# Effectively Prefetching Remote Memory with Leap
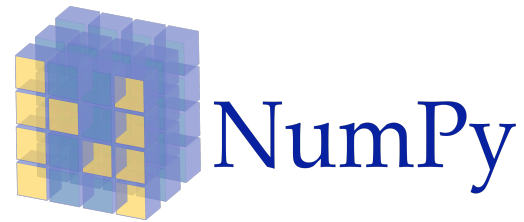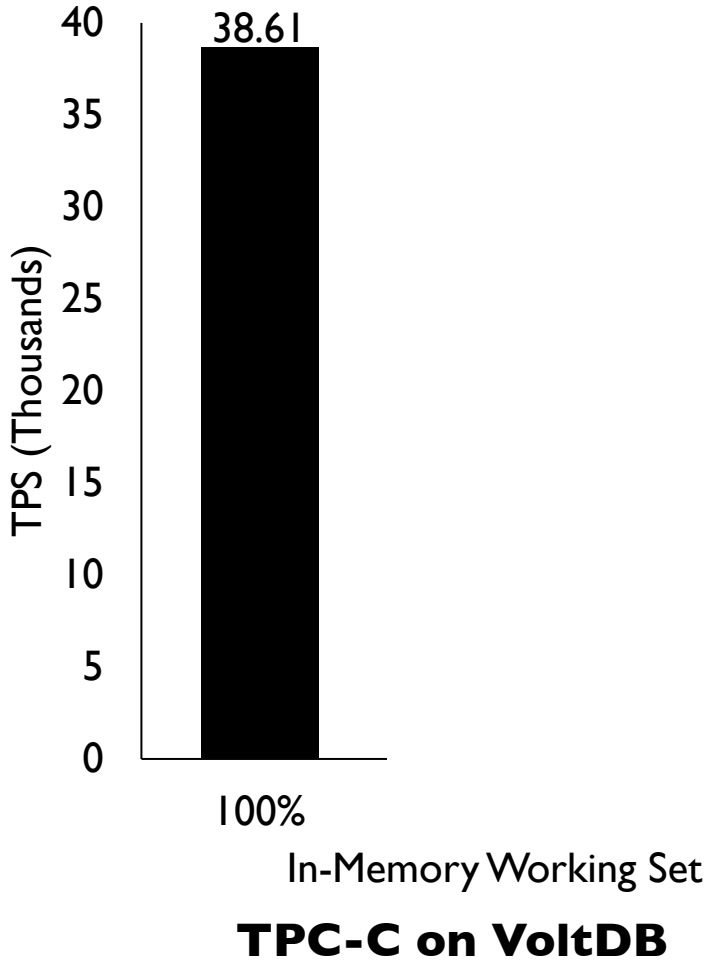
*Hasan Al Maruf* and *Mosharaf Chowdhury*
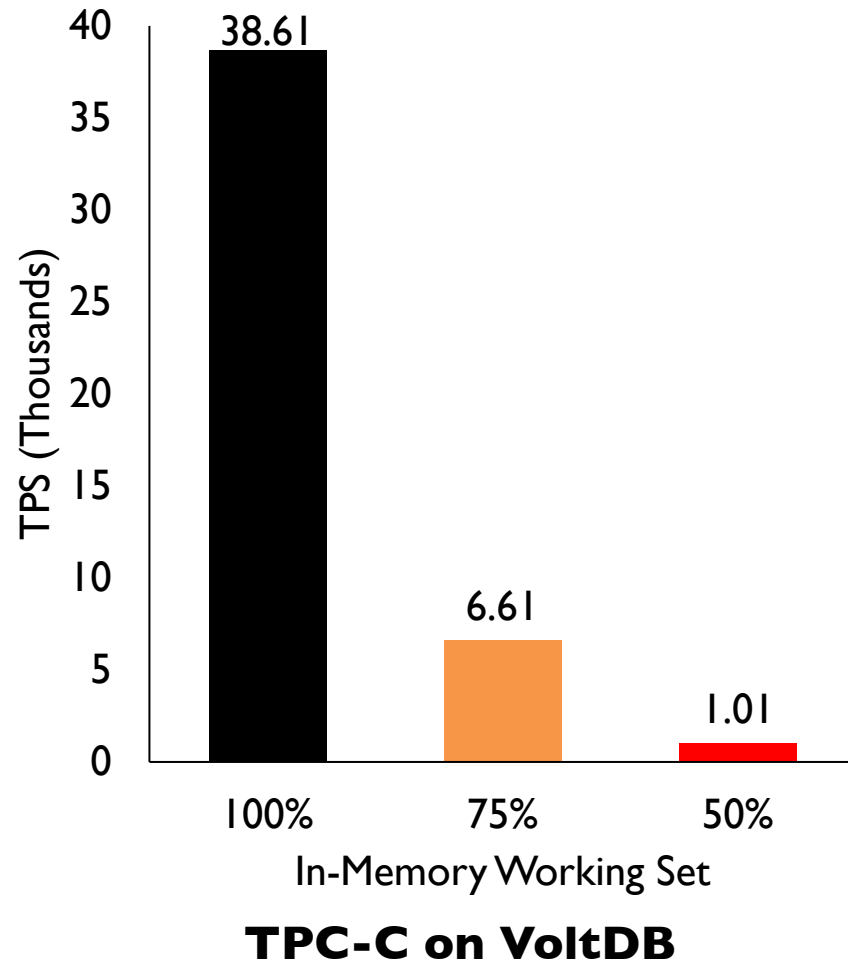
University of Michigan

1

# Memory-Intensive Applications

# Perform Great!



**TPC-C on VoltDB**

# Perform Great Until Memory Runs Out



TPC-C on VoltDB

PageRank on PowerGraph

# 50% Less Memory Causes Slowdown of …



TPC-C on VoltDB

PageRank on PowerGraph

**38X**

**4X**

# Between a Rock and a Hard Place

## Underallocation

Leads to severe performance loss

### vs.

## Overallocation

Leads to underutilization

30-40% in Google, Alibaba, and Facebook

# Memory Disaggregation

# Remote Memory Access

User-space Applications

Memory Disaggregation Frameworks

**Infiniswap**
(NSDI'17)

Remote
memory paging

**Remote Regions**
(ATC'18)

Remote file
abstraction

**LegoOS**
(OSDI'18)

Disaggregated
OS

Remote Memory

4KB page access latency
local vs. remote

**100 ns** vs. **4 μs**

# Remote Memory Access

| User-space Applications |
|---|

↕

**Memory Disaggregation Frameworks**

| **Infiniswap** (NSDI'17) Remote memory paging | **Remote Regions** (ATC'18) Remote file abstraction | **LegoOS** (OSDI'18) Disaggregated OS |
|---|---|---|

↕

**Remote Memory**

[1] P. X. Gao et al. "Network requirements for resource disaggregation" OSDI'16.

