

Efficient Vector Search on Disaggregated Memory with d-HNSW

Yi Liu*, Fei Fang*, Chen Qian
University of California Santa Cruz

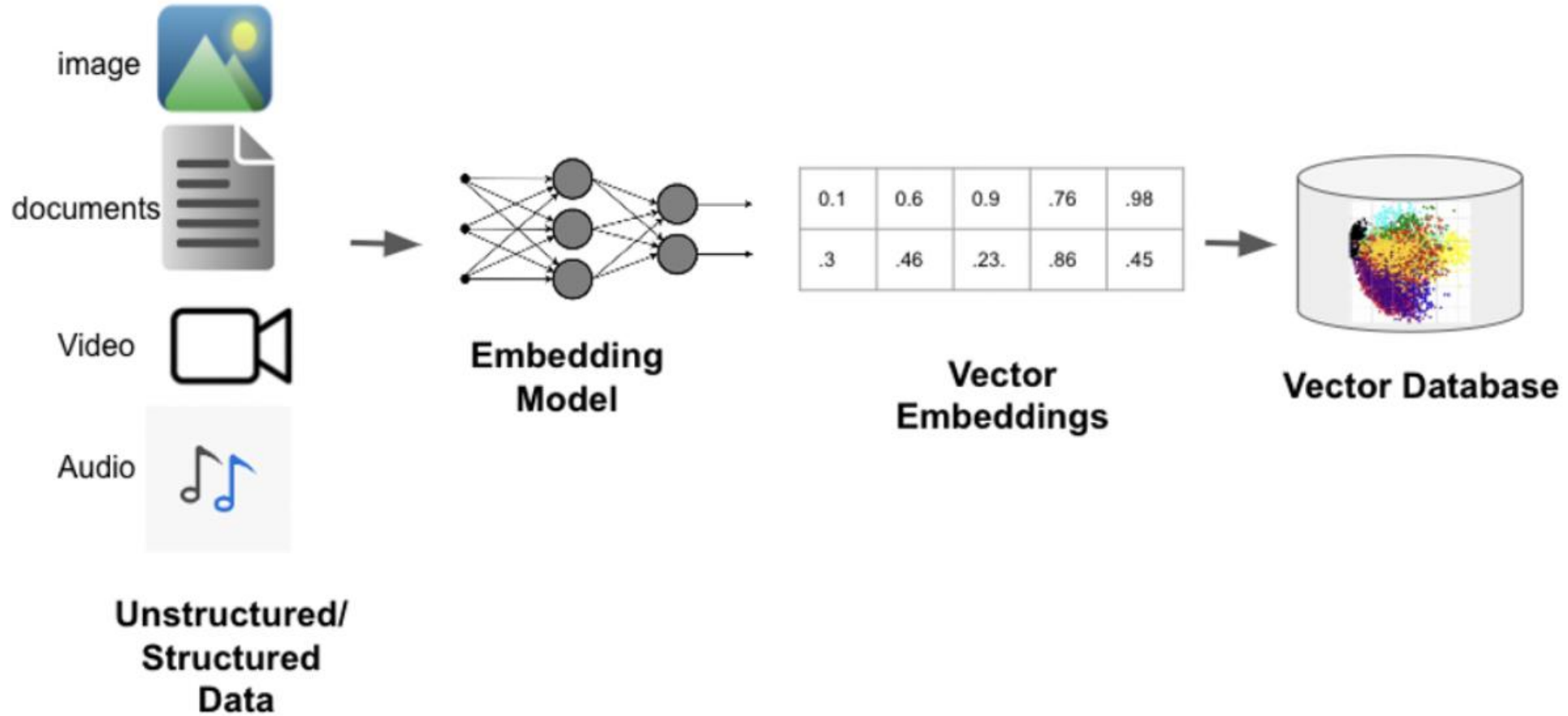


The 17th ACM Workshop on
Hot Topics in Storage and File Systems

July 10-11
Boston, MA



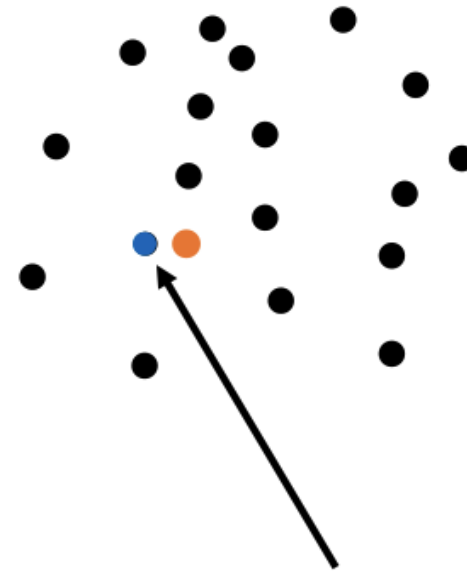
Vector DataBases



Vector Query

Given a query vector

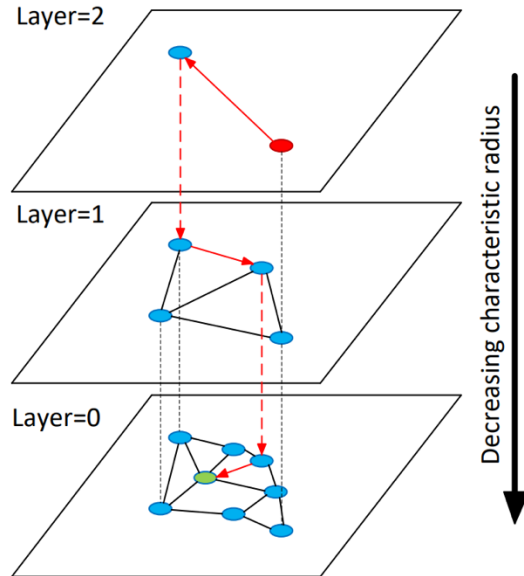
Return Top-k Nearest
vectors



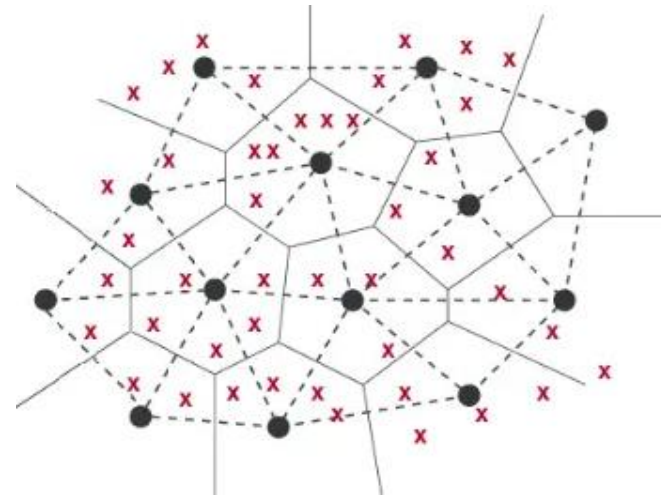
Top-1(k=1) Nearest Vector

Approximate Vector Query

Graph Index



Inverted Index

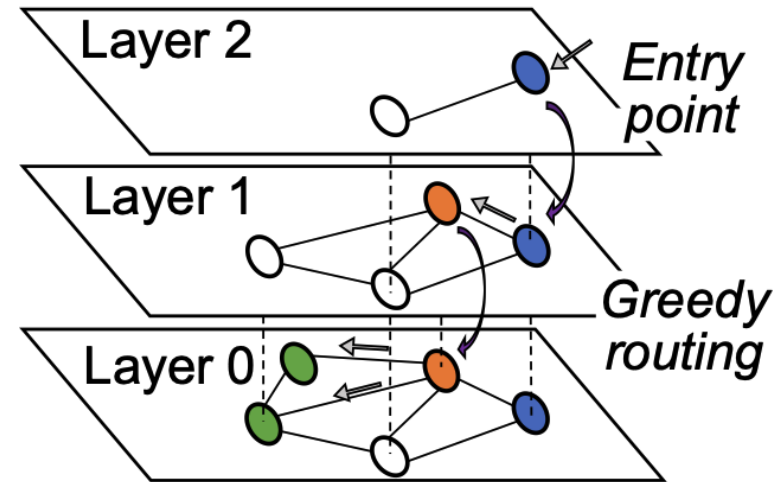


Hierarchical Navigable Small World (HNSW)

