

# Reaping the performance of fast NVM storage with uDepot

Kornilios Kourtis, Nikolas Ioannou, and Ioannis Koltsidas

IBM Research, Zurich

FAST<sup>19</sup>

# Key-Value (KV) stores

- ▶ Many applications require low-latency high-throughput KV storage
  - ▶ Flash-based SSDs not performant enough
  - ▶ Most are using DRAM KV-stores (Memcache, Redis)

# Key-Value (KV) stores

- ▶ Many applications require low-latency high-throughput KV storage
  - ▶ Flash-based SSDs not performant enough
  - ▶ Most are using DRAM KV-stores (Memcache, Redis)
- ▶ DRAM performance underutilized on commodity networks (e.g., 10GbE)
  - ▶ High-performance DRAM KV stores use: RDMA (RamCloud, FaRM), Direct NIC access (MICA), programmable NICs (KV-Direct).

# Key-Value (KV) stores

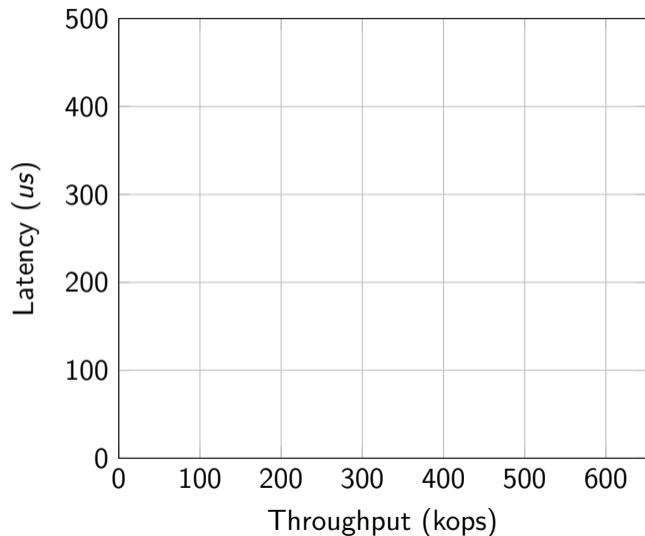
- ▶ Many applications require low-latency high-throughput KV storage
  - ▶ Flash-based SSDs not performant enough
  - ▶ Most are using DRAM KV-stores (Memcache, Redis)
- ▶ DRAM performance underutilized on commodity networks (e.g., 10GbE)
  - ▶ High-performance DRAM KV stores use: RDMA (RamCloud, FaRM), Direct NIC access (MICA), programmable NICs (KV-Direct).
- ▶ DRAM is not getting cheaper

# Fast NVM Devices

(FNDs)

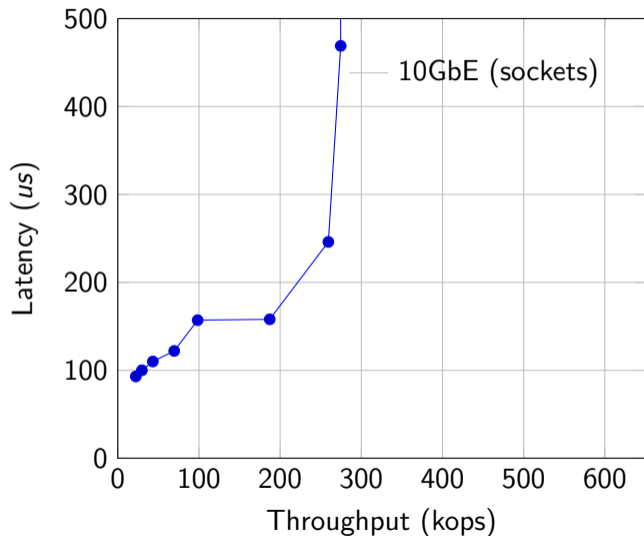
- ▶ new class of SSDs
  - ▶ Intel Optane (3DXP)
  - ▶ Samsung Z-SSD (Z-NAND)
- ▶ An order of magnitude better performance than Flash SSDs
- ▶ Significantly cheaper than DRAM
  - ▶ \$1.25 vs \$10 per GiB (Intel Optane)
  - ▶ smaller TCO (number of machines, energy, etc.)

# 10GiB network vs Flash SSD vs FND SSD



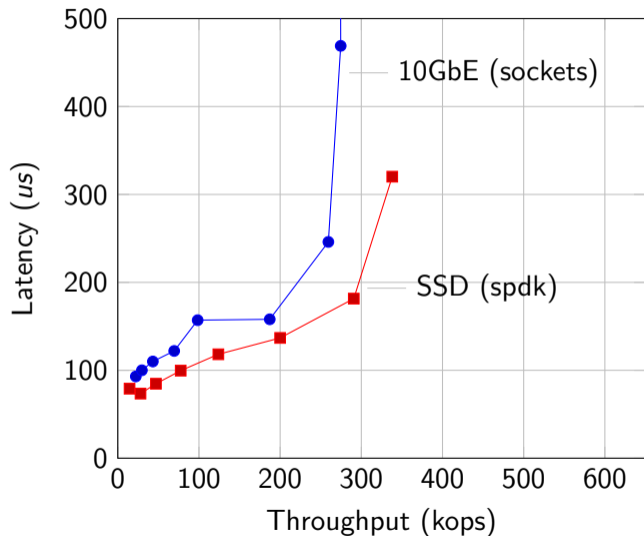
- ▶ queue depth (reqs in flight):  
1, 2, 4, 8, 16, ...

# 10GiB network vs Flash SSD vs FND SSD



- ▶ queue depth (reqs in flight): 1, 2, 4, 8, 16, ...
- ▶ *10GbE*: req:1b, res:4KiB (Intel X710, netperf)

