

Some Duplication with Kyle's talk on Deployment (After this)



Hardware



Memory

MongoDB revolves around memory mapped files



Operating System map files on the Filesystem to Virtual Memory

- •(200 gigs of MongoDB files creates 200 gigs of virtual memory)
- OS controls what data in RAM
- When a piece of data isn't found, a page fault occurs (Expensive + Locking!)
- OS goes to disk to fetch the data
- Indexes are part of the Regular Database files
- Deployment Trick: Pre-Warm your Database (Pre-Warming your cache) to prevent cold start slowdown



MongoDB will take advantage of multiple cores

For working set queries, CPU usage is typically low



Full Tablescans

- Surprise: Queries which don't hit indexes make heavy use of CPU & Disk
- Deployment Trick: Avoid counting & computing on the fly by caching & precomputing data



Map Reduce

- Currently Single Threaded; runs in parallel across shards
- Deployment Trick: Use the new aggregation output options



Working set is crucial!!!

- Working set should be, as much as possible, in memory
- your entire dataset need not be



Disks & I/O

 Disk I/O becomes your definer of performance in nonworking set queries



Surprise: Faster Disks is better than slow disks. More is also better

- RAID is good for a variety of reasons
- Our Recommendations ...



RAID 10 (Mirrored sets inside a striped set; minimum 4 disks)

- Improved write performance
- Survives single disk failure
- Downside: Needs double storage needs
- e.g. 4 20 gig disks gives you 40 gigs of usable space
- LVM of RAID 10 on EBS seems to smooth out performance and reliability best for MongoDB



RAID 5 or 6

- 1 or 2 additional disks required for parity
- Can survive 1 or 2 disk failures
- Implementations seem inconsistent, buyer beware



Flash (SSD)

- Expensive, but getting cheaper
- Significantly reduced seek time and increased I/O
 Throughput
- Random Writes and Sequential Reads are still a weak point



OS

- •For production: Use a 64 bit OS and a 64 bit MongoDB Build
- 32 Bit has a 2 gig limit; imposed by the operating systems for memory mapped files
- Clients can be 32 bit
- MongoDB Supports (little endian only)
- Linux, FreeBSD, OS X (on Intel, not PowerPC)
- Windows
- Solaris (Intel only, Joyent offers a cloud service which works for Mongo)



Put your journal on a separate spindle if possible



Server Status

Some tools for examining server status



MongoStat - free tool which comes with MongoDB

Shows I/O counters, time spent in locks, etc.

11:37:21

11:37:22

11:37:23

11:37:24

11:37:25

0 0 0

0 0 0

0 0 0

0 0 0

0 0 0

```
56694
                                                                                                                   57360
                                                                                                                             688
                                                                                   0
                                                                                                            56694
                                                                                                                   57360
                                                                                                                             688
                                                                                                                                         0
                                                                                                                                                0.3
                                                                                                            56694
                                                                                                                   57360
                                                                                                                             688
                                                    14
                                                                                                            56694
                                                                                                                   57360
                                                                                                                             688
                                                                                                                                         0
                                                                                                                                                0.4
                                                    14
                                                                                                            56694
                                                                                                                   57360
                                                                                                                             688
                                                                                                            56694
                                                                                                                   57360
                                                                                                                             688
aults/s locked % idx miss %
                                     a tiriw
                                               conn
                                                                                              1
                                                                        0
                                                                                                            56694
                                                                                                                   57360
                                                                                                                             688
                                                                                                                                         0
                                       0 0 0
                                                        11:37:13
                                                  48
                                                                        0
                                                                                                            56694
                                                                                                                   57360
                                                                                                                             688
                                                                                                                                         0
               0.3
                                       0 0 0
                                                        11:37:14
                                                                        0
                                                                                                            56694
                                                                                                                   57360
                                                                                                                             688
                                       0 0 0
                                                  48
                                                        11:37:15
                                                                                                            56694
                                                                                                                   57360
                                                                                                                             688
                                       0 0 0
                                                  48
                                                       11:37:16
                                                                        0
                                                                                                            56694
                                                                                                                   57360
                                                                                                                             688
                                       0 0 0
                                                       11:37:17
                                                  48
                                                                                                                   57360
                                                                                                                             688
                                       0 0 0
                                                                                                                   57360
                                                                                                                             688
                                                                                                                                         0
                                                  48
                                                       11:37:18
                                       0 0 0
                  0
                                                  48
                                                       11:37:19
       0
                  0
                                       0 0 0
                                                  48
                                                       11:37:20
```