similarity search

trade off: latency vs
accuracy

Robustness and high
performance

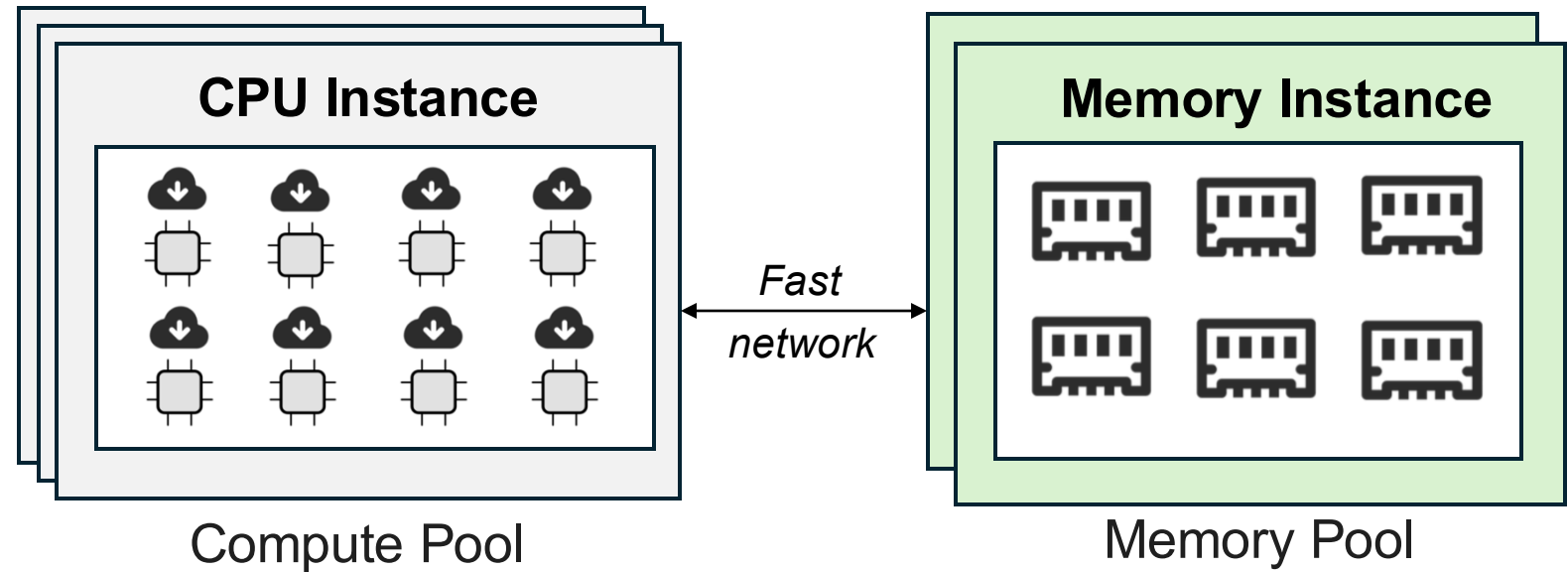


Disaggregated Memory Systems (DMS)

High resource
utilization

Flexible hardware
scalability

Efficient data sharing



RDMA-based Disaggregated Memory System

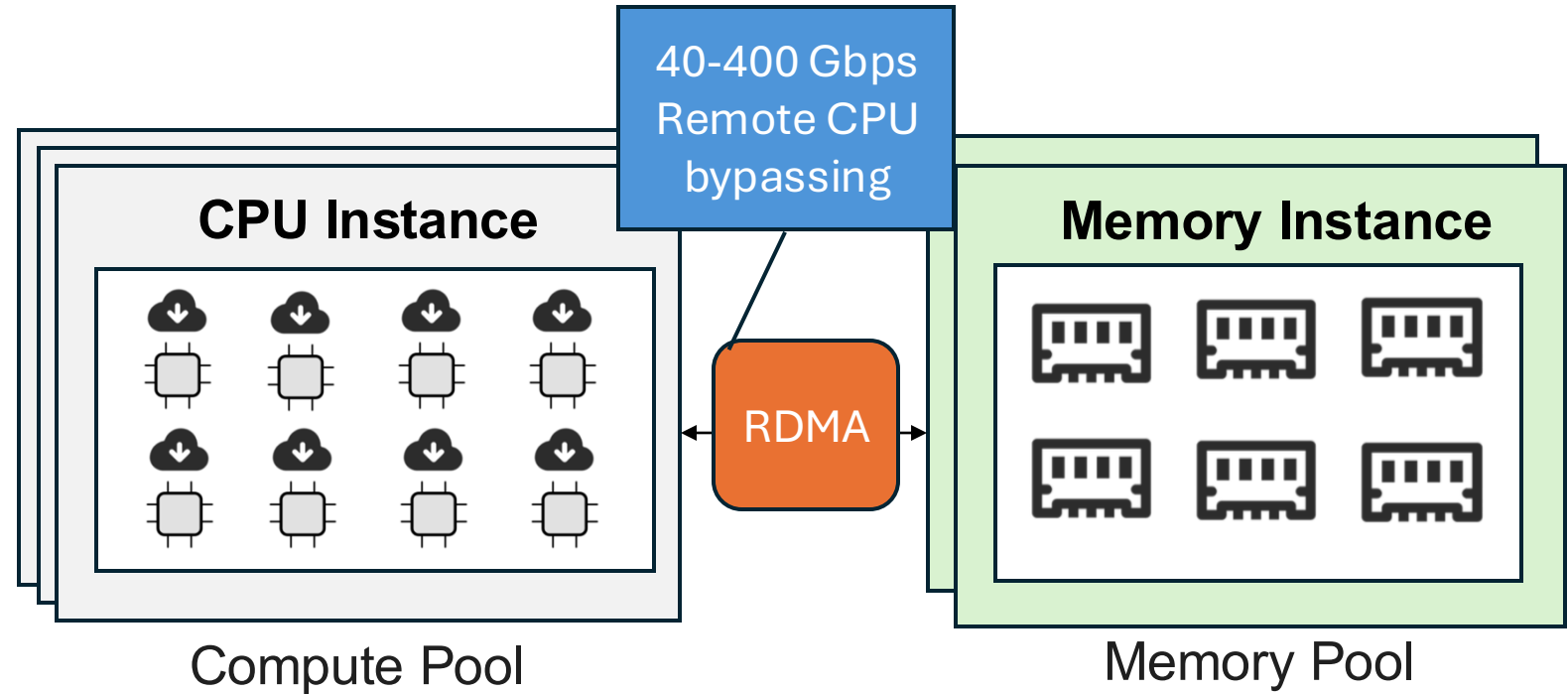
file system

key-value stores

transactional
databases

...

vector databases



Vector Database on Disaggregated Memory System

- Challenge 1: How to reduce network round trips?
- Challenge 2: How to enable one-side insertions?
- Challenge 3: How to support efficient batched operations?

