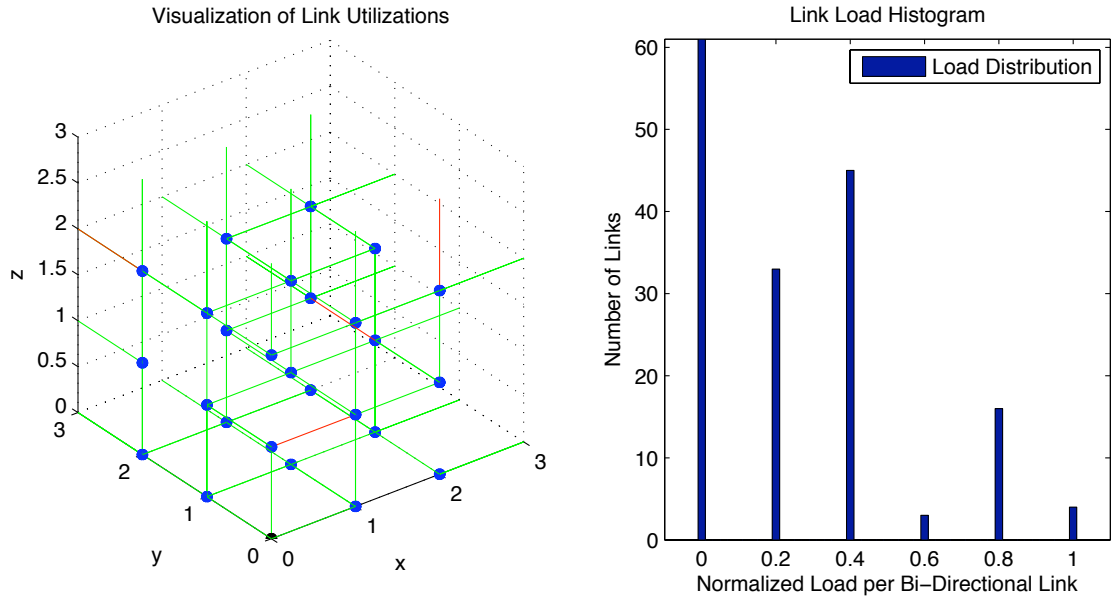




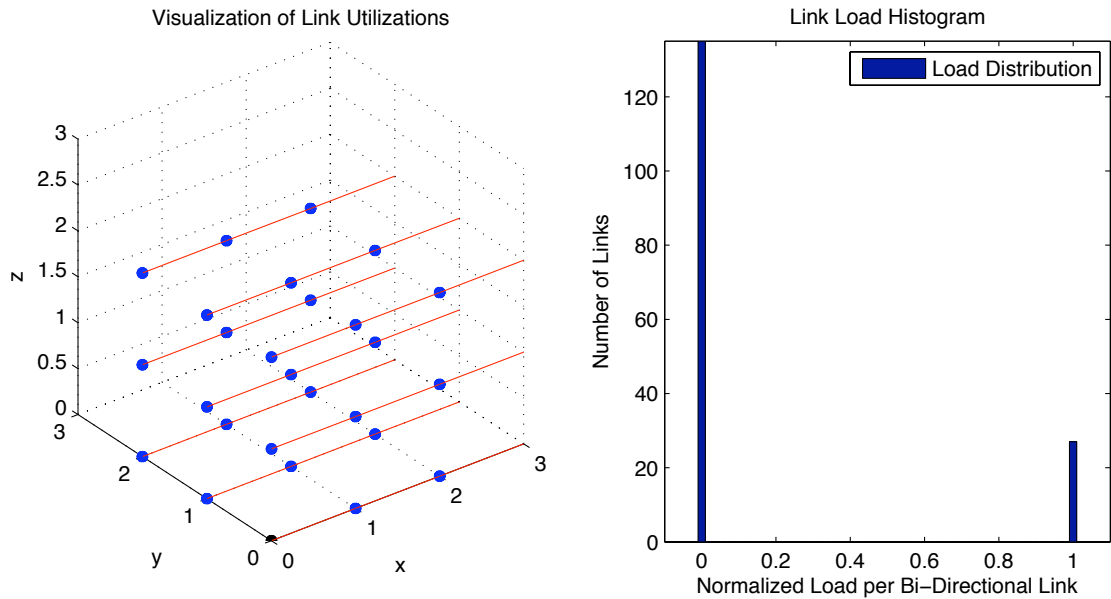
**Figure E.14:** Results of Minimal Oblivious Routing (MO) using the uniform random traffic demand (UR). For each node, a value of 0-9 was assigned as having 0-9 individual traffic demands.



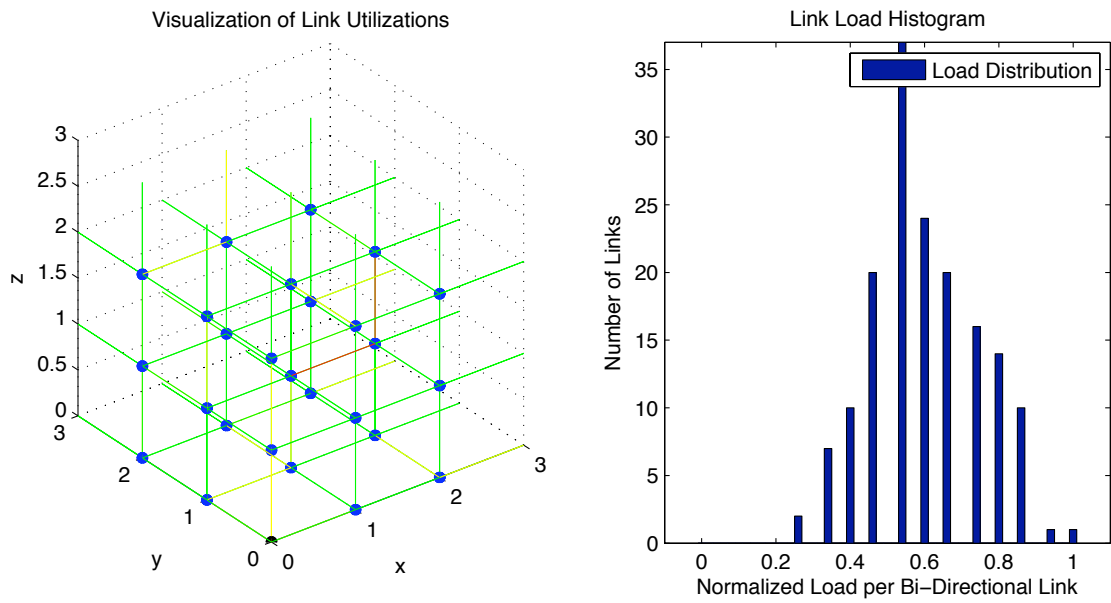
**Figure E.15:** Results of Minimal Oblivious Routing (MO) using the bit complement traffic demand (BC). Two individual traffic demands were assigned for possible values described in Table 2.1.



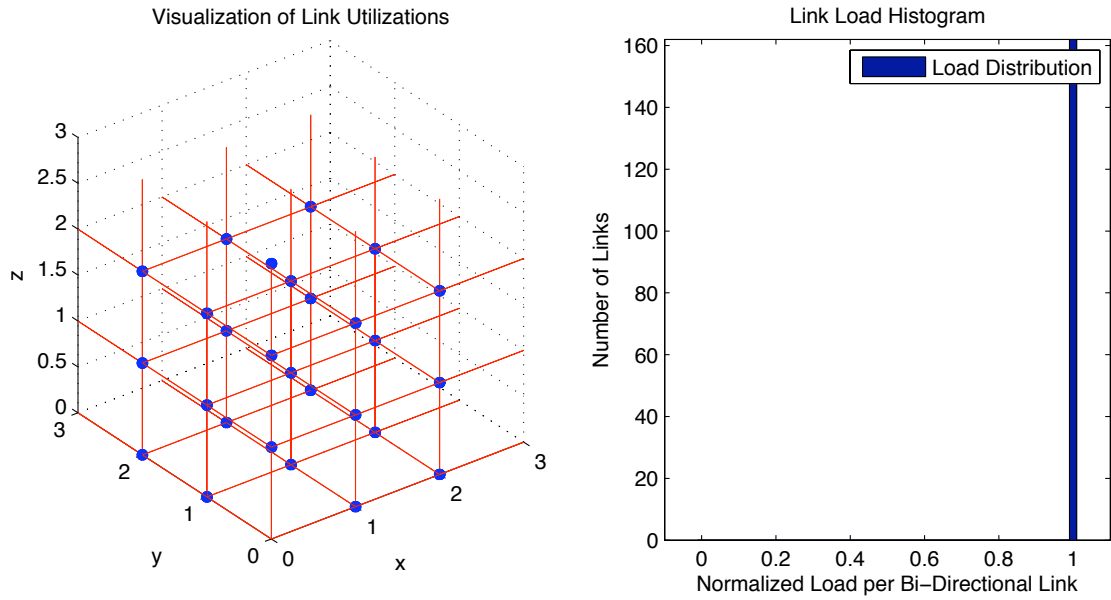
**Figure E.16:** Results of Minimal Oblivious Routing (MO) using the transpose traffic demand (TP). Two individual traffic demands were assigned for each permutation of all  $x, y, z$  values.



**Figure E.17:** Results of Minimal Oblivious Routing (MO) using the tornado traffic demand (TOR). Three individual traffic demands were assigned for possible values described in Table 2.1.

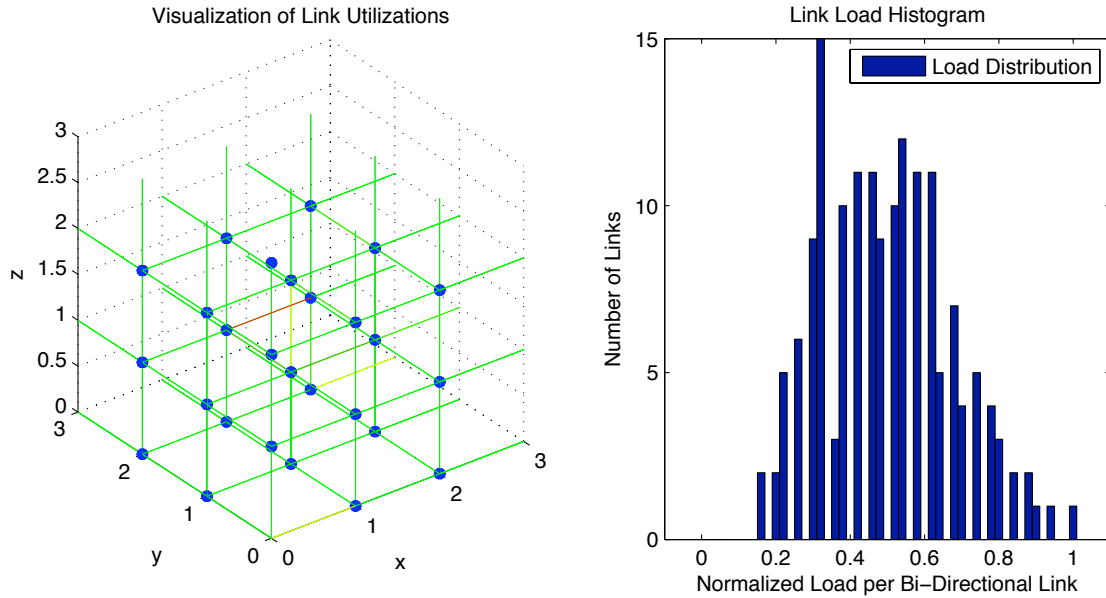


**Figure E.18:** Results of Minimal Oblivious Routing (MO) using the flood traffic demand (FL). Two individual traffic demands were assigned to each node from each node.

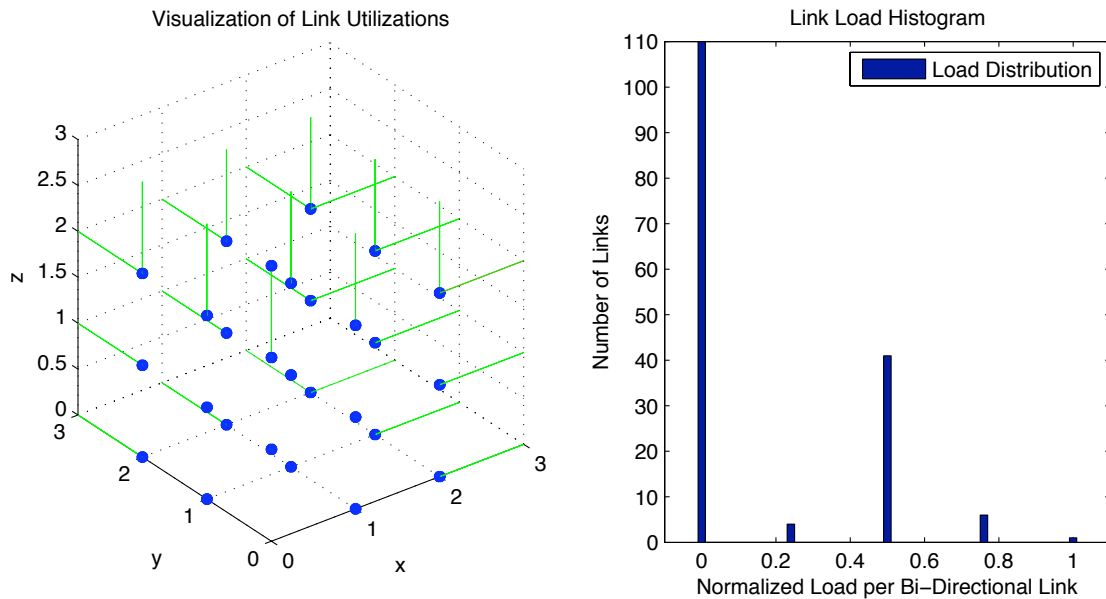


**Figure E.19:** Results of Minimal Adaptive Routing (MA) using the nearest neighbor traffic demand (NN). Each node had two individual traffic demands to each of its six neighbors.

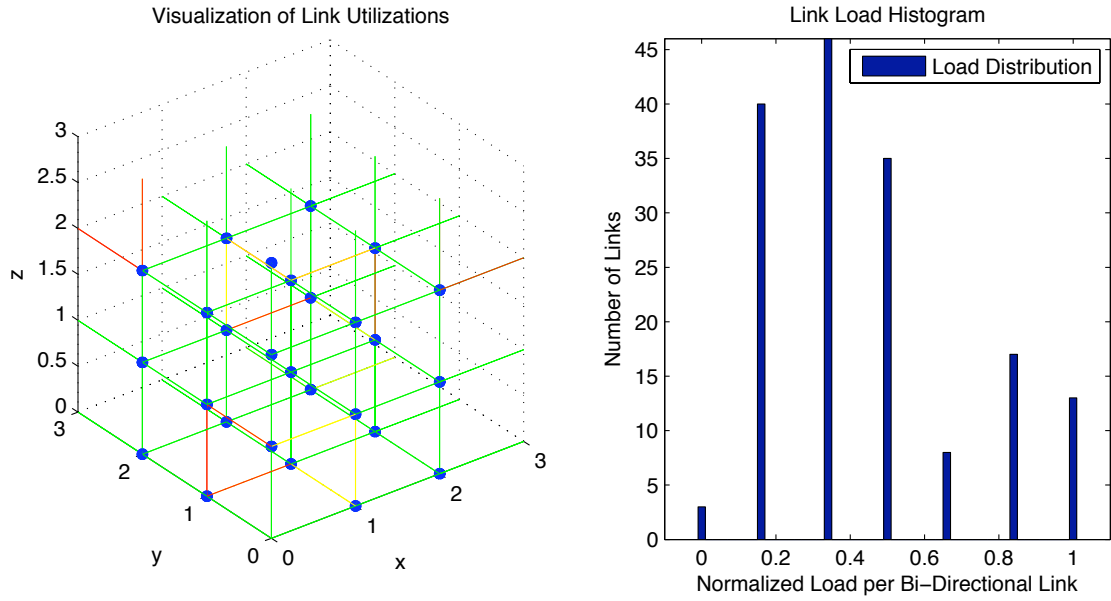
### E.1.3 Adaptive Algorithms



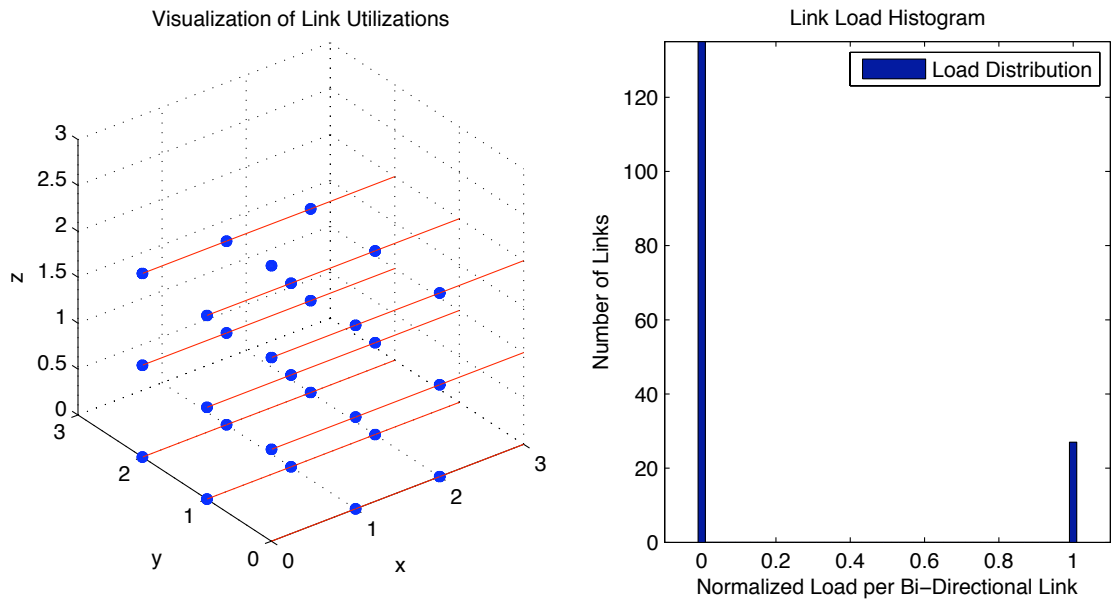
**Figure E.20:** Results of Minimal Adaptive Routing (MA) using the uniform random traffic demand (UR). For each node, a value of 0-9 was assigned as having 0-9 individual traffic demands.



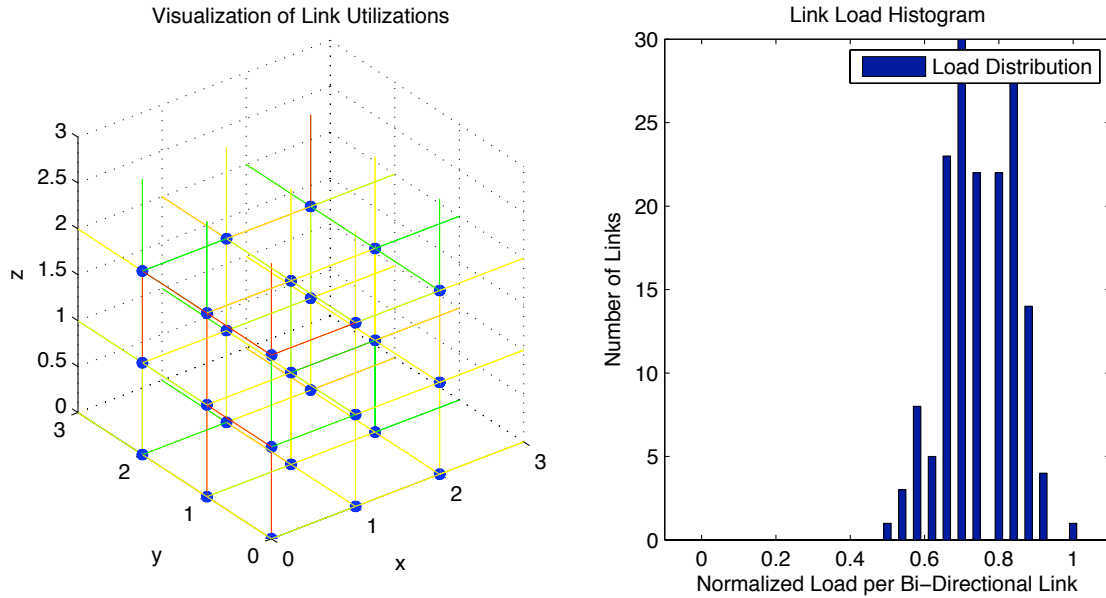
**Figure E.21:** Results of Minimal Adaptive Routing (MA) using the bit complement traffic demand (BC). Two individual traffic demands were assigned for possible values described in Table 2.1.



