

# Deep Dive into MongoDB from a developer's point of view

# Craig Lawrence

- Software Developer/Lead/Manager in cable/telecom for over 25 years
- Bachelor and Master's degree in Electrical Engineering
- Worked at Ericsson on OpenStream and MDMS since 2008
- Using MongoDB as a developer since 2011

# What's New in 2014?

- Browser: Chrome
- OS: Microsoft
- Device: Android

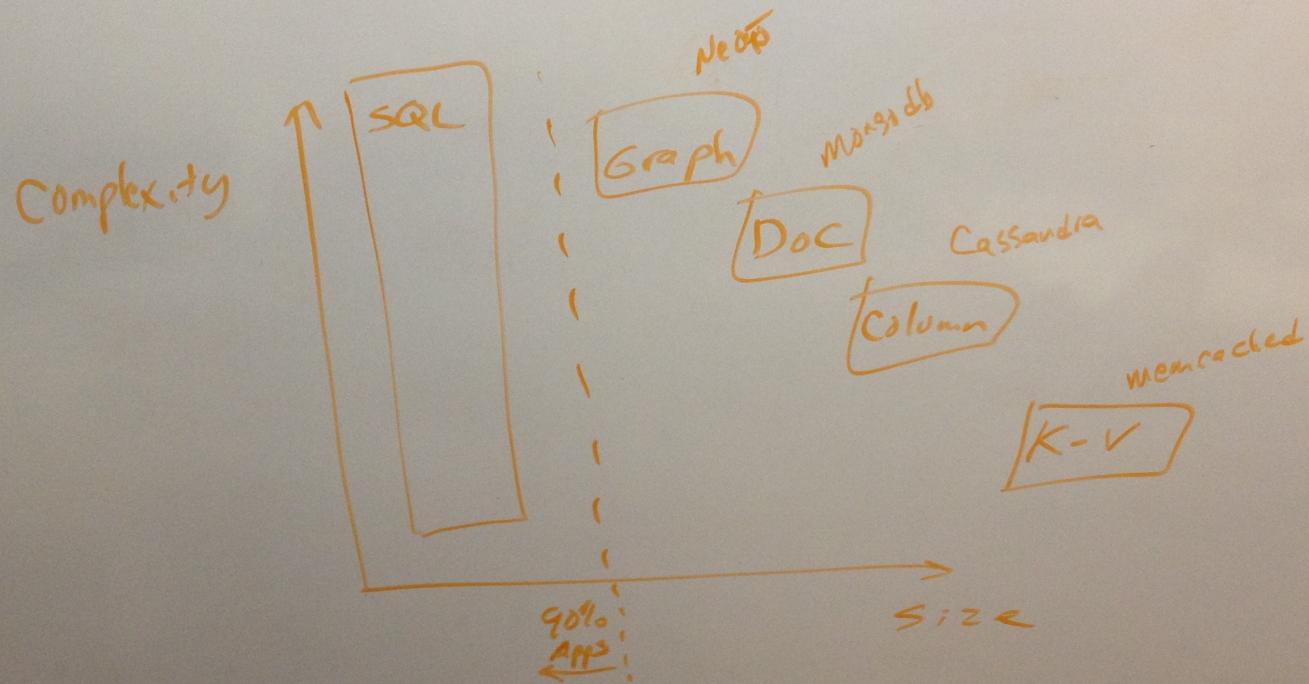
# w3schools.com stats

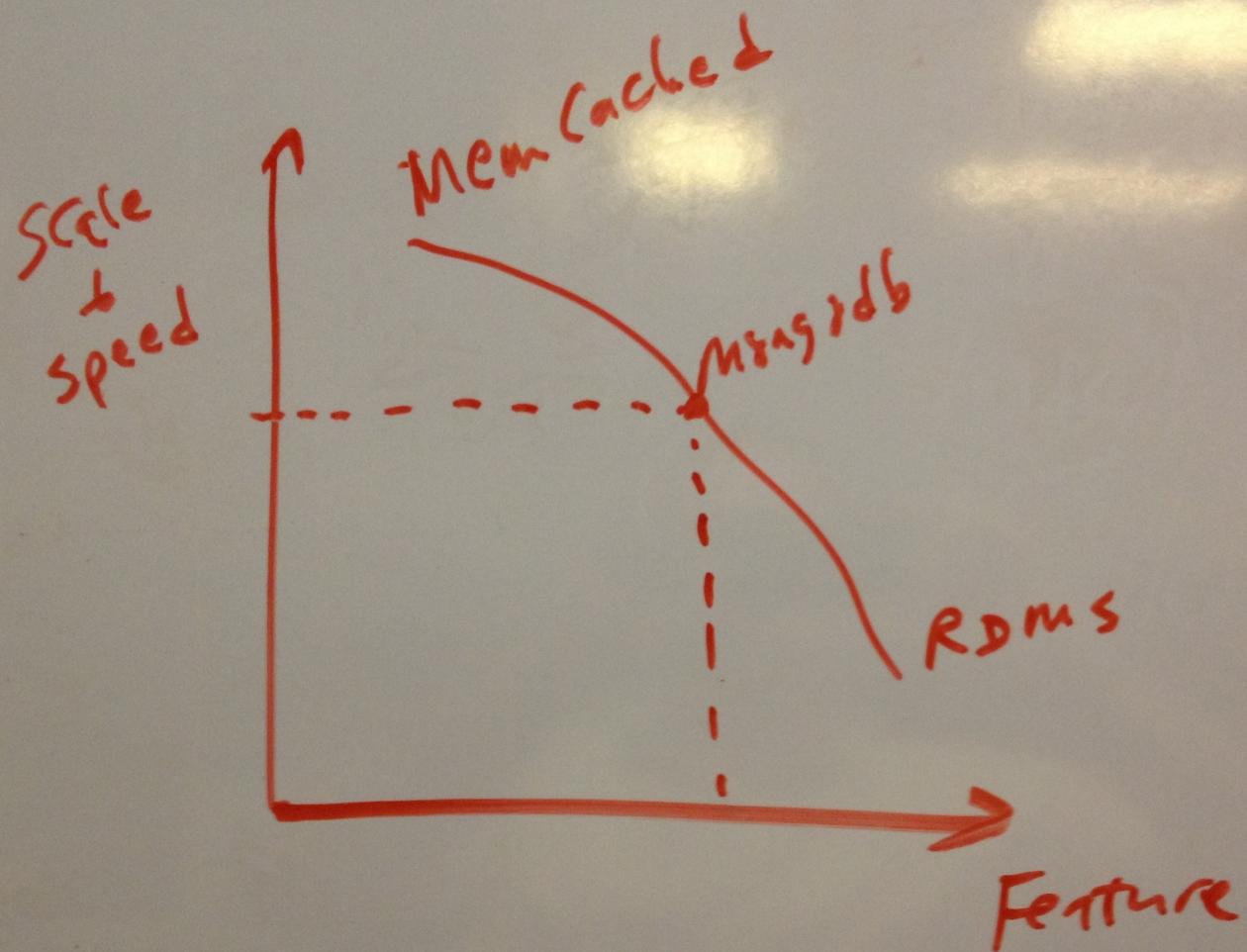
Browser	Chrome	60%
	FireFox	24%
	I.E.	10%
OS	Windows	81%
	Mac	10%
	Linux	6%
	Mobile	4%
Mobile	Android	2.3%
	iOS	1.1%
	Windows	0.4%

# What about Databases?

- Relational vs NoSQL
- MDMS uses MongoDB

## Complexity vs. Size





# Target Audience

- Software Developers
- Already familiar with Mongodb basics
- Already familiar with Linux

# MongoDB basics

- Store Document as BSON
- Docs can reference or embed other Docs
- Javascript Shell
- Replication for High Availability
- Sharding for Scalability

# Using Brew on Mac

```
clarence@acs:~/mongo$ brew list
