

Yang Zhou

Department of Computer Science
Harvard University
150 Western Ave, SEC 4.429
Allston, MA 02134, USA
+1 617 599 8532
yangzhou@g.harvard.edu

November 28, 2023

Dear Faculty Search Committee:

I am writing to apply for a tenure-track position as an assistant professor in your department. I am currently a PhD candidate in the Department of Computer Science at Harvard University and expect to complete my dissertation work by June 2024.

My research interests lie in the areas of computer systems, especially networking, operating systems, and distributed systems. I am particularly interested in full-stack optimizations for efficient and evolvable datacenter infrastructure, by codesigning networking stacks and data-center applications. My research has studied kernel eBPF and kernel-bypass techniques, and various applications like fault-tolerant far memory, Paxos consensus, distributed transactions, and microsecond-scale RPCs.

I have enclosed my curriculum vitae with a list of references, research statement, teaching statement, diversity statement, and three representative publications. These materials are also available online at <https://yangzhou1997.github.io/application/>.

I look forward to hearing from you.

Sincerely,

Yang Zhou

[This page intentionally left blank.]

Yang Zhou

yangzhou1997.github.io
yangzhou@g.harvard.edu ♦ +1 617 599 8532
150 Western Ave, SEC 4.429, Allston, MA 02134, USA

RESEARCH INTERESTS

Networked systems, operating systems, distributed systems, networking stacks, and network telemetry.

EDUCATION

Harvard University, Cambridge, MA, USA

Ph.D. in Computer Science

(Expected) June 2024

M.S. in Computer Science

November 2021

Thesis title: Codesigning Networking Stacks and Datacenter Applications for High Efficiency and Evolvability

Advisors: Minlan Yu and James Mickens

Peking University, Beijing, China

B.S. in Computer Science

July 2018

Thesis title: Towards Faster and More Accurate Data Stream Processing

Advisors: Tong Yang

WORK EXPERIENCE

Harvard University, Research Assistant

August 2018–Present

- *Kernel offloads*: Designed eBPF-based kernel offloads for distributed system protocols including Paxos (Electrode [1]) and serializable transactions (DINT [12]) to reduce kernel networking stack overhead. Implemented and evaluated atop unmodified Linux OSes, and achieved kernel-bypass-like throughput and latency.
- *μ s-scale RPCs*: Designed an efficient inter-server load balancing scheme for μ s-scale RPCs to achieve low tail latency and high goodput (Mew [11]). Implemented and evaluated for both kernel-bypass and kernel-based networking stacks.
- *SmartNIC architecture*: Designed and prototyped SGX-like trusted execution environments for network functions in SmartNICs under multi-tenant cloud environments (S-NIC [13]).

Google NetInfra Group and System Research Group, Student Researcher

June 2021–May 2023

- *Far memory*: Designed an efficient far memory system that leverages erasure-coding, remote memory compaction, one-sided RMAs, and offloadable parity calculations to achieve fast, storage-efficient fault tolerance (Carbink [2]). Implemented and evaluated using production networking stack.
- *Distributed runtime*: Designed an efficient fault-tolerant distributed runtime based on tasks and actors by leveraging the Chandy–Lamport consistent checkpointing algorithm and causal logging mechanism.
- *μ s-scale RPCs*: Identified and motivated the inter-server scheduling problem for μ s-scale RPCs (leading to Mew).

VMware Research, Research Intern

July 2020–September 2020

- *Geo-distributed data analytics*: Applied traffic redundancy elimination (TRE) technique to accelerate geo-distributed data analytics and save WAN traffic cost. Implemented atop Alluxio, an in-memory data cache system for analytics.

Facebook, Research Collaborator

November 2019–May 2020

- *Network telemetry*: Conducted extensive measurement and analysis on Facebook’s network telemetry system. Identified the importance of being evolvable and handling changes. Proposed a change cube abstraction to systematically track changes, and an intent-based layering design to confine and track changes (PCAT [3]).

SenseTime, Software Engineering Intern

March 2018–May 2018

- *Distributed storage*: Worked on Ceph storage setup, testing, maintenance, monitoring, and alerting.

Peking University, Research Assistant

April 2016–July 2018

- *Network telemetry*: Designed and implemented novel probabilistic data structures (e.g., sketches and Bloom filters) to optimize the memory usage, speed, and accuracy of network telemetry tasks (Cold Filter [4], Elastic Sketch [5], Pyramid Sketch [8], and more [6][15][19]).

PUBLICATIONS

Total 780 citations till November 2024 based on Google Scholar.

Conference Publications

- [1] **Yang Zhou**, Zezhou Wang, Sowmya Dharanipragada, and Minlan Yu.
Electrode: Accelerating Distributed Protocols with eBPF. [\[link\]](#)
USENIX NSDI 2023.
- [2] **Yang Zhou**, Hassan Wassel, Sihang Liu, Jiaqi Gao, James Mickens, Minlan Yu, Chris Kennelly, Paul Turner, David Culler, Hank Levy, and Amin Vahdat.
Carbink: Fault-Tolerant Far Memory. [\[link\]](#)
USENIX OSDI 2022.
- [3] **Yang Zhou**, Ying Zhang, Minlan Yu, Guangyu Wang, Dexter Cao, Eric Sung, and Starsky Wong.
Evolvable Network Telemetry at Facebook. [\[link\]](#)
USENIX NSDI 2022.
- [4] **Yang Zhou**, Tong Yang, Jie Jiang, Bin Cui, Minlan Yu, Xiaoming Li, and Steve Uhlig.
Cold Filter: A Meta-Framework for Faster and More Accurate Stream. Processing [\[link\]](#)
ACM SIGMOD 2018.
- [5] Tong Yang, Jie Jiang, Peng Liu, Qun Huang, Junzhi Gong, **Yang Zhou**, Rui Miao, Xiaoming Li, and Steve Uhlig.
Elastic Sketch: Adaptive and Fast Network-Wide Measurements. [\[link\]](#)
ACM SIGCOMM 2018.
- [6] Omid Alipourfard, Masoud Moshref, **Yang Zhou**, Tong Yang, and Minlan Yu.
A Comparison of Performance and Accuracy of Measurement Algorithms in Software. [\[link\]](#)
ACM Symposium on SDN Research (SOSR) 2018.
- [7] Xiangyang Gou, Chenxingyu Zhao, Tong Yang, Lei Zou, **Yang Zhou**, Yibo Yan, Xiaoming Li, and Bin Cui.
Single Hash: Use One Hash Function to Build Faster Hash Based Data Structures. [\[link\]](#)
IEEE International Conference on Big Data and Smart Computing (BigComp) 2018.
- [8] Tong Yang, **Yang Zhou**, Hao Jin, Shigang Chen, and Xiaoming Li.
Pyramid Sketch: A Sketch Framework for Frequency Estimation of Data Streams. [\[link\]](#)
VLDB 2017.
- [9] **Yang Zhou**, Peng Liu, Hao Jin, Tong Yang, Shoujiang Dang, and Xiaoming Li.
One Memory Access Sketch: A More Accurate and Faster Sketch for Per-Flow Measurement. [\[link\]](#)
IEEE Global Communications Conference (Globecom) 2017.
- [10] Junzhi Gong, Tong Yang, **Yang Zhou**, Dongsheng Yang, Shigang Chen, Bin Cui, and Xiaoming Li.
ABC: A Practicable Sketch Framework for Non-Uniform Multisets. [\[link\]](#)
IEEE International Conference on Big Data (BigData) 2017.

Papers Under Reviews

- [11] **Yang Zhou**, Hassan Wassel, James Mickens, Minlan Yu, and Amin Vahdat.
Mew: Efficient Inter-Server Load Balancing for Microsecond-Scale RPCs. [\[link\]](#)
September 2023.
- [12] **Yang Zhou**, Xingyu Xiang, Matthew Kiley, Sowmya Dharanipragada, and Minlan Yu.
DINT: Fast In-Kernel Distributed Transactions with eBPF. [\[link\]](#)
September 2023.
- [13] **Yang Zhou**, Mark Wilkening, James Mickens, and Minlan Yu.
SmartNIC Security Isolation in the Cloud with S-NIC. [\[link\]](#)
October 2023.

Workshop and Demo Publications

- [14] **Yang Zhou**, Hao Jin, Peng Liu, Haowei Zhang, Tong Yang, and Xiaoming Li.
Accurate Per-Flow Measurement with Bloom Sketch. [\[link\]](#)

Journal Publications

- [15] Zhuochen Fan, Gang Wen, Zhipeng Huang, **Yang Zhou**, Qiaobin Fu, Tong Yang, Alex X Liu, and Bin Cui.
On the Evolutionary of Bloom Filter False Positives - An Information Theoretical Approach to Optimizing Bloom Filter Parameters. [\[link\]](#)
IEEE Transactions on Knowledge & Data Engineering 2022.
- [16] Yuanpeng Li, Xiang Yu, Yilong Yang, **Yang Zhou**, Tong Yang, Zhuo Ma, and Shigang Chen.
Pyramid Family: Generic Frameworks for Accurate and Fast Flow Size Measurement. [\[link\]](#)
IEEE/ACM Transactions on Networking 2021.
- [17] Tong Yang, Jie Jiang, **Yang Zhou**, Long He, Jinyang Li, Bin Cui, Steve Uhlig, and Xiaoming Li.
Fast and Accurate Stream Processing by Filtering the Cold. [\[link\]](#)
The VLDB Journal 2019.
- [18] Tong Yang, Jie Jiang, Peng Liu, Qun Huang, Junzhi Gong, **Yang Zhou**, Rui Miao, Xiaoming Li, and Steve Uhlig.
Adaptive Measurements Using One Elastic Sketch. [\[link\]](#)
IEEE/ACM Transactions on Networking 2019.
- [19] **Yang Zhou**, Omid Alipourfard, Minlan Yu, and Tong Yang.
Accelerating Network Measurement in Software. [\[link\]](#)
ACM SIGCOMM Computer Communication Review 2018.

TALKS

-
- Electrode: Accelerating Distributed Protocols with eBPF
Duke University, ACE Center for Evolvable Computing, Google, USENIX NSDI
Columbia University
April 2023
March 2023
 - Carbink: Fault-Tolerant Far Memory
Cornell University
WORDS workshop
Microsoft Research Redmond, USENIX OSDI
Google
November 2023
November 2022
July 2022
March & June 2022
 - Evolvable Network Telemetry at Facebook
USENIX NSDI
Boston University, Meta
April 2022
March 2022
 - Cold Filter: A Meta-Framework for Faster and More Accurate Stream Processing
Harvard University
October 2018

MENTORING EXPERIENCE

-
- Matt Kiley, Harvard College undergraduate
Accelerating distributed transactions using eBPF and AF_XDP-based RPC systems. 2023
 - Yunxi Shen, Tsinghua University undergraduate
Resource-efficient job scheduling in data centers. 2023
 - Xingyu Xiang, Peking University undergraduate
Accelerating distributed transactions using eBPF. 2023
 - Zezhou Wang, Peking University undergraduate → University of Washington PhD
Accelerating Paxos using eBPF (NSDI 2023, [\[1\]](#)). 2022

TEACHING EXPERIENCE

-
- **Guest Lecture** on far memory, CS294-252: Architectures and Systems for Warehouse-Scale Computers, UC Berkeley
Nov 2023

- **Teaching Assistant** for Prof. Minlan Yu, CS145: Networking at Scale, Harvard University *Spring 2021*
- **Teaching Assistant** for Prof. Tong Yang, Algorithm Design and Analysis, Peking University *Fall 2018*

PATENTS

- **Yang Zhou**, Hassan Wassel, Minlan Yu, Hank Levy, David Culler, and Amin Vahdat. “Fault Tolerant Disaggregated Memory”. Pending (US20230185666A1), filed by Google in December 2022.

ACADEMIC HONORS

- Google Ph.D. Fellowship in Systems and Networking *2022*
- Finalist, Meta Ph.D. Fellowship in Networking *2022*
- Graduate Fellowship, Harvard University *2018*
- Excellent Bachelor Thesis (10/327), School of EECS, Peking University *2018*
- New Academic Star Award (1/193), School of EECS, Peking University *2018*
- Arawana Scholarship (2/193), Peking University *2017*
- Pinyou Hudong Scholarship, School of EECS, Peking University *2016*
- May Fourth Scholarship, Peking University *2015*

PROFESSIONAL ACTIVITIES

- PC Member: ACM SIGCOMM Poster/Demo 2023, IEEE INFOCOM Workshop on Networking Algorithms 2020.
- Reviewer (Conferences): ACM SIGKDD 2023.
- Reviewer (Journals): ACM Transactions on Modeling and Performance Evaluation of Computing Systems, IEEE/ACM Transactions on Networking, IEEE Journal on Selected Areas in Communications.
- Panelist: “Getting started with systems research” at Students@Systems 2022.

REFERENCES

Prof. Minlan Yu
Department of Computer Science
Harvard University
150 Western Ave, SEC 4.415
Allston, MA 02134, USA
+1 617 495 3986
minlanyu@g.harvard.edu

Prof. James Mickens
Department of Computer Science
Harvard University
150 Western Ave, SEC 4.416
Allston, MA 02134, USA
+1 617 384 8132
mickens@seas.harvard.edu

Dr. Amin Vahdat
Google Fellow and Vice President of Engineering
Google LLC
1600 Amphitheatre Parkway
Mountain View, CA 94042, USA
+1 650 390 7073
vahdat@google.com

Prof. Adam Belay
MIT CSAIL
32 Vassar St, 32-G996
Cambridge, MA 02139, USA
+1 617 253 0004
abelay@mit.edu

Dr. Ying Zhang
Senior Engineering Manager
Meta Platforms, Inc.
1 Hacker Way
Menlo Park, CA 94025, USA
+1 408 250 9961
zhangying@meta.com

Research Statement

Yang Zhou

I am a systems researcher spanning the areas of networking, operating systems, and distributed systems, focusing on datacenter environments. A datacenter centralizes hundreds of thousands of machines with high-speed networks, enables computations over huge amounts of data, and hosts popular services (e.g., Google search, Netflix streaming, ChatGPT) that impact billions of people’s lives.

To handle massive-scale data and computations, datacenter applications run across multiple machines over networks. Ideally, the underlying datacenter infrastructure should be efficient to maintain steady cloud revenues while meeting high user expectations, and be evolvable to handle the increasingly diverse and performance-hungry applications as well as heterogeneous hardware. However, application-level goals are having a growing **mismatch** with the goals of host networking stacks (involving NICs, kernels, transport layers, and threading) that play a core role in connecting machines in datacenters, causing severe efficiency and evolvability problems. For example, the most widely used kernel networking stack prioritizes security and isolation with separated kernel and user contexts, incurring prohibitive CPU overheads; meanwhile, emerging in-memory applications demand ultra-low latency and high throughput, preferring coalescing different contexts but losing isolation. Even though the networking stacks keep evolving, e.g., the modern kernel-bypass RDMA stacks, applications tend to be network-unaware and abuse (or even deplete) network resources. Such mismatch gets largely exacerbated in large-scale datacenters where networking stacks and applications are usually developed and maintained by disjoint groups of engineers, i.e., network vs. application engineers (due to their growing complexities and industrial organizational structures). This fundamental mismatch causes less efficient use of datacenter resources and hinders the scaling-out of diverse datacenter applications.

My research has focused on bridging the mismatch by **codesigning** low-level host networking stacks and high-level datacenter applications from a systems perspective. My codesign aims to realize high efficiency and agile evolvability for datacenter infrastructure, and it innovates in two directions: (1) customizing networking stacks based on application needs, and (2) redesigning applications to be network-aware and network-efficient. They have borne fruit for many important datacenter applications, including existing ones (e.g., consensus, distributed transactions) and emerging ones (e.g., far memory over networks, microsecond-scale RPCs). My Electrode [1], Dint [2], and Mew [3] safely inject Paxos, transactions, and RPC load balancing logics into the kernel networking stack respectively via eBPF. This not only achieves remarkable performance improvements (by avoiding kernel overheads) but also allows customizing and evolving the kernel stack based on application needs. My Carbink [4] enables network-aware fault tolerance for far memory with high network and memory efficiency, making it practically usable in datacenters with failures being the norm. Specific to evolvability, my PCAT [5] helps Facebook design an evolvable telemetry system to handle frequent changes in production networks.

My research methodology has been empiricism-guided *measuring, tailoring, and fitting* to analyze, optimize, and implement real-world systems—just like how tailors made clothes in the old times. First, I thoroughly measure to reason through the performance characteristics of various networking stack primitives and complex applications; I also draw on my two-year experiences in Google’s networking and system teams to uncover critical feature requirements in production systems. Second, I aggressively tailor unnecessary or overlapping operations in networking stacks and applications to optimize for high efficiency. Third, I strategically partition and fit applications to the right networking stack primitives to efficiently implement the entire system. This focus on **full-stack optimizations** defines my niche as a systems researcher.

Previous Work

CPU efficient distributed protocols with evolvable kernel networking via eBPF. In-memory distributed protocols such as consensus and distributed transactions are important building blocks for datacenter applications. They require intensive network IOs, while the widely-used kernel networking stack gives low IO performance due to high per-IO CPU overhead. Such mismatch has fostered a popular belief that kernel-bypass is the necessary key to high performance for these protocols. However, kernel-bypass is not a panacea: it essentially trades security, isolation, protection, maintainability, and debuggability for performance; it also burns one or more CPU cores for busy-polling even at low loads, which is usually hard to adopt in public cloud deployments due to per-core pricing [6]. As such,

I revisit the above popular belief and ask: is the current kernel networking stack really ill-suited for CPU-efficient distributed protocols, especially given many kernel advancements over decades?

I first measure the source of the high overhead for kernel networking stacks. When running a prior well-designed transaction protocol under a recent Linux kernel networking stack, I find that networking stack traversing dominates the overhead (64% vs. 16% on context switching and 12% on interrupt handling). This motivates me to aggressively tailor unnecessary components of the stack for specific distributed protocols, trading slight genericity loss for performance boosts. For example, the reliable transport along with complex queue disciplines, which incurs costly `sk_buff` maintenance and packet copies, could be cut; this is because (1) distributed protocols themselves can recover from packet loss with application-level timeouts, and (2) packet loss happens rarely within today’s well-engineered datacenter networks. To realize such tailoring, I leverage eBPF to *safely* offload protocol-specific request processing logic into the early stages of the kernel stack; this avoids going through the full stack and user space, removing most of the kernel overheads.

However, offloading complex distributed protocols into the kernel is challenging, because eBPF has a constrained programming model for kernel safety and liveness. To address this challenge, I strategically partition the distributed protocols to fit frequent critical paths into the kernel for high performance while complex rare paths into the user space for full functionalities. Take the classic Multi-Paxos protocol as an example. Electrode [1] offloads failure-free Multi-Paxos operations of broadcasting, acknowledging, and waiting-on-quorums into the kernel via eBPF; when failure happens, it runs complex failure-handling operations in the user space. I implement such partitioning for Multi-Paxos and two transaction protocols (version-based and lock-based) atop unmodified Linux kernels, and achieve remarkable performance boosts. For instance, Dint [2] for transaction offloading achieves up to $23\times$ higher throughput than kernel networking stacks, and $2.6\times$ higher than a recent DPDK-based kernel-bypass stack [7] (as the eBPF offloads directly work on raw ethernet packets, bypassing any socket connections). Owing to the kernel-friendliness and high performance, my eBPF offloading work has sparked interest in both industry (e.g., Meta, Intel) and academia (e.g., University of Washington, University of Michigan, NYU).

Looking further out, future kernel networking stacks should be evolvable in order to efficiently tackle increasingly diverse applications and heterogeneous hardware. My Electrode and Dint projects already demonstrate that eBPF can provide significant evolvability to kernel networking stacks for specific applications. I am now working on an evolvable generic RPC framework by implementing a reliable RPC transport in eBPF; it leverages efficient `AF_XDP` sockets to direct RPC requests to user-space applications for processing. The evolvability of this RPC framework manifests into three aspects: (1) customizing network transport protocols based on application types (e.g., video), (2) customizing the locations of transport layer offloads ranging from SmartNICs to hosts (as many SmartNICs directly support eBPF), and (3) application-informed request load balancing among CPU cores.

Network and memory efficient fault-tolerant far memory. In a datacenter, matching a particular application to just enough memory and CPUs is hard. A commodity server tightly couples memory and compute, hosting a fixed number of CPUs and RAM modules that are unlikely to exactly match the computational requirements of any particular application. Even if a datacenter contains a heterogeneous mix of server configurations, the load on each server (and thus the amount of available resources for a new application) changes dynamically as old applications exit and new applications arrive. Thus, even state-of-the-art cluster schedulers struggle to efficiently bin-pack a datacenter’s aggregate collection of CPUs and RAM. For example, Google [8] and Alibaba [9] report that the average server has only 60% memory utilization, with substantial variance across machines.

Disaggregated datacenter memory is a promising solution. It pairs a CPU with an arbitrary set of possibly-remote RAM modules, with a fast network interconnect keeping access latencies to far memory small. Much of the prior work in this space [10, 11] has a common limitation: a lack of fault tolerance. Unfortunately, in a datacenter containing hundreds of thousands of machines, faults are pervasive. Without fault tolerance, the failure rate of an application using far memory will be much higher than the failure rate of an application that only uses local memory; the reason is that the use of far memory increases the set of machines whose failure can impact an application.

Achieving both network and memory efficient fault-tolerant far memory is challenging. Conventional memory-efficient fault tolerance scheme applies erasure coding, and stripes a single memory page across multiple remote nodes with RMA-based swapping. For brevity, I use *span* to denote “memory page”. Assuming Reed-Solomon code with 4 data chunks and 2 parity chunks, the conventional scheme requires 6 RMAs per span swap-out and 4 RMAs per swap-in, incurring excessive network IO pressure on the networking stack. In Carbink, I tailor the excessive network IOs by eschewing the span-granularity erasure coding, and instead erasure code at the spanset granularity. A spanset consists of multiple spans with the same size, i.e., 4 data spans and 2 parity spans in our example, and gets swapped out together in a batch. This only requires averagely $(4 + 2)/4 = 1.5$ RMAs per span swap-out and a single RMA per swap-in, significantly improving network efficiency.

However, spanset-granularity erasure coding inevitably incurs memory fragmentation. This is because each span lives in exactly one place (either local memory or far memory), and swapping a span inside a spanset from far memory to local memory creates dead space (and thus fragmentation) in far memory. To address this problem, I design a pauseless defragmentation mechanism running off the swapping critical path, asynchronously reclaiming dead space for later swap-outs in the background. In contrast to the simple span swapping via RMA, this background defragmentation has complex two-phase commit procedures to guarantee crash consistency; therefore, I choose to implement it using more expressive RPCs. Carbink is implemented and evaluated atop Google’s datacenter infrastructure. Compared to a state-of-the-art fault-tolerant design that uses span-granularity erasure coding, Carbink has 29% lower tail latency and 48% higher application performance, with at most 35% higher far memory usage (due to asynchronous memory defragmentation). Carbink also results in a joint patent with Google.

CPU efficient load balancing for microsecond-scale RPCs. Datacenter applications are evolving into microservice architectures, with many small services connected via RPCs to serve user requests. To ensure responsiveness, these services require high throughput and low tail latency, reaching millions of operations/sec per server and microsecond-scale latency respectively. This creates a mismatch between existing RPC frameworks and application demands, in terms of efficiently load balancing microsecond-scale RPCs. Conventional Power-of-Two load balancing probes servers’ load too often (i.e., probing before each RPC) and hurts application throughput, as a load probing consumes comparable server CPUs as a microsecond-scale RPC. My measurement shows that it reduces the goodput (i.e., maximum throughput under tail latency SLO) by half compared to naively dispatching RPCs at random. On the other hand, probing too infrequently will result in stale estimates of load, resulting in suboptimal load balancing, the emergence of hot spots, and violated SLOs. To break this dilemma, Mew [3] tailors unnecessary load probings to just fulfill the staleness requirement that does not degrade tail latency. To do so, Mew performs probing statistically following an optimal probing frequency, obtained by running a gradient descent algorithm on the probing frequency vs. tail latency space.

However, there are more challenges in how to efficiently fit RPC load balancing into RPC frameworks. The first is what load signal to use that is general enough to capture different load levels of servers, and is strongly correlated to future RPC’s tail latency. Instead of using the conventional signal of CPU utilization, I use the low-level thread and packet queueing delay, because the former cannot differentiate between the ideal case of exactly-saturated CPUs and the bad case of overloaded CPUs. The second challenge is how to efficiently implement load probing, especially for kernel networking stacks with high overhead. My solution is leveraging eBPF to directly return load signal values in the kernel, without going through the full kernel stack or user space. With all the above designs, Mew is able to reduce RPC tail latency by $2\times$, while achieving $1.7\times$ higher goodput, over a state-of-the-art solution.

Other datacenter infrastructure research:

Evolvable and memory efficient network telemetry. As modern datacenter networks get larger and more complex, operators must rely on network telemetry systems for continuous monitoring, alerting, failure troubleshooting, etc. However, changes happen frequently in production networks (e.g., modifications to monitoring intent, advances of device APIs), impacting the reliability of network telemetry systems. To handle various changes, I helped Facebook develop their evolvable network telemetry system PCAT [5]. PCAT proposes to use a change cube abstraction to systematically track changes, and an intent-based layering design to confine and track changes. The overall result of PCAT is a change-aware network telemetry system that supports fast-evolving datacenter networks at Facebook.

Network telemetry also requires high efficiency for memory. Telemetry data must be stored in memory, at least temporarily, but memory is a precious resource. Network devices (e.g., NICs, switches) often have less than 100MB of memory; server memory is more plentiful, but should be mostly devoted to applications. My Cold Filter [12], Elastic Sketch [13], Pyramid Sketch [14], and more [15, 16, 17] design memory-efficient probabilistic data structures that can be updated at line rate, have low memory footprints, and high accuracy. At the time of this writing, Elastic Sketch is cited over 400 times by follow-up work across many academic research groups (e.g., CMU, Princeton, University of Pennsylvania, Technion, KTH). Some of them try to further optimize its memory usage, speed, or accuracy; some re-purpose its design for more telemetry tasks; and some leverage its implementation for P4 compiler research.

Secure hardware architecture for SmartNICs. Cloud providers are deploying various SmartNICs with wimpy-yet-power-efficient RISC cores to offload simple network functions such as network virtualization and traffic scheduling. Unfortunately, vast cloud tenants are barred from the efficiency benefits of SmartNICs, because they are not allowed to run their own customized functions on SmartNICs. The root cause is that modern SmartNICs provide little isolation between the network functions belonging to different tenants; these NICs also do not protect network functions from the datacenter-provided management OS running on the NIC. My S-NIC [18] project proposes minimal changes to

SmartNIC hardware, so that datacenters can provide offloaded functions with strong isolation, while preserving most of the total-cost-of-ownership benefits with minimal performance degradations. S-NIC's designs target various commodity multi-core SmartNICs, and explicitly isolate their IO subsystems and on-NIC accelerators.

Future Research

Building on my past experiences in networking, memory management, OS kernels, and datacenter applications, I am excited to apply my full-stack optimization approach with cross-layer codesign to the following problems.

Deployment-friendly approaches to memory efficiency via malloc queueing. Previous work on increasing memory efficiency is mostly *not* deployment-friendly, requiring modifying either OS kernels [10] or application code and third-party libraries [11]. In search of deployment-friendly approaches, I have a preliminary insight around separating the provisioning of average memory usage and bursty usage: application's peak memory usage is usually dominated by bursty, large memory allocations (e.g., temporarily loading a large file into memory); if one can time-interleave such allocations from different applications to avoid their memory peaks coinciding with each other, the overall memory provisioning can be reduced, thus improving memory efficiency. One way to implement time-interleaving is overwriting the `Malloc()` function to strategically delay memory allocations, which I believe is far more deployment-friendly than previous work. I call this approach malloc queueing, and it would mostly target batch processing applications whose performance is not sensitive to the incurred memory allocation delays.

eBPF for accelerators and more. eBPF programming language features verified safety and liveness, and has been widely applied to packet processing in kernels and SmartNICs. I intend to extend eBPF to manage heterogeneous hardware accelerators, and build a generic and easy-to-use programming interface between accelerators and application developers. Example accelerators include GPU and FPGA for massively parallel computing, and U2F (Universal 2nd Factor) keys for security. Through verification, this interface would enable strong safety and liveness guarantees for computations running on these accelerators. Besides hardware, I believe eBPF can shed light on more software applications. I intend to explore the following ones: (1) fast task scheduling (e.g., work stealing) for distributed computation framework like Ray [19], and (2) generic shared logs to support various distributed data structures [20]. Both applications would benefit from the efficient network IOs via kernel offloads, and require addressing challenges from the constrained programming model in eBPF.

Resource efficient machine learning. Machine learning (ML) workloads such as the training and inference of Large Language Models (LLMs) are extremely resource-hungry, requiring expensive accelerators like GPUs. I intend to take a full-stack approach to improve the resource efficiency of ML workloads, covering GPU memory efficiency and compute efficiency. One direction is applying far memory techniques to LLM training and inference by swapping to CPU memory. For performance, I plan to codesign far memory swapping with the memory access patterns of LLM weights and key-value cache, e.g., different access frequencies for different weights due to the attention mechanism in LLMs. Another direction is developing a unified GPU memory abstraction that allows easily accessing remote GPU memory over high-speed networks such as NVLink; this kind of GPU memory pooling would help reduce memory stranding and fragmentation caused by dynamic memory allocations in ML workloads. For performance, I plan to codesign such memory pooling with ML workload characteristics, e.g., allowing relaxed consistency. Finally, I am interested in fine-grained GPU kernel scheduling at the microsecond scale possibly with preemption; the goal is to efficiently multiplex GPU compute resources among multiple jobs without losing performance.

Datacenter-scale distributed runtime. A long-term goal of my research is to build a datacenter-scale distributed runtime to not only simplify application development but also increase the whole datacenter efficiency and evolvability. This distributed runtime sits between applications and datacenter resources: (1) for applications, it provides generic and stable interfaces to use compute, memory, storage, and accelerators, and customizable fault tolerance and recovery schemes based on application needs; (2) for resources, it eschews the conventional reservation-based provisioning strategy, and instead provisions resources in a best-effort manner to achieve high resource efficiency.

Today's datacenters have already provisioned network resources in a best-effort manner, and I plan to expand this strategy to cover more resources like compute, memory, storage, and accelerators. For these new best-effort resources, many networking techniques like congestion control can be applied to enable efficient fair sharing. However, unlike the network resources that are delay-tolerable for applications, other resources especially the memory are not (think of out-of-memory errors). To address this challenge, I intend to leverage techniques like far memory and malloc queueing to create a delay-tolerable memory abstraction, at the cost of lower resource utility than normal memory.

References

- [1] **Yang Zhou**, Zezhou Wang, Sowmya Dharanipragada, and Minlan Yu. Electrode: Accelerating Distributed Protocols with eBPF. In *Proceedings of USENIX NSDI*, 2023.
- [2] **Yang Zhou**, Xingyu Xiang, Matthew Kiley, Sowmya Dharanipragada, and Minlan Yu. DINT: Fast In-Kernel Distributed Transactions with eBPF. Under submission.
- [3] **Yang Zhou**, Hassan Wassel, James Mickens, Minlan Yu, and Amin Vahdat. Mew: Efficient Inter-Server Load Balancing for Microsecond-Scale RPCs. Under submission.
- [4] **Yang Zhou**, Hassan Wassel, Sihang Liu, Jiaqi Gao, James Mickens, Minlan Yu, Chris Kennelly, Paul Turner, David Culler, Hank Levy, and Amin Vahdat. Carbink: Fault-Tolerant Far Memory. In *Proceedings of USENIX OSDI*, 2022.
- [5] **Yang Zhou**, Ying Zhang, Minlan Yu, Guangyu Wang, Dexter Cao, Eric Sung, and Starsky Wong. Evolvable Network Telemetry at Facebook. In *Proceedings of USENIX NSDI*, 2022.
- [6] William Tu, Yi-Hung Wei, Gianni Antichi, and Ben Pfaff. Revisiting the Open vSwitch Dataplane Ten Years Later. In *Proceedings of ACM SIGCOMM*, 2021.
- [7] Joshua Fried, Zhenyuan Ruan, Amy Ousterhout, and Adam Belay. Caladan: Mitigating Interference at Microsecond Timescales. In *Proceedings of USENIX OSDI*, 2020.
- [8] Muhammad Tirmazi, Adam Barker, Nan Deng, Md E Haque, Zhijing Gene Qin, Steven Hand, Mor Harchol-Balter, and John Wilkes. Borg: the Next Generation. In *Proceedings of ACM EuroSys*, 2020.
- [9] Chengzhi Lu, Kejiang Ye, Guoyao Xu, Cheng-Zhong Xu, and Tongxin Bai. Imbalance in the Cloud: An Analysis on Alibaba Cluster Trace. In *IEEE International Conference on Big Data (Big Data)*, 2017.
- [10] Emmanuel Amaro, Christopher Branner-Augmon, Zhihong Luo, Amy Ousterhout, Marcos K Aguilera, Aurojit Panda, Sylvia Ratnasamy, and Scott Shenker. Can Far Memory Improve Job Throughput? In *Proceedings of ACM EuroSys*, 2020.
- [11] Zhenyuan Ruan, Malte Schwarzkopf, Marcos K Aguilera, and Adam Belay. AIFM: High-Performance, Application-Integrated Far Memory. In *Proceedings of USENIX OSDI*, 2020.
- [12] **Yang Zhou**, Tong Yang, Jie Jiang, Bin Cui, Minlan Yu, Xiaoming Li, and Steve Uhlig. Cold Filter: A Meta-Framework for Faster and More Accurate Stream Processing. In *Proceedings of ACM SIGMOD*, 2018.
- [13] Tong Yang, Jie Jiang, Peng Liu, Qun Huang, Junzhi Gong, **Yang Zhou**, Rui Miao, Xiaoming Li, and Steve Uhlig. Elastic Sketch: Adaptive and Fast Network-Wide Measurements. In *Proceedings of ACM SIGCOMM*, 2018.
- [14] Tong Yang, **Yang Zhou**, Hao Jin, Shigang Chen, and Xiaoming Li. Pyramid Sketch: A Sketch Framework for Frequency Estimation of Data Streams. *Proceedings of the VLDB Endowment*, 2017.
- [15] **Yang Zhou**, Omid Alipourfard, Minlan Yu, and Tong Yang. Accelerating Network Measurement in Software. *ACM SIGCOMM Computer Communication Review*, 2018.
- [16] Omid Alipourfard, Masoud Moshref, **Yang Zhou**, Tong Yang, and Minlan Yu. A Comparison of Performance and Accuracy of Measurement Algorithms in Software. In *Proceedings of ACM Symposium on SDN Research (SOSR)*, 2018.
- [17] Zhuochen Fan, Gang Wen, Zhipeng Huang, **Yang Zhou**, Qiaobin Fu, Tong Yang, Alex X Liu, and Bin Cui. On the Evolutionary of Bloom Filter False Positives - An Information Theoretical Approach to Optimizing Bloom Filter Parameters. *IEEE Transactions on Knowledge & Data Engineering*, 2022.
- [18] **Yang Zhou**, Mark Wilkening, James Mickens, and Minlan Yu. SmartNIC Security Isolation in the Cloud with S-NIC. Under submission.
- [19] Philipp Moritz, Robert Nishihara, Stephanie Wang, Alexey Tumanov, Richard Liaw, Eric Liang, Melih Elibol, Zongheng Yang, William Paul, Michael I. Jordan, and Ion Stoica. Ray: A Distributed Framework for Emerging AI Applications. In *Proceedings of USENIX OSDI*, 2018.
- [20] Mahesh Balakrishnan, Dahlia Malkhi, Ted Wobber, Ming Wu, Vijayan Prabhakaran, Michael Wei, John D Davis, Sriram Rao, Tao Zou, and Aviad Zuck. Tango: Distributed Data Structures Over a Shared Log. In *Proceedings of ACM SOSP*, 2013.

