# A Memory Pool Allocator for eBPF Applications

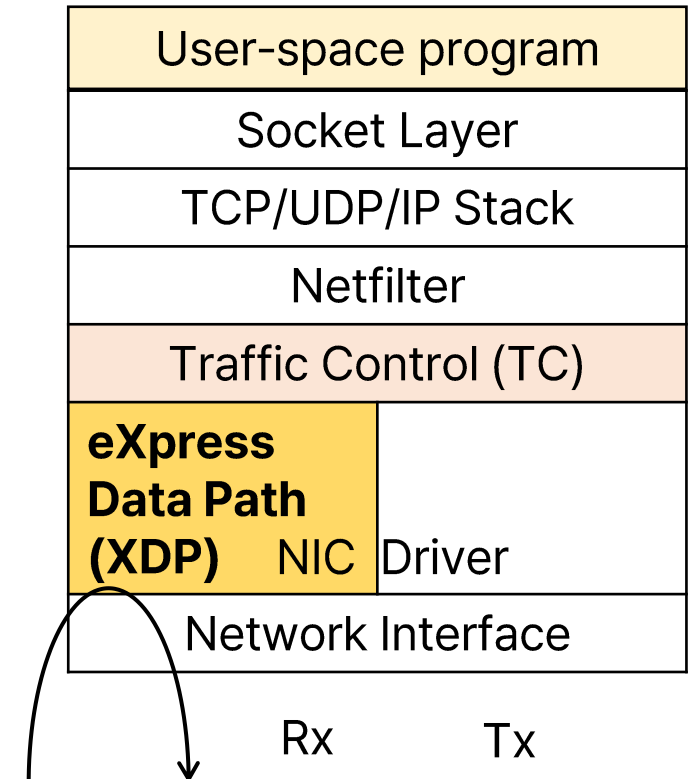**Gyuyeong Kim** and Dongsu Han

# In-Kernel Offloading with eBPF

eBPF enables kernel-level execution of application logic

- High performance by avoding networking stacks
- Safety guarantee through static verification

Key-value stores [BMC@NSDI'21, DINT-KV@NSDI'24]

Consensus protocol [Electrode@NSDI'23]

Lock managers [DINT-Lock@NSDI'24]

| User-space program | |
|---|---|
| Socket Layer | |
| TCP/UDP/IP Stack | |
| Netfilter | |
| Traffic Control (TC) | |
| **eXpress Data Path (XDP)** NIC | Driver |
| Network Interface | |

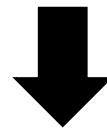Rx      Tx

# Memory Allocation Matters!

eBPF typically relies on **static memory allocation** for safety

Cannot allocate **variable-size** memory space

Runtime allocation is only possible for restricted cases

Critical for apps. with variable-size data or dynamic data structures
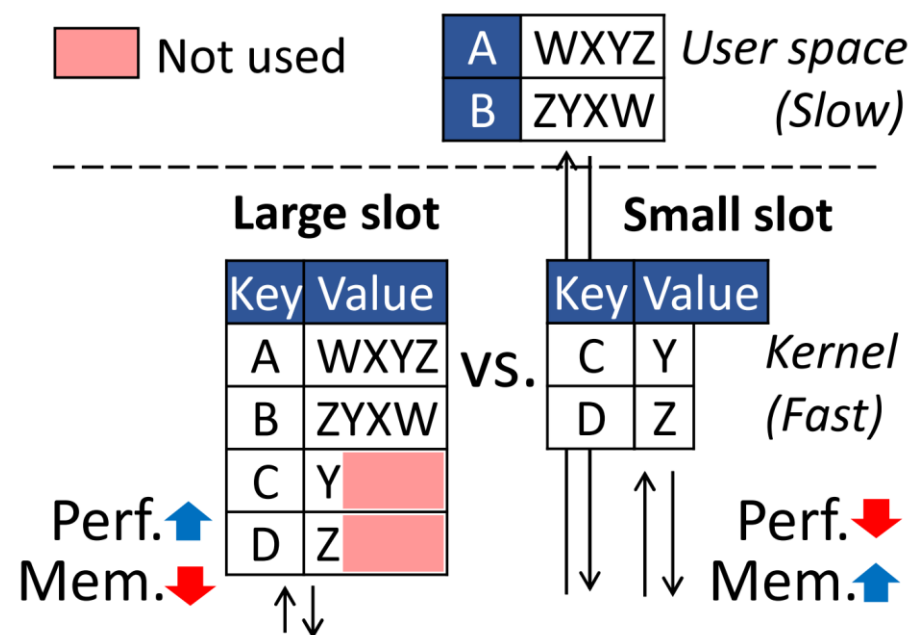


