

# Scalable Distributed Inverted List Indexes in Disaggregated Memory

MANUEL WIDMOSER

DANIEL KOCHER

NIKOLAUS AUGSTEN

University of Salzburg, Austria

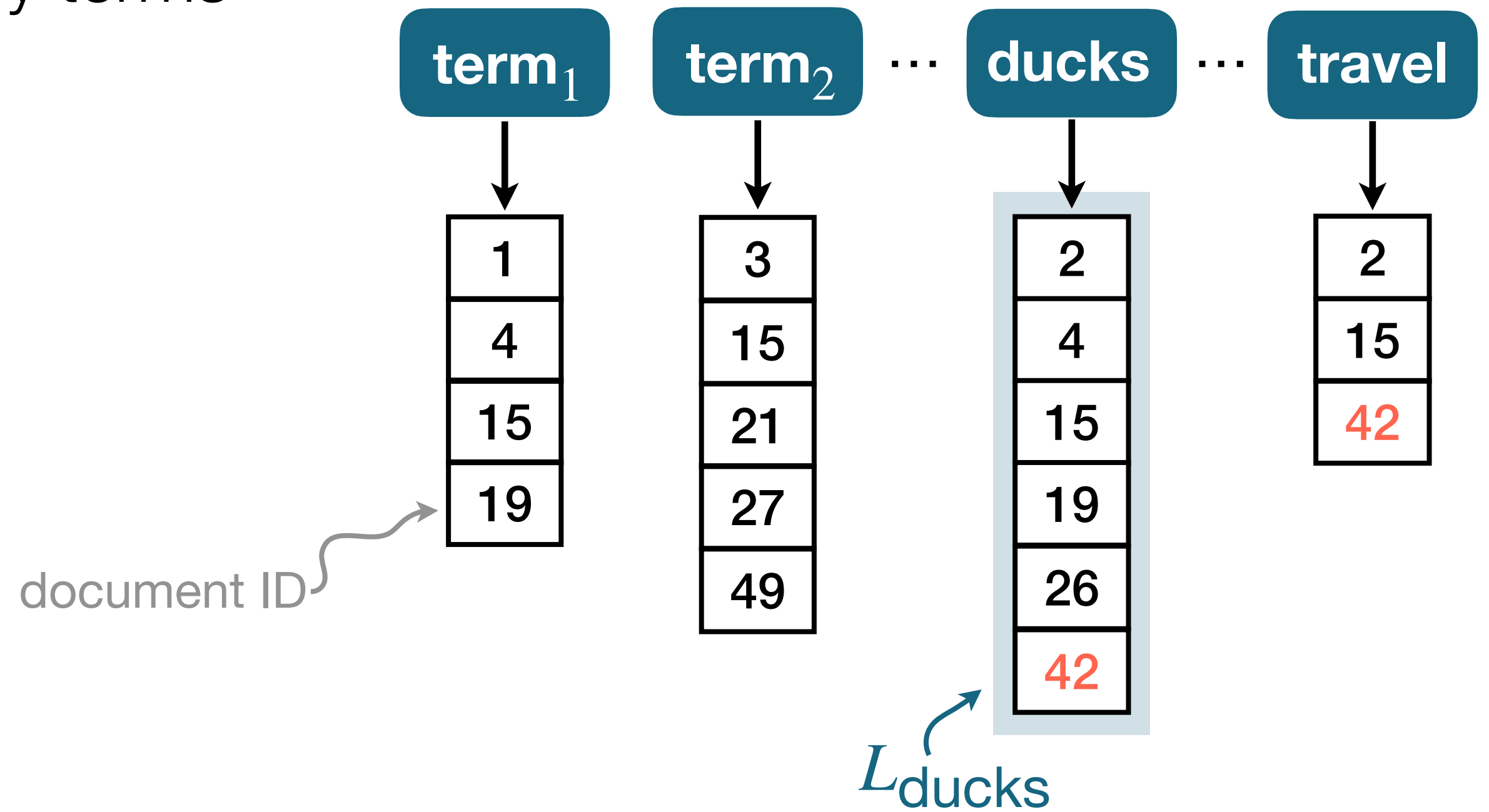
*June 13th, 2024*



# Inverted List Index

A **term** maps to an *ordered* list of documents that contain this term

**GOAL** Return all documents with common query terms



## APPLICATIONS

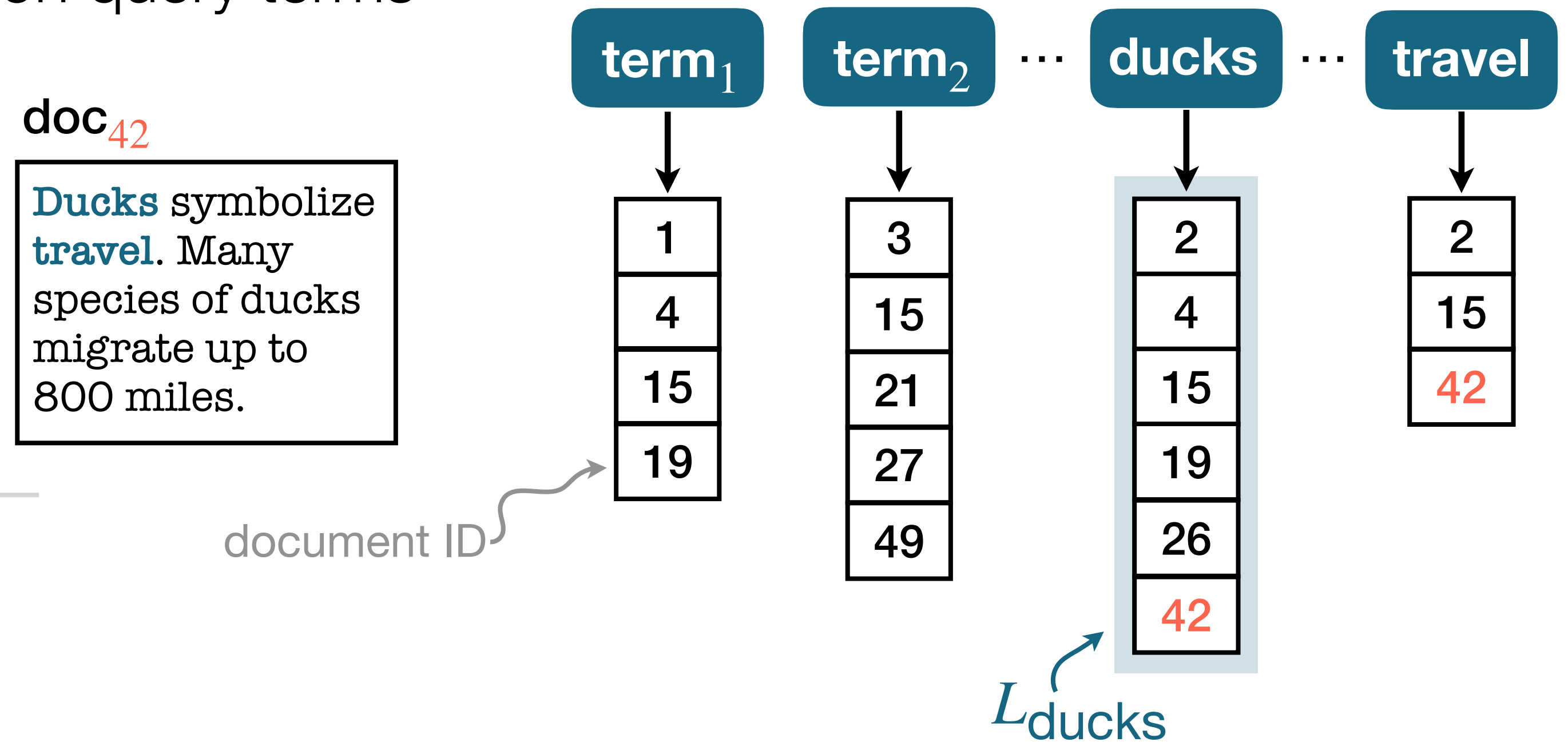
- Information retrieval
- DB query processing
- Graph analytics
- Similarity search

# Inverted List Index

A **term** maps to an *ordered* list of documents that contain this term

**GOAL** Return all documents with common query terms

- **INSERT:**
  - For each term in document **D**:  
Append D.id to list  $L_{\text{term}}$



## APPLICATIONS

- Information retrieval
- DB query processing
- Graph analytics
- Similarity search

# Inverted List Index

A **term** maps to an *ordered* list of documents that contain this term

**GOAL** Return all documents with common query terms

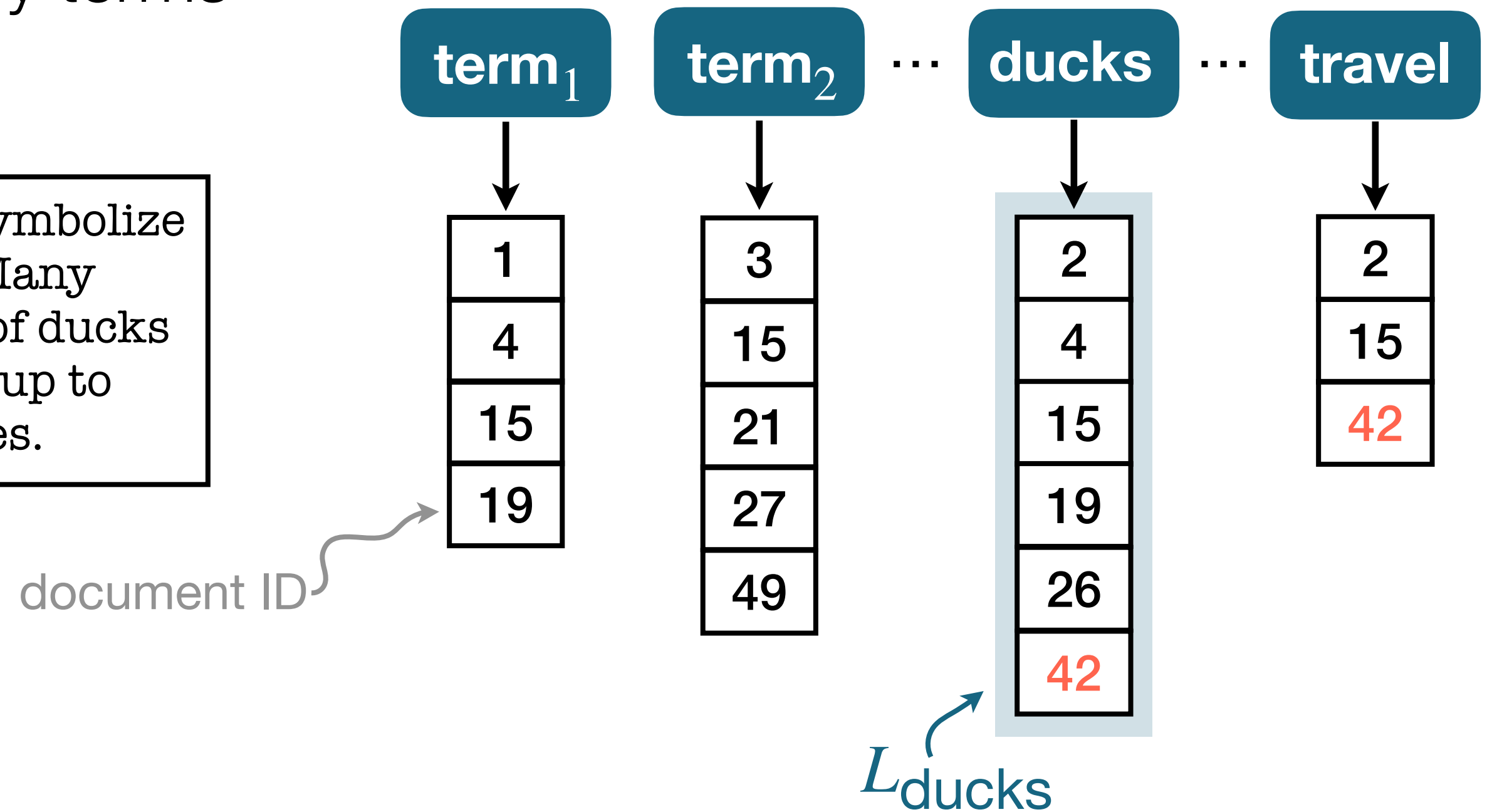
- **INSERT:**

- For each term in document **D**:  
Append D.id to list  $L_{\text{term}}$

doc<sub>42</sub>  
Ducks symbolize travel. Many species of ducks migrate up to 800 miles.

- **QUERY:**

- Query terms: **ducks** **travel**
- Result:  $L_{\text{ducks}} \cap L_{\text{travel}} = \{2, 15, 42\}$



## APPLICATIONS

- Information retrieval
- DB query processing
- Graph analytics
- Similarity search

# Why to Distribute Index Structures?

Two major reasons:

## MEMORY

---

- Very large datasets
- Spilling data to secondary (slow) storage is undesirable (e.g., index structure in an in-memory DBMS)

# Why to Distribute Index Structures?

Two major reasons:

## MEMORY

---

- Very large datasets
- Spilling data to secondary (slow) storage is undesirable (e.g., index structure in an in-memory DBMS)

## PERFORMANCE

---

- Scale out to massive amount of queries

# Why to Distribute Index Structures?

Two major reasons:

## MEMORY

---

- Very large datasets
- Spilling data to secondary (slow) storage is undesirable (e.g., index structure in an in-memory DBMS)

## PERFORMANCE

---

- Scale out to massive amount of queries

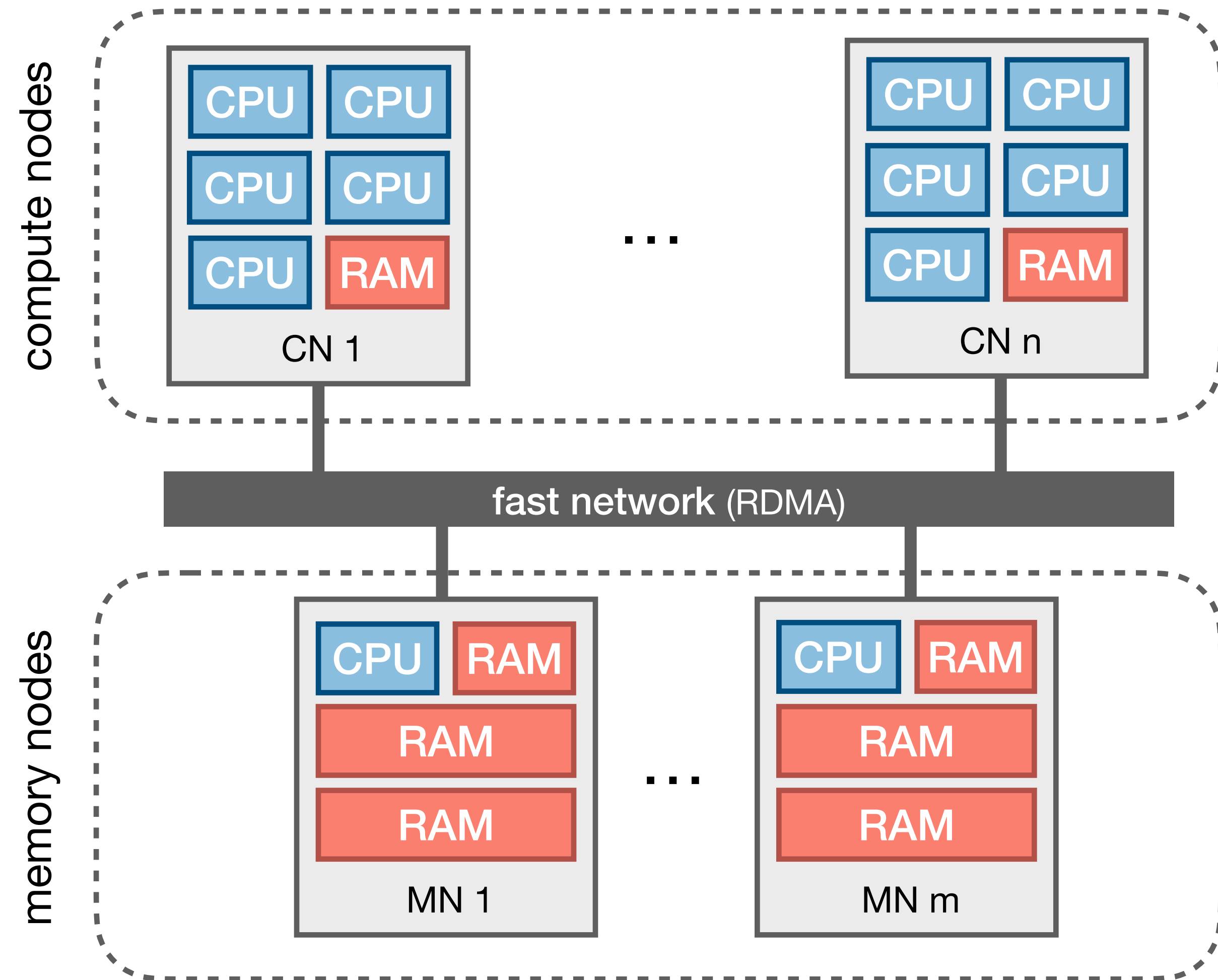
Architecture:

- Classic monolithic server architectures underutilize CPU and memory

CPU and RAM tightly coupled

# Disaggregated Memory (DM)

- High flexibility (elasticity)
- Cost-efficient resource utilization (sustainability)



