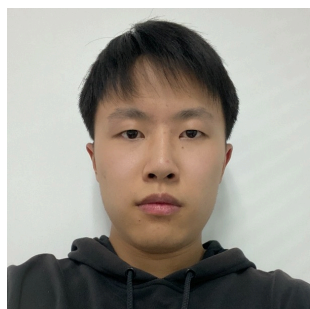


# DINT: Fast In-Kernel Distributed Transactions with eBPF



**\* Yang Zhou**

*Harvard University*



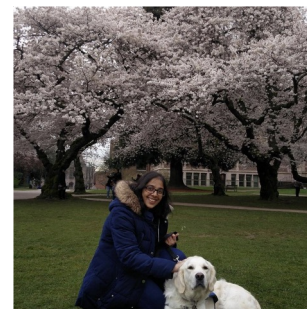
**\* Xingyu Xiang**

*Peking University*



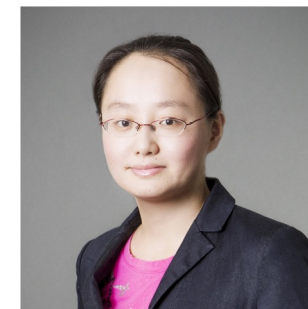
Matthew Kiley

*Harvard University*



Sowmya  
Dharanipragada

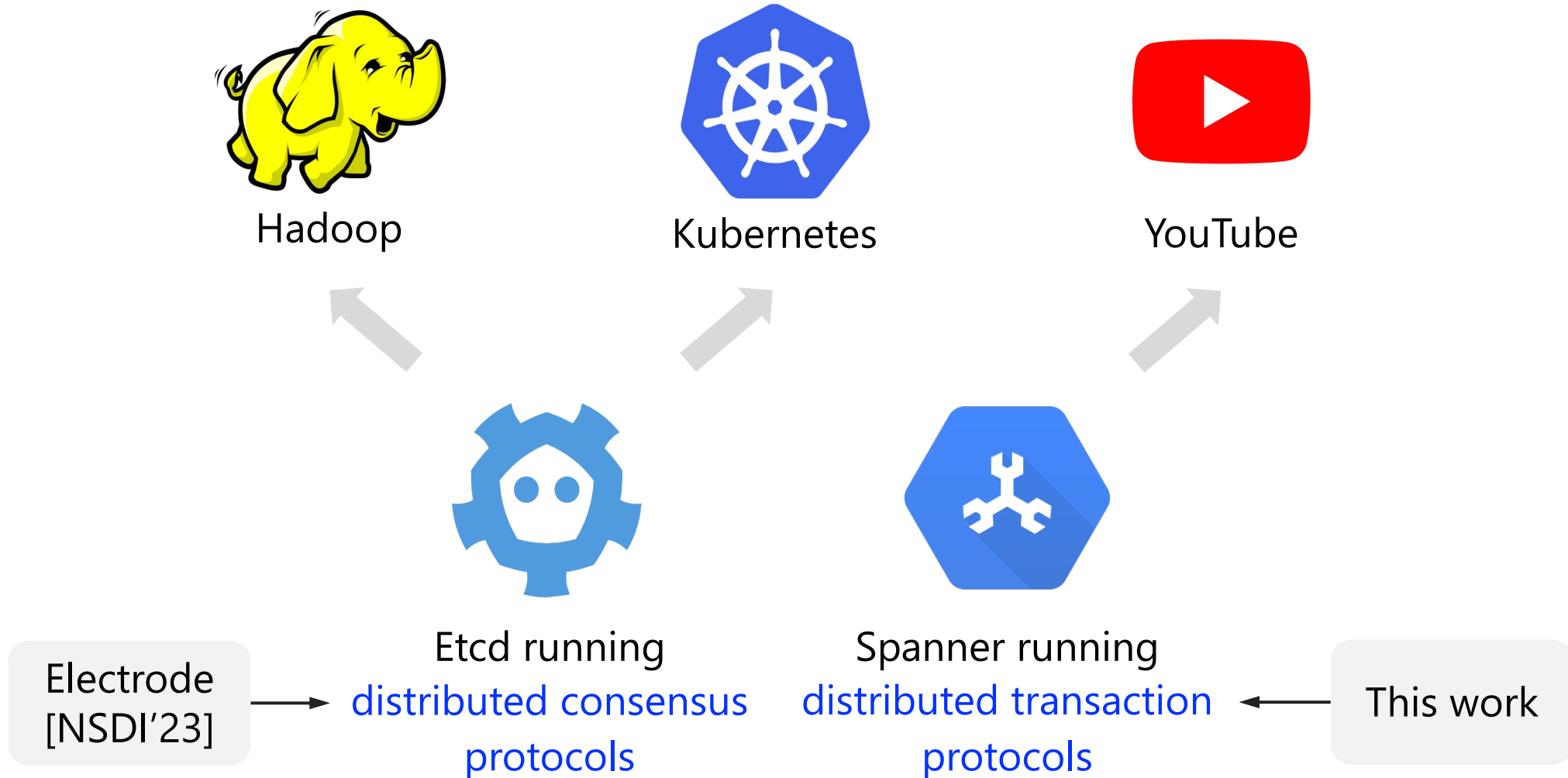
*Cornell University*



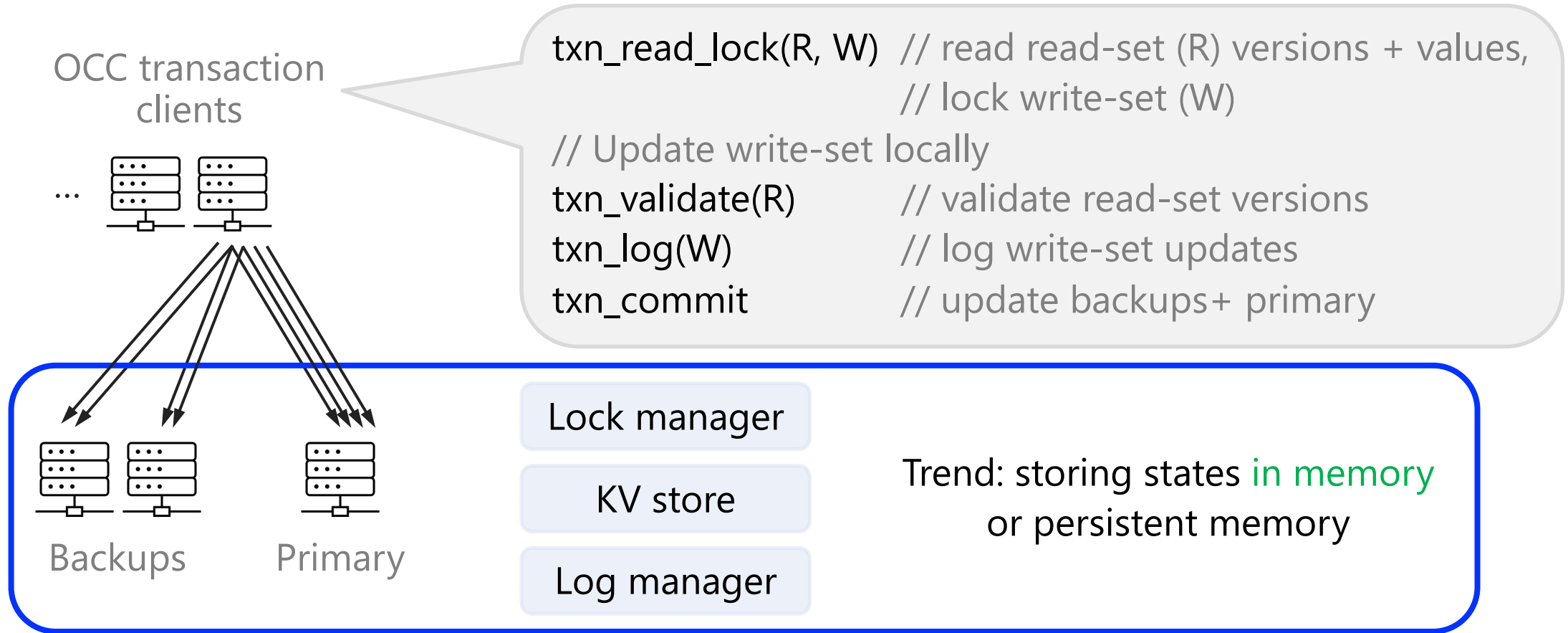
Minlan Yu

*Harvard University*

# Distributed Protocols for High Availability

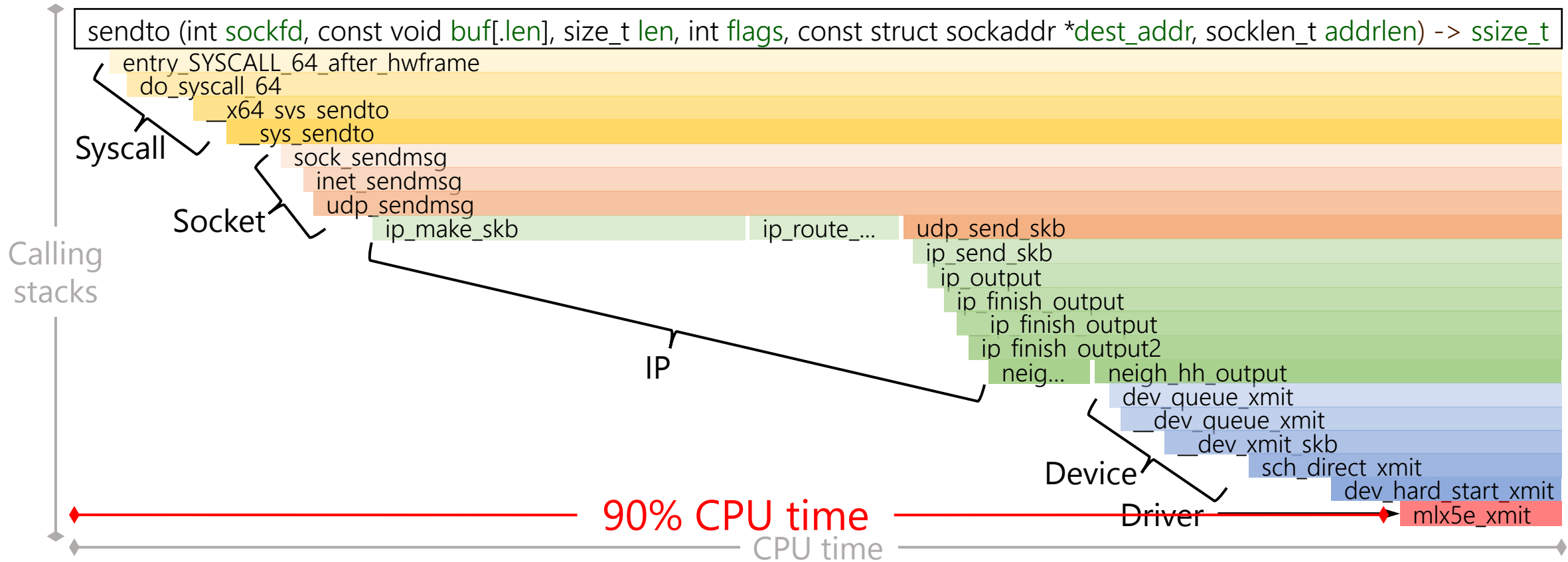


# Distributed Transactions inside a Datacenter

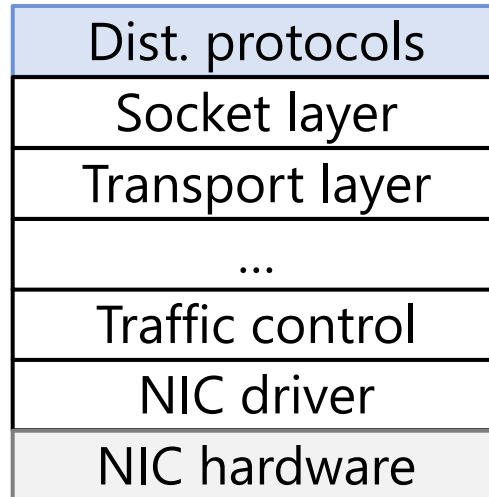


