

# Distributed Data Structures

## xFS: Serverless Network File System

February 3, 2004

### I. DDS

Goal: new persistent storage layer that is a better fit for Internet services

- o separate service logic from durability, availability, consistency issues
- o interface is that of a data structure, not a SQL query (intentionally navigational)
- o automate replication and partitioning across the cluster
- o automate recovery and high availability
- o CAP: CA (not P) -- depends on a machine room with a redundant network

Nit: abstract's perf numbers have too much precision! (should have been limited to 2-3 significant digits)

Alternatives for persistent storage:

- o DBMS: chooses C over A, interface is inherently slow: parse and plan each query.
- o Enterprise Java Beans (EJB): map Java objects onto tables. Typical use has an array of objects map to a table with one row for each object; row stores the serialized bytes for that object => expensive copy to get the object in and out of the database. Also, storage is never local and there is little control over partitioning and replication (which nodes). Should be 10x or more slower, but no reliable data.
- o File system: also has an expensive interface and very hard to provide fine-grain atomicity. (primary atomic operation is file rename, which implies copying the whole file). Fine choice for large objects, but not for small ones...

Internet requirements:

- o extreme scale: billion request/day
- o high availability
- o unknown but potential fast growth -- must be able to add capacity quickly
- o overload will occur -- need graceful degradation

Cluster properties:

- o Incremental scalability -- add nodes over time
- o Potential for high availability -- but needs to be written!
- o Natural parallelism for both I/O and CPUs
- o High BW, low-latency, \*partition free\* network (CAP)
- o Machine room properties: security, administration, reliable power, AC, networking

(these generally don't apply to P2P systems in practice)

#### DDS design decisions:

- o hash table interface: 64-bit keys, byte array values, put/get/delete elements, create/destroy hash tables
- o operations are atomic, sequences of operations (trans-actions) are not
- o two-phase commit across all replicas for consistency of writes
- o read any replica ( $\Rightarrow$  higher read throughput than write throughput)
- o event-driven design for high concurrency. this was a mistake and led to Capriccio. It is being rebuilt from scratch in C over Capriccio. (see Java problems below)
- o assume no network partitions (CAP)

#### Partitions:

- o key idea: break hash tables into many small fixed size "partitions"
- o replicate partitions (not tables); replica groups are a set of partitions with the same data
- o recovery: recover each partition independently
- o small partitions mean that you can lock the whole partition during recovery (other partitions operate independently)
- o also implies that you can lazily recover each partition, or you can be proactive.

#### Fault tolerance

- o DDS library: handles two-phase commit across replica group.
  - commit occurs when at least one replica receives a commit message
  - failure before this and all replicas will time out and discover there was no commit by contacting each other
  - failure after this and the replicas will learn about the commit from the member that recorded it, and will then all commit
- o Brick:
  - failure before commit will cause 2PC to fail and everyone aborts
  - failure after agreement to commit, but before commit: other replica's commit, this node will get the new value during recovery
- o Other service code:
  - if all durable state is in the DDS, this is recovered automatically on restart
  - soft state can be rebuilt (e.g. caches)
  - if local durable state (not in DDS), then service is on its own (e.g. might store some data in a DBMS also)

#### Data Partitioning Map:

- o maps HT key to partition id
- o uses a trie so that it is easy to split or merge partitions as needed (increment data repartitioning!)
- o replicated on all DDS clients, eventual consistency (stale data causes a repair and a

retry); staleness is detected by comparing a hash of the maps to the current values.

#### Replica Group Map:

- o maps partition id to the set of replicas
- o replicated on all DDS clients, eventual consistency as above
- o writes are 2PC to all replicas
- o reads go to one replica. which one? want to spread out reads in the same way to maximize effective cache space!

#### Recovery:

- o Node failure takes down all bricks on that node; brick failure takes down all partitions on that brick
- o a recovering brick must catch up all its partitions
- o key idea 1: allow a brick to say “no” -- this simplifies the code greatly! DDS library will retry in a bit. “NO” means that brick can be inconsistent for a while (but not too long)
- o all committed operations must use up-to-date maps (DP and RG), else retry
- o bricks catch up one partition at a time: just copy that partition from another replica.
- o Updates to a partition stop during recovery, but this is OK but partitions are small!
- o Must decide how proactive to be in recovery; OK to be completely lazy (wait for partition to be accessed, but hurts latency)
- o see recovery graph (figure 8)

#### Performance:

- o great scalability, availability, graceful degradation for reads (less so for writes)
- o Problems: Java GC can screw up performance! One a node gets behind it needs more memory (for growing queues), so it tends to get worse and worse! Need admission control, control over GC, or probably both. SEDA would probably fix this...
- o Events make this much worse! Event driven systems (in Java) create a lot of garbage, as events are just passed up through layer and then GC'd. Threads would have most of this state on the stack and thus create very little garbage... (this would also be a bad place for functional languages, which need to copy a great deal)
- o serious problems with Java: GC issues, extra copying, asynchronous I/O... We later fixed async I/O, but the other two issues remain. In general, Java is too high level for this kind of work...