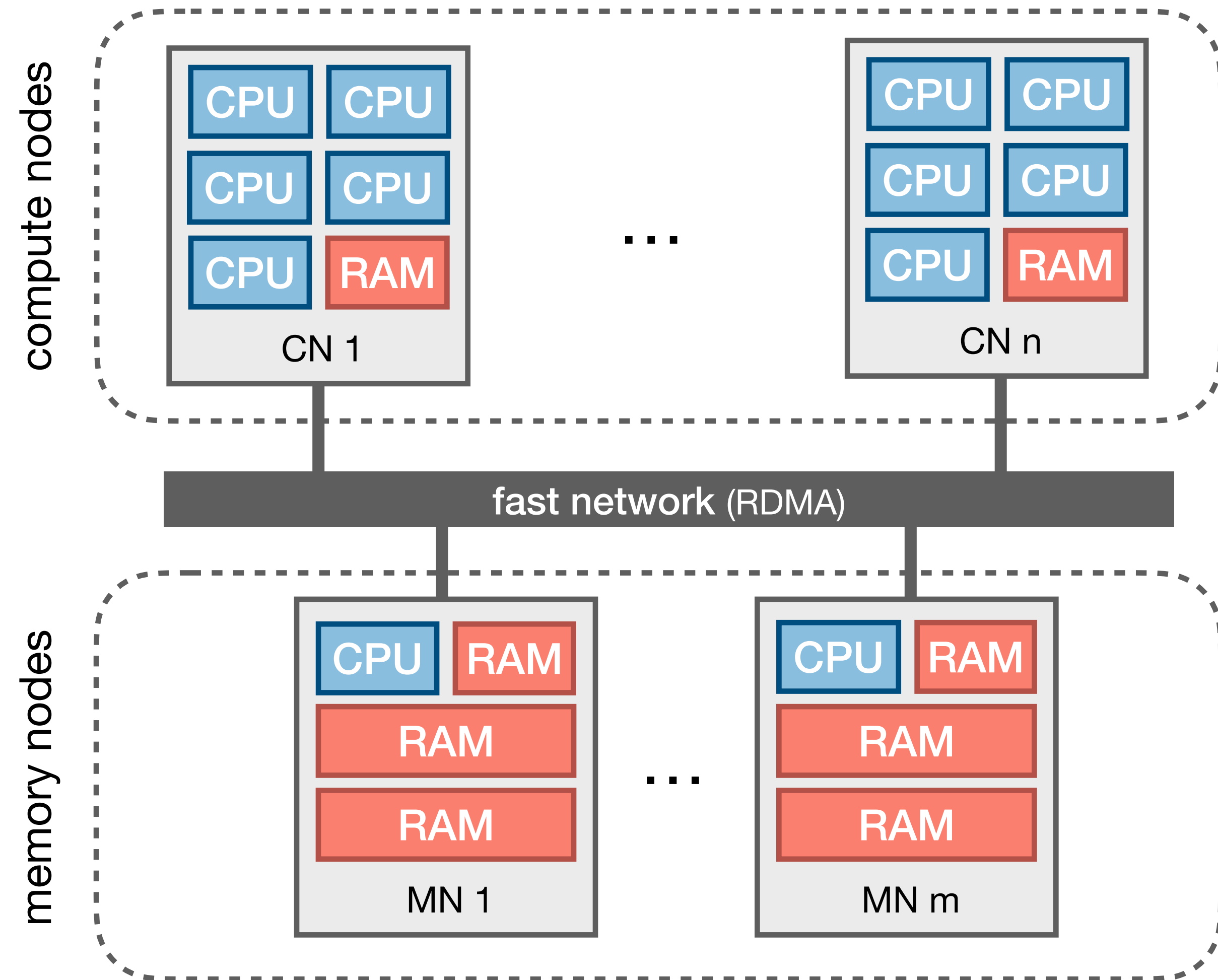


# Disaggregated Memory (DM)

- High flexibility (elasticity)
- Cost-efficient resource utilization (sustainability)

## DESIGN CHALLENGES

- Data must be accessed over the network
- MNs have near-zero computation power



# Disaggregated Memory (DM)

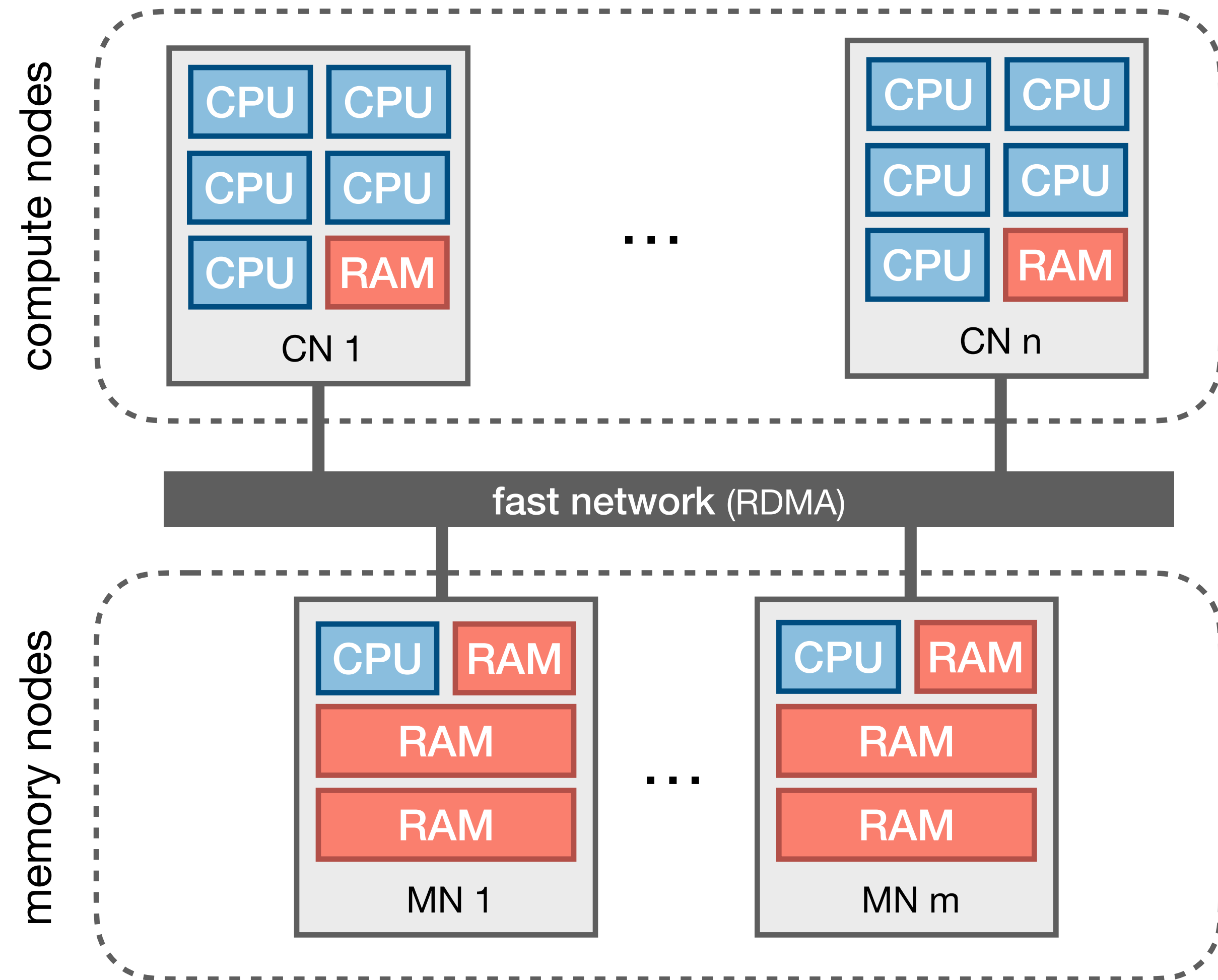
- High flexibility (elasticity)
- Cost-efficient resource utilization (sustainability)

## DESIGN CHALLENGES

- Data must be accessed over the network
- MNs have near-zero computation power

## RDMA

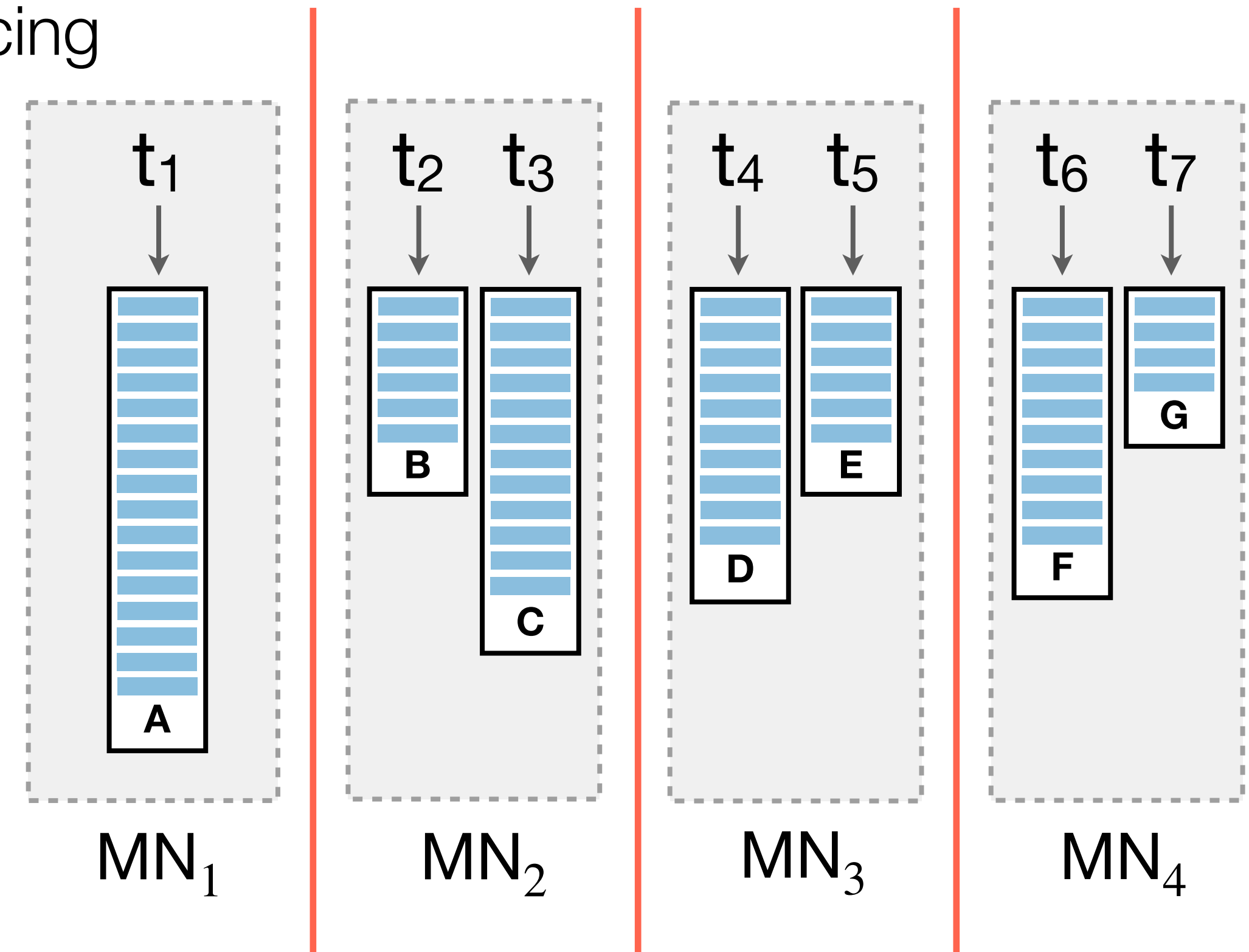
- Fully bypasses the remote CPU
- Low latency (single digit  $\mu s$ )



# TRADITIONAL DISTRIBUTION SCHEMES

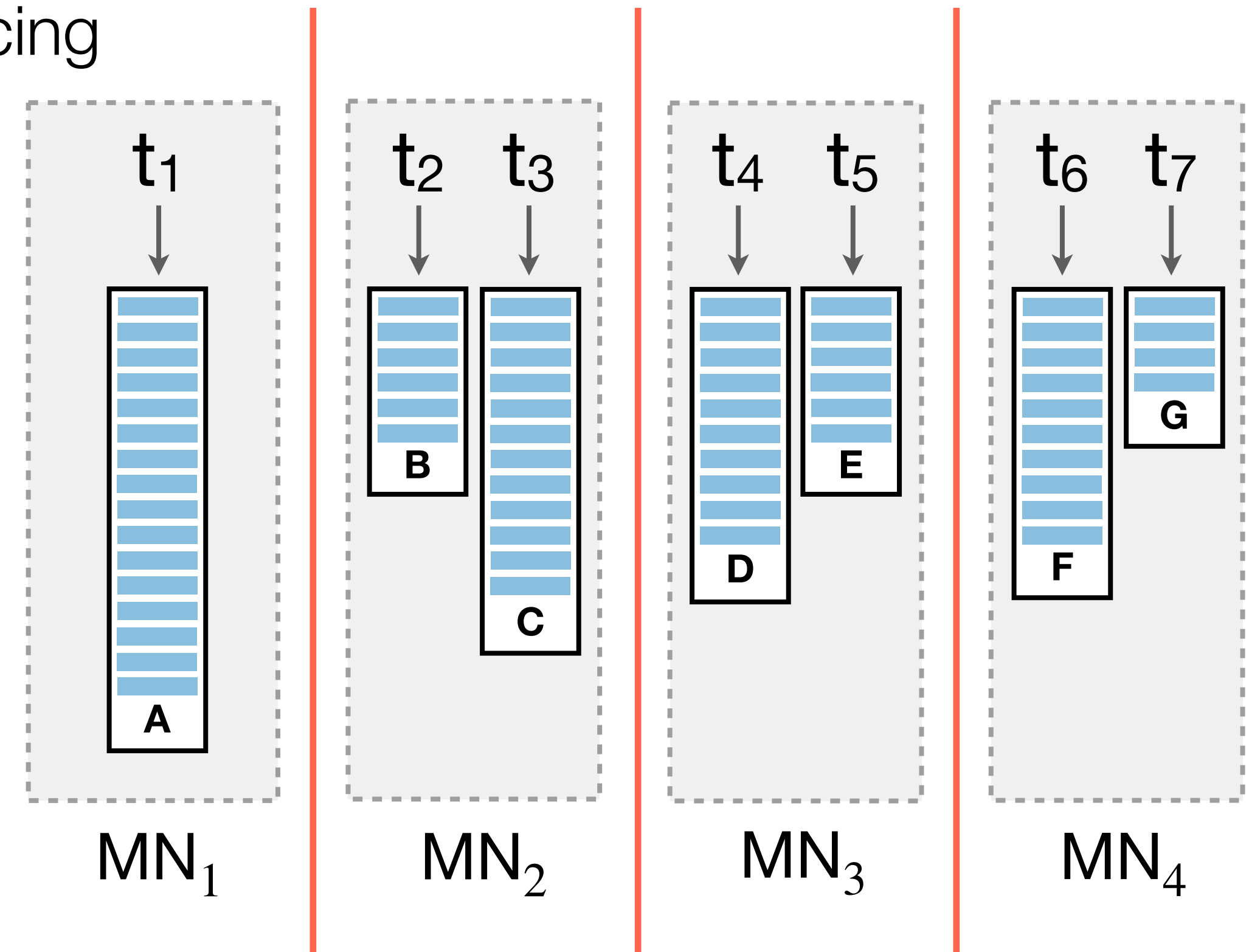
# Method 1: VERTICAL Partitioning

- Distribute complete lists across MNs using load balancing



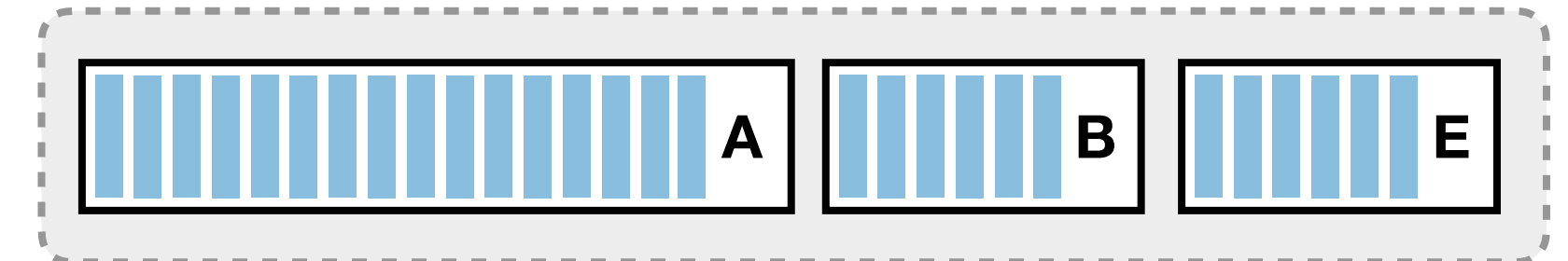
# Method 1: VERTICAL Partitioning

- Distribute complete lists across MNs using load balancing
- Query (on worker):
  - Read lists from remote memory to local buffer
  - Perform list operation



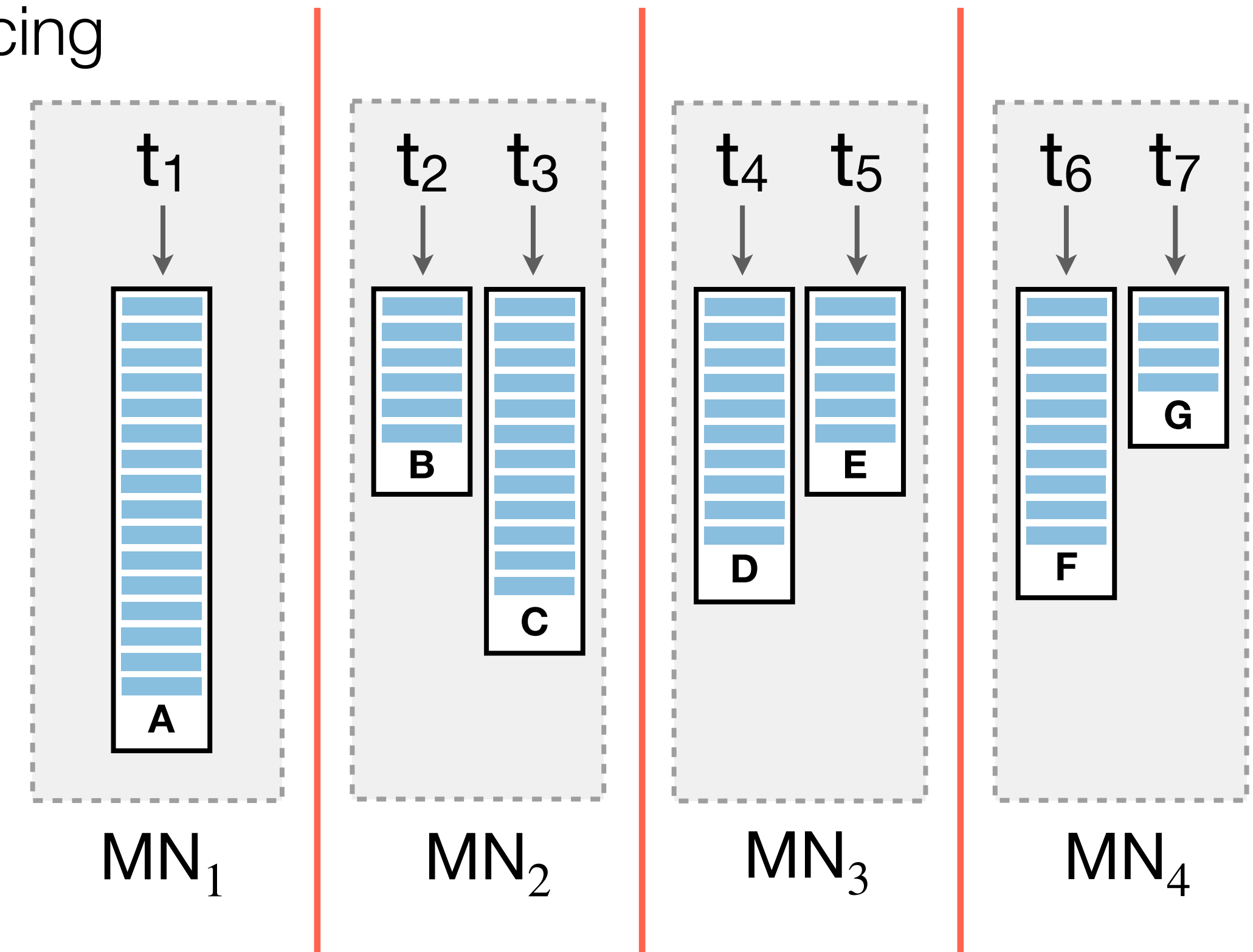
local buffer (on CN; per worker):

Query: {t<sub>1</sub>, t<sub>2</sub>, t<sub>5</sub>}



# Method 1: VERTICAL Partitioning

- Distribute complete lists across MNs using load balancing
- Query (on worker):
  - Read lists from remote memory to local buffer
  - Perform list operation

