

# WiscKey: Separating Keys from Values in SSD-Conscious Storage

Lanyue Lu, Thanumalayan Pillai,  
Andrea C. Arpaci-Dusseau, Remzi H. Arpaci-Dusseau

University of Wisconsin-Madison

# Key-Value Stores

# Key-Value Stores

Key-value stores are important

- web indexing, e-commerce, social networks
- various key-value stores
  - hash table, b-tree
  - log-structured merge trees (LSM-trees)

# Key-Value Stores

## Key-value stores are important

- web indexing, e-commerce, social networks
- various key-value stores
  - hash table, b-tree
  - log-structured merge trees (LSM-trees)

## LSM-tree based key-value stores are popular

- optimize for write intensive workloads
- widely deployed
  - BigTable and LevelDB at Google
  - HBase, Cassandra and RocksDB at FaceBook

# Why LSM-trees ?

# Why LSM-trees ?

## Good for hard drives

- batch and write sequentially
- high sequential throughput
- sequential access up to 1000x faster than random

# Why LSM-trees ?

## Good for hard drives

- batch and write sequentially
- high sequential throughput
- sequential access up to 1000x faster than random

## Not optimal for SSDs

- large write/read amplification
- wastes device resources

# Why LSM-trees ?

## Good for hard drives

- batch and write sequentially
- high sequential throughput
- sequential access up to 1000x faster than random

## Not optimal for SSDs

- large write/read amplification
  - wastes device resources
- unique characteristics of SSDs
  - fast random reads
  - internal parallelism



# Our Solution: Wisckey

# Our Solution: Wisckey

Separate keys from values

# Our Solution: WisckKey

Separate keys from values

→ decouple sorting and garbage collection

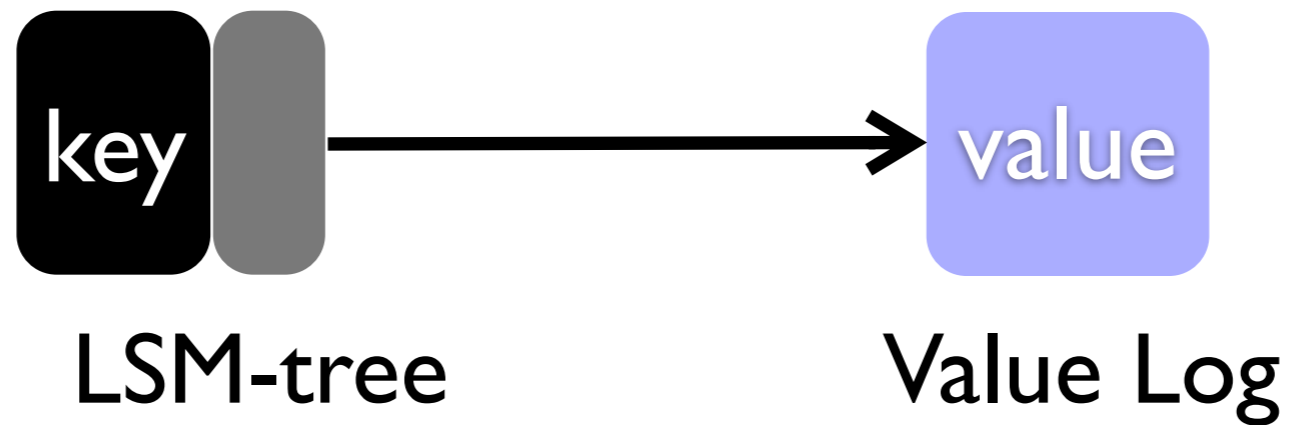


LSM-tree

# Our Solution: WisckKey

## Separate keys from values

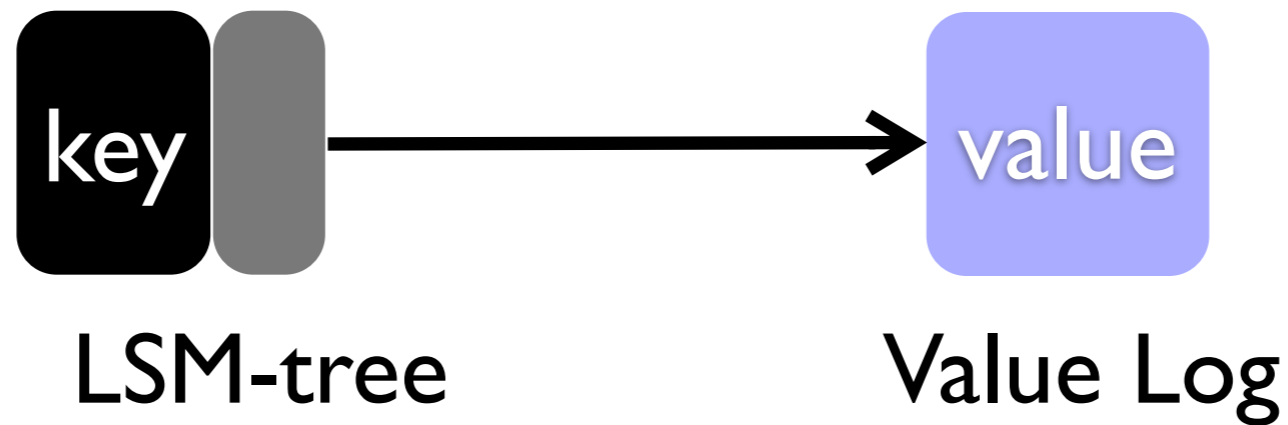
→ decouple sorting and garbage collection



# Our Solution: WisckKey

## Separate keys from values

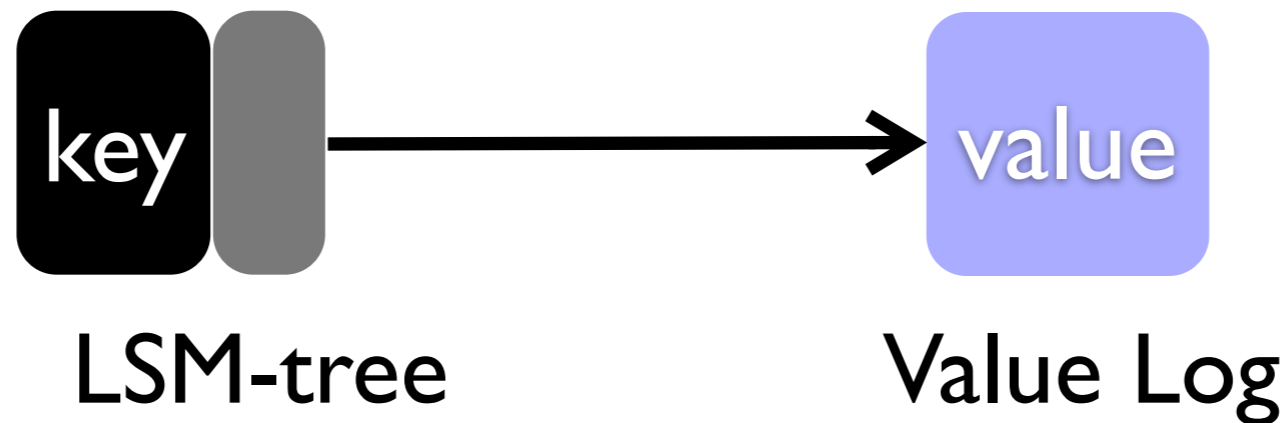
- decouple sorting and garbage collection
- harness SSD's internal parallelism for range queries



# Our Solution: WisckKey

## Separate keys from values

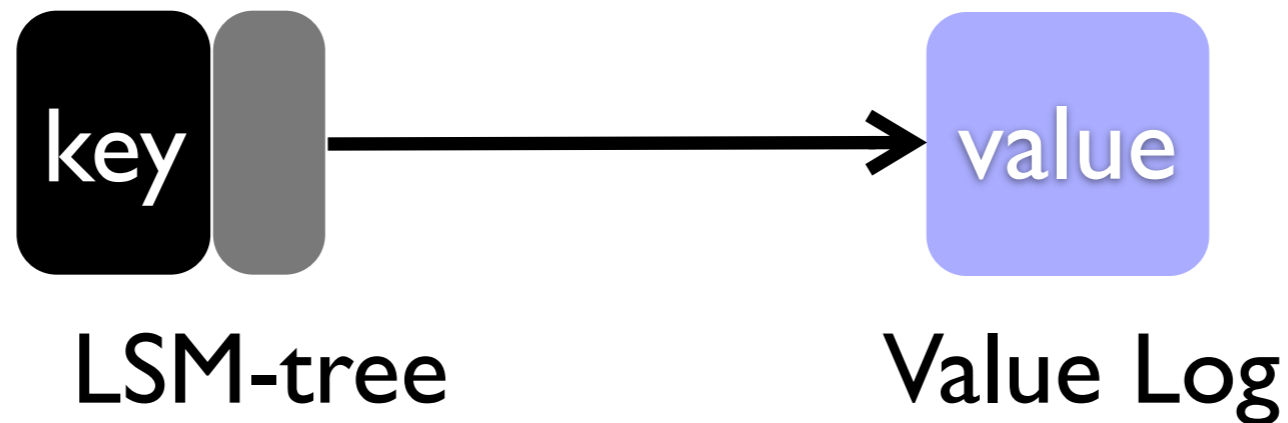
- decouple sorting and garbage collection
- harness SSD's internal parallelism for range queries
- online and light-weight garbage collection



# Our Solution: WisckKey

## Separate keys from values

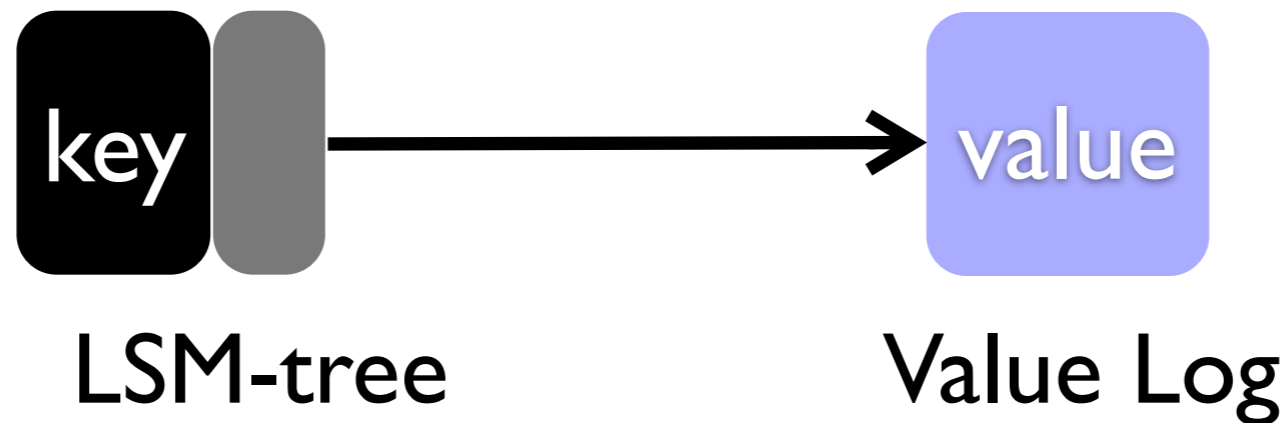
- decouple sorting and garbage collection
- harness SSD's internal parallelism for range queries
- online and light-weight garbage collection
- minimize I/O amplification and crash consistent



# Our Solution: WisckKey

## Separate keys from values

- decouple sorting and garbage collection
- harness SSD's internal parallelism for range queries
- online and light-weight garbage collection
- minimize I/O amplification and crash consistent



