# XRP: In-Kernel Storage Function with eBPF
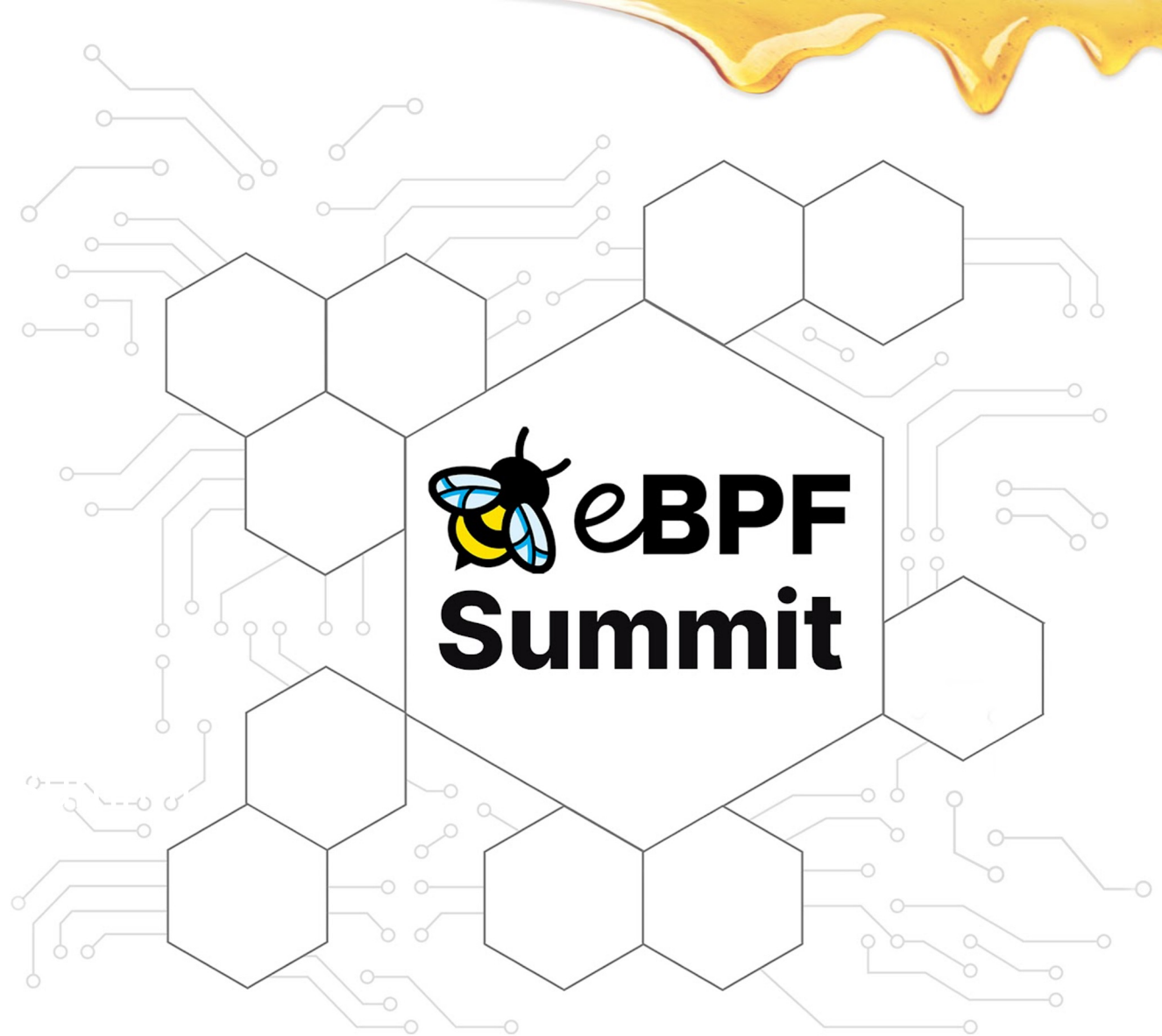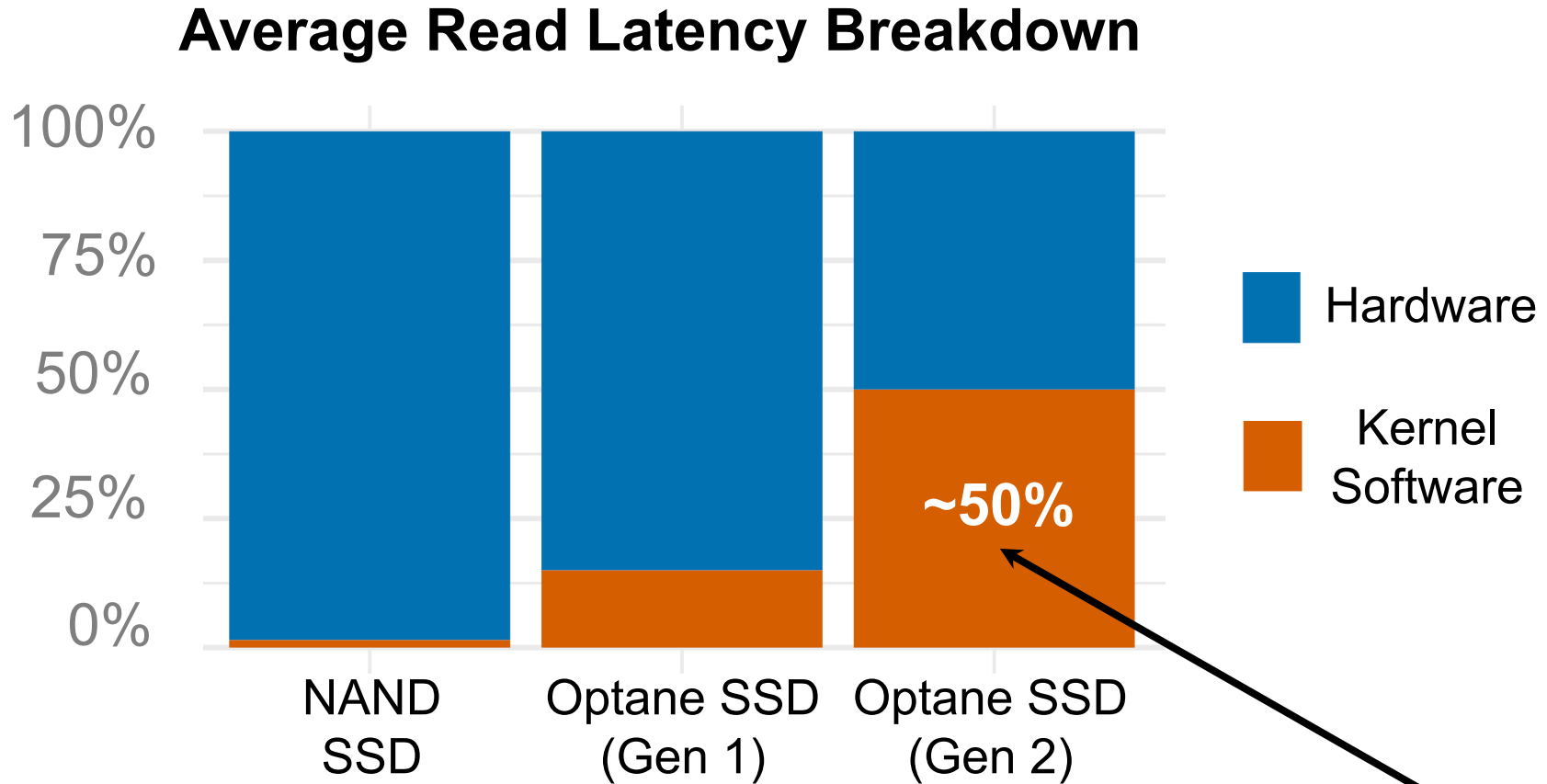
Yuhong Zhong

