

# Design Space Exploration of On-chip Ring Interconnection for a CPU-GPU Architecture

Jaekyu Lee<sup>1</sup>   Si Li<sup>2</sup>   Hyesoon Kim<sup>1</sup>   Sudhakar Yalamanchili<sup>2</sup>

<sup>1</sup>School of Computer Science  
Georgia Institute of Technology  
{jaekyu.lee, hyesoon}@cc.gatech.edu

<sup>2</sup>School of Electrical and Computer Engineering  
Georgia Institute of Technology  
sli@gatech.edu   sudha@ece.gatech.edu

## Abstract

Future chip multiprocessors (CMP) will only grow in core count and diversity in terms of frequency, power consumption, and resource distribution. Incorporating a GPU architecture into CMP, which is more efficient with certain types of applications, is the next stage in this trend. This heterogeneous mix of architectures will use an on-chip interconnection to access shared resources such as last-level cache tiles and memory controllers. The configuration of this on-chip network will likely have a significant impact on resource distribution, fairness, and overall performance.

The heterogeneity of this architecture inevitably exerts different pressures on the interconnection due to the differing characteristics and requirements of applications running on CPU and GPU cores. CPU applications are sensitive to latency, while GPGPU applications require massive bandwidth. This is due to the difference in the thread-level parallelism of the two architectures. GPUs use more threads to hide the effect of memory latency but require massive bandwidth to supply those threads. On the other hand, CPU cores typically running only one or two threads concurrently are very sensitive to latency.

This study surveys the impact and behavior of the interconnection network when CPU and GPGPU applications run simultaneously. This will shed light on other architectural interconnection studies on CPU-GPU heterogeneous architectures.

## 1. Introduction

The demand for more computational power never ends. Traditionally, growth in computational power was carried out by ever increasing clock frequency until the power wall was hit. To circumvent this barrier, future chip multiprocessors (CMPs) will only grow in core count and diversity in terms of frequency, power consumption, and resource distribution. The next stage in this trend involves heterogeneous architectures where each architecture is more power efficient at a subset of tasks. A general-purpose GPU (GPGPU) is one such example that is more power efficient for tasks involving massive thread-level parallelism. Incorporating a GPU architecture into CMPs is the next logical step. Recent examples include Intel's Sandy Bridge [21], AMD's Fusion [3], and NVIDIA's Denver [30] project. In these architectures, various on-chip resources are shared by CPU and GPU cores, including the last-level cache, memory controllers, and DRAM. Access to these shared resources is controlled by the on-chip interconnection, which has a significant impact on resource distribution, fairness, and overall performance.

To improve network and overall performance, many researchers have proposed a variety of mechanisms involving topologies, adaptive routing, scheduling, and arbitration policies. Many different topologies [6, 10, 13, 23, 35] have been proposed to improve performance. They tend to focus on either all CPU-based architectures or on specialized SoCs operating running a narrow spectrum of applications. Various adaptive routing algorithms are introduced to improve performance [18, 19, 25, 27]. A significant amount of work is proposed in router arbitration policies [11, 12]. Recently, proposals on heterogeneous interconnection configurations have been introduced [16, 26].

On-chip CPU-GPU heterogeneous architectures as well as their interconnections, however, are not as well studied. While we anticipate that they will have characteristics similar to more conventional CMP networks, we also expect additional complexities involving resource sharing mechanisms caused by the opposing memory demands exerted by applications running on the two architectures. CPU and GPU architectures possess fundamentally diametric network demands. CPU cores have high context switch overhead and rely on instruction parallelism, large caches, and speculative mechanisms to achieve high performance in serial execution. Since CPUs usually operate on a very small number of threads, when one is stalled on a memory access, it incurs a large penalty until that access is satisfied. At the opposite end of the spectrum, GPUs target data-parallel applications to achieve high throughput. They exchange large caches and other power-consuming mechanisms for more processing elements (PE) to execute on multiple data sets in parallel. Although a hardware-managed cache exists, GPUs mainly leverage the zero overhead, single-cycle context switch capability to remove latency introduced by long latency instructions and their dependencies. When a thread is blocked, the instruction scheduler context switches to the next available thread or group of threads. Due to the often massive number of potential threads waiting for execution, GPUs can execute many memory instructions concurrently and also in parallel, thereby achieving higher bandwidth requirements than CPUs.

In this paper, we evaluate the NoC behavior of this CPU-GPU heterogeneous architecture. However, we limit our study to the ring network. Although the ring network is not scalable with many cores, it is still relevant because most commercial processors currently employ a ring network with a reasonable number of CPU and GPU cores. Based on a network characterization of the Cell processor [2] and Intel's Larrabee [36] we believe that the ring network will be used at least for the next few years until the number of cores breaches a threshold of 10 or 12 cores.

In this study, we seek answers to the following questions: 1) How does the ring interconnection behave in CPU-GPU heterogeneous workloads? and 2) What is the best ring router configuration in heterogeneous workloads. We first study the impact of a variety of network resources and mechanisms on the system performance of CPU and GPGPU applications running separately, including the number of virtual channels, link width, link latency, and different placements. Then, we evaluate the resource sharing in the interconnection when both applications are running concurrently. In particular, we study virtual and physical channel partitioning between CPU and GPU cores, heterogeneous link configuration for each router, arbitrations, routing algorithms, and placements.

Based on the findings from empirical studies, we suggest an improved ring interconnection network in CPU-GPU heterogeneous architectures that improves performance by 22%, 19%, and 16% for one-CPU/one-GPU, two-CPU/one-GPU, and four-CPU/one-GPU workloads, respectively. We believe our studies will lead to further architectural studies in this area of on-chip interconnection for a CPU-GPU heterogeneous architecture.<sup>1</sup>

## 2. Background

### 2.1 CPU-GPU Heterogeneous Architecture

Intel’s Sandy Bridge [21], which was released in 2011, is the first commercial CPU-GPU heterogeneous architecture product. In the Sandy Bridge microarchitecture, the CPU and GPU cores share the last-level (L3) cache and memory controllers connected by a 256-bit wide ring network. Figure 1 shows the die picture of Sandy Bridge. Note that GPU cores can execute graphics as well as data-parallel GPGPU applications.

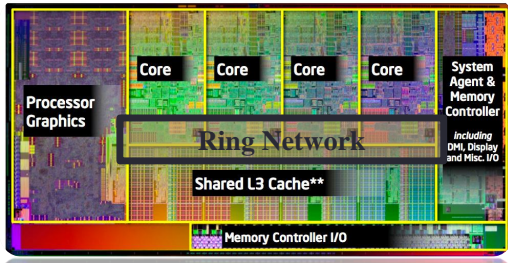


Figure 1. The die map of Intel’s Sandy Bridge.

AMD also released a different heterogeneous design, the Fusion architecture [3], that integrates more powerful GPU cores. Although CPU and GPU cores share caches and memory controllers, all communications are performed through the north bridge, not the generic interconnection. NVIDIA also disclosed its Denver project [30], which integrates ARM-based cores into the GPU. However, few details about this project have been released at the time of this writing.

Although the current GPUs in on-chip heterogeneous architectures are not as powerful as today’s high-performing GPUs, we project that more SIMD cores will be integrated in future generations. Therefore, we model our baseline heterogeneous architecture such that it has both high-performance CPU and GPU cores on-chip. Section 4.1 details the configuration of our baseline architecture.

### 2.2 A Comparison of CPU and GPU Cores

This section examines the differences between CPU and GPU cores. Modern high-performance CPU cores are typically based

on N-wide superscalar out-of-order cores. To reduce the penalty of the branch instructions, novel and often power-intensive branch prediction mechanisms are used. Large private caches (L1 and L2) are often employed to avoid long-latency access to off-chip memory. These cores are ideal for the serial execution of a small number of threads (1-4 way Simultaneous Multi-Threading (SMT)), so they have limited thread-level parallelism (TLP).

On the other hand, GPUs pack more processing elements in each core. Each GPU core is an in-order SIMD processor. Multiple threads execute the same instruction with different data sets per core. When branch directions within a batch of threads are diverged, the execution of each branch path is serialized. Currently, no branch prediction mechanism exists to reduce branch latency. The core context switches to other batches of threads until the branch is resolved. Similarly, to hide memory latencies, GPU cores utilize massive multi-threading. When a thread is stalled due to the memory instruction, the execution is switched to other available threads. Since GPU cores are designed to pay zero context overhead, this can happen on every instruction issued. Due to this high level of TLP, GPUs coalesce memory requests to reduce memory traffic when possible [29]. Additionally, GPUs are afforded single-cycle access to massive register files. Some GPGPU applications have frequent scatter-gather memory operations that hurt performance due to unaligned memory accesses. To mitigate this costly operation, GPUs often have special hardware to support the scatter-gather operation. Table 1 summarizes the differences between CPU and GPU cores.

Table 1. Comparison between CPU and GPU cores.

	CPU	GPU
Core	OOO superscalar	in-order SIMD
Branch Predictor	2-level	N/A
TLP	1-4 way SMT	abundant
Memory	Latency-limited	Bandwidth-limited
Latency tolerance	Caching	Caching, multi-threading
Miscellany		Scatter-gather operation

### 2.3 Network-on-Chip (NoC) Router Microarchitecture

In this section, we provide a brief background of the structure of an NoC router microarchitecture.