4KB page access latency
local vs. remote

**100 ns** vs. **4 μs**

Latency requirement for
preferable performance[1]

**3 μs**

**Existing frameworks can't achieve!**

# Remote Memory Access

User-space Applications

Memory Disaggregation Frameworks

**Infiniswap** (NSDI'17)

Remote memory paging

**Remote Regions** (ATC'18)

Remote file abstraction

**LegoOS** (OSDI'18)

Disaggregated OS

Remote Memory

data path overhead

variation in network latency

4KB page access latency local vs. remote

**100 ns** vs. **4 μs**

Latency requirement for preferable performance[1]

**3 μs**

**Existing frameworks can't achieve!**
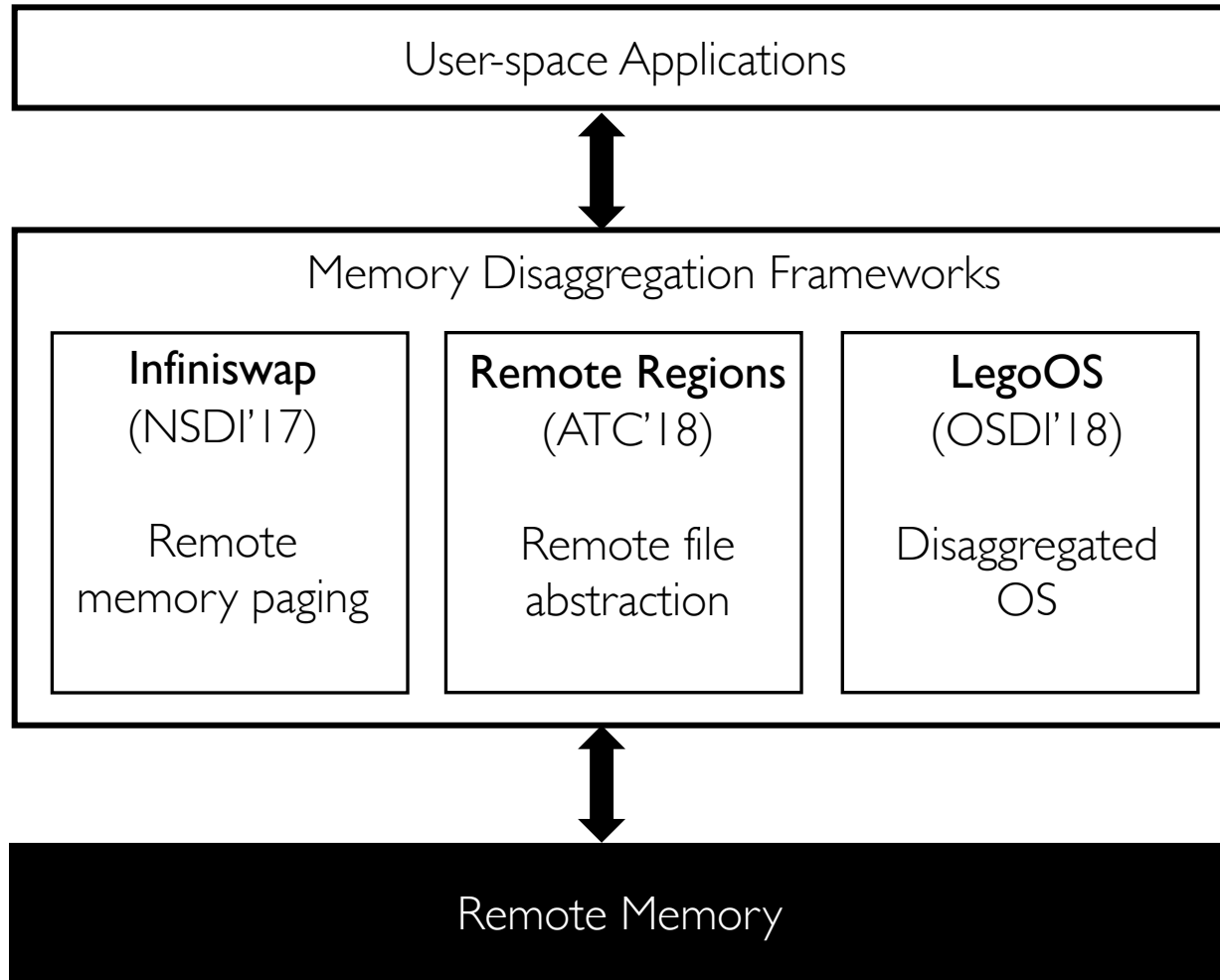
[1] P. X. Gao et al. "Network requirements for resource disaggregation" OSDI'16.