[This page intentionally left blank.]

Teaching Statement

Yang Zhou

I greatly enjoy the rewards of teaching and mentoring students. For me, the rewards consist of two significant parts: (1) the pride and fulfillment when my teaching helps students carry out their studies smoothly and when my mentored students grow into independent researchers, and (2) the interesting future research directions inspired or confirmed during teaching and mentoring. Driven by these rewards, I have taught as a teaching assistant and as a small-group “supervisor”, and mentored four undergraduates and five junior PhD students in their research. Based on my research background, I am qualified to teach undergraduate courses of computer networks, operating systems, distributed systems, and algorithms and data structures, and graduate courses of data center networking and dataplane operating systems (detailed later).

Mentoring Experience and Methodology

I have mentored four undergraduates on system research, and informally mentored five junior PhDs on their research ideas, internship applications, and study experience at Harvard. Among the four undergraduates, one (Zezhou Wang) published an NSDI’23 paper with me and went to University of Washington (UW) as a system PhD; two of them (Xingyu Xiang and Matt Kiley) co-authored an NSDI’24 submission with me, and are about to apply for system PhDs as well as the rest one. Such mentoring brings me enormous pride, e.g., seeing Zezhou gets into the UW PhD program. It inspires my future research—working with Zezhou on eBPF sparks two follow-up projects: one has become the NSDI’24 submission, and another is showing promising results. Below I summarize my mentoring methodology:

- *Building students’ confidence.* It is well-known that confidence is crucial for students, but how to build their confidence is challenging. One way I find helpful is respecting students’ thoughts by giving them enough freedom to try their thoughts while keeping an eye on the big agendas and goals. Another way is connecting them to experts upon entering a new field, avoiding the steep learning curves overwhelming or destroying their confidence. The experts, who could be the mentors themselves, would point out the proper materials or steps for quick ramp-ups.
- *Encouraging students to form their own opinions and tastes.* I encourage and anticipate students to form their own opinions about systems, develop their own tastes on promising research problems, and stick with them. I do not worry too much about if students’ opinions/tastes are wrong, as once they go deep into specific directions they believe, they will learn extensive experiences and insights to refine their previous opinions/tastes.
- *Collaborating widely.* Wide collaboration across industry and academia is especially beneficial for practical system research, and mentors should play the important role in connecting students with proper researchers in the wild. For example, my fault-tolerant far memory project Carbink was collaborated with Google via my co-advisor’s connections, and then inspired by Google’s desire for high availability. However, collaborating with industry usually requires teasing out real research challenges, while not being misled by massive engineering details; advisors should leverage their experience to help students (especially junior PhDs) navigate efficiently in this space. For another example, my eBPF-for-Paxos project Electrode would not be possible without the collaboration with Sowmya Dharanipragada who is a distributed system PhD at Cornell. Going forward, I would like to expand collaborations to theory, machine learning, architecture, programming languages, etc.

Teaching Experience and Philosophy

System course teaching: I was the teaching assistant (TA) for a computer system course, the Harvard CS145 Networking at Scale, along with an undergraduate TA. This course features eight P4-switch related projects, three of which are designed and developed by me including detailed guides and skeleton code. I held three one-hour sections covering network programming, background knowledge for projects, and handy tools for developing and debugging. Other duties include holding weekly office hours, answering students’ questions on forums, and grading projects. In addition to TA, I also had a guest lecture experience at UC Berkeley on far memory techniques in data centers, mainly facing junior graduate students from architecture areas. I started from common and accessible facts like resource utilization and

DRAM prices, then explained why data center operators have an interest in far memory, and finally discussed my work in this space.

Algorithm course teaching: I was the small-group supervisor for the Algorithm Design and Analysis course at Peking University as an undergraduate. This role requires supervising around 14 students in small classes, giving recitations, teaching advanced algorithms and data structures, preparing new problem sets and quizzes, and grading, all on a weekly basis. I extensively introduced non-textbook topics related to my undergraduate research of probabilistic data structures. Although time-consuming, being such a supervisor is truly gratifying, especially when students understand my research and try various optimizations as their final course projects. One student (Yicheng Jin) in my small class is now pursuing a computer science PhD at Duke University.

Introductory teaching: I taught non-CS audiences about the Internet from a computer science perspective during the English Language Program at Harvard. It was a slightly difficult yet fun experience especially when I told the audience that Internet data is transmitted in small packets: they were shocked and immediately asked why, and then I gave them detailed yet understandable explanations until they grasped the design philosophy behind it. This experience gave me a good sense of how to teach introductory courses in the future. Below I summarize my teaching philosophy:

- *Building safe and inclusive environments.* Students in the same class usually have different prior knowledge; thus it is important to create safe and inclusive environments to make students feel they are welcome to ask both the simplest questions and challenging ones. I got such first-hand experiences when I took my co-advisor James Mickens' CS263 System Security course: it has the most open class environment I have ever seen because of James' unique humor, and students ask so many interesting questions during the class. As a result, I personally learned so much security knowledge, though my research is on networked systems.
- *Focusing on hands-on experiences.* I believe the best way to learn computer systems is through reading, running, debugging, and hacking well-written codebases in a hands-on manner. My personal experience in learning dataplane operating systems exactly follows this pattern: after reading relevant papers, I could not understand how specific designs get implemented and contributed to the final performance; then I decided to read the codebase of a dataplane OS called Caladan [1], and run and debug it; finally, I built my own research prototype atop it. After the process, my understanding of dataplane OSes became much clearer, and I gradually began appreciating the merits of various designs in this space. For future system courses I teach, I would like to incorporate well-written teaching systems, such as the WeensyOS [2], into my agenda to help students gain hands-on experiences.
- *Promoting critical thinking on the pros and cons of techniques.* I learned this from the Harvard CS260r Projects and Close Readings in Software Systems—Serverless Computing by Eddie Kohler, where he discussed serverless computing research from a traditional system research perspective. He showed impressive critical thinking on the pros and cons of serverless computing, and helped us grasp the real novel components of this paradigm without deifying any new terms. I plan to apply a similar philosophy to my teaching, encouraging students to critically think about new techniques around us, such as the emerging LLM techniques.

Course plans: In addition to the aforementioned undergraduate courses based on textbook knowledge, I would like to hold two advanced graduate courses and a seminar course based on my research:

- *Data center networking:* I will discuss how modern data centers design and build high-performance network fabrics including topology, routing, congestion control, fault tolerance, load balancing, etc.
- *Dataplane operating systems:* I will discuss how the OS evolves to keep up with the fast hardware in data centers, including user-space networking, efficient threading, light-weight isolation, etc.
- *System seminar course:* I will invite a broad set of system researchers from both academia and industry to give talks on various system research topics, and foster potential collaborations with students.

References

- [1] The Caladan authors. Caladan opensource. <https://github.com/shenango/caladan>.
- [2] Eddie Kohler. Harvard CS61 Systems Programming and Machine Organization (2023): WeensyOS. <https://cs61.seas.harvard.edu/site/2023/WeensyOS/>.

Diversity Statement

Yang Zhou

I view DEI as the basic soil for growing humanity and excellence in society, including the academic community; it is about the daily respect for people regardless of their self-identifications, and self-introspection on “whether I want to be treated like what I treat others”. Everyone has the duty to foster DEI in her/his surroundings, because that eventually determines how the society will treat them in one day. Here, I would like to sample my and my family’s experiences of being underrepresented due to educational background, language, political affiliation, and ethnic origin, to motivate how I grow awareness of the challenges faced by underrepresented populations and the importance of DEI, and possible ways to foster DEI—some I have adopted and some I plan to do.

I am a first-generation college student, so my parents could hardly give me advice on how to succeed in college and in my PhD studies. However, I was lucky to receive tremendous emotional support from them. I was also fortunate to receive academic mentorship from a variety of professors and student peers. Thus, I am proud to be a faculty job applicant today, and I look forward to creating a sharing and inclusive environment in the classroom and in my research group.

As a first-generation immigrant to the US, one of the first challenges that I faced was mastering the English language. At Harvard, I greatly benefited from the university’s English Language Program (ELP), which offered weekly lectures by experienced English teachers, and recruited native English speakers from the university to serve as language partners. The ELP experience showed me how community building is a critical aspect of helping students integrate into challenging environments. As a professor, I hope to make students aware of programs like the ELP that target specific barriers to students’ success (e.g., language issues, or a lack of adequate high school preparation for college-level classes).

Fifty years ago, my uncle was denied admission to his dream civil aviation university, despite his excellent academic performance and physical fitness. He was rejected because his father (my grandfather) was a combat medic for the Chinese Nationalist Party—the party who had fought with the Communist Party of China that founded the People’s Republic of China. Such political discrimination prevented a whole generation of my uncles from participating in activities that were even slightly related to military service. I was told this experience at a very young age; thus, I have always known that the political environment of the past can influence personal outcomes in the present.

My mother and her family are Hui Chinese, one of the ethnic minorities that comprises 0.79% of the total Chinese population. Being an ethnic minority in China often results in discrimination by the majority Han population. For example, a popular stereotype is that Hui Chinese are thieves. Fortunately, my parents always taught me to not treat people by their ethnicity, race, or religion. As a result, I am always conscious of potential biases that may impact my interactions with others, and I hope to support DEI principles as a professor.

My Past Contributions to Advancing DEI

I have participated in various activities that supported DEI via mentoring and teaching.

Mentoring: During the summers of 2022 and 2023, I mentored four undergraduate students for research internships at Harvard: three came from non-US schools, with two being in the US for the first time. To help the students get familiar with systems research (and life in the US), I held weekly meetings with each student, talking about not only research but also various cultural acclimation challenges that I had experienced during my own PhD. At the time of this writing, one of them has co-authored a paper with me that was published at a premier system conference. This student was also accepted to the University of Washington as a computer science PhD student. The other three students have also decided to apply to systems PhD programs, including one that was hesitating for a long time before working with me. I also consistently (monthly) shared my research and internship experiences with five junior PhD students over the past two years. All of them are non-native English speakers and are non-white.

Occasionally, I received email inquiries from PhDs who are in other research areas or from underrepresented minorities; I often scheduled one-to-one meetings to learn about their difficulties or puzzles. For example, Jessica Quaye, originally from the Republic of Ghana in West Africa, was interested in system research though she is in an architecture research group. I had long meetings with her both in person and online, and introduced her to my co-advisor Minlan Yu to identify potential opportunities for collaboration and advising.

Besides one-to-one mentoring, I also participate in one-to-many panels to share my research experience with junior system PhDs. For example, I was a panelist for the “Getting started with systems research” panel [1] organized by Students@Systems in 2022. The video recording for the panel is freely accessible online to help systems PhD students regardless of their university or physical location.

Talking: Academic networking (e.g., talking to peer researchers at conferences) is crucial to the success of a PhD student. However, junior graduate students are often afraid of professional networking, e.g., due to fears about having little experience or being from less prodigious schools. However, I vividly remember how, at a conference, James Mickens (one of my co-advisers) stood in front of the door of a breakout room and publicly said “I am James, a Professor at Harvard, and you are welcome to talk to me!” This event inspired me to proactively interact with junior students during conferences, to talk about mutual research interests and identify potential collaboration opportunities. I also like to encourage poster presenters for their research, especially when there are no people who are currently engaging with their posters.

I also talk to undergraduates and high school students regarding computer science research. For example, in October 2022, I gave a research talk at a Harvard AM/CS/EE PhD recruitment event (accessible to all US universities) which targeted students “that hold membership in an underrepresented and/or historically minoritized group in STEM.” In 2022, I also gave talks at the Harvard SEAS Undergraduate Research Open House and the SEAS Research Showcase, targeting Harvard freshman and sophomore undergraduates. These talks were well-received, with several undergraduates in the audience later contacting my research lab to learn more about participation opportunities; I still mentor one of these undergraduates. Going back to the time when I was an undergraduate, I had the privilege to talk to juniors in my alma mater high school on why a computer science major is a good college major. Some of these students still contact me for advice.

Teaching: I make an explicit effort to help students with little prior exposure to computer science, and I try to promote inclusiveness during teaching. When I was the small-group “supervisor” for the Algorithm Design and Analysis course at Peking University, I realized that some students lacked high school experience with programming contests; these students often found it hard to catch up with peers who did have this experience. To help them, I wrote step-by-step, thorough explanations for the algorithms discussed in class, and I handed out these explanations after class. When TA’ing a course at Harvard University, I answered all questions that appeared in the Ed forum, no matter whether the questions were anonymous or not, to keep everyone’s learning progress on track.

My Future Plans for Fostering DEI

Going forward, as a faculty member, I plan to take the following actions:

- *Advising:* Actively recruiting underrepresented students, being attentive to any anti-DEI atmosphere in my research group, and explicitly adopting counter-measures to foster DEI with affirmative actions.
- *Connecting:* Reducing the barriers of students finding research opportunities by organizing mutual-connecting programs like UCB DARE [2]—matching students with faculty members for research.
- *Teaching:* Being attentive to any students with weaker prior knowledge in my classes, and helping them build confidence with support on a case-by-case basis.
- *Daily life:* Being kind to people I meet, no matter their age, color, disability, gender, ethnicity, politics, religion, education, language, and more. I believe “kindness is the ultimate nobility” [3].

References

- [1] Student@Systems. A panel on “Getting started with systems research”. <https://students-at-systems.org/pages/events/getting-started-with-systems-research.html>.
- [2] UC Berkeley. DARE: Diversifying Access to Research in Engineering. <https://dare.berkeley.edu/>.
- [3] Amin Vahdat. SIGCOMM Lifetime Achievement Award 2020 Keynote (48m44s): kindness is the ultimate nobility. https://youtu.be/Am_itCzkaE0?t=2924.

Electrode: Accelerating Distributed Protocols with eBPF

Yang Zhou*
Harvard University

Zezhou Wang*
Peking University

Sowmya Dharanipragada
Cornell University

Minlan Yu
Harvard University

Abstract

Implementing distributed protocols under a standard Linux kernel networking stack enjoys the benefits of load-aware CPU scaling, high compatibility, and robust security and isolation. However, it suffers from low performance because of excessive user-kernel crossings and kernel networking stack traversing. We present Electrode with a set of eBPF-based performance optimizations designed for distributed protocols. These optimizations get executed in the kernel before the networking stack but achieve similar functionalities as were implemented in user space (e.g., message broadcasting, collecting quorum of acknowledgments), thus avoiding the overheads incurred by user-kernel crossings and kernel networking stack traversing. We show that when applied to a classic Multi-Paxos state machine replication protocol, Electrode improves its throughput by up to 128.4% and latency by up to 41.7%.

1 Introduction

Distributed protocols such as Paxos [37] for state machine replication are important building blocks for highly-available distributed applications. For example, Google's Chubby [6] uses a variant of classic Paxos [37] and Multi-Paxos [36] to implement a highly-available lock service, powering their business-critical GFS [16] and Bigdata [7] applications. Google's globally-distributed database Spanner [8] and Microsoft's data center management tool Autopilot [22] also run Paxos protocols to maintain their high availability.

Existing high-performance implementation of distributed protocols tends to be radical and not readily-deployable. DPDK-based kernel-bypass approaches [27, 79] allow direct access to the underlying NIC hardware, but require application developers to build their own networking stack and maintain compatibility with the evolving kernel networking stack [75]. DPDK also dedicates CPU cores to busily poll the network interface for I/O competition, sacrificing CPU resources and wasting energy during low I/O loads. This is especially a problem for embedded devices [51, 60, 70] where CPU resources are rare. Other approaches co-design specialized distributed systems with niche network hardware including RDMA [11, 28, 76], FPGA [23], SmartNICs [66], and programmable switches [25]. These advanced hardware devices are not widely available in today's cloud environments, and systems built on top of them tend to be difficult to design, implement, and deploy [27].

Instead, we would prefer the widely-deployed and well-maintained standard kernel networking stack that also provides load-aware CPU scaling and strong security and isolation among different applications [5, 59]. However, implementing distributed protocols under the standard kernel networking stack often gives poor performance. The root causes are the high packet processing overhead in the kernel networking stack and heavy communications in distributed protocols. Our measurement shows that over half of CPU time is spent on the kernel networking stack in a typical Paxos deployment (§2); such overhead is mainly caused by user-kernel crossings (and associated context switches) and traversing the kernel networking stack. Moreover, when using a classic leader-based Multi-Paxos protocol [43, 54] to implement state machine replication, e.g., with five replicas, processing a single request would require the leader node to send/receive *fourteen* messages in total (see Figure 1a), suffering from the kernel stack overhead fourteen times¹.

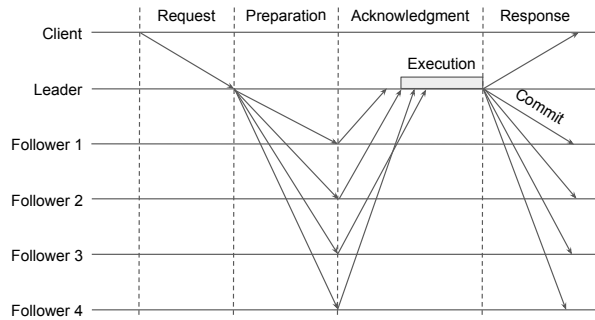
In this paper, we focus on accelerating Paxos protocols inside data centers by offloading protocol operations to the kernel via eBPF (i.e., extended Berkeley Packet Filter) [46, 49]. eBPF allows safely executing customized yet constrained functions inside the kernel at various locations. Similar to kernel bypass, the offloaded operations get executed immediately after the NIC driver receives the packet, without user-kernel crossing and kernel networking stack traversing. Unlike kernel bypass, eBPF is an OS-native mechanism such that eBPF offloaded operations do not sacrifice security and isolation properties while amenable to load-aware CPU scaling without busy-polling.

The key challenge is, given the constrained programming model of eBPF, *which parts of Paxos protocols to offload that can greatly reduce kernel stack overhead while being implementable and efficient in eBPF*. Note that the eBPF subsystem requires every offloaded function to be statically verified to guarantee kernel security, which only allows limited instructions, bounded loops, static memory allocation, etc.

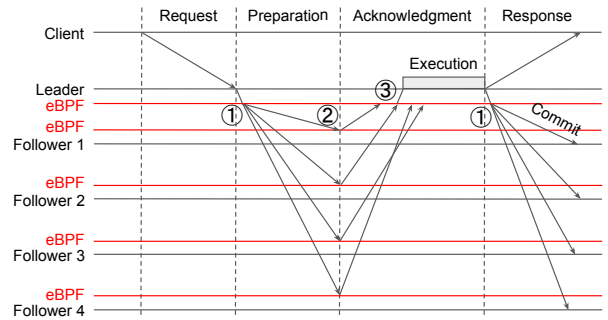
Our insight is that common operations of Paxos protocols, e.g., message broadcasting and waiting on quorums, incur large kernel stack overhead, but are naturally offloadable by existing eBPF programming capacity. For example, Paxos protocols require a leader node to broadcast preparation messages to follower nodes; if implemented using multiple `sendto()` syscalls conventionally, it would incur multiple user-kernel

*Equal contribution

¹Linux `io_uring` [1] can reduce user-kernel crossings, but cannot reduce kernel stack traversing (see §8 for details).



(a) The Multi-Paxos/Viewstamped Replication protocol.



(b) Electrode-accelerated Multi-Paxos/Viewstamped Replication.

Figure 1: Normal case execution of the leader-based Multi-Paxos/Viewstamped Replication protocol vs. Electrode-accelerated one with 5 replicas. Electrode offloads ①: message broadcasting (§4.1), ②: fast acknowledging (§4.2), and ③: wait-on-quorums (§4.3) to eBPF to reduce the kernel networking stack overhead.

crossings and kernel networking stack traversing. Instead, eBPF has a `bpf_clone_redirect()` [45] function that enables us to clone an in-kernel packet buffer multiple times and send them to different destinations; this eBPF-based message broadcasting only needs one user-kernel crossing and one kernel networking stack traversing. Besides broadcasting, we also utilize eBPF to reduce unnecessary wake-ups of user-space applications when waiting on quorums, and optimize how follower nodes handle preparation messages by early acknowledging before entering the kernel networking stack. The final result of these three eBPF-based optimizations is Electrode² (Figure 1b). When applying Electrode to a classic leader-based Multi-Paxos protocol, it achieves up to 128.4% higher throughput and 41.7% lower latency. This translates into up to 112.9% higher throughput and 19.3% lower latency for a Paxos-based transactional replicated key-value store.

Electrode has some limitations: it currently targets protocols implemented in UDP and relies on application-level retransmission to handle packet loss. This works well for Paxos protocols whose requests are usually small enough to fit into a single packet, and data center environments where packet loss is rare [28, 61].

2 Background

2.1 Consensus Protocols

Distributed protocols that coordinate and synchronize among a collection of nodes have become an indispensable part of the modern data center application stack. Storage systems in data centers replicate data for fault tolerance and availability. For instance, Berkeley-DB [55] uses a consensus protocol to replicate its logs over a set of distributed replicas. Transactional storage systems like H-Store [71] and Spanner commit their updates to multiple replicas in order to be more failure resilient. At the heart of most replication-based systems is a consensus protocol [36, 37, 43, 54] that ensures that operations execute in a consistent manner across all replicas.

Here, we consider a set of nodes either functioning as clients or replicas. Clients are the users of a particular application-level service hosted by a collection of replicas. It should also be noted here that clients could often just be other servers within the same data center. Clients submit requests to one or more replicas, which triggers a round of agreement to occur. Paxos is a common protocol that is used to obtain an agreement in the presence of node and network failures.

Since applications often need to reach agreements on many client requests, servers use agreement protocols like Paxos to implement a state machine-based abstraction that requires all the replicas to process the exact same set of client requests in the same order. This log-based state machine abstraction is often optimized by the use of a leader. In a leader-based protocol, all the instances of agreement on client requests are mediated through the leader and the leader also dictates the order of the log.

In Figure 1a, we have an example of VR (Viewstamped Replication), a leader-based Multi-Paxos protocol that uses Paxos for running agreements on individual requests. The leader here is responsible for ordering all client requests by assigning sequence numbers to them, and the followers (non-leader nodes) are responsible for responding to the leader and applying all the updates in the order in which they're sequenced by the leader.

The leader is also responsible for initiating agreement by sending out a preparation message to all the other replicas. The leader then waits for a quorum of acknowledgments from all the other replicas before broadcasting a commit message to all the replicas. A successful iteration of this two-round protocol ensures that all non-failed replicas have the client's request. And the sequence number assigned by the leader determines the order in which all the replicas process this client's request. This pattern of broadcasting and waiting on quorums is common in many distributed protocols [38, 39, 80].

To gain more insights into the performance of the Multi-Paxos/VR protocol under the standard Linux kernel networking stack, we measure the CPU time breakdown of the leader node, shown in Table 1. There is 44.7% +

²Electrode is a Pokémon that has a high speed score.

Function Name	Description	% CPU
<code>__libc_sendto()</code>	User function to send packets.	44.7
<code>sock_sendmsg()</code>	Kernel function to send packets.	32.2
<code>__alloc_skb()</code>	Allocate <code>sk_buff</code> for packets.	4.5
<code>dev_queue_xmit()</code>	Transmit <code>sk_buff</code> .	6.8
bookkeeping	For sock, IP, and UDP layers.	20.9
user-kernel crossing	Interrupt, mode switching, etc.	12.5
<code>__libc_recvfrom()</code>	User function to rcv packets.	11.8
<code>sock_recvmsg()</code>	Kernel function to rcv packets.	5.7
user-kernel crossing	Interrupt, mode switching, etc.	6.1

Table 1: CPU time breakdown for the leader node when running the Multi-Paxos/Viewstamped Replication protocol with 5 replicas. See §7 for measurement setup.

11.8% = 56.5% of time spent on the `__libc_sendto()` and `__libc_recvfrom()` functions, while 20.9% + 12.5% + 6.1% = 39.5% of time spent on user-kernel crossing and kernel networking stack bookkeeping. These numbers concrete our previous motivations that implementing distributed protocols under kernel networking stack incurs significant overhead on user-kernel crossings and kernel stack traversing (while eBPF can potentially save them).

2.2 eBPF and Hooks

BPF (i.e., Berkeley Packet Filter) [49] enables user-space applications to customize packet filtering in the kernel. A BPF program, written in some predicates on packet fields, is triggered by the kernel event that a packet arrives at a NIC driver. Once triggered, the BPF program will run inside a kernel virtual machine with limited registers and scratch memory, and a reduced instruction set [49]. For example, the well-known *tcpdump* [20] command-line packet analyzer is based on BPF.

eBPF extends the BPF by increasing the number of registers and adding stack memory. The increased number of registers and stack memory enable the eBPF program to execute more complex operations—the developers can use a C-like language to express customized operations. This C-like code is compiled into an eBPF bytecode by the Clang/LLVM toolchain and runs inside the kernel virtual machine via just-in-time compilation.

eBPF also introduces various powerful in-kernel data structures called *eBPF maps*, which, paired with various helper functions, are used to store and maintain states across multiple triggering of eBPF programs. Example eBPF maps include array, per-CPU arrays, queues, stacks, and hashMaps [46]. These maps are also used to communicate among different eBPF programs and between eBPF programs and user-space processes. Each eBPF map can be identified by a `map_path` through the file system, e.g., `/sys/fs/bpf/<map_name>`, and user-space processes can access a map based on its path.

The kernel events that can trigger eBPF programs are called *eBPF hooks*. There are many hooks existing in Linux kernels

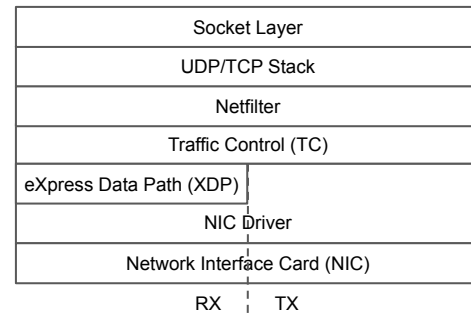


Figure 2: Linux kernel networking stacks and eBPF XDP/TC hooks.

and various device drivers, such as hooks in NIC drivers right after it receives a packet. User-space applications can attach eBPF programs to these eBPF hooks to customize the handling of corresponding kernel events.

Constrained programming model: An eBPF program needs to go through strict verification by an in-kernel eBPF verifier before attaching to an eBPF hook and running inside the kernel. The verification process does a static sanity check to make sure the eBPF program does not have out-of-bounds memory access (i.e., safety) and will always terminate (i.e., liveness). The verifier basically enumerates all possible cases of every conditional branch and loop to make sure every execution path meets the safety and liveness requirements. Because the verification tends to be time-consuming, each eBPF program can only contain up to 1 million instructions. For a larger eBPF program, the developer needs to split it into multiple smaller eBPF programs and uses *tail calls* to let one eBPF program call another one in a continuation manner.

Because of the strict verification process, dynamical memory allocation is not supported in eBPF programs; instead, eBPF programs can only rely on eBPF maps with capacity *specified statically* to maintain in-kernel states.

Due to these limitations, eBPF is commonly used in kernel tracing, profiling, and monitoring [3, 63] and L2-L4 low-level packet processing such as load balancing [14].

XDP (eXpress Data Path) [21, 64] technique implements an in-kernel eBPF hook that enables attached eBPF programs to process RX packets directly out of the NIC driver (Figure 2). Such processing gets triggered before any `sk_buff` [31] allocation or entering software socket queues, thus bypassing any higher-level networking stacks (e.g., UDP, TCP, Socket). XDP-based packet processing normally achieves comparable throughput and latency as DPDK-based kernel-bypass packet processing [21].

TC (Traffic Control) [47] is another important layer/hook which locates right after the XDP (Figure 2). In the TC layer, the `sk_buff` data structure has already been allocated by the kernel networking stack, thus the performance of TC-based packet processing will be slightly worse than XDP. However, the TC hook allows attached eBPF programs to process both RX and TX packets and manipulate the packet `sk_buff`. For

example, one can clone the `sk_buff` for a TX packet and thus implements packet broadcasting in the TC layer.

3 Electrode Overview

Electrode is a framework for offloading Paxos protocols under kernel networking stack to in-kernel eBPF programs to reduce user-kernel crossings and kernel networking stack traversing. Electrode has two goals in designing its eBPF offloads: 1) largely reducing kernel stack overhead to improve performance, and 2) carefully partitioning user- and kernel-space functionalities to keep offloads implementable and efficient inside the eBPF subsystem.

To achieve the first goal, Electrode carefully extracts generic and performance-critical fast-path operations from Paxos protocols to offload to the eBPF. As shown in Figure 1b, Electrode offloads message broadcasting (§4.1), fast acknowledging (§4.2), and wait-on-quorums (§4.3). These operations, if purely implemented in the user space, would involve many user-kernel crossings and kernel stack traversing, causing significant kernel stack overhead as shown in §2. Once implemented in the eBPF, message broadcasting allows the leader node to efficiently send preparation and commit messages to multiple follower nodes, by cloning and sending packets in the kernel; fast acknowledging enables follower nodes to buffer preparation messages in the kernel, and quickly respond to the leader node without involving user-space processes; wait-on-quorums lets the leader node eBPF program wait for a quorum number of acknowledgments from follower nodes, and only notify user-space processes once. Moreover, to simplify how user-space applications use these eBPF-based accelerations, Electrode further designs a set of user-space APIs (Table 2). Each API corresponds to one operation that Electrode offloads to the eBPF, and is used to invoke the offloaded function or retrieve eBPF processing results.

To achieve the second goal, Electrode keeps complicated slow-path operations of Paxos protocols in the user space. Specifically, Electrode leaves the procedures of failure recovery and handling message loss/reordering (i.e., gap agreement) to user-space applications, using similar mechanisms as VR [43] and NOPaxos [40]. These procedures involve accessing dynamic ranges of memory, which is hard to implement in eBPF under the static verification (see §8 for details).

Overall, Electrode has the following workflow: first, user-space applications attach eBPF programs to various hook locations corresponding to a network interface; then, user-space applications use Electrode APIs to invoke eBPF-offloaded functions or retrieve eBPF processing results; finally, the eBPF programs intercept and process target packets in the kernel without going through the networking stack or user-space applications (i.e., Paxos protocols in our case). Electrode targets accelerating the handling of messages that can fit into one ethernet packet (i.e., up to 9KB for jumbo frames). This is well-suited for locks, barriers, and configuration parameters [25, 78] that Paxos protocols commonly maintain. Non-

target packets still go through the stack and reach user-space applications, without impacting applications' other operations or protocol semantics.

Finally, we note that Electrode does not aim to offload every operation of Paxos protocols to the eBPF, because of eBPF's constrained programming model vs. the diverse set of operations that Paxos protocols and related services could have. For example, currently, Electrode does not offload client-facing request/response handling. There are two reasons: 1) Paxos clients normally serialize/deserialize their requests using widely-used libraries such as protocol buffers [19]; however, parsing or constructing protocol buffers is difficult in eBPF, because it involves complex pointer arithmetics and conditional branches which cannot easily pass the eBPF verifier. 2) client-facing requests/responses are normally embedded into application-level services like the Chubby lock service [6], but it is hard and inefficient to implement them in eBPF because of the strict eBPF verifier and the lack of dynamic memory allocation. We discuss more on Electrode's offloading decisions in §8.

4 Electrode Designs

4.1 Message Broadcasting in TC

In Paxos protocols, one-to-all message broadcasting is widely used. For example, 1) the leader node sends preparation messages to all follower nodes, and 2) (after receiving enough acknowledgments from followers) the leader node sends commit messages to all follower nodes.

To implement the above message broadcasting, the most common way is sending the same message multiple times in the user space to different destinations. However, the overhead (i.e., user-kernel crossing and kernel networking stack traversing) of this implementation on the leader node increases linearly as the number of followers increases, while the overhead on each follower node remains constant. Thus, the leader node essentially becomes the system bottleneck, e.g., Table 1 has shown that 44.7% of CPU time is spent on sending messages on the leader node.

An alternative implementation is to use IP multicast [42, 68, 77]. However, IP multicast normally requires support from the underlying network switches (e.g., storing a large number of multicast group-table entries for the whole network topology) [68, 77] or considerable modifications of the Linux networking stack [42].

Electrode approach: Electrode provides a flexible host-based broadcasting solution by utilizing eBPF on the TC hook. Here, we require the eBPF program that implements broadcasting operations to attach to the TC hook, because only the TC hook can intercept and process outgoing packets (§2.2). After attaching the eBPF program, user-space applications can call the `elec_broadcast()` function shown in Table 2 with specified `sock_fd`, message, and a list of destination IPs to broadcast the message to these destinations through the socket.

Function Name	Arguments	Output	Description
<code>elec_broadcast</code>	<code>sock_fd, message, {dst_ips}</code>	status	Broadcasts <message> to all destinations through <sock_fd>
<code>elec_poll_message</code>	<code>map_path</code>	messages	Polls buffered messages from an eBPF-maintained in-kernel ring buffer identified by <map_path>
<code>elec_check_quorum</code>	<code>received_message</code>	bool	Checks if <received_message> (acknowledgment) indicates quorum reaching

Table 2: Electrode user-space APIs.

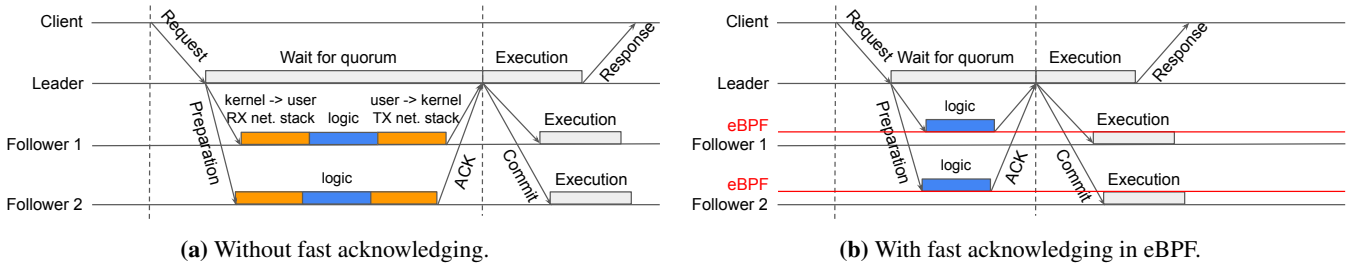


Figure 3: Fast acknowledging in eBPF reduces Paxos request latency. This example follows Figure 1, but omits followers 3 and 4 for brevity.

Under the hood, the eBPF program makes clones of the message packet using the `bpf_clone_redirect()` [45] helper function, modifies the destination addresses of cloned packets accordingly, and sends these packets out. The benefit of cloning packets and broadcasting in the kernel compared with sending the same message multiple times in the user space is that we only need to cross the user-kernel boundary and traverse the UDP and socket layer once.

Handling message loss: Electrode relies on application-level timeout and retransmission to handle message loss, similar to modern RPC-based applications [13, 69]. Specifically, if the leader node does not receive a response after a certain time of sending a request, it will resend the request; once a request experiences several timeouts, the leader node will mark the destination node as dead and start Paxos failure recovery. An alternative approach to handling message loss is doing retransmission in the kernel, which could save user-kernel context switching overheads, but such savings become marginal as packet loss happens rarely in data centers [28, 61]; it would also involve complex message buffer management in kernel/eBPF, hurting performance.

