Zipper: Latency-Tolerant Optimizations for High-Performance Buses

Shibo Chen chshibo@umich.edu University of Michigan Ann Arbor, Michigan, USA

Hailun Zhang hzhang664@wisc.edu University of Wisconsin Madison, Wisconsin, USA

Todd Austin austin@umich.edu University of Michigan Ann Arbor, Michigan, USA

ABSTRACT

As heterogeneous designs take over the world of hardware designs, the data bus plays a vital role in interconnecting hosts and accelerators. While past works have emphasized increasing communication bandwidth for data-hungry workloads, this work focuses on optimizing latency for latency-sensitive acceleration applications. We first study the pattern of various accelerator workloads and demonstrate that various optimization opportunities exist to reduce the communication latency overhead. To help developers exploit these opportunities, we introduce Zipper-a protocol optimization layer that reduces communication costs by enabling device and request level parallelism and exploiting data locality for existing bus standards. We applied Zipper to two domains and implemented the end-to-end system on a heterogeneous hardware platform with an integrated FPGA. Our physical experiments show that Zipper provides 8x speedup for one accelerator with 4.3% logic overhead and 1.5x speedup for another with 0.9% logic overhead.

CCS CONCEPTS

• Computer systems organization → Heterogeneous (hybrid) systems; • Hardware \rightarrow Buses and high-speed links; • Networks \rightarrow Network on chip.

KEYWORDS

Accelerator, HW/SW codesign

ACM Reference Format:

Shibo Chen, Hailun Zhang, and Todd Austin. 2025. Zipper: Latency-Tolerant Optimizations for High-Performance Buses. In 30th Asia and South Pacific Design Automation Conference (ASPDAC '25), January 20-23, 2025, Tokyo, Japan. ACM, New York, NY, USA, 8 pages. https://doi.org/10.1145/3658617. 3697546

1 INTRODUCTION

Data buses are vital in connecting heterogeneous components in today's hardware designs. While high-performance data buses have ramped up bandwidth over time, the access latency has not been scaling on par because link traversal scales poorly as the technology node shrinks [44]. A recent study [26] shows that the round-trip latency through the popular PCI Express Gen 3.0 [2] or Intel Ultra

ASPDAC '25, January 20-23, 2025, Tokyo, Japan

© 2025 Copyright held by the owner/author(s). Publication rights licensed to ACM. ACM ISBN 979-8-4007-0635-6/25/01

https://doi.org/10.1145/3658617.3697546

requests. We can exploit this locality to reduce data movement. To capitalize on latency-tolerant optimization opportunities in a real production system, we desire a generic and reusable solution that provides the following benefits:

Path Interconnect [21] languishes at the micro-second (μ s) scale. Many applications cannot tolerate µs-level latencies, leading to most

of these latencies being fully exposed [8]. As such, long-latency bus

transactions hinder the broad deployment of a wider spectrum of

proved physical designs, we observed two significant opportunities

for latency-tolerant optimization. First, there is parallelism at both

the request and the device level. We can enhance host and accel-

erator utilization by enabling out-of-order and parallel execution;

second, many compute kernels exhibit significant temporal locality:

the results of previous requests often become inputs for subsequent

Although communication latency is hard to reduce through im-

applications and accelerators in the production environment.

- Simple parallelism model: Data dependencies can exist between host and accelerator instructions or within accelerator instructions. Harvesting parallelism requires a dynamic scheduling mechanism that works across ISAs and device boundaries.
- Efficient data tracking: Since data is continually moved between host and accelerator, the system needs to precisely and efficiently track the location of the data to ensure functional correctness.
- Reduced design complexity: Due to the sheer size of possible accelerator designs and platforms, customizing APIs and compilers would be a heavy technology burden for ordinary system developers. Major vendors only provide functionlevel APIs [28, 43], which makes adding compiler support even more difficult. A portable and extensible solution is necessary to make latency optimizations accessible, scalable, and agnostic to the underlying bus standards.

To deliver these features, we propose Zipper. Working on top of existing data buses, Zipper is a dynamic protocol-level optimization layer that reduces bus transaction latency between heterogeneous devices connected through a high-performance data bus. It dynamically analyzes data dependencies, tracks data movement across devices, and exploits locality and parallelism as a program proceeds. Zipper uses a software-defined request scheduling approach that requires no modification to application logic, compilers, or data buses. Zipper's runtime library identifies temporal locality and parallel execution opportunities. Then, it schedules optimized accelerator requests to enable resource-constraint-aware parallelism and data reuse. The implementation details of this runtime library are hidden from developers, and the developers can use encapsulated data types as if they are host-native. For hardware, Zipper offers a small

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

request buffer to cache data and enables out-of-order execution of requests with data reuse.

In Section 4.1, two FPGA-based case studies are presented and shown to benefit significantly from Zipper. Our experiments show that Zipper provides uniformly good end-to-end application speedups, with as much as 8x speedup for one case study.

