SOSP 2024

# CHIME: A Cache-Efficient and High-Performance Hybrid Index on Disaggregated Memory

**Xuchuan Luo,** Jiacheng Shen, Pengfei Zuo, Xin Wang, Micheal R.Lyu, Yangfan Zhou

***School of Computer Science, Fudan University***
*National Key Laboratory of Parallel and Distributed Computing, China*
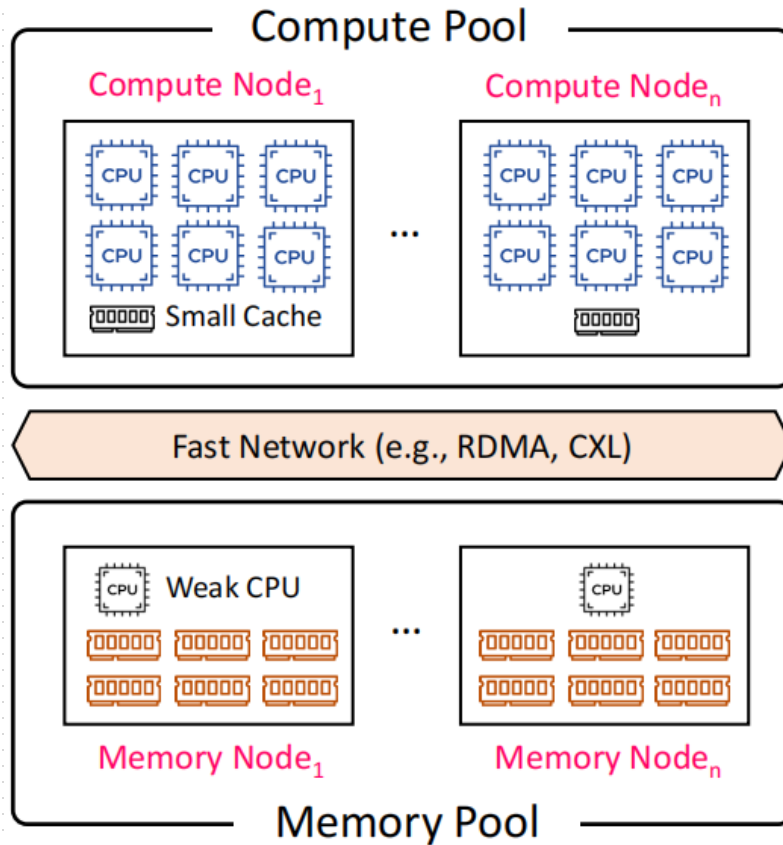*Duke Kunshan University  Huawei Cloud The Chinese University of Hong Kong*
*Shanghai Key Laboratory of Intelligent Information Processing, Shanghai, China*

Presented by Sen Han
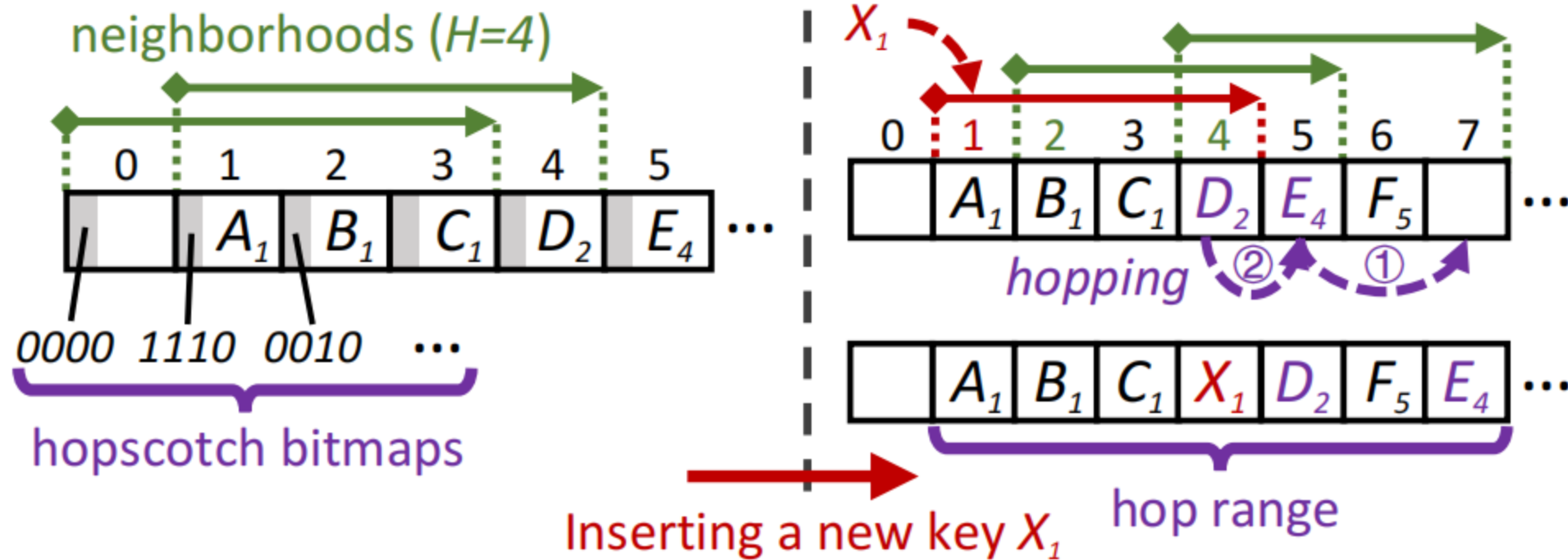
# Background-DM

## Disaggregated Memory(DM)



**Compute Pool**

Compute Node$_1$ ... Compute Node$_n$

CPU CPU CPU CPU CPU CPU
CPU CPU CPU CPU CPU CPU

Small Cache

Fast Network (e.g., RDMA, CXL)

