Breakfast of Champions: Towards Zero-Copy Serialization with NIC Scatter-Gather

Deepti Raghavan[†], Philip Levis[†], Matei Zaharia[†], Irene Zhang[†]

Abstract

Microsecond I/O will make data serialization a major bottleneck for datacenter applications. Serialization is fundamentally about data movement: serialization libraries coalesce and flatten in-memory data structures into a single transmittable buffer. CPU-based serialization approaches will hit a performance limit due to data movement overheads and be unable to keep up with modern networks.

We observe that widely deployed NICs possess scatter-gather capabilities that can be re-purposed to accelerate serialization's core task of coalescing and flattening in-memory data structures. It is possible to build a completely *zero-copy*, *zero-allocation* serialization library with commodity NICs. Doing so introduces many research challenges, including using the hardware capabilities efficiently for a wide variety of non-uniform data structures, making application memory available for zero-copy I/O, and ensuring memory safety.

CCS Concepts

• Networks → Programming interfaces.

Keywords

data serialization, kernel bypass networking, datacenters

ACM Reference Format:

Deepti Raghavan, Philip Levis, Matei Zaharia, Irene Zhang. 2021. Breakfast of Champions: Towards Zero-Copy Serialization with NIC Scatter-Gather. In *Workshop on Hot Topics in Operating Systems (HotOS '21), May 31–June 2, 2021, Ann Arbor, MI, USA.* ACM, New York, NY, USA, 7 pages. https://doi.org/10.1145/3458336.3465287

1 Introduction

The microsecond era is here [5]. As Figure 1 shows, datacenter applications today can achieve microsecond packet round-trip times, reaching single digit RTTs with

Permission to make digital or hard copies of part or all 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 third-party components of this work must be honored. For all other uses, contact the owner/author(s).

HotOS '21, May 31–June 2, 2021, Ann Arbor, MI, USA

© 2021 Copyright held by the owner/author(s).

ACM ISBN 978-1-4503-8438-4/21/05.

https://doi.org/10.1145/3458336.3465287

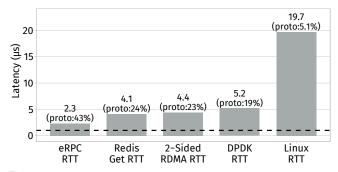


Figure 1: Reported RTTs of recent microsecond-scale systems, annotated with the percentage overhead that Protobuf serialization and deserialization of a single 1024 byte string (1.0 μ s) would add (shown in the dashed line). Redis RTT comes from Arrakis [27], eRPC from eRPC [17], while the RDMA, DPDK and Linux RTTs are measured on the Demikernel. [42].

kernel-bypass. At these latencies, everyday systems services, like data serialization, become unaffordable bottlenecks.

Data serialization [2,3,36–38] is important in datacenter applications. Many distributed applications [1,33,41], RPC libraries [13], and microservice deployments [10] rely on serialization as a communication primitive, but serialization already causes a big performance penalty. Google reported that Protobuf [37] accounted for 5% of its datacenter cycles [18] in 2015, and we expect the problem to worsen today.

Concretely, we find that Protobuf takes 1.0 µs to serialize and deserialize a simple data structure with a single 1024 byte-sized string. Figure 1 overlays this overhead. Protobuf serialization for this data structure adds a staggering 43% overhead to eRPC [17]. Each extra microsecond of serialization overhead significantly affects the throughput a server can achieve and the number of cores necessary to saturate the network.

The main problem is that general-purpose CPUs cannot perform serialization's core task efficiently enough. Serialization must move data, because there is fundamental tension between the application's optimal in-memory layout and the network's optimal on-the-wire layout for a data structure. Data structures often contain pointers (e.g., trees and graphs), so applications can easily modify data structures without having to re-allocate all the memory contiguously. Serialization coalesces these scattered pointers into a contiguous buffer for transmission. Performing this data movement in software will limit throughput in modern

networks, because it requires copying each field at least once and providing a buffer to store the final result.