We summarize the contributions of this work as follows:

- We detail the design of Zipper, a general protocol optimization layer to optimize bus transaction latency for heterogeneous system communication over existing high-performance buses. The optimization layer exploits locality and parallelism opportunities that are currently overlooked without making demands on the data bus or programmers.
- We show a real-world FPGA-based implementation of Zipper, built over Intel's QuickPath Interconnect (QPI) for CPUintegrated FPGAs. We dive into the hardware and runtime software components implemented for this platform.
- We present two real-world FPGA-based case studies with Zipper in a production environment and show significant application-level performance improvements (as much as 8x), with modest area overheads (of less than 5% increase in logic). Zipper exploits significant temporal locality and transaction parallelism in both studies, even for designs with only a 4-entry request scheduling window.

Zipper is open-sourced at https://github.com/zipper-bus-optimizations and ready to deploy in real production environments that connect accelerators using AXI, CCI-P, or CXL data buses.

2 DISCOVERING BUS OPTIMIZATION OPPORTUNITIES

This section discusses three key opportunities that we focus on to optimize bus communication latency.

2.1 Host-Accelerator Communication Convention

After the host connects to the accelerator, the developer creates a shared memory space between them to pass inputs and compute results. To kick off the kernel, the host writes inputs into the shared memory and issues instructions with metadata (*i.e.*, input starting address, result write back address, etc.) to the accelerator through, typically, Memory Mapped Input Output (MMIO). After receiving the instruction, the accelerator fetches the input data from the shared memory, writes back the result to the specified write-back address, and notifies the host. This is usually agnostic to physical implementations (*e.g.*, UPI [21], PCIe [2], Infinity [5], etc.) or data transfer protocols (*e.g.*, CCI-P [20], CXL [30], AXI [7], etc.).

2.2 Optimization Opportunities

In this section, we use a reduction algorithm over an abstract hardware-accelerated operator \otimes , as shown in Algorithm 1, to demonstrate existing latency-tolerant opportunities. In this algorithm, we want to calculate the product of the 2n inputs over a hardware accelerated operator \otimes and store the result to the write-back address. The accelerator is attached to the host system, which runs the algorithm by accessing high-performance data buses.

Figure 1 shows the data-dependence graph of Algorithm 1 and the optimizations to eliminate dependencies, exploit locality, and

ingoritania in refutuetion angoritania with operator 6. The
offload instruction sends the indicated operation to the at-
tached accelerator.
Data: An array of 2n elements: arr[2n], A writeback address
addr _{wr}
Result: y = summation of all elements in arr over special
operator ⊗
$i \leftarrow 0$; result $\leftarrow 1$;
while $i < (2n - 1)$ do
1: $a \leftarrow load(Mem[arr+i]);$
2: b ← load(Mem[arr+i+1]);
3: offload($a \otimes b$);
4: c ← fetch(a \otimes b);
5: offload(c⊗result);
6: result← fetch(c⊗result);
$7: i \leftarrow i+2;$
end
8: Mem[addr.ur] \leftarrow result

Algorithm 1. A reduction algorithm with operator \otimes The

enable parallelism. Starting with the unoptimized implementation in Figure 1a, the developer partitions the instructions based on the devices' capabilities: execute instruction 1, 2, 4, and 6-8 on the host, and offload instruction 3, 5 to the accelerator. A stock compiler cannot optimize cross-device dependencies; thus, the unoptimized system has to execute instructions sequentially in program order. As a result, the host always waits for the accelerator to complete computation and loads the result from the shared memory (performed by fetch instructions 4 and 6) before it can move on to the next instruction and/or use the result in subsequent requests.

2.2.1 Exploitable Temporal Locality. We notice that not every dependency is created equal. A cross-device data dependency is much costlier to resolve than a local data dependency due to the bus communication overhead. Based on this observation, we eliminate cross-device dependencies and replace them with local ones whenever possible. That is, instruction 5's two operands are the result from instruction 3 and its result from the last iteration. Therefore, as shown in Figure 1b, instruction 5 does not need to wait for instruction 4 on the host side to complete to get its input from the host. Rather, we can remove this cross-device dependency by directly forwarding the results from the previous requests 3 and 5. By relocating cross-device dependencies, we can avoid much inter-device communication and thus reduce communication overhead.

2.2.2 Device-level Parallelism. Instructions 4 and 6 block instructions that fetch results from the shared memory. After the dependencies have been relocated, instructions 4 and 6 can be moved off the critical path. As shown in Figure 1c, the host can continue execution while the accelerator is working on the received requests. Being non-blocking, the host can run ahead to fetch new data for future accelerator requests. As long as there is no data dependency across devices, the two devices can run in parallel and do not need to synchronize.

2.2.3 *Request-level Parallelism.* After relocating the data dependencies and enabling device-level parallelism, we can completely offload a sequence of requests to the accelerator. We can also extract request-level parallelism locally on the accelerator to maximize the



Figure 1: Latency-tolerant optimizations for the instruction sequence shown in Algorithm 1.

performance gain. In our example, since instruction 3 is independent of previous accelerator instructions, shown in Figure 1d, it can bypass previous requests or interleave with other requests. The only limitations would be the number of requests the accelerator can handle simultaneously and the accelerator's compute throughput.