Reduce Network Round Trips Challenge

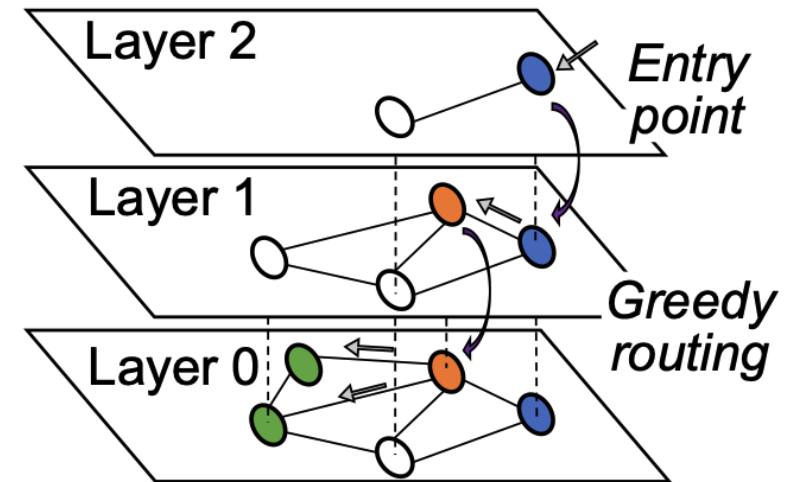
greedy algorithm

node represents
vector

unpredictable
traversal path

fetch each step

**excessive
round-trips**

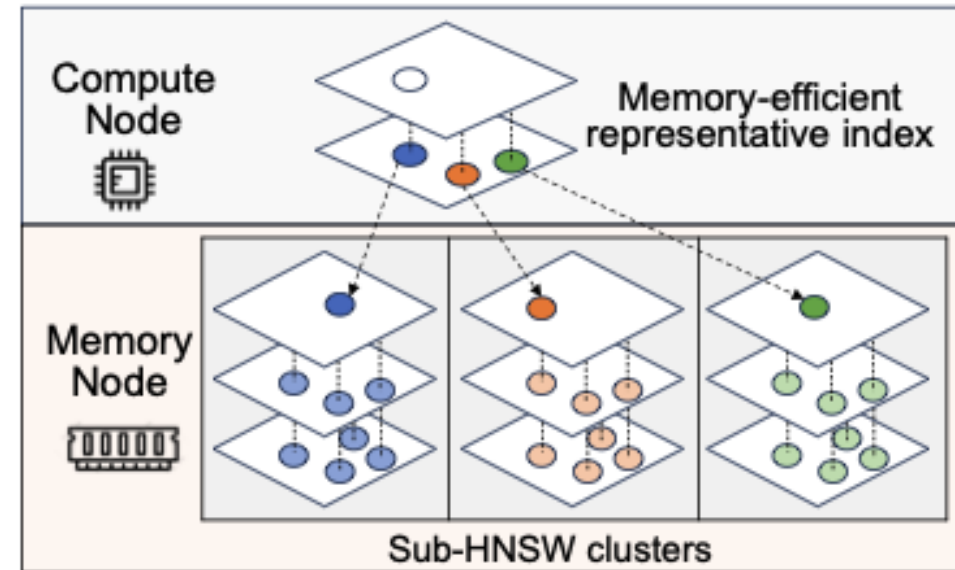


Reduce Network Round Trips Challenge

- Our approach: **Representative index caching**

a lightweight meta-HNSW

- **Minimizes network transfer** of irrelevant vectors
- **Reduces latency** and bandwidth usage
- **Preserves high recall**



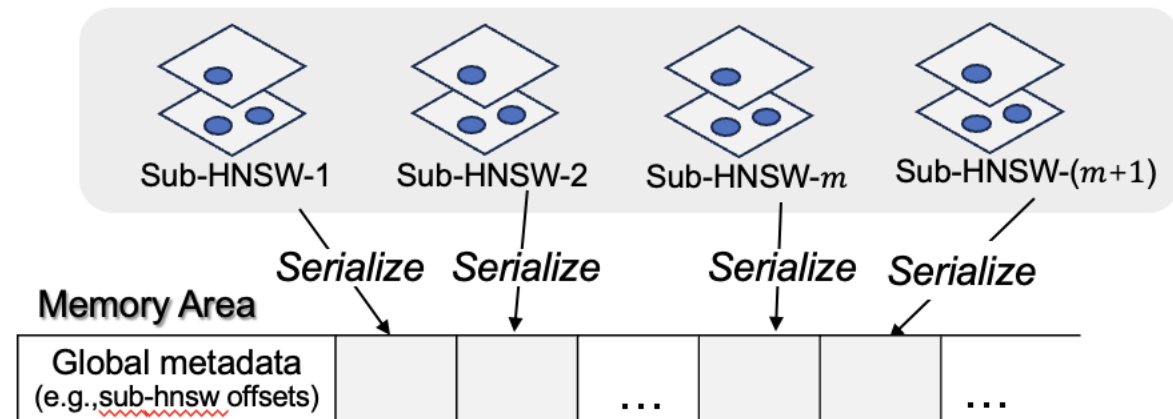
One-side Insertions Challenge

compactly serialized

new vector insertion

scatters vectors

High latency from fragmented memory access and multiple RDMA round-trips.