@zhong_yuhong
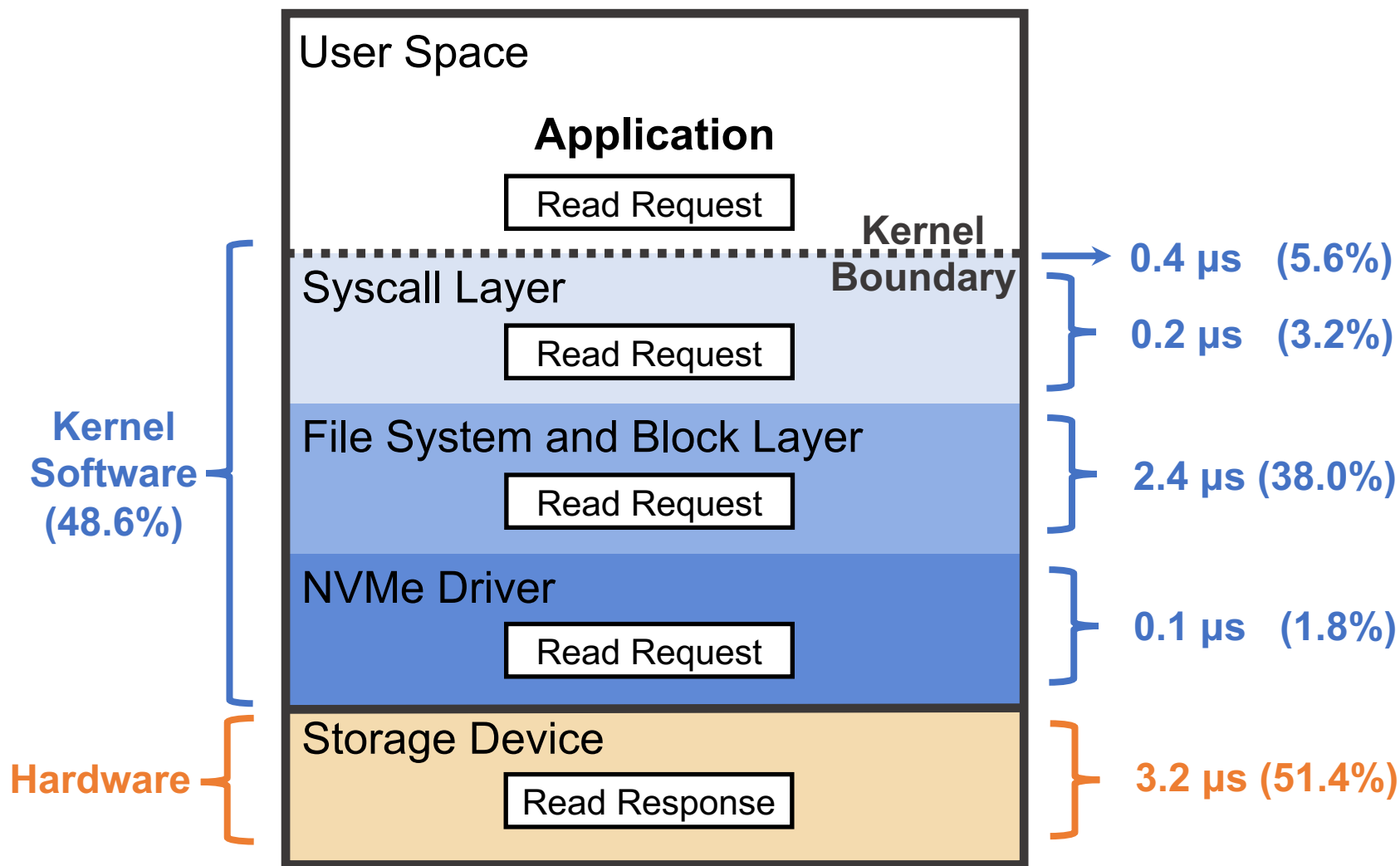
# Kernel Software is Becoming the Bottleneck for Storage

## Average Read Latency Breakdown



Kernel software overhead accounts for ~50% of read latency on Optane SSD Gen 2

Workload: Random 512B Read

# Where Does the Latency Come From?



User Space

**Application**

Read Request

**Kernel Boundary**

Syscall Layer

Read Request

0.4 μs (5.6%)

0.2 μs (3.2%)

**Kernel Software (48.6%)**

File System and Block Layer

Read Request

2.4 μs (38.0%)

NVMe Driver

Read Request

0.1 μs (1.8%)

**Hardware**

Storage Device

Read Response

3.2 μs (51.4%)