4.2 Fast Acknowledging in XDP

As shown in Figure 3a, a significant portion of Paxos request latency comes from the round-trip delay between the leader node and follower nodes. Note that the ACK messages in this figure mean Paxos protocol acknowledgments, not TCP acknowledgments. For Paxos protocols under the kernel networking stack, this round-trip delay includes not only physical propagation and transmission delay, but also the delay caused by the kernel networking stack (i.e., user-kernel crossing and networking stack traversing). As the fabric latency of nowadays data center network reaches a few tens of microseconds [48] or sub-ten microseconds [18, 27], the latency of the kernel networking stack, which is also around sub-ten microseconds [59], becomes non-negligible.

Electrode approach to reducing the Paxos request latency is to optimize the preparation handling in follower nodes by directly buffering the preparation messages into an in-kernel log and early acknowledging to the leader node. At the same time, user-space applications asynchronously poll and consume the buffered messages from the log, using the `elec_poll_message()` function shown in Table 2. Under the hood, the function calls a corresponding eBPF syscall to poll messages in batches, amortizing kernel crossing overhead. This asynchrony does not break the correctness of Paxos protocols because as long as a preparation message gets buffered into the log, it will be eventually processed by the user-space Paxos protocols, and the message processing order has been specified by the sequence number assigned by the leader node. Figure 3b shows that this approach removes *two* user-kernel crossings and networking stack traversing from the critical path of the Paxos request.

Note that not every preparation message can be handled using fast acknowledging; in some non-critical path cases (e.g., message loss/reordering, and node failure) where the eBPF program cannot handle because of its constrained programming model, our eBPF program can detect them and directly forward preparation messages to user-space Paxos protocols (detailed in §6).

In-kernel log implementation: The in-kernel log temporally stores incoming early-acknowledged preparation messages, which are polled and consumed by user-space applications concurrently. To implement this in-kernel log, we use a special eBPF map named `BPF_MAP_TYPE_RINGBUF` [30] (introduced from Linux kernel 5.8). This map implements an efficient multi-producer single-consumer (MPSC) ring buffer using shared memory and a lightweight spinlock, where we can have multiple writers in eBPF and one reader in user space. Based on our measurement, the time of pushing a preparation message into the ring buffer is roughly equal to memcpying this message, in cases without any lock contention. Note

that the in-kernel ring buffer also has a fixed size, because eBPF does not support dynamic memory allocation; in case it becomes full, the eBPF program can detect them and directly forward preparation messages to user-space applications.

4.3 Wait-on-Quorums in TC + XDP

Another common operation in Paxos protocols is the leader node waiting for a quorum number of acknowledgments (ACKs) from follower nodes (i.e., wait-on-quorums). Assume there are $2f + 1$ replicas including one leader node and $2f$ follower nodes. In most Paxos protocols, once the leader collects f ACKs from different follower nodes, the Paxos request is considered *committed*.

Conventionally, wait-on-quorums is implemented by the user-space applications that receive all ACKs and count towards the quorum number. However, each acknowledgment handling incurs the overhead of the user-kernel crossing and traversing the kernel networking layer. The total overhead of handling all ACKs is linear to the number of follower replicas (i.e., $2f$). Moreover, among these $2f$ ACKs, only the first f ones are required to commit a Paxos request.

Electrode approach: Electrode moves the leader-side wait-on-quorums operations to the eBPF, requiring only one user-kernel crossing and one networking stack traversing. Electrode maintains an array of bitsets (and other metadata) in eBPF, each of which indicates whether a Paxos request has reached the quorum. Electrode only forwards ACK messages that indicate reaching the quorum to the user-space applications, while dropping others. Electrode maps each Paxos request to a specific bitset by using the unique increasing sequence number assigned by the leader node (§2). Note that we use the bitset instead of a counter to check if the quorum gets reached; this is because a timed-out preparation request could cause duplicate ACK messages from follower nodes, and we want to avoid double counting.

Electrode maintains the bitset setting and clearing (i.e., zeroing out) operations through two eBPF programs hooked at TC and XDP layers, respectively. The TC-hooked eBPF program intercepts each outgoing preparation message and clears the indexed bitset, while the XDP-hooked eBPF program intercepts each incoming ACK message from follower nodes and sets the bit corresponding to the follower node's index in replicas.

As shown in Listing 1, the `tc_ebpf` function/program intercepts each outgoing preparation message and clears a specific bitset indexed by the sequence number in each message. Line 6 checks if it is the first time to intercept a preparation message corresponding to this Paxos request, by comparing the `seq` stored along this bitset and the `seq` extracted from the message; if so, it updates the stored `seq` in the array and clears the bitset that may have been used by previous Paxos requests (line 17-18).

The `xdp_ebpf` program intercepts each incoming ACK message, updates the indexed bitset, drops most of the ACK

```

1 # Processing outgoing preparation message
2 # pkt: the packet of the message
3 # seq: unique increasing sequence number (from pkt)
4 def tc_ebpf(pkt, seq):
5     idx = seq % array_length
6     if array[idx].seq != seq
7         array[idx].seq = seq
8         array[idx].bitset.clear()
9     forward(pkt) # to follower node
10
11 # Processing incoming ACK message
12 # pkt : the packet of the message
13 # seq : unique increasing sequence number (from pkt)
14 # node_i: follower node index (from pkt)
15 def xdp_ebpf(pkt, seq, node_i):
16     idx = seq % array_length
17     if array[idx].seq == seq
18         array[idx].bitset.set(node_i)
19         if array[idx].bitset.count() == f
20             pkt.mark_quorum_reach(true)
21             forward(pkt) # to user-space application
22         else: drop(pkt)
23     else: # bitset overwritten by tc_ebpf
24         pkt.mark_quorum_reach(false)
25         forward(pkt)

```

Listing 1: Maintaining the fixed-length bitset array to achieve wait-on-quorums in eBPF. Each bitset operation is also protected by a spinlock; we omit it here for simplicity.

packets, and only forwards packets to user-space applications that indicate reaching quorum or array overflow (explained in the next paragraph). Lines 17-18 check if this bitset corresponds to the `seq` in the ACK message, and set the proper bitset bit if so. Line 19 further checks if this ACK message reaches the quorum: if so, lines 20-21 will mark the packet as quorum-reaching and forward it to user-space applications; otherwise, line 22 just drops the packet. Once the user-space applications receive a quorum-reaching packet—checked by calling the `elec_check_quorum()` function shown in Table 2, it can directly consider this Paxos request as committed.

Handling array overflow: In some cases, a bitset might be overwritten by the `tc_ebpf` because of the fixed size of the bitset array. `xdp_ebpf` detects such array overflow in lines 17&23; once detected, lines 24-25 will mark the packet as *not-quorum-reaching* and forward it to user-space applications. Once the user-space applications receive a not-quorum-reaching packet, it resends the preparation messages to all follower nodes and waits for ACKs again. In practice, the leader node could limit the number of in-flight preparations while provisioning a large bitset array, such that the array overflow does not normally happen.

RSS: Electrode supports RSS (Receive-Side Scaling) which distributes incoming packets to different NIC queues and CPU cores. Specifically, Electrode has two receive-side optimizations: fast acknowledging and wait-on-quorums. For fast acknowledging, the eBPF programs in the follower node could maintain separate in-kernel ring buffers on different cores to avoid synchronization overhead during log append-

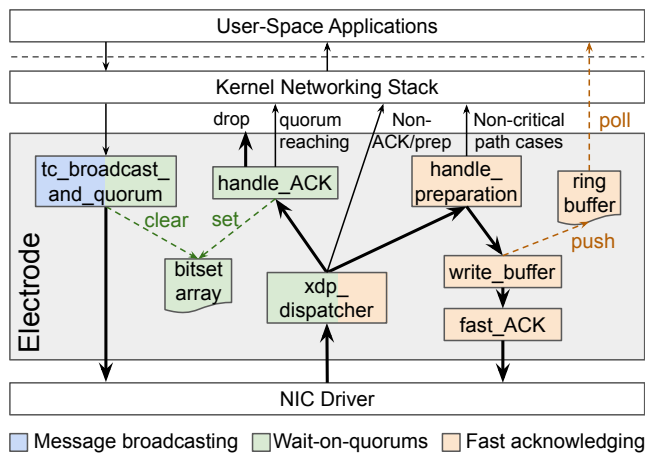


Figure 4: eBPF program structure of Electrode. The thickness of solid lines indicates traffic volume (the thicker, the higher).

ing, and use spinlocks to synchronize accesses to small shared in-kernel states (e.g., `ebpf_seq` in §6); the user-space applications asynchronously pull messages from all ring buffers, and process messages following the order specified by their embedded sequence numbers. For wait-on-quorums, the eBPF programs in the leader node could use atomic instructions to count how many ACKs it has received and check if the quorum is reached.

5 Electrode Implementation

Electrode is prototyped with six eBPF programs written in a restricted C language, and we utilize the Clang/LLVM toolchain for compiling source code to eBPF bytecode. These eBPF programs consist of 500 lines of C code in total. Application developers can also customize their own eBPF programs based on needs, e.g., only processing packets with a specific source port like [25]. Our prototype does not implement the RSS handling yet.

Figure 4 shows the structure of the six eBPF programs. One program can transfer its control flow to the next program via the eBPF tail call. We break the implementation into these six programs because of 1) avoiding breaking the instruction limits in the eBPF verifier (§2.2), and 2) modularity. In the following, we describe each program in detail.

- `tc_broadcast_and_quorum`: This program intercepts outgoing preparation messages. It implements the message broadcasting mechanism (§4.1) and the `tc_ebpf` function in Listing 1 for wait-on-quorums (§4.3). For broadcasting, we generate multiple clones of the preparation packets using the `bpf_clone_redirect()` [45] helper function.
- `xdp_dispatcher`: This program checks the types of incoming messages and calls corresponding message handlers. It only intercepts the ACK (only received on the leader node) and preparation (only received on follower nodes) messages, and calls the corresponding `handle_ACK` and `handle_preparation` programs. It directly forwards

other types of messages to user-space applications.

- `handle_ACK`: This program implements the `xdp_ebpf` function in Listing 1 for wait-on-quorums (§4.3). In common cases, it drops most ACK messages, and only forwards the quorum-reaching ACK messages to user-space applications.
- `handle_preparation`: This program implements various checks to detect non-critical path cases where it should forward messages to user-space applications (§4.2). In normal cases (mostly), it will call `write_buffer` to begin `fast_ACK`.
- `write_buffer`: This program stores message/packet data into an in-kernel log for user-space applications to poll and consume. As mentioned earlier, We use the eBPF ring buffer [30] to implement the log data structure. This program then calls the `fast_ACK` program.
- `fast_ACK`: This program reuses and modifies the received packet buffer to create an ACK packet and sent it out. This requires swapping the src-dst IP addresses and filling the corresponding fields of the ACK message.

6 Apply Electrode to Multi-Paxos

Optimizing throughput: We apply the eBPF-based message broadcasting (§4.1) and wait-on-quorums (§4.3) mechanisms to the leader node in the Multi-Paxos protocol. This implies two throughput optimizations: 1) when the leader node sends out preparation messages to follower nodes, it relies on eBPF to broadcast these messages instead of sending them one by one; and 2) when the leader node is waiting for a quorum number of ACK messages from follower nodes, it only needs to process the quorum-reaching ACK message while the other ACK messages are pruned/dropped by the eBPF program. These two optimizations largely reduce the number of user-kernel crossings and kernel networking stack traversing, thus alleviating the CPU bottleneck on the leader node and improving system throughput.

Optimizing latency: We apply the eBPF-based fast acknowledging mechanism (§4.2) to each follower node in the Multi-Paxos protocol. In normal cases (e.g., without packet loss/reordering, and all nodes are alive), the preparation messages from the leader node are quickly buffered and acknowledged by the eBPF program in the follower nodes, bypassing both the kernel networking stack and the user-space Multi-Paxos protocol. This reduces the commit latency of each Multi-Paxos request by twice the time of user-kernel crossing and kernel networking stack traversing.

Detecting non-critical path cases in fast acknowledging: As mentioned in §4.2, there are some non-critical path cases in fast acknowledging where the eBPF program must detect them and forward the incoming packets to the user-space Paxos protocols. To understand why non-critical path cases happen and how to detect them, we first elaborate on the Multi-Paxos/VR protocol shown in §2, following the literature [43]. In the Multi-Paxos protocol, the leader node assigns

each Multi-Paxos request a unique and strictly increasing sequence number, `seq`. Each replica including both the leader node and follower nodes maintains locally a view number, a status, and its last observed `seq`; each message sent by a replica will piggyback these three variables. The view number indicates which (leader) election epoch this replica is in; the status indicates if this replica is during a leader election (`status_viewchange`), recovering (`status_recovering`), or normal state (`status_normal`). This protocol requires a follower node to only process a preparation message if the node is in the normal state, and the message has a matched view and strictly increasing `seq`; otherwise, the follower node needs to drop the message, or execute a complex view-change or state-transfer procedure [43,54]. Therefore, the non-critical path cases for Multi-Paxos are:

1. the follower is during a leader election or recovering,
2. the follower receives a message with an unmatched view that is either (a) stale or (b) newer,
3. the follower receives a message with a non-strictly-increasing `seq` caused by message (a) loss/reordering or (b) duplication.

These cases only happen when replicas fail or join, or messages get lost/reordered, which is less common in data centers [27,61].

To detect these non-critical path cases in eBPF, we maintain an `ebpf_status`, an `ebpf_view`, and an `ebpf_seq` variable in the eBPF program using the eBPF map. In particular, these three variables can be updated by the user-space Multi-Paxos protocols to reflect the current protocol state. Listing 2 shows the detection pseudocode. Line 5 detects case 1, and line 6 detects case 2(a); for these two cases, the eBPF program needs to drop the packet. Line 7 detects cases 2(b) and 3(a), and forwards the packet to the user space to execute the view-change or state-transfer procedure. For case 3(b), i.e., `msg_seq < ebpf_seq + 1`, the eBPF program function replies an ACK (line 11), because it could be a re-transmitted preparation message due to timeout.

Handling the cases 2(a)&3(a) in fast acknowledging is tricky, because it (i.e., forwarding packets to the user space for processing) involves the concurrency between the user-space protocols and the kernel-space eBPF program, while eBPF only supports map-based communication *but not synchronization* between the user and kernel. Our approach is to let the user-space protocols *detach* the eBPF program from the hook while executing the view-change or state-transfer procedure. Specifically, once a user-space protocol receives a preparation message corresponding to the case 2(a) or 3(a), it detaches the eBPF program, then it finishes the view-change or state-transfer procedure, next it updates the `ebpf_status`, `ebpf_view`, and `ebpf_seq` properly, and finally it reattaches the eBPF program. This guarantees the cases 2(a)&3(a) are exclusively handled by the user-space protocol, avoiding the synchronization between the user and kernel. An alternative approach to achieving the same effect as eBPF detach-reattach

```

1 # pkt      : the packet of the preparation message
2 # msg_view: view piggybacked by the pkt
3 # msg_seq : unique increasing sequence number (from pkt)
4 def detect_non_crit_path_cases(pkt, msg_view, msg_seq):
5     if (ebpf_status != status_normal): drop(pkt)
6     if (msg_view < ebpf_view): drop(pkt)
7     if (msg_view > ebpf_view or msg_seq > ebpf_seq + 1):
8         forward(pkt)
9     if (msg_seq == ebpf_seq + 1):
10         append_log(++ebpf_seq, pkt)
11     reply_ack(pkt)

```

Listing 2: Detecting non-critical path cases during fast acknowledging for Multi-Paxos. Assume the protocol works in a single core, in line with prior Paxos work [40,44,61].

is to use an eBPF map with a branch testing before any Electrode logic. The first packet in the non-critical path can update this map atomically and let all following packets directly go to the user-space application (i.e., closing Electrode optimizations); later, the user-space application can update this map to reopen Electrode optimizations.

There are a few caveats: 1) After the user-space protocol detaches the eBPF program, it needs to poll the in-kernel ring buffer again, in case the eBPF program still appends a few messages to the ring buffer before detaching. Note that the eBPF map can outlive the eBPF program, as long as the user-space process holds a reference to it, because its lifetime is managed through reference counting [50]. 2) While the user-space protocol is setting the `ebpf_seq` value and is about to reattach the eBPF program, some preparation packets might just pass the eBPF hook location but have not been processed by the user-space protocol, e.g., queued in the socket layer. In this case, the user-space protocol actually has set a smaller `ebpf_seq` value in the map; once the eBPF program gets reattached, it will trigger more case 3(a) (lines 7&8). Our solution to this problem is: after the user-space protocol finishes the view-change or state-transfer procedure, it first sends a `stop_sending_preparation` message to the leader node to stop it from sending preparation messages, then it polls the socket to drain and process any queued packet, next it sets the proper `ebpf_seq` value, finally it sends a `resume_sending_preparation` message to the leader node to resume sending preparation messages, and reattaches the eBPF program. These two messages should be sent using reliable transport like TCP to handle packet loss.

Generalizability: Electrode’s eBPF-based optimizations are generic to many more distributed protocols, which normally consist of broadcasting and wait-on-quorums operations. More discussions can be found in Appendix A.

7 Evaluation

This section answers the following questions:

1. How do Electrode and each optimization improve the performance of the Multi-Paxos protocol (§7.1 and §7.2)?

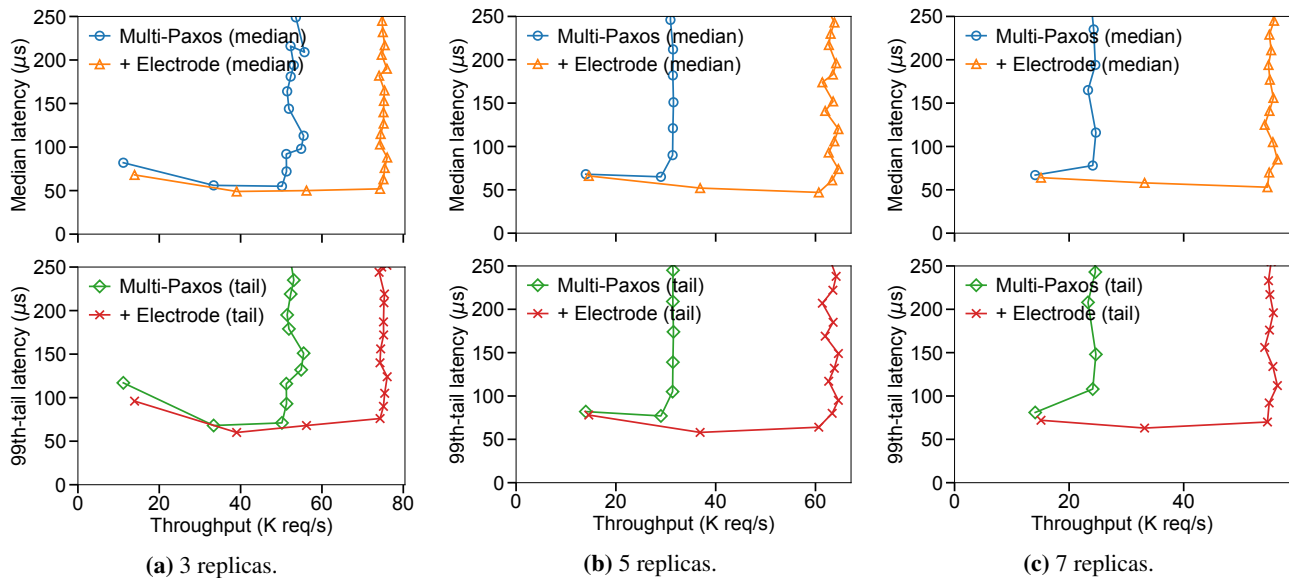


Figure 5: Performance comparison of the Multi-Paxos protocol vs. Electrode-accelerated one with different numbers of replicas.

2. How does Electrode improve the performance of real-world Paxos-based applications (§7.3)?
3. How does Electrode save kernel stack overhead (§7.4)?
4. How does Electrode compare to kernel-bypassing (§7.5)?

Testbed setup: We use eight x1170 servers from Cloud-Lab [12], each of which has a ten-core Intel E5-2640v4 CPU at 2.4 Ghz, 64GB memory, and a Mellanox ConnectX-4 25 Gbps NIC. Each server runs an unmodified Ubuntu 20.04 OS with kernel v5.8.0. All servers are connected using a two-level topology: five Mellanox 2410 as rack switches (each connecting to forty x1170 servers) and one Mellanox 2700 as the spine switch. One server is dedicated as the client server that generates Paxos requests, and other servers run the Paxos protocol with 3/5/7-replica configurations. By default, we configure each server to use one core for interrupt processing and another core for Paxos processing, following the performance optimizations in [41]. We disable irqbalance to avoid out-of-order packet deliveries as much as possible (which would hurt Paxos performance), in line with prior Paxos work [40, 44, 61]. Unlike prior Paxos work [32, 40, 61], we *do not* use IP multicast which requires specialized support from the network (§4.1).

Measurement methodology: The client server runs multiple Paxos/application clients, and each client sends Paxos/application requests in either a closed-loop or open-loop manner. In closed-loop experiments, each client sends the next request once it receives the response of the last request; we vary the number of clients and measure the corresponding throughput, and median and 99th-percentile tail latency, in line with prior Paxos work [40, 44, 61]. In open-loop experiments, each client sends requests one by one at a specific time interval, such that the overall request rate reaches a specified value; we use enough clients (i.e., they could saturate the Paxos servers), specify different request rates, and measure the corresponding

CPU utilization of each replica node.

Comparisons: We use the Multi-Paxos/VR protocol implementation in the SpecPaxos [61] open-sourced code [35] as the baseline, and optimize it using Electrode. We also run a transactional replicated key-value store similar to the one in SpecPaxos [61] atop the baseline Multi-Paxos protocol and Electrode-accelerated Multi-Paxos protocol. All implementation uses the standard UDP stack and socket layer from the Linux kernel.

7.1 Overall Results

Figure 5a, 5b, and 5c show the performance comparison of the Multi-Paxos protocol and the Electrode-accelerated one when using 3, 5, and 7 replicas, respectively. In each figure, we vary the number of clients sending Multi-Paxos requests in a closed-loop manner, and report throughput and median and 99th-percentile tail latency. All curves eventually hit a “hockey stick” in their median or tail latency growth when the system reaches its maximum throughput.

Throughput: the Electrode-accelerated Multi-Paxos protocol achieves 34.9%, 104.8%, and 128.4% higher maximum throughput than the original Multi-Paxos protocol under 3, 5, and 7 replicas, respectively. The large throughput improvements benefit from the eBPF-based broadcasting and wait-on-quorums which reduce the kernel stack overhead significantly on the leader node. With more replicas, the improvement becomes more significant. This is because, for each Multi-Paxos request, the leader node will send more preparation and commit messages, and handle more ACK messages; thus the eBPF-based broadcasting and wait-on-quorums can save more user-kernel crossings and kernel networking stack traversing.

Latency: the Electrode-accelerated Multi-Paxos protocol achieves 12.5%, 20.0%, and 25.6% lower median latency than the original Multi-Paxos protocol with 2 clients (before

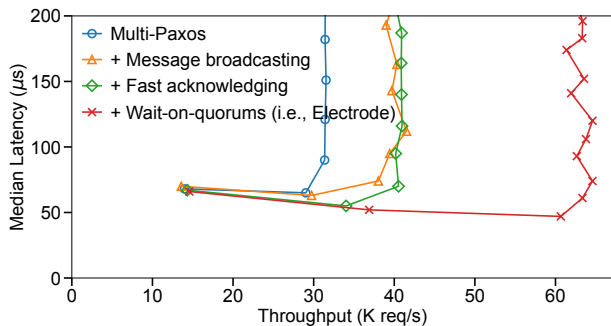


Figure 6: Performance impact of different optimizations for Electrode-accelerated Multi-Paxos protocol (with 5 replicas).

the “hockey stick”) under 3, 5, and 7 replicas, respectively; the corresponding tail latency is 11.8%, 24.7%, and 41.7% lower. The latency reduction mostly comes from the fast acknowledging in the follower nodes, which, for each Multi-Paxos request, saves the time of two user-kernel crossings, two kernel networking stack traversing, and one wake-up of the user-space process. With more replicas, the latency reduction becomes larger. This is because the fast acknowledging bypasses user-space process scheduling and avoids unpredictable scheduling delays [48] by the OS; for the original Multi-Paxos, with more follower nodes, such unpredictable scheduling delays would raise the chance of follower nodes straggling, thus increasing commit latency. Besides, for Multi-Paxos under 3/5 replicas and Electrode under 7 replicas, their latency curves first decline a bit and arrive at the lowest point, then rise and reach the “hockey stick”. This is because, under lower throughput, the Linux scheduler would schedule the Paxos process off the CPU more frequently, while under higher throughput, the Paxos process is mostly scheduled on the CPU.

7.2 Performance Gain Breakdown

Figure 6 shows the performance impact of different optimizations for the Electrode-accelerated Multi-Paxos protocol with 5 replicas. Similar to §7.1, we vary the number of clients sending Multi-Paxos requests in a closed-loop manner, and report the throughput and latency. eBPF-based message broadcasting improves the maximum throughput of the Multi-Paxos protocol by 31.7%; fast acknowledging further reduces the median latency by 4.3%-12.7% (before the “hockey stick”); finally, wait-on-quorums improves the maximum throughput by 57.7%. Overall, we find that the two throughput optimizations (i.e., eBPF-based message broadcasting and wait-on-quorums) have almost no impact on the median latency, while the latency optimization (i.e., fast acknowledging) does not nearly impact maximum throughput. This division of labor demonstrates good modularity of each optimization design in Electrode, and each design can be independently used to accelerate more distributed protocols as shown in Table 4.

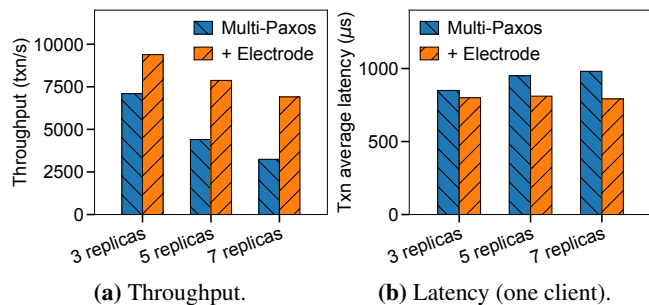


Figure 7: Performance comparison of a transactional key-value store atop the Multi-Paxos protocol vs. Electrode-accelerated one.

7.3 Application Performance

To demonstrate how Electrode can bring benefits to real-world Paxos-based applications, we run a transactional replicated key-value store (similar to the one in SpecPaxos [61]) atop the Multi-Paxos protocol and Electrode-accelerated one. This key-value store supports serializable transactions using two-phase commit and optimistic concurrency control (OCC). Clients use `BEGIN_TXN`, `COMMIT_TXN`, `ABORT_TXN`, `SET`, and `GET` operations to express transactions. We use a synthetic workload derived from the Retwis application [56]—an open-source Twitter clone. This workload consists of four types of transactions with different ratios, and each one issues different numbers of `GET` and `PUT` operations. The workload details can be found in Table 2 of [80]. We vary the number of clients that execute transactions in a closed-loop manner, and measure the maximum throughput these clients can achieve and the average latency under one client.

Figure 7a and 7b shows the maximum throughput and average latency of the key-value store atop the Multi-Paxos protocol vs. Electrode-accelerated one under different numbers of replicas, respectively. Overall, Electrode improves the key-value store throughput by 32.3%-112.9% and latency by 5.9%-19.3%. The improvement becomes larger with more replicas, due to the similar reasons described in §7.1. The latency of the key-value store atop the original Multi-Paxos gradually increases with more replicas, while Electrode-accelerated one’s remains relatively stable, because the former is more vulnerable to follower nodes straggling (§7.1).

7.4 CPU Utilization

One design goal of Electrode is to reduce the kernel networking stack overhead (§3) when implementing Paxos protocols. Thus, in this subsection, we study the impact of Electrode on CPU utilizations, which indicates how much kernel stack overhead gets reduced.

Figure 8a and 8b show the CPU utilization of the leader node and follower nodes, respectively, for the Multi-Paxos protocol and Electrode-accelerated one with different offered throughput. The experiments are done in an open-loop manner to control the offered throughput when measuring CPU utilization. The CPU utilization covers both the core handling

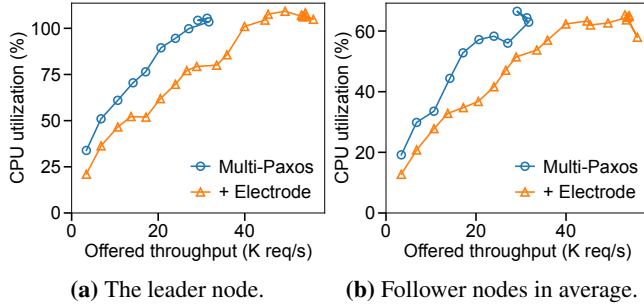


Figure 8: CPU utilization comparison of the Multi-Paxos protocol vs. Electrode-accelerated one (with 5 replicas).

interrupts and the core running Paxos. With higher offered throughput, the CPU utilization gradually increases, demonstrating the load-aware CPU scaling provided by the kernel networking stack (§1). We note that for DPDK-based Multi-Paxos protocol implementation, the CPU utilization would be always 100% because DPDK busily polls the network interface. Overall, Electrode reduces the CPU utilization by 22.7%-38.0% on the leader node and 16.0%-35.7% on the follower nodes, benefiting from the reduced user-kernel crossings and kernel stack traversing.

7.5 Comparison with Kernel-Bypassing

Electrode still handles client-facing requests/responses and initiates message broadcasting using the Linux kernel networking stack (§3); thus, it will achieve lower performance than pure kernel-bypassing approaches. This subsection compares the performance of Electrode with a kernel-bypassing baseline, aiming to reveal the performance upper bound of kernel-based approaches and identify the possible improvements for future work.

We choose Caladan [15] and use its high-performance DPDK-based UDP stack to implement our kernel-bypassing baseline. Similar to Caladan, our baseline dedicates one CPU core for packet polling and another core for running the Paxos protocol. We also configure the Caladan runtime to never idle the Paxos core even under low request load.

Table 3 compares the latency and throughput of kernel-based Multi-Paxos and the kernel-bypassing one. To exclude the latency incurred by the client-side kernel stack, we tested all three Paxos implementations with a request generator implemented using Caladan. Electrode achieves 1.4-1.6x lower latency and 2.0x higher throughput than vanilla Linux, but it still has 2.2x higher latency and 2.4x lower throughput compared to pure kernel-bypassing. The performance gap between Electrode and kernel-bypassing exists, because there are still substantial Paxos messages going through the kernel networking stack in Electrode. In particular, our profiling shows that, on the leader node, around 59.5% CPU time is spent on `__libc_sendto()` caused by frequent `dev_queue_xmit()` and `sk_buff` clones. Although eBPF-based broadcasting reduces a significant number of user-kernel crossings and sock-

	Lowest median/99p latency	Maximum throughput
Vanilla Linux	59/69 μ s	32 K req/s
Electrode	38/49 μ s	65 K req/s
Kernel-bypassing	17/22 μ s	154 K req/s

Table 3: Performance comparison of kernel-based Multi-Paxos vs. kernel-bypassing one (with 5 replicas).

/UDP/IP layer traversing, it cannot fundamentally optimize how the Linux kernel manages NICs and packet buffers. Finally, we note that Electrode’s goal is to provide generic eBPF-based accelerations for distributed protocol implementations that stick to kernel networking stacks because of compatibility, security, isolation, and elastic CPU scaling.

An additional evaluation regarding how the interrupt coalescing feature of modern NICs impacts Electrode is in Appendix B.

8 Discussion and Future Work

Electrode’s offloading decisions: Electrode decides to leave four components of the Multi-Paxos protocol to the user space: 1) failure recovery, 2) handling packet loss and reordering, 3) handling client-facing requests/responses, and 4) executing application-specific operations after reaching the consensus. The first two components involve complex operations on the log, e.g., scanning the log and sending inconsistent entries to other replicas, and inserting missing log entries received from others. These operations require accessing dynamic ranges of log entries, which would fail the eBPF static verification. The last two involve complex serialization/deserialization and application-level operations (see §3). We note that it is possible to offload these four components into eBPF by modifying the kernel eBPF subsystem or verifier—we leave this as future work.

How to improve the eBPF subsystem for offloading? Verifying memory accesses more smartly could make more application operations offloadable. The current eBPF verifier only allows accessing static ranges of memory, which hinders many applications with complex memory accessing behaviors. Another useful construct in eBPF would be dynamic memory allocation, which could ease the maintenance of more advanced data structures in eBPF. To avoid memory leaks, a possible solution could be enforcing Rust-style single-owner memory semantics.

io_uring [1] was recently introduced into the Linux kernel to support efficient batching of asynchronous I/Os via shared memory between the user and kernel space, thus reducing the overhead of frequent user-kernel crossings. Therefore, when implementing Paxos protocols using `io_uring`, it can help reduce the overhead of message broadcasting, which accounts for 12.5% of CPU time based on Table 1. However, each preparation and ACK message still goes through the

full Linux networking stack and wakes up user-space applications, incurring significant overhead; Electrode can be used together with `io_uring` to reduce such overhead. A recent work XRP [82] shares a similar view regarding `io_uring`.

Electrode on shared environments: Electrode requires attaching eBPF programs to the network interface, which then processes every packet accordingly. However, multiple Electrode applications might share the same NIC and attach different eBPF programs that might interfere with each other. We can use the SR-IOV (Single Root IO Virtualization) feature that is widely available in modern NICs [2, 9] to avoid such interference. SR-IOV virtualizes a physical network interface into multiple virtualized ones; the Electrode eBPF program can be attached to only one virtualized interface, without impacting others (e.g., used by non-Paxos applications). Besides SR-IOV, Electrode can also check the port numbers of incoming packets in eBPF, and only execute optimizations if the port numbers belong to target Paxos applications.

Accelerating leader-less consensus protocols using eBPF: Electrode targets at leader-based consensus protocols such as Paxos [37] and its variants [36, 43, 54], because they are the most-used ones by modern distributed applications [6, 8, 22]. Electrode's eBPF-based optimizations could also be applied to leader-less consensus protocols, e.g., EPaxos [52], Mencius [4], SD-Paxos [81], etc. For example, replicas in EPaxos could acknowledge preparation messages earlier in an eBPF program before entering the kernel networking stack, thus reducing latency. We leave the exploration of applying Electrode to leader-less consensus protocols as future work.

9 Related Work

Kernel-bypass and hardware offloading: Overheads of the monolithic kernel networking stack have spurred various attempts to design new kernel-bypassed networking stacks like mTCP [24], eRPC [27], Demikernel [79] and more [15, 29, 33, 48, 57, 67], which attempt to eliminate the kernel from the I/O datapath. But all of these solutions are not backward compatible with solutions that already use the standard kernel networking stack, and they incur more costs in terms of CPU cycles and energy during low I/O loads due to busy-polling. Electrode attempts to leverage eBPF to unclog some of the bottlenecks in the kernel networking stack for distributed protocols without completely having to shift to kernel-bypassed stacks.

Similarly, network offload solutions attempt to offload I/O-intensive operations to specialized hardware, e.g., RDMA [11, 28, 76], FPGA [23], SmartNICs [66], and programmable switches [10, 25]. But they come with limited interfaces for programmability and need custom hardware to be installed.

Co-designing distributed systems with networks: There have been attempts to optimize distributed systems by co-designing them with data center networks for improved performance. SpecPaxos [61] attempts to leverage the natural order of packet delivery in data centers to optimize the ordering of

messages needed for state machine replication. NoPaxos [40] uses in-network switches to sequence packets for a similar purpose. Eris [39] further applies in-network sequencing to distributed transactions to avoid coordination overhead. These are orthogonal ways to optimize distributed systems and can be used in conjunction with Electrode.

