Leif Walsh

Two Sigma Investments, LLC. leif.walsh@gmail.com @leifwalsh

April 16, 2015

	10/	2		-
LEIL	**		ъ	

-

• • • • • • • • • • • • •

Background

(日)

Data structures:

<ロト < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 < つ < つ </p>

Data structures:



Data structures:

- Provide retrieval of data.
 - Lookup(Key)
 - Pred(Key)
 - Succ(Key)



Data structures:

- Provide retrieval of data.
 - Lookup(Key)
 - Pred(Key)
 - Succ(Key)
- *Dynamic* data structures let you change the data.
 - Insert(Key, Value)
 - Delete(Key)



Image: A matrix and a matrix

-

[Aggarwal & Vitter '88]

DAM model

- Problem size *N*.
- Memory size *M*.

▲ロ▶▲圖▶▲≣▶▲≣▶ = 更 のQ@

[Aggarwal & Vitter '88]

DAM model

- Problem size *N*.
- Memory size *M*.
- Transfer data to/from memory in *blocks* of size *B*.



[Aggarwal & Vitter '88]

DAM model

- Problem size N.
- Memory size *M*.
- Transfer data to/from memory in *blocks* of size *B*.

Efficiency of operations is measured as the *number of block transfers*, a.k.a. IOPS.



A B-tree is an *external memory data structure*:

▲ロ▶▲圖▶▲≣▶▲≣▶ ■ のQの

A B-tree is an *external memory data structure*:



A B-tree is an *external memory data structure*:



- Balanced search tree.
- Fanout of B (block size / key size).

A B-tree is an *external memory data structure*:





- Balanced search tree.
- Fanout of *B* (block size / key size).

- Internal nodes < *M*.
- Leaf nodes > M.

A B-tree is an *external memory data structure*:





- Balanced search tree.
- Fanout of B (block size / key size).

- Internal nodes < *M*.
- Leaf nodes > M.
- Search: $O(\log_B N)$ I/Os
- Insert: $O(\log_B N)$ I/Os

Image: A matrix

[Brodal & Fagerberg '03]



[Brodal & Fagerberg '03]



Э.

OLAP

<ロト < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □

OLAP: Online Analytical Processing

April 16, 2015

9/33

OLAP: Online **Analytical** Processing Key idea: Analyze data collected in the past.

4 D > 4 B

OLAP: Online **Analytical** Processing Key idea: Analyze data collected in the past.

B-tree inserts are slow, but...logging and sorting are fast.

Image: A matrix

OLAP: Online **Analytical** Processing

Key idea: Analyze data collected in the past.

B-tree inserts are slow, but...logging and sorting are fast.

Plan: Log new data unsorted, then build indexes in large batches.

Write-optimization technique #1: OLAP

Merge sort:



◆□▶ ◆□▶ ◆三▶ ◆三▶ ● □ ● ●



<ロト < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □



Merge sort cost in DAM model is:

- Cost to scan through all the data once.
- Multiplied by the # of levels in the merge tree.



Merge sort cost in DAM model is:

- Cost to scan through all the data once. N/B
- Multiplied by the # of levels in the merge tree.



Merge sort cost in DAM model is:

• Cost to scan through all the data once.

N/B

• Multiplied by the # of levels in the merge tree.

 $\log_{M/B} N/B$



Merge sort cost in DAM model is:

- Cost to scan through all the data once. N/B
- Multiplied by the # of levels in the merge tree.

 $\log_{M/B} N/B$

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

$$O\left(N\log_B \frac{N}{M}\right)$$

Merge sort:

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

$$O\left(N\log_B\frac{N}{M}\right)$$

Merge sort:

$$O\left(rac{N}{B}\log_{M/B}rac{N}{B}
ight) pprox rac{2N}{B}$$

Typically, M/B is large, so only two passes are needed to sort.

-

$$O\left(N\log_B \frac{N}{M}\right) \geq N$$

Merge sort:

$$O\left(\frac{N}{B}\log_{M/B}\frac{N}{B}
ight) \approx \frac{2N}{B}$$

Typically, M/B is large, so only two passes are needed to sort.

Intuition: Each insert into a B-tree costs ~1 seek, while sorting is close to disk bandwidth.