ack          git          libpng        memcached      rabbitmq
autoconf     imagemagick  libtiff       mongodb       readline
automake     jasper       libtool       openssl        rebar
cscope       jmeter       libyaml       ossp-uuid     scons
ctags        jpeg         little-cms   pkg-config    sqlite
erlang       libevent     macvim       postgresql   wget
freetype     libgpg-error makedepend  protools     wxmac
gdbm        libksba      markdown     python       zeromq22
clarence@acs:~/mongo$ brew info mongodb | grep Cellar
/usr/local/Cellar/mongodb/1.8.1-x86_64 (16 files, 93M)
/usr/local/Cellar/mongodb/2.2.0-x86_64 (20 files, 170M)
/usr/local/Cellar/mongodb/2.4.6 (18 files, 320M) *
/usr/local/Cellar/mongodb/2.6.3 (17 files, 407M)
clarence@acs:~/mongo$ brew switch mongodb 2.4.6
Cleaning /usr/local/Cellar/mongodb/1.8.1-x86_64
Cleaning /usr/local/Cellar/mongodb/2.2.0-x86_64
Cleaning /usr/local/Cellar/mongodb/2.4.6
Cleaning /usr/local/Cellar/mongodb/2.6.3
14 links created for /usr/local/Cellar/mongodb/2.4.6
clarence@acs:~/mongo$ █
```

