

ARCHITECTURAL PRINCIPLES AND EXPERIMENTATION OF  
DISTRIBUTED HIGH PERFORMANCE VIRTUAL CLUSTERS

by

Andrew J. Younge

Submitted to the faculty of

Indiana University

in partial fulfillment of the requirements for the degree of

Doctor of Philosophy

Department of Computer Science

Indiana University

September 2016

Copyright © 2016 Andrew J. Younge

All Rights Reserved

INDIANA UNIVERSITY

GRADUATE COMMITTEE APPROVAL

of a dissertation submitted by

Andrew J. Younge

This dissertation has been read by each member of the following graduate committee and by majority vote has been found to be satisfactory.

---

Date

Geoffrey C. Fox, Ph.D, Chair

---

Date

Judy Qiu, Ph.D

---

Date

Thomas Sterling, Ph.D

---

Date

D. Martin Swany, Ph.D

INDIANA UNIVERSITY

As chair of the candidate's graduate committee, I have read the dissertation of Andrew J. Younge in its final form and have found that (1) its format, citations, and bibliographical style are consistent and acceptable and fulfill university and department style requirements; (2) its illustrative materials including figures, tables, and charts are in place; and (3) the final manuscript is satisfactory to the graduate committee and is ready for submission to the university library.

---

Date

Geoffrey C. Fox, Ph.D  
Chair, Graduate Committee

Accepted for the Department

---

Department Chair Name, Chair  
Computer Science Program

Accepted for the College

---

Dean Name, Associate Dean  
School of Informatics and Computing

## ABSTRACT

# ARCHITECTURAL PRINCIPLES AND EXPERIMENTATION OF DISTRIBUTED HIGH PERFORMANCE VIRTUAL CLUSTERS

Andrew J. Younge

Department of Computer Science

Doctor of Philosophy

With the advent of virtualization and Infrastructure-as-a-Service (IaaS), the broader scientific computing community is considering the use of clouds for their scientific computing needs. This is due to the relative scalability, ease of use, advanced user environment customization abilities, and the many novel computing paradigms available for data-intensive applications. However, a notable performance gap exists between IaaS and typical high performance computing (HPC) resources. This has limited the applicability of IaaS for many potential users, not only for those who look to leverage the benefits of virtualization with traditional scientific computing applications, but also for the growing number of big data scientists whose platforms are unable to build on HPC's advanced hardware resources.

Concurrently, we are at the forefront of a convergence in infrastructure between Big Data and HPC, in which a unified distributed computing archi-

tecture could provide computing and storage capabilities for both differing distributed systems use cases. There is the potential to leverage performance and advanced hardware from the HPC community, and provide it in a virtualized infrastructure using High Performance Virtual Clusters. This will not only enable a more diverse user environment within supercomputing applications, but also bring increased performance and capabilities to big data platform services.

This work proposes to bridge the gap between HPC and cloud infrastructure and enable infrastructure convergence through a framework for virtual clusters. It begins with an evaluation of current hypervisors and their viability to run HPC workloads within current infrastructure, which helps define existing performance gaps. Next, we uncover mechanisms to enable the use of specialized hardware available in many HPC resources, such as advanced accelerators like the Nvidia GPUs and high-speed, low-latency InfiniBand interconnects. The virtualized infrastructure that developed, which leverages such specialized HPC hardware and utilizes best-practices virtualization using KVM, supports advanced Molecular Dynamics simulations at near-native performance. These advances are incorporated into a framework for constructing distributed virtual clusters using the OpenStack cloud infrastructure. With high performance virtual clusters, we look to support a broad range of scientific computing challenges, from HPC simulations to big data analytics with a single, unified infrastructure.

## ACKNOWLEDGMENTS

Thanks, Mom! (Acknowledgements TBD)

# Contents

<b>Table of Contents</b>	viii
<b>List of Figures</b>	xi
<b>1 Introduction</b>	1
1.1 Overview . . . . .	1
1.2 Research Statement . . . . .	6
1.3 Research Challenges . . . . .	11
1.4 Outline . . . . .	15
<b>2 Related Research</b>	18
2.1 Virtualization . . . . .	18
2.1.1 Hypervisors and Containers . . . . .	21
2.2 Cloud Computing . . . . .	23
2.2.1 Infrastructure-as-a-Service . . . . .	28
2.2.2 Virtual Clusters . . . . .	38
2.2.3 The FutureGrid Project . . . . .	42
2.3 High Performance Computing . . . . .	45
2.3.1 Brief History of Supercomputing . . . . .	45
2.3.2 Distributed Memory Computation . . . . .	47
2.3.3 Exascale . . . . .	49
<b>3 Analysis of Virtualization Technologies for High Performance Computing Environments</b>	52
3.1 Abstract . . . . .	52
3.2 Introduction . . . . .	53
3.3 Related Research . . . . .	55
3.4 Feature Comparison . . . . .	57
3.4.1 Usability . . . . .	59
3.5 Experimental Design . . . . .	61
3.5.1 The FutureGrid Project . . . . .	61
3.5.2 Experimental Environment . . . . .	63
3.5.3 Benchmarking Setup . . . . .	63

3.6	Performance Comparison . . . . .	66
3.7	Discussion . . . . .	71
<b>4</b>	<b>Evaluating GPU Passthrough in Xen for High Performance Cloud Computing</b>	<b>75</b>
4.1	Abstract . . . . .	75
4.2	Introduction . . . . .	76
4.3	Virtual GPU Directions . . . . .	78
4.3.1	Front-end Remote API invocation . . . . .	79
4.3.2	Back-end PCI passthrough . . . . .	81
4.4	Implementation . . . . .	83
4.4.1	Feature Comparison . . . . .	84
4.5	Experimental Setup . . . . .	85
4.6	Results . . . . .	86
4.6.1	Floating Point Performance . . . . .	87
4.6.2	Device Speed . . . . .	88
4.6.3	PCI Express Bus . . . . .	90
4.7	Discussion . . . . .	92
4.8	Chapter Summary and Future Work . . . . .	95
<b>5</b>	<b>GPU-Passthrough Performance: A Comparison of KVM, Xen, VMWare ESXi, and LXC for CUDA and OpenCL Applications</b>	<b>96</b>
5.1	Introduction . . . . .	96
5.2	Related Work & Background . . . . .	98
5.2.1	GPU API Remoting . . . . .	98
5.2.2	PCI Passthrough . . . . .	98
5.2.3	GPU Passthrough, a Special Case of PCI Passthrough . . . . .	99
5.3	Experimental Methodology . . . . .	100
5.3.1	Host and Hypervisor Configuration . . . . .	100
5.3.2	Guest Configuration . . . . .	102
5.3.3	Microbenchmarks . . . . .	102
5.3.4	Application Benchmarks . . . . .	103
5.4	Performance Results . . . . .	104
5.4.1	SHOC Benchmark Performance . . . . .	106
5.4.2	GPU-LIBSVM Performance . . . . .	111
5.4.3	LAMMPS Performance . . . . .	113
5.4.4	LULESH Performance . . . . .	114
5.5	Lessons Learned . . . . .	116
5.6	Directions for Future Work . . . . .	118
<b>6</b>	<b>Supporting High Performance Molecular Dynamics in Virtualized Clusters using IOMMU, SR-IOV, and GPUDirect</b>	<b>119</b>
6.1	Introduction . . . . .	119

6.2	Background and Related Work . . . . .	122
6.2.1	GPU Passthrough . . . . .	123
6.2.2	SR-IOV and InfiniBand . . . . .	124
6.2.3	GPUDirect . . . . .	126
6.3	A Cloud for High Performance Computing . . . . .	126
6.4	Benchmarks . . . . .	127
6.5	Experimental Setup . . . . .	129
6.5.1	Node configuration . . . . .	129
6.5.2	Cluster Configuration . . . . .	131
6.6	Results . . . . .	131
6.6.1	LAMMPS . . . . .	132
6.6.2	HOOMD . . . . .	133
6.7	Discussion . . . . .	135
6.8	Chapter Summary . . . . .	136
<b>7</b>	<b>Virtualization advancements to support HPC applications</b>	<b>138</b>
7.1	Memory Page Table Optimizations . . . . .	139
7.2	Live Migration Mechanisms . . . . .	142
7.2.1	RDMA-enabled VM Migration . . . . .	144
7.2.2	Moving the Compute to the Data . . . . .	147
7.3	Fast VM Cloning . . . . .	148
7.4	Virtual Cluster Scheduling . . . . .	151
7.5	Chapter Summary . . . . .	154
<b>8</b>	<b>Conclusion</b>	<b>155</b>
8.1	Impact . . . . .	157
	<b>Bibliography</b>	<b>158</b>

# List of Figures

1.1	Data analytics and computing ecosystem compared (from [21]), with virtualization included . . . . .	5
1.2	High Performance Virtual Cluster Framework . . . . .	9
1.3	Architectural diagram for High Performance Virtual Clusters . . . . .	10
2.1	Virtual Machine Abstraction . . . . .	19
2.2	Hypervisors and Containers . . . . .	22
2.3	View of the Layers within a Cloud Infrastructure . . . . .	27
2.4	Eucalyptus Architecture . . . . .	32
2.5	OpenNebula Architecture . . . . .	36
2.6	Virtual Clusters on Cloud Infrastructure . . . . .	39
2.7	FutureGrid Participants, Network, and Resources . . . . .	43
2.8	Top 500 Development over time from 2002 to 2016 [8] . . . . .	48
3.1	A comparison chart between Xen, KVM, VirtualBox, and VMWare ESX	57
3.2	FutureGrid Participants and Resources . . . . .	62
3.3	Linpack performance . . . . .	68
3.4	Fast Fourier Transform performance . . . . .	69
3.5	Ping Pong bandwidth performance . . . . .	70
3.6	Ping Pong latency performance (lower is better) . . . . .	71
3.7	Spec OpenMP performance . . . . .	72
3.8	Benchmark rating summary (lower is better) . . . . .	73
4.1	GPU PCI passthrough within the Xen Hypervisor . . . . .	84
4.2	GPU Floating Point Operations per Second . . . . .	88
4.3	GPU Fast Fourier Transform . . . . .	89
4.4	GPU Matrix Multiplication . . . . .	90
4.5	GPU Stencil and S3D . . . . .	91
4.6	GPU Device Memory Bandwidth . . . . .	92
4.7	GPU Molecular Dynamics and Reduction . . . . .	93
4.8	GPU PCI Express Bus Speed . . . . .	94

5.1	SHOC Levels 0 and 1 relative performance on Bespin system. Results show benchmarks over or under-performing by 0.5% or greater. Higher is better. . . . .	107
5.2	SHOC Levels 1 and 2 relative performance on Bespin system. Results show benchmarks over or under-performing by 0.5% or greater. Higher is better. . . . .	108
5.3	SHOC Levels 0 and 1 relative performance on Delta system. Benchmarks shown are the same as Bespin's. Higher is better. . . . .	109
5.4	SHOC Levels 1 and 2 relative performance. Benchmarks shown are the same as Bespin's. Higher is better. . . . .	110
5.5	GPU-LIBSVM relative performance on Bespin system. Higher is better.	112
5.6	GPU-LIBSVM relative performance on Delta showing improved performance within virtual environments. Higher is better. . . . .	113
5.7	LAMMPS Rhodopsin benchmark relative performance for Bespin system. Higher is better. . . . .	114
5.8	LAMMPS Rhodopsin benchmark relative performance for Delta system. Higher is better. . . . .	115
5.9	LULESH relative performance on Bespin. Higher is better. . . . .	116
6.1	Node PCI Passthrough of GPUs and InfiniBand . . . . .	130
6.2	LAMMPS LJ Performance . . . . .	132
6.3	LAMMPS RHODO Performance . . . . .	134
6.4	HOOMD LJ Performance with 256k Simulation . . . . .	135
7.1	Transparent Huge Pages with KVM . . . . .	141
7.2	Adding extra specs to a VM flavor in OpenStack . . . . .	152

# <sup>1</sup> Chapter 1

## <sup>2</sup> Introduction

### <sup>3</sup> 1.1 Overview

<sup>4</sup> For years visionaries in computer science have predicted the advent of utility-based  
<sup>5</sup> computing. This concept dates back to John McCarthy's vision stated at the MIT  
<sup>6</sup> centennial celebrations in 1961.

<sup>7</sup> “If computers of the kind I have advocated become the computers of the  
<sup>8</sup> future, then computing may someday be organized as a public utility just  
<sup>9</sup> as the telephone system is a public utility... The computer utility could  
<sup>10</sup> become the basis of a new and important industry.“

<sup>11</sup> Only recently has the hardware and software become available to support the concept  
<sup>12</sup> of utility computing on a large scale.

<sup>13</sup> The concepts inspired by the notion of utility computing have combined with  
<sup>14</sup> the requirements and standards of Web 2.0 [1] to create Cloud computing [2–4].  
<sup>15</sup> Cloud computing is defined as, “A large-scale distributed computing paradigm that  
<sup>16</sup> is driven by economies of scale, in which a pool of abstracted, virtualized, dynamically-  
<sup>17</sup> scalable, managed computing power, storage, platforms, and services are delivered on

18 demand to external customers over the Internet.” This concept of Cloud computing  
19 is important to Distributed Systems because it represents a true paradigm shift [5]  
20 within the entire IT infrastructure. Instead of adopting the in-house services, client-  
21 server model, and mainframes, Clouds push resources out into abstracted services  
22 hosted *en masse* by larger organizations. This concept of distributing resources is  
23 similar to many of the visions of the Internet itself, which is where the “Clouds”  
24 nomenclature originated, as many people depicted the internet as a big fluffy cloud  
25 one connects to.

26 As the core of most cloud infrastructure lies virtualization, a computer architecture  
27 technology by which 1 or more Virtual Machines (VMs) are run on the same physical  
28 host. In doing this, a layer of abstraction is inserted between and around the hardware  
29 and Operating System (OS). Specifically, hardware resources such as CPUs, memory,  
30 and I/O devices, and software resources such as OS functionality and low level libraries  
31 are abstracted and provided to VMs directly. While virtualization has existed for  
32 many years, its availability with Intel x86 commodity hardware in conjunction with  
33 the rise of clouds has brought virtualization to the forefront of distributed systems.

34 While Cloud computing is changing IT infrastructure, it also has had a drastic  
35 impact on distributed systems itself, which has a different evolution. Gone are the  
36 IBM Mainframes of the 1980’s which dominated the enterprise landscape. While  
37 some mainframes still exist, they are used only for batch related processing tasks  
38 and are relatively unused for scientific applications as they are inefficient at Floating  
39 Point Operations. As such, they were replaced with Beowulf Clusters [6], Massively  
40 Parallel Processors (MPPs) and Supercomputers of the 90’s and 00’s. A novelty of  
41 these distributed memory systems is that instead of just one large machine, many  
42 machines are connected together and used to achieve a common goal, thereby maxi-  
43 mizing the overall speed of computation. Clusters represent a more commodity-based

<sup>44</sup> supercomputer, where off the shelf CPUs are used instead of the highly customized  
<sup>45</sup> and expensive processors and interconnects found in Supercomputers.

<sup>46</sup> Supercomputers and Clusters are best suited for large scale applications. These  
<sup>47</sup> HPC applications can even include “Grand Challenge” applications [7] and can repre-  
<sup>48</sup> sent a sizable amount of the scientific calculations done on large-scale Supercomputing  
<sup>49</sup> resources today. However, there exists a gap of many orders of magnitude between  
<sup>50</sup> leading-class high performance computing and what is available on the common labo-  
<sup>51</sup> ratory workshop. This gap, described here as mid-tier scientific computation is a fast  
<sup>52</sup> growing field that struggles to efficiently harness distributed systems while hoping to  
<sup>53</sup> minimize extensive development efforts. These mid-tier scientific endeavours need to  
<sup>54</sup> leverage distributed systems to effectively complete the calculations at hand, however  
<sup>55</sup> may not require the extreme scale provided by the latest machines at the peak of  
<sup>56</sup> the Top500 list [8]. Furthermore, it may not be feasible for small research teams to  
<sup>57</sup> effectively handle the development complexity of extreme-scale resources. This can  
<sup>58</sup> include scientific disciplines such as high energy physics [9], materials science [10],  
<sup>59</sup> bioinformatics [11], and climate research [12], to name a few.

<sup>60</sup> As more domain science turns to the aid of computational resources for con-  
<sup>61</sup> ducting novel scientific endeavours, there is a continuing and growing need for na-  
<sup>62</sup> tional cyberinfrastructure initiatives to support an increasingly diverse set of scientific  
<sup>63</sup> workloads. Substantial growth can be see in the number of computational resource  
<sup>64</sup> requests [13, 14] from many of the larger computational facilities. Concurrently, there  
<sup>65</sup> has also been an increase in accelerators and hybrid computing models capable of  
<sup>66</sup> quickly providing additional resources [15] beyond commodity clusters.

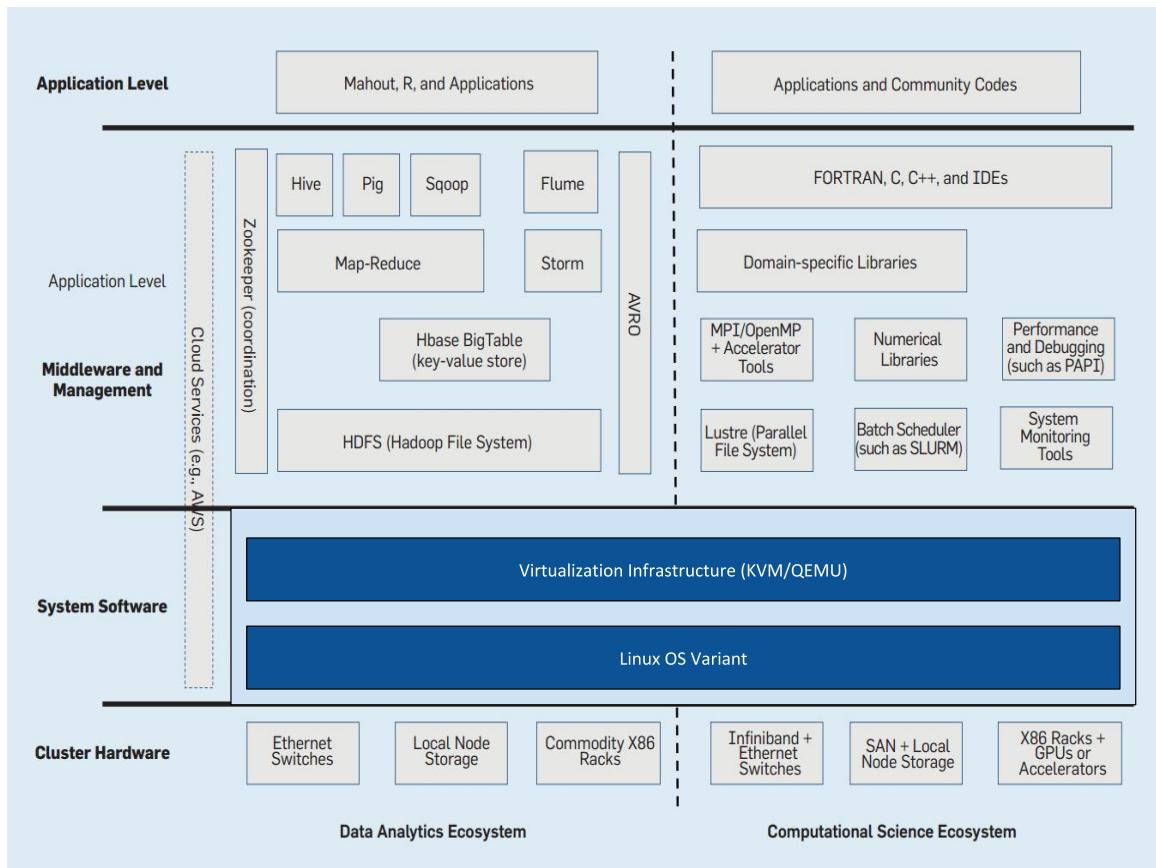
<sup>67</sup> Historically, application diversity was separated into High Performance Comput-  
<sup>68</sup> ing (HPC) and High Throughput Computing (HTC). With HTC, computational  
<sup>69</sup> problems can be separated into independent tasks that can execute in a pleasingly

70 parallel fashion, happily gaining any available resources and rarely require significant  
71 communication or synchronization between tasks. HPC often represents computa-  
72 tional problems that require significant communication and coordination to effectively  
73 produce results, usually with the use of a communication protocol such as MPI [16].  
74 Recently however, many big-data paradigms [17] have been introduced in distributed  
75 systems that represent new computational models for distributed computing, such as  
76 MapReduce [18] and the corresponding Apache Big Data stack [19, 20]. Supporting  
77 these different distributed computational paradigms requires a flexible infrastructure  
78 capable of providing computational resources for all possible models in a fast and  
79 efficient manner.

80 Currently we are at the forefront of a convergence within scientific computing be-  
81 tween HPC and big data computation [21]. This amalgamation of historically differing  
82 viewpoints of Distributed Systems looks to combine the performance characteristics of  
83 HPC and the pursuit towards Exascale with the data-centric and concurrency models  
84 found in Big Data and cloud services.

85 Much of the convergence effort has been focused on applications and platform  
86 services. Specifically, significant work towards convergent applications has been out-  
87 lined with the Big Data Ogres [22] and the corresponding HPC-ABDS model [23].  
88 This convergence can also be seen with efforts in bringing interconnect advances to  
89 classically big data platforms such as with InfiniBand and MapReduce [24]. However,  
90 the underlying hardware and OS environments are still something to be reconciled,  
91 and represents something that virtualization provides. It is expected that new big  
92 data efforts will continue to move in this direction [25], especially if virtualization can  
93 make traditionally prohibitive HPC hardware such as accelerators and high-speed in-  
94 terconnects readily available to cloud and big data platforms. As the deployment of  
95 such applications and platform services on virtualized infrastructure is well defined

and regarded within industry [26], this dissertation has focused the difficulty of running HPC applications on similar yet performance-enabled virtualized infrastructure. However, it is possible and hopeful that research regarding virtualization can also play a part in bringing advanced hardware and performance-focused considerations to Big Data applications, effectively cross-cutting the convergence with HPC. In Summary, success of the research could be defined by the ability to support the convergence between HPC and Big Data.



**Figure 1.1** Data analytics and computing ecosystem compared (from [21]), with virtualization included

To further illustrate where virtualization can play a part in HPC and Big Data convergence, we look at Figure 1.1 from Reed & Dongarra [21]. While the two ecosystems depicted are only representative and in no way exhaustive, they do show how

106 drastically different user environments exist and are reliant on differing hardware. If  
107 we insert a performance-oriented virtualization mechanism within the system soft-  
108 ware capable of handling the advanced cluster hardware, it could provide a single,  
109 comprehensive *convergent ecosystem* that supports both HPC and Big Data efforts  
110 at a critical level.

111 This work proposes the use of virtual clusters [27] to provide distinct, separated  
112 environments on similar resources using virtual machines. These virtual clusters,  
113 similar in design to the Beowulf clusters and commodity HPC systems, provide a dis-  
114 tributed memory environment, usually with a local resource management system [28],  
115 however they do so in virtual machines. However, past concerns with virtualization  
116 have limited the adoption of virtual clusters in many large scale cyberinfrastructure  
117 deployments. This has largely been due to the overhead of virtualization, whereby  
118 many scientific computations have experienced a significant and notable degradation  
119 in performance. In an ecosystem familiarized with HPC systems where performance  
120 is paramount, this has been an obstructive hurdle for many even mid-level scientific  
121 endeavors.

## 122 1.2 Research Statement

123 With the rise of cloud computing within the greater realm of distributed systems,  
124 there have been a number of scientific computing endeavors map well to cloud in-  
125 frastructure. This first includes the simple and most common practice of on demand  
126 computing where by users can rent-a-workstation [29]. Perhaps these resources are  
127 more powerful than a given researcher's laptop and used to run their scientific appli-  
128 cations or support greater laboratory collaborative efforts, such as a shared database  
129 or Web services. Second, we've seen virtualized cloud infrastructure support high

130 throughput computing very well. Often times pleasingly parallel applications, be it  
131 from high energy physics such as the LHC effort [9,30] or bioinformatics with BLAST  
132 alignment jobs [11], have proven to run with high efficiency in public and private cloud  
133 environments. The rise of public cloud infrastructure has also coincided with increase  
134 in big data computation and analytics. Many of these big data platform services  
135 have evolved complimentary to cloud infrastructure, and as such have a symbiotic  
136 relationship with virtualization technologies [31].

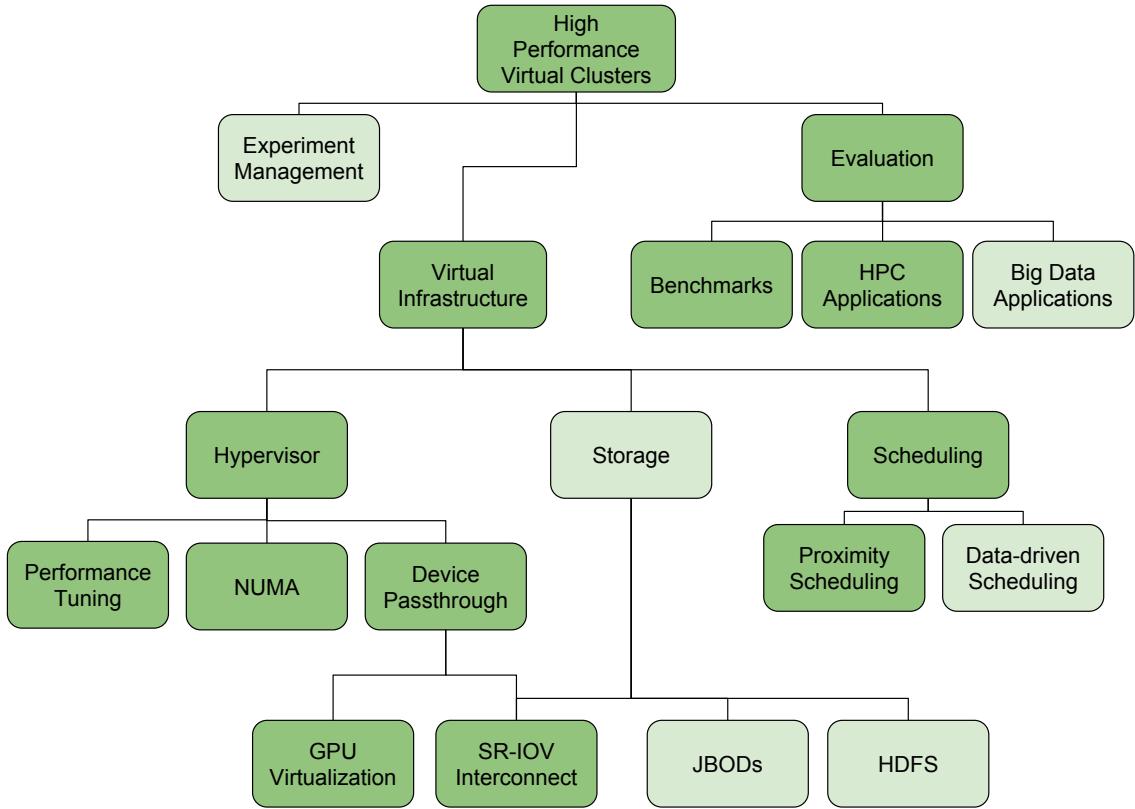
137 However, with tightly coupled, high performance distributed memory applications,  
138 the same endeavors that support leading class scientific efforts, run very poorly on  
139 virtualized cloud infrastructure [32]. This is due to a myriad of addressable reasons  
140 ranging from scheduling, abstraction overhead, or a lack of advanced hardware sup-  
141 port necessary for tightly coupled communication. This postulates the question on  
142 whether virtualization can in fact support such tightly coupled large scale applica-  
143 tions without an imposed significant performance penalty. Simply put, *the goal of this*  
144 *dissertation is to investigate the viability of mid-tier scientific applications supported*  
145 *in a virtualized infrastructure.*

146 Historically, mid-tier scientific applications are distributed memory HPC applica-  
147 tions that require more complex process communication mechanisms. These systems  
148 need far more performance than a single compute resources (such as a workstation)  
149 can provide. This could include hundreds or thousands of processes calculating and  
150 communicating concurrently on a cluster, perhaps using a messaging interface such  
151 as MPI. These applications are likely distinct either in application composition or  
152 operating parameters, from extreme-scale HPC applications that run at the highest  
153 end of the supercomputing resources today, which operate on petascale machines and  
154 beyond.

155 Given the current outlook on virtualization for supporting HPC applications, this

156 dissertation proposes a framework for High Performance Virtual Clusters that enable  
157 advanced computational workloads, including tightly coupled distributed memory ap-  
158 plications, to run with a high degree of efficiency in virtualized environments. This  
159 framework, outlined in Figure 1.2, illustrates the topics to be addressed to provide  
160 a supportive virtual cluster environment for high performance mid-tier scientific ap-  
161 plications. Areas marked in darker green indicate topics this dissertation may touch  
162 upon, whereas light green areas in Figure 1.2 identify outstanding considerations.  
163 We specifically identify mid-tier distributed memory parallel computations as a focal  
164 point for the computational challenges at hand as a way to separate from some of the  
165 latest efforts in towards Exascale [33–35] computing. While virtualization may in fact  
166 be able to play a role towards usable Exascale computing, such efforts fall outside the  
167 immediate scope of this dissertation.

168 In order to provide comprehensive high performance virtual clusters, we need to  
169 first look at a key area itself, the virtualized infrastructure itself. At the core, we  
170 have to consider the hypervisor or virtual machine monitor and the overhead and  
171 performance characteristics associated with it. This includes performance tuning  
172 considerations, NUMA effects, and advanced hardware device passthrough. Specifi-  
173 cally, device passthrough in the context of this manuscript refers to two major device  
174 types; GPUs and InfiniBand interconnects (the later using SR-IOV). The virtual in-  
175 frastructure also must consider scheduling as a major factor in performing efficient  
176 placement of workloads on the underlying host infrastructure, and in particular a  
177 Proximity scheduler is of interest [36]. Storage solutions in a virtualized environ-  
178 ment is an increasingly important aspect of this framework, as both HPC and big  
179 data solutions are continuing to prioritize I/O performance compared to computa-  
180 tion. Storage is also likely to be heavily dependent on interconnect considerations as  
181 well, as potentially provided by device passthrough, however such I/O considerations

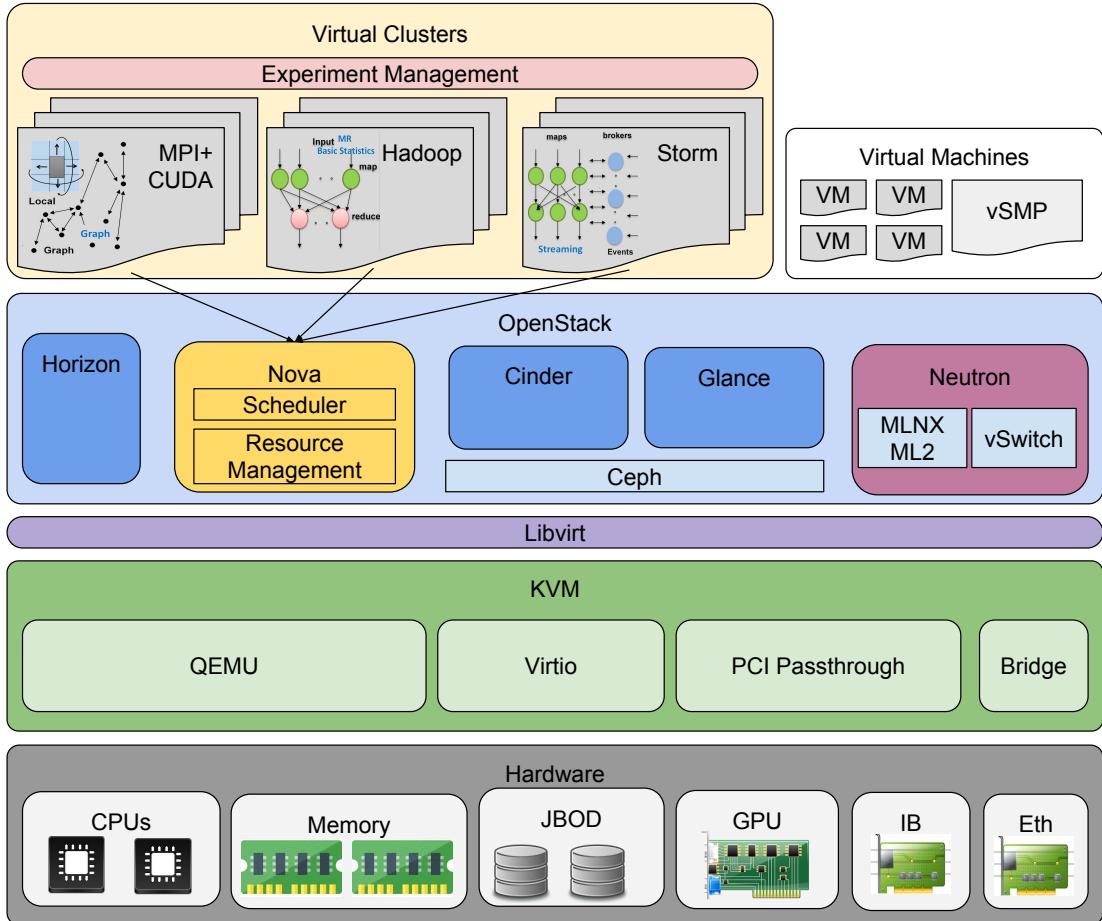


**Figure 1.2** High Performance Virtual Cluster Framework

<sup>182</sup> lie beyond this dissertation's immediate scope.

<sup>183</sup> However, simply providing an enhanced virtualized infrastructure may not guarantee that all implementations of high performance virtual clusters are performant. <sup>184</sup> Specifically, proposed infrastructures need to be properly evaluated in a systematic way through the use of a wide array of benchmarks, mini-applications, and full-scale scientific applications. This effort can further be broken down into three major problem sets; base level benchmarking tools, HPC applications, and big data applications. <sup>185</sup> Evaluating the stringent performance requirements of all three sets, when compared with bare metal (no virtualization) common solutions, will illuminate not only successful designs but also the focus areas that require more attention. As such, <sup>186</sup> we look to continually use these benchmarks and applications as a tool to measure

193 the viability of virtualization in this context.



**Figure 1.3** Architectural diagram for High Performance Virtual Clusters

194 Building from the historical virtual clusters in Grid computing, we see a new  
195 architectural model for high performance virtual clusters illustrated in 1.3 that is  
196 implied by this dissertation. Here, we leverage commodity hardware, as well as some  
197 advanced HPC hardware. While this notion could incorporate a wide array of dif-  
198 fering technologies, we focus here on GPU-based accelerators and InfiniBand inter-  
199 connects in conjunction with x86 CPUs. Atop this, we leverage KVM and QEMU to  
200 provide an advanced hypervisor for creating and hosting VMs with direct hardware  
201 involvement from the lower level. Moving up, the Libvirt API is leveraged due to

202 its hypervisor interoperability and popularity. Atop Libvirt is the OpenStack private  
203 cloud infrastructure. In Figure 1.3 we illustrate some (but not all) of OpenStack’s  
204 services including the Horizon UI, Cinder and Glance storage mechanisms, and the  
205 Neutron (previously Quantum) components. The Nova component of OpenStack is  
206 the point of focus for providing comprehensive VM management.<sup>1</sup> Atop OpenStack,  
207 we can create a wide array of virtual clusters and machines to support the wide rang-  
208 ing scientific computing ecosystems necessary. This includes application models such  
209 as tightly coupled MPI+CUDA HPC applications, to emerging big data analytics  
210 toolkits such as Apache Storm [37].

211 One of the higher-level aspects of providing high performance virtual clusters is  
212 the high level orchestration of the virtual clusters themselves. While this largely  
213 remains tangential to this immediate research, it is none the less a key aspect for  
214 a successful solution. Some effort has been put forth for virtual cluster experiment  
215 management [38], and many ongoing open sources solutions also offer compelling  
216 options, such as OpenStack Heat [39]. An example of a project delivering advanced  
217 orchestration mechanisms and a toolkit to aid in configurable virtual clusters on  
218 heterogeneous IaaS deployments is the Cloudmesh effort [40].

### 219 **1.3 Research Challenges**

220 The framework, architecture, and efforts described in this dissertation represent a  
221 movement forward in providing virtualized infrastructure to support a wide arrange  
222 of scientific applications. However, there still exist some challenges that will need

---

<sup>1</sup>While many of the features for nova’s additions in GPU Passthrough and SR-IOV InfiniBand support have been put together at USC/ISI as an OpenStack Nova fork (<https://libraries.io/github/usc-isi/nova>), much of the features have since been modified and matured by the OpenStack community in later releases and made available in upstream Nova.

223 to be addressed. This includes a sigma of virtualization being inherently slow and  
224 unable to support tightly coupled computations, limitations with advancing at scale,  
225 and even that containers may provide a better alternative. While this work hopes to  
226 move beyond these challenges, they none the less must be considered.

227 The notion that virtualization and Cloud infrastructure are not able to support  
228 parallel distributed memory applications has been characterized many times. One  
229 of the most prominent examples of this is the Department of Energy’s Magellan  
230 Project [41], where by the Magellan Final Report [42] states the following finding as  
231 a Key Finding:

232 **“Finding 2. Scientific applications with minimal communication  
233 and I/O are best suited for clouds.”**

234 We have used a range of application benchmarks and micro-benchmarks  
235 to understand the performance of scientific applications. Performance of  
236 tightly coupled applications running on virtualized clouds using commod-  
237 ity networks can be significantly lower than on clusters optimized for these  
238 workloads. This can be true at even mid-range computing scales. For ex-  
239 ample, we observed slowdowns of about 50x for PARATEC on the Amazon  
240 EC2 instances compared to Magellan bare-metal (non-virtualized) at 64  
241 cores and about 7x slower at 1024 cores on Amazon Cluster Compute in-  
242 stances, the specialized high performance computing offering. As a result,  
243 current cloud systems are best suited for high-throughput, loosely coupled  
244 applications with modest data requirements.“

245 These findings underscore how classical usage of virtualization in cloud infrastruc-  
246 ture has serious performance issues when running tightly coupled distributed memory  
247 applications. Many of these performance concerns are sound, given the limitation of a

248 number of virtualization overheads commonplace at the time, including shadow page  
249 tables, emulated Ethernet drivers, experimental hypervisors, and a complete lack of  
250 sophisticated hardware concurrently was commonplace in supercomputers and clus-  
251 ters. As a result, the advantages of virtualization, including on-demand resource  
252 allocation, live migration and advanced hybrid migration, and user-defined environ-  
253 ments, have not been able to effectively show their value in the context of the HPC  
254 community.

255 Other and related efforts within the scientific community too found limitations  
256 with HPC applications in public cloud environments. This includes the study by  
257 Jackson et. al [43] which illustrates how the Amazon Elastic Compute Cloud (EC2)  
258 creates a 6x performance impact compared to a local cluster, due to a large part  
259 on the limiting Gigabit Ethernet that benchmarks relied heavily on within the EC2  
260 system. Other studies also found similar results such as Ostermann [32], concluding  
261 that Amazon EC2 “is insufficient for scientific computing at large, though it still  
262 appeals to the scientists that need resources immediately and temporarily.” However,  
263 these studies are now outdated and do not take into account the advancements in  
264 virtualization and hardware availability detailed in this dissertation. Specifically, it is  
265 estimated that with the KVM hypervisor in a performance-tuned environment, using  
266 accelerators and most certainly a high-speed, low latency interconnect as detailed in  
267 Chapter 6, the results would be drastically different.

268 One limitation in research related to high performance virtual clusters is the degree  
269 to which applications can scale remains relatively unknown. While initial results  
270 with SR-IOV InfiniBand are promising, scaling is naturally hard to predict. While  
271 unfounded, it would be hypothetically possible that as the number of VM’s increases,  
272 tail-latency could also increase, causing notable slowdowns during distributed memory  
273 synchronization barriers. It is only when infrastructure comes available to support

274 high performance virtual clusters will scaling beyond to thousands of cores and beyond  
275 be feasible.

276 Another potential challenge, and perhaps also a strength, is the rising use of  
277 containers within industry, such as with efforts like Docker [44]. Recently, we've seen  
278 efforts at NERSC/LBNL adapt a container solution called Shifter with the SLURM  
279 resource manager on CRAY XC based systems [45]. Shifter's goal is to provide  
280 user defined images for NERSC's bioinformatics users, and it adapts remarkably well  
281 to the HPC environment. While further efforts are needed by Cray and NERSC  
282 to fully provide a container-bases solution on a large scale Supercomputer for all  
283 applications, its efforts are in many ways parallel to using virtualization. In Chapter 5,  
284 we specifically compare virtualization efforts with LXC [46], a popular Linux container  
285 solution and find performance to be comparable and largely near-native.

286 While the major concern with virtualization in the HPC community is perfor-  
287 mance issues, virtualization itself may not be fundamentally limited by the overhead  
288 that causes issues in running high performance computing applications. Recent im-  
289 provements in performance, along with increased usability of accelerators and high  
290 speed, low latency interconnects in virtualized environments as demonstrated in this  
291 dissertation, have made virtual clusters a more viable solution for mid-tier sci-  
292 entific applications. Furthermore, it is possible for virtualization technologies to bring  
293 enhanced usability and enable specialized runtime components to future HPC re-  
294 sources, adding significant value over today's supercomputing resources. This could  
295 potentially include infrastructure advances for higher level cloud platform services for  
296 supporting big data applications [23].

**297 1.4 Outline**

298 The rest of this dissertation is organized into chapters, each signifying the steps to  
299 move forward the notion of a high performance virtual cluster solution.

300 Chapter 2 investigates the related research surrounding both cloud computing  
301 and high performance computing. Within cloud computing, an introduction to cloud  
302 infrastructure, virtualization, and containers will all be discussed. This also includes  
303 details regarding virtual clusters as well as an overview of some national scale cloud  
304 infrastructure efforts that exist. Furthermore, we investigate the state of high per-  
305 formance computing and supercomputing, as well touch upon some of the current  
306 Exascale efforts.

307 Chapter 3 takes a look at the potential for virtualization, in its current state,  
308 for high performance computing. This includes a feature comparison for hardware  
309 availability of a few common hypervisors, specifically Xen, KVM, VirtualBox, and  
310 VMWare. Then, a few common HPC benchmarks are evaluated to determine what  
311 overhead exists and where in a single node configuration. This identifies how in some  
312 scenarios, virtualization adds only a minor overhead, whereas with other scenarios,  
313 overheads can be up to 50% compared to native configurations.

314 Chapter 4 starts to overcome one of the main limitations of virtualization for use  
315 in advanced scientific computing, specifically the lack of hardware availability. In this  
316 chapter, The Xen hypervisor is used to demonstrate the affect of GPU Passthrough,  
317 allowing for GPUs to be used in a guest VM. The efficiency of this method is briefly  
318 evaluated using two different hardware setups, and finds hardware can play a notable  
319 role in single node performance.

320 Chapter 5 continues where chapter 4 leaves off, by demonstrating that GPU  
321 passthrough is possible on many other hypervisors, specifically also KVM and VMWare,

322 and compares with one of the main containerization solutions, LXC. Here, the GPUs  
323 are evaluated using not only the SHOC GPU benchmark suite developed at Oak  
324 Ridge National Laboratory, but also a diverse mix of real-world applications to ex-  
325 amine how and where overhead exists with GPUs in VMs for each virtualization setup.  
326 Specifically, we find that with properly tuned hardware and NUMA-balanced configu-  
327 rations, that the KVM solution can perform at roughly near-native performance, with  
328 on average 1.5% overhead compared to no virtualization. This illustrates that with  
329 the combination correct hypervisor selection, careful tuning, and advanced hardware,  
330 scientific computations can be supported using virtualized hardware.

331 Chapter 6 takes the findings from the previous chapter to the next level. Specifi-  
332 cally, the lessons learned from successful KVM virtualization with GPUs is expanded  
333 and combined with a missing key component of supporting advanced parallel com-  
334 putations: a high speed, low latency interconnect, specifically InfiniBand. Using  
335 SR-IOV and PCI passthrough of QDR InfiniBand interconnect across a small cluster,  
336 it is demonstrated that two Molecular Dynamics simulations, both very commonly  
337 used in the HPC community, can be run at near-native performance in the designed  
338 virtual cluster.

339 Chapter 7 takes a look at the given situation of virtualization, and puts forth an  
340 argument for enhancements forthcoming in high performance virtual cluster solutions.  
341 Specifically, we look at the given state of the art, how virtual clusters can be used  
342 to provide an infrastructure to support the convergence between HPC and big data.  
343 Specifically, this chapter outlines and investigates potential next steps for virtualiza-  
344 tion, including the potential for advanced live migration techniques and VM cloning,  
345 which can be made available with the inclusion of a high-performance RDMA-capable  
346 interconnect.

347 Finally, this work concludes with an overall view of the current state of high

<sup>348</sup> performance virtualization, as well as it's potential to impact and support a wide  
<sup>349</sup> array of disciplines.

# <sup>350</sup> Chapter 2

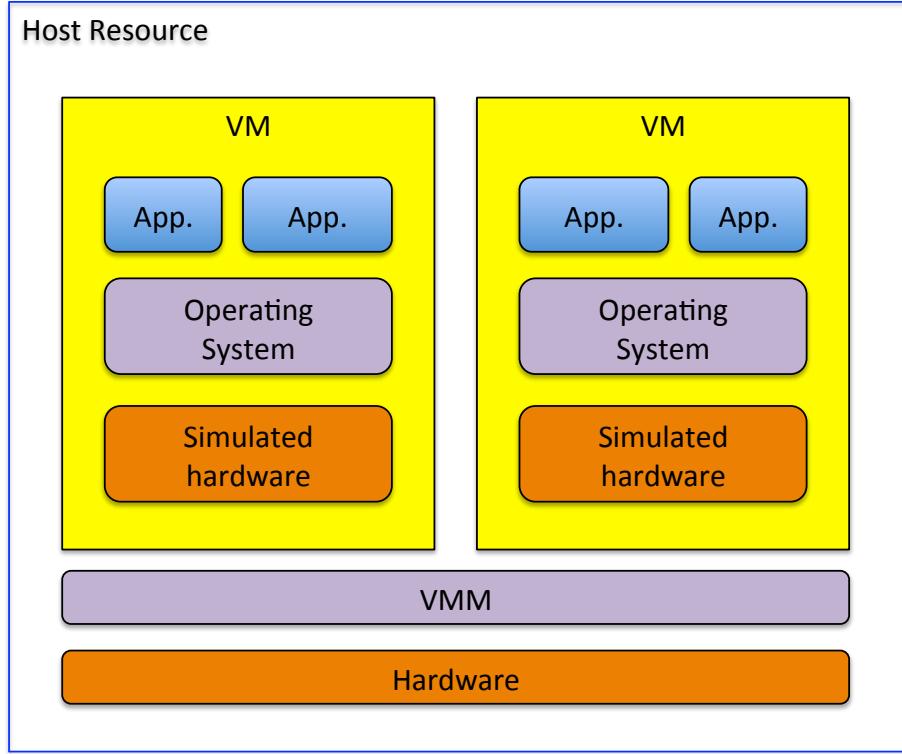
## <sup>351</sup> Related Research

<sup>352</sup> In order to accurately depict the research presented in this article, the topics within  
<sup>353</sup> Virtualization, Cloud computing, and High Performance Computing are reviewed in  
<sup>354</sup> detail.

### <sup>355</sup> 2.1 Virtualization

<sup>356</sup> Virtualization is a way to abstract the hardware and system resources from a operating  
<sup>357</sup> system. This is typically performed within a Cloud environment across a large  
<sup>358</sup> set of servers using a Hypervisor or Virtual Machine Monitor (VMM) which lies in  
<sup>359</sup> between the hardware and the Operating System (OS). From here, one or more virtualized  
<sup>360</sup> OSs can be started concurrently as seen in Figure 2.1, leading to one of the key  
<sup>361</sup> advantages of Cloud computing. This, along with the advent of multi-core processing  
<sup>362</sup> capabilities, allows for a consolidation of resources within any data center. It is the  
<sup>363</sup> Cloud's job to exploit this capability to its maximum potential while still maintaining  
<sup>364</sup> a given QoS. It should be noted that virtualization is not specific to Cloud computing.  
<sup>365</sup> IBM originally pioneered the concept in the 1960's with the M44/44X systems. It

<sup>366</sup> has only recently been reintroduced for general use on x86 platforms.



**Figure 2.1** Virtual Machine Abstraction

<sup>367</sup> However, virtualization is in the most general form, just another form of abstraction.  
<sup>368</sup> As such, there are in fact many levels to virtualization that exist [47].

- <sup>369</sup> • **ISA -** Virtualization can start from the instruction set architecture (ISA) level, where by an entire processor instruction set is emulated and provided. This may be useful for running software or services developed for one instruction set (say MIPS) but is needed to run on Intel x86 hardware. ISA level virtualization is usually emulated through an interpreter which translates source instructions to target instructions, however this can often be extremely inefficient. Dynamic binary translation can help aid in efficiency by translating blocks of source instructions to target instructions, however this still can be limiting.

- 377     ● **Hardware** - One of the most important technologies in cloud computing is  
378         hardware level virtualization [48, 49]. Hardware virtualization, in its most pure  
379         form, refers to the process of creating virtual abstraction to hardware plat-  
380         forms, operating systems, or software resources. This enables the creation of  
381         1 or more virtual machines (VMs) that are run concurrently on the same op-  
382         erating environment, be it hardware or some higher software. Here, a virtual  
383         hardware environment is generating for a VM, including providing virtual pro-  
384         cessors, memory, and I/O devices, allowing for a multiplexing of VMs to exist,  
385         as depicted in Figure 2.1. This layer also manages the physical hardware on  
386         behalf of a host OS as well as for guests. While the most common type of hard-  
387         ware virtualization is with the Xen Virtual Machine Monitor [48], this method  
388         has further separated into type 1 and type 2 hypervisors, as detailed further in  
389         Section 2.1.1.
  