Without high performance serialization libraries, applications are forced to hand-roll their own serialization or integrate custom hardware accelerators. Redis [31] improves CPU-based serialization by restricting its functionality, but cannot avoid the overhead required to move memory. The most complicated object Redis can serialize is a list. On the other hand, deploying and integrating custom hardware accelerators that do serialization [15, 28] can be difficult in today's datacenters as it requires extra coordination between network administrators, offload developers and application developers [22].

Our key observation is that while CPUs coalesce scattered memory regions inefficiently, widely deployed NICs already perform a similar function: scatter-gather. Scatter-gather was designed for high-performance computing, where applications frequently move large, statically-sized chunks of memory between servers. Kernel bypass exposes this NIC capability to the serialization library, but it is not obvious how to directly use it for serialization. Thus, this paper asks: How can we leverage NIC scatter-gather capabilities to build serialization libraries that keep up with modern networks?

The remainder of the paper describes why existing software serialization is inefficient (§2) and a simple use of NIC scatter-gather for serialization (§3). We finally discuss open research questions around building general-purpose serialization libraries with scatter-gather (§4) and related work (§5).

2 The Limits of Software Serialization

This section shows that CPU-based serialization cannot keep up with the peak packet processing throughput of kernel bypass I/O (§2.1), because CPU-based serialization cannot avoid certain data movement overheads (§2.2).

2.1 Software Serialization Hits a Performance Limit

To demonstrate the overhead of serialization, we benchmark three software serialization libraries [36–38] on DPDK and find that they only achieve up to 52% of DPDK's peak single core throughput. We only consider compilation-based serialization [2, 3, 36, 37] because dynamic type inference at runtime [20, 23] (e.g., Java serialization of arbitrary Java classes) adds unaffordable overheads. We use a data structure with a single 1024-byte string field. Although the data structure is so simple that serialization is theoretically unnecessary, it captures the minimal overhead for serialization today.

The experiment runs on 11 20-core dual socket Xeon Silver 4114 2.2 GHz servers, connected by Mellanox ConnectX-5 100 Gbps NICs and an Arista 7060CX 100 Gbps switch, with a minimum 450 ns of switching latency. We use concurrent, closed-loop clients to send a serialized message to the server, which deserializes, then re-serializes the same payload and

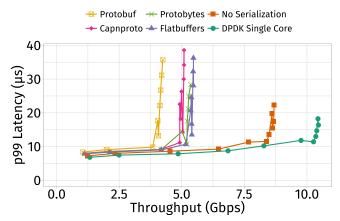


Figure 2: Measured achieved throughput and p99 latency for sets of 1 to 20 concurrent clients (across up to 10 separate machines) pinging a single-core serialization echo server with a message containing a single 1024-byte string. No software serialization library can keep up with the peak zero copy throughput without serialization, which is about 10.4 Gbps.

returns it to the client. We use a minimal UDP networking stack for DPDK based on LWIP [7].

We show the results in Figure 2. The "No Serialization" line removes serialization and gives the raw networking stack performance. Kernel bypass requires that packet memory lives in pinned, non-swappable pages, so the networking stack still copies application payloads into registered packet memory on transmission and copies packets into general memory on receive. The "DPDK Single Core" line removes these copies and represents the peak, zero-copy processing throughput possible with DPDK. We include another version of Protobuf, "Protobytes", where the payload is bytes, not a string, as Protobuf spends a significant amount of time in utf8-validation.

Experiment Results. FlatBuffers, the fastest serialization baseline, achieves only 5.4 Gbps, about 52% of DPDK's peak throughput of 10.4 Gbps (highest throughput measured under 15 μs of tail latency), due to two performance gaps. Serialization itself contributes the first 3 Gbps gap between FlatBuffers and No Serialization. Having the networking stack and serialization manage memory separately contributes the 2 Gbps gap between No Serialization and DPDK Single Core. Section 2.2 closely breaks down these gaps.