# MongoDB Files

```
clarence@mac:~/mongo/debug/data/shard01/db$ ls -ltr
total 2293760
-rw----- 1 clarence staff 134217728 Nov 14 17:57 MONACODB.1
-rw----- 1 clarence staff 134217728 Nov 14 17:57 SESSIONDB.1
-rw----- 1 clarence staff 134217728 Nov 14 17:57 SUBSCRIBERDB.1
-rw----- 1 clarence staff 134217728 Nov 14 17:57 PURCHASEDB.1
-rw----- 1 clarence staff 134217728 Nov 14 17:58 RENTALDB.1
-rw----- 1 clarence staff 16777216 Nov 18 10:12 SESSIONDB.ns
-rw----- 1 clarence staff 67108864 Nov 18 10:12 SESSIONDB.0
-rw----- 1 clarence staff 16777216 Nov 18 10:12 PURCHASEDB.ns
-rw----- 1 clarence staff 67108864 Nov 18 10:12 PURCHASEDB.0
-rw----- 1 clarence staff 16777216 Nov 18 10:13 RENTALDB.ns
-rw----- 1 clarence staff 67108864 Nov 18 10:13 RENTALDB.0
-rw----- 1 clarence staff 16777216 Nov 19 16:14 SUBSCRIBERDB.ns
-rw----- 1 clarence staff 67108864 Nov 19 16:15 SUBSCRIBERDB.0
-rw----- 1 clarence staff 16777216 Nov 24 17:22 MONACODB.ns
-rw----- 1 clarence staff 67108864 Nov 25 11:52 MONACODB.0
-rwxr-xr-x 1 clarence staff 0 Dec 12 10:53 mongod.lock
-rw----- 1 clarence staff 16777216 Dec 12 10:53 local.ns
-rw----- 1 clarence staff 67108864 Dec 12 10:53 local.0
drwxr-xr-x 2 clarence staff 68 Dec 12 10:53 journal
clarence@mac:~/mongo/debug/data/shard01/db$
```

# Indexing a Collection

- Faster Read, Slower Writes
- Document move more expensive
- Build index before you have data is faster
- Covered index (explain, hint)
- Unique, sparse, TTL, text, geospatial

```
db.assets.find({ _type: "TermLimitation"}).hint({_id: 1}).explain()
2 db.assets.find({ _type: "TermLimitation"})
3 .hint({_type: 1}).explain()
4 {
5 >     "cursor" : "BtreeCursor _type_1",
6 >     "isMultiKey" : false,
7 >     "n" : 38,
8 >     "nscannedObjects" : 38,
9 >     "nscanned" : 38,
10 >    "nscannedObjectsAllPlans" : 38,
11 >    "nscannedAllPlans" : 38,
12 >    "scanAndOrder" : false,
13 >    "indexOnly" : false,
14 >    "nYields" : 0,
15 >    "nChunkSkips" : 0,
16 >    "millis" : 0,
17 >    "indexBounds" : {
18 >        "_id" : [
19 >            [
20 >                {
21 >                    "$minElement" : 1
22 >                },
23 >                {
24 >                    "$maxElement" : 1
25 >                }
26 >            ]
27 >        ]
28 >    }
29 >    "server" : "macs:27018",
30 >    "millis" : 0
31 > }
```

# MONGODB Memory Usage

- MongoDB runs fast when working set is in memory.
- Working set: data accessed most frequently (indexes & subset of data)
- Resident memory (includes file system cache & mongod process memory)
- Restarts and dropping caches

# Determining Working Set

- Knowledge of data & indexes
- Average document size and index size
- Pages in memory & page Size
- Hard & Soft page faults under load

