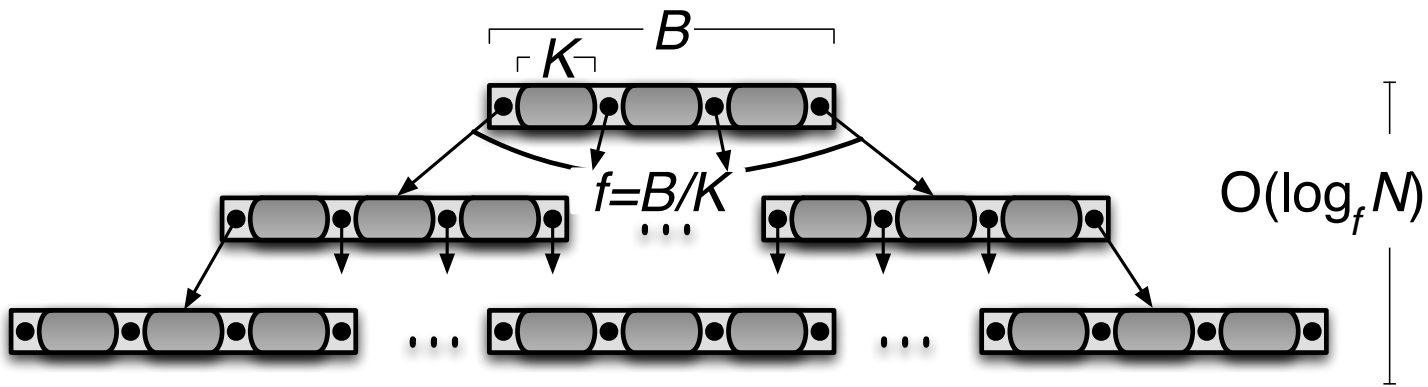# Performance Guarantees for B-trees with Different-Size Atomic Keys

**Michael A. Bender**
**Stony Brook**
**Tokutek, Inc.**
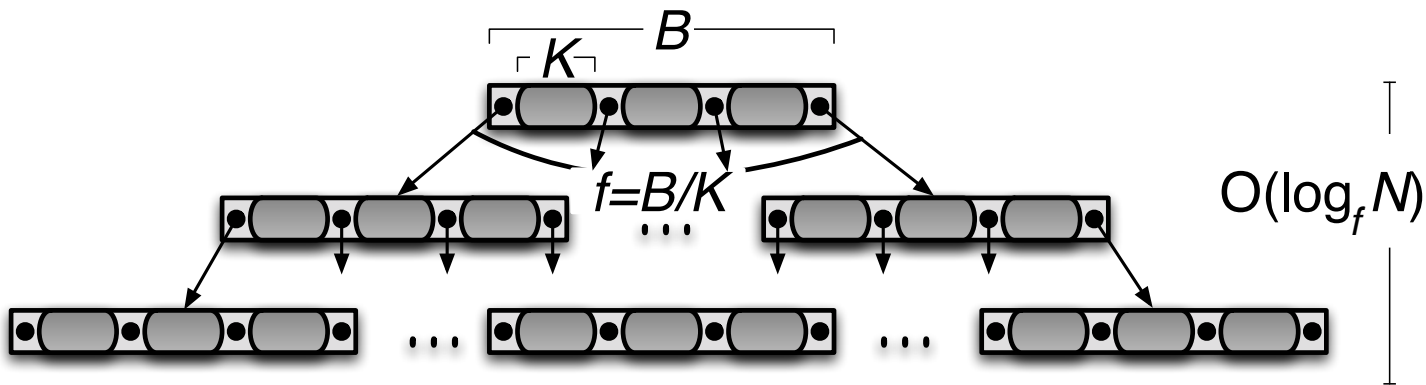
**Haodong Hu**
**Microsoft**

**Bradley C. Kuszmaul**
**MIT**
**Tokutek, Inc**

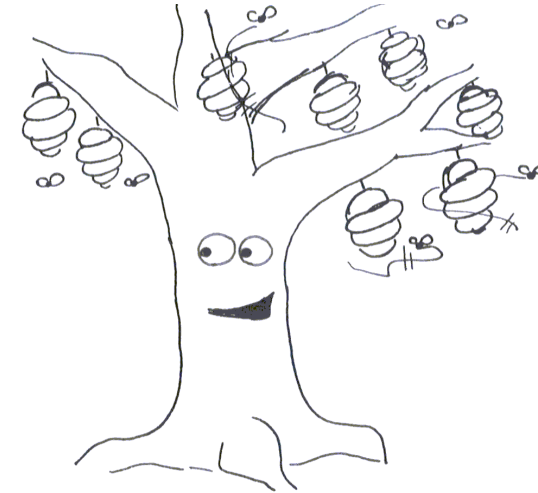**In B-trees in textbooks, all keys have the same size.**