2.2 Why is Software Serialization So Expensive?

The overhead of moving data on CPUs limits the performance of today's software serialization libraries. In-memory data structures often contain pointers, so serialization must flatten the data into a contiguous representation. Additionally, sometimes applications use serialization libraries to construct and transmit data structures on-demand to respond to application requests (e.g., returning the value of a range of specified keys in a key-value store).

| Step | Protobuf | Cap'n Proto |
|---------------------------|----------|-------------|
| Initialize Data Structure | 34 ns | 408 ns |
| Copy String Payload | 167 ns* | 80 ns* |
| Encode to Wire Format | 351 ns* | 53 ns |
| Decode from Wire Format | 491 ns* | 78 ns |
| Total Overhead | 1043 ns | 619 ns |

Table 1: Breakdown of steps to serialize and deserialize a message with a single 1024-byte-sized string field. Cap'n Proto's encode and decode are zero-copy because the inmemory buffer layout matches the eventual wire format, while Protobuf requires an expensive transformation to the wire format. Both libraries' copy-based overheads, marked by stars, scale with message size.

All current serialization libraries, no matter their final wire-format, pay the cost of the copies and allocations required for this data movement. Table 1 breaks down the serialization latencies from Figure 2 with Protobuf and Cap'n Proto (FlatBuffers behaves similarly to Cap'n Proto). After copying the field in ("Copy String Payload"), Protobuf performs an expensive transformation to the on-the-wire format. This transformation causes an additional allocation, copy and utf8-validation during "Encode", and corresponding costs during "Decode". Cap'n Proto's "Encode" and "Decode" are cheaper because the in-memory format matches the wire-format exactly, but even Cap'n Proto must allocate space for the serialized buffer ("Initialize Data Structure") and copy the payload in ("Copy String Payload") during transmission. For data structures with large payloads, data movement dominates serialization costs, while converting integers to network ordering, which few wire formats require, adds minimal cost.

The second performance gap in Figure 2 comes from the firm separation between the serialization library and networking stack. Modern kernel bypass stacks require that packets live in non-swappable, pinned memory, so they typically use their own buffers for I/O. Serialization libraries are unaware of the networking stack altogether, so there are inherently copies between the two. Completely eliminating the performance gap in Figure 2 would require tight integration between the serialization library, application and networking stack. This integration would involve agreeing on an interface, making pinned memory available, and coordinating ownership and memory safety of buffers. Fortunately, with kernel bypass, the networking stack, serialization library, and application are all in the same address space, so coordinating memory management may be possible (§4.3).

3 Leveraging the NIC for Serialization

Speeding up serialization requires reducing CPU data movement. Our key insight is that datacenter servers already

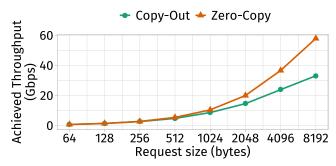


Figure 3: Achieved throughput for 16 clients pinging a single-core echo server with different message sizes (transmitted as a single chunk, without scatter-gather). The server either copies the payload out to a transmit buffer or uses zero-copy transmission. The difference between zero-copy and copy-out becomes visible only at 512 bytes. Note the log scale in the x-axis.

have a hardware accelerator for coalescing non-contiguous I/O regions: the NIC itself. Modern NICs have scatter-gather engines for high-performance computing, e.g., to optimize MPI communication primitives [8, 32]. Networking stacks [9, 34] have re-purposed scatter-gather to manage sending packets that are larger than the maximum packet buffer size. Serialization differs from these use cases because it needs to move potentially many fields whose size and placement dynamically depend on external data or user requests. This section describes the design of a prototype serialization library for the popular Mellanox CX-5 [25] NIC.

3.1 NIC Scatter-Gather Capabilities

Whether NIC scatter-gather can be used for high-performance serialization depends on its performance properties and restrictions. The section focuses on the Mellanox CX-5; other modern scatter-gather NICs with PCIe interconnects likely behave similarly (§4.1).

Given a list of I/O addresses, a CX-5 makes multiple PCIe requests to coalesce the memory into a single packet. The NIC supports up to 60 scattered memory chunks, but each chunk requires a NIC-to-PCIe round trip. The number of these round trips that can execute concurrently depends on hardware implementation details of the PCIe endpoint at the NIC and the CPU, which we currently do not have knowledge of. To understand this penalty, we ran an experiment where the DPDK echo server described in Section 2.1 transmits a pre-initialized payload of size 1024 bytes (no copies) equally divided into different numbers of chunks to a single client. The RTT increases from 6 µs to a 10.5 µs RTT when the message is sent as a single buffer, versus 60 scatter-gather chunks. Sending back the 1024-byte message as 16 chunks results in higher latency than using FlatBuffers to deserialize, reserialize and transmit the request (which requires copying the payload

```
struct ScatterGatherArray {
    size_t num_entries;
    void * ptrs[MAX_ENTRIES];
    size_t length[MAX_ENTRIES];
};
```

Listing 1: The scatter-gather array, the core abstraction for scatter-gather based serialization.

twice). These results suggest that, for a 1024-byte message, the "maximum" number of chunks should be fewer than 16.

