

Anatomy of LSM Memory Buffer

Shubham Kaushik

kaushiks@brandeis.edu



Brandeis
UNIVERSITY



Log-Structured Merge-tree

The Log-Structured Merge-Tree (LSM-Tree)

1996

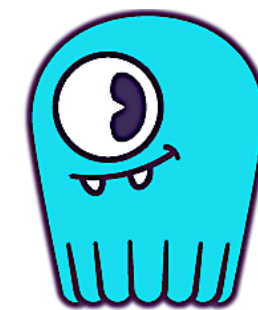
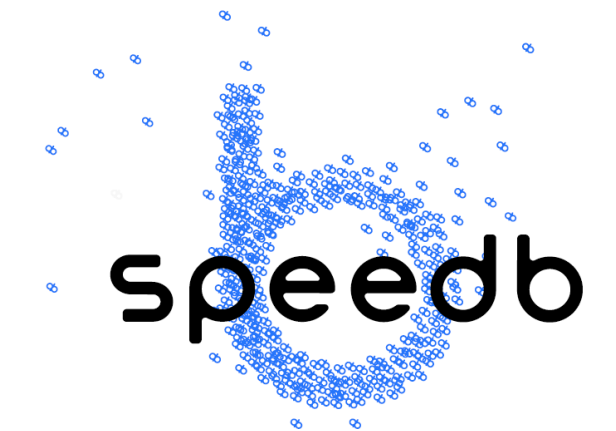
Patrick O'Neil¹, Edward Cheng²
Dieter Gawlick³, Elizabeth O'Neil¹
To be published: Acta Informatica



LSM-based Datastore

2004

2024



InnoDB

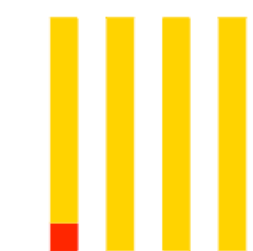


CockroachDB

SCYLLA



Couchbase



ClickHouse

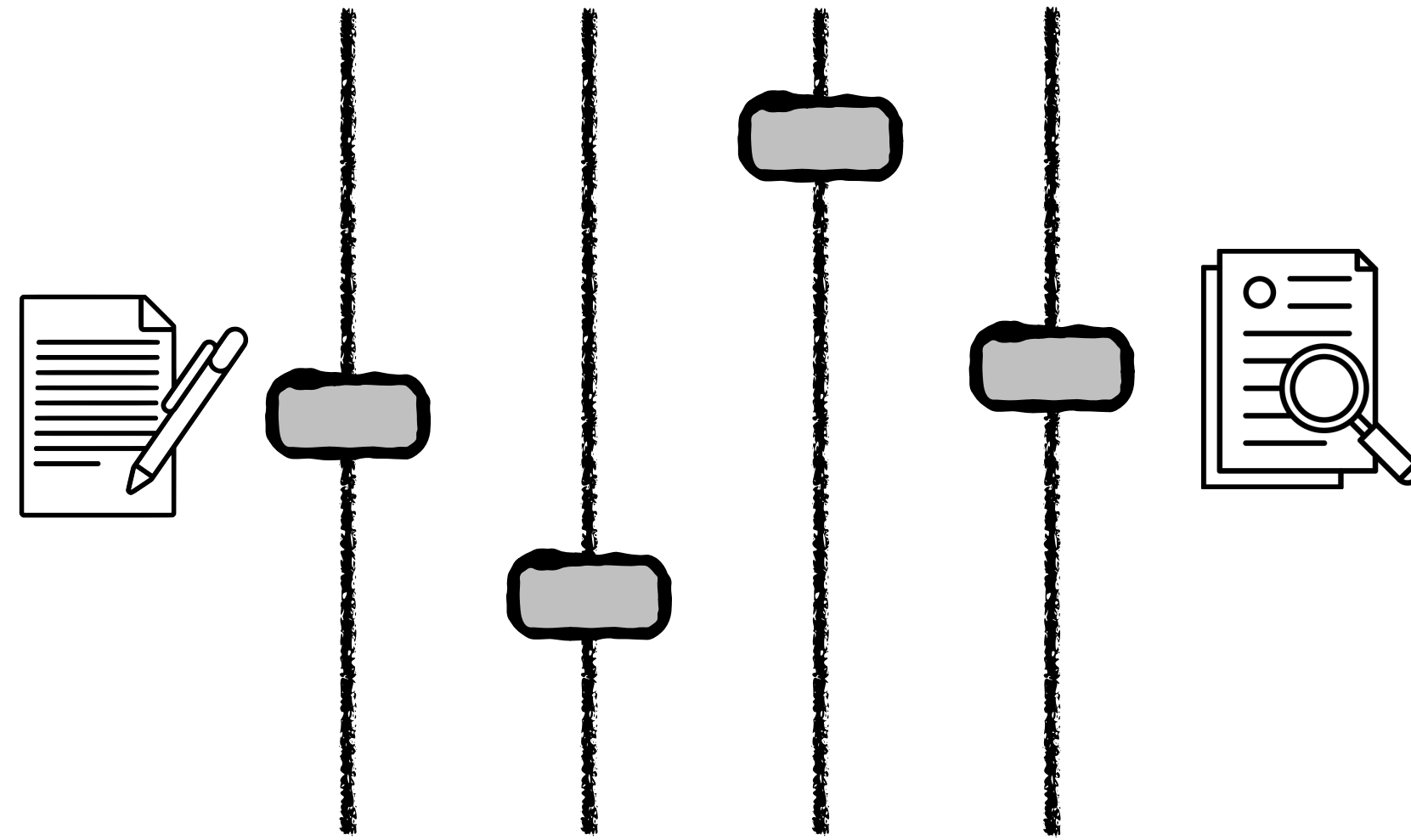


Brandeis
UNIVERSITY

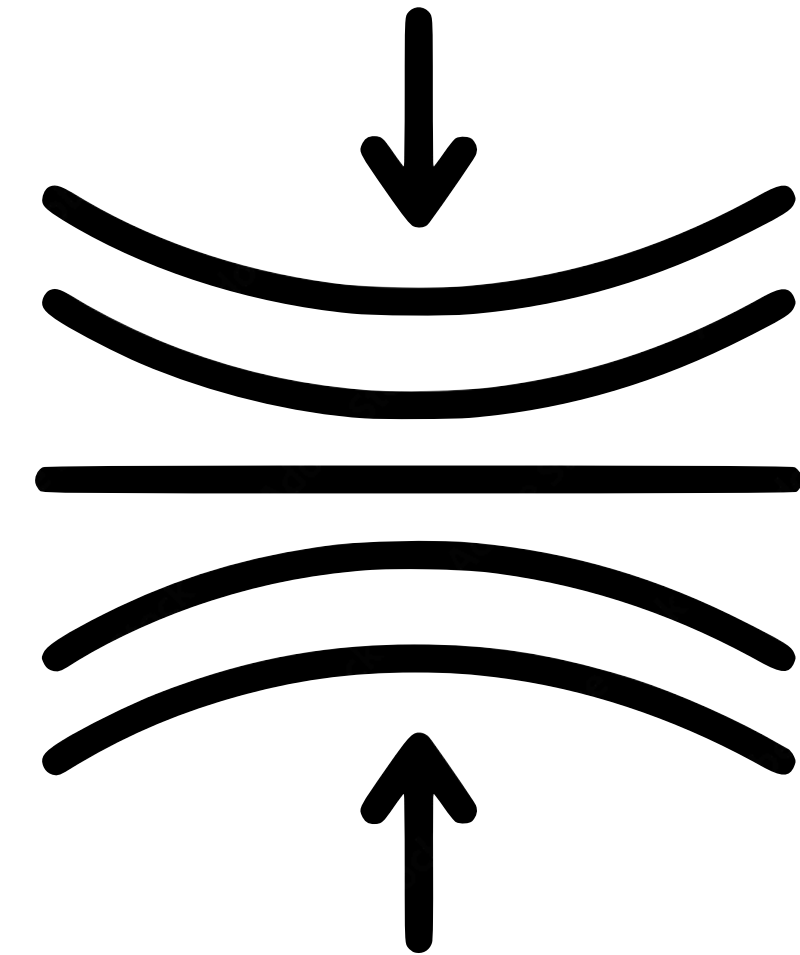
LSM-tree



fast **writes**

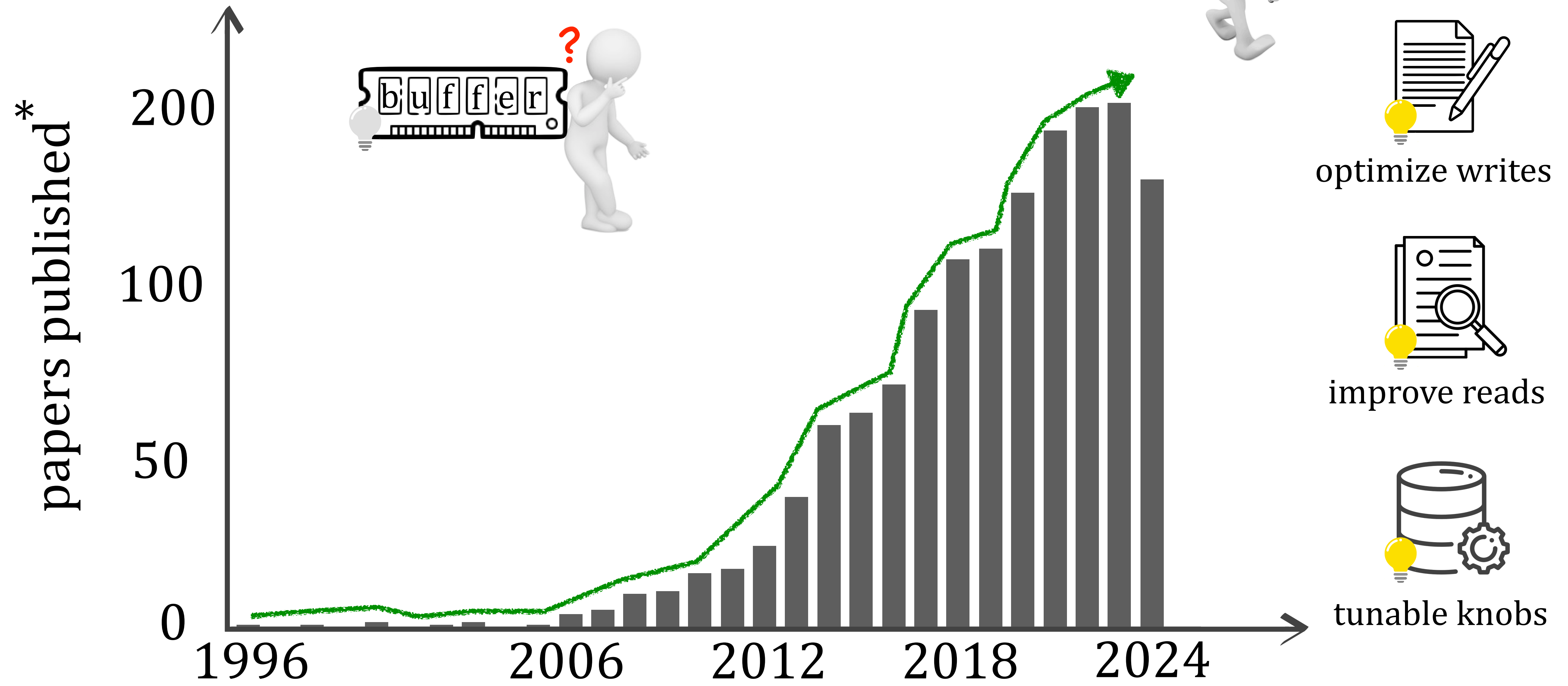


tunable read-write
performance



good **space**
utilization

Research Trend



optimize writes

improve reads

tunable knobs

* papers published

* data from Google scholar



Memory Buffer

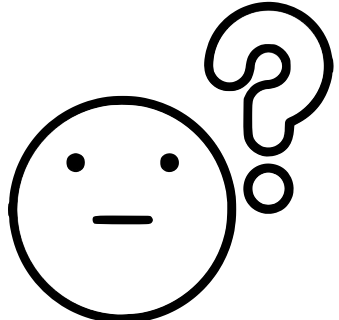
a temporary space that holds data in memory

Implementations of memory buffer

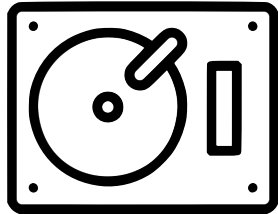
- 1 vector
- 2 skip-list
- 3 hash-hybrids
- 4 trie



LSM Basics



Why do we care?



L1



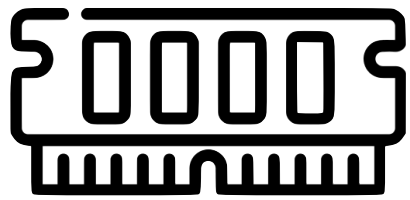
L2



L3



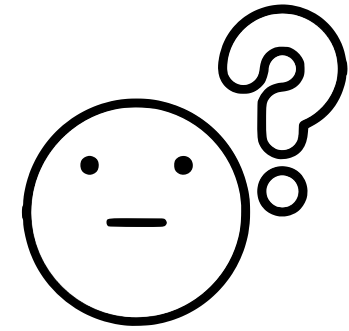
L4



buffer



Brandeis UNIVERSITY



Ideal buffer implementation?



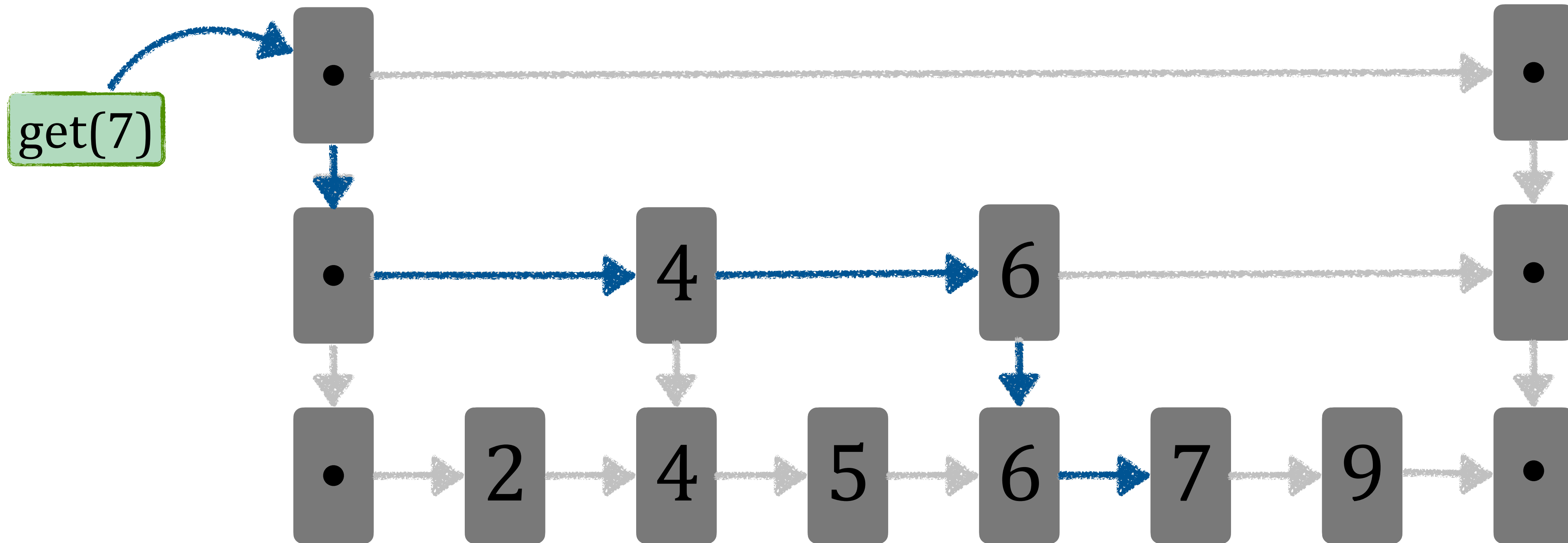
DEFAULT
skip-list



Brandeis
UNIVERSITY

N : entries in buffer
M : meta data

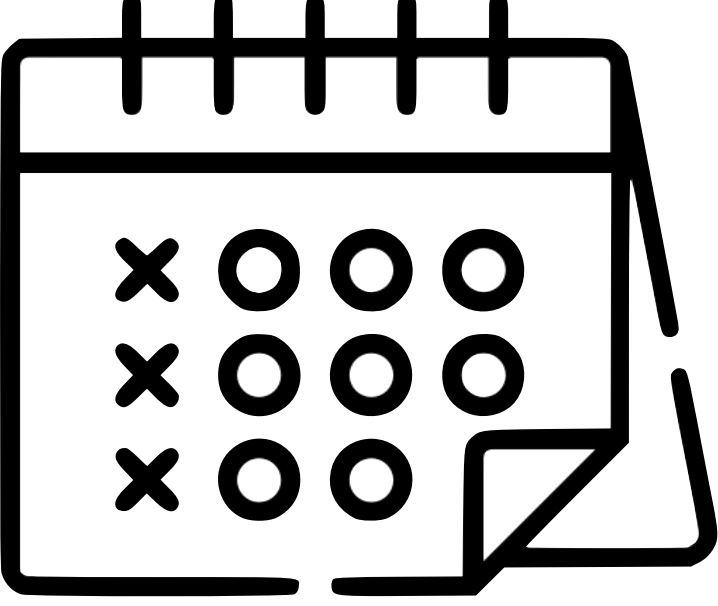
Implementation: Skip-list




- great for mixed w/l
- some extra space needed
- good for point queries

insert cost: $\mathcal{O}(\log N)$
space complexity: $\mathcal{O}(N + M)$
point query cost: $\mathcal{O}(\log N)$