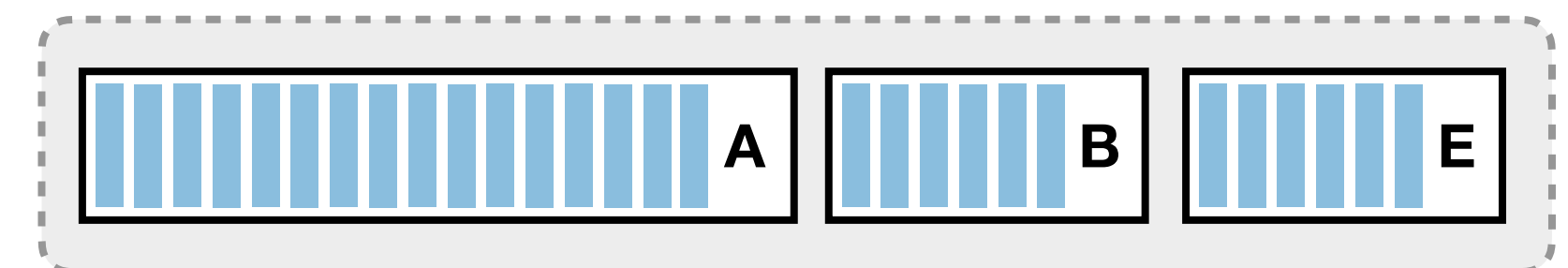


## SCALABILITY CHALLENGES

- C1** Network latency (waiting for long list)

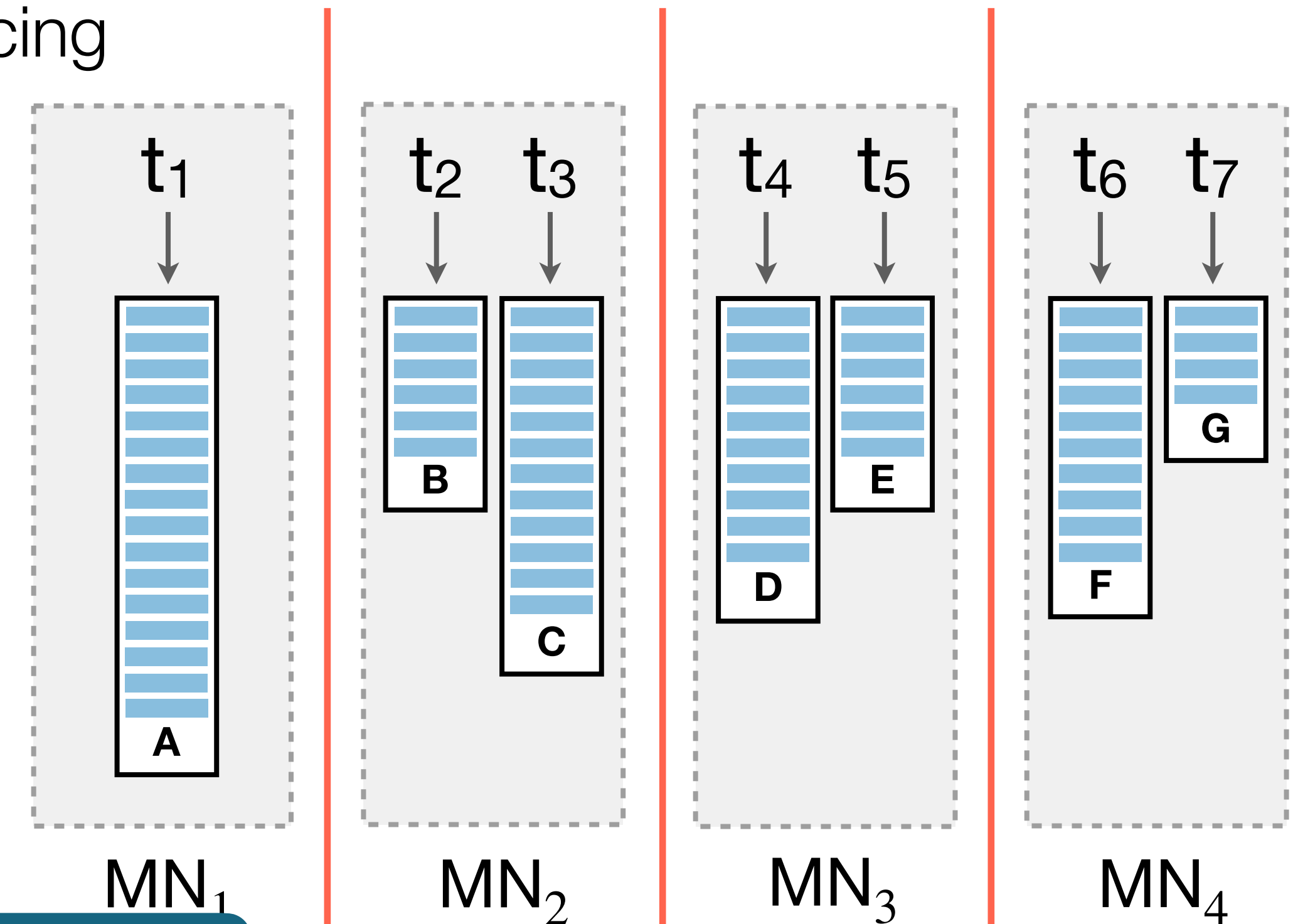
local buffer (on CN; per worker):

Query: {t<sub>1</sub>, t<sub>2</sub>, t<sub>5</sub>}



# Method 1: VERTICAL Partitioning

- Distribute complete lists across MNs using load balancing
- Query (on worker):
  - Read lists from remote memory to local buffer
  - Perform list operation



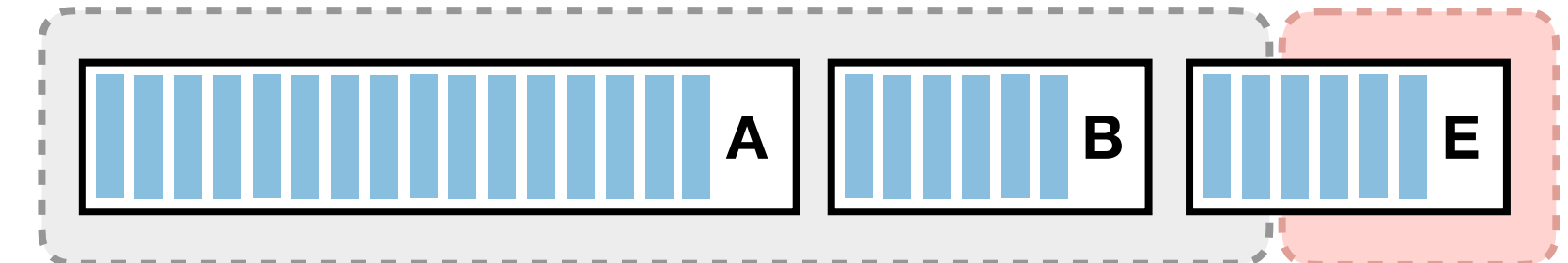
## SCALABILITY CHALLENGES

- C1** Network latency (waiting for long list)
- C2** Limited memory on CNs

**GOAL: fast list operation**  
(e.g.,  $k$ -way intersection)

local buffer (on CN; per worker):

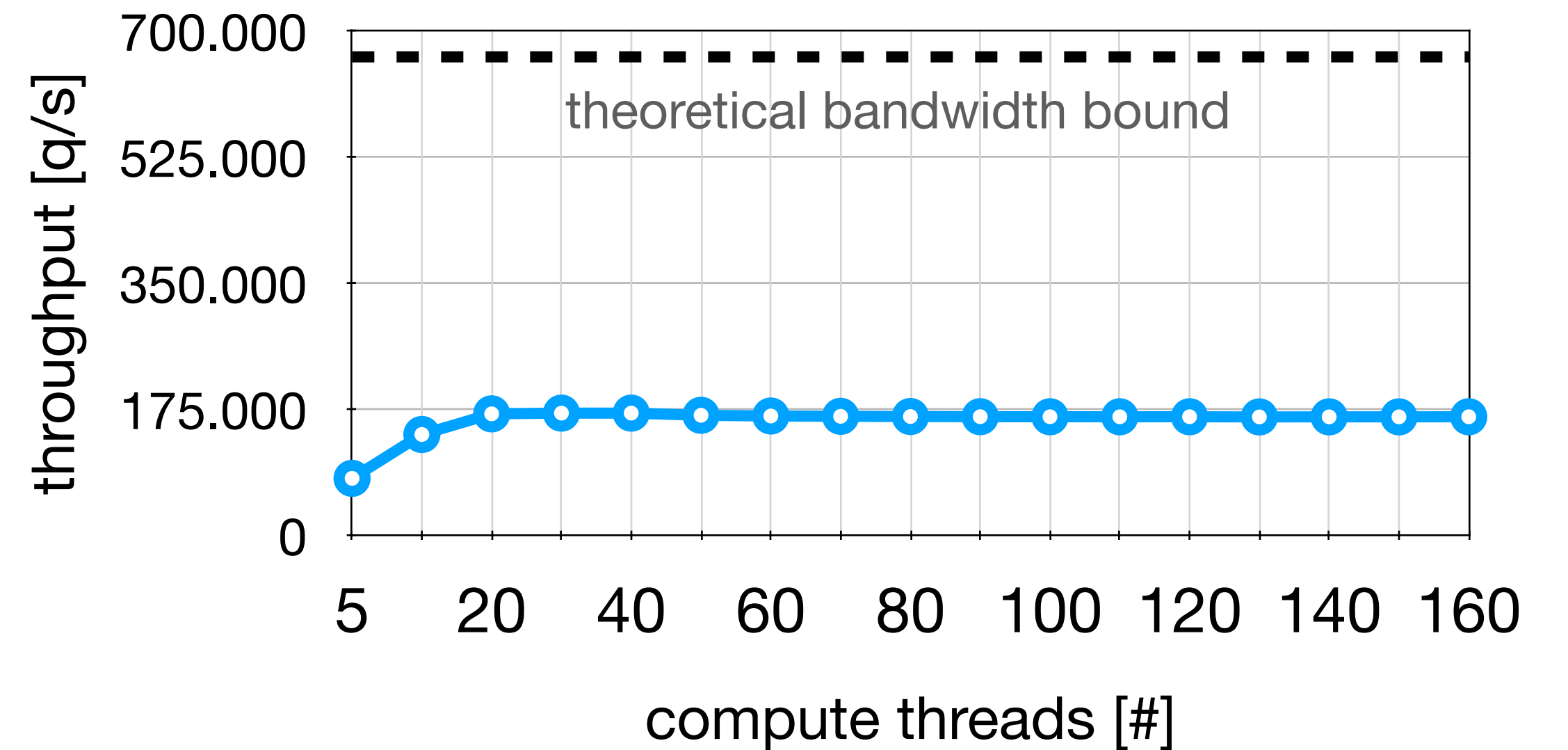
Query:  $\{t_1, t_2, t_5\}$



# Method 1: VERTICAL Partitioning

- Distribute complete lists across MNs using load balancing
- Query (on worker):
  - Read lists from remote memory to local buffer
  - Perform list operation

(higher is better)



## SCALABILITY CHALLENGES

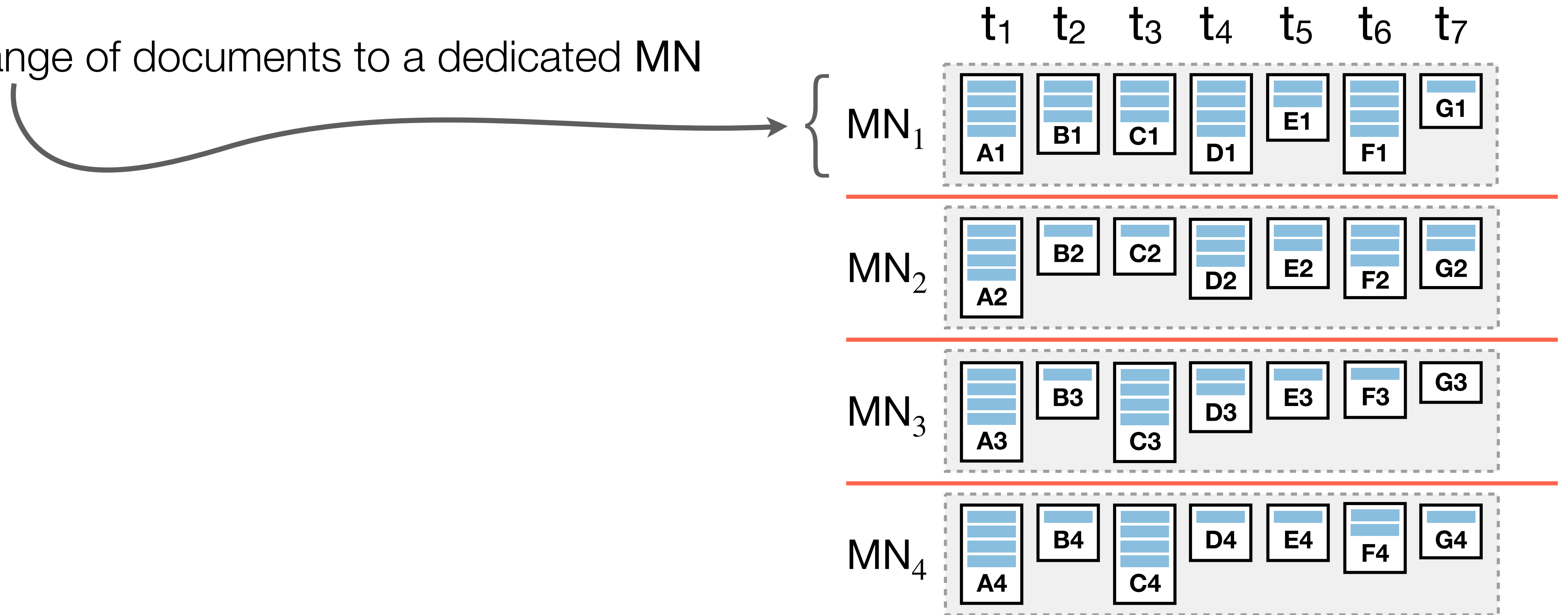
- C1** Network latency (waiting for long list)
- C2** Limited memory on CNs
- C3** Access skew

8KB lists (skewed accesses) / 5 CNs, 4 MNs



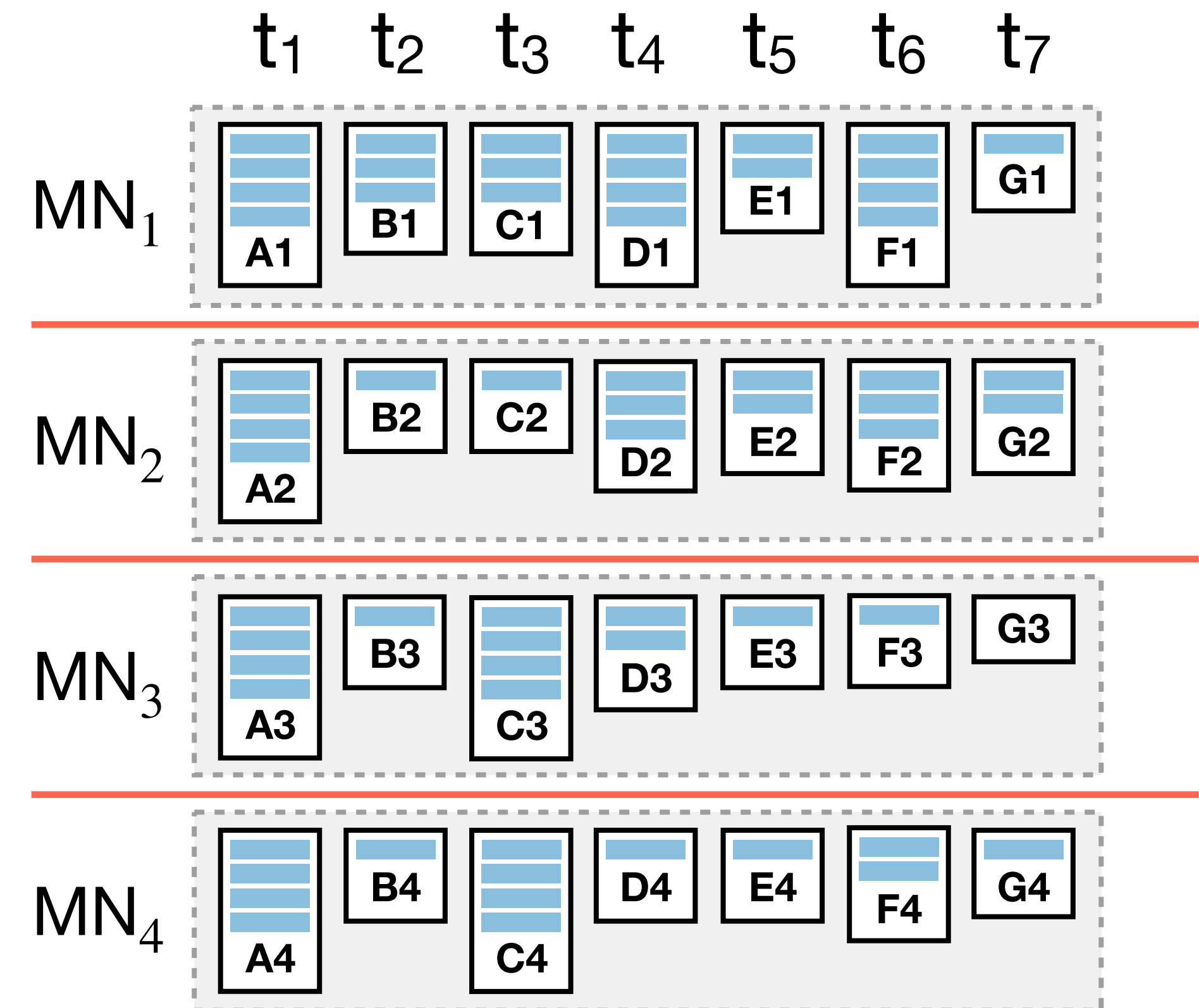
# Method 2: HORIZONTAL Partitioning

- Assign a range of documents to a dedicated MN



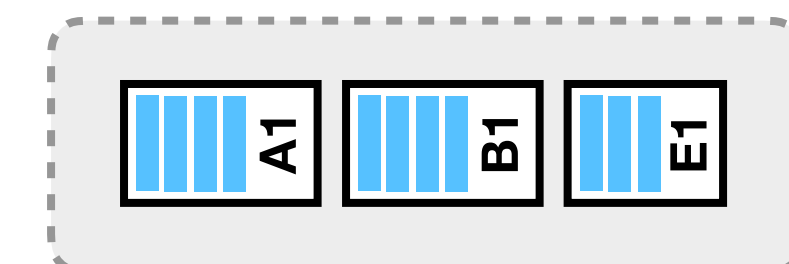
# Method 2: HORIZONTAL Partitioning

- Assign a range of documents to a dedicated MN
- Query (on worker):
  - For each MN:
    - Read partial lists to local buffer
    - Perform list operation



local buffer (on CN; per worker):

