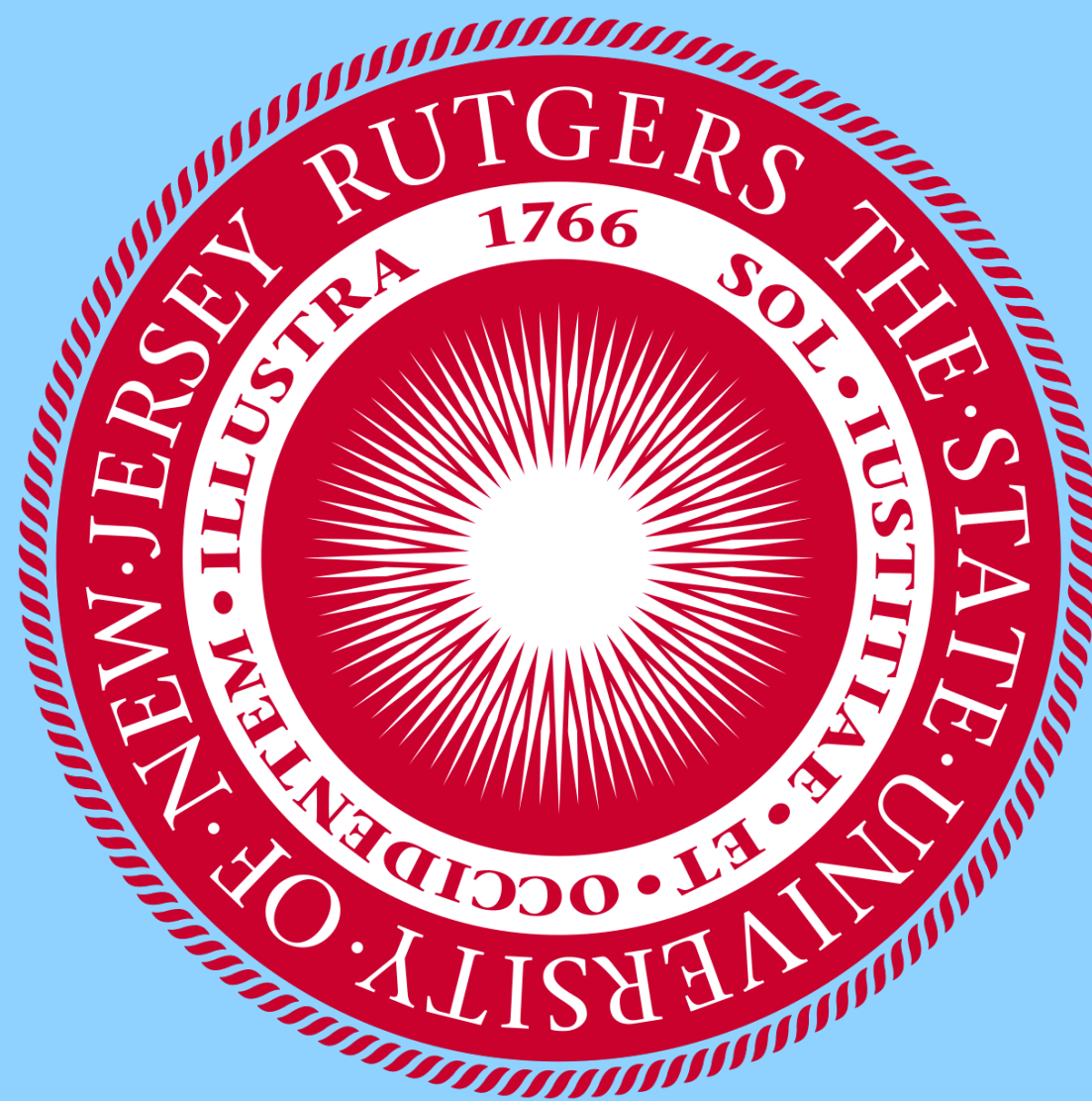


SplinterDB:

Closing the NVMe Bandwidth Gap

Alex Conway

Vijay Chidambaram
Martin Farach-Colton
Abihsihek Gupta
Richard Spillane
Amy Tai
Rob Johnson



ATC 2020

SplinterDB: A Key-Value Store for the Hard Cases

Key-Value Stores

What's hard?

Fast Storage

Our Approach

Use new data structures to lower IO amplification and CPU overhead while enabling concurrency

Key-Value Stores

What's hard?

Fast Storage

Small Key-Value Pairs

Our Approach

Use new data structures to lower IO amplification and CPU overhead while enabling concurrency

Keep key-value pairs sorted and packed into data blocks, delay merging as much as possible

Key-Value Stores

What's hard?

Fast Storage

Small Key-Value Pairs

Small Cache

Our Approach

Use new data structures to lower IO amplification and CPU overhead while enabling concurrency

Keep key-value pairs sorted and packed into data blocks, delay merging as much as possible

Make all data structures swappable in order to gracefully degrade under cache pressure

In this talk

Fast Storage
(NVMe)

SplinterDB

Data Structures

Flush-then-Compact

In this talk

Fast Storage
(NVMe)

SplinterDB

Data Structures

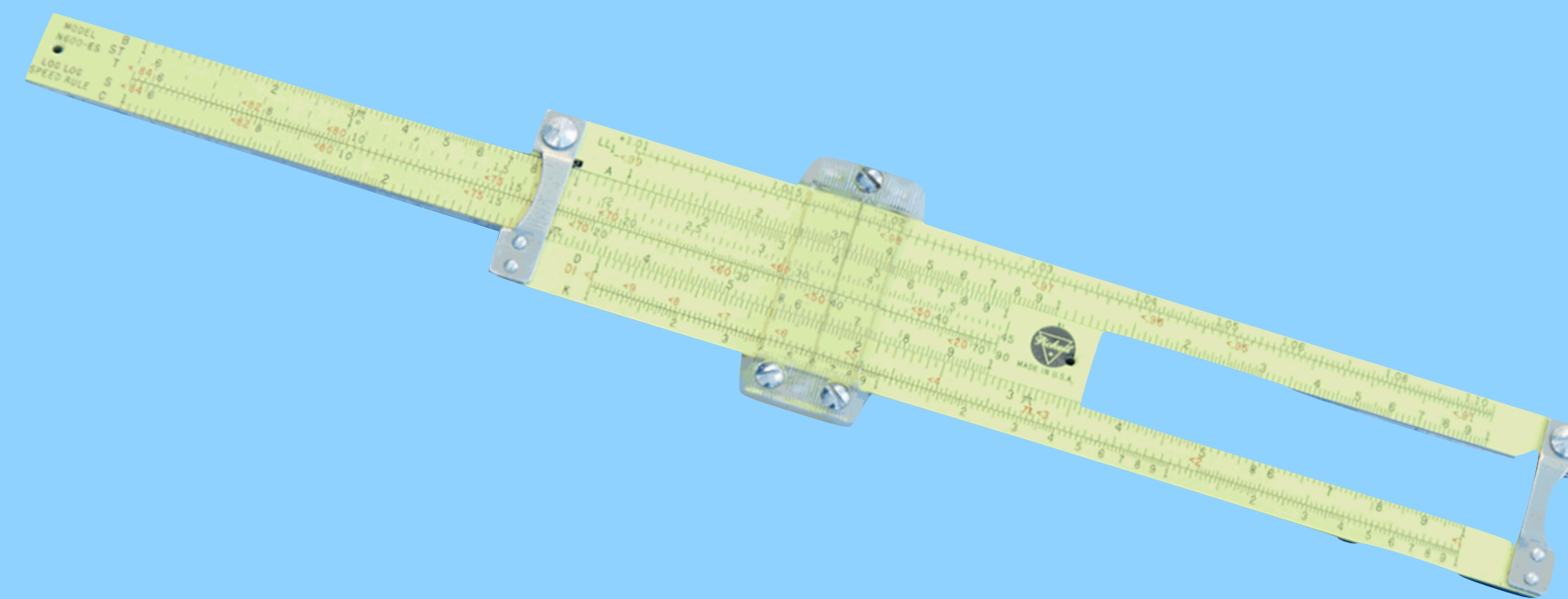
Flush-then-Compact

Back in the day...

People used to use...

Back in the day...

People used to use...



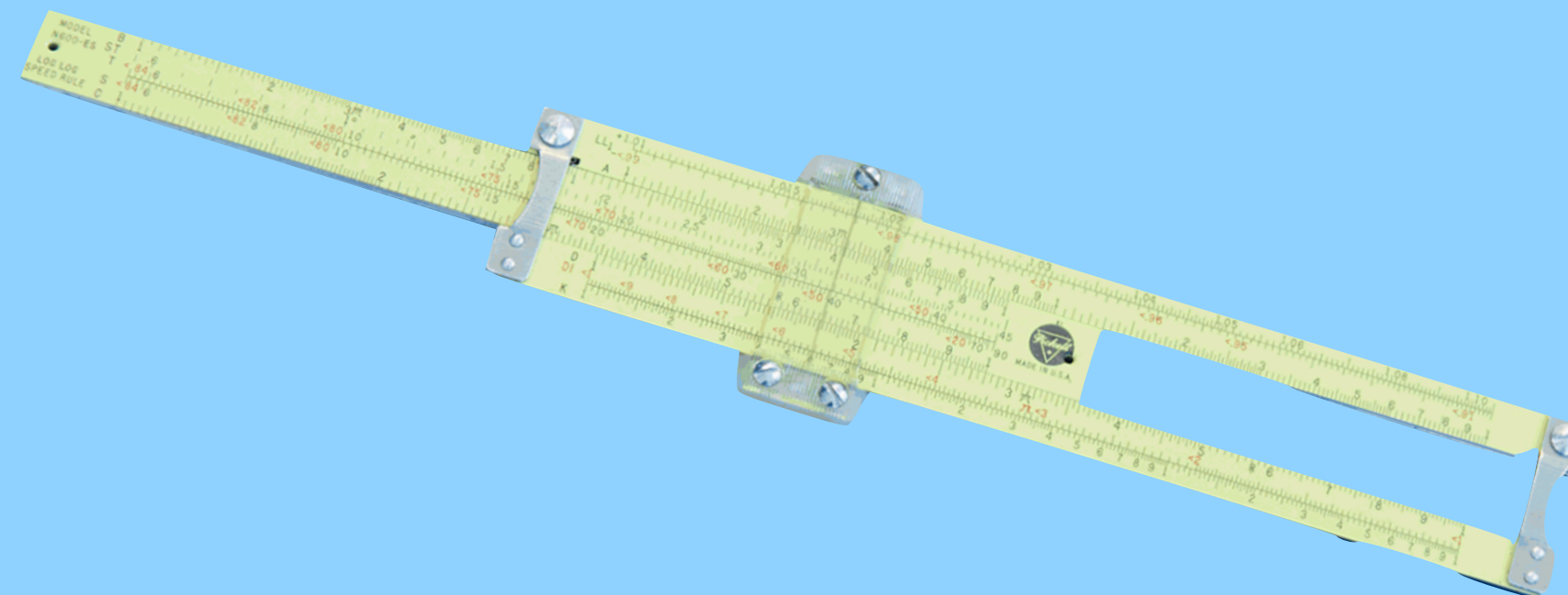
Slide rules

Back in the day...

People used to use...



VHS tapes



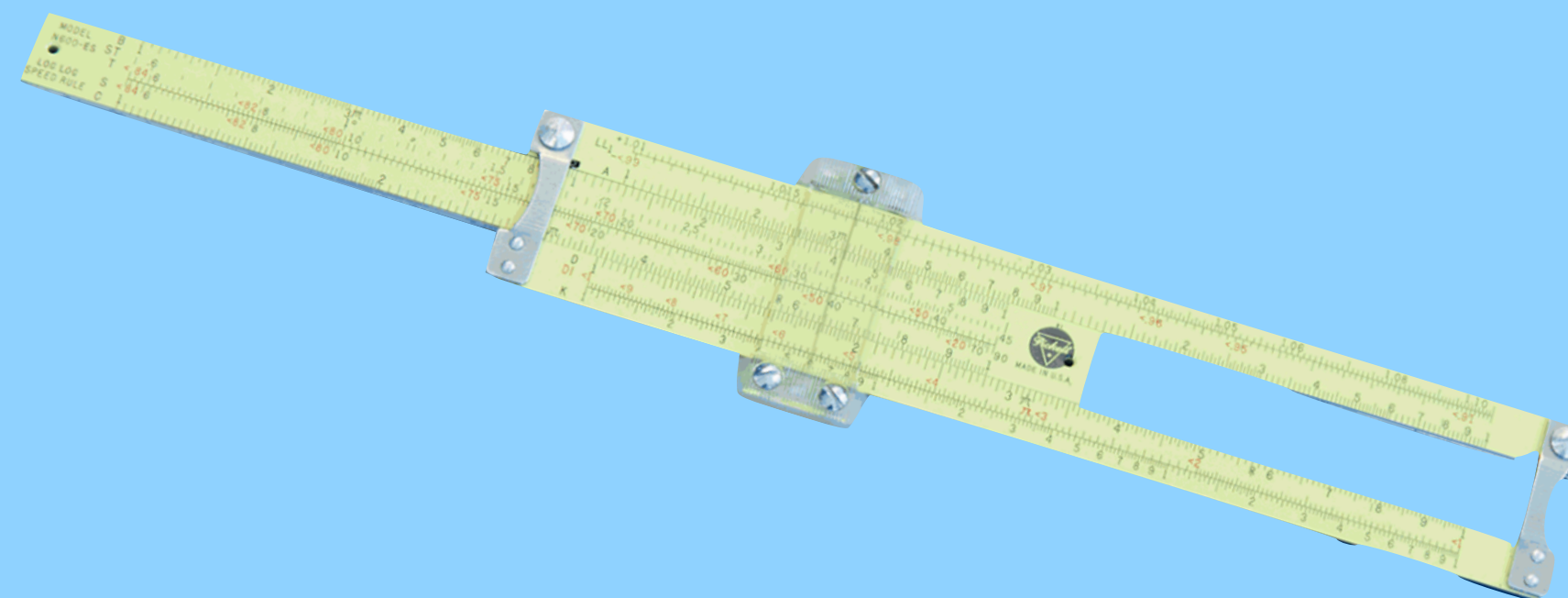
Slide rules

Back in the day...

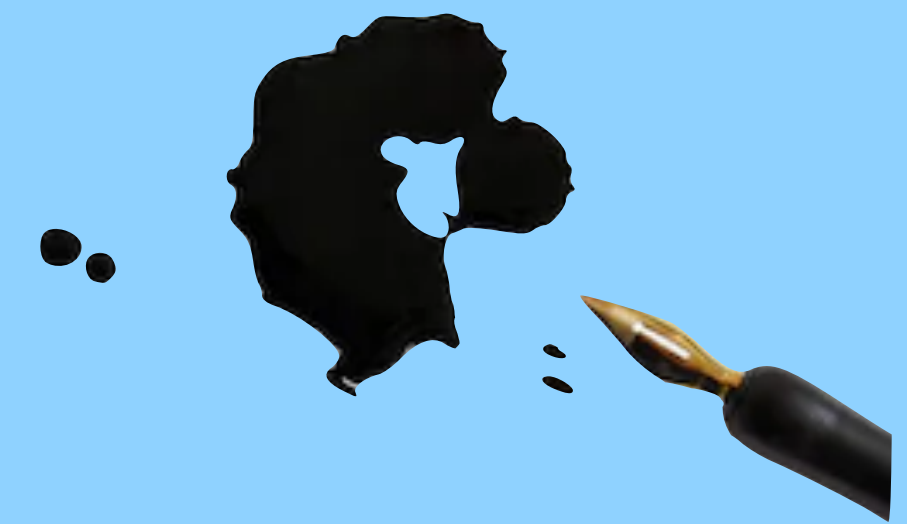
People used to use...



VHS tapes



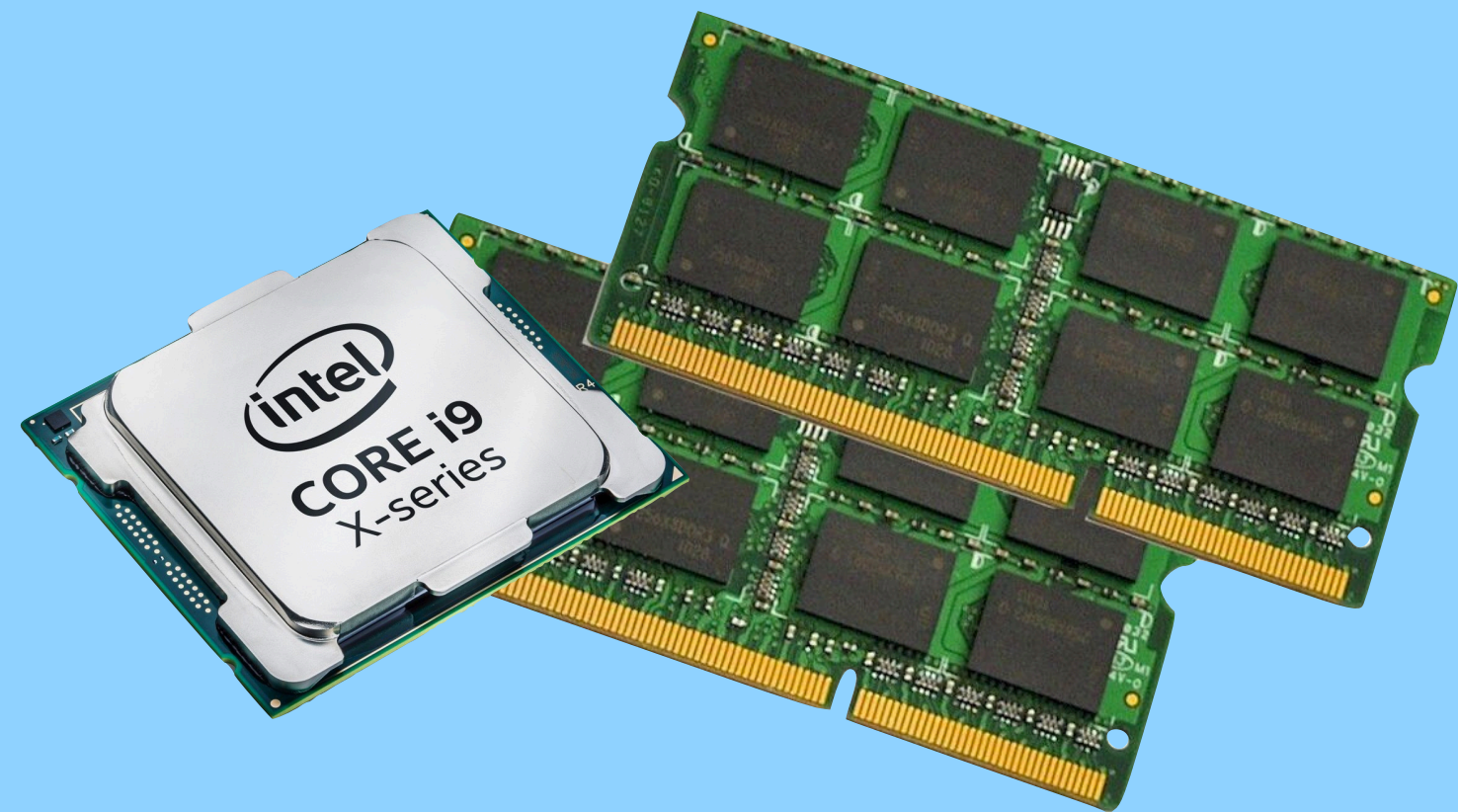
Slide rules



Fountain pens

Back in the day...

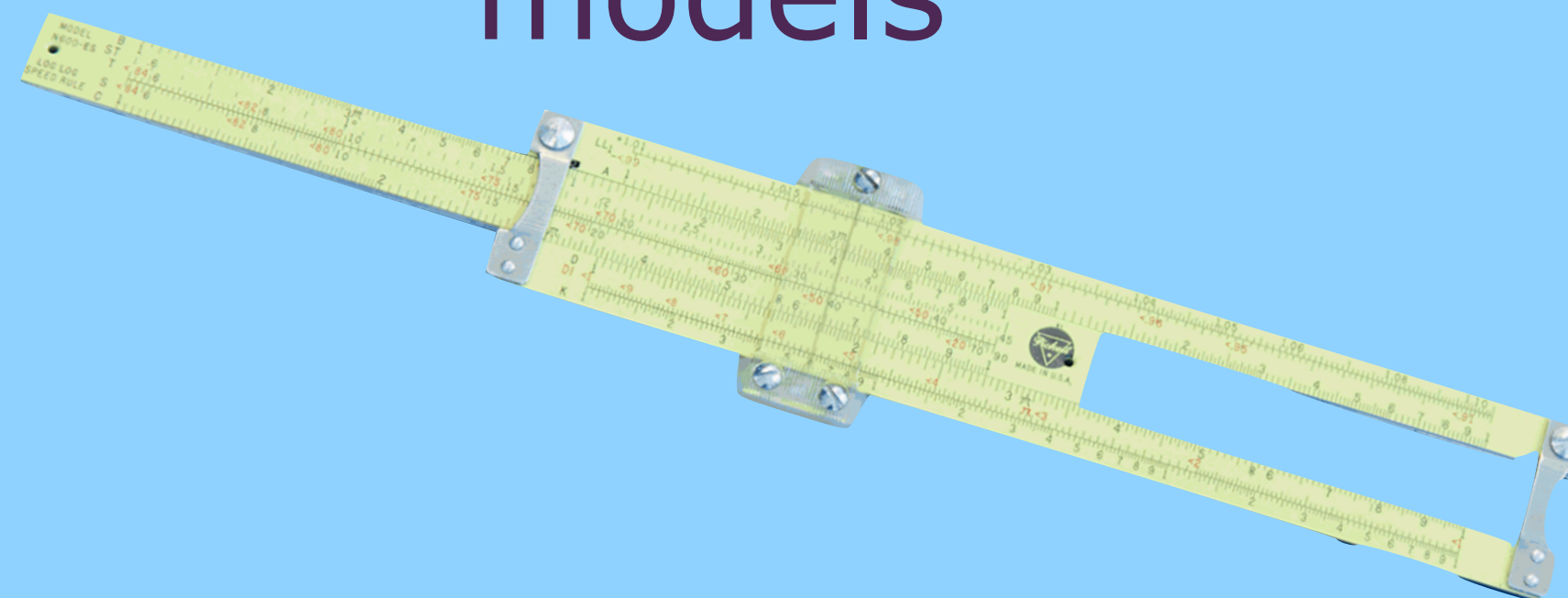
People used to use...



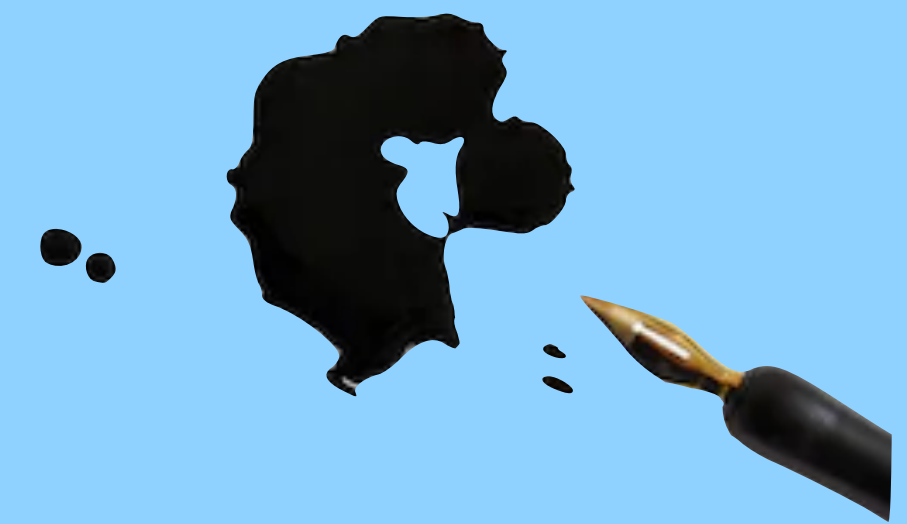
Different Performance models



VHS tapes



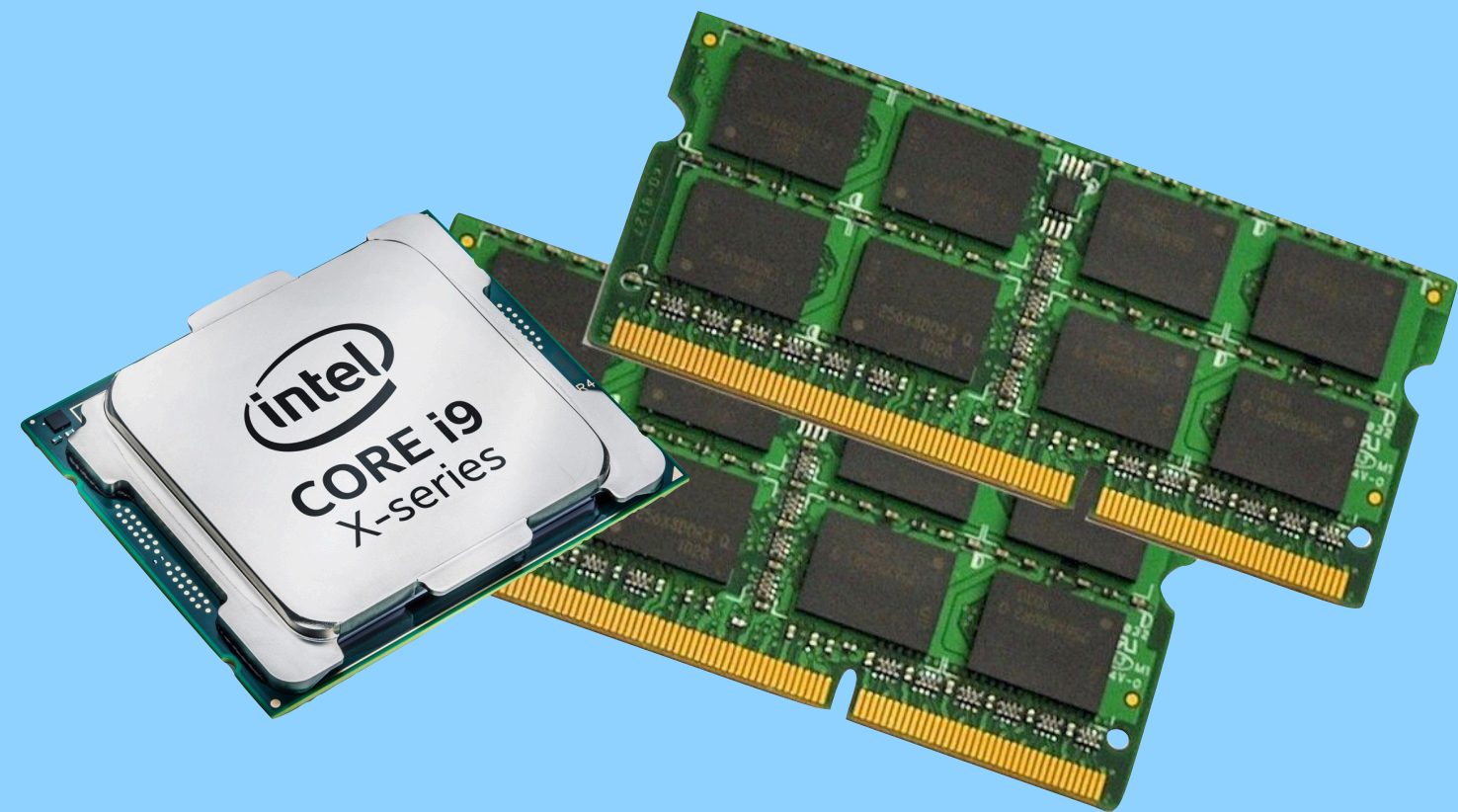
Slide rules



Fountain pens

Back in the day...

People used to use...



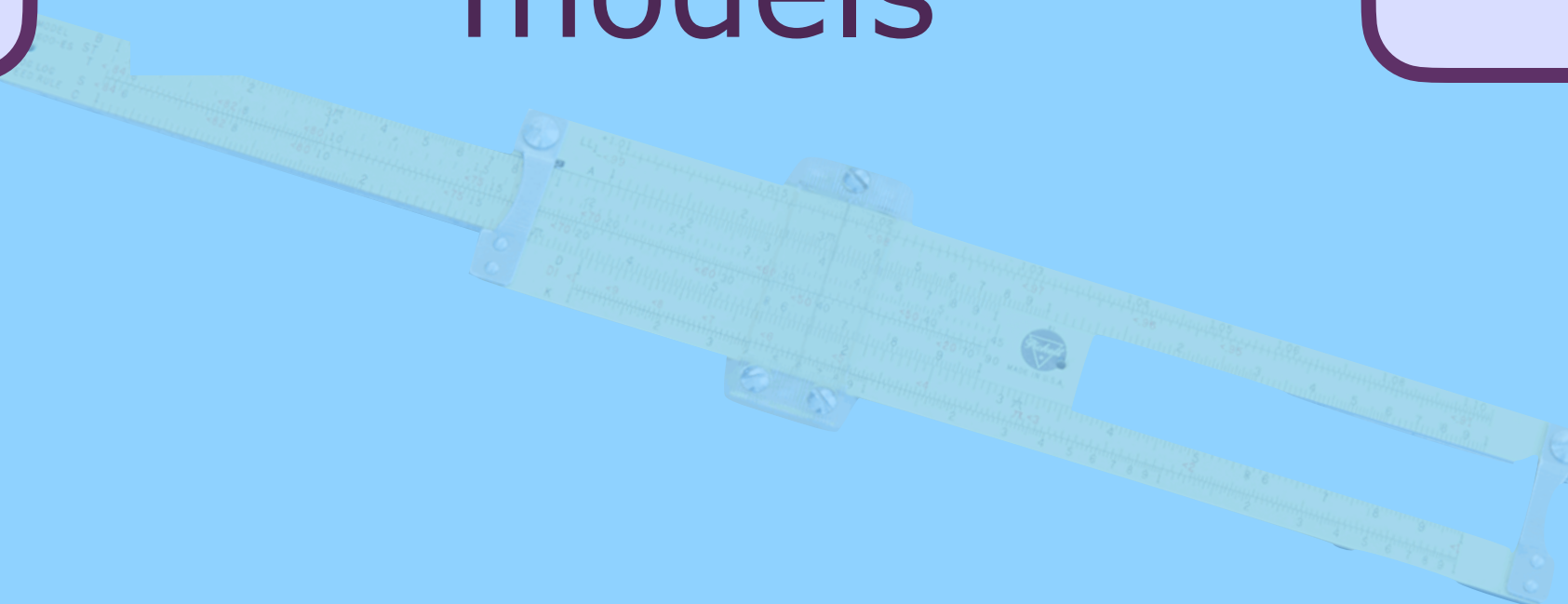
CPU

Different
Performance
models

IO



VHS tapes



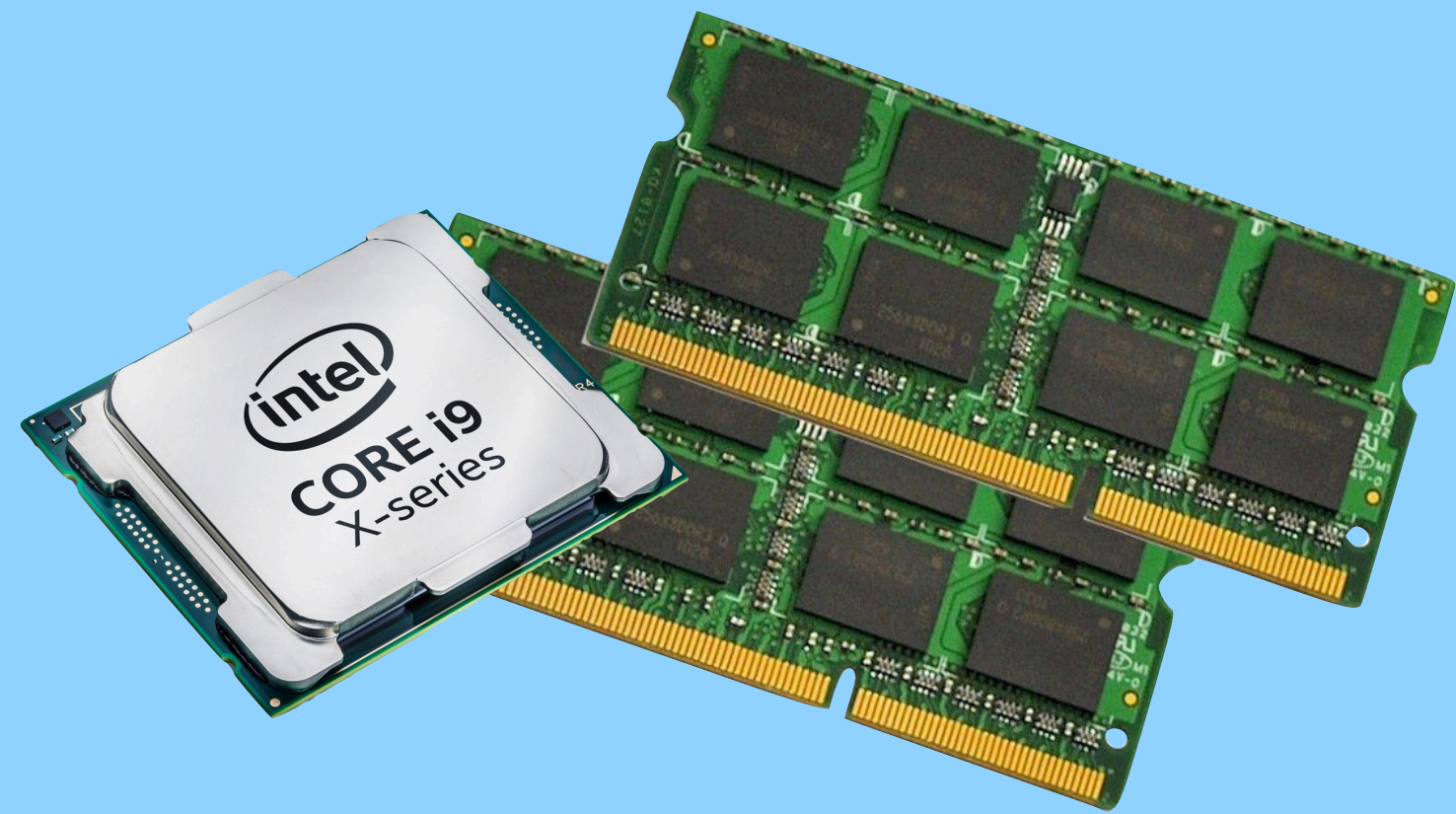
Slide rules



Fountain pens

Back in the day...

People used to use...



CPU



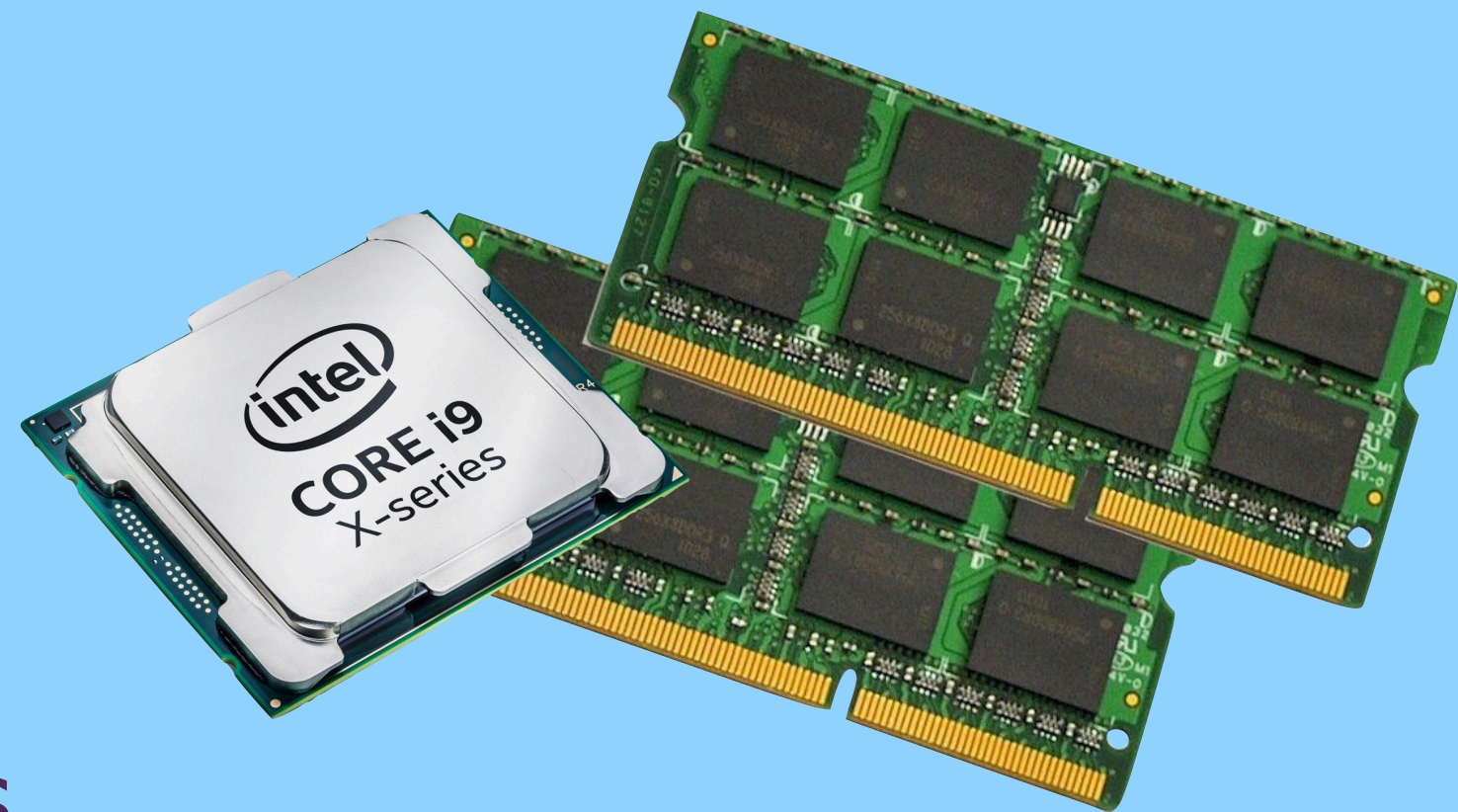
IO

Different
Performance
models

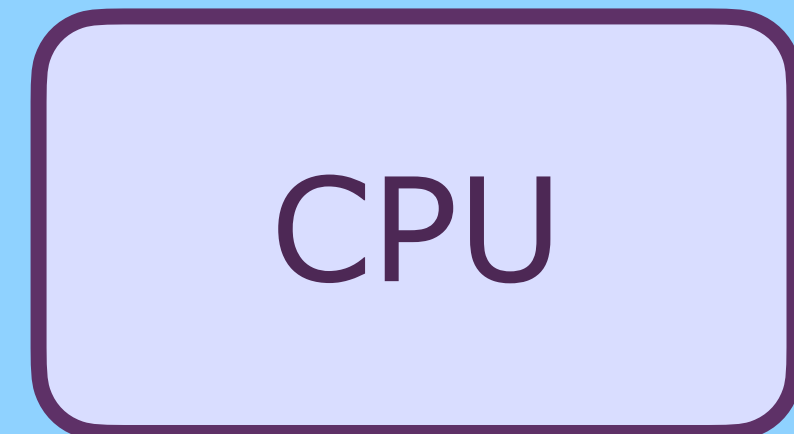
Look at *e.g.*
key-value stores

Back in the day...

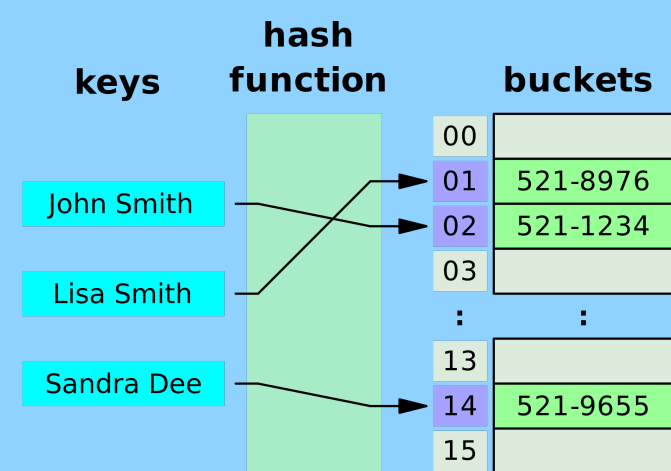
People used to use...



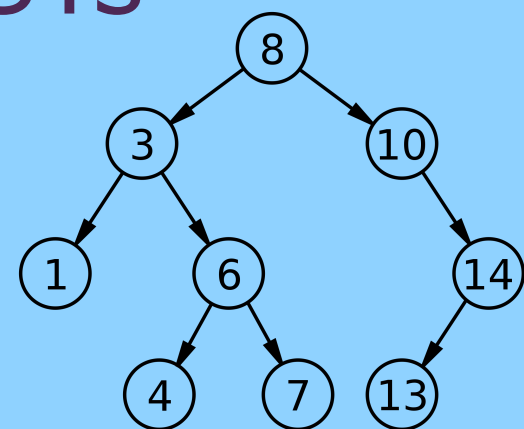
hash tables



Different Performance models



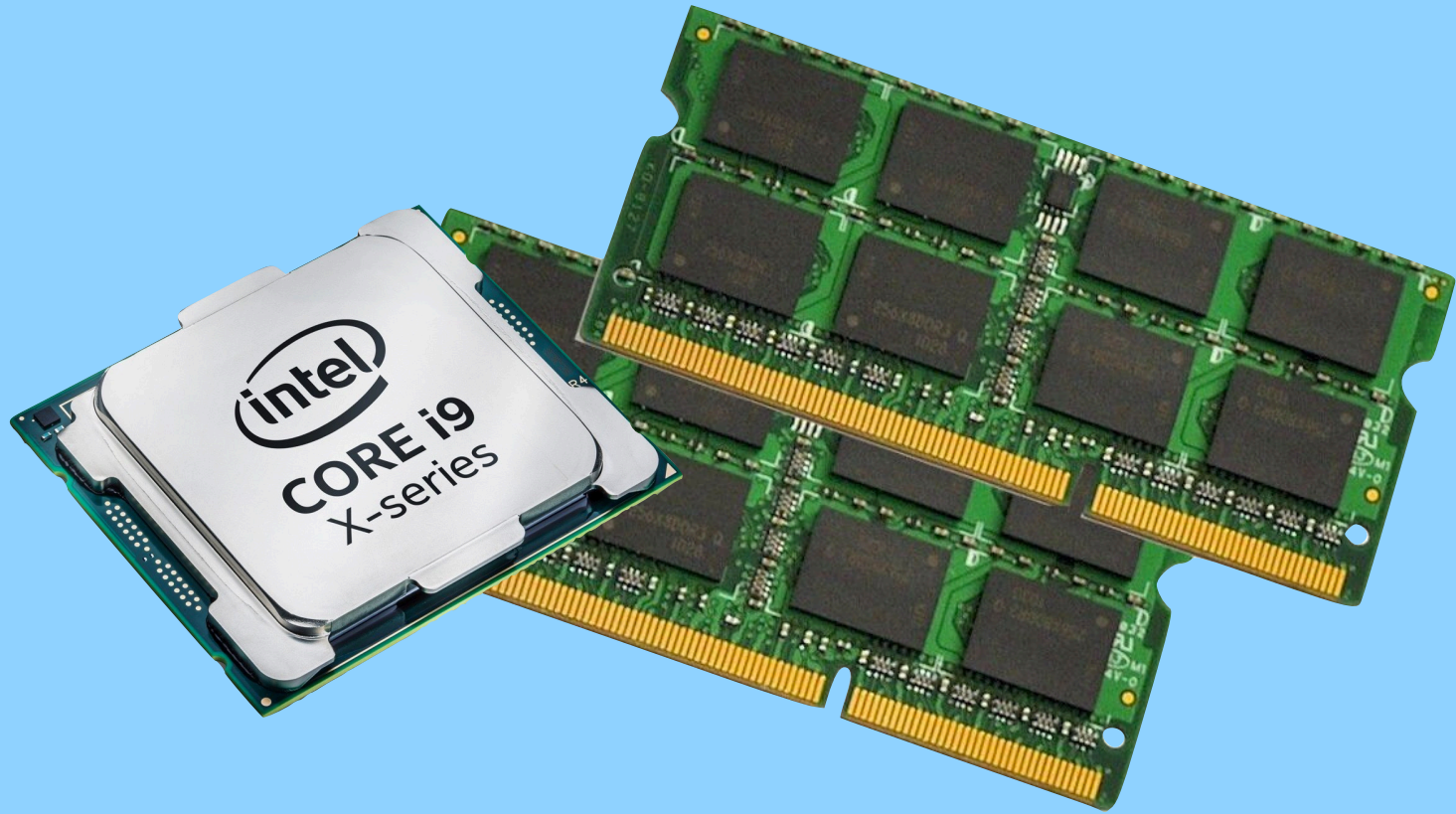
BSTs



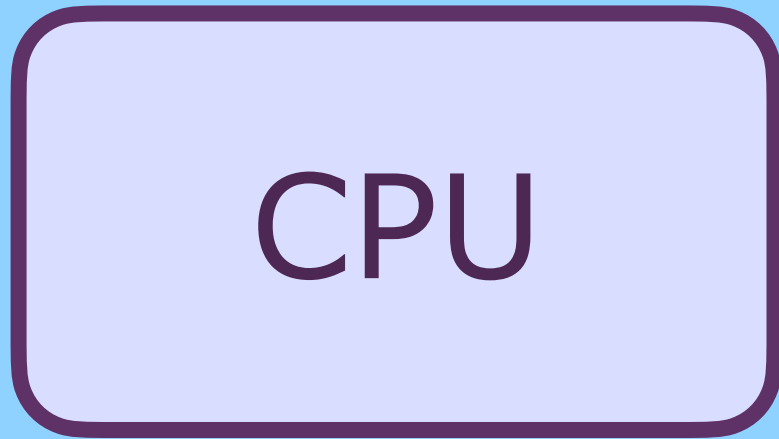
Look at *e.g.* key-value stores

Back in the day...

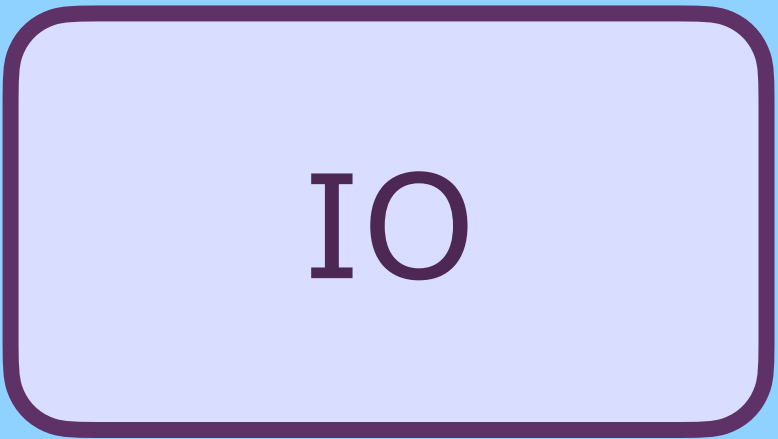
People used to use...



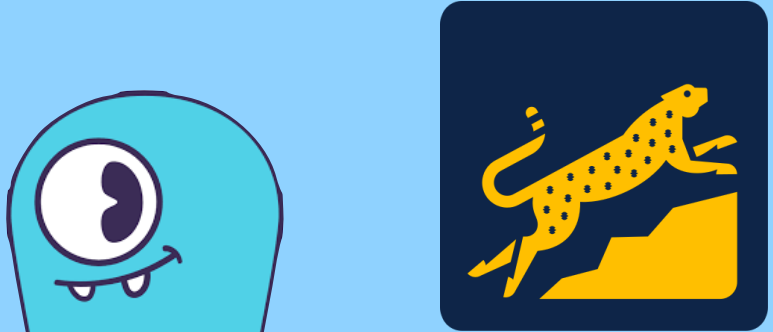
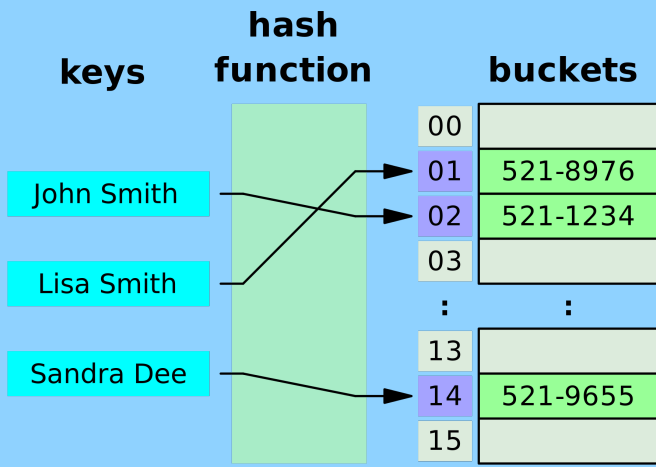
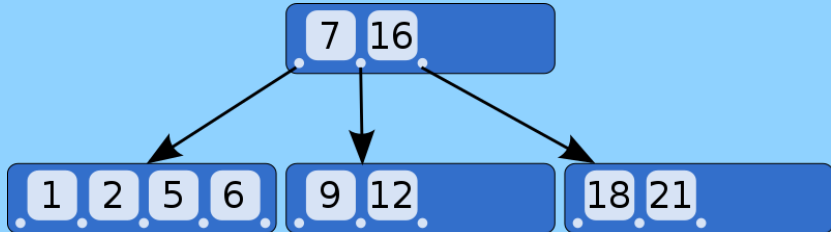
hash tables



Different Performance models

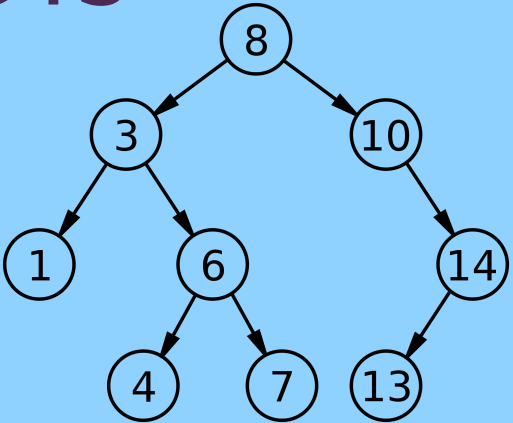


B-trees



SCYLLA. RocksDB

BSTs



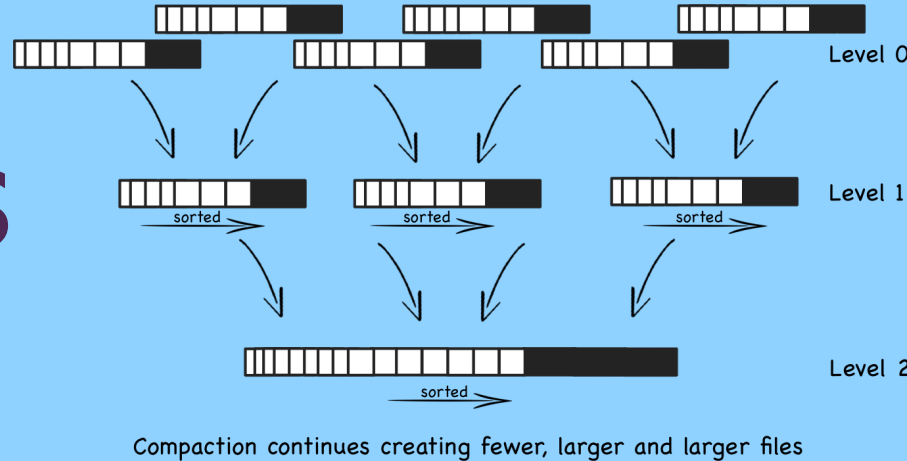
Memcached



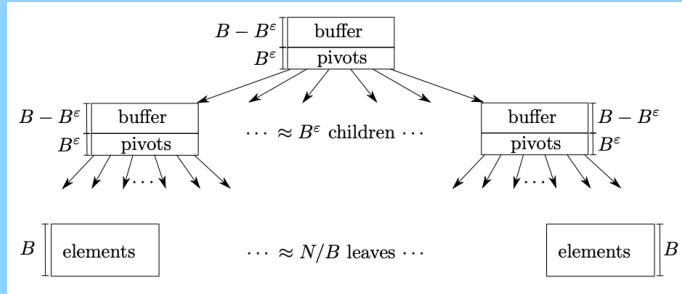
redis

Look at e.g. key-value stores

LSMs



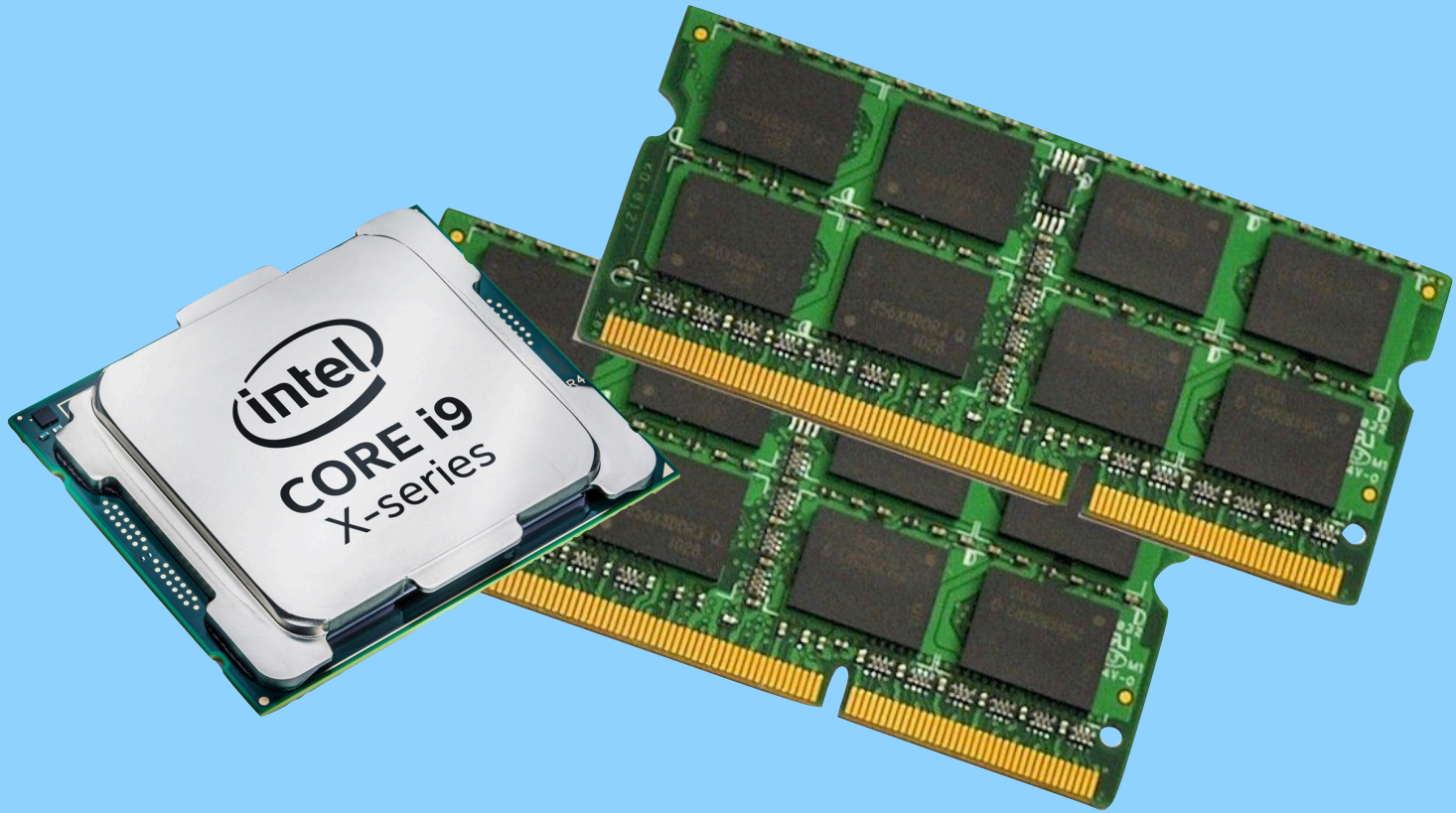
B^ε-trees



Back in the day...

Cycles per 64b word
at bandwidth

People used to use...



CPU

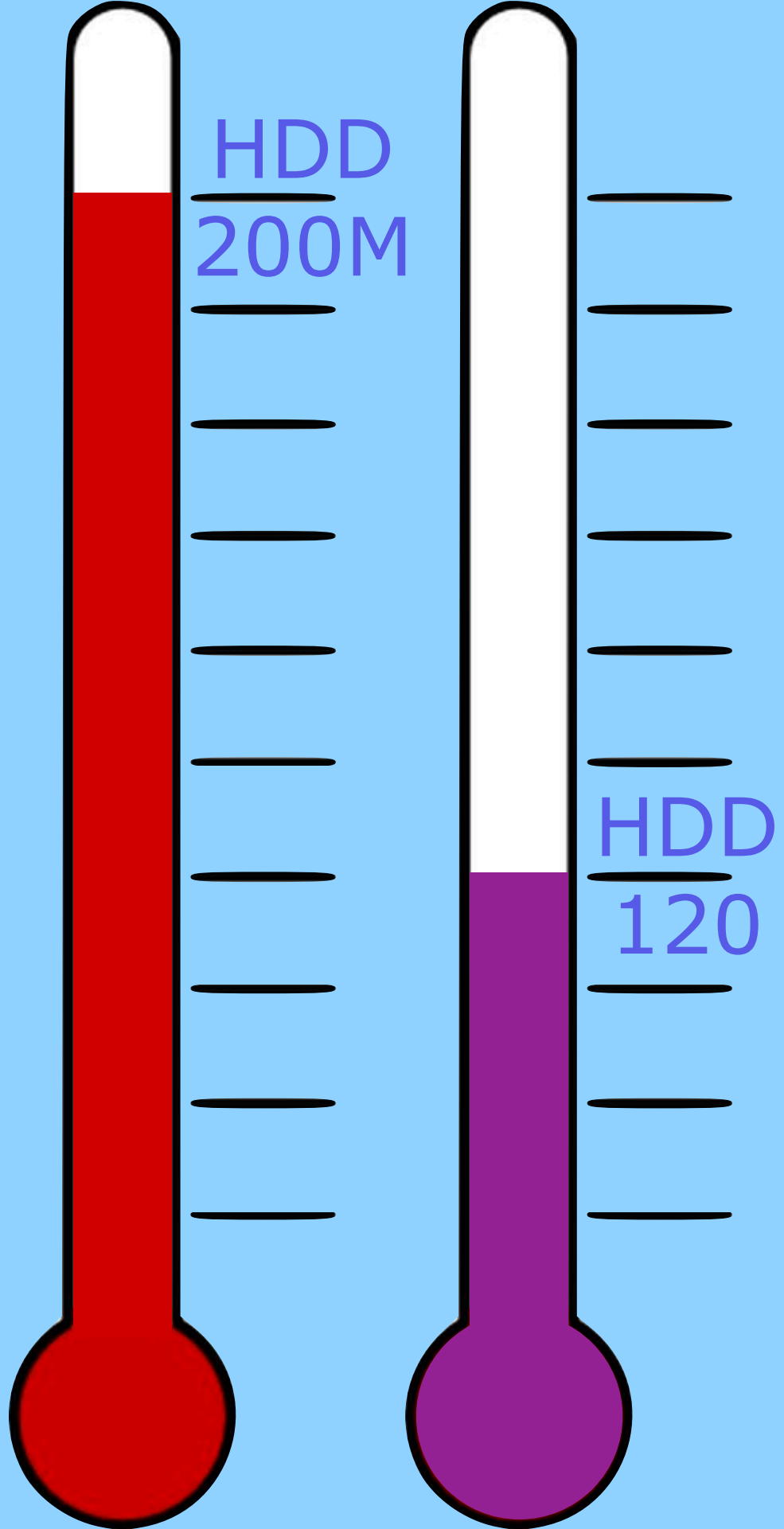


Hard Drive



IO

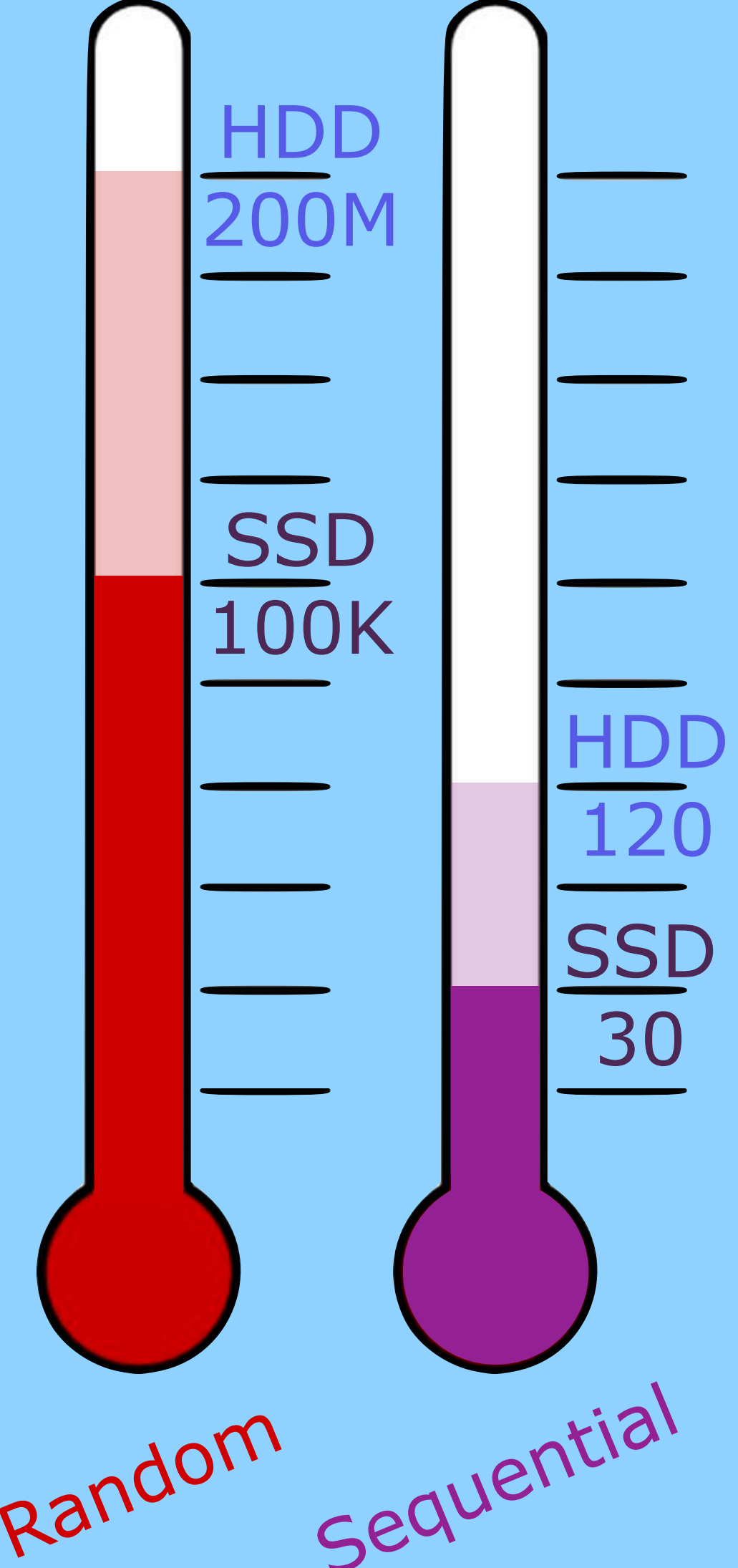
Different
Performance
models



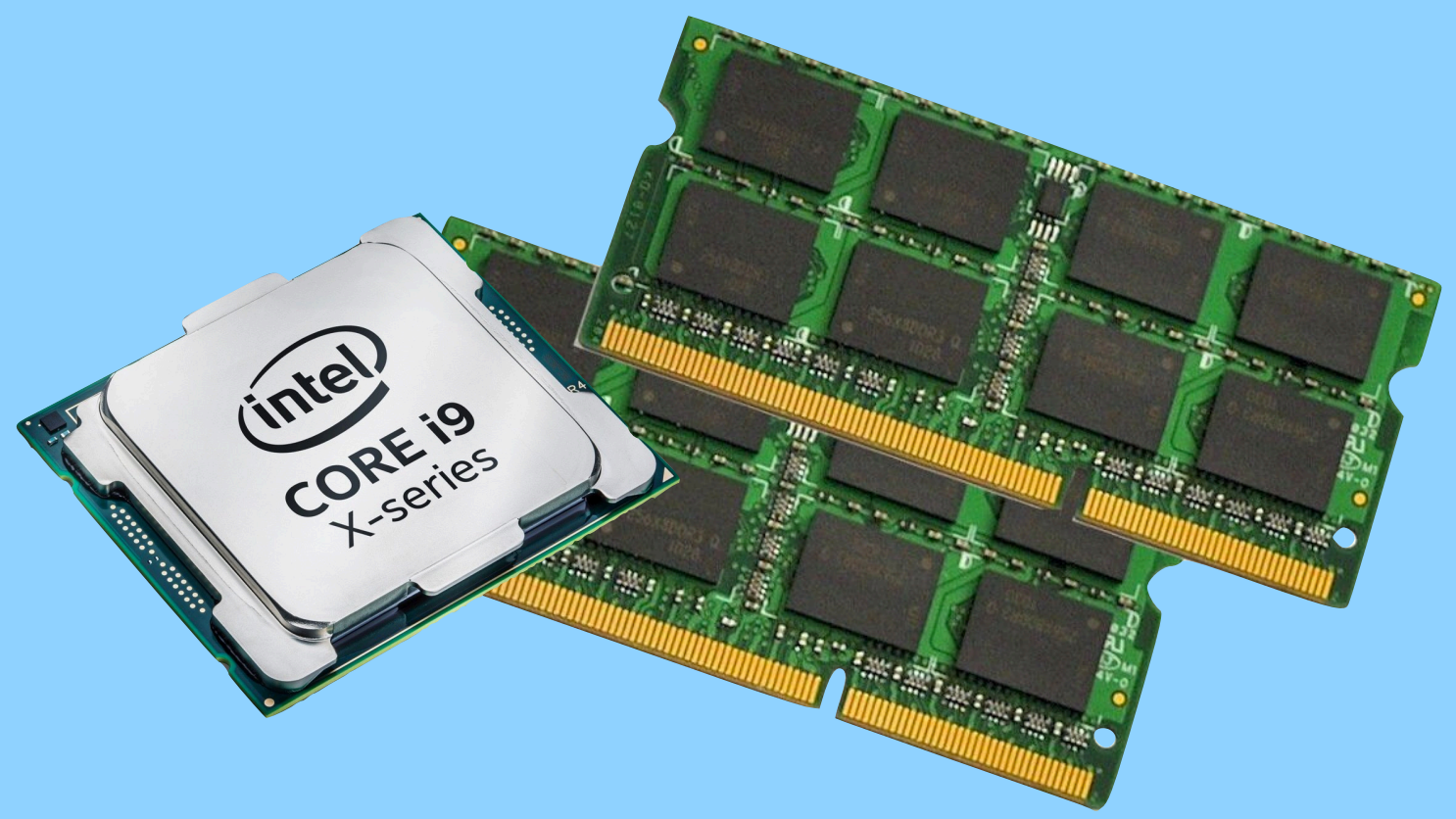
Random
Sequential

Back in the day...

Cycles per 64b word
at bandwidth

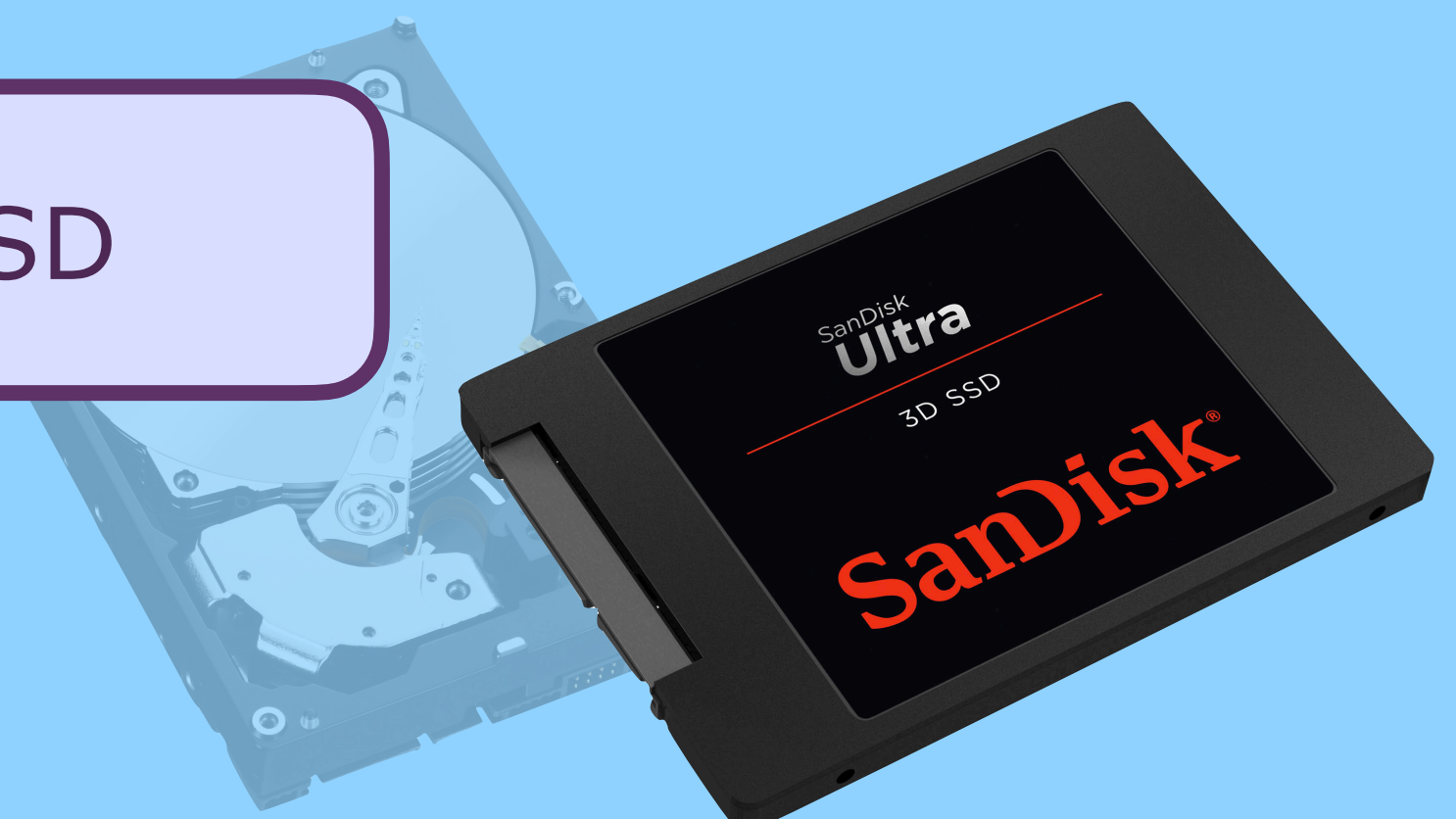


People used to use...



