Carnegie Mellon University

Trees Indexes (Part II)







Andy Pavlo Computer Science Carnegie Mellon Univ.

ADMINISTRIVIA

Project #1 is due Wednesday Sept 26th @ 11:59pm

Homework #2 is due Friday Sept 28th @ 11:59pm

Project #2 will be released on Wednesday Sept 26th. First checkpoint is due Monday Oct 8th.



TODAY'S AGENDA

Additional Index Usage Skip Lists Radix Trees Inverted Indexes



Most DBMSs automatically create an index to enforce integrity constraints.

- \rightarrow Primary Keys
- \rightarrow Unique Constraints
- \rightarrow Foreign Keys (?)









Most DBMSs automatically create an index to enforce integrity constraints.

- \rightarrow Primary Keys
- \rightarrow Unique Constraints
- \rightarrow Foreign Keys (?)

```
CREATE TABLE foo (
  id SERIAL PRIMARY KEY,
  val1 INT NOT NULL UNIQUE,
  val2 VARCHAR(32) UNIQUE
```

CREATE TABLE bar (id INT REFERENCES foo (val1), val VARCHAR(32));

DATABASE GROUP

PARTIAL INDEXES

Create an index on a subset of the entire table. This potentially reduces its size and the amount of overhead to maintain it.

One common use case is to partition indexes by date ranges. \rightarrow Create a separate index per month, year.







PARTIAL INDEXES

Create an index on a subset of the entire table. This potentially reduces its size and the amount of overhead to maintain it.

One common use case is to partition indexes by date ranges. \rightarrow Create a separate index per month, year.



| SELECT | b | FROM foo |
|--------|---|-------------|
| WHERE | а | = 123 |
| AND | С | = 'WuTang'; |



COVERING INDEXES

If all of the fields needed to process the query are available in an index, then the DBMS does not need to retrieve the tuple.

This reduces contention on the DBMS's buffer pool resources.





INDEX INCLUDE COLUMNS

Embed additional columns in indexes to support index-only queries. Not part of the search key.







INDEX INCLUDE COLUMNS

Embed additional columns in indexes to support index-only queries. Not part of the search key.





INDEX INCLUDE COLUMNS

Embed additional columns in indexes to support index-only queries. Not part of the search key.



The index does not need to store keys in the same way that they appear in their base table. 

The index does not need to store keys in the same way that they appear in their base table. CREATE INDEX idx_user_login
 ON users (login);



The index does not need to store keys in the same way that they appear in their base table. You can use expressions when declaring an index.

CREATE INDEX __user_login ON users __gin);



The index does not need to store keys in the same way that they appear in their base table. You can use expressions when declaring an index.

CREATE INDEX ._user_login ON users ...gin);

CREATE INDEX idx_user_login
 ON users (EXTRACT(dow FROM login));



The index does not need to store keys in the same way that they appear in their base table.

| SELECT | * FROM users | |
|--------|--------------------------|--|
| WHERE | EXTRACT (dow | |
| | ♥ FROM login) = 2 | |

You can use expressions when declaring an index.

CREATE INDEX __user_login ON users __gin);

CREATE INDEX idx_user_login
ON users (EXTRACT(dow FROM login));



The index does not need to store keys in the same way that they appear in their base table.

You can use expressions when declaring an index.

CREATE INDEX idx_user_login
ON users (EXTRACT(dow FROM login));

```
CREATE INDEX idx_user_login
   ON foo (login)
WHERE EXTRACT(dow FROM login) = 2;
```



OBSERVATION

The easiest way to implement a <u>dynamic</u> orderpreserving index is to use a sorted linked list. All operations have to linear search. \rightarrow Average Cost: O(N)





OBSERVATION

The easiest way to implement a <u>dynamic</u> orderpreserving index is to use a sorted linked list. All operations have to linear search. \rightarrow Average Cost: O(N)





OBSERVATION

The easiest way to implement a <u>dynamic</u> orderpreserving index is to use a sorted linked list. All operations have to linear search. \rightarrow Average Cost: O(N)





SKIP LISTS

Multiple levels of linked lists with extra pointers that **<u>skip</u>** over intermediate nodes.

Maintains keys in sorted order without requiring global rebalancing.

Skip Lists: A Probabilistic Alternative to Balanced Trees

Skip lists are a data structure that can be used in place of balanced trees. Skip lists use probabilistic balancing rather than strictly enforced balancing and as a result the algorithms for insertion and deletion in skip lists are much simpler and significantly faster than equivalent algorithms for balanced trees.

William Pugh

Binary trees can be used for representing abstract data types such a discinturies and extered lists. They were well well been the elements are inserted in a random order. Some sequences of operations, such as inserting the elements in order, produce degenerate data structures that give very poor performance. If it were possible is considered by the order of the list particular data structures that give very loop or performance. If it were possible is considered by the structure of the strucput sequences. It must cause queries sum by answered on vince, so mandowing permitting the ling stat is superclass. Marking the algorithms rearrange the tree as operations are performed to maintain central balance conditions and assume good perfor-

 \overline{Mg} for one as probabilities dimensive to balanced stress. So this are balanced stress. So this are balanced stress to work of the set of

Balancing a data structure probabilistically is easier than explicitly minimizing the balance. For many applications, skip lasts are a more natural representation than trees, also leading to simpler algorithms. This simplicitly of skip its algotrithm makes them easier to implement and provides significant constant factors peed improvements over halanced tree space efficient. They can easily be configured to require an average of 11-y goodness peed easily over leading and swarges of 11-y goodness peed easily over leading and the scale part of the scale part of the scale mode.