There is also a tradeoff between the cost of an additional PCIe request and simply copying the memory. Figure 3 shows an experiment that measures the difference in achieved throughput for 16 clients pinging the single-core DPDK echo server with messages of varying size consisting of a single buffer. The payload is either pre-initialized ("Zero-Copy") or copied into the packet ("Copy-Out"). The only discernible difference between copy-out and zero-copy starts at about 512 bytes. Additionally, entries much smaller than 256 bytes could hurt performance. When the NIC reads memory regions over PCIe, the PCIe controller sends back 256-byte-sized memory chunks (the chunk size is a hardware setting). Each chunk contains a header, so the header could dominate in the case of small payloads.

These results indicate that maximum performance on a CX-5 requires passing in I/O lists with entries that are at least 512 bytes large. The "maximum" number of entries in the I/O list depends on the size of each entry as well as how many concurrent DMAs can run. These tradeoffs preclude simple solutions, such as one scatter-gather operation per data structure field.

3.2 Integrating Networking and Serialization

Core Abstraction: Scatter-Gather Array. Our serialization library's core abstraction is the *scatter-gather array* abstraction, shown in Listing 1. Scatter-gather arrays point to application data in their original memory location. When applications call serialize, the library produces a scatter-gather array that can be passed to the networking stack instead of a single contiguous buffer. Transmitting scatter-gather arrays is conceptually similar to calling the writev system call [12] in Linux with an iovec data structure, except the Linux kernel still copies the iovec into a contiguous buffer before transmission. Section 4.3 discusses research challenges around ensuring application memory can be used for I/O directly.

Serialization API. Our prototype serialization library requires a zero-copy application interface. The generated setter functions store pointers to application memory directly, rather than moving the memory. Listing 2 shows the interface our library would produce for the simple data structure benchmarked in Section 2.1 and how an echo server could use the interface. However, the library only stores

```
message Object { optional string msg = 1; }
class ObjectGenerated {
    std::pair<char *, size_t> get_msg();
    void set_msg(const char *addr, size_t len);
    ScatterGatherArray serialize(size_t num_entries);
    void deserialize(const char *payload);
};
ObjectGenerated obj_recv, obj_send;
obj_recv.deserialize(connection.recv());
recved = obj_recv.get_msg();
obj_send.set_msg(recved.0, recved.1);
ScatterGatherArray sga_send = obj_send.serialize();
connection.send(sga_send);
```

Listing 2: Interface produced by our serialization library in C++, for the listed object schema (in Protobuf syntax), along with example code for an echo server. Unlike prior serialization interfaces, this interface uses zero-copy writes and reads. The serialization library avoids copying fields into a pre-allocated buffer and passes a scatter-gather array to the networking stack for transmission.

pointers for variable-sized values, such as strings, bytes or nested objects. Maintaining pointers to integer fields would not improve performance (storing the pointer to an integer takes about the same space as storing the integer itself), so the serialize function copies integers into the object header.

The header contains a bitmap to index which fields are present, followed by metadata for each field that is present. For the data structure in Listing 2, the corresponding scatter-gather array points to the object header in the first entry and to the string field in the second entry. The object header contains a bitmap that indexes whether the single field is present or not and an offset which points to the string field if it is present. The resulting wireformat is similar to Cap'n Proto's wireformat.

Our library can support nested objects and lists, like Cap'n Proto, FlatBuffers and Protobuf. To support a nested field, the object header contains an offset to the nested object's header (if present). To support a list, the header stores the length of the list and an offset to the actual list data. The final scatter-gather array contains the object header in the first entry (including any nested header data), and pointers to string or bytes fields in further entries from the top-level object as well as any nested objects or lists.