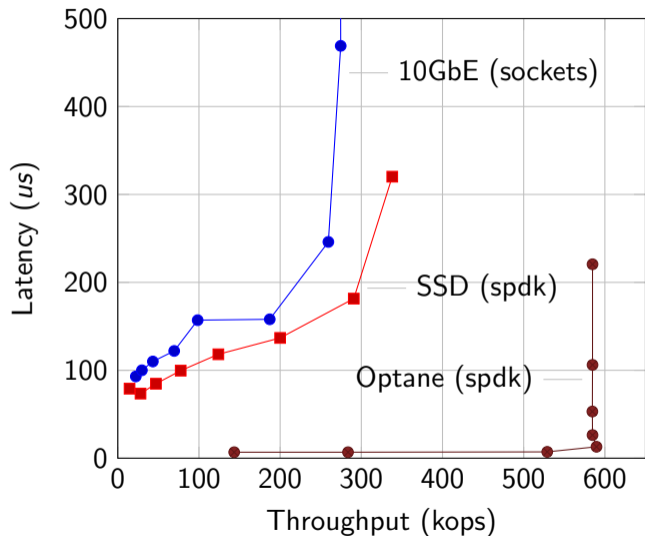
# 10GiB network vs Flash SSD vs FND SSD



- ▶ queue depth (reqs in flight): 1, 2, 4, 8, 16, ...
- ▶ *10GbE*: req:1b, res:4KiB (Intel X710, netperf)
- ▶ *SSD*: SPDK perf: 4KiB RDs (Flash NVMe)



# 10GiB network vs Flash SSD vs FND SSD



- ▶ queue depth (reqs in flight): 1, 2, 4, 8, 16, ...
  - ▶ *10GbE*: req:1b, res:4KiB (Intel X710, netperf)
  - ▶ *SSD*: SPDK perf: 4KiB RDs (Flash NVMe)
  - ▶ *Optane*: SPDK perf: 4KiB RDs ( $\approx 0.6\text{Mops/sec}$ ,  $\approx 10\text{us}$ )
- Storage no longer the bottleneck!

# KV store for DRAM FNDs

- ▶ reduced cost
- ▶ equivalent performance to DRAM KV store (at least, under commodity networks)

# KV store for DRAM FNDs

- ▶ reduced cost
- ▶ equivalent performance to DRAM KV store (at least, under commodity networks)

## **Existing KV stores cannot deliver FNDs' performance**

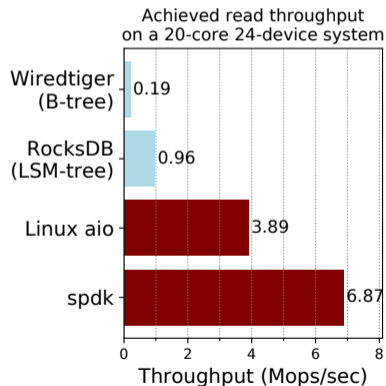
- ▶ Built for slower devices  
(e.g., use synchronous IO)
- ▶ Data structures with inherent IO amplification  
(LSM- or B-trees)
- ▶ Cache data in DRAM, limiting scalability
- ▶ Rich feature set (e.g., transactions, snapshots)

# KV store for DRAM FNDs

- ▶ reduced cost
- ▶ equivalent performance to DRAM KV store (at least, under commodity networks)

## Existing KV stores cannot deliver FNDs' performance

- ▶ Built for slower devices (e.g., use synchronous IO)
- ▶ Data structures with inherent IO amplification (LSM- or B-trees)
- ▶ Cache data in DRAM, limiting scalability
- ▶ Rich feature set (e.g., transactions, snapshots)

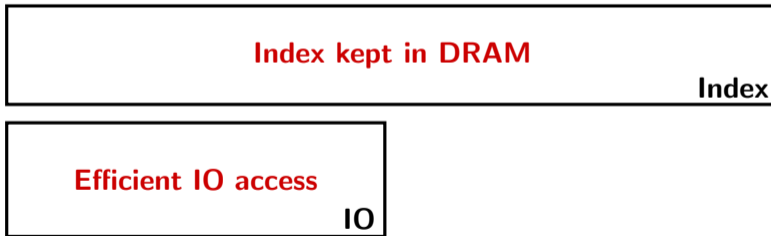


## **Deliver the performance of FNDs to the application**

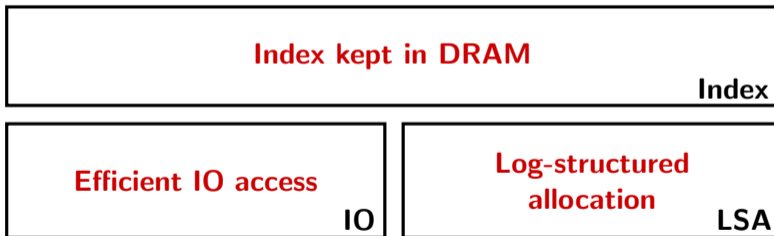
- ▶ minimize IO amplification
- ▶ scalability (cores, devices, capacity)
- ▶ bottom-up approach
  - ▶ basic interface: GET, PUT, DEL on variable-sized keys and values.