Distributed transactions are network IO-intensive

# Kernel Networking: High Kernel Overhead



# Kernel Networking: High Kernel Overhead



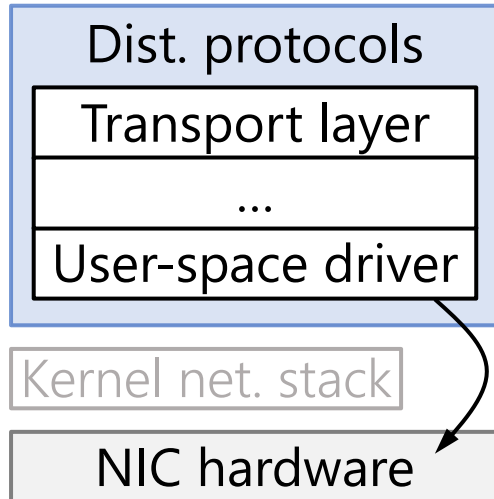
OCC distributed transactions

92% CPU time<sup>1</sup>



Only around  $\frac{1}{10}$  is on NIC driver

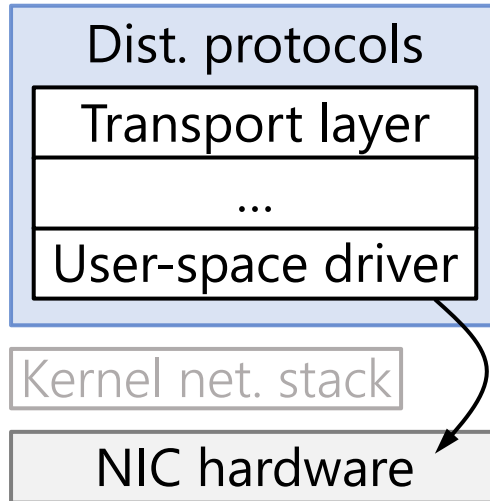
# Kernel Bypass: Not a Panacea



DPDK (Data Plane Development Kit) or RDMA:

- Customized networking stacks in user space or NIC
- Busy polling instead of costly interrupt

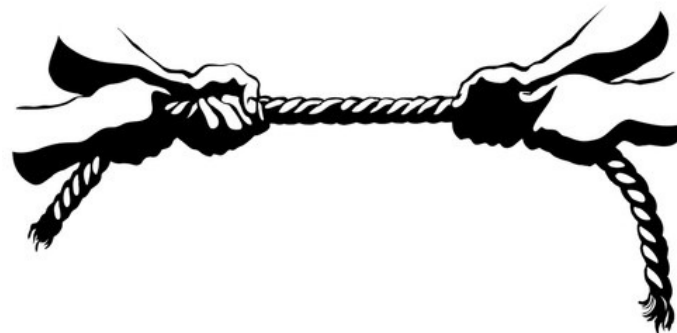
# Kernel Bypass: Not a Panacea



DPDK (Data Plane Development Kit) or RDMA:

- + High performance
- Dedicated resources (eg, busy-polling cores)
- Security vulnerabilities (user manages NICs)<sup>1,2</sup>

Kernel bypass:  
high performance

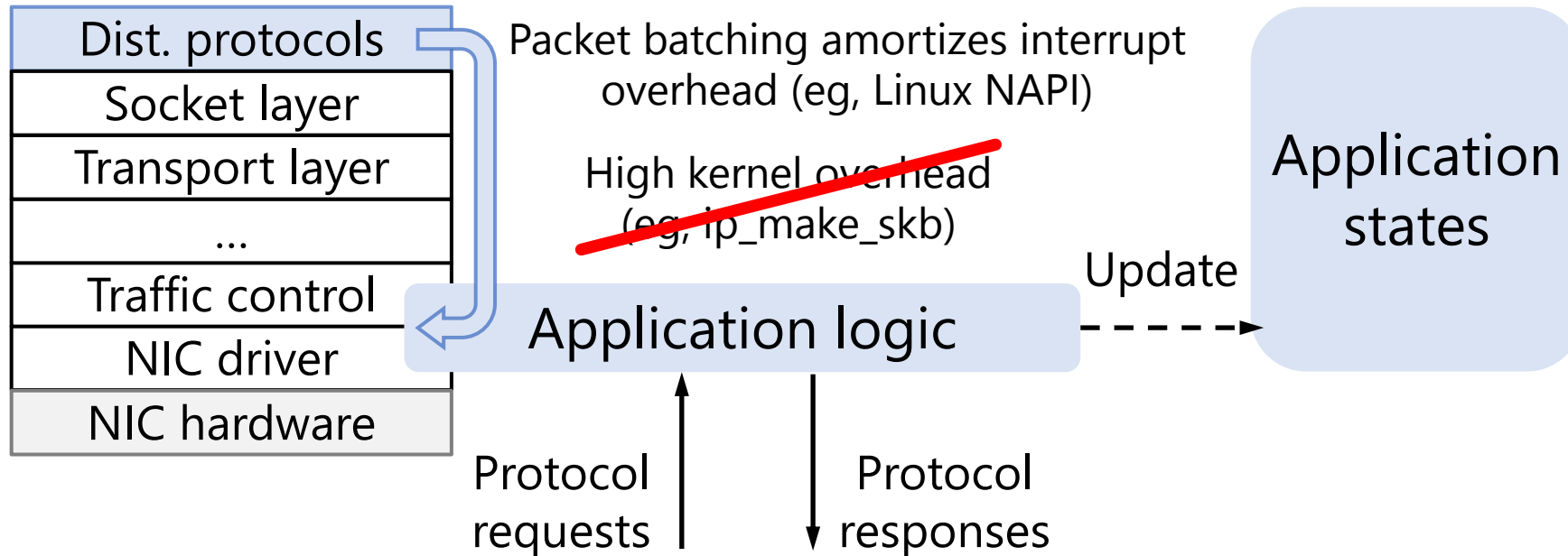


Kernel:  
resource sharing, security

[1] Bellovin, Steven M. "Security Problems in the TCP/IP Protocol Suite." SIGCOMM CCR 1989

[2] Smolyar et al. "Securing Self-Virtualizing Ethernet Devices." USENIX Security 2015

# DINT<sup>1</sup>: Application-Customized Networking Stacks



- + High performance
- + Resource sharing: interrupt-driven
- + Secure: kernel manages NICs



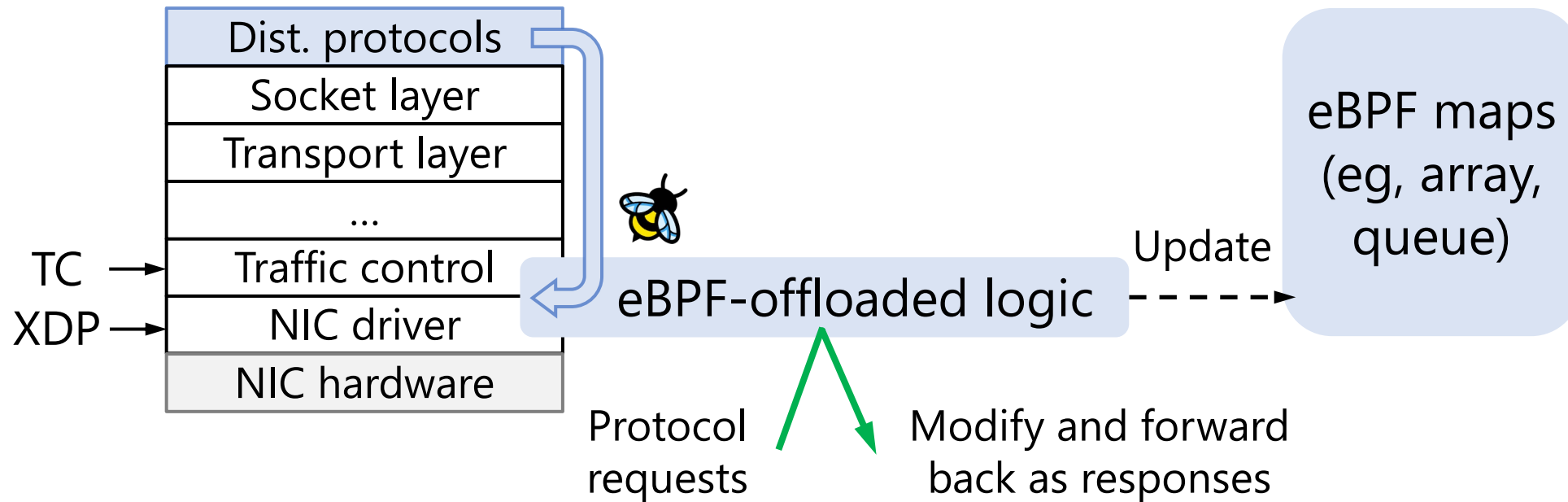
[1] DINT: an archaic word, meaning force and power