Query: {t<sub>1</sub>, t<sub>2</sub>, t<sub>5</sub>}

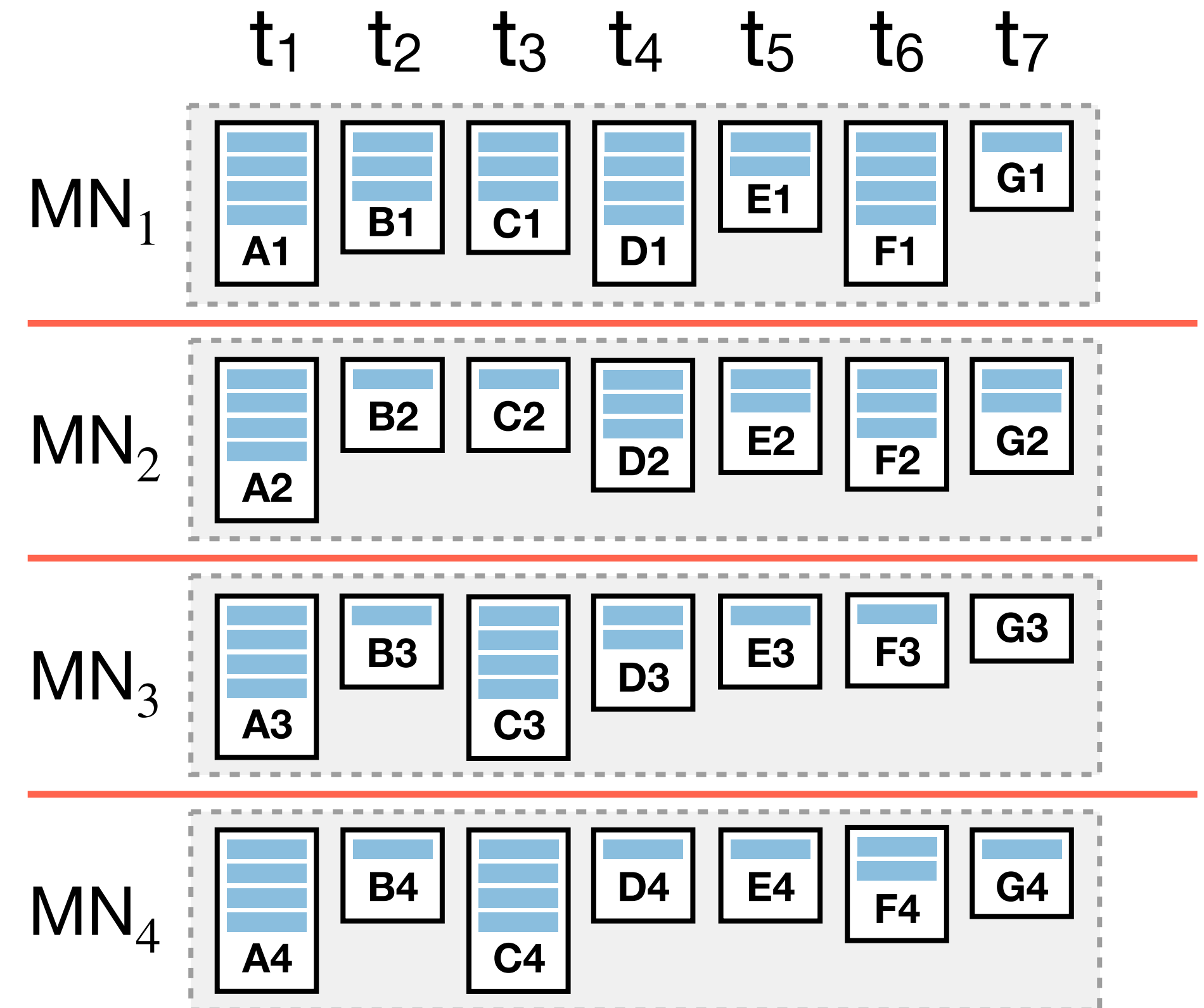


# Method 2: HORIZONTAL Partitioning

- Assign a range of documents to a dedicated MN
- Query (on worker):  
For each MN:
  - Read partial lists to local buffer
  - Perform list operation

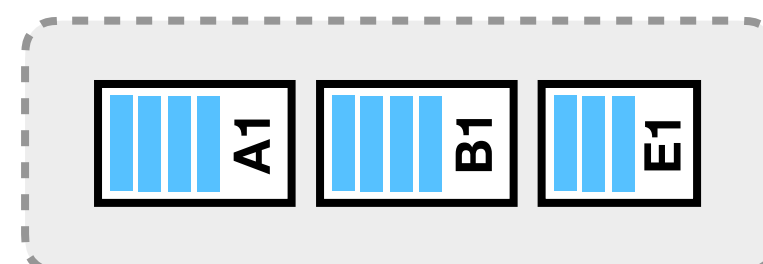
## SCALABILITY CHALLENGES

- C1** Network latency
- C2** Limited memory on CNs
- ~~**C3** Access skew~~



local buffer (on CN; per worker):

Query: {t<sub>1</sub>, t<sub>2</sub>, t<sub>5</sub>}



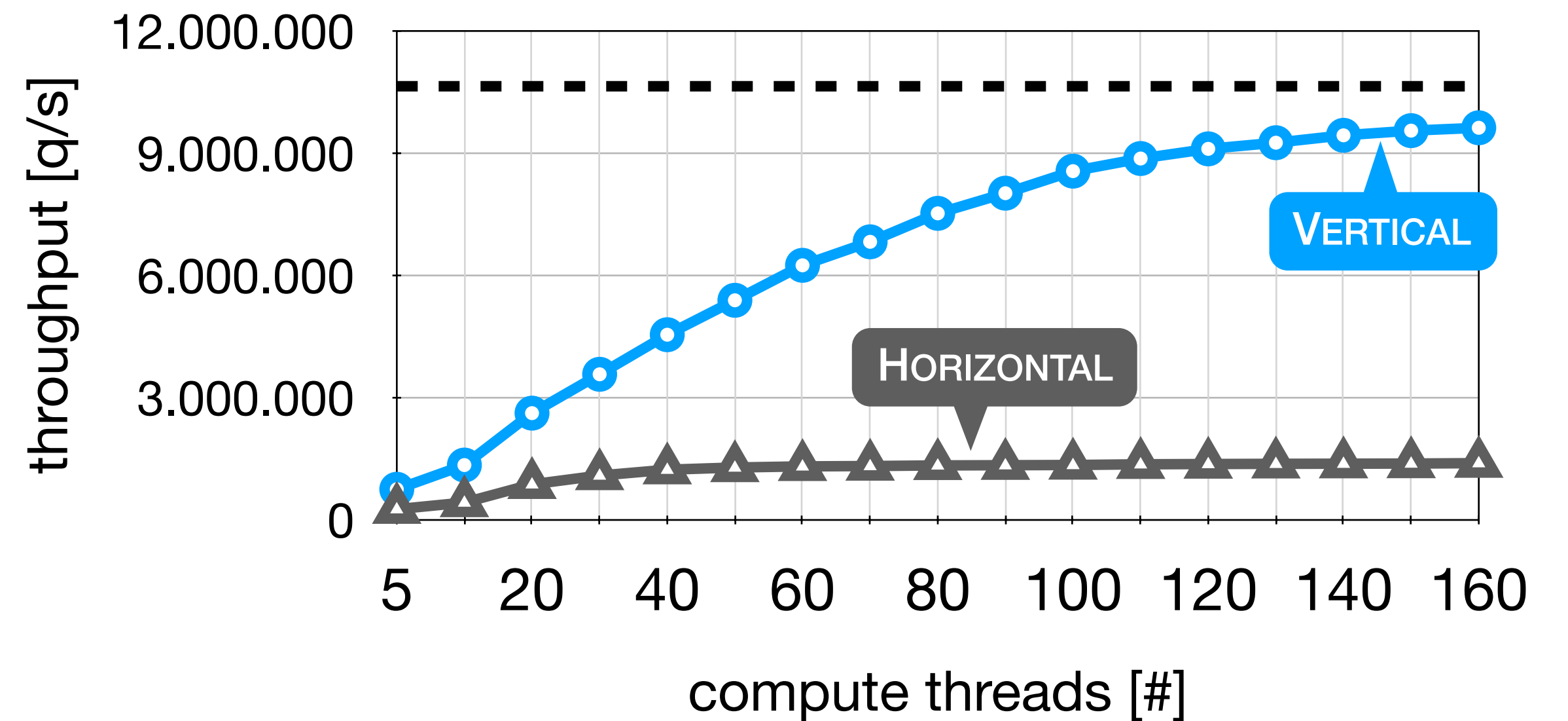
# Method 2: HORIZONTAL Partitioning

- Assign a range of documents to a dedicated MN
- Query (on worker):
  - For each MN:
    - Read partial lists to local buffer
    - Perform list operation

## SCALABILITY CHALLENGES

- C1** Network latency
- C2** Limited memory on CNs
- ~~**C3** Access skew~~
- C4** Network roundtrips

(higher is better)

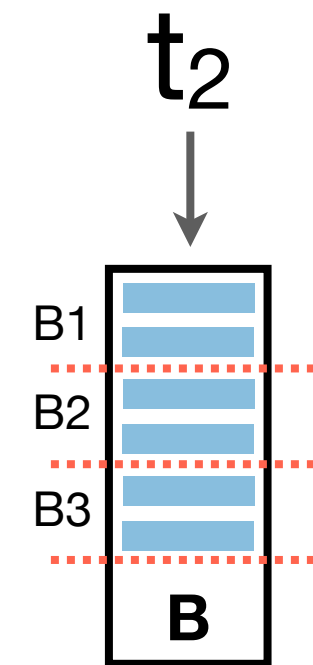


512B lists (uniform accesses) / 5 CNs, 4 MNs

# BLOCK-BASED SCHEME FOR INVERTED LISTS IN DM

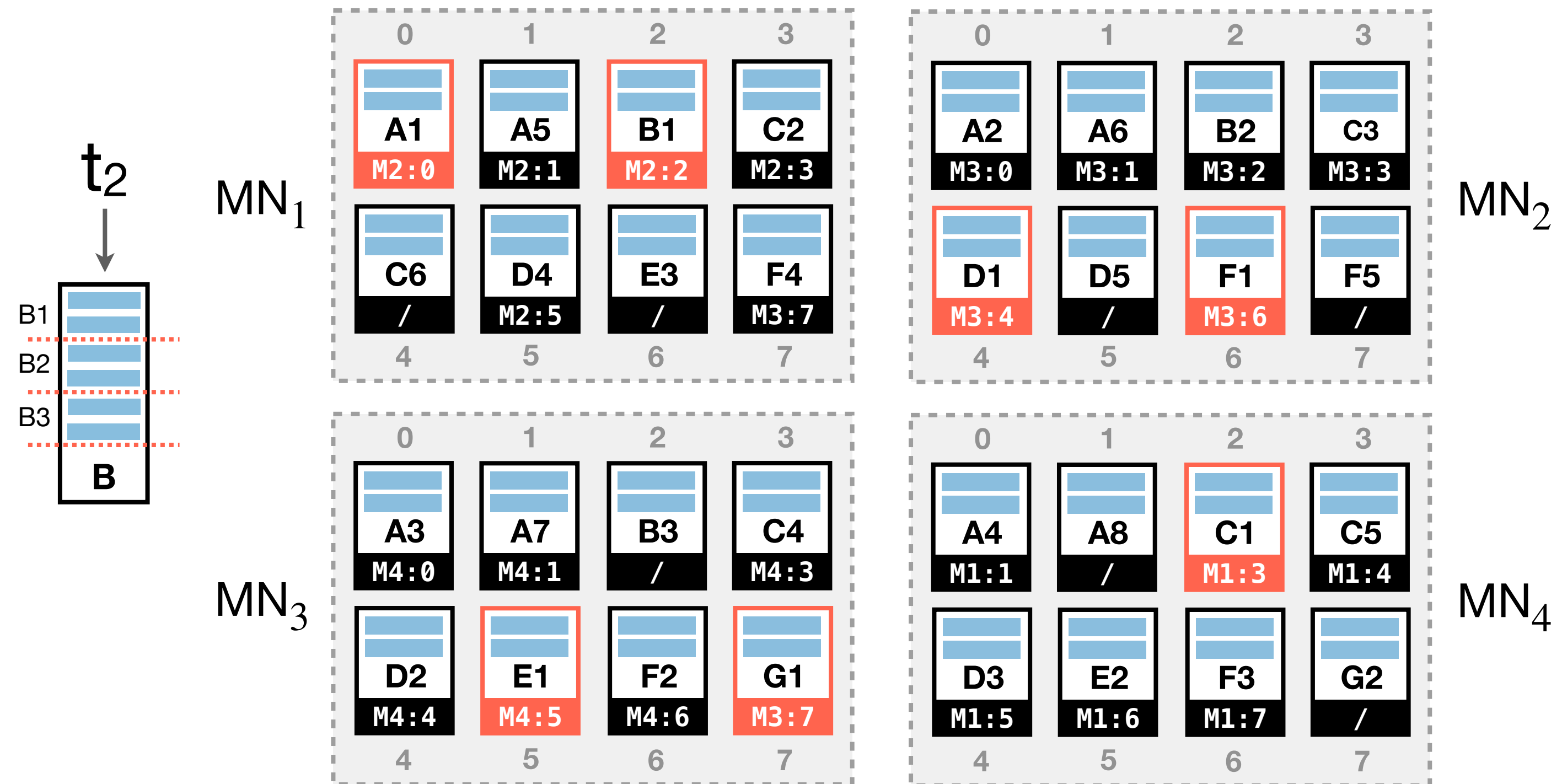
# Block-based Partitioning

- Divide lists into fix-sized blocks



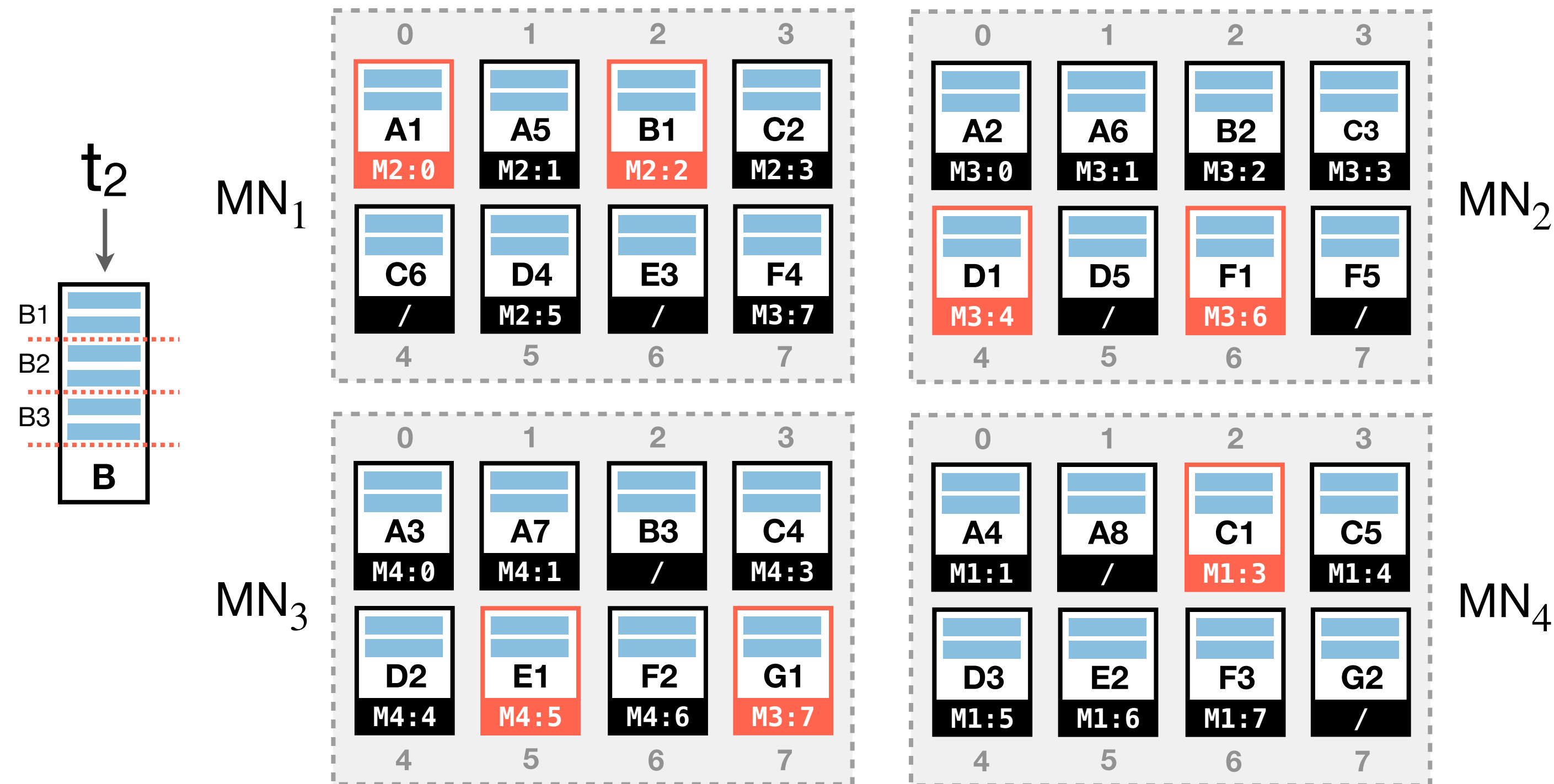
# Block-based Partitioning

- Divide lists into fix-sized blocks
- Distribute blocks among all MNs (e.g., round-robin, hash-based, ...)