**User-space apps.**

**eBPF applications**

Over-provisioned

# Trade-off Between Perf. And Mem. Efficiency

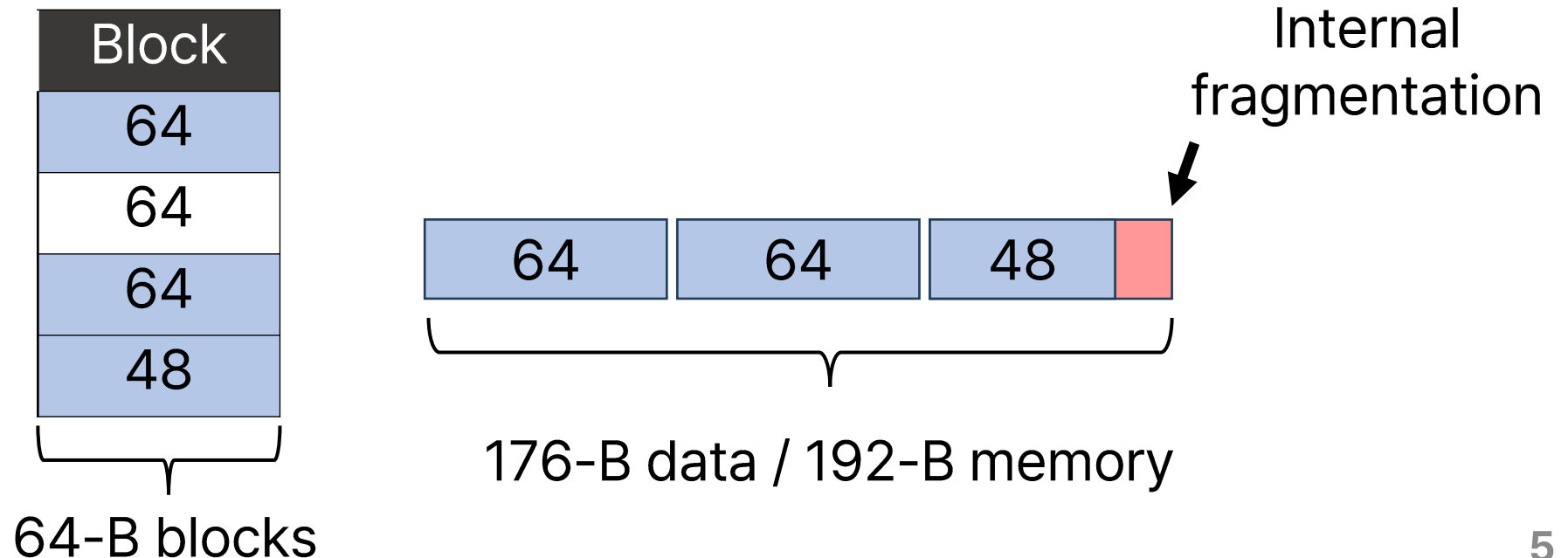A large slot provides better performance but wastes memory

A small slot saves memory but degrades performance



BMC@NSDI'21
DINT-KV@NSDI'24

# Kerby: A Memory Pool Allocator for eBPF

- **Idea:** dynamically manages a memory pool consisting of *fixed-size* memory blocks
  - Combines blocks to represent variable-size data
  - Fragmentation only happens in the tail block

Internal fragmentation

Block

| 64 |
| 64 |
| 64 |
| 48 |

64-B blocks

64 | 64 | 48 |

176-B data / 192-B memory

5

# Free List Management

The allocator should know which block is allocated or free

Bitmap-based free list?

We cannot call functions while holding a lock in eBPF for safety reasons

Block Index

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
|---|---|---|---|---|---|---|---|---|
| 1 | 0 | 0 | 1 | 1 | 0 | 1 | 0 | 1 |

Allocated?

```
lock();
Update_freelist();
unlock();
```

Prohibited!