# How to Guarantee Kernel Safety?

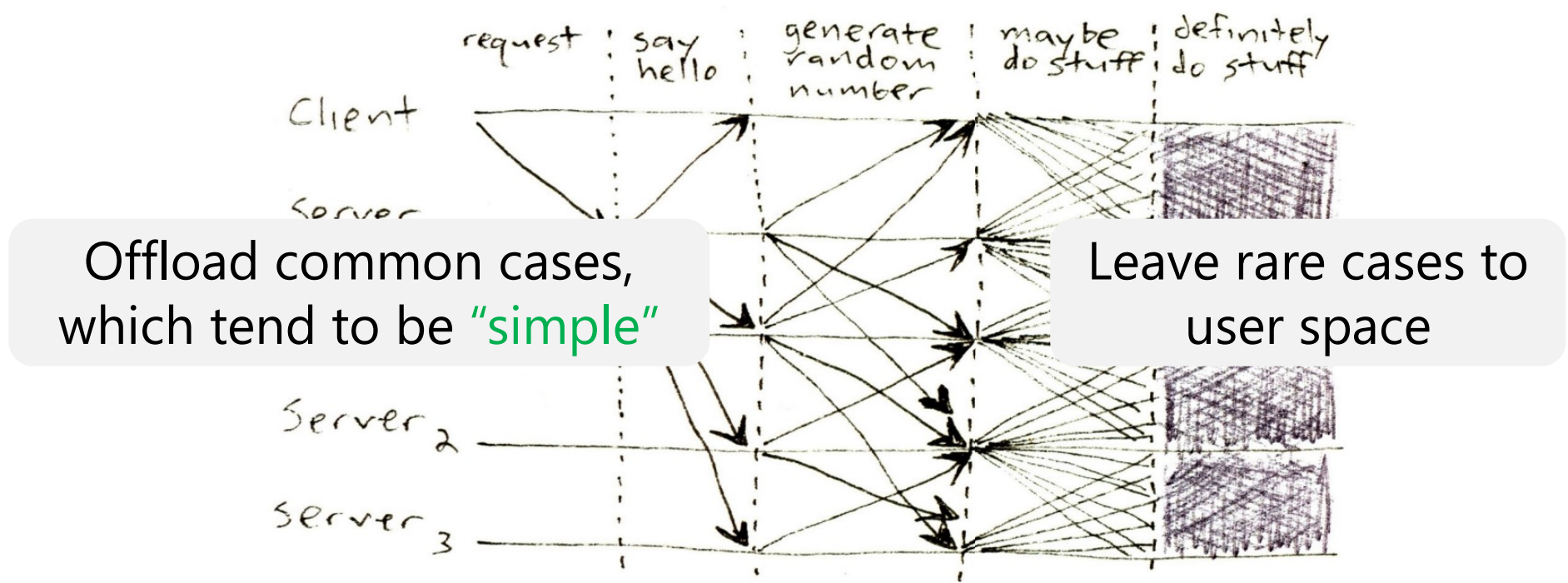
eBPF (extended Berkeley Packet Filter) to safely run programs in kernel at runtime

- Guaranteeing safety via static verification
- Originally for packet filtering and monitoring

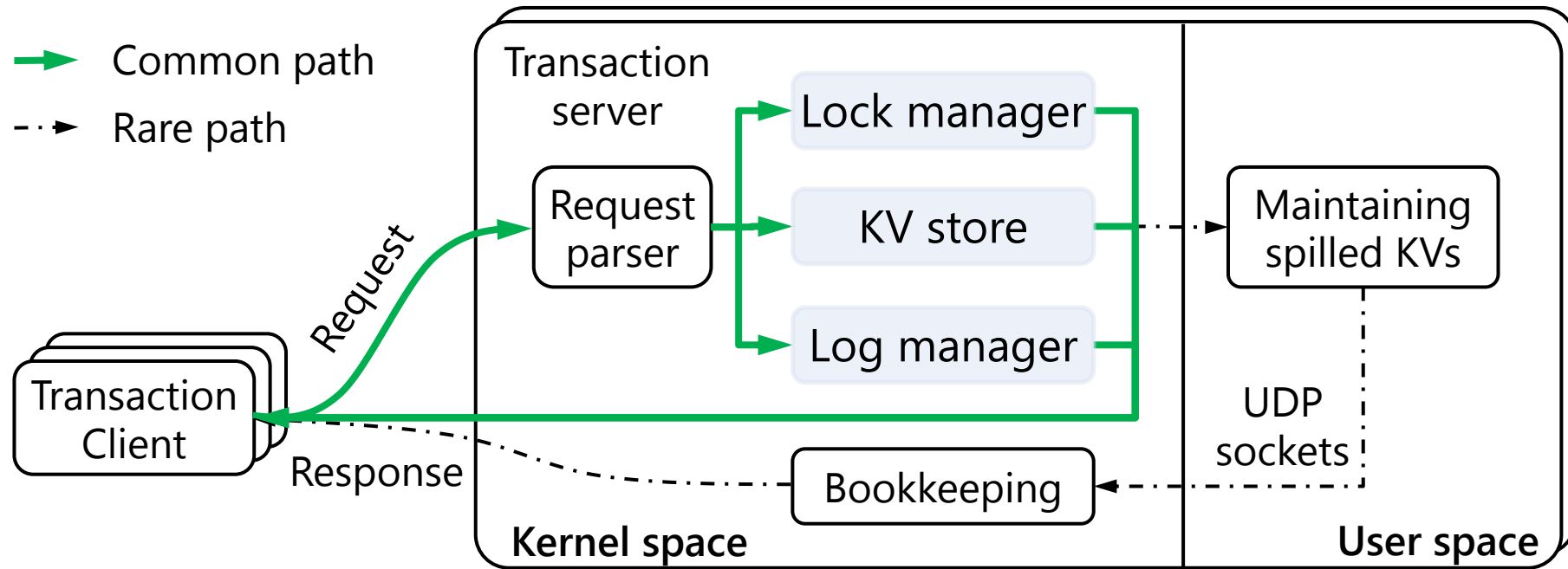


# Challenge of Kernel Offloads with eBPF

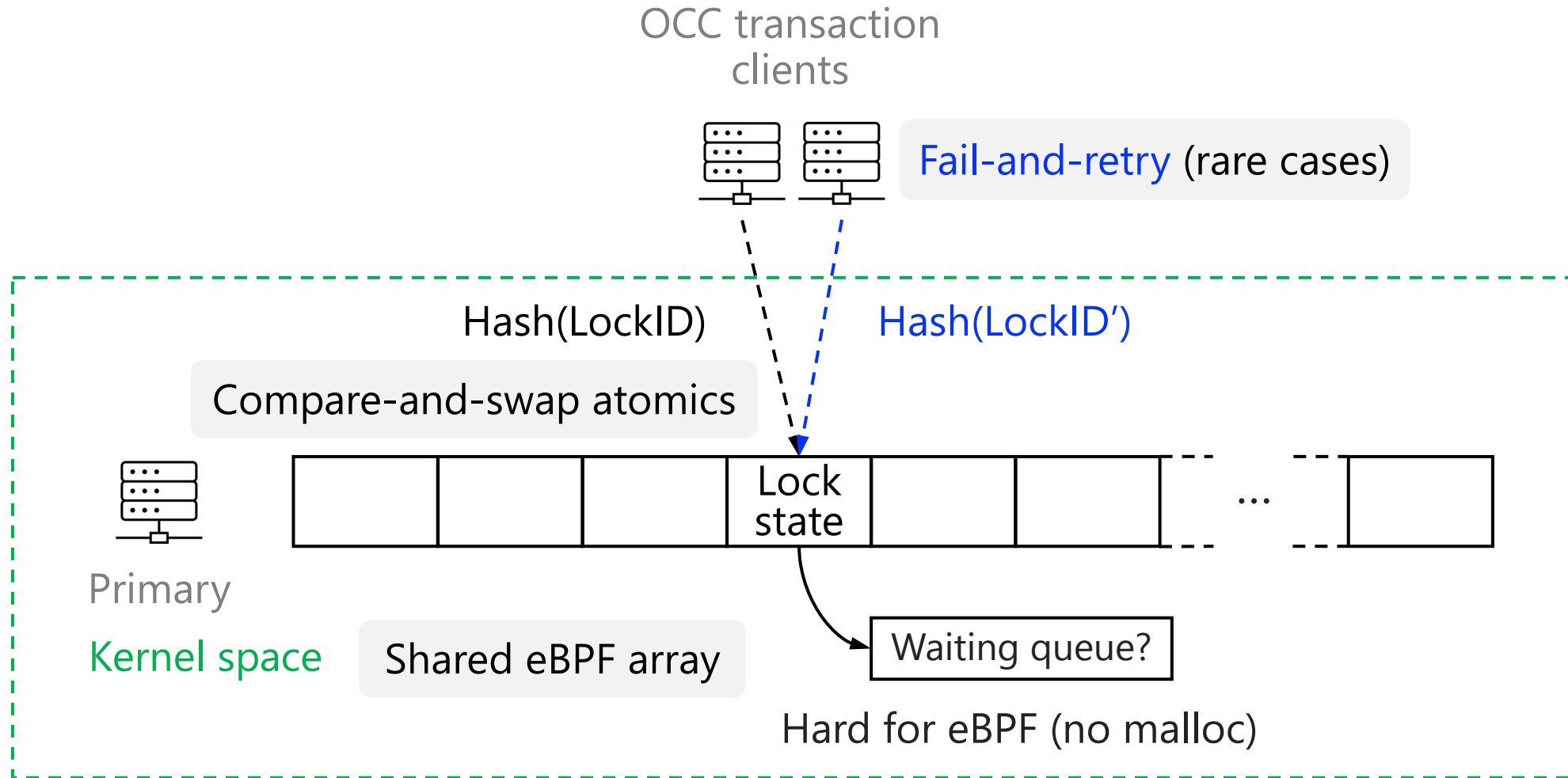
- eBPF programming model is **constrained** because of static verification
  - Limited # of instructions, bounded loops, static memory allocation
- Distributed protocols are **complex**
  - Some rare cases are too complex for eBPF: eg, failure, message loss, malloc