# Block-based Partitioning

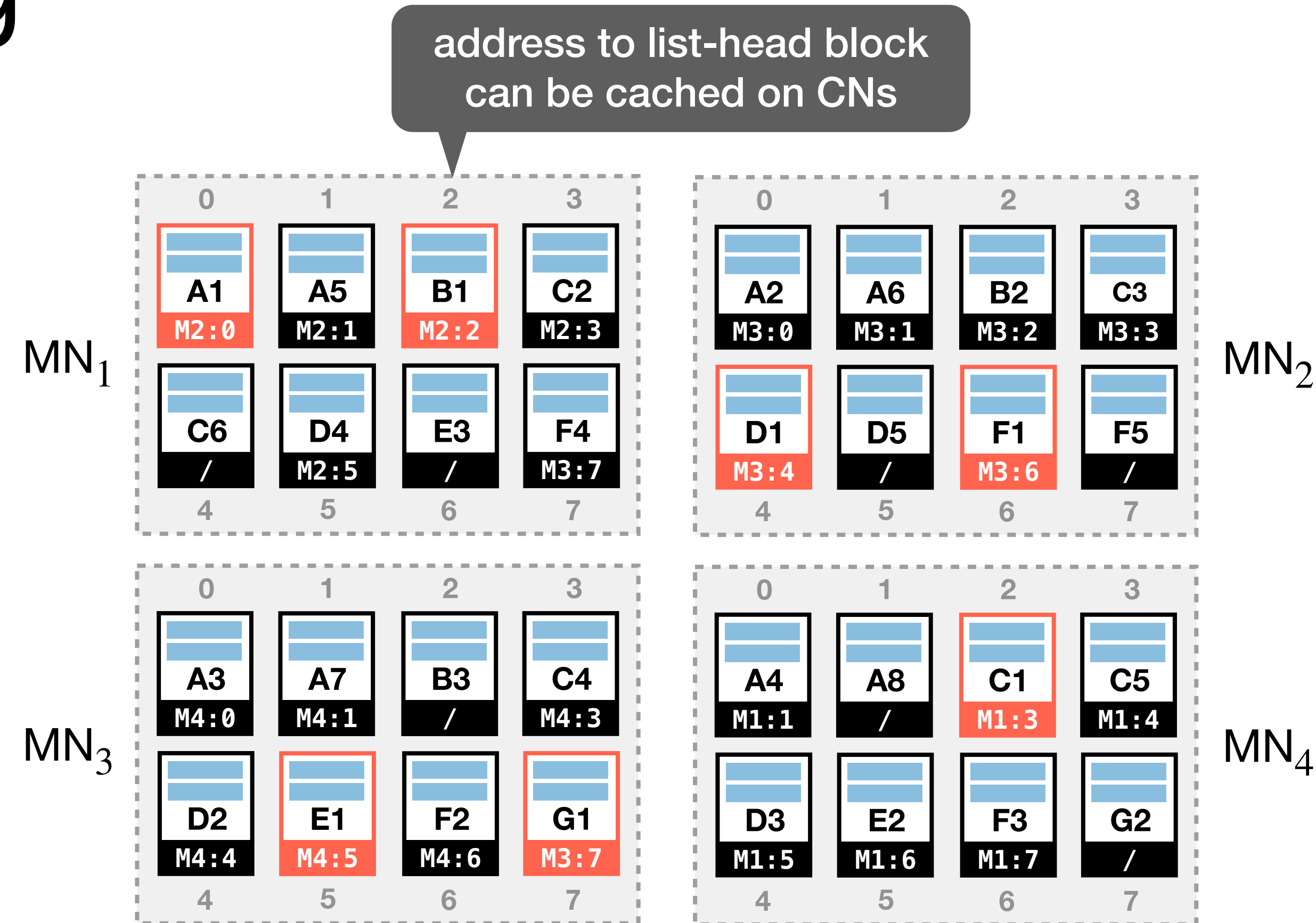
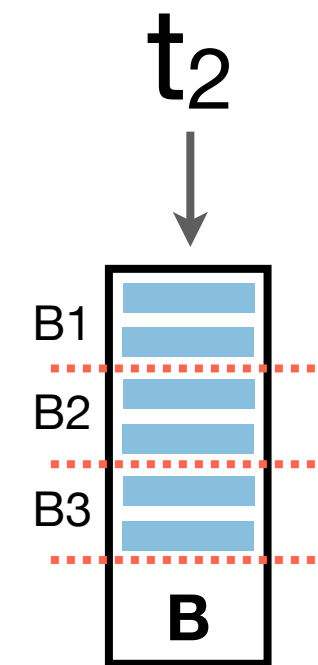
- Divide lists into fix-sized blocks
- Distribute blocks among all MNs (e.g., round-robin, hash-based, ...)
- Connect blocks via remote pointers





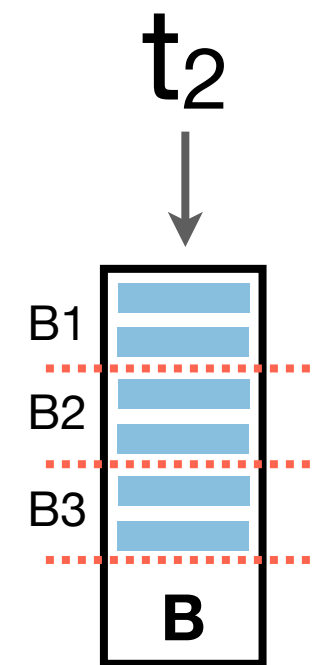
# Block-based Partitioning

- Divide lists into fix-sized blocks
- Distribute blocks among all MNs (e.g., round-robin, hash-based, ...)
- Connect blocks via remote pointers

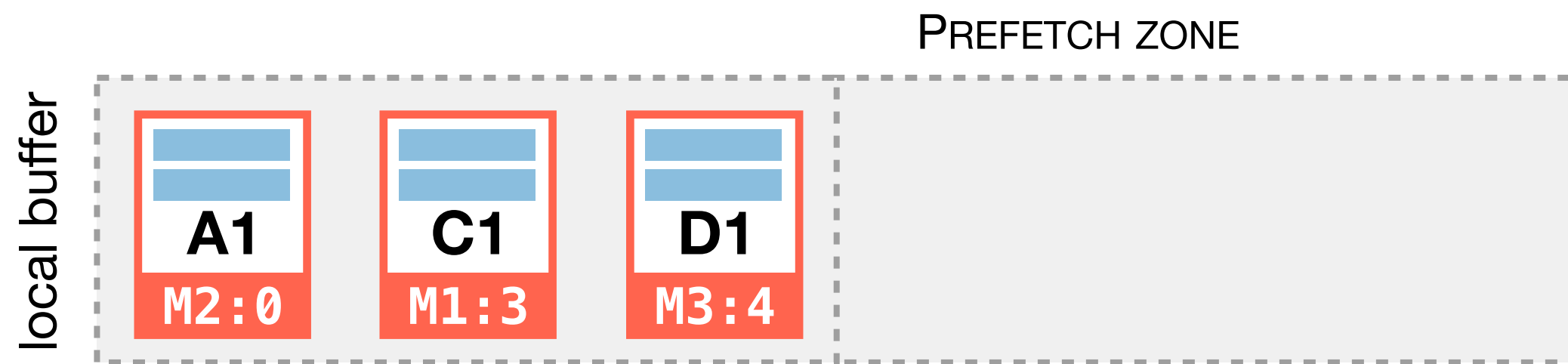
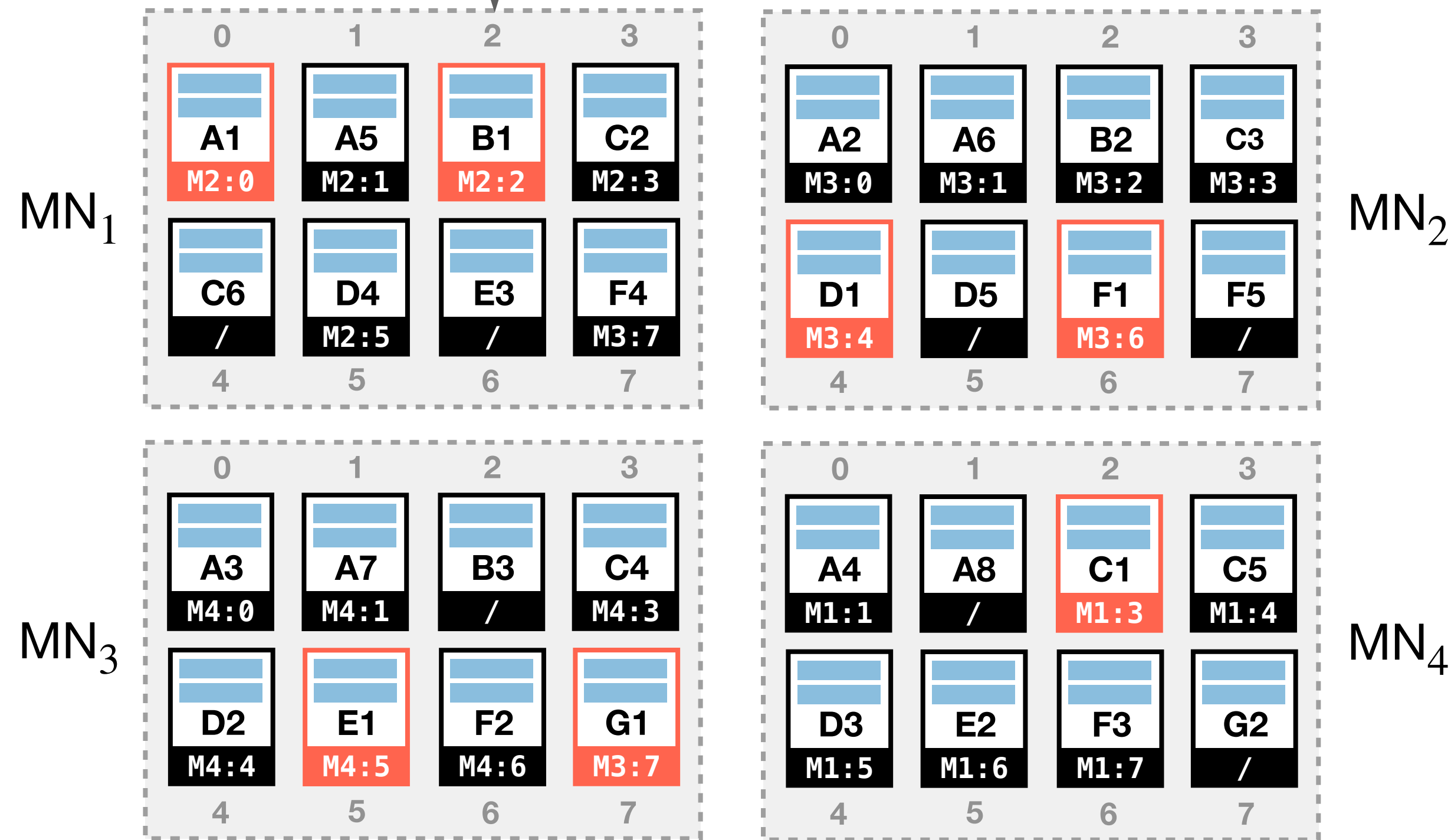


# Block-based Partitioning

- Divide lists into fix-sized blocks
- Distribute blocks among all MNs (e.g., round-robin, hash-based, ...)
- Connect blocks via remote pointers
- Query (on worker):
  - 1 Read list-head blocks to local buffer

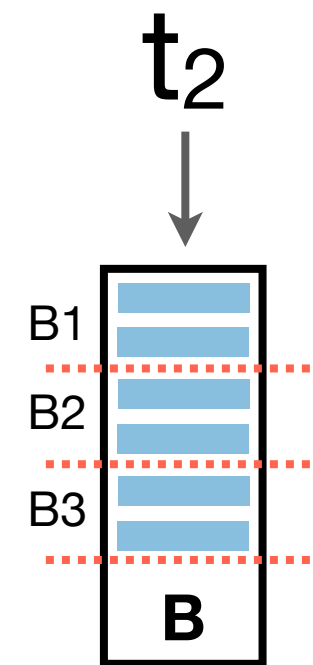


address to list-head block can be cached on CNs

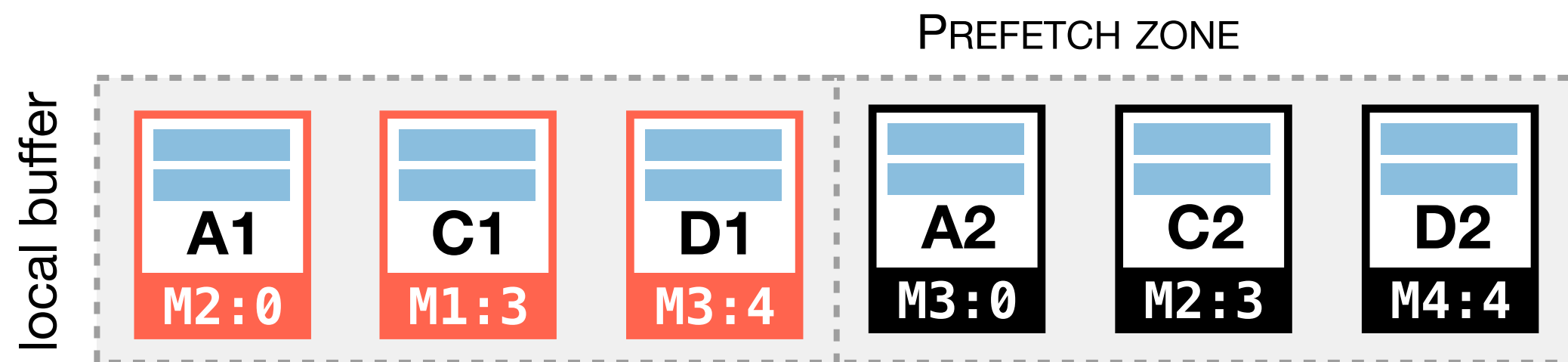
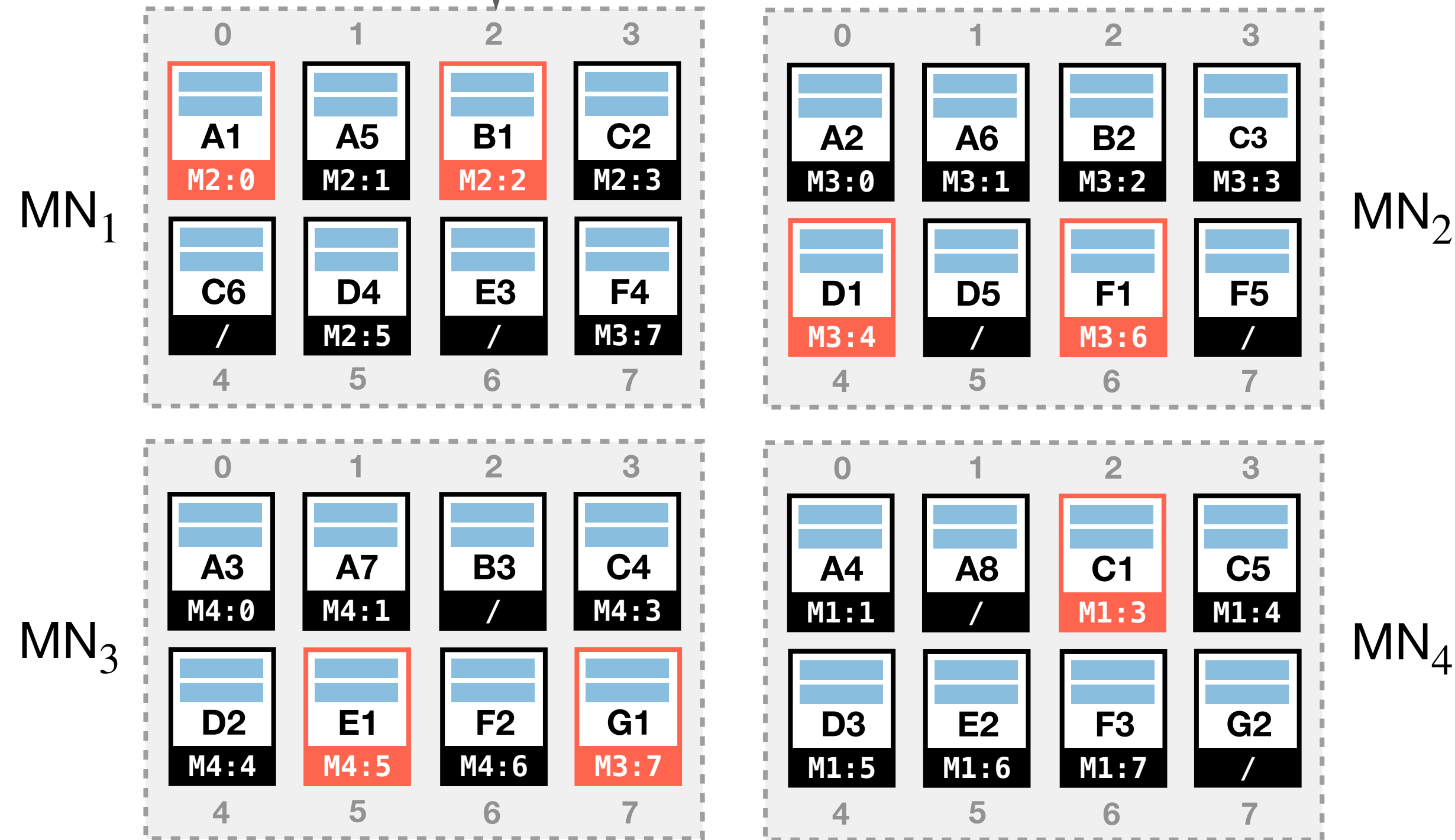


# Block-based Partitioning

- Divide lists into fix-sized blocks
- Distribute blocks among all MNs (e.g., round-robin, hash-based, ...)
- Connect blocks via remote pointers
- Query (on worker):
  - 1 Read list-head blocks to local buffer
  - 2 Prefetch subsequent blocks (issue RDMA\_READ)

