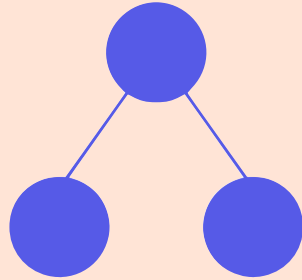# Lecture 10
# Write-Optimized Indexes

## Prashant Pandey

prashant.pandey@utah.edu
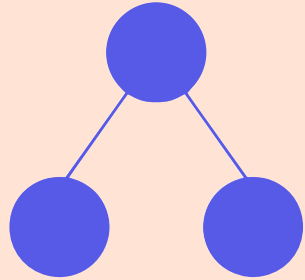
Slides taken from Prof. Alex Conway, Cornell Tech

# The Story of SplinterDB

Model the problem:
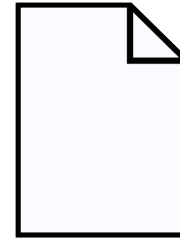external memory dictionary

# The Story of SplinterDB

Model the problem:
external memory dictionary

Metadata is fine-grained
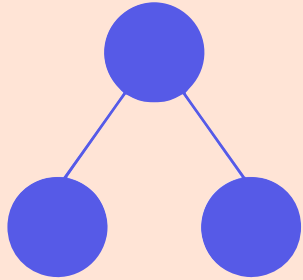
48 B

4 KiB

IO 4 KiB

# The Story of SplinterDB

Model the problem:
external memory dictionary

Metadata is fine-grained

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

Metadata is fine-grained

Model the problem:
external memory dictionary

4 KiB

48 B

IO 4 KiB



Here B is the number of items in an IO:
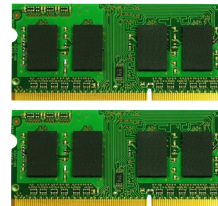B = 4 KiB / 48 B
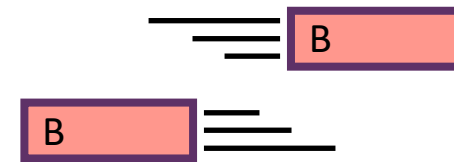
If the items were larger, the model wouldn't be as good

Internal Memory of size M

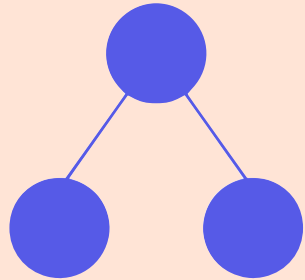A B-sized block can be read or written in 1 IO