## Performance of WisckKey

- 2.5x to 111x for loading, 1.6x to 14x for lookups



**Background**

**Key-Value Separation**

**Challenges and Optimizations**

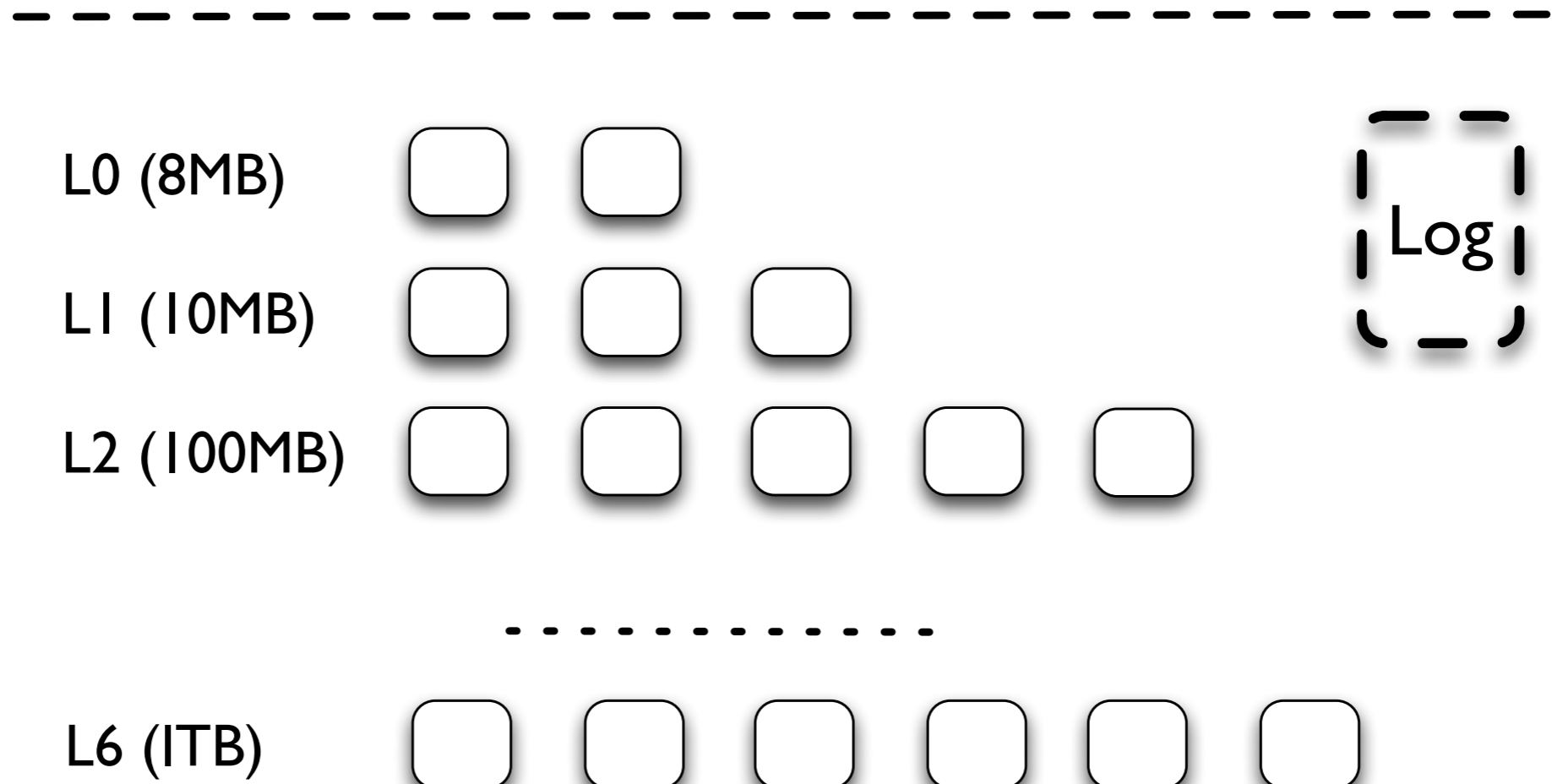
**Evaluation**

**Conclusion**

# LSM-trees: Insertion

memory

disk

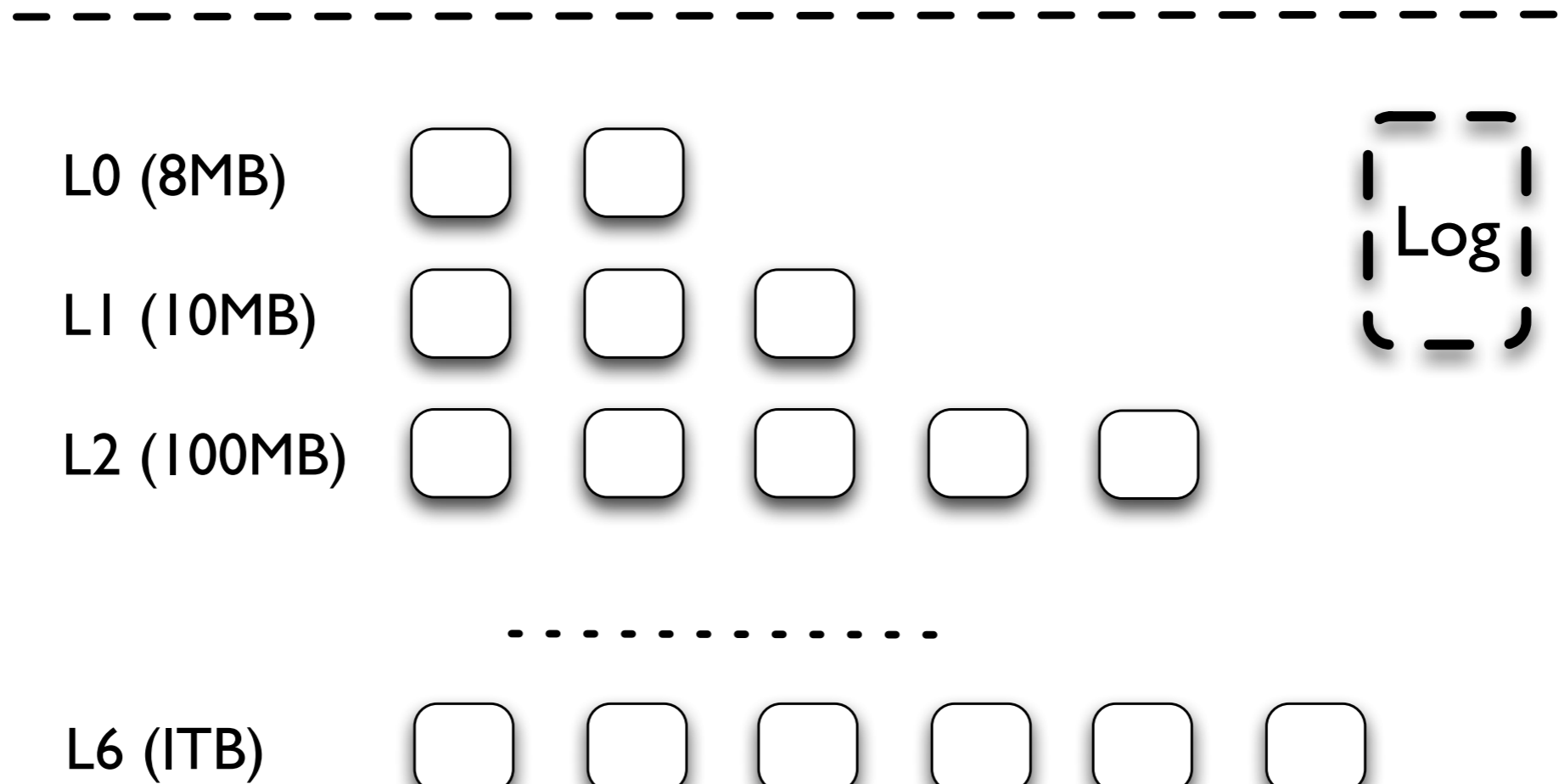


# LSM-trees: Insertion

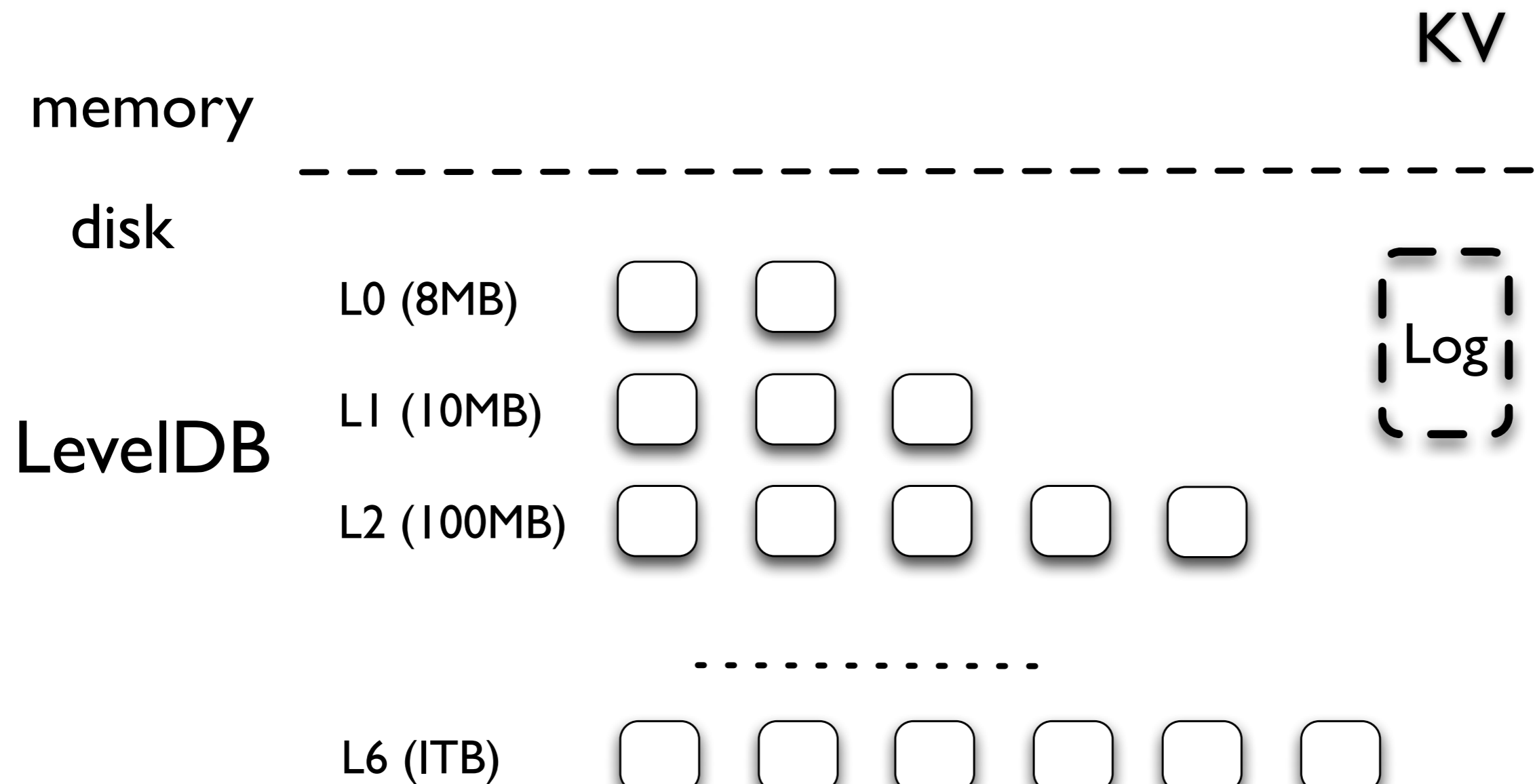
memory

disk

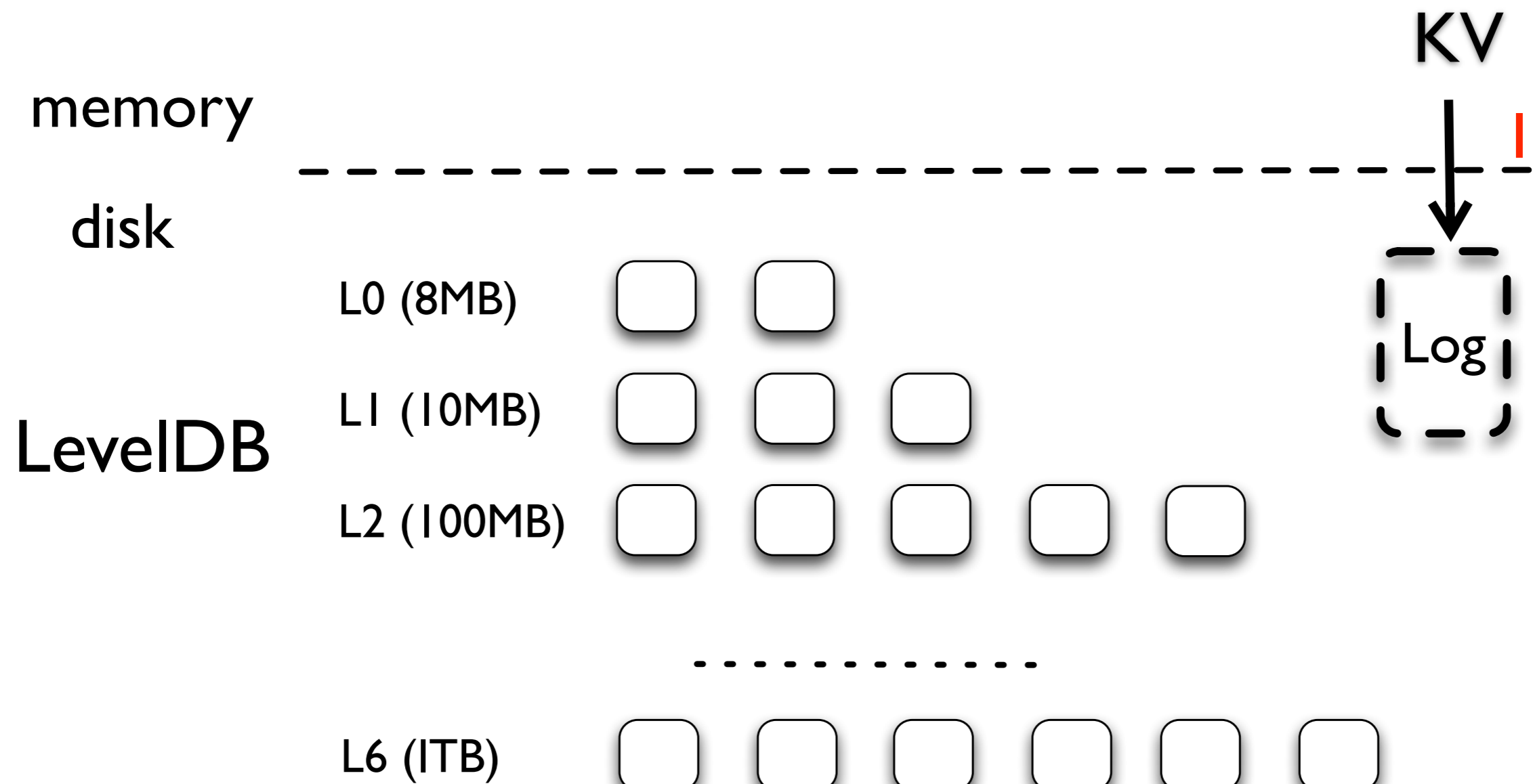
LevelDB



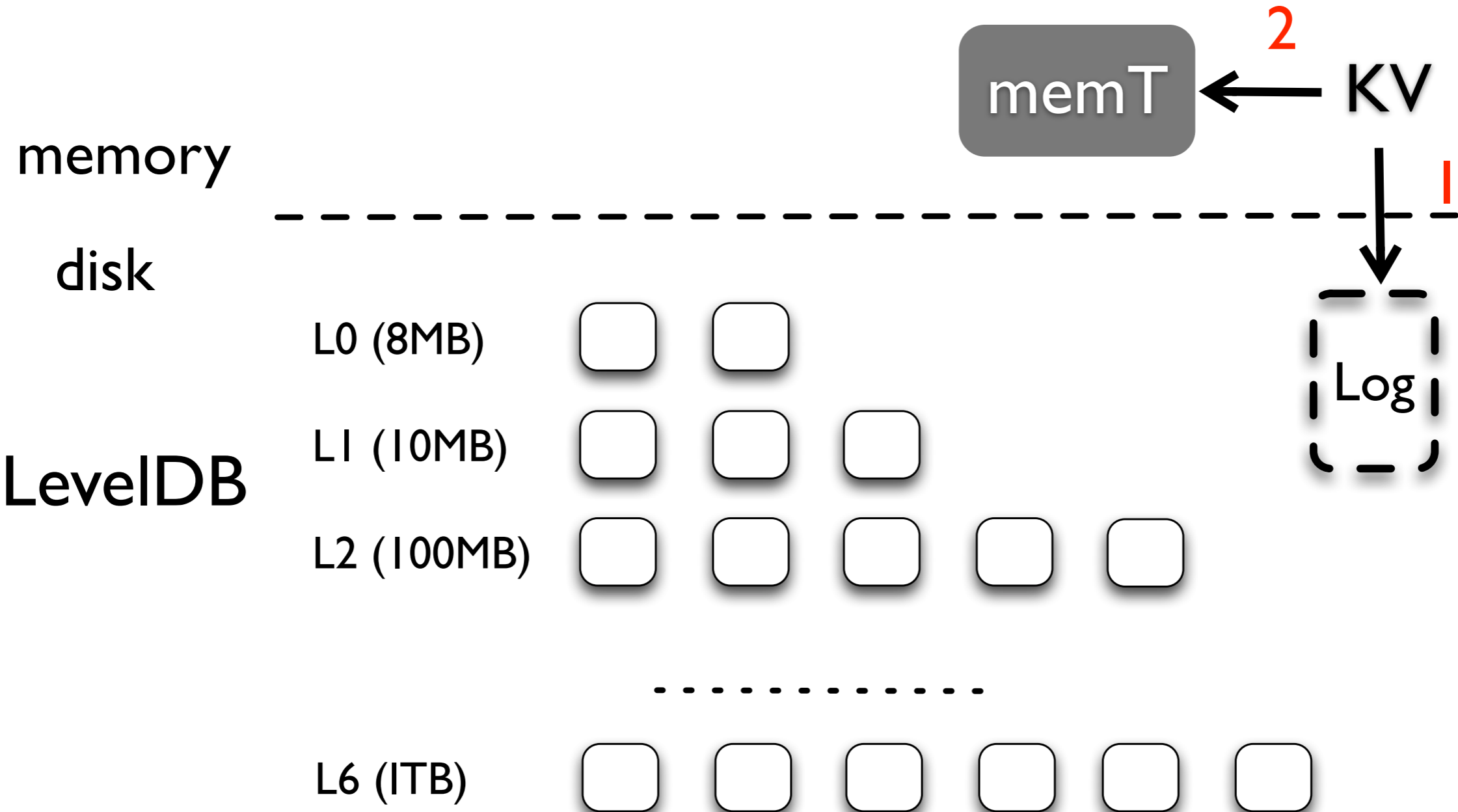
# LSM-trees: Insertion



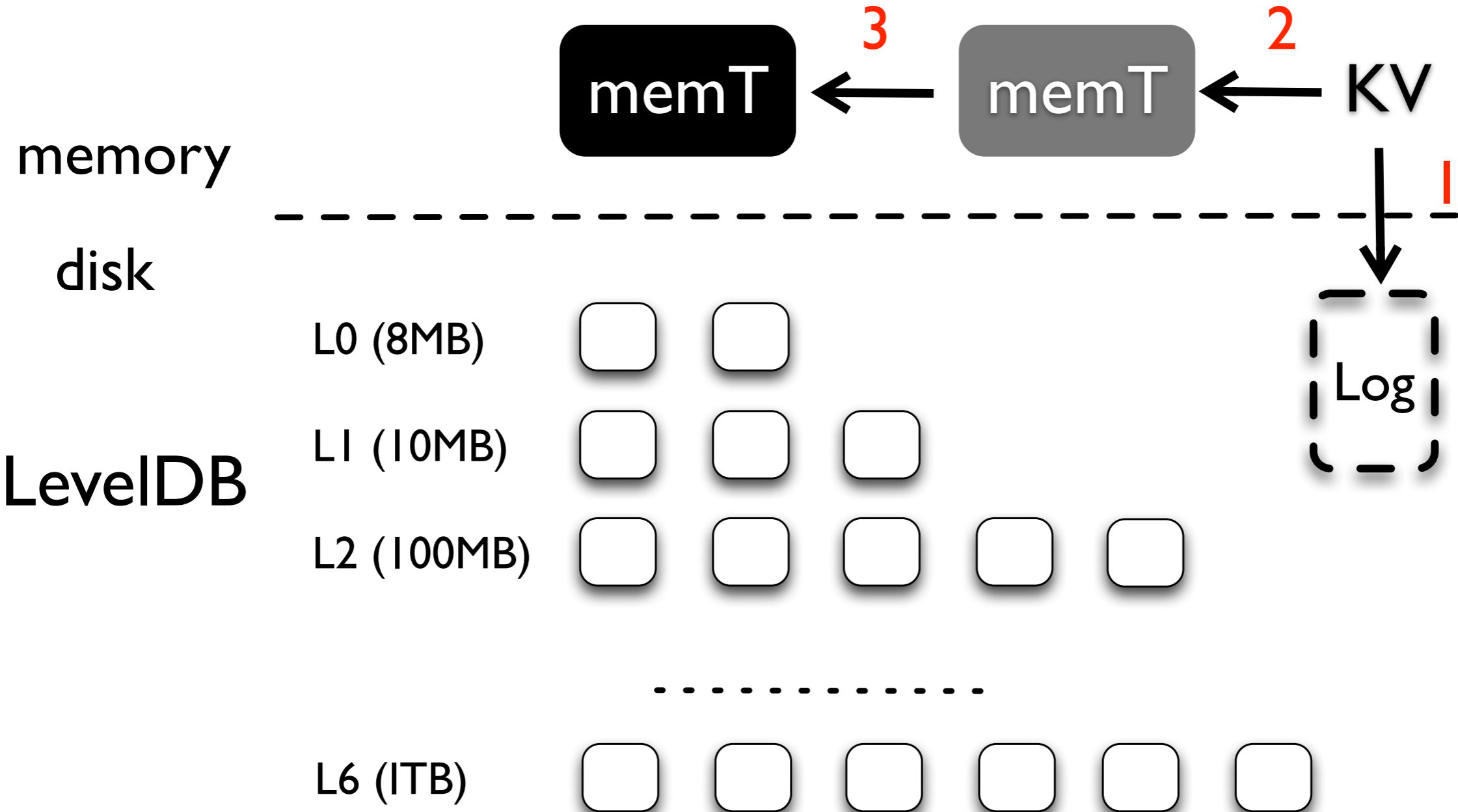
# LSM-trees: Insertion



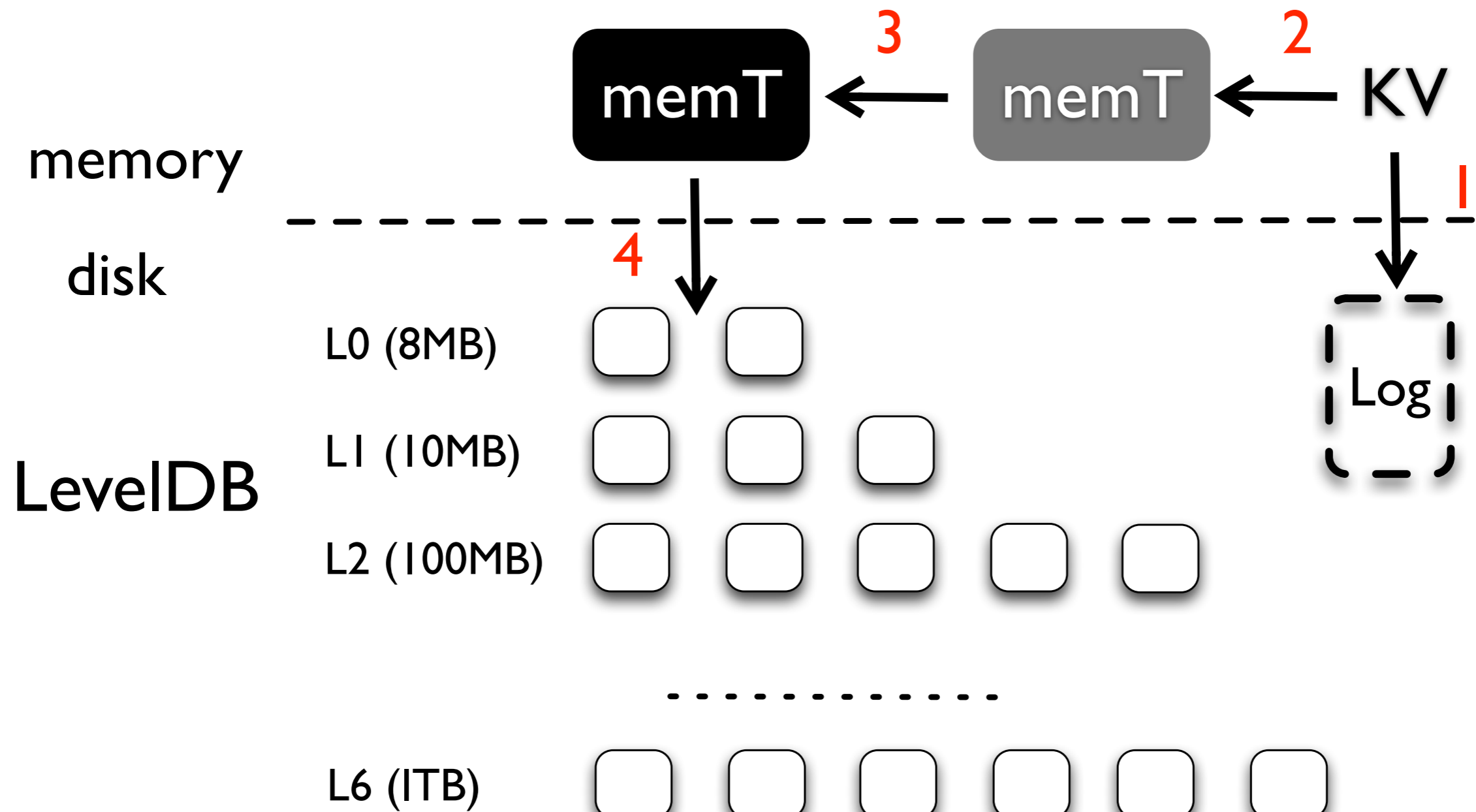
# LSM-trees: Insertion



# LSM-trees: Insertion

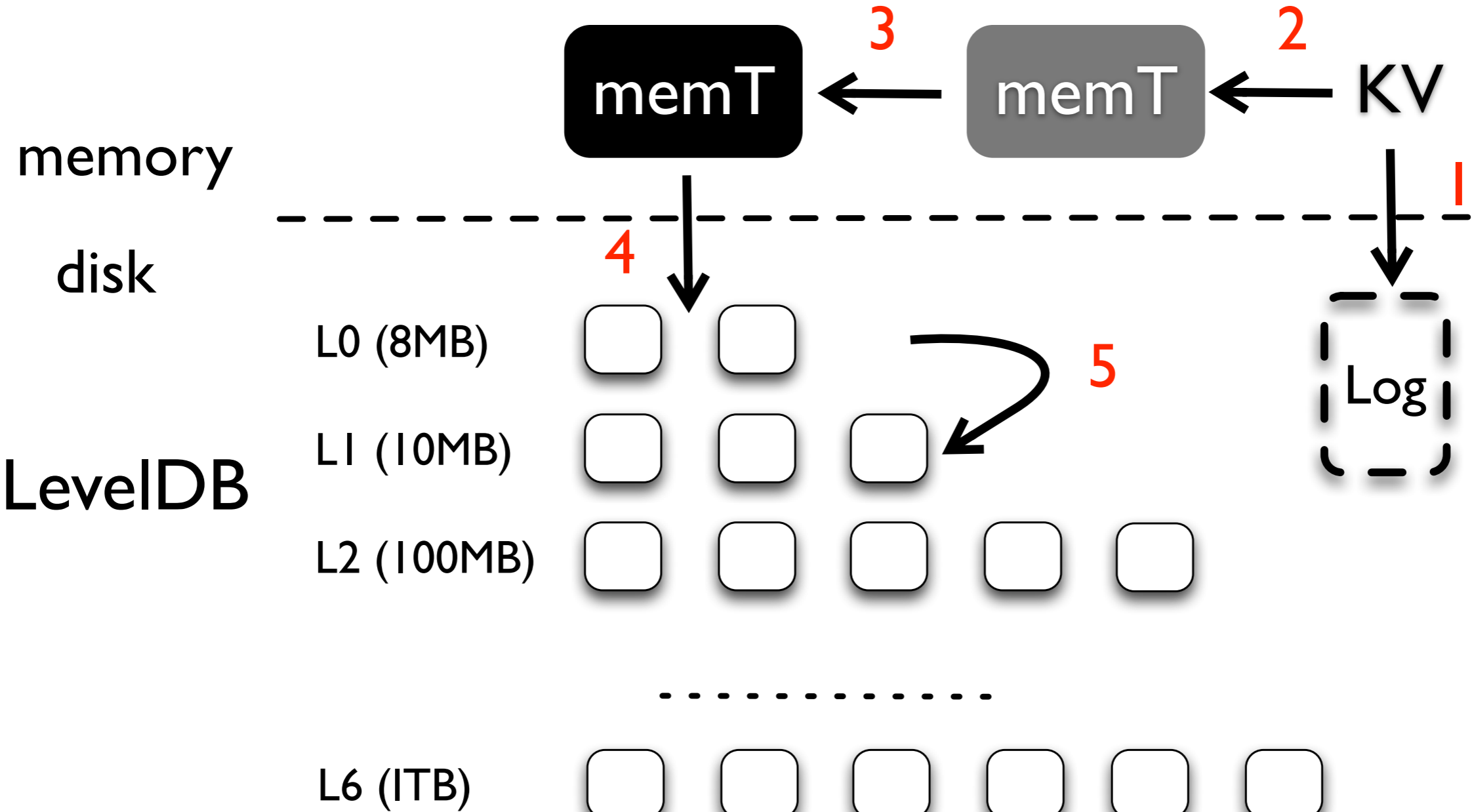


# LSM-trees: Insertion



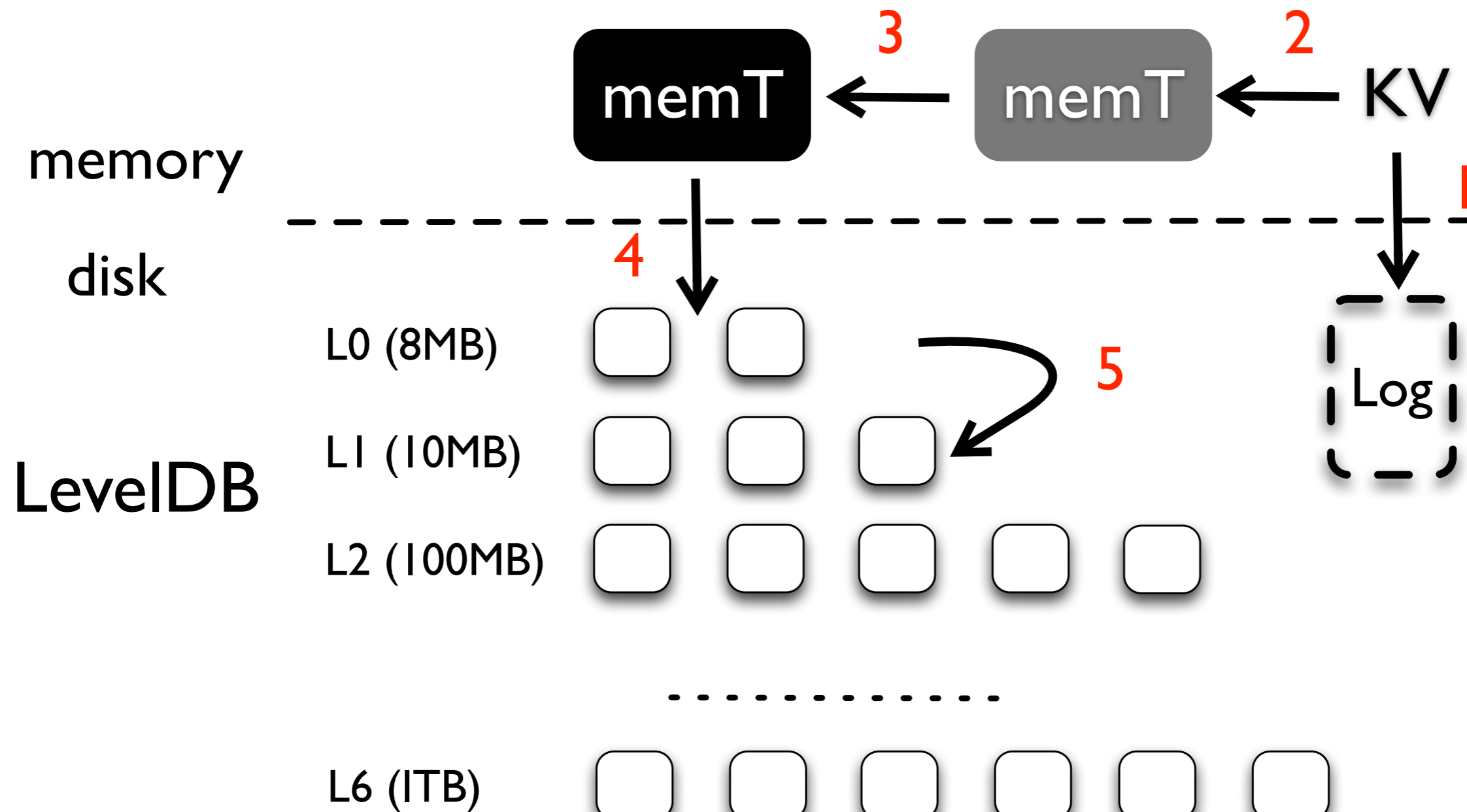


# LSM-trees: Insertion



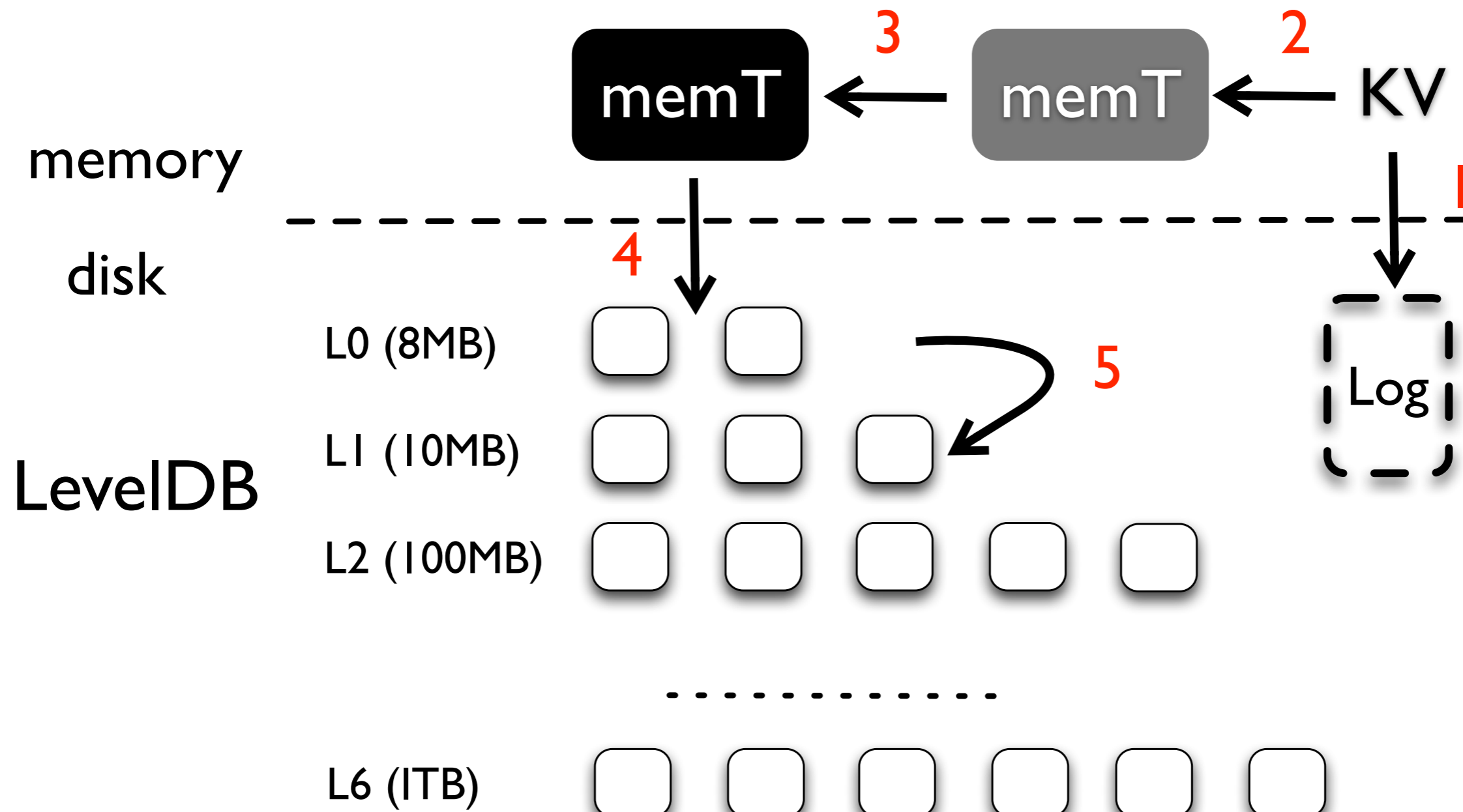
# LSM-trees: Insertion

1. Write sequentially
2. Sort data for quick lookups



# LSM-trees: Insertion

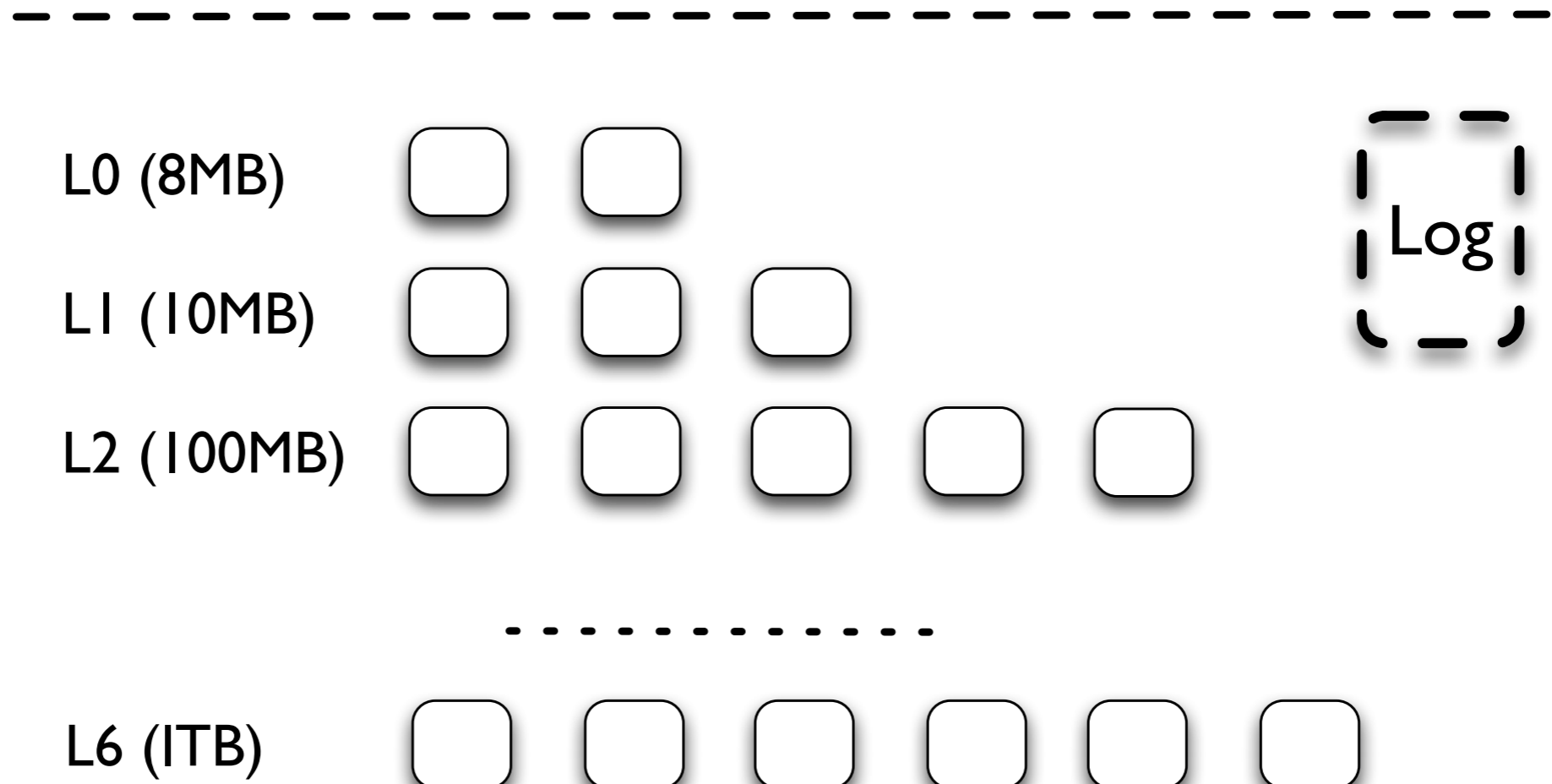
1. Write sequentially
2. Sort data for quick lookups
3. Sorting and garbage collection are coupled