## II. xFS

#### Key background:

- o RAID small write problem
- o LFS for basic servers
- o Zebra for network striping, but replace centralized manager with a distributed manager, and use a parallel cleaner and stripe groups (rather than stripe across all disks)

- o Multiprocessor cache consistency: aggressive locality exploitation

Three problems to solve:

- o scalable, distributed metadata
  - cache coherent
  - reconfigurable after failure (⇒ a level of indirection!)
- o Scalable way break storage servers into stripe groups
- o Scalable log cleaning

Philosophy: “anything, anywhere”. All parts of the file system can be anywhere:

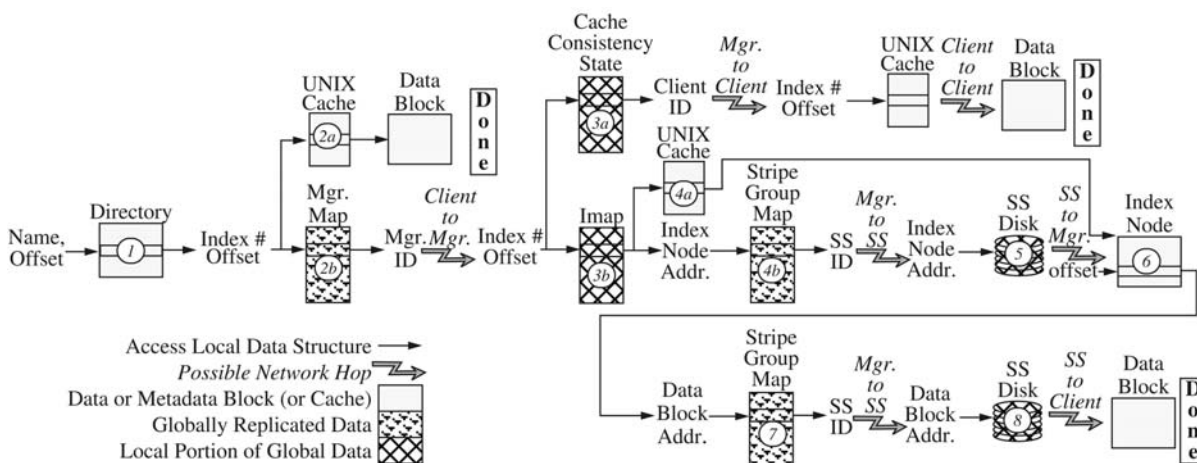
- o location independent
- o Any node can take over for any other
- o locality determines location...

Four basic modules:

- o client
- o storage server
- o manager: for each file (index number), tracks cache info and metadata (inode)
- o cleaner

Map 1: Manager Map:

- o globally replicated
- o map index number to manager (which stores the imap for that index number), i.e. which machine manages this file?



**Figure 3. Procedure to read a block.** The circled numbers refer to steps described in Section 3.2.1. The network hops are labelled as “possible” because clients, managers, and storage servers can run on the same machines. For example, xFS tries to co-locate the manager of a file on the same machine as the client most likely to use the file to avoid some of the network hops. “SS” is an abbreviation for “Storage Server.”

- o can use to balance metadata load

#### Map 2: Imap:

- o split across managers
- o maps index number to inode (i.e. where on disk is the inode?)

#### Map 3: File Directories:

- o maps file names to index numbers
- o stored in regular files (just like previous file systems)
- o How to pick new index numbers when create a file or directory? a) pick randomly for load balancing, b) pick locally for better locality. xFS does the latter.

#### Map 4: Stripe Group Map

- o partitions storage servers into groups
- o globally replicated
- o maps stripe-group IDs to a list of storage servers
- o Implies a disk address = (group ID, segment number, offset)
- o Groups can become obsolete, which allows GC or groups rather than having to update all of the disk addresses (everywhere)
- o Groups have better fault tolerance than striping without groups (why and at what cost?)

#### Caching

1. local block cache at client
2. someone else's cache (known by the manager)
3. inode cache at manager

#### Deltas:

- o critical section for atomic writes
- o Since we write data and metadata to different servers, we can't ensure that they happen atomically (as a group)
- o Deltas just allow the client to mark the boundaries of an atomic op, so that on replay it will ignore partial deltas

#### Cooperative caching: manager tracks list of readers

- o write causes invalidation of all readers, acquires write permission
- o read from different client converts writer to reader and adds new reader to list

#### Recovery:

- o starts at storage servers
- o Rebuild fragment ID to location map from checkpoint + roll forward

- o Stripe ID (and members) in every fragment
- o ifile stores imap on disk, but it has to be in a known place since it's needed for recovery (list it in a checkpoint)
- o Rollback to last checkpoint, load imap from ifile
- o Then roll forward through client logs to finish imap changes (using deltas)
- o Find consensus on which manager handles which clients logs...