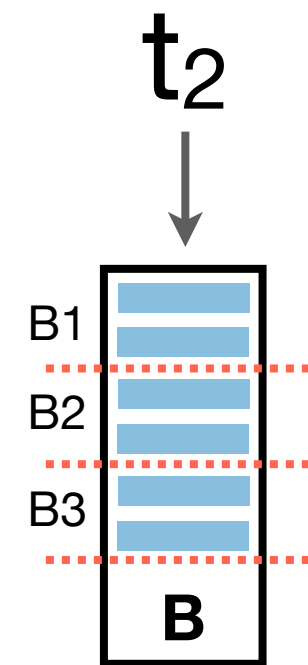


address to list-head block can be cached on CNs



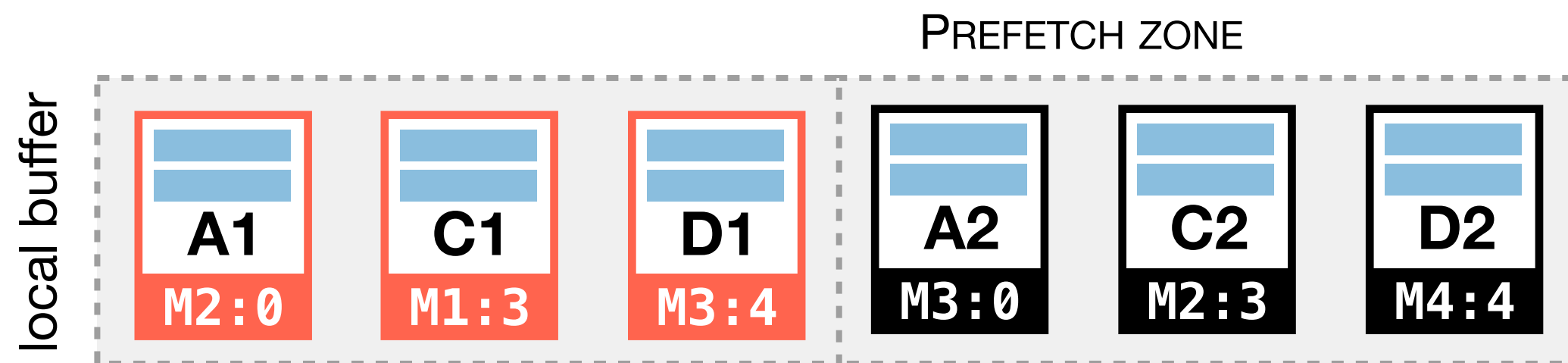
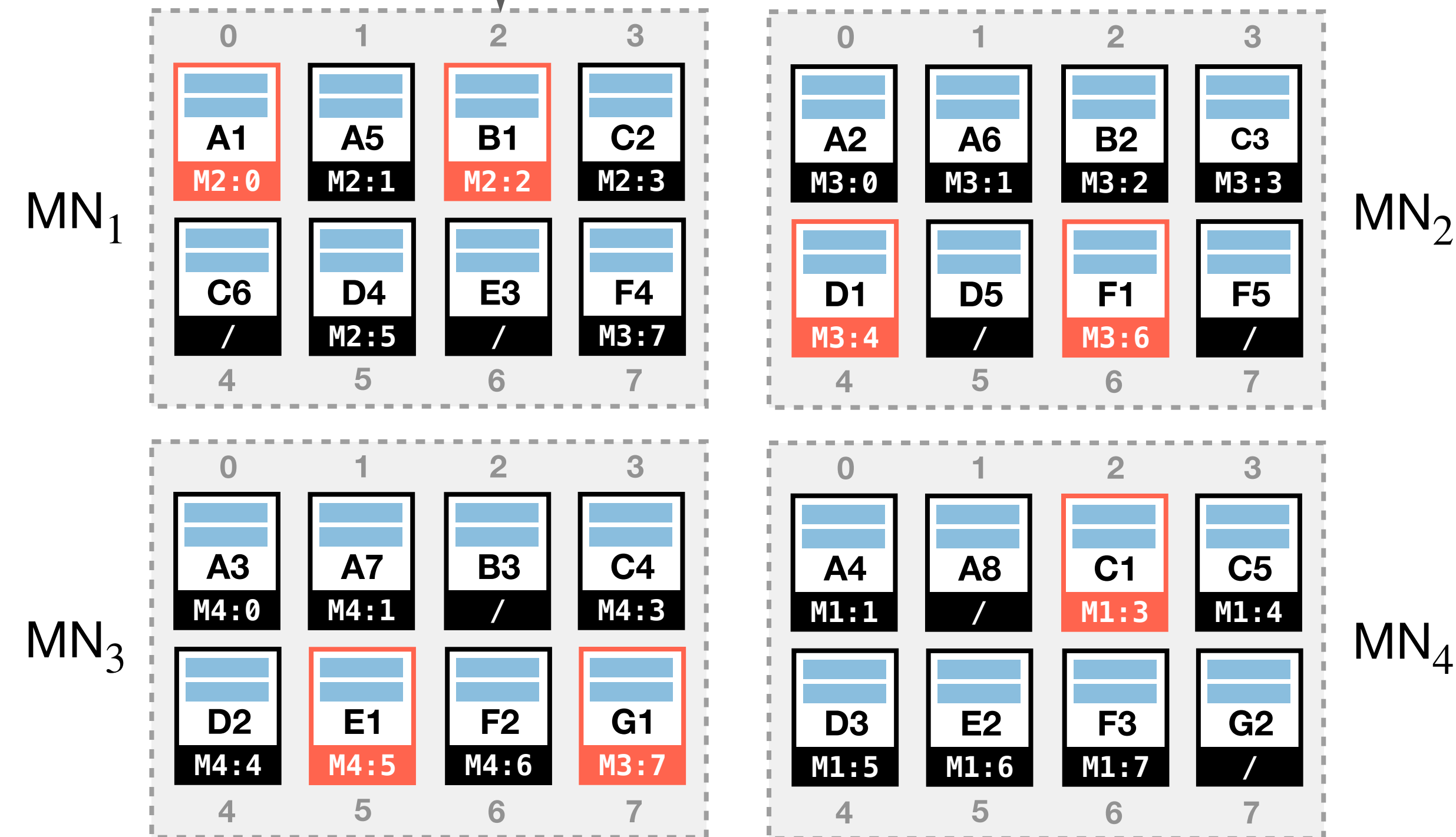
# Block-based Partitioning

- Divide lists into fix-sized blocks
- Distribute blocks among all MNs (e.g., round-robin, hash-based, ...)
- Connect blocks via remote pointers



- Query (on worker):
  - 1 Read list-head blocks to local buffer
  - 2 Prefetch subsequent blocks (issue RDMA\_READ)
  - 3 Perform the list operation — on block switch, go to 2

address to list-head block can be cached on CNs



# Addressing Scalability Challenges

## **C1** — NETWORK LATENCY

Network accesses and list operations are **fully interleaved**

# Addressing Scalability Challenges

**C1** — NETWORK LATENCY

Network accesses and list operations are **fully interleaved**

**C2** — LIMITED MEMORY ON CN

Local buffer size is **independent of the data** (list sizes)

# Addressing Scalability Challenges

**C1** — NETWORK LATENCY

Network accesses and list operations are **fully interleaved**

**C2** — LIMITED MEMORY ON CN

Local buffer size is **independent of the data** (list sizes)

**C3** — ACCESS SKEW

Skewed workloads do not overload individual MNs  
(due to the fine granularity of blocks)

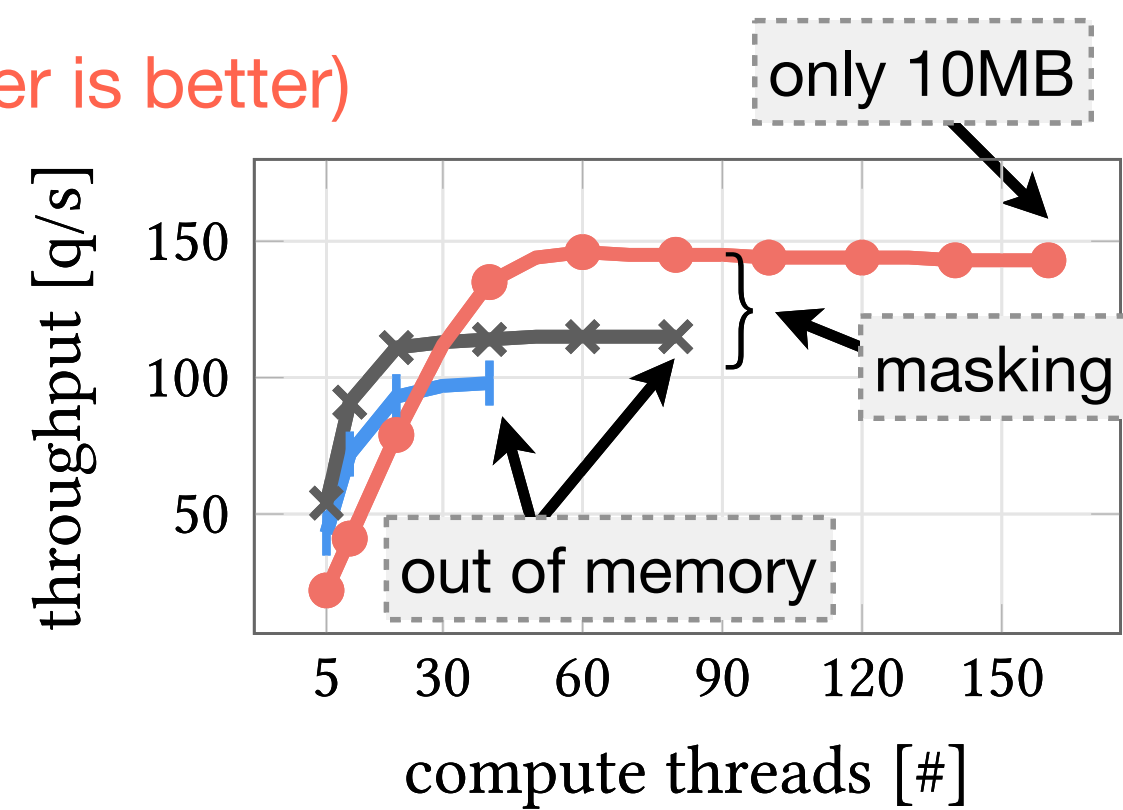
# Addressing Scalability Challenges

- C1** — NETWORK LATENCY  
Network accesses and list operations are **fully interleaved**
- C2** — LIMITED MEMORY ON CN  
Local buffer size is **independent of the data** (list sizes)
- C3** — ACCESS SKEW  
Skewed workloads do not overload individual MNs  
(due to the fine granularity of blocks)
- C4** — NETWORK ROUNDTRIPS  
Number of roundtrips is minimal for very small lists  
(otherwise, latency is masked with list operation)



# Scalability Experiments

(higher is better)



(a) CCNEWS

(information retrieval)

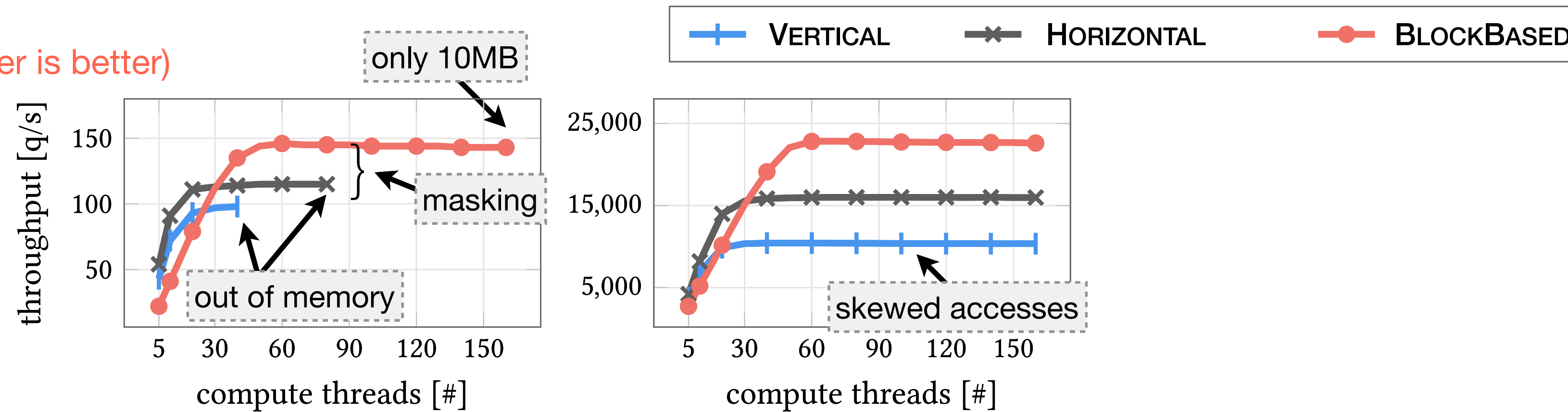
average query size: 6

## SETUP:

- 5 CNs (32 threads / 10GB RAM)
- 4 MNs (1 thread / 96GB RAM)
- Block size: 1KB

# Scalability Experiments

(higher is better)



**(a) CCNEWS**  
(information retrieval)

**(b) TWITTER**  
(graph analytics)

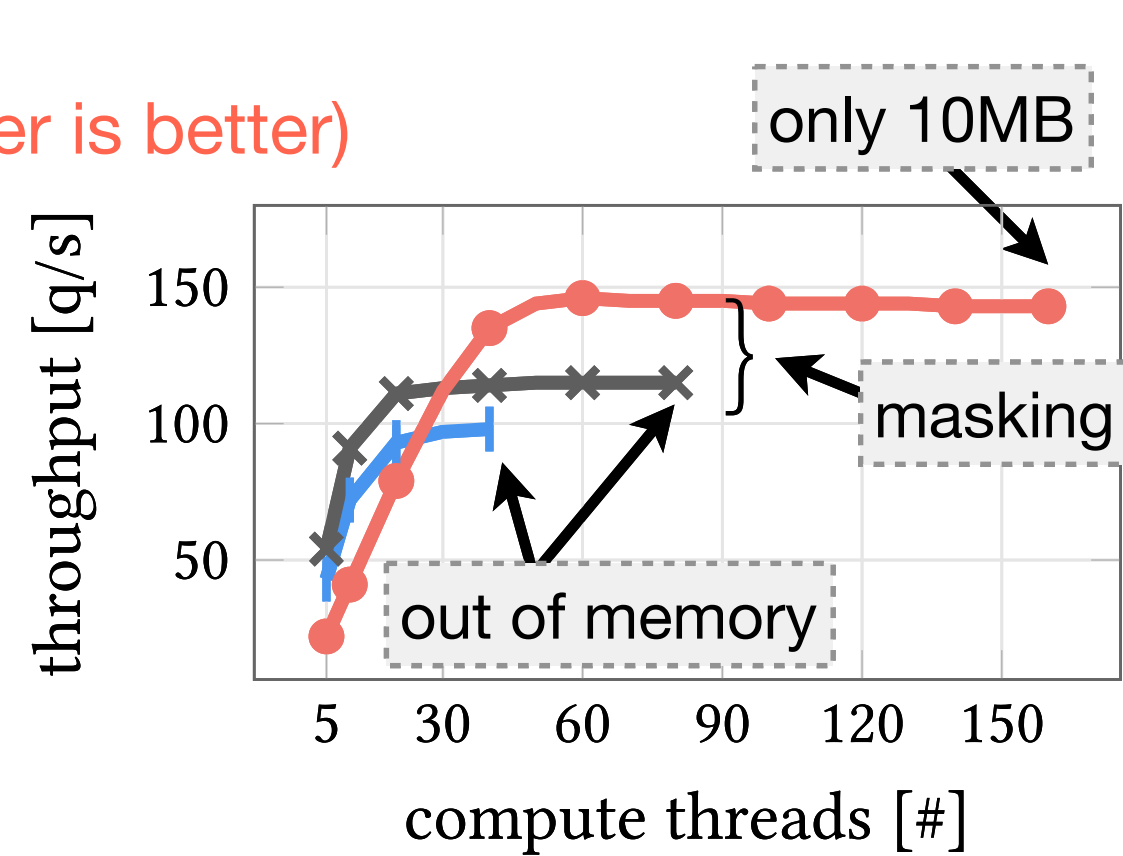
average query size: 6

## SETUP:

- 5 CNs (32 threads / 10GB RAM)
- 4 MNs (1 thread / 96GB RAM)
- Block size: 1KB

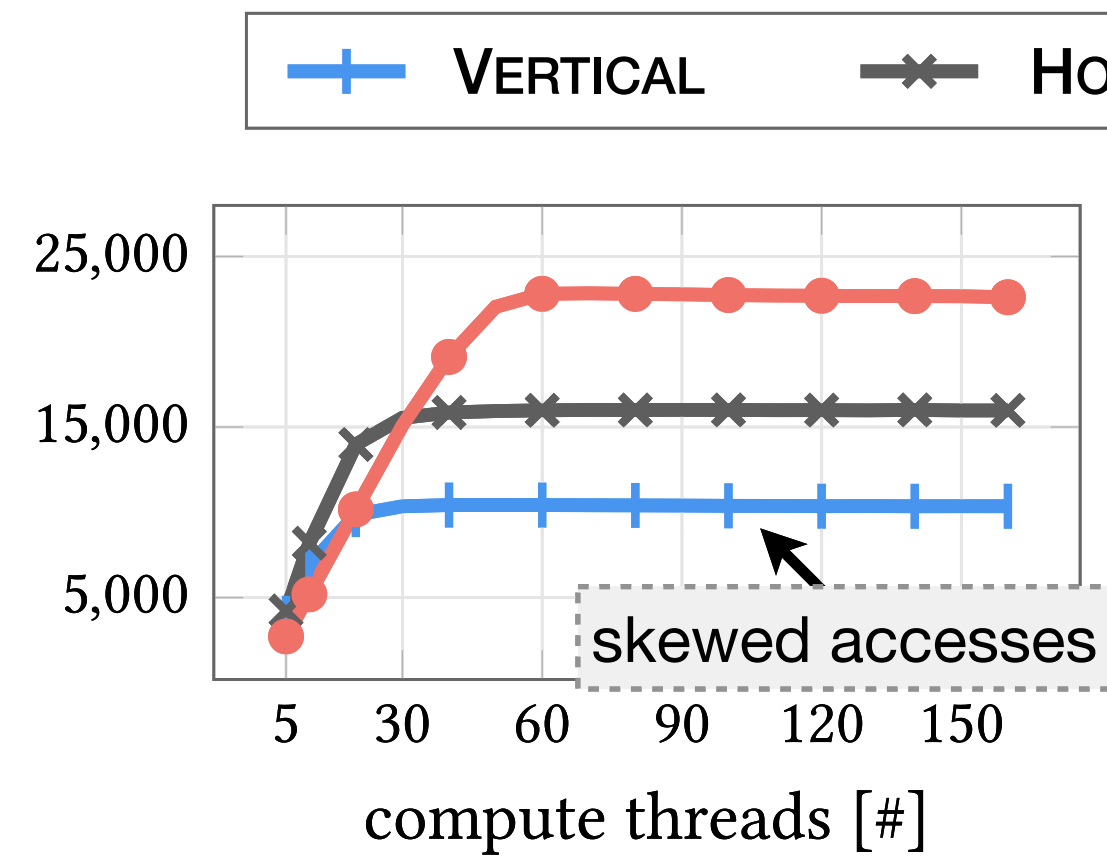
# Scalability Experiments

(higher is better)