# MONGODB General Info

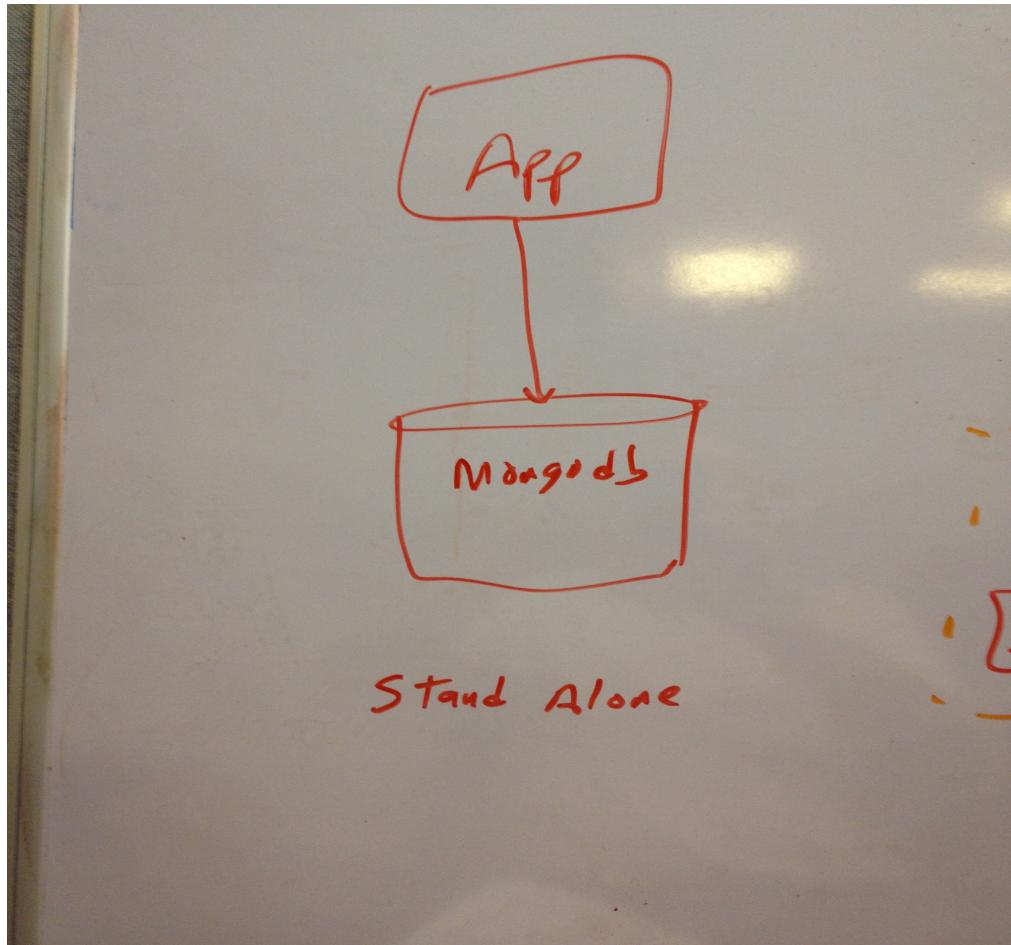
- Syncs Journal every 100msec
- Syncs Data every 60sec
- Preheating Data can be beneficial
- Follow file system/os config recommendations
- Does not track available disk space
- Mongodb is new and evolving