$f = B/K$
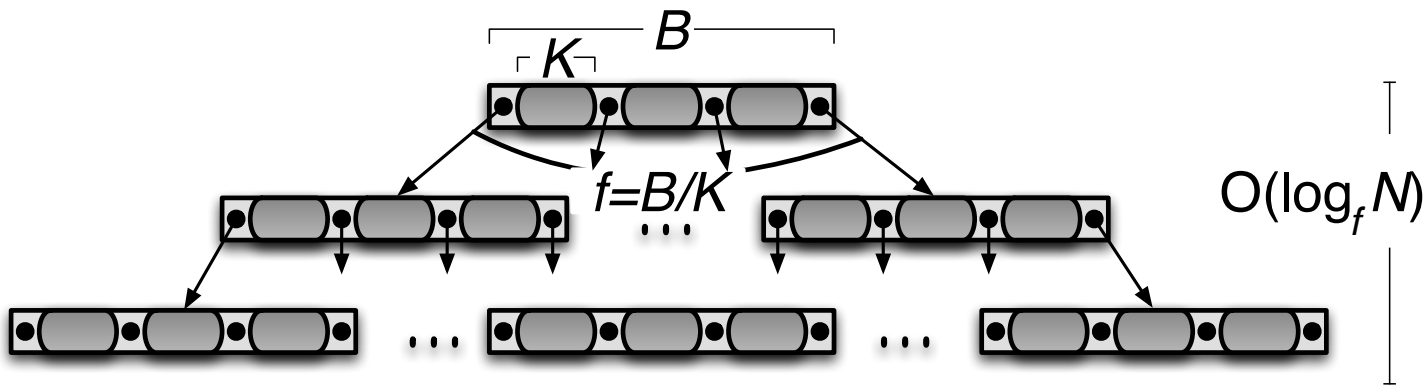
$O(\log_f N)$

Bender--B-trees with different-sized keys

**In B-trees in textbooks, all keys have the same size.**



$O(\log_f N)$

$f = B/K$

**Production B-trees support different-size keys, but with no nontrivial performance guarantees.**

**In B-trees in textbooks, all keys have the same size.**
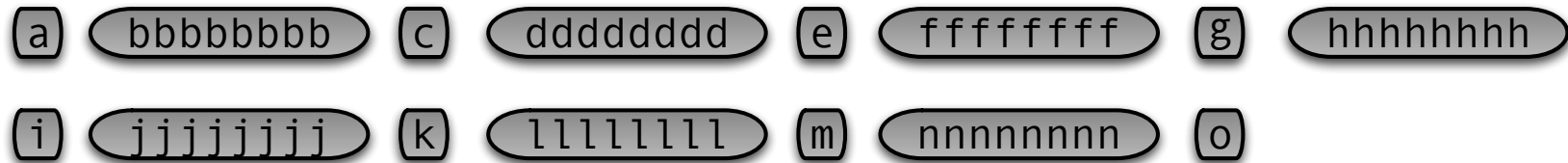


$$O(\log_f N)$$

$f = B/K$

**Production B-trees support different-size keys, but with no nontrivial performance guarantees.**

*This talk: give provably good guarantees in an only slightly modified B-tree.*

Bender--B-trees with different-sized keys

a  bbbbbbbb  c  dddddddd  e  ffffffff  g  hhhhhhhh

i  jjjjjjjj  k  llllllll  m  nnnnnnnn  o

**When the length-8 keys are pivots (and the block size is 8), the tree height is 4:**

```
                        hhhhhhhh

        dddddddd                      llllllll

   bbbbbbbb      ffffffff        jjjjjjjj      nnnnnnnn

   a       c    e       g        i       k    m        o
```

Bender--B-trees with different-sized keys

a  bbbbbbbb  c  dddddddd  e  ffffffff  g  hhhhhhhh

i  jjjjjjjj  k  llllllll  m  nnnnnnnn  o

**When the length-1 keys are pivots (and the block size is 8), the tree height is 2:**

a c e g i k m o

bbbbbbbb  dddddddd  ffffffff  hhhhhhhh  jjjjjjjj  llllllll  nnnnnnnn

*Choice of pivot affects the B-tree performance.*

Bender--B-trees with different-sized keys

**Cannot compare first byte of** `bbbbbbbb` **with** `c`.

**Cannot compare first byte of** `bbbbbbbb` **with** `c`.

**Only the comparison function understands the keys.**

**Keys are opaque. Need to send entire key to comparison function and store entire key in node.**

**Cannot compare first byte of** `bbbbbbbb` **with** `c`.

**Only the comparison function understands the keys.**

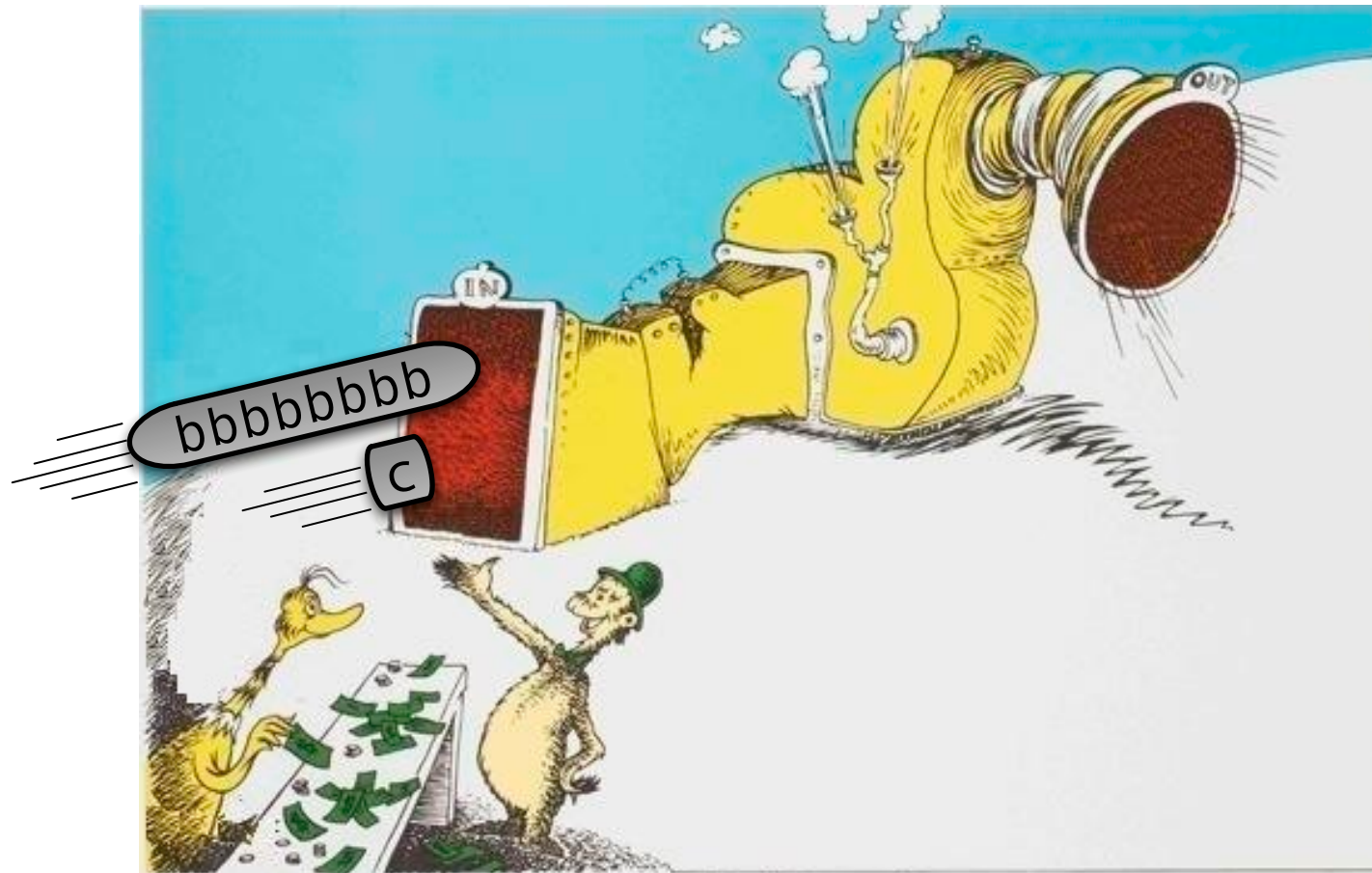**Keys are opaque. Need to send entire key to comparison function and store entire key in node.**