- 1	Leif	Wa	lsh

(assuming 100-1000 byte elements)

$$O\left(N\log_{B}\frac{N}{M}\right) \geq N \approx 10 - 100$$
kB/s = 100 elements/s

Merge sort:

$$O\left(\frac{N}{B}\log_{M/B}\frac{N}{B}
ight) \approx \frac{2N}{B}$$

Typically, M/B is large, so only two passes are needed to sort.

Intuition: Each insert into a B-tree costs ~1 seek, while sorting is close to disk bandwidth.

Fractal Troop		۸	16 201E	12/3
Flacial flees		AL	JHL 10. ZU15	·

人口人 人間人 人民人 人民人

Leif Wals

(assuming 100-1000 byte elements)

$$O\left(N\log_{B}\frac{N}{M}\right) \geq N \approx 10 - 100$$
kB/s = 100 elements/s

Merge sort:

$$O\left(\frac{N}{B}\log_{M/B}\frac{N}{B}\right) \approx \frac{2N}{B} \approx 50$$
 MB/s = 50k - 500k elements/s

Typically, M/B is large, so only two passes are needed to sort.

Intuition: Each insert into a B-tree costs ~1 seek, while sorting is close to disk bandwidth.

Fractal Trees	April 16, 2015	12/33

<ロト < 回 ト < 回 ト < 回 ト < 回 ト < 回 ト < 回 ト < 回 への</p>

• Log new data unindexed until you accumulate a lot of it (~10% of the data set).

- Log new data unindexed until you accumulate a lot of it (~10% of the data set).
- Sort the new data.

- Log new data unindexed until you accumulate a lot of it (~10% of the data set).
- Sort the new data.
- Use a merge pass through existing indexes to incorporate new data.
So, how does OLAP work?

- Log new data unindexed until you accumulate a lot of it (~10% of the data set).
- Sort the new data.
- Use a merge pass through existing indexes to incorporate new data.
- Use indexes to do analytics.

So, how does OLAP work?

- Log new data unindexed until you accumulate a lot of it (~10% of the data set).
- Sort the new data.
- Use a merge pass through existing indexes to incorporate new data.
- Use indexes to do analytics.

Moral: OLAP techniques can handle high insertion volume, but query results are *delayed*.

So, how does OLAP work?

- Log new data unindexed until you accumulate a lot of it (~10% of the data set).
- Sort the new data.
- Use a merge pass through existing indexes to incorporate new data.
- Use indexes to do analytics.

Moral: OLAP techniques can handle high insertion volume, but query results are *delayed*.

Write-optimization in external memory data structures

LSM-trees

▲ロト▲御と▲酒と▲酒と 酒 の()

Leif Walsh

The insight for LSM-trees starts by asking: how can we reduce the queryability delay in OLAP?

▲ロト ▲御 ▶ ▲ 唐 ▶ ▲ 唐 ▶ ● 唐 ■ ∽ � � (

The insight for LSM-trees starts by asking: how can we reduce the queryability delay in OLAP? The buffer is small, let's index it!

The insight for LSM-trees starts by asking: how can we reduce the queryability delay in OLAP? The buffer is small, let's index it!

- Inserts go into the "buffer B-tree".
- When the buffer gets full, we merge it with the "main B-tree".

Queries have to touch both trees and merge results, but results are available immediately.

The insight for LSM-trees starts by asking: how can we reduce the queryability delay in OLAP? The buffer is small, let's index it!

- Inserts go into the "buffer B-tree".
- When the buffer gets full, we merge it with the "main B-tree".

Queries have to touch both trees and merge results, but results are available immediately.

(This specific technique (which is not yet an LSM-tree) is used in InnoDB and is called the "change buffer".)

イロト イボト イヨト イヨト

<ロト < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > の Q (

- The buffer is in-memory, so inserts are fast.
- When we merge, we put many new elements in each leaf in the main B-tree (this *amortizes* the I/O cost to read the leaf).

- The buffer is in-memory, so inserts are fast.
- When we merge, we put many new elements in each leaf in the main B-tree (this *amortizes* the I/O cost to read the leaf).

Eventually, we reach a problem:

- The buffer is in-memory, so inserts are fast.
- When we merge, we put many new elements in each leaf in the main B-tree (this *amortizes* the I/O cost to read the leaf).

Eventually, we reach a problem:

- If the buffer gets too big, inserts get slow.
- If the buffer stays too small, the merge gets inefficient (back to $O(N \log_{\beta} N)$).

How can we fix this?