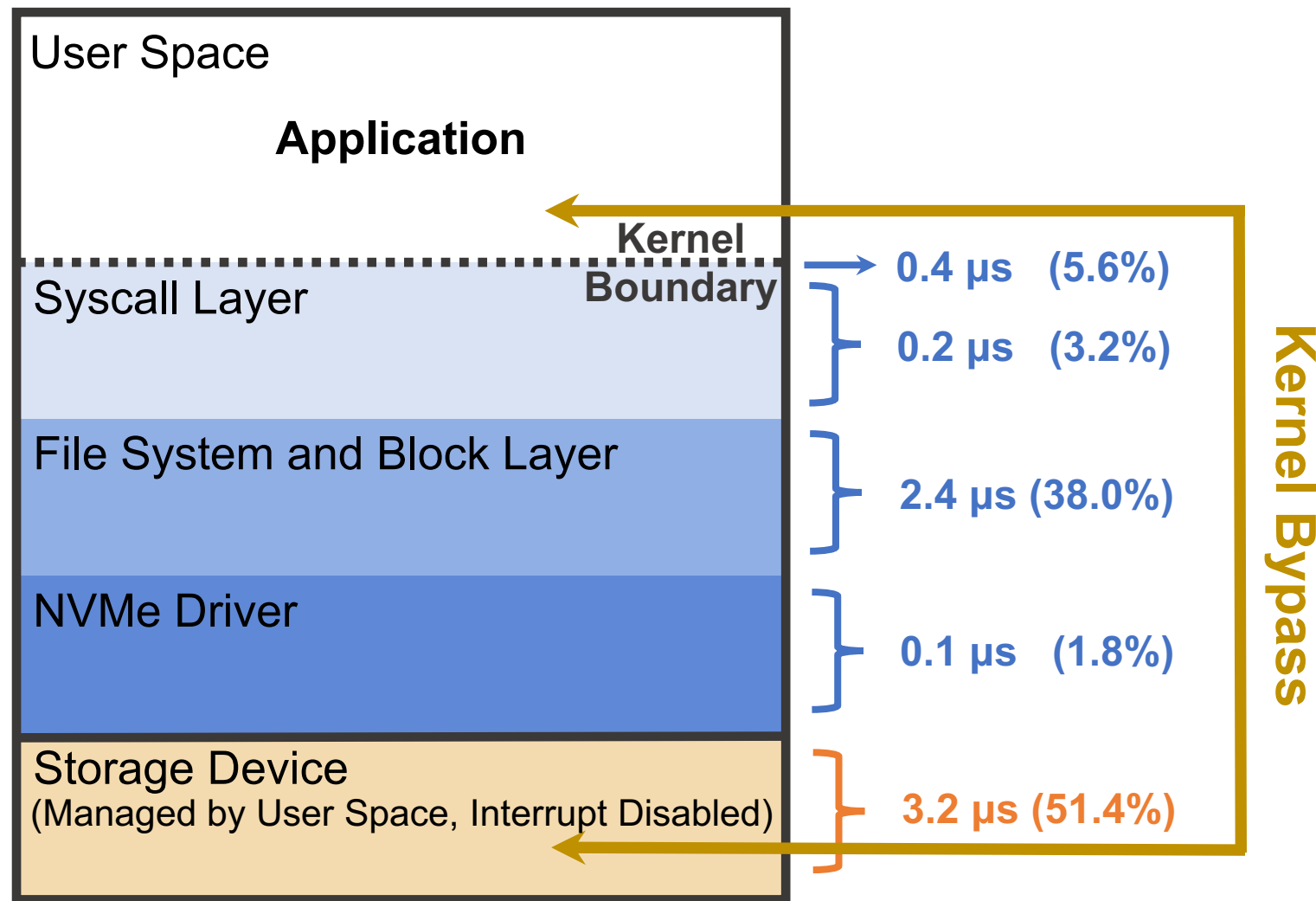Workload: Random 512B Read, Disk: Optane SSD P5800X

3

# Bypass Kernel to Eliminate Overhead

**Academic Work**

Demikernel (SOSP '21),
Shenango (NSDI '19),
Snap (SOSP '19),
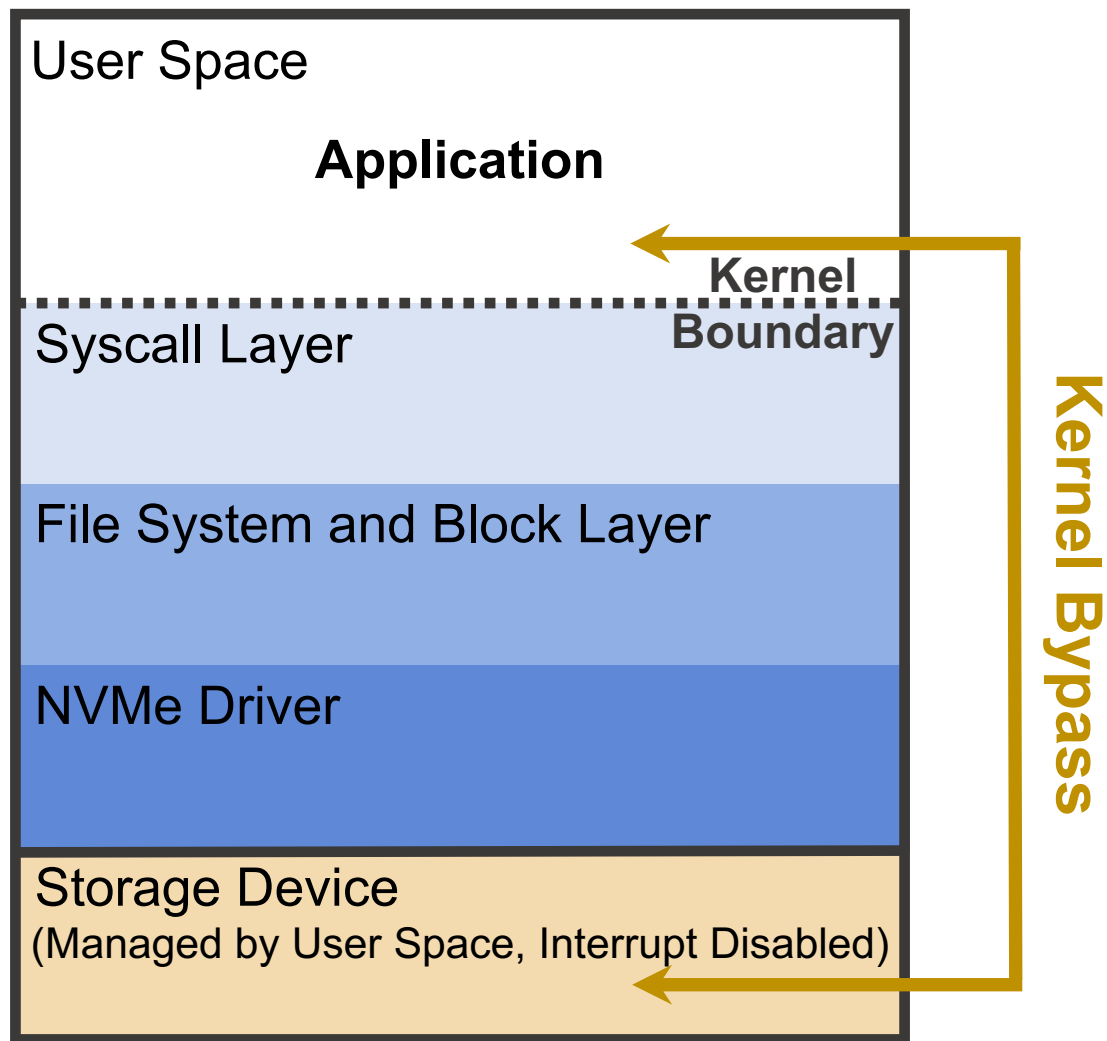IX (SOSP '17),
Arrakis (OSDI '14),
mTCP (NSDI '14),
...