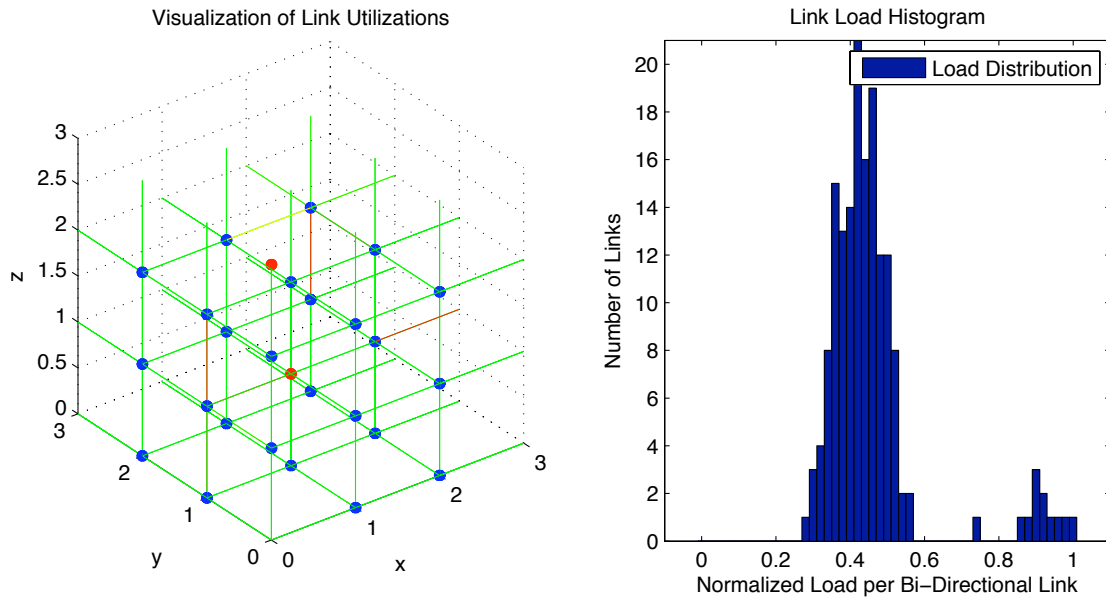
**Figure E.22:** Results of Minimal Adaptive Routing (MA) using the transpose traffic demand (TP). Two individual traffic demands were assigned for each permutation of all  $x, y, z$  values.



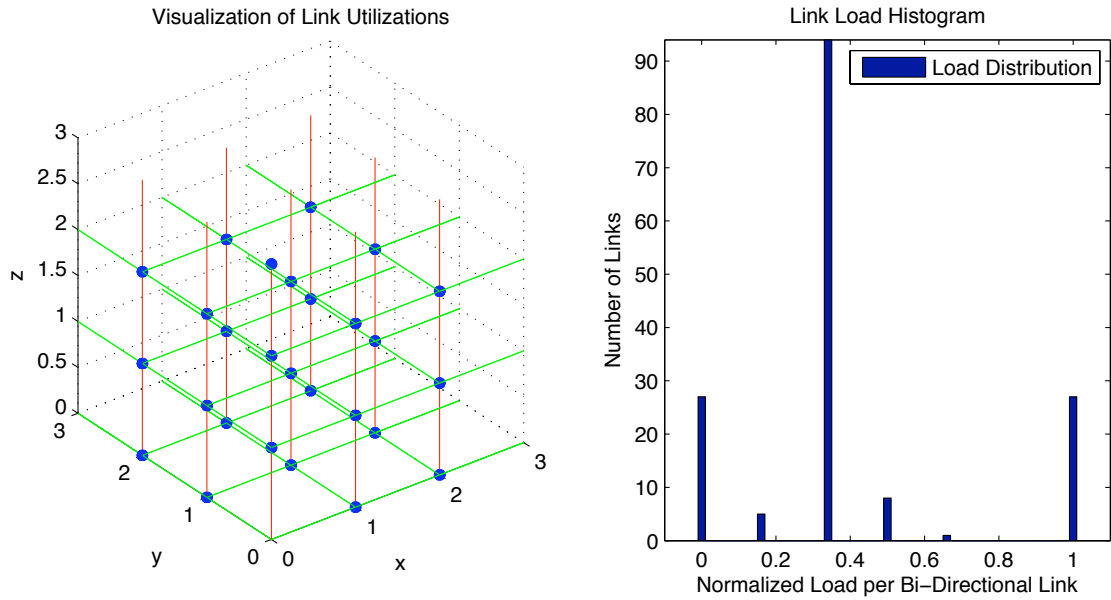
**Figure E.23:** Results of Minimal Adaptive Routing (MA) using the tornado traffic demand (TOR). Three individual traffic demands were assigned for possible values described in Table 2.1.



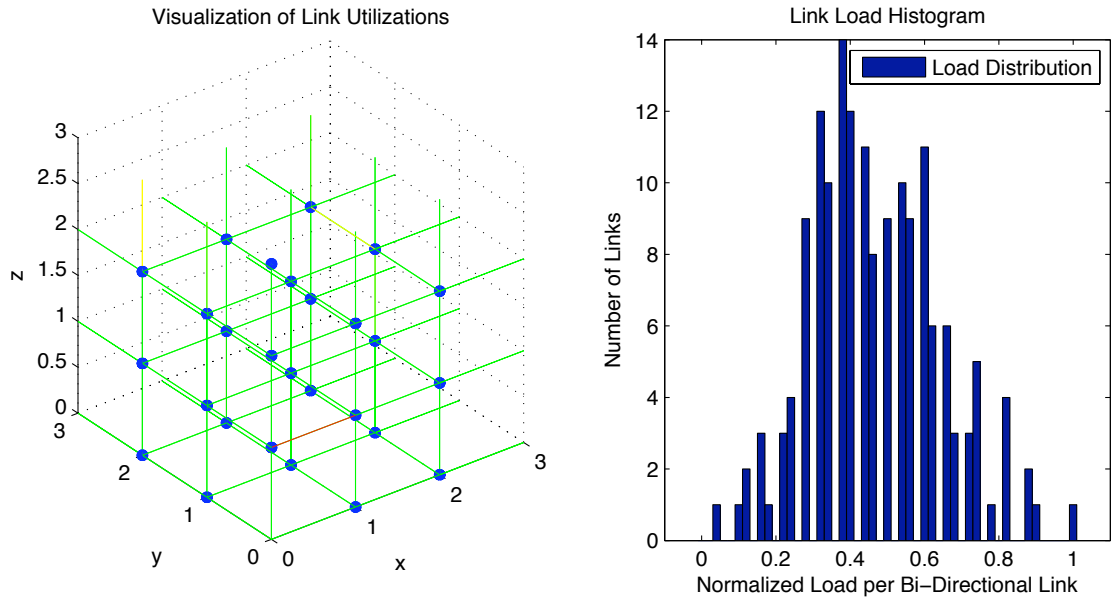
**Figure E.24:** Results of Minimal Adaptive Routing (MA) using the flood traffic demand (FL). Two individual traffic demands were assigned to each node from each node.



**Figure E.25:** Results of Minimal Adaptive Routing (MA) using the flood traffic demand (FL) with Hot-spots. Two individual traffic demands were assigned to each node from each node. Five percent of the nodes were made hot-spots, and they are marked as red in the figure.

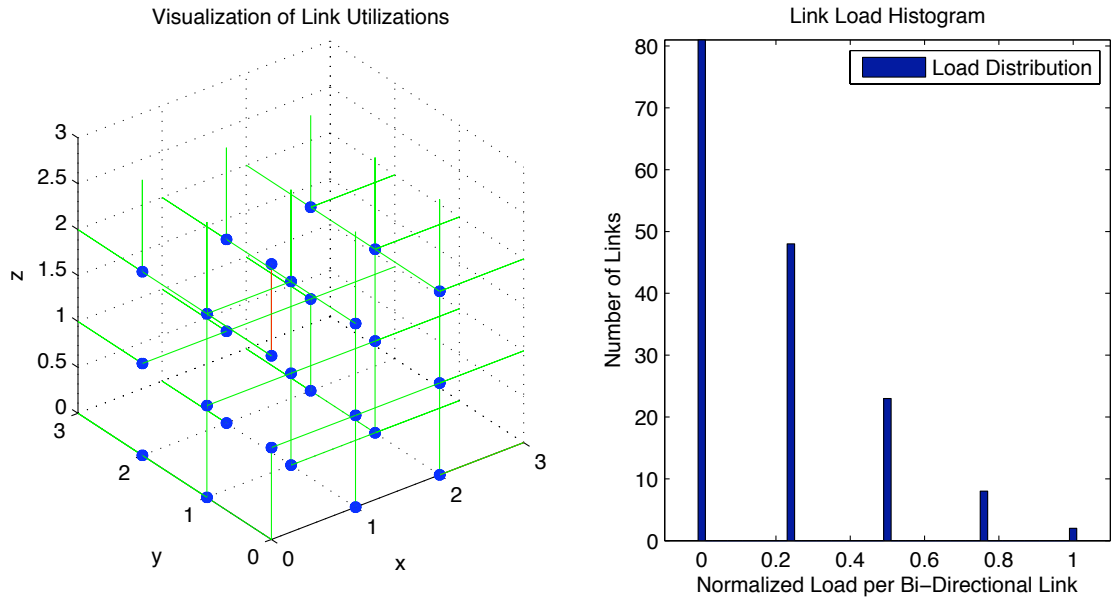


**Figure E.26:** Results of CQR Routing using the nearest neighbor traffic demand (NN). Each node had two individual traffic demands to each of its six neighbors.

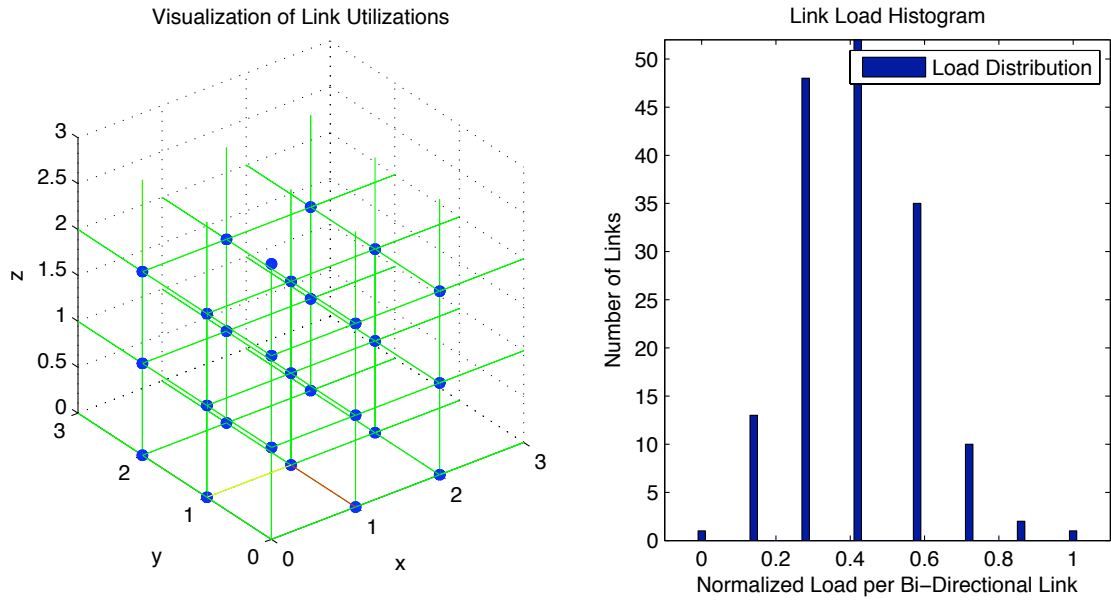


**Figure E.27:** Results of CQR Routing using the uniform random traffic demand (UR). For each node, a value of 0-9 was assigned as having 0-9 individual traffic demands.

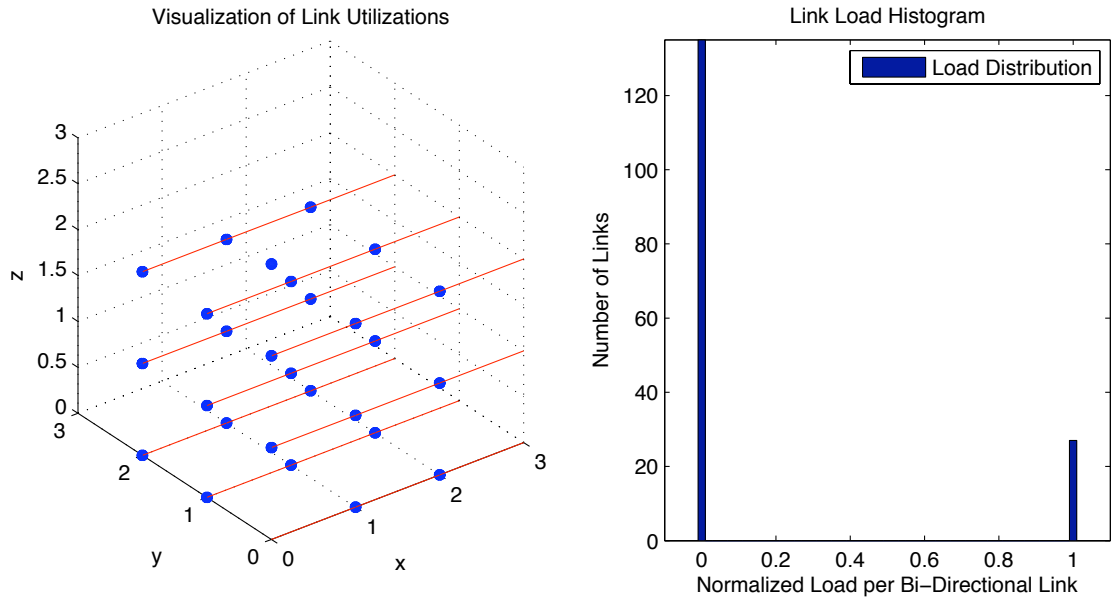




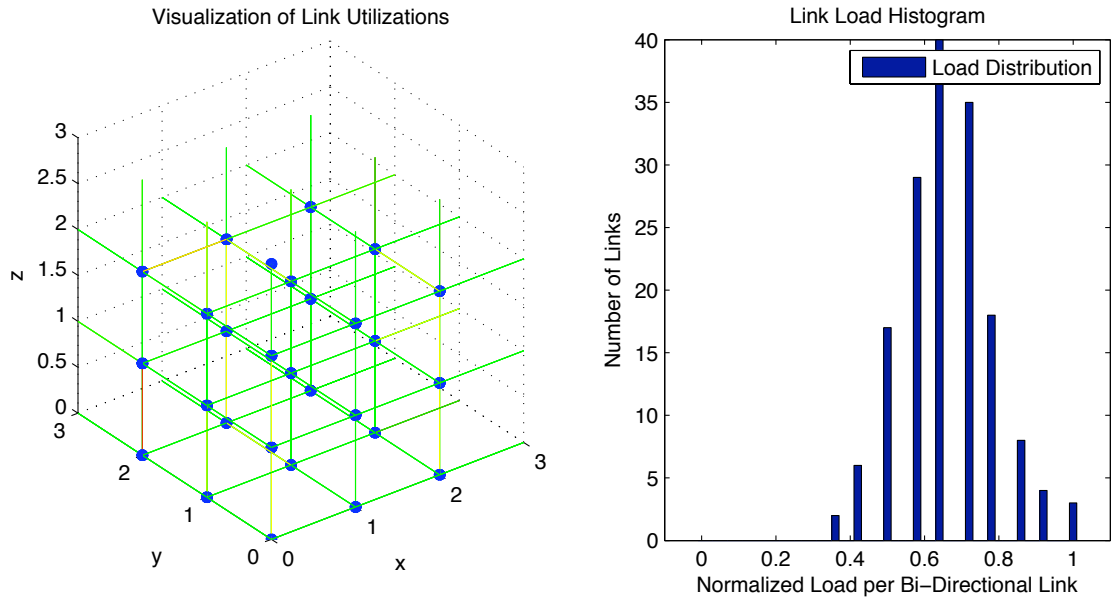
**Figure E.28:** Results of CQR Routing using the bit complement traffic demand (BC). Two individual traffic demands were assigned for possible values described in Table 2.1.



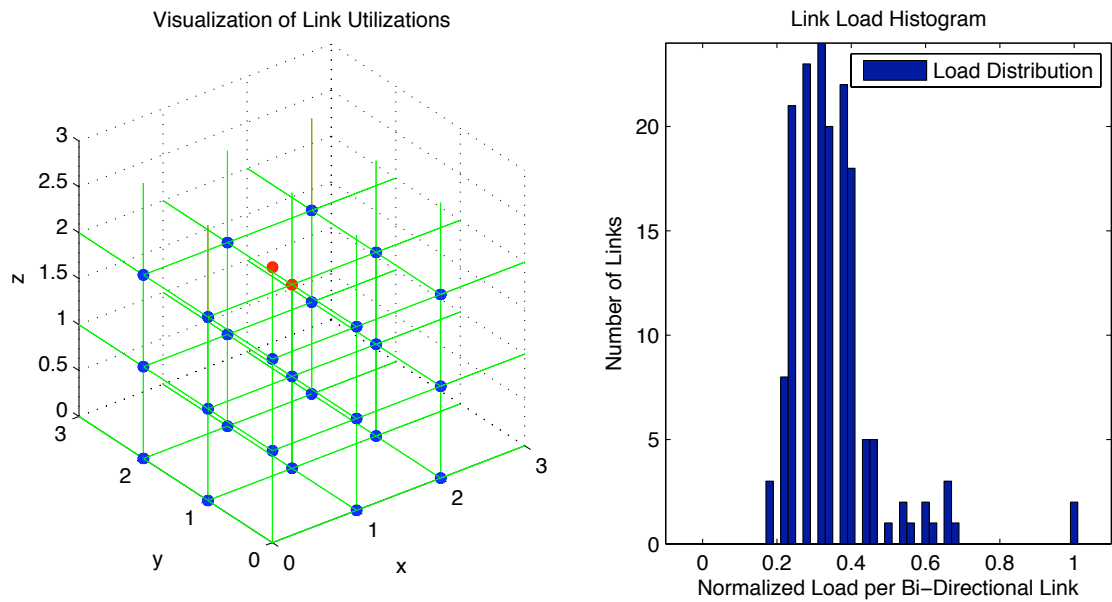
**Figure E.29:** Results of CQR Routing using the transpose traffic demand (TP). Two individual traffic demands were assigned for each permutation of all  $x, y, z$  values.



