# Successor Queries in Optimal External-Memory Dictionaries
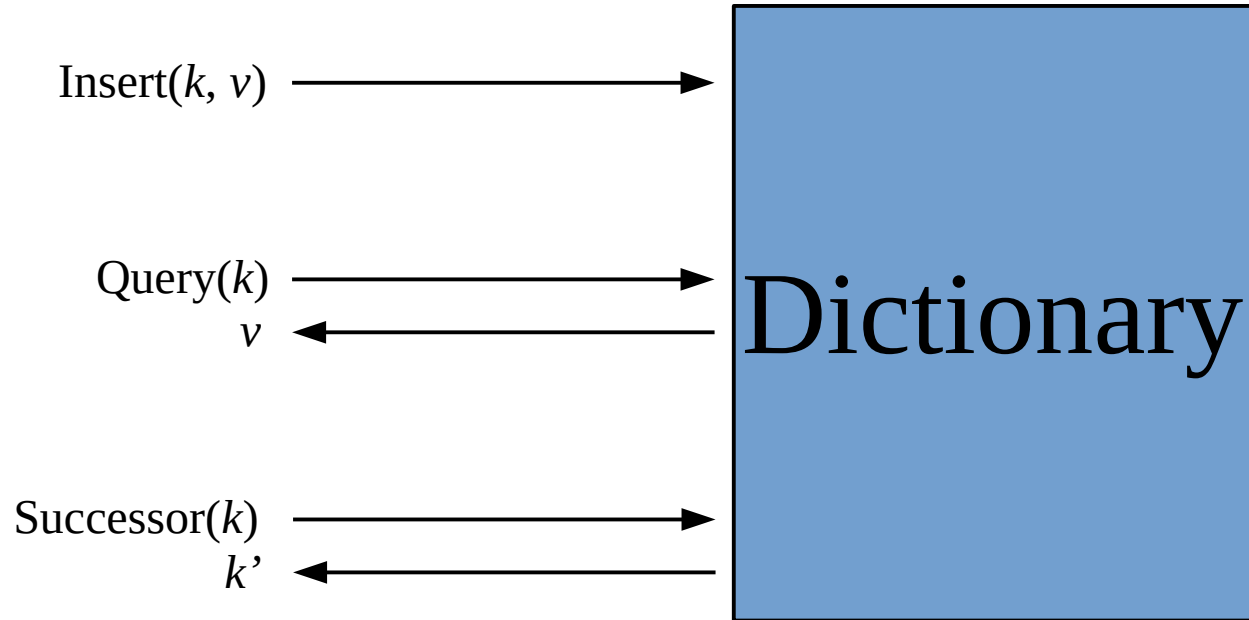
Rob Johnson

VMware Research Group

# What is a dictionary?

Insert($k$, $v$) ⟶ [Dictionary]