**In industry, the most common library is SPDK**

**Reduce read latency by 49%**

**User Space**

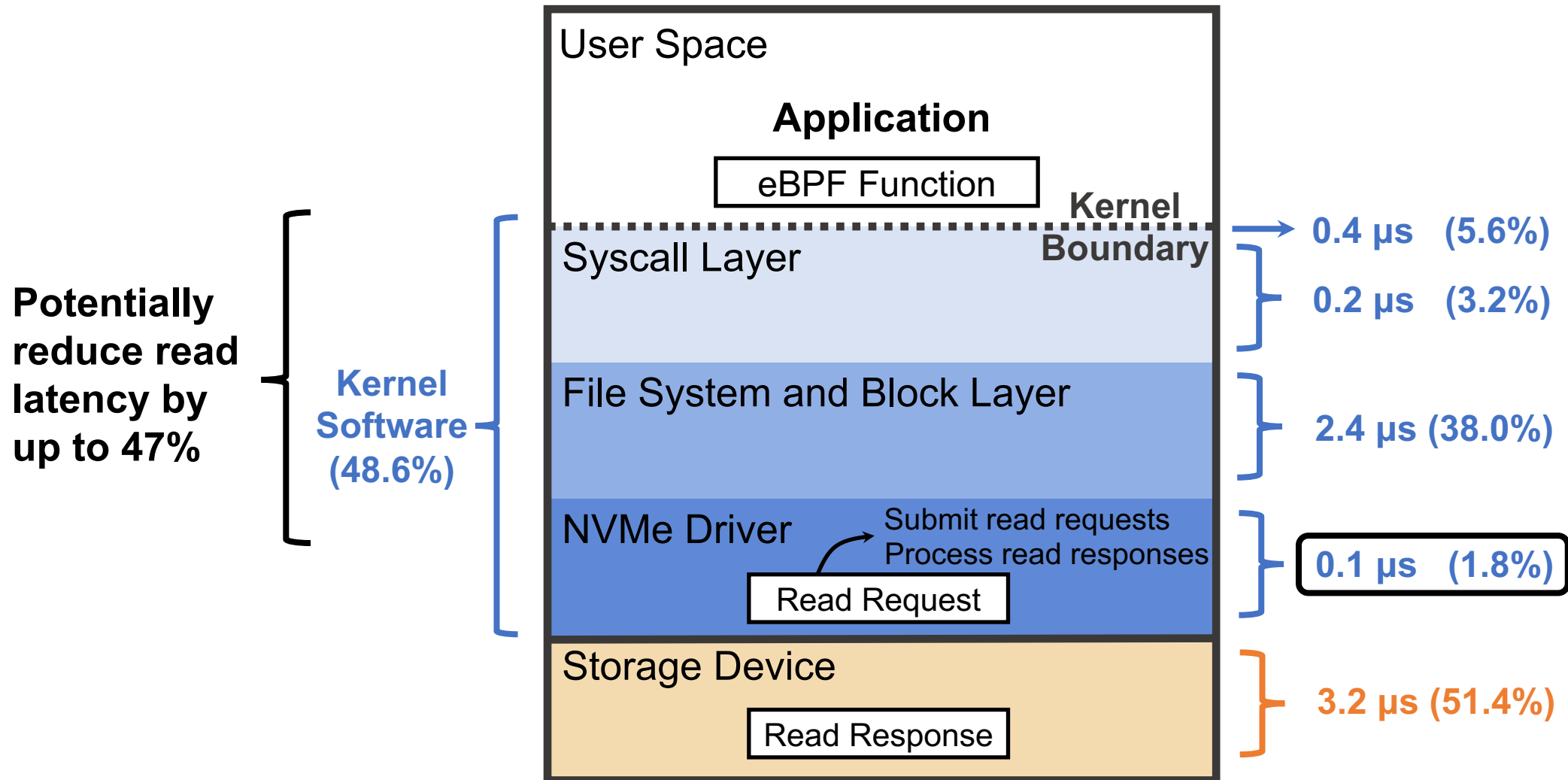**Application**

**Kernel Boundary**

Syscall Layer — 0.4 µs (5.6%) / 0.2 µs (3.2%)

File System and Block Layer — 2.4 µs (38.0%)

NVMe Driver — 0.1 µs (1.8%)

Storage Device
(Managed by User Space, Interrupt Disabled) — 3.2 µs (51.4%)

**Kernel Bypass**

# Kernel Bypass is Not a Panacea

**User Space**

**Application**

**Kernel Boundary**

Syscall Layer

File System and Block Layer

NVMe Driver

Storage Device
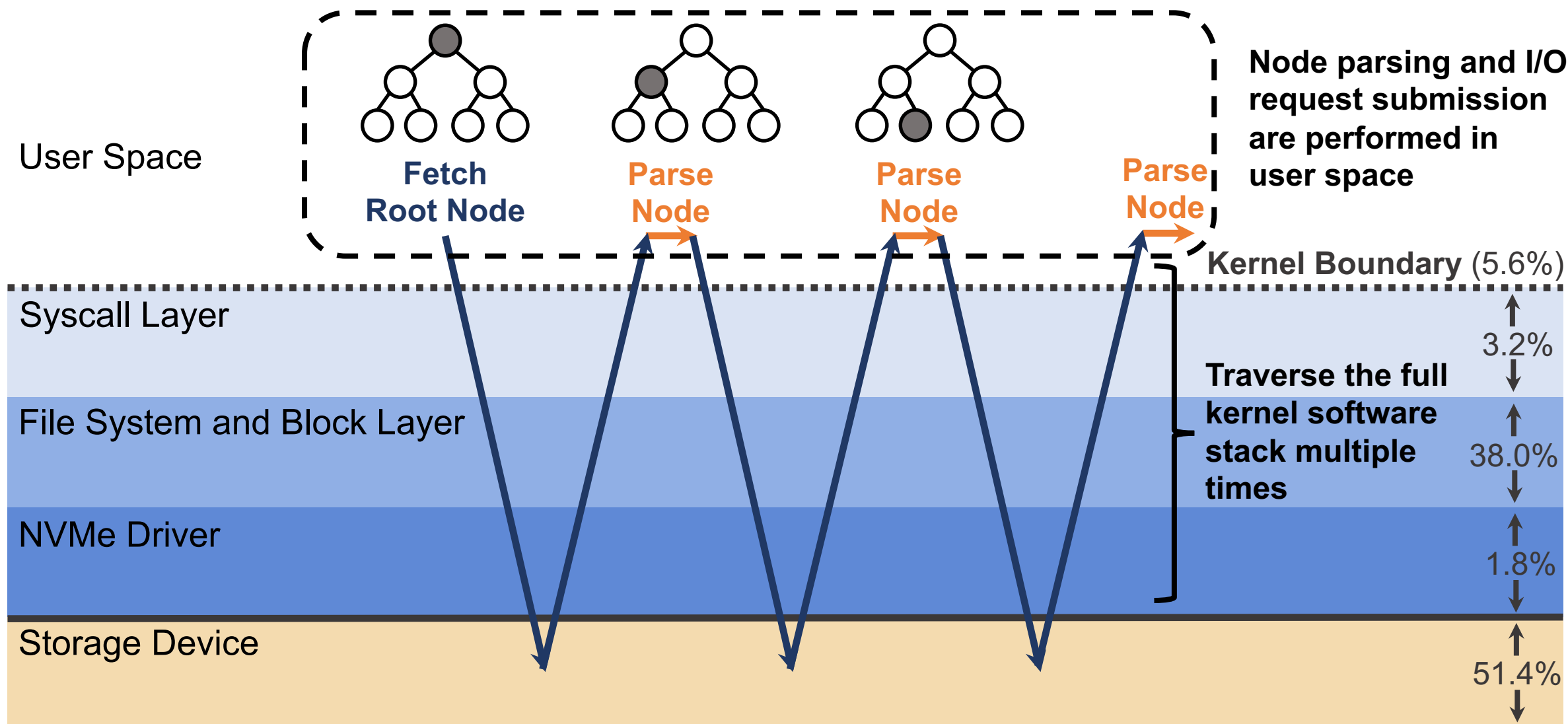(Managed by User Space, Interrupt Disabled)