DEFAULT
skip-list



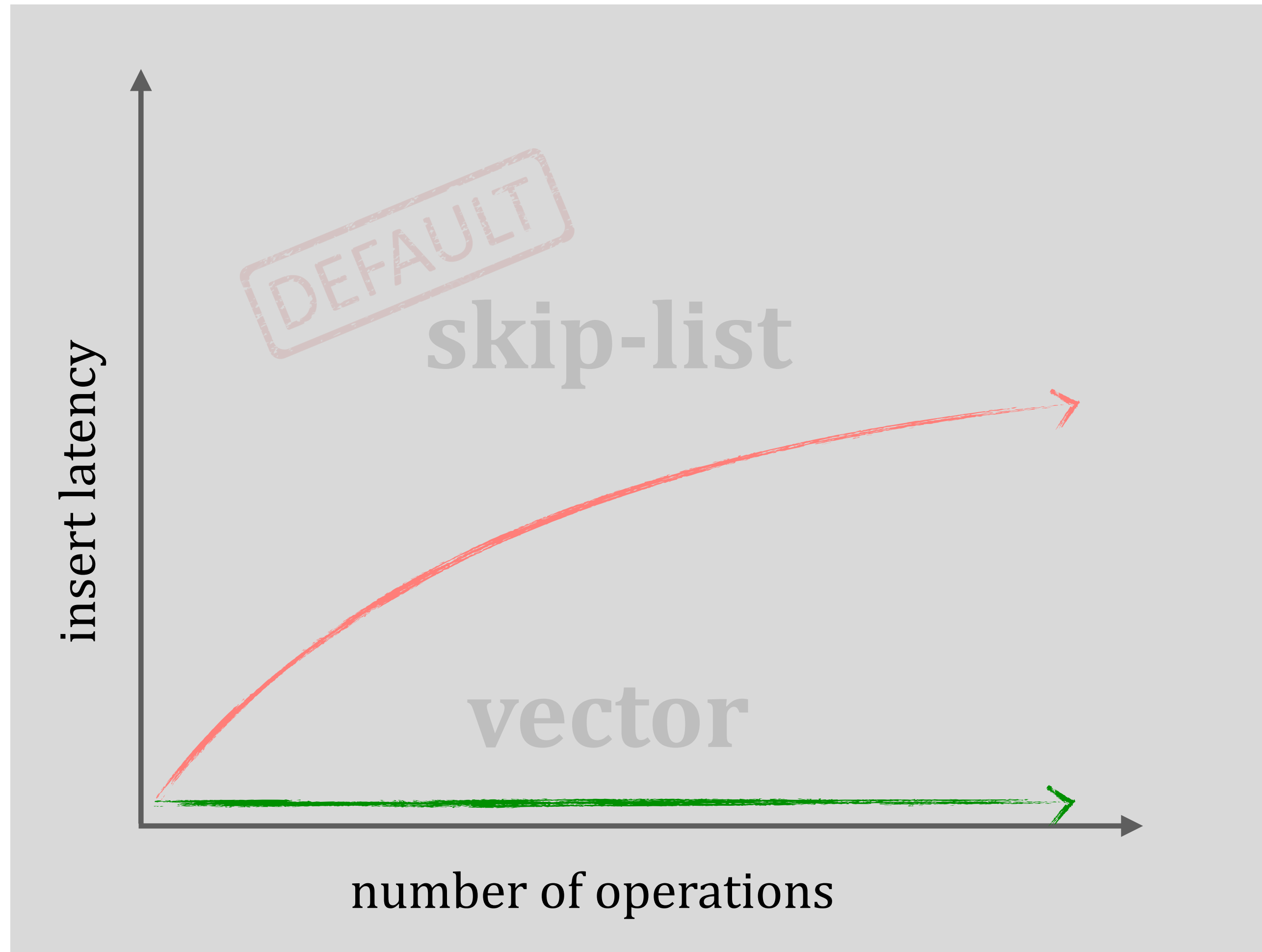
event logging



time-series

insert heavy



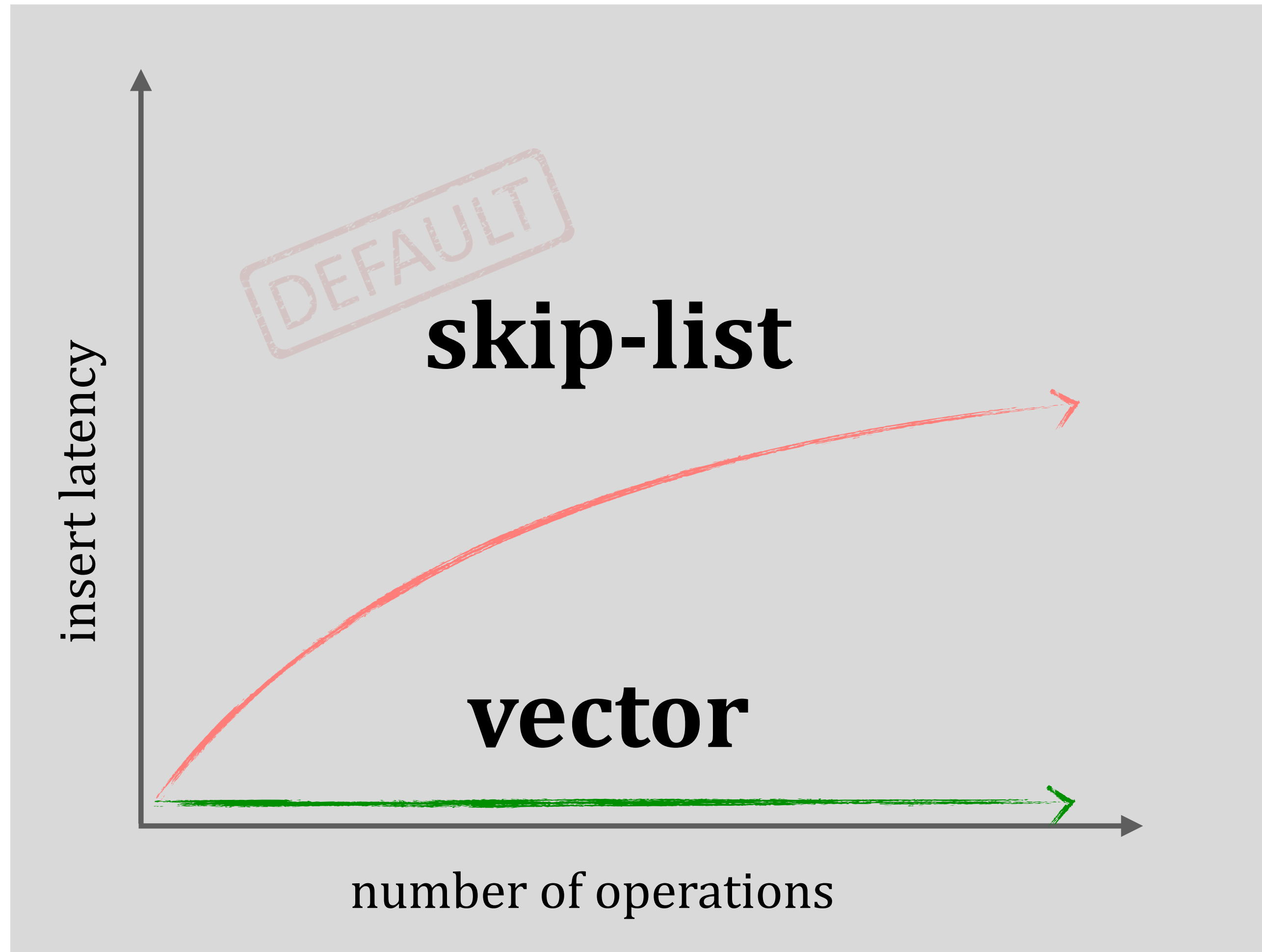


event logging

time-series

insert heavy





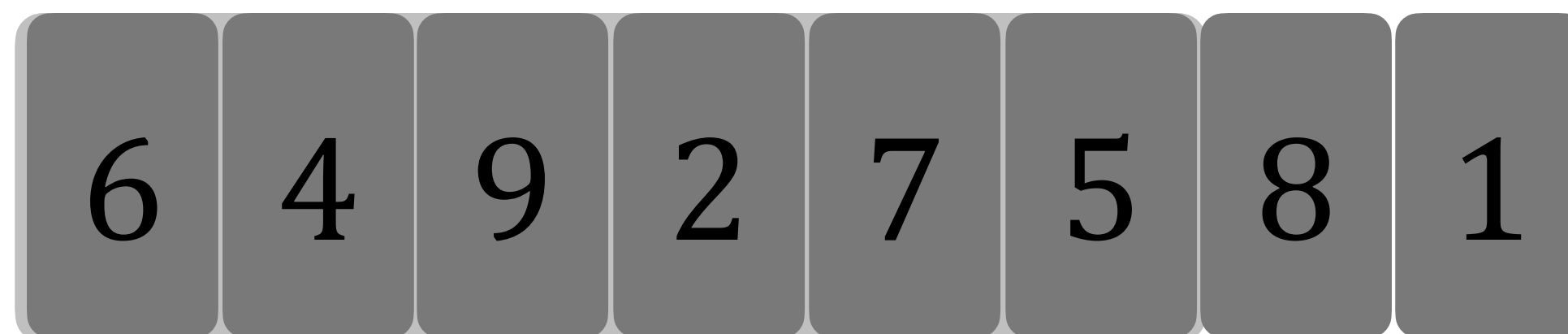
event logging

time-series

insert heavy

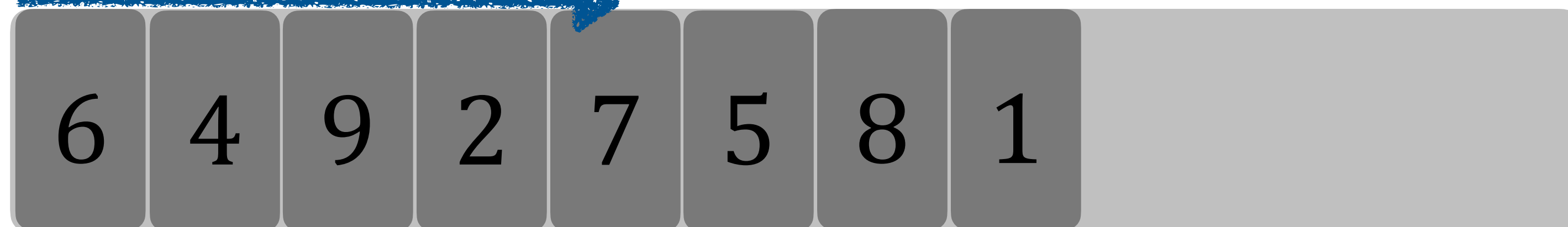


Implementation: Vector



dynamic

get(7)



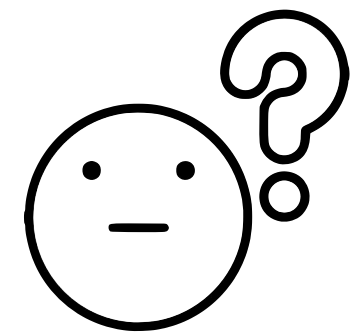
static

- great for insert-heavy w/l
- no extra space needed
- expensive point queries

insert cost: $\mathcal{O}(1)$

space complexity: $\mathcal{O}(N)$

point query cost: $\mathcal{O}(N)$



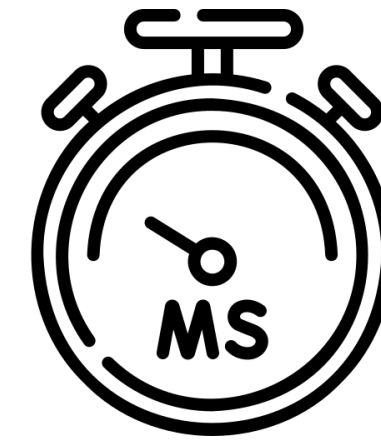
Ideal buffer implementation?

depends on a *workload composition*

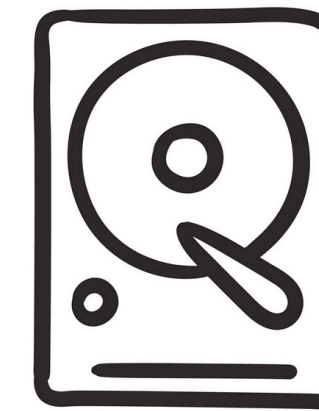
- 1 insert **heavy**
- 2 insert with **point** queries
- 3 insert with **range** queries
- 4 insert with **updates**

Research Focus

- 1 insert **heavy**
- 2 insert with **point** queries
- 3 insert with **range** queries
- 4 insert with **updates**



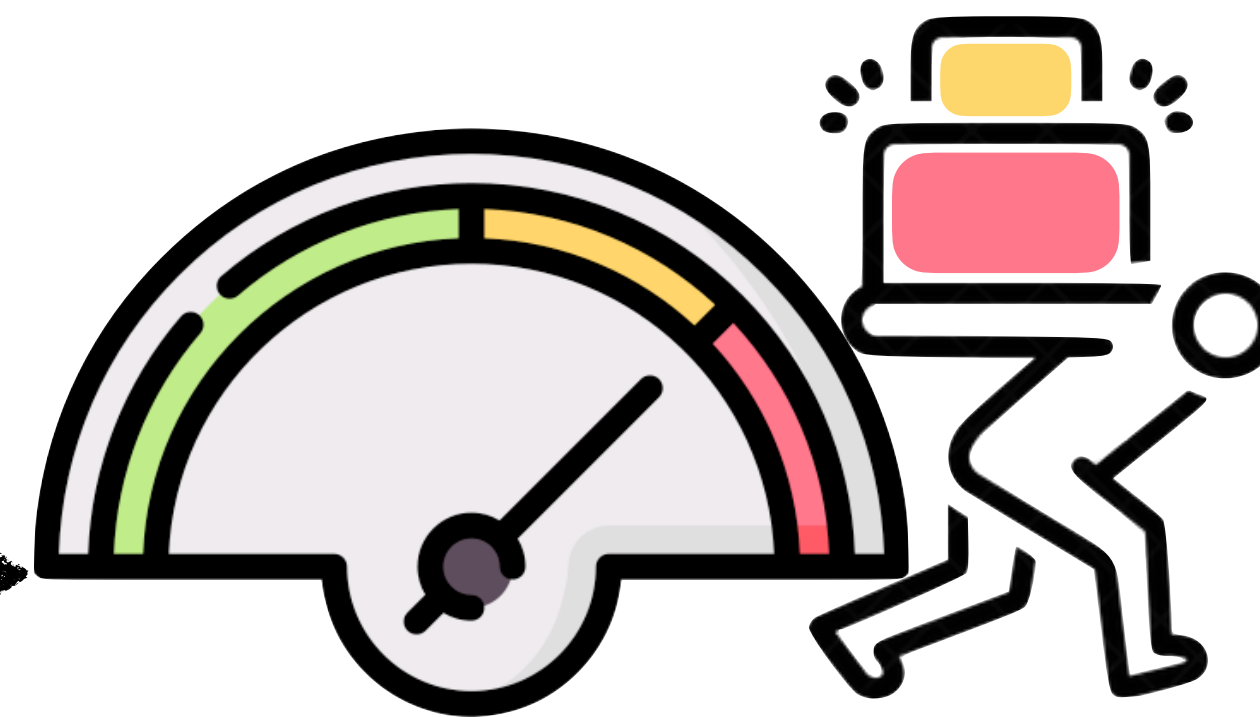
latency



memory footprint



throughput



benchmark

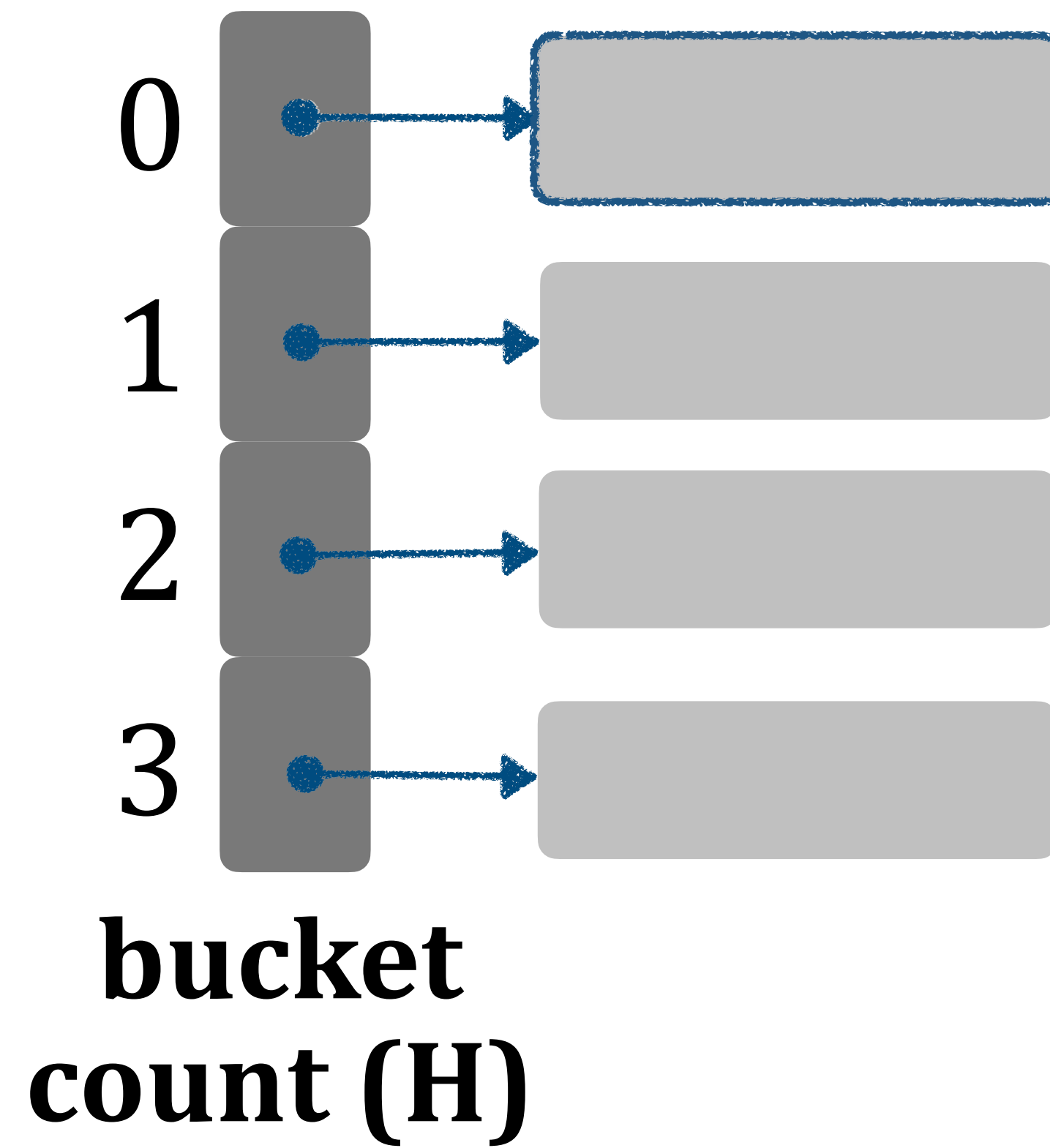
| | vector | skip-list | hash-hybrid | trie |
|---|--------|-----------|-------------|------|
| 1 | ?? | | | |
| 2 | | ?? | | |
| 3 | | | ?? | |
| 4 | | | | ?? |