# No Free List in Kerby

1. Monotonically-increasing index up to $2^{64}$-1

2. Access within the bound using a BPF hash map for the pool
   - Internally hashses index, resolves hash collisions using chaining
   - `BPF_F_NO_PREALLOC` flag enables runtime allocation
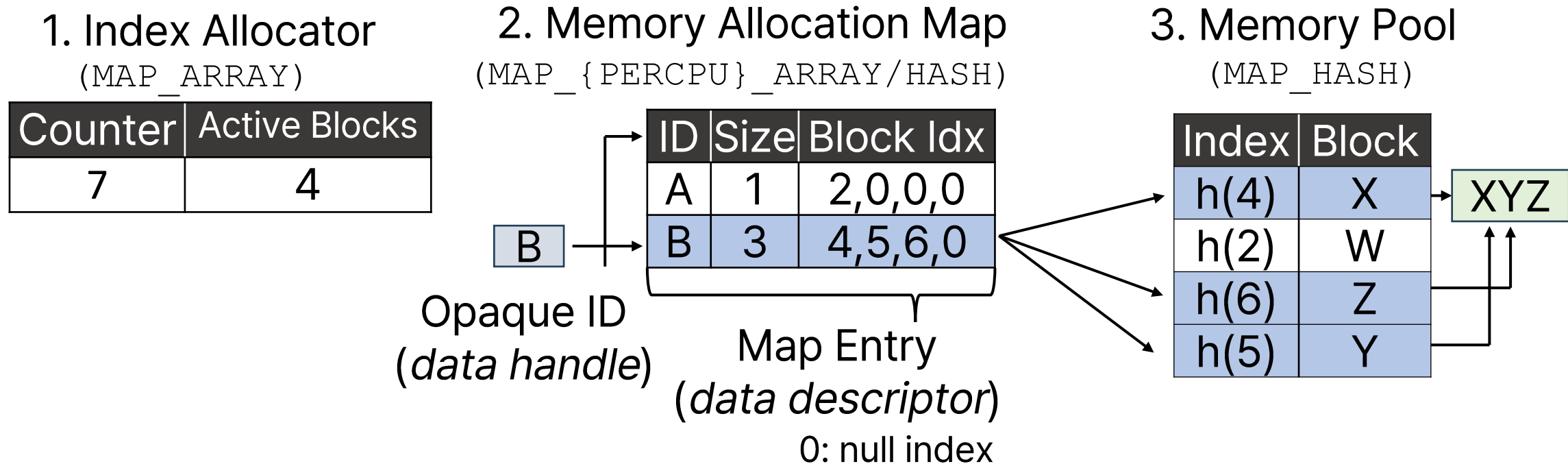
Index Allocator
(`MAP_ARRAY`)

| Counter |
|---------|
| 64-bit  |

Memory Pool
(`MAP_HASH`)

| Index | Block |
|-------|-------|
| h(2) | A |
| h(1) | D |
| h(3) | B |
| h($2^{64} - 1$) | C |

# Kerby Overview

Kerby consists of three components



1. Index Allocator (`MAP_ARRAY`)

| Counter | Active Blocks |
|---------|---------------|
| 7 | 4 |

2. Memory Allocation Map (`MAP_{PERCPU}_ARRAY/HASH`)

| ID | Size | Block Idx |
|----|------|-----------|
| A | 1 | 2,0,0,0 |
| B | 3 | 4,5,6,0 |

Opaque ID (*data handle*)

Map Entry (*data descriptor*)

0: null index

3. Memory Pool (`MAP_HASH`)

| Index | Block |
|-------|-------|
| h(4) | X |
| h(2) | W |
| h(6) | Z |
| h(5) | Y |