# LSM-trees: Lookup

memory

disk

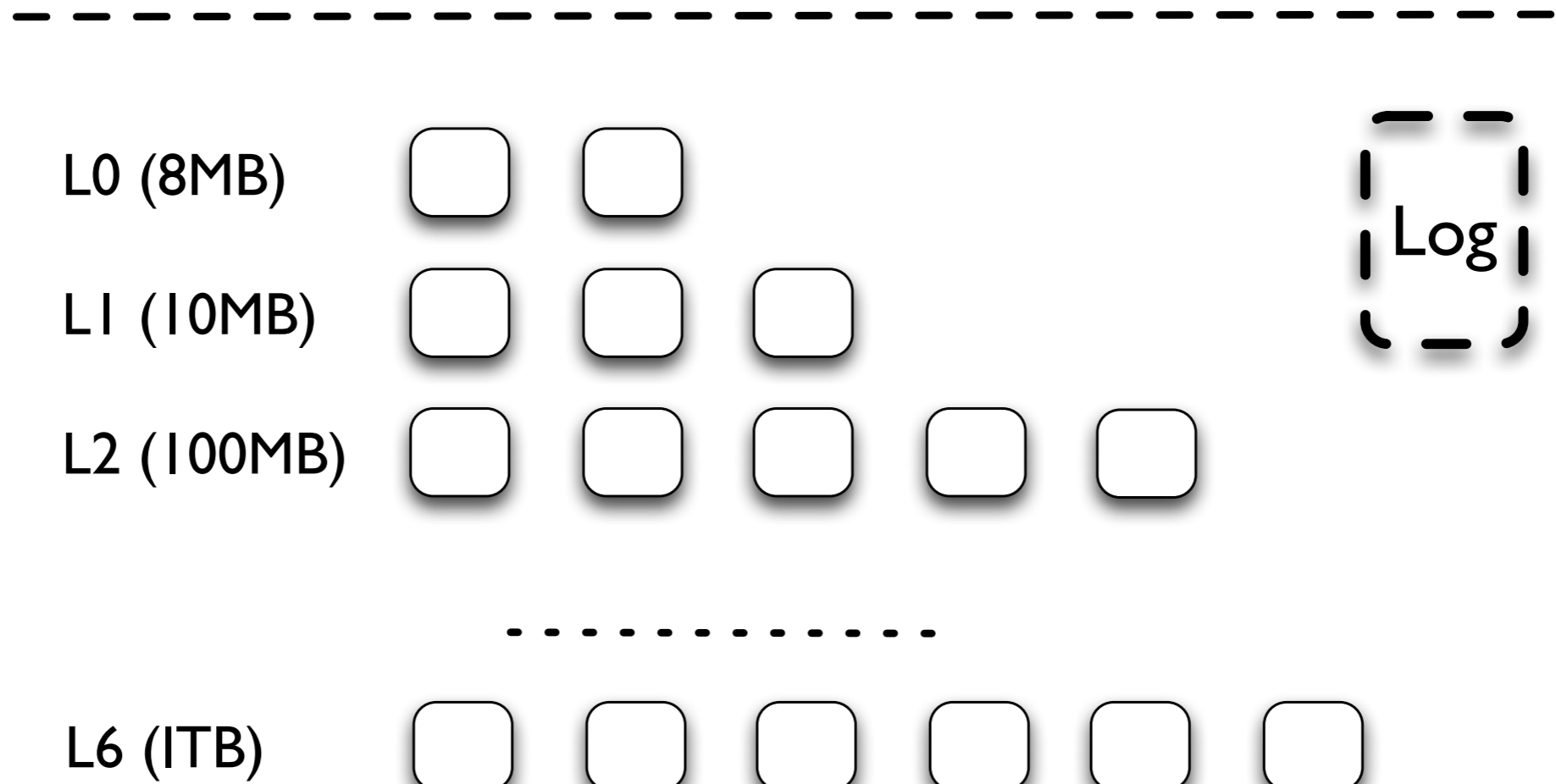


# LSM-trees: Lookup

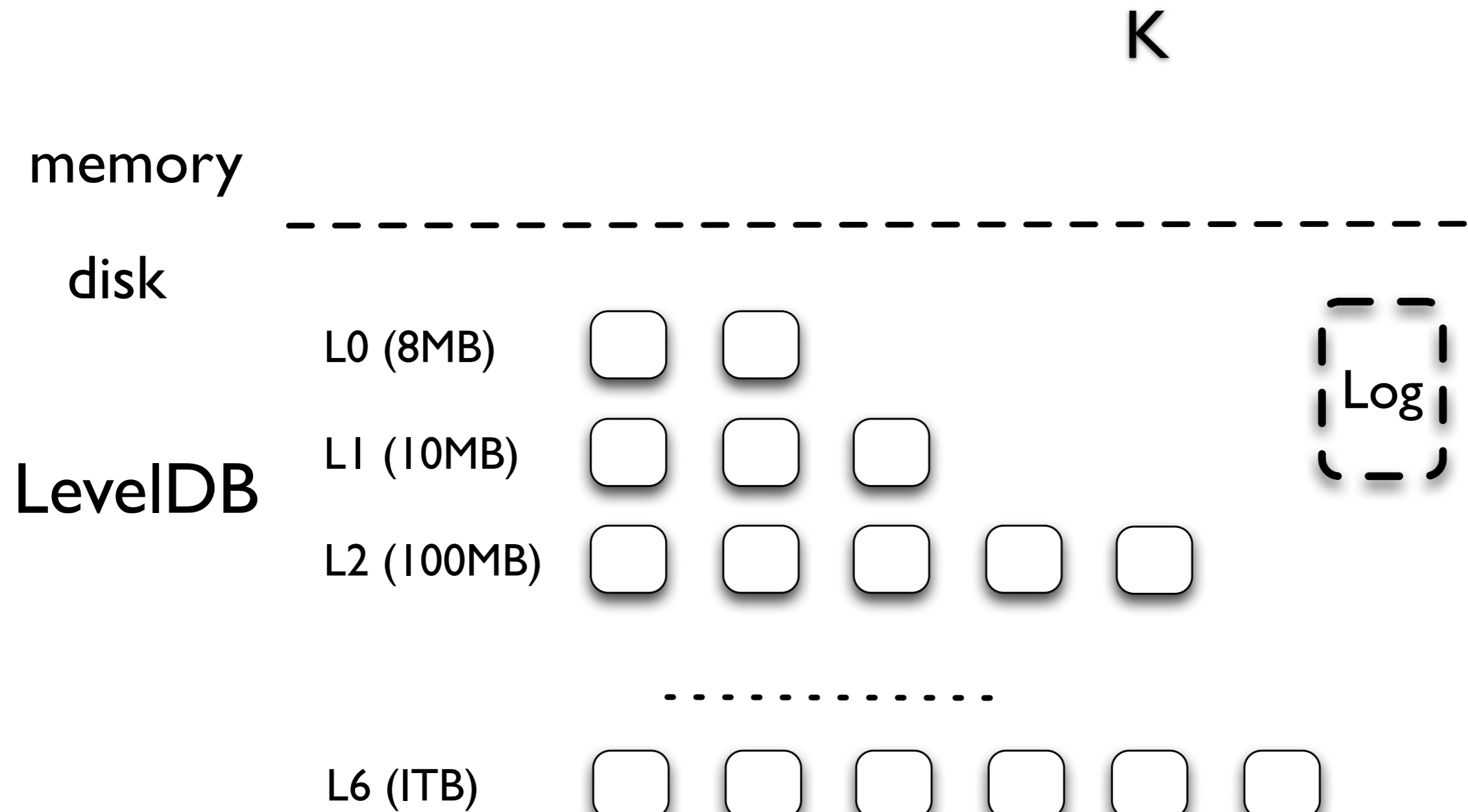
memory

disk

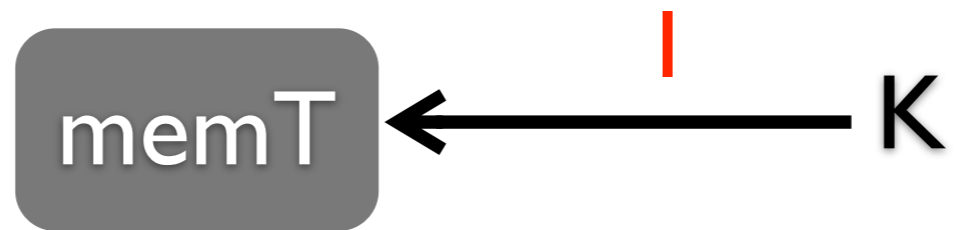
LevelDB



# LSM-trees: Lookup



# LSM-trees: Lookup

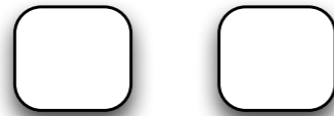


memory

disk

LevelDB

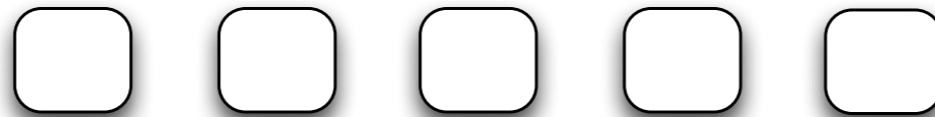
L0 (8MB)



L1 (10MB)



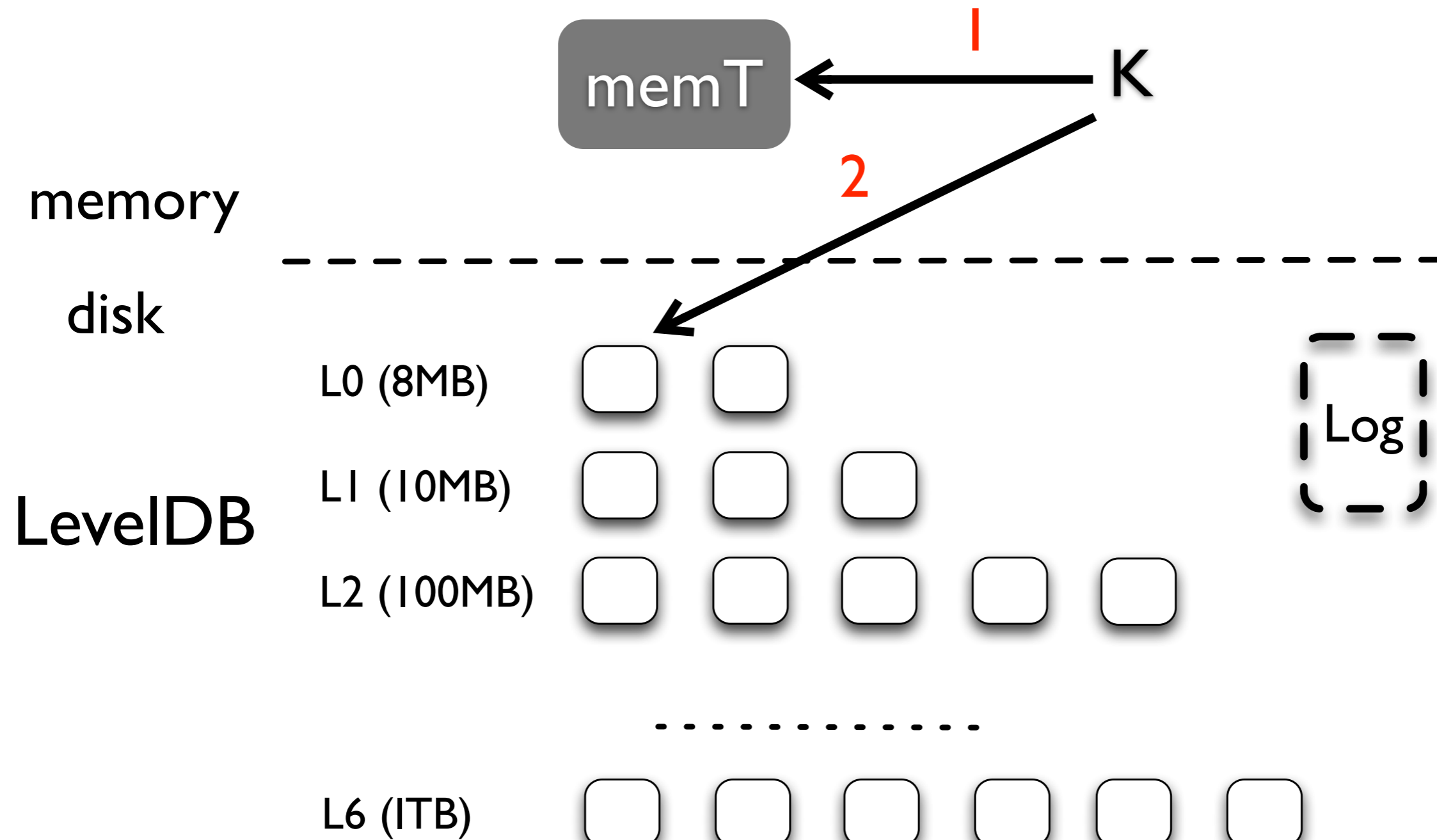
L2 (100MB)



L6 (1TB)