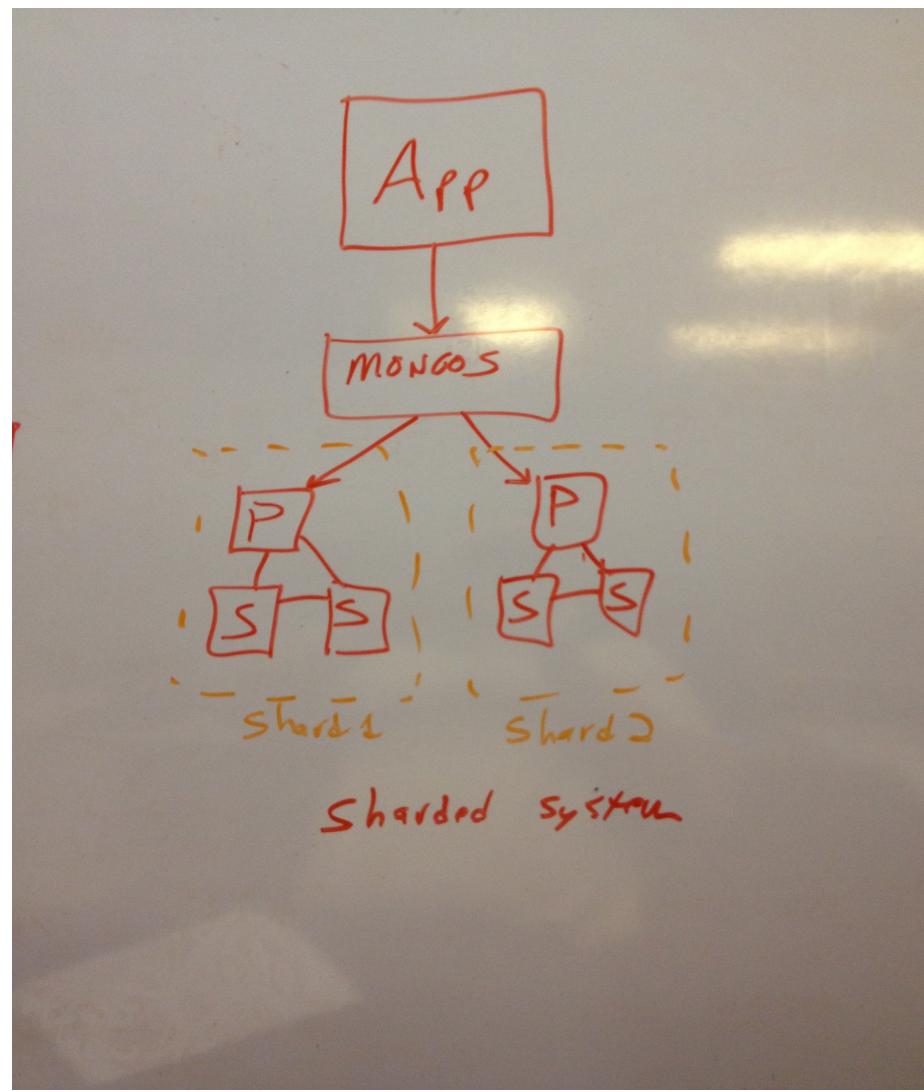
# MONGODB Sizing Considerations

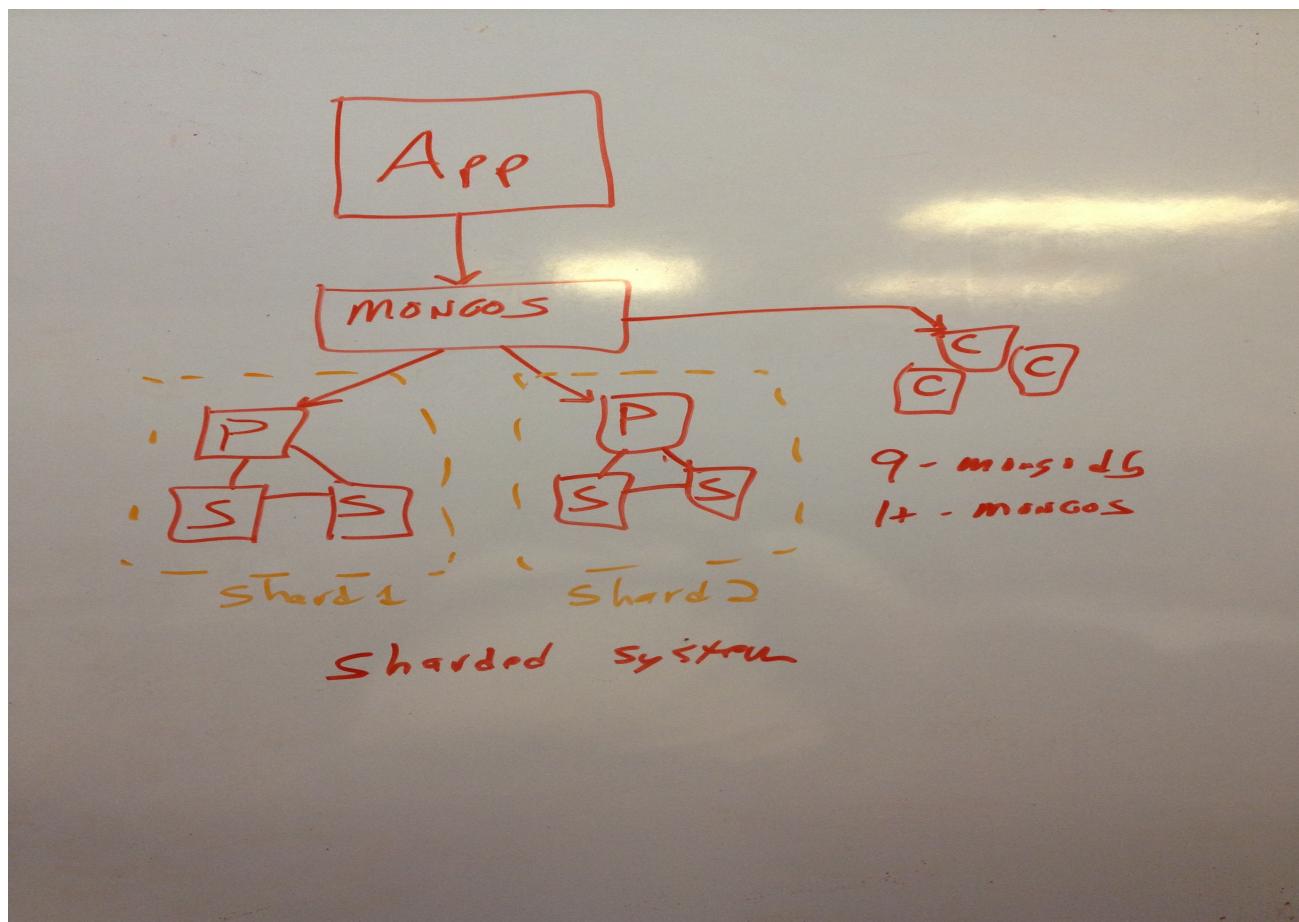
- Database size, number processors, I/O Benchmarking
- TCP connections, file descriptor limits
- Port limitations, connection churn, kernel connection parameters

