

CS 7270: Advanced Database Systems Fall 2025

Lecture 11

B^e-tree and SplinterDB

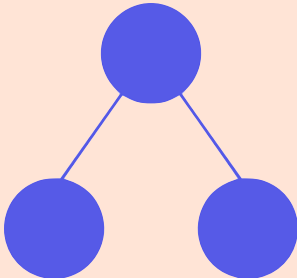
Prashant Pandey

p.pandey@northeastern.edu

Slides taken from Prof. Alex Conway, Cornell Tech

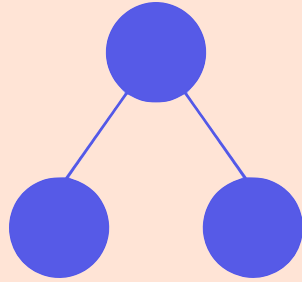
The Story of SplinterDB

Model the problem:
external memory dictionary



The Story of SplinterDB

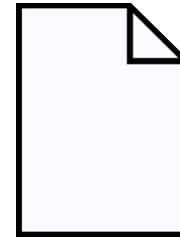
Model the problem:
external memory dictionary



48 B



4 KiB

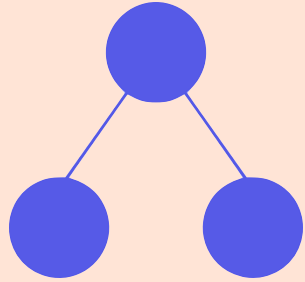


IO 4 KiB



The Story of SplinterDB

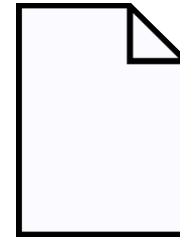
Model the problem:
external memory dictionary



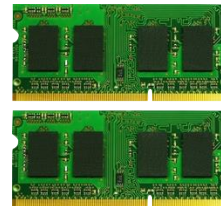
48 B



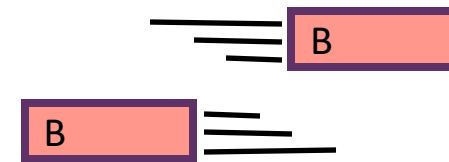
4 KiB



IO 4 KiB



Internal
Memory of size
M



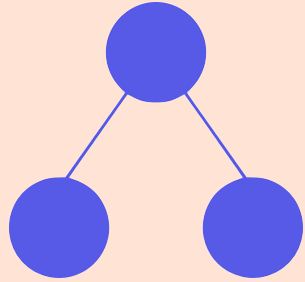
A B-sized block can be read or
written in 1 IO



External Memory Model

The Story of SplinterDB

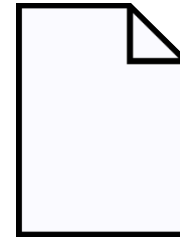
Model the problem:
external memory dictionary



48 B



4 KiB

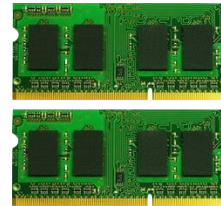


IO 4 KiB

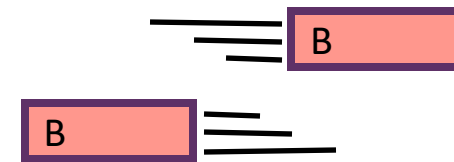


Here B is the number of
items in an IO:
 $B = 4 \text{ KiB} / 48 \text{ B}$

If the items were larger, the model
wouldn't be as good



Internal
Memory of size
M



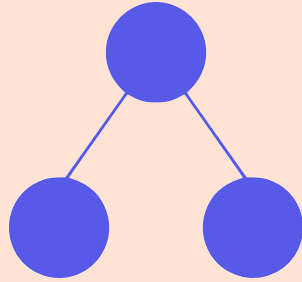
A B-sized block can be read or
written in 1 IO



External Memory Model

The Story of SplinterDB

Model the problem:
external memory dictionary



Two Flavors of
External-Memory Dictionary

Different lower bounds
(performance limits)

Comparison-Based Dictionaries

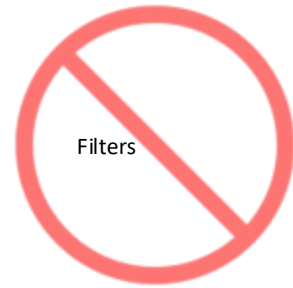
Comparison External Memory Model

user024299 < user082587
= user082587
>

Comparison-Based Dictionaries

Comparison External Memory
Model

user024299 <
= user082587
>

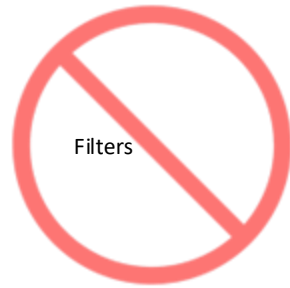


Comparison-Based Dictionaries

Comparison External Memory Model

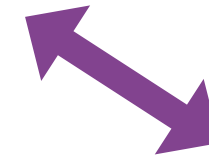
Brodal-Fagerberg Lower Bound

user024299 <
= user082587
>



Insertions in

$$O\left(\frac{\lambda}{B} \log_{\lambda} N\right)$$



Lookups in

$$\Omega(\log_{\lambda} N)$$

where λ is a tuning parameter

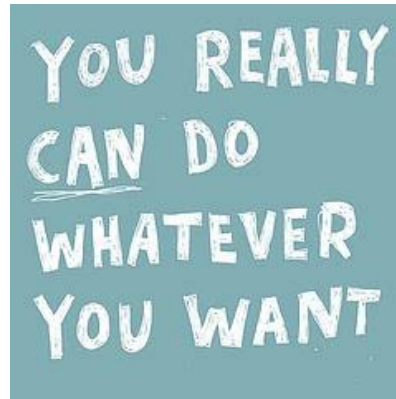
General Dictionaries

External Memory Model

General Dictionaries

External Memory Model

user024299



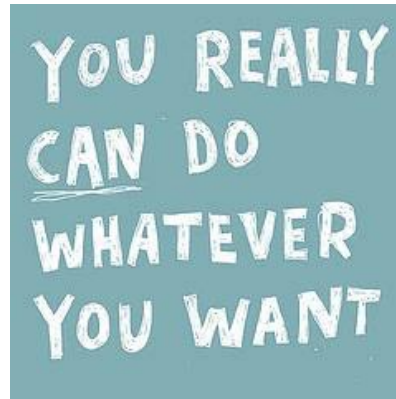
General Dictionaries

External Memory Model

user024299

Hashing

XXH (user024299)



General Dictionaries

External Memory Model

user024299

Hashing

XXH (user024299)

Filters

qf_insert (user024299)



General Dictionaries

External Memory Model

user024299

Hashing

XXH (user024299)

Filters

qf_insert (user024299)

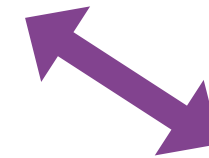


*Using hashing to solve the dictionary problem (in external memory),
Iacono, J., Pătraşcu, M. SODA '12*

Iacono-Pătraşcu Lower Bound

Insertions in

$$O\left(\frac{\lambda}{B} \log_{\lambda} N\right)$$



Lookups in

$$\Omega(\log_{\lambda} N)$$

where λ is a tuning parameter

Lower Bounds

Brodal-Fagerberg Lower Bound

Insertions in

$$O\left(\frac{\lambda}{B} \log_{\lambda} N\right)$$



Lookups in

$$\Omega(\log_{\lambda} N)$$

Iacono-Pătraşcu Lower Bound

Insertions in

$$O\left(\frac{\lambda}{B}\right)$$



Lookups in

$$\Omega(\log_{\lambda} N)$$

Comparison External Memory Model

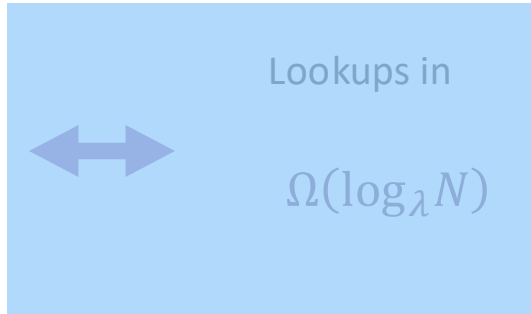
General External Memory Model

Lower Bounds

Brodal-Fagerberg Lower Bound

Insertions in

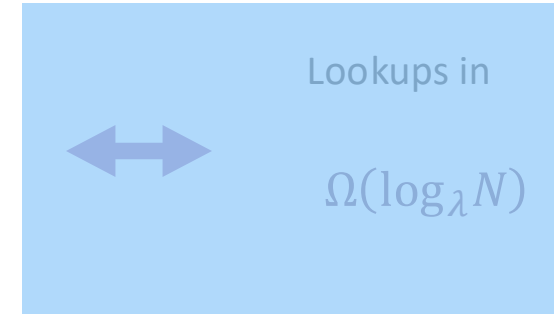
$$O\left(\frac{\lambda}{B} \log_{\lambda} N\right)$$



Iacono-Pătraşcu Lower Bound

Insertions in

$$O\left(\frac{\lambda}{B}\right)$$

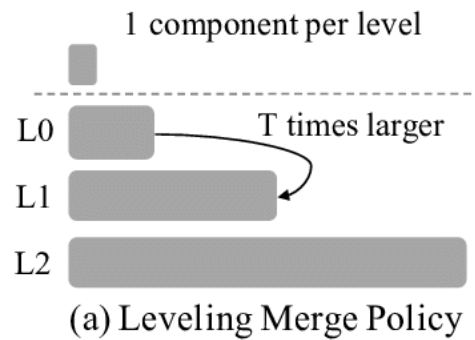
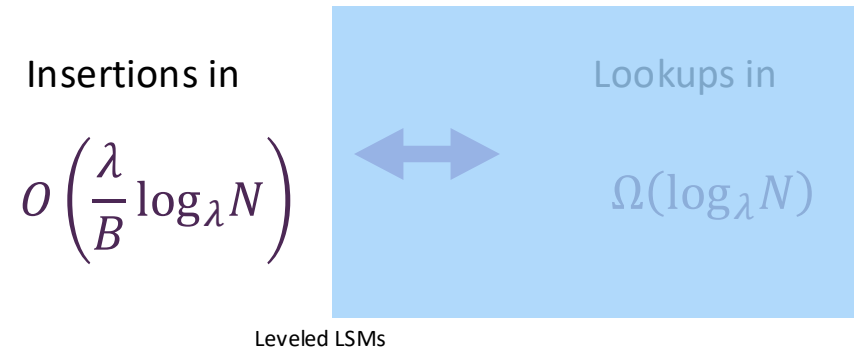


Comparison External Memory Model

General External Memory Model

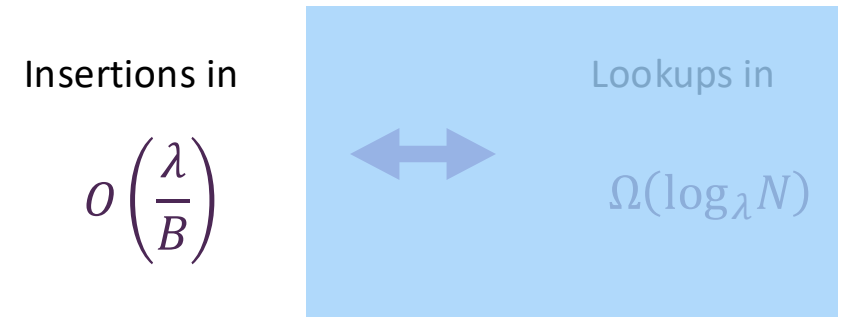
Lower Bounds

Brodal-Fagerberg Lower Bound



Comparison External Memory Model

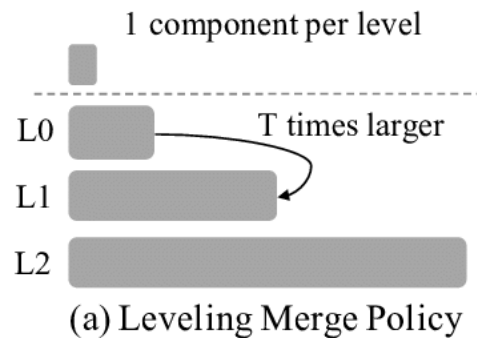
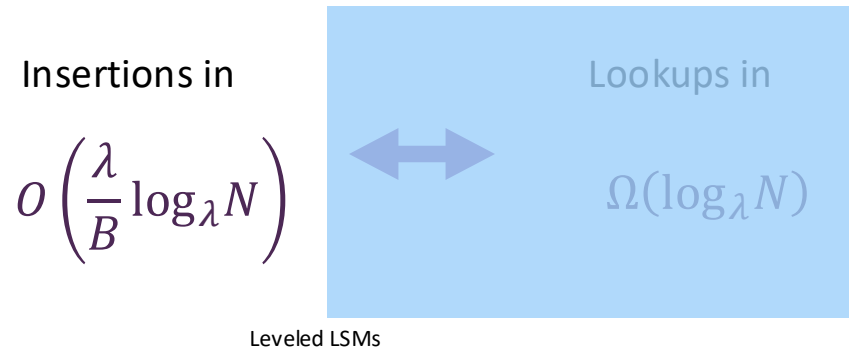
Iacono-Pătraşcu Lower Bound



General External Memory Model

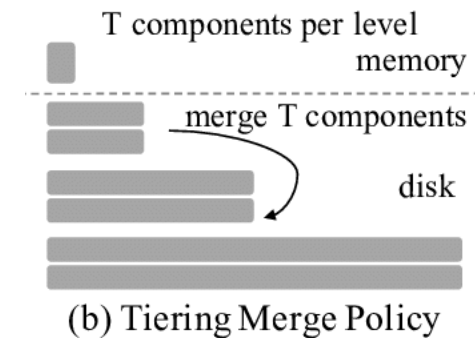
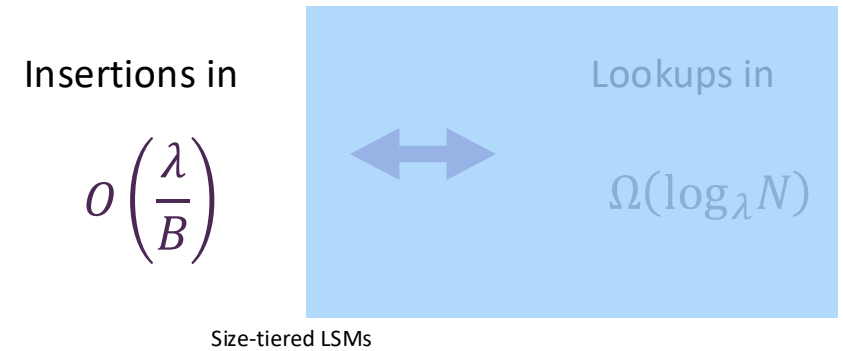
Lower Bounds

Brodal-Fagerberg Lower Bound



Comparison External Memory Model

Iacono-Pătrașcu Lower Bound



General External Memory Model

Lower Bounds

Brodal-Fagerberg Lower Bound

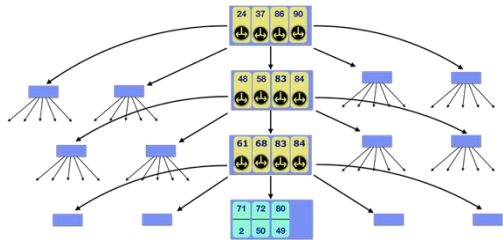
Insertions in

$$O\left(\frac{\lambda}{B} \log_{\lambda} N\right)$$



Lookups in

$$\Omega(\log_{\lambda} N)$$



B-Trees

$$(\lambda = B)$$

Comparison External Memory Model

Iacono-Pătrașcu Lower Bound

Insertions in

$$O\left(\frac{\lambda}{B}\right)$$



Lookups in

$$\Omega(\log_{\lambda} N)$$

General External Memory Model

Lower Bounds

Brodal-Fagerberg Lower Bound

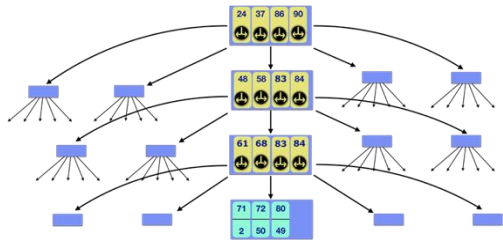
Insertions in

$$O\left(\frac{\lambda}{B} \log_{\lambda} N\right)$$



Lookups in

$$\Omega(\log_{\lambda} N)$$



B-Trees

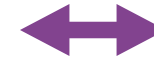
$$(\lambda = B)$$

Comparison External Memory Model

Iacono-Pătrașcu Lower Bound

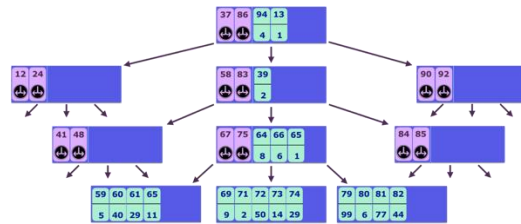
Insertions in

$$O\left(\frac{\lambda}{B}\right)$$



Lookups in

$$\Omega(\log_{\lambda} N)$$



B^ε-Trees

$$(\lambda = B^{\epsilon})$$

General External Memory Model

Lower Bounds

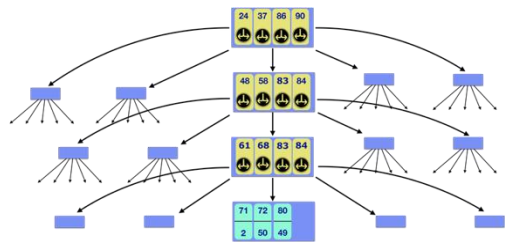
Brodal-Fagerberg Lower Bound

Insertions in

$$O\left(\frac{\lambda}{B} \log_{\lambda} N\right)$$

Lookups in

$$\Omega(\log_{\lambda} N)$$



B-Trees

$$(\lambda = B)$$

Comparison External Memory Model

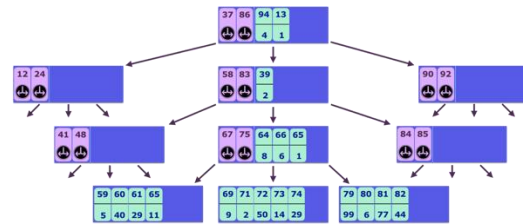
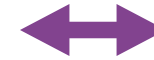
Iacono-Pătrașcu Lower Bound

Insertions in

$$O\left(\frac{\lambda}{B}\right)$$

Lookups in

$$\Omega(\log_{\lambda} N)$$



B^{ϵ} -Trees

$$(\lambda = B^{\epsilon})$$

General External Memory Model

Iacono-Patrascu Hash Table

BoA/BoT Hash Table

Lower Bounds

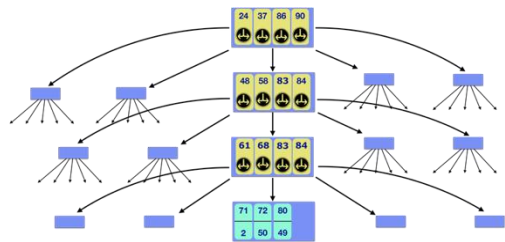
Brodal-Fagerberg Lower Bound

Insertions in

$$O\left(\frac{\lambda}{B} \log_{\lambda} N\right)$$

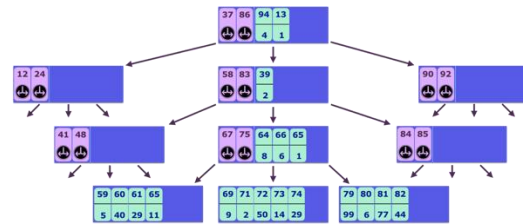
Lookups in

$$\Omega(\log_{\lambda} N)$$



B-Trees

$$(\lambda = B)$$



B^{ϵ} -Trees

$$(\lambda = B^{\epsilon})$$

Comparison External Memory Model

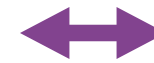
Iacono-Pătrașcu Lower Bound

Insertions in

$$O\left(\frac{\lambda}{B}\right)$$

Lookups in

$$\Omega(\log_{\lambda} N)$$



Iacono-Patrascu Hash Table

BoA/BoT
Hash Table

Optimal Hashing in External Memory, **Conway**, Farach-Colton, Shillane, ICALP 2018

Lower Bounds

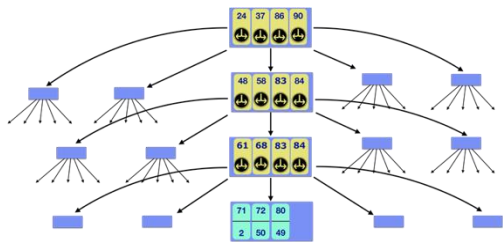
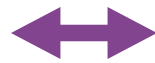
Brodal-Fagerberg Lower Bound

Insertions in

$$O\left(\frac{\lambda}{B} \log_{\lambda} N\right)$$

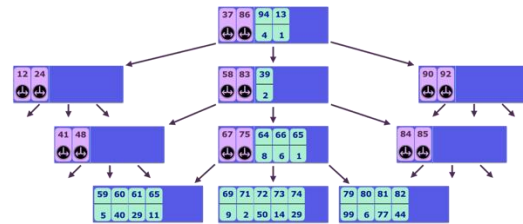
Lookups in

$$\Omega(\log_{\lambda} N)$$



B-Trees

$$(\lambda = B)$$



B^{ϵ} -Trees

$$(\lambda = B^{\epsilon})$$

Comparison External Memory Model

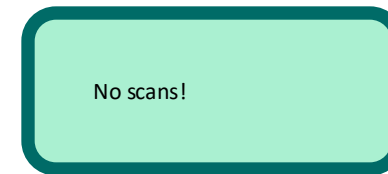
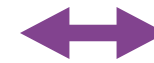
Iacono-Pătraşcu Lower Bound

Insertions in

$$O\left(\frac{\lambda}{B}\right)$$

Lookups in

$$\Omega(\log_{\lambda} N)$$



Iacono-Patrascu Hash Table

BoA/BoT Hash Table

General External Memory Model

Lower Bounds

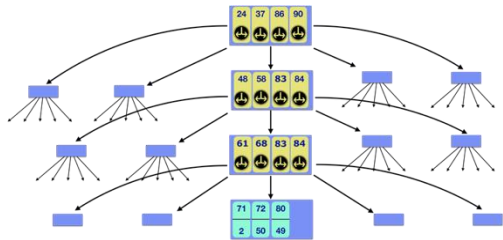
Brodal-Fagerberg Lower Bound

Insertions in

$$O\left(\frac{\lambda}{B} \log_{\lambda} N\right)$$

Lookups in

$$\Omega(\log_{\lambda} N)$$



B-Trees

$$(\lambda = B)$$

Comparison External Memory Model

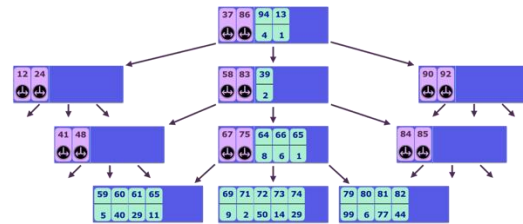
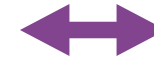
Iacono-Pătrașcu Lower Bound

Insertions in

$$O\left(\frac{\lambda}{B}\right)$$

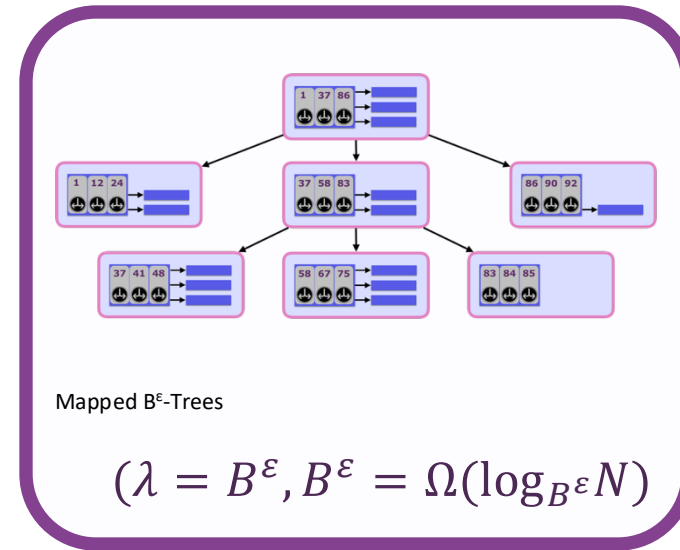
Lookups in

$$\Omega(\log_{\lambda} N)$$



B^{ϵ} -Trees

$$(\lambda = B^{\epsilon})$$



Mapped B^{ϵ} -Trees

$$(\lambda = B^{\epsilon}, B^{\epsilon} = \Omega(\log_{B^{\epsilon}} N))$$

General External Memory Model

Iacono-Patrascu Hash Table

BoA/BoT Hash Table

Lower Bounds

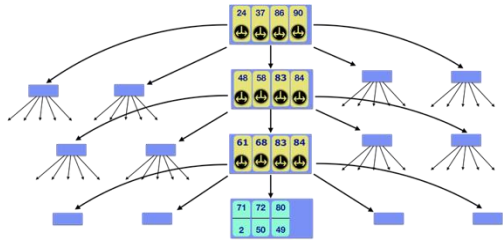
Brodal-Fagerberg Lower Bound

Insertions in

$$O\left(\frac{\lambda}{B} \log_{\lambda} N\right)$$

Lookups in

$$\Omega(\log_{\lambda} N)$$



B-Trees

$$(\lambda = B)$$

Comparison External Memory Model

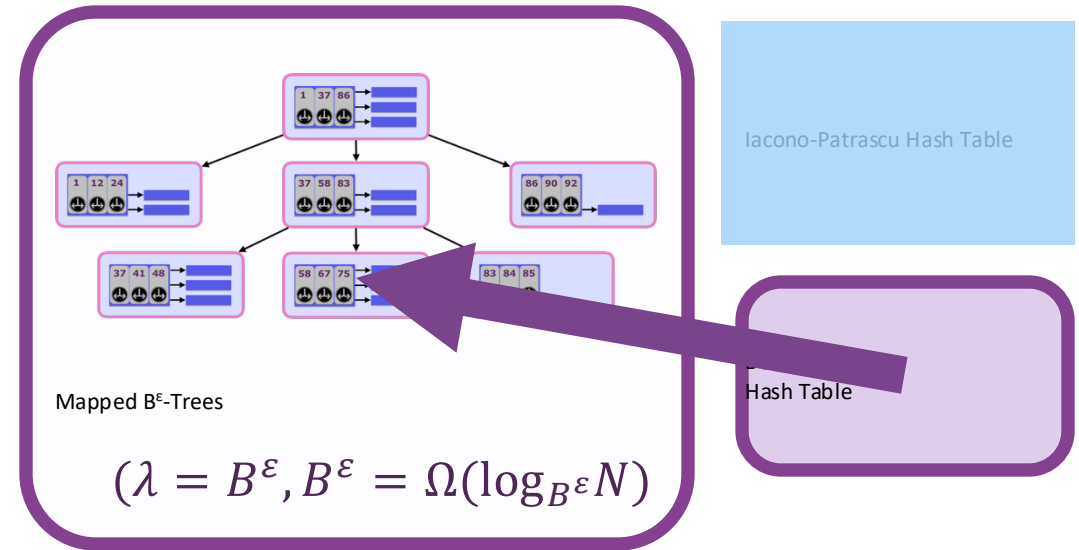
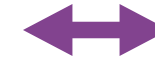
Iacono-Pătrașcu Lower Bound

Insertions in

$$O\left(\frac{\lambda}{B}\right)$$

Lookups in

$$\Omega(\log_{\lambda} N)$$



Mapped B^{ϵ} -Trees

$$(\lambda = B^{\epsilon}, B^{\epsilon} = \Omega(\log_{B^{\epsilon}} N))$$

General External Memory Model

I/O Amplification

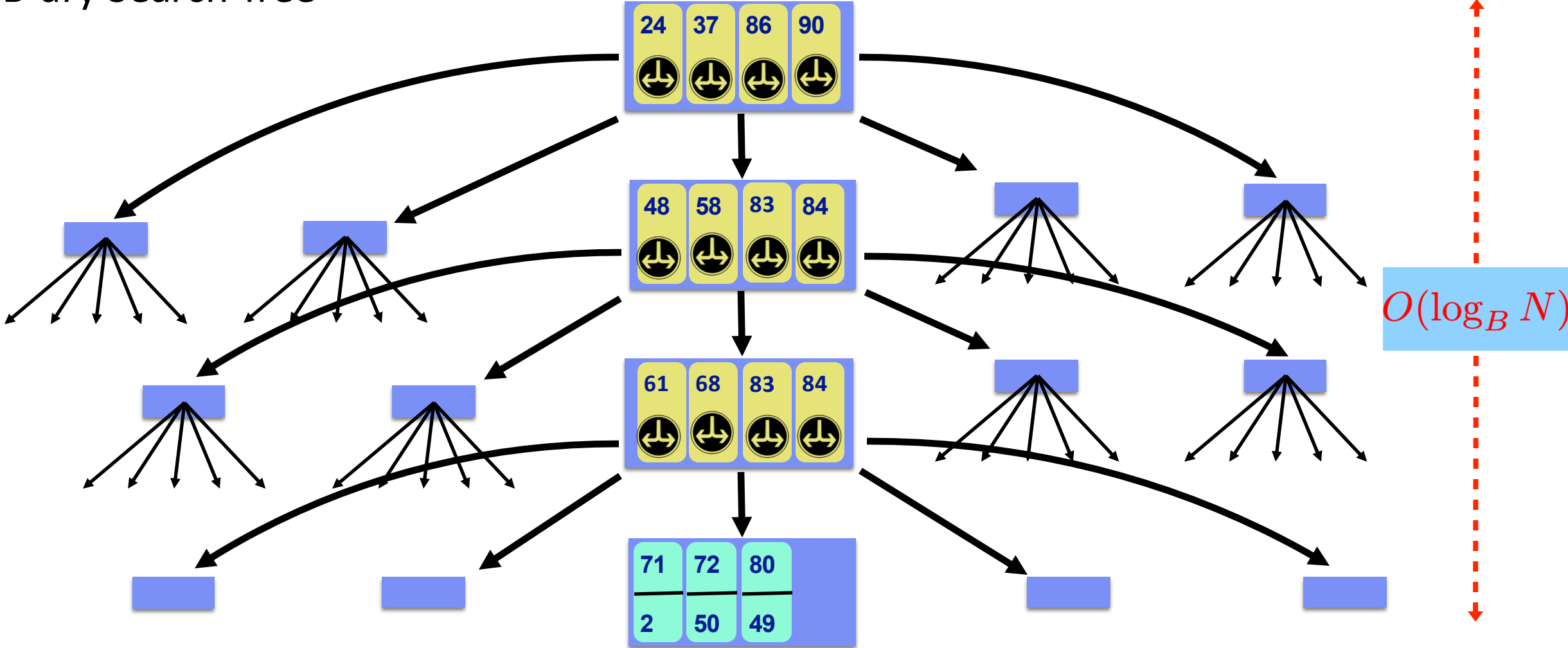
Read amplification is the ratio of the number of blocks read from the disk versus the number of blocks required to read the key-value pair.

Write amplification is the ratio of the number of blocks written to the disk versus the number of blocks required to write the key-value pair.