# uDepot architecture

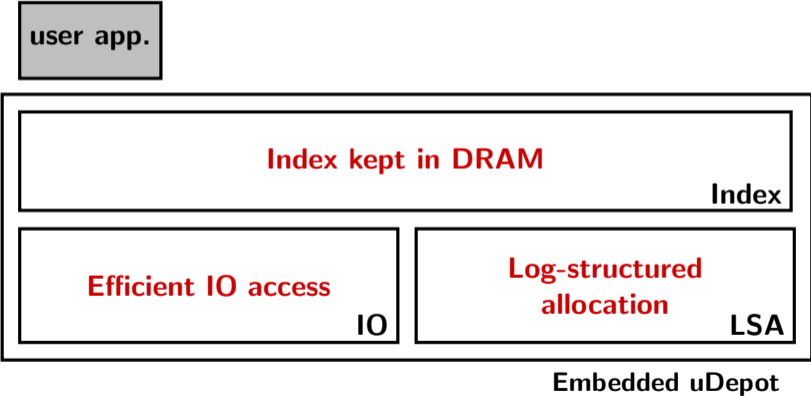




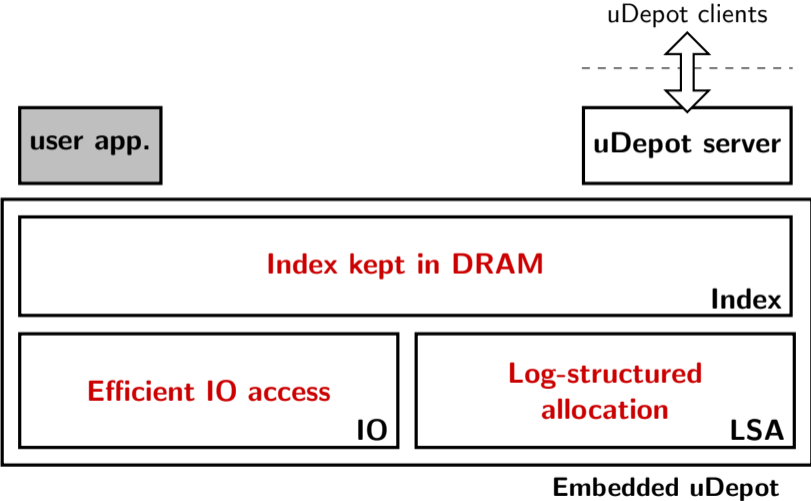




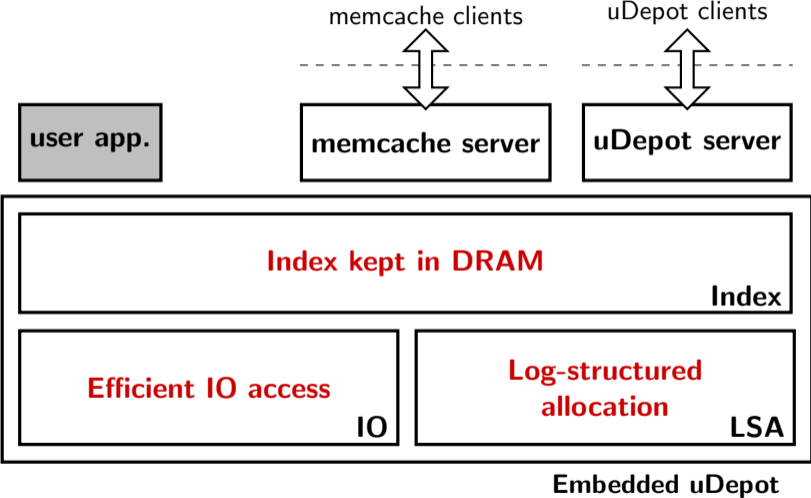
# uDepot architecture



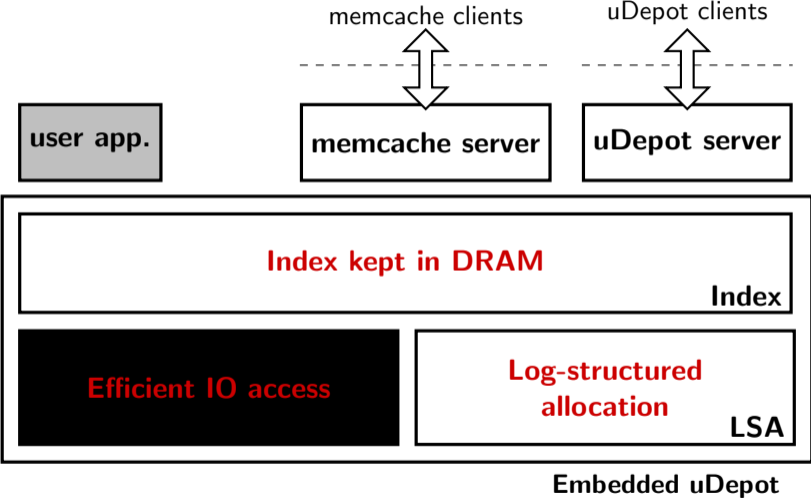
# uDepot architecture



# uDepot architecture



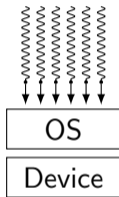
# uDepot architecture



Performance →

Performance →

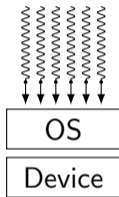
## Synchronous IO



- ▶ one thread per request
- ▶ synchronous (blocking) syscalls (e.g., `pread`)

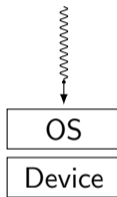
Performance →

## Synchronous IO



- ▶ one thread per request
- ▶ synchronous (blocking) syscalls (e.g., `pread`)

## Asynchronous IO

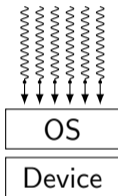


- ▶ issue IO requests
- ▶ receive IO completions
- ▶ e.g., Linux AIO



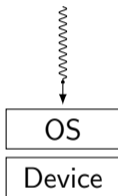
Performance →

## Synchronous IO



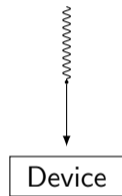
- ▶ one thread per request
- ▶ synchronous (blocking) syscalls (e.g., `pread`)

## Asynchronous IO



- ▶ issue IO requests
- ▶ receive IO completions
- ▶ e.g., Linux AIO

## User-space IO

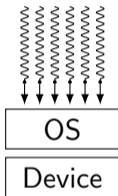


- ▶ Directly access the device
- ▶ Polling instead of interrupts
- ▶ e.g., SPDK

# IO Facilities

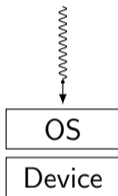
Performance

## Synchronous IO



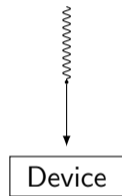
- ▶ one thread per request
- ▶ synchronous (blocking) syscalls (e.g., `pread`)
- ▶ uDepot Linux backend

## Asynchronous IO



- ▶ issue IO requests
- ▶ receive IO completions
- ▶ e.g., Linux AIO
- ▶ uDepot aio backend

## User-space IO

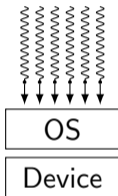


- ▶ Directly access the device
- ▶ Polling instead of interrupts
- ▶ e.g., SPDK
- ▶ uDepot spdk backend

# IO Facilities

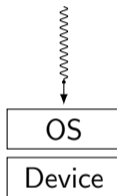
Performance →

## Synchronous IO



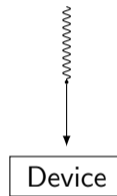
- ▶ one thread per request
- ▶ synchronous (blocking) syscalls (e.g., `pread`)
- ▶ uDepot Linux backend

## Asynchronous IO



- ▶ issue IO requests
- ▶ receive IO completions
- ▶ e.g., Linux AIO
- ▶ uDepot aio backend

## User-space IO



- ▶ Directly access the device
- ▶ Polling instead of interrupts
- ▶ e.g., SPDK
- ▶ uDepot spdk backend

TRT: a run-time system for async IO

# TRT: A task run-time system for asynchronous IO

## Goals

- ▶ programmer-friendly  
(e.g., avoid stack ripping)
- ▶ Handle multiple IO endpoints  
(e.g., SPDK and epoll)