- 390     ● **Operating System** - Moving up the latter of implementation with virtu-  
391         alization, we find OS level virtualization, where multiplexing happens at the  
392         level of the OS, rather than the hardware. Usually, this refers to isolating a  
393         filesystem and associated process and runtime effects in a single "chroot" or  
394         "container" environment at the user level, with kernel level operations being  
395         handled by a single OS kernel. This containerization allows for multiple user  
396         environments to exist concurrently without the complexity of hardware or ISA  
397         level virtualization. Containers are also describe in the next section in more  
398         detail.
  
- 399     ● **Library (API)** - Continuing, we move up to library or API level virtualization.  
400         Here, some API is separated and a provided to a user or programming environ-  
401         ment, and the communication between this API and the back-end library with

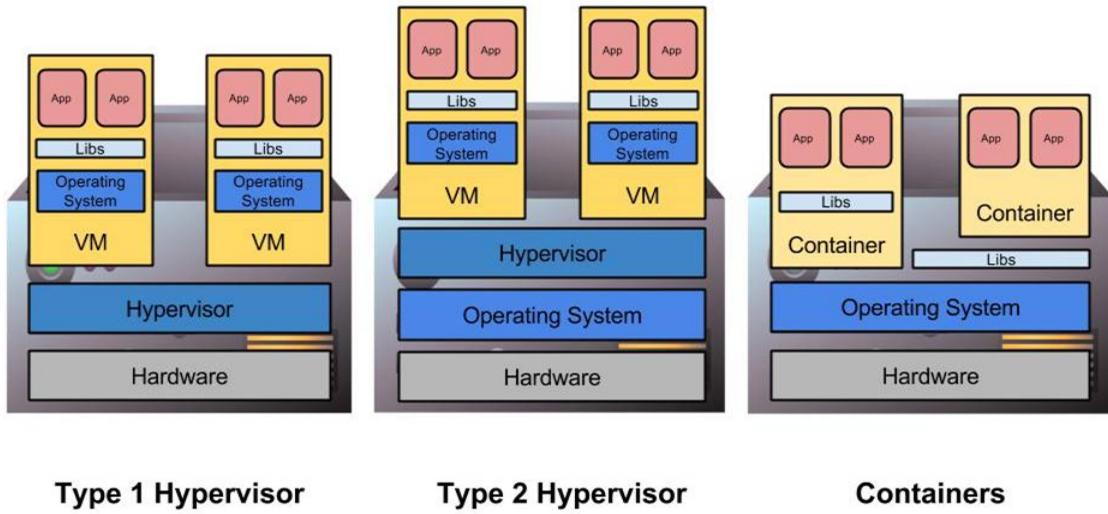
402 the computation is virtualized. This can give the effect that resources are local,  
403 when really they are remote and abstract. This is how GPUs are "virtualized"  
404 with tools such as rCUDA [50] and vCUDA [51]. This method's performance  
405 often can be very dependent on the underlying communications mechanisms, as  
406 well as complete virtualization of the whole API (which may be difficult). API  
407 virtualization also exists for entire OS libraries to translate between differing  
408 OS, without actual containerization.

409 • **Process -** Virtualization can also take the form at the process or user level,  
410 which can then be used to deliver a specialized or high level language that is the  
411 same across several different OSs. This is how the Java Virtual Machine (JVM)  
412 and Microsoft's .NET platform work. This type of application predominantly  
413 falls outside of the scope of this dissertation.

### 414 2.1.1 Hypervisors and Containers

415 While there are many types of virtualization, this dissertation predominately focuses  
416 on hardware and OS level virtualization. With hardware virtualization, a Virtual  
417 Machine Monitor, or hypervisor, is used. Abstractly, a hypervisor is a piece of software  
418 that creates virtual machines or *guests*, usually that model the underlying physical  
419 machine *host* system.

420 Hardware virtualization can actually be dissected in more detail to two different  
421 types of hypervisors, and they can be directly compared to containers as seen in  
422 Figure 2.2. With a Type 1 virtualization system, the hypervisor or VMM sits directly  
423 on the bare-metal hardware, below the OS. These native hypervisors provide direct  
424 control of the underlying hardware, and are controlled and operated usually through  
425 the use of a privileged VM. One example of a type 1 hypervisor is the Xen Virtual



**Figure 2.2** Hypervisors and Containers

426 Machine Monitor [48], which uses its VMM as well as a privileged Linux OS, called  
 427 Dom0, to create and manage other user VMs, or running DomU instances. The  
 428 VMM in this place provides the necessary hardware abstraction of CPU, memory,  
 429 and some I/O aspects, leaving the control aspects of the other DomUs to Dom0.  
 430 With a type 1 hypervisor, all virtualization functions are kept separate from control  
 431 and OS functionality, effectively making a cleaner design. This design could lead to  
 432 end application performance implications, as illustrated with the Palacios VMM [52].

433 Type 2 hypervisors utilize a different, and sometimes more convoluted design.  
 434 With a type 2 hypervisor, there is a "host" OS that, like with native OSs, sits di-  
 435 rectly atop hardware. This OS is just like any normal native OS. However, the OS  
 436 itself can abstract its own hardware, and provide and manage a VM, effectively as  
 437 an OS process. In this case, the hypervisor providing the abstraction is effectively  
 438 hosted within, atop, or as a module part of a given OS. There are many different  
 439 type 2 hypervisors, the most common of which is the Linux Kernel Virtual Machine  
 440 (KVM) [53]. KVM is often used in conjunction with QEMU, a ISA level hypervisor  
 441 to provide some basic ISA level virtualization and emulation capabilities. KVM is

442 simply provided as a Linux kernel module within a given host, and guest VMs are  
443 run as a single process on the host OS.

444 Both types of hypervisors are very distinct from OS level virtualization, also known  
445 as containers. With containers, there is a single OS, however instead of direct hard-  
446 ware abstraction,a single kernel is used to simultaneously run multiple user-space  
447 instances in a jailed-root environment. These environments may look and feel like  
448 a separate machine, but in fact are not. Often times the kernel itself provides re-  
449 source management tools to help control resource utilization and allocations. Linux  
450 containers (LXC) provide a great example of this, with their use of namespaces for  
451 filesystem control and cgroups for resource management. With the recent advent of  
452 Docker, which looks to control versioning and easy deployment of traceable contain-  
453 ers, this aspect of OS level virtualization has grown in popularity, however security  
454 and usability concerns still exist.

## 455 2.2 Cloud Computing

456 Cloud computing is one of the most explosively expanding technologies in the com-  
457 puting industry today. However it is important to understand where it came from, in  
458 order to figure out where it will be heading in the future. While there is no clear cut  
459 evolutionary path to Clouds, many believe the concepts originate from two specific  
460 areas: Grid Computing and Web 2.0.

461 Grid computing [54,55], in its practical form, represents the concept of connect-  
462 ing two or more spatially and administratively diverse clusters or supercomputers  
463 together in a federating manner. The term “the Grid” was coined in the mid 1990’s  
464 to represent a large distributed systems infrastructure for advanced scientific and en-  
465 gineering computing problems. Grids aim to enable applications to harness the full

466 potential of resources through coordinated and controlled resource sharing by scalable  
467 virtual organizations. While not all of these concepts carry over to the Cloud, the  
468 control, federation, and dynamic sharing of resources is conceptually the same as in  
469 the Grid. This is outlined by [3], as Grids and Clouds are compared at an abstract  
470 level and many concepts are remarkably similar. From a scientific perspective, the  
471 goals of Clouds and Grids are also similar. Both systems attempt to provide large  
472 amounts of computing power by leveraging a multitude of sites running diverse ap-  
473 plications concurrently in symphony. The only significant differences between Grids  
474 and Clouds exist in the implementation details, and the reproductions of them, as  
475 outlined later in this section.

476 The other major component, Web 2.0, is also a relatively new concept in the  
477 history of Computer Science. The term Web 2.0 was originally coined in 1999 in a  
478 futuristic prediction by Dracy DiNucci [56]: “The Web we know now, which loads  
479 into a browser window in essentially static screenfulls, is only an embryo of the Web  
480 to come. The first glimmerings of Web 2.0 are beginning to appear, and we are just  
481 starting to see how that embryo might develop. The Web will be understood not  
482 as screenfulls of text and graphics but as a transport mechanism, the ether through  
483 which interactivity happens. It will [...] appear on your computer screen, [...] on your  
484 TV set [...] your car dashboard [...] your cell phone [...] hand-held game machines  
485 [...] maybe even your microwave oven.” Her vision began to form, as illustrated in  
486 2004 by the O’Riley Web 2.0 conference, and since then the term has been a pivotal  
487 buzz word among the internet. While many definitions have been provided, Web 2.0  
488 really represents the transition from static HTML to harnessing the Internet and the  
489 Web as a platform in of itself.

490 Web 2.0 provides multiple levels of application services to users across the Inter-  
491 net. In essence, the web becomes an application suite for users. Data is outsourced

492 to wherever it is wanted, and the users have total control over what they interact  
493 with, and spread accordingly. This requires extensive, dynamic and scalable host-  
494 ing resources for these applications. This demand provides the user-base for much  
495 of the commercial Cloud computing industry today. Web 2.0 software requires ab-  
496 stracted resources to be allocated and relinquished on the fly, depending on the Web's  
497 traffic and service usage at each site. Furthermore, Web 2.0 brought Web Services  
498 standards [57] and the Service Oriented Architecture (SOA) [58] which outline the  
499 interaction between users and cyber infrastructure. In summary, Web 2.0 defined  
500 the interaction standards and user base, and Grid computing defined the underlying  
501 infrastructure capabilities.

502 Cloud computing [59] is one of the most explosively expanding technologies in  
503 the computing industry today. A Cloud computing implementation typically enables  
504 users to migrate their data and computation to a remote location with some varying  
505 impact on system performance [60]. This provides a number of benefits which could  
506 not otherwise be achieved:

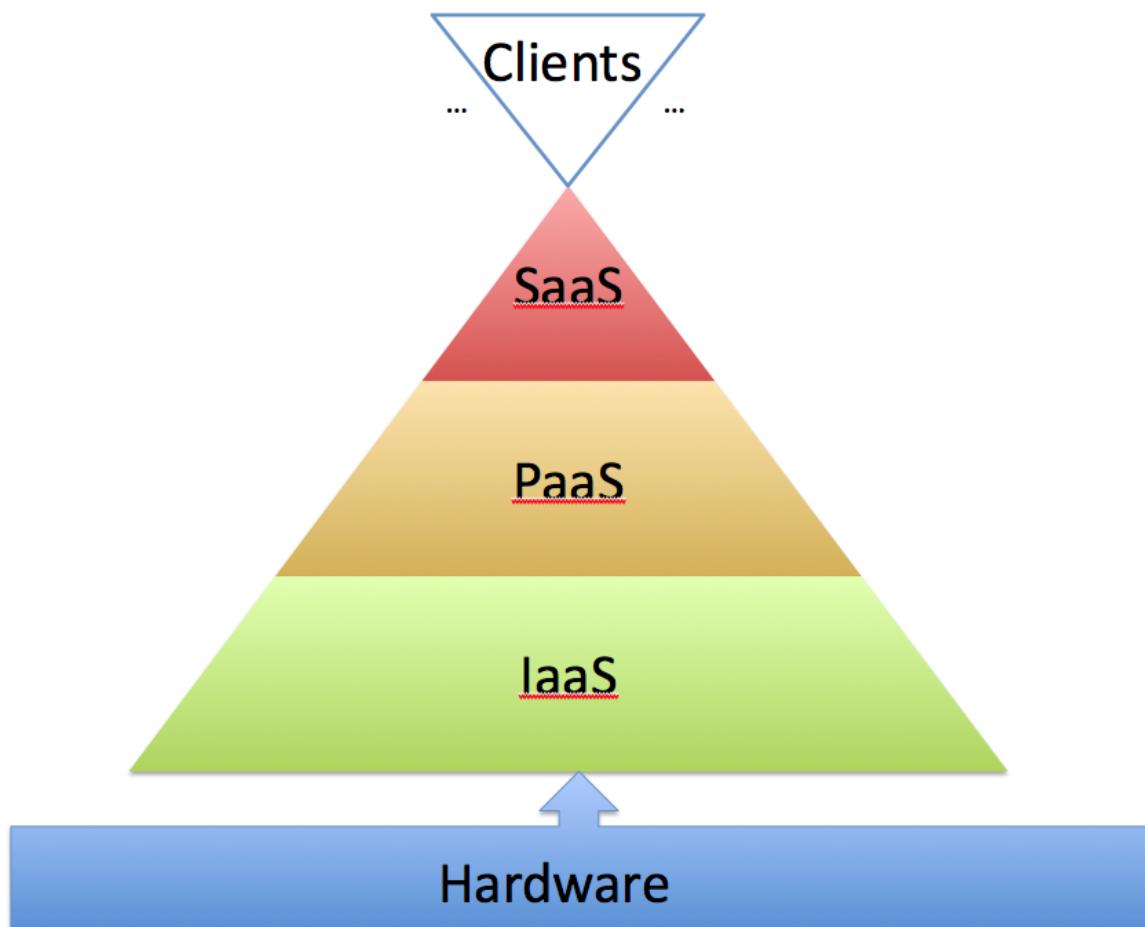
- 507     ● *Scalable* - Clouds are designed to deliver as much computing power as any  
508       user needs. While in practice the underlying infrastructure is not infinite, the  
509       cloud resources are projected to ease the developer's dependence on any specific  
510       hardware.
- 511     ● *Quality of Service (QoS)* - Unlike standard data centers and advanced com-  
512       puting resources, a well-designed Cloud can project a much higher QoS than  
513       traditionally possible. This is due to the lack of dependence on specific hard-  
514       ware, so any physical machine failures can be mitigated without the prerequisite  
515       user awareness.
- 516     ● *Specialized Environment* - Within a Cloud, the user can utilize customized tools

517 and services to meet their needs. This can be to utilize the latest library, toolkit,  
518 or to support legacy code within new infrastructure.

- 519 • *Cost Effective* - Users finds only the hardware required for each project. This  
520 reduces the risk for institutions potentially want build a scalable system, thus  
521 providing greater flexibility, since the user is only paying for needed infrastruc-  
522 ture while maintaining the option to increase services as needed in the future.
- 523 • *Simplified Interface* - Whether using a specific application, a set of tools or  
524 Web services, Clouds provide access to a potentially vast amount of computing  
525 resources in an easy and user-centric way. We have investigated such an interface  
526 within Grid systems through the use of the Cyberaide project [61, 62].

527 Many of the features noted above define what Cloud computing can be from a  
528 user perspective. However, Cloud computing in its physical form has many different  
529 meanings and forms. Since Clouds are defined by the services they provide and not  
530 by applications, an integrated as-a-service paradigm has been defined to illustrate the  
531 various levels within a typical Cloud, as in Figure 2.3.

- 532 • *Clients* - A client interacts with a Cloud through a predefined, thin layer of  
533 abstraction. This layer is responsible for communicating the user requests and  
534 displaying data returned in a way that is simple and intuitive for the user.  
535 Examples include a Web Browser or a thin client application.
- 536 • *Software-as-a-Service (SaaS)* - A framework for providing applications or soft-  
537 ware deployed on the Internet packaged as a unique service for users to consume.  
538 By doing so, the burden of running a local application directly on the client's  
539 machine is removed. Instead all the application logic and data is managed



**Figure 2.3** View of the Layers within a Cloud Infrastructure

540 centrally and to the user through a browser or thin client. Examples include  
541 Google Docs, Facebook, or Pandora.

- 542 • *Platform-as-a-Service (PaaS)* - A framework for providing a unique computing  
543 platform or software stack for applications and services to be developed on.  
544 The goal of PaaS is to alleviate many of the burdens of developing complex,  
545 scalable software by proving a programming paradigm and tools that make ser-  
546 vice development and integration a tractable task for many. Examples include  
547 Microsoft Azure and Google App Engine.
- 548 • *Infrastructure-as-a-Service (IaaS)* - A framework for providing entire computing

549 resources through a service. This typically represents virtualized Operating  
550 Systems, thereby masking the underlying complexity details of the physical  
551 infrastructure. This allows users to rent or buy computing resources on demand  
552 for their own use without needing to operate or manage physical infrastructure.  
553 Examples include Amazon EC2, and OpenStack, and is the major cloud focal  
554 point for this dissertation.

- 555 • *Physical Hardware* - The underlying set of physical machines and IT equipment  
556 that host the various levels of service. These are typically managed at a large  
557 scale using virtualization technologies which provide the QoS users expect. This  
558 is the basis for all computing infrastructure.

559 When all of these layers are combined, a dynamic software stack is created to  
560 focus on large scale deployment of services to users.

### 561 2.2.1 Infrastructure-as-a-Service

562 Today there are a number of Clouds that offer solutions for Infrastructure-as-a-Service  
563 (IaaS). There have been multiple comparison efforts between various IaaS service  
564 [4, 63–65] which provide insight to the similarities and differences between the long  
565 array of cloud infrastructure deployment solutions. However at a high level, IaaS can  
566 be split into 3 tiers, based on their availability.

- 567 • **Public** - Public IaaS is where the services and virtualizaton of hardware re-  
568 sources are provided over the internet. Usually this is in a centralized data  
569 centers whereby many users concurrently access these resources from across  
570 the globe, often at a pre-negotiated price point and service level agreement  
571 (SLA) [66]. Public clouds are the best utilization of economies of scale, by

572 selling hosting services en-masse and thereby offering competitive costs. The  
573 Amazon Elastic Compute Cloud (EC2) is a primary example of a public cloud.

- 574     ● **Private** - Private IaaS is where the cloud infrastructure is limited to within  
575       a distinct group, set of users, business, or virtual organization. Usually such  
576       private cloud infrastructure is also on a private, dedicated network that can  
577       either be separate or connected to other services. Data within private clouds  
578       are often more secure than those on a public cloud, and as such can be the choice  
579       for many users who have sensitive data, or who large-scale users who find better  
580       cost, performance, or QoS than what's provided within public clouds.
  
- 581     ● **Hybrid** - Hybrid cloud IaaS combines the computational power of both private  
582       and public IaaS, enabling users to keep data or costs within a private IaaS, but  
583       then "burst" usage to a public cloud on peak computational demands. Virtual  
584       Private Networks (VPN) can be useful to try to handle such a hybrid cloud not  
585       only for management and network addressing but also for security.

586 **Amazon EC2**

587 The Amazon Elastic Compute Cloud (EC2) [67], is probably the most popular of  
588 cloud infrastructure to date, and is used extensively in the IT industry. EC2 is the  
589 central component of Amazon Web Services platform. EC2 allows users to effectively  
590 rent virtual machines (called instances), hosted within Amazon's data centers, at a  
591 certain price point. Through their advanced UI or a RESTful API, users can start,  
592 stop, pause, migrate, and destroy instances exactly as needed, and to match the  
593 required computational tasks at hand.

594 Amazon EC2 predominantly relies on the Xen hypervisor to provide VMs on  
595 demand to users, with an equivalent compute unit equal to a 1.7Ghz Intel Xeon

596 processor. However, recent advancements, instance types, and upgrades to EC2 have  
597 increased this compute unit's power. Instance reservations have 3 types: On-demand,  
598 Reserved, and Spot. With On-demand instances, users pay by the hour for however  
599 long the desired instance is running. Users can instead rent reserved instances, where  
600 they pay a one-time (discounted) cost based on a pre-determined allocation time.  
601 There is also Spot pricing, where VMs are provisioned only when a given spot price is  
602 met, which is determined simply based on supply and demand within the EC2 system  
603 itself.

604 EC2 supports a wide array of user environments and setups. From an OS per-  
605 spective, this includes running Linux, Unix, and even Windows VM instances. EC2  
606 also provides instances with persistent storage through Elastic Block Storage (EBS)  
607 and Simple Storage Service (S3) object storage mechanisms. These tools are necessary  
608 for data persistence, as EC2 instances, as with most IaaS solutions, do not implic-  
609 itly persist data beyond the lifetime of the instance. EBS-rooted instances use an  
610 EBS volume as a root device and as such, persist there data beyond the lifetime of  
611 a given instance using the EBS volume as a persistent root partition. EC2 also of-  
612 fers elastic IPs, whereby public IP addresses are assigned to instances at boot (or in  
613 situ), however these elastic IPs do not require the DNS updates to propagate or an  
614 administrator to adjust the network.

615 These advanced features, coupled with the first-to-market viability and continual  
616 updates have made EC2 the largest cloud infrastructure today. While vendor lock-in  
617 is a concern (EC2 is not available for download or replication) and other alternatives  
618 exist such as Google's Compute Engine [68], the prevalence and support with EC will  
619 likely mean its status quo as the public cloud of choice will continue for the foreseeable  
620 future.

621 **Nimbus**

622 Nimbus [69, 70] is a set of open source tools that provide a private IaaS cloud com-  
623 putting solution. Nimbus is based on the concept of virtual workspaces previously  
624 introduced for Globus [70]. A virtual workspace is an abstraction of an execution  
625 environment that can be made dynamically available to authorized clients by using  
626 well-defined protocols. In this way, it can create customized environments by de-  
627 ploying virtual machines (VMs) among remote resources. To such an end, Nimbus  
628 provides a web interface called Nimbus Web. Its aim is to provide administrative and  
629 user functions in a friendly interface.

630 Within Nimbus, a storage cloud implementation called Cumulus [69] has been  
631 tightly integrated with the other central services, although it can also be used stan-  
632 dalone. Cumulus is compatible with the Amazon Web Services S3 REST API [71],  
633 but extends its capabilities by including features such as quota management. The  
634 Nimbus cloud client uses the Jets3t library [72] to interact with Cumulus. However,  
635 since it is compatible with S3 REST API, other interfaces like boto [73] or s2cmd [74]  
636 can also be used to interact with Nimbus.

637 Nimbus supports two resource management strategies. The first one is the default  
638 “resource pool” mode. In this mode, the service has direct control of a pool of  
639 virtual machine managers (VMM) nodes and it assumes it can start VMs. The other  
640 supported mode is called “pilot”. Here, the service makes requests, to a cluster’s  
641 Local Resource Management System (LRMS), to get a VMM available where deploy  
642 VMs.

643 Nimbus also provides an implementation of EC2’s interface that allows you to use  
644 clients developed for the real EC2 system on Nimbus based clouds.

645 **Eucalyptus**

646 Eucalyptus is a product from Eucalyptus Systems [75–77], that developed out of  
 647 a research project at the University of California, Santa Barbara. Eucalyptus was  
 648 initially aimed at bringing the cloud computing paradigm of computing to academic  
 649 super computers and clusters. Eucalyptus provides a Amazon Web Services (AWS)  
 650 complaint EC2 based web service interface for interacting with the Cloud service.

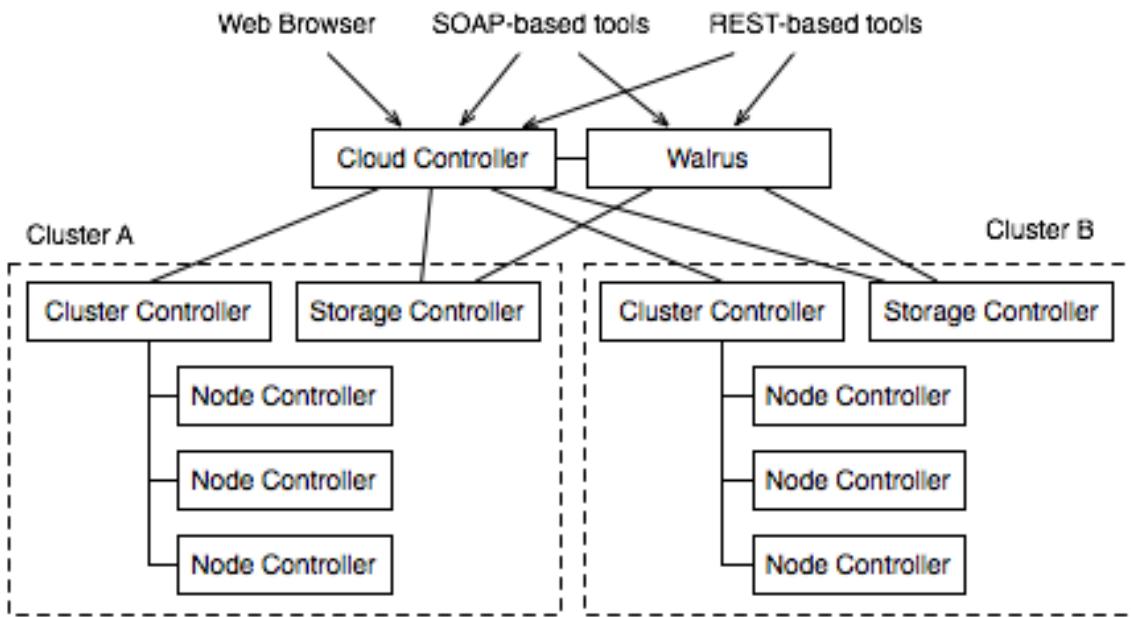


Figure 2.4 Eucalyptus Architecture

651 The architecture depicted in Figure 2.4 is based on a two level hierarchy of the  
 652 Cloud controller and the Cluster controller [78]. The Cluster Controller usually man-  
 653 ages the nodes within a single cluster and multiple such Cluster Controllers can be  
 654 used to connect to a single Cloud Controller. The Cloud Controller is responsible for  
 655 the resource management, scheduling and accounting aspects of the Cloud.

656 Being one of the first private cloud computing solutions, Eucalyptus has a focused  
 657 user interface. Much of Eucalyptus's design is based on the functionality of Amazon's  
 658 EC2 cloud solution, and the user interface is a prime example of that model. While

659 EC2 is a proprietary public cloud, it uses an open interface through the use of well  
660 designed Web Services which are open to all. Eucalyptus, looking to provide complete  
661 compatibility with EC2 to market the private cloud market, uses the same interface  
662 for all communication to the Cloud Controller. With Eucalyptus's interface being  
663 AWS complaint, it provides the same form of authentication that AWS supports,  
664 namely the shared key and PKI models.

665 While Eucalyptus can be controlled using the EC2 AMI tools, it also provides its  
666 own specific tool set; euca2ools. Euca2ools provides support for creating and man-  
667 aging keypairs, querying the cloud system, managing VMs, starting and terminating  
668 instances, network configuration, and block storage usage. The Eucalyptus system  
669 also provides a secure web front end to allow new users to create and manage account  
670 information, view available VMs, and download their security credentials.

671 As seen with the user interface, Eucalyptus takes many design queues from Ama-  
672 zons EC2 and the Image management system is no different. Eucalyptus stores images  
673 in Walrus, the block storage system that is analogous to the Amazon S3 service. As  
674 such, any user can bundle there own root filesystem, upload and then register this  
675 image and link that image with a particular kernel and ramdisk image. This image  
676 is uploaded into a user-defined bucket within Walrus, and can be retrieved anytime  
677 from any availability zone. This allows users to create specialty virtual appliances  
678 and deploy them within Eucalyptus with ease.

679 In 2014, Eucalyptus was acquired by Hewlett-Packard, which now maintains the  
680 HPE Helion Eucalyptus Cloud to have full compatibility with Amazon EC2. The  
681 most recent release of Helion eucalyptus is version 4.2.2 in early 2016.

682 **OpenStack**

683 OpenStack [79, 80], another private cloud infrastructure service, was introduced by  
684 Rackspace and NASA in July 2010. The project is trying to build an open source  
685 community spanning technologists, developers, researchers, and industry to share  
686 resources and technologies with the goal to create a massively scalable and secure  
687 cloud infrastructure. In tradition with other open source projects the entire software  
688 is open sources and limited to just open source API's such as Amazon.

689 Historically, OpenStack focuses on the development of two aspects of cloud com-  
690 puting to address compute and storage aspects with their OpenStack Compute and  
691 OpenStack Storage solutions. According to the documentation “OpenStack Com-  
692 pute is the internal fabric of the cloud creating and managing large groups of virtual  
693 private servers” and “OpenStack Object Storage is software for creating redundant,  
694 scalable object storage using clusters of commodity servers to store Terabytes or even  
695 petabytes of data.” However, OpenStack as a platform has evolved much more than  
696 its original efforts, and has created a wide array of new sub-projects.

697 As part of the computing support efforts OpenStack utilizes a cloud fabric con-  
698 troller known under the name Nova. The architecture for Nova is built on the concepts  
699 of shared-nothing and messaging-based information exchange. Hence most commu-  
700 nication in Nova are facilitated by message queues. To prevent blocking components  
701 while waiting for a response from others, deferred objects are introduced. Nova  
702 supports multiple scheduling paradigms, and includes plugins for a wide array of hy-  
703 pervisors, including Xen, KVM, and VMWare. The flexibility found within Nova  
704 is useful for supporting a wide array of cloud IaaS computational efforts. Recently,  
705 OpenStack has even looked to implement containers and bare-metal provisioning to  
706 keep on pace with the latest technologies.

707 The OpenStack Swift storage solution is build around a number of interacting com-

708 ponents and concepts including a Proxy Server, a Ring, Object Server, a Container  
709 Server, an Account Server, Replication, Updaters, and Auditors. This distributed  
710 architecture attempts to have no centralized components, to enable scalability and  
711 resiliency for data. Swift represents the long-term, object-based storage, similar to  
712 Amazon S3, and attempts to maintain rough API compatibility with S3. As Swift  
713 looks to use simple data replication as a main form of resiliency and fast read/write  
714 is rarely a priority, Swift is often built using commodity disk drives instead of most  
715 costly flash solutions.

716 With OpenStack Nova's increased prevalence, the number of auxiliary OpenStack  
717 projects has also increased to support Nova. While there are many other recent Open-  
718 Stack projects, these listed OpenStack efforts, along with Nova and Swift, represent  
719 the common core of a current OpenStack deployment. i

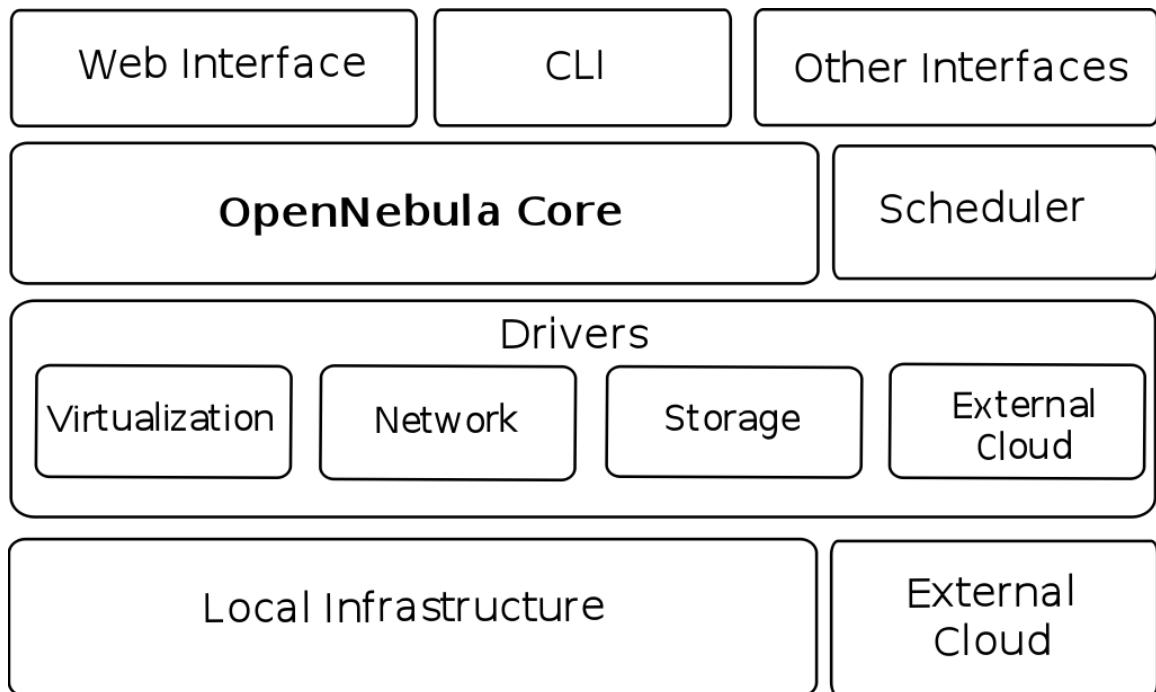
- 720     • **Cinder** for persistent block-level storage mechanisms to support VM instances  
721         and elastic block storage
- 722     • **Neutron** provides advanced networking and SDN solutions, IP addressing, and  
723         VLAN configuration
- 724     • **Glance** delivers comprehensive image management, including image discovery,  
725         registration, and delivery mechanisms
- 726     • **Keystone** identity and authentication service for all OpenStack services
- 727     • **Horizon**, a dashboard web-based UI framework, complimentary to the RESTful  
728         client API.

729 Currently, OpenStack exists as one of the largest ongoing private IaaS efforts,  
730 with over 500 companies contributing to the effort, and thousands of deployments.  
731 While releases have pushed forth approximately every 6 months, the latest current

732 release at the time of writing is *Mitaka*, which now includes full support for GPUs  
733 and SR-IOV interconnects. It is expected that OpenStack's prevalence in the cloud  
734 computing community will only increase in the next few years.

735 **OpenNebula**

736 OpenNebula [81, 82] is an open-source toolkit which allows to transform existing in-  
737 frastructure into an Infrastructure as a Service (IaaS) cloud with cloud-like interfaces.  
738 Figure 2.5 shows the OpenNebula architecture and their main components.



**Figure 2.5** OpenNebula Architecture

739 The architecture of OpenNebula has been designed to be flexible and modular to  
740 allow its integration with different storage and network infrastructure configurations,  
741 and hypervisor technologies. Here, the core is a centralized component that manage  
742 the virtual machine's (VM) full life cycle, including setting up networks dynamically  
743 for groups of VMs and managing their storage requirements, such as VM disk image

744 deployment or on-the-fly software environment creation. Another important compo-  
745 nent is the capacity manager, which governs the functionality provided by the core  
746 for scheduling. The default capacity scheduler is a requirement/rank matchmaker.  
747 However, it is also possible to develop more complex scheduling policies, through a  
748 lease model and advance reservations like Haizea [83]. The last main components are  
749 the access drivers. They provide an abstraction of the underlying infrastructure to  
750 expose the basic functionality of the monitoring, storage and virtualization services  
751 available in the cluster. Therefore, OpenNebula is not tied to any specific environ-  
752 ment and can provide a uniform management layer regardless of the virtualization  
753 platform.

754 Additionally, OpenNebula offers management interfaces to integrate the core's  
755 functionality within other data center management tools, such as accounting or mon-  
756 itoring frameworks. To this end, OpenNebula implements the libvirt API [84], an open  
757 interface for VM management, as well as a command line interface (CLI). A subset  
758 of this functionality is exposed to external users through a cloud interface. Due to  
759 its architecture, OpenNebula is able to adapt to organizations with changing resource  
760 needs, including the addition or failure of physical resources [65]. Some essential fea-  
761 tures to support changing environments are the live migration and the snapshotting  
762 of VMs [81]. Furthermore, when the local resources are insufficient, OpenNebula can  
763 support a hybrid cloud model by using cloud drivers to interface with external clouds.  
764 This lets organizations supplement the local infrastructure with computing capacity  
765 from a public cloud to meet peak demands, or implement high availability strategies.  
766 OpenNebula includes an EC2 driver, which can submit requests to Amazon EC2 and  
767 Eucalyptus [75], as well as an ElasticHosts driver [85].

768 Regarding the storage, an OpenNebula Image Repository allows users to easily  
769 specify disk images from a catalog without worrying about low-level disk configuration

770 attributes or block device mapping. Also, image access control is applied to the  
771 images registered in the repository, hence simplifying multi-user environments and  
772 image sharing. Nevertheless, users can also set up their own images.

773 **Others**

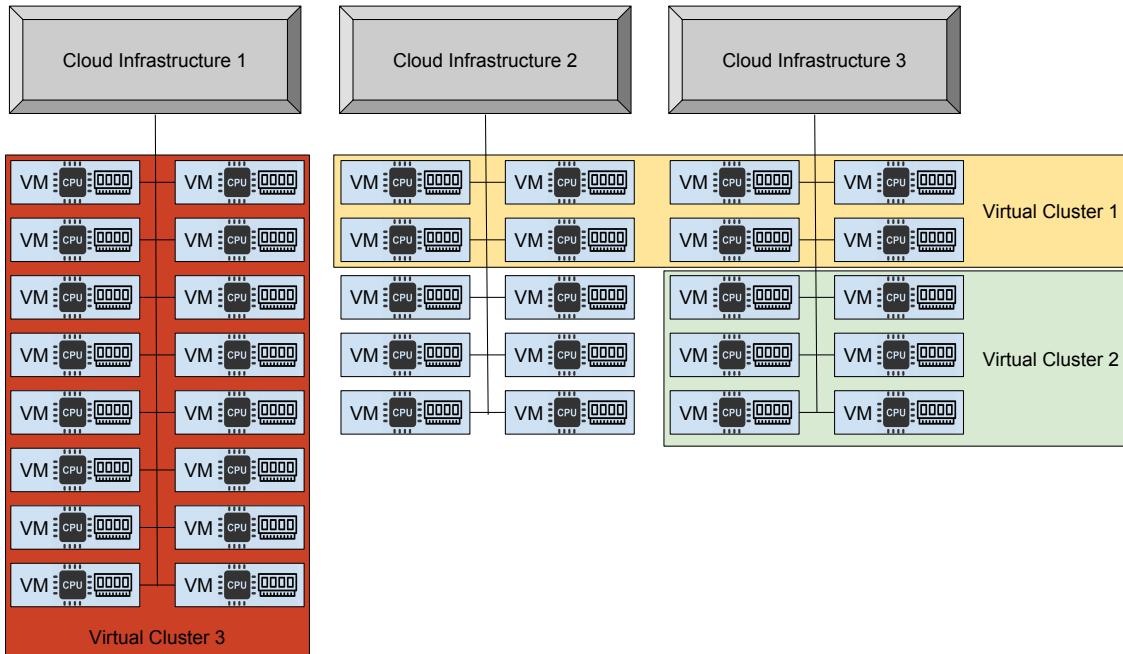
774 Other cloud specific projects exist such as In-VIGO [86], Cluster-on-Demand [87],  
775 and VMWare's own proprietary vCloud Air [88]. Each effort provides their own  
776 interpretation of private cloud services within a data center, often with the ability to  
777 interplay with public cloud offerings such as Amazon's EC2. Docker [44] also looks to  
778 provide IaaS capabilities with specialized and easily configurable containers, based on  
779 LXC and libcontainer solutions described in the previous section. While it is still to be  
780 determined how Docker and containers will change the private IaaS landscape, they  
781 do provide similar functionality for Linux users without some of the complexities of  
782 traditional virtualized IaaS.

783 **2.2.2 Virtual Clusters**

784 While virtualization and cloud IaaS provide many key advancements, this technology  
785 alone is not sufficient. Rather, a collective scheduling and management for virtual  
786 machines is required to piece together a working virtual cluster.

787 Let us consider a typical usage for a Cloud data center that is used in part to  
788 provide computational power for the Large Hadron Collider at CERN [89], a global  
789 collaboration from more than 2000 scientists of 182 institutes in 38 nations. Such a  
790 system would have a small number of experiments to run. Each experiment would  
791 require a very large number of jobs to complete the computation needed for the anal-  
792 ysis. Examples of such experiments are the ATLAS [90] and CMS [91] projects, which  
793 (combined) require Petaflops of computing power on a daily basis. Each job of an

794 experiment is unique, but the application runs are often the same. Therefore, virtual  
 795 machines are deployed to execute incoming jobs. There is a file server which provides  
 796 virtual machine templates. All typical jobs are preconfigured in virtual machine tem-  
 797 plates. When a job arrives at the head node of the cluster, a correspondent virtual  
 798 machine is dynamically started on a certain compute node within the cluster to ex-  
 799 ecute the job. While the LHC project and CERN's cloud effort is a formidable one,  
 800 it only covers pleasingly parallel HTC workloads, and often times HPC and big data  
 801 workloads can equally complex in differing ways.



**Figure 2.6** Virtual Clusters on Cloud Infrastructure

802 Cluster computing has become one of the core tools in distributed systems for  
 803 use in parallel computation. Cluster computing revolves around the desire to get more  
 804 computing power and better reliability by utilizing many computers together across a  
 805 network to achieve larger computational tasks. Clusters have manifested themselves  
 806 in many different ways, ranging from Beowulf clusters [6] which run using commodity

807 PCs to some of the TOP500 [8] supercomputing systems today. Virtual clusters rep-  
808 resents the growing need of users to effectively organize computational resources in an  
809 environment specific to their tasks at hand, instead of sharing a common architecture  
810 across many users. With the advent of modern virtualization, virtual clusters are  
811 deployed across a set of VMs in order to gain relative isolation and flexibility between  
812 disjoint virtual clusters. Virtual clusters, or a set of multiple cluster computing de-  
813 ployments on a single, larger physical cluster infrastructure, often have the following  
814 properties and attributes [47]:

- 815     ● Resources allocation based on a VM unit
- 816     ● Clusters built of many VMs together, or by provisioning physical nodes
- 817     ● Leverage local infrastructure management tools to provide a middleware solu-  
818         tion for virtual clusters
- 819             – Implementations could be a cloud IaaS such as OpenStack
- 820             – Some instances use a queueing system such as PBS
- 821     ● User experience based on virtual cluster management, not single VM manage-  
822         ment
- 823     ● Consolidates functionality on a smaller resource platform using multiple VMs
- 824     ● Can provide fault tolerance through VM migration and management
- 825     ● Can utilize dynamic scaling through the addition or deletion of VMs from the  
826         virtual cluster
- 827     ● Connection to back-end storage solution to provide virtual persistent storage

Given the properties, Virtual clusters can take on many forms, however a very simplified set of virtual clusters across cloud infrastructure are provided as a representation in Figure 2.6. This could lead to the simple provisioning of multiple disjoint OSs on a single physical resource. Virtual clusters generally have the ability to provide and manage their own user environment and tuned internal middleware. Virtual clusters may enable the separation of multiple tasks into separate VMs, which still in fact run on the same or similar underlying physical resources, effectively providing task isolation. Virtual clusters can be deployed to be persistence, stored, shared, or re-provisioned on demand. The size of a virtual cluster could potentially expand and contrast relative to the necessary resource requirements, taking advantage of elasticity found with virtualization. Furthermore, VM migration may enable fault tolerance in the event of physical machine errors if properly managed.

With virtual clusters, the capability to quickly deploy custom environments becomes critical. As such, efforts have been put forth to quickly configure and create VM images on-demand. This includes custom efforts with configuration engines such as CfEngine, Chef, Ansible, and others. Within FutureGrid, an image management system was defined to provide preconfigured VM images for cloud infrastructure using the BCFG2 engine [92].

Initially, virtual clusters were proposed for the use of Grid communities [27]. Specifically, Foster et al look to provide commodity clusters to various Virtual Organizations [93], whereby grid services can instantiate and deploy VMs. This design and implementation was further refined through the use of metadata and contextualization using appliances [94]. Some of these ideas have even come to take shape in larger scale supercomputing deployments, such as with SDSC Comet’s virtual cluster availability [95].

Virtual Clusters, require orchestration services to be able to organize, deploy,

854 manage, and re-play the desired user environment, and there have been a number of  
855 efforts to bring this orchestration to utility. One efforts within FutureGrid is with  
856 the experiment management design [38], which attempts to define how resources are  
857 connected to and monitored, as well as how VMs are stored in a repository and  
858 provisioned across multiple heterogeneous resources. This effort moved forward with  
859 Cloudmesh [40], which provides a simple client interface to access multiple cloud  
860 resources with a command-line shell interface.

861 Another virtual cluster orchestration effort has developed within the OpenStack  
862 private IaaS solution itself, named OpenStack Heat [39]. Heat provides a method  
863 by which you can deploy "stacks", which could essentially be virtual clusters, us-  
864 ing OpenStack infrastructure. Specifically, Heat provides a human readable and  
865 machine-accessible template for specifying environments and requirements, as well  
866 as a RESTful API. In submitting a Heat orchestration template to the API, heat  
867 will interpret and build the designed custom environment within a given OpenStack  
868 cloud deployment. Kubernetes [96], a related effort, is a infrastructure orchestration  
869 framework for managing containerized applications within Docker.

870 **2.2.3 The FutureGrid Project**

871 FutureGrid was a NSF-funded national-scale Grid and Cloud test-bed facility that  
872 included a number of computational resources across many distributed locations.  
873 This FutureGrid test-bed allowed users can evaluate differing systems for applica-  
874 bility with their given research task or application. These areas include computer  
875 science research topics ranging from authentication, authorization, scheduling, virtu-  
876 alization, middleware design, interface design and cybersecurity, to the optimization  
877 of grid-enabled and cloud-enabled computational schemes for Astronomy, Chemistry,  
878 Biology, Engineering, High Energy Physics, or Atmospheric Science. This project

879 started at an opportune time, when cloud infrastructure was still in its experimental  
 880 stages and its applicability to mid-tier scientific efforts were unknown.

881 The FutureGrid features a unique WAN network structure that lent itself to a  
 882 multitude of experiments specifically for evaluating middleware technologies and ex-  
 883 periment management services. This network can be dedicated to conduct experi-  
 884 ments in isolation, using a network impairment device for introducing a variety of  
 885 predetermined network conditions. Figure 2.7 depicts the geographically distributed  
 886 resources that are outlined in Table 2.1 in more detail. All network links within Fu-  
 887 tureGrid are dedicated 10GbE links with the exception of a shared 10GbE link to  
 888 TACC over the TeraGrid [97,98] network, enabling high-speed data management and  
 889 transfer between each partner site within FutureGrid.



**Figure 2.7** FutureGrid Participants, Network, and Resources

890 Although the total number of systems within FutureGrid is comparatively conser-  
 891 vative, they provide some heterogeneity to the architecture and are connected by the  
 892 high-bandwidth network links. One important feature to note is that most systems

**Table 2.1** FutureGrid hardware

System type	Name	CPUs	Cores	TFLOPS	RAM(GB)	Disk(TB)	Site
IBM iDataPlex	India	256	1024	11	3072	†335	IU
Dell PowerEdge	Alamo	192	1152	12	1152	15	TACC
IBM iDataPlex	Hotel	168	672	7	2016	120	UC
IBM iDataPlex	Sierra	168	672	7	2688	72	UCSD
Cray XT5m	Xray	168	672	6	1344	†335	IU
ScaleMP vSMP	Echo	32	192	3	5872	192	IU
Dell PoweEdge	Bravo	32	128	2	3072	192	IU
SuperMicro	Delta	32	192	‡20	3072	128	IU
IBM iDataPlex	Foxtrot	64	256	2	768	5	UF
Total		1112	4960	70	23056	1394	

†Indicates shared file system. ‡Best current estimate

893 can be dynamically provisioned, e.g. these systems can be reconfigured when needed  
 894 by special software that is part of FutureGrid with proper access control by users  
 895 and administrators. Therefore its believed that this hardware infrastructure can fully  
 896 accommodate the needs of an experiment management system.

897 As of Fall 2014, the FutureGrid project has ended. The computing resources  
 898 and facilities at Indiana University have continued on as FutureSystems, continuing  
 899 to provide a cloud, big data, and HPC testbed to approved researchers. Recently  
 900 with new projects as part of the digital Science Center, the FutureSystems effort has  
 901 added two new machines, *Romeo* and *Juliet*, presenting clusters with Intel Haswell  
 902 CPU architectures to be used for big data research.

903 **2.3 High Performance Computing**

904 **2.3.1 Brief History of Supercomputing**

905 Supercomputing itself can date back to some of the forefront of computing itself,  
906 especially if we consider the ENIAC [99], the first Turing Complete general purpose  
907 digital computer, also as the first supercomputer. ENIAC was first deployed to cal-  
908 culate artillery firing tables, but was later used during the Second World War for  
909 helping the Manhattan project's thermonuclear calculations and later dedicated to  
910 the University of Pennsylvania after the war.

911 The first properly termed supercomputer was the Control Data Corporation's 6600  
912 mainframe [100], first delivered to CERN in 1965. The CDC 6600 was notably faster  
913 than the IBM counterparts, and the first deployments were able to perform on the  
914 order of 1 MFLOP. Interestingly, the CPU design that came from Seymour Cray's  
915 CDC 6600 took advantage of a simplified yet fast CPU design with silicon-based  
916 transistors, which founded the basis of the RISC processor architecture.

917 Cray's efforts eventually lead him to start his own company, and in 1975 released  
918 the Cray 1 system [101]. The Cray 1 took the powerful aspects of vector processing  
919 and memory pipelining from the STAR architecture (developed later by CDC) and in-  
920 troduced scalar performance through splitting vectors and instruction chaining. This  
921 resulted in an overall performance of around 250 MFLOPS at peak, but realistically  
922 closer to 1300 MFLOPS for general applications. The Cray 1 system also helped push  
923 forward the integrated circuit design, which was finally performant to be used, and  
924 also required an entirely new freon based coolant system.

925 The Cray 1 system gave way to the Cray X-MP and Y-MP in the mid 1980s.  
926 These machines were shared-memory vector processors, with two processors in the  
927 X-MP and up to 8 processors for the later Y-MP systems. These first shared-memory

928 systems were aided by increased memory speeds. The X-MP machine was capable  
929 of 200 MFLOPS sustained and 400 MFLOPS peak performance, whereas the Y-MP  
930 variants were capable of over 2 GFLOPS.