#### 2.3.1 Structure of Router Architecture

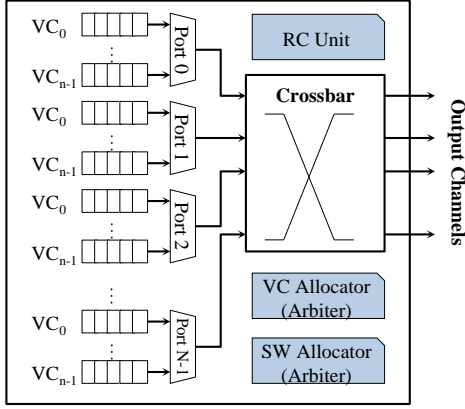
A router has  $N$  *input and output ports*. For example, a bi-directional ring network has three (local, left, and right) input and output ports. A 2D-mesh has five (local, north, south, east, and west) ports. Each input port has  $M$  *input buffers* or *virtual channels (VC)*. New packets from the local network interface are inserted into the VCs, and packets from other routers are inserted into their respective buffers. When a new packet (or flit) is inserted, the *routing computation* unit decides the output port to the next router. The *virtual channel allocator (arbiter)* assigns a virtual channel on that output port. The *switch arbitration* unit controls the *crossbar* to move a flit to the assigned output port. Each flit traverses a link from the output port of one router to the input port of the next. Figure 2 shows a diagram of the router microarchitecture.

#### 2.3.2 Pipeline stages

The flow of packets is pipelined in the NoC router. This is modeled in a five-stage pipeline model.

- Input buffering (IB): Flits received over a link or the source node are inserted into the buffer.
- Route computation (RC): Using the information in the header flit, the output port is determined.

<sup>1</sup>In this paper, we interchangeably use the term on-chip interconnection network and the network on chip (NoC).



**Figure 2.** The architecture of an NoC router.

- Virtual channel allocation (VCA): Using the output port information, a downstream virtual channel with available credits acquires a packet.
- Switch allocation (SA): To traverse to the output port, a packet needs an exclusive grant to access the cross bar from its input buffer to the output port.
- Switch traversal (ST): Once a switch is allocated to a packet, it can traverse to the output port over the crossbar.
- Link traversal (LT): A flit is moved to the next router through the link.

### 2.3.3 The Arbitration

As explained in Section 2.3.2, packets from the same or different input buffers compete against each other for a grant to a virtual channel or a switch to the output port. Simple policies are used to arbitrate between these packets: 1) round-robin: the winning virtual channel is chosen in a sequential manner and 2) oldest-first: all packets in the router are searched and the oldest request is scheduled. More sophisticated proposals in the literature are described in Section 6.

## 3. Problems and Design Space Exploration in NoCs of CPU-GPU Heterogeneous Architecture

This section describes the potential problems in designing the on-chip interconnection network in a CPU-GPU heterogeneous architecture.

### 3.1 Routing Algorithm

NoC routers typically employ a simple static routing algorithm to minimize latency and complexity. For example, x-y or shortest-distance algorithms are widely used. However, this may result in link congestion in the heterogeneous architecture. For example, in Figure 1, the GPGPU packets are not likely to use the upper link since the lower link offers the shortest distance from the GPU cores to the L3. The lower link is also used between the L3 and the memory controllers. Therefore, only CPU packets use the upper link, which is possibly under-utilized. While studies on other algorithms show improved network performance, they are limited to traffic generated by specialized or CPU-only applications [18, 19, 25].

### 3.2 Resource Contention and Partitioning

CPU and GPU packets compete to acquire resources in various places, especially virtual and physical channels. When the

resources are naively shared by both kinds of cores, higher-demanding cores will acquire the most resources, which are GPU cores. This is the same problem found in the LRU cache-replacement policy in the shared cache. To solve this problem, many researchers have proposed various static and dynamic cache partitioning mechanisms [34, 37]. Similarly, partitioning mechanisms can be applied to on-chip virtual and physical channels. As explained in Section 2.3.1, each port has multiple virtual channels. We can partition these virtual channels to each application. Similarly, if multiple physical channels exist, we can dedicate some channels to CPU cores and the other channels to GPU cores. If the interference exhibited by other applications is significant, resource partitioning would prevent interference and improve performance. However, this can lead to resource under-utilization if partitioning is not balanced with demand. Therefore, partitioning should be carefully applied to on-chip network resources.

### 3.3 Arbitration Policy

As described in Section 2.3.3, multiple arbiters exist in each router to coordinate packets from different ports. In a CPU-GPU heterogeneous architecture, due to the different network demands, arbitration between CPU and GPU packets is a non-trivial problem. At first glance, statically giving higher priority to CPU applications appears to be a reasonable solution since CPU applications are more latency sensitive. However, when CPU and GPGPU applications are both bandwidth-intensive, CPUs may be robbed of their fair share of the bandwidth. Therefore, the arbitration policy should also be carefully applied.

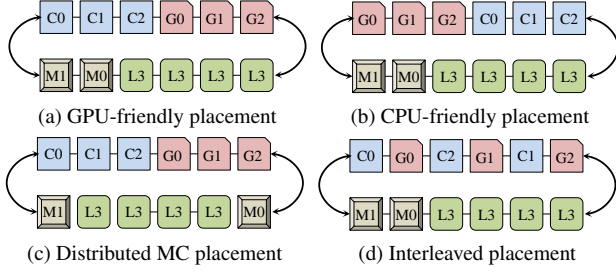
### 3.4 Homogeneous or Heterogeneous Link Configuration

A homogeneous router configuration has the practical benefit of easier implementation. If all NoC routers are identical, each router module can be duplicated with little or no individual adjustment. Since the requirements of CPU and GPU cores are very different, routers may require higher bandwidth interconnection in terms of the link width or larger buffers to effectively handle traffic from both applications. However, this may result in under-utilization of resources in a certain core. For example, if a wider link width is used, GPGPU applications may directly benefit from more bandwidth capability, but CPU applications may not because they do not require such a high bandwidth. Therefore, the utilization of CPU links will be low. A heterogeneous link configuration may work better in this situation but requires more complex implementation and may not perform as well in some bandwidth-intense situations. However, a heterogeneous configuration will require more design and implementation efforts compared to the homogeneous network. We leave this discussion to future work since this is beyond the scope of our study.

### 3.5 Placement

As explained in Section 3.1, any placement of these components – CPU, GPU, L3, MC – results in unbalanced utilization of on-chip interconnection resources for some scenarios or under-utilization for all situations. Figure 3 shows four possible examples of placement in the ring network. Among these examples, the placement of memory controllers (Figure 3 (d)) in many-core CMPs is studied by Abts et al. [1].

The first two examples are GPU- and CPU-friendly placements. Since all cache misses from a core need to reach L3 caches first, the distance between a core and a target L3 node may have a major impact on performance. Figure 3 (a) shows that the distance between the GPU cores and the L3 caches is shorter than the distance from the CPU cores. If there are more frequent accesses from the GPU cores to the L3 caches, this placement results in



**Figure 3.** Placement examples in the ring network.

better system performance. For the same reason, Figure 3 (b) is more beneficial for CPU applications.

In another configuration, each memory controller is placed at the end of the die in Figure 3 (c). If we can map the disjoint address range of the physical memory for the two types of cores (by the operating system), we can balance the link usage and the latency between each core to the L3 cache and traffic to the memory controllers will be reduced. This setup could effectively divide the chip into two halves, which would be the most beneficial when each half requires the same amount of bandwidth, but would otherwise result in a major imbalance in resource distribution.

Figure 3 (d) shows an interleaved placement, where CPU and GPU cores are interleaved. The possible benefit of this design is that it can balance the traffic in each direction from each application. In other designs, the traffic from each type of application tends to head in the same direction due to the shortest-distance routing algorithm. When too much traffic is headed in one direction, the application will slow down.

Although some placements are not practical in an actual implementation, this is beyond the scope of our study. We leave this discussion to future work.

## 4. Evaluation Methodology

### 4.1 Simulator

We use MacSim [17] for our simulations. For all simulations, we repeat early terminated applications until all applications have finished at least once. This is to model the resource contention uniformly across the duration of simulation, which is similar to the work in [22, 24, 34, 39].

**Table 2.** Processor configuration.

CPU	1-4 cores, 3.5GHz, 4-wide, out-of-order (OOO)
	gshare branch predictor
	8-way, 32KB L1 D/I cache, 2-cycle
	8-way 256KB L2 cache, 8-cycle
GPU	4 cores, 1.5GHz, in-order, 2-wide SIMD
	8-way, 32KB L1 D (2c), 4-way 4KB L1 I (1c)
	16KB s/w managed cache
L3 Cache	4 tiles (each tile: 32-way, 2MB), 64B line
Memory	DDR3-1333, 2 MCs (each 8 banks, 2 channels)
Controller	41.6GB/s BW, 2KB row buffer, FR-FCFS scheduler

Table 2 shows the processor configuration. We model our baseline CPU similarly to Intel’s Sandy Bridge [21] with GPU cores similar to NVIDIA Fermi [31].

Table 3 shows the configuration of the NoC router. To avoid the deadlock configuration, we use bubble routing [33]. We set the GPU-friendly placement in Figure 3 (a) as the baseline configuration, which is a configuration similar to Intel’s Sandy Bridge.

**Table 3.** NoC configuration.

Topology	Bi-directional ring network
Pipeline	5-stage (IB, RC, VCA, SA/ST, LT)
# VCs	4 per port (4-flit buffer)
# ports	3 (Local, Left, Right) per router
Link width	128 bits (16 B)
Link latency	2 cycles
Routing	Shortest distance
Flow control	credit-based, bubble routing [33]

### 4.2 Benchmarks

We use 29 SPEC 2006 CPU benchmarks and 31 CUDA GPGPU benchmarks from publicly available suites, including Nvidia CUDA SDK, Rodinia [9], Parboil [20], and ERCBench [8]. For the CPU workloads, Pinpoint [32] was used to select a representative simulation region with the reference input set. Most GPGPU applications run until completion.