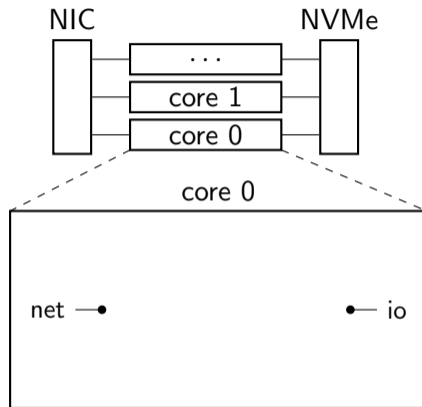
# TRT: A task run-time system for asynchronous IO

## Goals

- ▶ programmer-friendly (e.g., avoid stack ripping)
- ▶ Handle multiple IO endpoints (e.g., SPDK and epoll)

## TRT highlights

- ▶ avoid cross-core communication



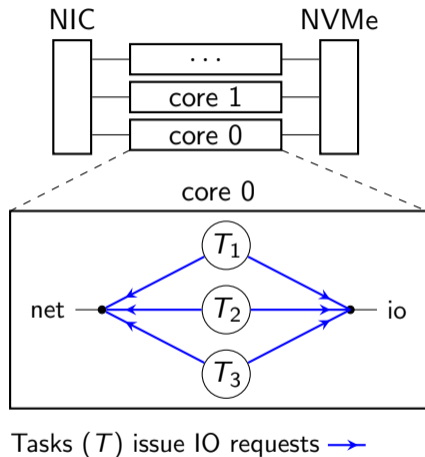
# TRT: A task run-time system for asynchronous IO

## Goals

- ▶ programmer-friendly (e.g., avoid stack ripping)
- ▶ Handle multiple IO endpoints (e.g., SPDK and epoll)

## TRT highlights

- ▶ avoid cross-core communication
- ▶ user-space tasks (aka green threads, aka co-routines)



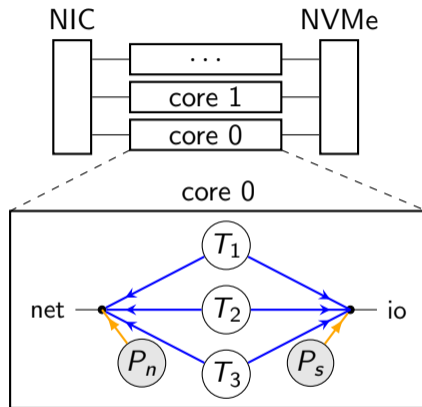
# TRT: A task run-time system for asynchronous IO

## Goals

- ▶ programmer-friendly (e.g., avoid stack ripping)
- ▶ Handle multiple IO endpoints (e.g., SPDK and epoll)

## TRT highlights

- ▶ avoid cross-core communication
- ▶ user-space tasks (aka green threads, aka co-routines)
- ▶ poller task for multiple IO backends



Tasks ( $T$ ) issue IO requests  $\rightarrow$   
Pollers ( $P$ ) poll for completions  $\rightarrow$

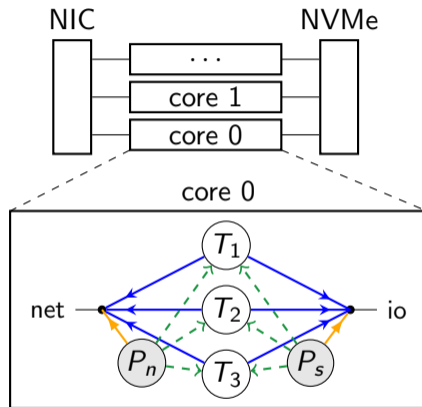
# TRT: A task run-time system for asynchronous IO

## Goals

- ▶ programmer-friendly (e.g., avoid stack ripping)
- ▶ Handle multiple IO endpoints (e.g., SPDK and epoll)

## TRT highlights

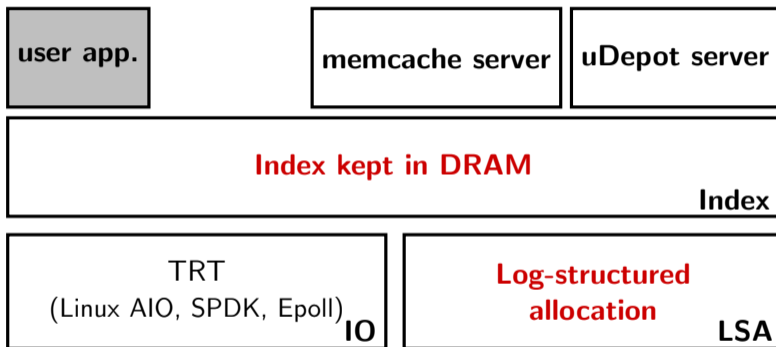
- ▶ avoid cross-core communication
- ▶ user-space tasks (aka green threads, aka co-routines)
- ▶ poller task for multiple IO backends



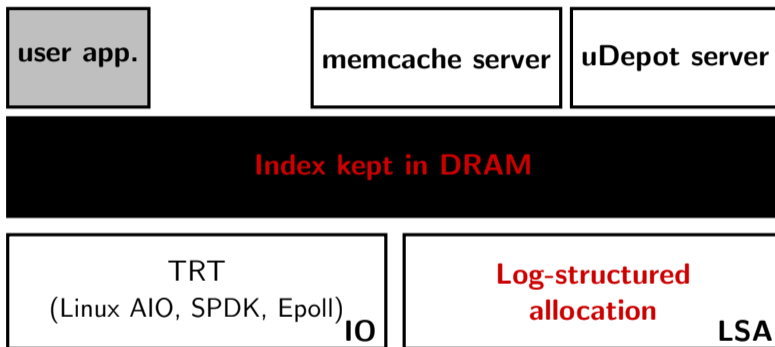
Tasks ( $T$ ) issue IO requests  $\rightarrow$   
Pollers ( $P$ ) poll for completions  $\rightarrow$   
Pollers notify tasks  $- - \rightarrow$



# uDepot architecture

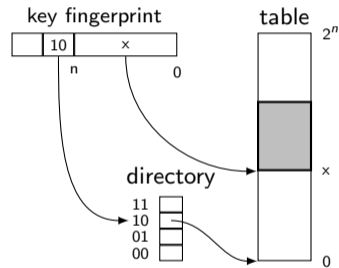


# uDepot architecture



# uDepot index

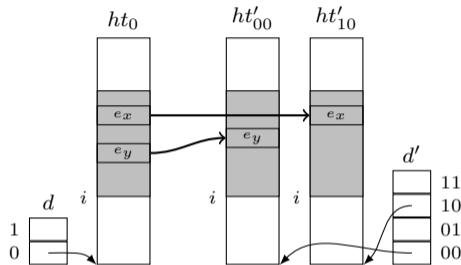
- ▶ Two-level hopscotch hash table
  - ▶ directory + tables
- ▶ 8 byte hash entry (cf. 6-byte for FAWN, FlashStore)
  - ▶ maintain KV size in the entry
  - ▶ larger storage addresses
- ▶ high-performance, scalable
- ▶ efficient resizing