Distributed protocols in data centers: Data centers have a variety of distributed protocols that are deployed for fault tolerance and data consistency. These include replication protocols like Mencius [4], EPaxos [52], chain replication [74], SDPaxos [81], and transaction protocols like TAPIR [80] and Meerkat [72]. Since many distributed protocols share similar patterns of communication like broadcasting and quorum responses, Electrode can be applied to speed up these distributed protocols as well.

eBPF applications: For a long time, eBPF was only used for packet filtering [49], monitoring [3, 63], and load balancing [14] because of its restricted programming model. Now, it is shown to be able to offload small yet critical operations to improve application performance. CCP [53] mentions that it may be possible to leverage the JIT feature of eBPF to gather datapath's congestion measurements for congestion control. BMC [17] uses eBPF to implement an in-kernel cache to accelerate UDP-based Memcached GET requests and achieves significant throughput improvement. Syrup [26] uses eBPF maps to share incoming request information across OS, networking stacks, and application runtimes to enable user-defined scheduling. SPRIGHT [65] employs fast eBPF-based packet forwarding to accelerate sidecar proxies in serverless computing. XRP [82] offloads storage functions (e.g., B-tree lookups) into the kernel using eBPF to reduce kernel storage stack overhead. SynCord [58] leverages eBPF to inject workload-specific and hardware-aware kernel lock policies specified by application developers. Electrode further demonstrates that eBPF can be used to accelerate distributed protocols under the kernel networking stack.

10 Conclusion

Electrode is a system that accelerates distributed protocols using safe in-kernel eBPF-based packet processing before the networking stack. Electrode retains the benefits of using the standard Linux networking stack (e.g., good maintenance, elastic CPU scaling, security, and isolation), while optimizing the performance-critical operations of distributed protocols (e.g., broadcasting, and wait-on-quorums) in a non-intrusive manner. When applying Electrode to a classic Multi-Paxos protocol, we achieve up to 128.4% higher throughput and 41.7% lower latency. We believe that the designs of eBPF-based optimizations in Electrode can motivate more research on improving networked application performance while maintaining the standard Linux networking stack.

Electrode code is available at <https://github.com/Electrode-NSDI23/Electrode>.

Acknowledgments

We thank our shepherd Adam Belay and the anonymous reviewers for their insightful comments. We thank Cloudlab [12] for providing us with the development and evaluation infrastructure. This work was supported in part by ACE, one of the seven centers in JUMP 2.0, a Semiconductor Research Corporation (SRC) program sponsored by DARPA. Yang Zhou is also supported by the Google PhD Fellowship.

References

- [1] Efficient IO with io_uring. https://kernel.dk/io_uring.pdf.
- [2] NVIDIA Corporation affiliates. Single Root IO Virtualization (SR-IOV) for Mellanox NICs. <https://docs.nvidia.com/networking/pages/viewpage.action?pageId=43718746>.
- [3] The Cilium Authors. Cilium: eBPF-Based Networking, Observability, Security. <https://cilium.io/>.
- [4] Catalonia-Spain Barcelona. Mencius: Building Efficient Replicated State Machines for WANs. In *Proceedings of USENIX OSDI*, 2008.
- [5] Adam Belay, George Prekas, Ana Klimovic, Samuel Grossman, Christos Kozyrakis, and Edouard Bugnion. IX: A Protected Dataplane Operating System for High Throughput and Low Latency. In *Proceedings of USENIX OSDI*, pages 49–65, 2014.
- [6] Mike Burrows. The Chubby Lock Service for Loosely-Coupled Distributed Systems. In *Proceedings of USENIX OSDI*, pages 335–350, 2006.
- [7] Fay Chang, Jeffrey Dean, Sanjay Ghemawat, Wilson C Hsieh, Deborah A Wallach, Mike Burrows, Tushar Chandra, Andrew Fikes, and Robert E Gruber. Bigtable: A Distributed Storage System for Structured Data. *ACM Transactions on Computer Systems (TOCS)*, 26(2):1–26, 2008.
- [8] James C Corbett, Jeffrey Dean, Michael Epstein, Andrew Fikes, Christopher Frost, Jeffrey John Furman, Sanjay Ghemawat, Andrey Gubarev, Christopher Heiser, Peter Hochschild, et al. Spanner: Google’s Globally Distributed Database. *ACM Transactions on Computer Systems (TOCS)*, 31(3):1–22, 2013.
- [9] Intel Corporation. Single Root IO Virtualization (SR-IOV) for Intel NICs. <https://www.intel.com/content/www/us/en/support/articles/000005722/ethernet-products.html>.
- [10] Huynh Tu Dang, Daniele Sciascia, Marco Canini, Fernando Pedone, and Robert Soulé. Netpaxos: Consensus at Network Speed. In *Proceedings of ACM SIGCOMM Symposium on Software Defined Networking Research (SOSR)*, pages 1–7, 2015.
- [11] Aleksandar Dragojević, Dushyanth Narayanan, Edmund B Nightingale, Matthew Renzelmann, Alex Shamis, Anirudh Badam, and Miguel Castro. No Compromises: Distributed Transactions with Consistency, Availability, and Performance. In *Proceedings of ACM SOSP*, pages 54–70, 2015.
- [12] Dmitry Duplyakin, Robert Ricci, Aleksander Maricq, Gary Wong, Jonathon Duerig, Eric Eide, Leigh Stoller, Mike Hibler, David Johnson, Kirk Webb, et al. The Design and Operation of CloudLab. In *Proceedings of USENIX ATC*, pages 1–14, 2019.
- [13] Facebook. Facebook’s Branch of Apache Thrift, Including a New C++ Server. <https://github.com/facebook/fbthrift/blob/main/thrift/doc/cpp/cpp2.md#options>.
- [14] Facebook. Katran: A High-Performance Layer 4 Load Balancer. <https://github.com/facebookincubator/katran>.
- [15] Joshua Fried, Zhenyuan Ruan, Amy Ousterhout, and Adam Belay. Caladan: Mitigating Interference at Microsecond Timescales. In *Proceedings of USENIX OSDI*, pages 281–297, 2020.
- [16] Sanjay Ghemawat, Howard Gobioff, and Shun-Tak Leung. The Google File System. In *Proceedings of ACM SOSP*, pages 29–43, 2003.
- [17] Yoann Ghigoff, Julien Sopena, Kahina Lazri, Antoine Blin, and Gilles Muller. BMC: Accelerating Memcached using Safe In-kernel Caching and Pre-stack Processing. In *Proceedings of USENIX NSDI*, pages 487–501, 2021.
- [18] Dan Gibson, Hema Hariharan, Eric Lance, Moray McLaren, Behnam Montazeri, Arjun Singh, Stephen Wang, Hassan MG Wassel, Zhehua Wu, Sunghwan Yoo, et al. Aquila: A unified, low-latency fabric for datacenter networks. In *Proceedings of USENIX NSDI*, pages 1249–1266, 2022.
- [19] Google. Protocol Buffers. <https://developers.google.com/protocol-buffers/>.
- [20] The Tcpdump Group. tcpdump. <https://www.tcpdump.org/>.
- [21] Toke Høiland-Jørgensen, Jesper Dangaard Brouer, Daniel Borkmann, John Fastabend, Tom Herbert, David

- Ahern, and David Miller. The eXpress Data Path: Fast Programmable Packet Processing in the Operating System Kernel. In *Proceedings of ACM CoNEXT*, pages 54–66, 2018.
- [22] Michael Isard. Autopilot: Automatic Data Center Management. *ACM SIGOPS Operating Systems Review*, 41(2):60–67, 2007.
- [23] Zsolt István, David Sidler, Gustavo Alonso, and Marko Vukolic. Consensus in a Box: Inexpensive Coordination in Hardware. In *Proceedings of USENIX NSDI*, pages 425–438, 2016.
- [24] EunYoung Jeong, Shinae Wood, Muhammad Jamshed, Haewon Jeong, Sunghwan Ihm, Dongsu Han, and KyoungSoo Park. mTCP: a Highly Scalable User-level TCP Stack for Multicore Systems. In *Proceedings of USENIX NSDI*, pages 489–502, 2014.
- [25] Xin Jin, Xiaozhou Li, Haoyu Zhang, Nate Foster, Jeongkeun Lee, Robert Soulé, Changhoon Kim, and Ion Stoica. NetChain: Scale-Free Sub-RTT Coordination. In *Proceedings of USENIX NSDI*, pages 35–49, 2018.
- [26] Kostis Kaffes, Jack Tigar Humphries, David Mazières, and Christos Kozyrakis. Syrup: User-Defined Scheduling Across the Stack. In *Proceedings of ACM SOSP*, pages 605–620, 2021.
- [27] Anuj Kalia, Michael Kaminsky, and David Andersen. Datacenter RPCs can be General and Fast. In *Proceedings of USENIX NSDI*, pages 1–16, 2019.
- [28] Anuj Kalia, Michael Kaminsky, and David G Andersen. FaSST: Fast, Scalable and Simple Distributed Transactions with Two-Sided (RDMA) Datagram RPCs. In *Proceedings of USENIX OSDI*, pages 185–201, 2016.
- [29] Antoine Kaufmann, Tim Stamler, Simon Peter, Naveen Kr Sharma, Arvind Krishnamurthy, and Thomas Anderson. TAS: TCP Acceleration as an OS Service. In *Proceedings of EuroSys*, pages 1–16, 2019.
- [30] The Linux kernel development community. BPF Ring Buffer. <https://www.kernel.org/doc/html/latest/bpf/ringbuf.html>.
- [31] The Linux kernel development community. struct sk_buff. <https://docs.kernel.org/networking/skbuff.html>.
- [32] Marios Kogias and Edouard Bugnion. Hovercraft: Achieving Scalability and Fault-tolerance for microsecond-scale Datacenter Services. In *Proceedings of EuroSys*, pages 1–17, 2020.
- [33] Marios Kogias, George Prekas, Adrien Ghosn, Jonas Fietz, and Edouard Bugnion. R2P2: Making RPCs First-Class Datacenter Citizens. In *Proceedings of USENIX ATC*, pages 863–880, 2019.
- [34] Hsiang-Tsung Kung and John T Robinson. On Optimistic Methods for Concurrency Control. *ACM Transactions on Database Systems (TODS)*, 6(2):213–226, 1981.
- [35] UW Systems Lab. Speculative Paxos Open Source. <https://github.com/UWSysLab/specpaxos>.
- [36] Leslie Lamport. Paxos Made Simple. *ACM SIGACT News (Distributed Computing Column)* 32, 4 (Whole Number 121, December 2001), pages 51–58, 2001.
- [37] Leslie Lamport. The Part-Time Parliament. In *Concurrency: the Works of Leslie Lamport*, pages 277–317, 2019.
- [38] Leslie Lamport, Dahlia Malkhi, and Lidong Zhou. Vertical Paxos and Primary-Backup Replication. In *Proceedings of ACM PODC*, pages 312–313, 2009.
- [39] Jialin Li, Ellis Michael, and Dan RK Ports. Eris: Coordination-Free Consistent Transactions Using In-Network Concurrency Control. In *Proceedings of ACM SOSP*, pages 104–120, 2017.
- [40] Jialin Li, Ellis Michael, Naveen Kr Sharma, Adriana Szekeres, and Dan RK Ports. Just Say NO to Paxos Overhead: Replacing Consensus with Network Ordering. In *Proceedings of USENIX OSDI*, pages 467–483, 2016.
- [41] Jialin Li, Naveen Kr Sharma, Dan RK Ports, and Steven D Gribble. Tales of the Tail: Hardware, OS, and Application-level Sources of Tail Latency. In *Proceedings of ACM SoCC*, pages 1–14, 2014.
- [42] John C Lin and Sanjoy Paul. RMTP: A Reliable Multicast Transport Protocol. In *Proceedings of IEEE INFOCOM*, volume 96. Citeseer, 1996.
- [43] Barbara Liskov and James Cowling. Viewstamped Replication Revisited. 2012.
- [44] Xuhao Luo, Weihai Shen, Shuai Mu, and Tianyin Xu. DepFast: Orchestrating Code of Quorum Systems. In *Proceedings of USENIX ATC*, pages 557–574, 2022.
- [45] Linux Programmer’s Manual. bpf-helpers(7). <https://man7.org/linux/man-pages/man7/bpf-helpers.7.html>.
- [46] Linux Programmer’s Manual. bpf(2). <https://man7.org/linux/man-pages/man2/bpf.2.html>.

- [47] Linux Programmer's Manual. tc-bpf(8). <https://man7.org/linux/man-pages/man8/tc-bpf.8.html>.
- [48] Michael Marty, Marc de Kruijf, Jacob Adriaens, Christopher Alfeld, Sean Bauer, Carlo Contavalli, Michael Dalton, Nandita Dukkipati, William C Evans, Steve Gribble, et al. Snap: A Microkernel Approach to Host Networking. In *Proceedings of ACM SOSP*, pages 399–413, 2019.
- [49] Steven McCanne and Van Jacobson. The BSD Packet Filter: A New Architecture for User-level Packet Capture. In *USENIX winter*, volume 46, 1993.
- [50] Paul E McKenney. Overview of Linux-Kernel Reference Counting. *N2167*, pages 07–0027, 2007.
- [51] Henrique Moniz, Nuno Ferreira Neves, and Miguel Correia. Turquois: Byzantine Consensus in Wireless Ad hoc Networks. In *2010 IEEE/IFIP International Conference on Dependable Systems & Networks (DSN)*, pages 537–546. IEEE, 2010.
- [52] Iulian Moraru, David G Andersen, and Michael Kaminsky. There is More Consensus in Egalitarian Parliaments. In *Proceedings of ACM SOSP*, pages 358–372, 2013.
- [53] Akshay Narayan, Frank Cangialosi, Deepti Raghavan, Prateesh Goyal, Srinivas Narayana, Radhika Mittal, Mohammad Alizadeh, and Hari Balakrishnan. Restructuring Endpoint Congestion Control. In *Proceedings of ACM SIGCOMM*, pages 30–43, 2018.
- [54] Brian M Oki and Barbara H Liskov. Viewstamped Replication: A New Primary Copy Method to Support Highly-Available Distributed Systems. In *Proceedings of ACM PODC*, pages 8–17, 1988.
- [55] Michael A Olson, Keith Bostic, and Margo I Seltzer. Berkeley DB. In *Proceedings of USENIX ATC, FREENIX Track*, pages 183–191, 1999.
- [56] VMware Inc. or its affiliates. Spring Data Redis - Retwis-J. <https://docs.spring.io/spring-data/data-keyvalue/examples/retwisj/current/>.
- [57] Amy Ousterhout, Joshua Fried, Jonathan Behrens, Adam Belay, and Hari Balakrishnan. Shenango: Achieving High CPU Efficiency for Latency-Sensitive Datacenter Workloads. In *Proceedings of USENIX NSDI*, pages 361–378, 2019.
- [58] Sujin Park, Diyu Zhou, Yuchen Qian, Irina Calciu, Taesoo Kim, and Sanidhya Kashyap. Application-Informed Kernel Synchronization Primitives. In *Proceedings of USENIX OSDI*, pages 667–682, 2022.
- [59] Simon Peter, Jialin Li, Irene Zhang, Dan RK Ports, Doug Woos, Arvind Krishnamurthy, Thomas Anderson, and Timothy Roscoe. Arrakis: The Operating System is the Control Plane. *ACM Transactions on Computer Systems (TOCS)*, 33(4):1–30, 2015.
- [60] Valentin Poirot, Beshr Al Nahas, and Olaf Landsiedel. Paxos Made Wireless: Consensus in the Air. In *EWSN*, pages 1–12, 2019.
- [61] Dan RK Ports, Jialin Li, Vincent Liu, Naveen Kr Sharma, and Arvind Krishnamurthy. Designing Distributed Systems Using Approximate Synchrony in Data Center Networks. In *Proceedings of USENIX NSDI*, pages 43–57, 2015.
- [62] Ravi Prasad, Manish Jain, and Constantinos Dovrolis. Effects of Interrupt Coalescence on Network Measurements. In *International Workshop on Passive and Active Network Measurement*, pages 247–256. Springer, 2004.
- [63] The IO Visor Project. BPF Compiler Collection (BCC). <https://github.com/iovisor/bcc>.
- [64] The IO Visor Project. eXpress Data Path (XDP). <https://www.iovisor.org/technology/xdp>.
- [65] Shixiong Qi, Leslie Monis, Ziteng Zeng, Ian-chin Wang, and KK Ramakrishnan. SPRIGHT: Extracting the Server From Serverless Computing! High-Performance eBPF-Based Event-Driven, Shared-Memory Processing. In *Proceedings of ACM SIGCOMM*, pages 780–794, 2022.
- [66] Henry N Schuh, Weihao Liang, Ming Liu, Jacob Nelson, and Arvind Krishnamurthy. Xenic: SmartNIC-Accelerated Distributed Transactions. In *Proceedings of ACM SOSP*, pages 740–755, 2021.
- [67] ScyllaDB. SeaStar High Performance Server-Side Application Framework. <https://github.com/scylladb/seastar>.
- [68] Muhammad Shahbaz, Lalith Suresh, Jennifer Rexford, Nick Feamster, Ori Rottenstreich, and Mukesh Hira. Elmo: Source Routed Multicast for Public Clouds. In *Proceedings of ACM SIGCOMM*, pages 458–471. 2019.
- [69] Gráinne Sheerin. gRPC and Deadlines. <https://grpc.io/blog/deadlines/>.
- [70] Alberto Spina, Julie McCann, Michael Breza, and Anandha Gopalan. *Reliable Distributed Consensus for Low-Power Multi-Hop Networks*. PhD thesis, Master's thesis, Imperial College London, 2019.
- [71] Michael Stonebraker, Samuel Madden, Daniel J. Abadi, Stavros Harizopoulos, Nabil Hachem, and Pat Helland.

- The End of an Architectural Era: (It's Time for a Complete Rewrite). In *Proceedings of VLDB*, page 1150–1160. VLDB Endowment, 2007.
- [72] Adriana Szekeres, Michael Whittaker, Jialin Li, Naveen Kr Sharma, Arvind Krishnamurthy, Dan RK Ports, and Irene Zhang. Meerkat: Multicore-Scalable Replicated Transactions Following the Zero-Coordination Principle. In *Proceedings of EuroSys*, pages 1–14, 2020.
- [73] Amy Tai, Igor Smolyar, Michael Wei, and Dan Tsafir. Optimizing Storage Performance with Calibrated Interrupts. *ACM Transactions on Storage (TOS)*, 18(1):1–32, 2022.
- [74] Robbert Van Renesse and Fred B Schneider. Chain Replication for Supporting High Throughput and Availability. In *Proceedings of USENIX OSDI*, volume 4, 2004.
- [75] Ed. W. Eddy. RFC 9293: Transmission Control Protocol (TCP). <https://datatracker.ietf.org/doc/html/rfc9293>.
- [76] Xingda Wei, Zhiyuan Dong, Rong Chen, and Haibo Chen. Deconstructing RDMA-Enabled Distributed Transactions: Hybrid is Better! In *Proceedings of USENIX OSDI*, pages 233–251, 2018.
- [77] IJsbrand Wijnands, E Rosen, Andrew Dolganow, Tony Przygienda, and Sam Aldrin. RFC 8279: Multicast Using Bit Index Explicit Replication (BIER). <https://www.rfc-editor.org/rfc/rfc8279>.
- [78] Zhuolong Yu, Yiwen Zhang, Vladimir Braverman, Mosharaf Chowdhury, and Xin Jin. NetLock: Fast, Centralized Lock Management Using Programmable Switches. In *Proceedings of ACM SIGCOMM*, pages 126–138, 2020.
- [79] Irene Zhang, Amanda Raybuck, Pratyush Patel, Kirk Olynyk, Jacob Nelson, Omar S Navarro Leija, Ashlie Martinez, Jing Liu, Anna Kornfeld Simpson, Sujay Jayakar, et al. The Demikernel Datapath OS Architecture for Microsecond-Scale Datacenter Systems. In *Proceedings of ACM SOSP*, pages 195–211, 2021.
- [80] Irene Zhang, Naveen Kr Sharma, Adriana Szekeres, Arvind Krishnamurthy, and Dan RK Ports. Building Consistent Transactions with Inconsistent Replication. *ACM Transactions on Computer Systems (TOCS)*, 35(4):1–37, 2018.
- [81] Hanyu Zhao, Quanlu Zhang, Zhi Yang, Ming Wu, and Yafei Dai. SDPaxos: Building Efficient Semi-Decentralized Geo-Replicated State Machines. In *Proceedings of ACM SoCC*, pages 68–81, 2018.
- [82] Yuhong Zhong, Haoyu Li, Yu Jian Wu, Ioannis Zarkadas, Jeffrey Tao, Evan Mesterhazy, Michael Makris, Junfeng Yang, Amy Tai, Ryan Stutsman, et al. XRP: In-Kernel Storage Functions with eBPF. In *Proceedings of USENIX OSDI*, pages 375–393, 2022.

Types	Protocols	Applying message broadcasting	Applying fast acknowledging	Applying wait-on-quorums
Replication	Primary-backup	The primary broadcasts requests to backups.	Each backup buffers messages in the kernel and quickly responds to the primary.	The primary waits for responses from all backups.
	Chain	None	Each replica (except for the last one) buffers write requests in the kernel and forwards them to the next replica.	None
Concurrency control	Two-phase locking	A transaction coordinator broadcasts LOCK and UNLOCK requests to all shards.	Each shard maintains a lock table in the kernel and directly handles lock acquiring and releasing.	A transaction coordinator waits for responses from all shards.
	OCC	None	Each shard checks in the kernel if the committing transaction's timestamp conflicts with all other running ones.	None
Atomic commitment	Two-phase commit	A transaction coordinator broadcasts PREPARE and COMMIT requests to all shards.	Each shard buffers PREPARE messages in the kernel and responds to the coordinator, and handles COMMIT requests by polling the buffered messages.	A transaction coordinator waits for responses from all shards

Table 4: Applying Electrode to more distributed protocols.

APPENDIX

A Electrode Generalizability

Table 4 summarizes how the classic replication, concurrency control, and atomic commitment protocols can leverage Electrode optimizations. For example, the primary-back replication, two-phase locking, and two-phase commit protocols follow the pattern of sending requests to multiple nodes and waiting for a quorum number of responses; thus they naturally fit well with the eBPF-based message broadcasting and wait-on-quorums. Together with the above protocols, the chain replication [74] and opportunistic concurrency control (OCC) [34] protocols include some critical-yet-simple operations like storing messages in memory, maintaining a lock table, and checking timestamp conflicts; these operations are also suitable for offloading to the eBPF following the fast acknowledging mechanism.

B Impact of Interrupt Coalescing

During benchmarking, we noticed that the interrupt coalescing [62] (IC) feature of modern NICs has a big impact on the measured performance. In IC, after an incoming packet triggers an interrupt, the kernel networking stack waits until a threshold of packets arrives or a timeout gets triggered, aiming to amortize the interrupt cost. In our scenarios, we find it significantly hurts latency and performance predictability in our settings; similar results are also reported in [73]. Thus, in all our experiments, we disable IC by default.

Figure 9 shows the performance impact of IC on the Multi-Paxos protocol and Electrode-accelerated one, by varying the number of open-loop clients. With IC, load-latency curves become unpredictable with two “hockey stick”s. The second “hockey stick” is because the extremely high load triggers coalescing/batching much more packets in one interrupt. Overall, IC does not nearly impact the maximum throughput for the Multi-Paxos protocol and Electrode-accelerated one, but it increases the latency by 57.4%-129.2% and 9.1%-246.8% with 1-3 clients (before the first “hockey stick”). Moreover,

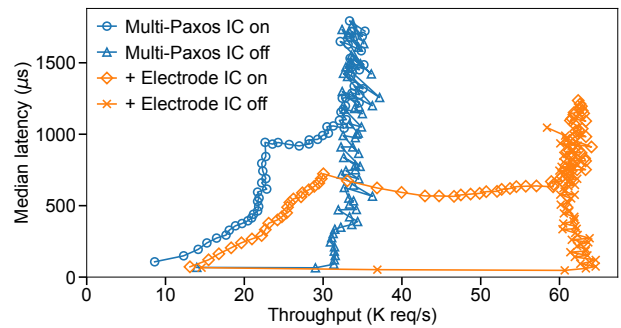


Figure 9: Performance impact of interrupt coalescing (IC) on the Multi-Paxos protocol vs. Electrode-accelerated one (with 5 replicas).

enabling IC decreases the one-client throughput by 38.3% and 10.1% for the original Multi-Paxos and Electrode-accelerated one, respectively.

Electrode performance with IC: Electrode accelerates the maximum throughput of the Multi-Paxos protocol by 81.4% and latency by 32.7% with 1 client (before the first “hockey stick”) when IC is on.

[This page intentionally left blank.]

Carbink: Fault-Tolerant Far Memory

Yang Zhou^{†*} Hassan M.G. Wassel[‡] Sihang Liu^{§*} Jiaqi Gao[†] James Mickens[†] Minlan Yu^{†‡}
Chris Kennelly[‡] Paul Turner[‡] David E. Culler[‡] Henry M. Levy^{||‡} Amin Vahdat[‡]

[†]*Harvard University* [‡]*Google* [§]*University of Virginia* ^{||}*University of Washington*

Abstract

Far memory systems allow an application to transparently access local memory as well as memory belonging to remote machines. Fault tolerance is a critical property of any practical approach for far memory, since machine failures (both planned and unplanned) are endemic in datacenters. However, designing a fault tolerance scheme that is efficient with respect to both computation and storage is difficult. In this paper, we introduce Carbink, a far memory system that uses erasure-coding, remote memory compaction, one-sided RMAs, and offloadable parity calculations to achieve fast, storage-efficient fault tolerance. Compared to Hydra, a state-of-the-art fault-tolerant system for far memory, Carbink has 29% lower tail latency and 48% higher application performance, with at most 35% higher memory usage.

1 Introduction

In a datacenter, matching a particular application to just enough memory and CPUs is hard. A commodity server tightly couples memory and compute, hosting a fixed number of CPUs and RAM modules that are unlikely to exactly match the computational requirements of any particular application. Even if a datacenter contains a heterogeneous mix of server configurations, the load on each server (and thus the amount of available resources for a new application) changes dynamically as old applications exit and new applications arrive. Thus, even state-of-the-art cluster schedulers [51, 52] struggle to efficiently bin-pack a datacenter’s aggregate collection of CPUs and RAM. For example, Google [52] and Alibaba [34] report that the average server has only ~60% memory utilization, with substantial variance across machines.

Memory is a particularly vexing resource for two reasons. First, for several important types of applications [19, 20, 33, 54], the data set is too big to fit into the RAM of a single machine, even if the entire machine is assigned to a single application instance. Second, for these kinds of applications, alleviating memory pressure by swapping data between RAM and storage [14] would lead to significant application slowdowns, because even SSD accesses are orders of magnitude slower than RAM accesses. For example, Google runs a graph

analysis engine [28] whose data set is dozens of GBs in size. This workload runs 46% faster when it shuffles data purely through RAM instead of between RAM and SSDs.

Disaggregated datacenter memory [2, 5, 15, 16, 22, 44, 46] is a promising solution. In this approach, a CPU can be paired with an arbitrary set of possibly-remote RAM modules, with a fast network interconnect keeping access latencies to far memory small. From a developer’s perspective, far memory can be exposed to applications in several ways. For example, an OS can treat far RAM as a swap device, transparently exchanging pages between local RAM and far RAM [5, 22, 46]. Alternatively, an application-level runtime like AIFM [44] can expose remotable pointer abstractions to developers, such that pointer dereferences (or the runtime’s detection of high memory pressure) trigger swaps into and out of far memory.

Much of the prior work on disaggregated memory [2, 44, 55] has a common limitation: a lack of fault tolerance. Unfortunately, in a datacenter containing hundreds of thousands of machines, faults are pervasive. Many of these faults are planned, like the distribution of kernel upgrades that require server reboots, or the intentional termination of a low-priority task when a higher-priority task arrives. However, many server faults are unpredictable, like those caused by hardware failures or kernel panics. Thus, any *practical* system for far memory has to provide a scalable, fast mechanism to recover from unexpected server failures. Otherwise, the failure rate of an application using far memory will be much higher than the failure rate of an application that only uses local memory; the reason is that the use of far memory increases the set of machines whose failure can impact an application [8].

Some prior far-memory systems do provide fault tolerance via replication [5, 22, 46]. However, replication-based approaches suffer from high storage overheads. Hydra [29] uses erasure coding, which has smaller storage penalties than replication. However, Hydra’s coding scheme stripes a single memory page across multiple remote nodes. This means that a compute node requires multiple network fetches to reconstruct a page; furthermore, computation over that page cannot be outsourced to remote memory nodes, since each node contains only a subset of the page’s bytes.

*Contributed to this work during internships at Google.

In this paper, we present Carbink,¹ a new framework for far memory that provides efficient, high-performance fault recovery. Like (non-fault-tolerant) AIFM, Carbink exposes far memory to developers via application-level remoteable pointers. When Carbink’s runtime must evict data from local RAM, Carbink writes erasure-coded versions of that data to remote memory nodes. The advantage of erasure coding is that it provides equivalent redundancy to pure replication, while avoiding the double or triple storage overheads that replication incurs. However, straightforward erasure coding is a poor fit for the memory data created by applications written in standard programming languages like C++ and Go; those applications allocate variable-sized memory objects, but erasure coding requires equal-sized blocks. To solve this problem, Carbink eschews the object-granularity swapping strategy of AIFM, and instead swaps at the granularity of *spans*. A single span consists of multiple memory pages that contain objects with similar sizes. Carbink’s runtime asynchronously and transparently moves local objects within the spans in local memory, grouping cold objects together and hot objects together. When necessary, Carbink batch-evicts cold spans, calculating parity bits for those spans at eviction time, and writing the associated fragments to remote memory nodes. Carbink utilizes one-sided remote memory accesses (RMAs) to efficiently perform swapping activity, minimizing network utilization. Unlike Hydra, Carbink’s erasure coding scheme allows a compute node to fetch a far memory region using a single network request.

In Carbink, each span lives in exactly one place: the local RAM of a compute node, or the far RAM of a memory node. Thus, swapping a span from far RAM to local RAM creates dead space (and thus fragmentation) in far RAM. Carbink runs pauseless defragmentation threads in the background, asynchronously reclaiming space to use for later swap-outs.

We have implemented Carbink atop our datacenter infrastructure. Compared to Hydra, Carbink has up to 29% lower tail latency and 48% higher application performance, with at most 35% more remote memory usage. Unlike Hydra, Carbink also allows computation to be offloaded to remote memory nodes.

In summary, this paper has four contributions:

- a span-based approach for solving the size mismatch between the granularity of erasure coding and the size of the objects allocated by compute nodes;
- new algorithms for defragmenting the RAM belonging to remote memory nodes that store erasure-encoded spans;
- an application runtime that hides spans, object migration within spans, and erasure coding from application-level developers; and
- a thorough evaluation of the performance trade-offs made by different approaches for adding fault tolerance to far memory systems.

¹Carbink is a Pokémon that has a high defense score.

2 Background

Recent work on far memory has used one of two approaches. The first approach modifies the OS that runs applications, exploiting the fact that preexisting OS abstractions already decouple application-visible in-memory data from the backing storage hierarchy. For example, INFINISWAP [22], Fastswap [5], and LegoOS [46] leverage virtual memory support to swap application memory to far RAM instead of a local SSD or hard disk. Applications use standard language-level pointers to interact with memory objects; behind the scenes, the OS swaps pages between local RAM and far RAM, e.g., in response to page faults for non-locally-resident pages. In contrast, the remote region approach [2] exposes far memory via file system abstractions. Applications name remote memory regions using standard filenames, and interact with regions using standard file operations like `open()` and `read()`.

Exposing far memory via OS abstractions is attractive because it requires minimal changes to application-level code. However, invasive kernel changes are needed; such changes require substantial implementation effort, and are difficult to maintain as other parts of the kernel evolve.

The second far-memory approach requires more help from application-level code. For example, AIFM [44] uses a modified C++ runtime to hide the details of managing far memory. The runtime provides special pointer types whose dereferencing may trigger the swapping of a remote C++-level object into local RAM. AIFM’s runtime tracks object hotness using GC-style read/write barriers, and uses background threads to swap out cold local objects when local memory pressure is high. To synchronize the local memory accesses generated by application threads and runtime threads, AIFM embeds a variety of metadata bits (e.g., `present`, `isBeingEvicted`) in each smart pointer, leveraging an RCU-like scheme [36] to protect concurrent accesses to a pointer’s referenced object.

Listing 1 provides an example of how applications use AIFM’s smart pointers. Like AIFM, Carbink exposes far memory via smart pointers, but unlike AIFM, Carbink provides fault tolerance.

3 Carbink Design

Figure 1 depicts the high-level architecture of Carbink. **Compute nodes** execute single-process (but potentially multi-threaded) applications that want to use far memory. **Memory nodes** provide far memory that compute nodes use to store application data that cannot fit in local RAM. A logically-centralized **memory manager** tracks the liveness of compute nodes and memory nodes. The manager also coordinates the assignment of far memory **regions** to compute nodes. When a memory node wants to make a local memory region available to compute nodes, the memory node *registers* the region with the memory manager. Later, when a compute node requires far memory, the compute node sends an *allocation* request to the memory manager, who then assigns a registered, unallo-

```

RemUniquePtr<Node> rem_ptr = AIFM::MakeUnique<Node>();
{
    DerefScope scope;
    Node* normal_ptr = rem_ptr.Deref(scope);
    computeOverNodeObject(normal_ptr);
} // Scope is destroyed; Node object can be evicted.

```

Listing 1: Example of how AIFM applications interact with far memory. In the code above, the application first allocates a Node object that is managed by a particular RemUniquePtr. Such a remote unique pointer represents a pointer to an object that (1) can be swapped between local and far memory, and (2) can only be pointed to by a single application-level pointer. The code then creates a new scope via an open brace, declares a DerefScope variable, and invokes the RemUniquePtr’s Deref() method, passing the DerefScope variable as an argument. Deref() essentially grabs an RCU lock on the remotable memory object, and returns a normal C++ pointer to the application. After the application has finished using the normal pointer, the scope terminates and the destructor of the DerefScope runs, releasing the RCU lock and allowing the object to be evicted from local memory.

cated region. Upon receiving a *deallocation* message from a compute node, the memory manager marks the associated region as available for use by other compute nodes. A memory node can ask the memory manager to *deregister* a previously registered (but currently unallocated) region, withdrawing the region from the global pool of far memory.

Carbink does not require participating machines to use custom hardware. For example, any machine in a datacenter can be a memory node if that machine runs the Carbink memory host daemon. Similarly, any machine can be a compute node if that node’s applications use the Carbink runtime.

From the perspective of an application developer, the Carbink runtime allows a program to dynamically allocate and deallocate memory objects of arbitrary size. As described in Section 3.2, programs access those objects through AIFM-like remotable pointers [44]. When applications dereference pointers that refer to non-local (i.e., swapped-out) objects, Carbink pulls the desired objects from far memory. Under the hood, Carbink’s runtime manages objects using **spans** (§3.3) and **spansets** (§3.4). A span is a contiguous run of memory pages; a single region allocated by a compute node contains one or more spans. Similar to slab allocators like Facebook’s jemalloc [17] and Google’s TCMalloc [21, 24], Carbink rounds up each object allocation to the bin size of the relevant span, and aligns each span to the page size used by compute nodes and memory nodes. Carbink swaps far memory into local memory at the granularity of a span; however, when local memory pressure is high, Carbink swaps local memory out to far memory at the granularity of a spanset (i.e., a collection of spans of the same size). In preparation for

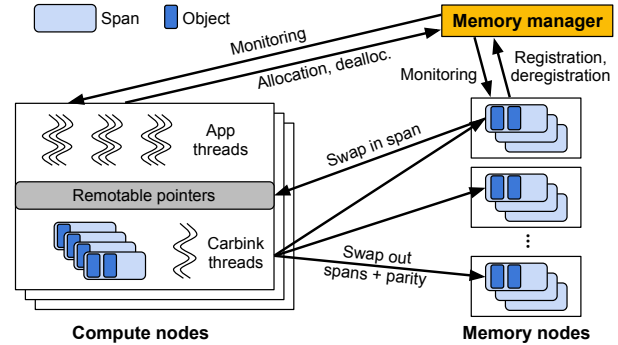


Figure 1: Carbink’s high-level architecture.

swap-outs, background threads on compute nodes group cold objects into cold spans, and bundle a group of cold spans into a spanset; at eviction time, the threads generate erasure-coding parity data for the spanset, and then evict the spanset and the parity data to remote nodes. As we discuss in Sections 3.4 and 3.5, this approach simplifies memory management and fault tolerance.