XYZ
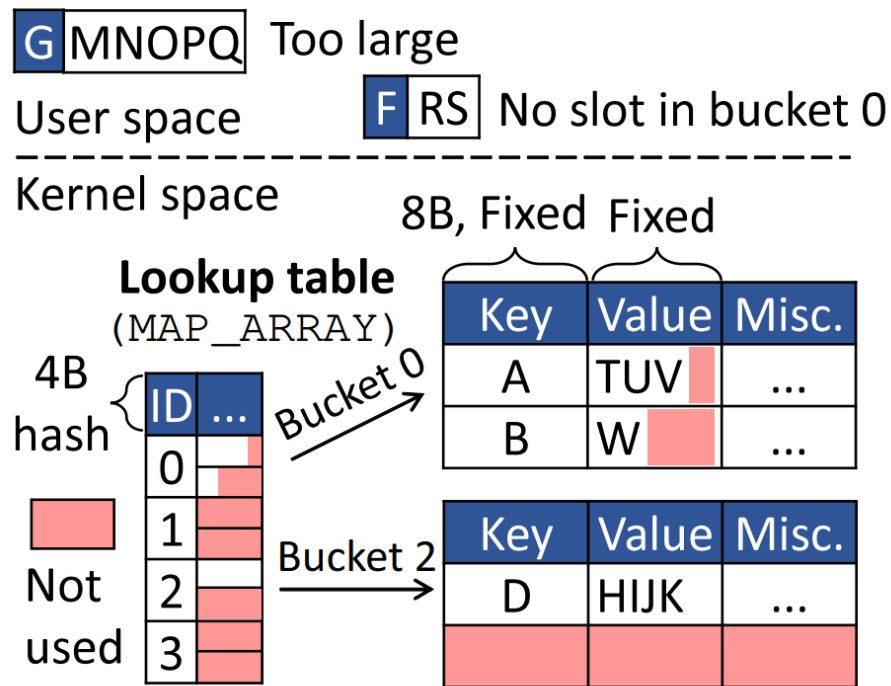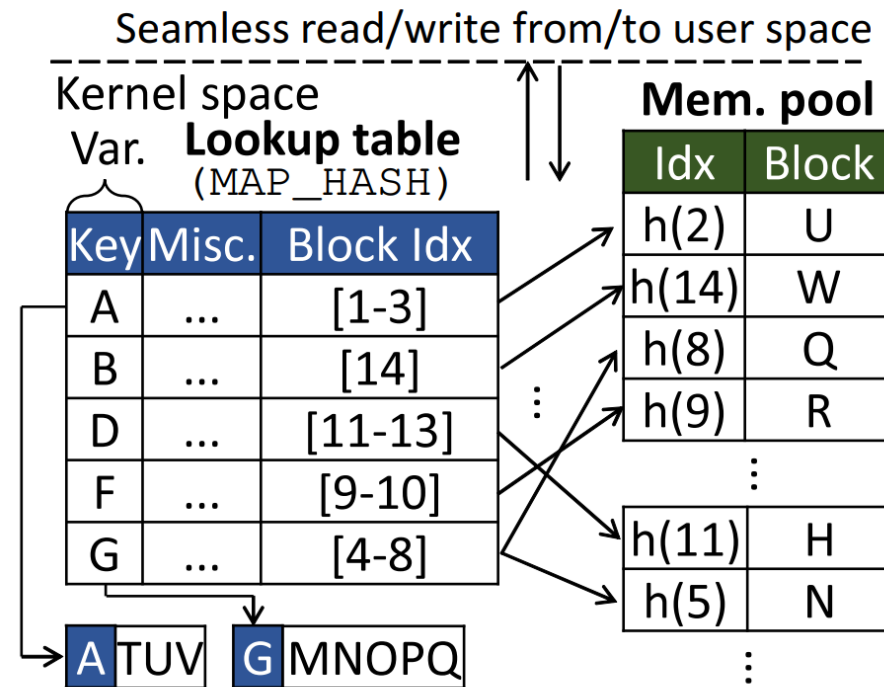
# Key-Value Stores: Design

