# Performance of Fractal-Tree Databases
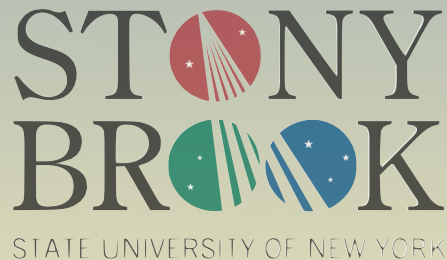
**Michael A. Bender**

**Problem: maintain a dynamic dictionary on disk.**
**Motivation: file systems, databases, etc.**

Michael Bender --  Performance of Fractal-Tree Databases

**Problem: maintain a dynamic dictionary on disk.**
**Motivation: file systems, databases, etc.**

**State of the art (algorithmic perspective):**

- B-tree [Bayer, McCreight 72]
- cache-oblivious B-tree [Bender, Demaine, Farach-Colton 00]
- buffer tree [Arge 95]
- buffered-repository tree [Buchsbaum,Goldwasser,Venkatasubramanian,Westbrook 00]
- $B^\varepsilon$ tree [Brodal, Fagerberg 03]
- log-structured merge tree [O'Neil, Cheng, Gawlick, O'Neil 96]
- string B-tree [Ferragina, Grossi 99]
- etc, etc!

Michael Bender -- Performance of Fractal-Tree Databases

**Problem: maintain a dynamic dictionary on disk.**
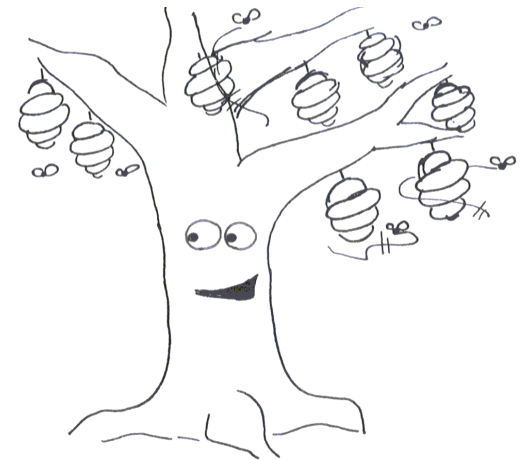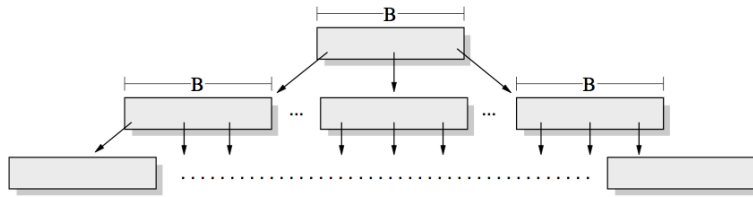**Motivation: file systems, databases, etc.**

**State of the art (algorithmic perspective):**

- B-tree [Bayer, McCreight 72]
- cache-oblivious B-tree [Bender, Demaine, Farach-Colton 00]
- buffer tree [Arge 95]
- buffered-repository tree [Buchsbaum,Goldwasser,Venkatasubramanian,Westbrook 00]
- $B^\varepsilon$ tree [Brodal, Fagerberg 03]
- log-structured merge tree [O'Neil, Cheng, Gawlick, O'Neil 96]
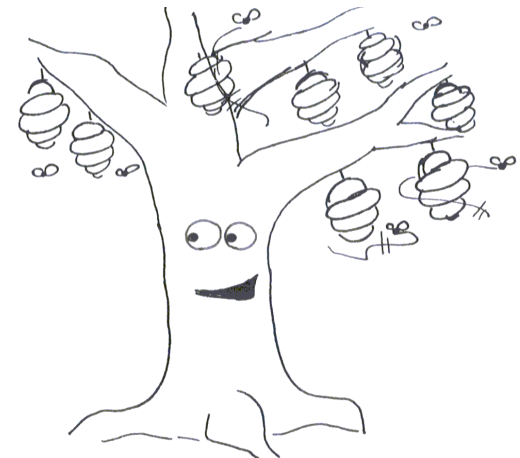- string B-tree [Ferragina, Grossi 99]
- etc, etc!

**State of the practice:**

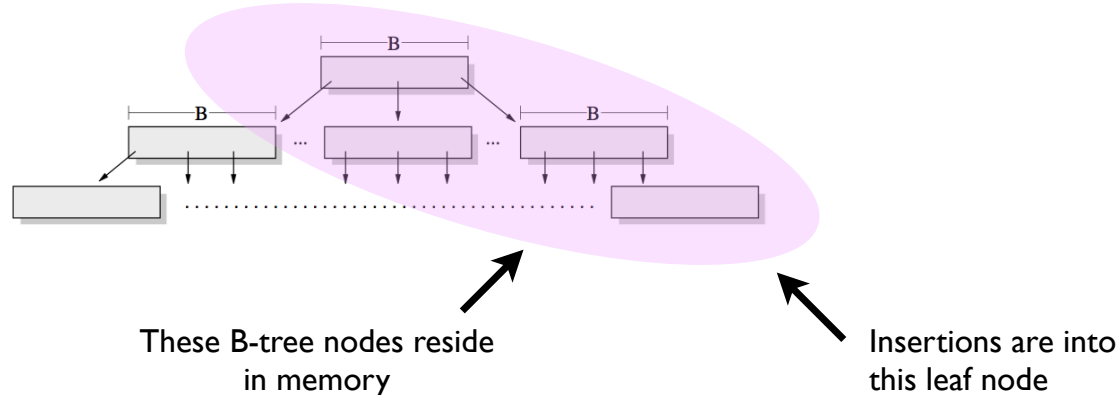- B-trees + industrial-strength features/optimizations

STONY BROOK
STATE UNIVERSITY OF NEW YORK

Michael Bender --  Performance of Fractal-Tree Databases

*Tokutek*™

# B-trees are Fast at Sequential Inserts

# B-trees are Fast at Sequential Inserts

**Sequential inserts in B-trees have near-optimal data locality**



These B-tree nodes reside in memory

Insertions are into this leaf node

- One disk I/O per leaf (which contains many inserts).
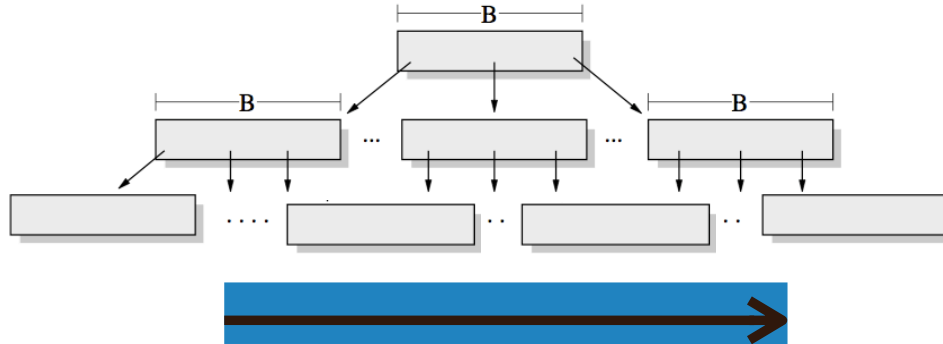- Sequential disk I/O.
- Performance is disk-bandwidth limited.

Michael Bender -- Performance of Fractal-Tree Databases

*Tokutek*™

# B-Trees Are Slow at Ad Hoc Inserts

**High entropy inserts (e.g., random) in B-trees have poor data locality**
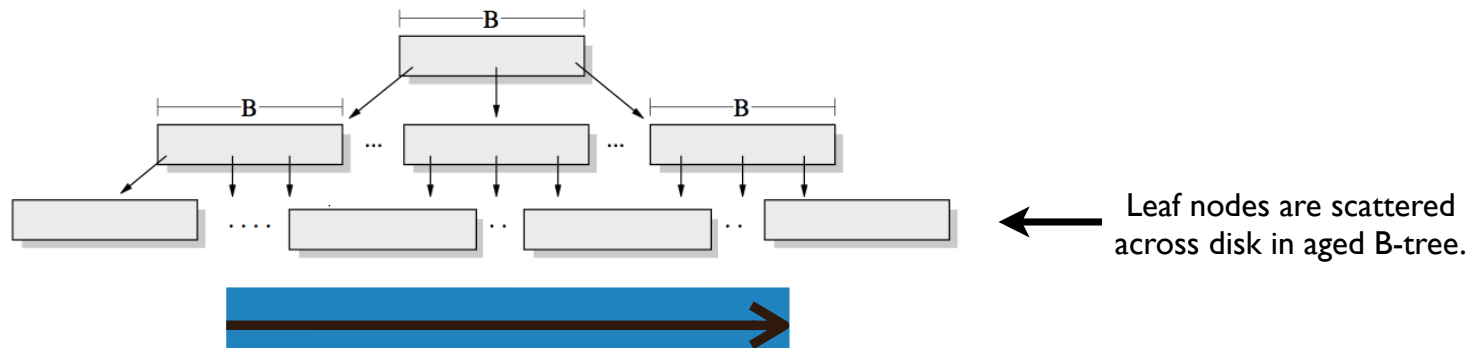
These B-tree nodes reside in memory

- Most nodes are not in main memory.
- Most insertions require a random disk I/O.
- Performance is disk-seek limited.
- ≤ 100 inserts/sec/disk (≤ 0.05% of disk bandwidth).

Michael Bender -- Performance of Fractal-Tree Databases

**Range queries in newly built B-trees have good locality**

Michael Bender -- Performance of Fractal-Tree Databases

# B-trees Have a Similar Story for Range Queries



Leaf nodes are scattered across disk in aged B-tree.

**Range queries in newly built B-trees have good locality**

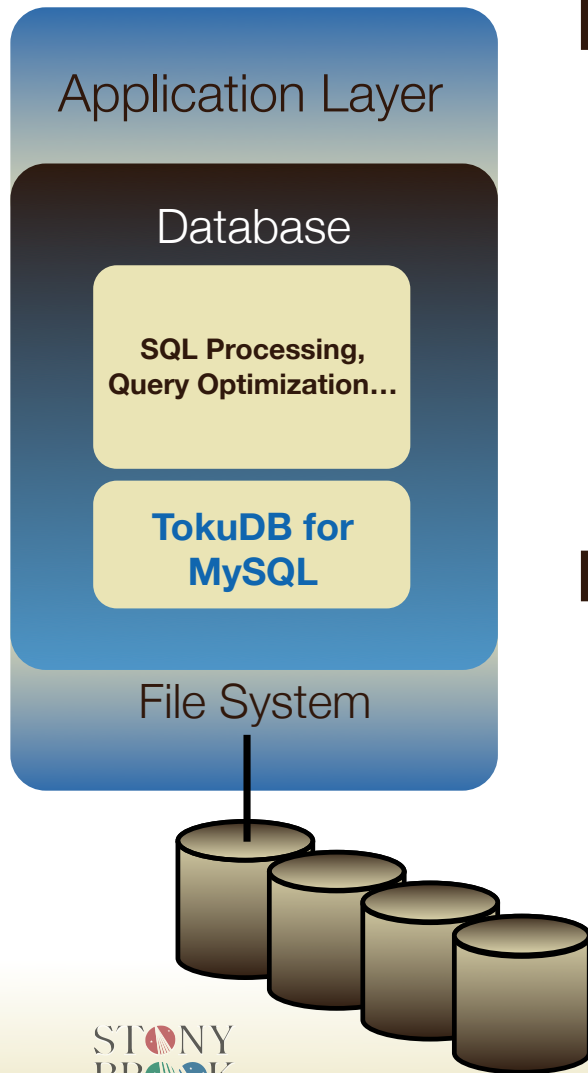**Range queries in aged B-trees have poor locality**

- Leaf blocks are scattered across disk.
- For page-sized nodes, as low as 1% disk bandwidth.