Query($k$) ⟶
$v$ ⟵

Successor($k$) ⟶
$k'$ ⟵

Scan($k$, $L$) = return the $L$ successors of $k$

# Dictionary Performance Trade-Offs
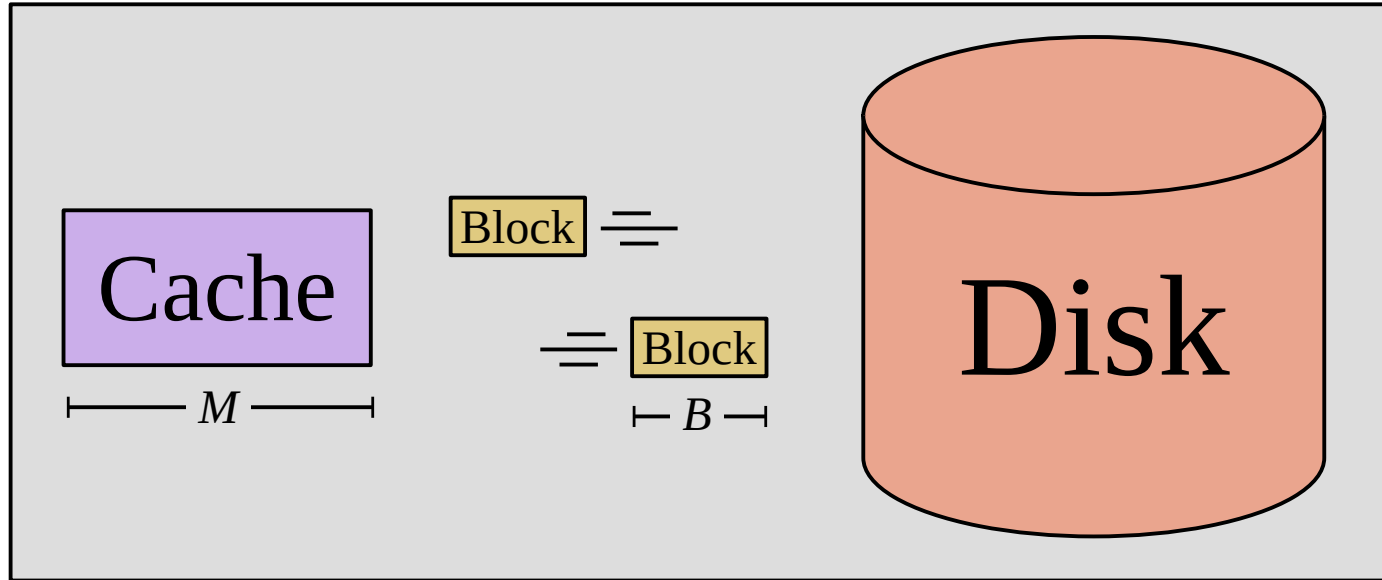
Inserts ⟷ Point Queries & Scans

# The research program

- Make insertions as fast as possible

- While preserving fast point queries

- And "reasonable" successor queries.
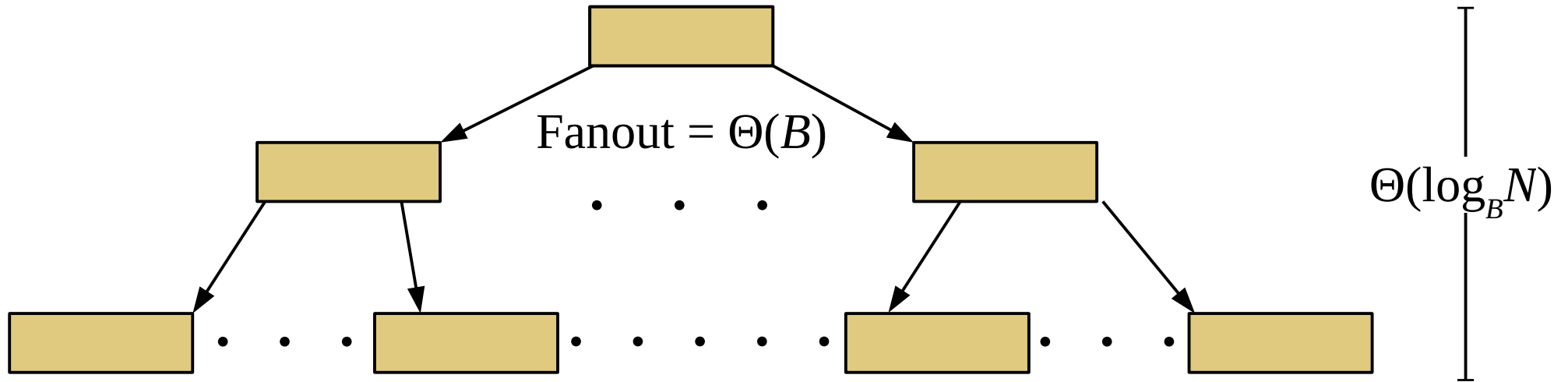
# The Disk Access Machine (DAM) Model

[Aggarwal & Vitter '88]



Algorithm design goal: minimize number of block transfers

# B-trees were long thought to be optimal
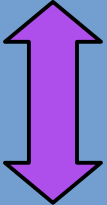
[Bayer & McCreight '70]

Fanout = $\Theta(B)$

$\Theta(\log_B N)$

Insertions
Queries    } $\Theta(\log_B N)$ I/Os        Scans: $\Theta(L/B + \log_B N)$ I/Os
Successors

# The Brodal-Fagerberg bounds [Brodal & Fagerberg '03]

# The Brodal-Fagerberg bounds [Brodal & Fagerberg '03]

Atomic key comparison-based bounds
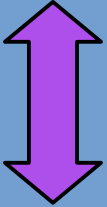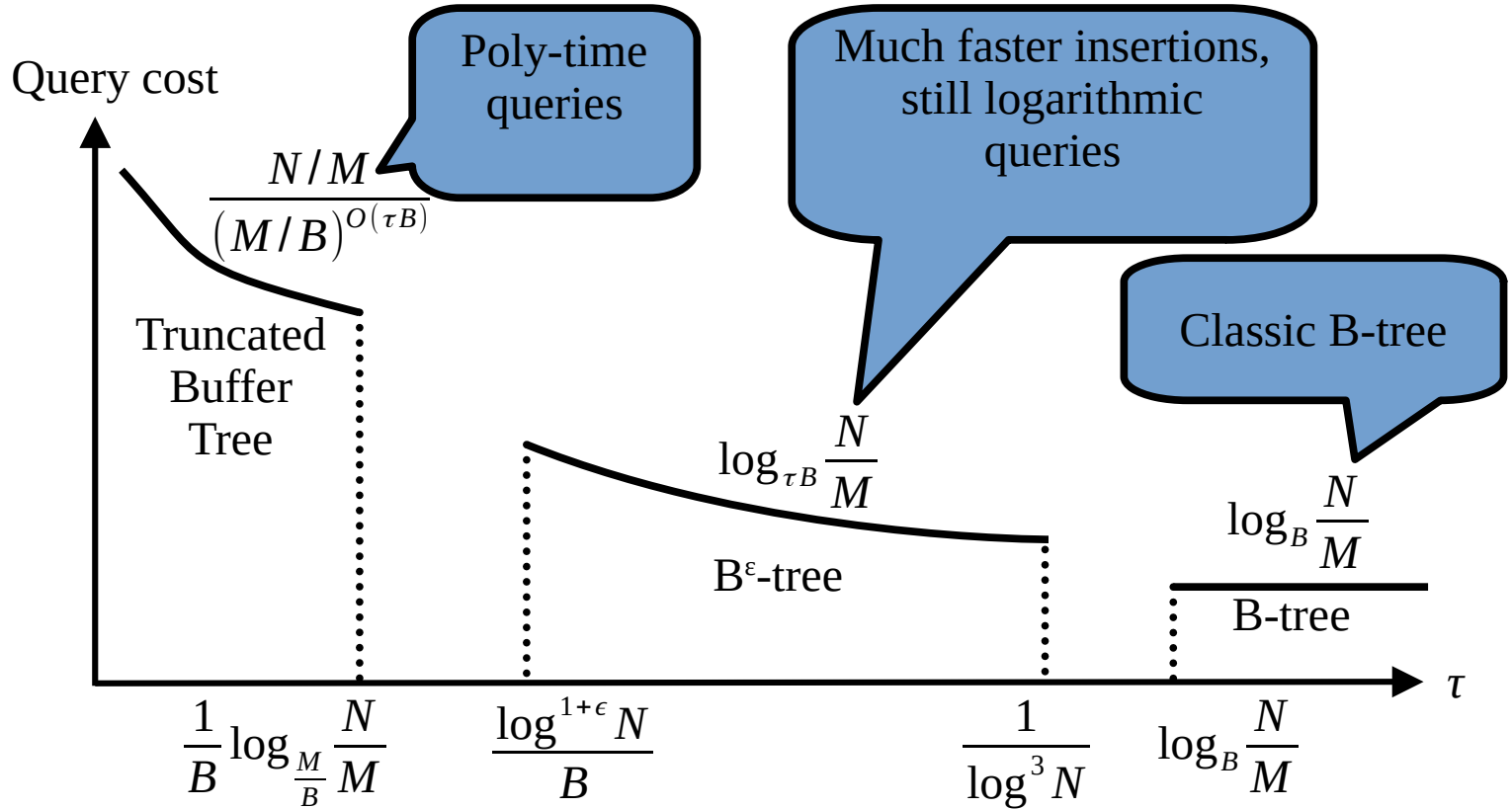
No hashing