insert/s query/s update/s delete/s getmore/s command/s flushes/s

vsize

mapped

faults/s locked %

Similarly, iostat ships on most Linux machines (or can be installed)

- •iostat [args] <seconds per poll>
- -x for extended report
- Disk can be a bottleneck in large datasets where working set > ram
- •~200-300Mb/s on XL EC2 instances, but YMMV (EBS is slower)
- On Amazon Latency spikes are common, 400-600ms
 (No, this is not a good thing)



avg-cpu:				%iowait 0.35							OSt	at
Device:		rrqm/s	wrqm/s	r/s	w/s	rsec/s	wsec/s	avgrq-sz	avgqu-sz	await	svctm	%util
sda1		1.22	3.84	0.70	2.87	34.46	53.67	24.68	0.05	15.04	4.80	1.71
sda2		0.01	0.01	0.00	0.00	0.08	0.09	31.59	0.00	27.94	12.18	0.01
	N	0/	O/	W4 4 to	4 L	0/4 47 -	- - KB/t	disk0 tps MB/s	diskl KB/t tps	MB/s	disk2 KB/t tps	cpu MB/s



db.serverStatus()

```
> db.serverStatus()
          "host": "b-mac",
          "version": "1.8.0",
          "process" : "mongod",
          "uptime" : 845952,
          "uptimeEstimate" : 355978,
          "localTime": ISODate("2011-04-26T17:04:40.938Z"),
          "globalLock": {
                    "totalTime": 845952742184,
                    "lockTime" : 5154824,
                    "ratio": 0.0000060935129623101254,
                    "currentQueue": {
                               "total" : o,
                               "readers" : o,
                               "writers" : o
                    },
                    "activeClients" : {
                               "total" : o,
                               "readers": o,
                               "writers": o
```

"writers" : o

db, serverStatus()

```
"mem" : {
           "bits": 64,
           "resident" : 7,
           "virtual" : 28917,
           "supported": true,
           "mapped" : 13245
},
"connections" : {
           "current": 1,
           "available" : 203
},
"extra_info" : {
           "note": "fields vary by platform"
},
"indexCounters" : {
           "btree" : {
                      "accesses" : 37,
                      "hits" : 37,
                      "misses": o,
                      "resets": o,
                      "missRatio": o
```



"misses": 0, db.serverS resets": 0, missRatio

```
"backgroundFlushing" : {
          "flushes" : 5999,
          "total_ms" : 83634,
          "average_ms" : 13.941323553925654,
          "last_ms" : 2,
          "last_finished" : ISODate("2011-04-26T17:04:02.582Z")
"cursors" : {
          "totalOpen": o,
          "clientCursors_size" : o,
          "timedOut": 21
},
"network" : {
          "bytesIn" : 1614262,
          "bytesOut" : 19368250,
          "numRequests" : 6395
"opcounters" : {
          "insert" : 2359,
```

```
"getmore": 2776,
db, serverStatus()
      "asserts" : {
               "regular" : o,
               "warning": o,
               "msg": o,
               "user": 8,
               "rollovers" : o
     },
      "writeBacksQueued": false,
      "dur" : {
               "commits": o,
               "journaledMB": o,
               "writeToDataFilesMB" : o,
               "commitsInWriteLock": o,
               "earlyCommits": o,
               "timeMs": {
                         "dt" : 3008,
                         "prepLogBuffer" : o,
                         "writeToJournal": o,
                         "writeToDataFiles": o,
                         "remapPrivateView": o
```

7

Filesystems



All data & namespace files are stored in the 'data' directory (- dbpath) •You can create symbolic links to keep different databases