## Cache-Oblivious Streaming B-tree [Bender, Farach-Colton, Fineman, Fogel,  Kuszmaul, Nelson 07]

- Replacement for Traditional B-tree

- High entropy inserts/deletes run up to 100x faster

- No aging --> always fast range queries

- Streaming B-tree is cache-oblivious
  - ▸ Good data locality without memory-specific parameterization.

Michael Bender --  Performance of Fractal-Tree Databases

**Tokutek**™

## Application Layer

### Database

SQL Processing,
Query Optimization…

**TokuDB for MySQL**

File System

# Fractal Tree™ database

- TokuDB is a *storage engine* for MySQL
  - ▸ A storage engine is a structure that stores on-disk data.
  - ▸ Traditionally a storage engine is a B-tree.
- MySQL is an open-source database
  - ▸ Most installations of any database
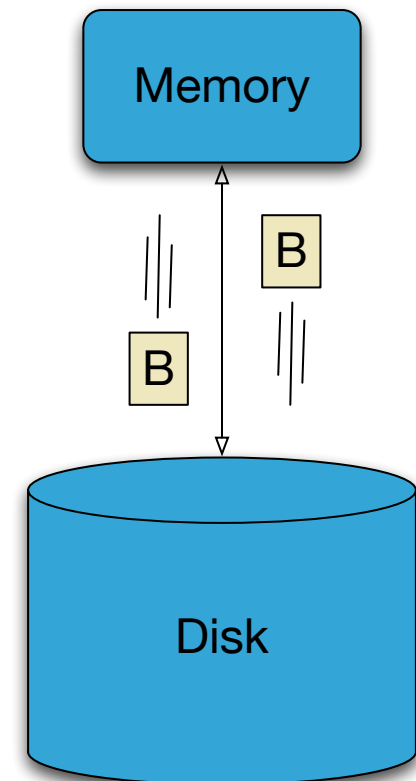- Built in context of our startup Tokutek.

# Performance

- 10x-100x faster index inserts
- No aging
- Faster queries in important cases

Michael Bender -- Performance of Fractal-Tree Databases

**Tokutek**™

## Minimize # of block transfers per operation

## Disk-Access Machine (DAM) [Aggrawal, Vitter 88]

- Two-levels of memory.
- Two parameters:

    block-size $B$, memory-size $M$.



Memory

B

B

Disk
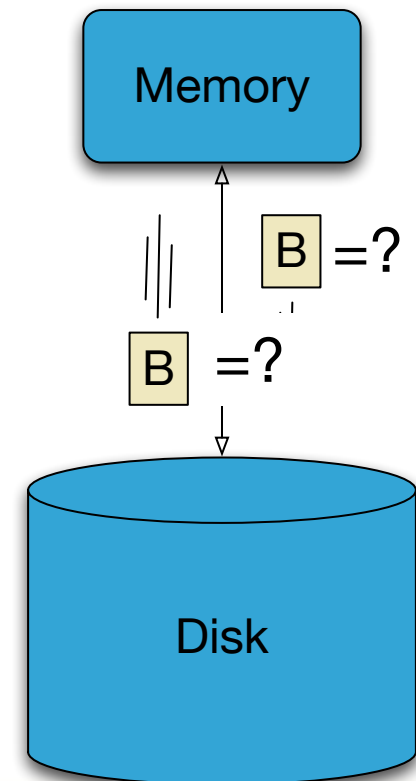
# Algorithmic Performance Model

**Minimize # of block transfers per operation**

**Disk-Access Machine (DAM)** [Aggrawal, Vitter 88]
- Two-levels of memory.
- Two parameters:

    block-size **B**, memory-size **M**.

**Cache-Oblivious Model (CO)** [Frigo, Leiserson, Prokop, Ramachandran 99]
- Parameters **B** and **M** are unknown to the algorithm or coder.
- (Of course, used in proofs.)

Memory

B =?

B =?

Disk

# Fractal Tree Inserts (and Deletes)

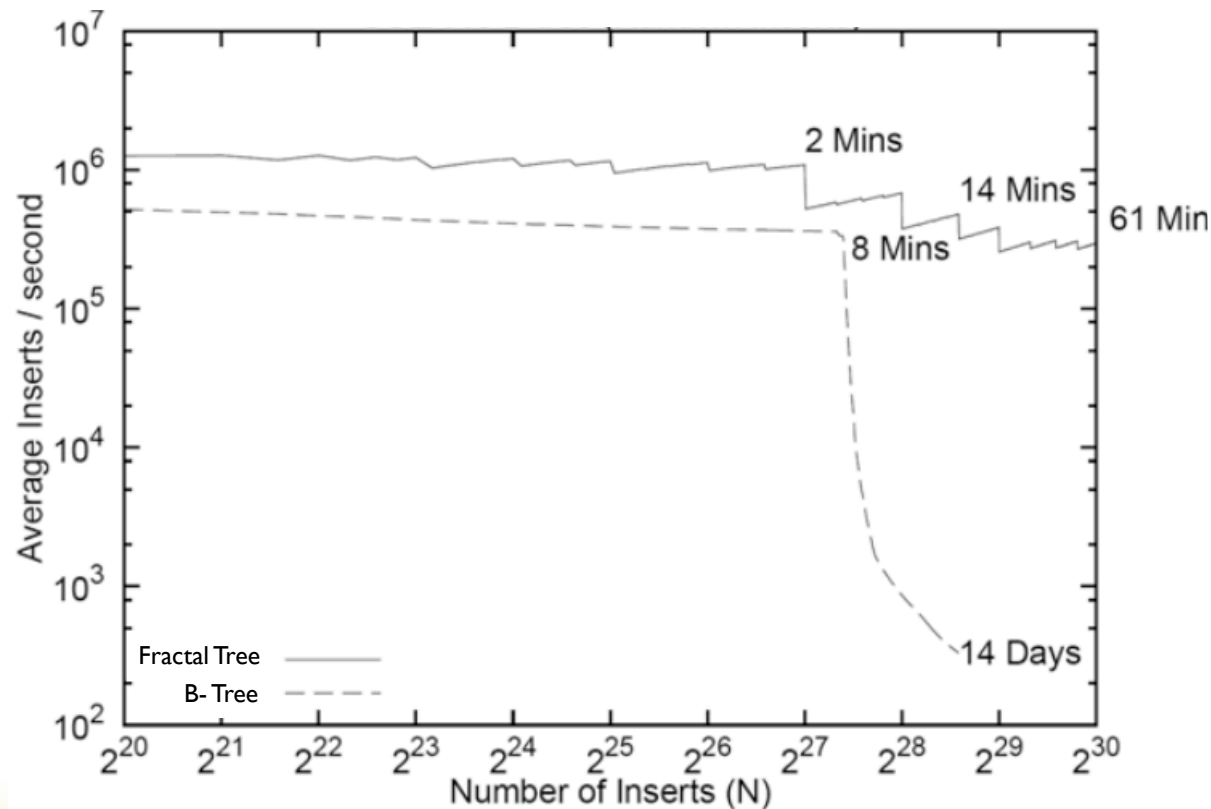| | B-tree | Streaming B-tree |
|---|---|---|
| **Insert** | $O(\log_B N) = O\left(\frac{\log N}{\log B}\right)$ | $O\left(\frac{\log N}{B}\right)$ |

## Example: $N$=1 billion, $B$=4096

- 1 billion 128-byte rows (128 gigabytes)
  - ▶ $\log_2$ (1 billion) = 30
- Half-megabyte blocks that hold 4096 rows each
  - ▶ $\log_2$ (4096) = 12

Michael Bender --  Performance of Fractal-Tree Databases

# Fractal Tree Inserts (and Deletes)

| | B-tree | Streaming B-tree |
|---|---|---|
| Insert | $O(\log_B N) = O\left(\frac{\log N}{\log B}\right)$ | $O\left(\frac{\log N}{B}\right)$ |

## Example: *N*=1 billion, *B*=4096

- 1 billion 128-byte rows (128 gigabytes)
  - ▸ $\log_2$ (1 billion) = 30
- Half-megabyte blocks that hold 4096 rows each
  - ▸ $\log_2$ (4096) = 12

- B-trees require $\dfrac{\log N}{\log B}$ = 30/12 = 3 disk seeks (modulo caching, insertion pattern)
- Streaming B-trees require $\dfrac{\log N}{B}$ = 30/4096 = 0.007 disk seeks
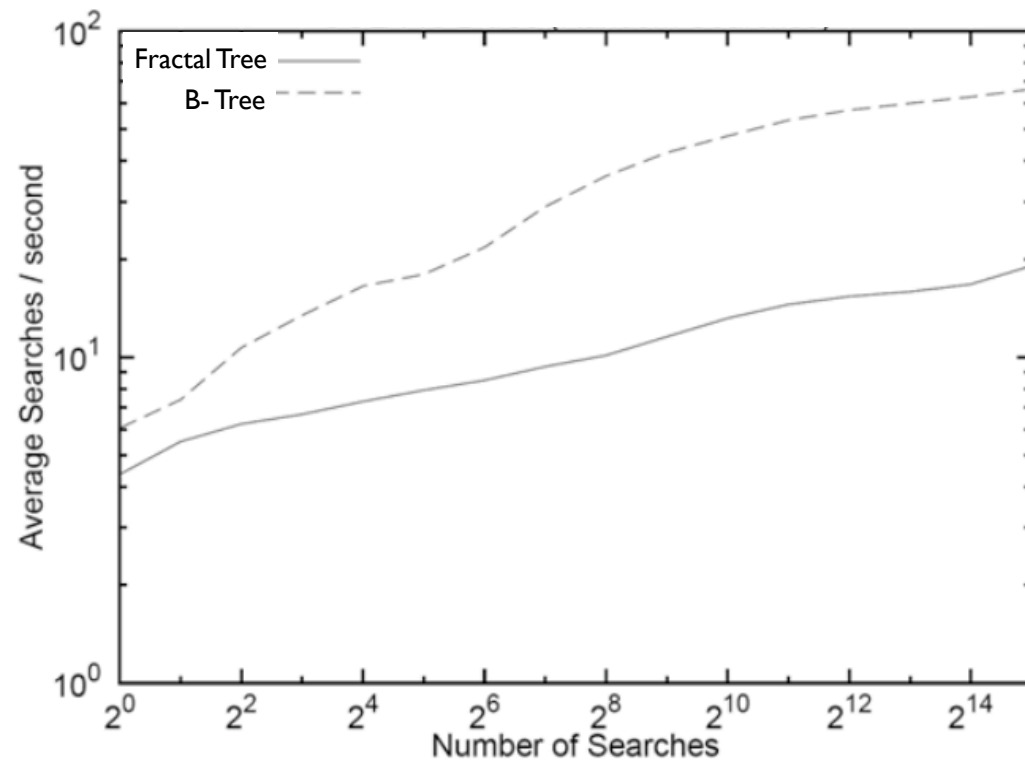
## Random Inserts into Fractal Tree ("streaming B-tree") and B-tree (Berkeley DB)