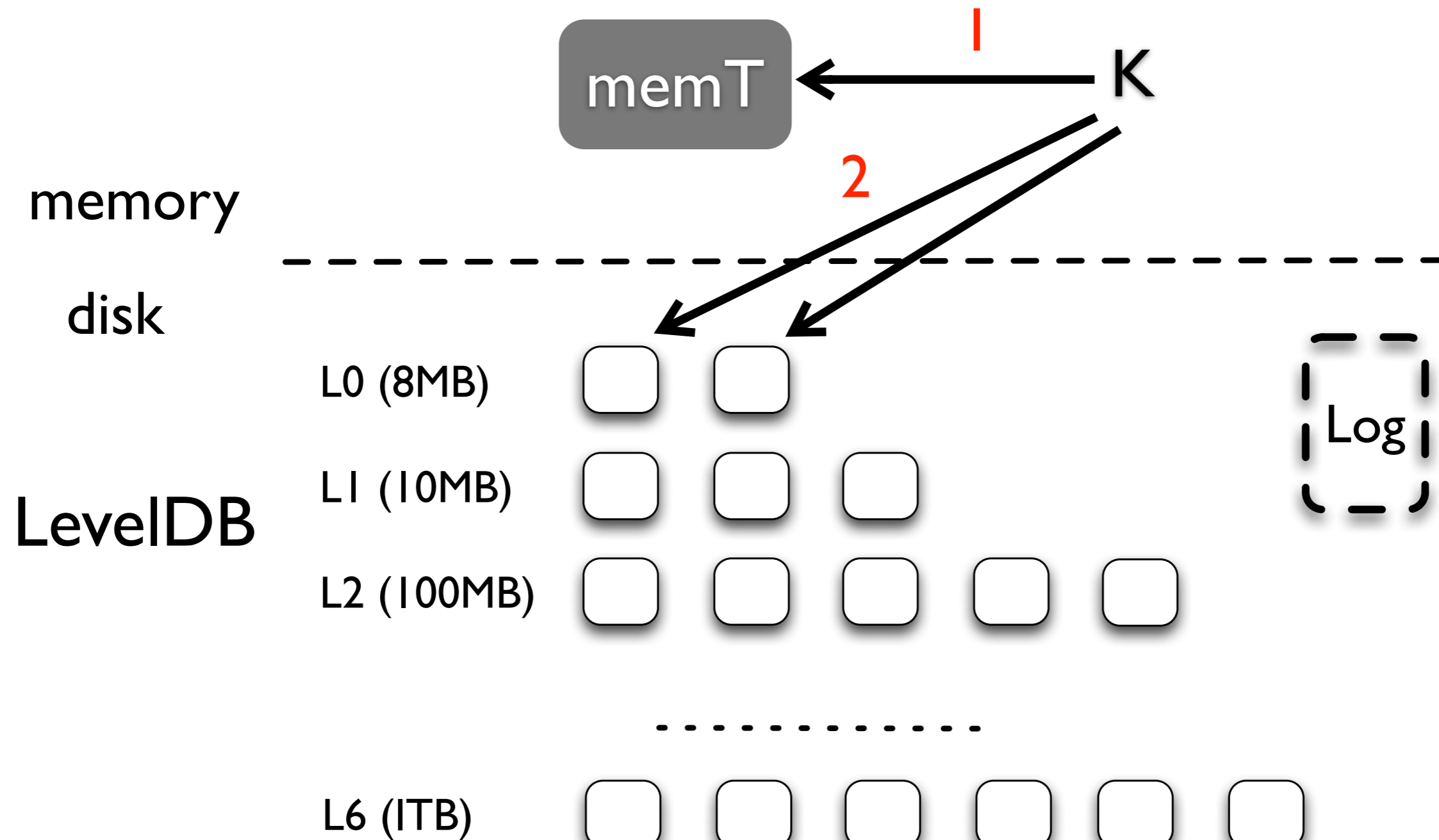


# LSM-trees: Lookup

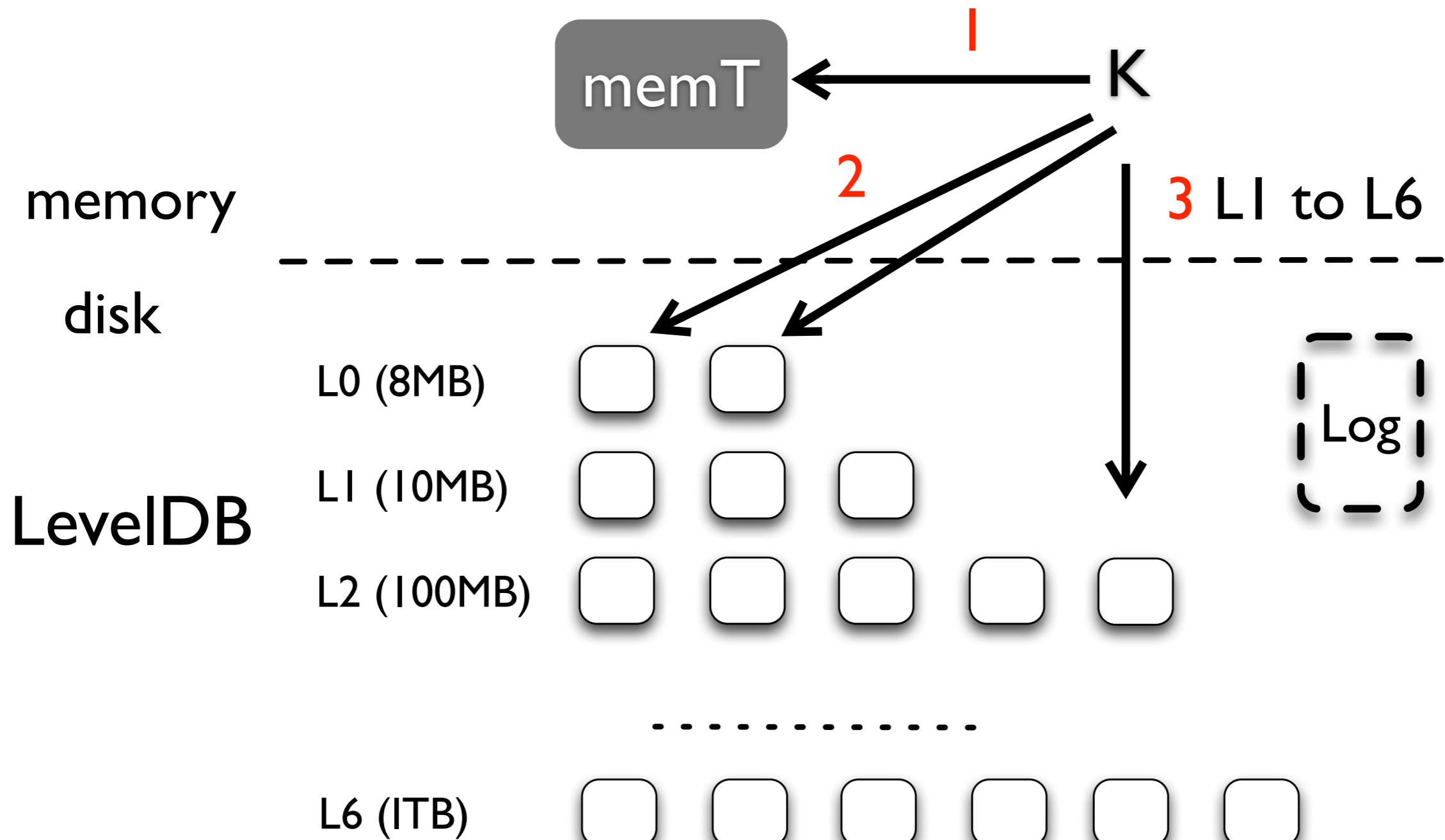




# LSM-trees: Lookup

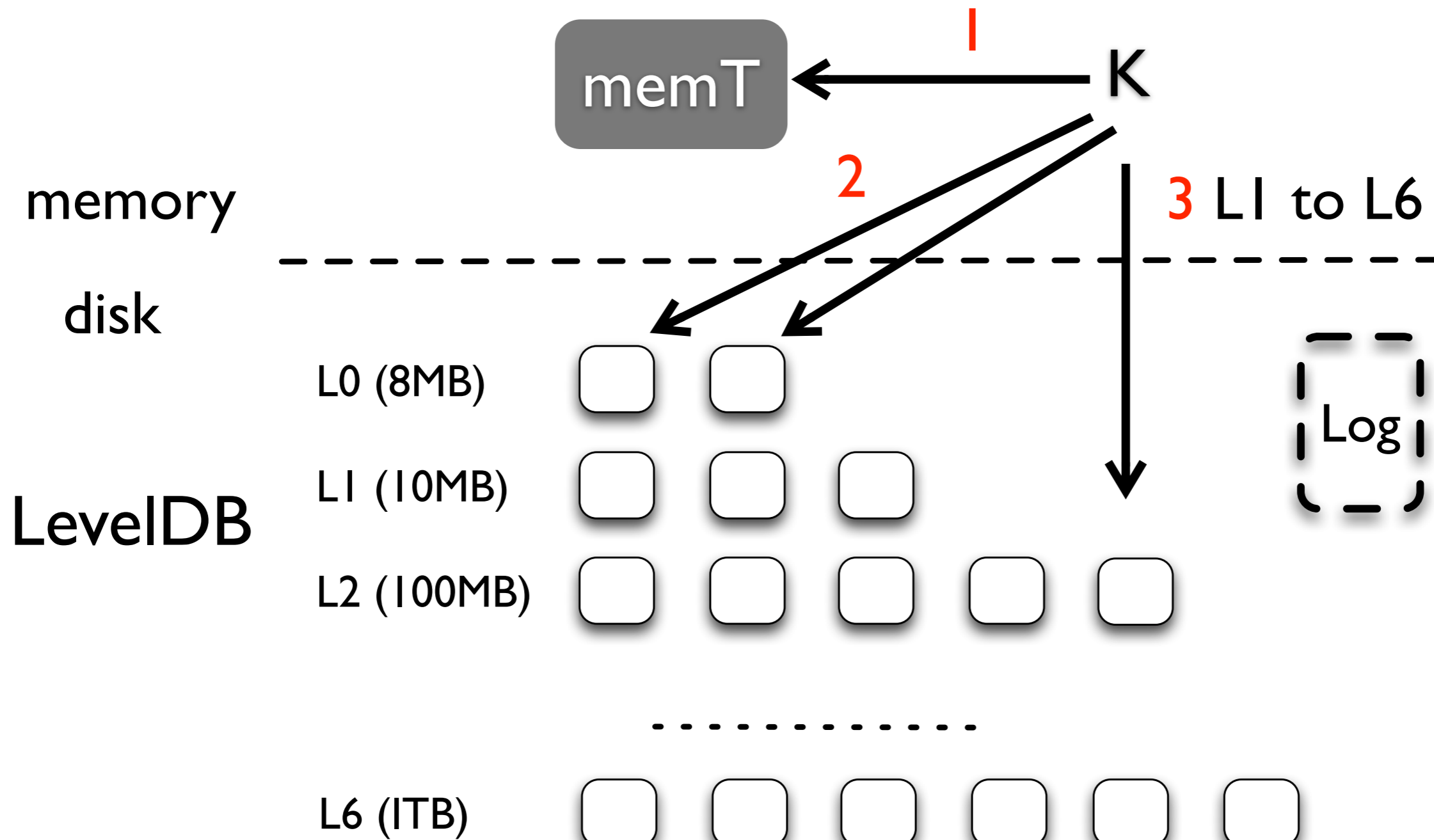


# LSM-trees: Lookup



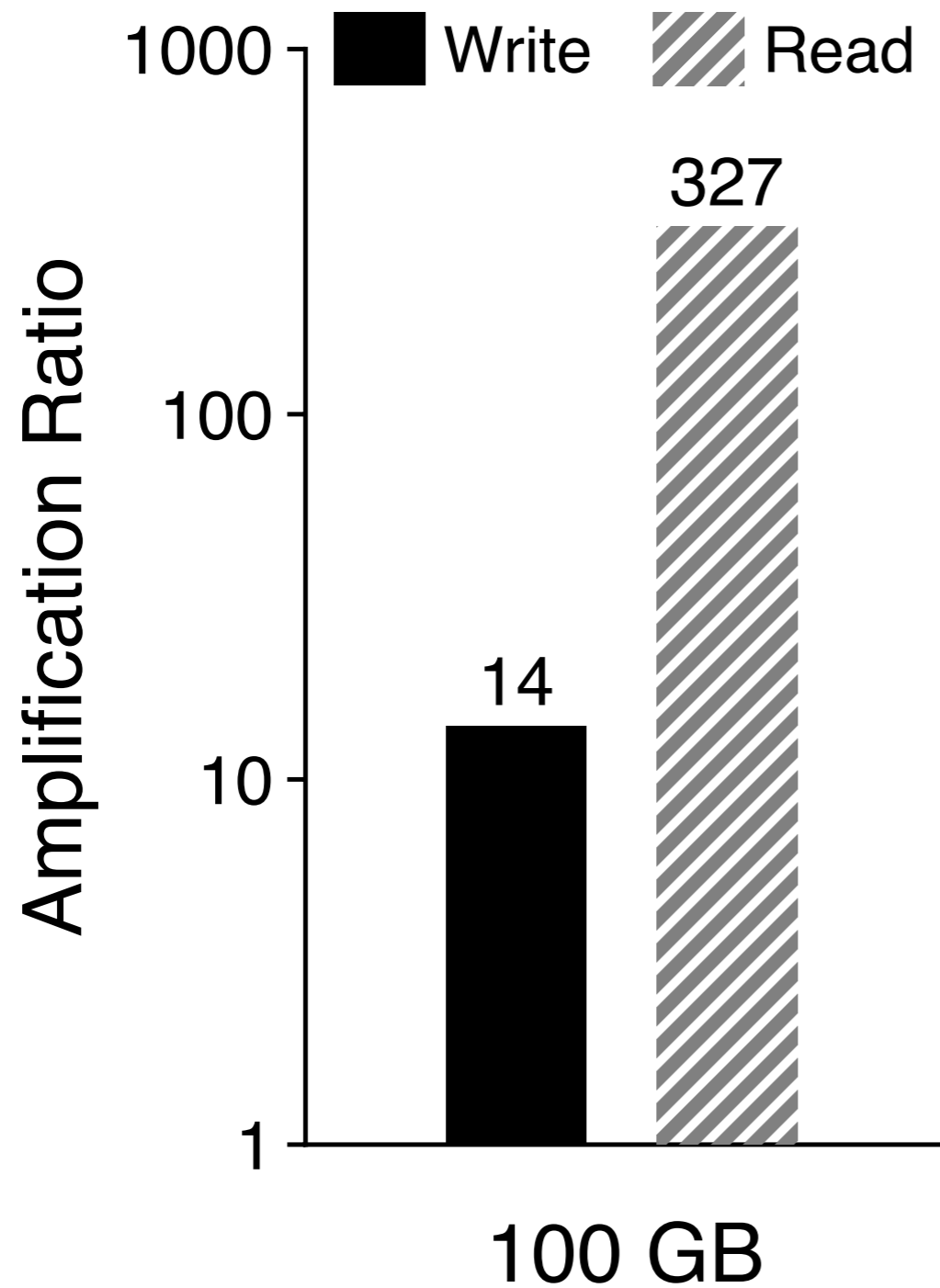
# LSM-trees: Lookup

1. Random reads
2. Travel many levels for a large LSM-tree

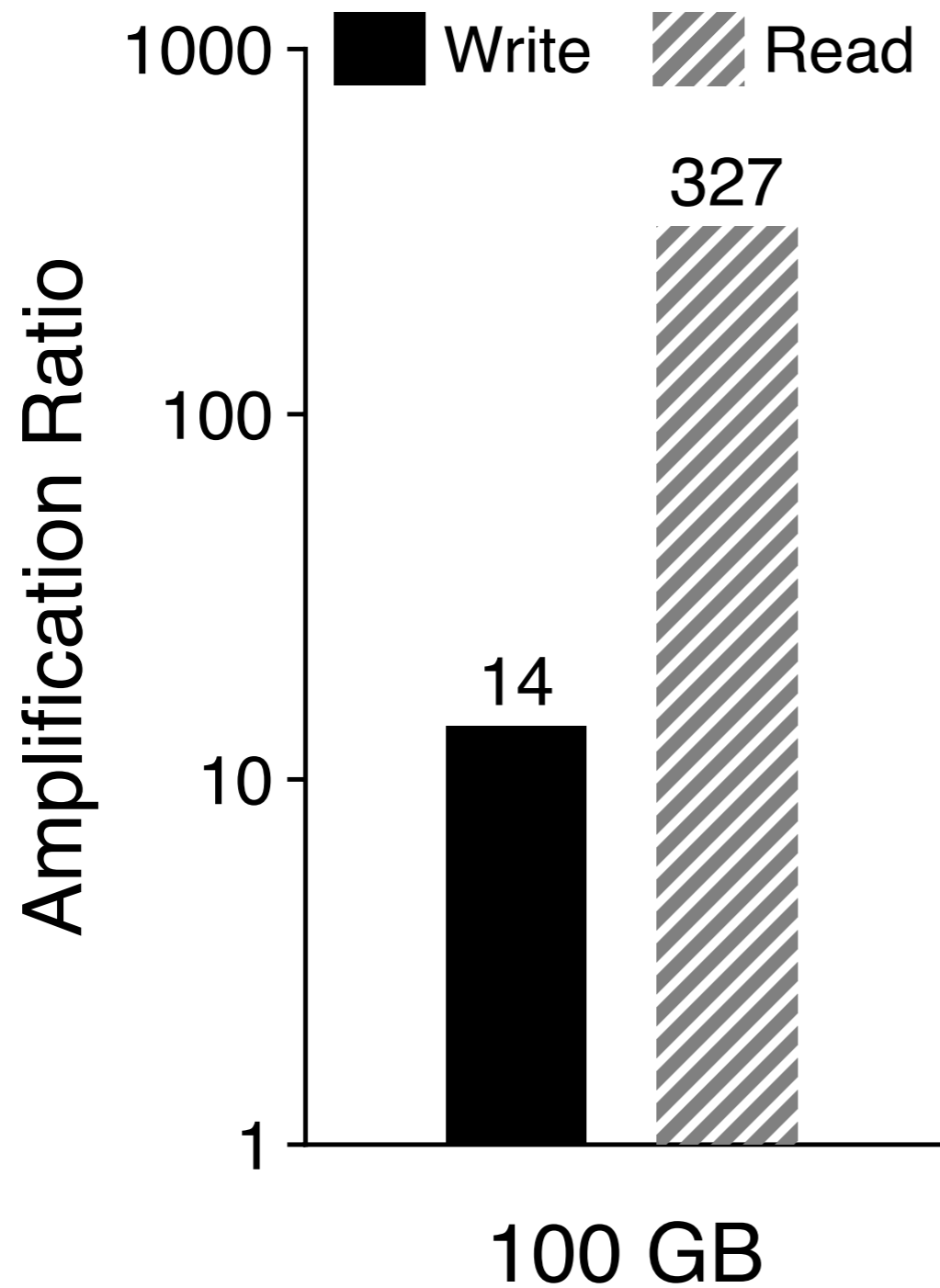


# I/O Amplification in LSM-trees

# I/O Amplification in LSM-trees



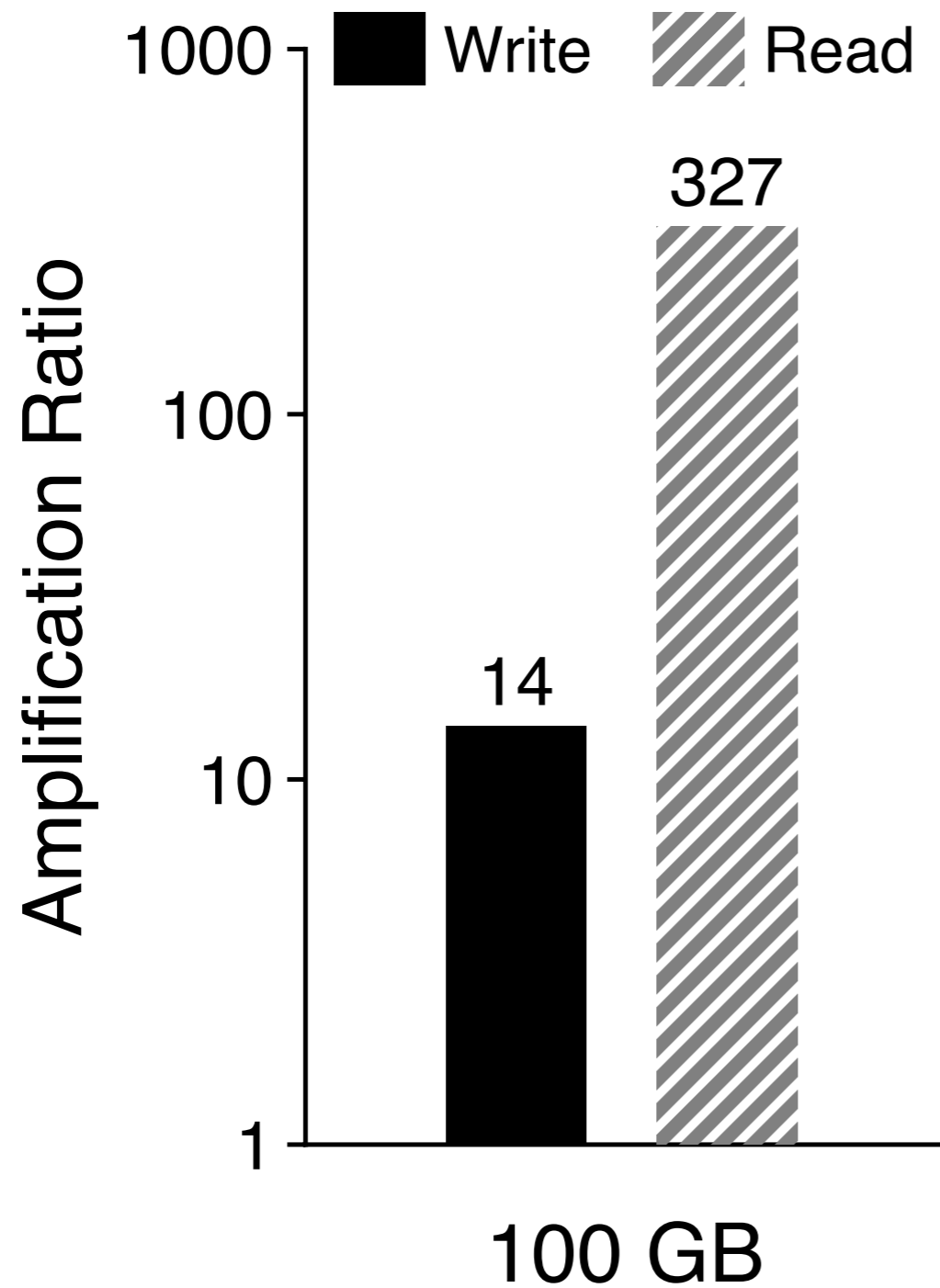
# I/O Amplification in LSM-trees



Random load:  
a 100GB database

Random lookup:  
100,000 lookups

# I/O Amplification in LSM-trees



Random load:  
a 100GB database

Random lookup:  
100,000 lookups

**Problems:**

large write amplification

large read amplification

Background

**Key-Value Separation**

**Challenges and Optimizations**

**Evaluation**

**Conclusion**



# Key-Value Separation

# Key-Value Separation

Main idea: only keys are required to be sorted

# Key-Value Separation

Main idea: only keys are required to be sorted

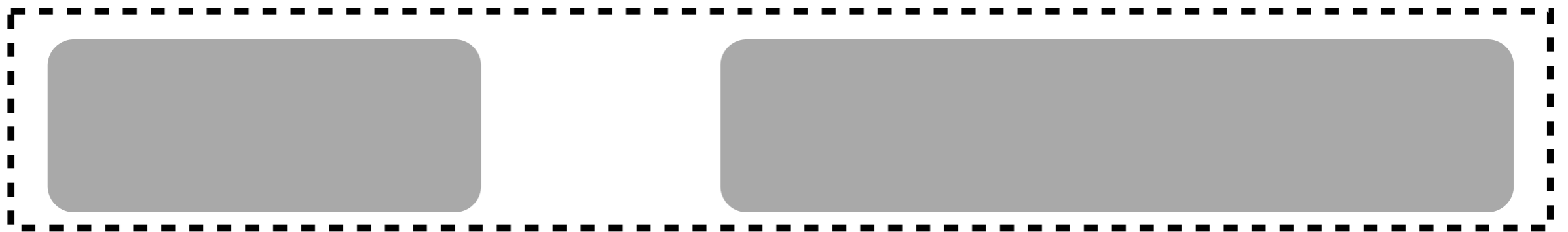
Decouple sorting and garbage collection

# Key-Value Separation

Main idea: only keys are required to be sorted

Decouple sorting and garbage collection

SSD device



LSM-tree

Value Log

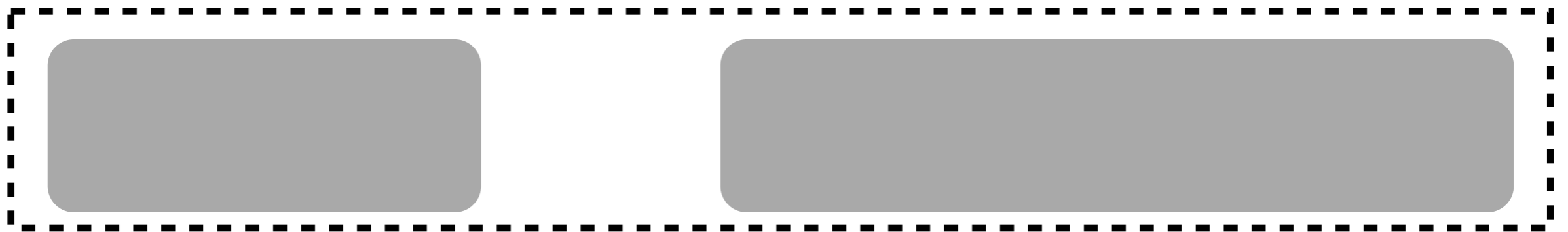
# Key-Value Separation

Main idea: only keys are required to be sorted

Decouple sorting and garbage collection



SSD device



LSM-tree

Value Log

# Key-Value Separation

Main idea: only keys are required to be sorted

Decouple sorting and garbage collection



SSD device



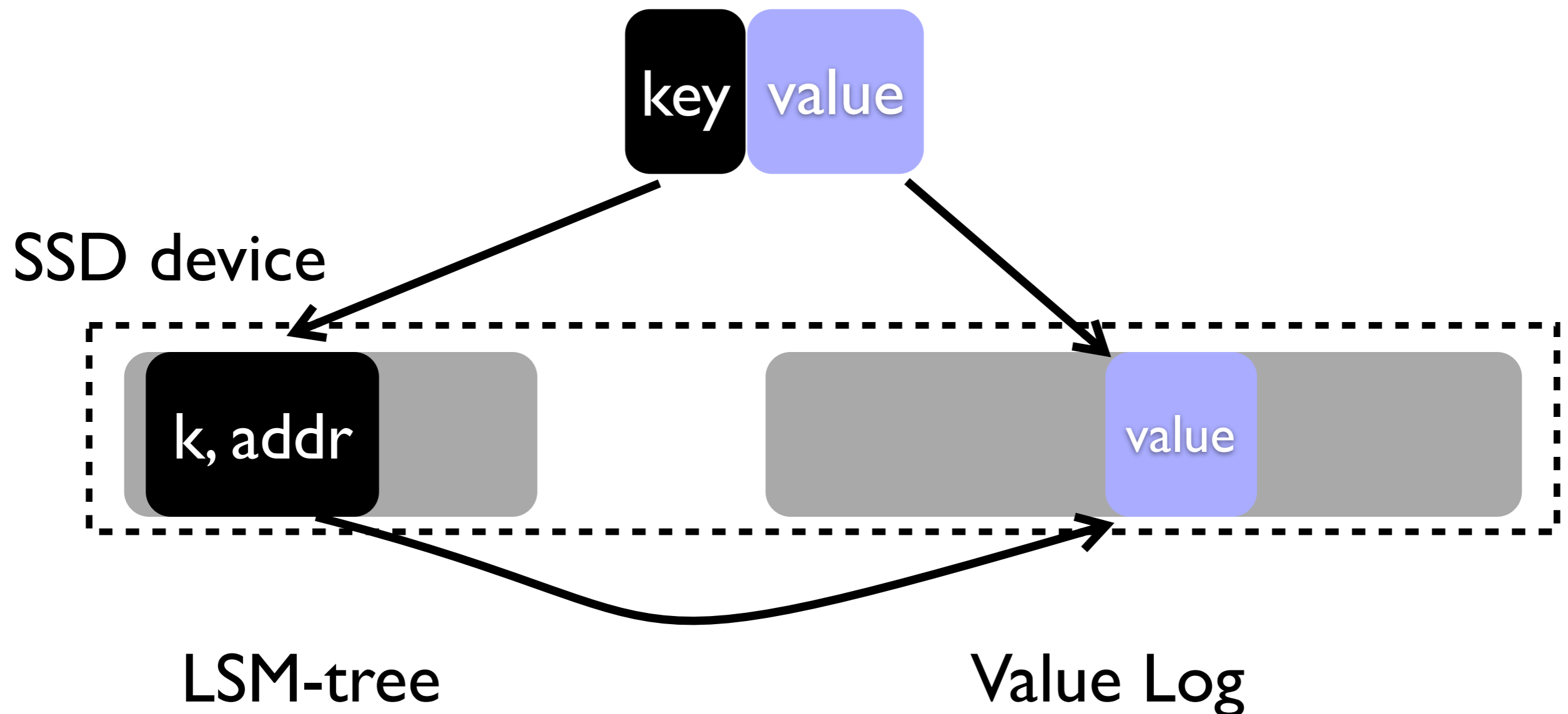
LSM-tree

Value Log

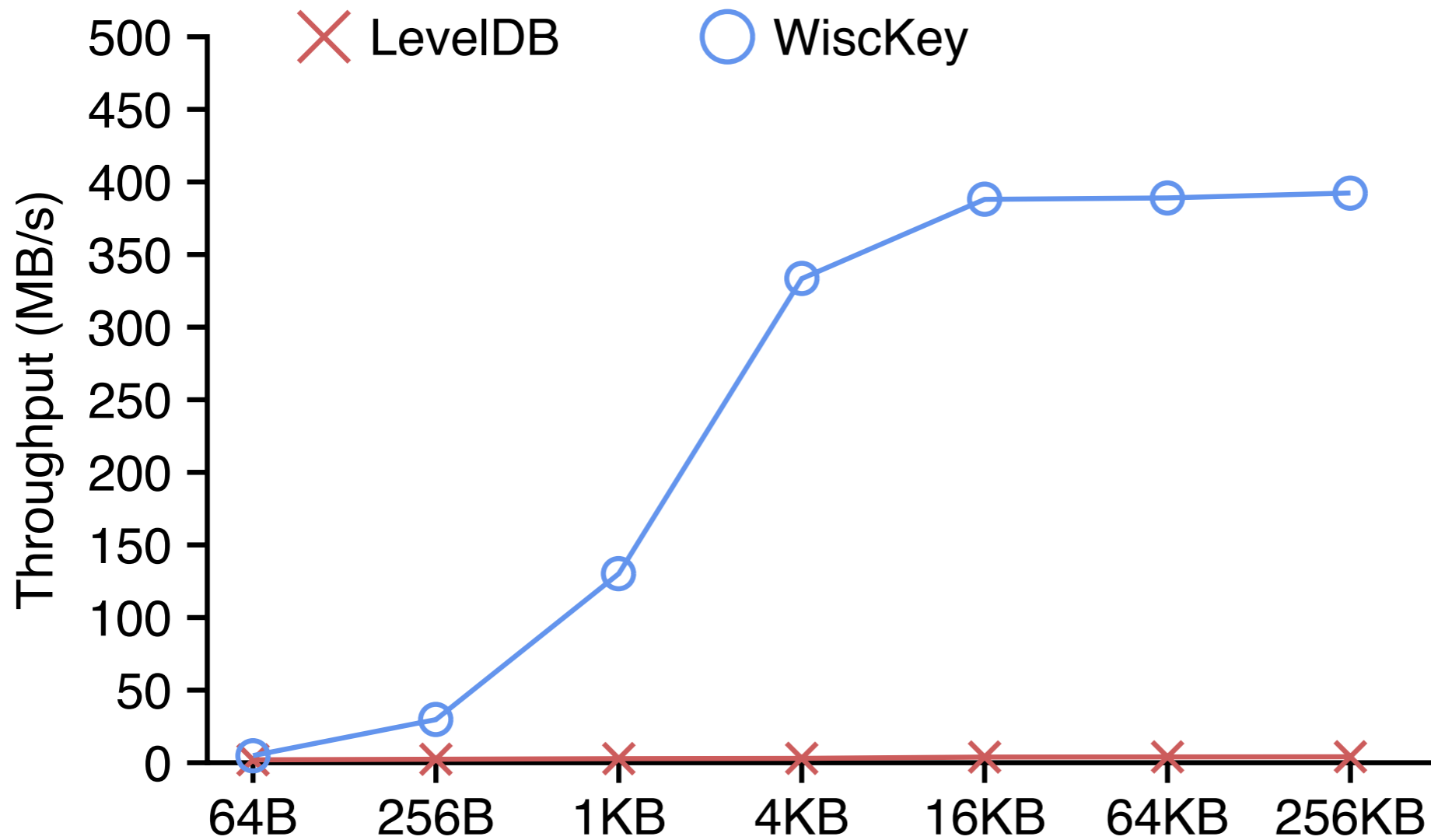
# Key-Value Separation

Main idea: only keys are required to be sorted

Decouple sorting and garbage collection



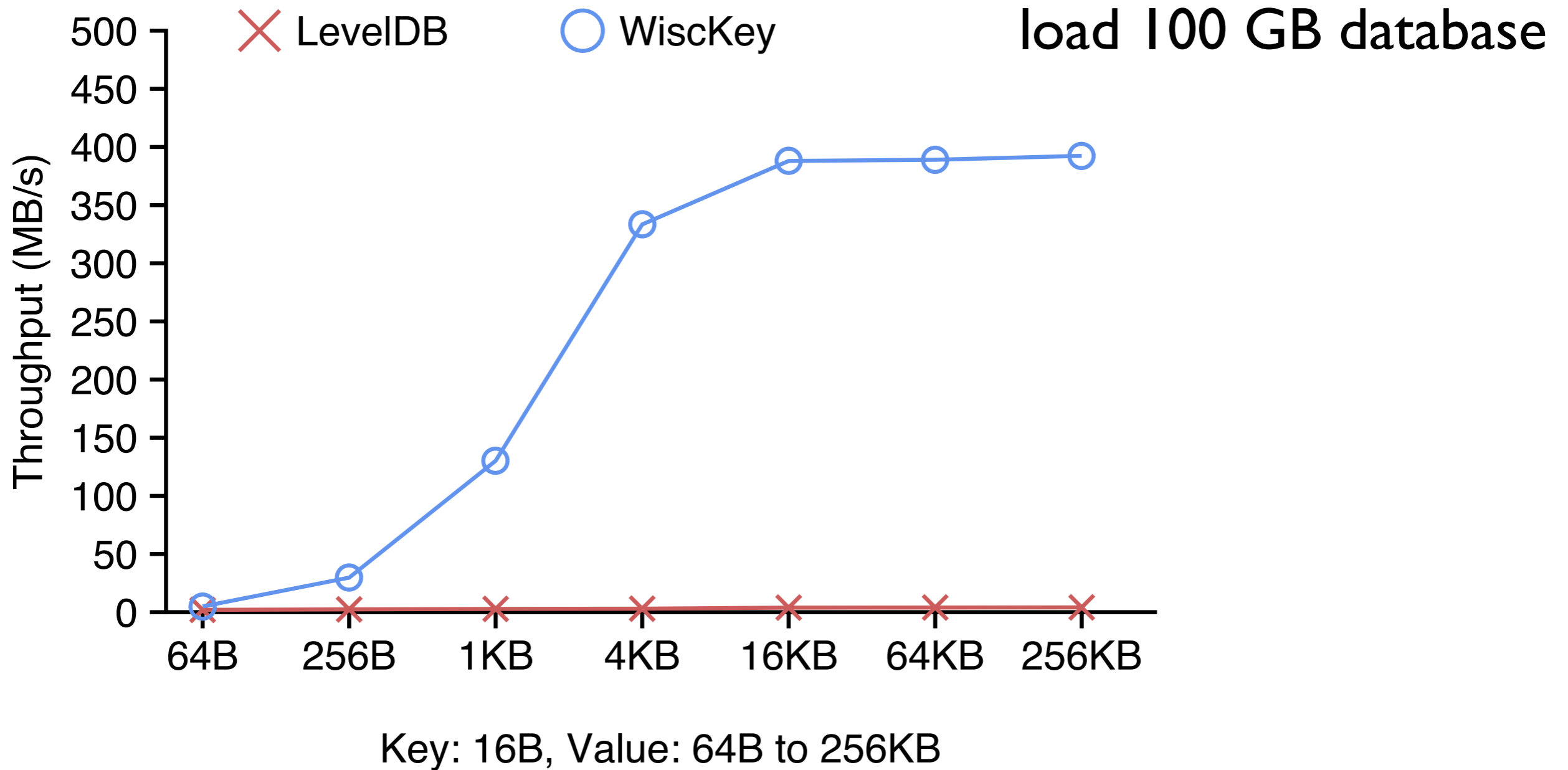
# Random Load



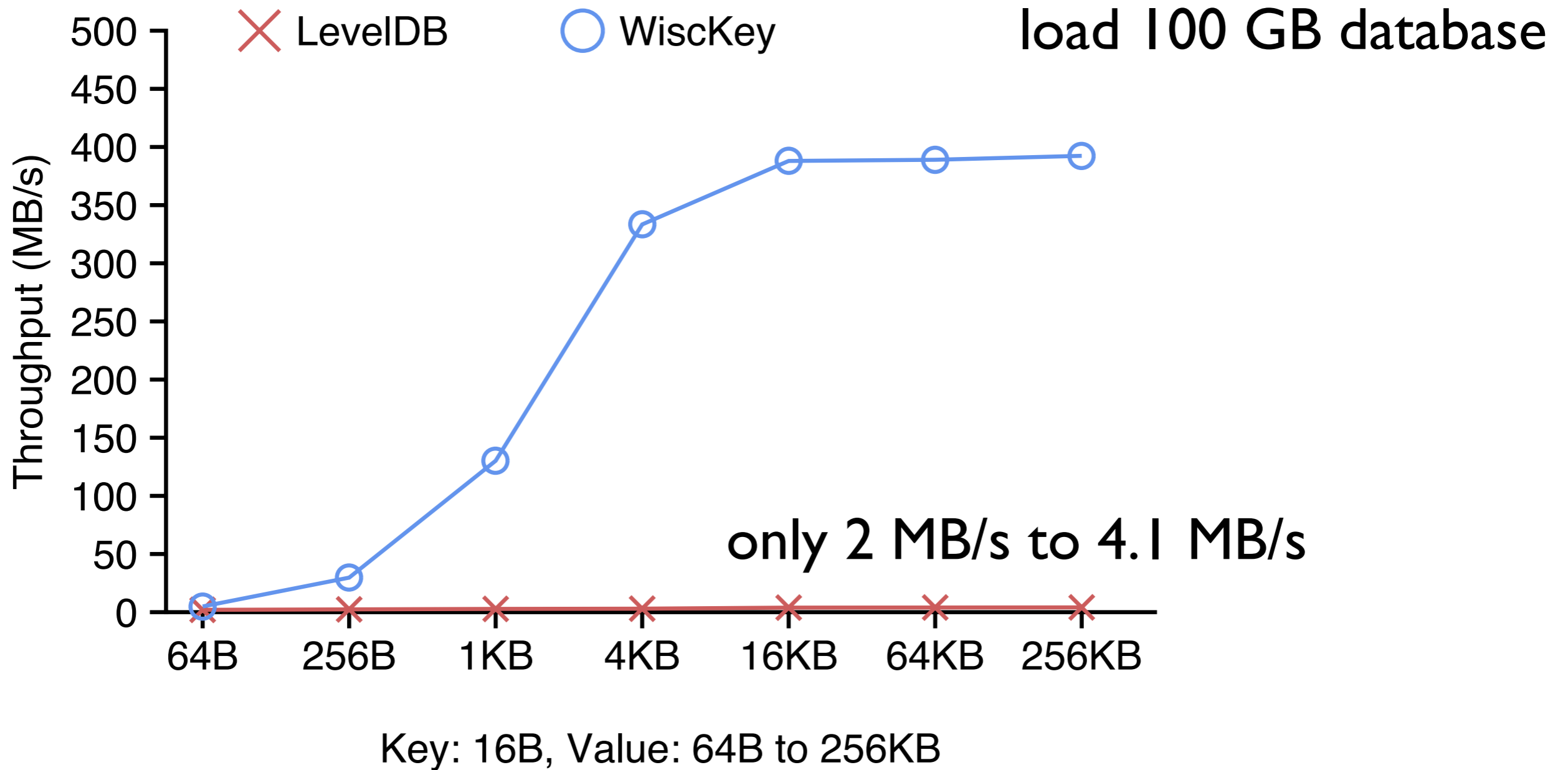
Key: 16B, Value: 64B to 256KB