**Cannot compare first byte of** `bbbbbbbb` **with** `c`.
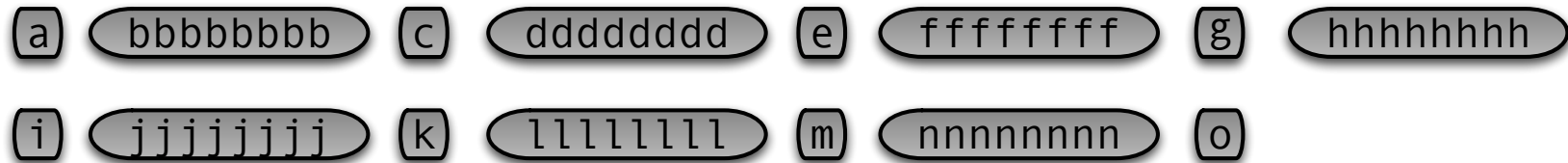
**Only the comparison function understands the keys.**

**Keys are opaque. Need to send entire key to comparison function and store entire key in node.**

**String B-tree techniques don't work** [Ferragina, Grossi 98].

**Front compression (prefix compression) doesn't work** [Bayer, Unterauer 77] [Wagner 73].

STONY BROOK

Bender--B-trees with different-sized keys

*Tokutek* ™

a bbbbbbbb c dddddddd e ffffffff g hhhhhhhh

i jjjjjjjj k llllllll m nnnnnnnn o

*earliest* vs. *latest*: order determined by comparison function.

*shortest* versus *longest*: how many bytes to store key

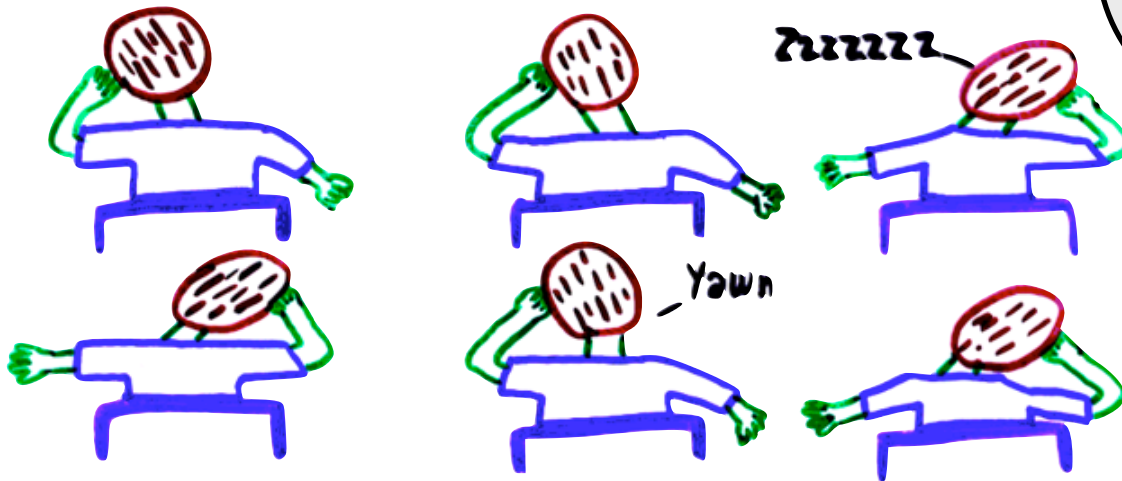(Words to avoid: *small*, *large*, *minimum*.)

# So if I say....

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

- Two-levels of memory.
- Two parameters: block-size $B$, memory-size $M$.

## Performance metric:

- Minimize # of block transfers

Memory

B

B

Disk

**_Example:_ _N_ keys with average size <2.**

- _N/B_ keys with size _B_ and _N-N/B_ keys with size 1.



1    _B_    _B_-1

# Choice of Pivot Matters For Variable-Size Keys

**_Example:_ _N_ keys with average size <2.**

- *N/B* keys with size *B* and *N-N/B* keys with size 1.



**Size 1 keys as pivots: optimal.**



$O(\log_B N)$

Bender--B-trees with different-sized keys

# Choice of Pivot Matters For Variable-Size Keys

**Example: *N* keys with average size <2.**

- *N/B* keys with size *B* and *N-N/B* keys with size 1.



**Size *B* keys as pivots: $O(\log B)$ factor worse.**



$O(\log_2 N)$

Bender--B-trees with different-sized keys

Comparison cost:
# of transfers to bring keys into memory.

**Let $K$ be the *average* key size.**

**Goal: $O(\log_{B/K} N)$ memory transfers per operation.**

- Generalizes what happens if keys all have the same size $K$.

Bender--B-trees with different-sized keys

**Let *K* be the *average* key size.**

**Goal: $O(\log_{B/K} N)$ memory transfers per operation.**

- Generalizes what happens if keys all have the same size *K*.

**Unfortunately, we cannot get this for worst-case searches, but we'll get it in expectation.**

Bender--B-trees with different-sized keys

**Example: *N* keys with average size *K<2*.**

- *N/B* keys with size *B* and *N-N/B* keys with size 1.