Michael Bender -- Performance of Fractal-Tree Databases

**Point searches ~3.5x slower (N=$2^{30}$)**

- Searches/sec improves as more of data structure fits in cache)

Michael Bender -- Performance of Fractal-Tree Databases

**Small specification changes affect complexity**

**E.g., duplicate keys**

- Slow: Return an error when a duplicate key is inserted
  ▸ Hidden search

- Fast: Overwrite duplicates or maintain all versions
  ▸ No hidden search

# Asymmetry Between Inserts and Key Searches

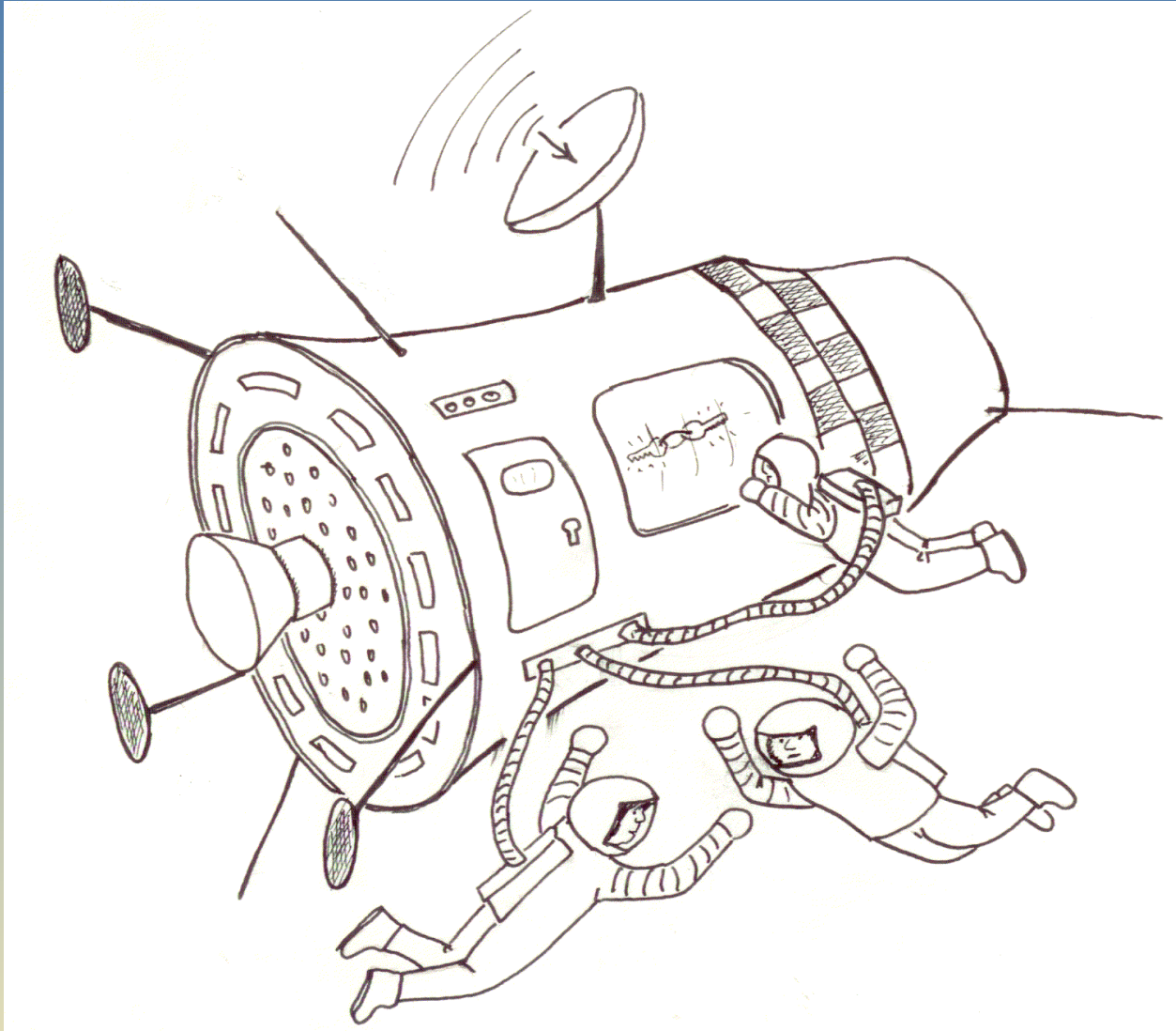## Small specification changes affect complexity

## E.g., duplicate keys

- Slow: Return an error when a duplicate key is inserted
  - ▸ Hidden search
- Fast: Overwrite duplicates or maintain all versions
  - ▸ No hidden search

## E.g. deletes

- Slow: Return number of elements deleted
  - ▸ Hidden search
- Fast: Delete without feedback
  - ▸ No hidden search

# Asymmetry Between Inserts and Key Searches

**Small specification changes affect complexity**

**E.g., duplicate keys**

- Slow: Return an error when a duplicate key is inserted
  ‣ Hidden search

- Fast: Overwrite duplicates or maintain all versions
  ‣ No hidden search

**E.g. deletes**

- Slow: Return number of elements deleted
  ‣ Hidden search

- Fast: Delete without feedback
  ‣ No hidden search

**Next slide: extra difficulty of key searches**
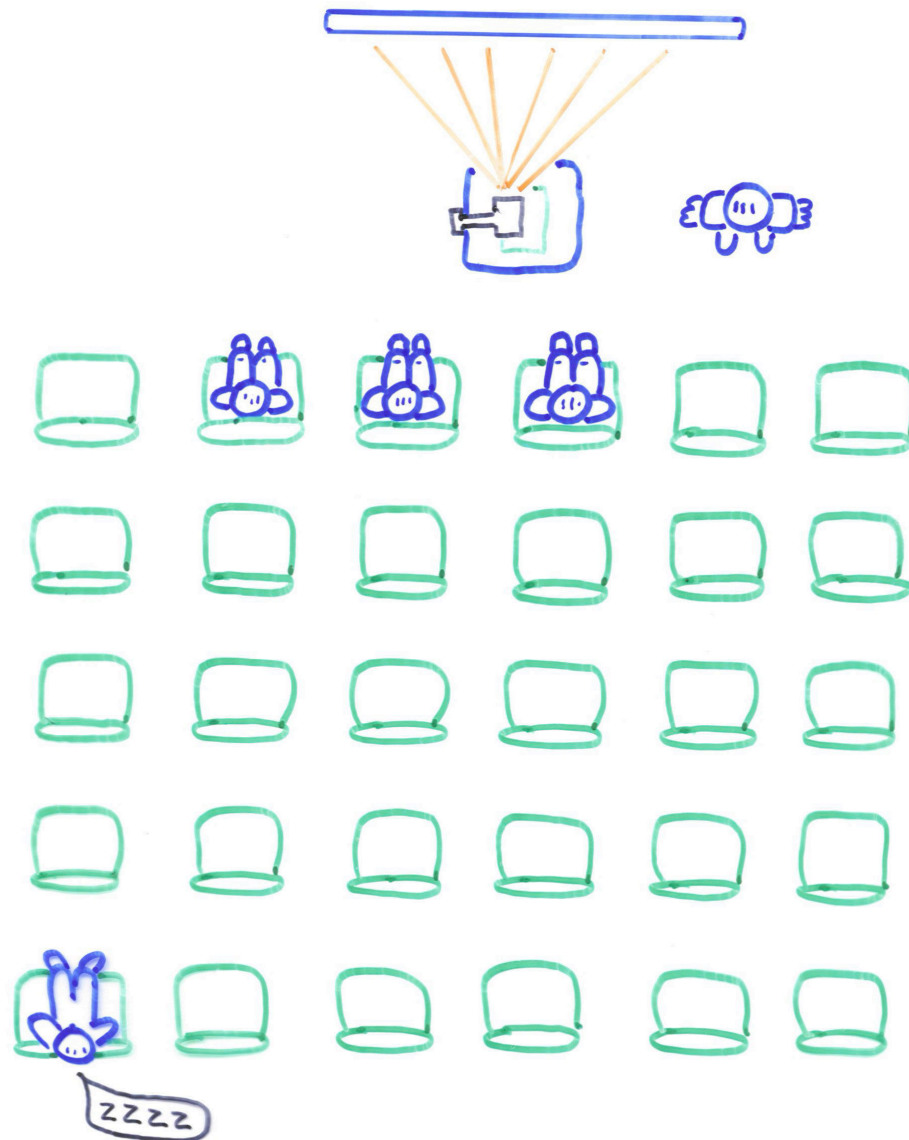
# Extra Difficulty of Key Searches

# Asymmetry Between Inserts and Key Searches

**Inserts/point query asymmetry has impact on**

- *System design.* How to redesign standard mechanisms (e.g., concurrency-control mechanism).

- *System use.* How to take advantage of faster inserts (e.g., to enable faster queries).

