

Document downloaded from:

<http://hdl.handle.net/10251/103729>

This paper must be cited as:

Prades Gasulla, J.; Varghese, B.; Reaño González, C.; Silla Jiménez, F. (2017). Multi-tenant virtual GPUs for optimising performance of a financial risk application. *Journal of Parallel and Distributed Computing*. 108:28-44. doi:10.1016/j.jpdc.2016.06.002



The final publication is available at

<https://doi.org/10.1016/j.jpdc.2016.06.002>

Copyright Elsevier

Additional Information

Multi-Tenant Virtual GPUs for Optimising Performance of a Financial Risk Application

Javier Prades^a, Blesson Varghese (Corresponding Author)^b, Carlos Reaño^a,
Federico Silla^a

^a*Department of Computer Engineering,
Universitat Politècnica de València, Spain*

^b*School of Electronics, Electrical Engineering and Computer Science,
Queen's University Belfast, UK*

Abstract

Graphics Processing Units (GPUs) are becoming popular accelerators in modern High-Performance Computing (HPC) clusters. Installing GPUs on each node of the cluster is not efficient resulting in high costs and power consumption as well as underutilisation of the accelerator. The research reported in this paper is motivated towards the use of few physical GPUs by providing cluster nodes access to remote GPUs on-demand for a financial risk application. We hypothesise that sharing GPUs between several nodes, referred to as multi-tenancy, reduces the execution time and energy consumed by an application. Two data transfer modes between the CPU and the GPUs, namely concurrent and sequential, are explored. The key result from the experiments is that multi-tenancy with few physical GPUs using sequential data transfers lowers the execution time and the energy consumed, thereby improving the overall performance of the application.

Keywords: GPU virtualisation, Acceleration-as-a-Service, rCUDA, multi-tenancy, energy efficiency

Email addresses: japraga@gap.upv.es (Javier Prades), varghese@qub.ac.uk (Blesson Varghese (Corresponding Author)), carregon@gap.upv.es (Carlos Reaño), fsilla@disca.upv.es (Federico Silla)

URL: www.blessonv.com (Blesson Varghese (Corresponding Author))

1. Introduction

Hardware accelerators are achieving a prominent role in modern High-Performance Computing (HPC) clusters for making applications faster. This is evidenced by four out of top ten supercomputers listed on Top500 (<http://top500.org>) and the top ten supercomputers listed on Green500 (<http://www.green500.org>) in November 2015 have employed hardware accelerators, such as Graphics Processing Units (GPU). Incorporating GPUs in large clusters allows for heterogeneity, thus making it possible for an application to exploit the regular processor as well as the accelerator [1, 2].

Clusters can now be set up to employ a small number of GPUs by providing applications shared access to remote GPUs on-demand [3, 4]. Such a set up is feasible on a limited budget because not only are a few GPUs used to provide acceleration, but also the energy consumed is well justified since the GPUs are well utilised in the cluster [5, 6]. This is possible as a result of maturing GPU virtualisation technologies that facilitate virtual GPUs (vGPUs) in a cluster. An application can request Acceleration-as-a-Service[7] from one or many vGPUs. One vGPU can reside on a physical GPU (pGPU), referred to as *single tenancy*, but is limiting in that multiple applications cannot make use of the same pGPU since it is exclusively locked for a single application. When multiple vGPUs reside on the same pGPU, otherwise known as *multi-tenancy*, either the same application has access to a pool of vGPUs on the same pGPU or multiple applications can share the same pGPU. We hypothesise that using multi-tenancy can improve the performance of an application.

Numerous challenges arise when multiple GPUs are shared across a cluster for an application, of which three are considered in this paper. The challenges are addressed in this paper by exploring remote CUDA (rCUDA) [8], a GPU virtualisation framework, for improving the performance of a real-world case study employed in the financial industry. The application typically runs in a cluster environment, but can hugely benefit from GPU acceleration for deriving important risk metrics in real-time. The benefit of executing the application

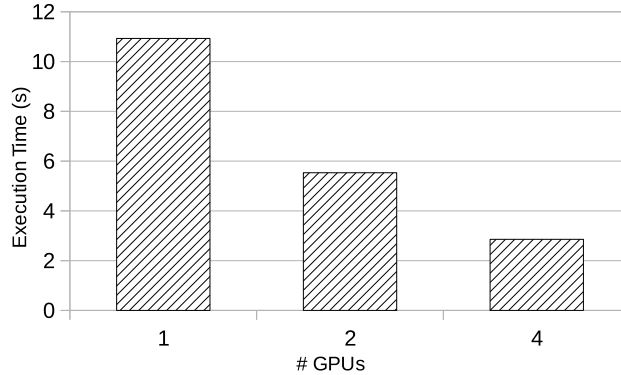


Figure 1: Execution time of the financial application on multiple local GPUs

31 on multiple physical GPUs is shown in Figure 1. We hypothesise that using a
 32 large number of vGPUs can further optimise application performance. However,
 33 the following three challenges and research questions arise, which are addressed
 34 in this paper: (i) Data will need to be transferred from the CPU to the vG-
 35 PUs for computations. However, data transfer will be restricted by bottlenecks
 36 due to limited bandwidth which affects the overall scalability of the applica-
 37 tion. Hence, “What data transfer approaches can mitigate the effect of data
 38 bottlenecks?” (ii) Multi-tenancy may degrade application performance since
 39 the underlying hardware resource is shared. This results in increased execution
 40 time and consequently higher energy consumption. Hence, “How can vGPUs be
 41 shared effectively to optimise application performance and energy consumed?”
 42 (iii) Using multi-tenancy an application can be deployed in multiple ways. For
 43 example, an application can be executed on 2 vGPUs residing on 1 pGPU or 8
 44 vGPUs residing on 1 pGPU. These possibilities significantly increase with mul-
 45 tiple pGPUs. Each deployment option consumes different amounts of energy
 46 and impacts the overall execution time. Hence, “Can performance and energy
 47 of an application be estimated in the multi-tenancy approach?”

48 To address the above challenges we propose two data transfer approaches,
 49 namely concurrent and sequential, for transferring data with the aim of mitigat-
 50 ing the effect of data bottlenecks. In the context of the financial application, the

51 sequential data transfer approach is expected to improve performance since data
52 transfers from the CPU to the GPU and GPU computations can be overlapped
53 for multiple pGPUs. The approach is further extended for overlapping the data
54 movement and computation time for multiple vGPUs on the same pGPU result-
55 ing in a further improvement in performance of the application. The key result
56 is that the financial application can be executed under two seconds for deriving
57 risk metrics in an energy efficient manner on the same hardware compared to
58 single tenancy thus confirming our initial hypothesis. Performance and energy
59 consumed by the application are modelled to determine the combination of vG-
60 PUs on a pGPU that can maximise performance and GPU utilisation and at
61 the same time minimise the energy consumed.

62 The key contributions of this research are: (i) investigating the lack of scala-
63 bility due to data transfer from CPU to the GPU in the context of the financial
64 risk application, (ii) proposing two approaches to transfer data, namely concur-
65 rent and sequential, (iii) evaluating the above data transfer approaches in the
66 context of single-tenancy for overlapping computations and data transfer of mul-
67 tiple pGPUs, (iv) developing an approach that exploits multi-tenancy for over-
68 lapping computations and data transfer of multiple virtual GPUs on the same
69 physical GPU to optimise the performance of the application, (v) evaluating the
70 performance of the application, considering execution time, GPU utilisation and
71 energy consumed by the application, and (vi) developing a mathematical model
72 to derive deployment options for the application by estimating performance and
73 energy of different combinations of virtual GPUs mapped onto physical GPUs.

74 The remainder of this paper is organised as follows. Section 2 highlights re-
75 lated work in the area of HPC solutions for GPU virtualisation and financial risk
76 applications. Section 3 briefly presents the rCUDA framework. Section 4 con-
77 siders a financial risk application for evaluating the feasibility of multi-tenancy
78 for improving performance. Section 5 presents the platform, experiments per-
79 formed and the key results obtained. Section 6 concludes this paper.

80 2. Related Work

81 High Performance Computing (HPC) solutions are exploited in the financial
82 risk industry to accelerate the underlying computations of applications. This
83 reduces overall execution times making such applications fit for real-time use.
84 Solutions range from small scale clusters [9, 10] to large supercomputers [11, 12].
85 More recently, hardware accelerators with multi-core and many-core processors
86 are employed. For example, financial risk applications are accelerated on Cell
87 BE processors [13, 14], FPGAs [15, 16] and GPUs [17, 18].

88 HPC clusters offering heterogeneous solutions by using hardware accelera-
89 tors, such as GPUs, along with processors on nodes are feasible [1, 2]. Clusters
90 can be set up to incorporate a GPU on each node. This is not an efficient solu-
91 tion for accelerating an application because of the relatively high cost of GPUs,
92 high power consumption of nodes using GPUs and the under utilisation of GPUs
93 (applications do not require acceleration of GPUs during their entire execution).
94 However, a more efficient solution would be if nodes executing an application
95 can access GPUs when required. This can be facilitated by GPU virtualisation.
96 Currently there are no solutions available for the financial risk industry to har-
97 ness the potential of GPU virtualisation. In this paper, we investigate the use
98 of virtual GPUs for a financial risk application.

99 The mechanism of GPU virtualisation allows nodes of a cluster that do not
100 own a physical GPU for accelerating computations of applications that run on
101 it to remotely access GPUs. Acceleration is obtained as a service seamlessly
102 to a requesting node without being aware of accessing remote GPUs. A single
103 application (running on a Virtual Machine (VM) or on a node of a cluster
104 without a hardware accelerator) benefits from the acceleration of a remotely
105 located single GPU or multiple GPUs to reduce execution time. The rate of
106 GPU utilisation can be increased since multiple applications can access the same
107 GPU. This in turn reduces the number of GPUs that need to be installed in
108 a cluster, and reduces the cost spent on energy consumption, cooling, physical
109 space and maintenance, usually referred to as the Total Cost of Ownership

110 (TCO). Furthermore, the source code of an application usually does not need
111 any modification to reap the benefits of virtual GPUs.

112 GPU virtualisation is usually applied at the high-level Application Program-
113 ming Interface (API) of GPUs because low level protocols used to interact with
114 accelerators are proprietary and, additionally, not publicly available. Hence,
115 APIs such as CUDA [19] or OpenCL [20] are used. In this paper we use CUDA
116 (Compute Unified Device Architecture) for an application that is used in the
117 financial risk industry.

118 There are several remote GPU virtualization frameworks supporting CUDA.
119 GridCuda [21] supports CUDA 3.2, although it is not publicly available. vCUDA [22]
120 supports the CUDA 3.2 and implements an unspecified subset of the CUDA
121 runtime API. The communication protocol between the node that executes the
122 application and the remote GPU has a considerable overhead, because of the
123 costs incurred during encoding and decoding, which results in a noticeable drop
124 of overall performance. GViM [23] is based on CUDA 1.1 and does not imple-
125 ment the entire runtime API. Furthermore, GViM is designed to be used on
126 VMs so that applications executed on them can access GPUs located in the
127 real host; GViM does not support the access of GPUs in remote nodes. gVir-
128 tuS [24] supports CUDA 2.3 and again implements only a small portion of the
129 runtime API. For example, in the case of the memory management module, it
130 implements only 17 out of the 37 available functions. Although it is intended
131 mainly to be used by VMs for accessing real GPUs located in the same node, it
132 facilitates TCP/IP communications between clients and servers, thus allowing
133 the access to GPUs located in other nodes. DS-CUDA [25] supports CUDA 4.1
134 and includes specific communication support for InfiniBand Verbs, thus reduc-
135 ing the overhead of communications between the node executing the application
136 and the node owning the GPU. However, DS-CUDA is limited in that it does not
137 allow data transfers with pinned memory and supports maximum data transfer
138 of 32 MB.

139 The rCUDA framework [8] is binary compatible with CUDA 6.5 and im-
140 plements the entire CUDA Runtime and Driver APIs (with the exception of

141 graphics functions). It provides support for the libraries included within CUDA,
 142 such as cuBLAS or cuFFT. In addition, a number of underlying interconnection
 143 technologies are supported by making use of a set of runtime-loadable, network-
 144 specific communication modules (currently TCP/IP and InfiniBand). Concur-
 145 rent virtualization services are made available to remote clients simultaneously
 146 demanding GPU acceleration by managing an independent GPU context for
 147 each client. rCUDA performs better than other publicly available GPU virtu-
 148 alisation frameworks (considered in Section 3) and is therefore chosen for this
 149 research.