Tables 4 and 5 show the characteristic of the evaluated network-intensive CPU and GPGPU benchmarks, respectively. The two metrics, MPKI (Miss Per Kilo Instruction) and IPKC ((Packet) Injection Per Kilo Cycles), are strongly correlated because cache misses will introduce more traffic into the network. For this reason we only consider IPKC when categorizing the benchmarks into two groups: network-intensive (Group N, IPKC greater than 10 (CPU), 37.5 (GPU)) and non-network-intensive (Group C). We demonstrate that GPGPU benchmarks in general generate higher network traffic than CPU benchmarks. This is mainly due to the high number of concurrent threads running in a GPU core, which serves to hide memory latency by overlapping memory accesses. Subsequently, GPU cores generate a higher intensity of network traffic.

**Table 4.** CPU benchmark characteristics.

Benchmark	Suite	MPKI/Core	IPKC/Core
bzip2	Int	0.4	10.2
gcc	Int	1.2	10.9
mcf	Int	43.6	29.9
libquantum	Int	26.8	20.8
omnetpp	Int	10.2	21.5
astar	Int	7.6	17.8
bwaves	FP	22.3	61.6
milc	FP	31.1	49.4
zeusmp	FP	5.8	23.7
cactusADM	FP	6.2	31.6
leslie3d	FP	24.7	71.6
soplex	FP	13.9	34.4
GemsFDTD	FP	18.9	51.2
lbm	FP	18.0	71.1
wrf	FP	15.2	55.2
sphinx3	FP	0.6	31.2

Section 5.1 evaluates the network behavior of N-type applications running individually in the arrangement of Figure 3 (a). Each application was tested by running one copy on compatible cores of the network. For the heterogeneous configuration experiments, we randomly choose combinations of N-type CPU and GPGPU applications. Table 6 describes the workload we evaluated.

### 4.3 Evaluation Metric

We use the geometric mean (Eq. 1) of the speedup of each application (Eq. 2) as the main evaluation metric.

$$speedup = geomean(speedup_{(0 \text{ to } n-1)}) \quad (1)$$

$$speedup_i = \frac{IPC_i}{IPC_i^{baseline}} \quad (2)$$

**Table 5.** GPGPU benchmark characteristics (MPKI and IPKC is the average of each core).

Benchmark	Suite	MPKI/Core	IPKC/Core
BlackS	SDK	25.6	153.2
ConvS	SDK	0.0	39.6
Dct8x8	SDK	0.1	42.7
Histog	SDK	5.4	39.8
ImageD	SDK	0.1	62.7
MonteC	SDK	0.0	47.1
Reduct	SDK	123.5	164.5
SobolQ	SDK	22.7	152.1
Scalar	SDK	0.4	80.7
backPr	Rodinia	3.4	39.7
cfid	Rodinia	323.9	112.9
neares	Rodinia	0.1	81.7
bfs	Rodinia	10.6	95.4
needle	Rodinia	10.2	65.0
SHA1	ERCBench	4.5	47.1
fft	parboil	0.2	56.4
stencil	parboil	16.7	134.3

**Table 6.** Workload description.

Workload	# CPU	# GPU	# combinations	Reference
W-1CPU	1	1	13	Section 5.2
W-2CPU	2	1	10	Section 5.3
W-4CPU	4	1	10	Section 5.3

We occasionally use the weighted speedup metric defined in Eq. 3.

$$weighted\_speedup = \sum_i \frac{IPC_i^{shared}}{IPC_i^{alone}} \quad (3)$$

## 5. Results

In this section, we first evaluate each application running singly in the system. Then we evaluate combinations of applications in the following order: one-CPU/one-GPU, two-CPU/one-GPU, and four-CPU/one-GPU applications. In each evaluation, we measure the impact of the number of virtual channels, link width, physical channel partition, link latency, arbitration policy, and network placement configuration. Note that we do not evaluate different routing algorithms since the ring network has only two possibilities of the route decision (left or right).

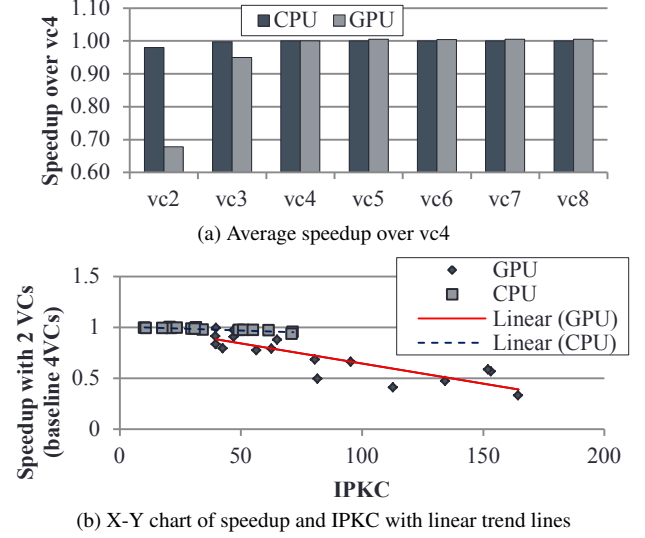
### 5.1 Single Application Analysis

In this section, we first evaluate each application in isolation (CPU or GPGPU application only) to analyze its characteristics without interference. Although previous studies exist [4, 5] on the effect of the on-chip network for GPGPU applications, we again present the data to correlate with our other experiments. Each CPU application is tested by running it on one CPU core while all other cores remain idle. Similarly, GPGPU applications run all GPU cores while CPU cores remain idle.

#### 5.1.1 Different Number of Virtual Channels

Figure 4 (a) shows the result when we vary the number of virtual channels from two to eight. All results are normalized to the baseline configuration (4 VCs). CPU applications suffered less performance loss than GPU applications with a small number of VCs. With two VCs, the maximum performance loss is 6.6% for the lbm benchmark, which is the most network-intensive benchmark in Table 4. On the other hand, six GPGPU benchmarks show more than a 40% degradation. When these applications are sharing the on-chip interconnection, we expect the inter-application interference to be a serious problem. We expect GPGPU applications to

have a considerably greater impact on other applications when running concurrently on the network. However, from this experiment, guaranteeing a small number of VCs (one or two) to CPU applications is sufficient to maintain CPU application performance. Section 5.2.2 evaluates the effect of virtual channel partitioning in a multi-application environment.



**Figure 4.** Evaluation of different # of virtual channels.

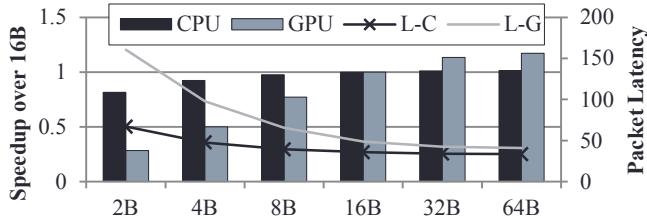
To seek the correlation between the slowdown of having smaller VCs and IPKC, we show an x-y chart in Figure 4 (b). Although both applications show linear regression lines, CPU applications do not show much variance. Applications with higher IPKC show greater degradation in general.

#### 5.1.2 Different Link Width

Figure 5 shows the impact of link widths for each type of application compared to the baseline of the 16B link width. Compared to the impact of varying the number of virtual channels, link width has a greater effect on the performance of both applications. GPGPU benchmarks generally show more sensitivity to link widths than CPU benchmarks, especially in Reduct, Scalar, cfd, and bfs benchmarks. This is potentially due to the large number of memory requests made in a batch by a huge number of concurrently running threads. GPGPU applications can take advantage of the wider links and prevent network congestion. On the other hand, CPU benchmarks do not show significant improvement with wider links (32B: 1.1%, 64B: 1.6% on average) compared to the baseline. However, lower network bandwidth in GPGPU applications induces a significant latency increase, thereby hurting performance excessively. The L-G line in Figure 5 shows that the average network latency of GPGPU applications is increased by 3.3 times, while CPU applications (L-C line) show a relatively small increase (1.87 times) with the 2B link. As a result, GPGPU applications show 72% slowdown over the 16B link, but CPU applications only show a 19% slowdown.

From this observation, we conclude that having a wider link is very helpful to GPGPU applications, but not to CPU applications. This confirms that GPGPU benchmarks are more bandwidth-limited, while CPU benchmarks are not. As a result, this observation leads us to study heterogeneous link configurations in Section 5.2.6.

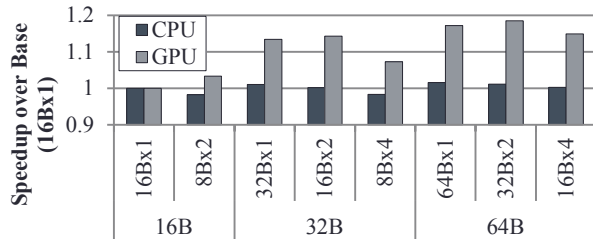




**Figure 5.** Evaluation of different link width (L-C: CPU network latency, L-G: GPU network latency).

### 5.1.3 Different Channel Configuration

In this section, we evaluate the impact of different widths and number of physical channels. The purpose of this evaluation is based on the intuition that a wider link can be beneficial to reduce the latency in the cases of larger packet requests, but multiple channels can be better utilized for more general traffic when smaller packets become a contributing portion of traffic.



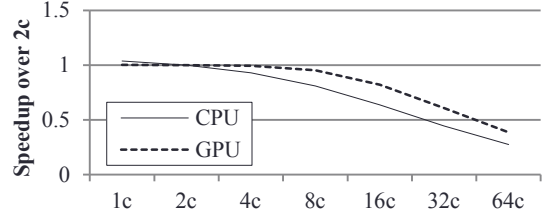
**Figure 6.** Evaluation of different physical channel.