# Overview of Talk

External-memory dictionaries

Performance limitations of B-trees

Fractal-Tree data structure (Streaming B-tree)
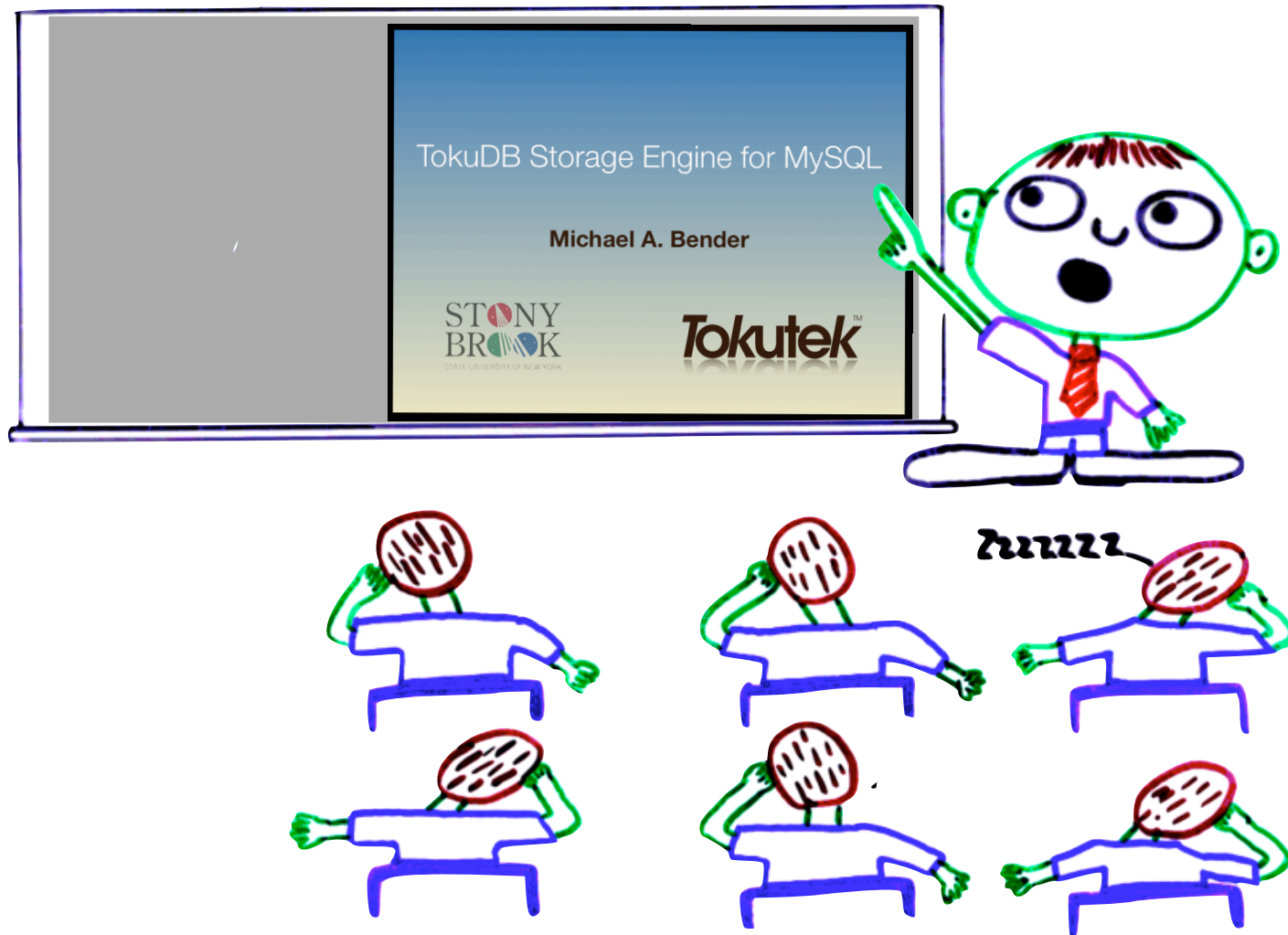
Search/point-query asymmetry

**Impact of search/point-query asymmetry on database use**
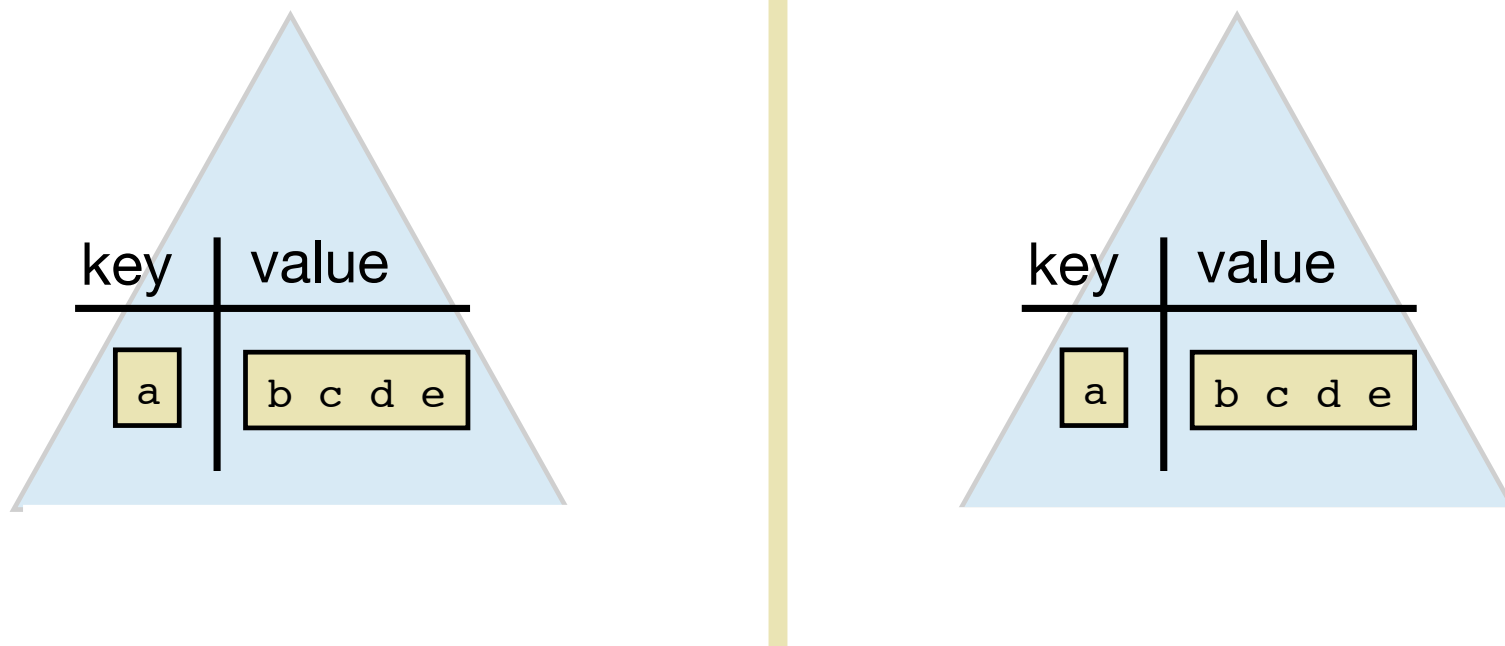
**How to build a streaming B-tree**

**Impact of search/point-query asymmetry on system design**

**Scaling into the future**

Michael Bender -- Performance of Fractal-Tree Databases

# Search/point-query asymmetry affecting database use

| key | value |
| --- | --- |
| a | b c d e |

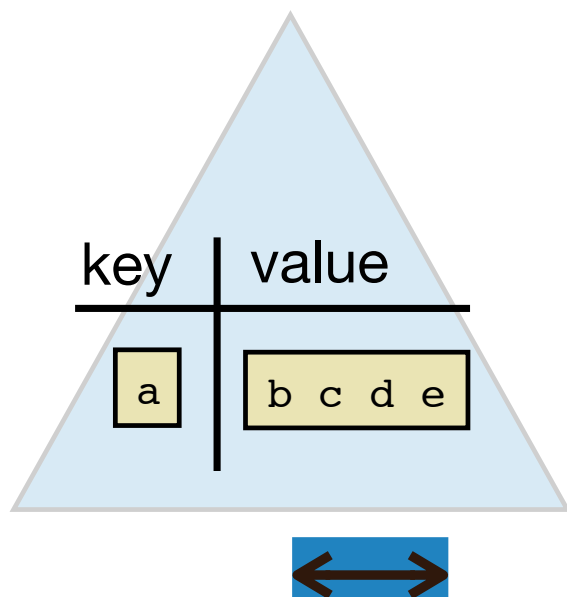| key | value |
| --- | --- |
| a | b c d e |

Data maintained in rows and stored in B-trees.

# How B-trees Are Used in Databases

**Select via Index**

`select` d `where 270 ≤` a `≤ 538`

key | value

a | b c d e

**Select via Table Scan**

`select` d `where 270 ≤` e `≤ 538`

key | value

a | b c d e

Data maintained in rows and stored in B-trees.

**Selecting via an index can be slow, if it is coupled with point queries.**

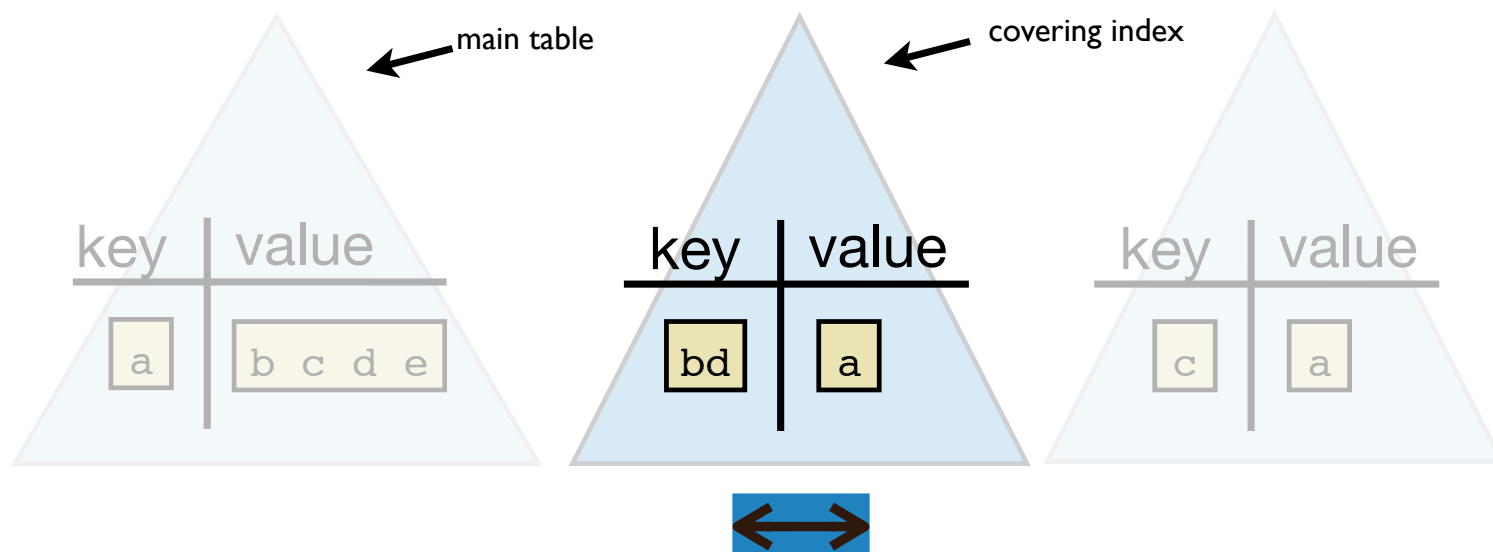`select` **d** `where` **270 ≤ b ≤ 538**

Michael Bender -- Performance of Fractal-Tree Databases

*Tokutek* ™

## Covering index can speed up selects

- Key contains all columns necessary to answer query.