150 3. rCUDA

151 The rCUDA framework, otherwise referred to as remote CUDA, is used in the
 152 research presented in this paper. As shown in Figure 2, the rCUDA framework is
 153 a client-server architecture. Numerous *Clients* executing applications that can
 154 benefit from hardware acceleration can concurrently access *Servers* that have
 155 physical GPUs on them. The client makes use of the remote GPU to accelerate
 156 part of the software code of the application, referred to as kernel, running on it.
 157 The framework transparently handles the data management and the execution
 158 management; the transfer of data between the local memory of the client, the
 159 local memory of the server and the GPU memory, and the remote execution of
 160 the kernel.

161 Figure 3 shows the hardware and software stack of the client and the rCUDA

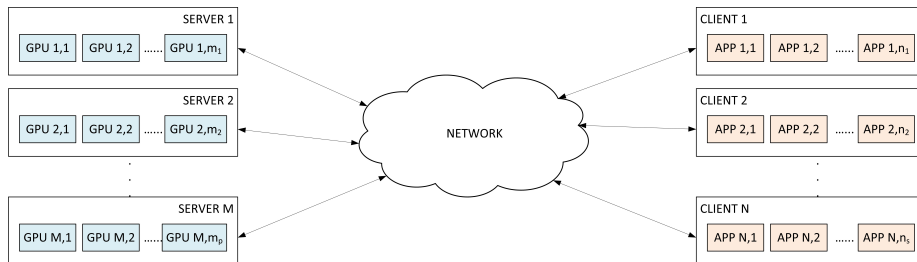


Figure 2: Distributed acceleration architecture facilitated by rCUDA

162 server. The client nodes that execute the application (shown in Figure 2), make
 163 use of the rCUDA Client Library, which is a wrapper around the CUDA Runtime
 164 and Driver APIs. The library is responsible for (i) intercepting calls made by
 165 the application to a CUDA device, (ii) processing them for forwarding the calls
 166 to the remote rCUDA server, and (iii) retrieving the results of the calls from the
 167 rCUDA server. On the other hand, each GPU server has an rCUDA daemon
 168 running on it which receives CUDA requests and executes them on the physical
 169 GPU.

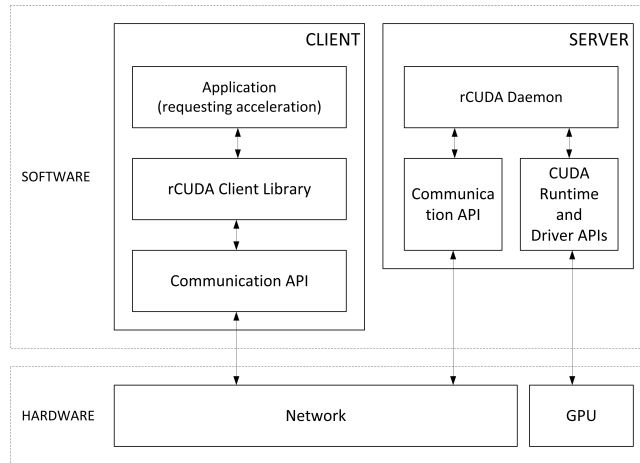


Figure 3: rCUDA client and server software/hardware stack

170 An efficient communication protocol is developed for seamless execution be-
 171 tween rCUDA clients and servers. This protocol, using either regular TCP/IP
 172 sockets or the InfiniBand Verbs API when this high performance interconnect is
 173 available in the cluster, is designed to provide lightweight support to the remote
 174 CUDA operations provided by the external accelerator. The CUDA commands
 175 intercepted by the rCUDA client wrapper are encapsulated into messages in
 176 the form of one or more packets that travel across the network towards the
 177 rCUDA server. The format of the messages depends on the specific CUDA
 178 command transported. In general, the messages have low network overheads.
 179 Every CUDA command forwarded to the remote GPU server is followed by

180 a response message, which acknowledges the success/failure of the operation
181 requested on the remote server.

182 Figure 4 shows an example of the communication between the rCUDA client
183 and the rCUDA daemon executing on the remote server. In this example, the
184 following steps occur:

185 *Step 1 - Initialise:* The client establishes connection with the remote server
186 automatically, and the request for acceleration services is intercepted and the
187 GPU kernel along with related information such as statically allocated variables
188 are sent to the server.

189 *Step 2 - Allocate Memory:* Based on the client request device memory is al-
190 located on the GPU for data that will be required by the GPU kernel. The
191 `cudaMalloc` requests are intercepted by the client and forwarded to the remote
192 server.

193 *Step 3 - Transfer Data to Device:* All data required by the kernel is transferred
194 from the host to the remote device.

195 *Step 4 - Execute Kernel:* The GPU kernel is executed remotely on the rCUDA
196 server.

197 *Step 5 - Transfer Data to Host:* After the execution of the kernel on the remote
198 server the data is transmitted back to the host.

199 *Step 6 - Release Memory:* The memory allocated on the remote device is re-
200 leased.

201 *Step 7 - Quit:* In this final step the client application stops communicating with
202 the remote server. The rCUDA daemon executing on the server stops servicing
203 the execution and releases the resources associated with the execution.

204 Figure 5 compares the performance of publicly available GPU virtualisation
205 frameworks, namely DS-CUDA, gVirtuS and rCUDA by using the `bandwidthTest`
206 benchmark from the NVIDIA CUDA Samples [26]. Our choice of selecting

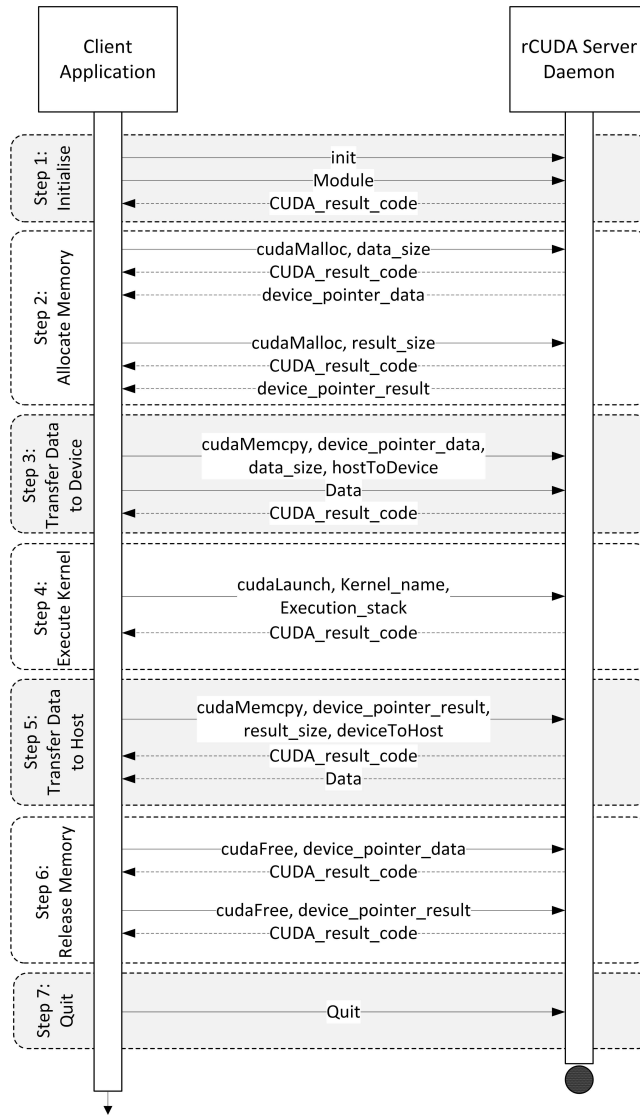


Figure 4: Communication sequence between a client and the rCUDA server daemon

207 rCUDA for this research is based on its superior performance over other frame-
 208 works as shown in the figure. The performance of CUDA 6.5 is used as the
 209 baseline reference. Bandwidth is used as a measure for comparing performance
 210 since it is a limiting factor for data transfers between host (CPU) memory and
 211 device (GPU) memory (data size can be in the order of MB) and affects the

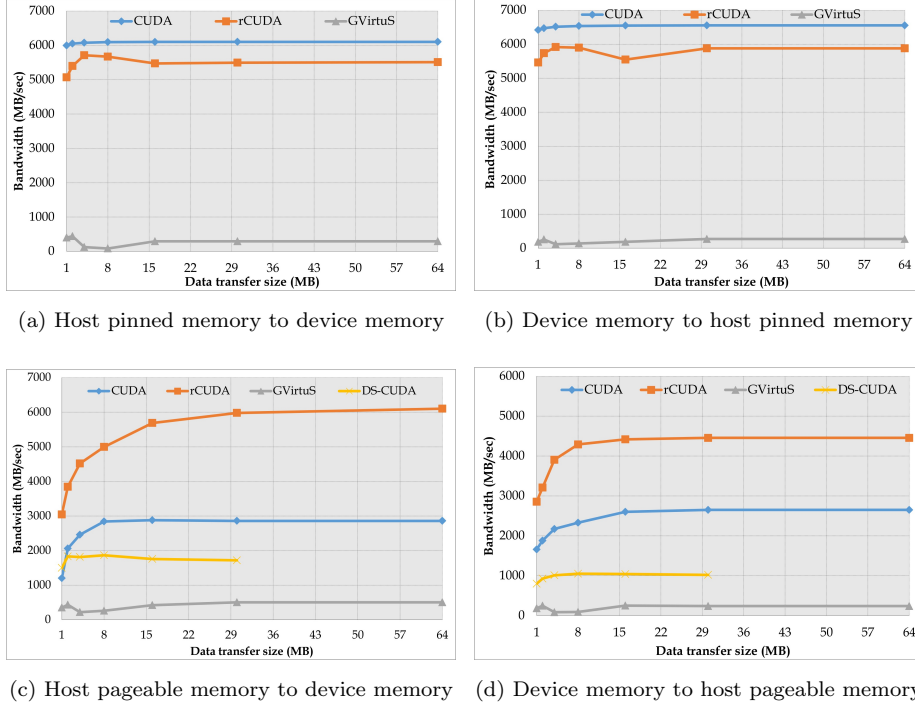


Figure 5: Comparison of bandwidth for pinned memory and pageable memory of rCUDA, DS-CUDA and gVirtuS using CUDA as a baseline reference (DS-CUDA does not support pinned memory)

212 performance of the virtualisation frameworks. Other metrics such as latency
 213 are less relevant in this context.

214 The test-bed employed for carrying out the bandwidth performance exper-
 215 iments is presented later in Section 5.1. Virtual Machine (VMs) were not em-
 216 ployed to simplify the experiments. The bandwidth test was run on a native
 217 domain and the server side of the virtualisation framework used was executed in
 218 a remote node. The InfiniBand FDR network technology was used to connect
 219 both nodes. The rCUDA and DS-CUDA frameworks made use of the InfiniBand
 220 Verbs API and gVirtuS made use of TCP/IP over InfiniBand since it cannot
 221 take advantage of the InfiniBand Verbs API.

222 The three virtualisation frameworks support different versions of CUDA

223 which had to be used for obtaining the bandwidth benchmarks. DS-CUDA
224 is compatible with CUDA 4.1, gVirtuS supports CUDA 2.3 and rCUDA sup-
225 ports CUDA 6.5. In our experience, employing different CUDA versions has
226 minimal impact on bandwidth performance and therefore no additional noise
227 was introduced by using different versions.

228 The following observations are made from Figure 4. Firstly, CUDA achieves
229 highest performance when pinned memory is used (refer Figure 5a and Fig-
230 ure 5b), achieving nearly a bandwidth of 6000 MB/s. The bandwidth is however
231 reduced for copies using pageable memory (refer Figure 5c and Figure 5d).

232 Secondly, Figure 5 shows that rCUDA outperforms DS-CUDA and gVirtuS.
233 For copies using pageable memory rCUDA even performs better than CUDA;
234 this has been previously reported, which is due to the use of an efficient pipelined
235 communication between rCUDA clients and servers based on the use of internal
236 and pre-allocated pinned memory buffers [8]. rCUDA and DS-CUDA support
237 InfiniBand Verbs API and therefore have access to large bandwidths which are
238 available on this interconnect. However, DS-CUDA has relatively poor perfor-
239 mance when compared to rCUDA. Therefore, it is assumed that both frame-
240 works manage the InfiniBand interconnect differently. DS-CUDA neither sup-
241 ports memory copies larger than 32MB nor pinned memory. The performance
242 of gVirtuS is significantly lower than the other frameworks. It may be immedi-
243 ately inferred that this is because TCP/IP is used and has a lower bandwidth in
244 comparison to InfiniBand Verbs. However, using the `iperf` tool [27], TCP/IP
245 over InfiniBand FDR provides approximately 1190 MB/s, which is a noticeably
246 larger bandwidth than the one achieved by gVirtuS. Therefore, the poor per-
247 formance of gVirtuS may be due to the inefficient handling of communication.