**Kernel Bypass**

✅ Does not incur the overhead of the kernel storage stack

❌ No fine-grained access control

❌ Requires busy polling for completion

❌ Processes cannot yield CPU when waiting for I/O

❌ CPU cycles are wasted when I/O utilization is low

❌ CPU cannot be shared efficiently among multiple processes
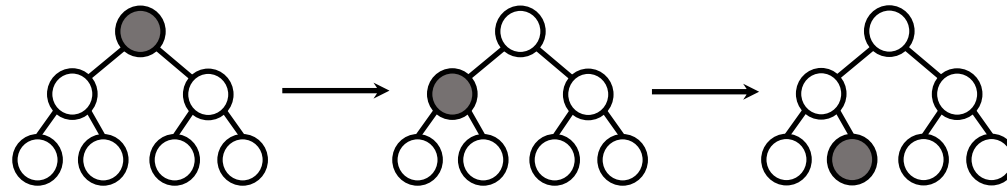
# Move Application Logic Into the Kernel

**Potentially reduce read latency by up to 47%**

**Kernel Software (48.6%)**

| User Space |
| --- |
| **Application** |
| eBPF Function |

**Kernel Boundary**

| Syscall Layer |
| --- |
| File System and Block Layer |
| NVMe Driver — Submit read requests / Process read responses / Read Request |
| Storage Device — Read Response |

0.4 μs   (5.6%)

0.2 μs   (3.2%)

2.4 μs (38.0%)

0.1 μs   (1.8%)

3.2 μs (51.4%)