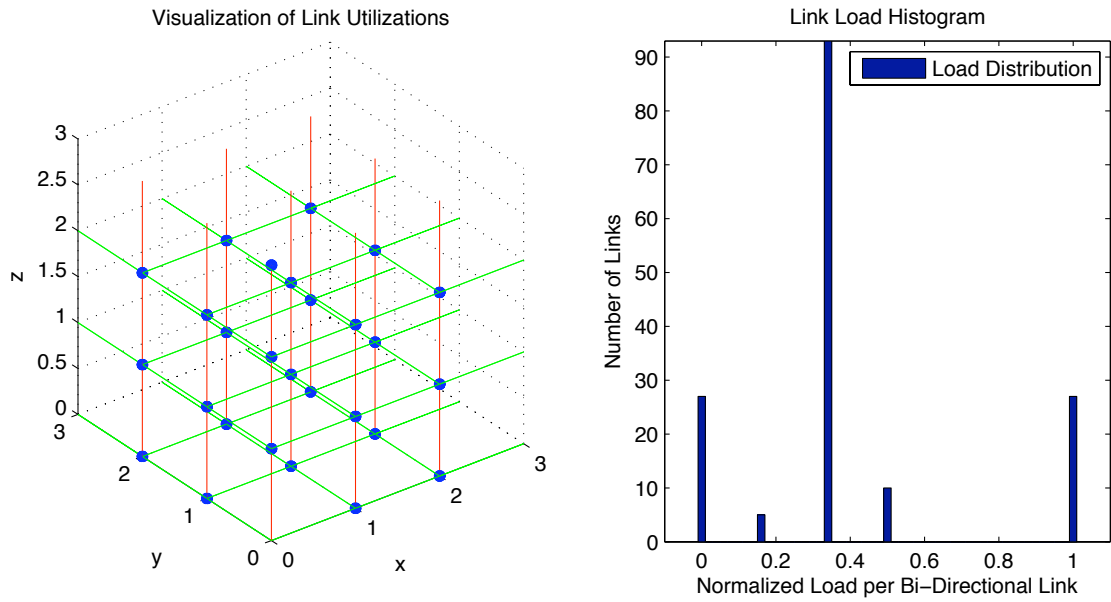
**Figure E.30:** Results of CQR Routing using the tornado traffic demand (TOR). Three individual traffic demands were assigned for possible values described in Table 2.1.



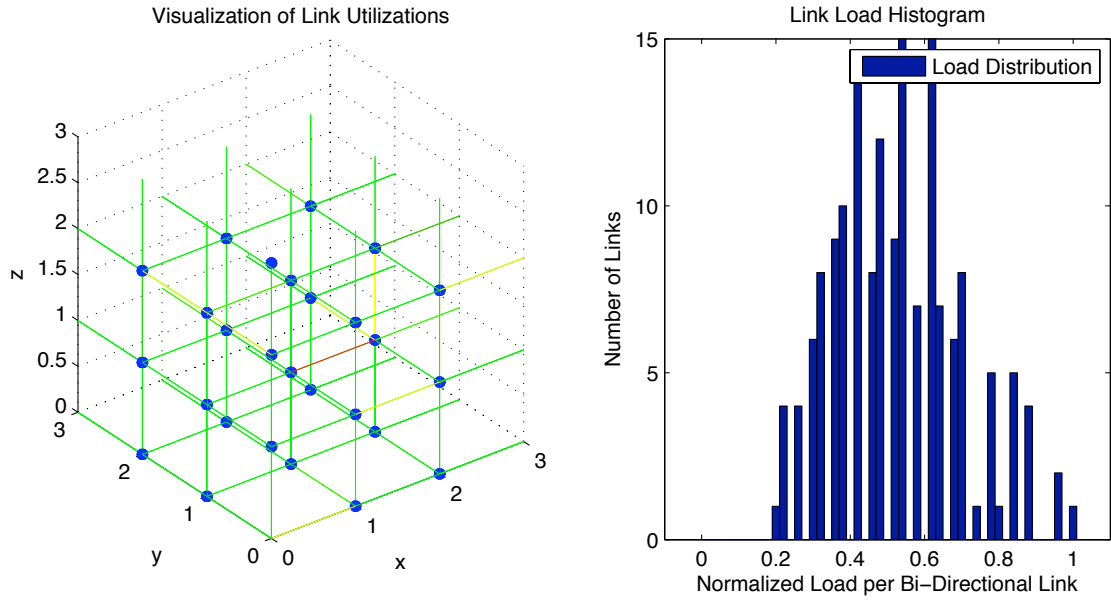
**Figure E.31:** Results of CQR Routing using the flood traffic demand (FL). Two individual traffic demands were assigned to each node from each node.



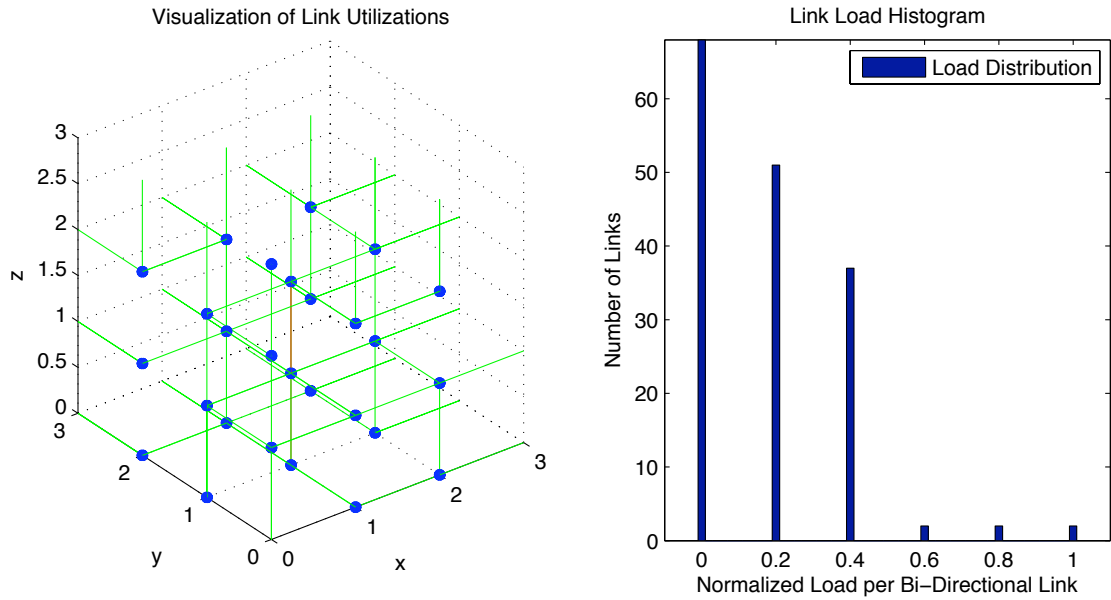
**Figure E.32:** Results of CQR Routing using the flood traffic demand (FL) with Hot-spots. Two individual traffic demands were assigned to each node from each node. Five percent of the nodes were made hot-spots, and they are marked as red in the figure.



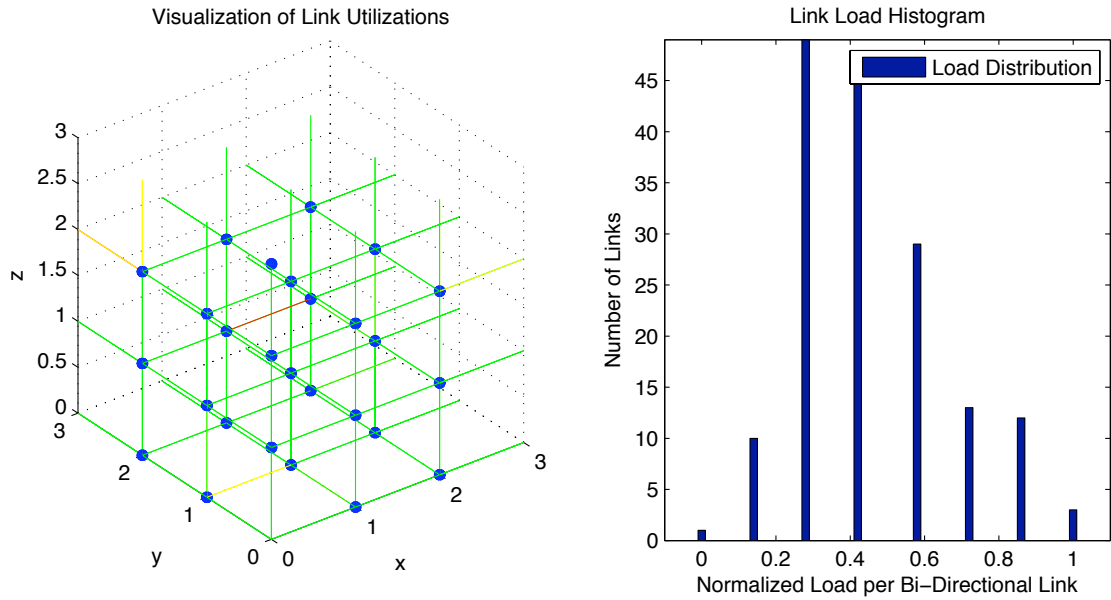
**Figure E.33:** Results of Enhanced CQR Routing using the nearest neighbor traffic demand (NN). Each node had two individual traffic demands to each of its six neighbors.



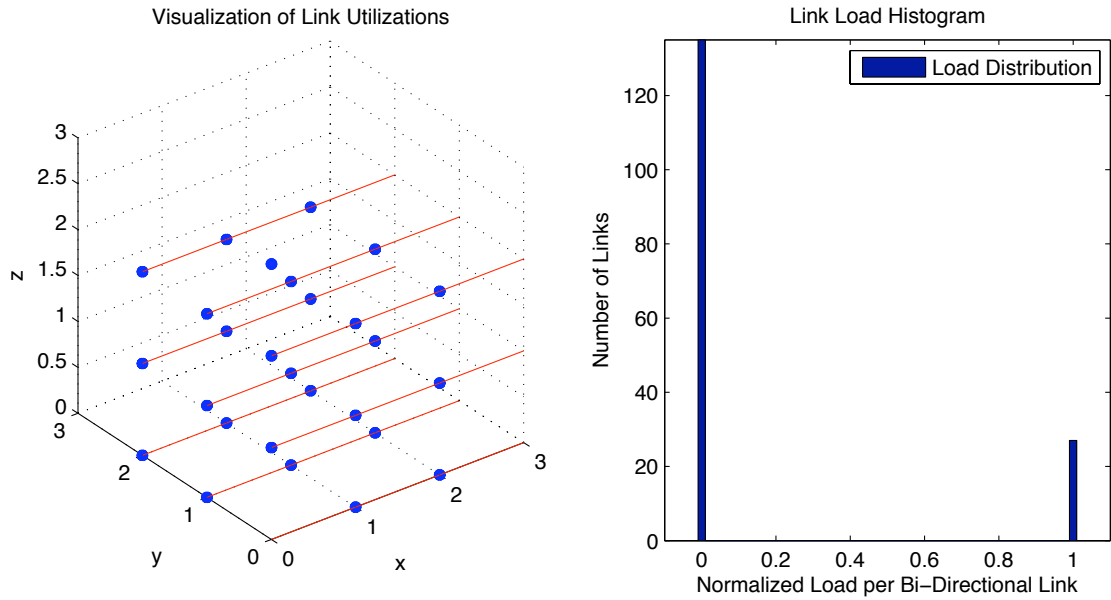
**Figure E.34:** Results of Enhanced CQR Routing using the uniform random traffic demand (UR). For each node, a value of 0-9 was assigned as having 0-9 individual traffic demands.



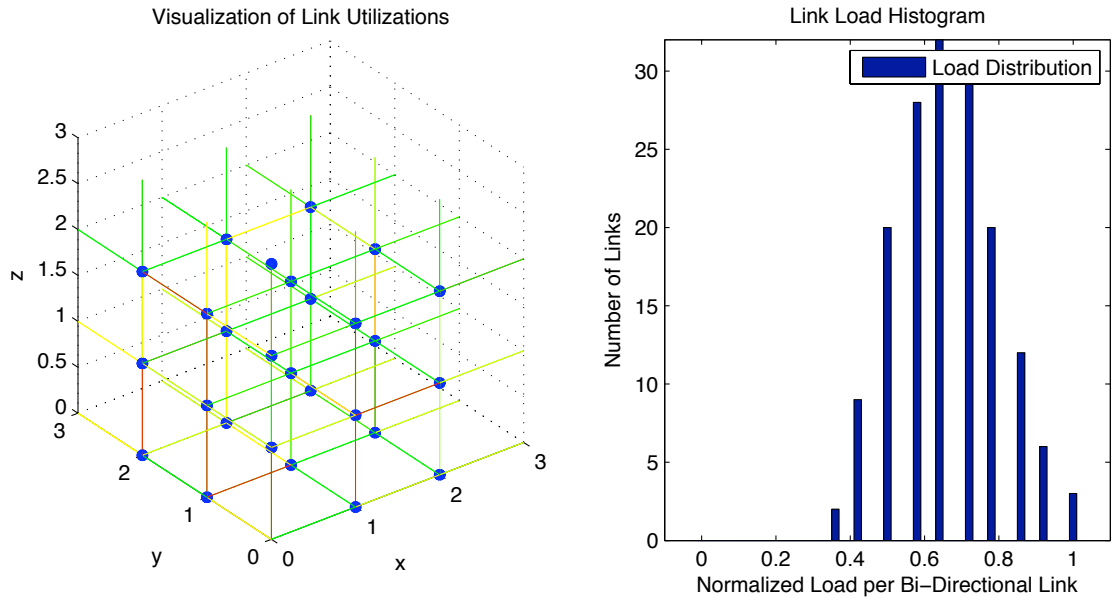
**Figure E.35:** Results of CQR Routing using the bit complement traffic demand (BC). Two individual traffic demands were assigned for possible values described in Table 2.1.



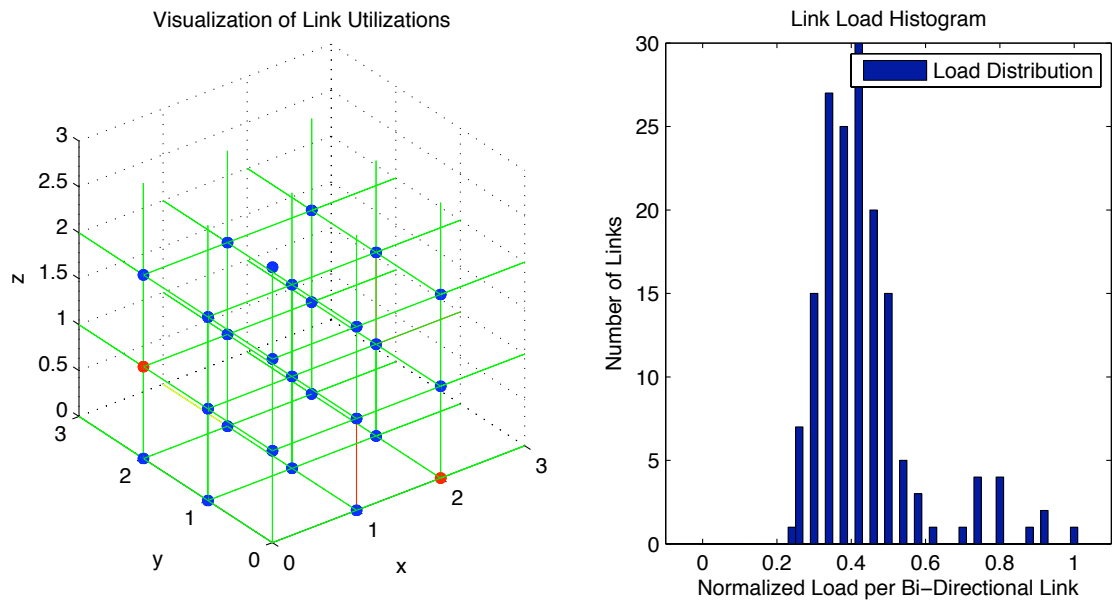
**Figure E.36:** Results of Enhanced CQR Routing using the transpose traffic demand (TP). Two individual traffic demands were assigned for each permutation of all  $x, y, z$  values.



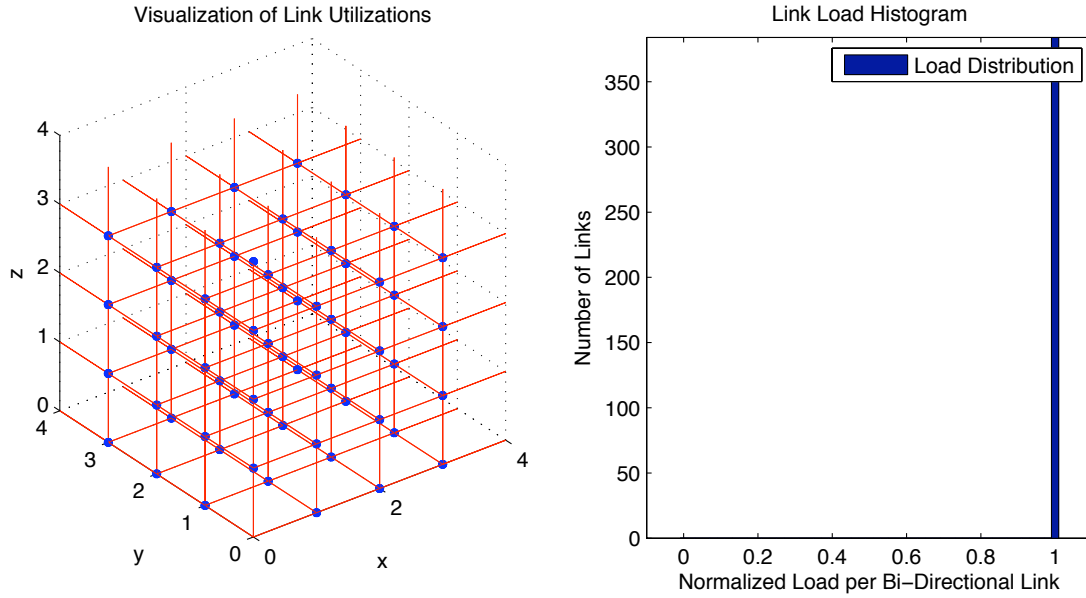
**Figure E.37:** Results of Enhanced CQR Routing using the tornado traffic demand (TOR). Three individual traffic demands were assigned for possible values described in Table 2.1.



**Figure E.38:** Results of Enhanced CQR Routing using the flood traffic demand (FL). Two individual traffic demands were assigned to each node from each node.



**Figure E.39:** Results of Enhanced CQR Routing using the flood traffic demand (FL) with Hot-spots. Two individual traffic demands were assigned to each node from each node. Five percent of the nodes were made hot-spots, and they are marked as red in the figure.

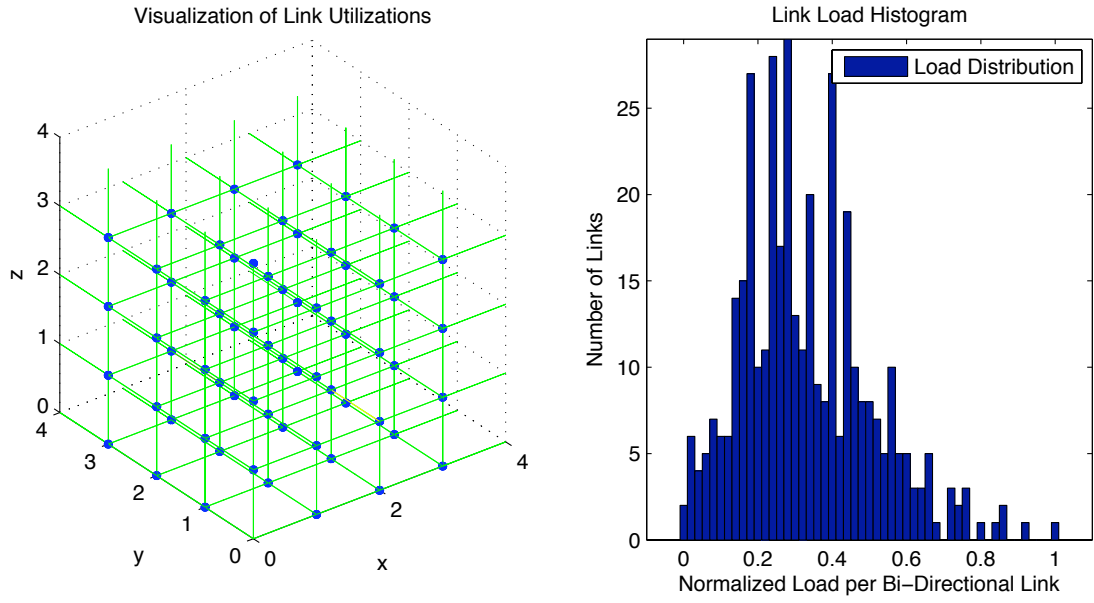


**Figure E.40:** Results of Dimension Ordered Routing (DOR) using the nearest neighbor traffic demand (NN). Each node had two individual traffic demands to each of its six neighbors.