Weak CPU ... CPU

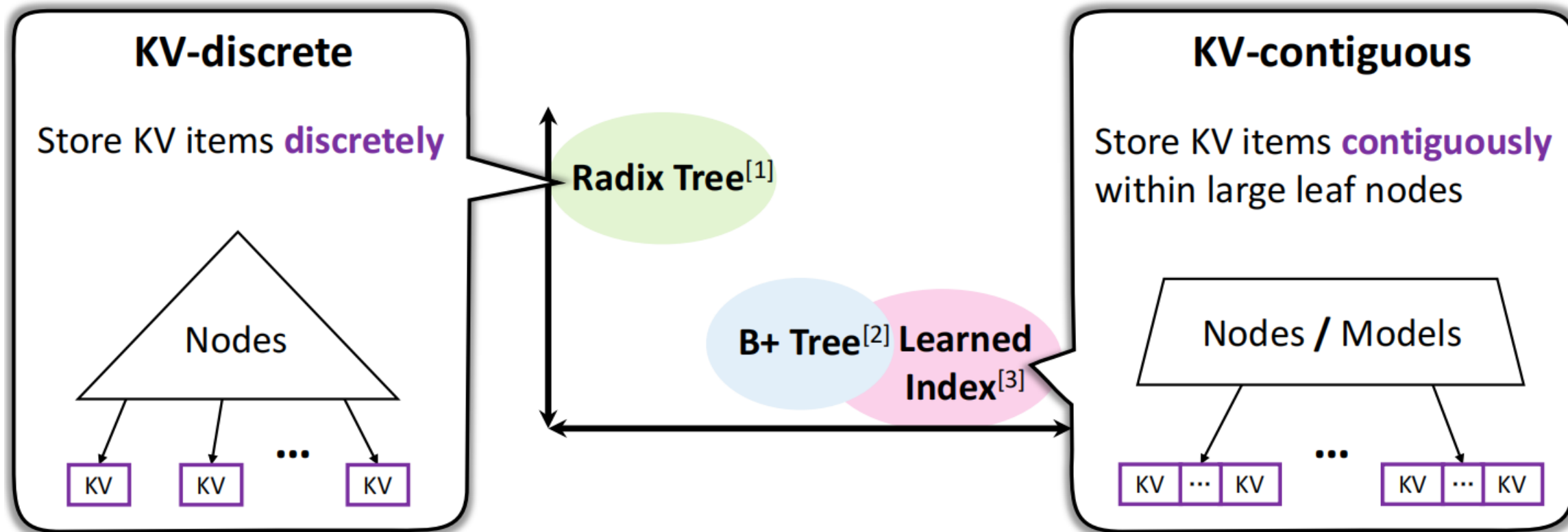Memory Node$_1$ ... Memory Node$_n$

**Memory Pool**

☺ **Benefits:**
✓ Resource utilization
✓ Elasticity

## HopScotch Hashing

## Existing range indexes on DM can be classified into two types:

[1] Xuchuan Luo et al. SMART: A high-performance adaptive radix tree for disaggregated memory. OSDI 2023.
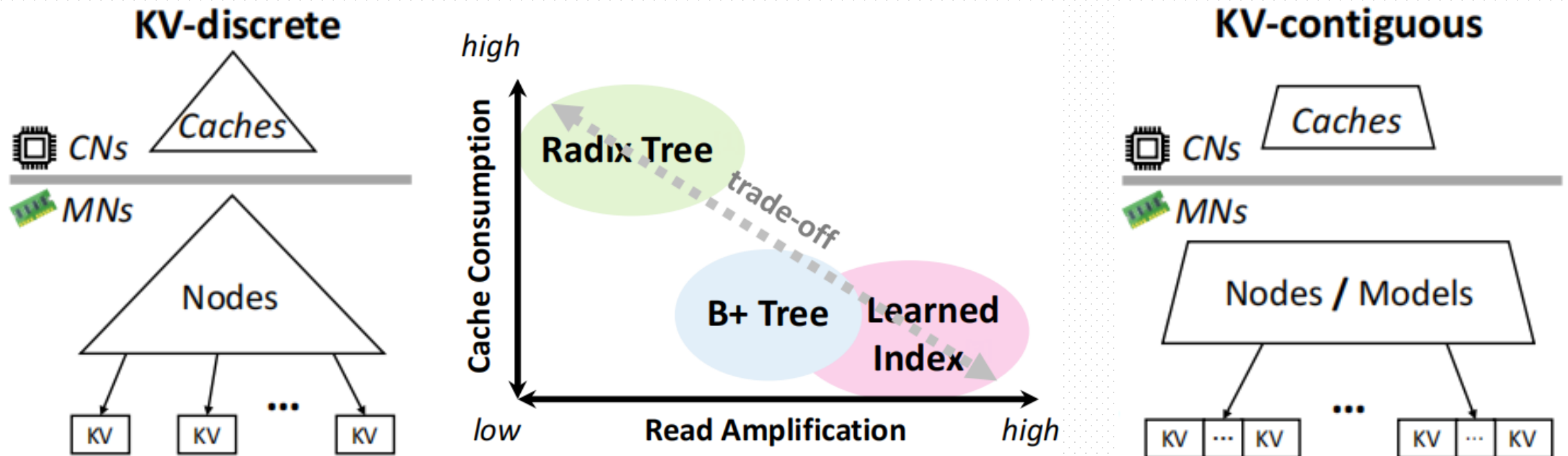[2] Qing Wang et al. Sherman: A write-optimized distributed B+ tree index on disaggregated memory. SIGMOD 2022.
[3] Pengfei Li et al. ROLEX: A scalable RDMA-oriented learned key-value store for disaggregated memory. FAST 2023.
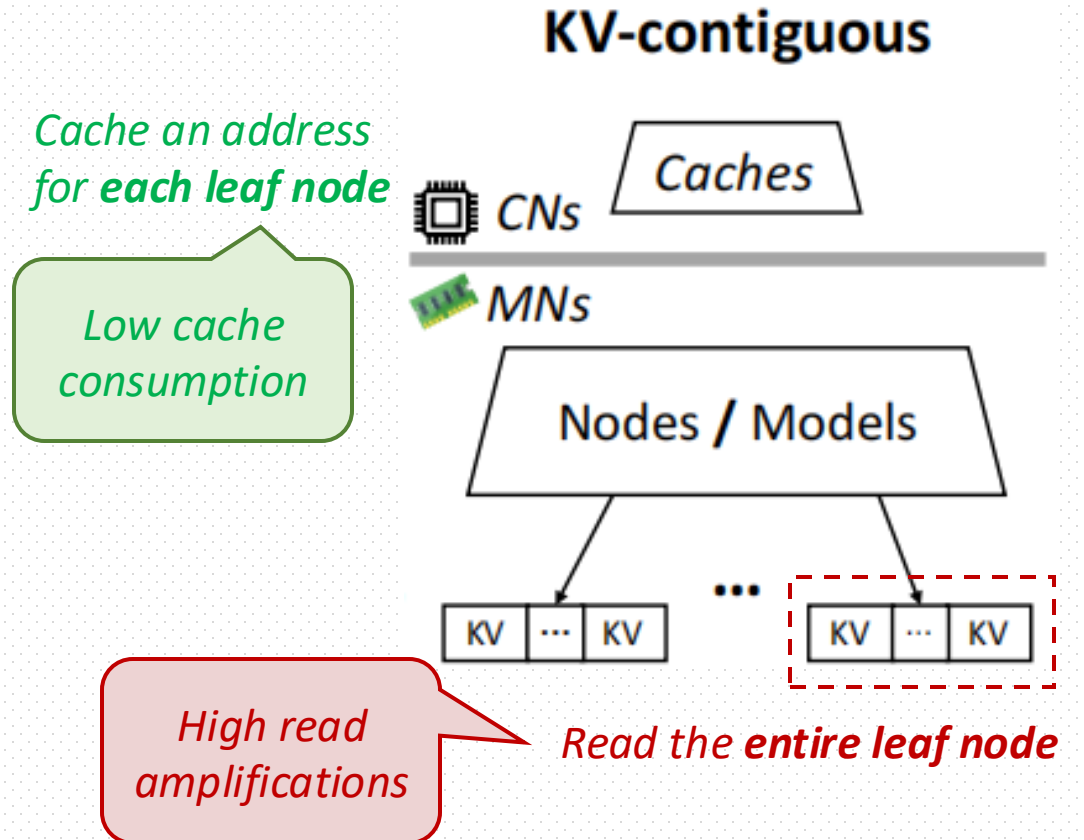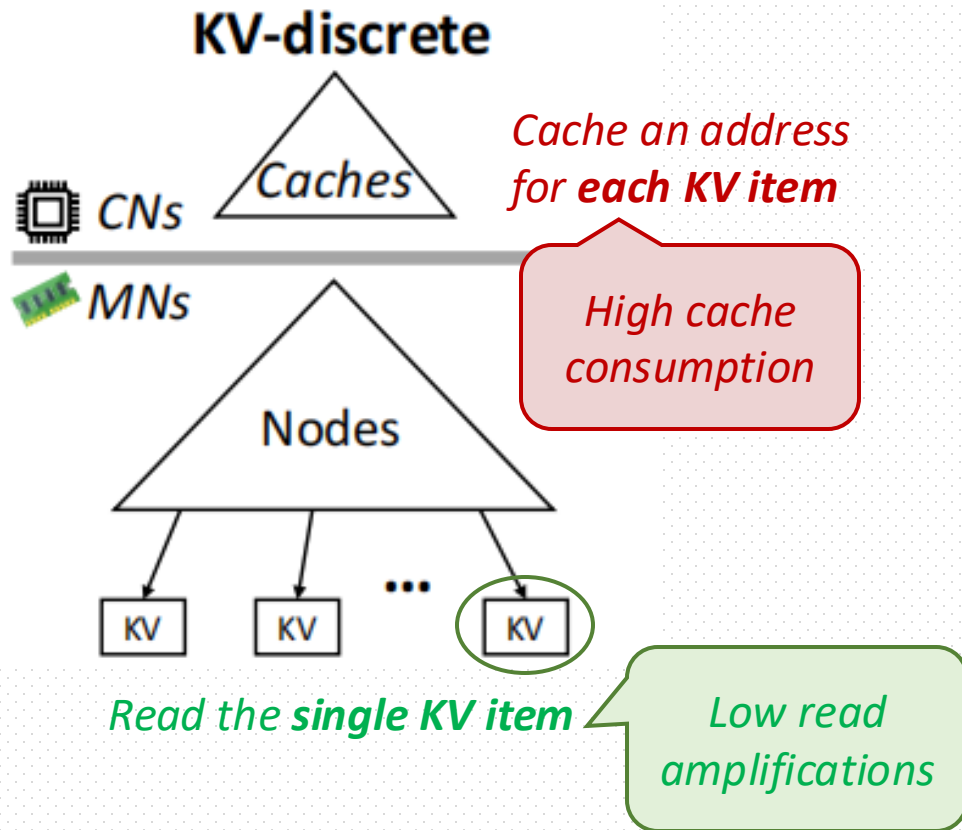
# A Trade-off for Range indexes on DM