Deserialization API. Deserialization requires turning the received payload back into a pointer-based data structure. This requires linearly scanning through all of the possible fields in the object schema, checking if they are present in the bitmap, and recasting each field offset into a pointer. While linearly scanning through all the fields may add overhead for a data structure with a large number of fields, deserialization

could be "lazily" evaluated if the library changed its wire format slightly. If the object header stored information for all fields, instead of only fields that are present, the compiler would know the location of any field's header information ahead of time. Deserialization could then be a constant-time operation and the library could lazily recover the pointer for any given field when the programmer calls get_field.

Zero-copy deserialization (without copies) causes the application to take ownership of data allocated in the networking stack's packet buffers, which the networking stack might need to reclaim later. Additionally, unless the application uses in-place updates when writing data from received packets (e.g., a put request in Redis), the deserialized data might need to be "re-scattered" into specific in-memory data structures, which requires copies. A fully integrated serialization library and networking stack would need to deal with memory safety and reclamation on the deserialization path (§4.4).

3.3 Prototype Implementation

We implemented this approach for the echo server workload for the data structure in Listing 2 in C++ on top of the same UDP networking stack for DPDK used in Section 2.1. We modified the DPDK datapath to produce a linked list of mbuf packet data structures given the scatter-gather array. The first mbuf contains the packet header with the serialization header copied in. The further mbuf's point to the payloads referenced by the scatter-gather array using DPDK's attach_extbuf API. To comply with kernel bypass I/O memory requirements, the server directly initializes the data structure payload from preregistered memory. However, Section 4.3 discusses strategies to ensure application memory addresses can be used for I/O.

The prototype implementation achieves about 9.15 Gbps (highest throughput measured under 15 µs of tail latency). The prototype's performance improves on all the serialization libraries and the 1-copy ("No Serialization") baseline, but falls about 1.2 Gbps short of the optimal DPDK throughput. We speculate this gap comes from inefficient use of scatter-gather entries (allocating an entire mbuf for just the packet header and object header). Nonetheless, this prototype shows that leveraging NIC scatter-gather is a promising way to accelerate serialization.

4 Open Research Challenges

Many challenges remain in building general-purpose and usable serialization libraries that leverage NIC scatter-gather. This section covers four areas of future work.

4.1 NIC Support for Scatter-Gather

Building a scatter-gather based serialization library requires modeling the performance trade-offs of scatter-gather, which can vary across NICs as well as device drivers. Modeling scatter-gather in current NICs gives insight into how future NIC designs can better support scatter-gather based serializations. Section 3.1 shows that our PCIe-connected NIC adds overhead for transferring small payloads, so scatter-gather can only help for data structures with large enough payloads. Eliminating the PCIe interconnect in the NIC [24] could change these tradeoffs and make scatter-gather beneficial for data structures with smaller payloads. Additionally, understanding how to manage the number of concurrent PCIe requests would help model the time required to transmit any given scatter-gather array.

4.2 Using Scatter-Gather Efficiently

Translating application data structures into scatter-gather arrays that work efficiently with a specific NIC requires optimizing the memory layout of the scatter-gather array. Data structures could vary in size (many fields or few fields), shape (differently-sized fields) and complexity (contain nested objects). Naively creating one scatter-gather entry per data structure field could add overhead, so the serialization library must modify the memory layout of the scatter-gather array before handing it to the NIC. This optimization encompasses coalescing some fields into larger buffers and keeping some fields as separate entries, given a model of scatter-gather performance.

4.3 Accessing Application Memory for Zero-Copy I/O

A completely zero-copy serialization solution requires using arbitrary application memory for I/O, which raises issues related to programming effort and memory fragmentation. Kernel bypass requires that any memory used for I/O lives in pinned and backed pages, because the virtual to physical mappings of this memory must remain the same during the program lifetime. As a result, pinning an entire application's memory for kernel bypass I/O could lead the OS to allocate large amounts of memory that the application will never use. For memory-intensive datacenter workloads, this could impact the performance of other processes or even the ability for other applications to share infrastructure. Thus, the networking stack and serialization library must understand which application memory will be used for I/O and must be pinned.