<ロト < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > の Q (

<ロト < 回 ト < 回 ト < 回 ト < 回 ト < 回 ト < 回 への()</p>











Each level is twice as large as the previous level, for some value of 2.





Each level is twice as large as the previous level, for some value of 2. We'll use 2.

How do queries work?

How do queries work?



A D K A B K A B K A B K

How do queries work?



Search cost is:

$$\log_B B + \ldots + \log_B \frac{N}{8} + \log_B \frac{N}{4} + \log_B \frac{N}{2} + \log_B N$$

ъ

イロトイロトイラトイラ

How do queries work?



Search cost is:

$$\log_B B + \ldots + \log_B \frac{N}{8} + \log_B \frac{N}{4} + \log_B \frac{N}{2} + \log_B N$$
$$= \frac{1}{\log B} \left(1 + \ldots + \lg(N) - 3 + \lg(N) - 2 + \lg(N) - 1 + \lg(N)\right)$$

▲□▶▲□▶▲目▶▲目▶ 目 のへぐ

How do queries work?



Search cost is:

$$\log_{B} B + \dots + \log_{B} \frac{N}{8} + \log_{B} \frac{N}{4} + \log_{B} \frac{N}{2} + \log_{B} N$$
$$= \frac{1}{\log B} (1 + \dots + \lg(N) - 3 + \lg(N) - 2 + \lg(N) - 1 + \lg(N))$$

How do queries work?



Search cost is:

$$\log_{B} B + \dots + \log_{B} \frac{N}{8} + \log_{B} \frac{N}{4} + \log_{B} \frac{N}{2} + \log_{B} N$$
$$= \frac{1}{\log B} (1 + \dots + \lg(N) - 3 + \lg(N) - 2 + \lg(N) - 1 + \lg(N)) = O(\log N \cdot \log_{B} N)$$

How much do inserts cost?

How much do inserts cost? Cost to flush a tree T_i of size X is O(X/B).

▲ロト ▲御 ▶ ▲ 唐 ▶ ▲ 唐 ▶ ● 唐 ■ ∽ � � (

How much do inserts cost? Cost to flush a tree T_j of size X is O(X/B). Cost per element to flush T_j is O(1/B).



How much do inserts cost? Cost to flush a tree T_j of size X is O(X/B). Cost per element to flush T_j is O(1/B).



Each element moves $\leq \log N$ times.

How much do inserts cost? Cost to flush a tree T_j of size X is O(X/B). Cost per element to flush T_j is O(1/B).



Each element moves $\leq \log N$ times.

Total amortized insert cost per element is $O\left(\frac{\log N}{B}\right)$.

Write-optimization in external memory data structures

Fractal Trees

The pain in LSM-trees is doing a full $O(\log_B N)$ search in each level.

<ロ><四><日><日><日><日><日><日><日><日><日><日><日<</p>

The pain in LSM-trees is doing a full $O(\log_B N)$ search in each level. We use *fractional cascading* to reduce the search per level to O(1).

Image: A matrix

The pain in LSM-trees is doing a full $O(\log_B N)$ search in each level. We use *fractional cascading* to reduce the search per level to O(1).

The idea is that once we've searched T_i , we know where the key would be in T_i , and we can use that information to guide our search of T_{i+1} .

The pain in LSM-trees is doing a full $O(\log_B N)$ search in each level. We use *fractional cascading* to reduce the search per level to O(1).

The idea is that once we've searched T_i , we know where the key would be in T_i , and we can use that information to guide our search of T_{i+1} .

Let's examine the leaves of two consecutive levels of the LSM-tree...

Write-optimization technique #3: Fractal Trees

Add *forwarding pointers* from leaves in T_i to leaves in T_{i+1} (but remove the redundant ones that point to the same leaf):



Write-optimization technique #3: Fractal Trees

Add *forwarding pointers* from leaves in T_i to leaves in T_{i+1} (but remove the redundant ones that point to the same leaf):



Now, from a leaf node in T_i , we can jump forward to some of the leaves in T_{i+1} without searching the whole tree at T_{i+1} .

Leif Walsh	Fractal Trees	April 16, 2015	22/33
Add *ghost pointers* to leaves not pointed to in T_{i+1} in leaves in T_i :



▲□▶▲□▶▲□▶▲□▶ ▲□ ● のへの

Add *ghost pointers* to leaves not pointed to in T_{i+1} in leaves in T_i :



Now *every* leaf in T_{i+1} can be reached by a pointer in a leaf node in T_i .