**There is a trade-off between read amplifications and cache consumption:**

# A Trade-off for Range indexes on DM

**There is a trade-off between read amplifications and cache consumption:**



KV-discrete

CNs — Caches

Cache an address for **each KV item**

High cache consumption

MNs

Nodes

KV  KV  ...  KV

Read the **single KV item**

Low read amplifications

KV-contiguous

Cache an address for **each leaf node**

CNs — Caches

Low cache consumption

MNs

Nodes / Models

KV ... KV    KV ... KV

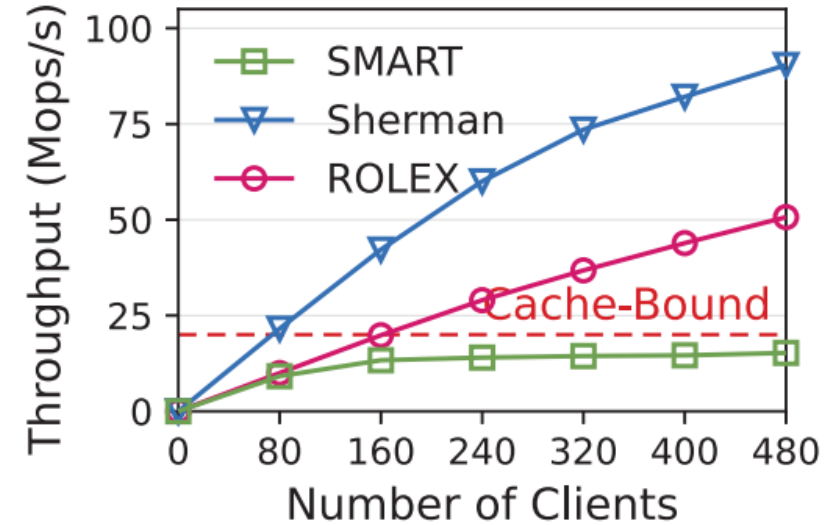High read amplifications

Read the **entire leaf node**

6

**There is a trade-off between read amplifications and cache consumption:**
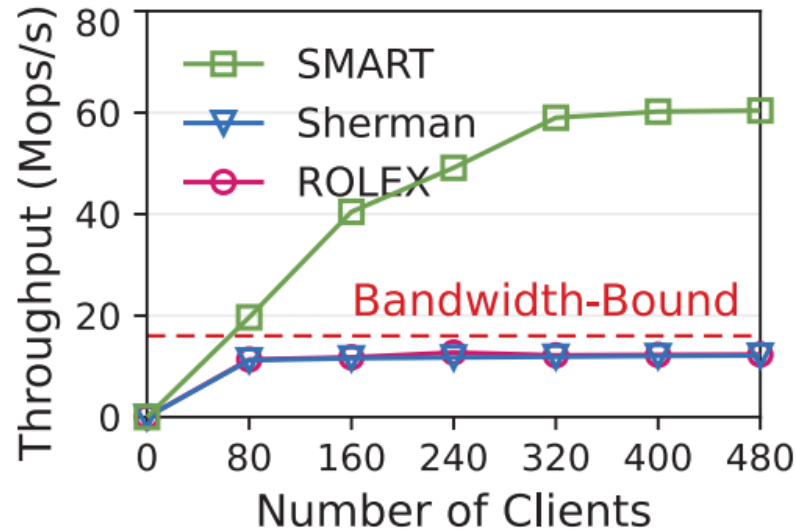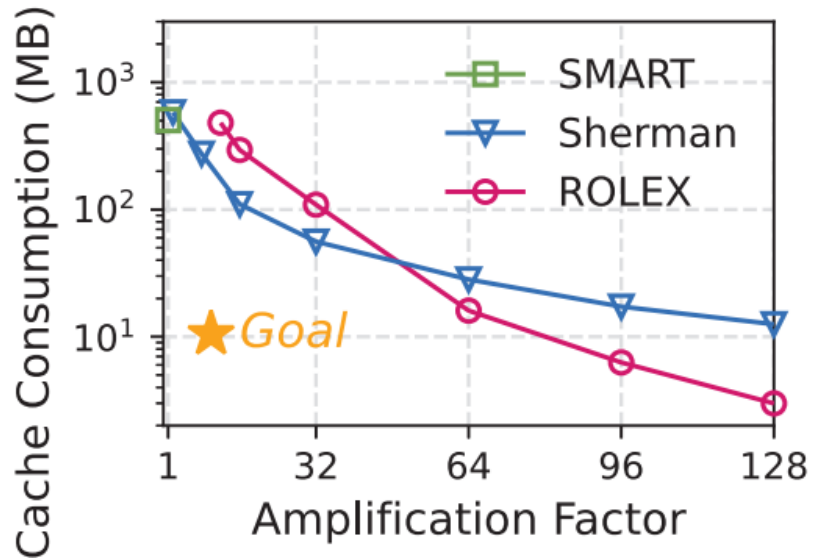
# A Trade-off for Range indexes on DM

**Root Cause:** The alignment between keys and memory addresses:



KV-discrete

$Key_2$

$Key_1$ ... $Key_{64}$

$min <= Key_i < max$ (i=1,2,...,64)

not aligned

$Addr_i$ is arbitrary

high

Cache Consumption

Radix Tree

trade-off

B+ Tree    Learned Index

low    Read Amplification    high

KV-contiguous

Leaf node:

$Key_1$ | $Key_2$ | ... | $Key_{64}$

$min <= Key_i < max$ (i=1,2,...,64)

aligned

$Addr_1 <= Addr_i < Addr_1 + leaf\_size$

## Use a KV-contiguous index (e.g., B+ tree) with hash-table-based leaf nodes

# Widely Choose a Suitable Hashing Scheme

## Remote access

*Few accesses:*

- *Simple associativity*
- *Hopscotch hashing[1]*
- *FaRM[2]*
- *RACE[3]*

## Read amplifications & Space efficiency



> *Hopscotch hashing best fits the requirements*

[1] Maurice Herlihy et al. Hopscotch hashing. DISC 2008.
[2] Aleksandar Dragojevic et al. FaRM: Fast Remote Memory. NSDI 2014.
[3] Pengfei Zuo et al. One-sided RDMA-conscious extendible hashing for disaggregated memory. ATC 2021.

# Viable Idea

## Use a hybrid index combining a B+ tree with hopscotch hashing

USTC, CHINA
ADSLAB

## Various granularities in reads/writes complicate optimistic synchronization

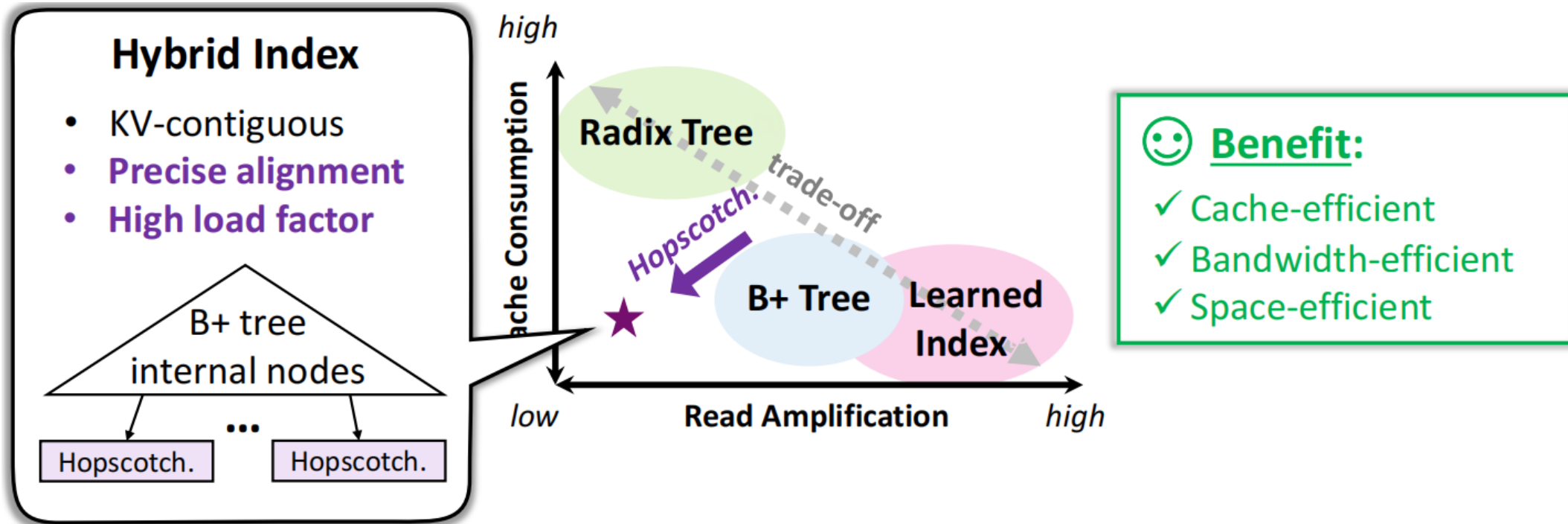Clients

Caches

CNs

MNs

B+ tree
internal nodes

Hopscotch. • • • Hopscotch.

Read a neighborhood
to **search a key**

neighborhood

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
|---|---|---|---|---|---|---|---|
|   | $A_1$ | $B_1$ | $C_1$ | $D_2$ | $E_4$ | $F_5$ |   | ...

*hopping*

*Insert a key $X_1$*

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
|---|---|---|---|---|---|---|---|
|   | $A_1$ | $B_1$ | $C_1$ | $X_1$ | $D_2$ | $F_5$ | $E_4$ | ...

hop range

Write an entire node
to **split the node**

Write a hop range
to **insert a key**

## Various granularities in reads/writes complicate optimistic synchronization



Clients

Caches

CNs

MNs

B+ tree internal nodes

Hopscotch.    ...    Hopscotch.

Read a neighborhood to **search a key**

Check

*Invisible to **fine-grained** neighborhood reads*

Data Block：

| ... | Version/checksum | ...... |

*Maintain*

*Difficult to maintain for **various** hop range writes*

Write an entire node to **split the node**

Write a hop range to **insert a key**

## Metadata for B+ trees and hopscotch hashing induces extra remote accesses

**Hopscotch hash still incurs read amplifications compared with reading KV**



*Search a key*

*Fetch the neighborhood*

| Metadata | | | $A_1$ | $B_1$ | $C_1$ | $X_1$ | |
|----------|--|--|-------|-------|-------|-------|--|

*Still need to fetch all items within the neighborhood*

Clients

Caches

CNs

MNs

B+ tree internal nodes

Hopscotch.   ...   Hopscotch.