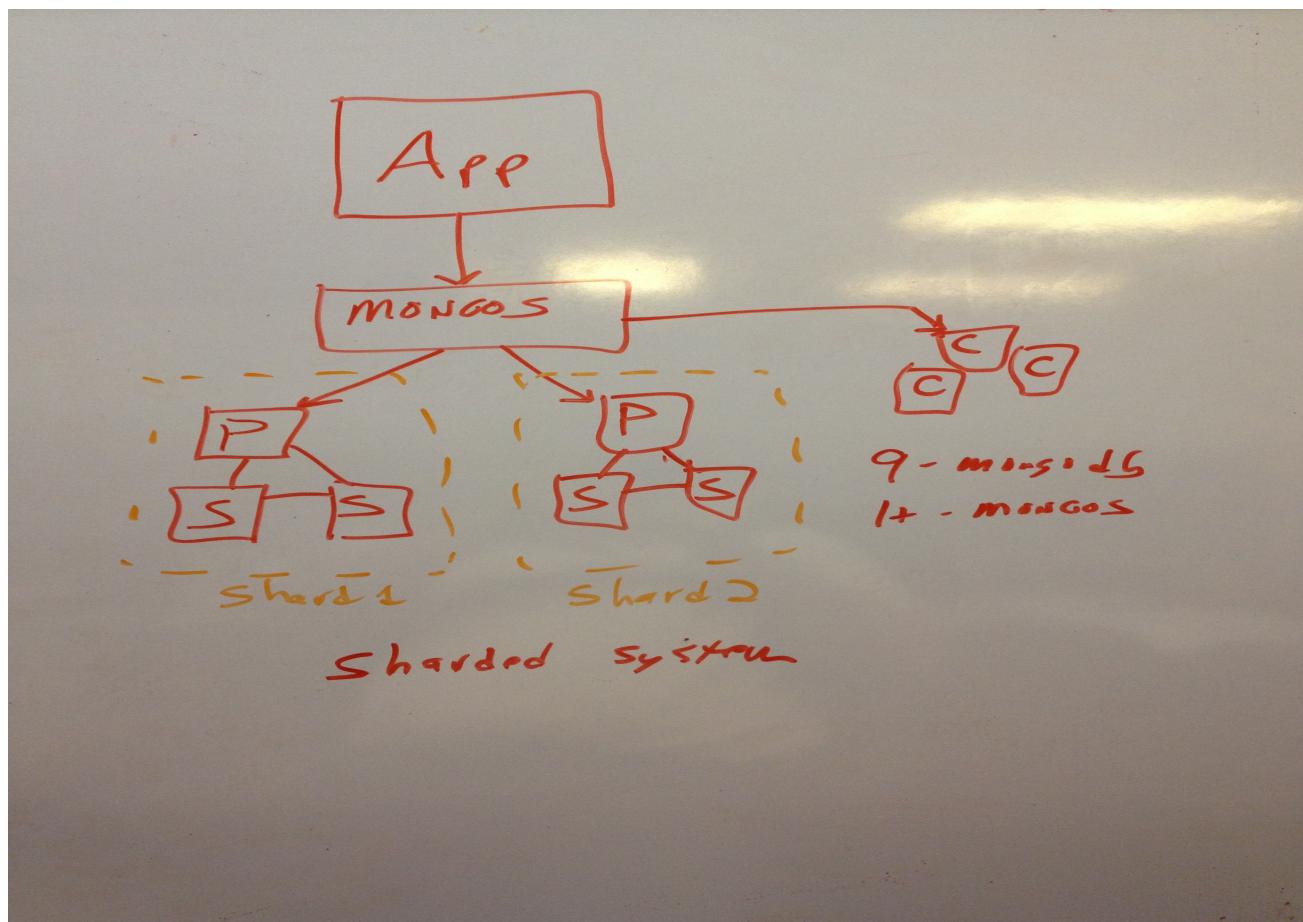
# Development Environment

- Typical Standalone Mongod
- Could be 3 system replica set or a sharded system with mongos
- Likely running on a single dev box
- On Mac Brew and mlaunch are nice