CPU

SSD

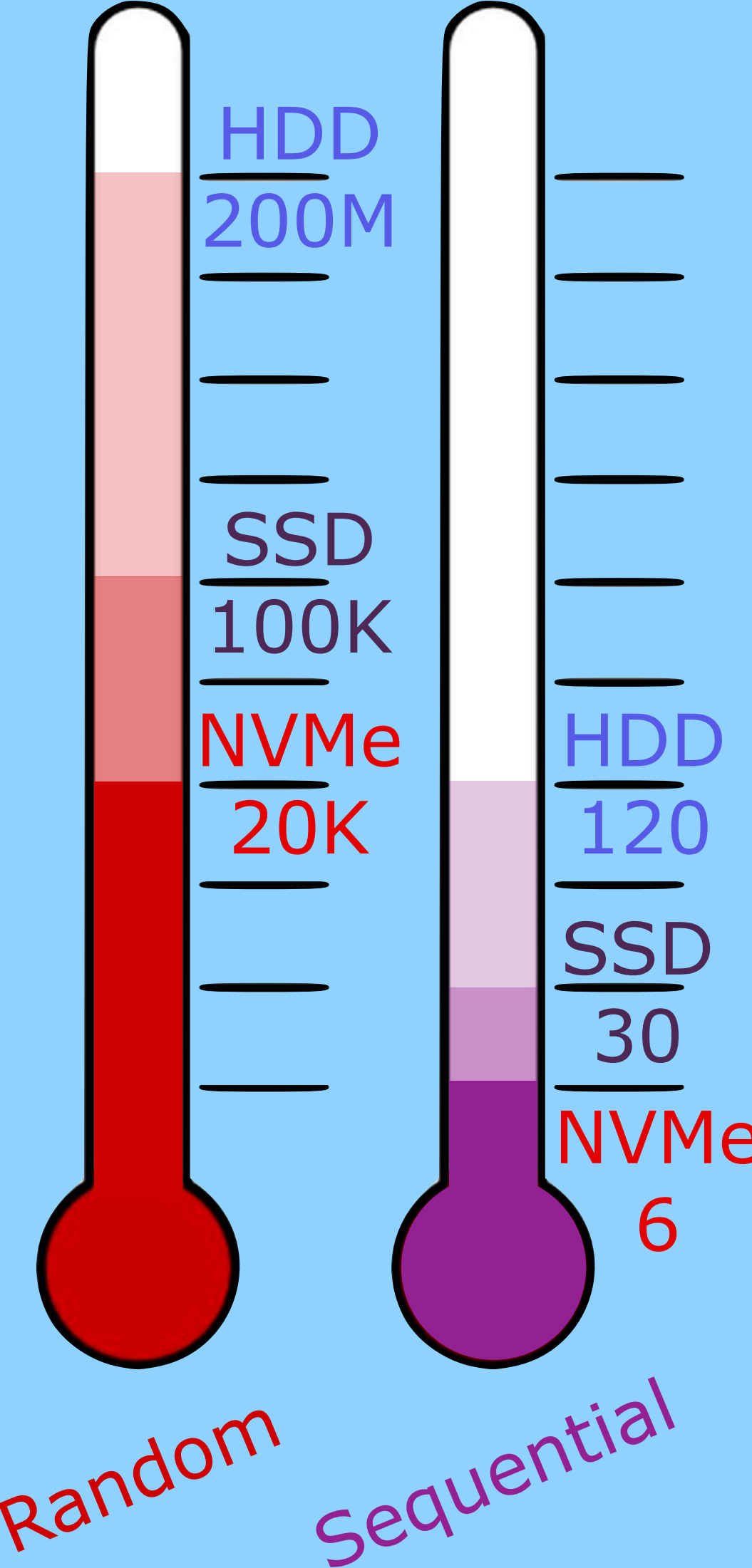


Different
Performance
models

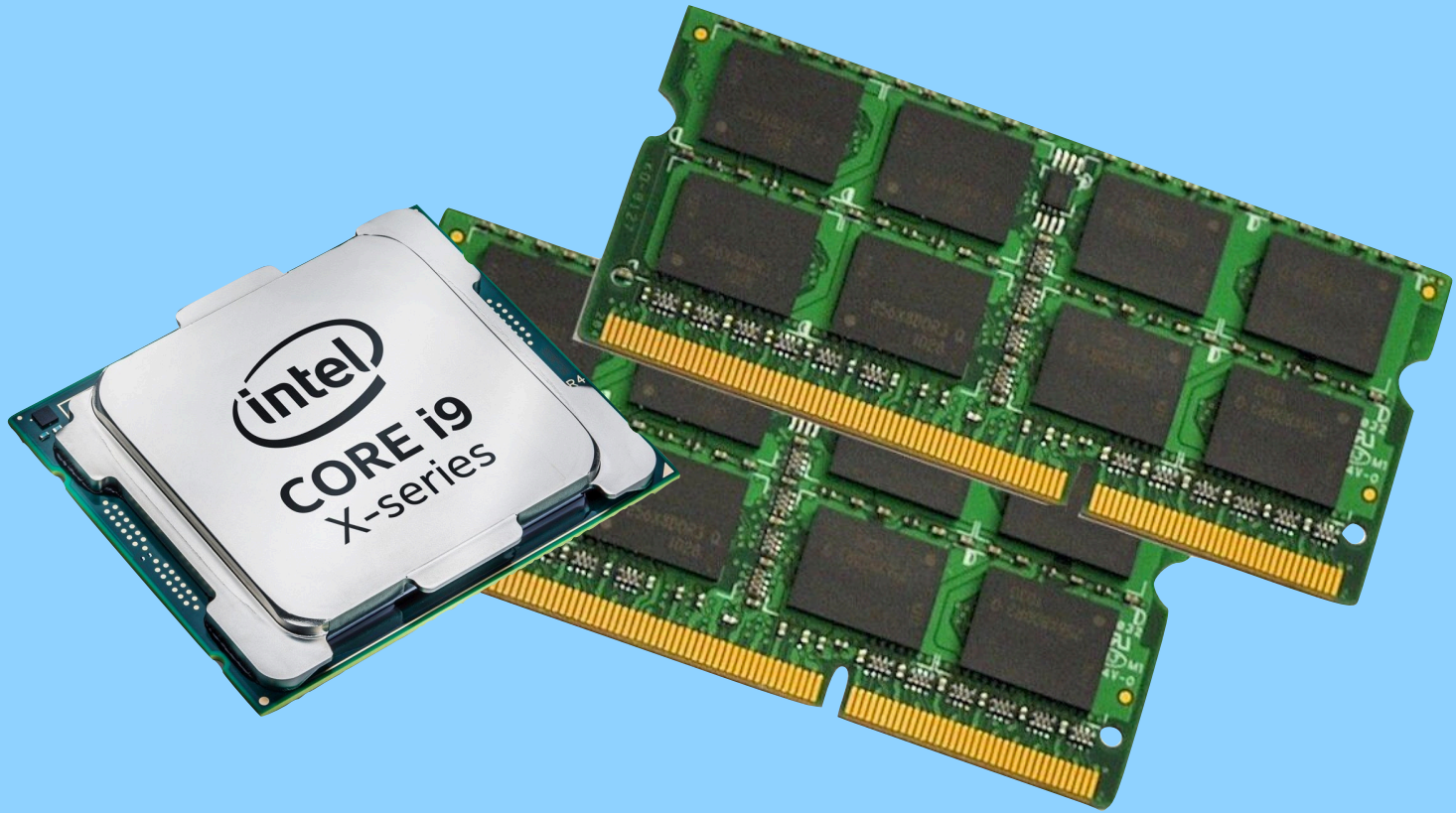
IO

Back in the day...

Cycles per 64b word
at bandwidth



People used to use...



CPU



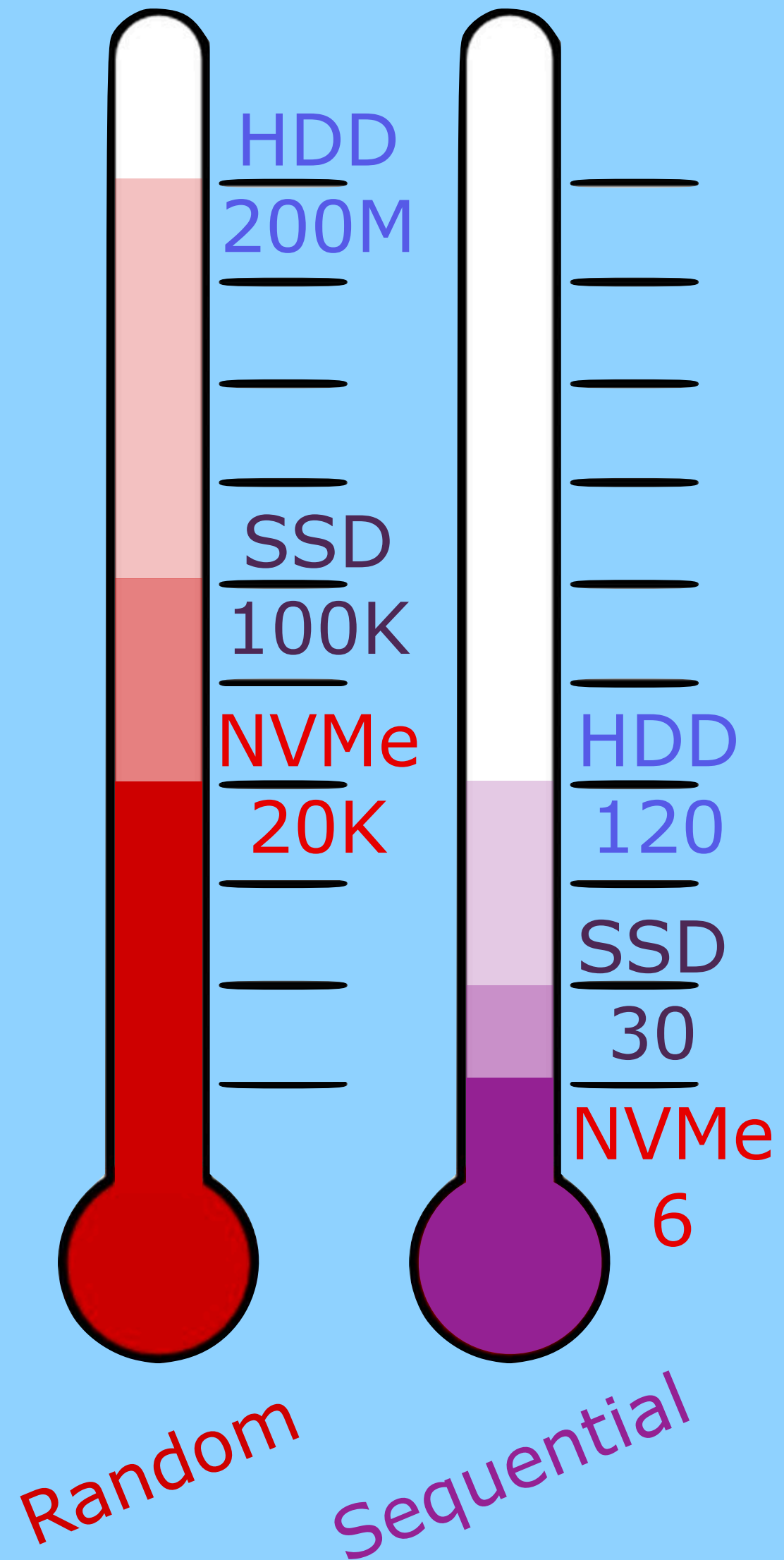
NVMe

Different
Performance
models

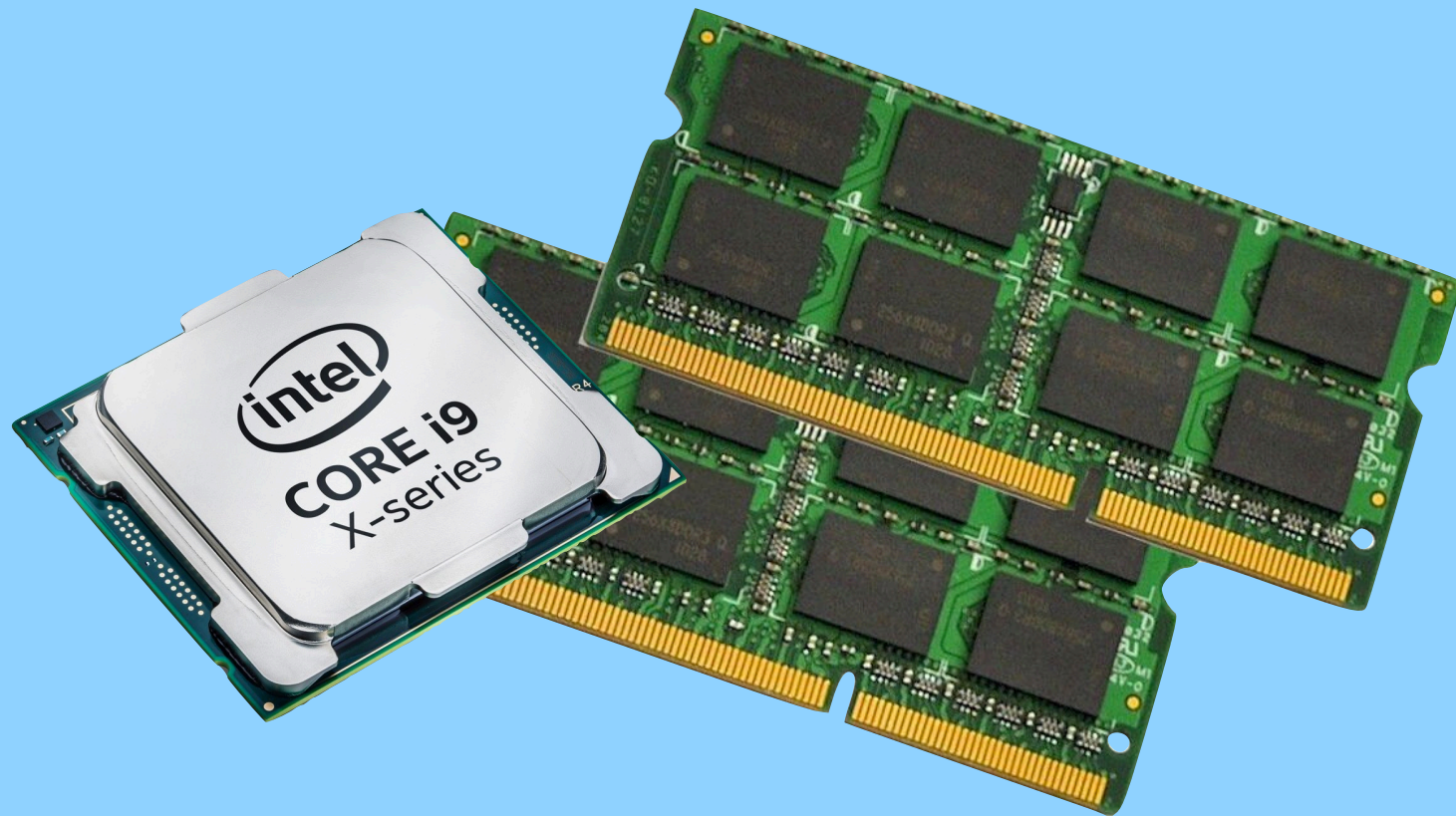
IO

Back in the day...

Cycles per 64b word
at bandwidth



People used to use...



NVMe

CPU

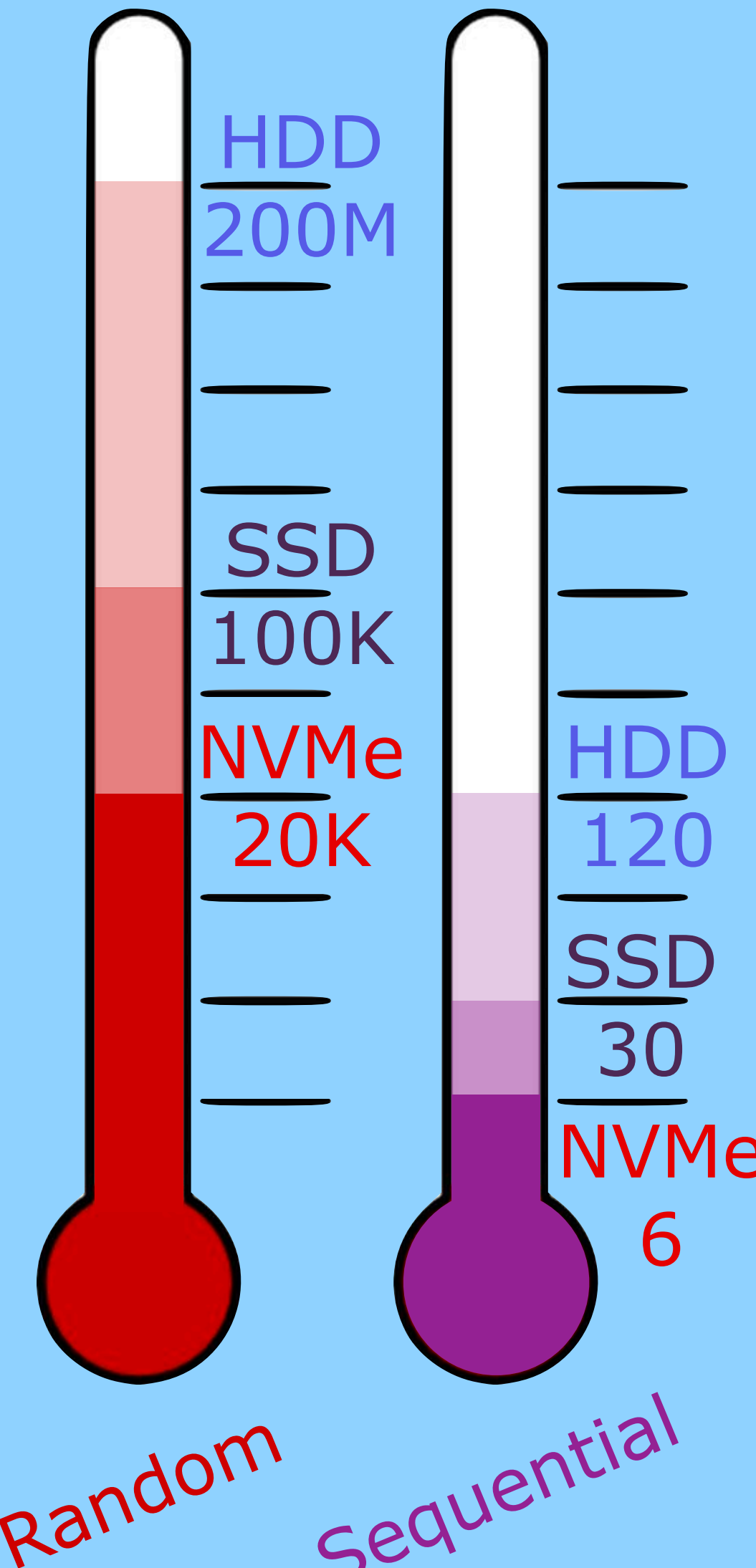
Different
Performance
models

IO

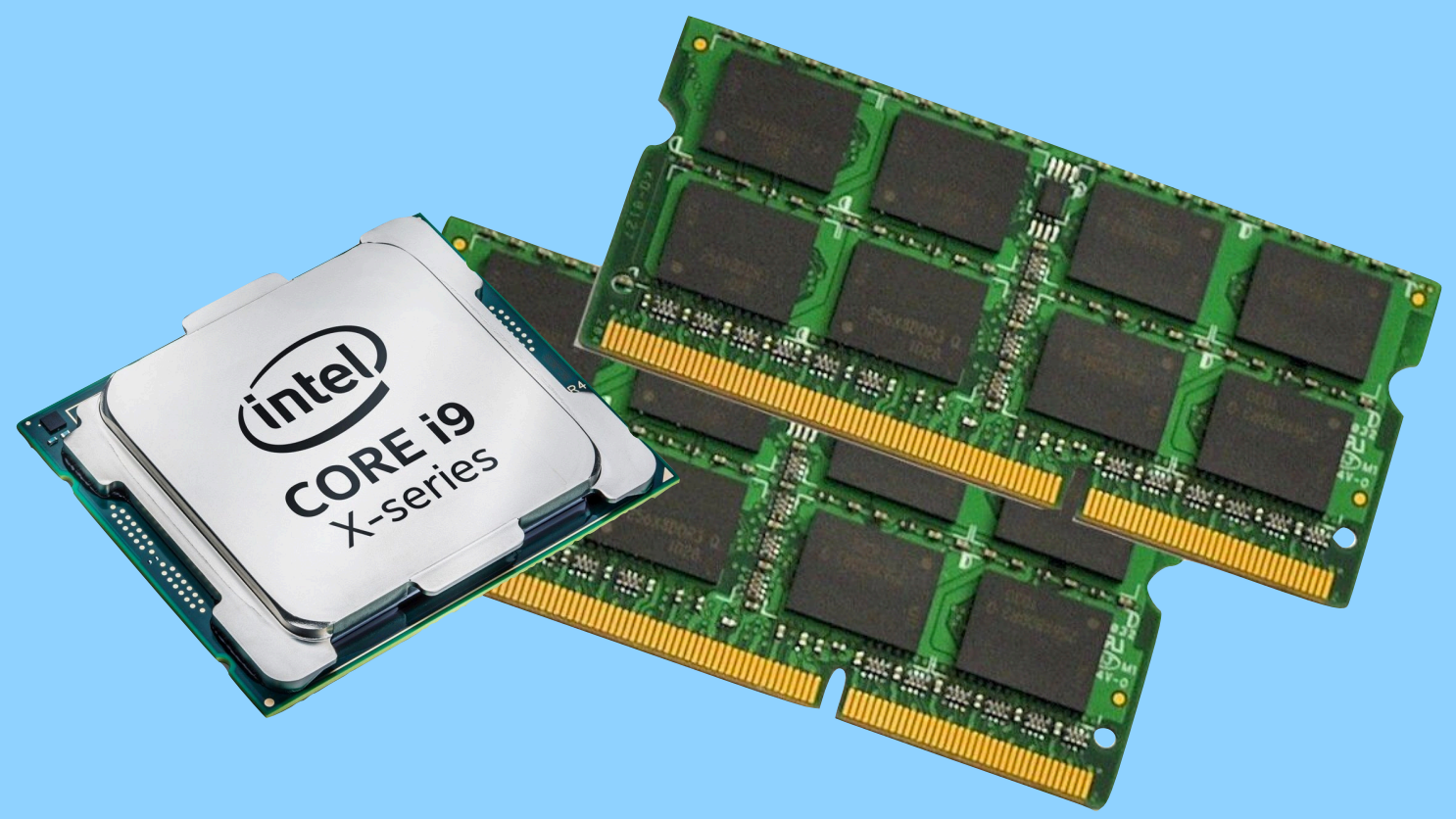
For NVMe, need data
structures to optimize both

Back in the day...

Cycles per 64b word
at bandwidth



People used to use...



NVMe

CPU

Different
Performance
models

IO

For NVMe, need data
structures to optimize both

~~Write Amplification~~

Work Amplification

In this talk

Fast Storage
(NVMe)

SplinterDB

Data Structures

Flush-then-Compact

SplinterDB

SplinterDB is a key-value store
which handles these tough cases:

Fast Storage

Small Key-Value Pairs

Small Cache

SplinterDB

SplinterDB is a key-value store
which handles these tough cases:

Fast Storage



Small Key-Value Pairs



Small Cache



SplinterDB

SplinterDB is a key-value store which handles these tough cases:

Fast Storage



Small Key-Value Pairs

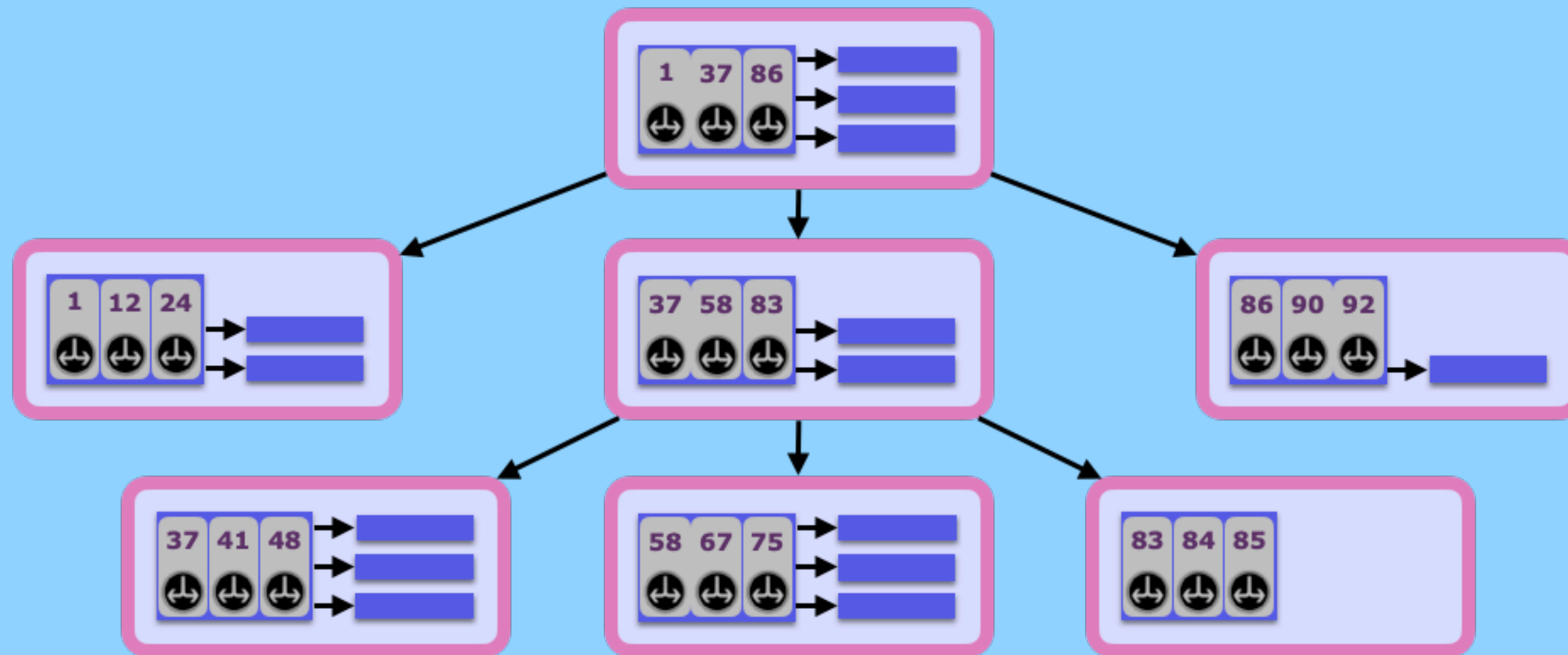


Small Cache



SplinterDB

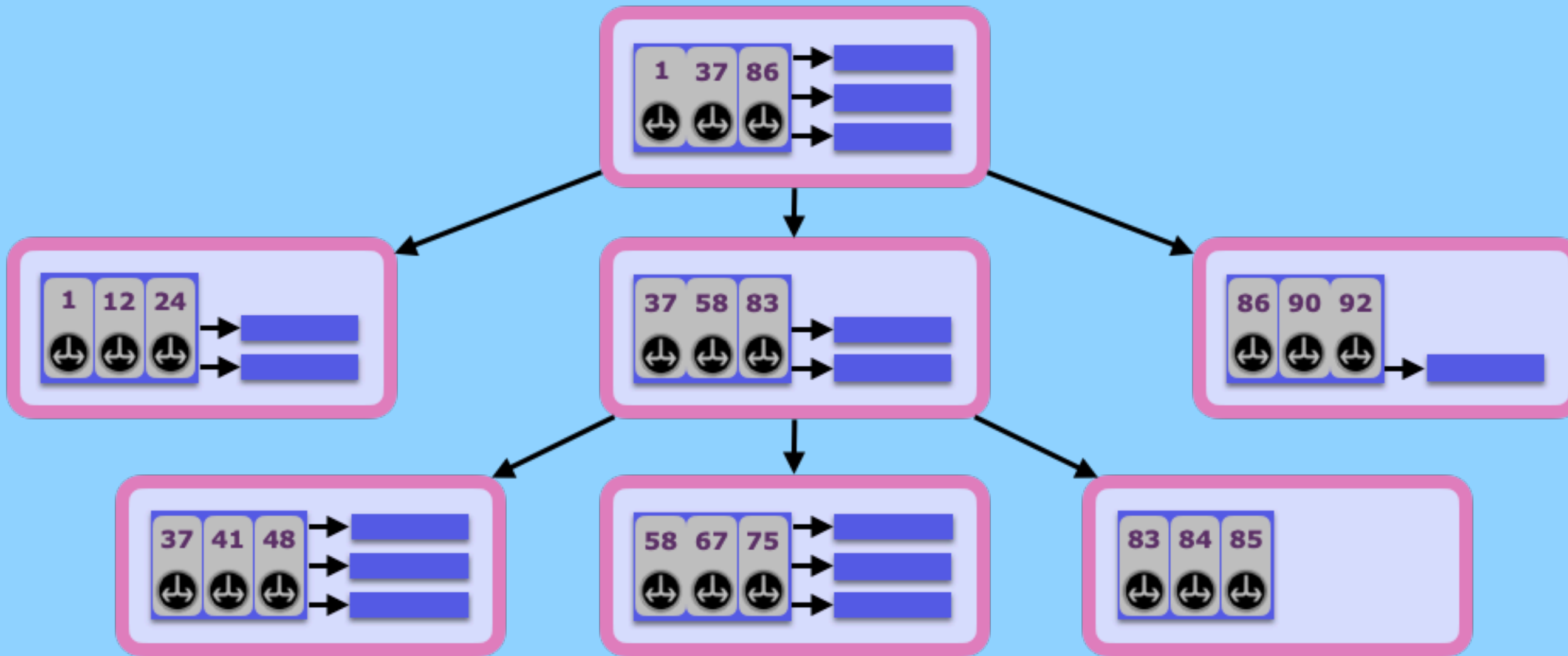
SplinterDB is a key-value store which handles these tough cases:



Size-Tiered B^ϵ -Tree

SplinterDB

SplinterDB is a key-value store which handles these tough cases:



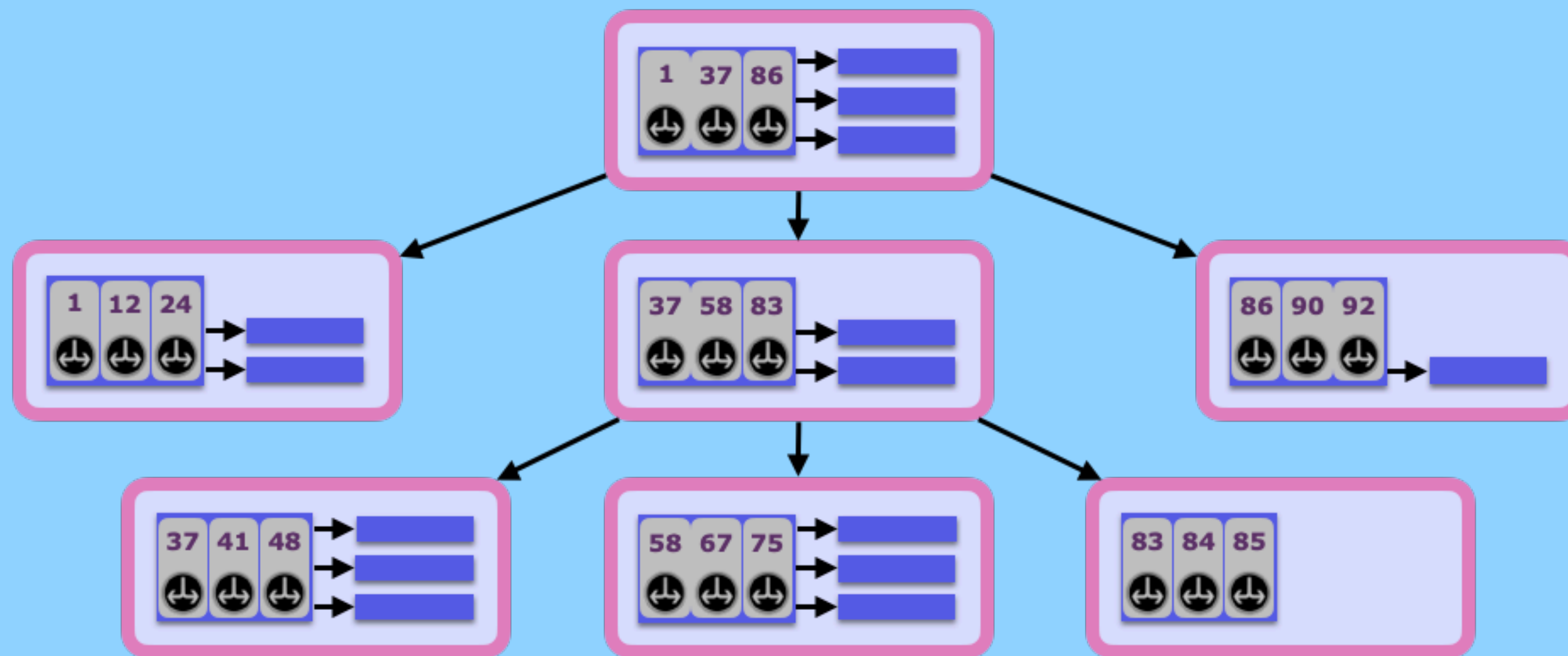
Size-Tiered B^ϵ -Tree

Reducing Work



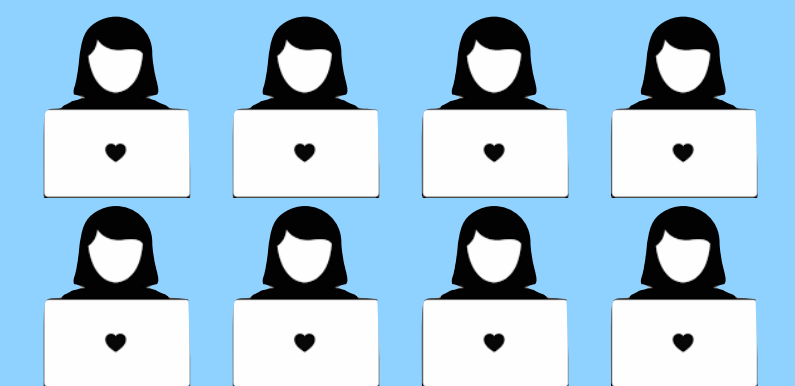
SplinterDB

SplinterDB is a key-value store which handles these tough cases:



Size-Tiered B^ϵ -Tree

Reducing Work



Concurrency

How Does SplinterDB Perform?

State of the Art: RocksDB

RocksDB



RocksDB is a high performance embedded database for key-value data. It is a fork of LevelDB by Facebook optimized to exploit many central processing unit (CPU) cores, and make efficient use of fast storage, such as solid-state drives (SSD), for input/output (I/O) bound workloads.



- Released 2012, LevelDB traces back to 2004
- Built and maintained by full-time engineering team
- Continuous performance improvements

SplinterDB Performance

32 2Ghz cores

Intel Optane 905P

Block-addressable
NVMe

24B keys
100B values

Small KV-pairs

4GiB RAM
80GiB dataset

Small cache
(using cgroup)

SplinterDB Performance

32 2Ghz cores

Intel Optane 905P

Block-addressable
NVMe

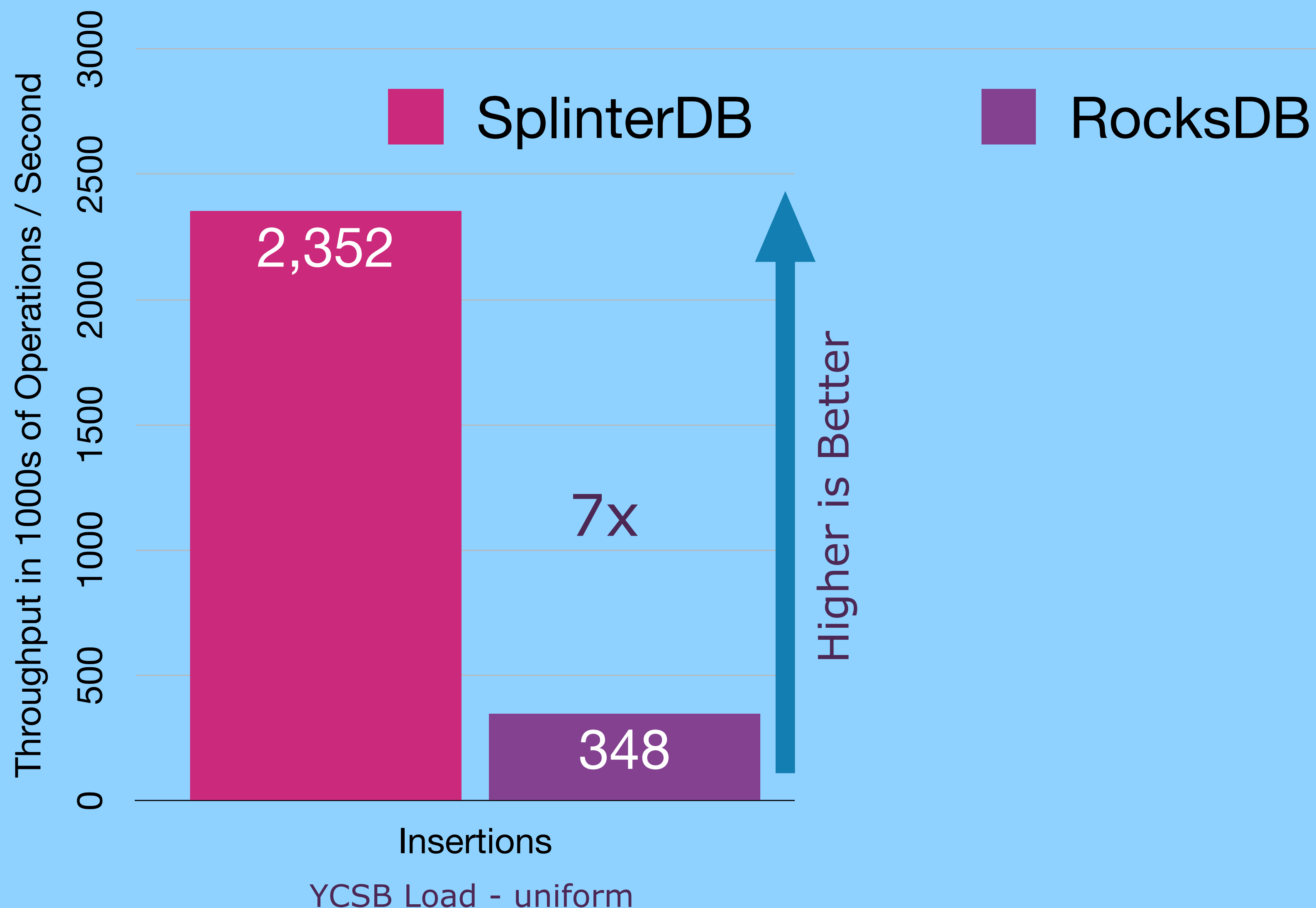
24B keys
100B values

Small KV-pairs

4GiB RAM
80GiB dataset

Small cache
(using cgroup)

Basic Operation Throughput



SplinterDB Performance

32 2Ghz cores

Intel Optane 905P

Block-addressable
NVMe

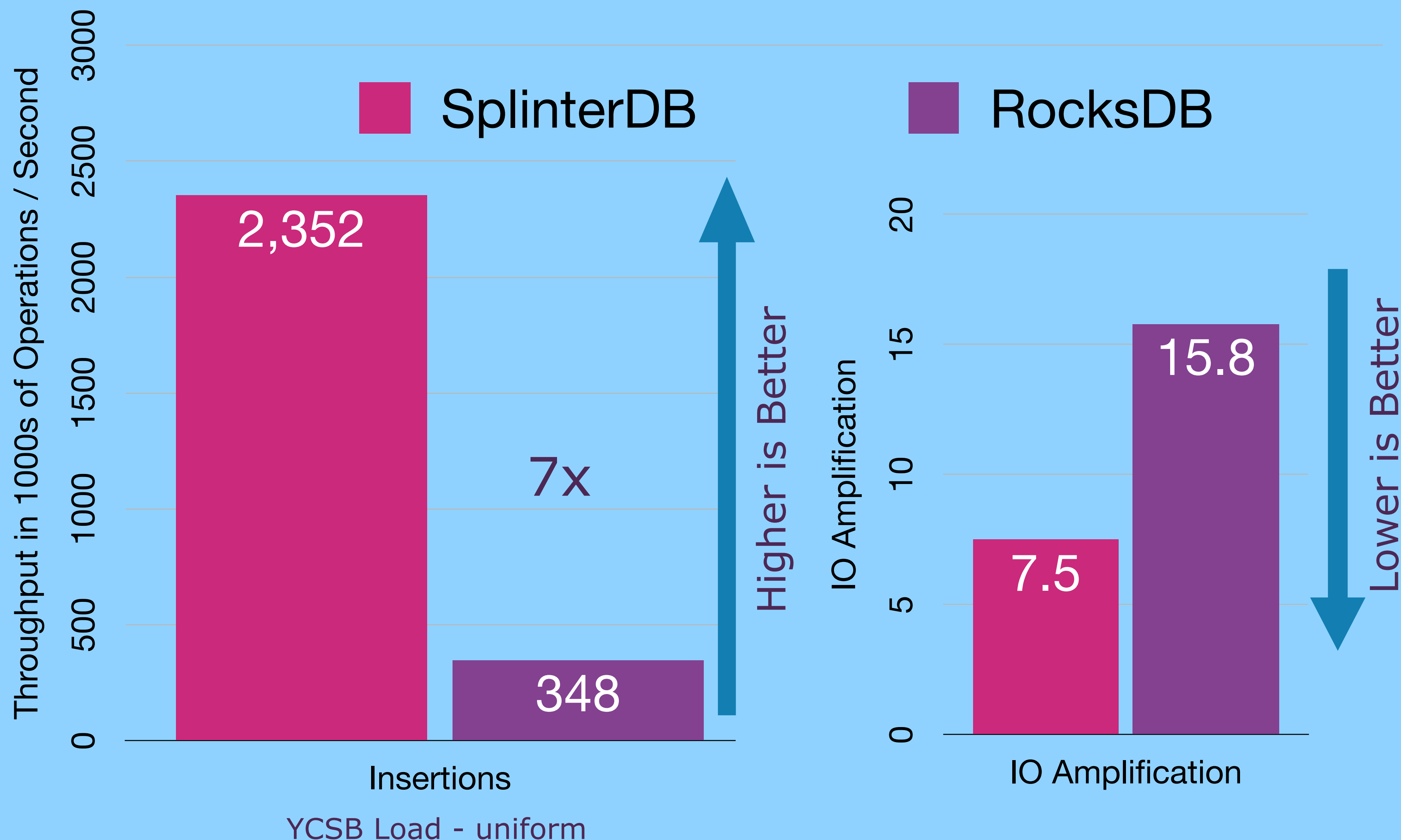
24B keys
100B values

Small KV-pairs

4GiB RAM
80GiB dataset

Small cache
(using cgroup)

Basic Operation Throughput



SplinterDB Performance

32 2Ghz cores

Intel Optane 905P

Block-addressable
NVMe

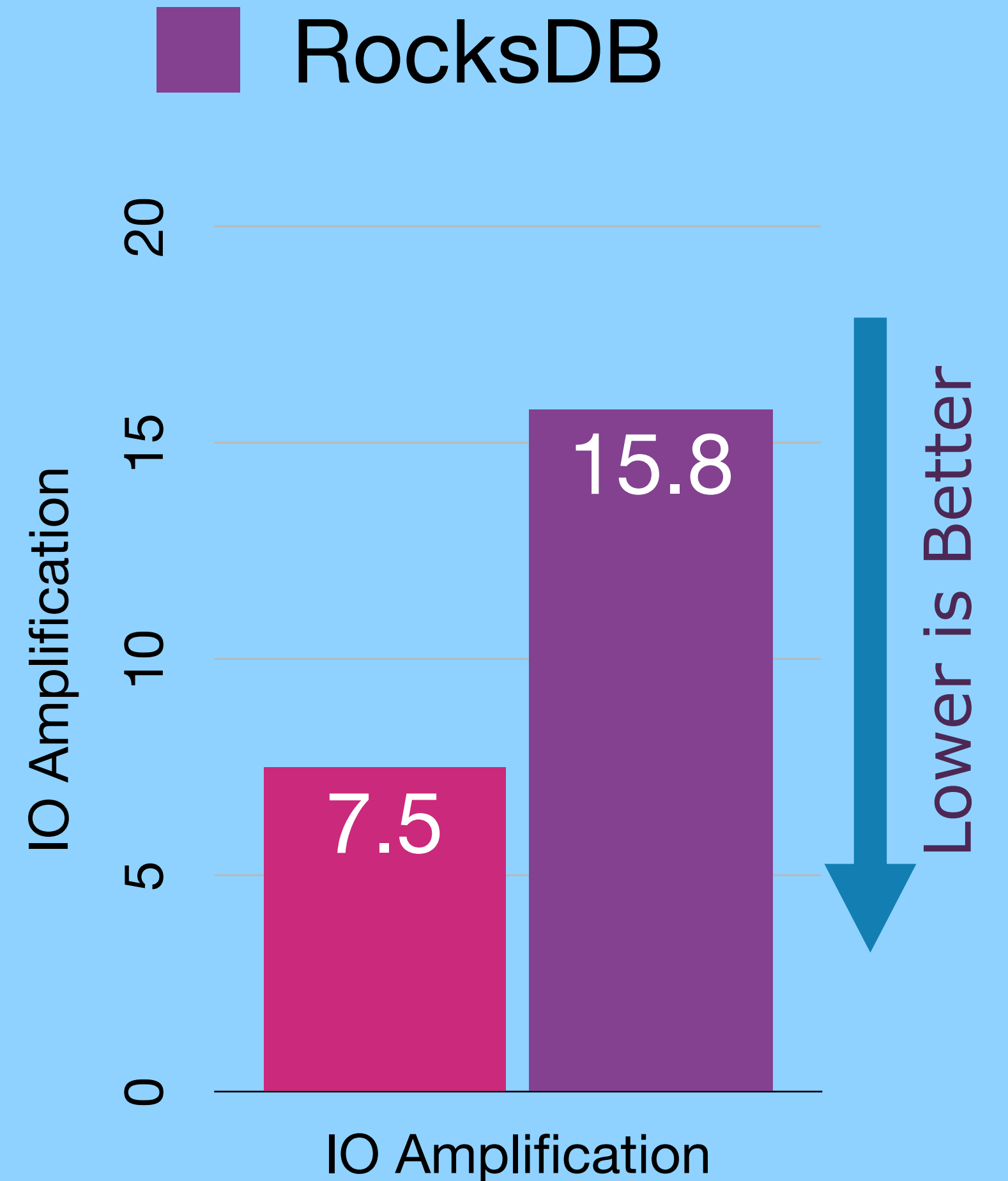
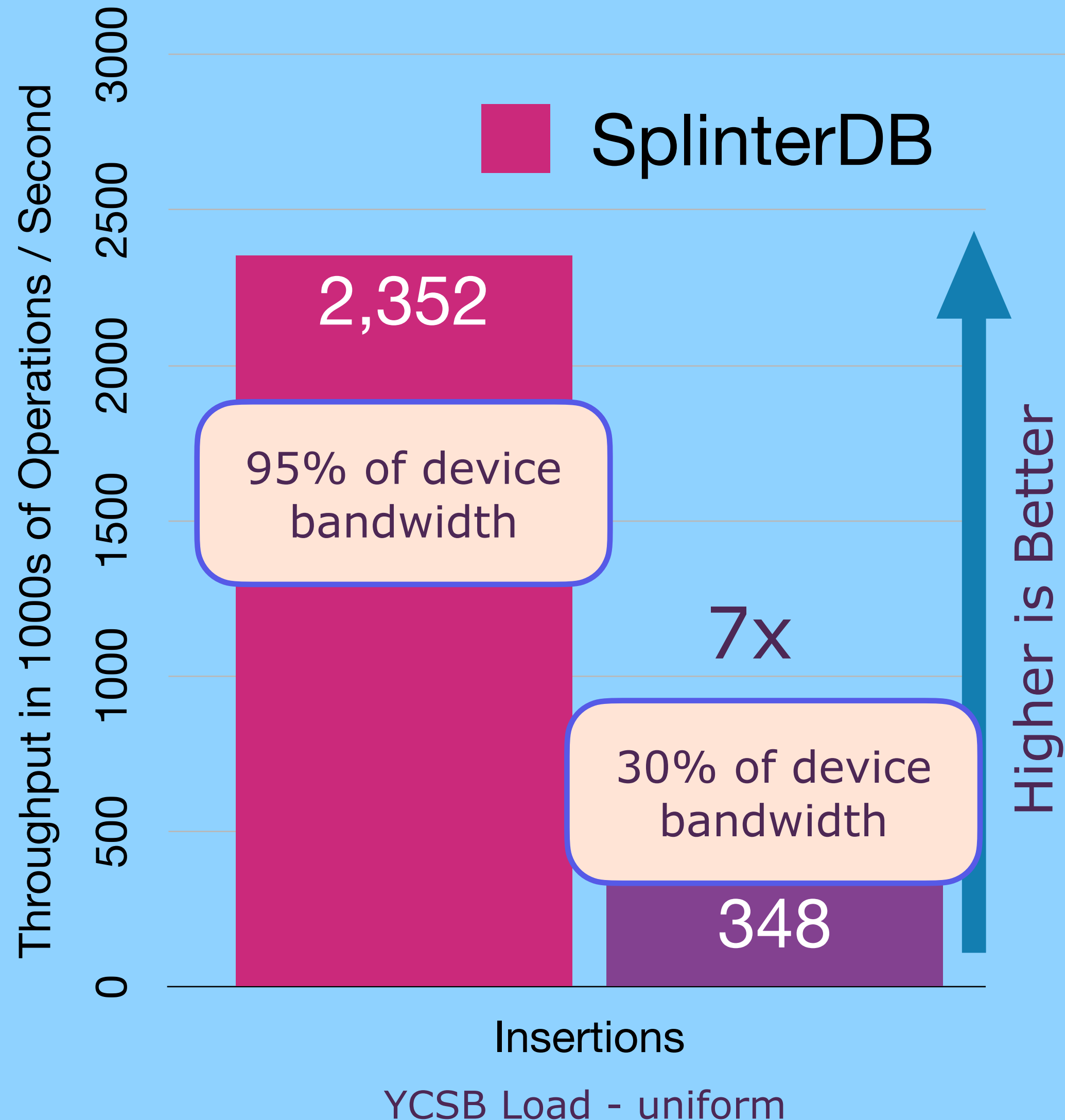
24B keys
100B values

Small KV-pairs

4GiB RAM
80GiB dataset

Small cache
(using cgroup)

Basic Operation Throughput



SplinterDB Performance

32 2Ghz cores

Intel Optane 905P

Block-addressable
NVMe

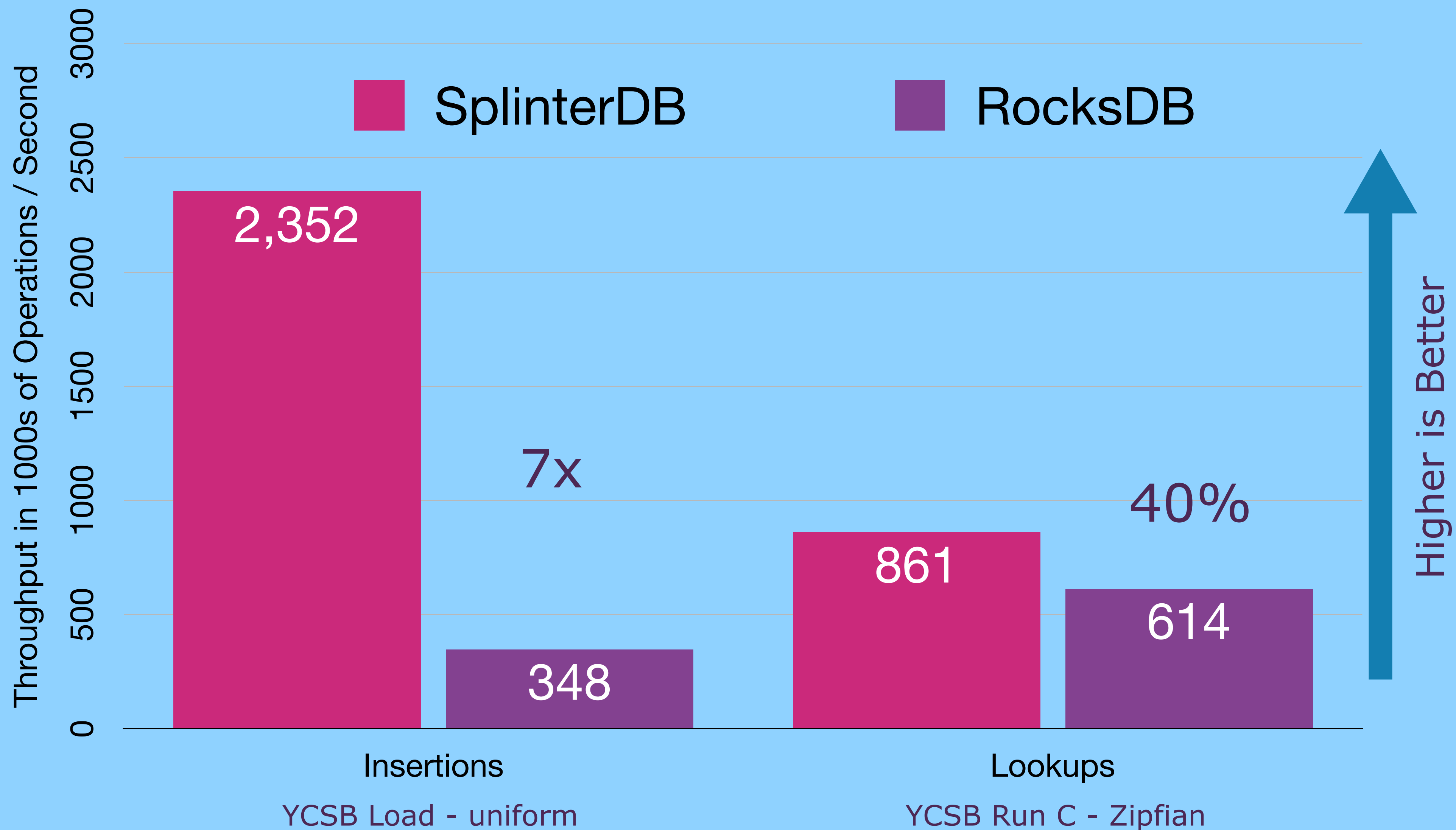
24B keys
100B values

Small KV-pairs

4GiB RAM
80GiB dataset

Small cache
(using cgroup)

Basic Operation Throughput



SplinterDB Performance

32 2Ghz cores

Intel Optane 905P

Block-addressable
NVMe

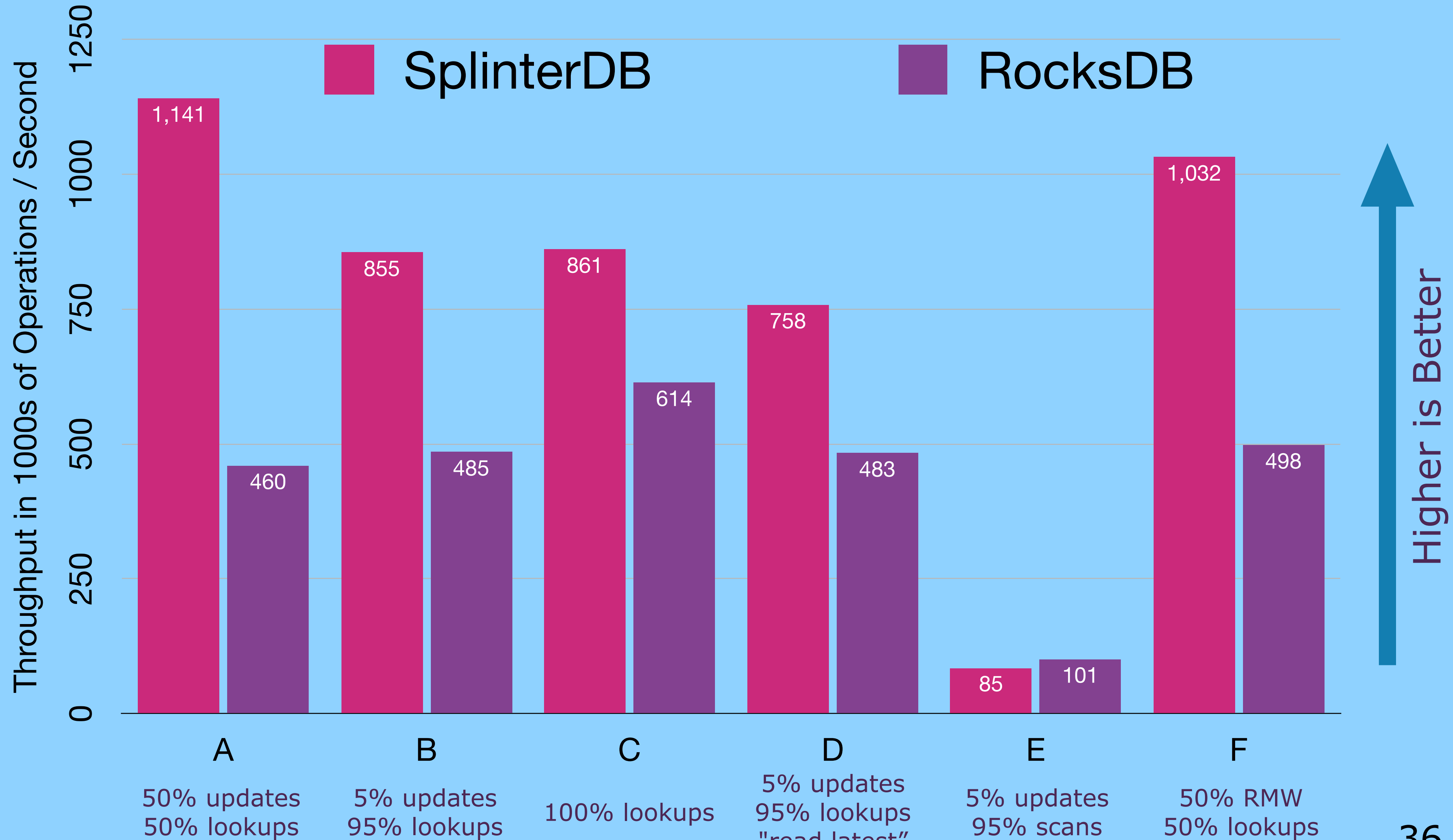
24B keys
100B values

Small KV-pairs

4GiB RAM
80GiB dataset

Small cache
(using cgroup)

YCSB Application Benchmark



In this talk

Fast Storage
(NVMe)