248 **4. Financial Risk Application**

249 A candidate application that can benefit from Acceleration-as-a-Service (AaaS)
250 in HPC clusters is investigated in this section. We present such an application
251 employed in the financial risk industry, referred to as *‘Aggregate Risk Analy-*

252 *sis'* [28] for validating the feasibility of our proposed multi-tenancy approach.
253 The analysis of financial risk is underpinned by a simulation that is computa-
254 tionally intensive. Typically, this analysis is periodically performed on a routine
255 basis on production clusters to derive important risk metrics. Such a set up is
256 sufficient when the analysis does not need to be performed outside routine.

257 Risk metrics will need to be obtained in real-time, such as in an online pricing
258 scenario, in addition to routine executions. In such settings, a number of input
259 parameters to the analysis will need to be varied to satisfy the customer. This
260 generates a large number of requests to execute the analysis multiple times based
261 on the complexity of the client's portfolio. It may not be feasible to furnish all
262 these requests generated by single or multiple clients; it will be impossible to
263 quickly obtain a large set of resources on an in-house cluster already provisioned
264 for executing other routine jobs. Here, GPUs can play an important role in
265 furnishing a large number of requests.

266 While GPUs can provide a feasible solution, employing a large number of
267 GPUs to furnish bursts of requests will be expensive. As considered in Section 1
268 virtual GPUs are pragmatic and cost effective to minimise under utilisation. In
269 this context, we leverage the acceleration offered by virtual GPUs in an HPC
270 cluster to develop a faster application fit for use in real-time settings. The
271 rCUDA framework suits such an application because minimal changes need to
272 be brought about to the production cluster and the acceleration required for the
273 analysis is obtained as a service from a remote host. The analysis has previously
274 been investigated in the context of many-core architectures [29], but we believe
275 virtual GPUs can be a better option.

276 Aggregate risk analysis is performed on a portfolio of risk held by an insurer
277 or reinsurer and provides actuaries and decision makers with millions of alter-
278 nate views of catastrophic events, such as earthquakes, that can occur and the
279 order in which they can occur in a year. To obtain millions of alternate views,
280 millions of trials are simulated with each trial comprising a set of possible future
281 earthquake events and the probable loss for each trial is estimated.

282 *4.1. Input and Output Data*

283 Three data tables are required for the analysis, which are as follows:

284 i. *Year Event Table*, which is a database of pre-simulated occurrences of
 285 events from a catalogue of stochastic events that is denoted as *YET*. Each
 286 record in a *YET* called a ‘trial’, denoted as T_i , represents a possible sequence
 287 of event occurrences for any given year. The sequence of events is defined by
 288 an ordered set of tuples containing the ID of an event and the time-stamp of its
 289 occurrence in that trial $T_i = \{(E_{i,1}, t_{i,1}), \dots, (E_{i,k}, t_{i,k})\}$.

290 The set is ordered by ascending time-stamp values. A typical *YET* may
 291 comprise thousands to millions of trials, and each trial may have approximately
 292 between 800 to 1500 ‘event time-stamp’ pairs, based on a global event catalogue
 293 covering multiple perils. The representation of the *YET* is shown in Equation 1,
 294 where $i = 1, 2, \dots$ and $k = 1, 2, \dots, 1500$.

$$YET = \{T_i = \{(E_{i,1}, t_{i,1}), \dots, (E_{i,k}, t_{i,k})\}\} \quad (1)$$

295 ii. *Event Loss Tables*, which is a representation of collections of specific
 296 events and their corresponding losses with respect to an exposure set denoted
 297 as *ELT*. Each record in an *ELT* is denoted as $EL_i = \{E_i, l_i\}$ and the financial
 298 terms associated with the *ELT* are represented as a tuple $\mathcal{I} = (\mathcal{I}_1, \mathcal{I}_2, \dots)$.

A typical aggregate analysis may comprise 10,000 *ELTs*, each containing
 10,000-30,000 event losses with exceptions even up to 2,000,000 event losses. The
ELTs can be represented as shown in Equation 2, where $i = 1, 2, \dots, 30,000$.

$$ELT = \left\{ \begin{array}{l} EL_i = \{E_i, l_i\}, \\ \mathcal{I} = (\mathcal{I}_1, \mathcal{I}_2, \dots) \end{array} \right\} \quad (2)$$

299 iii. *Portfolio*, which is denoted as *PF* and contains a group of Programs, *P*
 300 represented as $PF = \{P_1, P_2, \dots, P_n\}$ with $n = 1, 2, \dots, 10$.

301 Each Program in turn covers a set of Layers, denoted as *L*, cover a collection
 302 of *ELTs* under a set of layer terms. A single layer L_i is composed of two at-
 303 tributes. Firstly, the set of *ELTs* $\mathcal{E} = \{ELT_1, ELT_2, \dots, ELT_j\}$, and secondly,
 304 the Layer Terms, denoted as $\mathcal{T} = (\mathcal{T}_1, \mathcal{T}_2, \dots)$.

A typical Layer covers approximately 3 to 30 individual *ELTs* and is represented as shown in Equation 3, where $j = 1, 2, \dots, 30$.

$$L = \left\{ \begin{array}{l} \mathcal{E} = \{ELT_1, ELT_2, \dots, ELT_j\}, \\ \mathcal{T} = (\mathcal{T}_1, \mathcal{T}_2, \dots) \end{array} \right\} \quad (3)$$

305 The output of the analysis is a loss value associated with each trial of the
 306 *YET*. A reinsurer can derive important portfolio risk metrics such as the Prob-
 307 able Maximum Loss (PML) [30] and the Tail Value-at-Risk (TVaR) [31] which
 308 are used for both internal risk management and reporting to regulators and
 309 rating agencies. Furthermore, these metrics flow into a final stage of the risk
 310 analytics pipeline, namely Enterprise Risk Management, where liability, asset,
 311 and other forms of risks are combined and correlated to generate an enterprise
 312 wide view of risk.

313 4.2. Algorithm and GPU Implementation

314 Given the above three inputs, Aggregate Risk Analysis is shown in Algo-
 315 rithm 1. The data tables, *YET*, *ELT* and *PF*, are loaded into host (CPU)
 316 memory. The analysis is performed for each Layer and for each Trial in the
 317 *YET* and a Year Loss Table (*YLT*) is produced. In this paper, we assume a
 318 Portfolio comprising one Program and one Layer, and therefore the for loops of
 319 lines 1 and 2 iterate once. If there are N available devices (GPUs), then the
 320 *YET* is split to N smaller *YETs*, represented as YET_i , where $i = 1, 2, \dots, N$.

321 There are two functions that facilitate device execution. The first function
 322 `TransferDataToDevice` copies YET_i and the *ELT* to the device memory as
 323 shown in Algorithm 2.

324 The second function `LaunchDeviceKernel` executes the function on the de-
 325 vice as shown in Algorithm 3. Each event of a trial and its corresponding event
 326 loss in the set of *ELTs* associated with the Layer is determined. A set of con-
 327 tractual financial terms (\mathcal{I}) are applied to each loss value of the Event-Loss pair
 328 extracted from an *ELT* to the benefit of the layer. The event loss for each event
 329 occurrence in the trial, combined across all *ELTs* associated with the layer, are
 330 subject to further financial terms (\mathcal{T}) [28].

Algorithm 1: Aggregate Risk Analysis

Input : YET, ELT, PF

Output: YLT

```
1 for each Program,  $P$ , in  $PF$  do
2   for each Layer,  $L$ , in  $P$  do
3     Split  $YET$  to  $YET_i$ , where  $i = 1, 2, \dots, N$ 
4     for each  $i$  do
5       TransferDataToDevice ( $i, YET_i, ELT$ )
6       LaunchDeviceKernel ( $i$ )
7     end
8   end
9 end
10 Populate  $YLT$  from  $YLT_i$ , where  $i = 1, 2, \dots, N$ 
11 return
```

331 Two occurrence terms, namely (i) Occurrence Retention, \mathcal{T}_{OccR} , which is the
332 retention or deductible of the insured for an individual occurrence loss, and (ii)
333 Occurrence Limit, \mathcal{T}_{OccL} , which is the limit of coverage the insurer will pay for
334 occurrence losses in excess of the retention are applied. Occurrence terms are
335 applicable to individual event occurrences independent of any other occurrences
336 in the trial. The event losses net of occurrence terms are then accumulated into
337 a single aggregate loss for the given trial. The occurrence terms are applied as
338 $l_T = \min(\max(l_T - \mathcal{T}_{OccR}), \mathcal{T}_{OccL})$.

Algorithm 2: TransferDataToDevice Function

Input : i

```
1 Select device  $i$ 
2 Copy  $YET_i, ELT$  to device  $i$ 
3 return
```

Algorithm 3: LaunchDeviceKernel Function

Input : i

Output: YLL_i

```
1 Select device  $i$ 
2 for each Trial,  $T$ , in  $YET_i$  do
3   for each Event,  $E$ , in  $T$  do
4     for each  $ELT$  covered by  $L$  do
5       Lookup  $E$  in the  $ELT$  and find corresponding loss,  $l_E$ 
6       Apply Financial Terms to  $l_E$ 
7        $l_T \leftarrow l_T + l_E$ 
8     end
9     Apply Financial Terms to  $l_T$ 
10  end
11 end
12 return
```

339 Two aggregate terms, namely (i) Aggregate Retention, \mathcal{T}_{AggR} , which is the
340 retention or deductible of the insured for an annual cumulative loss, and (ii)
341 Aggregate Limit, \mathcal{T}_{AggL} , which is the limit or coverage the insurer will pay for
342 annual cumulative losses in excess of the aggregate retention are applied. Ag-
343 gregate terms are applied to the trial’s aggregate loss for a layer. The aggregate
344 loss net of the aggregate terms is referred to as the trial loss or the year loss.
345 The aggregate terms are applied as $l_T = \min(\max(l_T - \mathcal{T}_{AggR}), \mathcal{T}_{AggL})$.

346 A single thread is employed for the computations of each trial of the ap-
347 plication. $ELTs$ corresponding to a Layer were implemented as direct access
348 tables to facilitate fast lookup of losses corresponding to events. Each ELT is
349 implemented as an independent table; therefore, in a read cycle, each thread
350 independently looks up its events from the $ELTs$. All threads within a block
351 access the same ELT . The device global memory stores all data required for the
352 analysis. Chunking, which refers to processing a block of events of fixed size (or

353 chunk size) for the efficient use of shared memory is employed to optimise the
354 implementation; the computations related to the events in a trial and for apply-
355 ing financial terms benefit from chunking. The financial terms are stored in the
356 streaming multi-processor’s constant memory. In this case, chunking reduces
357 the number of global memory update and global read operations.

358 In this paper, the implementation of fine-grain parallelism in `LaunchDeviceKernel`
359 is not the focus. Instead, the optimisation of performance and efficiency of re-
360 source utilisation by managing the two functions, namely `TransferDataToDevice`
361 and `LaunchDeviceKernel` on virtual GPUs is considered and reported in the
362 next section.

363 5. Evaluation

364 In this section we optimise the performance of the financial risk application
365 to reduce its execution time such that real-time response can be achieved. To
366 this end we present (i) the hardware platform on which the experiments are
367 performed and, (ii) the use of the remote GPU virtualisation framework, and
368 (iii) an approach for transferring data from a CPU to GPUs with the aim of
369 reducing the execution time.

370 5.1. Platform

371 The experimental platform employed in this research comprises 1027GR-
372 TRF Supermicro nodes. Each node contains two Intel Xeon E5-2620 v2 proces-
373 sors, each with six cores, operating at 2.1 GHz and 32 GB of DDR3 SDRAM
374 memory at 1600 MHz. Each node has a Mellanox ConnectX-3 VPI single-port
375 InfiniBand adapter (InfiniBand FDR) as well as a Mellanox ConnectX-2 VPI
376 single-port adapter (InfiniBand QDR). The nodes are connected either by a Mel-
377 lanox switch MTS3600 with QDR compatibility (a maximum rate of 40Gb/s)
378 or by a Mellanox Switch SX6025, which is compatible with InfiniBand FDR (a
379 maximum rate of 56Gb/s). One NVIDIA Tesla K20 GPU is available for accel-
380 eration on each node. Additionally, one SYS7047GR-TRF Supermicro server

Table 1: Scalability of the financial risk application when executed using CUDA

	No. of GPUs		
	1 GPU	2 GPUs	4 GPUs
Total execution time	10.928	5.53	2.857
Normalised execution time	1	0.506	0.261
Execution time with perfect scalability	10.928	5.464	2.732
Offset with respect to perfect scalability	0	0.066	0.125
% offset with respect to perfect scalability	0	1.2%	4.57%

381 with identical processors was populated with 4 NVIDIA Tesla K20 GPUs and
 382 128 GB of DDR3 SDRAM memory at 1600MHz, to serve as a local server for
 383 the purpose of comparison. The CentOS 6.4 operating system was used, and
 384 the Mellanox OFED 2.4-1.0.4 (InfiniBand drivers and administrative tools) was
 385 used at the servers along with CUDA 6.5.