N : entries in buffer
M : meta data

Implementation: Hash Hybrids

prefix
length (X)

$$f(k)$$

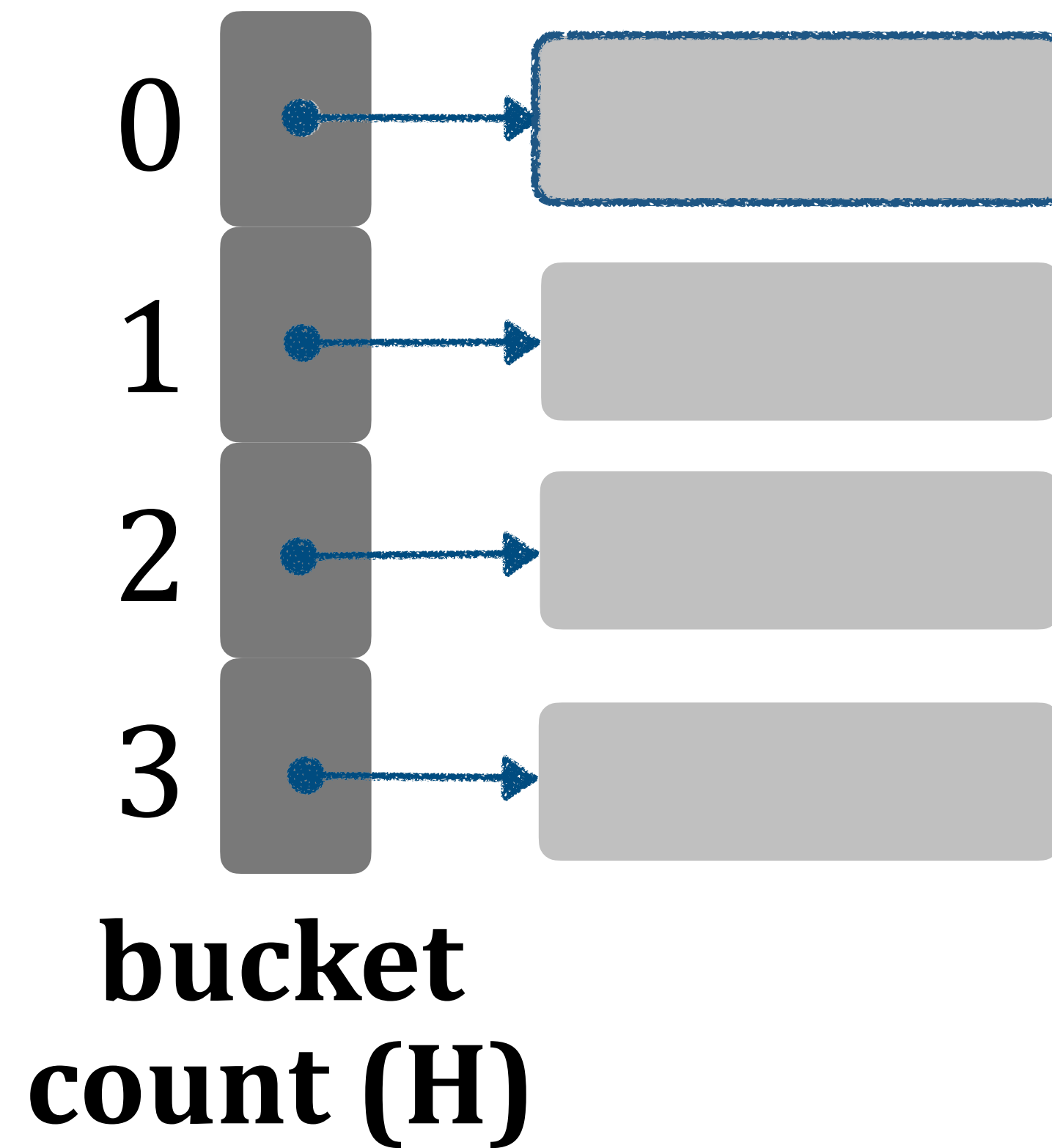


N : entries in buffer
M : meta data

Implementation: Hash Hybrids

prefix
length (X)

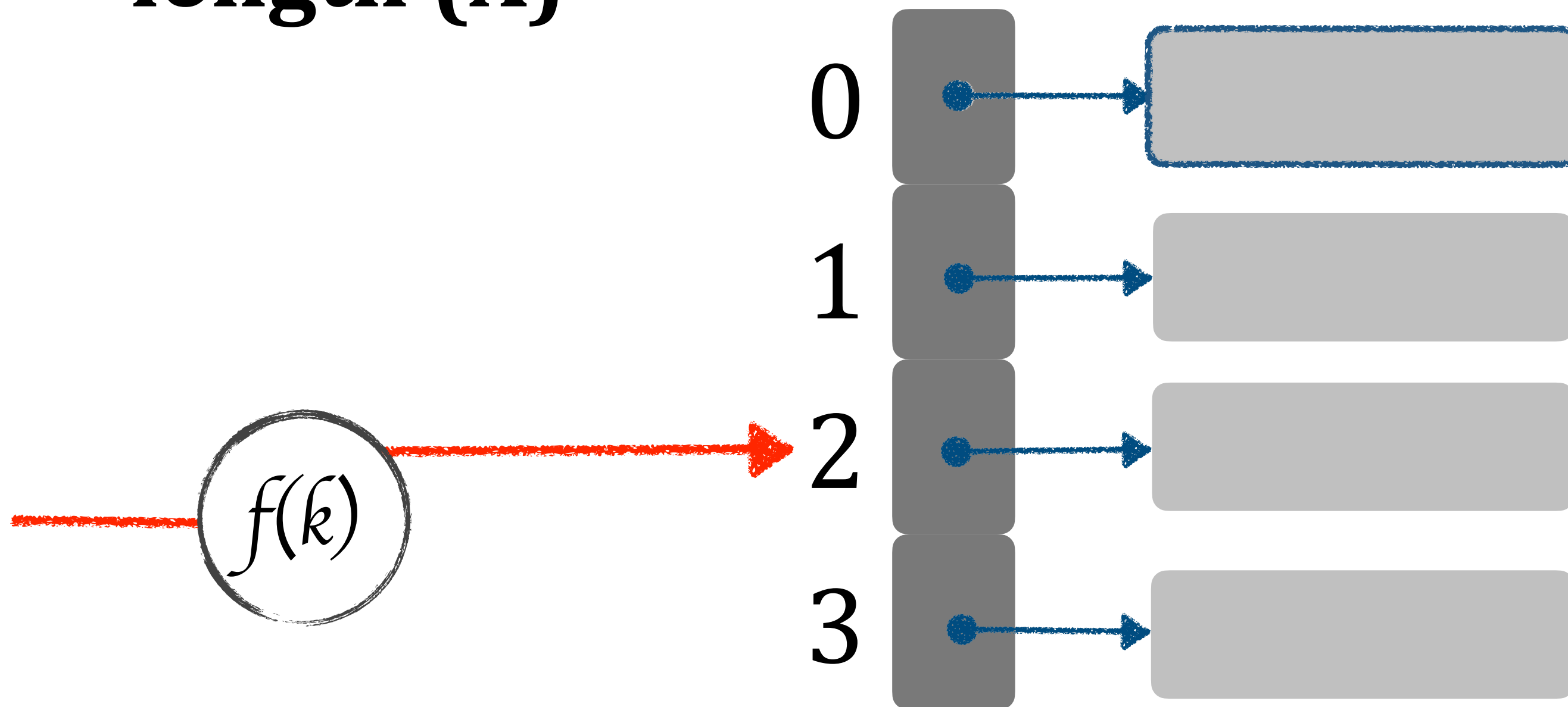
$$f(k)$$



N : entries in buffer
M : meta data

Implementation: Hash Hybrids

prefix
length (X)



key value
12301 abc
X = 3

bucket
count (H)

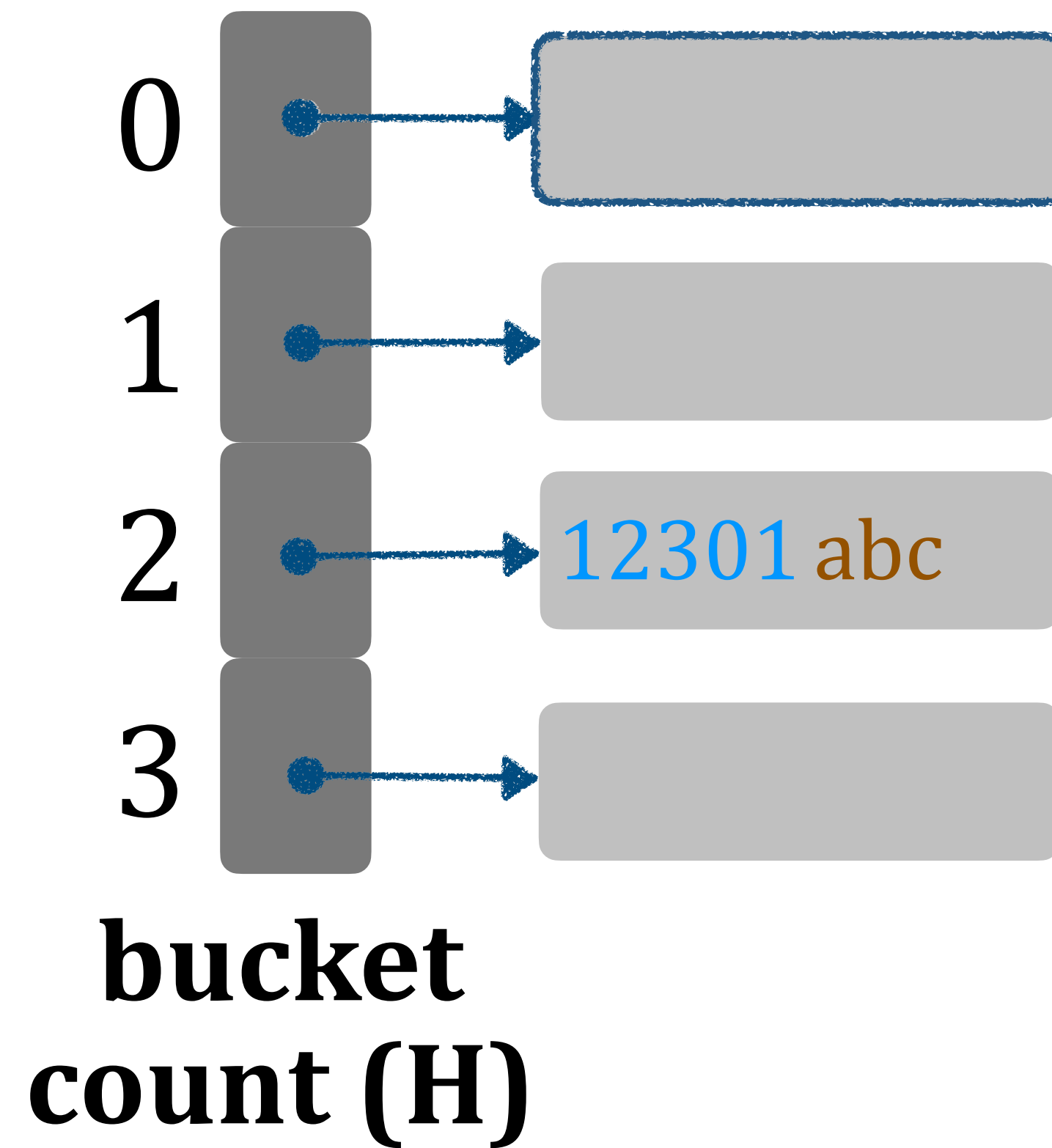


N : entries in buffer
M : meta data

Implementation: Hash Hybrids

prefix
length (X)

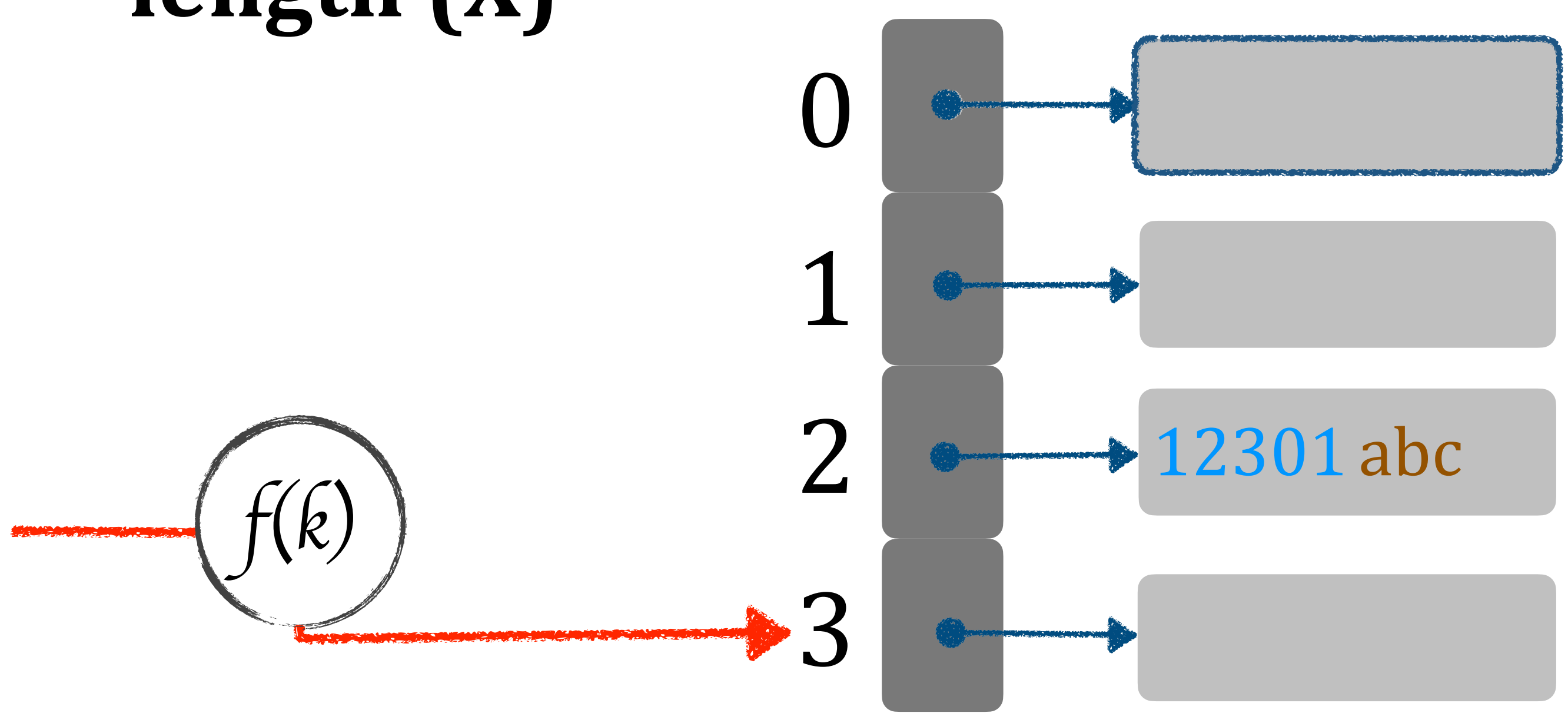
$$f(k)$$



N : entries in buffer
M : meta data

Implementation: Hash Hybrids

**prefix
length (X)**



| key | value | bucket |
|-------|-------|--------|
| 12401 | def | 3 |

X = 3

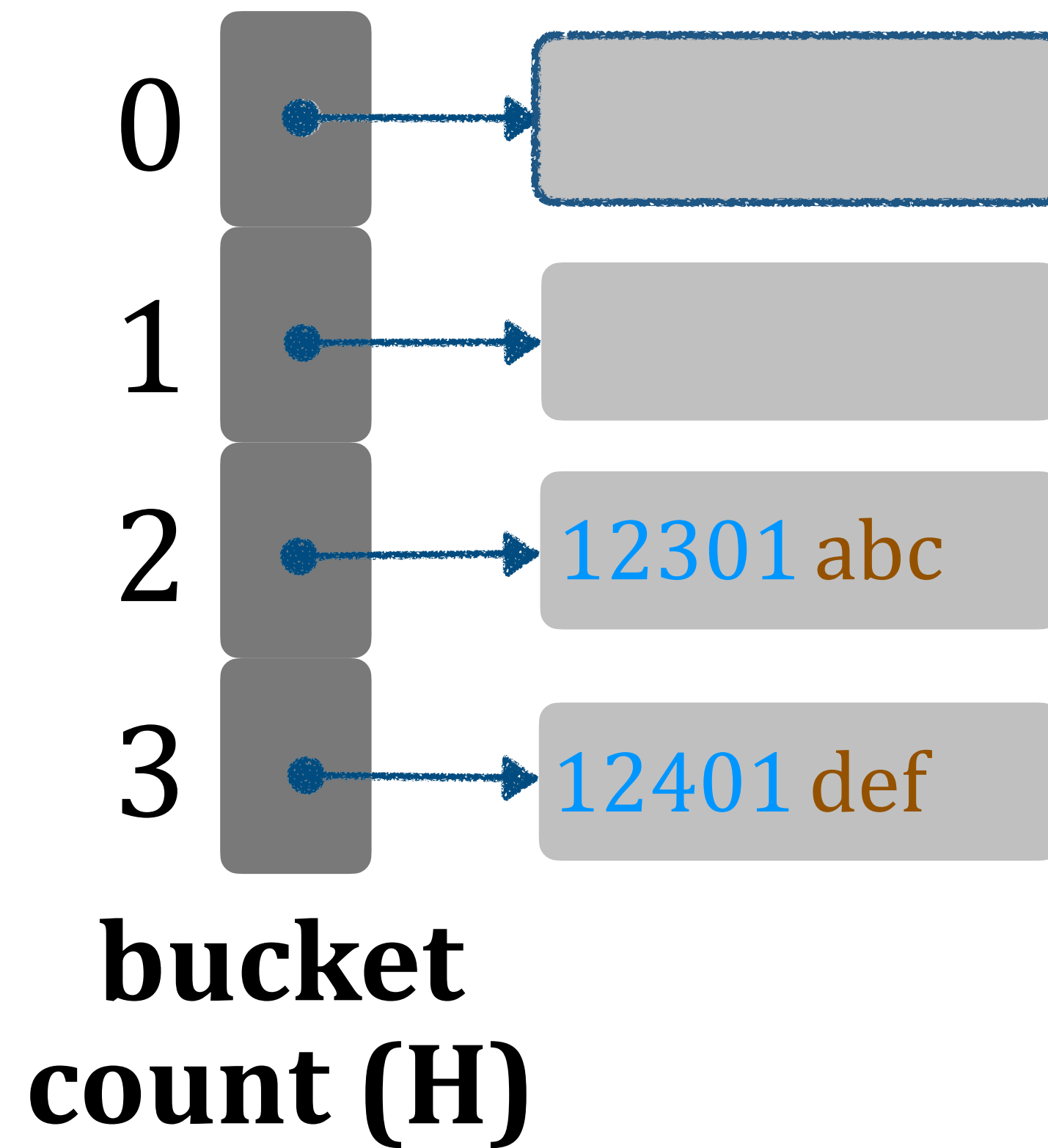


N : entries in buffer
M : meta data

Implementation: Hash Hybrids

prefix
length (X)