DINT-KV supports fixed-size items
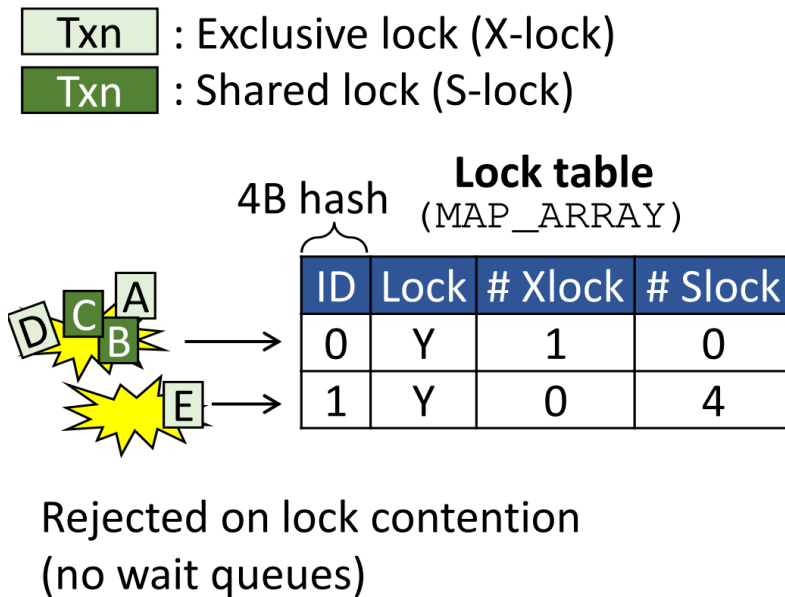
Kerby-KV supports variable-size items
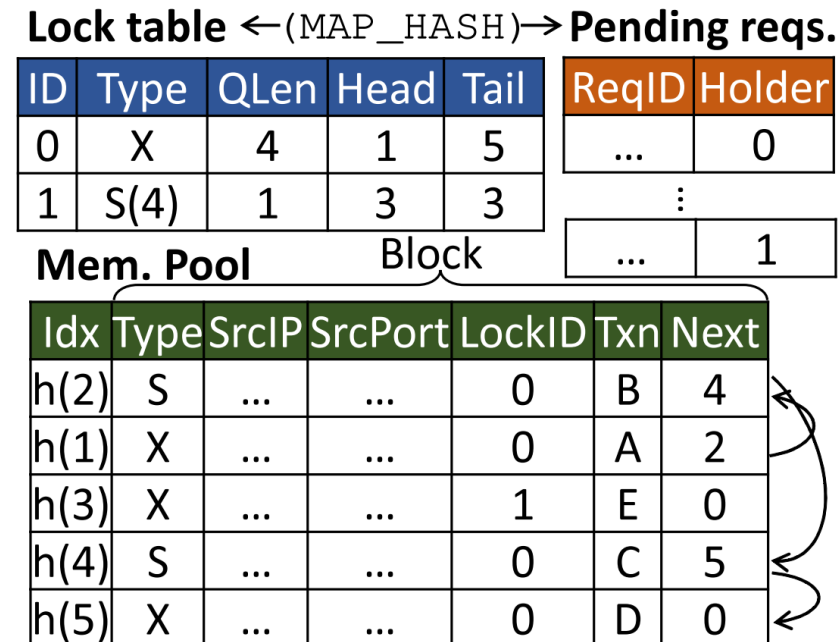


(a) DINT-KV

(b) Kerby-KV

# Lock Managers: Design

DINT-Lock lacks wait queues, cannot handle lock contention

Kerby-Lock supports dynamic variable-length wait queues by linking memory blocks (dynamic linked-list)



| Txn | : Exclusive lock (X-lock) |
| Txn | : Shared lock (S-lock) |

**Lock table** (MAP_ARRAY)

4B hash

| ID | Lock | # Xlock | # Slock |
|----|------|---------|---------|
| 0 | Y | 1 | 0 |
| 1 | Y | 0 | 4 |

Rejected on lock contention
(no wait queues)

(a) DINT-Lock

**Lock table** ←(MAP_HASH)→ **Pending reqs.**

| ID | Type | QLen | Head | Tail |
|----|------|------|------|------|
| 0 | X | 4 | 1 | 5 |
| 1 | S(4) | 1 | 3 | 3 |