# Growing the uDepot index

## Operations:

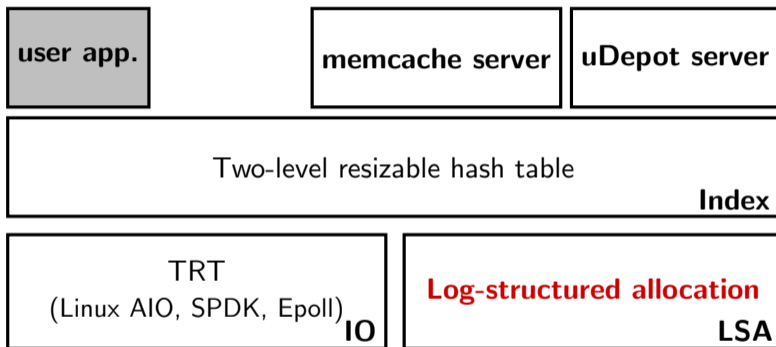
- ▶ double the size of the directory
- ▶ migrate entries to new tables



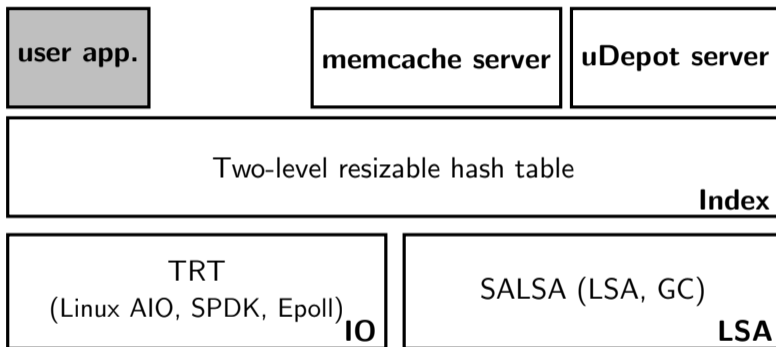
## Minimal disruption

- ▶ **unobstructed reads**
- ▶ **no IO required**: information in the hash entry to reconstruct the fingerprint
- ▶ **incremental**: each operation migrates a bounded number of entries

# uDepot architecture



# uDepot architecture



# Evaluation

1. Performance of uDepot index (with and without resize)
2. Embedded uDepot performance vs device performance
3. uDepot server vs other NVMe stores (Aerospike, ScyllaDB)
4. uDepot memcache vs DRAM-based memcache implementations

1. Performance of uDepot index (with and without resize)
2. Embedded uDepot performance vs device performance
3. uDepot server vs other NVMe stores (Aerospike, ScyllaDB)
4. uDepot memcache vs DRAM-based memcache implementations

## Experiment

- ▶ vs libcuckoo (better performance, see paper)
- ▶ Here: How is tail latency affected by the grow operation?
- ▶ ubench: perform 50M (no grow) and 1B (4 grows) inserts and lookups



## uDepot index latency

percentile	lookup/50M	lookup/1B	insert/50M	insert/1B
50%	0.2 $\mu$ s	0.3 $\mu$ s	0.2 $\mu$ s	0.4 $\mu$ s
99%	1.1 $\mu$ s	1.2 $\mu$ s	0.6 $\mu$ s	1.0 $\mu$ s
99.9%	1.9 $\mu$ s	2.0 $\mu$ s	1.6 $\mu$ s	9.2 $\mu$ s
99.99%	11.0 $\mu$ s	8.9 $\mu$ s	7.5 $\mu$ s	1168.0 $\mu$ s

## uDepot index latency

percentile	lookup/50M	lookup/1B	insert/50M	insert/1B
50%	0.2 $\mu$ s	0.3 $\mu$ s	0.2 $\mu$ s	0.4 $\mu$ s
99%	1.1 $\mu$ s	1.2 $\mu$ s	0.6 $\mu$ s	1.0 $\mu$ s
99.9%	1.9 $\mu$ s	2.0 $\mu$ s	1.6 $\mu$ s	9.2 $\mu$ s
<b>99.99%</b>	11.0 $\mu$ s	8.9 $\mu$ s	7.5 $\mu$ s	<b>1168.0 <math>\mu</math>s</b>

# Evaluation

1. Performance of uDepot index (with and without resize)
2. Embedded uDepot performance vs device performance
3. uDepot server vs other NVMe stores (Aerospike, ScyllaDB)
4. uDepot memcache vs DRAM-based memcache implementations

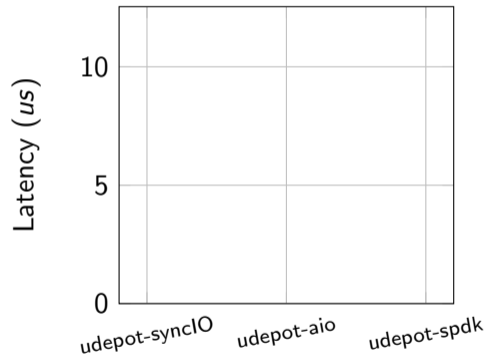
1. Performance of uDepot index (with and without resize)
2. Embedded uDepot performance vs device performance
3. uDepot server vs other NVMe stores (Aerospike, ScyllaDB)
4. uDepot memcache vs DRAM-based memcache implementations

## Experiment

- ▶ uDepot ubench: perform 10M uDepot PUTs and GETs
  - ▶ multiple backends (how different backends behave)
- ▶ vs. dev ubench: fio and SPDK perf
- ▶ same workload: 4KiB

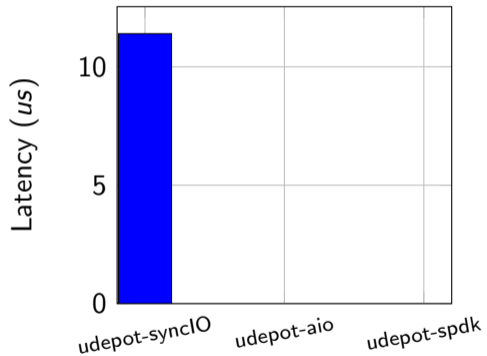
# Embedded uDepot: Efficiency (1 core / 1 Optane)