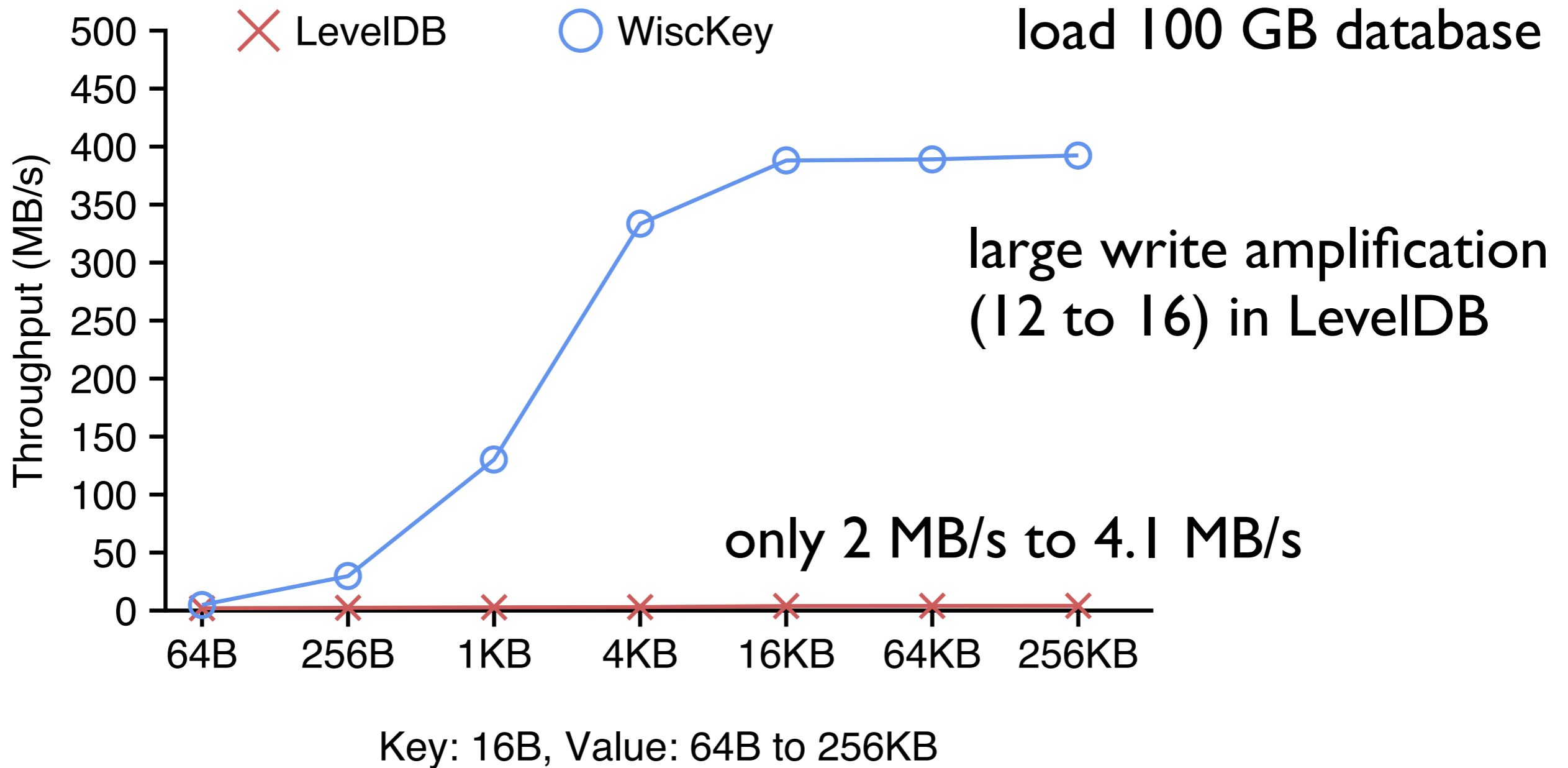
# Random Load



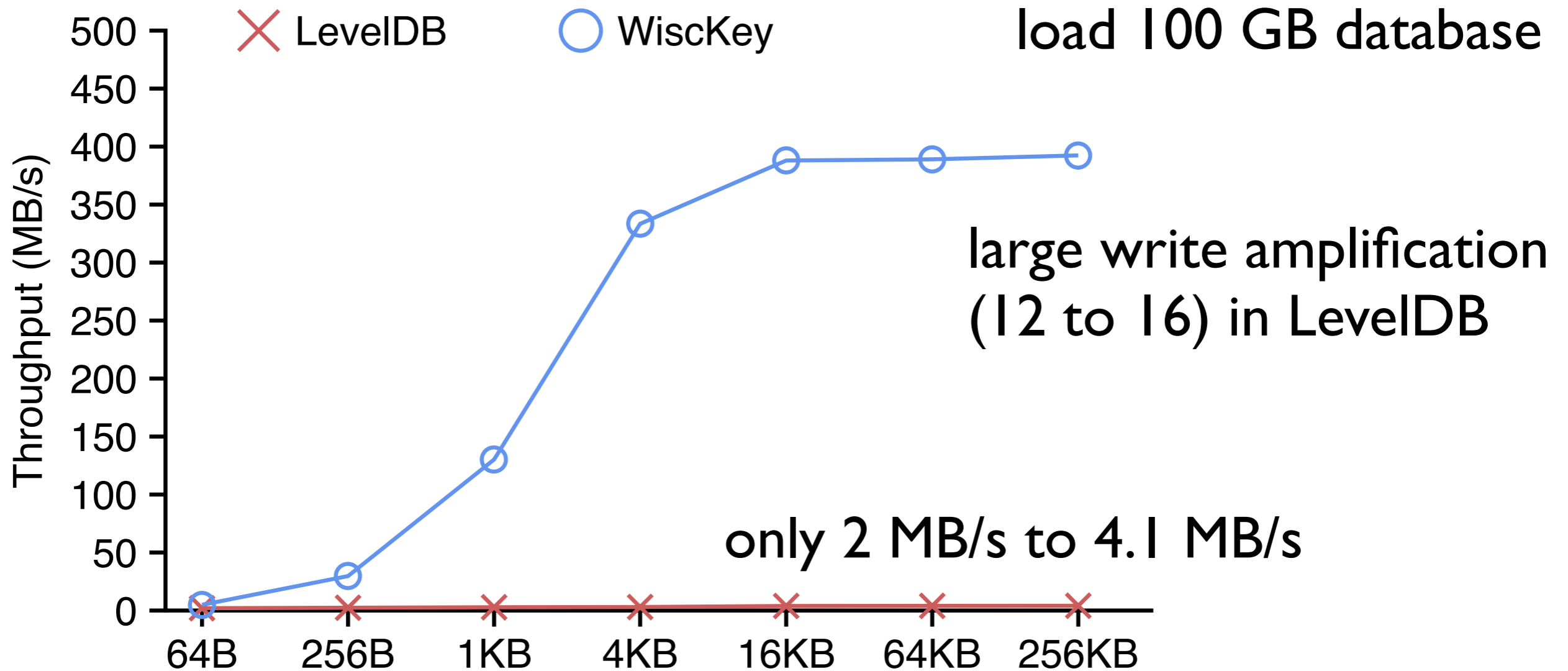
# Random Load



# Random Load



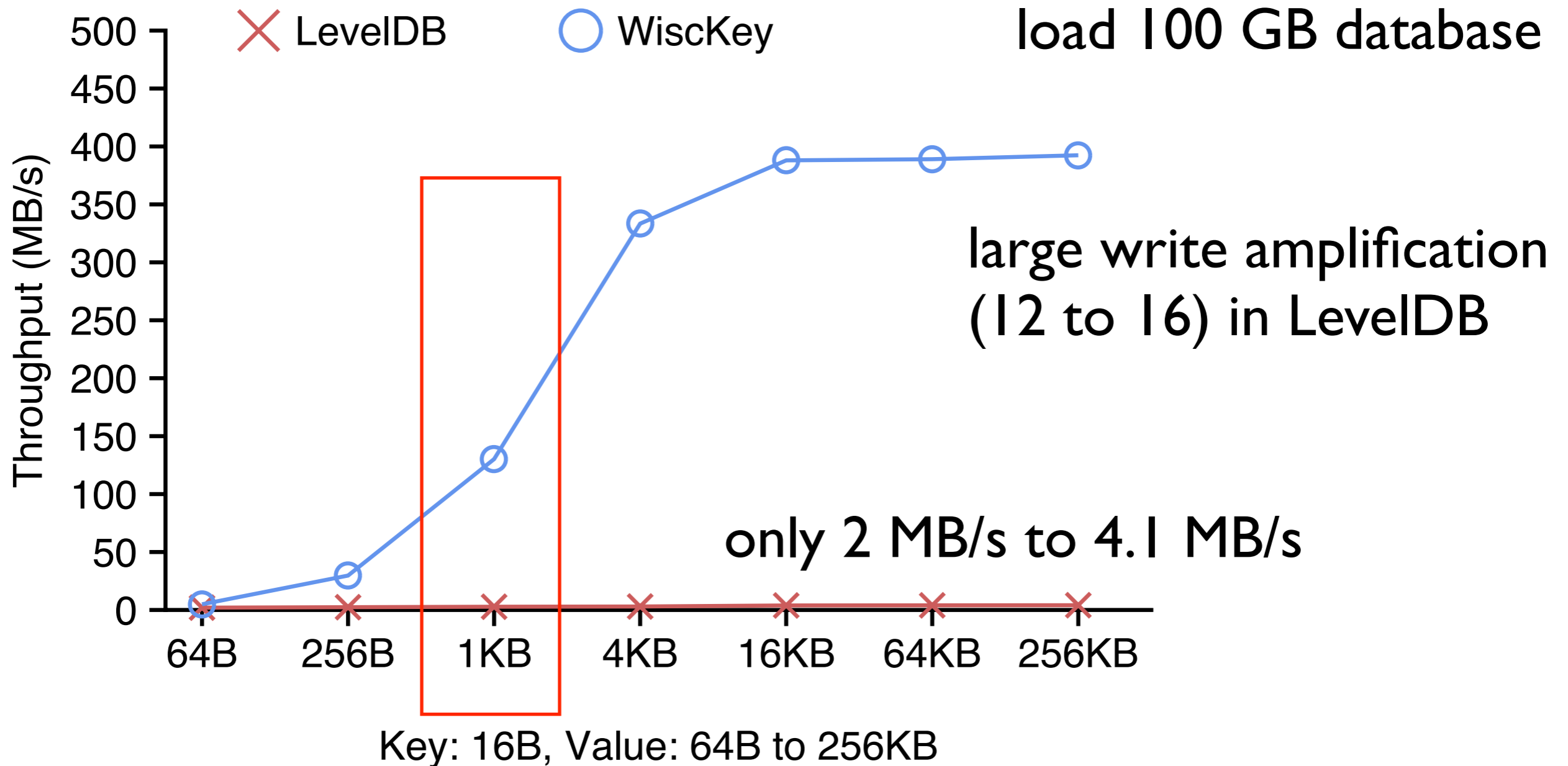
# Random Load



Key: 16B, Value: 64B to 256KB

Small write amplification in WiscKey due to key-value separation (up to **111x** in throughput)

# Random Load



Small write amplification in WiscKey due to key-value separation (up to **111x** in throughput)

# LevelDB

**limits of files**

**num of files**

**L0**

**9**

**L1 (5)**

**30**

**L2 (50)**

**365**

**L3 (500)**

**2184**

**L4 (5000)**

**15752**

**L5 (50000)**

**23733**

**L6 (500000)**

**0**

# LevelDB

## limits of files

## num of files

L0	9	Large LSM-tree:
L1 (5)	30	Intensive compaction
L2 (50)	365	→ repeated reads/writes
L3 (500)	2184	→ stall foreground I/Os
L4 (5000)	15752	Many levels
L5 (50000)	23733	→ travel several levels for each lookup
L6 (500000)	0	

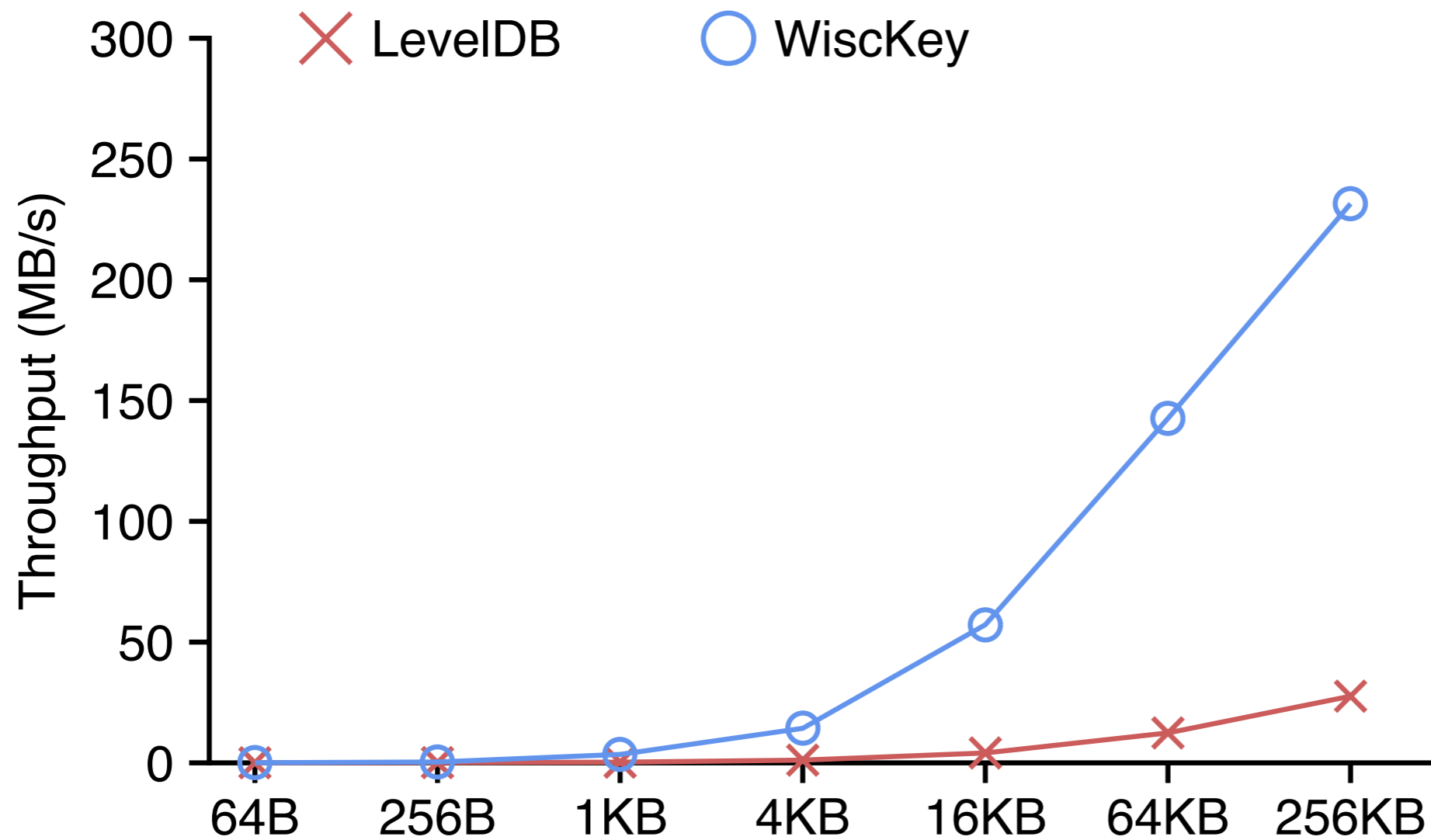
<b>limits of files</b>	<b>LevelDB num of files</b>	<b>WiscKey num of files</b>
<b>L0</b>	<b>9</b>	<b>7</b>
<b>L1 (5)</b>	<b>30</b>	<b>11</b>
<b>L2 (50)</b>	<b>365</b>	<b>127</b>
<b>L3 (500)</b>	<b>2184</b>	<b>460</b>
<b>L4 (5000)</b>	<b>15752</b>	<b>0</b>
<b>L5 (50000)</b>	<b>23733</b>	<b>0</b>
<b>L6 (500000)</b>	<b>0</b>	<b>0</b>



	LevelDB	WiscKey
limits of files	num of files	num of files
L0	9	7
L1 (5)	30	11
L2 (50)	365	127
L3 (500)	2184	460
L4 (5000)	15752	0
L5 (50000)	23733	0
L6 (500000)	0	0

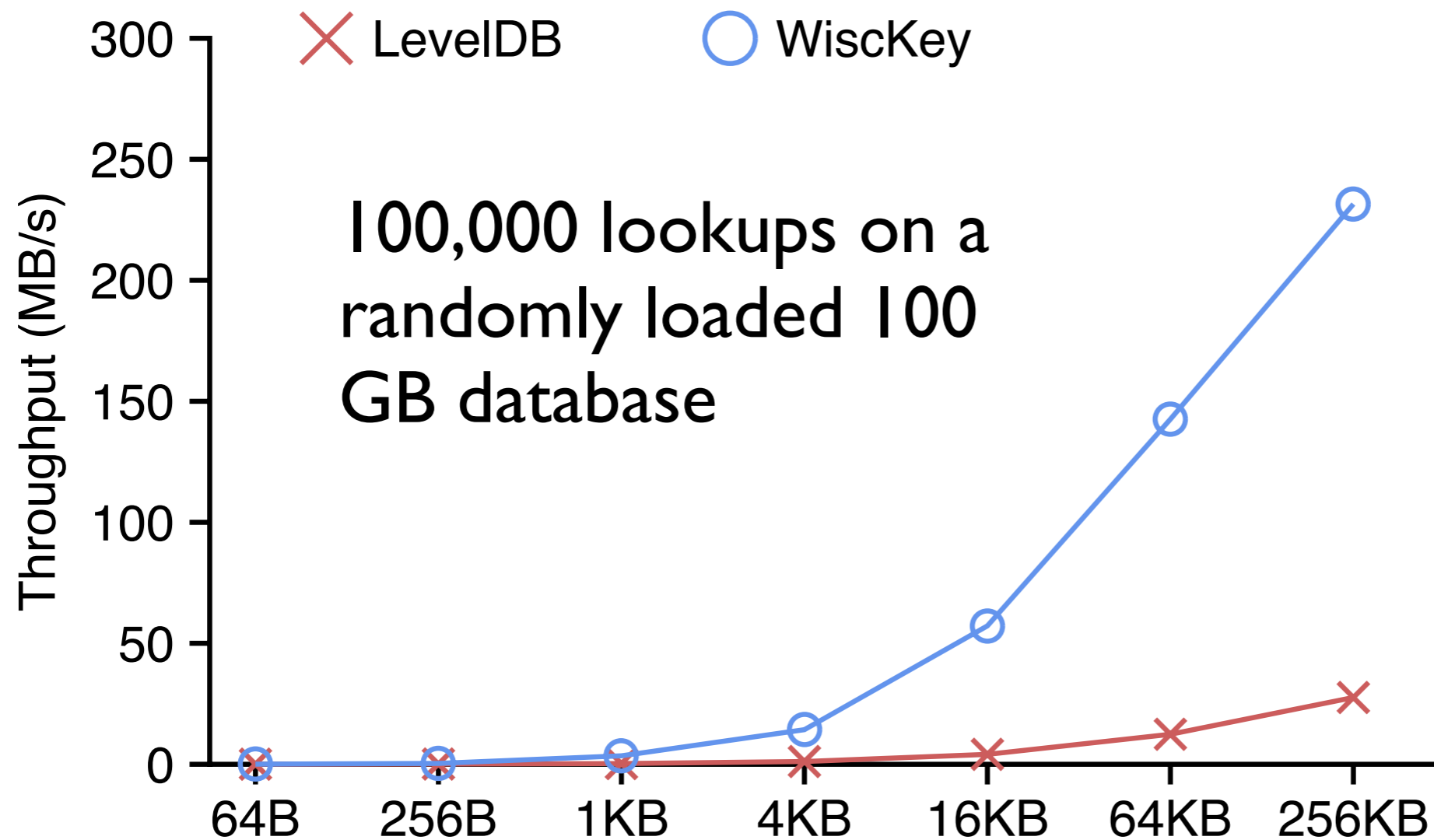
**Small LSM-tree:** less compaction, fewer levels to search, and better caching

# Random Lookup



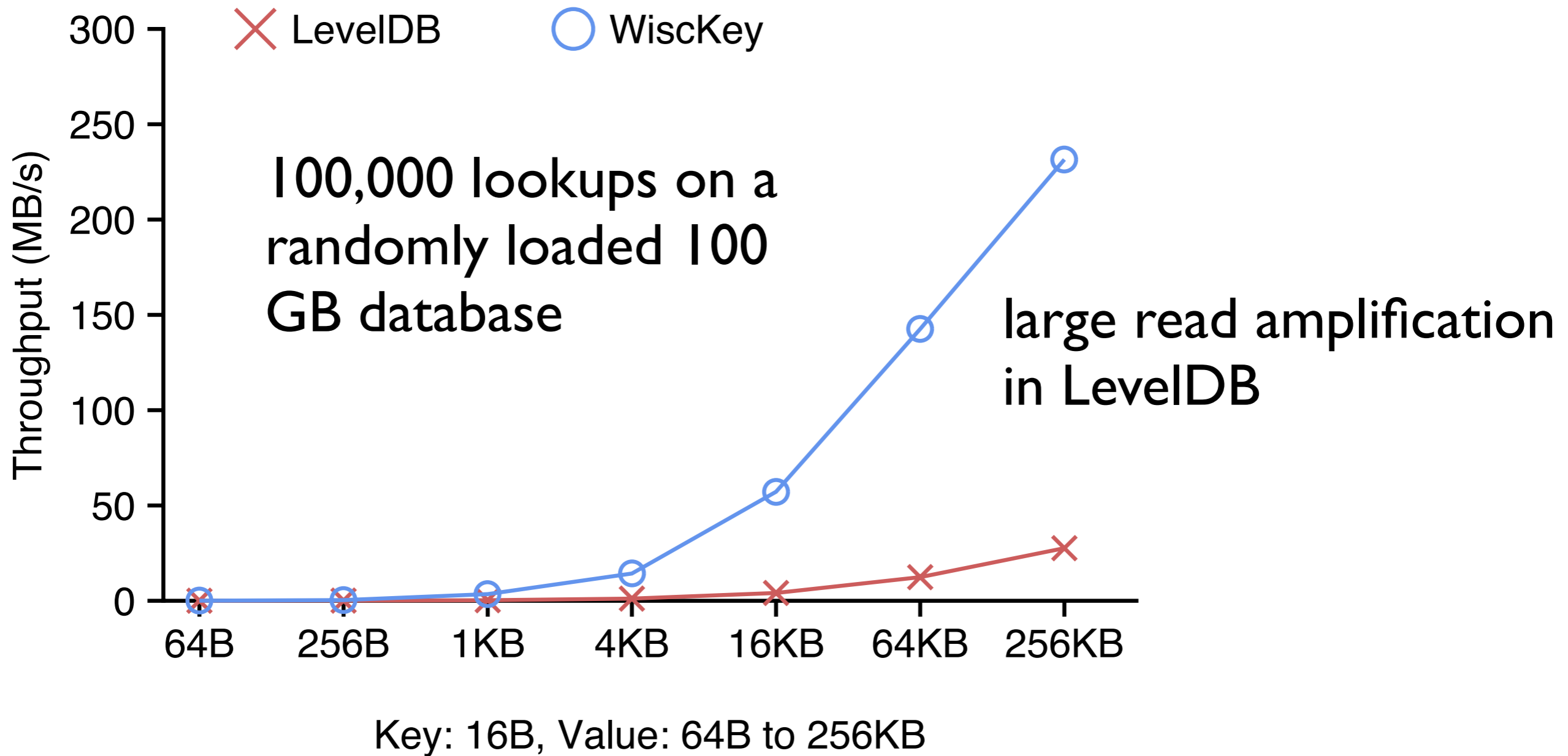
Key: 16B, Value: 64B to 256KB

# Random Lookup

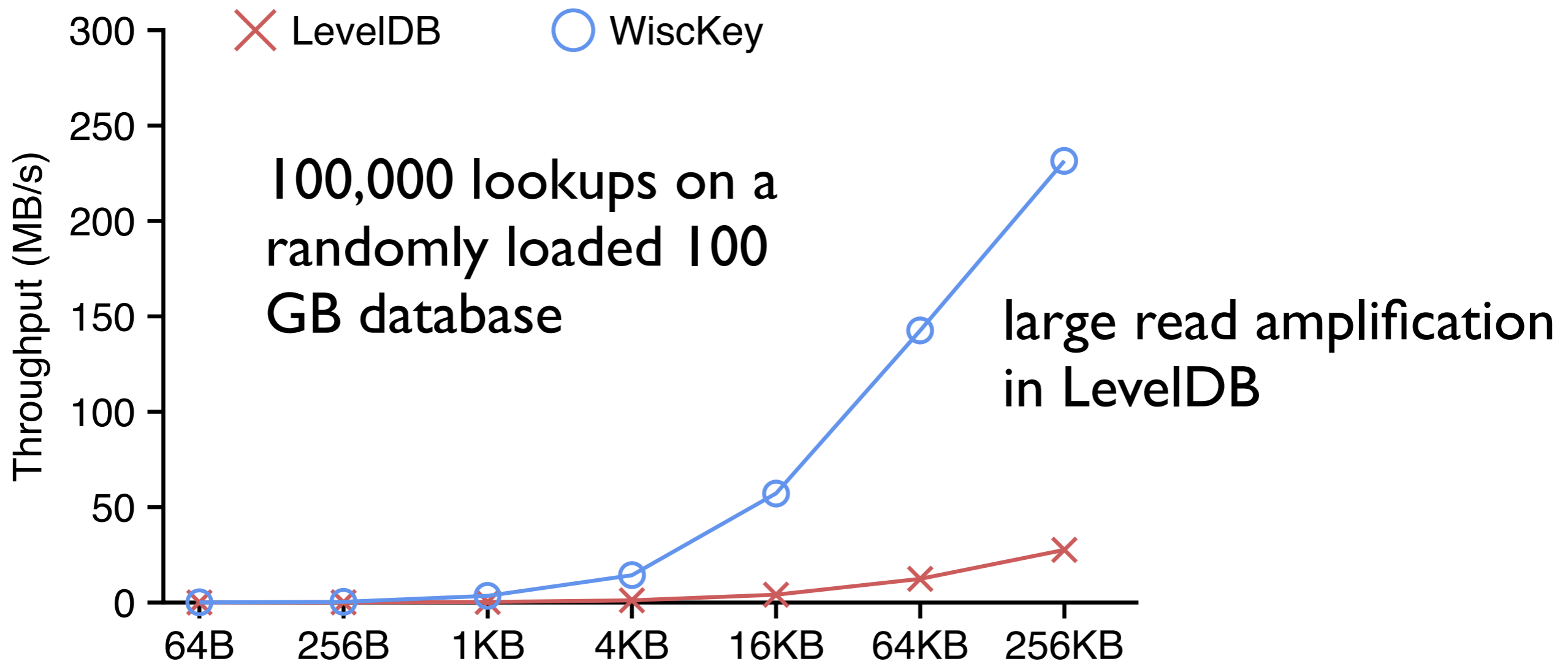


Key: 16B, Value: 64B to 256KB

# Random Lookup



# Random Lookup



Key: 16B, Value: 64B to 256KB

Smaller LSM-tree in WiscKey leads to better lookup performance (1.6x - 14x)

