

CS 561: Data Systems Architectures

Introduction to Indexing: Trees, Tries, Hashing, Bitmap Indexes, Database Cracking

Prof. Manos Athanassoulis

https://bu-disc.github.io/CS561/

Recap: Key-Value Stores

<key, value>

put(key, value)
stores value and associates with key

get(key) returns the associated value

delete(key) deletes the value associated with the key

get_range (key_start,key_end)
get_set(key1, key2, ...)

how to organize keys/values?

depends on the workload!





Recap: Key-Value Stores

inserts and point queries?

inserts, point queries, and range queries?



log-structured merge tree



LSM-Trees

A quick review of LSM-Trees and what is expected for the systems project















performance & cost trade-offs



bigger filters → fewer false positives memory space vs. read cost

tuning *reads* vs. *updates*



Merge Policies

Tiering write-optimized



read-optimized











 $O(log_T(N) \cdot e^{-M/N})$ $O(T \cdot log_T(N) \cdot e^{-M/N})$ lookup cost: false runs levels levels positive rate per level positive rate

false



lookup cost:

$$O(T \cdot log_T(N) \cdot e^{-M/N})$$

update cost:

$$O(log_T(N))$$

 $O(log_T(N) \cdot e^{-M/N})$



merges per level

lookup cost: $O(T \cdot log_T(N) \cdot e^{-M/N})$ $O(log_T(N) \cdot e^{-M/N})$ update cost: $O(log_T(N))$ $O(T \cdot log_T(N))$



lookup cost:

$$O(\log_T(N) \cdot e^{-M/N}) = O(\log_T(N) \cdot e^{-M/N})$$

update cost:

$$O(\log_T(N)) = O(\log_T(N))$$

for size ratio T ~

lookup cost: $O(T \cdot log_T(N) \cdot e^{-M/N})$ $O(log_T(N) \cdot e^{-M/N})$ update cost: $O(log_T(N))$ $O(T \cdot log_T(N))$

for size ratio T \land





performance & cost trade-offs



Details on Bloom filters

Inserting into a Bloom filter





Inserting into a Bloom filter





Probing a Bloom filter (true negative)





Probing a Bloom filter (false positive)



what is the probability of a false positive?






















Bloom filter false positive





Bloom filter false positive





Bloom filter false positive (derivation details)

let's focus on the term:
$$\left(1 - \frac{1}{m}\right)^n$$

assuming $\alpha = \frac{m}{n}$, and for large m, n :
 $\left(1 - \frac{1}{m}\right)^n = \left(1 - \frac{1}{\alpha \cdot n}\right)^n = \left(1 + \frac{-1/\alpha}{n}\right)^n \approx e^{-1/\alpha} = e^{-n/m}$, because $\lim_{n \to \infty} \left(1 + \frac{x}{n}\right)^n = e^x$

hence, the probability that <u>all **k** bits</u> are found set by inserting previous keys is given by

$$\left(1 - \left(1 - \frac{1}{m}\right)^{n \cdot k}\right)^k = \left(1 - \left(\left(1 - \frac{1}{m}\right)^n\right)^k\right)^k = \left(1 - \left(e^{-n/m}\right)^k\right)^k = \left(1 - \left(e^{-k \cdot n/m}\right)^k\right)^k$$



Bloom filter false positive

false positive
$$p = (1 - e^{-k \cdot n/m})^k$$

how to minimize?

it can be shown (not easy):

the optimal number of hash functions k, that minimize the false positive is:

$$k = \frac{m}{n} \cdot \ln(2)$$

Rule of thumb: k is a number, often between 2 and 10.



Bloom filter false positive

Combining
$$p = (1 - e^{-k \cdot n/m})^k$$
 and $k = \frac{m}{n} \cdot ln(2)$
we get: $e^{-\frac{m}{n} \cdot (ln(2))^2}$

details:

$$p = \left(1 - e^{-\frac{m}{n} \cdot ln(2) \cdot \frac{n}{m}}\right)^{\frac{m}{n} \cdot ln(2)} = \left(1 - e^{-ln(2)}\right)^{\frac{m}{n} \cdot ln(2)} = \left(1 - \frac{1}{2}\right)^{\frac{m}{n} \cdot ln(2)} = \left(\frac{1}{2}\right)^{\frac{m}{n} \cdot ln(2)} = \left(\frac{1}{2}\right)^{\frac{m}{n} \cdot ln(2)}$$

using twice that
$$1/2 = e^{-ln(2)}$$
, $p = (e^{-ln(2)})^{\frac{m}{n} \cdot ln(2)} = e^{-\frac{m}{n} \cdot ln(2) \cdot ln(2)} = e^{-\frac{m}{n} \cdot (ln(2))^2}$



key-value stores vs. indexes

What is an index?

Auxiliary structure to quickly find rows based on arbitrary attribute

Special form of <key, value>

indexed attribute

position/location/rowID/primary key/...