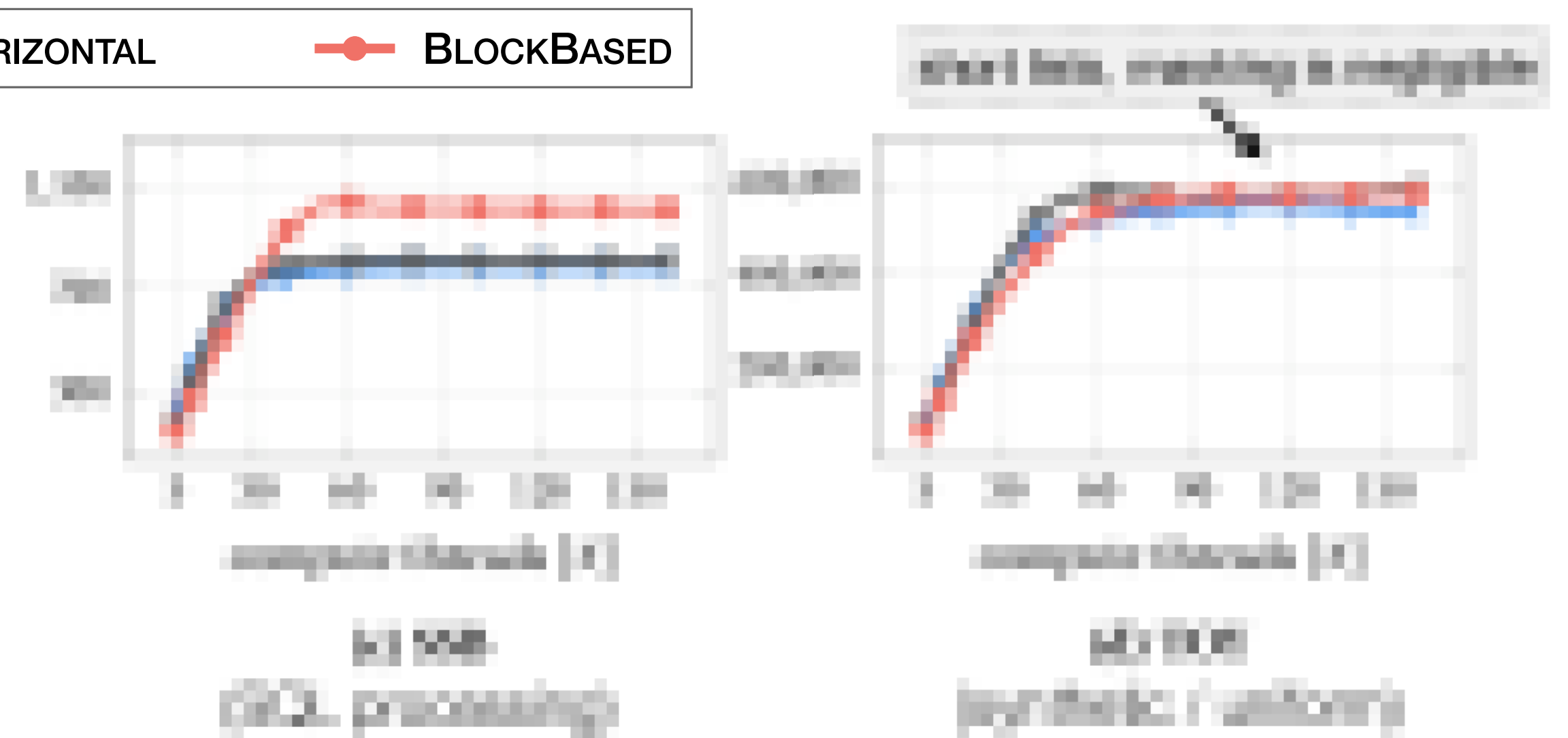


**(a) CCNEWS**  
(information retrieval)

average query size: 6



**(b) TWITTER**  
(graph analytics)



MORE RESULTS IN THE PAPER

## SETUP:

- 5 CNs (32 threads / 10GB RAM)
- 4 MNs (1 thread / 96GB RAM)
- Block size: 1KB

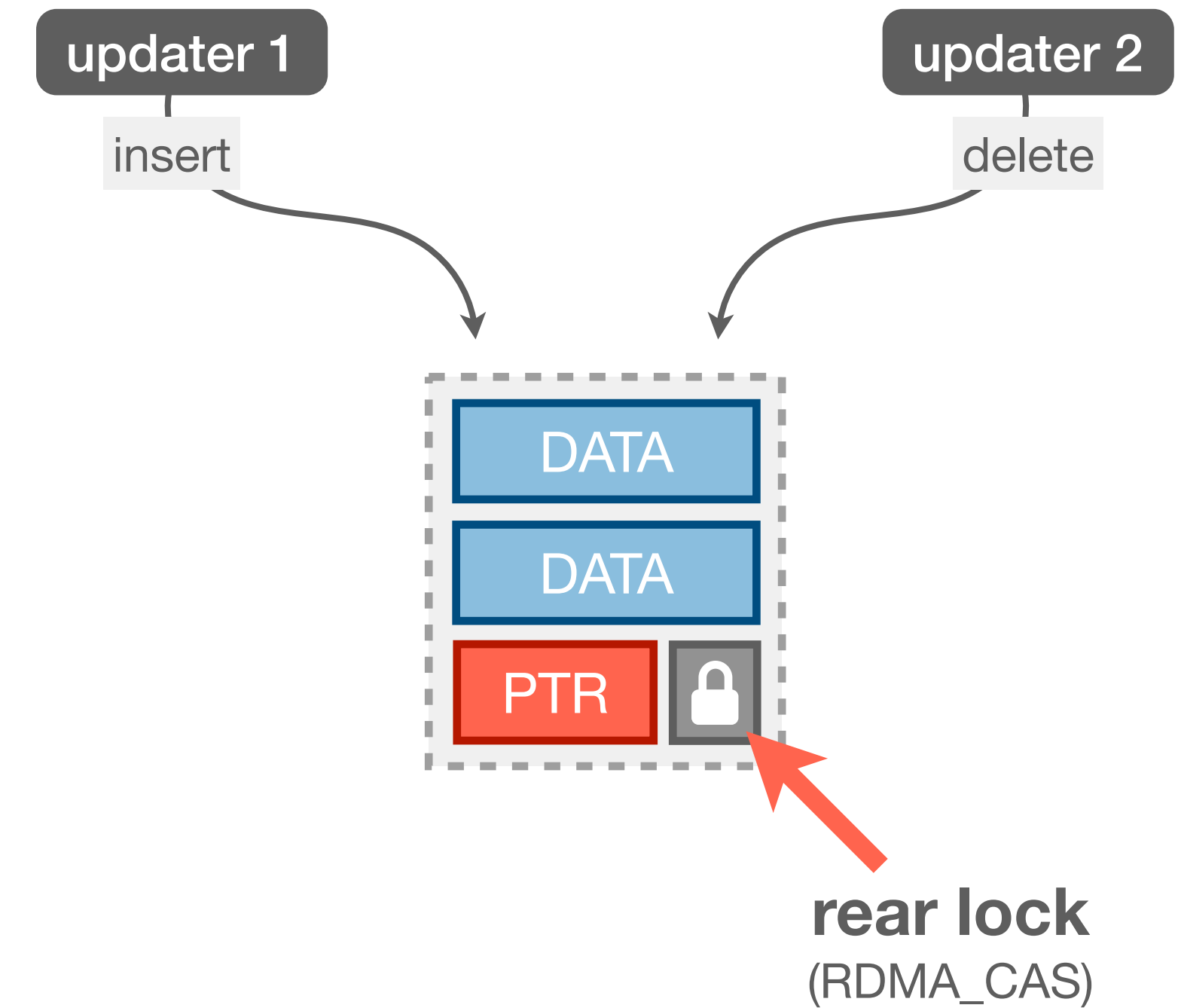
# INCREMENTAL INDEX UPDATES

# Index Updates

Protocol for incremental concurrent index updates

## MAIN CHALLENGES

- Write-write conflicts
  - Rear write locks



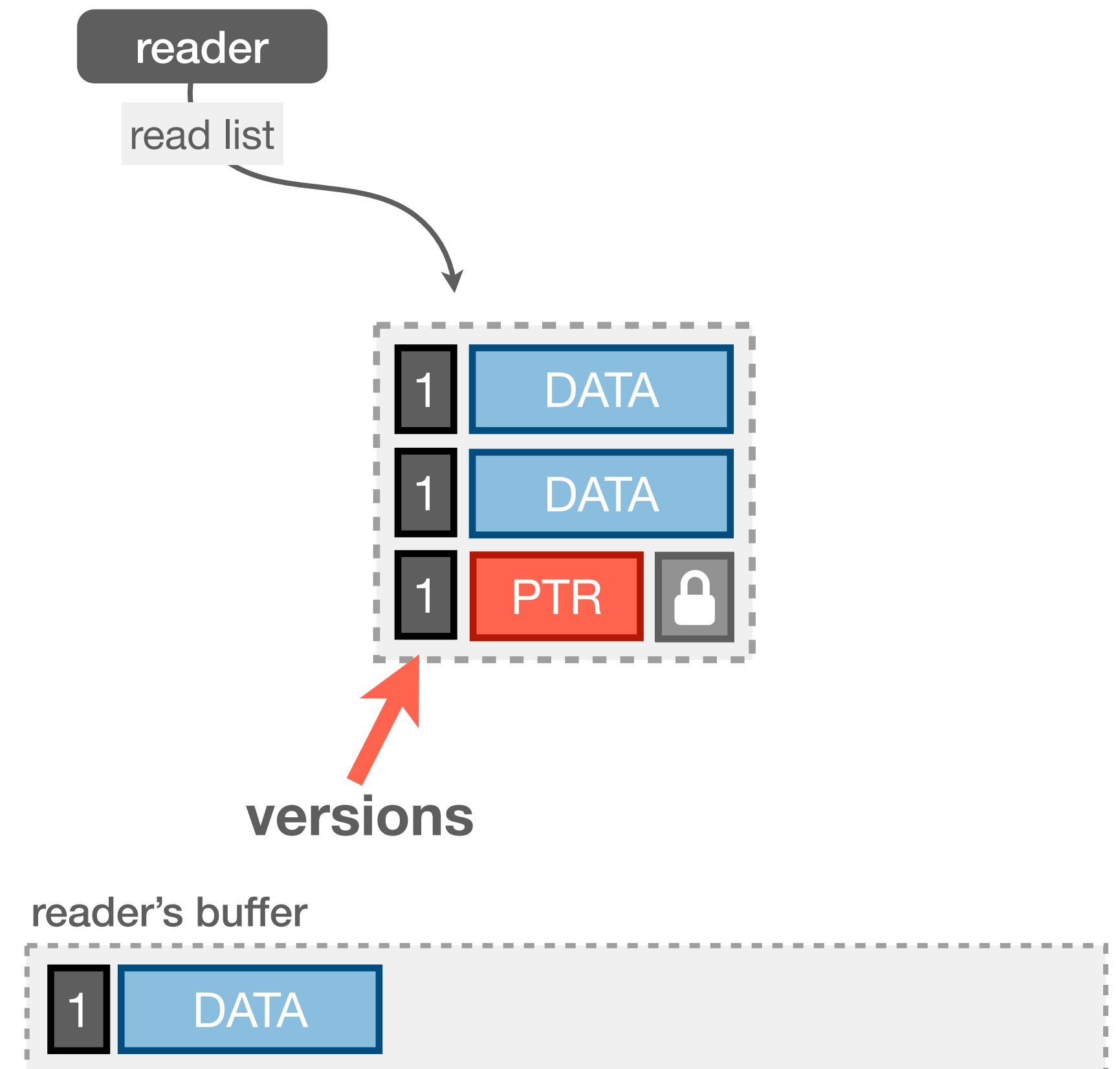
# Index Updates

Protocol for incremental concurrent index updates

## MAIN CHALLENGES

- Write-write conflicts
  - Rear write locks
- Read-write conflicts
  - Verification through versioning (lock-free)

OPTIMISTIC CONCURRENCY CONTROL



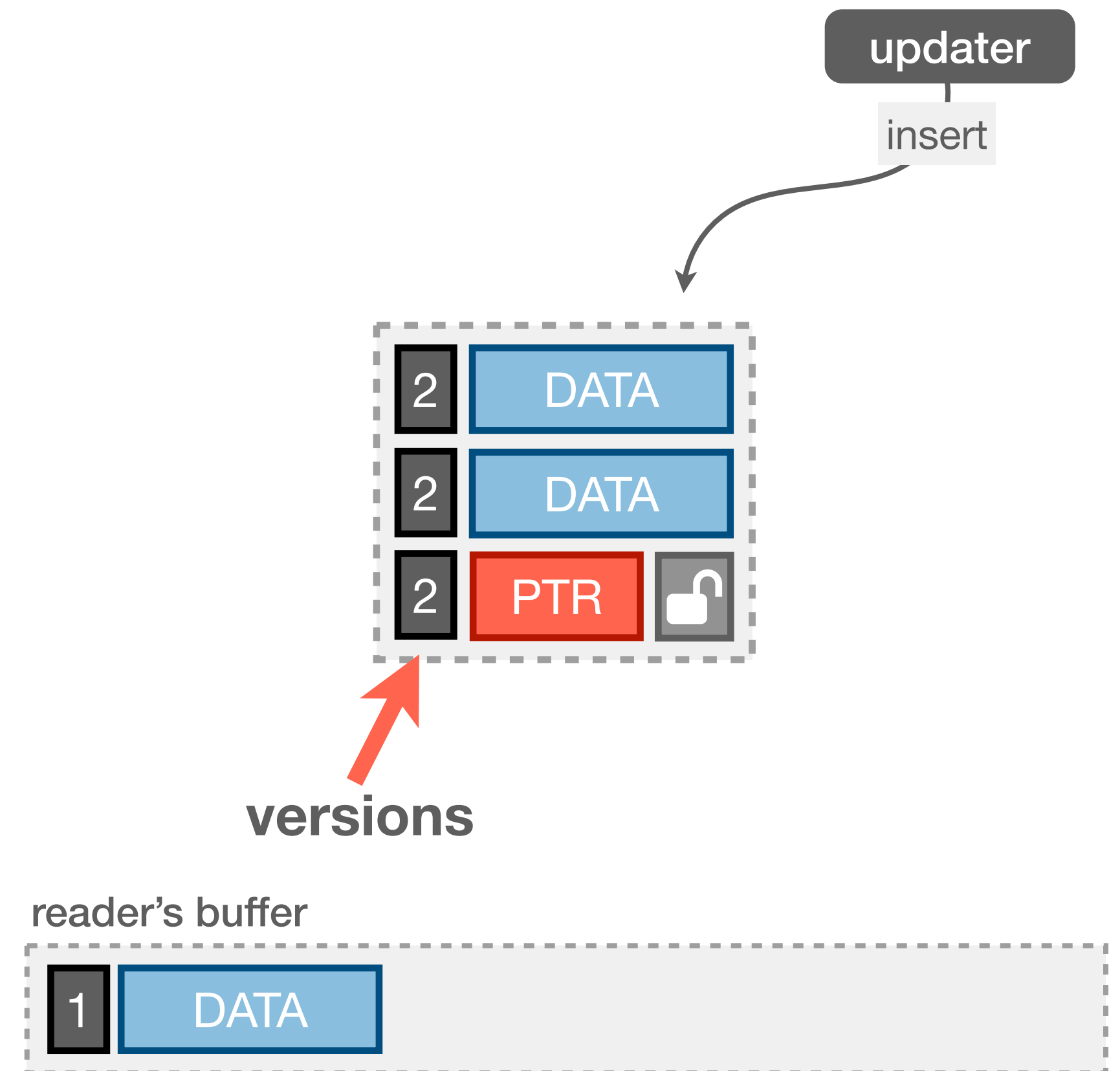
# Index Updates

Protocol for incremental concurrent index updates

## MAIN CHALLENGES

- Write-write conflicts
  - Rear write locks
- Read-write conflicts
  - Verification through versioning (lock-free)

OPTIMISTIC CONCURRENCY CONTROL



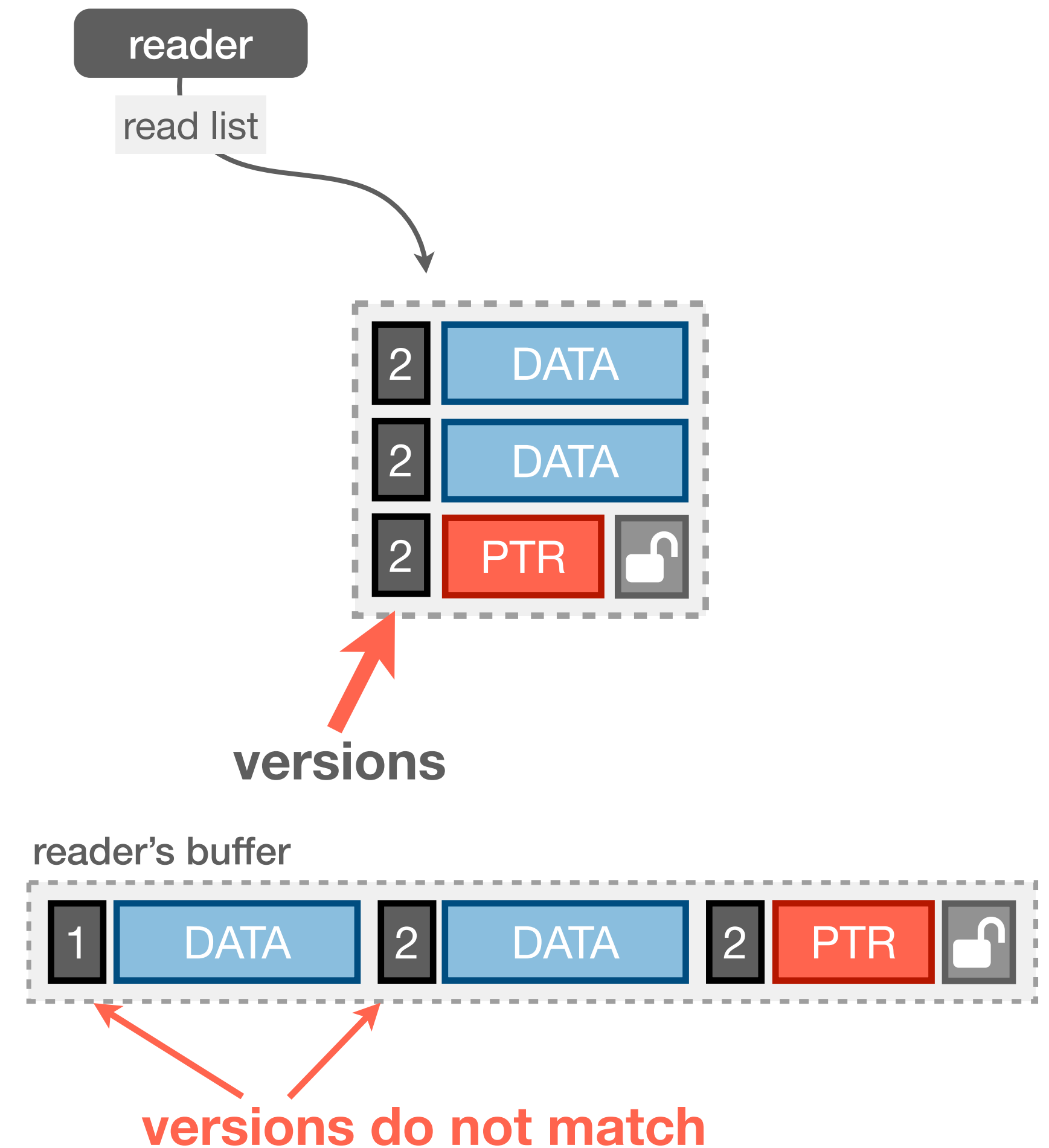
# Index Updates

Protocol for incremental concurrent index updates

## MAIN CHALLENGES

- Write-write conflicts
  - Rear write locks
- Read-write conflicts
  - Verification through versioning (lock-free)

OPTIMISTIC CONCURRENCY CONTROL





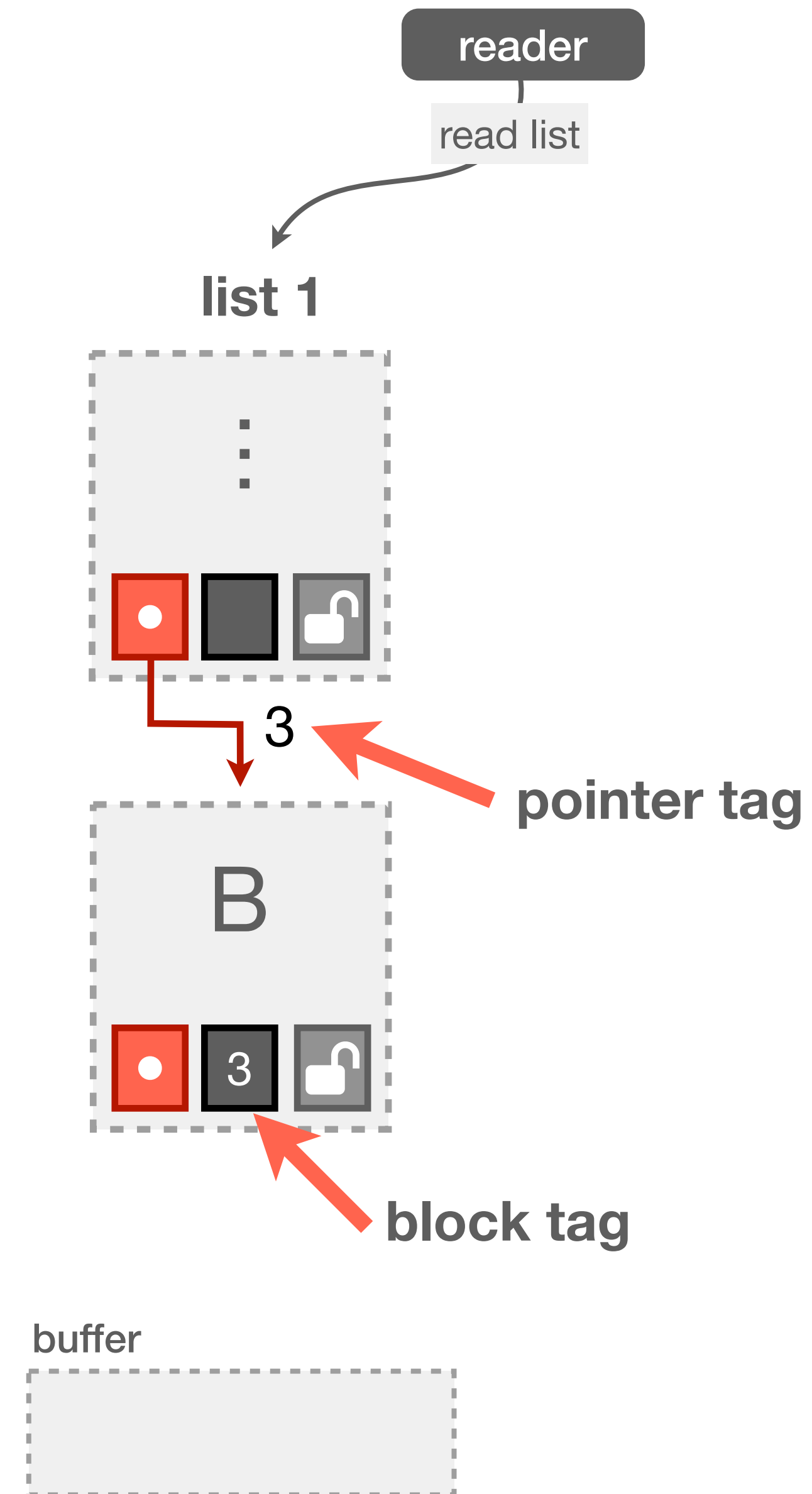
# Index Updates

Protocol for incremental concurrent index updates

## MAIN CHALLENGES

- Write-write conflicts
  - Rear write locks
- Read-write conflicts
  - Verification through versioning (lock-free)
- Block re-allocations
  - Pointer/block tagging

