

# Inside LiveJournal's Backend

or,  
*“holy hell that's a lot of hits!”*

November 2004

**Brad Fitzpatrick <brad@danga.com>**  
**Lisa Phillips <lisa@danga.com>**

**Danga Interactive**  
danga.com / livejournal.com

**SOME RIGHTS RESERVED**



This work is licensed under the Creative Commons **Attribution-NonCommercial-ShareAlike** License. To view a copy of this license, visit <http://creativecommons.org/licenses/by-nc-sa/1.0/> or send a letter to Creative Commons, 559 Nathan Abbott Way, Stanford, California 94305, USA.

# Administrivia

- Question Policy
  - Anytime... interrupt!
  - also at end
- Slides online:
  - <http://www.danga.com/words/>

# The Plan

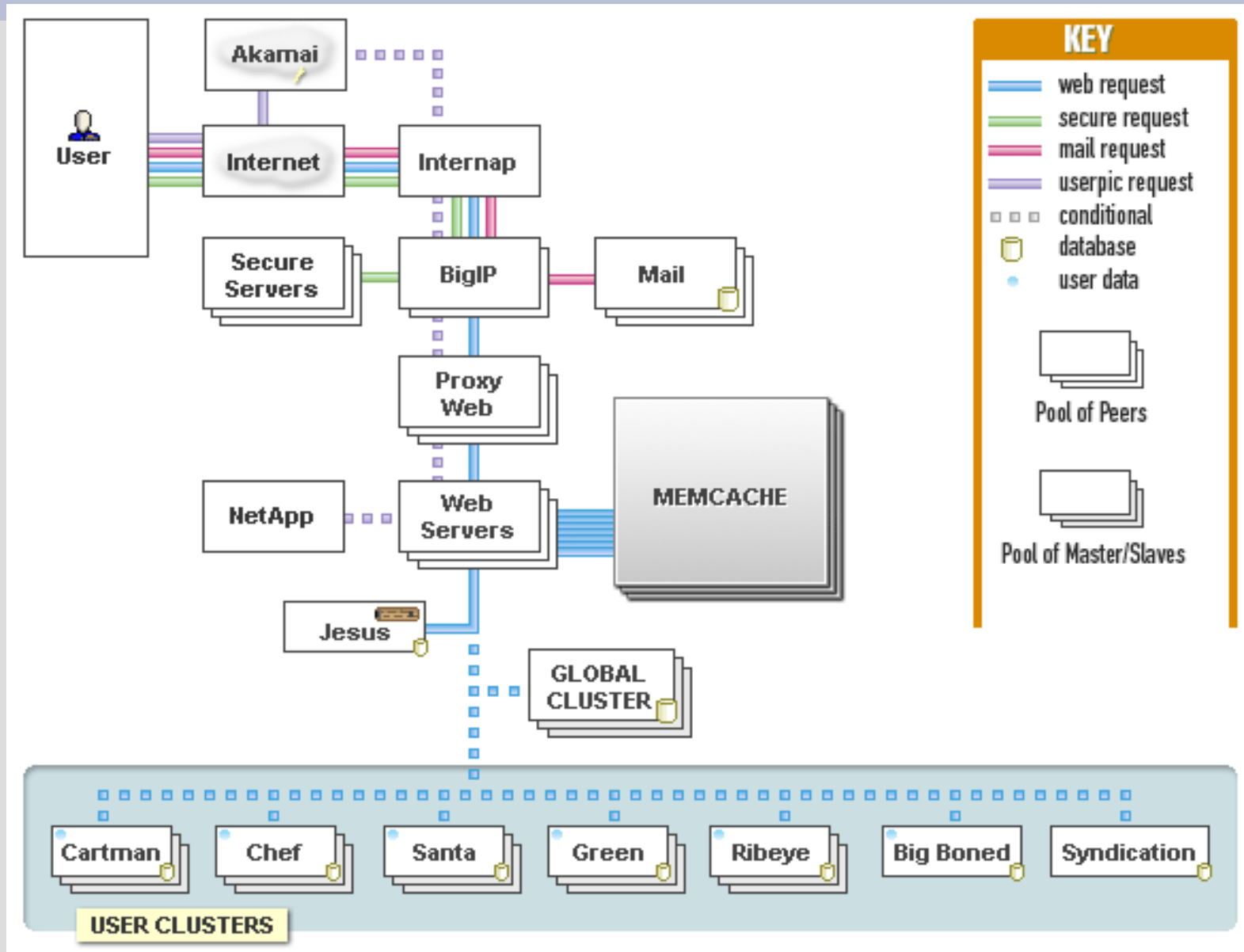
- LiveJournal overview
- Scaling history
- Perlbal
  - load balancer
- memcached
  - distributed caching
- MogileFS
  - distributed filesystem
- Wrap-up
  - Monitoring
  - Software/Architecture overview
- Future

# LiveJournal Overview

- college hobby project, Apr 1999
- blogging, forums
- aggregator, social-networking ('friends')
- 5+ million accounts; ~half active
- 50M+ dynamic page views/day. 1k+/s at peak hours (old data)
- why it's interesting to you...
  - 100+ servers
  - lots of Open Source:
    - existing open source:
      - Linux, Debian, Apache, perl, mod\_perl, MySQL, ...
    - our open source
      - memcached, perlbal, mogilefs, livejournal server, ...

# LiveJournal Backend

(as of a few months ago)

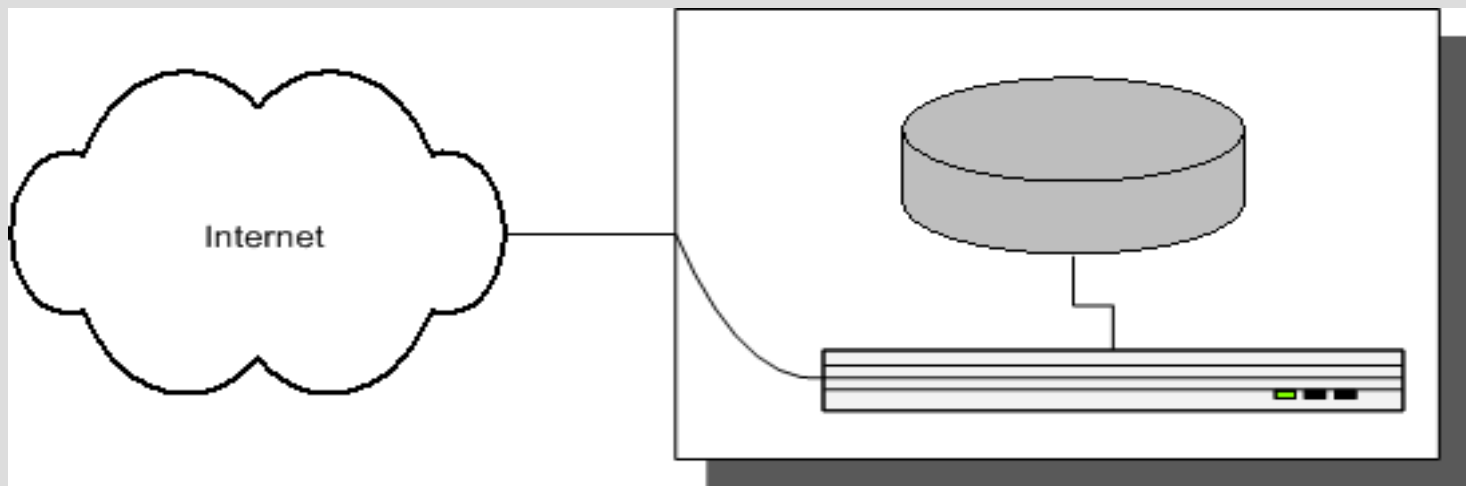


# Backend Evolution

- From 1 server to 100+....
  - where it hurts
  - how to fix
- Learn from this!
  - don't repeat my mistakes
  - can implement much of our design on a single server

# One Server

- shared server (killed it)
- dedicated server (killed it)
  - still hurting, but could tune it
  - learned Unix pretty quickly
  - CGI to FastCGI
- Simple



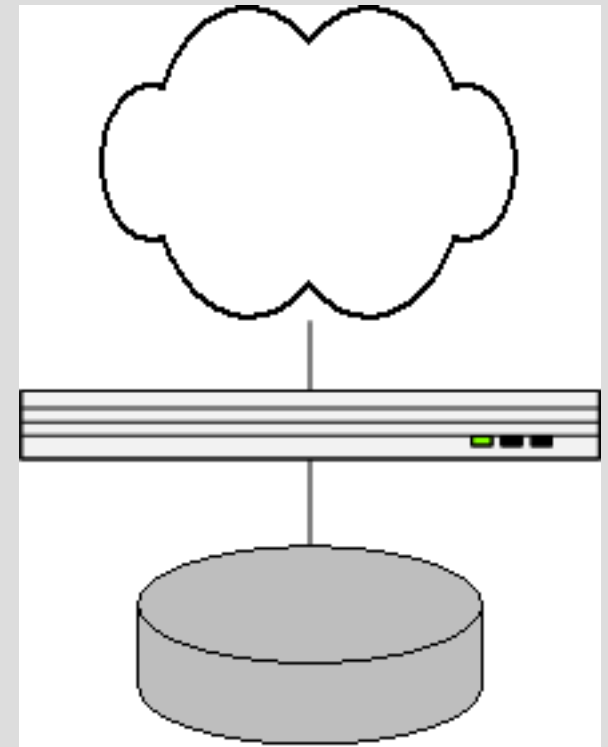
# One Server - Problems

- Site gets slow eventually.
  - reach point where tuning doesn't help
- single point of failure
- Need servers
  - start “paid accounts”



# Two Servers

- Paid account revenue buys:
  - Kenny: 6U Dell web server
  - Cartman: 6U Dell database server
    - bigger / extra disks
- Network simple
  - 2 NICs each
- Cartman runs MySQL on internal network

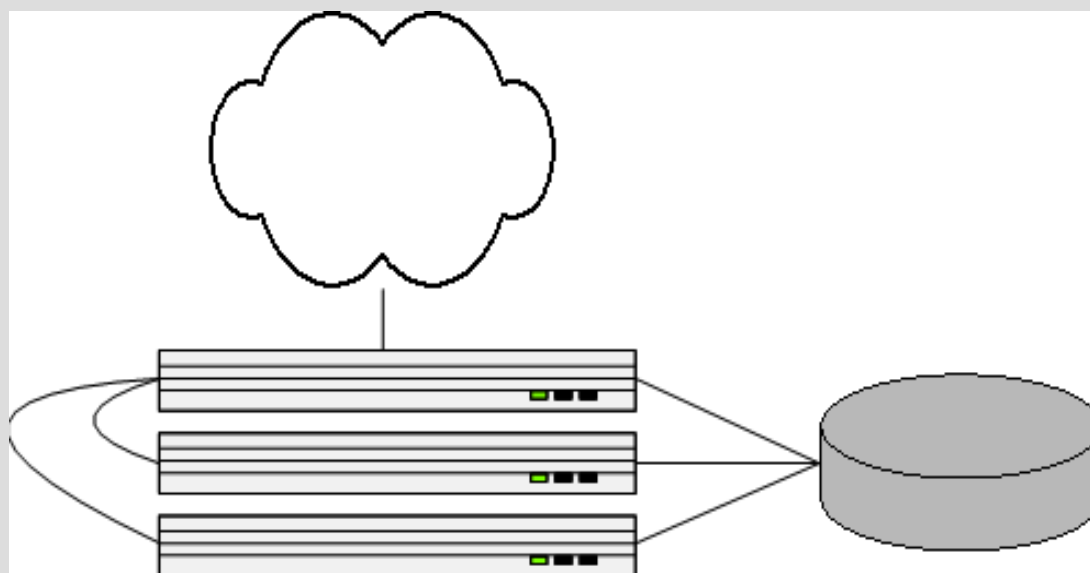


# Two Servers - Problems

- Two points of failure
- No hot or cold spares
- Site gets slow again.
  - CPU-bound on web node
  - need more web nodes...