GET: median latency for qd=1



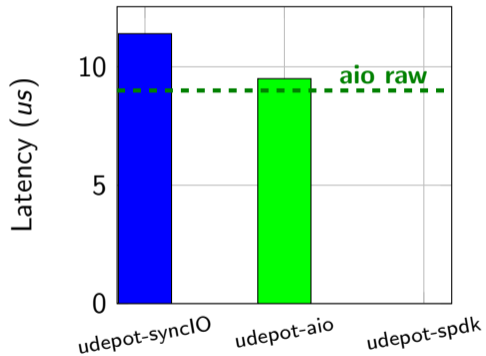
# Embedded uDepot: Efficiency (1 core / 1 Optane)

GET: median latency for qd=1



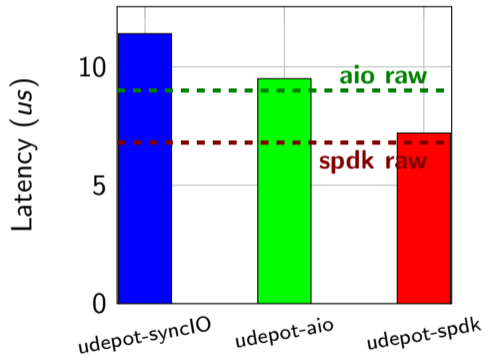
# Embedded uDepot: Efficiency (1 core / 1 Optane)