931 Also concurrently during the 1980s the advent of distributed memory architectures  
932 for supercomputing were starting to emerge. Specifically work on Caltech's Cosmic  
933 Cube, also known as a Hypercube, by Seitz and Fox [102,103], started the movement of  
934 concurrent or parallel computing. The Cosmic Cube leveraged new VLSI techniques  
935 and assembled Intel 8086/87 processors together with a novel hypercube interconnect  
936 which required no switching, creating one of the first truly parallel computers. SIMD  
937 programming and computation was done through a novel message passing architec-  
938 ture, instead of shared variables. This design was first commercialized with Intel's  
939 iPSC, and later contributing directly to designs in the Intel Paragon, ASCI Red, and  
940 Cray T3D/E systems.

941 While concurrent and parallel processor supercomputers continued into the 90s  
942 with aforementioned hypercube designs and the IBM Thinking Machines [104], a new  
943 commodity-based strategy emerged with Becker and Sterling's Beowulf clusters [6].  
944 Effectively, a Beowulf cluster is simply a cluster of commodity x86 machines linked  
945 together with a simplified LAN network. Beowulf clusters often (but not always)  
946 run Linux OS and leverage Ethernet solutions, and are programmed using a message  
947 passing construct such as MPI [16]. This allows for the building of massively parallel  
948 systems with relatively low cost and investment. Many specialized clusters today  
949 still utilize commodity x86 hardware and Linux OSs similar to the original Beowulf  
950 systems.

951 Concurrency and parallel computation on supercomputing resources has only  
952 flourished since. This has been even more pronounced as CPU clock frequencies  
953 stabilized and multi-core architectures took hold, driving the need for concurrency

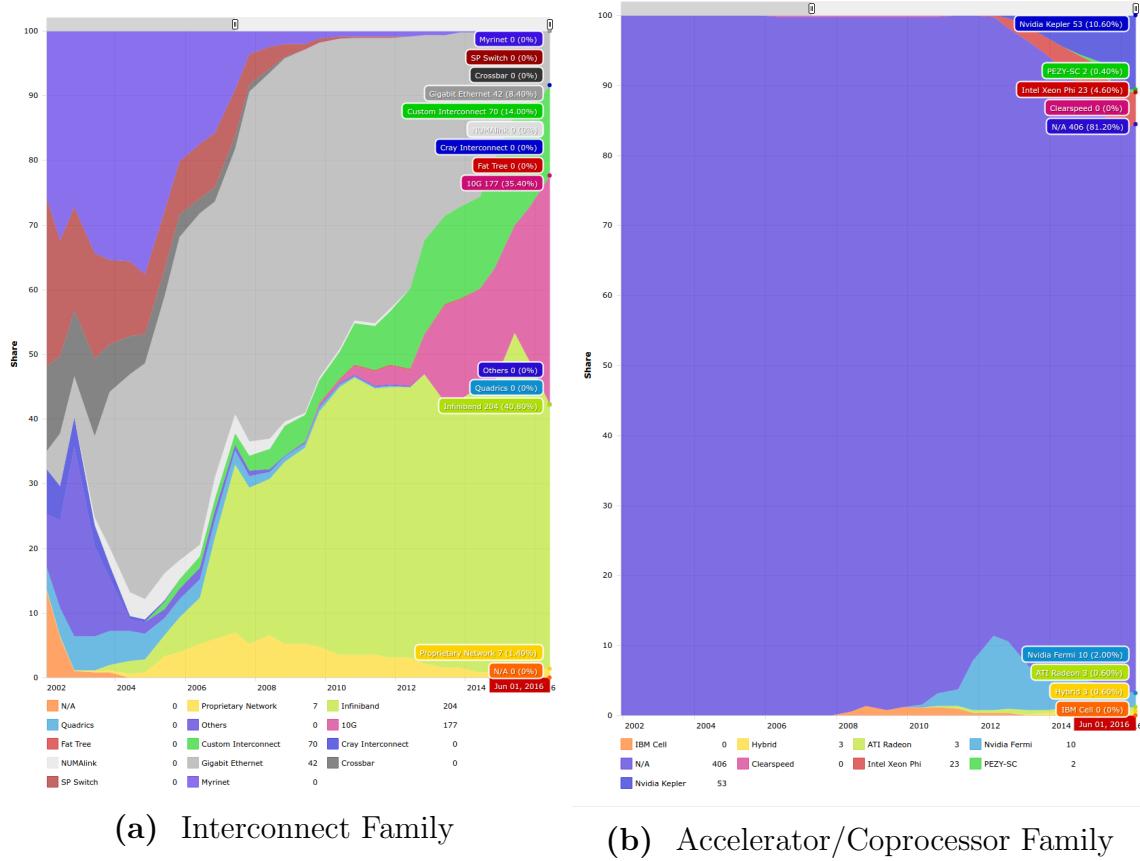
954 not only at an increased rate for supercomputing, but also even for commodity sys-  
955 tems. While commodity single-CPU, multi-core systems often look towards exploiting  
956 shared memory parallelism, distributed memory architectures have become a way of  
957 life for high performance computing.

### 958 2.3.2 Distributed Memory Computation

959 Abstractly, distributed memory architectures consist of multiple instances of a pro-  
960 cessor, memory, and an interconnect which allows each instance to perform inde-  
961 pendent computations and communicate over the interconnect. These interconnects  
962 could be built using point-to-point, or through more complex and advanced switching  
963 hardware, building a larger topology. While a simple example could involve several  
964 commodity PCs connected through Ethernet switch, this model scales to the latest  
965 supercomputing resources of today with millions of cores [105].

966 Programming such distributed memory systems is a nontrivial task that the  
967 greater HPC community has been wrangling for years. In the early days, this took  
968 the form of either the Message Passing Interface (MPI) [16] or PVM [106], however  
969 MPI has been far more successful and dominates the market for distributed memory  
970 parallel computation. MPI is a standardized message passing system, which defines  
971 the semantics and syntax for writing parallel programs in C, C++, or Fortran. MPI  
972 has even been implemented in other languages such as Java [107]. As MPI is a stan-  
973 dardization, there are many implementations that exist, including OpenMPI [108],  
974 MPICH [109], and MVAPICH [110], to name a few. Many MPI-enabled applications  
975 have been shown to be, with proper and careful design, the most efficient way to run  
976 a parallel application across a large subset of tightly coupled distributed resources,  
977 and can often represent the status-quo for HPC applications today.

978 More recently, the MPI programming model has been modified or joined with



**Figure 2.8** Top 500 Development over time from 2002 to 2016 [8]

other models. This change or deviation largely revolves around the hardware that has changed within HPC resources themselves to meet the need for more computational power. This includes hybrid MPI + OpenMP models which become useful as multi-core and many-core technologies becomes increasingly available [111], or a MPI+CUDA model which interleaves message passing with GPU utilization [112]. As some of the latest supercomputing resources have been deployed specifically with Nvidia GPUs for GPGPU programming such as the ORNL’s Titan system [113], the need for distributed memory computation with GPU computing has increased.

To get an idea of some of the hardware advances over the past few years, it is useful to examine the Top 500 [8], a comprehensive list of the top 500 supercomputers known and their key characteristics. Looking at the past decade in HPC, we can see

990 some trends emerge within hardware. Specifically looking at Interconnect and Copro-  
991 cessor architectures in Figure 2.8, we see a notable jump in the number of deployed  
992 system that are using both InfiniBand and GPUs in the past decade. Specifically, In-  
993 finiBand usage has increased to roughly 40% of the total number of deployed systems  
994 and remained somewhat stable as an interconnect family. Concurrently, the use of  
995 accelerators has increased from only 1 in 500 systems a decade ago, to almost a 20%  
996 of the top 500 systems that are using coprocessors. Within that factor, the majority  
997 of such accelerator-equipped systems have been using GPUs, with a concurrent in-  
998 crease in the Intel Xeon Phi coprocessor as well. While these do not represent all of  
999 supercomputing nor exclusively the most high end of systems, they do represent ad-  
1000 vanced hardware that has been increasingly common in the last decade, yet relatively  
1001 underutilized in comparison within cloud infrastructure. As such, much of the effort  
1002 in this dissertation focuses on, but is not limited to, these two technology families.

1003 **2.3.3 Exascale**

1004 As the forefront of supercomputing moves beyond the latest petaflop machines in the  
1005 past few years [113], the HPC community is setting their sights on the next signifi-  
1006 cant milestone: Exascale. Exascale computing refers broadly to performing roughly  
1007 one exaFLOPS, or  $10^{18}$  floating point operations per second. However, exascale itself  
1008 is far more than just a theoretical FLOPS goal, but instead a set of new comput-  
1009 ing advancements and challenges that requires reaching computational power at such  
1010 magnitude. While FLOPs are often used as the ubiquitous yard stick for supercom-  
1011 puting with the LINPACK benchmark [114], other efforts have taken hold to classify  
1012 systems under a different set of parameters [115, 116], with the loose understanding  
1013 that these may incorporate a richer application set destined for exascale systems. This  
1014 could include for instance integer calculations at a similar scale to satisfy defence and

1015 intelligence perspectives, or graph processing at a

1016 With exascale, there are a number of barriers that exist with current technolo-  
1017 gies that must be overcome to reach exascale. The exascale Computing Study [34]  
1018 specifically lists 4 major focal areas:

1019 1. Energy and Power Challenge

1020 • Describes the physical difficulties in providing the amount of power needed  
1021 to drive a sufficiently large exascale system. The US DOE estimates the  
1022 maximum power envelope for a deployed first exascale system to be within  
1023 20-40MW. Extrapolating current technology power utilization shows an  
1024 order of magnitude more energy utilization than the specificity power en-  
1025 envelop. As such, new architectures and conversation techniques will need  
1026 to be investigated.

1027 2. Memory and Storage Challenge

1028 • This challenge illustrates the problem that has grown in relation to the  
1029 memory wall, defined by the exponential difference between processor and  
1030 memory performance, as well as the storage capacity limits to support  
1031 calculations at the level of performance necessary. This challenge incorpo-  
1032 rates not only main memory limitations, but also tertiary storage issues as  
1033 well.

1034 3. Concurrency and Storage Challenge

1035 • This challenge is born from the recent limit in CPU clock rates as a way  
1036 to gain performance. Instead, performance must be gained through paral-  
1037 lelism. The depth of this challenge is especially profound when we consider  
1038 parallelism on the order of millions, if not billions of threads.

1039        4. Resiliency Challenge

- 1040            • The resiliency challenge defines the necessity of computation to recover  
1041            and continue in the event of a fault or fluctuation. As parallelism and  
1042            the number of individualized components substantially increases in a path  
1043            towards exascale, the mean time to failure of any given component also  
1044            increases.

1045        While current exascale efforts are as wide as they are varying, not only with  
1046        concepts, architectures, and runtime systems, but also with deployment plans and  
1047        expectations between future deployments. Of particular interest in current exas-  
1048        scale research is in Operating System and runtime (OS/R) developments to support  
1049        new extreme-scale applications in an efficient manner. Two examples of novel OS  
1050        approaches are the Hobbes project [117] and the ARGO Exascale Operating Sys-  
1051        tem [118]. These OS efforts, along with novel programming models for exascale such  
1052        as ParalleX [119] look to fundamentally change the relationship between HPC hard-  
1053        ware architectures and the libraries and applications to be leveraged on such future  
1054        exascale deployments.

1055        It is possible that virtualization itself may have an impact in OS and runtime ser-  
1056        vices in exascale [117]. While some of the work herein may be tangentially of utility  
1057        to such efforts, the immediate goal of this dissertation is not to investigate the appli-  
1058        cability of virtualization for exascale systems, but rather to enable the diversification  
1059        of HPC towards cloud infrastructure.

1060 **Chapter 3**

1061 **Analysis of Virtualization**

1062 **Technologies for High Performance**

1063 **Computing Environments**

1064 **3.1 Abstract**

1065 As Cloud computing emerges as a dominant paradigm in distributed systems, it is  
1066 important to fully understand the underlying technologies that make Clouds possible.  
1067 One technology, and perhaps the most important, is virtualization. Recently virtual-  
1068 ization, through the use of hypervisors, has become widely used and well understood  
1069 by many. However, there are a large spread of different hypervisors, each with their  
1070 own advantages and disadvantages. This chapter provides an in-depth analysis of  
1071 some of today's commonly accepted virtualization technologies from feature com-  
1072 parison to performance analysis, focusing on the applicability to High Performance  
1073 Computing environments using FutureGrid resources. The results indicate virtualiza-  
1074 tion sometimes introduces slight performance impacts depending on the hypervisor

1075 type, however the benefits of such technologies are profound and not all virtualization  
1076 technologies are equal.

1077 **3.2 Introduction**

1078 Cloud computing [59] is one of the most explosively expanding technologies in the  
1079 computing industry today. A Cloud computing implementation typically enables  
1080 users to migrate their data and computation to a remote location with some varying  
1081 impact on system performance [60]. This provides a number of benefits which could  
1082 not otherwise be achieved.

1083 Such benefits include:

- 1084     ● *Scalability* - Clouds are designed to deliver as much computing power as any  
1085       user needs. While in practice the underlying infrastructure is not infinite, the  
1086       cloud resources are projected to ease the developer's dependence on any specific  
1087       hardware.
- 1088     ● *Quality of Service (QoS)* - Unlike standard data centers and advanced com-  
1089       puting resources, a well-designed Cloud can project a much higher QoS than  
1090       traditionally possible. This is due to the lack of dependence on specific hard-  
1091       ware, so any physical machine failures can be mitigated without the prerequisite  
1092       user awareness.
- 1093     ● *Customization* - Within a Cloud, the user can utilize customized tools and  
1094       services to meet their needs. This can be to utilize the latest library, toolkit, or  
1095       to support legacy code within new infrastructure.
- 1096     ● *Cost Effectiveness* - Users finds only the hardware required for each project.  
1097       This reduces the risk for institutions potentially want build a scalable system,

1098 thus providing greater flexibility, since the user is only paying for needed in-  
1099 frastructure while maintaining the option to increase services as needed in the  
1100 future.

- 1101 • *Simplified Access Interfaces* - Whether using a specific application, a set of  
1102 tools or Web services, Clouds provide access to a potentially vast amount of  
1103 computing resources in an easy and user-centric way.

1104 While Cloud computing has been driven from the start predominantly by the in-  
1105 dustry through Amazon [67], Google [120] and Microsoft [121], a shift is also occurring  
1106 within the academic setting as well. Due to the many benefits, Cloud computing is  
1107 becoming immersed in the area of High Performance Computing (HPC), specifically  
1108 with the deployment of scientific clouds [122] and virtualized clusters [27].

1109 There are a number of underlying technologies, services, and infrastructure-level  
1110 configurations that make Cloud computing possible. One of the most important  
1111 technologies is virtualization. Virtualization, in its simplest form, is a mechanism to  
1112 abstract the hardware and system resources from a given Operating System. This is  
1113 typically performed within a Cloud environment across a large set of servers using a  
1114 Hypervisor or Virtual Machine Monitor (VMM), which lies in between the hardware  
1115 and the OS. From the hypervisor, one or more virtualized OSs can be started concur-  
1116 rently, leading to one of the key advantages of Cloud computing. This, along with the  
1117 advent of multi-core processors, allows for a consolidation of resources within any data  
1118 center. From the hypervisor level, Cloud computing middleware is deployed atop the  
1119 virtualization technologies to exploit this capability to its maximum potential while  
1120 still maintaining a given QoS and utility to users.

1121 The rest of this chapter is as follows: First, we look at what virtualization is,  
1122 and what current technologies currently exist within the mainstream market. Next

we discuss previous work related to virtualization and take an in-depth look at the features provided by each hypervisor. We follow this by outlining an experimental setup to evaluate a set of today's hypervisors on a novel Cloud test-bed architecture. Then, we look at performance benchmarks which help explain the utility of each hypervisor and the feasibility within an HPC environment. We conclude with our final thoughts and recommendations for using virtualization in Clouds for HPC.

### 3.3 Related Research

While the use of virtualization technologies has increased dramatically in the past few years, virtualization is not specific to the recent advent of Cloud computing. IBM originally pioneered the concept of virtualization in the 1960's with the M44/44X systems [123]. It has only recently been reintroduced for general use on x86 platforms. Today there are a number of public Clouds that offer IaaS through the use of virtualization technologies. The Amazon Elastic Compute Cloud (EC2) [124] is probably the most popular Cloud and is used extensively in the IT industry to this day. Nimbus [125] and Eucalyptus [75] are popular private IaaS platforms in both the scientific and industrial communities. Nimbus, originating from the concept of deploying virtual workspaces on top of existing Grid infrastructure using Globus, has pioneered scientific Clouds since its inception. Eucalyptus has historically focused on providing an exact EC2 environment as a private cloud to enable users to build an EC2-like cloud using their own internal resources. Other scientific Cloud specific projects exist such as OpenNebula [126], In-VIGO [127], and Cluster-on-Demand [87], all of which leverage one or more hypervisors to provide computing infrastructure on demand. In recent history, OpenStack [128] has also come to light from a joint collaboration between NASA and Rackspace which also provide compute and storage

1147 resources in the form of a Cloud.

1148 While there are currently a number of virtualization technologies available today,  
1149 the virtualization technique of choice for most open platforms over the past 5 years has  
1150 typically been the Xen hypervisor [48]. However more recently VMWare ESX [129]  
1151 <sup>1</sup>, Oracle VirtualBox [130] and the Kernel-based Virtual Machine (KVM) [53] are  
1152 becoming more commonplace. As these look to be the most popular and feature-  
1153 rich of all virtualization technologies, we look to evaluate all four to the fullest extent  
1154 possible. There are however, numerous other virtualization technologies also available,  
1155 including Microsoft's Hyper-V [131], Parallels Virtuozzo [132], QEMU [133], OpenVZ  
1156 [134], Oracle VM [135], and many others. However, these virtualization technologies  
1157 have yet to seen widespread deployment within the HPC community, at least in their  
1158 current form, so they have been placed outside the scope of this work.

1159 In recent history there have actually been a number of comparisons related to  
1160 virtualization technologies and Clouds. The first performance analysis of various hy-  
1161 pervisors started with, unsurprisingly, the hypervisor vendors themselves. VMWare  
1162 has happy to put out its own take on performance in [136], as well as the original  
1163 Xen article [48] which compares Xen, XenoLinux, and VMWare across a number of  
1164 SPEC and normalized benchmarks, resulting in a conflict between both works. From  
1165 here, a number of more unbiased reports originated, concentrating on server consol-  
1166 idation and web application performance [129, 137, 138] with fruitful yet sometimes  
1167 incompatible results. A feature base survey on virtualization technologies [139] also  
1168 illustrates the wide variety of hypervisors that currently exist. Furthermore, there  
1169 has been some investigation into the performance within HPC, specifically with In-  
1170 finiBand performance of Xen [140] and rather recently with a detailed look at the  
1171 feasibility of the Amazon Elastic Compute cloud for HPC applications [43], however

---

<sup>1</sup>Due to the restrictions in VMWare's licensing agreement, benchmark results are unavailable.

1172 both works concentrate only on a single deployment rather than a true comparison  
1173 of technologies.

1174 As these underlying hypervisor and virtualization implementations have evolved  
1175 rapidly in recent years along with virtualization support directly on standard x86  
1176 hardware, it is necessary to carefully and accurately evaluate the performance impli-  
1177 cations of each system. Hence, we conducted an investigation of several virtualization  
1178 technologies, namely Xen, KVM, VirtualBox, and in part VMWare. Each hypervisor  
1179 is compared alongside one another with base-metal as a control and (with the exception  
1180 of VMWare) run through a number of High Performance benchmarking tools.

## 1181 3.4 Feature Comparison

1182 With the wide array of potential choices of virtualization technologies available, its  
1183 often difficult for potential users to identify which platform is best suited for their  
1184 needs. In order to simplify this task, we provide a detailed comparison chart between  
1185 Xen 3.1, KVM from RHEL5, VirtualBox 3.2 and VMWWare ESX in Figure 2.

	Xen	KVM	VirtualBox	VMWare
<b>Para-virtualization</b>	Yes	No	No	No
<b>Full virtualization</b>	Yes	Yes	Yes	Yes
<b>Host CPU</b>	x86, x86-64, IA-64	x86, x86-64,IA64,PPC	x86, x86-64	x86, x86-64
<b>Guest CPU</b>	x86, x86-64, IA-64	x86, x86-64,IA64,PPC	x86, x86-64	x86, x86-64
<b>Host OS</b>	Linux, UNIX	Linux	Windows, Linux, UNIX	Proprietary UNIX
<b>Guest OS</b>	Linux, Windows, UNIX	Linux, Windows, UNIX	Linux, Windows, UNIX	Linux, Windows, UNIX
<b>VT-x / AMD-v</b>	Opt	Req	Opt	Opt
<b>Cores supported</b>	128	16	32	8
<b>Memory supported</b>	4TB	4TB	16GB	64GB
<b>3D Acceleration</b>	Xen-GL	VMGL	Open-GL	Open-GL, DirectX
<b>Live Migration</b>	Yes	Yes	Yes	Yes
<b>License</b>	GPL	GPL	GPL/proprietary	Proprietary

**Figure 3.1** A comparison chart between Xen, KVM, VirtualBox, and VMWare ESX

1186 The first point of investigation is the virtualization method of each VM. Each

1187 hypervisor supports full virtualization, which is now common practice within most  
1188 x86 virtualization deployments today. Xen, originating as a para-virtualized VMM,  
1189 still supports both types, however full virtualization is often preferred as it does  
1190 not require the manipulation of the guest kernel in any way. From the Host and  
1191 Guest CPU lists, we see that x86 and, more specifically, x86-64/amd64 guests are all  
1192 universally supported. Xen and KVM both support Itanium-64 architectures for full  
1193 virtualization (due to both hypervisors dependency on QEMU), and KVM also claims  
1194 support for some recent PowerPC architectures. However, we concern ourselves only  
1195 with x86-64 features and performance, as other architectures are out of the scope of  
1196 this manuscript. Of the x86-64 platforms, KVM is the only hypervisor to require  
1197 either Intel VT-X or AMD-V instruction sets in order to operate. VirtualBox and  
1198 VMWare have internal mechanisms to provide full virtualization even without the  
1199 virtualization instruction sets, and Xen can default back to para-virtualized guests.

1200 Next, we consider the host environments for each system. As Linux is the pri-  
1201 mary OS type of choice within HPC deployments, its key that all hypervisors sup-  
1202 port Linux as a guest OS, and also as a host OS. As VMWare ESX is meant to be  
1203 a virtualization-only platform, it is built upon a specially configured Linux/UNIX  
1204 proprietary OS specific to its needs. All other hypervisors support Linux as a host  
1205 OS, with VirtualBox also supporting Windows, as it was traditionally targeted for  
1206 desktop-based virtualization. However, as each hypervisor uses VT-X or AMD-V in-  
1207 structions, each can support any modern OS targeted for x86 platforms, including all  
1208 variants of Linux, Windows, and UNIX.

1209 While most hypervisors have desirable host and guest OS support, hardware sup-  
1210 port within a guest environment varies drastically. Within the HPC environment,  
1211 virtual CPU (vCPU) and maximum VM memory are critical aspects to choosing the  
1212 right virtualization technology. In this case, Xen is the first choice as it supports up

1213 to 128 vCPUs and can address 4TB of main memory in 64-bit modes, more than  
1214 any other. VirtualBox, on the other hand, supports only 32 vCPUs and 16GB of  
1215 addressable RAM per guest OS, which may lead to problems when looking to deploy  
1216 it on large multicore systems. KVM also faces an issue with the number of vCPU  
1217 supported limited to 16, recent reports indicate it is only a soft limit [141], so deploy-  
1218 ing KVM in an SMP environment may not be a significant hurdle. Furthermore, all  
1219 hypervisors provide some 3D acceleration support (at least for OpenGL) and support  
1220 live migration across homogeneous nodes, each with varying levels of success.

1221 Another vital juxtaposition of these virtualization technologies is the license agree-  
1222 ments for its applicability within HPC deployments. Xen, KVM, and VirtualBox are  
1223 provided for free under the GNU Public License (GPL) version 2, so they are open  
1224 to use and modification by anyone within the community, a key feature for many  
1225 potential users. While VirtualBox is under GPL, it has recently also offered with  
1226 additional features under a more proprietary license dictated by Oracle since its ac-  
1227 quirement from Sun last year. VMWare, on the other hand, is completely proprietary  
1228 with an extremely limited licensing scheme that even prevents the authors from will-  
1229 fully publishing any performance benchmark data without specific and prior approval.  
1230 As such, we have neglected VMWare form the remainder of this chapter. Whether  
1231 going with a proprietary or open source hypervisor, support can be acquired (usually  
1232 for an additional cost) with ease from each option.

### 1233 3.4.1 Usability

1234 While side by side feature comparison may provide crucial information about a poten-  
1235 tial user's choice of hypervisor, that may also be interested in its ease of installation  
1236 and use. We will take a look at each hypervisor from two user perspectives, a systems  
1237 administrator and normal VM user.

1238 One of the first things on any system administrator's mind on choosing a hypervi-  
1239 sor is the installation. For all of these hypervisors, installation is relatively painless.  
1240 For the FutureGrid support group, KVM and VirualBox are the easiest of the all  
1241 tested hypervisors to install, as there are a number of supported packages available  
1242 and installation only requires the addition of one or more kernel modules and the sup-  
1243 port software. Xen, while still supported in binary form by many Linux distributions,  
1244 is actually much more complicated. This is because Xen requires a full modification  
1245 to the kernel itself, not just a module. Loading a new kernel into the boot process  
1246 which may complicate patching and updating later in the system's maintenance cycle.  
1247 VMWare ESX, on the other hand, is entirely separate from most other installations.  
1248 As previously noted, ESX is actually a hypervisor and custom UNIX host OS com-  
1249 bined, so installation of ESX is likewise to installing any other OS from scratch. This  
1250 may be either desirable or adverse, depending on the system administrator's usage of  
1251 the systems and VMWare's ability to provide a secure and patched environment.

1252 While system administrators may be concerned with installation and maintenance,  
1253 VM users and Cloud developers are more concerned with daily usage. The first thing  
1254 to note about all of such virtualiation technologies is they are supported (to some  
1255 extent) by the libvirt API [142]. Libvirt is commonly used by many of today's IaaS  
1256 Cloud offerings, including Nimbus, Eucalyptus, OpenNebula and OpenStack. As  
1257 such, the choice of hypervisor for Cloud developer's is less of an issue, so long as  
1258 the hypervisor supports the features they desire. For individual command line usage  
1259 of each tool, it varies quite a bit more. Xen does provide their own set of tools for  
1260 controlling and monitoring guests, and seem to work relatively well but do incur a  
1261 slight learning curve. KVM also provides its own CLI interface, and while it is often  
1262 considered less cumbersome it provides less advanced features directly to users, such as  
1263 power management or quick memory adjustment (however this is subject to personal

opinion). One advantage of KVM is each guest actually runs as a separate process within the host OS, making it easy for a user to manage and control the VM inside the host if KVM misbehaves. VirtualBox, on the other hand, provides the best command line and graphical user interface. The CLI, is especially well featured when compared to Xen and KVM as it provides clear, decisive and well documented commands, something most HPC users and system administrators alike will appreciate. VMWare provides a significantly enhanced GUI as well as a Web-based ActiveX client interface that allows users to easily operate the VMWare host remotely. In summary, there is a wide variance of interfaces provided by each hypervisor, however we recommend Cloud developers to utilize the libvirt API whenever possible.

## 3.5 Experimental Design

In order to provide an unaltered and unbiased review of these virtualization technologies for Clouds, we need to outline a neutral testing environment. To make this possible, we have chosen to use FutureGrid as our virtualization and cloud test-bed.

### 3.5.1 The FutureGrid Project

FutureGrid (FG) [143] provides computing capabilities that enable researchers to tackle complex research challenges related to the use and security of Grids and Clouds. These include topics ranging from authentication, authorization, scheduling, virtualization, middleware design, interface design and cybersecurity, to the optimization of Grid-enabled and cloud-enabled computational schemes for researchers in astronomy, chemistry, biology, engineering, atmospheric science and epidemiology.

The test-bed includes a geographically distributed set of heterogeneous computing systems, a data management system that will hold both metadata and a growing

library of software images necessary for Cloud computing, and a dedicated network allowing isolated, secure experiments, as seen in Figure 3.2. The test-bed supports virtual machine-based environments, as well as operating systems on native hardware for experiments aimed at minimizing overhead and maximizing performance. The project partners are integrating existing open-source software packages to create an easy-to-use software environment that supports the instantiation, execution and recording of grid and cloud computing experiments.



**Figure 3.2** FutureGrid Participants and Resources

One of the goals of the project is to understand the behavior and utility of Cloud computing approaches. However, it is not clear at this time which of these toolkits will become the users' choice toolkit. FG provides the ability to compare these frameworks with each other while considering real scientific applications [38]. Hence, researchers are be able to measure the overhead of cloud technology by requesting linked experiments on both virtual and bare-metal systems, providing valuable information that help decide which infrastructure suits their needs and also helps users that want to

1301 transition from one environment to the other. These interests and research objectives  
1302 make the FutureGrid project the perfect match for this work. Furthermore, we expect  
1303 that the results gleaned from this chapter will have a direct impact on the FutureGrid  
1304 deployment itself.

### 1305 3.5.2 Experimental Environment

1306 Currently, one of FutureGrid's latest resources is the *India* system, a 256 CPU IBM  
1307 iDataPlex machine consisting of 1024 cores, 2048 GB of ram, and 335 TB of storage  
1308 within the Indiana University Data Center. In specific, each compute node of India  
1309 has two Intel Xeon 5570 quad core CPUs running at 2.93Ghz, 24GBs of Ram, and a  
1310 QDR InfiniBand connection. A total of four nodes were allocated directly from India  
1311 for these experiments. All were loaded with a fresh installation of Red Hat Enterprise  
1312 Linux server 5.5 x86\_64 with the 2.6.18-194.8.1.el5 kernel patched. Three of the four  
1313 nodes were installed with different hypervisors; Xen version 3.1, KVM (build 83), and  
1314 VirtualBox 3.2.10, and the forth node was left as-is to act as a control for bare-metal  
1315 native performance.

1316 Each guest virtual machine was also built using Red Hat EL server 5.5 running  
1317 an unmodified kernel using full virtualization techniques. All tests were conducted  
1318 giving the guest VM 8 cores and 16GB of ram to properly span a compute node.  
1319 Each benchmark was run a total of 20 times, with the results averaged to produce  
1320 consistent results, unless indicated otherwise.

### 1321 3.5.3 Benchmarking Setup

1322 As this chapter aims to objectively evaluate each virtualization technology from a  
1323 side-by-side comparison as well as from a performance standpoint, the selection of

1324 benchmarking applications is critical.

1325 The performance comparison of each virtual machine is based on two well known  
1326 industry standard performance benchmark suites; HPCC and SPEC. These two  
1327 benchmark environments are recognized for their standardized reproducible results in  
1328 the HPC communit, and the National Science Foundation (NSF), Department of En-  
1329 ergy (DOE), and DARPA are all sponsors of the HPCC benchmarks. The following  
1330 benchmarks provide a means to stress and compare processor, memory, inter-process  
1331 communication, network, and overall performance and throughput of a system. These  
1332 benchmarks were selected due to their importance to the HPC community sinse they  
1333 are often directly correlated with overall application performance [144].

1334 **HPCC Benchmarks**

1335 The HPCC Benchmarks [145, 146] are an industry standard for performing bench-  
1336 marks for HPC systems. The benchmarks are aimed at testing the system on multiple  
1337 levels to test their performance. It consists of 7 different tests:

- 1338     ● *HPL* - The Linpack TPP benchmark measures the floating point rate of exe-  
1339         cution for solving a linear system of equations. This benchmark is perhaps the  
1340         most important benchmark within HPC today, as it is the basis of evaluation  
1341         for the Top 500 list [8].
- 1342     ● *DGEMM* - Measures the floating point rate of execution of double precision real  
1343         matrix-matrix multiplication.
- 1344     ● *STREAM* - A simple synthetic benchmark program that measures sustainable  
1345         memory bandwidth (in GB/s) and the corresponding computation rate for sim-  
1346         ple vector kernel.

- 1347     ● *PTRANS* - Parallel matrix transpose exercises the communications where pairs  
1348       of processors communicate with each other simultaneously. It is a useful test of  
1349       the total communications capacity of the network.
  
- 1350     ● *RandomAccess* - Measures the rate of integer random updates of memory (GUPS).
  
- 1351     ● *FFT* - Measures the floating point rate of execution of double precision complex  
1352       one-dimensional Discrete Fourier Transform (DFT).
  
- 1353     ● *Communication bandwidth and latency* - A set of tests to measure latency and  
1354       bandwidth of a number of simultaneous communication patterns; based on b\_eff  
1355       (effective bandwidth benchmark).

1356     This benchmark suite uses each test to stress test the performance on multiple  
1357     aspects of the system. It also provides reproducible results which can be verified by  
1358     other vendors. This benchmark is used to create the Top 500 list [8] which is the list  
1359     of the current top supercomputers in the world. The results that are obtained from  
1360     these benchmarks provide an unbiased performance analysis of the hypervisors. Our  
1361     results provide insight on inter-node PingPong bandwidth, PingPong latency, and  
1362     FFT calculation performance.

### 1363     **SPEC Benchmarks**

1364     The Standard Performance Evaluation Corporation (SPEC) [147, 148] is the other  
1365     major standard for evaluation of benchmarking systems. SPEC has several different  
1366     testing components that can be utilized to benchmark a system. For our benchmark-  
1367     ing comparison we will use the SPEC OMP2001 because it appears to represent a  
1368     vast array of new and emerging parallel applications while simultaneously providing  
1369     a comparison to other SPEC benchmarks. SPEC OMP continues the SPEC tradi-

1370 tion of giving HPC users the most objective and representative benchmark suite for  
1371 measuring the performance of SMP (shared memory multi-processor) systems.

- 1372     ● The benchmarks are adapted from SPEC CPU2000 and contributions to its  
1373       search program.
- 1374     ● The focus is to deliver systems performance to real scientific and engineering  
1375       applications.
- 1376     ● The size and runtime reflect the needs of engineers and researchers to model  
1377       large complex tasks.
- 1378     ● Two levels of workload characterize the performance of medium and large sized  
1379       systems.
- 1380     ● Tools based on the SPEC CPU2000 toolset make these the easiest ever HPC  
1381       tests to run.
- 1382     ● These benchmarks place heavy demands on systems and memory.

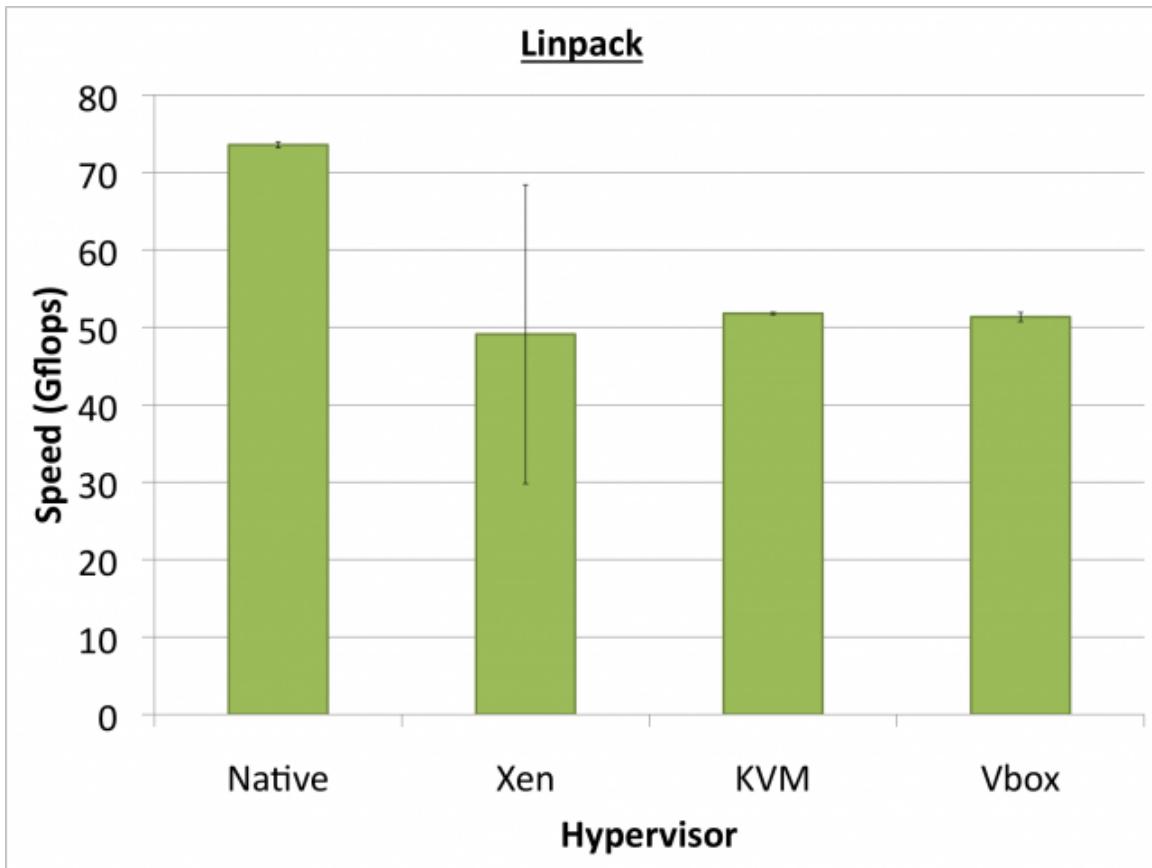
1383 **3.6 Performance Comparison**

1384 The goal of this chapter is to effectively compare and contrast the various virtual-  
1385 ization technologies, specifically for supporting HPC-based Clouds. The first set of  
1386 results represent the performance of HPCC benchmarks. Each benchmark was run  
1387 a total of 20 times, and the mean values taken with error bars represented using the  
1388 standard deviation over the 20 runs. The benchmarking suite was built using the Intel  
1389 11.1 compiler, uses the Intel MPI and MKL runtime libraries, all set with defaults  
1390 and no optimizations whatsoever.

1391 We open first with High Performance Linpack (HPL), the de-facto standard for  
1392 comparing resources. In Figure 3.3, we can see the comparison of Xen, KVM, and  
1393 Virtual Box compared to native bare-metal performance. First, we see that native  
1394 is capable of around 73.5 Gflops which, with no optimizations, achieves 75% of the  
1395 theoretical peak performance. Xen, KVM and VirtualBox perform at 49.1, 51.8 and  
1396 51.3 Gflops, respectively when averaged over 20 runs. However Xen, unlike KVM  
1397 and VirtualBox, has a high degree of variance between runs. This is an interesting  
1398 phenomenon for two reasons. First, this may impact performance metrics for other  
1399 HPC applications and cause errors and delays between even pleasingly-parallel appli-  
1400 cations and add to reducer function delays. Second, this wide variance breaks a key  
1401 component of Cloud computing providing a specific and predefined quality of service.  
1402 If performance can sway as widely as what occurred for Linpack, then this may have  
1403 a negative impact on users.

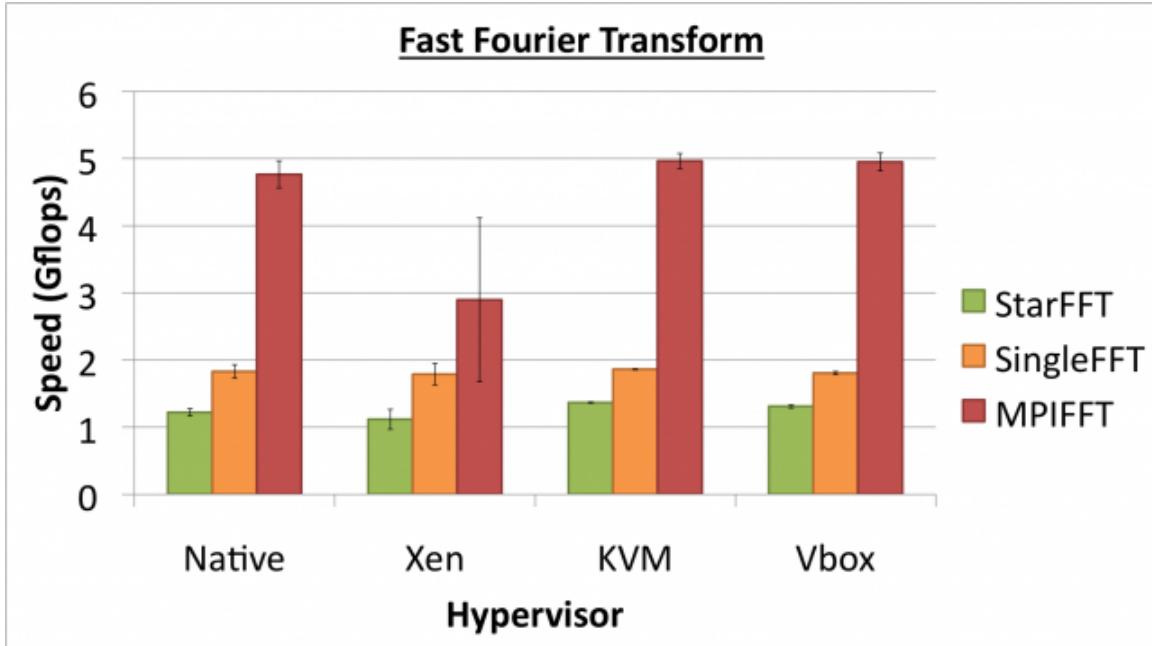
1404 Next, we turn to another key benchmark within the HPC community, Fast Fourier  
1405 Transforms (FFT). Unlike the synthetic Linpack benchmark, FFT is a specific, pur-  
1406 poseful benchmark which provides results which are often regarded as more relative  
1407 to a user's real-world application than HPL. From Figure 3.4, we can see rather dis-  
1408 tinct results from what was previously provided by HPL. Looking at Star and Single  
1409 FFT, its clear performance across all hypervisors is roughly equal to bare-metal per-  
1410 formance, a good indication that HPC applications may be well suited for use on  
1411 VMs. The results for MPI FFT also show similar results, with the exception of Xen,  
1412 which has a decreased performance and high variance as seen in the HPL benchmark.  
1413 Our current hypothesis is that there is an adverse affect of using Intel's MPI runtime  
1414 on Xen, however the investigation is still ongoing.

1415 Another useful benchmark illustrative of real-world performance between bare-  
1416 metal performance and various hypervisors are the ping-pong benchmarks. These



**Figure 3.3** Linpack performance

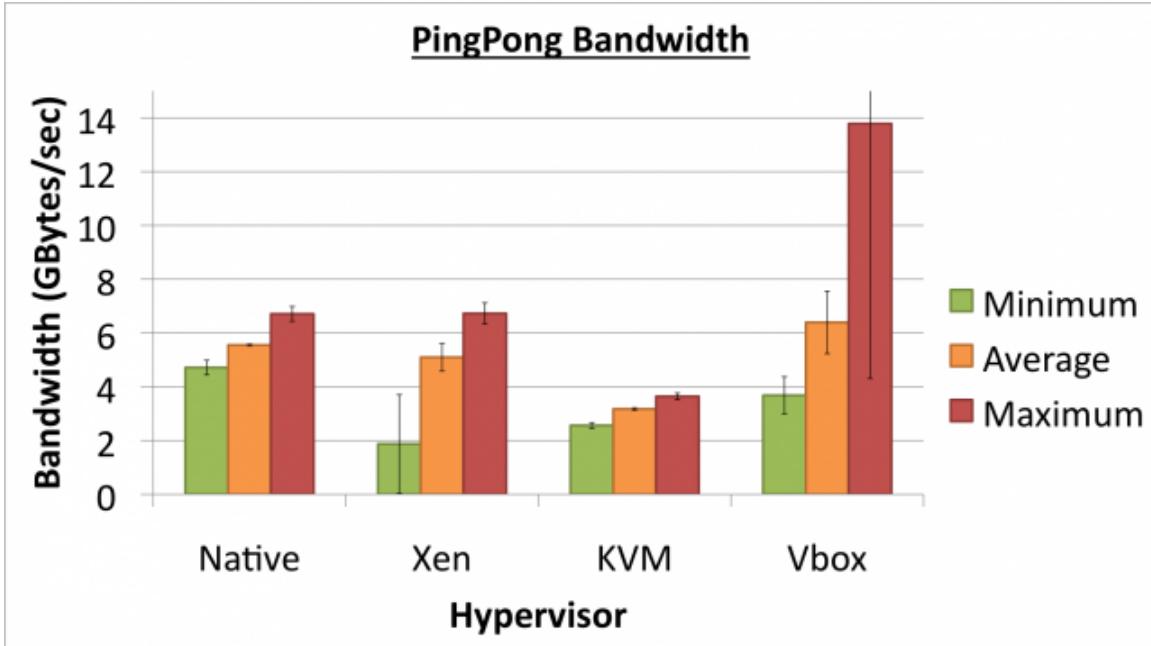
1417 benchmarks measure the bandwidth and latency of passing packets between multiple  
 1418 CPUs. With this experiment, all ping-pong latencies are kept within a given node,  
 1419 rather than over the network. This is done to provide further insight into the CPU and  
 1420 memory overhead within each hypervisor. From Figure 3.5 the intranode bandwidth  
 1421 performance is uncovered, with some interesting distinctions between each hypervi-  
 1422 sor. First, Xen performs, on average, close to native speeds, which is promising for  
 1423 the hypervisor. KVM, on the other hand, shows consistent overhead proportional  
 1424 to native performance across minimum, average, and maximum bandwidth. Virtu-  
 1425 alBox, on the other hand, performs well, in fact too well to the point that raises  
 1426 alarm. While the minimum and average bandwidths are within native performance,  
 1427 the maximum bandwidth reported by VirtualBox is significantly greater than native



**Figure 3.4** Fast Fourier Transform performance

1428 measurements, with a large variance. After careful examination, it appears this is  
 1429 due to how VirtualBox assigns its virtual CPUs. Instead of locking a virtual CPU to  
 1430 a real CPU, a switch may occur which could benefit on the off-chance the two CPU's  
 1431 in communication between a ping-pong test could in fact be the same physical CPU.  
 1432 The result would mean the ping-pong packet would remain in cache and result in a  
 1433 higher perceived bandwidth than normal. While this effect may be beneficial for this  
 1434 benchmark, it may only be an illusion towards the real performance gleaned from the  
 1435 VirtualBox hypervisor.

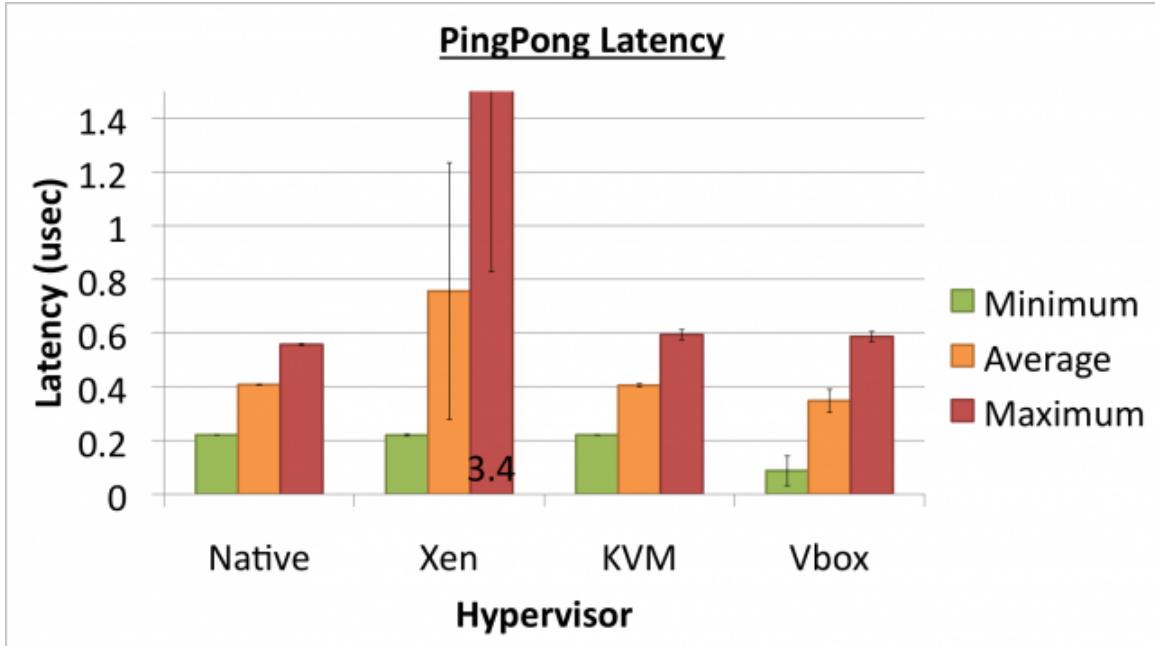
1436 The Bandwidth may in fact be important within the ping-pong benchmark, but  
 1437 the latency between each ping-pong is equally useful in understanding the perfor-  
 1438 mance impact of each virtualization technology. From Figure 3.6, we see KVM and  
 1439 VirtualBox have near-native performance; another promising result towards the util-  
 1440 ity of hypervisors within HPC systems. Xen, on the other hand, has extremely high  
 1441 latencies, especially at for maximum latencies, which in turn create a high variance



**Figure 3.5** Ping Pong bandwidth performance

<sup>1442</sup> within the average latency within the VM's performance.