3 ARCHITECTING ZIPPER OPTIMIZATIONS

To enable optimizations discussed in Section 2, we propose *Zipper*. Zipper is a protocol layer that resides between the physical bus and the application logic, and thus, it does not require any changes to the compiler, compute kernel, or the underlying data bus. Zipper is a drop-in optimization that significantly reduces the exposed latency that is common for high-performance buses, essentially widening the applicability of these emerging bus technologies.

3.1 Overview of Zipper

Zipper uses a set of communication semantics that captures the locality and dependency information to connect the host and accelerator. Zipper adds a request buffer table to the accelerator that tracks the status of operands and caches recent request results. On the host side, a runtime library analyzes data dependency and catches the data reuse opportunities by observing and tracking accelerator requests. The runtime library also manages communication between the host and the accelerator and hides tedious implementation details from software developers. The rest of this section will describe Zipper's communication protocol, hardware structure, and runtime library.

3.2 Host-Accelerator Communication Protocol

In Zipper, the host sends requests to the accelerator through MMIO while communicating input operands and results with the accelerator through shared memory. The number of fields and the bits in each field can vary depending on the use case. As a rule of thumb, each request should include Instruction, Write-back Address, and Operand Information. Each request can have multiple operands. Each operand may reside in the shared memory or the Zipper hardware structure.

The shared memory is the communication channel between the host-accelerator for input operands and results. We partition the



Figure 2: Zipper hardware structure and life cycle of an accelerator request.

shared memory into an operand partition and a result partition. The input operands are continuously placed in the operand partition and wrapped over to reuse the old memory when it reaches capacity. The result partition maintains the bijection with the acceleratorside buffer table entries. For input operands smaller than the size of one memory request granularity, Zipper packs multiple operands for different requests into one cache line to reduce the number of accesses to the memory. Zipper software attaches a version bit to each operand when issuing the request. The Zipper accelerator verifies the freshness of the operand by matching the version bit with the version bit it receives from the runtime library.

3.3 Zipper Hardware Structure

Zipper hardware is on the accelerator side and handles requests it receives from the host. It consists of four parts, shown in Figure 2: a request buffer table, an execution scheduler, a memory controller, and the accelerator. The memory controller is platform-specific, and the compute kernel is user-specific. Zipper does not make intrusive modifications to these two components to work.

Zipper uses the request buffer table and the execution scheduler to enable request-level parallelism. In ①, when Zipper hardware receives a request from its software counterpart, typically through MMIO, it will first store the request information in the request buffer table. The request buffer table stores and tracks all the details on pending and recently completed requests: the instruction, the status of each operand, the write-back address, etc. The request buffer table can be sized as needed. We will discuss the impact of buffer size in Section 5.3. The execution scheduler decides which request is ready for execution and can dispatch instructions out-oforder. The scheduling logic prioritizes older requests when multiple requests are ready to be dispatched.

The request buffer table caches recent results until new requests overwrite the entries. Zipper hardware then reconstructs the data dependency chain based on the information embedded in the requests. For each accelerator instruction, Zipper fetches their values based on the information provided in the request. (2) If the operand is in memory, Zipper issues a read request to the memory controller and marks it as "in fetch" to avoid duplicated access. If the operand comes from a prior request, Zipper either fetches the value if it is ready in the buffer table or waits until the prior request has been completed. ④ Once all operands are resolved, Zipper marks this request as ready to be dispatched. (5) When the computation is done, Zipper stores the results back in the buffer table and writes the results into their corresponding write-back address in the memory, shown in (7). (6) If there are pending requests whose inputs are dependent on the newly completed request, Zipper directly forwards the value when the result is ready.

3.4 Zipper Runtime Library

Zipper provides a non-blocking host-side interface that tolerates multiple pending requests within the accelerator's resource limitations. On the host side, Zipper conducts dependency analysis, request scheduling, and result fetching with a software runtime library. Zipper provides packaged data classes to host applications as if they were host-native types. These data classes encapsulate overridden functions and necessary metadata. This approach provides flexibility and dynamic scheduling capabilities without compiler modification. Figure 3 shows Zipper's software data structures and corresponding updates when behaving different functions.

The Zipper runtime library maintains two data structures to track data objects and communicate with the accelerator: class objects and result lists. The class objects track the status of the requests and the results if the requests are complete.

The result lists track all the software data objects associated with each hardware buffer table entry. When Zipper fetches the results back to the host or clears a table entry, it iterates through the list and updates all relevant data objects. This enables Zipper to track multiple in-flight requests.

We use a code snippet shown in Figure 2 to demonstrate the operations of the Zipper runtime library. The code first calculates an accelerator request and its result *a* with input from the host and then calculates another variable *b*, reusing *a*'s value. After these two accelerator requests, it re-assigns *b* to *a*. Lastly, it retrieves the *a*'s value from the accelerator back to the host.