Background

Key-Value Separation

**Challenges and Optimizations**

- Parallel range query
- Garbage collection
- LSM-log

**Evaluation**

**Conclusion**

# Parallel Range Query

# Parallel Range Query

SSD read performance

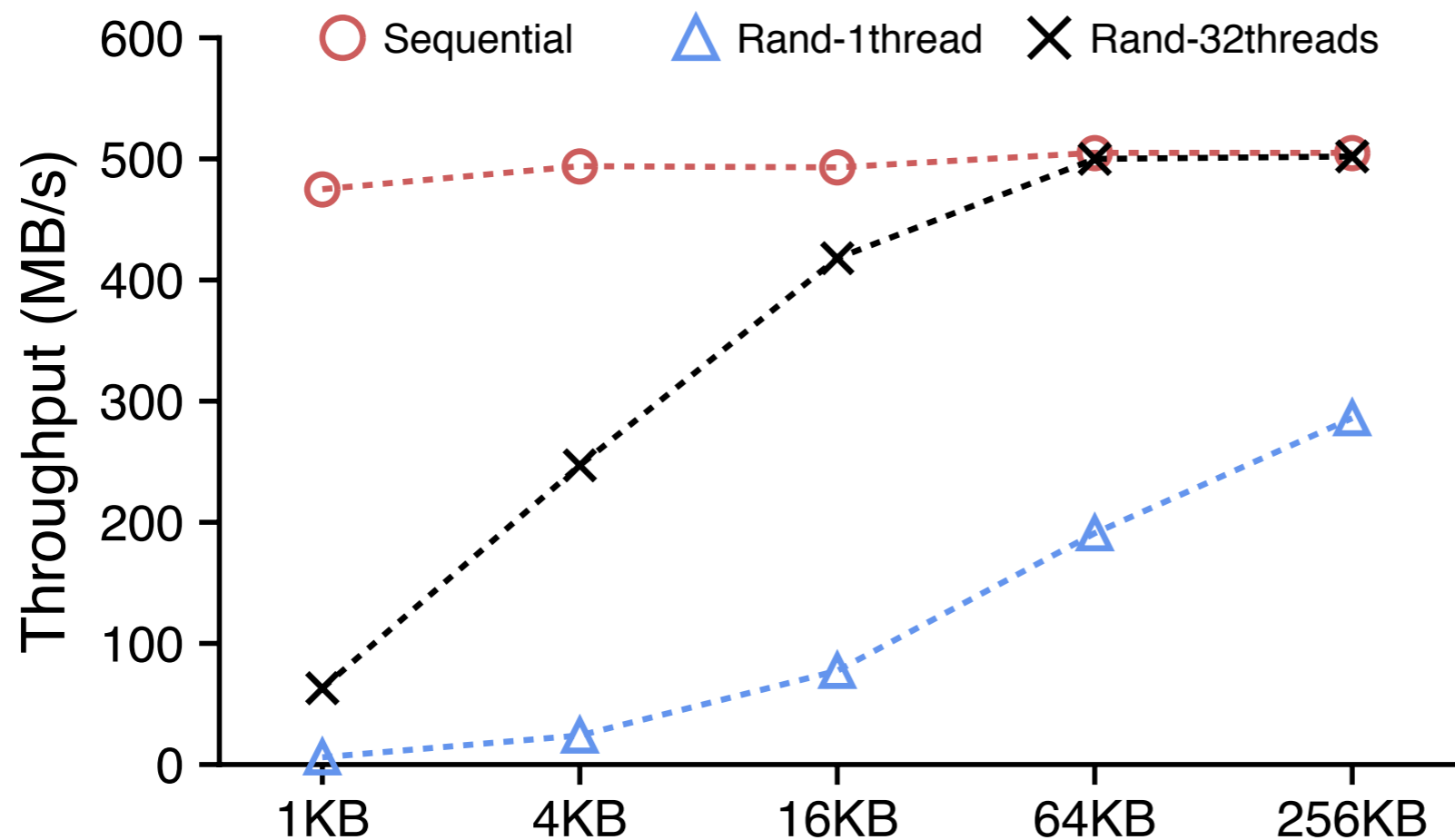
→ sequential, random, parallel



# Parallel Range Query

## SSD read performance

→ sequential, random, parallel



SSD: Samsung 840  
EVO 500GB

Reads on a 100GB  
file on ext4

Request size: 1KB to 256KB

# Parallel Range Query

# Parallel Range Query

## Challenge

- sequential reads in LevelDB
- read keys and values separately in WiscKey

# Parallel Range Query

## Challenge

- sequential reads in LevelDB
- read keys and values separately in WiscKey

## Parallel range query

- leverage parallel random reads of SSDs

# Parallel Range Query

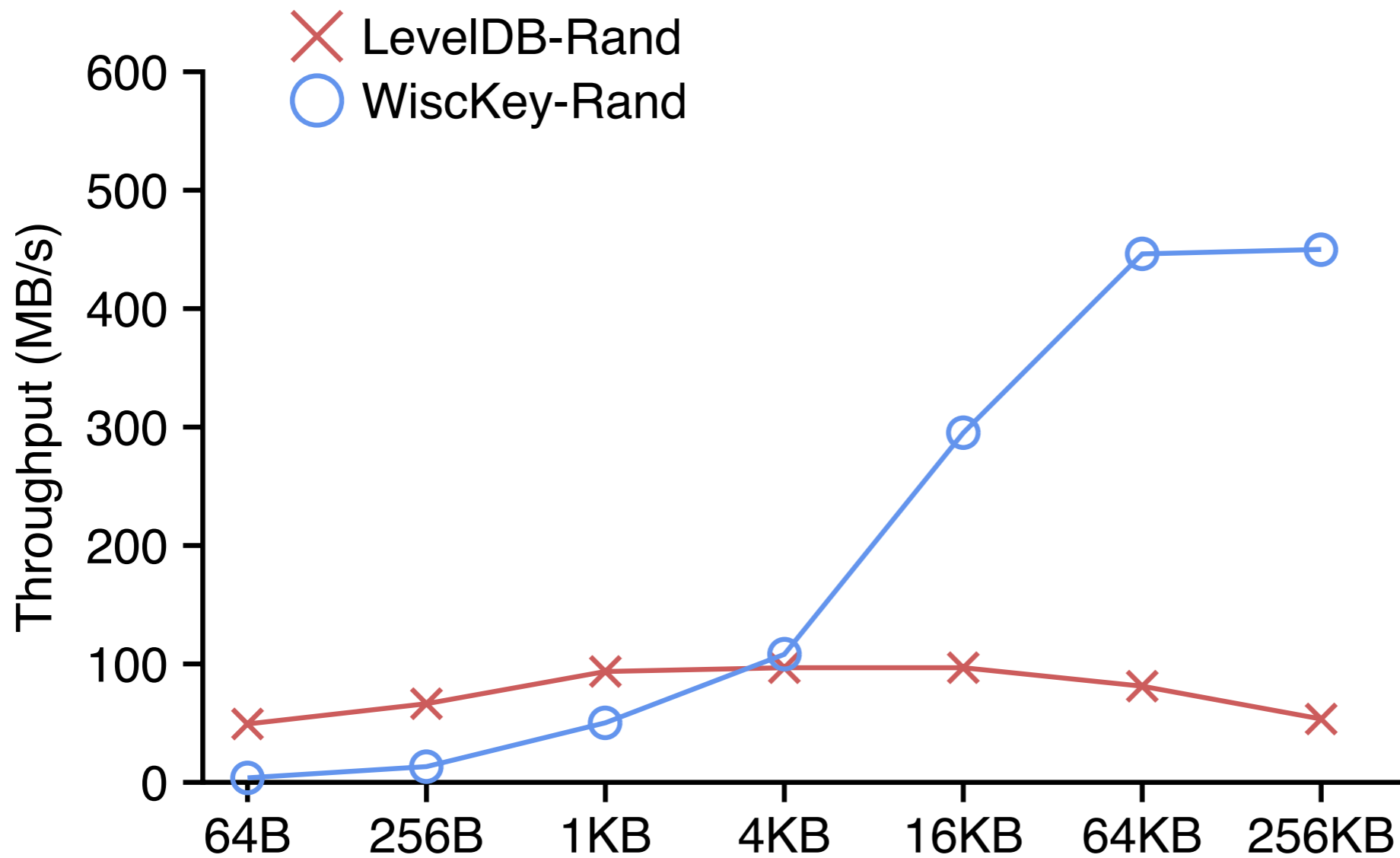
## Challenge

- sequential reads in LevelDB
- read keys and values separately in WiscKey

## Parallel range query

- leverage parallel random reads of SSDs
- prefetch key-value pairs in advance
  - range query interface: seek(), next(), prev()
  - detect a sequential pattern
  - prefetch concurrently in background

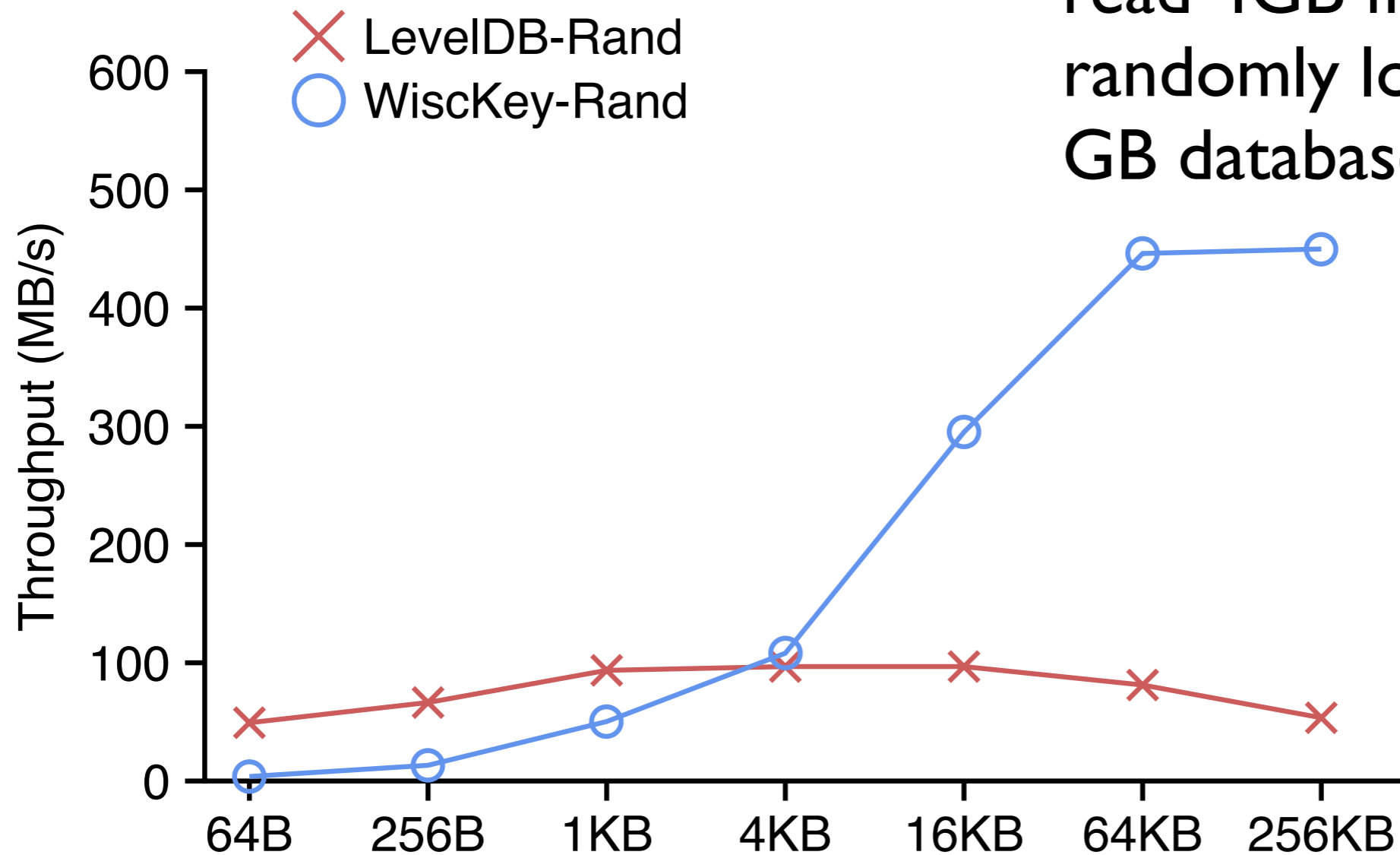
# Range Query



Key: 16B, Value: 64B to 256KB

# Range Query

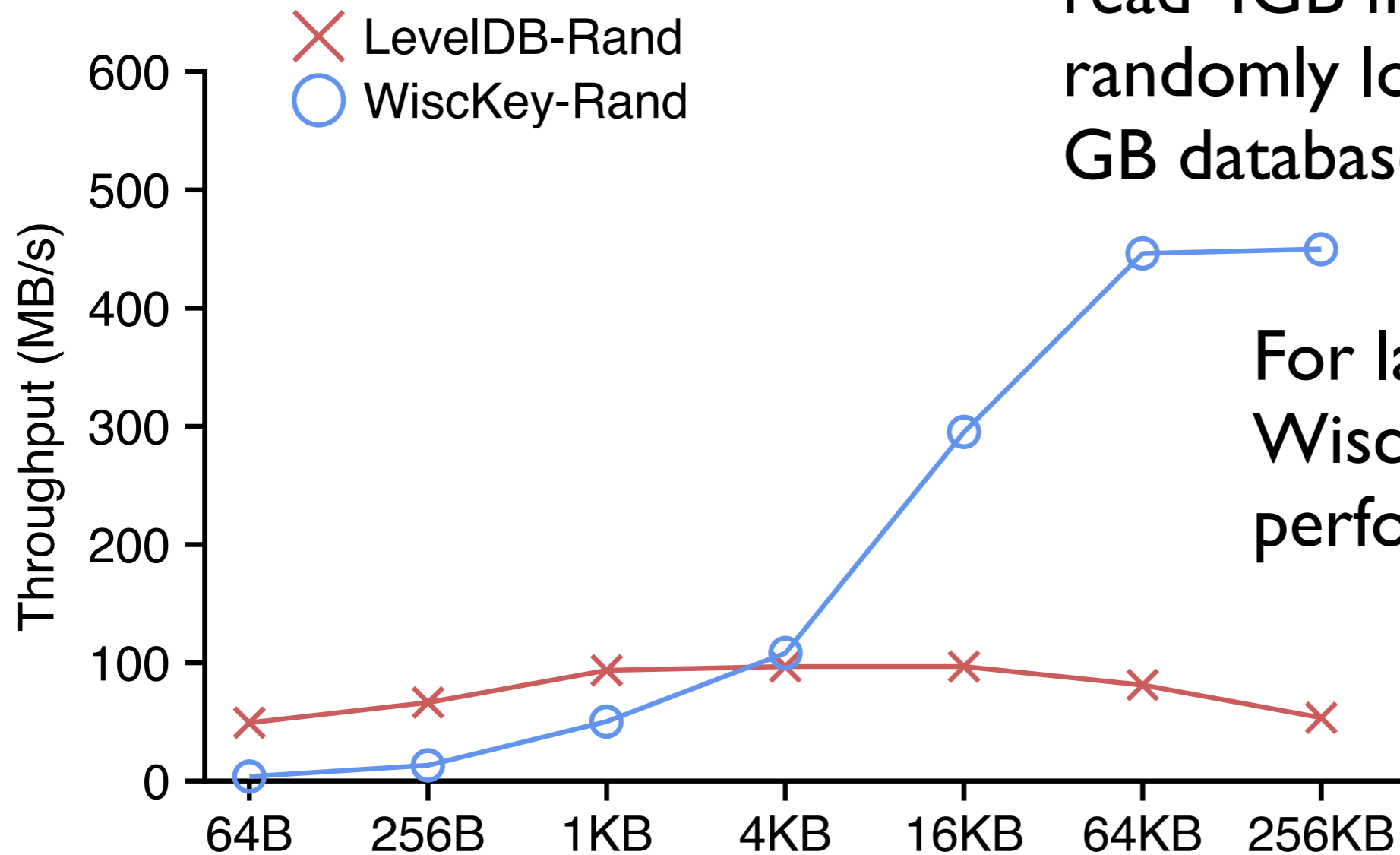
read 4GB from a  
randomly loaded 100  
GB database



Key: 16B, Value: 64B to 256KB

# Range Query

read 4GB from a  
randomly loaded 100  
GB database



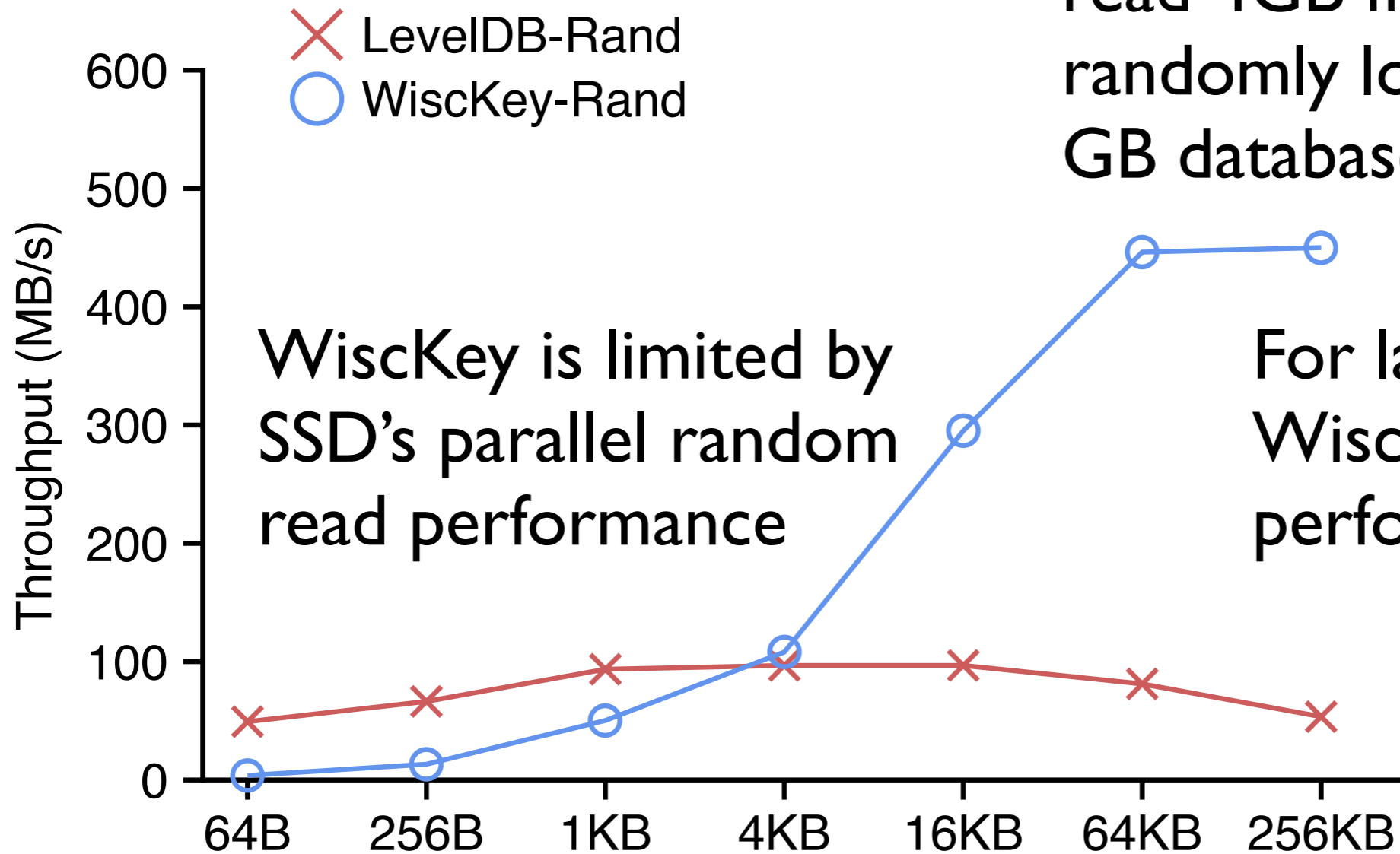
For large kv pairs,  
WiscKey can  
perform better