# Life of a Page

# Where Does the Time Go?

# Where Does the Time Go?

# Design Goal

1. Increase cache hit
   - faster path serves more page faults

2. Reduce the latency of the slow path
   - remove unnecessary block-layer operations for RDMA

# **Leap**

*Online remote memory prefetcher*

Identifies memory access patterns to prefetch pages in a
- fast,
- cache-efficient, and
- resilient manner

**without modifying any**
- applications, or
- hardware

# Life of a Page



Cache Hit

0.27 us

Cache Miss

2.1 us

10.04 us

21.88 us

RDMA: 4.3 us

Process 1   Process 2   ...   Process N      User Space

Page Fault

Memory Management Unit (MMU)

MMU Page Cache      Kernel Space

Device Mapping Layer

Generic Block Layer      bio

I/O Scheduler

Request queue processing:
Insertion, Merging,
Sorting, Staging and Dispatch

Request Queue

Dispatch Queue

Block Device Driver

Remote Memory

# Life of a Page w/ Leap

Cache Hit

0.27 us

2.1 us

Cache Miss

RDMA: 4.3 us

User Space

Process 1    Process 2    …    Process N

Page Fault

Memory Management Unit (MMU)

MMU Page Cache

Kernel Space

Remote Memory

# Prefetching in Linux

Reads ahead pages sequentially

Based only on the last page access

*too aggressive on seq*: cache pollution

*too conservative off seq*: brings nothing

Does not distinguish between processes

Cannot detect thread-level access irregularities

# Prefetching Techniques

| Approach | Low Computational Complexity | Low Memory Overhead | Unmodified Application | HW/SW Independence | Temporal Locality | Spatial Locality | Low Cache Pollution |
|----------|------------------------------|---------------------|------------------------|--------------------|--------------------|------------------|---------------------|

# Prefetching Techniques

| Approach | Low Computational Complexity | Low Memory Overhead | Unmodified Application | HW/SW Independence | Temporal Locality | Spatial Locality | Low Cache Pollution |
|----------|------------------------------|---------------------|------------------------|--------------------|--------------------|-------------------|----------------------|
| Next N-Line | Yes | Yes | Yes | Yes | No | Yes | No |
| Stride | Yes | Yes | Yes | Yes | No | Yes | No |

# Prefetching Techniques

| Approach | Low Computational Complexity | Low Memory Overhead | Unmodified Application | HW/SW Independence | Temporal Locality | Spatial Locality | Low Cache Pollution |
|---|---|---|---|---|---|---|---|
| Next N-Line | Yes | Yes | Yes | Yes | No | Yes | No |
| Stride | Yes | Yes | Yes | Yes | No | Yes | No |
| Instruction Prefetch | No | No | No | No | Yes | Yes | No |

# Prefetching Techniques

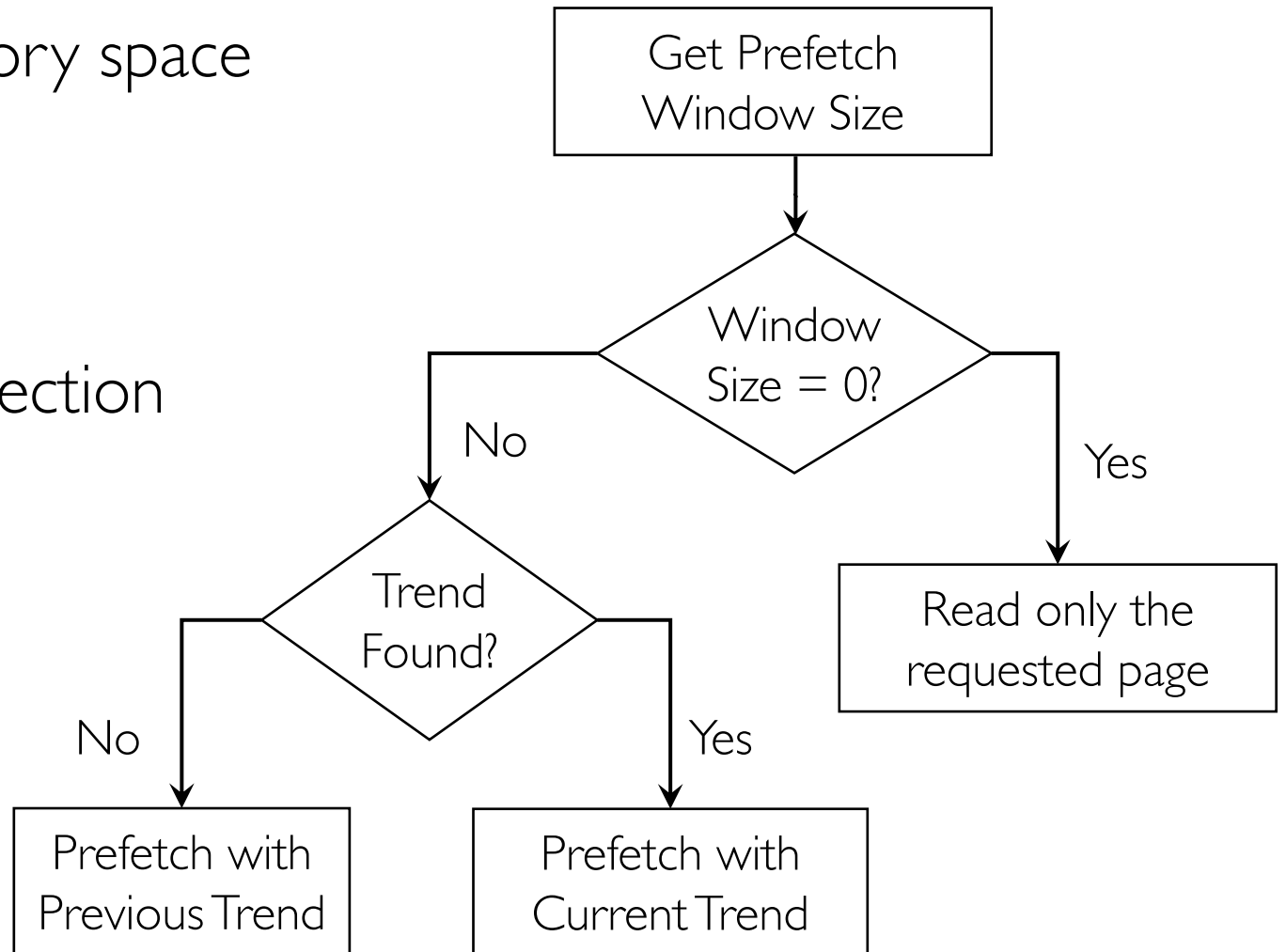| Approach | Low Computational Complexity | Low Memory Overhead | Unmodified Application | HW/SW Independence | Temporal Locality | Spatial Locality | Low Cache Pollution |
|---|---|---|---|---|---|---|---|
| Next N-Line | Yes | Yes | Yes | Yes | No | Yes | No |
| Stride | Yes | Yes | Yes | Yes | No | Yes | No |
| Instruction Prefetch | No | No | No | No | Yes | Yes | No |
| Linux Read-Ahead | Yes | Yes | Yes | Yes | Yes | Yes | No |