Algorithm 2 shows the process of Zipper runtime library fetching results and sending new requests to the accelerator.

3.4.1 Issuing New Requests. Figure 3a shows Zipper issuing a new accelerator request to the accelerator. ① Within the context shown in the figure, Zipper registers object a into an available slot in result lists. ② Zipper will store the input operands m and n in the shared memory and send their relative location to the accelerator. In the last step ④, Zipper updates a's validity as false and marks it to

Algorithm 2: Procedures to issue new requests in Zipper						
runtime library. BUF_TBL = buffer table. MEM = memory.						
Data: A list of <i>n</i> input operands <i>ops</i> and Instruction <i>inst</i>						
Result: An object <i>r</i> that tracks the result						
$r.valid \leftarrow false;$						
$r.inAccl \leftarrow true;$						
$r.location \leftarrow nextSlot + +;$						
if resultLists[nextSlot].occupied then						
for $obj \leftarrow resultLists[nextSlot]$ do						
obj.fetchResult;						
end						
end						
Request req;						
$req.inst \leftarrow inst;$						
for $i \leftarrow (0 \rightarrow n-1)$ do						
if ops[i].inAccl then						
$req.ops[i].mode \leftarrow BUF_TBL;$						
$req.ops[i].location \leftarrow ops[i].location;$						
end						
else						
$req.ops[i].mode \leftarrow MEM;$						
$req.ops[i].location \leftarrow nextMemOperandSlot;$						
$nextMemOperandSlot \leftarrow (nextMemOperandSlot +$						
1) %MaxNumOfSlots;						
end						
end						
send req;						
return r;						

be inside the accelerator at location 3. After this step, the program can continue onto the next host instruction or accelerator request.

3.4.2 Enabling Accelerator-Side Caching. In Figure 3b, Zipper makes another accelerator request. Since no empty slot is available, Zipper first clears the oldest entry as shown in Step (①). Zipper forces each object mapped to slot 1 to fetch its value to the host memory if they have not already and update them as not in the accelerator's buffer anymore. During the analysis stage, Zipper detects variable *a* is at location 3 of the accelerator buffer and its value can be reused, so Zipper will not write *a* to the shared memory nor need to fetch *a*'s value back. Instead, Zipper instructs the hardware to get *a*'s value directly from buffer table slot 3. In this way, Zipper detects the relocation opportunities on the host and utilizes the hardware buffer to exploit them. We then append *b* to the result lists and update its metadata similar to what we did to *a* in the last request.

3.4.3 Object Reassignment. When reassigning an object to track another object, as in Figure 3c, Zipper changes the data structure to reflect this reassignment. We reassign b to variable a. ① Zipper copies a's metadata to b and moves b away from its original slot in the result lists to the same slot as a. Similarly, Zipper removes the object from the result list when the object is getting deleted.

3.4.4 Lazy Fetch. Zipper never proactively retreives results until the value is needed. As the code execution progresses, the host eventually asks for the value of *a* to proceed, shown in Figure 3d. In this case, Zipper fetches *a*'s result from its tracking location 3. If the

Zipper: Latency-Tolerant Optimizations for High-Performance Buses







result is not ready, the host will stall due to hard dependency. Once Zipper fetches the value from the shared memory, it will update *a*'s value and its metadata. Zipper will also update all the objects that are tracking location 3. However, *a*'s value is still in the accelerator buffer for future use as no new request evicts *a* yet.

4 REAL-WORLD EXPERIMENTAL SETUP

This section discusses the case studies and hardware setup that are representative of the production environment. We conducted our experiments on Intel HARPv2 [13] with an in-package FPGA. The system contains a 64K FPGA-side coherent cache. The Zipperaugmented software runs on the Intel Xeon E5-2699v4 at 2.2 GHz, and the Zipper-enabled hardware kernels run on the Arria10 FPGA. The host and the FPGA are connected with Intel QuickPath Interconnect using Core Cache Interface.

4.1 Real-World Case Studies

We evaluated Zipper on two applications that rely on CPU and accelerator to compute and are highly sensitive to the communication latency between the two devices:

(1) We replaced the floating point representation in the NASA Parallel Benchmark (*NPB*) [11] with a Posit32 number representation. NPB is a high-performance scientific computing benchmark, including algorithms that require both performance and precisions, such as Fast Fourier Transform (FFT) [18] and MultiGrid (MG) [42]. Posit is a 32-bit number format that achieves better precision than floating points but currently lacks native hardware support. All posit computations are computed with a hardware kernel.

(2) We implemented sequestered encyption (SE) based hardware isolation support for the integer subset of VIP-Bench [10]. VIP-Bench is a set of basic algorithms and applications (*i.e.*, bubble-sort [15], Tiny Encryption Algorithm (TEA) [41],*etc.*) implemented in a data-oblivious manner where only the SE hardware enclave can see the plaintext values of the secrets [9]. We prototyped the SE enclave on an FPGA, and all privacy-enhanced operators are offloaded to the SE enclave.