386 5.2. Application Scalability

387 As presented in Section 1 the use of multiple GPUs reduces the execution
 388 time of the application by evenly distributing computations across the GPUs
 389 assigned to the application. However, a closer look at the performance as shown
 390 in Figure 1 highlights that the scalability of the application as the number of
 391 GPUs increases is sub-linear. Table 1 is the result of executing the application
 392 on the Supermicro SYS7047GR-TRF server using CUDA with up to four GPUs.
 393 The normalised execution time indicates that perfect scalability is not achieved.
 394 For example, when two GPUs are used the normalised execution time should be
 395 0.5 instead of 0.506 and similarly when four GPUs are employed 0.25 is expected
 396 as against 0.261. The offset of execution time with respect to perfect scalability
 397 as a reference increases with the number of GPUs involved in the computations.

398 To account for sub-linear scalability further investigations were carried out.
 399 The time taken for computations on the GPUs and the time taken for trans-
 400 ferring data to the GPUs (1, 2, and 4 GPUs) were considered as shown in

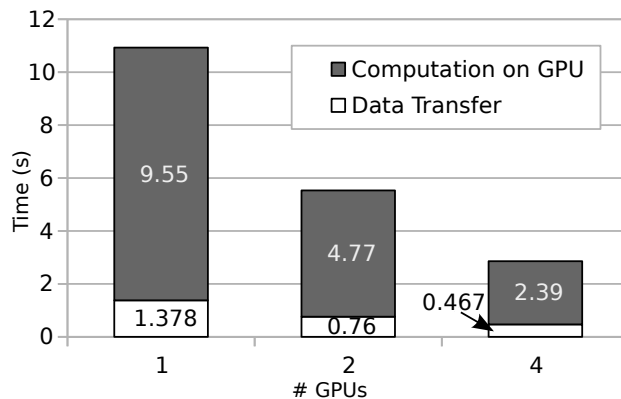


Figure 6: Computation and data transfer times for the financial risk application when executed on single and multiple GPUs with CUDA

401 Figure 6. The GPU computations take most of the execution time of the ap-
 402 plication (87.39%, 86.25%, and 63.65% of the total application execution time
 403 when 1, 2, and 4 GPUs are used respectively). The GPU computations scale in a
 404 perfect manner as the number of GPUs available to the application is increased.
 405 However, the time taken for data transfer does not scale well and accounts for
 406 12.6%, 13.74%, and 16.34% of total execution time when 1, 2, and 4 GPUs are
 407 used, respectively.

408 At first glance, it can be assumed that the increase in data transfer time
 409 may be due to the lower communication bandwidth of CUDA for transfers of
 410 small chunks of data (refer Figure 5c and Figure 5d). When pageable memory is
 411 transferred the attained bandwidth for data smaller than 10 MB is significantly
 412 reduced. Therefore, given that the size of input data transferred to each GPU
 413 is progressively reduced as the number of GPUs increases, then the input data
 414 may be smaller than 10 MB and thus the effective bandwidth for moving data
 415 to the GPUs is reduced in practice. However, in the case of our application the
 416 initial data size is 4 GB and when this data is shared among four GPUs the
 417 data transferred to each GPU is larger than 10 MB. Hence, the data transfer to
 418 the GPUs is performed at full bandwidth.

419 A closer look at the application reveals that the *YET* data structure (4 GB)

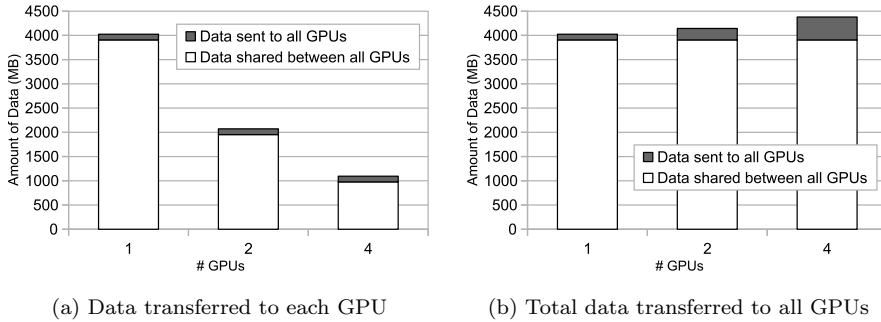


Figure 7: Amount of data transferred during the execution of the financial risk application

420 presented in Section 4 is uniformly split between the GPUs for computations.
 421 However, the *ELTs* and *PF* data structures (120 MB and 4 MB) are not split
 422 between the GPUs, instead are transferred fully to each GPU. Consequently, the
 423 total data movement to GPUs increases which is shown in Figure 7. Excluding
 424 the *ELTs*, the data that is not split between the GPUs is less than 10 MB
 425 resulting in a lower bandwidth for transferring this data requiring an additional
 426 2.6 milliseconds. However, this cannot fully account for sub-linear performance.

427 One important reason for the degradation of performance is data transfers
 428 to all GPUs are concurrently performed. Although each GPU is located in a
 429 different PCIe link, all data is extracted from main memory, which results in a
 430 bottleneck. This memory bottleneck is highlighted in Figure 8, which shows the
 431 bandwidth attained for each individual data copy when several data transfers
 432 are carried out concurrently to different destination GPUs by a single memory
 433 controller.

434 We summarise that for the financial risk application executing on multiple
 435 GPUs data transfers do not scale perfectly as the computations for two reasons.
 436 Firstly, there are input data structures that cannot be split between the GPUs
 437 and need to be copied onto each GPU creating an overhead. Secondly, concu-
 438 rent data transfers from the CPU main memory to GPUs result in a bottleneck
 439 at the memory controller.



Figure 8: Attained bandwidth when concurrent data transfers to GPUs are performed. Source data is located in the same memory bank.

440 *5.3. Reducing Execution Time Using rCUDA*

441 Current servers are constrained in the number of GPUs that can be accom-
 442 modated on them¹. We believe remote GPU virtualisation (in this research
 443 rCUDA is employed) is an appropriate mechanism to make a large number of
 444 GPUs available to an application. Figure 9a and Figure 9b present the perfor-
 445 mance of the application using the QDR InfiniBand and the FDR InfiniBand
 446 networks respectively for up to 16 GPUs.

447 Figure 9 indicates that the computation times when using rCUDA on 1, 2,
 448 and 4 GPUs are the same as shown in Figure 6 using CUDA. This is expected
 449 given that the computation time on the GPU is independent of whether it is
 450 on the same node as the application or on a remote node. With increasing
 451 number of GPUs there is perfect scalability. When 16 GPUs are employed, the
 452 computation time is less than one second (0.62 seconds) making it possible to
 453 do an industry size simulation in real-time.

¹Manufacturers, such as Cirrascale and Supermicro, have integrated up to 8 GPU cards in a single server. However, these are exceptions and costly options. Moreover, there are performance bottlenecks since the GPUs are usually grouped as a set of four cards that share a single PCIe x16 link with a processor socket. This results in slower communication between main memory and the GPUs. Performance is further degraded when a GPU card comprises multiple devices.

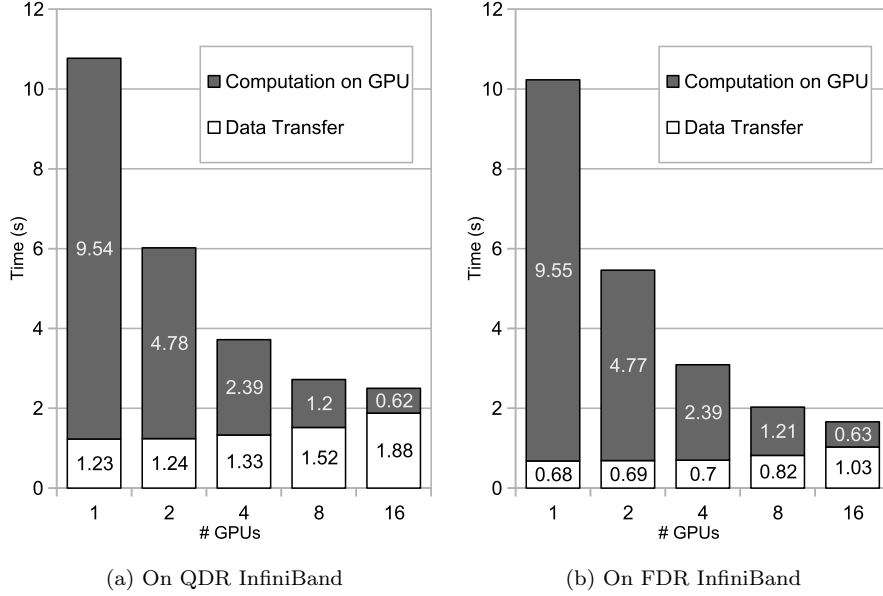


Figure 9: Scalability of the financial risk application when executed with rCUDA.

454 Two observations are made regarding data transfers. Firstly, when one re-
 455 mote GPU is used, the data transfer time using rCUDA is better than using
 456 CUDA (CUDA requires 1.378 seconds whereas rCUDA takes 1.23 seconds with
 457 QDR InfiniBand and 0.68 seconds with FDR InfiniBand). This lower transfer
 458 time as considered in Figure 5c is because rCUDA obtains more bandwidth
 459 than CUDA by using pageable memory. The improvement of communication
 460 performance is seen in Figure 9b for 2 GPUs.

461 Secondly, data transfer using rCUDA follows a different trend to CUDA. For
 462 CUDA the data transfer times to each GPU reduced as the number of GPUs
 463 increased (refer Figure 6). On the contrary, rCUDA time increases when both
 464 QDR and FDR InfiniBand are used. This is not surprising since the reasons for
 465 sub-linear scalability of data transfer time considered in the previous section is
 466 applicable for rCUDA. In this case, the bandwidth bottleneck is the InfiniBand
 467 card in the cluster node executing the application, which is a single communi-
 468 cation link for all the GPUs. This bottleneck is highlighted in Figure 10.

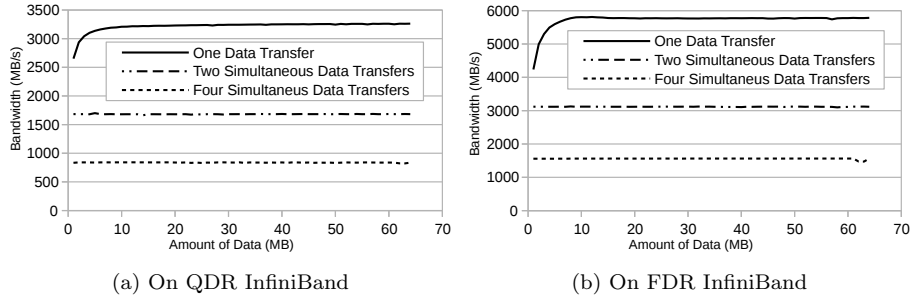


Figure 10: Bandwidth attained for multiple data transfers concurrently to different remote GPUs using rCUDA.

469 Figure 10 shows the bandwidth achieved for individual data transfer to a
 470 different remote GPU when multiple transfers are executed concurrently. The
 471 bandwidth for each transfer is proportional to the number of data movement
 472 operations in progress. In addition to the previous observations that result in
 473 an increase of data transfer times, there are a large number of `cudaMalloc()`
 474 functions that are invoked prior to the data transfer (the memory allocation
 475 time is included in the data transfer time). In rCUDA, memory allocations for
 476 a large number of data structures on remote GPUs requires 2.7 milliseconds with
 477 FDR InfiniBand (compared to 1.7 milliseconds in CUDA on a local GPU) and
 478 2.67 milliseconds with QDR InfiniBand (lower time due to low latency, despite
 479 reduced bandwidth [32]). Therefore, when a large number of GPUs are used
 480 by an application the time required for memory allocations can increase up to
 481 43.2 milliseconds for 16 remote GPUs; this is 4.2% of the total data transfer
 482 time.

483 The use of rCUDA allows to leverage a large number of GPUs to speed up
 484 the application despite poor performance for data transfers. The total execution
 485 time is reduced from 2.86 seconds when using local GPUs on CUDA to 1.66
 486 seconds when using remote GPUs on rCUDA. Reducing the total execution
 487 time enables the application to provide a solution in real-time.