SplinterDB

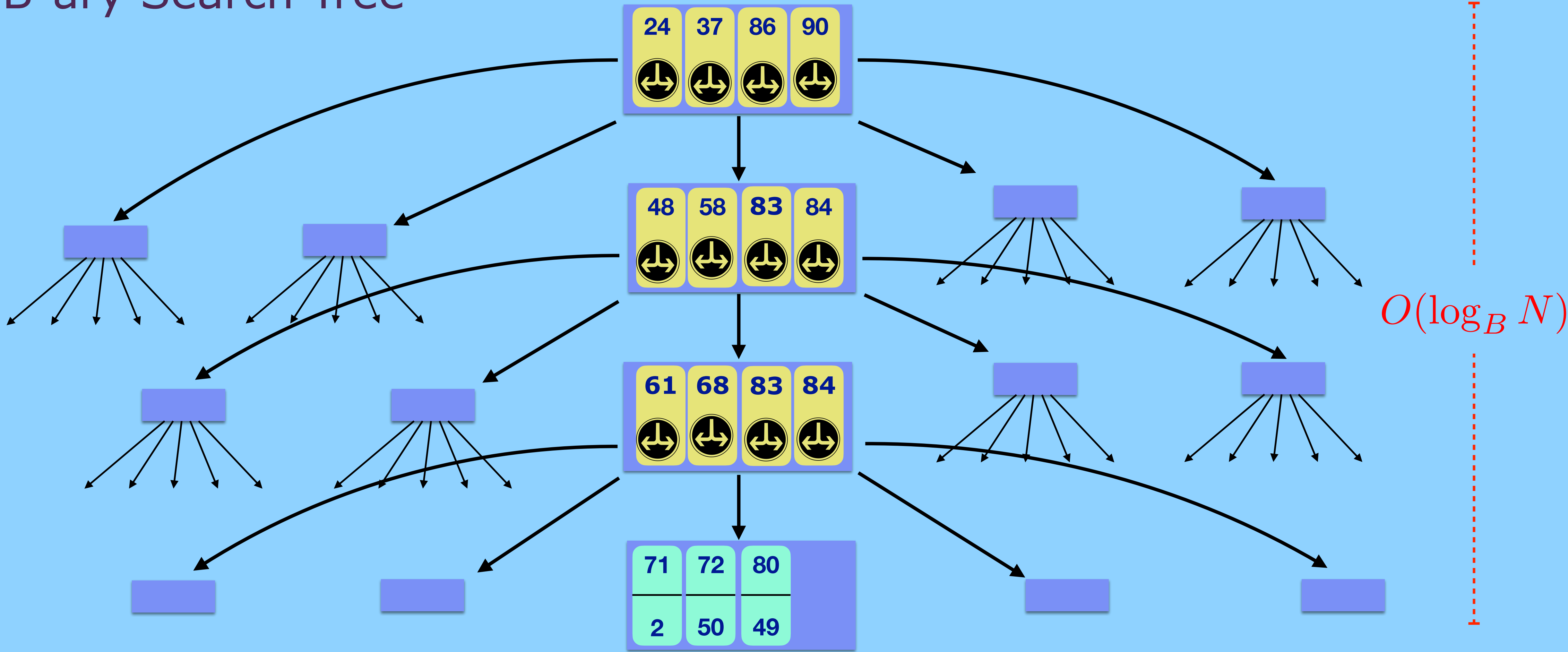
Data Structures

Flush-then-Compact

B-Trees

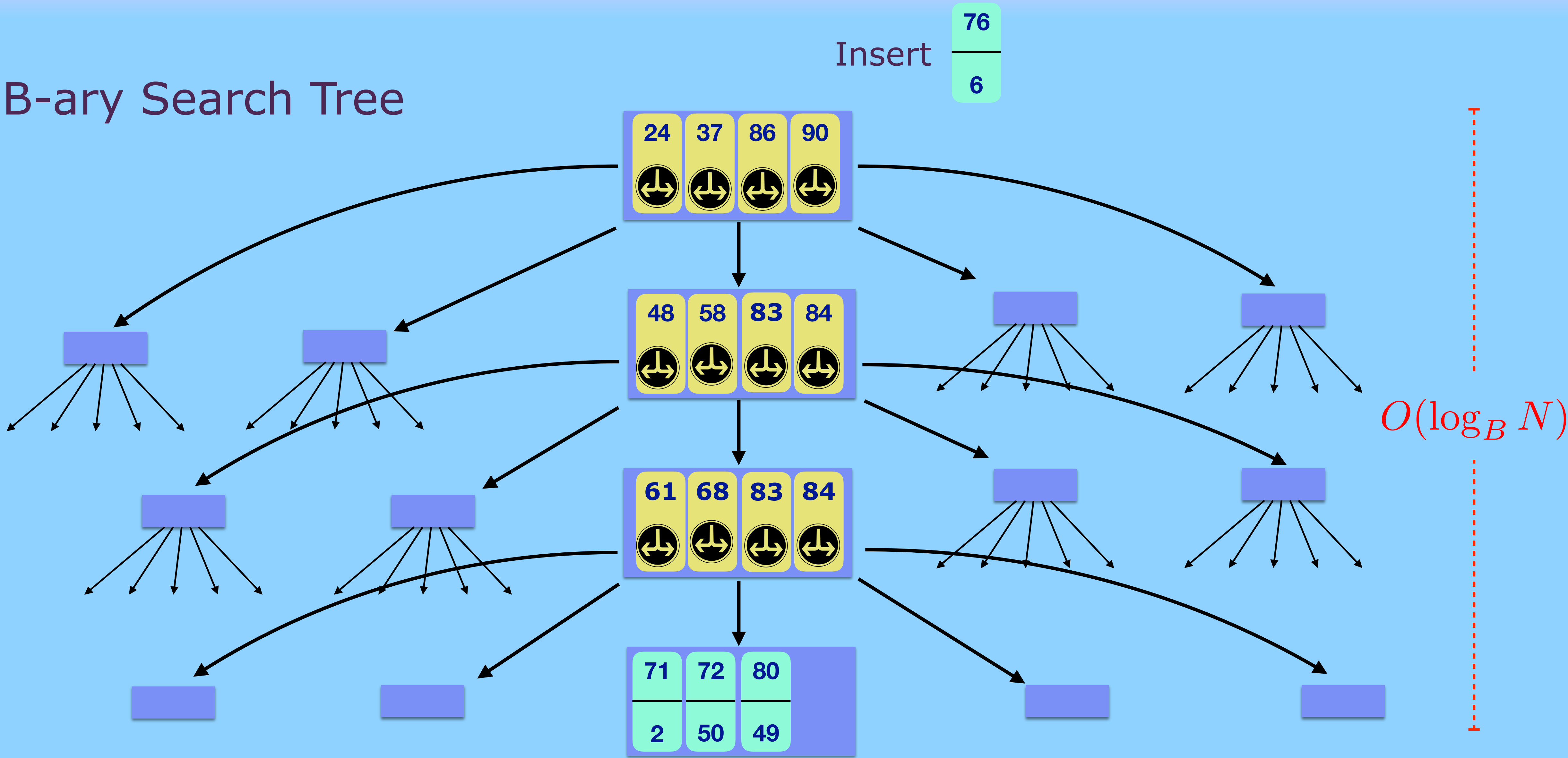
B-Trees

B-ary Search Tree



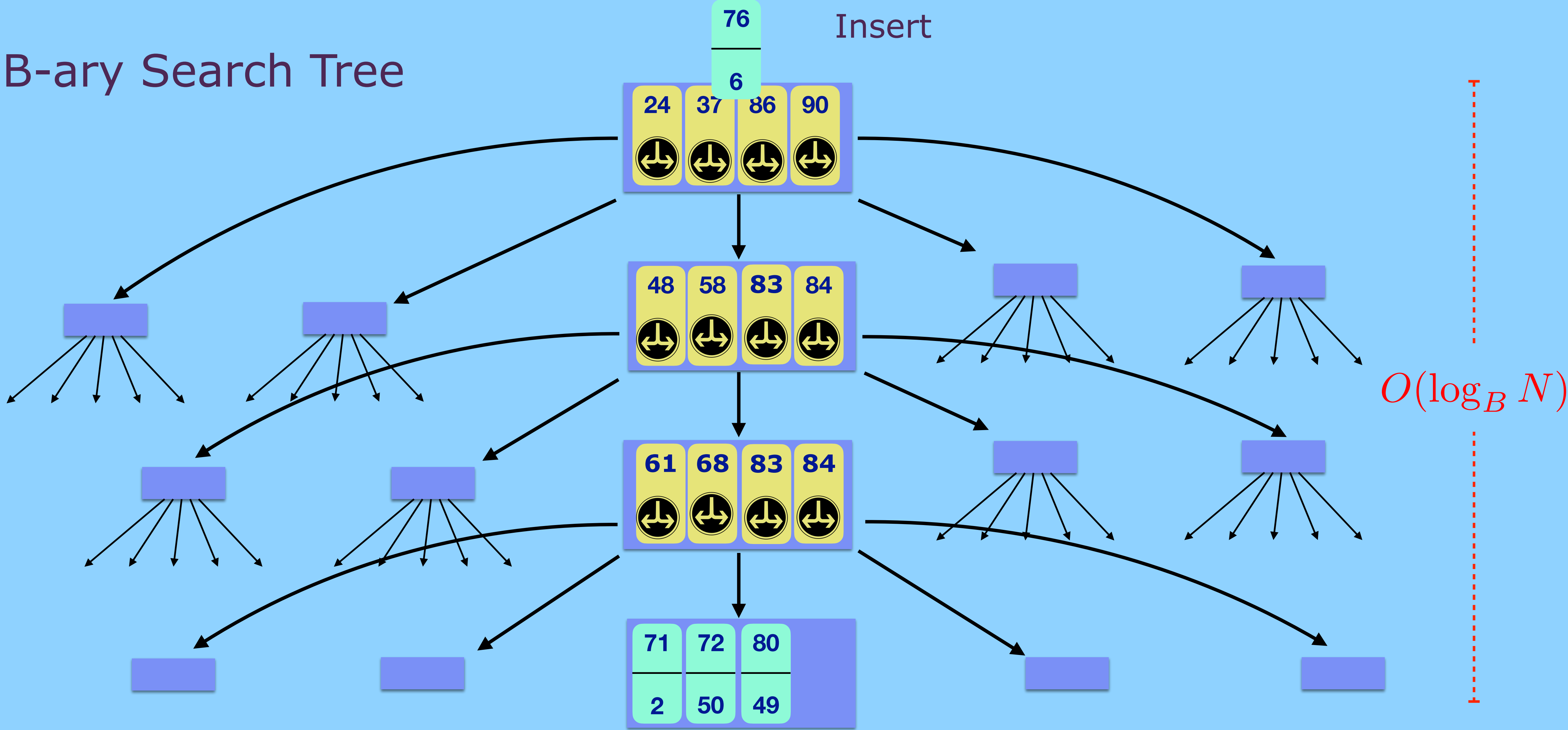
B-Trees

B-ary Search Tree



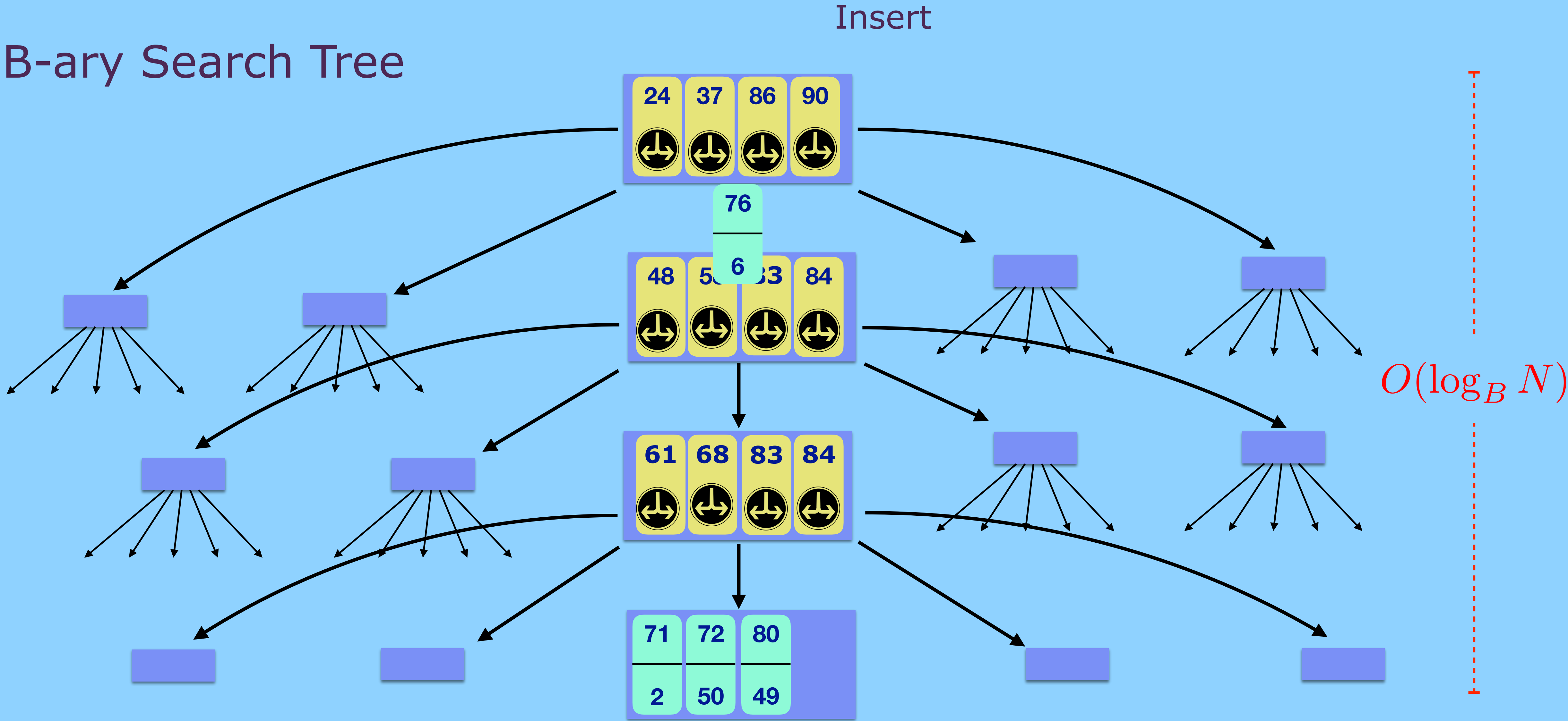
B-Trees

B-ary Search Tree



B-Trees

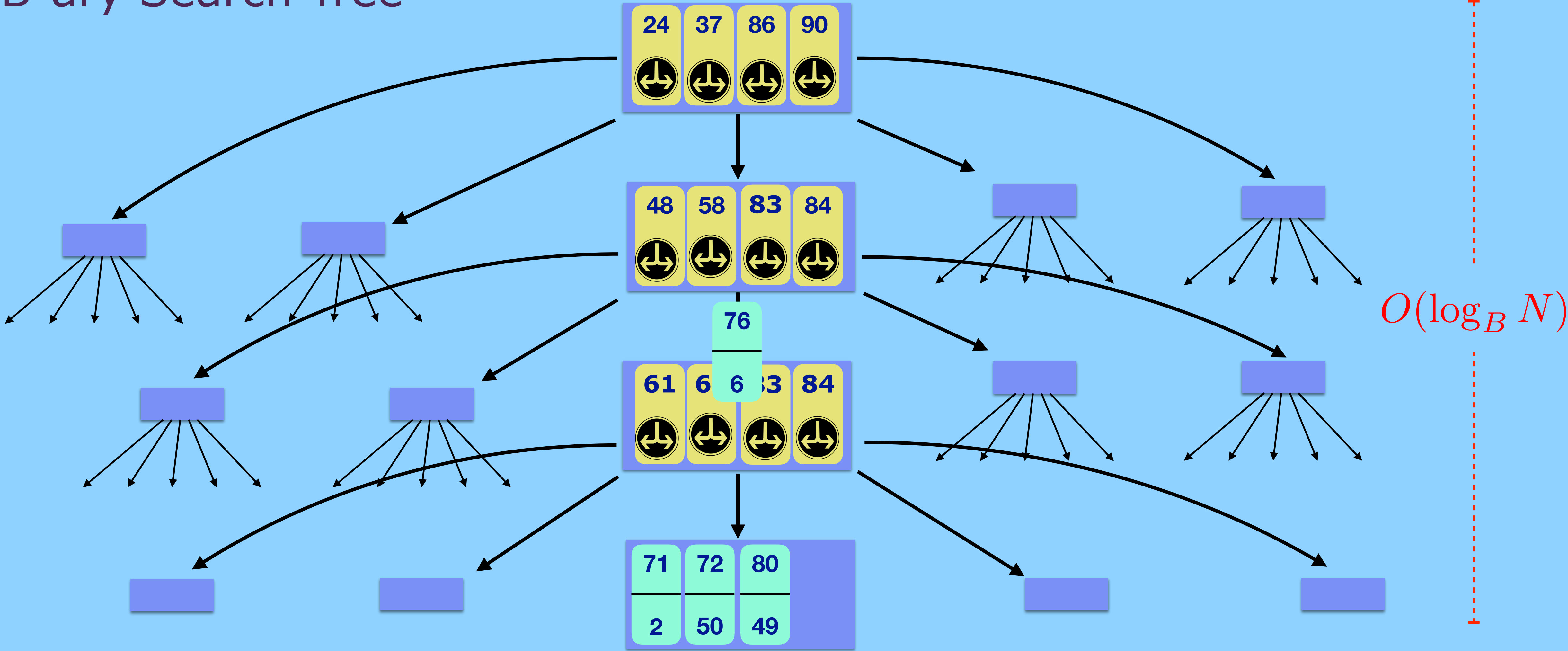
B-ary Search Tree



B-Trees

B-ary Search Tree

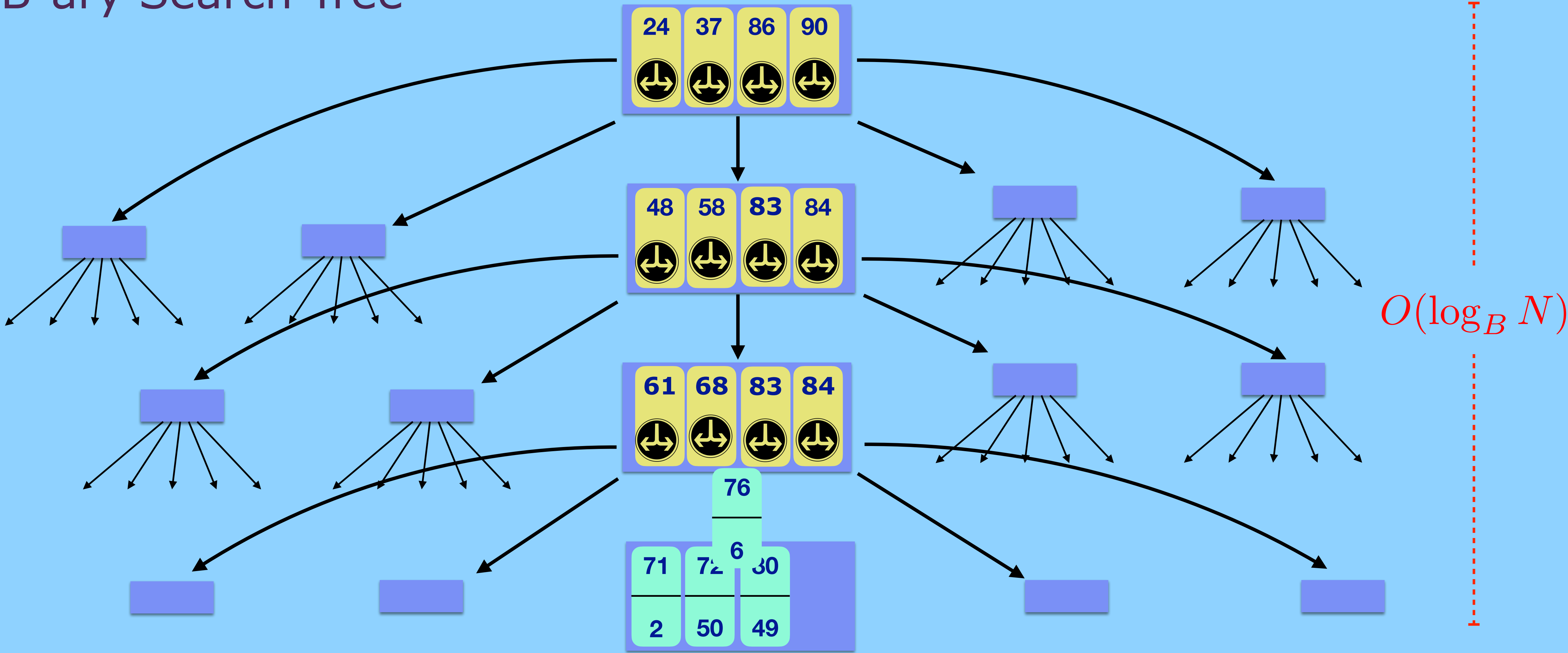
Insert



B-Trees

B-ary Search Tree

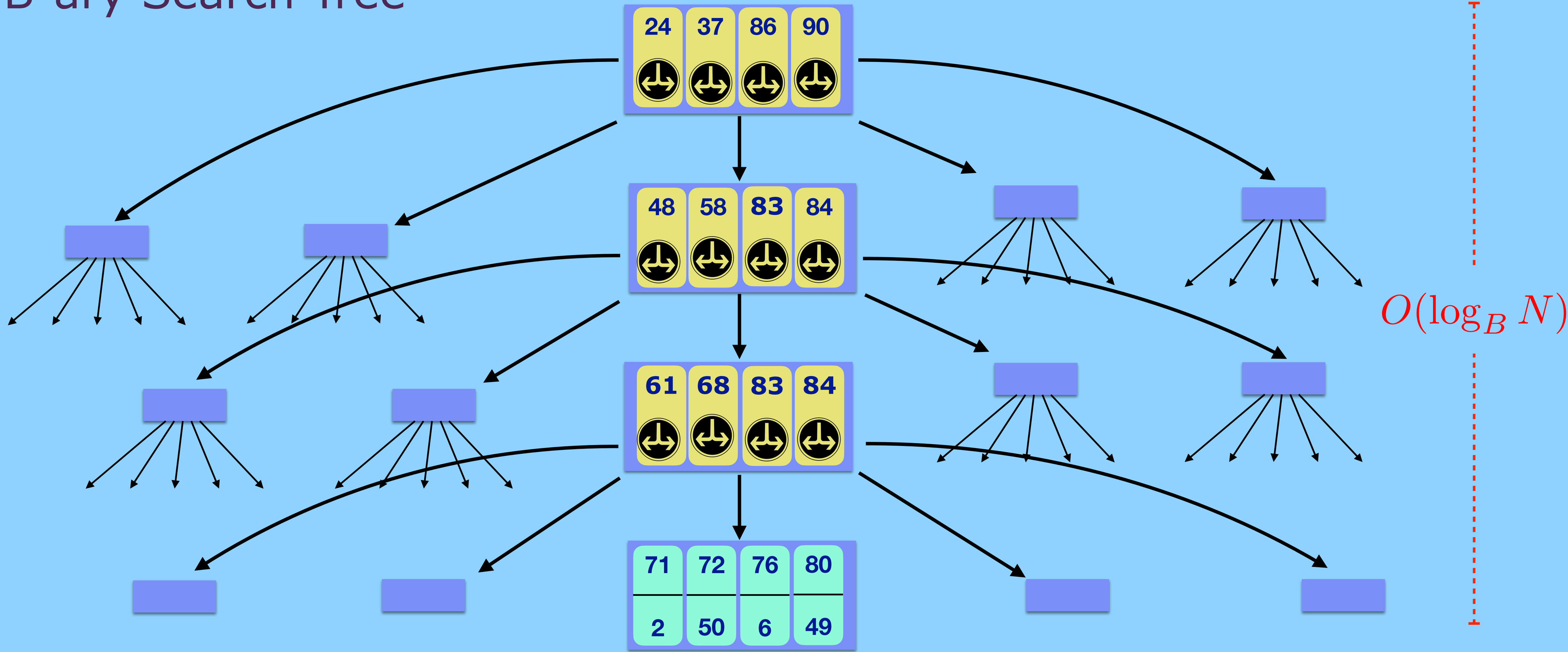
Insert



B-Trees

B-ary Search Tree

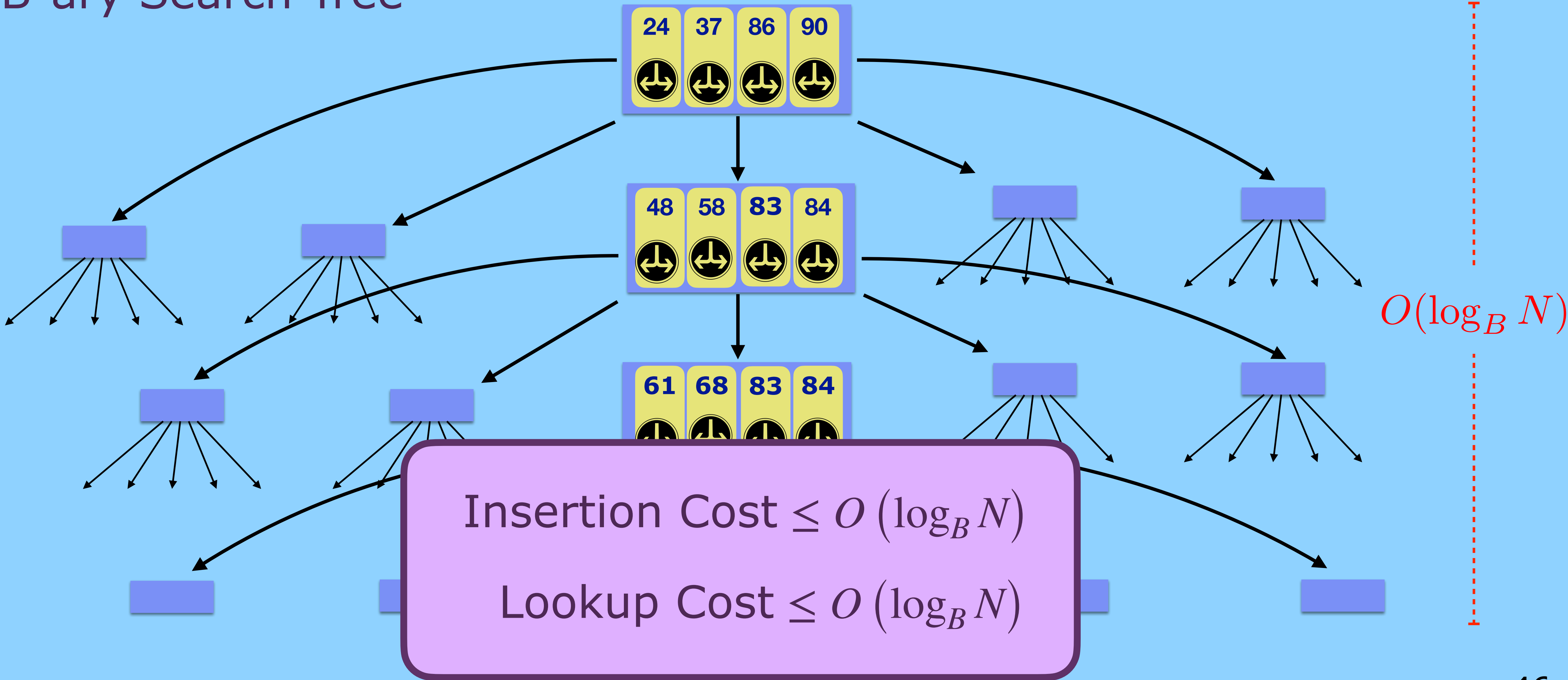
Insert



B-Trees

B-ary Search Tree

Insert

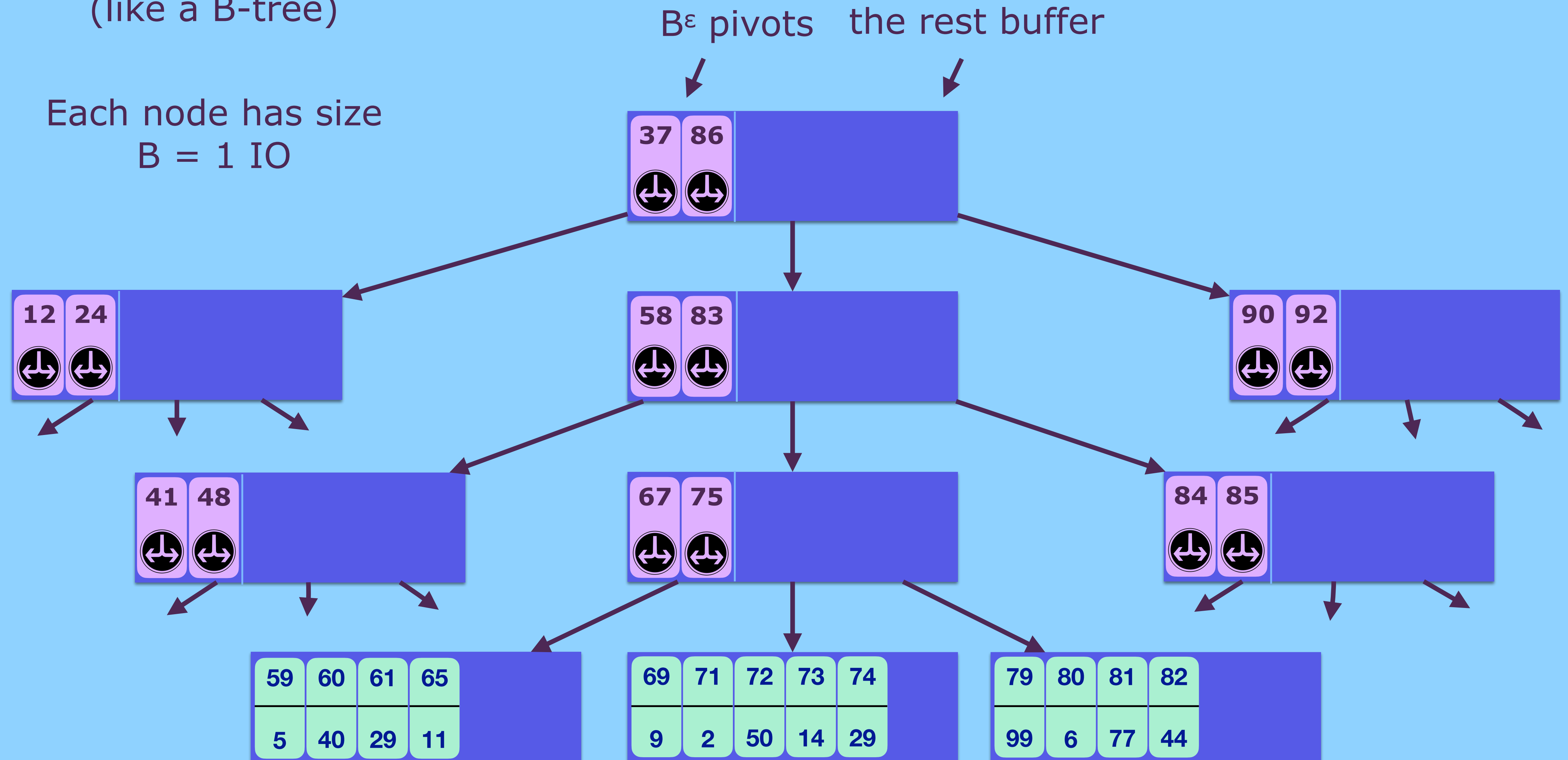


B^ϵ -Trees

B ϵ -Trees

A B ϵ -tree is a search tree
(like a B-tree)

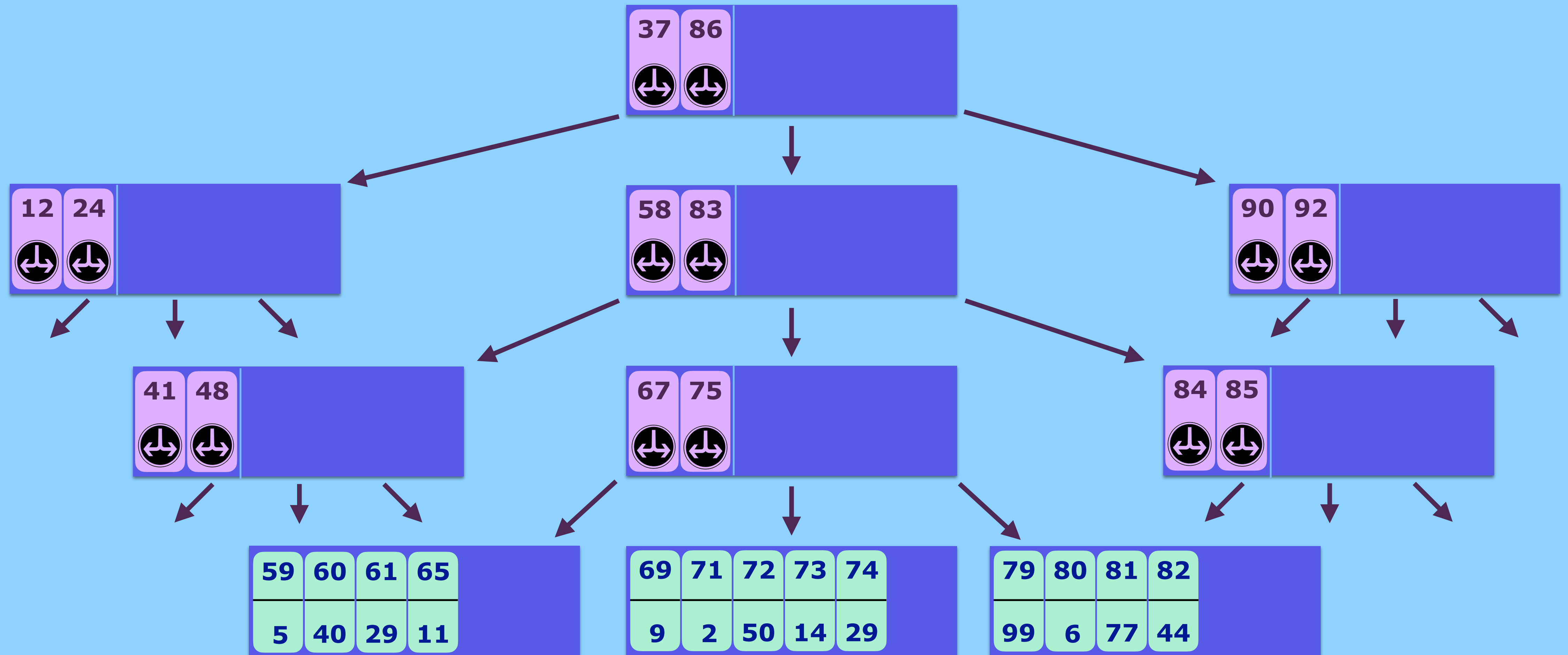
Each node has size
 $B = 1 \text{ IO}$



Insertions in B^ε -Trees

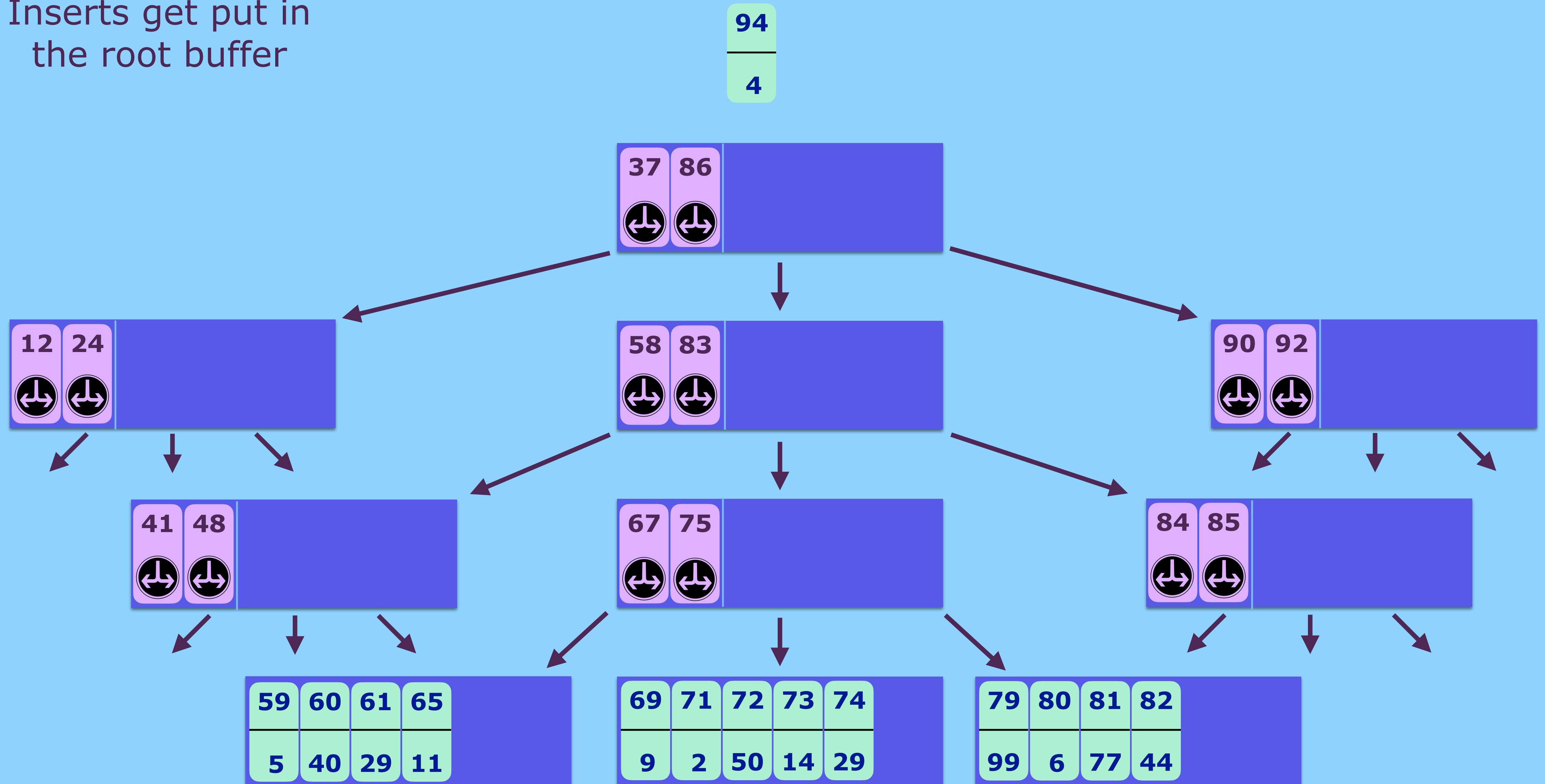
B ϵ -Trees

Inserts get put in
the root buffer



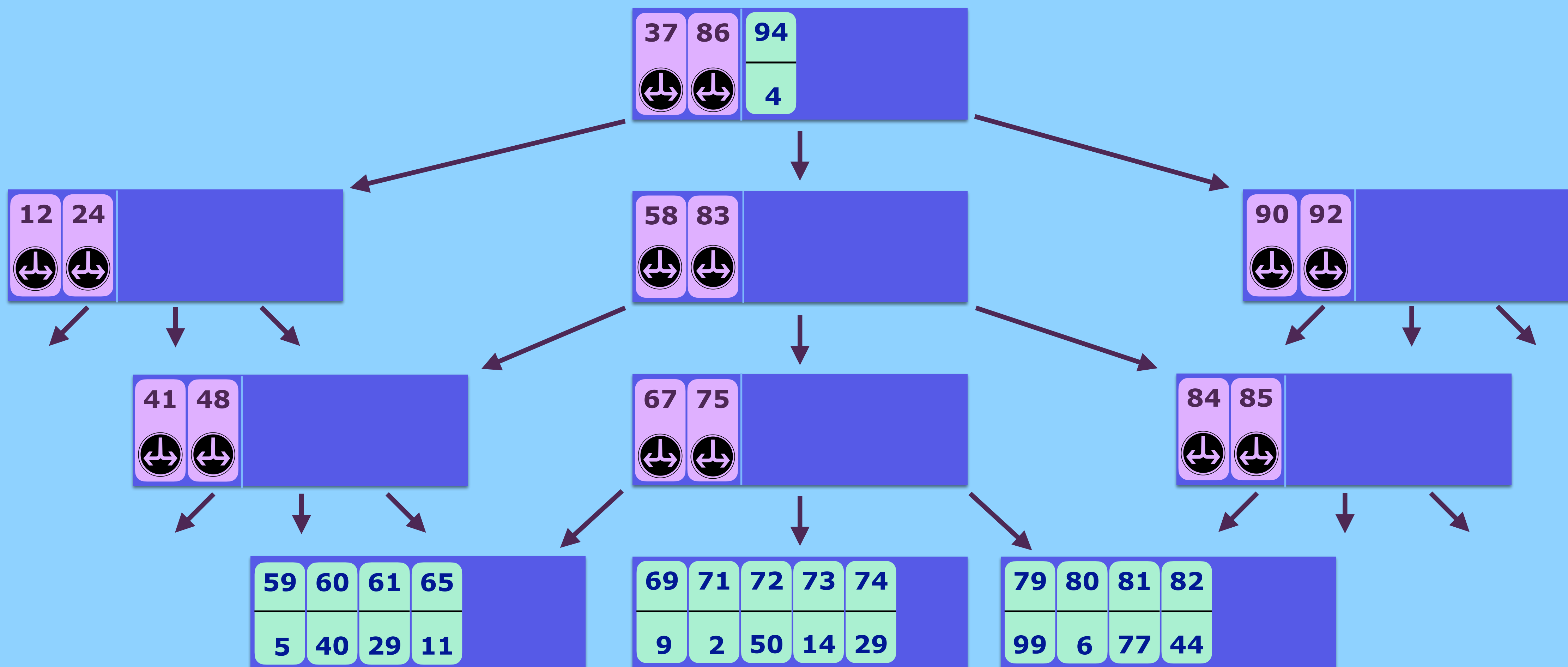
B ϵ -Trees

Inserts get put in
the root buffer



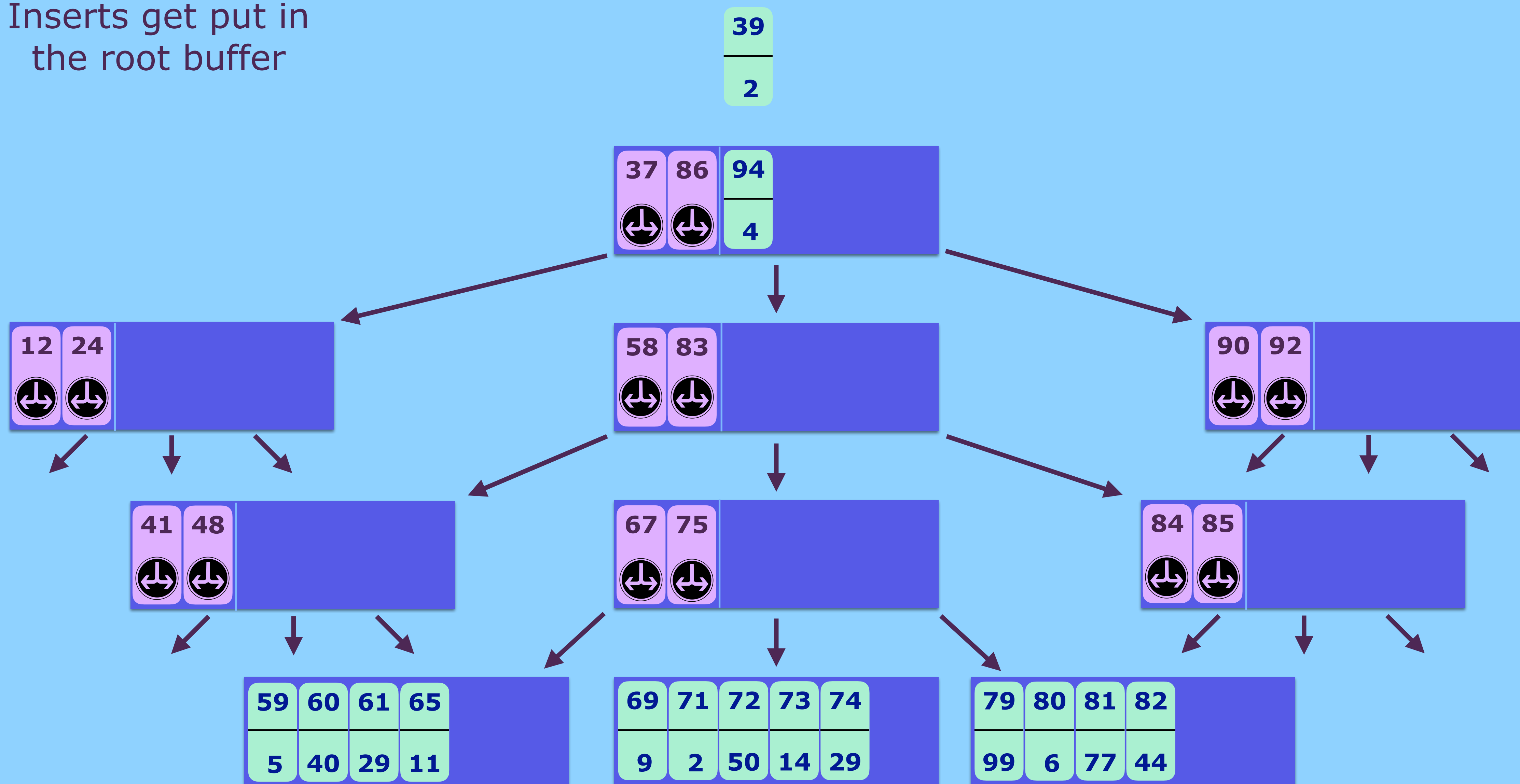
B ϵ -Trees

Inserts get put in
the root buffer



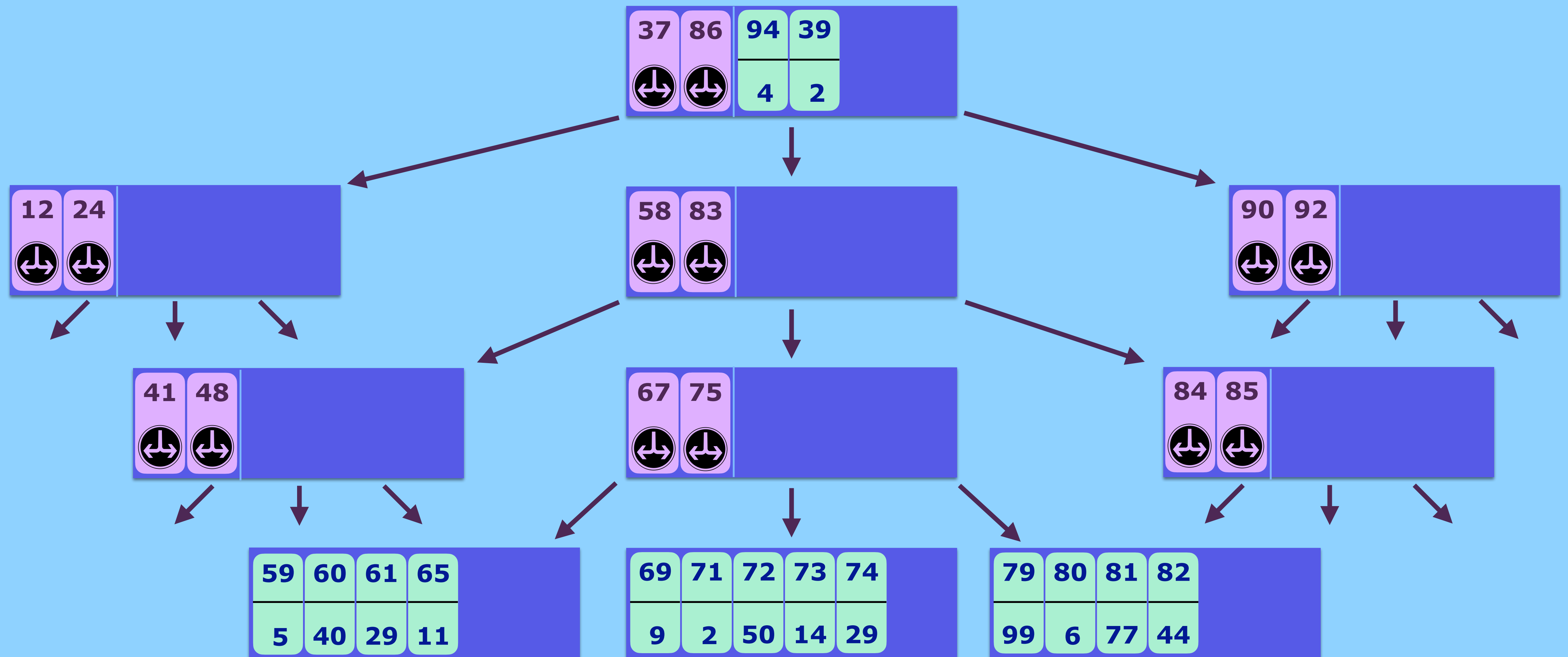
B ϵ -Trees

Inserts get put in
the root buffer



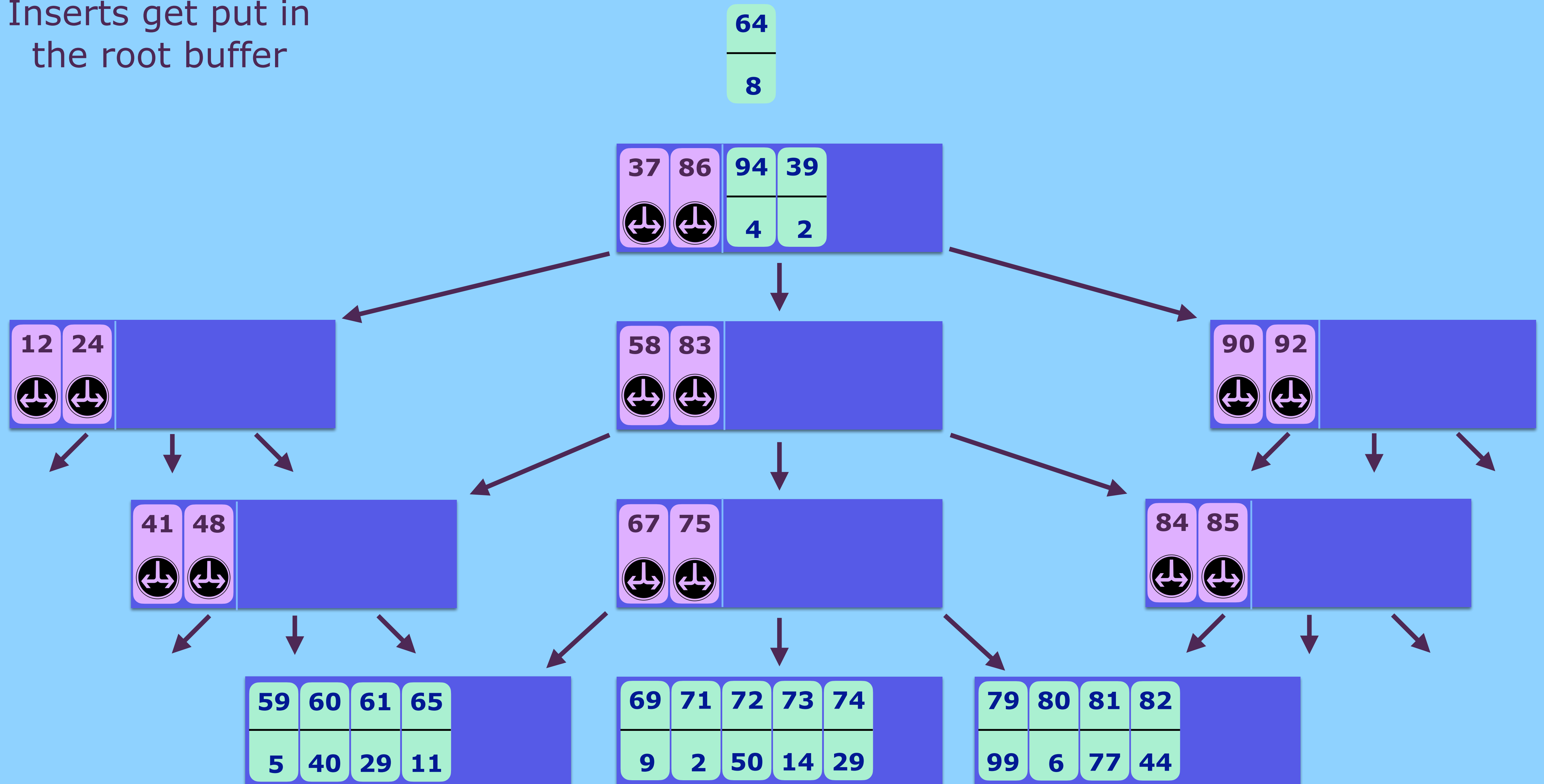
B ϵ -Trees

Inserts get put in
the root buffer



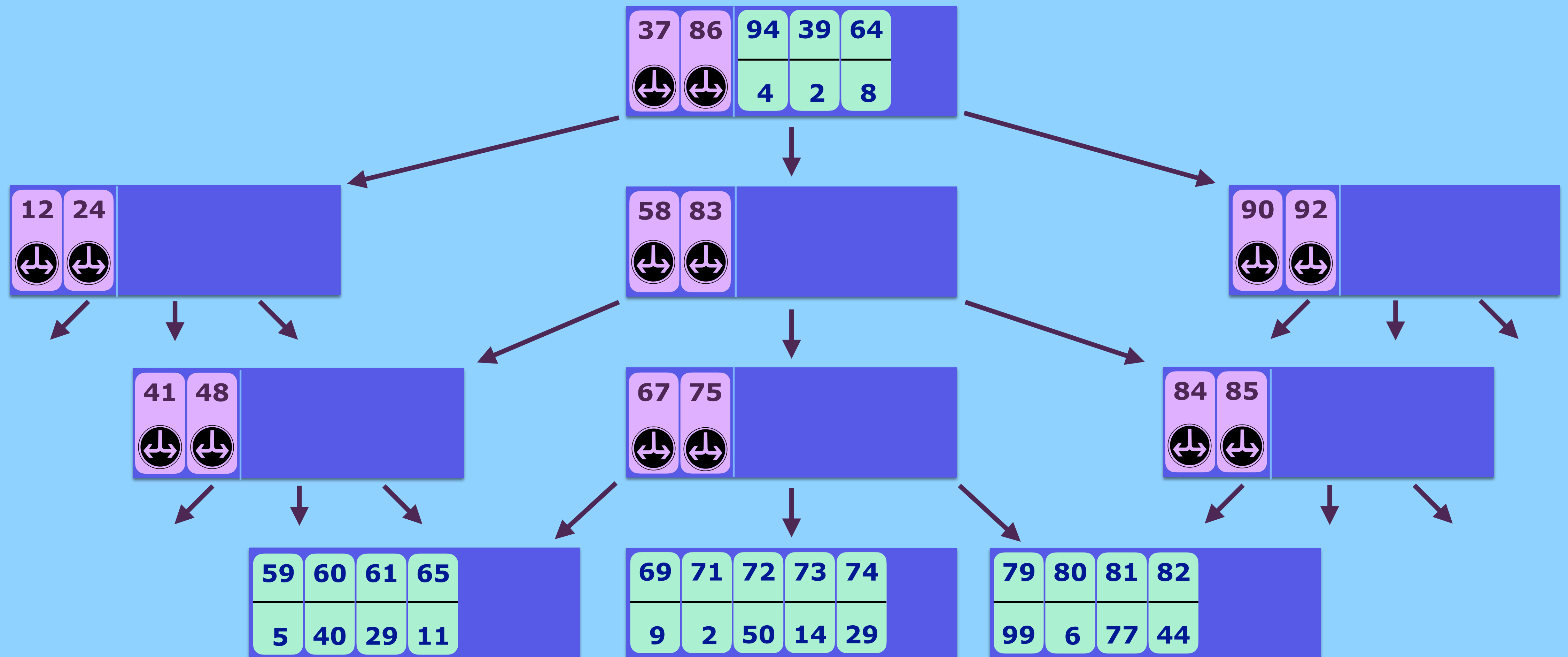
B ϵ -Trees

Inserts get put in
the root buffer



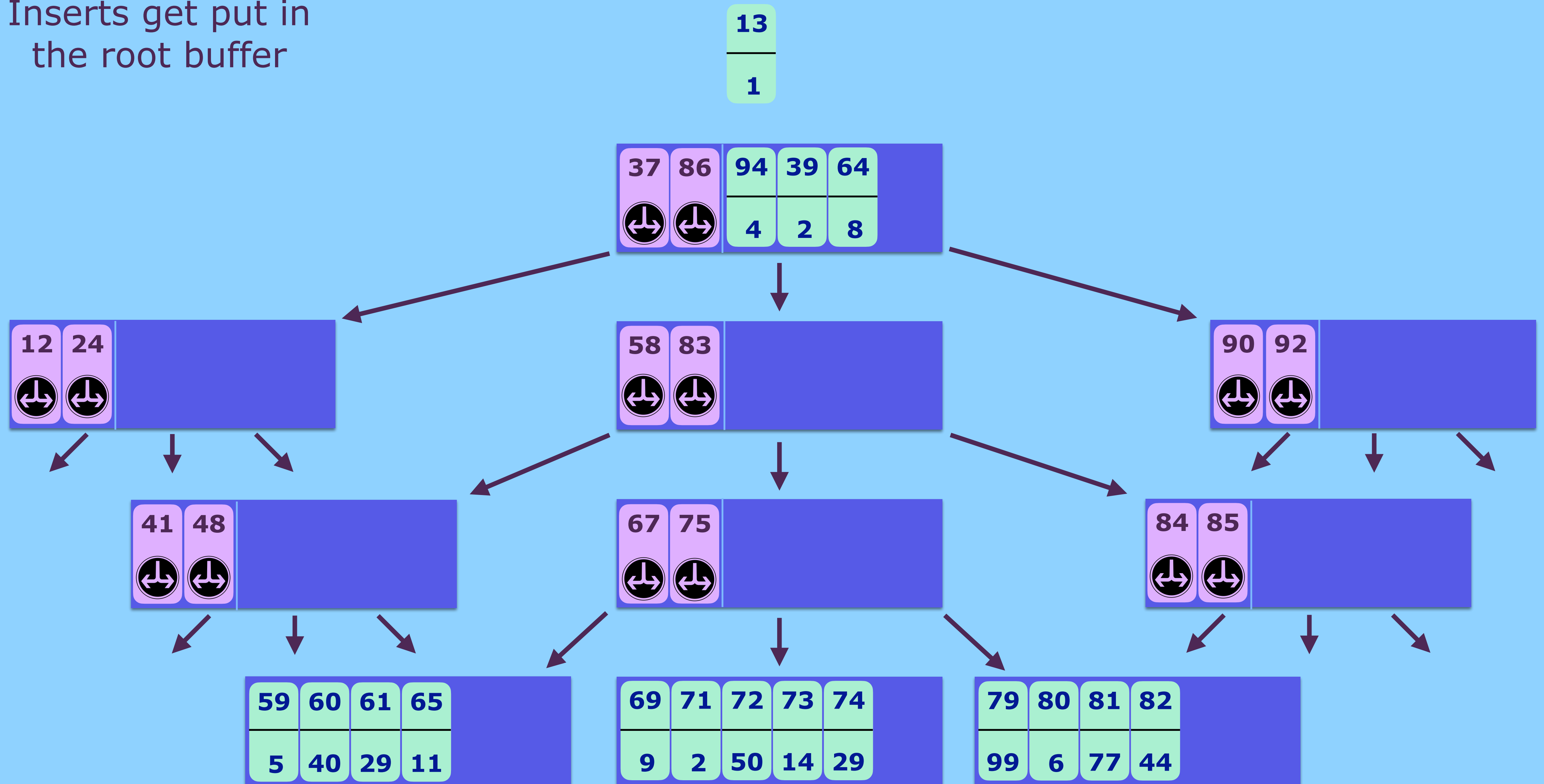
B ϵ -Trees

Inserts get put in
the root buffer



B ϵ -Trees

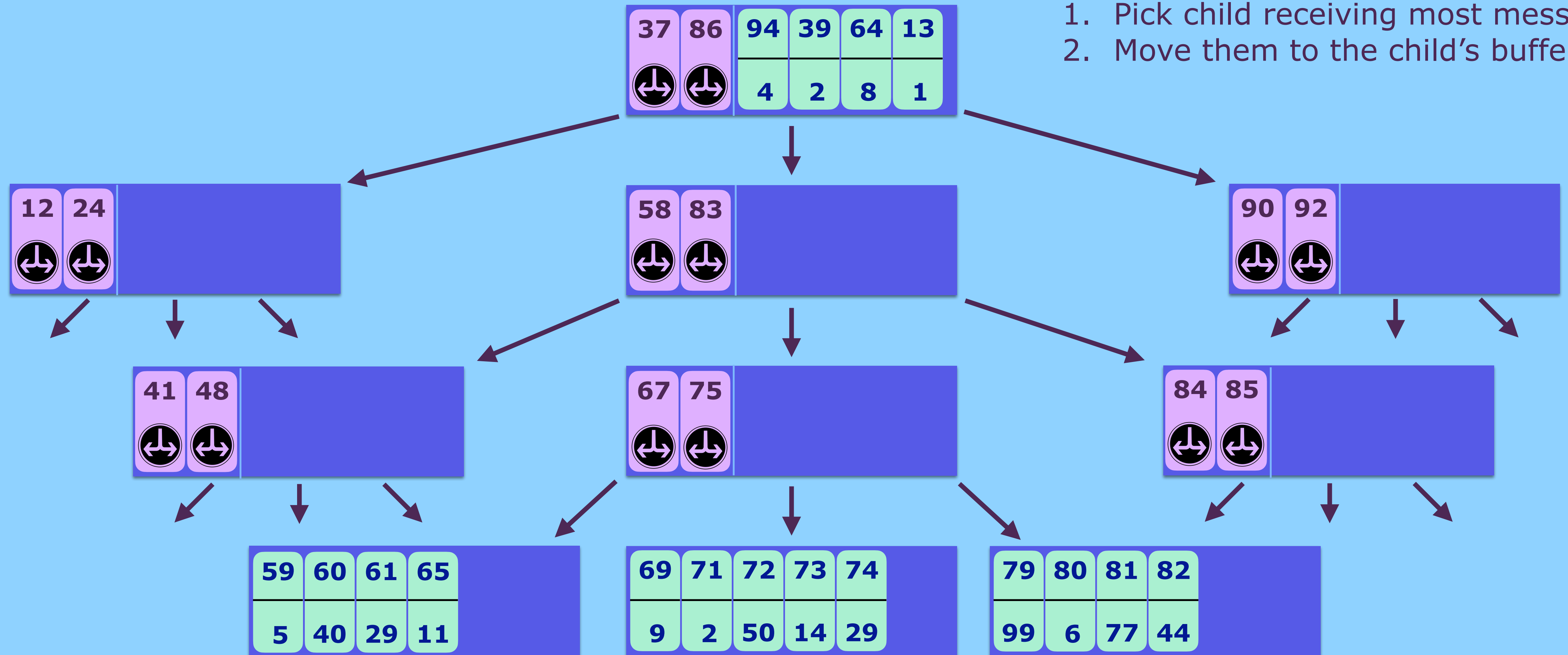
Inserts get put in
the root buffer



B ϵ -Trees

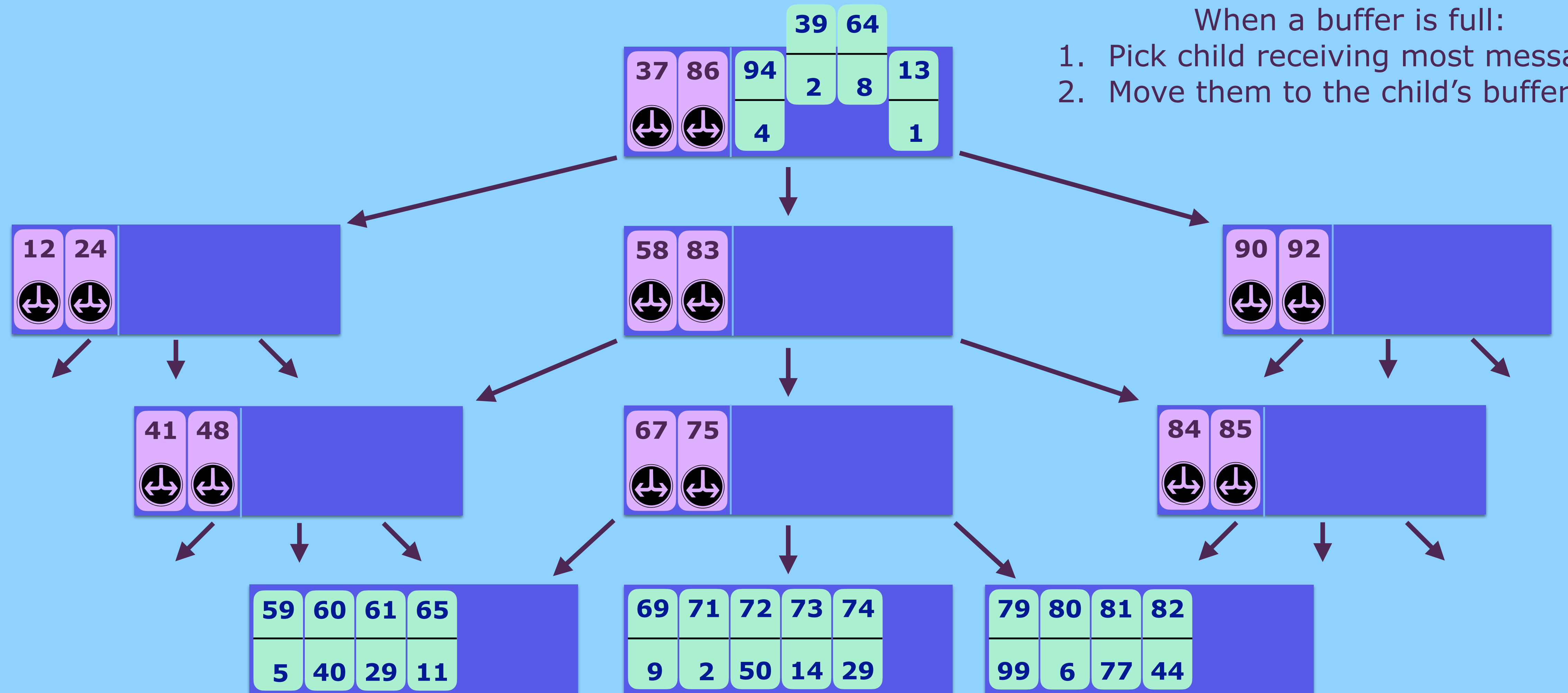
Inserts get put in the root buffer

- When a buffer is full:
1. Pick child receiving most messages
 2. Move them to the child's buffer



B ϵ -Trees

Inserts get put in the root buffer

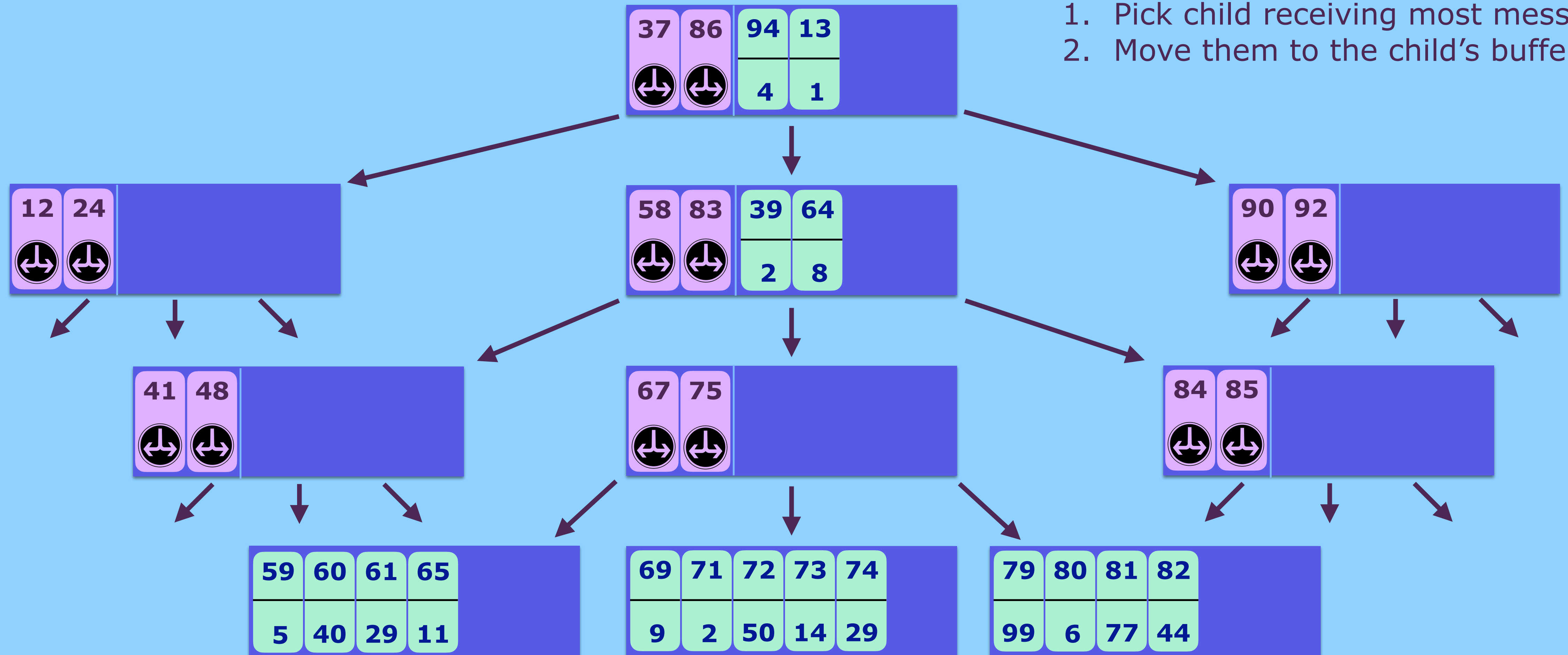


- When a buffer is full:
1. Pick child receiving most messages
 2. Move them to the child's buffer

B ϵ -Trees

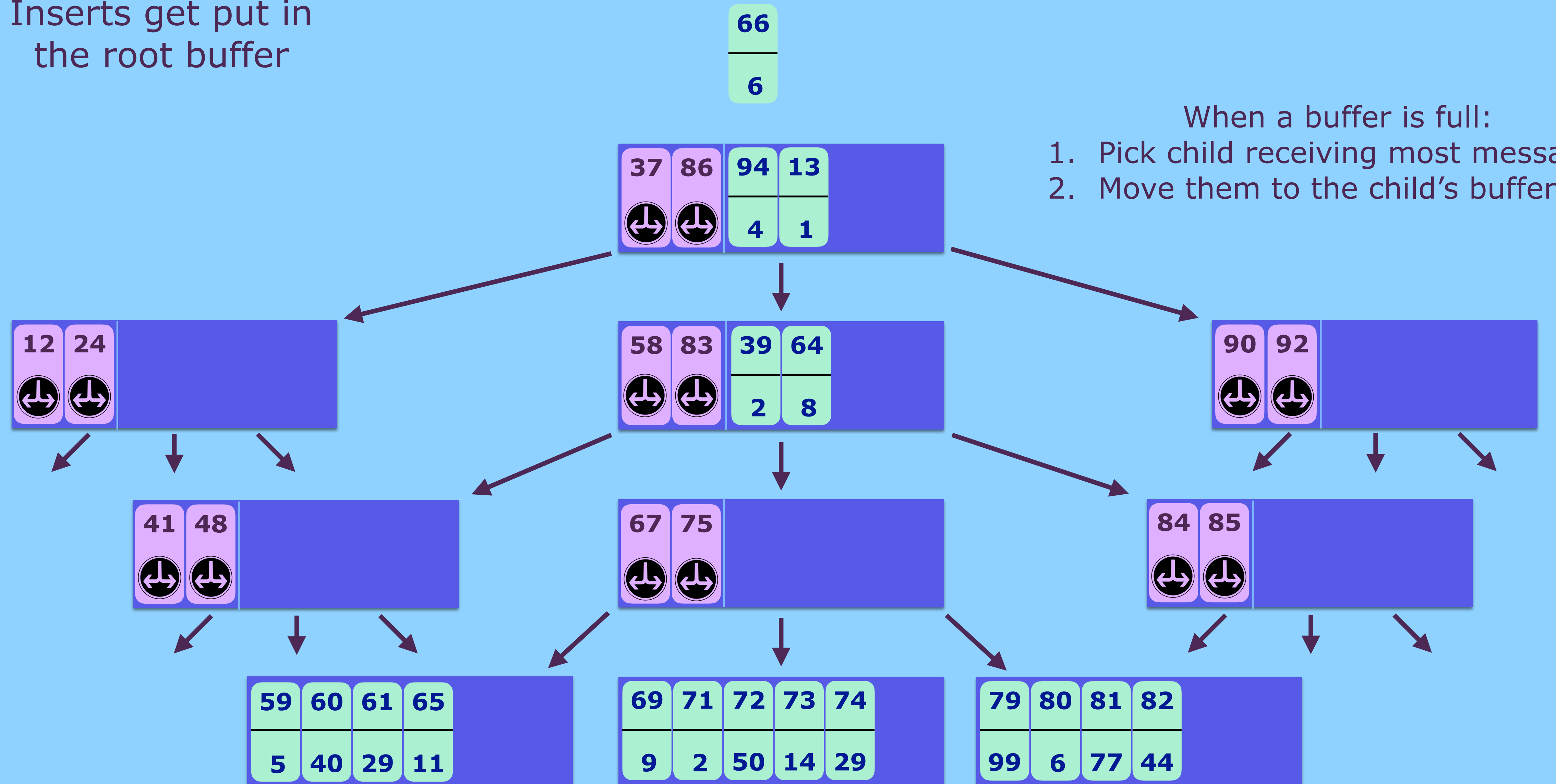
Inserts get put in the root buffer

- When a buffer is full:
1. Pick child receiving most messages
 2. Move them to the child's buffer



B ϵ -Trees

Inserts get put in the root buffer

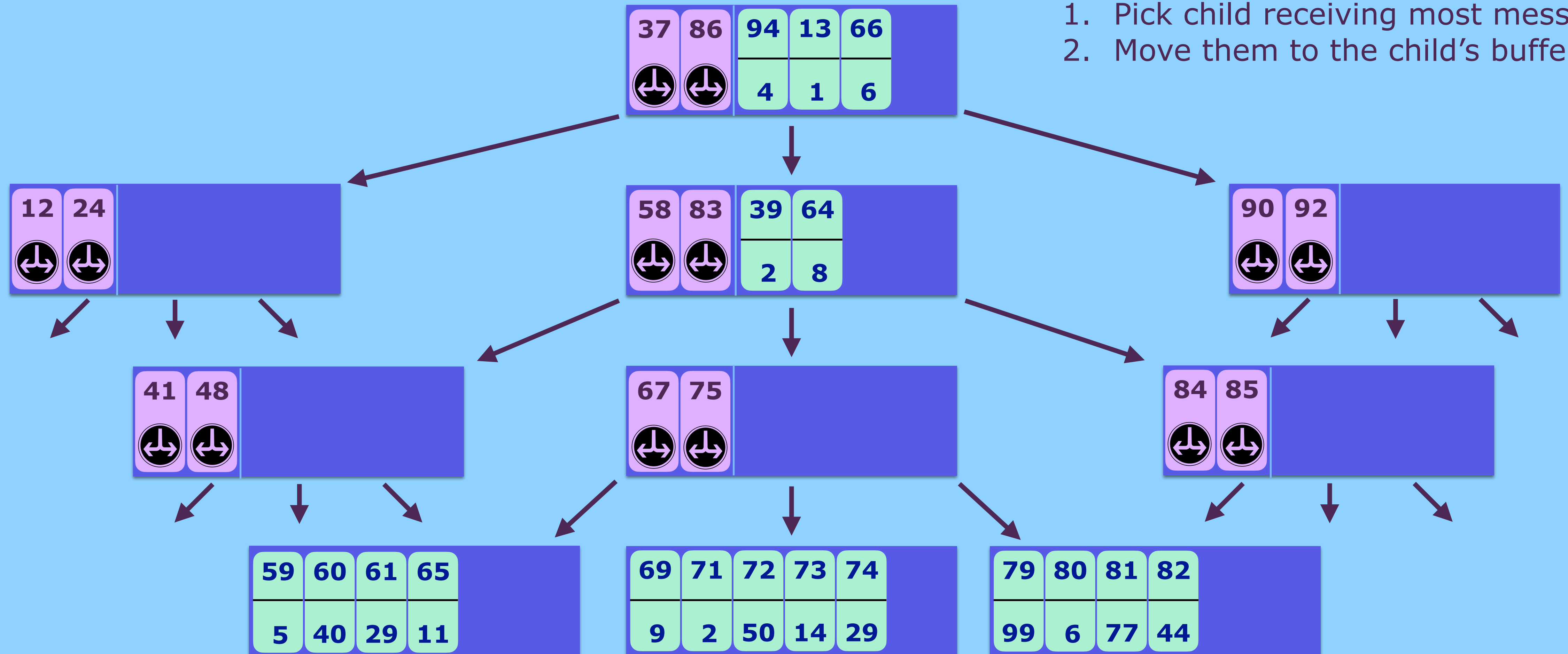


- When a buffer is full:
1. Pick child receiving most messages
 2. Move them to the child's buffer

B ϵ -Trees

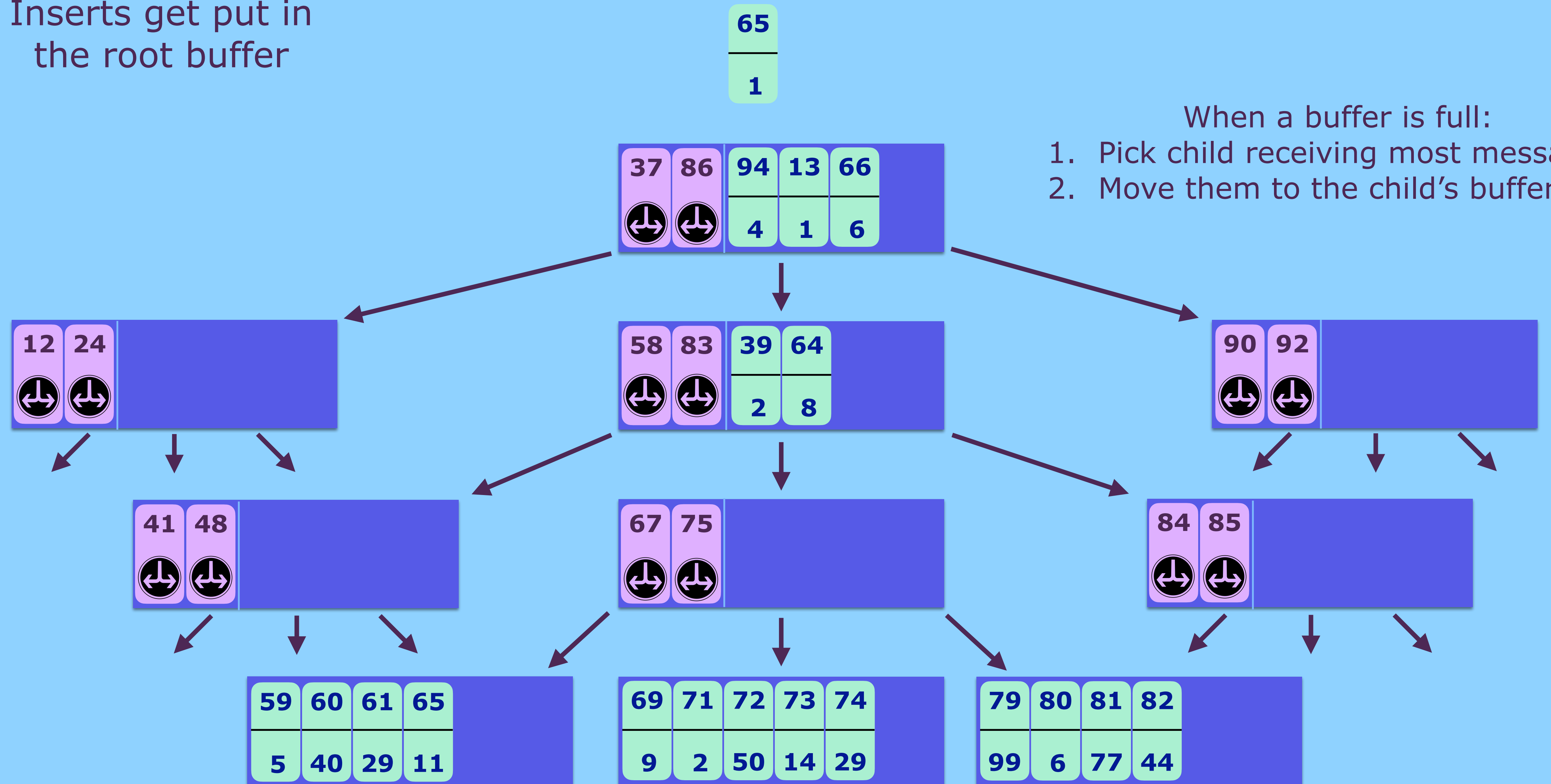
Inserts get put in the root buffer

- When a buffer is full:
1. Pick child receiving most messages
 2. Move them to the child's buffer



B ϵ -Trees

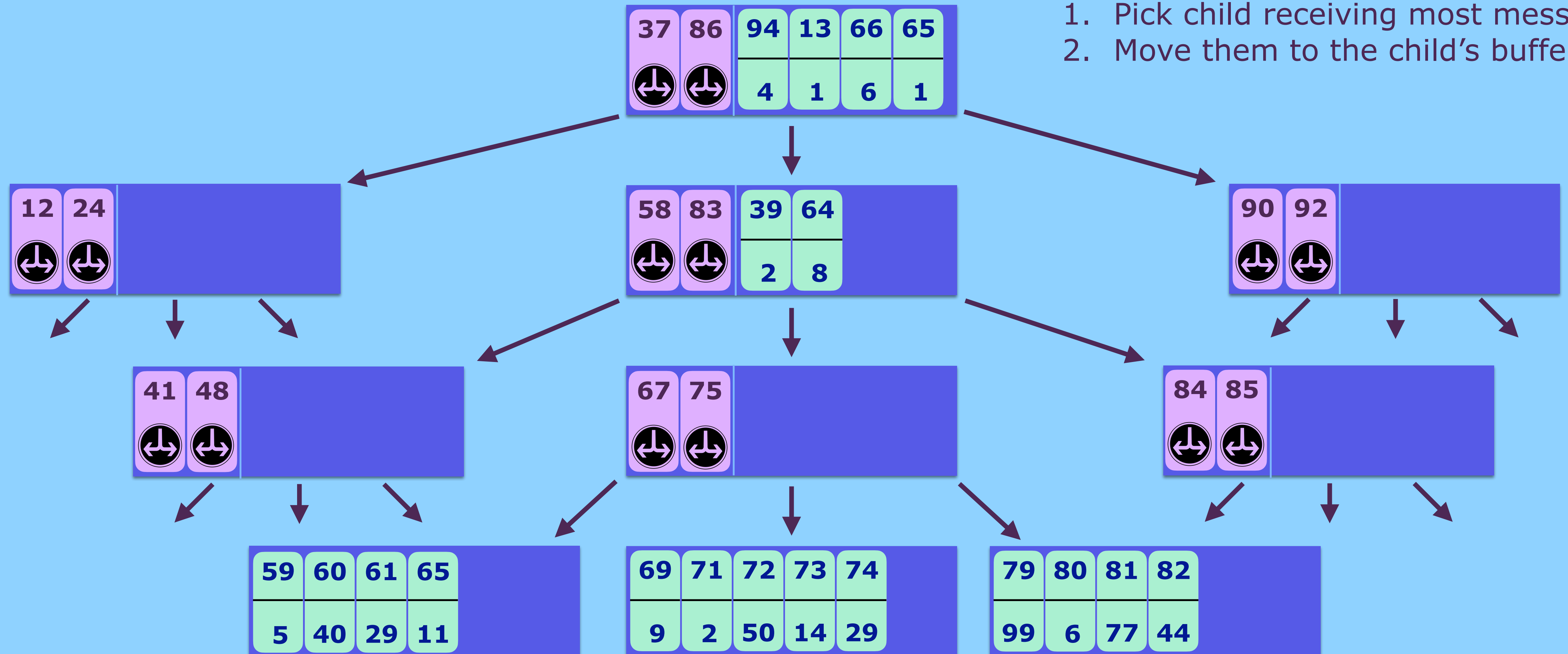
Inserts get put in the root buffer



B ϵ -Trees

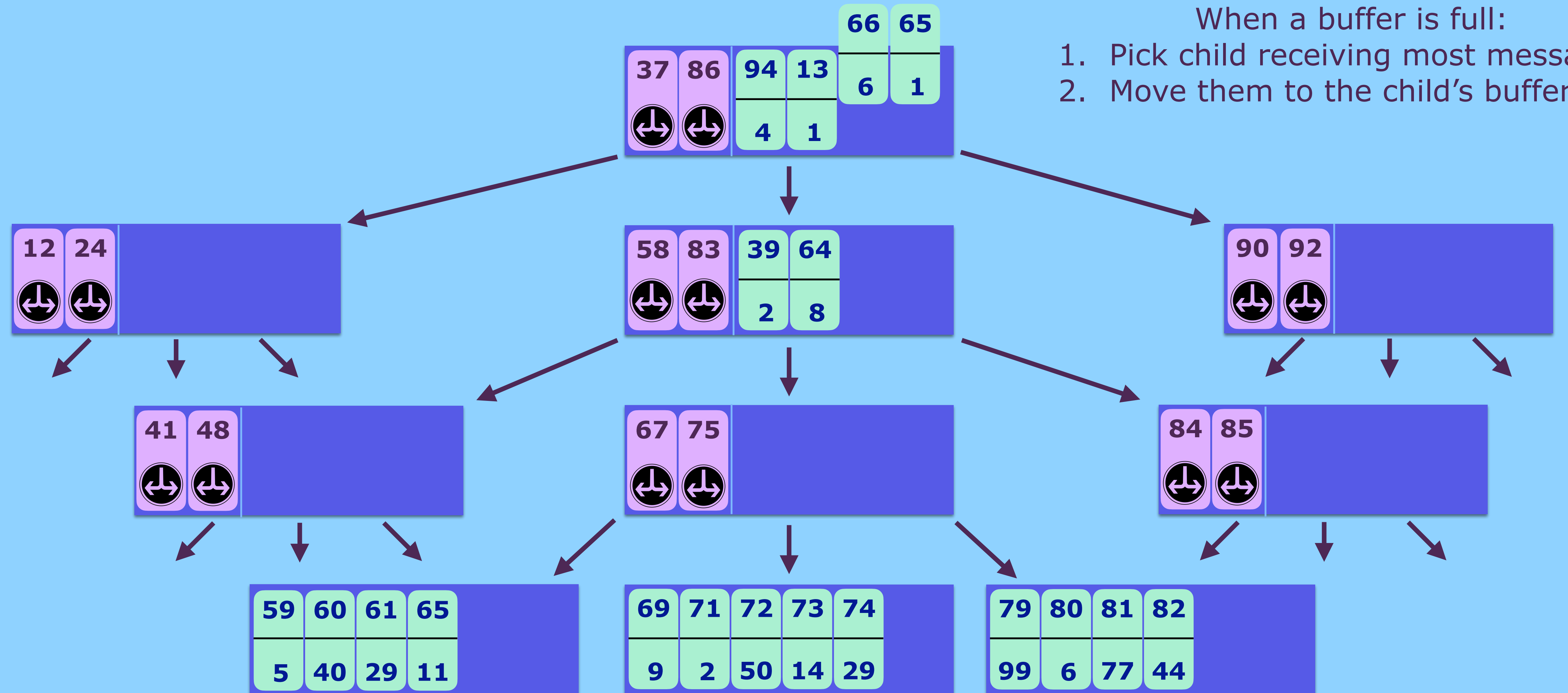
Inserts get put in the root buffer

- When a buffer is full:
1. Pick child receiving most messages
 2. Move them to the child's buffer