These two latency-sensitive applications represent interesting privacy and HPC acceleration opportunities that benefit greatly from fine-grained offloading. For optimal performance-area tradeoffs, we used an 8-entry buffer table design for the Posit32 accelerator and a 2-entry design for the SE enclave design. In the baseline design, each request is issued and executed sequentially.

5 EXPERIMENTAL EVALUATION

This section provides an analysis of the performance speedup, area overhead, impact of various optimizations, request buffer table sizes, workload profiles, and other relevant design aspects.

5.1 Performance Speedup and Logic Overhead

Figure 4 shows the relative performance of Zipper over the baseline design. The figure also provides insights into each feature's contribution to the overall performance. On average, Zipper provides 8x speedup for NPB with Posit32 and 1.5x for VIP-Bench with the SE. Dependency relocation provides the most significant



Figure 4: Relative performance of Zipper and various de-featured Zipper over the baseline. RLP = Request-level Parallelism.



Figure 5: Comparison of the number of bus transactions by accelerator between Zipper, de-featured Zipper, and the baseline.

speedup, while request-level parallelism and memory coalescing each provide a smaller but noticeable speedup on top of each other.

We synthesized our design with Intel Quartus Pro 16.0.0.211 onto the targeted FPGA platform. We compute the logic overhead by taking the added Zipper logic over the accelerator and the existing bus control logic. The logic overhead of Zipper is only 4.3% for the 8-entry Zipper Posit32 design over the baseline design and 0.9% for the 2-entry Zipper SE enclave design over the baseline design.

5.2 Accelerator Memory Access

Zipper's performance benefits greatly from reducing accelerator memory demands by exploiting temporal locality and memory coalescing. Figure 5 shows the percentage of the bus transactions Zipper and other de-featured design options make over the baseline design.

For NPB with Posit32, Zipper reduces the accelerator's bus transactions by 77% from the baseline. Request-level parallelism enables out-of-order execution but does not reduce any memory access. Since Zipper can pack 8 input operands into one cache line, memory coalescing reduces 63% of bus transactions over the baseline. Dependency relocation exploits temporal locality and data reuse, reducing 34% of bus transactions over the baseline.

Since operands are larger in the VIP-Bench with SE enclave design, Zipper cannot pack the operands as tight as with Posit32 numbers. Therefore, it is more likely that the operands for the same request span over two cache lines, which leads to more memory access and fewer opportunities for memory coalescing. Zipper reduces 46% of the bus transactions while memory coalescing and dependency relocation reduce 37% and 27% of the bus transactions over the baseline SE enclave design, respectively.

5.3 Impact of Hardware Buffer Size

To study the optimal number of buffers for different workloads, we analyzed the distance of the data dependency chain in Zipper requests or the number of entries we need to provide for efficient dependency relocation.

Our experiment results show that 91% and 92% of the temporal locality can be captured with only four buffer entries for NPB and VIP-Bench, respectively. Table 1 shows the number of requests that can be processed in parallel, the percentage of results required to be fetched back into the host memory, and the average time distance (*in microseconds*) between the host issuing request and the hosting using the request result. As we increase the number of buffer entries, Zipper can exploit more parallelism while facing diminishing returns. As Zipper harvests more operand reuse with larger buffers, the results that need to be fetched decrease as more request dependencies get relocated. The average time distance also increases as fewer results are fetched back to the host, giving the host more time to execute host-side codes in parallel. Note that Zipper: Latency-Tolerant Optimizations for High-Performance Buses

Application	Window Size			
	0	2	4	8
	Exploitable parallelism			
NPB w/ Posit	1	1.89	3.53	6.17
VIP-Bench w/ SE	1	1.87	2.69	3.94
	Percentage of results to be fetched back			
NPB w/ Posit	100%	45%	26%	22%
VIP-Bench w/ SE	100%	55%	21%	10%
	Distance between issue and use			
NPB w/ Posit	251.86	543.59	840.71	821.11
VIP-Bench w/ SE	65	136.68	1656.9	2243.7

Table 1: Zipper characteristics under different instruction window sizes for two applications on average.



Figure 6: Impact of different number of buffer table entries on performance and area for Zipper. ALM = Adaptive Logic Module.

this analysis assumes the system has perfect knowledge of the instruction dependencies during runtime. In practice, Zipper always fetches results back when the buffer entry gets recycled to ensure correctness. A larger buffer gives more time to continue execution until it needs to recycle a buffer entry.

We then analyzed the performance and logic overhead of various sizes. The logic overhead is measured by the usage of adaptive logic modules (ALM), the basic logic units in Intel FPGA families. We construct our experiments around the buffer size of 4. The results are shown in Figure 6. The logic overhead increases exponentially as the size of the buffer increases because of the logic needed for scheduling and storage for operands and results. For NPB with Posit32, the speedup increases logarithmically as we put more entries in the buffer table. However, VIP-Bench with SE Enclave's performance only increases slightly with more buffer entries. The difference is attributed to the latency of each compute kernel. The Posit32 kernel takes two cycles to complete an instruction, while the SE Enclave kernel takes 24 cycles for each instruction. The SE Enclave is more compute-bound to the kernel itself.