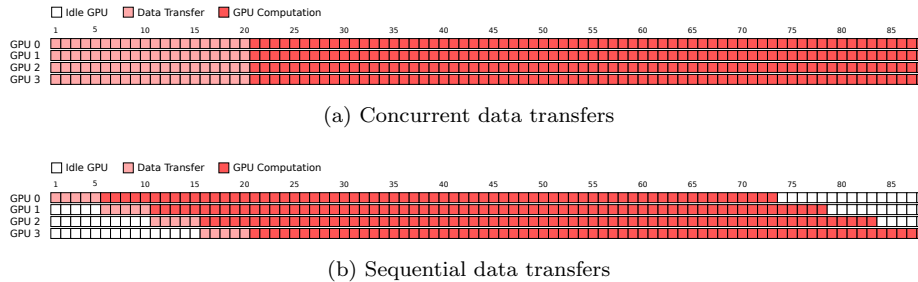


Figure 11: Communication approaches for transferring data to GPUs.

488 *5.4. Mitigating the Impact of Data Transfers in rCUDA*

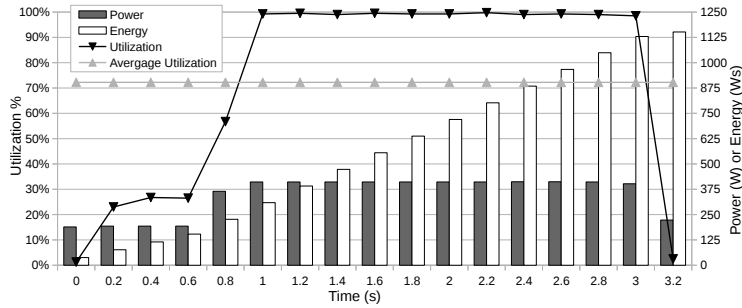
489 In this section, we consider two data transfer modes, namely concurrent and
 490 sequential, and further develop an approach based on multi-tenant GPUs in
 491 rCUDA.

492 *5.4.1. Concurrent vs Sequential Data Transfers*

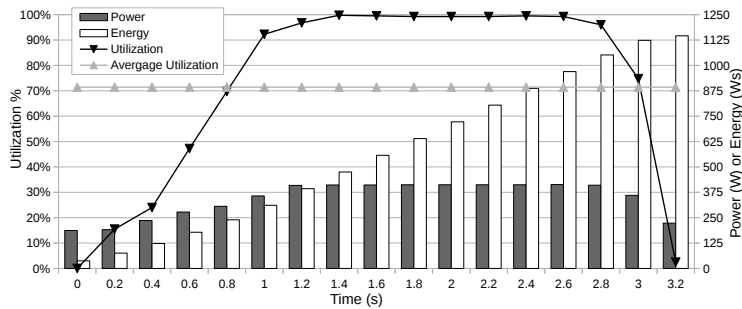
493 Figure 11a shows the life cycle of execution of a real application using rCUDA
 494 with four remote GPUs and FDR InfiniBand. Each cell represents execution
 495 time of 35 milliseconds. This corresponds to the four GPU execution shown in
 496 Figure 9b. The same amount of data is moved to the four GPUs concurrently
 497 by interleaving across the network and the remote GPUs start computations at
 498 the same time approximately. However, from Figure 10 it was noted that the
 499 bandwidth achieved is inversely proportional to the number of multiple data
 500 transfers concurrently performed which results in degrading performance.

501 An alternate method is shown in Figure 11b. Data to the first GPU is trans-
 502 ferred without sharing the bandwidth for the remaining three data streams.
 503 Since there is no competition for bandwidth it only takes a quarter of the time
 504 required when data is concurrently transferred (shown in Figure 11a). Compu-
 505 tations on the first GPU start while data is transferred to the second GPU. In
 506 this manner, data transfer is performed on fully available network bandwidth.
 507 This is referred to as the sequential data transfer method.

508 Data is transferred at full network bandwidth and there is an overlap with
 509 GPU computations in the sequential data transfer approach. However, it is



(a) Concurrent

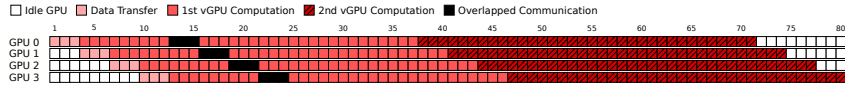


(b) Sequential

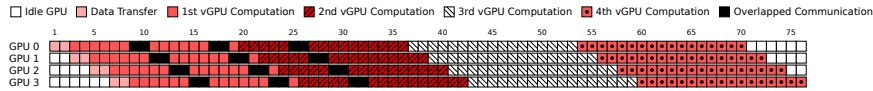
Figure 12: GPU utilisation, power and energy consumption of concurrent and sequential data transfers to GPUs considered in Figure 11

510 noted that the execution time is not reduced since the fourth GPU begins its
 511 computations when it would in concurrent data transfers. Figure 12 shows the
 512 GPU utilisation, power and energy consumption of concurrent and sequential
 513 data transfers to GPUs. The average values of the four GPUs considered in
 514 Figure 11 are used. The Y-axis on the left indicates GPU utilisation and the
 515 Y-axis on the right shows power (in Watts) and energy (in Watts per second,
 516 denoted as Ws in the figure) consumed. The power and energy of GPUs are
 517 measured instead of the cluster since multiple GPU configurations (n GPUs
 518 per node) could be employed, which results in different energy measurements.
 519 There are no gains in the energy consumed and very little difference in GPU
 520 utilisation for both concurrent and sequential transfers.

521 Regardless, in this research sequential data transfer is foundational in devel-



(a) 2 vGPUs per GPU



(b) 4 vGPUs per GPU

Figure 13: Sequential data copies with several vGPUs per GPU.

522 oping an optimised approach for executing the application using remote GPUs
 523 which is based on multi-tenancy of virtual GPUs.

524 5.4.2. Multi-tenancy Approach

525 The key concept of the multi-tenancy approach is based on the fact that
 526 current GPUs perform kernel executions and DMA (Direct Memory Access)
 527 operations concurrently. If it were possible to move data to a GPU the same
 528 time it was executing a kernel, there could be gains in further improving the
 529 performance of the executing application.

530 This can be facilitated by a multi-tenancy approach in which a number of
 531 remote GPUs (or virtual GPUs referred to as vGPUs) reside on or are mapped
 532 onto the same physical GPU (pGPU)². Figure 13 shows the concept of multi-
 533 tenancy when 2 and 4 vGPUs are mapped to a pGPU.

534 When 2 vGPUs are mapped on to a pGPU as shown in Figure 13a 8 GPUs
 535 are available to the application (4 pGPUs are used). Input data will be split

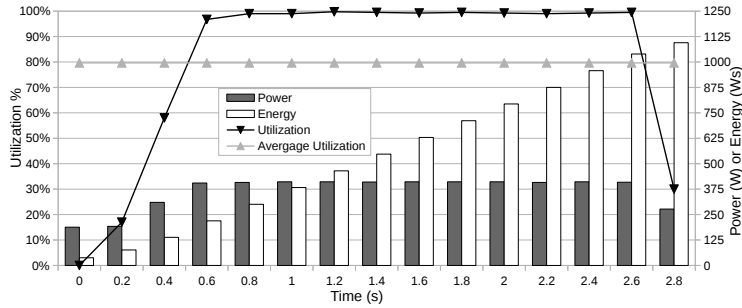
²Multi-tenancy is achieved on rCUDA by setting two environment variables prior to application execution, namely `RCUDA_DEVICE_COUNT` and `RCUDA_DEVICE_j`. The first variable indicates the number of GPUs accessible to the application. The second variable indicates the cluster node in which the j^{th} GPU is located. For example, “`export RCUDA_DEVICE_COUNT=2`” when 2 GPUs are assigned to the application and “`export RCUDA_DEVICE_0=192.168.0.1`” and “`export RCUDA_DEVICE_1=192.168.0.2`”. The server of the `RCUDA_DEVICE_j` variables need to point to the same node. Hence, the application does not require to be modified to accommodate multi-tenancy using rCUDA.

536 such that 8 GPUs will be used for computations. The initial data transfer is
537 shown as “*Data Transfer*” followed by computations by the first vGPU labelled
538 as “*1st vGPU Computation*”. After transferring data in the 12th time step, there
539 are four more vGPUs that will require their input data. Data transferred to
540 the remaining four vGPUs beginning at time step 13 are overlapped with the
541 computations of the first four vGPUs. Since two vGPUs are mapped onto a
542 single pGPU, computations of both vGPUs cannot progress in parallel as they
543 belong to different GPU contexts. Therefore, the NVIDIA driver executes them
544 sequentially (using as many GPU resources required by each kernel). So the
545 second kernel must wait until the execution of the first kernel is completed.

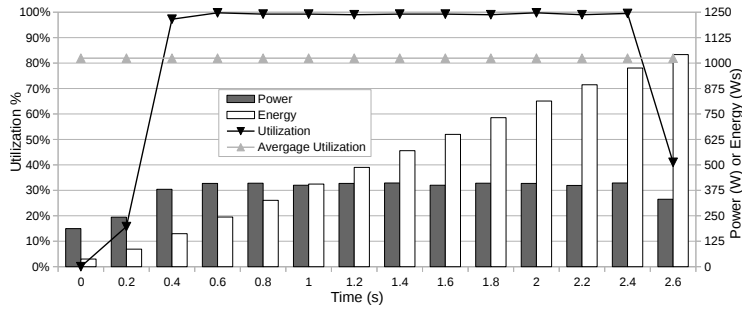
546 Two key observations are made from multi-tenant executions. Firstly, the
547 total execution time has reduced in contrast to the execution life cycle presented
548 in Figure 11b although the same hardware resources are used. The application
549 completed execution in time step 80 using 2 vGPUs per pGPU compared to
550 time step 88 when no multi-tenancy is employed. The time that each GPU
551 computes is exactly the same. The time saved is because of the overlap between
552 computations and data transfers of multiple vGPUs on the same pGPU. In
553 Figure 11b data transfers overlapped with computations of other pGPUs but
554 there were no overlaps on the same GPU.

555 Secondly, the data transfer time takes longer when more vGPUs are em-
556 ployed. In Figure 11b, data is transferred completely to all GPUs at time step
557 20, whereas in Figure 13a, the input data arrives at time step 24. The reasons
558 for longer data transfer times have been considered in the previous section. De-
559 spite the larger data transfer time, the total execution time gains since there is
560 an overlap between computation and data movement.

561 Figure 13b shows the use of 16 vGPUs mapped on to 4 pGPUs. The execu-
562 tion time is further reduced due to the larger overlap between computation and
563 data transfers when compared to 2 vGPUs residing on a single pGPU. Again
564 the time for computing is the same on each physical GPU but the data copying
565 time has increased. The overall execution time is further reduced to 76 time
566 steps.



(a) 2 vGPUs per pGPU



(b) 4 vGPUs per pGPU

Figure 14: GPU utilisation, power and energy consumption of the multi-tenancy approach considered in Figure 13.

567 Multi-tenancy can be analysed from the perspective of energy required to
 568 complete the execution of the application. Figure 14 shows the energy consumed
 569 during the execution of the application along with the utilization of the
 570 physical GPU. The multi-tenancy energy consumption is lower than sequential
 571 communications without an overlap between data transfers and computations
 572 on the same GPU seen in Figure 12. The energy consumed is 1145 Watts per
 573 second without using multi-tenancy and 1094 and 1041 Watts per second when
 574 2 and 4 vGPUs are tenants on a pGPU, respectively. It is observed that GPU
 575 utilisation increases in the multi-tenancy approach. The average GPU utilisation
 576 rises from 71.44% without multi-tenancy up to 79.65% for 2 vGPUs per
 577 pGPU and up to 81.93% when 4 vGPUs are mapped on to a pGPU.