<sup>1443</sup> While the HPCC benchmarks provide a comprehensive view for many HPC applications including Linpack and FFT using MPI, performance of intra-node SMP applications using OpenMP is also investigated. Figure 3.7 illustrates SPEC OpenMP performance across the VMs we concentrate on, as well as baseline native performance. <sup>1446</sup> First, we see that the combined performance over all 11 applications executed 20 times yields the native testbed with the best performance at a SPEC score of 34465. <sup>1447</sup> KVM performance comes close with a score of 34384, which is so similar to the native performance that most users will never notice the difference. Xen and VirtualBox both perform notably slower with scores of 31824 and 31695, respectively, however <sup>1451</sup> this is only an 8% performance drop compared to native speeds. Further results can <sup>1452</sup> be found on the SPEC website [149]. <sup>1453</sup>



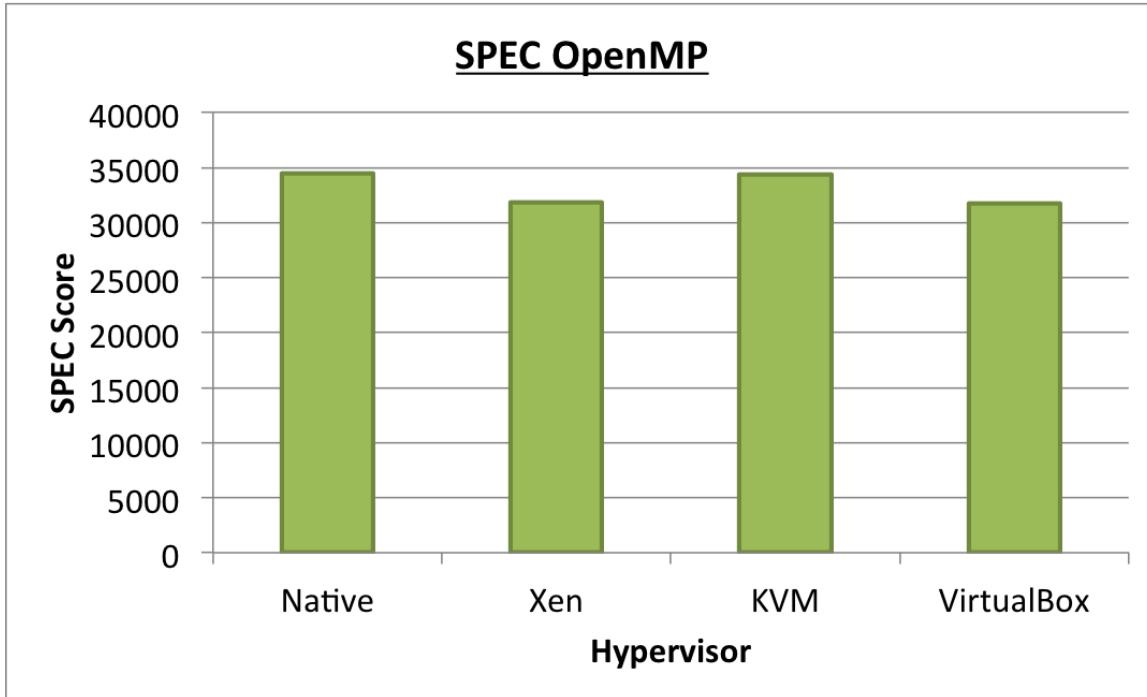
**Figure 3.6** Ping Pong latency performance (lower is better)

## 1454 3.7 Discussion

1455 The primary goal of this chapter is to evaluate the viability of virtualization within  
 1456 HPC. After our analysis, the answer seems to be a resounding "yes." However, we  
 1457 also hope to select the best virtualization technology for such an HPC environment.

1458 In order to do this, we combine the feature comparison along with the performance  
 1459 results, and evaluate the potential impact within the FutureGrid testbed.

1460 From a feature standpoint, most of today's virtualization technologies fit the bill  
 1461 for at least small scale deployment, including VMWare. In short, each support Linux  
 1462 x86\_64 platforms, use VT-X technology for full virtualization, and support live mi-  
 1463 gration. Due to VMWare's limited and costly licensing, it is immediately out of  
 1464 contention for most HPC deployments. From a CPU and memory standpoint, Xen  
 1465 seems to provide the best expandability, supporting up to 128 cpus and 4TB of ad-  
 1466 dressable RAM. So long as KVM's vCPU limit can be extended, it too shows promise



**Figure 3.7** Spec OpenMP performance

as a feature-full virtualization technology. One of Virtualbox's greatest limitations was the 16GB maximum memory allotment for individual guest VMs, which actually limited us from giving VMs more memory for our performance benchmarks. If this can be fixed and Oracle does not move the product into the proprietary market, VirtualBox may also stand a chance for deployment in HPC environments.

From the benchmark results previously described, the use of hypervisors within HPC-based Cloud deployments is mixed batch. Figure 3.8 summarizes the results based on a 1-3 rating, 1 being best and 3 being worst. While Linpack performance seems to take a significant performance impact across all hypervisors, the more practical FFT benchmarks seem to show little impact, a notably good sign for virtualization as a whole. The ping-pong bandwidth and latency benchmarks also seem to support this theory, with the exception of Xen, who's performance continually has wide fluctuations throughout the majority of the benchmarks. OpenMP performance

	Xen	KVM	VirtualBox
Linpack	3	1	2
FFT	3	1	2
Bandwidth	2	3	1
Latency	3	2	1
OpenMP	2	1	3
Total Rating	13	8	9

**Figure 3.8** Benchmark rating summary (lower is better)

1480 through the SPEC OMP benchmarking suite also shows promising results for the use  
 1481 of hypervisors in general, with KVM taking a clear lead by almost matching native  
 1482 speeds.

1483 While Xen is typically regarded as the most widely used hypervisor, especially  
 1484 within academic clouds and grids, it's performance has shown lack considerably when  
 1485 compared to either KVM or VirtualBox. In particular, Xen's wide and unexplained  
 1486 fluctuations in performance throughout the series of benchmarks suggests that Xen  
 1487 may not be the best choice for building a lasting quality of service infrastructure upon.

1488 From Figure 3.8, KVM rates the best across all performance benchmarks, making it  
 1489 the optimal choice for *general* deployment in an HPC environment. Furthermore,  
 1490 this work's illustration of the variance in performance among each benchmark and  
 1491 the applicability of each benchmark towards new applications may make possible the  
 1492 ability to preemptively classify applications for accurate prediction towards the ideal  
 1493 virtualized Cloud environment. We hope to further investigate this concept through  
 1494 the use of the FutureGrid experiment management framework at a later date.

1495 In summary, it is the authors' projection that KVM is the best overall choice for

1496 use within HPC Cloud environments. KVM's feature-rich experience and near-native  
1497 performance makes it a natural fit for deployment in an environment where usability  
1498 and performance are paramount. Within the FutureGrid project specifically, we hope  
1499 to deploy the KVM hypervisor across our Cloud platforms in the near future, as it  
1500 offers clear benefits over the current Xen deployment. Furthermore, we expect these  
1501 findings to be of great importance to other public and private Cloud deployments, as  
1502 system utilization, Quality of Service, operating cost, and computational efficiency  
1503 could all be improved through the careful evaluation of underlying virtualization  
1504 technologies.

1505 **Chapter 4**

1506 **Evaluating GPU Passthrough in  
1507 Xen for High Performance Cloud  
1508 Computing**

1509 **4.1 Abstract**

1510 With the advent of virtualization and Infrastructure-as-a-Service (IaaS), the broader  
1511 scientific computing community is considering the use of clouds for their technical  
1512 computing needs. This is due to the relative scalability, ease of use, advanced user  
1513 environment customization abilities clouds provide, as well as many novel comput-  
1514 ing paradigms available for data-intensive applications. However, there is concern  
1515 about a performance gap that exists between the performance of IaaS when com-  
1516 pared to typical high performance computing (HPC) resources, which could limit the  
1517 applicability of IaaS for many potential scientific users.

1518 Most recently, general-purpose graphics processing units (GPGPUs or GPUs) have  
1519 become commonplace within high performance computing. We look to bridge the

1520 gap between supercomputing and clouds by providing GPU-enabled virtual machines  
1521 (VMs) and investigating their feasibility for advanced scientific computation. Specif-  
1522 ically, the Xen hypervisor is utilized to leverage specialized hardware-assisted I/O  
1523 virtualization and PCI passthrough in order to provide advanced HPC-centric Nvidia  
1524 GPUs directly in guest VMs. This methodology is evaluated by measuring the per-  
1525 formance of two Nvidia Tesla GPUs within Xen VMs and comparing to bare-metal  
1526 hardware. Results show PCI passthrough of GPUs within virtual machines is a vi-  
1527 able use case for many scientific computing workflows, and could help support high  
1528 performance cloud infrastructure in the near future.

1529 **4.2 Introduction**

1530 Cloud computing [4] has established itself as a prominent paradigm within the realm  
1531 of Distributed Systems [150] in a very short period of time. Clouds are an internet-  
1532 based solution that provide computational and data models for utilizing resources,  
1533 which can be accessed directly by users on demand in a uniquely scalable way. Cloud  
1534 computing functions by providing a layer of abstraction on top of base hardware  
1535 to enable a new set of features that are otherwise intangible or intractable. These  
1536 benefits and features include Scalability, a guaranteed Quality of Service (QoS), cost  
1537 effectiveness, and direct user customization via a simplified user interface [60].

1538 While the origin of cloud computing is based in industry through solutions such  
1539 as Amazon's EC2 [151], Google's MapReduce [152], and Microsoft's Azure [153], the  
1540 paradigm has since become integrated in all areas of science and technology. Most  
1541 notably, there is an increasing effort within the High Performance Computing (HPC)  
1542 community to leverage the utility of clouds for advanced scientific computing to solve  
1543 a number of challenges still standing in the field. This can be clearly seen in large-

scale efforts such as the FutureGrid project [154], the Magellan project [42], and through various other Infrastructure-as-a-Service projects including OpenStack [155], Nimbus [70], and Eucalyptus [156].

Within HPC, there has also been a notable movement toward dedicated accelerator cards such as general purpose graphical processing units (GPGPUs, or GPUs) to enhance scientific computation problems by upwards of two orders of magnitude. This is accomplished through dedicated programming environments, compilers, and libraries such as CUDA [157] from Nvidia as well as the OpenCL effort [158]. When combining GPUs in an otherwise typical HPC environment or supercomputer, major gains in performance and computational ability have been reported in numerous fields [159, 160], ranging from Astrophysics to Bioinformatics. Furthermore, these gains in computational power have also reportedly come at an increased performance-per-watt [161], a metric that is increasingly important to the HPC community as we move closer to exascale computing [162] where power consumption is quickly becoming the primary constraint.

With the advent of both clouds and GPUs within the field of scientific computing, there is an immediate and ever-growing need to provide heterogeneous resources, most immediately GPUs, within a cloud environment in the same scalable, on-demand, and user-centric way that many cloud users are already accustomed to [163]. While this task alone is nontrivial, it is further complicated by the high demand for performance within HPC. As such, it is performance that is paramount to the success of deploying GPUs within cloud environments, and thus is the central focus of this work.

The rest of this chapter is organized as follows. First, in Section 2, we discuss the related research and the options currently available for providing GPUs within a virtualized cloud environment. In Section 3, we discuss the methodology for providing GPUs directly within virtual machines. In Section 4 we outline the evaluation of the

1570 given methodology using two different Nvidia Tesla GPUs and compare to the best-  
1571 case native application in Section 5. Then, we discuss the implications of these results  
1572 in Section 6 and consider the applicability of each method within a production cloud  
1573 system. Finally, we conclude with our findings and suggest directions for future work.

### 1574 4.3 Virtual GPU Directions

1575 Recently, GPU programming has been a primary focus for numerous scientific com-  
1576 puting applications. Significant progress has been accomplished in many different  
1577 workloads, both in science and engineering, based on parallel abilities of GPUs for  
1578 floating point operations and very high on-GPU memory bandwidth. This hardware,  
1579 coupled with CUDA and OpenCL programming frameworks, has led to an explosion  
1580 of new GPU-specific applications. In some cases, GPUs outperform even the fastest  
1581 multicore counterparts by an order of magnitude [164]. In addition, further research  
1582 could leverage the per-node performance of GPU accelerators with the high speed,  
1583 low latency interconnects commonly utilized in supercomputers and clusters to create  
1584 a hybrid GPU + MPI class of applications. The number of distributed GPU appli-  
1585 cations is increasing substantially in supercomputing, usually scaling many GPUs  
1586 simultaneously [165].

1587 Since the establishment of cloud computing in industry, research groups have  
1588 been evaluating its applicability to science [3]. Historically, HPC and Grids have  
1589 been on similar but distinct paths within distributed systems, and have concentrated  
1590 on performance, scalability, and solving complex, tightly coupled problems within  
1591 science. This has led to the development of supercomputers with many thousands  
1592 of cores, high speed, low latency interconnects, and sometimes also coprocessors and  
1593 FPGAs [166, 167]. Only recently have these systems been evaluated from a cloud

1594 perspective [42]. An overarching goal exists to provide HPC Infrastructure as its own  
1595 service (HPCaaS) [168], aiming to classify and limit the overhead of virtualization, and  
1596 reducing the bottlenecks classically found in CPU, memory, and I/O operations within  
1597 hypervisors [43, 169]. Furthermore, the transition from HPC to cloud computing  
1598 becomes more complicated when we consider adding GPUs to the equation.

1599 GPU availability within a cloud is a new concept that has sparked a large amount  
1600 of interest within the community. The first successfully deployment of GPUs within  
1601 a cloud environment was the Amazon EC2 GPU offering. A collaboration between  
1602 Nvidia and Citrix also exists to provide cloud-based gaming solutions to users using  
1603 the new Kepler GPU architecture [170]. However, this is currently not targeted  
1604 towards HPC applications.

1605 The task of providing a GPU accelerator for use in a virtualized cloud environment  
1606 is one that presents a myriad of challenges. This is due to the complicated nature of  
1607 virtualizing drivers, libraries, and the heterogeneous offerings of GPUs from multiple  
1608 vendors. Currently, two possible techniques exist to fill the gap in providing GPUs  
1609 in a cloud infrastructure: back-end I/O virtualization, which this chapter focuses on,  
1610 and Front-end remote API invocation.

#### 1611 4.3.1 Front-end Remote API invocation

1612 One method for using GPUs within a virtualized cloud environment is through front-  
1613 end library abstractions, the most common of which is remote API invocation. Also  
1614 known as API remoting or API interception, it represents a technique where API  
1615 calls are intercepted and forwarded to a remote host where the actual computation  
1616 occurs. The results are then returned to the front-end process that spawned the  
1617 invocation, potentially within a virtual machine. The goal of this method is to provide  
1618 an emulated device library where the actual computation is offloaded to another

1619 resource on a local network.

1620 Front-end remote APIs for GPUs have been implemented by a number of differ-  
1621 ent technologies for different uses. To solve the problem of graphics processing in  
1622 VMs, VMWare [171] has developed a device-emulation approach that emulates the  
1623 Direct3D and OpenGL calls to leverage the host OS graphics processing capabili-  
1624 ties to provide a 3D environment within a VM. API interception through the use of  
1625 wrapper binaries has also been implemented by technologies such as Chromium [172],  
1626 and Blink. However these graphics processing front-end solutions are not suitable for  
1627 general purpose scientific computing, as they do not expose interfaces that CUDA or  
1628 OpenCL can use.

1629 Currently, efforts are being made to provide a front-end remote API invocation  
1630 solutions for the CUDA programming architecture. vCUDA [51] was the first of such  
1631 technologies to enable transparent access of GPUs within VMs by API call inter-  
1632 ception and redirection of the CUDA API. vCUDA substitutes the CUDA runtime  
1633 library and supports a transmission mode using XMLRPC, as well as a sharing mode  
1634 that is built on VMRPC, a dedicated remote procedure call architecture for VMM  
1635 platforms. This share model can leads to better performance, especially as the volume  
1636 of data increases, although there may be limitations in VMM interoperability.

1637 Like vCUDA, gVirtuS uses API interception to enable transparent CUDA, OpenCL,  
1638 and OpenGL support for Xen, KVM, and VMWare virtual machines [173]. gVirtuS  
1639 uses a front-end/back-end model to provide a VMM-independent abstraction layer  
1640 to GPUs. Data transport from gVirtuS' front-end to the back-end is accomplished  
1641 through a combination of shared memory, sockets, or other hypervisor-specific APIs.  
1642 gVirtuS' primary disadvantage is in its decreased performance in host-to-device and  
1643 device-to-host data movement due to overhead of data copies to and from its shared  
1644 memory buffers. Recent work has also enabled the dynamic sharing of GPUs by

1645 leveraging the gVirtus back-end system with relatively good results [174], however  
1646 process-level GPU resource sharing is outside the scope of this manuscript.

1647 rCUDA [50], a recent popular remote CUDA framework, also provides remote API  
1648 invocation to enable VMs to access remote GPU hardware by using a sockets based  
1649 implementation for high-speed near-native performance of CUDA based applications.  
1650 rCUDA recently added support for using InfiniBand’s high speed, low latency network  
1651 to increase performance for CUDA applications with large data volume requirements.  
1652 rCUDA version 4.1 also supports the CUDA runtime API version 5.0, which supports  
1653 peer device memory access and unified addressing. One drawback of this method  
1654 is that rCUDA cannot implement the undocumented and hidden functions within  
1655 the runtime framework, and therefore does not support all CUDA C extensions.  
1656 While rCUDA provides some support tools, native execution of CUDA programs  
1657 is not possible and programs need to be recompiled or rewritten to use rCUDA.  
1658 Furthermore, like gVirtuS and many other solutions, performance between host-to-  
1659 device data movement is only as fast as the underlying interconnect, and in the best  
1660 case with native RDMA InfiniBand, is roughly half as fast as native PCI Express  
1661 usage when using the standard QDR InfiniBand.

#### 1662 4.3.2 Back-end PCI passthrough

1663 Another approach to using a GPU in a virtualized environment is to provide a VM  
1664 with direct access to the GPU itself, instead of relying on a remote API. This chapter  
1665 focuses on such an approach. Devices on a host’s PCI-express bus are virtualized  
1666 using directed I/O virtualization technologies recently implemented by chip manu-  
1667 facturers, and then direct access is relinquished upon request to a guest VM. This  
1668 can be accomplished using the VT-d and IOMMU instruction sets from Intel and  
1669 AMD, respectively. This mechanism, typically called PCI passthrough, uses a mem-

1670 memory management unit (MMU) to handle direct memory access (DMA) coordination  
1671 and interrupt remapping directly to the guest VM, thus bypassing the host entirely.  
1672 With host involvement being nearly non-existent, near-native performance of the PCI  
1673 device within the guest VM can be achieved, which is an important characteristic for  
1674 using a GPU within a cloud infrastructure.

1675 PCI passthrough itself has recently become a standard technique for many other  
1676 I/O systems such as storage or network controllers. However, GPUs (even from  
1677 the same vendor) have additional legacy VGA compatibility issues and non-standard  
1678 low-level interface DMA interactions that make direct PCI passthrough nontrivial.  
1679 VMWare has started use of a vDGA system for hardware GPU utilization, however it  
1680 remains in tech preview and only documentation for Windows VMs is present [171]. In  
1681 our experimentation, we have found that the Xen hypervisor provides a good platform  
1682 for performing PCI passthrough of GPU devices to VMs due to its open nature,  
1683 extensive support, and high degree of reconfigurability. Work with Xen in [175] gives  
1684 hints at good performance for PCI passthrough in Xen, however further evaluation  
1685 with independent benchmarks is needed when looking at scientific computing with  
1686 GPUs.

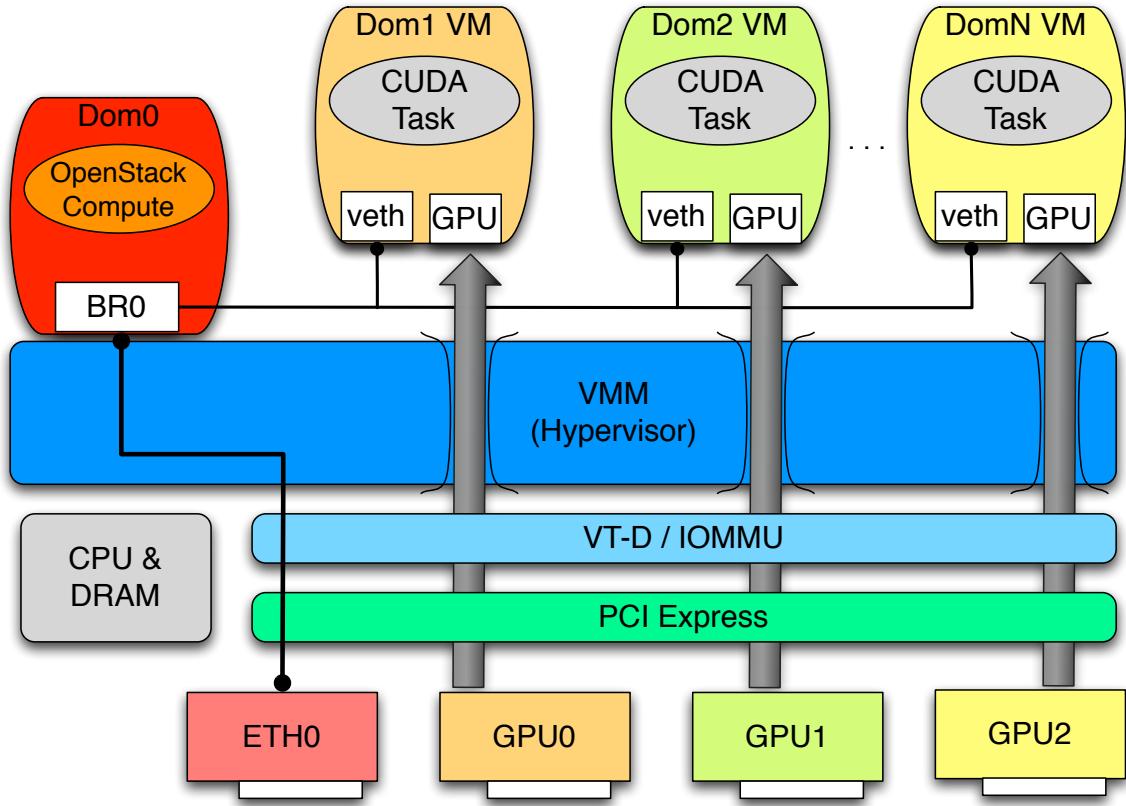
1687 Today's GPUs can provide a variety of frameworks for application programmers to  
1688 use. Two common solutions are CUDA and OpenCL. CUDA, or the Compute Unified  
1689 Device Architecture, is a framework for creating and running parallel applications on  
1690 Nvidia GPUs. OpenCL provides a more generic and open framework for parallel  
1691 computation on CPUs and GPUs, and is available for a number of Nvidia, AMD, and  
1692 Intel GPUs and CPUs. While OpenCL provides a more robust and portable solution,  
1693 many HPC applications utilize the CUDA framework. As such, we focus only on  
1694 Nvidia based CUDA-capable GPUs as they offer the best support for a wide array of  
1695 programs, although this work is not strictly limited to Nvidia GPUs.

**1696 4.4 Implementation**

1697 In this chapter we use a specific host environment to enable PCI passthrough. First,  
1698 we start with the Xen 4.2.2 hypervisor on Centos 6.4 and a custom 3.4.50-8 Linux  
1699 kernel with Dom0 Xen support. Within the Xen hypervisor, GPU devices are seized  
1700 upon boot and assigned to the xen-pciback kernel module. This process blocks the  
1701 host devices from loading the Nvidia or nouveau drivers, keeping the GPUs uninitialized  
1702 and therefore able to be assigned to DomU VMs.

1703 Xen, like other hypervisors, provides a standard method of passing through PCI  
1704 devices to guest VMs upon creation. When assigning a GPU to a new VM, Xen loads  
1705 a specific VGA BIOS to properly initialize the device enabling DMA and interrupts  
1706 to be assigned to the guest VM. Xen also relinquishes control of the GPU via the  
1707 xen-pciback module. From there, the Linux Nvidia drivers are loaded and the device  
1708 is able to be used as expected within the guest. Upon VM termination, the xen-  
1709 pciback module re-seizes the GPU and the devices can be re-assigned to new VMs in  
1710 the future.

1711 This mechanism of PCI passthrough for GPUs can be implemented using multiple  
1712 devices per host, as illustrated in Figure 6.1. Here, we see how the device's connection  
1713 to the VM totally bypasses the Dom0 host as well as the Xen VMM, and is managed  
1714 by VT-d or IOMMU to access the PCI-Express bus which the GPUs utilize. This is in  
1715 contrast to other common virtual device uses, where hardware is emulated by the host  
1716 and shared across all guests. This is the common usage for Ethernet controllers and  
1717 input devices to enable users to interact with VMs as they would with native hosts,  
1718 unlike the bridged model shown in the figure. The potential downside of this method  
1719 is there can be a 1:1 or 1:N mapping of VMs to GPUs only. A M:1 mapping where  
1720 multiple VMs use a GPU is not possible. However, almost all scientific applications



**Figure 4.1** GPU PCI passthrough within the Xen Hypervisor

1721 environments using GPUs generally do not share GPUs between processes or other  
 1722 nodes, as doing so would cause unpredictable and serious performance degradation.  
 1723 As such, this GPU isolation within a VM can be considered an advantage in many  
 1724 contexts.

#### 1725 4.4.1 Feature Comparison

1726 Using the GPU PCI passthrough technique described previously has a number of  
 1727 advantages compared to front-end API implementations. First, it allows for an op-  
 1728 erating environment that more closely relates to native bare-metal usage of GPUs.  
 1729 Essentially, a VM provides a nearly identical infrastructure to clusters and supercom-  
 1730 puters with integrated GPUs. This lowers the learning curve for many researchers,

and even enables researchers to potentially use other tools within a VM that might not be supported within a supercomputer or cluster. Unlike with remote API implementations, users don't need to recompile or modify their code, as the GPUs are essentially local to the data. Further comparing to remote API implementations, using PCI passthrough within a VM allows users to leverage any GPU framework available, OpenCL or CUDA, and any higher level programming frameworks such as within Matlab or Python.

Through the use of advanced scheduling techniques within cloud infrastructure, we can also take advantage of PCI passthrough implementation for efficiency purposes. Users could request VMs with GPUs which get scheduled for creation on machines that provide such resources, but can also request normal VMs as well. The scheduler can correctly map VM requirement requests to heterogeneous hardware. This enables large scale resources to support a wide array of scientific computing applications without the added cost of putting GPUs in all compute nodes.

## 4.5 Experimental Setup

In this chapter back-end GPU PCI passthrough to virtual machines using the Xen hypervisor is detailed, however proper evaluation of the performance of such method needs to be properly considered. As such, we ran an array of benchmarks that evaluate the performance of this method compared to the same hardware running native bare-metal GPU code without any virtualization. We focus our tests on single-node performance to best understand low level overhead.

To evaluate the effectiveness of GPU-enabled VMs within Xen, two different machines were used to represent two generations of Nvidia GPUs. The first system at Indiana University consists of dual-socket Intel Xeon X5660 6-core CPUs at 2.8Ghz

1755 with 192GB DDR3 RAM, 8TB RAID5 array, and two Nvidia Tesla "Fermi" C2075  
1756 GPUs. The system at USC/ISI uses a dual-socket Intel Xeon E5-2670 8-core CPUs  
1757 at 2.6Ghz with 48GB DDR3 RAM, 3 600GB SAS disk drives, but with the latest  
1758 Nvidia Tesla "Kepler" K20m GPU supplied by Nvidia. These machines represent  
1759 the present Fermi series GPUs along with the recently release Kepler series GPUs,  
1760 providing a well-rounded experimental environment. Native systems were installed  
1761 with standard Centos 6.4 with a stock 2.6.32-279 Linux kernel. Xen host systems  
1762 were still loaded with Centos 6.4 but with a custom 3.4.53-8 Linux kernel and Xen  
1763 4.2.22. Guest VMs were also Centos 6.4 with N-1 processors and 24GB of memory  
1764 and 1 GPU passed through in HVM full virtualization mode. Using both IU and  
1765 USC/ISI machine configurations in native and VM modes represent the 4 test cases  
1766 for our work.

1767 In order to evaluate the performance, the SHOC Benchmark suite [176] was used  
1768 to extensively evaluate performance across each test platform. The SHOC bench-  
1769 marks were chosen because they provide a higher level of evaluation regarding GPU  
1770 performance than the sample applications provided in the Nvidia SDK, and can also  
1771 evaluate OpenCL performance in similar detail. The benchmarks were compiled us-  
1772 ing the NVCC compiler in CUDA5 driver and library, along with OpenMPI 1.4.2 and  
1773 GCC 4.4. Each benchmark was run a total of 20 times, with the results given as an  
1774 average of all runs.

## 1775 4.6 Results

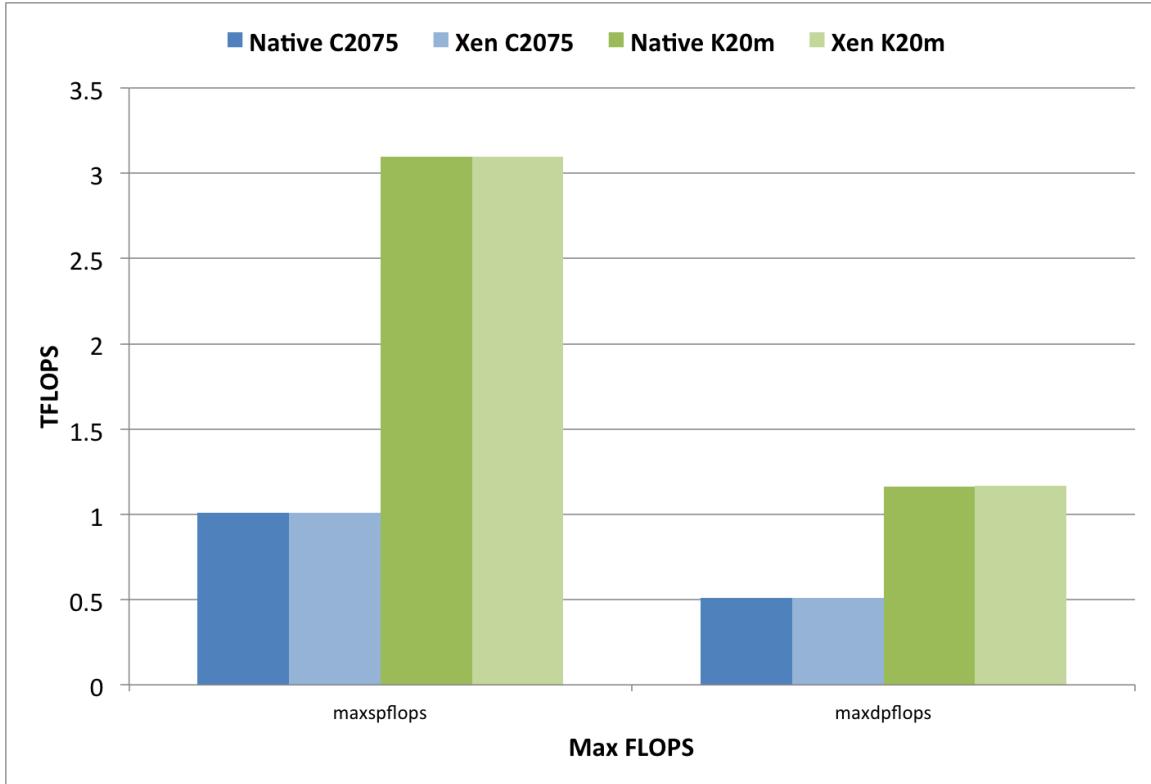
1776 Results of all benchmarks are compressed into three subsections: floating point oper-  
1777 ations, device bandwidth and pci bus performance. Each represents a different level  
1778 of evaluation for GPU-enabled VMs compared to bare-metal native GPU usage.

**4.6.1 Floating Point Performance**

Flops, or floating point operations per second, are a common measure of computational performance, especially within scientific computing. The Top 500 list [8] has been the relative gold standard for evaluating supercomputing performance for more than two decades. With the advent of GPUs in some of the fastest supercomputers today, including GPUs identical to those used in our experimentation, understanding the performance relative to this metric is imperative.

Figure 4.2 shows the raw peak FLOPs available using each GPU in both native and virtualized modes. First, we observe the advantage of using the Kepler series GPUs over Fermi, with peak single precision speeds tripling in each case. Even for double precision FLOPs, there is roughly a doubling of performance with the new GPUs. However its most interesting that there is less than a 1% overhead when using GPU VMs compared to native case for pure floating point operations. The K20m-enabled VM was able to provide over 3 Teraflops, a notable computational feat for any single node.

Figures 4.3 and 4.4 show Fast Fourier Transform (FFT) and the Matrix Multiplication implementations across both test platforms. For all benchmarks that do not take into account the PCI-Express (pcie) bus transfer time, we again see near-native performance using the Kepler GPUs and Fermi GPUs when compared to bare metal cases. Interestingly, overhead within Xen VMs is consistently less than 1%, confirming the synthetic MaxFlops benchmark above. However, we do see some performance impact when calculating the total FLOPs with the pcie bus in the equation. This performance decrease ranges significantly for the C2075-series GPU, roughly about a 15% impact for FFT and a 5% impact for Matrix Multiplication. This overhead in pcie runs is not as pronounced for the Kepler K20m test environment, with near-native performance in all cases (less than 1%).

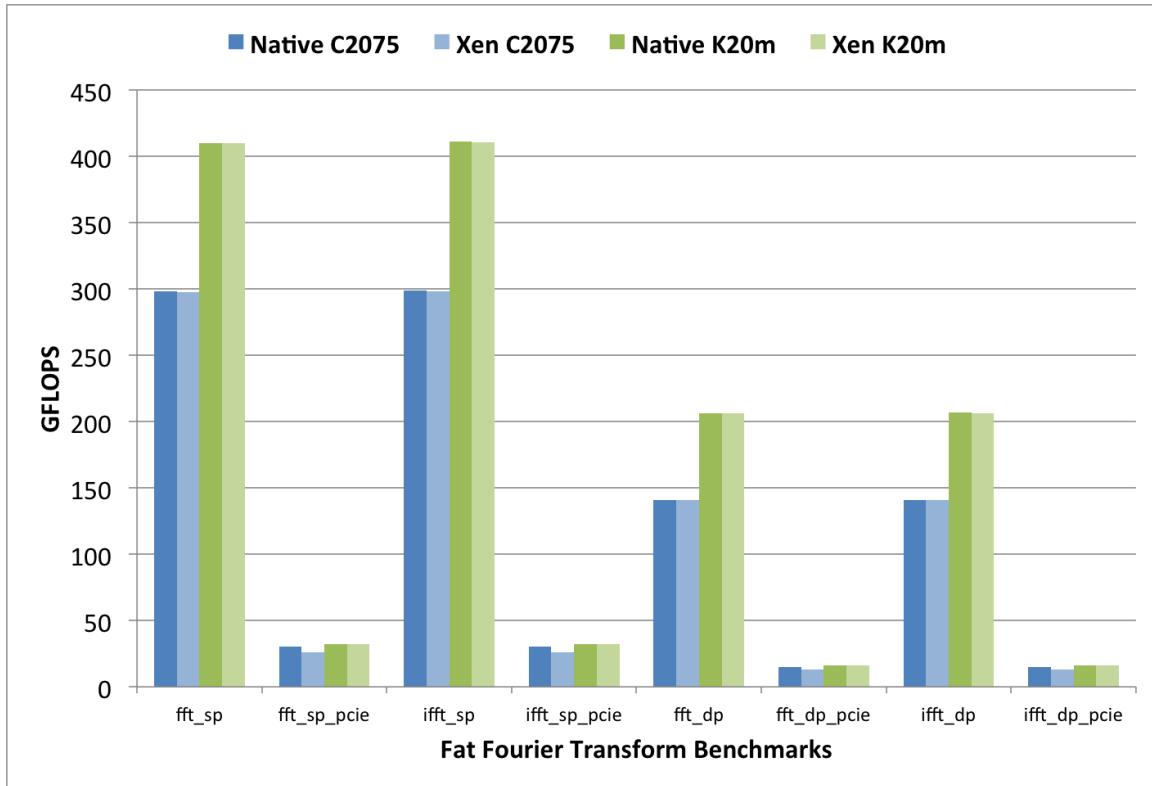


**Figure 4.2** GPU Floating Point Operations per Second

1805 Other FLOP-based benchmarks are used to emulate higher level applications.  
 1806 Stencil represents a 9-point stencil operation applied to a 2D data set, and the S3D  
 1807 benchmark is a computationally-intensive kernel from the S3D turbulent combustion  
 1808 simulation program [177]. In Figure 4.5, we see that both the Fermi C2075 and Kepler  
 1809 K20m GPUs performing well compared to the native base case, showing the overhead  
 1810 of virtualization is low. The C2075-enabled VMs experience slightly more overhead  
 1811 when compared to native performance again for pcie runs, but overhead is at most  
 1812 7% for the S3D benchmark.

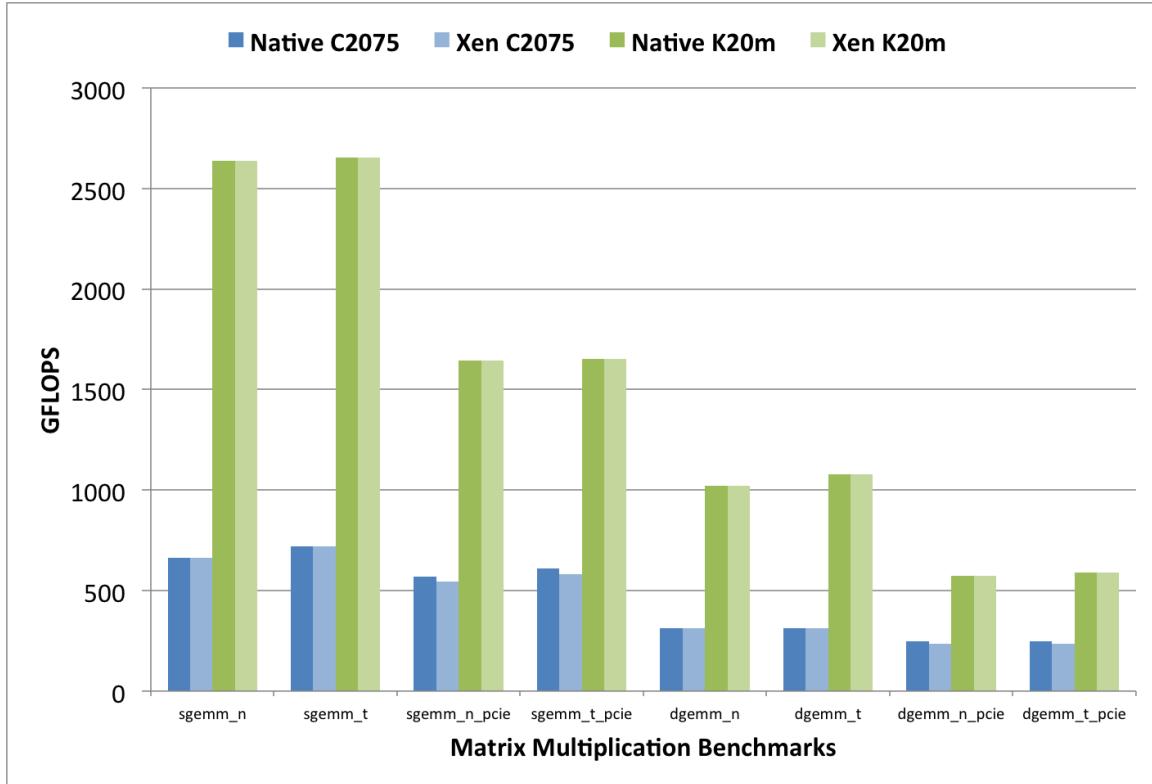
### 1813 4.6.2 Device Speed

1814 While floating point operations allow for the proposed solution to relate to many tradi-  
 1815 tional HPC applications, they are just one facet of GPU performance within scientific



**Figure 4.3** GPU Fast Fourier Transform

1816 computing. Device speed, measured in both raw bandwidth and additional bench-  
 1817 marks, provides a different perspective towards evaluating GPU PCI passthrough in  
 1818 Xen. Figure 4.6 illustrates device level memory access of various GPU device mem-  
 1819 ory structures. With both Nvidia GPUs, virtualization has little to no impact on the  
 1820 performance of inter-device memory bandwidth. As expected the Kepler K20m out-  
 1821 performed the C2075 VMs and there was a higher variance between runs with both  
 1822 native and VM cases. Molecular Dynamics and Reduction benchmarks in Figure 4.7  
 1823 perform again at near-native performance without the pcie bus taken into account.  
 1824 However the overhead observed increases to 10-15% when the PCI-Express bus is  
 1825 considered when looking at the Fermi C2075 VMs.

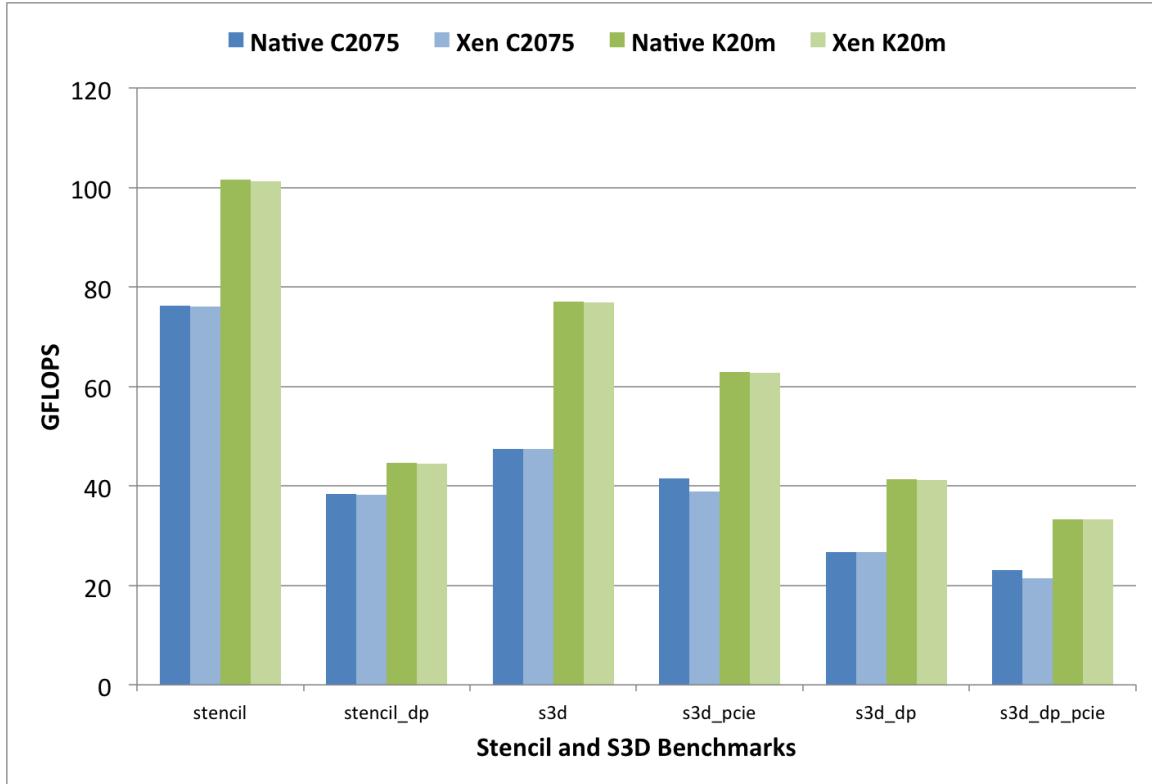


**Figure 4.4** GPU Matrix Multiplication

### 1826 4.6.3 PCI Express Bus

1827 Upon evaluating PCI passthrough of GPUs, it is apparent that the PCI express bus  
 1828 is subject to the greatest potential for overhead, as was observed in the Fermi C2075  
 1829 benchmarks. The VT-d and IOMMU chip instruction sets interface directly with the  
 1830 PCI bus to provide operational and security related mechanisms for each PCI device,  
 1831 thereby ensuring proper function in a multi-guest environment but potentially intro-  
 1832 ducing some overhead. As such, it is imperative to investigate any and all overhead  
 1833 at the PCI Express bus.

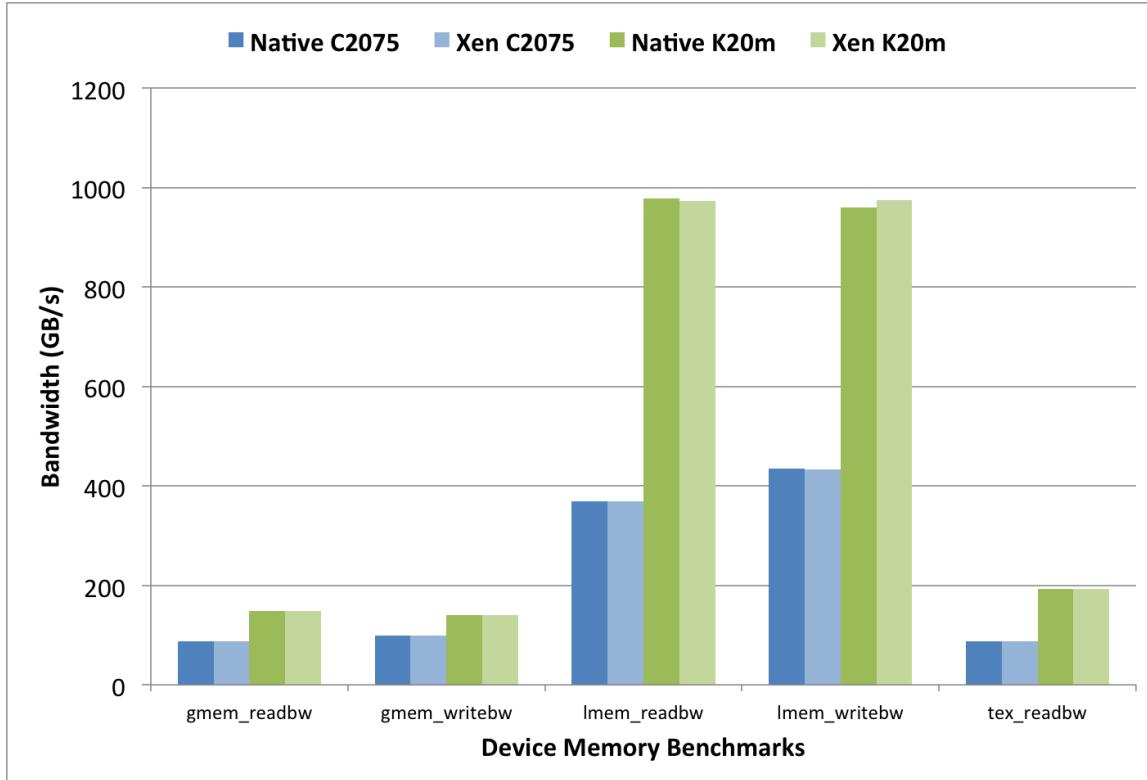
1834 Figure 4.8 looks at maximum PCI bus speeds for each experimental implemen-  
 1835 tation. First, we see a consistent overhead in the Fermi C2075 VMs, with a 14.6%  
 1836 performance impact for download (to-device) and a 26.7% impact in readback (from-



**Figure 4.5** GPU Stencil and S3D

device). However, the Kepler K20 VMs do not experience the same overhead in the PCI-Express bus, and instead perform at near-native performance. These results indicate that the overhead is minimal in the actual VT-d mechanisms, and instead due to other factors.

Upon further examination, we have identified the performance degradation within the Fermi-based VM experimental setup is due to the testing environment's NUMA node configuration with the Westmere-based CPUs. With the Fermi nodes, the GPUs are connected through the PCI-Express bus from CPU Socket #1, yet the experimental GPU VM was run using CPU Socket #0. This meant that the VM's PCI-Express communication involved more host involvement with Socket #1, leading to a notable decrease in read and write speeds across the PCI-Express Bus. Such architectural limitations seem to be relieved in the next generation with a Sandy-Bridge CPU ar-

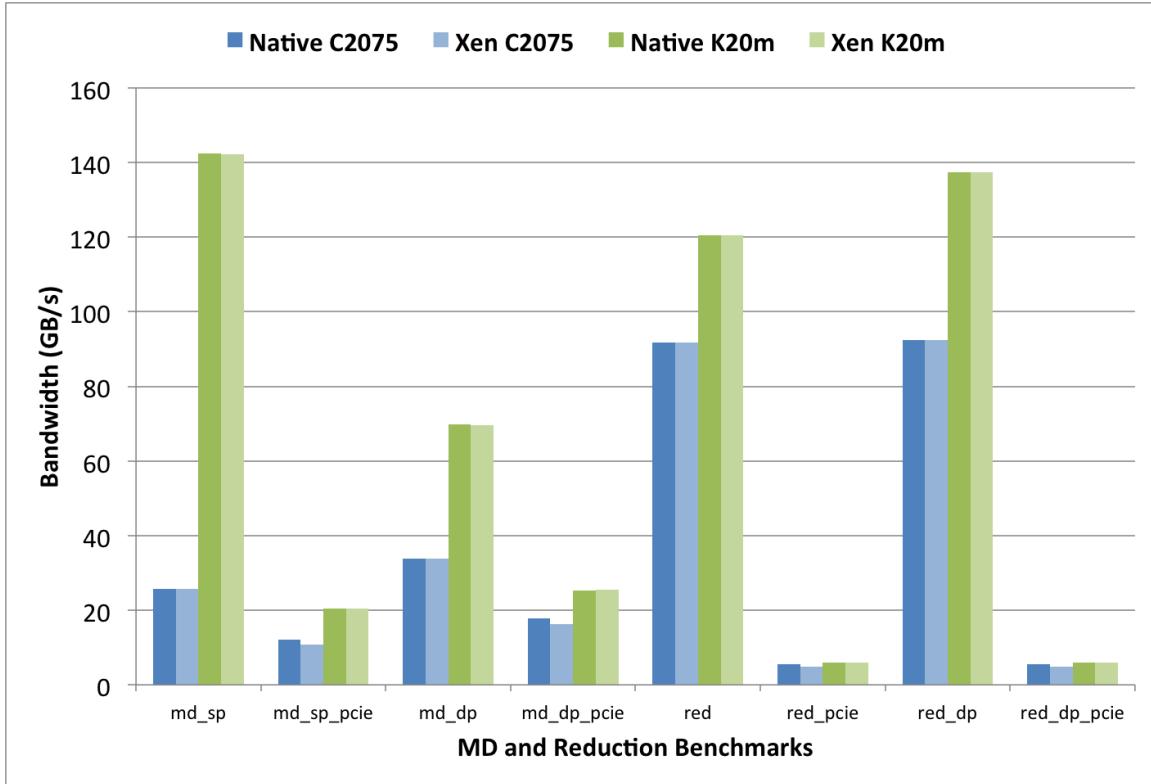


**Figure 4.6** GPU Device Memory Bandwidth

1849 chitecture and the Kepler K20m experimental setup. In summary, GPU PCI-Express  
 1850 performance in a VM is sensitive to the host CPU's NUMA architecture and care is  
 1851 needed to mitigate the impact, either by leveraging new architectures or by proper us-  
 1852 age of Xen's VM core assignment features. Furthermore, the overhead in this system  
 1853 diminishes significantly when using the new Kepler GPUs by Nvidia.