# Prefetching Techniques

| Approach | Low Computational Complexity | Low Memory Overhead | Unmodified Application | HW/SW Independence | Temporal Locality | Spatial Locality | Low Cache Pollution |
|---|---|---|---|---|---|---|---|
| Next N-Line | Yes | Yes | Yes | Yes | No | Yes | No |
| Stride | Yes | Yes | Yes | Yes | No | Yes | No |
| Instruction Prefetch | No | No | No | No | Yes | Yes | No |
| Linux Read-Ahead | Yes | Yes | Yes | Yes | Yes | Yes | No |
| Leap | Yes | Yes | Yes | Yes | Yes | Yes | Yes |

# Leap Prefetcher

Linear-time and constant memory space

Two main components:
- Trend detection
- Prefetch window size detection



Get Prefetch Window Size

Window Size = 0?

No

Yes

Trend Found?

No

Yes

Read only the requested page

Prefetch with Previous Trend
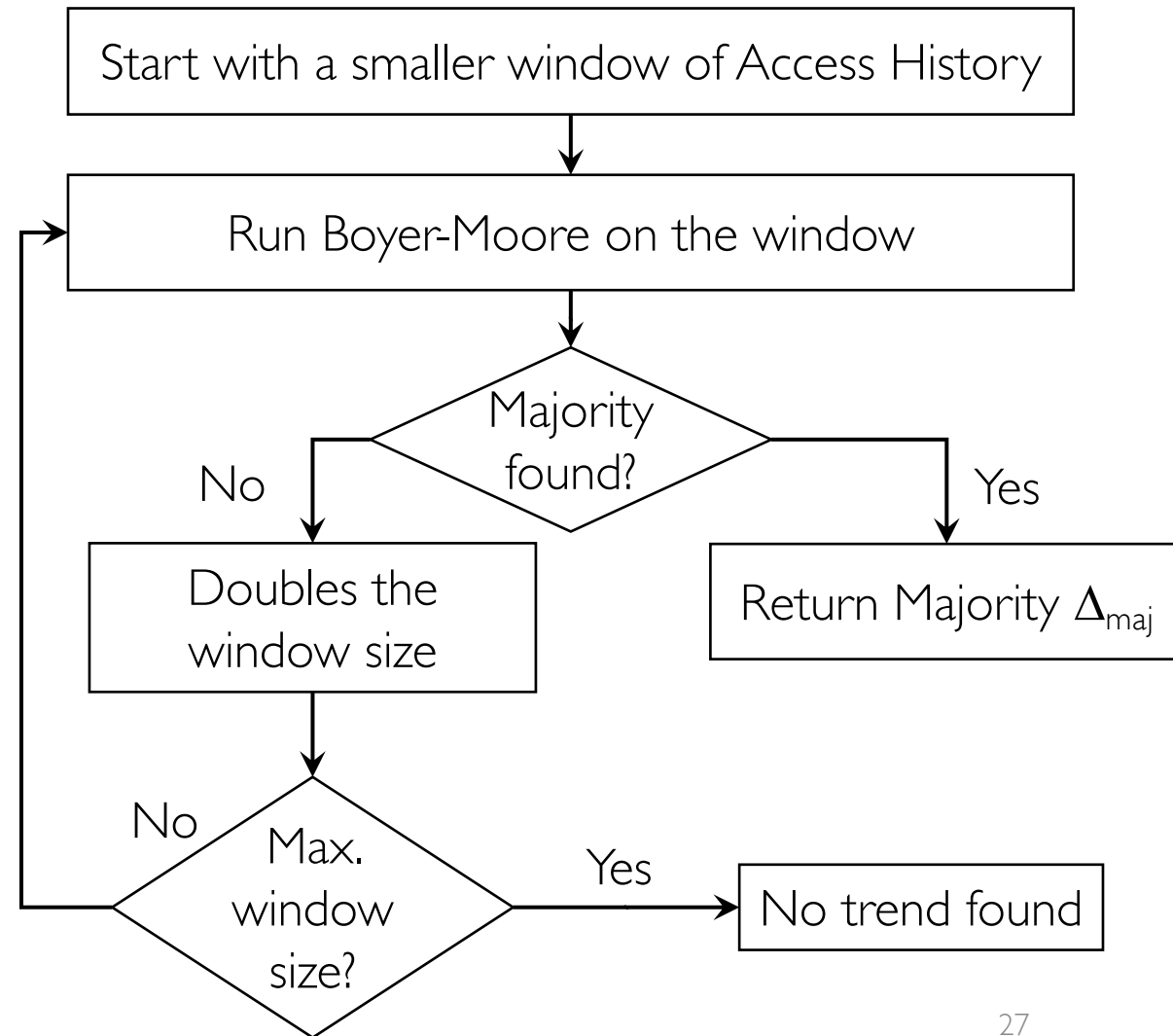
Prefetch with Current Trend

# Trend Detection
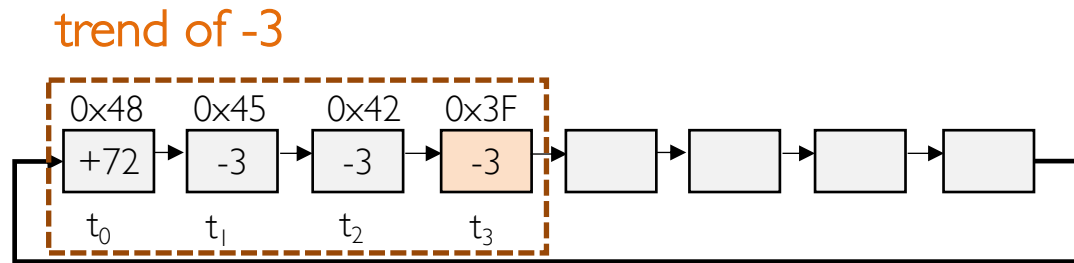
Flexible to short term irregularity
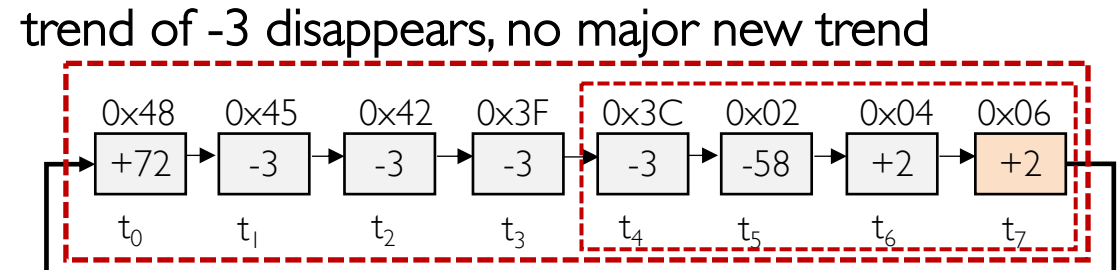
Identifies the majority element in access history