# Four Servers

- Buy two more web nodes (1U this time)
  - Kyle, Stan
- Overview: 3 webs, 1 db
- Now we need to load-balance!
  - Kept Kenny as gateway to outside world
  - mod\_backhand amongst 'em all



# mod\_backhand

- web nodes broadcasting their state
  - free/busy apache children
  - system load
  - ...
- internally proxying requests around
  - network cheap

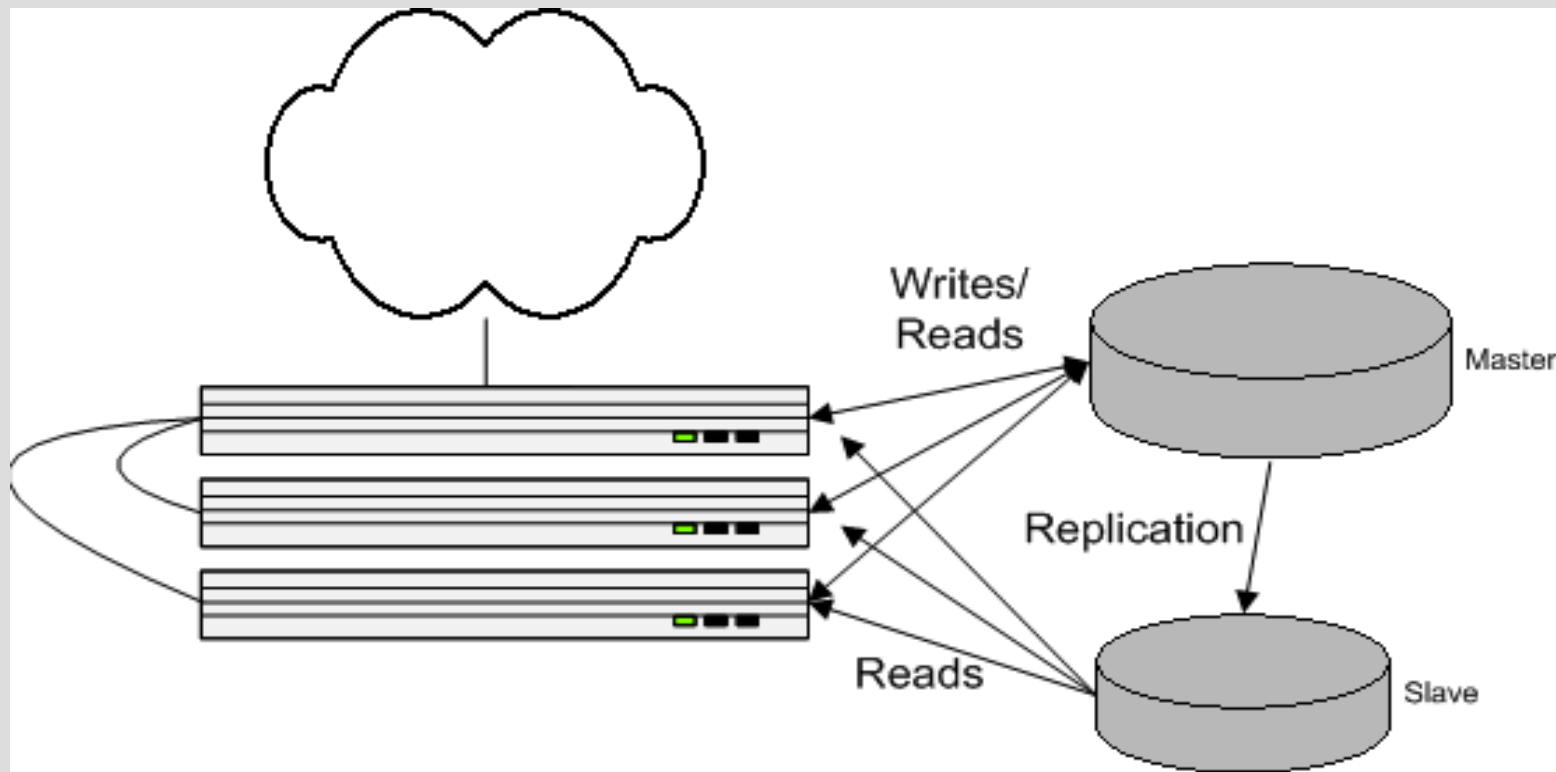
# Four Servers - Problems

- Points of failure:
  - database
  - kenny (but could switch to another gateway easily when needed, or used heartbeat, but we didn't)
- Site gets slow...
  - IO-bound
  - need another database server ...
  - ... how to use another database?

# Five Servers

## introducing MySQL replication

- We buy a new database server
- MySQL replication
- Writes to Cartman (master)
- Reads from both

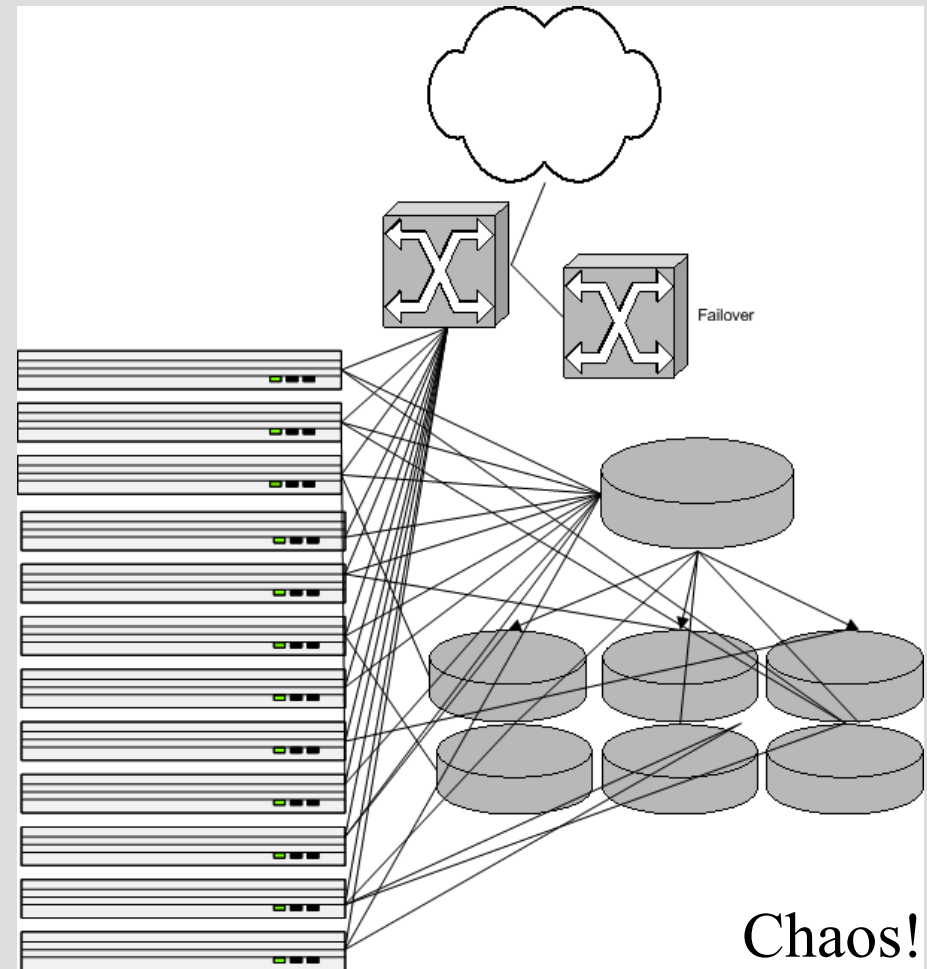


# Replication Implementation

- `get_db_handle() : $dbh`
  - existing
- `get_db_reader() : $dbr`
  - transition to this
  - weighted selection
- **permissions: slaves select-only**
  - mysql option for this now
- **be prepared for replication lag**
  - easy to detect in MySQL 4.x
  - user actions from `$dbh`, not `$dbr`

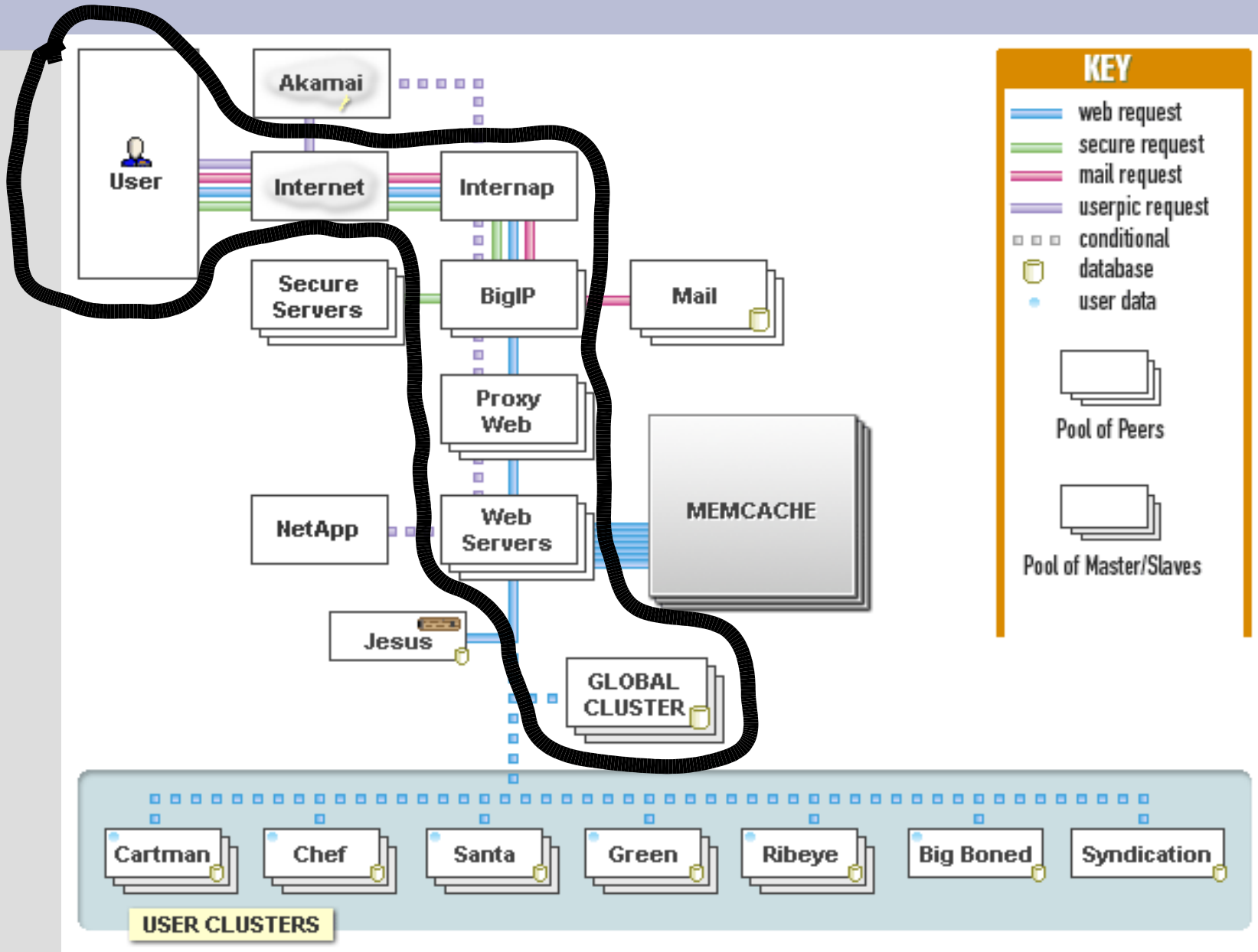
# More Servers

- Site's fast for a while,
- Then slow
- More web servers,
- More database slaves,
- ...
- IO vs CPU fight
- BIG-IP load balancers
  - cheap from usenet
  - LVS would work too
  - nowadays: wackamole