$B$         1

search: $\Theta(\log_2 N/B)$      search: $\Theta(\log_B N)$

**Example: *N* keys with average size *K*<2.**

- *N/B* keys with size *B* and *N-N/B* keys with size 1.



$B$

$1$

search: $\Theta(\log_2 N/B)$      search: $\Theta(\log_B N)$

$B$

$O(\log_B N)$

$O(\log_2 N/B)$

Bender--B-trees with different-sized keys

**Example: $N$ keys with average size $K<2$.**

- *$N/B$ keys with size $B$ and $N-N/B$ keys with size 1.*



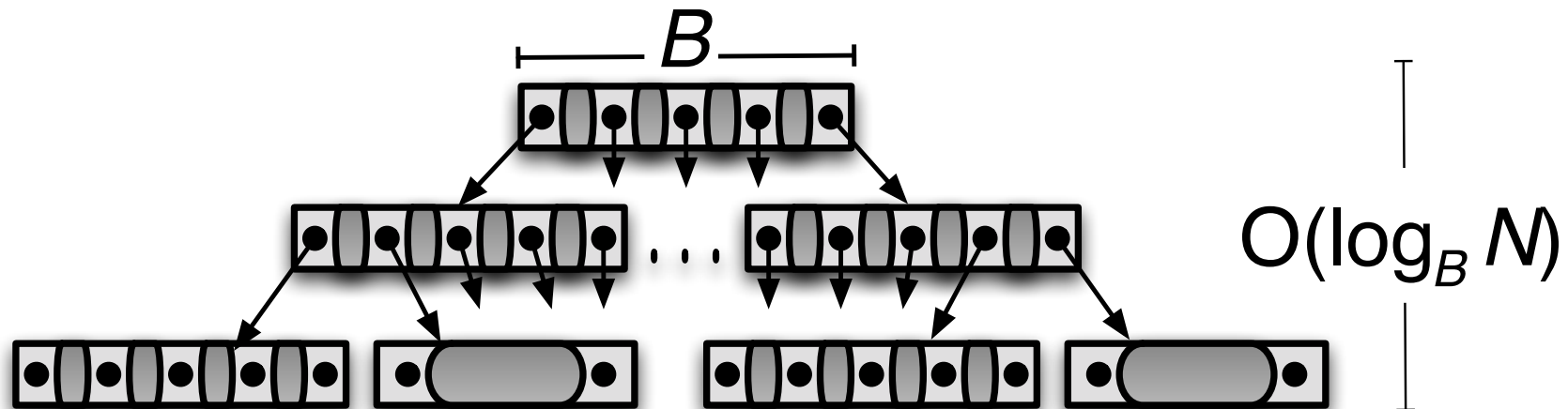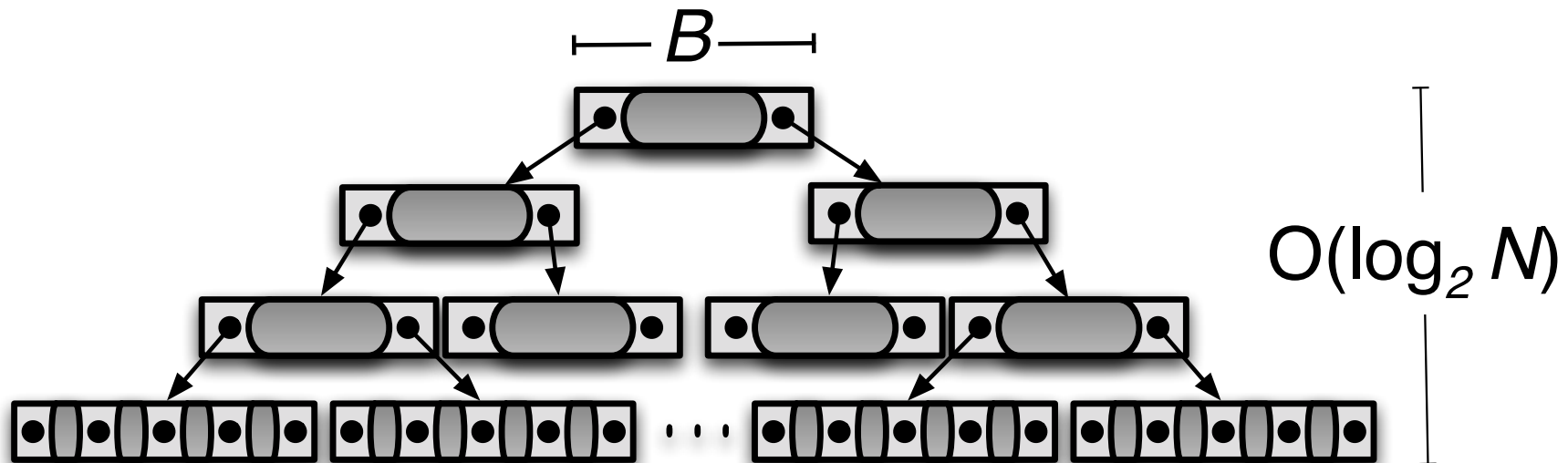search: $\Theta(\log_2 N/B)$     search: $\Theta(\log_B N)$



$O(\log_B N)$

$O(\log_2 N/B)$

But we are ok in expectation: $(1\text{-}1/B)\log_B N + (1/B)\log_2 N = O(\log_B N)$.

# Why We Cannot Attain Good Worst-Case Bounds

**Example: *N* keys with average size *K*<2.**

- *N/B* keys with size *B* and *N-N/B* keys with size 1.



$B$                          $1$

search: $\Theta(\log_2 N/B)$         search: $\Theta(\log_B N)$

**Related work: how to optimize B-tree height**

[Vaishnavi, Kriegel, Wood 80] [Gotleib 81] [Huang, Vishwanathan 90] [Becker 94]

- static (no inserts/deletes)
- DP-based
- (far from our target guarantee)

# Why We Cannot Attain Good Worst-Case Bounds

**Extreme example: _N_ keys with average size _K_<2.**

- 1 key with size _M_ and _N-1_ keys with size 1.