# Challenge Summary

1. Complicated Optimistic Synchronization

2. Extra Metadata Accesses

3. Read Amplifications of Hopscotch Hashing

# Three-level Optimistic Synchronization

## Synchronization Overview

*Readers*

*Writers*

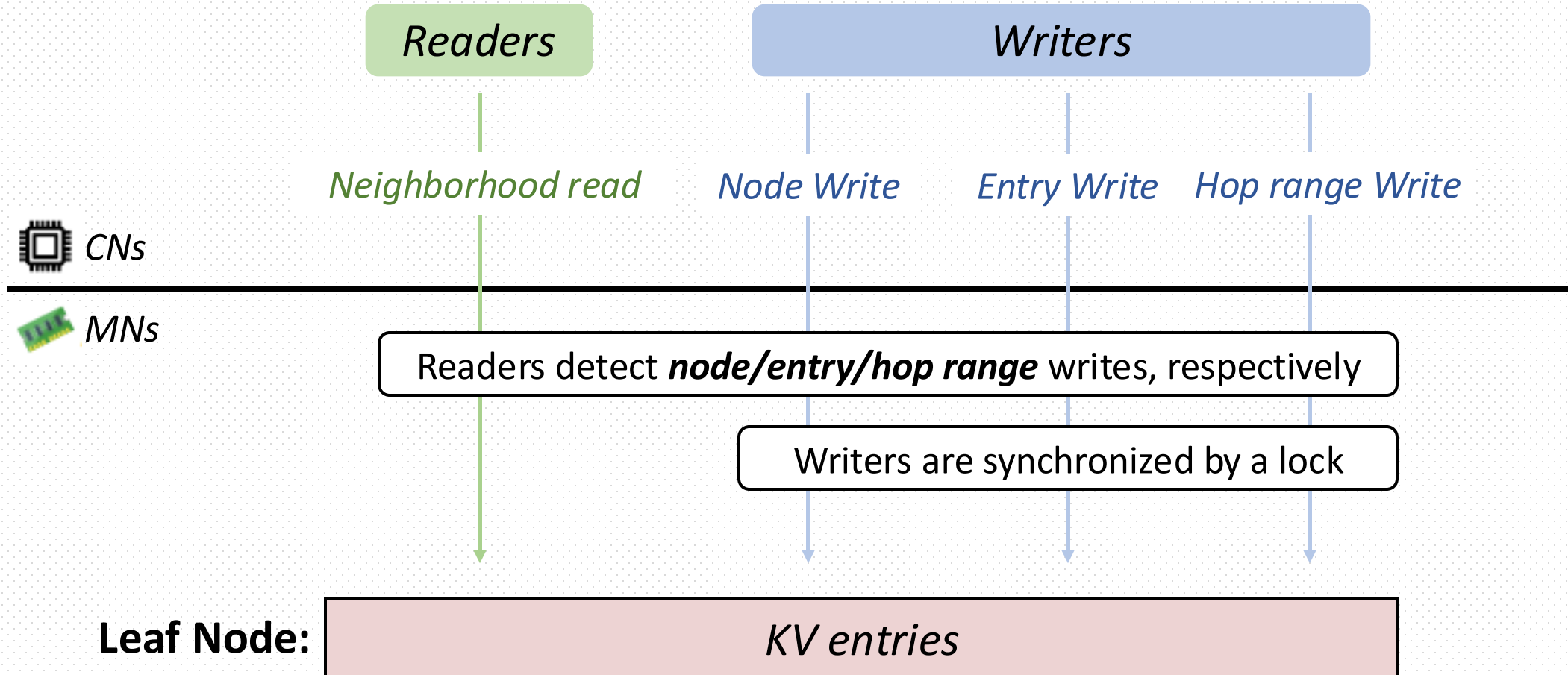*Neighborhood read*     *Node Write*     *Entry Write*     *Hop range Write*

CNs

MNs

Readers detect **node/entry/hop range** writes, respectively

Writers are synchronized by a lock

**Leaf Node:**     *KV entries*

# Three-level Optimistic Synchronization

## Level 1 & Level 2: Detect the Node Write & Entry Write



*<Search a key>*

*Reading a neighborhood*

**Node:**

**Cache Line**

**Entry:**

**Write the entire Node**

**Write an entry**

*<Split the Node>*

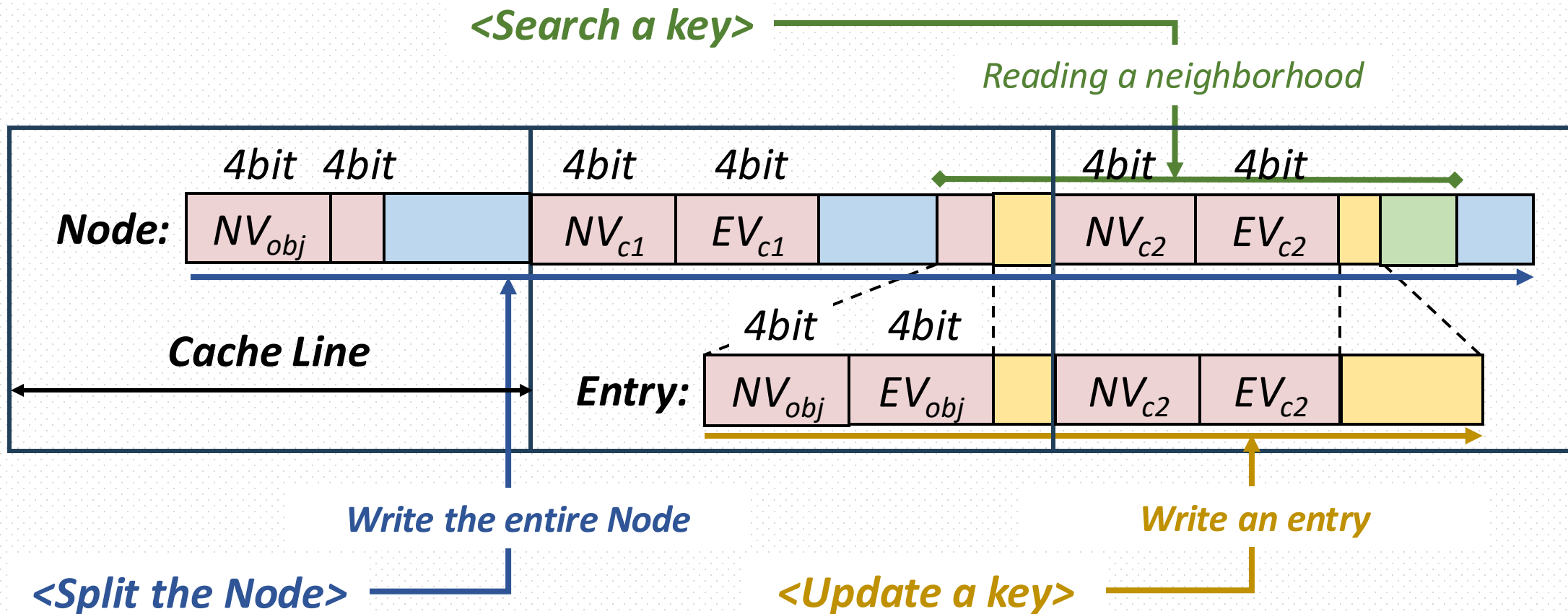*<Update a key>*

# Three-level Optimistic Synchronization

## Level 1 & Level 2: Detect the Node Write & Entry Write

➡ **_Solution:_** _Use two-level cache line versioning_

**<Search a key>**

_Reading a neighborhood_

| 4bit | 4bit | | 4bit | 4bit | | | 4bit | 4bit |
|---|---|---|---|---|---|---|---|---|

**Node:** $NV_{obj}$ $NV_{c1}$ $EV_{c1}$ $NV_{c2}$ $EV_{c2}$

**Cache Line**

| | 4bit | 4bit | | | |
|---|---|---|---|---|---|

**Entry:** $NV_{obj}$ $EV_{obj}$ $NV_{c2}$ $EV_{c2}$

**Write the entire Node**

**Write an entry**

**<Split the Node>**

**<Update a key>**

# Three-level Optimistic Synchronization

**Level 1 & Level 2: Detect the Node Write & Ent**
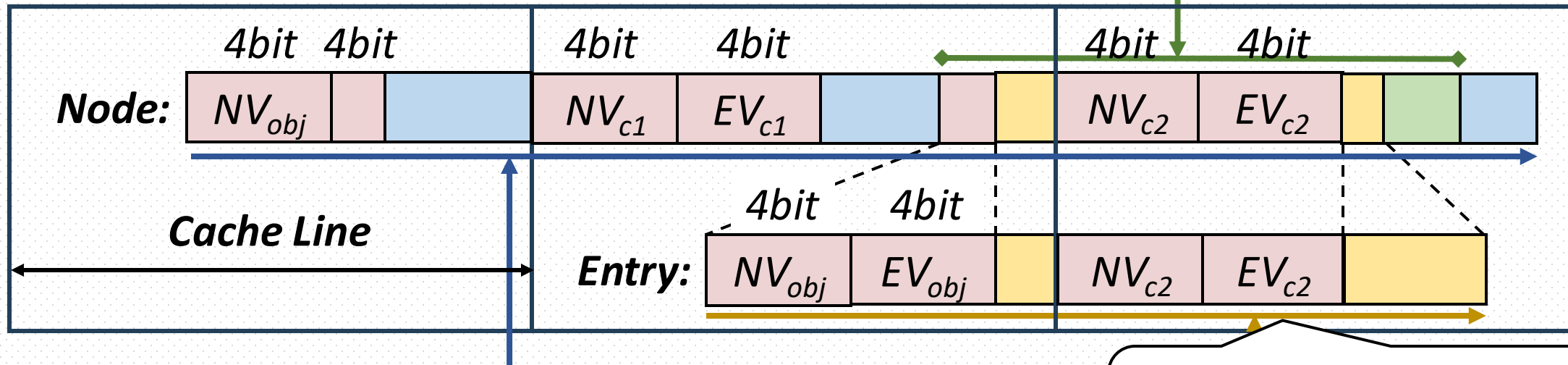
➡️ **_Solution:_** _Use two-level cache line versioning_

Readers:
- Check the **node write** via NVs
- Check the **entry write** via EVs

**_\<Search a key\>_**

_Reading a neighborhood_



**Node:**
| 4bit | 4bit | | 4bit | 4bit | | 4bit | 4bit |
| $NV_{obj}$ | | | $NV_{c1}$ | $EV_{c1}$ | | $NV_{c2}$ | $EV_{c2}$ |