select **d** where 270 ≤ **b** ≤ 538



main table

covering index

| key | value |
|-----|-------|
| a | b c d e |

| key | value |
|-----|-------|
| bd | a |

| key | value |
|-----|-------|
| c | a |

Michael Bender -- Performance of Fractal-Tree Databases

**People often don't use these indexes.**
**They use simplistic schema.**

- Sequential inserts via autoincrement key
- Few indexes, few covering indexes

Autoincrement key
(effectively a timestanp) →

| key | value |
|-----|-------|
| t   | a b c d e |

**Then insertions are fast but queries are slow.**

Michael Bender -- Performance of Fractal-Tree Databases

**People often don't use these indexes.
They use simplistic schema.**

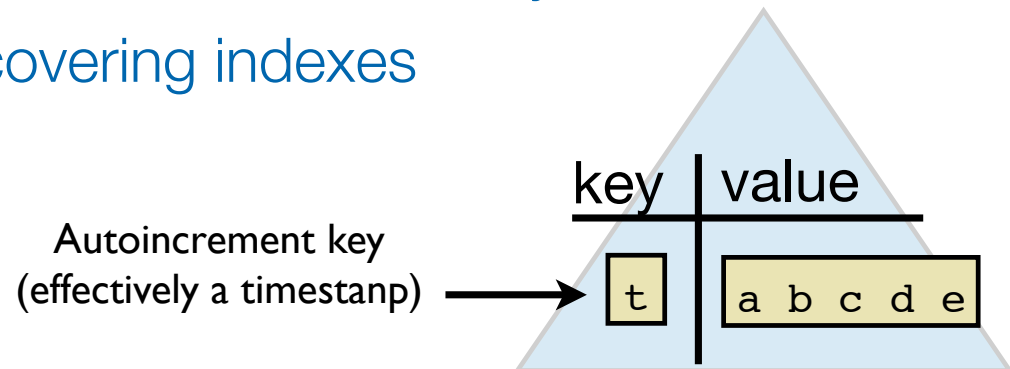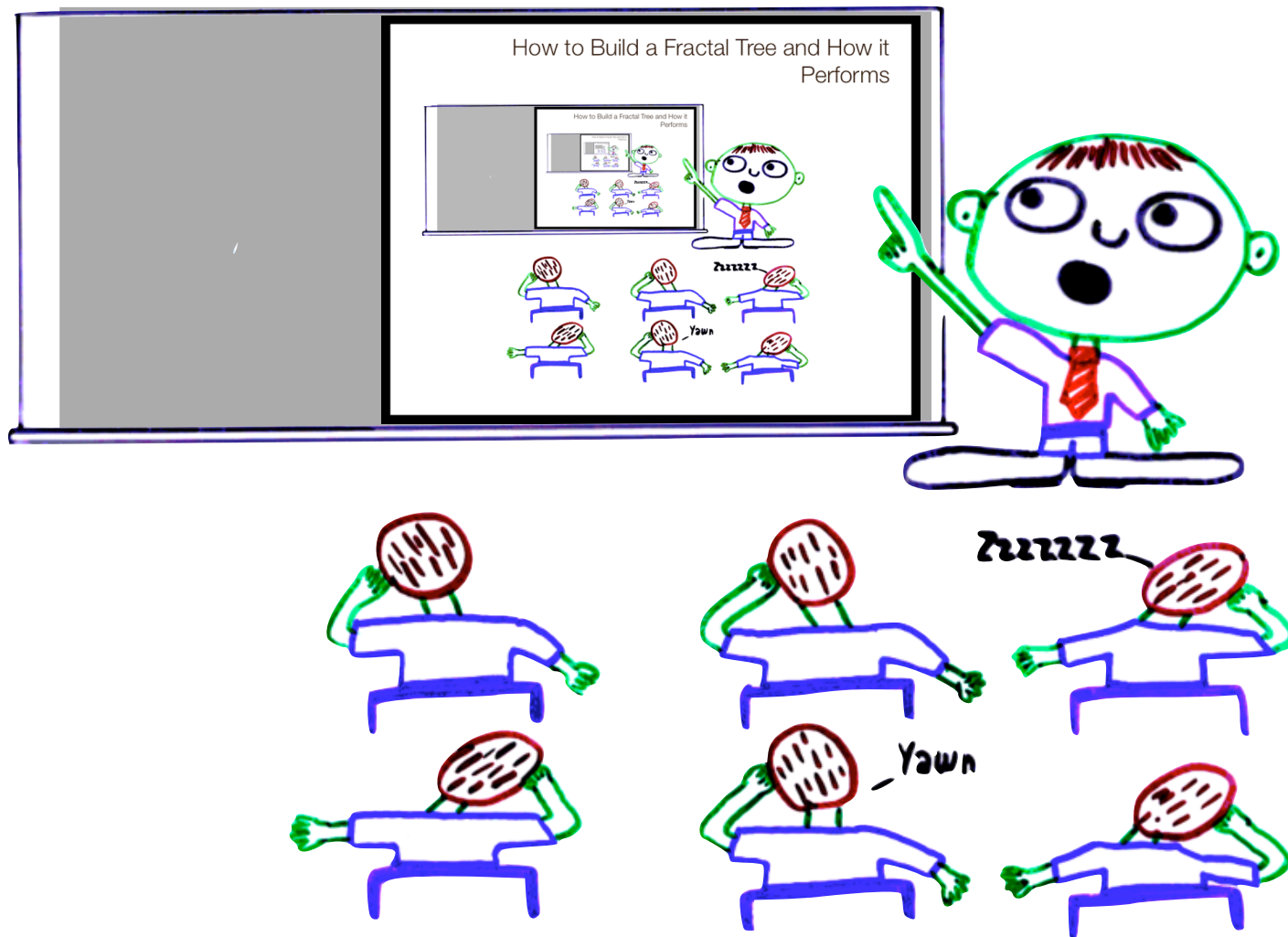- Sequential inserts via autoincrement key

- Few indexes, few covering indexes

Autoincrement key
(effectively a timestanp) →

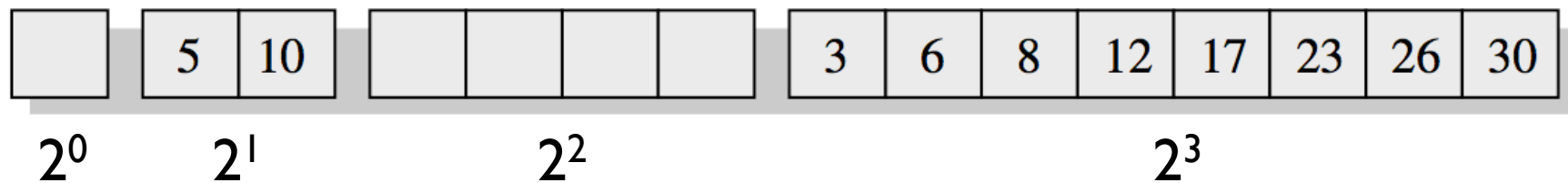| key | value |
|-----|-------|
| t | a b c d e |

**Then insertions are fast but queries are slow.**

**Adding sophisticated indexes helps queries**

- B-trees cannot afford to maintain them.
  Fractal Trees can.

Michael Bender -- Performance of Fractal-Tree Databases
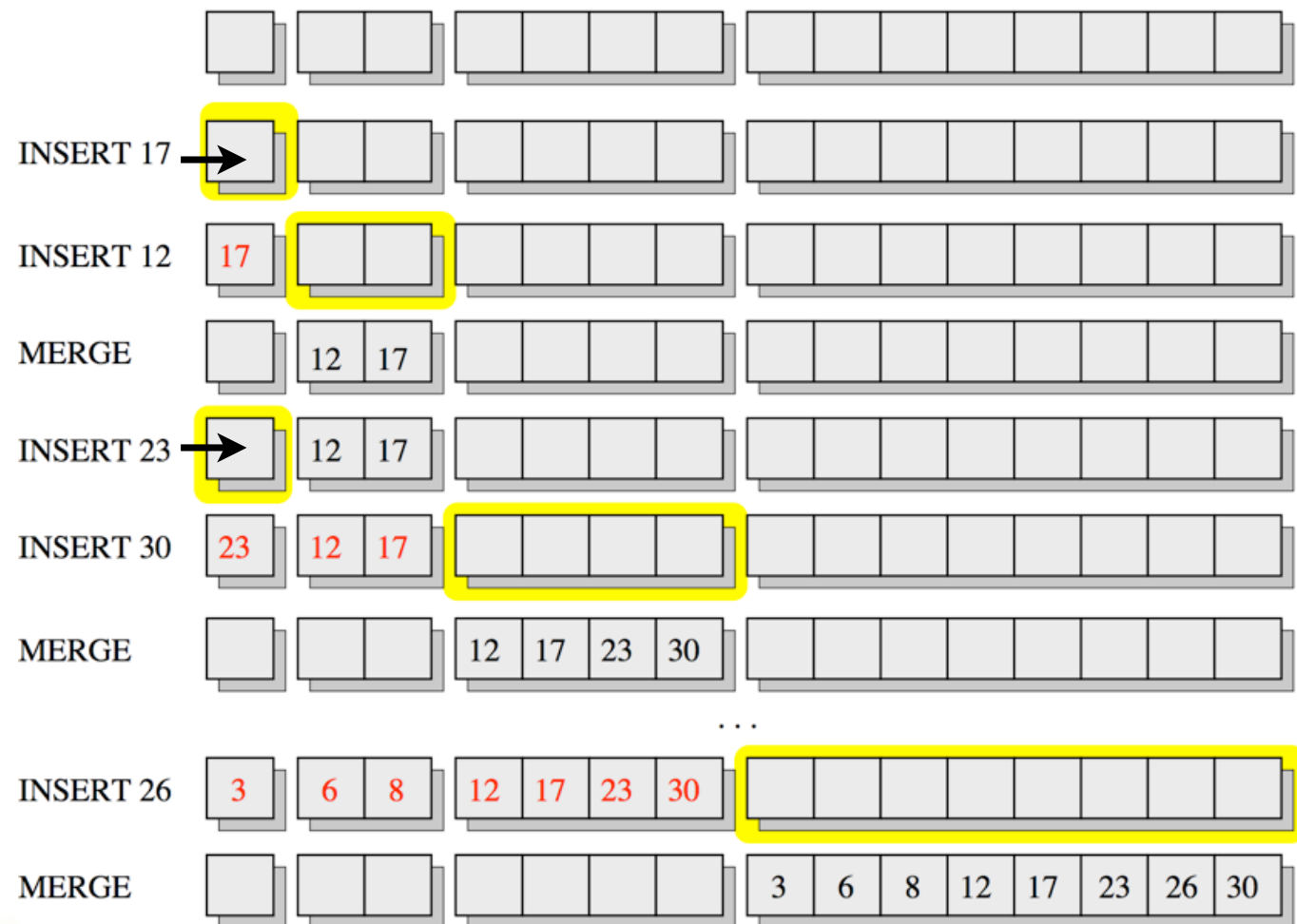
# How to Build a Fractal Tree and How it Performs

# Simplified (Cache-Oblivious) Fractal Tree

| | | 5 | 10 | | | | | 3 | 6 | 8 | 12 | 17 | 23 | 26 | 30 |

$2^0$      $2^1$          $2^2$                $2^3$

## $O((\log N)/B)$ insert cost  &  $O(\log^2 N)$ search cost

- Sorted arrays of exponentially increasing size.