B-Trees

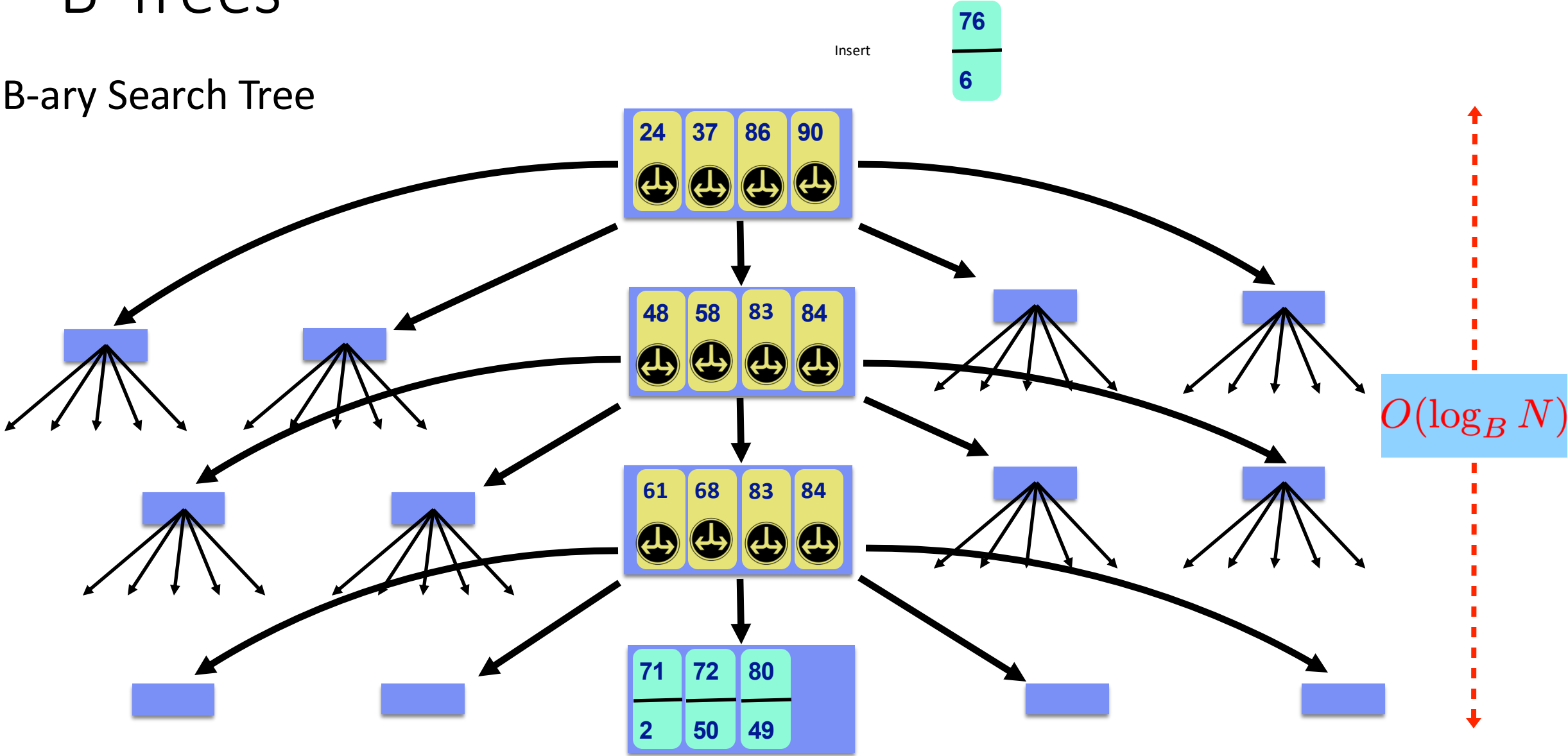
B-Trees

B-ary Search Tree



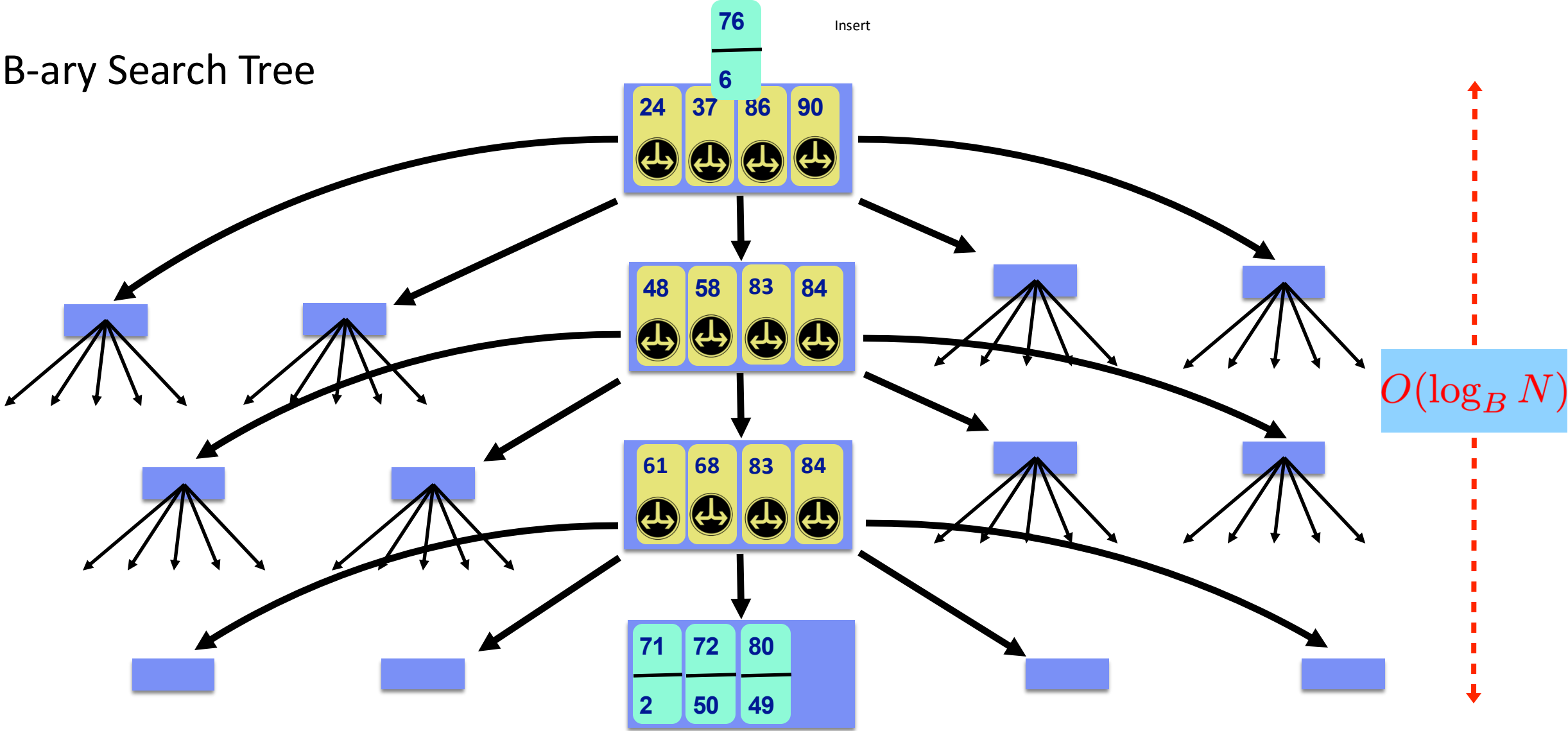
B-Trees

B-ary Search Tree



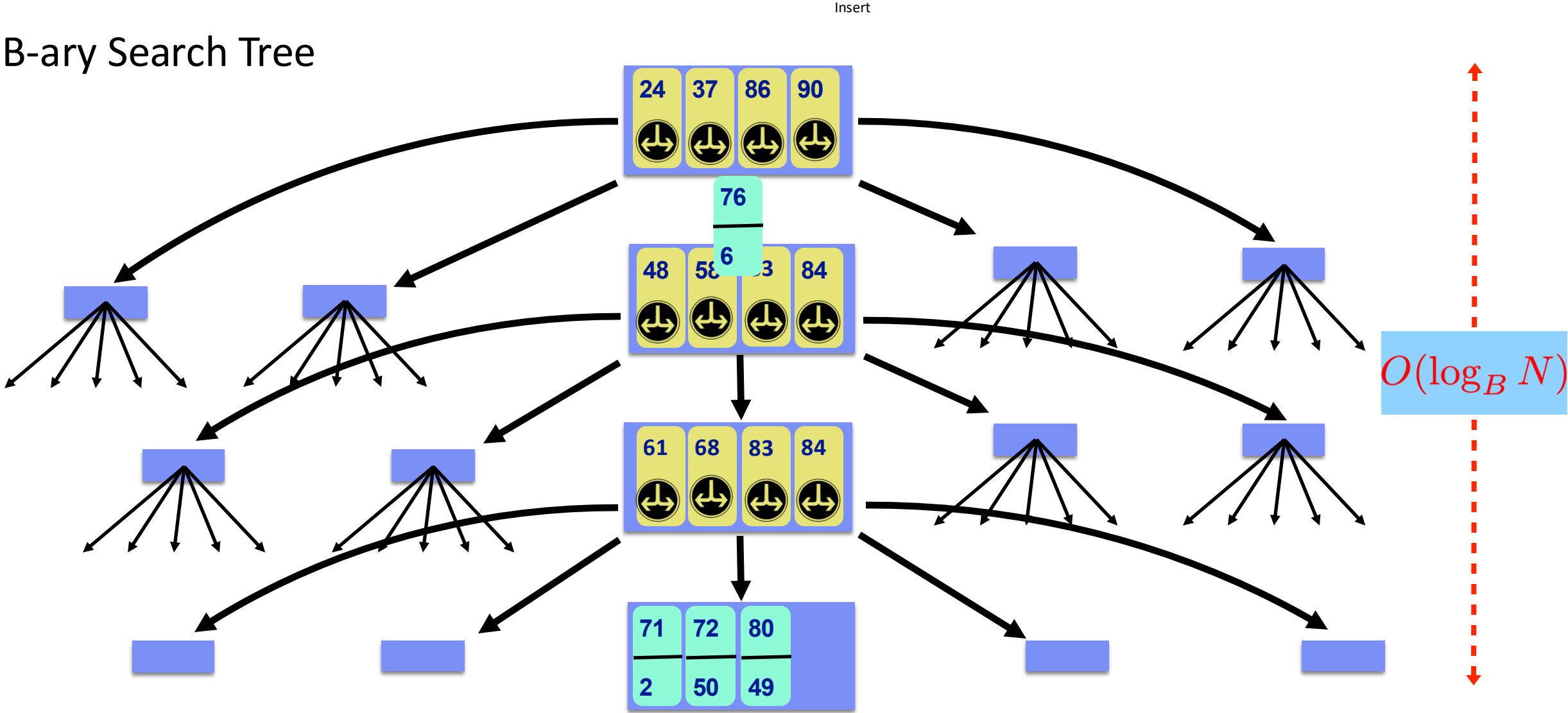
B-Trees

B-ary Search Tree



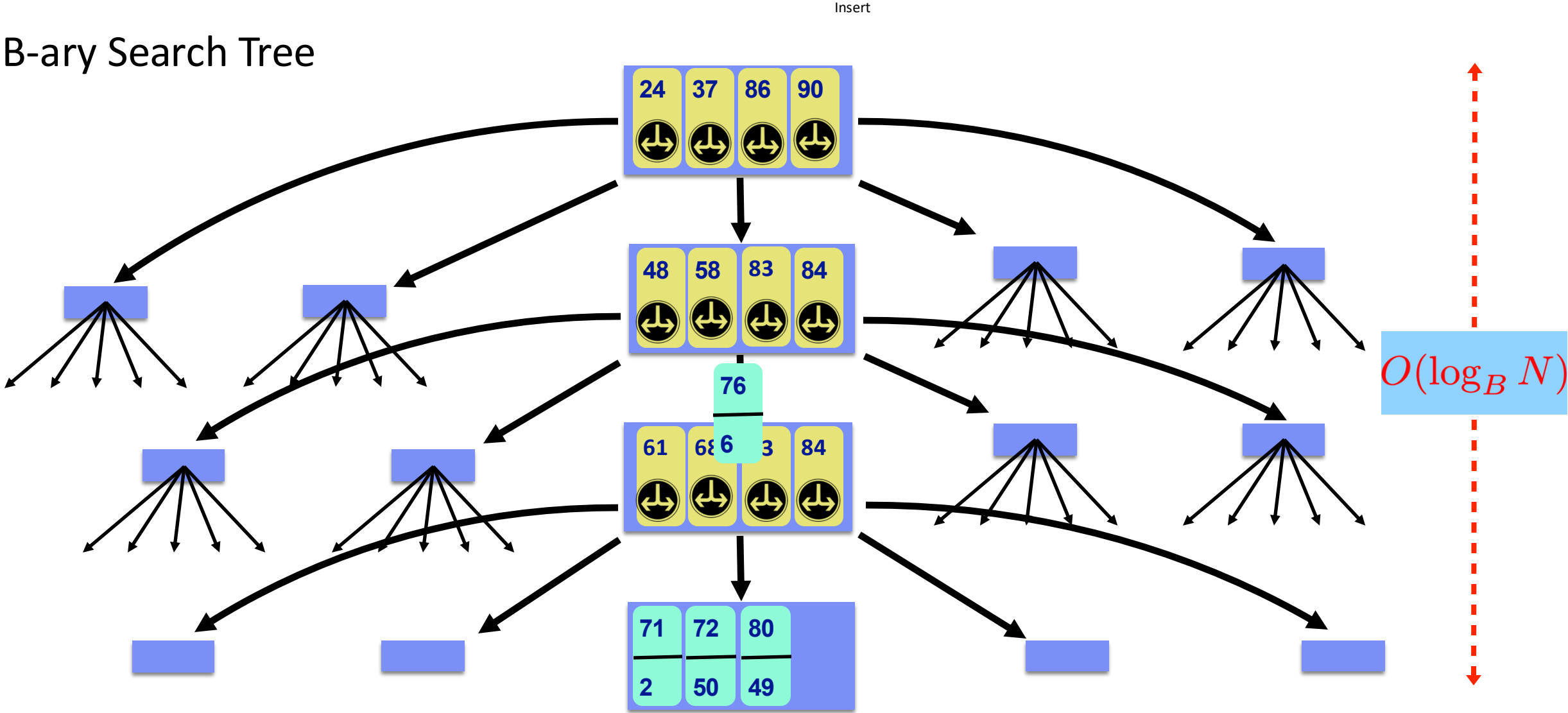
B-Trees

B-ary Search Tree



B-Trees

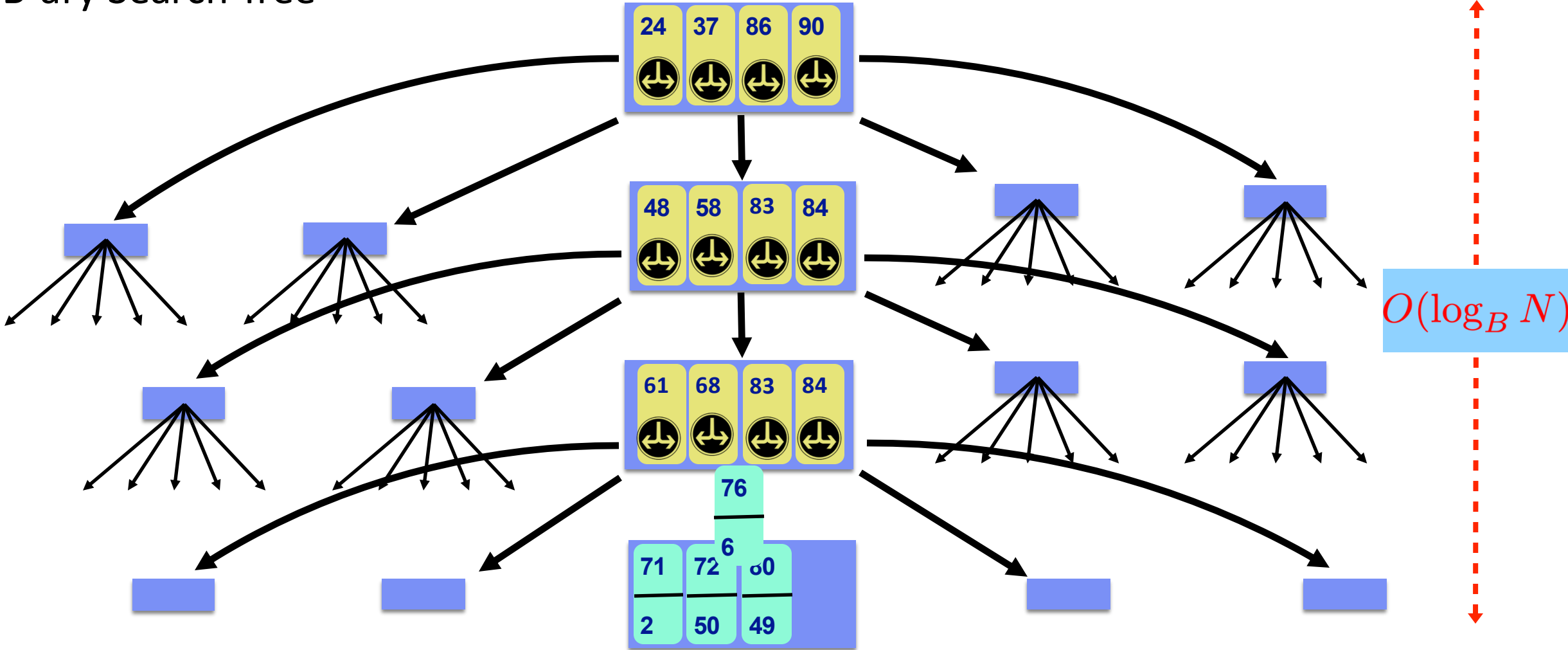
B-ary Search Tree



B-Trees

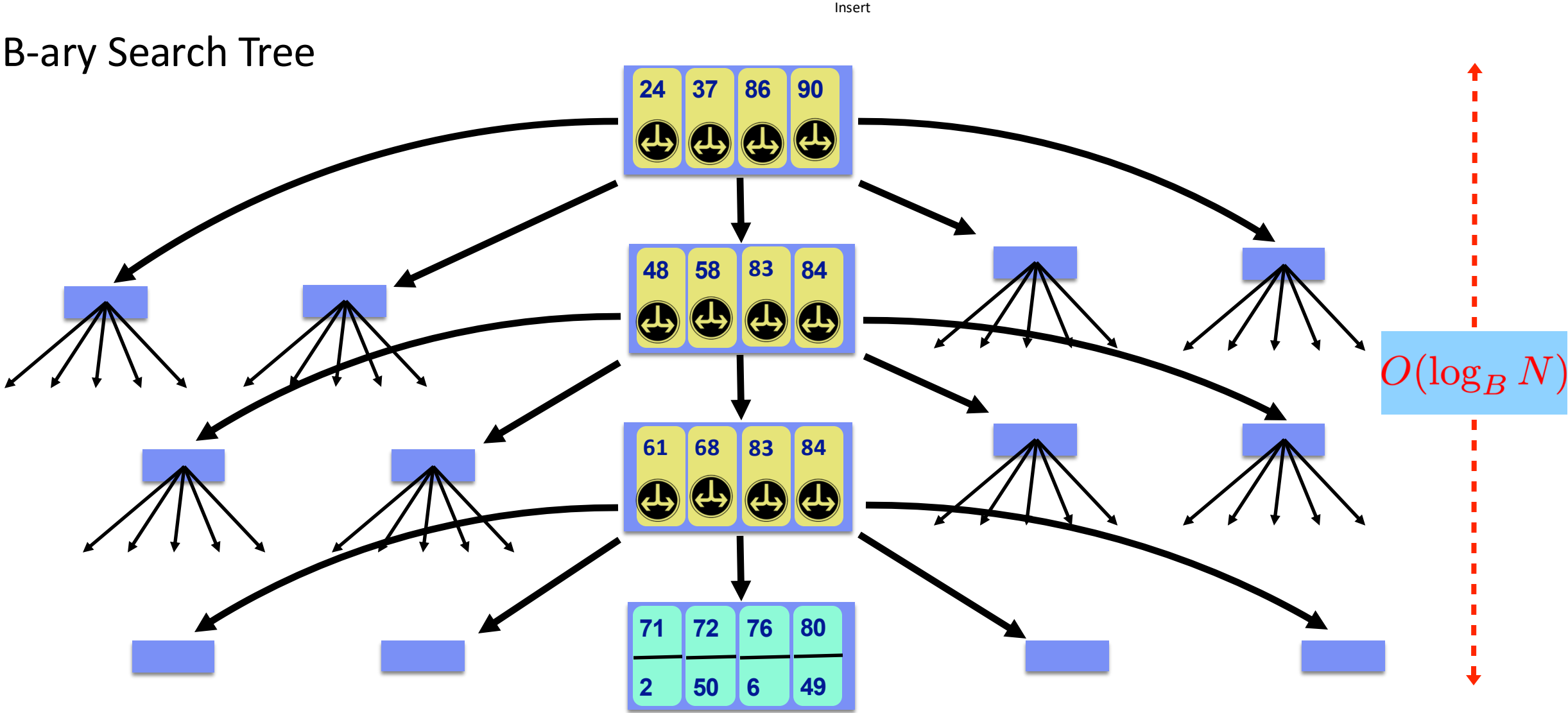
B-ary Search Tree

Insert



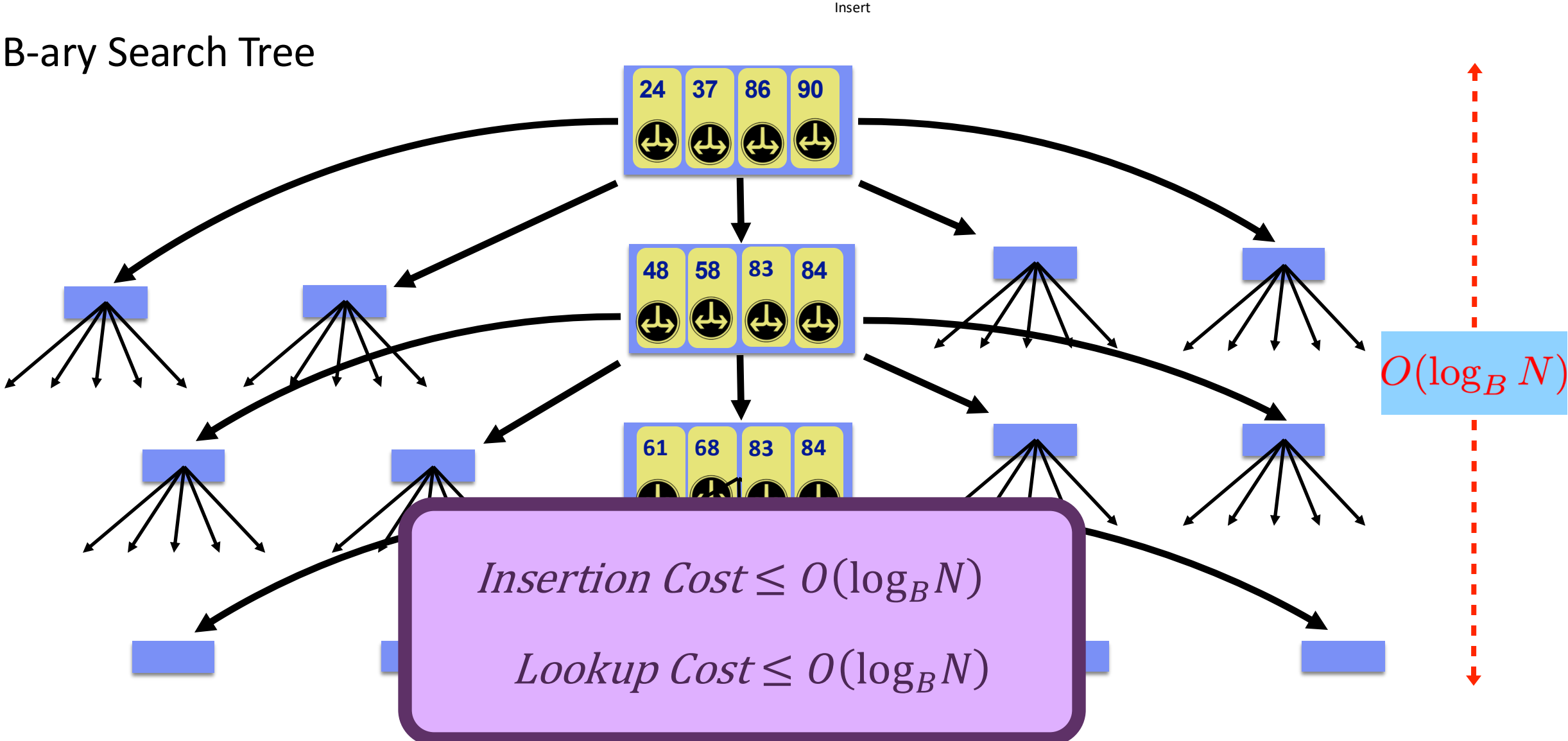
B-Trees

B-ary Search Tree



B-Trees

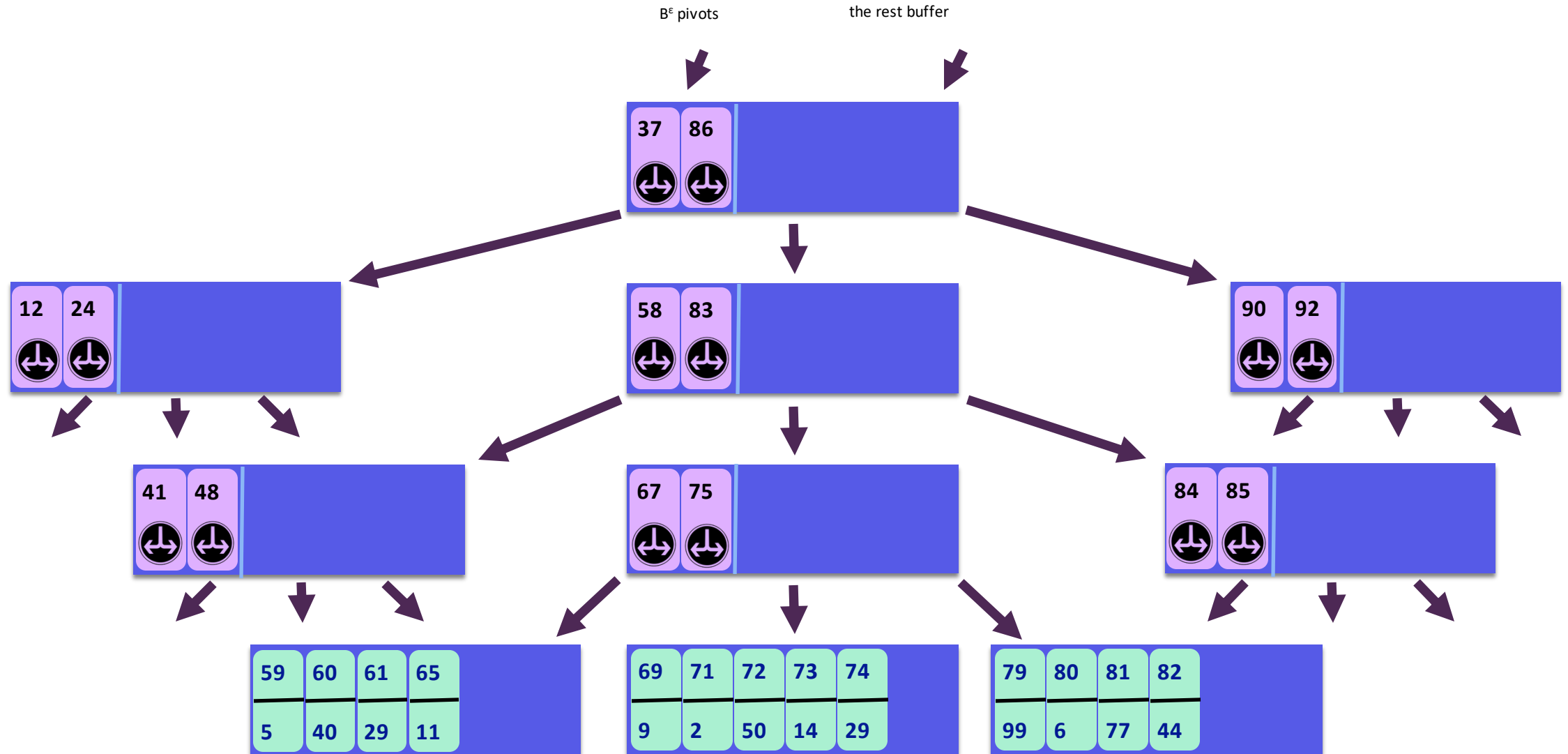
B-ary Search Tree



B^ϵ -Trees

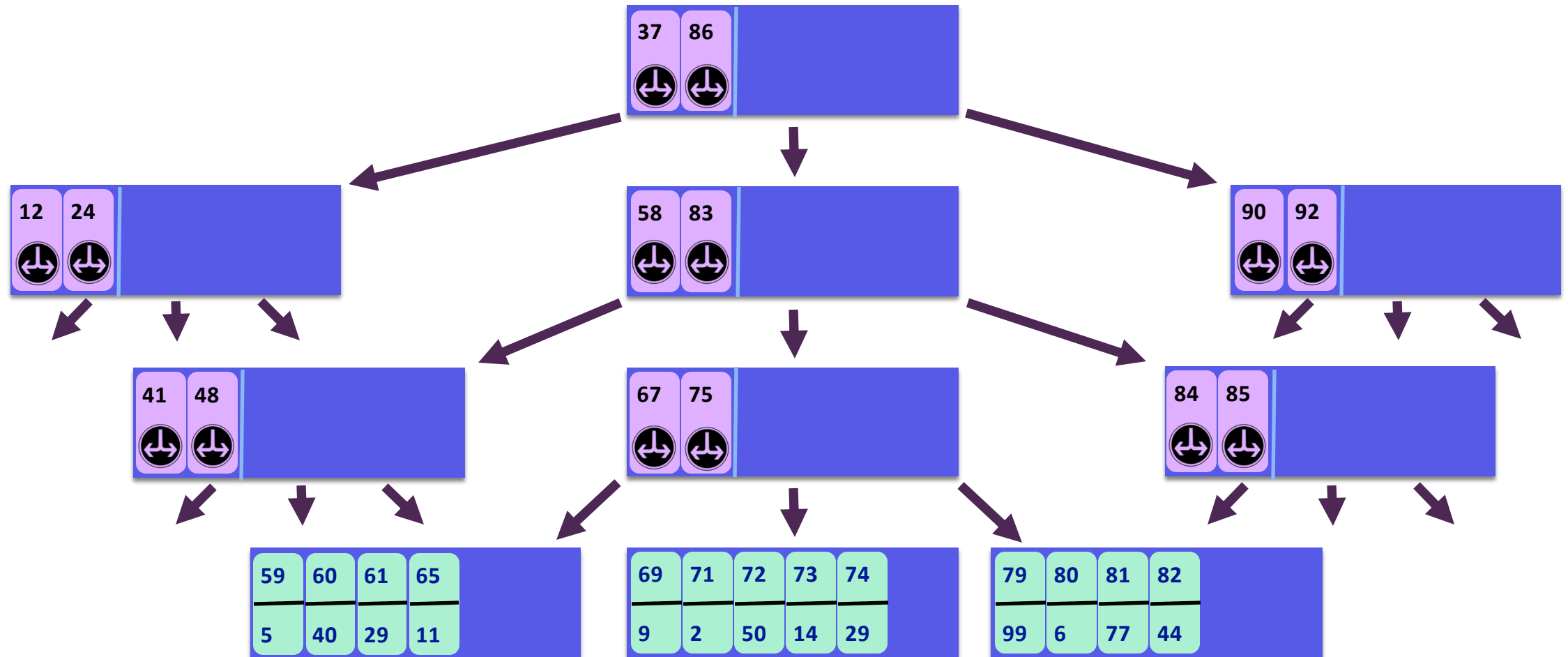
B^ε-Trees

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



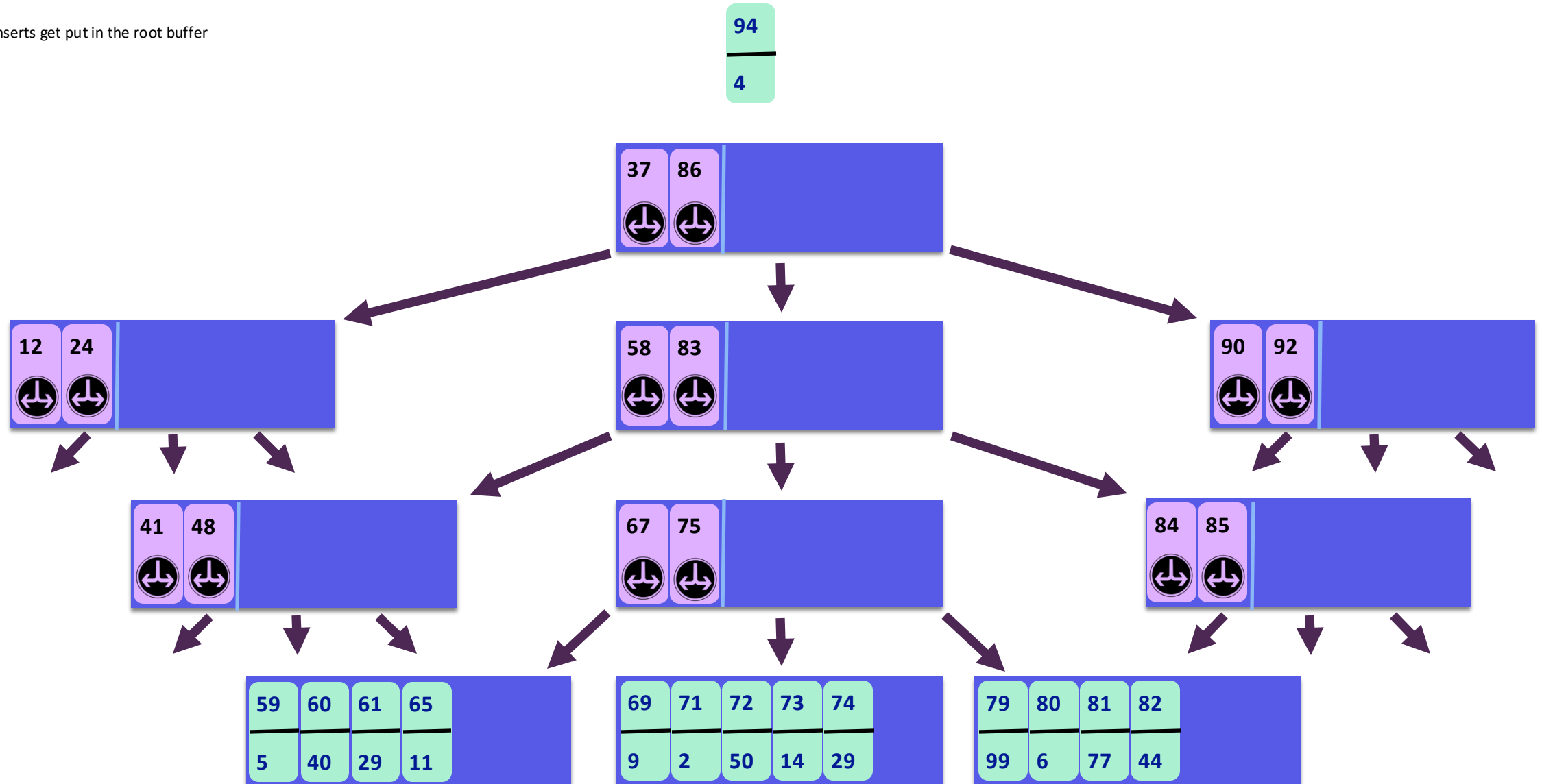
B^ε-Trees

Inserts get put in the root buffer



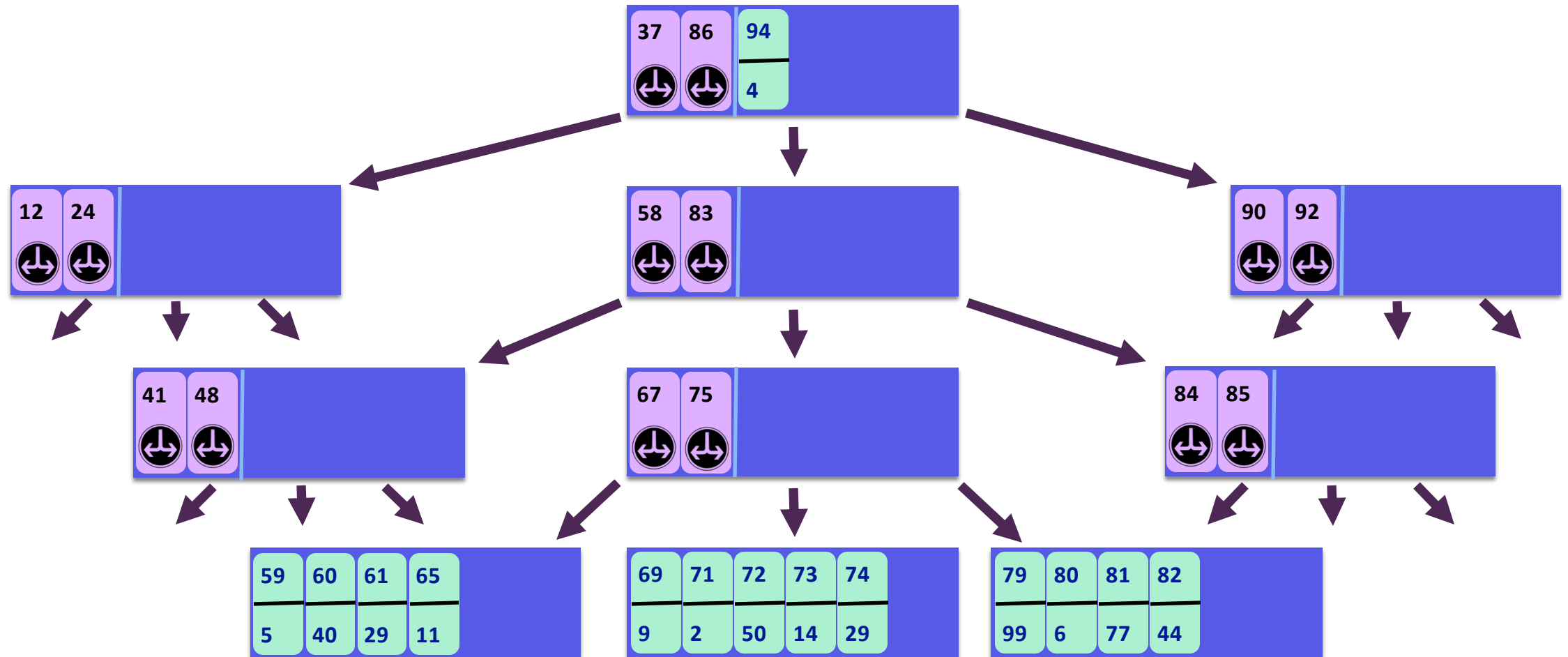
B^ε-Trees

Inserts get put in the root buffer



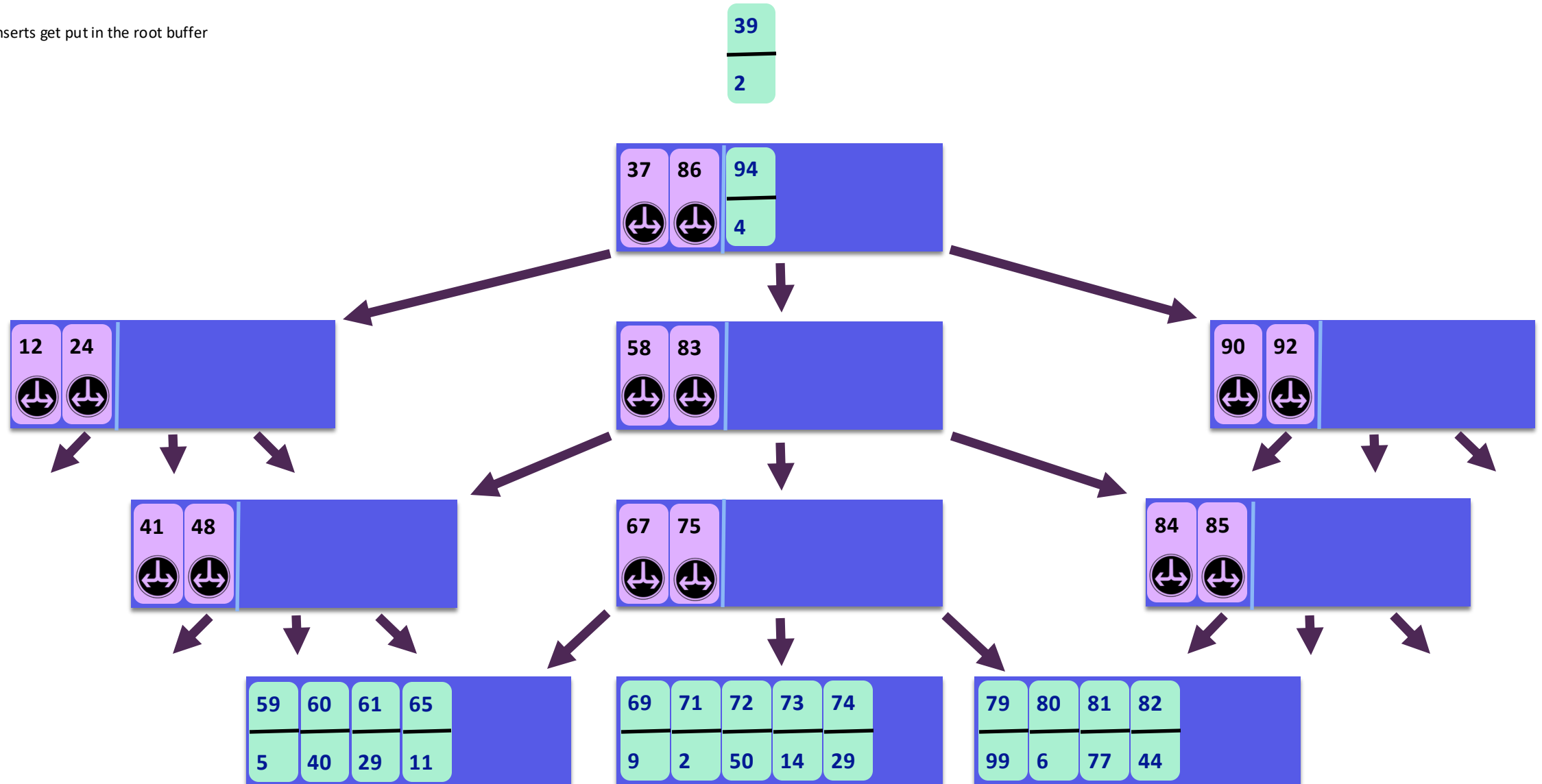