B ϵ -Trees

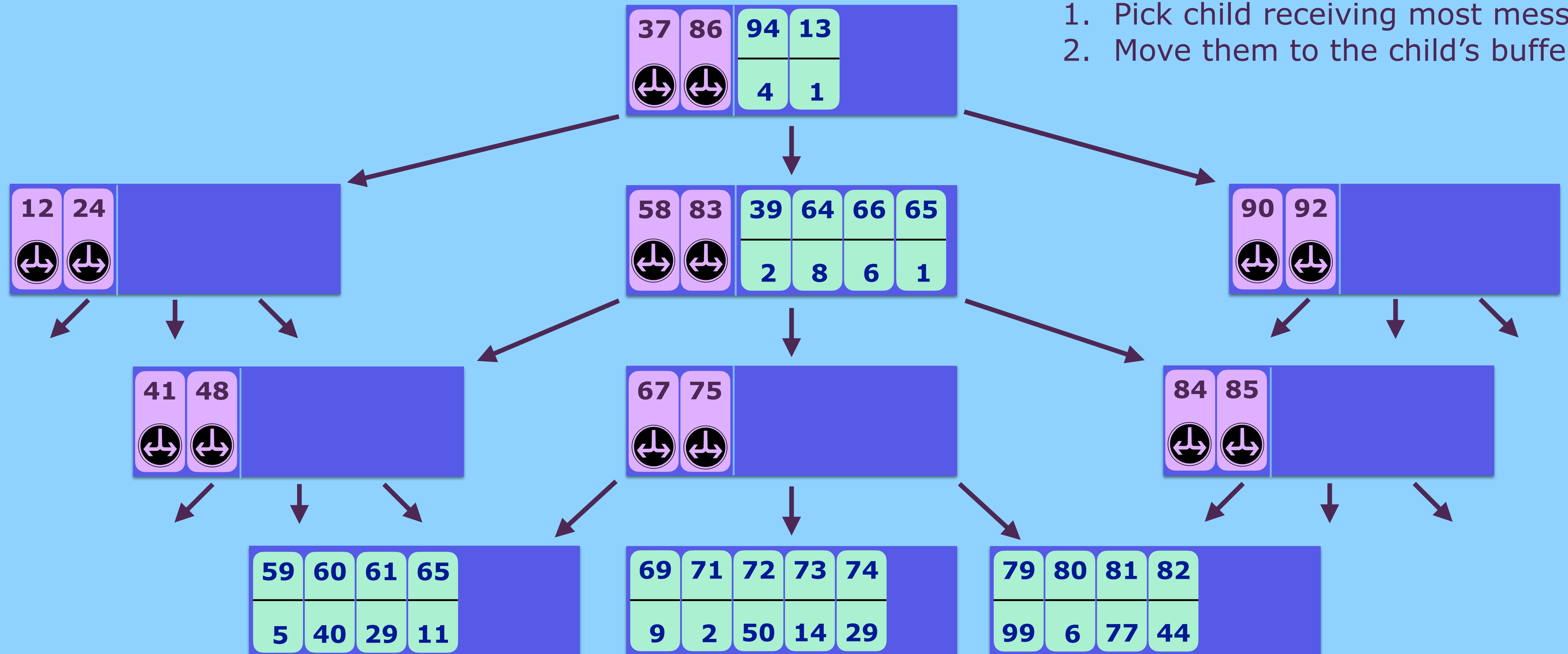
Inserts get put in the root buffer



B ϵ -Trees

Inserts get put in the root buffer

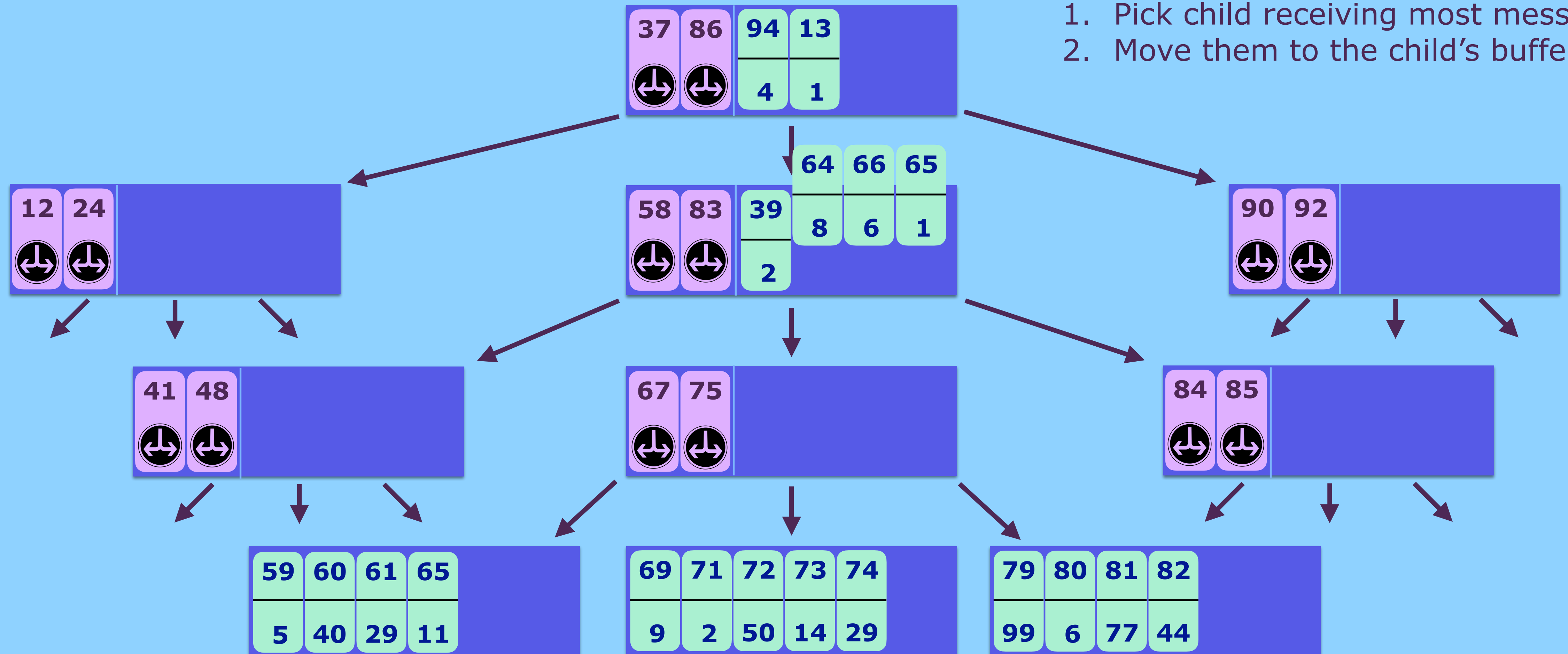
- When a buffer is full:
1. Pick child receiving most messages
 2. Move them to the child's buffer



B ϵ -Trees

Inserts get put in the root buffer

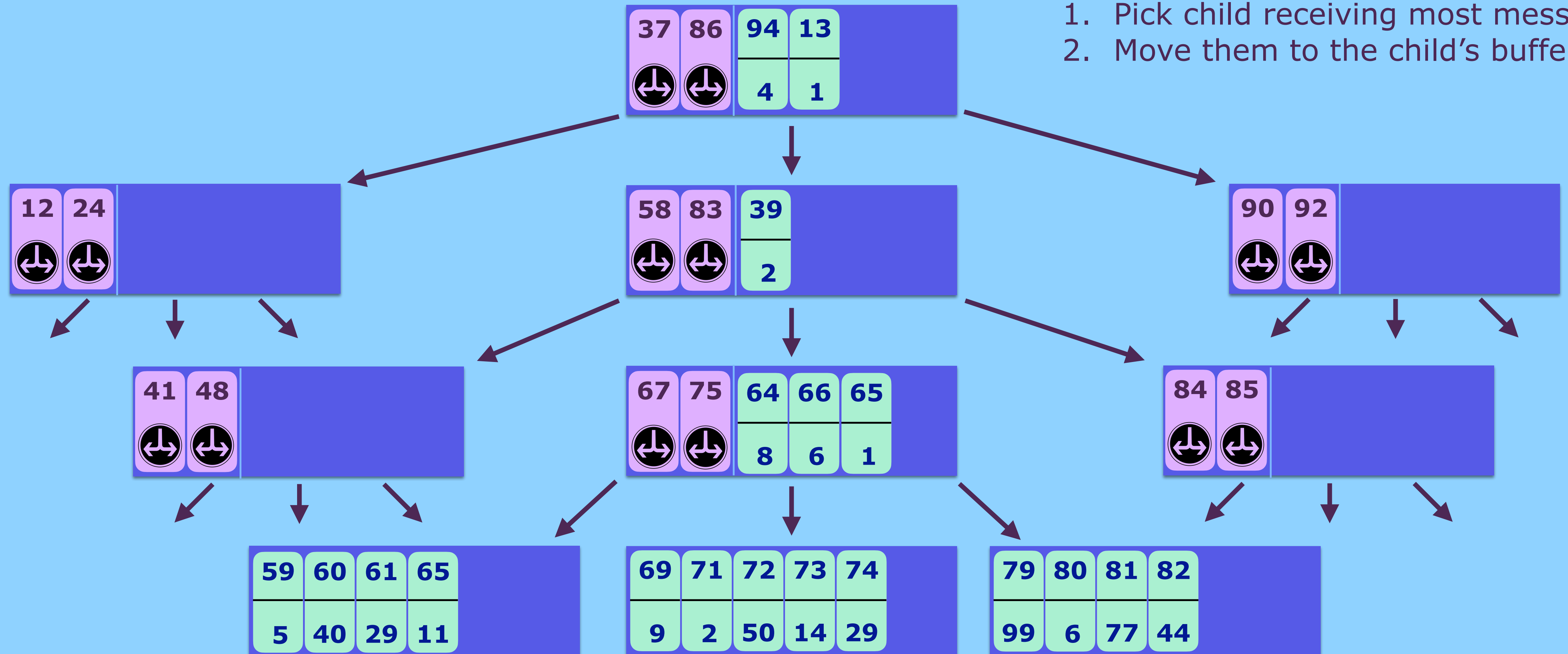
- When a buffer is full:
1. Pick child receiving most messages
 2. Move them to the child's buffer



B ϵ -Trees

Inserts get put in the root buffer

- When a buffer is full:
1. Pick child receiving most messages
 2. Move them to the child's buffer

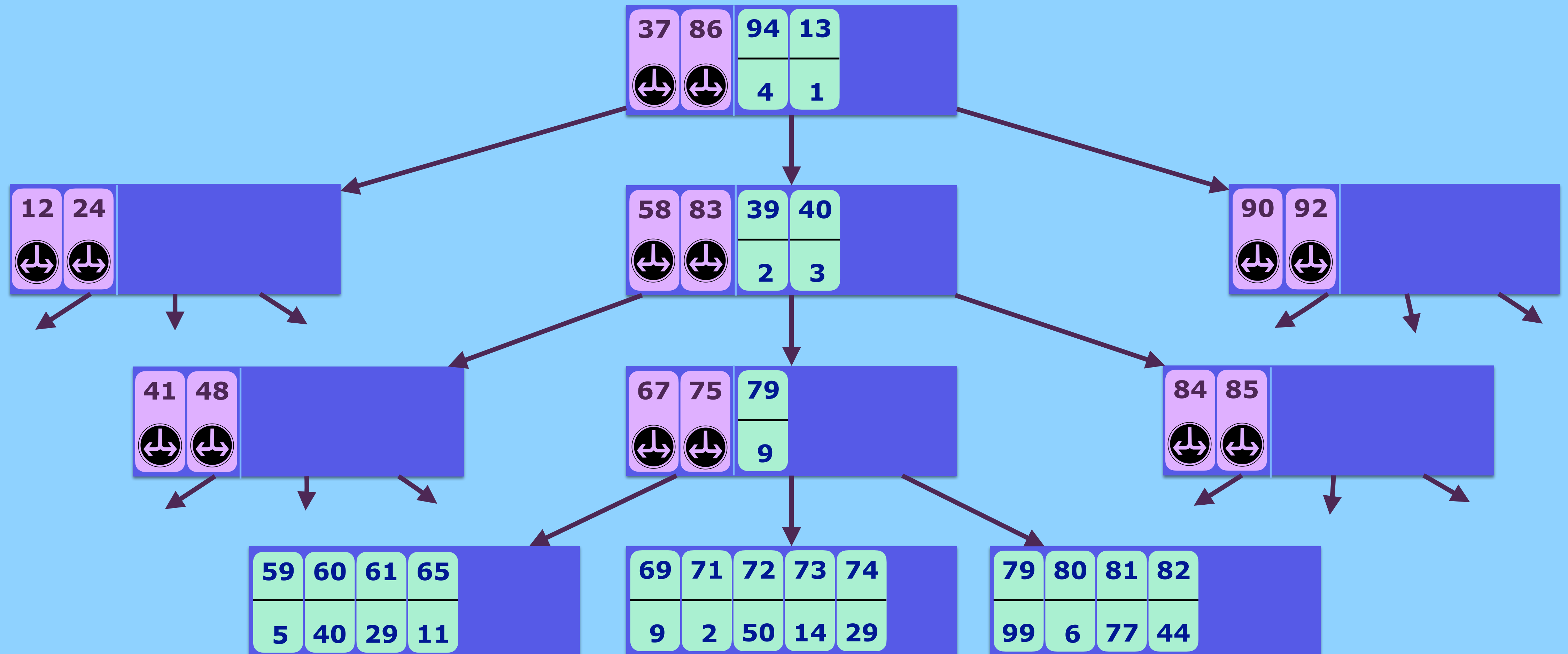


Lookups in B^ϵ -Trees

B ϵ -Trees

Lookups follow pivots, but check buffers along the way

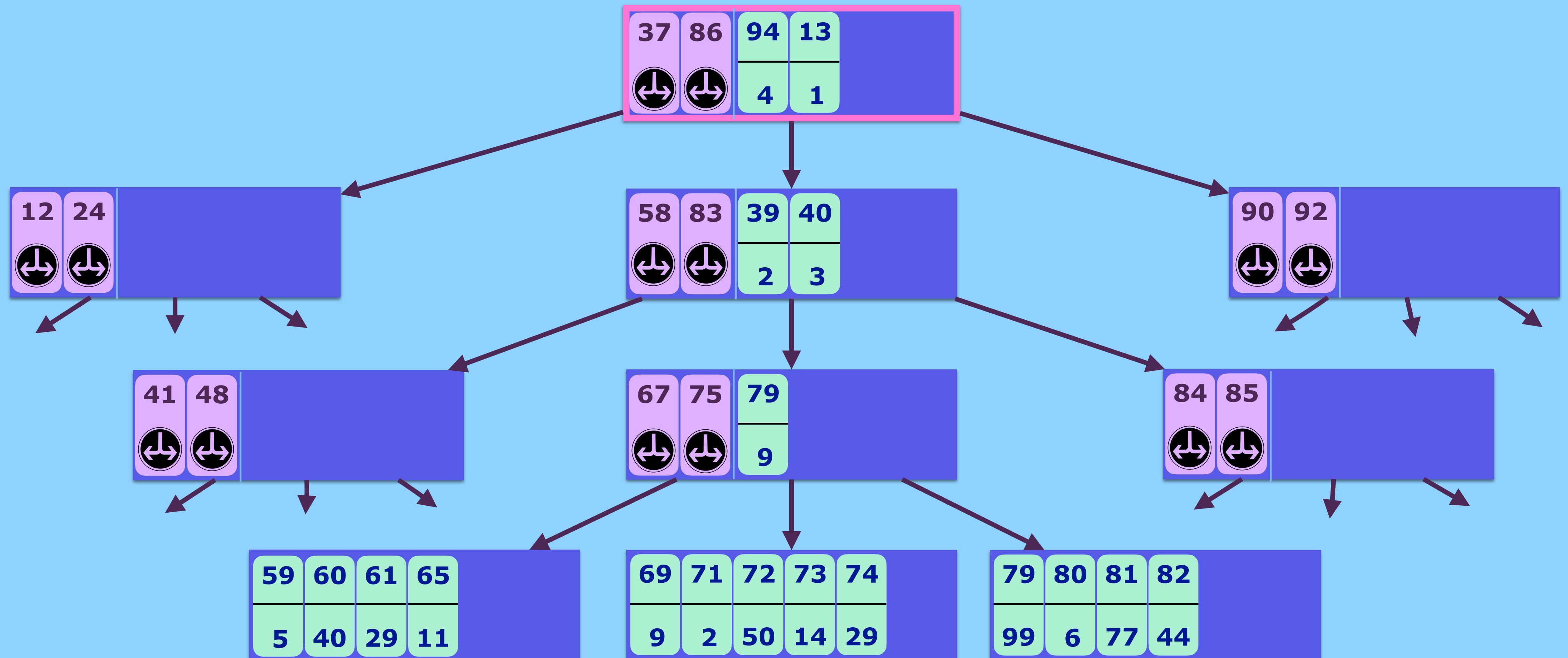
Query(71)



B ϵ -Trees

Lookups follow pivots, but check buffers along the way

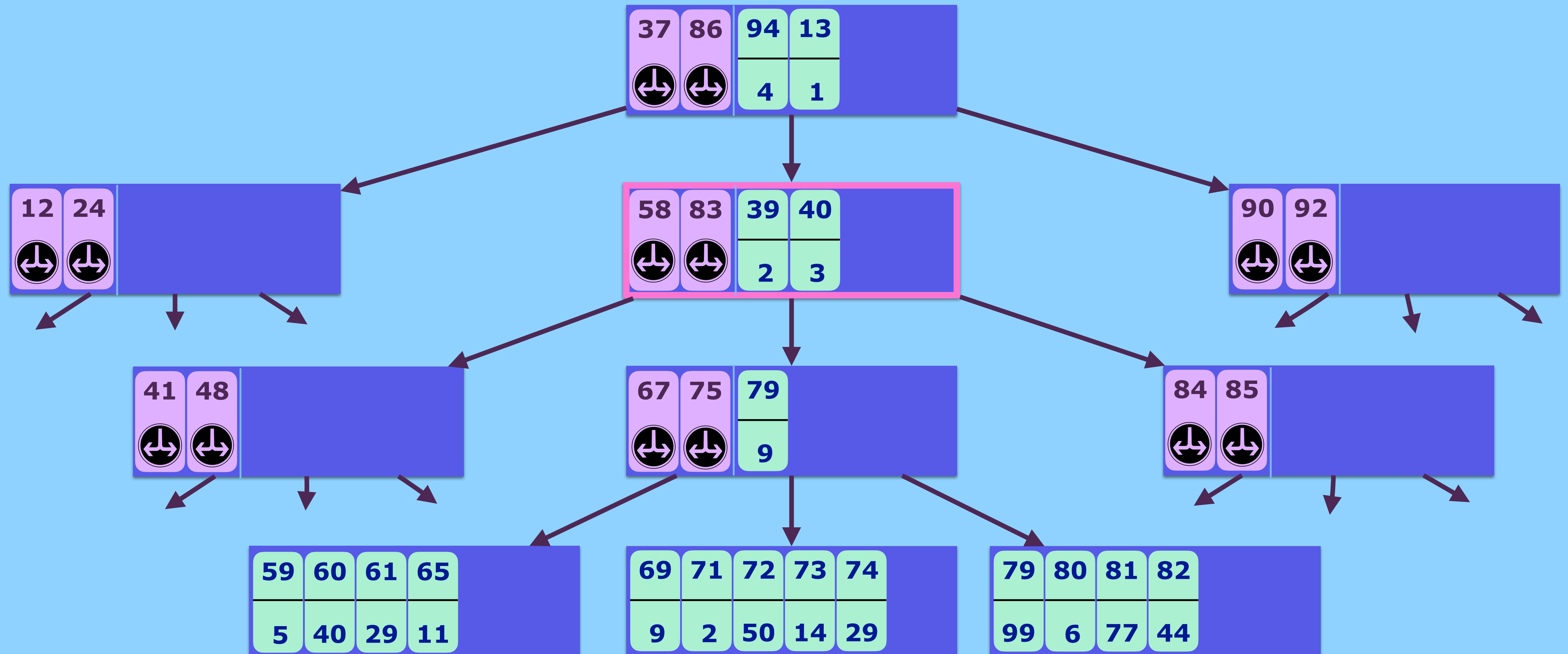
Query(71)



B ϵ -Trees

Lookups follow pivots, but check buffers along the way

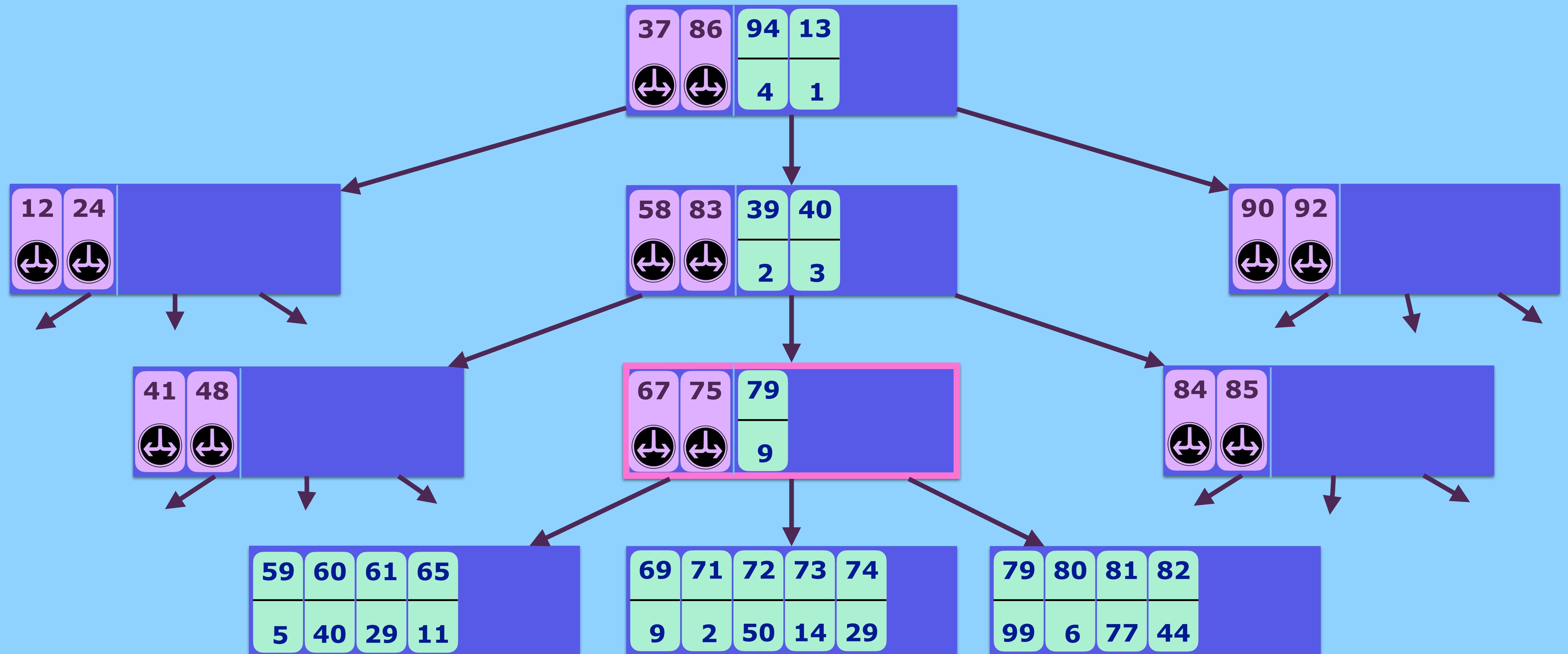
Query(71)



B ϵ -Trees

Lookups follow pivots, but check buffers along the way

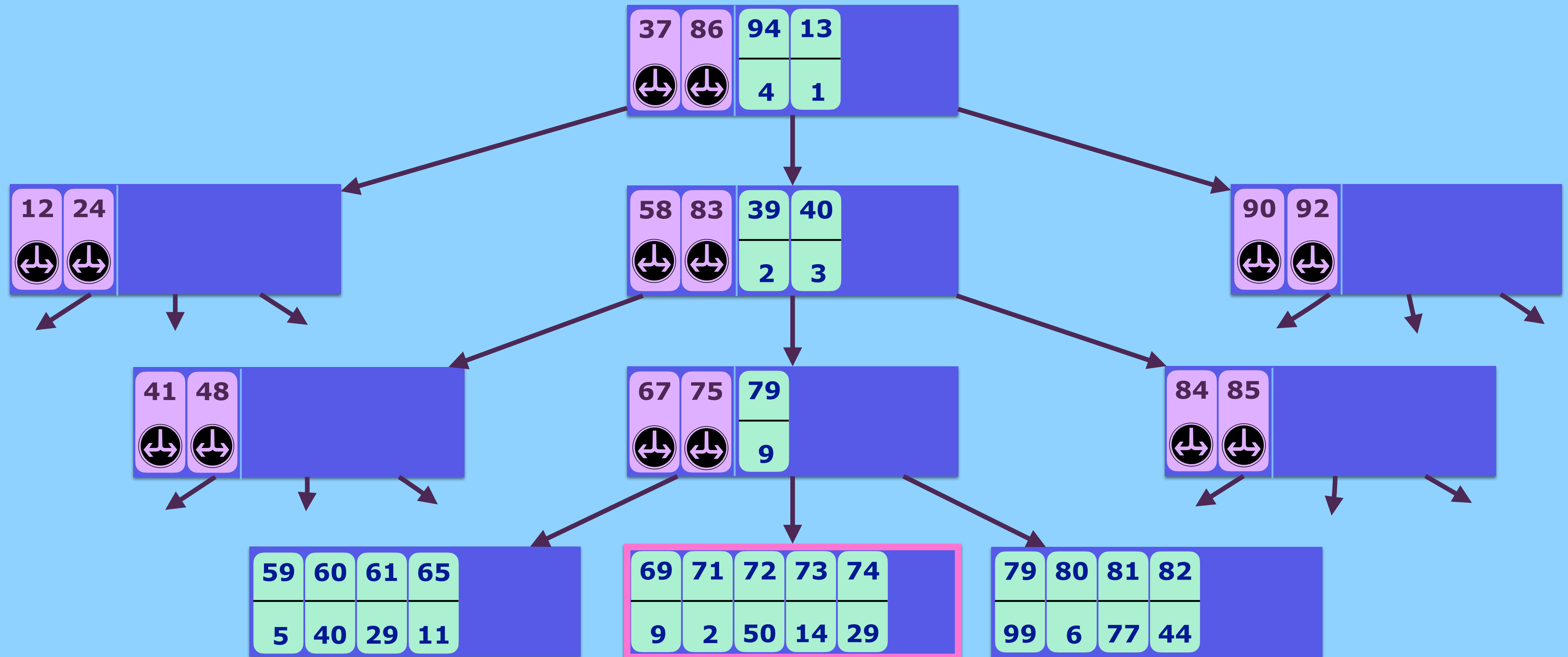
Query(71)



B ϵ -Trees

Lookups follow pivots, but check buffers along the way

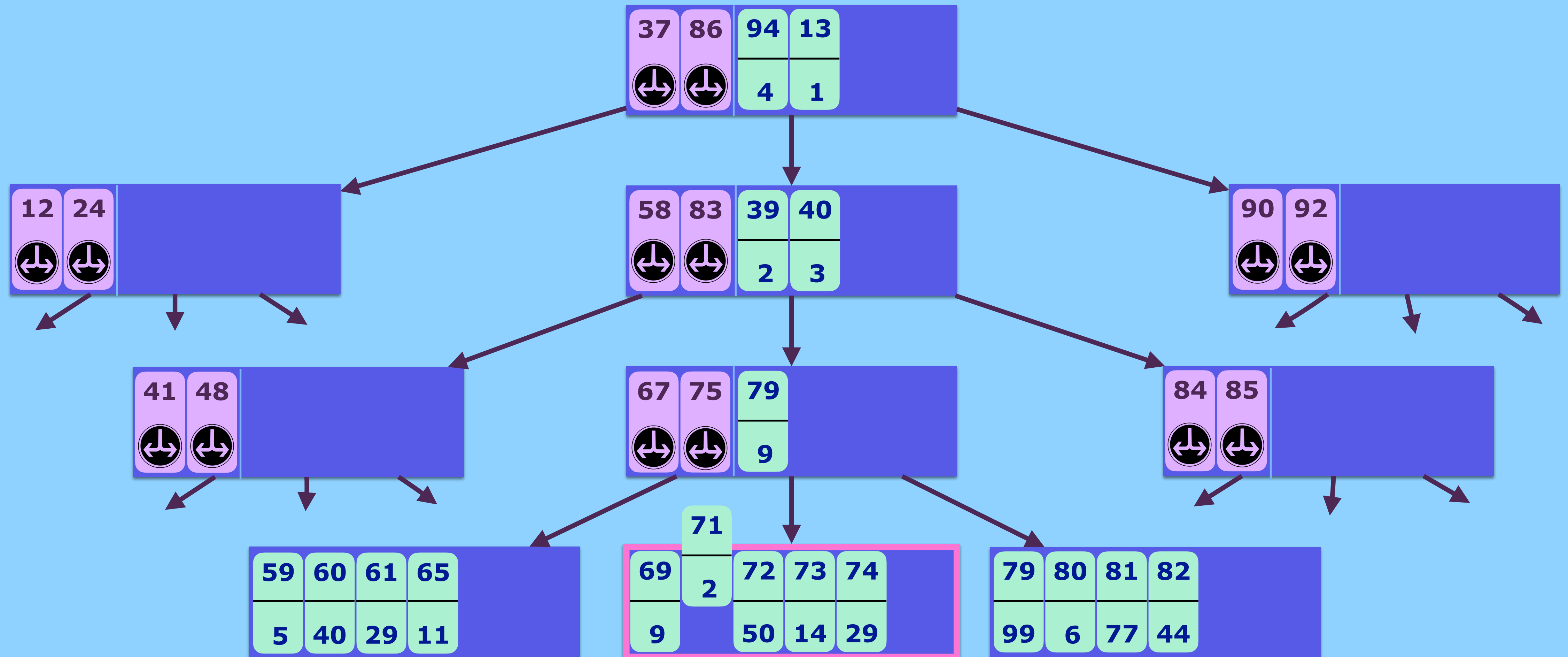
Query(71) \rightarrow 2



B ϵ -Trees

Lookups follow pivots, but check buffers along the way

Query(71) \rightarrow 2



Work
~~Write~~ Amplification
in B^ε -Trees

Work Amplification in B^ϵ -Trees

Recall: Insertions in B^ϵ -trees

65	72	80
11	50	6



To add new data to a B^ϵ -tree node, the node must be rewritten

58	83	39	64	66
		2	8	6

Work Amplification in B^ϵ -Trees

Recall: Insertions in B^ϵ -trees



To add new data to a B^ϵ -tree node, the node must be rewritten

58	83	39	64	66	65	72	80
		2	8	6	11	50	6

Work Amplification in B^ϵ -Trees

Recall: Insertions in B^ϵ -trees

To add new data to a B^ϵ -tree node, the node must be rewritten

58	83	39	64	66	65	72	80
		2	8	6	11	50	6



Therefore, any messages already in the node get written out again

Work Amplification in B^ϵ -Trees

Recall: Insertions in B^ϵ -trees

98	44
1	3

To add new data to a B^ϵ -tree node, the node must be rewritten

58	83	39	64	66	65	72	80
		2	8	6	11	50	6

Therefore, any messages already in the node get written out again



Work Amplification in B^ϵ -Trees

Recall: Insertions in B^ϵ -trees

To add new data to a B^ϵ -tree node, the node must be rewritten

Therefore, any messages already in the node get written out again

And again



58	83	39	64	66	65	72	80	98	44
		2	8	6	11	50	6	1	3

Work Amplification in B^ϵ -Trees

Recall: Insertions in B^ϵ -trees

28	91
24	43

To add new data to a B^ϵ -tree node, the node must be rewritten

58	83	39	64	66	65	72	80	98	44
		2	8	6	11	50	6	1	3

Therefore, any messages already in the node get written out again

And again

Work Amplification in B^ϵ -Trees



Recall: Insertions in B^ϵ -trees

To add new data to a B^ϵ -tree node, the node must be rewritten

Therefore, any messages already in the node get written out again

And again

And again

58	83	39	64	66	65	72	80	98	44	28	91
		2	8	6	11	50	6	1	3	24	43

Work Amplification in B^ϵ -Trees

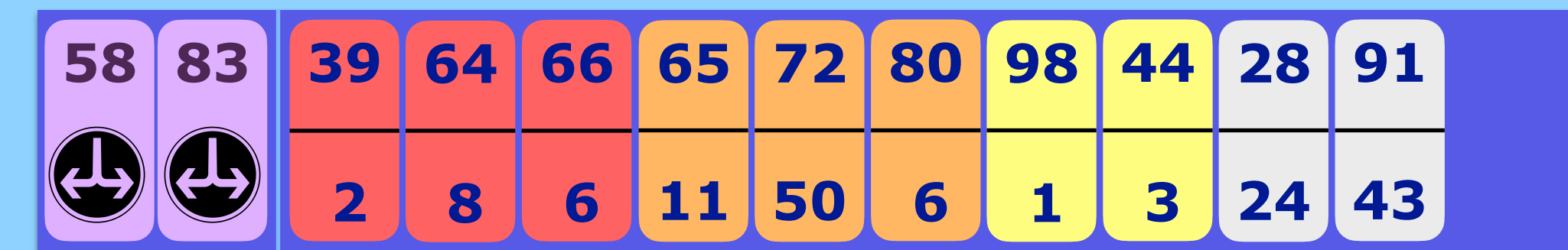
Recall: Insertions in B^ϵ -trees

To add new data to a B^ϵ -tree node, the node must be rewritten

Therefore, any messages already in the node get written out again

And again

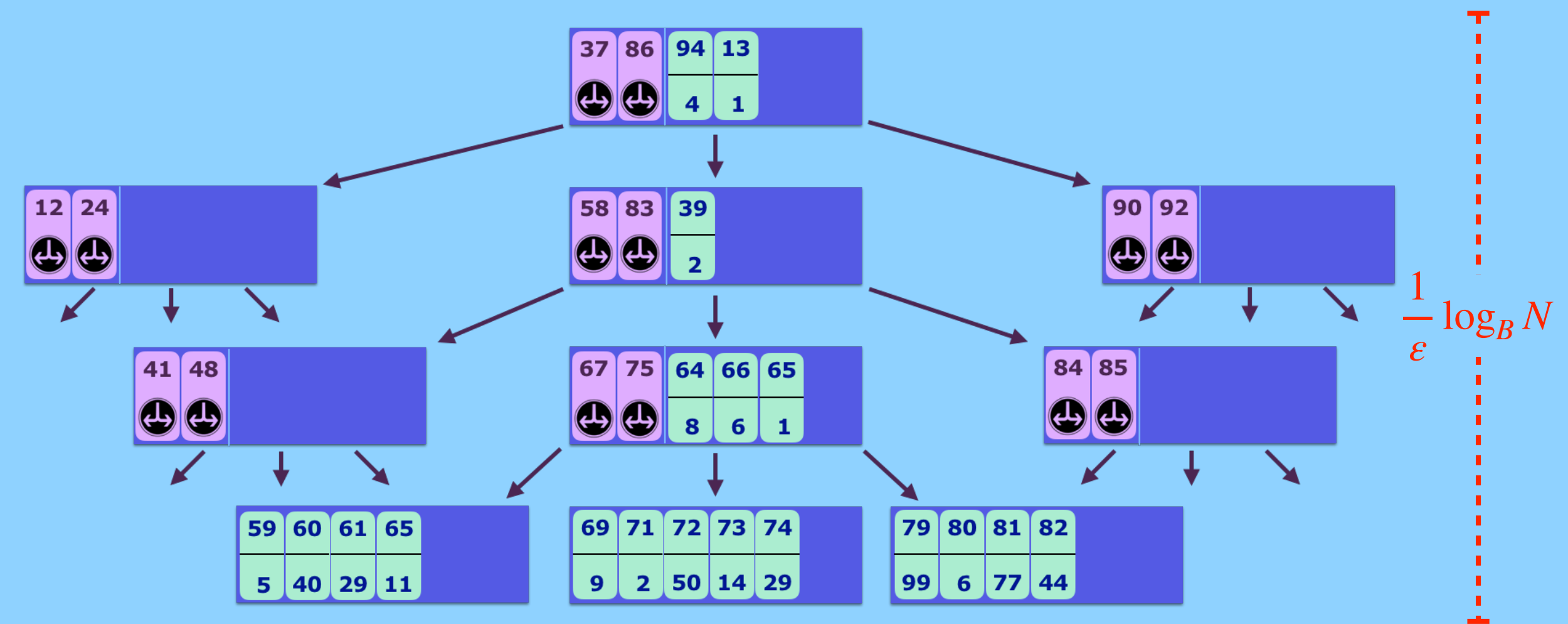
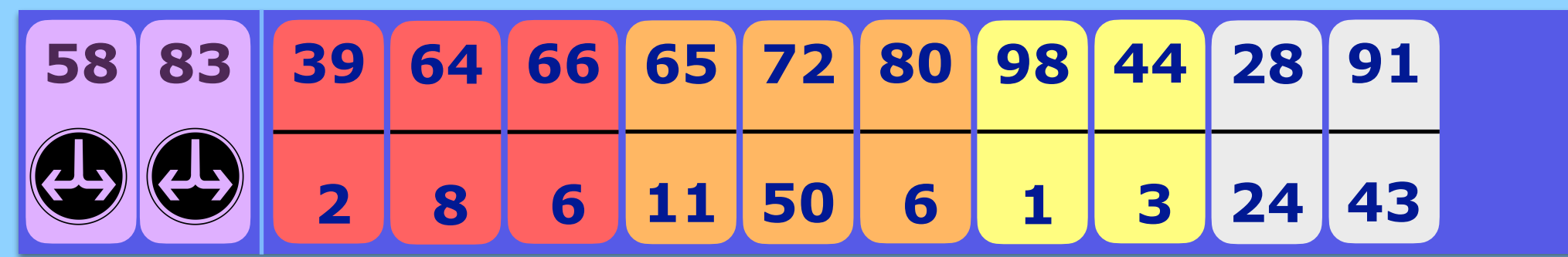
And again



In the worst case, the average message is rewritten $B^\epsilon/2$ times in each node

Work Amplification in B^ϵ -Trees

Recall: Insertions in B^ϵ -trees



In the worst case, the average message is rewritten $B^\epsilon/2$ times in each node

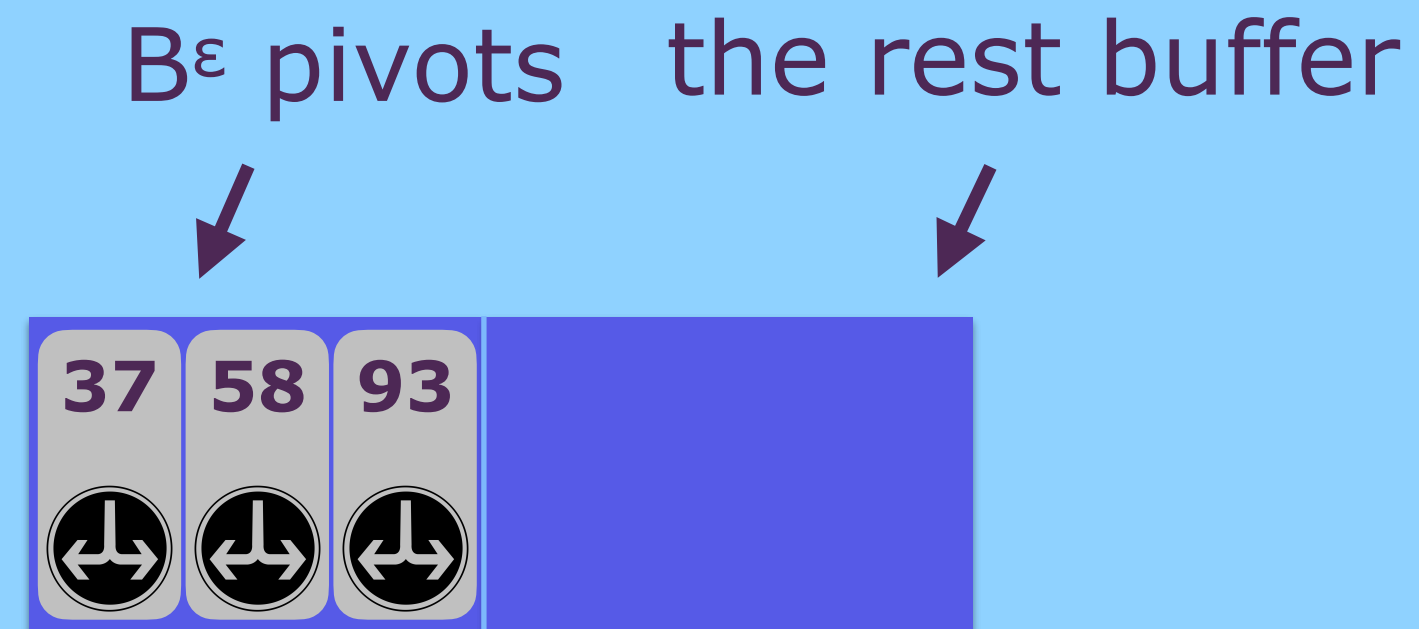
$$B^\epsilon\text{-Tree Work Amplification} = O(B^\epsilon \times \log_{B^\epsilon} N)$$

Size-Tiered B^ϵ -Trees (SplinterDB)

Size-Tiered B^ϵ -Trees

A Size-Tiered B^ϵ -tree is a B^ϵ -tree where the buffer is stored discontinuously

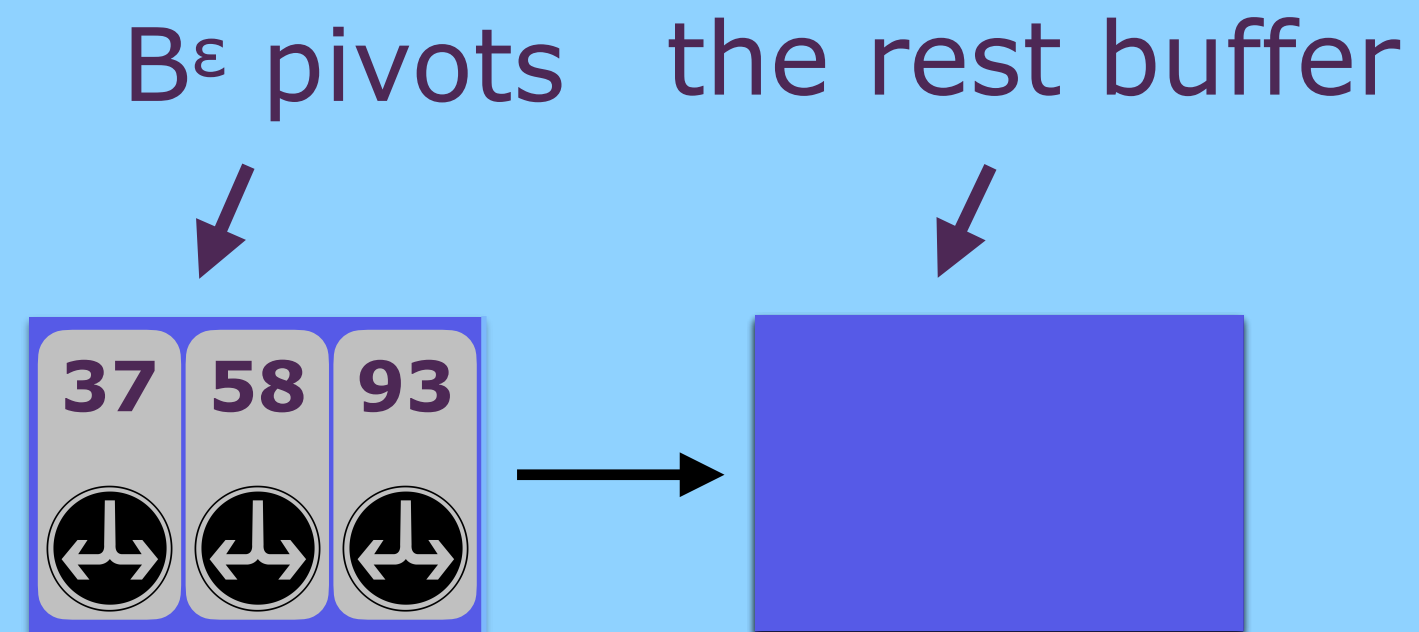
Recall:
a B^ϵ -tree node has
pivots and a buffer



Size-Tiered B^ϵ -Trees

A Size-Tiered B^ϵ -tree is a B^ϵ -tree where the buffer is stored discontinuously

Recall:
a B^ϵ -tree node has
pivots and a buffer

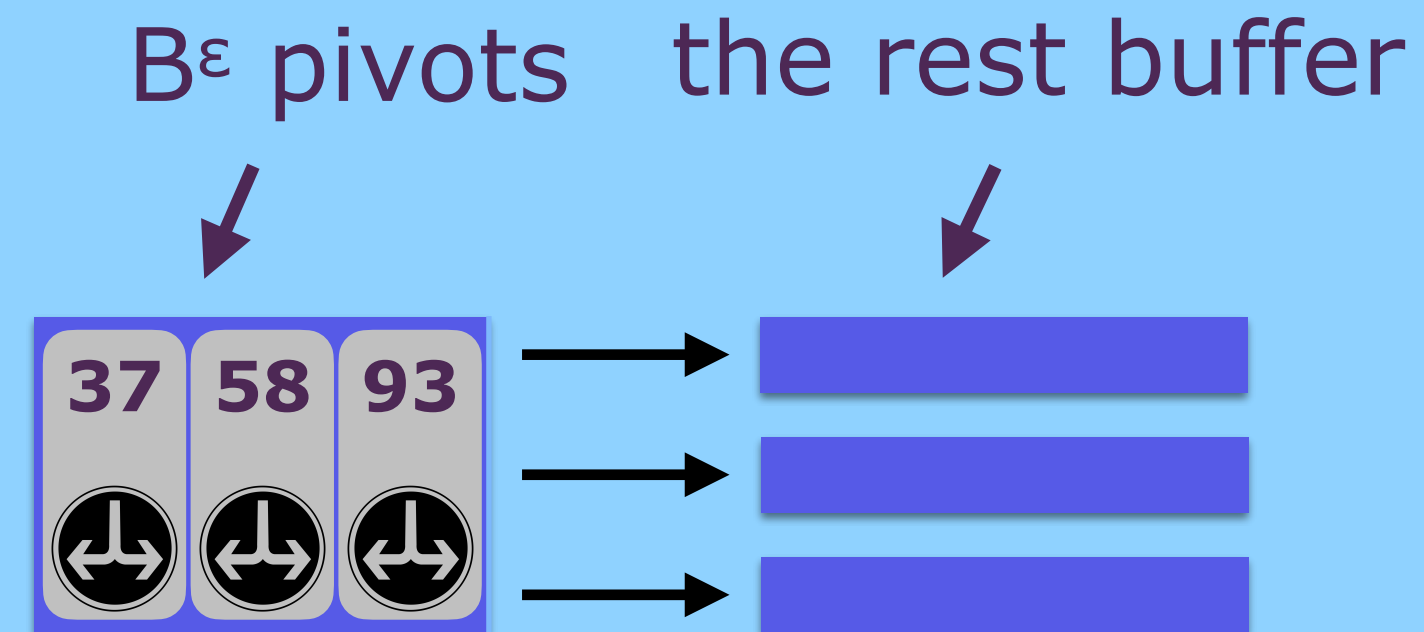


In an STB^ϵ -tree, the
buffer is stored
separately

Size-Tiered B^ϵ -Trees

A Size-Tiered B^ϵ -tree is a B^ϵ -tree where the buffer is stored discontinuously

Recall:
a B^ϵ -tree node has
pivots and a buffer



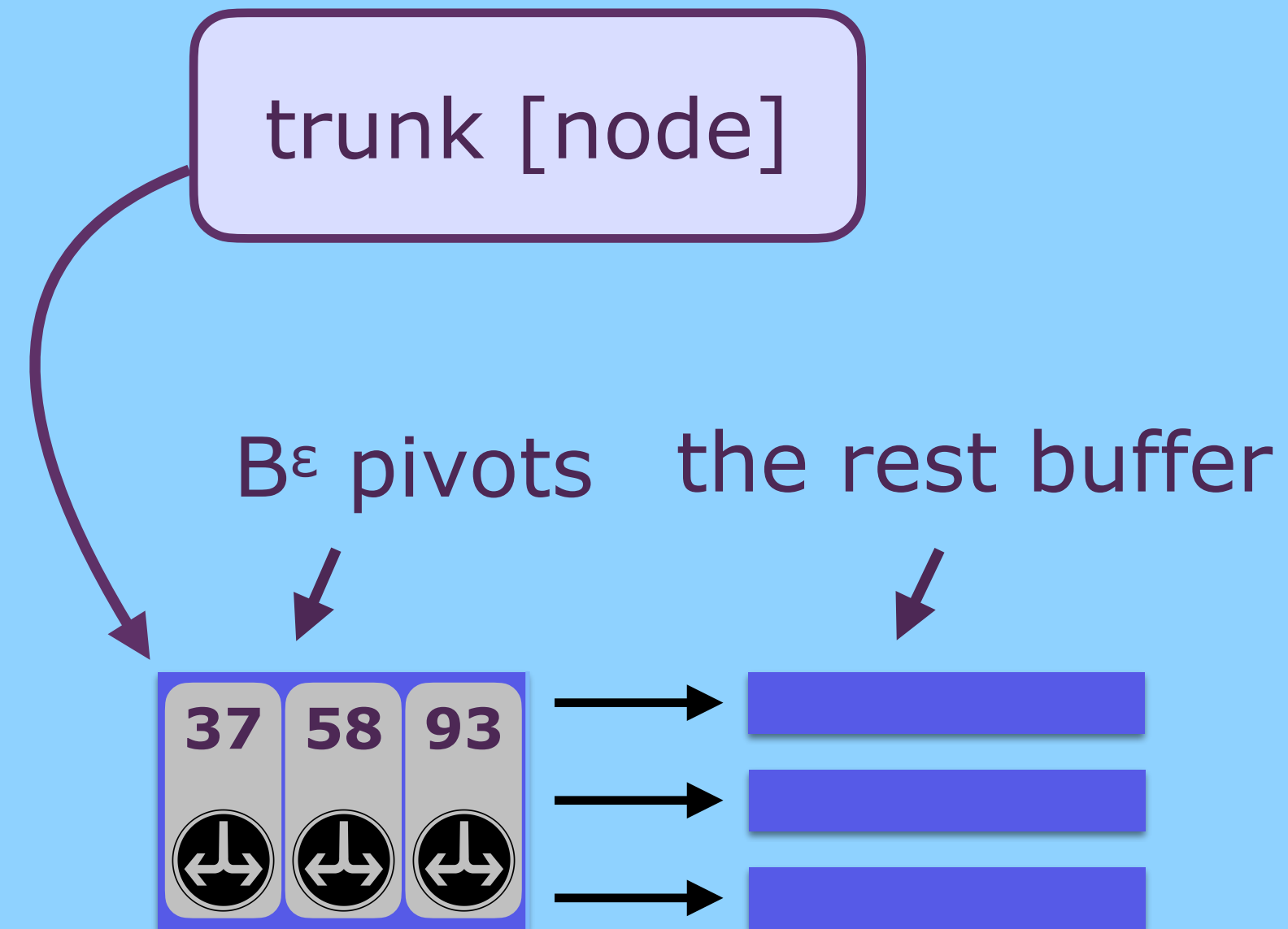
and in several
discontiguous pieces

In an STB^ϵ -tree, the
buffer is stored
separately

Size-Tiered B^ϵ -Trees

A Size-Tiered B^ϵ -tree is a B^ϵ -tree where the buffer is stored discontinuously

Recall:
a B^ϵ -tree node has pivots and a buffer



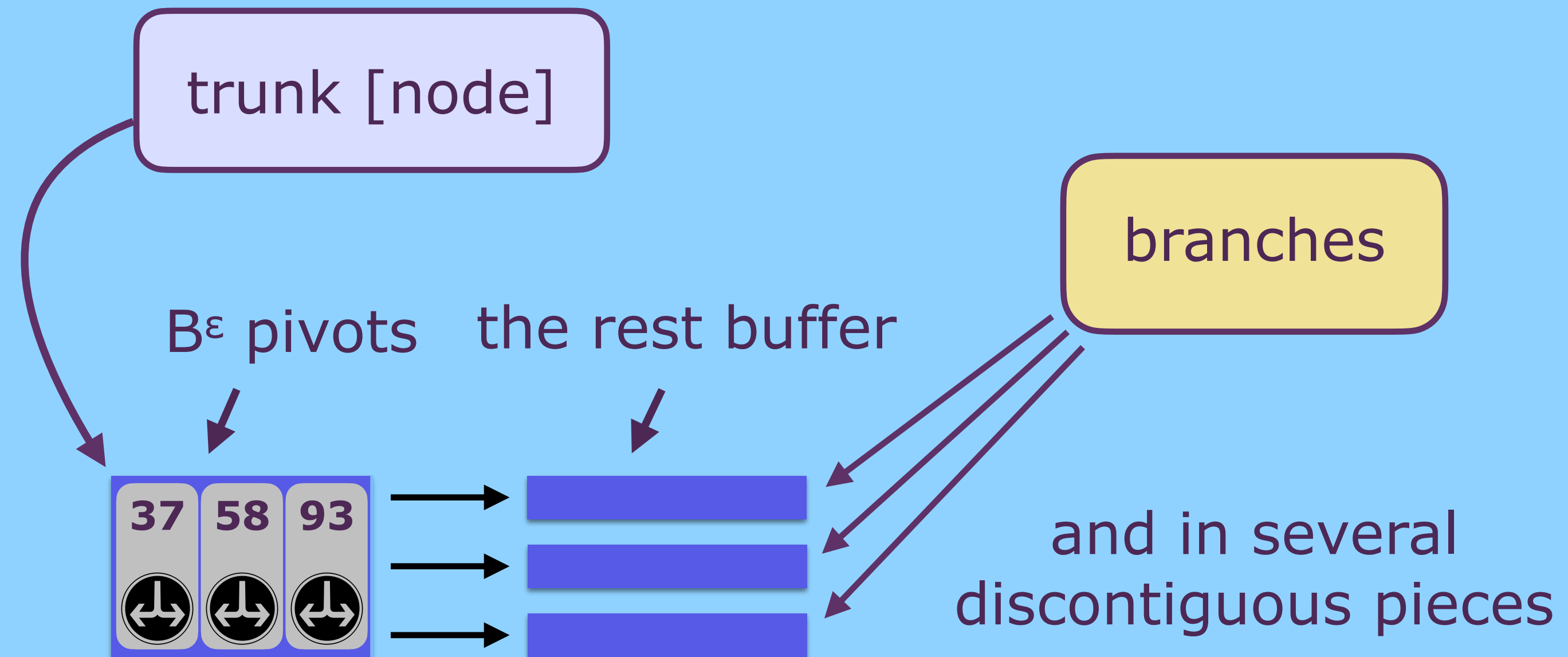
and in several
discontiguous pieces

In an STB^ϵ -tree, the
buffer is stored
separately

Size-Tiered B^ϵ -Trees

A Size-Tiered B^ϵ -tree is a B^ϵ -tree where the buffer is stored discontinuously

Recall:
a B^ϵ -tree node has pivots and a buffer



In an STB^ϵ -tree, the buffer is stored separately

Insertions in Size-Tiered B^ε -Trees

Size-Tiered B^ϵ -Trees

A Size-Tiered B^ϵ -tree is a B^ϵ -tree where the buffer is stored discontinuously

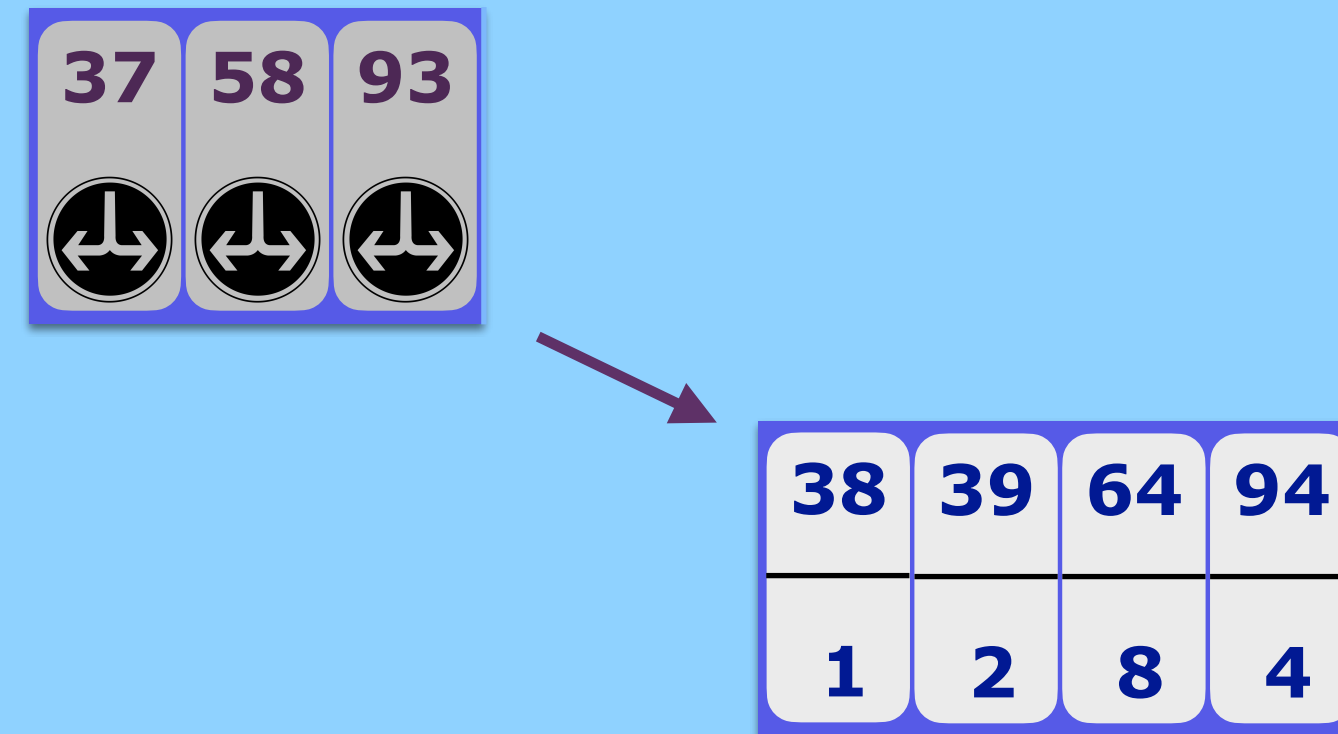
38	39	64	94
1	2	8	4

37	58	93
		

When new data is flushed into the trunk node...

Size-Tiered B^ϵ -Trees

A Size-Tiered B^ϵ -tree is a B^ϵ -tree where the buffer is stored discontinuously

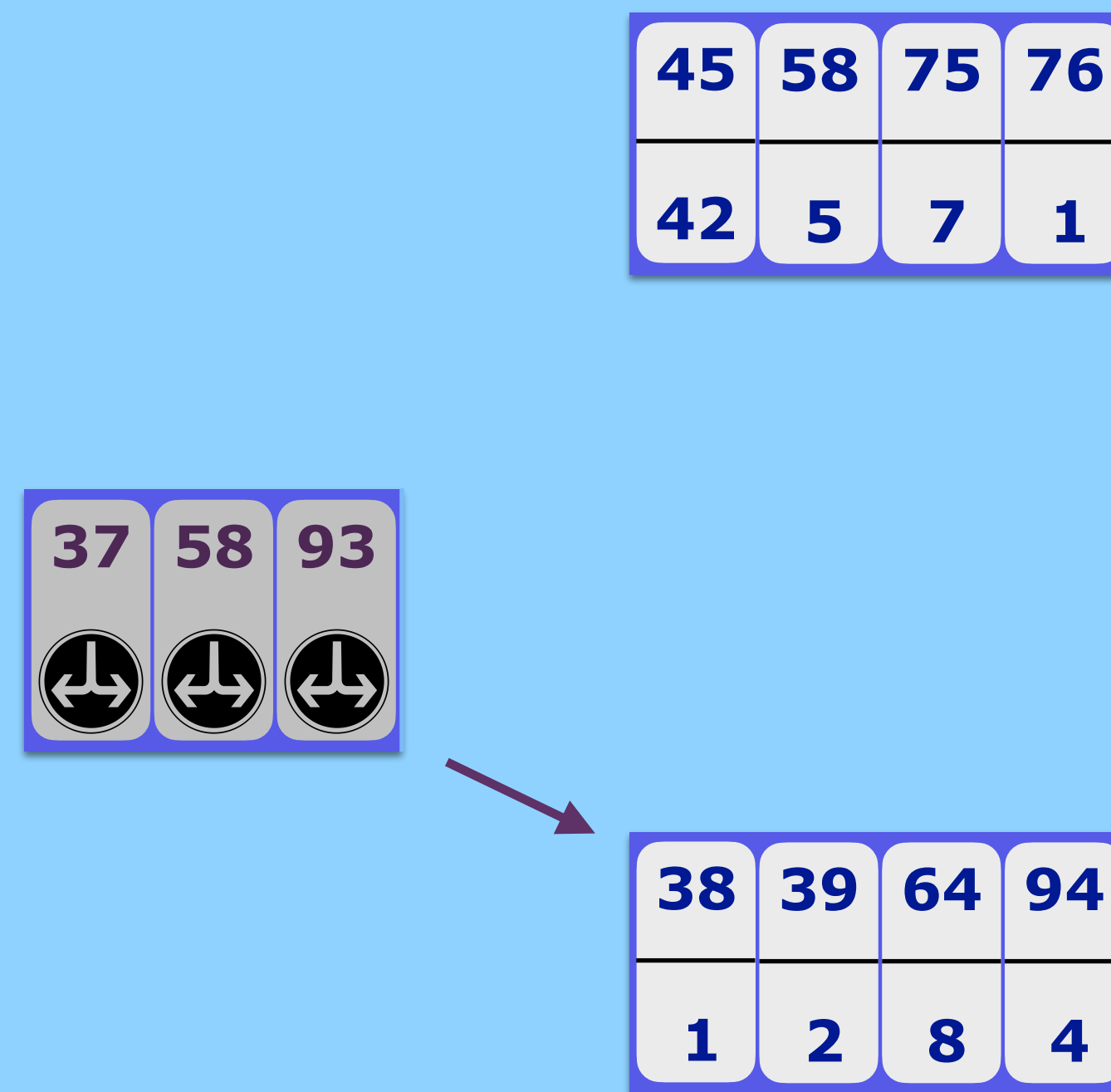


When new data is flushed into the trunk node...

...it is added as a new branch

Size-Tiered B^ϵ -Trees

A Size-Tiered B^ϵ -tree is a B^ϵ -tree where the buffer is stored discontinuously

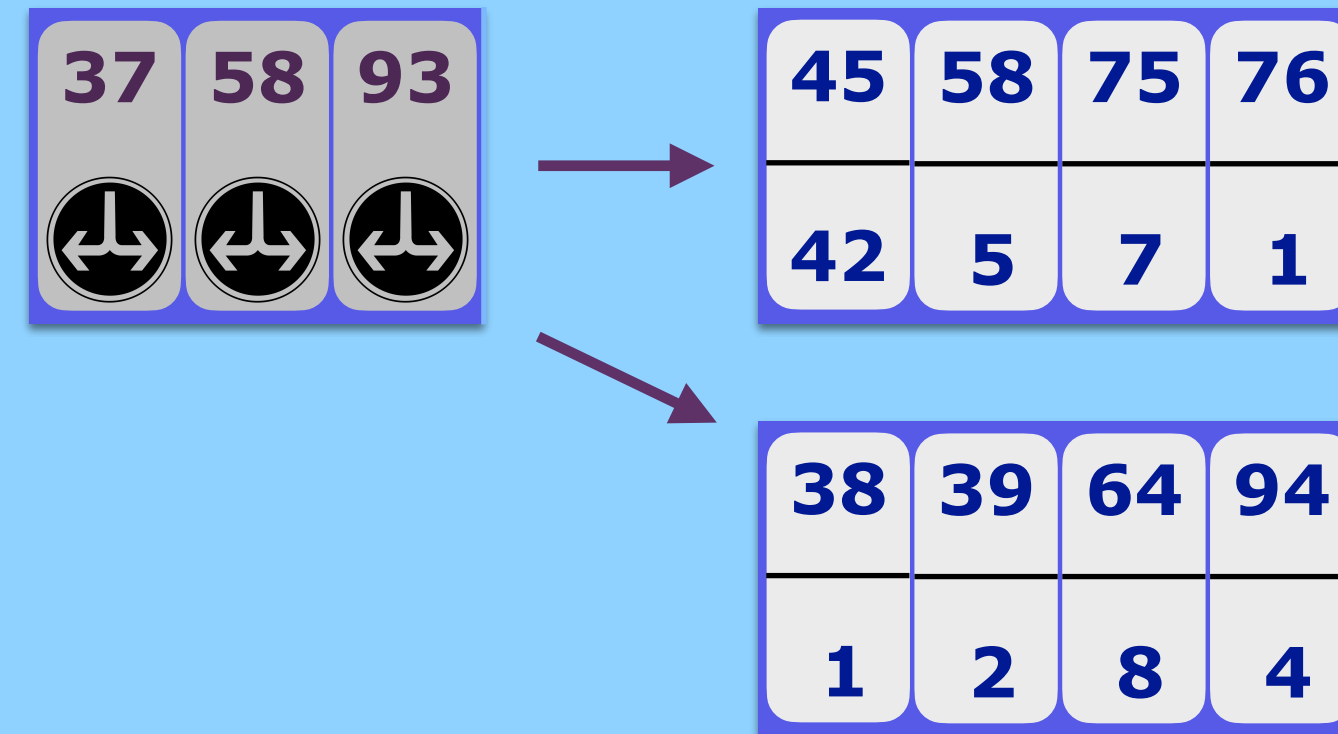


When new data is flushed into the trunk node...

...it is added as a new branch

Size-Tiered B^ϵ -Trees

A Size-Tiered B^ϵ -tree is a B^ϵ -tree where the buffer is stored discontinuously



When new data is flushed into the trunk node...

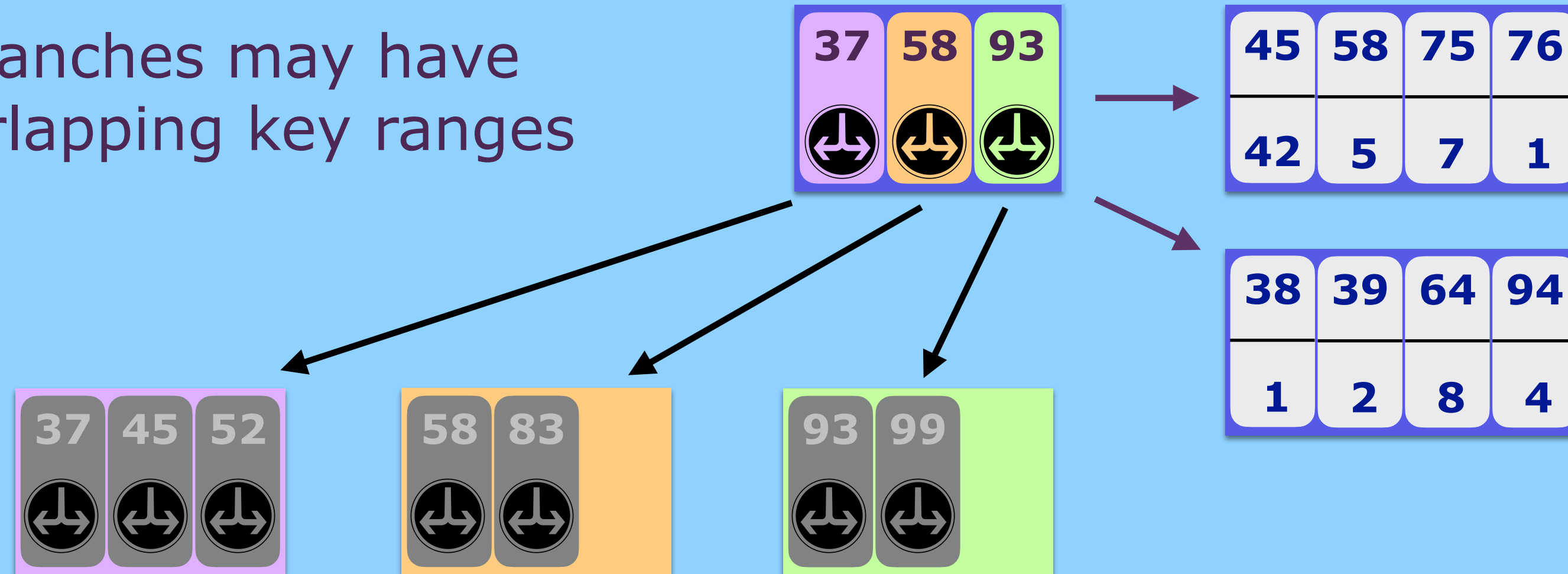
...it is added as a new branch

The old branches do not need to be rewritten

Size-Tiered B^ϵ -Trees

A Size-Tiered B^ϵ -tree is a B^ϵ -tree where the buffer is stored discontinuously

Branches may have overlapping key ranges



When new data is flushed into the trunk node...

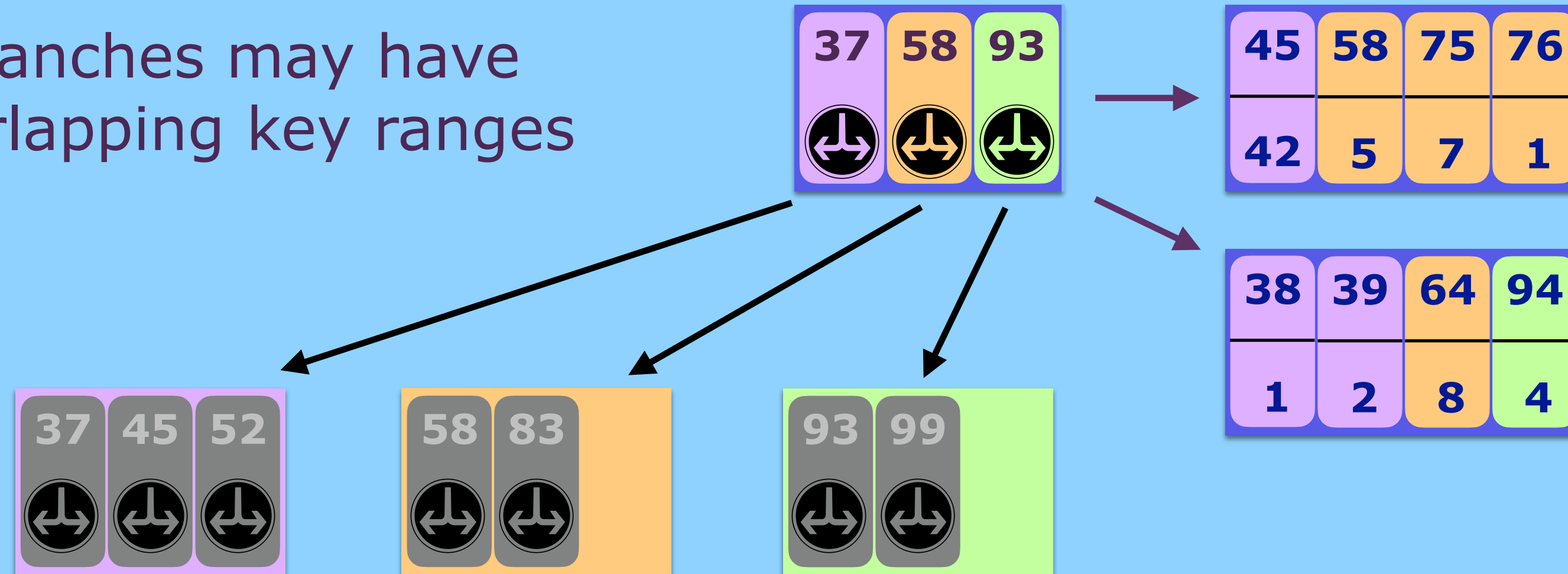
...it is added as a new branch

The old branches do not need to be rewritten

Size-Tiered B ϵ -Trees

A Size-Tiered B ϵ -tree is a B ϵ -tree where the buffer is stored discontinuously

Branches may have overlapping key ranges



When new data is flushed into the trunk node...

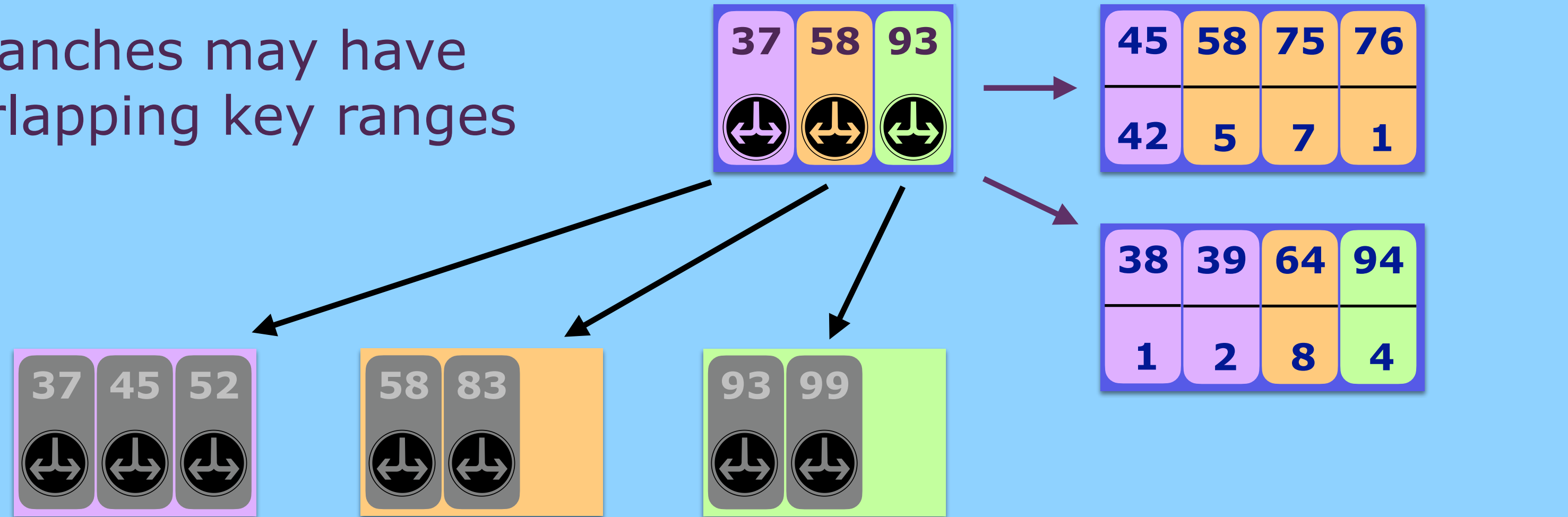
...it is added as a new branch

The old branches do not need to be rewritten

Size-Tiered B^ϵ -Trees

A Size-Tiered B^ϵ -tree is a B^ϵ -tree where the buffer is stored discontinuously

Branches may have overlapping key ranges



When new data is flushed into the trunk node...