# Mlaunch for dev environment

- mlaunch --init --single --sharded 3 --config 1
- mlaunch start
- mlaunch list
- mlaunch stop

```
clavrence@macs:~/mongo/debug$ mlaunch list
PROCESS          PORT      STATUS      PID
debug            10        running     -
mongos           27017    running     -
homework         27021    running     -
config server   27021    running     -
shard01          27018    running     -
shard02          27019    running     -
shard03          27020    running     -

```

# Production Deployment

- 3 member replica Set with voting
- Primary, Secondary, Secondary
- Primary, Secondary, Arbitor
- Maybe extras for analytics or backups
- Often distribute replicas across different sites
- Sharding used when necessary

# Distributed Replica Sets

- Writes to a remote primary may have longer network delays
- Replica Set chaining can be used to help reduce network delays for updating secondaries

# Sharding

- Each Shard normally a replica set
- Sharding requires 3 config servers (1 for dev)
- Sharded writes fail when a config server is down or unavailable
- Unsharded collection will live on one shard
- Sharded collections are split across shards using a shard key
- Shard keys need to be chosen carefully

# Shard Balancing

- A load balancer moves chunks of data across all shards
- Chunks are a general abstraction. They are not data size nor number of documents nor load balancing.
- Balancing chunks is not the same as balancing load or data

# Sharding Considerations

- Mongos has responsibility to combine queries that span different shards
- Jumbo chunks and empty chunks create problems
- Orphan chunks sometimes show up on secondary reads or direct mongod queries (see `mongodb/support-tools/orphanage.js`)
- `AllChunkInfo`  
[`\(https://gist.github.com/normgraham/8716338\)`](https://gist.github.com/normgraham/8716338)



normgraham / allChunkInfo.js

Created on Jan 30



Display detail and aggregate chunk info for a sharded collection. To run, load the following code for allChunkInfo() into a mongos process in your cluster and run allChunkInfo( "dbname.collname" ).

[allChunkInfo.js](#)[Raw](#)

```
1 var allChunkInfo = function(ns){
2     var chunks = db.getSiblingDB("config").chunks.find({"ns" : ns}).sort({min:1}); //this will return all c
3     //some counters for overall stats at the end
4     var totalChunks = 0;
5     var totalSize = 0;
6     var totalEmpty = 0;
7     print("ChunkID,Shard,ChunkSize,ObjectsInChunk"); // header row
8     // iterate over all the chunks, print out info for each
9     chunks.forEach(
10         function printChunkInfo(chunk) {
11
12             var db1 = db.getSiblingDB(chunk.ns.split(".")[0]); // get the database we will be running the comma
13             var key = db.getSiblingDB("config").collections.findOne({_id:chunk.ns}).key; // will need this for
14             // dataSize returns the info we need on the data, but using the estimate option to use counts is le
15             var dataSizeResult = db1.runCommand({datasize:chunk.ns, keyPattern:key, min:chunk.min, max:chunk.ma
16             // printjson(dataSizeResult); // uncomment to see how long it takes to run and status
17             print(chunk._id+","+chunk.shard+","+dataSizeResult.size+","+dataSizeResult.numObjects);
18             totalSize += dataSizeResult.size;
19             totalChunks++;
20             if (dataSizeResult.size == 0) { totalEmpty++ }; //count empty chunks for summary
21         }
22     )
23     print("*****Summary Chunk Information*****");
24     print("Total Chunks: "+totalChunks);
25     print("Average Chunk Size (bytes): "+(totalSize/totalChunks));
26     print("Empty Chunks: "+totalEmpty);
27     print("Average Chunk Size (non-empty): "+(totalSize/(totalChunks-totalEmpty)));
28 }
```