## 1854 4.7 Discussion

1855 This chapter evaluates the use of general purpose GPUs within cloud computing in-  
 1856 frastructure, primarily targeted towards advanced scientific computing. The method  
 1857 of PCI passthrough of GPUs directly to a guest virtual machine running on a tuned  
 1858 Xen hypervisor shows initial promise for an ubiquitous solution in cloud infrastruc-

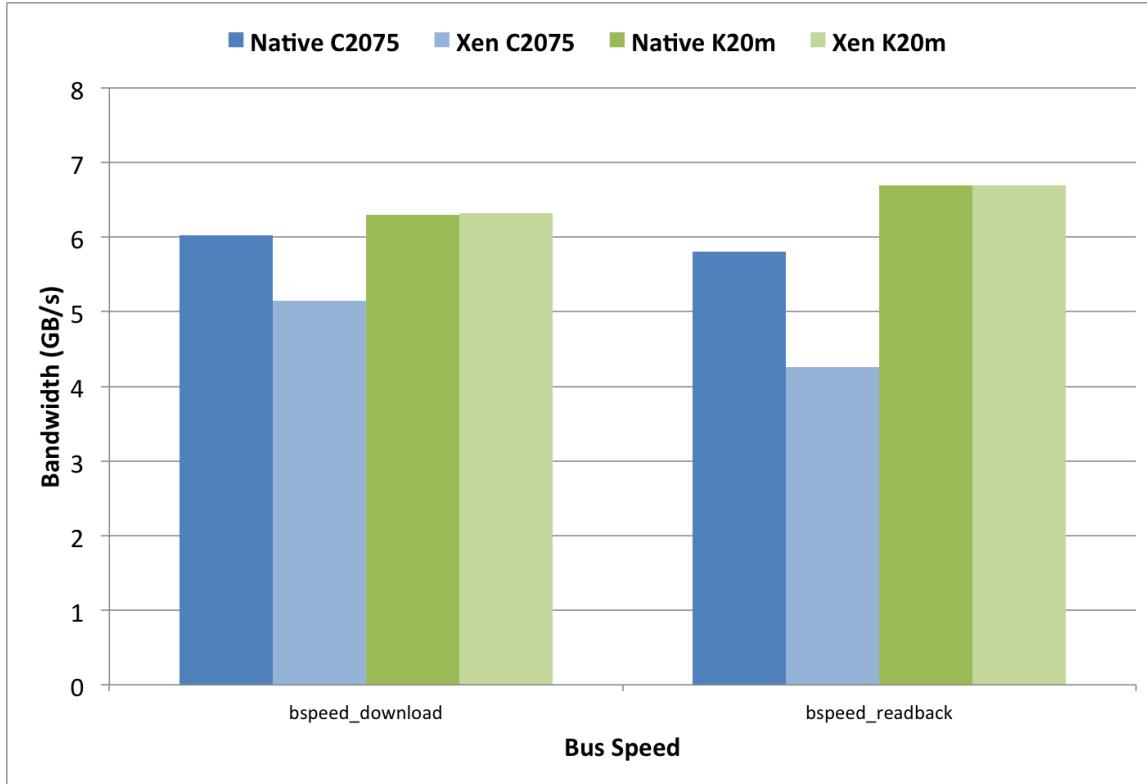


**Figure 4.7** GPU Molecular Dynamics and Reduction

ture. In evaluating the results in the previous section, a number of points become clear.

First, we can see that there is a small overhead in using PCI passthrough of GPUs within VMs, compared to native bare-metal usage, which represents the best possible use case. As with all abstraction layers, some overhead is usually inevitable as a necessary trade-off to added feature sets and improved usability. The same is true for GPUs within Xen VMs. The Kepler K20m GPU-enabled VMs operated at near-native performance for all runs, with a 1.2% reduction at worst in performance. The Fermi based C2075 VMs experience more overhead due to the NUMA configuration that impacted PCI-Express performance. However, when the overhead of the PCI-Express bus is not considered, the C2075 VMs perform at near-native speeds.

GPU PCI passthrough also has the potential to perform better than other front-



**Figure 4.8** GPU PCI Express Bus Speed

end API solutions that are applicable within VMs, such as gVirtus or rCUDA. This is because such solutions are designed to communicate via a network interconnect such as 10Gb Ethernet or QDR InfiniBand [178], which introduces an inherent bottleneck. Even with the theoretically optimal configuration of rCUDA using QDR InfiniBand, the maximum theoretical bus speed is 40Gbs, which is comparably less than the measured 54.4Gps real-world performance measured between host-to-device transfers with GPU-enabled Kepler VMs.

Overall, it is our hypothesis that the overhead in using GPUs within Xen VMs as described in this chapter will largely go unnoticed by most mid-level scientific computing applications. This is especially true when using the latest Sandy-Bridge CPUs with the Kepler series GPUs. We expect many mid-tier scientific computing groups to benefit the most from the ability to use GPUs in a scientific cloud infras-

1883 tructure. Already this has been confirmed in [179], where similar a methodology has  
1884 been leveraged specifically for Bioinformatics applications in the cloud.

1885 **4.8 Chapter Summary and Future Work**

1886 The ability to use GPUs within virtual machines represents a leap forward for sup-  
1887 porting advanced scientific computing within cloud infrastructure. The method of  
1888 direct PCI passthrough of Nvidia GPUs using the Xen hypervisor offers a clean, re-  
1889 producible solution that can be implemented within many Infrastructure-as-a-Service  
1890 (IaaS) deployments. Performance measurements indicate that the overhead of pro-  
1891 viding a GPU within Xen is minimal compared to the best-case native use, however  
1892 NUMA inconsistencies can impact performance. The New Kepler-based GPUs oper-  
1893 ate with a much lower overhead, making those GPUs an ideal choice when designing  
1894 a new GPU IaaS system.

1895 Next steps for this work could involve providing GPU-based PCI passthrough  
1896 within the OpenStack nova IaaS framework. This will enable research laboratories  
1897 and institutions to create new private or national-scale cloud infrastructure that have  
1898 the ability to support new scientific computing challenges. Other hypervisors could  
1899 also leverage GPU PCI passthrough techniques and warrant in-depth evaluation in  
1900 the future. Furthermore, we hope to integrate this work with advanced interconnects  
1901 and other heterogeneous hardware and provide a parallel high performance cloud  
1902 infrastructure to enable mid-tier scientific computing.

1903 **Chapter 5**

1904 **GPU-Passthrough Performance: A  
1905 Comparison of KVM, Xen,  
1906 VMWare ESXi, and LXC for  
1907 CUDA and OpenCL Applications**

1908 **5.1 Introduction**

1909 As scientific workloads continue to demand increasing performance at greater power  
1910 efficiency, high performance architectures have been driven towards heterogeneity and  
1911 specialization. Intel's Xeon Phi, and GPUs from both NVIDIA and AMD represent  
1912 some of the most common accelerators, with each capable of delivering improved  
1913 performance and power efficiency over commodity multi-core CPUs.

1914 Infrastructure-as-a-Service (IaaS) clouds have the potential to democratize access  
1915 to the latest, fastest, and most powerful computational accelerators. This is true  
1916 of both public and private clouds. Yet today's clouds are typically homogeneous

1917 without access to even the most commonly used accelerators. Historically, enabling  
1918 virtual machine access to GPUs and other PCIe devices has proven complex and  
1919 error-prone, with only a small subset of GPUs being certified for use within a few  
1920 commercial hypervisors. This is especially true for NVIDIA GPUs, likely the most  
1921 popular for scientific computing, but whose drivers have always been closed source.

1922 Given the complexity surrounding the choice of GPUs, host systems, and hypervi-  
1923 sors, it is perhaps no surprise that Amazon is the only major cloud provider offering  
1924 customers access to GPU-enabled instances. All of this is starting to change, however,  
1925 as open source and other freely available hypervisors now provide sufficiently robust  
1926 PCI passthrough functionality to enable GPU and other accelerator access whether  
1927 in the public or private cloud.

1928 Today, it is possible to access GPUs at high performance within all of the major  
1929 hypervisors, merging many of the advantages of cloud computing (e.g. custom images,  
1930 software defined networking, etc.) with the accessibility of on-demand accelerator  
1931 hardware. Yet, no study to date has systematically compared the performance of PCI  
1932 passthrough across all major cloud hypervisors. Instead, alternative solutions have  
1933 been proposed that attempt to virtualize the GPU [180] , but sacrifice performance.

1934 In this chapter, we characterize the performance of both NVIDIA Fermi and  
1935 Kepler GPUs operating in PCI passthrough mode in VMWare VSphere, Linux KVM,  
1936 Xen, and Linux Containers (LXC). Through a series of microbenchmarks as well as  
1937 scientific and Big Data applications, we make two contributions:

1938 1. We demonstrate that PCI passthrough at high performance is possible for GPUs  
1939 across 4 major hypervisors.

1940 2. We describe the lessons learned through our performance analysis, as well as  
1941 the relative advantages and disadvantages of each hypervisor for GPU support.

**1942 5.2 Related Work & Background**

1943 GPU virtualization and GPU-passthrough are used within a variety of contexts, from  
1944 high performance computing to virtual desktop infrastructure. Accessing one or more  
1945 GPUs within a virtual machine is typically accomplished by one of two strategies: 1)  
1946 via API remoting with device emulation; or 2) using PCI passthrough.

**1947 5.2.1 GPU API Remoting**

1948 rCUDA, vCUDA, GVIM, and gVirtuS are well-known API remoting solutionsFun-  
1949 damentally, these approaches operate similarly by splitting the driver into a front-  
1950 end/back-end model, where calls into the interposed CUDA library (front-end) are  
1951 sent via shared memory or a network interface to the back-end service that executes  
1952 the CUDA call on behalf of the virtual machine. Notably, this technique is not limited  
1953 to CUDA, but can be used to decouple OpenCL, OpenGL, and other APIs from their  
1954 local GPU or accelerator.

1955 The performance of API-remoting depends largely on the application and the  
1956 remoting solution's implementation. Bandwidth and latency-sensitive benchmarks and  
1957 applications will tend to expose performance bottlenecks more than compute-intensive  
1958 applications. Moreover, solutions that rely on high speed networks, such as Infini-  
1959 band, will compete with application-level networking for bandwidth.

**1960 5.2.2 PCI Passthrough**

1961 Input/Output Memory Management Units, or IOMMUs, play a fundamental roll in  
1962 the PCI-passthrough virtualization mechanism. Like traditional MMUs that provide  
1963 a virtual memory address space to CPUs [181], an IOMMU serves the fundamental  
1964 purpose of connecting a direct memory access (DMA) capable I/O bus to main mem-

1965 ory. The IOMMU unit, typically within the chipset, maps device virtual addresses to  
1966 physical memory addresses. This process also has the added improvement of guaran-  
1967 teeing device isolation by blocking rogue DMA and interrupt requests [182], with a  
1968 slight overhead, especially in early implementations [183].

1969 Currently two major IOMMU implementations exist, VT-d and AMD-Vi by In-  
1970 tel and AMD, respectively. Both specifications provide DMA remapping to enable  
1971 PCI-passthrough as well as other features such as interrupt remapping, hypervisor  
1972 snooping, and security control mechanisms to ensure proper and efficient hardware  
1973 utilization. PCI passthrough has been studied within the context of networking [184],  
1974 storage [185], and other PCI-attached devices; however, GPUs have historically lagged  
1975 behind other devices in their support for virtual machine passthrough.

### 1976 **5.2.3 GPU Passthrough, a Special Case of PCI Passthrough**

1977 While generic PCI passthrough can be used with IOMMU technologies to pass through  
1978 many PCI-Express devices, GPUs represent a special case of PCI devices, and a spe-  
1979 cial case of PCI passthrough. In traditional usage, GPUs usually serve as VGA  
1980 devices primarily to render screen output, and while the GPUs used in this study do  
1981 not render screen out, the function still exists in legacy. In GPU-passthrough, another  
1982 VGA device (such as onboard graphics built into the motherboard, or a baseboard  
1983 management controller) is necessary to serve as the primary display for the host, as  
1984 well as providing emulated VGA devices for each guest VM. Most GPUs also have a  
1985 video BIOS that requires full initialization and reset functions, which is often difficult  
1986 due to the proprietary nature of the cards and their drivers.

1987 Nevertheless, for applications that require native or near-native GPU performance  
1988 across the full spectrum of applications with immediate access to the latest GPU  
1989 drivers and compilers, GPU passthrough solutions are preferable to API remoting.

1990 Today, Citrix Xenserver, open source Xen [186], and VMWare ESXi [187], and most  
1991 recently KVM all support GPU passthrough. To our knowledge, no one has system-  
1992 atically characterized the performance of GPU passthrough across a range of hyper-  
1993 visors, across such a breadth of benchmarks, and across multiple GPU generations as  
1994 we do.

1995 **5.3 Experimental Methodology**

1996 **5.3.1 Host and Hypervisor Configuration**

1997 We used two hardware systems, named Bespin and Delta, to evaluate four hypervisors.  
1998 The Bespin system at USC/ISI represents Intel’s Sandy Bridge microarchitecture with  
1999 a Kepler class K20 GPU. The Delta system, provided by the FutureGrid project [188],  
2000 represents the Westmere microarchitecture with a Fermi class C2075 GPU. Table 5.1  
2001 provides the major hardware characteristics of both systems. Note that in addition,  
2002 both systems include 10 gigabit Ethernet, gigabit Ethernet, and either FDR or QDR  
2003 Infiniband. Our experiments do not emphasize networking, and we use the gigabit  
2004 ethernet network for management only.

2005 A major design goal of these experiments was to reduce or eliminate NUMA effects  
2006 (non-uniform memory access) on the PCI passthrough results in order to facilitate  
2007 fair comparisons across hypervisors and to reduce experimental noise. To this end,  
2008 we configured our virtual machines and containers to execute only on the NUMA  
2009 node containing the GPU under test. We acknowledge that the NUMA effects on  
2010 virtualization may be interesting in their own right, but they are not the subject of  
2011 this set of experiments.

2012 We use Bespin and Delta to evaluate three hypervisors and one container-based  
2013 approach to GPU passthrough. The hypervisors and container system, VMWare

**Table 5.1** Host hardware configurations

	Bespin	Delta
CPU (cores)	2x E5-2670 (16)	2x X5660 (12)
Clock Speed	2.6 GHz	2.6 GHz
RAM	48 GB	192 GB
NUMA Nodes	2	2
GPU	1x K20m	2x C2075

2014 ESXi, Xen, KVM, and LXC, are summarized in Table 5.2. Note that each virtual-  
 2015 ization solution imposes its own unique requirements on the base operating system.  
 2016 Hence, Xen’s Linux kernel refers to the Domain-0 kernel, whereas KVM’s Linux kernel  
 2017 represents the actual running kernel hosting the KVM hypervisor. Linux Containers  
 2018 share a single kernel between the host and guests, and VMWare ESXi does not rely  
 2019 on a Linux kernel at all.

2020 Similarly, hypervisor requirements prevented us from standardizing on a single  
 2021 host operating system. For Xen and KVM, we relied on the Arch Linux 2013.10.01  
 2022 distribution because it provides easy access to the mainline Linux kernel. For our  
 2023 LXC tests, we use CentOS 6.4 because its shared kernel was identical to the base  
 2024 CentOS 6.4 kernel used in our testing. VMWare has a proprietary software stack.  
 2025 All of this makes comparison challenging, but as we describe in Section 5.3.2, we are  
 2026 running a common virtual machine across all experiments.

**Table 5.2** Host Hypervisor/Container Configuration

Hypervisor	Linux Kernel	Linux Version
KVM	3.12	Arch 2013.10.01
Xen 4.3.0-7	3.12	Arch 2013.10.01
VMWare ESXi 5.5.0	N/A	N/A
LXC	2.6.32-358.23.2	CentOS 6.4

### 2027 5.3.2 Guest Configuration

2028 We treat each hypervisor as its own system, and compare virtual machine guests  
 2029 to a base CentOS 6.4 system. The base system and the guests are all composed of  
 2030 CentOS 6.4 installation with a 2.6.32-358.23.2 stock kernel and CUDA version 5.5.  
 2031 Each guest is allocated 20 GB of RAM and a full CPU socket (either 6 or 8 CPU  
 2032 cores). Bespin experiments received 8 cores and Delta experiments received 6 cores.  
 2033 VMs were restricted to a single NUMA node. On the Bespin system, the K20m GPU  
 2034 was attached to NUMA node 0. On the Delta system, the C2075 GPU was attached  
 2035 to NUMA node 1. Hence VMs ran on NUMA node 0 for the Bespin experiments,  
 2036 and node 1 for the Delta experiments.

### 2037 5.3.3 Microbenchmarks

2038 Our experiments are composed of a mix of microbenchmarks and application-level  
 2039 benchmarks, as well as a combination of CUDA and OpenCL benchmarks. The  
 2040 SHOC benchmark suite provides a series of microbenchmarks in both OpenCL and  
 2041 CUDA [189]. For this analysis, we focus on the OpenCL benchmarks in order to  
 2042 exercise multiple programming models. Benchmarks range from low-level peak Flops

2043 and bandwidth measurements, to kernels and mini-applications.

2044 **5.3.4 Application Benchmarks**

2045 For our application benchmarks, we have chosen the LAMMPS molecular dynam-  
2046 ics simulator [190], the GPU-LIBSVM [191], and the LULESH shock hydrodynamics  
2047 simulator [192]. These represent a range of computational characteristics, from com-  
2048 putational physics to big data analytics, and are representative of GPU-accelerated  
2049 applications in common use.

2050 **LAMMPS** The Large-scale Atomic/Molecular Parallel Simulator (LAMMPS), is a  
2051 parallel molecular dynamics simulator [190, 193] used for production MD simulation  
2052 on both CPUs and GPUs [194]. LAMMPS has two packages for GPU support, the  
2053 USER-CUDA and GPU packages. With the USER-CUDA package, each GPU is used  
2054 by a single CPU, whereas the GPU package allows multiple CPUs to take advantage  
2055 of a single GPU. There are performance trade-offs with both approaches, but we chose  
2056 to use the GPU package in order to stress the virtual machine by exercising multiple  
2057 CPUs. Consistent with the existing GPU benchmarking approaches, our results are  
2058 based on the Rhodopsin protein.

2059 **GPU-LIBSVM** LIBSVM is a popular implementation [195] of the machine learn-  
2060 ing classification algorithm support vector machine (SVM). GPU-accelerated LIB-  
2061 SVM [191] enhances LIBSVM by providing GPU-implementations of the kernel ma-  
2062 trix computation portion of the SVM algorithm for radial basis kernels. For bench-  
2063 marking purposes we use the NIPS 2003 feature extraction gisette data set. This data  
2064 set has a high dimensional feature space and large number of training instances, and  
2065 these qualities are known to be computational intensive to generate SVM models.

2066 The GPU-accelerated SVM implementation shows dramatic improvement over the  
2067 CPU-only implementation.

2068 **LULESH** Hydrodynamics is widely used to model continuum properties and inter-  
2069 actions in materials when there is an applied force [196]. Hydrodynamics applications  
2070 consume approximately one third of the runtime of data center resource throughout  
2071 the U.S. DoD (Department of Defense). The Livermore Unstructured Lagrange Ex-  
2072 plicit Shock Hydro (LULESH) was developed by Lawrence Livermore National Lab  
2073 as one of five challenge problems in the DARPA UHPC program. LULESH is widely  
2074 used as a proxy application in the U.S. DOE (Department of Energy) co-design effort  
2075 for exascale applications [192].

## 2076 5.4 Performance Results

2077 We characterize GPGPU performance within virtual machines across two hardware  
2078 systems, 4 hypervisors, and 3 application sets. We begin with the SHOC benchmark  
2079 suite before describing the GPU-LIBSVM, LAMMPS, and LULESH results. All  
2080 benchmarks are run 20 times and averaged. Results are scaled with respect to a base  
2081 CentOS 6.4 system for both systems. That is, we compare virtualized Bespin per-  
2082 formance to non-virtualized Bespin performance, and virtualized Delta performance  
2083 to non-virtualized Delta performance. Values less than 1 indicate that the base sys-  
2084 tem outperformed the virtual machine, while values greater than 1 indicate that the  
2085 virtual machine outperformed the base system. In cases where we present geometric  
2086 means across multiple benchmarks, the means are taken over these scaled values, and  
2087 the semantics are the same: less than 1 indicates overhead in the hypervisor, greater  
2088 than 1 indicates a performance increase over the base system.

**Table 5.3** SHOC overheads for Bespin (K20) expressed as geometric means of scaled values within a level, while maximum overheads are expressed as a percentage.

		Bespin (K20)							
		KVM		Xen		LXC		VMWare	
		Mean	Max	Mean	Max	Mean	Max	Mean	Max
L0	0.999	1.57	0.997	3.34	1.00	1.77	1.00	1.90	
L1	0.998	1.23	0.998	1.39	1.00	1.47	1.00	0.975	
L2	0.998	0.48	0.995	0.846	0.999	1.90	0.999	1.20	
Bespin PCIe-only									
		KVM		Xen		LXC		VMWare	
		Mean	Max	Mean	Max	Mean	Max	Mean	Max
L0	0.997	0.317	0.999	0.143	0.995	0.981	0.995	1.01	
L1	0.998	0.683	0.997	1.39	1.00	0.928	1.01	0.975	
L2	0.998	0.478	0.996	0.846	1.00	0.247	1.01	0.133	

**Table 5.4** SHOC overheads for Delta (c2075) expressed as geometric means of scaled values within a level, while maximum overheads are expressed as a percentage.

		Delta (C2075)							
		KVM		Xen		LXC		VMWare	
		Mean	Max	Mean	Max	Mean	Max	Mean	Max
L0	1.01	0.031	0.969	12.7	1.00	0.073	1.00	4.95	
L1	1.00	1.45	0.959	24.0	1.00	0.663	0.933	36.6	
L2	1.00	0.101	0.982	4.60	1.00	0.016	0.962	7.01	
Delta PCIe-only									
		KVM		Xen		LXC		VMWare	
		Mean	Max	Mean	Max	Mean	Max	Mean	Max
L0	1.04	0.029	0.889	12.7	1.00	0.01	0.995	4.37	
L1	1.00	1.45	0.914	20.5	0.999	0.380	0.864	36.6	
L2	1.00	0.075	0.918	4.60	1.00	N/A	0.869	7.01	

**2089 5.4.1 SHOC Benchmark Performance**

2090 SHOC splits its benchmarks into 3 Levels, named Levels 0 through 2. Level 0 repre-  
2091 sents device-level characteristics, peak Flop/s, bandwidth, etc. Level 1 characterizes  
2092 computational kernels: FFT and matrix multiplication, among others. Finally, Level  
2093 2 includes “mini-applications,” in this case an implementation of the S3D, a compu-  
2094 tational chemistry application.

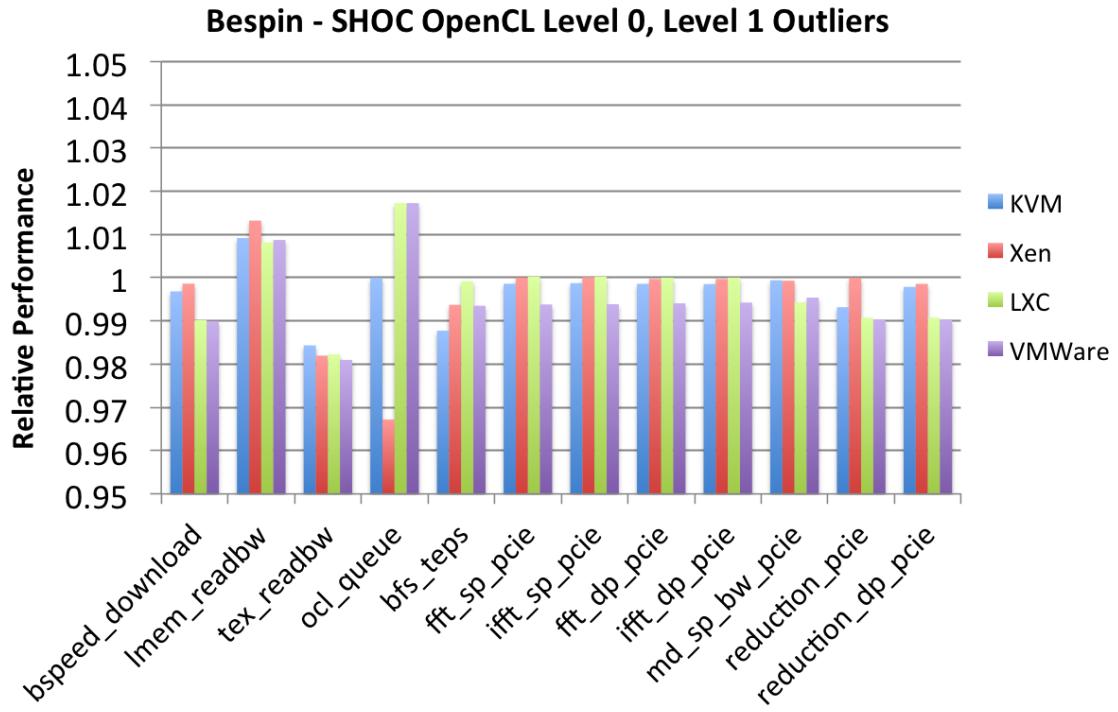
2095 Because the SHOC OpenCL benchmarks report more than 70 individual mi-  
2096 crobenchmarks, space does not allow us to show each benchmark individually. In-  
2097 stead, we start with a broad overview of SHOC’s performance across all benchmarks,  
2098 hypervisors, and systems. We then discuss in more detail those benchmarks that  
2099 either outperformed or underperformed the Bespin (K20) system by 0.50% or more.  
2100 We call these benchmarks outliers. As we will show, those outlier benchmarks iden-  
2101 tified on the Bespin system, also tend to exhibit comparable characteristics on the  
2102 Delta system as well, but the overhead is typically higher.

2103 In Tables 5.3 and 5.4, we provide geometric means for each SHOC level across each  
2104 hypervisor and system. We also include the maximum overhead for each hypervisor  
2105 at each level to facilitate comparison across hypervisors and systems. Finally, we  
2106 provide a comparable breakdown of only the PCIe SHOC benchmarks, based on the  
2107 intuition that PCIe-specific benchmarks will likely result in higher overhead.

2108 At a high level, we immediately notice that in the cases of KVM and LXC, both  
2109 perform very near native across both the Bespin and Delta platforms. On average,  
2110 these systems are almost indistinguishable from their non-virtualized base systems.  
2111 So much so, that experimental noise occasionally boosts performance slightly above  
2112 their base systems.

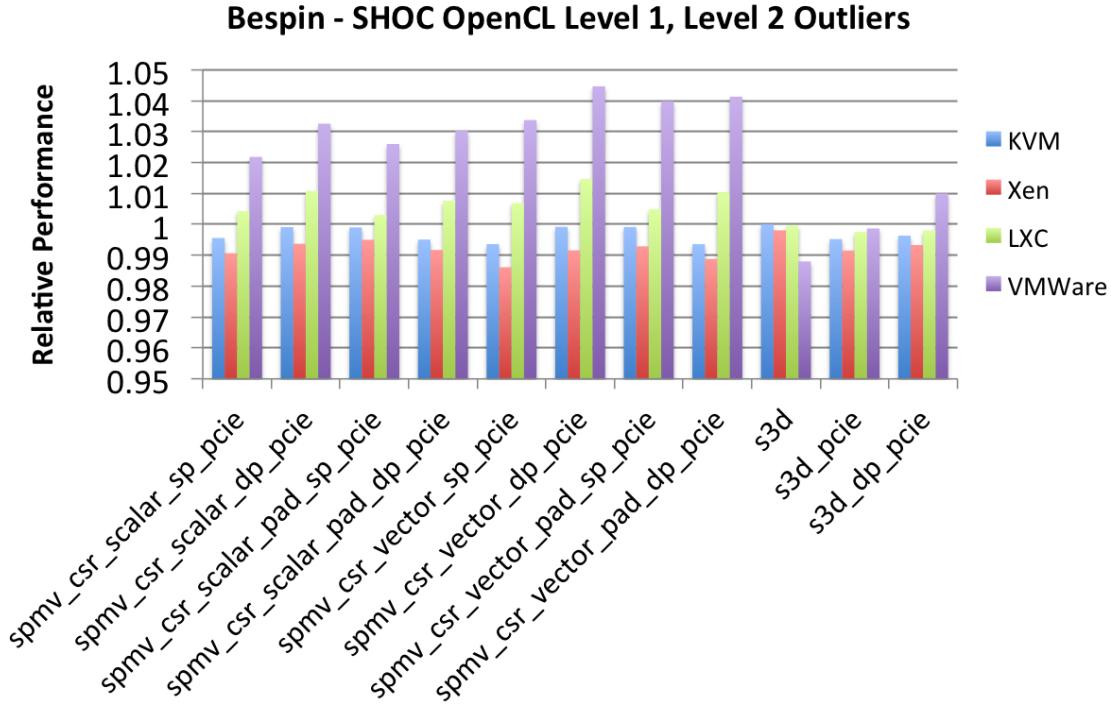
2113 This is in sharp contrast to the Xen and VMWare hypervisors, which perform  
2114 well on the Bespin system, but poorly on the Delta system in some cases. This is

particularly evident when looking at the maximum overheads for Xen and VMWare across both systems. In this case, we see that on the Bespin system, Xen's maximum overhead of 3.34% is dwarfed by Delta's maximum Xen overhead of 24.0%. VMWare exhibits similar characteristics, resulting in a maximum overhead of 1.9% in the case of the Bespin system, and a surprising 36.6% in the case of the Delta system. We provide a more in-depth discussion of these overheads below.



**Figure 5.1** SHOC Levels 0 and 1 relative performance on Bespin system. Results show benchmarks over or under-performing by 0.5% or greater. Higher is better.

Of the Level 0 benchmarks, only four exhibited notable overhead in the case of the Bespin system: bspeed\_download, lmem\_readbw, tex\_readbw, and ocl\_queue. These are shown in Figure 5.1. These benchmarks represent device-level characteristics, including host-to-device bandwidth, onboard memory reading, and OpenCL kernel queue delay. Of the four, only bspeed\_download incurs a statistically significant

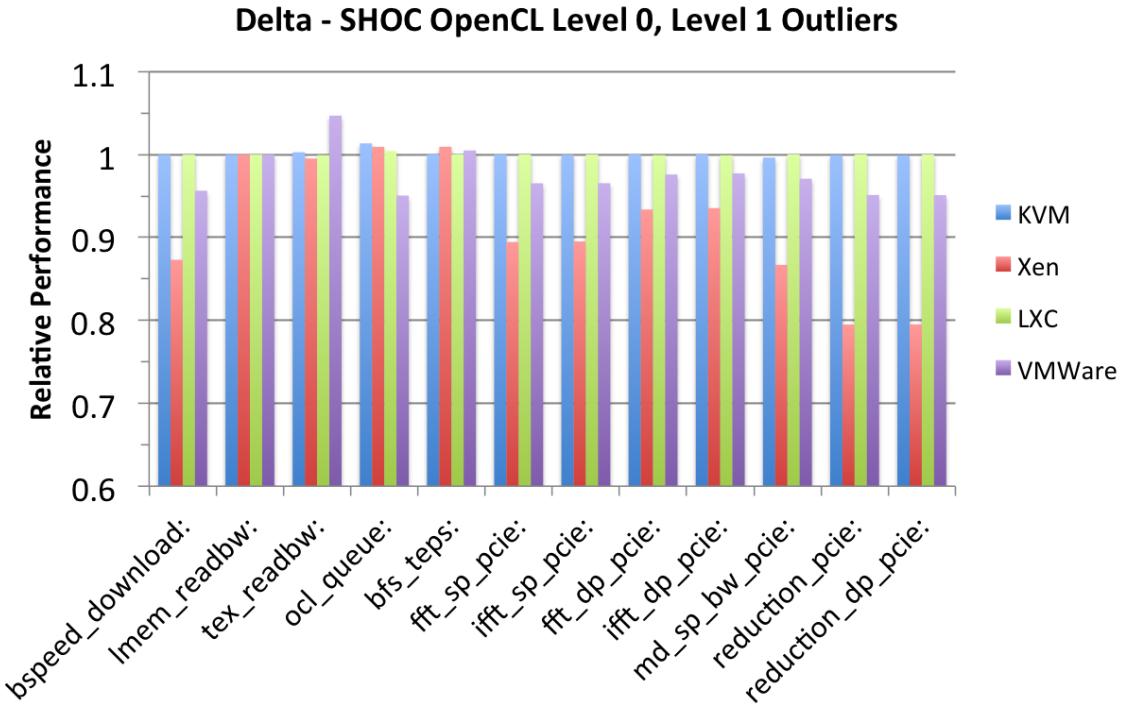


**Figure 5.2** SHOC Levels 1 and 2 relative performance on Bespin system. Results show benchmarks over or under-performing by 0.5% or greater. Higher is better.

overhead. The remainder perform within one standard deviation of the base, despite an overhead of greater than 0.5%.

bspeed\_download is representative of the most likely source of virtualization overhead, data movement across the PCI-Express bus. The PCIe lies at the boundary of the virtual machine and the physical GPU, and requires interrupt remapping, IOMMU interaction, etc. in order to enable GPU passthrough into the virtual machine. Despite this, in Figure 5.1 we see a maximum of 1% overhead for bspeed\_download in VMWare, and less than 0.5% overhead for both KVM and Xen.

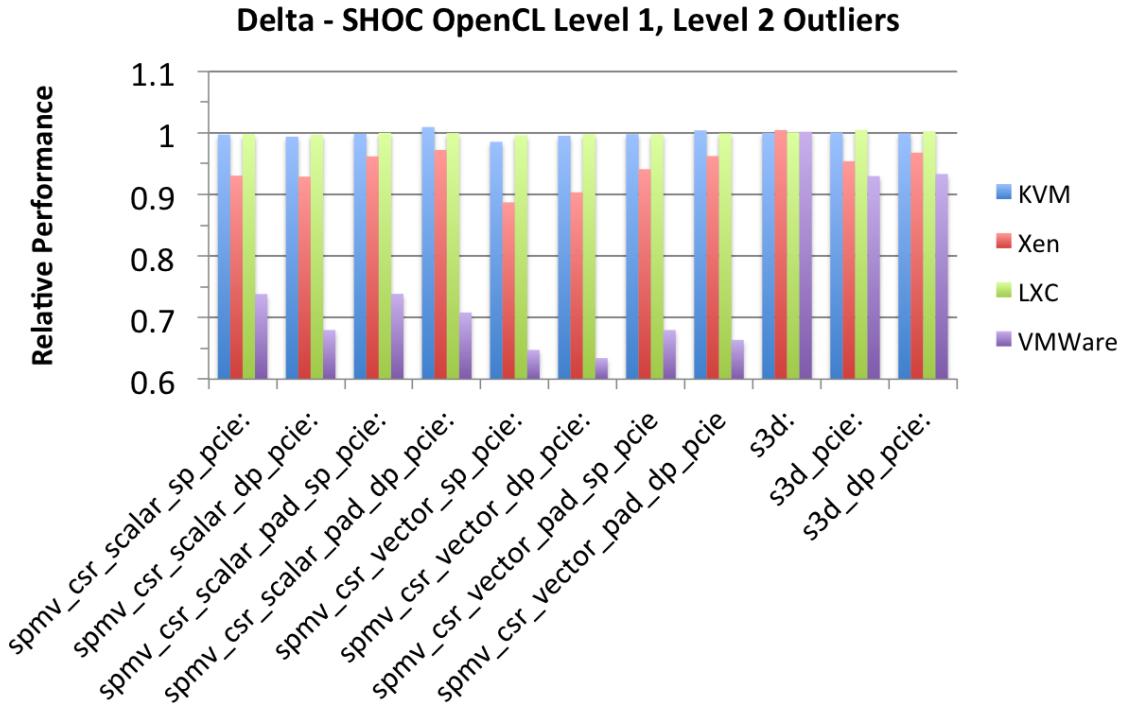
The remainder of Figure 5.1 includes a series of SHOC's Level 1 benchmarks, representing computational kernels. This includes BFS, FFT, molecular dynamics, and reduction kernels. Notably, nearly all of the benchmarks exhibiting overhead are the



**Figure 5.3** SHOC Levels 0 and 1 relative performance on Delta system. Benchmarks shown are the same as Bespin's. Higher is better.

2137 PCIe portion of SHOC's benchmarks. This is unsurprising, since the Level 0 benchmarks  
 2138 suggest PCIe bandwidth as the major source overhead. Still, results remain  
 2139 consistent with the bspeed\_download overhead observed in the Level 0 benchmarks,  
 2140 further suggesting that host/device data movement is the major source of overhead.

2141 In Figure 5.2 we present the remaining SHOC Level 1 results as well as the SHOC  
 2142 Level 2 results (S3D). In these results we begin to see an interesting trend, namely  
 2143 that VMWare consistently outperforms the base system in the Spmv and S3D mi-  
 2144 crobenchmarks on the Bespin system. We believe this to be a performance regression  
 2145 in CentOS 6.4, rather than a unique improvement due to the VMWare ESXi hyper-  
 2146 visor. When running the benchmarks on a non-virtualized Bespin system with Arch  
 2147 Linux, the VMWare ESXi performance gains were erased. An interesting finding of  
 2148 this was that Spmv was unique in this way - no other benchmarks were affected by



**Figure 5.4** SHOC Levels 1 and 2 relative performance. Benchmarks shown are the same as Bespin's. Higher is better.

2149 this performance issue.

2150 Turning to the Delta system, in Figures 5.3 and 5.4, we show the same benchmarks  
 2151 for the Delta system as was shown in Figures 5.1 and 5.2. Again, we find that  
 2152 the same benchmarks are responsible for most of the overhead on the Delta system.  
 2153 This is unsurprising, since PCIe was shown to be the source of the bulk of the over-  
 2154 head. A major difference in the case of the Delta system, however, is the amount  
 2155 of overhead. While the Bespin system saw overheads of approximately 1%, Delta's  
 2156 overhead routinely jumps above 35%, especially in the case of the Spmv benchmark  
 2157 for VMWare.

2158 On further examination, we determined that Xen was unable to activate IOMMU  
 2159 large page tables on the Delta system. KVM successfully allocated 4k, 2M, and 1G  
 2160 page table sizes, while Xen was limited to size 4k page tables. The Bespin system

2161 was able to take advantage of 4k, 2M, and 1G page sizes on both KVM and Xen. It  
2162 appears that this issue is correctable and does not represent a fundamental limitation  
2163 to the Xen hypervisor on the Nehalem/Westmere microarchitecture. While as a  
2164 closed source product, we have limited insight into VMWare ESXi, we speculate that  
2165 VMWare may be experiencing a similar issue on the Delta system and not on our  
2166 Bespin system.

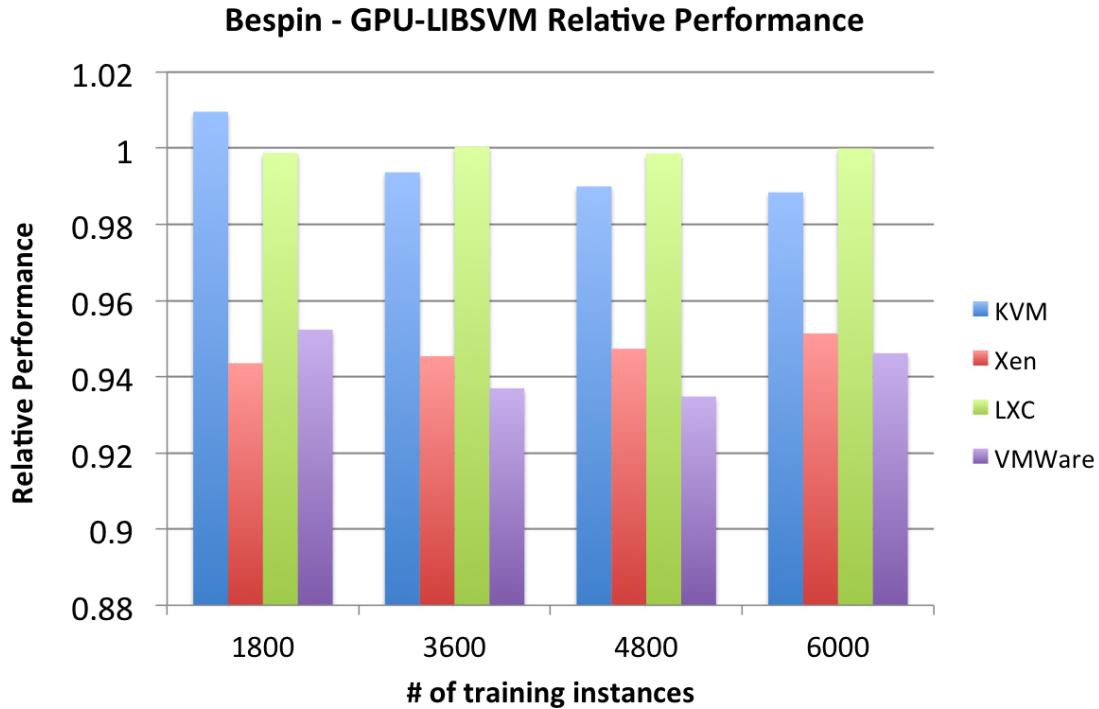
2167 In light of this, we broadly find that virtualization overhead across hypervisors and  
2168 architectures is minimal. Questions remain as to the source of the exceptionally high  
2169 overhead in the case of Xen and VMWare on the Delta system, but because KVM  
2170 shows no evidence of this overhead, we believe the Westmere/Fermi architecture to  
2171 be suitable for VGA passthrough in a cloud environment. In the case of the Bespin  
2172 system, it is clear that VGA passthrough can be achieved across hypervisors with  
2173 virtually no overhead.

2174 A surprising finding is that LXC showed little performance advantage over KVM,  
2175 Xen, and VWMare. While we expected near-native performance from LXC we did  
2176 not expect the hardware-assisted hypervisors to achieve such high performance. Still,  
2177 LXC carries some advantages. In general, its maximum overheads are comparable to  
2178 or less than KVM, Xen, and VMWare, especially in the case of the Delta system.

2179 **5.4.2 GPU-LIBSVM Performance**

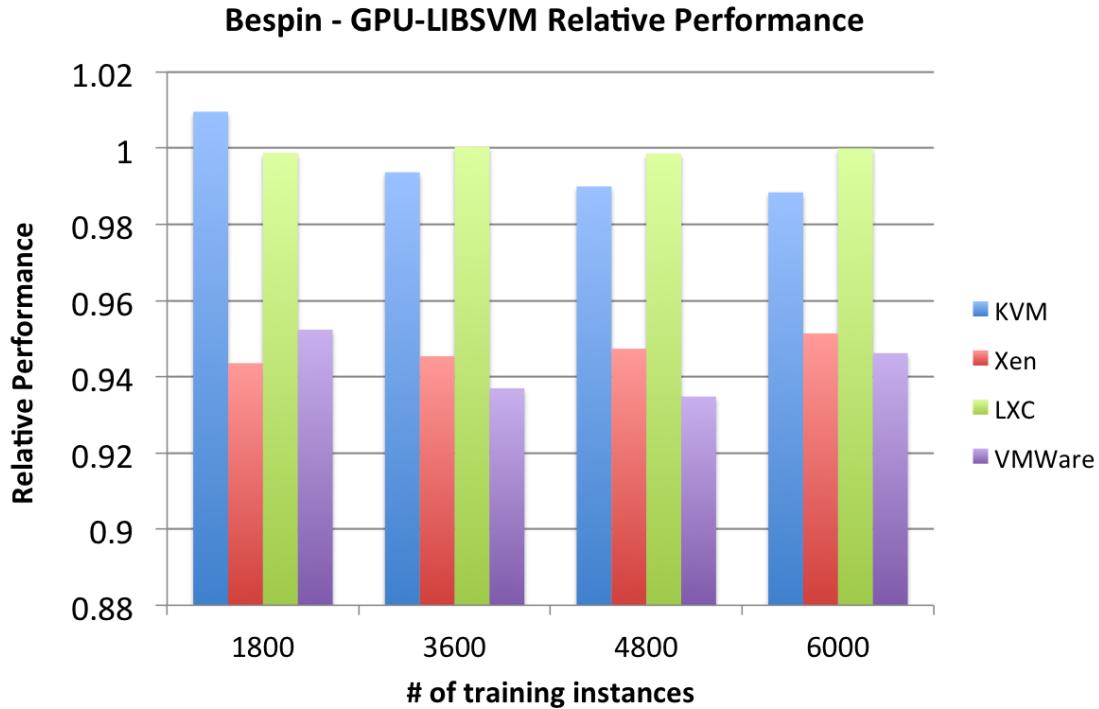
2180 In Figures 5.5 and 5.6, we display our GPU-LIBSVM performance results for the  
2181 Bespin (K20m) and Delta (C2075) systems. Using the NIPS 2003 gisette data set,  
2182 we show the performance across 4 problems sizes, ranging from 1800 to 6000 training  
2183 instances. The NIPS 2003 gisette dataset is a standard dataset that was used as a  
2184 part of the NIPS 2003 feature selection challenge.

2185 KVM again performs well across both the Delta and Bespin systems. In the case



**Figure 5.5** GPU-LIBSVM relative performance on Bespin system. Higher is better.

of the Delta system, in fact, KVM, significantly outperforms the base system. We determined this to be caused by KVM's support for transparent hugepages. When sufficient memory is available, transparent hugepages may be used to back the entirety of the guest VM's memory. Hugepages have previously been shown to improve TLB performance for KVM guests, and have been shown to occasionally boost the performance of a KVM guests beyond its host [197]. After disabling hugepages on the KVM host for the Delta system, performance dropped to 80–87% of the base system, suggesting that memory optimizations such as transparent hugepages can substantially improve the performance of virtualized guests under some circumstances. LXC and VMWare perform close to the base system, while Xen achieves between 72–90% of the base system's performance. We speculate that this may be related to Xen's inability to enable page sizes larger than 4k.

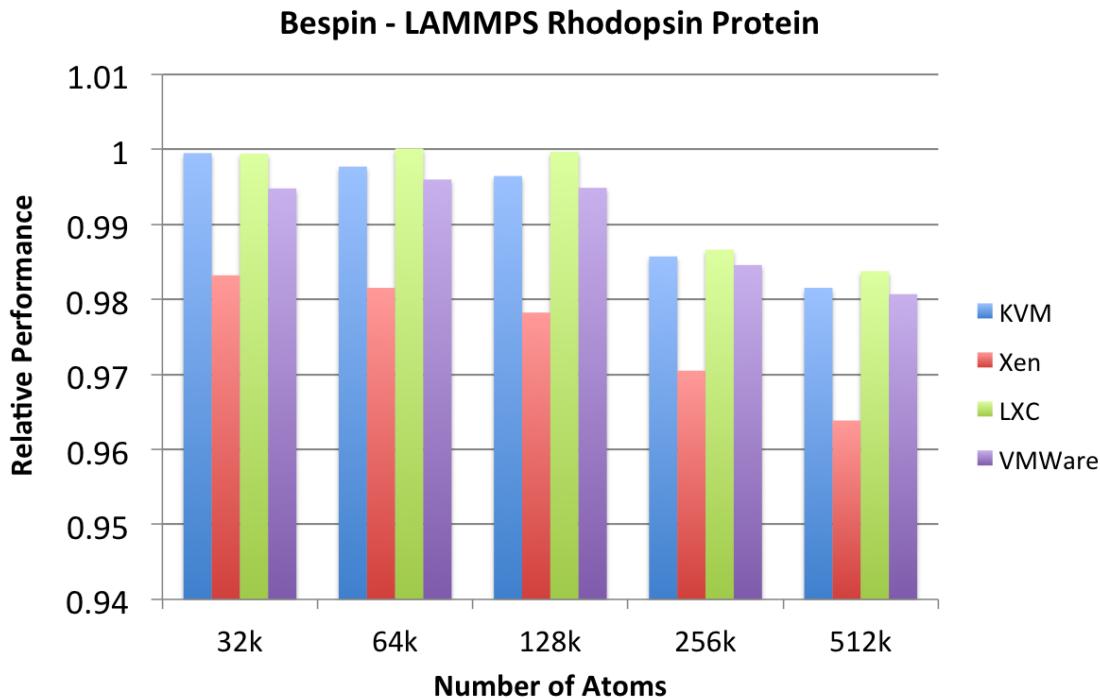


**Figure 5.6** GPU-LIBSVM relative performance on Delta showing improved performance within virtual environments. Higher is better.