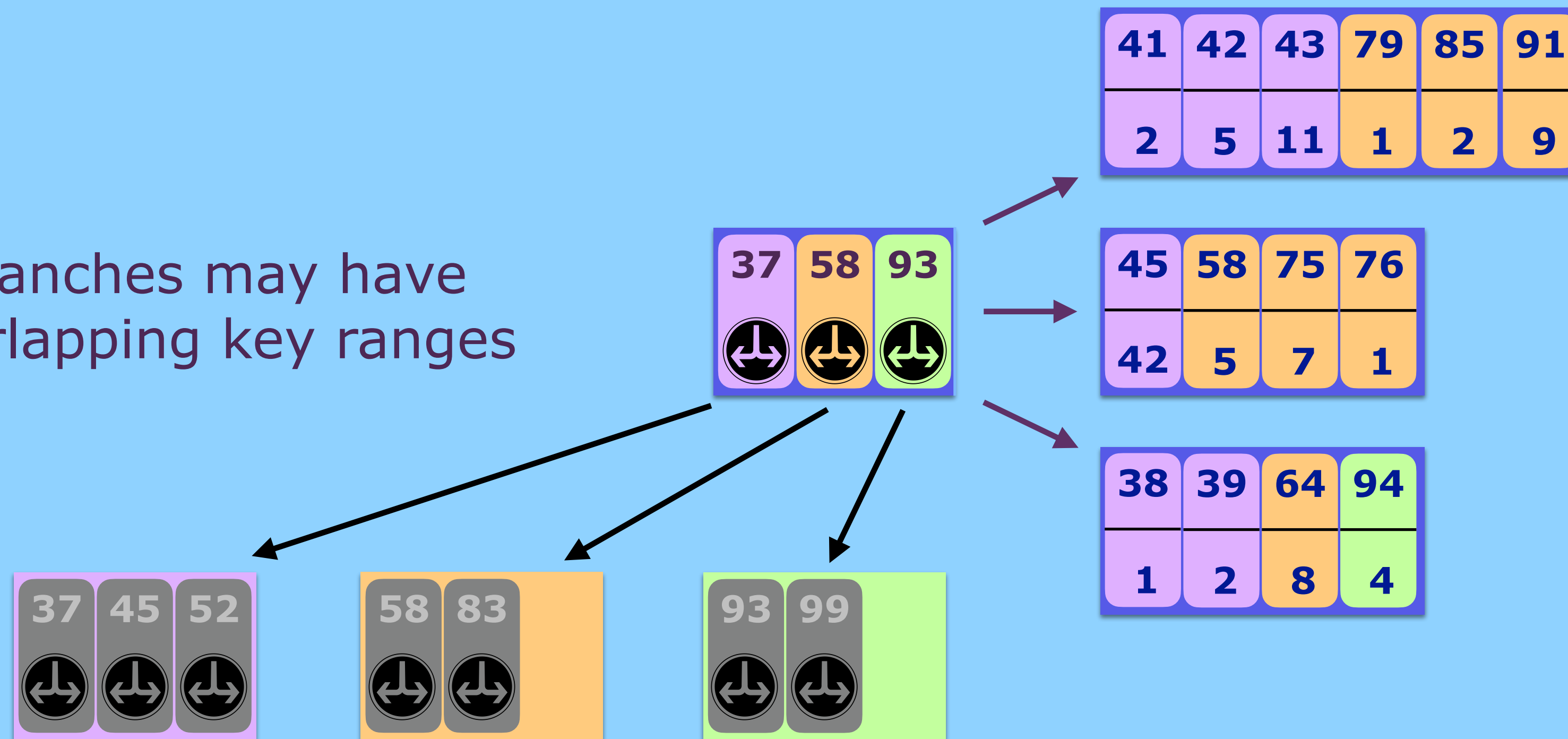
...it is added as a new branch

The old branches do not need to be rewritten

Size-Tiered B^ϵ -Trees

A Size-Tiered B^ϵ -tree is a B^ϵ -tree where the buffer is stored discontinuously

Branches may have overlapping key ranges



When new data is flushed into the trunk node...

...it is added as a new branch

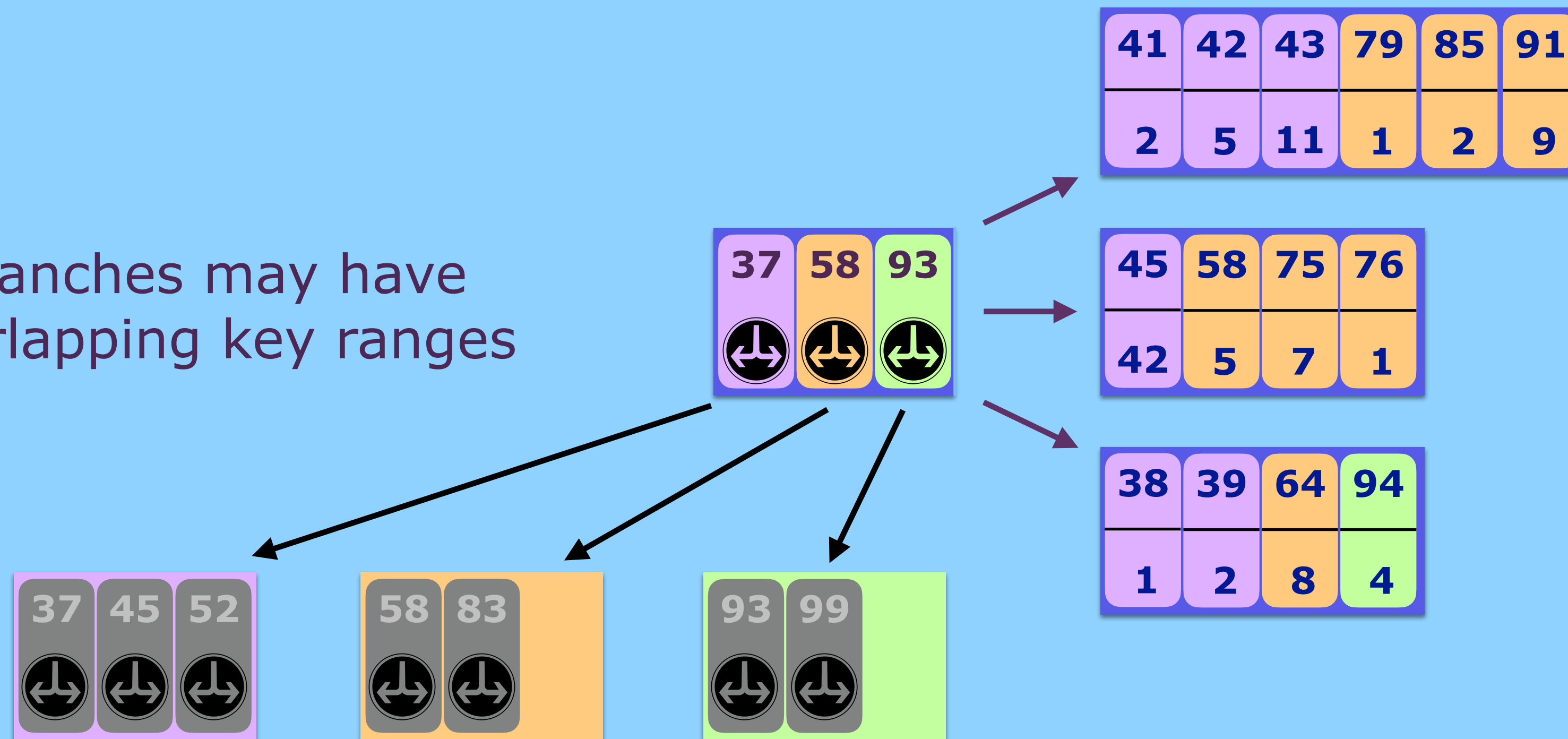
The old branches do not need to be rewritten

Size-Tiered B ϵ -Trees

A Size-Tiered B ϵ -tree is a B ϵ -tree where the buffer is stored discontinuously

The fullness threshold is:
Fanout \times Average Buffer Size

Branches may have overlapping key ranges



When new data is flushed into the trunk node...

...it is added as a new branch

The old branches do not need to be rewritten

Size-Tiered B ϵ -Trees

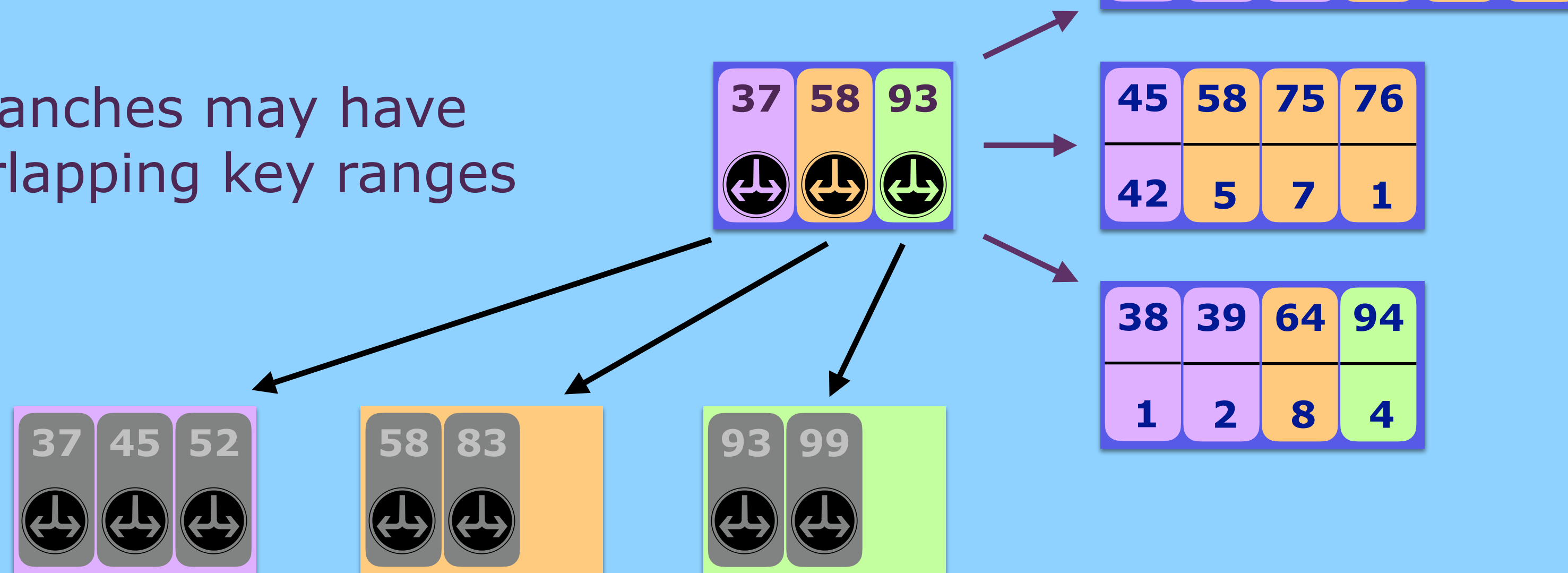
A Size-Tiered B ϵ -tree is a B ϵ -tree where the buffer is stored discontinuously

The fullness threshold is:
Fanout \times Average Buffer Size

When the node is full:

1. Pick child receiving most messages
2. Merge them into a new branch for the child

Branches may have overlapping key ranges



When new data is flushed into the trunk node...

...it is added as a new branch

The old branches do not need to be rewritten

Size-Tiered B ϵ -Trees

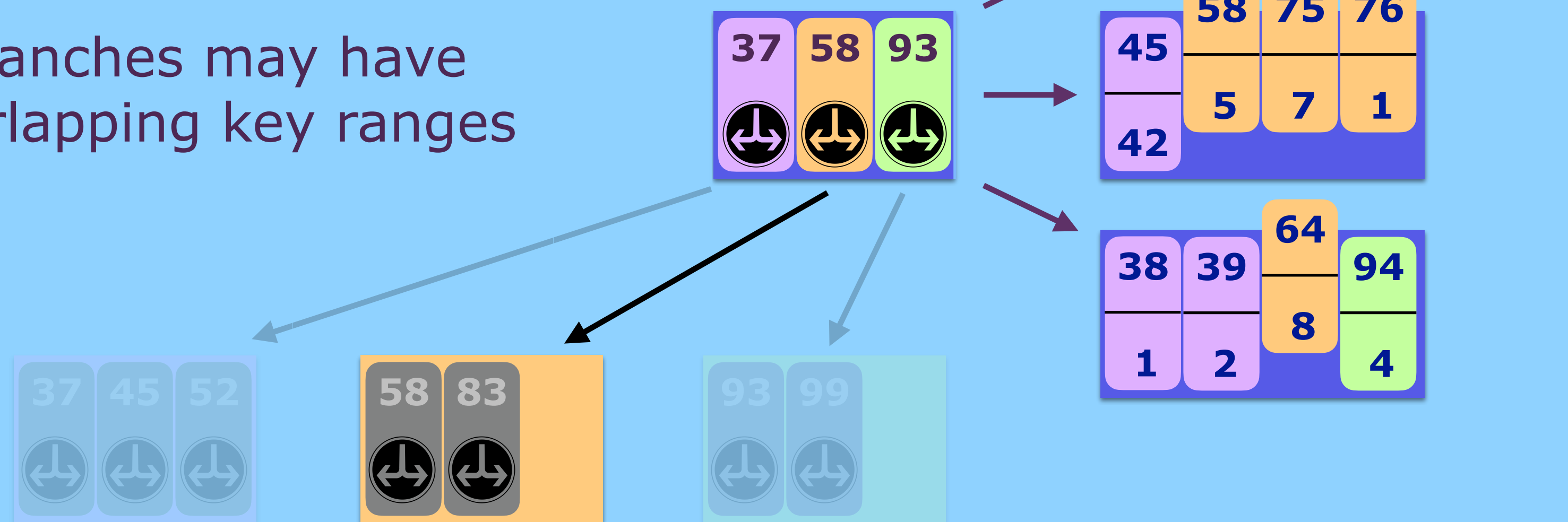
A Size-Tiered B ϵ -tree is a B ϵ -tree where the buffer is stored discontinuously

The fullness threshold is:
Fanout \times Average Buffer Size

When the node is full:

1. Pick child receiving most messages
2. Merge them into a new branch for the child

Branches may have overlapping key ranges



When new data is flushed into the trunk node...

...it is added as a new branch

The old branches do not need to be rewritten

Size-Tiered B ϵ -Trees

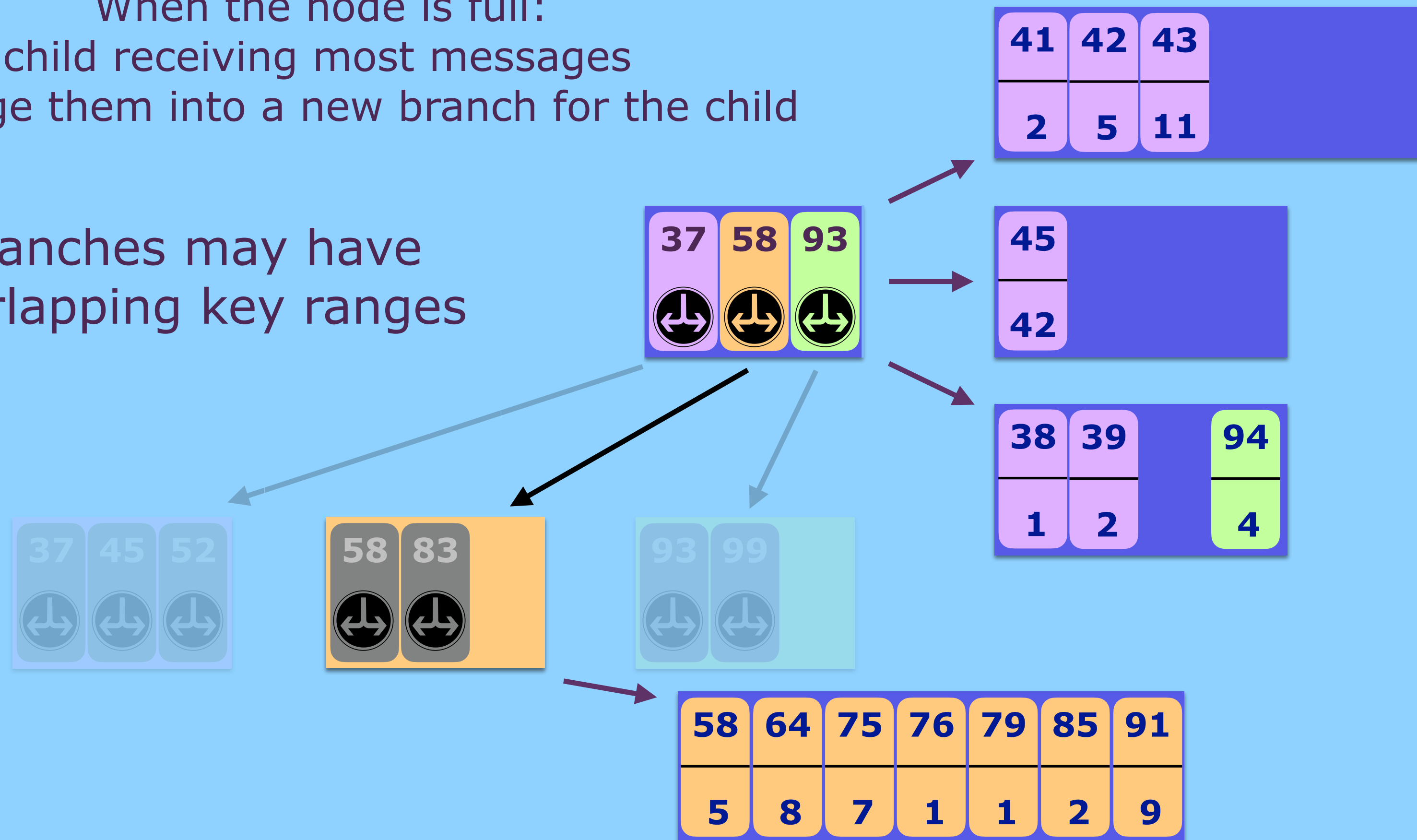
A Size-Tiered B ϵ -tree is a B ϵ -tree where the buffer is stored discontinuously

The fullness threshold is:
Fanout \times Average Buffer Size

When the node is full:

1. Pick child receiving most messages
2. Merge them into a new branch for the child

Branches may have overlapping key ranges



When new data is flushed into the trunk node...

...it is added as a new branch

The old branches do not need to be rewritten

Size-Tiered B ϵ -Trees

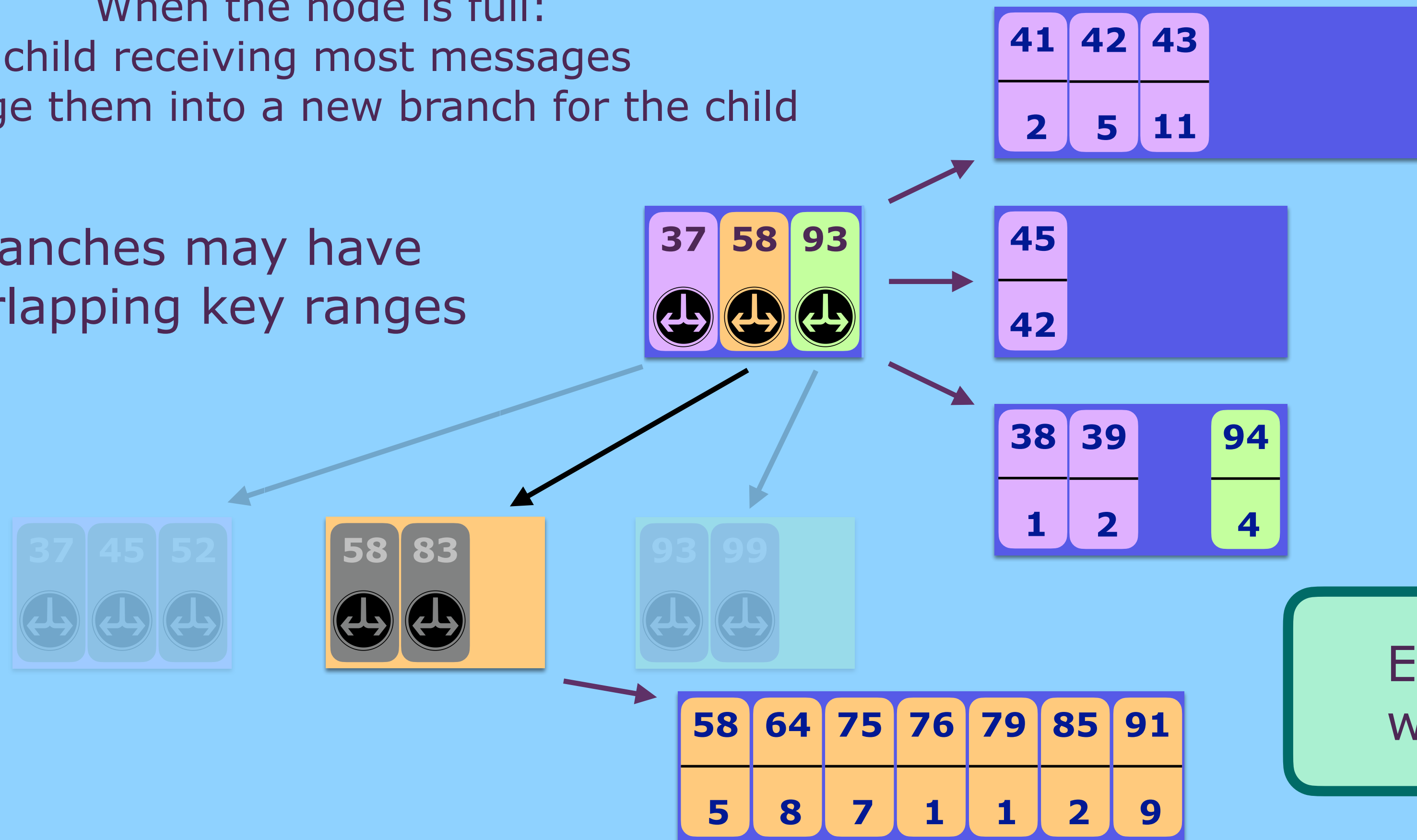
A Size-Tiered B ϵ -tree is a B ϵ -tree where the buffer is stored discontinuously

The fullness threshold is:
Fanout \times Average Buffer Size

When the node is full:

1. Pick child receiving most messages
2. Merge them into a new branch for the child

Branches may have overlapping key ranges



When new data is flushed into the trunk node...

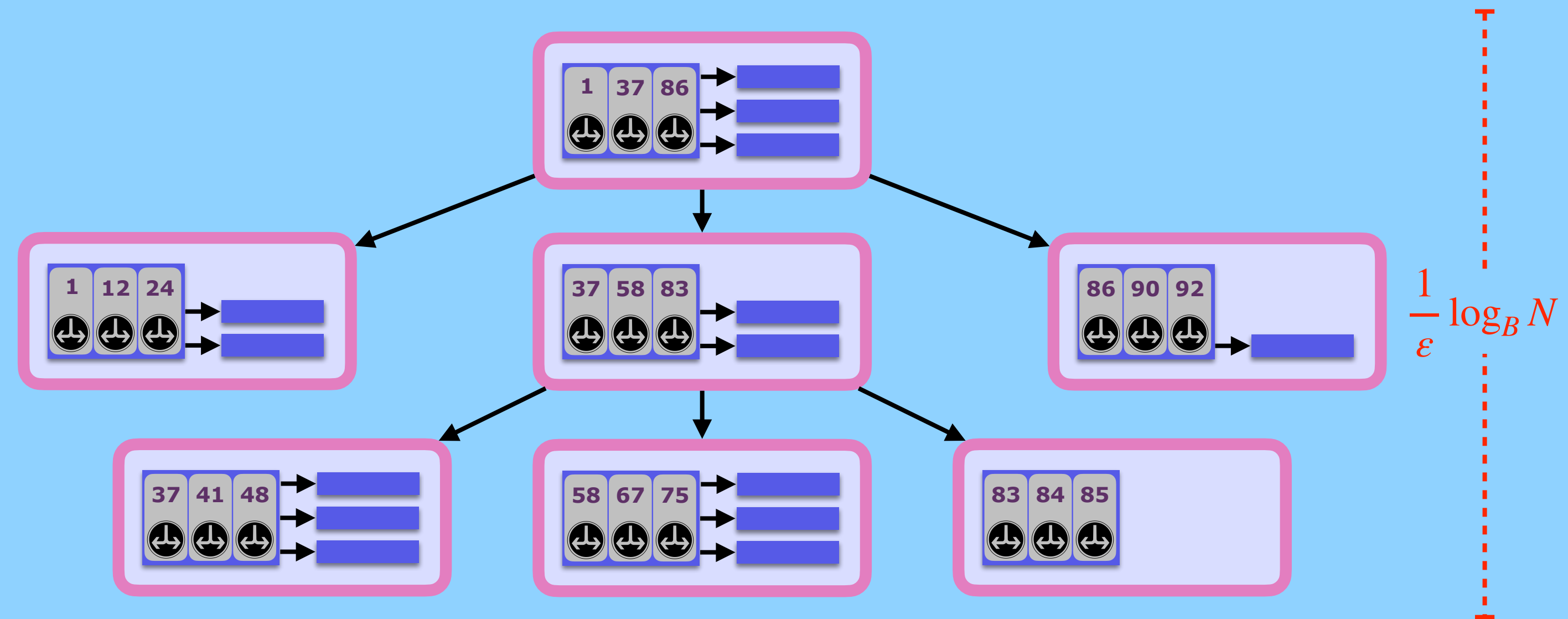
...it is added as a new branch

The old branches do not need to be rewritten

Each key-value pair is read/written once per trunk node

Size-Tiered B^ϵ -Trees

A Size-Tiered B^ϵ -tree is a B^ϵ -tree where the buffer is stored discontinuously



Each key-value pair is read/
written once per trunk node

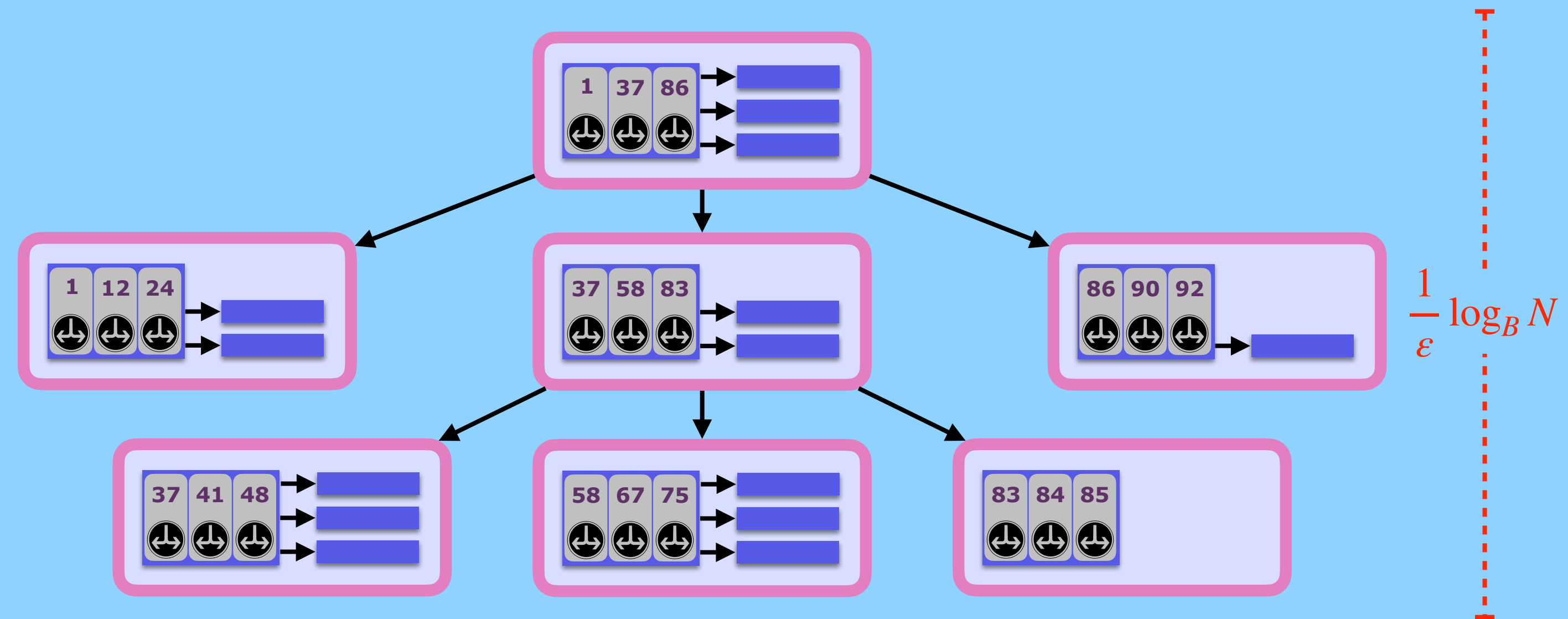
Size-Tiered B^ϵ -Trees

A Size-Tiered B^ϵ -tree is a B^ϵ -tree where the buffer is stored discontinuously

Work Amplification

B^ϵ -Tree: $O(B^\epsilon \times \log_{B^\epsilon} N)$

Size-Tiered B^ϵ -Tree: $O(\log_{B^\epsilon} N)$



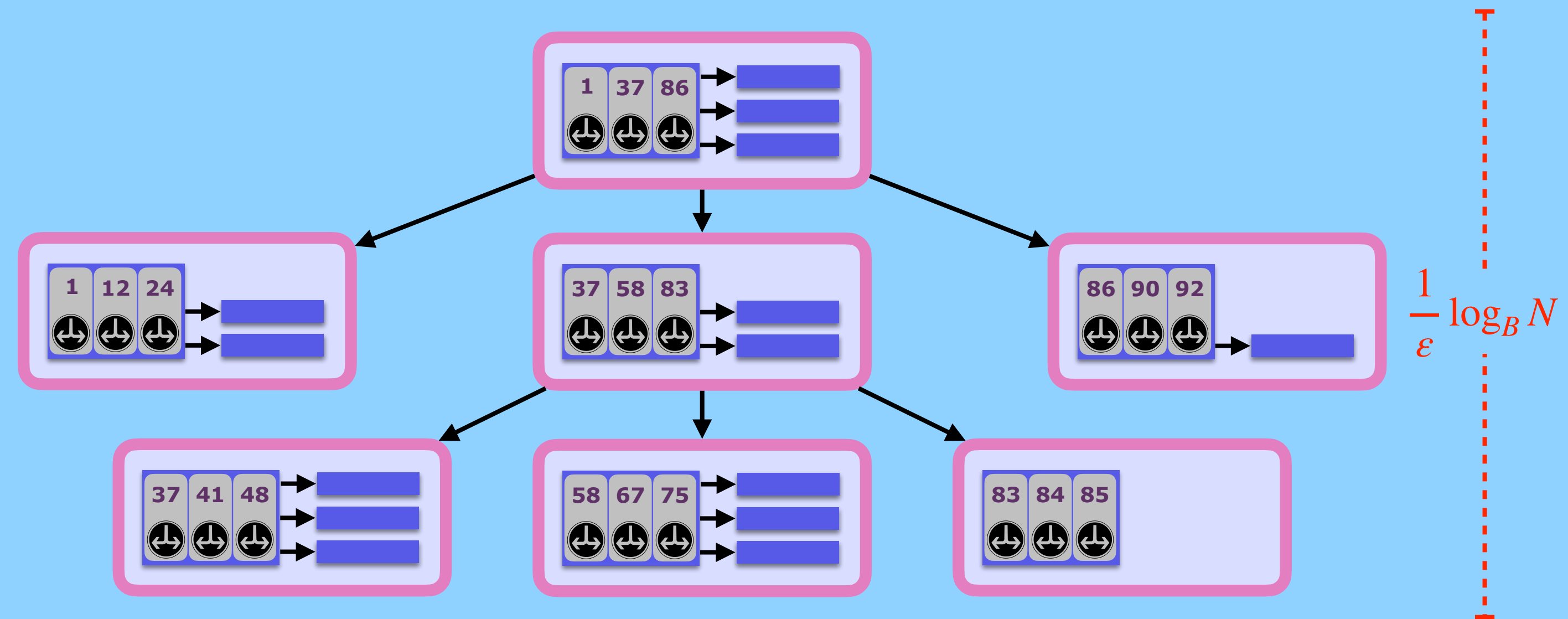
Each key-value pair is read/
written once per trunk node

Size-Tiered B^ϵ -Trees

A Size-Tiered B^ϵ -tree is a B^ϵ -tree where the buffer is stored discontinuously

Work Amplification

B^ϵ -Tree: $O(B^\epsilon \times \log_{B^\epsilon} N)$ $B^\epsilon \times$ less
 Size-Tiered B^ϵ -Tree: $O(\log_{B^\epsilon} N)$



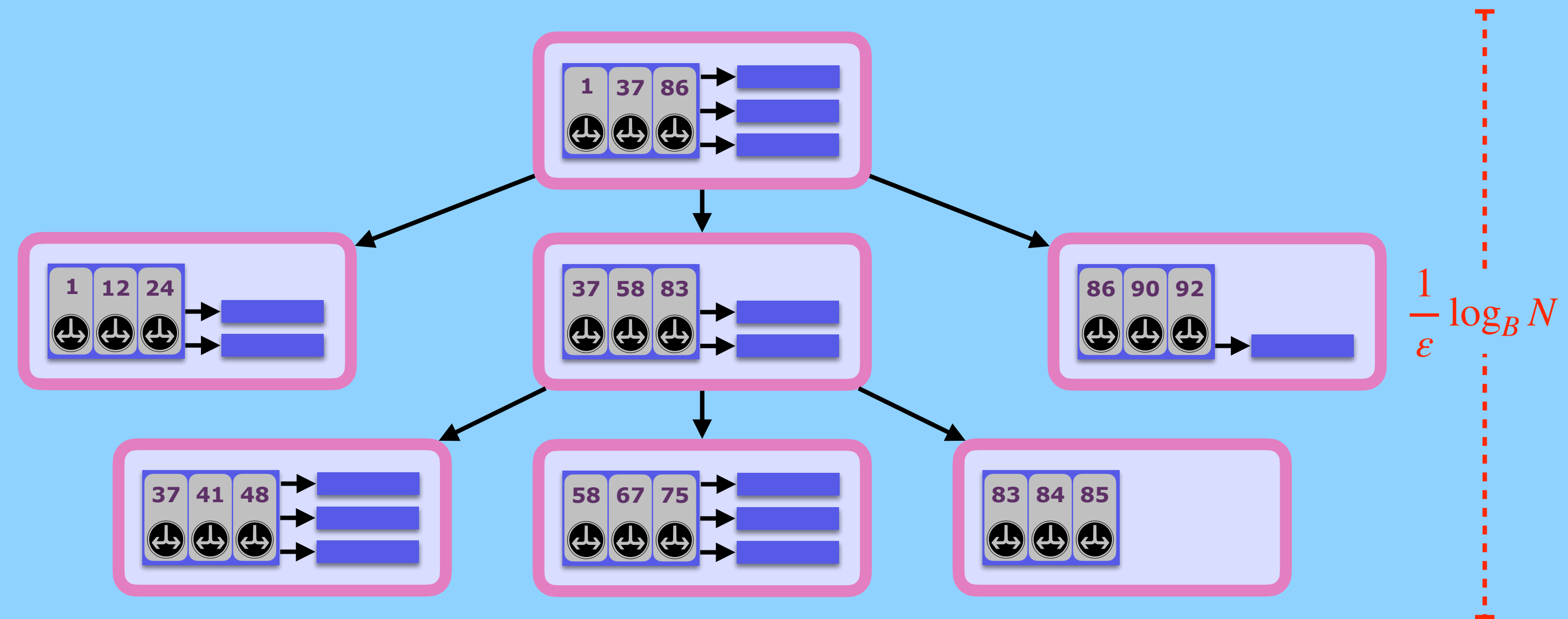
Each key-value pair is read/written once per trunk node

Size-Tiered B^ϵ -Trees

A Size-Tiered B^ϵ -tree is a B^ϵ -tree where the buffer is stored discontinuously

Work Amplification

B^ϵ -Tree: $O(B^\epsilon \times \log_{B^\epsilon} N)$ $B^\epsilon \times$ less
Size-Tiered B^ϵ -Tree: $O(\log_{B^\epsilon} N)$ 😊



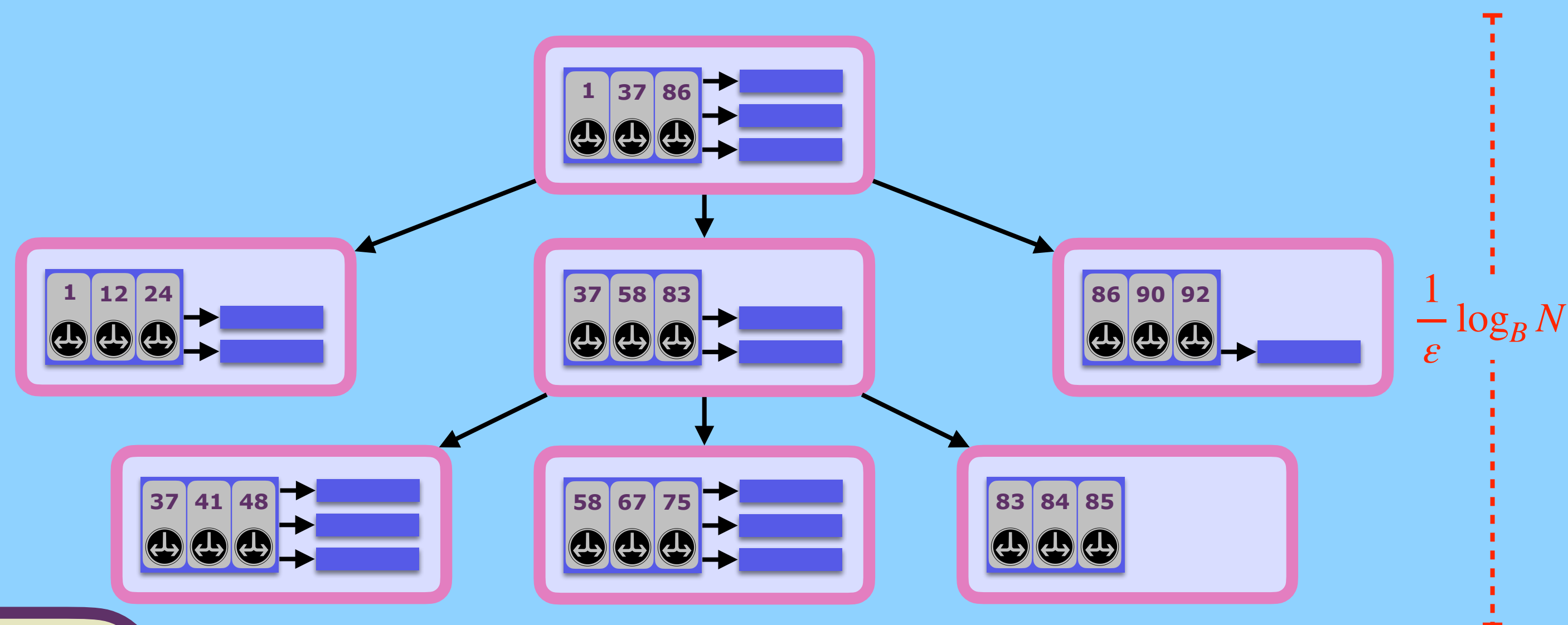
Each key-value pair is read/
written once per trunk node

Size-Tiered B^ϵ -Trees

A Size-Tiered B^ϵ -tree is a B^ϵ -tree where the buffer is stored discontinuously

Work Amplification

B^ϵ -Tree: $O(B^\epsilon \times \log_{B^\epsilon} N)$ $B^\epsilon \times$ less 😊
 Size-Tiered B^ϵ -Tree: $O(\log_{B^\epsilon} N)$ 😊



Insertion Cost

B^ϵ -Tree: $= O\left(\frac{1}{B} B^\epsilon \times \log_{B^\epsilon} \frac{N}{M}\right)$ $B^\epsilon \times$ less 😊
 Size-Tiered B^ϵ -Tree: $= O\left(\frac{1}{B} \log_{B^\epsilon} \frac{N}{M}\right)$ 😊

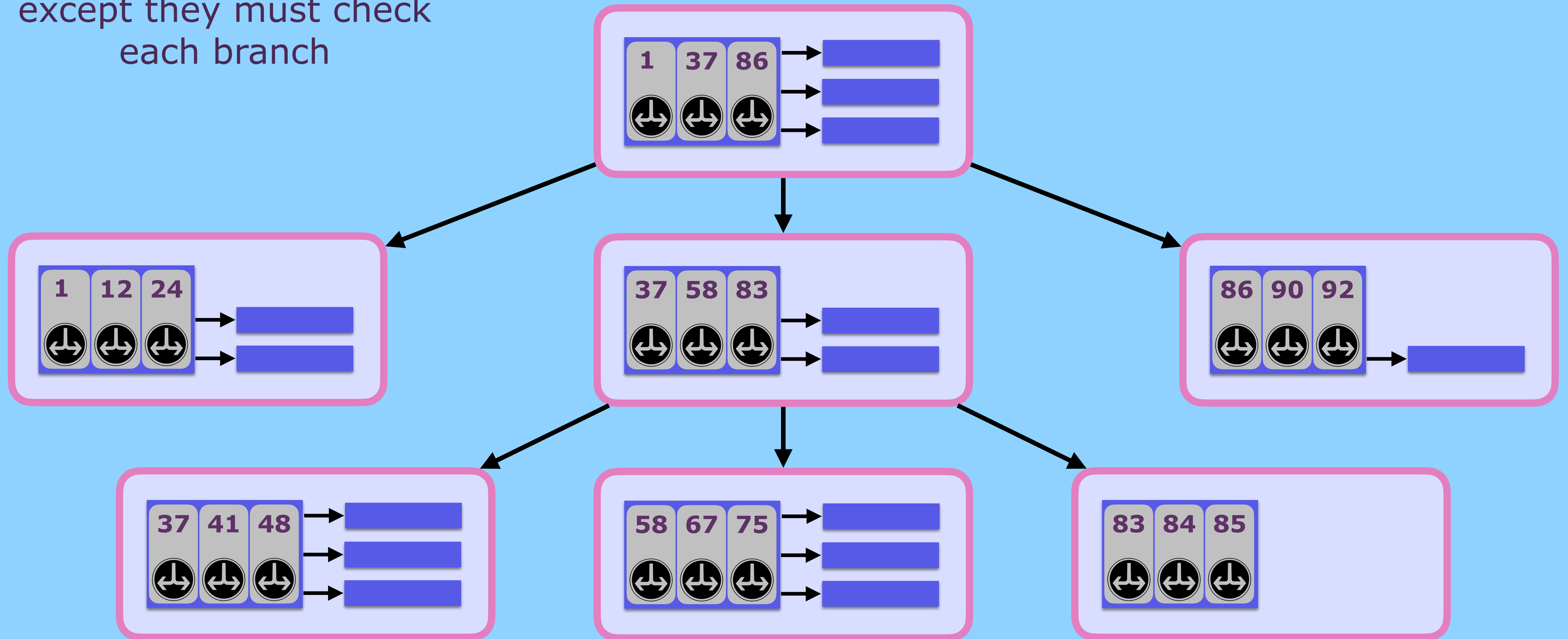
Each key-value pair is read/written once per trunk node

Lookups in Size-Tiered B^ϵ -Trees

Size-Tiered B^ϵ -Trees

Lookups in a STB^ϵ -tree are like lookups in a B^ϵ -tree, except they must check each branch

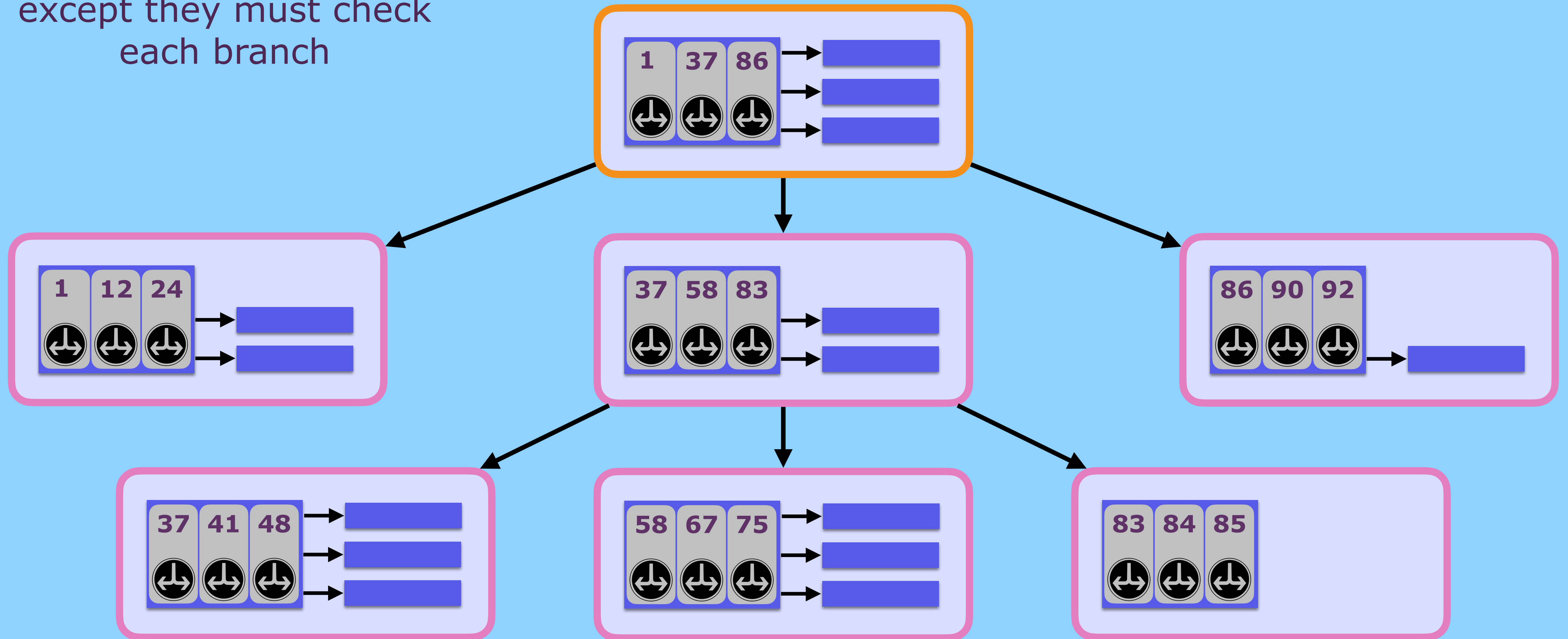
Query(71)



Size-Tiered B^ϵ -Trees

Lookups in a STB^ϵ -tree are like lookups in a B^ϵ -tree, except they must check each branch

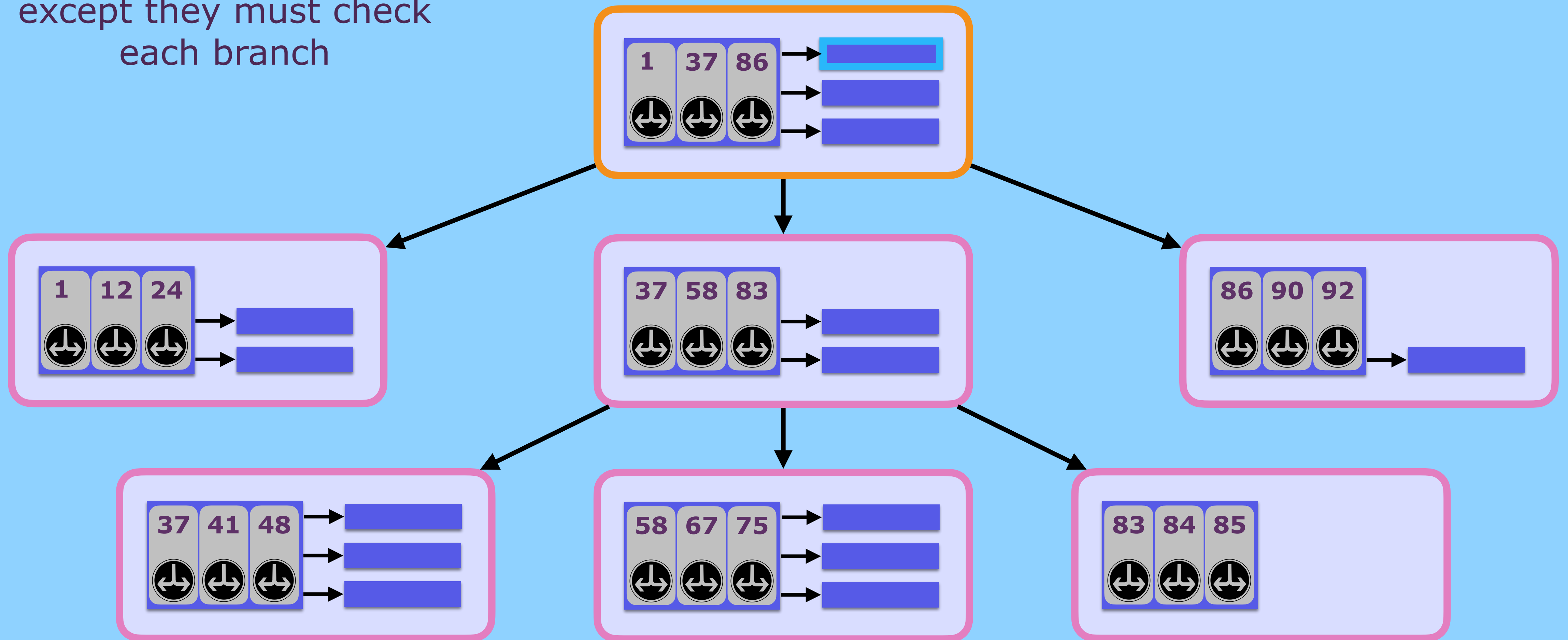
Query(71)



Size-Tiered B^ϵ -Trees

Lookups in a STB^ϵ -tree are like lookups in a B^ϵ -tree, except they must check each branch

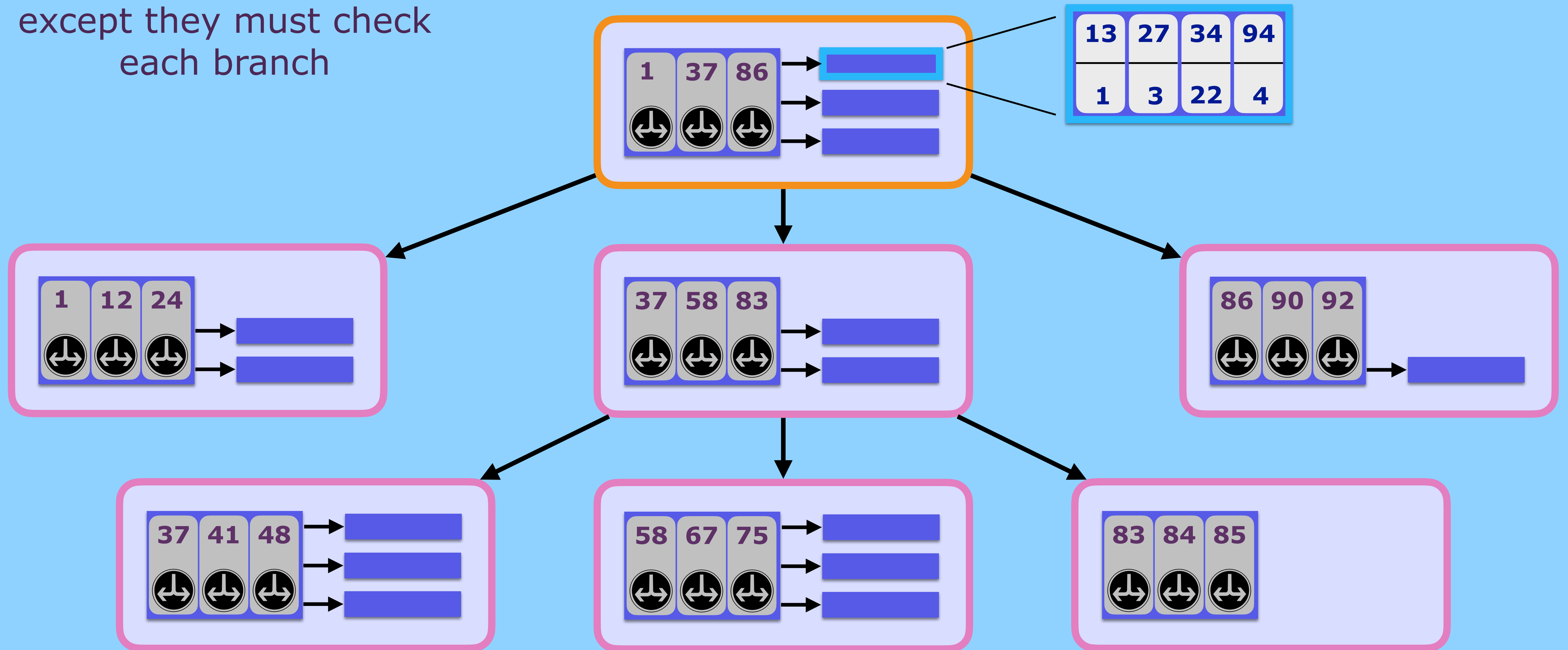
Query(71)



Size-Tiered B^ϵ -Trees

Lookups in a STB^ϵ -tree are like lookups in a B^ϵ -tree, except they must check each branch

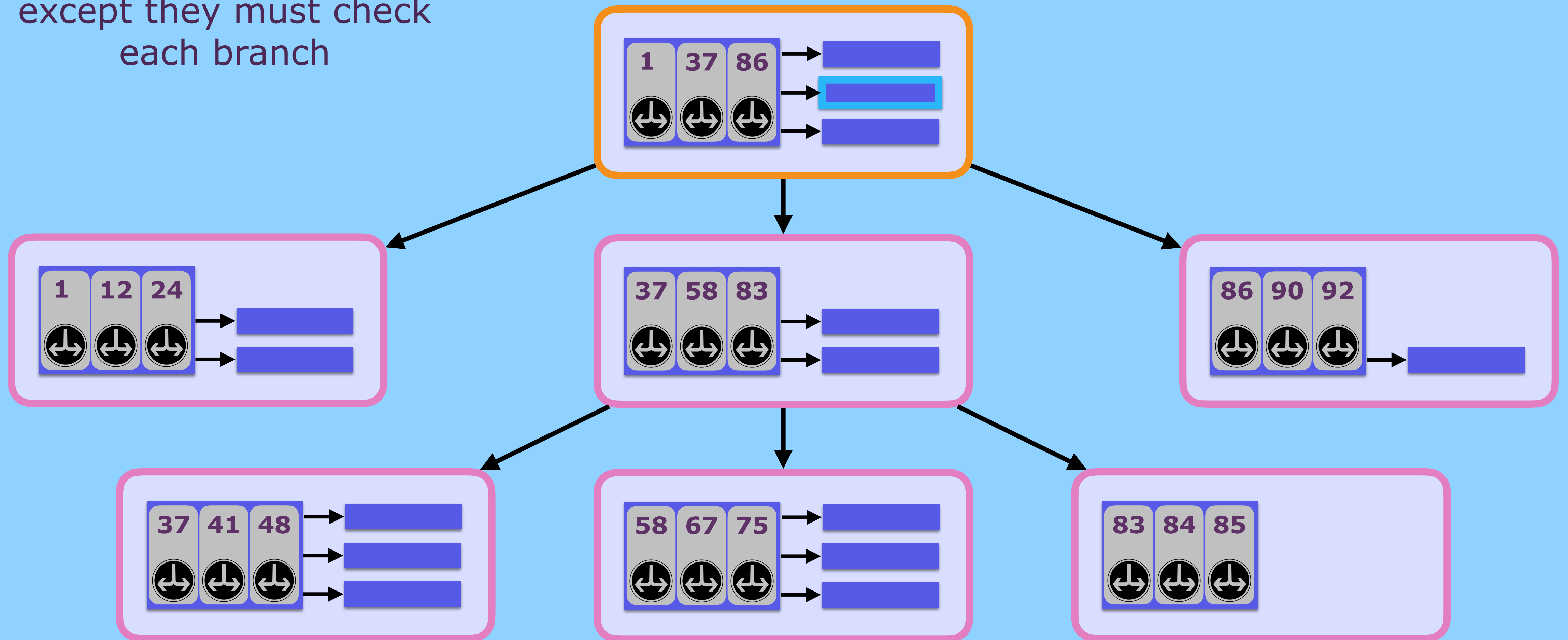
Query(71)



Size-Tiered B^ϵ -Trees

Lookups in a STB^ϵ -tree are like lookups in a B^ϵ -tree, except they must check each branch

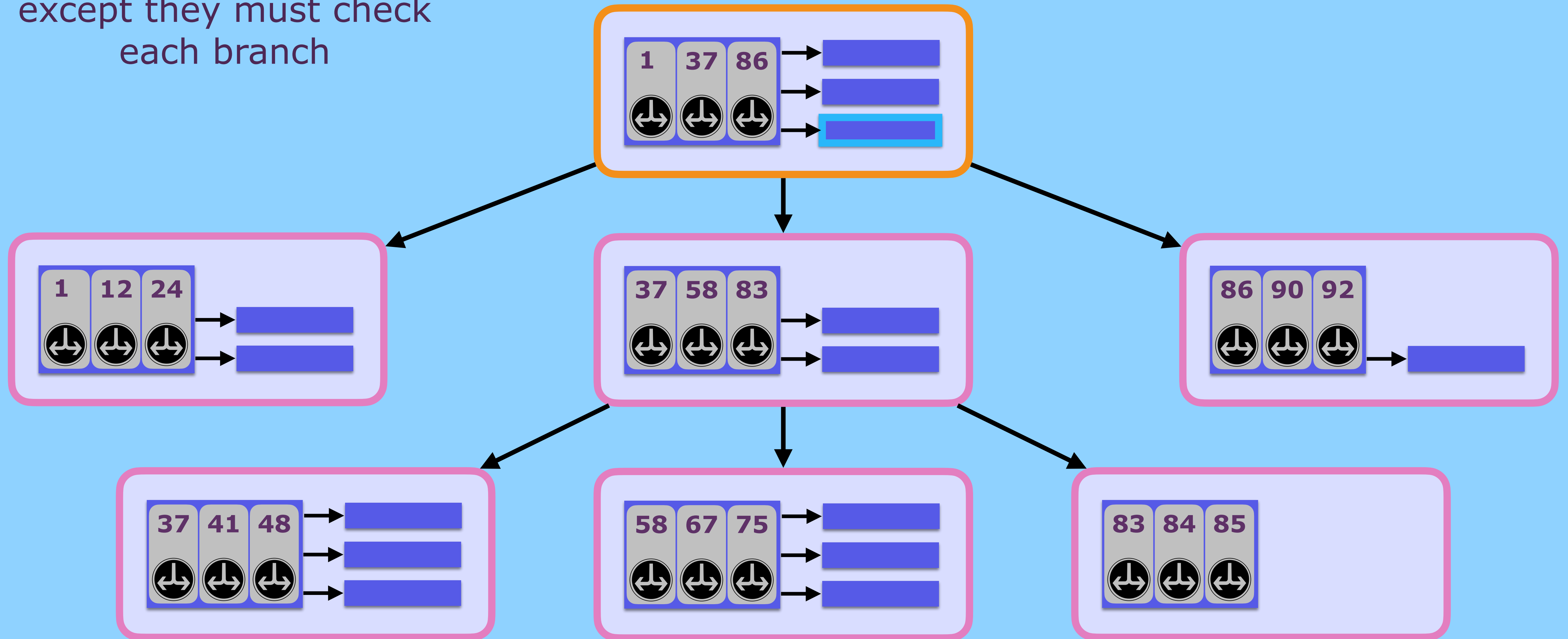
Query(71)



Size-Tiered B^ϵ -Trees

Lookups in a STB^ϵ -tree are like lookups in a B^ϵ -tree, except they must check each branch

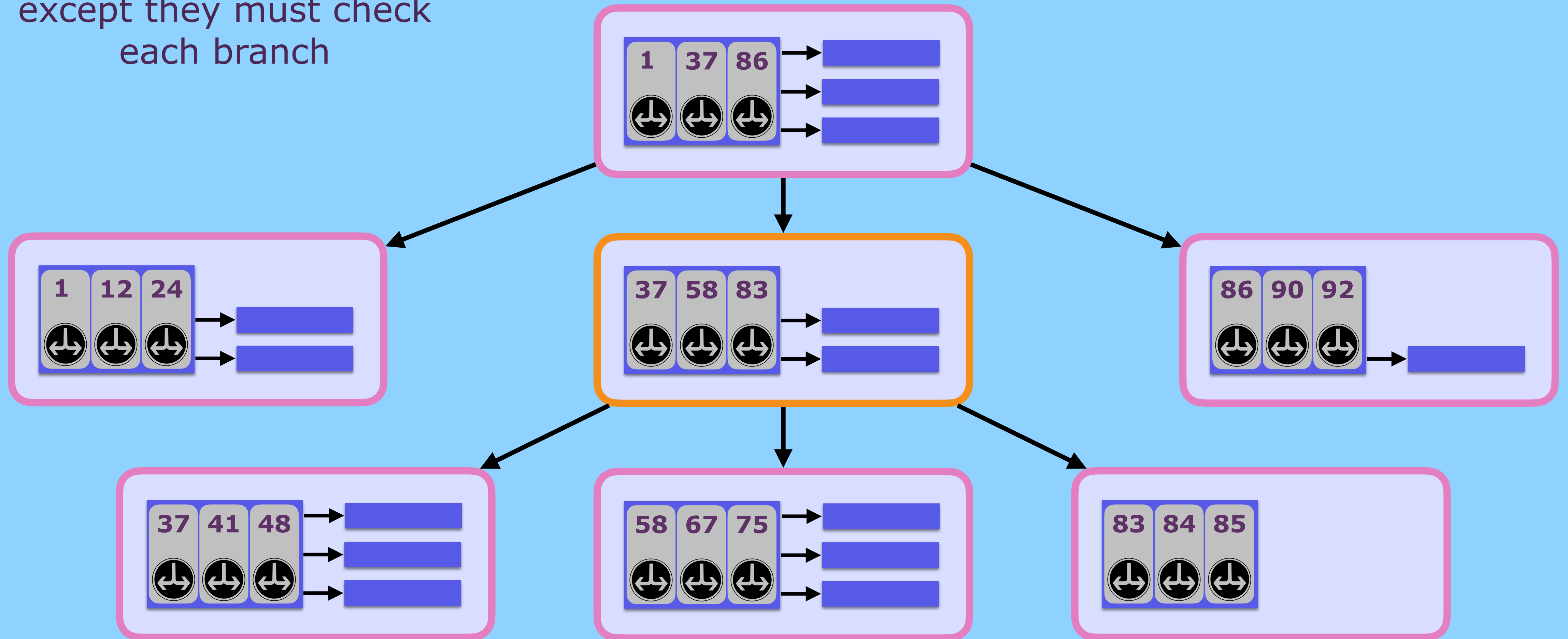
Query(71)



Size-Tiered B^ϵ -Trees

Lookups in a STB^ϵ -tree are like lookups in a B^ϵ -tree, except they must check each branch

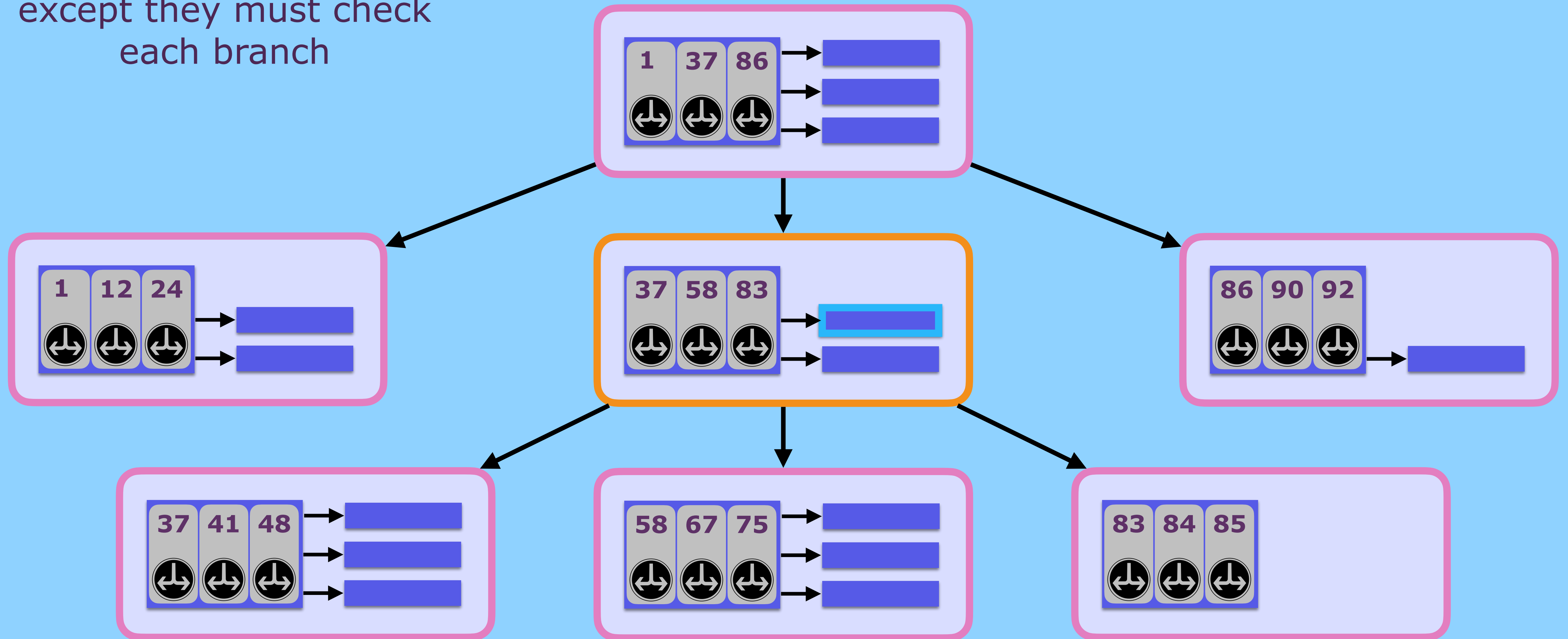
Query(71)



Size-Tiered B^ϵ -Trees

Lookups in a STB^ϵ -tree are like lookups in a B^ϵ -tree, except they must check each branch

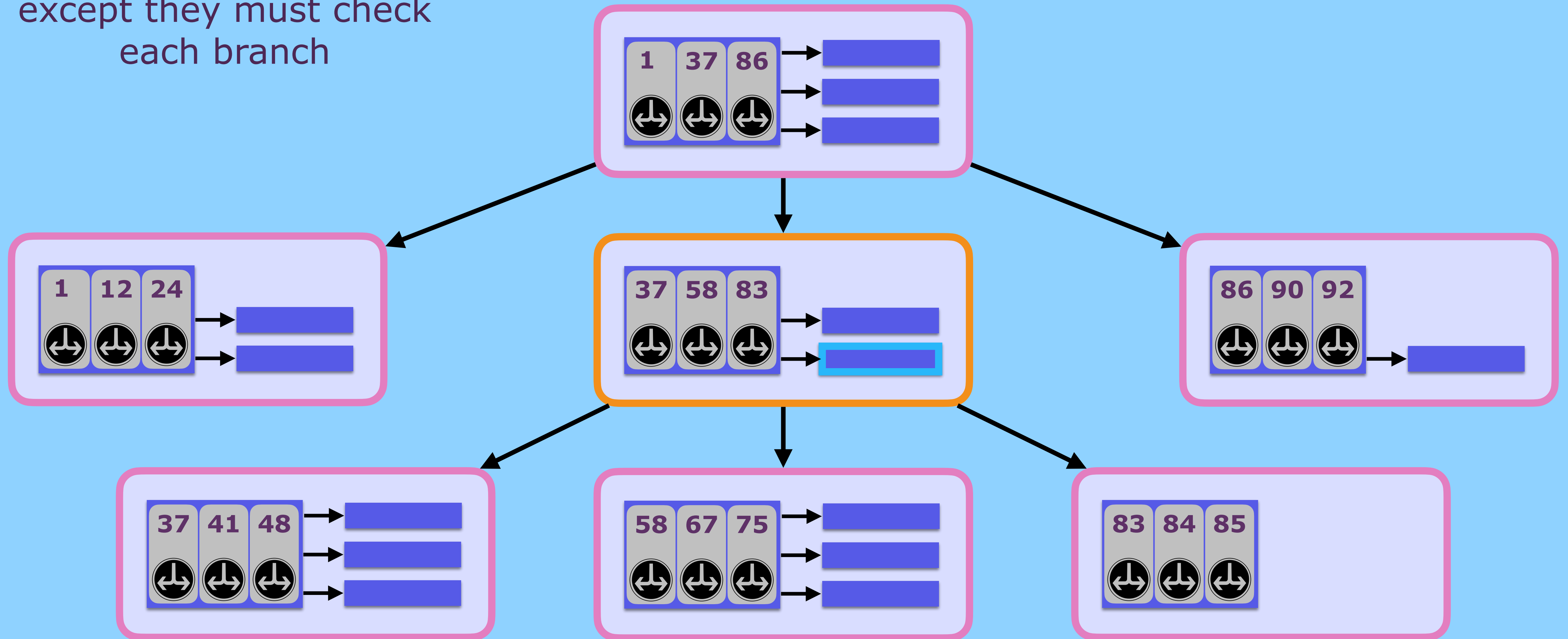
Query(71)



Size-Tiered B^ϵ -Trees

Lookups in a STB^ϵ -tree are like lookups in a B^ϵ -tree, except they must check each branch

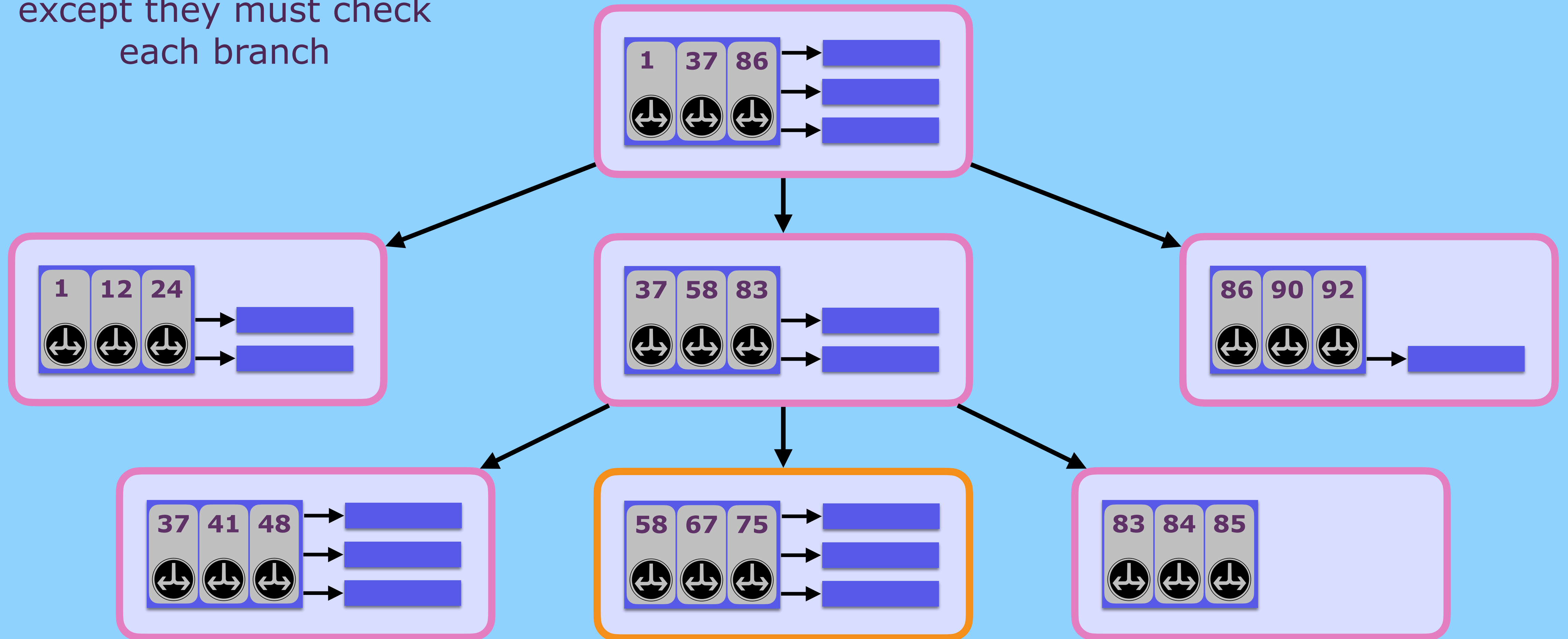
Query(71)



Size-Tiered B^ϵ -Trees

Lookups in a STB^ϵ -tree are like lookups in a B^ϵ -tree, except they must check each branch