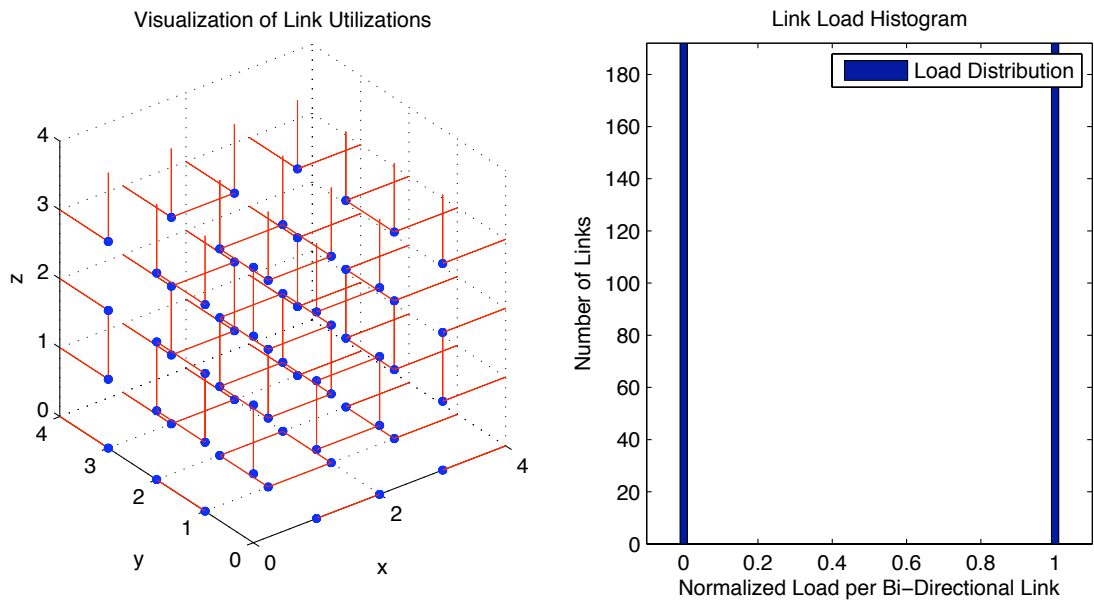
## E.2 Results from the 4-ary 3-cube Topology

### E.2.1 Static Algorithms

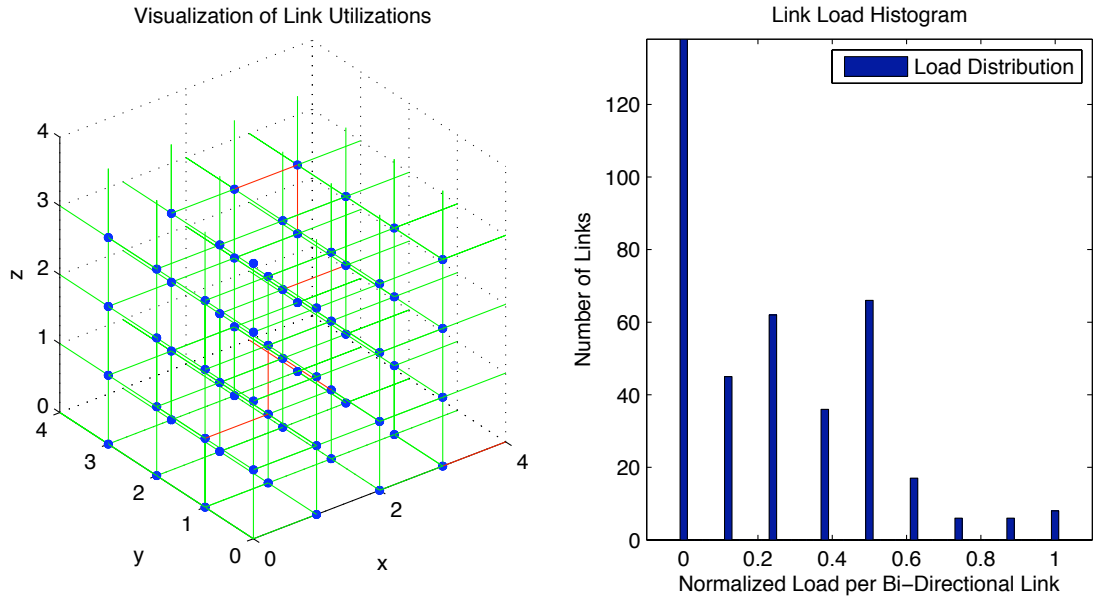




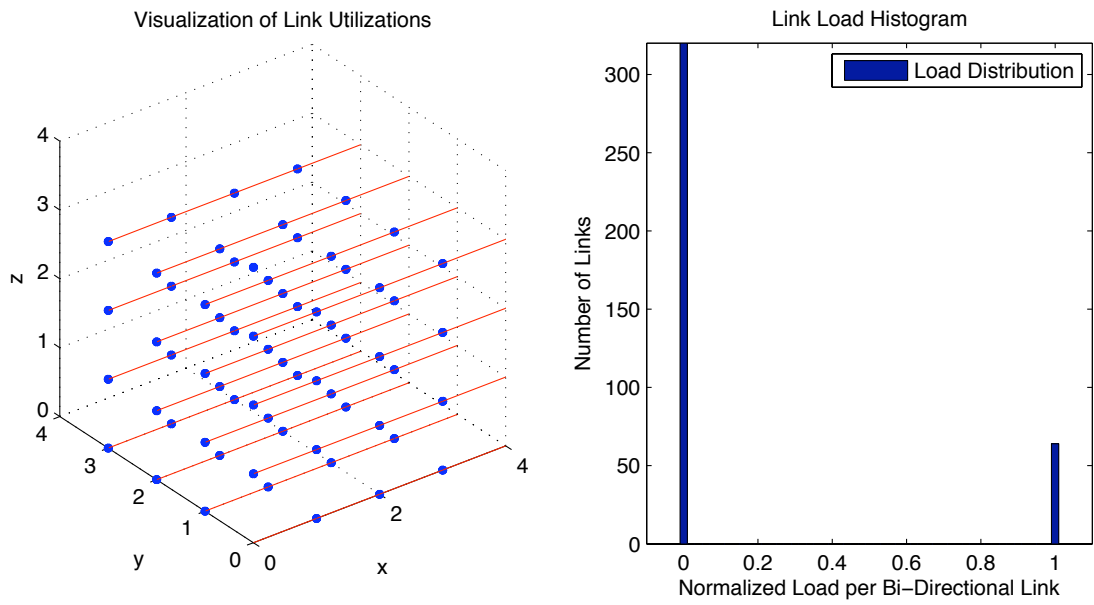
**Figure E.41:** Results of Dimension Ordered Routing (DOR) using the uniform random traffic demand (UR). For each node, a value of 0-9 was assigned as having 0-9 individual traffic demands.



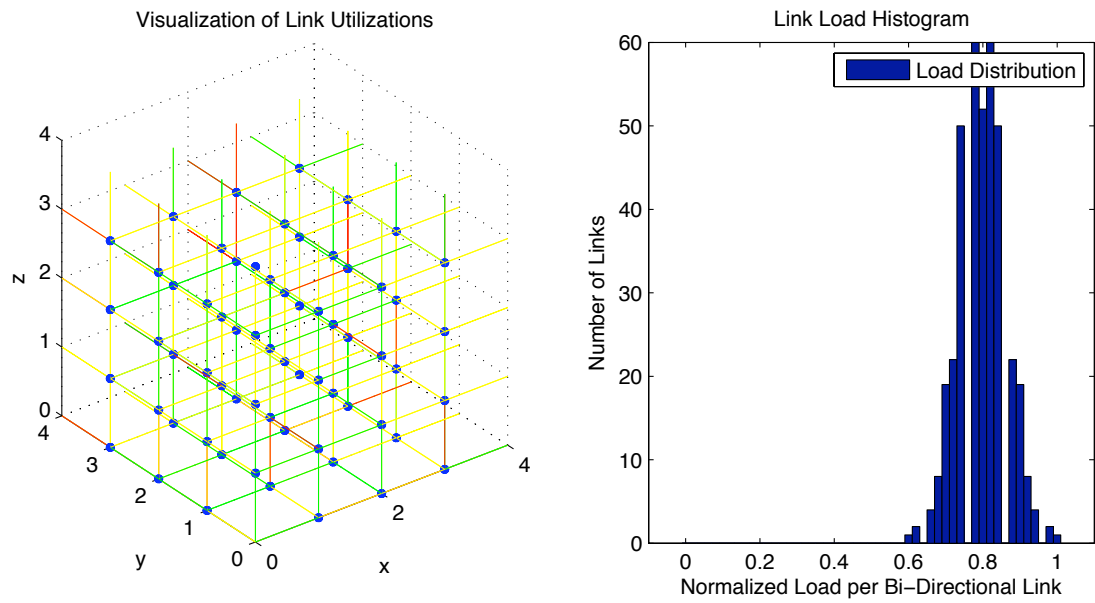
**Figure E.42:** Results of Dimension Ordered Routing (DOR) using the bit complement traffic demand (BC). Two individual traffic demands were assigned for possible values described in Table 2.1.



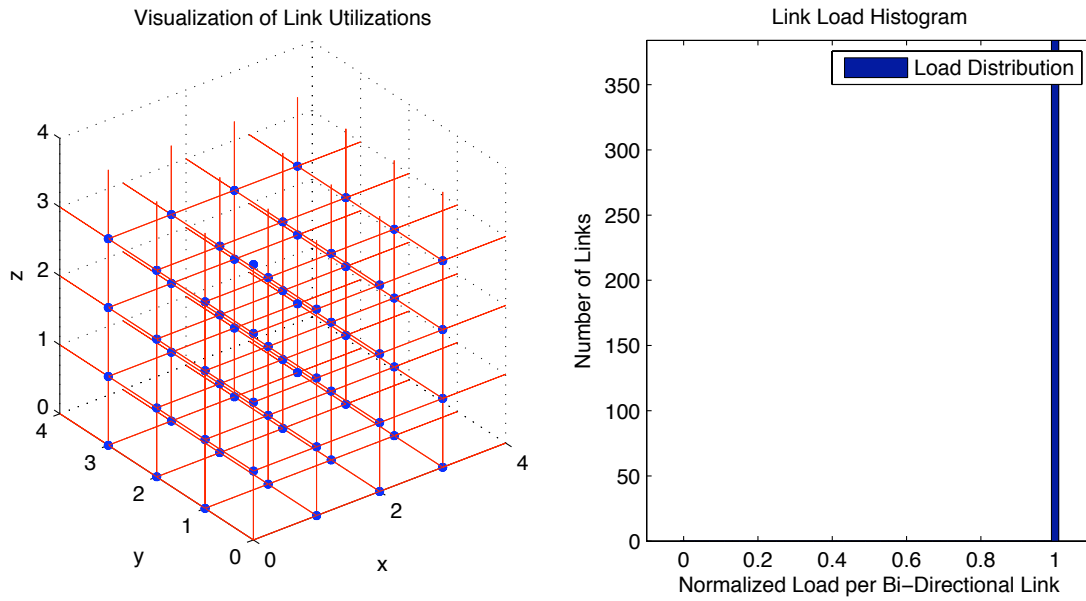
**Figure E.43:** Results of Dimension Ordered Routing (DOR) using the transpose traffic demand (TP). Two individual traffic demands were assigned for each permutation of all  $x, y, z$  values.



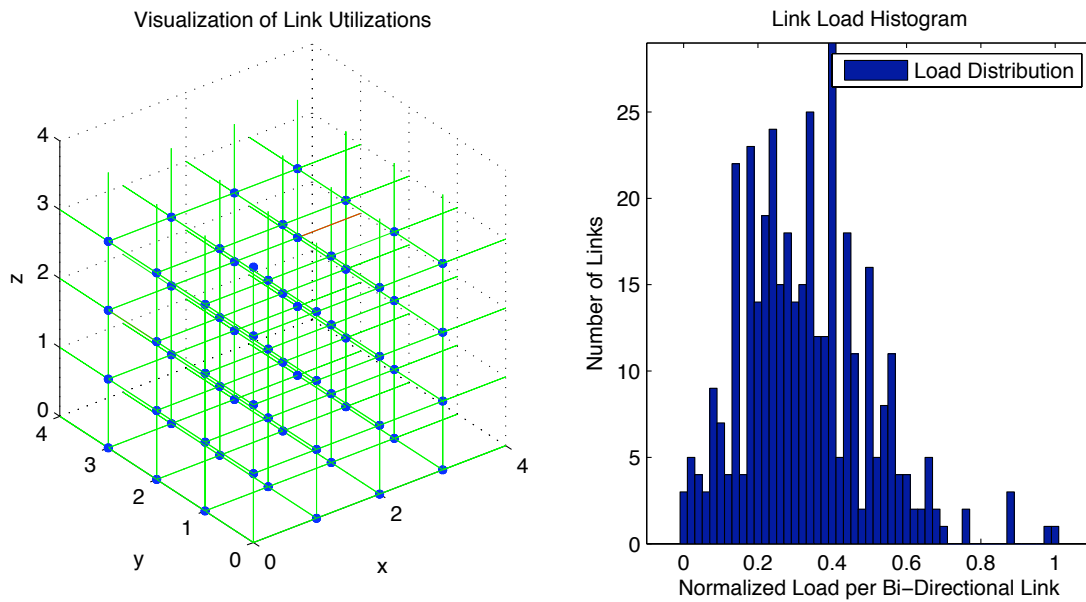
**Figure E.44:** Results of Dimension Ordered Routing (DOR) using the tornado traffic demand (TOR). Three individual traffic demands were assigned for possible values described in Table 2.1.



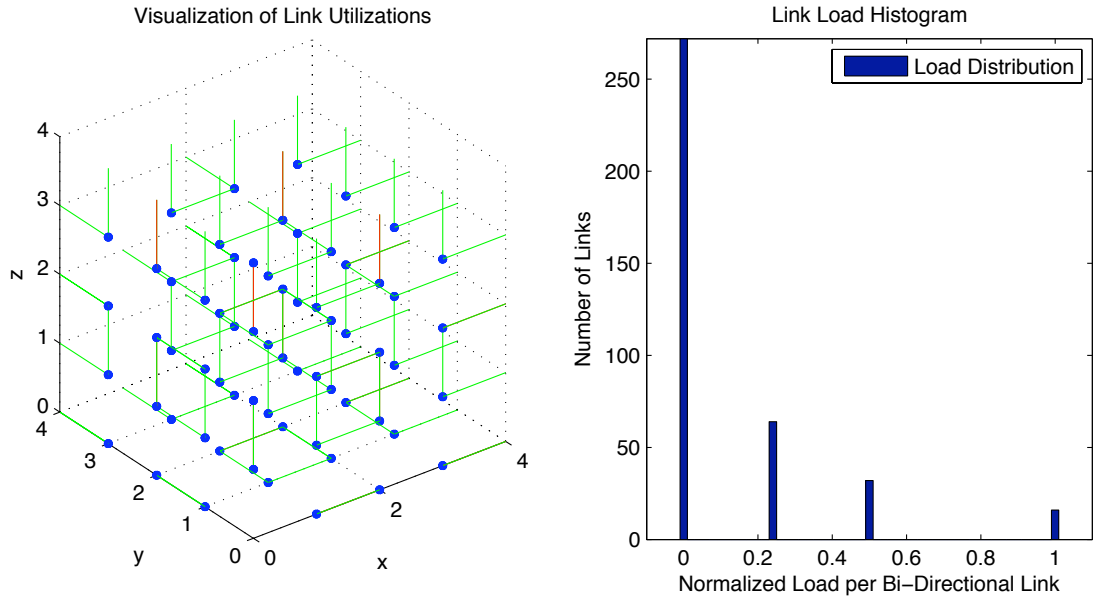
**Figure E.45:** Results of Dimension Ordered Routing (DOR) using the flood traffic demand (FL). Two individual traffic demands were assigned to each node from each node.



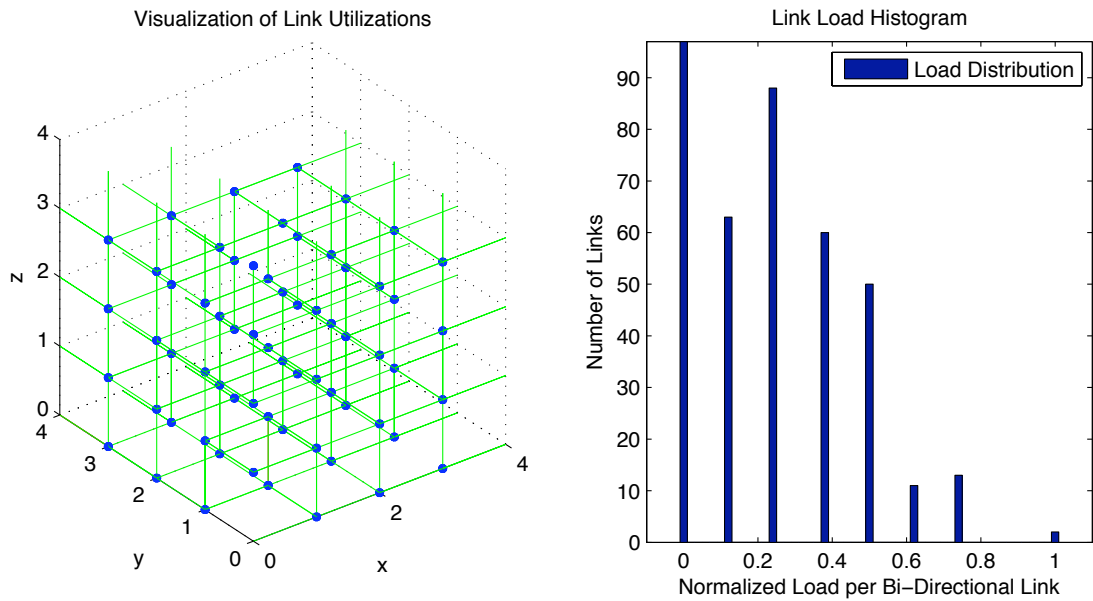
**Figure E.46:** Results of Direction Ordered Routing (DIR) using the nearest neighbor traffic demand (NN). Each node had two individual traffic demands to each of its six neighbors.



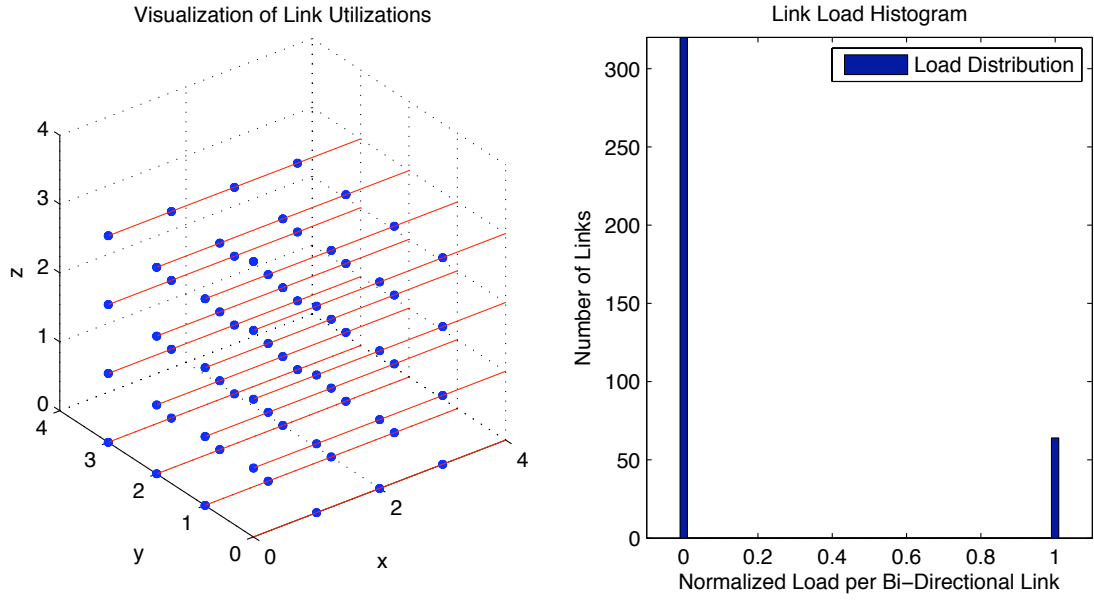
**Figure E.47:** Results of Direction Ordered Routing (DIR) using the uniform random traffic demand (UR). For each node, a value of 0-9 was assigned as having 0-9 individual traffic demands.