B^ε-Trees

Inserts get put in the root buffer



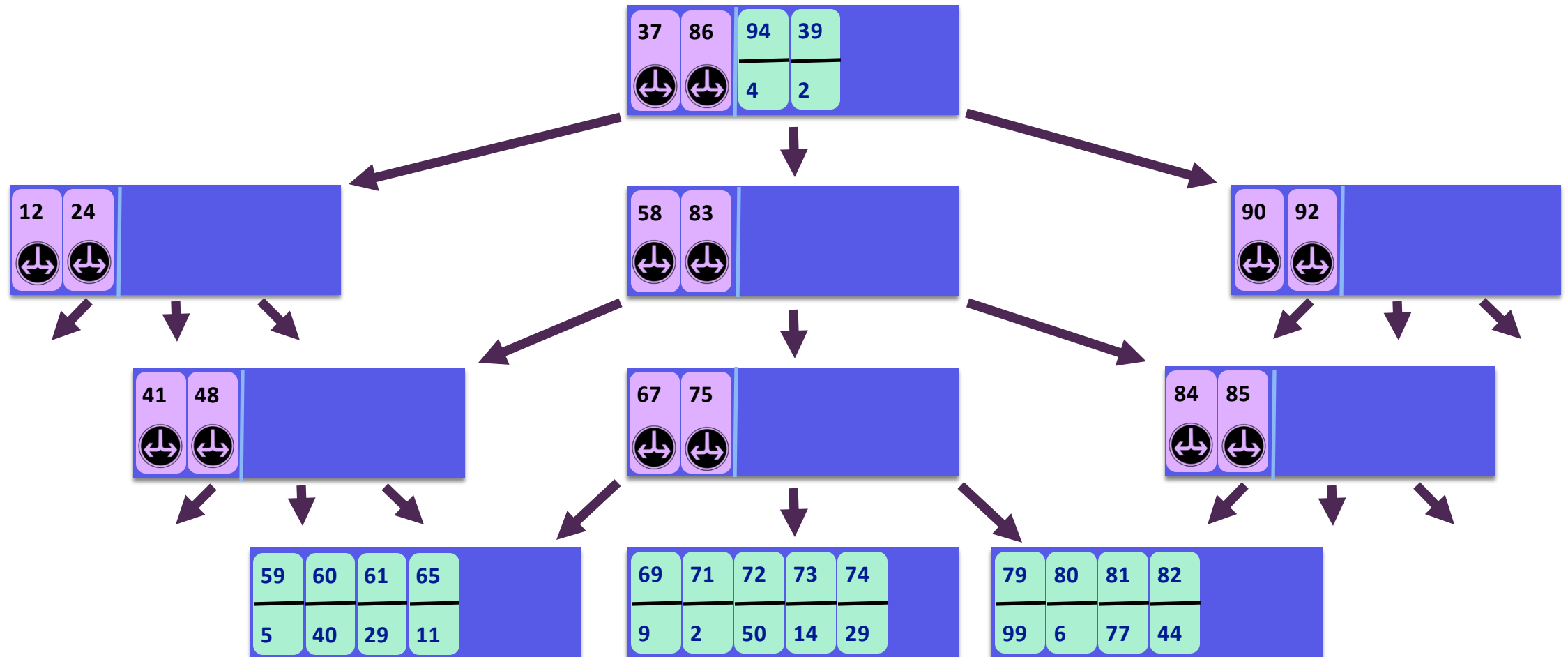
B^ε-Trees

Inserts get put in the root buffer



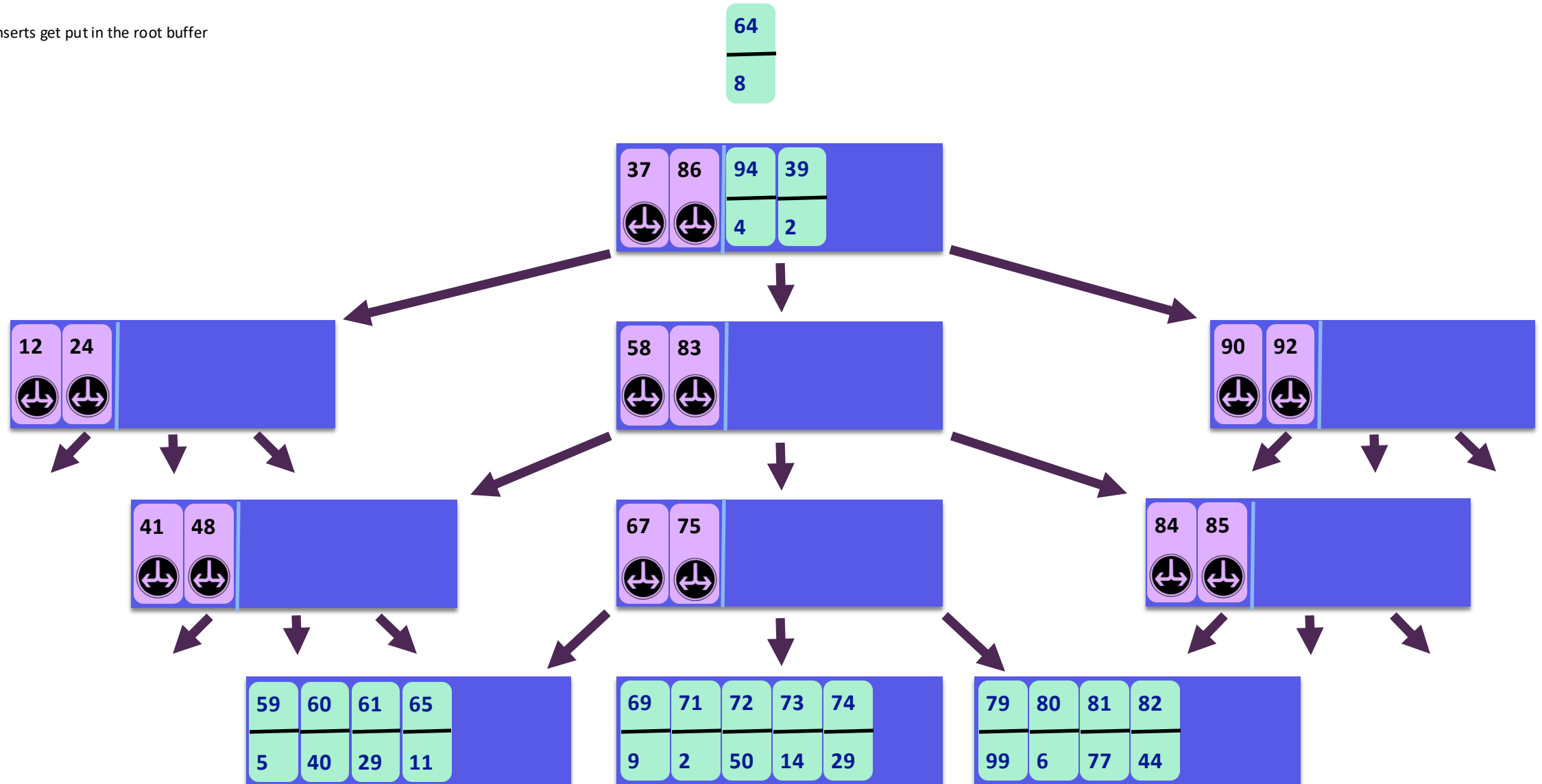
B^ε-Trees

Inserts get put in the root buffer



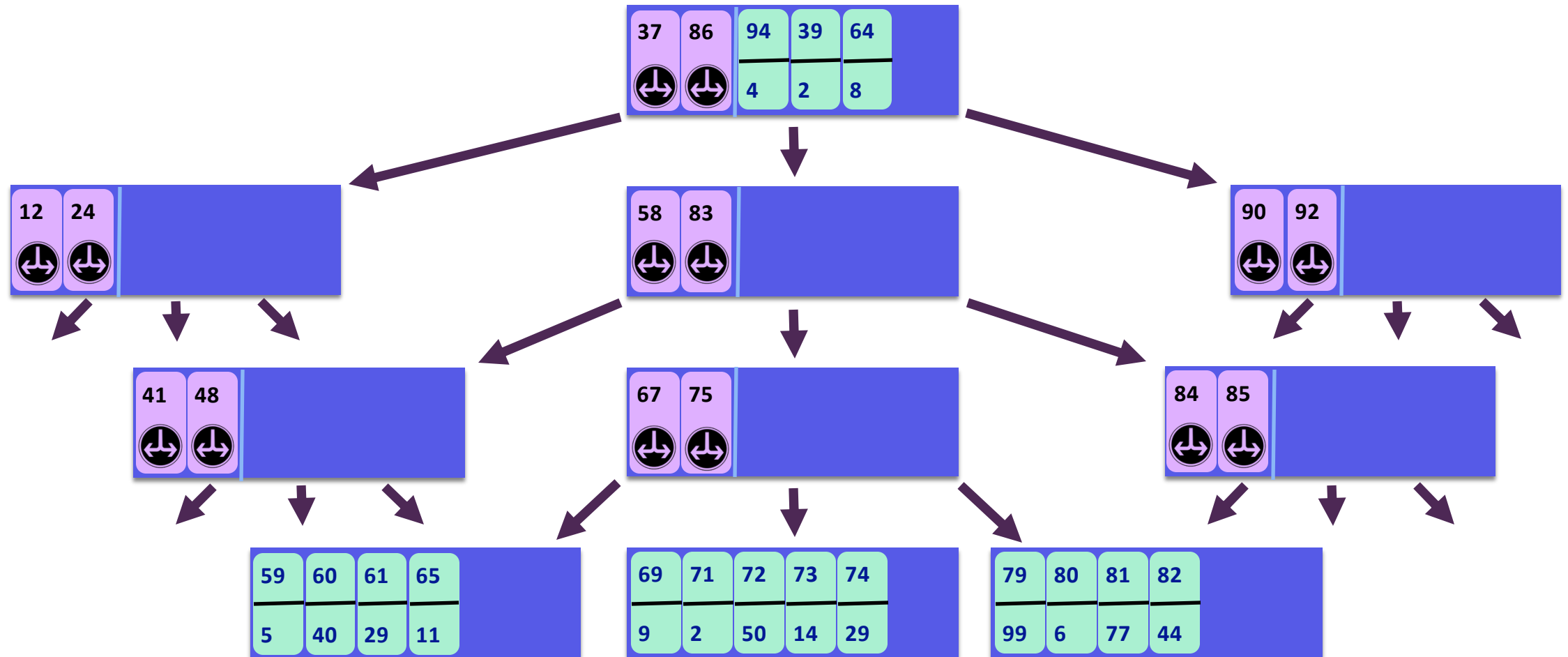
B^ε-Trees

Inserts get put in the root buffer



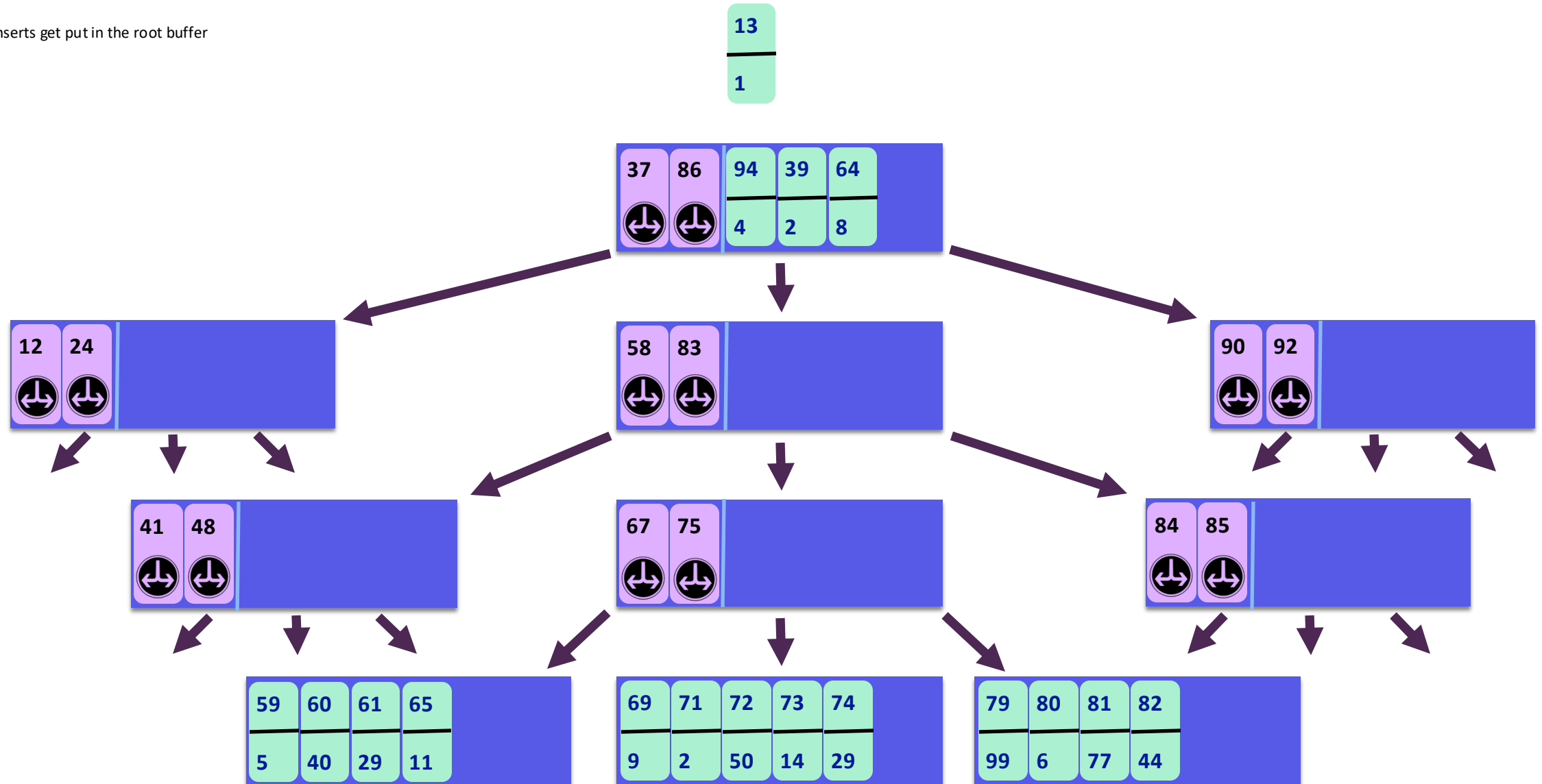
B^ε-Trees

Inserts get put in the root buffer



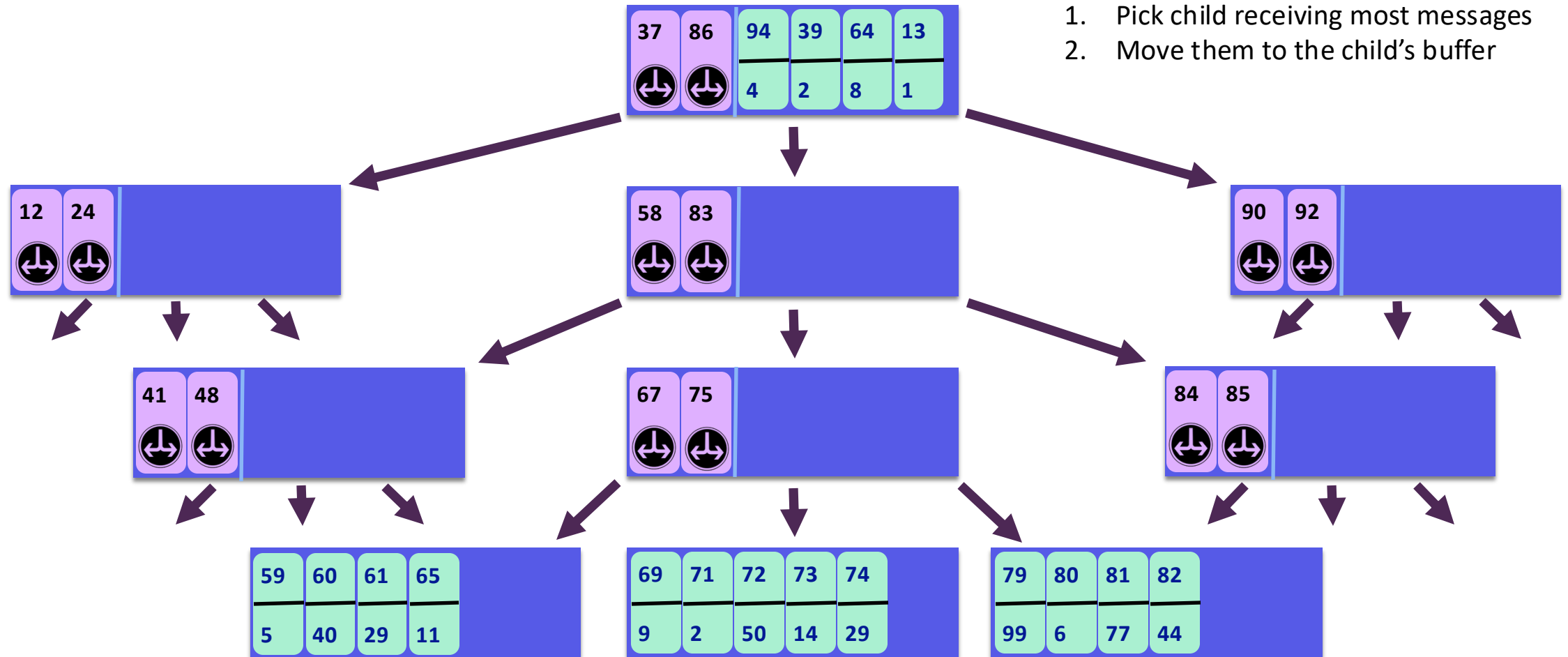
B^ε-Trees

Inserts get put in the root buffer



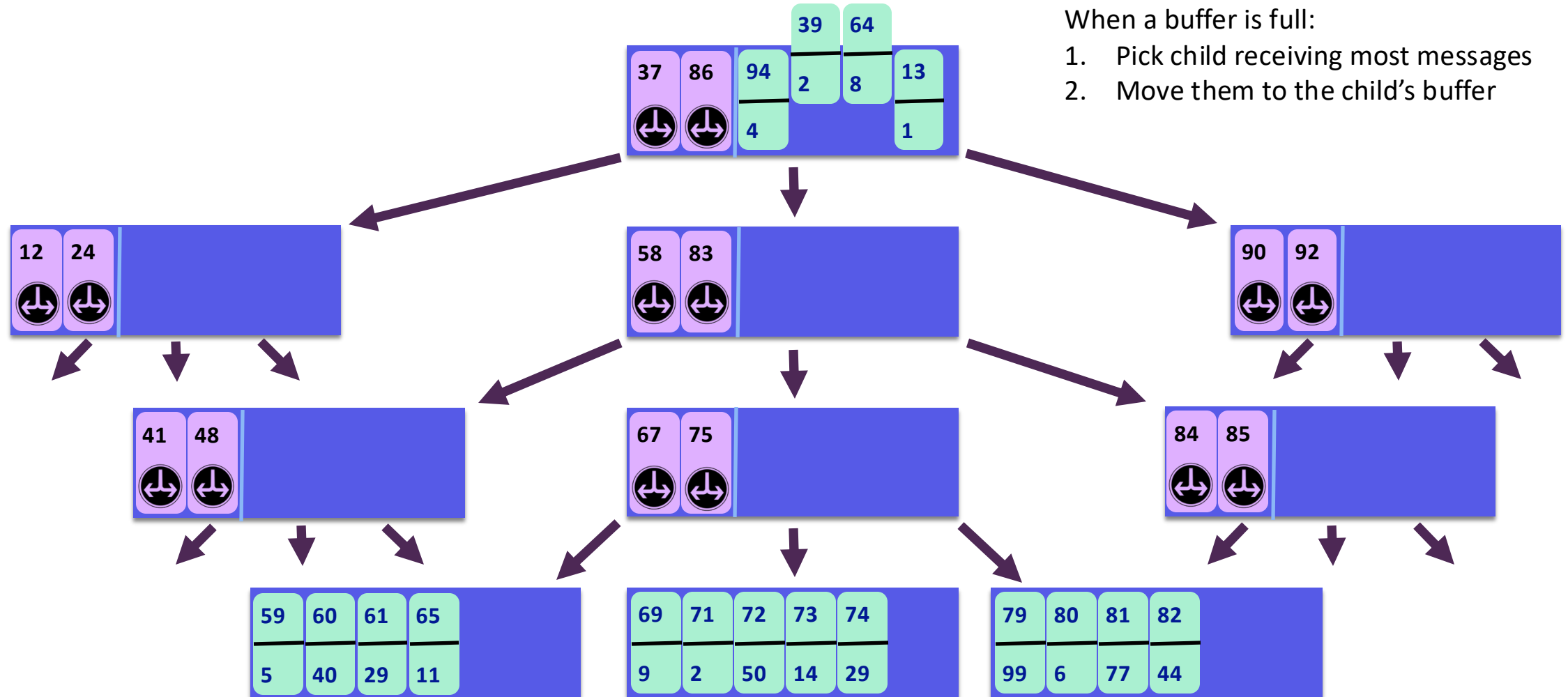
B^ε-Trees

Inserts get put in the root buffer



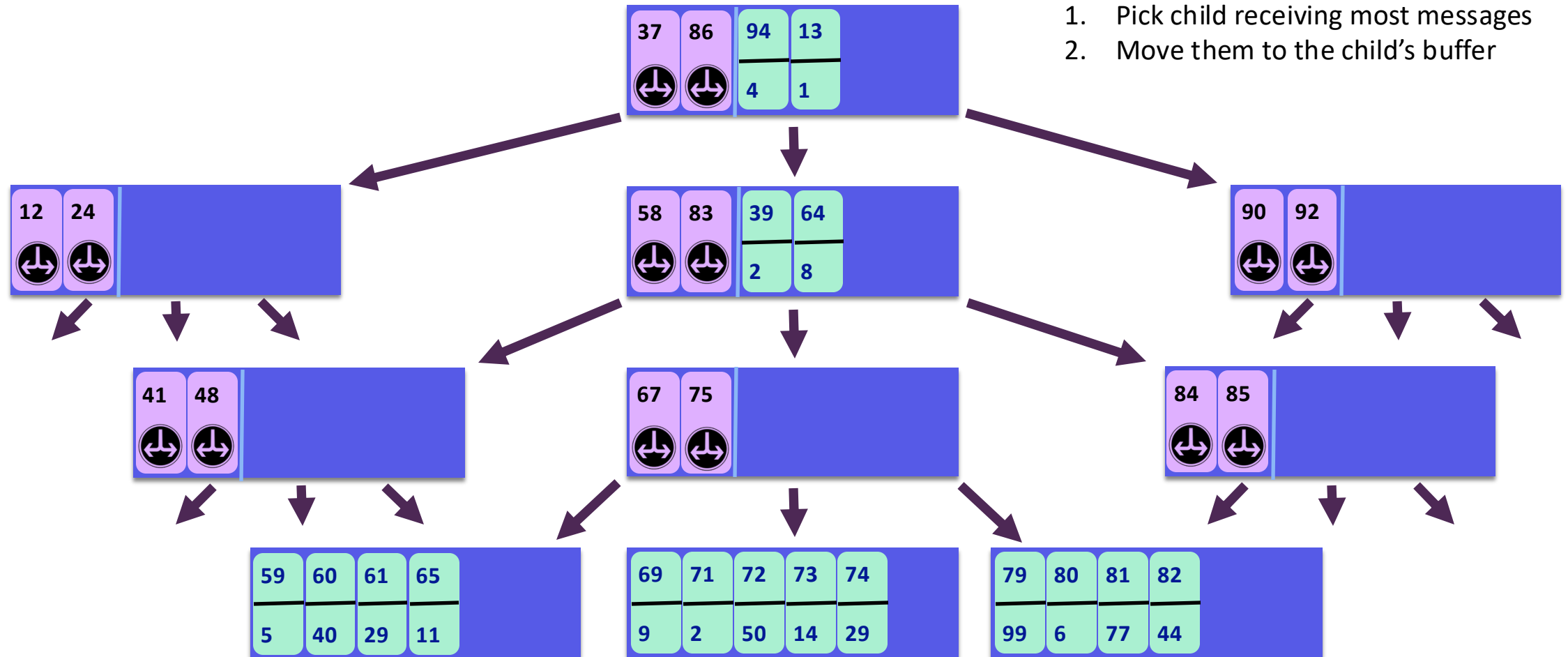
B^ε-Trees

Inserts get put in the root buffer



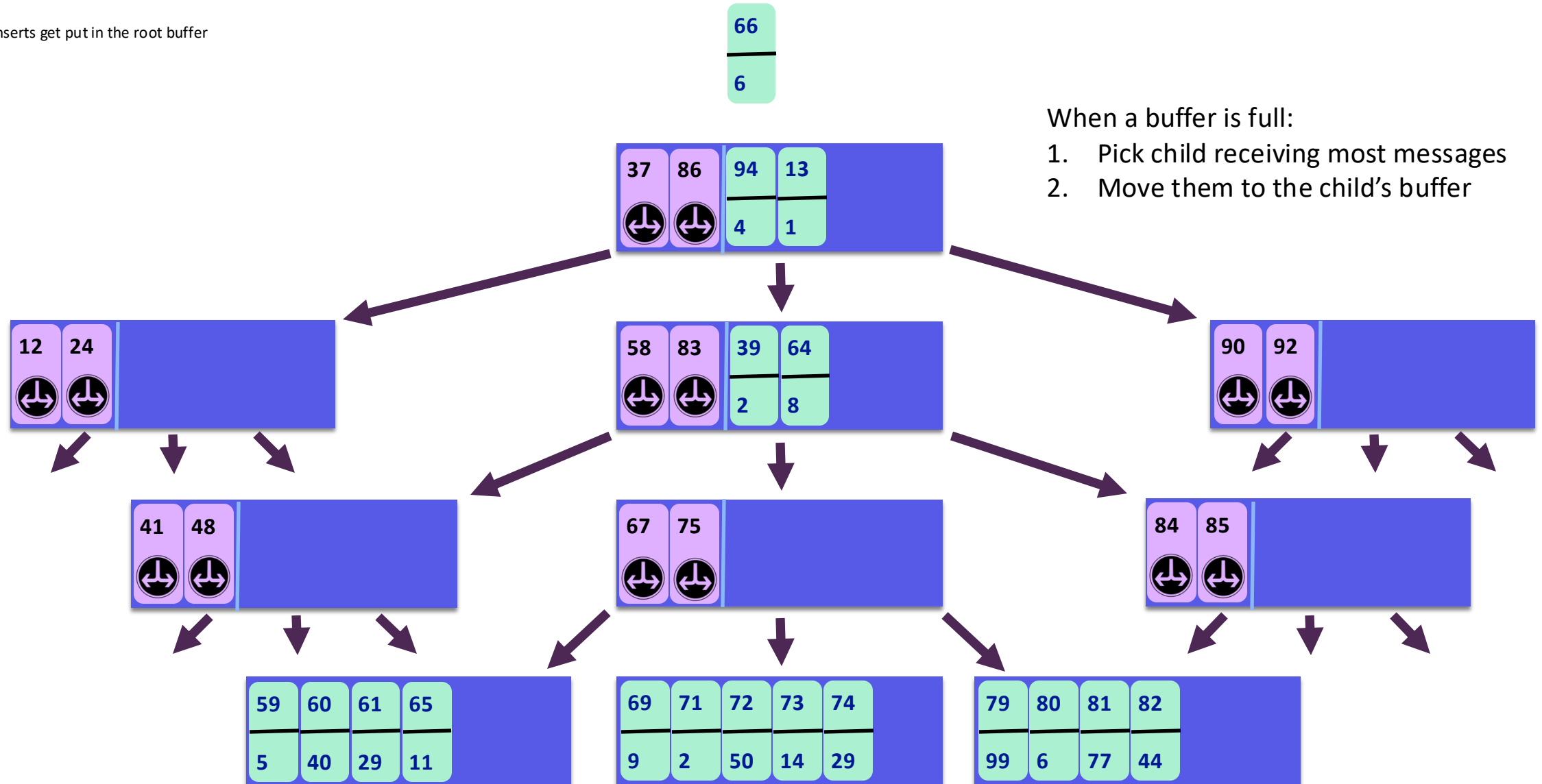
B^ε-Trees

Inserts get put in the root buffer



B^ε-Trees

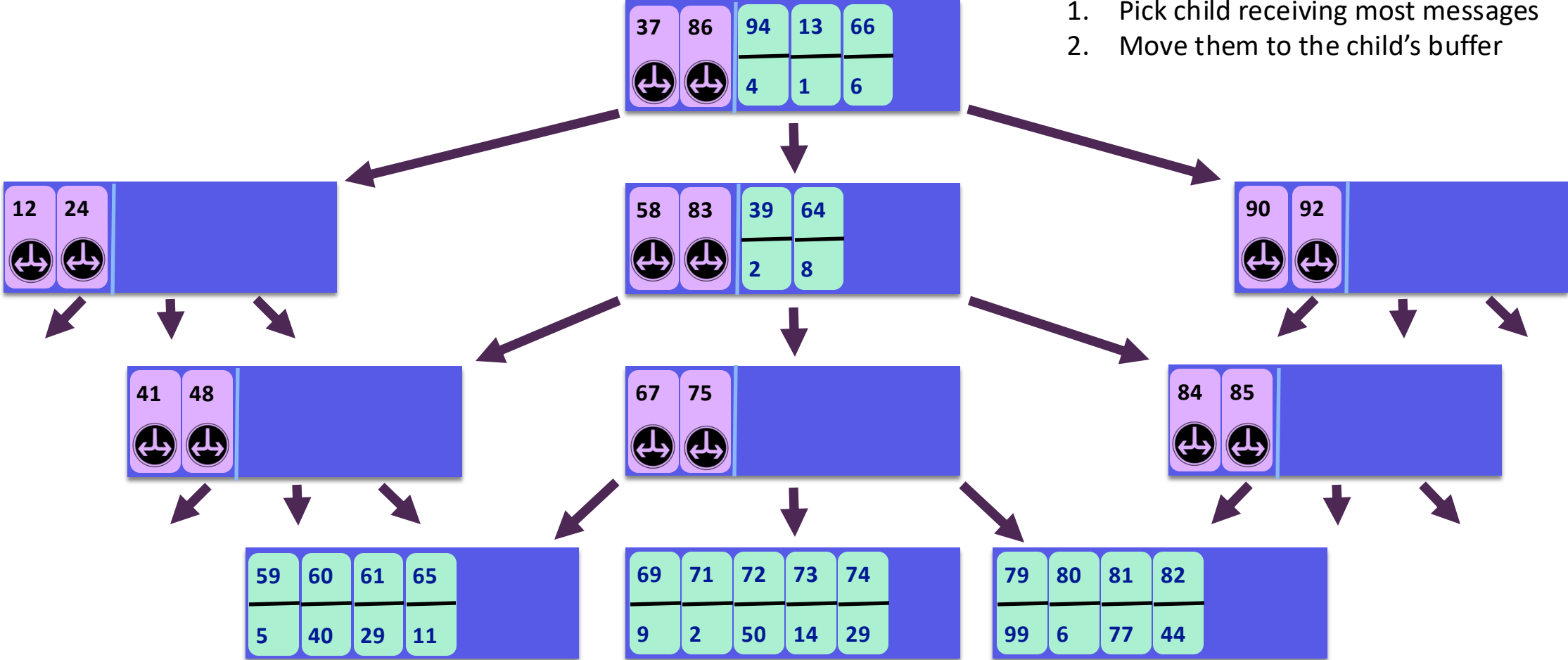
Inserts get put in the root buffer



B-Trees

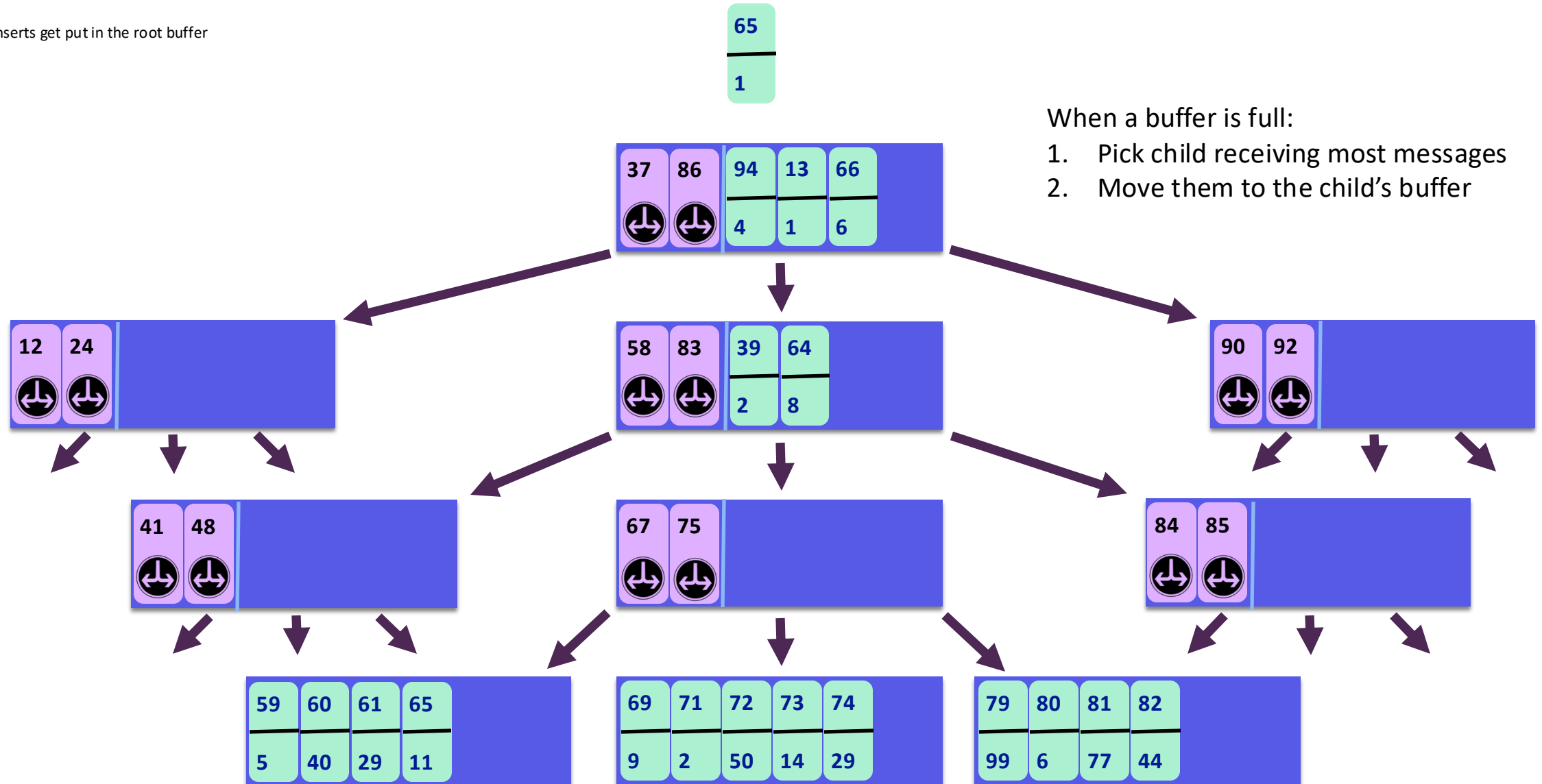
Inserts get put in the root buffer