Inserts: $\tau$

Queries: $f(\tau)$

Query cost

$$\frac{N/M}{(M/B)^{O(\tau B)}}$$

Poly-time queries

Much faster insertions, still logarithmic queries

Classic B-tree

Truncated Buffer Tree

$$\log_{\tau B} \frac{N}{M}$$

$B^{\epsilon}$-tree

$$\log_B \frac{N}{M}$$

B-tree

$\tau$

$$\frac{1}{B} \log_{\frac{M}{B}} \frac{N}{M}$$

$$\frac{\log^{1+\epsilon} N}{B}$$

$$\frac{1}{\log^3 N}$$

$$\log_B \frac{N}{M}$$

*Some conditions apply

# B$^\varepsilon$-tree asymptotics

| Fanout | Insertions | Point queries | Scans |
|---|---|---|---|
| $F$ | $\dfrac{F \log_F N}{B}$ | $\log_F N$ | $L/B + \log_F N$ |
| $B^\varepsilon$ | $\dfrac{\log_B N}{\epsilon B^{1-\epsilon}}$ | $\dfrac{\log_B N}{\epsilon}$ | $L/B + \dfrac{\log_B N}{\epsilon}$ |
| $\sqrt{B}$ | $\dfrac{\log_B N}{\sqrt{B}}$ | $\log_B N$ | $L/B + \log_B N$ |
| $\dfrac{\tau B}{\log N}$ | $\tau$ | $\log_{\tau B} N$ | $L/B + \log_{\tau B} N$ |