# Where we're at...

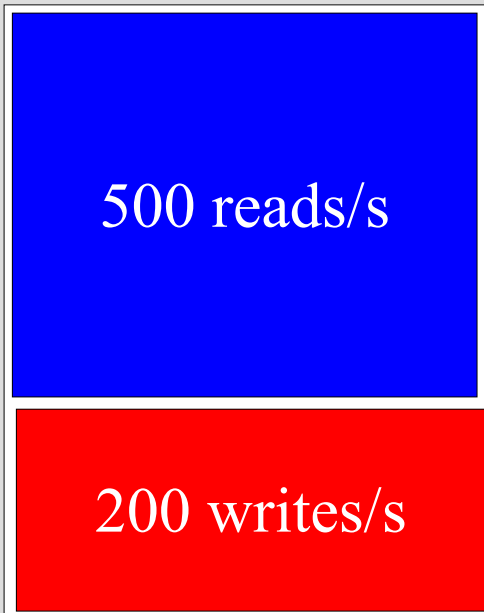


# Problems with Architecture

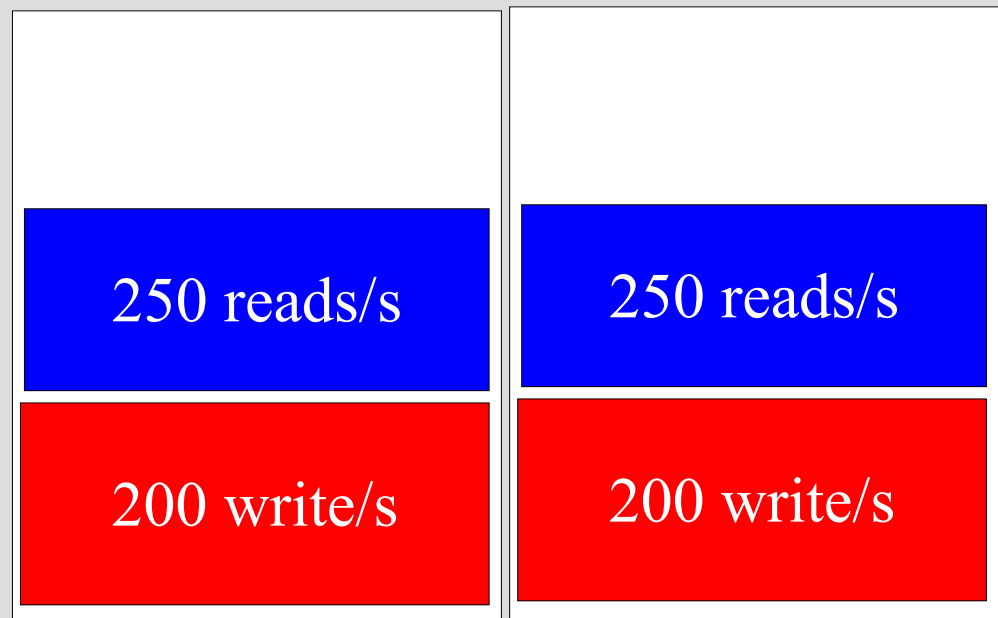
or,  
*"This don't scale..."*

- Slaves upon slaves doesn't scale well...
  - only spreads reads

w/ 1 server

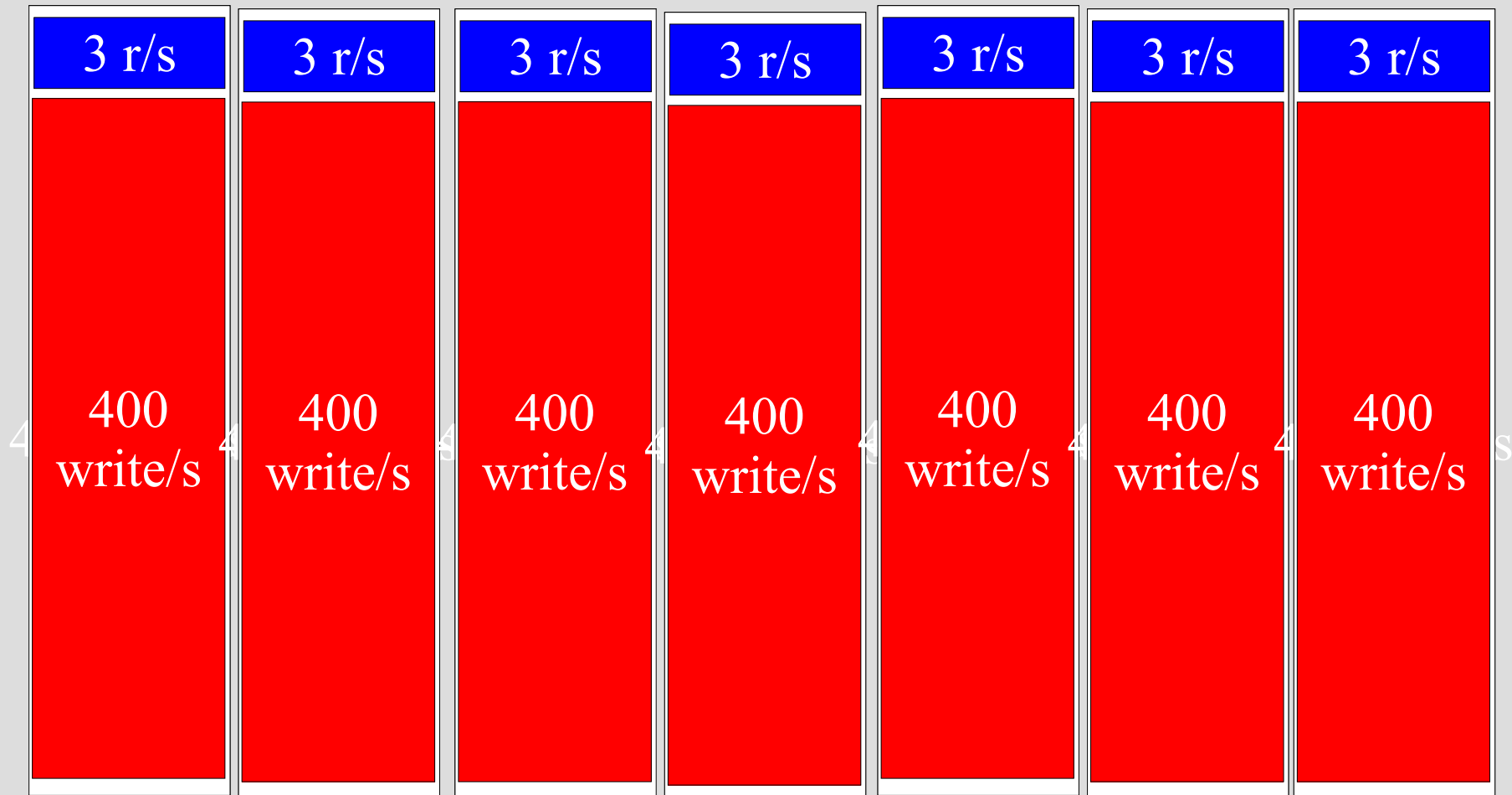


w/ 2 servers



# Eventually...

- databases eventual consumed by writing



# Not to mention,

- Database master is point of failure
- Reparenting slaves on master failure tricky at best
  - (without downtime)

# Spreading Writes

- Our database machines already did RAID
- We did backups
- So why put user data on 6+ slave machines?  
(~12+ disks)
  - overkill redundancy
  - wasting time writing everywhere

# Introducing User Clusters

- Already had `get_db_handle()` vs `get_db_reader()`
- Specialized handles:
- Partition dataset
  - can't join. don't care. never join user data w/ other user data
- Each user assigned to a cluster number
- Each cluster has multiple machines
  - writes self-contained in cluster (writing to 2-3 machines, not 6)

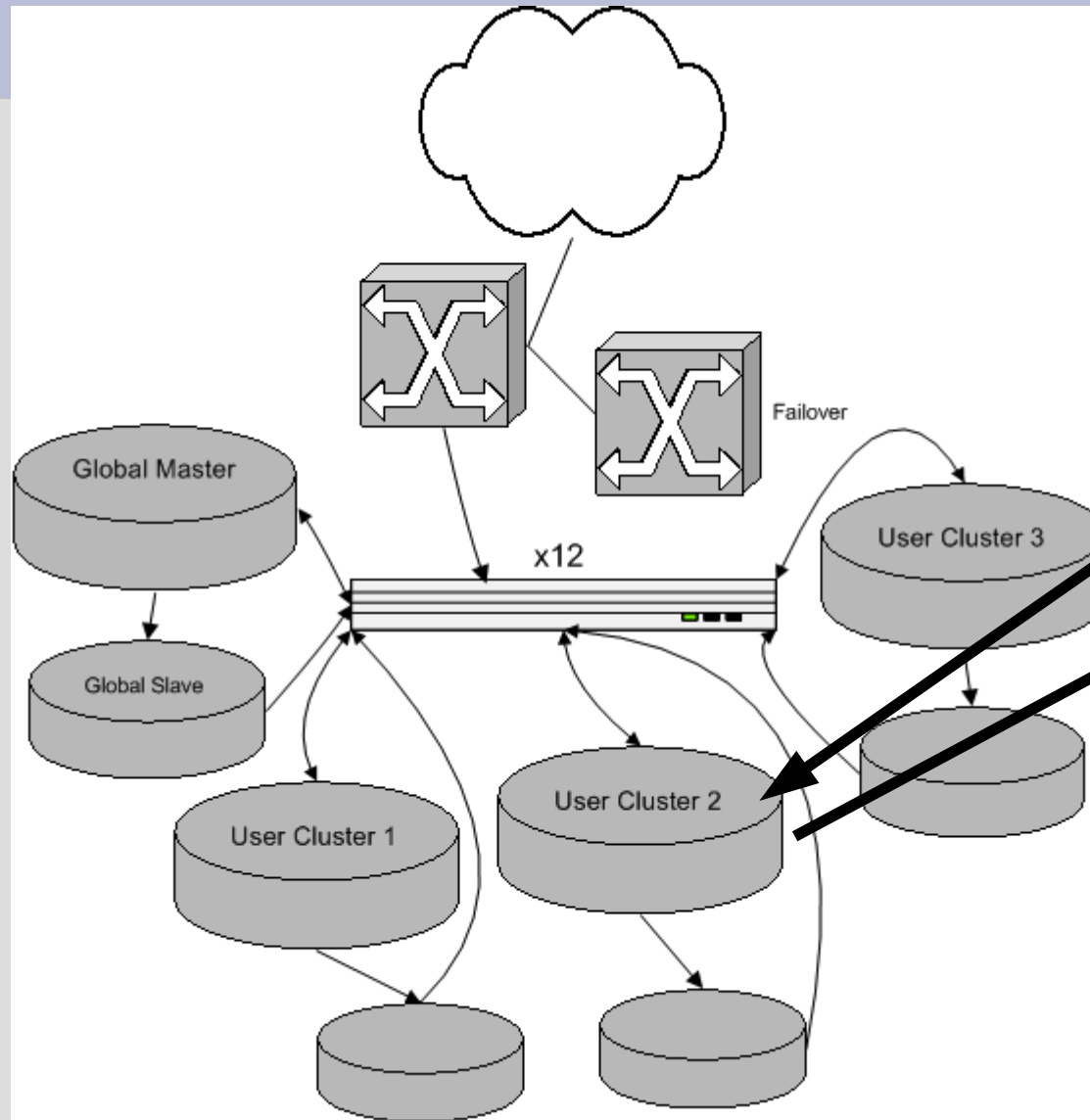
# User Cluster Implementation

- `$u = LJ::load_user("brad")`
  - hits global cluster
  - `$u` object contains its clusterid
- `$dbcm = LJ::get_cluster_master($u)`
  - writes
  - definitive reads
- `$dbcr = LJ::get_cluster_reader($u)`
  - reads

# User Clusters

```
SELECT userid,  
clusterid FROM  
user WHERE  
user='bob'
```

```
userid: 839  
clusterid: 2
```



```
SELECT ....  
FROM ...  
WHERE  
userid=839 ...
```

OMG i like  
totally hate  
my parents  
they just  
dont  
understand me  
and i h8 the  
world omg lol  
rofl \*! :^-  
^^;

- almost resembles today's architecture

add me as a  
friend!!!