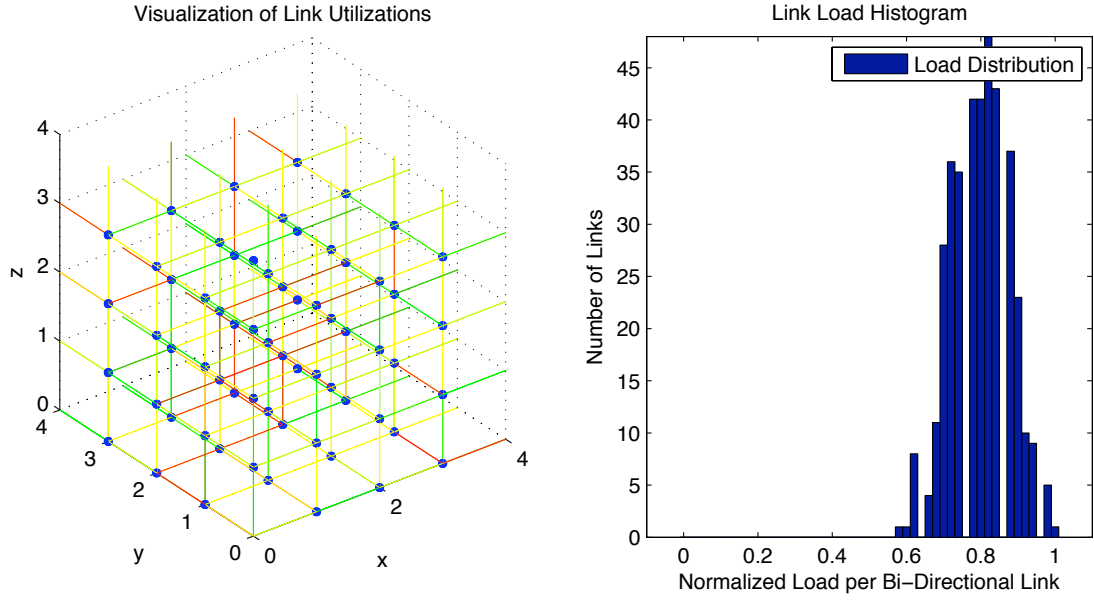
**Figure E.48:** Results of Direction Ordered Routing (DIR) using the bit complement traffic demand (BC). Two individual traffic demands were assigned for possible values described in Table 2.1.



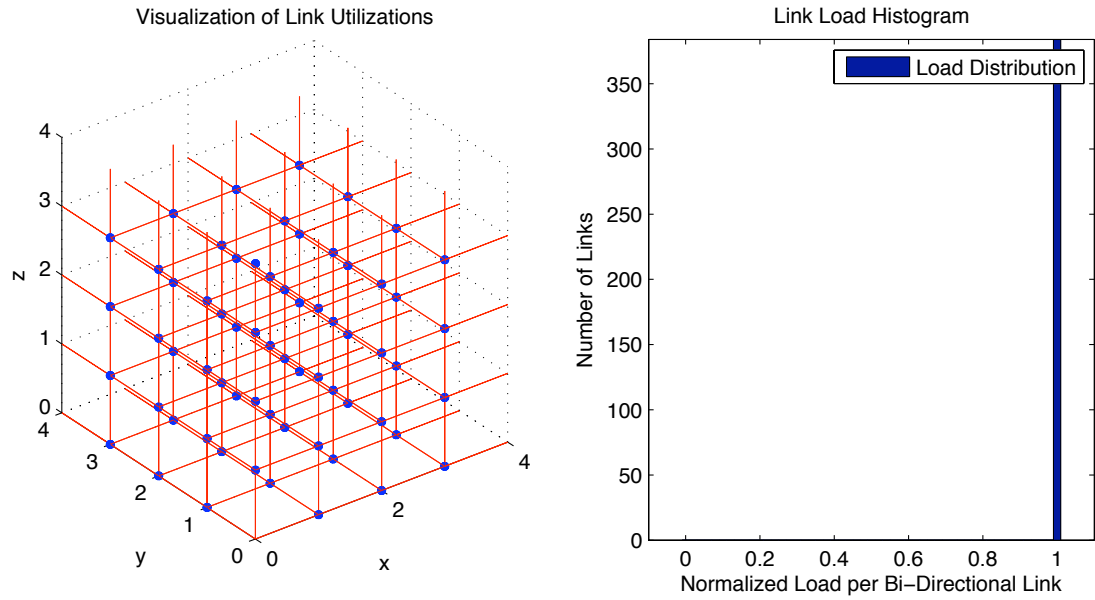
**Figure E.49:** Results of Direction Ordered Routing (DIR) using the transpose traffic demand (TP). Two individual traffic demands were assigned for each permutation of all  $x, y, z$  values.



**Figure E.50:** Results of Direction Ordered Routing (DIR) using the tornado traffic demand (TOR). Three individual traffic demands were assigned for possible values described in Table 2.1.

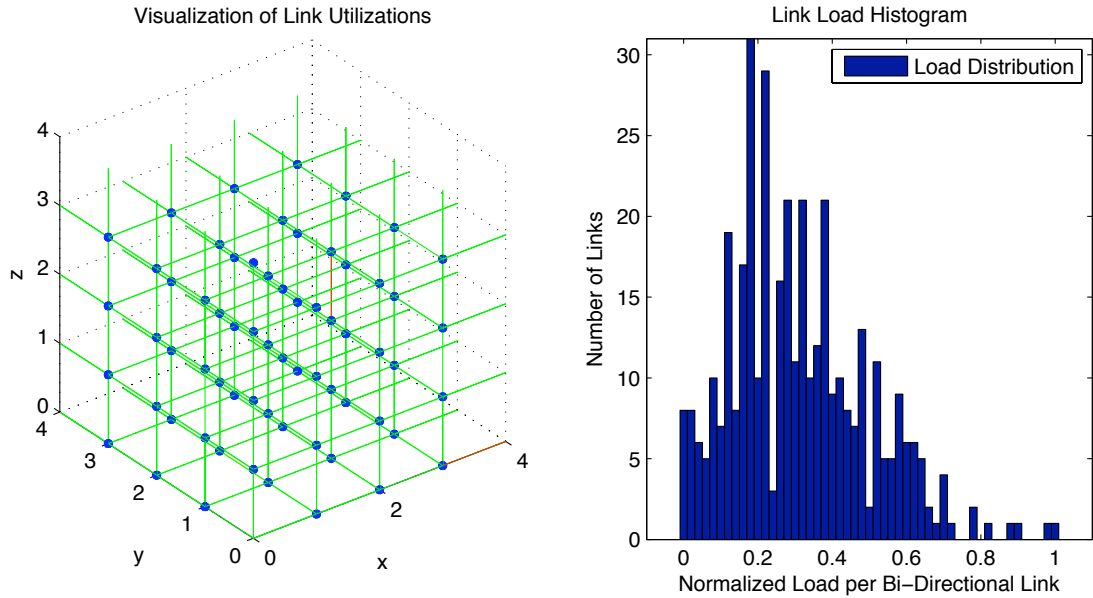


**Figure E.51:** Results of Direction Ordered Routing (DIR) using the flood traffic demand (FL). Two individual traffic demands were assigned to each node from each node.

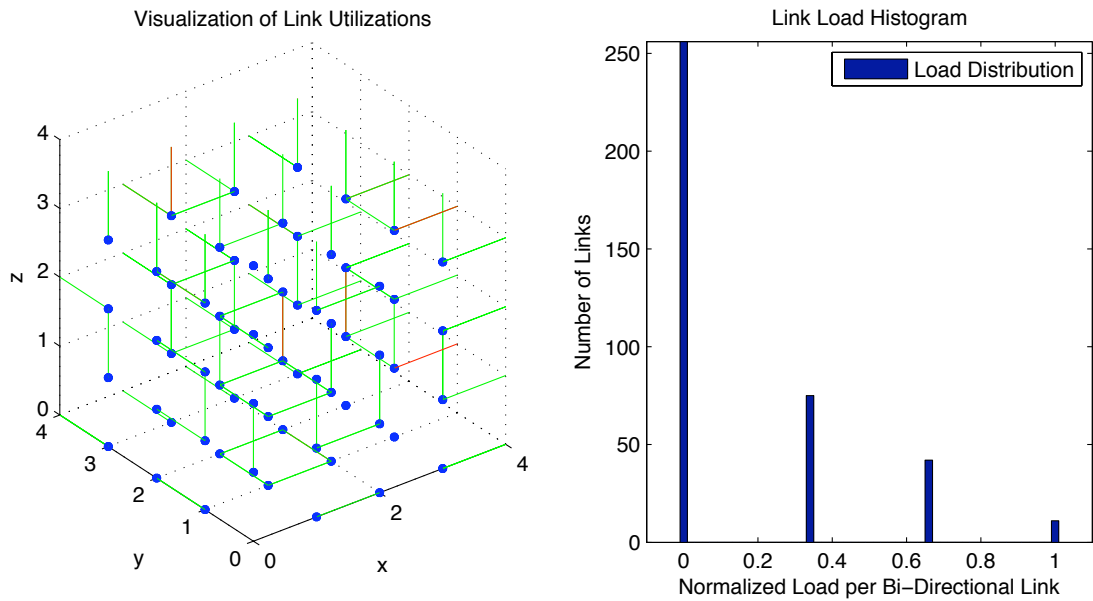


**Figure E.52:** Results of Minimal Oblivious Routing (MO) using the nearest neighbor traffic demand (NN). Each node had two individual traffic demands to each of its six neighbors.

## E.2.2 Oblivious Algorithms

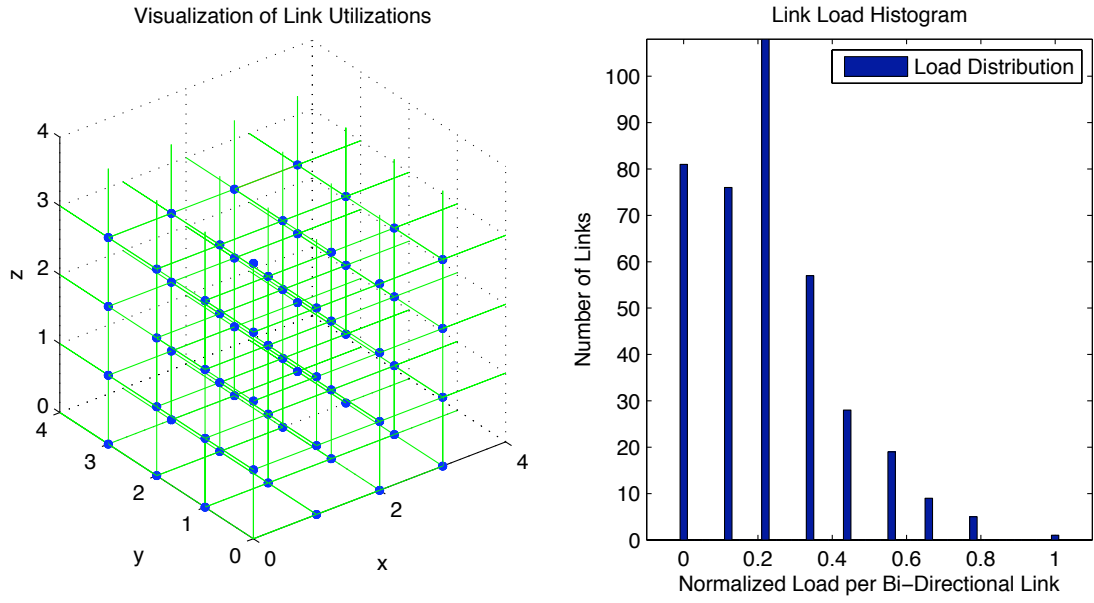


**Figure E.53:** Results of Minimal Oblivious Routing (MO) using the uniform random traffic demand (UR). For each node, a value of 0-9 was assigned as having 0-9 individual traffic demands.

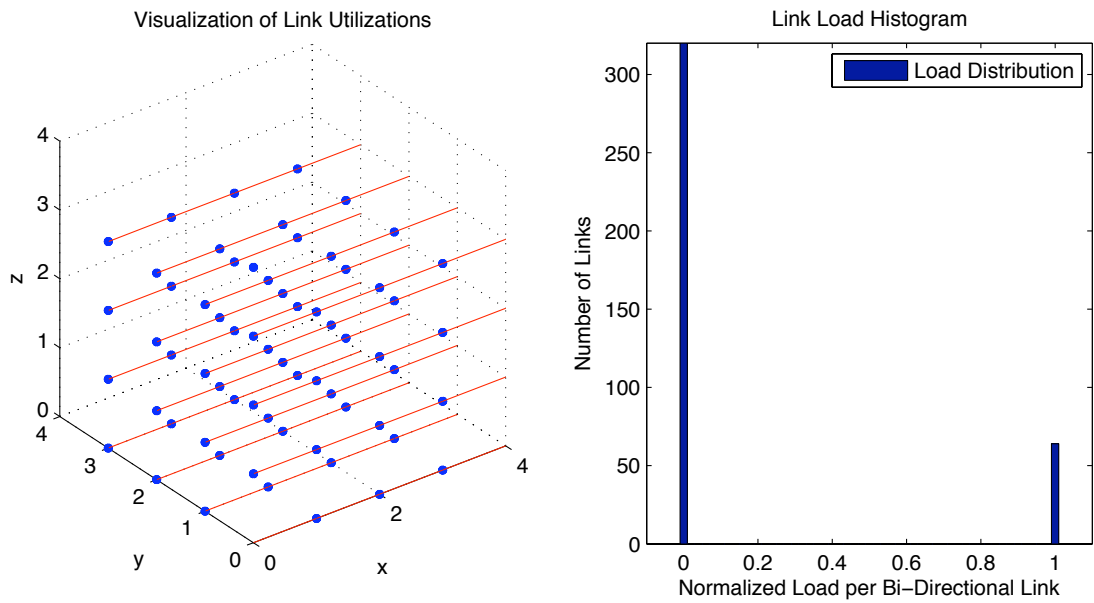


**Figure E.54:** Results of Minimal Oblivious Routing (MO) using the bit complement traffic demand (BC). Two individual traffic demands were assigned for possible values described in Table 2.1.

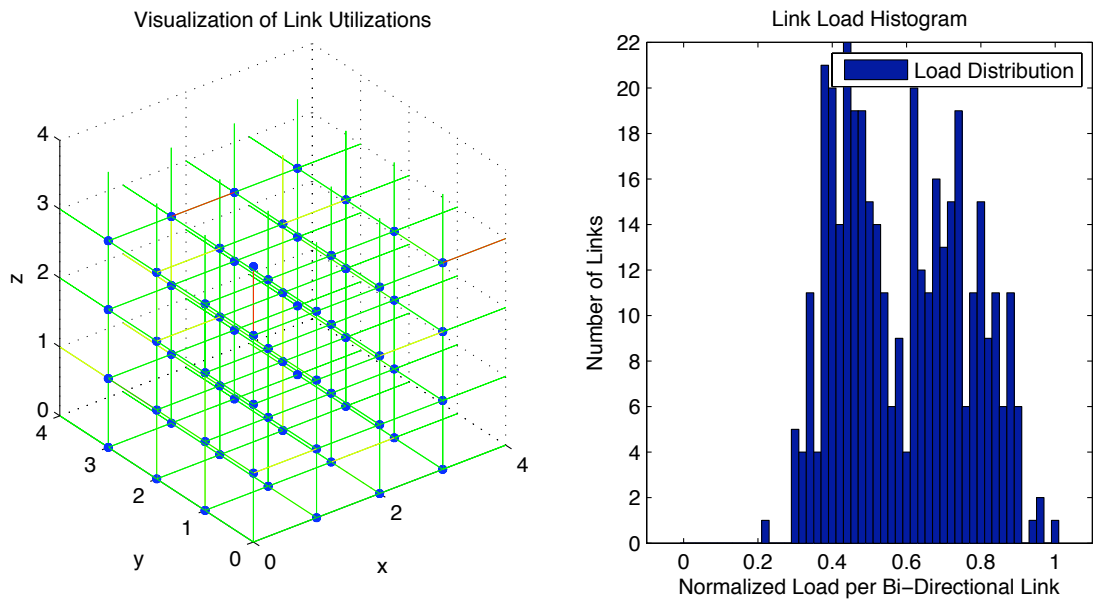




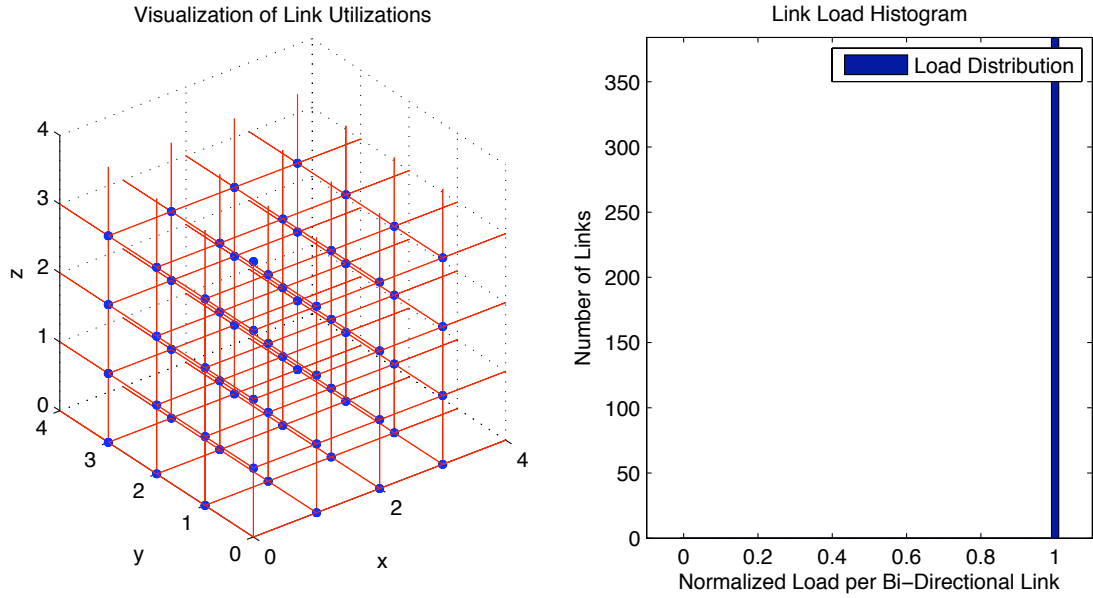
**Figure E.55:** Results of Minimal Oblivious Routing (MO) using the transpose traffic demand (TP). Two individual traffic demands were assigned for each permutation of all  $x, y, z$  values.



**Figure E.56:** Results of Minimal Oblivious Routing (MO) using the tornado traffic demand (TOR). Three individual traffic demands were assigned for possible values described in Table 2.1.