Carbink disallows cross-application memory sharing. This approach is a natural fit for our target applications, and has the advantage of simplifying failure recovery and avoiding the need for expensive coherence traffic [46].

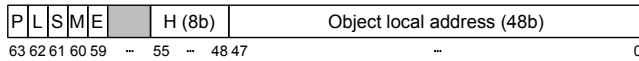
3.1 Failure Model

Carbink implements the logically-centralized memory manager as a replicated state machine [1, 45]. Thus, Carbink assumes that the memory manager will not fail. Carbink assumes that memory nodes and compute nodes may experience fail-stop faults. Carbink does not handle Byzantine failures or partial network failures.

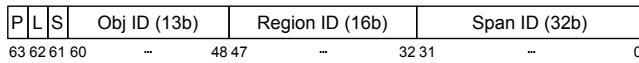
The memory manager tracks the liveness of compute nodes and memory nodes via heartbeats. When a compute node fails, the memory manager instructs the memory nodes to deallocate the relevant spans; if applications desire, they can use an application-level fault tolerance scheme like checkpointing to ensure that application-level data is recoverable. When a memory node fails, the memory manager deregisters the node’s regions from the global pool of far memory. However, erasure-coding recovery of the node’s regions is initiated by a compute node when the compute node unsuccessfully tries to read or write a span belonging to the failed memory node. If an application thread on a compute node tries to read a span that is currently being recovered, the read will use Carbink’s degraded read protocol (§3.5), reconstructing the span using data from other spans and parity blocks.

3.2 Remotable Pointers

Like AIFM, Carbink exposes far memory through C++-level smart pointers. However, as shown in Figure 2, Carbink uses a different pointer encoding to represent span information.



(a) Local object.



(b) Far object.

Field	Meaning
Present	Is the object in local RAM or far RAM?
Lock	Is the object (spin)locked by a thread?
Shared	Is the pointer a unique pointer or a shared pointer?
Moving	Is the object being moved by a background thread?
Evicting	Is the object being evicted by a background thread?
Hotness	Is the object frequently accessed?

(c) Field semantics.

Figure 2: Carbink’s `RemUniquePtr` representation. In contrast to AIFM [44], Carbink does not embed information about a data structure ID or an object size. Instead, Carbink embeds span metadata (namely, a Region ID and a Span ID) to associate a pointed-to object with its backing span.

A Carbink `RemUniquePtr` has the same size as a traditional `std::unique_ptr` (i.e., 8 bytes). The **Present** bit indicates whether the pointed-to object resides in local RAM. The **Shared** bit indicates whether a pointer implements unique-pointer semantics or shared-pointer semantics; the former only allows a single reference to the pointed-to object. The **Lock**, **Moving**, and **Evicting** bits are used to synchronize object accesses between application threads and Carbink’s background threads (§3.6). The **Hotness** byte is consulted by the background threads when deciding whether an object is cold (and thus a priority for eviction).

If an object is local, the local virtual address of the object is directly embedded in the pointer. If an object has been evicted, the pointer describes how to locate the object. In particular, the **Obj ID** indicates the location of an object within a particular span; the **Span ID** identifies that span; and the **Region ID** denotes the far memory region that contains the span.

Carbink supports two smart pointer types: `RemUniquePtr`, which only allows one reference to the underlying object, and `RemSharedPtr`, which allows multiple references. When moving or evicting an object, Carbink’s background threads need a way to locate and update the smart pointer(s) which reference the object. To do so, Carbink uses AIFM’s approach of embedding a “reverse pointer” in each object; the reverse pointer points to the object’s single `RemUniquePtr`, or to the first `RemSharedPtr` that references the object. An individual `RemSharedPtr` is 16 bytes large, with the last 8 bytes storing a pointer that references the next `RemSharedPtr` in the list. Thus, Carbink’s runtime can find all of an object’s `RemSharedPtrs` by discovering the first one via the object’s reverse pointer, and then iterating across the linked list.

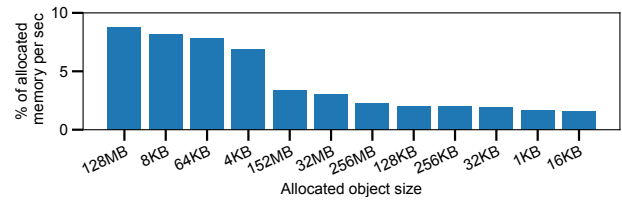


Figure 3: Allocation sizes in our production workloads.

3.3 Span-Based Memory Management

Local memory management: A span is a contiguous set of pages that contain objects of the same size class. Carbink supports 86 different size classes, and aligns each span on an 8KB boundary; Carbink borrows these configuration parameters from TCMalloc [21, 24], which observed these parameters to reduce internal fragmentation. When an application allocates a new object, Carbink tries to round the object size up to the nearest size class and assign a free object slot from an appropriate span. If the object is bigger than the largest size class, Carbink rounds the object size up to the nearest 8KB-aligned size, and allocates a dedicated span to hold it.

To allocate spans locally, Carbink uses a *local page heap*. The page heap is an array of free lists, with each list tracking 8KB-aligned free spans of a particular size (e.g., 2MB, 4MB, etc.). If Carbink cannot find a free span big enough to satisfy an allocation request, Carbink allocates a new span, using `mmap()` to request 2MB huge pages from the OS.

Allocating and deallocating via the page heap is mutex-protected because application threads may issue concurrent allocations or deallocations. To reduce contention on the page heap, each thread reserves a private (i.e., thread-local) cache of free spans for each size class. Carbink also maintains a global cache of free lists, with each list having its own spinlock. When a thread wants to allocate a span whose size can be handled by one of Carbink’s predefined size classes, the thread first tries to allocate from the thread-local cache, then the global cache, and finally the page heap. For larger allocation requests, threads allocate spans directly from the page heap.

Carbink associates each span with several pieces of metadata, including an integer that describes the span’s size class, and a bitvector that indicates which object slots are free. To map a locally-resident object to its associated span metadata, Carbink uses a two-level radix tree called the *local page map*. The lookup procedure is similar to a page table walk: the first 20 bits of an object’s virtual address index into the first-level radix tree table, and the next 15 bits index into a second-level table. The same mapping approach allows Carbink to map the virtual address of a locally-resident span to its metadata.

Far memory management: On a compute node, local spans contain a subset of an application’s memory state. The rest of that state is stored in far spans that live in far memory regions. Recall from Figure 2b that a Carbink pointer to a non-local object embeds the object’s Region ID and Span ID.

To allocate or deallocate a region, a compute node sends a request to the memory manager. A single Carbink region is 1GB or larger, since Carbink targets applications whose total memory requirements are hundreds or thousands of GBs. Upon successfully allocating a region, the compute node updates a *region table* which maps the Region ID of the allocated region to the associated far memory node.

A compute node manages far spans and far regions using additional data structures that are analogous to the ones that manage local spans. A *far page heap* handles the allocation and deallocation of far spans belonging to allocated regions. A *far page map* associates a far Span ID with metadata that (1) names the enclosing region (as a Region ID) and (2) describes the offset of the far span within that region.

Each application thread has a private far cache; Carbink also maintains a global far cache that is visible to all application threads. To swap out a local span of size s , a compute node must first use the far page heap (or a far cache if possible) to allocate a free far span of size s . Similarly, after a compute node swaps in a far span, the node deallocates the far span, returning the far span to its source (either the far page heap or a far cache).

Span filtering and swapping: The Carbink runtime executes *filtering threads* that iterate through the objects in locally-resident spans and move those objects to different local spans. Carbink’s object shuffling has two goals.

- First, Carbink wants to create *hot spans* (containing only hot objects) and *cold spans* (containing only cold ones); when local memory pressure is high, Carbink’s *eviction threads* prefer to swap out spansets containing cold spans. Carbink tracks object hotness using GC-style read/write barriers [4, 23]. Thus, by the time that a filtering thread examines an object, the Hotness byte in the object’s pointer (see Figure 2) has already been set. Upon examining the Hotness byte, a filtering thread updates the byte using the CLOCK algorithm [12].
- Second, object shuffling allows Carbink to garbage-collect dead objects by moving live objects to new spans and then deallocating the old spans. During eviction, Carbink utilizes efficient one-sided RMA writes to swap spansets out to far memory nodes; this approach allows Carbink to avoid software-level overheads (e.g., associated with thread scheduling) on the far node.

From the application’s perspective, object movement and spanset eviction are transparent. This transparency is possible because each object embeds a reverse pointer (§3.2) that allows filtering threads and evicting threads to determine which smart pointers require updating.

Carbink swaps far memory into local memory at the granularity of a span. As with swap-outs, Carbink uses one-sided RMAs for swap-ins. Swapping at the granularity of a span simplifies far memory management, since compute nodes only have to remember how spans map to memory nodes (as opposed to how the much larger number of *objects* map to

memory nodes). However, swapping in at span granularity instead of object granularity has a potential disadvantage: if a compute node swaps in a span containing multiple objects, but only uses a small number of those objects, then the compute node will have wasted network bandwidth (to fetch the unneeded objects) and CPU time (to update the remotable pointers for those unneeded objects). We collectively refer to these penalties as *swap-in amplification*.

To reduce the likelihood of swap-in amplification, Carbink’s filtering and eviction threads prioritize the scanning and eviction of spansets containing large objects. The associated spans contain fewer objects per span; thus, swapping in these spans will reduce the expected number of unneeded objects. Figure 3 shows that, for our production workloads, large objects occupy the majority of memory. Moreover, most hot objects are small; for example, in our company’s geo-distributed database [13], roughly 95% of accesses involve objects smaller than 1.8KB. As a result, an eviction scheme which prioritizes large-object spansets is well-suited for our target applications.

In Carbink, a local span has a three-state lifecycle. A span is first *created* due to a swap-in or local allocation. The span transitions to the *filtering* state upon being examined by filtering threads. Once filtering completes, those spans transition to the *evicting* state when evicting threads begin to swap out spansets. The transition from created to filtering to evicting is fixed, and determines which Carbink runtime threads race with application threads at any given moment (§3.6).

3.4 Fault Tolerance via Erasure Coding

Erasure coding provides data redundancy with lower storage overhead than traditional replication. However, the design space for erasure coding schemes is more complex. Carbink seeks to minimize both average and long-tail access penalties for far objects; per our fault model (§3.1), Carbink also wants to efficiently recover from the failure of memory nodes. Achieving these goals forced us to make careful decisions involving coding granularity, parity recalculation, and cross-node transport protocols.

Coding granularity: To motivate Carbink’s decision to erasure-code at the spanset granularity, first consider an approach that erasure-codes individual spans. In this approach, to swap out a span, a compute node breaks the span into data fragments, generates the associated parity fragments, and then writes the entire set of fragments (data+parity) to remote nodes. During the swap-in of a span, a compute node must fetch multiple fragments to reconstruct the target span.

This scheme, which we call EC-Split, is used by Hydra [29]. With EC-Split, handling the failure of memory nodes during swap-out or swap-in is straightforward: the compute node who is orchestrating the swap-out or swap-in will detect the memory node failure, select a replacement memory node, trigger span reconstruction, and then restart the swap-in or

Schemes	EC data fragment size	Network transport	Parity computation	Defragmentation
EC-Split (Hydra [29])	Span chunk	RMA in & out	Local	N/A
EC-2PC	Full span	RMA in, RPC out (+updating parity via 2PC)	Remote	N/A
EC-Batch Local (Carbink)	Full span	RMA in & out	Local	Remote compaction
EC-Batch Remote (Carbink)	Full span	RMA in & out (+parallel 2PC for compaction)	Local (swap-out)+ Remote (compaction)	Remote compaction

Table 1: The erasure-coding approaches that we study.

swap-out. The disadvantage of EC-Split is that, to reconstruct a single span, a compute node must contact multiple memory nodes to pull in all of the needed fragments. This requirement to contact multiple memory nodes makes the swap-in operation vulnerable to stragglers (and thus high tail latency²). This requirement also frequently prevents a compute node from offloading computation to memory nodes; unless a particular object is small, the object will span multiple fragments, meaning that no single memory node will have a complete local copy of the object.

An alternate approach is to erasure-code across a group of equal-sized spans. We call such a group a *spanset*. In this approach, each span in the spanset is treated as a fragment, with parity data computed across all of the spans in the set. To reconstruct a span, a compute node merely has to contact the single memory node which stores the span. Carbink uses this approach to minimize tail latencies.

Parity updating: Erasure-coding at the spanset granularity but swapping in at the span granularity does introduce complications involving parity updates. The reason is that swapping in a span s leaves an invalid, span-sized hole in the backing spanset; the hole must be marked as invalid because, when s is later swapped out, s will be swapped out as part of a new spanset. The hole created by swapping in s causes fragmentation in the backing spanset. Determining how to garbage-collect the hole and update the relevant parity information is non-trivial. Ideally, a scheme for garbage collection and parity updating would not incur overhead on the critical path of swap-ins or swap-outs. An ideal scheme would also allow parity recalculations to occur at either compute nodes or memory nodes, to enable opportunistic exploitation of free CPU resources on both types of nodes.

Cross-node transport protocols: In systems like RAM-Cloud [39], machines use RPCs to communicate. RPCs involve software-level overheads on both sides of a communication. Carbink avoids these overheads by using one-sided RMA, avoiding unnecessary thread wakeups on the receiver. However, in and of itself, RMA does not automatically solve the consistency issues that arise when offloading parity calculations to remote nodes (§3.4.2).

Throughout the paper, we compare Carbink’s erasure-coding approach to various alternatives.

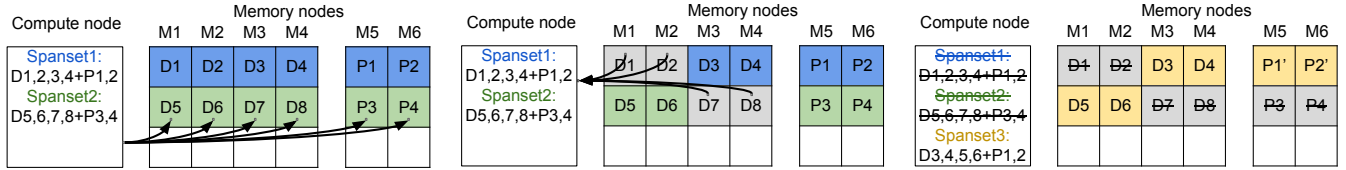
²Hydra [29] and EC-Cache [42] try to minimize straggler-induced latencies by contacting $k + \Delta$ memory nodes instead of the minimum k , using the first k responses to reconstruct an object. This approach increases network traffic and compute-node CPU overheads.

- **EC-Split** is Hydra’s approach, which erasure-codes at the span granularity, swaps data using RMA, and synchronously recalculates parity at compute nodes when swap-outs occur. Fragmentation within an erasure-coding group never occurs, as a span is swapped in and out as a full unit.
- **EC-2PC** erasure-codes using spansets, and uses RMA to swap in at the span granularity. During a swap-out (which happens at the granularity of a span), EC-2PC writes the updated span to the backing memory node; the memory node then calculates the updates to the parity fragments, and sends the updates to the relevant memory nodes which store the parity fragments. To provide crash consistency for the update to the span and the parity fragments, EC-2PC implements a two-phase commit protocol using RPCs. There is no fragmentation within an erasure-coding group because swap-ins and swap-outs both occur at the span granularity.
- **EC-Batch Local** and **EC-Batch Remote** are the approaches used by Carbink. Both schemes erasure-code at spanset granularity, using RMA for swap-in as well as swap-out. Swap-ins occur at the granularity of a span, but swap-outs occur at the granularity of spansets (§3.4.1); thus, both EC-Batch approaches deallocate a span’s backing area in far memory upon swapping that span into a compute node’s local RAM. The result is that swap-ins create dead space on a remote memory node. Both EC-Batch schemes reclaim dead space and recalculate parity data using asynchronous garbage collection. EC-Batch Local always recalculates parity on compute nodes, whereas EC-Batch Remote can recalculate parity on compute nodes or memory nodes. When EC-Batch Remote offloads parity computations to remote nodes, it employs a parallel commit scheme that avoids the latencies of traditional two-phase commit (§3.4.2).

Table 1 summarizes the various schemes. We now discuss EC-Batch Local and Remote in more detail.

3.4.1 EC-Batch: Swapping

Swapping out: In both varieties of EC-Batch, a spanset contains multiple spans of the same size. At swap-out time, a compute node writes a *batch* (i.e., a spanset and its parity fragments) to a memory node. Figure 4a shows an example. In that example, the compute node has two spansets: spanset1 (consisting of data spans $\langle D1, D2, D3, D4 \rangle$ and parity fragments $\langle P1, P2 \rangle$), and spanset2 (containing data spans $\langle D5, D6, D7, D8 \rangle$ and parity fragments $\langle P3, P4 \rangle$). Carbink uses Reed-Solomon codes [43] to create parity data, and prioritizes the eviction of spansets that contain cold spans



(a) Swapping out spans and parity in a batch.

(b) Swapping in individual spans.

(c) Compacting spansets to reclaim space.

Figure 4: EC-Batch swapping-out, swapping-in, and far compaction.

(§3.3). Neither variant of EC-Batch updates spansets in place, so eviction may require a compute node to request additional far memory regions from the memory manager.

Swapping in: When an application tries to access an object that is currently far, the Carbink runtime inspects the application pointer and extracts the Span ID (see Figure 2b). The runtime consults the far page map (§3.3) to discover which remote node holds the span. Finally, the runtime initiates the appropriate RMA operation to swap in the span.

However, swapping in at the span granularity creates *remote fragmentation*. In Figure 4b, the compute node in the running example has pulled four spans ($D1$, $D2$, $D7$, and $D8$) into local memory. Any particular span lives exclusively in local memory or far memory; thus, the swap-ins of the four spans creates dead space on the associated remote memory nodes. If Carbink wants to fill (say) $D1$'s dead space with a new span $D9$, Carbink must update parity fragments $P1$ and $P2$. For a Reed-Solomon code, those parity fragments will depend on both $D1$ and $D9$.

There are two strawman approaches to update $P1$ and $P2$:

- The compute node can read $D1$ into local memory, generate the parity information, and then issue writes to $P1$ and $P2$.
- Alternatively, the compute node can send $D9$ to memory node $M1$, and request that $M1$ compute the new parity data and update $P1$ and $P2$.

The second approach requires a protocol like 2PC to guarantee the consistency of data fragments and parity fragments; without such a protocol, if $M1$ fails after updating $P1$, but before updating $P2$, the parity information will be out-of-sync with the data fragments.

The first approach, in which the compute node orchestrates the parity update, avoids the inconsistency challenges of the second approach. If a memory node dies in the midst of a parity update, the compute node will detect the failure, pick a new memory node to back the parity fragment, and retry the parity update. If the compute node dies in the midst of the parity update, then the memory manager will simply deallocate all regions belonging to the compute node (§3.1).

Unfortunately, both approaches require a lot of network bandwidth to fill holes in far memory. To reclaim one vacant span, the first approach requires three span-sized transfers—the compute node must read $D1$ and then write $P1$ and $P2$. The second approach requires two span-sized transfers to update $P1$ and $P2$. To reduce these network overheads, Carbink performs *remote compaction*, as described in the next section.

3.4.2 EC-Batch: Remote Compaction

Carbink employs *remote compaction* to defragment far memory using fewer network resources than the two strawmen above. On a compute node, Carbink executes several *compaction threads*. These threads look for “matched” spanset pairs; in each pair, the span positions containing dead space in one set are occupied in the other set, and vice versa. For example, the two spansets in Figure 4b are a matched pair. Once the compaction threads find a matched pair, they create a new spanset whose data consists of the live spans in the matched pair (e.g., $\langle D3, D4, D5, D6 \rangle$ in Figure 4b). The compaction threads recompute and update the parity fragments $P1'$ and $P2'$ using techniques that we discuss in the next paragraph. Finally, the compaction threads deallocate the dead spaces in the matched pair (e.g., $\langle D1, D2, D7, D9, P3, P4 \rangle$ in Figure 4b), resulting in a situation like the one shown in Figure 4c. Carbink's compaction can occur in the background, unlike the synchronous parity updates of EC-2PC which place consensus activity on the critical path of swap-outs.

So, how should compaction threads update parity information? Carbink uses Reed-Solomon codes over the Galois field $GF(2^8)$. The new parity data to compute in Figure 4c is therefore represented by the following equations on $GF(2^8)$:

$$P1' - P1 = A_{1,1}(D5 - D1) + A_{2,1}(D6 - D2)$$

$$P2' - P2 = A_{1,2}(D5 - D1) + A_{2,2}(D6 - D2)$$

where $A_{i,j}$ ($i \in \{0, 1, 2, 3\}, j \in \{0, 1\}$) are fixed coefficient vectors in the Reed-Solomon code. Carbink provides two approaches for updating the parity information.

- In EC-Batch Local, the compute node that triggered the swap-out orchestrates the updating of parity data. In the running example, the compute node asks $M1$ to calculate the span delta $D5 - D1$, and asks $M2$ to calculate the span delta $D6 - D2$. After retrieving those updates, the compute node determines the parity deltas (i.e., $P1' - P1$ and $P2' - P2$) and pushes those deltas to the parity nodes $M5$ and $M6$.
- In EC-Batch Remote, the compute node offloads parity recalculation and updating to memory nodes. In the running example, the compute node asks $M1$ to calculate the span delta $D5 - D1$, and $M2$ to calculate the span delta $D6 - D2$. The compute node also asks $M1$ and $M2$ to calculate partial parity updates (e.g., $A_{1,1}(D5 - D1)$ and $A_{1,2}(D5 - D1)$ on $M1$). $M1$ and $M2$ are then responsible for sending the partial parity updates to the parity nodes. For example, $M1$ sends $A_{1,1}(D5 - D1)$ to $M5$, and $A_{1,2}(D5 - D1)$ to $M6$.

In EC-Batch Local, recovery from memory node failure is orchestrated by the compute node in a straightforward way, as in EC-Split (§3.4). In EC-Batch Remote, a compute node performs remote compaction by offloading parity updates to memory nodes. The compute node ensures fault tolerance for an individual compaction via 2PC. However, the compute node aggressively issues compaction requests in parallel. Two compactions (i.e., two instance of the 2PC protocol) are safe to concurrently execute if the compactions involve different spansets; the prepare and commit phases of the two compactions can partially or fully overlap.

On a compute node, Carbink’s runtime can monitor the CPU load and network utilization of remote memory nodes. The runtime can default to remote compaction via EC-Batch Local, but opportunistically switch to EC-Batch Remote if spare resources emerge on memory nodes. During a switch to a different compaction mode, Carbink allows all in-flight compactions to complete before issuing new compactions that use the new compaction mode.

The strawmen defragmentation schemes in Section 3.4.1 require two or three span-sized network transfers to recover one dead span. In the context of Figure 4, EC-Batch Local recovers four dead spans using four span-sized network transfers. EC-Batch Remote requires four span-sized network transfers (plus some small messages generated by the consistency protocol) to recover four dead spans.

3.5 Failure Recovery

Carbink handles two kinds of memory node failures: planned and unplanned. Planned failures are scheduled by the cluster manager [51, 52] to allow for software updates, disk reformatting, and so on. Unplanned failures happen unexpectedly, and are caused by phenomena like kernel panics, defective hardware, and power disruptions.

Planned failures: When the cluster manager decides to schedule a planned failure, the manager sends a warning notification to the affected memory nodes. When a memory node receives such a warning, the memory node informs the memory manager. In turn, the memory manager notifies any compute nodes that have allocated regions belonging to the soon-to-be-offline memory node. Those compute nodes stop swapping-out to the memory node, but may continue to swap-in from the node as long as the node is still alive. Meanwhile, the memory manager orchestrates the migration of regions from the soon-to-be-offline memory node to other memory nodes. When a particular region’s migration has completed, the memory manager informs the relevant compute node, who then updates the local mapping from Region ID to backing memory node. At some point during this process, the memory manager may also request non-failing memory nodes to contribute additional regions to the global pool of far memory.

Unplanned Failures: On a compute node, the Carbink runtime is responsible for detecting the unplanned failure of a

memory node. The runtime does so via connection timeouts or more sophisticated leasing protocols [15, 16]. Upon detecting an unplanned failure, the runtime spawns background threads to reconstruct the affected spans using erasure coding. The runtime is also responsible for allowing application threads to read spans whose recovery is in-flight.

Span reconstruction: To reconstruct the spans belonging to a failed memory node M_{fail} , a compute node first requests a new region from the memory manager. Suppose that the new region is provided by memory node M_{new} . The compute node iterates through each lost spanset associated with M_{fail} ; for each spanset, the compute node tells M_{new} which external spans and parity fragments to read in order to erasure-code-restore M_{fail} ’s data. As the relevant spans are restored, a compute node can still swap in and remotely compact those spans. However, the swap-in and remote compaction activity will have to synchronize with recovery activity (§3.6).

In EC-Batch Local, when a compute node detects a memory node failure, the compute node cancels all in-flight compactions involving that node. A compute node using EC-Batch Remote does the same; however, for each canceled compaction, the compute node must also instruct the surviving memory nodes in the 2PC group to cancel the transaction.

The data and parity for a swapped-out spanset reside on multiple memory nodes. As a compute node recovers from the failure of one of the nodes in that group, another node in the group may fail. As long as the number of failed nodes does not exceed the number of parity nodes, Carbink can recover the spanset. The reason is that all of the information needed to recover is stored on a compute node, e.g., in the far page heap (§3.3). Due to space limitations, we omit a detailed explanation of how Carbink deals with concurrent failures.

Degraded reads: During the reconstruction of an affected span, application threads may try to swap in the span. The runtime handles such a fetch using a *degraded read* protocol. For example, consider Figure 4a. Suppose that $M1$ fails unexpectedly, and while the Carbink runtime is recovering $M1$ ’s spans ($D1$ and $D5$), an application thread tries to read an object residing in $D1$. The runtime will swap in data spans $D2$, $D3$, and $D4$, as well as parity fragment $P1$, and then reconstruct $D1$ via erasure coding. Degraded reads ensure that the failure of a memory node merely slows down an application instead of blocking it. In Section 5.3, we show that application performance only drops for 0.6 seconds, and only suffers a throughput degradation of 36% during that time.

Network bandwidth consumption: During failure recovery, Carbink consumes the same amount of network bandwidth as Hydra. For example, suppose that both Hydra and Carbink use RS4.2 encoding and have 4 spans, with a span stored on each of 4 memory nodes. In Hydra, a single node failure will lose four 1/4th spans. Reconstructing each 1/4th span will require the reading of four 1/4th span/parity regions from the surviving nodes, resulting in an aggregate network bandwidth requirement of 1 full span. So, reconstructing four 1/4th spans

will require an aggregate network bandwidth of 4 full spans. In Carbink, the failure of a single memory node results in the loss of 1 full span. To recover that span, Carbink (like Hydra) must read 4 span/parity regions.

3.6 Thread Synchronization

On a compute node, the main kinds of Carbink threads are applications threads (which read objects, write objects, and swap in spans), filtering threads (which move objects within local spans), and eviction threads (which reclaim space by swapping local spansets to far memory). At any given time, a span may be in one of two concurrency regimes (§3.3): the span is either accessible to application threads and filtering threads, or to application threads and eviction threads. In both regimes, Carbink has to synchronize how the relevant threads update Carbink’s smart pointers (§3.2).

At a high level, Carbink uses an RCU locking scheme that is somewhat reminiscent of AIFM’s approach [44]. Due to space restrictions, we merely sketch the design. Carbink optimizes for the common case in which a span is only being accessed by an application thread. In this common case, an application thread grabs an RCU read lock on the pointer via the pointer’s `Deref()` method, as shown in Listing 1. The thread sees that either (1) the **Present** bit is not set, in which case the Carbink runtime issues an RMA read to swap in the appropriate span; (2) alternatively, the thread sees that the **Present** bit is set, but the **M** and **E** bits are unset. In the second case, `Deref()` can just return a normal pointer back to the application. The application can be confident that concurrent filtering or evicting threads will not move or evict the object, because those threads cannot touch the object until application-level threads have released their RCU read locks via the `DerefScope` destructor (Listing 1).

The more complicated scenarios arise when the **Present** bit is set and either the **M** or **E** bit are set as well. In this case, the (say) **M** bit has been set because the filtering thread set the bit and then called `SyncRCU()` (i.e., the RCU write waiting lock). The concurrent application thread and filtering thread essentially race to acquire the pointer’s spinlock; if the application thread (i.e., `Deref()`) wins, it makes a copy of the object, clears **M**, releases the spinlock, and returns the address of the object copy to the application. Otherwise, if the filtering thread wins, it moves the object, clears **M**, and releases the spinlock. The losing thread has to retry the desired action. An analogous situation occurs if the **E** bit is set.

Carbink’s eviction and remote compaction threads directly poll the network stack to learn about RMA completions and RPC completions. An application thread which has issued an RMA swap-in operation will yield, but a dedicated RMA poller thread detects when application RMAs have completed and awakens the relevant application threads. Polling avoids the overheads of context switching to new threads and notifying old threads that network events have occurred.

During recovery (§3.5), Carbink spawns additional threads to orchestrate the reconstruction of spans. Those threads acquire per-spanset mutexes which are also acquired by threads performing swap-ins, swap-outs, and remote compactions.

4 Implementation

Our Carbink prototype contains 14.3K lines of C++. It runs atop unmodified OSes, using standard POSIX abstractions for kernel-visible threads and synchronization. The runtime leverages the PonyExpress user-space network stack [35]. On a compute node, all threads in a particular application (both application-defined threads and Carbink-defined threads) execute in the same process. On a memory node, a Carbink daemon exposes far memory via RMAs or RPCs. We use Intel ISA-L v2.30.0 [25] for Reed-Solomon erasure coding.

Our current prototype has a simplified memory manager that is unreplicated, does not handle planned failures, and statically assigns memory nodes to compute nodes. Implementing the full version of the memory manager will be conceptually straightforward, since we can use off-the-shelf libraries for replicated state machines [1, 45] and cluster management [51, 52]. We also note that the experiments in §5 are insensitive to the performance of the memory manager, regardless of whether the manager is replicated or not. The reason is that memory allocations and deallocations (which must be routed through the memory manager) are rare and are not on the critical path of steady-state compute node operations like swap-in and swap-out.

To better understand the performance overheads of Carbink’s erasure-coding approach, we built an AIFM-like [44] far memory system. That system uses remotable pointers like Carbink, but swaps in and out at the granularity of objects, and provides no fault tolerance. Like Carbink, it leverages the PonyExpress [35] user-space network stack. Our AIFM clone is 5.8K lines of C++.

5 Evaluation

In this section, we answer the following questions:

1. What is the latency, throughput, and remote memory usage of EC-Batch compared with the other fault tolerance schemes (§5.1 and §5.2)?
2. How does an unplanned memory node failure impact the performance of Carbink applications (§5.3)?
3. How does the performance of Carbink’s span-based memory organization compare to the performance of an AIFM-like object-level approach (§5.4)?

Testbed setup: We deployed eight machines in the same rack, including one compute node and seven memory nodes; one of the memory nodes was used for failover. Each machine was equipped with dual-socket 2.2 GHz Intel Broadwell processors and a 50 Gbps NIC.

Fault tolerance schemes: Using the Carbink runtime, we compared our proposed EC-Batch schemes to four ap-

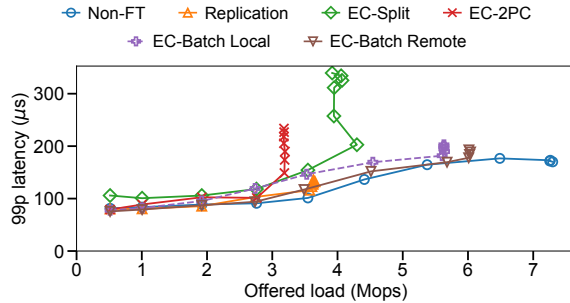


Figure 5: Microbenchmark load-latency curves.

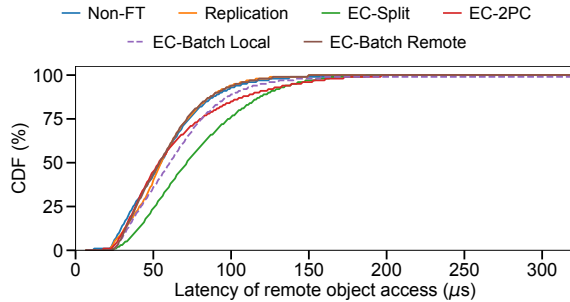


Figure 6: Latency distribution of remote object accesses in the microbenchmark under an offered load of 2 Mops.

proaches: Non-FT (a non-fault-tolerant scheme that used RMA to swap spans), Replication (which replicated spans on multiple nodes), EC-Split (the approach used by Hydra [29]), and EC-2PC (Table 1). We configured all fault tolerance schemes to tolerate up to two memory node failures. So, the Replication scheme replicated each swapped-out span on three memory nodes, whereas the EC schemes used six memory nodes—four held data, and two held RS4.2 parity bits [43]. EC-Batch spawned two compaction threads by default.

As mentioned in Section 4, we also built an AIFM-like far memory system. This system did not provide fault tolerance, but it provided a useful comparison with our Non-FT Carbink version.

Carbink borrows the span sizes that are used by TCMalloc (§3.3). These parameters have been empirically observed to reduce internal fragmentation. In our evaluation, EC-Batch (both Local and Remote) grouped four equal-size spans into a spanset, swapping out at the granularity of a spanset. Increasing spanset sizes would allow Carbink to issue larger batched RMAs, improving network efficiency. However, spansets whose evictions are in progress must be locked in local memory while RMAs complete; thus, larger spanset sizes would delay the reclamation of larger portions of local memory.

5.1 Microbenchmarks

To get a preliminary idea of Carbink’s performance, we created a synthetic benchmark that wrote 15 million 1 KB objects (totalling 15 GB) to a remotable array. The compute node’s local memory had space to store 7.5 GB of objects (i.e., half of the total set). By default, the compute node spawned 128

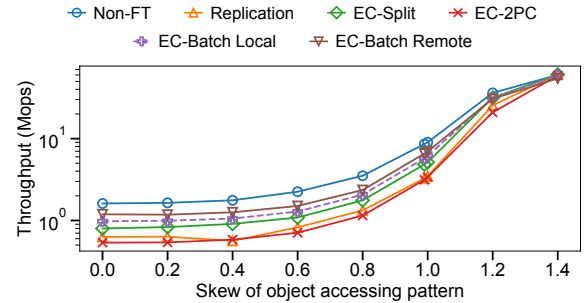


Figure 7: Impact of skew on throughput.

threads on 32 logical cores to access objects; the access pattern had a Zipfian-distributed [41] skew of 0.99. Such skews are common in real workloads for key/value stores [7].