SKIP LISTS

We might need to examine every node of the list when searching a linked list ($F(gure \ 1a)$). If the list is stored in sorted order and every other node of the list also has a pointer to the node two ahead it in the list ($F(gure \ 1b)$), we have to examine no more than $\lceil n/2 \rceil + 1$ nodes (where n is the length of the list). Also giving every fourth node a pointer four abaud (Figure 1c) requires that no more than (m+4) + 2 nodes be cannined. If every $(2^k)^{th}$ node has a pointer 2^k nodes ahead (Figure Id), the number of nodes that must be examined can be reduced to logg n^k while only doubling the number of pointers. This data structure could be used for fast searching, but insertion and deletion would be impractical.

A stort has the forward pottern is called a level analytic (the event (2)) mode that has a forward pottern is such as the form (eV). (F) would be the potter 2 mode shared, then iterative (the event (2)) and the stort potter potter is and a stort. This would be possible of the trend or mode share are possed another, but in the store proposed in the trend or mode source characterization would be the trend or the trend of the trend

SKIP LIST ALGORITHMS

This section gives algorithms to search for, insert and delete elements in a dictionary or symbol table. The Search opera tion returns the contents of the value associated with the de sired key or failure if the key is not present. The Insert opera tion associates a specified key with a new value (inserting the key if it had not already been present). The Delete operation deletes the specified key. It is easy to support additional oper ations such as "find the minimum key" or "find the next key" Each element is represented by a node, the level of which is chosen randomly when the node is inserted without regard for the number of elements in the data structure. A level i node has (forward pointers, indexed 1 through i. We do not need to store the level of a node in the node. Levels are capped at some appropriate constant MaxLevel. The level of a list is the maximum level currently in the list (or 1 if the list is empty). The header of a list has forward pointers at levels on through MaxLevel. The forward pointers of the header at levels higher than the current maximum level of the list point





SKIP LISTS

A collection of lists at different levels

- \rightarrow Lowest level is a sorted, singly linked list of all keys
- \rightarrow 2nd level links every other key
- \rightarrow 3rd level links every fourth key
- \rightarrow In general, a level has half the keys of one below it

To insert a new key, flip a coin to decide how many levels to add the new key into. Provides approximate **O(log n)** search times.

















12
















































First **logically** remove a key from the index by setting a flag to tell threads to ignore.

Then **physically** remove the key once we know that no other thread is holding the reference.

























SKIP LISTS

Advantages:

- \rightarrow Uses less memory than a typical B+Tree if you don't include reverse pointers.
- \rightarrow Insertions and deletions do not require rebalancing.

Disadvantages:

- \rightarrow Not disk/cache friendly because they do not optimize locality of references.
- \rightarrow Reverse search is non-trivial.



RADIX TREE

Represent keys as individual digits. This allows threads to examine prefixes one-by-one instead of comparing entire key.

- \rightarrow The height of the tree depends on the length of keys.
- \rightarrow Does not require rebalancing
- \rightarrow The path to a leaf node represents the key of the leaf
- \rightarrow Keys are stored implicitly and can be reconstructed from paths.





































































Not all attribute types can be decomposed into binary comparable digits for a radix tree.

- → **Unsigned Integers:** Byte order must be flipped for little endian machines.
- → **Signed Integers:** Flip two's-complement so that negative numbers are smaller than positive.
- → **Floats:** Classify into group (neg vs. pos, normalized vs. denormalized), then store as unsigned integer.
- → **Compound:** Transform each attribute separately.



Int Key: **168496141** Hex Key: **0A 0B 0C 0D**













IN-MEMORY TABLE INDEXES

Processor: 1 socket, 10 cores w/ 2×HT Workload: 50m Random Integer Keys (64-bit)







OBSERVATION

The tree indexes that we've discussed so far are useful for "point" and "range" queries:

- \rightarrow Find all customers in the 15217 zip code.
- \rightarrow Find all orders between June 2018 and September 2018.

They are **<u>not</u>** good at keyword searches:

 \rightarrow Find all Wikipedia articles that contain the word "Pavlo"



WIKIPEDIA EXAMPLE



WIKIPEDIA EXAMPLE

If we create an index on the content attribute, what does that actually do?

CREATE INDEX idx_rev_cntnt
 ON revisions (content);

This doesn't help our query. Our SQL is also not correct...

SELECT pageID FROM revisions
WHERE content LIKE '%Pavlo%';



INVERTED INDEX

An *inverted index* stores a mapping of words to records that contain those words in the target attribute.

- \rightarrow Sometimes called a *full-text search index*.
- \rightarrow Also called a *concordance* in old (like really old) times.

The major DBMSs support these natively. There are also specialized DBMSs.



CMU 15-445/645 (Fall 2018)

elasticsearch

Sphinx

QUERY TYPES

Phrase Searches

 \rightarrow Find records that contain a list of words in the given order.

Proximity Searches

 \rightarrow Find records where two words occur within *n* words of each other.

Wildcard Searches

→ Find records that contain words that match some pattern (e.g., regular expression).


DESIGN DECISIONS

Decision #1: What To Store

- \rightarrow The index needs to store at least the words contained in each record (separated by punctuation characters).
- \rightarrow Can also store frequency, position, and other meta-data.

Decision #2: When To Update

 \rightarrow Maintain auxiliary data structures to "stage" updates and then update the index in batches.



CMU 15-445/645 (Fall 2018)

CONCLUSION

B+Trees are still the way to go for tree indexes.

Inverted indexes are covered in <u>CMU 11-442</u>.

We did not discuss geo-spatial tree indexes:

- \rightarrow Examples: R-Tree, Quad-Tree, KD-Tree
- \rightarrow This is covered in <u>CMU 15-826</u>.



CMU 15-445/645 (Fall 2018)

NEXT CLASS

How to make indexes thread-safe!



CMU 15-445/645 (Fall 2018)