Pinning memory on demand in the networking stack seems promising but would hinder performance on the packet-processing fast path. On-demand pinning would tell the networking stack which data needs to be pinned, but would add the overhead of a system call to packet transmission. Some NICs have additional penalties to consider. Mellanox NICs require memory registration, so the device can do address translation. However, the NIC can only hold a fixed number of address mappings. Fetching a mapping, done when the first address in a newly mapped region is transmitted, adds a 1 µs latency penalty. If the networking stack registers too many regions, some mappings might fall out of the NIC memory, causing an effect similar to a cache miss.

A new class of *kernel bypass-aware memory allocators* [40, 42] could enable zero-copy dataflows, but raises research challenges related to application integration and memory fragmentation. They could pin large regions of memory beforehand and allocate "dataplane" memory directly into these regions, while allocating "control" memory into a normal heap. To do this transparently, allocators would need to understand which data needs to be registered with minimal programming effort, perhaps with some sort of compiler-based control flow analysis [4]. To enable multitenancy and minimize interference with other processes, the allocators need to to minimize memory fragmentation and understand how to give up unused memory back to the OS.

4.4 Providing Zero-Copy I/O with Memory Safety

A zero-copy serialization stack must provide memory safety, in the form of write and free protection during transmission, and a memory management scheme on the deserialization path. As the Demikernel paper [42] suggests, the memory allocator could provide free protection by adding a reference count to any buffers that are transmitted.

However, providing transparent, efficient write protection from concurrent memory accesses between the NIC and CPU is an open problem. Relying on Linux write protection would add the overhead of a page fault to kernel bypass applications [11]. The networking stack could adopt techniques from recent work [6] to use cache invalidation to detect when addresses are being overwritten and accordingly respond, but this requires custom hardware. Relying on a memory-safe language such as Rust to build the serialization library and networking stack would not protect against read-write races between the NIC hardware and CPU.

On the deserialization path, the networking stack may need to eventually reclaim application buffers (e.g., if an application uses an in-place update to write a value from a received packet). If the application does not free received buffers in time, the networking stack could run out of memory.

5 Related Work

Serialization Acceleration. Many libraries attempt to improve CPU-based serialization by optimizing their wire format [36,38], employing SIMD parallelism for decoding [21], or reducing the overhead of type inference in dynamic serialization [20,23]. These approaches do not remove the fundamental cost required to move memory in software. As a result, recent research proposes offloading serialization to custom accelerators [15, 28, 39] or directly within SSDs for storage [35]. Unlike these accelerators, the scatter-gather functionality already exists in widely used NICs.

Kernel Bypass Systems. Our work is enabled by recent kernel bypass I/O frameworks that expose NIC interfaces directly

to applications in userspace [14,30,34] to eliminate OS level packet processing overheads. Many recent kernel bypass networking stacks [26,27,29,42] build on top of these interfaces to provide APIs to applications while offering low latency, optimized thread scheduling, or zero-copy I/O. eRPC [17] offers general-purpose RPC for commodity networking hardware, and zero-copy networking. None of these systems directly offer general-purpose, zero-copy, data structure serialization as a programming primitive, which requires scatter-gather.

Scatter-Gather Capabilities. High-performance computing applications have used scatter-gather to optimize MPI all-to-all communication primitives [8], or provide zero-copy communication over MPI derived datatypes [32]. Kesavan, et al. [19] uses scatter-gather to measure when zero-copy helps an in-memory database, but does not consider serialization of arbitrary data structures. Derecho [16], a recent SMR system, uses scatter-gather to provide zero-copy I/O for scattered data structures, but relies on specific layouts of data structures provided by their memory allocator. We propose designing general-purpose serialization for application data in arbitrary memory layouts.

6 Conclusion

As link speeds have increased, servers have less cycles to process packets. Object serialization is a core component of datacenter systems, but it cannot keep up with modern networks. We identify that CPU-based software serialization is inherently inefficient, as it relies the CPU to perform data movement. We propose using a hardware capability already present in widely deployed NICs to accelerate serialization: NIC scatter-gather functionality. Our prototype shows that by leveraging NIC scatter-gather to offload data movement from the CPU to the NIC, it is possible to build a zero-copy and zero-allocation serialization library. We identify several areas of future work: better hardware support for scatter-gather, using scatter-gather efficiently, providing transparent memory registration, and ensuring memory safety with zero-copy.