Key: 16B, Value: 64B to 256KB



# Range Query

read 4GB from a randomly loaded 100 GB database

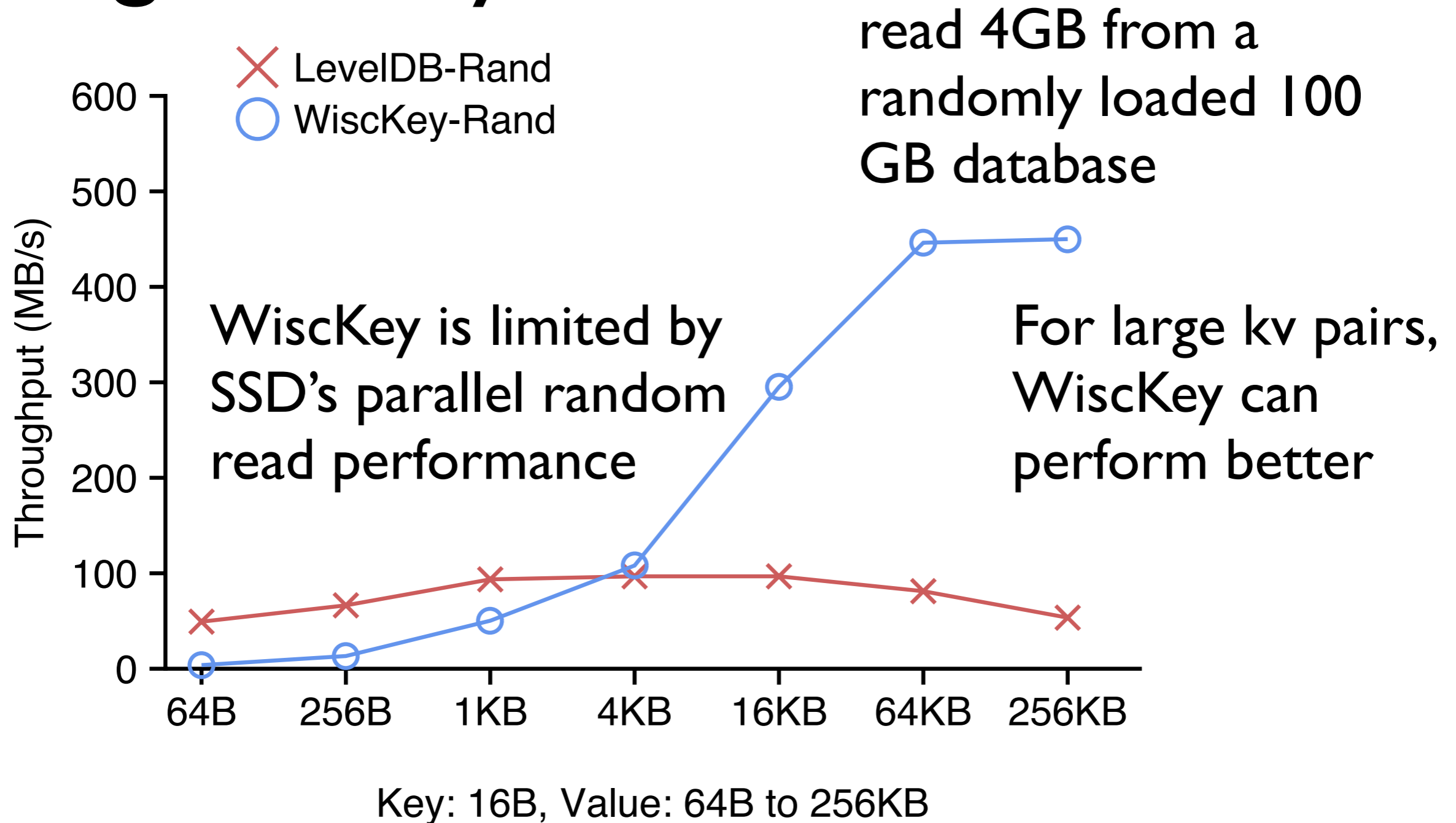


WiscKey is limited by SSD's parallel random read performance

For large kv pairs, WiscKey can perform better

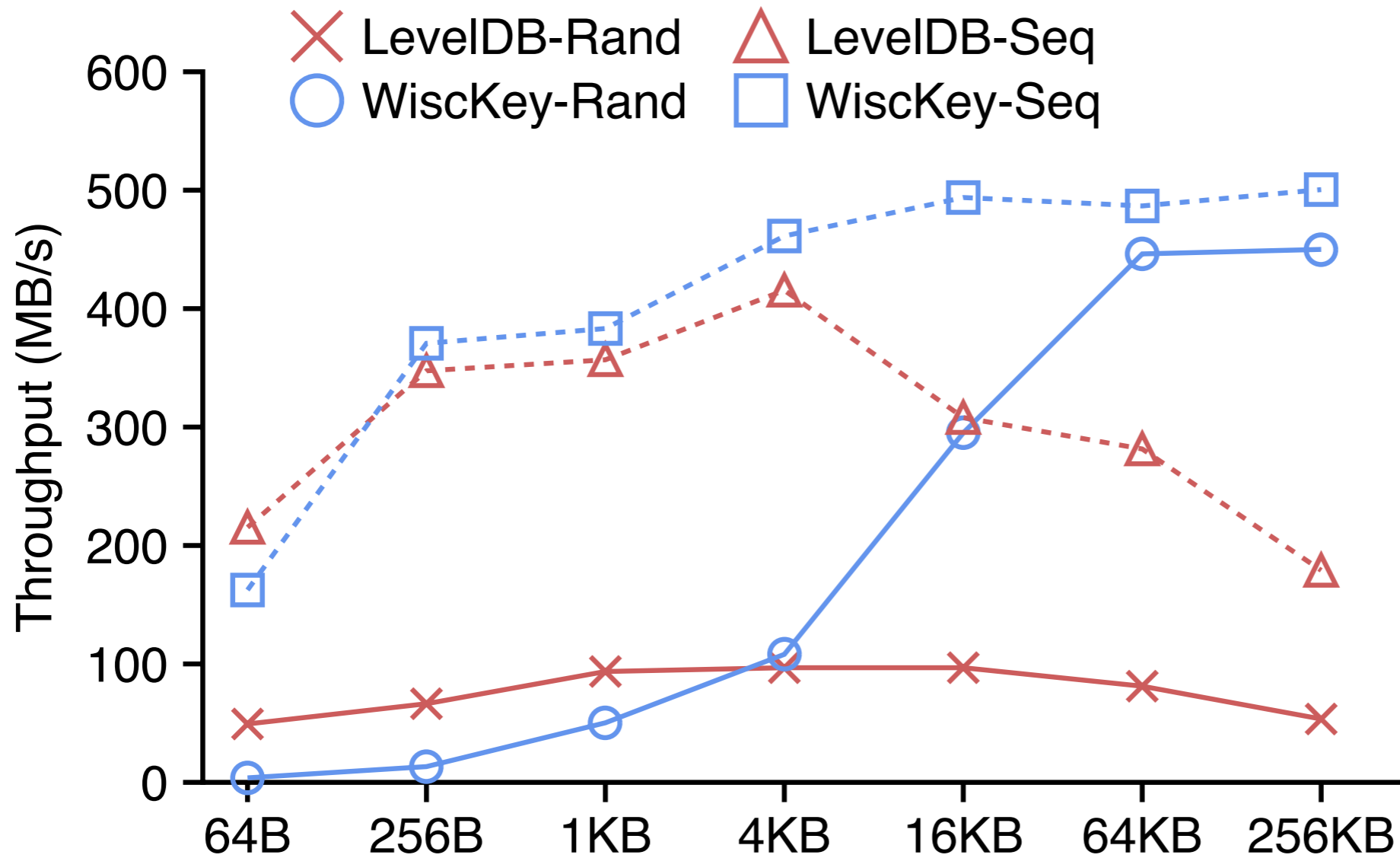
Key: 16B, Value: 64B to 256KB

# Range Query



Better for large kv pairs, but worse for small kv pairs on an unsorted database

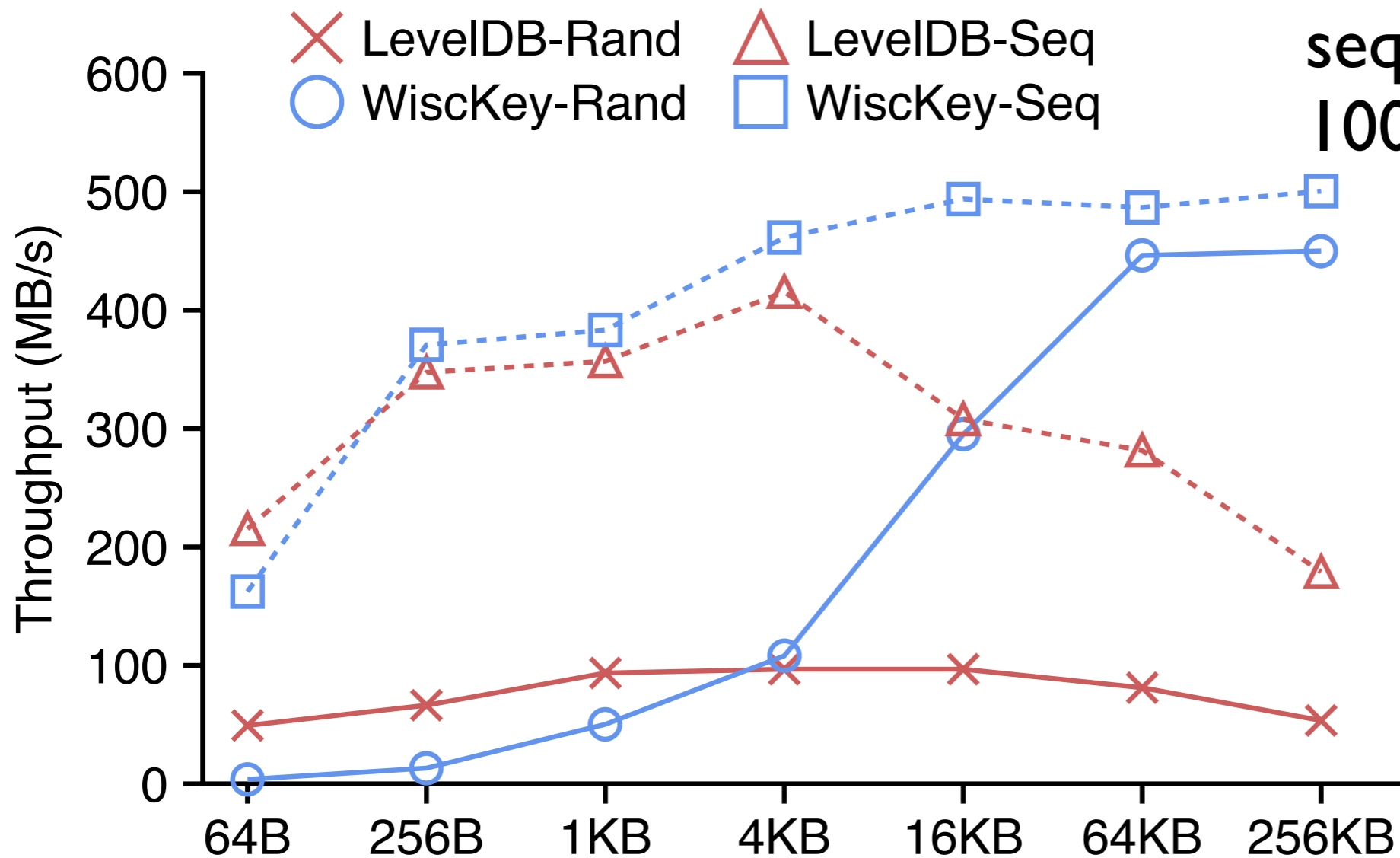
# Range Query



Key: 16B, Value: 64B to 256KB

# Range Query

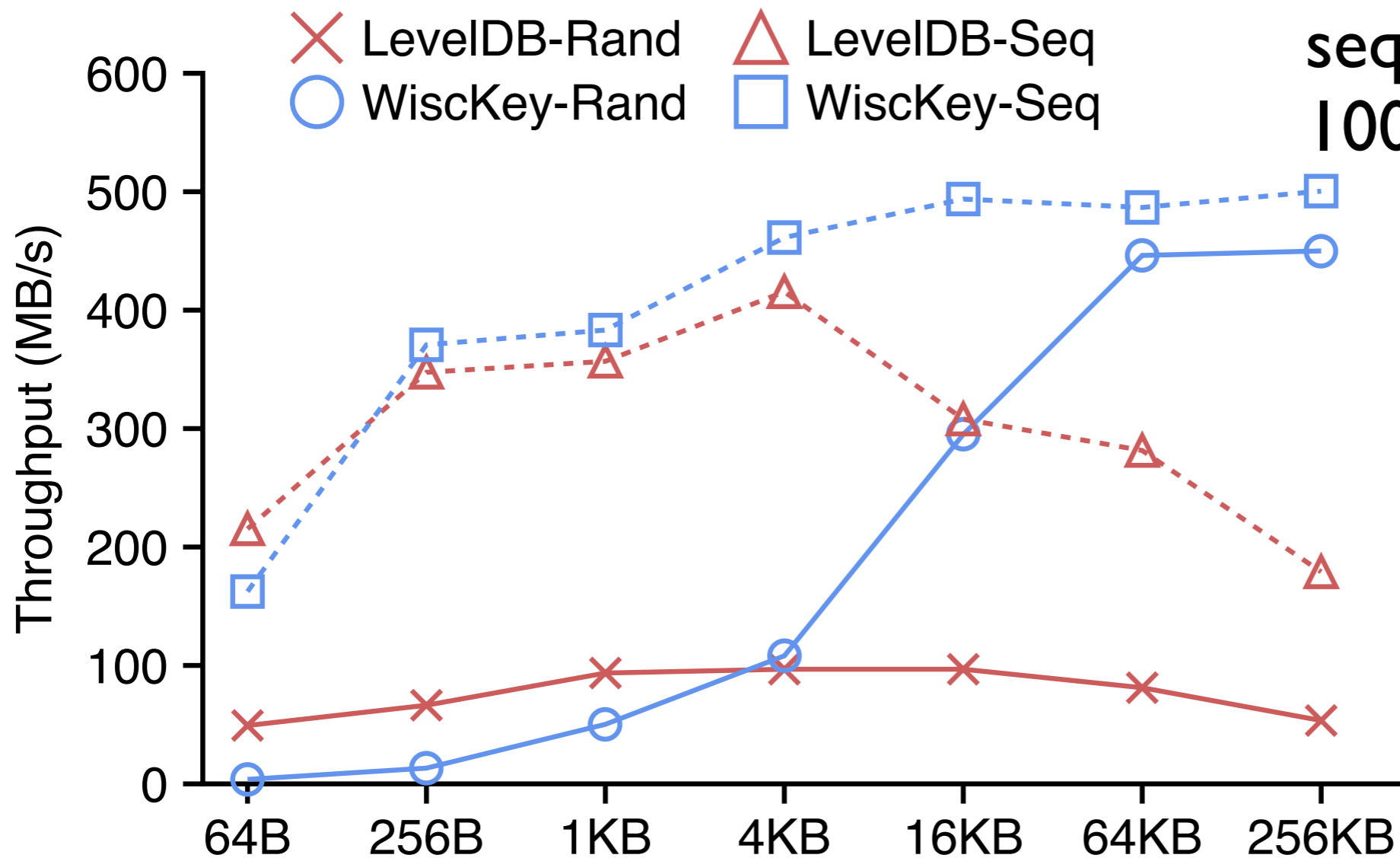
read 4GB from a  
sequentially loaded  
100 GB database



Key: 16B, Value: 64B to 256KB

# Range Query

read 4GB from a  
sequentially loaded  
100 GB database

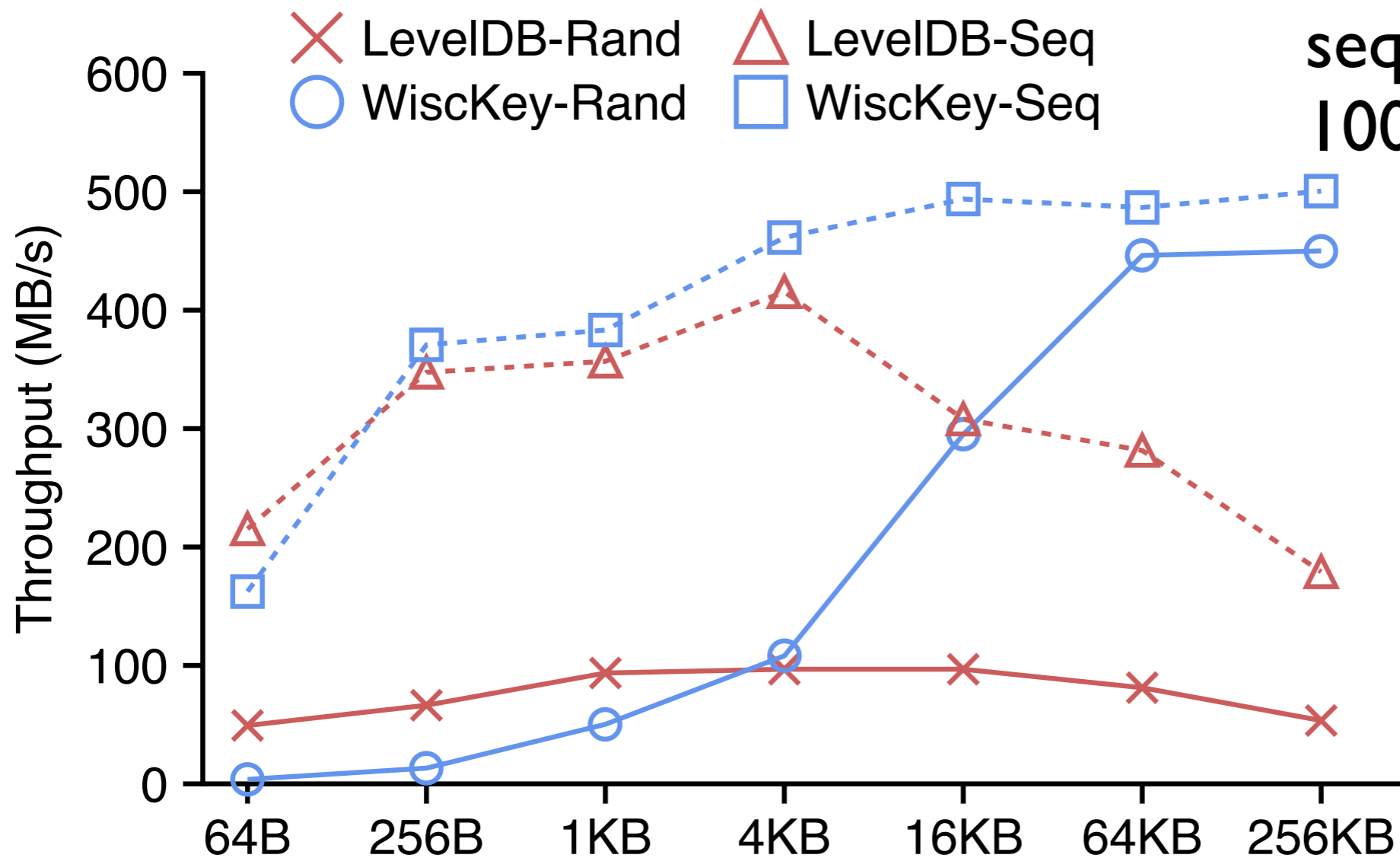


Key: 16B, Value: 64B to 256KB

Both WiscKey  
and LevelDB  
read sequentially

# Range Query

read 4GB from a sequentially loaded 100 GB database



Key: 16B, Value: 64B to 256KB

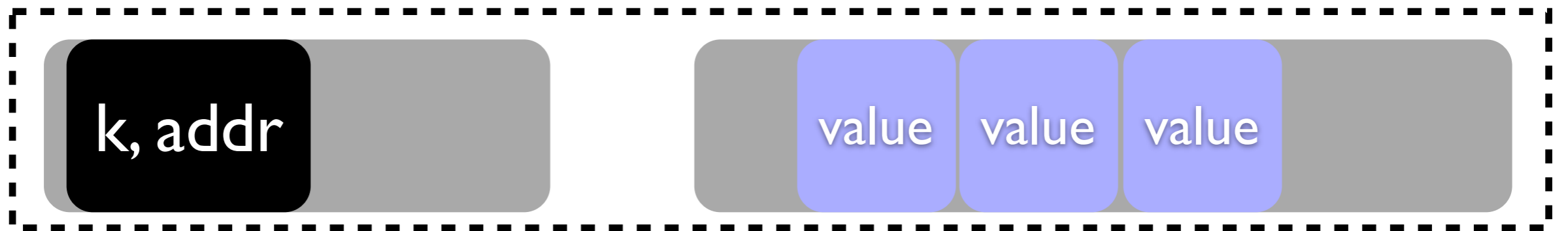
Both WiscKey and LevelDB read sequentially

Sorted databases help WiscKey's range query

# Optimizations

# Optimizations

SSD device



LSM-tree

Value Log

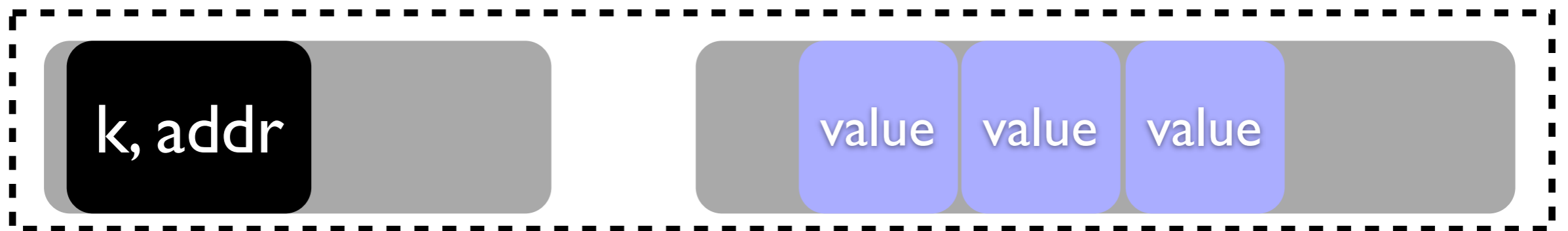


# Optimizations

Online and light-weight garbage collection

→ append (ksize, vsize, key, value) in value log

SSD device



LSM-tree

Value Log

# Optimizations

Online and light-weight garbage collection

→ append (ksize, vsize, key, value) in value log

SSD device



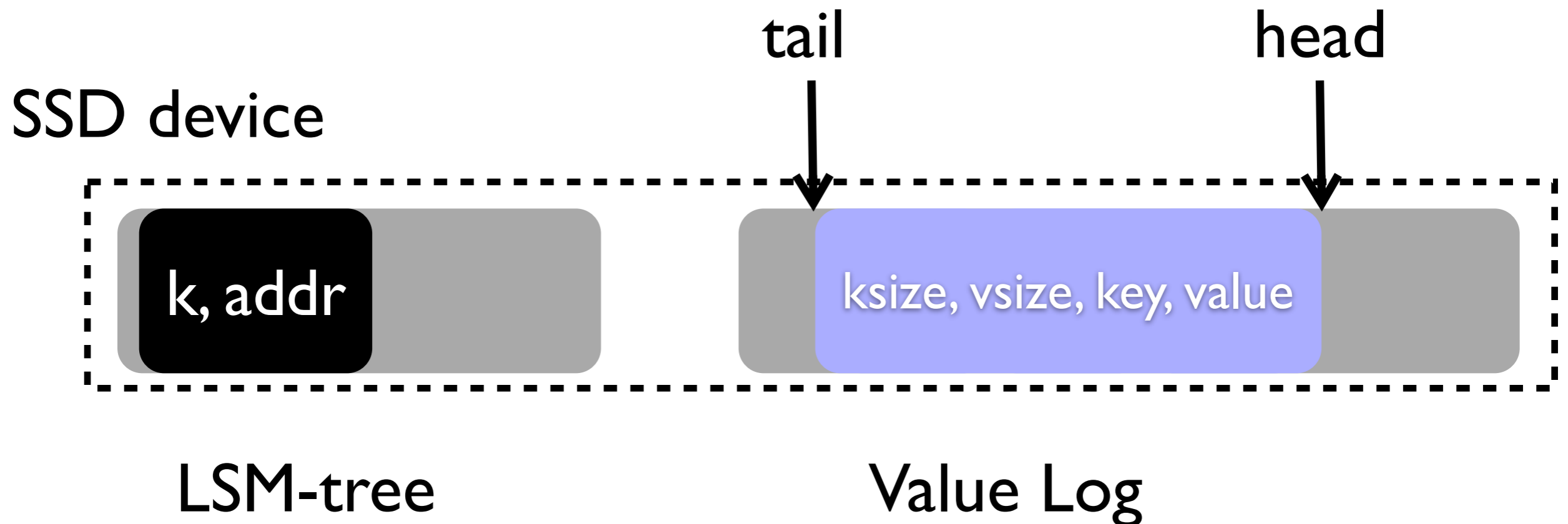
LSM-tree

Value Log

# Optimizations

Online and light-weight garbage collection

→ append (ksize, vsize, key, value) in value log

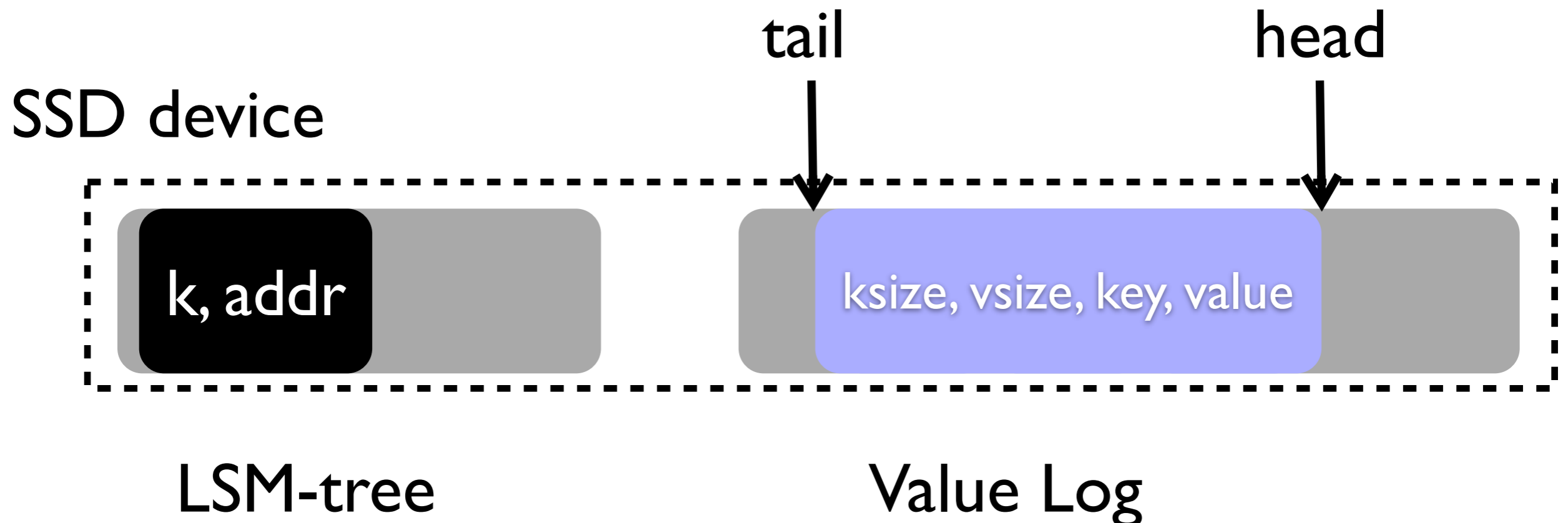


# Optimizations

Online and light-weight garbage collection

→ append (ksize, vsize, key, value) in value log

Remove LSM-tree log in WiscKey



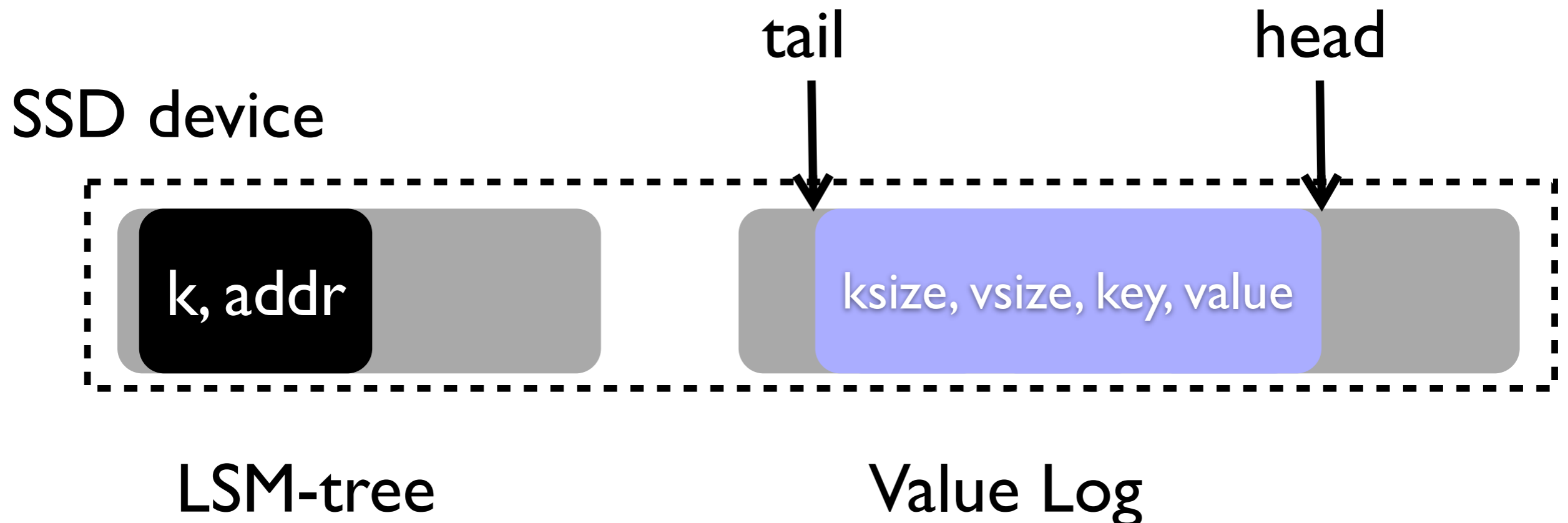
# Optimizations

Online and light-weight garbage collection

→ append (ksize, vsize, key, value) in value log

Remove LSM-tree log in WiscKey

→ store head in LSM-tree periodically



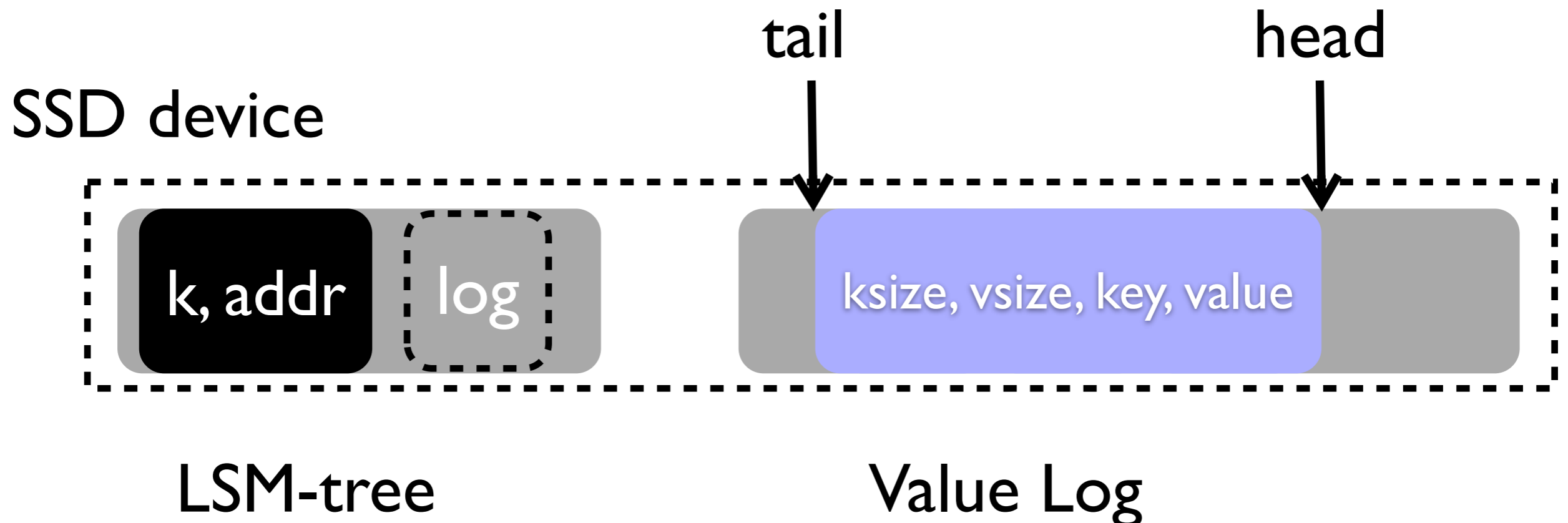
# Optimizations

## Online and light-weight garbage collection

- append (ksize, vsize, key, value) in value log

## Remove LSM-tree log in WiscKey

- store head in LSM-tree periodically
- scan the value log from the head to recover