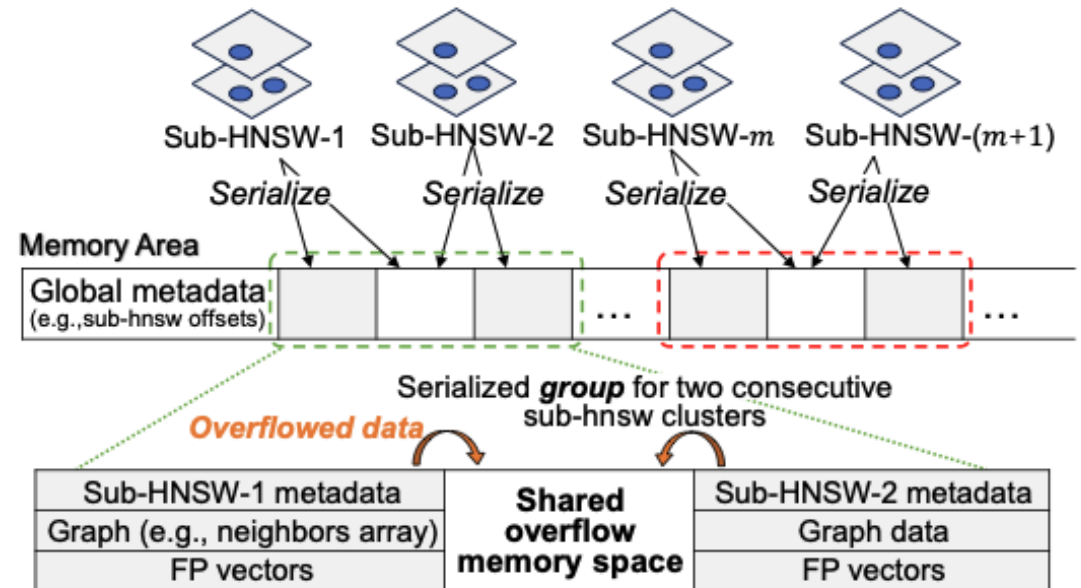
One-side Insertions Challenge

- Our approach: **RDMA-friendly graph index storage layout in remote memory.**

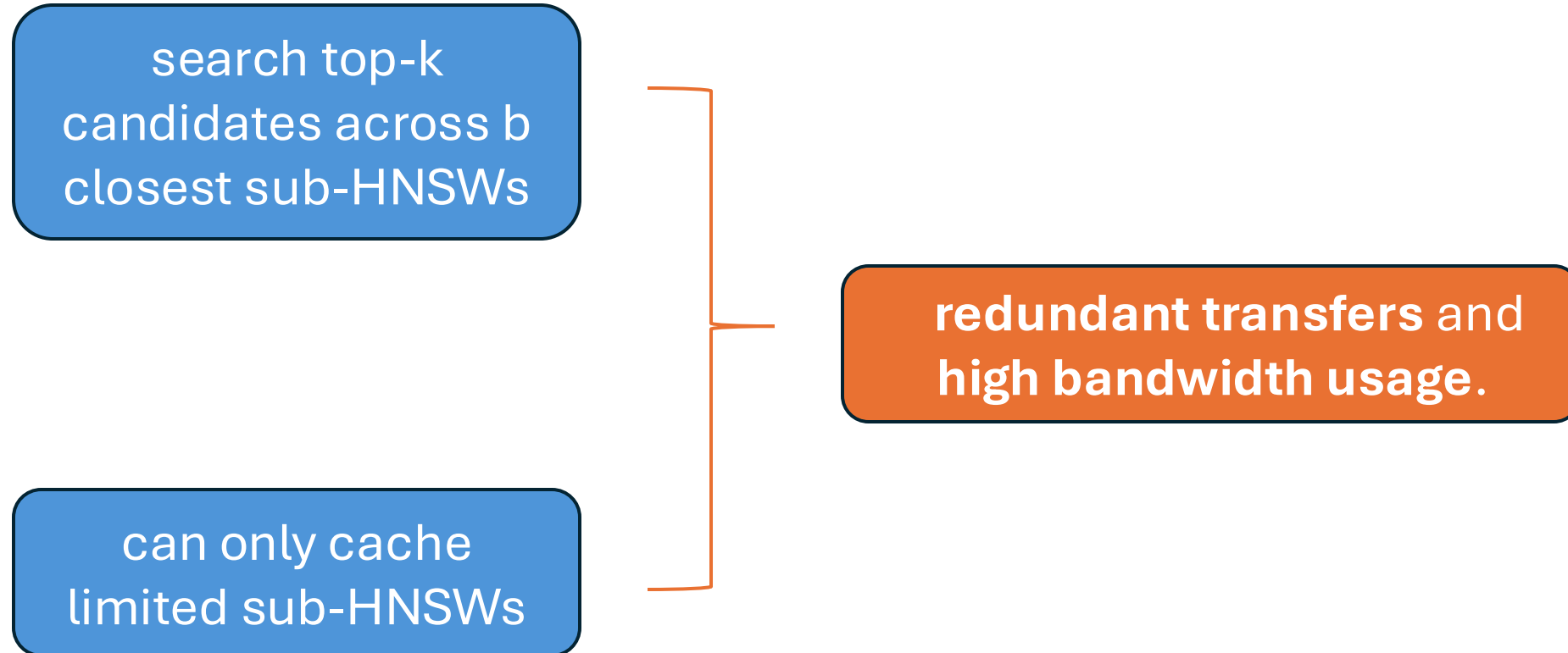
Internal gap

Shared overflow
memory space

- Enables **single-round RDMA reads** for query
- **Avoid fragmented memory access**
- **Preserves high throughput** under insertions
- **Balances insertions between two sub-HNSW**