7 Acknowledgements

We thank the anonymous HotOS reviewers, Akshay Narayan, Amy Ousterhout, Anirudh Sivaraman, Anuj Kalia, Jacob Nelson, Kostis Kaffes, Qian Li, Shoumik Palkar, and the members of the Stanford Future Data and SING Research groups for their invaluable feedback. This research was supported in part by affiliate members and other supporters of the Stanford DAWN project—Ant Financial, Facebook, Google, Infosys, NEC, and VMware—as well as Toyota Research Institute, Northrop Grumman, Cisco, SAP, and the NSF under CAREER grant CNS-1651570 and Graduate Research Fellowship grant DGE-1656518. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the authors and do not necessarily reflect the views of the National Science Foundation. Toyota Research Institute ("TRI") provided funds to assist the authors with their research but this article solely reflects the opinions and conclusions of its authors and not TRI or any other Toyota entity.

References

- [1] Apache Software Foundation. Hadoop. https://hadoop.apache.org.
- [2] Apache Software Foundation. Apache avro. https://avro.apache.org/, 2012.
- [3] Apache Sofware Foundation. Apache thrift. https://thrift.apache.org/ download. 2017.
- [4] K. Ashcraft and D. Engler. Using programmer-written compiler extensions to catch security holes. In *IEEE Symposium on Security and Privacy*, 2002.
- [5] L. Barroso, M. Marty, D. Patterson, and P. Ranganathan. Attack of the killer microseconds. *Communicatons of the ACM*, 2017.
- [6] I. Calciu, I. Puddu, A. Kolli, A. Nowatzyk, J. Gandhi, O. Mutlu, and P. Subrahmanyam. Project pherry: Fpga acceleration for remote memory. In *HotOS*, 2019.
- [7] lwIP A Lightweight TCP/IP stack Summary. https://savannah.nongnu. org/projects/lwip/.
- [8] A. Gainaru, R. L. Graham, A. Polyakov, and G. Shainer. Using infiniband hardware gather-scatter capabilities to optimize mpi all-to-all. In EuroMPI 2016. 2016.
- [9] A. Gallatin, J. Chase, and K. Yocum. Trapeze/ip: Tcp/ip at near-gigabit speeds. In ATC, 1999.
- [10] Y. Gan, Y. Zhang, D. Cheng, A. Shetty, P. Rathi, N. Katarki, A. Bruno, J. Hu, B. Ritchken, B. Jackson, et al. An open-source benchmark suite for microservices and their hardware-software implications for cloud & edge systems. In ASPLOS, 2019.
- [11] mprotect(2) linux manual page. https://man7.org/linux/manpages/man2/mprotect.2.html.
- [12] writev(2) linux man page. https://linux.die.net/man/2/writev.
- [13] gRPC Authors. grpc: A high-performance, open source universal rpc framework. https://grpc.io/.
- [14] Storage performance development kit. https://spdk.io/.
- [15] J. Jang, S. J. Jung, S. Jeong, J. Heo, H. Shin, T. J. Ham, and J. W. Lee. A specialized architecture for object serialization with applications to big data analytics. In *ISCA*, 2020.
- [16] S. Jha, J. Behrens, T. Gkountouvas, M. Milano, W. Song, E. Tremel, R. V. Renesse, S. Zink, and K. P. Birman. Derecho: Fast state machine replication for cloud services. ACM Transactions on Computer Systems, 2019.
- [17] A. Kalia, M. Kaminsky, and D. Andersen. Datacenter rpcs can be general and fast. In NSDI, 2019.
- [18] S. Kanev, J. P. Darago, K. Hazelwood, P. Ranganathan, T. Moseley, G.-Y. Wei, and D. Brooks. Profiling a warehouse-scale computer. In ISCA, 2015.
- [19] A. Kesavan, R. Ricci, and R. Stuntsman. To copy or not to copy: Making in-memory databases fast on modern nics. https://rstutsman.github.io/papers/copy-not-to-copy.pdf.
- [20] Kyro. https://github.com/EsotericSoftware/kryo, Accessed January 23, 2021.
- [21] G. Langdale and D. Lemire. Parsing gigabytes of json per second. The VLDB Journal, 2019.
- [22] A. Narayan, A. Panda, M. Alizadeh, H. Balakrishnan, A. Krishnamurthy, and S. Shenker. Bertha: Tunneling through the network api. In *HotNets*, 2020
- [23] K. Nguyen, L. Fang, C. Navasca, G. Xu, B. Demsky, and S. Lu. Skyway: Connecting managed heaps in distributed big data systems. In ASPLOS, 2018.
- [24] S. Novakovic, A. Daglis, E. Bugnion, B. Falsafi, and B. Grot. Scale-out numa. In ASPLOS, 2014.
- [25] Nvidia. Connectx-5. advanced offload capabilities for the most demanding applications. https://www.nvidia.com/enus/networking/ethernet/connectx-5/.
- [26] A. Ousterhout, J. Fried, J. Behrens, A. Belay, and H. Balakrishnan. Shenango: Achieving high CPU efficiency for latency-sensitive datacenter workloads. In NSDI, 2019.