#### 2198 5.4.3 LAMMPS Performance

2199 In Figures 5.7 and 5.8, we show the LAMMPS Rhodopsin protein simulation results.  
 2200 LAMMPS is unique among our benchmarks, in that it exercises both the GPU and  
 2201 multiple CPU cores. In keeping with the LAMMPS benchmarking methodology,  
 2202 we execute the benchmark using 1 GPU and 1–8 CPU cores on the Bespin system,  
 2203 selecting the highest performing configuration. In the case of the Delta system, we  
 2204 execute the benchmark on 1–6 cores and 1 GPU, selecting the highest performing  
 2205 configuration.

2206 Overall, LAMMPS performs well across both hypervisors and systems. Surpris-  
 2207 ingly, LAMMPS showed better efficiency on the Delta system than the Bespin system,  
 2208 achieving greater than 98% efficiency across the board, while Xen on the Bespin sys-



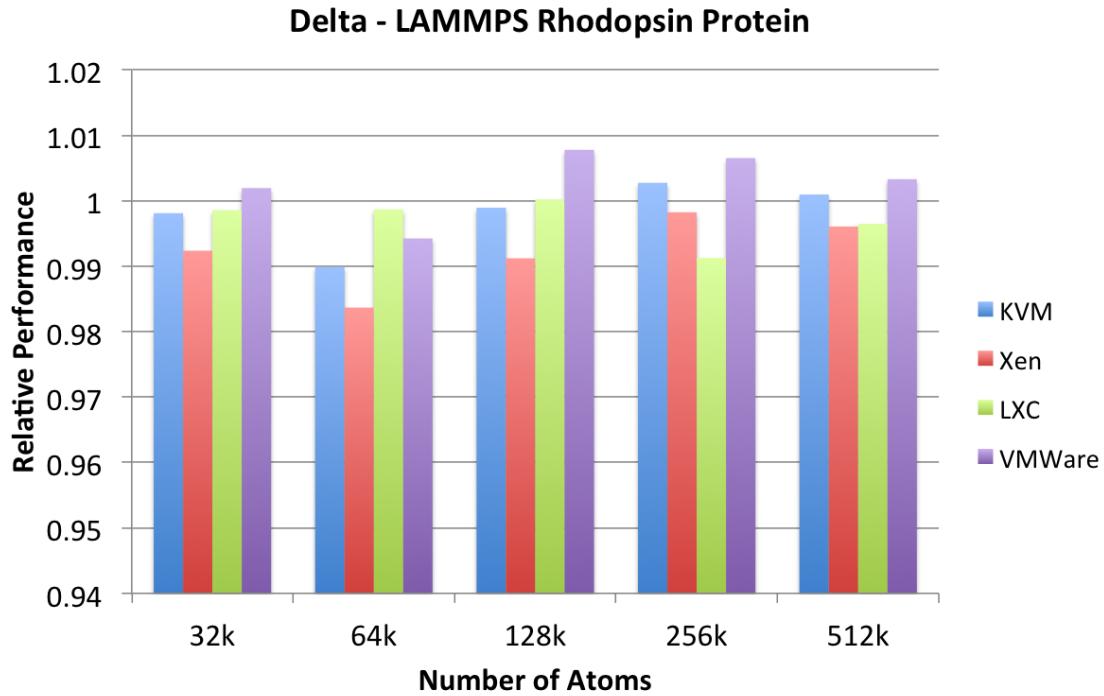
**Figure 5.7** LAMMPS Rhodopsin benchmark relative performance for Bespin system. Higher is better.

tem occasionally drops as low as 96.5% efficiency.

This performance is encouraging because it suggests that even heterogeneous CPU + GPU code is capable of performing well in a virtualized environment. Unlike SHOC, GPU-LIBSVM, and LULESH, LAMMPS fully exercises multiple host CPUs, splitting work between one or more GPUs and one or more CPU cores. This has the potential to introduce additional performance overhead, but the results do not bear this out in the case of LAMMPS.

#### 5.4.4 LULESH Performance

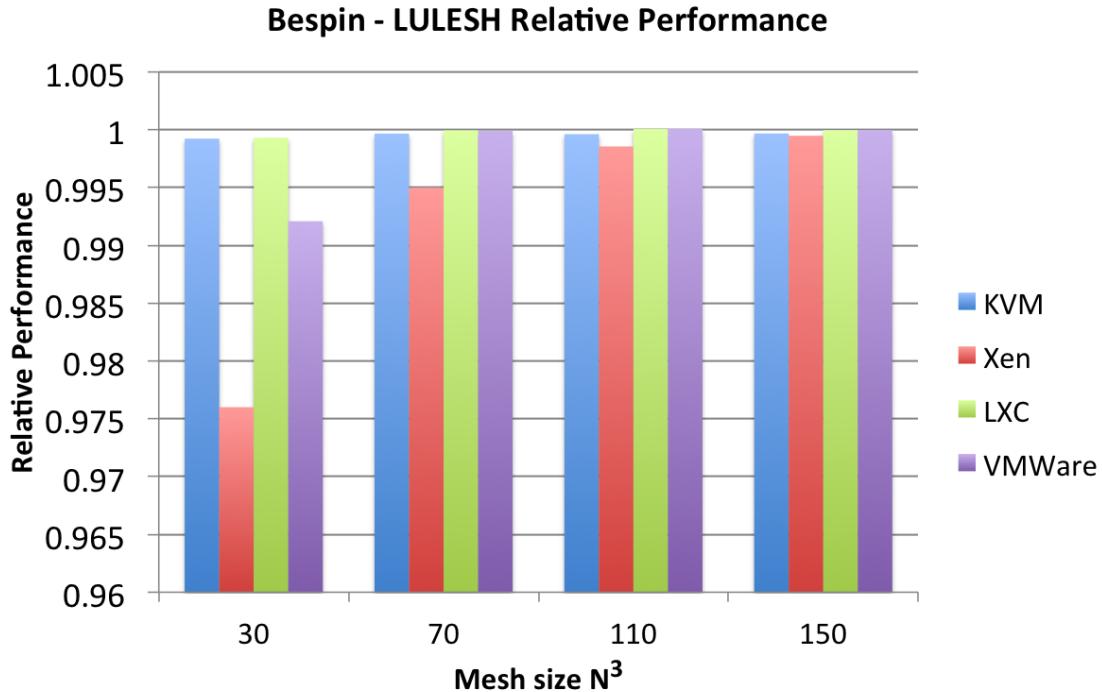
In Figure 5.9 we present our LULESH results for problem sizes ranging from mesh resolutions of  $N = 30$  to  $N = 150$ . LULESH is a highly compute-intensive simulation, with limited data movement between the host/virtual machine and the GPU, making



**Figure 5.8** LAMMPS Rhodopsin benchmark relative performance for Delta system. Higher is better.

it ideal for GPU acceleration. Consequently, we would expect little overhead due to virtualization. We show LULESH results only on the Bespin system, because differences in the code bases between the Kepler and Fermi implementations led to unsound comparisons.

While, overall, we see very little overhead, there is a slight scaling effect that is most apparent in the case of the Xen hypervisor. As the mesh resolution ( $N^3$ ) increases from  $N = 30$  to  $N = 150$ , we see that the Xen overhead decreases until Xen performs on-par with KVM, LXC, and VMWare.



**Figure 5.9** LULESH relative performance on Bespin. Higher is better.

## 2228 5.5 Lessons Learned

2229 Virtualizing performance-critical workloads has always proven controversial, whether  
 2230 the workload is CPU-only [169] or CPU with GPU. From our Westmere results, we  
 2231 can see that this criticism is in part legitimate, at times resulting in a performance  
 2232 penalty of greater than 10%, especially in the case of the VMWare ESXi and Xen  
 2233 hypervisors. We believe much of this to be fixable, especially in the case of the Xen  
 2234 hypervisor.

2235 At the same time, however, we have shown that the Sandy Bridge processor genera-  
 2236 tion has nearly erased those performance overheads, suggesting that old arguments  
 2237 against virtualization for performance-critical tasks should be reconsidered. In light of  
 2238 this, the primary lesson from this study is that VGA-passthrough to virtual machines  
 2239 is achievable at low overhead, and across a variety of hypervisors and virtualization

2240 platforms. Virtualization performance remains inconsistent across hypervisors for the  
2241 Westmere generation of processors, but starting with the Sandy Bridge architecture,  
2242 performance and consistency increase dramatically. In the case of the Sandy Bridge  
2243 architecture, even the lowest performing hypervisor, open source Xen, typically per-  
2244 forms within 95% of the base case.

2245 This study has also yielded valuable insight into the merits of each hypervisor.  
2246 KVM consistently yielded near-native performance across the full range of bench-  
2247 marks. Its support for transparent hugepages resulted in slight performance boosts  
2248 over-and-above even the base CentOS system in the case of the Delta system.

2249 VMWare's performance proved inconsistent across architectures, performing well  
2250 in the case of Bespin, and relatively poorly in the case of the Delta system. Because  
2251 hypervisor configurations were identical across systems, we can only speculate that  
2252 VMWare's performance is aided by the virtualization improvements offered by the  
2253 Sandy Bridge microarchitecture.

2254 The Xen hypervisor was consistently average across both architectures, perform-  
2255 ing neither poorly nor extraordinarily well in any individual benchmark. Xen and  
2256 VMWare ESXi are the only two hypervisors from this study that officially support  
2257 VGA passthrough. As a result, PCI passthrough support in both Xen and VMWare is  
2258 more robust than KVM. We expect that this advantage will not last long, as commer-  
2259 cial solutions targeting PCI passthrough in KVM are becoming common, particularly  
2260 with regard to SR-IOV and networking adapters.

2261 Linux Containers (LXC), consistently performed closest to the native case. This,  
2262 of course, is not surprising given that LXC guests share a single kernel with their  
2263 hosts. This performance comes at the cost of both flexibility and security, however.  
2264 LXC is less flexible than its full virtualization counterparts, offering support for only  
2265 Linux guests. More importantly, LXC device passthrough has security implications

2266 for multi-GPU systems. In the case of a multi-GPU-enabled host with NVIDIA  
2267 hardware, both GPUs must be passed to the LXC guest in order to initialize the  
2268 driver. This limitation may be addressable in future revisions to the NVIDIA driver.

## 2269 **5.6 Directions for Future Work**

2270 In this chapter we have characterized the performance of 4 common hypervisors across  
2271 two generations of GPUs and two host microarchitectures, and across 3 sets of bench-  
2272 marks. We showed the dramatic improvement in virtualization performance between  
2273 the Fermi/Westmere and the Kepler/Sandy Bridge system, with the Sandy Bridge  
2274 system typically performing within 1% of the base system. Finally, this study sought  
2275 to characterize the GPU and CPU+GPU performance with carefully tuned hypervi-  
2276 sor and guest configurations, especially with respect to NUMA. Improvements must  
2277 be made to today’s hypervisors in order to improve virtual NUMA support. Finally,  
2278 cloud infrastructure, such as OpenStack, must be capable of automatically allocating  
2279 virtual machines in a NUMA-friendly manner in order to achieve acceptable results  
2280 at cloud-scale.

2281 The next step in this work is to move beyond the single node to show that clus-  
2282 ters of accelerators can be efficiently used with minimal overhead. This will require  
2283 studies in high speed networking, particularly SR-IOV-enabled ethernet and Infini-  
2284 band. Special attention is needed to ensure that latencies remain tolerable within  
2285 virtual environments. Some studies have begun to examine these issues [198], but  
2286 open questions remain.

2287 **Chapter 6**

2288 **Supporting High Performance**  
2289 **Molecular Dynamics in Virtualized**  
2290 **Clusters using IOMMU, SR-IOV,**  
2291 **and GPUDirect**

2292 **6.1 Introduction**

2293 At present we stand at the inevitable intersection between High Performance Com-  
2294 puting (HPC) and clouds. Various platform tools such as Hadoop and MapReduce,  
2295 among others, have already percolated into data intensive computing within HPC [22].  
2296 In addition, there are efforts to support traditional HPC-centric scientific computing  
2297 applications in virtualized cloud infrastructure. There are a multitude of reasons for  
2298 supporting parallel computation in the cloud [59], including features such as dynamic  
2299 scalability, specialized operating environments, simple management interfaces, fault  
2300 tolerance, and enhanced quality of service, to name a few. The growing importance

2301 of supporting advanced scientific computing using virtualized infrastructure can be  
2302 seen by a variety of new efforts, including the NSF-funded Comet resource part of  
2303 XSEDE at San Diego Supercomputer Center [199].

2304 Nevertheless, there exists a past notion that virtualization used in today's cloud  
2305 infrastructure is inherently inefficient. Historically, cloud infrastructure has also done  
2306 little to provide the necessary advanced hardware capabilities that have become al-  
2307 most mandatory in supercomputers today, most notably advanced GPUs and high-  
2308 speed, low-latency interconnects. The result of these notions has hindered the use  
2309 of virtualized environments for parallel computation, where performance must be  
2310 paramount.

2311 A growing effort is currently underway that looks to systematically identify and  
2312 reduce any overhead in virtualization technologies. This effort has, thus far, proven  
2313 to be a qualified success [169, 200], though further research is needed to address  
2314 issues of scalability and I/O. Thus, we see a constantly diminishing overhead with  
2315 virtualization, not only with traditional cloud workloads [201] but also with HPC  
2316 workloads. While virtualization will almost always include some additional overhead  
2317 in relation to its dynamic features, the eventual goal for supporting HPC in virtualized  
2318 environments is to minimize what overhead exists whenever possible. To advance  
2319 the placement of HPC applications on virtual machines, new efforts are emerging  
2320 which focus specifically on key hardware now commonplace in supercomputers. By  
2321 leveraging new virtualization tools such as IOMMU device passthrough and SR-IOV,  
2322 we can now support the such advanced hardware as the latest Nvidia Tesla GPUs [202]  
2323 as well as InfiniBand fabrics for high performance networking and I/O [203, 204].

2324 With the advances in hypervisor performance coupled with the newfound avail-  
2325 ability of HPC hardware in virtual machines analogous to the most powerful super-  
2326 computers used today, we see can see the possibility of a high performance cloud in-

2327 frastructure using virtualization. While our previous efforts in this area have focused  
2328 on single-node advancements, it is now imperative to ensure real-world applications  
2329 can also operate in distributed environments as found in today’s cluster and cloud  
2330 infrastructures.

2331 Efforts to improve power efficiency and performance in data centers has led to  
2332 more heterogeneous architectures. That move toward heterogeneity has, in turn, led  
2333 to support for heterogeneity in the cloud. For example, Amazon EC2 supports GPU  
2334 accelerators in EC2 [205], and OpenStack supports heterogeneity using flavors [206].  
2335 These advancements in cloud-level support for heterogeneity combined with better  
2336 support for high-performance virtualization makes the use of cloud for HPC much  
2337 more feasible for a wider range of applications and platforms.

2338 In this paper we describe background a related work. Then, we describe a heteroge-  
2339 neous cloud platform, based on OpenStack. This effort has been under development  
2340 at USC/ISI since 2011 [163]. We describe our work towards integrating GPU and  
2341 InfiniBand support into OpenStack, and we describe the heterogeneous scheduling  
2342 additions that are necessary to support not only attached accelerators, but any cloud  
2343 composed of heterogeneous elements.

2344 We then demonstrate running two molecular dynamics simulations, LAMMPS  
2345 and HOOMD, in a virtual infrastructure complete with both Kepler GPUs and QDR  
2346 InfiniBand. Both HOOMD and LAMMPS are used extensively in some of the world’s  
2347 fastest supercomputers and represent example simulations that HPC supports today.  
2348 We show that these applications are able to run at near-native speeds within a com-  
2349 pletely virtualized environment, demonstrating just small performance impacts that  
2350 are usually acceptable by many users. Furthermore, we demonstrate the ability of  
2351 such a virtualized environment to support cutting edge software tools such as RDMA  
2352 GPUDirect, illustrating how cutting-edge HPC technologies are also possible in a

2353 virtualized environment.

2354 Following these efforts, we hope to ensure upstream infrastructure projects such  
2355 as OpenStack [128, 207] are able to make effective and quick use of these features,  
2356 allowing users to build private cloud infrastructure to support high performance dis-  
2357 tributed computational workloads.

## 2358 6.2 Background and Related Work

2359 Virtualization technologies and hypervisors have been seen widespread deployment  
2360 in support of a vast array of applications. This ranges from public commercial Cloud  
2361 deployments such as Amazon EC2 [208, 209], Microsoft Azure [210], and Google’s  
2362 Cloud Platform [211] to private deployments within colocation facilities, corporate  
2363 data centers, and even national scale cyber infrastructure initiatives. All these sup-  
2364 port look to support various use cases and applications such as web servers, ACID  
2365 and BASE databases, online object storage, and even distributed systems, to name a  
2366 few.

2367 The use of virtualization and hypervisors specifically support various HPC so-  
2368 lutions has been studied with mixed results. In [169], it is found that there is a  
2369 great deal of variance between hypervisors when running various distributed memory  
2370 and MPI applications, finding that KVM overall performed well across an array of  
2371 HPC benchmarks. Furthermore, some applications may not fit well into default  
2372 virtualized environments, such as High Performance Linpack [200]. Other studies  
2373 have specifically looked at interconnect performance in virtualization and found the  
2374 best-case scenario to be lacking [212] with up to 60% performance penalties with  
2375 conventional techniques.

2376 Recently, various CPU architectures have added support for I/O virtualization

2377 mechanisms in the CPU ISA through the use of an I/O memory management unit  
2378 (IOMMU). Often, this is referred to as PCI Passthrough, as it enabled devices on the  
2379 PCI-Express bus to be passed directly to a specific virtual machine (VM). Specific  
2380 hardware implementations include Intel’s VT-d [213], AMD’s IOMMU [214] from  
2381 x86\_64 architectures, and even more recently ARM System MMU [215]. All of these  
2382 implementations effectively look to aid in the usage of DMA-capable hardware to be  
2383 used within a specific virtual machine. Using these features, a wide array of hardware  
2384 can be utilized directly within VMs and enable fast and efficient computation and  
2385 I/O capabilities.

2386 With PCI Passthrough, a PCI device is handed directly to a running (or booting)  
2387 VM, thereby relinquishing control of the device within the host entirely. This is  
2388 different from typical VM usage where hardware is emulated in the host and used  
2389 in a guest VM, such as with bridged ethernet adapters or emulated VGA devices.  
2390 Performing PCI Passthrough requires the host to seize the device upon boot using a  
2391 specialized driver to effectively block normal driver initialization. In the instance of  
2392 the KVM hypervisor, this is done using the *vfio* and *pci\_stub* drivers. Then, this driver  
2393 relinquishes control to the VM, whereby normal device drivers initiate the hardware  
2394 and enable the device for use by the guest OS.

### 2395 6.2.1 GPU Passthrough

2396 Nvidia GPUs comprise the single most common accelerator in the Nov 2014 Top 500  
2397 List [8] and represent an increasing shift towards accelerators for HPC applications.  
2398 Historically, GPU usage in a virtualized environment has been difficult, especially  
2399 for scientific computation. Various front-end remote API implementations have been  
2400 developed to provide CUDA and OpenCL libraries in VMs, which translate library  
2401 calls to a back-end or remote GPU. One common use case of this is rCUDA [50], which

2402 provides a front-end CUDA API within a VM or any compute node, and then sends  
2403 the calls via Ethernet or InfiniBand to a separate node with 1 or more GPUs. While  
2404 this method is valid, it has the drawback of relying on the interconnect itself and the  
2405 bandwidth available, which can be especially problematic on Ethernet. Furthermore,  
2406 as this method consumes bandwidth, it can leave little remaining for MPI or RDMA  
2407 routines, thereby constructing a bottleneck for some MPI+CUDA applications that  
2408 depend on inter-process communication.

2409 Recently efforts have been seen to support such GPU accelerators within VMs  
2410 using IOMMU technologies, with implementations now available with KVM [202],  
2411 Xen [216] and VMWare [217]. These efforts have shown that GPUs can achieve up  
2412 to 99% of their bare metal performance when passed to a virtual machine using PCI  
2413 Passthrough. VMWare specifically shows how the such PCI Passthrough solutions  
2414 perform well and are likely to outperform front-end Remote API solutions such as  
2415 rCUDA within a VM [217]. While these works demonstrate PCI Passthrough perfor-  
2416 mance across a range of hypervisors and GPUs, they have been limited to investigating  
2417 single node performance until now.

2418 **6.2.2 SR-IOV and InfiniBand**

2419 With almost all parallel HPC applications, the interconnect fabric which enables fast  
2420 and efficient communication between processors becomes a central requirement to  
2421 achieving good performance. Specifically, a high bandwidth link is needed for dis-  
2422 tributed processors to share large amounts of data across the system. Furthermore,  
2423 low latency becomes equally important for ensuring quick delivery of small mes-  
2424 sage communications and resolving large collective barriers within many parallelized  
2425 codes. One such interconnect, InfiniBand, has become the most common implemen-  
2426 tation used within the Top500 list. However previously InfiniBand was inaccessible

2427 to virtualized environments.

2428 Supporting I/O interconnects in VMs has been aided by Single Root I/O Virtu-  
2429 alization (SR-IOV), whereby multiple virtual PCI functions are created in hardware  
2430 to represent a single PCI device. These virtual functions (VFs) can then be passed  
2431 to a VM and used as by the guest as if it had direct access to that PCI device. SR-  
2432 IOV allows for the virtualization and multiplexing to be done within the hardware,  
2433 effectively providing higher performance and greater control than software solutions.

2434 SR-IOV has been used in conjunction with Ethernet devices to provide high perfor-  
2435 mance 10Gb TCP/IP connectivity within VMs [218], offering near-native bandwidth  
2436 and advanced QoS features not easily obtained through emulated Ethernet offerings.  
2437 Currently Amazon EC2 offers a high performance VM solution utilizing SR-IOV en-  
2438 abled 10Gb Ethernet adapters. While SR-IOV enabled 10Gb Ethernet solutions offers  
2439 a big forward in performance, Ethernet still does not offer the high bandwidth or low  
2440 latency typically found with InfiniBand solutions.

2441 Recently SR-IOV support for InfiniBand has been added by Mellanox in the Con-  
2442 nectX series adapters. Initial evaluation of SR-IOV InfiniBand within KVM VMs  
2443 has proven has found point-to-point bandwidth to be near-native, but up to 30%  
2444 latency overhead for very small messages [203, 219]. However, even with the noted  
2445 overhead, this still signifies up to an order of magnitude difference in latency between  
2446 InfiniBand and Ethernet with VMs. Furthermore, advanced configuration of SR-IOV  
2447 enabled InfiniBand fabric is taking shape, with recent research showing up to a 30%  
2448 reduction in the latency overhead [204]. However, real application performance has  
2449 not yet been well understood until now.

2450 **6.2.3 GPUDirect**

2451 NVIDIA’s GPUDirect technology was introduced to reduce the overhead of data  
2452 movement across GPUs [220, 221]. GPUDirect supports both networking as well as  
2453 peer-to-peer interfaces for single node multi-GPU systems. The most recent imple-  
2454 mentation of GPUDirect, version 3, adds support for RDMA over InfiniBand for  
2455 Kepler-class GPUs.

2456 The networking component of GPUDirect relies on three key technologies: CUDA  
2457 5 (and up), a CUDA-enabled MPI implementation, and a Kepler-class GPU (RDMA  
2458 only). Both MVAPICH and OpenMPI support GPUDirect. Support for RDMA over  
2459 GPUDirect is enabled by the MPI library, given supported hardware, and does not  
2460 depend on application-level changes to a user’s code.

2461 In this paper, our GPUDirect work focuses on GPUDirect v3 for multi-node  
2462 RDMA support. We demonstrate scaling for up to 4 nodes connected via QDR  
2463 InfiniBand and show that GPUDirect RDMA improves both scalability and overall  
2464 performance by approximately 9% at no cost to the end user.

2465 **6.3 A Cloud for High Performance Computing**

2466 With support for GPU Passthrough, SR-IOV, and GPUDirect, we have the building  
2467 blocks for a high performance, heterogeneous cloud. In addition, other common  
2468 accelerators (e.g. Xeon Phi [222]) have similarly been demonstrated in virtualized  
2469 environments. Our vision is of a heterogeneous cloud, supporting both high speed  
2470 networking and accelerators for tightly coupled applications.

2471 To this end we have developed a heterogeneous cloud based on OpenStack [128].  
2472 In our previous work, we have demonstrated the ability to rapidly provision GPU,  
2473 bare metal, and other heterogeneous resources within a single cloud [163]. Building

2474 on this effort we have added support for GPU passthrough to OpenStack as well as  
2475 SR-IOV support for both ConnectX-2 and ConnectX-3 Infiniband devices. Mellanox  
2476 separately supports an OpenStack InfiniBand networking plugin for OpenStack’s Neu-  
2477 tron service [223], however the Mellanox plugin depends on the ConnectX-3 adapter.  
2478 Our institutional requirements depend on ConnectX-2 SR-IOV support, requiring  
2479 an independent implementation.

2480 OpenStack supports services for networking (Neutron), compute (Nova), identity  
2481 (Keystone), storage (Cinder, Swift), and others. Our work focuses entirely on the  
2482 compute service.

2483 Scheduling is implemented at two levels: the cloud-level and the node-level. In our  
2484 earlier work, we have developed a cloud-level heterogeneous scheduler for OpenStack,  
2485 allowing scheduling based on architectures and resources [163]. In this model, the  
2486 cloud-level scheduler dispatches jobs to nodes based on resource requirements (e.g.  
2487 Kepler GPU) and node-level resource availability.

2488 At the node, a second level of scheduling occurs to ensure that resources are  
2489 tracked and not over-committed. Unlike traditional cloud paradigms, devices passed  
2490 into VMs cannot be over-committed. We treat devices, whether GPUs or InfiniBand  
2491 virtual functions, as schedulable resources. Thus, it is the responsibility of the indi-  
2492 vidual node to track resources committed and report availability to the cloud-level  
2493 scheduler. For reporting, we piggyback on top of OpenStack’s existing reporting  
2494 mechanism to provide a low overhead solution.

## 2495 6.4 Benchmarks

2496 We selected two molecular dynamics (MD) applications for evaluation in this study:  
2497 LAMMPS and HOOMD [224, 225]. These MD simulations are chosen to represent a

2498 subset of advance parallel computation for a number of fundamental reasons:

- 2499     ● MD simulations provide a practical representation of N-Body simulations, which  
2500         is one of the major computational *Dwarfs* [226] in parallel and distributed com-  
2501         puting.
- 2502     ● MD simulations are one of the most widely deployed applications on large scale  
2503         supercomputers today.
- 2504     ● Many MD simulations have a hybrid MPI+CUDA programming model, which  
2505         has often become commonplace in HPC as the use of accelerators increases.

2506     As such, we look to LAMMPS and HOOMD to provide a real-world example for  
2507     running cutting-edge parallel programs on virtualized infrastructure. While these  
2508     applications by no means represent all parallel scientific computing efforts (as justi-  
2509         fied by the 13 Dwarfs defined in [226]), we hope these MD simulators offer a more  
2510         pragmatic viewpoint than traditional synthetic HPC benchmarks such as High Per-  
2511         formance Linpack.

2512     **LAMMPS** The Large-scale Atomic/Molecular Parallel Simulator is a well-understood  
2513     highly parallel molecular dynamics simulator. It supports both CPU and GPU-  
2514     based workloads. Unlike many simulators, both MD and otherwise, LAMMPS is  
2515     heterogeneous. It will use both GPUs and multicore CPUs concurrently. For this  
2516     study, this heterogeneous functionality introduces additional load on the host, allow-  
2517     ing LAMMPS to utilize all available cores on a given system. Networking in LAMMPS  
2518     is accomplished using a typical MPI model. That is, data is copied from the GPU  
2519     back to the host and sent over the InfiniBand fabric. No RDMA is used for these  
2520     experiments.

2521 **HOOMD-blue** The Highly Optimized Object-oriented Many-particle Dynamics  
2522 – Blue Edition is a particle dynamics simulator capable of scaling into the thou-  
2523 sands of GPUs. HOOMD supports executing on both CPUs and GPUs. Unlike  
2524 LAMMPS, HOOMD is homogeneous and does not support mixing of GPUs and  
2525 CPUs. HOOMD supports GPUDirect using a CUDA-enabled MPI. In this paper  
2526 we focus on HOOMD’s support for GPUDirect and show its benefits for increasing  
2527 cluster sizes.

## 2528 6.5 Experimental Setup

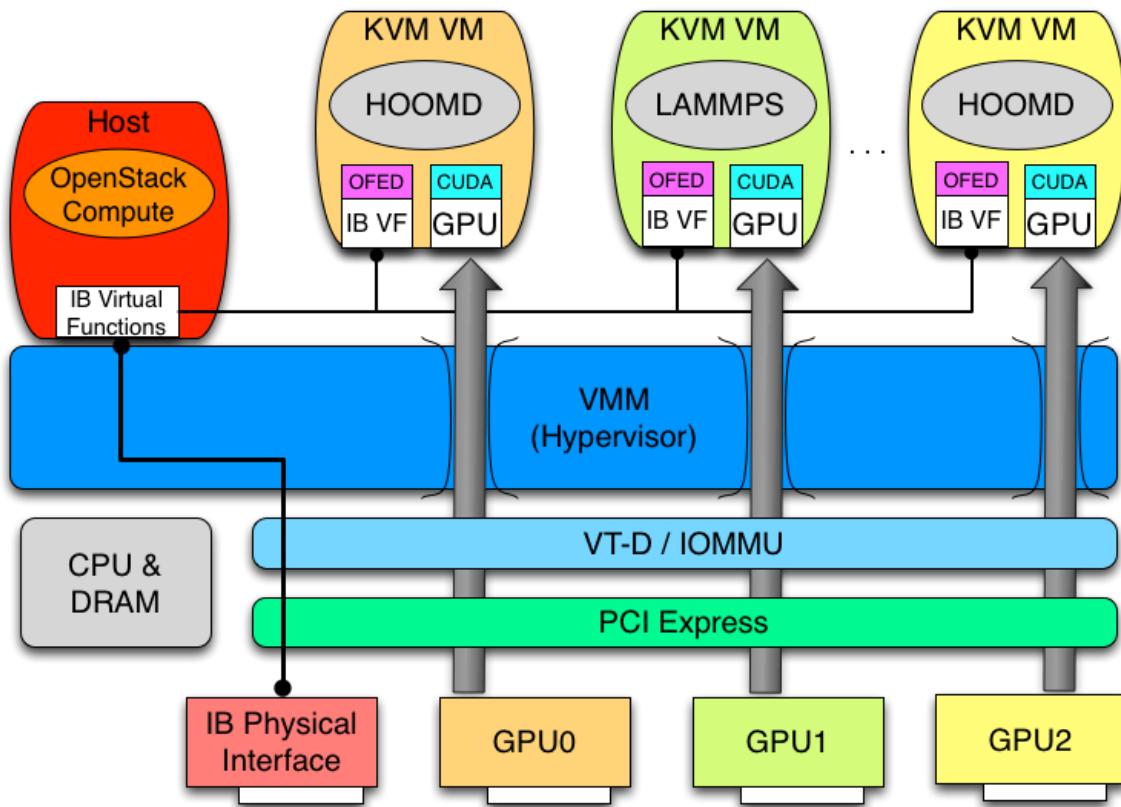
2529 Using two molecular dynamics tools, LAMMPS [224] and HOOMD [225], we demon-  
2530 strate a high performance *system*. That is, we combine PCI passthrough for Nvidia  
2531 Kepler-class GPUs with QDR Infiniband SR-IOV and show that high performance  
2532 molecular dynamics simulations are achievable within a virtualized environment.

2533 For the first time, we also demonstrate Nvidia GPUDirect technology within such  
2534 a virtual environment. Thus, we look to not only illustrate that virtual machines  
2535 provide a flexible high performance infrastructure for scaling scientific workloads in-  
2536 cluding MD simulations, but also that the latest HPC features and programming  
2537 environments are also available in this same model.

### 2538 6.5.1 Node configuration

2539 To support the use of Nvidia GPUs and InfiniBand within a VM, specific and exact  
2540 host configuration is needed. This node configuration is illustrated in Figure 6.1.  
2541 While our implementation is specific to the KVM hypervisor, this setup represents a  
2542 design that can be hypervisor agnostic.

2543 Each node in the testbed uses CentOS 6.4 with a 3.13 upstream Linux kernel for



**Figure 6.1** Node PCI Passthrough of GPUs and InfiniBand

2544 the host OS, along with the latest KVM hypervisor, QEMU 2.1, and the *vfio* driver.  
 2545 Each Guest VM runs CentOS 6.4 with a stock 2.6.32-358.23.2 kernel. A Kepler GPU  
 2546 is passed through using PCI Passthrough and directly initiated within the VM via  
 2547 the Nvidia 331.20 driver and CUDA release 5.5. While this specific implementation  
 2548 used only a single GPU, it is also possible to include as many GPUs as one can fit  
 2549 within the PCI Express bus if desired. As the GPU is used by the VM, an on-board  
 2550 VGA device was used by the host and a standard Cirrus VGA was emulated in the  
 2551 guest OS.

2552 With using SR-IOV, the OFED drivers version 2.1-1.0.0 are used with Mellanox  
 2553 ConnectX-3 VPI adapter with firmware 2.31.5050. The host driver initiates 4 VFs,  
 2554 one of which is passed through to the VM where the default OFED mlnx\_ib drivers

2555 are loaded.

### 2556 6.5.2 Cluster Configuration

2557 Our test environment is composed of 4 servers each with a single Nvidia Kepler-  
2558 class GPU. Two servers are equipped with K20 GPUs, while the other two servers  
2559 are equipped with K40 GPUs, demonstrating the potential for a more heterogeneous  
2560 deployment. Each server is composed of 2 Intel Xeon E5-2670 CPUs, 48GB of DDR3  
2561 memory, and Mellanox ConnectX-3 QDR InfiniBand. CPU sockets and memory are  
2562 split evenly between the two NUMA nodes on each system. All InfiniBand adapters  
2563 use a single Voltaire 4036 QDR switch with a software subnet manager for IPoIB  
2564 functionality.

2565 For these experiments, both the GPUs and InfiniBand adapters are attached to  
2566 NUMA node 1 and both the guest VMs and the base system utilized identical software  
2567 stacks. Each guest was allocated 20 GB of RAM and a full socket of 8 cores, and  
2568 pinned to NUMA node 1 to ensure optimal hardware usage. While all VMs are  
2569 capable of login via the InfiniBand IPoIB setup, a 1Gb Ethernet network was used  
2570 for all management and login tasks.

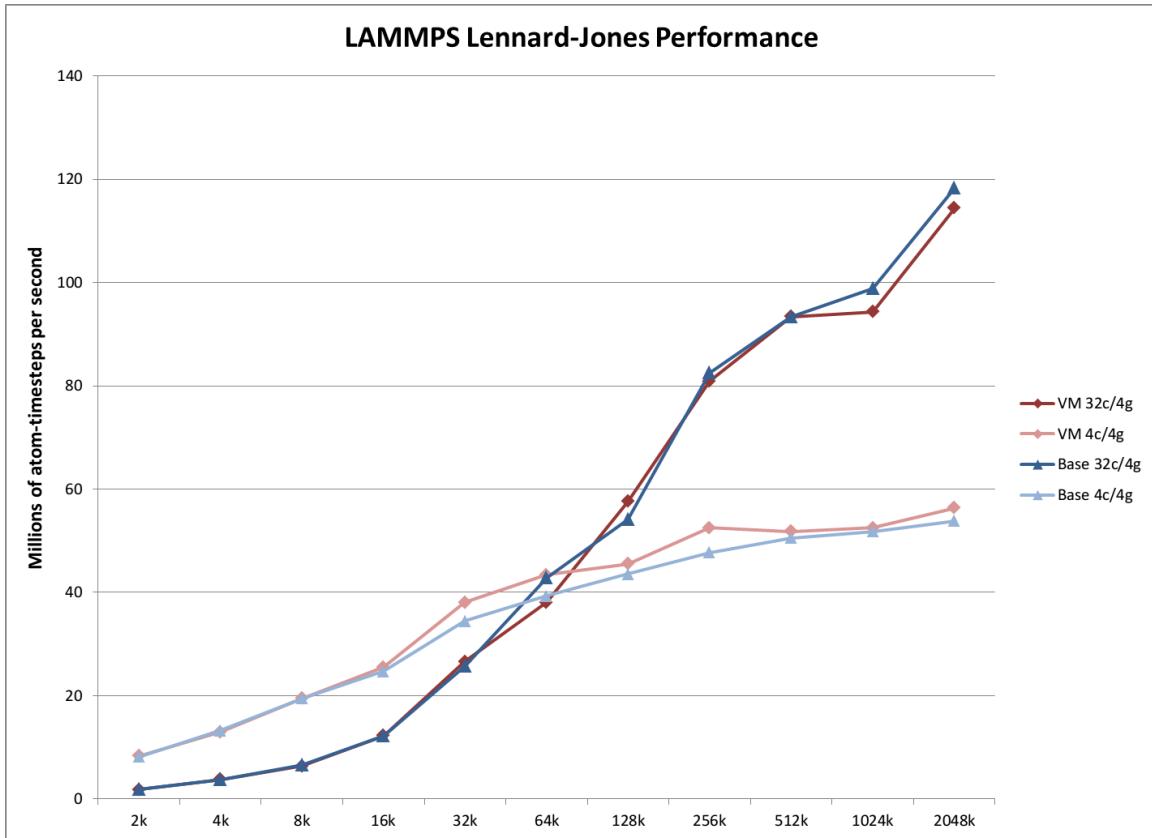
2571 For a fair and effective comparison, we also use a native environment without any  
2572 virtualization. This native environment employs the same hardware configuration,  
2573 and like the Guest OS runs CentOS 6.4 with the stock 2.6.32-358.23.2 kernel.

## 2574 6.6 Results

2575 In this section, we discuss the performance of both the LAMMPS and HOOMD  
2576 molecular dynamics simulation tools when running within a virtualized environment.  
2577 Specifically, we scale each application to 32 cores and 4 GPUs, both in a native bare-

2578 metal and virtualized environments. Each application set was run 10 times, with the  
 2579 results averaged accordingly.

2580 **6.6.1 LAMMPS**



**Figure 6.2** LAMMPS LJ Performance

2581 Figure 6.2 shows one of the most common LAMMPS algorithms used; the Lennard-  
 2582 Jones potential (LJ). This algorithm is deployed in two main configurations - a 1:1  
 2583 core to GPU mapping, and a 8:1 core to GPU mapping. With the LAMMPS GPU  
 2584 implementation, a delicate balance between GPUs and CPUs is required to find the  
 2585 optimal ratio for fastest computation, however here we just look at the two most  
 2586 obvious choices. With small problem sizes, the 1:1 mapping outperforms the more  
 2587 complex core deployment, as the problem does not require the additional complexity

2588 provided with multi-core solution. As expected the multi-core configuration quickly  
2589 offers better performance for larger problem sizes, achieving roughly twice the perfor-  
2590 mance with all 8 available cores. This is largely due to the availability of all 8 cores  
2591 to keep the GPU running 100% with continual computation.

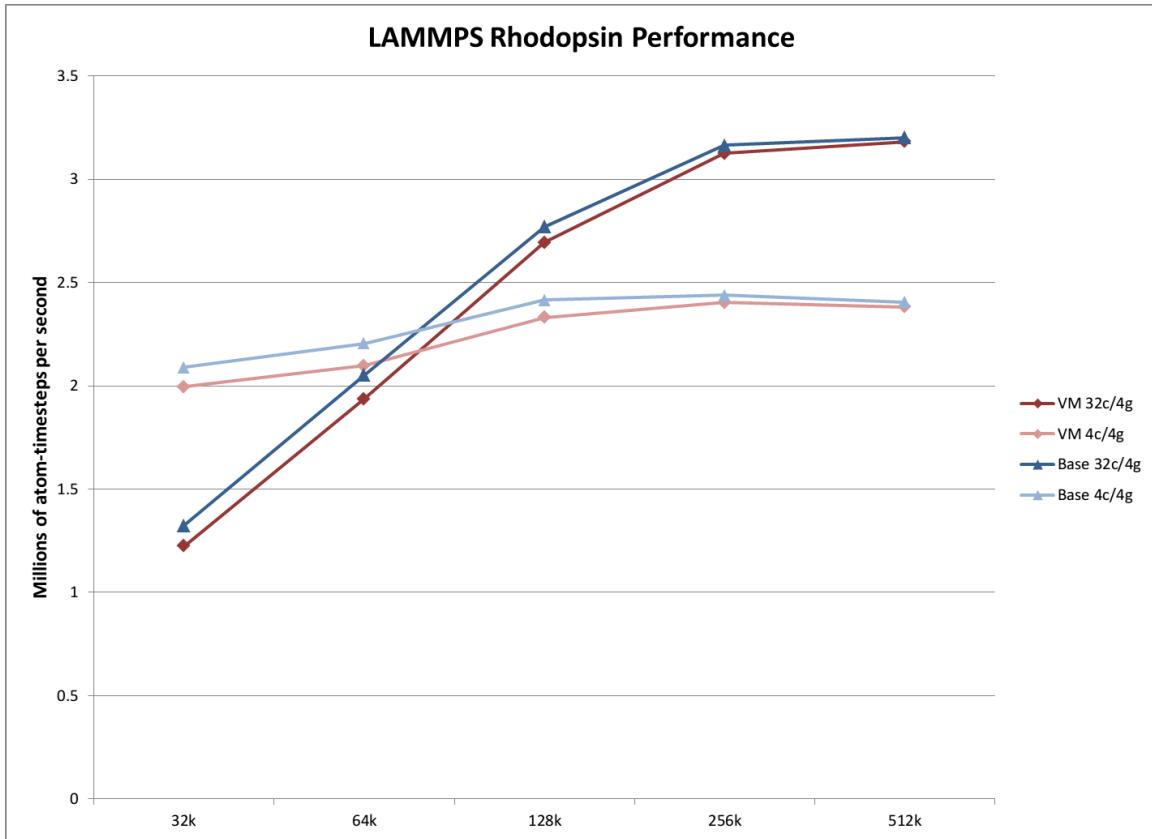
2592 The important factor for this manuscript is the relative performance of the virtu-  
2593 alized environment. From the results, it is clear the VM solution performs very well  
2594 compared to the best-case native deployment. For the multi-core configuration across  
2595 all problem sizes, the virtualized deployment averaged 98.5% efficiency compared to  
2596 native. The single core per GPU deployment reported better-than native perfor-  
2597 mance at 100% native. This is likely due to caching effects, but further investigation  
2598 is needed to fully identify this occurrence.

2599 Another common LAMMPS algorithm, the Rhodopsin protein in solvated lipid  
2600 bilayer benchmark (Rhodo), was also run with results given in Figure 6.3. As with  
2601 the LJ runs, we see the multi-core to GPU configuration resulting in higher computa-  
2602 tional performance for the larger problem sizes compared to the single core per GPU  
2603 configuration, as expected.

2604 Again, the overhead of the virtualized configuration remains low across all con-  
2605 figurations and problem sizes, with an average 96.4% efficiency compared to native.  
2606 Interestingly enough, we also see the performance gap decrease as the problem size  
2607 increases, with the 512k problem size in yielding 99.3% of native performance. This  
2608 finding leads us to extrapolate that a virtualized MPI+CUDA implementation would  
2609 scale to a larger computational resource with similar success.

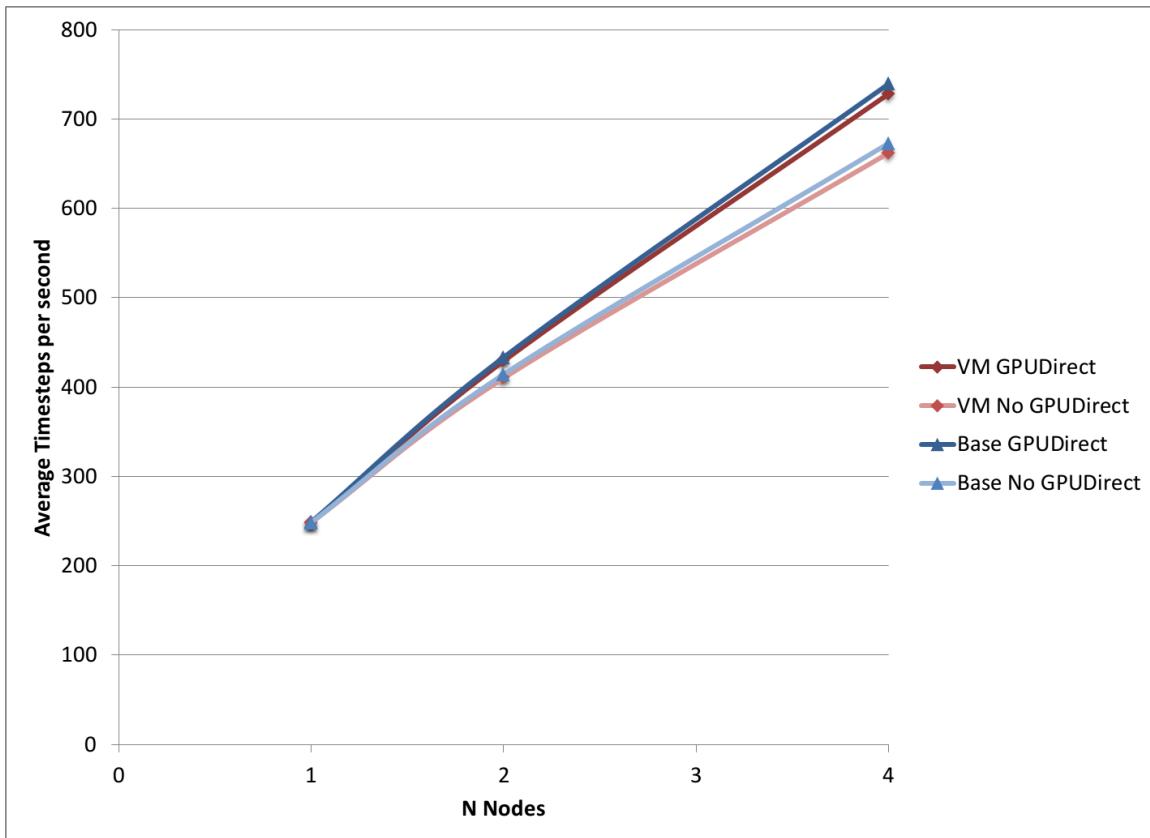
2610 **6.6.2 HOOMD**

2611 In Figure 6.4 we show the performance of a Lennard-Jones liquid simulation with 256K  
2612 particles running under HOOMD. HOOMD includes support for CUDA-aware MPI



**Figure 6.3** LAMMPS RHODO Performance

2613 implementations via GPUDirect. The MVAPICH 2.0 GDR implementation enables  
 2614 a further optimization by supporting RDMA for GPUDirect. From Figure 6.4 we  
 2615 can see that HOOMD simulations, both with and without GPUDirect, perform very  
 2616 near-native. The GPUDirect results at 4 nodes achieve 98.5% of the base system's  
 2617 performance. The non-GPUDirect results achieve 98.4% efficiency at 4 nodes. These  
 2618 results indicate the virtualized HPC environment is able to support such complex  
 2619 workloads. While the effective testbed size is relatively small, it indicates that such  
 2620 workloads may scale equally well to hundreds or thousands of nodes.



**Figure 6.4** HOOMD LJ Performance with 256k Simulation

## 2621 6.7 Discussion

2622 From the results, we see the potential for running HPC applications in a virtualized  
 2623 environment using GPUs and InfiniBand interconnect fabric. Across all LAMMPS  
 2624 runs with ranging core configurations, we found only a 1.9% overhead between the  
 2625 KVM virtualized environment and native. For HOOMD, we found a similar 1.5%  
 2626 overhead, both with and without GPU Direct. These results go against conventional  
 2627 wisdom that HPC workloads do not work in VMs. In fact ,we show two N-Body  
 2628 type simulations programmed in an MPI+CUDA implementation perform at roughly  
 2629 near-native performance in tuned KVM virtual machines.

2630 With HOOMD, we see how GPUDirect RDMA shows a clear advantage over

2631 the non-GPUDirect implementation, achieving a 9% performance boost in both the  
2632 native and virtualized experiments. While GPUDirect's performance impact has been  
2633 well evaluated previously [220], it is the author's belief that this manuscript represents  
2634 the first time GPUDirect has been utilized in a virtualized environment.

2635 Another interesting finding of running LAMMPS and HOOMD in a virtualized  
2636 environment is as workload scales from a single node to 32 cores, the overhead does  
2637 not increase. These results lend credence to the notion that this solution would also  
2638 work for a much larger deployment. Specifically, it would be possible to expand  
2639 such computational problems to a larger deployment in FutureGrid [227], Chameleon  
2640 Cloud [228], or even the planned NSF Comet machine at SDSC, scheduled to provide  
2641 up to 2 Petaflops of computational power. Effectively, these results help support  
2642 the theory that a majority of HPC computations can be supported in virtualized  
2643 environment with minimal overhead.

## 2644 6.8 Chapter Summary

2645 With the advent of cloud infrastructure, the ability to run large-scale parallel sci-  
2646 entific applications has become possible but limited due to both performance and  
2647 hardware availability concerns. In this work we show that advanced HPC-oriented  
2648 hardware such as the latest Nvidia GPUs and InfiniBand fabric are now available  
2649 within a virtualized infrastructure. Our results find MPI + CUDA applications such  
2650 as molecular dynamics simulations run at near-native performance compared to tra-  
2651 ditional non-virtualized HPC infrastructure, with just an averaged 1.9% and 1.5%  
2652 overhead for LAMMPS and HOOMD, respectively. Moving forward, we show the  
2653 utility of GPUDirect RDMA for the first time in a cloud environment with HOOMD.  
2654 Effectively, we look to pave the way for large-scale virtualized cloud Infrastructure

2655 to support a wide array of advanced scientific computation commonly found running  
2656 on many supercomputers today. Our efforts leverage these technologies and provide  
2657 them in an open source Infrastructure-as-a-Service framework using OpenStack.