# DINT Overall Architecture



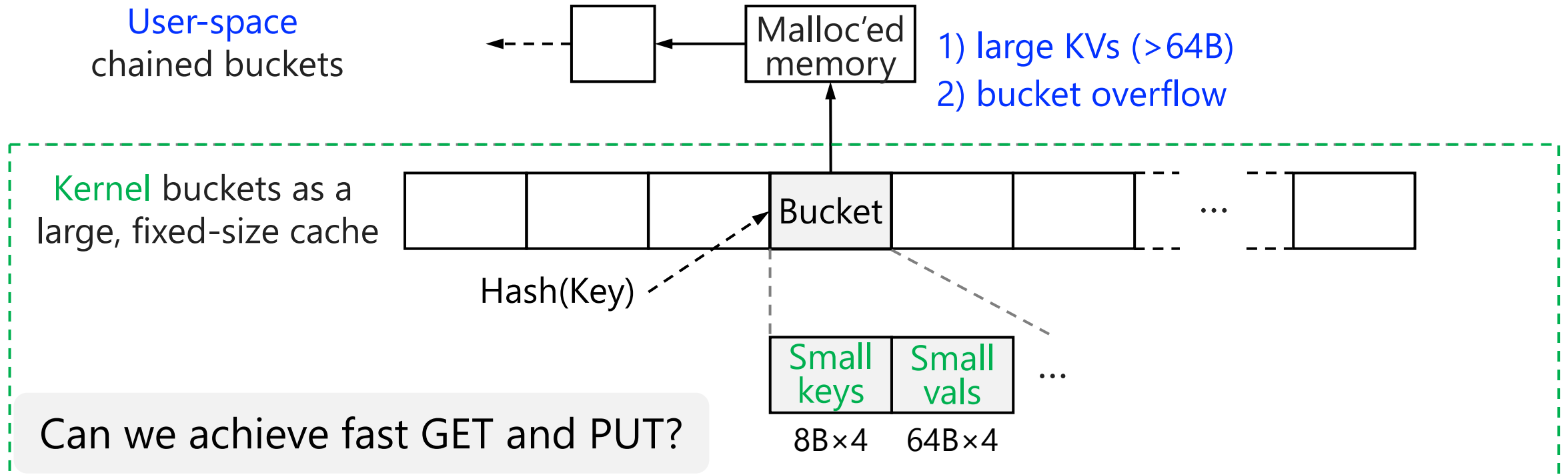
# Offloading Lock Manager



# Offloading KV Store

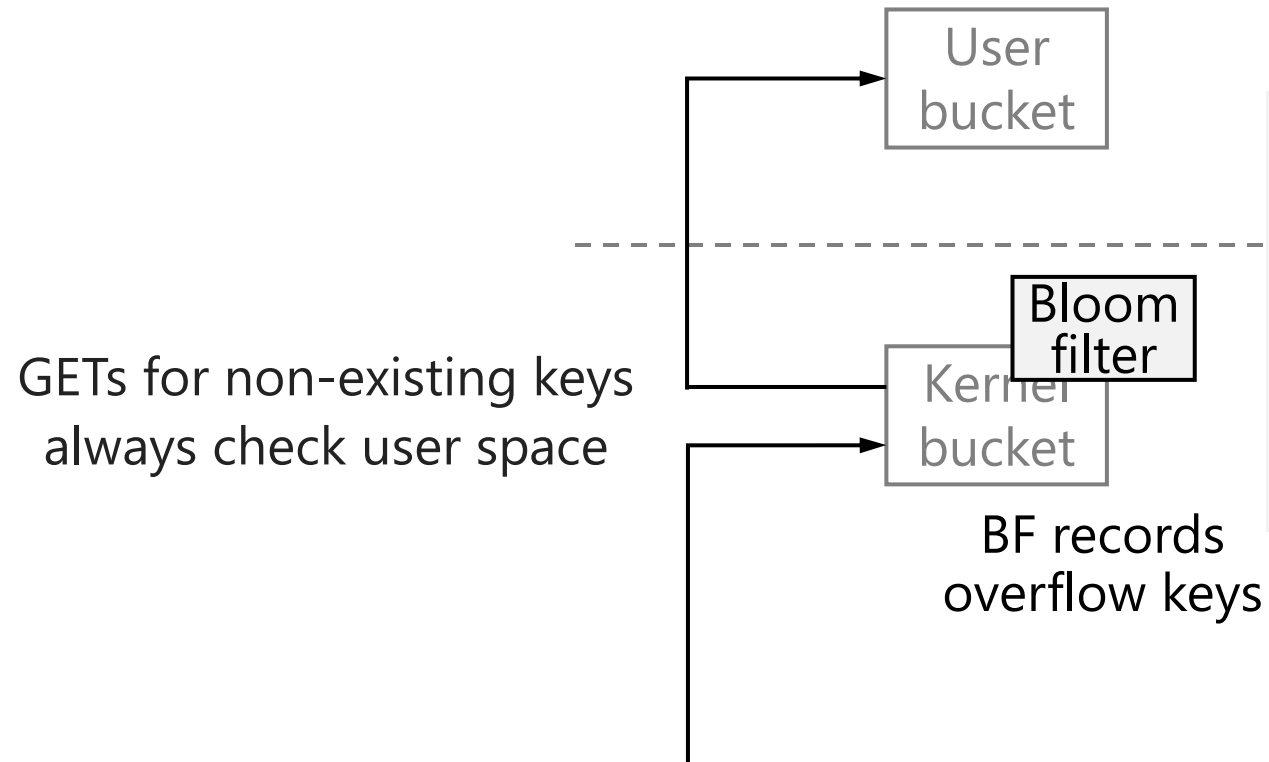
Common cases: most KVs are small in typical workloads

- Dozens of bytes in transactional workloads (eg, TATP, SmallBank)
- Statically-allocated eBPF map to store small KVs and avoid malloc



# Offloading KV Store

How to achieve **fast GET** especially for non-existing keys?