- [27] S. Peter, J. Li, I. Zhang, D. R. K. Ports, D. Woos, A. Krishnamurthy, T. Anderson, and T. Roscoe. Arrakis: The operating system is the control plane. In OSDI, 2014.
- [28] A. Pourhabibi, S. Gupta, H. Kassir, M. Sutherland, Z. Tian, M. P. Drumond, B. Falsafi, and C. Koch. Optimus prime: Accelerating data transformation in servers. In ASPLOS, 2020.
- [29] G. Prekas, M. Kogias, and E. Bugnion. Zygos: Achieving low tail latency for microsecond-scale networked tasks. In SOSP, 2017.
- [30] A rdma protocol specification. http://rdmaconsortium.org/, 2009.
- [31] redis labs. Redis. https://redis.io/.
- [32] G. Santhanaraman, J. Wu, W. Huang, and D. K. Panda. Designing zerocopy message passing interface derived datatype communication over infiniband: Alternative approaches and performance evaluation. The International Journal of High Performance Computing Applications, 2005.
- [33] R. Taft, I. Sharif, A. Matei, N. VanBenschoten, J. Lewis, T. Grieger, K. Niemi, A. Woods, A. Birzin, R. Poss, P. Bardea, A. Ranade, B. Darnell, B. Gruneir, J. Jaffray, L. Zhang, and P. Mattis. Cockroachdb: The resilient geo-distributed sql database. In SIGMOD, 2020.
- [34] Dpdk: Data plane development kit. https://www.dpdk.org/.
- [35] H.-W. Tseng, Q. Zhao, Y. Zhou, M. Gahagan, and S. Swanson. Morpheus: Creating application objects efficiently for heterogeneous computing. In ISCA, 2016.
- [36] W. Van Oortmerssen. Flatbuffers: a memory efficient serialization library. https://opensource.googleblog.com/2014/06/flatbuffers-memoryefficient.html, 2014.
- [37] K. Varda. Protocol buffers: Google's data interchange form. https://opensource.googleblog.com/2008/07/protocol-buffers-googles-data.html, 2008.
- [38] K. Varda. Cap'n proto. https://capnproto.org/, 2020 (Accessed October 22, 2020).
- [39] A. Wolnikowski, S. Ibanez, J. Stone, C. Kim, R. Manohar, and R. Soulé. Zerializer: Towards zero-copy serialization. In *HotOS*, 2021.
- [40] B. Yi, J. Xia, L. Chen, and K. Chen. Towards zero copy dataflows using rdma. In SIGCOMM Posters and Demos, 2017.
- [41] M. Zaharia, M. Chowdhury, T. Das, A. Dave, J. Ma, M. McCauly, M. J. Franklin, S. Shenker, and I. Stoica. Resilient distributed datasets: A faulttolerant abstraction for in-memory cluster computing. In NSDI, 2012.
- [42] I. Zhang, J. Liu, A. Austin, M. L. Roberts, and A. Badam. I'm not dead yet! the role of the operating system in a kernel-bypass era. In *HotOS*, 2019.