Object access throughput and tail latency: Figure 5 shows the 99th-percentile latency with various object access loads. All of the fault-tolerant schemes eventually hit a “hockey stick” in tail latency growth when the schemes could no longer catch up with the offered load. EC-Batch Remote had the highest sustained throughput (6.0 Mops), which was 40% higher than the throughput of the state-of-the-art EC-Split (4.3 Mops). EC-Batch Local achieved 5.6 Mops, which was 30% higher than EC-Split. EC-Split had worse performance because it had to issue four RMA requests to swap in one span; thus, EC-Split quickly became bottlenecked by network IO. In contrast, EC-Batch only issued one RMA request per swap-in.

EC-Batch Remote had 18%-29% lower tail latency than EC-Split under the same load (before reaching the “hockey-stick”). The reason was that EC-Split’s larger number of RMAs per swap-in left EC-Split more vulnerable to stragglers [29]. Also recall that EC-Batch can support computation offloading [3, 27, 44, 57], which is hard with EC-Split (§3.4).

EC-2PC had the worst throughput because it relied on costly RPCs and 2PC protocols to swap out spans. Thus, EC-2PC could not reclaim local memory as fast as other schemes. The Replication scheme was bottlenecked by network bandwidth, since every swap-out incurred a $3\times$ network write penalty; in contrast, EC-based schemes used RS4.2 erasure coding to reduce the write penalty to $1.5\times$.

Latency distribution of remote object accesses: Figure 6 shows the latency of accessing remote objects under 2 Mops of offered load. With this low offered load, Replication and EC-Batch Remote achieved similar access latencies as Non-FT because none of the schemes were bottlenecked by network bandwidth. EC-Batch Local had slightly higher remote access latencies. However, EC-Split had significantly higher access latencies (e.g., at the median and tail) than EC-Batch Local and Remote; the reason was that EC-Split issued four times as many network IOs and thus was more sensitive to stragglers. EC-2PC’s tail latency was slightly higher than that of EC-Batch Local and Remote due to the overhead of costly RPCs and 2PC traffic.

Impact of skewness: Figure 7 shows how the skewness of object accesses impacted throughput. EC-Batch Remote and

	# Compaction threads	Norm. remote mem usage	Avg. # remote logical cores	Avg. BW (Gbps)
EC-Batch Local	1	2.54	0.23	1.27
	2	2.35	0.53	1.64
	3	2.28	0.56	1.76
EC-Batch Remote	1	1.89	1.97	2.98
	2	1.83	2.10	3.15
	3	1.74	2.27	3.40
W/o compaction	0	3.03	—	—

Table 2: Remote resource usage in the microbenchmark. The remote memory usage is normalized with respect to the usage of Non-FT. The number of remote logical cores and the network bandwidth are averaged across all six memory nodes.

Local performed best due to their more efficient swapping approaches. However, the throughput of all schemes increased with higher skewness. The reason is that high skewness led to a smaller working set and thus a higher likelihood that hot objects were locally resident. In these scenarios, schemes with faster swapping were not rewarded as much.

Remote resource usage with compaction: Table 2 shows the impact of compaction on the average memory, CPU, and bandwidth usage per memory node. Without compaction, EC-Batch used $3.03\times$ remote memory (normalized with respect to Non-FT memory consumption). With two local compaction threads, EC-Batch Remote’s memory overhead reduced to $1.83\times$. The memory reduction was at the expense of 2.1 cores and 3.15 Gbps bandwidth on each memory node. With more compaction threads, Carbink could further reduce memory usage at the cost of higher CPU and bandwidth utilization. That being said, we note that the synthetic microbenchmark application represented an extreme case of remote CPU and network usage, since the workload accessed objects without actually computing on them.

EC-Batch Remote vs. Local: EC-Batch Remote had higher throughput and lower tail latency than EC-Batch Local (Figure 5). This was because EC-Batch Local’s compaction required (1) local CPUs for parity computation and (2) network bandwidth for transferring span deltas and parity updates, leaving fewer local resources for application threads and RMA reads. Because of EC-Batch Remote’s faster compaction, EC-Batch Remote also used 28%-34% less remote memory than EC-Batch Local (Table 2). However, EC-Batch Remote consumed more remote CPUs (2.10 vs. 0.53 cores) and more network bandwidth (3.15 vs. 1.64 Gbps) than Local. In practice, the Carbink runtime could transparently switch between EC-Batch Remote and Local based on an application developer’s policy about resource/performance trade-offs.

5.2 Macrobenchmarks

We evaluated Carbink using two memory-intensive applications that would benefit from remote memory: an in-memory transactional key-value store, and a graph processing algorithm. The two applications exhibited different patterns of

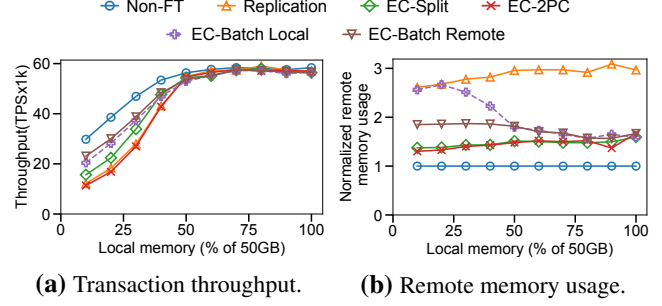


Figure 8: Transactional KV-store evaluation.

object accesses, and had different working set behaviors.

Transactional KV-store: This application implemented a transactional in-memory B-tree, exposing it via a key/value interface similar to that of MongoDB [37]. Each remotable object was a 4 KB value stored in a B-tree leaf. The application spawned 128 threads, and each thread processed 20 K transactions. The compute node provisioned 32 logical cores, with the application overlapping execution of the threads for higher throughput [26, 38, 44, 56]. Each transaction contained three reads and three writes, similar to the TPC-A benchmark [53]. Each update created a new version of a particular key’s value; asynchronously, the application trimmed old versions. The maximum working set size during the experiment was roughly 50 GB.

Throughput: Figure 8a shows the KV-store throughput when varying the size of local memory (normalized as a fraction of the maximum working set size). In scenarios with less than 50% local memory, EC-Batch Remote achieved higher transactions per second (TPS) than all other fault tolerance schemes. For example, TPS for EC-Batch Remote was 1.5%-48% higher than that of EC-Split; this was because EC-Batch only needed one RMA request to swap in a span. EC-Batch Remote was at most 29% slower than Non-FT, mainly due to the additional parity update required for fault tolerance. EC-Batch Local was at most 13% slower than EC-Batch Remote. EC-2PC performed the worst among EC schemes.

All schemes achieved similar throughput when the local memory size was above 50%. The reason was that the *average* working set size of the workload was only half the size of the *maximum* memory usage. The maximum memory usage only occurred when the B-Tree had fallen very behind in culling old versions of objects.

Remote memory usage: Figure 8b plots remote memory usage as a function of local memory sizes; remote memory usage is normalized with respect to that of Non-FT. Compared to EC-Split, EC-Batch Remote and Local used up to 35% and 93% more remote memory, respectively. EC-Batch schemes defragmented remote memory using compaction, but when local memory space was less than 50%, remote compaction could not immediately defragment the spanset holes created by frequent span swap-ins. As local memory grew larger, span fetching became less frequent, making it

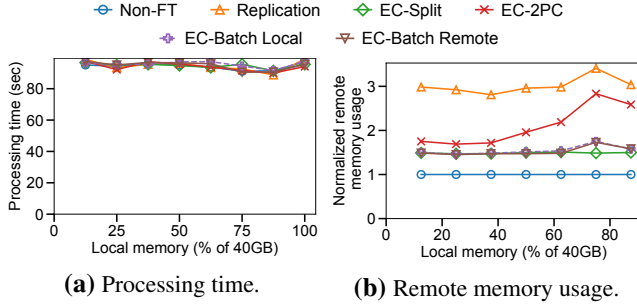


Figure 9: Graph processing evaluation.

easier for remote compaction to reclaim space. In this less hectic environment, EC-Batch’s remote memory usage was similar to that of the other erasure-coding schemes.³

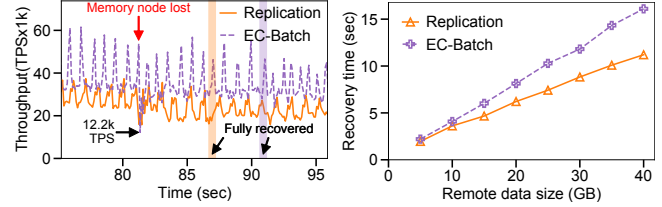
Graph processing: We implemented a connected-components algorithm [50] that found all sets of linked vertices in a graph. This kind of algorithm is critical to various Google services. We evaluated the algorithm using the Friendster graph [30] which contained 65 million vertexes and 1.8 billion edges. In the graph analysis code, each vertex’s adjacency list was referenced via remotable pointers. The total size of the objects stored in Carbink was roughly 40 GB. The application used 80 application threads that ran atop 80 logical cores. In our experimental results, the reported processing times exclude graph loading, since graph loading is dominated by disk latencies.

Figure 9a shows that all schemes had similar processing times as Non-FT, regardless of the local memory size. The reason was that the graph application had a high compute-to-network ratio—the application fetched all neighbors associated with each vertex and then spent non-trivial time enumerating each neighbor and computing on them. As a result of this good spatial locality and high “think time,” the graph application did not incur frequent data swapping, and thus avoided fault tolerance overhead that the KV-store could not.

Figure 9b shows that EC-Batch Local and Remote had similar remote memory usage as EC-Split: 15%-39% lower than EC-2PC and roughly 50% lower than Replication. All EC-based schemes had lower remote memory overheads than Replication because the erasure coding only incurred a $1.5\times$ space overhead for the extra parity data.

EC-2PC used more memory than EC-Batch because the graph workload randomly fetched diverse-sized spans. The random fetch sizes reflected the fact that different vertices had different sizes for their adjacency lists. This lack of span size locality hindered dead space reclamation, since EC-2PC had to wait longer for all of the spans in an erasure-coding group to be swapped in. EC-Batch avoided this problem by bundling equal-sized spans into the same spanset and using remote compaction.

³The remote memory usage of triple-replication was slightly less than $3\times$ the usage of Non-FT because Non-FT could swap out memory faster during periods of high local memory pressure.



(a) KV-store TPS over time. (b) Microbenchmark recovery.

Figure 10: Failure recovery evaluation.

5.3 Failure Recovery

We measured the recovery time for an unplanned memory node failure in the KV-store, the graph processor, and the microbenchmark application. For the graph application, all schemes achieved similar processing time during unplanned failures; thus, in the text below, we focus on the KV-store and the microbenchmark.

Transactional KV-store: Figure 10a shows the KV-store throughput of Replication and EC-Batch Local, with a data point collected every 100 ms before and after an unplanned memory node failure. Upon detecting the failure, EC-Batch Local immediately reconstructed the lost data on a pre-configured failover memory node. We gave the KV-store 15 GB of local memory, equivalent to 30% of the 50 GB maximum working set size.

The throughput of both schemes fluctuated sinusoidally because the KV-store frequently tried to swap in remote objects, but the swap-ins sometimes had to synchronously block until eviction threads could reclaim enough local memory. After a memory node failed, EC-Batch needed 0.6 seconds to restore normal throughput, while replication needed 0.3 seconds. This is because, during failure recovery, an EC-Batch read that targeted an affected span used the degraded read protocol which uses more bandwidth than a normal read (§3.5); in contrast, a Replication read that targeted an affected span consumed the same amount of bandwidth as a read during non-failure-recovery. During recovery, the throughput of Replication and EC-Batch dropped an average of 35% and 36% respectively.

EC-Batch required 9.7 seconds to fully regenerate the lost data on the failover node, taking $1.7\times$ longer than Replication. This difference arose because, in EC-Batch, the new memory node read $4\times$ span/parity information involving the lost data and computed erasure codes to reconstruct the lost data. In contrast, Replication lost more data per memory node, but only read one copy of the lost data. Note that with EC-Batch, degraded reads mostly happened during the first second of failure recovery; the skewed workload meant that a small number of objects were the targets of most reads, and once a hot object was pulled into local memory (perhaps by a degraded read), the object would not generate additional degraded reads.

Microbenchmark: Figure 10b shows recovery times as a function of the remote data size. The recovery time of EC-Batch increased almost linearly with the remote data size,

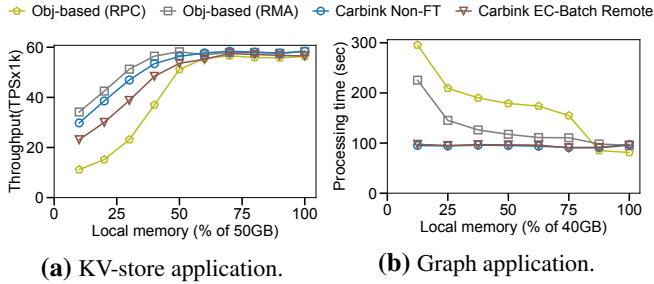


Figure 11: Application performance: AIFM-like object-based systems and Carbink.

with 0.6 GB/s recovery speed. This speed was 12%-44% slower than Replication due to the larger amount of recovery information that EC-Batch had to transfer around the network, and the computational overhead of generating erasure codes.

Prior work [10, 29, 58] also found that, during recovery, erasure-coding schemes had longer recovery times and worse performance degradation than replication schemes. However, this drawback only happens for unplanned failures which, in our production environment, are rare compared to planned failures; in an erasure-coding scheme, handling a *planned* failure just requires simple copying of the information on a departing memory node, and does not incur additional work to find parity information or recompute erasure coding. Thus, in our deployment setting where unplanned failures are rare, erasure-coding schemes (which have lower memory utilization than replication schemes) are very attractive.

5.4 Comparison with AIFM-like Systems

We compared span-based swapping in Carbink with the object-based approach used in AIFM [44]. We implemented two AIFM-like systems using our threading and network stack (§4). The first system used RPCs to swap individual objects, with the remote memory nodes tracking the object-to-remote-location mapping (as done in AIFM). Our second object-granularity swapping system used more-efficient RMAs to swap objects, and had compute nodes track the mapping between objects and their remote locations; recall that RMA is one-sided, so compute nodes could not rely on memory nodes to synchronously update mappings during swaps. Like the original AIFM, neither system provided fault tolerance.

Transactional KV-store: Figure 11a shows that, if local memory was too small to hold the average working set, Non-FT Carbink had 45%-167% higher throughput than the AIFM-like system with RPC. The reason is that, when local memory pressure was high, more swapping occurred, and the better efficiency of RMAs over RPCs became important. However, Non-FT Carbink achieved 5.6%-15% lower throughput than the object-based system with RMA. This was due to swap-in amplification. For example, Non-FT Carbink might swap in an 8KB span but only use one 4KB object in the span; this never happens in a system that swaps at an object granularity.

Graph processing: Figure 11b shows the graph application’s processing time. When the local memory size was below 87.5%, Carbink performed 18%-58% faster than the object-based system with RMA. This is because, in the graph workload, 4% of large objects occupied 50% of the overall data set. Carbink prioritized swapping out large cold objects (§3.3), keeping most small objects in local memory and reducing the miss rate for those objects. In contrast, the object-based systems did not consider object sizes when swapping, leading to an increased miss rate for small objects. Note that, with larger local memories, all schemes had similar performance; indeed, when all objects fit into local memory, the object-based system with RPC slightly outperformed the rest because it did not require a dedicated core to poll for RMA completions.

6 Discussion

EC-Batch for paging-based systems: Carbink uses EC-Batch to transparently expose far memory via remotable pointers. However, EC-Batch can also be used to expose far memory via OS paging mechanisms [5, 22, 46]. In a traditional paging-based approach for far memory, a compute node swaps in and out at the granularity of a page. However, a compute node can use EC-Batch to treat each page as a span, such that pages are swapped out at the “pageset” granularity, and pages are swapped in at the page granularity.

Custom one-sided operations: EC-Batch requires memory nodes to calculate span deltas and parity updates (§3.4.2). In our Carbink prototype, memory nodes use separate threads to execute these calculations. However, memory nodes could instead implement them as custom one-sided operations in the network stack, such that the network stack itself performs the calculations, avoiding the need to context-switch to external threads. This approach has been used in prior work [6, 9, 35, 47, 48] to avoid thread scheduling overheads.

Designing the memory manager: We used a centralized manager because such a manager (1) simplified our overall design, and (2) made it easier to drive memory utilization high (because a centralized manager will have a global, accurate view of memory allocation metadata). A similarly-centralized memory manager is used by the distributed transaction system FaRM [16]. If the centralized manager became unavailable, Carbink could fall back to a decentralized memory allocation scheme like the one used by Hydra [29] or INFINISWAP [22].

The state maintained by the memory manager is not large. With 1 GB regions, we expect up to 500 regions in a typical memory node (similar to FaRM [16]). With thousands of memory nodes, the memory manager just needs to store a few MBs of state for region assignments.

Fault tolerance for compute nodes: In Carbink, a compute node does not share memory with other compute nodes. Thus, a Carbink application can checkpoint its own state without fear of racing with other compute nodes that modify the state being checkpointed. Checkpoint data could be placed in a

	Fast s/o	Low mem	Fast s/i	Interface	Coding granularity
On-disk repl.	✗	✓	✓	Various	–
In-memory repl.	✓	✗	✓	Various	–
Hydra [29]	✓	✓	✗	Paging	Split 4KB pages
Cocytus [10]	✓	✓	✗	KV-store	Across 4KB pages
BCStore [31]	✓	✓	✗	KV-store	Across objs
Hybrid [32]	✗	✗	✓	KV-store	Split 4KB pages
Carbink	✓	✓	✓	Remotable pointers	Across spans

Table 3: Comparison of existing fault-tolerant approaches for far memory. “Fast s/o” indicates whether a system can swap out at network/memory speeds. “Low mem” means that a system has relatively low memory pressure. “Fast s/i” refers to whether a system can swap in at network/memory speeds.

non-Carbink store, obviating the need to track how checkpointed spans move across Carbink memory nodes during compaction and invalidation. Alternatively, Carbink itself could store checkpoints, e.g., in the fault-tolerant address space of a well-known Carbink application whose sole purpose is to store checkpoints.

7 Related Work

Fault tolerance for far memory: Many far memory systems do not provide fault tolerance [2, 44, 55]. Of the systems that do, most replicate swapped-out data to local disks or remote ones [5, 22, 46]. Unfortunately, this approach forces application performance to bottleneck on disk bandwidth or disk IOPs during bursty workloads or failure recovery [29]. This behavior is unattractive, since a primary goal of a far memory system is to have applications run at *memory* speeds as much as possible.

Like Carbink, Hydra [29] is a far memory system that provides fault tolerance by writing erasure-coded local memory data to far RAM. Hydra uses the EC-Split coding approach that we describe in Section 3.4. As we demonstrate in Section 5, Carbink’s erasure-coding scheme provides better application performance in exchange for somewhat higher memory consumption. Carbink’s coding scheme also enables the offloading of computations to far memory nodes. Such offloading can significantly improve the performance of various applications [3, 27, 44, 57].

Fault tolerance for in-memory transactions and KV-stores: In-memory transaction systems typically provide fault tolerance by replicating data across the memory of multiple nodes [15, 16, 26]. These approaches suffer from the classic disadvantages of replication: double or triple storage overhead, and the associated increase in network traffic.

Recent in-memory KV-stores use erasure coding to provide fault tolerance. For example, Cocytus [10] and BCStore [31] only rely on in-memory replication to store small instances of metadata; object data is erasure-coded using a default page size of 4KB. Cocytus erasure-codes using a scheme that resembles EC-2PC (§3.4). To reduce the network utilization of

a Cocytus-style approach, a BCStore compute node buffers outgoing writes; this approach allows the node to batch the computation of parity fragments (and thus issue fewer updates to remote data and parity regions). Batching reduces network overhead at the cost of increasing write latency.

Both Cocytus and BCStore rely on two-sided RPCs to manipulate far memory. RPCs incur software-level overheads involving thread scheduling and context switching on remote nodes. To avoid these costs, Carbink eschews RPCs for one-side RMA operations. Carbink also issues fewer parity updates than Cocytus; whereas Cocytus uses expensive 2PC to update parity information during every write, Carbink defers parity updates until compaction occurs on remote nodes (§3.4.2). Carbink’s compaction approach is also more efficient than that of BCStore. BCStore’s compaction algorithm performs actual copying of data objects on memory nodes, whereas Carbink compaction just manipulates span pointers inside of spanset metadata.

A far memory system could use both replication and erasure coding [32]. For example, during a Hydra-style swap-out, a span would be erasure-coded and the fragments written to memory nodes; however, a full replica of the span would also be written out. Relative to Carbink, this hybrid approach would have lower reconstruction costs (assuming that the full replica did not live on the failed node). However, Carbink would have lower memory overheads because no full replica of a span would be stored. Carbink would also have faster swap-outs, because swap-outs in the hybrid scheme would require an EC-2PC-like mechanism to ensure consistency.

Table 3 summarizes the strengths and weaknesses of the various systems discussed above.

Memory compaction: In Carbink, the far memory regions used by a program become fragmented as spans are swapped in. Memory compaction is a well-studied topic in the literature about “moving” garbage collectors for managed languages (e.g., [11, 18, 49]). Moving garbage collection is also possible for C/C++ programs; Mesh [40] represents the state-of-the-art. With respect to this prior work, Carbink’s unique challenge is that the compaction algorithm (§3.4.2) must compose well with an erasure coding scheme that governs how objects move between local memory and far memory.

8 Conclusion

Carbink is a far memory system that provides low-latency, low-overhead fault tolerance. Carbink erasure-codes data using a span-centric approach that does not expose swap-in operations to stragglers. Whenever possible, Carbink uses efficient one-sided RMAs to exchange data between compute nodes and memory nodes. Carbink also uses novel compaction techniques to asynchronously defragment far memory. Compared to Hydra, a state-of-the-art fault-tolerant system for far memory, Carbink has 29% lower tail latency and 48% higher application performance, with at most 35% higher memory usage.

Acknowledgments

We thank our shepherd Luís Rodrigues and the anonymous reviewers for their insightful comments. We also thank Kim Keeton and Jeff Mogul for their comments on early drafts of the paper, and Maria Mickens for her comments on a later draft. Yang Zhou and Minlan Yu were supported in part by NSF CNS-1955422 and CNS-1955487.

References

- [1] Ittai Abraham, Dahlia Malkhi, Kartik Nayak, Ling Ren, and Maofan Yin. Sync HotStuff: Simple and Practical Synchronous State Machine Replication. In *Proceedings of IEEE Symposium on Security and Privacy*, pages 106–118, 2020.
- [2] Marcos K. Aguilera, Nadav Amit, Irina Calciu, Xavier Deguillard, Jayneel Gandhi, Stanko Novakovic, Arun Ramanathan, Pratap Subrahmanyam, Lalith Suresh, Kiran Tati, and et al. Remote Regions: A Simple Abstraction for Remote Memory. In *Proceedings of USENIX ATC*, pages 775–787, 2018.
- [3] Marcos K. Aguilera, Kimberly Keeton, Stanko Novakovic, and Sharad Singhal. Designing Far Memory Data Structures: Think Outside the Box. In *Proceedings of ACM HotOS*, pages 120–126, 2019.
- [4] Shoaib Akram, Jennifer B. Sartor, Kathryn S. McKinley, and Lieven Eeckhout. Write-Rationing Garbage Collection for Hybrid Memories. *ACM SIGPLAN Notices*, 53(4):62–77, 2018.
- [5] Emmanuel Amaro, Christopher Branner-Augmon, Zhihong Luo, Amy Ousterhout, Marcos K. Aguilera, Aurojit Panda, Sylvia Ratnasamy, and Scott Shenker. Can Far Memory Improve Job Throughput? In *Proceedings of ACM EuroSys*, pages 1–16, 2020.
- [6] Emmanuel Amaro, Zhihong Luo, Amy Ousterhout, Arvind Krishnamurthy, Aurojit Panda, Sylvia Ratnasamy, and Scott Shenker. Remote Memory Calls. In *Proceedings of ACM HotNets*, pages 38–44, 2020.
- [7] Berk Atikoglu, Yuehai Xu, Eitan Frachtenberg, Song Jiang, and Mike Paleczny. Workload Analysis of a Large-Scale Key-Value Store. *ACM SIGMETRICS Performance Evaluation Review*, 40(1):53–64, 2012.
- [8] Cristina Băescu and Bryan Ford. Immunizing Systems from Distant Failures by Limiting Lamport Exposure. In *Proceedings of ACM HotNets*, pages 199–205, 2021.
- [9] Matthew Burke, Shannon Joyner, Adriana Szekeres, Jacob Nelson, Irene Zhang, and Dan R.K. Ports. PRISM: Rethinking the RDMA Interface for Distributed Systems. In *Proceedings of USENIX SOSP*, pages 228–242, 2021.
- [10] Haibo Chen, Heng Zhang, Mingkai Dong, Zhaoguo Wang, Yubin Xia, Haibing Guan, and Binyu Zang. Efficient and Available In-Memory KV-Store with Hybrid Erasure Coding and Replication. *ACM Transactions on Storage (TOS)*, 13(3):1–30, 2017.
- [11] Jon Coppeard. Compacting Garbage Collection in SpiderMonkey. <https://hacks.mozilla.org/2015/07/compacting-garbage-collection-in-spidermonkey/>, 2015.
- [12] Fernando J. Corbato. A Paging Experiment with the Multics System. Technical report, MASSACHUSETTS INST OF TECH CAMBRIDGE PROJECT MAC, 1968.
- [13] James C. Corbett, Jeffrey Dean, Michael Epstein, Andrew Fikes, Christopher Frost, Jeffrey John Furman, Sanjay Ghemawat, Andrey Gubarev, Christopher Heiser, Peter Hochschild, and et al. Spanner: Google’s Globally Distributed Database. *ACM Transactions on Computer Systems (TOCS)*, 31(3):1–22, 2013.
- [14] Jeffrey Dean and Sanjay Ghemawat. MapReduce: Simplified Data Processing on Large Clusters. *Communications of the ACM*, 51(1):107–113, 2008.
- [15] Aleksandar Dragojević, Dushyanth Narayanan, Miguel Castro, and Orion Hodson. FaRM: Fast Remote Memory. In *Proceedings of USENIX NSDI*, pages 401–414, 2014.
- [16] Aleksandar Dragojević, Dushyanth Narayanan, Edmund B. Nightingale, Matthew Renzelmann, Alex Shamis, Anirudh Badam, and Miguel Castro. No Compromises: Distributed Transactions with Consistency, Availability, and Performance. In *Proceedings of ACM SOSP*, pages 54–70, 2015.
- [17] Jason Evans. A Scalable Concurrent malloc (3) Implementation for FreeBSD. In *Proceedings of BSDCan Conference*, 2006.
- [18] Robert R. Fenichel and Jerome C. Yochelson. A LISP Garbage-Collector for Virtual-Memory Computer Systems. *Communications of the ACM*, 12(11):611–612, 1969.
- [19] Joseph E. Gonzalez, Yucheng Low, Haijie Gu, Danny Bickson, and Carlos Guestrin. Powergraph: Distributed Graph-Parallel Computation on Natural Graphs. In *Proceedings of USENIX OSDI*, pages 17–30, 2012.
- [20] Joseph E. Gonzalez, Reynold S. Xin, Ankur Dave, Daniel Crankshaw, Michael J. Franklin, and Ion Stoica. Graphx: Graph Processing in a Distributed Dataflow Framework. In *Proceedings of USENIX OSDI*, pages 599–613, 2014.

- [21] Google. TCMalloc Open Source. <https://github.com/google/tcmalloc>.
- [22] Juncheng Gu, Youngmoon Lee, Yiwen Zhang, Mosharaf Chowdhury, and Kang G. Shin. Efficient Memory Disaggregation with INFINISWAP. In *Proceedings of USENIX NSDI*, pages 649–667, 2017.
- [23] Xianglong Huang, Stephen M Blackburn, Kathryn S. McKinley, J. Eliot B. Moss, Zhenlin Wang, and Perry Cheng. The Garbage Collection Advantage: Improving Program Locality. *ACM SIGPLAN Notices*, 39(10):69–80, 2004.
- [24] Andrew Hamilton Hunter, Chris Kennelly, Paul Turner, Darryl Gove, Tipp Moseley, and Parthasarathy Ranganathan. Beyond Malloc Efficiency to Fleet Efficiency: A Hugepage-Aware Memory Allocator. In *Proceedings of USENIX OSDI*, pages 257–273, 2021.
- [25] Intel. Intel Intelligent Storage Acceleration Library. <https://github.com/intel/isa-1>.
- [26] Anuj Kalia, Michael Kaminsky, and David G. Andersen. FaSST: Fast, Scalable and Simple Distributed Transactions with Two-Sided RDMA Datagram RPCs. In *Proceedings of USENIX OSDI*, pages 185–201, 2016.
- [27] Dario Korolija, Dimitrios Koutsoukos, Kimberly Keeton, Konstantin Taranov, Dejan Milojević, and Gustavo Alonso. Farview: Disaggregated Memory with Operator Off-loading for Database Engines. In *Proceedings of Conference on Innovative Data Systems Research*, 2022.
- [28] Jakub Łacki, Vahab Mirrokni, and Michał Włodarczyk. Connected Components at Scale via Local Contractions. *arXiv preprint arXiv:1807.10727*, 2018.
- [29] Youngmoon Lee, Hasan Al Maruf, Mosharaf Chowdhury, Asaf Cidon, and Kang G. Shin. Mitigating the Performance-Efficiency Tradeoff in Resilient Memory Disaggregation. *arXiv preprint arXiv:1910.09727*, 2019.
- [30] Jure Leskovec. Friendster Social Network Dataset. <https://snap.stanford.edu/data/com-Friendster.html>.
- [31] Shenglong Li, Quanlu Zhang, Zhi Yang, and Yafei Dai. BCStore: Bandwidth-Efficient In-Memory KV-store with Batch Coding. *Proceedings of IEEE International Conference on Massive Storage Systems and Technology*, 2017.
- [32] Yuzhe Li, Jiang Zhou, Weiping Wang, and Yong Chen. RE-Store: Reliable and Efficient KV-Store with Erasure Coding and Replication. In *Proceedings of IEEE International Conference on Cluster Computing*, pages 1–12, 2019.
- [33] Yucheng Low, Joseph E. Gonzalez, Aapo Kyrola, Danny Bickson, Carlos E. Guestrin, and Joseph Hellerstein. Graphlab: A New Framework for Parallel Machine Learning. *arXiv preprint arXiv:1408.2041*, 2014.
- [34] Chengzhi Lu, Kejiang Ye, Guoyao Xu, Cheng-Zhong Xu, and Tongxin Bai. Imbalance in the Cloud: An Analysis on Alibaba Cluster Trace. In *Proceedings of IEEE International Conference on Big Data*, pages 2884–2892, 2017.
- [35] Michael Marty, Marc de Kruijf, Jacob Adriaens, Christopher Alfeld, Sean Bauer, Carlo Contavalli, Michael Dalton, Nandita Dukkipati, William C. Evans, Steve Gribble, and et al. Snap: A Microkernel Approach to Host Networking. In *Proceedings of ACM SOSP*, pages 399–413, 2019.
- [36] Paul E. McKenney and John D. Slingwine. Read-Copy Update: Using Execution History to Solve Concurrency Problems. In *Proceedings of Parallel and Distributed Computing and Systems*, pages 509–518, 1998.
- [37] MongoDB Inc. MongoDB Open Source. <https://github.com/mongodb/mongo>.
- [38] Amy Ousterhout, Joshua Fried, Jonathan Behrens, Adam Belay, and Hari Balakrishnan. Shenango: Achieving High CPU Efficiency for Latency-Sensitive Datacenter Workloads. In *Proceedings of USENIX NSDI*, pages 361–378, 2019.
- [39] John Ousterhout, Arjun Gopalan, Ashish Gupta, Ankita Kejriwal, Collin Lee, Behnam Montazeri, Diego Ongaro, Seo Jin Park, Henry Qin, Mendel Rosenblum, and et al. The RAMCloud Storage System. *ACM Transactions on Computer Systems (TOCS)*, 33(3):1–55, 2015.
- [40] Bobby Powers, David Tench, Emery D. Berger, and Andrew McGregor. Mesh: Compacting Memory Management for C/C++ Applications. In *Proceedings of ACM PLDI*, pages 333–346, 2019.
- [41] David M.W. Powers. Applications and Explanations of Zipf’s Law. In *Proceedings of New Methods in Language Processing and Computational Natural Language Learning*, 1998.
- [42] KV Rashmi, Mosharaf Chowdhury, Jack Kosaian, Ion Stoica, and Kannan Ramchandran. EC-Cache: Load-Balanced, Low-Latency Cluster Caching with Online Erasure Coding. In *Proceedings of USENIX OSDI*, pages 401–417, 2016.
- [43] Irving S. Reed and Gustave Solomon. Polynomial Codes over Certain Finite Fields. *Journal of the Society for Industrial and Applied Mathematics*, 8(2):300–304, 1960.

- [44] Zhenyuan Ruan, Malte Schwarzkopf, Marcos K. Aguilera, and Adam Belay. AIFM: High-Performance, Application-Integrated Far Memory. In *Proceedings of USENIX OSDI*, pages 315–332, 2020.
- [45] Fred B. Schneider. Implementing Fault-Tolerant Services Using the State Machine Approach: A Tutorial. *ACM Computing Surveys (CSUR)*, 22(4):299–319, 1990.
- [46] Yizhou Shan, Yutong Huang, Yilun Chen, and Yiying Zhang. LegoOS: A Disseminated, Distributed OS for Hardware Resource Disaggregation. In *Proceedings of USENIX OSDI*, pages 69–87, 2018.
- [47] David Sidler, Zeke Wang, Monica Chiosa, Amit Kulkarni, and Gustavo Alonso. StRoM: Smart Remote Memory. In *Proceedings of ACM EuroSys*, pages 1–16, 2020.
- [48] Arjun Singhvi, Aditya Akella, Maggie Anderson, Rob Cauble, Harshad Deshmukh, Dan Gibson, Milo M.K. Martin, Amanda Strominger, Thomas F. Wenisch, and Amin Vahdat. CliqueMap: Productionizing an RMA-Based Distributed Caching System. In *Proceedings of ACM SIGCOMM*, pages 93–105, 2021.
- [49] SUN Microsystems. Memory Management in the Java HotSpot Virtual Machine, 2006.
- [50] Michael Sutton, Tal Ben-Nun, and Amnon Barak. Optimizing Parallel Graph Connectivity Computation via Subgraph Sampling. In *Proceedings of IEEE International Parallel and Distributed Processing Symposium*, pages 12–21, 2018.
- [51] Chunqiang Tang, Kenny Yu, Kaushik Veeraraghavan, Jonathan Kaldor, Scott Michelson, Thawan Kooburat, Aravind Anbudurai, and et al. Twine: A Unified Cluster Management System for Shared Infrastructure. In *Proceedings of USENIX OSDI*, pages 787–803, 2020.
- [52] Muhammad Tirmazi, Adam Barker, Nan Deng, Md E. Haque, Zhijing Gene Qin, Steven Hand, Mor Harchol-Balter, and John Wilkes. Borg: the Next Generation. In *Proceedings of ACM EuroSys*, pages 1–14, 2020.
- [53] Transaction Processing Performance Council (TPC). TPC-A. <http://tpc.org/tpca/default5.asp>.
- [54] Volt Active Data. VoltDB. <https://www.voltdb.com/>.
- [55] Chenxi Wang, Haoran Ma, Shi Liu, Yuanqi Li, Zhenyuan Ruan, Khanh Nguyen, Michael D. Bond, Ravi Netravali, Miryung Kim, and Guoqing Harry Xu. Semeru: A Memory-Disaggregated Managed Runtime. In *Proceedings of USENIX OSDI*, pages 261–280, 2020.
- [56] Xingda Wei, Zhiyuan Dong, Rong Chen, and Haibo Chen. Deconstructing RDMA-Enabled Distributed Transactions: Hybrid Is Better! In *Proceedings of USENIX OSDI*, pages 233–251, 2018.
- [57] Jie You, Jingfeng Wu, Xin Jin, and Mosharaf Chowdhury. Ship Compute or Ship Data? Why Not Both? In *Proceedings of USENIX NSDI*, pages 633–651, 2021.
- [58] Zhe Zhang, Amey Deshpande, Xiaosong Ma, Eno Thereska, and Dushyanth Narayanan. Does Erasure Coding Have a Role to Play in My Data Center? *Microsoft Research Technical Report*, 2010.