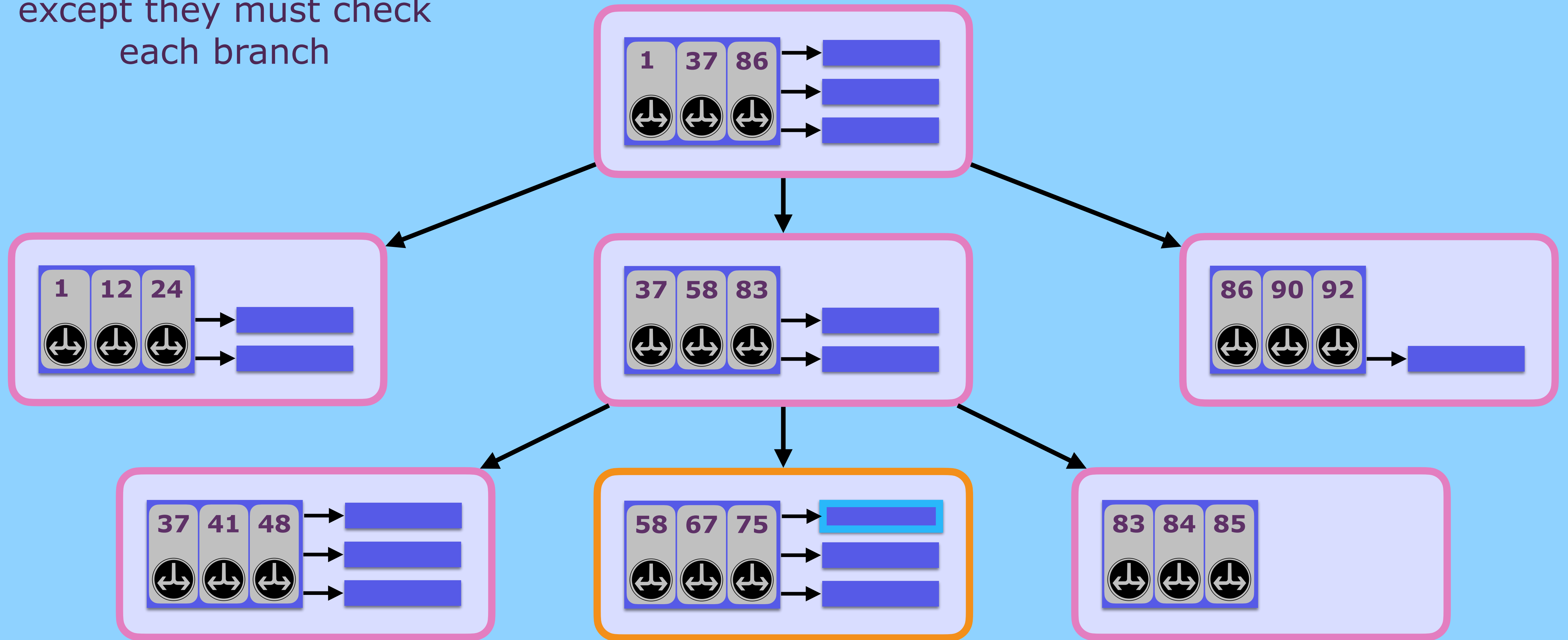
Query(71)



Size-Tiered B^ϵ -Trees

Lookups in a STB^ϵ -tree are like lookups in a B^ϵ -tree, except they must check each branch

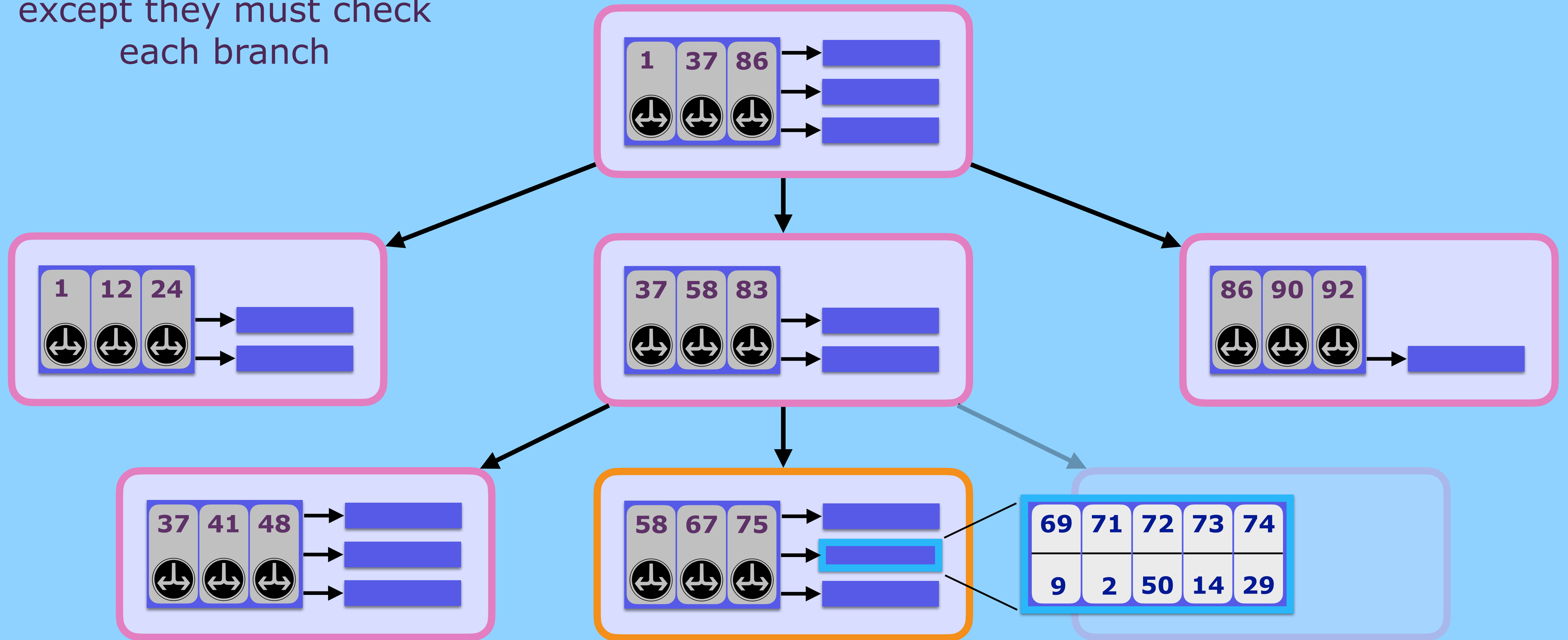
Query(71)



Size-Tiered B^ϵ -Trees

Lookups in a STB^ϵ -tree are like lookups in a B^ϵ -tree, except they must check each branch

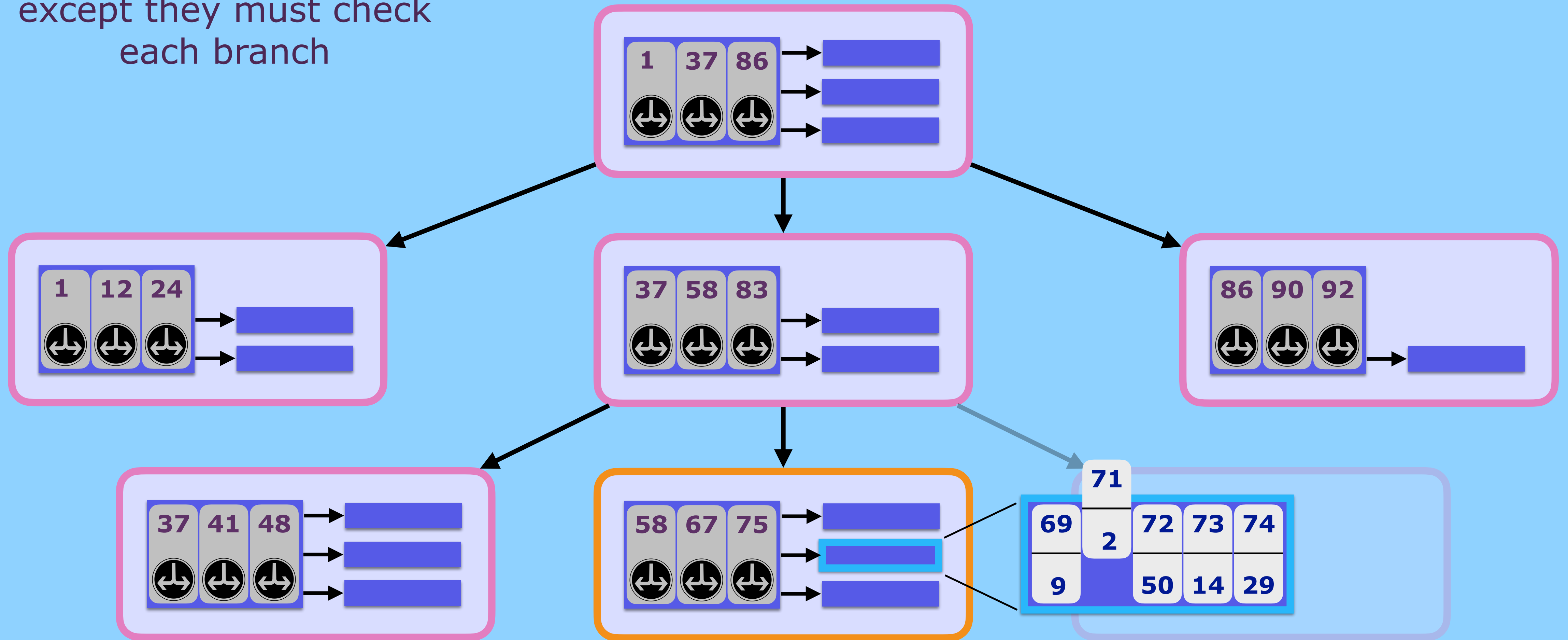
Query(71)



Size-Tiered B^ϵ -Trees

Lookups in a STB^ϵ -tree are like lookups in a B^ϵ -tree, except they must check each branch

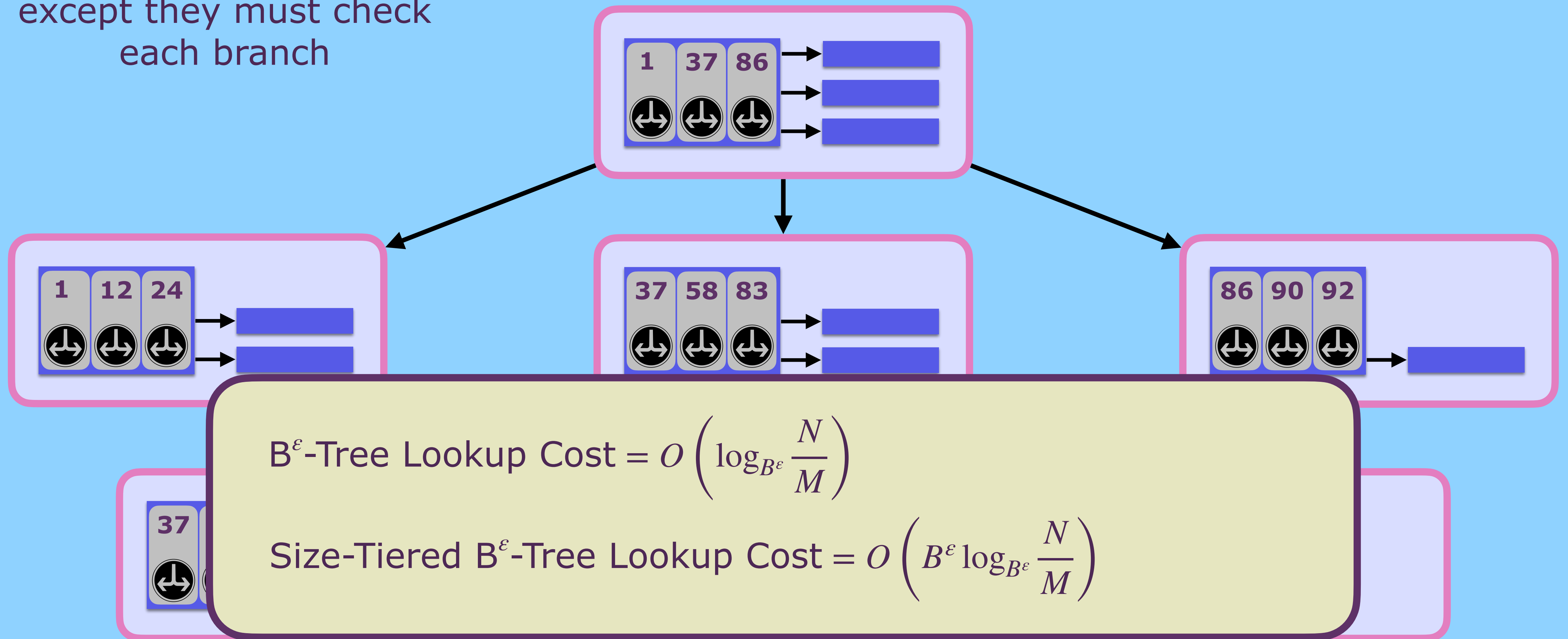
Query(71) \rightarrow 2



Size-Tiered B^ϵ -Trees

Lookups in a STB^ϵ -tree are like lookups in a B^ϵ -tree, except they must check each branch

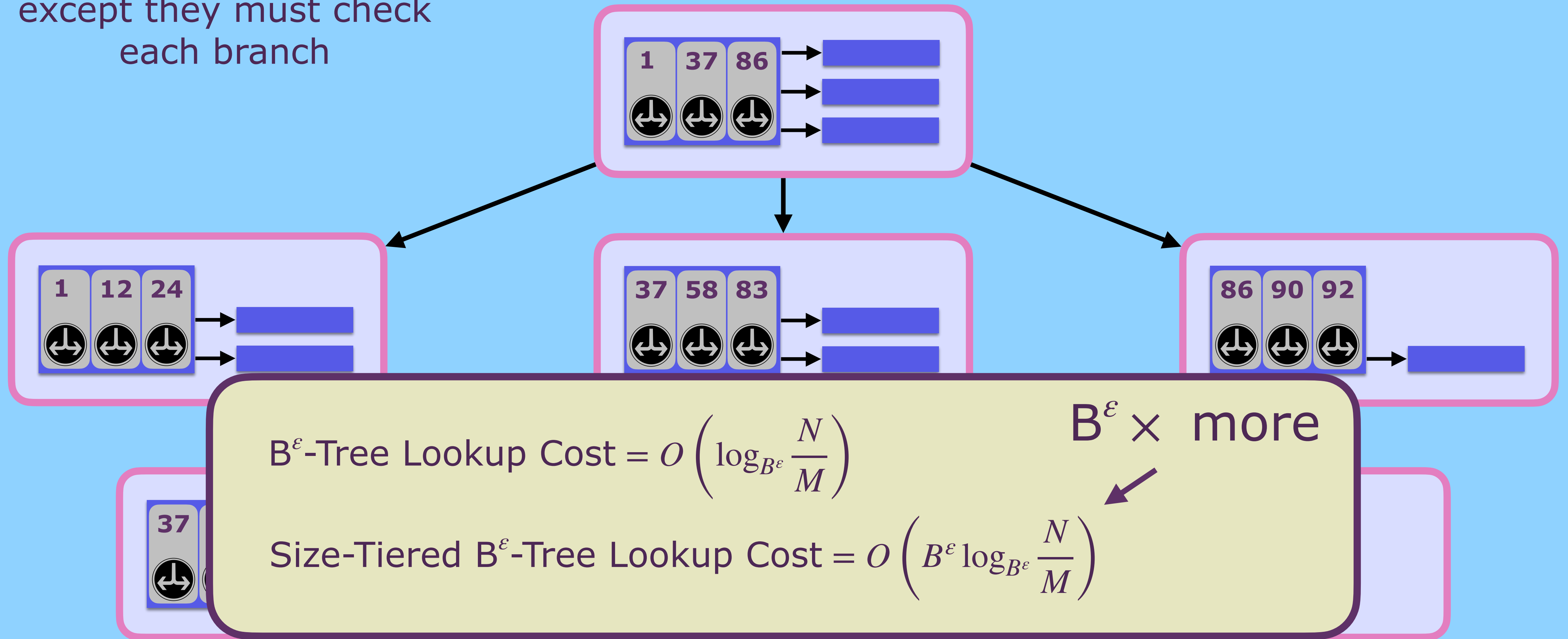
Query(71)



Size-Tiered B^ϵ -Trees

Lookups in a STB^ϵ -tree are like lookups in a B^ϵ -tree, except they must check each branch

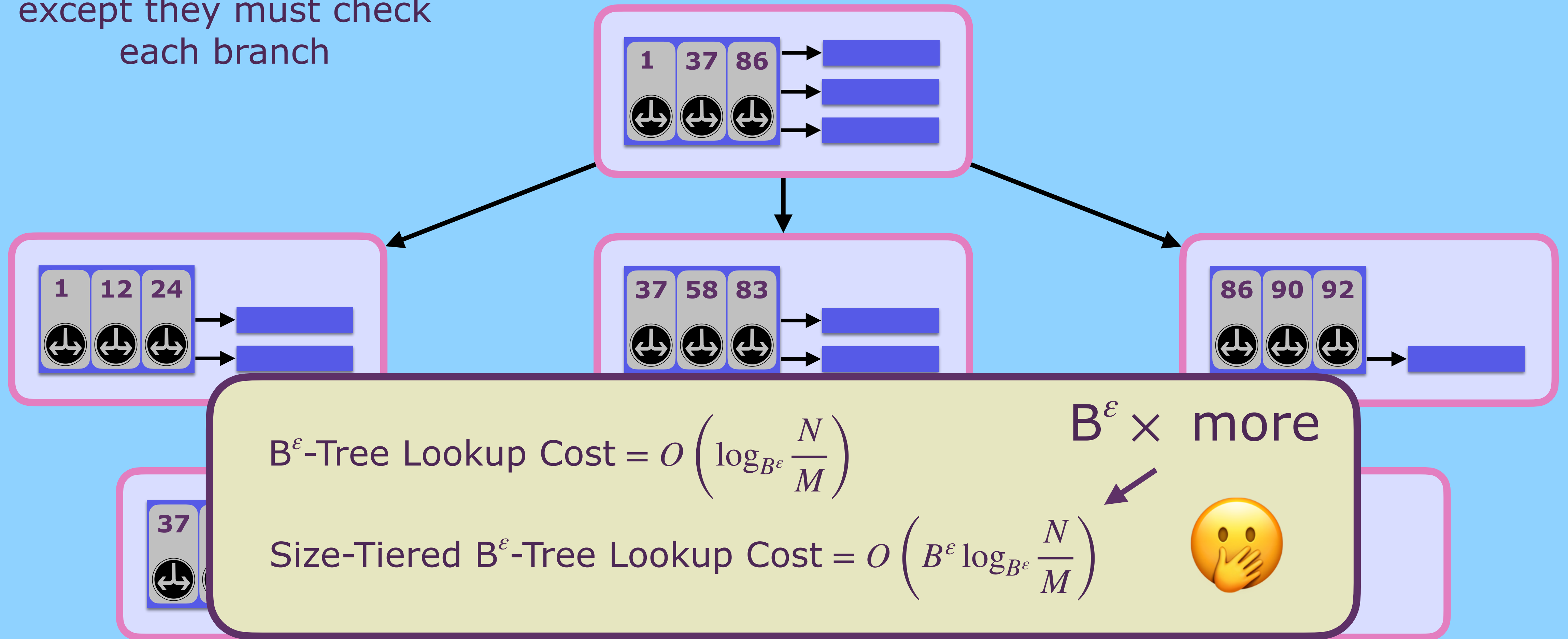
Query(71)



Size-Tiered B^ϵ -Trees

Lookups in a STB^ϵ -tree are like lookups in a B^ϵ -tree, except they must check each branch

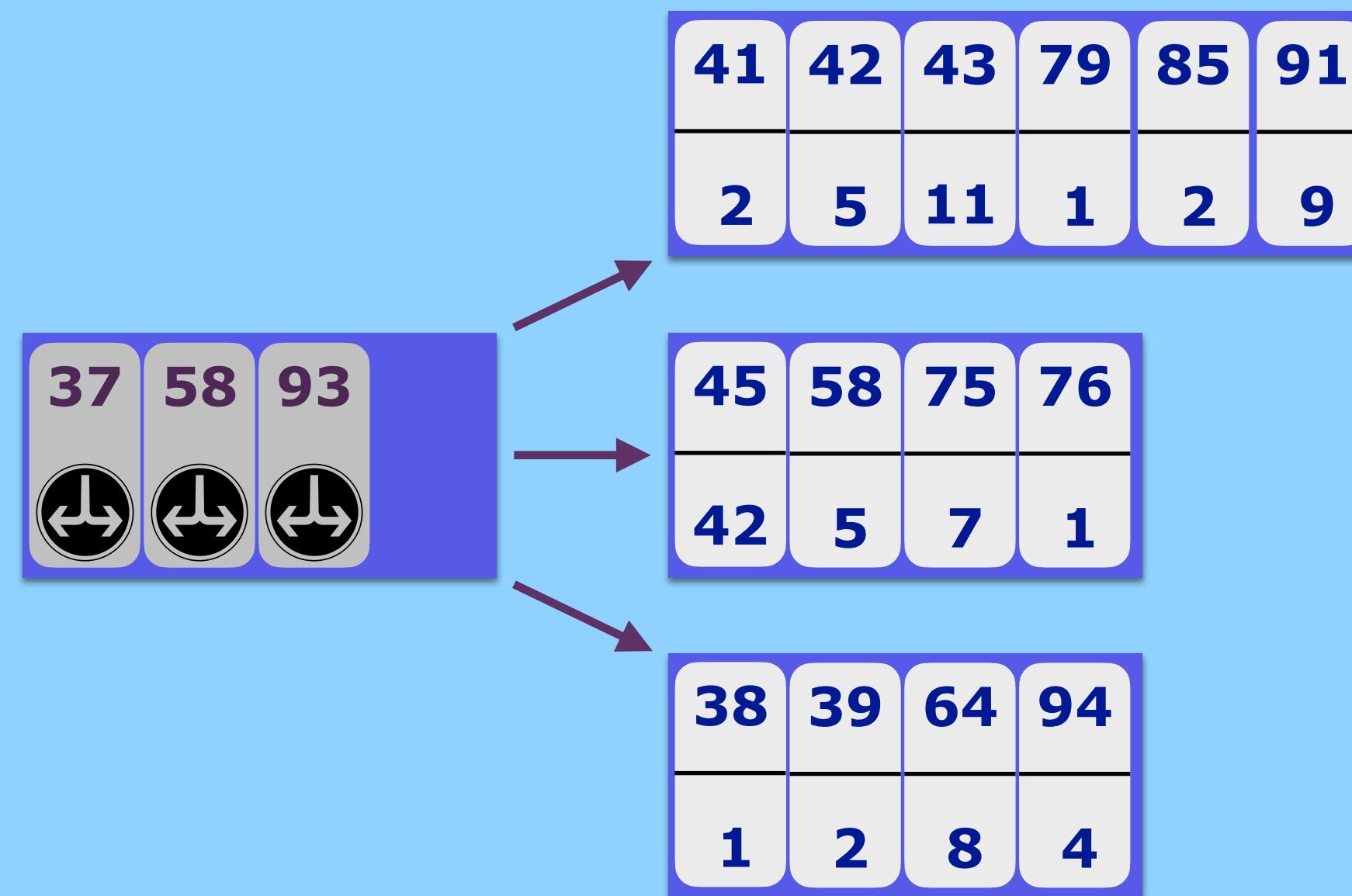
Query(71)



Fixing Lookups in Size-Tiered B^ϵ -Trees

Fixing Lookups in Size-Tiered B_ϵ -Trees

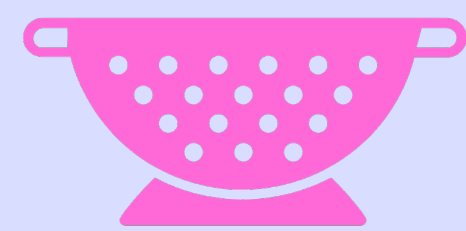
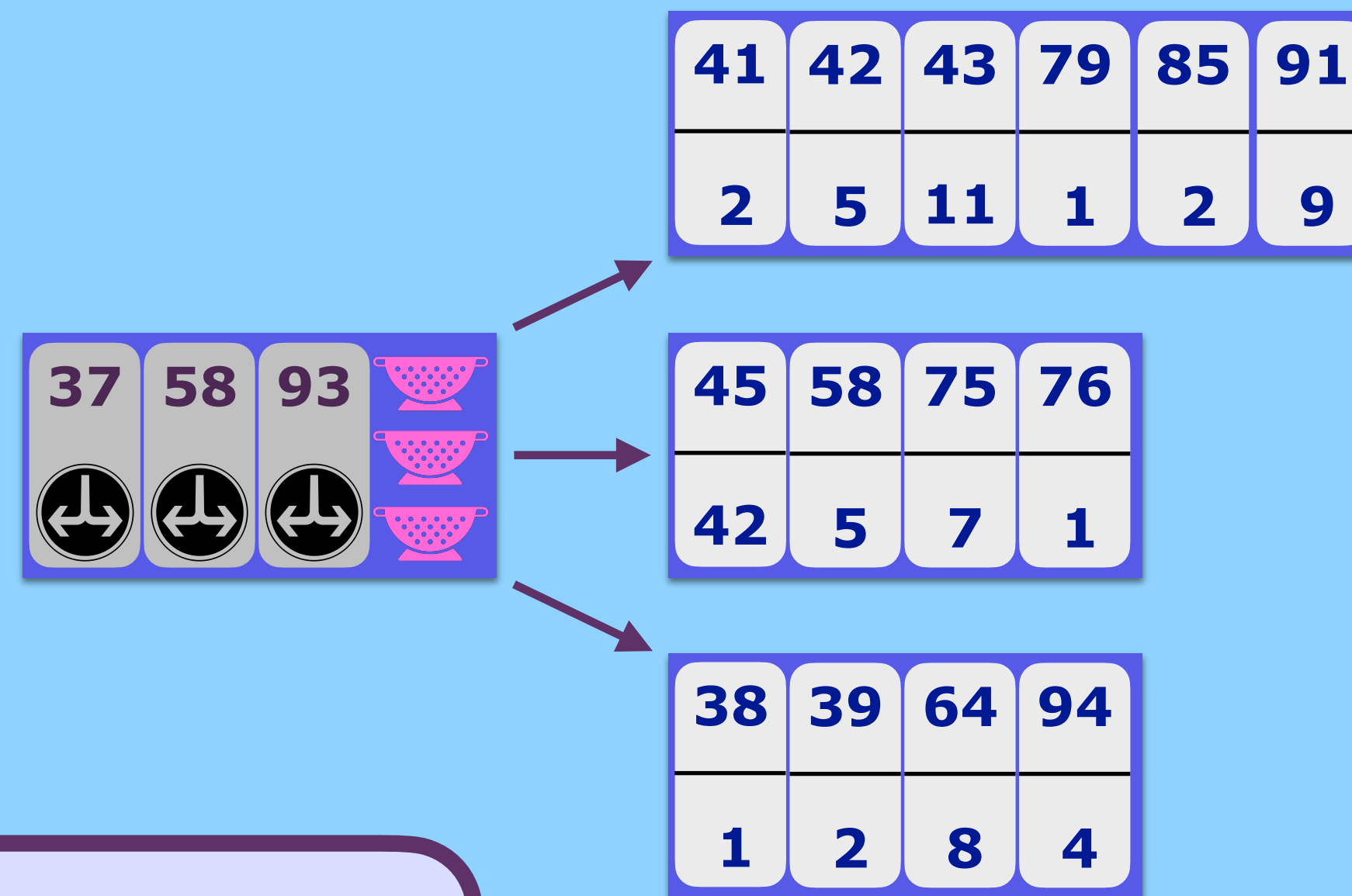
The problem is that each node has multiple branches



Fixing Lookups in Size-Tiered B_ϵ -Trees

The problem is that each node has multiple branches

Idea: use filters to avoid searching them



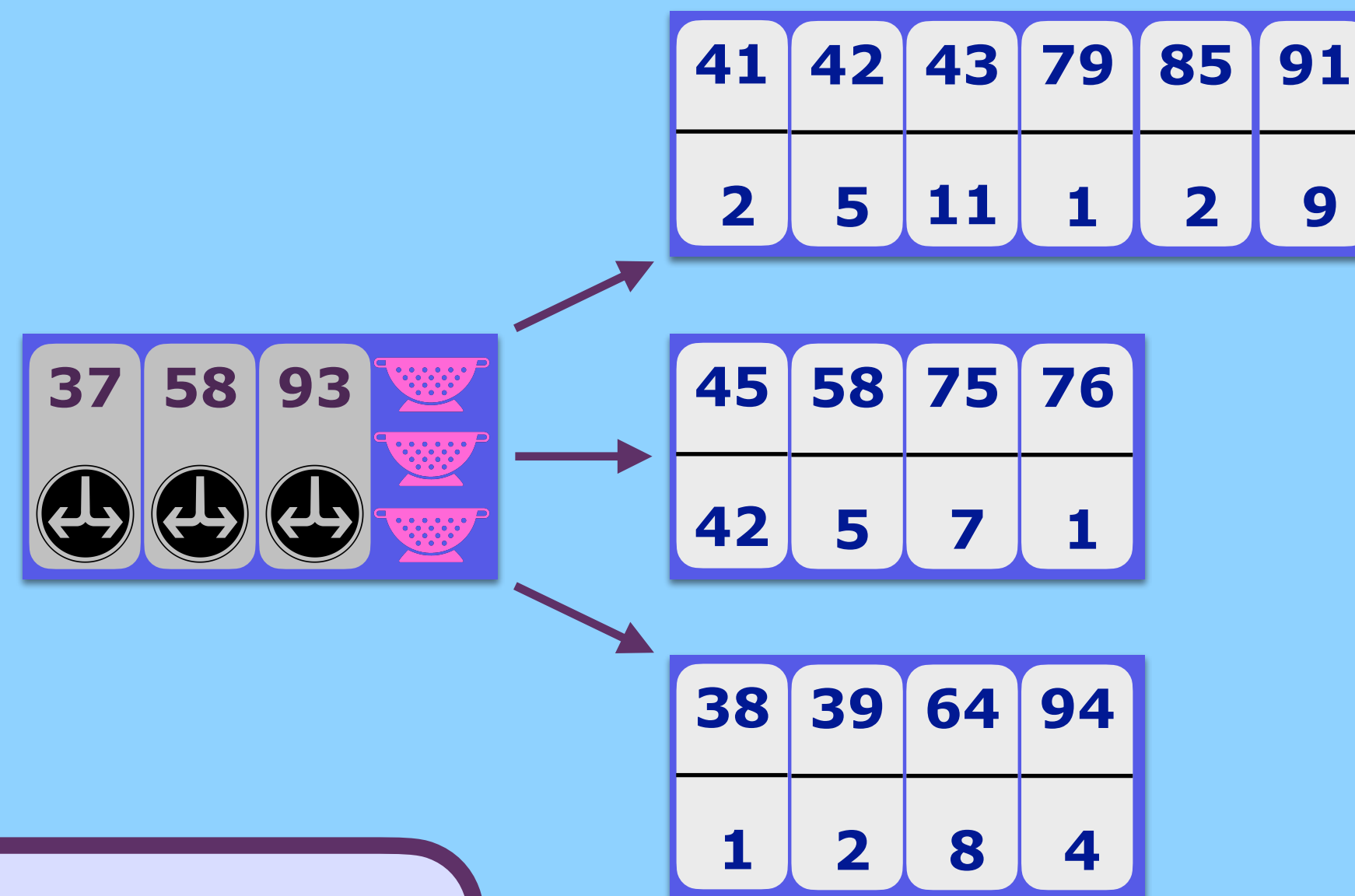
A filter is a probabilistic data structure with answers membership with no false negatives

Examples: Bloom, cuckoo, quotient

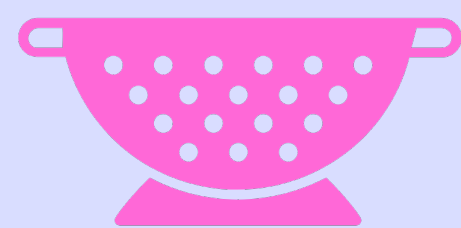
Fixing Lookups in Size-Tiered B^ϵ -Trees

The problem is that each node has multiple branches

Idea: use filters to avoid searching them



Now a lookup will only search those branches which contain the key (plus rare false positives)



A filter is a probabilistic data structure with answers membership with no false negatives

Examples: Bloom, cuckoo, quotient

Fixing Lookups in Size-Tiered B^ϵ -Trees

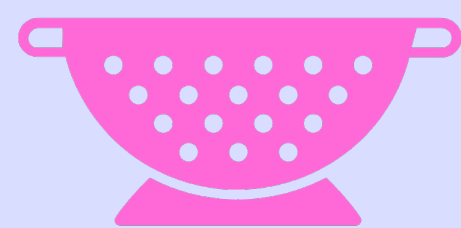
Query(64)

The problem is that each node has multiple branches

Idea: use filters to avoid searching them



Now a lookup will only search those branches which contain the key (plus rare false positives)



A filter is a probabilistic data structure with answers membership with no false negatives

Examples: Bloom, cuckoo, quotient

Fixing Lookups in Size-Tiered B^ϵ -Trees

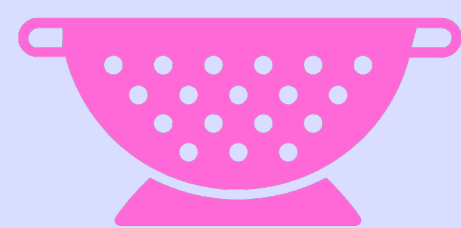
Query(64)

The problem is that each node has multiple branches



Idea: use filters to avoid searching them

Now a lookup will only search those branches which contain the key (plus rare false positives)



A filter is a probabilistic data structure with answers membership with no false negatives

Examples: Bloom, cuckoo, quotient

Fixing Lookups in Size-Tiered B^ϵ -Trees

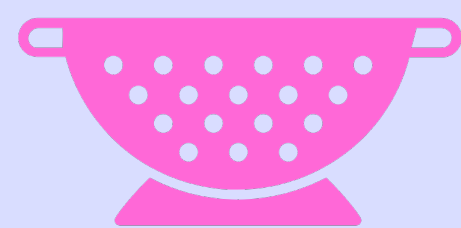
The problem is that each node has multiple branches

Idea: use filters to avoid searching them

Query(64)



Now a lookup will only search those branches which contain the key (plus rare false positives)



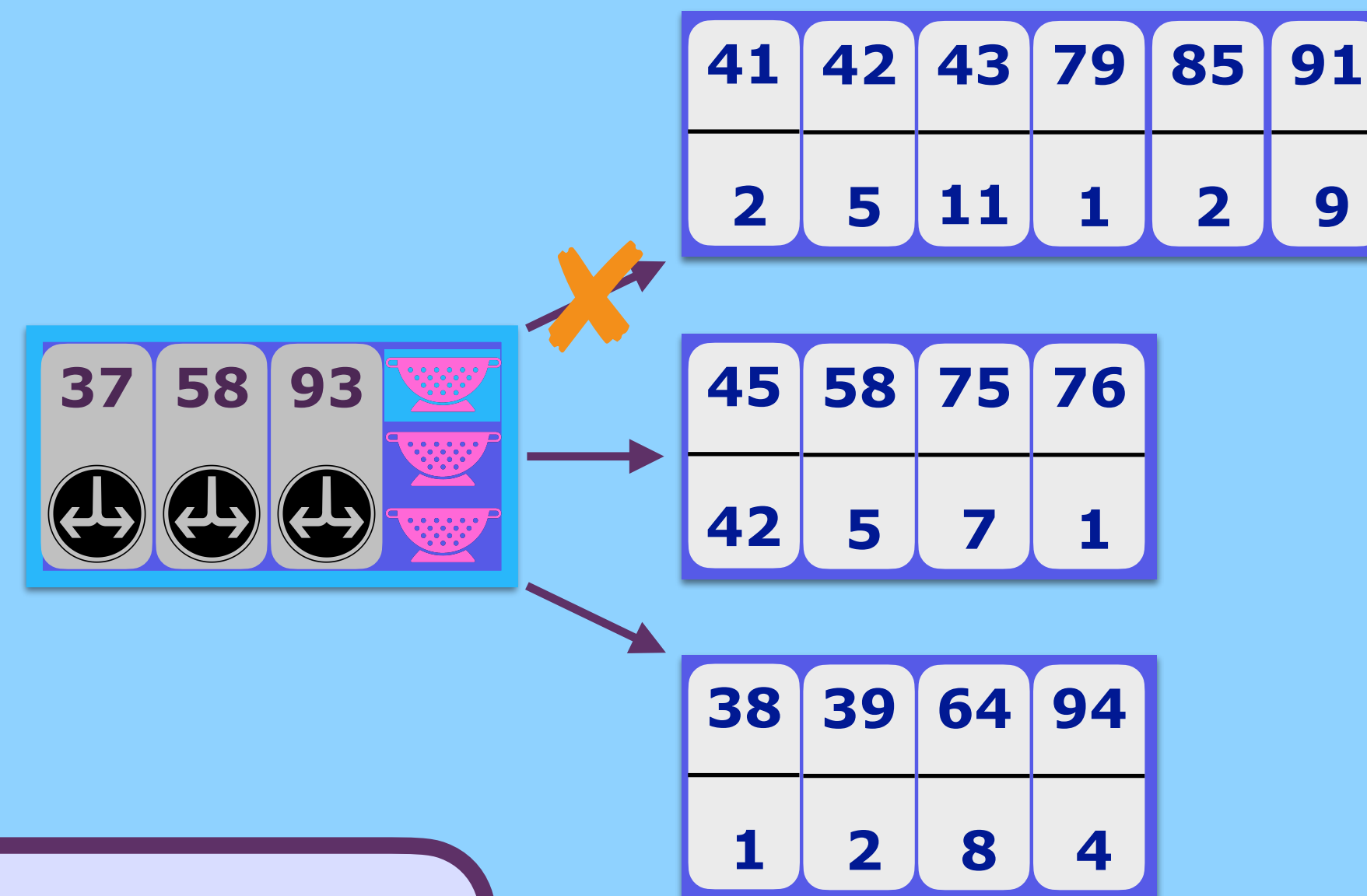
A filter is a probabilistic data structure with answers membership with no false negatives

Examples: Bloom, cuckoo, quotient

Fixing Lookups in Size-Tiered B^ϵ -Trees

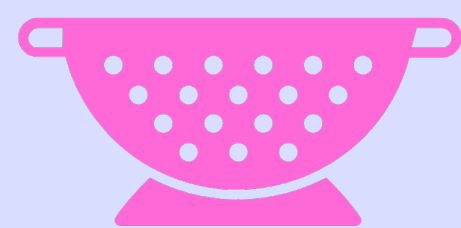
Query(64)

The problem is that each node has multiple branches



Idea: use filters to avoid searching them

Now a lookup will only search those branches which contain the key (plus rare false positives)



A filter is a probabilistic data structure with answers membership with no false negatives

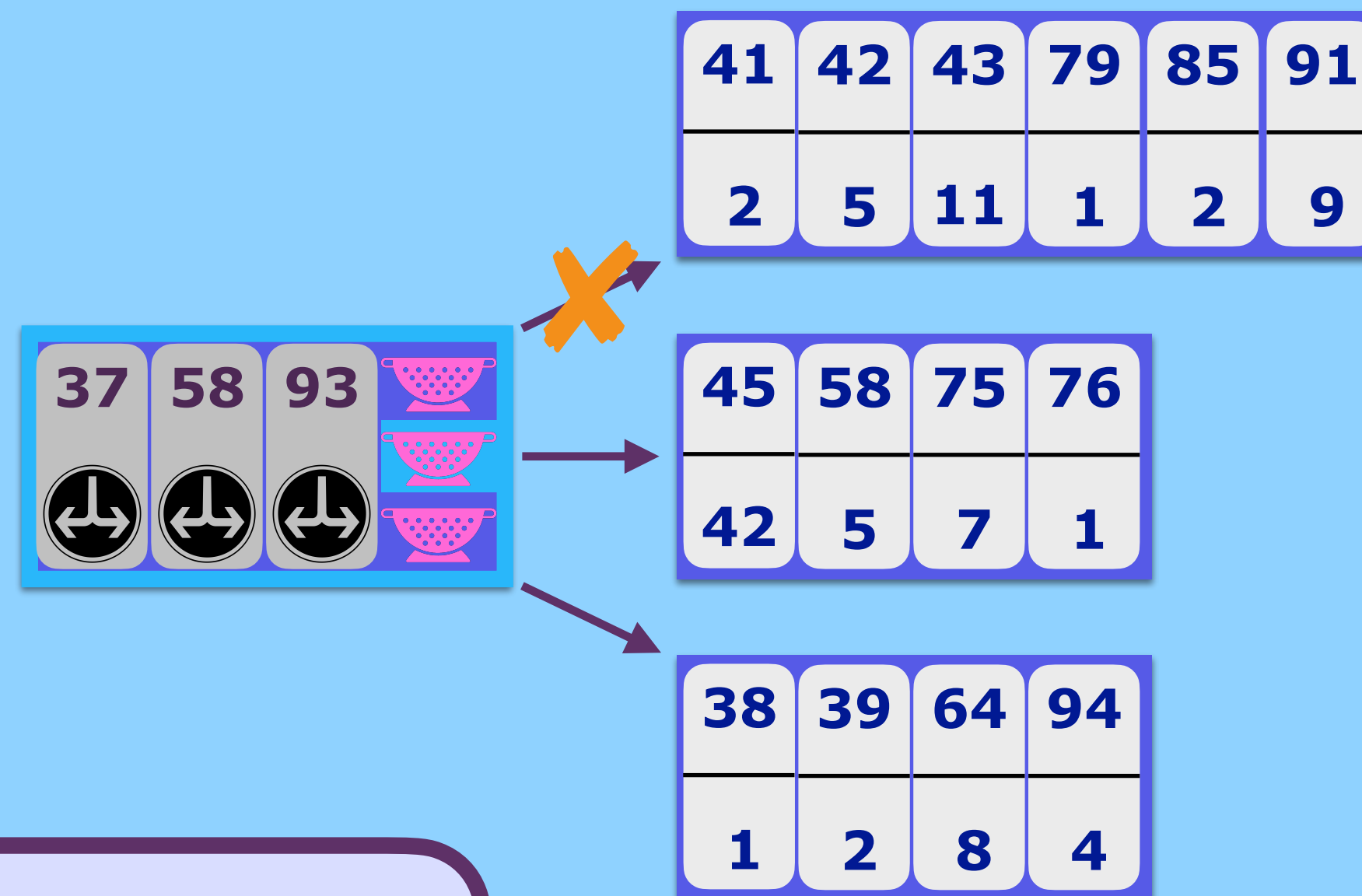
Examples: Bloom, cuckoo, quotient

Fixing Lookups in Size-Tiered B^ϵ -Trees

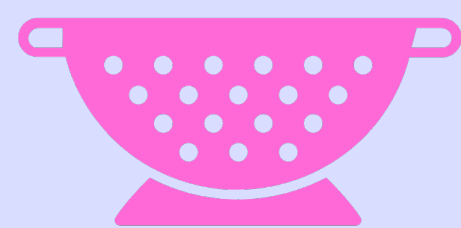
The problem is that each node has multiple branches

Idea: use filters to avoid searching them

Query(64)



Now a lookup will only search those branches which contain the key (plus rare false positives)



A filter is a probabilistic data structure with answers membership with no false negatives

Examples: Bloom, cuckoo, quotient

Fixing Lookups in Size-Tiered B^ϵ -Trees

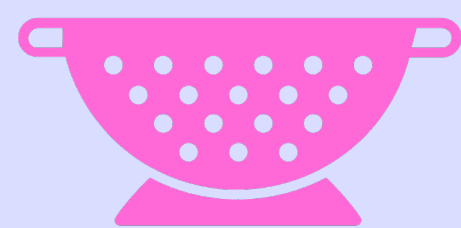
The problem is that each node has multiple branches

Idea: use filters to avoid searching them

Query(64)



Now a lookup will only search those branches which contain the key (plus rare false positives)



A filter is a probabilistic data structure with answers membership with no false negatives

Examples: Bloom, cuckoo, quotient

Fixing Lookups in Size-Tiered B^ϵ -Trees

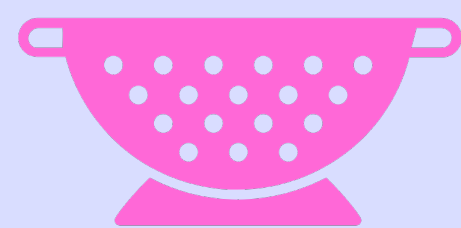
The problem is that each node has multiple branches

Idea: use filters to avoid searching them

Query(64)



Now a lookup will only search those branches which contain the key (plus rare false positives)



A filter is a probabilistic data structure with answers membership with no false negatives

Examples: Bloom, cuckoo, quotient

Fixing Lookups in Size-Tiered B^ϵ -Trees

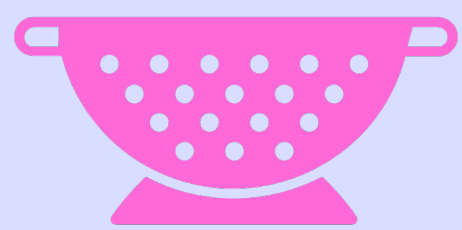
The problem is that each node has multiple branches

Idea: use filters to avoid searching them

Query(64)



Now a lookup will only search those branches which contain the key (plus rare false positives)



A filter is a probabilistic data structure with answers membership with no false negatives

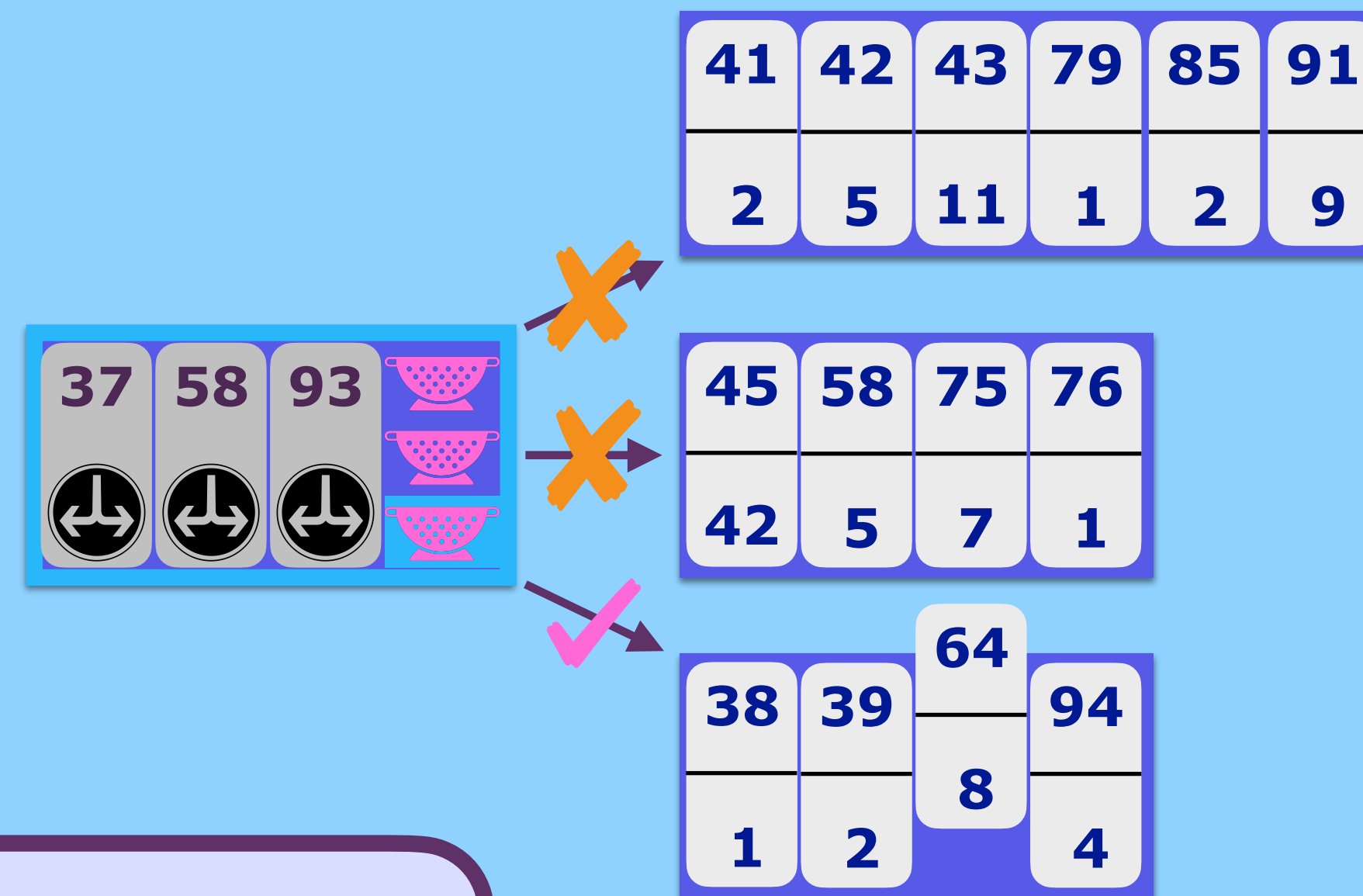
Examples: Bloom, cuckoo, quotient

Fixing Lookups in Size-Tiered B^ϵ -Trees

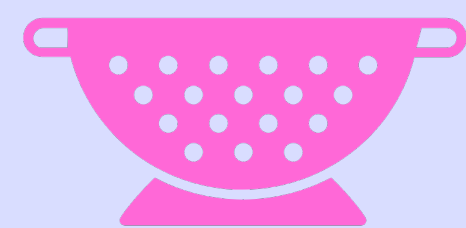
Query(64) \rightarrow 8

The problem is that each node has multiple branches

Idea: use filters to avoid searching them



Now a lookup will only search those branches which contain the key (plus rare false positives)



A filter is a probabilistic data structure with answers membership with no false negatives

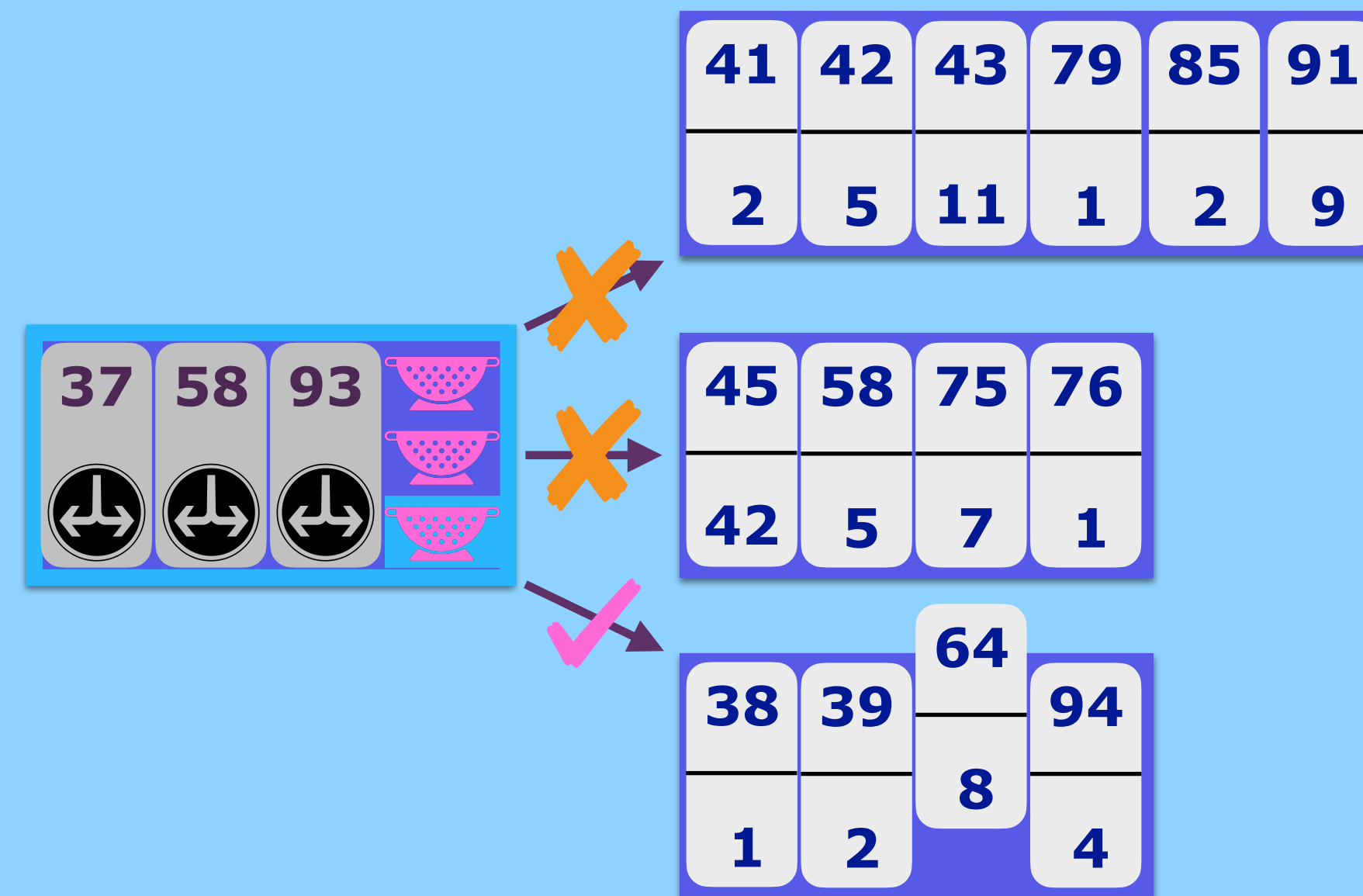
Examples: Bloom, cuckoo, quotient

Fixing Lookups in Size-Tiered B^ϵ -Trees

Query(64) \rightarrow 8

The problem is that each node has multiple branches

Idea: use filters to avoid searching them



Now a lookup will only search those branches which contain the key (plus rare false positives)

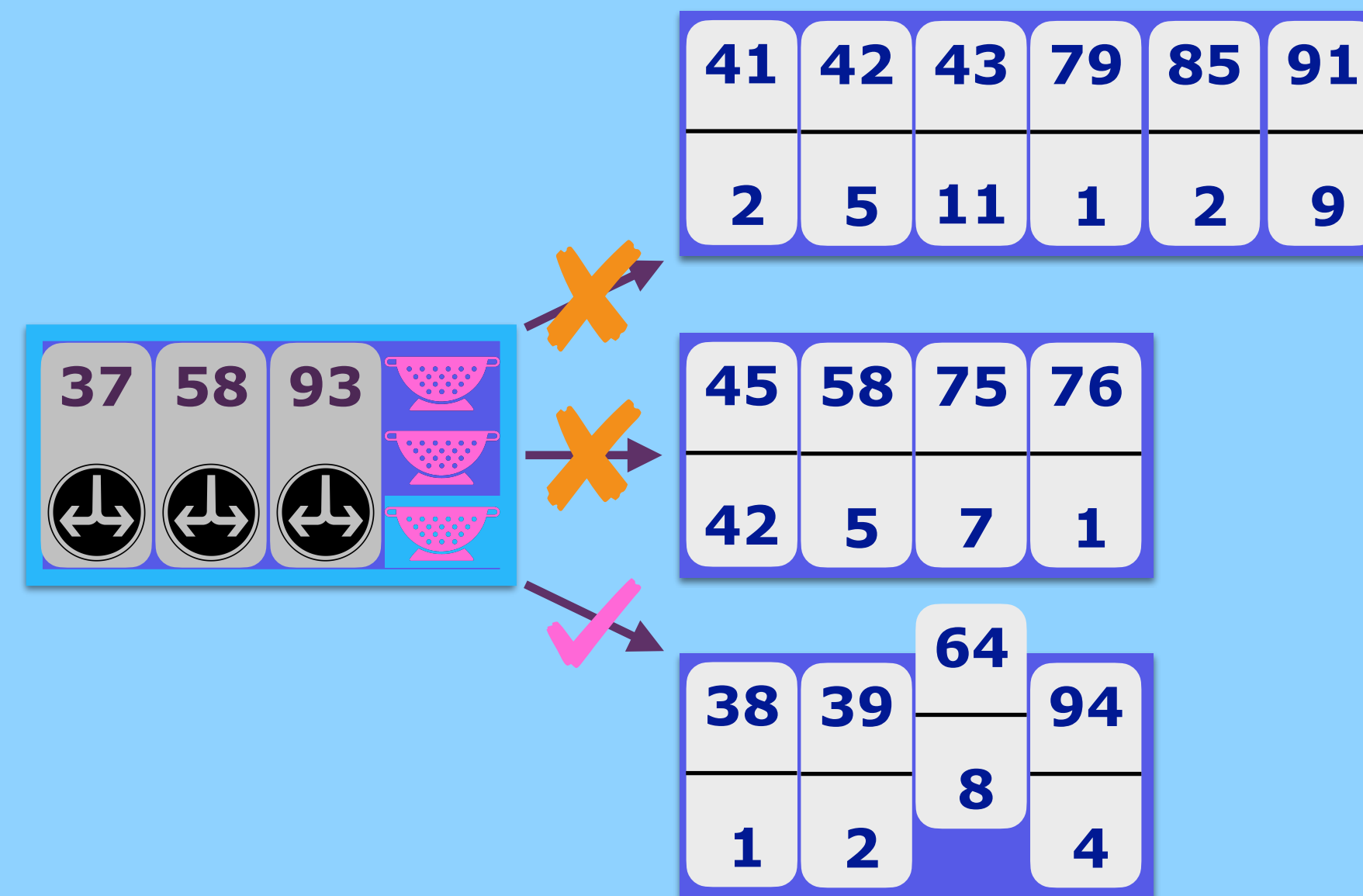
$$\text{False Positive Rate} \leq O\left(\frac{\epsilon}{B^\epsilon \log_B N}\right)$$

Fixing Lookups in Size-Tiered B^ϵ -Trees

Query(64) \rightarrow 8

The problem is that each node has multiple branches

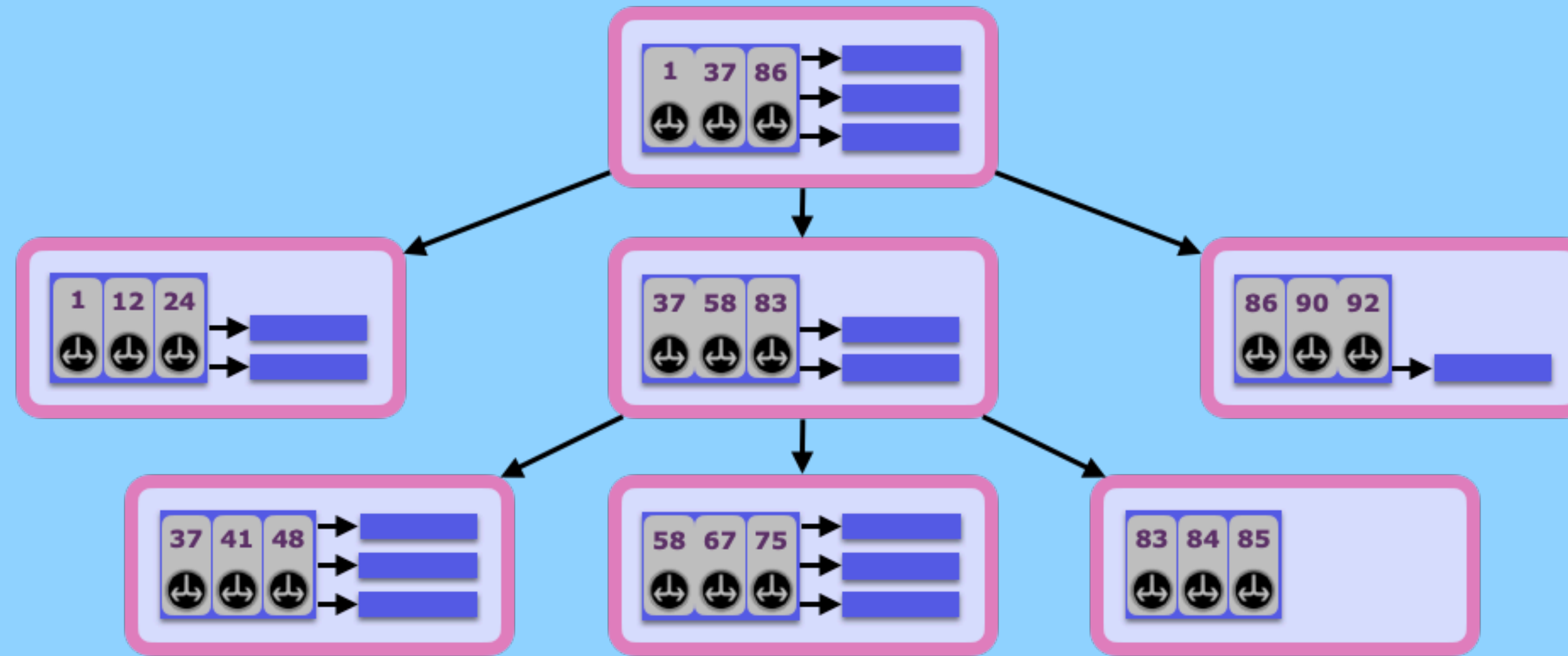
Idea: use filters to avoid searching them



Now a lookup will only search those branches which contain the key (plus rare false positives)

$$\text{False Positive Rate} \leq O\left(\frac{\epsilon}{B^\epsilon \log_B N}\right) \Rightarrow \text{Lookups in } O(1) \text{ IOs}$$

Size-Tiered B ϵ -Tree



Less compaction

Less IO

Less CPU

Low lookup cost

Scans

See the text!

Short — more expensive

Long — disk bandwidth

In this talk

Fast Storage
(NVMe)