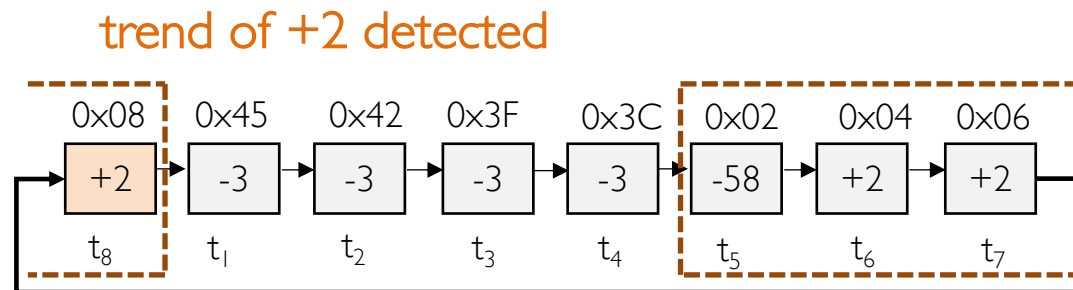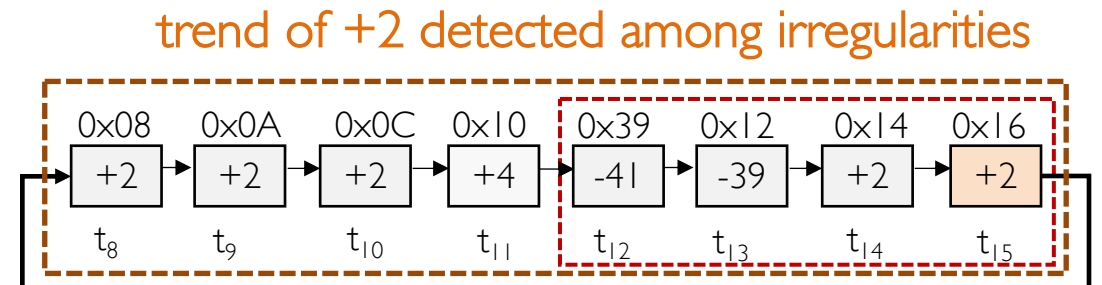Regular trends can be found within recent accesses

Start with a smaller window of Access History

Run Boyer-Moore on the window

Majority found?

No

Yes

Doubles the window size

Return Majority $\Delta_{maj}$

Max. window size?

No

Yes

No trend found

# Trend Detection Example

trend of -3

| 0x48 | 0x45 | 0x42 | 0x3F | | | | |
|------|------|------|------|--|--|--|--|
| +72 | -3 | -3 | -3 | | | | |
| $t_0$ | $t_1$ | $t_2$ | $t_3$ | | | | |

(a) at time t3

trend of -3 disappears, no major new trend

| 0x48 | 0x45 | 0x42 | 0x3F | 0x3C | 0x02 | 0x04 | 0x06 |
|------|------|------|------|------|------|------|------|
| +72 | -3 | -3 | -3 | -3 | -58 | +2 | +2 |
| $t_0$ | $t_1$ | $t_2$ | $t_3$ | $t_4$ | $t_5$ | $t_6$ | $t_7$ |

(b) at time t7

trend of +2 detected

| 0x08 | 0x45 | 0x42 | 0x3F | 0x3C | 0x02 | 0x04 | 0x06 |
|------|------|------|------|------|------|------|------|
| +2 | -3 | -3 | -3 | -3 | -58 | +2 | +2 |
| $t_8$ | $t_1$ | $t_2$ | $t_3$ | $t_4$ | $t_5$ | $t_6$ | $t_7$ |

(c) at time t8

trend of +2 detected among irregularities

| 0x08 | 0x0A | 0x0C | 0x10 | 0x39 | 0x12 | 0x14 | 0x16 |
|------|------|------|------|------|------|------|------|
| +2 | +2 | +2 | +4 | -41 | -39 | +2 | +2 |
| $t_8$ | $t_9$ | $t_{10}$ | $t_{11}$ | $t_{12}$ | $t_{13}$ | $t_{14}$ | $t_{15}$ |

(d) at time t15

# Prefetch Window Size Detection

Cache hit indicates prefetch utilization

High cache hit: increase prefetch window aggressively

No cache hit
- *trend availability: increase prefetch window gradually*
- *no trend: decrease prefetch window gradually*
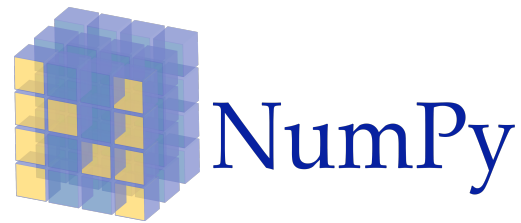
Gradual slow down helps during sudden changes