[This page intentionally left blank.]

Evolvable Network Telemetry at Facebook

Yang Zhou[†] Ying Zhang[‡] Minlan Yu[†] Guangyu Wang[‡] Dexter Cao[‡] Eric Sung[‡] Starsky Wong[‡]
[†]*Harvard University* [‡]*Facebook*

Abstract

Network telemetry is essential for service availability and performance in large-scale production environments. While there is recent advent in novel measurement primitives and algorithms for network telemetry, a challenge that is not well studied is *Change*. Facebook runs fast-evolving networks to adapt to varying application requirements. Changes occur not only in the data collection and processing stages but also when interpreted and consumed by applications. In this paper, we present PCAT, a production change-aware telemetry system that handles changes in fast-evolving networks. We propose to use a change cube abstraction to systematically track changes, and an intent-based layering design to confine and track changes. By sharing our experiences with PCAT, we bring a new aspect to the monitoring research area: improving the adaptivity and evolvability of network telemetry.

1 Introduction

Network telemetry is an integral component in modern, large-scale network management software suites. It provides visibility to fuel all other applications for operation and control. At Facebook, we built a telemetry system that has been the cornerstone for continuous monitoring of our production networks over a decade. It collects device-level data and events from hundreds of thousands of heterogeneous devices, millions of device interfaces, and billions of counters, covering IP and optical equipments in datacenter, backbone and edge networks. In addition to data retrieval, our telemetry system performs device-level and network-wide processing that generates time-series data streams and derives real-time states. The system serves a wide range of applications such as alerting, failure troubleshooting, configuration verification, traffic engineering, performance diagnosis, and asset tracking.

While our telemetry system can adopt algorithm and system proposals from the research community (e.g., [18, 27, 48, 50]), a remaining open challenge is *Change*. Changes happen frequently in our network hardware and software to meet the soaring application demands and traffic growth [16]. These changes have a significant impact on the network telemetry system. First, we have to collect data on increasingly heterogeneous devices. This is exaggerated as we introduce in-house built FBOSS [13], which allows switches to update as frequently as software. Second, we have growing applications (e.g., [1]) that rely on real-time, comprehensive, and accurate data from network telemetry systems. These applications introduce diverse and changing requirements for the telemetry system on the types of data they need, data collection

frequency, and the reliability and performance of collection methods.

The changes this paper considers include not only the network events from the monitored data, but also those updates to the telemetry system itself: modification to monitoring intent, advance of device APIs, adjustment of frequency configurations, mutation of processing, and restructure of storage formats. Without explicitly tracking them in our network telemetry system, we struggle to mitigate their impact to network reliability. For example, a switch vendor may change a packet counter format when it upgrades a switch version without notifying Facebook operators. This format change implicitly affects many counters in our telemetry database (e.g., aggregated packet counters), leading to adverse impact to downstream alerting systems and traffic engineering decisions. This example highlights several challenges: (1) Production telemetry is a complex system with many components (e.g., data collection, normalization, aggregation) from many teams (e.g., vendors, data processing team, database team, application teams). A change at one component can lead to many changes or even errors at other components. As a result, when telemetry data changes, it is difficult to discern legitimate data changes from semantic changes. (2) Sometimes, we only detect the error passively when traffic engineering team notices congestion. Yet, we cannot diagnose it easily because the error involves many data. Even worse, it may only affect a small portion of vendor devices due to phased updates. Section 2 shares more such examples.

In this paper, we propose to treat changes as first-class citizens by introducing PCAT, a Production *Change-Aware* Telemetry system. PCAT includes three key designs:

First, inspired by the database community [8], we introduce the *change cube* abstraction for telemetry to explicitly track the time, entities, property, and components for each change, and a set of primitives to explore changes systematically. Using change cubes and their primitives, we conduct the first comprehensive study on change characterization in a production telemetry system (Section 3). Our results uncover the magnitudes and the diversity of changes in production, which can be used for future telemetry and reliability research.

Second, we re-architect our telemetry system to be change-aware and evolvable. In the first version of our telemetry system, we have to modify configurations and code at many devices every time a vendor changes the counter semantics or collection methods, or an application changes monitoring intents. To constrain the impact of changes, i.e., the number of affected components, PCAT includes an intent-based lay-

ering design (Section 4) which separates monitoring intents from data collection and supports change cubes across layers. PCAT enables change attribution by allowing network engineers with rich network domain knowledge to define intents while having software engineers building distributed data collection infrastructure with high reliability and scalability. PCAT then compiles intents to vendor-agnostic intermediate representation (IR) data model, and subsequently to vendor-specific collection models, and job models. The intent-driven layering design reduces the number of cascading changes by 54%-100%, and enables systematically tracking changes through the monitoring process.

Third, we build several change-aware applications that explore the dependencies across change cubes to improve application efficiency and accuracy. For example, Toposyncer is our *topology derivation* service that builds on telemetry data and serves many other applications. We transformed Toposyncer to subscribe to change cubes based on derivation dependencies and greatly reduce topology derivation delay by up to 118s. We leverage correlation dependencies across change cubes to enable troubleshooting and validation.

The main contribution of this paper is to bring the community's attention to a new aspect of telemetry systems—how to adapt to changes from network devices, configurations, and applications. We also share our experiences of building change-aware telemetry systems and applications that can be useful to other fast-evolving systems.

2 Motivation

To keep up with new application requirements and traffic growth, data center networks are constantly evolving [16]. As a result, changes happen frequently across all the components in telemetry systems, ranging from device-level changes, collection configuration changes, to changes in the applications that consume telemetry data.

Our first generation of production telemetry system was not built to systematically track changes. This brings significant challenges for telemetry data collection at devices, integration of telemetry system components, debugging network incidents, and building efficient applications. In this section, we share our experiences of dealing with changes in our telemetry system and discuss the system design and operational challenges for tracking changes.

2.1 Bringing changes to first-class citizens

We motivate the needs of treating changes as first-class citizen in network telemetry with a few examples.

1. Build trustful telemetry data. Many management applications rely on telemetry data to make decisions. However, in production, telemetry data is always erroneous, incomplete, or inconsistent due to frequent changes of devices and configurations. Moreover, there are constant failures in large-scale networks (e.g., network connection issues, device overload, message loss, system instability). Therefore, applications need

to know which time range and data source are trustful and how to interpret and use the data. This requires tracking changes for each telemetry data value and semantics.

For example, we collect device counters at various scopes (e.g., interfaces, queues, linecards, devices, circuits, clusters). These counters may have different semantics with device hardware and software upgrades or network re-configurations. For example, we have a counter for 90th percentile CPU usage within a time window of a switch. When we change the switch architecture to multiple subswitches [13], we set the counter as the average of 90th percentile CPU of subswitches. However, our alert on this counter cannot catch single sub-switch CPU spikes that caused bursty packet drops. We need to know when to change the alerts based on counter changes.

2. Track API changes across telemetry components. Our telemetry system consists of multiple data processing components, which are independently developed by different vendors and teams. When one component changes its interfaces, many other components may get affected without notice. There are no principled ways to handle such changes across telemetry components. For example, vendor-proprietary monitoring interfaces often get changed without an explicit notification or detailed specification. This is because telemetry interfaces are traditionally viewed as secondary compared to other major features. However today cloud providers heavily rely on telemetry data for decisions in a fine-grained and continuous manner. If we do not update data processing logic based on device-level changes, the inconsistency may cause bugs and monitoring service exceptions.

In one incident, a routing controller had a problem of unbalanced traffic distribution, caused by incomplete input topology: a number of circuits were missing from the derived topology. This took the routing team and the topology team over three days to diagnose. The root cause was an earlier switch software upgrade that changed the linecard version from integer (e.g., 3) to string (e.g., 3.0.0). Such a simple format change was not compatible with the post-processing code that aggregated the linecard information into a topology. Thus, we missed several linecards in the topology, which then mislead TE decision and cause congestion in the network. This is not a one-off case, given many vendors and software versions coexist in our continuously evolving networks.

3. Debug with change-aware data correlation. As telemetry components keep evolving, it is hard to attribute a problem to a change using data correlation without explicitly tracking changes and their impacts. For example, when we fail to get a counter, the problem can come from data collection at the device, the network transfer, or both.

In production, we make changes in small phases: first canary on a few devices serving non-critical applications, then gradually on more devices to minimize disruptions to the network [13]. In one incident, there were a small number of devices with “empty data” errors for a power counter. The errors increased gradually and ultimately went beyond 1% threshold

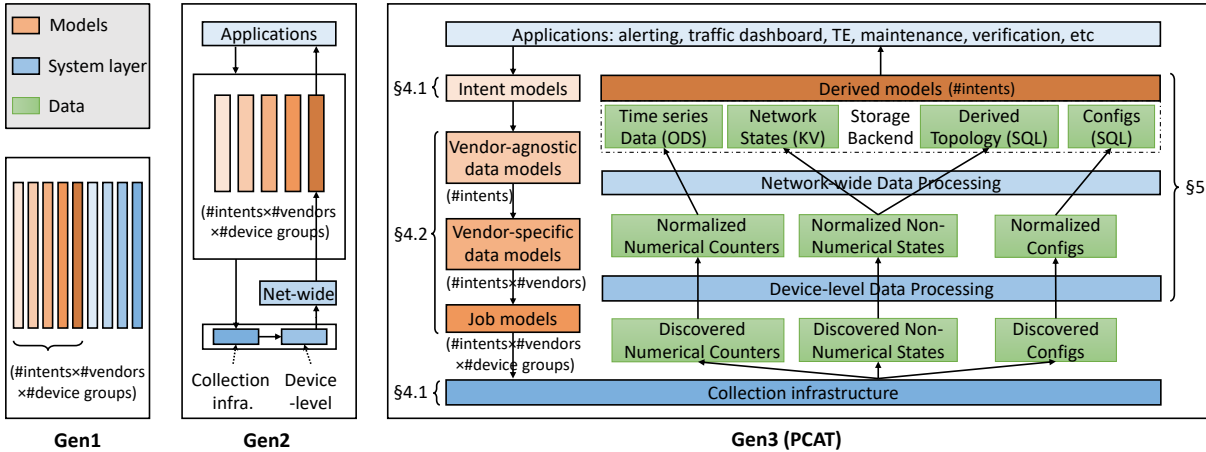


Figure 1: Generations towards change-aware telemetry.

after two weeks and triggered an alarm. This problem was difficult to troubleshoot due to its small percentage. We manually explored the changes through correlation: checking whether there were code changes before the failure, whether the failed devices shared a common region (indicating regional failure), a common vendor, or on common data types. We tried many dimensions of correlation and finally found the errors were mostly related to power and environment counters. The root cause was a vendor changing its format but the processing code could not recognize it. This example shows a tedious manual process of data correlation to debug problems because of gradual change rollouts. To improve debugging, we need to use changes to guide data correlation.

2.2 Lessons from Previous Generations

We now discuss our previous two generations of telemetry systems prior to PCAT and their limitations in handling changes. **Gen1: Monolithic collection script.** In a nutshell, a telemetry system is a piece of code that collects data using APIs from the devices. Our first generation is naturally a giant script that codifies what counters to collect. It hardcodes the collection method, polling frequency, post-processing logic, and where to store the data. Figure 1 illustrates Gen1 as intertwined models and system blocks. It runs as multiple cron jobs, each collecting data from different groups of devices. This design is intuitive to implement but is not change-friendly. If a vendor changes the format of a counter, we need to sweep through the entire script to change the processing logic accordingly, and redeploy the new code to all monitors. It has high maintenance burdens as it relies on expert’s deep understanding of the code to make changes. Further, tracking changes relies on version control system in the form of code differences, which do not reveal the intent directly.

Gen2: Decoupled counter definition from collection process. As our network expanded, the hulking script in Gen1 became hard to manage. We moved to Gen2, which separates the monitoring model (i.e., what counter to collect) from the actual collection code, shown as orange and blue boxes in Fig-

ure 1. The separation allows us to track changes to data types separately from the collection logic (e.g., sending requests, handling connections). However, the intent is still mingled with the vendor-specific counter definition. For instance, one may want to collect the “packet drops per interface”. One needs to specify the exact SNMP MIB entry name and the specific API command. A low-level format change would result in updates on all model definitions. Moreover, the data collection system includes both the collection infrastructure and data processing logic. The data processing logics scatter across many places, e.g., when the data is collected locally at the collector, or before it is put into the storage. To change a piece of processing logic, we have to change many such places, which is cumbersome to track. In addition, when a piece of data is changed or is absent, tracing back on what causes the change is manual and tedious.

2.3 Challenges and PCAT Overview

Our experiences of previous two generations indicate three main challenges in handling changes: change abstraction, attribution, and exploration. To address these challenges, we build our Gen3 telemetry system – PCAT.

Change abstraction. In Gen1 and Gen2, changes were not stored structurally. They exist either as diffs in code reviews in Gen1, or logs to temporary files in Gen2. Without a uniform representation, each application needs to develop ad-hoc scripts to parse each data source. This leads to not only duplicate efforts but also missing changes or mis-interpretations. A uniform and generic change abstraction allows hundreds of engineers to publish and subscribe to changes to boost reliable collaboration without massive coordination overhead. In §3, we propose a generic abstraction called *change cube* to tackle this challenge.

Change attribution. The second challenge is the turmoil to ascribe the intent of the scattering changes. The solution involves a surgically architectural change to a multi-layer design, shown in Figure 1 and elaborated below.

Data collection. The first step is to collect data from de-

vices, called *discovered data*. There are three types: numeric counters, non-numeric states, and configurations (see Table 4 in Appendix). We use different protocols for collecting different data and for different devices: SNMP [10], XML, CLI, Syslog [28], and Thrift for our in-house switches [13].

Device-level data processing (normalization). The data is different in formats and semantics across devices, vendors, and switch software. This makes it difficult for applications to parse and aggregate the data from different devices. We use a device-level data processing layer to parse the raw data to a unified format across devices, vendors, and switch software.

Network-wide data processing. Next, we aggregate device-level (normalized) counters, states, and configurations into network-wide storage systems for applications to query. The normalized non-numerical states (as network states) are stored in a key-value store. We build a tool called **Toposyncer** which constructs *derived topology* from normalized non-numerical states. For example, from per-device data, we can construct the device, its chassis, linecard, as well as cross-device links.

Data consumer applications. There are many critical network applications that consume PCAT data. Network health monitoring and failure detection use monitoring data to detect and react to faults. Network control relies on real-time data for making routing and load balancing decisions [2, 38]. Maintenance and verification use telemetry data to compare network states before and after any network operations.

There are several advantages of the new design compared to previous generations. First, compared to Gen2, Gen3 dissects a monolithic data definition into three different types, each focusing on defining one aspect of the monitoring. The separation brings better scalability and manageability. We describe the details in §4. Second, we not only care about tracking changes in data format and code, but also need to attribute changes to the right teams (i.e., who/what authored the change). Change attribution builds the trust of the data for applications. It facilitates collaboration across teams towards transparent and verifiable system development. Gen3's intent-based layering design lets each team play by their strength and work together seamlessly. Specifically, the network engineers can leverage their rich domain knowledge and focus on intent definition, while software engineers focus on scaling the distributed collection system.

Change exploration. Many designs and operations require a clear understanding of the relations amongst changes. For example, to debug why a piece of data is missing, we always find the last time the data appears and check what has changed since then. We may find one change to be the cause, which could be caused by another change somewhere else. Similarly, when receiving a change of an interface state, we need to reflect the change on the derived topology and upper-layer applications. It motivates us to develop primitives for change exploration that serves many applications. We demonstrate the usage in real-time topology derivation in §5.

3 Changes in Facebook Network Telemetry

In this section, we define the change cube concept and explain how they are generated in this system, together with extensive measurement results by composing queries on top of the change cubes.

3.1 Change Cube Definition

To systematically handle changes in network telemetry, we leverage the concept of *change cubes*. Change cubes are used in databases [8] to tackle the data change exploration problem by efficiently identifying, quantifying, and summarizing changes in data values, aggregation, and schemas. Change cube defines a set of schemas for changes and provides a set of query primitives. However, changes in network telemetry are different from those in databases in two aspects: (1) Network telemetry generates streaming data with constant value changes, so the change cubes in network telemetry do not care about value changes but only changes in schema and data aggregation. (2) Network telemetry has frequent changes due to fast advances of hardware and software that result in data semantics changes.

Change cubes. We define a change cube to be a tuple $\langle \text{Time}, \text{Entity}, \text{Property}, \text{Type}, \text{Dependency} \rangle$. We summarize each field of the change cube in Table 1 and explain below.

- *Time* dimension captures when the change happens. It depends on the granularity we detect changes, e.g., seconds, minutes, or days.
- *Entity* represents a measurement object, e.g. a switch, a linecard, as well as the models that describe what to measure and how.
- *Property* contains the fields or attributes of the entity that get changed. For example, a loopback IP address of a switch, an ingress packet drop of an interface.
- *Layer* dictates the layer or component in the telemetry system (in Figure 1) where changes happen. We discuss how we land in these choices in §4.
- *Dependency* dimension contains a list of other changes that this change is correlated with. Each item in the list is a $\langle \text{ChangeCube}, \text{Dependency Type} \rangle$ pair. We support two dependency types: correlation dependency and derivation dependency. Derivation dependency means that a lower-layer change causes an upper-layer change. Correlation dependency means two changes on correlated entities or properties.

Primitives on change cubes. Next, we introduce the operators on the change cube, which are used to explore the change sets. We leverage the operators proposed in [8] but redefine and expand them in the context of telemetry systems.

- $\text{Sort}_f(C)$ applies function f to a set of change cubes C , on one or a few dimensions to a comparable value, and uses it to generate an ordered list of C . In our problem, sort is mainly used with time to focus on the most recent changes.
- $\text{Slice}_p(C)$ means selecting a subset of C where the predicate

Dimension	Sub category	Examples
Time	Multiple time granularities	Second, minute, hour, day
Entity	Intent model	High-level intent, e.g., packet drops at spine switches
	Vendor-agnostic data model	Counter scope, unit
	Vendor-specific data model	Format, API
	Job model	Collection channel, frequency, protocol
Property	Derived model	Derived network switch
	Model fields	IP address, network type
Layer	Change attributes	LoC, reason
	Application	Adding alert to detect a new failure type
	Network-wide processing	Topology discovery code logic
	Device-level processing	Normalization rule
Dependency	Collection infrastructure	Codebase for collection tasks
	Correlation dependency	BGP session and interface status
	Derivation dependency	Circuit is derived from two interfaces' data

Table 1: Change Cube Definition.

- p is true. It is used to filter an entity or a property value.
- $Split_a(C)$ partitions C to multiple subsets by attribute a . An example is to split the changes by the layer to group changes according to where they occur. A reverse operator to $Split$ is $Union$, which combines multiple change sets.
 - $Rank_f(P_C)$ After we split C to multiple sets P_C , we further analyze these sets and rank them based on a function, e.g., cube size, the time span, the volatility.
 - $TraceUp(c)$ and $TraceDown(c)$. These two operators are used with the *Dependency* field, which are new compared to [8]. The former traces the changes that the current change c depends on, and the latter traces the changes that depend on the c . They are useful for debugging through layers and validation across data.

Explicitly tracking changes in a structured representation eases the diagnosis process. Considering the second example in §2.1, when the switch software is updated, it populates a change cube to the database, indicating the API's return result has changed. Consequently, it triggers another change cube at the counter model level on this specific CPU counter. This change cube in turn propagates through the monitoring stack to job changes and retrieved data changes. The applications using the CPU counters can subscribe to such data change, which can then be notified immediately. The chain of change notifications eliminates the post-mortem debugging after the counter change causes application errors.

3.2 Changes in PCAT

Leveraging *change cubes*, we provide the first systematic study of changes and their impact on network telemetry systems. We populate change cubes of PCAT using multiple ways. For the data stored in database, we leverage our database change pub/sub infrastructure [39]. We subscribe to the telemetry objects' change log and translate them to change cubes. For code changes in collection infrastructure, data processing logics (both device-level and network-wide), and

Queries	Formulas
Q1 (Fig. 2a)	$Sort_{Time(Week)}(Slice_{layer="application"}(C))$
Q2 (Fig. 2a)	$Sort_{Time(Week)}(Slice_{entity="vendor-agnostic data model"}(C))$
Q3 (Fig. 2c)	$\sum_c c.LoC, c \in Split_{Time(Week)}(Slice_{layer="application"}(C))$
Q4 (Fig. 3a)	$Sort_{Time(Day)}(Slice_{entity="job model" \& property="frequency"}(C))$
Q5 (Fig. 3b)	$Split_{network type}(Slice_{entity="job model" \& property="frequency"}(C))$
Q6 (Fig. 3c)	$Split_{network type}(Slice_{entity="job model" \& property="channel"}(C))$
Q7 (Fig. 4a)	$Split_{blueprint type}(Slice_{layer="application" \& reason="blueprint"}(C))$
Q8 (Fig. 4b)	$Split_{vendor}(Slice_{layer="application" \& reason="new model"}(C))$

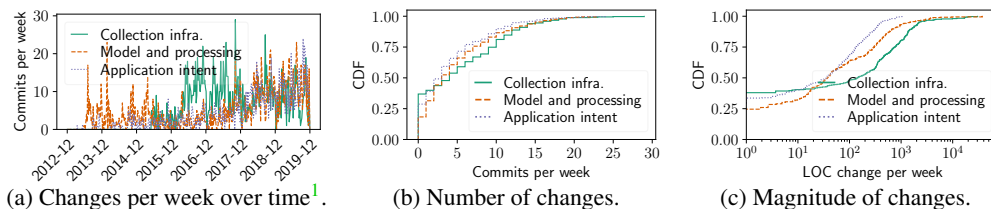
Table 2: Queries used in §3.2.

applications, we parse the logs in the code version control system to generate change cubes. Intent model, data model (both vendor-agnostic and -specific), and job model changes are codified and thus tracked through code changes [41]. They can be populated using the same way as other code changes. We store all change cubes to a separate database called *ChangeDB* and develop APIs to explore these changes.

We analyze changes from the perspectives of devices, collection configurations, and application intents, over seven years (2012-2019). Our results below uncover surprisingly frequent changes and quantify the diverse causes of changes.

3.2.1 Change Overview

Change frequency. We first quantify the code changes of our monitoring system. We map one code commit to one change cube, involving multiple lines of code across multiple files. We group the changes into three categories according to where they happen in Figure 1: collection infrastructure (bottom layer), data & job models and processing (middle two layers), and applications, representing the infrastructure, data, and intent respectively. We construct queries using the primitives defined earlier. We put the actual query to generate the figures in Table 2. Q1 uses *Slice* to filter the changes in application layer, and sorts the changes by time. We replace the “application” with other values for changes in other layers. Q1



(a) Changes per week over time¹.

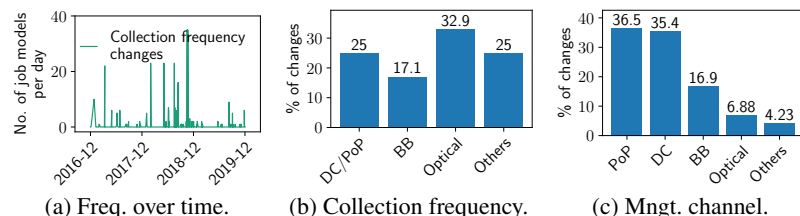
(b) Number of changes.

(c) Magnitude of changes.

Figure 2: Change characteristics.

Change reasons	%
Collection infrastructure	67.9
Adding new devices	17.8
Topology processing	8.30
Data format	4.86
Counter processing	1.19

Table 3: Change categorization by change reasons.

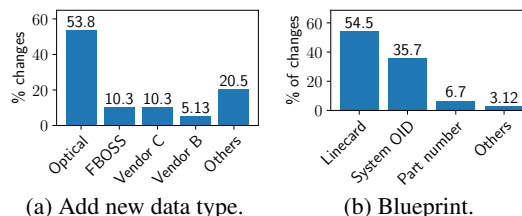


(a) Freq. over time.

(b) Collection frequency.

(c) Mngt. channel.

Figure 3: Collection configuration changes².



(a) Add new data type.

(b) Blueprint.

Figure 4: Network configuration changes.

can be compiled into the following SQL: `SELECT COUNT(*) FROM ChangeDB WHERE layer = "application" GROUP BY time_week ORDER BY time_week.`

Figure 2a shows the number of changes per week. We find that *types of changes vary greatly as the telemetry system scales*. More model and processing changes occur at the beginning (the year of 2013), as we begin by adding more counters to monitor. When the number of counters reaches a certain scale (the year of 2016), we realize the infrastructure needs better scalability. Thus there are more changes to refactor the collection infrastructure. Application intent follows the same trend as data changes, as adding new data is often driven by the need from applications.

Cumulatively across time, we show the average numbers of weekly changes of three categories in Figure 2b. They are on the same order of magnitude, with slightly more infrastructure changes. It can be as high as 25-30 changes per week. Note that each change is deployed on many switches and the changes it introduces to the network is significant.

Change magnitude. We quantify the magnitude of changes in terms of Lines of Code (LoC) using query Q3. While most code changes are not big, some changes could touch multiple lines due to consolidation of processing logic and refactoring. This is obtained by first getting a slice of changes of a given category, splitting the changes into weeks, and summing up the LoC property. Figure 2c shows that *collection infrastructure has larger changes and the application has changes with larger volatility*. We can dig into the volatility of each change set by computing its variance and use the *Rank* primitive. Both figures show there exists a significant number of large changes. For example, there are 27 weeks with more than 1000 LoC changes for collection infrastructure. However, as the industry's trend is to move away from monolithic changes

to many small incremental changes [13], we expect to have more frequent small changes going forward.

Change reason categorization. We analyze the breakdown by reason of change using $Split_{reason}(C)$, which is obtained by parsing the commit log text and adding it to the ChangeDB. Table 3 shows that one major reason is collection infrastructure changes (67.9%). Adding new devices to the network is the second dominant reason (17.8%). Topology processing changes occupy 8.3%. The fourth reason is adjusting the data formats of collection models (4.86%). Lastly, 1.19% come from the device-level counter processing code.

3.2.2 Device-Level Changes

In our large-scale networks, we constantly add new vendors and devices to leverage a rich set of functions and to minimize the risk of single-vendor failures. The number of devices increased 19.0 times and the number of vendors increased 4.7 times as observed by PCAT in six years. Even with the same vendor, we gradually increase the chassis types, which have different combinations of linecard slots and port speeds. More choices of chassis types allow us to have fine-grained customization to our network needs. Furthermore, the number of chassis types grows from 26 to 129 (4.4 \times). In addition, our in-house software switch has tens of code changes daily and deploys once every few days [13].

3.2.3 Collection Configuration Changes

Collection frequency. Applications adjust the collection frequency to balance between data freshness and collection overhead. We first analyze collection changes by counting daily changes of collection frequencies over time, using query Q4. Figure 3a shows that there are constant collection frequency changes over time, with more frequent changes near December 2018 – because of tuning collection frequencies for newly-added optical devices. We analyze collection frequency changes by applying *Slice* on both the entity and the property. Interestingly, Figure 3b shows that *optical devices*

¹The collection and processing infrastructure were not merged into the codebase before 2015-04; so its commits are non-trackable before that.

²The “Other” contains some changes that are hard to classify programmatically. The same applies to Figure 4.

change frequency more often (32.9%) because they cannot sustain high-frequency data polling and thus require more careful frequency tuning.

Management channel changes. PCAT collects data from the management interfaces at devices. As our management network evolves, we frequently reconfigure management interfaces (e.g., IP addresses, in-band vs. out-of-band interfaces). Backbone and PoP devices have multiple out-of-band network choices for high failure resiliency. Figure 3c breaks the IP preference changes into PoPs, DCs, and Backbones. *PoPs have more frequent channel changes (36.5%)* because PoPs are in remote locations and thus have more variant network conditions. Selecting the right channel is important to keep the device reachable during network outages so that we can mitigate the impact quickly.

3.2.4 Application Intent Changes

Data type changes. PCAT supports an increasing number of diverse applications over years. Applications may add new types of data to collect (e.g., to debug new types of failures), or remove some outdated data. Figure 4a shows how different vendors add new data types. Optical device vendors add more data (53.8%) because we recently start building our own optical management software and thus need more counter types. Indeed, optical devices generally have more types of low-level telemetry data compared to IP devices, e.g., power levels, signal-to-noise ratio. They are also less uniformed across vendors than IP devices.

Hardware blueprint changes. Hardware blueprint specifies the internal components (chassis, linecards) of each switch and determines what data to collect. Figure 4b shows the percentage of changes for hardware blueprints such as linecard map, system Object Identifier (OID) map, part number map, and others (e.g., OS regex map). These changes are due to network operations such as device retrofit and migration. They may cause data misinterpretation if not treated carefully.

4 Change Tracking in Telemetry System

In this section, we describe the layering design of our current intent-based telemetry system to help track changes.

4.1 Towards change-aware telemetry

Intent modeling. We use a thrift-based modeling language that empowers network engineers to easily specify their monitoring intents. Compared to other intent language proposed in academia [19, 31, 32], our language puts more emphasis on device state in addition to traffic flows, and defines actions in addition to monitoring. Our language contains three components shown in Figure 5.

- *Scope* captures both the device-level scope (e.g., Backbone Router) and network-level scope, (e.g., DC fabric network).
- *Monitor* specifies what to monitor in a vendor-agnostic way. For example, an intent could be capturing packet discard for the gold-level traffic class, which will get translated

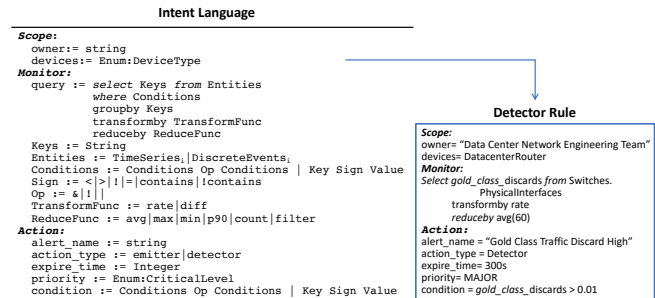


Figure 5: Intent model.

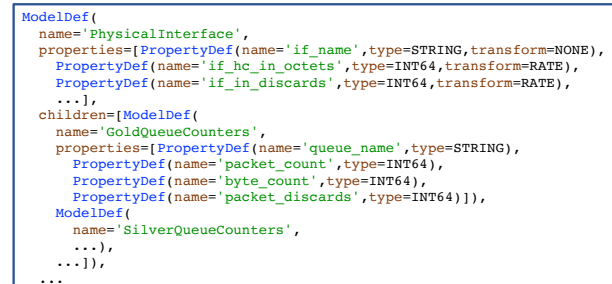


Figure 6: Data model.

to a specific SNMP MIB entry or particular counters. In the left part of Figure 5, we describe the SQL-like query language. The *keys* are monitored metrics and the *entities* are time-series data streams and discrete events tables. We also support data aggregation functions such as *avg*, *count*, *filter*, which aggregate samples over time and devices.

- *Action* includes two types: Emitter and Detector. Emitters subscribe to *discrete network events* that are pushed from devices, and define actions upon receiving these events. Detectors allow us to write formulas for various time-series data, and set up a threshold for the formula value as the alerting condition. A detector example is shown in the right part of Figure 5; it defines a detector based on the key *gold_class_discards* which captures the packet drops for gold-class traffic on a physical interface. The discard is transformed to rate, and aggregated every 60 seconds. The alert is triggered if the threshold is greater than 0.01.

The intent model hides low-level changes. Vendors may change the queue drop counter names, or the mapping between queues and gold-class traffic may change. The intent configuration remains unchanged in both cases.

Runtime system. We handle heterogeneous intents with homogeneous software infrastructure. Thanks to separation, software engineers can focus on the runtime execution system to solve the hard system building problems: scheduling, load balancing, scalability, and reliability. The runtime execution system collects data from devices according to the model, which includes a distributed set of engines and a centralized controller to distribute jobs and collect data from these engines. The centralized controller fetches the latest collection and job models, combines with device information in our database, and generates a sequence of jobs to be executed



Figure 7: Collection method and job model.

with a given deadline. It dispatches jobs to designated engines based on load and latency. The engine executes the collection command, performs device-level processing, and sends data back to the corresponding storage. The system is heavily engineered to tackle the reliability and scalability challenges.

4.2 Change reduction w/ vendor-agnostic IR

Next, we zoom into the intermediate layer between the intent and the runtime. The high-level monitoring intent is translated to the intermediate representation data model, which gets mapped to the vendor-dependent collection model, and finally is materialized to the job model on each device. We emphasize that how the modular design principle is translated to different models in order to *limit the impact of changes*.

Vendor-agnostic intermediate representation (IR) data model. The data model is created based on the *keys* field in the intent model. It specifies data schema in the following way, as shown in Figure 6.

- *Hierarchical.* We choose a tree structure as an intermediate representation, called the *model tree*. An example is shown in Figure 7. An *AggregateInterface* model has multiple child models, e.g., *PhysicalInterface*, *BGPSessions*. A *PhysicalInterface* also has multiple child models. Models are like templates waiting to be filled in. When they are materialized by actual monitoring data, we call them **objects**. By organizing the materialized objects in the same hierarchy as the model tree and adding a dummy root to connect up the top-level objects, we get a materialized **object tree**. The models define the data to be collected, which is derived from the *keys* field in the intent model.
- *Typed.* The data model defines the types of data to make interpretation of the data easier, e.g., *if_hc_in_octets* is the incoming traffic in octets.
- *Processing instruction.* It also defines basic processing primitives to go with the data using the *transform* field, e.g., computing a per-second rate from consecutive absolute counts. Both the type and the processing instruction are determined by the intent. Placing all the processing logic in a separated blob makes it much easier to track the changes

in processing logic.

Vendor-dependent collection model. The IR model is further compiled down to vendor-specific counter names, specific commands to use, e.g., a CLI command, Thrift function name. Figure 7 shows two collection methods for the *GoldQueueCounters* data model: CLI and thrift. In each implementation, we define the collection API and the post-processing function in the *parser* field. We show an example of the CLI parser function that matches the regex in the output of a command on the vendor1 device. Creating this layer of model separately allows us a place to capture all changes due to vendor format and API changes, which are quite common.

Vendor-dependent job model. The job model combines the collection method with a concrete set of devices, shown in Figure 7. The *implementation* field matches with what is defined in the collection method. Instead of defining a job spec for each device, we group devices and apply the same job spec for all of them. Figure 7 uses *DeviceFilter* to define device role (e.g., rack switches), OS type, region, device state, etc. Job models are the input to the runtime execution to handle job scheduling and manage job completion. Job model captures the system aspects of changes. It can be adapted according to performance and scalability requirement, which is independently controlled from the intent or data specifications.

5 Change Exploration

Once PCAT collects data based on monitoring intents, we run device-level and network-level processing to report the data back to applications. Below, we build a few change-aware applications by exploring dependencies across change cubes.