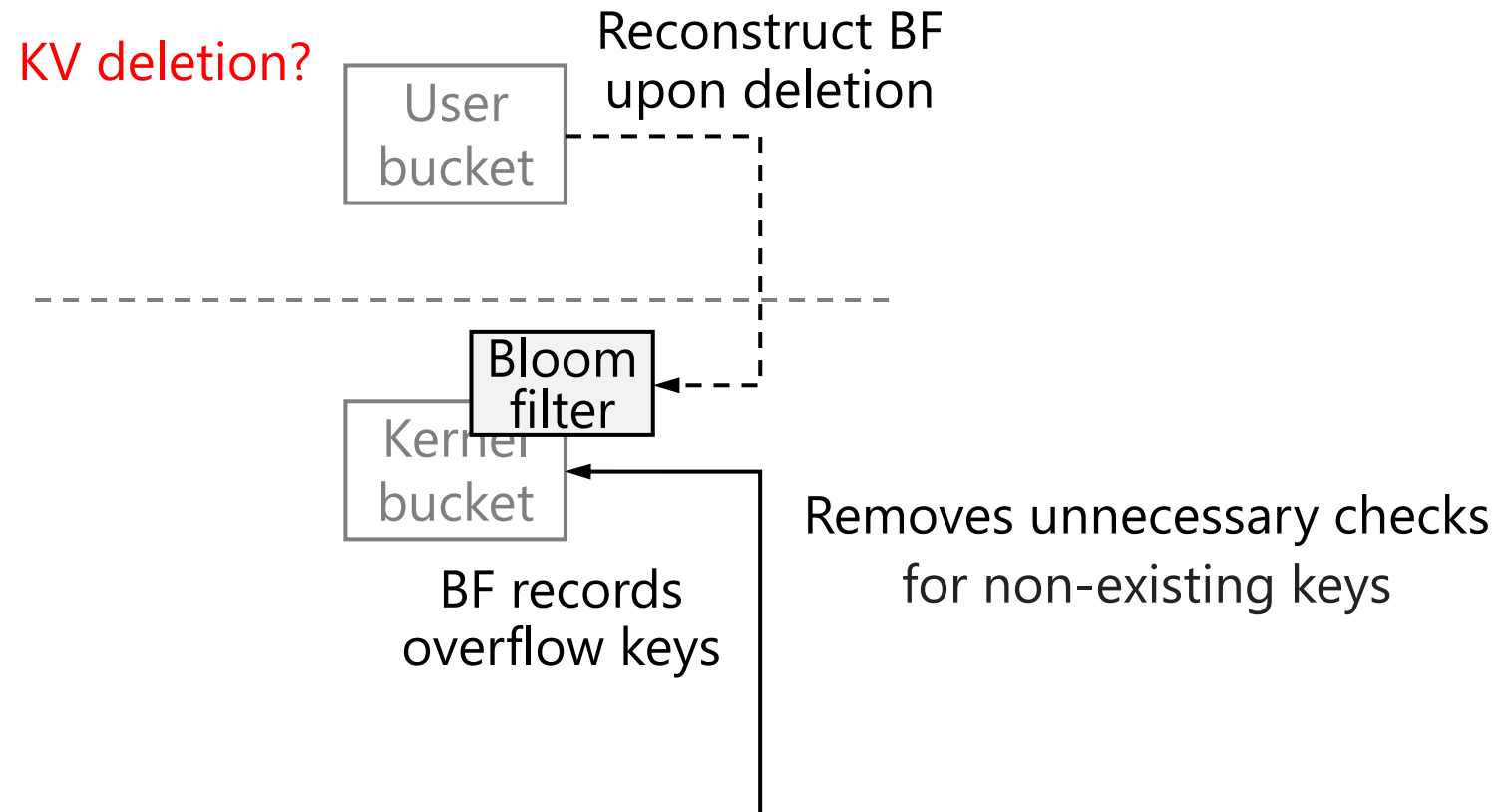
Bloom Filter<sup>1</sup>:

- An approximate data structure that quickly tells whether a key is in a set
- Using a **fixed size** of bitset array

[1] Bloom, Burton H. "Space/Time Trade-offs in Hash Coding with Allowable Errors." Communications of the ACM 13.7 (1970): 422-426

# Offloading KV Store

How to achieve **fast GET** especially for non-existing keys?



# Offloading KV Store

How to achieve **fast GET** especially for non-existing keys?

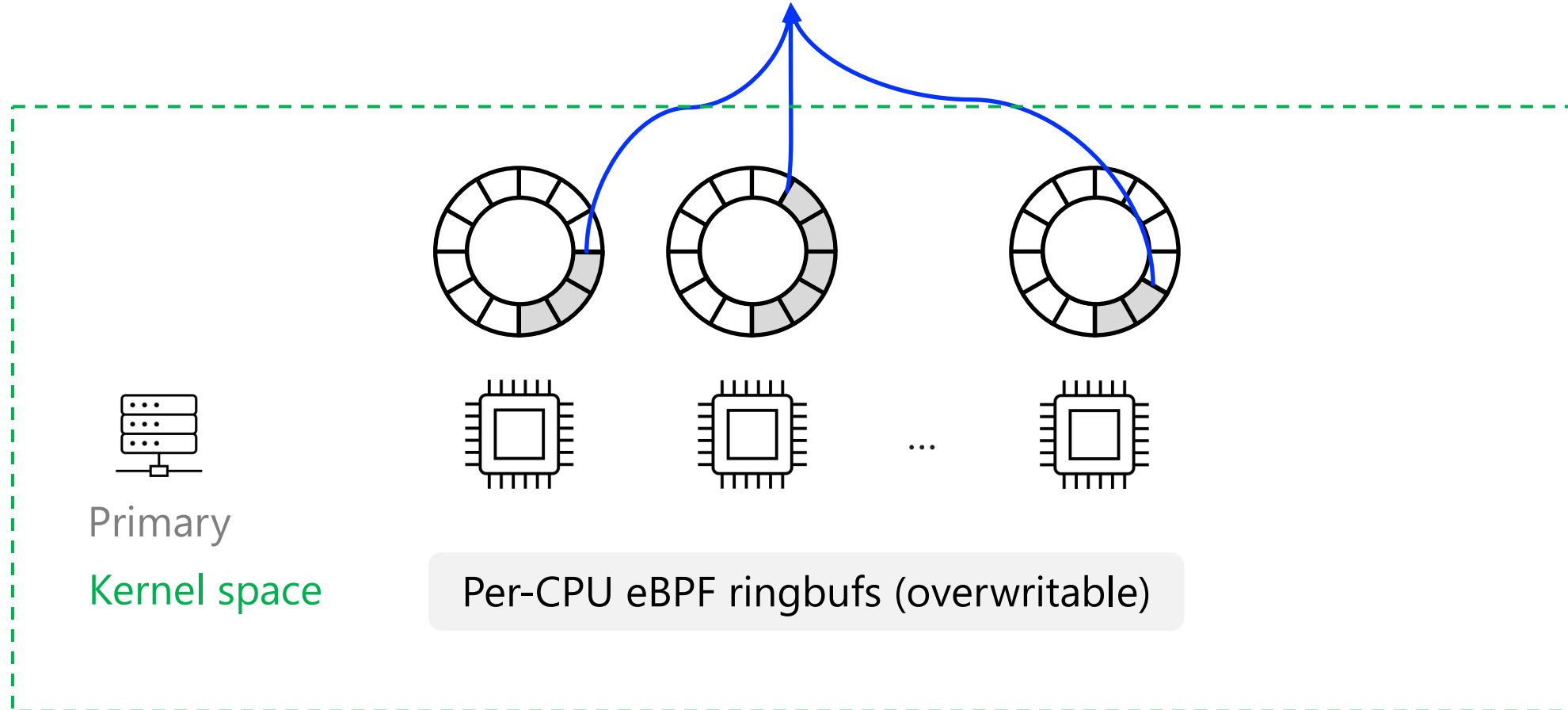
Bloom filter to record overflow keys

- Write-back cache for **fast PUT**
- Lock sharing for **fast locking**
- Per-core circular logs for **fast logging**
- Piggybacking states on packets for **fast user-kernel synchronization**
- ...