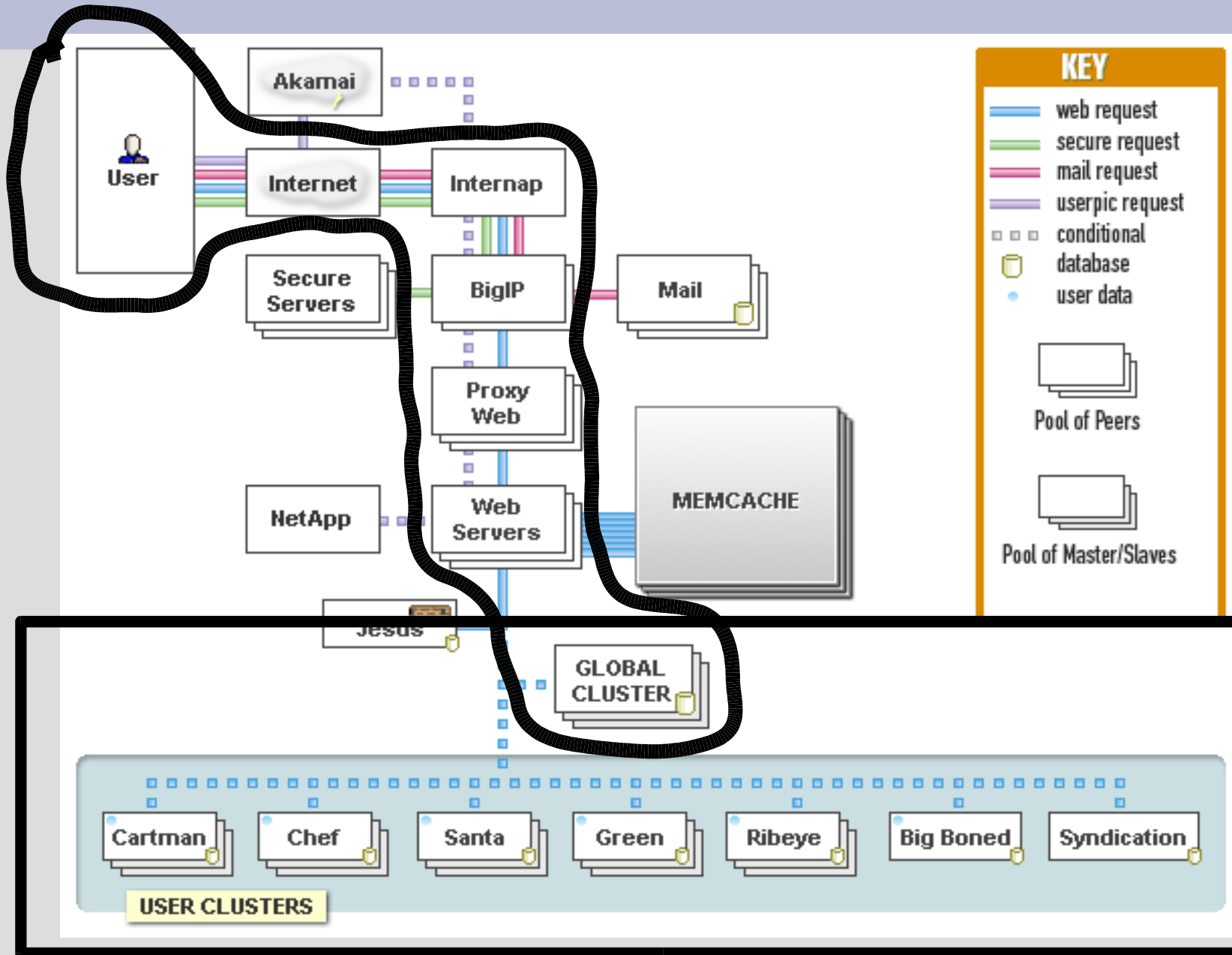
# User Cluster Implementation

- per-user numberspaces
  - can't use AUTO\_INCREMENT
  - avoid it also on final column in multi-col index: (MyISAM-only feature)
    - CREATE TABLE foo (uid INT, postid INT AUTO\_INCREMENT, PRIMARY KEY (userid, postid))
- moving users around clusters
  - very, very paranoid mover
  - user-moving harness
    - job server that coordinates, distributed long-lived user-mover clients who ask for tasks
  - balancing disk I/O
  - balance disk space
    - archive inactive users to space-efficient MyISAM

# DBI::Role – DB Load Balancing

- Our library on top of DBI
  - GPL; not packaged anywhere but our cvs
- Returns handles given a role name
  - master (writes), slave (reads)
  - directory (innodb), ...
  - cluster<n>{,slave,a,b}
  - Can cache connections within a request or forever
- Verifies connections from previous request
- Realtime balancing of DB nodes within a role
  - web / CLI interfaces (not part of library)
  - dynamic reweighting when node down

# Where we're at...



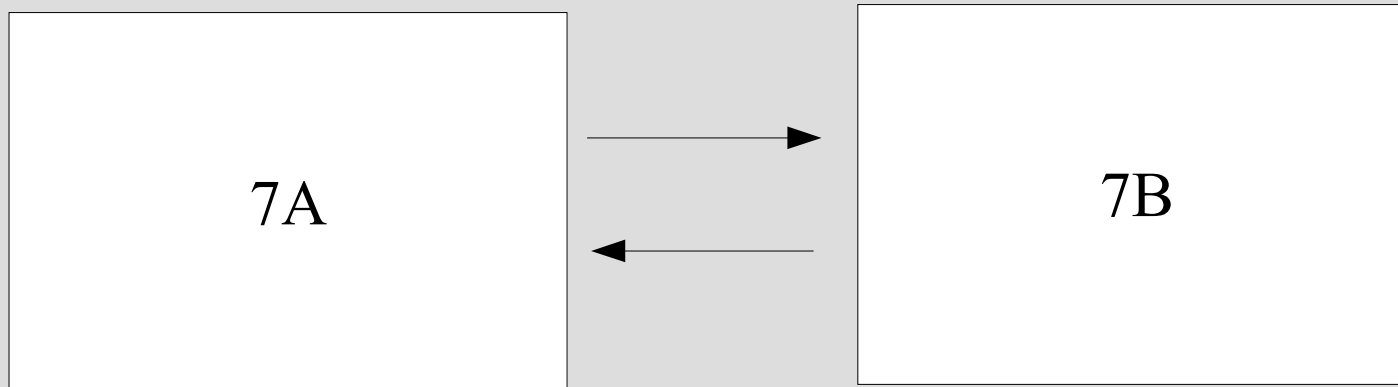
# Points of Failure

- 1 x Global master
  - lame
- $n$  x User cluster masters
  - $n$  x lame.
- Slave reliance
  - one dies, others reading too much

*Solution?*

# Master-Master Clusters!

- two identical machines per cluster
  - both “good” machines
- do all reads/writes to one at a time, both replicate from each other
- intentionally only use half our DB hardware at a time to be prepared for crashes
- easy maintenance by flipping active node
- backup from inactive node

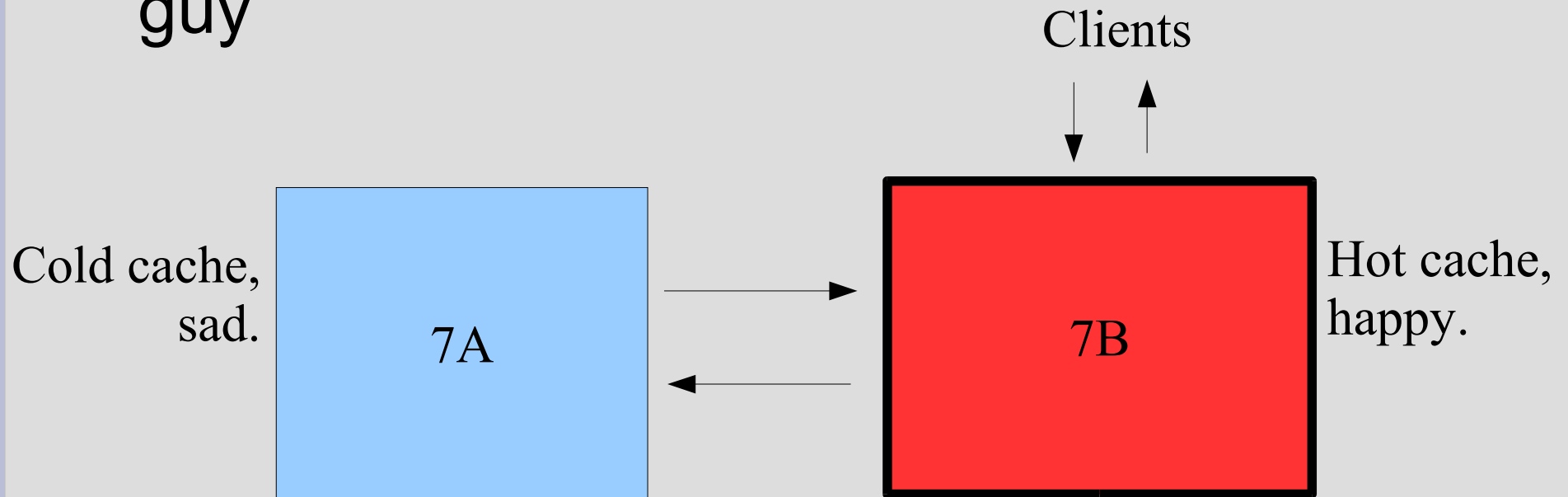


# Master-Master Prereqs

- failover can't break replication, be it:
  - automatic
    - be prepared for flapping
  - by hand
    - probably have other problems if swapping, don't need more breakage
- fun/tricky part is number allocation
  - same number allocated on both pairs
  - avoid AUTO\_INCREMENT
  - cross-replicate, explode.
  - do your own sequence generation w/ locking, 3<sup>rd</sup> party arbitrator, odd/even, centralized, etc...