6 LIMITATIONS AND FUTURE WORK

While Zipper demonstrates tremendous performance improvement over the baseline, there are additional improvements we can explore as the continuation of this line of work:

- **Request Reordering**: Zipper leverages the optimization opportunities that applications present. However, there would be more temporal locality by reordering the requests and exploiting operator commutativity and associativity. To achieve this, Zipper can issue requests in batches after requests within a scheduling window have been optimized.
- Multi-Agent Cooperation: We considered the scenario with only one accelerator in this work. Multiple accelerators can cooperate to complete the computation in a more complex system. Zipper poises well to enable such extensions as the developer can use optimized scheduling algorithms to dispatch requests to different accelerators.

7 RELATED WORKS

With the emergence of heterogeneous and large-scale systems, communication latency between nodes has become a key focus.

There are four major approaches to tolerate latency: prefetching [1, 4, 6, 22, 23, 29, 31, 38], caching [12, 17, 25, 32, 36, 39], multithreading [3, 14, 34, 40], and relocating [16, 19, 24, 27, 33, 35, 37]. Prefetching predicts the memory access pattern and issues memory accesses before the data is used. This technique does not apply to the challenge tackled in this paper because accelerator requests often rely on host-side data-based control flow, making it hard to issue in advance. Caching keeps data closer to the compute by exploiting spatial and temporal locality. Being tailored specifically to CPU-accelerator interactions, Zipper is more flexible and areaefficient than cache-based designs. Multithreading hides access latency by allocating the hardware resources to another thread while waiting for the long-latency operation to complete. However, its benefits diminish when the operation is at or below the microsecond level due to context switch overhead. Moreover, multithreading relies on having enough threads to schedule and focuses on the throughput. In comparison, Zipper does not rely on switching to other work to occupy the host and significantly speeds up the end-to-end latency. Relocating (i.e., in-memory/near-memory computing) is a design philosophy that moves compute closer to the data. However, even if placed near the memory, the system still needs to tolerate the latency between the host and the accelerator. As a result, this challenge is not directly addressed by relocation.

8 CONCLUSIONS

This paper details Zipper, a bus latency optimization framework for latency-sensitive accelerated applications. By carefully tracking CPU-accelerator dependencies, Zipper can exploit device- and request-level parallelism and temporal locality to significantly reduce exposed bus transaction latency. Zipper is implemented as a protocol optimization layer over an existing bus interface. Minimal system or programmer support is required, as Zipper uses runtime library support for dynamic scheduling and an additional hardware structure for executing the requests from the host. Zipper is deployed on Intel's HARPv2 platform, where two real-world accelerated applications are examined with and without Zipper optimizations. Zipper achieves a 1.5x-8x speedup with low logic overheads for the two case studies presented. This work demonstrates that protocol latency optimizations have significant promise to reduce the exposed latency of high-performance buses and widen their applicability to future application acceleration opportunities.