Assuming $\tau \geq \dfrac{\log^{1+\epsilon} N}{B}$

# The Iacono-Pătrașcu bounds
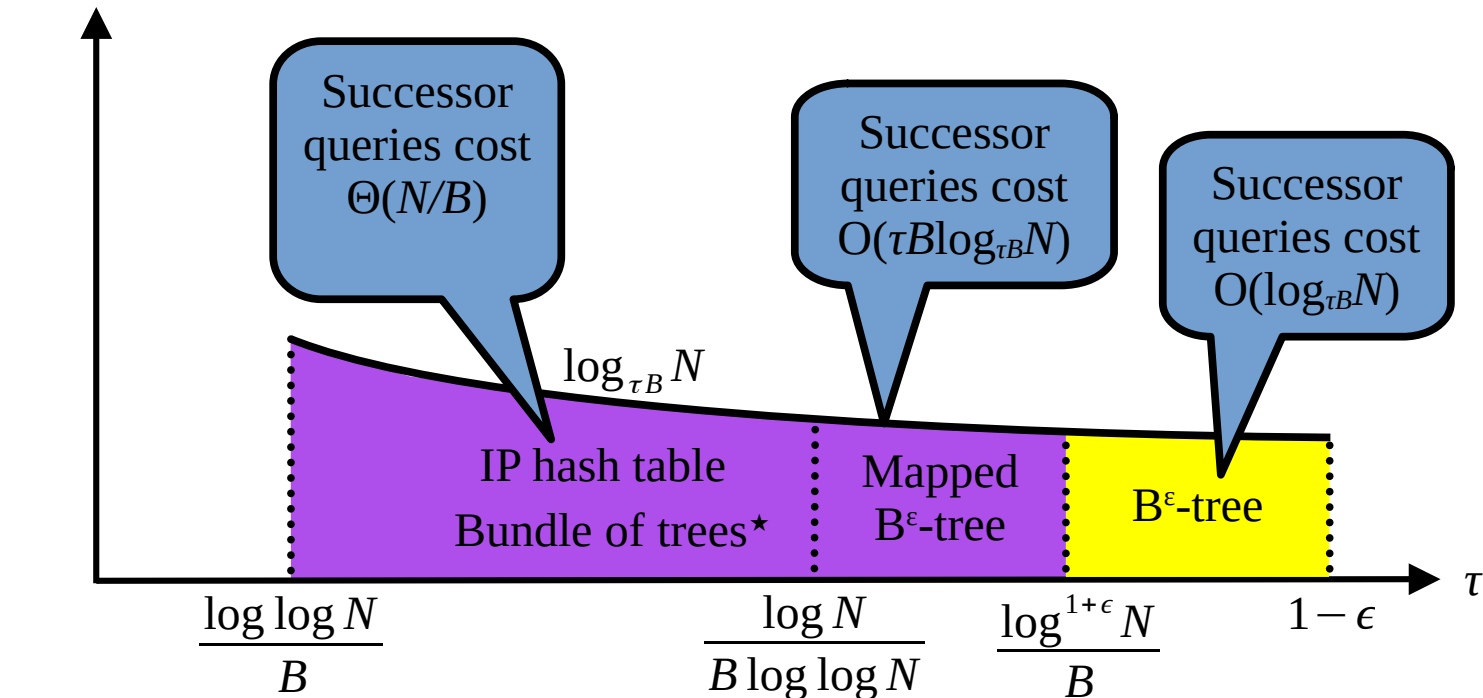
[Iacono & Pătrașcu '11]

Non-atomic key model

Hashing allowed

Inserts: $\tau$

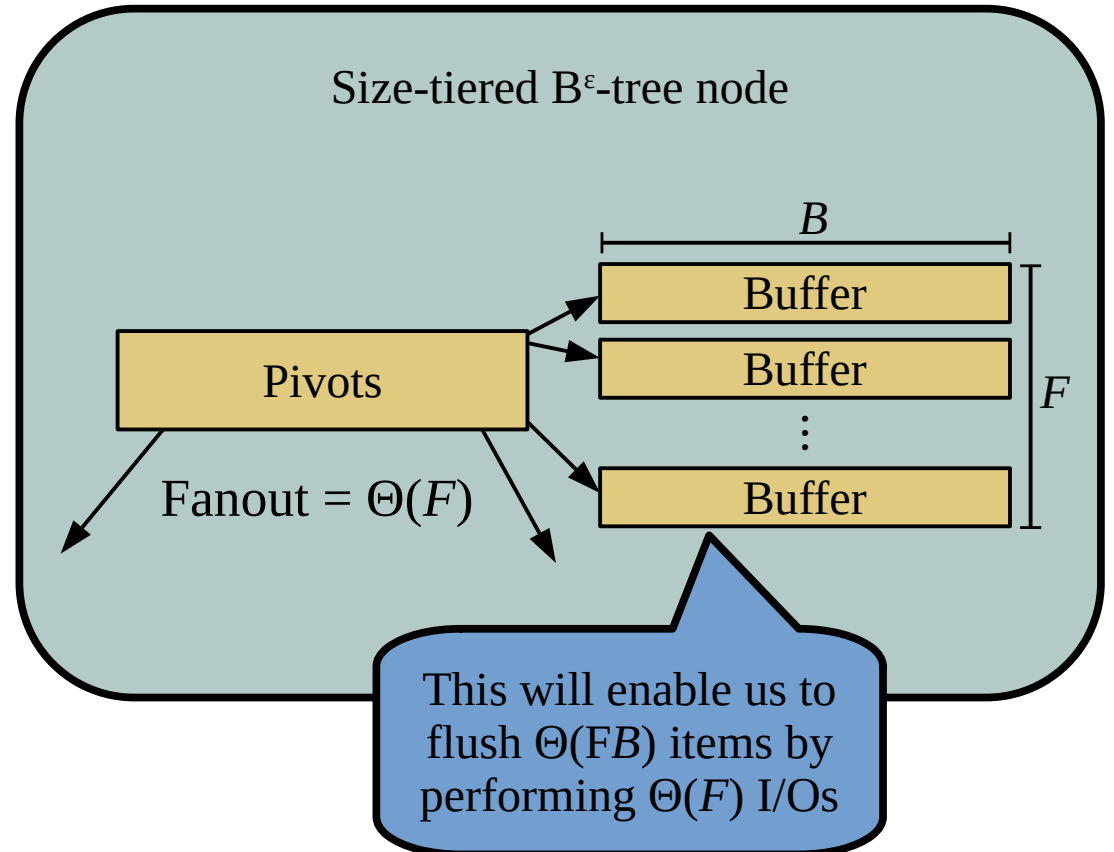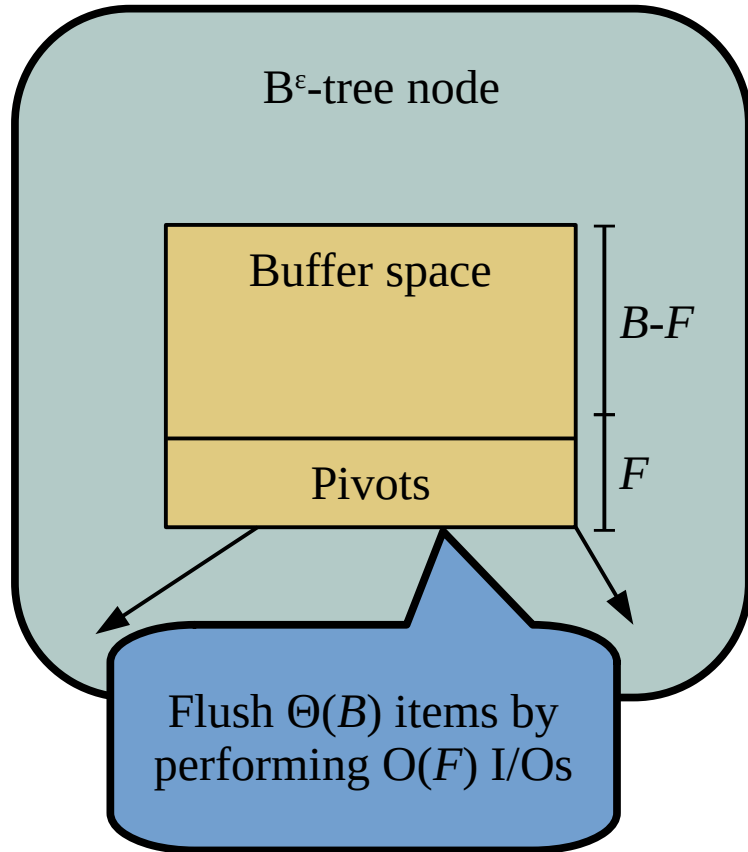Queries: $\log_{\tau B} N$

Query cost

Successor queries cost $\Theta(N/B)$

Successor queries cost $O(\tau B \log_{\tau B} N)$

Successor queries cost $O(\log_{\tau B} N)$

$\log_{\tau B} N$

IP hash table
Bundle of trees★

Mapped $B^\epsilon$-tree

$B^\epsilon$-tree

$\tau$

$\dfrac{\log \log N}{B}$

$\dfrac{\log N}{B \log \log N}$

$\dfrac{\log^{1+\epsilon} N}{B}$

$1 - \epsilon$

★[Conway, Farach-Colton, Shilane '18]