# Cold Co-Master

- inactive pair isn't getting reads
- after switching active machine, caches full, but not useful (few min to hours)
- switch at night, or
- sniff reads on active pair, replay to inactive guy

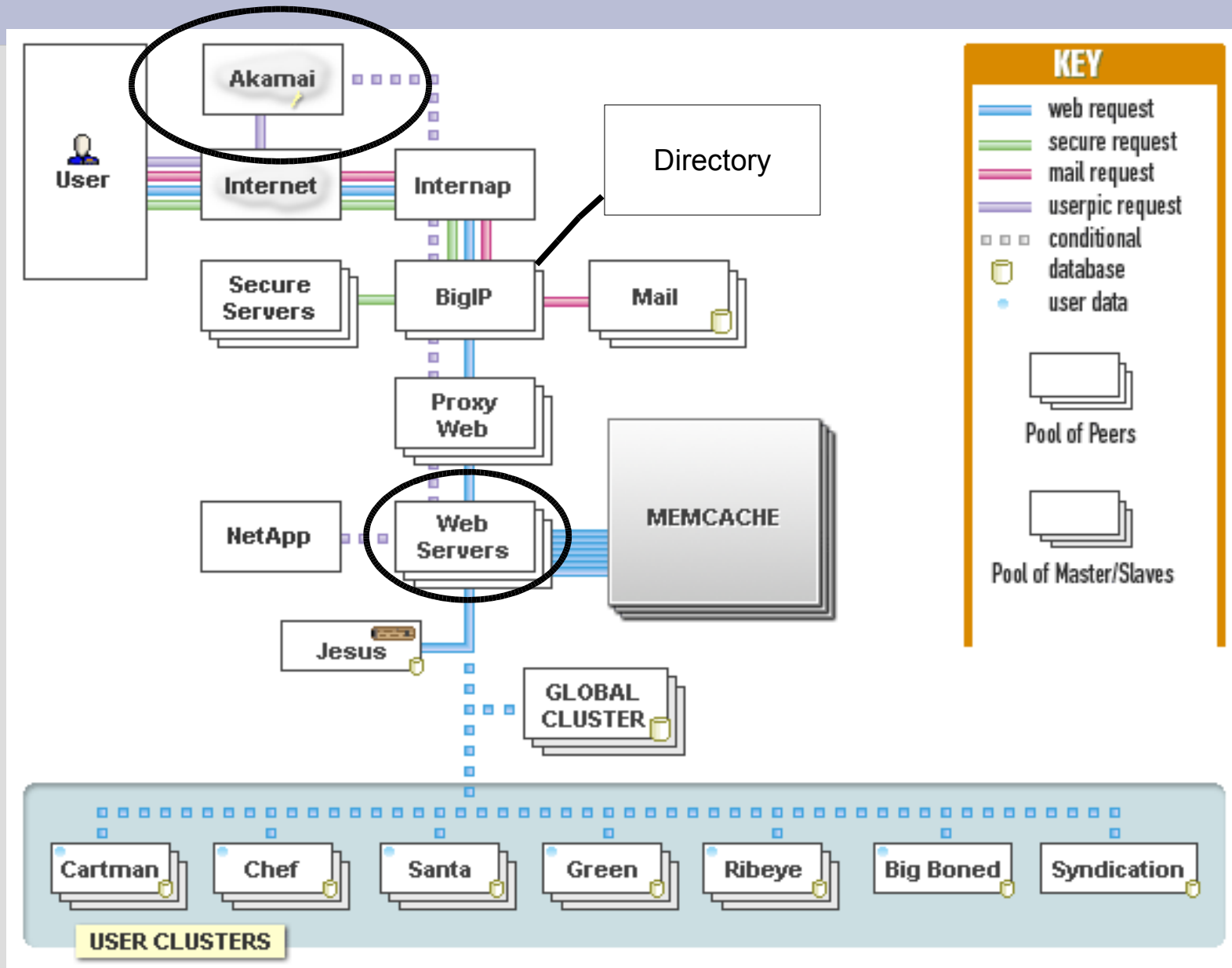


# Summary Thus Far

- Dual BIG-IPs (or LVS+heartbeat, or..)
- ~40 web servers
- 1 “global cluster”:
  - non-user/multi-user data
  - what user is where?
  - master-slave (lame)
    - point of failure; only cold spares
    - pretty small dataset (<4 GB)
      - future: MySQL Cluster!
        - in memory, shared-nothing, 99.999% uptime
- bunch of “user clusters”:
  - master-slave (old ones)
  - master-master (new ones)
- ...



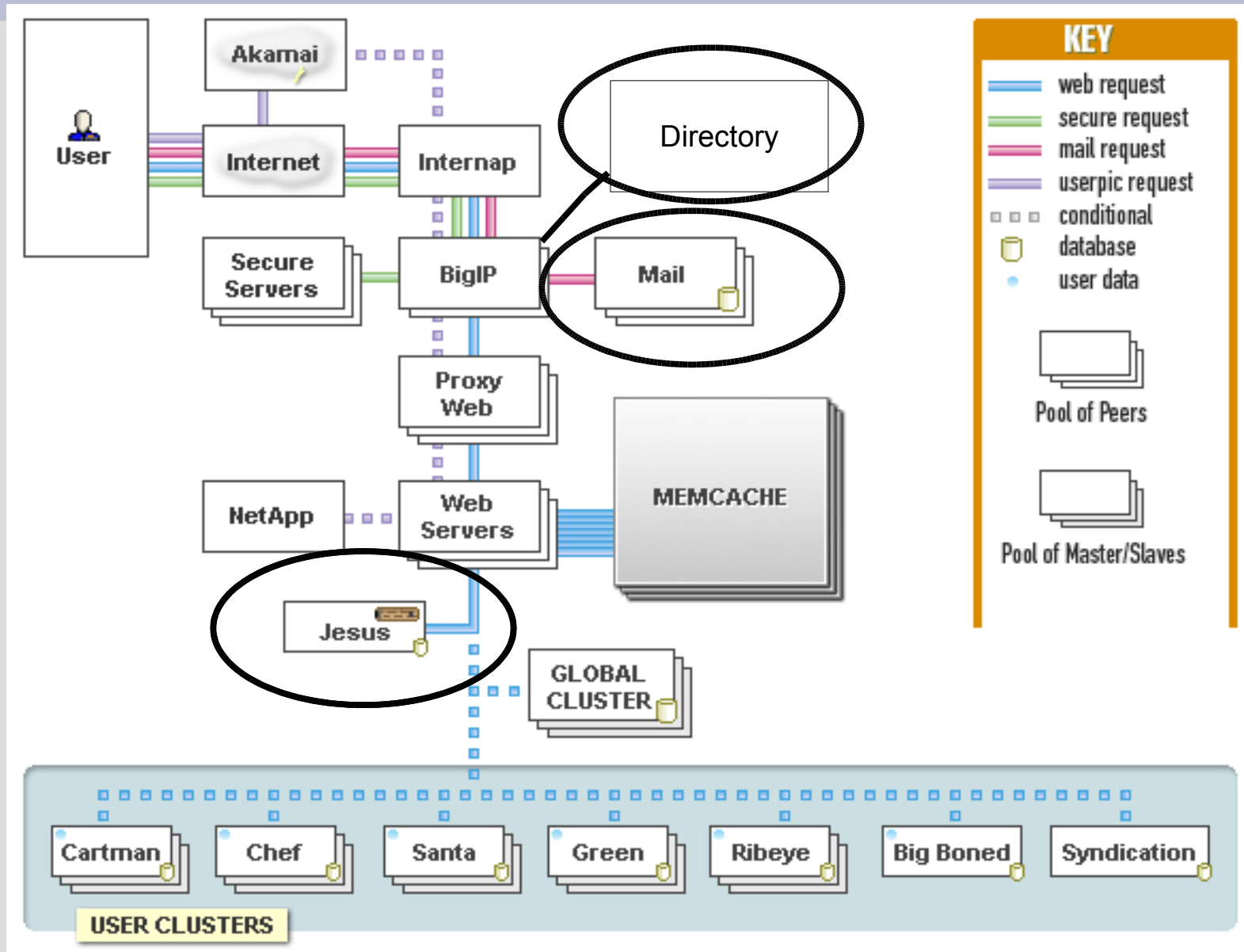
# Static files...



# Dynamic vs. Static Content

- static content
  - images, CSS
  - TUX, epoll-thttpd, etc. w/ thousands conns
  - boring, easy
- dynamic content
  - session-aware
    - site theme
    - browsing language
  - security on items
  - deal with heavy (memory hog) processes
  - exciting, harder

# Misc MySQL Machines



# MyISAM vs. InnoDB

- We use both
- MyISAM
  - fast for reading xor writing,
  - bad concurrency, compact,
  - no foreign keys, constraints, etc
  - easy to admin
  - logs
- InnoDB
  - ACID
  - wonderful concurrency
    - long slow queries while updates continue
    - directory server
  - detects hardware failures (bad memory/disks)

# MyISAM & InnoDB memory requirements

- index vs. data caching
- MyISAM caches indexes in memory
  - 32-bit machines w/ 8GB of memory: 2 GB of indexes in userspace, ~4GB of data cache in kernel buffer cache
- InnoDB
  - primary key is clustered index
  - indexes and data cached in userspace
  - kernel caching can be useless/harmful
    - O\_DIRECT helps a lot
      - caveat: linux 2.6 problems being fixed: XFS race, other filesystems have concurrency issues (one thread per file)
    - alternatively, raw partitions
  - begs for 64-bit

# MyISAM to InnoDB

- don't run both on same machine
  - InnoDB starves MyISAM disk-wise
  - separate caches which fight
    - one big cache better than two small ones
- MyISAM concurrency hack:
  - multiple dbs per machine. lame.

# Postfix & MySQL

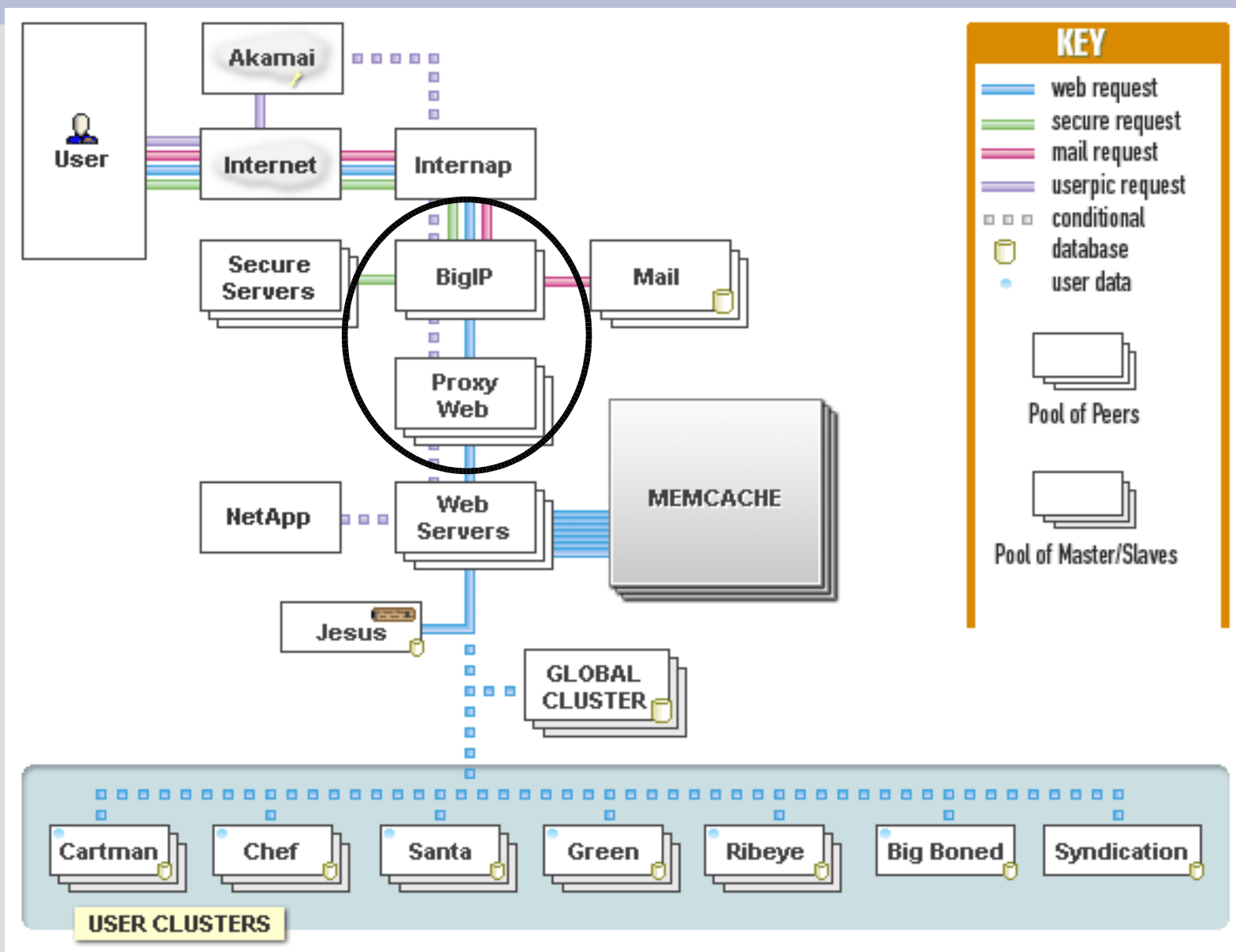
- 4 postfix servers
  - load balance incoming connections w/ BIG-IP
  - each runs tiny MySQL install
    - replicates one table (email\_aliases)
- Incoming mail uses mysql map type
  - To: brad@livejournal.com
  - SELECT email FROM email\_aliases WHERE alias='brad@livejournal.com'
- Don't have rebuild huge DBM files every few minutes