2658 **Chapter 7**

2659 **Virtualization advancements to**

2660 **support HPC applications**

2661 Many of the previous chapters have focused on identifying and decreasing the per-  
2662 formance gap that exists with running HPC workloads in virtualized infrastructure,  
2663 as compared to native HPC environments without virtualization. While this is a  
2664 significant technical challenge to overcome, the use of virtualization itself has the  
2665 potential to offer additional benefits to HPC infrastructure. In this chapter, we look  
2666 to identify research, design, and future implementation of advanced virtualization  
2667 techniques to enable a new class of HPC infrastructure with added performance and  
2668 features that has yet to be realized. It may be possible for these advancements, once  
2669 implemented to have an impact not only on cloud infrastructure, but also a dedicated  
2670 HPC environment for added usability and performance.

## 2671 7.1 Memory Page Table Optimizations

2672 As we've seen both in Chapter 3 as well as in other supported literature [42], virtual-  
2673 ization of memory structures becomes a point of contention and potential overhead.  
2674 This is often due to the extensive effort a hypervisor has to perform in order to  
2675 translate memory addresses from guest-virtual addresses to host-virtual addresses,  
2676 and then again to machine-physical addresses. Historically, this was handled using  
2677 Shadow Page tables [229], which eliminate the need for emulation of physical memory  
2678 inside the VM by creating a page table mapping from guest virtual to machine mem-  
2679 ory. However, these page tables are not walkable by hardware such as a transition  
2680 lookaside buffer (TLB), and as such guest OS page tables require updating of the  
2681 shadow page table. This can be costly not only in the additional management, but  
2682 also notably by the cost of VMexit and VMentry calls.

2683 Recently, Intel and AMD have implemented Extended Page Tables (EPT) and  
2684 nested paging, respectively. With EPT, support is added to allow the TLB hardware  
2685 to keep track of both guest pages and hypervisor pages concurrently, effectively re-  
2686 moving the need for shadow page table. The downside of EPT and nested paging  
2687 occurs when there is a TLB miss (the page isn't in the TLB), and as such each miss  
2688 requires a walk through each VMM nested paging, effectively creating TLB miss cost  
2689 of 16 table walks (instead of just 3-4) for 4k pages. While many applications find  
2690 this TLB miss additional cost less much less than that of managing shadow page  
2691 tables, it still can lead to a significant gap in performance between non-virtualized  
2692 applications, especially as VM count or an application's memory footprint increase.

2693 One potential way to decrease the chance of a TLB miss (and therefore the cost  
2694 of a miss) is by using a larger page size. By default, x86 hardware uses 4k pages  
2695 sizes, but newer hardware can support 2M and 1G page sizes as well, effectively

2696 named *transparent huge pages* or THP. Using the KVM hypervisor with transparent  
2697 hugepages enabled, we can create guest VMs backed entirely 2M hugepages [197]. We  
2698 can also enable transparant hugepage support within the guest as well, to have the  
2699 entire guest OS (including kernel and modules) using 2M pages.

2700 The result of THP-enabled guest VMs can be significant. With huge pages on Intel  
2701 x86 CPUs with EPT, there exists an entirely separate TLB for hugepages as well. This  
2702 will naturally alleviate TLB pressure and therefore reduce TLB contention between  
2703 guest and host operating systems. Because 2M pages provide larger addressable  
2704 memory, the size of the page tables themselves are also decreased. Effectively, this  
2705 reduces the TLB miss cost from 4 to 3 page table walks, which when handling a VM  
2706 TLB miss, then requires only 15 walks to the default 24. This in effect can have a  
2707 significant improvement in VM performance, as hypothesized in [197].

2708 To evaluate the effect of 2M transparent huge pages on guest performance, we  
2709 leverage the same KVM setup in Chapter 5 on the Bespin hardware. Specifically, THP  
2710 is switched both in the host as well as the guest OS kernels, and the same LibSVM  
2711 application using GPUs is re-run. The libSVM GPU application can have significant  
2712 memory requirements, as large chunks of the support vector machine datasets are  
2713 stored in CPU memory and tranferred in a sudo-random order to and from GPU  
2714 memory, making it an ideal application to use for evaluating THP.

2715 The results of running libSVM in a THP-enabled VM, a VM with no THP, and  
2716 natively without virtualization are displayed in Figure 7.1. Comparing first just  
2717 the KVM results without THP to the native solution, we can see the impact of  
2718 THP the overal application runtime, especially at larger problem sizes (6000 training  
2719 sets). However, when TLB is enabled in the guest and host, we actually see the  
2720 KVM VM solution *outperform* the native solution. This is because guest privileged  
2721 OS memory used to buffer to/from GPU memory is backed by 2M pages, instead



**Figure 7.1** Transparent Huge Pages with KVM

of the normal 4k pages as in the native solution. The result is less TLB misses during application runtime compared to 4k TLB misses in a native solution, resulting in improved performance. While this is likely a special case for THP usage with the libSVM application, the fact that a VM can even occasionally outperform a native runtime is a noteworthy accomplishment. This also underscores the need for careful tuning and best-practices for hypervisors when supporting advanced scientific workloads.

**2729 7.2 Live Migration Mechanisms**

2730 Migration of VMs represents one of the fundamental advantages to virtualization, and  
2731 also one of the greatest challenges to efficiency. With VM migration, the complete VM  
2732 state is copied from a source to a unallocated destination host, where disk, memory,  
2733 and network connections are kept intact. For disk continuity, a distributed and/or  
2734 shared filesystem is utilized, most commonly but not exclusively NFS, where both  
2735 the source and destination hosts have access to the VM disk. Network continuity is  
2736 preserved so long as the destination guest is within the same LAN and generates an  
2737 unsolicited ARP reply to maintain the original IP after migration. VM vCPU states  
2738 and machine states are recorded from the source and quickly sent to the destination  
2739 host when the VM is paused. For live migration, the source VM is paused only  
2740 after all state and memory contents are copied to the destination. The last of the  
2741 dirtied memory pages are copied over, and the newly formed destination VM is then  
2742 un-paused. This pause and transfer time represents the entirety of a VM downtime  
2743 during live migration, and is often at or under 100 milliseconds across commodity  
2744 Ethernet networks (given a low memory utilization).

2745 VM memory transfers are often the main performance consideration for overhead  
2746 during live migration. This is not only due to the potentially large amount of memory  
2747 to be sent across the network, but the veracity at which the memory is changed. This  
2748 is defined directly by the amount of main memory allocated (or in use) by the source  
2749 VM. However, as a VM's memory is sent, the VM is still running and therefore memory  
2750 pages can be dirtied, creating the need for any written page to be retransmitted.  
2751 Given a small network and a memory bound processes running within a VM, this  
2752 can be an infinitely long process of page dirtying. Many live migration strategies  
2753 provide an iterative timeout mechanism to avoid this infinite state, but this will lead

2754 to increased downtime during migration.

2755 The copying of memory pages for live migration can take multiple implemen-  
2756 tations. Three common options are summarized below:

- 2757     ● Pre-copy Migration: All memory pages are transmitted to the destination before  
2758         the VM is paused. The hypervisor will note and track all dirtied memory pages,  
2759         and retransmit those pages in iterative rounds. The rounds end when either no  
2760         dirtied pages exist or a max iteration count has been reached. The VM state  
2761         is then transmitted and resumed on the destination. This method was the first  
2762         live migration technique used in [230].
- 2763     ● Post-copy Migration: The Vm state is immediately paused and sent to the  
2764         destination hypervisor, and immediately resumed. If the new destination VM  
2765         generates a page fault, the VM is paused, and faulted pages are transmitted  
2766         across the network on demand from the source and the VM resumed. This  
2767         methodology is proposed for use in the Xen hypervisor by Hines et al [231].
- 2768     ● Hybrid-copy Migration: Provides a compromise solution to memory paging.  
2769         First, single copy of the VM memory pages, or a subset of known-necessary  
2770         memory pages are sent to the destination. Then the source VM is paused,  
2771         its VM state sent to the destination, and resumed on the destination. Known  
2772         dirtied source pages, or missing pages are then copied to the destination upon a  
2773         triggered page fault utilizing the same mechanism as post-copy migration. An  
2774         example of hybrid migration can be found via Lu et al [232].

2775 While pre-copy migration is the traditional and most used live migration tech-  
2776 nique, there are opportunities about to implement other migration techniques to  
2777 advance the mobility of distributed computing in high performance virtual clusters  
2778 with virtualization.

**2779 7.2.1 RDMA-enabled VM Migration**

2780 Currently, most live migration in production environments occur over TCP/IP con-  
2781 nections, largely due to the prevalence of commodity Ethernet connections within  
2782 cloud infrastructure. However even if RDMA-capable interconnects are available in  
2783 such infrastructure as described with InfiniBand in Chapter 6, live migration still  
2784 usually occurs over TCP/IP. For the case of InfiniBand, this is via IP over InfiniBand  
2785 (IPoIB) [233], which can be an inefficient use of the interconnect bandwidth and add  
2786 extra latency [234].

2787 The time it takes to migrate the memory contents from a source to destination  
2788 VM can be significantly decreased by using an RDMA based mechanisms. Huang et  
2789 al first provided a proposed pre-copy method for RDMA-based migration using the  
2790 Xen hypervisor [235]. Specifically, they found a 80% decrease in migration time with  
2791 RDMAwrite operations. This speedup is largely due to the removal of overhead nec-  
2792 essary for processing TCP/IP communications, largely in CPU utilization for copying  
2793 buffers, packet processing, and the included context switch overhead when competing  
2794 for resources in a CPU-bound application (which are common in HPC environments).

2795 The live migration algorithm proposed in [235] uses the standard pre-copy mecha-  
2796 nism that sends the entire memory contents across the network as RDMA operations,  
2797 then iteratively copies dirtied pages before switching the running states. As discussed  
2798 in the previous section, a post-copy migration strategy may have benefits for quick  
2799 VM migration in high performance virtual clusters, or even for VM cloning as de-  
2800 scribed later in Section 7.3. Furthermore, efforts in Chapter 3 have found that the  
2801 Xen hypervisor is not best suited for HPC workloads. As such, there is a need to  
2802 redefine the use of RDMA for VM migration using a hybrid post-copy mechanism in  
2803 a high performance hypervisor.

- 2804     ● Transfer initial CPU state, registers
- 2805     ● Start the page table pages translation process: MFN to PFN and use copy-base
- 2806        approach
- 2807     ● Concurrently, allocate remote memory on destination VM.
- 2808     ● Set up other machine state settings in destination
- 2809     ● Start destination VM, pause source VM.
- 2810     ● Initiate RDMAwrite of entire memory contents from source to destination
- 2811     ● As page faults occur in destination VM, catch faults and perform RDMAread
- 2812        requesting pages.

2813        Currently efforts are underway to provide post-copy live migration in KVM/QEMU,  
2814        using the `migrate_set_capability x-postcopy-ram` on mode within KVM [236].  
2815        This method uses the Linux *userfaultfd* kernel mechanisms from a kernel 4.3 or newer.  
2816        At the start, all memory blocks are registered as *userfaultfd* so all faults cause the  
2817        running thread to pause. In kernel space, the missing page is requested from the  
2818        sender, which prioritized over other pages being sent and is returned and mapped  
2819        to the destination guest memory space and the thread or process is unpause. This  
2820        mechanism operates asynchronously, so that multiple outstanding page faults will not  
2821        stop other executables within the VM. As such, the `xpostcopy-ram` KVM extention  
2822        has been identified as an ideal place to implement an RDMA based implementation.

2823        To provide RDMA functionality, the InfiniBand interconnect could first be used  
2824        as a proof-of-concept. This choice is due to InfiniBand's rise in popularity, increased  
2825        prevelance in virtualized systems with SR-IOV as noted in Chapter 6, and RMDA  
2826        functionality. However, other interconnect options exist that may be better suited

2827 for enhanced functionality and performance, such as Intel's new Omnipath [237], or  
2828 an Ethernet solution such as RoCE [238], to name a few alternatives. While RDMA  
2829 can be managed through the kernel level, it is often used in user-level APIs, such  
2830 as MPI, or in the case of InfiniBand, ibVerbs. However, it may be best to select a  
2831 interconnect-agnostic middleware for RDMA, such as Photon [239].

2832 RDMA semantics can be used to either read or write contents of remote mem-  
2833 ory, in this case VM guest memory pages. However before such operations can take  
2834 place, the target side of the operation must register the remote memory buffers and  
2835 send the remote key to the initiator, effectively providing the DMA addressing to  
2836 be used. Beyond the increased bandwidth and decreased latency benefits provided  
2837 by an advanced interconnect with InfiniBand, RDMA also allows VM memory to be  
2838 sent without involving the OS. This is due to InfiniBand's zero-copy, kernel bypass  
2839 mechanisms, and utilizing asynchronous operations. This keeps the CPU load down,  
2840 and focused on the computation at hand rather than the I/O transfer, as it often has  
2841 to when utilizing a TCP/IP stack.

2842 Selecting the proper RDMA based communication operations to use can make a  
2843 notable difference in the overall performance of the migration. Some RDMA oper-  
2844 ations are largely a one-sided movement, which can have performance impact and  
2845 should be carefully considered in using for post-copy live migration. The RDMAread  
2846 operation requires more effort at the destination host, while RDMA write operation  
2847 puts burden on the source host. One way to determine which method is best is by  
2848 evaluating both source and sink CPU loads. However this method likely will not  
2849 be as obvious when we consider VMs will be not be running in an over-subscribed  
2850 virtualized environment, but rather a highly optimized one.

2851 When the destination VM page faults, it is necessary to retrieve the page with  
2852 as little latency as possible, as the running thread is paused. Using the method

2853 developed in [236], the *userfaultfd* will trigger an RDMAread to retreive the missing  
2854 page. RDMAread requires no interaction from the source VM, as RDMAread is one  
2855 sided and the memory buffers have been passed in the setup phase. To further enable  
2856 quick return time for the missing page, enabling polling on the source host may further  
2857 reduce latency, however verification of this will be necessary. When the RDMAread  
2858 operation is finished, the running thread will resume and execution will continue.

2859 **7.2.2 Moving the Compute to the Data**

2860 While pre-copy migration is dominant in the live migration techniques of nearly every  
2861 mainstream hypervisor today, it is proposed that for HPC applications, post-copy  
2862 migration could provide some key new advances. One particular use case would be to  
2863 send a lightweight VM to directly act on a large or set of large datasets, and return a  
2864 slimmed down result set. This would reduce the requirement of transmitting the data  
2865 across a network entirely, and potentially speed up data access latency drastically, as  
2866 on result information in the form of memory pages and VM state are transmitted.

2867 With post-copy migration, one could move the computation at hand close to a  
2868 data source in significantly less time than full pre-copy live migration. This data  
2869 source, and lightweight VM sink, could potentially be something similar to a Burst  
2870 Buffer system [240, 241] or a classic HPC I/O node with a distribute filesystem such  
2871 as Lustre or GPFS [242]. This data source could even potentially be a remote sci-  
2872 entific instrument completely separate from the HPC infrastructure itself, especially  
2873 if the network at hand is capable of RoCE [238] or iWARP [243]. A VM would ini-  
2874 tiate post-copy live migration, transmitting only the necessary cpu state, registers,  
2875 and non-paged memory rather than the full VM memory state. Once migrated, the  
2876 VM could connect to a (now local) I/O or storage device, accessing data fast and  
2877 perofrming the necessary calculations. The VM could potentially even forgo the copy

2878 of the majority of its memory. During this time, only the necessary memory pages  
2879 required to complete the immediate calculation would generate a fault and trigger  
2880 their transmission from the source. The VM could even return the result (rather  
2881 than a very large dataset) to the original source VM.

2882 This post-copy live migration technique for remote data computation avoids the  
2883 cost of spawning a whole new job and/or process with associated running parameters,  
2884 as well as the extremely high cost of a full VM live migration using the pre-copy  
2885 method. However, one potential downside of post-copy live migration would be the  
2886 non-deterministic runtime, as it would be unknown how much remote memory paging  
2887 would be required. This could lead to more time spent with the destination VM in  
2888 a paused state awaiting remote memory pages, rather than if the entire VM memory  
2889 contents were transmitted completely. Careful analysis of memory usage, or a hybrid  
2890 copy method based on predetermined memory sections could help overcome this issue,  
2891 but require a more advanced migration architecture.

## 2892 7.3 Fast VM Cloning

2893 In many distributed system environments, concurrency is achieved through the use  
2894 of homogeneous compute nodes that handle the bulk of the computational load in  
2895 parallel. This can take many forms, including master/slave configurations, or even a  
2896 traditional HPC cluster with identical compute nodes, often used to support Single  
2897 Process Multiple Data (SPMD) computational models [244]. With high performance  
2898 virtual clusters, there is a need to efficiently deploy and manage near identical virtual  
2899 machines for distributed computation.

2900 In Snowflock [245,246], the notion of VM cloning is given. Specifically, Lagar et al.  
2901 define the notion of VM Fork, where VMs are treated similar to a fork system call for

2902 processes. This process is conceptually similar to VM migration, with the exception  
2903 that the source VM is not destroyed after the migration. Starting with a master  
2904 VM, an impromptu cluster can be created across a network using the Xen hypervisor.  
2905 Snowflock specifically uses Multicast to linearly scale out VM creation to many hosts,  
2906 only coping a minimal state and then remotely coping memory pages when requested.  
2907 This fetched memory on-demand is similar in principal to the post-copy live migration  
2908 technique described in the previous section. This is further augmented with blocktap-  
2909 based virtual disks with Copy-on-Write (CoW) functionally, delivering CoW slices for  
2910 each child VM.

2911 While Snowflock provides an excellent framework for VM cloning, it is not suitable  
2912 for the current implementation. First, it uses Xen, which in previous research de-  
2913 scribed in Chapter 3 has found to have limited performance for HPC workloads [169].  
2914 Second, SnowFlock is designed for Ethernet and IP based networks, which have signif-  
2915 icantly higher latency and lower bandwidth when compared to InfiniBand solutions.  
2916 Developing analogous mechanisms with a high performance hypervisor such as KVM  
2917 or Palacios [52] to use RDMA for VM memory paging. RDMA could even be utilized  
2918 in conjunction with CoW disc blocks.

2919 Potential VM Cloning RDMA mechanism:

- 2920     ● Prepare parent VM state, including registers and info.
- 2921     ● Prepare CoW disk images.
- 2922     ● Create large buffers for all VM memory on child hosts.
- 2923     ● Send vmstate via RDMAwrite or IB Send to N child nodes, where N is the clone  
2924       size.
- 2925     ● Resume/start child VMs in tandem.

2926     ● Parent VM set up RDMA multicast (unreliable connection) to send memory  
2927        pages to all childs efficiently

2928     ● Child clone VM joins RDMA mlticast.

2929     ● If child page faults, perform RDMAread operation for specific page.

2930     This method allows for efficient cloning of VMs based on a running parent VM.

2931     First, the VM state and images are copied to all child nodes to recieve the VM,  
2932     and blank memory is allocated. Then each child VM is started, and page faults are  
2933     handled via RDMA read. This will result in an initial slowdown, but as pages are  
2934     received the child VMs will start to run. The rest of the memory is eventually sent  
2935     via multicast to all VMs simultaniously. As multicast is unreliable, delivery failure  
2936     will just trigger a page fault and subsequent page sending via RDMA. It allows for  
2937     direct page fault handling, while allowing child VMs to start and run immediately.  
2938     As RDMA multicast mechanisms are relatively questionable, more investigation is  
2939     needed to evaluate the feasibility of this situation.

2940     It is expected for post-copy VM cloning to also work most efficiently if used in  
2941     conjunction with guest VMs backed with hugepages. Transferring memory in 2M  
2942     chunks will more effectively utilize network bandwdith by eliminating send/receive  
2943     overhead. This also will hide the latency found in SR-IOV enabled interconnects  
2944     for small messages. Furthermore, it will reduce the overhead of page fault handling  
2945     mechanisms, as less overall pages will fault and be transferred. While hugepages are  
2946     expected to improve VM cloning efficency, emperical testing will still be necessary to  
2947     properly evaluate their viability.

2948     While a VM fork mechanism leveraging post-copy live migration in KVM will  
2949     quickly spool up cloned VMs, the eventual memory transfer will eventually fail to scale  
2950     past the network's capacity. This could happen if hundreds or thousands of child clone

2951 VMs are being started simultaneously, as is likely in large scale deployments. As such,  
2952 a higherarchical distribution may be necessary. One possible method for this would be  
2953 a two-stage cloning mechanism, where child VMs are cloned one to each cabinet, and  
2954 the entire memory contents copied using the pre-copy migration mechanism. From  
2955 there, further cloning occurs to deploy many cloned VMs to individual nodes within  
2956 the cabinet. Organization of mid-tier cloned VMs would likely be determined based  
2957 on RDMA fabric configurations within cabinets, as this is a network-bound process.

2958 **7.4 Virtual Cluster Scheduling**

2959 Historically, cloud infrastructure has taken a simplistic approach when it comes to VM  
2960 and workload scheduling. Often, a simple round-robin or greedy [247] scheduling  
2961 algorithm is appropriated. With round robin scheduling, a simple list of host ma-  
2962 chines are used and iterated over as VM requests are made. This essentially scatters  
2963 the VMs without regard to their locality, and oversubscription becomes commonplace  
2964 regardless of the number of requested VMs or their interconnection. A greedy algo-  
2965 rithm can help keep spacial locality, but focuses specifically on over-subscription as  
2966 well to help consolidate VM allocations. While this oversubscription aspect is very  
2967 good for public cloud providers such as Amazon EC2, it becomes counterproductive  
2968 for high performance virtual clusters.

2969 With high performance virtual clusters, a number of VM instances are needed  
2970 that can gain near-native performance. Furthermore, these VMs will be running  
2971 tightly coupled applications, and as such need to be allocated in a way in which  
2972 communication latency is minimized and bandwidth is maximized for all VMs. This  
2973 will help insure the entire virtual clusters can perform optimally. However, given  
2974 off-the-shelf private cloud providers or, worse still, public cloud infrasturcture, these

2975 requirements are at best opaque to the user, and at worst extremely suboptimal.

2976 One way in which high performance virtual clusters can operate efficiently is  
 2977 by defining specific instance types, or *flavors* within OpenStack, that define what  
 2978 resources a VM has allocated. Given previous research on the NUMA effects of  
 2979 VMs [248], these specialized flavors should be defined to fit within a NUMA socket.  
 2980 Furthermore, CPU pinning should be used to specifically keep the VM within the  
 2981 NUMA socket itself. Using libvirt, a common API utilized in many cloud infrastruc-  
 2982 ture deployments (including OpenStack), we can specify CPU pinning directly in the  
 2983 XML config.

```
2984 <cputune>
2985   <vcpu pin vcpu="0" cpuset="0"/>
2986   <vcpu pin vcpu="1" cpuset="1"/>
2987   <vcpu pin vcpu="2" cpuset="4"/>
2988   <vcpu pin vcpu="3" cpuset="5"/>
2989   <vcpu pin vcpu="4" cpuset="6"/>
2990   <emulator pin cpuset="2"/>
2991 </cputune>
```

<input type="checkbox"/>	Key	Value	Actions
<input type="checkbox"/>	pci_passthrough:labels	[ "gpu", "infiniband" ]	<button>Edit</button> <button>More ▾</button>

Displaying 1 item

Figure 7.2 Adding extra specs to a VM flavor in OpenStack

2992 With OpenStack, this NUMA configuration can be executed through the KVM

2993 Nova plugin and a scheduling filter, which defines how to place VMs effectively. Fur-  
2994 themore, the instance flavor can also specify the addition of

2995 `instance_type_extra_specs`

2996 within Nova, whereby specialized hardware such as GPUs and InfiniBand intercon-  
2997 nects (as described in Chapters 5 and 6) can be passed through to the VMs directly.

2998 Once defined, the same specialized high performance flavors in OpenStack simply have  
2999 to add the labels (such as 'gpu' or 'infiniband') to the flavor to gain the requested  
3000 hardware, as illustrated in Figure 7.2 with OpenStack Horizon's UI interface. The  
3001 implementation put forth in the OpenStack Havana build has since been upgraded by  
3002 Intel with Sr-IOV support and additional scheduling filter additions, and is available  
3003 in the latest OpenStack releases [249].

```
3004 pci_passthrough_devices=[{"label":"gpu","address":"0000:08:00.0"},  
3005  {"label":"infiniband","address":"0000:21:00.0"}]  
3006 instance_type_extra_specs={'pci_passthrough:labels': '["gpu"]'}  
3007 instance_type_extra_specs={'pci_passthrough:labels': '["infiniband"]'}
```

3008 To provide a space for high performance virtual clusters, we need a scheduling  
3009 mechanism within a cloud infrastructure to support the effective and proper plate-  
3010 ment beyond controlling for NUMA characteristics. While there are many effective  
3011 scheduling algorithms for workload placement within an environment, a Proximity  
3012 Scheduler [36], as defined in OpenStack, may work well. Specifically, one could either  
3013 define or compute the underlying cloud infrastructure network topology. This could  
3014 be as simple as a YAML file that defines an underlying InfiniBand interconnect 2-1 Fat  
3015 Tree topology, or a more complex solution that utilizes network performance metrics  
3016 to measure bandwidth and latency between disjoint nodes to build a weighted prox-  
3017 mity graph. As it is possible that such network parameters could change due to

3018 other usage or to reconfiguration, a method of periodically monitoring and updating  
3019 this proximity network using measurement tools [250] may also help keep an effec-  
3020 tive proximity metric between hosts. With a metric, one can then apply a proximity  
3021 scheduler to handle high performance virtual cluster allocation requests effectively  
3022 and in a way that will be far more optimal than a simple round robin scheduling  
3023 mechanism. The OpenStack private cloud IaaS framework is proposed for this ef-  
3024 fort, however, further development is needed to bring such a scheudling mechanism  
3025 to fruition.

3026 The use of service level agreements (SLAs) within cloud infrastructure allocation is  
3027 a well studied aspect [251]. It is also possible that certain SLAs could be incorporated  
3028 into a private cloud infrasturcture such as OpenStack to simultaniously gaurantee  
3029 performance for virtual clusters with the above defined methods, while concurrently  
3030 offering "classical" VM workload scheduling for HTC, big data, or other cloud usage  
3031 models. This would provide the same user experience yet support diverse workloads  
3032 and performance expectations as defined by a given SLA

## 3033 7.5 Chapter Summary

3034 In summary, we expect the combination of transparant huge pages, an RDMA-capable  
3035 interconnect multiplexed in hardware for use by both guest and hosts, a high perfor-  
3036 mant hypervisor, and post-copy migration and cloning mechanisms to enable a novel  
3037 architecture for high perofrmance virtual clusters. These mechanisms, if implemented  
3038 and combined with an efficient OS, such as Kitten [252], could help enable a new run-  
3039 time system to support extreme scale distirbuted memory computations with in-situ  
3040 analysis [253,254] for large scale scientific systems.

3041 **Chapter 8**

3042 **Conclusion**

3043 With the advent of virtualization and the availability of virtual machines through  
3044 cloud infrastructure, a paradigm shift in distributed systems has occurred. Many  
3045 services and applications once deployed on workstations, private servers, personal  
3046 computers, and even some supercomputers have migrated to a cloud infrastructure.  
3047 The reasons for this change are vast, however these reasons are not enough to support  
3048 all computational challenges within such a virtualized infrastructure.

3049 The use of tightly coupled, distributed memory applications common in High  
3050 Performance Computing communities, has seen a number of problems and compli-  
3051 cations when deployed in virtualized infrastructure. While the reasons for this can  
3052 be vast, many challenges stem from the performance impact and overhead associated  
3053 with virtualization, along with a lack of hardware necessary to support such concur-  
3054 rent environments. This dissertation looks to evaluate virtualization for supporting  
3055 mid-tier scientific HPC applications and provide potential solutions to these issues.

3056 From the beginning, this dissertation proposes the advent of high performance  
3057 virtual clusters to support a wide array of scientific computation, including mid-tier  
3058 HPC applications, as well as a framework for building such an environment. This

3059 framework aims at identifying virtualization overhead and finding solutions and best  
3060 practices with performant hypervisors, providing support for advanced accelerators  
3061 and interconnects to enable new class of applications, and evaluating potential meth-  
3062 ods using benchmarks and real-world applications.

3063 Chapter 2 studied the related research necessary for defining not only the context  
3064 for virtualization and cloud computing, but also virtual clusters, and their history  
3065 through supercomputing. Chapter 3 looked to study the applicability of various hy-  
3066 pervisors for supporting common HPC workloads through the use of benchmarks from  
3067 a single-node aspect. This found challenges and some solutions to these workloads,  
3068 and identified missing gaps that exist.

3069 Chapter 4 started the investigation of the utility of GPUs to support mid-tier  
3070 scientific applications using the Xen hypervisor. This chapter provided a proof-of-  
3071 concept that with proper configuration and utilizing the latest in hardware support,  
3072 GPU passthrough was possible and a viable model for supporting CUDA-enabled  
3073 applications, a fast-growing application set. Chapter 5 provides an in-depth com-  
3074 parison of multiple hypervisors using the SHOC GPU benchmark suite, as well as  
3075 GPU-enabled HPC applications. Here we discover our KVM implementation per-  
3076 forms at near-native speeds and allows for effective GPU utilization.

3077 Chapter 6 takes the lessons learned with KVM in GPU passthrough and adds in  
3078 SR-IOV InfiniBand support, a critical tool for supporting tightly coupled distributed  
3079 memory applications, to build a small virtual cluster. This environment supports  
3080 two class-leading Molecular Dynamics simulations, LAMMPS and HOOMD-blue, and  
3081 shows how both applications can not only perform at near-native speeds, but also  
3082 leverage the latest HPC technologies such as GPUDirect for efficient GPU-to-GPU  
3083 communication. This framework is also enveloped in an OpenStack environment.

3084 Chapter 7 is an introspective look at other advancements that can be made in

3085 virtualization to support high performance virtual clusters. Specifically, this chap-  
3086 ter details the utility of virtual clusters backed with hugepages, added support for  
3087 specialized live migration techniques leveraging high speed RDMA-capable intercon-  
3088 nects, VM cloning for fast deployment of virtual clusters themselves, and scheduling  
3089 considerations for integration of high performance virtual clusters in OpenStack.

3090 *TODO: Answer the research question, can we support HPC with virtual clusters?*  
3091 *How could this help with big data convergence?*

## 3092 8.1 Impact

3093 TBD

## 3094      Bibliography

- 3095      [1] B. Alexander, “Web 2.0: A New Wave of Innovation for Teaching and Learn-  
3096      ing?” *Learning*, vol. 41, no. 2, pp. 32–44, 2006.
- 3097      [2] R. Buyya, C. S. Yeo, and S. Venugopal, “Market-oriented cloud computing:  
3098      Vision, hype, and reality for delivering it services as computing utilities,” in  
3099      *Proceedings of the 10th IEEE International Conference on High Performance*  
3100      *Computing and Communications (HPCC-08, IEEE CS Press, Los Alamitos,*  
3101      *CA, USA)*, 2008, pp. 5–13.
- 3102      [3] I. Foster, Y. Zhao, I. Raicu, and S. Lu, “Cloud Computing and Grid Comput-  
3103      ing 360-Degree Compared,” in *Grid Computing Environments Workshop, 2008.*  
3104      *GCE’08*, 2008, pp. 1–10.
- 3105      [4] M. Armbrust, A. Fox, R. Griffith, A. D. Joseph, R. Katz, A. Konwinski,  
3106      G. Lee, D. Patterson, A. Rabkin, I. Stoica, and M. Zaharia, “Above  
3107      the clouds: A berkeley view of cloud computing,” University of  
3108      California at Berkeley, Tech. Rep., February 2009. [Online]. Available:  
3109      <http://berkeleyclouds.blogspot.com/2009/02/above-clouds-released.html>
- 3110      [5] T. Kuhn, *The structure of scientific revolutions*. University of Chicago press  
3111      Chicago, 1970.

- 3112 [6] T. Sterling, *Beowulf cluster computing with Linux*. The MIT Press, 2001.
- 3113 [7] T. Hoare and R. Milner, “Grand challenges for computing research,” *The Computer Journal*, vol. 48, no. 1, pp. 49–52, 2005.
- 3115 [8] J. Dongarra, H. Meuer, and E. Strohmaier, “Top 500 supercomputers,” website,  
3116 November 2013.
- 3117 [9] P. Buncic, C. A. Sanchez, J. Blomer, L. Franco, A. Harutyunian, P. Mato, and  
3118 Y. Yao, “Cernvm—a virtual software appliance for lhc applications,” in *Journal  
of Physics: Conference Series*, vol. 219, no. 4. IOP Publishing, 2010, p. 042003.
- 3120 [10] L.-W. Wang, “A survey of codes and algorithms used in nersc material science  
3121 allocations,” *Lawrence Berkeley National Laboratory*, 2006.
- 3122 [11] K. Menon, K. Anala, S. G. Trupti, and N. Sood, “Cloud computing: Applications  
3123 in biological research and future prospects,” in *Cloud Computing Technologies, Applications  
3124 and Management (ICCCTAM), 2012 International Conference on*. IEEE, 2012, pp. 102–107.
- 3126 [12] Q. He, S. Zhou, B. Kobler, D. Duffy, and T. McGlynn, “Case study  
3127 for running hpc applications in public clouds,” in *Proceedings of the 19th  
3128 ACM International Symposium on High Performance Distributed Computing*,  
3129 ser. HPDC ’10. New York, NY, USA: ACM, 2010, pp. 395–401. [Online].  
3130 Available: <http://doi.acm.org/10.1145/1851476.1851535>
- 3131 [13] J. Towns, T. Cockerill, M. Dahan, I. Foster, K. Gaither, A. Grimshaw, V. Hazle-  
3132 wood, S. Lathrop, D. Lifka, G. D. Peterson *et al.*, “Xsede: accelerating scientific  
3133 discovery,” *Computing in Science & Engineering*, vol. 16, no. 5, pp. 62–74, 2014.

- 3134 [14] K. Antypas, “Nersc-6 workload analysis and benchmark selection process,”  
3135           *Lawrence Berkeley National Laboratory*, 2008.
- 3136 [15] J. S. Vetter, R. Glassbrook, J. Dongarra, K. Schwan, B. Loftis, S. McNally,  
3137           J. Meredith, J. Rogers, P. Roth, K. Spafford *et al.*, “Keeneland: Bringing het-  
3138           erogeneous gpu computing to the computational science community,” *Comput-  
3139           ing in Science and Engineering*, vol. 13, no. 5, pp. 90–95, 2011.
- 3140 [16] P. Pacheco, *Parallel Programming with MPI*, ser. ISBN. Morgan Kaufmann,  
3141           October 1996, no. 978-1-55860-339-4.
- 3142 [17] D. Agrawal, S. Das, and A. El Abbadi, “Big data and cloud computing: cur-  
3143           rent state and future opportunities,” in *Proceedings of the 14th International  
3144           Conference on Extending Database Technology*. ACM, 2011, pp. 530–533.
- 3145 [18] J. Dean and S. Ghemawat, “Mapreduce: simplified data processing on large  
3146           clusters,” *Communications of the ACM*, vol. 51, no. 1, pp. 107–113, 2008.
- 3147 [19] S. Kamburugamuve, G. Fox, D. Leake, and J. Qiu, “Survey of apache big data  
3148           stack,” Ph.D. dissertation, Ph. D. Qualifying Exam, Dept. Inf. Comput., Indi-  
3149           ana Univ., Bloomington, IN, 2013.
- 3150 [20] M. Chen, S. Mao, and Y. Liu, “Big data: a survey,” *Mobile Networks and  
3151           Applications*, vol. 19, no. 2, pp. 171–209, 2014.
- 3152 [21] D. A. Reed and J. Dongarra, “Exascale computing and big data,” *Communi-  
3153           cations of the ACM*, vol. 58, no. 7, pp. 56–68, 2015.
- 3154 [22] S. Jha, J. Qiu, A. Luckow, P. K. Mantha, and G. C. Fox, “A tale of two  
3155           data-intensive paradigms: Applications, abstractions, and architectures,” in  
3156           *Proceedings of the 3rd International Congress on Big Data*, 2014.

- 3157 [23] J. Qiu, S. Jha, A. Luckow, and G. C. Fox, “Towards hpc-abds: An initial high-  
3158 performance big data stack,” *Building Robust Big Data Ecosystem ISO/IEC  
3159 JTC 1 Study Group on Big Data*, pp. 18–21, 2014.
- 3160 [24] M. W. ur Rahman, N. S. Islam, X. Lu, J. Jose, H. Subramoni, H. Wang, and  
3161 D. K. D. Panda, “High-performance rdma-based design of hadoop mapreduce  
3162 over infiniband,” in *Parallel and Distributed Processing Symposium Workshops  
3163 PhD Forum (IPDPSW), 2013 IEEE 27th International*, May 2013, pp. 1908–  
3164 1917.
- 3165 [25] S. Ekanayake, S. Kamburugamuve, and G. Fox, “Spidal: High performance data  
3166 analytics with java and mpi on large multicore hpc clusters,” in *Proceedings of  
3167 24th High Performance Computing Symposium (HPC 2016)*, 2016.
- 3168 [26] F. Tian and K. Chen, “Towards optimal resource provisioning for running  
3169 mapreduce programs in public clouds,” in *Cloud Computing (CLOUD), 2011  
3170 IEEE International Conference on*. IEEE, 2011, pp. 155–162.
- 3171 [27] I. Foster, T. Freeman, K. Keahy, D. Scheftner, B. Sotomayer, and X. Zhang,  
3172 “Virtual clusters for grid communities,” *Cluster Computing and the Grid, IEEE  
3173 International Symposium on*, vol. 0, pp. 513–520, 2006.
- 3174 [28] K. Czajkowski, I. Foster, N. Karonis, C. Kesselman, S. Martin, W. Smith, and  
3175 S. Tuecke, “A resource management architecture for metacomputing systems,”  
3176 in *Workshop on Job Scheduling Strategies for Parallel Processing*. Springer,  
3177 1998, pp. 62–82.
- 3178 [29] D. Kondo, B. Javadi, P. Malecot, F. Cappello, and D. P. Anderson, “Cost-  
3179 benefit analysis of cloud computing versus desktop grids,” in *Parallel & Dis-*

- 3180                 tributed Processing, 2009. IPDPS 2009. IEEE International Symposium on.  
3181                 IEEE, 2009, pp. 1–12.
- 3182 [30] T. Bell, B. Bompastor, S. Bukowiec, J. C. Leon, M. Denis, J. van Eldik, M. F.  
3183                 Lobo, L. F. Alvarez, D. F. Rodriguez, A. Marino *et al.*, “Scaling the cern  
3184                 openstack cloud,” in *Journal of Physics: Conference Series*, vol. 664, no. 2.  
3185                 IOP Publishing, 2015, p. 022003.
- 3186 [31] T. Gunarathne, T.-L. Wu, J. Qiu, and G. Fox, “Mapreduce in the clouds for  
3187                 science,” in *Cloud Computing Technology and Science (CloudCom), 2010 IEEE*  
3188                 *Second International Conference on*. IEEE, 2010, pp. 565–572.
- 3189 [32] S. Ostermann, A. Iosup, N. Yigitbasi, R. Prodan, T. Fahringer, and D. Epema,  
3190                 “A performance analysis of ec2 cloud computing services for scientific comput-  
3191                 ing,” in *International Conference on Cloud Computing*. Springer, 2009, pp.  
3192                 115–131.
- 3193 [33] J. Dongarra *et al.*, “The international exascale software project roadmap,”  
3194                 *International Journal of High Performance Computing Applications*, p.  
3195                 1094342010391989, 2011.
- 3196 [34] K. Bergman, S. Borkar, D. Campbell, W. Carlson, W. Dally, M. Denneau,  
3197                 P. Franzon, W. Harrod, K. Hill, J. Hiller, S. Karp, S. K. and Dean Klein,  
3198                 R. Lucas, M. Richards, A. Scarpelli, S. Scott, A. Snavely, T. Sterling, R. S.  
3199                 Williams, and K. Yelick, “Exascale computing study: Technology challenges  
3200                 in achieving exascale systems,” *Defense Advanced Research Projects Agency*  
3201                 *Information Processing Techniques Office (DARPA IPTO), Tech. Rep*, vol. 15,  
3202                 2008.

- 3203 [35] J. Shalf, S. Dosanjh, and J. Morrison, “Exascale computing technology chal-  
3204 lenges,” in *International Conference on High Performance Computing for Com-*  
3205 *putational Science*. Springer, 2010, pp. 1–25.
- 3206 [36] J. Suh, “Proximity scheduler in openstack,” Webpage. [Online]. Available:  
3207 <https://wiki.openstack.org/wiki/ProximityScheduler>
- 3208 [37] S. Kamburugamuve, S. Ekanayake, M. Pathirage, and G. Fox, “Towards high  
3209 performance processing of streaming data in large data centers,” in *HPBDC*  
3210 *2016 IEEE International Workshop on High-Performance Big Data Comput-*  
3211 *ing in conjunction with The 30th IEEE International Parallel and Distributed*  
3212 *Processing Symposium (IPDPS 2016), Chicago, Illinois USA, Friday, 2016.*
- 3213 [38] G. von Laszewski, G. C. Fox, F. Wang, A. J. Younge, A. Kulshrestha, and  
3214 G. Pike, “Design of the FutureGrid Experiment Management Framework,”  
3215 in *Proceedings of Gateway Computing Environments 2010 at Supercomputing*  
3216 *2010*. New Orleans, LA: IEEE, Nov 2010. [Online]. Available: <http://grids.ucs.indiana.edu/ptliupages/publications/vonLaszewski-10-FG-exp-GCE10.pdf>
- 3217 [39] S. Drake and O. development team, “Heat: OpenStack Orchestration,”  
3218 Webpage. [Online]. Available: <https://wiki.openstack.org/wiki/Heat>
- 3219 [40] G. Von Laszewski, F. Wang, H. Lee, H. Chen, and G. C. Fox, “Accessing  
3220 multiple clouds with cloudmesh,” in *Proceedings of the 2014 ACM international*  
3221 *workshop on Software-defined ecosystems*. ACM, 2014, pp. 21–28.
- 3222 [41] “The magellan project,” Webpage. [Online]. Available: <http://magellan.alcf.anl.gov/>
- 3223 [42] K. Yelick, S. Coghlan, B. Draney, and R. S. Canon, “The Magellan Report on  
3224 Cloud Computing for Science,” U.S. Department of Energy Office of Science
- 3225
- 3226

- 3227        Office of Advanced Scientific Computing Research (ASCR), Tech. Rep., Dec.  
3228        2011.
- 3229        [43] K. Jackson, L. Ramakrishnan, K. Muriki, S. Canon, S. Cholia, J. Shalf,  
3230           H. Wasserman, and N. Wright, “Performance Analysis of High Performance  
3231           Computing Applications on the Amazon Web Services Cloud,” in *2nd IEEE*  
3232           *International Conference on Cloud Computing Technology and Science*. IEEE,  
3233           2010, pp. 159–168.
- 3234        [44] D. Merkel, “Docker: lightweight linux containers for consistent development  
3235           and deployment,” *Linux Journal*, vol. 2014, no. 239, p. 2, 2014.
- 3236        [45] D. M. Jacobsen and R. S. Canon, “Contain this, unleashing docker for hpc,”  
3237           *Proceedings of the Cray User Group*, 2015.
- 3238        [46] M. G. Xavier, M. V. Neves, F. D. Rossi, T. C. Ferreto, T. Lange, and C. A.  
3239           De Rose, “Performance evaluation of container-based virtualization for high  
3240           performance computing environments,” in *2013 21st Euromicro International*  
3241           *Conference on Parallel, Distributed, and Network-Based Processing*. IEEE,  
3242           2013, pp. 233–240.
- 3243        [47] K. Hwang, J. Dongarra, and G. C. Fox, *Distributed and cloud computing: from*  
3244           *parallel processing to the internet of things*. Morgan Kaufmann, 2013.
- 3245        [48] P. Barham, B. Dragovic, K. Fraser, S. Hand, T. L. Harris, A. Ho, R. Neuge-  
3246           bauer, I. Pratt, and A. Warfield, “Xen and the art of virtualization,” in *Pro-*  
3247           *ceedings of the 19th ACM Symposium on Operating Systems Principles*, New  
3248           York, U. S. A., Oct. 2003, pp. 164–177.
- 3249        [49] “Vmware esx server,” [Online], <http://www.vmware.com/de/products/vi/esx/>.

- 3250 [50] J. Duato, A. J. Pena, F. Silla, J. C. Fernández, R. Mayo, and E. S. Quintana-  
3251 Orti, “Enabling CUDA acceleration within virtual machines using rCUDA,” in  
3252 *High Performance Computing (HiPC), 2011 18th International Conference on*.  
3253 IEEE, 2011, pp. 1–10.
- 3254 [51] L. Shi, H. Chen, J. Sun, and K. Li, “vCUDA: GPU-accelerated high-  
3255 performance computing in virtual machines,” *Computers, IEEE Transactions*  
3256 on, vol. 61, no. 6, pp. 804–816, 2012.
- 3257 [52] J. Lange, K. Pedretti, T. Hudson, P. Dinda, Z. Cui, L. Xia, P. Bridges, A. Gocke,  
3258 S. Jaconette, M. Levenhagen *et al.*, “Palacios and kitten: New high performance  
3259 operating systems for scalable virtualized and native supercomputing,” in *Par-*  
3260 *allel & Distributed Processing (IPDPS), 2010 IEEE International Symposium*  
3261 *on*. IEEE, 2010, pp. 1–12.
- 3262 [53] A. Kivity, Y. Kamay, D. Laor, U. Lublin, and A. Liguori, “kvm: the Linux  
3263 virtual machine monitor,” in *Proceedings of the Linux Symposium*, vol. 1, 2007,  
3264 pp. 225–230.
- 3265 [54] I. Foster, C. Kesselman, and S. Tuecke, “The Anatomy of the Grid: Enabling  
3266 Scalable Virtual Organizations,” *Intl. J. Supercomputer Applications*, vol. 15,  
3267 no. 3, 2001.
- 3268 [55] I. Foster, C. Kesselman *et al.*, “The Physiology of the Grid: An Open Grid  
3269 Services Architecture for Distributed Systems Integration,” Argonne National  
3270 Laboratory, Chicago, Tech. Rep., Jan. 2002.
- 3271 [56] D. DiNucci, “Fragmented future,” *AllBusiness-Champions of Small Business*,  
3272 1999. [Online]. Available: <http://www.cdinucci.com/Darcy2/articles/Print/Printarticle7.html>

- 3274 [57] A. Arkin, S. Askary, S. Fordin, W. Jekeli, K. Kawaguchi, D. Orchard,  
3275 S. Pogliani, K. Riemer, S. Struble, P. Takacs-Nagy, I. Trickovic, and S. Zimek,  
3276 “Web Services Choreography Interface,” Jun. 2002, version 1.0. [Online].  
3277 Available: <http://wwws.sun.com/software/xml/developers/wsci/index.html>
- 3278 [58] D. Krafzig, K. Banke, and D. Slama, *Enterprise SOA: Service-Oriented Archi-*  
3279 *tecture Best Practices (The Coad Series)*. Prentice Hall PTR Upper Saddle  
3280 River, NJ, USA, 2004.
- 3281 [59] M. Armbrust, A. Fox, R. Griffith, A. Joseph, R. Katz, A. Konwinski, G. Lee,  
3282 D. Patterson, A. Rabkin, I. Stoica *et al.*, “A view of cloud computing,” *Com-*  
3283 *munications of the ACM*, vol. 53, no. 4, pp. 50–58, 2010.
- 3284 [60] L. Wang, G. von Laszewski, A. J. Younge, X. He, M. Kunze, and J. Tao,  
3285 “Cloud Computing: a Perspective Study,” *New Generation Computing*, vol. 28,  
3286 pp. 63–69, Mar 2010. [Online]. Available: <http://cyberaide.googlecode.com/svn/trunk/papers/10-cloudcomputing-NGC/vonLaszewski-10-NGC.pdf>
- 3288 [61] G. von Laszewski, A. J. Younge, X. He, K. Mahinthakumar, and  
3289 L. Wang, “Experiment and Workflow Management Using Cyberaide  
3290 Shell,” in *Proceedings of the 4th International Workshop on Workflow*  
3291 *Systems in e-Science (WSES 09) with 9th IEEE/ACM International*  
3292 *Symposium on Cluster Computing and the Grid (CCGrid 09)*. IEEE,  
3293 May 2009. [Online]. Available: <http://cyberaide.googlecode.com/svn/trunk/papers/09-gridshell-ccgrid/vonLaszewski-ccgrid09-final.pdf>
- 3295 [62] G. von Laszewski, F. Wang, A. J. Younge, X. He, Z. Guo, and M. Pierce,  
3296 “Cyberaide JavaScript: A JavaScript Commodity Grid Kit,” in *Proceedings*  
3297 *of the Grid Computing Environments 2007 at Supercomputing 2008*. Austin,

- 3298 TX: IEEE, Nov 2008. [Online]. Available: <http://cyberaide.googlecode.com/svn/trunk/papers/08-javascript/vonLaszewski-08-javascript.pdf>
- 3299
- 3300 [63] G. von Laszewski, J. Diaz, F. Wang, and G. C. Fox, “Comparison of multiple  
3301 cloud frameworks,” in *Cloud Computing (CLOUD), 2012 IEEE 5th Interna-*  
3302 *tional Conference on*, June 2012, pp. 734–741.
- 3303 [64] B. P. Rimal, E. Choi, and I. Lumb, “A taxonomy and survey of cloud computing  
3304 systems,” *INC, IMS and IDC*, pp. 44–51, 2009.
- 3305 [65] B. Sotomayor, R. S. Montero, I. M. Llorente, and I. Foster, “Virtual infrastruc-  
3306 ture management in private and hybrid clouds,” *IEEE Internet Computing*,  
3307 vol. 13, no. 5, pp. 14–22, 2009.
- 3308 [66] S. A. Baset, “Cloud slas: present and future,” *ACM SIGOPS Operating Systems  
3309 Review*, vol. 46, no. 2, pp. 57–66, 2012.
- 3310 [67] “Amazon Elastic Compute Cloud.” [Online]. Available: <http://aws.amazon.com/ec2/>
- 3311
- 3312 [68] S. Krishnan and J. L. U. Gonzalez, “Google compute engine,” in *Building Your  
3313 Next Big Thing with Google Cloud Platform*. Springer, 2015, pp. 53–81.
- 3314 [69] “Nimbus Project,” <http://www.nimbusproject.org>. Last access Mar. 2011.
- 3315 [70] K. Keahey, I. Foster, T. Freeman, and X. Zhang, “Virtual workspaces: Achiev-  
3316 ing quality of service and quality of life in the grid,” *Scientific Programming  
3317 Journal*, vol. 13, no. 4, pp. 265–276, 2005.
- 3318 [71] “Amazon web services s3 rest api,” <http://awsdocs.s3.amazonaws.com/S3/latest/s3-api.pdf>. Last Access Mar. 2011.
- 3319

- 3320 [72] “Jets3t project,” <http://bitbucket.org/jmurty/jets3t/wiki/Home>. Last Access  
3321 Mar. 2011.
- 3322 [73] “Boto project,” <http://code.google.com/p/boto>. Last Access Mar. 2011.
- 3323 [74] “S3tools project,” <http://s3tools.org/s3cmd>. Last Access Mar. 2011.
- 3324 [75] D. Nurmi, R. Wolski, C. Grzegorczyk, G. Obertelli, S. Soman, L. Youseff, and  
3325 D. Zagorodnov, “The Eucalyptus Open-source Cloud-computing System,” *Pro-*  
3326 *ceedings of Cloud Computing and Its Applications*, 2008.
- 3327 [76] “Eucalyptus open-source cloud computing infrastructure,” an Overview, Euca-  
3328 lyptus Systems, Inc. 2009.
- 3329 [77] “Eucalyptus,” <http://www.eucalyptus.com>, Last Access Mar. 2011.
- 3330 [78] I. H. Shaon, “Eucalyptus and its components,” Webpage.  
3331 [Online]. Available: <https://mdshaonimran.wordpress.com/2011/11/26/eucalyptus-and-its-components>
- 3332
- 3333 [79] “Openstack. cloud software,” <http://openstack.org>, Last Access Mar. 2011.
- 3334 [80] O. Sefraoui, M. Aissaoui, and M. Eleuldj, “Openstack: toward an open-source  
3335 solution for cloud computing,” *International Journal of Computer Applications*,  
3336 vol. 55, no. 3, 2012.
- 3337 [81] “Opennebula project,” <http://www.opennebula.org>. Last access Mar. 2011.
- 3338 [82] I. M. Llorente, R. Moreno-Vozmediano, and R. S. Montero, “Cloud computing  
3339 for on-demand grid resource provisioning,” *Advances in Parallel Computing*,  
3340 vol. 18, no. 5, pp. 177–191, 2009.