578 In short, multi-tenancy allows for data transfers to be overlapped with com-
 579 putations on the same GPUs thereby reducing total execution time of the fi-

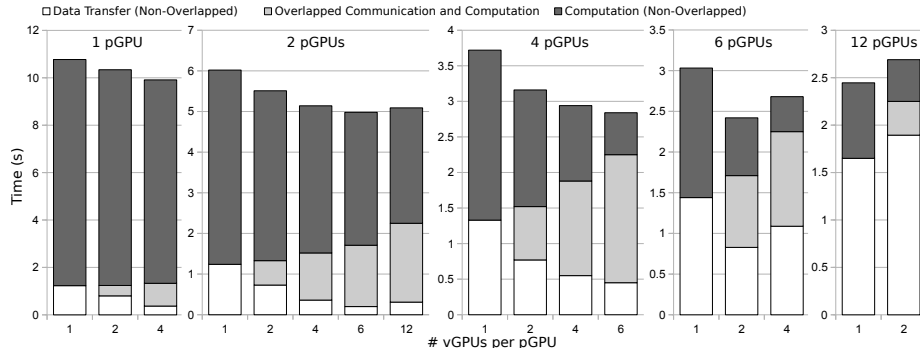


Figure 15: Application performance for different combinations of pGPUs and vGPUs using QDR InfiniBand

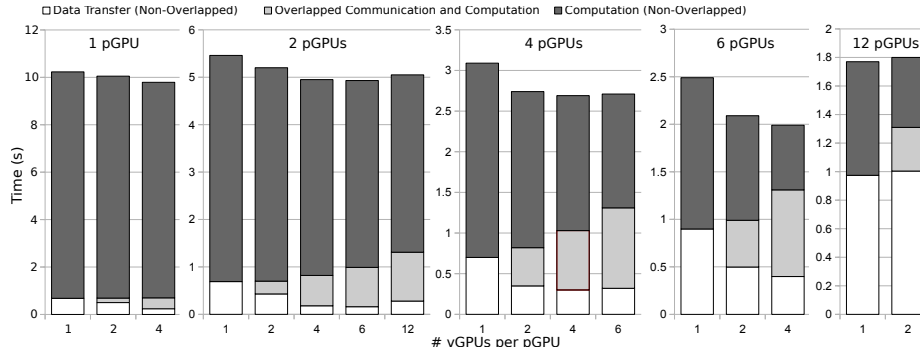


Figure 16: Application performance for different combinations of pGPUs and vGPUs using FDR InfiniBand

580 nancial risk application. Furthermore, the energy required to execute the appli-
 581 cation is reduced and the GPU utilisation is increased.

582 5.5. Performance Analysis Using Multi-tenancy

583 An analysis of the application performance as measured by execution time is
 584 presented in this section. The cluster nodes in our experimental set up have 12
 585 cores (up to 24 threads with hyper-threading) and therefore we use a maximum
 586 of 24 vGPUs (to avoid any noise due to CPU overhead). Up to 12 pGPUs will
 587 be used to map the vGPUs.

588 Figure 15 and Figure 16 show the time taken for data transfer and compu-
 589 tation for varying pGPUs when the rCUDA framework is used over QDR and

590 FDR InfiniBand. The ‘Overlapped data transfer and computation’ label denotes
591 that data transfers and computation are carried out concurrently on the same
592 pGPU. The behaviour of the application is as expected. Multi-tenancy with
593 sequential transfers allows for overlapping computations and data movement on
594 the same pGPU, thus reducing the execution time. When QDR InfiniBand is
595 used, time for data transfer without overlaps with communication is reduced up
596 to 70%, 84%, 66%, and 42% when vGPUs are mapped to 1, 2, 4, and 6 pGPUs,
597 respectively. In the case of FDR InfiniBand, the same time is 65%, 77%, 57%,
598 and 56%. Consequently, the total power consumed is reduced but not indicated
599 on the graph.

600 It is noted that when 12 pGPUs are used the data transfer times are not
601 reduced further because (i) the execution time decreases with more pGPUs,
602 and (ii) the data transfer time increases when more vGPUs are used allowing
603 for little overlap between data transfers and computation on the same pGPU.
604 This necessitates the need for determining the effective combination of pGPUs
605 and vGPUs by estimating application performance both in terms of execution
606 time and energy consumption.

607 *5.6. Modelling Multi-tenancy for Performance and Energy Estimation*

608 An important challenge is to automatically determine the best multi-tenancy
609 configuration for a deployment that can maximise performance (minimising ex-
610 ecution time), but at the same time minimise the energy consumed.

611 *5.6.1. Performance Model*

612 We firstly consider a basic model to account for execution time of the ap-
613 plication when sequential data transfers are used with rCUDA, but without
614 exploiting multi-tenancy. Subsequently, the model is optimised to take multi-
615 tenancy into account. The model is then applied in the context of the hardware
616 (NVIDIA Tesla K20 GPUs with QDR and FDR InfiniBand) we have employed
617 in this research.

The total execution time depends on: (i) time for transferring data and (ii) time for computing on the GPUs as shown in Equation 4, which inherently depends on the number of GPUs (pGPUs or vGPUs) available to the application.

$$TotalExecutionTime = T_{transfer}(\#GPUs) + T_{computation}(\#GPUs) \quad (4)$$

Since there is perfect scalability for the computation times on the GPU (Section 5.2 and Section 5.3), the time required for computations by a given number of GPUs can be obtained as shown in Equation 5.

$$T_{computation}(\#GPUs) = ComputationTime_{1pGPU} / \#GPUs \quad (5)$$

618 The time to transfer the input data to all GPUs is shown in Equation 6.
 619 The time taken to allocate memory on each GPU using `cudaMalloc()` and the
 620 time for moving small and large data structures to the GPUs are taken into
 621 account. Different data sizes achieve varying network bandwidth (Figure 5c).
 622 To simplify the equation, the time to transfer data structures smaller than 100
 623 bytes is denoted as $T_{small_transfers}$ ³

$$\begin{aligned} T_{transfer}(\#GPUs) = \#GPUs * (T_{cudaMalloc} + T_{small_transfers} \\ + T_{transfer_4MB} + T_{transfer_120MB}) \\ + T_{transfer_AGB} \end{aligned} \quad (6)$$

624 When multi-tenancy is taken into account there is an overlap between data
 625 transfers and computations on the same pGPU which reduces the total execu-
 626 tion time. As shown in Figure 13a, when 2 vGPUs are mapped onto a single
 627 pGPU, the time for data transfer is the time taken to move the first chunks
 628 of data to the pGPUs (until the completion of time step 12). The time for

³Data structures smaller than 100 bytes achieve the same bandwidth and are therefore grouped together. The InfiniBand frame size is typically 2 KB, which will be sent to the GPU in all cases where data is smaller than 100 bytes.

629 moving the remaining data chunks are not accounted for since it is overlapped
 630 by computation time. This is captured in Equation 7.

$$\begin{aligned}
 ExecTime_Multitenancy_{fully_overlapped} &= T_{transfer}(\#vGPUs) / vGPUs_per_pGPU \\
 &+ vGPUs_per_pGPU * T_{computation}(\#vGPUs)
 \end{aligned}
 \tag{7}$$

631 If a very large number of vGPUs are used, then all data transfer times
 632 may not be overlapped with computation times. This can happen when the
 633 computation on the vGPU is not long enough to overlap data transfers to the
 634 pGPU and the computations on it. In this case, the total execution time depends
 635 on the time required to copy data to all the vGPUs and is shown in Equation 8.

$$\begin{aligned}
 ExecTime_Multitenancy_{not_fully_overlapped} &= T_{transfer}(\#vGPUs) \\
 &+ T_{computation}(\#vGPUs)
 \end{aligned}
 \tag{8}$$

636 As shown in Equation 9 the maximum value from Equation 7 and Equa-
 637 tion 8 determines whether the application has significant overlaps between data
 638 transfer and computations.

$$\begin{aligned}
 ExecTime_Multitenancy &= MAX(ExecTime_Multitenancy_{fully_overlapped}, \\
 &ExecTime_Multitenancy_{not_fully_overlapped})
 \end{aligned}
 \tag{9}$$

639 Table 2 shows actual values of the model for the experimental platform used
 640 in this research.

641 Figure 17 and Figure 18 use these values in Equation 9 for 1 to 16 pGPUs
 642 and up to 12 vGPUs per pGPU. The combinations of pGPUs and vGPUs that
 643 require the lowest execution time can be explored in this space. The estimated
 644 execution times are grouped for 1 to 4 pGPUs, 5 to 8 pGPUs, 9 to 12 pGPUs,
 645 and 13 to 16 pGPUs. In Figure 17a and Figure 18a, for one pGPU up to 4
 646 vGPUs can be used. The total memory on the Tesla K20 devices is 4799 MB

Table 2: Time in seconds for GPU memory allocation and data transfer tasks of the financial risk application

Parameter	QDR	FDR
$ComputationTime_{1pGPU}$	9.55	
$T_{cudaMalloc}$	0.00267	0.0027
$T_{small_transfers}$	0.0048	0.0028
$T_{transfer_4MB}$	0.00133	0.00079
$T_{transfer_120MB}$	0.036	0.0205
$T_{transfer_4GB}$	1.171	0.67

647 (from the `nvidia-smi` command), which is exhausted by more than 4 vGPUs
648 (total memory size consumed by the application on 4 vGPUs is 4484 MB). It is
649 inferred from the figures that a large number of vGPU has detrimental effect on
650 performance due to the overheads in data movements. Using QDR InfiniBand
651 the model predicts a saturation sooner than FDR InfiniBand because of the
652 overhead of data transfers due to a lower bandwidth available on the QDR
653 network. The optimal deployment configuration of the application is 7 pGPUs
654 with 2 vGPUs per pGPU and 9 pGPUs with 2 vGPUs per pGPU using QDR
655 InfiniBand and FDR InfiniBand respectively.

656 5.6.2. Energy Model

657 The amount of energy required to execute the application is modelled in this
658 section. From Figure 13 it is inferred that a GPU can be in the following four
659 different states: (1) idle, (2) receive data, but no computations, (3) receive data
660 and compute simultaneously, and (4) compute, but no data to receive.

661 Power is measured by querying `nvidia-smi` every 200 milliseconds. The
662 power required by the GPU in the first two states is the same. The NVIDIA
663 Tesla K20 device requires 47 Watts while idling⁴ and receiving data. The GPU

⁴The idle state in Figure 13 is distinguished from the commonly known “idle” state. In

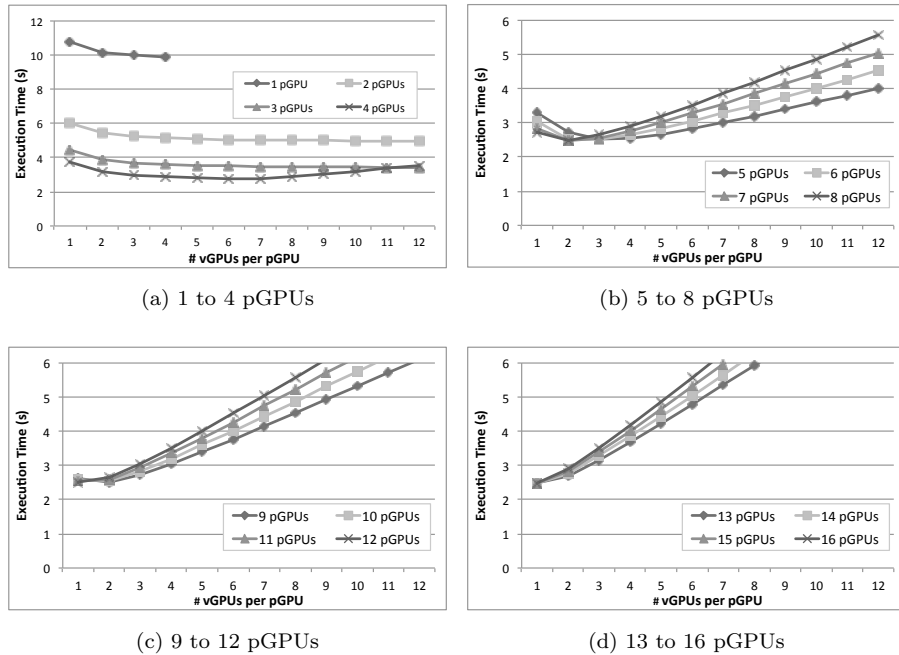


Figure 17: Results from performance model for QDR InfiniBand

664 requires 102 Watts in the last two states.

665 Using the above power readings for the four GPU states along with total
 666 execution time obtained from Equation 9 an energy model is developed as shown
 667 in Equation 10. The energy required by the GPU for computations (time spent
 668 on computations is obtained from Equation 5) is eliminated to obtain the energy
 669 spent in the first and second states. The computation time on the pGPUs is

Figure 13, the GPU has already been assigned to the application and therefore has been initialised by the GPU driver (this requires approximately 1.3 seconds in CUDA). After initialisation, the GPU does not perform any task, but actively waits for commands. In the commonly known “idle” state, the GPU is not assigned to an application and is not initialised by the driver. In this state, the Tesla K20 GPU requires 25 Watts.