**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

                              >
```

# Comparison-Based Dictionaries

Comparison External Memory
Model

<
user024299           =           user082587
>

Hashing          Filters

# Comparison-Based Dictionaries

## Comparison External Memory Model

<

user024299                =          user082587

>

Hashing

Filters

*Lower bounds for external memory dictionaries*,
Brodal, G., Fagerberg, R. SODA '03

## Brodal-Fagerberg Lower Bound

Insertions in

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

Lookups in

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

where $\lambda$ is a tuning parameter

# General Dictionaries

External Memory Model

# General Dictionaries

External Memory Model

user024299



YOU REALLY CAN DO WHATEVER YOU WANT

# General Dictionaries

External Memory Model

user024299

Hashing

XXH(user024299)

YOU REALLY CAN DO WHATEVER YOU WANT

# 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., Pătrașcu, M. SODA '12

## Iacono-Pătrașcu Lower Bound

Insertions in

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

Lookups in

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

where $\lambda$ is a tuning parameter

# Lower Bounds

Brodal-Fagerberg Lower Bound

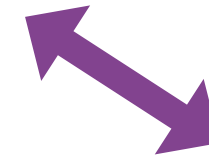Iacono-Pătrașcu Lower Bound

Insertions in

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

Lookups in

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

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

Iacono-Pătrașcu Lower Bound

Insertions in

Lookups in

Insertions in

Lookups in

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



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

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



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



B-Trees

$B^\varepsilon$-Trees

$$(\lambda = B)$$

$$(\lambda = B^\varepsilon)$$

Comparison External Memory Model

General External Memory Model

# Lower Bounds

Insertions in

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

Lookups in

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

Insertions in

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

Lookups in

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



B-Trees

$$(\lambda = B)$$



B$^\varepsilon$-Trees

$$(\lambda = B^\varepsilon)$$

Iacono-Patrascu Hash Table

BoA/BoT
Hash Table

Comparison External Memory Model

General External Memory Model

# Lower Bounds

Brodal-Fagerberg Lower Bound

Iacono-Pătrașcu Lower Bound

Insertions in

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

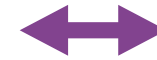Lookups in

$$\Omega(\log_\lambda N)$$
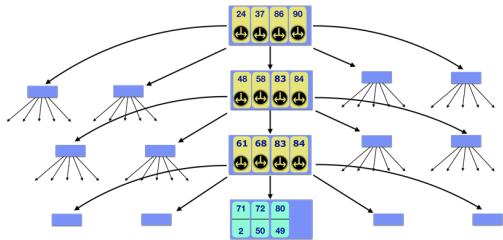
Insertions in

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

Lookups in

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



B-Trees

$$(\lambda = B)$$

$B^\varepsilon$-Trees

$$(\lambda = B^\varepsilon)$$

Iacono-Patrascu Hash Table

BoA/BoT
Hash Table

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

Comparison External Memory Model

# Lower Bounds

Insertions in

Lookups in

Insertions in

Lookups in

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

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

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

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



B-Trees

Bᵋ-Trees

Iacono-Patrascu Hash Table

No scans!

BoA/BoT
Hash Table

$$(\lambda = B)$$

$$(\lambda = B^\varepsilon)$$

Comparison External Memory Model

General External Memory Model

SCHOOL OF COMPUTING
UNIVERSITY OF UTAH

# Lower Bounds

Brodal-Fagerberg Lower Bound

Iacono-Pătrașcu Lower Bound

Insertions in $\qquad$ Lookups in $\qquad$ Insertions in $\qquad$ Lookups in

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



B-Trees

$$(\lambda = B)$$

B$^\varepsilon$-Trees

$$(\lambda = B^\varepsilon)$$

Iacono-Patrascu Hash Table

BoA/BoT
Hash Table

Mapped B$^\varepsilon$-Trees

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

Comparison External Memory Model

General External Memory Model

**SCHOOL OF COMPUTING**
UNIVERSITY OF UTAH
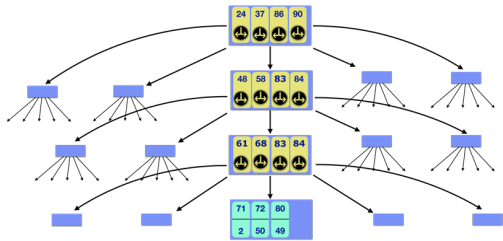
# Lower Bounds



Brodal-Fagerberg Lower Bound

Insertions in

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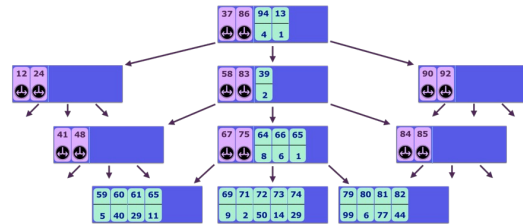
Lookups in

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

B-Trees

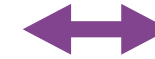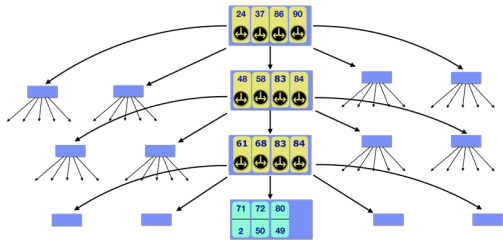$$(\lambda = B)$$

$B^\varepsilon$-Trees

$$(\lambda = B^\varepsilon)$$

Comparison External Memory Model

Iacono-Pătrașcu Lower Bound

Insertions in

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

Lookups in

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

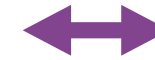Iacono-Patrascu Hash Table

Mapped $B^\varepsilon$-Trees

Hash Table

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

General External Memory Model

SCHOOL OF COMPUTING
UNIVERSITY OF UTAH

# 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



$$O(\log_B N)$$

# B-Trees

## B-ary Search Tree



Insert

| 76 |
|----|
| 6  |

$$O(\log_B N)$$

# B-Trees

B-ary Search Tree



$$O(\log_B N)$$

# B-Trees

## B-ary Search Tree

Insert

$$O(\log_B N)$$

# B-Trees

B-ary Search Tree

Insert



$$O(\log_B N)$$

# B-Trees

## B-ary Search Tree

Insert



$$O(\log_B N)$$

SCHOOL OF COMPUTING
UNIVERSITY OF UTAH

# B-Trees

## B-ary Search Tree

Insert



$$O(\log_B N)$$

# B-Trees

B-ary Search Tree

Insert

| 24 | 37 | 86 | 90 |

| 48 | 58 | 83 | 84 |

| 61 | 68 | 83 | 84 |

$O(\log_B N)$

**Insertion Cost $\leq O(\log_B N)$**

**Lookup Cost $\leq O(\log_B N)$**

# B$^\varepsilon$-Trees

# Bᵋ-Trees

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

Bᵋ pivots

the rest buffer

| 37 | 86 | | |

| 12 | 24 | | |

| 58 | 83 | | |

| 90 | 92 | | |

| 41 | 48 | | |

| 67 | 75 | | |

| 84 | 85 | | |

| 59 | 60 | 61 | 65 |
|----|----|----|----|