**Cache Line**

**Entry:**
| 4bit | 4bit | | | |
| $NV_{obj}$ | $EV_{obj}$ | | $NV_{c2}$ | $EV_{c2}$ |

**_Write the entire Node_**

**_\<Split the Node\>_**

**_\<Update a key\>_**

Writers:
- Increment NVs during a **node write**
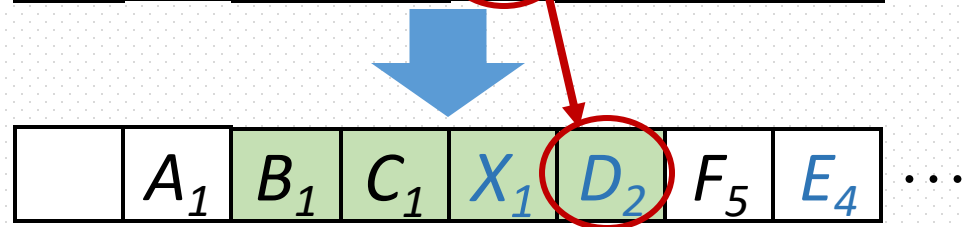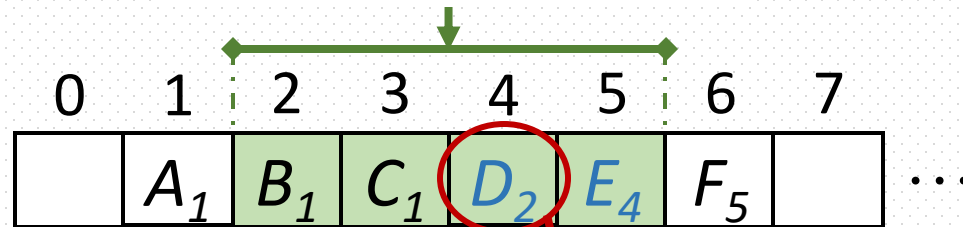- Increment EVs during an **entry write**

# Three-level Optimistic Synchronization

**Level 3: Detect the hop range write**     **Problem:** *Location changes of hooped items*

*<Search a key>*

*Reading a neighborhood*

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
|---|---|---|---|---|---|---|---|
|  | $A_1$ | $B_1$ | $C_1$ | $D_2$ | $E_4$ | $F_5$ |  |

|  | $A_1$ | $B_1$ | $C_1$ | $X_1$ | $D_2$ | $F_5$ | $E_4$ |

*Writing the hop range*

*<Insert a key>*

# Three-level Optimistic Synchronization
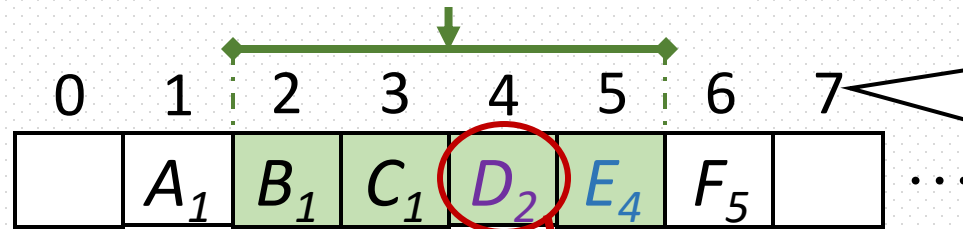
## Level 3: Detect the hop range write

**Problem:** *Location changes of hoped items*

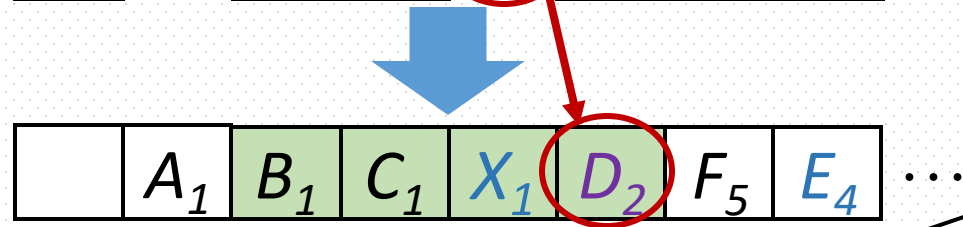➡ ***Solution*:** Reuse the hopscotch bitmaps

*<Search a key>*

*Reading a neighborhood*

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
|---|---|---|---|---|---|---|---|
|   | $A_1$ | $B_1$ | $C_1$ | $D_2$ | $E_4$ | $F_5$ |   |

...

Readers:
- **Re-construct** the bitmap according to fetched keys
- Retry if the two bitmaps cannot match

| | | | | | | | |
|---|---|---|---|---|---|---|---|
|   | $A_1$ | $B_1$ | $C_1$ | $X_1$ | $D_2$ | $F_5$ | $E_4$ |

...

Writers:
Hopscotch hashing **has maintained a bitmap inside each neighborhood** to track the occupancy status

*Writing the hop range*

*<Insert a key>*

**Metadata for hopscotch hashing**

**Problem:** Vacancy bitmaps induce extra accesses



*<Insert a key>*

CNs

MNs

**2. Fetch a vacancy bitmap**          **3. Fetch the hop range**          *1.Lock the Node*
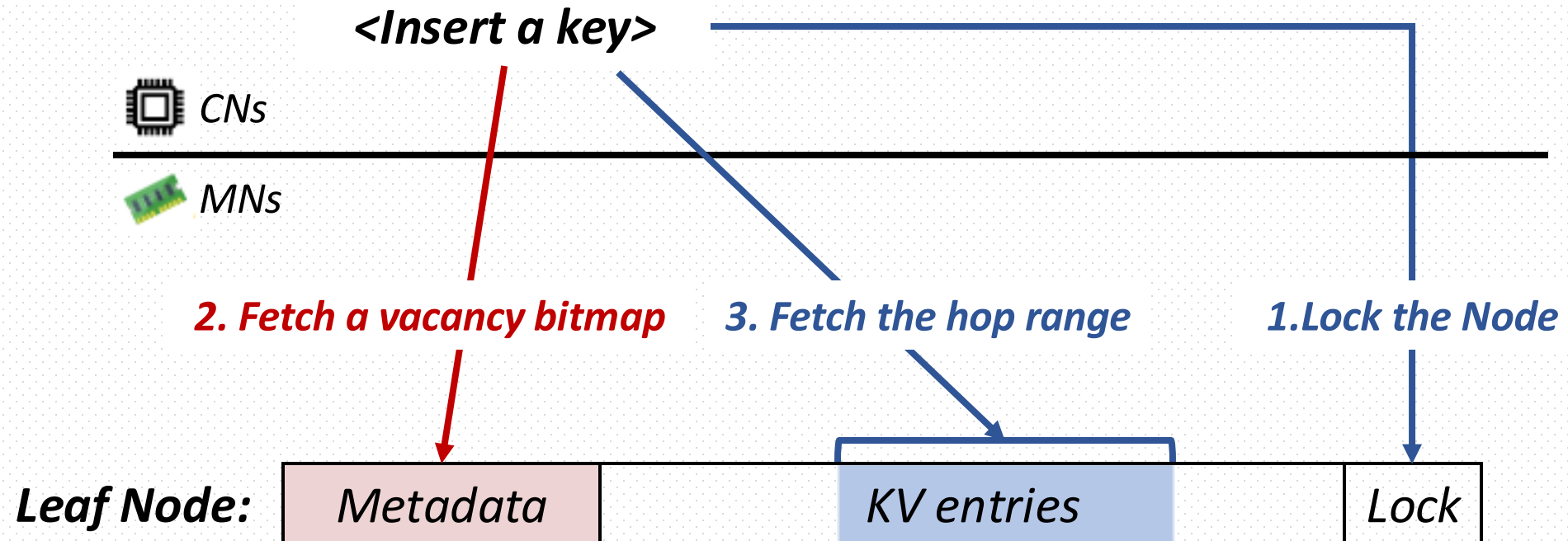
*Leaf Node:*     Metadata          KV entries          Lock
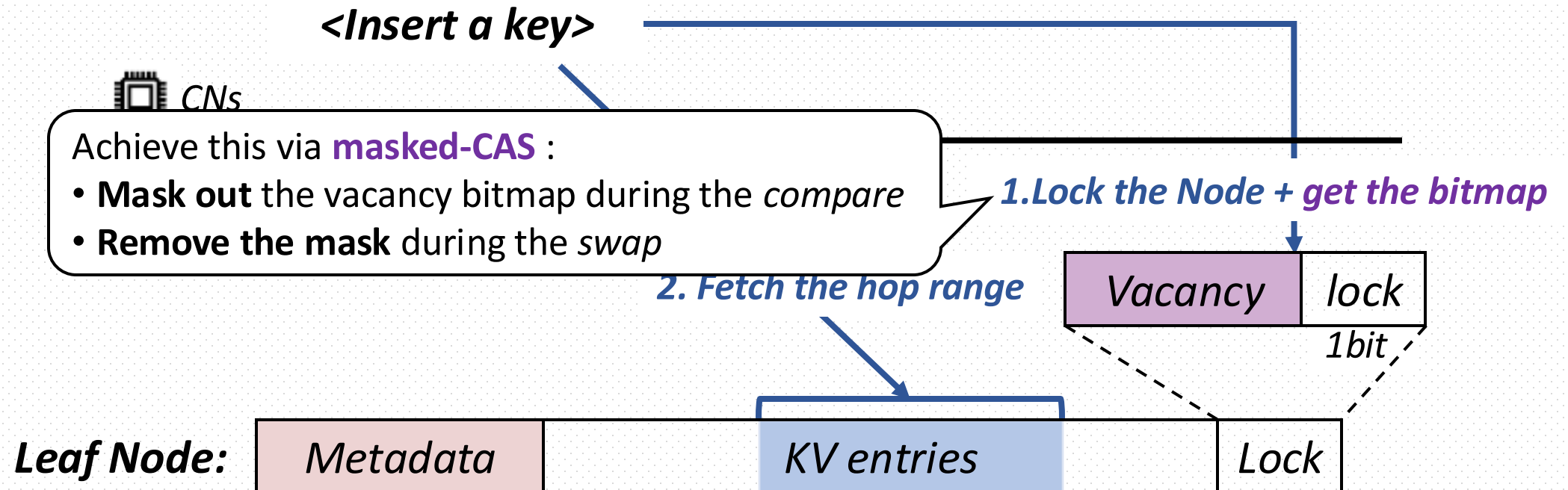
# Access-Aggregated Metadata Management