Efficient Batched Operations Challenge



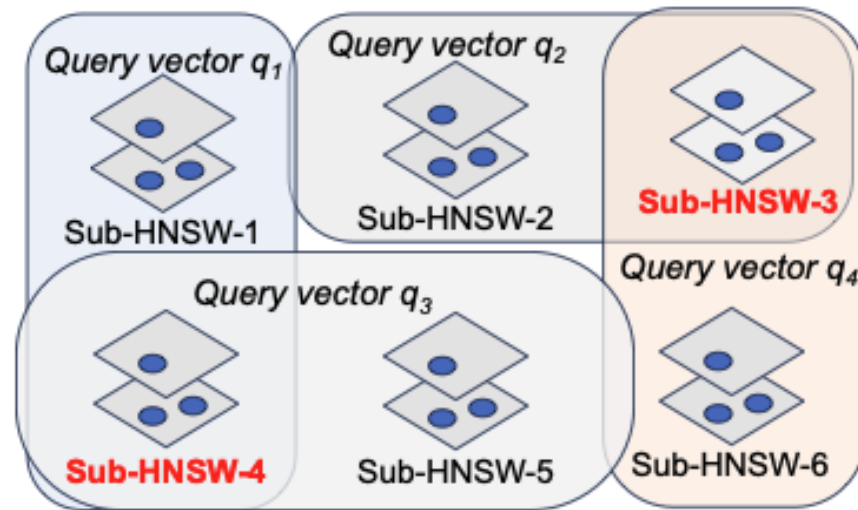
Efficient Batched Operations Challenge

- Our approach: **Query-aware batched data loading**

required sub-
HNSW **only once**
per batch

doorbell
batching

- **Minimizes redundant data transfers**
- **Reduces network round-trips** with batching
- **Preserves cache efficiency**