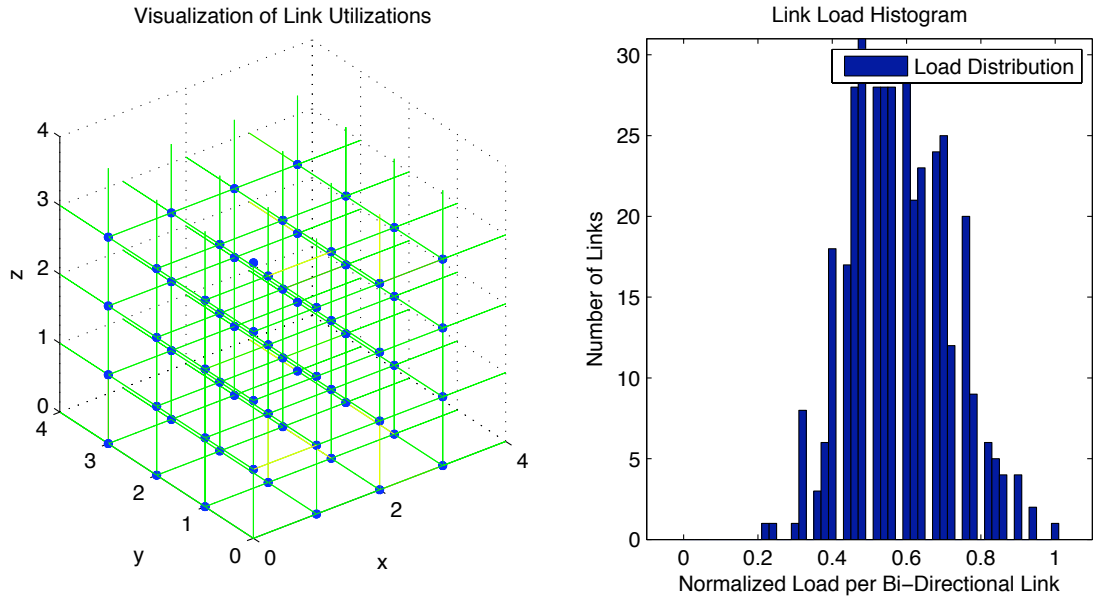


**Figure E.57:** Results of Minimal Oblivious Routing (MO) using the flood traffic demand (FL). Two individual traffic demands were assigned to each node from each node.

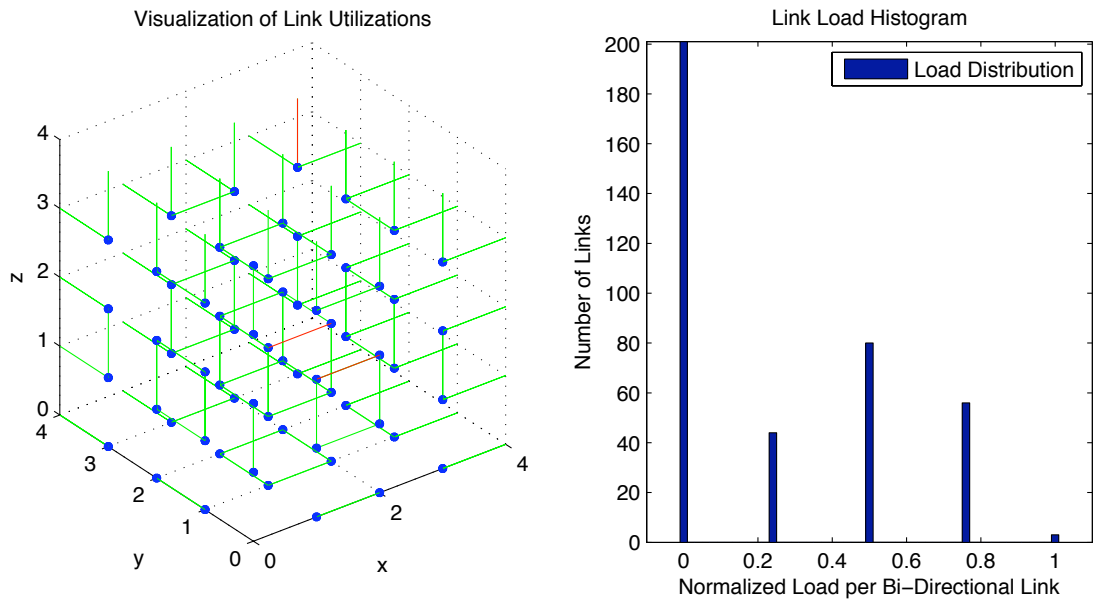


**Figure E.58:** Results of Minimal Adaptive Routing (MA) using the nearest neighbor traffic demand (NN). Each node had two individual traffic demands to each of its six neighbors.

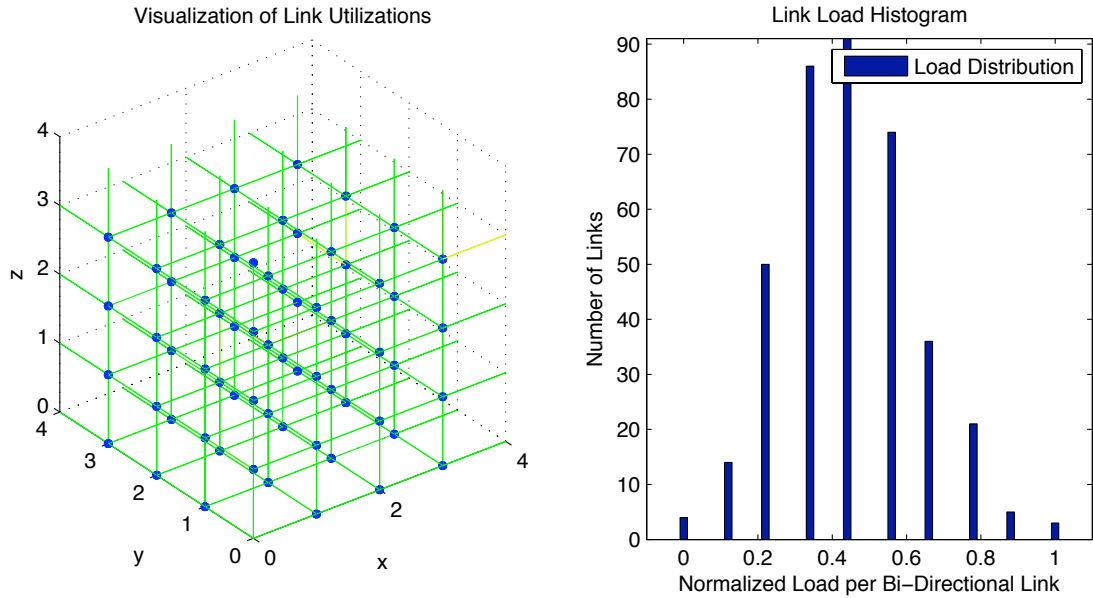
### E.2.3 Adaptive Algorithms



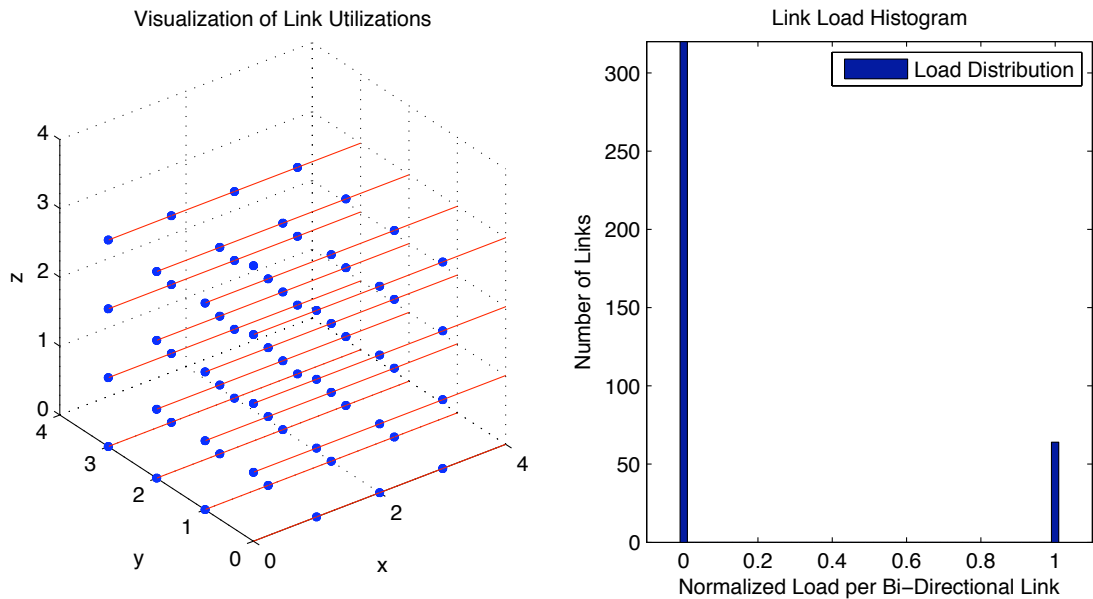
**Figure E.59:** Results of Minimal Adaptive Routing (MA) using the uniform random traffic demand (UR). For each node, a value of 0-9 was assigned as having 0-9 individual traffic demands.



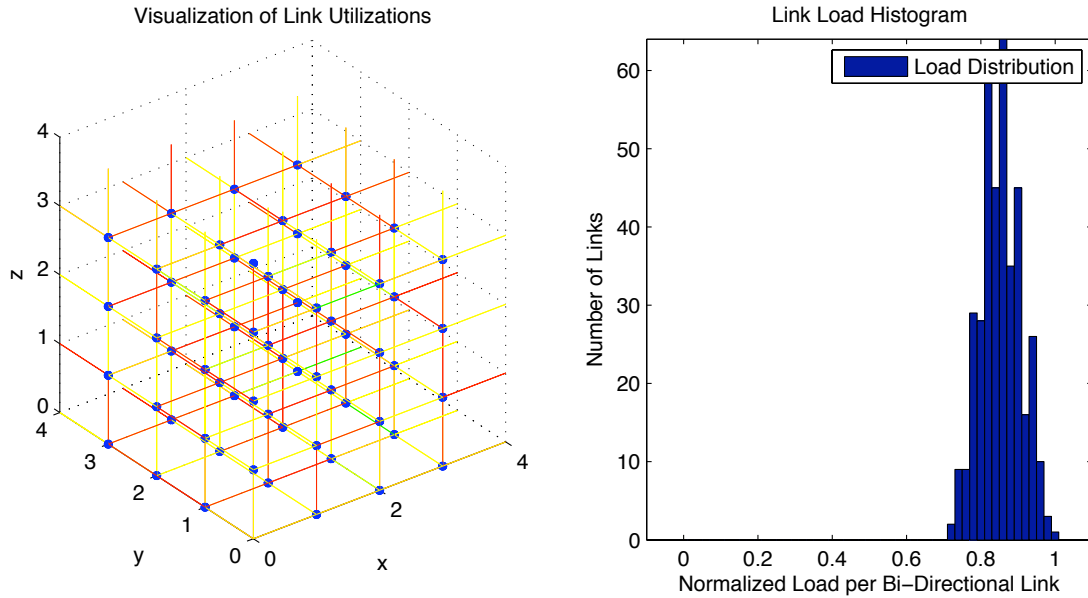
**Figure E.60:** Results of Minimal Adaptive Routing (MA) using the bit complement traffic demand (BC). Two individual traffic demands were assigned for possible values described in Table 2.1.



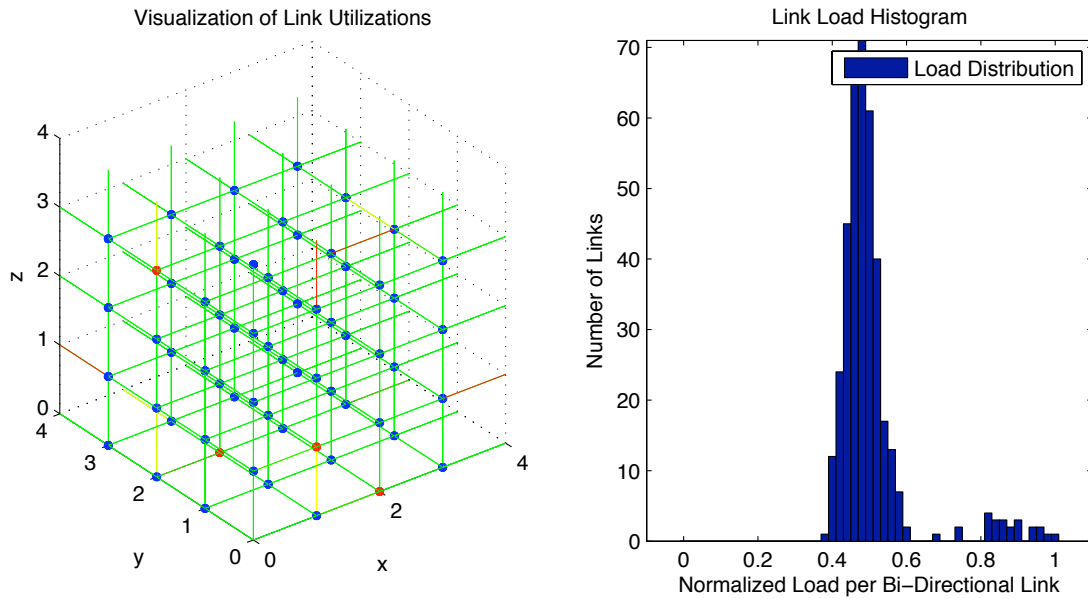
**Figure E.61:** Results of Minimal Adaptive Routing (MA) using the transpose traffic demand (TP). Two individual traffic demands were assigned for each permutation of all  $x, y, z$  values.



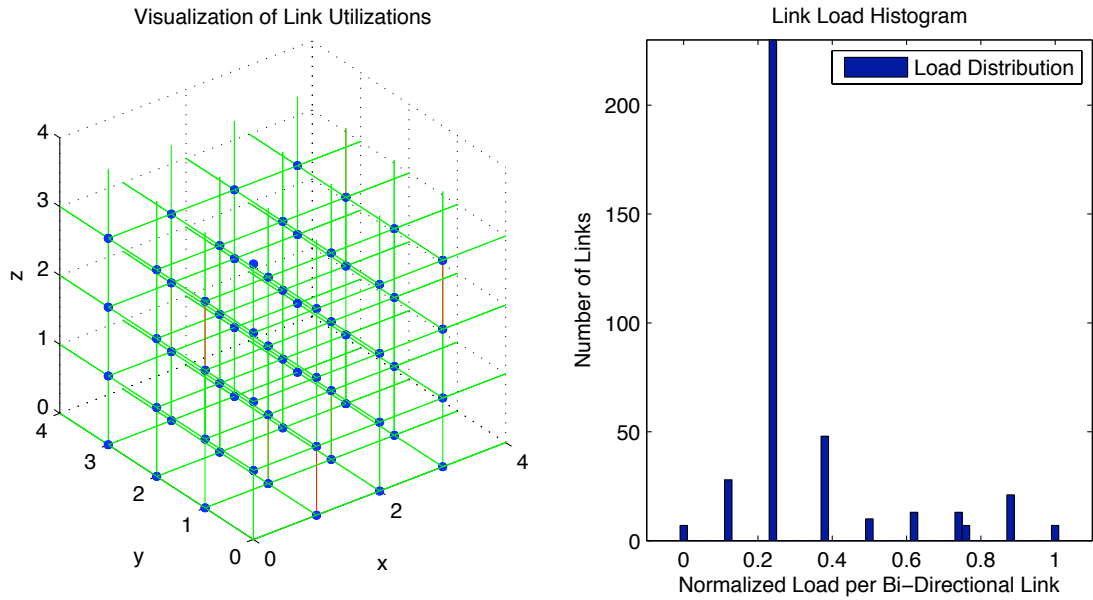
**Figure E.62:** Results of Minimal Adaptive Routing (MA) using the tornado traffic demand (TOR). Three individual traffic demands were assigned for possible values described in Table 2.1.



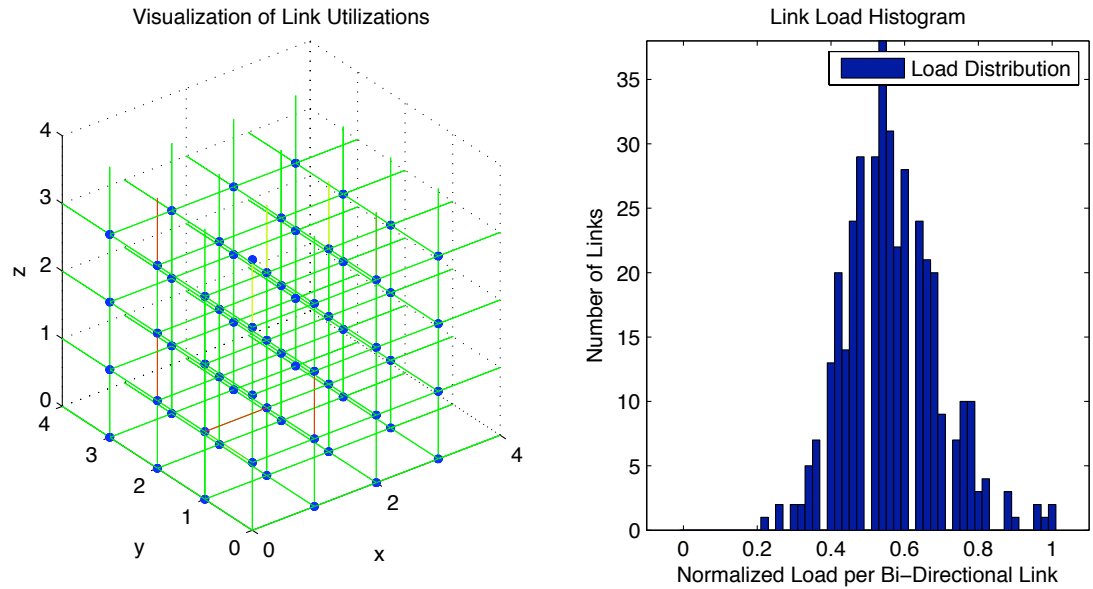
**Figure E.63:** Results of Minimal Adaptive Routing (MA) using the flood traffic demand (FL). Two individual traffic demands were assigned to each node from each node.



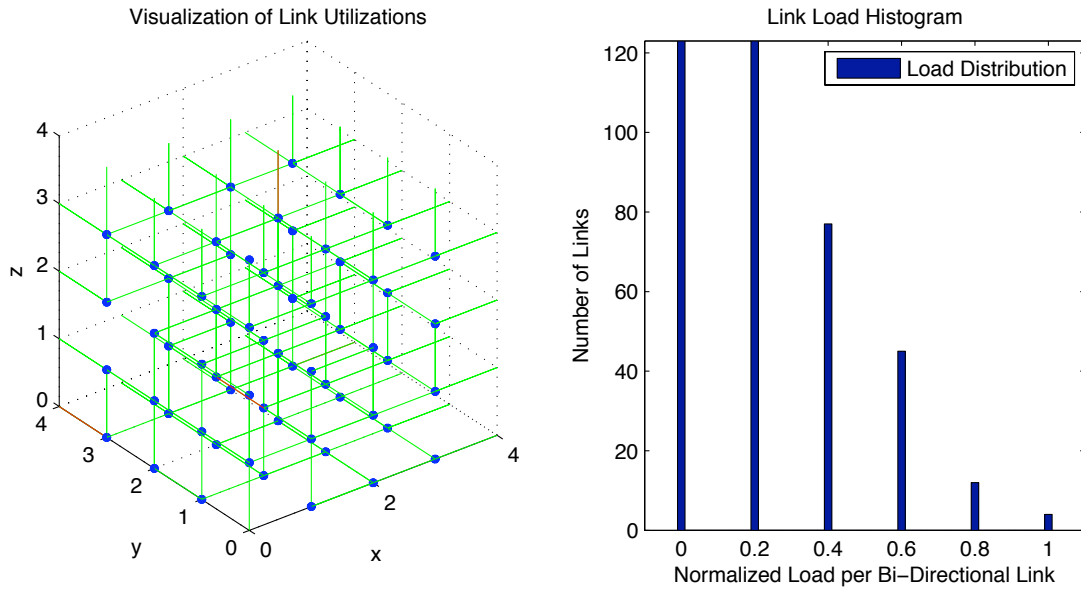
**Figure E.64:** Results of Minimal Adaptive Routing (MA) using the flood traffic demand (FL) with Hot-spots. Two individual traffic demands were assigned to each node from each node. Five percent of the nodes were made hot-spots, and they are marked as red in the figure.