# Logging to MySQL

- mod\_perl logging handler
- new table per hour
  - MyISAM
- Apache access logging off
  - diskless web nodes, PXE boot
  - apache error logs through syslog-ng
- **INSERT DELAYED**
  - increase your insert buffer if querying
- minimal/no indexes
  - table scans are fine
- background job doing log analysis/rotation



# Load Balancing!



# Load Balancing Problem Overview

- slow clients (hogging mod\_perl/php)
  - even DSL/Cable is “slow”
  - need to spoon-feed clients
    - who will buffer?
- heterogeneous hardware and response latencies
  - load balancing algorithms
  - unlucky, clogged nodes
- dealing with backend failures
- The “Listen Backlog Problem”
  - is proxy/client talking to kernel or apache?
- live config changes

# Two proxy / load balancing layers

- 1: IP-level proxy
  - little or no buffering
    - 1 or 2 machines
      - hot spare, stateful failover
    - finite memory
  - Gbps+ switching
- 2: HTTP-level proxy
  - more machines
  - buffer here

# Proxy layer 1: IP-level

- Options:
  - Commercial:
    - BIG-IP, Alteon, Foundry, etc, etc...
  - Open Source:
    - Linux Virtual Server, Wackamole\*
- load balance methods:
  - round robin, weighted round robin
  - least connections
- some have L7 capabilities
  - useful, but still need another proxy layer...

# Proxy layer 2: HTTP-level

- Options:
  - mod\_proxy
    - “typical” setup with mod\_perl
    - to one host by default
    - mod\_rewrite + external map program (prg:) with mod\_proxy dest ([P])
      - broadcast Apache free/idle status from Apache scoreboard
      - flakey
    - “proxy connect error” to clients
  - pound
  - mod\_backhand
  - Squid
  - plb (pure load balancer)
- Frustrated, needy, we wrote our own...

# Perlbal

- Perl
- uses Linux 2.6's epoll
- single threaded, event-based
- console / HTTP remote management
  - live config changes
- handles dead nodes
- static webserver mode
  - sendfile(), async stat() / open()
- plug-ins
  - GIF/PNG altering
- ...

# Perlbal: Persistent Connections

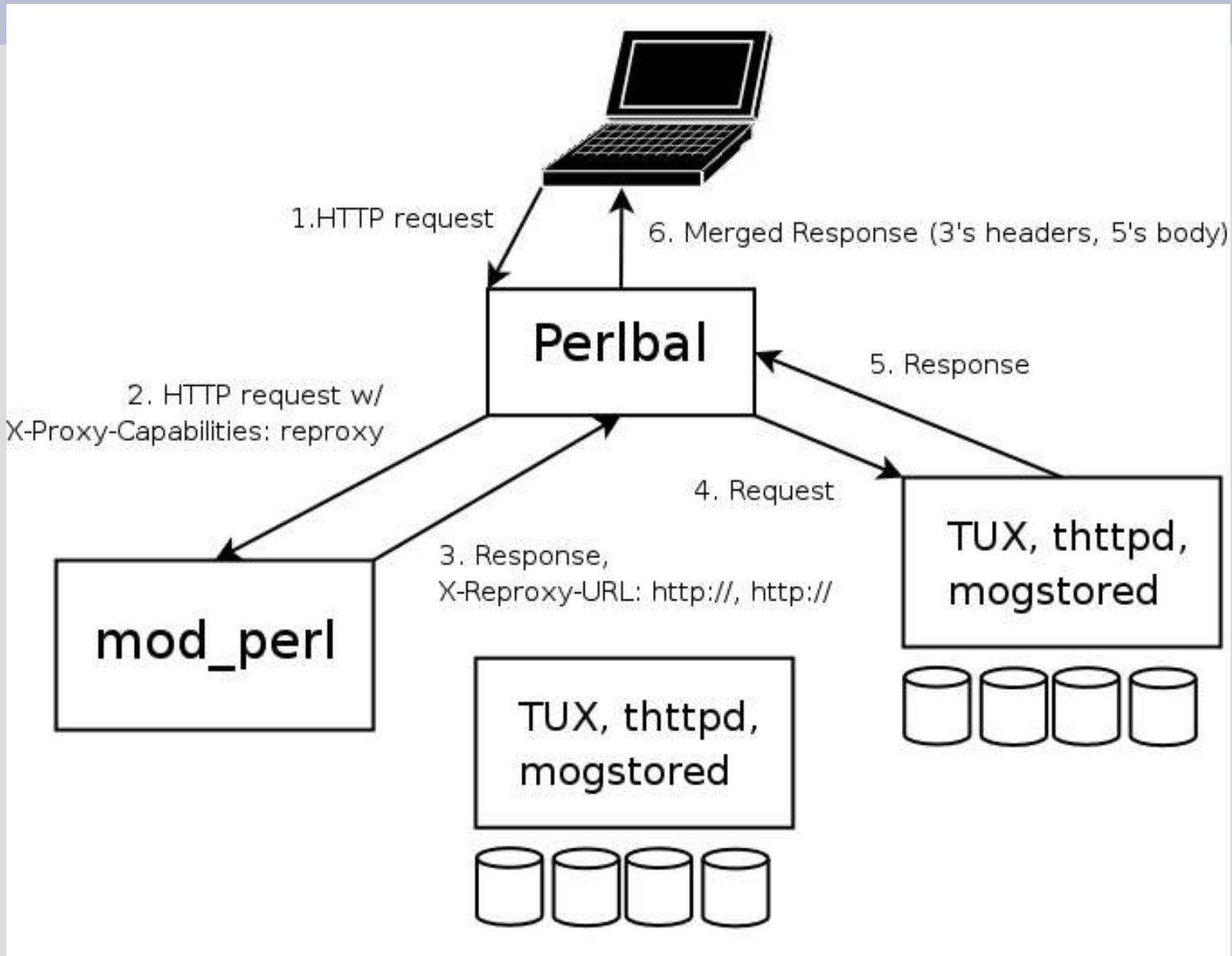
- persistent connections
  - perlbal to backends (mod\_perls)
  - know exactly when a connection is ready for a new request
    - keeps backends busy
    - connection known good
      - tied to mod\_perl, not kernel
- verifies new connections
  - one new pending connect per backend
  - verifies backend connection
    - OPTIONS request w/ keep-alive
    - response quick for apache
- multiple queues
  - free vs. paid user queues

# Perlbal: cooperative large file serving

- large file serving w/ mod\_perl bad...
  - buffering
- internal redirects
  - to URLs (plural) or file path
    - (hence Perlbal's web server mode)
  - client sees no HTTP redirect
- The path:
  - Perlbal advertises “X-Proxy-Capability: reproxy” to backend
  - backend (mod\_perl) does path trans & auth, sees proxy capability, sends URL/path back in header, not response
    - let mod\_perl do hard stuff, not push bytes around



# Internal redirect picture



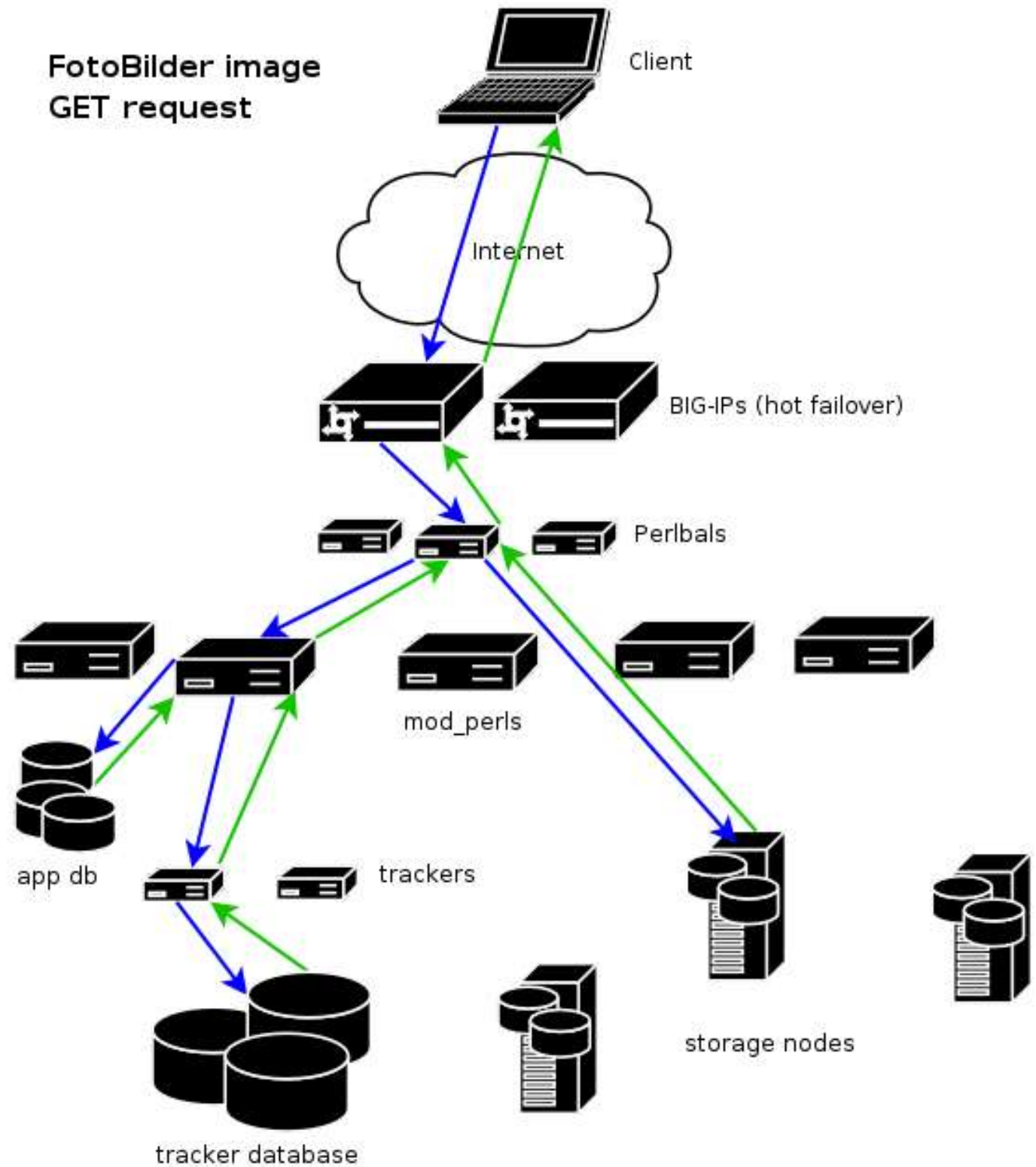
# MogileFS: distributed filesystem

- looked into Lustre, GFS, scared of in-development status
- MogileFS main ideas:
  - files belong to classes
    - classes: minimum replica counts (thumbnails == 1)
  - track what devices (disks) files are on
    - states: up, temp\_down, dead
  - keep replicas on devices on different hosts
    - Screw RAID! (for this, for databases it's good.)
  - multiple tracker databases
    - all share same MySQL cluster database
  - big, cheap disks (12 x 250GB SATA in 3U)
  - dumb storage nodes