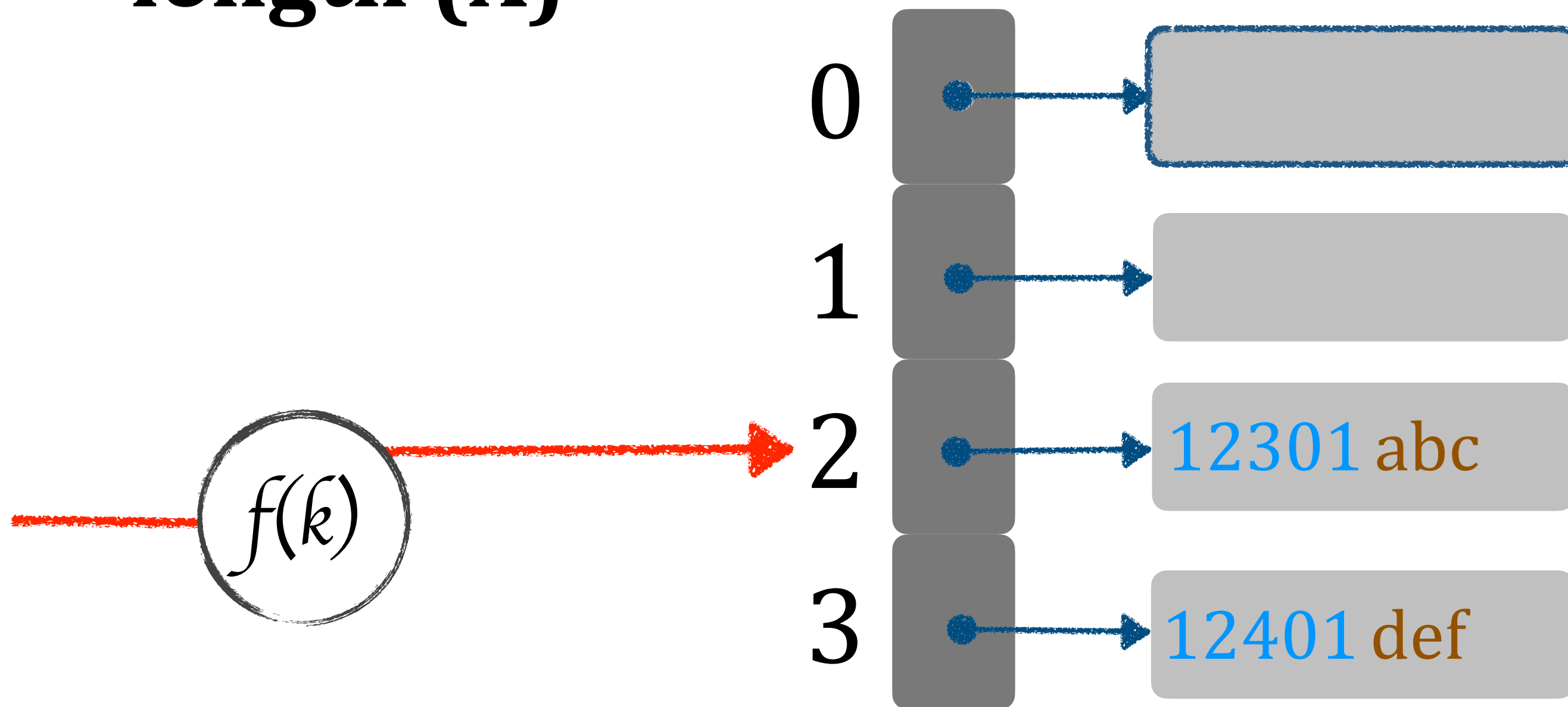
$$f(k)$$



N : entries in buffer
M : meta data

Implementation: Hash Hybrids

prefix
length (X)



key value
12311 xyz
X = 3

bucket
count (H)

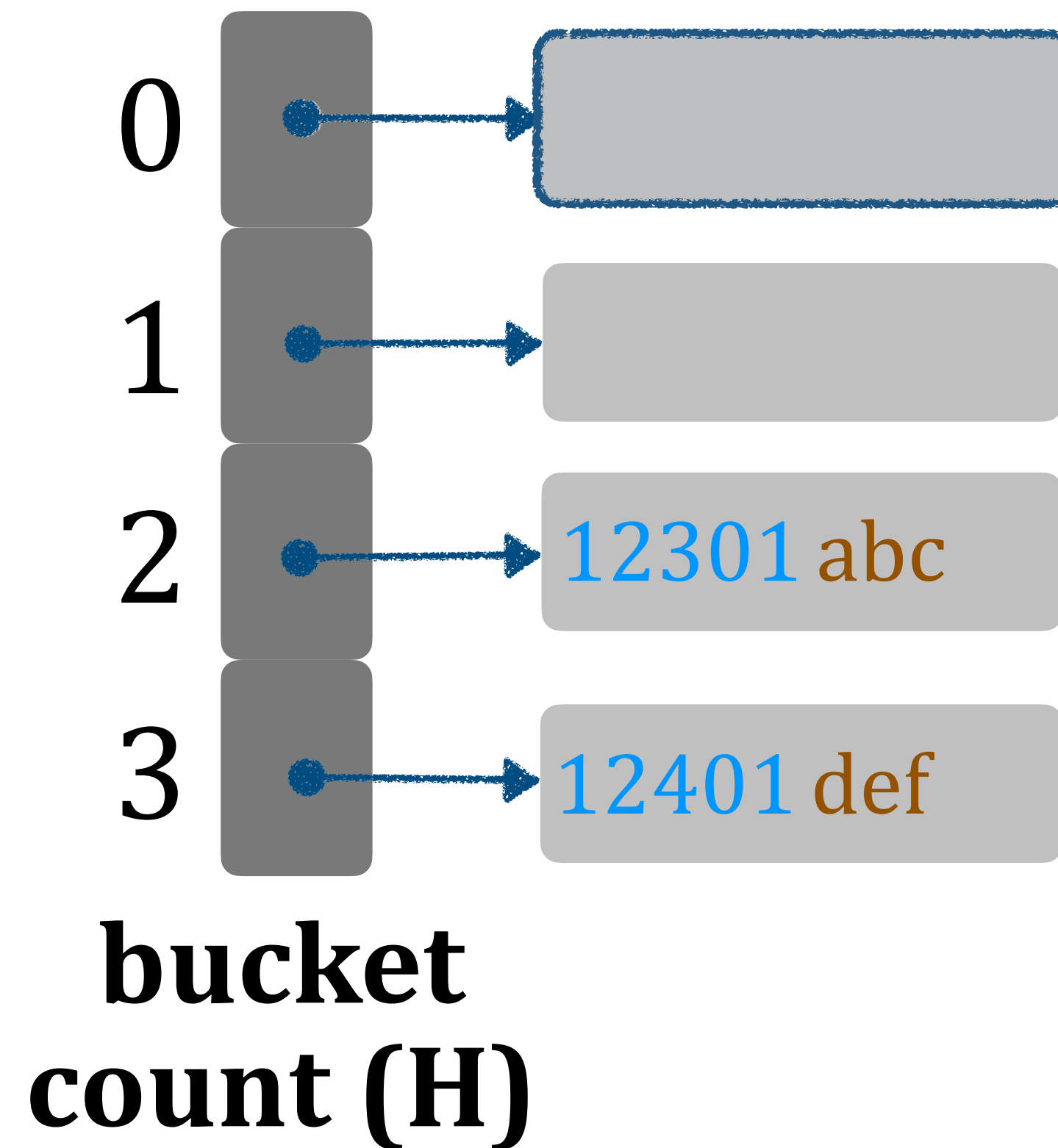


N : entries in buffer
M : meta data

Implementation: Hash Hybrids

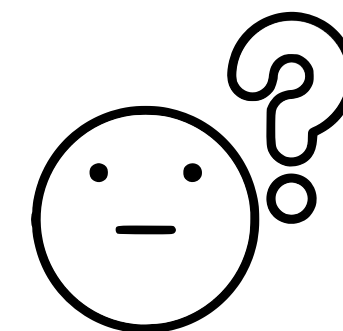
prefix
length (X)

$f(k)$



12311 xyz

collision



What do we do for overflow?

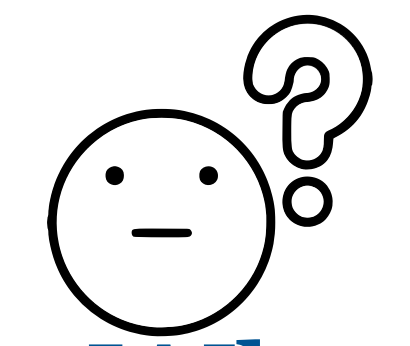
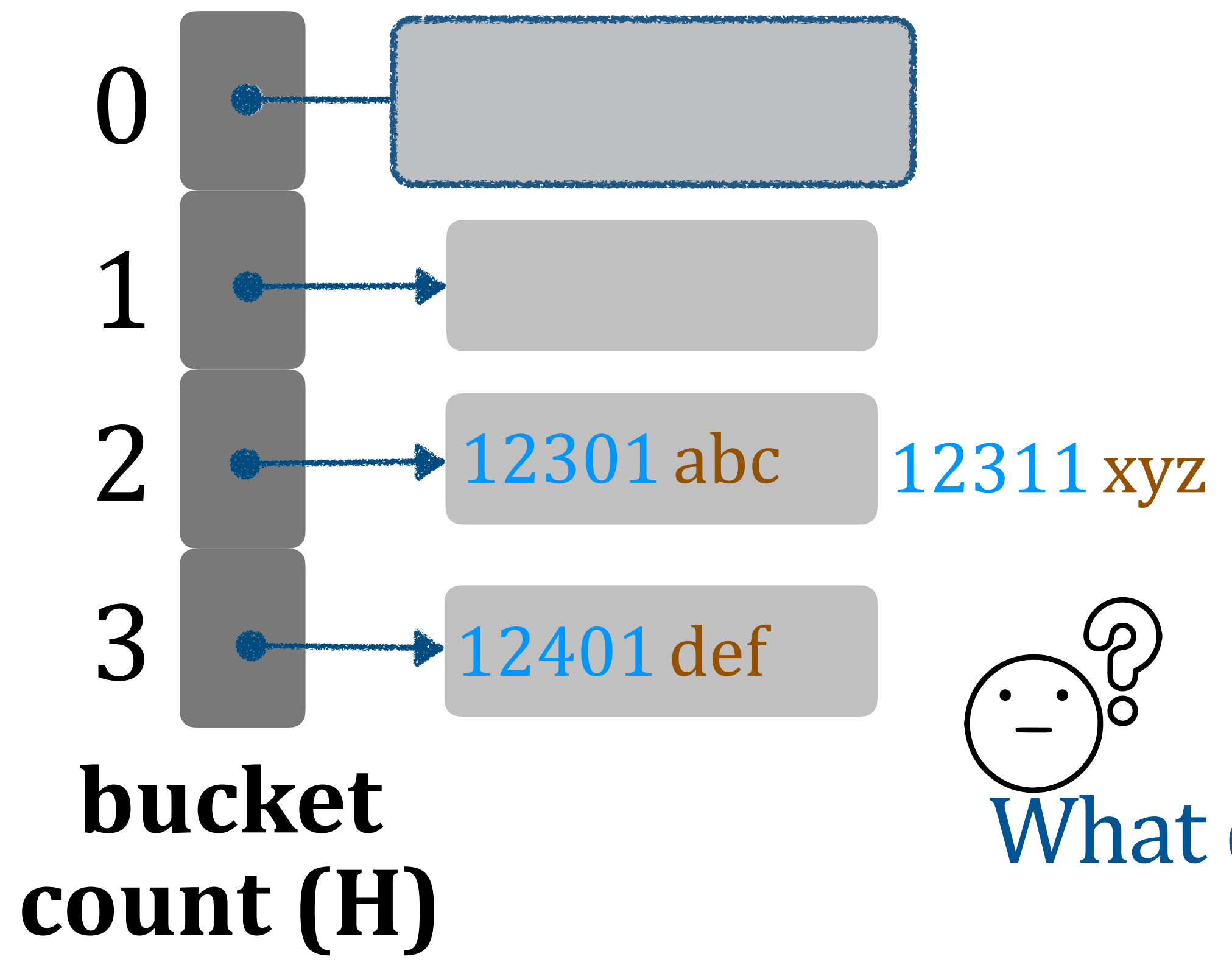
N : entries in buffer
M : meta data

Implementation: Hash Hybrids

prefix length (X)

linked-list
skip-list
vector

$$f(k)$$



What do we do for overflow?



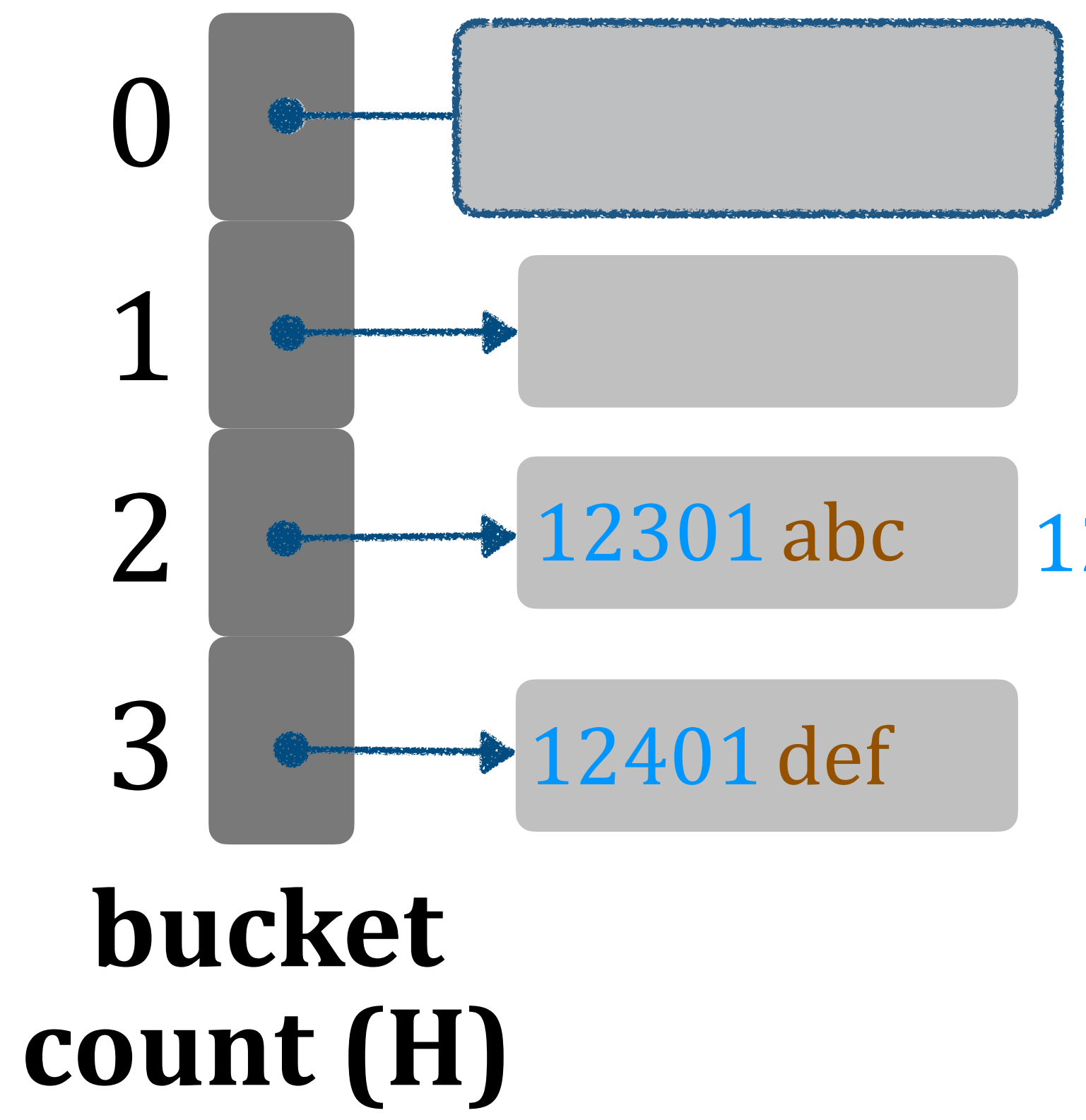
N : entries in buffer
M : meta data

Implementation: Hash Hybrids

prefix
length (X)

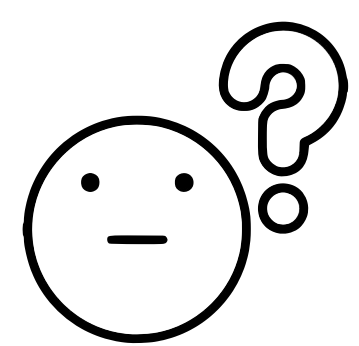
linked-list
skip-list
vector

$f(k)$



12311 xyz

collision



What if we don't store prefix?