- 3341 [83] B. Sotomayor, R. S. Montero, I. M. Llorente, and I. Foster, “Capacity leasing in cloud systems using the opennebula engine,” *Cloud Computing and its*  
3342       *Applications*, 2008.
- 3343
- 3344 [84] “Libvirt API webpage,” [http:// libvirt.org](http://libvirt.org). Last access Mar. 2011.
- 3345 [85] “Elastichosts webpage,” <http://www.elastichosts.com>. Last Access Mar. 2011.
- 3346 [86] S. Adabala, V. Chadha, P. Chawla, R. Figueiredo, J. Fortes, I. Krsul, A. Mat-  
3347 sunaga, M. Tsugawa, J. Zhang, M. Zhao *et al.*, “From virtualized resources  
3348 to virtual computing grids: the in-vigo system,” *Future Generation Computer  
3349 Systems*, vol. 21, no. 6, pp. 896–909, 2005.
- 3350 [87] J. Chase, D. Irwin, L. Grit, J. Moore, and S. Sprenkle, “Dynamic virtual clus-  
3351       ters in a grid site manager,” in *12th IEEE International Symposium on High  
3352       Performance Distributed Computing, 2003. Proceedings*, 2003, pp. 90–100.
- 3353 [88] I. VMware, “VMware vCloud Air,” Webpage. [Online]. Available: <http://vcloud.vmware.com>
- 3354
- 3355 [89] CERN, “LHC Computing Grid Project,” Web Page, Dec. 2003. [Online].  
3356       Available: <http://lcg.web.cern.ch/LCG/>
- 3357 [90] J. Luo, A. Song, Y. Zhu, X. Wang, T. Ma, Z. Wu, Y. Xu, and L. Ge, “Grid sup-  
3358       porting platform for AMS data processing,” *Lecture notes in computer science*,  
3359       vol. 3759, p. 276, 2005.
- 3360 [91] “CMS,” Web Page. [Online]. Available: <http://cms.cern.ch/>
- 3361 [92] J. Diaz, A. J. Younge, G. von Laszewski, F. Wang, and G. C. Fox, “Grappling  
3362       Cloud Infrastructure Services with a Generic Image Repository,” in *Proceedings  
3363       of Cloud Computing and Its Applications (CCA 2011)*, Argonne, IL, Mar 2011.

- 3364 [93] I. Foster, "The anatomy of the grid: Enabling scalable virtual organizations,"  
3365 in *Euro-Par 2001 Parallel Processing: 7th International Euro-Par Conference*.  
3366 Springer, 2001, pp. 1–5. [Online]. Available: [www.globus.org/alliance/publications/papers/anatomy.pdf](http://www.globus.org/alliance/publications/papers/anatomy.pdf)
- 3368 [94] K. Keahey and T. Freeman, "Contextualization: Providing one-click virtual  
3369 clusters," in *eScience, 2008. eScience'08. IEEE Fourth International Conference on*. IEEE, 2008, pp. 301–308.
- 3371 [95] R. L. Moore, C. Baru, D. Baxter, G. C. Fox, A. Majumdar, P. Papadopoulos,  
3372 W. Pfeiffer, R. S. Sinkovits, S. Strande, M. Tatineni, R. P. Wagner, N. Wilkins-  
3373 Diehr, and M. L. Norman, "Gateways to discovery: Cyberinfrastructure  
3374 for the long tail of science," in *Proceedings of the 2014 Annual Conference on Extreme  
3375 Science and Engineering Discovery Environment*, ser. XSEDE '14. New York, NY, USA: ACM, 2014, pp. 39:1–39:8. [Online]. Available:  
3376 <http://doi.acm.org/10.1145/2616498.2616540>
- 3378 [96] D. Bernstein, "Containers and cloud: From lxc to docker to kubernetes," *IEEE  
3379 Cloud Computing*, vol. 1, no. 3, pp. 81–84, 2014.
- 3380 [97] F. Berman, "From TeraGrid to knowledge grid," *Communications of the ACM*,  
3381 vol. 44, no. 11, pp. 27–28, 2001.
- 3382 [98] C. Catlett, "The philosophy of TeraGrid: building an open, extensible, dis-  
3383 tributed TeraScale facility," in *2nd IEEE/ACM International Symposium on Cluster  
3384 Computing and the Grid, 2002*, 2002, pp. 8–8.
- 3385 [99] H. H. Goldstine and A. Goldstine, "The electronic numerical integrator and  
3386 computer (eniac)," *Mathematical Tables and Other Aids to Computation*, vol. 2,  
3387 no. 15, pp. 97–110, 1946.

- 3388 [100] J. E. Thornton, "Design of a computer - the control data 6600," 1970.
- 3389 [101] R. M. Russell, "The cray-1 computer system," *Communications of the ACM*,  
3390 vol. 21, no. 1, pp. 63–72, 1978.
- 3391 [102] C. L. Seitz, "The cosmic cube," *Communications of the ACM*, vol. 28, no. 1,  
3392 pp. 22–33, 1985.
- 3393 [103] G. C. Fox, S. W. Otto, and A. J. Hey, "Matrix algorithms on a hypercube i:  
3394 Matrix multiplication," *Parallel computing*, vol. 4, no. 1, pp. 17–31, 1987.
- 3395 [104] W. D. Hillis, *The connection machine*. MIT press, 1989.
- 3396 [105] X. Zhang, C. Yang, F. Liu, Y. Liu, and Y. Lu, "Optimizing and scaling hpcg  
3397 on tianhe-2: early experience," in *International Conference on Algorithms and*  
3398 *Architectures for Parallel Processing*. Springer, 2014, pp. 28–41.
- 3399 [106] A. Geist, *PVM: Parallel virtual machine: a users' guide and tutorial for net-*  
3400 *worked parallel computing*. MIT press, 1994.
- 3401 [107] B. Carpenter, V. Getov, G. Judd, A. Skjellum, and G. C. Fox, "Mpj: Mpi-like  
3402 message passing for java," 2000.
- 3403 [108] E. Gabriel, G. E. Fagg, G. Bosilca, T. Angskun, J. J. Dongarra, J. M. Squyres,  
3404 V. Sahay, P. Kambadur, B. Barrett, A. Lumsdaine *et al.*, "Open mpi: Goals,  
3405 concept, and design of a next generation mpi implementation," in *European*  
3406 *Parallel Virtual Machine/Message Passing Interface Users Group Meeting*.  
3407 Springer, 2004, pp. 97–104.
- 3408 [109] W. Gropp, E. Lusk, and A. Skjellum, *Using MPI: portable parallel programming*  
3409 *with the message-passing interface*. MIT press, 1999, vol. 1.

- 3410 [110] M. J. Koop, T. Jones, and D. K. Panda, “Mvapich-aptus: Scalable high-  
3411 performance multi-transport mpi over infiniband,” in *IEEE International Sym-  
3412 posium on Parallel and Distributed Processing, 2008. IPDPS 2008.* IEEE,  
3413 2008, pp. 1–12.
- 3414 [111] R. Rabenseifner, G. Hager, and G. Jost, “Hybrid mpi/openmp parallel program-  
3415 ming on clusters of multi-core smp nodes,” in *2009 17th Euromicro international  
3416 conference on parallel, distributed and network-based processing.* IEEE, 2009,  
3417 pp. 427–436.
- 3418 [112] D. A. Jacobsen, J. C. Thibault, and I. Senocak, “An mpi-cuda implementation  
3419 for massively parallel incompressible flow computations on multi-gpu clusters,”  
3420 in *48th AIAA aerospace sciences meeting and exhibit*, vol. 16, 2010, p. 2.
- 3421 [113] A. S. Bland, J. Wells, O. E. Messer, O. Hernandez, and J. Rogers, “Titan:  
3422 Early experience with the cray xk6 at oak ridge national laboratory,” *Cray  
3423 User Group*, 2012.
- 3424 [114] J. J. Dongarra, P. Luszczek, and A. Petitet, “The linpack benchmark: past,  
3425 present and future,” *Concurrency and Computation: practice and experience*,  
3426 vol. 15, no. 9, pp. 803–820, 2003.
- 3427 [115] R. C. Murphy, K. B. Wheeler, B. W. Barrett, and J. A. Ang, “Introducing the  
3428 graph 500,” *Cray Users Group (CUG)*, 2010.
- 3429 [116] M. A. Heroux and J. Dongarra, “Toward a new metric for ranking high perfor-  
3430 mance computing systems,” *Sandia Report, SAND2013-4744*, vol. 312, 2013.
- 3431 [117] R. Brightwell, R. Oldfield, A. B. Maccabe, and D. E. Bernholdt, “Hobbes:  
3432 Composition and virtualization as the foundations of an extreme-scale os/r,”

- 3433       in *Proceedings of the 3rd International Workshop on Runtime and Operating*  
3434       *Systems for Supercomputers.* ACM, 2013, p. 2.
- 3435       [118] S. Perarnau, R. Gupta, and P. Beckman, “Argo: An exascale operating system  
3436       and runtime,” in *Poster Proceedings from IEEE/ACM Supercomputing 2015,*  
3437       2015.
- 3438       [119] T. Sterling, D. Kogler, M. Anderson, and M. Brodowicz, “SLOWER: A perfor-  
3439       mance model for Exascale computing,” *Supercomputing Frontiers and Innova-*  
3440       *tions*, vol. 1, pp. 42–57, Sep 2014.
- 3441       [120] E. Ciurana, *Developing with Google App Engine.* Springer, 2009.
- 3442       [121] D. Chappell, “Introducing windows azure,” Microsoft, Inc, Tech. Rep., 2009.
- 3443       [122] K. Keahey, R. Figueiredo, J. Fortes, T. Freeman, and M. Tsugawa, “Science  
3444       clouds: Early experiences in cloud computing for scientific applications,” *Cloud*  
3445       *Computing and Applications*, vol. 2008, 2008.
- 3446       [123] R. Creasy, “The origin of the VM/370 time-sharing system,” *IBM Journal of*  
3447       *Research and Development*, vol. 25, no. 5, pp. 483–490, 1981.
- 3448       [124] “Amazon elastic compute cloud,” [Online], <http://aws.amazon.com/ec2/>.
- 3449       [125] K. Keahey, I. Foster, T. Freeman, X. Zhang, and D. Galron, “Virtual  
3450       Workspaces in the Grid,” *Lecture Notes in Computer Science*, vol. 3648,  
3451       pp. 421–431, 2005. [Online]. Available: [http://workspace.globus.org/papers/VW\\_EuroPar05.pdf](http://workspace.globus.org/papers/VW_EuroPar05.pdf)
- 3453       [126] J. Fontan, T. Vazquez, L. Gonzalez, R. S. Montero, and I. M. Llorente, “Open-  
3454       NEbula: The Open Source Virtual Machine Manager for Cluster Computing,”

- 3455       in *Open Source Grid and Cluster Software Conference*, San Francisco, CA, USA,  
3456       May 2008.
- 3457 [127] S. Adabala, V. Chadha, P. Chawla, R. Figueiredo, J. Fortes, I. Krsul, A. Mat-  
3458       sunaga, M. Tsugawa, J. Zhang, M. Zhao, L. Zhu, and X. Zhu, “From virtualized  
3459       resources to virtual computing Grids: the In-VIGO system,” *Future Generation  
3460       Comp. Syst.*, vol. 21, no. 6, pp. 896–909, 2005.
- 3461 [128] Rackspace, “Openstack,” WebPage, Jan 2011. [Online]. Available: <http://www.openstack.org/>
- 3462
- 3463 [129] P. Padala, X. Zhu, Z. Wang, S. Singhal, and K. Shin, “Performance evaluation  
3464       of virtualization technologies for server consolidation,” HP Laboratories, Tech.  
3465       Rep., 2007.
- 3466 [130] J. Watson, “Virtualbox: bits and bytes masquerading as machines,” *Linux  
3467       Journal*, vol. 2008, no. 166, p. 1, 2008.
- 3468 [131] D. Leinenbach and T. Santen, “Verifying the Microsoft Hyper-V Hypervisor  
3469       with VCC,” *FM 2009: Formal Methods*, pp. 806–809, 2009.
- 3470 [132] I. Parallels, “An introduction to os virtualization and paral-  
3471       lels virtuozzo containers,” Parallels, Inc, Tech. Rep., 2010. [On-  
3472       line]. Available: [http://www.parallels.com/r/pdf/wp/pvc/Parallels\\_Virtuozzo\\_Containers\\_WP\\_an\\_introduction\\_to\\_os\\_EN.pdf](http://www.parallels.com/r/pdf/wp/pvc/Parallels_Virtuozzo_Containers_WP_an_introduction_to_os_EN.pdf)
- 3473
- 3474 [133] D. Bartholomew, “Qemu: a multihost, multitarget emulator,” *Linux Journal*,  
3475       vol. 2006, no. 145, p. 3, 2006.
- 3476 [134] J. N. Matthews, W. Hu, M. Hapuarachchi, T. Deshane, D. Dimatos, G. Hamil-  
3477       ton, M. McCabe, and J. Owens, “Quantifying the performance isolation prop-

- 3478        erties of virtualization systems,” in *Proceedings of the 2007 workshop on Ex-*
- 3479        perimental computer science, ser. ExpCS ’07. New York, NY, USA: ACM,
- 3480        2007.
- 3481 [135] Oracle, “Performance evaluation of oracle vm server virtualization software,”
- 3482        Oracle, Whitepaper, 2008. [Online]. Available: <http://www.oracle.com/us/technologies/virtualization/oraclevm/026997.pdf>
- 3483
- 3484 [136] K. Adams and O. Agesen, “A comparison of software and hardware techniques
- 3485        for x86 virtualization,” in *Proceedings of the 12th international conference on Architectural support for programming languages and operating systems*. ACM,
- 3486
- 3487        2006, pp. 2–13, vMware.
- 3488 [137] Y. Koh, R. Knauerhase, P. Brett, M. Bowman, Z. Wen, and C. Pu, “An analy-
- 3489        sis of performance interference effects in virtual environments,” in *Performance Analysis of Systems & Software, 2007. ISPASS 2007. IEEE International Symposium on*. IEEE, 2007, pp. 200–209.
- 3490
- 3491
- 3492 [138] S. Rixner, “Network virtualization: Breaking the performance barrier,” *Queue*,
- 3493        vol. 6, no. 1, p. 36, 2008.
- 3494 [139] S. Nanda and T. Chiueh, “A survey of virtualization technologies,” Tech. Rep.,
- 3495
- 3496 [140] A. Ranadive, M. Kesavan, A. Gavrilovska, and K. Schwan, “Performance im-
- 3497        plications of virtualizing multicore cluster machines,” in *Proceedings of the 2nd workshop on System-level virtualization for high performance computing*. ACM,
- 3498
- 3499        2008, pp. 1–8.
- 3500 [141] R. Harper and K. Rister, “Kvm limits arbitrary or architectural?” IBM Linux
- 3501        Technology Center, Jun 2009.

- 3502 [142] M. Bolte, M. Sievers, G. Birkenheuer, O. Niehorster, and A. Brinkmann, “Non-  
3503 intrusive virtualization management using libvirt,” in *Design, Automation Test  
3504 in Europe Conference Exhibition (DATE), 2010*, 2010, pp. 574 –579.
- 3505 [143] “FutureGrid,” Web Page, 2009. [Online]. Available: <http://www.futuregrid.org>
- 3506 [144] J. J. Dujmovic and I. Dujmovic, “Evolution and evaluation of spec bench-  
3507 marks,” *SIGMETRICS Perform. Eval. Rev.*, vol. 26, no. 3, pp. 2–9, 1998.
- 3508 [145] P. Luszczek, D. Bailey, J. Dongarra, J. Kepner, R. Lucas, R. Rabenseifner,  
3509 and D. Takahashi, “The HPC Challenge (HPCC) benchmark suite,” in *SC06  
3510 Conference Tutorial*. Citeseer, 2006.
- 3511 [146] J. Dongarra and P. Luszczek, “Reducing the time to tune parallel dense linear  
3512 algebra routines with partial execution and performance modelling,” University  
3513 of Tennessee Computer Science Technical Report, Tech. Rep., 2010.
- 3514 [147] K. Dixit, “The SPEC benchmarks,” *Parallel Computing*, vol. 17, no. 10-11, pp.  
3515 1195–1209, 1991.
- 3516 [148] SPEC, “Standard performance evaluation corporation,” Webpage, Jan 2011.  
3517 [Online]. Available: <http://www.spec.org/>
- 3518 [149] R. Henschel and A. J. Younge, “First quarter 2011 spec omp results,” Webpage,  
3519 Mar 2011. [Online]. Available: <http://www.spec.org/omp/results/res2011q1/>
- 3520 [150] A. S. Tanenbaum and M. Van Steen, *Distributed systems*. Prentice Hall, 2002,  
3521 vol. 2.
- 3522 [151] Amazon, “Elastic Compute Cloud.” [Online]. Available: <http://aws.amazon.com/ec2/>

- 3524 [152] J. Dean and S. Ghemawat, "MapReduce: simplified data processing on large  
3525 clusters," *Commun. ACM*, vol. 51, no. 1, pp. 107–113, 2008.
- 3526 [153] "Windows azure platform," [Online], <http://www.microsoft.com/azure/>.
- 3527 [154] G. von Laszewski, G. C. Fox, F. Wang, A. J. Younge, A. Kulshrestha, G. G.  
3528 Pike, W. Smith, J. Vockler, R. J. Figueiredo, J. Fortes *et al.*, "Design of the  
3529 futuregrid experiment management framework," in *Gateway Computing Envi-  
ronments Workshop (GCE), 2010.* IEEE, 2010, pp. 1–10.
- 3530
- 3531 [155] OpenStack, "Openstack compute administration manual," 2013.
- 3532 [156] D. Nurmi, R. Wolski, C. Grzegorczyk, G. Obertelli, S. Soman, L. Youseff, and  
3533 D. Zagorodnov, "The Eucalyptus open-source cloud-computing system," in  
3534 *Proceedings of Cloud Computing and Its Applications*, October 2008. [Online].  
3535 Available: <http://eucalyptus.cs.ucsb.edu/wiki/Presentations>
- 3536 [157] C. Nvidia, "Programming guide," 2008.
- 3537 [158] J. E. Stone, D. Gohara, and G. Shi, "OpenCL: A parallel programming standard  
3538 for heterogeneous computing systems," *Computing in science & engineering*,  
3539 vol. 12, no. 3, p. 66, 2010.
- 3540 [159] S. Che, M. Boyer, J. Meng, D. Tarjan, J. W. Sheaffer, and K. Skadron,  
3541 "A performance study of general-purpose applications on graphics processors  
3542 using CUDA," *Journal of Parallel and Distributed Computing*, vol. 68,  
3543 no. 10, pp. 1370 – 1380, 2008, k-means implementation on CUDA  
3544 with 72x speedup. Compares to 4-threaded CPU version with 30x  
3545 speedup. [Online]. Available: <http://www.sciencedirect.com/science/article/B6WKJ-4SVV8GS-2/2/f7a1dccceb63cbbfd25774c6628d8412>
- 3546

- 3547 [160] Z. Liu and W. Ma, "Exploiting computing power on graphics processing unit,"  
3548 vol. 2, Dec. 2008, pp. 1062–1065.
- 3549 [161] S. Hong and H. Kim, "An integrated GPU power and performance model,"  
3550 in *ACM SIGARCH Computer Architecture News*, vol. 38. ACM, 2010, pp.  
3551 280–289.
- 3552 [162] P. Kogge, K. Bergman, S. Borkar, D. Campbell, W. Carson, W. Dally, M. Den-  
3553 neau, P. Franzon, W. Harrod, K. Hill *et al.*, "Exascale computing study: Tech-  
3554 nology challenges in achieving exascale systems," 2008.
- 3555 [163] S. Crago, K. Dunn, P. Eads, L. Hochstein, D.-I. Kang, M. Kang, D. Mod-  
3556 ium, K. Singh, J. Suh, and J. P. Walters, "Heterogeneous cloud computing," in  
3557 *Proceedings of the Workshop on Parallel Programming on Accelerator Clusters*  
3558 (*PPAC*). *Cluster Computing (CLUSTER), 2011 IEEE International Confer-  
3559 ence on*. IEEE, 2011, pp. 378–385.
- 3560 [164] J. Sanders and E. Kandrot, *CUDA by example: an introduction to general-*  
3561 *purpose GPU programming*. Addison-Wesley Professional, 2010.
- 3562 [165] V. V. Kindratenko, J. J. Enos, G. Shi, M. T. Showerman, G. W. Arnold, J. E.  
3563 Stone, J. C. Phillips, and W.-m. Hwu, "GPU clusters for high-performance  
3564 computing," in *Cluster Computing and Workshops, 2009. CLUSTER'09. IEEE*  
3565 *International Conference on*. IEEE, 2009, pp. 1–8.
- 3566 [166] K. J. Barker, K. Davis, A. Hoisie, D. J. Kerbyson, M. Lang, S. Pakin, and J. C.  
3567 Sancho, "Entering the petaflop era: the architecture and performance of Road-  
3568 runner," in *Proceedings of the 2008 ACM/IEEE conference on Supercomputing*.  
3569 IEEE Press, 2008, p. 1.

- 3570 [167] S. Craven and P. Athanas, "Examining the viability of FPGA supercomputing,"  
3571       *EURASIP Journal on Embedded systems*, vol. 2007, no. 1, pp. 13–13, 2007.
- 3572 [168] G. Shainer, T. Liu, J. Layton, and J. Mora, "Scheduling strategies for HPC as a  
3573       service (HPCaaS)," in *Cluster Computing and Workshops, 2009. CLUSTER'09.*  
3574       *IEEE International Conference on.* IEEE, 2009, pp. 1–6.
- 3575 [169] A. J. Younge, R. Henschel, J. T. Brown, G. von Laszewski, J. Qiu, and G. C.  
3576       Fox, "Analysis of Virtualization Technologies for High Performance Computing  
3577       Environments," in *Proceedings of the 4th International Conference on Cloud*  
3578       *Computing (CLOUD 2011).* Washington, DC: IEEE, July 2011.
- 3579 [170] W. Wade, "How NVIDIA and Citrix are driving the future of virtualized visual  
3580       computing," [Online], <http://blogs.nvidia.com/blog/2013/05/22/synergy/>.
- 3581 [171] S. Long, "Virtual machine graphics acceleration deployment guide," VMWare,  
3582       Tech. Rep., 2013.
- 3583 [172] G. Humphreys, M. Houston, R. Ng, R. Frank, S. Ahern, P. D. Kirchner, and  
3584       J. T. Klosowski, "Chromium: a stream-processing framework for interactive  
3585       rendering on clusters," in *ACM Transactions on Graphics (TOG)*, vol. 21.  
3586       ACM, 2002, pp. 693–702.
- 3587 [173] G. Giunta, R. Montella, G. Agrillo, and G. Coviello, "A GPGPU  
3588       transparent virtualization component for high performance computing clouds,"  
3589       in *Euro-Par 2010 - Parallel Processing*, ser. Lecture Notes in Computer  
3590       Science, P. D'Ambra, M. Guarracino, and D. Talia, Eds. Springer  
3591       Berlin Heidelberg, 2010, vol. 6271, pp. 379–391. [Online]. Available:  
3592       [http://dx.doi.org/10.1007/978-3-642-15277-1\\_37](http://dx.doi.org/10.1007/978-3-642-15277-1_37)

- 3593 [174] K. Diab, M. Rafique, and M. Hefeeda, “Dynamic sharing of GPUs in cloud  
3594 systems,” in *Parallel and Distributed Processing Symposium Workshops PhD*  
3595 *Forum (IPDPSW), 2013 IEEE 27th International*, 2013, pp. 947–954.
- 3596 [175] C.-T. Yang, H.-Y. Wang, and Y.-T. Liu, “Using pci pass-through for gpu vir-  
3597 tualization with cuda,” in *Network and Parallel Computing*. Springer, 2012,  
3598 pp. 445–452.
- 3599 [176] A. Danalis, G. Marin, C. McCurdy, J. S. Meredith, P. C. Roth, K. Spafford,  
3600 V. Tippurajju, and J. S. Vetter, “The scalable heterogeneous computing (SHOC)  
3601 benchmark suite,” in *Proceedings of the 3rd Workshop on General-Purpose*  
3602 *Computation on Graphics Processing Units*. ACM, 2010, pp. 63–74.
- 3603 [177] E. R. Hawkes, R. Sankaran, J. C. Sutherland, and J. H. Chen, “Direct numerical  
3604 simulation of turbulent combustion: fundamental insights towards predictive  
3605 models,” in *Journal of Physics: Conference Series*, vol. 16. IOP Publishing,  
3606 2005, p. 65.
- 3607 [178] F. Silla, “rCUDA: share and aggregate GPUs in your cluster,” Nov. 2012, mel-  
3608 lanox Booth Presaentation.
- 3609 [179] H. Jo, J. Jeong, M. Lee, and D. H. Choi, “Exploiting GPUs in virtual machine  
3610 for biocloud,” *BioMed research international*, vol. 2013, 2013.
- 3611 [180] J. Duato, A. Pena, F. Silla, R. Mayo, and E. Quintana-Orti, “rCUDA: Re-  
3612 ducing the number of gpu-based accelerators in high performance clusters,”  
3613 in *High Performance Computing and Simulation (HPCS), 2010 International*  
3614 *Conference on*, June 2010, pp. 224–231.
- 3615 [181] B. L. Jacob and T. N. Mudge, “A look at several memory management units,  
3616 TLB-refill mechanisms, and page table organizations,” in *Proceedings of the*

- 3617        *Eighth International Conference on Architectural Support for Programming*  
3618        *Languages and Operating Systems.* ACM, 1998, pp. 295–306.
- 3619 [182] B.-A. Yassour, M. Ben-Yehuda, and O. Wasserman, “Direct device assignment  
3620        for untrusted fully-virtualized virtual machines,” IBM Research, Tech. Rep.,  
3621        2008.
- 3622 [183] M. Ben-Yehuda, J. Xenidis, M. Ostrowski, K. Rister, A. Bruemmer, and  
3623        L. Van Doorn, “The price of safety: Evaluating IOMMU performance,” in  
3624        *Ottawa Linux Symposium*, 2007.
- 3625 [184] J. Liu, “Evaluating standard-based self-virtualizing devices: A performance  
3626        study on 10 GbE NICs with SR-IOV support,” in *Parallel Distributed Processing*  
3627        (*IPDPS*), 2010 IEEE International Symposium on, April 2010, pp. 1–12.
- 3628 [185] V. Jujjuri, E. Van Hensbergen, A. Liguori, and B. Pulavarty, “VirtFS-a virtu-  
3629        alization aware file system passthrough,” in *Ottawa Linux Symposium*, 2010.
- 3630 [186] C. Yang, H. Wang, W. Ou, Y. Liu, and C. Hsu, “On implementation of GPU  
3631        virtualization using PCI pass-through,” in *4th IEEE International Conference*  
3632        *on Cloud Computing Technology and Science Proceedings (CloudCom)*, 2012.
- 3633 [187] M. Dowty and J. Sugerman, “GPU virtualization on VMware’s hosted I/O  
3634        architecture,” *ACM SIGOPS Operating Systems Review*, vol. 43, no. 3, pp. 73  
3635        – 82, 2009.
- 3636 [188] “FutureGrid,” Web page. [Online]. Available: <http://portal.futuregrid.org>
- 3637 [189] A. Danalis, G. Marin, C. McCurdy, J. S. Meredith, P. C. Roth, K. Spaf-  
3638        ford, V. Tipparaju, and J. S. Vetter, “The Scalable Heterogeneous Computing

- 3639        (SHOC) benchmark suite,” in *Proceedings of the 3rd Workshop on General-*  
3640        *Purpose Computation on Graphics Processing Units*, ser. GPGPU ’10. ACM,  
3641        2010, pp. 63–74.
- 3642        [190] “LAMMPS molecular dynamics simulator,” <http://lammps.sandia.gov/>, [On-  
3643        line; accessed Jan. 2, 2014].
- 3644        [191] A. Athanasopoulos, A. Dimou, V. Mezaris, and I. Kompatsiaris, “GPU ac-  
3645        celeration for support vector machines,” in *Proc 12th International Workshop*  
3646        *on Image Analysis for Multimedia Interactive Services (WIAMIS 2001)*, April  
3647        2011.
- 3648        [192] “LULESH shock hydrodynamics application,” <https://codesign.llnl.gov/lulesh.php>, [Online; accessed Jan. 2, 2014].
- 3650        [193] S. Plimpton, “Fast parallel algorithms for short-range molecular dynamics,”  
3651        *Journal of Computational Physics*, vol. 117, no. 1, pp. 1 – 19, 1995.
- 3652        [194] W. M. Brown, “GPU acceleration in LAMMPS,” <http://lammps.sandia.gov/workshops/Aug11/Brown/brown11.pdf>, [Online; accessed Jan. 2, 2014].
- 3654        [195] C.-C. Chang and C.-J. Lin, “LIBSVM: A library for support vector machines,”  
3655        *ACM Trans. Intell. Syst. Technol.*, vol. 2, no. 3, pp. 27:1–27:27, May 2011.
- 3656        [196] “Hydrodynamics Challenge Problem,” <https://codesign.llnl.gov/pdfs/spec-7.pdf>, [Online; accessed Jan. 2, 2014].
- 3658        [197] A. Arcangeli, “Transparent hugepage support,” <http://www.linux-kvm.org/wiki/images/9/9e/2010-forum-thp.pdf>, 2010, [Online; accessed Jan. 29, 2014].
- 3660        [198] J. Jose, M. Li, X. Lu, K. C. Kandalla, M. D. Arnold, and D. K. Panda, “SR-IOV  
3661        support for virtualization on infiniband clusters: Early experience,” in *Cluster*

- 3662        *Computing and the Grid, IEEE International Symposium on.* IEEE Computer  
3663        Society, 2013, pp. 385–392.
- 3664        [199] R. L. Moore, C. Baru, D. Baxter, G. C. Fox, A. Majumdar, P. Papadopoulos,  
3665        W. Pfeiffer, R. S. Sinkovits, S. Strande, M. Tatineni *et al.*, “Gateways to  
3666        discovery: Cyberinfrastructure for the long tail of science,” in *Proceedings of*  
3667        *the 2014 Annual Conference on Extreme Science and Engineering Discovery*  
3668        *Environment.* ACM, 2014, p. 39.
- 3669        [200] P. Luszczek, E. Meek, S. Moore, D. Terpstra, V. M. Weaver, and J. Dongarra,  
3670        “Evaluation of the hpc challenge benchmarks in virtualized environments,” in  
3671        *Proceedings of the 2011 International Conference on Parallel Processing - Volume 2*, ser. Euro-Par’11. Berlin, Heidelberg: Springer-Verlag, 2012, pp. 436–  
3672        445.  
3673
- 3674        [201] N. Huber, M. von Quast, M. Hauck, and S. Kounev, “Evaluating and model-  
3675        ing virtualization performance overhead for cloud environments.” in *CLOSER*,  
3676        2011, pp. 563–573.
- 3677        [202] J. P. Walters, A. J. Younge, D.-I. Kang, K.-T. Yao, M. Kang, S. P. Crago,  
3678        and G. C. Fox, “GPU-Passthrough Performance: A Comparison of KVM, Xen,  
3679        VMWare ESXi, and LXC for CUDA and OpenCL Applications,” in *Proceedings*  
3680        *of the 7th IEEE International Conference on Cloud Computing (CLOUD 2014).*  
3681        Anchorage, AK: IEEE, 2014.
- 3682        [203] J. Jose, M. Li, X. Lu, K. C. Kandalla, M. D. Arnold, and D. K. Panda, “Sr-iov  
3683        support for virtualization on infiniband clusters: Early experience,” in *Cluster,*  
3684        *Cloud and Grid Computing (CCGrid), 2013 13th IEEE/ACM International*  
3685        *Symposium on.* IEEE, 2013, pp. 385–392.

- 3686 [204] M. Musleh, V. Pai, J. P. Walters, A. J. Younge, and S. P. Crago, “Bridging  
3687 the Virtualization Performance Gap for HPC using SR-IOV for InfiniBand,”  
3688 in *Proceedings of the 7th IEEE International Conference on Cloud Computing*  
3689 (*CLOUD 2014*). Anchorage, AK: IEEE, 2014.
- 3690 [205] “Aws high performance computing,” <http://aws.amazon.com/hpc/>, last Access  
3691 Nov. 2014.
- 3692 [206] “Openstack flavors,” <http://docs.openstack.org/openstack-ops/content/flavors.html>, last Access Nov. 2014.
- 3694 [207] K. Pepple, *Deploying OpenStack*. O’Reilly Media, Inc., 2011.
- 3695 [208] S. Hazelhurst, “Scientific computing using virtual high-performance computing:  
3696 a case study using the amazon elastic computing cloud,” in *Proceedings of the  
3697 2008 annual research conference of the South African Institute of Computer  
3698 Scientists and Information Technologists on IT research in developing countries:  
3699 riding the wave of technology*. ACM, 2008, pp. 94–103.
- 3700 [209] E. Amazon, “Amazon elastic compute cloud (amazon ec2),” *Amazon Elastic  
3701 Compute Cloud (Amazon EC2)*, 2010.
- 3702 [210] R. Jennings, *Cloud Computing with the Windows Azure Platform*. John Wiley  
3703 & Sons, 2010.
- 3704 [211] “Google cloud platform,” <https://cloud.google.com/>, last Access Nov. 2014.
- 3705 [212] L. Ramakrishnan, R. S. Canon, K. Muriki, I. Sakrejda, and N. J. Wright,  
3706 “Evaluating interconnect and virtualization performance forhigh performance  
3707 computing,” *SIGMETRICS Perform. Eval. Rev.*, vol. 40, no. 2, pp. 55–60,  
3708 Oct. 2012. [Online]. Available: <http://doi.acm.org/10.1145/2381056.2381071>

- 3709 [213] M. Righini, "Enabling intel virtualization technology features and benefits,"  
3710 Intel Corporation, Tech. Rep., 2010.
- 3711 [214] AMD, "AMD i/o virtualization technology (IOMMU) specification," AMD Cor-  
3712 poration, Tech. Rep., 2009.
- 3713 [215] A. Limited, "Arm system memory management unit architecture specification,"  
3714 ARM Limited, Tech. Rep., 2013.
- 3715 [216] A. J. Younge, J. P. Walters, S. Crago, and G. C. Fox, "Evaluating GPU  
3716 Passthrough in Xen for High Performance Cloud Computing," in *High-*  
3717 *Performance Grid and Cloud Computing Workshop at the 28th IEEE Inter-*  
3718 *national Parallel and Distributed Processing Symposium*, IEEE. Pheonix, AZ:  
3719 IEEE, 05/2014 2014.
- 3720 [217] L. Vu, H. Sivaraman, and R. Bidarkar, "Gpu virtualization for high  
3721 performance general purpose computing on the esx hypervisor," in *Proceedings*  
3722 *of the High Performance Computing Symposium*, ser. HPC '14. San Diego,  
3723 CA, USA: Society for Computer Simulation International, 2014, pp. 2:1–2:8.  
3724 [Online]. Available: <http://dl.acm.org/citation.cfm?id=2663510.2663512>
- 3725 [218] J. Liu, "Evaluating standard-based self-virtualizing devices: A performance  
3726 study on 10 gbe nics with sr-iov support," in *Parallel Distributed Processing*  
3727 (*IPDPS*), 2010 IEEE International Symposium on, April 2010, pp. 1–12.
- 3728 [219] T. P. P. D. L. Ruivo, G. B. Altayo, G. Garzoglio, S. Timm, H. Kim, S.-Y.  
3729 Noh, and I. Raicu, "Exploring infiniband hardware virtualization in opennebula  
3730 towards efficient high-performance computing." in *CCGRID*, 2014, pp. 943–948.
- 3731 [220] "NVIDIA GPUDirect," <https://developer.nvidia.com/gpudirect>, [Online; ac-  
3732 cessed Nov. 24, 2014].

- 3733 [221] G. Shainer, A. Ayoub, P. Lui, T. Liu, M. Kagan, C. R. Trott, G. Scantlen, and  
3734 P. S. Crozier, “The development of mellanox/nvidia gpudirect over infinibanda  
3735 new model for gpu to gpu communications,” *Computer Science-Research and*  
3736 *Development*, vol. 26, no. 3-4, pp. 267–273, 2011.
- 3737 [222] “Getting Xen working for Intel(R) Xeon Phi(tm)  
3738 Coprocessor,” <https://software.intel.com/en-us/articles/getting-xen-working-for-intelr-xeon-phitm-coprocessor>, [Online; accessed  
3739 Nov. 24, 2014].
- 3741 [223] “Mellanox Neutron Plugin,” <https://wiki.openstack.org/wiki/Mellanox-Neutron>, [Online; accessed Nov. 24, 2014].
- 3743 [224] S. Plimpton, P. Crozier, and A. Thompson, “Lammps-large-scale  
3744 atomic/molecular massively parallel simulator,” *Sandia National Laboratories*, 2007.
- 3746 [225] J. Anderson, A. Keys, C. Phillips, T. Dac Nguyen, and S. Glotzer, “Hoomd-  
3747 blue, general-purpose many-body dynamics on the gpu,” in *APS Meeting Abstracts*, vol. 1, 2010, p. 18008.
- 3749 [226] K. Asanovic, R. Bodik, B. C. Catanzaro, J. J. Gebis, P. Husbands, K. Keutzer,  
3750 D. A. Patterson, W. L. Plishker, J. Shalf, S. W. Williams *et al.*, “The land-  
3751 scape of parallel computing research: A view from berkeley,” Technical Report  
3752 UCB/EECS-2006-183, EECS Department, University of California, Berkeley,  
3753 Tech. Rep., 2006.
- 3754 [227] G. Fox, G. von Laszewski, J. Diaz, K. Keahey, J. Fortes, R. Figueiredo,  
3755 S. Smallen, W. Smith, and A. Grimshaw, “Futuregrida reconfigurable testbed

- 3756 for cloud, hpc and grid computing,” *Contemporary High Performance Comput-*  
3757 *ing: From Petascale toward Exascale, Computational Science. Chapman and*  
3758 *Hall/CRC*, 2013.
- 3759 [228] K. Keahey, J. Mambretti, D. K. Panda, P. Rad, W. Smith, and  
3760 D. Stanzione, “Nsf chameleon cloud,” website, November 2014. [Online].  
3761 Available: <http://www.chameleoncloud.org/>
- 3762 [229] M. Rosenblum and T. Garfinkel, “Virtual machine monitors: Current technol-  
3763 ogy and future trends,” *Computer*, vol. 38, no. 5, pp. 39–47, 2005.
- 3764 [230] C. Clark, K. Fraser, S. Hand, J. G. Hansen, E. Jul, C. Limpach, I. Pratt, and  
3765 A. Warfield, “Live migration of virtual machines,” in *Proceedings of the 2nd*  
3766 *conference on Symposium on Networked Systems Design & Implementation-*  
3767 *Volume 2*. USENIX Association, 2005, pp. 273–286.
- 3768 [231] M. R. Hines and K. Gopalan, “Post-copy based live virtual machine migration  
3769 using adaptive pre-paging and dynamic self-ballooning,” in *Proceedings of the*  
3770 *2009 ACM SIGPLAN/SIGOPS international conference on Virtual execution*  
3771 *environments*. ACM, 2009, pp. 51–60.
- 3772 [232] P. Lu, A. Barbalace, and B. Ravindran, “Hsg-lm: Hybrid-copy speculative  
3773 guest os live migration without hypervisor,” in *Proceedings of the*  
3774 *6th International Systems and Storage Conference*, ser. SYSTOR ’13.  
3775 New York, NY, USA: ACM, 2013, pp. 2:1–2:11. [Online]. Available:  
3776 <http://doi.acm.org/10.1145/2485732.2485736>
- 3777 [233] J. Chu and V. Kashyap, “Transmission of IP over InfiniBand (IPoIB),” IETF,  
3778 Tech. Rep., 2006.

- 3779 [234] W. Yu, N. S. Rao, P. Wyckoff, and J. S. Vetter, “Performance of rdma-capable  
3780 storage protocols on wide-area network,” in *2008 3rd Petascale Data Storage  
3781 Workshop*. IEEE, 2008, pp. 1–5.
- 3782 [235] W. Huang, Q. Gao, J. Liu, and D. K. Panda, “High performance virtual machine  
3783 migration with rdma over modern interconnects,” in *2007 IEEE International  
3784 Conference on Cluster Computing*. IEEE, 2007, pp. 11–20.
- 3785 [236] D. Gilbert, “Post copy live migration,” Webpage, 2015. [Online]. Available:  
3786 <http://wiki.qemu.org/Features/PostCopyLiveMigration>
- 3787 [237] M. S. Birrittella, M. Debbage, R. Huggahalli, J. Kunz, T. Lovett, T. Rimmer,  
3788 K. D. Underwood, and R. C. Zak, “Intel omni-path architecture: Enabling  
3789 scalable, high performance fabrics,” in *2015 IEEE 23rd Annual Symposium on  
3790 High-Performance Interconnects*, Aug 2015, pp. 1–9.
- 3791 [238] M. Beck and M. Kagan, “Performance evaluation of the rdma over ethernet  
3792 (roce) standard in enterprise data centers infrastructure,” in *Proceedings of  
3793 the 3rd Workshop on Data Center-Converged and Virtual Ethernet Switching*.  
3794 International Teletraffic Congress, 2011, pp. 9–15.
- 3795 [239] E. Kissel and M. Swany, “Photon: Remote memory access middleware for high-  
3796 performance runtime systems,” in *2016 IEEE International Parallel and Dis-  
3797 tributed Processing Symposium Workshops (IPDPSW)*. IEEE, 2016, pp. 1736–  
3798 1743.
- 3799 [240] J. Lofstead, I. Jimenez, and C. Maltzahn, “Consistency and fault tolerance con-  
3800 siderations for the next iteration of the doe fast forward storage and io project,”  
3801 in *2014 43rd International Conference on Parallel Processing Workshops*, Sept  
3802 2014, pp. 61–69.

- 3803 [241] C. G. Wright Jr, “Trinity burst buffer-architecture and design,” Los Alamos  
3804 National Laboratory (LANL), Tech. Rep., 2015.
- 3805 [242] F. B. Schmuck and R. L. Haskin, “Gpfss: A shared-disk file system for large  
3806 computing clusters.” in *FAST*, vol. 2, 2002, pp. 231–244.
- 3807 [243] M. J. Rashti and A. Afsahi, “10-gigabit iwarps ethernet: comparative perfor-  
3808 mance analysis with infiniband and myrinet-10g,” in *2007 IEEE International*  
3809 *Parallel and Distributed Processing Symposium*. IEEE, 2007, pp. 1–8.
- 3810 [244] P. Duclos, F. Boeri, M. Auguin, and G. Giraudon, “Image processing on a  
3811 simd/spmd architecture: Opsila,” in *Pattern Recognition, 1988., 9th Interna-*  
3812 *tional Conference on*, Nov 1988, pp. 430–433 vol.1.
- 3813 [245] H. A. Lagar-Cavilla, J. A. Whitney, A. M. Scannell, P. Patchin, S. M. Rumble,  
3814 E. De Lara, M. Brudno, and M. Satyanarayanan, “Snowflock: rapid virtual  
3815 machine cloning for cloud computing,” in *Proceedings of the 4th ACM European*  
3816 *conference on Computer systems*. ACM, 2009, pp. 1–12.
- 3817 [246] H. A. Lagar-Cavilla, J. A. Whitney, R. Bryant, P. Patchin, M. Brudno,  
3818 E. de Lara, S. M. Rumble, M. Satyanarayanan, and A. Scannell, “Snowflock:  
3819 Virtual machine cloning as a first-class cloud primitive,” *ACM Transactions on*  
3820 *Computer Systems (TOCS)*, vol. 29, no. 1, p. 2, 2011.
- 3821 [247] A. J. Younge, G. von Laszewski, L. Wang, and G. C. Fox, “Providing a Green  
3822 Framework for Cloud Based Data Centers,” in *The Handbook of Energy-Aware*  
3823 *Green Computing*, I. Ahmad and S. Ranka, Eds. Chapman and Hall/CRC  
3824 Press, 2011, ch. 17, in press.
- 3825 [248] D. P. Berrange, “Openstack performance optimization,” in *KVM Forum 2014*,  
3826 2014.

- 3827 [249] Y. Jiang and Y. He, “Openstack pci passthrough,” Webpage, Intel, Inc, 2016.  
3828 [Online]. Available: [https://wiki.openstack.org/wiki/Pci\\_passthrough](https://wiki.openstack.org/wiki/Pci_passthrough)
- 3829 [250] A. Hanemann, J. W. Boote, E. L. Boyd, J. Durand, L. Kudarimoti, R. Lapacz,  
3830 D. M. Swany, S. Trocha, and J. Zurawski, “Perfsonar: A service oriented archi-  
3831 tecture for multi-domain network monitoring,” in *International Conference on*  
3832 *Service-Oriented Computing*. Springer, 2005, pp. 241–254.
- 3833 [251] D. Serrano, S. Bouchenak, Y. Kouki, T. Ledoux, J. Lejeune, J. Sopena,  
3834 L. Arantes, and P. Sens, “Towards qos-oriented sla guarantees for online  
3835 cloud services,” in *Cluster, Cloud and Grid Computing (CCGrid), 2013 13th*  
3836 *IEEE/ACM International Symposium on*. IEEE, 2013, pp. 50–57.
- 3837 [252] K. Pedretti, “Kitten: A lightweight operating system for ultrascaling supercom-  
3838 puters,” 2011.
- 3839 [253] J. Ahrens, D. Rogers, and B. Springmeyer, “Visualization and data analysis at  
3840 the exascale,” *National Nuclear Security Administration (NNSA) Accelerated*  
3841 *Strategic Computing (ASC) Exascale Environment Planning Process*, vol. 2,  
3842 2010.
- 3843 [254] V. Vishwanath, M. Hereld, and M. E. Papka, “Toward simulation-time data  
3844 analysis and i/o acceleration on leadership-class systems,” in *Large Data Anal-*  
3845 *ysis and Visualization (LDAV), 2011 IEEE Symposium on*. IEEE, 2011, pp.  
3846 9–14.