*Some conditions apply

# Size-tiered $B^\varepsilon$-trees
## Faster inserts, slower queries
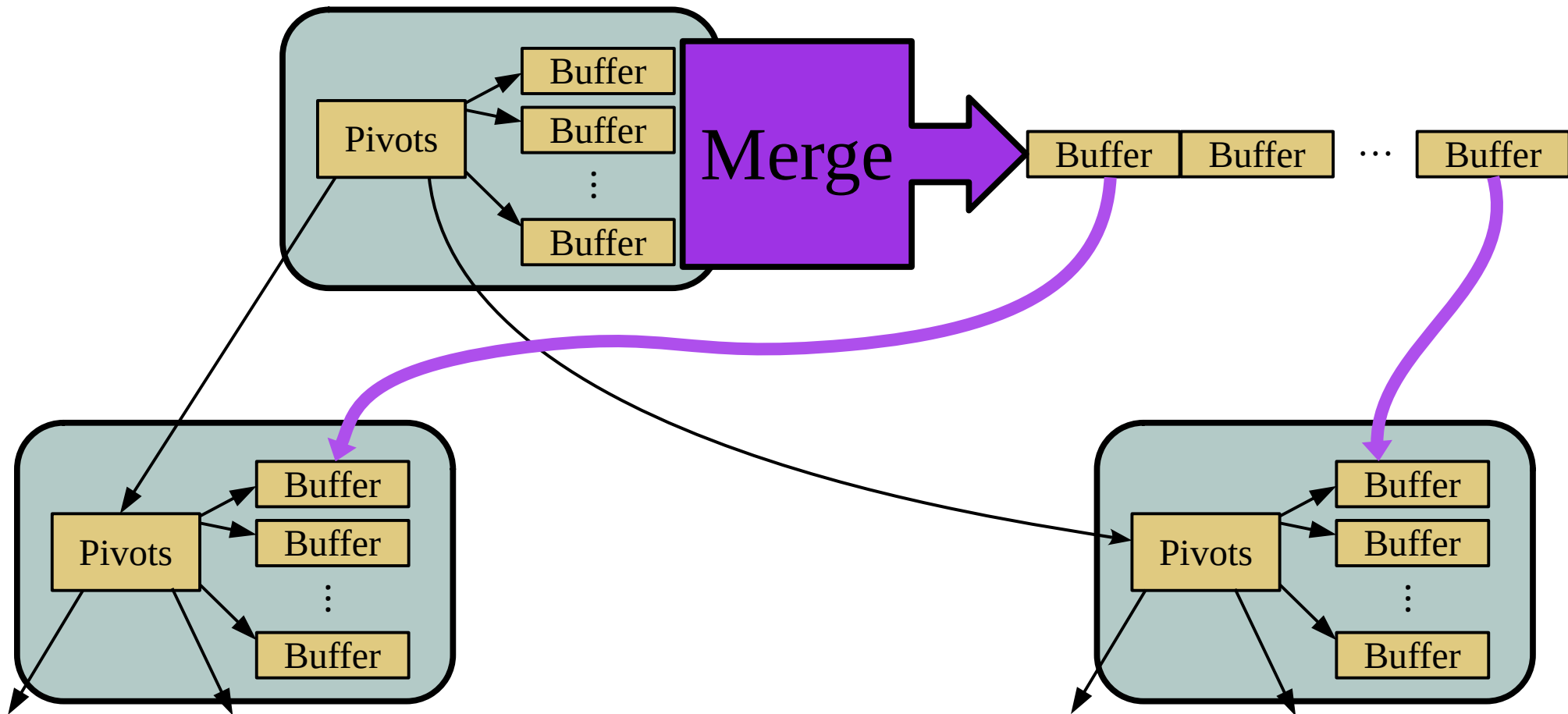
[Jagadish, Narayan, Seshadri, Sudarshan, Kanneganti '97]

# Size-tiered Bᵋ-tree nodes

# Flushes in size-tiered Bᵋ-trees

# Analysis: Flushes in size-tiered B$^\varepsilon$-trees

# Analysis: Flushes in size-tiered B$^\varepsilon$-trees

# Analysis: Insertions in size-tiered B$^\varepsilon$-trees



Flush $\Theta(FB)$ items using $\Theta(F)$ I/Os

Fanout $\Theta(F)$

Insertions: $\Theta\left(\dfrac{\log_F N}{B}\right)$

$\Theta\left(\log_F N\right)$

Pivots | Buffer Buffer ⋮ Buffer

# Analysis: I/O costs in size-tiered B$^\varepsilon$-trees

Flush $\Theta(FB)$ items using $\Theta(F)$ I/Os

Fanout $\Theta(F)$

Insertions: $\Theta\left(\dfrac{\log_F N}{B}\right)$

$\Theta\left(\log_F N\right)$

Point queries: $\Theta\left(F \log_F N\right)$

Scans: $\Theta\left(L/B + F \log_F N\right)$

Pivots
Buffer
Buffer
Buffer

# Maplets and mapped B$^\varepsilon$-trees

Fixing queries

# Maplets

- Maplets extend filters from **sets** to **maps**
  - maplet_query($k$) → { $v_1, v_2, …, v_\ell$ }

- Maplets save space by allowing false positives
  - False positives are extra values in a query result
  - False-positive rate = E[# of extra values]

- Basic implementation:
  - Store a ordered linear-probing hash table of *(h(k), v)* pairs
  - Compress table using **quotienting** [Knuth 1973]

maplet_query("Knuth")

$h$("Knuth") = 28

Result: {7, 9}

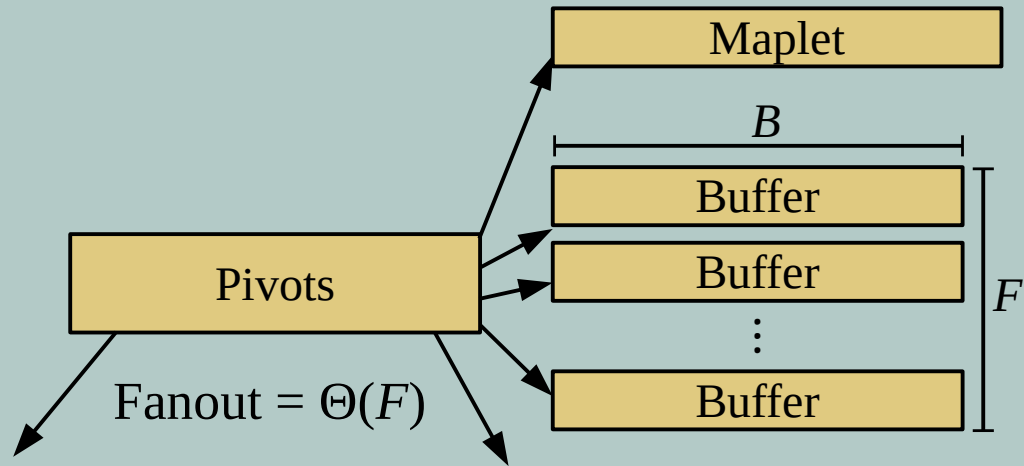Queries need O(1) I/Os w.h.p.

| $i$ | $h(k)$ | $v$ |
|---|---|---|
| 0 | 05 | 5 |
| 1 | - | - |
| 2 | 21 | 3 |
| 3 | 28 | 7 |
| 4 | 28 | 9 |
| 5 | - | - |
| 6 | 67 | 2 |