# Evaluation

*Deploy and evaluate over 56 Gbps InfiniBand network*
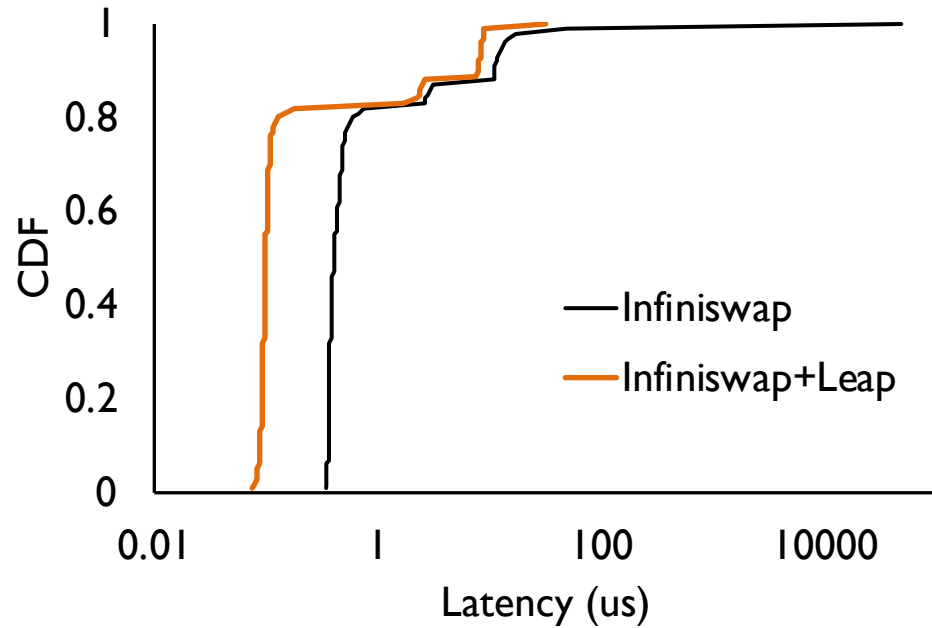
Memory Disaggregation Frameworks