# Offloading Log Manager

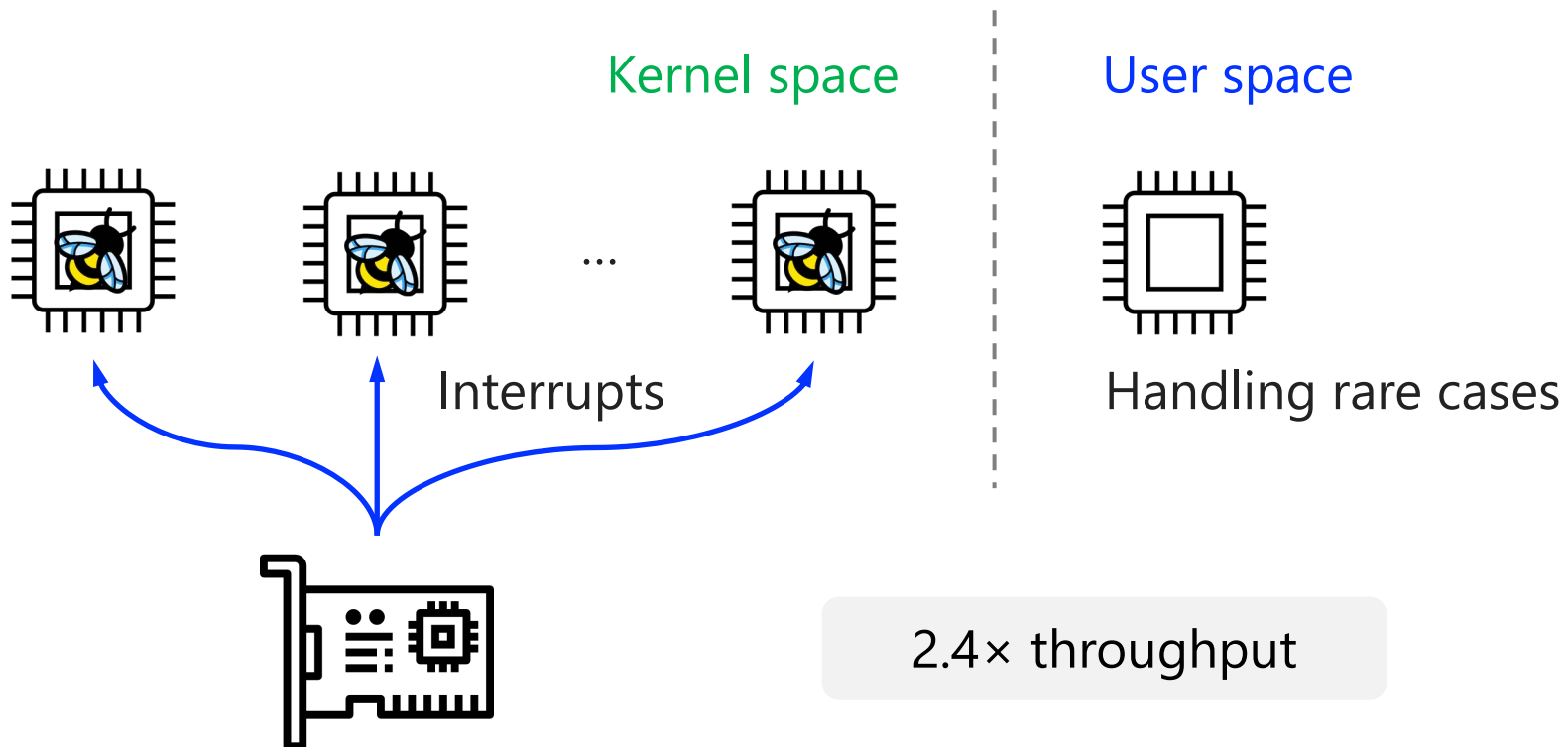
User space replays the logs on failure (rare cases)



# System Optimization: Interrupt Scheduling

Separating interrupt-handling cores and rare-case handling cores

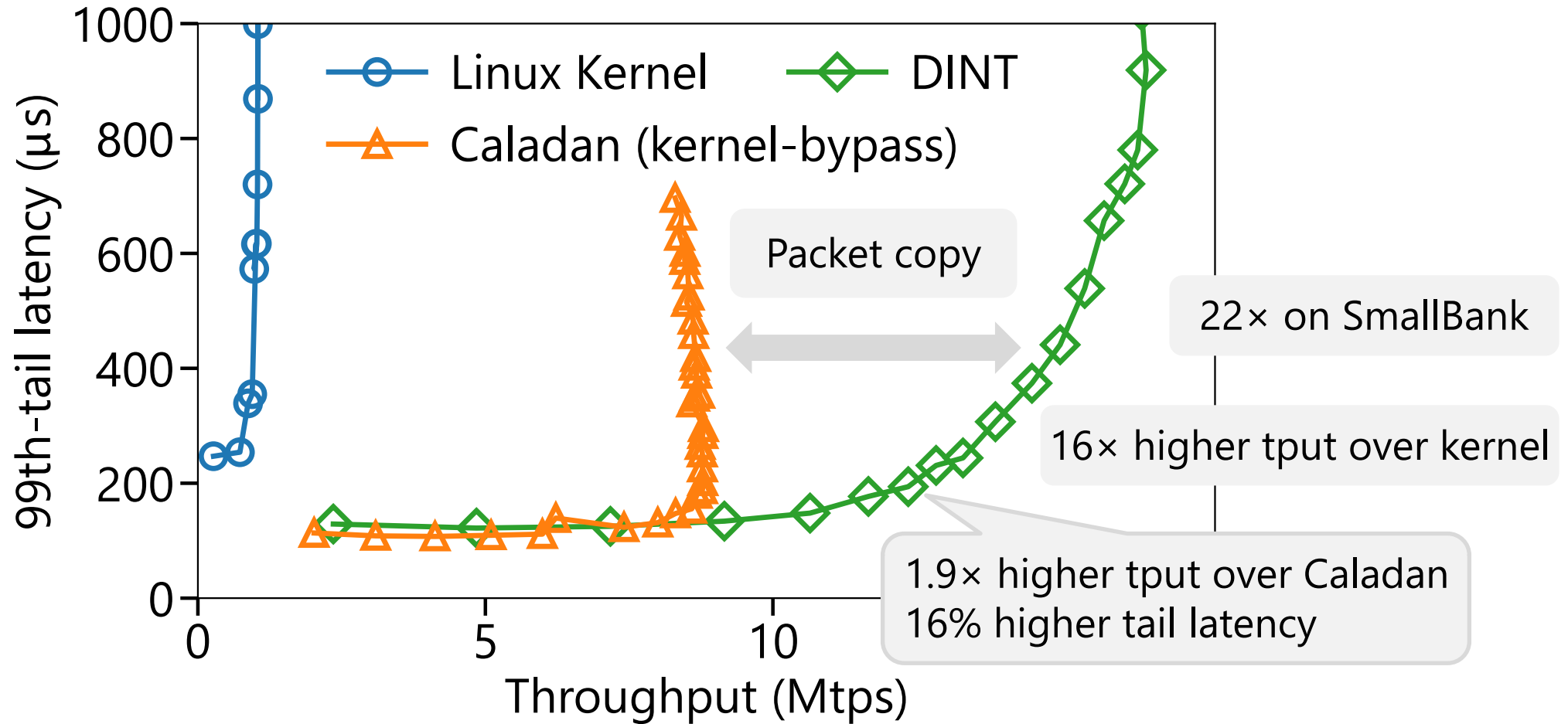
↳ Avoiding user-kernel context switching overhead