$M$      1

Read: $\Theta(1+M/B)$      search: $\Theta(\log_B N)$

But we are ok in expectation: $(1-1/N)\log_B N + (1/N)\,M/B = O(\log_B N)$.

**Extreme example: *N* keys with average size *K*<2.**

- 1 key with size *M* and *N-1* keys with size 1.



$M$      $1$

Read: $\Theta(1+M/B)$      search: $\Theta(\log_B N)$

But we are ok in expectation: $(1-1/N)\log_B N + (1/N)\,M/B = O(\log_B N)$.

**These two examples have different flavors.**

- LB for first example is based on tree structure.
- LB for for second example is based on reading the key.
- Both motivate why we consider expectation.

**Static atomic-key B-tree (only searches)**

- Expected leaf search cost: $O(\lceil K/B \rceil \log_{1+\lceil B/K \rceil} N)$
- Linear construction cost for sorted data: $O(NK/B)$

## Static atomic-key B-tree (only searches)

- Expected leaf search cost: $O(\lceil K/B \rceil \log_{1+\lceil B/K \rceil} N)$ ← Captures $K = O(B)$ and $K \leq \Omega(B)$

- Linear construction cost for sorted data: $O(NK/B)$

## Static atomic-key B-tree (only searches)

- Expected leaf search cost: $O(\lceil K/B \rceil \log_{1+\lceil B/K \rceil} N)$ ← Captures $K=O(B)$ and $K \leq \Omega(B)$

- Linear construction cost for sorted data: $O(NK/B)$ ← Scan bound since total length= $NK$

STONY BROOK

Bender--B-trees with different-sized keys

*Tokutek*

## Static atomic-key B-tree (only searches)

- Expected leaf search cost: $O(\lceil K/B \rceil \log_{1+\lceil B/K \rceil} N)$ ← Captures $K=O(B)$ and $K \leq \Omega(B)$

- Linear construction cost for sorted data: $O(NK/B)$ ← Scan bound since total length= $NK$

## Dynamic atomic-key B-tree

- Expected leaf search cost : $O(\lceil K/B \rceil \log_{1+\lceil B/K \rceil} N)$

- Cost to insert/delete/search for key $L$ of *random* rank (amort): $O(\lceil K/B \rceil \log_{1+\lceil B/K \rceil} N + |L|/B)$

- Cost to insert/delete/search for key of *arbitrary* rank: ***modification cost is dominated by search cost.***

Bender--B-trees with different-sized keys

## Static atomic-key B-tree (only searches)

- Expected leaf search cost: $O(\lceil K/B \rceil \log_{1+\lceil B/K \rceil} N)$ <span style="color:darkred">← Captures $K=O(B)$ and $K \leq \Omega(B)$</span>

- Linear construction cost for sorted data: $O(NK/B)$ <span style="color:darkred">← Scan bound since total length= $NK$</span>

## Dynamic atomic-key B-tree

- Expected leaf search cost : $O(\lceil K/B \rceil \log_{1+\lceil B/K \rceil} N)$

- Cost to insert/delete/search for key $L$ of *random* rank (amort):
  $O(\lceil K/B \rceil \log_{1+\lceil B/K \rceil} N + |L|/B)$ <span style="color:darkred">← $O(|L|/B)$ is cost to read $L$ into memory</span>

- Cost to insert/delete/search for key of *arbitrary* rank:
  ***modification cost is dominated by search cost.***

## Static atomic-key B-tree (only searches)

- Expected leaf search cost: $O(\lceil K/B \rceil \log_{1+\lceil B/K \rceil} N)$ ← Captures $K=O(B)$ and $K \leq \Omega(B)$

- Linear construction cost for sorted data: $O(NK/B)$ ← Scan bound since total length= $NK$

## Dynamic atomic-key B-tree

- Expected leaf search cost : $O(\lceil K/B \rceil \log_{1+\lceil B/K \rceil} N)$

- Cost to insert/delete/search for key $L$ of *random* rank (amort):
$O(\lceil K/B \rceil \log_{1+\lceil B/K \rceil} N + |L|/B)$ ← $O(|L|/B)$ is cost to read $L$ into memory

- Cost to insert/delete/search for key of *arbitrary* rank:
***modification cost is dominated by search cost.*** ← important

Bender--B-trees with different-sized keys

## Static atomic-key B-tree (only searches)

- Expected leaf search cost: $O(\lceil K/B \rceil \log_{1+\lceil B/K \rceil} N)$ ← Captures *K*=O(*B*) and *K*≤ Ω(*B*)

- Linear construction cost for sorted data: $O(NK/B)$ ← Scan bound since total length= *NK*

## Dynamic atomic-key B-tree

- Expected leaf search cost : $O(\lceil K/B \rceil \log_{1+\lceil B/K \rceil} N)$

- Cost to insert/delete/search for key *L* of *random* rank (amort):
$O(\lceil K/B \rceil \log_{1+\lceil B/K \rceil} N + |L|/B)$ ← O(|L|/B) is cost to read *L* into memory

- Cost to insert/delete/search for key of *arbitrary* rank:
***modification cost is dominated by search cost.*** ← important

## Optimal static atomic-key B-tree:

- $O(BN^3)$ operations
- Applies even for nonuniform search probabilities

## Static atomic-key B-tree (only searches)

- Expected leaf search cost: $O(\lceil K/B \rceil \log_{1+\lceil B/K \rceil} N)$ ⟵ *Captures K=O(B) and K≤ Ω(B)*

- Linear construction cost for sorted data: $O(NK/B)$ ⟵ *Scan bound since total length= NK*

## Dynamic atomic-key B-tree

- Expected leaf search cost : $O(\lceil K/B \rceil \log_{1+\lceil B/K \rceil} N)$

- Cost to insert/delete/search for key *L* of *random* rank (amort): $O(\lceil K/B \rceil \log_{1+\lceil B/K \rceil} N + |L|/B)$ ⟵ *O(|L|/B) is cost to read L into memory*

- Cost to insert/delete/search for key of *arbitrary* rank:
  ***modification cost is dominated by search cost.*** ⟵ *important*