- Arrays are completely full or completely empty (depends on the bit representation of # of elmts).

- Insert into the smallest array.
  Merge arrays to make room.

Michael Bender -- Performance of Fractal-Tree Databases

# Analysis of Simplified Fractal Tree

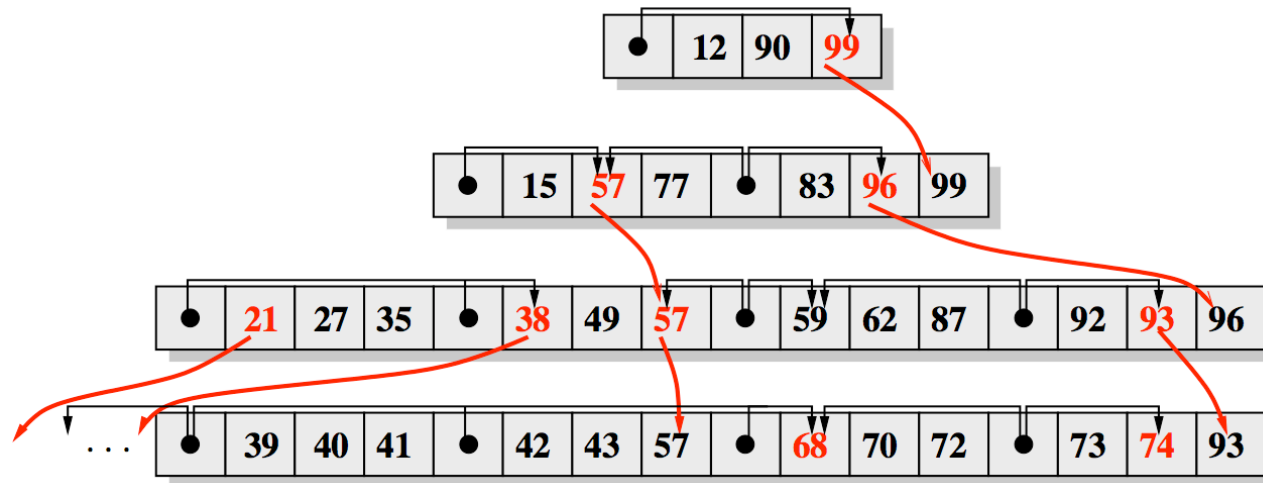| 17 | | 5 | 10 | | 13 | 41 | 57 | 90 | | 3 | 6 | 8 | 12 | 17 | 23 | 26 | 30 |
|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|

## Insert Cost:

- cost to flush buffer of size $X$ = O($X/B$)

- cost per element to flush buffer = O(1/$B$)

- max # of times each element is flushed = log $N$

- insert cost = O((log $N$))/$B$) amortized memory transfers

## Search Cost

- Binary search at each level

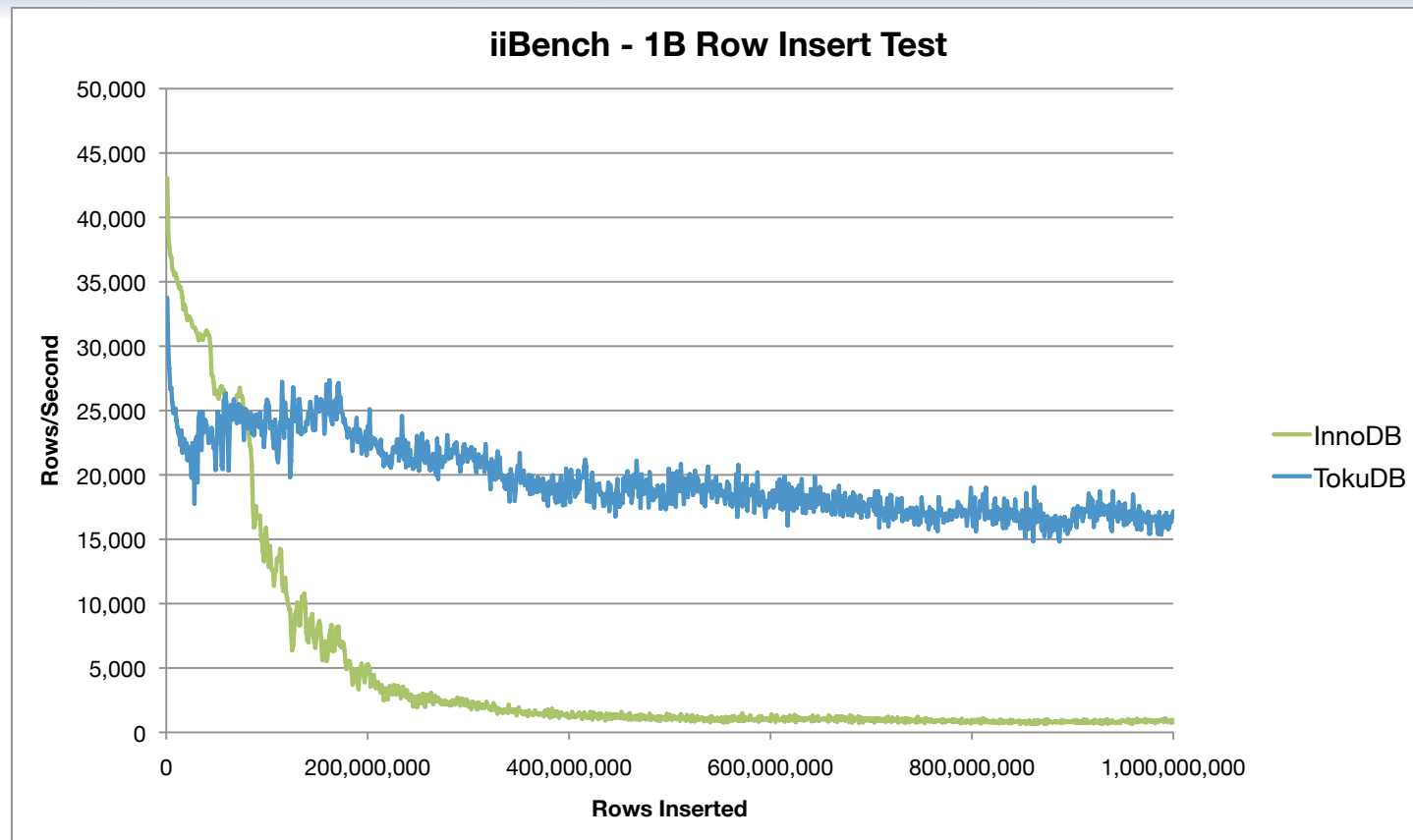- log($N/B$) + log($N/B$) - 1 + log($N/B$) - 2 + ... + 2 + 1
  = O(log$^2$($N/B$))

## O(log (*N/B*)) search cost

- Some redundancy of elements between levels
- Arrays can be partially full
- Horizontal and vertical pointers to redundant elements
- (Fractional Cascading)