5.1 Change-driven Topology Derivation

Toposyncer is our topology generation service, part of the collection infrastructure (see network-wide data processing in Figure 1). It creates *derived topology* from normalized device-level data (i.e., in vendor-agnostic format) (Figure 8). For example, from per-device data (e.g., interface counters, BGP sessions), *Toposyncer* constructs the device, its chassis, linecard, as well as cross-device links.

Toposyncer overview *Toposyncer* has four processes: (1) *Sync_device* constructs nodes with multiple sub-components: sub-switches, chassis, line cards (line 2-11 in Figure 9). It also derives device-level attributes such as power and temperature, control and management plane settings. (2) *Sync_port* derives physical and logical interfaces on each node and their settings (IP address, speed, QoS) (line 12-13). (3) *Sync_circuit* constructs cross-device circuits. A circuit is modeled as an entity with two endpoints, pointing to the interface of each end's router [41]. For each interface, it searches for all possible neighbors based on various protocol data, e.g., LLDP, MAC table, LACP table. In case some data source is incomplete, we search all data sources, independently identify all possible neighbors from each data source, and consolidate the results.

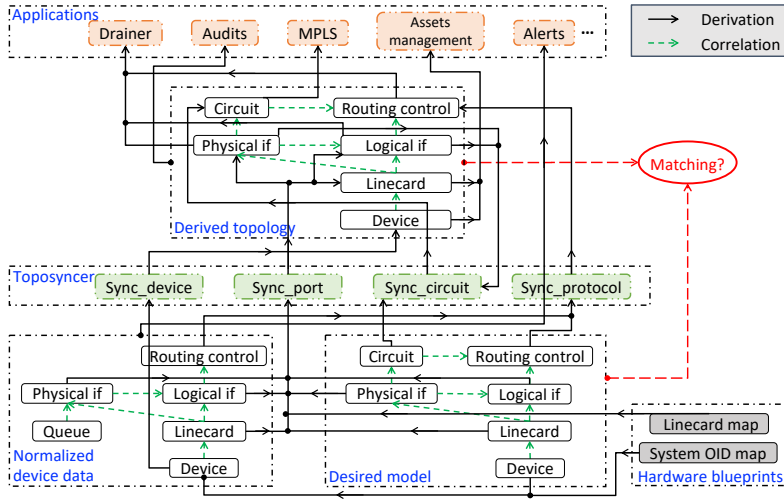


Figure 8: Data dependency graph. Toposyncer consumes the device data, desired model and hardware blueprints to generate derived data. PCAT verifies the derived data with the desired model.

(4) *Sync_protocol* creates the protocol layers on top of the circuits, such as OSPF areas, BGP sessions and their states.

Toposyncer uses two additional data sources as templates to guide the construction: the desired model which defines the operator’s intent topology and hardware blueprints which include hardware specifications, as shown at the bottom in Figure 8. Figure 9 shows the process. It uses desired device data (names, IP addresses) to decide what device to derive (line 2). Then, it uses the hardware blueprints and desired data to handle ambiguity. For instance, to figure out “what is this chassis”, it first checks the discovered chassis name in raw data from the device. But often the discovered name is not uniquely mapped to a chassis but to several possible chassis versions, e.g., two versions with 4 linecards, one version with 8 linecards. Toposyncer cross-checks with hardware blueprints and picks the best match³ (line 6-8). This process is similar to other topology services [29, 41], but we focus on derived models and how we populate them automatically from telemetry data.

Improve Toposyncer with change cubes. Our first implementation of Toposyncer did not utilize changes. It ran periodically against the latest snapshots of collected data at a fixed frequency (e.g., 15 minutes). This method leads to stale derived data, which affects the freshness and accuracy of upper-layer applications. Another challenge is debugging. When a piece of data (e.g., a circuit) is missing in the derived topology, it is hard to find out whether it is because of a raw data change, a normalized data change, a desired model change, a hardware blueprint change, or other reasons. We tackle these problems using change cubes and the dependency primitives below.

³When the guess is wrong, it exhibits as a discrepancy between desired and derived topology. We add alarming to detect such differences and involve humans to manually investigate.

```

1: procedure DERIVETOPOLOGY(Collection, Desired,
   HdwTemps, DependencyG)
2:   for  $d \in \text{Desired.Devices}$  do
3:     DeviceObj, dep = sync_device(Collection, d)
4:     DependencyG.add(dep)
5:     blueprint = getBlueprint(HdwTemps, d)
6:     for  $\text{chassis\_temp} \in \text{blueprint}$  do
7:       derived_chassis = sync_chassis(Collection,
         chassis_temp)
8:       for  $\text{linecard\_model} \in \text{chassis\_temp}$  do
9:         derived_chassis.add(sync_linecard(Collection,
           linecard_model))
10:      DeviceObj.addchassis(derived_chassis)
11:      DeviceObj.add(DeviceObj)
12:   for  $d \in \text{DeviceObjs}$  do
13:     derived_ifaces = sync_port(Collection, d)
14:     for  $\text{iface} \in \text{derived\_ifaces}$  do
15:       neighbors.add(findNeighbors(iface, Collection))
16:     circuits = sync_circuits(iface, neighbors)
17:     sync_protocol(Collection, circuits)
18: procedure UPDATETOPOLOGY(UpdateQ, DependencyG)
19:   while  $\text{UpdateQ} \neq \emptyset$  do
20:      $\text{change}_i = \text{UpdateQ.pop}()$ 
21:      $\text{dependent\_objs} = \text{change}_i.\text{Dependency}$ 
22:      $\text{update\_func} = \text{findFunc}(\text{dependent\_objs})$ 
23:      $\text{update\_func}(\text{dependent\_objs}, \text{change}_i)$ 

```

Figure 9: Toposyncer algorithm

Build change cubes. We generate change cubes for normalized data, desired model, hardware blueprint, as well as Toposyncer code changes, shown as each dotted box in Figure 8. We generate these cubes by parsing database transaction logs and model/code changes from version control system logs and publish them to ChangeDB. For example, when an operator changes the configuration of an SNMP MIB for a device, we generate a record to the DB.

Derivation dependency. We populate the derivation dependency across change cubes A and B if we derive data A from data B. In the above example, the MIB change will result in multiple change cubes of job models. We build the dependency between the MIB config change and the rest of job model changes. Figure 8 shows derivation dependency in solid arrows across objects in different layers (each large dash box representing a layer). The dependency exists between data objects as well as between code and data.

Subscription to change cubes. Toposyncer subscribes to the change cubes and invokes corresponding processing logic accordingly, shown in line 19-23. For example, *sync_port* subscribes to the device data (e.g., Thrift_Fboss_Linecards), *Snmp_entPhysicalTable*, and a hardware blueprint (i.e., linecard map). If the hardware blueprint changes, i.e., the same linecard name is mapped to a different hardware blueprint, the change cube will be published and *sync_port* triggers its function *sync_phy_iface* function on the impacted interfaces. Similar pub/sub relation is also built between applications and derived data. For instance, as shown in Figure 8, a drainer application subscribes to interface status and the routing control messages to determine if it is safe to perform an interface drain operation.

5.2 Improve Trust on Data Quality

Real-world telemetry data may contain dirty or missing data. By exploring the change history, one can better judge whether the current values are trustworthy. Observing patterns of data changes can help predict the occurrences of future changes and identify missing changes.

Correlation dependency. Amongst normalized device data or derived data, data have relationships between them, shown as dash arrows within each large dash box in Figure 8. The relationship represents the physical dependency across objects, such as “contain”, “connect”, “originate”. Previous topology-modeling works Robotron [41] and MALT [29] focus on the desired model and the correlation dependency in it. The desired model is built for the purpose of capturing topology intents and generating configurations. Here we use the model together with change cubes to verify if the actual topology’s change is legit by comparing it against the desired models. A change cube generated from the desired object should have a matching change cube in the derived object, and vice versa. This can be done with a *Slice* on the entity, *Sort* by time, and compare the *Entity* of the changes.

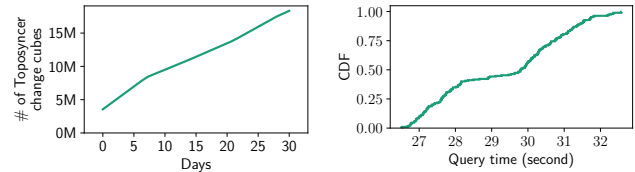
This correlation dependency can also be used for cross-layer validation of data quality. We implement *if-then* validation rules based on the correlation dependency on change cubes. We give two examples below. One use case is hard-stop fault detection. One rule is that *if* the logical interface fails (i.e., a specific change cube on a logical interface), *then* the routing session going through it will also fail (i.e., another change cube on the routing session must exist). If we observe significant errors at the lower layer but no upper failure, it indicates a measurement issue. In another use case, the aggregate interface consists of multiple physical interfaces. *If* a member physical interface reports packet errors, *then* the packet errors from aggregated interface should be larger than or equal to the physical interface errors. If the rule is not satisfied, it indicates some issues. These cross-layer dependencies can help us detect change-induced problems more quickly.

6 Evaluation

This section evaluates how the layer design of PCAT has helped with change tracking and how much benefit the change cube method has brought to use cases.

6.1 Change tracking implementation

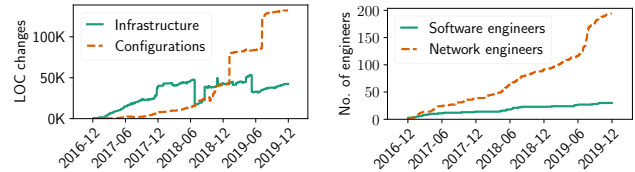
First, we examine whether tracking all the changes is even feasible in a production environment. We show the change cube data volume grows with time in Figure 10a. Drawing from the experiences of Facebook’s data infrastructure team, we employ a two-tier storage solution. We have an in-memory database to hold the change data for the most recent 30 days and have a disk-based SQL database for longer historical data. At the same time, the change data is published to our publish/subscribe system [4] for real-time propagation. Next, we



(a) # of change cubes over time.

(b) Query distribution time.

Figure 10: Scalability and performance of change tracking.



(a) LoC changes over time.

(b) Engineers over time.

Figure 11: Separating configurations with telemetry infra.

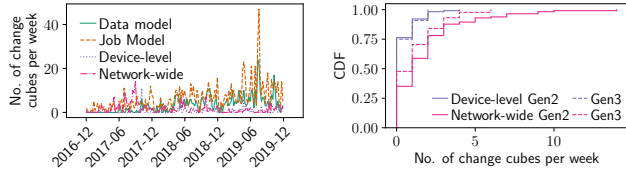
evaluate the performance of exploration using the primitives defined in §3.1, which is implemented using SQL statements. Figure 10b shows the query distribution time for data stored on disk, most of which centers around 27-32 seconds, due to the large data volume. For shorter duration of data in memory, it takes less than one second.

6.2 Benefits of separation

Analyzing change data over time helps us evaluate the long-term benefit of the layer design. We show it from three aspects. **Decoupled evolvement of configurations and infrastructure.** We categorize changes broadly to configuration changes vs. infrastructure changes. We quantify the magnitude of the change using the Lines of Code (LoC) change. Figure 11a shows that the changes for configurations are 3.1 times more than core collection infrastructure changes. The sudden jumps for configurations in January 2019 are due to adding a large set of optical devices, which was not monitored by PCAT. The second increase around July 2019 is due to the migration to Gen3, resulting in a large number of new models added. The result shows that we increase the monitoring scope by configuration layer changes with a stable infrastructure.

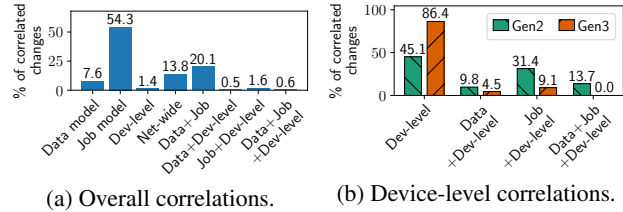
Scaling with divided responsibility. The separation in software systems has a long-term impact on the organization growth and people aspects. In Figure 11b, we analyze the change authors and categorize by their roles. It shows the number of network engineers who have made changes to configurations is increasing at a much faster pace than software engineers, with 7.2 times more people recently. The increase around June 2019 is due to both migration to Gen3 and adding more optical devices to monitor. It is clear that both of these changes are carried out by network engineers. It shows that PCAT enables network engineers to work on different network types while a small number of software engineers maintain infrastructure. It will boost a healthy collaboration environment where each team can play by their strength.

Confining the impact of changes. We use the number of change cubes as an approximate of the volume of changes.



(a) Change cubes across time. (b) CDF of change frequency.

Figure 12: Change cube frequency.



(a) Overall correlations. (b) Device-level correlations.

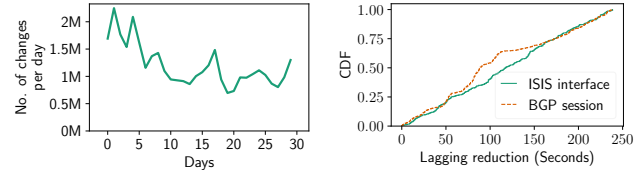
Figure 13: Correlated changes.

Figure 12a shows its trends across time for data models, job models, device-level processing, and network-wide processing. The maximum numbers range from 15 to 50 for different categories. The models (data and job) have more changes due to the frequent intent changes. The infrastructure layers (device and network-wide processing) are more stable. Recently there are more data model and job model changes, because of Gen2-to-Gen3 migration. To directly illustrate the benefit of modular design in Gen3 (§4.2), Figure 12b compares the frequency of change cubes for device-level and network-wide processing in Gen2 and Gen3 (after 2019-02). We observe that the average change frequency for network-wide processing in Gen3 is 38.1% lower than Gen2, while device-level remains similar. This means the modular design in Gen3 further prevents the changes in lower data model and job model layers from impacting upper processing layers, confining the impact of lower-layer changes. Note that we discounted the changes due to Gen2-to-Gen3 migration to have a fair comparison.

Reducing correlated changes. We find change cubes that occur close in time as correlated changes (e.g., data and job models are modified in the same commit). We show that PCAT’s way of separating layers and models has helped reduce correlated changes. We first present the breakdown of different correlation combinations in both generations in Figure 13a. The largest combination is data and job, accounting for 20.1%. It is because adding new devices requires adding both data and job models. There are a small fraction of changes that require updating data, job, and device-level processing all together. Most of them are due to adding some specific counters that require special processing. Figure 13b further breaks down all correlated changes related to device-level processing for Gen2 and Gen3. It shows that Gen3 has significantly reduced the correlated changes by 54.1%, 71.0%, and 100.0% (i.e., the second-to-last bar pairs) accordingly.

6.3 Benefits of change-driven Toposyncer

The first benefit is explicitly tracking changes in a centralized manner. Figure 14a shows the magnitude of the changes over



(a) No. of changes to Toposyncer. (b) Lagging reduction.

Figure 14: Change-driven Toposyncer.

time to Toposyncer. Note that this is much higher than the changes presented earlier, since it includes the changes of raw data for network states.

The second benefit lies in the efficiency and accuracy improvement to applications. We evaluate it using the lagging time, i.e. the time between the change happening and when changes are reflected in derived topology by Toposyncer. Figure 14b shows the topology derivation is much more timely: reducing 118.76s lagging time for ISIS interface updates, and 108.93s for BGP session state updates, averagely.

7 Lessons and Future Directions

We discuss our lessons from building PCAT and the opportunities for future research.

Efficiency vs. adaptivity. We work closely with vendors to improve the efficiency of data collection primitives at switches (similar to academic work on reducing memory usage and collection overhead with high accuracy [24, 26, 46]). However, pursuing efficiency brings us challenges on adaptivity. Different devices have different programming capabilities and resource constraints to adopt efficient algorithms. Introducing new primitives also adds diversity and dynamics to upper layers in the telemetry system. For example, we work with vendors to support a sophisticated micro-burst detection on hardware. However, if only a subset of switches supports this new feature, applications need complex logic to handle detected and missed micro-bursts. Thus we have a higher bar for adopting efficient algorithms due to adaptivity concerns.

To support diverse data collection algorithms, we need a full-stack solution with universal collection interfaces at switches and change-aware data processing and aggregation algorithms. Recent efforts on standardizing switch interfaces such as OpenConfig [3] is a great first step but does not put enough emphasis on standardizing telemetry interfaces. Recent trends on open-box switches (e.g., FBOSS [13]) bring new opportunities to develop adaptive telemetry primitives.

Trustful network telemetry. Telemetry becomes the foundation for many network management applications. Thus we need to know which data at which time period is trustful. However, building a trustful telemetry system is challenging in an evolving environment with many changes of devices, network configurations, and monitoring intents. Fast evolution also introduces more misconfigurations and software bugs. Explicitly tracking change cubes and exploring their dependencies in PCAT is only the first step.

We need more principled approaches for *telemetry verification and validation* across monitoring intents, data models, and collection jobs. Compared with configuration verification work [7, 35], telemetry verification requires quantifying the impact of changes to the measurement results. One opportunity is that we can leverage cross-validations across multiple counters covering the same or related network states or across aggregated statistics over time. For example, we collect power utilizations (watts) from both switches and power distribution units (PDUs). In this way, we can validate the correctness of these utilization counters by comparing the PDU value with the sum of switch values.

Telemetry systems are complex time-series databases. We can leverage provenance techniques [12] to support change tracking, data integration, and troubleshooting. One challenge is that we cannot build a full provenance system due to vendor-proprietary code and network domain-specific data aggregation algorithms. There are also unexpected correlation dependencies across data.

Integration between telemetry and management applications. Our production networks are moving towards self-driving network management with a full measure-control loop. PCAT shows that changes bring a new complexity to the measure-control loop. Control decisions not only affect the network state that telemetry system captures but also the telemetry system itself. For example, an interface change may affect a counter scope. A traffic engineering control change may affect data aggregation because traffic traverses through different switches. These telemetry data changes in turn affect control decisions. We need to identify solutions that can feed control-induced changes directly into the telemetry systems.

Another question is how to present large-scale multi-layer telemetry data to control applications. Rather than providing a unified data stream, control applications can benefit from deciding what time, at what granularity, frequency, and availability level for data collection and the resulting overhead and accuracy in the telemetry system. One lesson we learned is to have the telemetry data available when it is mostly needed. For example, the network's aggregated egress traffic counter, which is collected at the edge PoPs, is a strong indicator of the business healthiness. To ensure its high availability, we need to give control applications the option to transfer the counter on more expensive out-of-band overlay networks. Moreover, we may extend the intent model to explicitly express the reliability-cost tradeoffs and adapt the tradeoffs during changes. We also need new algorithms and systems that can automatically integrate data at different granularities, frequencies, and device scopes to feed in control applications.

8 Related Work

Network evolution. Several existing works have also pointed out the importance of considering changes. Both Robotron [41] and MALT [29] discuss it in the context of topology modeling, but miss the practical challenges of net-

work monitoring. [16] discusses network availability during changes, while we focus on telemetry systems during changes.

Other monitoring techniques. PCAT is a passive approach. Active measurement injects packets into the network [14, 17, 18, 34, 49], and they are complements to passive measurement. The design principle of PCAT to handle changes can be applied to existing monitoring systems [20, 25, 36, 44, 48, 50], languages and compilers [9, 19, 32, 33]. PCAT also benefits from recent software-defined measurement frameworks [25, 27, 32, 46, 48]. For example, similar to OpenSketch [48], PCAT frees network engineers from configuring different measurement tasks manually. PCAT's intent model design borrows ideas from the query language in Marple [32]. There are many memory-efficient monitoring algorithms [22, 23, 26, 30, 46] that focus on the expressiveness and performance of network monitoring. They provide adaptivity but only to a limited type of new queries, resource changes, or network condition changes. Here, PCAT focuses on a broader set of adaptivity (e.g., adaptive to counter semantics changes, data format changes, and more).

Dependency in network management. Dependency graph has been widely used for root cause localization [5, 6, 37, 42, 43, 47, 50]. Statesman [40] captures domain-specific dependencies among network states. We share some similarities but use dependency to tackle the change propagation.

Techniques from database and software engineering. Data provenance [12, 15] encodes causal relations between data and tables in metadata. Several works [11, 45] apply provenance to network diagnosis. [8] proposes the change cube concept and applies it to real-world datasets. All the above works focus on data face-value. On the other hand, software engineering community studies the problem of how a change in one source code propagates to impact other code [21, 51]. Ours looks at changes from telemetry systems from both data, configurations, and code.

9 Conclusion

This paper presents the practical challenge of a monitoring system to support an evolving network in Facebook. We propose explicitly tracking changes with change cubes and exploring changes with a set of primitives. We present extensive measurements to illustrate its prevalence and complexity in production, then share experiences in building a change-aware telemetry system. We hope to inspire more research on adaptive algorithms and evolvable systems in telemetry.

Acknowledgments

We thank our shepherd Chuanxiong Guo and the anonymous reviewers for their insightful comments. Yang Zhou and Minlan Yu are supported in part by NSF grant CNS-1834263.

References

- [1] Express backbone. <https://engineering.fb.com/data-center-engineering/building-express-backbone-facebook-s-new-long-haul-network/>.
- [2] Introducing proxygen facebook c++ http framework. <https://code.fb.com/production-engineering/introducing-proxygen-facebook-s-c-http-framework>.
- [3] OpenConfig YANG model. <http://www.openconfig.net/projects/models/>.
- [4] Scribe. <https://github.com/facebookarchive/scribe>.
- [5] Behnaz Arzani, Selim Ciraci, Luiz Chamon, Yibo Zhu, Hongqiang Harry Liu, Jitu Padhye, Boon Thau Loo, and Geoff Outhred. 007: Democratically finding the cause of packet drops. In *15th {USENIX} Symposium on Networked Systems Design and Implementation ({NSDI} 18)*, pages 419–435, 2018.
- [6] Paramvir Bahl, Ranveer Chandra, Albert Greenberg, Srikanth Kandula, David A Maltz, and Ming Zhang. Towards highly reliable enterprise network services via inference of multi-level dependencies. *ACM SIGCOMM Computer Communication Review*, 37(4):13–24, 2007.
- [7] Ryan Beckett, Aarti Gupta, Ratul Mahajan, and David Walker. A general approach to network configuration verification. In *Proceedings of the Conference of the ACM Special Interest Group on Data Communication*, pages 155–168, 2017.
- [8] Tobias Bleifuß, Leon Bornemann, Theodore Johnson, Dmitri V Kalashnikov, Felix Naumann, and Divesh Srivastava. Exploring change: a new dimension of data analytics. *Proceedings of the VLDB Endowment*, 12(2):85–98, 2018.
- [9] Kevin Borders, Jonathan Springer, and Matthew Burnside. Chimera: A declarative language for streaming network traffic analysis. In *Presented as part of the 21st {USENIX} Security Symposium ({USENIX} Security 12)*, pages 365–379, 2012.
- [10] Jeffrey D Case, Mark Fedor, Martin L Schoffstall, and James Davin. Simple network management protocol (snmp). Technical report, 1990.
- [11] Ang Chen, Yang Wu, Andreas Haeberlen, Wenchao Zhou, and Boon Thau Loo. The good, the bad, and the differences: Better network diagnostics with differential provenance. In *Proceedings of the 2016 ACM SIGCOMM Conference*, pages 115–128. ACM, 2016.
- [12] James Cheney, Laura Chiticariu, Wang-Chiew Tan, et al. Provenance in databases: Why, how, and where. *Foundations and Trends® in Databases*, 1(4):379–474, 2009.
- [13] Sean Choi, Boris Burkov, Alex Eckert, Tian Fang, Saman Kazemkhani, Rob Sherwood, Ying Zhang, and Hongyi Zeng. Fboss: building switch software at scale. In *Proceedings of the 2018 Conference of the ACM Special Interest Group on Data Communication (SIGCOMM)*, pages 342–356. ACM, 2018.
- [14] Cisco. Ip slas configuration guide, cisco ios release 12.4t. <http://www.cisco.com/c/en/us/td/docs/ios-xml/ios/ipsla/configuration/12-4t/sla-12-4t-book.pdf>.
- [15] Mahmoud Elkhodr, Belal Alsinglawi, and Mohammad Alshehri. Data provenance in the internet of things. In *2018 32nd International Conference on Advanced Information Networking and Applications Workshops (WAINA)*, pages 727–731. IEEE, 2018.
- [16] Ramesh Govindan, Ina Minei, Mahesh Kallahalla, Bikash Koley, and Amin Vahdat. Evolve or die: High-availability design principles drawn from googles network infrastructure. In *Proceedings of the 2016 ACM SIGCOMM Conference*, pages 58–72. ACM, 2016.
- [17] Nicolas Guilbaud and Ross Cartlidge. Google localizing packet loss in a large complex network. Nanog57, Feb 2013.
- [18] Chuanxiong Guo, Lihua Yuan, Dong Xiang, Yingnong Dang, Ray Huang, Dave Maltz, Zhaoyi Liu, Vin Wang, Bin Pang, Hua Chen, et al. Pingmesh: A large-scale system for data center network latency measurement and analysis. In *Proceedings of the 2015 ACM Conference on Special Interest Group on Data Communication*, pages 139–152, 2015.
- [19] Arpit Gupta, Rob Harrison, Marco Canini, Nick Feamster, Jennifer Rexford, and Walter Willinger. Sonata: Query-driven streaming network telemetry. In *Proceedings of the 2018 Conference of the ACM Special Interest Group on Data Communication*, pages 357–371, 2018.
- [20] Nikhil Handigol, Brandon Heller, Vimalkumar Jeyakumar, David Mazières, and Nick McKeown. I know what your packet did last hop: Using packet histories to troubleshoot networks. In *NSDI*, volume 14, pages 71–85, 2014.
- [21] Ahmed E Hassan and Richard C Holt. Predicting change propagation in software systems. In *20th IEEE International Conference on Software Maintenance, 2004. Proceedings.*, pages 284–293. IEEE, 2004.

- [22] Qun Huang, Xin Jin, Patrick PC Lee, Runhui Li, Lu Tang, Yi-Chao Chen, and Gong Zhang. Sketchvisor: Robust network measurement for software packet processing. In *Proceedings of the Conference of the ACM Special Interest Group on Data Communication*, pages 113–126, 2017.
- [23] Qun Huang, Patrick PC Lee, and Yungang Bao. Sketchlearn: Relieving user burdens in approximate measurement with automated statistical inference. In *Proceedings of the 2018 Conference of the ACM Special Interest Group on Data Communication*, pages 576–590, 2018.
- [24] Qun Huang, Haifeng Sun, Patrick PC Lee, Wei Bai, Feng Zhu, and Yungang Bao. Omnimon: Re-architecting network telemetry with resource efficiency and full accuracy. In *Proceedings of the Annual conference of the ACM Special Interest Group on Data Communication on the applications, technologies, architectures, and protocols for computer communication*, pages 404–421, 2020.
- [25] Yuliang Li, Rui Miao, Changhoon Kim, and Minlan Yu. Flowradar: A better netflow for data centers. In *13th {USENIX} Symposium on Networked Systems Design and Implementation ({NSDI} 16)*, pages 311–324, 2016.
- [26] Zaoxing Liu, Ran Ben-Basat, Gil Einziger, Yaron Kassner, Vladimir Braverman, Roy Friedman, and Vyas Sekar. Nitrosketch: Robust and general sketch-based monitoring in software switches. In *Proceedings of the ACM Special Interest Group on Data Communication*, pages 334–350. 2019.
- [27] Zaoxing Liu, Antonis Manousis, Gregory Vorsanger, Vyas Sekar, and Vladimir Braverman. One sketch to rule them all: Rethinking network flow monitoring with univmon. In *Proceedings of the 2016 ACM SIGCOMM Conference*, pages 101–114, 2016.
- [28] Chris Lonvick. The bsd syslog protocol. Technical report, 2001.
- [29] Jeffrey C Mogul, Drago Goricanec, Martin Pool, Anees Shaikh, Douglas Turk, Bikash Koley, and Xiaoxue Zhao. Experiences with modeling network topologies at multiple levels of abstraction. In *17th {USENIX} Symposium on Networked Systems Design and Implementation ({NSDI} 20)*, pages 403–418, 2020.
- [30] Masoud Moshref, Minlan Yu, Ramesh Govindan, and Amin Vahdat. Scream: Sketch resource allocation for software-defined measurement. In *Proceedings of the 11th ACM Conference on Emerging Networking Experiments and Technologies*, pages 1–13, 2015.
- [31] Masoud Moshref, Minlan Yu, Ramesh Govindan, and Amin Vahdat. Trumpet: Timely and precise triggers in data centers. In *Proceedings of the 2016 ACM SIGCOMM Conference*, pages 129–143, 2016.
- [32] Srinivas Narayana, Anirudh Sivaraman, Vikram Nathan, Prateesh Goyal, Venkat Arun, Mohammad Alizadeh, Vimalakumar Jeyakumar, and Changhoon Kim. Language-directed hardware design for network performance monitoring. In *Proceedings of the Conference of the ACM Special Interest Group on Data Communication*, pages 85–98, 2017.
- [33] Srinivas Narayana, Mina Tahmasbi, Jennifer Rexford, and David Walker. Compiling path queries. In *13th {USENIX} Symposium on Networked Systems Design and Implementation ({NSDI} 16)*, pages 207–222, 2016.
- [34] Yanghua Peng, Ji Yang, Chuan Wu, Chuanxiong Guo, Chengchen Hu, and Zongpeng Li. detector: a topology-aware monitoring system for data center networks. In *2017 USENIX Annual Technical Conference (USENIX ATC 17)*, pages 55–68. USENIX Association, 2017.
- [35] Santhosh Prabhu, Kuan Yen Chou, Ali Kheradmand, Brighten Godfrey, and Matthew Caesar. Plankton: Scalable network configuration verification through model checking. In *17th {USENIX} Symposium on Networked Systems Design and Implementation ({NSDI} 20)*, pages 953–967, 2020.
- [36] Arjun Roy, Hongyi Zeng, Jasmeet Bagga, George Porter, and Alex C Snoeren. Inside the social network’s (data-center) network. In *ACM SIGCOMM Computer Communication Review*, volume 45, pages 123–137. ACM, 2015.
- [37] Arjun Roy, Hongyi Zeng, Jasmeet Bagga, and Alex C Snoeren. Passive realtime datacenter fault detection and localization. In *14th {USENIX} Symposium on Networked Systems Design and Implementation ({NSDI} 17)*, pages 595–612, 2017.
- [38] Brandon Schlinker, Hyojeong Kim, Timothy Cui, Ethan Katz-Bassett, Harsha V Madhyastha, Italo Cunha, James Quinn, Saif Hasan, Petr Lapukhov, and Hongyi Zeng. Engineering egress with edge fabric: Steering oceans of content to the world. In *Proceedings of the Conference of the ACM Special Interest Group on Data Communication*, pages 418–431. ACM, 2017.
- [39] Yogeshwer Sharma, Philippe Ajoux, Petchean Ang, David Callies, Abhishek Choudhary, Laurent Demailly, Thomas Fersch, Liat Atsmon Guz, Andrzej Kotulski, Sachin Kulkarni, Sanjeev Kumar, Harry Li, Jun Li, Evgeniy Makeev, Kowshik Prakasam, Robbert Van Renesse, Sabyasachi Roy, Pratyush Seth, Yee Jiun Song, Benjamin Wester, Kaushik Veeraraghavan, and Peter

- Xie. Wormhole: Reliable pub-sub to support geo-replicated internet services. In *12th USENIX Symposium on Networked Systems Design and Implementation (NSDI 15)*, Oakland, CA, May 2015. USENIX Association.
- [40] Peng Sun, Ratul Mahajan, Jennifer Rexford, Lihua Yuan, Ming Zhang, and Ahsan Arefin. A network-state management service. In *Proceedings of the 2014 ACM conference on SIGCOMM*, pages 563–574, 2014.
- [41] Yu-Wei Eric Sung, Xiaozheng Tie, Starsky HY Wong, and Hongyi Zeng. Robotron: Top-down network management at facebook scale. In *Proceedings of the 2016 ACM SIGCOMM Conference*, pages 426–439. ACM, 2016.
- [42] Praveen Tammana, Rachit Agarwal, and Myungjin Lee. Simplifying datacenter network debugging with pathdump. In *12th {USENIX} Symposium on Operating Systems Design and Implementation ({OSDI} 16)*, pages 233–248, 2016.
- [43] Cheng Tan, Ze Jin, Chuanxiong Guo, Tianrong Zhang, Haitao Wu, Karl Deng, Dongming Bi, and Dong Xiang. Netbouncer: active device and link failure localization in data center networks. In *16th {USENIX} Symposium on Networked Systems Design and Implementation ({NSDI} 19)*, pages 599–614, 2019.
- [44] Mea Wang, Baochun Li, and Zongpeng Li. *sFlow: Towards resource-efficient and agile service federation in service overlay networks*. IEEE, 2004.
- [45] Yang Wu, Ang Chen, and Linh Thi Xuan Phan. Zeno: Diagnosing performance problems with temporal provenance. In *16th {USENIX} Symposium on Networked Systems Design and Implementation ({NSDI} 19)*, pages 395–420, 2019.
- [46] Tong Yang, Jie Jiang, Peng Liu, Qun Huang, Junzhi Gong, Yang Zhou, Rui Miao, Xiaoming Li, and Steve Uhlig. Elastic sketch: Adaptive and fast network-wide measurements. In *Proceedings of the 2018 Conference of the ACM Special Interest Group on Data Communication*, pages 561–575. ACM, 2018.
- [47] Da Yu, Yibo Zhu, Behnaz Arzani, Rodrigo Fonseca, Tianrong Zhang, Karl Deng, and Lihua Yuan. dshark: a general, easy to program and scalable framework for analyzing in-network packet traces. In *16th {USENIX} Symposium on Networked Systems Design and Implementation ({NSDI} 19)*, pages 207–220, 2019.
- [48] Minlan Yu, Lavanya Jose, and Rui Miao. Software defined traffic measurement with opensketch. In *NSDI*, volume 13, pages 29–42, 2013.
- [49] Hongyi Zeng, Peyman Kazemian, George Varghese, and Nick McKeown. Automatic test packet generation. In *Proceedings of the 8th international conference on Emerging networking experiments and technologies*, pages 241–252. ACM, 2012.
- [50] Yibo Zhu, Nanxi Kang, Jiaxin Cao, Albert Greenberg, Guohan Lu, Ratul Mahajan, Dave Maltz, Lihua Yuan, Ming Zhang, Ben Y Zhao, et al. Packet-level telemetry in large datacenter networks. In *Proceedings of the 2015 ACM Conference on Special Interest Group on Data Communication*, pages 479–491, 2015.
- [51] Thomas Zimmermann, Andreas Zeller, Peter Weissgerber, and Stephan Diehl. Mining version histories to guide software changes. *IEEE Transactions on Software Engineering*, 31(6):429–445, 2005.

APPENDIX

The first step of PCAT is to collect data from devices, which we call discovered data. There are three types of data including numeric counters, non-numeric states, and configurations. Table 4 shows the examples for each category.

Types	Categories & examples	Impact of software upgrades
Counters	<i>Device utilization</i> : CPU&memory utilization, routing table size, etc	Ambiguity between percentage and absolute values.
	<i>Device internal status</i> : Interface error counter, power supply temperature, fan speeds, linecard version, optical CRC error counter, etc	XML format gets changed; linecard version format changes from integer to string.
	<i>Packet processing counters</i> : Packet drops, errors, queue length, etc	Ambiguity of interface stats meaning.
	<i>Protocol counters</i> : BGP neighbor received routes, etc	General empty data error.
States	<i>Interface state</i> : Interface up, down, drained, configured IP address, MAC address, etc	Hex-decimal change causes MAC address retrieving error.
	<i>Protocol state</i> : BGP neighbor state, etc	State meaning ambiguity.
Configs	BGP policy, queuing algorithm, etc	Raw config format changed.

Table 4: Different discovered data in PCAT.