## Optimal static atomic-key B-tree:

- $O(BN^3)$ operations ⟵ *RAM not external memory. (Won't discuss in talk.)*

- Applies even for nonuniform search probabilities

## Static atomic-key B-tree (only searches)  ←To discuss next

- Expected leaf search cost: $O(\lceil K/B \rceil \log_{1+\lceil B/K \rceil} N)$

- Linear construction cost for sorted data: $O(NK/B)$

## Dynamic atomic-key B-tree

- Expected leaf search cost : $O(\lceil K/B \rceil \log_{1+\lceil B/K \rceil} N)$

- Cost to insert/delete/search for key *L* of *random* rank (amort): $O(\lceil K/B \rceil \log_{1+\lceil B/K \rceil} N + |L|/B)$

- Cost to insert/delete/search for key of *arbitrary* rank: ***modification cost is dominated by search cost.***

## Optimal static atomic-key B-tree:

- $O(BN^3)$ operations

- Applies even for nonuniform search probabilities

**Greedy construction algorithm**

- Greedily select pivot elements for the root node
- Proceed recursively on all subtrees of the root.
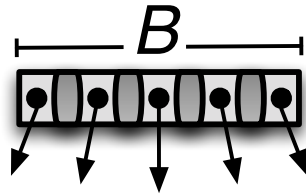
**Intuition**

- Pick small keys in root to maximize fanout.
- Pick evenly distributed keys to reduce the search space.

**To prove**
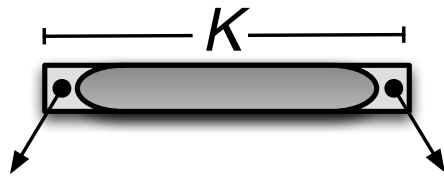
- Root has a good structure.
- Recursive substructures achieve good performance, even though subtrees may have different average key sizes.

Bender--B-trees with different-sized keys

# Root Structure of Static Atomic Key B-tree

**Case 1: $K=O(B)$. Root has size $O(B)$ and fanout $\Theta(B/K)$.**



**Case 2: $K= \Omega(B)$. Root has size $O(K)$ and fanout 2.**

# Root Structure of Static Atomic Key B-tree

**Case 1: $K$=O($B$). Root has size O($B$) and fanout $\Theta$($B$/$K$).**



*Overall search cost:* $O(\lceil K/B \rceil \log_{1+\lceil B/K \rceil} N)$

**Case 2: $K$= $\Omega$($B$). Root has size O($K$) and fanout 2.**



*Overall search cost:* $O(\lceil K/B \rceil \log_{1+\lceil B/K \rceil} N)$

Bender--B-trees with different-sized keys

# Root Structure of Static Atomic Key B-tree

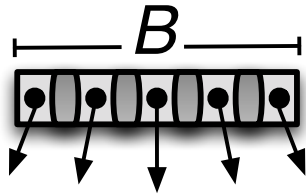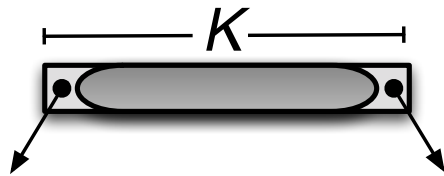**Case 1: $K = O(B)$. Root has size $O(B)$ and fanout $\Theta(B/K)$.**



*Overall search cost:*

$$O(\underbrace{\lceil K/B \rceil}_{=1} \log_{\underbrace{1 + \lceil B/K \rceil}_{\approx \Theta(B/K)}} N)$$

> cost of reading root

> ave fanout

**Case 2: $K = \Omega(B)$. Root has size $O(K)$ and fanout $2$.**



*Overall search cost:*

$$O(\underbrace{\lceil K/B \rceil}_{\approx K/B} \log_{\underbrace{1 + \lceil B/K \rceil}_{=2}} N)$$

> cost of reading root

> ave fanout

Bender--B-trees with different-sized keys

*Useful Lemma*: **Consider $N$ #s whose average is $K$. Divide into $f$ groups of equal cardinality (each has $N/f$). Take the min in each group (say $K_i$ is min of group $i$). Then average of these minima is at most the overall average $K$ (i.e., $(K_1+K_2+...+K_f)/f \le K$).**

**Ex.  4  4  3  4 | 1  2  2  2 | 3  2  5  4 | 3  1  4  4**

**Ave $K$=3.**

STONY
BROOK

*Tokutek*™

*Useful Lemma*: **Consider $N$ #s whose average is $K$. Divide into $f$ groups of equal cardinality (each has $N/f$). Take the min in each group (say $K_i$ is min of group $i$). Then average of these minima is at most the overall average $K$ (i.e., $(K_1+K_2+...+K_f)/f \leq K$).**

**Ex.  4  4  ③  4 | ①  2  2  2 | 3  ②  5  4 | 3  ①  4  4**

**Ave $K$=3. Ave of mins =1.75.**

# Root Structure of Static Atomic Key B-tree

*Useful Lemma*: **Consider *N* #s whose average is *K*. Divide into *f* groups of equal cardinality (each has *N/f*). Take the min in each group (say $K_i$ is min of group *i*). Then average of these minima is at most the overall average *K* (i.e., $(K_1+K_2+...+K_f)/f \le K$).**

**Ex.**  4  4  ③  4  ①  2  2  2  3  ②  5  4  3  ①  4  4

**Ave *K*=3. Ave of mins =1.75.**

*Note: only true because groups have the same size.*

*Enough structure to bound the size and fanout of the root.*

1. **Divide keys into** $f = \max\left\{3, \left\lfloor \frac{B}{K} \right\rfloor\right\}$ **equal-size groups.**

