#### Log-Structured Merge Trees (LSMs) Bloom Filters

14-848 (Cloud Infrastruture)

# Scenario

- A system is needed to support high-throughput updates
- The total data volume is larger than the main memory budget
- Writes to secondary storage occur more quickly and efficiently when batched than when written individually.
  - For example, writing a whole block of data at a time amortizes disk seek and rotational delay
- Sorting and indexing data in main memory can be done relatively efficiently causing relatively little delay.

# Collect and Batch Updates In Memory

- Collect updates in memory
  - Sort them somehow
    - Sorted string list
    - Tree, Etc.
- Updates possible to in-memory values
  - But, once a value is written to disk, it stays written
  - Queries will need to find all records and merge
  - Tombstone deletes

# Spill From Memory To Disk

- As memory budget approaches full, spill them to disk
  - Write out entire sorted string table
  - Write out a subtree, then remove and prune it in memory
- Each dump from memory to disk forms a "run" of some kind
  - Runs are time ordered

# Merge, Idea #1

- Possibility #1:
  - Merge portions of in-memory data structure into on-disk data structure as spilled
  - Common when pruning in-memory trees and merging into on-disk trees
  - Slows the freeing of memory

# Merge Idea #2

- Possibility #2:
  - Dump from memory into new "run", i.e. data structure in secondary storage
  - Maintain in-memory Bloom Filters, one per disk run, to support queries
  - Upon query, check Bloom Filters
  - Then check on-disk runs only where Bloom filters indicated possible match
- Merge updates disk data structures in background
  - By similar tree pruning, if tree
  - By merging files into new files if tables
  - Delete then update Bloom Filter, since false positives aren't fatal

# Merge Idea #3

- Compaction occurs as part of the merges
  - Deletes Tombstoned records
  - Merges multiple updates into one
  - Recovers storage from merged updates and deleted values

# Log-Structured Merge Trees (LSMs)

- When we spill subtrees or branches from an in-memory tree into a tree in secondary storage, this strategy is known as a Log-Structured Merge Tree (LSM) Tree
  - The in-memory tree is often known as  $C_0$  and the tree in secondary storage is known as  $C_1$ .
  - If there are more levels of trees (not within a tree), they are known as  $C_1$ ,  $C_2$ , etc.

## Memtable, SSIndex, SSTable

- Common idiom in practice
  - Memtable in main memory contains sorted values and likely sorted <key, offset> index.
  - When spilled to disk is divided into SSTables and SSIndexes written separately
  - Indexes or Bloom Filters kept in memory
  - Merging in background when threshold met in terms of number of tables, etc.
  - Merges perform compaction
  - Write-Ahead logs used to aid recovery.
- Used in some form by Cassandra, Hbase, LevelDB, BigTable, Etc.



- Overall strategy
  - Fill memory
  - Spill to disk
  - Search disk runs until they can be merged
  - Use Bloom Filters to minimize unproductive searches
  - Updated in-memory, but merge independent changes once on disk.
- The overall strategy is sound even if it...
  - Does not involve trees, for example by using sorted string tables, and
  - Even if it leaves a forest of data structures to be searched after consulting a Bloom Filter.

#### New Scenario

- Multiple (May be a large number), independent key-value stores, such as LSMs
  - Each is in secondary storage and relatively slow to access/search.
  - Values may not be contained within zero or more of the stores
    - Never present
    - Written Once
    - Written and Over-Written
    - Written and Lazy Deleted by "Tombstoning"
    - Multiple partial updates that need to be unioned to form up-to-date record

## Goal

- Minimize wasted time/energy searching data stores which do not have requested key.
- Add minimal overhead w.r.t. main memory foot-print, processor load, etc.

# Idea and Challenges

- Idea: Keep an in-memory index
- Challenge: The on-disk data structures are most-likely indexes or keep relatively small metadata and are still too large to fit in memory, hence the need to keep them on disk
  - Large data values are normally kept in separate data stores and indexes keep references to their location, etc.
- Revised idea: Need an extremely dense index, e.g. ideally one or a few bits per item

## Idea #1

#### • Idea:

- Keep an in-memory hash table: <key, store-ID>
- Good:
  - Near constant time ability to discover which stores to search
- Challenges:
  - Need to manage collision
  - Lists get to be expensive in processor and memory
    - Could be limiting

## Idea #2

- Idea #2:
  - Keep one hash table per external store, check each
  - One bit per address: Present or not
  - Search only data structures where present
- Good:
  - Smaller
  - Fast to check, even given multiple data structures
- Challenges:
  - Managing collision means complexity
  - Reducing collision means a much bigger table, which means more memory wasted

## Idea #3

- Idea:
  - Ignore Collision
- Good:
  - Can have small tables with little complexity
- Bad:
  - Could thing something is present in secondary store when it is really a colliding record
    - Waste time looking just to find nothing
- Thought:
  - If collisions aren't super common, this is probably okay.

### Bloom Filter

- This device is called a Bloom Filter
  - First conceived by Burton Howard Bloom in 1970

#### Parameters

- How many keys, max?
- Memory budget?
  - One bit per address
- How many hash functions? One bit per hash function consumed per key
  - Less overlap
- What false positive rate is tolerable?
- Based upon the probabilities, we could reduce this to an equation. But, we'll take a pass. It is well studied.
  - See Wikipedia for a discussion: https://en.wikipedia.org/wiki/Bloom\_filter