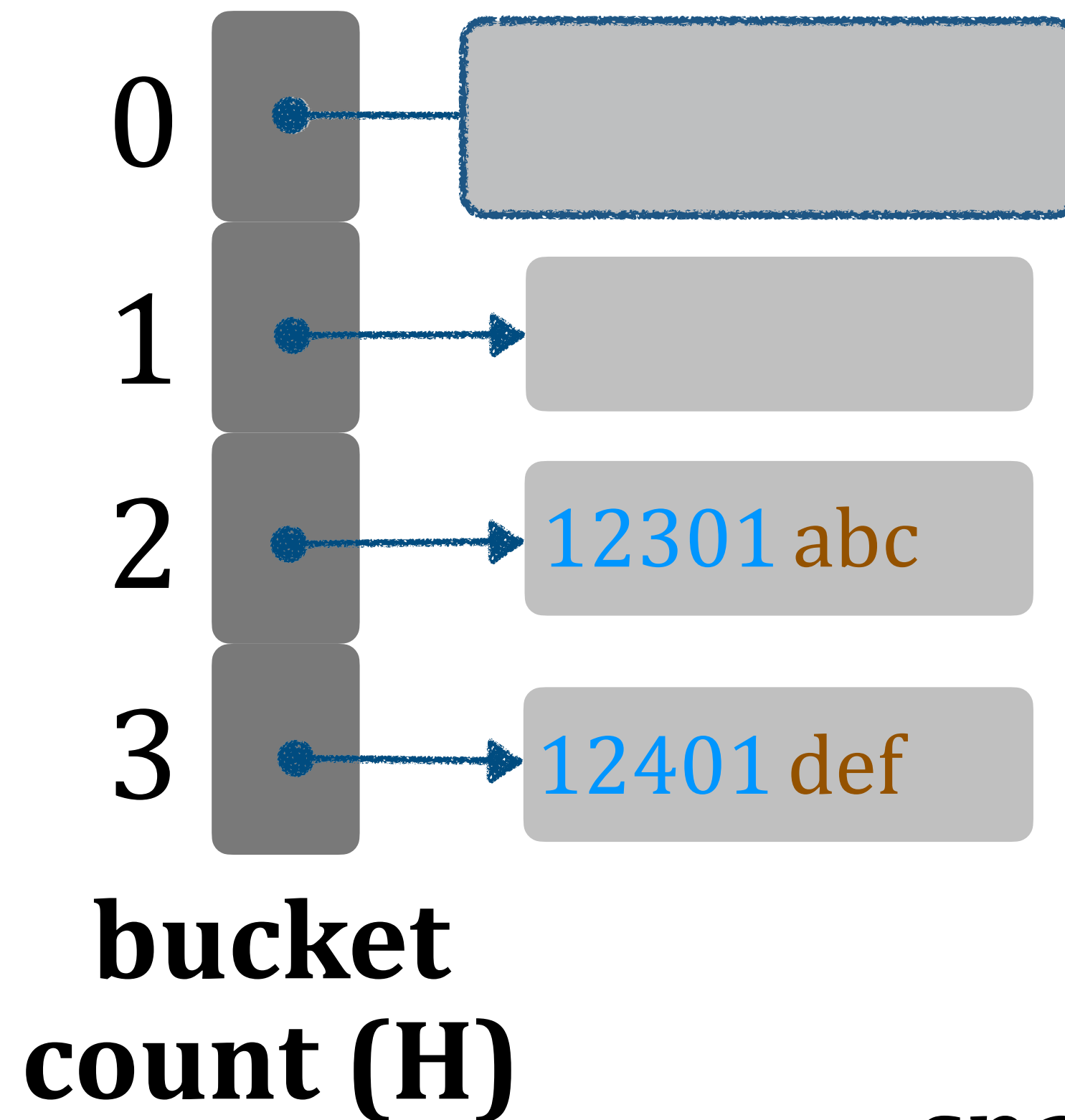


N : entries in buffer
M : meta data

Implementation: Hash Hybrids

prefix
length (X)

$f(k)$



linked-list
skip-list
vector

- great for mixed w/l
- more meta space needed
- good for point queries

insert cost: $\mathcal{O}(\log N/H)$

space complexity: $\mathcal{O}(N + M)$

point query cost: $\mathcal{O}(\log N/H)$

N : entries in buffer
M : meta data

Analyze: Hash Hybrids

prefix length (X)

bucket count (H)

any data structure

hash-table

$f(k_{i-N}) \rightarrow k_{iV}, \dots, k_{NV}$ ●

● $f(k_i) \rightarrow kv$

X = 0, H > 0

or

X = key size,
H = 1

○
○
○
 $f(k_N) \rightarrow kv$

X = key size,
H = N

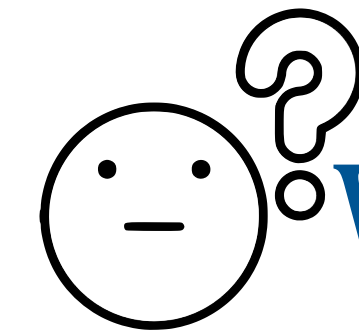


N : entries in buffer
M : meta data

Analyze: Hash Hybrids

prefix length (X)

bucket count (H)



What could be the possible values for X and H?

any data structure

$f(k_{i-N}) \rightarrow k_{iV}, \dots, k_{NV}$

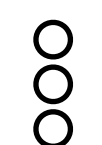
$X = 0, H > 0$

or

$X = \text{key size}, H = 1$

hash-table

$f(k_i) \rightarrow kv$



$f(k_N) \rightarrow kv$

$X = \text{key size}, H = N$



X ??
H ??

b is a finite set of prefixes $\{b_1, \dots, b_m\}$ for range queries

same prefix for start and end key



N : entries in buffer
M : meta data

Analyze: Hash Hybrids

prefix length (X)

bucket count (H)

 What could be the possible values for X and H?

any data structure

$f(k_{i-N}) \rightarrow k_{iV}, \dots, k_{NV}$ ●

$X = 0, H > 0$

or

$X = \text{key size}, H = 1$

hash-table

$f(k_i) \rightarrow kv$

⋮

$f(k_N) \rightarrow kv$

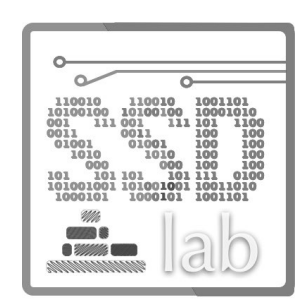
$X = \text{key size}, H = N$



$X = \text{len}(b_1)$
 $H > \text{len}(b)$

b is a finite set of prefixes $\{b_1, \dots, b_m\}$ for range queries

same prefix for start and end key

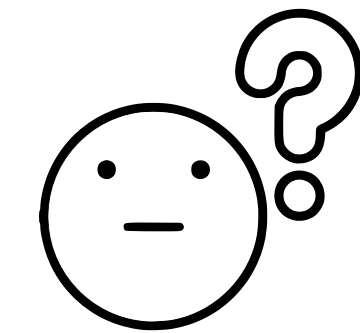


N : entries in buffer
M : meta data

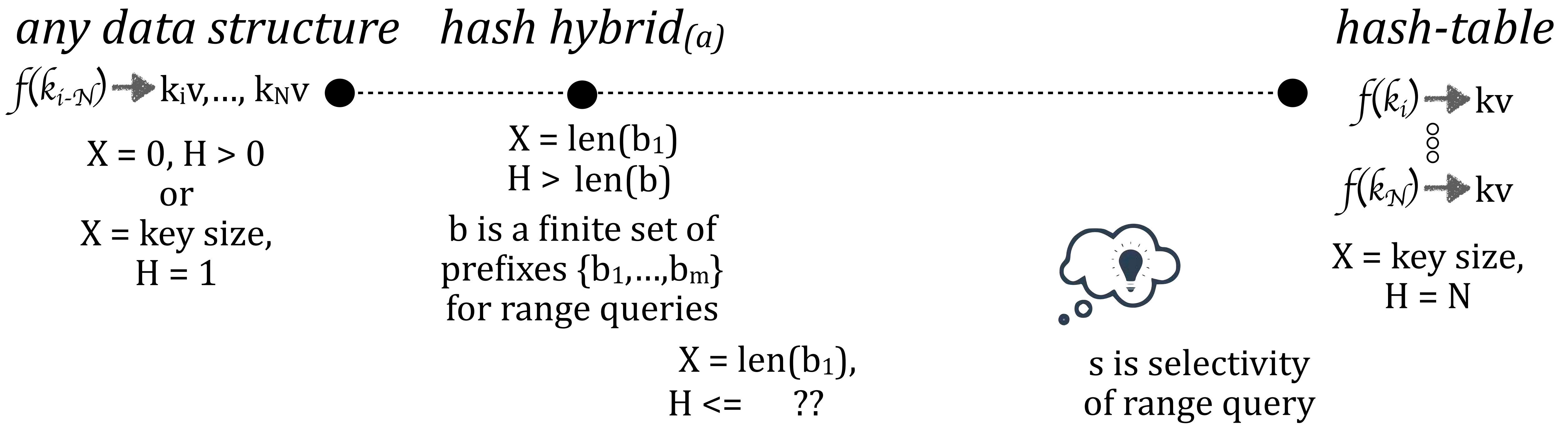
Analyze: Hash Hybrids

prefix length (X)

bucket count (H)



What could be the possible H value?



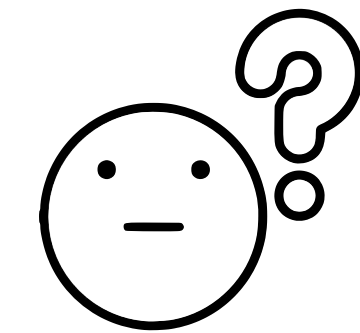
same prefix for start and end key

N : entries in buffer
M : meta data

Analyze: Hash Hybrids

prefix length (X)

bucket count (H)



What could be the possible H value?

any data structure

hash hybrid_(a)

hash-table

$f(k_{i-N}) \rightarrow k_{iV}, \dots, k_{NV}$



$f(k_i) \rightarrow kv$

\vdots
 $f(k_N) \rightarrow kv$

$X = 0, H > 0$

or

$X = \text{key size}, H = 1$

$X = \text{len}(b_1)$

$H > \text{len}(b)$

b is a finite set of prefixes $\{b_1, \dots, b_m\}$ for range queries



$X = \text{len}(b_1),$

$H \leq \text{size}(\text{buffer}) / (N \times s)$

s is selectivity of range query

same prefix for start and end key

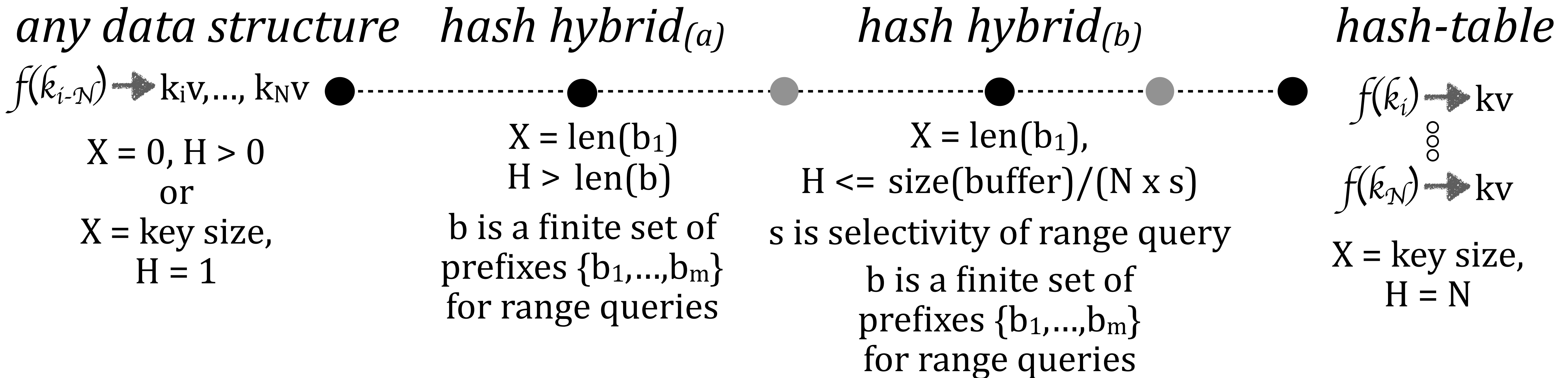


N : entries in buffer
M : meta data

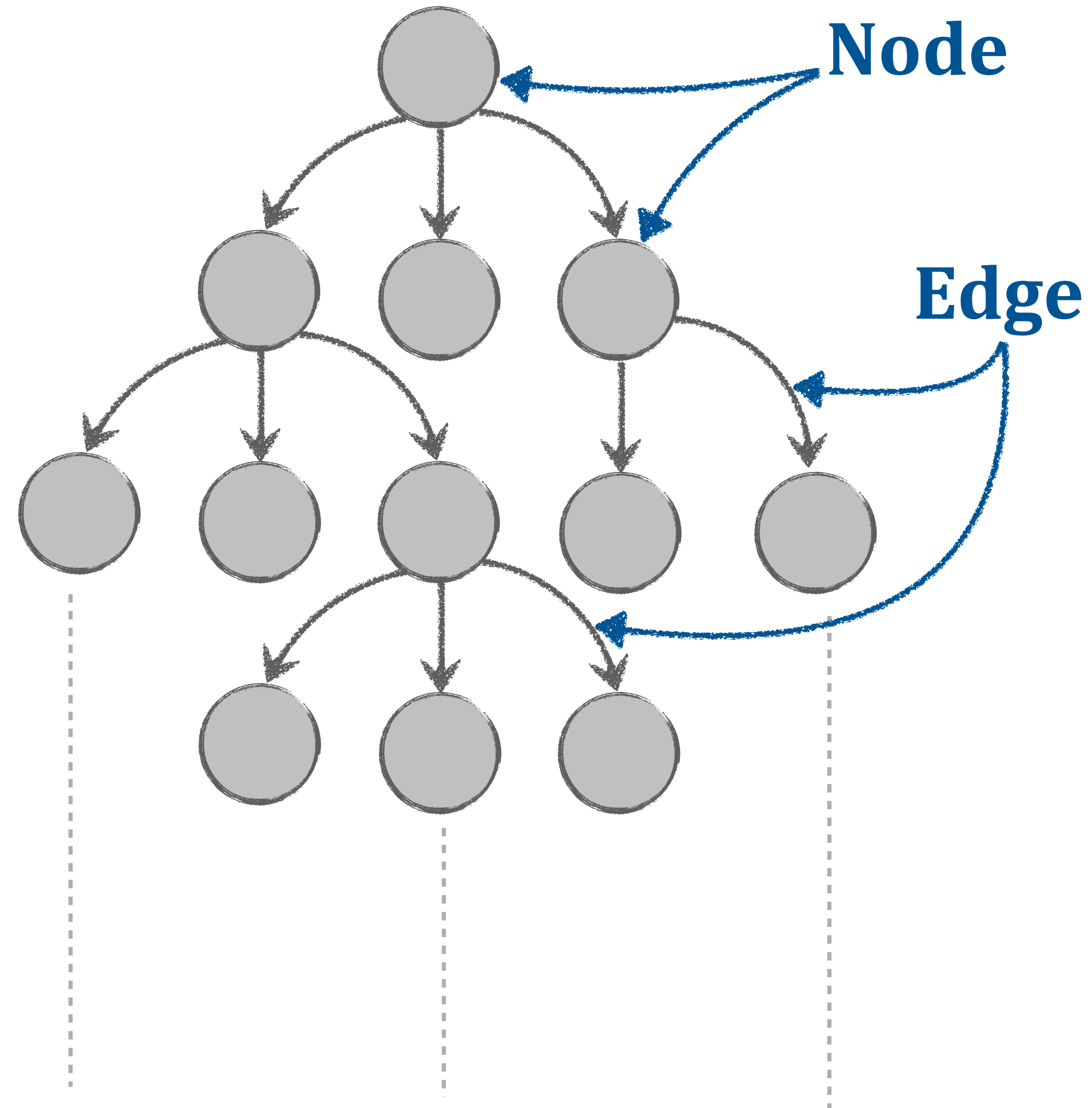
Analyze: Hash Hybrids

prefix length (X)

bucket count (H)