- 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

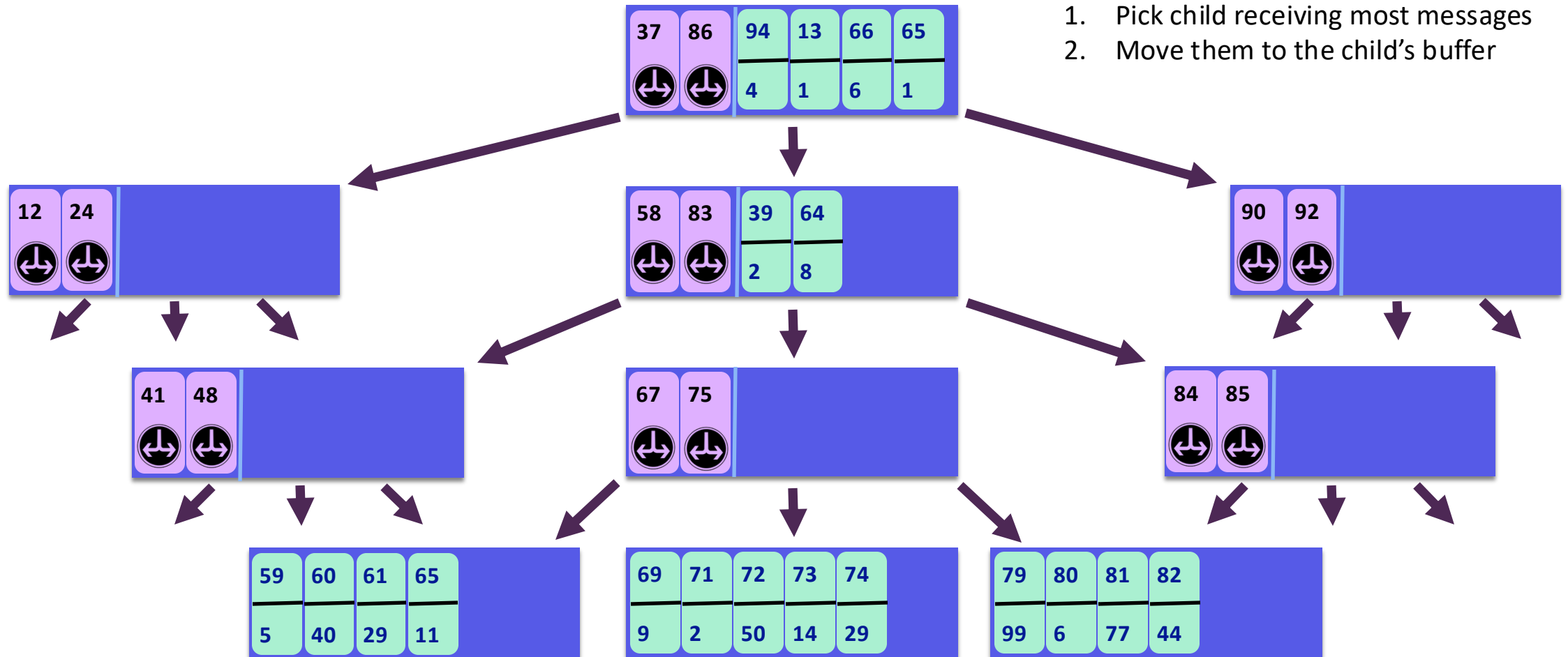


B^ε-Trees

Inserts get put in the root buffer

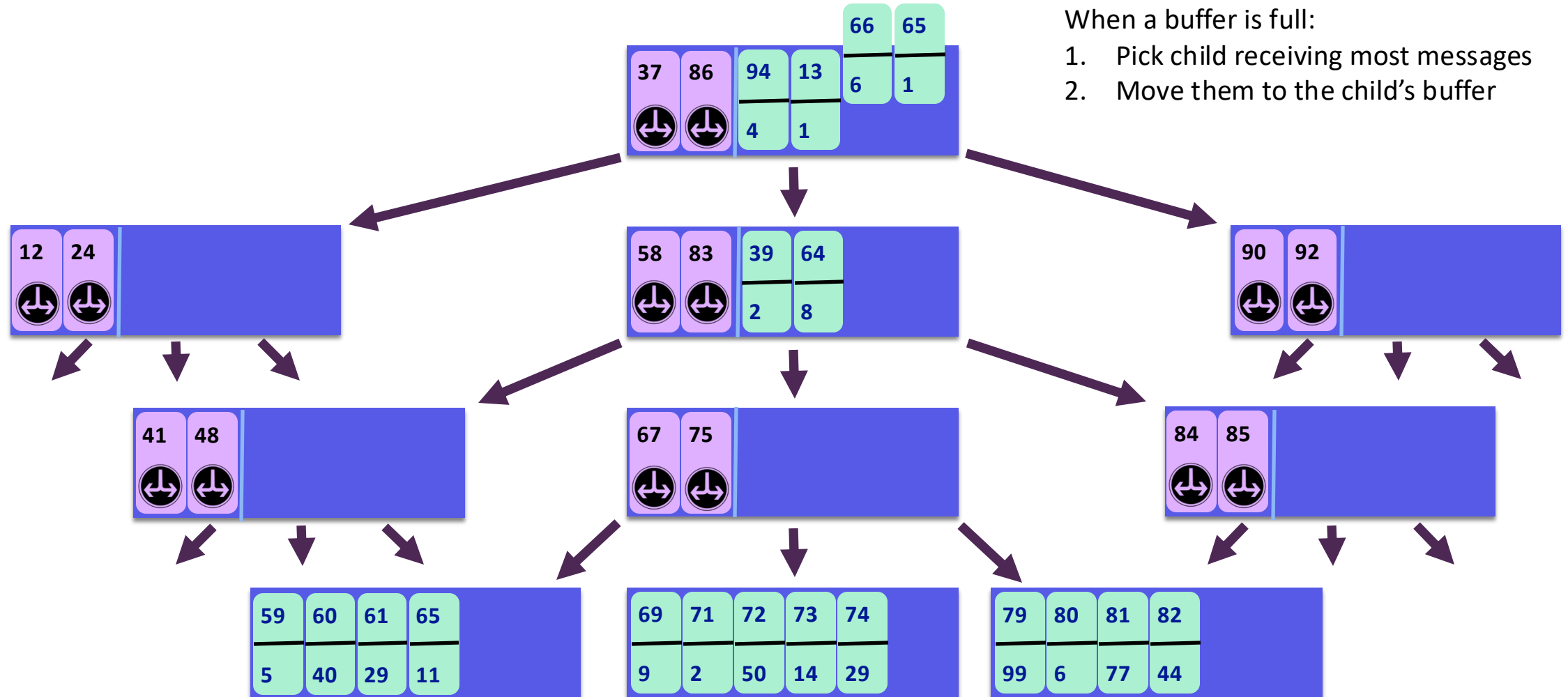
When a buffer is full:

1. Pick child receiving most messages
2. Move them to the child's buffer



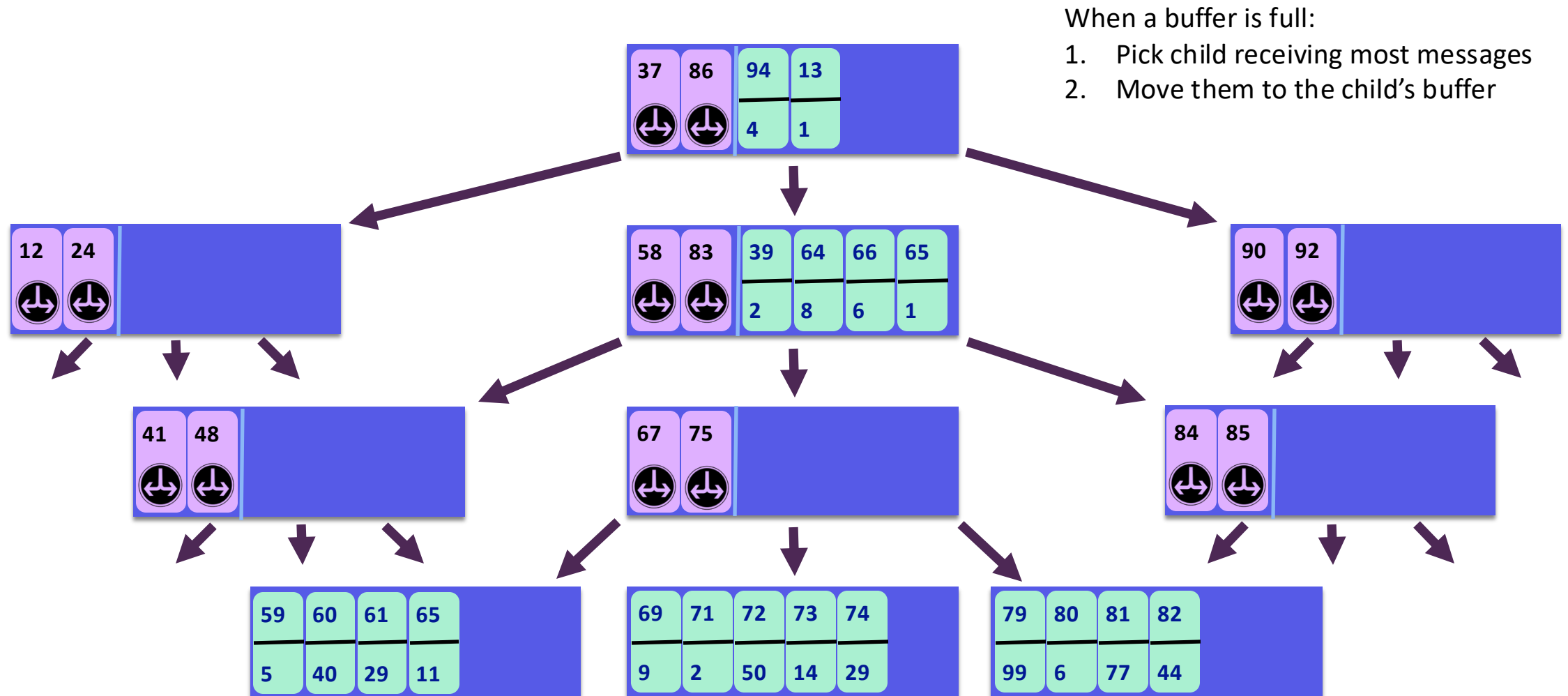
B^ε-Trees

Inserts get put in the root buffer



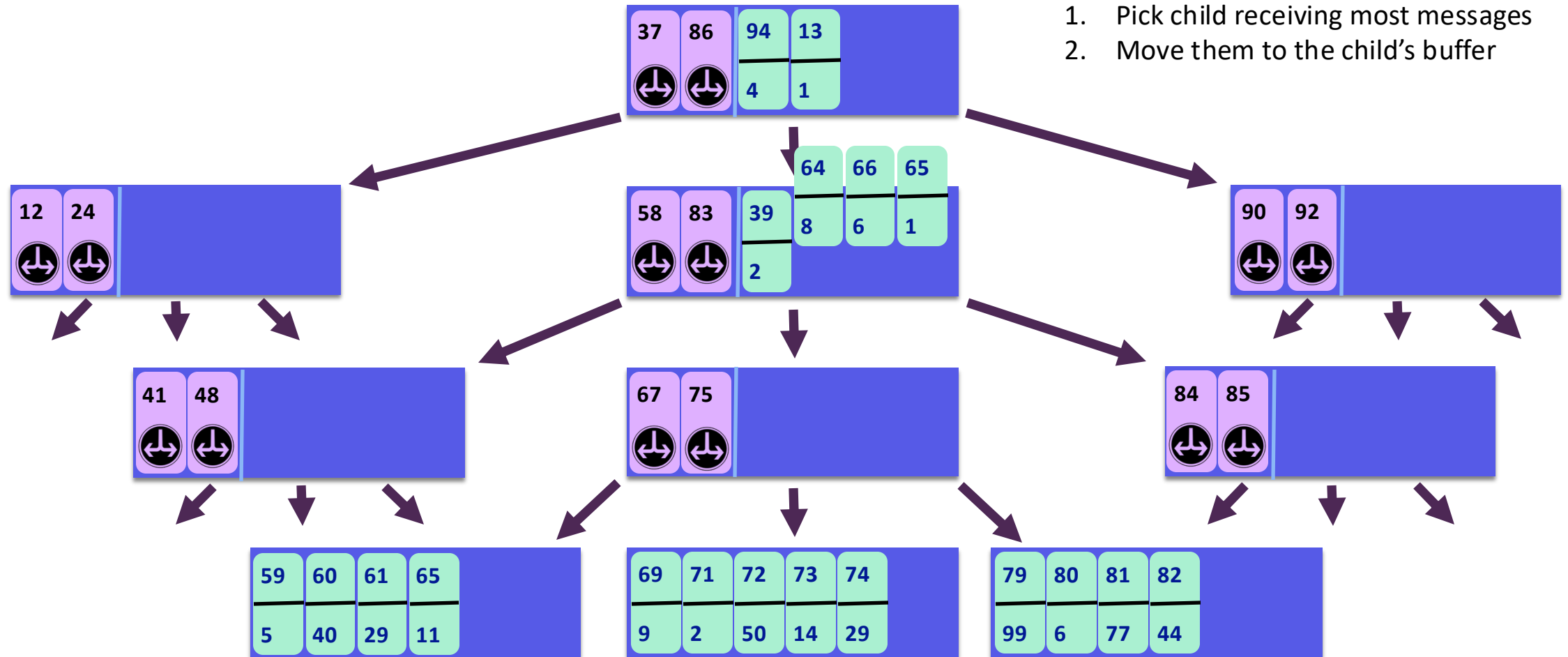
B^ε-Trees

Inserts get put in the root buffer



B^ε-Trees

Inserts get put in the root buffer



When a buffer is full:

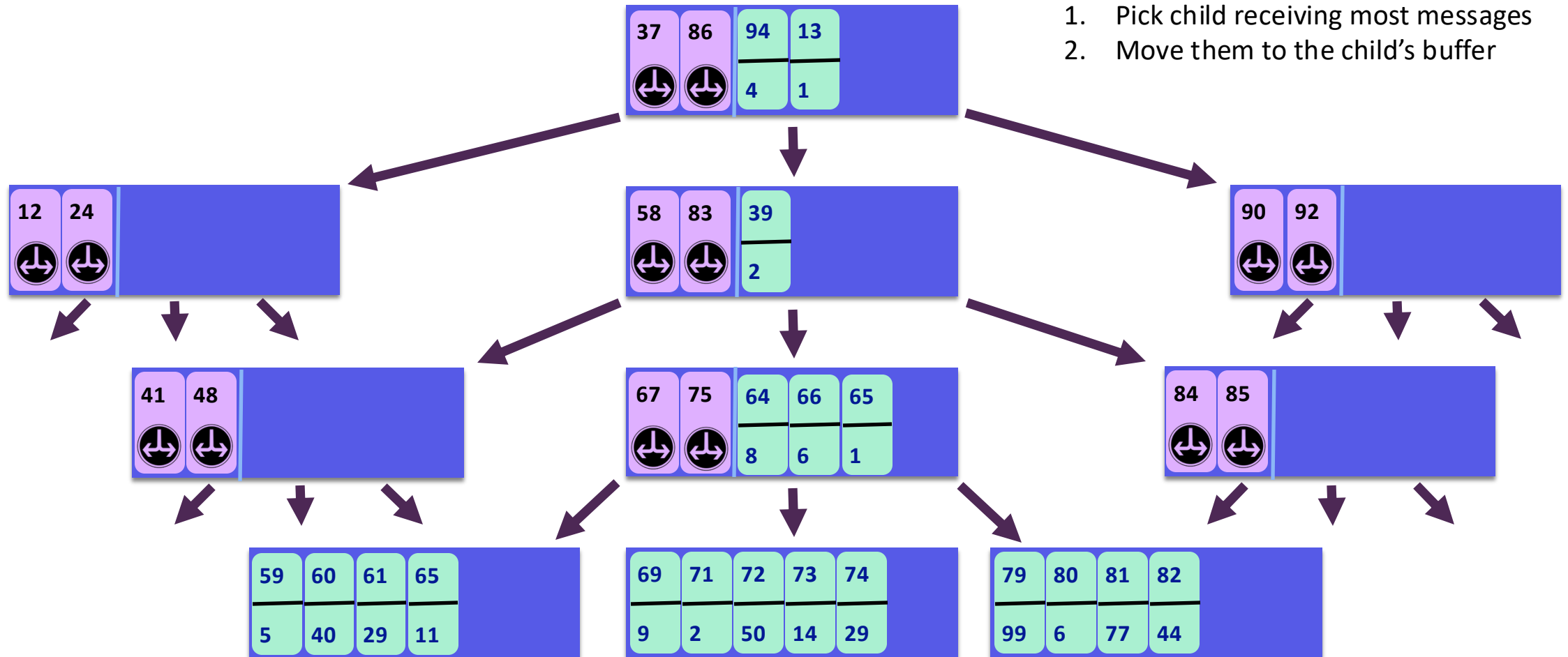
1. Pick child receiving most messages
2. Move them to the child's buffer

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

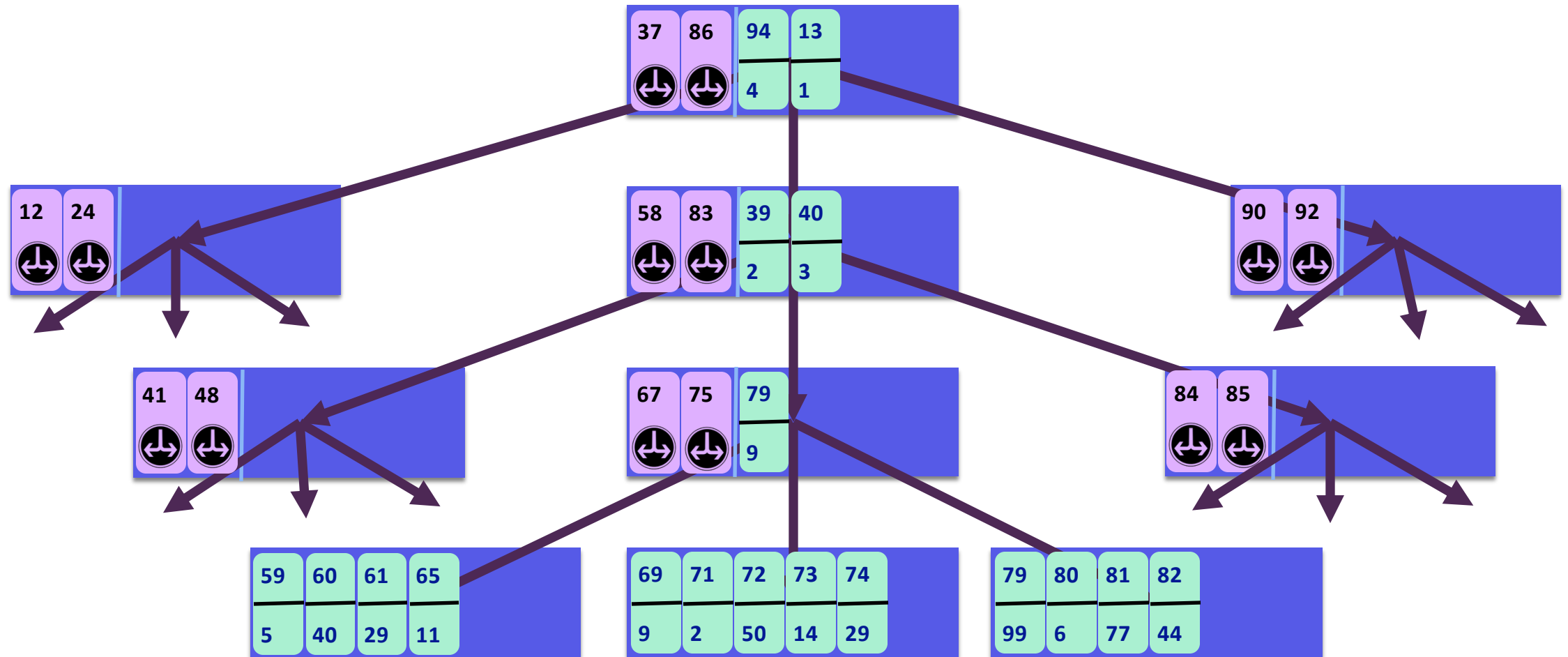


Lookups in B^ϵ -Trees

B^ε-Trees

Query(71)

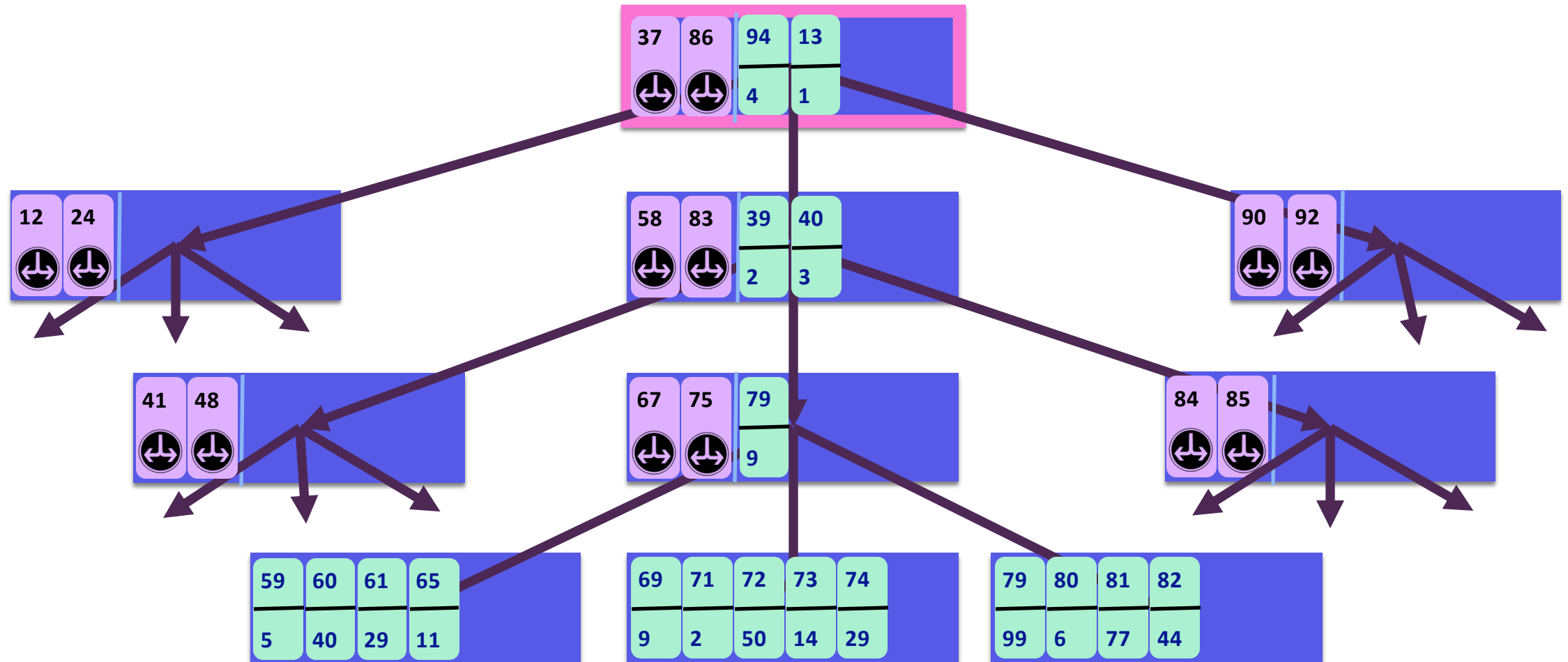
Lookups follow pivots, but check buffers along the way



B^ε-Trees

Query(71)