# Mapped B$^\varepsilon$-tree nodes

# Queries in mapped B$^\varepsilon$-trees

# Queries in mapped B$^\varepsilon$-trees



Mapped buffered B-tree node

Maplet

$B$

Buffer

Buffer

$F$

⋮

Buffer

Pivots

Fanout = $\Theta(F)$
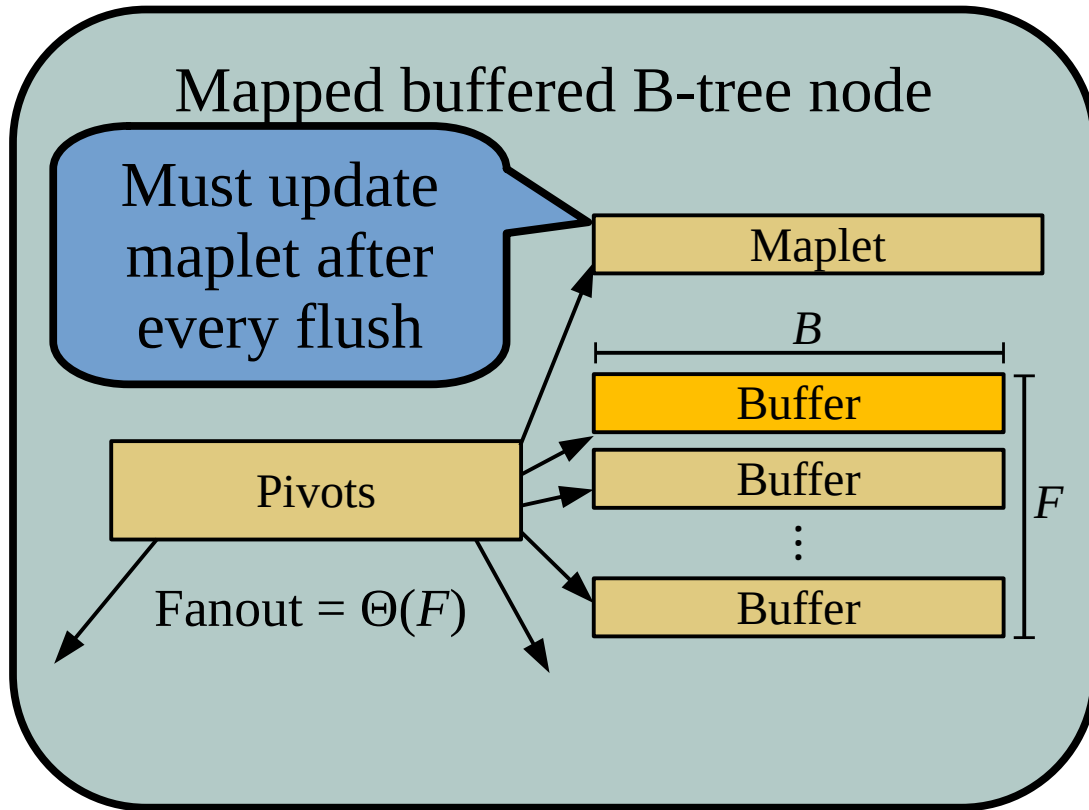
Scans: $\Theta\left(L/B + F\log_F N\right)$

Maplet: $k \rightarrow$ { buffers containing $k$ }

Queries access O(1) blocks in each node

Point queries: $\Theta\left(\log_F N\right)$

# Flushes in mapped B$^\varepsilon$-trees



Mapped buffered B-tree node

Must update maplet after every flush

Maplet

$B$

Buffer

Buffer

Buffer

$\vdots$

Buffer

$F$

Pivots

Fanout = $\Theta(F)$
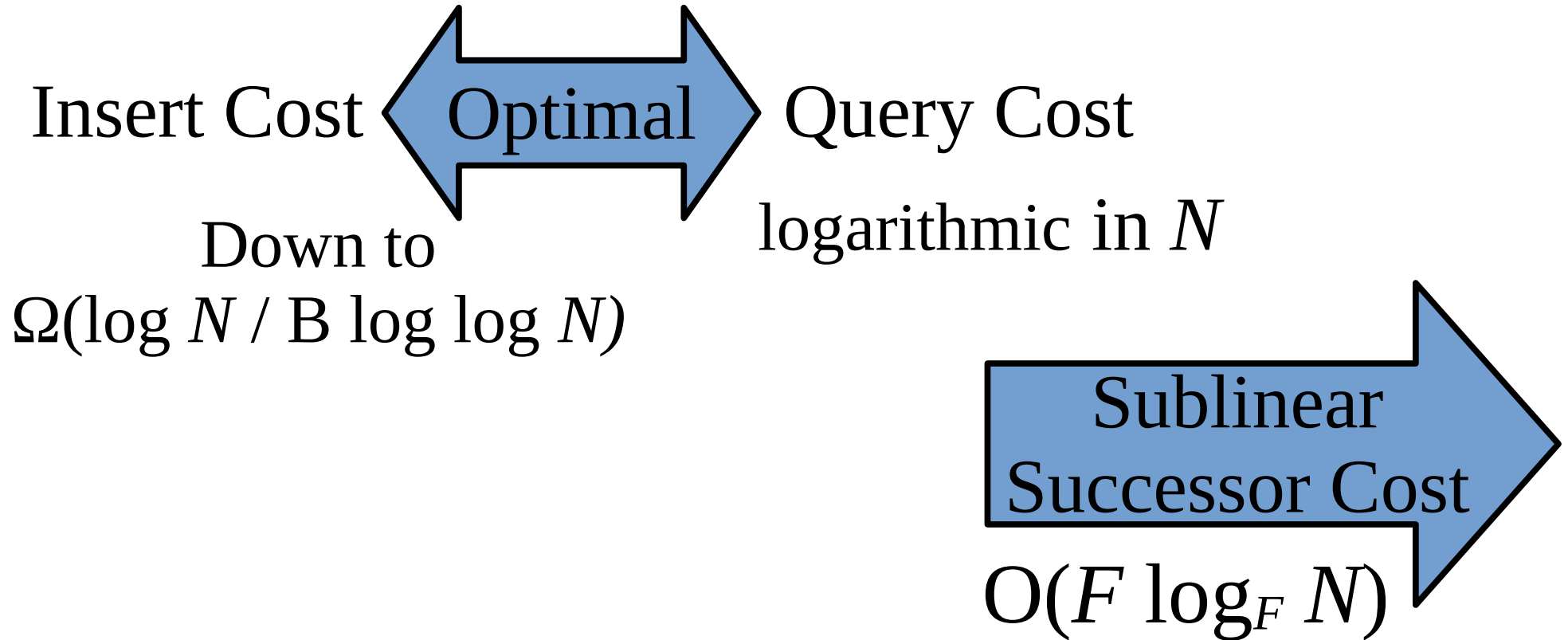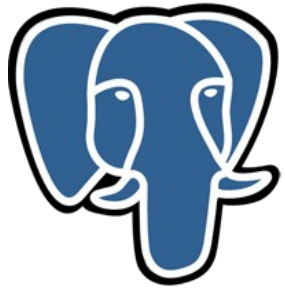
Maplet: $k \rightarrow \{$ buffers containing $k$ $\}$