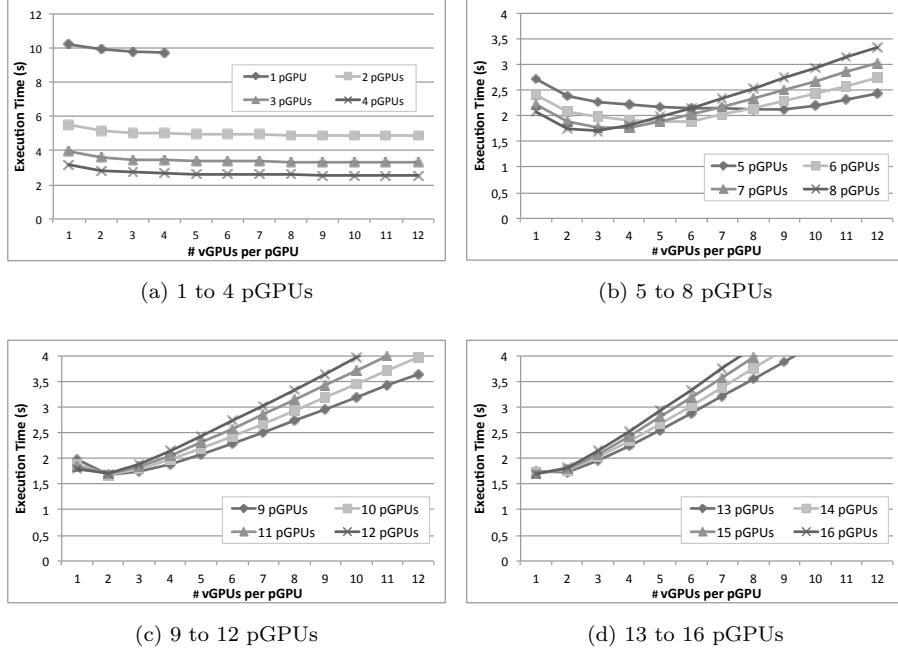


Figure 18: Results from performance model for FDR InfiniBand

670 $vGPUs_per_pGPUs * T_{computation}(\#vGPUs).$

$$\begin{aligned}
 TotalEnergy = \#pGPUs * (T_{computation}(\#pGPUs) * 102\ Watts + \\
 (ExecTime_Multitenancy - T_{computation}(\#pGPUs)) * 47\ Watts)
 \end{aligned}
 \tag{10}$$

671 Figure 19 and Figure 20 present the results of the energy model from Equa-
 672 tion 10. It is noted that an energy efficient deployment is obtained using 4
 673 vGPUs on 1 pGPU for both QDR InfiniBand and FDR InfiniBand. This is as
 674 expected given that the least amount of hardware is employed. However, there
 675 is a trade off since the lowest execution times are not obtained in this configura-
 676 tion. In Figure 21 and Figure 22, an alternate space (*energy * execution time*)
 677 is explored to find configurations that can maximise performance and minimise
 678 energy consumption.

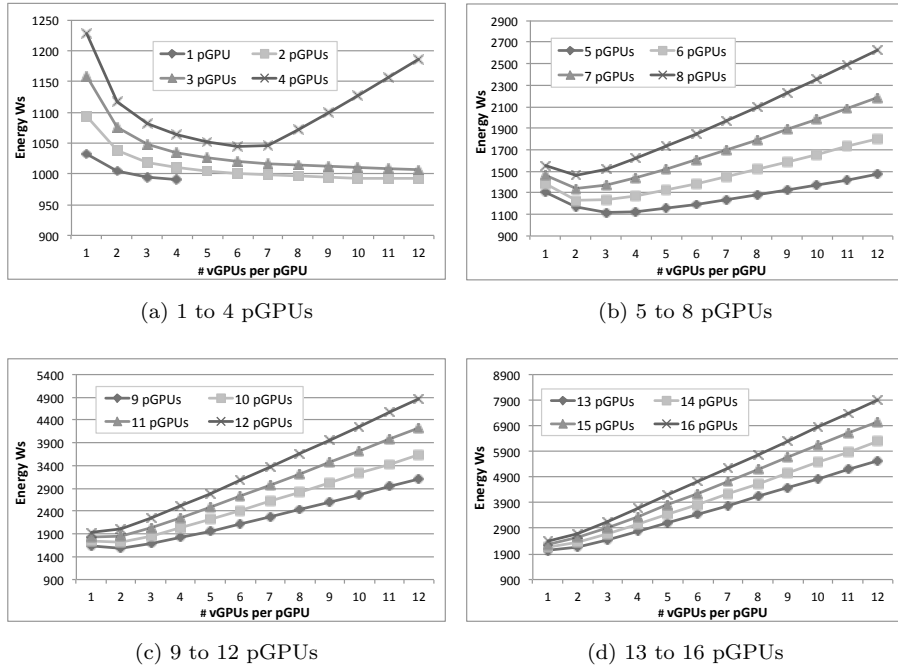


Figure 19: Results from energy model for QDR InfiniBand

679 *5.7. Generality of Proposed Approaches*

680 The financial risk application chosen in this paper is embarrassingly parallel
 681 and is representative of one class of workloads that execute in high-performance
 682 computing environments. The research challenges which were initially posed are
 683 hence relevant to a wide range of applications that aim to exploit vGPUs, particu-
 684 larly in the context of multi-tenant vGPUs on a single pGPU. The approaches
 685 we have proposed as solutions to mitigate the challenges can be broadly applied
 686 to the benefit of these applications.

687 Typically, when accelerators are employed for applications the data neces-
 688 sary for computations needs to be transferred from the host to the memory of
 689 vGPUs before computations can be actually performed. In the face of limited
 690 bandwidth for data transfers, linear scalability of the application will be affected
 691 degrading the overall performance of the application. However, by using our pro-
 692 posed approach of sequential data transfer performance can be improved; data

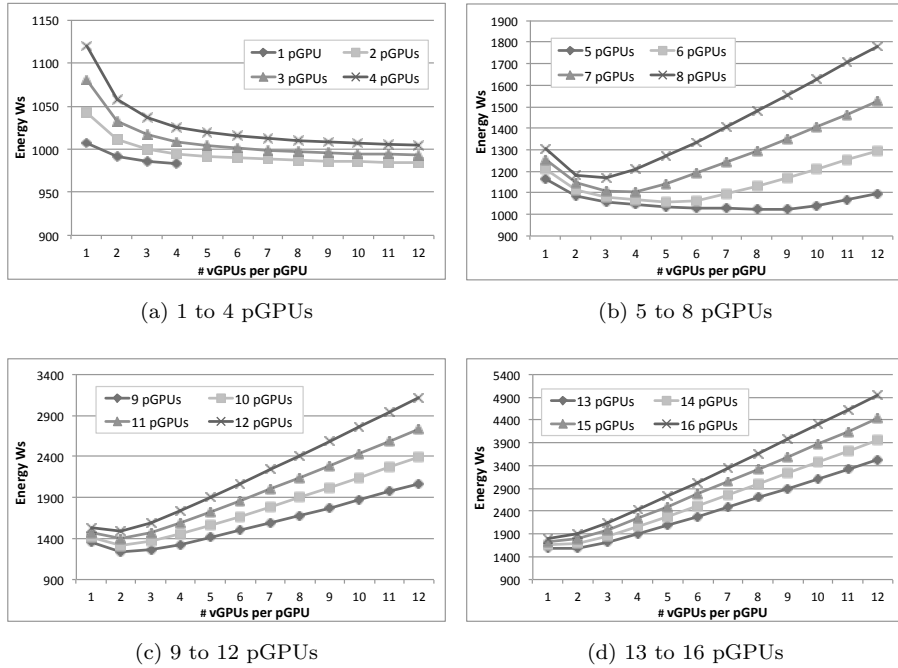
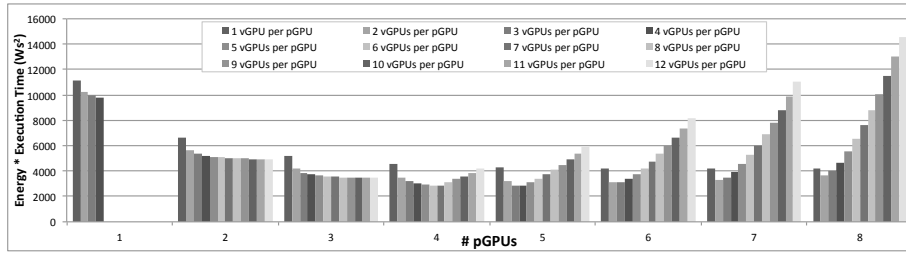


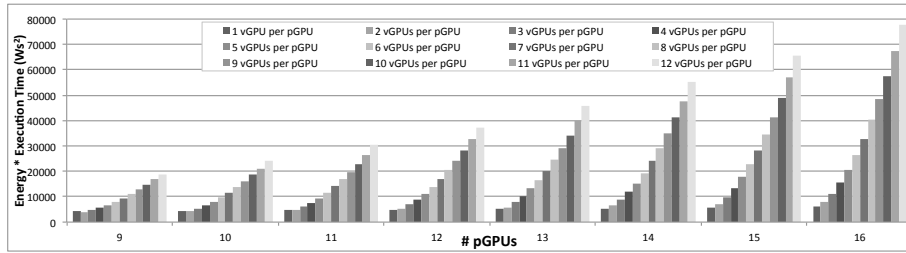
Figure 20: Results from energy model for FDR InfiniBand

693 transfers from the host to the GPU and GPU computations can be overlapped
 694 for multiple pGPUs. Performance can be further improved by incorporating our
 695 approach for overlapping data transfers and computation on multiple vGPUs
 696 which reside on the same pGPU. Such an approach effectively shares vGPUs to
 697 optimise an application's execution time and energy consumption.

698 There are multiple deployment options for an application when multi-tenancy
 699 is exploited. Each application will have its own best combination of vGPUs
 700 that need to be mapped onto a pGPU for best performance. Here our offline
 701 approach of modelling performance both in terms of energy and performance
 702 for estimations will be a handy method that can be broadly applied for other
 703 applications.

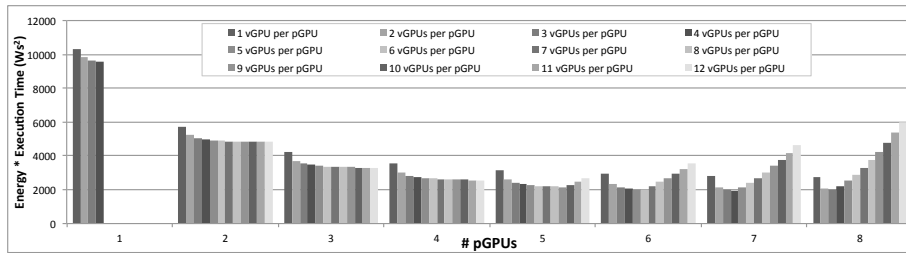


(a) 1 to 8 pGPUs

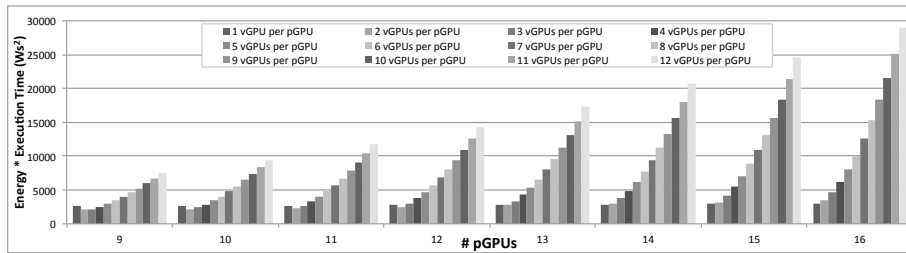


(b) 9 to 16 pGPUs

Figure 21: Combined space of energy and execution time using QDR InfiniBand



(a) 1 to 8 pGPUs



(b) 9 to 16 pGPUs

Figure 22: Combined space of energy and execution time using FDR InfiniBand

704 **6. Conclusions**

705 In this paper, we have demonstrated the benefits of virtual GPUs for an
706 application. Single tenancy (using one virtual GPU on a single physical GPU)
707 and multi-tenancy (using a number of virtual GPUs on a physical GPU) were
708 explored in this context. Concurrent and sequential data transfer models were
709 considered. We hypothesised that multi-tenancy can improve the performance
710 of the application. To validate the hypothesis the application was executed
711 using rCUDA (remote CUDA), a framework that virtualises GPUs in a High-
712 Performance Computing (HPC) cluster and provides remote GPUs to nodes
713 that require acceleration on demand. Experimental results indicate that multi-
714 tenant virtual GPUs with sequential data transfers optimise the performance of
715 the application with less hardware when compared to single tenancy.

716 This research highlights that multi-tenant virtual GPUs can improve per-
717 formance of an application. To achieve this we brought together the concepts
718 of virtual GPUs and multi-tenancy in a single framework. The contribution of
719 this research is to leverage multi-tenancy in the context of virtual GPUs within
720 the rCUDA framework. Further, we have demonstrated this concept using a
721 real world financial risk application of industrial use to optimise performance in
722 terms of metrics, namely execution time, energy consumption and GPU utilisation.
723 Given the application our research explores data transfer approaches with
724 the aim of improving performance and how it is affected by memory and band-
725 width bottlenecks. The experimental results provide insight that would not be
726 apparent without a thorough evaluation. For example, it may be assumed that
727 concurrent data transfers would improve performance, but the effect of memory
728 and bandwidth limitations make sequential data transfers more appealing. The
729 offline performance model is derived by making use of the experimental results
730 which determines the configuration of the vGPU mapping on the pGPU for
731 maximising performance.