**Metadata for hopscotch hashing**

**Problem:** Vacancy bitmaps induce extra accesses
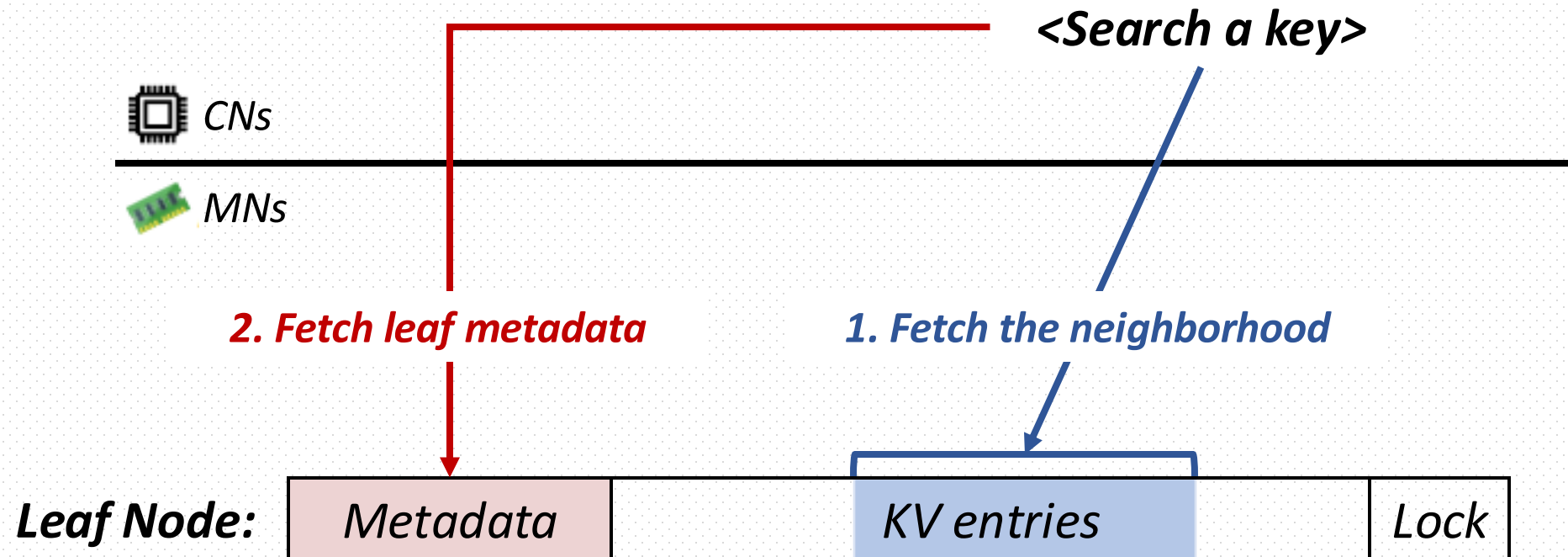
➡️ **Solution:** *Piggyback the vacancy bitmap*

*<Insert a key>*

CNs

Achieve this via **masked-CAS** :
- **Mask out** the vacancy bitmap during the *compare*
- **Remove the mask** during the *swap*

**1.Lock the Node + get the bitmap**

**2. Fetch the hop range**

| Vacancy | lock |

*1bit*

**Leaf Node:** | Metadata | | KV entries | | Lock |

**Metadata for the B+ tree**

**Problem:** Leaf metadata induce extra accesses



*<Search a key>*

CNs

MNs

**2. Fetch leaf metadata**

**1. Fetch the neighborhood**

*Leaf Node:* | Metadata | KV entries | Lock |

## Metadata for the B+ tree

**Problem:** Leaf metadata induce extra accesses

➡ **Solution:** *Replicate the leaf metadata*

**<Search a key>**

CNs

MNs

- Insert a leaf metadata **replica** at the position of every neighborhood size
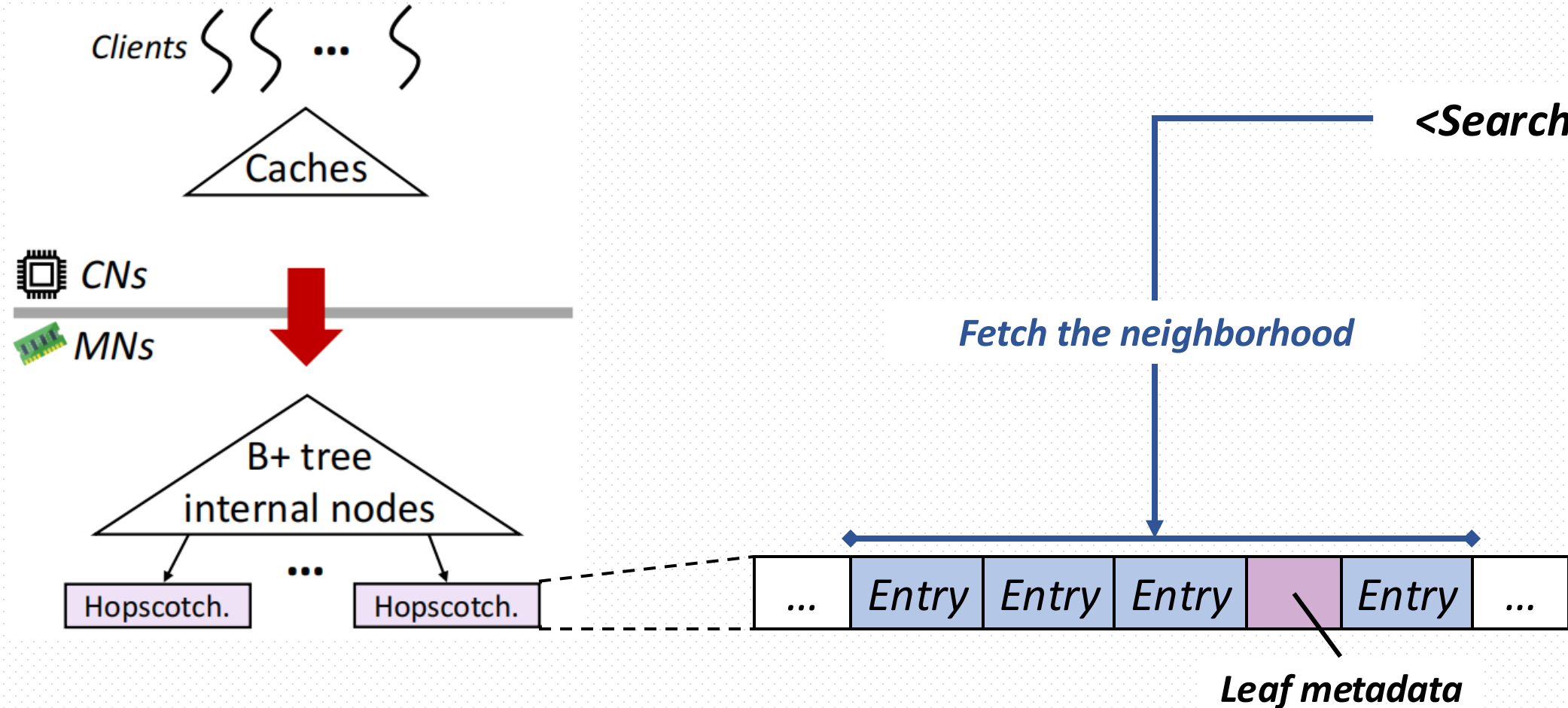
*Leaf metadata*

1. **Fetch the neighborhood + metadata**

**Leaf Node:** | | | KV entries | | | Lock |

**USTC, CHINA ADSLAB**

**Metadata for the B+ tree**

**Problem:** Still need to fetch all items within the neighborhood



Clients

Caches

CNs

MNs

B+ tree internal nodes

Hopscotch. ... Hopscotch.

*<Search a key>*

*Fetch the neighborhood*

... | *Entry* | *Entry* | *Entry* | | *Entry* | ...