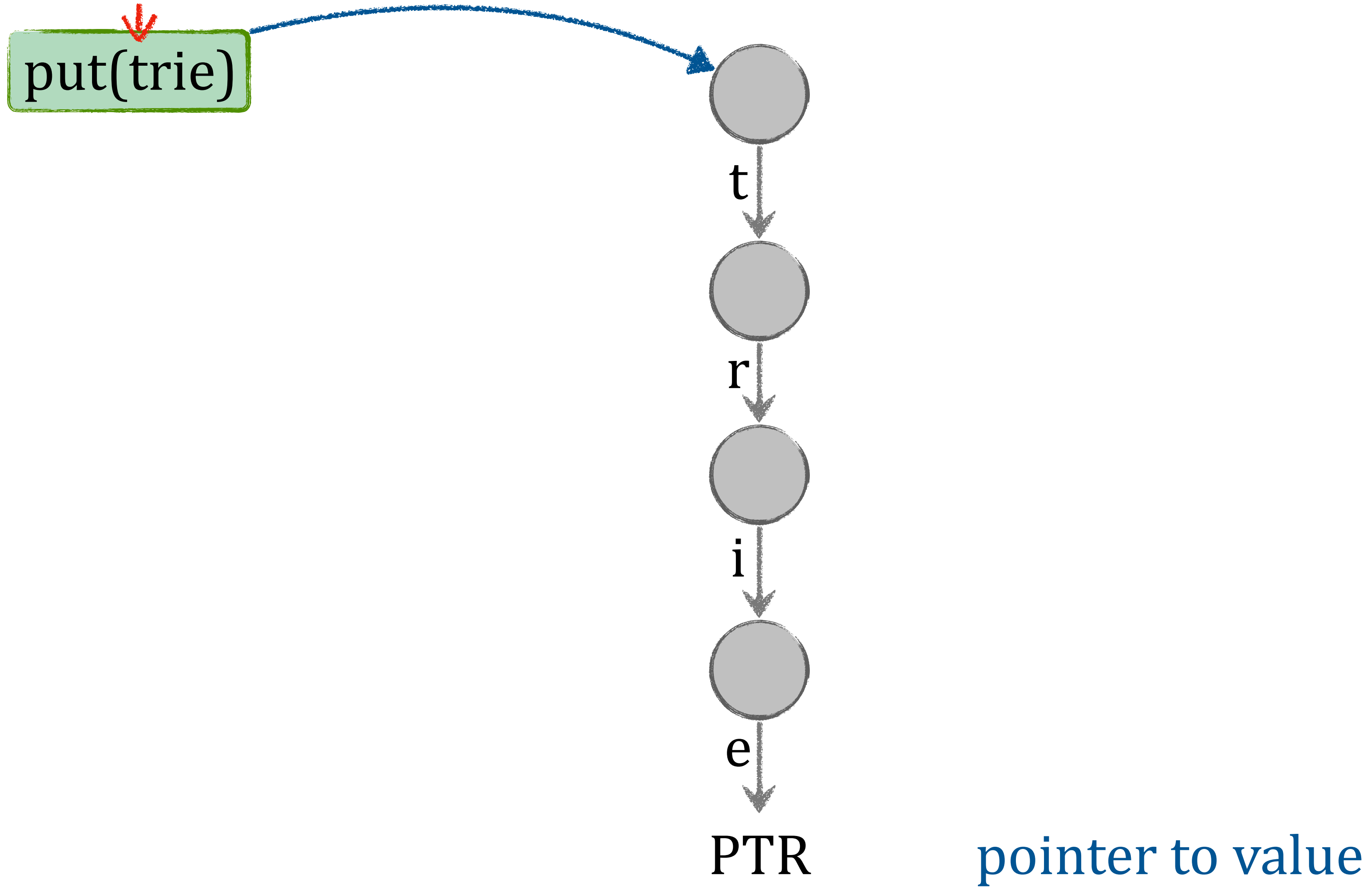


Implementation: Trie



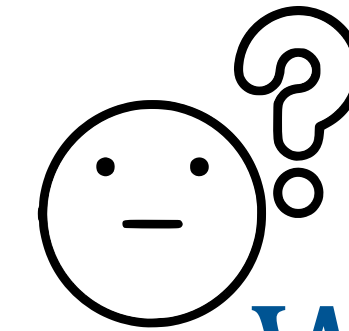
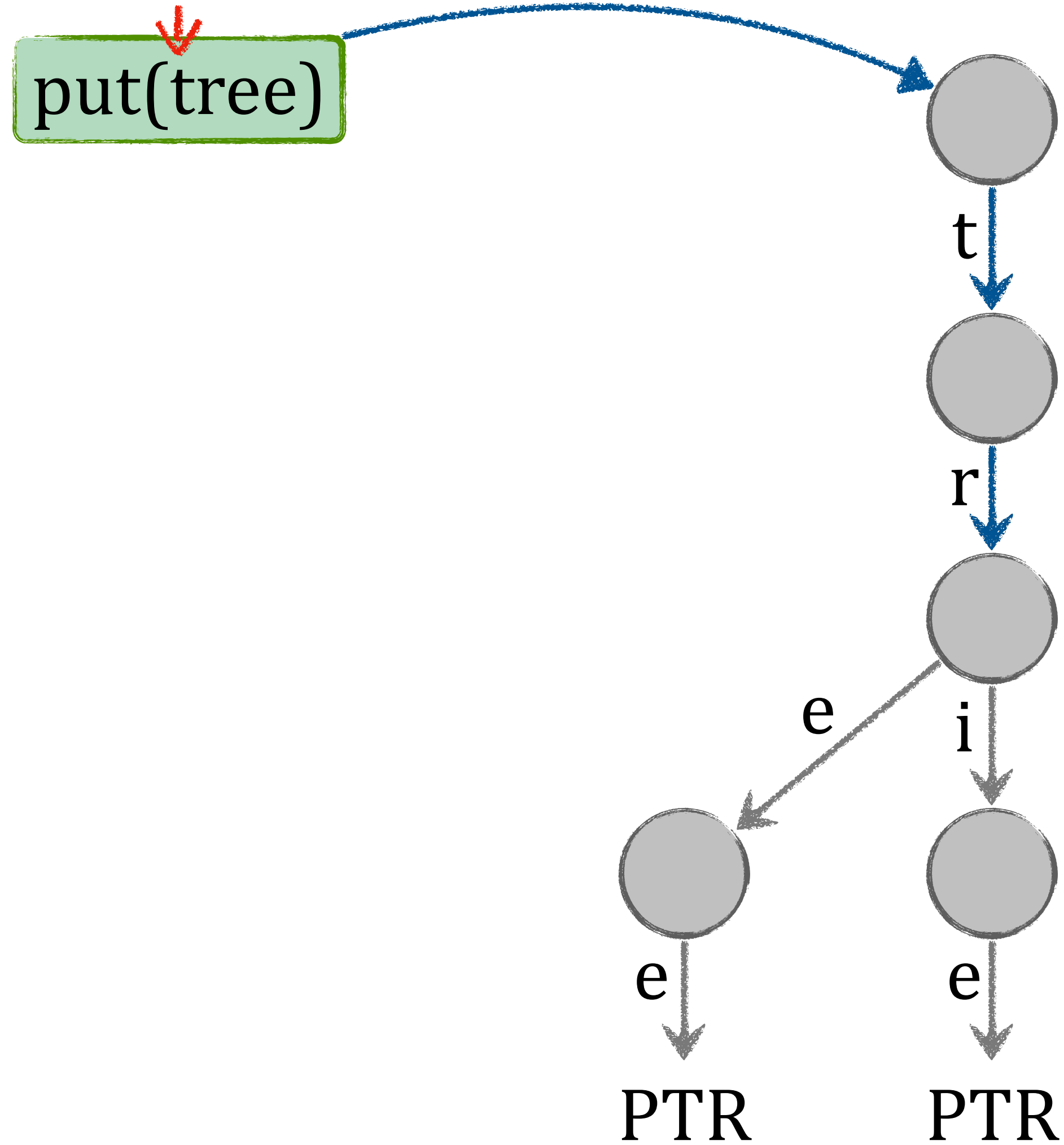
K : max size of key
M : meta data
N : entries in buffer

Put Implementation: Trie



K : max size of key
M : meta data
N : entries in buffer

Put Implementation: Trie



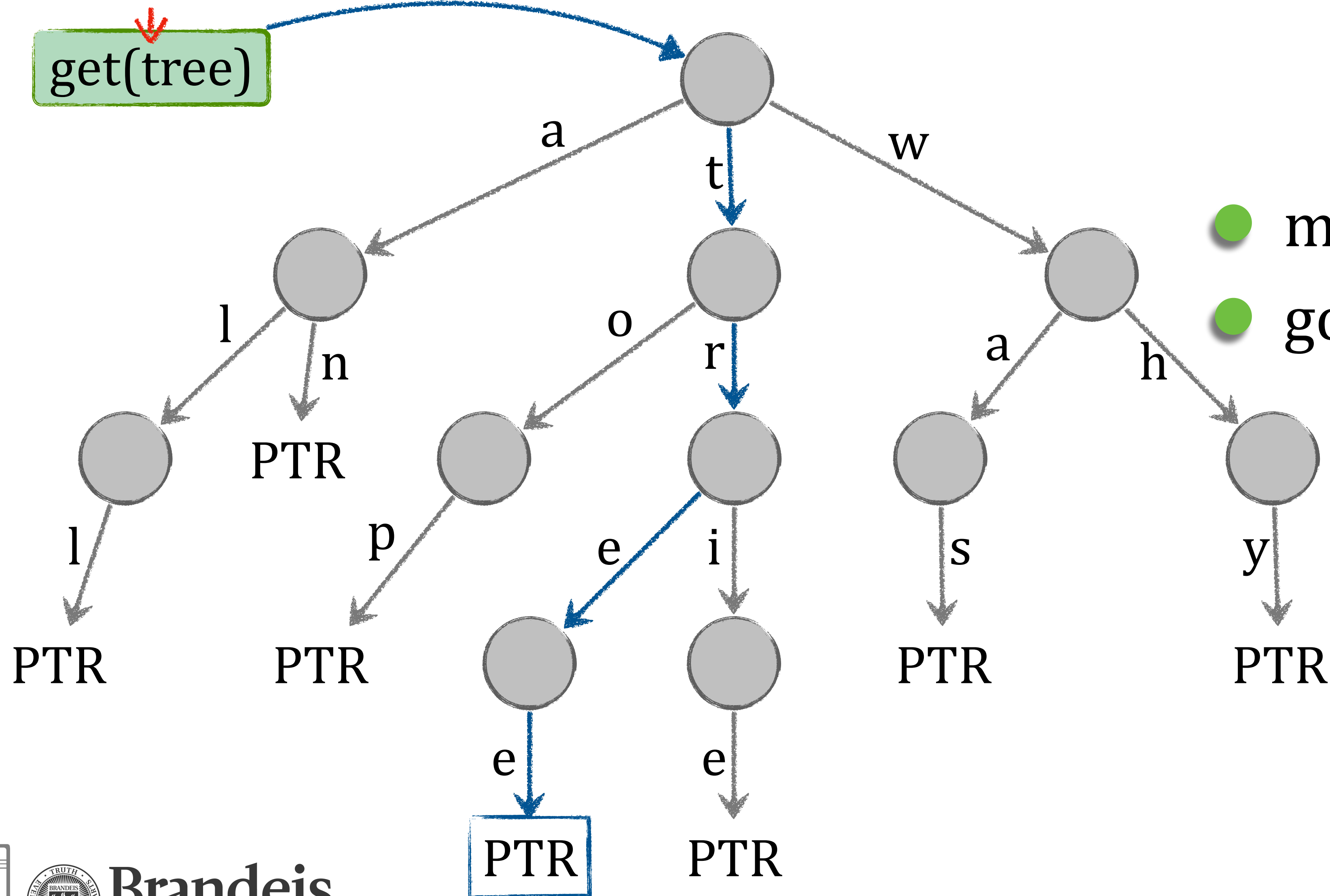
What's the difference?

- more space optimized



K : max size of key
M : meta data
N : entries in buffer

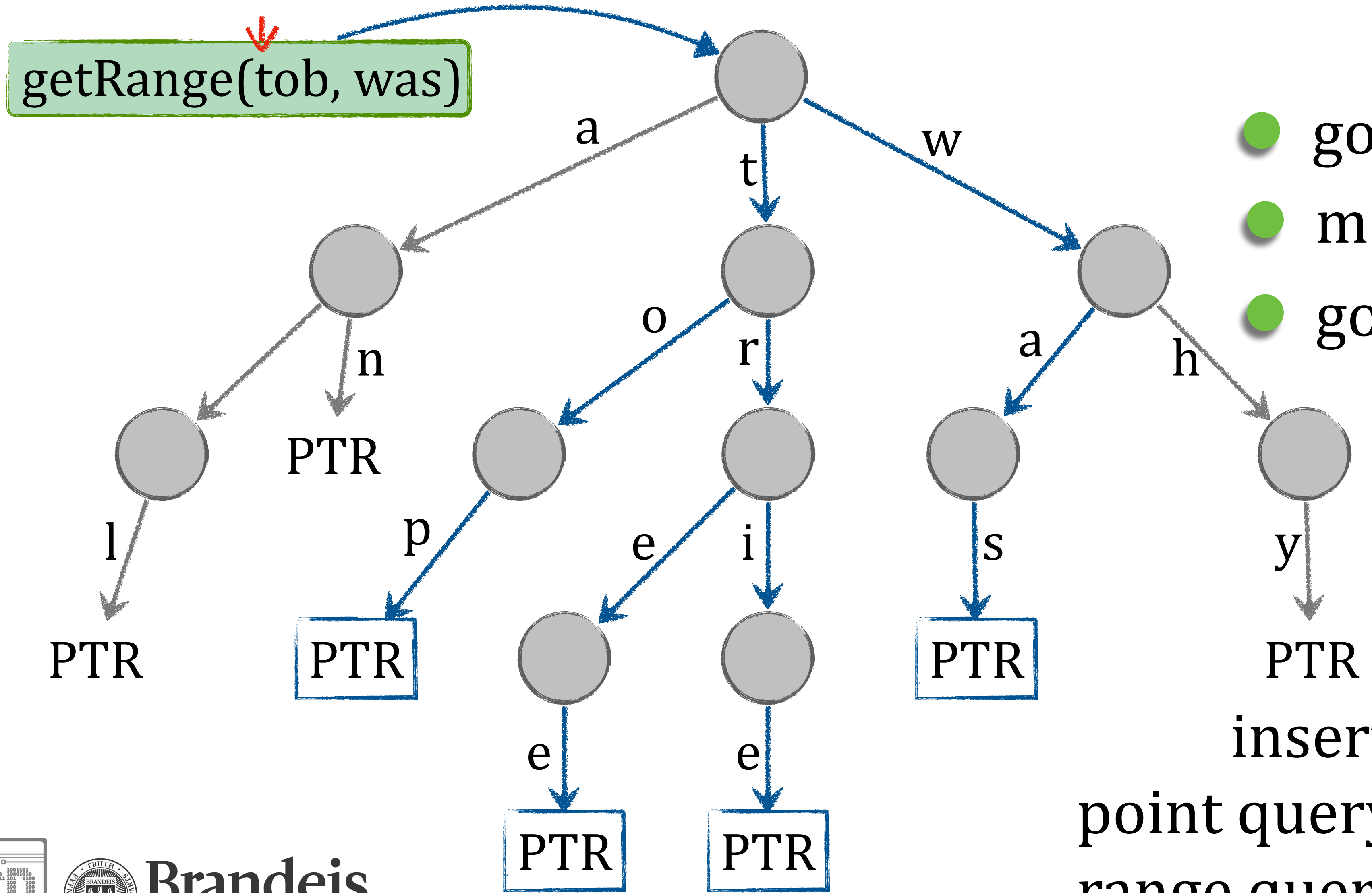
Get Implementation: Trie



- more space optimized
- good for point queries

K : max size of key
M : meta data
N : entries in buffer

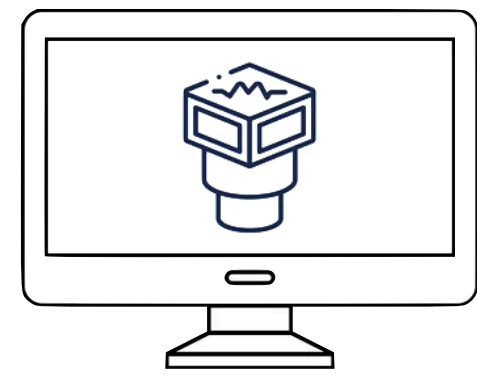
GetRange Implementation: Trie



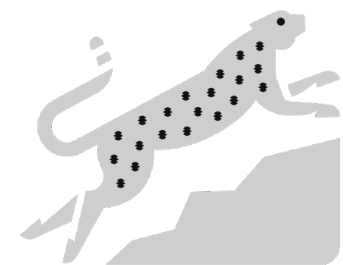
- good for range queries
- more space optimized
- good for point queries

insert cost: $\mathcal{O}(K)$
point query cost: $\mathcal{O}(K)$
range query cost: $\mathcal{O}(K \times N \times s)$

Experiments



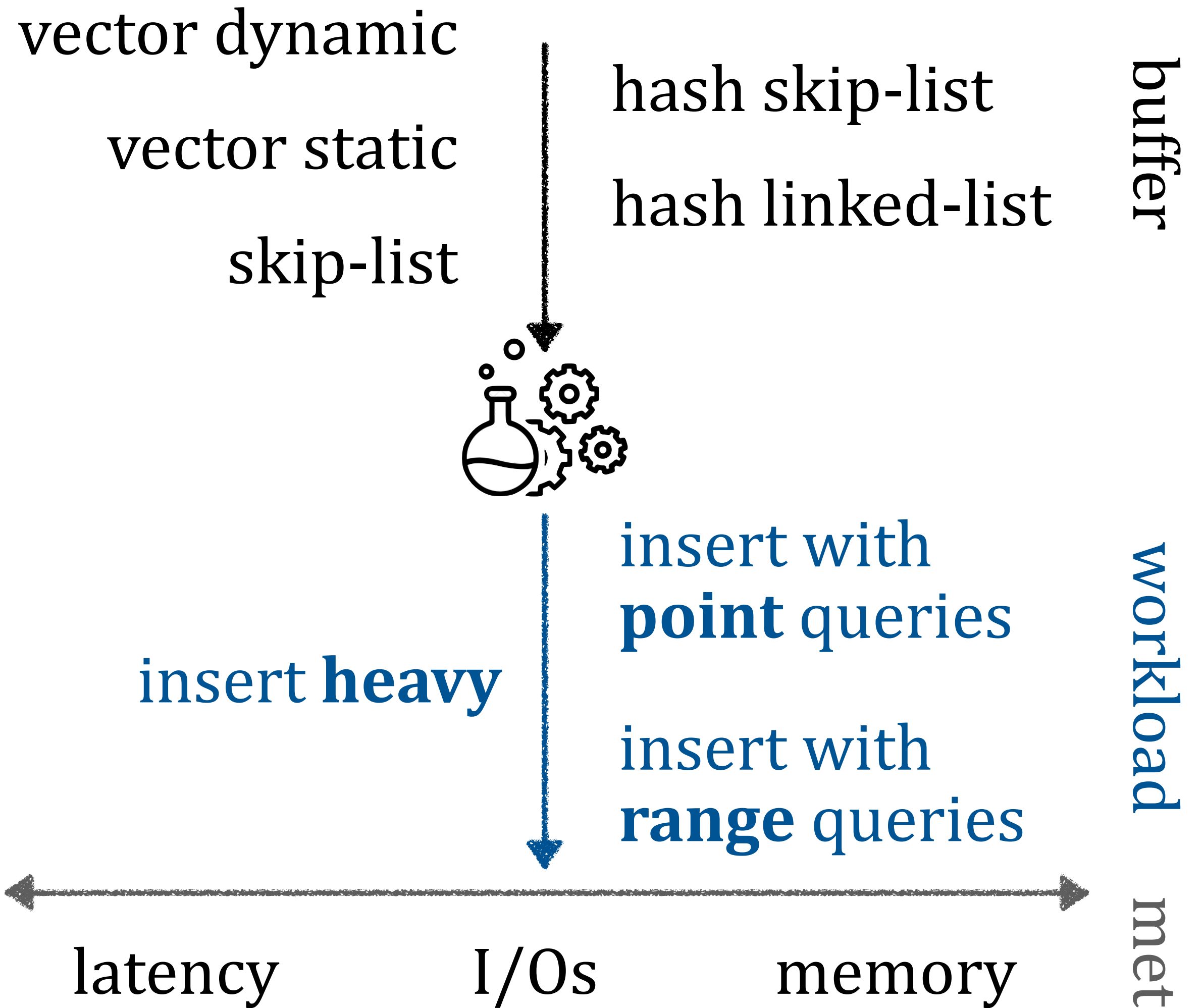
2 Intel Gold 6126 vCPUs
192 GB RAM
240 GB SSD
Ubuntu 20.04 LTS



RocksDB v9.0.0

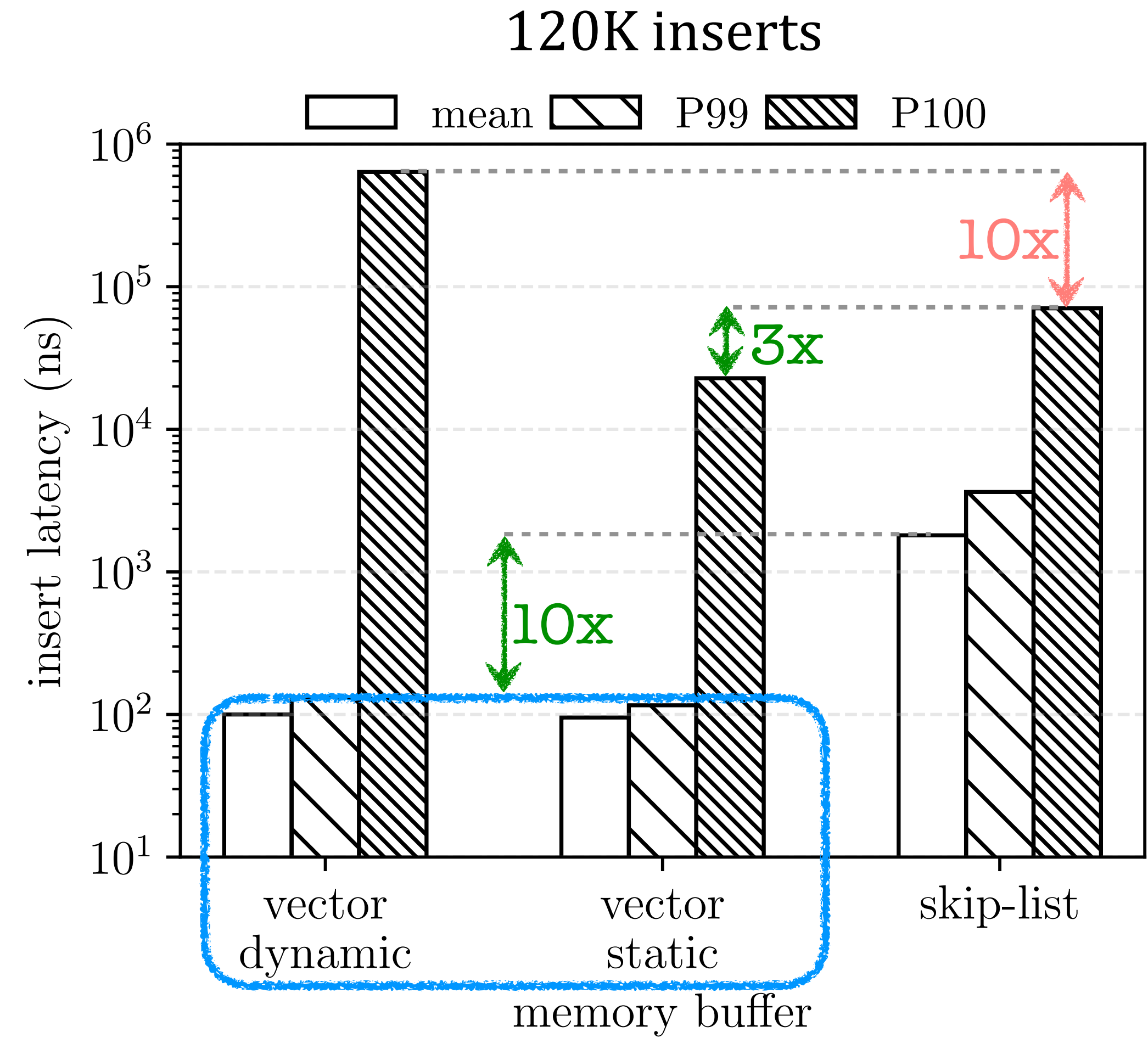


| | |
|-------------|------|
| size ratio | 4 |
| buffer size | 16MB |
| page size | 4KB |
| entry size | 64B |
| key size | 10 |



