

GDC MapReduce Notes

Josh Snider

October 16, 2015

FP Review

Data structures are not modified, new data structures are made instead. This means we don't need to lock data structures. It also means that since we don't have side effects, we make everything a separate thread and put them together at the end.

Functional programming also allows you to use functions as arguments as the name implies. Functional programming is also good for generic yet typesafe programming, but languages like OCaml can still screw this up.

Standard functional programming functions:

- Map - Apply a function to each element of a list and return a list made by concatenating the results together.
- Fold - Track the result of calling f on each element and an accumulator, return the final accumulator. Can be either a left or right fold.

These are basically the Map and Reduce in MapReduce.

MapReduce

Motivated by large scale data processing, parallelize across thousands of computers, and make it easy for programmers to use. MapReduce has built-in fault tolerance and network monitoring. The user basically has two functions to provide

- map (in_key, in_value) -> (out_key, intermediate_value) list - Input consists of records from various sources.
- reduce (out_key, intermediate_value list) -> out_value list

After the map phase happens, we combine all the intermediate values for a given key into one list and then start the reduce phase. It's good practice for reduce to only have out_value. We use barrier synchronization, to prevent people from going to the reduce step while people are still mapping.

Optimization

Each map function runs parallelly and the reduction for each key runs in parallel. This means we have a bottleneck waiting for mapping to finish before we start reduction. How do we optimize this?

- Slow-moving “map” tasks are run on multiple hosts and then we take the first one to finish.
- Google considered making reducers lazy, but that has some design flaws that made it not used.
- We can have “combiners” which run mini-reduce phases before the actual reduce in order to save bandwidth. If our reduce is both associative and commutative, then we can use it as our combine as well.

Written in C++, with Java and Python bindings. The dividing program tries to divide up map tasks so that mappers are on the same computer as the data. Tasks are chunked in 64MB which is the same size as the Google File System’s chunks.

Ensuring fault tolerance

- We give up on tasks if certain key-value pairs reliably crash.
- Mappers that fail before reporting results are redone.
- Reducers don’t claim to be done, until their results are reliably backed up.

Impact

Has had great results at Google and is a shining jewel of functional programming. Greatly simplifies large-scale computations. Lets programmer focus on problem and let library handle messy details.

File Systems

- A filesystem is a system to join data in a tree of files, by writing them to disk. Sometimes, they can support remote files or local caching.
- Folders are the same as namespaces.
- A “distributed filesystem” is one that can reach files on other machines.
- Any decent “distributed filesystem” guarantees that many people can access many files simultaneously. It needs to be performant at the same time and maintain data integrity / consistency.

NFS

- Made by Sun in the 80s. Pretty standard Unix FS. We mount remote drives onto local host.
- Original was completely stateless. Higher-level protocols handle that stuff.
- NFS defines a virtual filesystem (like an interface, not an implementation).
- NFS locking is done with leases. Clients request locks. Servers tell them if that succeeds or not and take back locks that aren’t renewed.
- When we close a file, we push changes. When we open a file, we pull changes.

- An NFS volume is managed by a single server. This makes concurrency easier, but puts stress on it.
- POSIX compliance makes it portable, but very generalized.

GFS

- Google needed a way to redundantly store data on cheap and unreliable systems. Also, wanted to optimize it for Googly purposes.
- Assumptions:
 - Components suck.
 - Modest number of HUGE files.
 - Files are write-once. Mostly appending to data.
 - Want to read large serial chunks.
- Design Decisions:
 - 64 MB fixed-size chunks.
 - Each thing replicated on 3+ chunk servers.
 - Master coordinates access and stores metadata.
 - No data caching (files are too big).
 - Interfaces customized for Google, but POSIX-esque.
- Master server is a single point of failure. We solve this by having backups of the master called “shadow masters”.
- What metadata does the master have?
 - File and chunk namespaces.
 - Mapping from files to chunks.
 - Locations of chunk’s replicas.
 - An operations log for metadata updates.
- Except for the operations log, the metadata is super small and can stay in RAM. The operations log needs to be stored on disk for robustness reasons.
- We can use the operations log to restore the master to a good state if it becomes corrupted.
- GFS Mutation = write or append. Our goal for mutations is to minimize master involvement.
- We use a lease mechanism. The master picks a replica as primary and gives it a lease for mutations, the primary defines a serial order for the changes, and the replicas follow it.
- Atomic record append - GFS lets us append a record to a chunk atomically. Message is at-least-once, so we may do it multiple times, and order may differ between replicas.
- So, we have a relaxed consistency model, except for our master’s metadata. Google’s fine with that, but it’s something to keep in mind.

- What are the master's responsibilities?
 - Metadata storage
 - Namespace management/locking.
 - Monitor system health.
 - Chunk creation, re-replication, and rebalancing.
 - Garbage collection - We delete things by renaming them to something hidden, then we do the actual clean up when we're not busy.
 - Prompt people with out of date chunks to update them.
- This is a highly fault tolerant system, with high availability and specified data integrity.
- In conclusion, sometimes simple solutions are good and you should expect hard drives to break.