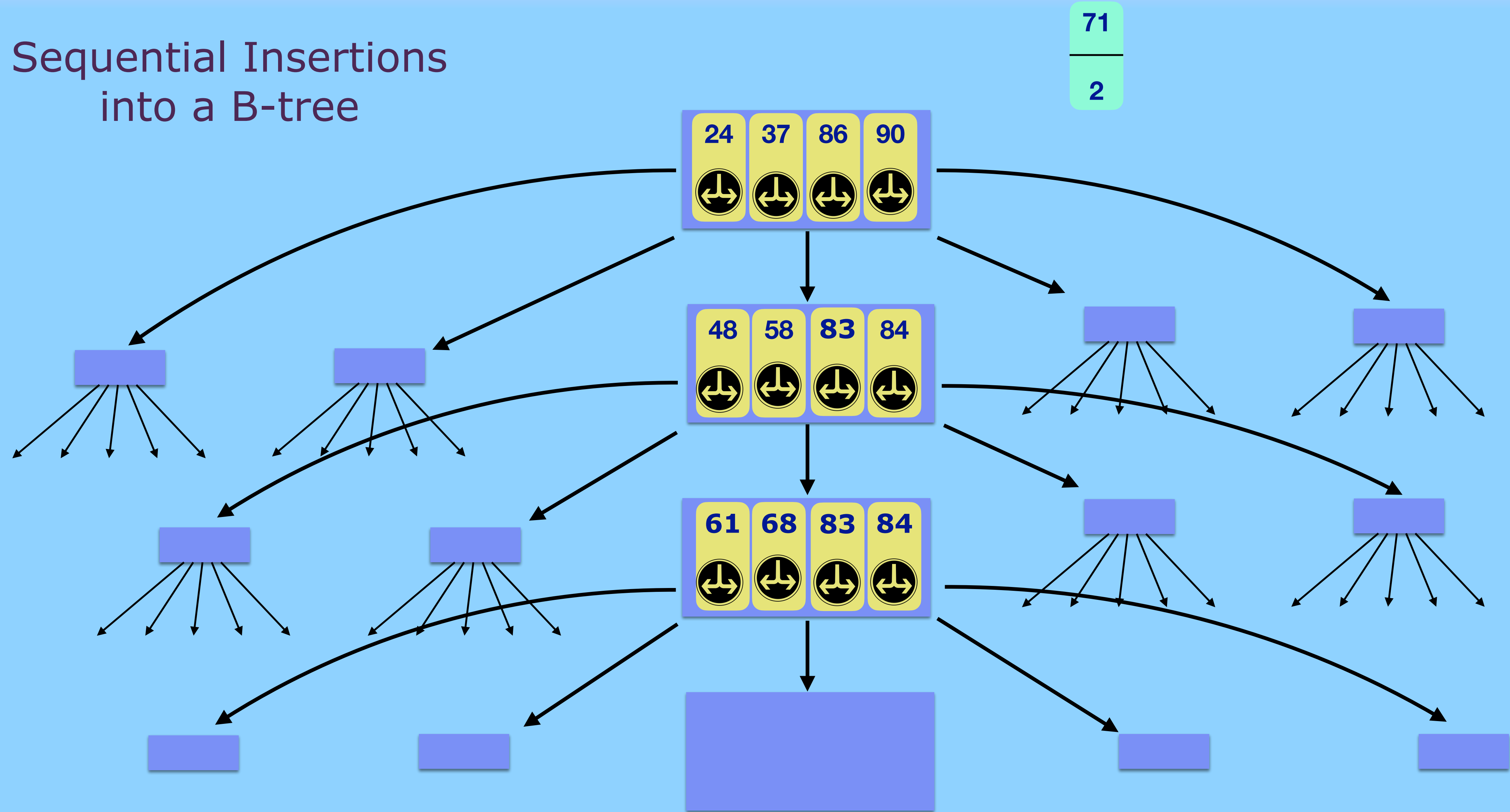
SplinterDB

Data Structures

Flush-then-Compact

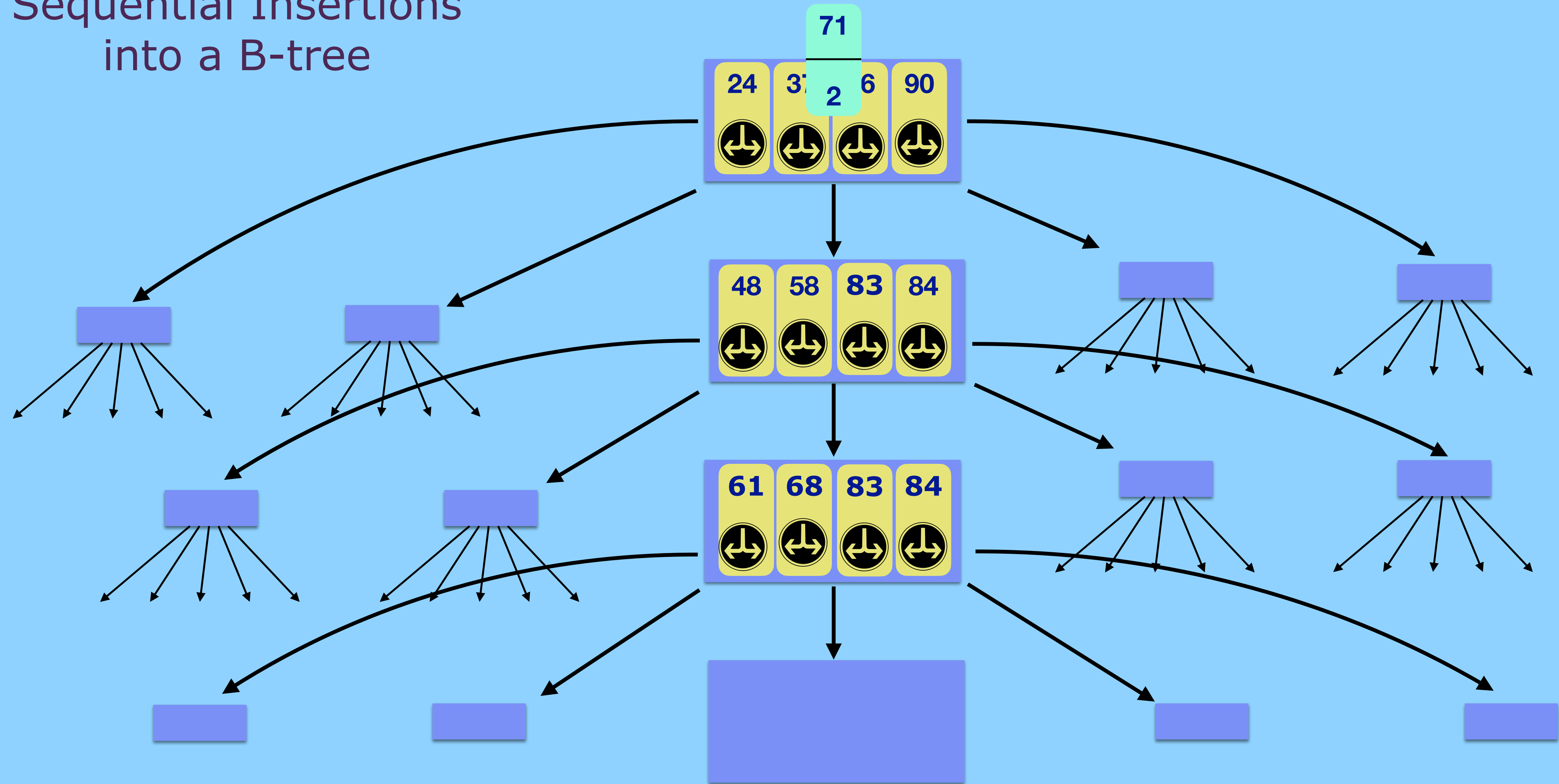
Flush-Then-Compact

Sequential Insertions
into a B-tree



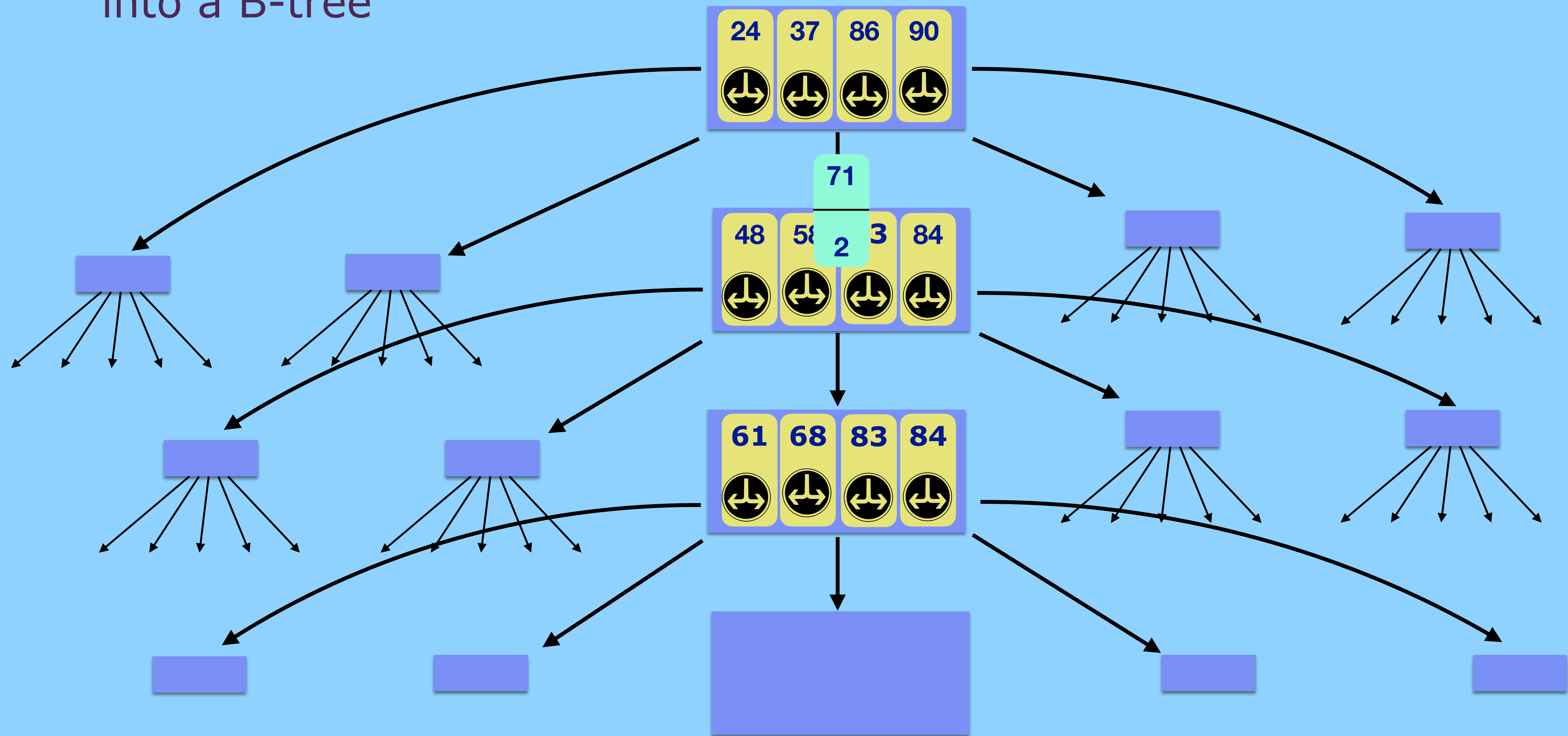
Flush-Then-Compact

Sequential Insertions
into a B-tree



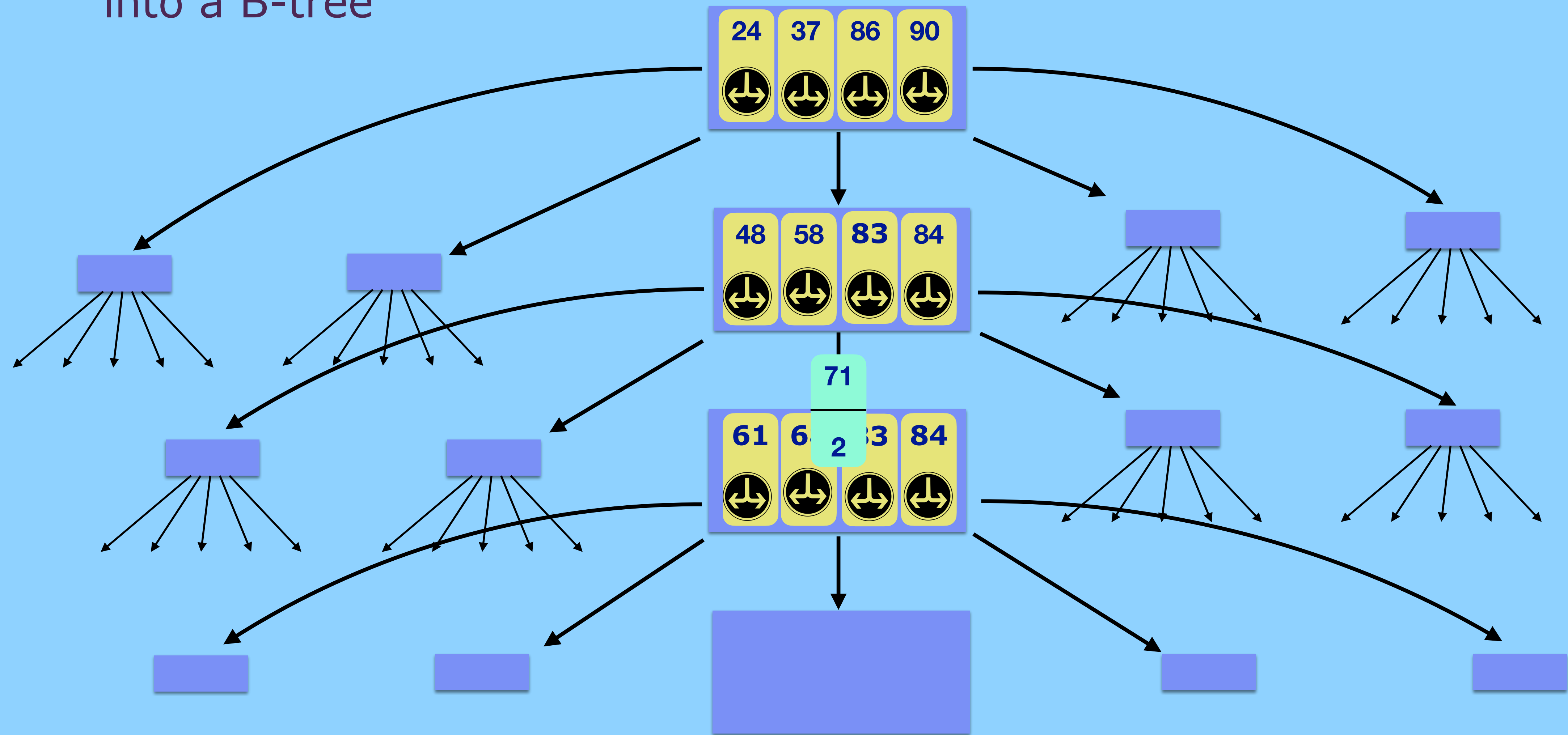
Flush-Then-Compact

Sequential Insertions
into a B-tree



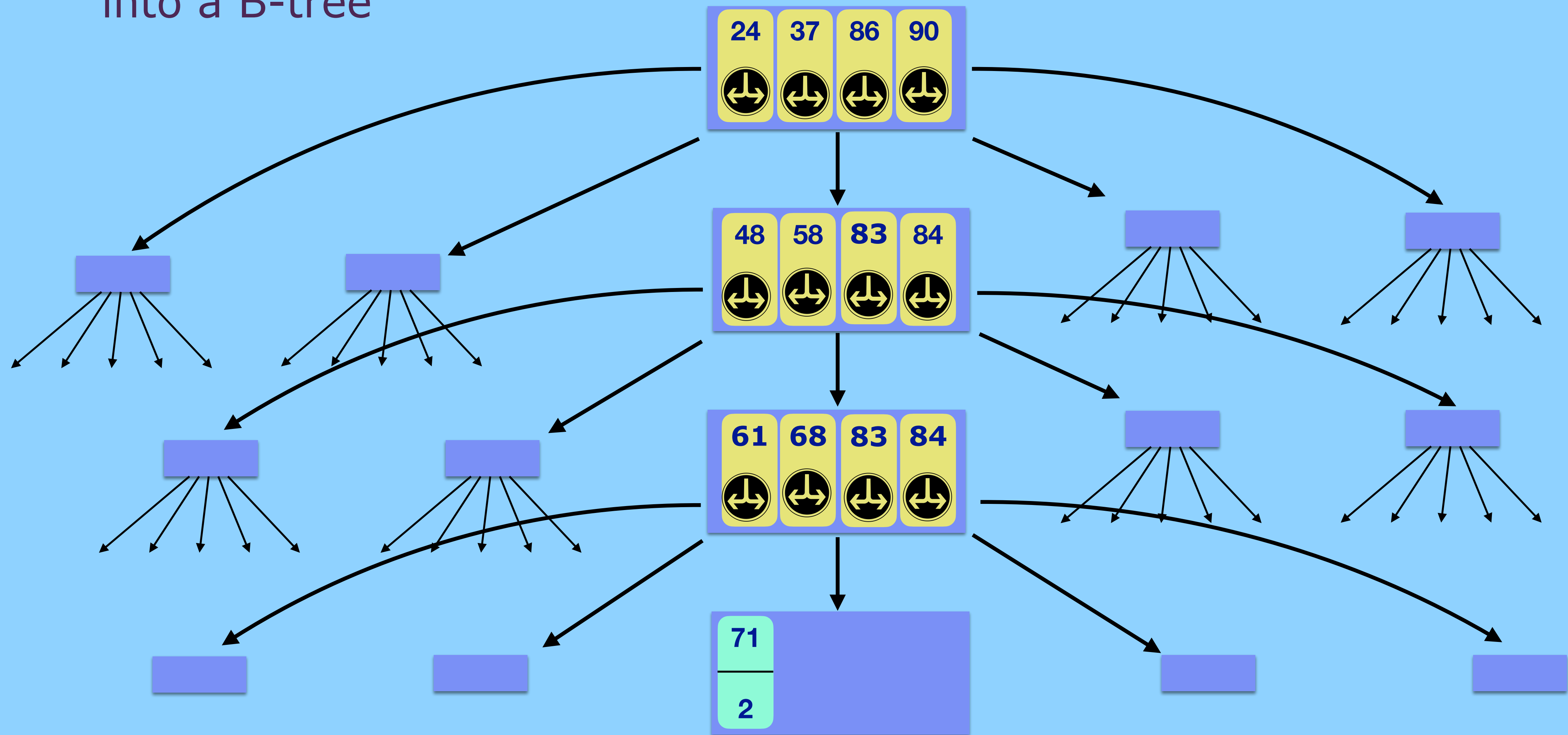
Flush-Then-Compact

Sequential Insertions
into a B-tree



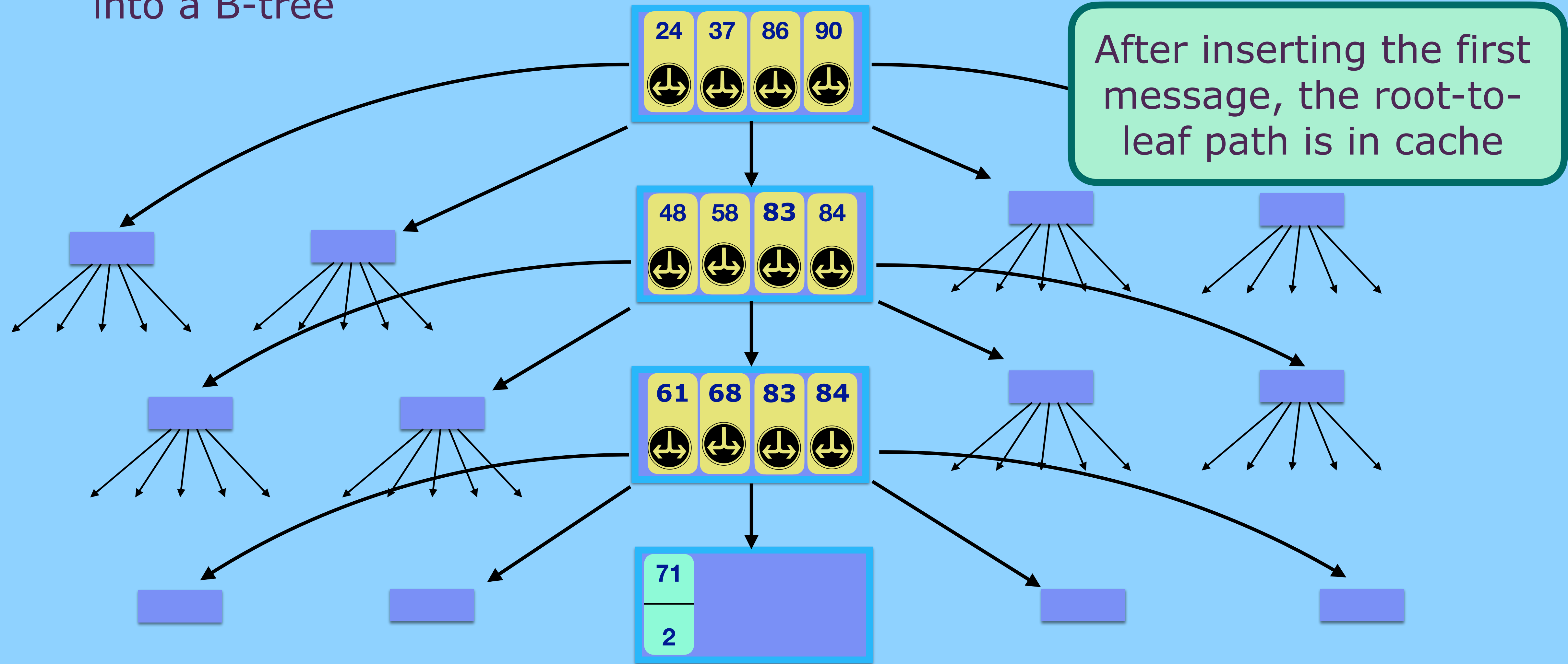
Flush-Then-Compact

Sequential Insertions
into a B-tree



Flush-Then-Compact

Sequential Insertions
into a B-tree

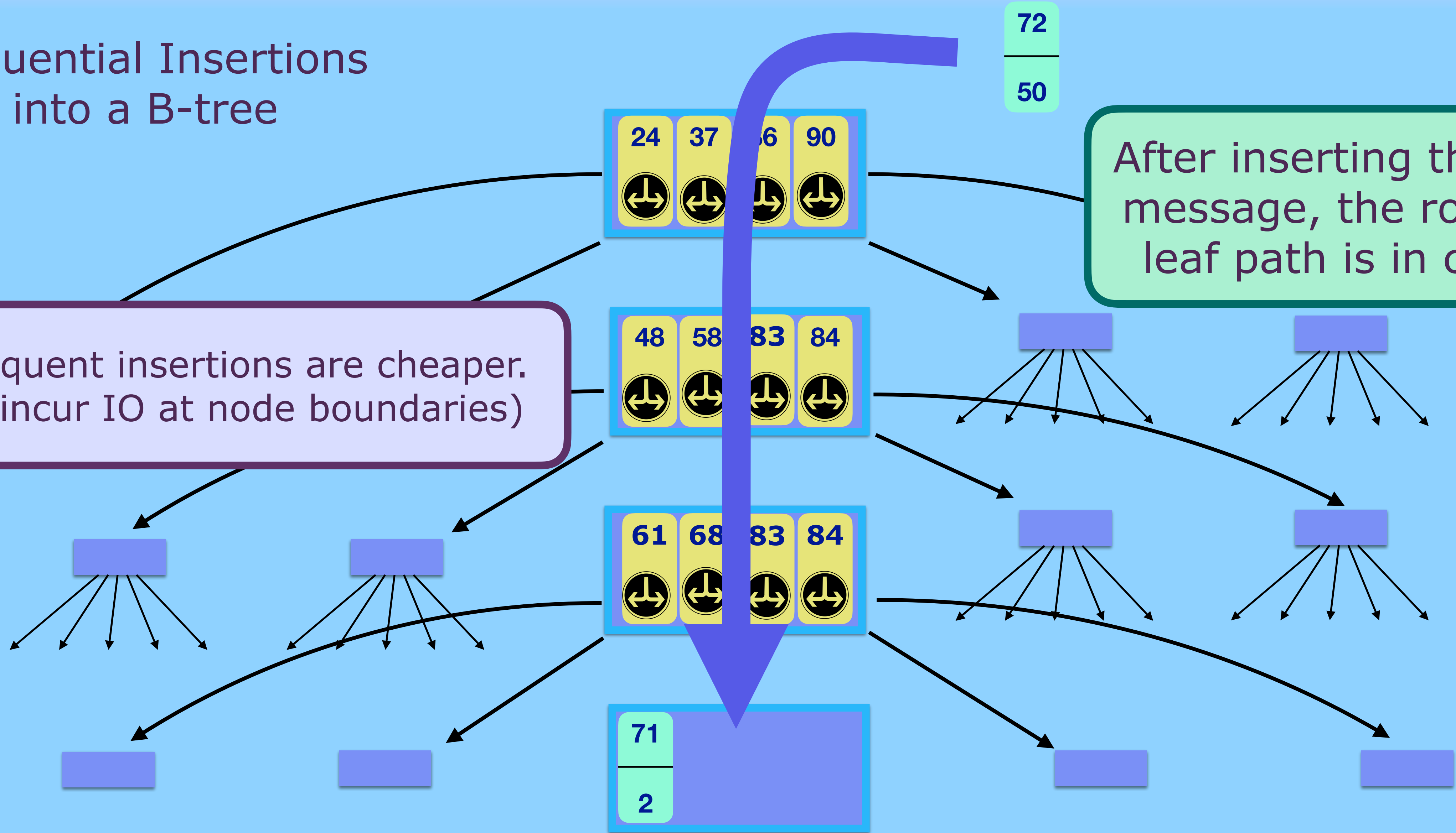


Flush-Then-Compact

Sequential Insertions
into a B-tree

Subsequent insertions are cheaper.
(only incur IO at node boundaries)

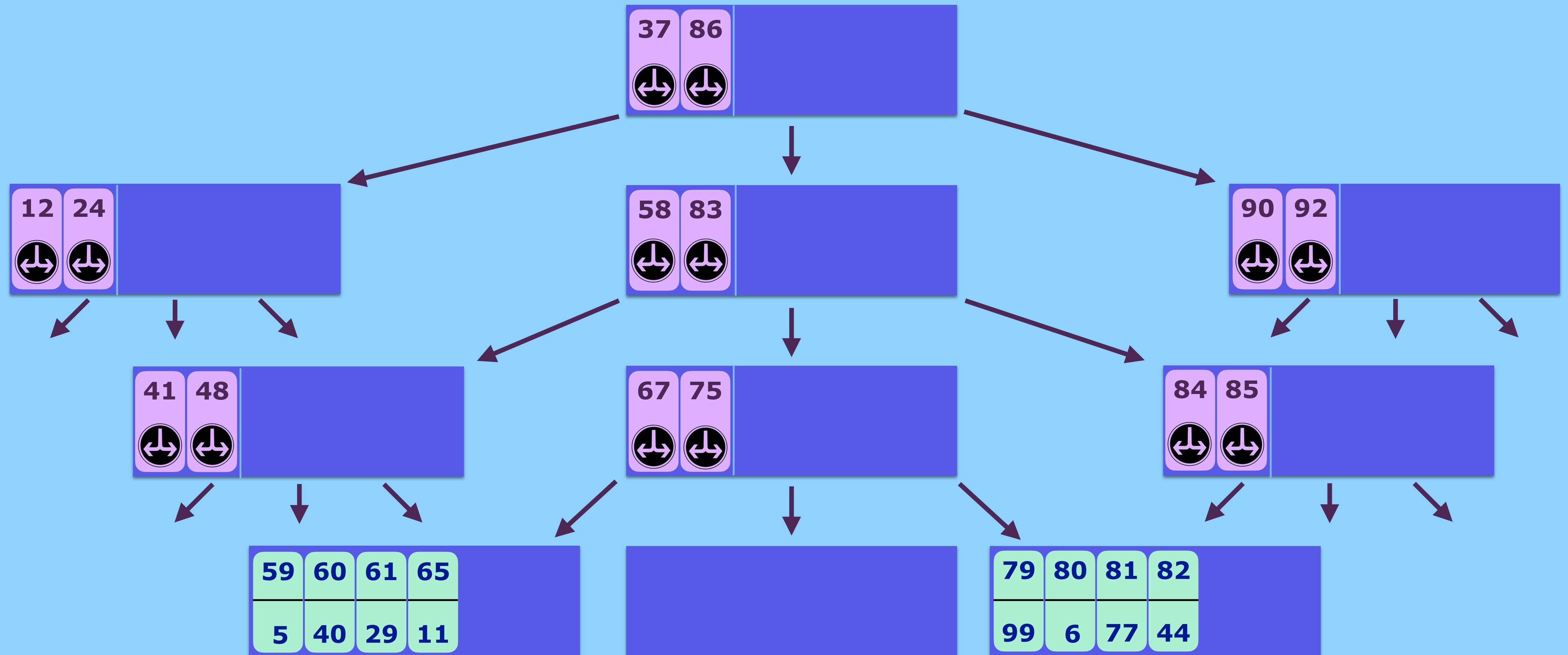
After inserting the first message, the root-to-leaf path is in cache



Flush-Then-Compact

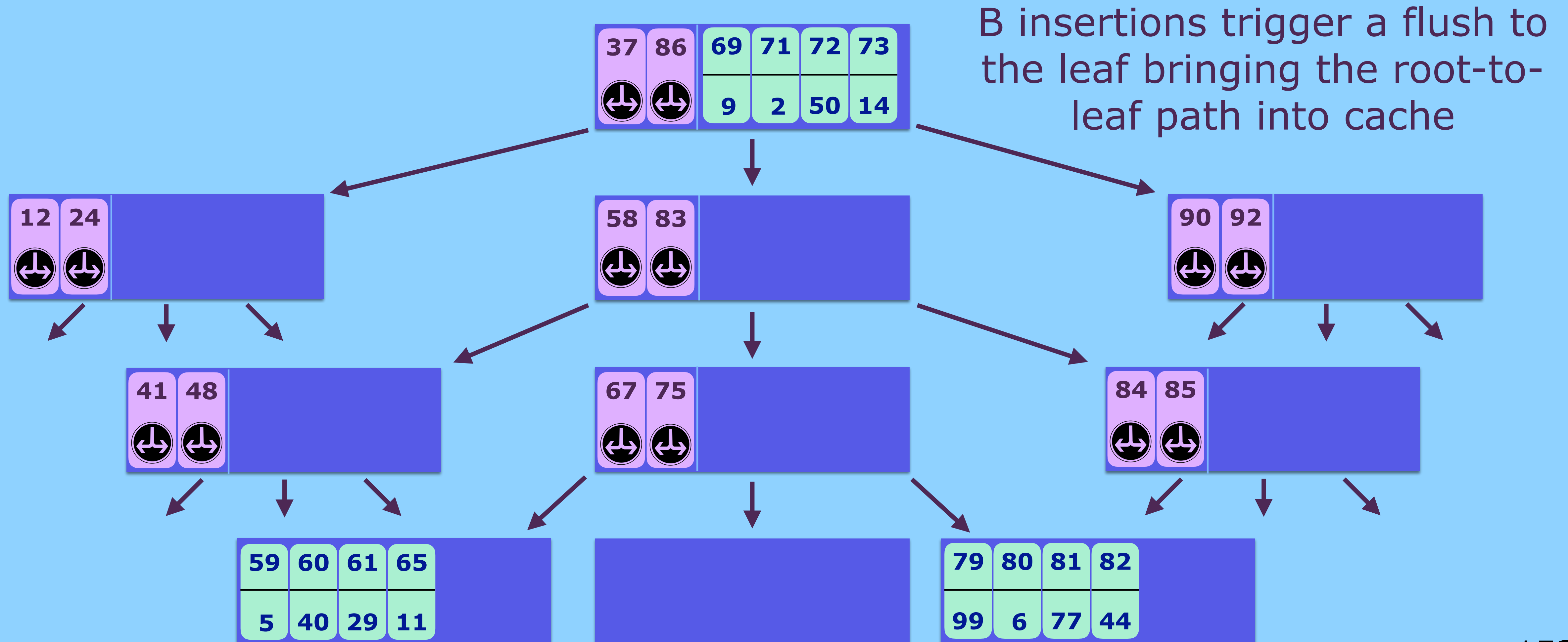
Sequential Insertions
into a B^ϵ -tree

69	71	72	73
9	2	50	14



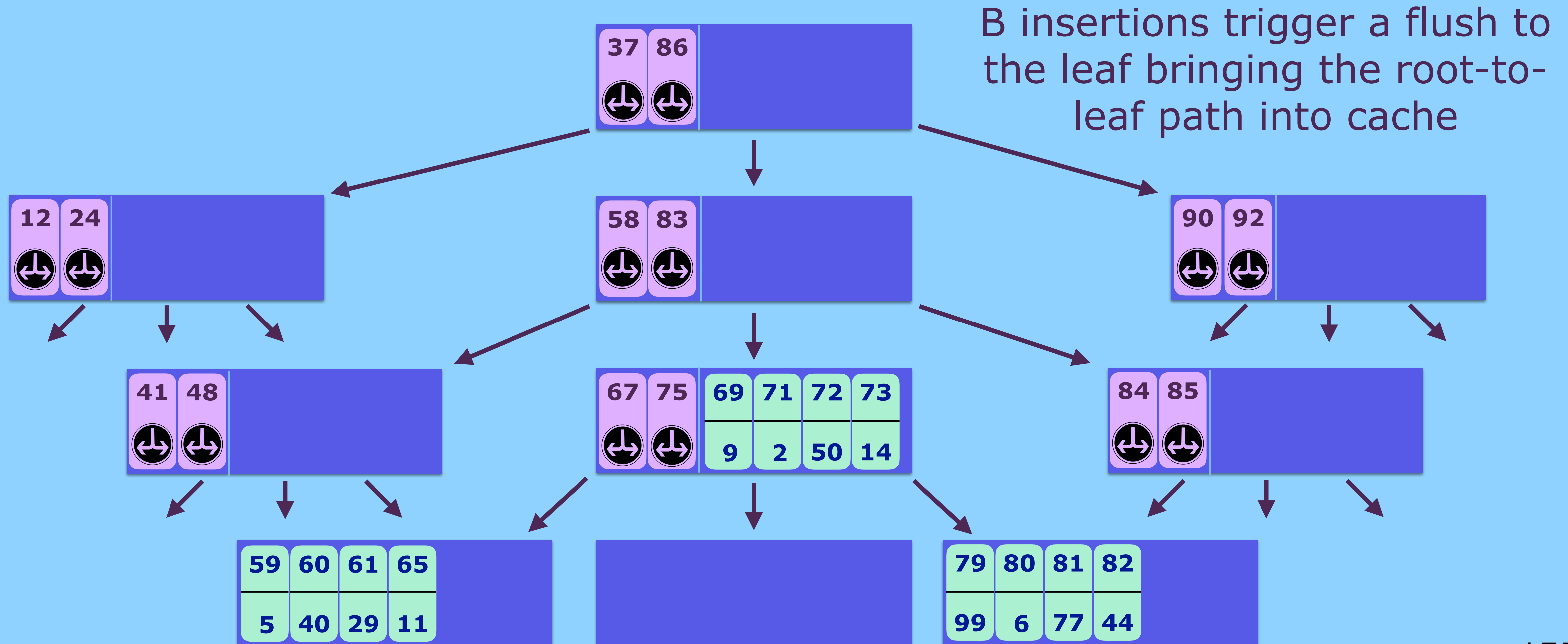
Flush-Then-Compact

Sequential Insertions
into a B^ϵ -tree



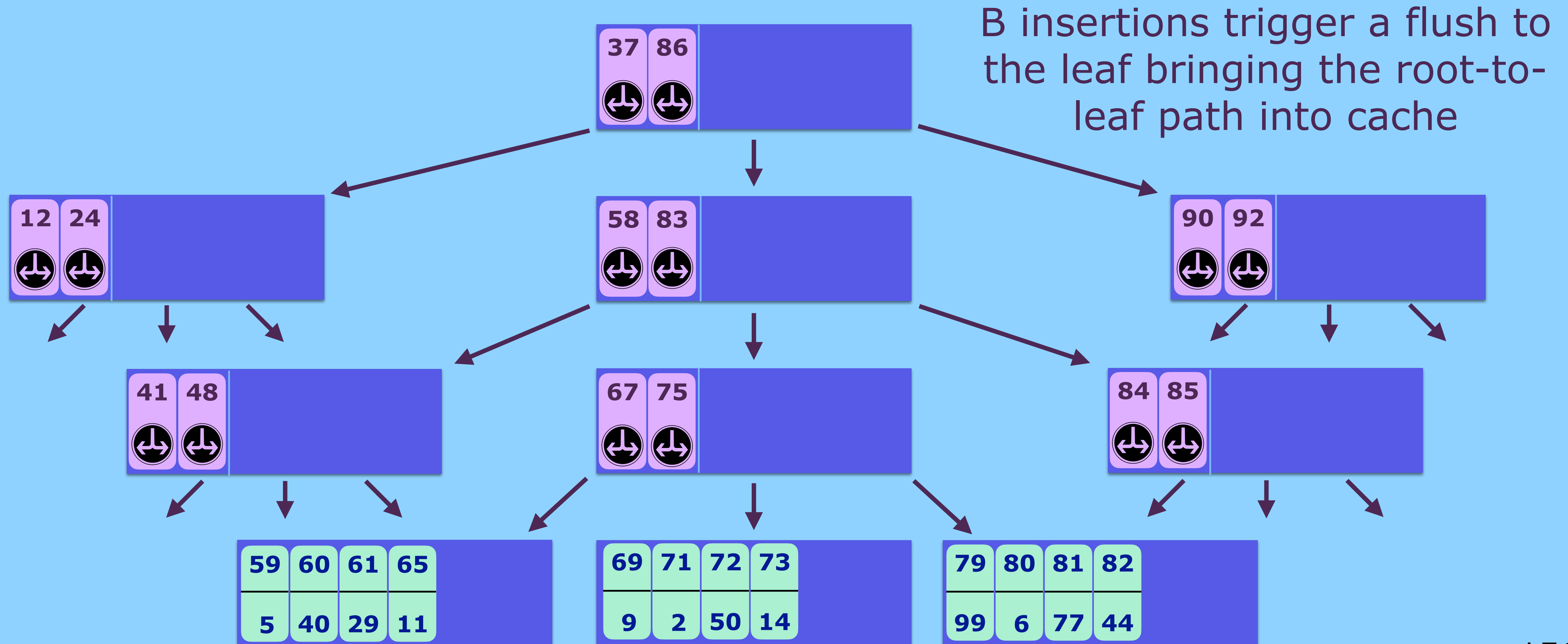
Flush-Then-Compact

Sequential Insertions
into a B^ϵ -tree



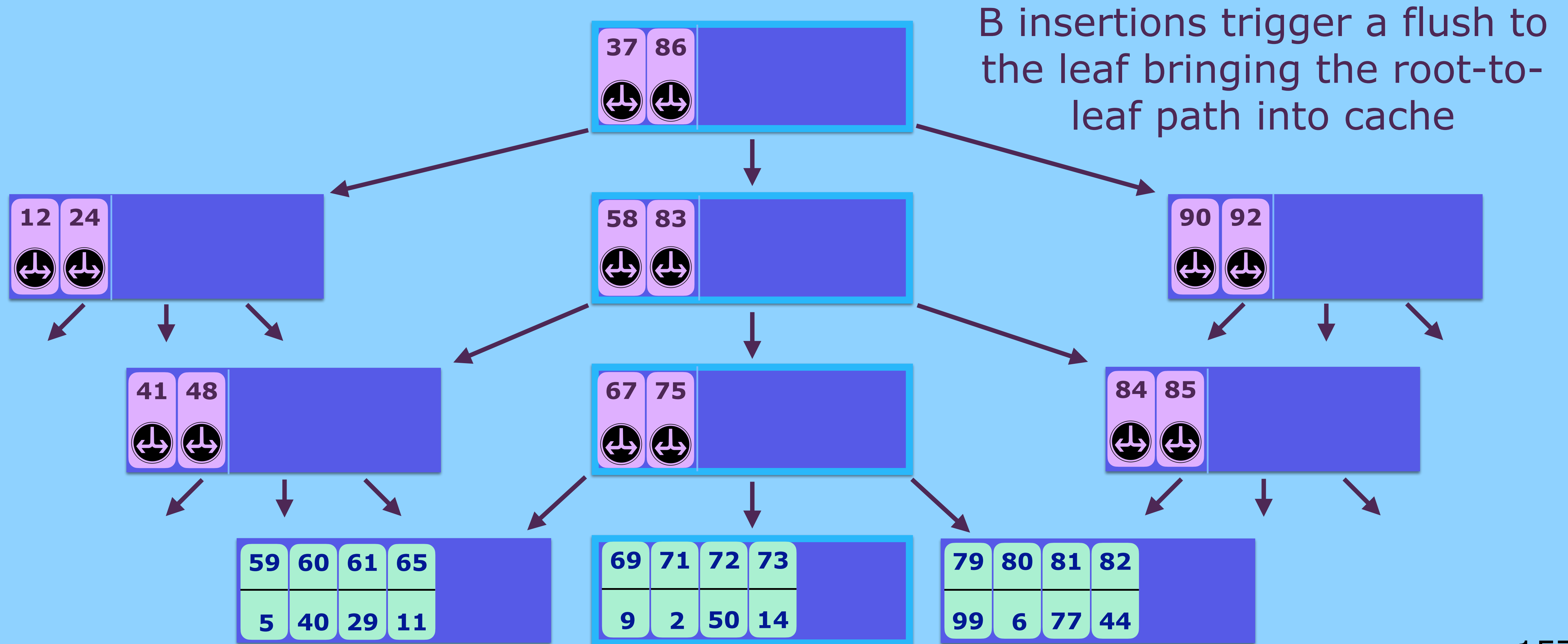
Flush-Then-Compact

Sequential Insertions
into a B^ϵ -tree



Flush-Then-Compact

Sequential Insertions
into a B^ϵ -tree



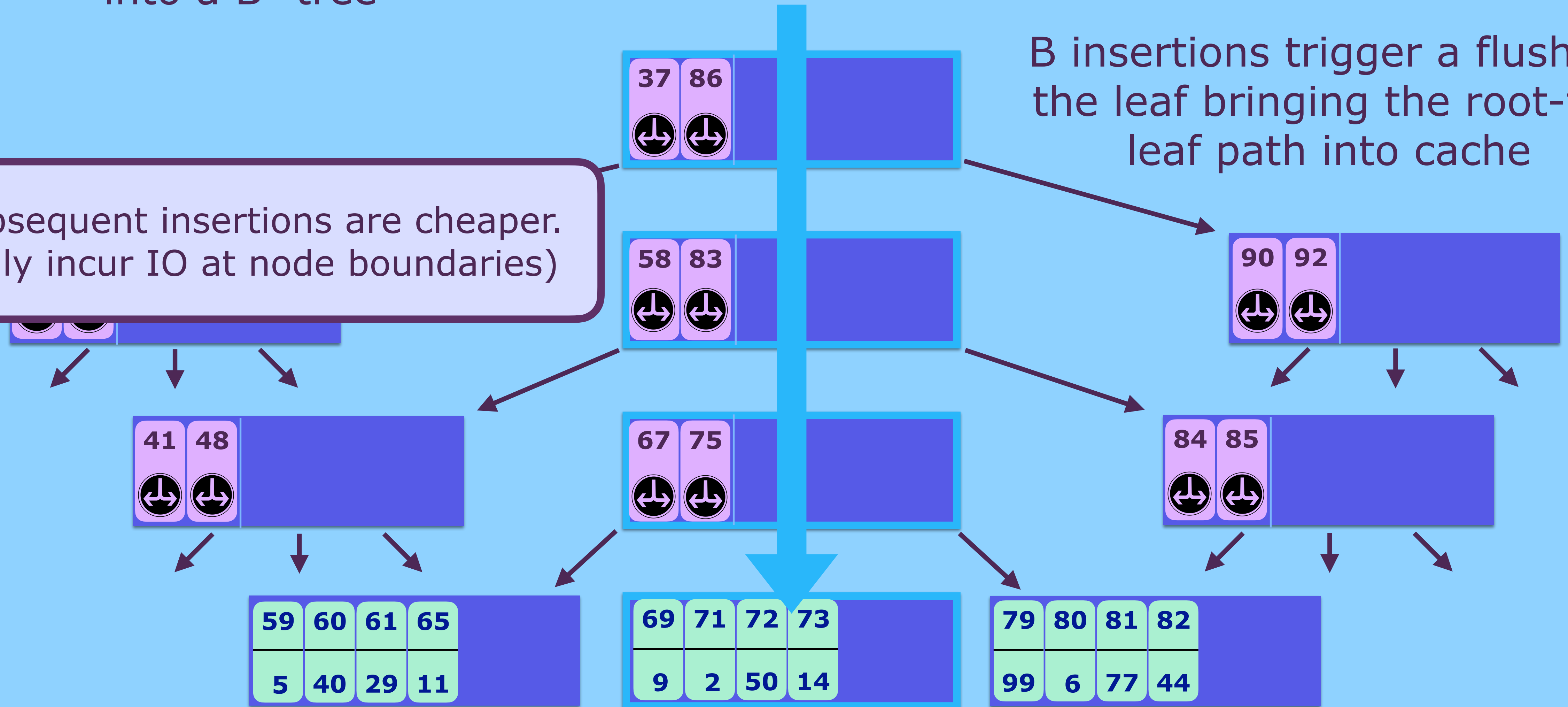
Flush-Then-Compact

Sequential Insertions
into a B^ϵ -tree

74	75	76	77
1	2	3	4

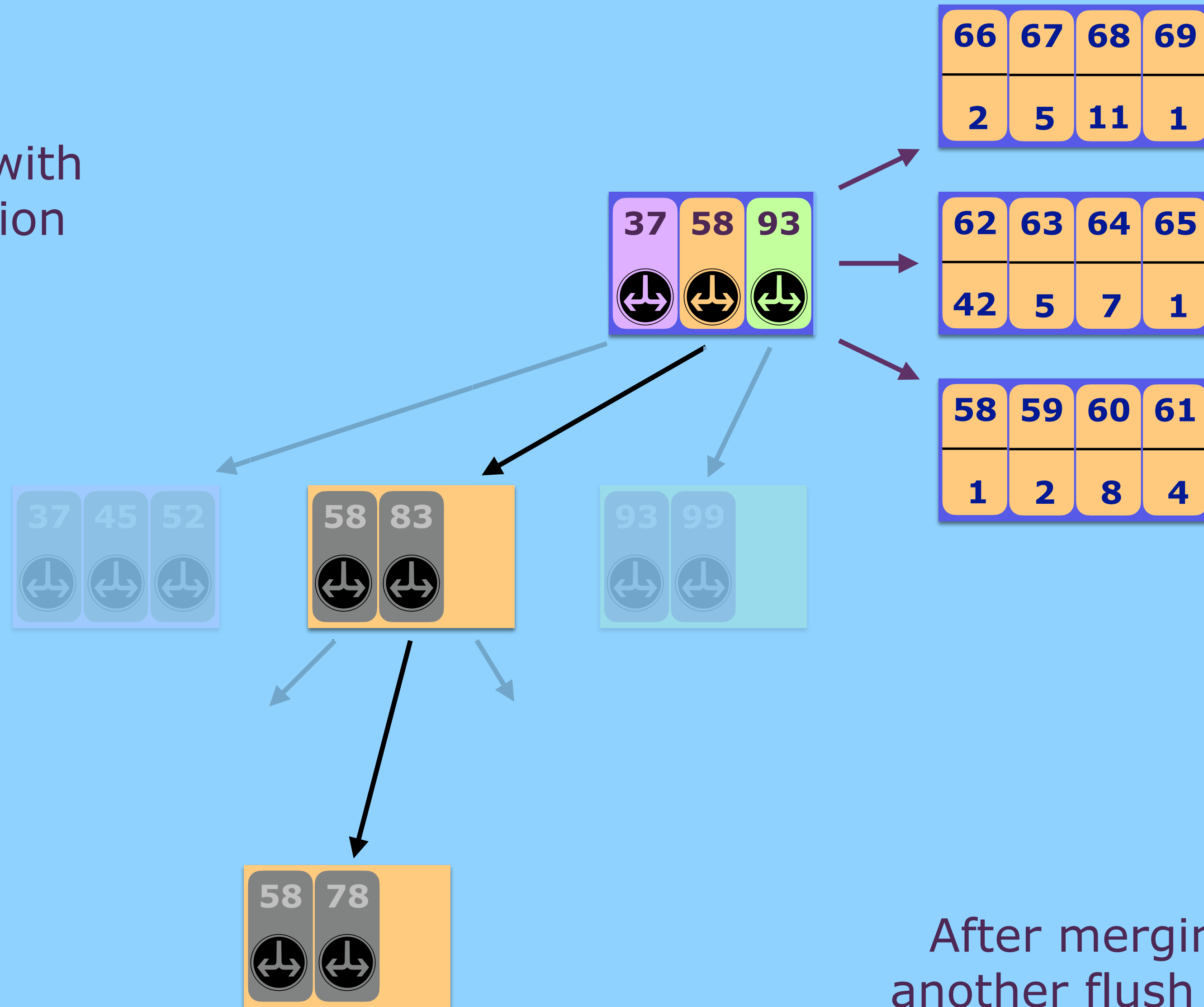
B insertions trigger a flush to the leaf bringing the root-to-leaf path into cache

Subsequent insertions are cheaper.
(only incur IO at node boundaries)



Flush-Then-Compact

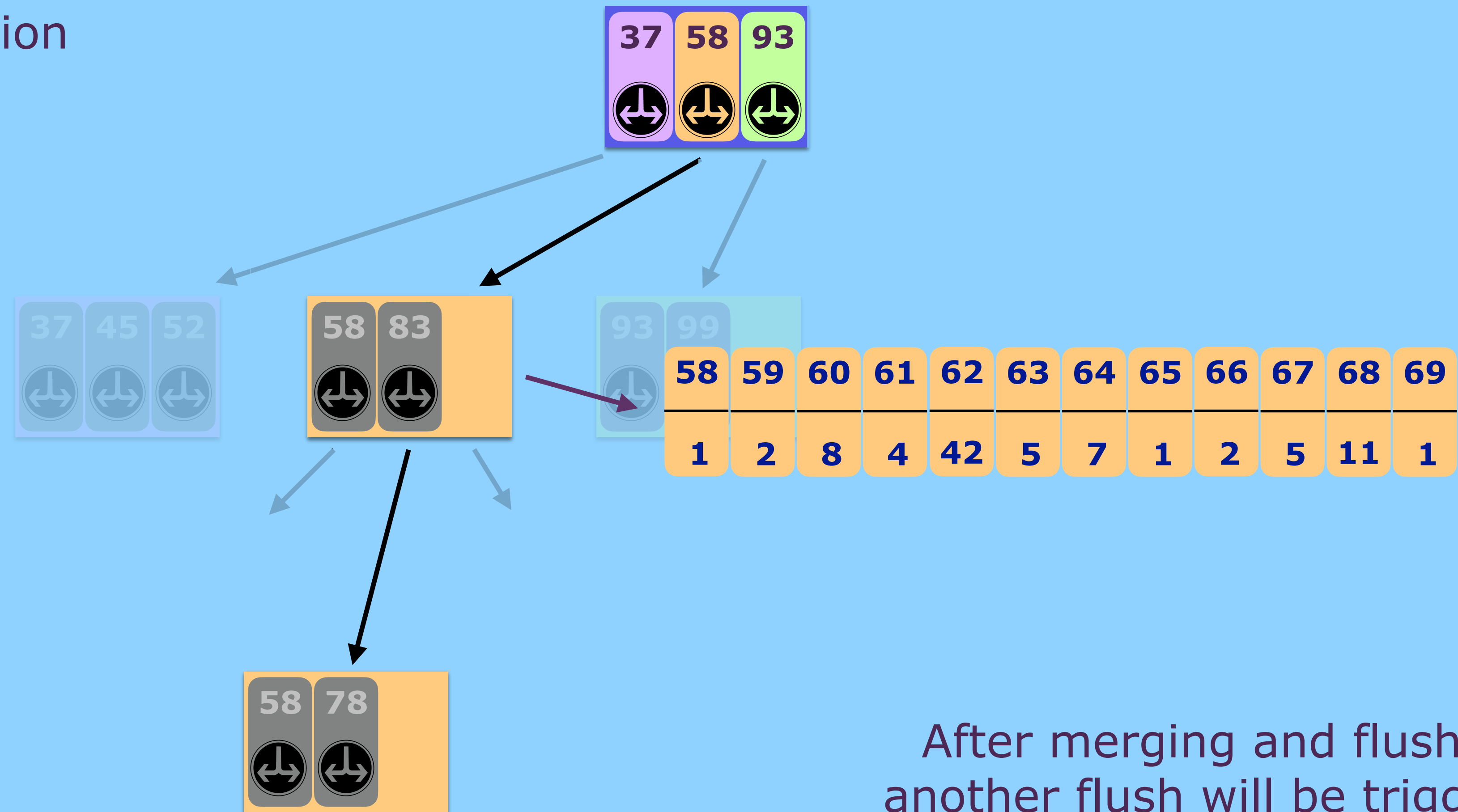
Want:
Sequential insertions with
lower work amplification



After merging and flushing
another flush will be triggered

Flush-Then-Compact

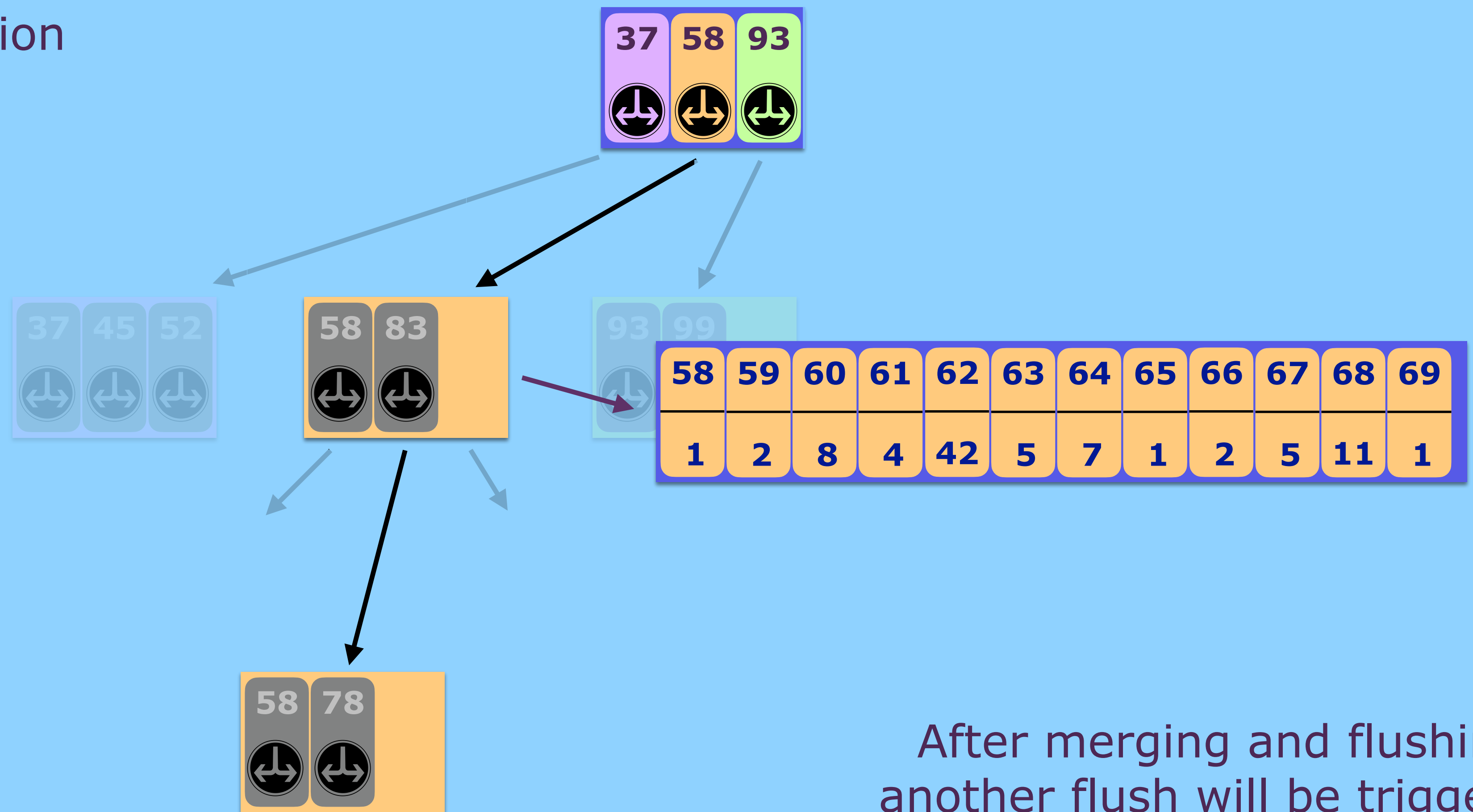
Want:
Sequential insertions with
lower work amplification



After merging and flushing
another flush will be triggered

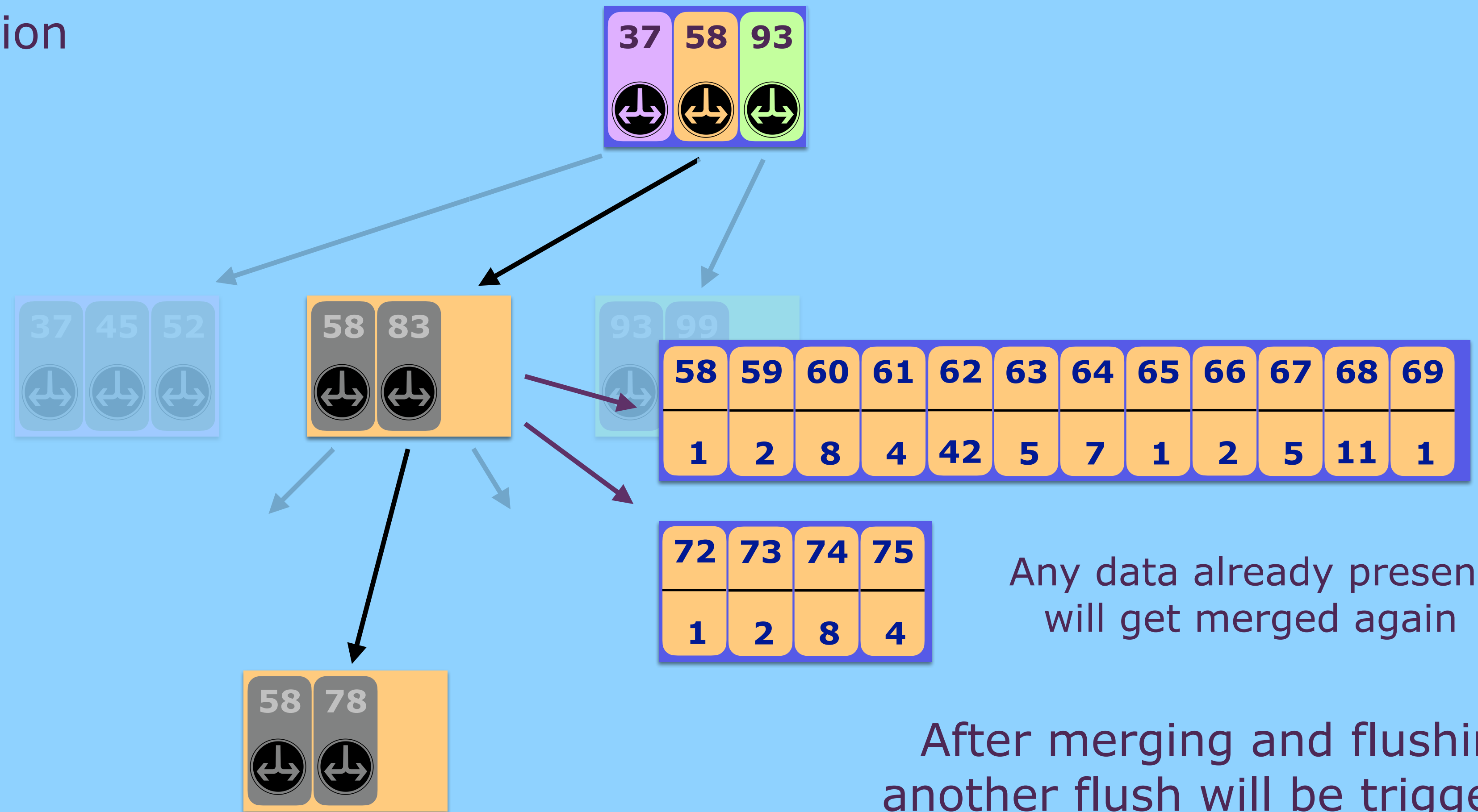
Flush-Then-Compact

Want:
Sequential insertions with
lower work amplification



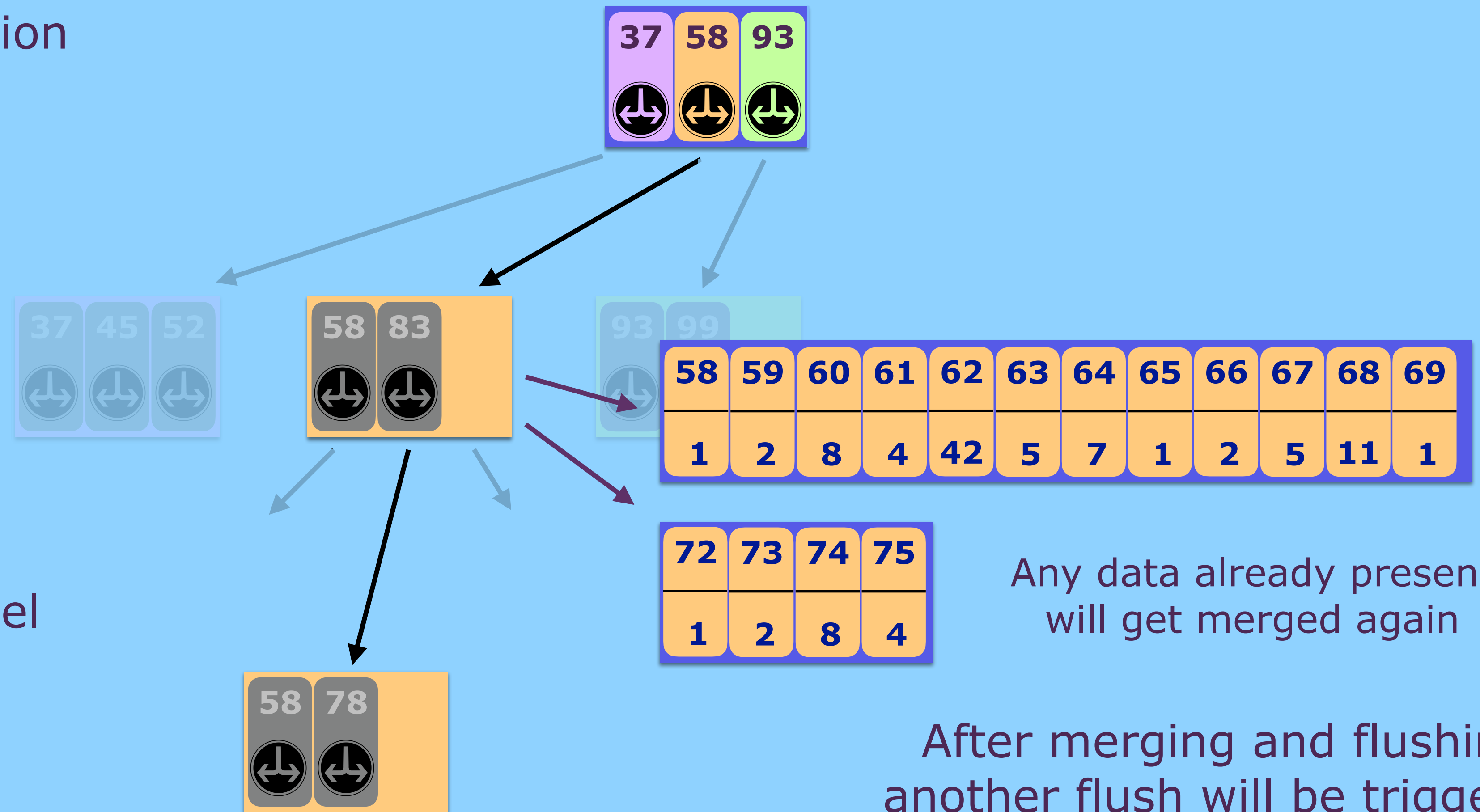
Flush-Then-Compact

Want:
Sequential insertions with
lower work amplification



Flush-Then-Compact

Want:
Sequential insertions with
lower work amplification



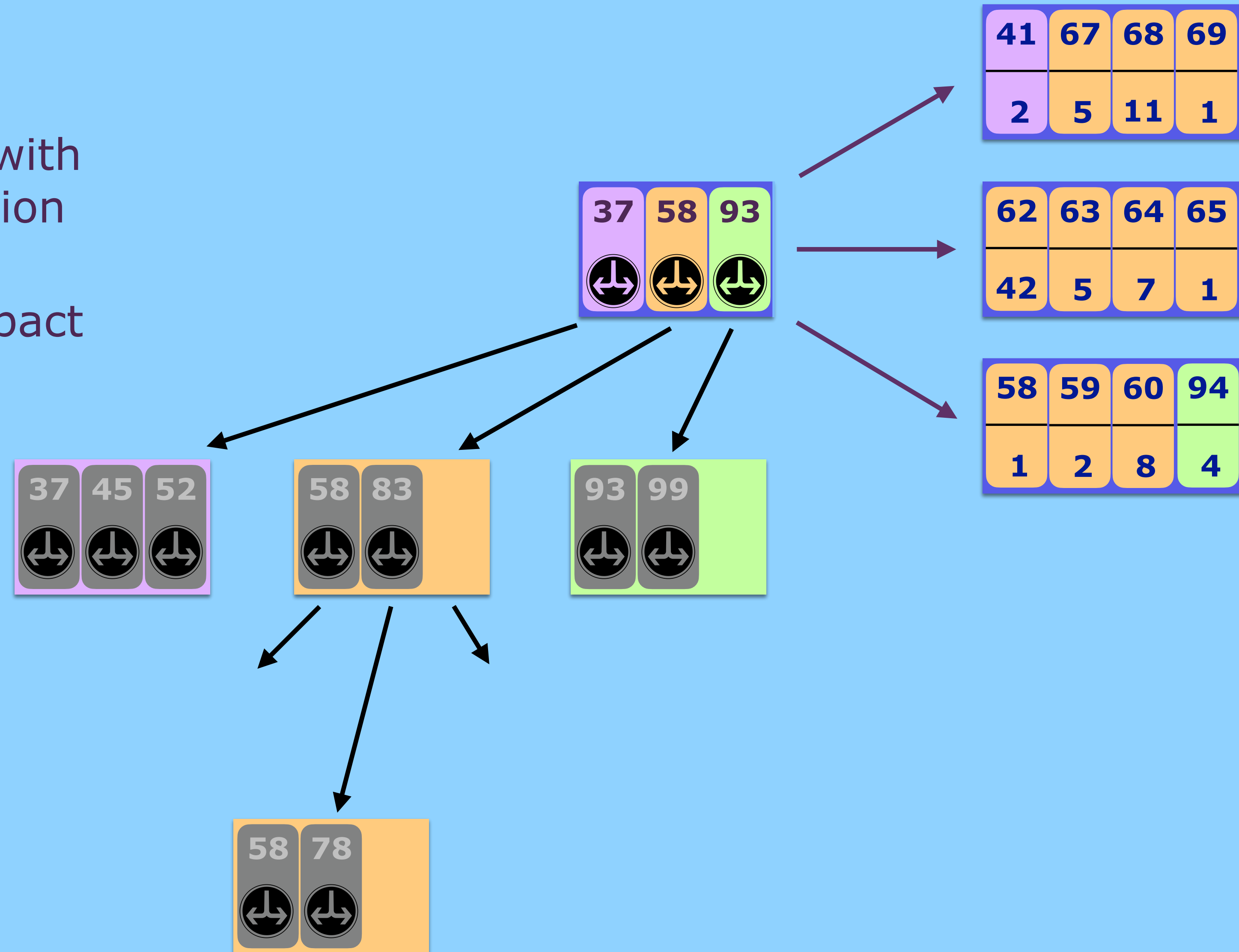
Can still end up
merging on each level

After merging and flushing
another flush will be triggered

Flush-Then-Compact

Want:
Sequential insertions with
lower work amplification

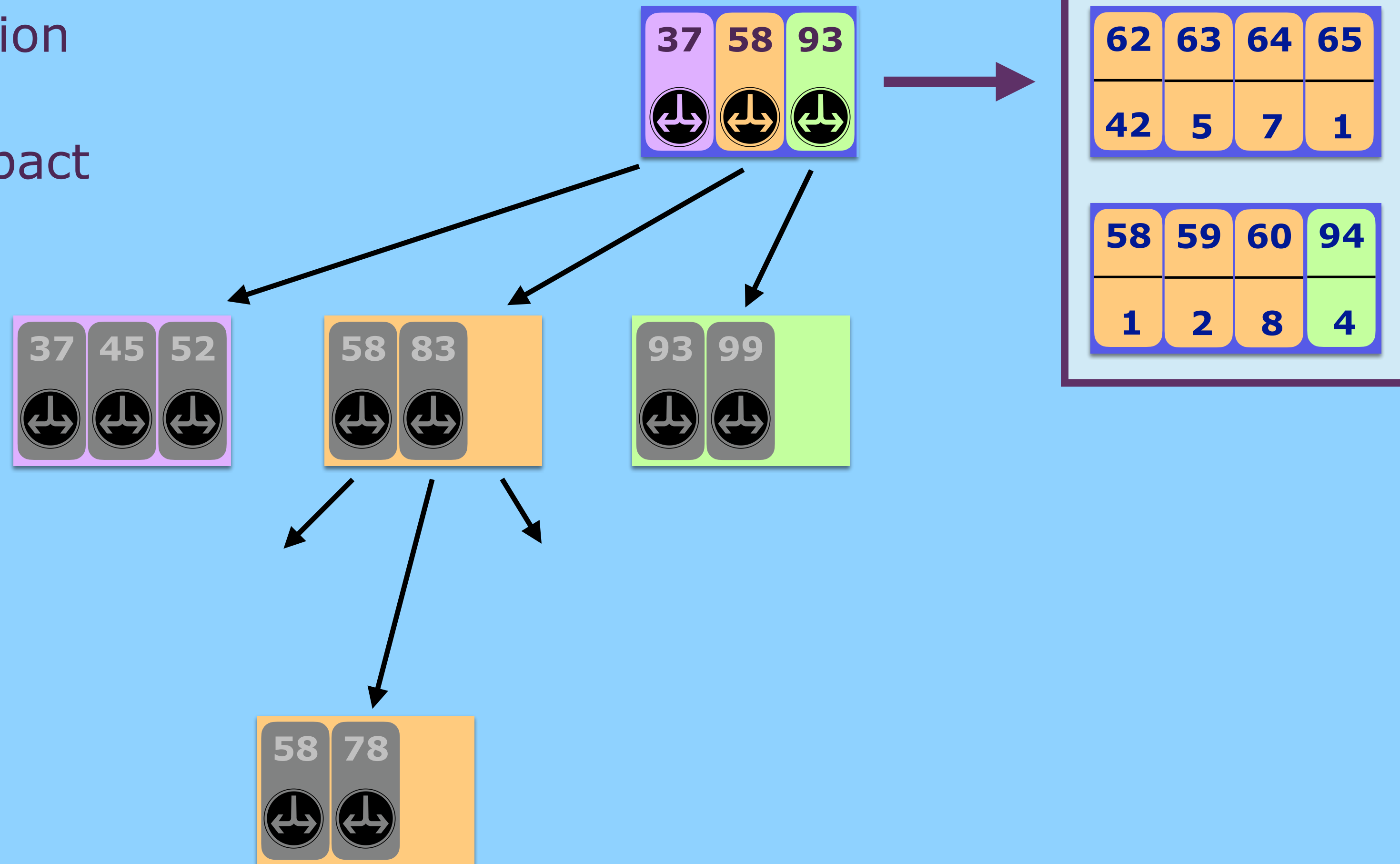
Idea: Flush-then-compact



Flush-Then-Compact

Want:
Sequential insertions with
lower work amplification

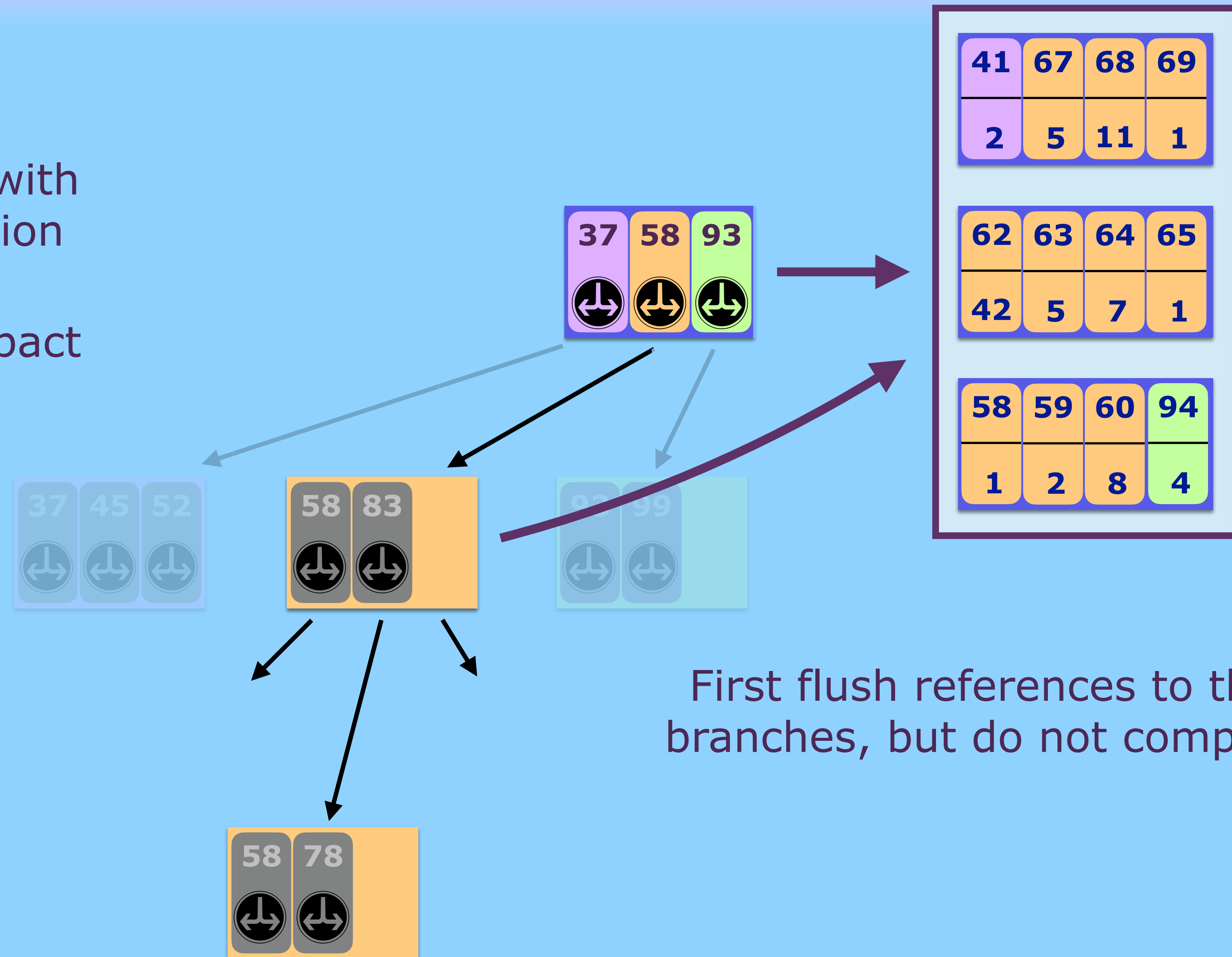
Idea: Flush-then-compact



Flush-Then-Compact

Want:
Sequential insertions with
lower work amplification

Idea: Flush-then-compact

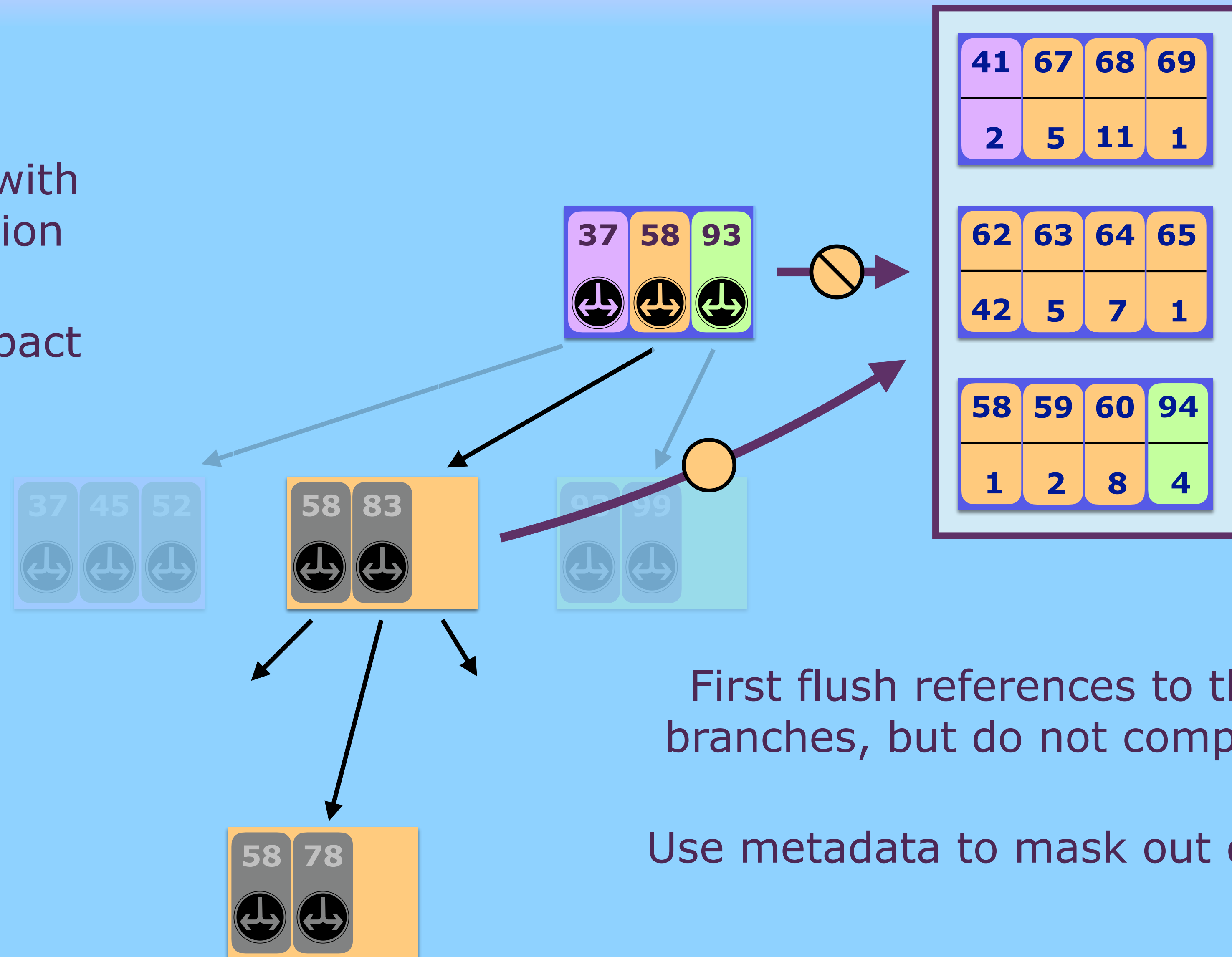


First flush references to the
branches, but do not compact

Flush-Then-Compact

Want:
Sequential insertions with
lower work amplification

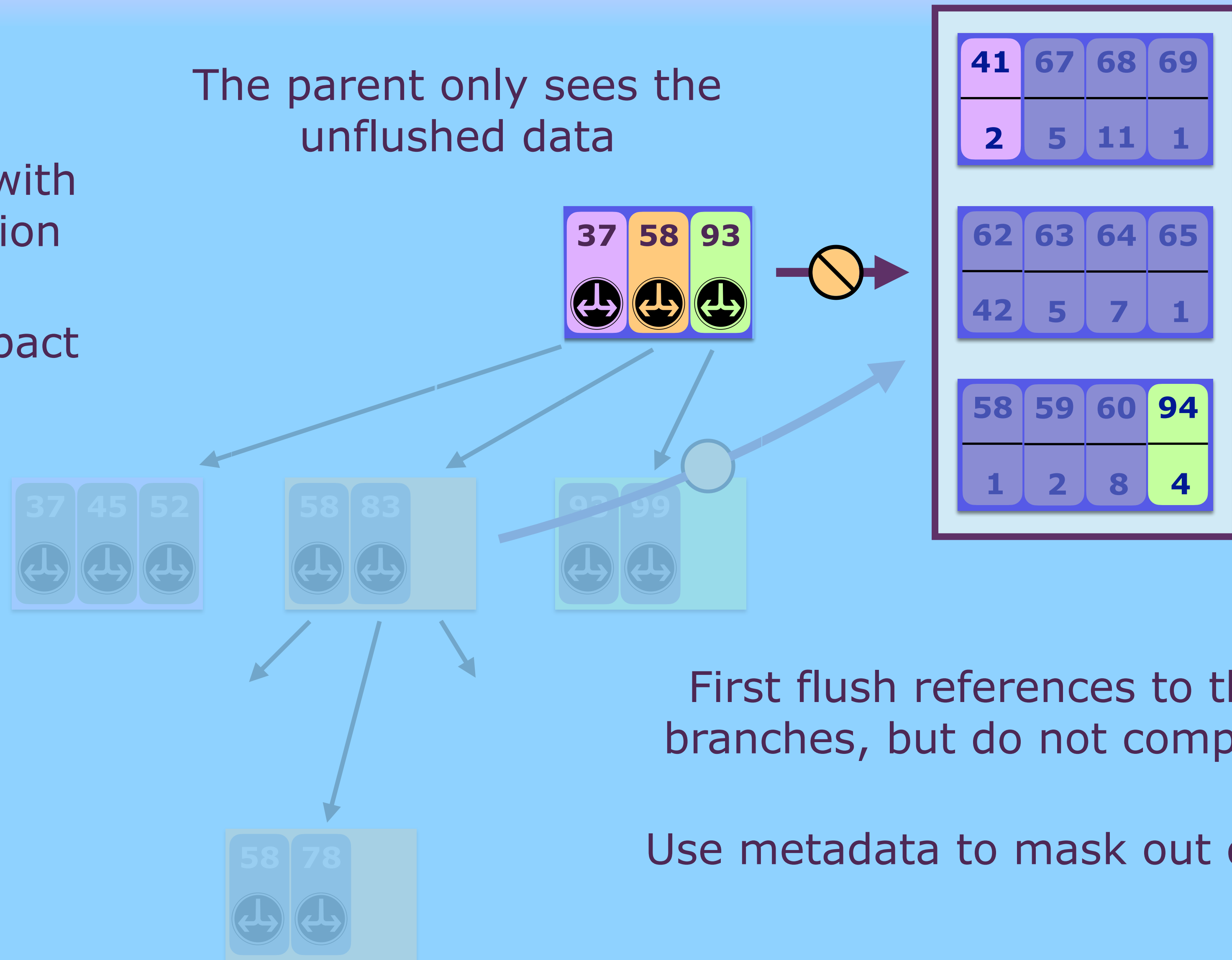
Idea: Flush-then-compact



Flush-Then-Compact

Want:
Sequential insertions with
lower work amplification

Idea: Flush-then-compact



Flush-Then-Compact

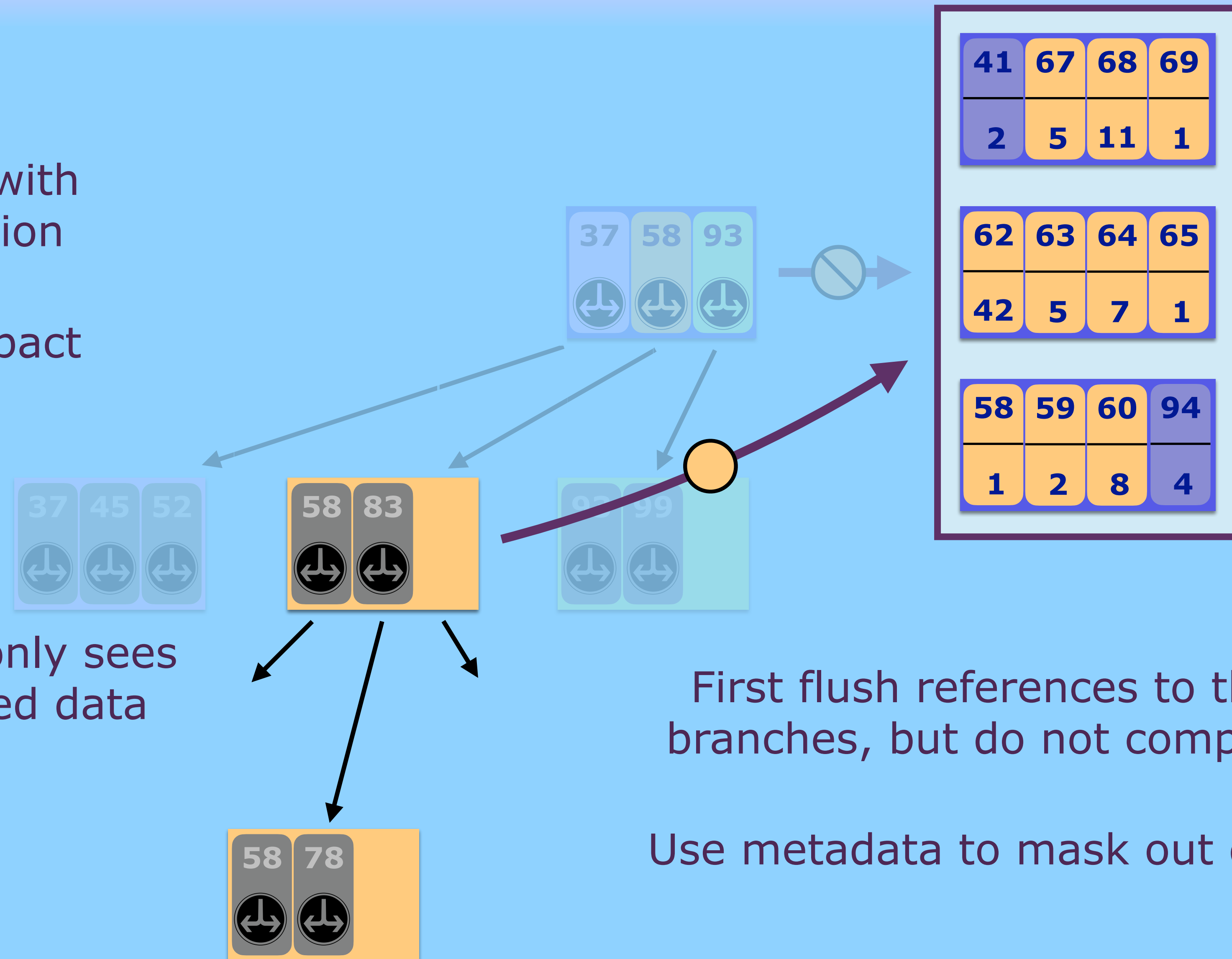
Want:
Sequential insertions with
lower work amplification