		◆□▶★圖▶★園▶★園▶ ▲園	୬୯୯
Leif Walsh	Fractal Trees	April 16, 2015	23/33

[Bender, Farach-Colton, Fineman, Fogel, Kuszmaul, & Nelson '07]

After searching T_i for a missing element c, we look left and right for forwarding or ghost pointers, and follow them down to look at O(1) leaves in T_{i+1} .



[Bender, Farach-Colton, Fineman, Fogel, Kuszmaul, & Nelson '07]

After searching T_i for a missing element c, we look left and right for forwarding or ghost pointers, and follow them down to look at O(1) leaves in T_{i+1} .



This way, search is only $O(\log_R N)$ (in our example, R = 2).

	Leif Walsh
--	------------

[Bender, Farach-Colton, Fineman, Fogel, Kuszmaul, & Nelson '07]

The internal node structure in each level is now redundant, so we can represent each level as an array. We can forget about the B-tree structure above the leaves in each level! This is called a *Cache-Oblivious Lookahead Array*. The amortized analysis says our inserts are fast, but we flush a very large level to the next one, we might see a big stall. Concurrent merge algorithms exist, but we can do better.

The amortized analysis says our inserts are fast, but we flush a very large level to the next one, we might see a big stall. Concurrent merge algorithms exist, but we can do better.

We break each level's array into chunks that can be flushed independently. Each chunk flushes to a small region of a few chunks in the next level down, found using its forwarding pointers.

The amortized analysis says our inserts are fast, but we flush a very large level to the next one, we might see a big stall. Concurrent merge algorithms exist, but we can do better.

We break each level's array into chunks that can be flushed independently. Each chunk flushes to a small region of a few chunks in the next level down, found using its forwarding pointers.

Now we have a tree again!



▲□▶▲圖▶▲圖▶▲圖▶ 圖 のQ(

Easier to manage an LRU cache of blocks.

- Easier to manage an LRU cache of blocks.
- Ø More flexible with "hotspots", or non-uniform workload distributions.

- Easier to manage an LRU cache of blocks.
- Ø More flexible with "hotspots", or non-uniform workload distributions.
- Solution Flushes are O(1), so easier to reduce latency and increase concurrency with client work.

- Easier to manage an LRU cache of blocks.
- Ø More flexible with "hotspots", or non-uniform workload distributions.
- Solution Flushes are O(1), so easier to reduce latency and increase concurrency with client work.
- Easier to implement a concurrent checkpoint algorithm with small flushes.

- Easier to manage an LRU cache of blocks.
- Ø More flexible with "hotspots", or non-uniform workload distributions.
- Solution Flushes are O(1), so easier to reduce latency and increase concurrency with client work.
- Easier to implement a concurrent checkpoint algorithm with small flushes.
- Enables good tradeoffs for queries, and allows that computation to be cached without inducing I/O (this is enough complexity for a whole other talk).

- Modified B-tree-like dynamic (inserts, updates, deletes) data structure that supports point and range queries.
- Inserts are a factor $B/\log B$ (typically 10-100x in practice) faster than a B-tree: $O\left(\frac{\log N}{B}\right) < O\left(\frac{\log N}{\log B}\right).$
- Searches are a factor $\log B / \log R$ slower than a B-tree: $O\left(\frac{\log N}{\log R}\right) > O\left(\frac{\log N}{\log B}\right)$.
- To amortize flush costs over many elements, we want each block we write to be large (~4MB), much larger than typical B-tree blocks (~16KB). These compress well.

TokuDB for MySQL, TokuMX for MongoDB:

- Faster indexed insertions.
- Hot schema changes.
- Compression.
- Read-free replication on secondaries.
- Fast (no read before write) updates with messages in buffers.
- ACID transactions.
- Mixed workload concurrency.
- Faster sharding migrations (TokuMX).



▲□▶▲圖▶▲≣▶▲≣▶ ≣ のQ@

Benchmarks



100M inserts into a collection with 3 secondary indexes

April 16, 2015 31/33

æ

イロト イヨト イヨト イヨト

Benchmarks



◆□▶ ◆□▶ ◆三▶ ◆三▶ ● 三 のへの

		~-		
Lei		٧c	an a	51

Leif Walsh @leifwalsh Downloads: www.tokutek.com/downloads Docs: docs.tokutek.com Slides: bit.ly/1au1uvr

A D K A B K A B K A B K