| 5  | 40 | 29 | 11 |

| 69 | 71 | 72 | 73 | 74 |
|----|----|----|----|----|
| 9  | 2  | 50 | 14 | 29 |

| 79 | 80 | 81 | 82 |
|----|----|----|----|
| 99 | 6  | 77 | 44 |

SCHOOL OF COMPUTING
UNIVERSITY OF UTAH

55

# Bᵉ-Trees

Inserts get put in the root buffer

SCHOOL OF COMPUTING
UNIVERSITY OF UTAH

# Bᵋ-Trees

Inserts get put in the root buffer

# Bᵋ-Trees

Inserts get put in the root buffer

# Bᵋ-Trees

Inserts get put in the root buffer

# Bᵋ-Trees

Inserts get put in the root buffer

# Bᵋ-Trees

Inserts get put in the root buffer

61

# Bᵋ-Trees

Inserts get put in the root buffer

# Bᵋ-Trees

Inserts get put in the root buffer

# Bᵋ-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

# Bᵋ-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

SCHOOL OF COMPUTING
UNIVERSITY OF UTAH

# Bᵋ-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

# Bᵋ-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-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$^\varepsilon$-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

# Bᵋ-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

# Bᵉ-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

# Bᵋ-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

# Bᵋ-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

SCHOOL OF COMPUTING
UNIVERSITY OF UTAH

# B$^\varepsilon$-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

SCHOOL OF COMPUTING
UNIVERSITY OF UTAH

# Lookups in B$^\varepsilon$-Trees

# Bᵋ-Trees

Query(71)

Lookups follow pivots, but check buffers along the way

# B$^\varepsilon$-Trees

Query(71)

Lookups follow pivots, but check buffers along the way

# B$^\varepsilon$-Trees

Query(71)

Lookups follow pivots, but check buffers along the way

# Bᵋ-Trees



Query(71)

Lookups follow pivots, but check buffers along the way

79

# Bᵋ-Trees

Lookups follow pivots, but check buffers along the way

Query(71) → 2

# Bᵋ-Trees

Query(71) → 2

Lookups follow pivots, but check buffers along the way

# Insertions in B$^\varepsilon$-Trees are more expensive than they look

# Insertions in B$^\varepsilon$-Trees are more expensive than they look

# (Also most LSMS)

# Insertions in B$^\varepsilon$-Trees are more expensive than they look
## Recall: Insertions in B$^\varepsilon$-trees

| 65 | 72 | 80 |
|----|----|----|
| 11 | 50 | 6  |

# Insertions in B$^\varepsilon$-Trees are more expensive than they look

## Recall: Insertions in B$^\varepsilon$-trees

| 65 | 72 | 80 |
|----|----|----|
| 11 | 50 | 6  |

| 58 | 83 | 39 | 64 | 66 |
|----|----|----|----|----|
|    |    | 2  | 8  | 6  |

Read the node

# Insertions in B$^\varepsilon$-Trees are more expensive than they look

## Recall: Insertions in B$^\varepsilon$-trees

| 65 | 72 | 80 |
|----|----|----|
| 11 | 50 | 6  |

**Merge the data**

| 58 | 83 | 39 | 64 | 66 |
|----|----|----|----|----|
|    |    | 2  | 8  | 6  |

**Read the node**

SCHOOL OF COMPUTING
UNIVERSITY OF UTAH

# Insertions in B$^\varepsilon$-Trees are more expensive than they look

## Recall: Insertions in B$^\varepsilon$-trees



Merge the data

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

Read the node

# Insertions in B$^\varepsilon$-Trees are more expensive than they look

## Recall: Insertions in B$^\varepsilon$-trees



Merge the data

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

Read the node

Write the node

SCHOOL OF COMPUTING
UNIVERSITY OF UTAH

# Insertions in B$^\varepsilon$-Trees are more expensive than they look

## Recall: Insertions in B$^\varepsilon$-trees

Merge the data

**CPU Work** = O(old + new)

**Volume of IO** = O(old + new)

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

Read the node

Write the node

SCHOOL OF COMPUTING
UNIVERSITY OF UTAH

# 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  |

**Read the node**

**Write the node**

**CPU Work** = O(old + new)

**Volume of IO** = O(old + new)

Older data gets written over and over again

# Insertions in B$^\varepsilon$-Trees are more expensive than they look

## Recall: Insertions in B$^\varepsilon$-trees

**Merge the data**

**CPU Work** = O(old + new)

**Volume of IO** = O(old + new)

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

Older data gets written over and over again

**Read the node**

**Write the node**

# Insertions in B$^\varepsilon$-Trees are more expensive than they look

## Recall: Insertions in B$^\varepsilon$-trees



**CPU Work** = O(old + new)

**Volume of IO** = O(old + new)

Older data gets written over and over again

Merge the data

Read the node

Write the node

SCHOOL OF COMPUTING
UNIVERSITY OF UTAH

# Insertions in B$^\varepsilon$-Trees are more expensive than they look

## Recall: Insertions in B$^\varepsilon$-trees



Merge the data

**CPU Work** = O(old + new)

**Volume of IO** = O(old + new)

Older data gets written over and over again

Read the node

Write the node

# Insertions in Bᵉ-Trees are more expensive than they look

**Merge the data**

**Read the node**

**Write the node**

**CPU Work** = O(old + new)

**Volume of IO** = O(old + new)

Older data gets written over and over again

Up to $B^\varepsilon$ times per node!

# Size-Tiered B$^\varepsilon$-Trees

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

# Size-Tiered B$^\varepsilon$-Trees

A Size-Tiered B$^\varepsilon$-tree is a B$^\varepsilon$-tree where the buffer is stored discontiguously

B$^\varepsilon$ pivots          the rest buffer

| 37 | 58 | 93 | |
|----|----|----|----|

Recall:
a B$^\varepsilon$-tree node has pivots and a buffer

SCHOOL OF COMPUTING
UNIVERSITY OF UTAH

# Size-Tiered B$^\varepsilon$-Trees

A Size-Tiered B$^\varepsilon$-tree is a B$^\varepsilon$-tree where the buffer is stored discontiguously

Recall:
a B$^\varepsilon$-tree node has pivots and a buffer

B$^\varepsilon$ pivots

the rest buffer



In an STB$^\varepsilon$-tree, the buffer is stored separately

SCHOOL OF COMPUTING
UNIVERSITY OF UTAH

# Size-Tiered B$^\varepsilon$-Trees

A Size-Tiered B$^\varepsilon$-tree is a B$^\varepsilon$-tree where the buffer is stored discontiguously