Maplet maintenance not a bottleneck to optimality or in practice

**Theorem:** Mapped B$^\varepsilon$-tree meets IP lower bound when $F=\Omega(\log N / \log \log N)$.

# Summary of theoretical results

Insert Cost $\longleftrightarrow$ Optimal $\longleftrightarrow$ Query Cost

Down to logarithmic in $N$

$\Omega(\log N \,/\, B \log \log N)$

Sublinear Successor Cost $\longrightarrow$

$O(F \log_F N)$

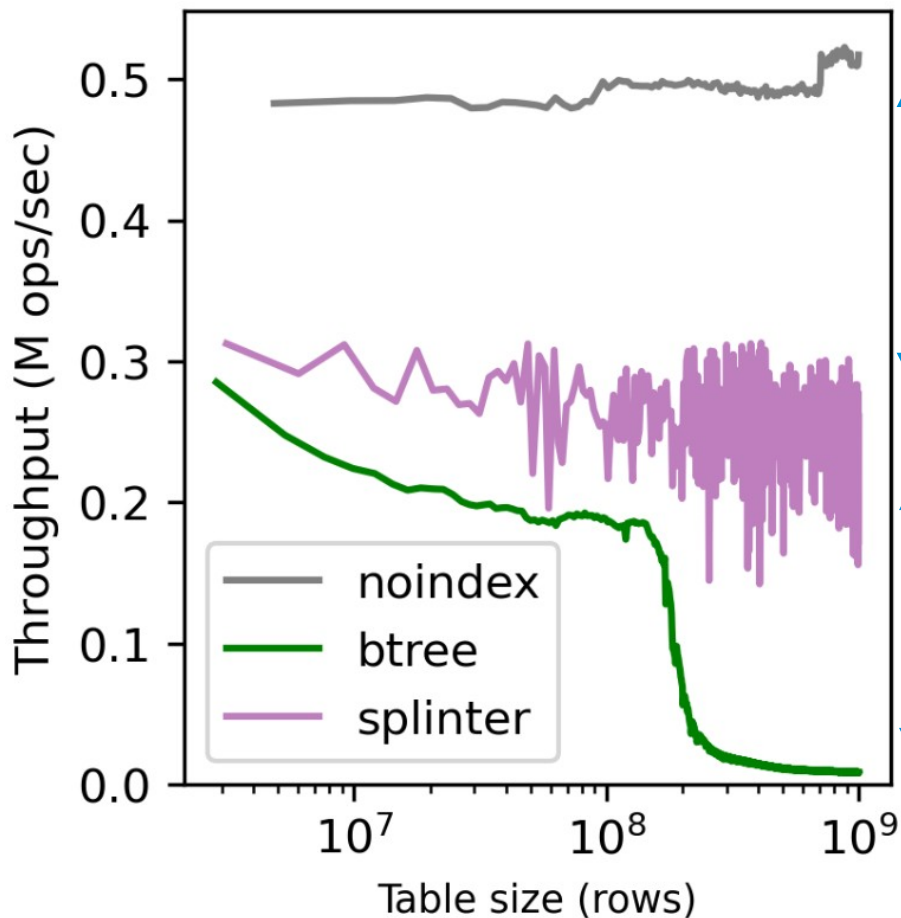# Empirical performance measurements



B-tree

Mapped B$^\varepsilon$-tree

# Random inserts

AWS i4i.16xlarge:
64 CPUs, fast local storage

1 client inserting

1 billion random rows

Higher throughput is
better.