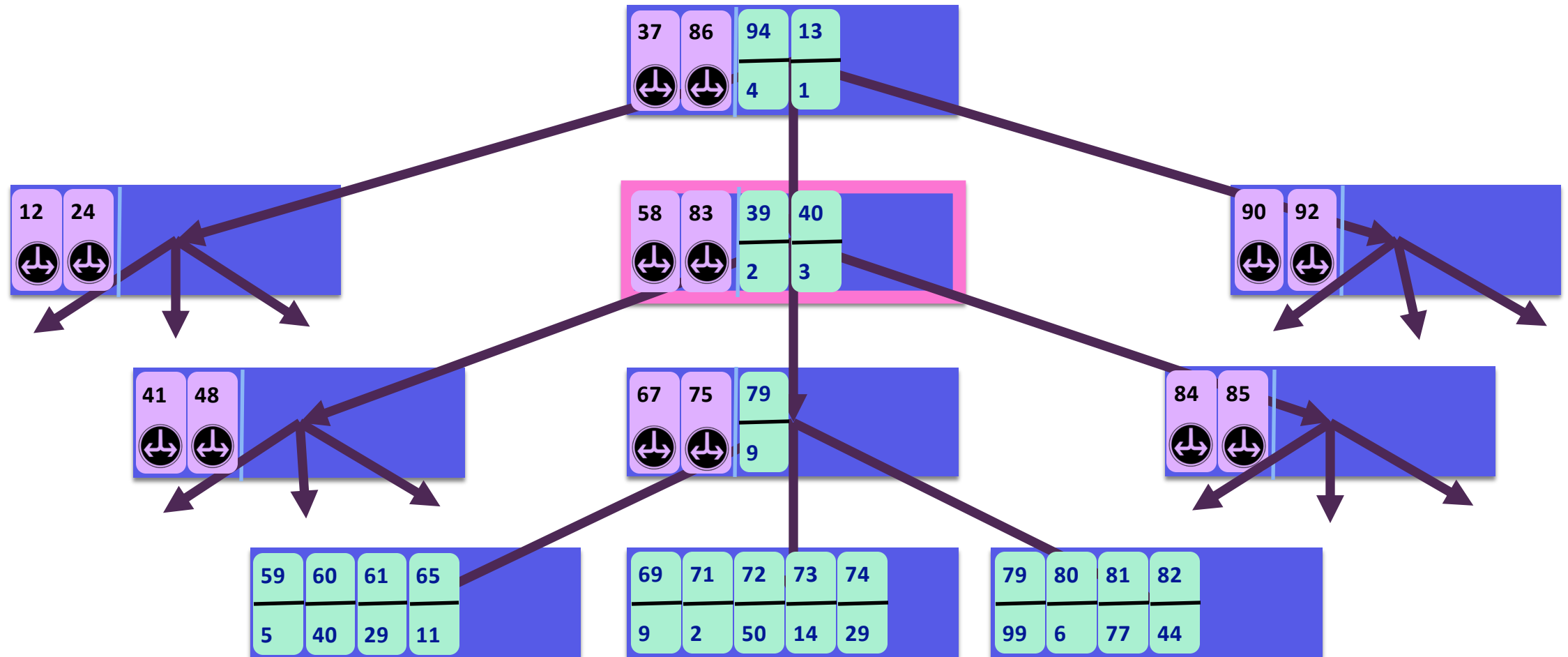
Lookups follow pivots, but check buffers along the way



B^ε-Trees

Query(71)

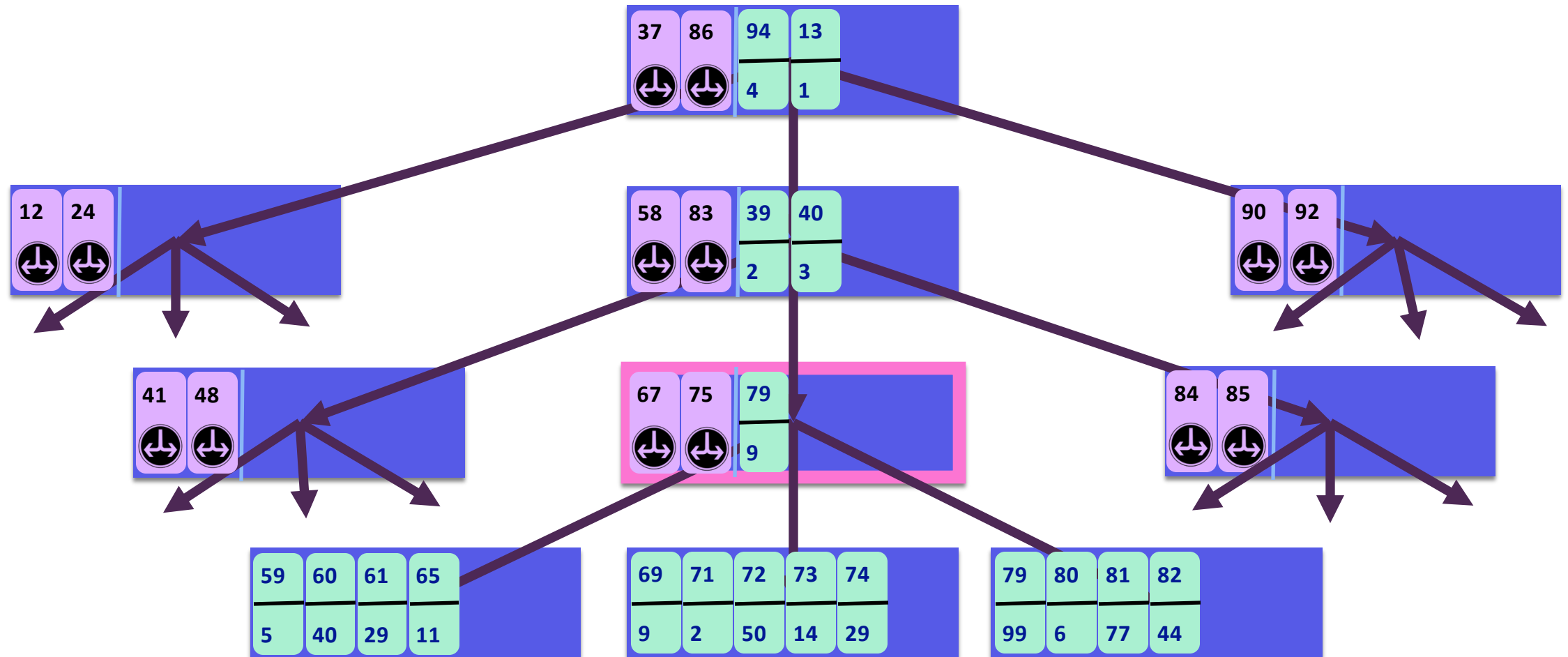
Lookups follow pivots, but check buffers along the way



B^ε-Trees

Query(71)

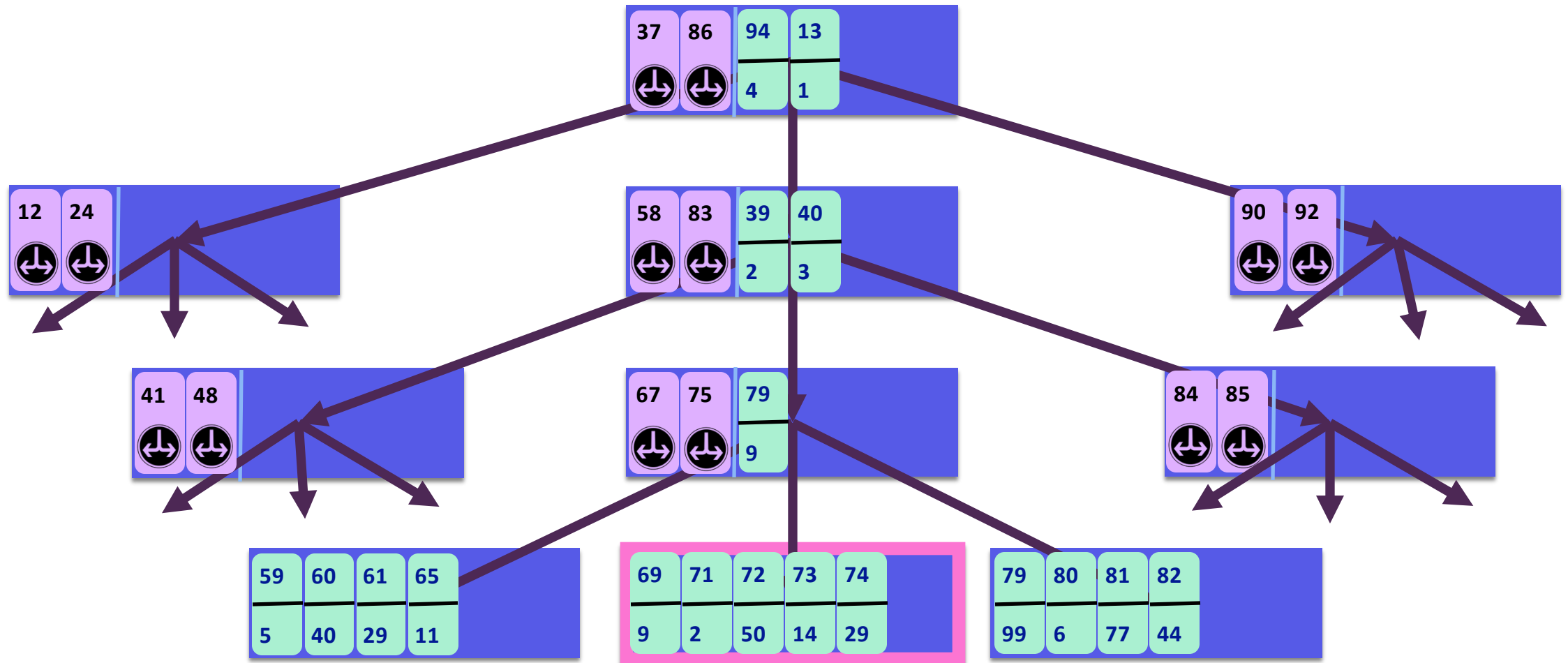
Lookups follow pivots, but check buffers along the way



B^ε-Trees

Lookups follow pivots, but check buffers along the way

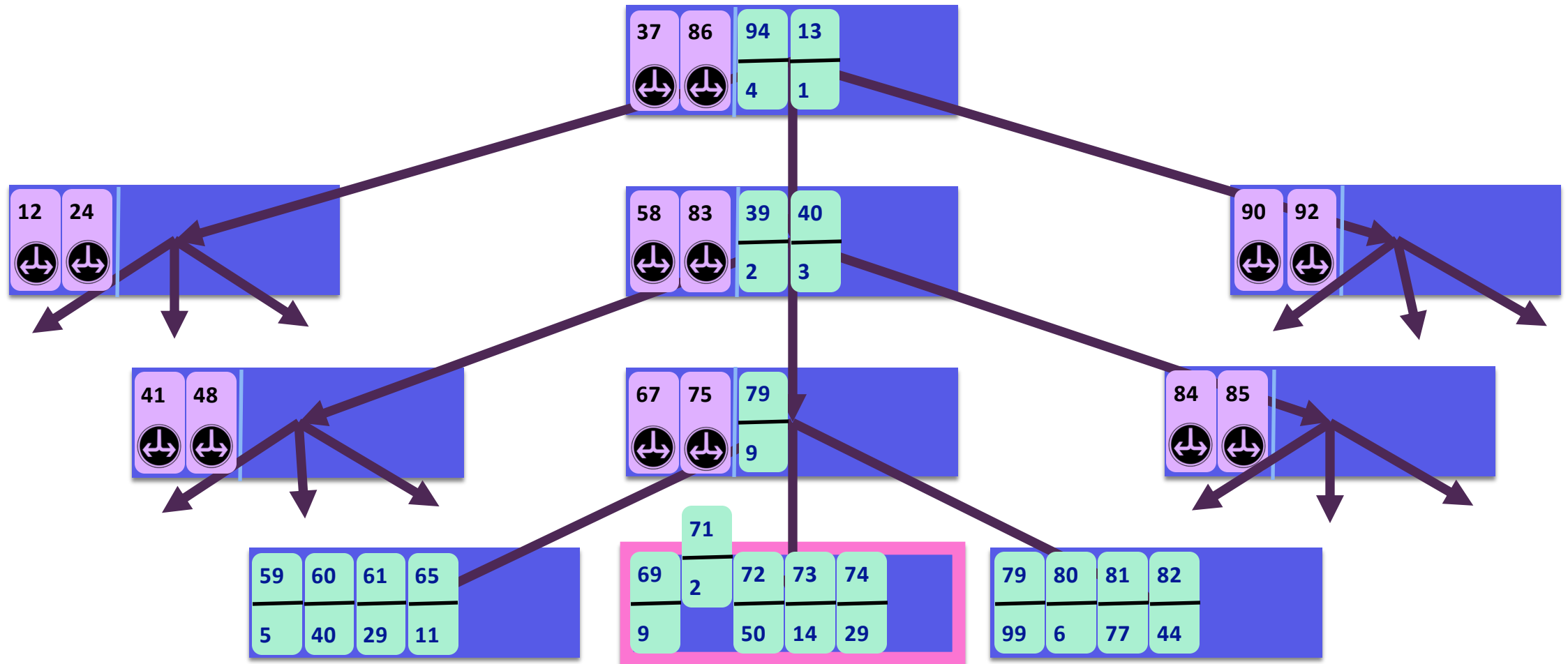
Query(71) → 2



B^ε-Trees

Lookups follow pivots, but check buffers along the way

Query(71) → 2



Insertions in B^ϵ -Trees are more expensive than they look

Insertions in B^ϵ -Trees are more expensive than they look

Recall: Insertions in B^ϵ -trees

65	72	80
11	50	6



Insertions in B^ϵ -Trees are more expensive than they look

Recall: Insertions in B^ϵ -trees

65	72	80
11	50	6

58	83	39	64	66
↕	↕	2	8	6



Read the node



Insertions in B^E -Trees are more expensive than they look

Recall: Insertions in B^E -trees

65	72	80
11	50	6

Merge the data



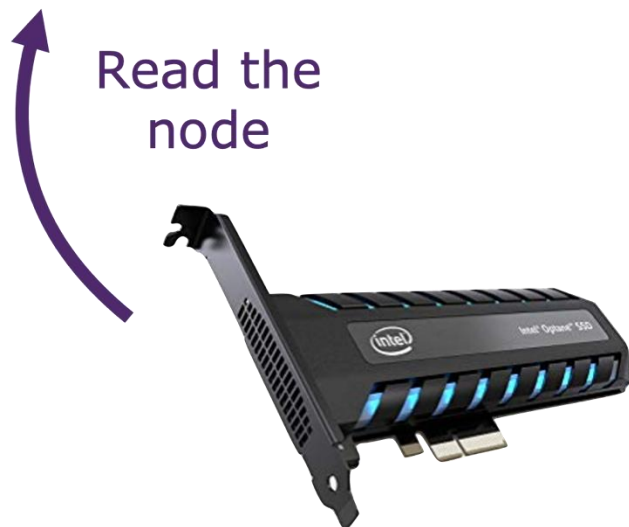
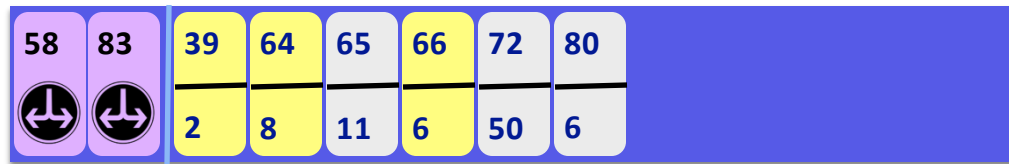
58	83	39	64	66
⬇	⬇	2	8	6

Read the node



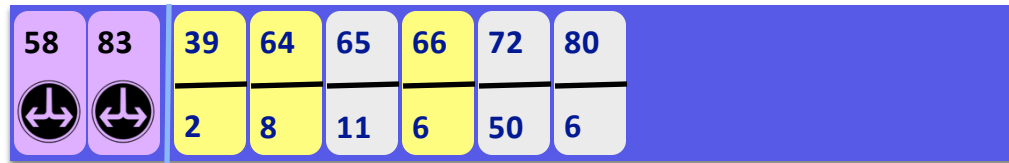
Insertions in B^ϵ -Trees are more expensive than they look

Recall: Insertions in B^ϵ -trees



Insertions in B^ϵ -Trees are more expensive than they look

Recall: Insertions in B^ϵ -trees



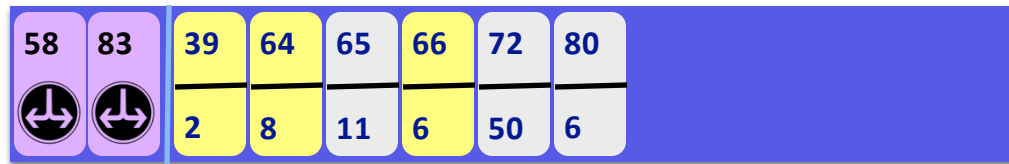
Insertions in B^ϵ -Trees are more expensive than they look

Recall: Insertions in B^ϵ -trees



CPU Work = $O(\text{old} + \text{new})$

Volume of IO = $O(\text{old} + \text{new})$



Insertions in B^ϵ -Trees are more expensive than they look

Recall: Insertions in B^ϵ -trees

98	44
1	3

Merge the data



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

CPU Work = $O(\text{old} + \text{new})$

Volume of IO = $O(\text{old} + \text{new})$

Older data gets written over and over again

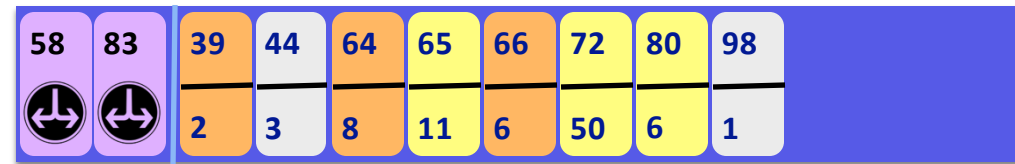
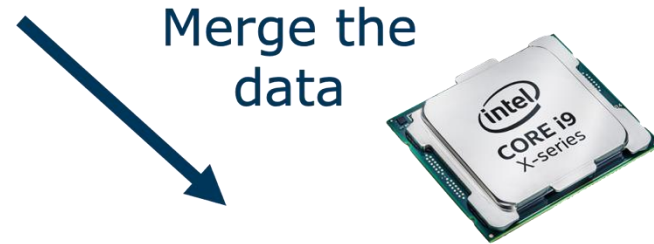
Read the node

Write the node



Insertions in B^ϵ -Trees are more expensive than they look

Recall: Insertions in B^ϵ -trees



CPU Work = $O(\text{old} + \text{new})$

Volume of IO = $O(\text{old} + \text{new})$

Older data gets written over and over again

Insertions in B^ϵ -Trees are more expensive than they look

Recall: Insertions in B^ϵ -trees

28	91
24	43

Merge the data



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

CPU Work = $O(\text{old} + \text{new})$

Volume of IO = $O(\text{old} + \text{new})$

Older data gets written over and over again

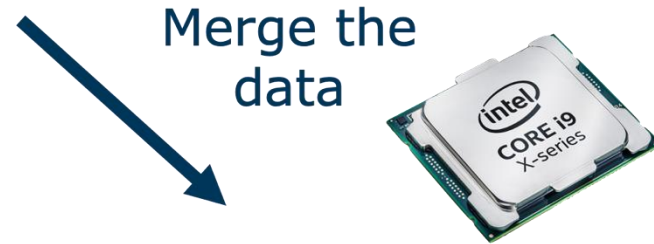
Read the node

Write the node



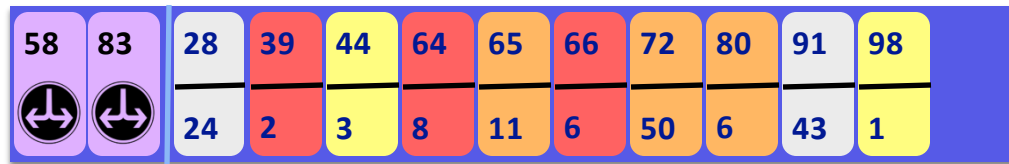
Insertions in B^ϵ -Trees are more expensive than they look

Recall: Insertions in B^ϵ -trees



CPU Work = $O(\text{old} + \text{new})$

Volume of IO = $O(\text{old} + \text{new})$

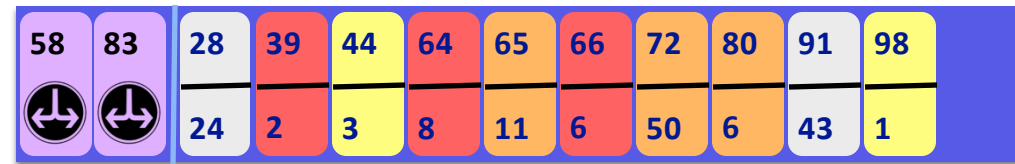


Older data gets written over and over again



Insertions in B^ϵ -Trees are more expensive than they look

Recall: Insertions in B^ϵ -trees



CPU Work = $O(\text{old} + \text{new})$

Volume of IO = $O(\text{old} + \text{new})$

Older data gets written over and over again

Up to B^ϵ times per node!

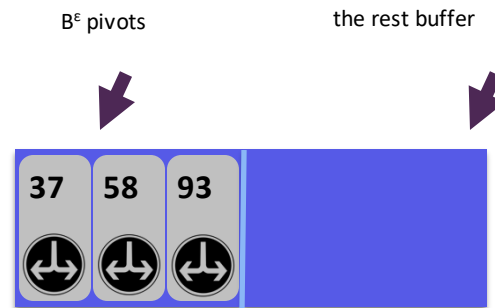
Size-Tiered B^ϵ -Trees

SplinterDB: Closing the Bandwidth Gap for NVMe Key-Value Stores
Conway, Gupta, Chidambaram, Farach-Colton, Spillane, Tai, Johnson,
ATC 2020

Size-Tiered B^ϵ -Trees

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

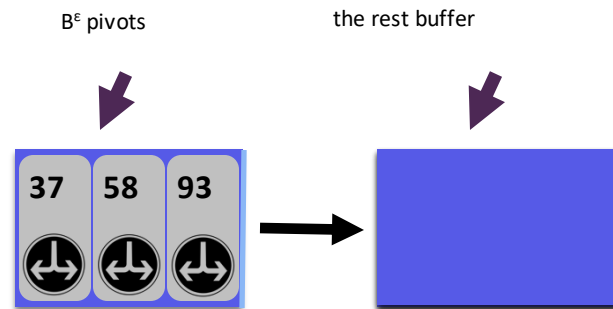
Recall:
a B^ϵ -tree node has pivots and a buffer



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

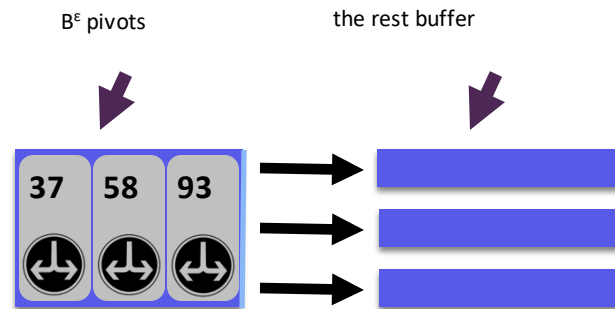


In an STB^ϵ -tree, the buffer is stored separately

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

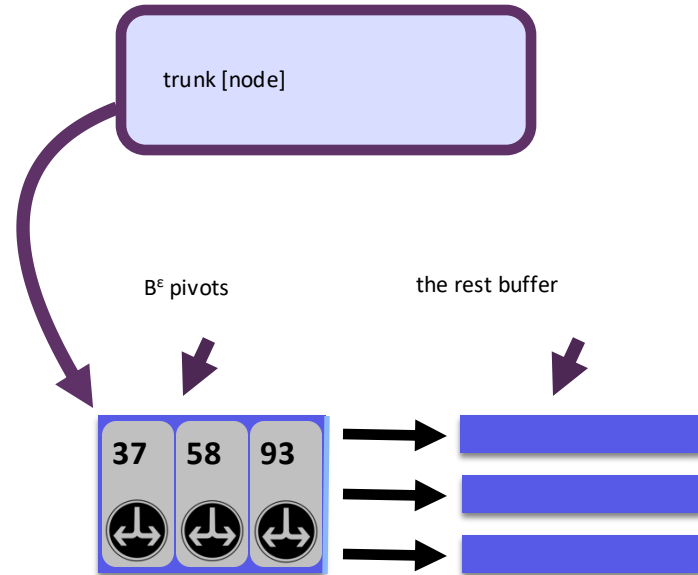


and in several discontinuous pieces

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



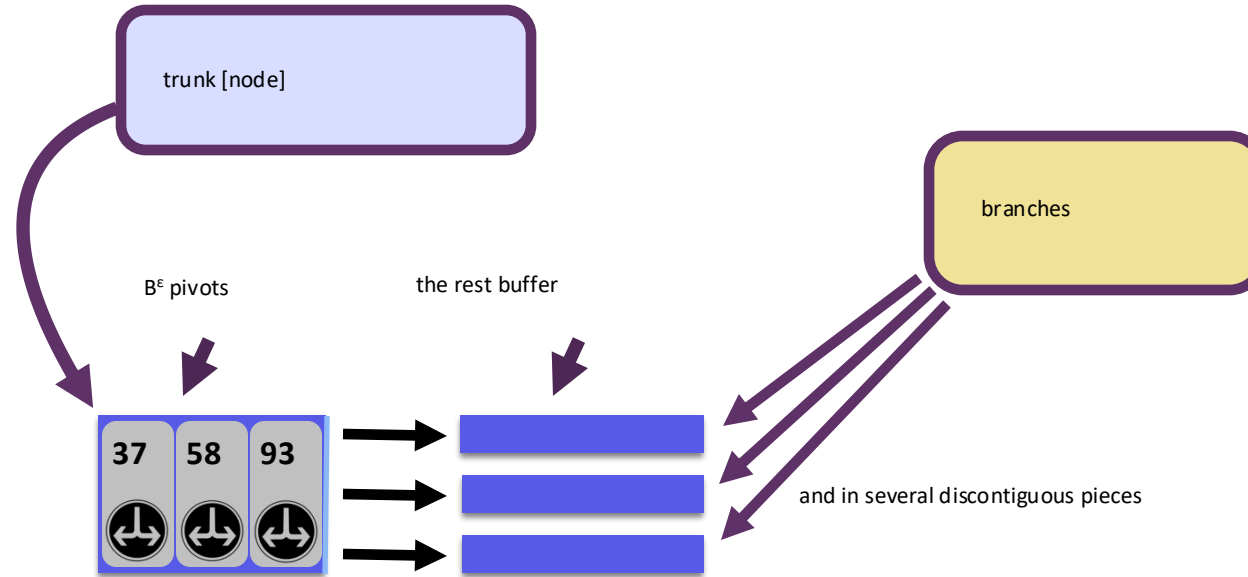
Recall:
a B^ϵ -tree node has pivots and a buffer

and in several discontinuous 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

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

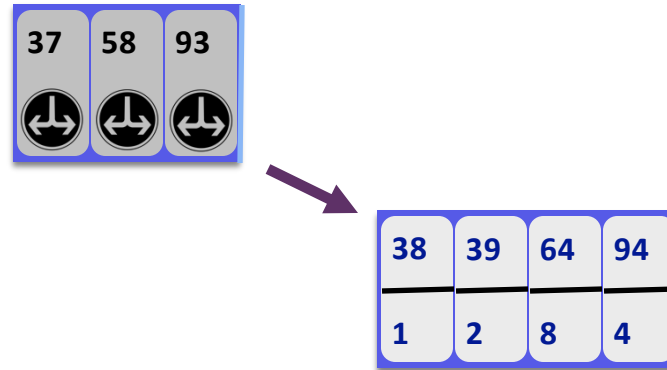
38	39	64	94
1	2	8	4

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

37	58	93
		

Size-Tiered B^ε-Trees

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

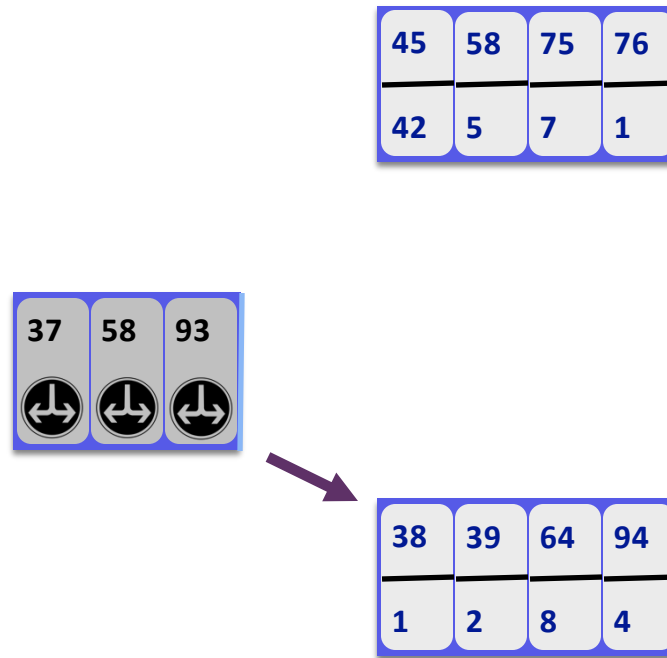


When new data is flushed into the trunk node...

...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

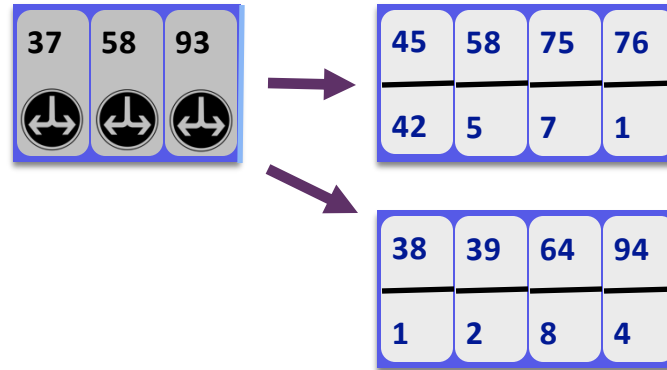


When new data is flushed into the trunk node...

...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...

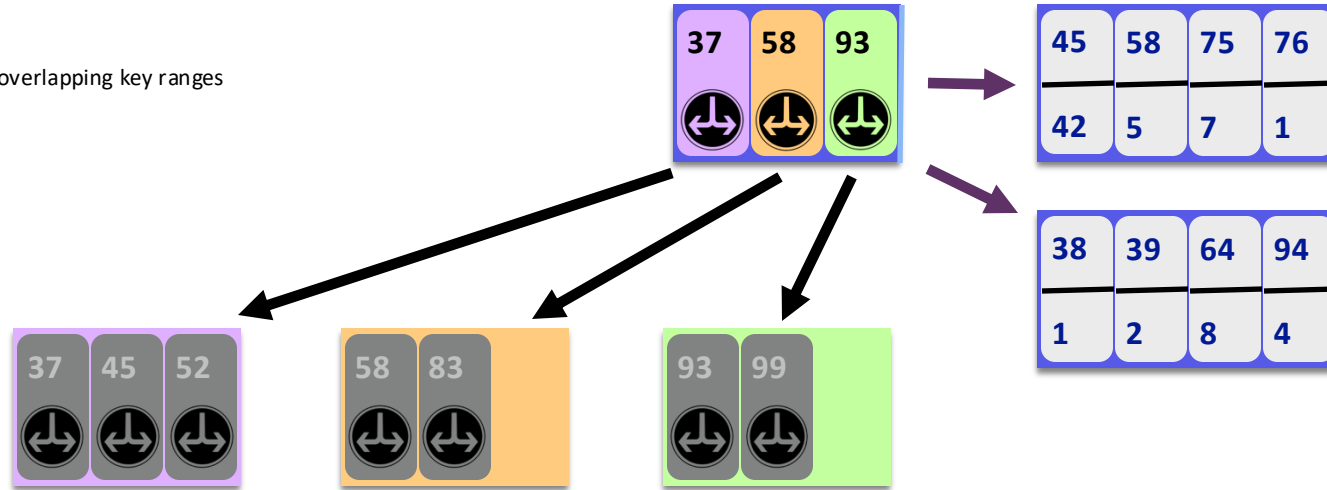
...it is added as a new branch

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...

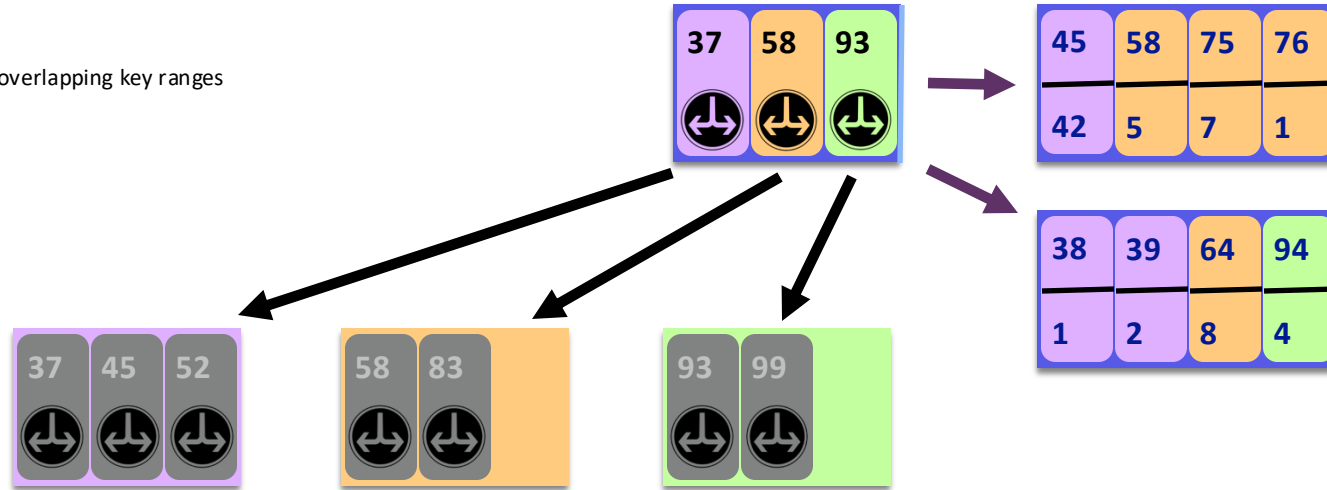
...it is added as a new branch

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...