ASPDAC '25, January 20-23, 2025, Tokyo, Japan

REFERENCES

- Sam Ainsworth and Timothy M Jones. 2017. Software Prefetching for Indirect Memory Accesses. In 2017 IEEE/ACM International Symposium on Code Generation and Optimization (CGO). IEEE, 305–317.
- [2] Jasmin Ajanovic. 2009. PCI Express 3.0 Overview. In Hot Chips Symposium. 1–61.
- [3] Haitham Akkary and Michael A Driscoll. 1998. A Dynamic Multithreading Processor. In Proceedings. 31st Annual ACM/IEEE International Symposium on Microarchitecture. IEEE, 226–236.
- [4] Hasan Al Maruf and Mosharaf Chowdhury. 2020. Effectively Prefetching Remote Memory with Leap. In 2020 USENIX Annual Technical Conference (USENIX ATC 20). 843–857.
- [5] AMD. 2019. AMD Infinity Architecture: The Foundation of the Modern Datacenter. https://www.amd.com/system/files/documents/LE-70001-SB-InfinityArchitecture.pdf.
- [6] DW Anderson, FJ Sparacio, and Robert M Tomasulo. 1967. The IBM System/360 Model 91: Machine Philosophy and Instruction-Handling. *IBM Journal of Research* and Development 11, 1 (1967), 8–24.
- [7] ARM. 2021. AMBA AXI and ACE Protocol Specification. Version H.c. https: //developer.arm.com/documentation/ihi0022/hc/?lang=en.
- [8] Luiz Barroso, Mike Marty, David Patterson, and Parthasarathy Ranganathan. 2017. Attack of the Killer Microseconds. *Commun. ACM* 60, 4 (mar 2017), 48–54. https://doi.org/10.1145/3015146
- [9] Lauren Biernacki, Meron Zerihun Demissie, Kidus Birkayehu Workneh, Fitsum Assamnew Andargie, and Todd Austin. 2022. Sequestered Encryption: A Hardware Technique for Comprehensive Data Privacy. In 2022 IEEE International Symposium on Secure and Private Execution Environment Design (SEED). 73–84. https://doi.org/10.1109/SEED55351.2022.00014
- [10] Lauren Biernacki, Meron Zerihun Demissie, Kidus Birkayehu Workneh, Galane Basha Namomsa, Plato Gebremedhin, Fitsum Assamnew Andargie, Brandon Reagen, and Todd Austin. 2021. VIP-Bench: A Benchmark Suite for Evaluating Privacy-Enhanced Computation Frameworks. In 2021 International Symposium on Secure and Private Execution Environment Design (SEED). IEEE, 139–149.
- [11] Steven W. D. Chien, Ivy B. Peng, and Stefano Markidis. 2020. Posit NPB: Assessing the Precision Improvement in HPC Scientific Applications. In *Parallel Processing and Applied Mathematics*, Roman Wyrzykowski, Ewa Deelman, Jack Dongarra, and Konrad Karczewski (Eds.). Springer International Publishing, Cham, 301–310.
- [12] Jongsok Choi, Kevin Nam, Andrew Canis, Jason Anderson, Stephen Brown, and Tomasz Czajkowski. 2012. Impact of Cache Architecture and Interface on Performance and Area of FPGA-based Processor/Parallel-accelerator Systems. In 2012 IEEE 20th International Symposium on Field-Programmable Custom Computing Machines. IEEE, 17–24.
- [13] Ian Cutress. 2018. Intel Shows Xeon Scalable Gold 6138P with Integrated FPGA, Shipping to Vendors. (2018). https://www.anandtech.com/show/12773/intelshows-xeon-scalable-gold-6138p-with-integrated-fpga-shipping-to-vendors
- [14] Susan J Eggers, Joel S Emer, Henry M Levy, Jack L Lo, Rebecca L Stamm, and Dean M Tullsen. 1997. Simultaneous Multithreading: A Platform for Next-Generation Processors. *IEEE MICRO* 17, 5 (1997), 12–19.
- [15] Edward H. Friend. 1956. Sorting on Electronic Computer Systems. J. ACM 3, 3 (July 1956), 134–168. https://doi.org/10.1145/320831.320833
- [16] Maya Gokhale, Bill Holmes, and Ken Iobst. 1995. Processing in Memory: The Terasys Massively Parallel PIM Array. Computer 28, 4 (1995), 23–31.
- [17] James R Goodman. 1983. Using Cache Memory to Reduce Processor-Memory Traffic. In Proceedings of the 10th annual international symposium on Computer architecture. 124–131.
- [18] M. Heideman, D. Johnson, and C. Burrus. 1984. Gauss and the History of the Fast Fourier Transform. *IEEE ASSP Magazine* 1, 4 (1984), 14–21. https://doi.org/10. 1109/MASSP.1984.1162257
- [19] Daniele Ielmini and H-S Philip Wong. 2018. In-memory Computing with Resistive Switching Devices. Nature electronics 1, 6 (2018), 333–343.
- [20] Intel. 2019. Intel Acceleration Stack for Intel® Xeon® CPU with FPGAs Core Cache Interface (CCI-P) Reference Manual. https://www.intel.com/content/ www/us/en/docs/programmable/683193/current/acceleration-stack-for-cpuwith-fpgas.html.
- [21] Intel. 2019. Intel® Xeon® Processor Scalable Family Technical Overview. https://www.intel.com/content/www/us/en/developer/articles/technical/xeonprocessor-scalable-family-technical-overview.html.
- [22] Saba Jamilan, Tanvir Ahmed Khan, Grant Ayers, Baris Kasikci, and Heiner Litz. 2022. APT-GET: Profile-guided Timely Software Prefetching. In Proceedings of the Seventeenth European Conference on Computer Systems. 747–764.
- [23] Adwait Jog, Onur Kayiran, Asit K Mishra, Mahmut T Kandemir, Onur Mutlu, Ravishankar Iyer, and Chita R Das. 2013. Orchestrated Scheduling and Prefetching for GPGPUs. In Proceedings of the 40th Annual International Symposium on Computer Architecture. 332–343.
- [24] Liu Ke, Udit Gupta, Benjamin Youngjae Cho, David Brooks, Vikas Chandra, Utku Diril, Amin Firoozshahian, Kim Hazelwood, Bill Jia, Hsien-Hsin S. Lee, Meng Li, Bert Maher, Dheevatsa Mudigere, Maxim Naumov, Martin Schatz,

Mikhail Smelyanskiy, Xiaodong Wang, Brandon Reagen, Carole-Jean Wu, Mark Hempstead, and Xuan Zhang. 2020. RecNMP: Accelerating Personalized Recommendation with Near-Memory Processing. In 2020 ACM/IEEE 47th Annual International Symposium on Computer Architecture (ISCA). 790–803. https: //doi.org/10.1109/ISCA45697.2020.00070