| ReqID | Holder |
|-------|--------|
| ... | 0 |
| : | |
| ... | 1 |

**Mem. Pool**  Block

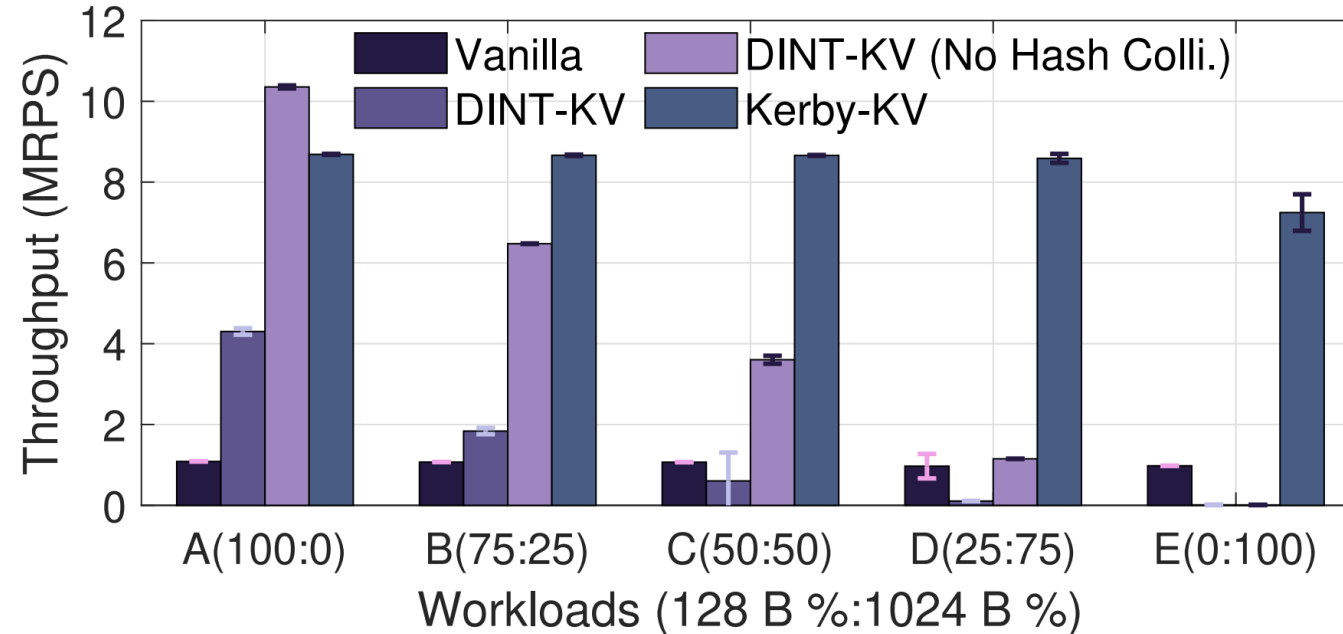| Idx | Type | SrcIP | SrcPort | LockID | Txn | Next |
|-----|------|-------|---------|--------|-----|------|
| h(2) | S | ... | ... | 0 | B | 4 |
| h(1) | X | ... | ... | 0 | A | 2 |
| h(3) | X | ... | ... | 1 | E | 0 |
| h(4) | S | ... | ... | 0 | C | 5 |
| h(5) | X | ... | ... | 0 | D | 0 |