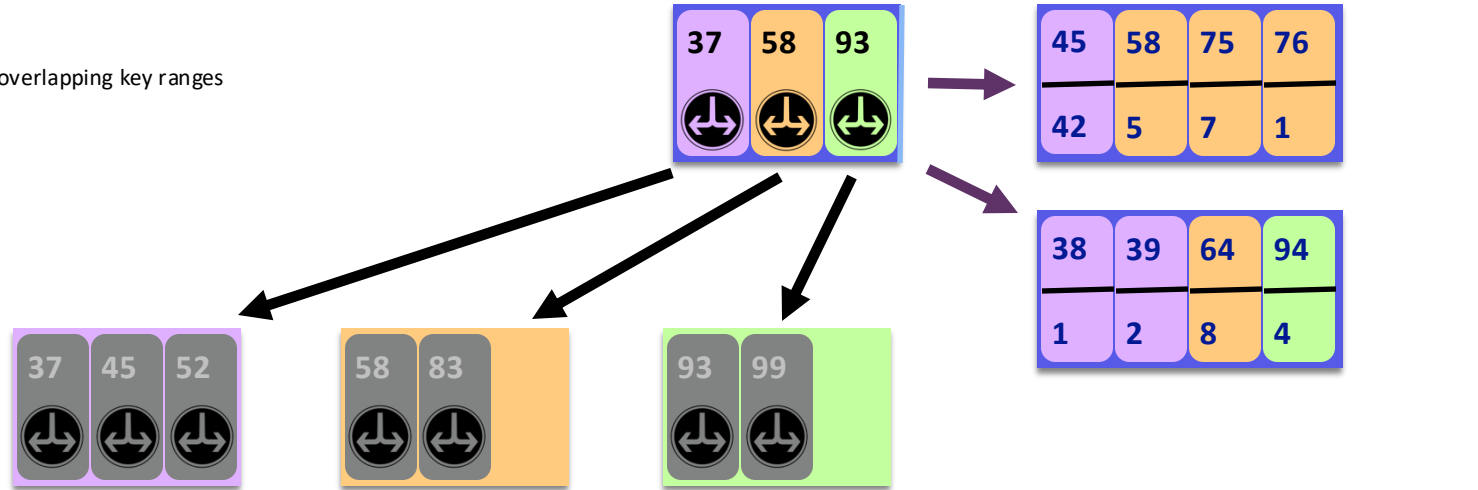
...it is added as a new branch

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



41	42	43	79	85	91
2	5	11	1	2	9

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

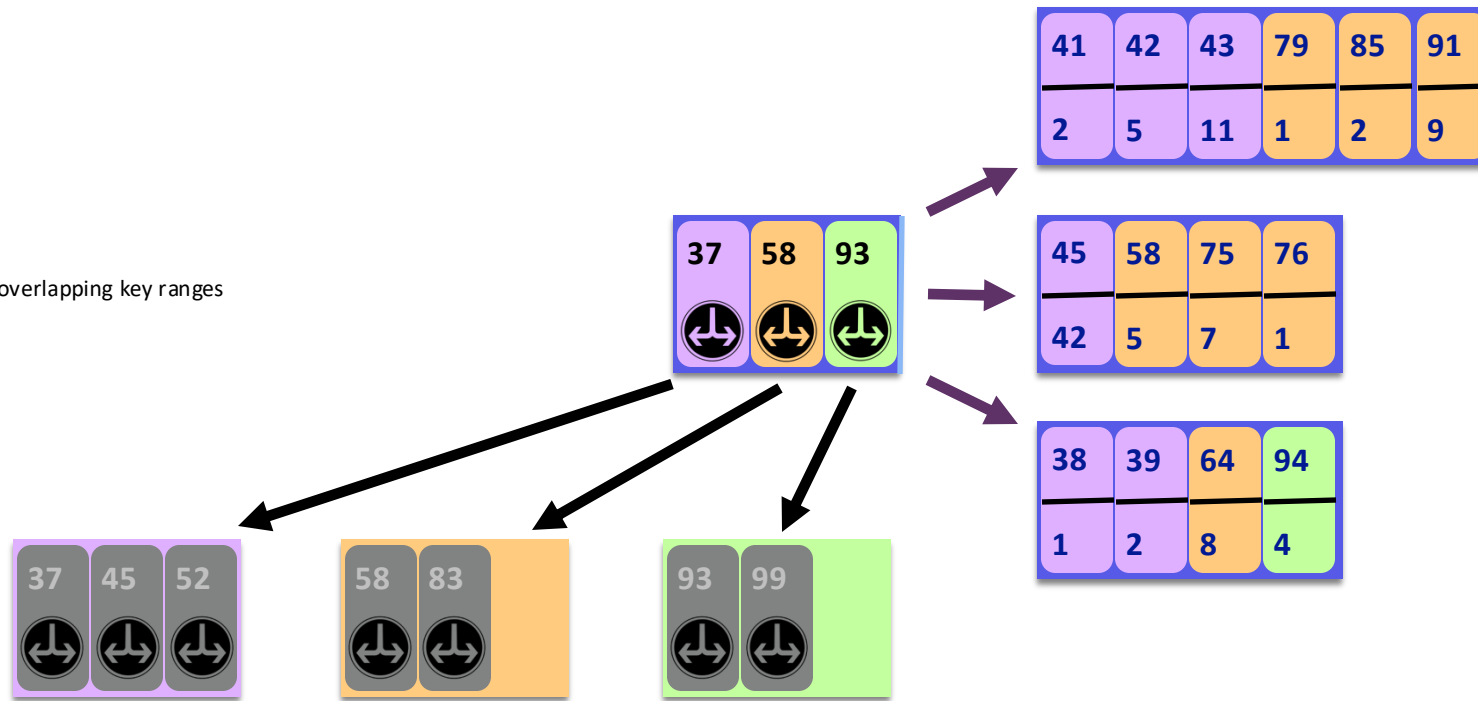
...it is added as a new branch

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

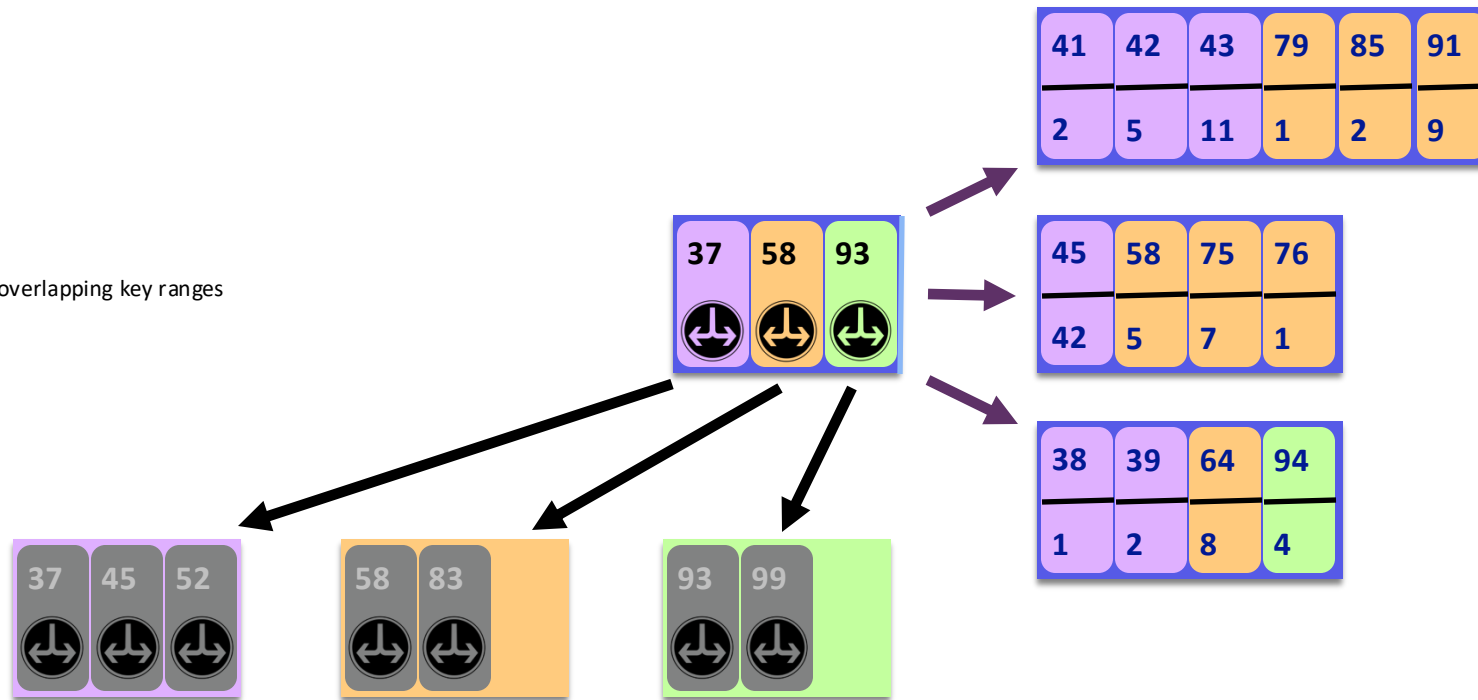
Size-Tiered B^ε-Trees

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

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

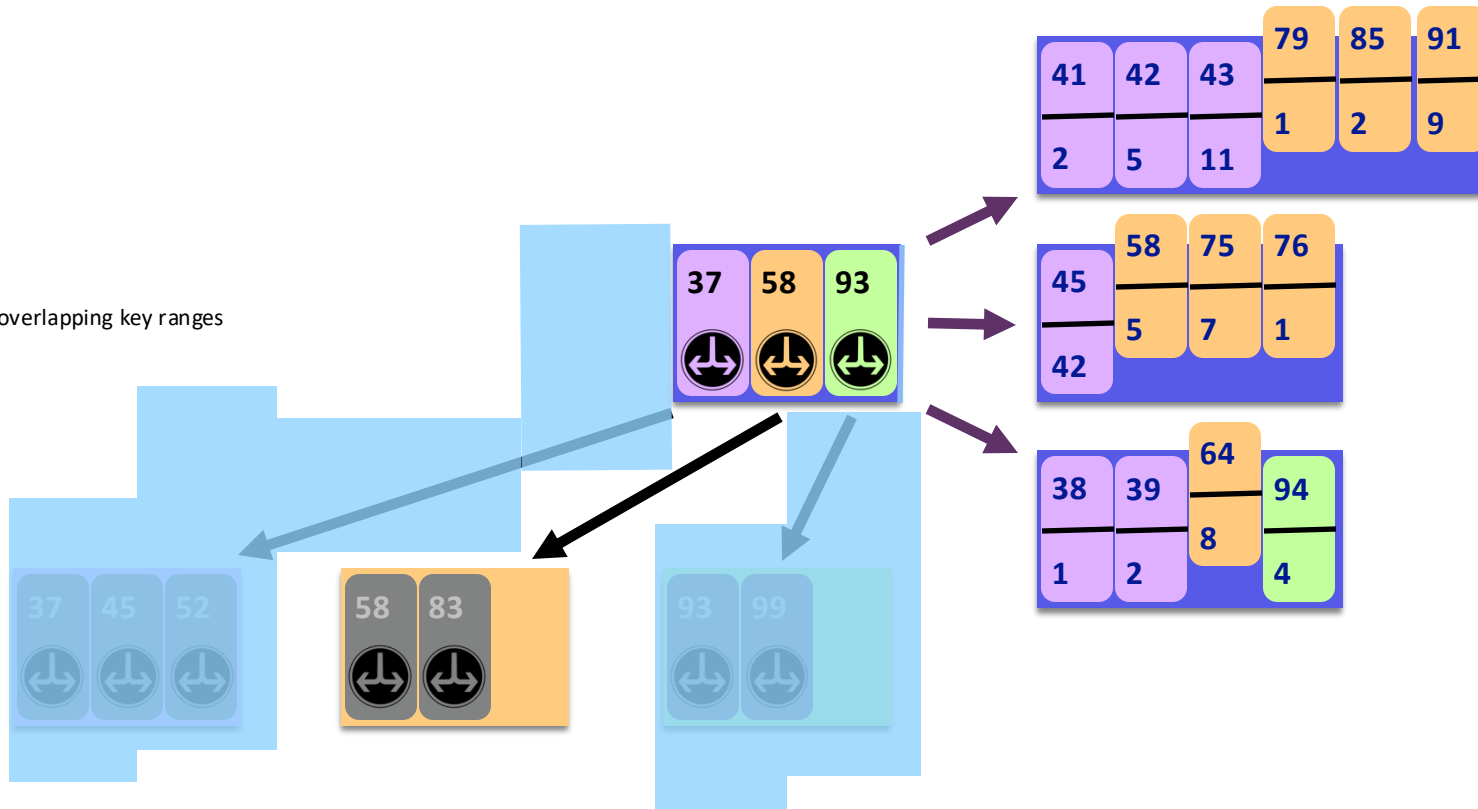
Size-Tiered B^ε-Trees

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

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

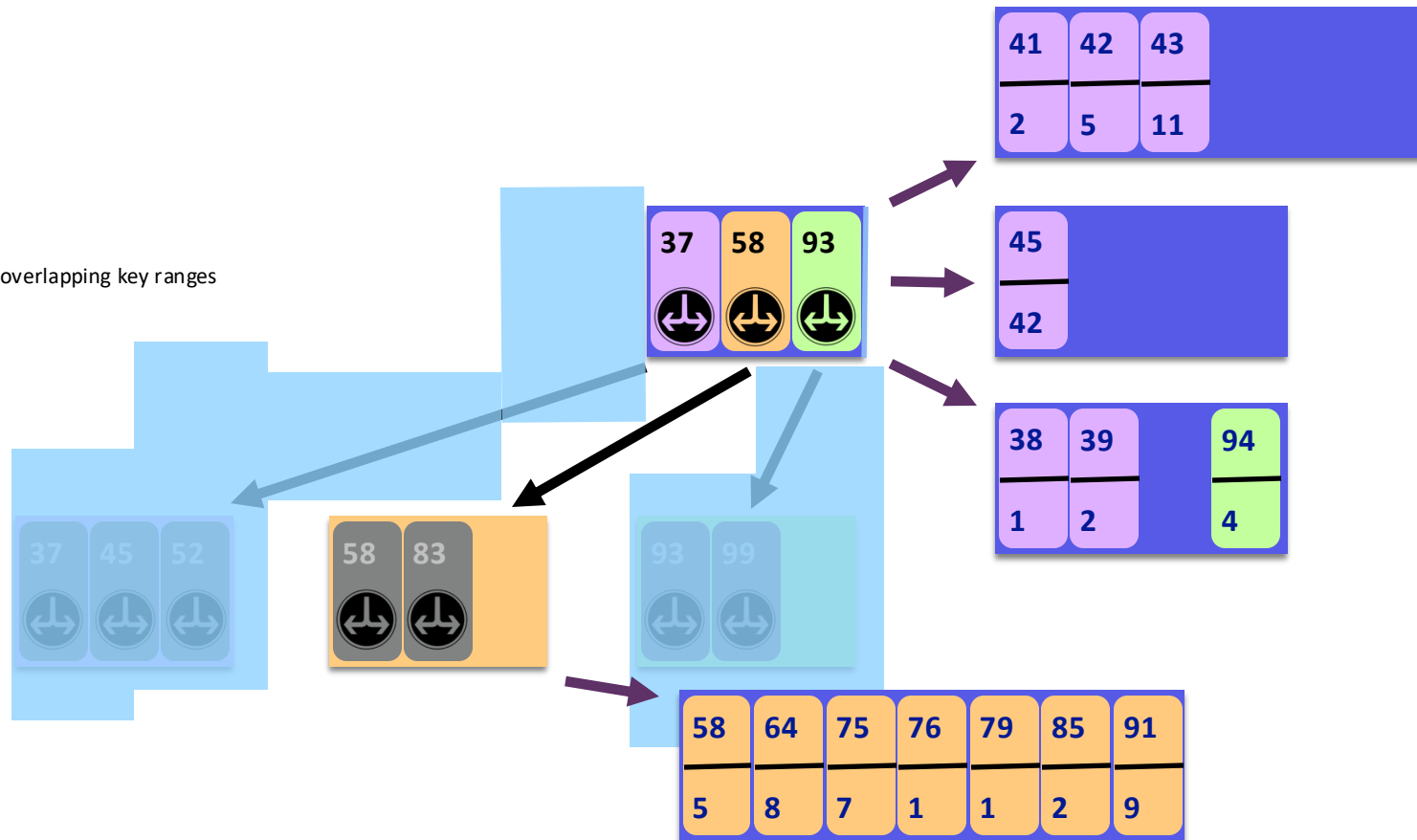
Size-Tiered B^ε-Trees

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

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

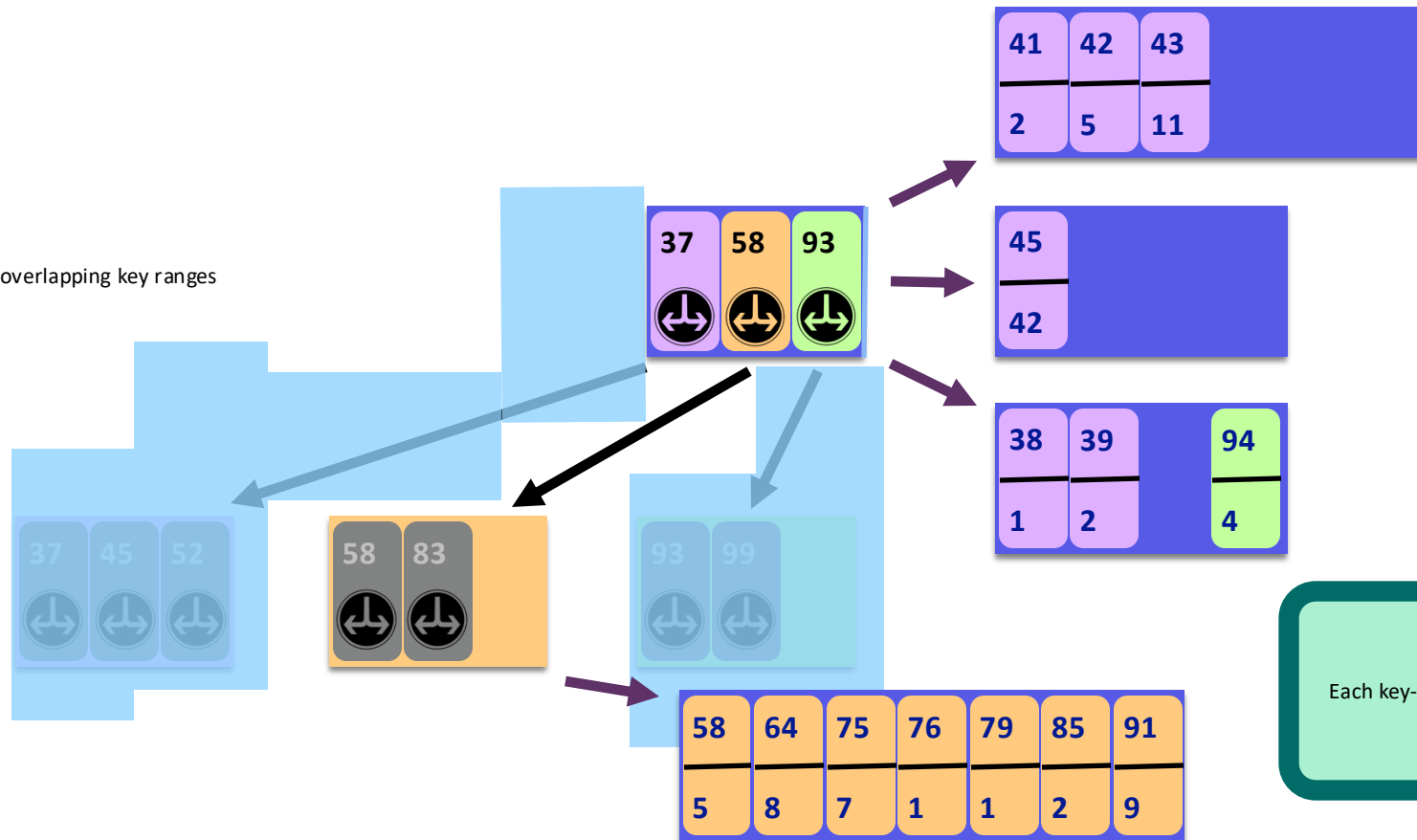
Size-Tiered B^ε-Trees

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

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

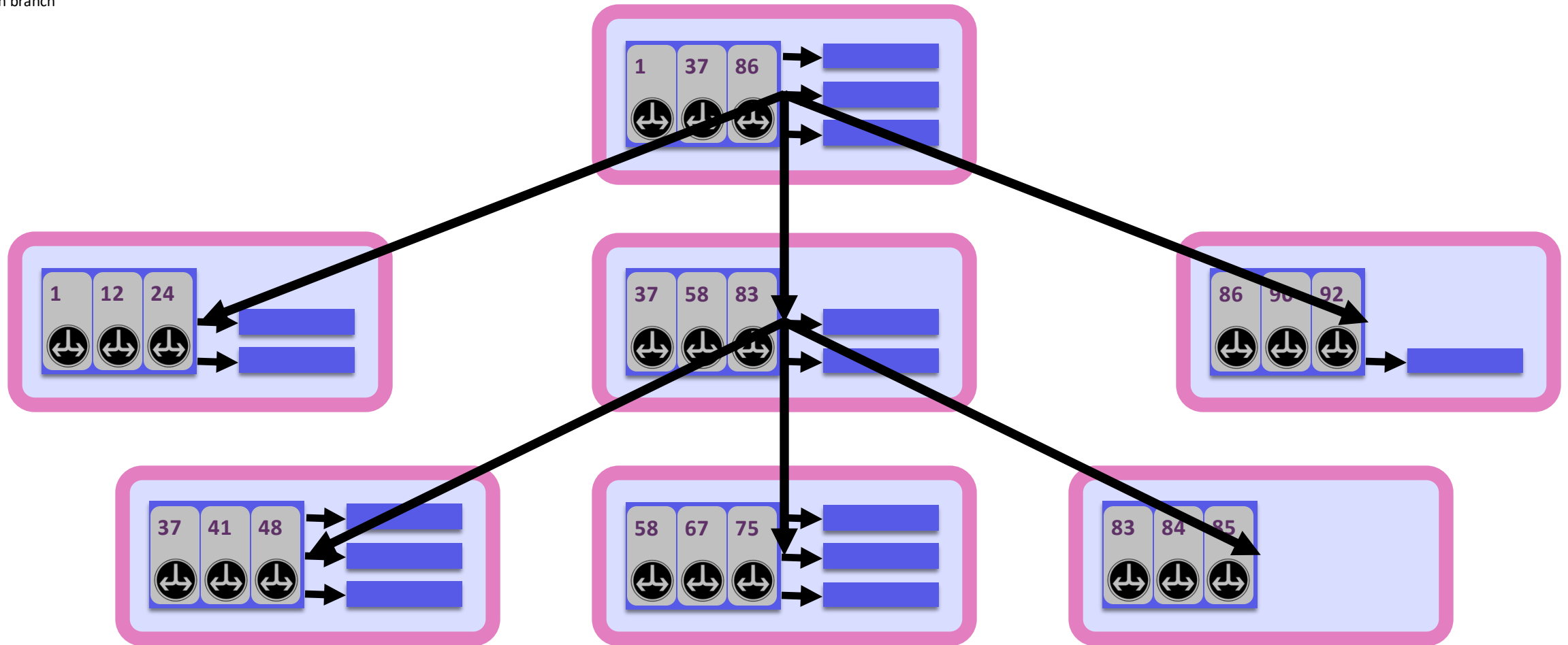
Each key-value pair is read/written once per trunk node

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

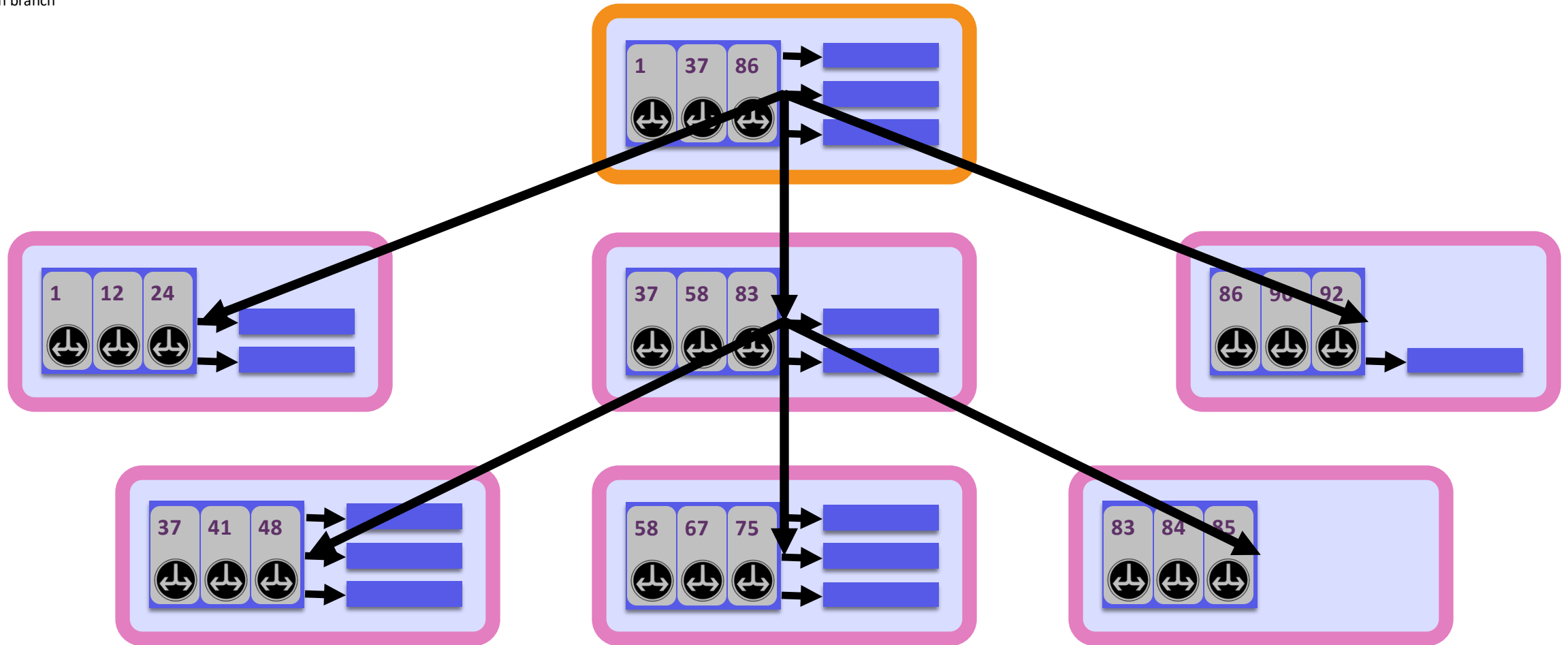
Query(71)



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

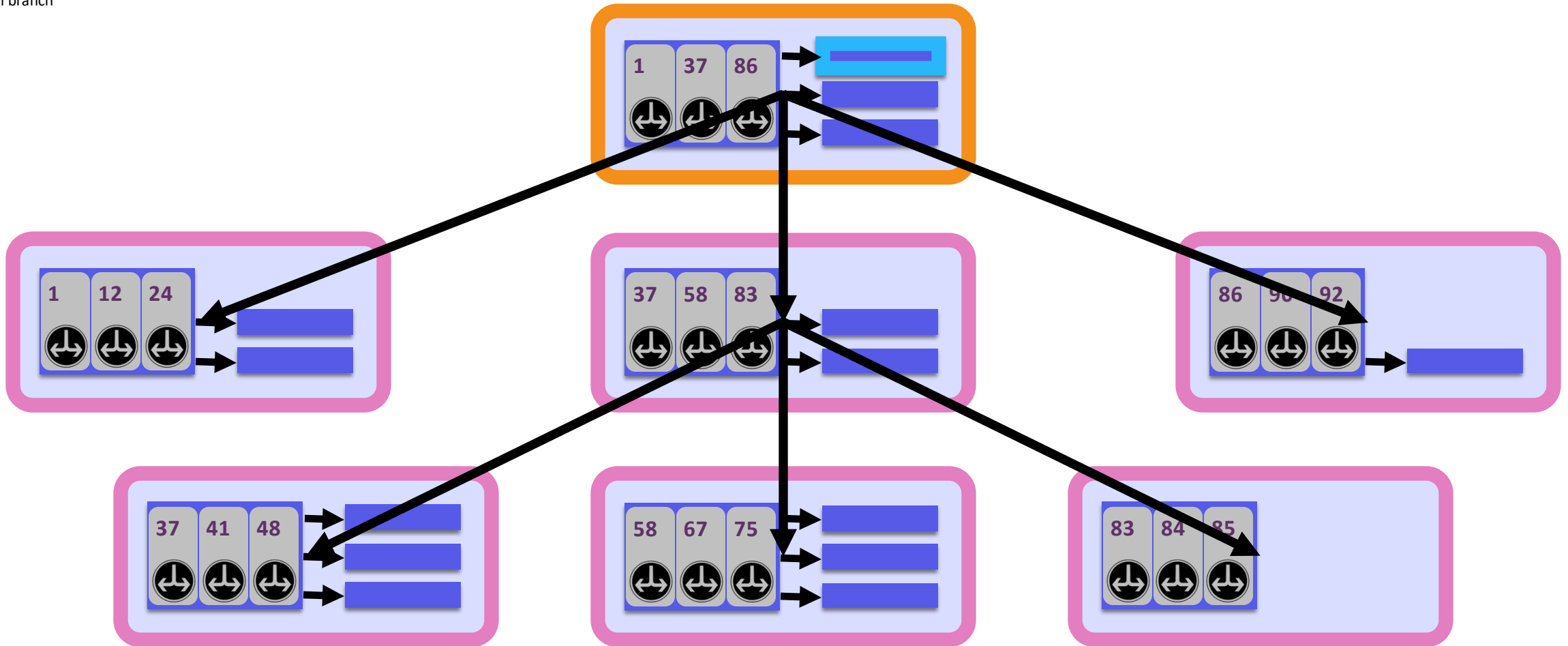
Query(71)



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

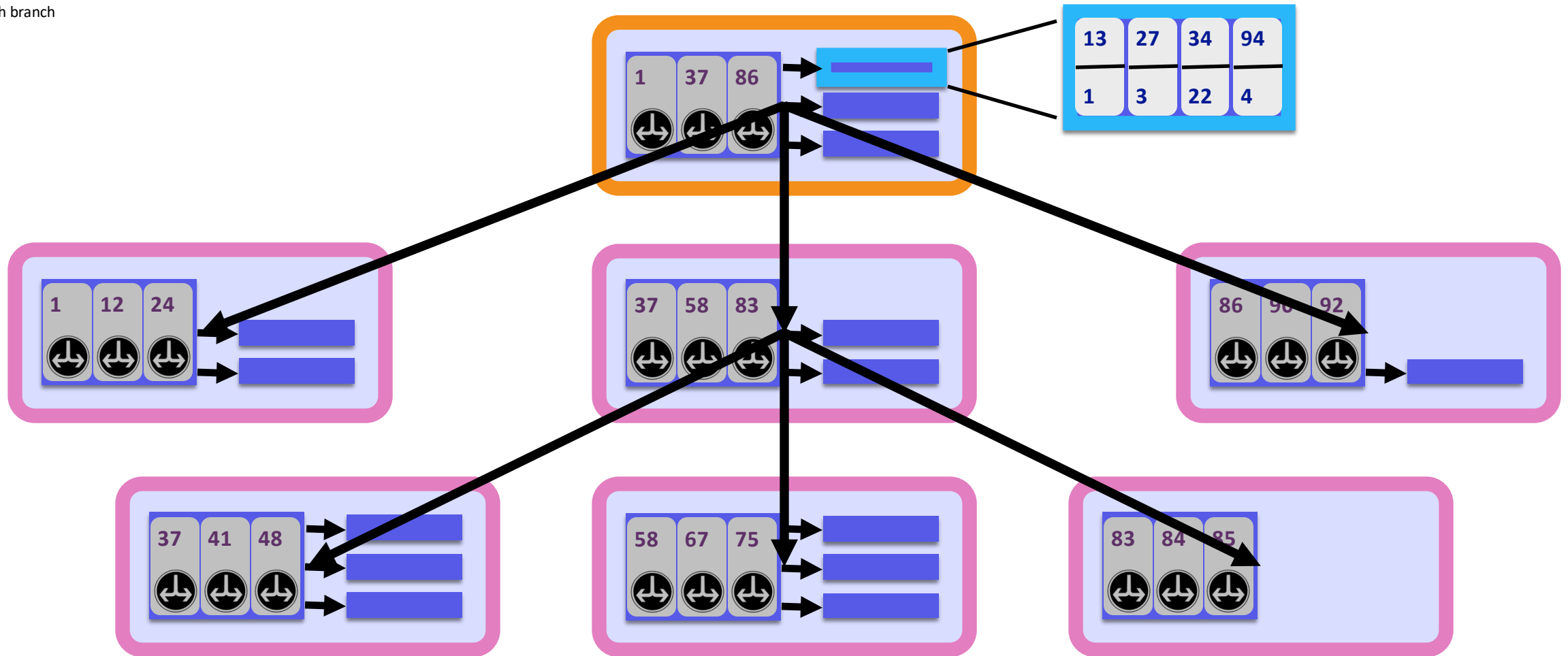
Query(71)