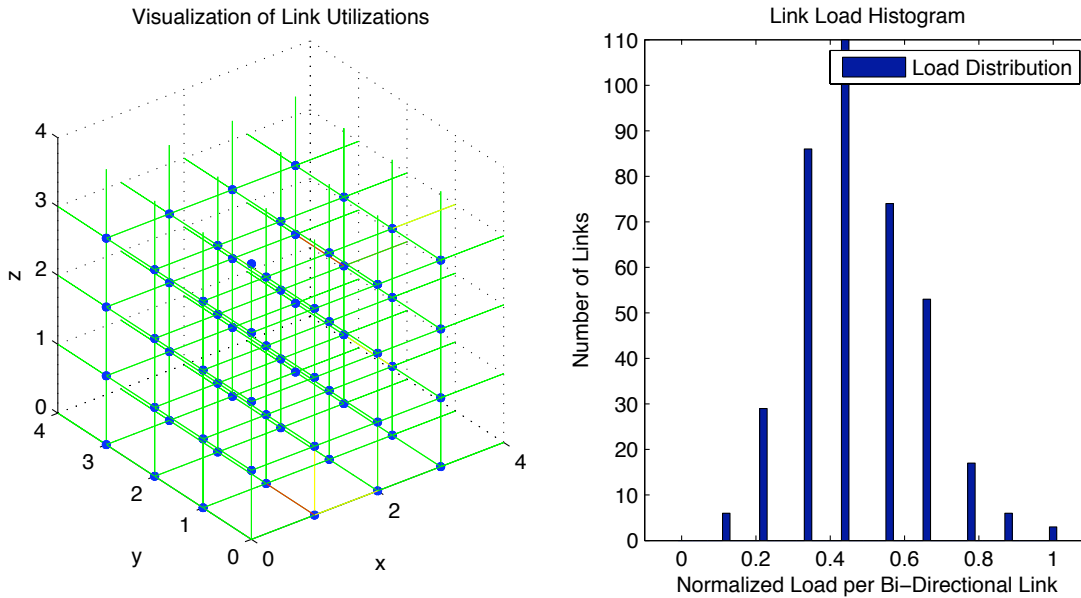
**Figure E.65:** Results of CQR Routing using the nearest neighbor traffic demand (NN). Each node had two individual traffic demands to each of its six neighbors.



**Figure E.66:** Results of CQR Routing using the uniform random traffic demand (UR). For each node, a value of 0-9 was assigned as having 0-9 individual traffic demands.

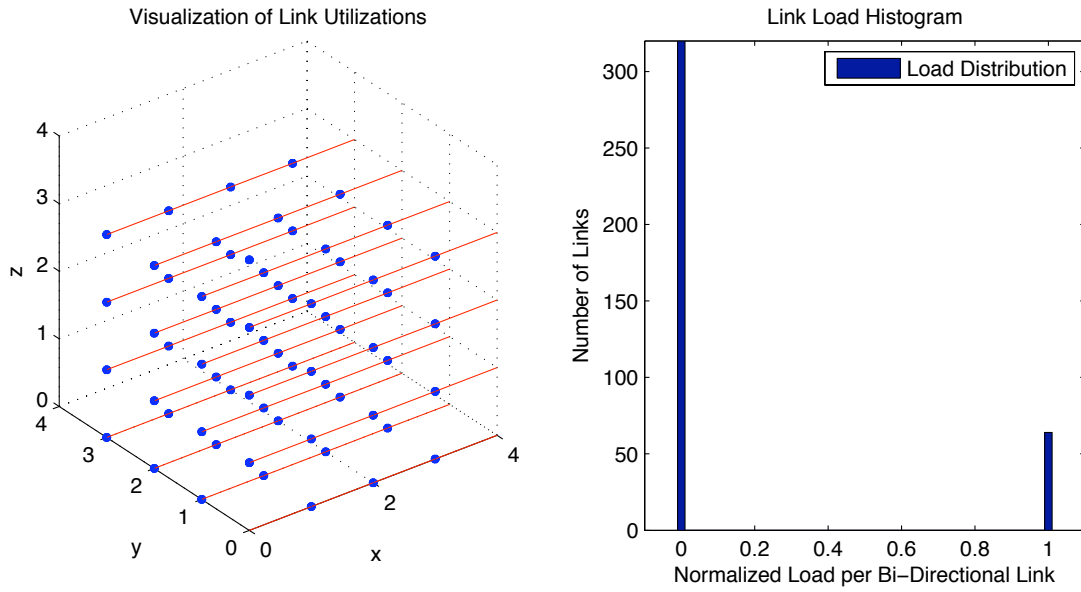


**Figure E.67:** Results of CQR Routing using the bit complement traffic demand (BC). Two individual traffic demands were assigned for possible values described in Table 2.1.

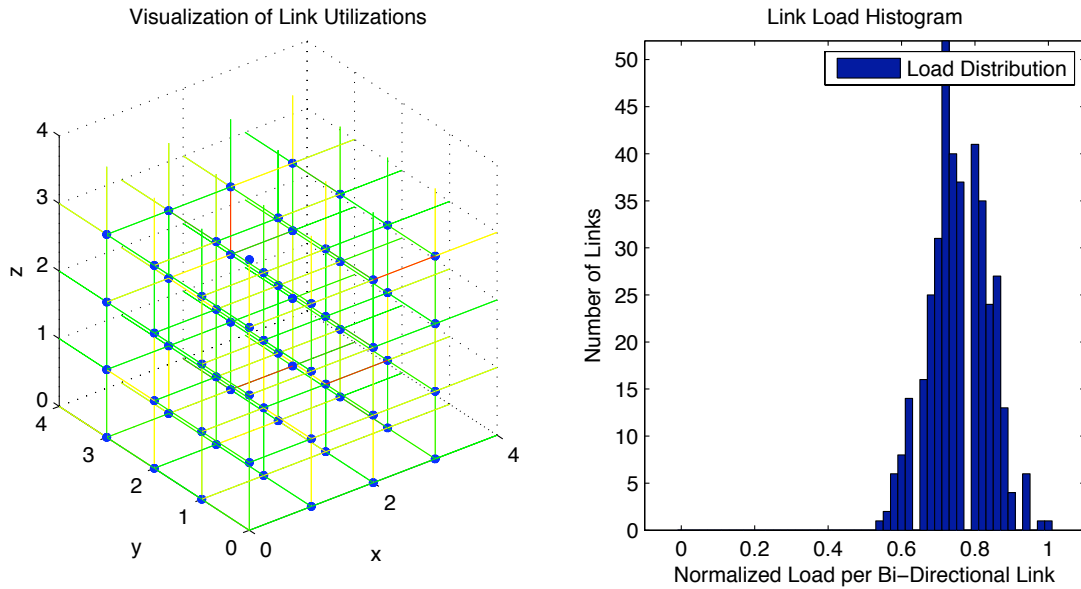


**Figure E.68:** Results of CQR Routing using the transpose traffic demand (TP). Two individual traffic demands were assigned for each permutation of all  $x, y, z$  values.

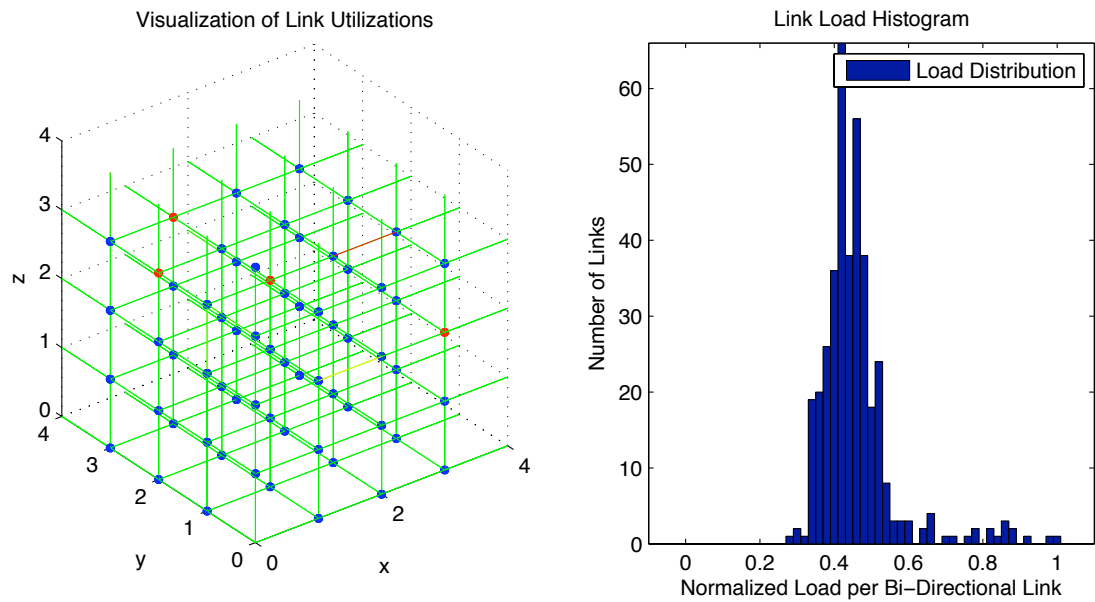




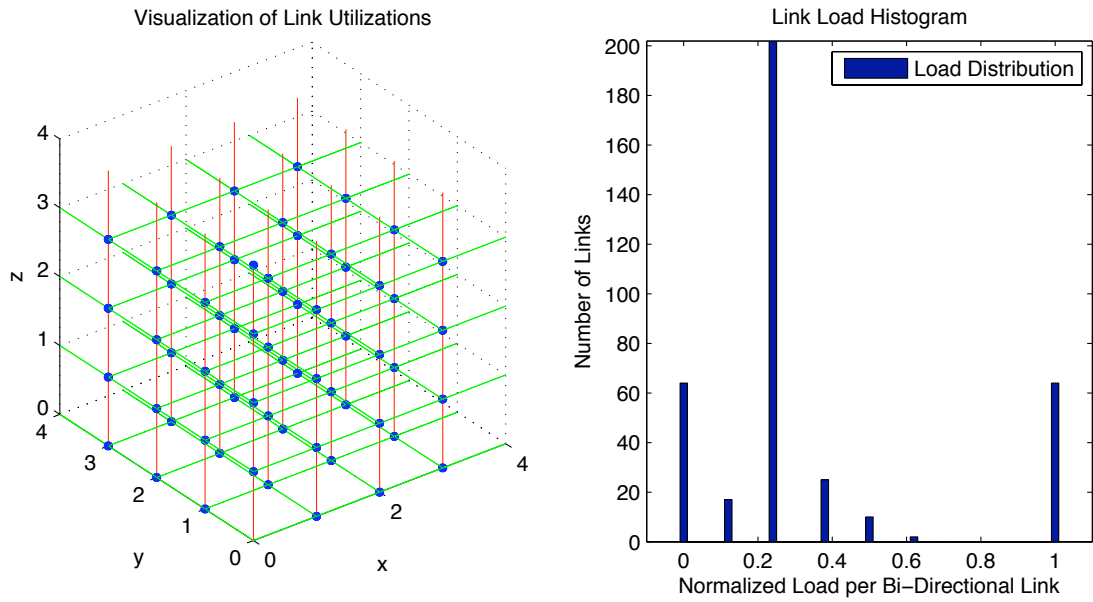
**Figure E.69:** Results of CQR Routing using the tornado traffic demand (TOR). Three individual traffic demands were assigned for possible values described in Table 2.1.



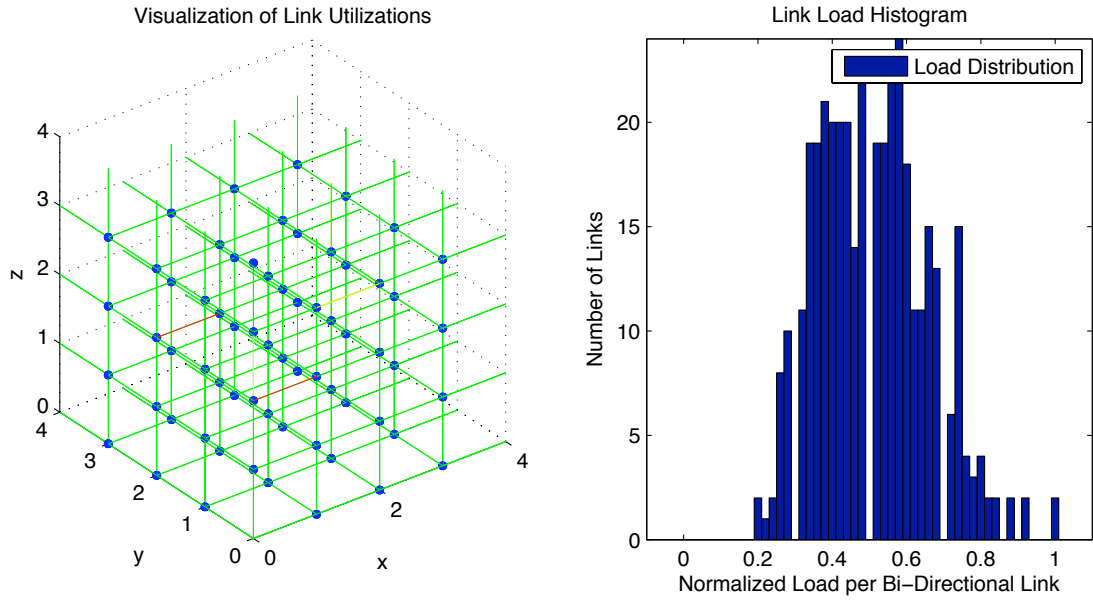
**Figure E.70:** Results of CQR Routing using the flood traffic demand (FL). Two individual traffic demands were assigned to each node from each node.



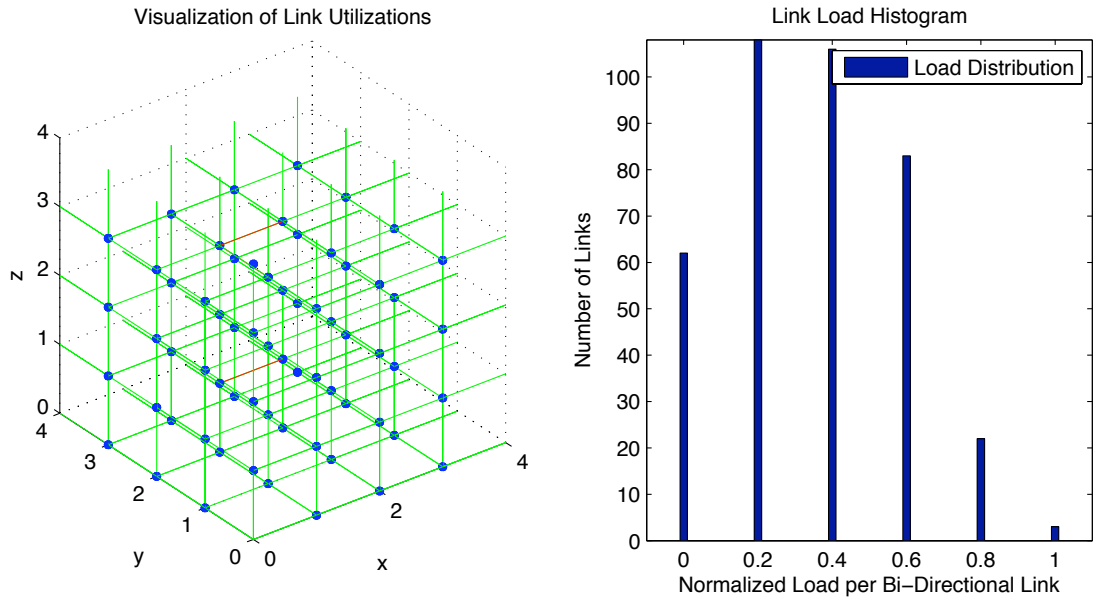
**Figure E.71:** Results of CQR Routing using the flood traffic demand (FL) with Hot-spots. Two individual traffic demands were assigned to each node from each node. Five percent of the nodes were made hot-spots, and they are marked as red in the figure.



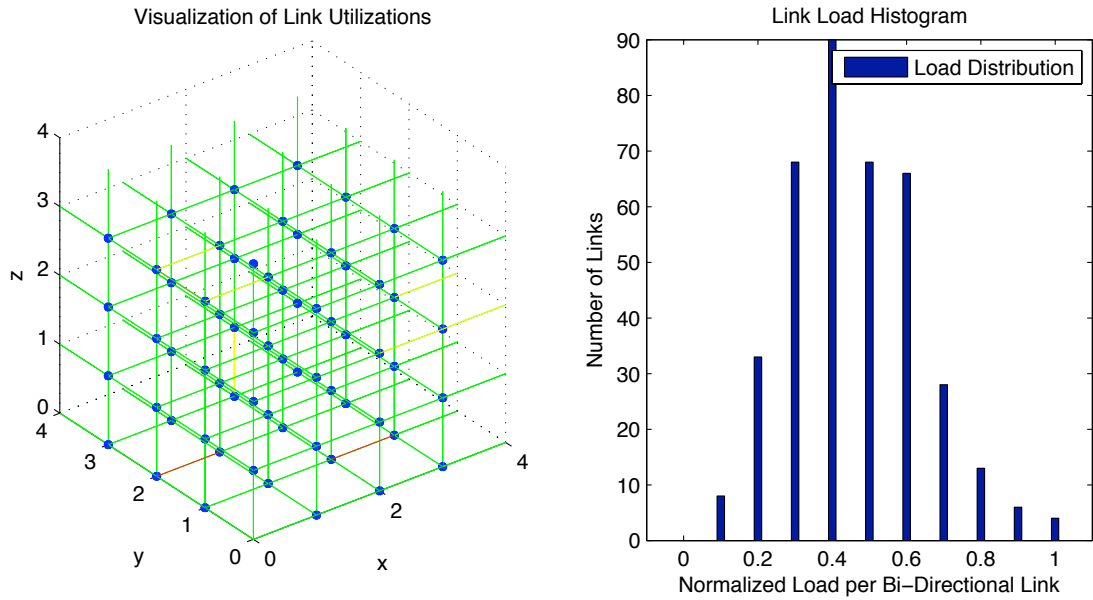
**Figure E.72:** Results of Enhanced CQR Routing using the nearest neighbor traffic demand (NN). Each node had two individual traffic demands to each of its six neighbors.



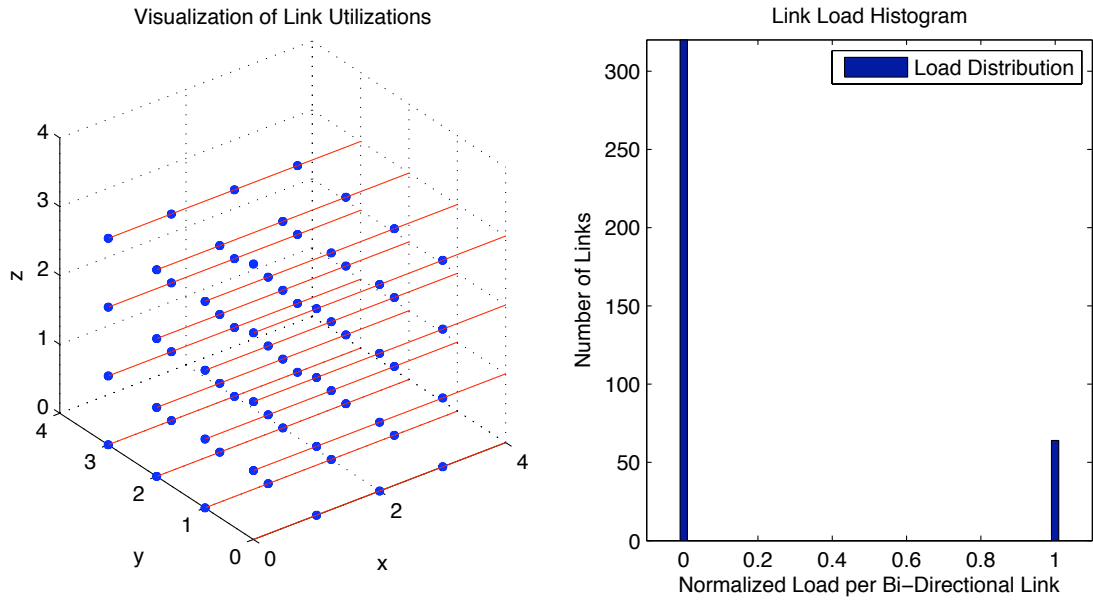
**Figure E.73:** Results of Enhanced CQR Routing using the uniform random traffic demand (UR). For each node, a value of 0-9 was assigned as having 0-9 individual traffic demands.