GET: median latency for qd=1



# Embedded uDepot: Efficiency (1 core / 1 Optane)

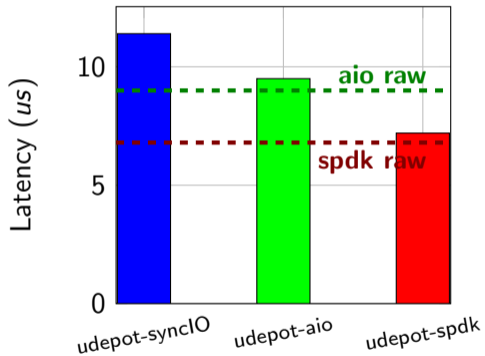
GET: median latency for qd=1



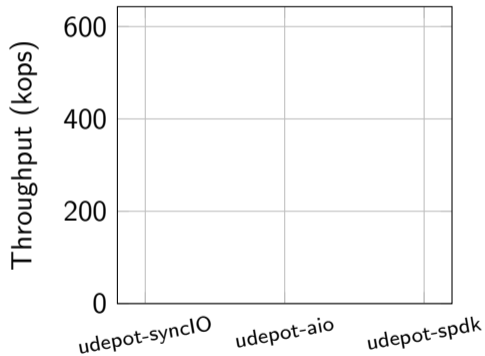


# Embedded uDepot: Efficiency (1 core / 1 Optane)

GET: median latency for qd=1

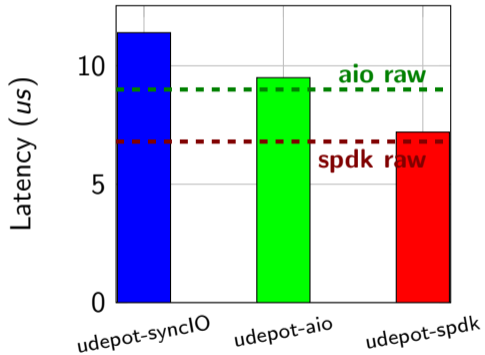


GET: throughput for qd=128

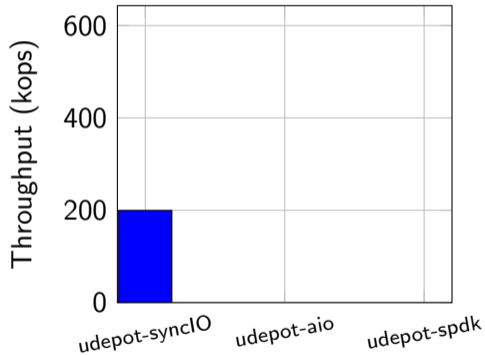


# Embedded uDepot: Efficiency (1 core / 1 Optane)

GET: median latency for qd=1

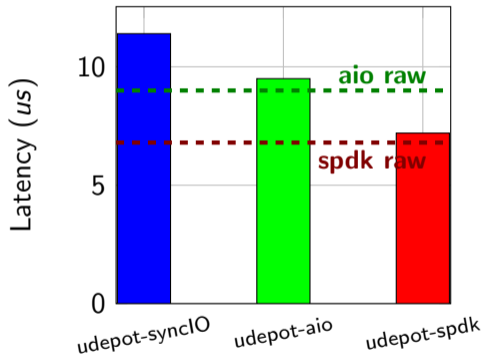


GET: throughput for qd=128

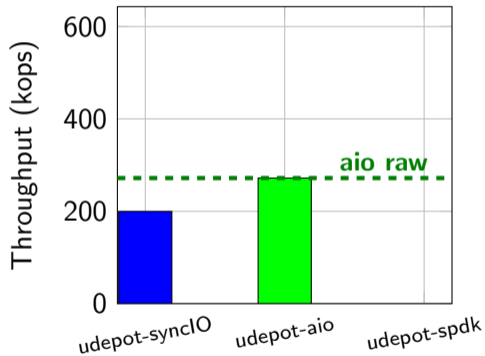


# Embedded uDepot: Efficiency (1 core / 1 Optane)

GET: median latency for qd=1

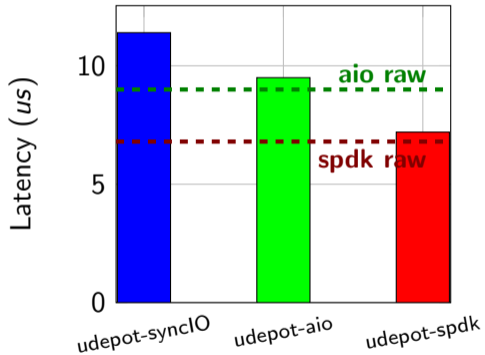


GET: throughput for qd=128

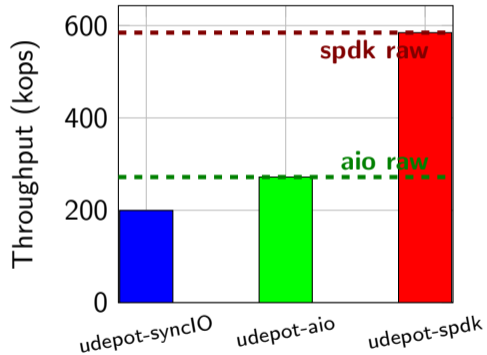


# Embedded uDepot: Efficiency (1 core / 1 Optane)

GET: median latency for qd=1

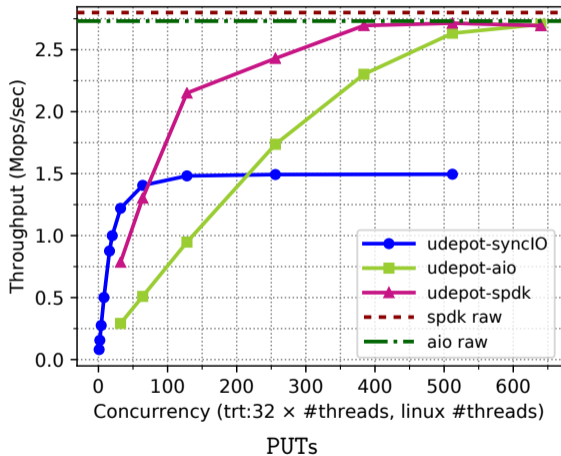
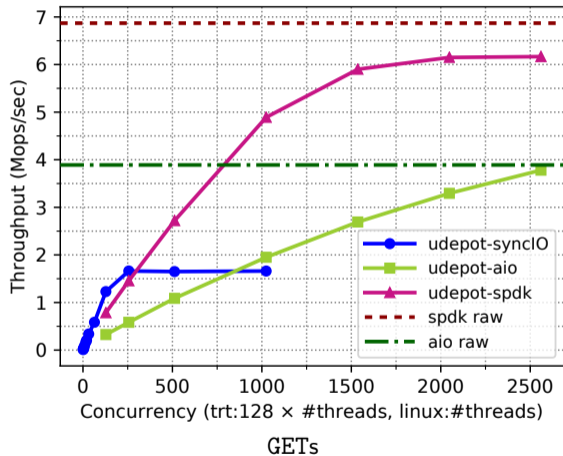


GET: throughput for qd=128

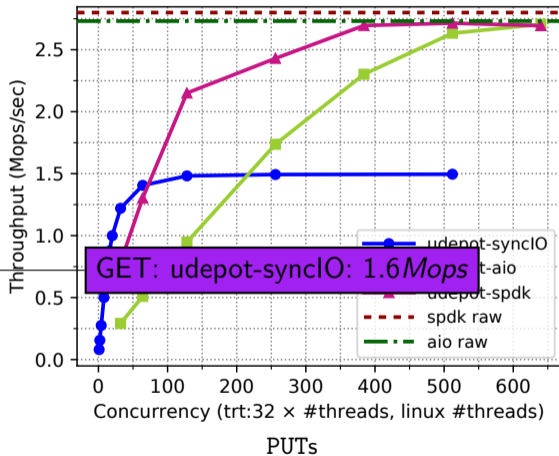
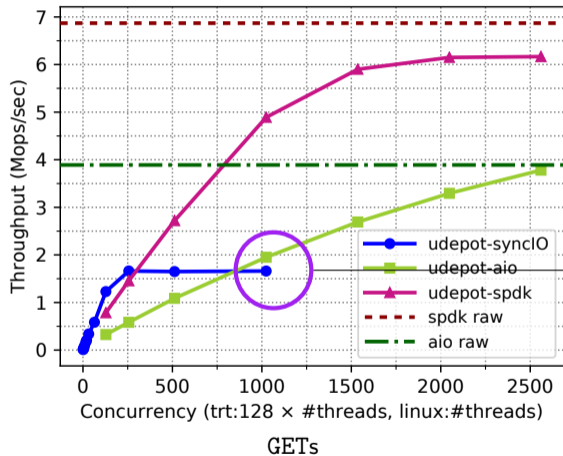


# Embedded uDepot: Scalability (20 cores, 24 Flash NVMeS)

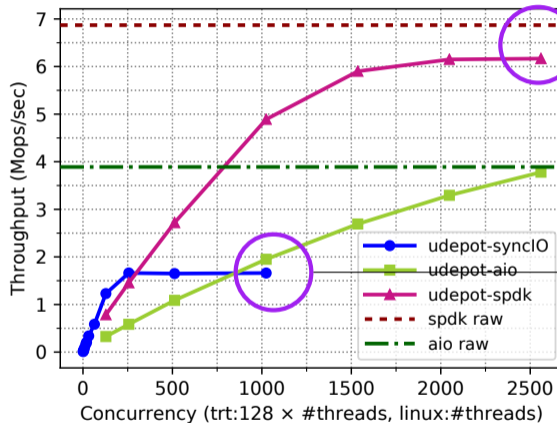
# Embedded uDepot: Scalability (20 cores, 24 Flash NVMes)



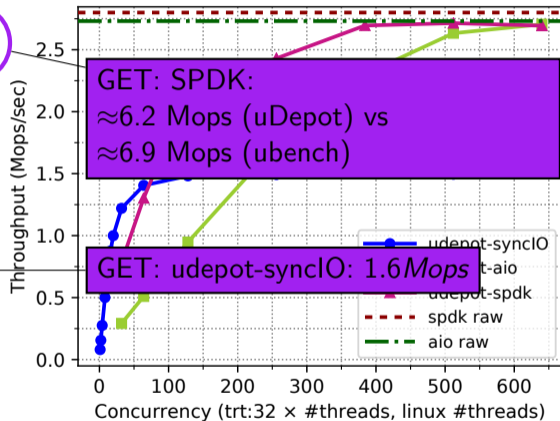
# Embedded uDepot: Scalability (20 cores, 24 Flash NVMe)



# Embedded uDepot: Scalability (20 cores, 24 Flash NVMe)



GETs



GET: SPDK:  
≈6.2 Mops (uDepot) vs  
≈6.9 Mops (ubench)

GET: udepot-syncIO: 1.6 Mops

PUTs



# Evaluation

1. Performance of uDepot index (with and without resize)
2. Embedded uDepot performance vs device performance
3. uDepot server vs other NVMe stores (Aerospike, ScyllaDB)
4. uDepot memcache vs DRAM-based memcache implementations