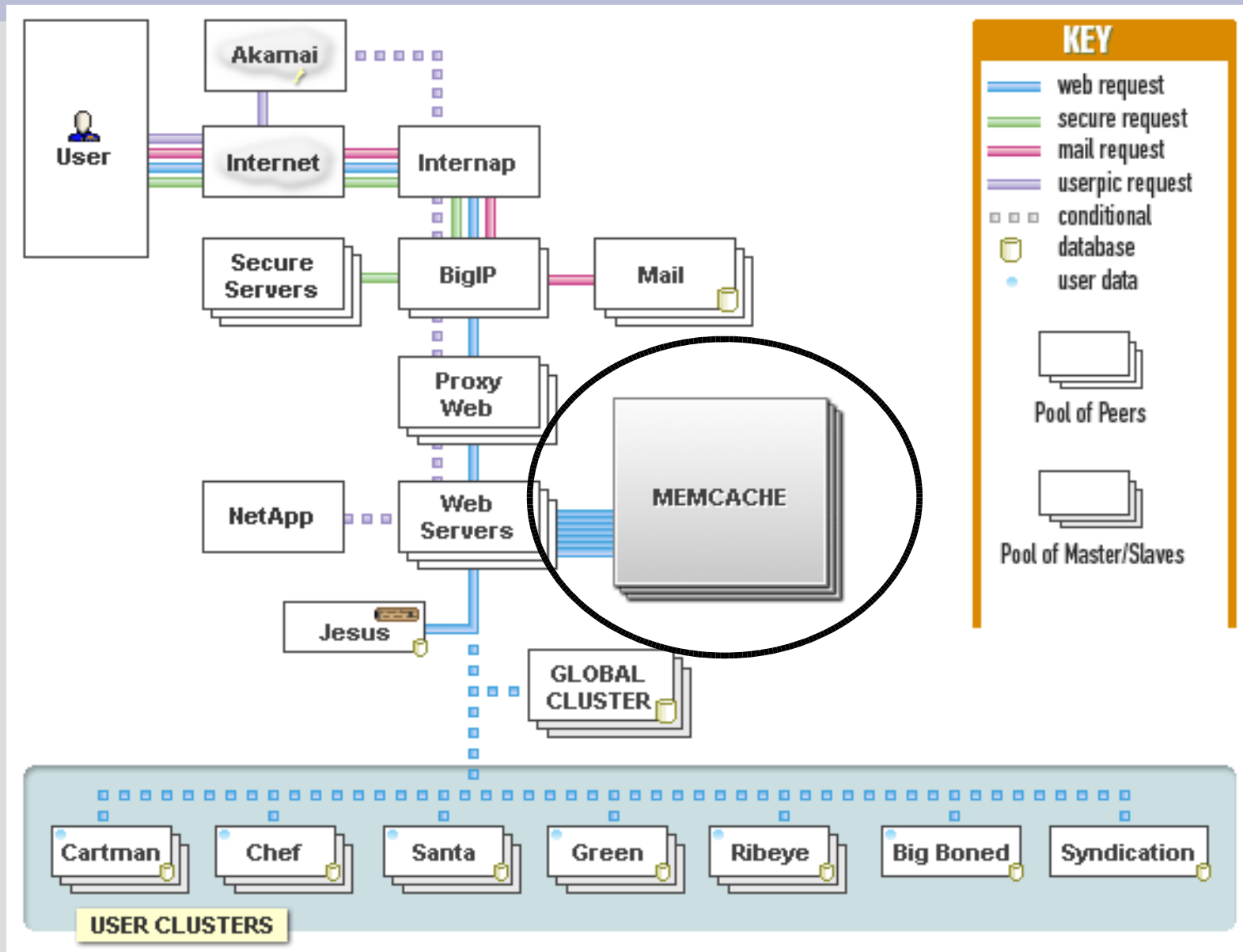
# MogileFS components

- clients
  - small, simple Perl library
  - FUSE filesystem driver (unfinished)
- trackers
  - interface between client protocol and MySQL Cluster
- MySQL Cluster
  - in memory, multiple machines
- Storage nodes
  - NFS or HTTP transport
    - [Linux] NFS *incredibly* problematic
  - HTTP transport is Perlbal with PUT & DELETE enabled

# Large file GET request



# Caching!



# Caching

- caching's key to performance
- can't hit the DB all the time
  - MyISAM: major r/w concurrency problems
  - InnoDB: good concurrency
    - not as fast as memory
  - MySQL has to parse your queries all the time
    - better with new MySQL 4.1 binary protocol
- Where to cache?
  - mod\_perl caching (address space per apache child)
  - shared memory (limited to single machine, same with Java/C#/Mono)
  - MySQL query cache: flushed per update, small max size
  - HEAP tables: fixed length rows, small max size

# memcached

<http://www.danga.com/memcached/>

- our Open Source, distributed caching system
- run instances wherever there's free memory
- no “master node”
- clients distribute requests
- In use by:
  - LiveJournal, Slashdot, Wikipedia, Meetup, mail systems, etc...
- protocol simple and XML-free; clients for:
  - perl, java, php(x3), python, ruby, C(?)...

# How memcached works

- requests hashed out amongst instance “buckets”
  - $\text{CRC32}(\text{“key”}) = 383472 \% \text{num\_buckets} = 6$
  - bucket 23 ... server 10.1.0.23: send: “key” = “value”

3 hosts, 7 buckets; 512 MB = 1 bucket (arbitrary)

10.1.0.18  
1024 MB; buckets 0-1

10.1.0.20  
2048 MB; buckets 2-5

10.1.0.23  
512 MB; bucket 6

weather = dismal  
tu:29323 = 1091029955

key = value



# memcached – speed

- C
  - prototype Perl version proved concept, too slow
- async IO, event-driven, single-threaded
- libevent (epoll, kqueue, select, poll...)
  - run-time mode selection
- lockless, refcounted objects
- slab allocator
  - glibc malloc died after 7~8 days
    - variable sized allocations, long life = difficult
  - slabs: no address space fragmentation ever.
- $O(1)$  operations
  - hash table, LRU cache
- multi-server parallel fetch (can't do in DBI)

# LiveJournal and memcached

- 10 unique hosts
  - none dedicated, whatever has extra memory
- 28 instances (512 MB = 1 bucket)
- 30 GB of cached data
- 90-93% hit rate
  - not necessarily 90-93% less queries:
    - FROM foo WHERE id IN (1, 2, 3)
    - would be 3 memcache hits; 1 mysql query
  - 90-93% potential disk seeks?
- 12 GB machine w/ five 2GB instances
  - left-over 'big' machines from our learn-to-scale-out days
- ~100,000 queries/second at peaks

# What to Cache

- Everything?
- Start with stuff that's hot
- Look at your logs
  - query log
  - update log
  - slow log
- Control MySQL logging at runtime
  - can't
    - (been bugging them)
  - sniff the queries! Net::Pcap
- count
  - add identifiers: `SELECT /* name=foo */`

# Caching Disadvantages

- more code
  - using
  - populating
  - invalidating
  - easy, if your API is clean
- conceptually lame
  - database should do it
    - kinda.
    - database doesn't know object lifetimes
      - putting memcached between app and DB doesn't work
- more stuff to admin
  - but memcached is easy
  - one real option: memory to use

# memcached TODO

- Very little
  - It Works.
  - have memcached processes w/ 190 day uptimes
- use tmpfs/ramfs? maybe.
  - eliminate 3GB limit on 32-bit machines
  - use even less CPU (sendfile from tmpfs)
- new memory allocator? maybe.
- virtual buckets & memcached bucket manager(s). definitely.
  - grow/shrink memcached farm at run-time
  - solves flapping problem if client re-hashes and can't detect old data

# MySQL Persistent Connection Woes

- connections == threads == memory
- max threads
  - limit max memory
- with 10 user clusters:
  - Bob is on cluster 5
  - Alice on cluster 6
  - Do you need Bob's DB handles alive while you process Alice's request?
- Major wins by disabling persistent conns
  - still use persistent memcached conns
  - db hits are rare (well, 14,000 queries/second)
  - mysql conns quick (opposed to, say, Oracle)
    - watch out for local port exhaustion

# Monitoring

- Cricket, Nagios
  - lots of custom Nagios plugins
- Interactive real-time tools...

lj@grimace:~\$

lj@grimace:~\$ dbcheck.pl

111 gm2	repl:	- <41>	conn:	0/ 214	UTC (master, oldids)
22 blue	111				()
17 gs1	111 repl:	121	conn:	1/ 55	PST (s2slave)
8 jesus	12 repl:	0	conn:	1/ 92	PST (slave)
48 mackey	111 repl:	0	conn:	0/ 1	PDT (slow)
18 red	111 repl:	16932	conn:	5/ 80	PST (slave)
14 bebe	111 repl:	0	conn:	0/ 24	PDT (email)
10 garrison	111 repl:	XXXXXXXX	conn:	0/ 1	PDT (directory)
7 gobbles	111 repl:	0	conn:	0/ 14	PST (email)
162 kyle	111 repl:	-	conn:	0/ 1	PDT ()
183 terrance	111 repl:	0	conn:	0/ 16	PDT ()
134 tender	135 repl:	1788 <17>	conn:	8/ 24	PDT (cluster10a)
135 loin	134 repl:	3858 <40>	conn:	15/ 25	PDT (cluster10b)
12 chef	repl:	- < 2>	conn:	0/ 16	PDT (cluster20, cluster20slave)
19 orange	12 repl:	0	conn:	1/ 9	UTC (cluster20slave)
21 bigboned	repl:	-	conn:	0/ 4	UTC (cluster7)
11 santa	repl:	- < 4>	conn:	0/ 15	PDT (cluster30, cluster30slave)
20 yellow	11 repl:	0	conn:	0/ 6	PDT (cluster30slave)
23 c04m	repl:	- <14>	conn:	1/ 20	PDT (cluster40, cluster40movemaster, cluster40slave)
24 green	23 repl:	0	conn:	3/ 10	UTC (cluster40slave)
26 indigo	25 repl:	0 < 5>	conn:	0/ 8	PST (cluster50b, cluster50slave)
25 ribeye	repl:	0 < 7>	conn:	4/ 19	PDT (cluster50, cluster50a, cluster50slave)
185 beef	184 repl:	0 <86>	conn:	5/ 49	UTC (cluster70b)
184 roast	185 repl:	510382 <82>	conn:	12/ 58	UTC (cluster70a)
112 pork	113 repl:	15269 < 7>	conn:	1/ 2	PST (cluster80a)
113 chop	112 repl:	0 < 9>	conn:	2/ 17	PST (cluster80b)
160 mad	161 repl:	11549 <22>	conn:	76/ 80	UTC (cluster90a)
161 cow	160 repl:	0 <26>	conn:	1/ 16	UTC (cluster90b)

## ERRORS:

\* Can't connect to blue

\* Slave not running: garrison:

lj@grimace:~\$ █



```
cow: 485.9 q/s
SUM: 10555.2 q/s

gm2: 1875.7 q/s
gs1: 603.1 q/s
jesus: 639.5 q/s
mackey: 95.7 q/s
red: 856.6 q/s
bebe: 114.9 q/s
garrison: 1.4 q/s
gobbles: 108.3 q/s
kyle: 0.6 q/s
terrance: 103.2 q/s
tender: 650.3 q/s
loin: 628.8 q/s
chef: 551.7 q/s
orange: 224.1 q/s
bigboned: 17.5 q/s
santa: 384.9 q/s
yellow: 192.6 q/s
c04m: 542.7 q/s
green: 227.6 q/s
indigo: 165.8 q/s
ribeye: 490.0 q/s
beef: 241.2 q/s
roast: 407.0 q/s
pork: 33.2 q/s
chop: 590.6 q/s
mad: 356.2 q/s
cow: 478.7 q/s
SUM: 10581.9 q/s
```

lj@grimace:~\$ mogcheck.pl

Checking mogilefsd availability...