*Leaf metadata*

# Hotness-Aware Speculative Read

**Metadata for the B+ tree**

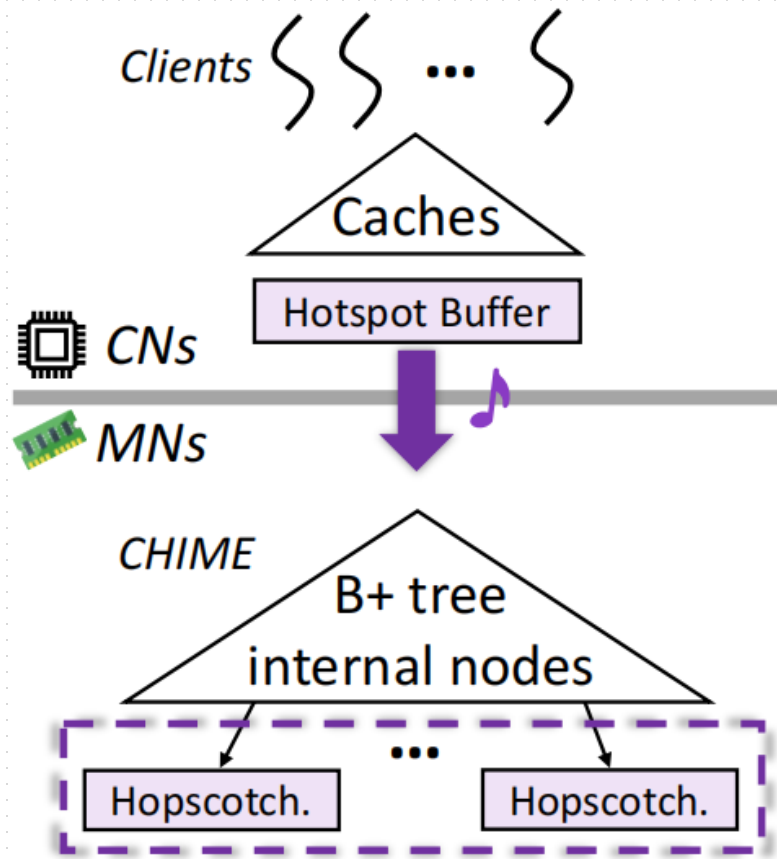**Problem:** Still need to fetch all items within the neighborhood

➡️ **Solution:** *Speculatively read the target entry*



*<Search a key>*

**Each buffer entry:**

| Leaf address | Key index | fingerprint | counter |
|---|---|---|---|

*Speculatively read an entry*

| ... | Entry | Entry | Entry | | Entry | ... |
|---|---|---|---|---|---|---|

*hotspot*

*Leaf metadata*

# Optimization Summary

1. Complicated Optimistic Synchronization

*Solution 1:* *Three-Level Optimistic Synchronization*

2. Extra Metadata Accesses

*Solution 2:* *Access-Aggregated Metadata Management*

3. Read Amplifications of Hopscotch Hashing

*Solution 3*: *Hotness-Aware Speculative Read*

# Evaluation

## Workloads and Parameters

- YCSB workloads
- 8-byte keys and 8-byte values
- Limit the cache size to 100 MB per CN
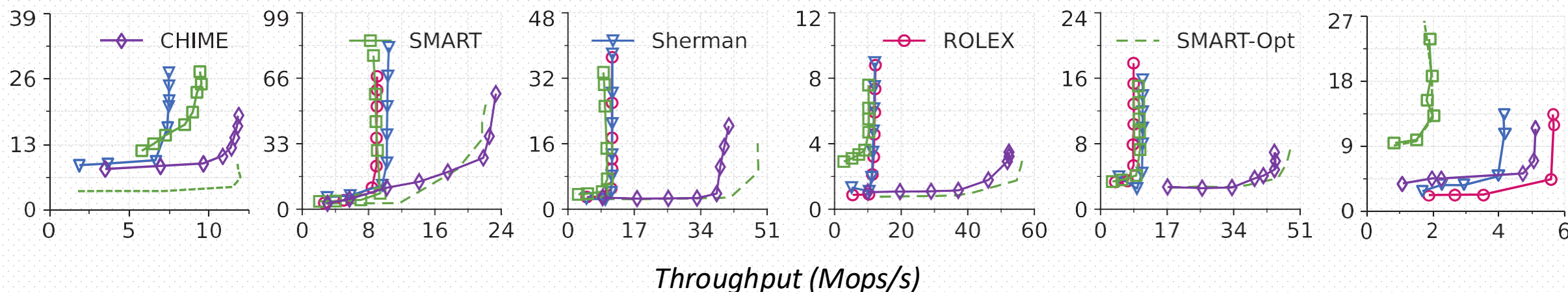
## Comparisons

- SMART[*OSDI' 23*]
  - The latest radix tree design on DM
- Sherman[*SIGMOD' 22*]
  - The classic B+ tree design on DM
- ROLEX[*FAST' 23*]
  - The latest learned index on DM
- SMART-Opt[*Optimal case*]
  - SMART with sufficient caches

# Performance Comparison

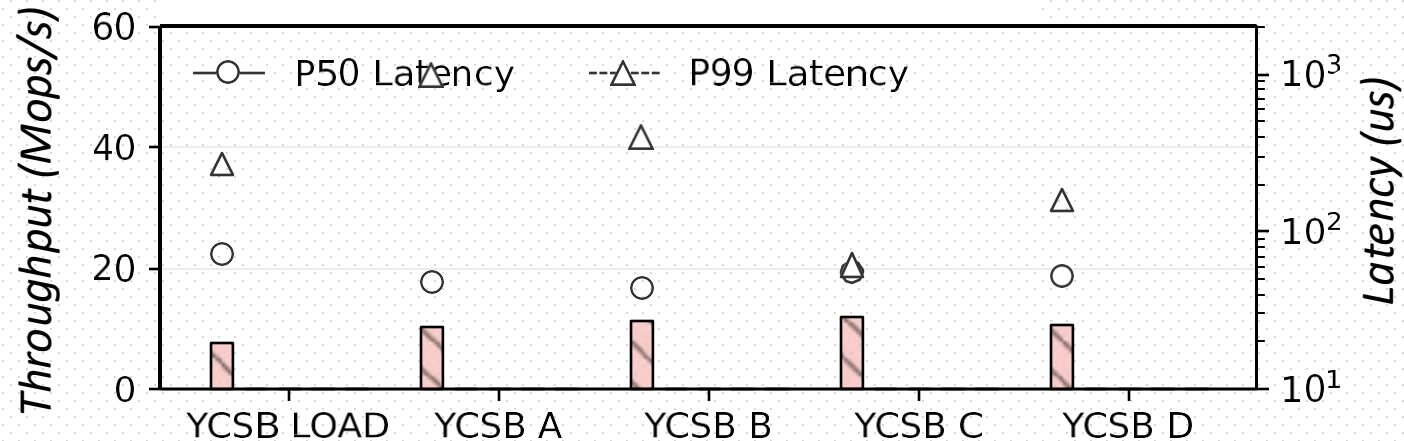| YCSB LOAD | YCSB A | YCSB B | YCSB C | YCSB D | YCSB E |
|---|---|---|---|---|---|
| 100% insert | 50% read 50% update | 95% read 5% update | 100% read | 95% read 5% insert | 95% scan 5% insert |



- CHIME achieves:
  - Up to 4.3x higher throughput than Sherman and ROLEX
  - Up to 5.1x higher throughput than SMART
  - A close performance to the optimal case, with up to 8.7x lower cache consumption *(57.6 MB vs. 503.6 MB)*

# Factor Analysis

| YCSB LOAD | YCSB A | YCSB B | YCSB C | YCSB D |
|---|---|---|---|---|
| 100% insert | 50% read<br>50% update | 95% read<br>5% update | 100% read | 95% read<br>5% insert |

Sherman



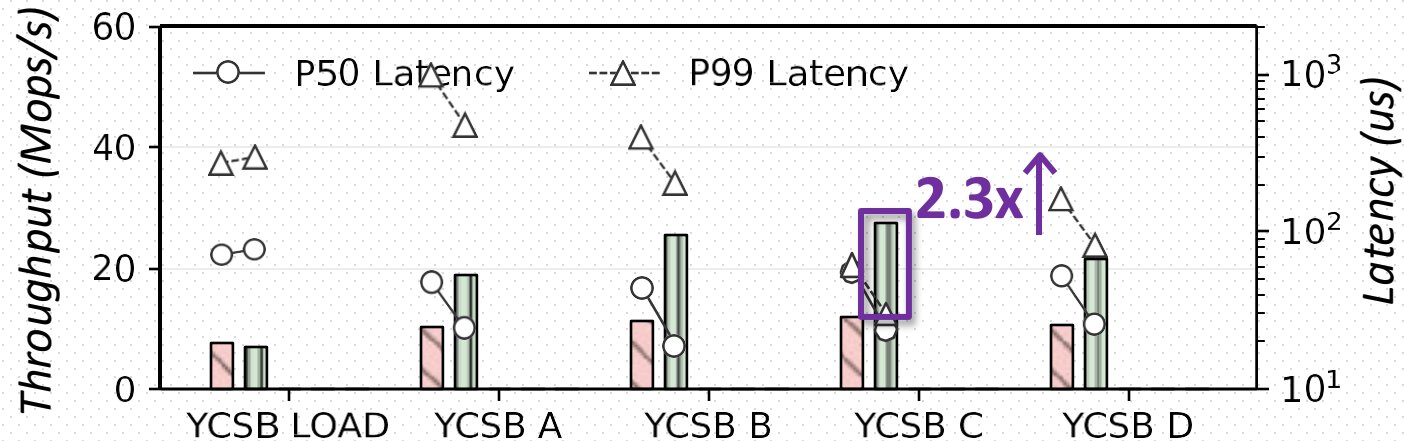- Start with Sherman and apply each proposed technique one by one