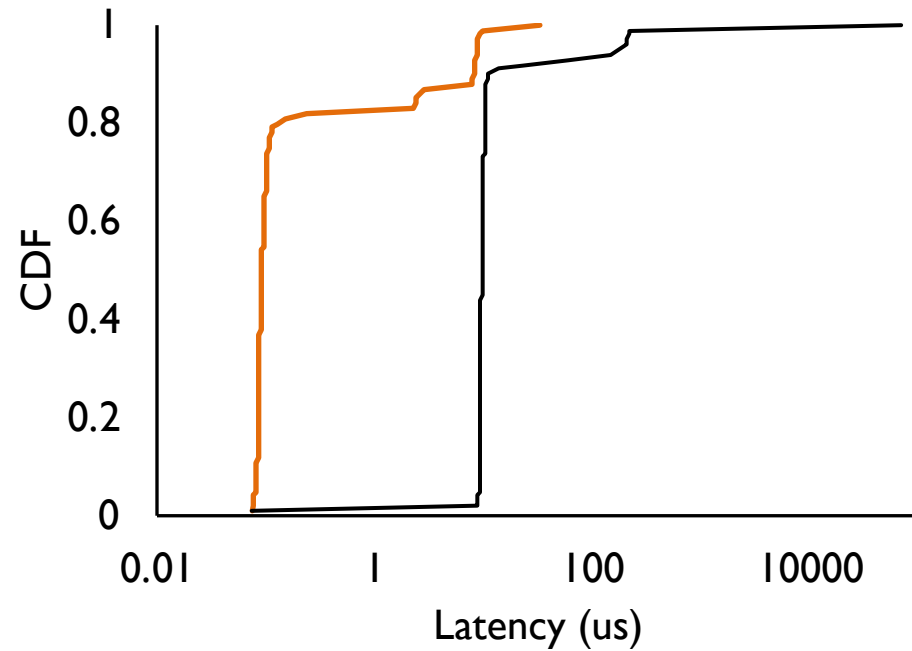*Disaggregated VMM:* **Infiniswap**

*Disaggregated VFS:* **Remote Regions**

# Lowers Remote Page Access Latency by…
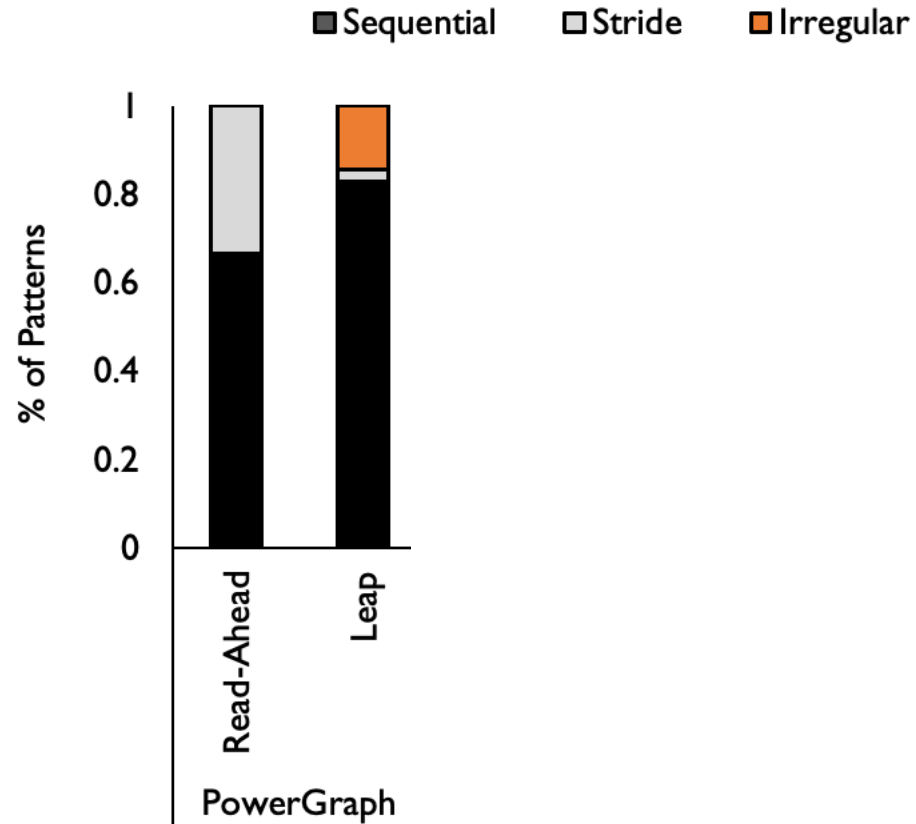
**Sequential Access**



**4X**

**Stride Access**



**104X**

# Efficient Pattern Detection

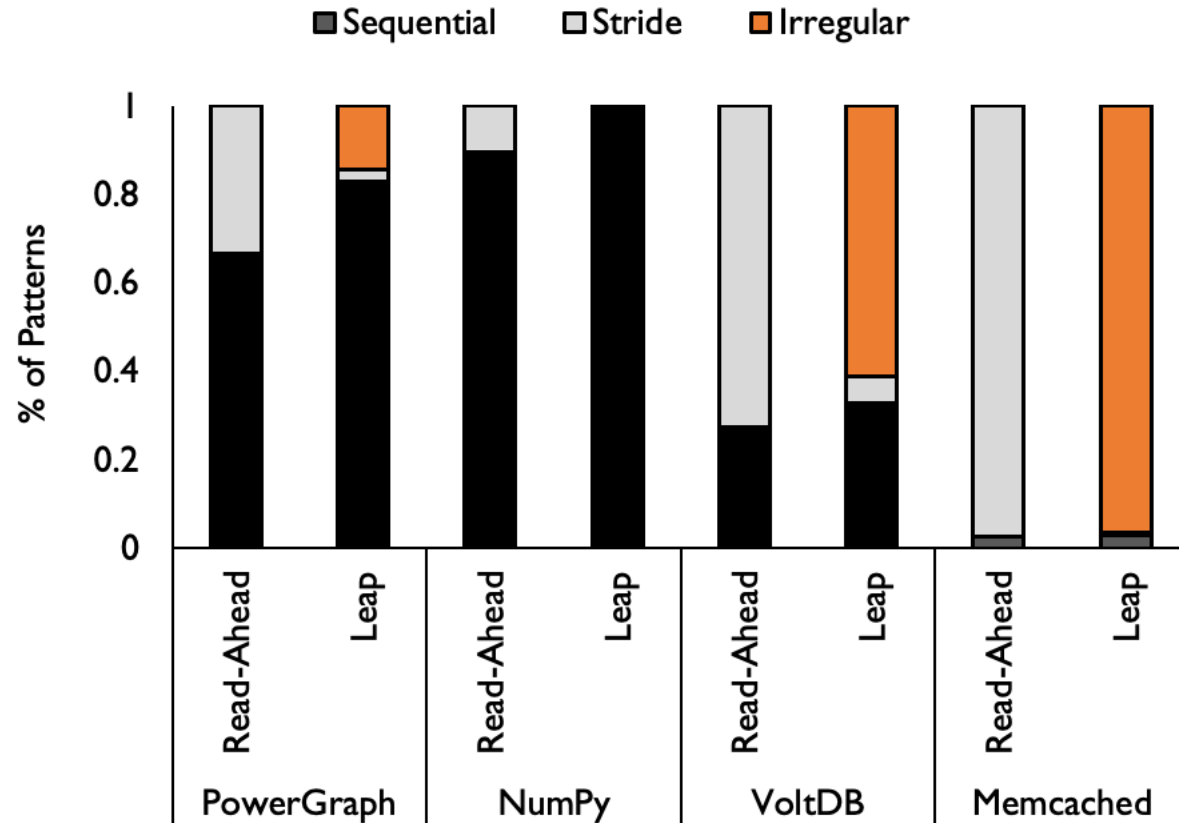Sequential ▪  Stride ▫  Irregular ▪

% of Patterns

Read-Ahead   Leap

PowerGraph

Detects **29.70%** more sequential accesses

Detects most of the irregularity
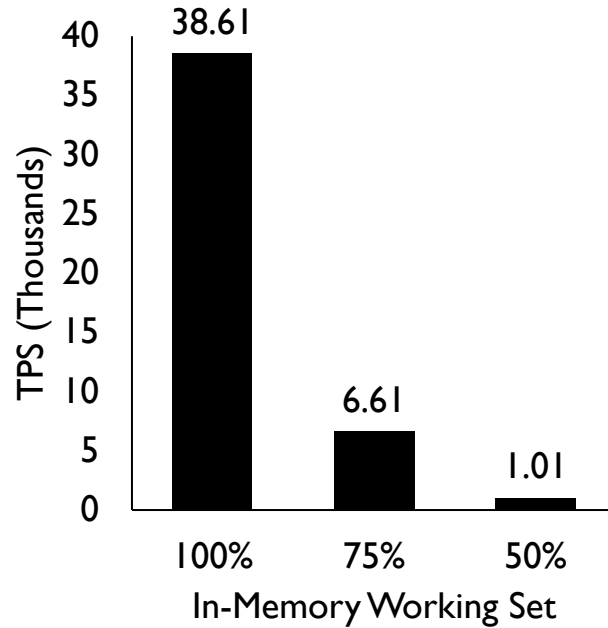
# Efficient Pattern Detection



Detects **29.70%** more sequential accesses

Detects most of the irregularity

During irregularities, doing nothing helps the most

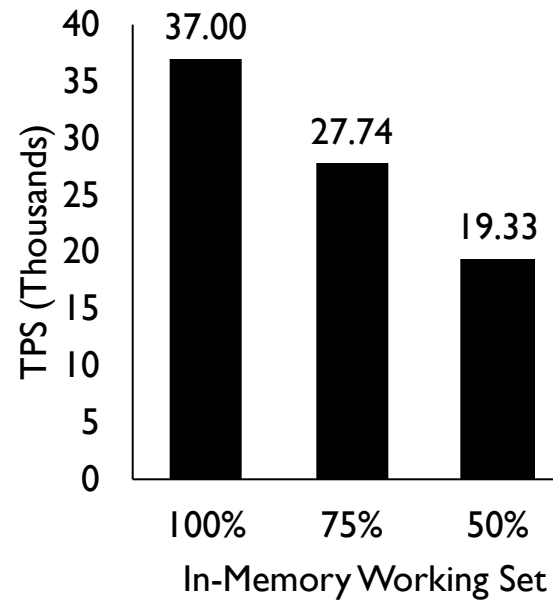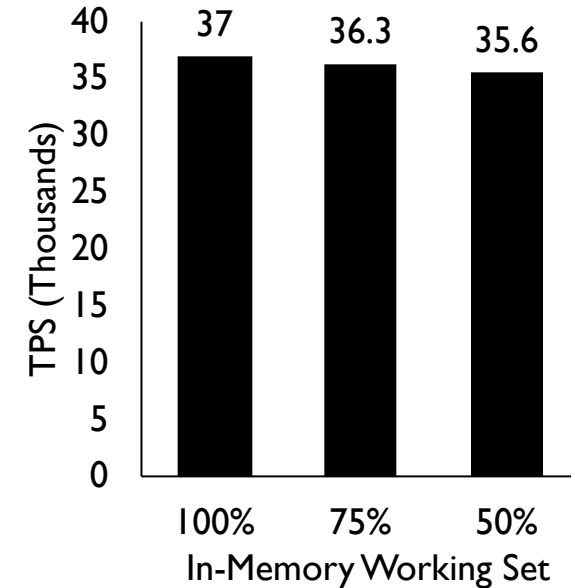# Perform Great Even After Memory Runs Out