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

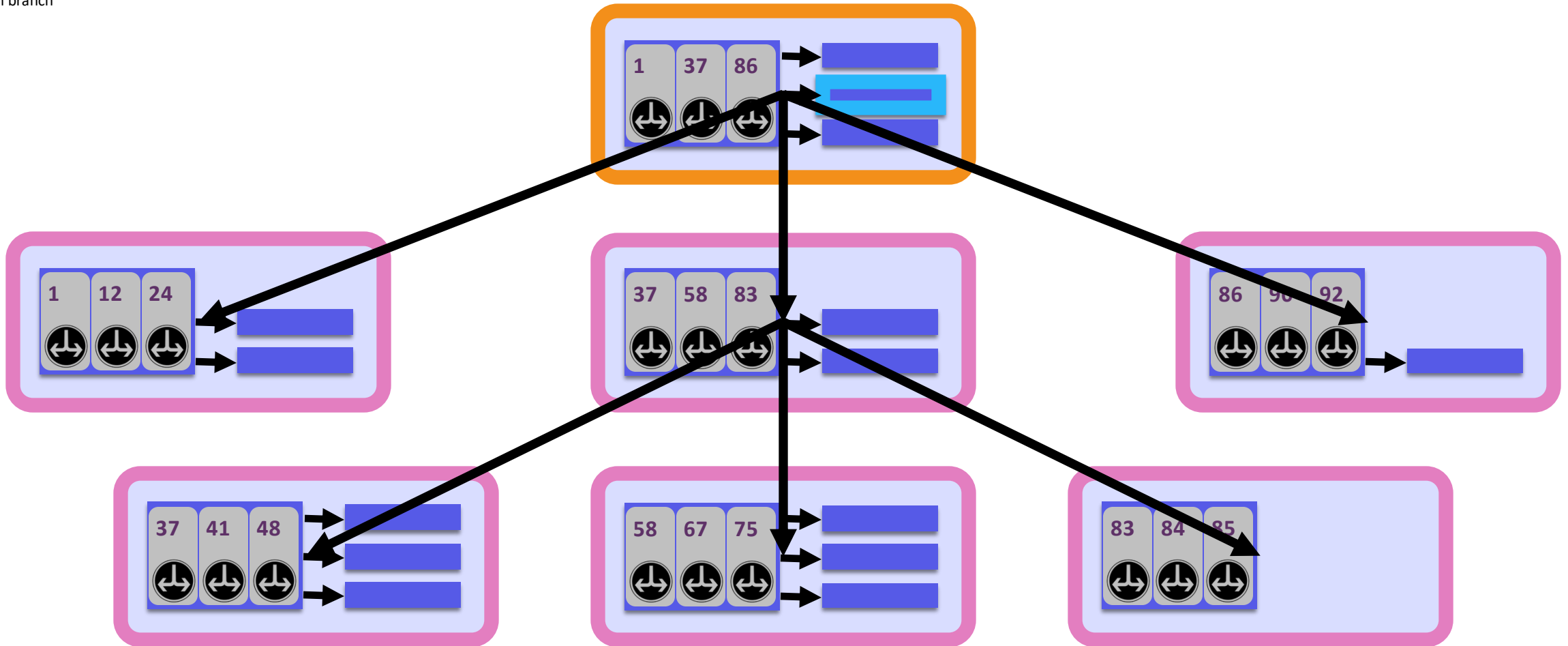
Query(71)



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

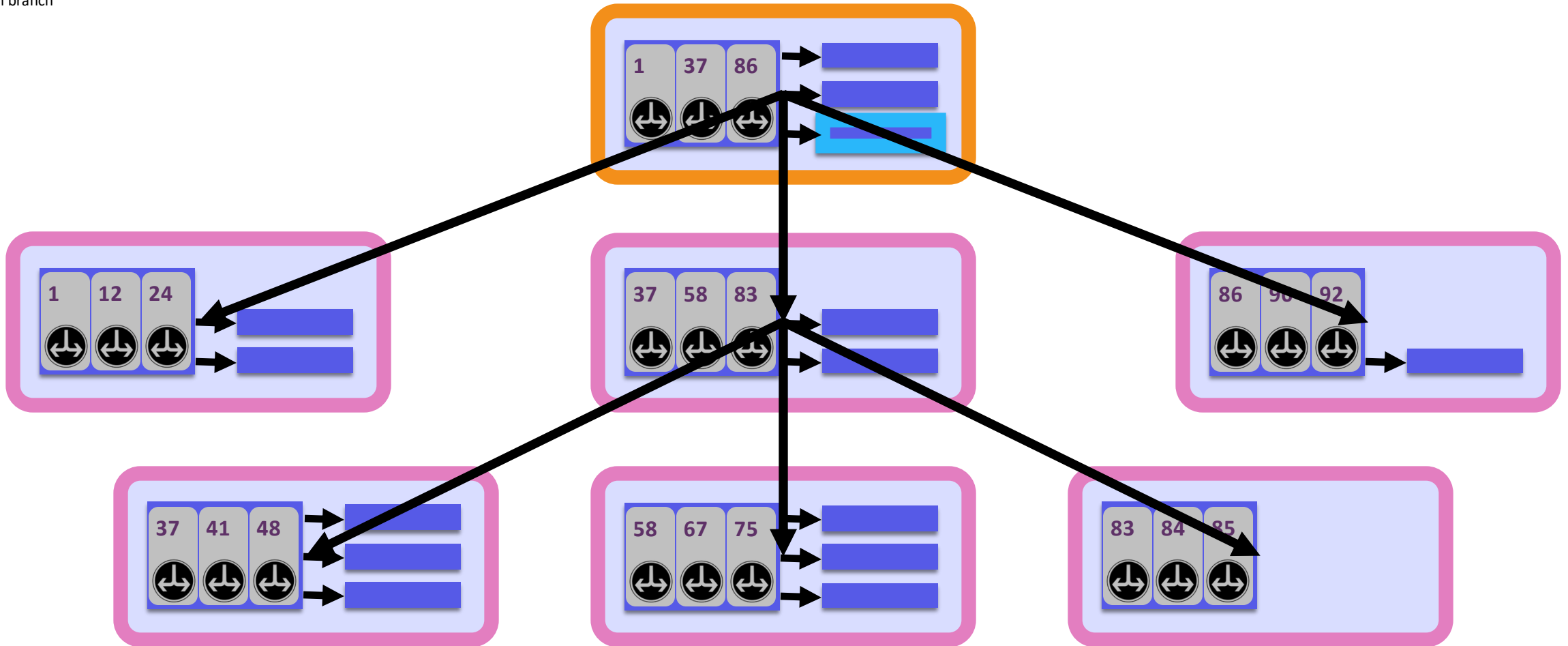
Query(71)



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

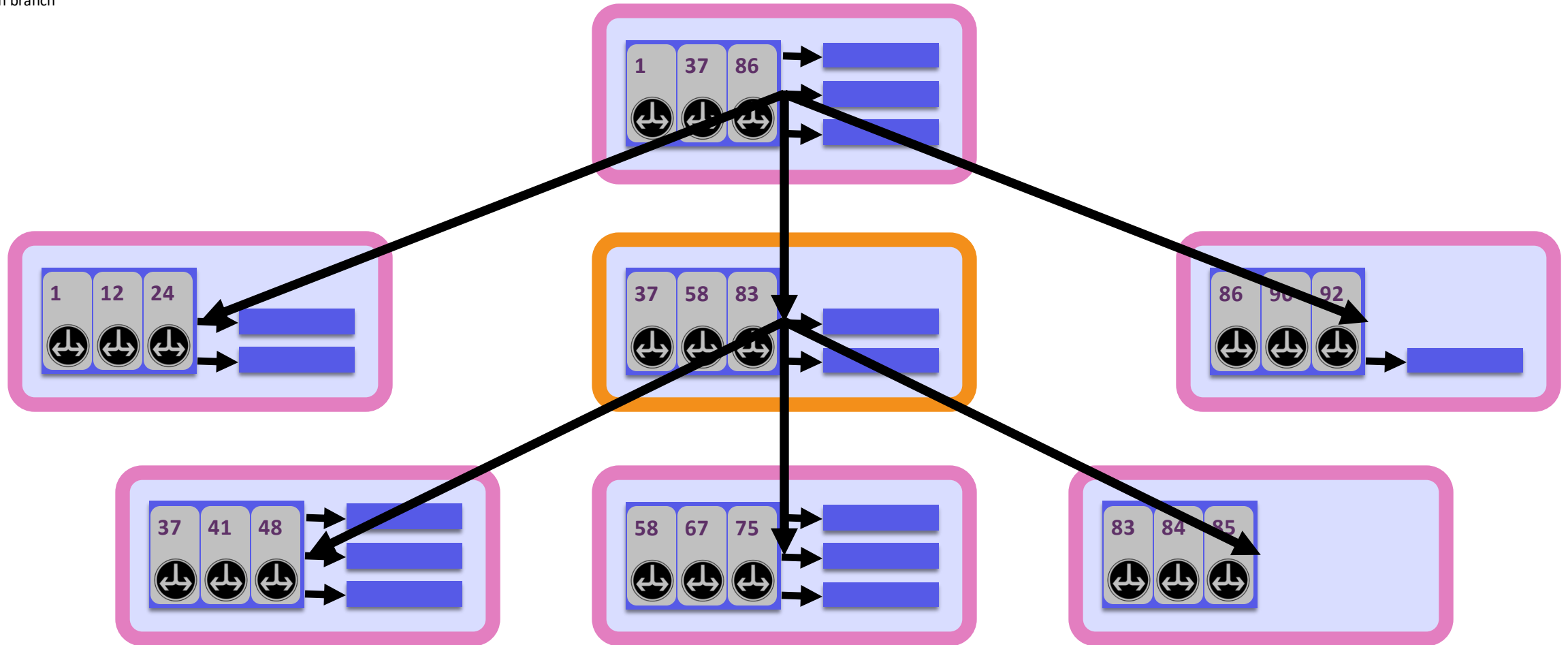
Query(71)



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

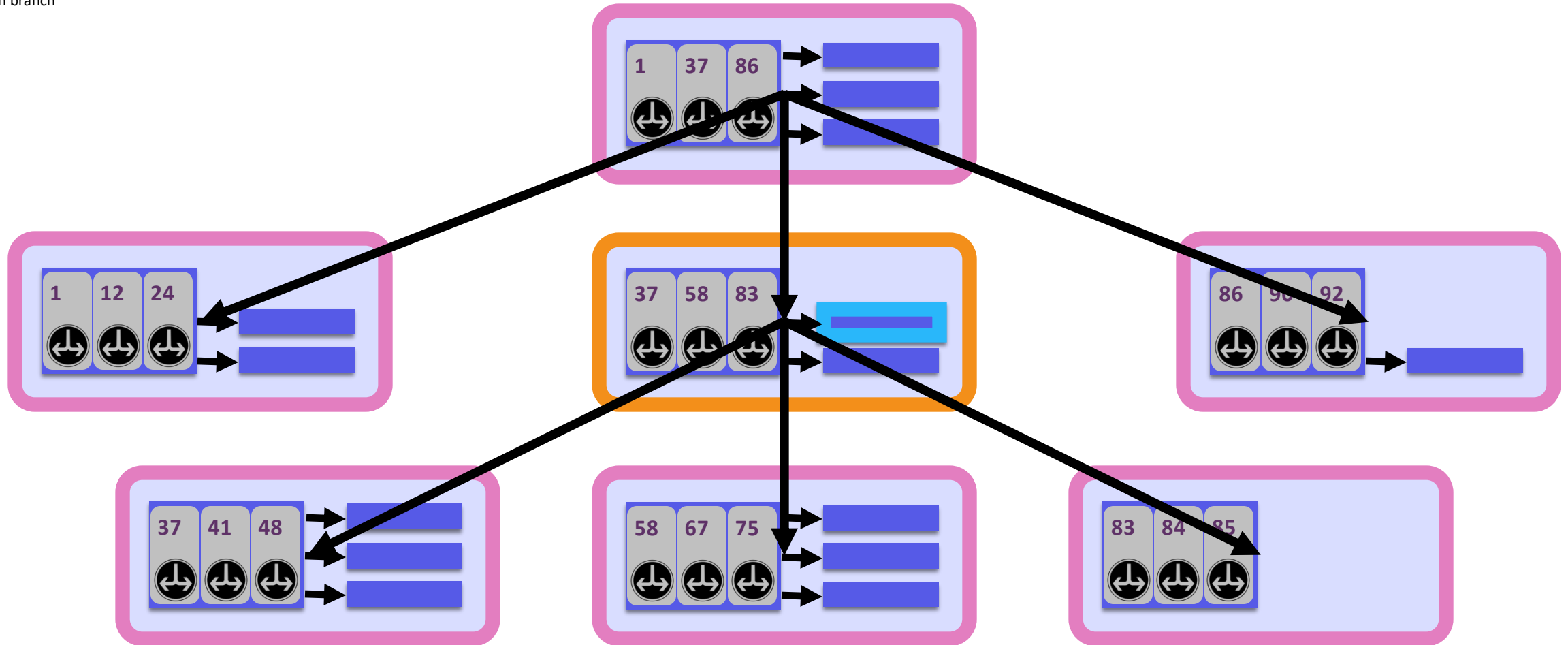
Query(71)



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

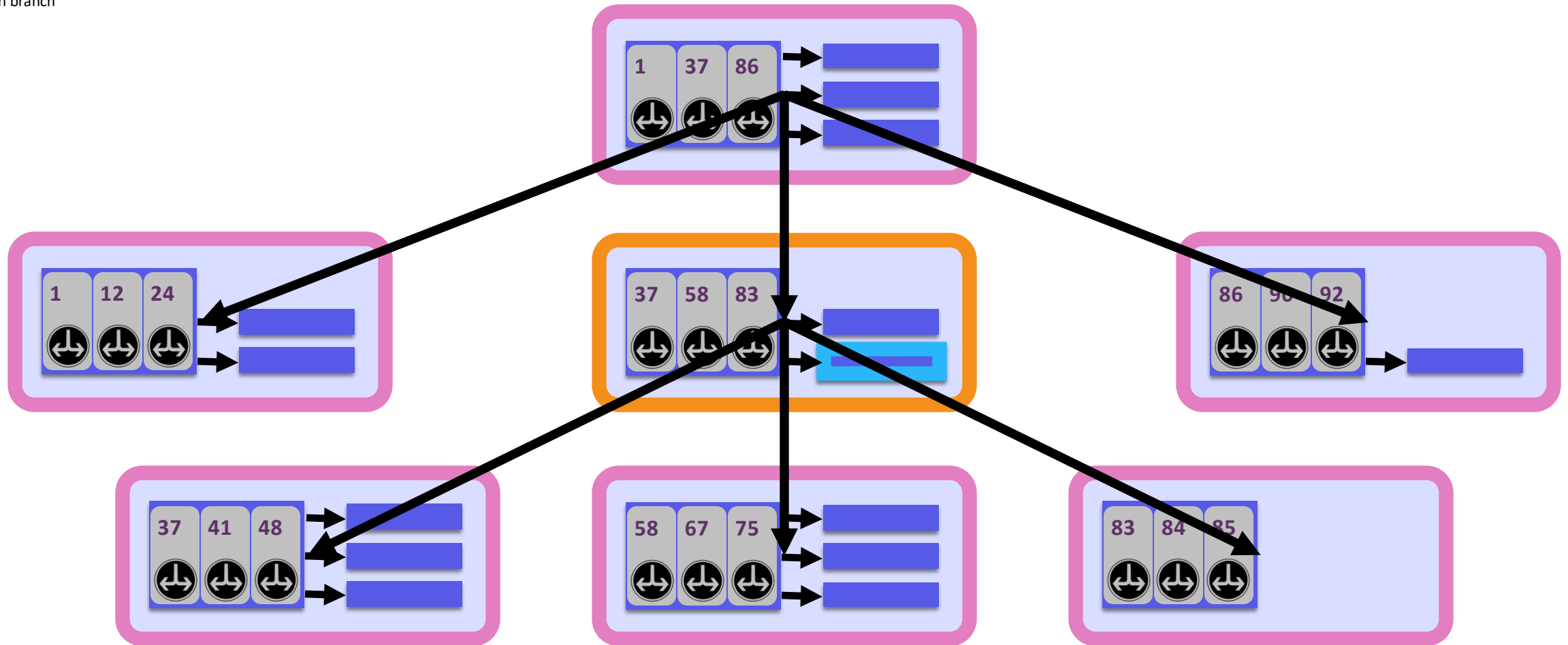
Query(71)



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

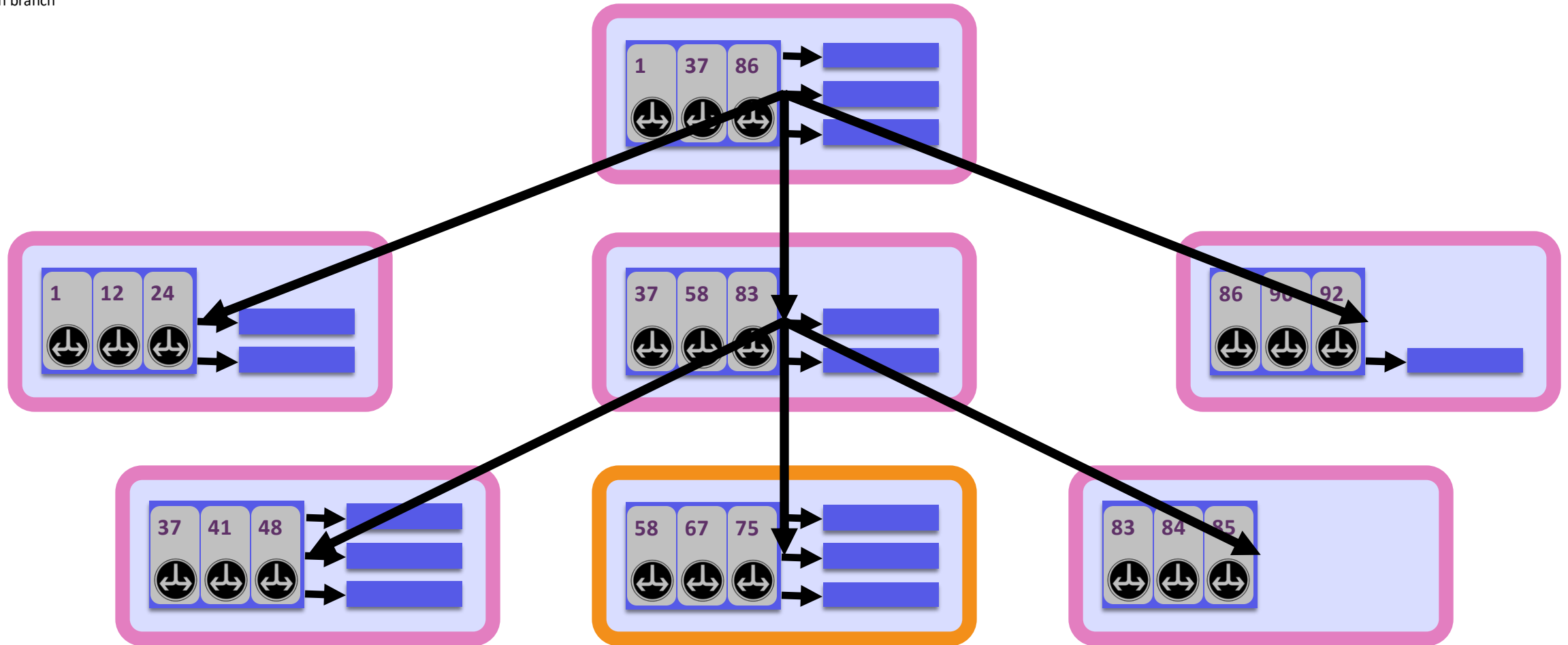
Query(71)



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

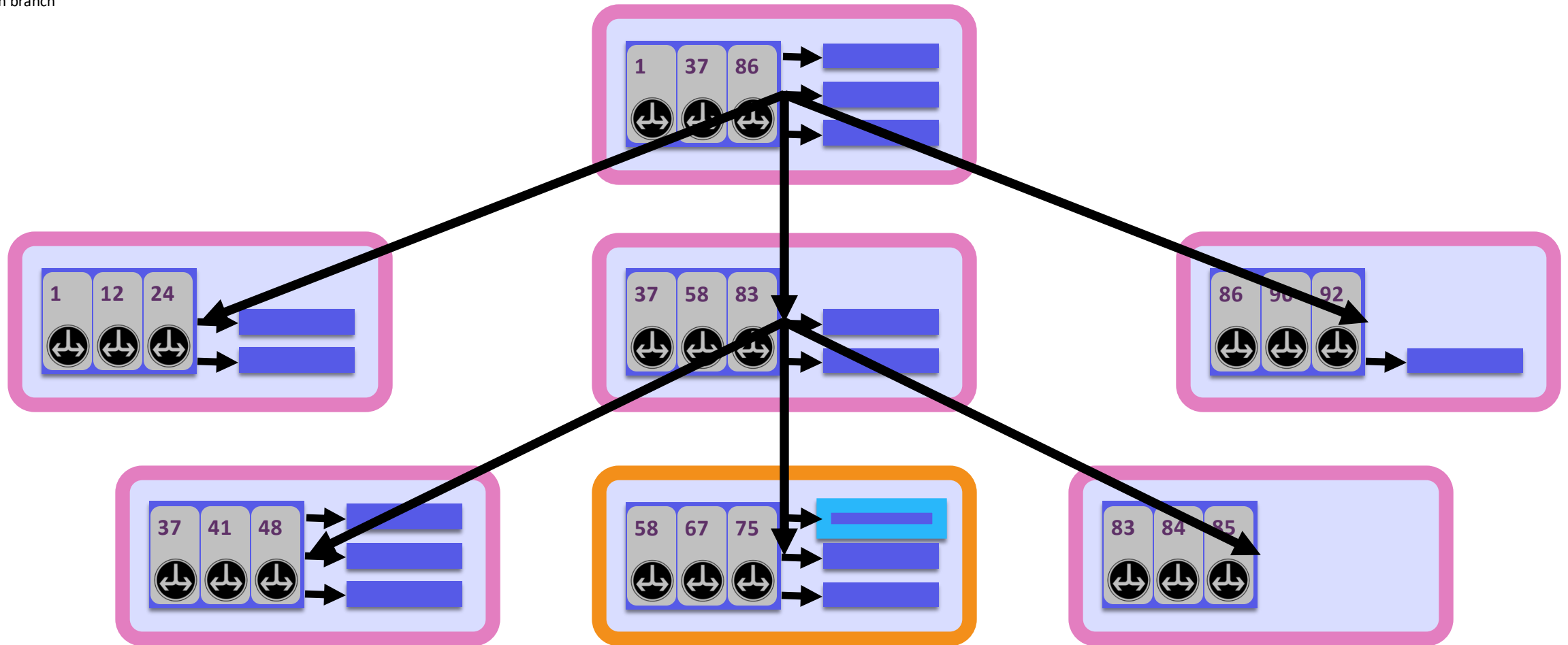
Query(71)



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

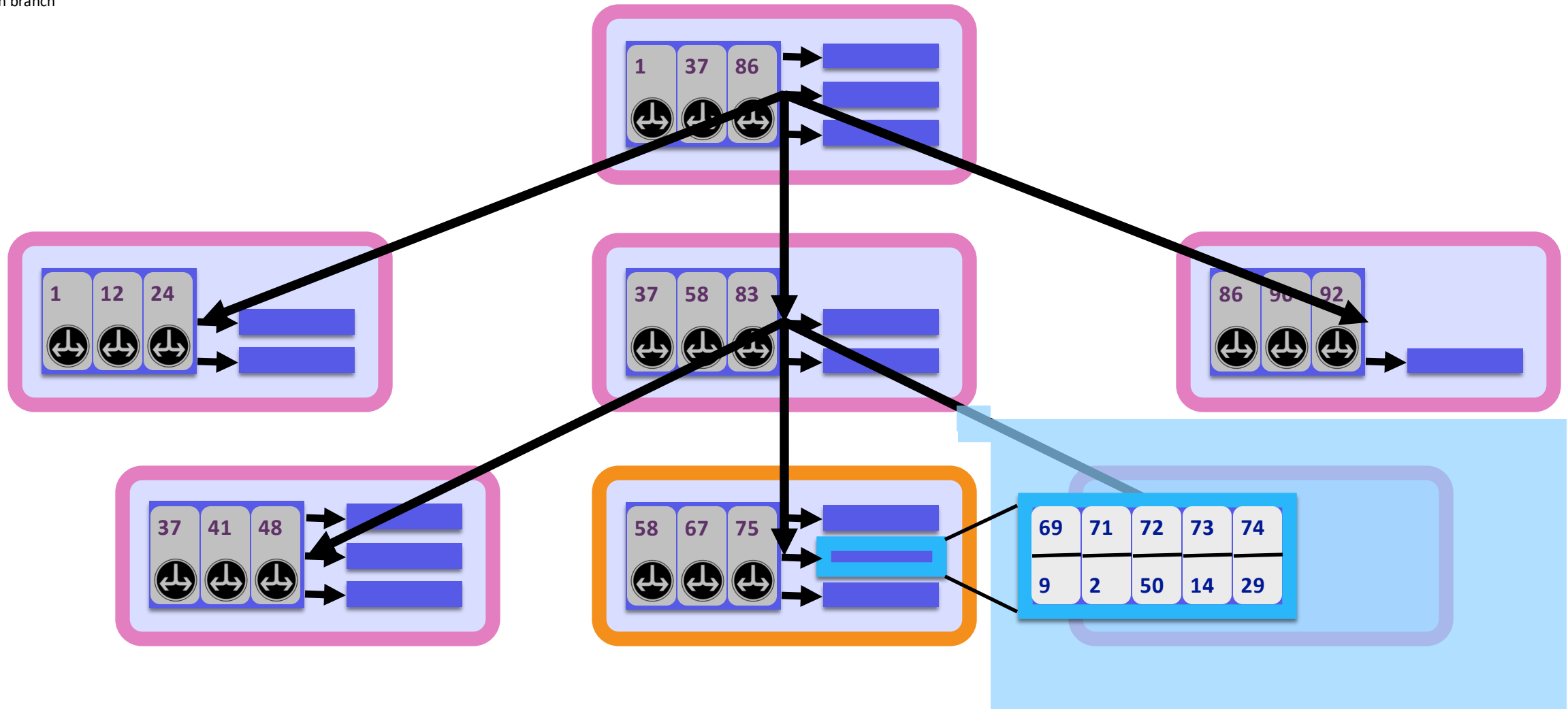
Query(71)



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

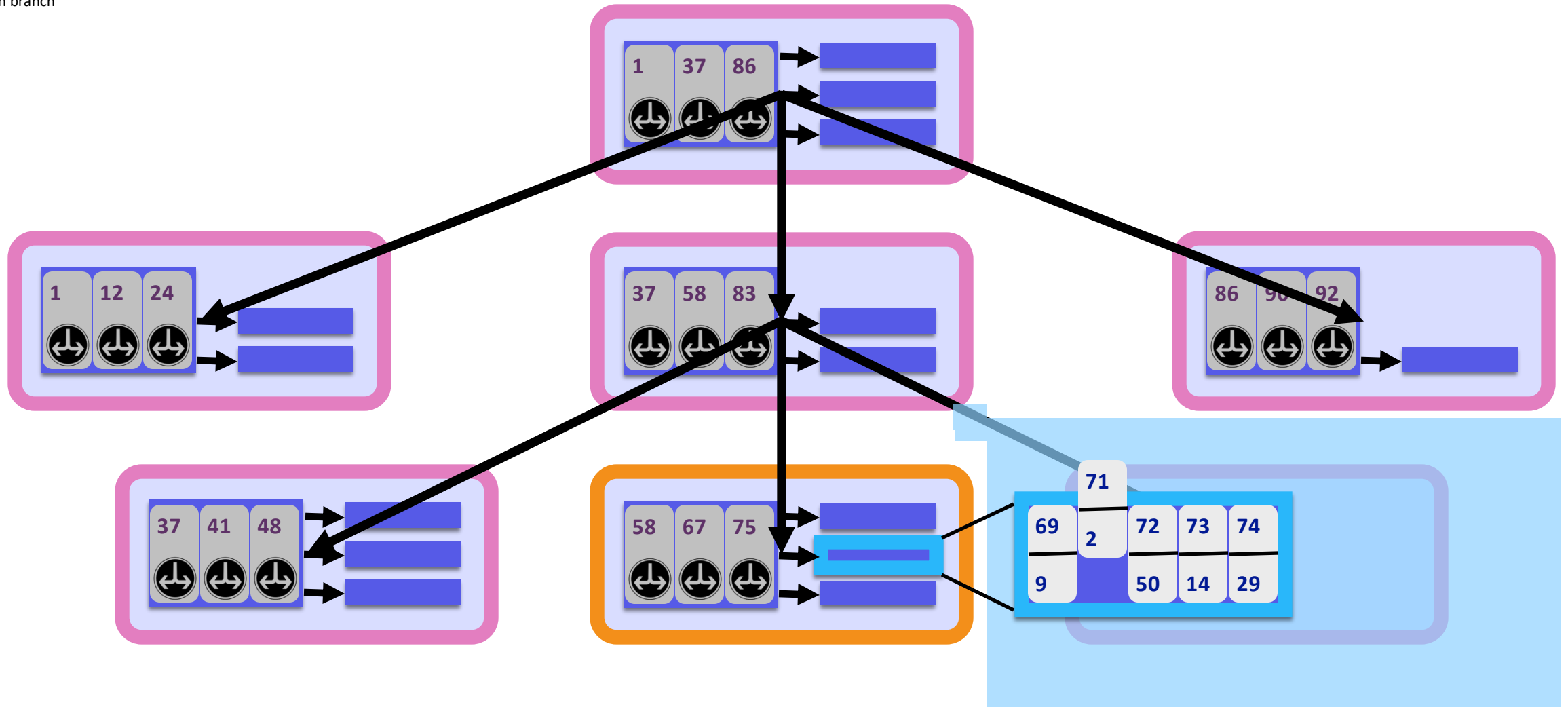
Query(71)