# Factor Analysis

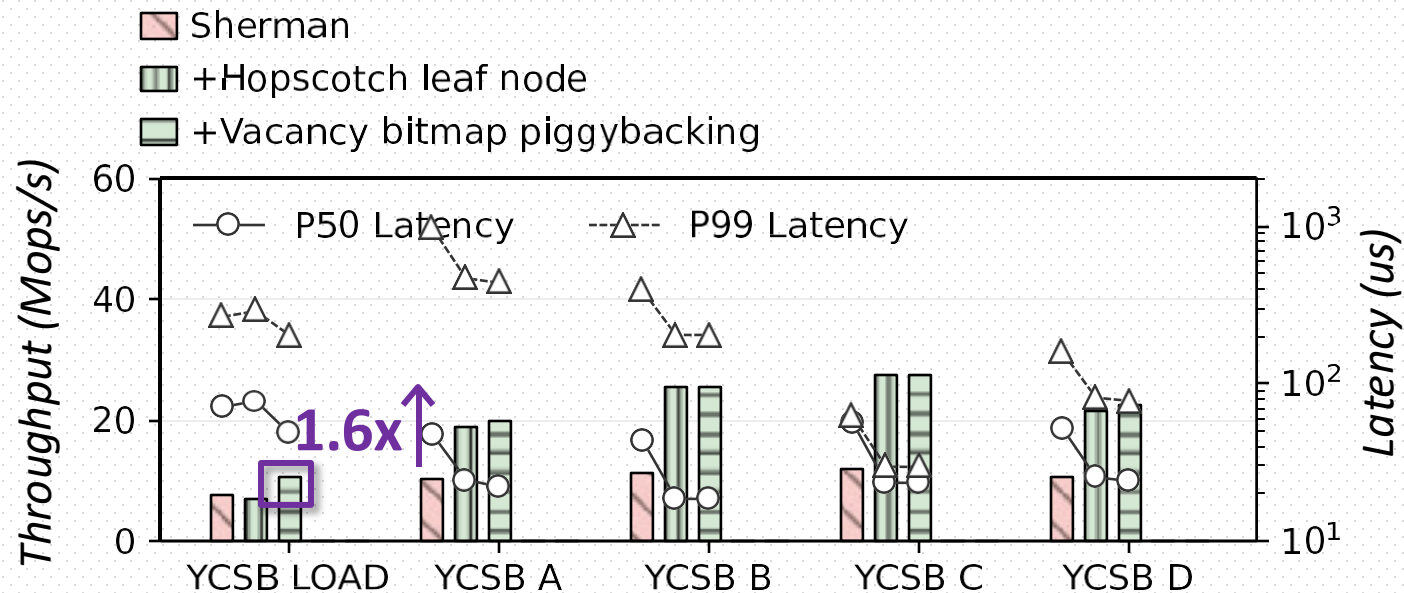| YCSB LOAD | YCSB A | YCSB B | YCSB C | YCSB D |
|---|---|---|---|---|
| 100% insert | 50% read 50% update | 95% read 5% update | 100% read | 95% read 5% insert |



- The *hopscotch leaf node* enables fetching the neighborhood rather than the entire leaf node

33

# Factor Analysis

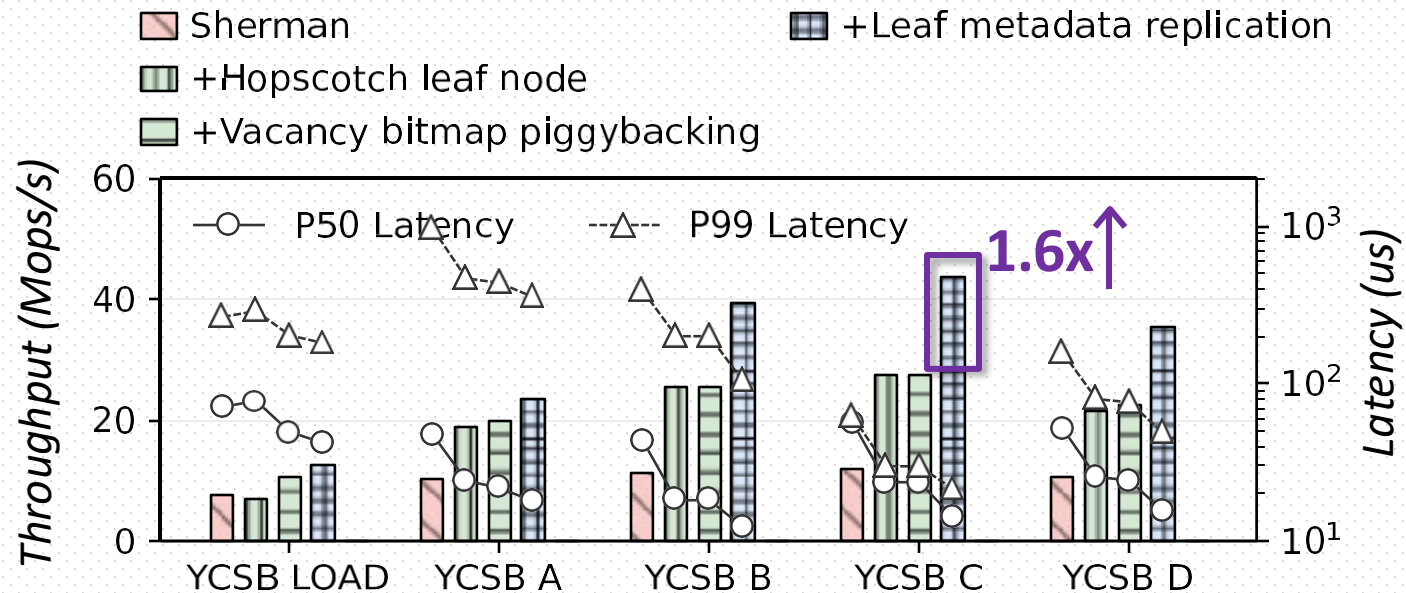| YCSB LOAD | YCSB A | YCSB B | YCSB C | YCSB D |
|:---:|:---:|:---:|:---:|:---:|
| 100% insert | 50% read<br>50% update | 95% read<br>5% update | 100% read | 95% read<br>5% insert |



- The *vacancy bitmap piggybacking* enables fetching the hop range rather than the entire leaf node, without inducing extra remote accesses

34

# Factor Analysis



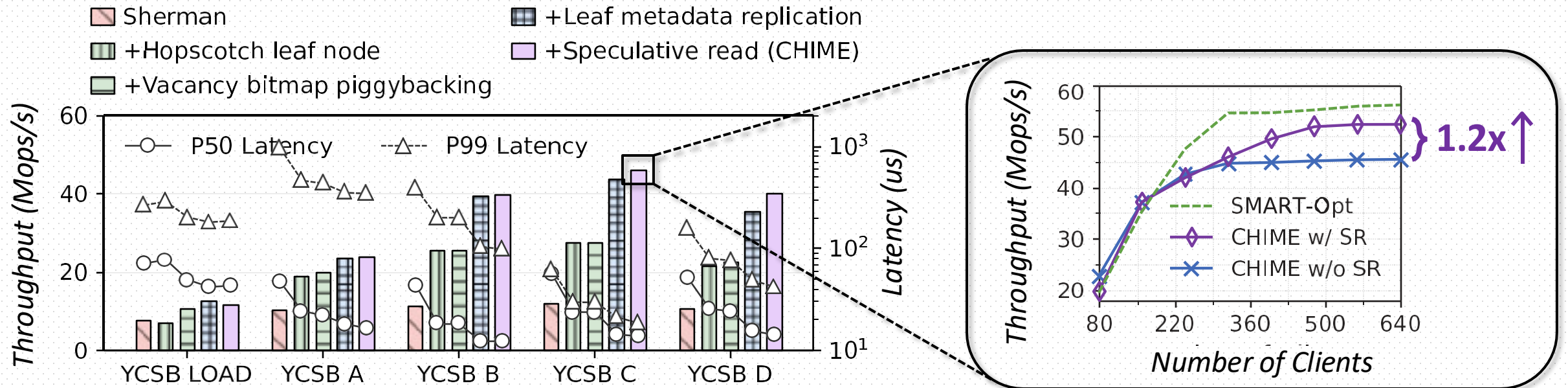| YCSB LOAD | YCSB A | YCSB B | YCSB C | YCSB D |
|---|---|---|---|---|
| 100% insert | 50% read 50% update | 95% read 5% update | 100% read | 95% read 5% insert |



- The *leaf metadata replication* avoids the extra remote accesses of fetching in-header leaf metadata

# Factor Analysis



| YCSB LOAD | YCSB A | YCSB B | YCSB C | YCSB D |
|---|---|---|---|---|
| 100% insert | 50% read 50% update | 95% read 5% update | 100% read | 95% read 5% insert |



- The *speculative read* enables greedily fetching the target entry rather than the entire neighborhood

36

# Conclusion

- This paper identifies the **trade-off** between read amplifications and cache consumption for range indexes on DM

- We propose **CHIME**, a hybrid index combining the B+ tree with hopscotch hashing to break the trade-off:
  - Three-level optimistic synchronization
  - Access-aggregated metadata management
  - Hotness-aware speculative read

- CHIME outperforms the state-of-the-art range indexes on DM by up to **5.1x in throughput** with the same cache size and achieves similar performance with up to **8.7x lower cache consumption**

# CHIME: A Cache-Efficient and High-Performance Hybrid Index on Disaggregated Memory

**Xuchuan Luo,** Jiacheng Shen, Pengfei Zuo, Xin Wang, Micheal R.Lyu, Yangfan Zhou

***School of Computer Science, Fudan University***
*National Key Laboratory of Parallel and Distributed Computing, China*
*Duke Kunshan University  Huawei Cloud The Chinese University of Hong Kong*
*Shanghai Key Laboratory of Intelligent Information Processing, Shanghai, China*

Thanks you for your attention