# B+ Tree Index Lookup from User Space



**Node parsing and I/O request submission are performed in user space**

**Kernel Boundary (5.6%)**

**Traverse the full kernel software stack multiple times**

User Space

Fetch Root Node

Parse Node

Parse Node

Parse Node

Syscall Layer — 3.2%

File System and Block Layer — 38.0%

NVMe Driver — 1.8%

Storage Device — 51.4%

# B+ Tree Index Lookup With an eBPF Function

**A Chain of Dependent Read Requests:**



User Space

**Kernel Boundary** (5.6%)

Syscall Layer

*Fetch Root Node*

*Return Leaf Node*

3.2%

File System and Block Layer

**Only traverse the full kernel software stack once**

38.0%

NVMe Driver

*Parse Node*

eBPF Function

1.8%

Storage Device

**Reduce the latency of the intermediate I/O by up to 47%**

51.4%

8

# Chains of Dependent Read Requests are Very Common
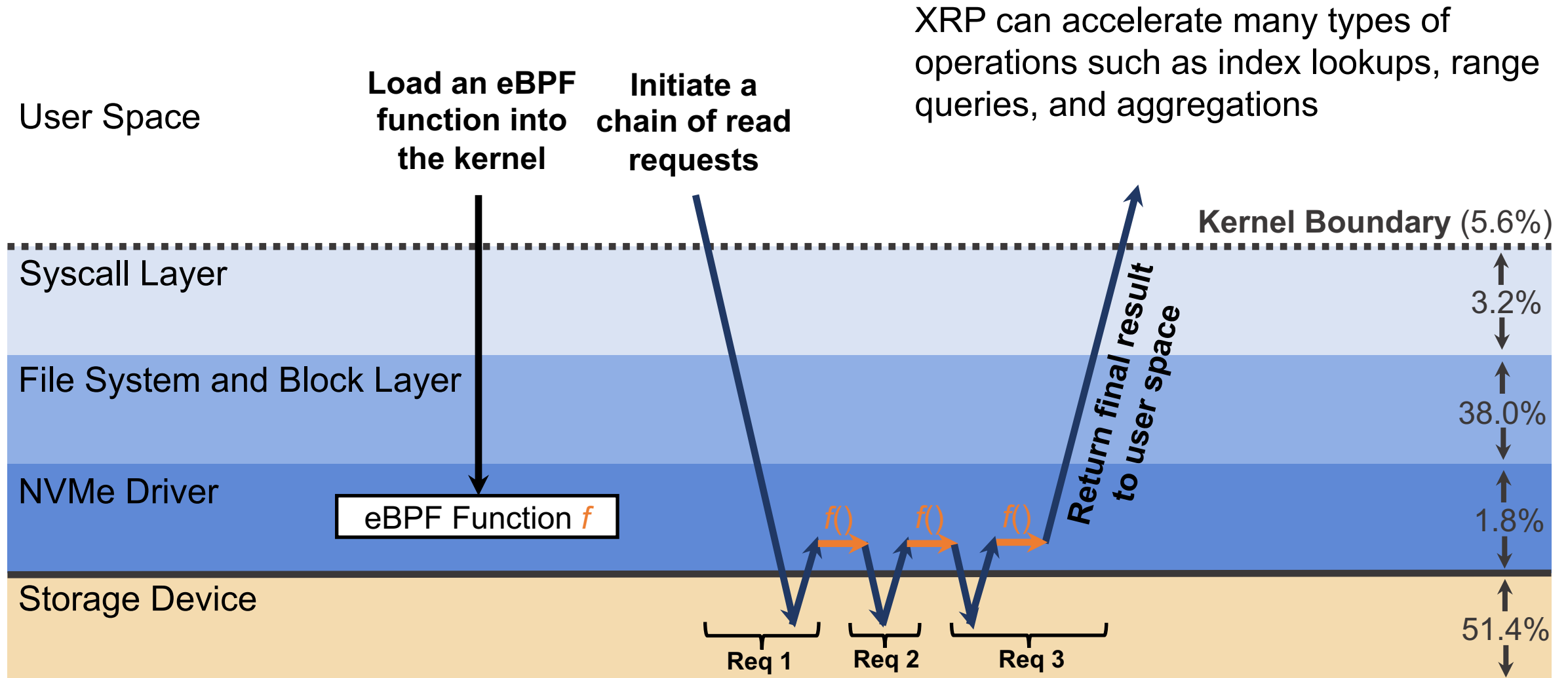
## B-Tree    LSM Tree

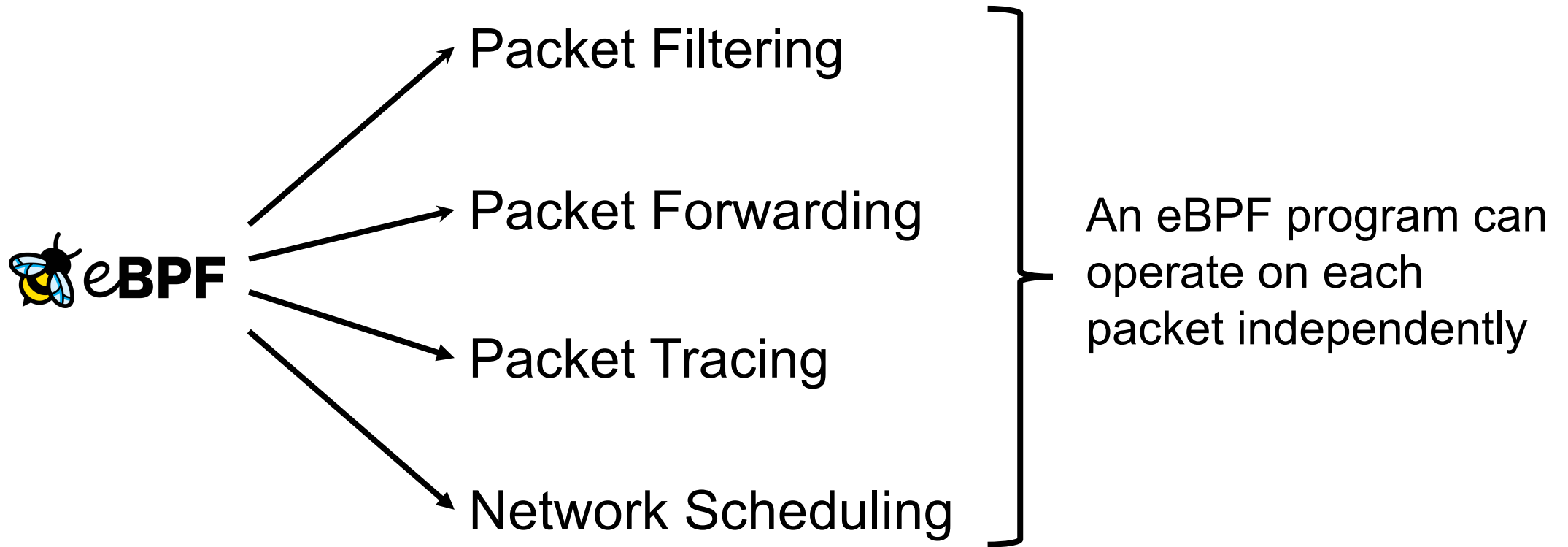Issue dependent read requests to perform lookups

**Goal: Build a framework for storage engines to accelerate dependent read requests using in-kernel functions**

# **XRP**: A Framework for In-Kernel Storage Functions With eBPF

User Space

**Load an eBPF function into the kernel**

**Initiate a chain of read requests**

XRP can accelerate many types of operations such as index lookups, range queries, and aggregations

**Kernel Boundary** (5.6%)

Syscall Layer

3.2%

File System and Block Layer

38.0%

NVMe Driver

eBPF Function *f*

*f()*    *f()*    *f()*

*Return final result to user space*

1.8%

Storage Device

**Req 1**    **Req 2**    **Req 3**

51.4%

10

# eBPF is Widely Used in Networking

Packet Filtering

Packet Forwarding

Packet Tracing

Network Scheduling

An eBPF program can operate on each packet independently

**However, a storage eBPF program needs to traverse a large on-disk data structure in a stateful way**

# Adopting eBPF in Storage is Challenging

**XRP is the first system that adopts eBPF to reduce the kernel software overhead for storage**
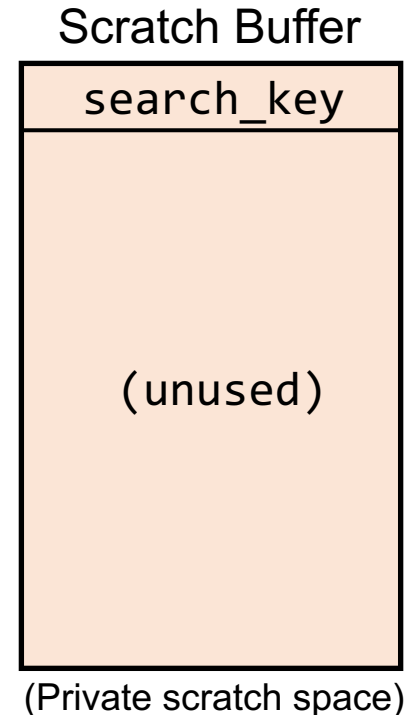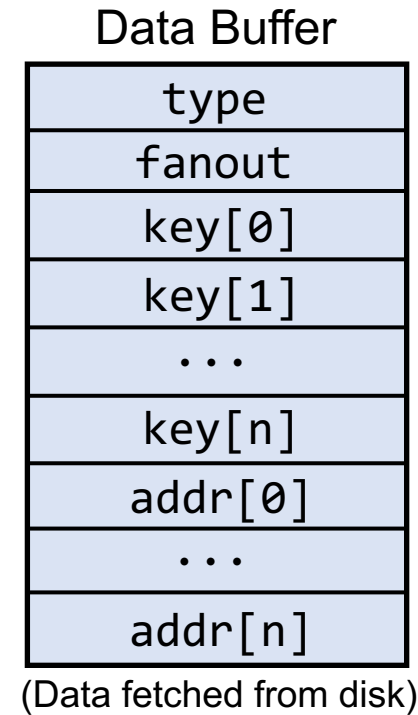
Key research challenges:

- Translating file offsets in the NVMe driver

- Augmenting the BPF verifier to support storage use cases

- Resubmitting NVMe requests

- Interaction with application-level caches

# eBPF Can Traverse Different Types of Data Structures
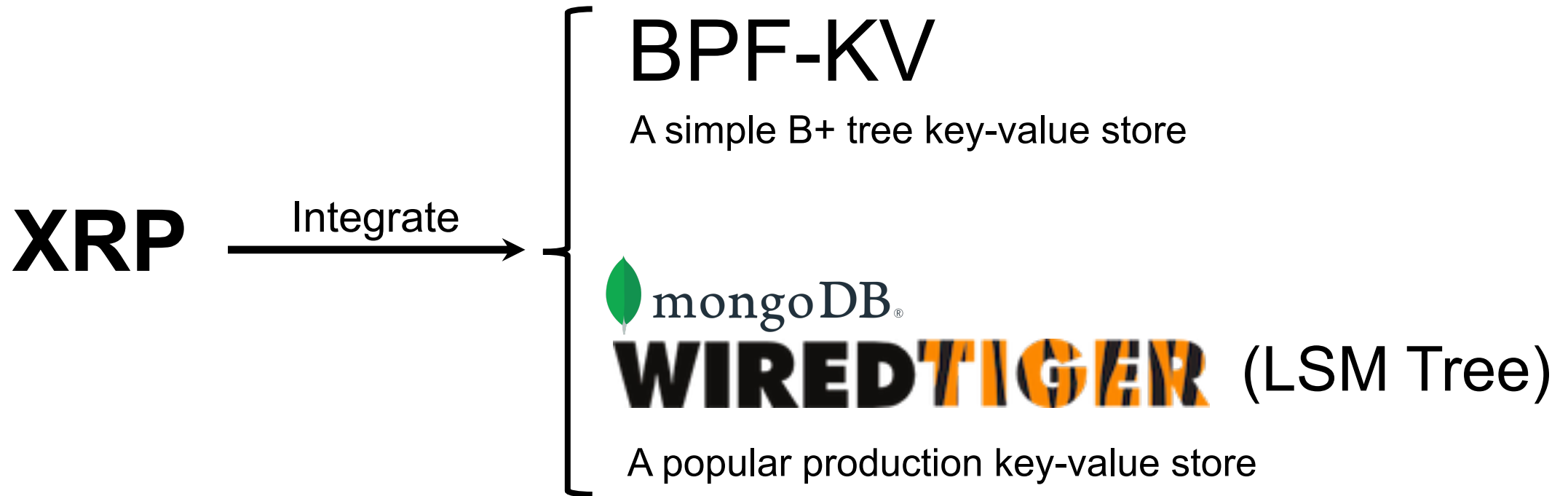
```
u32 btree_lookup(struct bpf_xrp *context) {
    struct node *n = (struct node *) context->data;
    u64 search_key = *(u64 *) context->scratch;
    if (node->type == LEAF) {
        context->done = true;
        return 0;
    }
    int i;
    for (i = 1; i < MIN(n->fanout, MAX_FANOUT); ++i) {
        if (search_key < n->key[i]) break;
    }
    context->done = false;
    context->next_addr[0] = n->addr[i - 1];
    return 0;
}
```