Evaluation

vector is 10x faster than skip-list for an *insert heavy* w/l



insert heavy

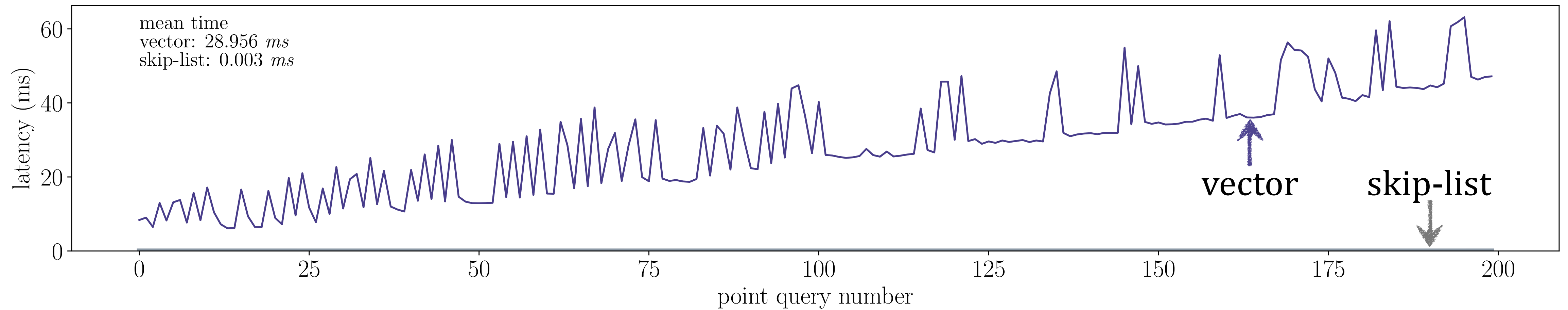


Evaluation

vector is 10x faster than skip-list for an *insert heavy* w/l



140K inserts, 200 interleaved point queries




insert with point queries

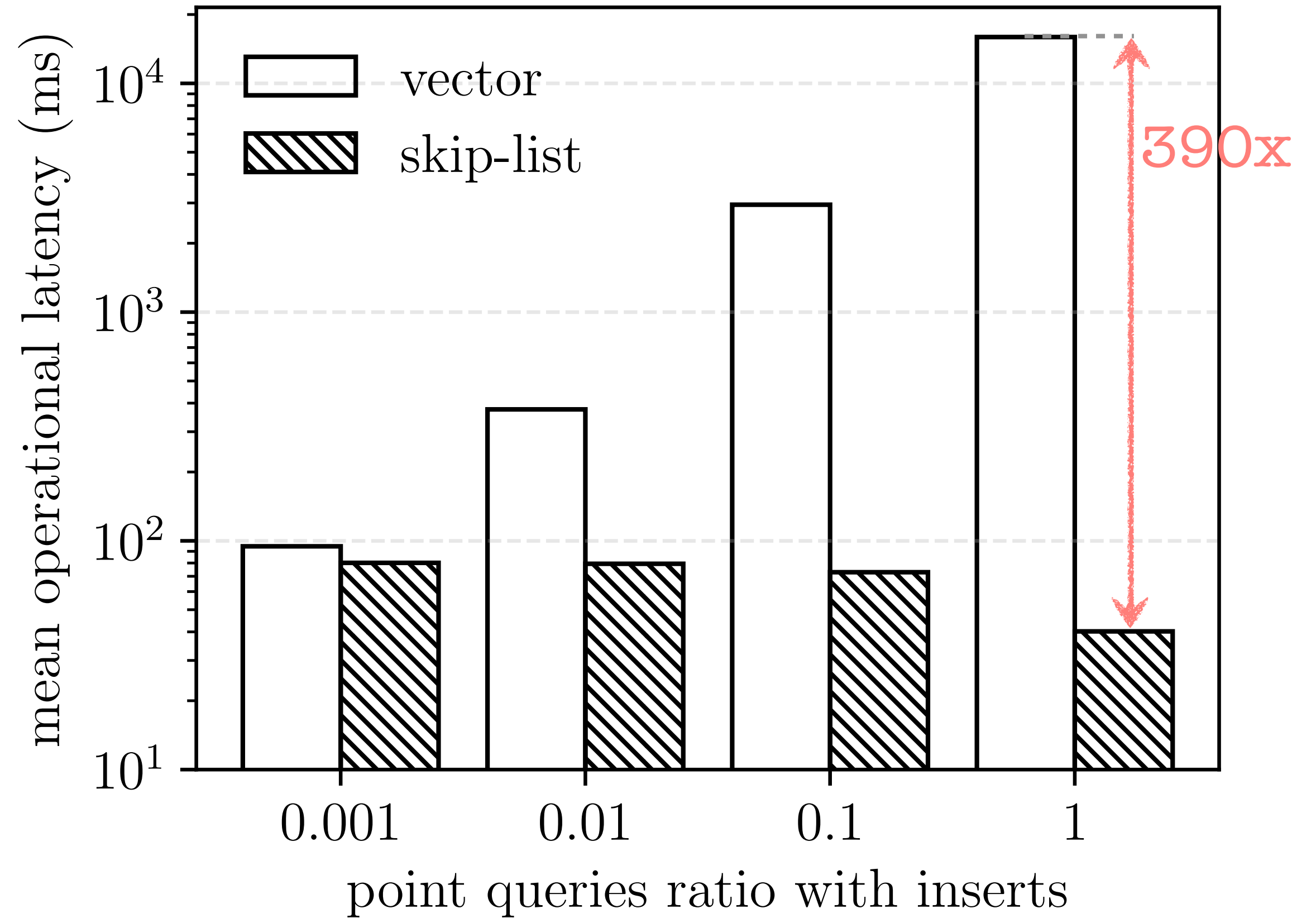


Evaluation

vector is 10x faster than skip-list for an *insert heavy* w/l 

vector is worst choice for w/l's with point queries 

120K inserts, interleaved point queries



insert with point queries

Evaluation

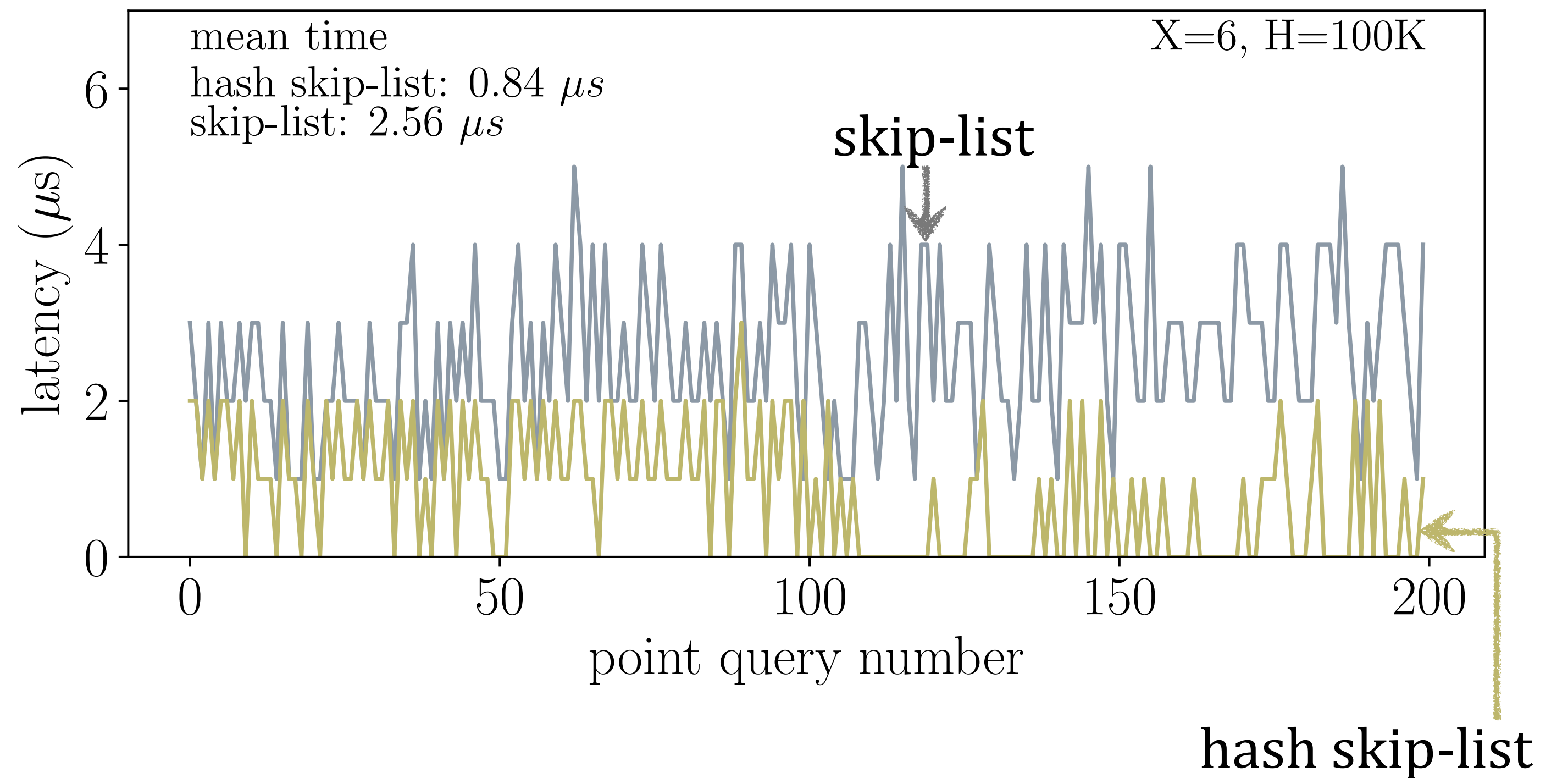
vector is 10x faster than skip-list for an *insert heavy* w/l



vector is worst choice for w/l's with point queries



140K inserts, 200 interleaved point queries





insert with point queries




Brandeis
UNIVERSITY

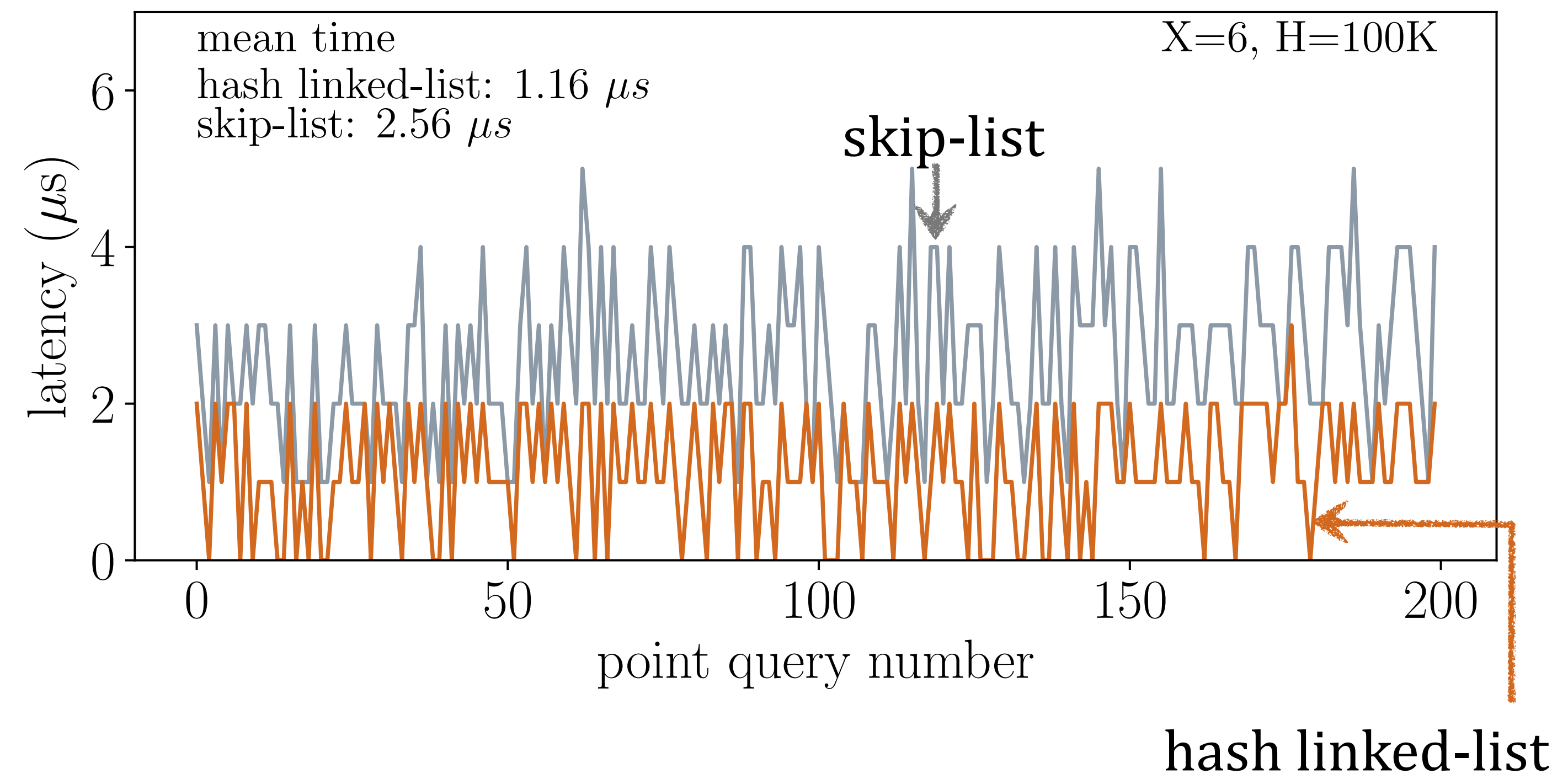
Evaluation

vector is 10x faster than skip-list for an *insert heavy* w/l 

vector is worst choice for w/l's with point queries 

hash buffers are 2-3x faster than skip-list for point queries w/l's 

140K inserts, 200 interleaved point queries




insert with point queries

Evaluation

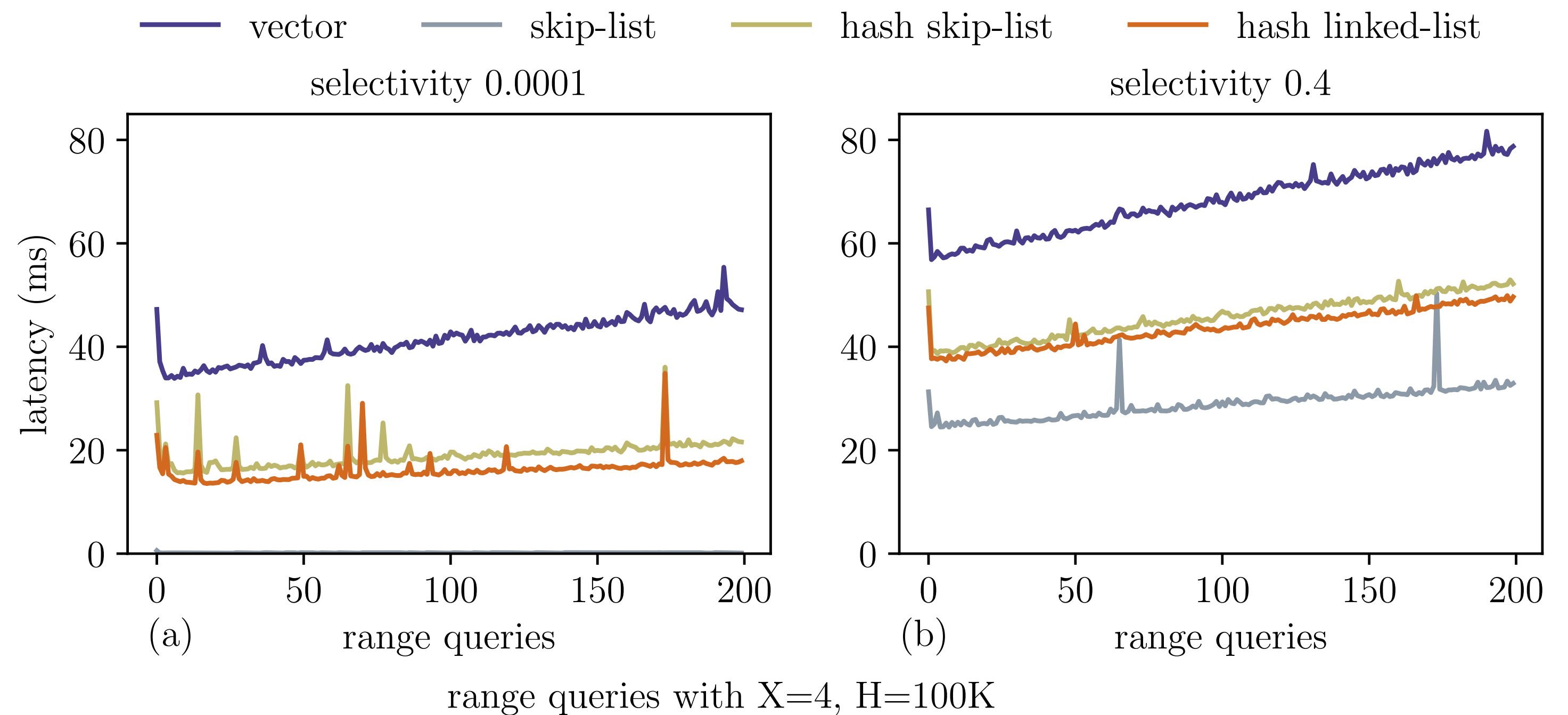
vector is 10x faster than skip-list for an *insert heavy* w/l 

vector is worst choice for w/ls with point queries 

hash buffers are 2-3x faster than skip-list for point queries w/ls 

skip-list outperforms for w/ls with range queries 

140K inserts, 200 interleaved range queries



insert with range queries



Evaluation

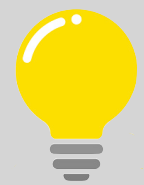
vector is 10x faster than skip-list for an *insert heavy* w/l