Figure 6 shows the results. For CPU applications, wider links improve performance by only a small percentage. However, having smaller, multiple links degrades performance. On the other hand, GPGPU applications show improvement with two physical channels (16Bx2 and 32Bx2). However, having even more channels increases the packet latency significantly (2x longer latency for data packets) while not fully utilizing all channels. As a result, the benefit of wider links decreases. This finding indicates that the width and number of physical channels should be well balanced to support various applications that have different latency/bandwidth requirements.

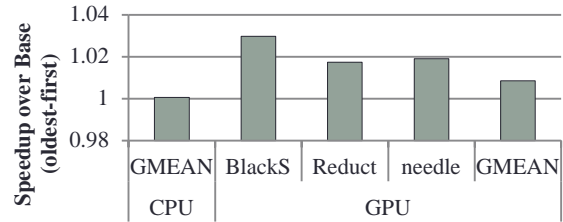
### 5.1.4 Different Link Latency

Figure 7 shows the different link latency evaluations with two-cycle as the baseline latency. We double the link latency from one to 64 cycles. Although both types of applications are sensitive to link latency, we observe that the degree of sensitivity is much higher in CPU benchmarks. Even with small changes (from two to eight cycles), most CPU benchmarks suffer from the increase in latency. We observe 7%, 19%, and 36% degradations on average for four-cycle, eight-cycle, and 16-cycle configurations, respectively. However, the performance of many GPGPU benchmarks did not degrade significantly until the 16-cycle latency configuration. On average the performance degradations are 5% and 18% for GPGPU benchmarks for eight-cycle and 16-cycle configurations, respectively.

As we explained in Section 2.1, CPU benchmarks are known to be latency-sensitive and GPGPU benchmarks are bandwidth-limited. This simulation result confirms this tendency.



**Figure 7.** Evaluation of different link latency.



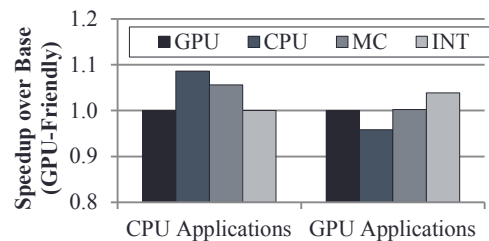
**Figure 8.** Round-robin arbitration (baseline: oldest-first).

### 5.1.5 Different Arbitration Policy

In this section, we evaluate two simple arbitration policies: round-robin arbitration and oldest-first policy while fixing other configurations the same.<sup>2</sup> Figure 8 shows the results. Not surprisingly, neither type of application shows significant changes. Since there are three input ports with few virtual channels in the ring network, not enough bottleneck is introduced by network complexity for an arbitration to resolve. CPU applications never show more than a 0.5% variance. Only a few GPGPU benchmarks, BlackS, Reduct, and needle, show more than a 1.5% variance.

### 5.1.6 The Effect of Network Placement Configuration

Figure 9 shows the effect of different placement policies. We evaluate four different placements, as shown in Figure 3.



**Figure 9.** Different placement results (MC: distributed memory controller, INT: interleaved).

For each type of application, placing the cores closer to the L3 caches improves performance by reducing the round trip latency between the cores and the memory system (the L3 caches and the memory controllers). The CPU-friendly placement improves the performance of CPU applications by more than 8% on average. On the other hand, the CPU-friendly placement degrades the performance of GPGPU applications by 4.1% compared to the GPU-friendly placement.

<sup>2</sup> We evaluate other static arbitration policies with heterogeneous workloads in Sections 5.2.3 and 5.3.4.

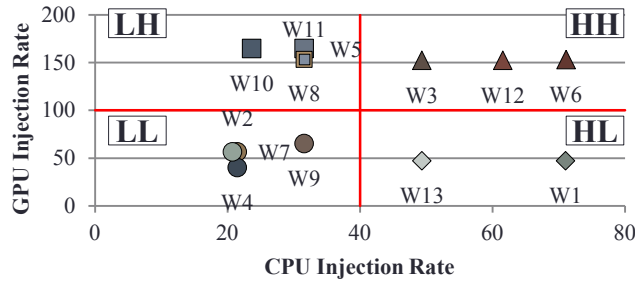
The MC placement improves the performance of CPU applications by 5.6%. The performance gain is partially from reducing the through-traffic of the memory controllers. In other configurations, the shortest distance from Core 0 to the L3 caches is through the memory controllers. Therefore, even though the destination of a packet is not a memory controller but an L3 cache node, all packets from and to the core must travel through these memory controllers. As a result, this path is always busy and tends to congest traffic. However, by placing the memory controllers on both sides of the cache nodes, through traffic is reduced.

On the other hand, there are cases in GPGPU applications where the MC placement either hurts or improves performance (from -7% to 6%). For negative cases, this placement introduces an extra node in the critical path between GPU cores and L3 cache nodes (M0 in Figure 3 (c)), resulting in extra latency and higher congestion (near M0) as traffic injected by the memory controller is also in the critical path. Hence, in some applications these negative effects offset the benefits of reducing congestion to/from the memory controllers explained earlier. Overall, for GPGPU applications, the benefit of the MC placement is washed out.

For the INT placement, note that CPU-friendly and INT placements are identical for CPU applications since we run only one application on Core 0 (C0), and this placement configuration does not impact the latency from this core to memory nodes. For GPGPU applications, we observe improvements of 3.9% on average due to an increase in path diversity, as the cores will take both sides of the ring instead of favoring only one side due to the shortest-distance routing algorithm described in Section 3.1.

## 5.2 Multiple-Application Experiments

In this section, we look at multi-application experiments to analyze the impact of inter-application interference. Each test randomly chooses one N-type CPU and GPGPU application to run them concurrently (the W-1CPU workloads in Table 6). Figure 10 shows an x-y chart of all W-1CPU workloads in terms of the IPKC characteristic. We split this into four regions (CPU Injection/GPU Injection, HH – High CPU and High GPU, HL, LH, and LL) based on the CPU and GPU injection rates.



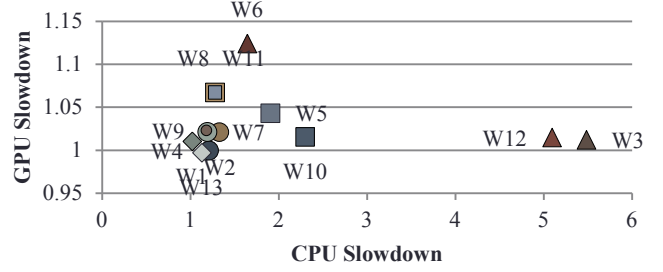
**Figure 10.** W-1CPU workload characterization (IPKC: injection per kilo cycles).

The following factors are measured for their impact: router buffer partitioning, arbitration policy, network placement, physical channel partitioning, and heterogeneous link configuration.

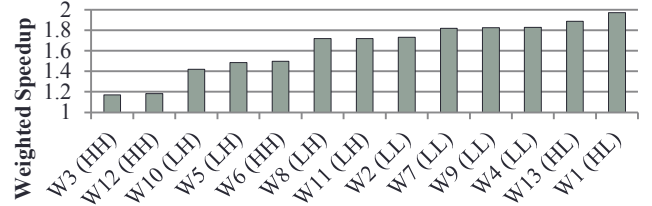
### 5.2.1 Interference with GPGPU applications

We first explain how applications interfere with each other. Figure 11 (a) shows an x-y chart of the slowdown of each application compared to when they are running alone and Figure 11 (b) shows

the weighted speedup<sup>3</sup> of each workload (sorted in ascending order).



(a) Slowdown X-Y chart ( $\triangle$ : HH,  $\diamond$ : HL,  $\square$ : LH,  $\circ$ : LL)



(b) Weighted speedup (Workload and Group Id in Figure 10))

**Figure 11.** Interference with GPGPU applications.

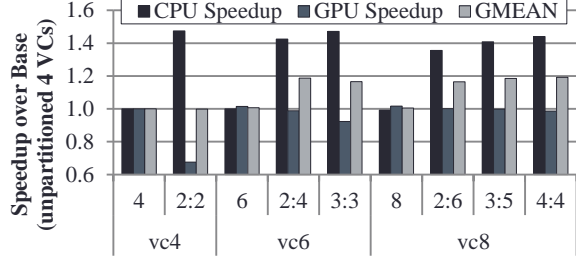
We can observe that significant interference is caused by GPGPU applications. In Figure 11 (a), most GPGPU applications do not show more than a 5% slowdown, while only three CPU applications show less than a 20% slowdown. In Figure 11 (b), GPGPU applications in poorly performing combinations (from the leftmost, W3, W12, W10, W5, and W6) all belong to the HH or LH region, while the best performing ones (from the rightmost, W1 and W13) belong to the HL region in Figure 10. We expect that although GPGPU applications will experience more interference with CPU applications as the number of concurrently running CPU applications increases, GPGPU applications are more likely to interfere with CPU applications.

### 5.2.2 Router buffer partitioning

As Tables 4 and 5 show, GPGPU benchmarks exhibit more frequent network injections. If we do not partition the router buffer space, the GPGPU packets will occupy more space and unnecessarily degrade the performance of CPU applications. This is a similar problem to the traditional LRU replacement policy in CMPs. Since the naive LRU policy does not partition cache space, higher-demanding applications will occupy more cache space. To solve this problem, many researchers have proposed static and dynamic cache partitioning mechanisms [34, 37]. Similarly, NoC virtual channels can be statically or dynamically partitioned; in other words, a few virtual channels can be dedicated to CPU packets and other channels to GPU packets. In this section, we evaluate the effect of virtual channel partitioning.