	Data Organization	Comments	
B+ Trees	Sorted & partitioned	Partition <i>k-ways</i> recursively	
LSM Trees	Log-structured & Sorted	Optimizes <i>insertion</i>	
Radix Trees	Radix	Partition using the <i>key radix</i> representation	
Hash Indexes	Hash	Partition by <i>hashing the key</i>	
Bitmap Indexes	None	Succinctly represent all rows with a key	
Scan Accelerators	None	Metadata to <i>skip accesses</i>	







B+ Trees

Search begins at root, and key comparisons direct it to a leaf. Search for 5*, 15*, all data entries >= 24* ...



Based on the search for 15^{*}, we <u>know</u> it is not in the tree!



















































LSM-trees

Is LSM-tree an index?



LSM-trees





- Hash Indexes
- Bitmap Indexes

Scan Accelerators



Radix Trees (special case of tries and prefix B-Trees)

Idea: use common prefixes for internal nodes to reduce size/height!









Bitmap Indexes

Scan Accelerators



Hash Indexes (static hashing)

#primary bucket pages fixed, allocated sequentially, never de-allocated; overflow pages if needed

h(k) mod M = bucket to insert data entry with key *k* (M: #buckets)



what if I have skew in the data set (or a bad hash function)?

Hash Indexes (static hashing)

#primary bucket pages fixed, allocated sequentially, never de-allocated; overflow pages if needed

h(k) mod M = bucket to insert data entry with key *k* (M: #buckets)





Bitmap Indexes

Scan Accelerators



Bitmap Indexes





Bitvectors

Can leverage fast Boolean operators

Bitwise AND/OR/NOT faster than looping over meta data

Bitmap Indexes



what about updates?





Scan Accelerators



Scan Accelerators

Zonemaps

Search for 25

				[00.0_]	
Z1: [32,72]	Z2: [13,45]	Z3: [1,10]	Z4: [21,100]	Z5: [28,35]	Z6: [5,12]



Scan Accelerators

Zonemaps

Search for 25 Search for [5,11]



Scan Accelerators

Zonemaps

Search for 25 Search for [5,11] Search for [31,46]

Z1: [32,72]	Z2: [13,45]	Z3: [1,10]	Z4: [21,100]	Z5: [28,35]	Z6: [5,12]
-------------	-------------	------------	--------------	-------------	------------


Scan Accelerators

Zonemaps

Search for 25 Search for [5,11] Search for [31,46]

Z1: [32,72]	Z2: [13,45]	Z3: [1,10]	Z4: [21,100]	Z5: [28,35]	Z6: [5,12]
-------------	-------------	------------	--------------	-------------	------------



Scan Accelerators

Zonemaps

							Search for 25
	Z1: [32,72]	Z2: [13,45]	Z3: [1,10]	Z4: [21,100]	Z5: [28,35]	Z6: [5,12]	Search for [5,11] Search for [31,46]
:r							

if data were sorted:

Z1: [1,15]	Z2: [16,30]	Z3: [31,50]	Z4: [50,67]	Z5: [68,85]	Z6: [85,100]
Z1: [1,15]	Z2: [16,30]	Z3: [31,50]	Z4: [50,67]	Z5: [68,85]	Z6: [85,100]

Search for 25 Search for [5,11] Search for [31,46]



Scan Accelerators

Zonemaps

						Search for 25
71. [20 70]	72. [12 /15]	72. [1 10]	7/1.[21 100]	75. [28 35]	76. [5 12]	Search for [5,11]
21. [32,72]	22.[13,45]	23. [1,10]	24. [21,100]		20. [3,12]	Search for [31,46]

if data were sorted:

Z1: [1,15]	Z2: [16,30]	Z3: [31,50]	Z4: [50,67]	Z5: [68,85]	Z6: [85,100]
------------	-------------	-------------	-------------	-------------	--------------

Search for 25 Search for [5,11] Search for [31,46]

what if data is perfectly uniformly distributed?



Z1: [1,99]Z2: [2,95]Z3: [1,100]Z4: [2,100]Z5: [3,97]Z6: [2,99]
--

What are the possible *index designs*?





idea: there is an *ideal* data organization

what is it (for a collection of keys – e.g., a column)? *sorted*!

we can reach it *eventually* if we use the *workload as a hint*







	search < 15		search < 90		
32		11		11	
19		6		6	
11	~ 1C	12	~ 1E	12	
6	< 12 -	32	< 15	32	
123		19		19	
55		123		123	
12		55	> 00	55	
78		78	> 90 -	78	









what about updates/inserts?



Project Implementation

What to plan for the implementation (1/3)

Durable Database (open/close without losing state)

Components:

Memory buffer (array, hashtable, B+ tree) Files (sorted levels/tiers) Fence pointers (**Zonemaps**) Bloom filters



What to plan for the implementation (2/3)

Durable Database (open/close without losing state)

Components:

Memory buffer (search, read, write, unpin) Priority data structure Eviction policy



What to plan for the implementation (3/3)

API + basic testing and benchmarking available at:

LSM Implementation:

https://github.com/BU-DiSC/cs561_templatedb

with a Reference Bloom filter implementation

Bufferpool Implementation:

https://github.com/BU-DiSC/cs561_templatebufferpool





CS 561: Data Systems Architectures

Introduction to Indexing: Trees, Tries, Hashing, Bitmap Indexes, Database Cracking

Prof. Manos Athanassoulis

https://bu-disc.github.io/CS561/