Idea: Flush-then-compact

The child only sees
the flushed data

First flush references to the
branches, but do not compact

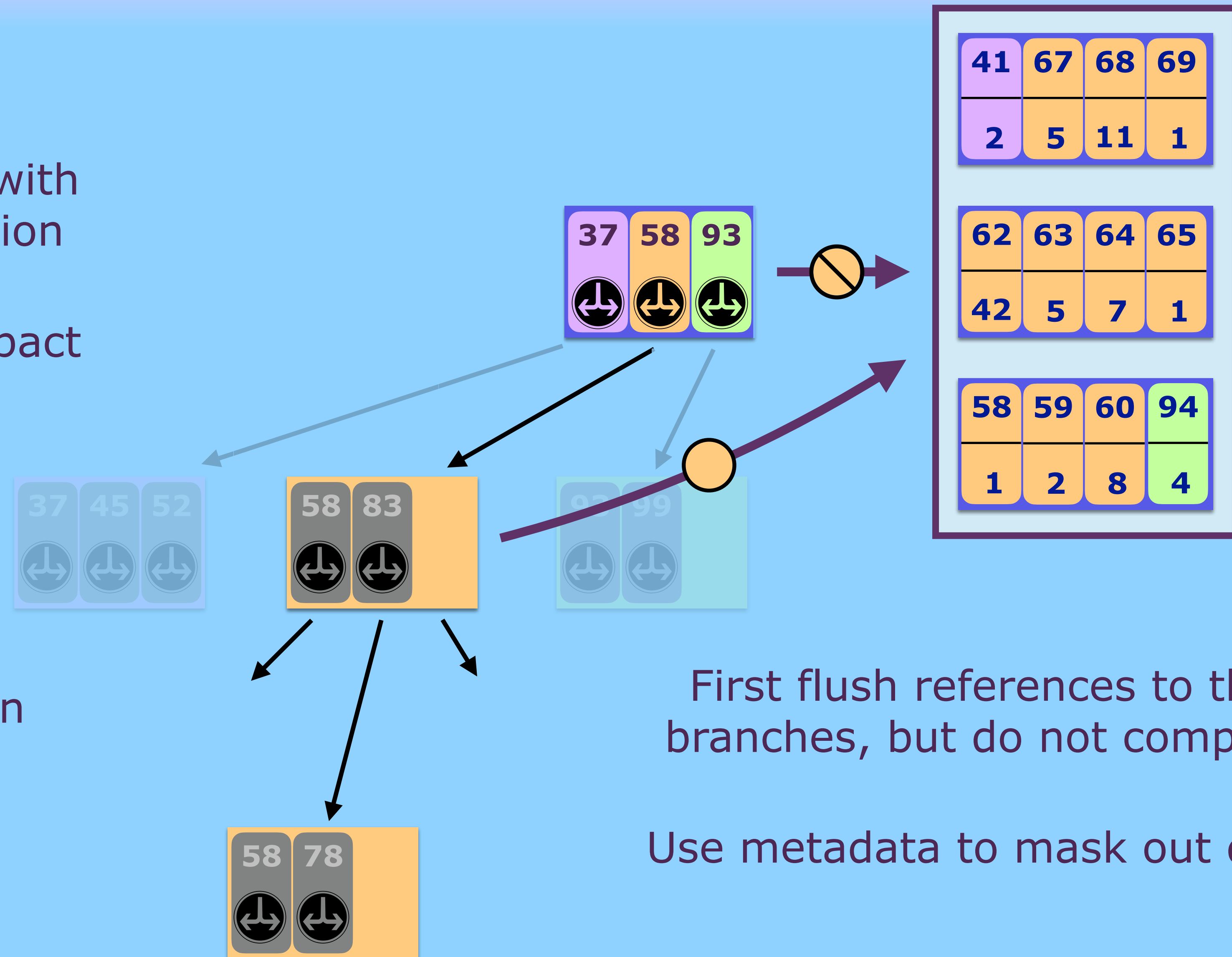
Use metadata to mask out data



Flush-Then-Compact

Want:
Sequential insertions with
lower work amplification

Idea: Flush-then-compact



Then can flush again

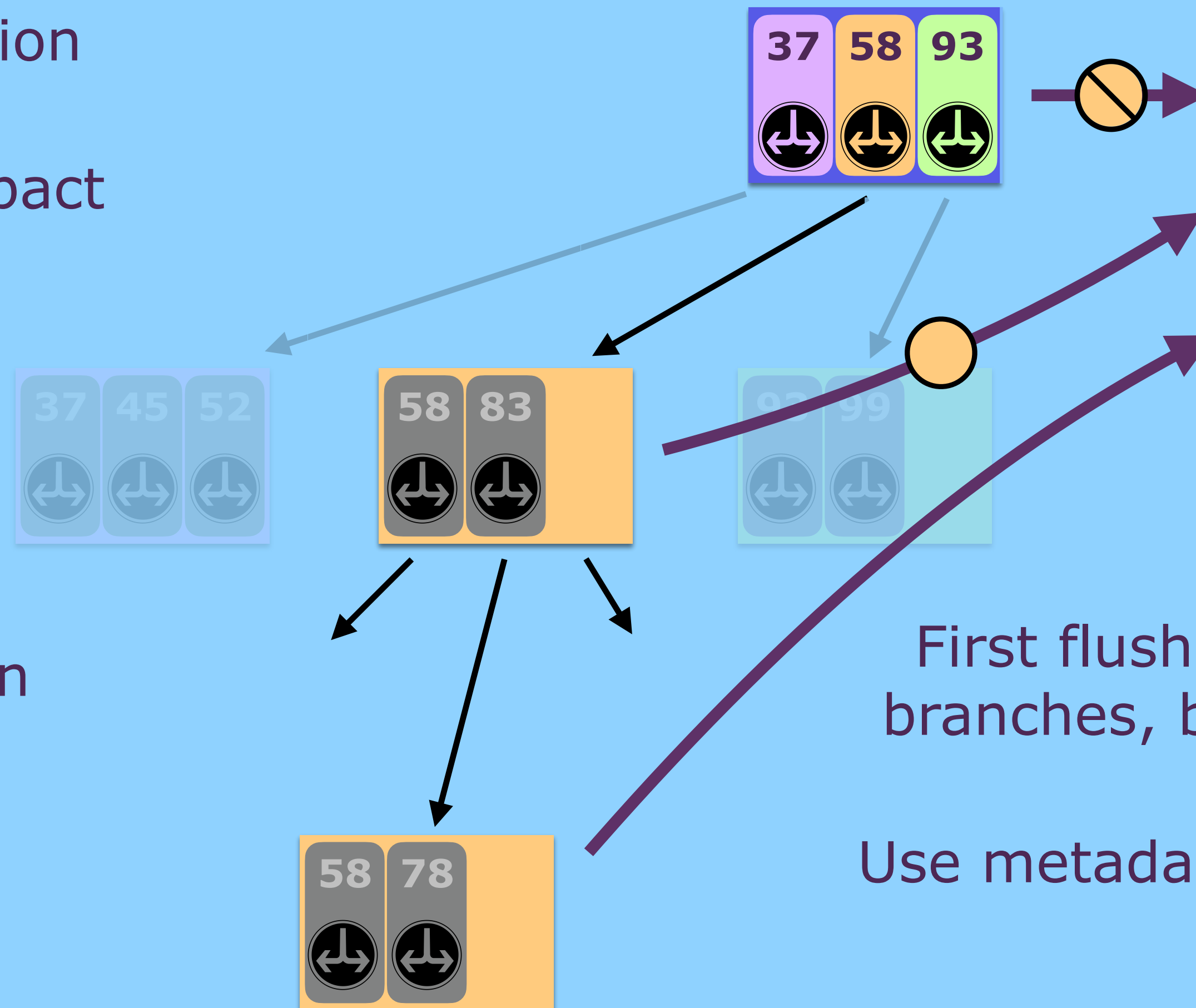
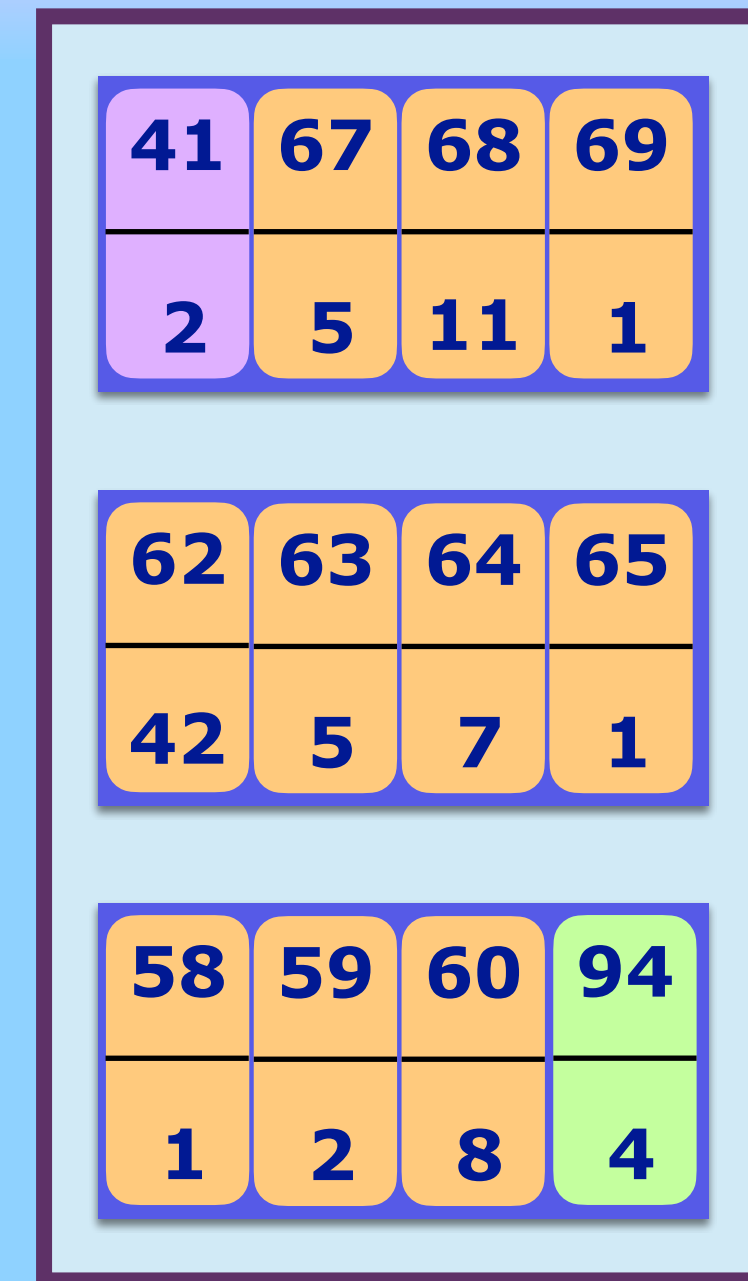
First flush references to the
branches, but do not compact

Use metadata to mask out data

Flush-Then-Compact

Want:
Sequential insertions with
lower work amplification

Idea: Flush-then-compact



Then can flush again

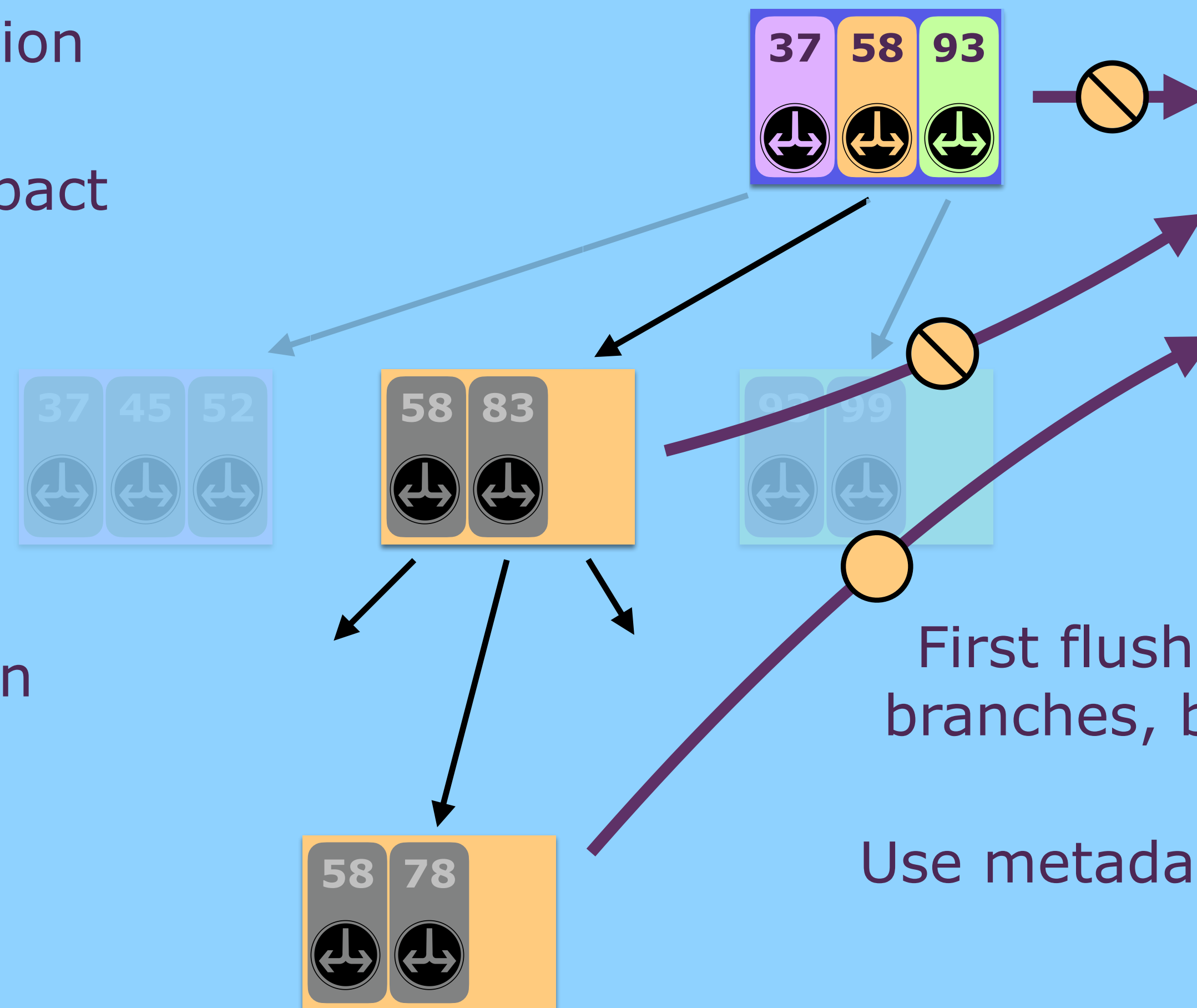
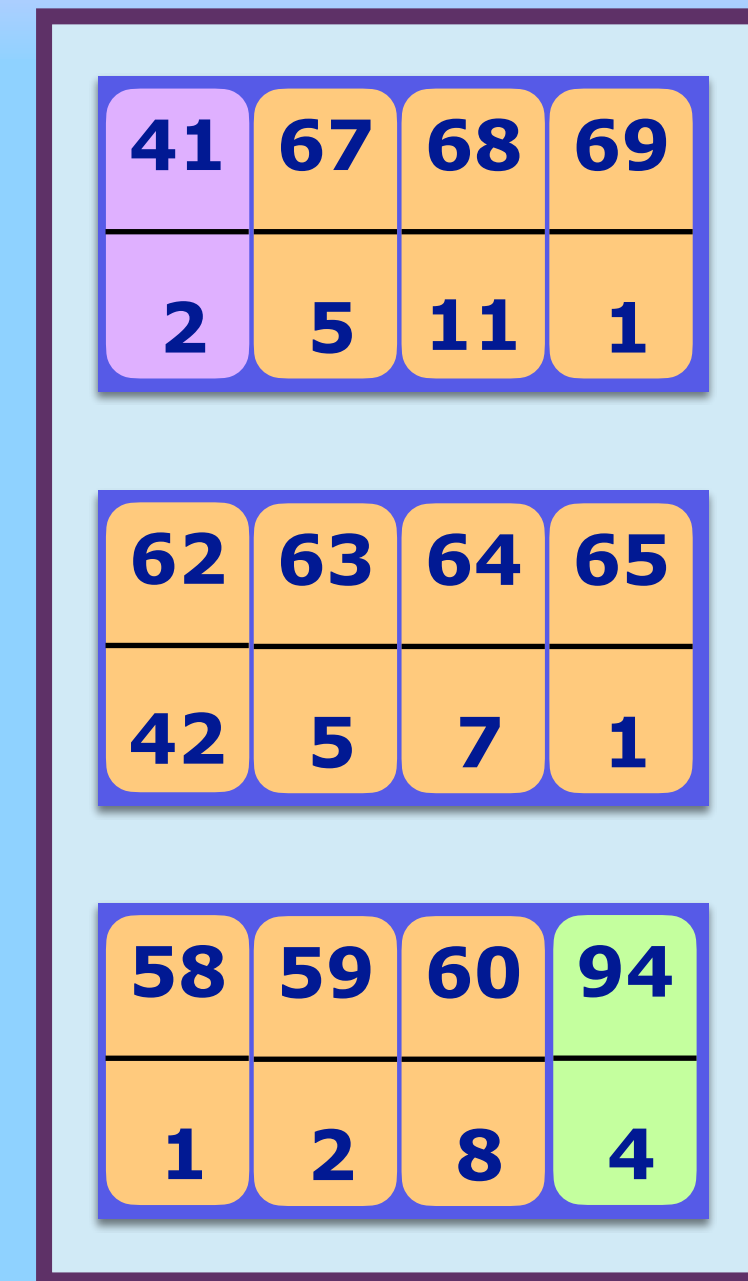
First flush references to the
branches, but do not compact

Use metadata to mask out data

Flush-Then-Compact

Want:
Sequential insertions with
lower work amplification

Idea: Flush-then-compact



Then can flush again

First flush references to the
branches, but do not compact

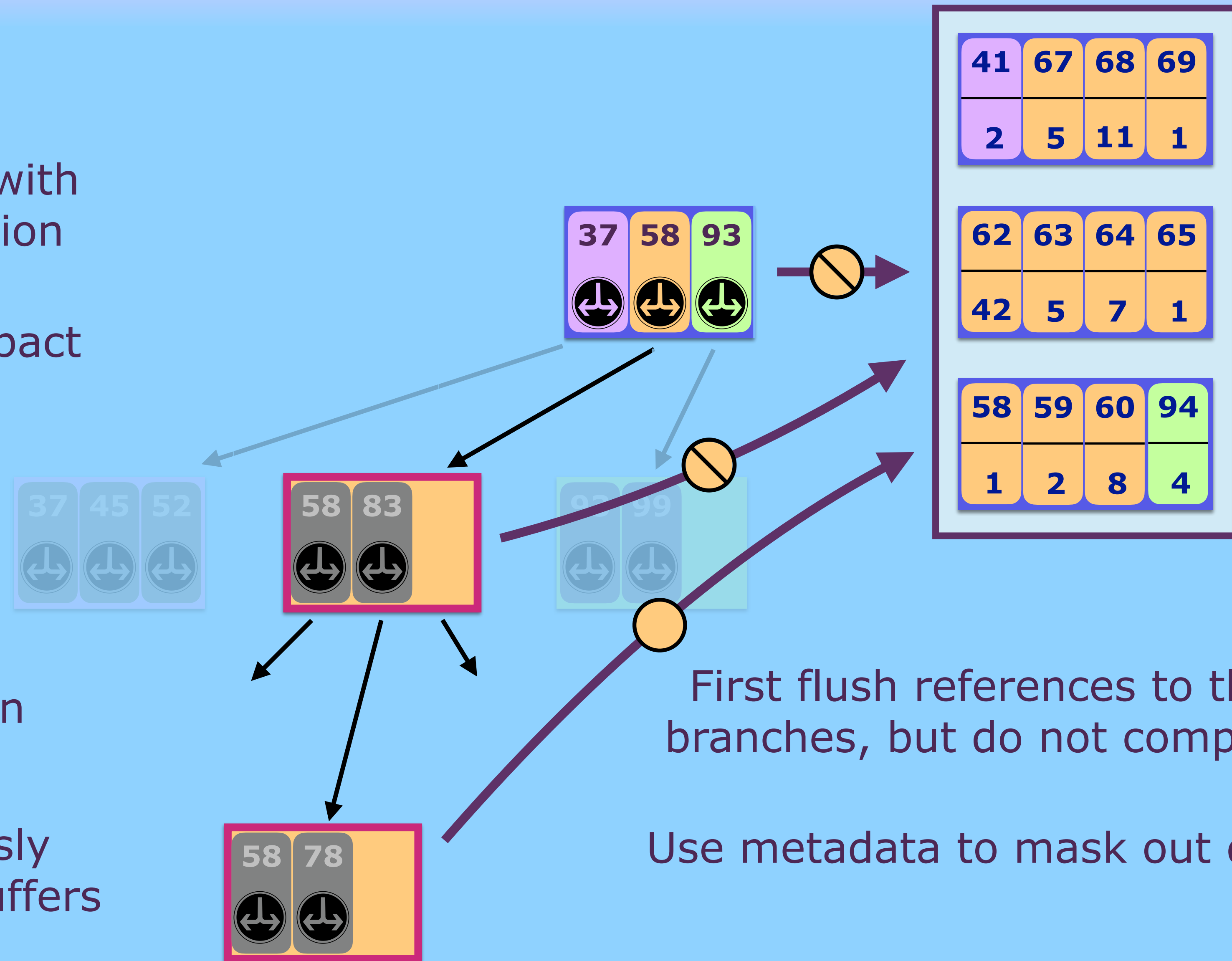
Use metadata to mask out data

Flush-Then-Compact

Want:

Sequential insertions with
lower work amplification

Idea: Flush-then-compact



Then can flush again

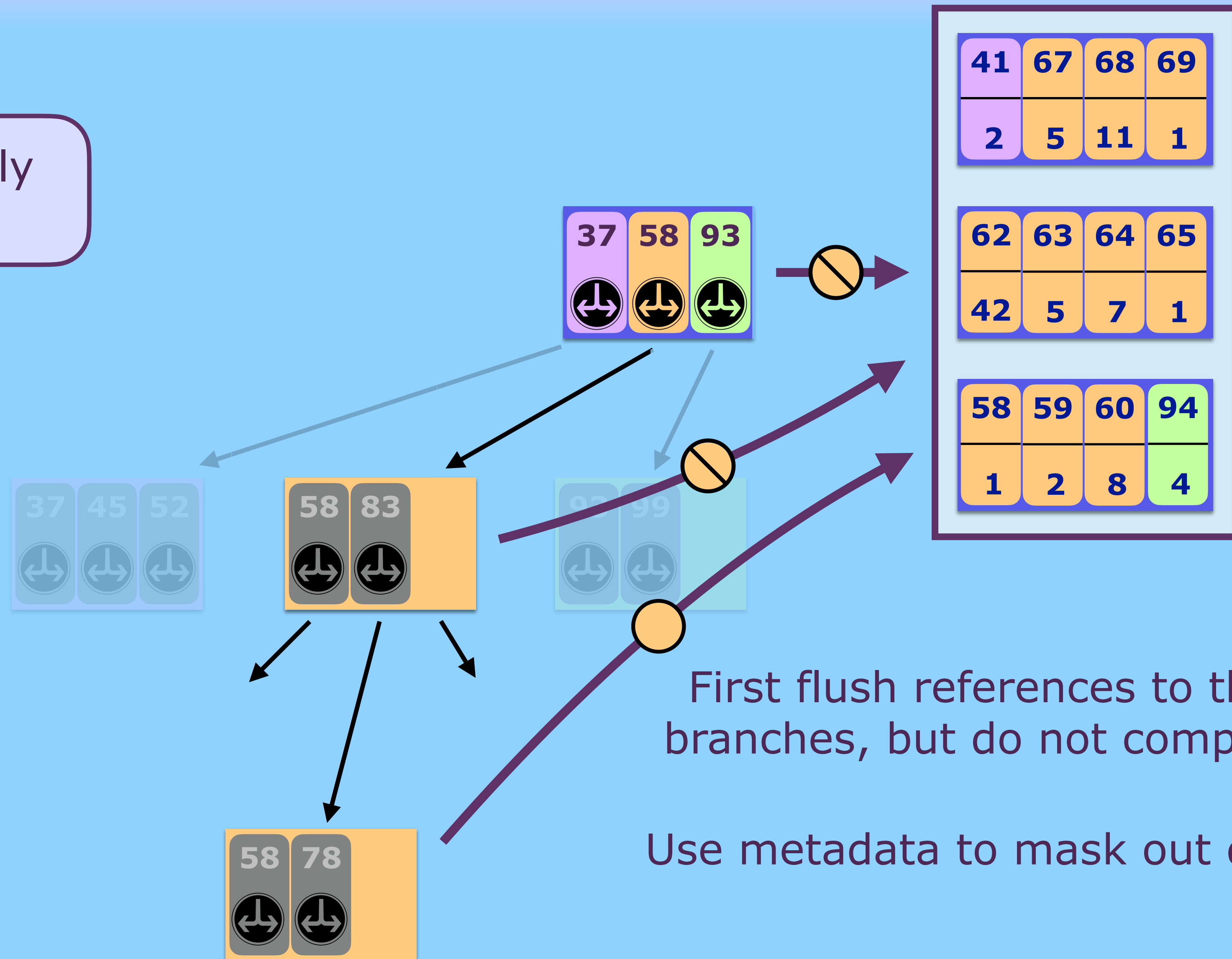
Finally, asynchronously
compact the flushed buffers
in each node

First flush references to the
branches, but do not compact

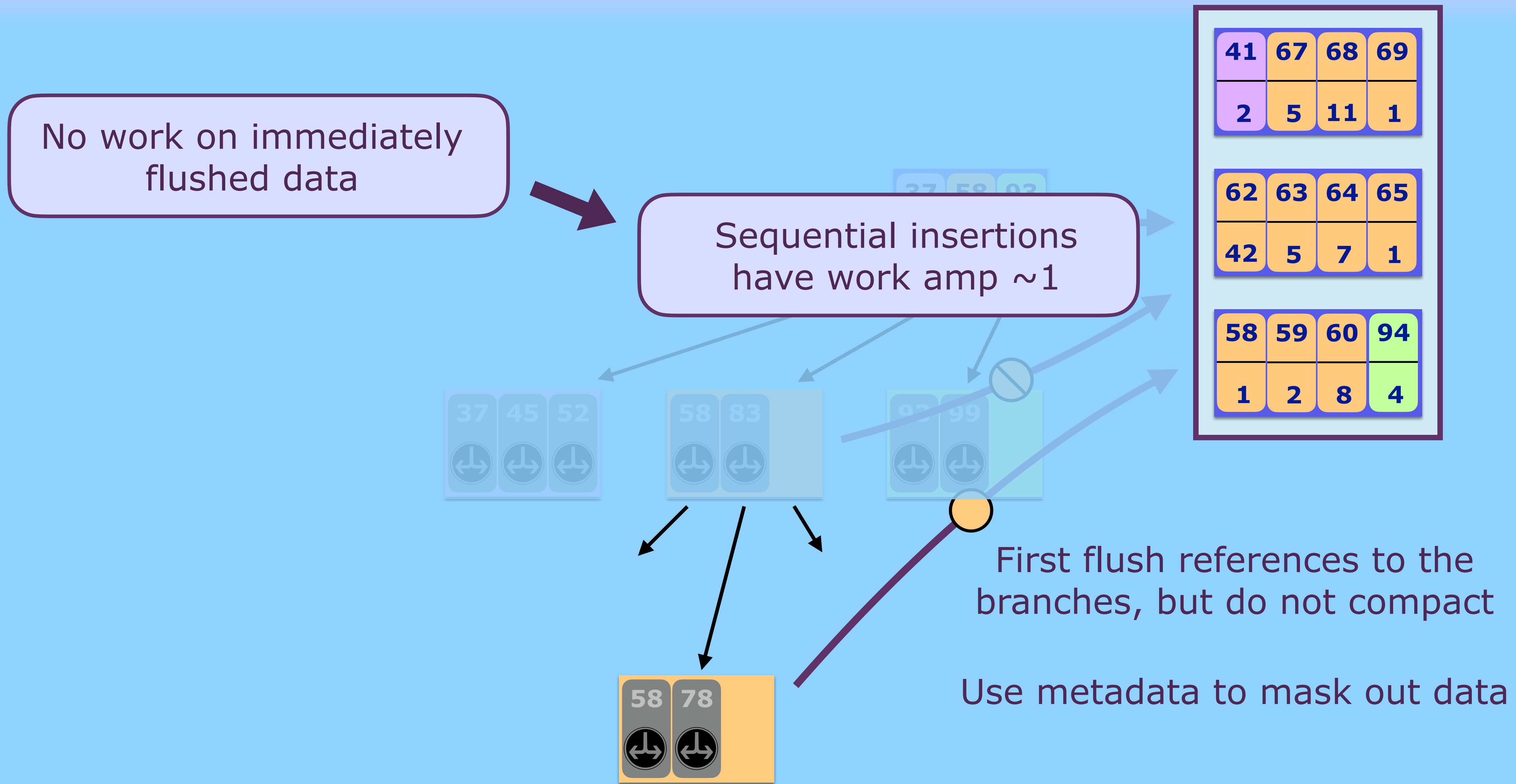
Use metadata to mask out data

Flush-Then-Compact

No work on immediately flushed data



Flush-Then-Compact



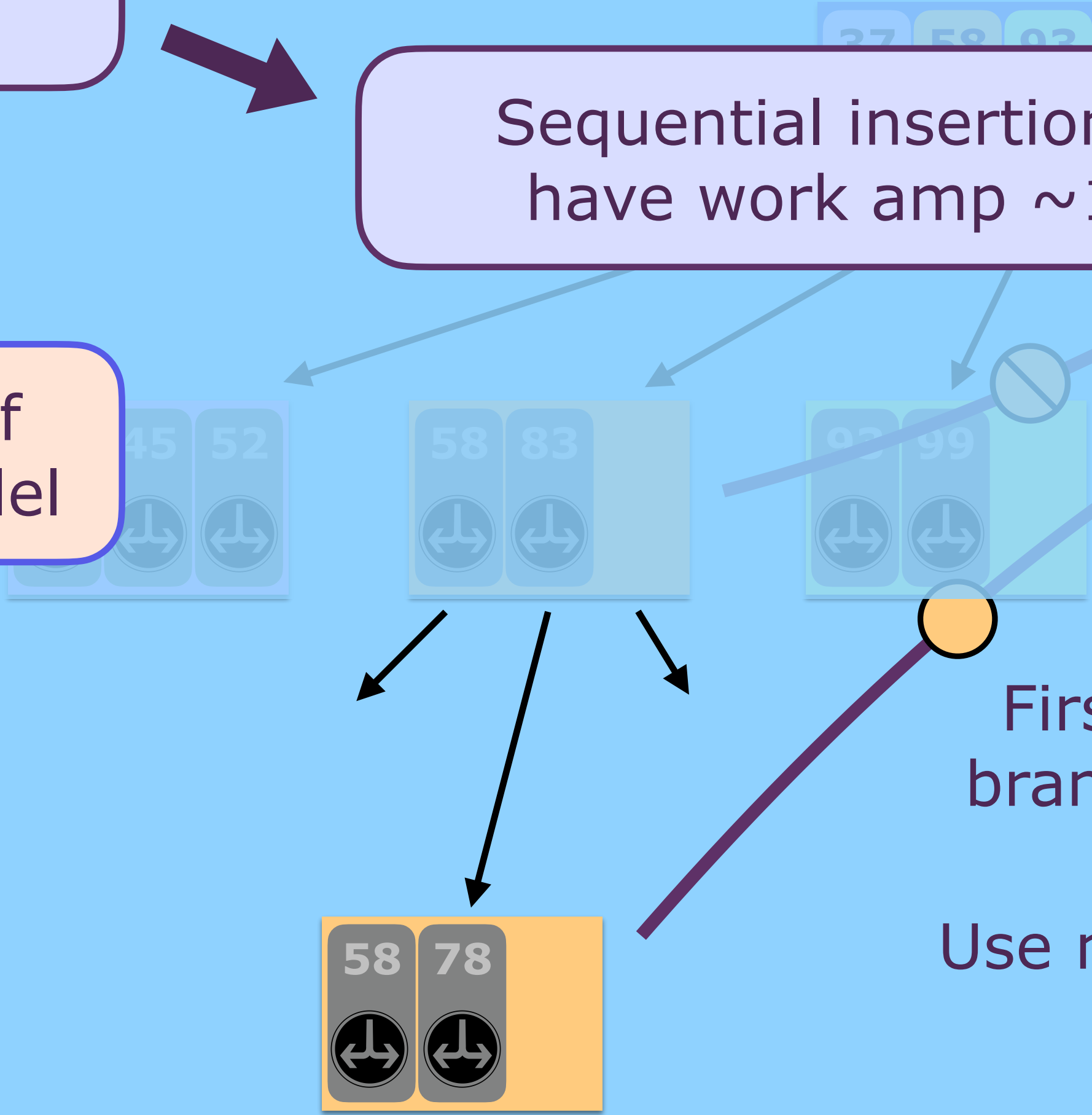
Flush-Then-Compact

No work on immediately flushed data

Sequential insertions have work amp ~ 1

Break a serial chain of compactions into parallel

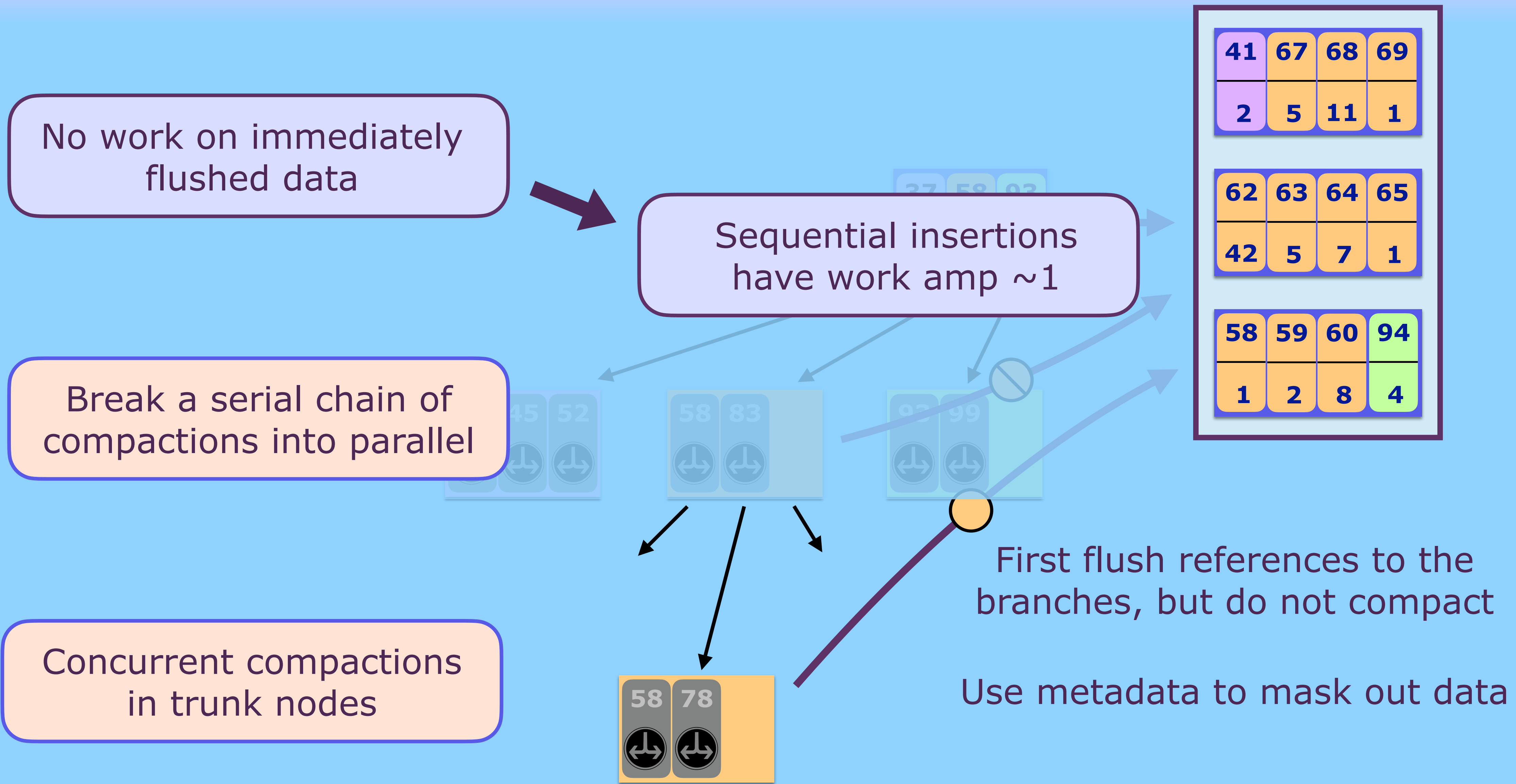
41	67	68	69
2	5	11	1
62	63	64	65
42	5	7	1
58	59	60	94
1	2	8	4



First flush references to the branches, but do not compact

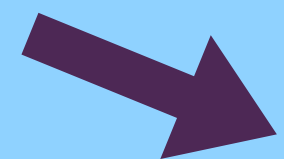
Use metadata to mask out data

Flush-Then-Compact



Flush-Then-Compact

No work on immediately flushed data



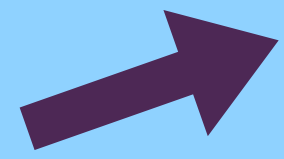
Sequential insertions have work amp ~ 1

Break a serial chain of compactions into parallel

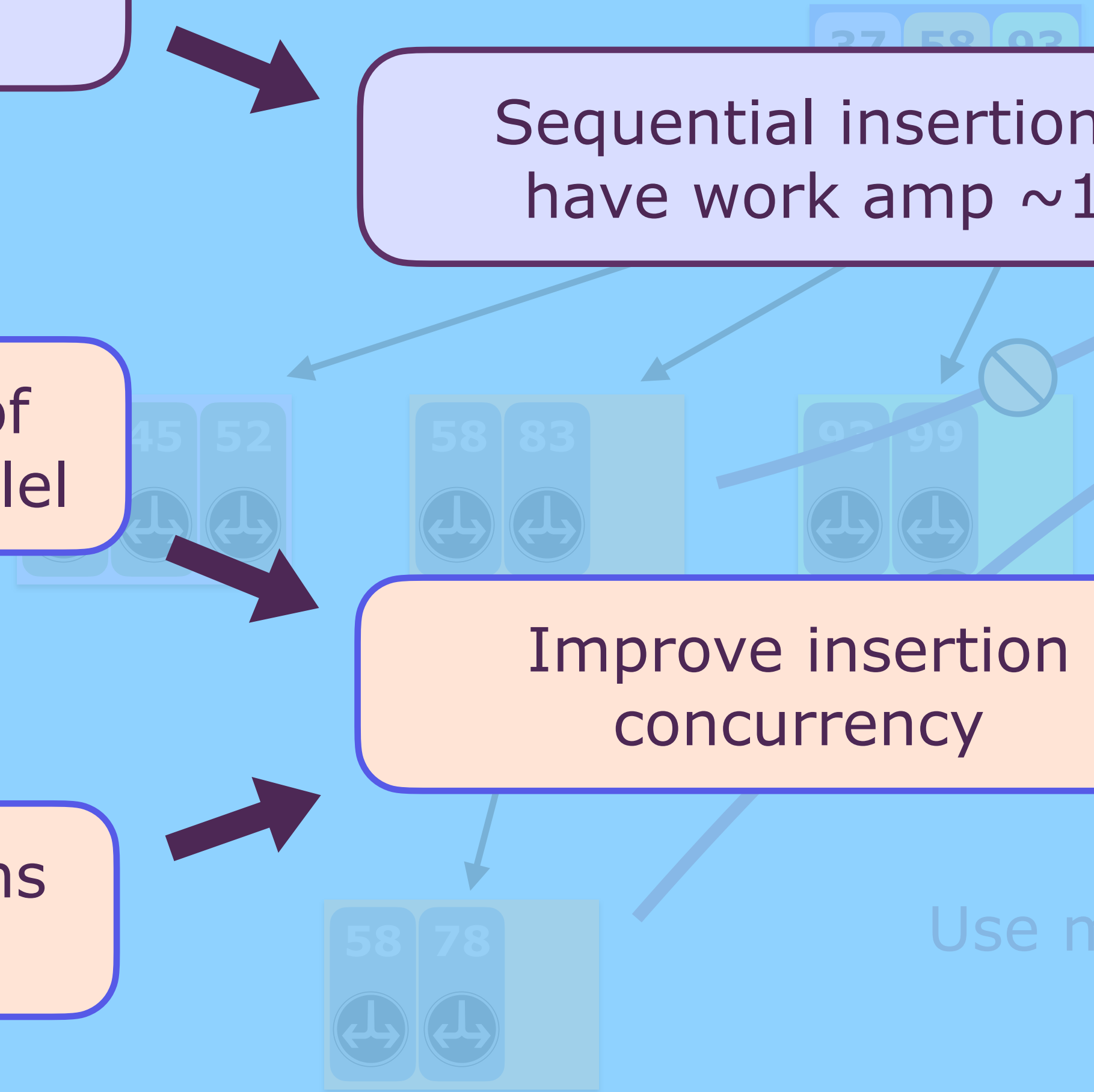


Improve insertion concurrency

Concurrent compactions in trunk nodes



41	67	68	69
2	5	11	1
62	63	64	65
42	5	7	1
58	59	60	94
1	2	8	4

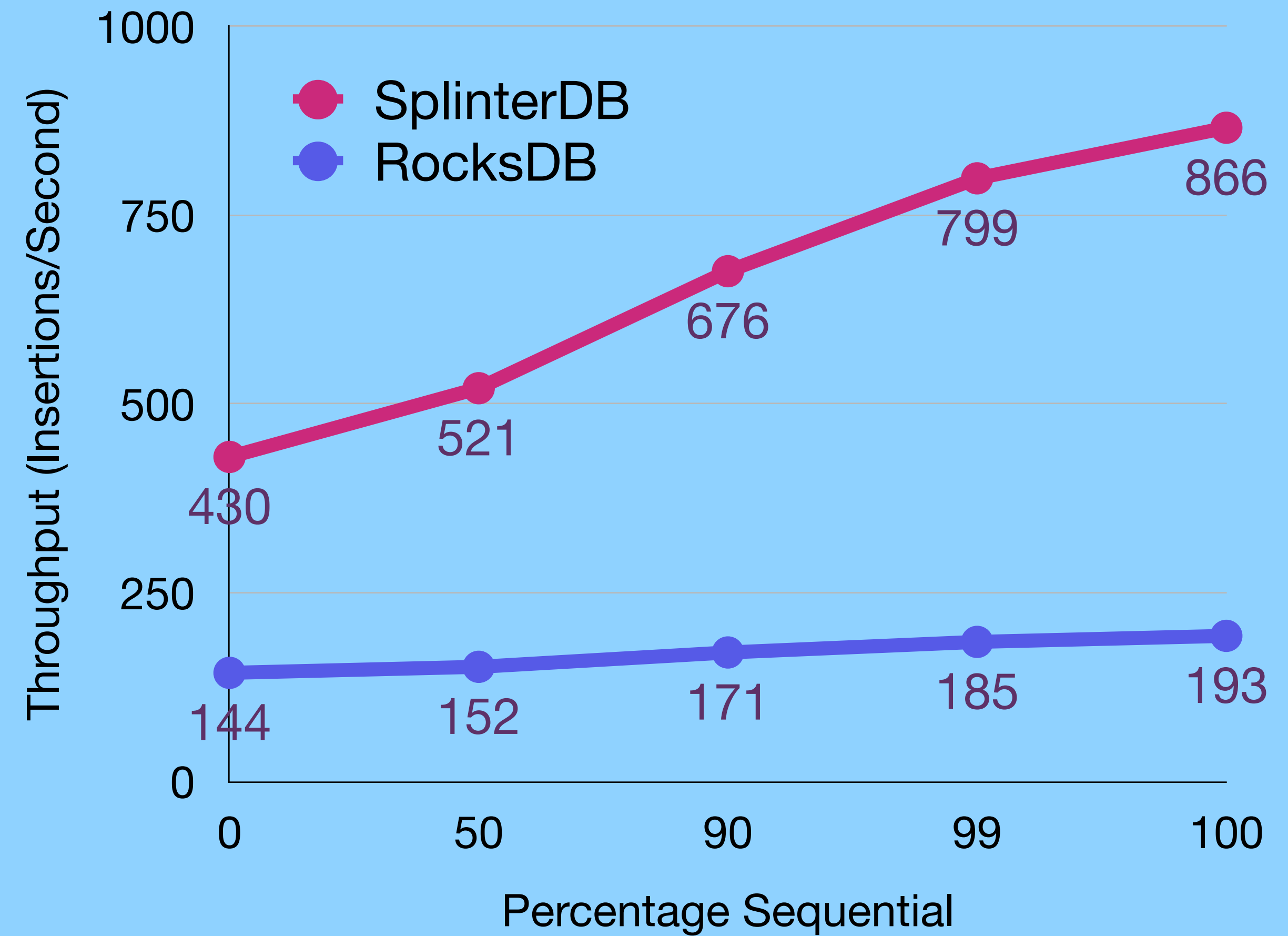


with references to the
, but do not compact

Use metadata to mask out data

Flush-Then-Compact

Run a single-threaded workload with a percentage sequential insertions and the rest random

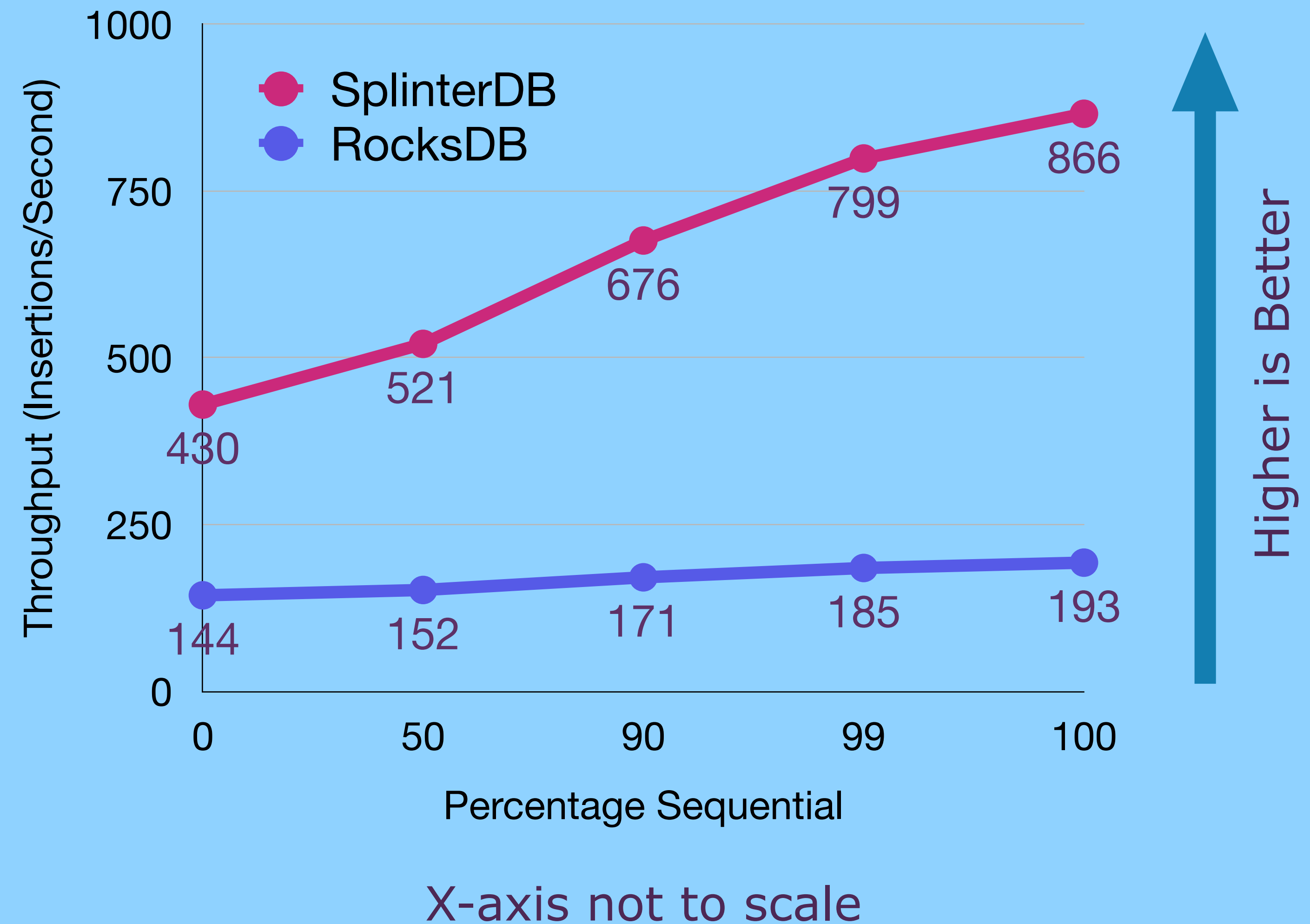


X-axis not to scale

Flush-Then-Compact

Run a single-threaded workload with a percentage sequential insertions and the rest random

Because of flush-then-compact, SplinterDB smoothly increases throughput as the workload gets more sequential

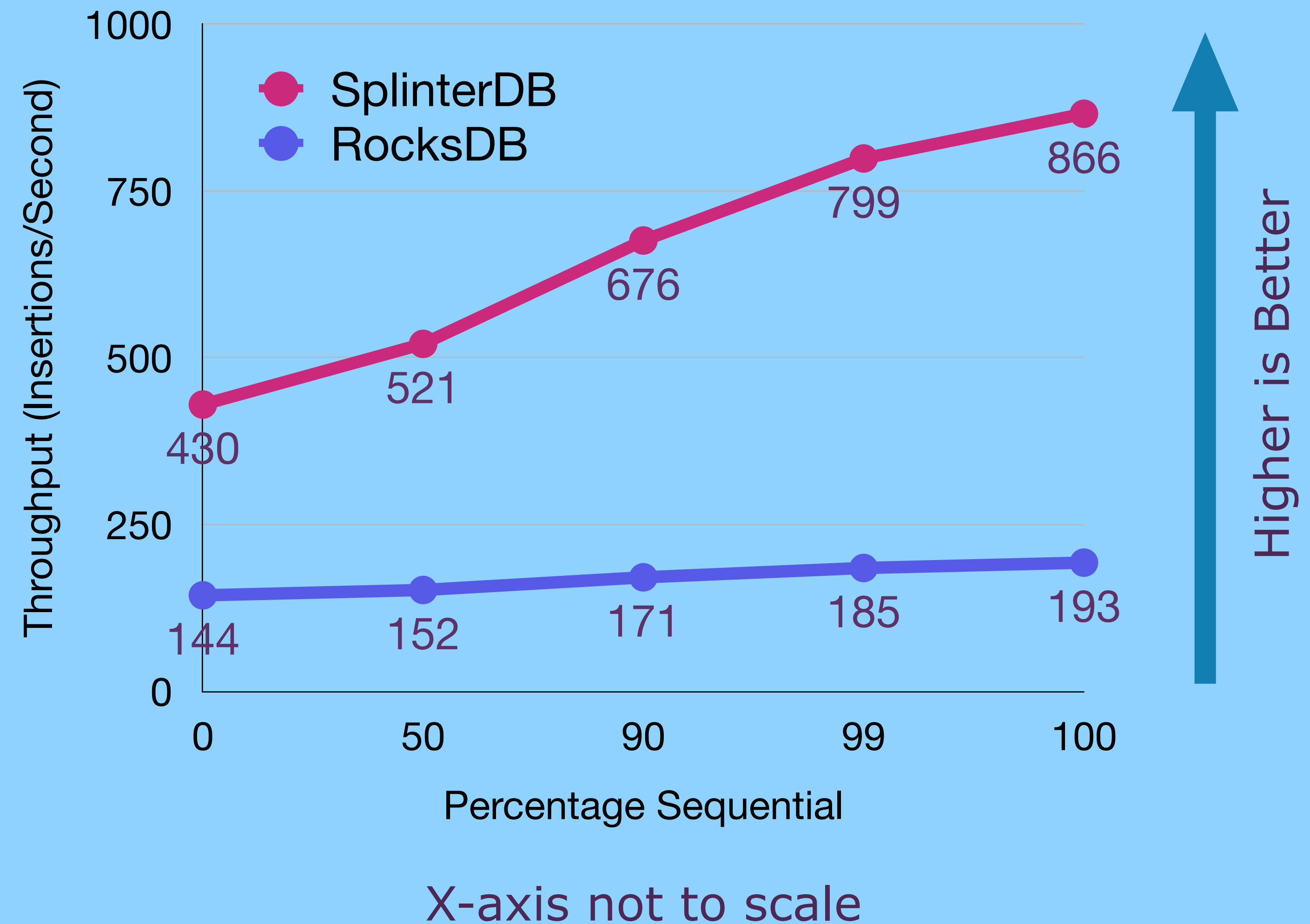


Flush-Then-Compact

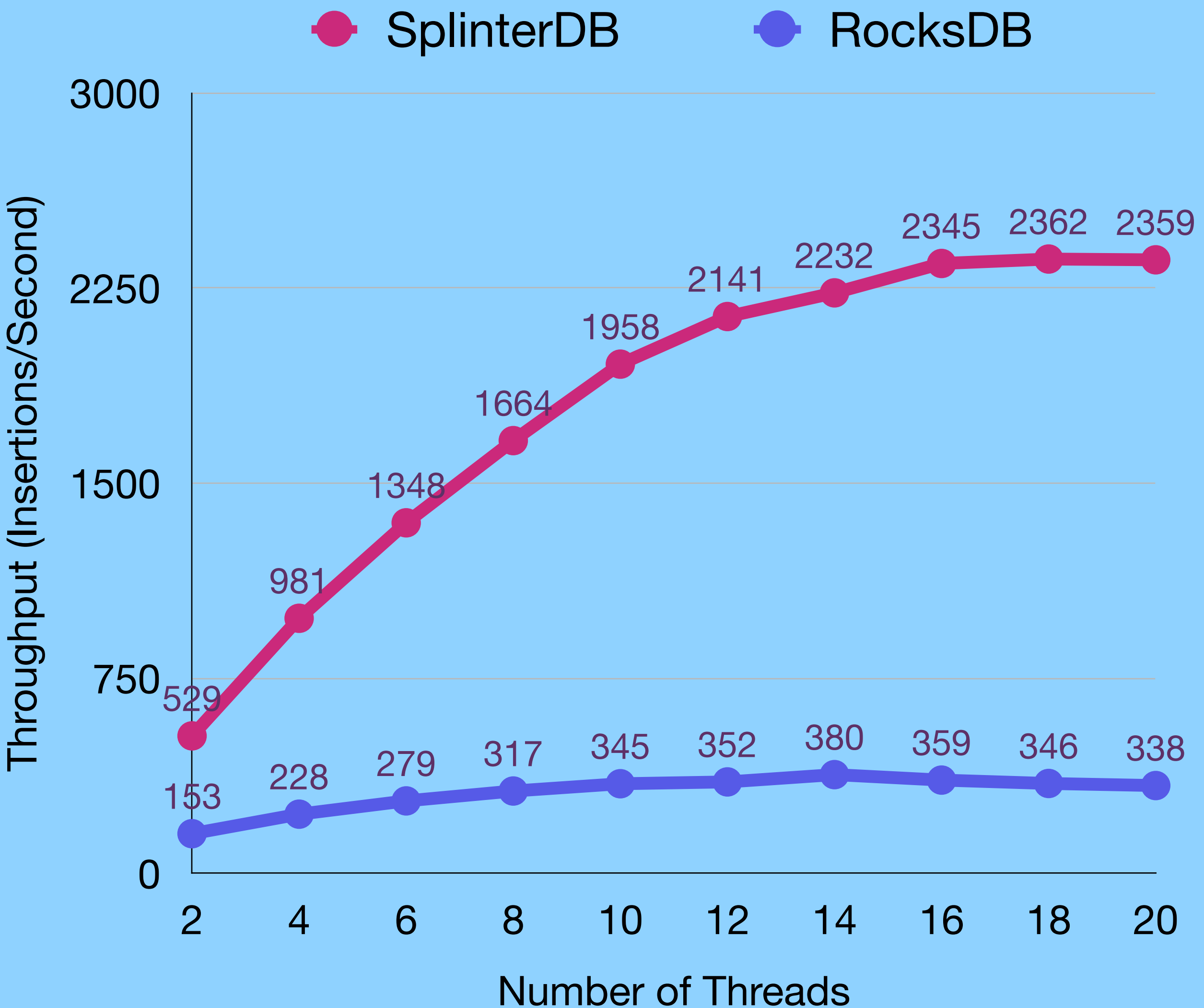
Run a single-threaded workload with a percentage sequential insertions and the rest random

Because of flush-then-compact, SplinterDB smoothly increases throughput as the workload gets more sequential

RocksDB improves, but at a much lower rate



Flush-then-Compact

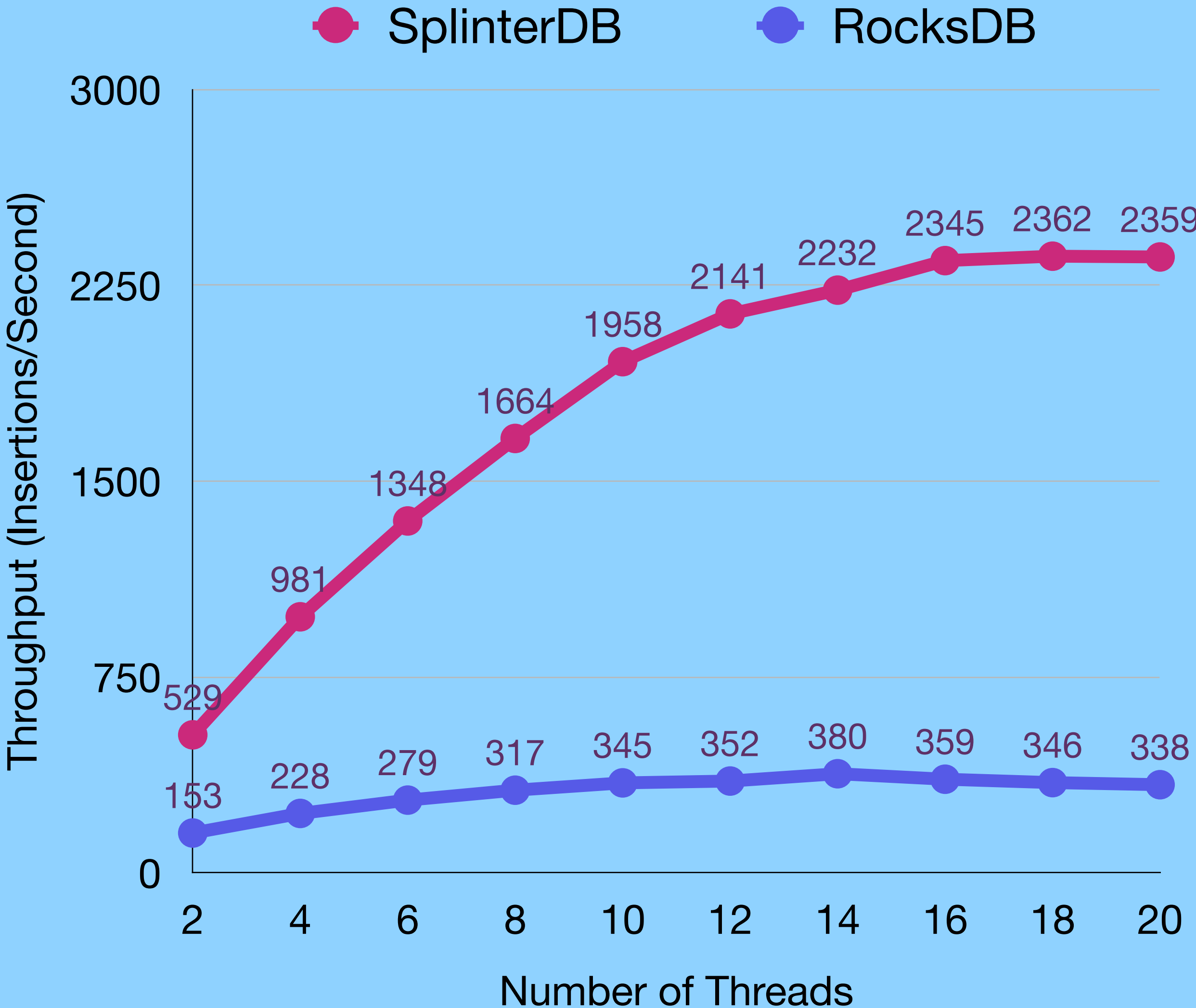


Insertions in SplinterDB scale well



Higher is Better

Flush-then-Compact

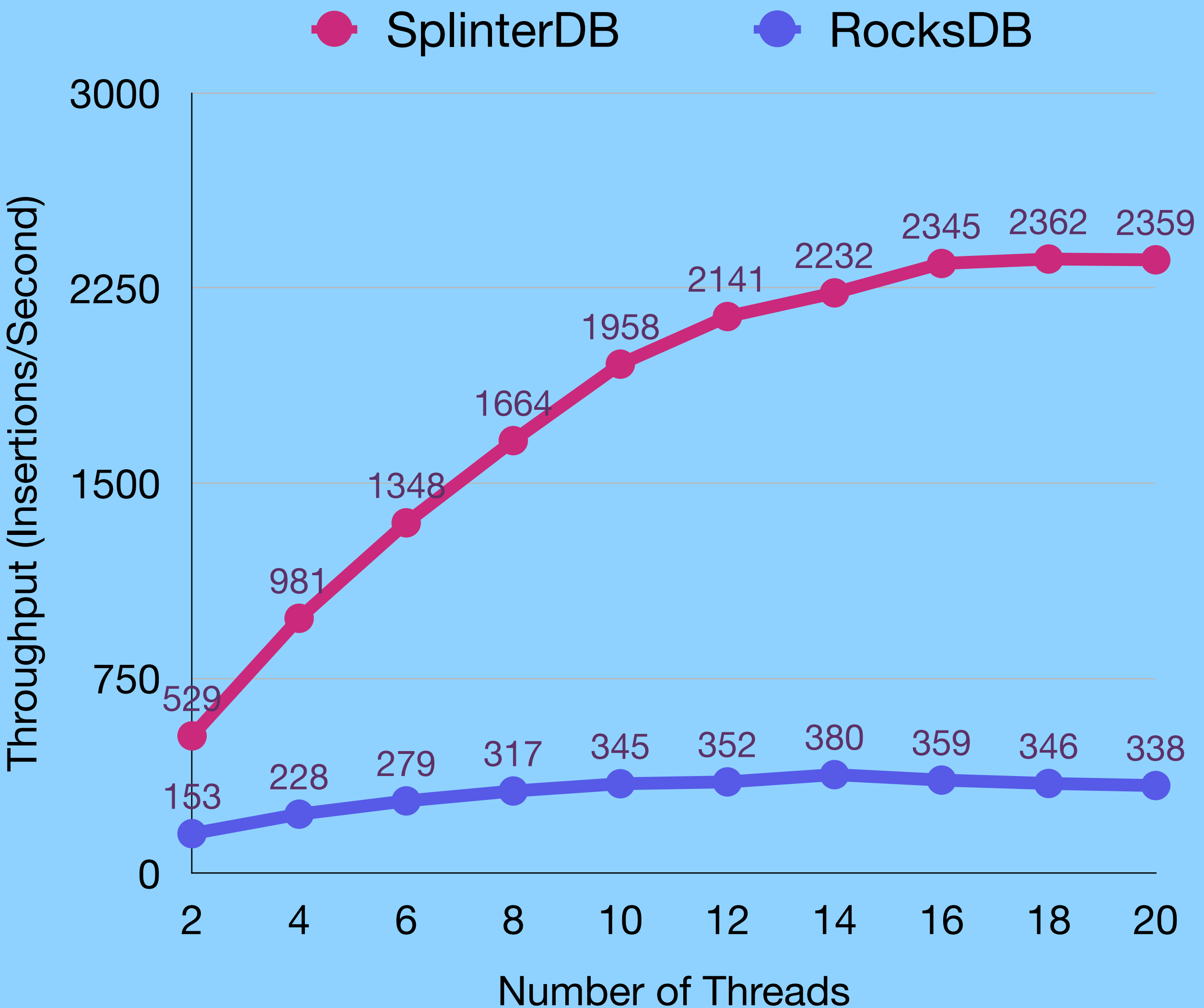


Insertions in SplinterDB scale well

At 12 threads, SplinterDB has 7x the throughput of 1 thread

Higher is Better

Flush-then-Compact



Insertions in SplinterDB scale well

At 12 threads, SplinterDB has 7x the throughput of 1 thread

At 12+ threads, SplinterDB uses 85%+ of the device bandwidth

Conclusion

Conclusion

SplinterDB is a key-value store which handles these tough cases:

Fast Storage



Small Key-Value Pairs



Small Cache



Conclusion

SplinterDB is a key-value store which handles these tough cases:

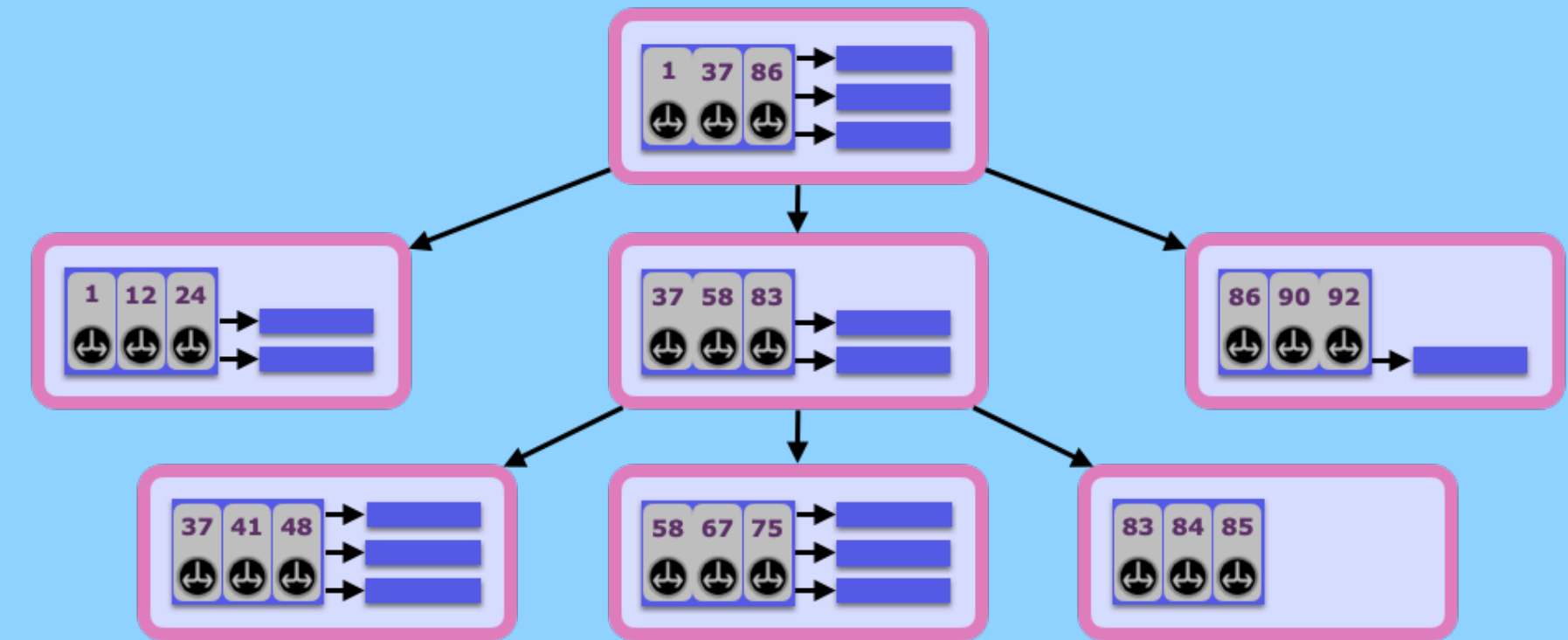
Fast Storage



Small Key-Value Pairs



Small Cache



Size-Tiered B^ϵ -Tree

Conclusion

SplinterDB is a key-value store which handles these tough cases:

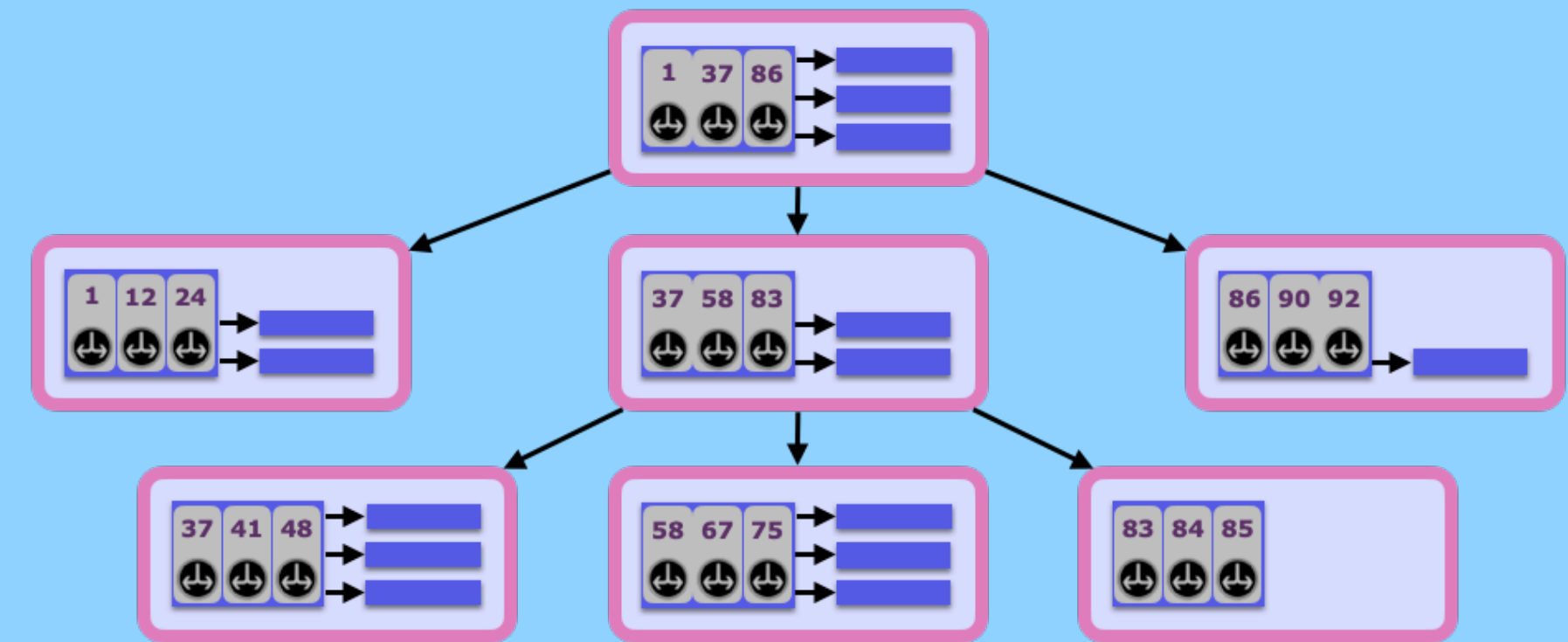
Fast Storage



Small Key-Value Pairs



Small Cache



Size-Tiered B^ϵ -Tree

Flush-then-Compact

Thank you!!!

Alex Conway

ajhconway.com

aconway@vmware.com