B$^\varepsilon$ pivots

the rest buffer

Recall:
a B$^\varepsilon$-tree node has pivots and a buffer

| 37 | 58 | 93 |

and in several discontiguous pieces

In an STB$^\varepsilon$-tree, the buffer is stored separately

# Size-Tiered B$^\varepsilon$-Trees

A Size-Tiered B$^\varepsilon$-tree is a B$^\varepsilon$-tree where the buffer is stored discontiguously

trunk [node]

B$^\varepsilon$ pivots

the rest buffer

Recall:
a B$^\varepsilon$-tree node has pivots and a buffer

| 37 | 58 | 93 |

and in several discontiguous pieces

In an STB$^\varepsilon$-tree, the buffer is stored separately

# Size-Tiered B$^\varepsilon$-Trees

A Size-Tiered B$^\varepsilon$-tree is a B$^\varepsilon$-tree where the buffer is stored discontiguously

trunk [node]

branches

B$^\varepsilon$ pivots

the rest buffer

| 37 | 58 | 93 |

Recall:
a B$^\varepsilon$-tree node has pivots and a buffer

and in several discontiguous pieces

In an STB$^\varepsilon$-tree, the buffer is stored separately

# Insertions in Size-Tiered B$^{\varepsilon}$-Trees

# Size-Tiered Bᵋ-Trees

A Size-Tiered Bᵋ-tree is a Bᵋ-tree where the buffer is stored discontiguously

When new data is flushed into the trunk node…

# Size-Tiered Bᵋ-Trees

A Size-Tiered Bᵋ-tree is a Bᵋ-tree where the buffer is stored discontiguously



When new data is flushed into the trunk node…

# Size-Tiered Bᵉ-Trees

A Size-Tiered $B^\varepsilon$-tree is a $B^\varepsilon$-tree where the buffer is stored discontiguously

When new data is flushed into the trunk node…

…it is added as a new branch

# Size-Tiered Bᵋ-Trees

A Size-Tiered Bᵋ-tree is a Bᵋ-tree where the buffer is stored discontiguously



When new data is flushed into the trunk node…

…it is added as a new branch

# Size-Tiered B$^\epsilon$-Trees

A Size-Tiered B$^\epsilon$-tree is a B$^\epsilon$-tree where the buffer is stored discontiguously

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

...it is added as a new branch

| 37 | 58 | 93 |
|----|----|----|

| 45 | 58 | 75 | 76 |
|----|----|----|----|
| 42 | 5  | 7  | 1  |

The old branches do not need to be rewritten

| 38 | 39 | 64 | 94 |
|----|----|----|----|
| 1  | 2  | 8  | 4  |

# Size-Tiered Bᵋ-Trees

A Size-Tiered Bᵋ-tree is a Bᵋ-tree where the buffer is stored discontiguously

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$^\varepsilon$-Trees

A Size-Tiered B$^\varepsilon$-tree is a B$^\varepsilon$-tree where the buffer is stored discontiguously

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

SCHOOL OF COMPUTING
UNIVERSITY OF UTAH

# Size-Tiered B$^\varepsilon$-Trees

A Size-Tiered B$^\varepsilon$-tree is a B$^\varepsilon$-tree where the buffer is stored discontiguously

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 discontiguously

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$^\varepsilon$-Trees

A Size-Tiered B$^\varepsilon$-tree is a B$^\varepsilon$-tree where the buffer is stored discontiguously

When the node is full:
1. Pick child receiving most messages
2. Merge them into a new branch for the child

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

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

...it is added as a new branch

Branches may have overlapping key ranges

| 37 | 58 | 93 |
|----|----|----|

| 45 | 58 | 75 | 76 |
|----|----|----|----|
| 42 | 5  | 7  | 1  |

The old branches do not need to be rewritten

| 37 | 45 | 52 |
|----|----|----|

| 58 | 83 |
|----|----|

| 93 | 99 |
|----|----|

| 38 | 39 | 64 | 94 |
|----|----|----|----|
| 1  | 2  | 8  | 4  |

# Size-Tiered Bᵋ-Trees

A Size-Tiered Bᵋ-tree is a Bᵋ-tree where the buffer is stored discontiguously

When the node is full:
1. Pick child receiving most messages
2. Merge them into a new branch for the child



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

Branches may have overlapping key ranges

# Size-Tiered B$^\varepsilon$-Trees

A Size-Tiered B$^\varepsilon$-tree is a B$^\varepsilon$-tree where the buffer is stored discontiguously

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 Bᵋ-Trees

A Size-Tiered Bᵋ-tree is a Bᵋ-tree where the buffer is stored discontiguously

When the node is full:
1. Pick child receiving most messages
2. Merge them into a new branch for the child

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

Branches may have overlapping key ranges

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



114

# Lookups in Size-Tiered B$^\varepsilon$-Trees

# Size-Tiered B$^{\varepsilon}$-Trees

Query(71)

Lookups in a STB$^{\varepsilon}$-tree are like lookups in a B$^{\varepsilon}$-tree, except they must check each branch

# Size-Tiered Bᵋ-Trees

Query(71)

Lookups in a STBᵋ-tree are like lookups in a Bᵋ-tree, except they must check each branch

# Size-Tiered B$^\varepsilon$-Trees

Query(71)

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

# Size-Tiered B$^\varepsilon$-Trees

Query(71)

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

# Size-Tiered Bᵋ-Trees

Query(71)

Lookups in a STBᵋ-tree are like lookups in a Bᵋ-tree, except they must check each branch

# Size-Tiered Bᵋ-Trees

Query(71)

Lookups in a STBᵋ-tree are like lookups in a Bᵋ-tree, except they must check each branch

# Size-Tiered Bᵋ-Trees

Query(71)

Lookups in a STBᵋ-tree are like lookups in a Bᵋ-tree, except they must check each branch

# Size-Tiered B$^\varepsilon$-Trees

Query(71)

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

# Size-Tiered B<sup>ε</sup>-Trees



Query(71)

Lookups in a STB<sup>ε</sup>-tree are like lookups in a B<sup>ε</sup>-tree, except they must check each branch

# Size-Tiered B$^\varepsilon$-Trees

Query(71)

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

# Size-Tiered Bᵋ-Trees

Query(71)

Lookups in a STBᵋ-tree are like lookups in a Bᵋ-tree, except they must check each branch

# Size-Tiered Bᵉ-Trees



Query(71)

Lookups in a STBᵉ-tree are like lookups in a Bᵉ-tree, except they must check each branch