**Disk** — TPC-C on VoltDB: **38X**
(100%: 38.61, 75%: 6.61, 50%: 1.01)

**Infiniswap** — TPC-C on VoltDB: **2X**
(100%: 37.00, 75%: 27.74, 50%: 19.33)

**Infiniswap + Leap** — TPC-C on VoltDB: **≈1X**
(100%: 37, 75%: 36.3, 50%: 35.6)

# Perform Great Even After Memory Runs Out



**Disk**

TPS (Thousands) vs In-Memory Working Set
- 100%: 38.61
- 75%: 6.62
- 50%: 1.01
- 25%: Fails

**TPC-C on VoltDB**

**Fails**

**Infiniswap**

TPS (Thousands) vs In-Memory Working Set
- 100%: 37.00
- 75%: 27.74
- 50%: 19.33
- 25%: 1.5

**TPC-C on VoltDB**

**24X**

**Infiniswap + Leap**

TPS (Thousands) vs In-Memory Working Set
- 100%: 37
- 75%: 36.3
- 50%: 35.6
- 25%: 15.6

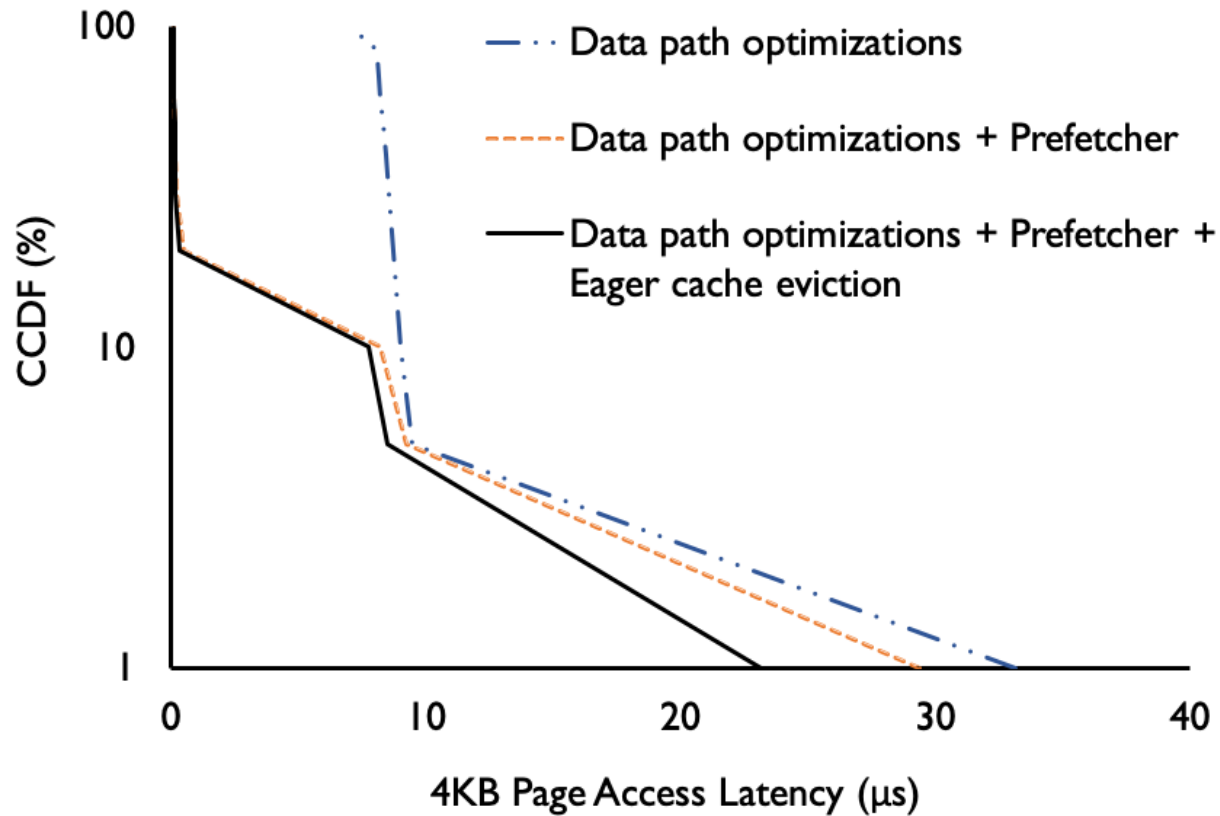**TPC-C on VoltDB**

**2.4X**

# Benefit Breakdown of Leap's Components



**Data path optimizations**: single-μs latency till 95th percentile

**Prefetcher**: sub-μs latency till 85th percentile

**Eager cache eviction**: improves the 99th percentile latency by 22%

# Future Work

1. Thread-specific prefetching for multiple concurrent streams
   - memory is managed at the process level
   - this requires significant changes in virtual memory subsystem

2. Optimized remote I/O interface
   - load balancing,
   - fault-tolerance,
   - data locality, and
   - application-specific isolation in remote memory

# Leap

## Lightweight and efficient data path for remote memory

source code available at https://github.com/SymbioticLab/leap

Online prefetcher with a leaner data path and eager cache eviction policy to improve
- cache hit,
- remote I/O latency, and
- application-level performance

without modifying any
- application, or
- hardware

# Thank You!

source code available at https://github.com/SymbioticLab/leap