2. **Pick shortest key in each group.**

3. **Store these keys in root (except 1st & last groups).**

aaaa  bbbb  ccc  dddd  e  ff  gg  hh  iii  jj  kkkk  llll  mmm  n  oooo  pppp

**Ex:** *B* = 12, *K* = 3.  **So** *f*=4.

Bender--B-trees with different-sized keys

**1. Divide keys into** $f = \max \left\{ 3, \left\lfloor \frac{B}{K} \right\rfloor \right\}$ **equal-size groups.**

2. Pick shortest key in each group.

3. Store these keys in root (except 1st & last groups).

aaaa bbbb ccc dddd e ff gg hh iii jj kkkk llll mmm n oooo pppp

**Ex:** *B* = 12, *K* = 3. **So** *f*=4.

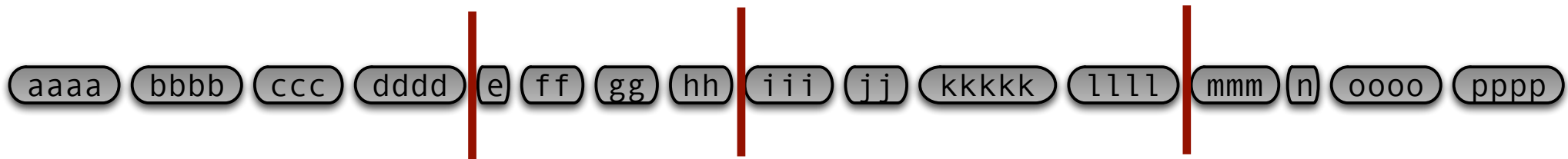Bender--B-trees with different-sized keys

1. **Divide keys into** $f = \max\left\{3, \left\lfloor \frac{B}{K} \right\rfloor\right\}$ **equal-size groups.**

2. **Pick shortest key in each group.**

3. Store these keys in root (except 1st & last groups).



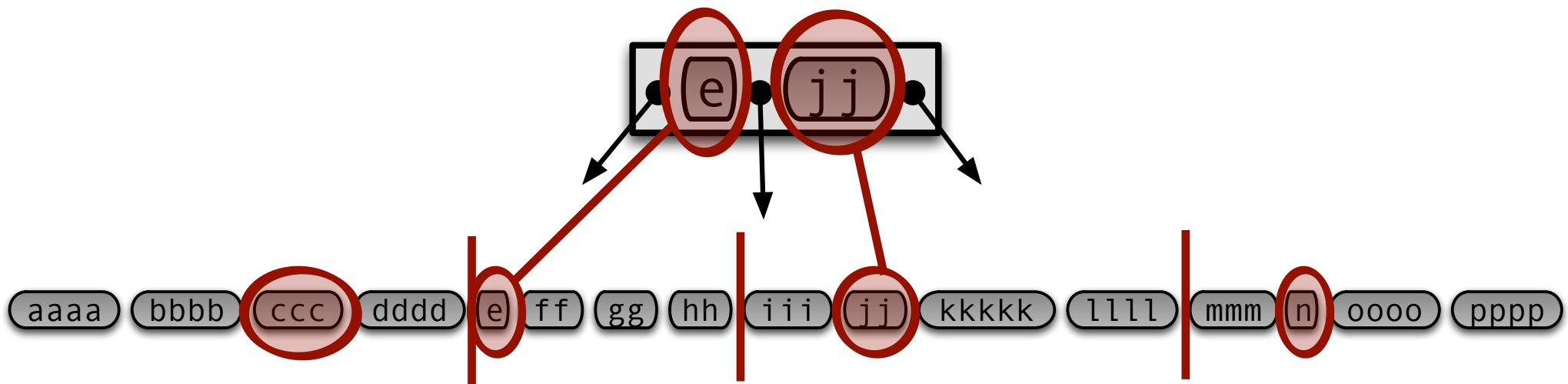**Ex:** *B* = 12, *K* = 3.  **So** *f*=4.

1. **Divide keys into** $f = \max \left\{ 3, \left\lfloor \frac{B}{K} \right\rfloor \right\}$ **equal-size groups.**

2. **Pick shortest key in each group.**

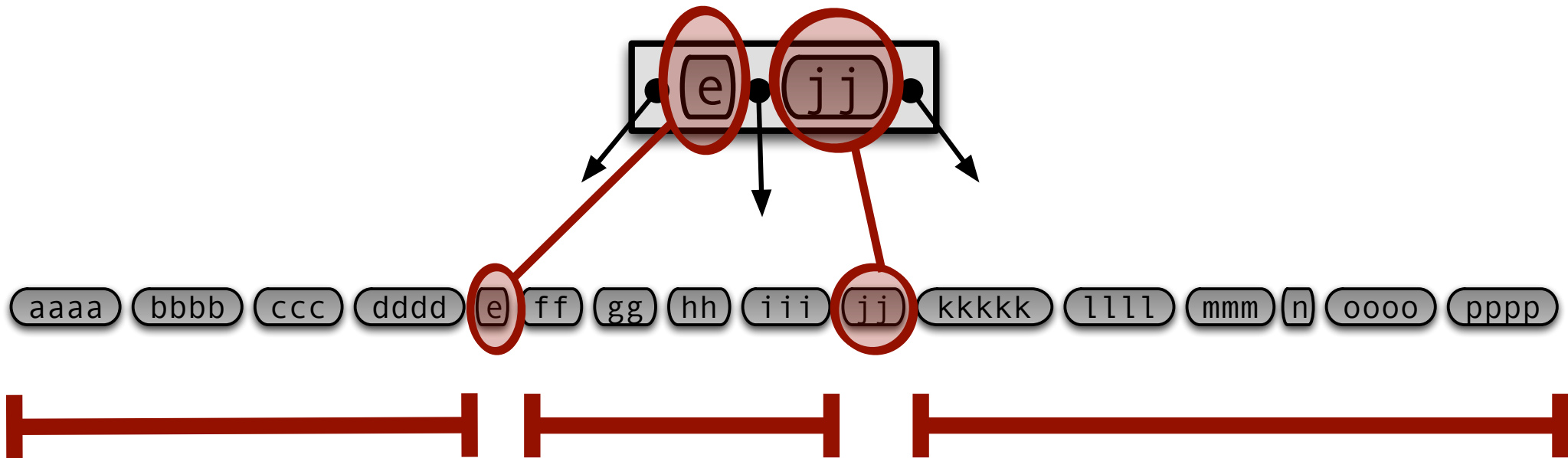3. **Store these keys in root (except 1st & last groups).**



**Ex:** *B* = 12, *K* = 3. **So** *f*=4.

**Proceed recursively in each subtree.**
**(Different value of $K$ (and thus $f = \max\left\{3, \left\lfloor \frac{B}{K} \right\rfloor\right\}$)**
**in each subtree.**