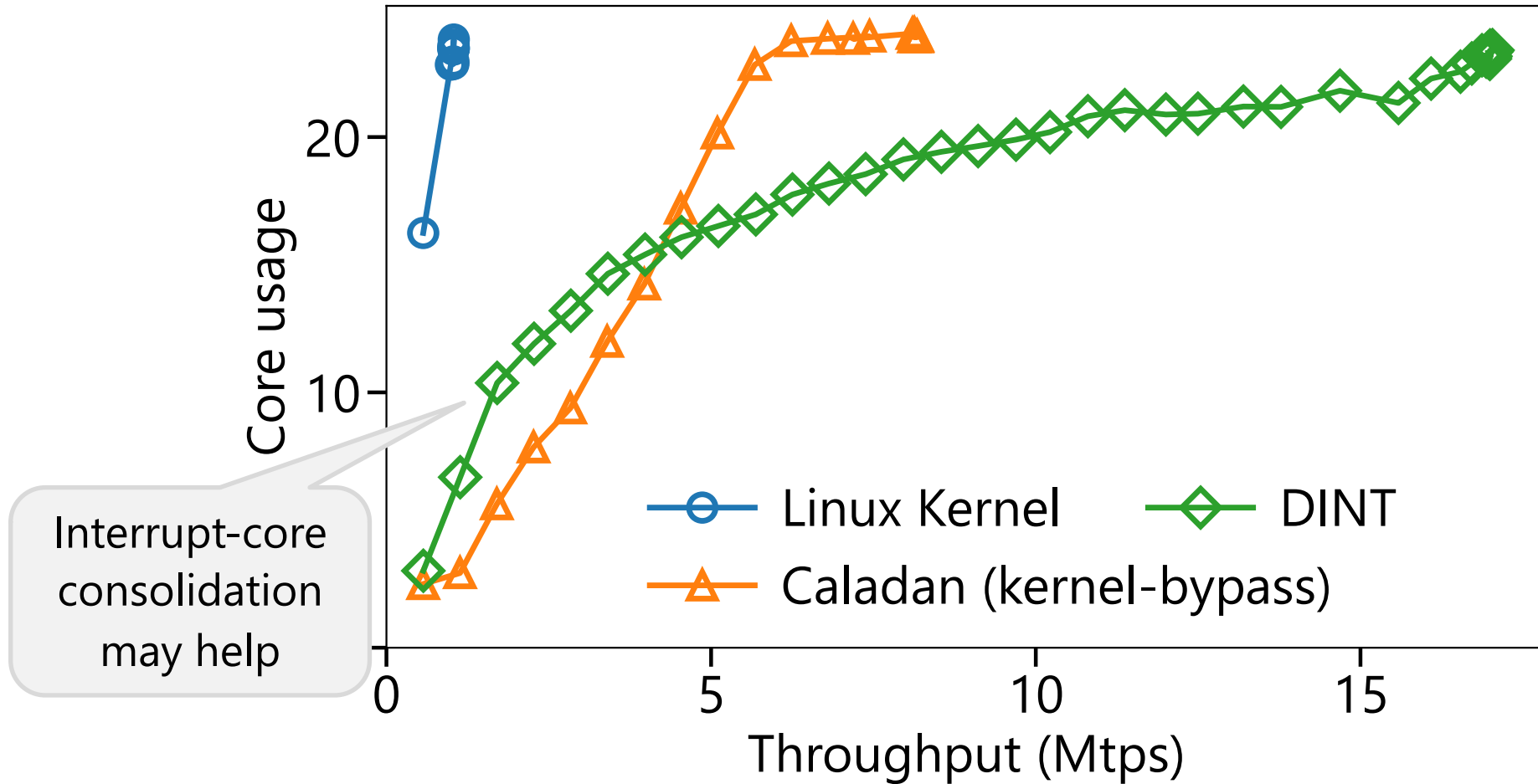
# Implementation & Evaluation

- 2.1K lines eBPF and 4.3K C++ for OCC and 2PL distributed transactions
  - 3-way primary-backup replication, 3-way sharding
  - 6.1K lines of C++ for baselines
- Experiment setup:
  - CloudLab r650 (10 clients, 3 servers) running **unmodified** Linux kernel 6.1.0
  - TATP for OCC, SmallBank for 2PL
- Open source: <https://github.com/DINT-NSDI24/DINT>

# Tail Latency vs. Throughput



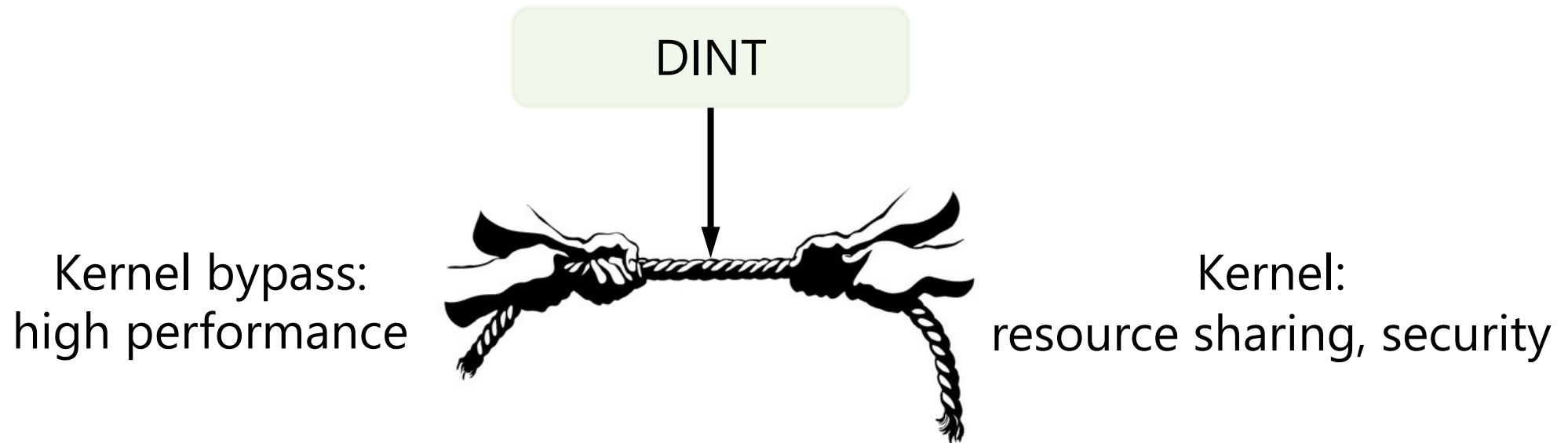
# CPU Utilization vs. Throughput



# DINT Conclusion

We enable application-customized kernel networking stacks with:

- eBPF offloads for common cases, while user space for rare cases
- distributed transaction offloads, but generalizable to many distributed protocols



Thank you!

[yangzhou@g.harvard.edu](mailto:yangzhou@g.harvard.edu)