- [25] PAH Knoben. 2021. Software Caching for Tree-based Algorithms on Accelerator Cards. Master's thesis. University of Twente.
- [26] Yanqiang Liu, Jiacheng Ma, Zhengjun Zhang, Linsheng Li, Zhengwei Qi, and Haibing Guan. 2021. MEGATRON: Software-Managed Device TLB for Shared-Memory FPGA Virtualization. In 2021 58th ACM/IEEE Design Automation Conference (DAC). 1213–1218. https://doi.org/10.1109/DAC18074.2021.9586197
- [27] Zhiyu Liu, Irina Calciu, Maurice Herlihy, and Onur Mutlu. 2017. Concurrent Data Structures for Near-Memory Computing. In Proceedings of the 29th ACM Symposium on Parallelism in Algorithms and Architectures. 235–245.
- [28] Enno Luebbers, Song Liu, and Michael Chu. [n. d.]. Simplify Software Integration for FPGA Accelerators with OPAE. ([n. d.]). http://eulerproject.com/assets/files/ open-programmable-acceleration-engine-paper.pdf
- [29] Sparsh Mittal. 2016. A Survey of Recent Prefetching Techniques for Processor Caches. ACM Computing Surveys (CSUR) 49, 2 (2016), 1–35.
- [30] Patrick Patrick Kennedy. 2022. Compute Express Link CXL Latency–How Much is Added at HC34. https://www.servethehome.com/compute-express-link-cxllatency-how-much-is-added-at-hc34/#:~:text=The%20CXL%20Consortium% 20is%20using, 170%2D250ns%20for%20CXL%20memory.&text=If%20CXL% 20seems%20to%20be, with%20Q2%202022%20Wind%2DDown.
- [31] R Hugo Patterson, Garth A Gibson, Eka Ginting, Daniel Stodolsky, and Jim Zelenka. 1995. Informed Prefetching and Caching. In Proceedings of the fifteenth ACM symposium on Operating systems principles. 79–95.
- [32] Christian Pinto, Yiannis Gkoufas, Andrea Reale, Seetharami Seelam, and Steven Eliuk. 2018. Hoard: A Distributed Data Caching System to Accelerate Deep Learning Training on the Cloud. arXiv preprint arXiv:1812.00669 (2018).
- [33] Carlos Ríos, Nathan Youngblood, Zengguang Cheng, Manuel Le Gallo, Wolfram HP Pernice, C David Wright, Abu Sebastian, and Harish Bhaskaran. 2019. In-memory Computing on a Photonic Platform. *Science Advances* 5, 2 (2019), eaau5759.
- [34] Amir Roth and Gurindar S Sohi. 2001. Speculative Data-Driven Multithreading. In Proceedings HPCA Seventh International Symposium on High-Performance Computer Architecture. IEEE, 37–48.
- [35] Fabian Schuiki, Michael Schaffner, Frank K Gürkaynak, and Luca Benini. 2018. A Scalable Near-Memory Architecture for Training Deep Neural Networks on Large In-Memory Datasets. *IEEE Trans. Comput.* 68, 4 (2018), 484–497.
- [36] Yakun Sophia Shao, Sam Xi, Viji Srinivasan, Gu-Yeon Wei, and David Brooks. 2015. Toward Cache-Friendly Hardware Accelerators. In HPCA Sensors and Cloud Architectures Workshop (SCAW). 1–6.
- [37] Gagandeep Singh, Lorenzo Chelini, Stefano Corda, Ahsan Javed Awan, Sander Stuijk, Roel Jordans, Henk Corporaal, and Albert-Jan Boonstra. 2018. A Review of Near-Memory Computing Architectures: Opportunities and Challenges. In 2018 21st Euromicro Conference on Digital System Design (DSD). IEEE, 608–617.
- [38] Alan Jay Smith. 1978. Sequential Program Prefetching in Memory Hierarchies. Computer 11, 12 (1978), 7–21.
- [39] Alan Jay Smith. 1982. Cache Memories. ACM Computing Surveys (CSUR) 14, 3 (1982), 473–530.
- [40] Lawrence Spracklen and Santosh G Abraham. 2005. Chip multithreading: Opportunities and challenges. In 11th International Symposium on High-Performance Computer Architecture. IEEE, 248–252.
- [41] David J. Wheeler and Roger M. Needham. 1995. TEA, a Tiny Encryption Algorithm. In *Fast Software Encryption*, Bart Preneel (Ed.). Springer Berlin Heidelberg, Berlin, Heidelberg, 363–366.
- [42] R. Wienands and W. Joppich. 2005. Practical Fourier Analysis for Multigrid Methods. Taylor & Francis. https://books.google.com/books?id=IOSux5GxacsC
- [43] Xilinx. [n. d.]. Xilinx Runtime Library (XRT). ([n. d.]). https://xilinx.github.io/ XRT/master/html/index.html
- [44] Greg Yeric. 2015. Moore's Law at 50: Are We Planning for Retirement. In 2015 IEEE International Electron Devices Meeting (IEDM). 1.1.1–1.1.8. https://doi.org/ 10.1109/IEDM.2015.7409607

Received 12 July 2024; accepted 8 September 2024