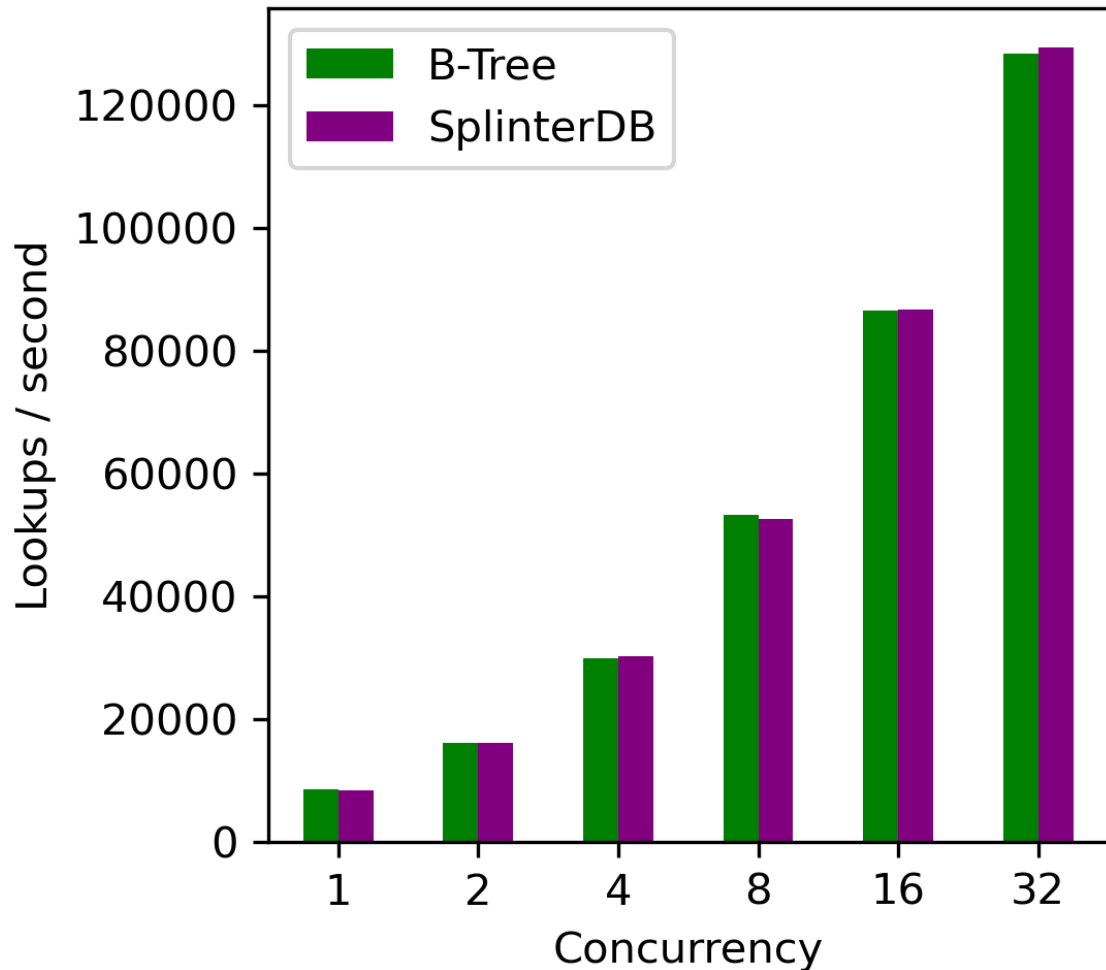SplinterDB is ~50% slower
than no-index at all

SplinterDB throughput
fluctuates due to flushes

SplinterDB is **18x**
faster than B-Tree

# Random Point Queries

AWS i4i.16xlarge:
64 CPUs, fast local storage

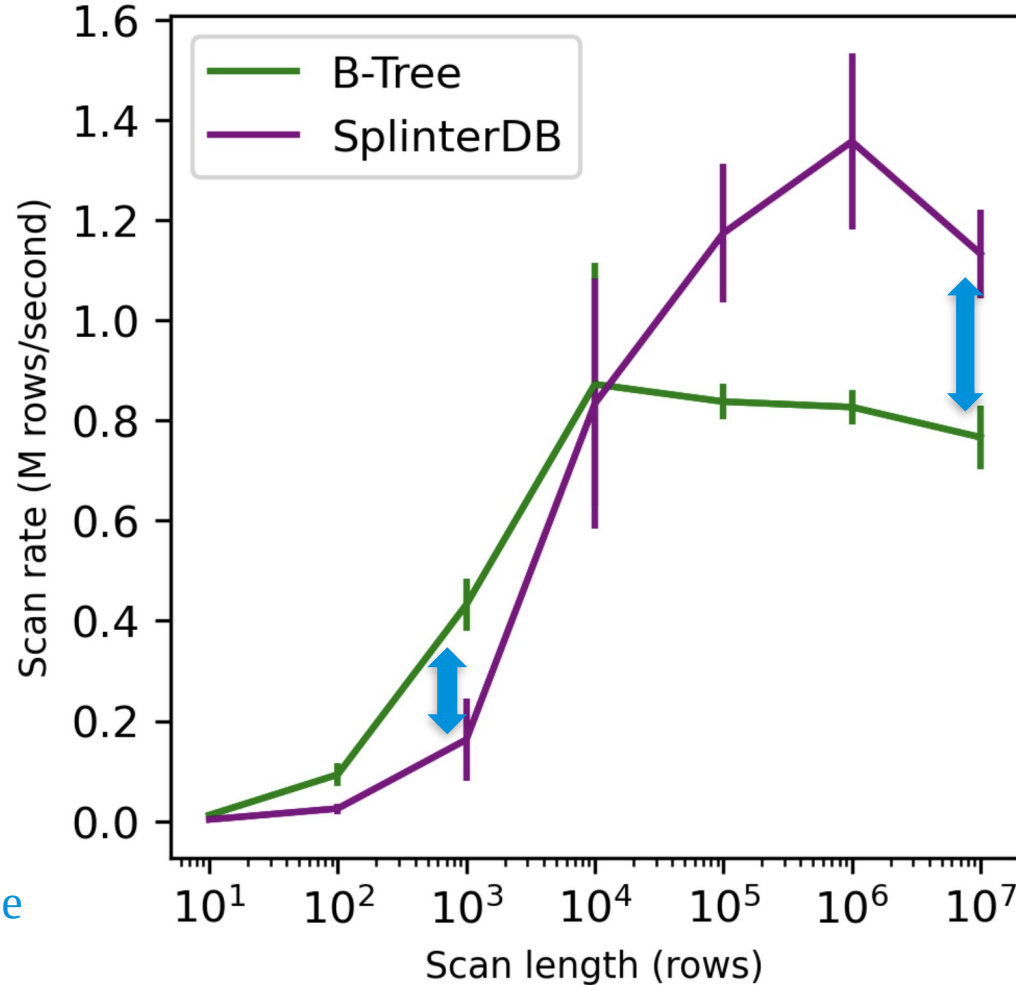1 Billion row table

with unique index.

# Scans



AWS i4i.16xlarge:
64 CPUs, fast local storage

Scans on a
1 billion row table,
from a cold cache

Scanning <= 1000 rows,
Splinter adds 1-5ms extra
latency compared to B-Tree

Scanning >= 100k
rows, Splinter is ~50%
faster than B-Tree.

# Open Questions

Can we get insertion costs below $\log N / B \log \log N$ while keeping sublinear successor queries?

- Maplets are not the bottleneck

Can we improve successor query costs?

- Range maplets?

- Fractional cascading?

Successor lower bounds?