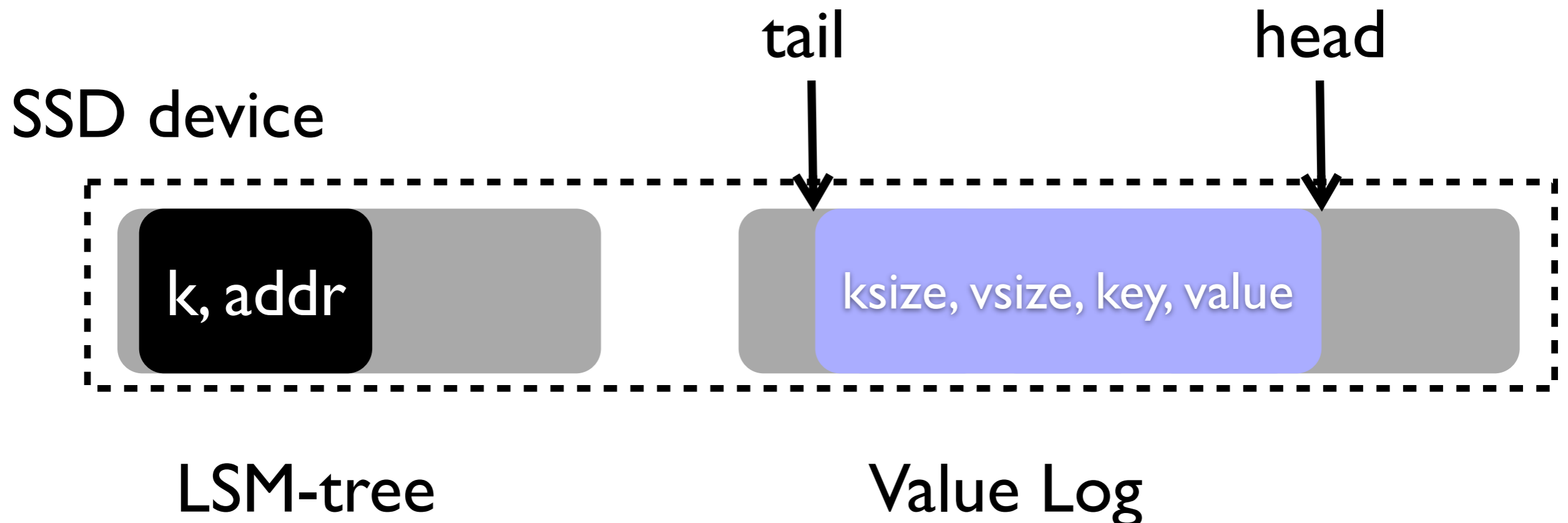
# Optimizations

## Online and light-weight garbage collection

- append (ksize, vsize, key, value) in value log

## Remove LSM-tree log in WiscKey

- store head in LSM-tree periodically
- scan the value log from the head to recover



# WiscKey Implementation



# WiscKey Implementation

## Based on LevelDB

- a separate vLog file for values
- modify I/O paths to separate keys and values
- leverages most of high-quality LevelDB source code

# WiscKey Implementation

## Based on LevelDB

- a separate vLog file for values
- modify I/O paths to separate keys and values
- leverages most of high-quality LevelDB source code

## Range query

- thread pool launches queries in parallel
- detect sequential pattern with the Iterator interface

# WiscKey Implementation

## Based on LevelDB

- a separate vLog file for values
- modify I/O paths to separate keys and values
- leverages most of high-quality LevelDB source code

## Range query

- thread pool launches queries in parallel
- detect sequential pattern with the Iterator interface

## File-system support

- fadvise to predeclare access patterns
- hole-punching to free space

Background

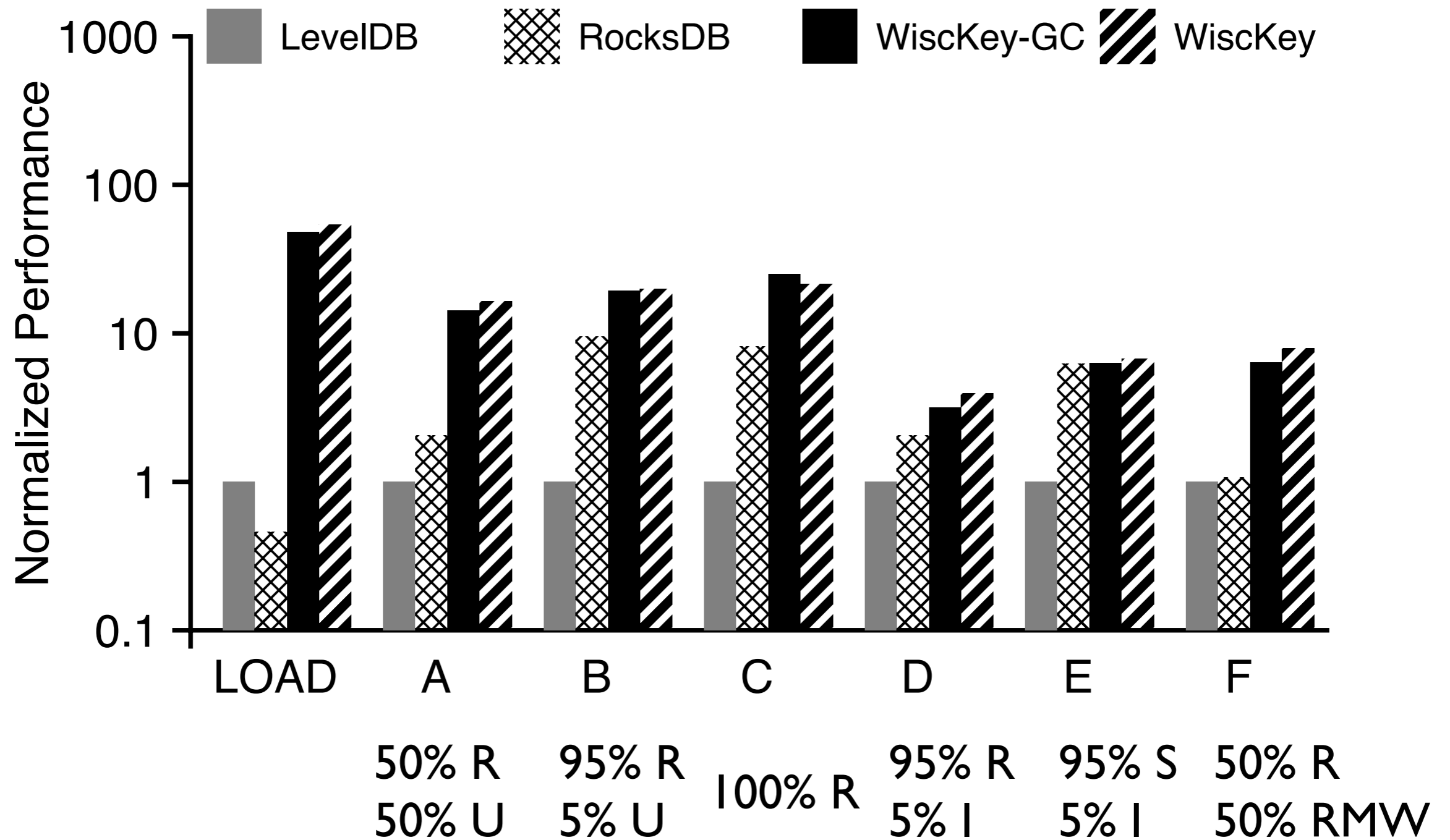
Key-Value Separation

Challenges and Optimizations

**Evaluation**

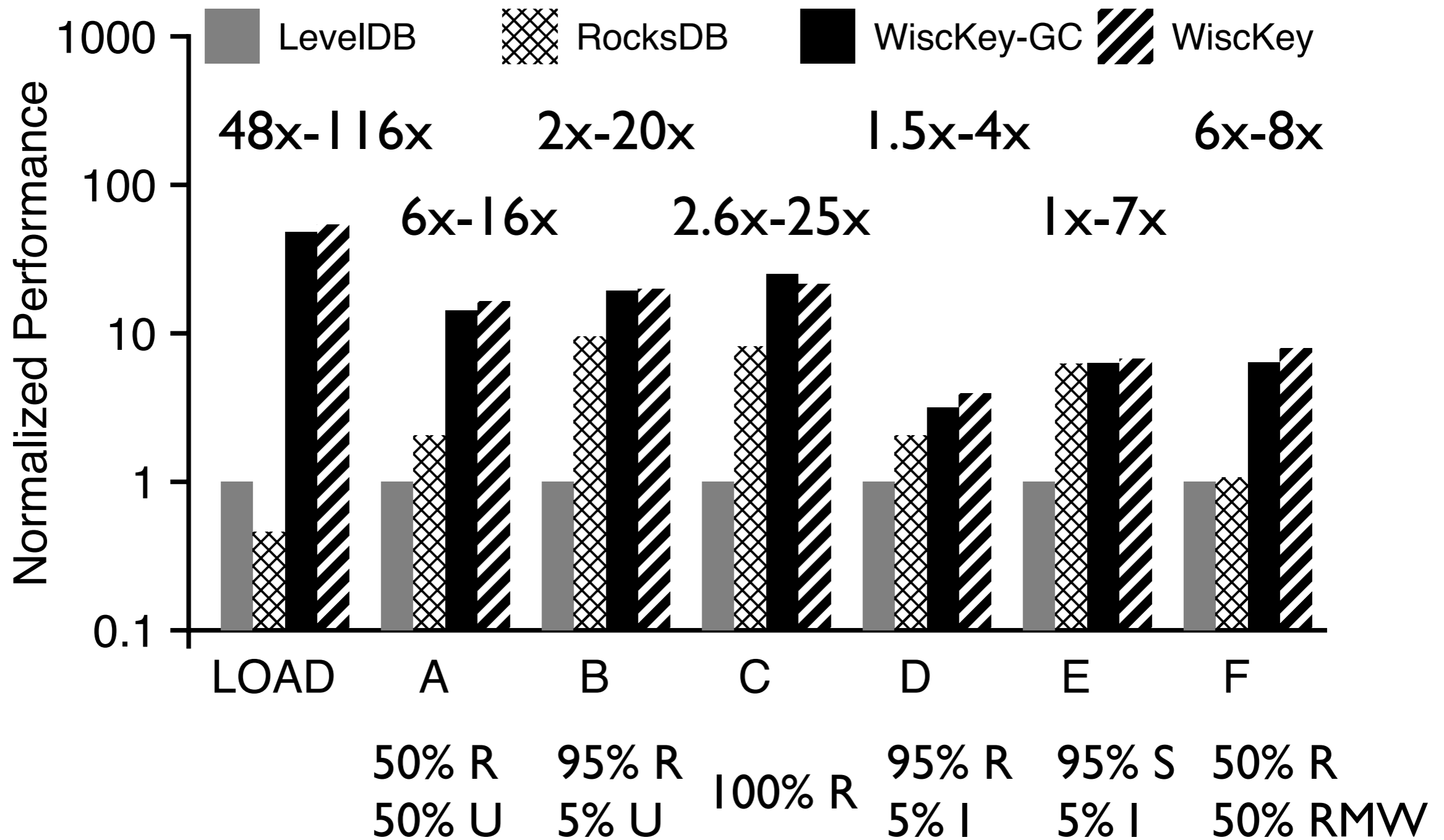
**Conclusion**

# YCSB Benchmarks



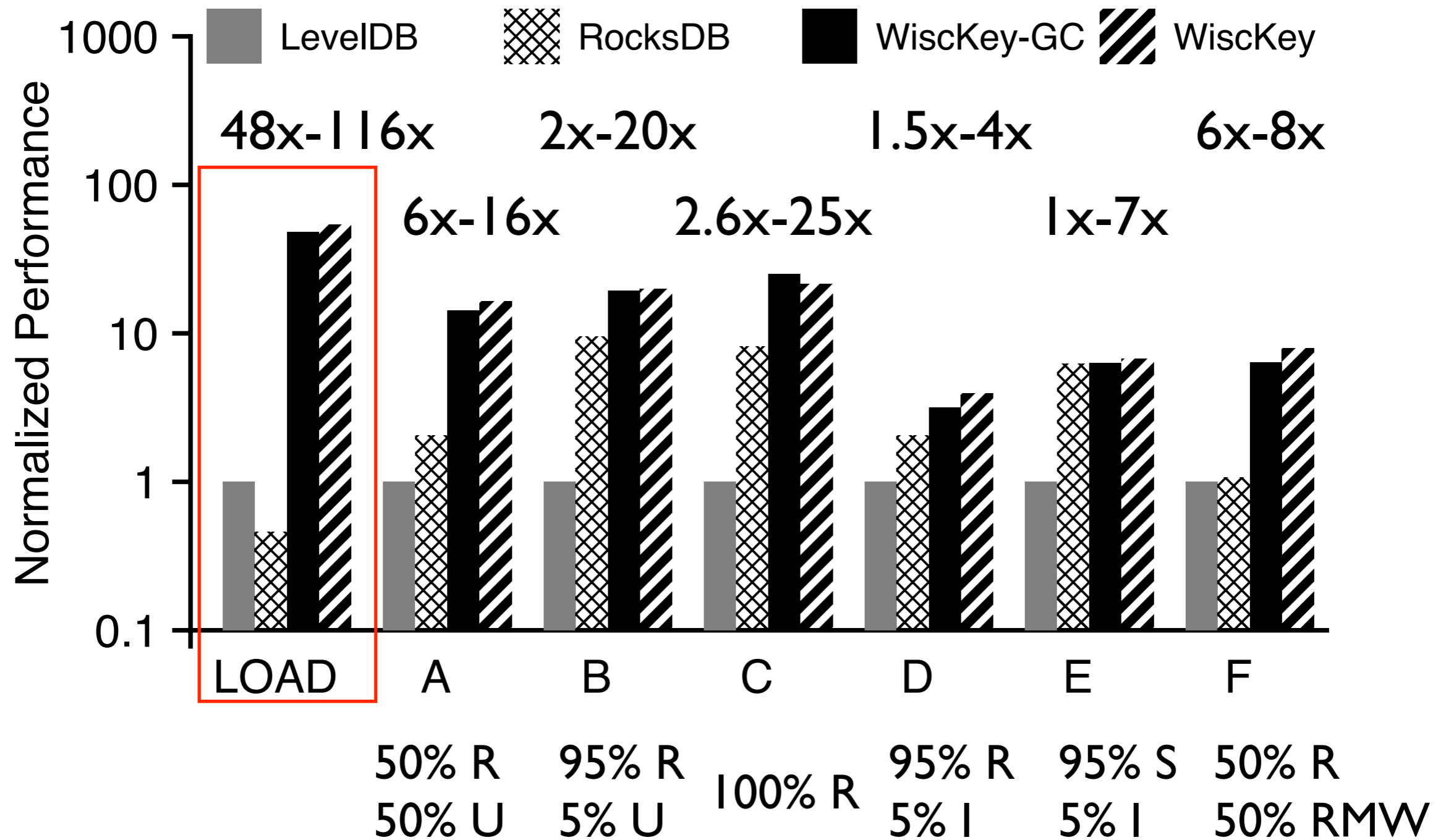
Key size: 16B, Value size: 1KB

# YCSB Benchmarks



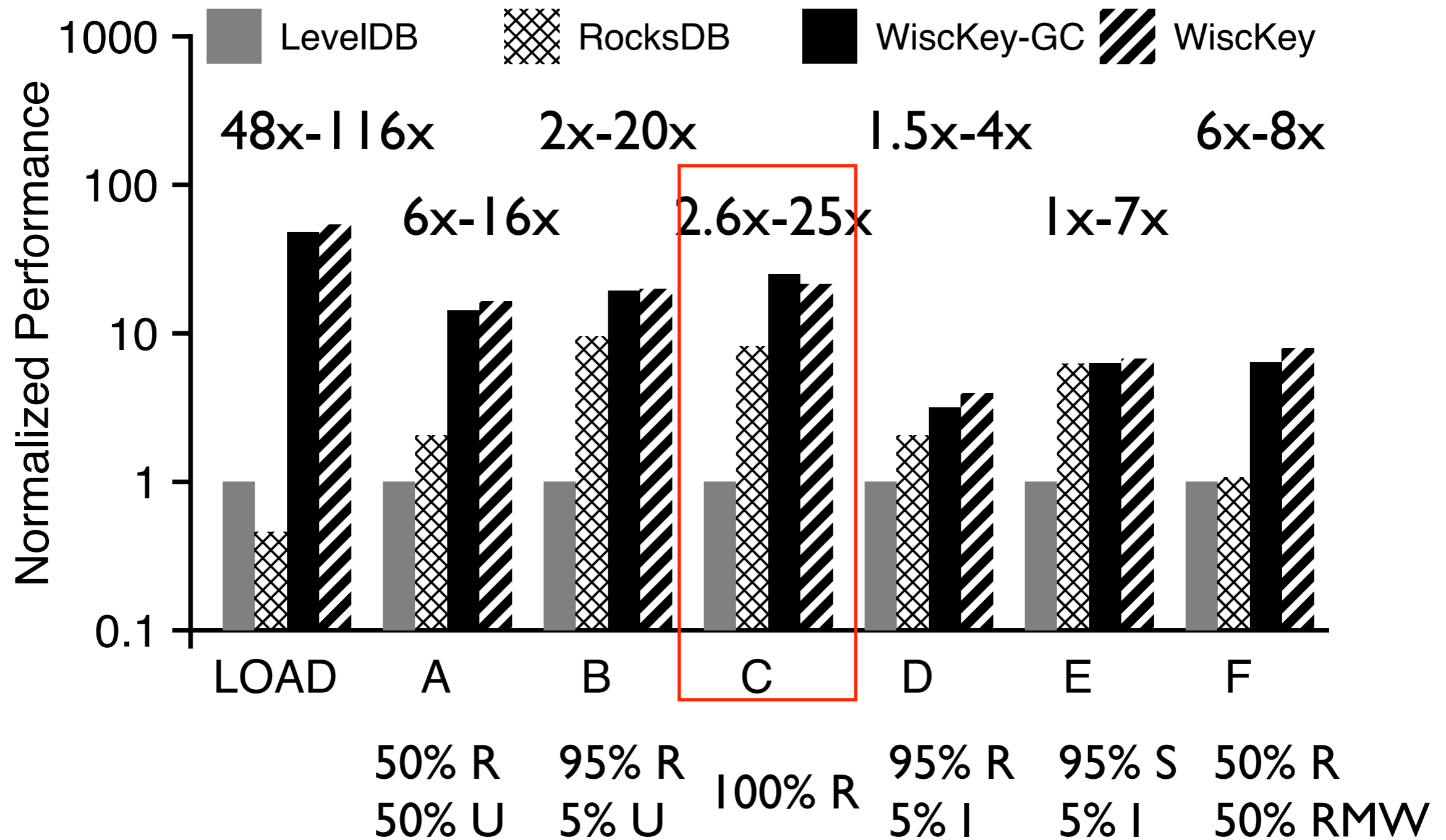
Key size: 16B, Value size: 1KB

# YCSB Benchmarks



Key size: 16B, Value size: 1KB

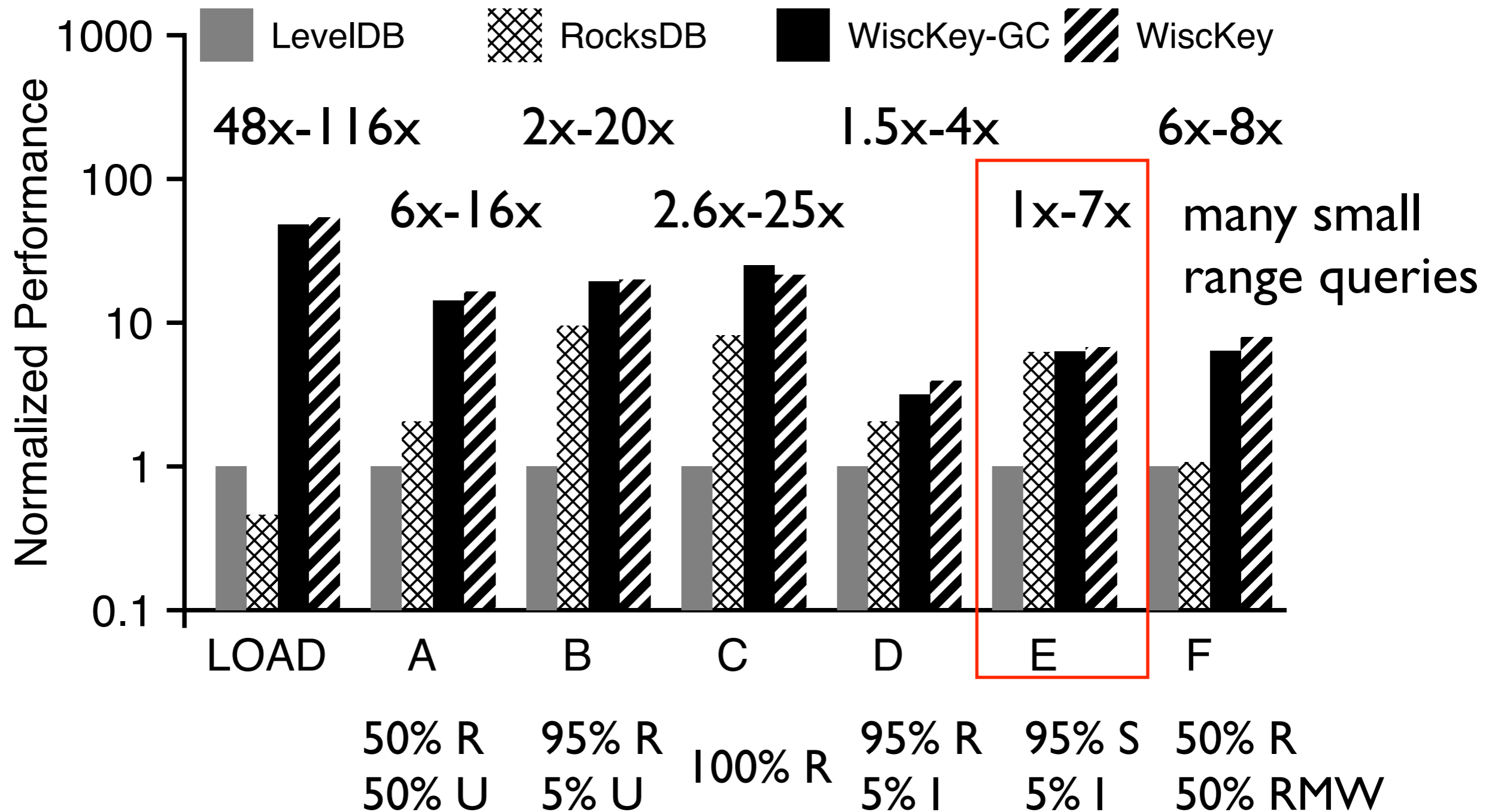
# YCSB Benchmarks



Key size: 16B, Value size: 1KB



# YCSB Benchmarks



Key size: 16B, Value size: 1KB

# Conclusion

# Conclusion

WiscKey: a LSM-tree based key-value store

- **decouple** sorting and garbage collection by **separating keys from values**
- SSD-conscious designs
- significant performance gain

# Conclusion

WiscKey: a LSM-tree based key-value store

- **decouple** sorting and garbage collection by **separating keys from values**
- SSD-conscious designs
- significant performance gain

Transition to new storage hardware

- understand and leverage existing software
- explore new designs to utilize the new hardware
- get the best of two worlds