(b) Kerby-Lock

10

# Kerby APIs

We expose developer-friendly APIs

Three primitives and extra helpers
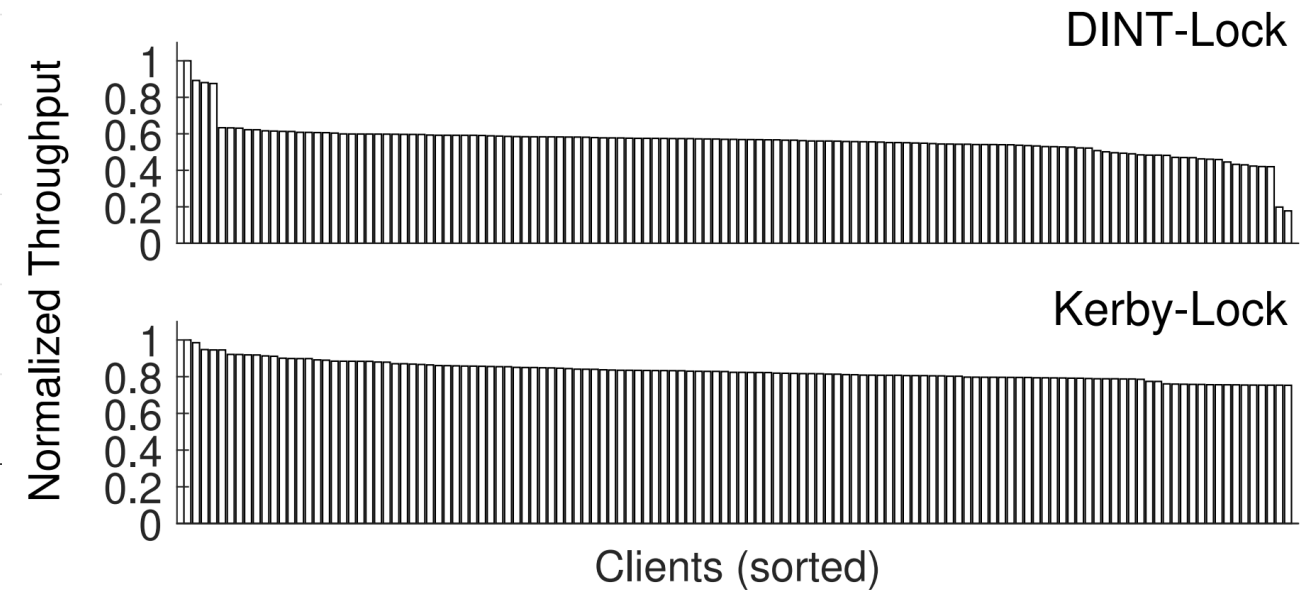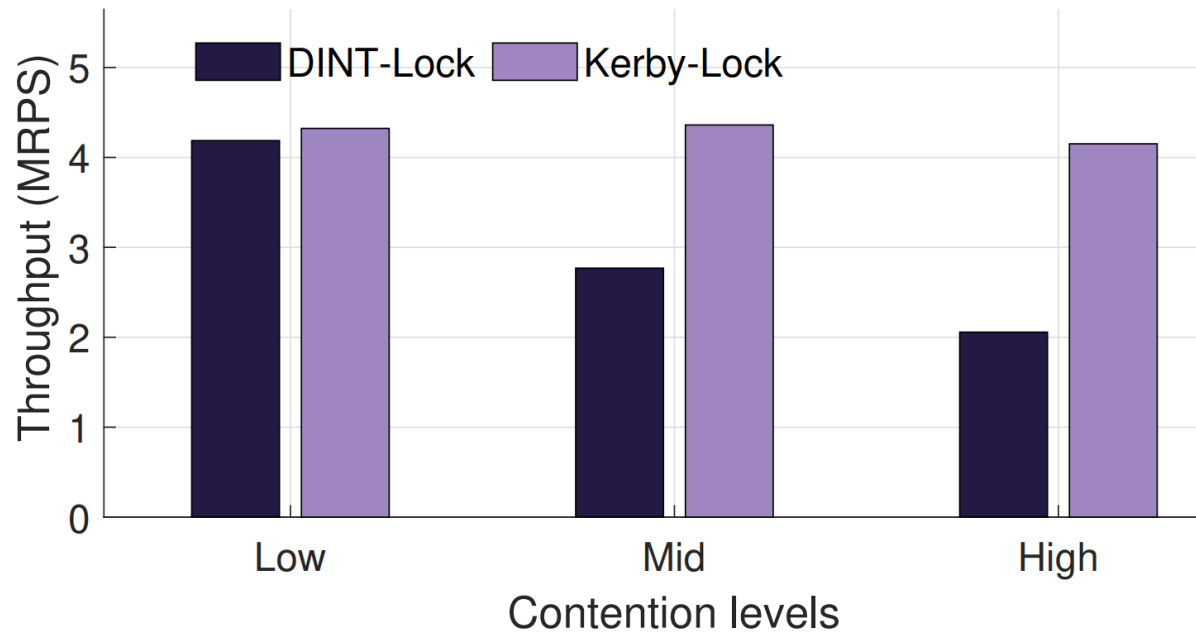
Using only BPF maps and helpers, no kernel modifications

| Abstract Signature | Category | Role |
|---|---|---|
| malloc(size, idx[]) | Primitive | Allocate memory |
| free(idx[]) | Primitive | Release memory |
| resize(size, idx[]) | Primitive | Adjust allocation |
| pool_update(idx[], values) | Helper | Write blocks |
| alloc_update(key, idx[], size) | Helper | Update alloc. map |
| alloc_del(key) | Helper | Remove alloc. map |

# Key-Value Stores: Evaluation Results



Kerby-KV performs the best
DINT-KV spills 1024-B items to user space

# Lock Managers: Evaluation Results



**Kerby-Lock can handle high lock contention and achieves fairness**

# Conclusion

- Kerby is a dynamic eBPF memory pool allocator to overcome the limitation of static memory allocation in eBPF

- Kerby enables us to implement high-performance and memory-efficient eBPF applications through memory pooling