10.0.0.81:7001 ... responding.

10.0.0.82:7001 ... responding.

Device information...

hostname	device	age	size(G)	used	free	use%	delay
sto1	dev1	42s	224,319	16,246	208,073	7.24%	0.005s
sto1	dev2	42s	229,161	10,641	218,520	4.64%	0.037s
sto1	dev3	42s	229,161	10,531	218,629	4.60%	0.023s
sto1	dev4	42s	229,161	10,586	218,575	4.62%	0.005s
sto1	dev5	42s	229,161	10,533	218,628	4.60%	0.004s
sto1	dev6	42s	229,161	10,657	218,503	4.65%	0.004s
sto1	dev7	42s	229,161	10,551	218,609	4.60%	0.006s
sto1	dev8	42s	229,161	10,619	218,542	4.63%	0.005s
sto1	dev9	42s	229,161	10,555	218,605	4.61%	0.007s
sto1	dev10	42s	229,161	10,474	218,687	4.57%	0.004s
sto1	dev11	42s	229,161	10,642	218,519	4.64%	0.005s
sto1	dev12	42s	229,161	10,520	218,641	4.59%	0.004s
sto1	dev13	42s	229,161	10,626	218,535	4.64%	0.005s
sto1	dev14	42s	229,161	10,558	218,603	4.61%	0.005s
sto2	dev15	37s	224,319	10,599	213,721	4.72%	0.006s
sto2	dev16	37s	229,161	10,569	218,592	4.61%	0.005s
sto2	dev17	37s	229,161	10,624	218,537	4.64%	0.005s
sto2	dev18	37s	229,161	10,675	218,486	4.66%	0.005s
sto2	dev19	37s	229,161	10,513	218,647	4.59%	0.012s
sto2	dev20	37s	229,161	10,619	218,542	4.63%	0.005s
sto2	dev21	37s	229,161	10,458	218,702	4.56%	0.004s
sto2	dev22	37s	229,161	10,597	218,563	4.62%	0.004s
sto2	dev23	37s	229,161	10,478	218,682	4.57%	0.004s
sto2	dev24	37s	229,161	10,630	218,531	4.64%	0.004s
sto2	dev25	37s	229,161	10,577	218,584	4.62%	0.004s
sto2	dev26	37s	229,161	10,522	218,638	4.59%	0.004s
sto2	dev27	37s	229,161	10,467	218,693	4.57%	0.023s
sto2	dev28	37s	229,161	10,575	218,585	4.61%	0.006s
	total		6406,817	301,643	6105,174	4.71%	0.211s

lj@grimace:~\$

```
10.0.0.29: free 1, active 14 c=111
10.0.0.28: free 7, active 8 c=146
10.0.0.27: free 5, active 8 c=156
10.0.0.23: free 0, active 15 c=131
10.0.0.22: free 1, active 18 c=152
----- Pool: int_web (@1100711245) ----- ct: 2977
10.0.0.97: free 7, active 11 c= 84
10.0.0.96: free 6, active 13 c=174
10.0.0.95: free 13, active 5 c= 81
10.0.0.94: ??
10.0.0.88: free 8, active 14 c=156
10.0.0.87: ??
10.0.0.86: free 9, active 8 c= 77
10.0.0.85: free 14, active 4 c= 58
10.0.0.84: free 13, active 5 c= 71
10.0.0.83: free 9, active 8 c= 65
10.0.0.80: free 7, active 13 c= 81
10.0.0.77: free 10, active 6 c= 73
10.0.0.76: free 12, active 8 c= 93
10.0.0.75: free 15, active 3 c= 28
10.0.0.74: free 11, active 6 c= 48
10.0.0.72: free 13, active 4 c= 37
10.0.0.71: free 8, active 7 c= 89
10.0.0.70: free 12, active 10 c= 77
10.0.0.69: free 8, active 11 c=113
10.0.0.64: free 12, active 9 c=100
10.0.0.63: free 9, active 11 c=155
10.0.0.62: free 11, active 6 c= 98
10.0.0.61: free 9, active 10 c= 73
10.0.0.59: free 9, active 11 c= 84
10.0.0.58: free 14, active 6 c=120
10.0.0.57: free 0, active 21 c=112
10.0.0.56: free 7, active 8 c= 60
10.0.0.53: free 7, active 7 c= 47
10.0.0.52: free 8, active 6 c= 76
10.0.0.51: free 7, active 7 c= 34
10.0.0.50: free 3, active 9 c= 83
10.0.0.42: free 2, active 13 c= 74
10.0.0.41: free 2, active 11 c= 71
10.0.0.29: free 1, active 13 c= 63
10.0.0.28: free 7, active 8 c= 94
10.0.0.27: free 9, active 5 c= 65
10.0.0.23: free 1, active 14 c= 68
10.0.0.22: free 1, active 18 c= 94
```



```
lj@grimace:~$ watch-queues.pl 1
```

```
Wed Nov 17 09:04:50 2004: [free users queued: 2, oldest: 0s] [paid users queued: 3, age: 1s]
Wed Nov 17 09:04:52 2004: [free users queued: 0, oldest: 0s] [paid users queued: 0, age: 0s]
Wed Nov 17 09:04:53 2004: [free users queued: 2, oldest: 0s] [paid users queued: 3, age: 0s]
Wed Nov 17 09:04:54 2004: [free users queued: 0, oldest: 0s] [paid users queued: 0, age: 0s]
Wed Nov 17 09:04:55 2004: [free users queued: 2, oldest: 0s] [paid users queued: 3, age: 0s]
Wed Nov 17 09:04:56 2004: [free users queued: 0, oldest: 0s] [paid users queued: 0, age: 0s]
Wed Nov 17 09:04:57 2004: [free users queued: 5, oldest: 0s] [paid users queued: 1, age: 0s]
```

```
lj@grimace:~$
```

```
lj@grimace:~$ watch-perlbal.pl 1
```

```
Wed Nov 17 09:05:03 2004: [w47:1 - 020, 0000, 1640] [w47:2 - 015, 0000, 1362] [w48:1 - 019, 0000, 1359] [w48:2 - 021, 0000, 1608] [w49:1 - 014, 0000, 1442] [w49:2 - 020, 0000, 1546] [w50:1 - 024, 0000, 1419] [w50:2 - 022, 0000, 1383] [w51:1 - 015, 0000, 1490] [w51:2 - 019, 0000, 1481]
Wed Nov 17 09:05:05 2004: [w47:1 - 024, 0000, 1383] [w47:2 - 024, 0001, 1148] [w48:1 - 023, 0000, 1131] [w48:2 - 017, 0000, 1330] [w49:1 - 027, 0001, 1247] [w49:2 - 016, 0000, 1310] [w50:1 - 029, 0000, 1223] [w50:2 - 015, 0000, 1110] [w51:1 - 029, 0000, 1242] [w51:2 - 023, 0000, 1291]
Wed Nov 17 09:05:06 2004: [w47:1 - 016, 0000, 1457] [w47:2 - 020, 0000, 1233] [w48:1 - 016, 0000, 1208] [w48:2 - 020, 0000, 1404] [w49:1 - 030, 0000, 1337] [w49:2 - 024, 0000, 1364] [w50:1 - 019, 0000, 1292] [w50:2 - 014, 0000, 1182] [w51:1 - 017, 0000, 1323] [w51:2 - 015, 0000, 1378]
Wed Nov 17 09:05:07 2004: [w47:1 - 022, 0000, 1543] [w47:2 - 027, 0000, 1322] [w48:1 - 027, 0000, 1283] [w48:2 - 033, 0000, 1479] [w49:1 - 030, 0000, 1395] [w49:2 - 027, 0000, 1459] [w50:1 - 020, 0000, 1375] [w50:2 - 022, 0000, 1256] [w51:1 - 023, 0000, 1409] [w51:2 - 028, 0003, 1449]
Wed Nov 17 09:05:08 2004: [w47:1 - 018, 0000, 1644] [w47:2 - 025, 0000, 1406] [w48:1 - 022, 0000, 1359] [w48:2 - 021, 0000, 1569] [w49:1 - 026, 0000, 1484] [w49:2 - 019, 0000, 1561] [w50:1 - 021, 0000, 1443] [w50:2 - 018, 0000, 1324] [w51:1 - 024, 0003, 1444] [w51:2 - 024, 0000, 1525]
Wed Nov 17 09:05:09 2004: [w47:1 - 026, 0000, 1359] [w47:2 - 024, 0000, 1150] [w48:1 - 019, 0000, 1110] [w48:2 - 027, 0000, 1253] [w49:1 - 025, 0000, 1207] [w49:2 - 026, 0000, 1322] [w50:1 - 031, 0000, 1171] [w50:2 - 022, 0000, 1078] [w51:1 - 019, 0000, 1190] [w51:2 - 025, 0000, 1273]
```

```
lj@grimace:~$ █
```

# Software Overview

- BIG-IPs
- Debian
  - Linux 2.4 (phasing out)
  - Linux 2.6
- mod\_perl
- MySQL
  - MyISAM, InnoDB
- Perlbal
- MogileFS
- Nagios, Cricket, ...

# Non-Technical Problems

- dealing w/ vendors
  - how much can they milk from you
  - fruit baskets
  - 6-month latency on returning calls, if ever
  - ... commoditize their stuff!
  - we like [siliconmechanics.com](http://siliconmechanics.com) (local, honest)
- asset management
  - `servers.yaml`
    - atrophied often until used it for generating configs, became useful and maintained
- incident logging
  - used to keep it in our head, then too many machines

# Misc Technical Problems

- few 64-bit issues
  - old MySQL codepaths (ISAM) from '97 not 64-bit safe
  - NUMA code crashing, XFS race, ...
- lame hardware raid
  - closed specs, hard to monitor
    - MegaRAID in Linux 2.6
  - prefer software except for battery-backed write-back caches
    - investigated solid state disks for ext3/xfs/innodb journals
- finding blocking (block-watcher.pl)
  - application notes latency on services, reports
  - lame, tedious (begs for DTrace)

# The Future

- finish MyISAM to InnoDB transition for user clusters
  - used to be “issues” in early days, but we're fairly happy now, esp. w/ 64-bit
- phase out old master-slave clusters
  - be fully master-master active/standby
- continue moving stuff off global DB
- MySQL Cluster or automatic master-election of 3 machines for global
  - MySQL Cluster very cool (distributed, in memory db), but the MySQL-NDB bridge immature





# Thank you!

Questions to...  
**brad@danga.com**  
**lisa@danga.com**

Slides linked off:  
**<http://www.danga.com/words/>**