## Static atomic-key B-tree (only searches)

- Expected leaf search cost: $O(\lceil K/B \rceil \log_{1+\lceil B/K \rceil} N)$

- Linear construction cost for sorted data: $O(NK/B)$

## Dynamic atomic-key B-tree    ⟵    To discuss next

- Expected leaf search cost : $O(\lceil K/B \rceil \log_{1+\lceil B/K \rceil} N)$

- Cost to insert/delete/search for key *L* of *random* rank (amort):
  $O(\lceil K/B \rceil \log_{1+\lceil B/K \rceil} N + |L|/B)$

- Cost to insert/delete/search for key of *arbitrary* rank:
  ***modification cost is dominated by search cost.***

## Optimal static atomic-key B-tree:

- $O(BN^3)$ operations

- Applies even for nonuniform search probabilities

***One idea:*** **Groups need not be of equal cardinality. Within constant factors is good enough. We don't need the shortest key as a pivot. We can choose a key whose length is at most twice the average in that group.**

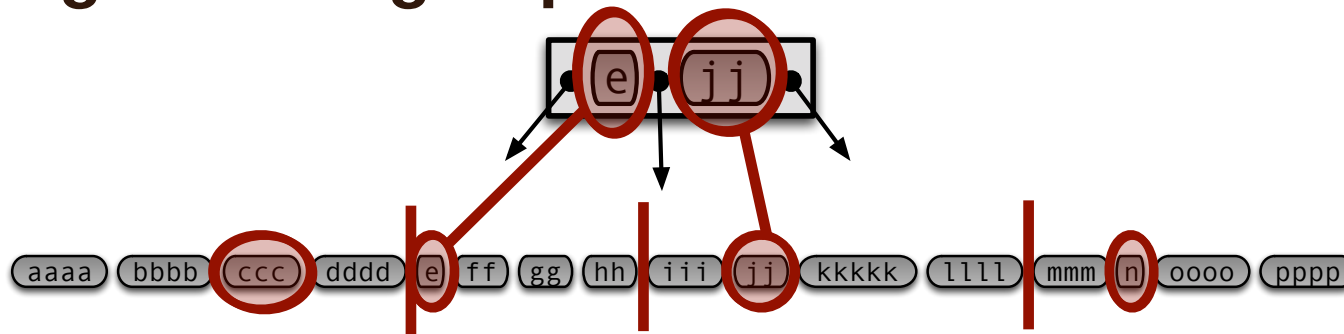***One idea:*** **Groups need not be of equal cardinality. Within constant factors is good enough. We don't need the shortest key as a pivot. We can choose a key whose length is at most twice the average in that group.**



*Thus, > half the keys in a group could be a pivot.*

*The shortest key can remain as pivot even if the group grows or shrinks by a constant factor.*

*Structure isn't brittle.*

**Second idea: Insert elements directly into leaves. Rebuild entire subtrees whether there have been "too many" inserts/deletes. (Don't bother splitting and merging. )**

$$\text{amortized update cost} = \frac{\text{rebuild cost}}{\text{\# updates between rebuilds}}$$

**Problem: standard technique chokes because value of *K* changes over time. (Value of *K* during rebuild is different from value of *K* at actual insert/delete.) Can be fixed. (Ask after talk.)**

**People want performance guarantees**

- I have a startup Tokutek. One of the things customers like most about our product TokuDB is its predictability.

Bender--B-trees with different-sized keys

ximum

**People want performance guarantees**
- I have a startup Tokutek. One of the things customers like most about our product TokuDB is its predictability.

**People want performance guarantees for B-trees**
- The manual for Oracle Berkeley DB manual claims that BDB runs in $O(\log_{B/K} N)$ transfers.
- As our results show, this folk theorem is incorrect.
- But the claim helps motivate the guarantees we achieve.

Bender--B-trees with different-sized keys

**People want performance guarantees**

- I have a startup Tokutek. One of the things customers like most about our product TokuDB is its predictability.

**People want performance guarantees for B-trees**

- The manual for Oracle Berkeley DB manual claims that BDB runs in $O(\log_{B/K} N)$ transfers.
- As our results show, this folk theorem is incorrect.
- But the claim helps motivate the guarantees we achieve.

**Will our theoretical guarantees have practical value?**

- Maybe B-trees empirically perform predictably enough.

STONY
BROOK

*Tokutek* ™

**People want performance guarantees**
- I have a startup Tokutek. One of the things customers like most about our product TokuDB is its predictability.

**People want performance guarantees for B-trees**
- The manual for Oracle Berkeley DB manual claims that BDB runs in $O(\log_{B/K} N)$ transfers.
- As our results show, this folk theorem is incorrect.
- But the claim helps motivate the guarantees we achieve.

**Will our theoretical guarantees have practical value?**
- Maybe B-trees empirically perform predictably enough.
- Maybe B-trees are so unpredictable already (e.g., because of memory cliffs) that our guarantees are second-order effects.

**People want performance guarantees**

- I have a startup Tokutek. One of the things customers like most about our product TokuDB is its predictability.

**People want performance guarantees for B-trees**

- The manual for Oracle Berkeley DB manual claims that BDB runs in $O(\log_{B/K} N)$ transfers.
- As our results show, this folk theorem is incorrect.
- But the claim helps motivate the guarantees we achieve.

**Will our theoretical guarantees have practical value?**

- Maybe B-trees empirically perform predictably enough.
- Maybe B-trees are so unpredictable already (e.g., because of memory cliffs) that our guarantees are second-order effects.
- Lots to explore.....

Bender--B-trees with different-sized keys

## People want performance guarantees

- I have a startup Tokutek. One of the things customers like most about our product TokuDB is its predictability.

## People want performance guarantees for B-trees

- The manual for Oracle Berkeley DB manual claims that BDB runs in $O(\log_{B/K}N)$ transfers.
- As our results show, this folk theorem is incorrect.
- But the claim helps motivate the guarantees we achieve.

## Will our theoretical guarantees have practical value?

- Maybe B-trees empirically perform predictably enough.
- Maybe B-trees are so unpredictable already (e.g., because of memory cliffs) that our guarantees are second-order effects.
- Lots to explore.....

STONY BROOK

Tokutek™