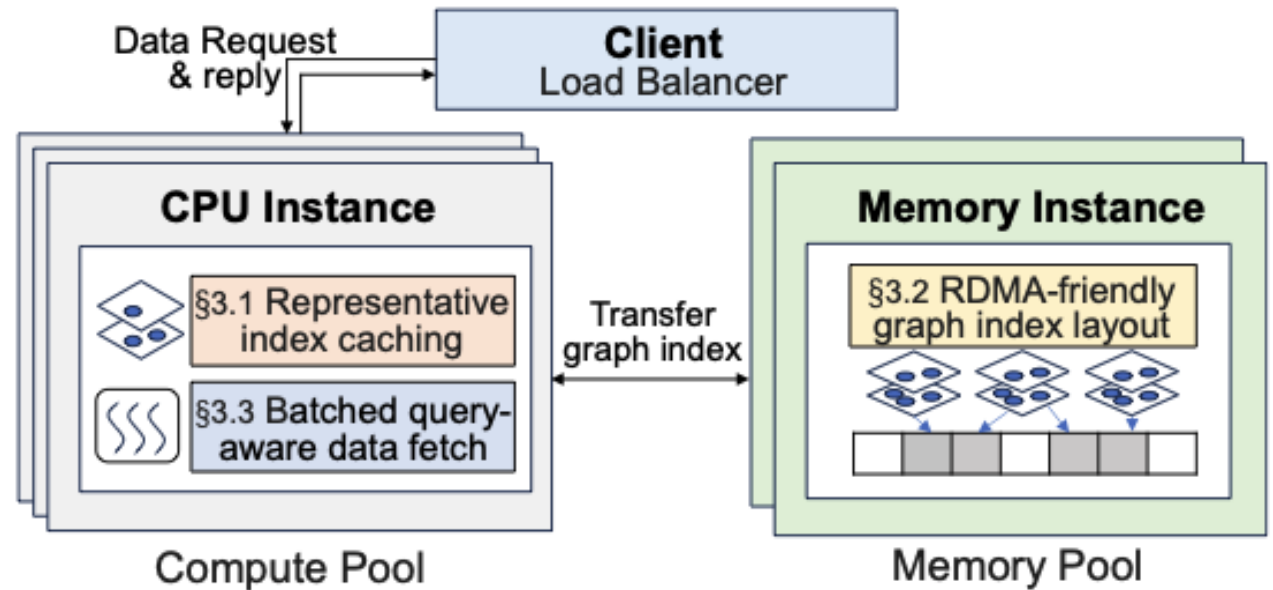


Overall design

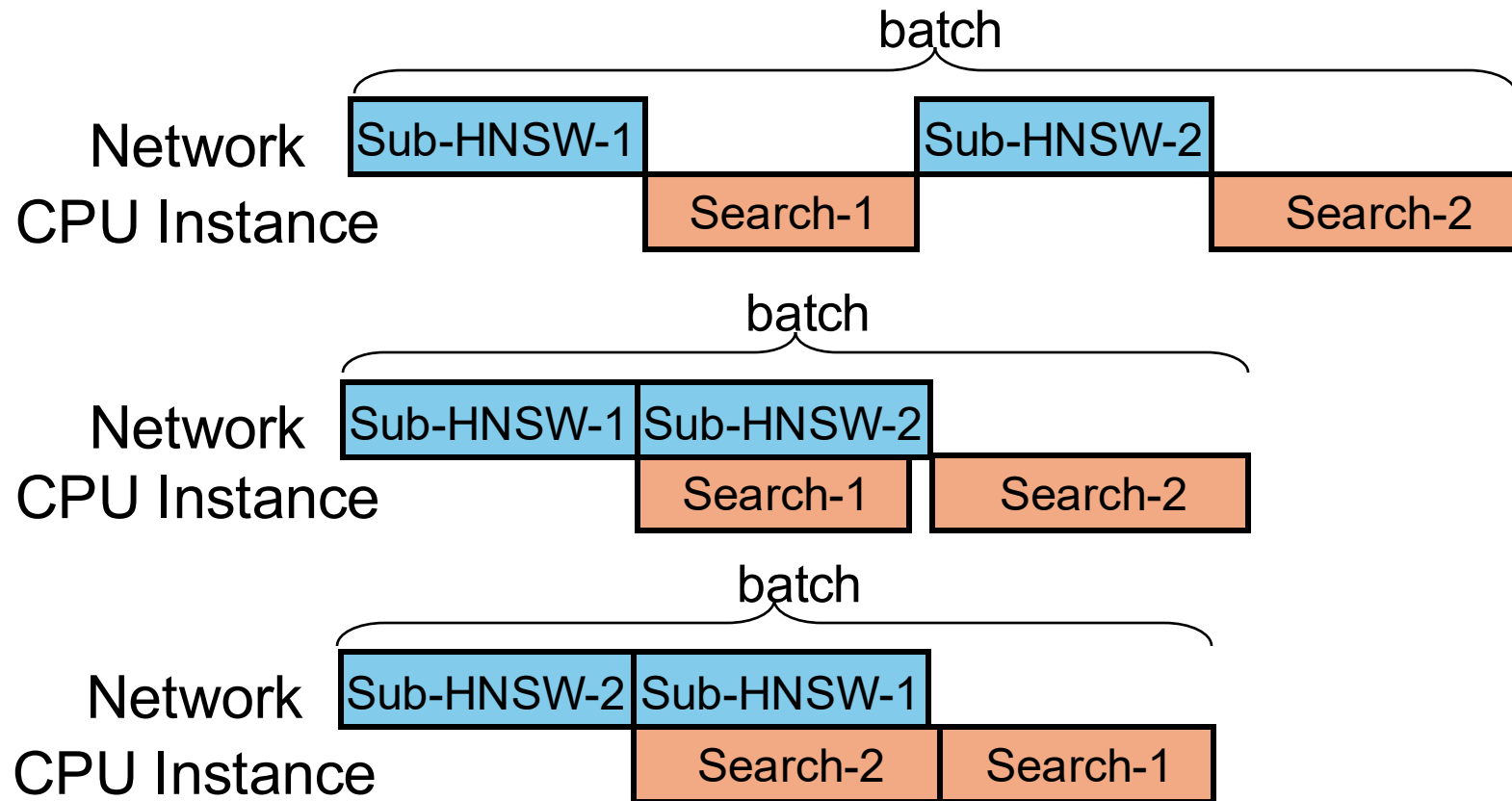
Representative index caching

RDMA-friendly graph index storage layout in remote memory.