# Size-Tiered B$^\varepsilon$-Trees

Query(71) → 2

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

128

# Size-Tiered B$^\varepsilon$-Trees

Query(71)

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

| 1 | 37 | 86 |
|---|----|----|

| 1 | 12 | 24 |
|---|----|----|

| 37 | 58 | 83 |
|----|----|----|

| 86 | 90 | 92 |
|----|----|----|

| 37 |
|----|

$$B^{\varepsilon}-Tree\ Lookup\ Cost = O\left(\log_{B^\varepsilon}\frac{N}{M}\right)$$

$$Size-Tiered\ B^{\varepsilon}-Tree\ Lookup\ Cost = O\left(B^{\varepsilon}\log_{B^\varepsilon}\frac{N}{M}\right)$$

# Size-Tiered B$^\varepsilon$-Trees

Query(71)

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



$$B^\varepsilon - Tree\ Lookup\ Cost = O\left(\log_{B^\varepsilon}\frac{N}{M}\right)$$

$$B^\varepsilon \times more$$

$$Size - Tiered\ B^\varepsilon - Tree\ Lookup\ Cost = O\left(B^\varepsilon\log_{B^\varepsilon}\frac{N}{M}\right)$$

# Size-Tiered B$^\varepsilon$-Trees

Query(71)

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



$$B^\varepsilon - Tree\ Lookup\ Cost = O\left(\log_{B^\varepsilon} \frac{N}{M}\right)$$

$$B^\varepsilon \times more$$

$$Size - Tiered\ B^\varepsilon - Tree\ Lookup\ Cost = O\left(B^\varepsilon \log_{B^\varepsilon} \frac{N}{M}\right)$$

# Fixing Lookups (almost)

SCHOOL OF COMPUTING
UNIVERSITY OF UTAH

# 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

| 37 | 58 | 93 | 🔽 |
| --- | --- | --- | --- |

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

| 45 | 58 | 75 | 76 |
| --- | --- | --- | --- |
| 42 | 5 | 7 | 1 |

| 38 | 39 | 64 | 94 |
| --- | --- | --- | --- |
| 1 | 2 | 8 | 4 |

A filter is a probabilistic data structure with 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

| 37 | 58 | 93 | |

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

| 45 | 58 | 75 | 76 |
|----|----|----|----|
| 42 | 5  | 7  | 1  |

| 38 | 39 | 64 | 94 |
|----|----|----|----|
| 1  | 2  | 8  | 4  |

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 (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)

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

| 37 | 58 | 93 |
|----|----|----|

| 45 | 58 | 75 | 76 |
|----|----|----|----|
| 42 | 5  | 7  | 1  |

Idea: use filters to avoid searching them

| 38 | 39 | 64 | 94 |
|----|----|----|----|
| 1  | 2  | 8  | 4  |

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

Idea: use filters to avoid searching them

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

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

| 37 | 58 | 93 | 🥄 |
|----|----|----|----|

| 45 | 58 | 75 | 76 |
|----|----|----|----|
| 42 | 5  | 7  | 1  |

| 38 | 39 | 64 | 94 |
|----|----|----|----|
| 1  | 2  | 8  | 4  |

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

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

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

| 37 | 58 | 93 | 🍳 |

| 45 | 58 | 75 | 76 |
|----|----|----|----|
| 42 | 5  | 7  | 1  |

Idea: use filters to avoid searching them

| 38 | 39 | 64 | 94 |
|----|----|----|----|
| 1  | 2  | 8  | 4  |

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

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

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

| 37 | 58 | 93 | |

| 45 | 58 | 75 | 76 |
|----|----|----|----|
| 42 | 5  | 7  | 1  |

Idea: use filters to avoid searching them

| 38 | 39 | 64 | 94 |
|----|----|----|----|
| 1  | 2  | 8  | 4  |

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



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

| 37 | 58 | 93 | |
|----|----|----|---|

| 45 | 58 | 75 | 76 |
|----|----|----|----|
| 42 | 5 | 7 | 1 |

| 38 | 39 | 64 | 94 |
|----|----|----|----|
| 1 | 2 | 8 | 4 |

Now a lookup will only search those branches which contain the key (plus rare false positives)
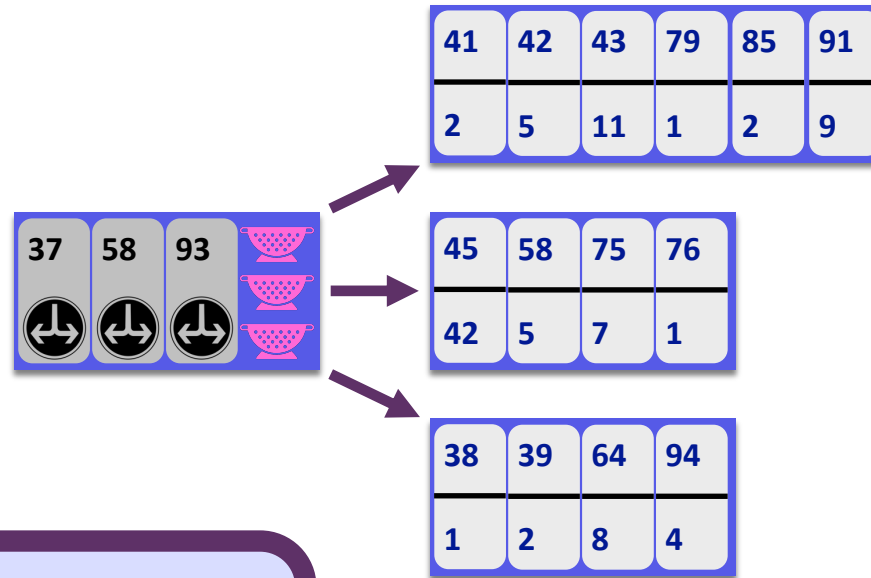
Idea: use filters to avoid searching them

A filter is a probabilistic data structure with answers membership with no false negatives
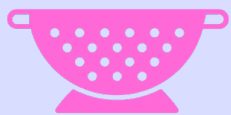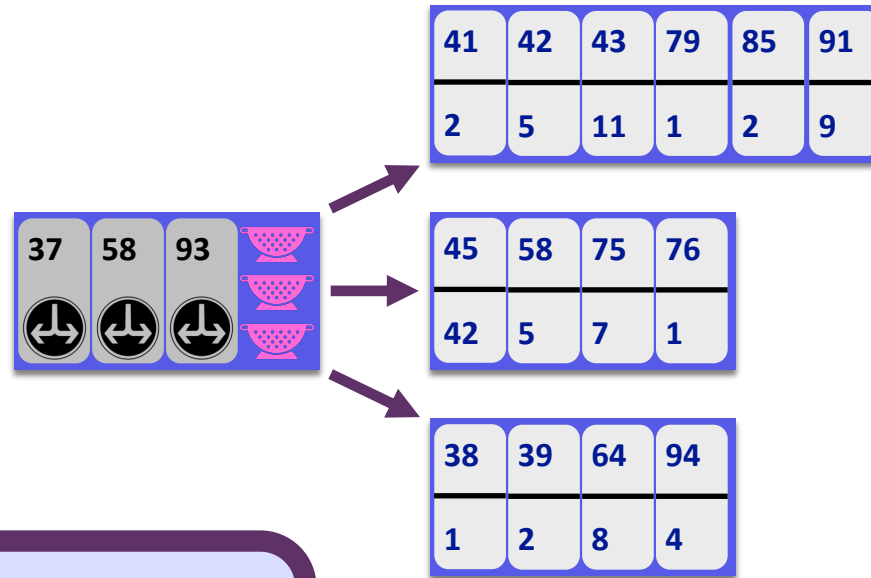
Examples: Bloom, cuckoo, quotient
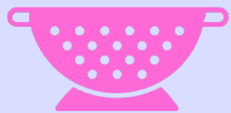
# Fixing Lookups (almost)

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)

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

| 45 | 58 | 75 | 76 |
|----|----|----|----|
| 42 | 5  | 7  | 1  |

| 37 | 58 | 93 |
|----|----|----|

| 38 | 39 | 64 | 94 |
|----|----|----|----|
| 1  | 2  | 8  | 4  |

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

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

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

| 37 | 58 | 93 | 🧺 |
|----|----|----|----|

| 45 | 58 | 75 | 76 |
|----|----|----|----|
| 42 | 5  | 7  | 1  |

Idea: use filters to avoid searching them

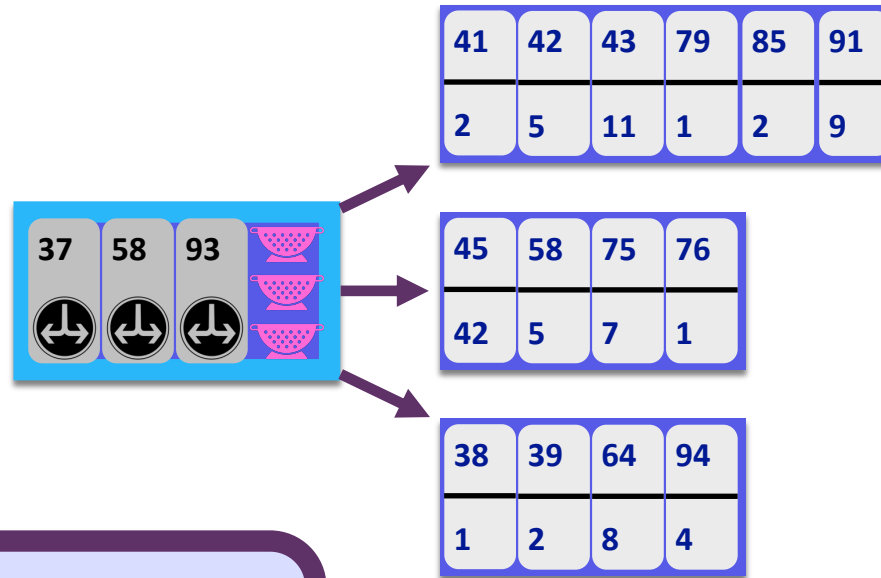| 38 | 39 | 64 | 94 |
|----|----|----|----|
| 1  | 2  | 8  | 4  |

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

Examples: Bloom, cuckoo, quotient
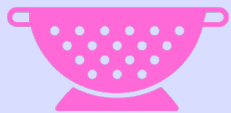
# Fixing Lookups (almost)

Query(64)

The problem is that each node has multiple branches

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

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

| 37 | 58 | 93 | |
|----|----|----|--|

| 45 | 58 | 75 | 76 |
|----|----|----|----|
| 42 | 5  | 7  | 1  |

Idea: use filters to avoid searching them

| 38 | 39 | 64 | 94 |
|----|----|----|----|
| 1  | 2  | 8  | 4  |

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

Idea: use filters to avoid searching them

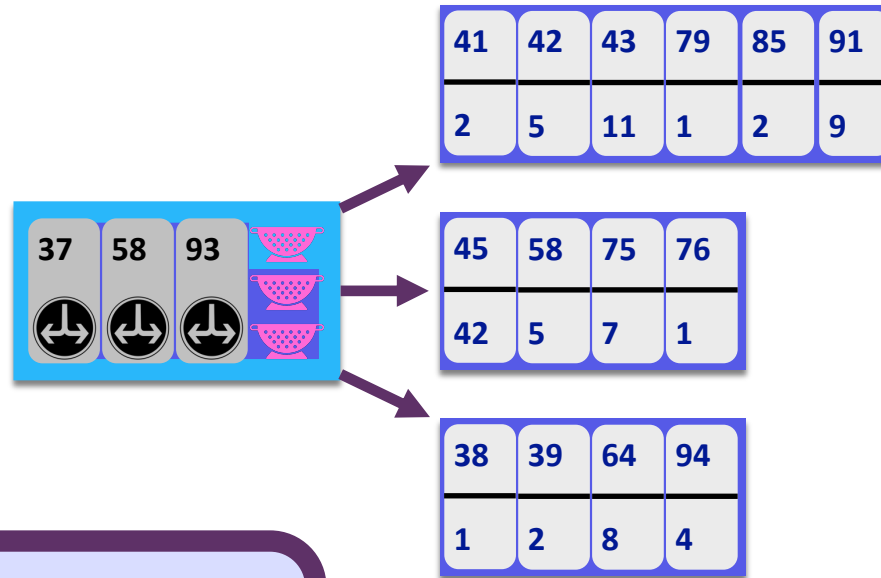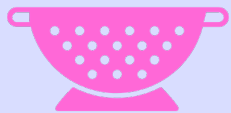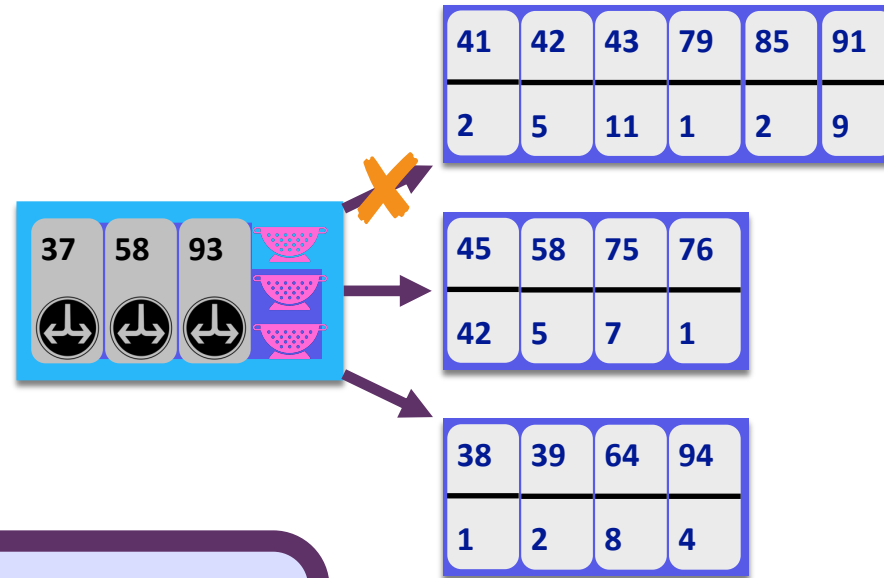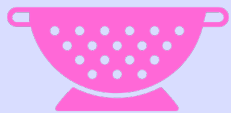Now a lookup will only search those branches which contain the key (plus rare false positives)

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

| 37 | 58 | 93 | |
|----|----|----| |

| 45 | 58 | 75 | 76 |
|----|----|----|----|
| 42 | 5  | 7  | 1  |

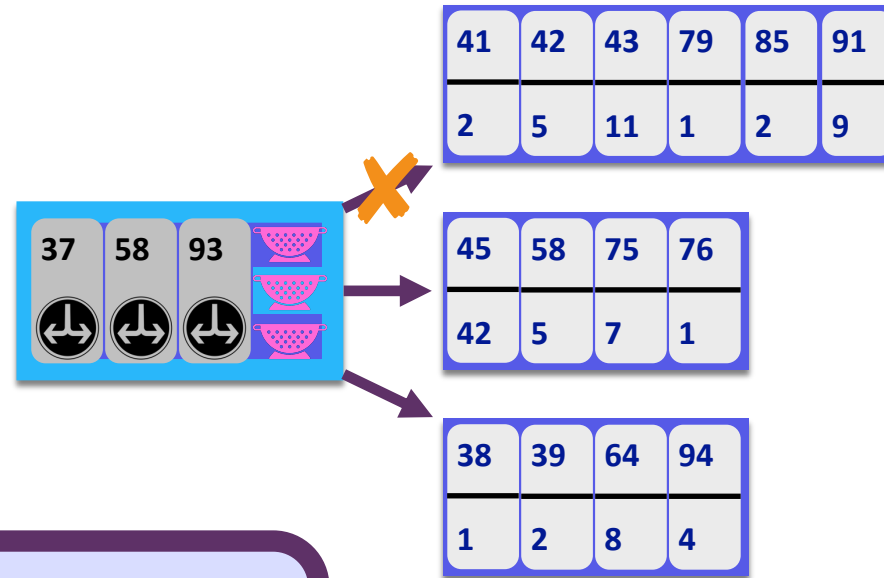| 38 | 39 | 64 | 94 |
|----|----|----|----|