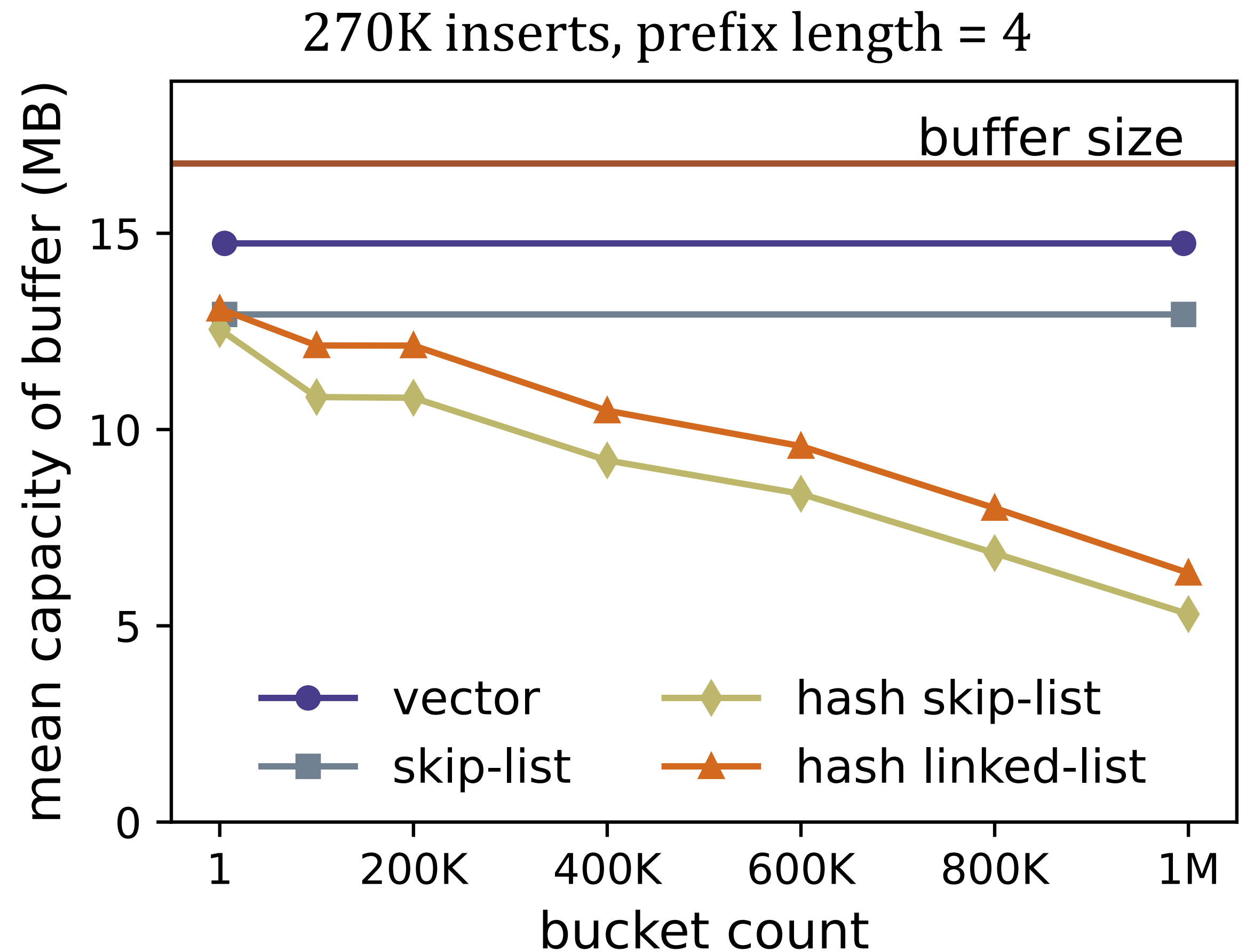
vector is worst choice for w/l's with point queries



hash buffers are 2-3x faster than skip-list for point queries w/l's



skip-list outperforms for w/l's with range queries




memory footprint



Evaluation

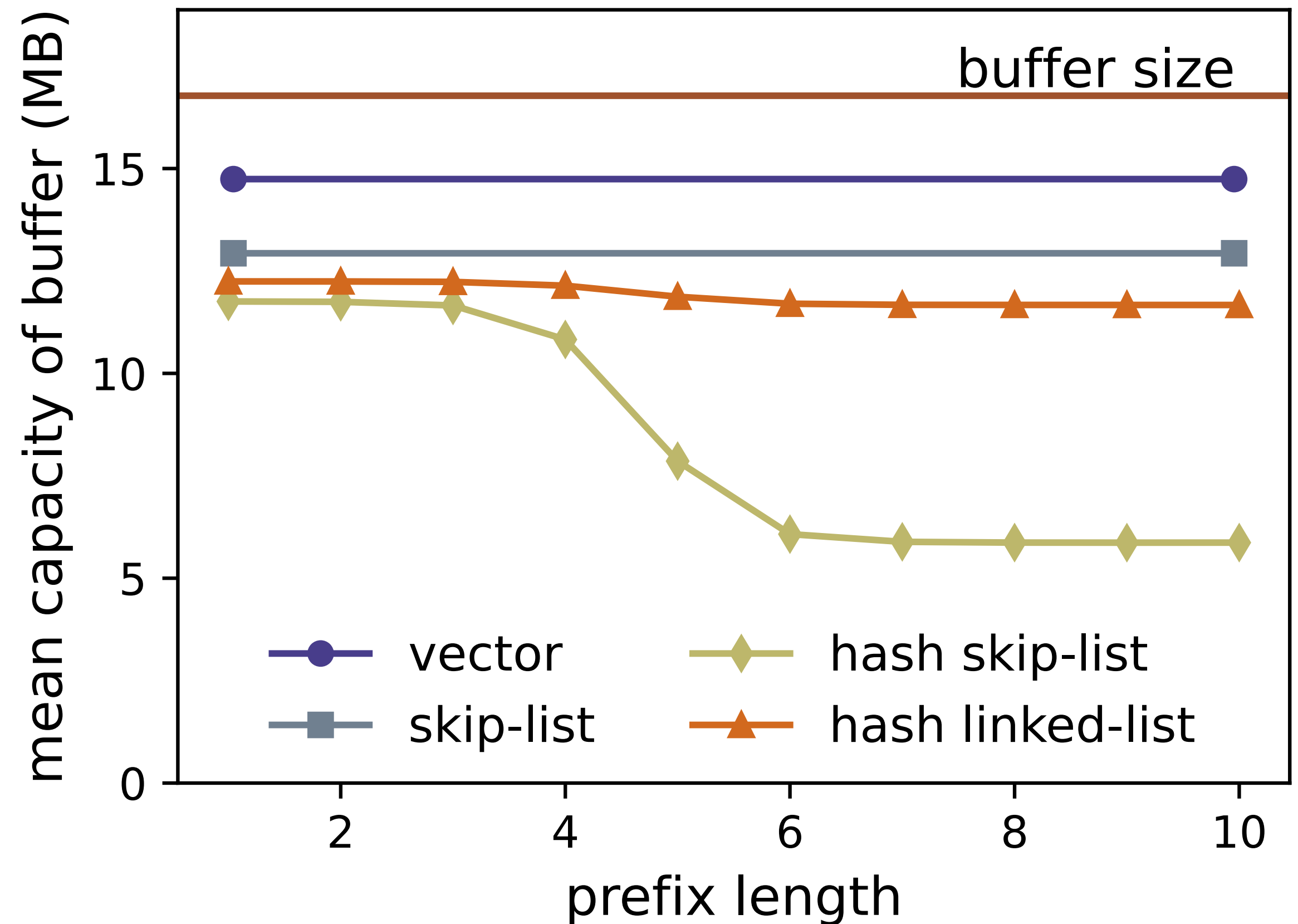
vector is 10x faster than skip-list for an *insert heavy* w/l 

vector is worst choice for w/l's with point queries 

hash buffers are 2-3x faster than skip-list for point queries w/l's 

skip-list outperforms for w/l's with range queries 

270K inserts, bucket count = 100K



memory footprint



Evaluation

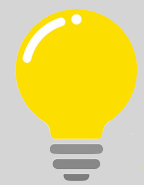
vector is 10x faster than skip-list for an *insert heavy* w/l



vector is worst choice for w/l's with point queries



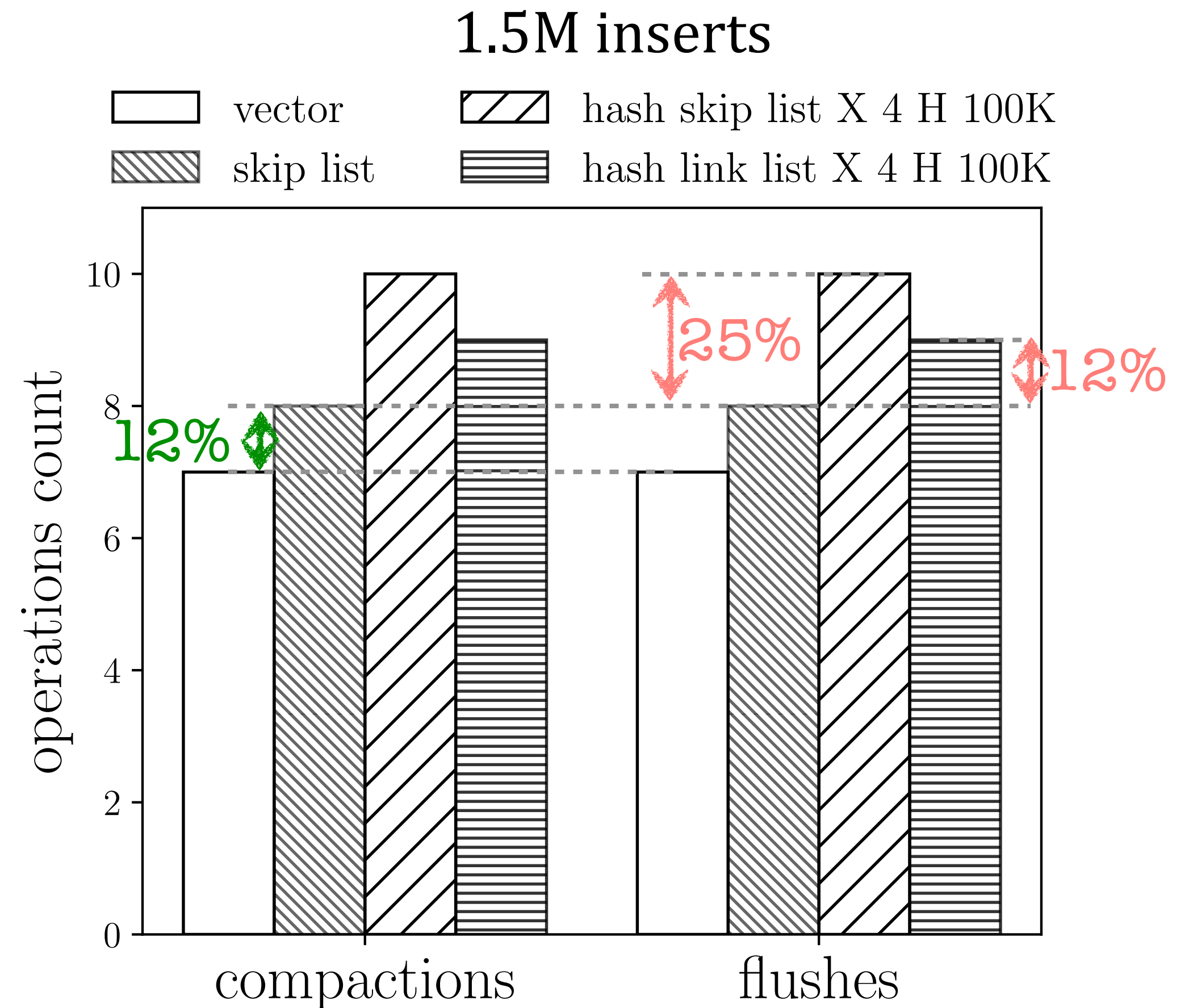
hash buffers are 2-3x faster than skip-list for point queries w/l's



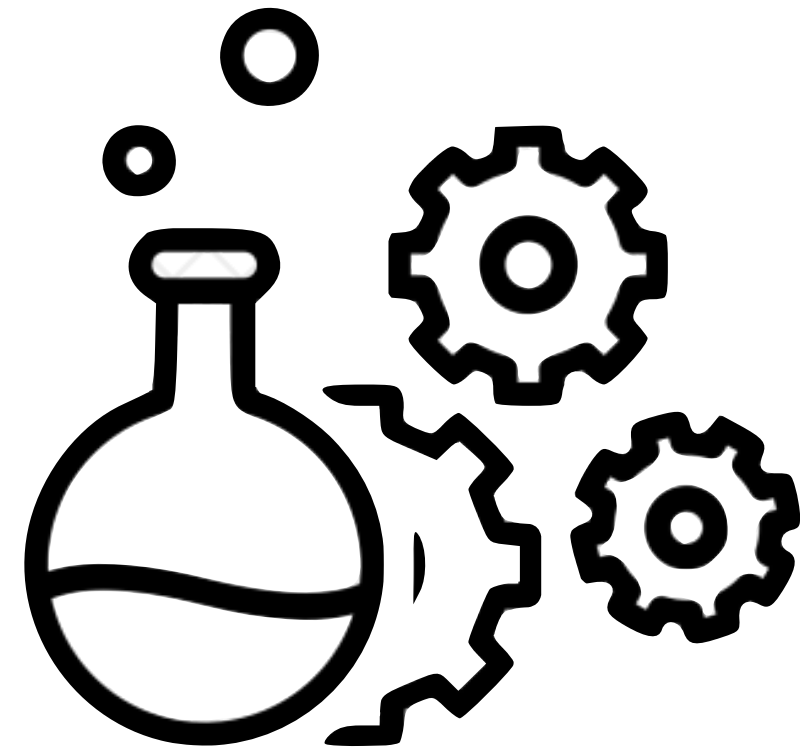
skip-list outperforms for w/l's with range queries



vector perform 12% less flushes; more I/Os in hash based buffers

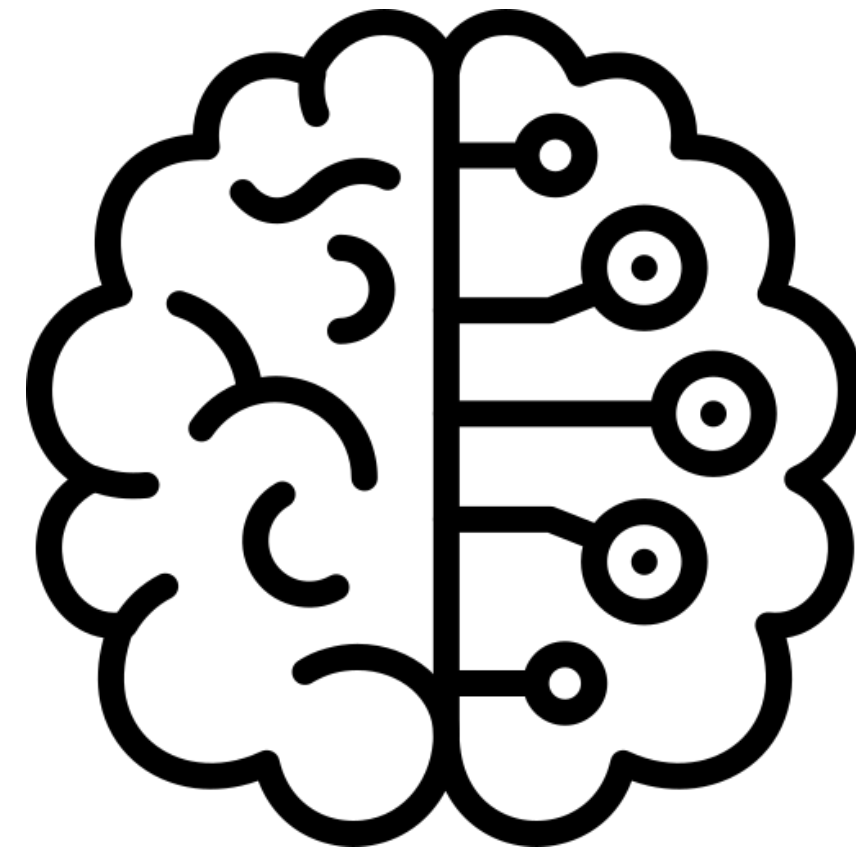


Ongoing Work



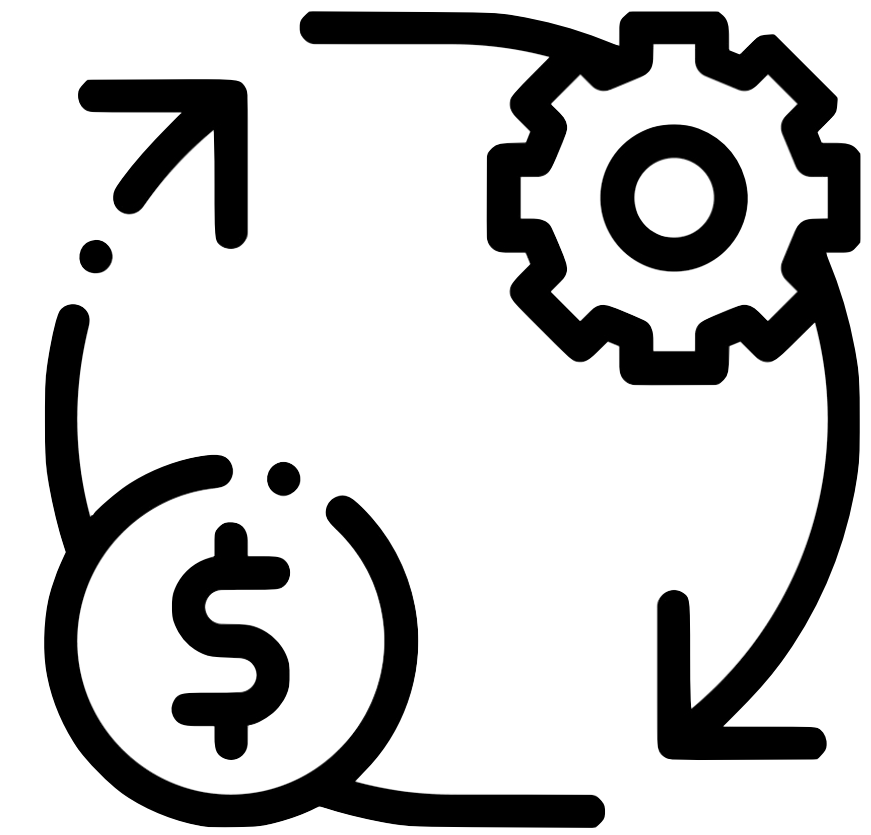
explore other
implementations

tree, trie, heap



finding best
buffer

using machine learning



switching
buffer

on the fly when w/l changes

Anatomy of LSM Memory Buffer

Thank You!

Questions?



Brandeis
UNIVERSITY