Query-aware batched data loading



Pipeline Parallelism

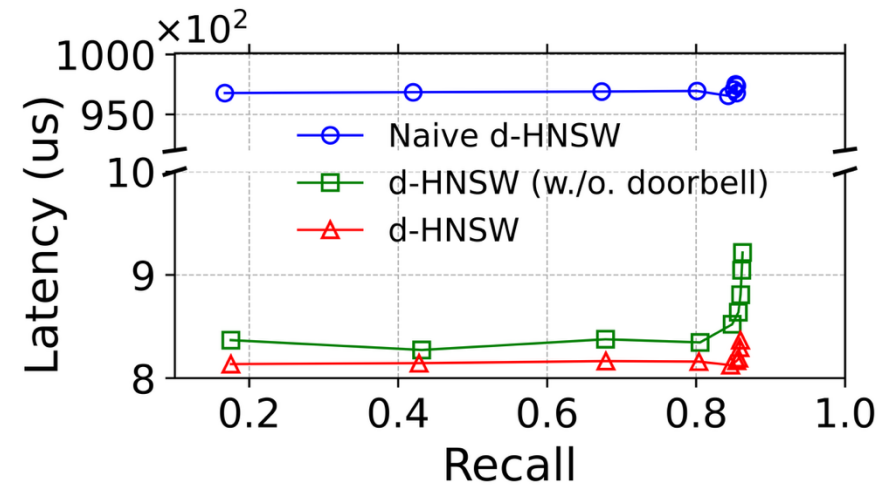
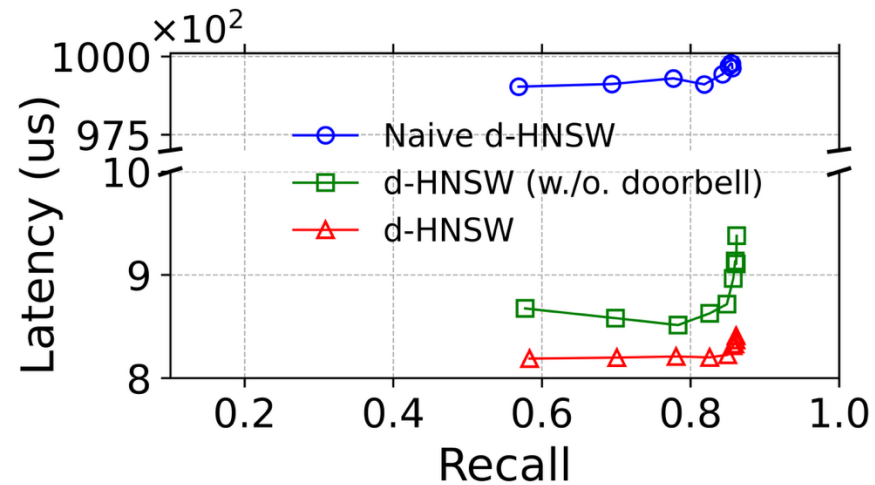