| 1  | 2  | 8  | 4  |

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

Examples: Bloom, cuckoo, quotient
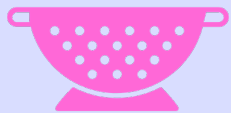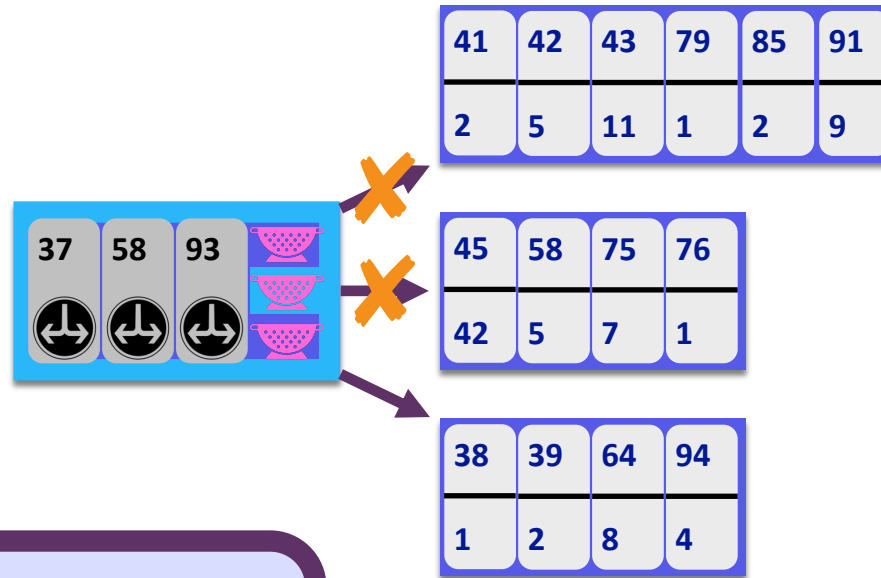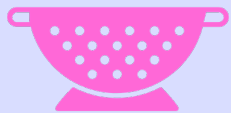
# Fixing Lookups (almost)

Query(64) → 8

The problem is that each node has multiple branches

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

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

| 37 | 58 | 93 | |

| 45 | 58 | 75 | 76 |
|----|----|----|----|
| 42 | 5  | 7  | 1  |

Idea: use filters to avoid searching them

| 38 | 39 | 64 | 94 |
|----|----|----|----|
| 1  | 2  | 8  | 4  |

$$False\ Positive\ Rate \leq O\left(\frac{\varepsilon}{B^{\varepsilon}\log_{B}N}\right)$$

# Fixing Lookups (almost)

Query(64) → 8

The problem is that each node has multiple branches

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

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

| 37 | 58 | 93 | |
|----|----|----|--|

| 45 | 58 | 75 | 76 |
|----|----|----|----|
| 42 | 5  | 7  | 1  |

Idea: use filters to avoid searching them

|    |    | 64 |    |
|----|----|----|----|
| 38 | 39 |    | 94 |
|    |    | 8  |    |
| 1  | 2  |    | 4  |

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

$\Rightarrow$ Lookups in O(1) IOs

# *Really* Fixing Lookups in Size-Tiered B$^\varepsilon$-Trees

# *Really* Fixing Lookups in Size-Tiered Bᵋ-Trees

Querying all these filters is expensive



Root-to-leaf path

# *Really* Fixing Lookups in Size-Tiered Bᵉ-Trees

Querying all these filters is expensive

In practice, we see 15-40 filter lookups per point query



Multiple filters per node

Root-to-leaf path

# *Really* Fixing Lookups in Size-Tiered B$^{\varepsilon}$-Trees

Querying all these filters is expensive

In practice, we see 15-40 filter lookups per point query



We could hope to amortize against IO

Multiple filters per node

BUT...

Root-to-leaf path

# *Really* Fixing Lookups in Size-Tiered B$^\varepsilon$-Trees

Querying all these filters is expensive

In practice, we see 15-40 filter lookups per point query



Multiple filters per node

We could hope to amortize against IO

BUT…

High Memory/Hot Queries

No IO, performance limited by CPU

Root-to-leaf path

# *Really* Fixing Lookups in Size-Tiered Bᵉ-Trees

Querying all these filters is expensive

In practice, we see 15-40 filter lookups per point query



We could hope to amortize against IO

Multiple filters per node

BUT...

High Memory/Hot Queries — No IO, performance limited by CPU

Medium Memory — 1 IO per query,
CPU cost of filter lookups ⇒ more threads

Root-to-leaf path

# *Really* Fixing Lookups in Size-Tiered B$^\varepsilon$-Trees

Querying all these filters is expensive

In practice, we see 15-40 filter lookups per point query



We could hope to amortize against IO

Multiple filters per node

BUT…

High Memory/Hot Queries — No IO, performance limited by CPU

Medium Memory — 1 IO per query, CPU cost of filter lookups $\Rightarrow$ more threads

Low Memory — Filters paged out to storage, Lookup performance degrades

Root-to-leaf path

# Maplets

# Maplets

A maplet is a filter which can also store small values

# Maplets

A maplet is a filter which can also store small values

Is X in the set?



no          yes

Filter

# Maplets

A maplet is a filter which can also store small values

Is X in the set?

Is X in the set?



no          yes

no          yes, 4

Filter                    Maplet

# Maplets

A maplet is a filter which can also store small values

Is X in the set?

Is X in the set?



no          yes

no          yes, 4          yes, 3, 4 and 7

Filter

Maplet

SCHOOL OF COMPUTING
UNIVERSITY OF UTAH

# Maplets

A maplet is a filter which can also store small values

Is X in the set?

Is X in the set?

No false negatives, same false positive guarantee

no          yes

no          yes, 4          yes, 3, 4 and 7

Filter

Maplet

# Maplets

A maplet is a filter which can also store small values

Is X in the set?

Is X in the set?

No false negatives, same false positive guarantee

Same memory footprint as multiple filters

no          yes

no          yes, 4          yes, 3, 4 and 7

Filter                    Maplet

# Maplets

A maplet is a filter which can also store small values

Is X in the set?



Filter

no          yes

Is X in the set?



Maplet

no        yes, 4        yes, 3, 4 and 7

No false negatives, same false positive guarantee

Same memory footprint as multiple filters

Lookups same cost as 1 quotient filter:
2 cache line misses
1 IO

SCHOOL OF COMPUTING
UNIVERSITY OF UTAH

# Mapped $B^\varepsilon$-Trees

# Mapped $B^\varepsilon$-Trees

Replace individual filters with a single maplet

# Mapped $B^\varepsilon$-Trees

Replace individual filters with a single maplet

# Mapped $B^\varepsilon$-Trees

Replace individual filters with a single maplet

Use the values to store which buffers contain matching keys

# Mapped $B^\varepsilon$-Trees

Query(64)

Replace individual filters with a single maplet

Use the values to store which buffers contain matching keys



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

| 37 | 58 | 93 |
|----|----|----|

| 45 | 58 | 75 | 76 |
|----|----|----|----|
| 42 | 5  | 7  | 1  |

| 38 | 39 | 64 | 94 |
|----|----|----|----|
| 1  | 2  | 8  | 4  |

# Mapped $B^\varepsilon$-Trees

Query(64)

Replace individual filters with a single maplet

Use the values to store which buffers contain matching keys

| 37 | 58 | 93 |
|----|----|----|

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

| 45 | 58 | 75 | 76 |
|----|----|----|----|
| 42 | 5  | 7  | 1  |

| 38 | 39 | 64 | 94 |
|----|----|----|----|
| 1  | 2  | 8  | 4  |

1

# Mapped $B^{\varepsilon}$-Trees

Replace individual filters with a single maplet

Use the values to store which buffers contain matching keys

# Using Maplets to Manage Space

# Using Maplets to Manage Space

Size-tiering can lead to redundant data, wasting space

# Using Maplets to Manage Space

Size-tiering can lead to redundant data, wasting space