# Evaluation

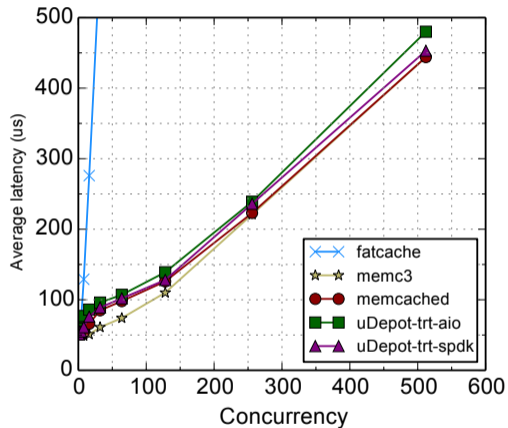
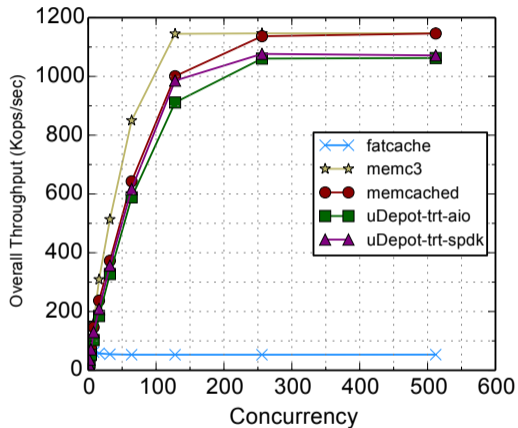
1. Performance of uDepot index (with and without resize)
2. Embedded uDepot performance vs device performance
3. uDepot server vs other NVMe stores (Aerospire, ScyllaDB) [see paper]
4. uDepot memcache vs DRAM-based memcache implementations

1. Performance of uDepot index (with and without resize)
2. Embedded uDepot performance vs device performance
3. uDepot server vs other NVMe stores (Aerospire, ScyllaDB) [see paper]
4. uDepot memcache vs DRAM-based memcache implementations

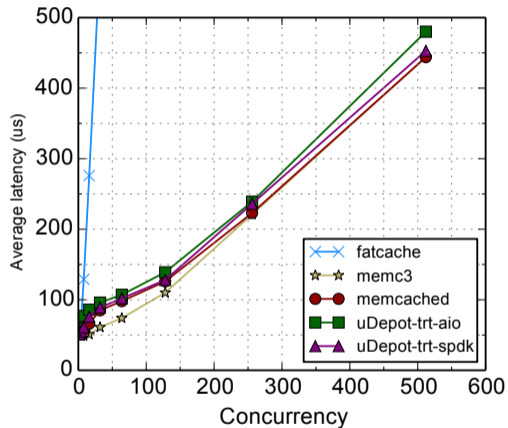
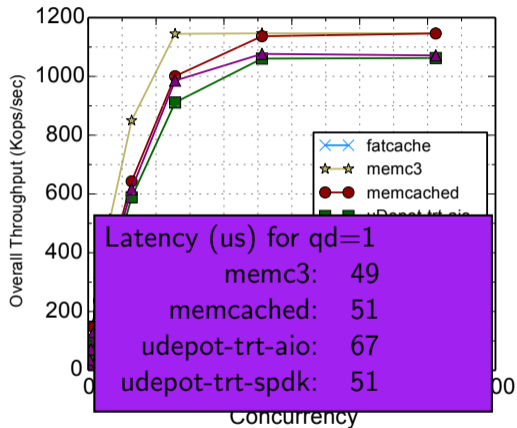
## Experiment

- ▶ memslap benchmark
- ▶ default workload: 1KiB objects, 10%/90% PUT/GET
- ▶ 2 Optane SSDs, 20 cores, 10GbE
- ▶ vs Memcached (expected performance), Memc3 (optimized memcached), Fatcache (traditional SSD impl.)

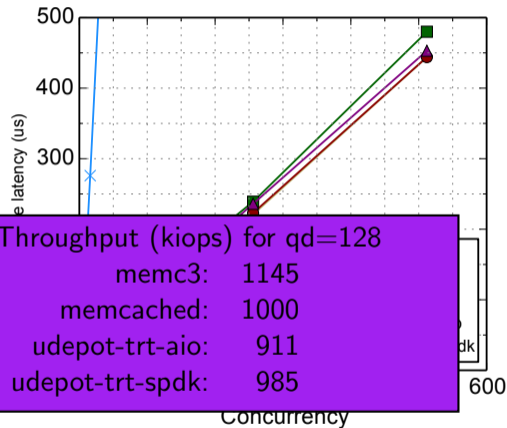
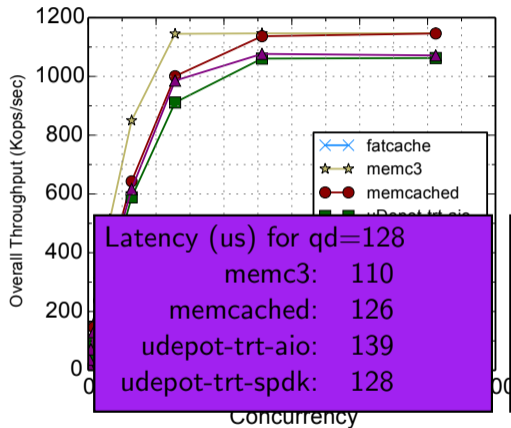
# uDepot memcache



# uDepot memcache



# uDepot memcache



- ▶ Fast NVMe devices offer an attractive cost-performance tradeoff between DRAM and Flash SSDs
  - ▶ They shift the bottleneck from the storage to the network

- ▶ Fast NVMe devices offer an attractive cost-performance tradeoff between DRAM and Flash SSDs
  - ▶ They shift the bottleneck from the storage to the network
- ▶ uDepot: a KV store that delivers their performance
  - ▶ low latency, high throughput
  - ▶ uDepot memcache has comparable performance with DRAM memcache



- ▶ Fast NVMe devices offer an attractive cost-performance tradeoff between DRAM and Flash SSDs
  - ▶ They shift the bottleneck from the storage to the network
- ▶ uDepot: a KV store that delivers their performance
  - ▶ low latency, high throughput
  - ▶ uDepot memcache has comparable performance with DRAM memcache
- ▶ Experimental Cloud service based on uDepot memcache implementation
  - ▶ try it out (for free):  
<https://cloud.ibm.com/catalog/services/data-store-for-memcache>

Thank you! Questions?