&lt;scr

HTTPS

http

You can  
SSH.

```
Fri Dec 12 11:40:06.586 ReferenceError: subscriber is not defined
mongos> sh.status()
--- Sharding Status ---
sharding version: {
  "_id" : 1,
  "version" : 3,           Sharding
  "minCompatibleVersion" : 3,
  "currentVersion" : 4,    requires 3 config servers (1 for dev)
  "clusterId" : ObjectId("546687e2f7607a7475c98504")
}
shards: [
  {
    "_id" : "shard0000", "host" : "Craig-Lawrences-MacBook-Pro.local:27018"
  },
  {
    "_id" : "shard0001", "host" : "Craig-Lawrences-MacBook-Pro.local:27019"
  },
  {
    "_id" : "shard0002", "host" : "Craig-Lawrences-MacBook-Pro.local:27020"
  }
]
databases: [
  {
    "_id" : "admin", "partitioned": false, "primary" : "config"
  },
  {
    "_id" : "MONACODB", "partitioned" : false, "primary" : "shard0000"
  },
  {
    "_id" : "BOOKMARKDB", "partitioned" : false, "primary" : "shard0000"
  },
  {
    "_id" : "CDNDB", "partitioned" : false, "primary" : "shard0000"
  },
  {
    "_id" : "NADB", "partitioned" : false, "primary" : "shard0000"
  },
  {
    "_id" : "PURCHASEDB", "partitioned" : false, "primary" : "shard0000"
  },
  {
    "_id" : "QAMDB", "partitioned" : false, "primary" : "shard0000"
  },
  {
    "_id" : "RENTALDB", "partitioned" : false, "primary" : "shard0000"
  },
  {
    "_id" : "SESSIONDB", "partitioned" : true, "primary" : "shard0000"
  },
  {
    "_id" : "SOPDB", "partitioned" : false, "primary" : "shard0000"
  },
  {
    "_id" : "SUBSCRIBERDB", "partitioned" : true, "primary" : "shard0000"
  }
]
SUBSCRIBERDB.subscribers
  shard key: { "sms_id" : 1, "subscriber_id" : 1 }
  chunks:
    shard0001 271
    shard0000 1
28 |     { "sms_id" : { "$minKey" : 1 }, "subscriber_id" : { "$minKey" : 1 } } --> { "sms_id" : "1", "subscriber_id" : { "$minKey" : 1 } }
{ "_id" : "test", "partitioned" : false, "primary" : "shard0001" }
{ "_id" : "SESSIONDB", "partitioned" : false, "primary" : "shard0001" }
{ "_id" : "MONACODB_TEST_7", "partitioned" : false, "primary" : "shard0002" }
{ "_id" : "MONACODB_TEST_1", "partitioned" : false, "primary" : "shard0002" }
{ "_id" : "MONACODB_TEST_5", "partitioned" : false, "primary" : "shard0002" }
{ "_id" : "MONACODB_TEST_2", "partitioned" : false, "primary" : "shard0002" }
{ "_id" : "MONACODB_TEST_4", "partitioned" : false, "primary" : "shard0002" }
{ "_id" : "MONACODB_TEST_3", "partitioned" : false, "primary" : "shard0002" }
{ "_id" : "MONACODB_TEST_0", "partitioned" : false, "primary" : "shard0002" }
{ "_id" : "MONACODB_TEST_6", "partitioned" : false, "primary" : "shard0002" }
{ "_id" : "SUBSCRIBER", "partitioned" : false, "primary" : "shard0001" }
```

# Example AllChunkInfo

```
allChunkInfo("SUBSCRIBERDB.subscribers")
```

```
ChunkID,Shard,ChunkSize,ObjectsInChunk
```

```
SUBSCRIBERDB.subscribers-sms_id_MinKeysubscriber_id_MinKey,shard0001,432,1
```

```
SUBSCRIBERDB.subscribers-sms_id_ "1"subscriber_id_ "subscriber_0",shard0000,2512,8
```

```
*****Summary Chunk Information*****
```

```
Total Chunks: 2
```

```
Average Chunk Size (bytes): 1472
```

```
Empty Chunks: 0
```

```
Average Chunk Size (non-empty): 1472
```

# Other Sharding Options

- Tag Based Sharding (pinning)
- Hash Based Sharding (scales writes, may slow queries)

# MONGODB Driver Options

- Options vary with different drivers and with different versions of a driver
- Connection timeout, connections per host
- Connections use 1MB per connection and file descriptors too)
- TODO: Show mongoid.yml

# Connection Strategy

- Primaries have greatest number of connections
- A 3 shard system may have 20 connection per mongos connection
- Consider maximum possible primary connections (1MB/per, ulimits, 20k typical mongod max)
- Consider limiting mongos and drivers to protect system from reaching overload

# MONGODB Failover

- Journaling
- Oplog delays
- Write Concerns
- Read Concerns
- Triggering Rollback

# Mongod Tools

- mlaunch (for dev use only)
- mloginfo (log file summary info, see [blog.rueckstiess.com](http://blog.rueckstiess.com))
- mlogfilter
- mplotqueries
- ([github.com/rueckstiess/mtools](https://github.com/rueckstiess/mtools))
- pip install mtools (pymongo, matplotlib)

# Example Data for mloginfo

- `mloginfo mtools_example.log --connections  
--distinct --queries --restarts --rsstate`
- `mlogfilter mtools_example.log --from "2014  
Jun 21 00:12:59" --to "+1s" --slow 1000  
--operation remove --namespace  
"grilled.indigo.papaya"`
- `mplotqueries mtools_example.log --type  
scatter --group namespace --group-limit 10`