MAX_FANOUT ensures for loop is bounded

**Data Buffer**

| type |
|------|
| fanout |
| key[0] |
| key[1] |
| ... |
| key[n] |
| addr[0] |
| ... |
| addr[n] |

(Data fetched from disk)

**Scratch Buffer**

| search_key |
|------------|
| (unused) |

(Private scratch space)

# Integration of XRP and Key-Value Stores

**XRP** → Integrate →

## BPF-KV
A simple B+ tree key-value store

mongoDB

**WIREDTIGER** (LSM Tree)

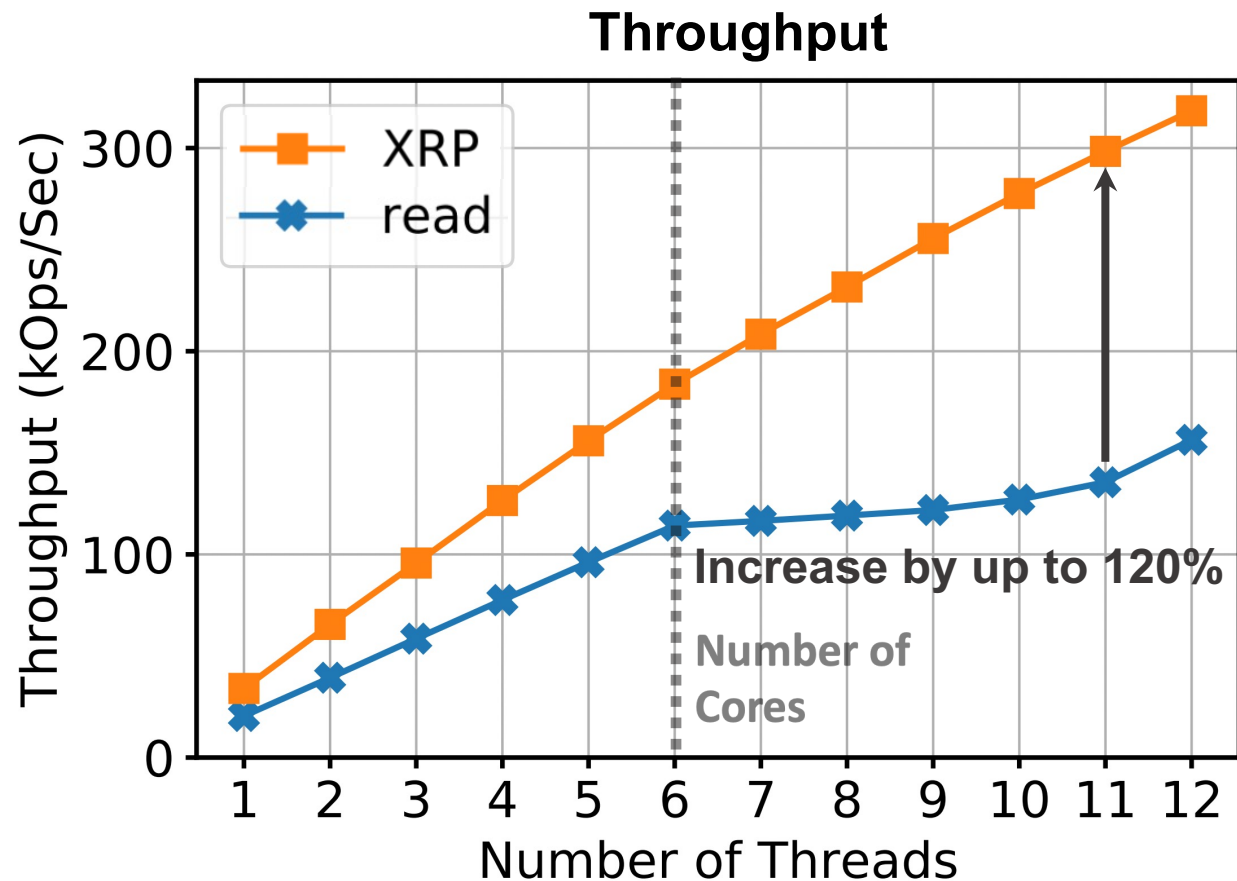A popular production key-value store

# Evaluation

- What is the performance benefit of XRP?

- How does XRP compare to kernel bypass?

- What types of operations can XRP support?

- Can XRP accelerate a production key-value store?

See our paper: XRP: In-Kernel Storage Functions with eBPF (OSDI '22)
https://www.usenix.org/conference/osdi22/presentation/zhong