# Why The Previous Data Structure is a Simplification

- Need concurrency-control mechanisms
- Need crash safety
- Need transactions, logging+recovery
- Need better search cost
- Need to store variable-size elements
- Need better amortization
- Need to be good for random and sequential inserts
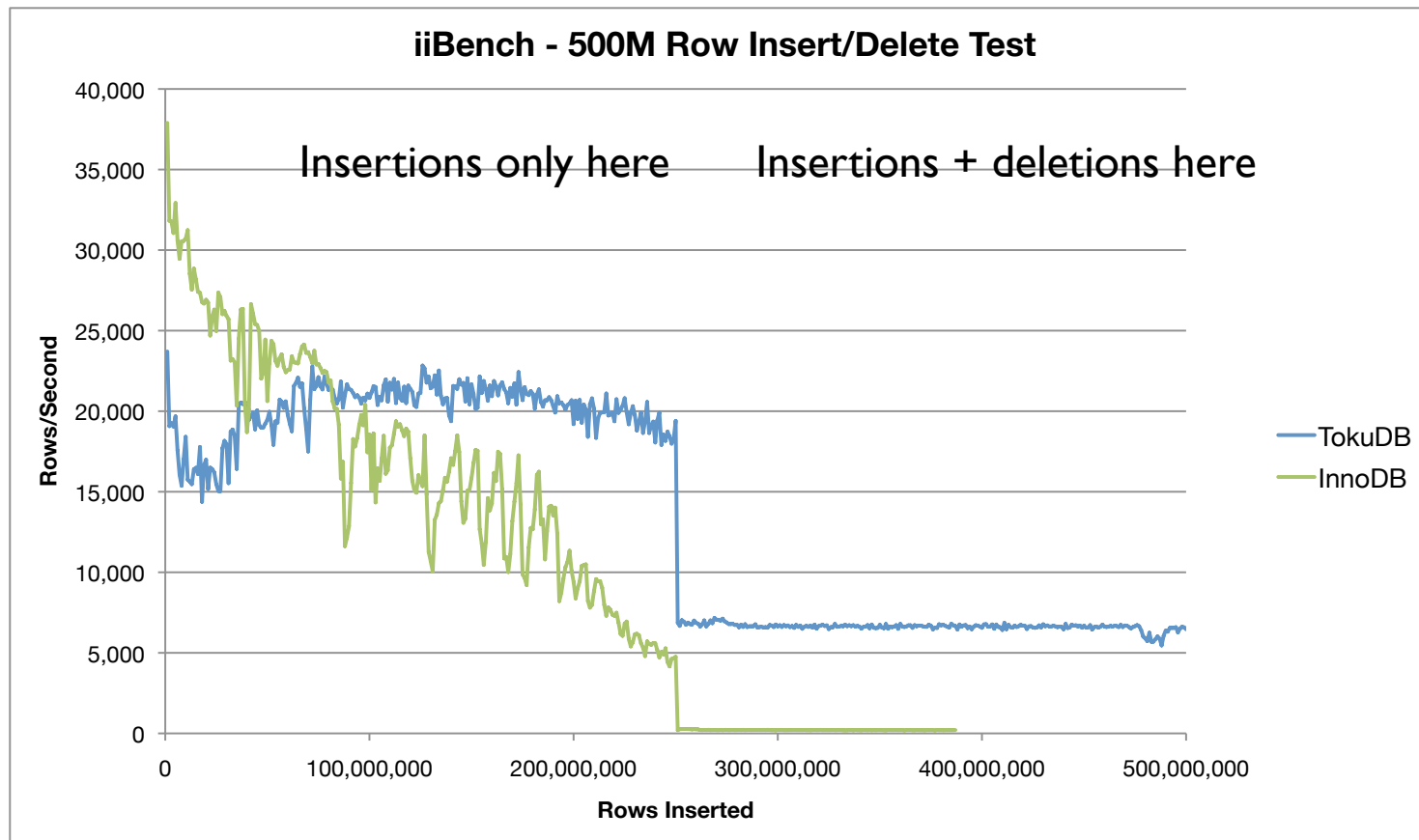- Need to support multithreading.
- Need compression
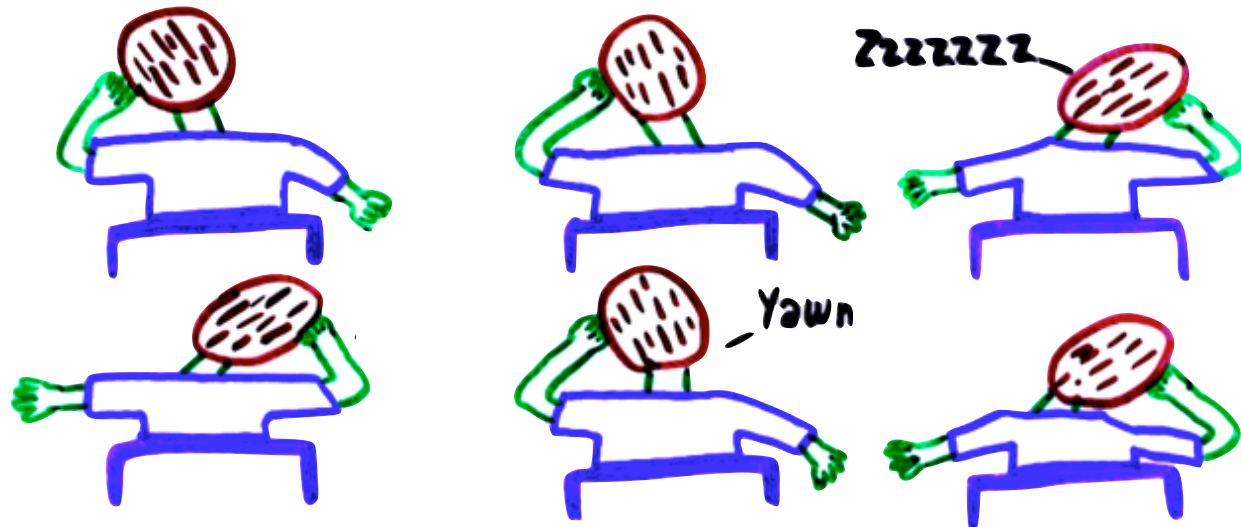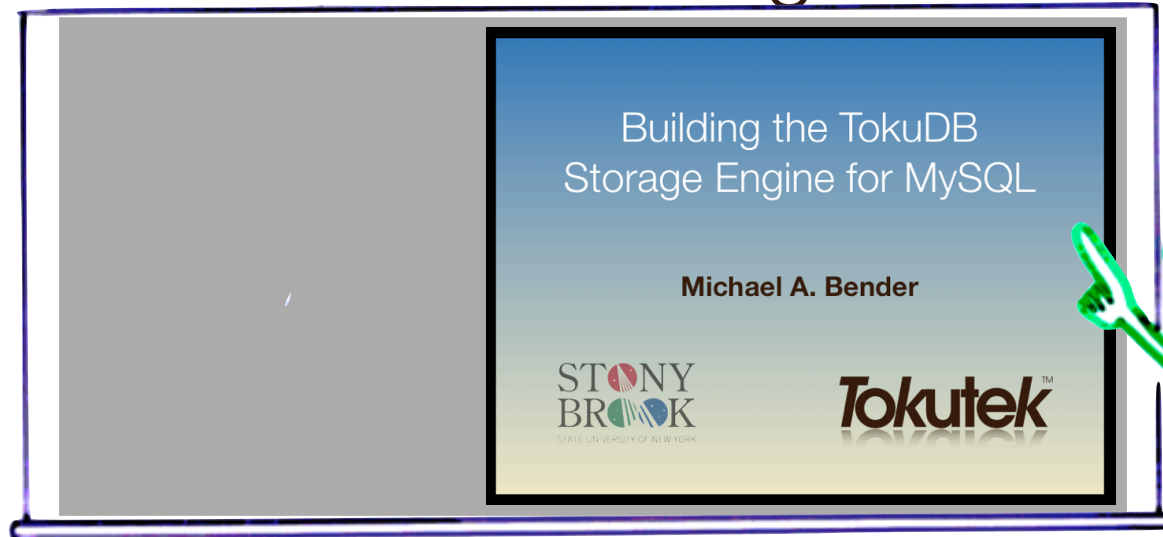
# iiBench Insertion Benchmark



**iiBench - 1B Row Insert Test**

Legend: InnoDB, TokuDB

Y-axis: Rows/Second (0 to 50,000)
X-axis: Rows Inserted (0 to 1,000,000,000)

**Fractal Trees scale with disk bandwidth not seek time.**

- In fact, now we are compute bound, so cannot yet take full advantage of more cores or disks. (This will change.)

iiBench - 500M Row Insert/Delete Test

# Search/point query asymmetry when building Fractal-Tree Database

# Building TokuDB Storage Engine for MySQL

## Engineering to do list

- Need concurrency-control mechanisms
- Need crash safety
- Need transactions, logging+recovery
- Need better search cost
- Need to store variable-size elements
- Need better amortization
- Need to be good for random and sequential inserts
- Need to support multithreading.
- Need compression

# Building TokuDB Storage Engine for MySQL

## Engineering to do list

- **Need concurrency-control mechanisms**
- Need crash safety
- Need transactions, logging+recovery
- Need better search cost
- Need to store variable-size elements
- Need better amortization
- Need to be good for random and sequential inserts
- Need to support multithreading.
- Need compression

Michael Bender -- Performance of Fractal-Tree Databases

## Transactions

- Sequence of durable operations.
- Happen atomically.

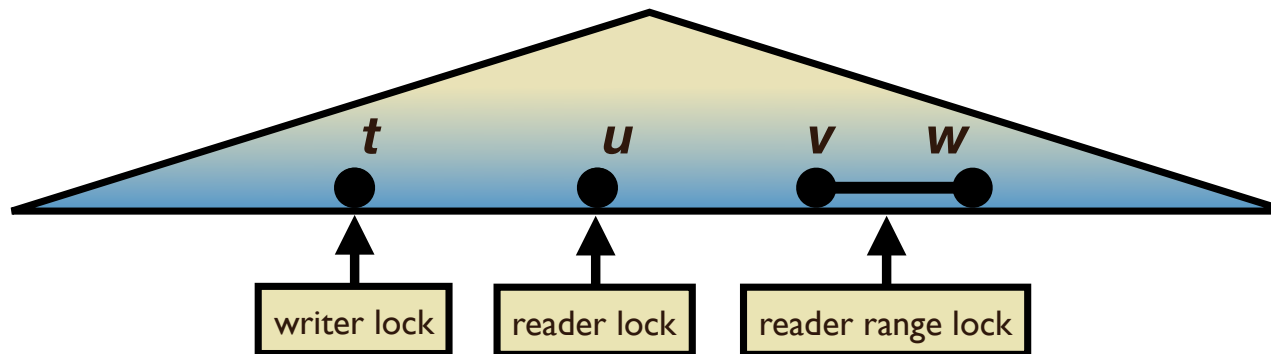## Atomicity in TokuDB via pessimistic locking

- **readers lock:** *A* and *B* can both read row *x* of database.
- **writers lock:** if *A* writes to row *x*, *B* cannot read *x* until *A* completes.
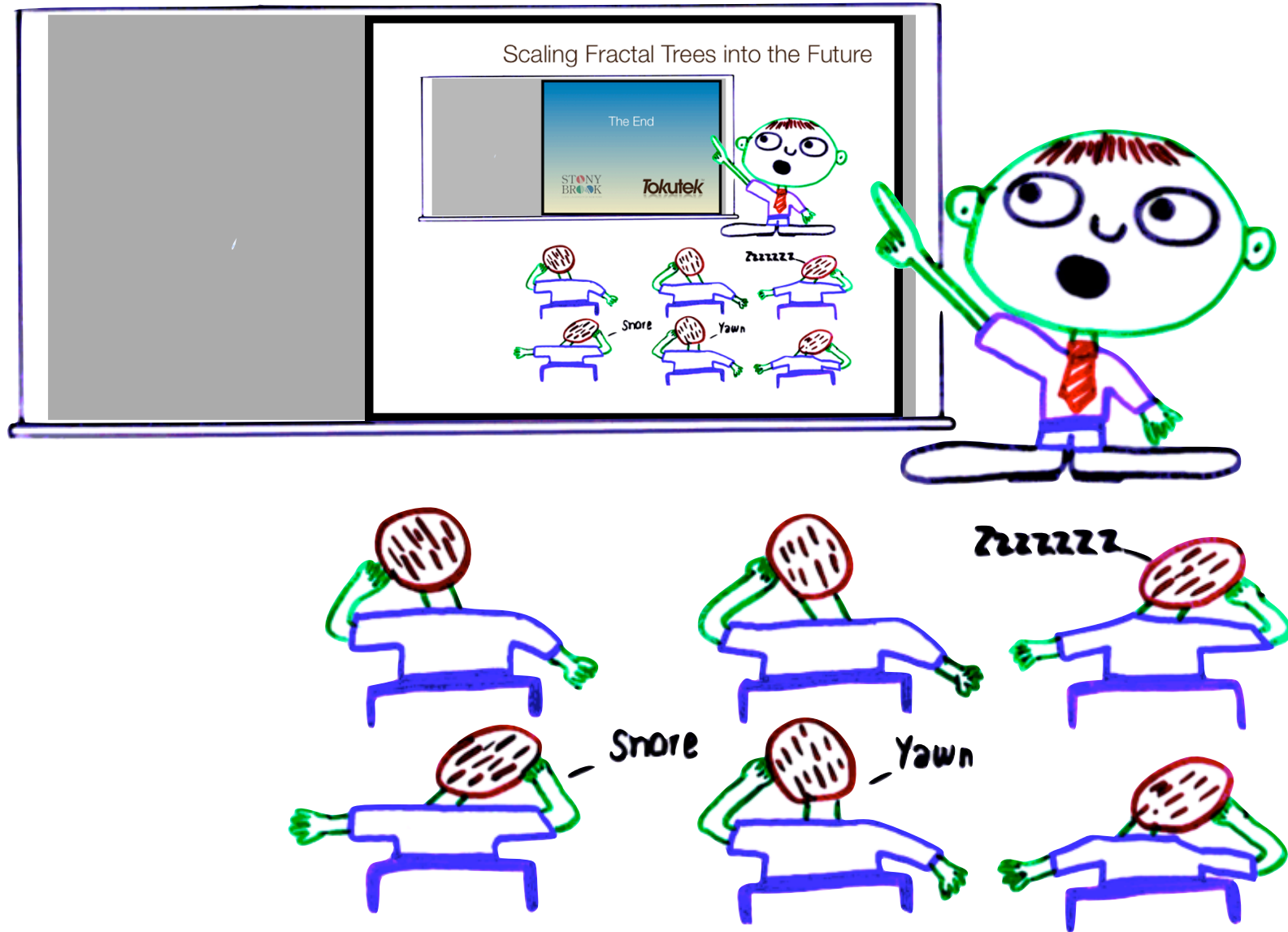
## B-tree implementation: maintain locks in leaves