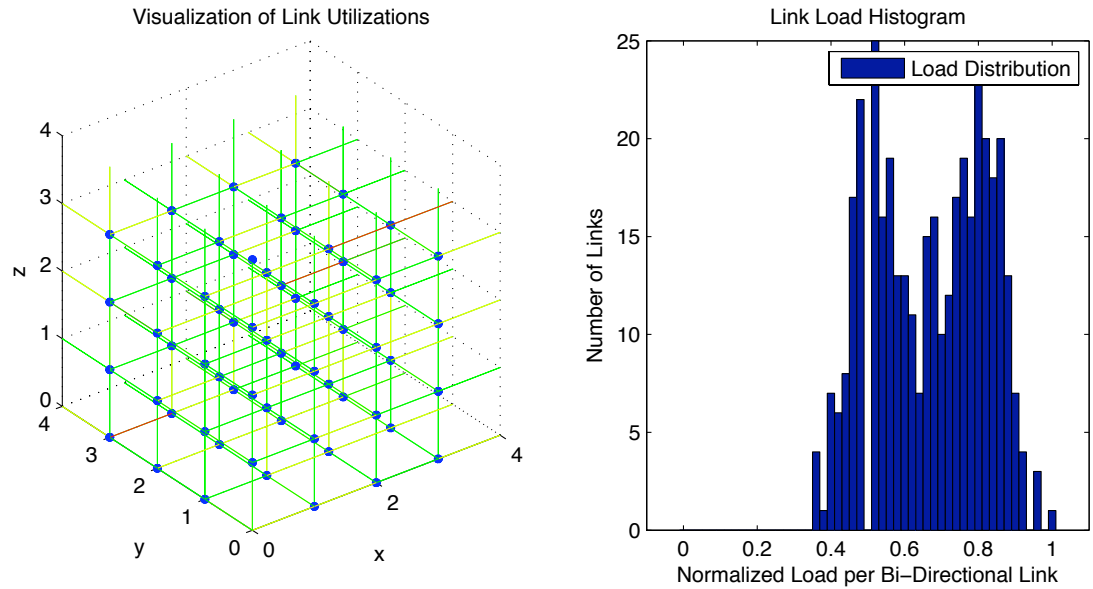
**Figure E.74:** Results of CQR Routing using the bit complement traffic demand (BC). Two individual traffic demands were assigned for possible values described in Table 2.1.



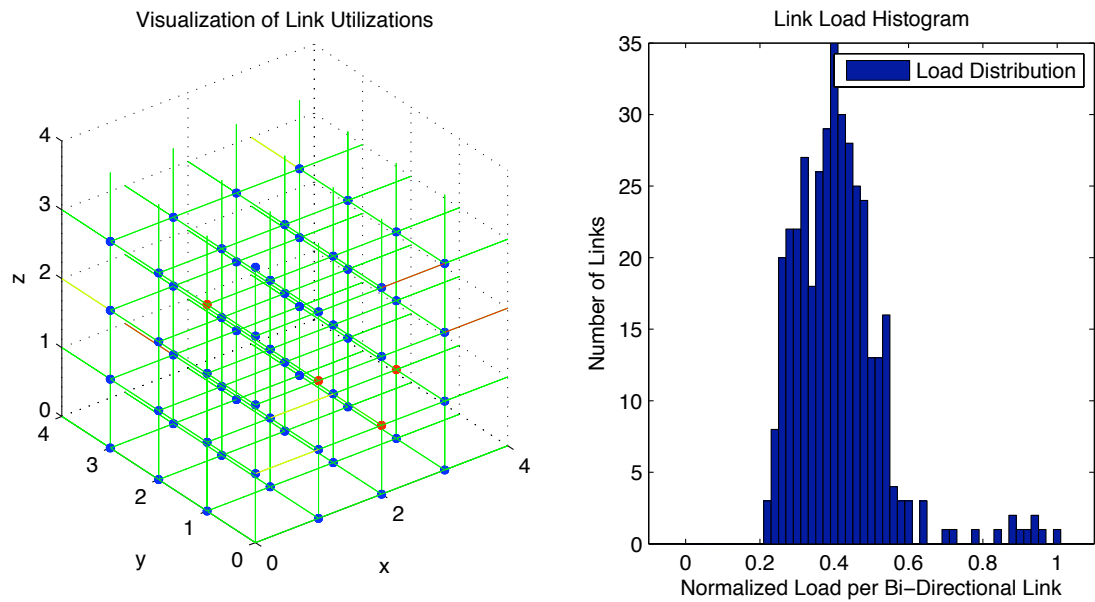
**Figure E.75:** Results of Enhanced CQR Routing using the transpose traffic demand (TP). Two individual traffic demands were assigned for each permutation of all  $x, y, z$  values.



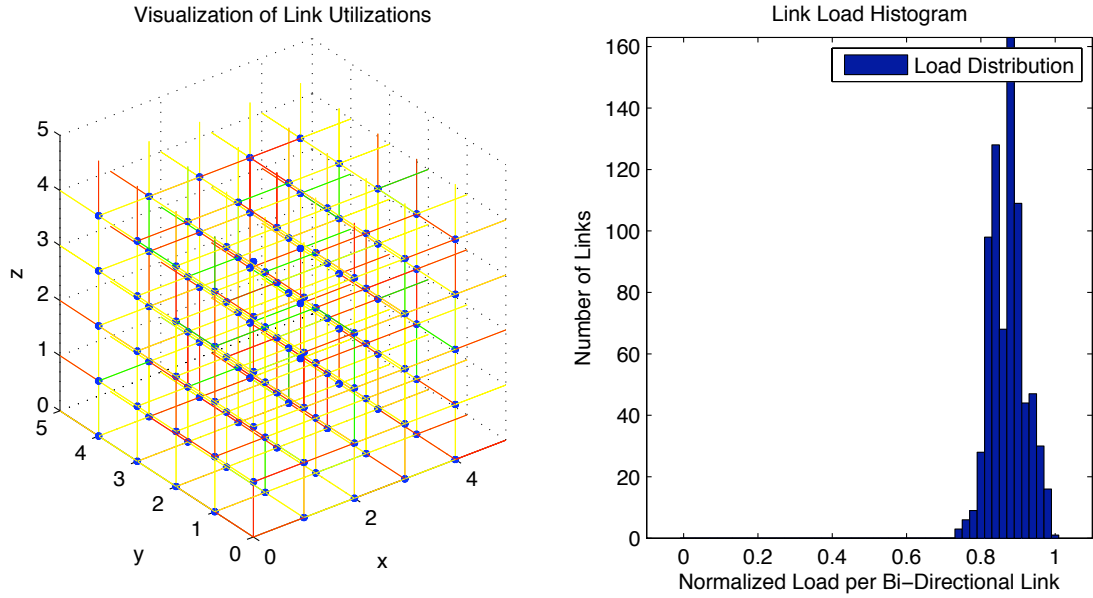
**Figure E.76:** Results of Enhanced CQR Routing using the tornado traffic demand (TOR). Three individual traffic demands were assigned for possible values described in Table 2.1.



**Figure E.77:** Results of Enhanced CQR Routing using the flood traffic demand (FL). Two individual traffic demands were assigned to each node from each node.

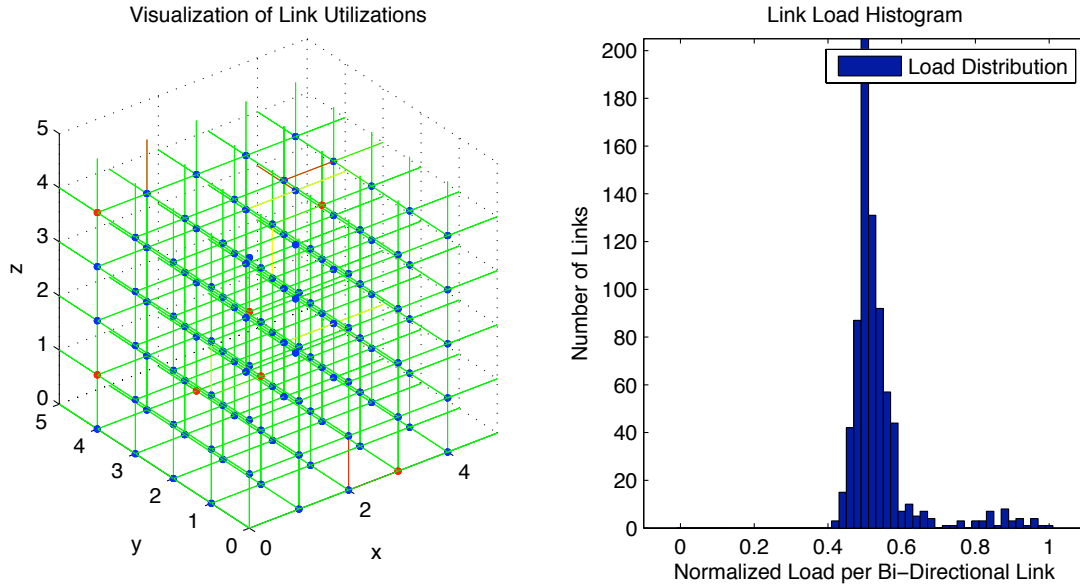


**Figure E.78:** Results of Enhanced CQR Routing using the flood traffic demand (FL) with Hot-spots. Two individual traffic demands were assigned to each node from each node. Five percent of the nodes were made hot-spots, and they are marked as red in the figure.

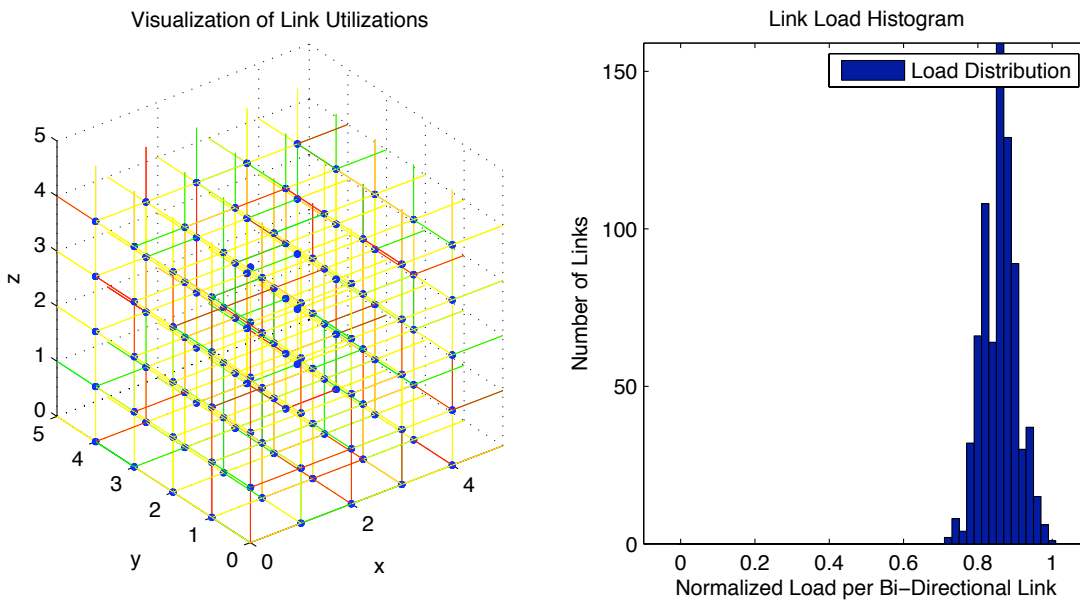


**Figure E.79:** Results of Minimal Adaptive Routing (MA) using the flood traffic demand (FL). Two individual traffic demands were assigned to each node from each node.

### E.3 Results from the 5-ary 3-cube Topology

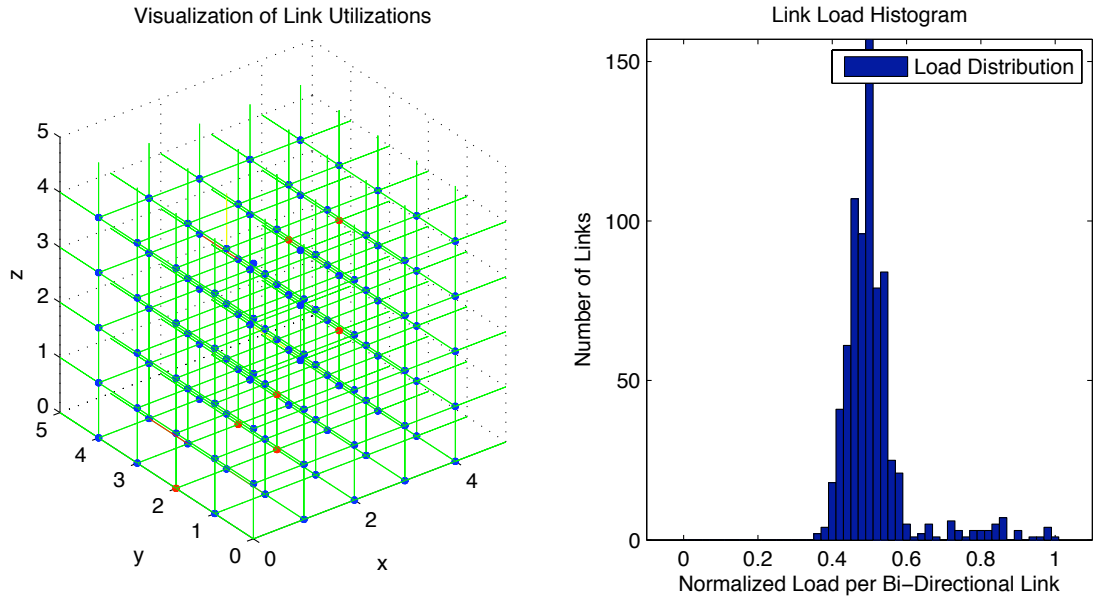


**Figure E.80:** Results of Minimal Adaptive Routing (MA) using the flood traffic demand (FL) with Hot-spots. Two individual traffic demands were assigned to each node from each node. Five percent of the nodes were made hot-spots, and they are marked as red in the figure.

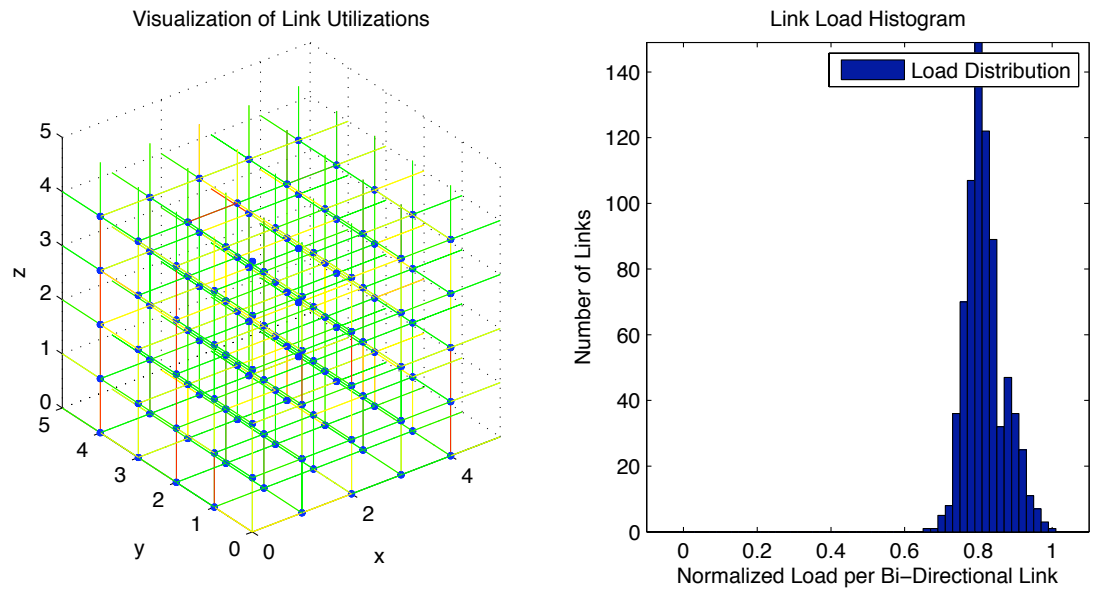


**Figure E.81:** Results of CQR Routing using the flood traffic demand (FL). Two individual traffic demands were assigned to each node from each node.

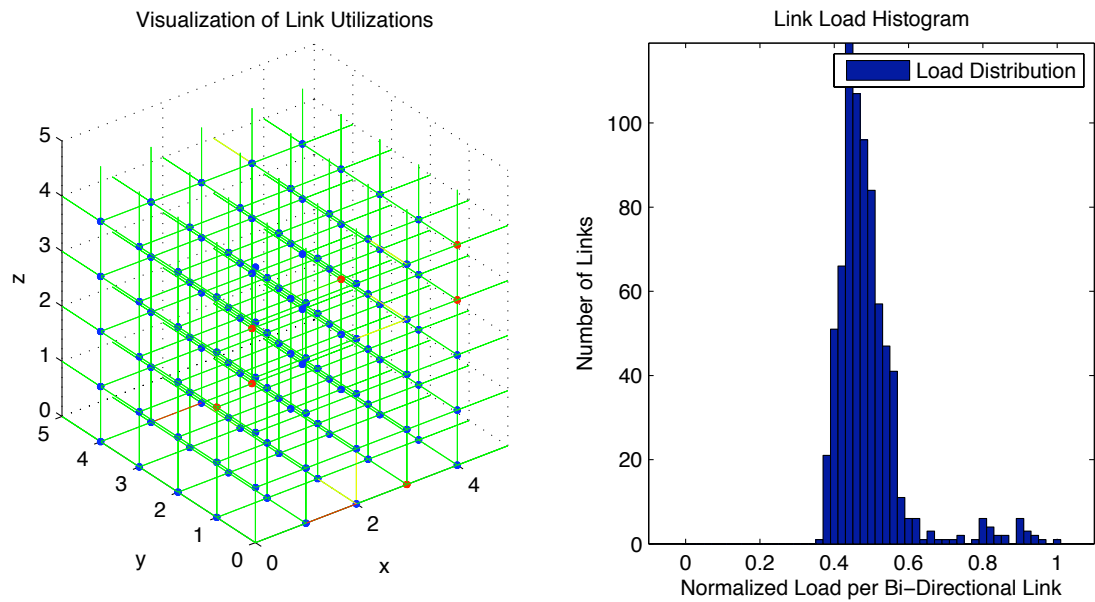




**Figure E.82:** Results of CQR Routing using the flood traffic demand (FL) with Hot-spots. Two individual traffic demands were assigned to each node from each node. Five percent of the nodes were made hot-spots, and they are marked as red in the figure.



**Figure E.83:** Results of Enhanced CQR Routing using the flood traffic demand (FL). Two individual traffic demands were assigned to each node from each node.



**Figure E.84:** Results of Enhanced CQR Routing using the flood traffic demand (FL) with Hot-spots. Two individual traffic demands were assigned to each node from each node. Five percent of the nodes were made hot-spots, and they are marked as red in the figure.