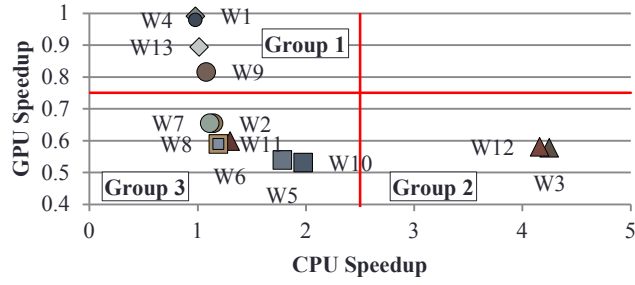
We compare the baseline unpartitioned virtual channels with various static partitioning configurations, as shown in Figure 12. We observe that increasing the number of virtual channels is not helpful. GPGPU packets occupy most of the available buffer space. Increasing to six and eight virtual channels results in minor improvements in GPGPU applications, while the performance of CPU

<sup>3</sup>Eq. 3. Higher is better. The ideal weighted speedup is 2 if no inter-application interference is exhibited.



**Figure 12.** Router buffer partitioning results (ex. 6 is 6 unpartitioned VCs; 2:4 is 2 CPU VCs, 4 GPU VCs).

applications stays the same. Compared to the unmanaged baseline, allocating at least some fixed number of virtual channels to CPU applications significantly improves performance by more than 35%, while GPGPU applications show 8% and 1% degradations with three and four dedicated virtual channels, respectively.



**Figure 13.** VC partitioning x-y chart (baseline: shared VC,  $\Delta$ : HH,  $\diamond$ : HL,  $\square$ : LH,  $\circ$ : LL in Figure 10).

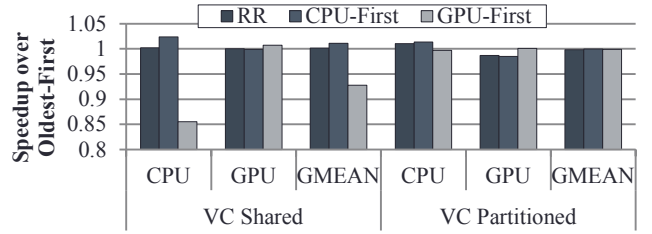
For a deeper analysis, we show an x-y chart of CPU and GPU speedups in the 2:2 (VC4) virtual channel partitioning configuration for all W-1CPU workloads in Figure 13. We pick this configuration because both CPU and GPGPU applications with this configuration show representative behavior from virtual channel partitioning. We observe that three groups exist in Figure 13: group 1 (top left region, W1, W4, W9, W13: moderate CPU and GPU speedup), group 2 (bottom right region, W3, W12: excessive CPU and GPU), and group 3 (others: moderate CPU and excessive GPU). In group 1, GPGPU applications have low injection rate (HL or LL group in Figure 10), so the interference by GPGPU applications is limited as well. Therefore, CPU and GPGPU applications effectively share VCs. As a result, both CPU and GPGPU applications do not show significant variances with VC partitioning. CPU applications in group 2 have higher injection rate (greater than 50 IPKC), so the performance of CPU applications is significantly improved with VC partitioning. Almost all CPU applications in group 3 have low injection rate (LH or LL group), so VC partitioning hurts the performance of GPGPU applications. In other configurations, we can observe that both CPU and GPGPU applications show similar trends, but we do not see as much degradation as in the 2:2 (VC4) configuration for GPGPU applications.

From this observation, we see that when CPU and GPGPU applications run concurrently in the shared on-chip network, guaranteeing the minimum buffer space to CPU applications would improve the overall system performance. However, providing additional buffer space does not result in additional performance increases.

### 5.2.3 Arbitration Policy

In heterogeneous architectures, CPU and GPU packets compete for buffer space and switch arbitration. To see the effect of arbitration policies, we try two static policies: CPU-first and GPU-first policies. To prevent starvation, we implement a form of batching similar to that in [12, 28]. We compare these two static policies and the round-robin policy with the oldest-first arbitration.

Figure 14 shows the results using the shared and partitioned virtual channel configurations. First, as seen in Section 5.1.5, there is no significant difference between oldest-first and round-robin policies, regardless of virtual channel configuration. However, for the two static policies, the VC configuration affects the result significantly. Since VC partitioning guarantees the minimal service for each type, the effect of different policies decreases. As a result, the three policies behave similarly (less than 1% delta) with VC partitioning.

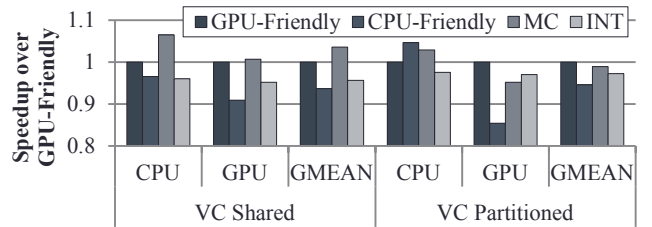


**Figure 14.** Different arbitration policy results (CPU and GPU indicate the speedup of each application).

However, with the shared VCs, the CPU-first policy slightly improves the performance of CPU applications without degrading that of GPGPU applications. On the other hand, the GPU-first policy degrades the performance of CPU applications by more than 15% on average while improving GPGPU performance by only 0.5%. Looking at the geometric mean, CPU-first improves performance by 1.5%, but GPU-first degrades it by 7.3% on average. This experiment again confirms that CPU packets are latency-sensitive, so we need to prioritize CPU packets.

### 5.2.4 Network Placement Configuration

We evaluate different placement policies on heterogeneous workloads, as shown in Figure 15. We again perform experiments with different VC configurations. With the shared VC, even CPU-friendly placement degrades CPU applications. This is because the lengthened distance from the GPU cores to both L3 caches and memory controllers increases system-level traffic congestion. As a result, each CPU packet is penalized by this congestion. The MC placement slightly improves the performance of both applications (3.5% on average), while the INT placement degrades both (-4.4% on average).



**Figure 15.** Different placement results for heterogeneous workloads (MC: distributed memory controller, INT: interleaved).

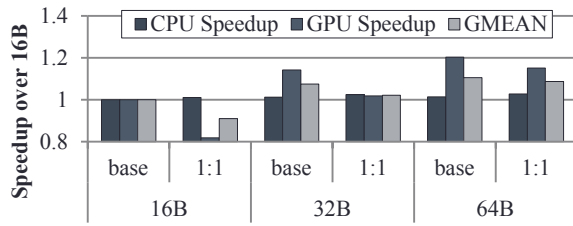


However, with the partitioned VC, we observe the expected behavior. The CPU-friendly placement shows better performance for the CPU applications (4.6%), but it degrades performance of the GPGPU applications even more (-14.6%). This degrades overall performance. The MC placement improves CPU application slightly by partially reducing the distance to the memory controllers, but it worsens GPGPU applications. The INT placement degrades the performance of both applications.

From this experiment, we observe that GPGPU applications have more impact on the network. The overall performance gain can be acquired by not penalizing GPGPU applications.

### 5.2.5 Physical Channel Partitioning

In this section, we evaluate physical channel partitioning. Similar to router buffer partitioning, if multiple physical channels exist in the router, we can partition channels to each type of application.

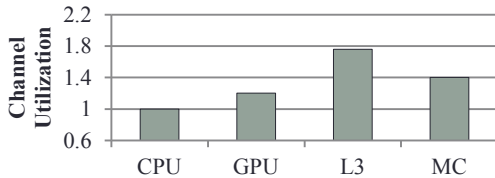


**Figure 16.** Physical channel partitioning results (16B-base: unpartitioned 1 channel, 16B-1:1: 2x8B channels and one channel is dedicated for CPU and the other is for GPGPU application).

Figure 16 shows the results. Similar to VC partitioning, the CPU applications benefit from a dedicated channel. However, a significant performance loss in GPGPU applications results. This indicates that GPGPU applications require a wider physical channel than CPU applications. When the channel is partitioned, we observe that the utilization of the GPU channel is slightly increased while that of the CPU shows very low utilization (around 5%). Therefore, to obtain better channel utilization for GPGPU applications, a wider channel instead of multiple channels is more effective.

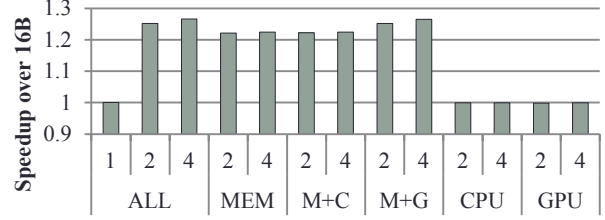
### 5.2.6 Heterogeneous Link Configuration

In this section, we evaluate heterogeneous link configurations. We first categorize routers into three groups: CPU-router (4 routers), GPU-router (4 routers), and memory-router (6 routers: 4 L3 cache and 2 memory controller routers). Our baseline uses one 128-bit (16B) link between each node. We vary the number of physical channels in each group.



**Figure 17.** Physical channel utilization (relative to CPU routers).

First, we show the physical channel utilization of each router normalized to the CPU routers in Figure 17. L3 and MC routers obviously have much more traffic, thereby utilizing channels more compared to processor routers (CPU and GPU). Even if the number



**Figure 18.** Heterogeneous link configuration results (ALL: all routers, MEM: vary memory router link only, and M+C: memory and CPU routers. 1, 2, or 4 in x-axis indicates the number of physical channels).

of concurrently running CPU applications increases, we expect that a similar trend will be observed.