- Insert row *t*
- Search for row *u*
- Search for row *v* and put a cursor
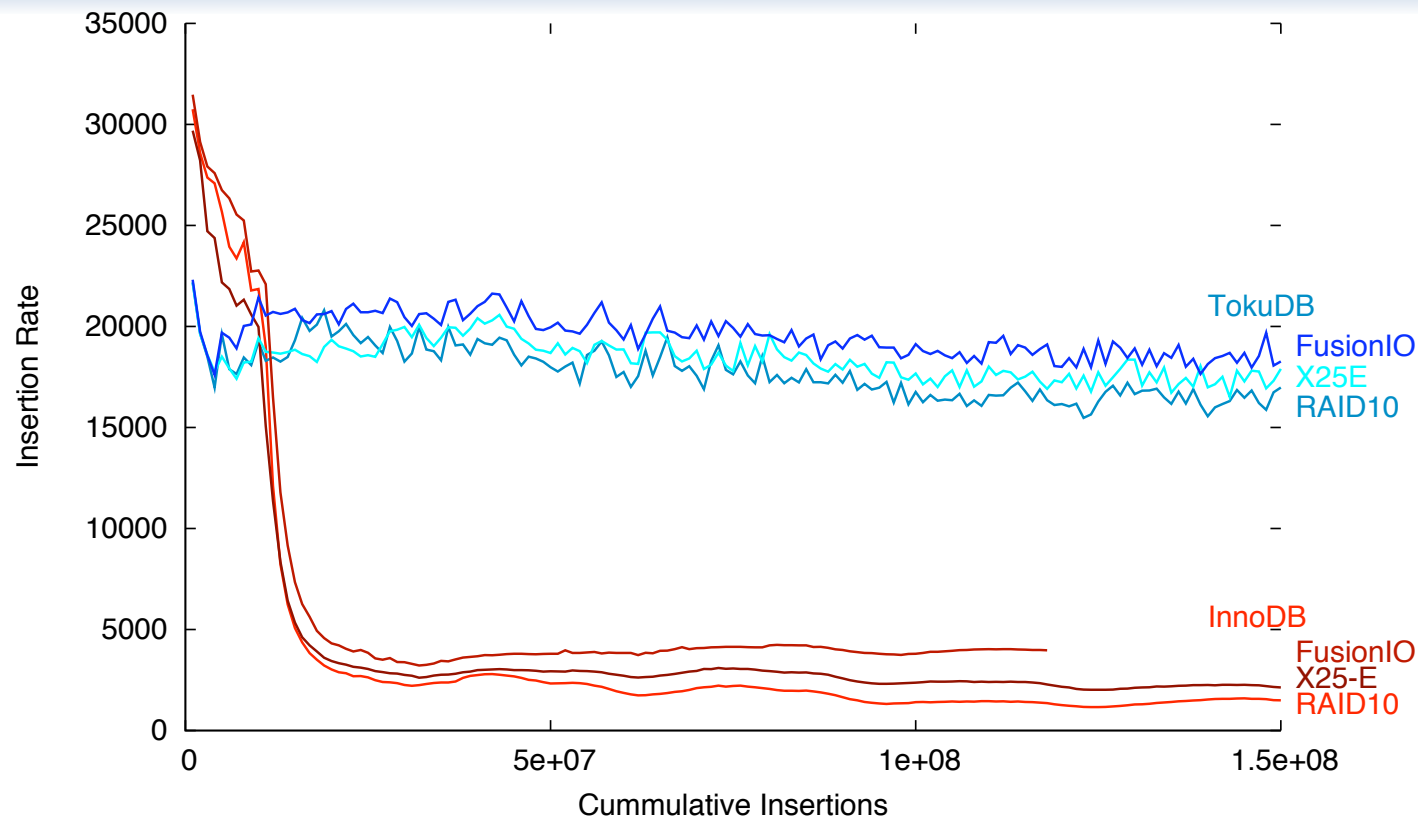- Increment cursor. Now cursor points to row *w*.



***Doesn't work for Fractal Trees: maintaining locks involves implicit searches on writes.***

Michael Bender -- Performance of Fractal-Tree Databases

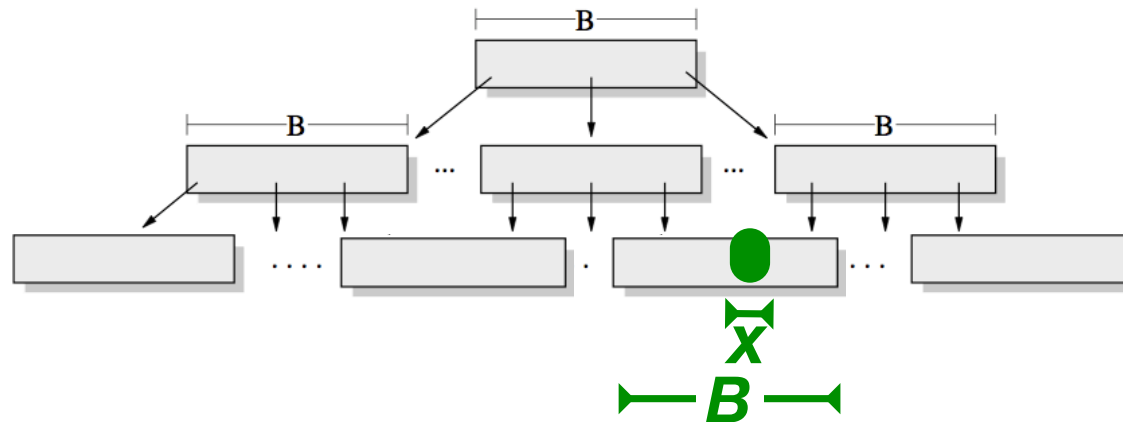# Scaling Fractal Trees into the Future

**B-trees are slow on SSDs, probably b/c they waste bandwidth.**

- When inserting one row, a whole block (much larger) is written.
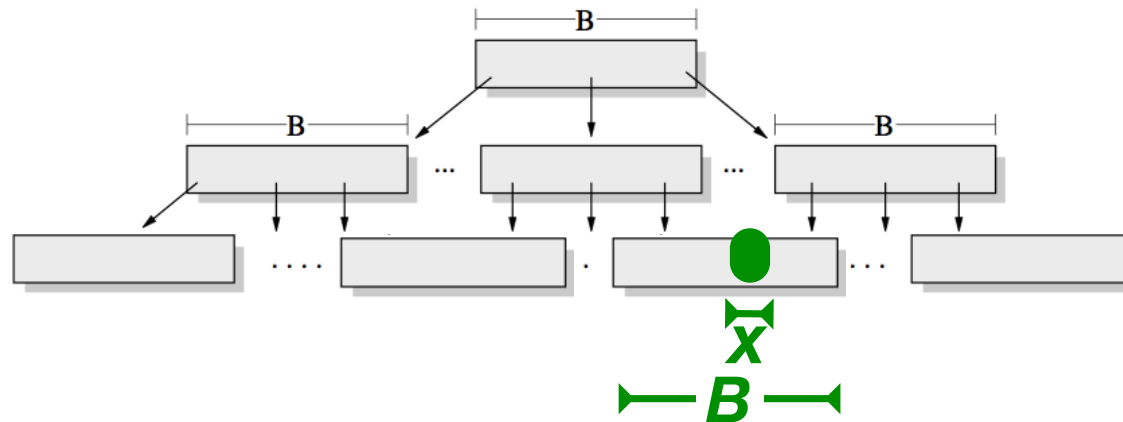
# B-tree Inserts Are Slow on SSDs

**Inserting an element of size *x* into a B-tree dirties a leaf block of size *B*.**



**We can write keys of size *x* into a B-tree using at most a O(*x/B*) fraction of disk bandwidth.**

**Inserting an element of size $x$ into a B-tree dirties a leaf block of size $B$.**



**We can write keys of size $x$ into a B-tree using at most a $O(x/B)$ fraction of disk bandwidth.**

*Fractal trees do efficient inserts on SSDs because they transform random I/O into sequential I/O.*

# Disk Hardware Trends

## Disk capacity will continue to grow quickly

| Year | Capacity | Bandwidth |
|------|----------|-----------|
| 2008 | 2 TB | 100MB/s |
| 2012 | 4.5 TB | 150MB/s |
| 2017 | 67 TB | 500MB/s |

## but seek times will change slowly.

- Bandwidth scales as square root of capacity.

Source: http://blocksandfiles.com/article/4501

# Fractal Trees Enable Compact Systems

## B-trees require capacity, bandwidth, and random I/O

- B-tree based systems achieve large random I/O rates by using more spindles and lower capacity disks.

## Fractal Trees require only capacity & bandwidth

- Fractal Trees enable the use of high-capacity disks.

# Fractal Trees Enable Big Disks

**B-trees require capacity, bandwidth, and seeks.**

**Fractal trees require only capacity and bandwidth.**

**Today, for a 50TB database,**

- Fractal tree with 25 2TB disks gives 500K ins/s.
- B-tree with 25 2TB disks gives 2.5K ins/s.
- B-tree with 500 100GB disks gives 50K ins/s but costs $, racks, and power.

**In 2017, for a 1500TB database:**

- Fractal tree with 25 67TB disks gives 2500K ins/s.
- B-tree with 25 67TB disks gives 2.5K ins/s.

**B-trees need spindles, and spindle density increases slowly.**

*Tokutek*™

# Using Big Disks Also Saves Energy

## Power consumption of disks

- Enterprise 80 to 160 GB disk runs at 4W (idle power).
- Enterprise 1-2 TB disk runs at 8W (idle power).

## Data centers/server farms use 80-160 GB disks

- Use many small-capacity disks, not large ones.

## Using large disks may save factor >10 in Storage Costs

- Other considerations modify this factor
  ▸ e.g., CPUs necessary to drive disks, scale-out infrastructure, cooling, etc.
  ▸ Metric: e.g., Watts/MB versus Inserts/Joule