

Lessons from Giant-Scale Services

February 19, 2004

Experience paper on how to build and operate very large Internet sites...

Key ideas:

- o Load management
- o Partitioning vs. Replication, load redirection
- o Availability metrics: yield and harvest, MTTR emphasis
- o Online evolution
- o Graceful degradation
- o the DQ principle

Basics:

- o Data center: power, AC, networking, no people, minimal cables
- o Extreme symmetry simplifies everything
- o FRU = “field replaceable unit” equals whole node (worry about subpieces offline)
- o management backplane is a good idea for staging, management under duress
- o data center contract limits temperature, power, networking variations

Load Management

- o key idea: need a **highly available name**
- o DNS: map names to IP addresses
 - remap when IP address disappears
 - mediocre solution: DNS expiration takes a while, not all browsers obey expiration correctly
- o L4 switches
 - forward incoming TCP connections to the “up” (IP, port) addresses
 - load balance based on number of open connections
 - manage the “up” set automatically (by detecting dropped connections/resets)
 - (but better to be proactive! why?)
 - come in pairs with automatic hot fail-over (avoid single point of failure)
 - L7 switches switch based on URLs
- o smart clients: clients manages the name mapping and failover directly
 - this is the best solution, as it handles disaster recovery well (redirect to new data center)

- not built into HTTP, but can do it if the client is an applet or program

Availability Metrics

- o focus on MTTR not MTBF -- faster debugging cycle, more stable
- o yield = better version of uptime
- o harvest reflects potential for incomplete data used for answers
- o good designer plans how faults affect MTTR, yield and harvest!

DQ Principle

- o claim: data per query * queries/sec == constant
- o represents the total data flowing through the system (MB/s like bandwidth)
- o only true is system is running near capacity, otherwise you can increase D or Q
- o Lots of uses: spec hardware, eval software changes, capacity planning, failover planning, etc.

Replication vs. Partitioning

- o replication typically viewed as “better”: maintains 100% harvest during a fault, but 50% yield (if it was at capacity!)
- o ...but partition maintains 100% yield, 50% harvest
- o There is no “better” -- just different optimizations
- o Load redirection problem: not enough to replicate data, must have replicated **capacity** (sufficient DQ points)
- o See Table 1 for the overload factors\
- o For write-intensive traffic replication costs more than partitioning, but for read mostly they are essentially the same!
- o => partition until you reach a convenient size, and then replicate the whole set. AOL Caches partition within a rack (used to be 5 machines), and then replicate racks for capacity
- o Can replicate only the important data, and then ensure that lost harvest does not include the important stuff
- o Randomization makes worst case data loss same as the average case.

Graceful Degradation

- o large sites on open networks will get overloaded
 - Schwab uses managers to handle the overload of phone calls during a market event
- o correlated failures, although rare, can reduce capacity and cause overload
- o key insight: to handle overload we can either limit Q (admission control), or reduce D

(to increase Q) to handle more capacity at some loss in quality

- o best answer is usually a combination! deny the expensive queries and increase caching and reduce harvest
 - cost-based admission control
 - priority or value-based AC (for good customers or financially important transactions)

Disaster Tolerance

- o key idea: estimate correlated failures due to a disaster, typically one whole data center
- o figure out replication/partitioning for that set of failures
- o need load redirection to *outside* the data center => can't use L4 switch (both of them fail as well); ideally smart clients, else DNS
- o plan on overload due to redirected load, and handle via graceful degradation!

Online evolution

- o quality: standards go down under fast evolution! also lower for online services (vs normal software) why?
- o staging: extra space to store two versions of the software (or data) makes it easy to switch back and forth; important to automate "revert"!
- o three ways to upgrade
 - fast upgrade
 - rolling upgrade
 - big flip

Moving sites!

- o Inktomi moved the data center twice while it was online!
- o Many upgrades to hardware, OS, schema, protocols, ...