# Using Maplets to Manage Space

Size-tiering can lead to redundant data, wasting space

# Using Maplets to Manage Space

Size-tiering can lead to redundant data, wasting space

# Using Maplets to Manage Space

Size-tiering can lead to redundant data, wasting space

Compaction can recover disk space when there are many updates

# Using Maplets to Manage Space

Compaction saves little space when there is little redundant data



So we don't want to waste time compacting branches with few updates

SCHOOL OF COMPUTING
UNIVERSITY OF UTAH

# Using Maplets to Manage Space

Maplets can tell us how much redundant data there is

# Using Maplets to Manage Space

Maplets can tell us how much redundant data there is



41 → {0,1,2}
42 → {1}
43 → {0,2}

...

# Using Maplets to Manage Space

Maplets can tell us how much redundant data there is



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

| 41 | 79 | 85 | 98 |
|----|----|----|----|
| 42 | 5  | 7  | 1  |

| 41 | 43 | 64 | 91 |
|----|----|----|----|
| 1  | 2  | 8  | 4  |

| 37 | 58 | 93 |
|----|----|----|

$41 \rightarrow \{0,1,2\}$
$42 \rightarrow \{1\}$
$43 \rightarrow \{0,2\}$

...

Lots of multiple entries

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

| 45 | 58 | 75 | 76 |
|----|----|----|----|
| 42 | 5  | 7  | 1  |

| 38 | 39 | 64 | 94 |
|----|----|----|----|
| 1  | 2  | 8  | 4  |

| 37 | 58 | 93 |
|----|----|----|

# Using Maplets to Manage Space



Maplets can tell us how much redundant data there is

# Using Maplets to Manage Space

Maplets can tell us how much redundant data there is



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

| 41 | 79 | 85 | 98 |
|----|----|----|----|
| 42 | 5  | 7  | 1  |

| 41 | 43 | 64 | 91 |
|----|----|----|----|
| 1  | 2  | 8  | 4  |

37  58  93

41 → {0,1,2}
42 → {1}
43 → {0,2}

...

Lots of multiple entries

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

| 45 | 58 | 75 | 76 |
|----|----|----|----|
| 42 | 5  | 7  | 1  |

| 38 | 39 | 64 | 94 |
|----|----|----|----|
| 1  | 2  | 8  | 4  |

37  58  93

41 → {2}
42 → {2}
43 → {2}

...

Few multiple entries

Maplet

# SplinterDB Adaptive Space Reclamation

SplinterDB maintains a heap of trunk nodes, sorted by estimated amount of redundant data



Update estimate every time we rebuild maplet

# SplinterDB Adaptive Space Reclamation

SplinterDB maintains a heap of trunk nodes, sorted by estimated amount of redundant data

Whenever disk usage gets too high, SplinterDB initiates compaction on top node of the heap.

Goal: maximal gains, minimal pains

37  58  93

Update estimate every time we rebuild maplet

# SplinterDB A...



28 2Ghz cores

Intel Optane 905P

24B keys 100B values

25GiB RAM
80GiB dataset

Legend: SplinterDB+Maplets, SplinterDB, RocksDB

Chart axes: Updates/Second (1M 2M 3M) vs Space Efficiency (60% 80% 100%)

Data labels: ∞, 120 GiB, 112 GiB, 104 GiB, 96 GiB, 0 GiB

faster & smaller

100% uniform updates

SCHOOL OF COMPUTING
UNIVERSITY OF UTAH

# Flush-Then-Compact

# Flush-Then-Compact

Sequential insertions into a B-tree

71
—
2

24 37 86 90

48 58 83 84

61 68 83 84

SCHOOL OF COMPUTING
UNIVERSITY OF UTAH

# 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



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

# Flush-Then-Compact

Sequential Insertions into a B-tree



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

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

# Flush-Then-Compact

Sequential Insertions into a $B^{\varepsilon}$-tree

| 69 | 71 | 72 | 73 |
|----|----|----|----|
| 9  | 2  | 50 | 14 |

201

# Flush-Then-Compact

Sequential Insertions into a
B$^\varepsilon$-tree



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

# Flush-Then-Compact

Sequential Insertions into a
B$^\varepsilon$-tree



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

# Flush-Then-Compact

Sequential Insertions into a
B$^\varepsilon$-tree

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

# Flush-Then-Compact

Sequential Insertions into a
$B^{\varepsilon}$-tree



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

# Flush-Then-Compact

Sequential Insertions into a
$B^{\varepsilon}$-tree



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

# Flush-Then-Compact

Sequential Insertions into a
B$^{\varepsilon}$-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

| 37 | 86 |
|----|----|

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

| 58 | 83 |
|----|----|

| 90 | 92 |
|----|----|

| 41 | 48 |
|----|----|

| 67 | 75 |
|----|----|

| 84 | 85 |
|----|----|

| 59 | 60 | 61 | 65 |
|----|----|----|----|
| 5  | 40 | 29 | 11 |

| 69 | 71 | 72 | 73 |
|----|----|----|----|
| 9  | 2  | 50 | 14 |

| 79 | 80 | 81 | 82 |
|----|----|----|----|
| 99 | 6  | 77 | 44 |

# Flush-Then-Compact

Want:
Cheap sequential insertions



After merging and flushing another flush will be triggered

# Flush-Then-Compact

Want:
Cheap sequential insertions



After merging and flushing another flush will be triggered

# Flush-Then-Compact

Want:
Cheap sequential insertions



After merging and flushing another flush will be triggered

# Flush-Then-Compact

Want:
Cheap sequential insertions



Any data already present will get merged again

After merging and flushing another flush will be triggered

# Flush-Then-Compact

Want:
Cheap sequential insertions

| 37 | 58 | 93 |
|----|----|----|

| 58 | 83 | |

| 37 | 45 | 52 |

| 93 | 99 |

| 58 | 59 | 60 | 61 | 62 | 63 | 64 | 65 | 66 | 67 | 68 | 69 |
|----|----|----|----|----|----|----|----|----|----|----|----|
| 1  | 2  | 8  | 4  | 42 | 5  | 7  | 1  | 2  | 5  | 11 | 1  |

Can still end up merging on each level

| 72 | 73 | 74 | 75 |
|----|----|----|----|
| 1  | 2  | 8  | 4  |

Any data already present will get merged again

| 58 | 78 | |

After merging and flushing another flush will be triggered

# Flush-Then-Compact

Want:
Cheap sequential insertions

Idea: Flush-then-compact

# Flush-Then-Compact

Want:
Cheap sequential insertions

Idea: Flush-then-compact

# Flush-Then-Compact

Want:
Cheap sequential insertions

Idea: Flush-then-compact



First flush references to the branches, but do not compact

# Flush-Then-Compact

Want:
Cheap sequential insertions

Idea: Flush-then-compact



First flush references to the branches, but do not compact

Use metadata to mask out data

# Flush-Then-Compact

The parent only sees the unflushed data

Want:
Cheap sequential insertions

Idea: Flush-then-compact



First flush references to the branches, but do not compact

Use metadata to mask out data

# Flush-Then-Compact

Want:
Cheap sequential insertions

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



| 41 | 67 | 68 | 69 |
|----|----|----|----|
| 2  | 5  | 11 | 1  |

| 62 | 63 | 64 | 65 |
|----|----|----|----|
| 42 | 5  | 7  | 1  |

| 58 | 59 | 60 | 94 |
|----|----|----|----|
| 1  | 2  | 8  | 4  |

37  58  93

37  45  52

58  83

58  78

**SCHOOL OF COMPUTING**
UNIVERSITY OF UTAH

218

# Flush-Then-Compact

Want:
Cheap sequential insertions

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:
Cheap sequential insertions

Idea: Flush-then-compact

Then can flush again



First flush references to the branches, but do not compact

Use metadata to mask out data

220

# Flush-Then-Compact

Want:
Cheap sequential insertions

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:
Cheap sequential insertions

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

| 37 | 58 | 93 |
|----|----|----|

| 58 | 83 |
|----|----|

| 37 | 45 | 52 |
|----|----|----|

| 58 | 78 |
|----|----|

First flush references to the branches, but do not compact

Use metadata to mask out data

| 41 | 67 | 68 | 69 |
|----|----|----|----|
| 2 | 5 | 11 | 1 |

| 62 | 63 | 64 | 65 |
|----|----|----|----|
| 42 | 5 | 7 | 1 |

| 58 | 59 | 60 | 94 |
|----|----|----|----|
| 1 | 2 | 8 | 4 |

223

# Flush-Then-Compact

No work on immediately flushed data

Sequential insertions have write amp ~1

| 41 | 67 | 68 | 69 |
|----|----|----|----|
| 2 | 5 | 11 | 1 |

| 62 | 63 | 64 | 65 |
|----|----|----|----|
| 42 | 5 | 7 | 1 |

| 58 | 59 | 60 | 94 |
|----|----|----|----|
| 1 | 2 | 8 | 4 |

| 37 | 45 | 52 |
|----|----|----|

| 58 | 83 |
|----|----|

| 99 |
|----|

First flush references to the branches, but do not compact

| 58 | 78 |
|----|----|

Use metadata to mask out data

# Flush-Then-Compact

No work on immediately flushed data

Sequential insertions have write 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

| 58 | 78 |
|----|----|

Use metadata to mask out data

SCHOOL OF COMPUTING
UNIVERSITY OF UTAH

# Flush-Then-Compact

No work on immediately flushed data

Sequential insertions have write amp ~1

Break a serial chain of compactions into parallel

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  |

| 58 | 78 |
|----|----|

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

Sequential insertions have write amp ~1
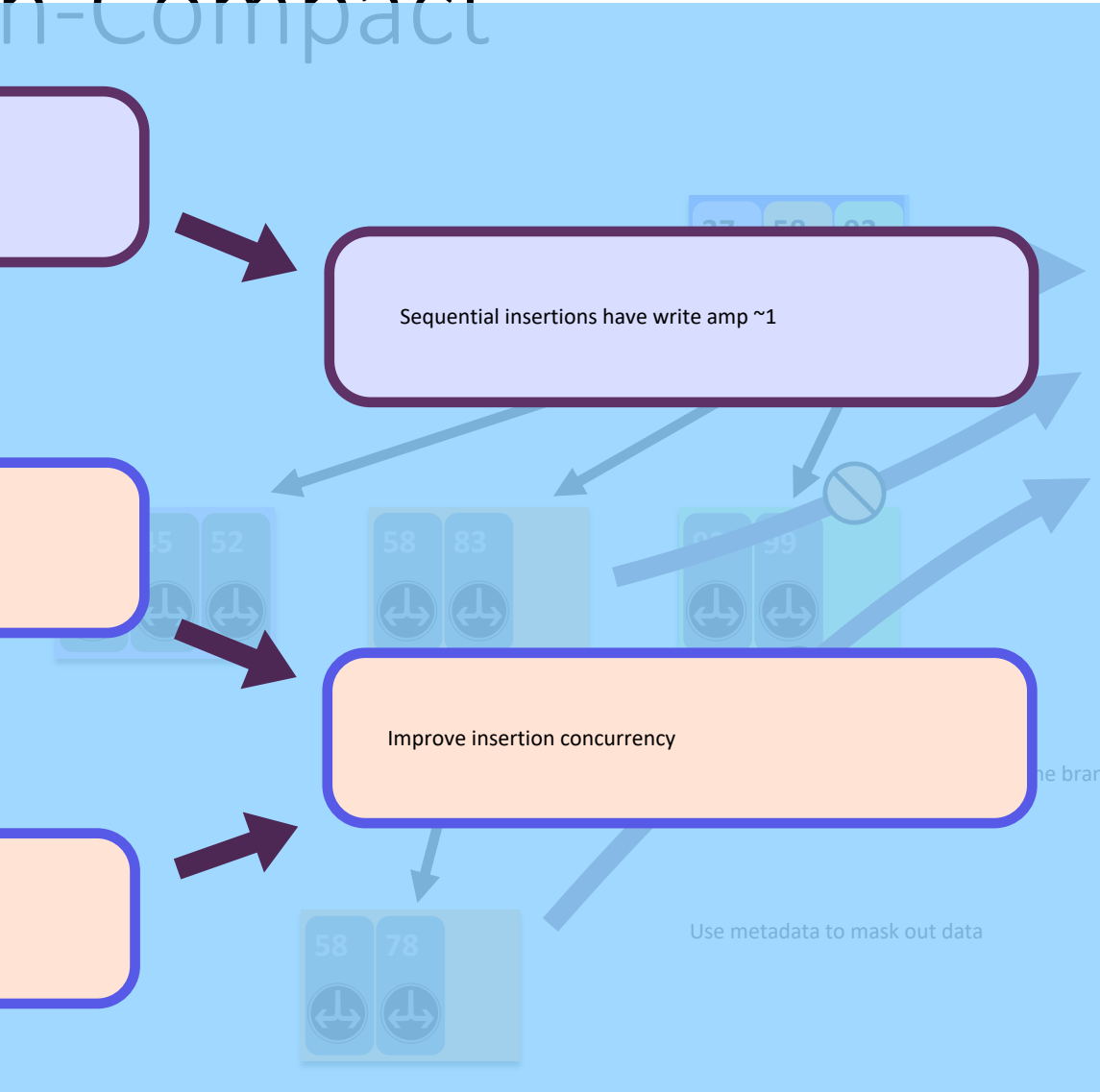
Break a serial chain of compactions into parallel

Improve insertion concurrency

Concurrent compactions in trunk nodes

the branches, but do not compact

Use metadata to mask out data

| 41 | 67 | 68 | 69 |
|----|----|----|----|
| 2  | 5  | 11 | 1  |

| 62 | 63 | 64 | 65 |
|----|----|----|----|
| 42 | 5  | 7  | 1  |

| 58 | 59 | 60 | 94 |
|----|----|----|----|
| 1  | 2  | 8  | 4  |

# Flush-Then-Compact

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



SplinterDB

Throughput
(Insertions/Sec)

1000

750

500

250

0

866

799

676

521

430

193

185

171

152

144

Percentage Sequential

0    50    90    99    100

Higher is Better

X-axis not to scale

228

# 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



**SplinterDB**

Throughput (Insertions/Sec)

1000

750

500

250

0

430
521
676
799
866

144
152
171
185
193

Percentage Sequential
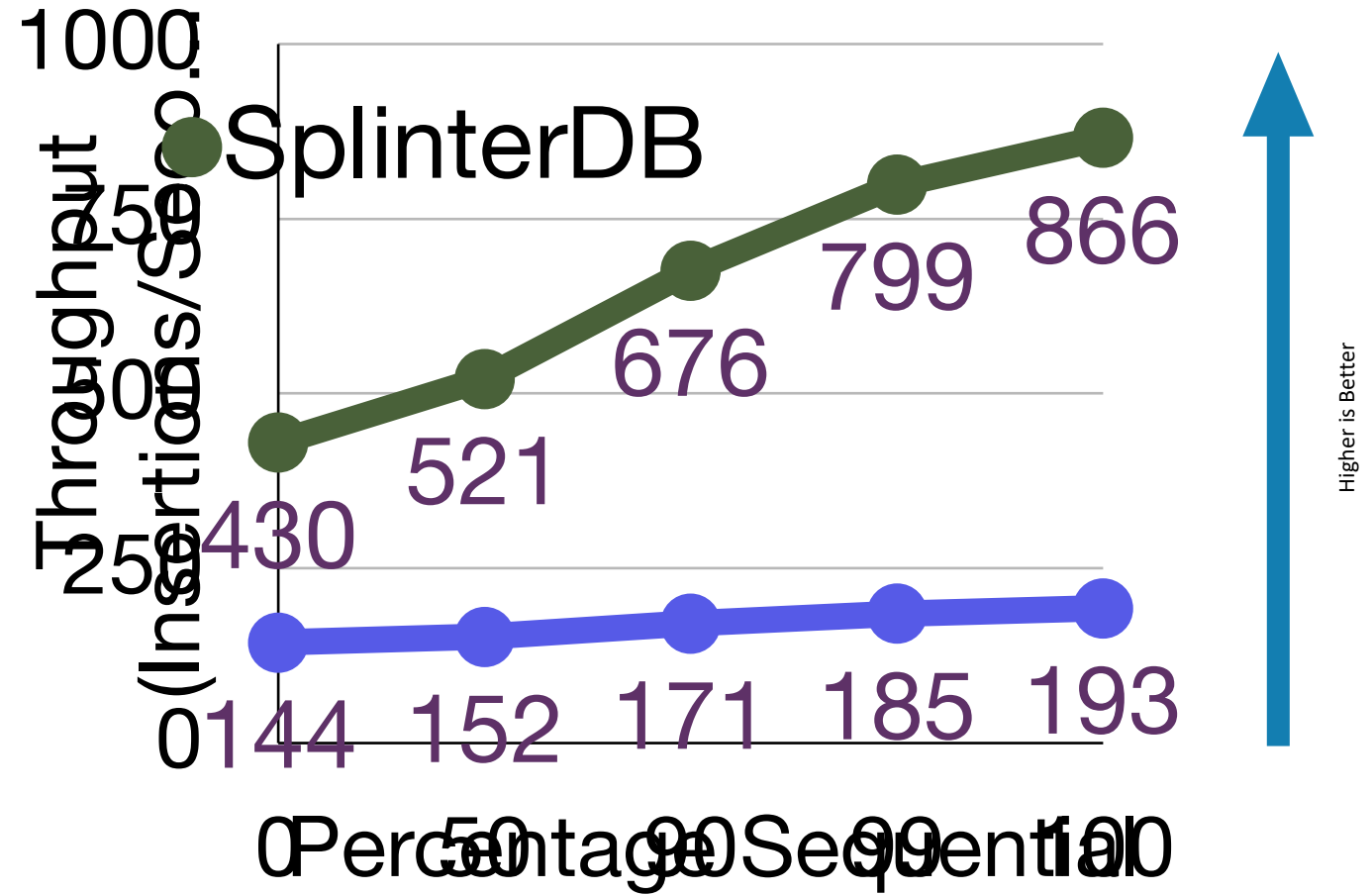
0    50    90    99    100

Higher is Better

X-axis not to scale

229

# 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
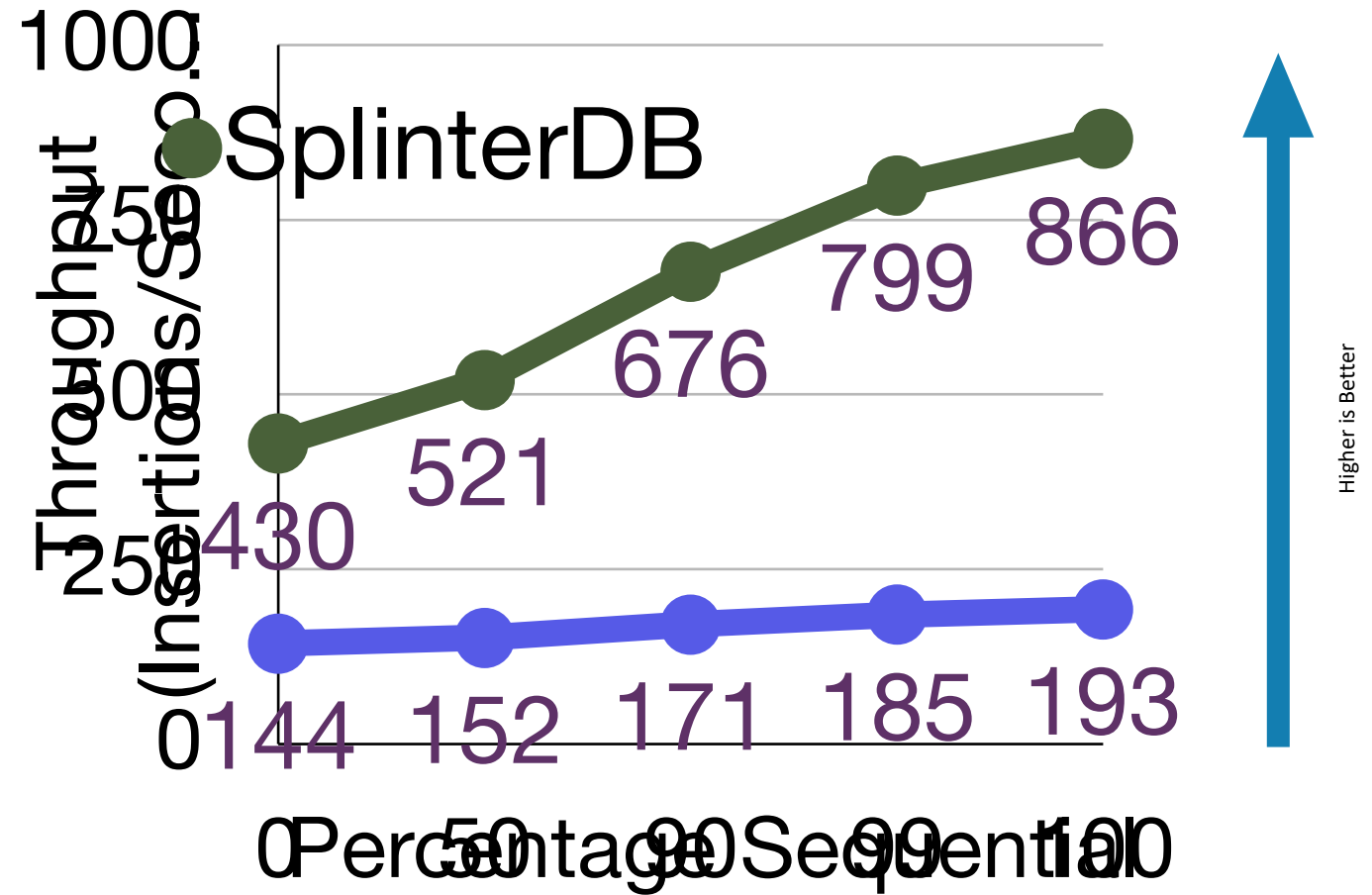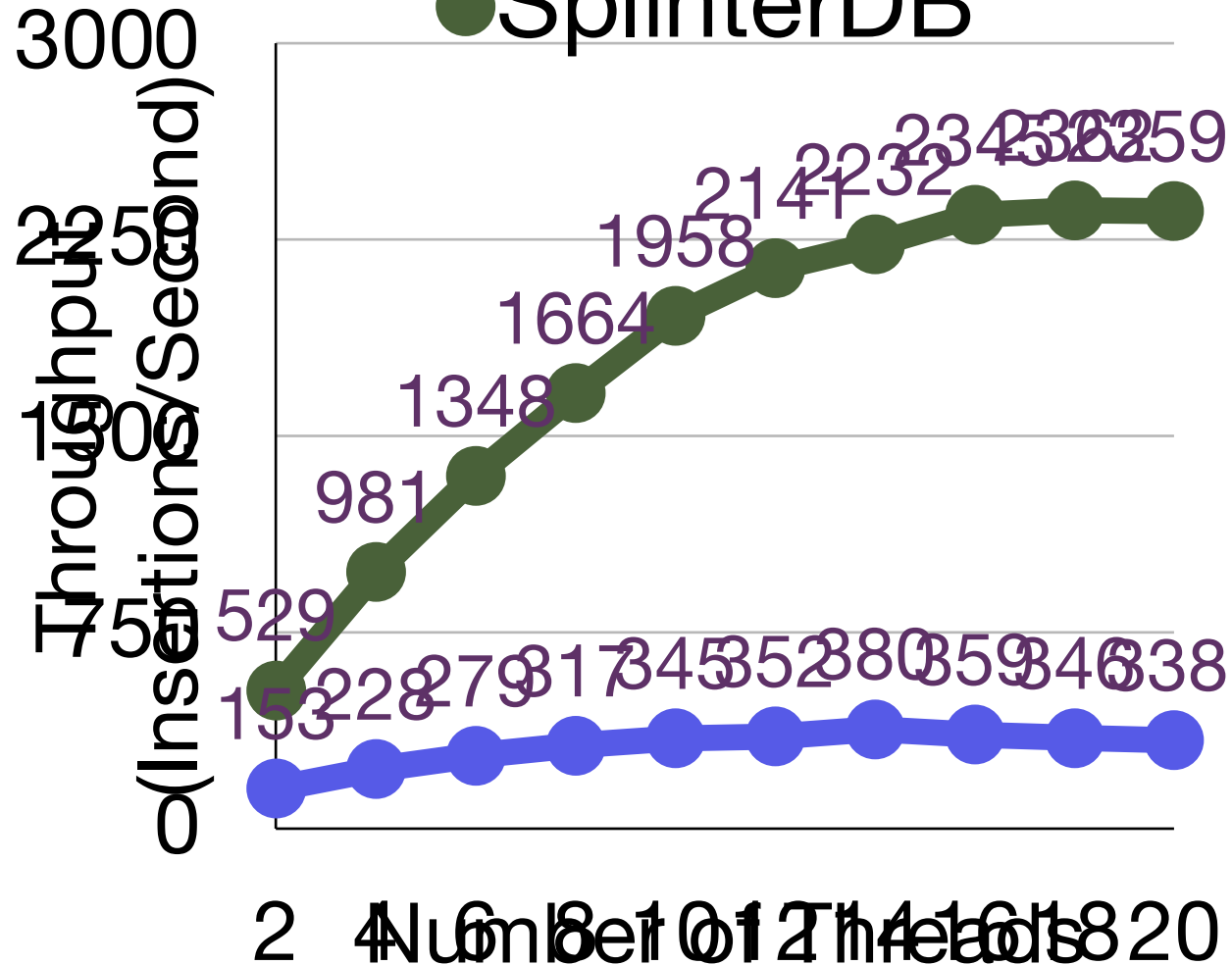
RocksDB improves, but at a much lower rate



SplinterDB

Throughput (Insertions/Sec)

1000

750

500

250

0

430
521
676
799
866

144
152
171
185
193

Percentage Sequential
0    50    90    99    100

Higher is Better

X-axis not to scale

# Flush-then-Compact



SplinterDB

Throughput (Insertions/Second)

3000

2250

1500

750

0

2 4 6 8 10 12 14 16 18 20

Number of Threads

Insertions in SplinterDB scale well

Higher is Better

SplinterDB values: 529, 981, 1348, 1664, 1958, 2141, 2232, 2342, 2363, 2359

Other values: 153, 228, 279, 317, 345, 352, 380, 359, 346, 338

231

# Flush-then-Compact



SplinterDB

Throughput(Insertions/Second)

3000

2250

1500

750

0

2 4 6 8 10 12 14 16 18 20
Number of Threads

529
981
1348
1664
1958
2147
2232
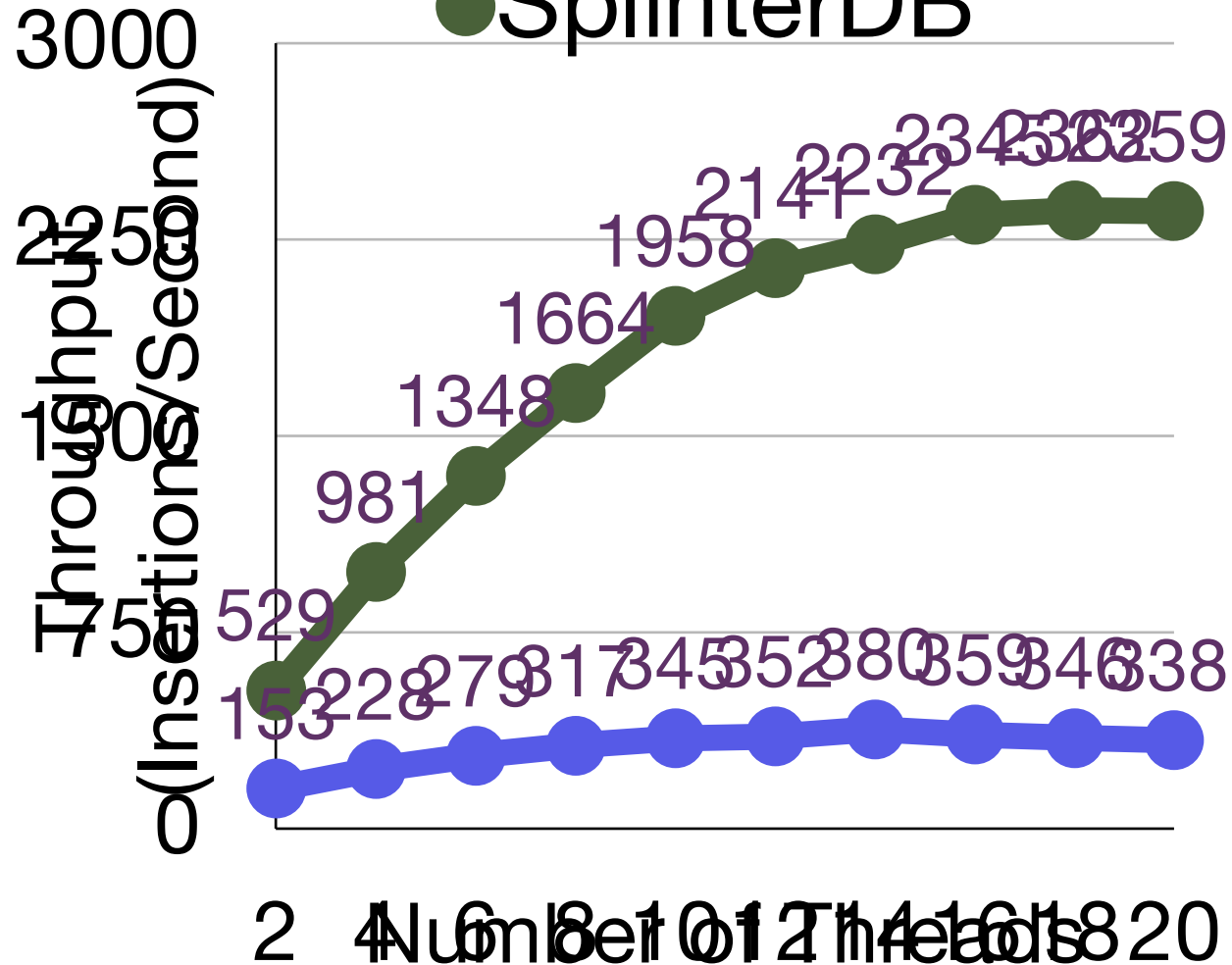2342
2363
2359

153
228
279
317
345
352
380
359
346
38

Higher is Better

Insertions in SplinterDB scale well

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

# Flush-then-Compact



**SplinterDB**

Throughput(Insertions/Second)

3000

2250

1500

750

0

SplinterDB data points: 529, 981, 1348, 1664, 1958, 2147, 2232, 2324, 2359, 2332, 2359

Lower line data points: 153, 228, 279, 317, 345, 352, 380, 359, 346, 38

Number of Threads

2  4  6  8  10  12  14  16  18  20
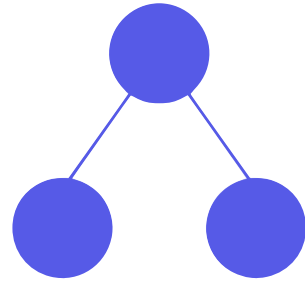
Higher is Better

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

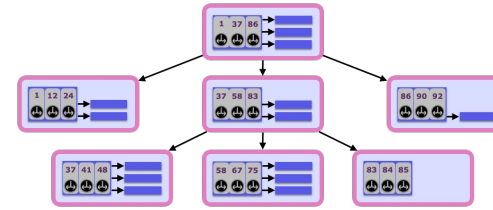Model the problem:
external memory dictionary

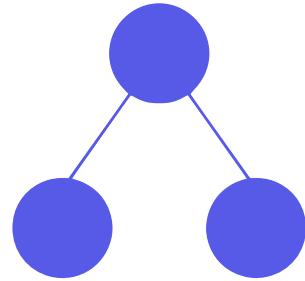Mapped $B^{\varepsilon}$-tree

Theory

Systems



vSAN needed a new way of storing metadata
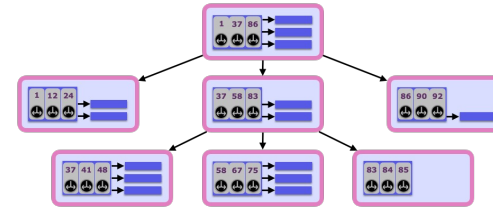
SplinterDB

# Conclusion

Model the problem:
external memory dictionary

Mapped $B^\varepsilon$-tree



Theory

Systems

SplinterDB is in vSAN 8.0

vSAN needed a new way of storing metadata

Open-source at
https://github.com/vmware/splinterdb

SplinterDB