# XRP Nearly Eliminates the Kernel Software Overhead

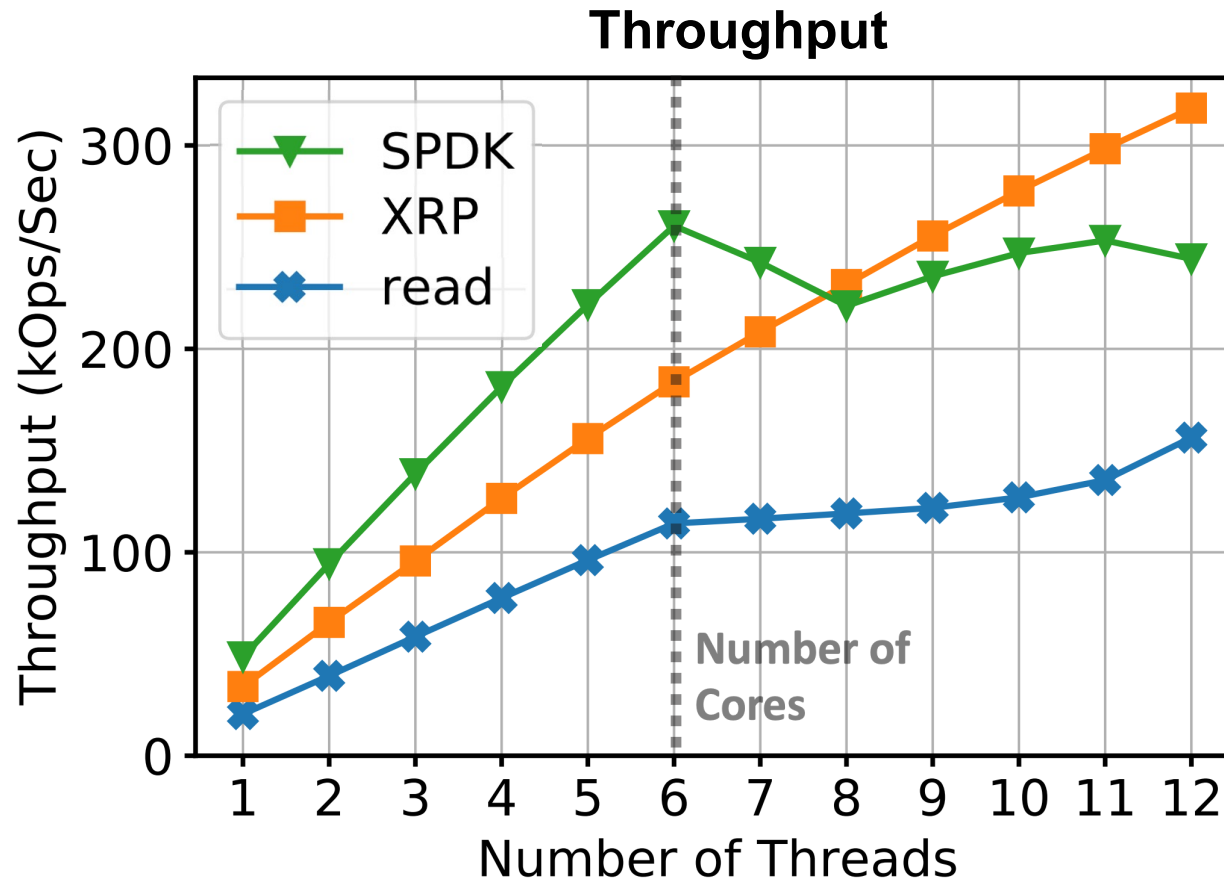Multi-threaded throughput in BPF-KV with uniform random 512B read:

**Throughput**



XRP can scale well even if the number of threads exceeds the number of cores

This is because XRP alleviates the CPU contention by reducing the CPU overhead per IO request

# XRP Handles CPU Contention, SPDK Not So Much

Multi-threaded throughput in BPF-KV with uniform random 512B read:

**Throughput**



SPDK fails to scale beyond 6 threads because SPDK threads cannot yield CPU when waiting for I/O to complete
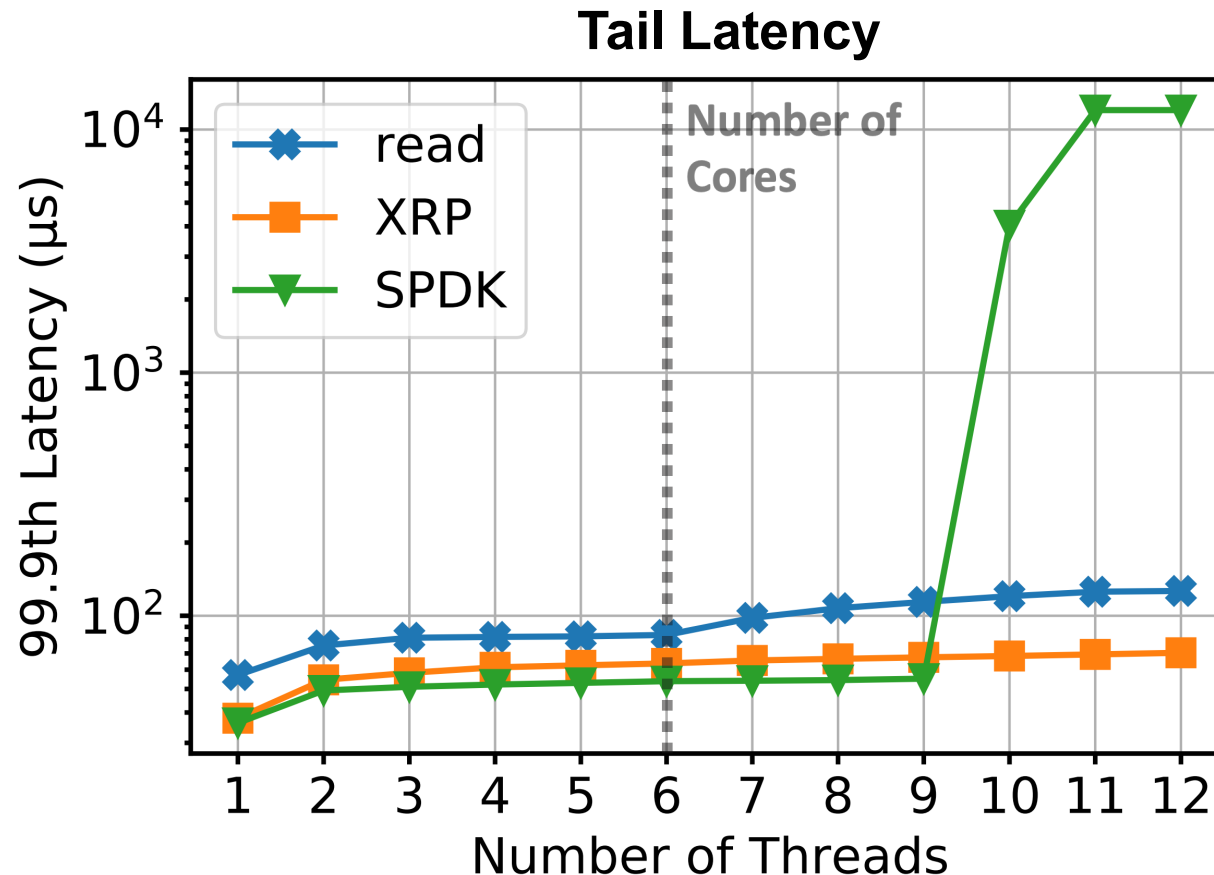
**XRP provides performance that is close to/better than SPDK without sacrificing isolation and CPU efficiency**

Each thread represents a different storage application on the same machine

# XRP Handles CPU Contention, SPDK Not So Much

Multi-threaded tail latency in BPF-KV with uniform random 512B read:

**Tail Latency**



Compared to read, XRP improves tail latency by up to 45%

Tail latency of SPDK spikes to ~10 ms when the number of threads is greater than the number of cores by more than 50%

# Conclusions

- XRP is the first system to use eBPF to accelerate common storage functions

- XRP captures most of the performance benefit of kernel bypass, without sacrificing CPU utilization and access control

We are actively integrating XRP with other popular key-value stores including RocksDB

Paper: https://www.usenix.org/conference/osdi22/presentation/zhong

Source Code: http://xrp-project.com/

Email: yz@cs.columbia.edu