- Kerby can enrich existing and future eBPF applications
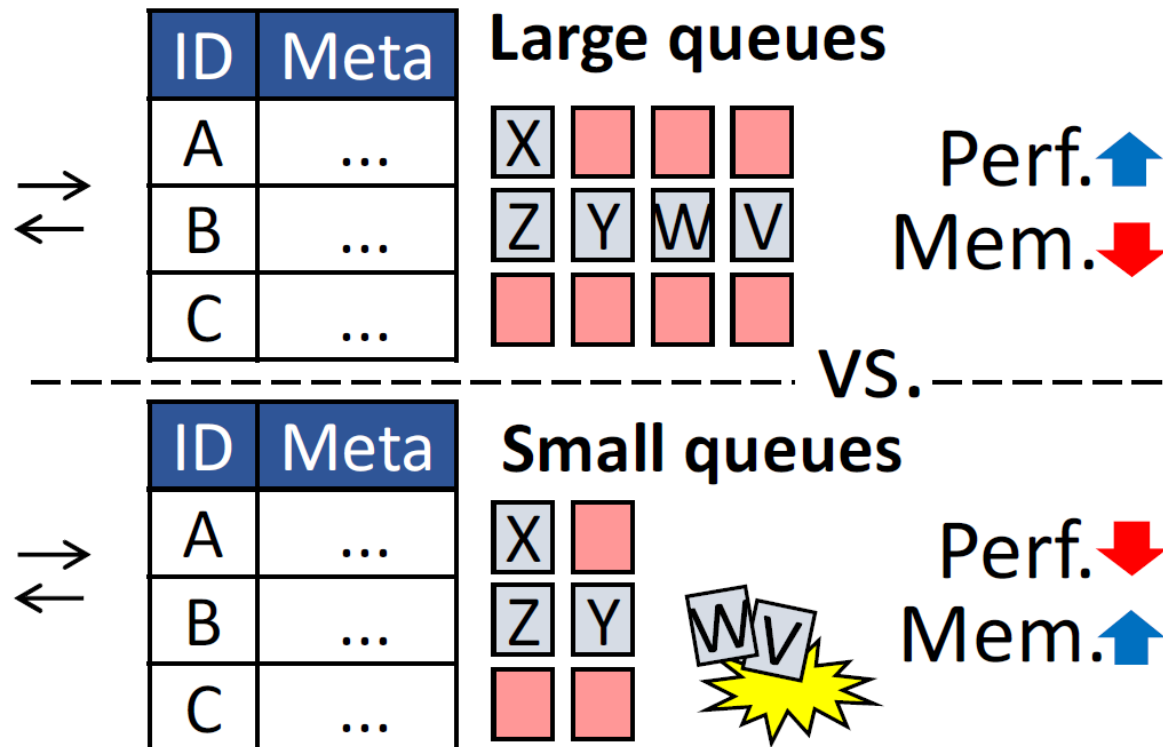
# Thank you!

Questions?

# Trade-off in Lock Managers

DINT-Lock@NSDI'24 lacks wait queues

We may implement wait queues using BPF queue maps

But, still has a trade-off between perf. and mem. efficiency

# Allocation Workflow

Allocator step

Application steps

Mem. Pool

## Idx Allocator

| Idx | Active |
|-----|--------|
| 2→4 | 0→2 |

| Idx | Block |
|------|-------|
| h(2) | AB |
| | |
| | |
| | |
| h(3) | C |

## Alloc. Map

| ID | Size | BlckIdx |
|----|------|---------|
| A | 3 B | 2,3,0,0 |
| | | |

❸ AllocUpdate(A,[2,3],3)

❶ malloc(3,[])

❷ PoolUpdate([2,3],ABC)

# Deallocation Workflow



Application step

Alloc. Map

| ID | Size | BlckIdx |
|----|------|---------|
| A | 3 | **X** 3,0,0 |
| | | |

Mem. Pool

| Idx | Block |
|-----|-------|
| h( **X** | AB |
| | |
| | |
| | |
| h( **X** | C |

Allocator step

Idx Allocator

| Idx | Active |
|-----|--------|
| 4 | 2→0 |

❶ AllocDel(A)          ❷ free([2,3])

# Resizing Workflow

**Allocator step**

**Application steps**

Idx Allocator

| Idx | Active |
|-----|--------|
| 5→6 | 3→4 |

Mem. Pool

| Idx | Block |
|------|-------|
| h(2) | DE |
| h(5) | HJ |
| h(4) | XY |
| h(3) | FG |

Alloc. Map

| ID | Size | BlckIdx |
|----|------|---------|
| A | 6 B | 2,3,5,0 |
| B | 2 B | 4,0,0,0 |

❸ AllocUpdate(A,[2,3,5],6)

❶ resize(6, [2,3])    ❷ PoolUpdate([2,3,5],DEFGHJ)