Implementation and Evaluation Step

- Implementation
 - ~12K LoC C++
- Testbed
 - Dell PowerEdge R650: 2×36-core Intel Xeon Platinum CPUs, 250GB RAM, 1.6TB NVMe SSD, **Mellanox ConnectX-6 100Gb NIC**
 - 3 as computing
 - 1 as memory node
- Datasets
 - SIFT1M
 - GIST1M

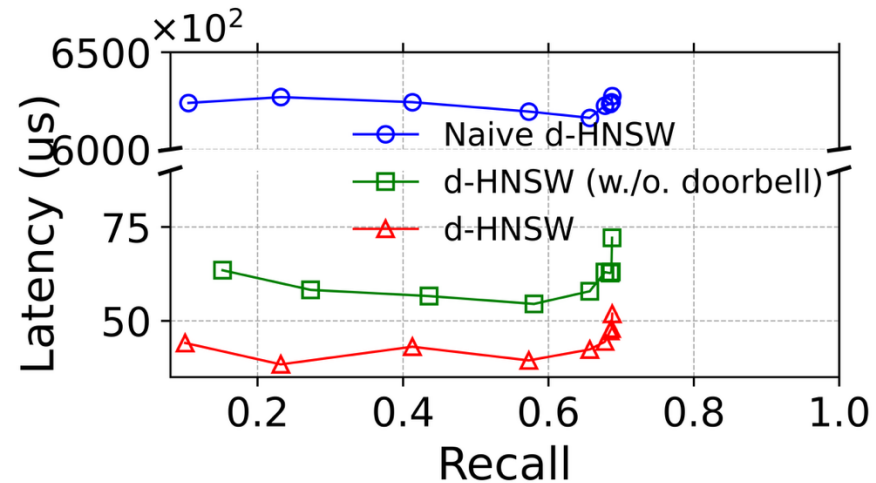
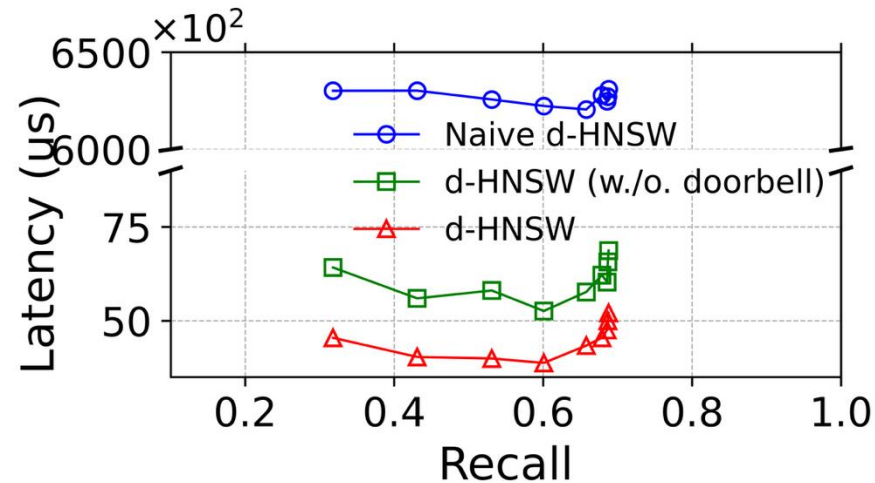
Evaluation

- d-HNSW reduces latency by up to $117\times$ and $1.12\times$ compared to naive d-HNSW and d-HNSW without doorbell



Evaluation

- d-HNSW achieves up to $121\times$ and $1.30\times$ lower latency compared to naive d-HNSW and d-HNSW without doorbell.



Evaluation

- d-HNSW outperforms Naïve d-HNSW and d-HNSW without doorbell in both network and sub-HNSW search latency.

Scheme	Network	Sub-HNSW	Meta-HNSW
Naive d-HNSW	90271.2 μ s	6564.5 μ s	13.52 μ s
d-HNSW (w./o. doorbell)	607.5 μ s	287.0 μ s	9.97 μ s
d-HNSW	<u>527.6μs</u>	269.2 μ s	9.75 μ s

Table 1: Latency breakdown for SIFT1M@1 with ef-Search as 48.

Scheme	Network	Sub-HNSW	Meta-HNSW
Naive d-HNSW	422.9ms	35.3ms	61.1 μ s
d-HNSW (w./o. doorbell)	2.9ms	1.27ms	52.6 μ s
d-HNSW	<u>1.3ms</u>	1.48ms	46.9 μ s

Table 2: Latency breakdown for GIST1M@1 with efSearch as 48.

Conclusion

- We present d-HNSW: the first RDMA-based vector similarity search engine for disaggregated memory system.
- d-HNSW enhances vector request throughput and minimizes data transfer overhead by implementing an RDMA-friendly data layout for memory nodes.
- d-HNSW optimizes batched vector queries by eliminating redundant vector transfers for batched vector queries.

Thank you!