Figure 18 shows the results of various link configurations. We first evaluate the homogeneous link configuration (ALL). Having more physical channels is always beneficial, but there is a diminishing return in performance after two channels. We observe performance improvements of 25% and 27% on average with two and four physical channels, respectively. Then, we evaluate heterogeneous configurations by varying each type of router (MEM, CPU, GPU, M+C, and M+G). While increasing the number of channels for CPU or GPU routers does not help improve performance, increasing memory-router channels has a direct effect on performance (22.1% and 22.4% with two and four physical channels, respectively). As in Figure 17, memory routers are mostly busy during the entire execution, but CPU and GPU routers are not. Because there are five different data flows in the ring network – 1) CPU to L3, 2) GPU to L3, 3) L3 to Memory controller (MC), 4) MC to L3, and 5) L3 to CPU (or GPU) – most traffic goes through the memory routers. By allocating wider links to only the memory routers rather than all routers, we can fully utilize these links without powering under-utilized links for the CPU and GPU routers. In addition, giving GPU routers additional channels, on top of wider memory router links, shows an additional boost in performance. Since GPUs are major sources of network traffic, without more channels between them, GPU-routers will become the new bottleneck.

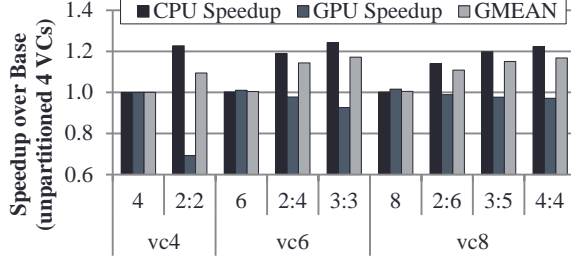
### 5.3 Multiple CPU Application Experiments

In this section, we evaluate multiple CPU applications and one GPGPU application running simultaneously. We repeat the same set of experiments as in Section 5.2.

#### 5.3.1 Router buffer partitioning

Figure 19 shows the results of VC partitioning in the W-4CPU workloads (4 CPUs + 1 GPGPU). Note that the W-2CPU (2 CPUs + 1 GPGPU) workloads are omitted, as they show roughly identical results as W-4CPU results. W-2CPU and W-4CPU workload data show very similar trends as W-1CPU (1 CPU + 1 GPGPU) experiments. Dedicated (at least a few) virtual channels to CPU applications are shown to be helpful, but the performance of GPGPU applications degrades with less than three VCs. This again indicates significant traffic injected by GPGPU applications and interference by GPGPU applications. VC partitioning will reduce this interference.

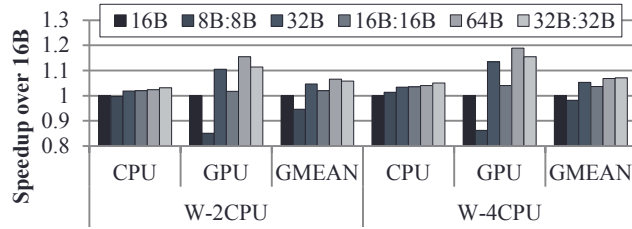
From the various VC-related experiments, our conclusions of the ideal VC configuration are that 1) VCs should be partitioned, especially for CPU applications; 2) However two to three VCs are sufficient; and 3) a GPGPU application requires at least three VCs, but having more does not help. Therefore, the ideal VC



**Figure 19.** Router buffer partitioning results (4 CPUs + 1 GPGPU workloads).

configuration will be five virtual channels: two are dedicated for CPU and the other three are for GPGPU applications.

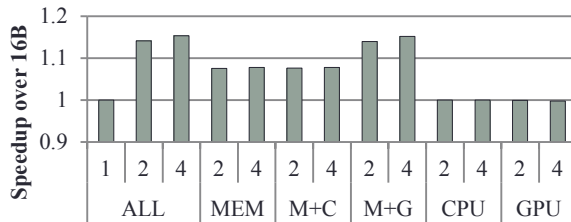
### 5.3.2 Physical Channel Partitioning



**Figure 20.** Physical channel partitioning results (16B-base: unpartitioned 1 channel, 16B-1:1: 2x8B channels and one channel is dedicated for CPU and the other is for GPGPU application).

Figure 20 shows the results for both W-2CPU and W-4CPU workloads. As the number of concurrently running CPU applications increases, the benefit of having a separate physical channel for CPU applications decreases since the channel itself is shared by more applications. However, GPGPU applications still suffer from narrower links, as they depend on bandwidth. As a result, we always observe a system performance degradation with a partitioned physical channel configuration.

### 5.3.3 Heterogeneous Link Configuration



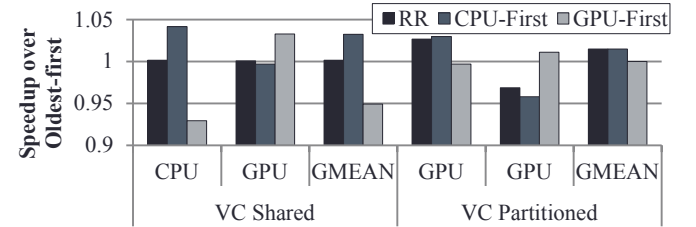
**Figure 21.** Heterogeneous link configuration results (ALL: all routers, MEM: vary memory router link only, and M+C: memory and CPU routers. 1, 2, or 4 in x-axis indicates the number of physical channels).

Figure 21 shows the results for the W-4CPU workloads. Again, the W-2CPU workloads are omitted due to their similarity to the W-4CPU workloads. Compared to the W-1CPU workloads, there is no significant difference. Multiple physical channels for the memory and the GPU routers proved to be beneficial. However, increasing

the number of physical channels for only CPU or GPU routers is not beneficial because half of the traffic is cache-miss requests. The size of request traffic is small (1 flit) and the request traffic would not occupy multiple physical channels. Therefore, having multiple physical channels does not improve performance. However, for GPU routers, the additional channels for the memory routers adds additional improvements.

### 5.3.4 Arbitration Policy

Figure 22 shows the results from different arbitration policies on both shared (CPU and GPU packets share all VCs) and partitioned VCs (2 VCs are dedicated for each type) for the W-4CPU workloads. Round-robin is comparable to the Oldest-First policy. Although the CPU-First policy consistently shows better performance, the benefit decreases with the partitioned VC. On average, CPU-First shows 3.2% and 1.5% improvements with the shared and partitioned VC, respectively. On the other hand, GPU-First degrades CPU applications in the shared VC configuration.



**Figure 22.** Different arbitration policy results.

Throughout multiple-workload evaluations (W-1CPU, W-2CPU, and W-4CPU), the effect of different arbitration policies is not so important, especially if virtual channels are partitioned to each type of application. This is not entirely unexpected since there is a smaller number of input and output ports in the ring network. We expect that the role of intelligent arbitration becomes significant with 2D topologies.

### 5.3.5 Placement

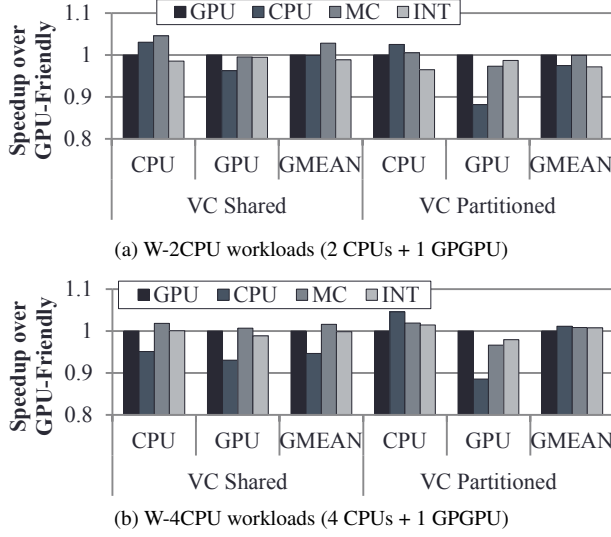
Figure 23 shows different placement evaluations for W-2CPU and W-4CPU workloads with different virtual channel configurations (shared and partitioned). With the shared VC, the MC placement shows the best results. This is mainly because traffic from/to the memory controllers is distributed to both sides of the chip, thereby reducing the congestion in those routers. With the partitioned VC, the effect of different placement policies is a bit reduced. GPU-friendly and MC placements perform the best.

### 5.4 Scalability

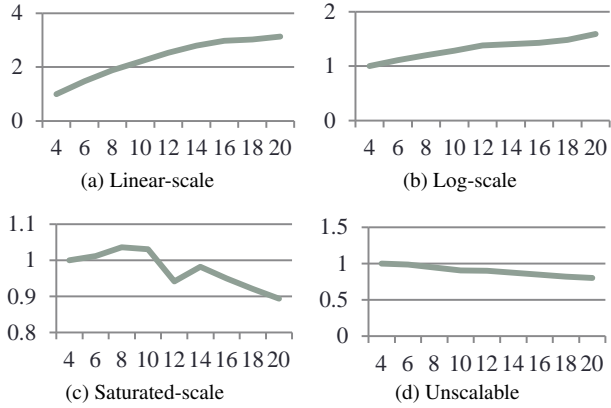
The last evaluation of our study is the scalability of the ring network. The ring network is known to handle a small number of nodes. In this section, we stress the ring network to show how many cores can be used. We run GPGPU applications and scale the number of cores from four to 20.

Figure 24 shows four different types of applications: linear-scale, log-scale, saturated-scale (performance saturated after N cores), and unscalable. Linear and log-scale applications are less network intensive and show performance improvement because increasing the number of cores does not saturate the interconnection network.