- You can create symbolic links to keep different databases on different disks
- Best to aggregate your IO across multiple disks
- File Allocation

```
$ 1s -sk /data/db/
16384 foo.ns
double in size { 65536 foo.0 (up to 2 gigs) { 131072 foo.1 16384 bar.ns ....
```



Logfiles

- •--logpath <file>
- Rotation can be requested of MongoDB...
- •db.runCommand("logRotate")
- •kill -SIGUSR1 <mongod pid>
- killall -SIGUSR1 mongod
- •Won't work for ./mongod > [file] syntax



Filesystems

- MongoDB is filesystem neutral
- ext3, ext4 and XFS are most used
- •BUT....
- ext4, XFS or any other filesystem with posix_fallocate() are preferred and best



EC2

- Many distros default to ext3 (but Amazon AMI now uses ext4 by default)
- For best performance reformat to EXT4 / XFS
- Make sure you use a recent version of EXT4
- Striping (MDADM / LVM) aggregates I/O
- See previous recommendations about RAID 10



Monitoring



Maintenance

- When doing a lot of updates or deletes....
- Compaction may be needed occasionally on indices and datafiles
- •db.repair()
- •Replica Sets:
- Rolling repairs, start nodes up with --repair param



Scale out

read

shard1

rep_a1

rep_b1

rep_c2

shard2

rep_a2

rep_b2

rep_c2

shard3

rep_a3

rep_b3

rep_c3

mongos / config server

mongos / config server

mongos / config server

write

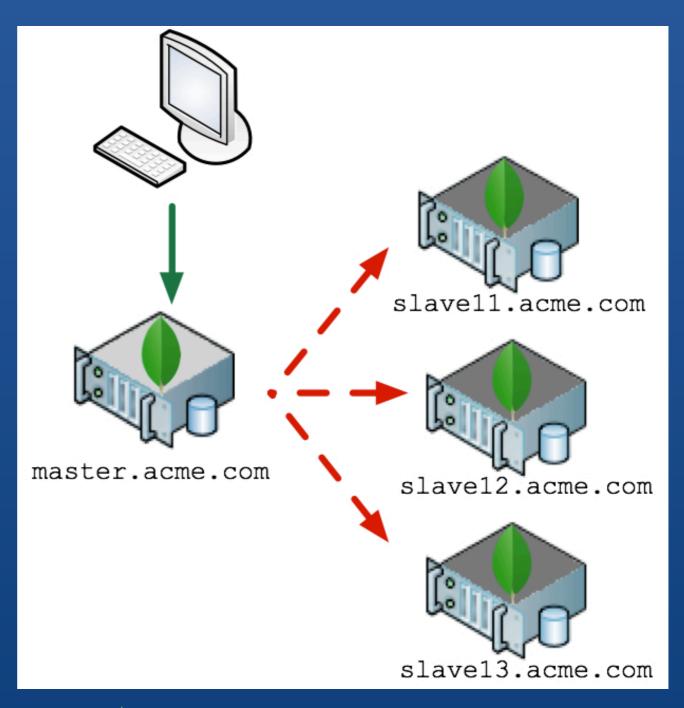


Backups



Best driven from a slave

- •Eliminates impact on master during backup
- Hidden Nodes in1.8





mongodump / mongorestore

- binary, compact object dump
- each consistent object is written
- NOT necessarily consistent from start to finish (Unless you lock the database)
- mongorestore to restore binary dump
- database doesn't have to be up to restore, can use dbpath



filelock / fsync

- lock: blocks writes
- •db.runCommand({fsync: 1, lock: 1})
- fsync to flush buffers to disk
- backup
- •then, unlock
- •db.\$cmd.sys.unlock.findOne();



With Journaling, you can run an LVM or EBS snapshot and recover later without locking

- EBS Can disappear (See: last week)
- S3 for longer term backups
- USE AMAZON AVAILABILITY ZONES
- •DR / HA





We're Hiring!
brendan@10gen.com
(twitter: @rit)

conferences, appearances, and meetups http://www.10gen.com/events







LinkedIn
http://linkd.in/joinmongo