OPTIMISTIC CONCURRENCY CONTROL



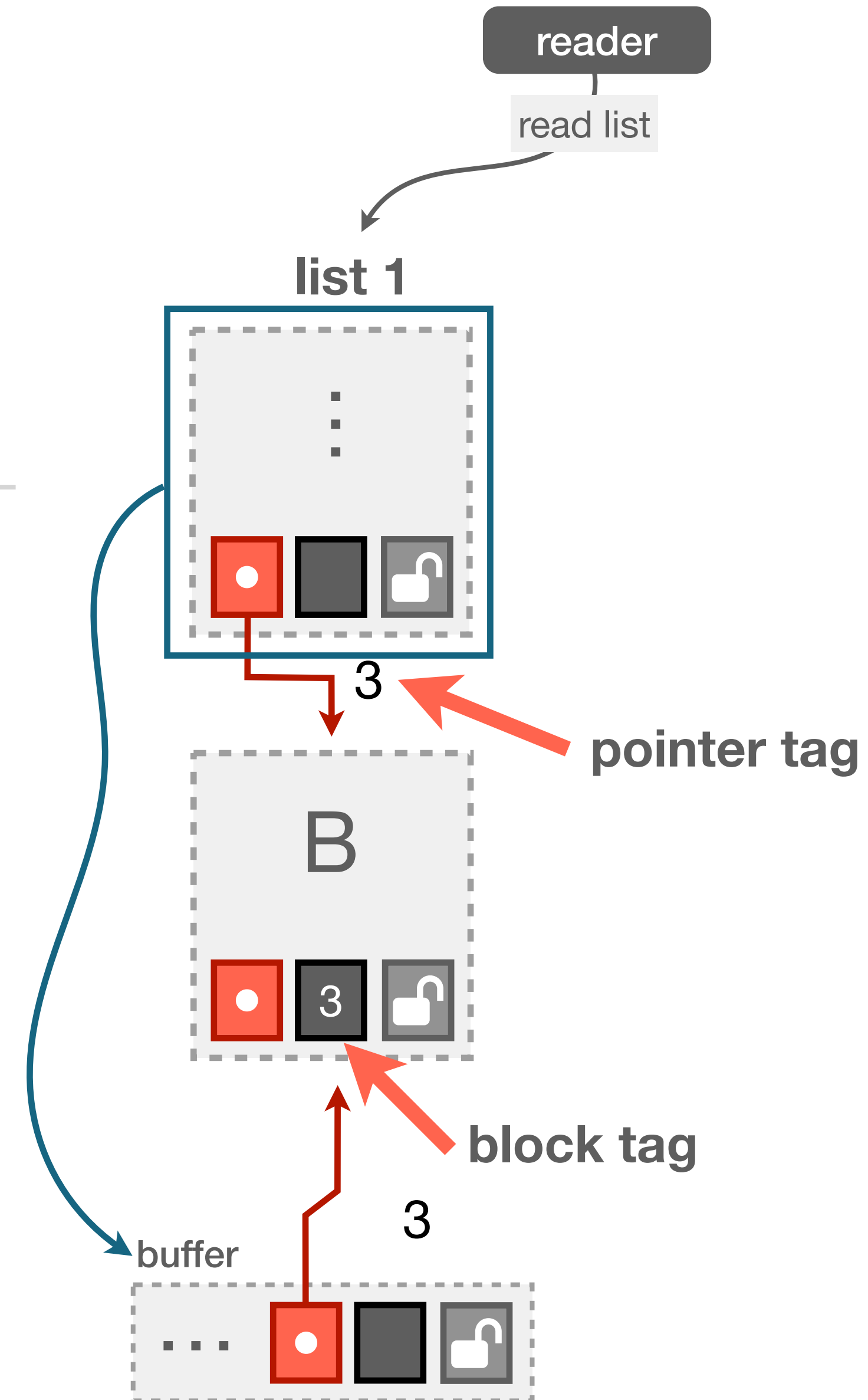
# Index Updates

Protocol for incremental concurrent index updates

## MAIN CHALLENGES

- Write-write conflicts
  - Rear write locks
- Read-write conflicts
  - Verification through versioning (lock-free)
- Block re-allocations
  - Pointer/block tagging

OPTIMISTIC CONCURRENCY CONTROL



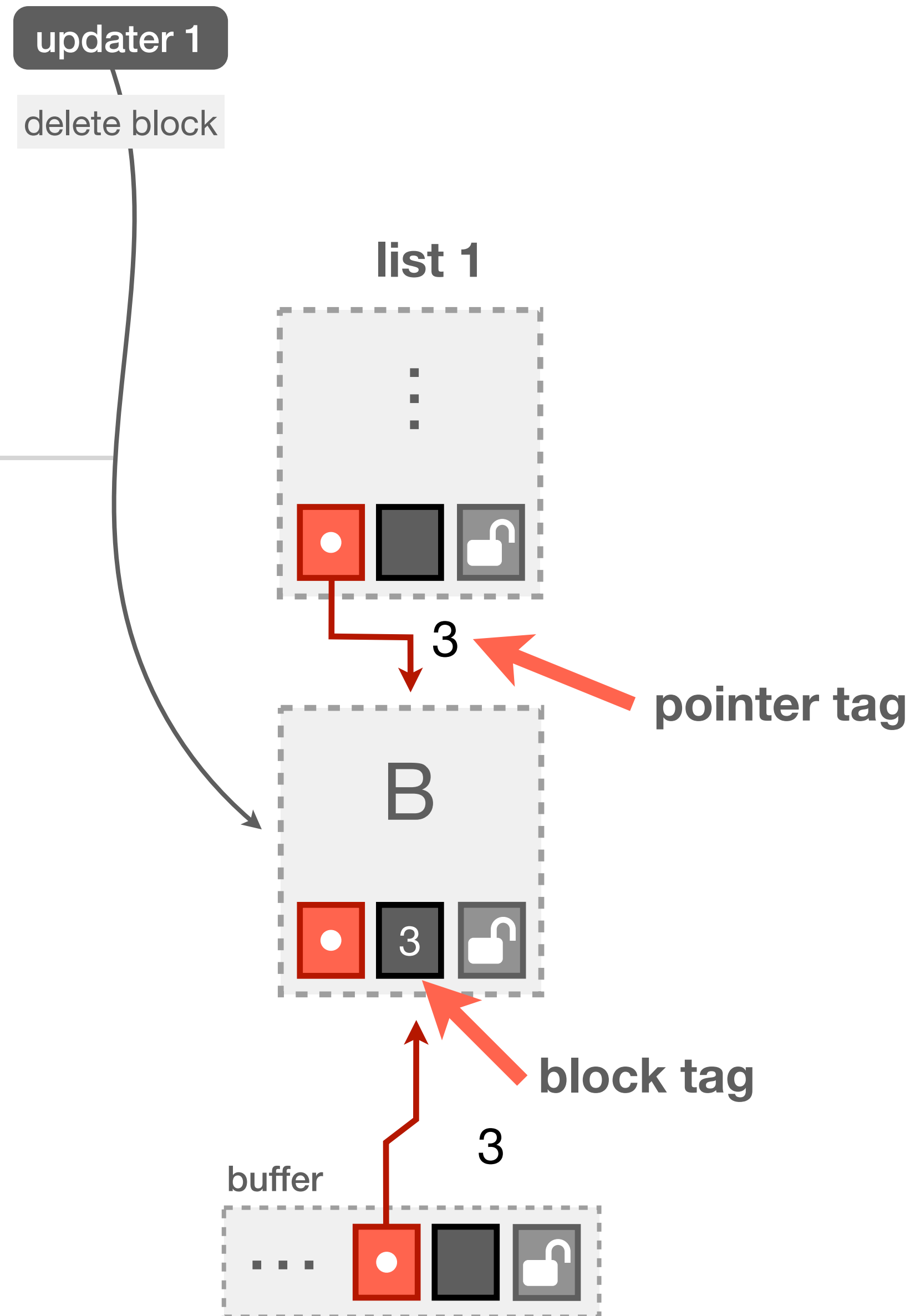
# Index Updates

Protocol for incremental concurrent index updates

## MAIN CHALLENGES

- Write-write conflicts
  - Rear write locks
- Read-write conflicts
  - Verification through versioning (lock-free)
- Block re-allocations
  - Pointer/block tagging

OPTIMISTIC CONCURRENCY CONTROL



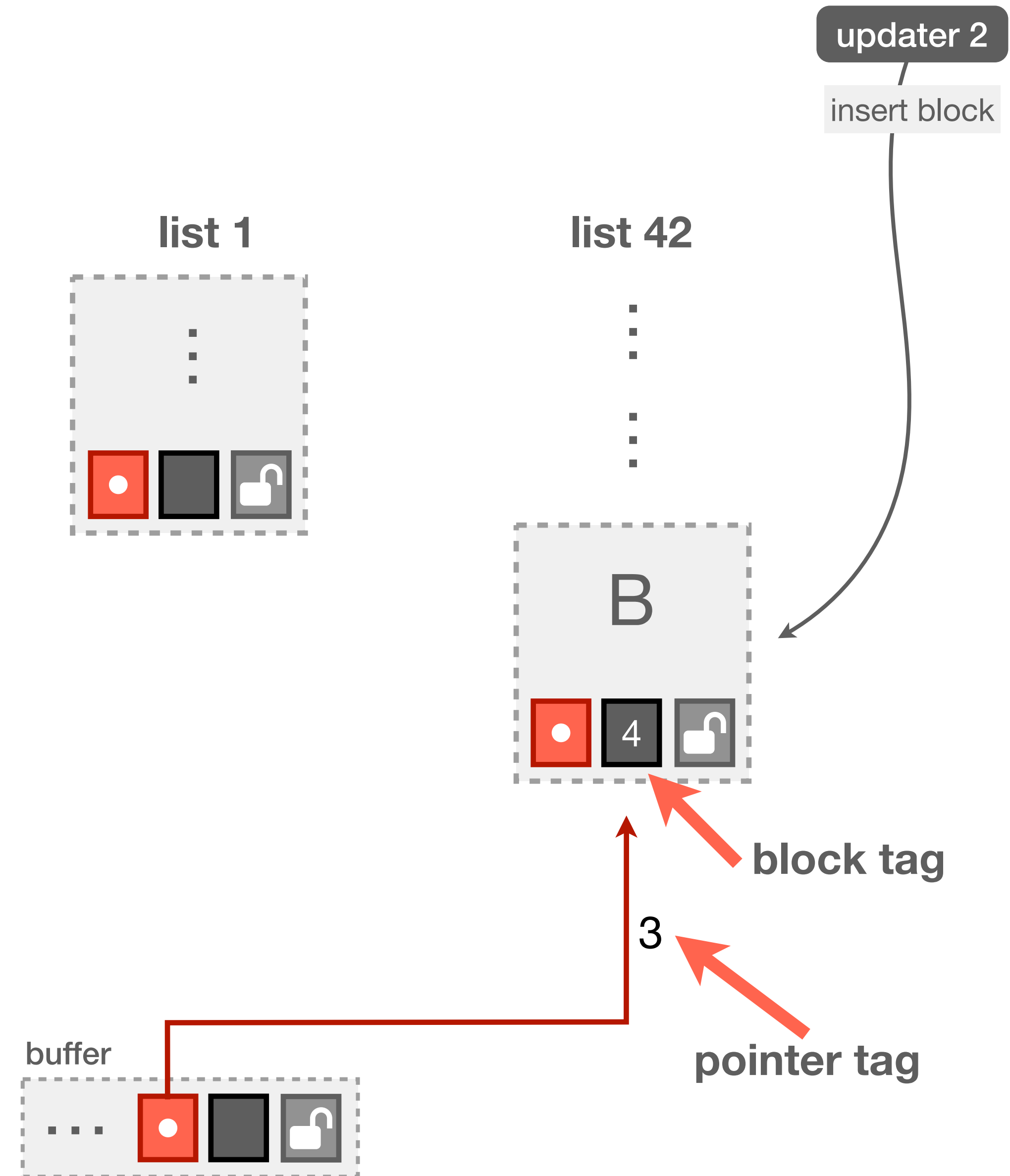
# Index Updates

Protocol for incremental concurrent index updates

## MAIN CHALLENGES

- Write-write conflicts
  - Rear write locks
- Read-write conflicts
  - Verification through versioning (lock-free)
- Block re-allocations
  - Pointer/block tagging

OPTIMISTIC CONCURRENCY CONTROL



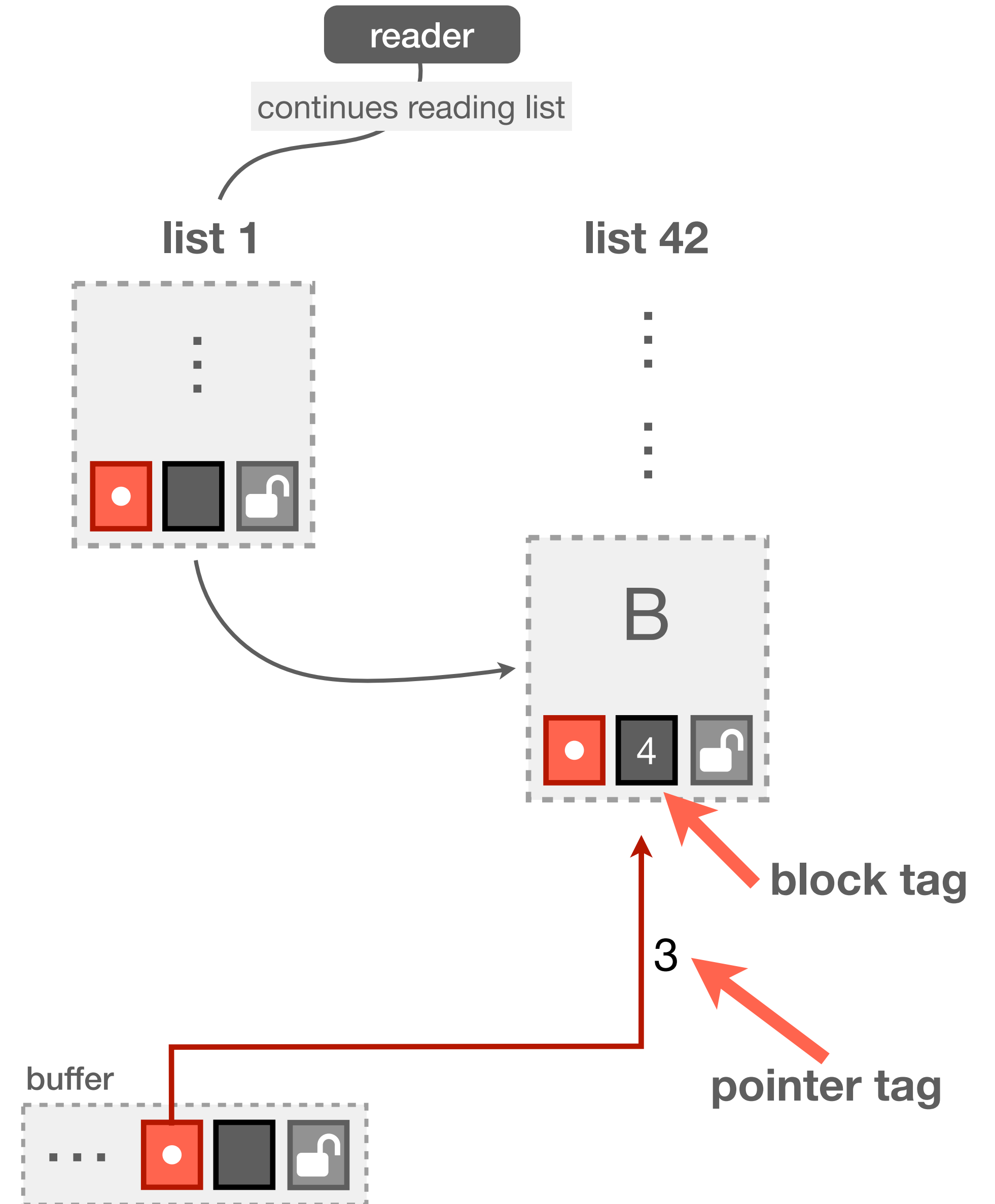
# Index Updates

Protocol for incremental concurrent index updates

## MAIN CHALLENGES

- Write-write conflicts
  - Rear write locks
- Read-write conflicts
  - Verification through versioning (lock-free)
- Block re-allocations
  - Pointer/block tagging

OPTIMISTIC CONCURRENCY CONTROL



# Conclusion

- Identified key performance bottlenecks of traditional partitioning schemes

not efficient under memory disaggregation

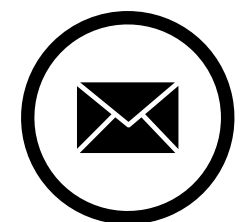
- Proposed a scalable inverted list index design for disaggregated memory
- Developed a protocol to support fast concurrent index updates:
  - Lock-free read queries
  - Fine-grained write locks for updates
- Scalability evaluation on real-world datasets

# SCALABLE DISTRIBUTED INVERTED LIST INDEXES IN DISAGGREGATED MEMORY

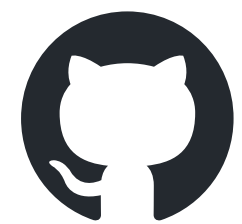
**MANUEL WIDMOSER**

DANIEL KOCHER

NIKOLAUS AUGSTEN



`manuel.widmoser@plus.ac.at`



`github.com/DatabaseGroup/rdma-inverted-index`



link to paper