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

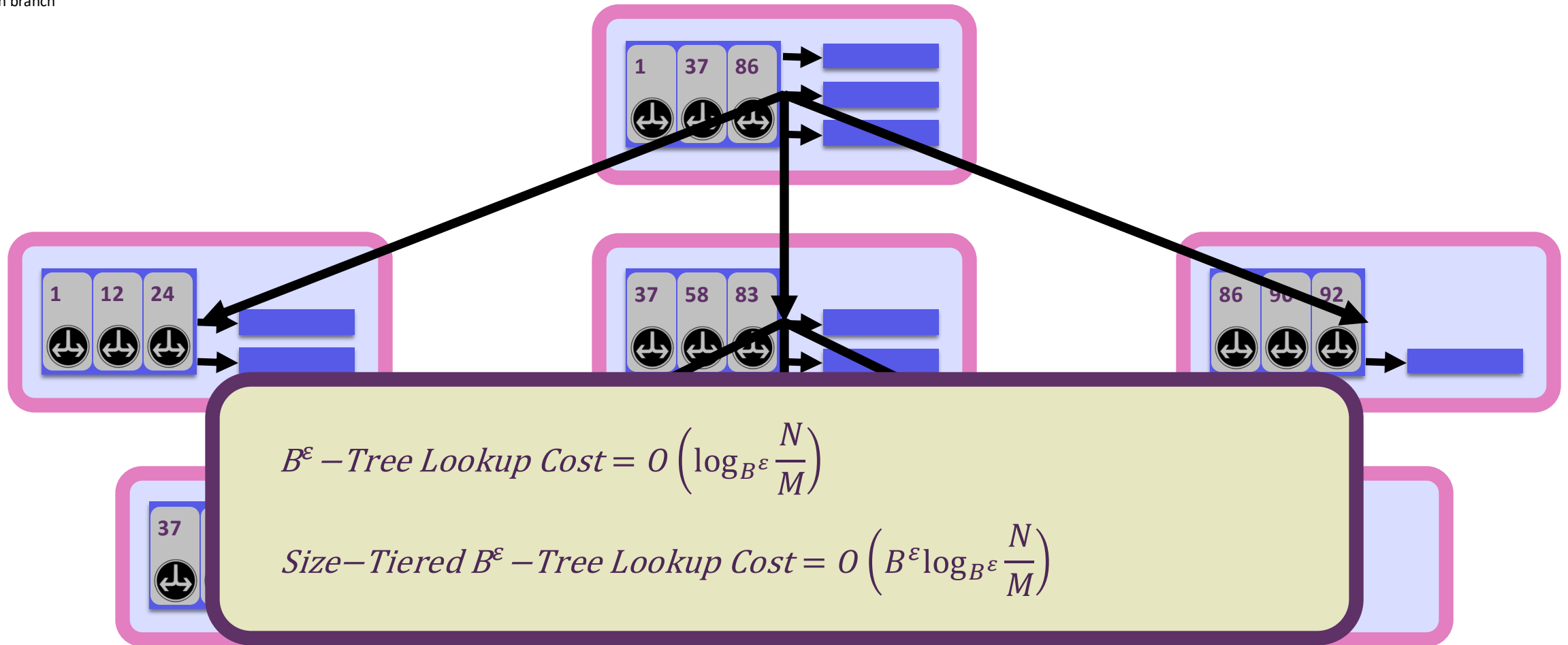
Query(71) → 2



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

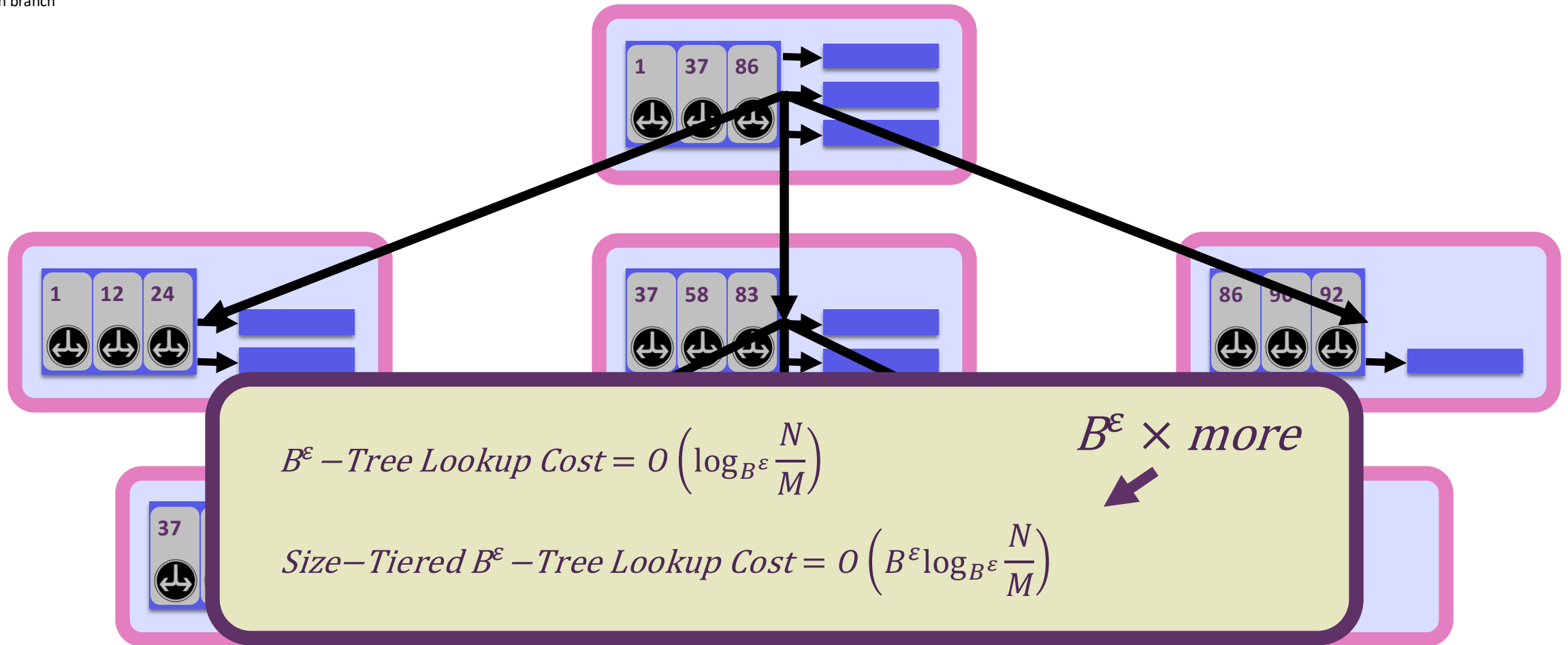
Query(71)



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

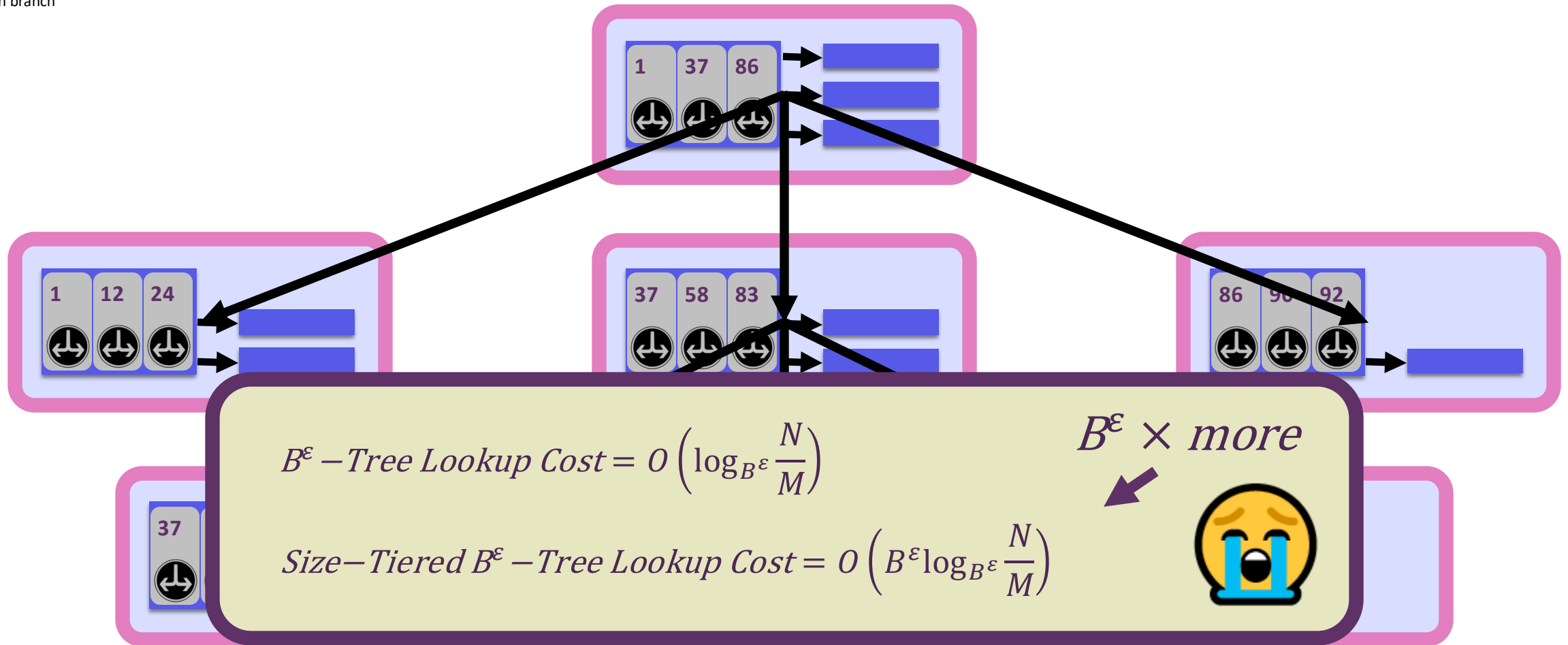
Query(71)



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

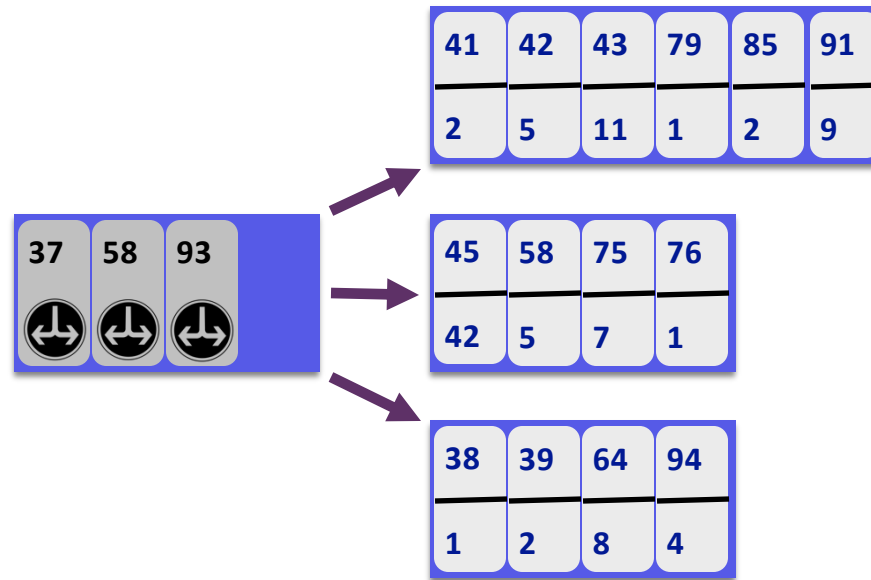
Query(71)



Fixing Lookups (almost)

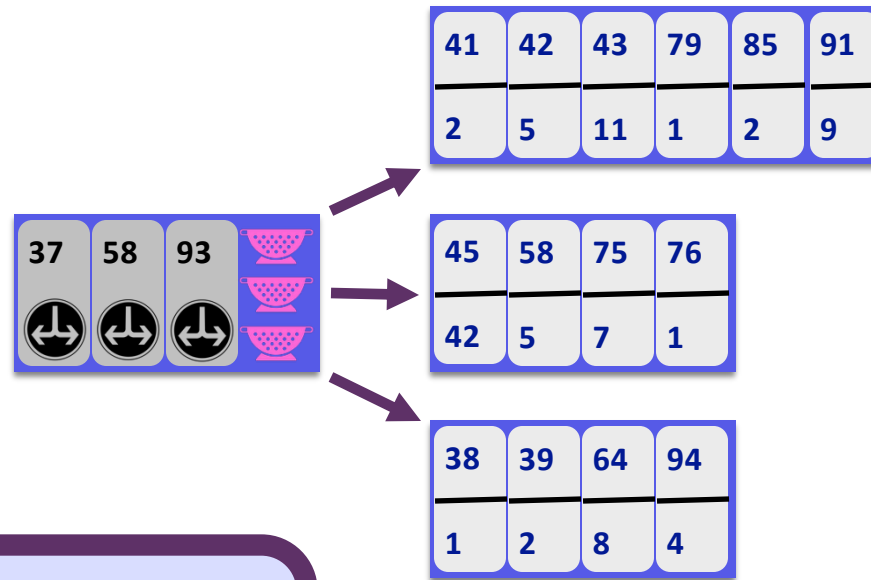
Fixing Lookups (almost)

The problem is that each node has multiple branches



Fixing Lookups (almost)

The problem is that each node has multiple branches



Idea: use filters to avoid searching them

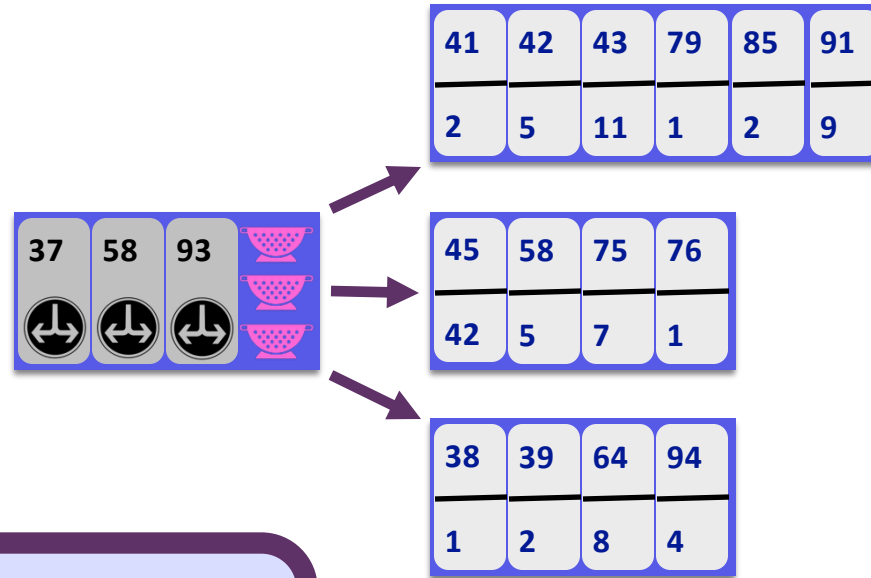


A filter is a probabilistic data structure which answers membership with no false negatives

Examples: Bloom, cuckoo, quotient

Fixing Lookups (almost)

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)



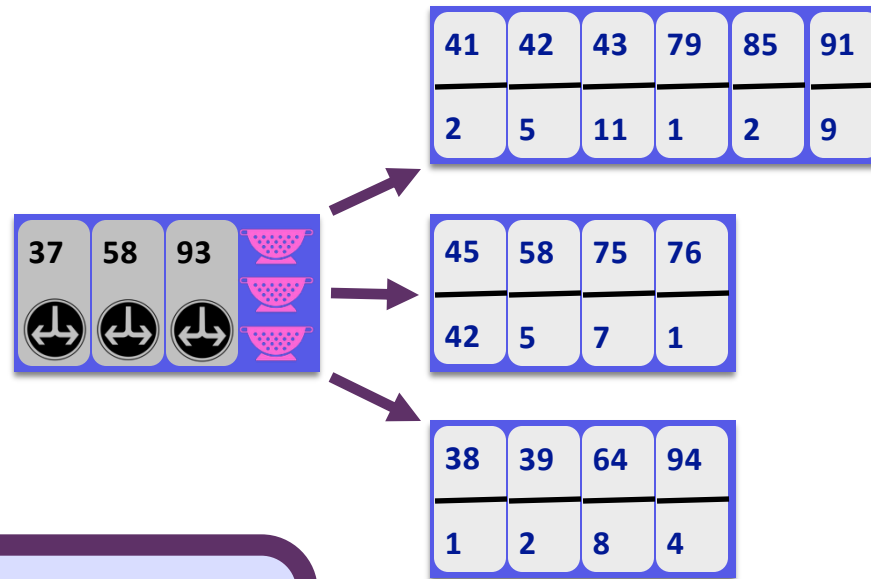
A filter is a probabilistic data structure which answers membership with no false negatives

Examples: Bloom, cuckoo, quotient

Fixing Lookups (almost)

Query(64)

The problem is that each node has multiple branches



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

Idea: use filters to avoid searching them



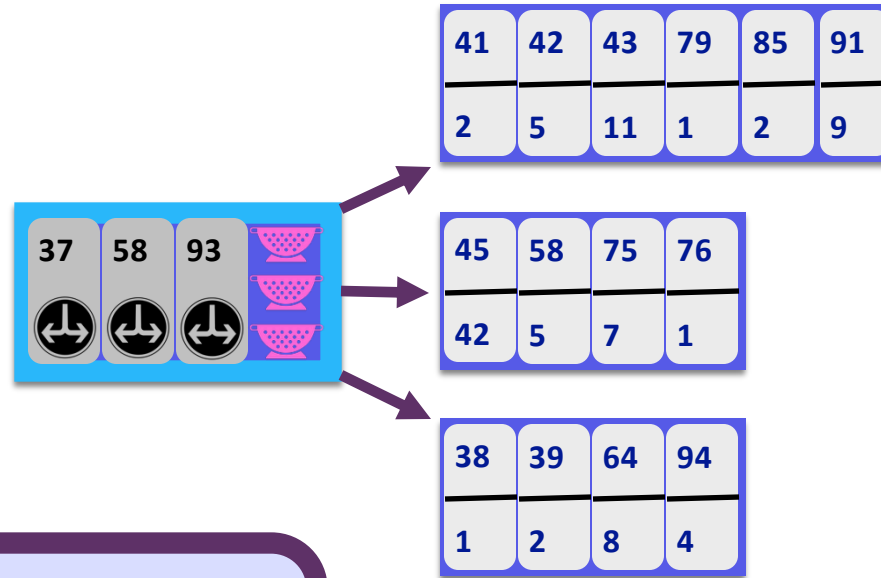
A filter is a probabilistic data structure which answers membership with no false negatives

Examples: Bloom, cuckoo, quotient

Fixing Lookups (almost)

Query(64)

The problem is that each node has multiple branches



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

Idea: use filters to avoid searching them



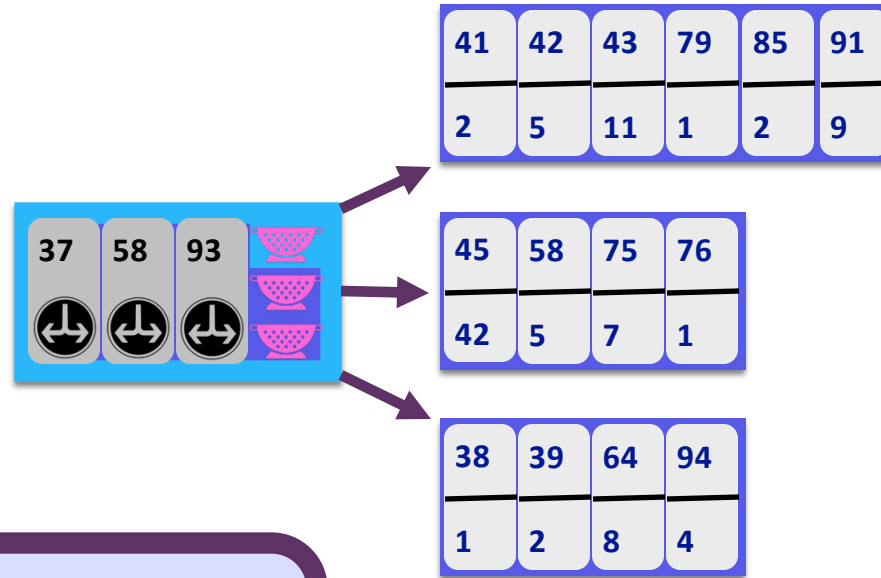
A filter is a probabilistic data structure which answers membership with no false negatives

Examples: Bloom, cuckoo, quotient

Fixing Lookups (almost)

Query(64)

The problem is that each node has multiple branches



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

Idea: use filters to avoid searching them



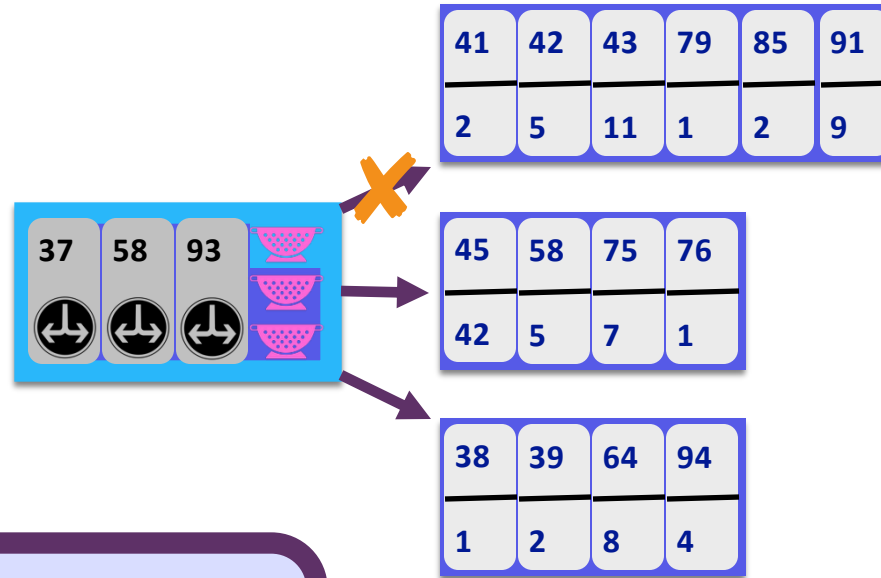
A filter is a probabilistic data structure which answers membership with no false negatives

Examples: Bloom, cuckoo, quotient

Fixing Lookups (almost)

Query(64)

The problem is that each node has multiple branches



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

Idea: use filters to avoid searching them



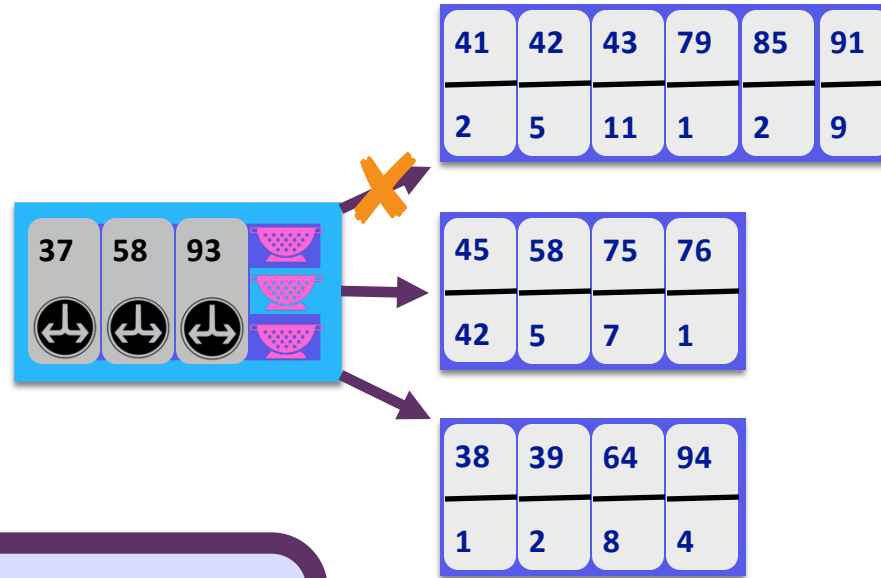
A filter is a probabilistic data structure with answers membership with no false negatives

Examples: Bloom, cuckoo, quotient

Fixing Lookups (almost)

Query(64)

The problem is that each node has multiple branches



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

Idea: use filters to avoid searching them



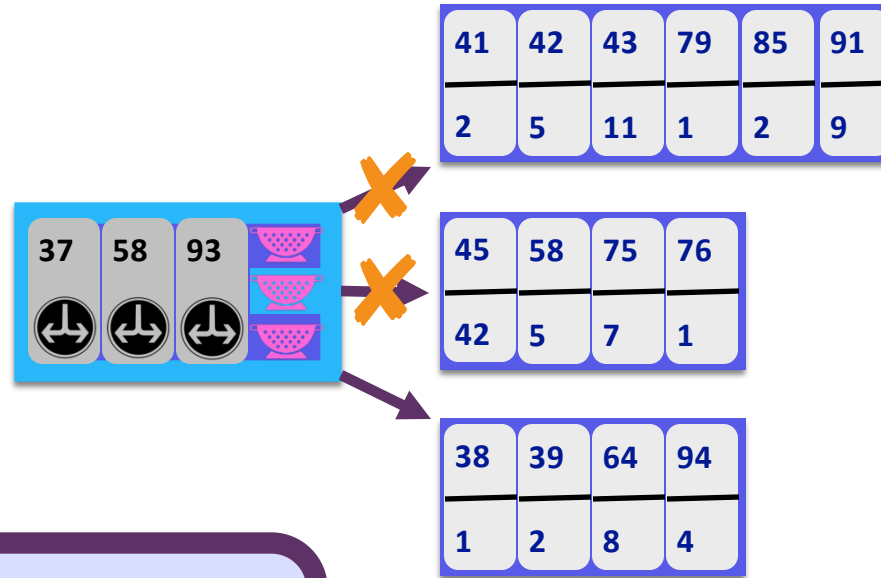
A filter is a probabilistic data structure with answers membership with no false negatives

Examples: Bloom, cuckoo, quotient

Fixing Lookups (almost)

Query(64)

The problem is that each node has multiple branches



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

Idea: use filters to avoid searching them

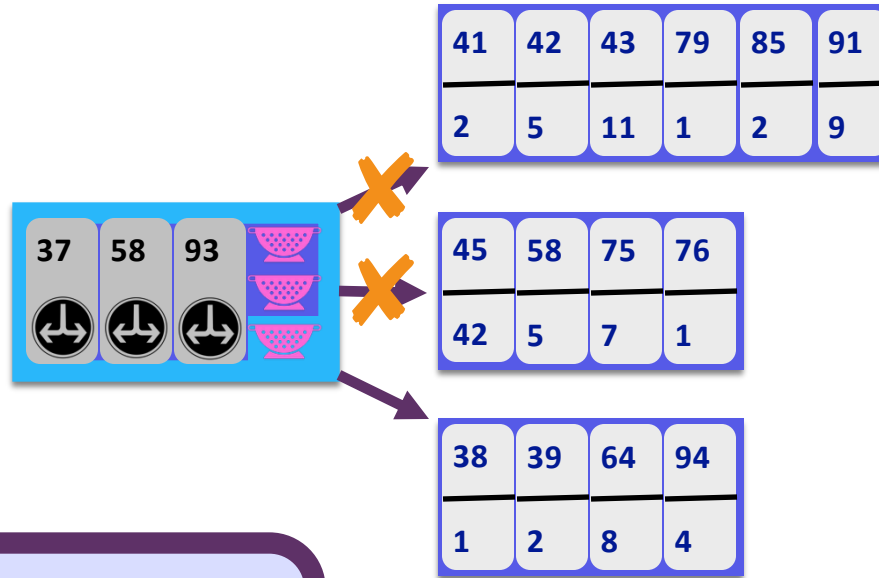


A filter is a probabilistic data structure with answers membership with no false negatives

Examples: Bloom, cuckoo, quotient

Query(64)

The problem is that each node has multiple branches



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

Idea: use filters to avoid searching them



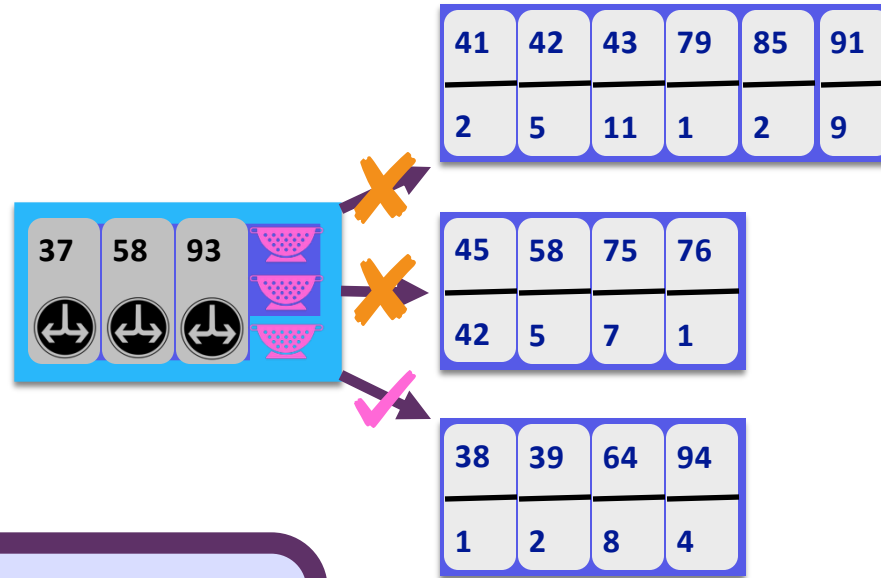
A filter is a probabilistic data structure with answers membership with no false negatives

Examples: Bloom, cuckoo, quotient

Fixing Lookups (almost)

Query(64)

The problem is that each node has multiple branches



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

Idea: use filters to avoid searching them



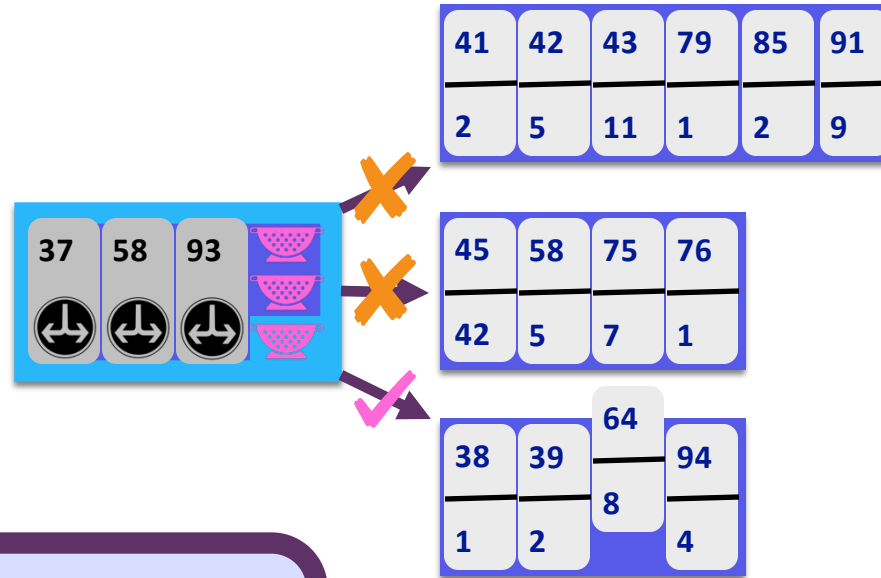
A filter is a probabilistic data structure with answers membership with no false negatives

Examples: Bloom, cuckoo, quotient

Fixing Lookups (almost)

Query(64) → 8

The problem is that each node has multiple branches



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

Idea: use filters to avoid searching them



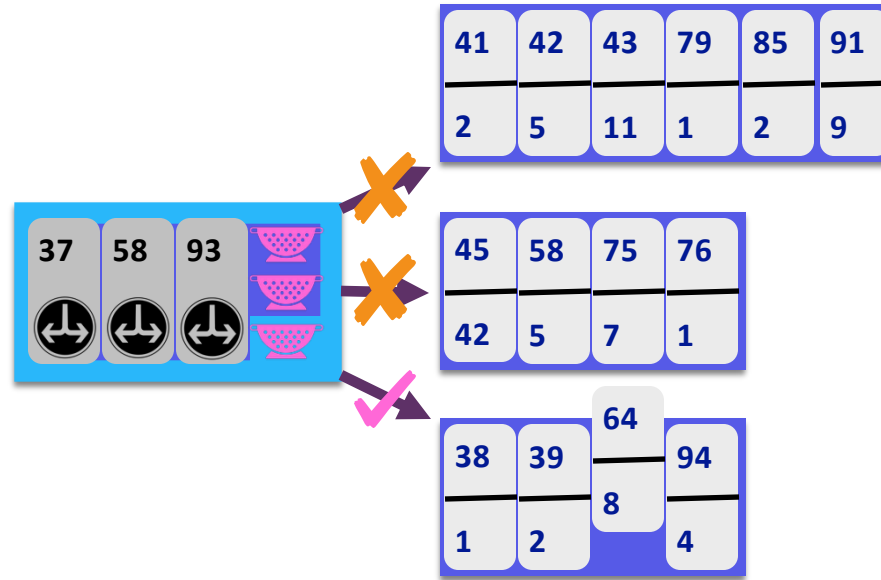
A filter is a probabilistic data structure with answers membership with no false negatives

Examples: Bloom, cuckoo, quotient

Fixing Lookups (almost)

Query(64) → 8

The problem is that each node has multiple branches



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

Idea: use filters to avoid searching them

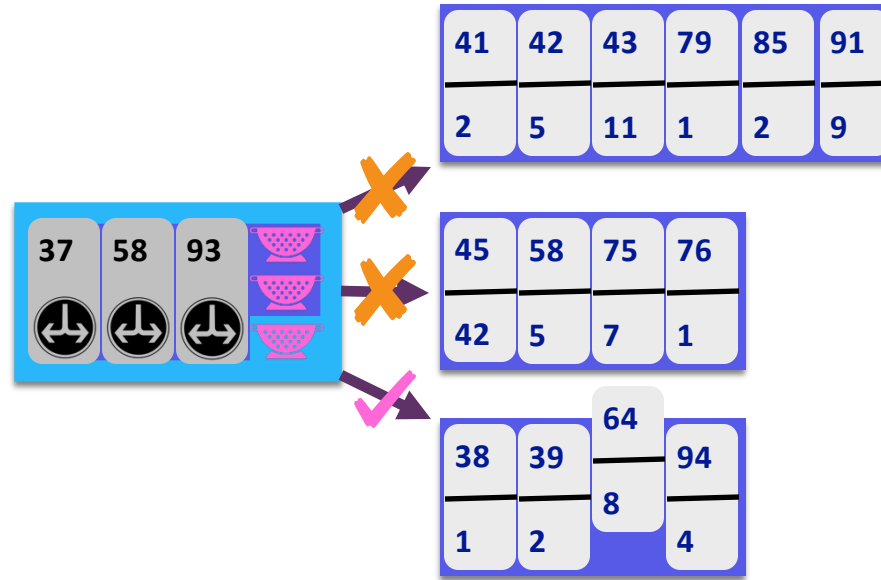
$$\text{False Positive Rate} \leq O\left(\frac{\epsilon}{B^\epsilon \log_B N}\right)$$

Fixing Lookups (almost)

Query(64) → 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)

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



Lookups in O(1) IOs

Conclusion

- B^e-trees are asymptotically faster than B-trees for insertions.
- They are appropriate for OLTP workloads
- Size-tiered B^e-trees help reduce write amplification
- Filter data structure can help reduce read amplification