732 **References**

- 733 [1] K. H. Tsoi, W. Luk, Axel: A heterogeneous cluster with FPGAs and GPUs,
734 in: Proceedings of the 18th Annual ACM/SIGDA International Symposium
735 on Field Programmable Gate Arrays, 2010, pp. 115–124.
- 736 [2] F. Song, J. Dongarra, A scalable framework for heterogeneous GPU-based
737 clusters, in: Proceedings of the 24th Annual ACM Symposium on Paral-
738 lelism in Algorithms and Architectures, 2012, pp. 91–100.
- 739 [3] M. Becchi, K. Sajjapongse, I. Graves, A. Procter, V. Ravi, S. Chakrad-
740 har, A virtual memory based runtime to support multi-tenancy in clusters
741 with GPUs, in: Proceedings of the 21st International Symposium on High-
742 Performance Parallel and Distributed Computing, 2012, pp. 97–108.
- 743 [4] D. Sengupta, R. Belapure, K. Schwan, Multi-tenancy on GPGPU-based
744 servers, in: Proceedings of the 7th International Workshop on Virtualiza-
745 tion Technologies in Distributed Computing, 2013, pp. 3–10.
- 746 [5] Y. Jiao, H. Lin, P. Balaji, W. Feng, Power and performance character-
747 ization of computational kernels on GPU, in: Proceedings of the 2010
748 IEEE/ACM Int’L Conference on Green Computing and Communications
749 & Int’L Conference on Cyber, Physical and Social Computing, 2010, pp.
750 221–228.
- 751 [6] S. Iserte, A. Castello, R. Mayo, E. Quintana-Orti, F. Silla, J. Duato, C. Re-
752 ano, J. Prades, Slurm support for remote GPU virtualization: Implemen-
753 tation and performance study, in: Computer Architecture and High Per-
754 formance Computing (SBAC-PAD), 2014 IEEE 26th International Sympo-
755 sium on, 2014, pp. 318–325.
- 756 [7] B. Varghese, J. Prades, C. Reano, F. Silla, Acceleration-as-a-service: Ex-
757 ploiting virtualised GPUs for a financial application, in: Proceedings of the
758 11th IEEE International Conference on eScience, 2015, pp. 47–56.

- 759 [8] A. J. Pena, C. Reano, F. Silla, R. Mayo, E. S. Quintana-Orti, J. Duato,
760 A complete and efficient CUDA-sharing solution for HPC clusters, *Parallel*
761 *Computing* 40 (2014) 574–588.
- 762 [9] A. Srinivasan, Parallel and distributed computing issues in pricing financial
763 derivatives through quasi monte carlo, in: *Parallel and Distributed Process-*
764 *ing Symposium.*, Proceedings International, IPDPS 2002, Abstracts and
765 CD-ROM, 2002.
- 766 [10] K. Huang, R. Thulasiram, Parallel algorithm for pricing american asian
767 options with multi-dimensional assets, in: *High Performance Computing*
768 *Systems and Applications*, 2005. HPCS 2005. 19th International Sympo-
769 sium on, 2005, pp. 177–185.
- 770 [11] C. Bekas, A. Curioni, I. Fedulova, Low cost high performance uncertainty
771 quantification, in: *Proceedings of the 2nd Workshop on High Performance*
772 *Computational Finance*, 2009.
- 773 [12] D. Daly, K. D. Ryu, J. Moreira, Multi-variate finance kernels in the blue
774 gene supercomputer, in: *High Performance Computational Finance*, 2008.
775 WHPCF 2008. Workshop on, 2008, pp. 1–7.
- 776 [13] V. Agarwal, L.-K. Liu, D. Bader, Financial modeling on the cell broadband
777 engine, in: *Parallel and Distributed Processing*, 2008. IPDPS 2008. IEEE
778 International Symposium on, 2008, pp. 1–12.
- 779 [14] C. Docan, M. Parashar, C. Marty, Advanced risk analytics on the cell
780 broadband engine (2009) 1–8.
- 781 [15] A. Irturk, B. Benson, N. Laptev, R. Kastner, FPGA acceleration of mean
782 variance framework for optimal asset allocation, in: *High Performance*
783 *Computational Finance*, 2008. WHPCF 2008. Workshop on, 2008, pp. 1–8.
- 784 [16] D. Thomas, Acceleration of financial monte-carlo simulations using FP-
785 GAs, in: *High Performance Computational Finance (WHPCF)*, 2010 IEEE
786 Workshop on, 2010, pp. 1–6.

- 787 [17] L. Abbas-Turki, S. Vialle, B. Lapeyre, P. Mercier, Pricing derivatives on
788 graphics processing units using monte carlo simulation, *Concurrency and*
789 *Computation: Practice and Experience* 26 (9) (2014) 1679–1697.
- 790 [18] D. M. Dang, C. C. Christara, K. R. Jackson, An efficient graphics process-
791 ing unit-based parallel algorithm for pricing multi-asset american options,
792 *Concurrency and Computation: Practice and Experience* 24 (8) (2012)
793 849–866.
- 794 [19] CUDA API Reference Manual 6.5 (2014).
- 795 [20] OpenCL 2.0 Specification (2013).
- 796 [21] T. Y. Liang, Y. W. Chang, Gridcuda: A grid-enabled CUDA programming
797 toolkit, in: *Advanced Information Networking and Applications (WAINA),*
798 *2011 IEEE Workshops of International Conference on*, 2011, pp. 141–146.
- 799 [22] L. Shi, H. Chen, J. Sun, vCUDA: GPU accelerated high performance com-
800 puting in virtual machines, in: *Parallel & Distributed Processing*, 2009.
801 *IPDPS 2009. IEEE International Symposium on*, 2009, pp. 1–11.
- 802 [23] V. Gupta, A. Gavrilovska, K. Schwan, H. Kharche, N. Tolia, V. Talwar,
803 P. Ranganathan, GViM: GPU-accelerated virtual machines, in: *Proceed-*
804 *ings of the 3rd ACM Workshop on System-level Virtualization for High*
805 *Performance Computing*, 2009, pp. 17–24.
- 806 [24] G. Giunta, R. Montella, G. Agrillo, G. Coviello, A GPGPU transparent
807 virtualization component for high performance computing clouds, in: *Pro-*
808 *ceedings of the 16th international Euro-Par conference on Parallel process-*
809 *ing*, 2010, pp. 379–391.
- 810 [25] M. Oikawa, A. Kawai, K. Nomura, K. Yasuoka, K. Yoshikawa, T. Narumi,
811 DS-CUDA: A middleware to use many GPUs in the cloud environment, in:
812 *Proceedings of the 2012 SC Companion: High Performance Computing,*
813 *Networking Storage and Analysis*, 2012, pp. 1207–1214.

- 814 [26] NVIDIA, The NVIDIA GPU Computing SDK Version 5.5 (2013).
- 815 [27] iperf3: A TCP, UDP, and SCTP network bandwidth measurement tool,
816 <https://github.com/esnet/iperf> (2015).
- 817 [28] A. Bahl, O. Baltzer, A. Rau-Chaplin, B. Varghese, Parallel simulations for
818 analysing portfolios of catastrophic event risk, in: High Performance Com-
819 puting, Networking, Storage and Analysis (SCC), 2012 SC Companion:,
820 2012, pp. 1176–1184.
- 821 [29] B. Varghese, The hardware accelerator debate: A financial risk case study
822 using many-core computing, *Computers & Electrical Engineering* 46 (2015)
823 157–175.
- 824 [30] G. Woo, Natural catastrophe probable maximum loss, *British Actuarial*
825 *Journal* 8 (2002) 943–959.
- 826 [31] A. A. Gaivoronski, G. Pflug, Value-at-risk in portfolio optimization: Prop-
827 erties and computational approach 7 (2) (2005) 1–31.
- 828 [32] C. Reano, R. Mayo, E. Quintana-Orti, F. Silla, J. Duato, A. Pena, Influence
829 of InfiniBand FDR on the performance of remote GPU virtualization, in:
830 IEEE International Conference on Cluster Computing, 2013, pp. 1–8.

831 **Figure Captions**

832 **Figure 1:** Execution time of the financial application on multiple GPUs

833 **Figure 2:** Distributed acceleration architecture facilitated by rCUDA

834 **Figure 3:** rCUDA client and server software/hardware stack

835 **Figure 4:** Communication sequence between a client and the rCUDA server
836 daemon

837 **Figure 5:** Comparison of bandwidth for pinned memory and pageable memory
838 of rCUDA, DS-CUDA and gVirtuS using CUDA as a baseline reference (DS-
839 CUDA does not support pinned memory)

840 Figure 5(a): Host pinned memory to device memory

841 Figure 5(b): Device memory to host pinned memory

842 Figure 5(c): Host pageable memory to device memory

843 Figure 5(d): Device memory to host pageable memory

844 **Figure 6:** Computation and data transfer times for the financial risk application
845 when executed on single and multiple GPUs with CUDA

846 **Figure 7:** Amount of data transferred during the execution of the financial risk
847 application

848 Figure 7(a): Data transferred to each GPU

849 Figure 7(b): Total data transferred to all GPUs

850 **Figure 8:** Attained bandwidth when concurrent data transfers to GPUs are
851 performed. Source data is located in the same memory bank.

852 **Figure 9:** Scalability of the financial risk application when executed with
853 rCUDA

854 Figure 9(a): On QDR InfiniBand

855 Figure 9(b): On FDR InfiniBand

856 **Figure 10:** Bandwidth attained for multiple data transfers concurrently to
857 different remote GPUs using rCUDA

858 Figure 10(a): On QDR InfiniBand

859 Figure 10(b): On FDR InfiniBand

860 **Figure 11:** Communication approaches for transferring data to GPUs.

861 Figure 11(a): Concurrent data transfers

862 Figure 11(b): Sequential data transfers

863 **Figure 12:** GPU utilisation, power and energy consumption of concurrent and
864 sequential data transfers to GPUs considered in Figure 11

865 Figure 12(a): Concurrent

866 Figure 12(b): Sequential

867 **Figure 13:** Sequential data copies with several vGPUs per GPU.

868 Figure 13(a): 2 vGPUs per GPU

869 Figure 13(b): 4 vGPUs per GPU

870 **Figure 14:** GPU utilisation, power and energy consumption of the multi-
871 tenancy approach considered in Figure 13.

872 Figure 14(a): 2 vGPUs per pGPU

873 Figure 14(b): 4 vGPUs per pGPU

874 **Figure 15:** Application performance for different combinations of pGPUs and
875 vGPUs using QDR InfiniBand.

876 **Figure 16:** Application performance for different combinations of pGPUs and
877 vGPUs using FDR InfiniBand.

878 **Figure 17:** Results from performance model for QDR InfiniBand

879 Figure 17(a): 1 to 4 pGPUs

880 Figure 17(b): 5 to 8 pGPUs

881 Figure 17(c): 9 to 12 pGPUs

882 Figure 17(d): 13 to 16 pGPUs

883 **Figure 18:** Results from performance model for FDR InfiniBand

884 Figure 18(a): 1 to 4 pGPUs

885 Figure 18(b): 5 to 8 pGPUs

886 Figure 18(c): 9 to 12 pGPUs

887 Figure 18(d): 13 to 16 pGPUs

888 **Figure 19:** Results from energy model for QDR InfiniBand

889 Figure 19(a): 1 to 4 pGPUs

890 Figure 19(b): 5 to 8 pGPUs

891 Figure 19(c): 9 to 12 pGPUs

892 Figure 19(d): 13 to 16 pGPUs

893 **Figure 20:** Results from energy model for FDR InfiniBand

894 Figure 20(a): 1 to 4 pGPUs

895 Figure 20(b): 5 to 8 pGPUs

896 Figure 20(c): 9 to 12 pGPUs

897 Figure 20(d): 13 to 16 pGPUs

898 **Figure 21:** Combined space of energy and execution time using QDR Infini-
899 Band

900 Figure 21(a): 1 to 8 pGPUs

901 Figure 21(b): 9 to 16 pGPUs

902 **Figure 22:** Combined space of energy and execution time using FDR Infini-
903 Band

904 Figure 22(a): 1 to 8 pGPUs

905 Figure 22(b): 9 to 16 pGPUs

906 **Table Captions**

907 **Table 1:** Scalability of the financial risk application when executed using CUDA

908 **Table 2:** Time in seconds for GPU memory allocation and data transfer tasks
909 of the financial risk application