The other two types show that the performance flat lines at some point due to the congestion in the network. Especially, the performance of unscalable benchmarks degrades for even a small number of cores (Figure 24 (d)). These are the most network-intensive benchmarks in Table 5 (BlackS, Reduct, SobolQ, cfd, bfs).



**Figure 23.** Placement (MC: distributed memory controller, INT: interleaved).



**Figure 24.** Scalability test (x-axis: # cores, y-axis: speedup over 4-core).

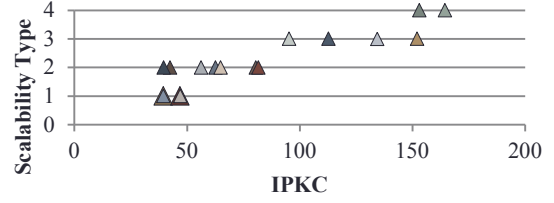
Figure 25 shows the correlation between the scalability of a ring network and the IPKC. As we can expect, in general, higher IPKC applications (especially above 100 IPKC) do not show good scalability.

In sum, the ring network is not scalable when it comes to significantly memory/network-intensive applications. However, we observe that the ring network is still a good candidate to handle moderate memory/network-intensive applications with a moderate number of cores until the on-chip interconnection bandwidth can handle them.

## 5.5 Summary of Findings

In this section, we summarize the findings from our evaluations.

1. As is widely known, CPU benchmarks are latency-sensitive and GPGPU benchmarks are bandwidth-limited. From the empirical data, we confirm that this applies to the on-chip interconnection under various circumstances.



**Figure 25.** Scalability (y-axis: scalability type, 1: linear, 2: log, 3: saturated, 4: unscalable)

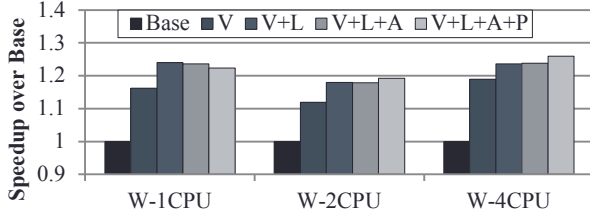
2. When CPU and GPGPU applications share the same NoC, a significant interference by GPGPU applications exists. To prevent this interference, we evaluate two resource partitioning schemes: virtual and physical channels. Having enough dedicated virtual channels for each type improves performance. On the other hand, multiple narrower physical channels rather than a single wider channel significantly degrades the performance of GPGPU applications. However, they do not improve CPU application performance due to lower link utilization.
3. A heterogeneous link configuration shows a promising result. Since the traffic is the highest in the last-level cache and the memory controllers (to/from CPU and GPU cores), those routers become traffic hotspots. By adding more physical channels, we can achieve similar performance compared to having multiple channels for all routers.
4. Prioritizing CPU request packets yields the best performance among other policies, but the benefit is very small. Generally, router arbitration has minimal effect on performance in a ring network.
5. GPU cores should be located close to the L3 caches and the memory controllers to avoid congestion, which eventually affects other applications as well. Moreover, separating memory controllers to both sides of the chip would reduce the traffic and improve performance in some cases.
6. The ring network is not scalable and even saturates with only a small number of cores (4 to 6) for some benchmarks. However, we observe that the ring network is still a good medium to handle a moderate number of cores until on-chip interconnection bandwidth can handle them.

From these findings, we suggest an optimal router option that is a combination of the best performance with minimal hardware resources for each individual experiment in Table 7, and Figure 26 shows the evaluation results. Please note that the suggested configuration requires more hardware resources (one more virtual channel per port in a router and more physical channels for memory and GPU routers) than the baseline. The detailed tradeoff between more hardware resources and performance/power improvements remains in future work.

**Table 7.** Putting it all together.

	Base	Suggested option
Virtual Channel	shared 4 VCs	partitioned 5 VCs (2: CPU, 3: GPU)
Physical Channel	1x16B	1x16B for cpu 2x16B for mem and gpu
Arbitration	Oldest-first	CPU-first
Placement	GPU-Friendly	MC or GPU-Friendly

Greater improvements mostly come from virtual channel partitioning and heterogeneous link configurations, while other optimizations are less critical. Our suggested router configuration improves performance by 22%, 19%, and 26% for W-1CPU, W-2CPU, and W-4CPU workloads, respectively.



**Figure 26.** Suggested on-chip interconnection configuration results (V: VC partitioning, L: Heterogeneous link configuration, A: Arbitration, P: Placement)

## 6. Related Work

### 6.1 On-chip Ring Network

Although the ring network has been extensively studied in the past, we limit our discussion to the on-chip ring network in this section. Bononi and Concer [6] studied and compared various on-chip network topologies, including the ring, in the SoC (System on Chip) domain. Bourduas and Zilic [7] proposed a hybrid ring/mesh on-chip network. A conventional 2D-mesh network has a large communication radius. To reduce the communication cost, the network is partitioned into several sub meshes and the ring connects these partitions. Ainsworth and Pinkston [2] performed a case study of the Cell Broadband Engine’s Element Interconnect Bus (EIB), which consists of four ring networks for data and a shared command bus that connects 12 elements.

### 6.2 CPU-GPU Heterogeneous Architecture Research

Since the CPU-GPU heterogeneous architecture was introduced recently, not many studies are available in the literature. In particular, the resource-sharing problem is not well discussed. To the best of our knowledge, the study by Lee and Kim [24] is the only study on resource sharing in this type of architecture. The authors recently studied the cache-sharing behaviors in a CPU-GPU heterogeneous architecture and proposed TLP-aware cache management schemes. However, some work on utilizing idle cores to boost performance has been done. Woo and Lee proposed Compass [38], which utilizes idle GPU cores for various prefetching algorithms in heterogeneous architectures.

### 6.3 Heterogeneous NoC Configuration

Much research has been done on heterogeneous NoCs involving many-core CPUs. Heterogeneous network configurations (HeteroNoC), by Mishra et al. [26], proposed asymmetric resource allocations (buffers and link bandwidth) to exploit non-uniform demand on a mesh topology. HeteroNoC showed that routers along the diagonals provided performance improvement over homogeneous resource distribution. Grot et al. [16] proposed Kilo-NOC, with shared resources isolated into QoS-enabled regions to minimize the network complexity. Kilo-NOC also reduces area and energy, in non-QoS regions by using a MECS- (Multidrop Express Channels) [15] based network with elastic buffer and novel virtual channel allocation that reduces VC buffer requirements by eight times over MECS with minimal latency impact.

### 6.4 NoC Prioritization

Application-aware prioritization [11] computes the network demand of applications at intervals by looking at a number of metrics such as private cache misses per instruction, average outstanding L1 misses in MSHRs, and average stall cycles per packet. This produces a ranking of an application, and all packets of one application are prioritized over another, resulting in a coarse granularity

of control. To prevent application starvation, a batching framework is implemented that prioritizes all packets of one time quantum over another, regardless of source application. Aéria [12] predicts the available latency (slack) of any packet by the number of outstanding L1 misses and prioritizes low-slack (critical) packets over packets with higher slack when they are within the same batching interval.

### 6.5 NoC Routing

Ma et al. [25] proposed destination-based adaptive routing (DBAR), a network with a novel congestion information network with a low wiring overhead. Like RCA [14], DBAR uses the virtual channel buffer status of nodes on the same dimension to route around congested paths. In addition, DBAR ignores nodes outside the potential path to eliminate interference and provides dynamic isolation from outside regions of the network.

Bufferless routing [27] showed substantial energy savings from removing input buffers by deflecting incoming packets from congested output ports. This routing algorithm managed performance similar to other buffered routing algorithms but only at low traffic.

## 7. Conclusion and Future Work

In this paper, we explore a broad design space in the on-chip ring interconnection network for a CPU-GPU heterogeneous architecture. We observe that this type of heterogeneous architecture has been adopted by major players in the industry and will become the mainstream processor type in subsequent generations. We observe that the interference exhibited by other applications, mostly by GPGPU applications, is significant and can be detrimental to system performance if not properly managed. We examine resource partitioning schemes for virtual and physical channels. Virtual channel partitioning improves performance, but physical channel partitioning degrades it because of link under-utilization. Heterogeneous link configurations, different arbitration policies, and placement configurations have been considered in this paper as well. The heterogeneous link configuration shows effectiveness, but other configurations have less benefit. From numerous experimental results, we suggest an optimal router configuration that combines the best of individual experiments for this architecture, which improves performance by 22%, 19%, and 16% for W-1CPU, W-2CPU, and W-4CPU workloads, respectively.

In future work, we will evaluate other topologies, including two-dimensional mesh and torus. Also, we will study various resource partitioning mechanisms for the on-chip network.

## References

- [1] D. Abts, N. D. E. Jerger, J. Kim, D. Gibson, and M. H. Lipasti. Achieving predictable performance through better memory controller placement in many-core CMPs. In *Proc. of the 31st annual Int’l. Symp. on Computer Architecture (ISCA)*, pages 451–461. ACM, 2009.
- [2] T. W. Ainsworth and T. M. Pinkston. On characterizing performance of the cell broadband engine element interconnect bus. In *Proc. of the 1st ACM/IEEE Int’l Symp. on Network-on-Chip (NOCS)*, pages 18–29, Washington, DC, USA, 2007. IEEE Computer Society.
- [3] AMD. Fusion. <http://sites.amd.com/us/fusion/apu/Pages/fusion.aspx>, 2011.
- [4] A. Bakhoda, G. L. Yuan, W. W. L. Fung, H. Wong, and T. M. Aamodt. Analyzing cuda workloads using a detailed GPU simulator. In *Proc. of the 2009 IEEE Int’l. Symp. on Performance Analysis of Systems and Software (ISPASS)*, pages 163–174. IEEE, 2009.
- [5] A. Bakhoda, J. Kim, and T. M. Aamodt. Throughput-effective on-chip networks for manycore accelerators. In *Proc. of the 43rd Int’l. Symp. on Microarchitecture (MICRO)*, pages 421–432. IEEE, 2010.
- [6] L. Bononi, N. Concer, M. D. Grammatikakis, M. Coppola, and R. Locatelli. Noc topologies exploration based on mapping and simulation



- models. In *Proc. of the 10th Euromicro Conf. on Digital System Design Architectures, Methods, and Tools (DSD)*, pages 543–546. IEEE, 2007.
- [7] S. Bourduas and Z. Zilic. A hybrid ring/mesh interconnect for network-on-chip using hierarchical rings for global routing. In *Proc. of the 1st ACM/IEEE Int'l Symp. on Network-on-Chip (NOCS)*, pages 195–204, Washington, DC, USA, 2007. IEEE Computer Society.
- [8] D. Chang, C. Jenkins, P. Garcia, S. Gilani, P. Aguilera, A. Nagarajan, M. Anderson, M. Kenny, S. Bauer, M. Schulte, and K. Compton. ERCBench: An open-source benchmark suite for embedded and reconfigurable computing. In *Proc. of 20th Intl. Conf. on Field Programmable Logic and Applications (FPL)*, pages 408–413, 2010.
- [9] S. Che, M. Boyer, J. Meng, D. Tarjan, J. W. Sheaffer, S.-H. Lee, and K. Skadron. Rodinia: A benchmark suite for heterogeneous computing. In *Proc. of the 2009 IEEE Int'l Symp. on Workload Characterization (IISWC)*, pages 44–54. IEEE, 2009.
- [10] N. Concer, S. Iamundo, and L. Bononi. aEqualized: A novel routing algorithm for the spidernoc network on chip. In *Proc. of Design, Automation, and Test in Europe (DATE)*, pages 749–754, Leuven, Belgium, 2009. European Design and Automation Association.
- [11] R. Das, O. Mutlu, T. Moscibroda, and C. R. Das. Application-aware prioritization mechanisms for on-chip networks. In *Proc. of the 42nd Int'l. Symp. on Microarchitecture (MICRO)*, pages 280–291, New York, NY, USA, 2009. ACM.
- [12] R. Das, O. Mutlu, T. Moscibroda, and C. R. Das. Aergia: exploiting packet latency slack in on-chip networks. In *Proc. of the 32nd annual Int'l. Symp. on Computer Architecture (ISCA)*, pages 106–116, New York, NY, USA, 2010. ACM.
- [13] V. Dumitriu and G. Khan. Throughput-oriented noc topology generation and analysis for high performance socs. *IEEE Transactions on Very Large Scale Integration (VLSI) Systems*, 17(10):1433–1446, 2009.
- [14] P. Gratz, B. Grot, and S. W. Keckler. Regional congestion awareness for load balance in networks-on-chip. In *Proc. of the 14th Int'l. Symp. on High Performance Computer Architecture (HPCA)*, pages 203–214, Washington, DC, USA, 2008. IEEE Computer Society.
- [15] B. Grot, J. Hestness, S. W. Keckler, and O. Mutlu. Express cube topologies for on-chip interconnects. In *Proc. of the 15th Int'l. Symp. on High Performance Computer Architecture (HPCA)*, pages 163–174, Washington, DC, USA, 2009. IEEE Computer Society.
- [16] B. Grot, J. Hestness, S. W. Keckler, and O. Mutlu. Kilo-noc: a heterogeneous network-on-chip architecture for scalability and service guarantees. In *Proc. of the 33rd annual Int'l. Symp. on Computer Architecture (ISCA)*, pages 401–412, New York, NY, USA, 2011. ACM.
- [17] HPArch. MacSim simulator. <http://code.google.com/p/macsim/>, 2012.
- [18] J. Hu and R. Marculescu. Exploiting the routing flexibility for energy/performance aware mapping of regular noc architectures. In *Proc. of Design, Automation, and Test in Europe (DATE)*, pages 688–693, Washington, DC, USA, 2003. IEEE Computer Society.
- [19] J. Hu and R. Marculescu. DyAD: smart routing for networks-on-chip. In S. Malik, L. Fix, and A. B. Kahng, editors, *Proc. of the 41st annual Design Automation Conference (DAC)*, pages 260–263, New York, NY, USA, 2004. ACM.
- [20] IMPACT. Parboil benchmark suite. <http://impact.crhc.illinois.edu/parboil.php>.
- [21] Intel. Sandy Bridge. <http://software.intel.com/en-us/articles/sandy-bridge/>, 2011.
- [22] A. Jaleel, K. B. Theobald, S. C. Steely, Jr., and J. Emer. High performance cache replacement using re-reference interval prediction (RRIP). In *Proc. of the 32nd annual Int'l. Symp. on Computer Architecture (ISCA)*, pages 60–71, New York, NY, USA, 2010. ACM.
- [23] J. Kim, W. J. Dally, and D. Abts. Flattened butterfly: a cost-efficient topology for high-radix networks. In *Proc. of the 29th annual Int'l. Symp. on Computer Architecture (ISCA)*, pages 126–137, New York, NY, USA, 2007. ACM.
- [24] J. Lee and H. Kim. TAP: A TLP-aware cache management policy for a CPU-GPU heterogeneous architecture. In *Proc. of the 18th Int'l. Symp. on High Performance Computer Architecture (HPCA)*, pages 91–102, Washington, DC, USA, 2012. IEEE.
- [25] S. Ma, N. D. E. Jerger, and Z. Wang. Dbar: an efficient routing algorithm to support multiple concurrent applications in networks-on-chip. In *Proc. of the 33rd annual Int'l. Symp. on Computer Architecture (ISCA)*, pages 413–424, New York, NY, USA, 2011. ACM.
- [26] A. K. Mishra, N. Vijaykrishnan, and C. R. Das. A case for heterogeneous on-chip interconnects for cmps. In *Proc. of the 33rd annual Int'l. Symp. on Computer Architecture (ISCA)*, pages 389–400, New York, NY, USA, 2011. ACM.
- [27] T. Moscibroda and O. Mutlu. A case for bufferless routing in on-chip networks. In *Proc. of the 31st annual Int'l. Symp. on Computer Architecture (ISCA)*, pages 196–207, New York, NY, USA, 2009. ACM.
- [28] O. Mutlu and T. Moscibroda. Parallelism-aware batch scheduling: Enhancing both performance and fairness of shared dram systems. In *Proc. of the 30th annual Int'l. Symp. on Computer Architecture (ISCA)*, pages 63–74, Washington, DC, USA, 2008. IEEE Computer Society.
- [29] *CUDA Programming Guide, V4.0*. NVIDIA.
- [30] NVIDIA. Project denver. <http://blogs.nvidia.com/2011/01/project-denver-processor-to-usher-in-new-era-of-computing/>, .
- [31] NVIDIA. Fermi: Nvidia's next generation cuda compute architecture. <http://www.nvidia.com/fermi/>, .
- [32] H. Patil, R. Cohn, M. Charney, R. Kapoor, A. Sun, and A. Karunanidhi. Pinpointing representative portions of large intel@itanium@programs with dynamic instrumentation. In *Proc. of the 37th Int'l. Symp. on Microarchitecture (MICRO)*, pages 81–92, Washington, DC, USA, 2004. IEEE Computer Society.
- [33] V. Puente, C. Izu, R. Beivide, J. A. Gregorio, F. Vallejo, and J. M. Prellezo. The adaptive bubble router. *J. Parallel Distrib. Comput.*, 61(9):1180–1208, 2001.
- [34] M. K. Qureshi and Y. N. Patt. Utility-based cache partitioning: A low-overhead, high-performance, runtime mechanism to partition shared caches. In *Proc. of the 39th Int'l. Symp. on Microarchitecture (MICRO)*, pages 423–432, Washington, DC, USA, 2006. IEEE Computer Society.
- [35] D. Sanchez, G. Michelogiannakis, and C. Kozyrakis. An analysis of on-chip interconnection networks for large-scale chip multiprocessors. *ACM Transactions on Architecture and Code Optimization (TACO)*, 7(1):4, 2010.
- [36] L. Seiler, D. Carmean, E. Sprangle, T. Forsyth, M. Abrash, P. Dubey, S. Junkins, A. Lake, J. Sugerman, R. Cavin, R. Espasa, E. Grochowski, T. Juan, and P. Hanrahan. Larrabee: a many-core x86 architecture for visual computing. *ACM Transactions on Graphics (TOG)*, 27(3), 2008.
- [37] G. Suh, S. Devadas, and L. Rudolph. A new memory monitoring scheme for memory-aware scheduling and partitioning. In *Proc. of the 8th Int'l. Symp. on High Performance Computer Architecture (HPCA)*, pages 117–128, Washington, DC, USA, feb. 2002. IEEE Computer Society.
- [38] D. H. Woo and H.-H. S. Lee. Compass: a programmable data prefetcher using idle gpu shaders. In *Proc. of the 15th Int'l. conference on Architectural support for programming languages and operating systems (ASPLOS)*, pages 297–310, New York, NY, USA, 2010. ACM.
- [39] Y. Xie and G. H. Loh. PIPP: promotion/insertion pseudo-partitioning of multi-core shared caches. In *Proc. of the 31st annual Int'l. Symp. on Computer Architecture (ISCA)*, pages 174–183, New York, NY, USA, 2009. ACM.