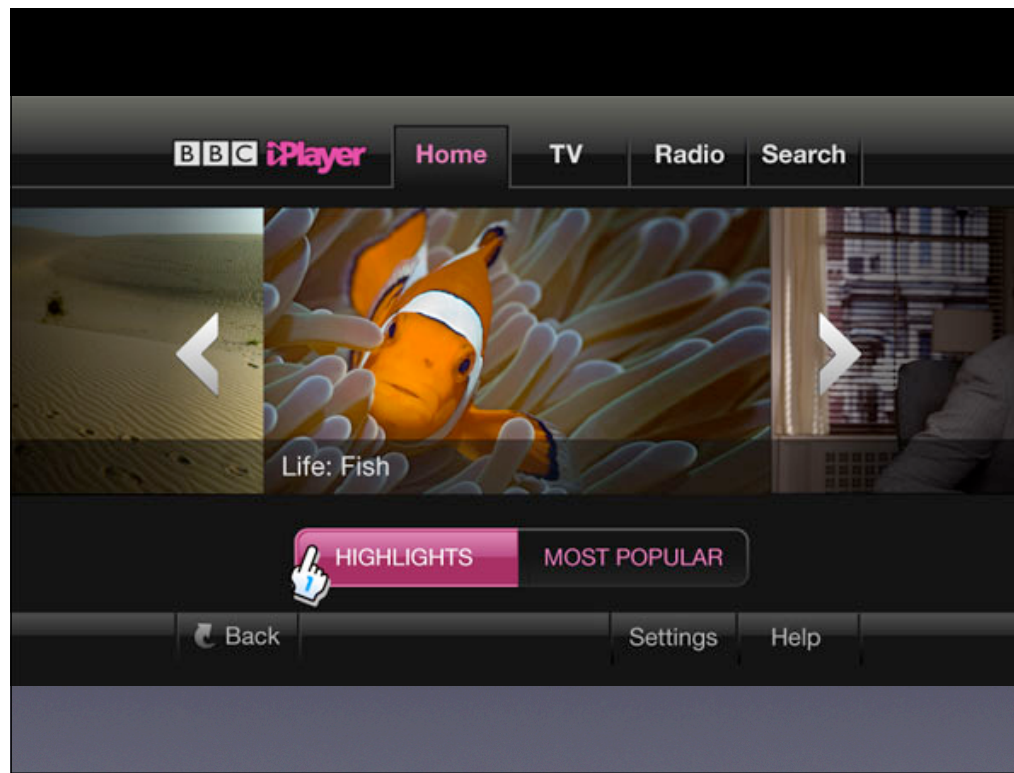# Enda Farrell



Software architect, product owner and lead developer for the BBC's usage of CouchDB

Auntie on the Couch

what CouchDB is, how to use it, and what it is like at a large scale

A little context before I start. I expect most of you have come across the BBC and it's web site. It's big, there are popular parts and there are obscure parts. Which are backed in some way by CouchDB?

So - which "little" sites are using CouchDB? Might I have ever come across them? Do they matter? Would anyone notice if they disappeared? ;-) Someone might ;-)

## What is CouchDB?

- ... is a document-oriented database that can be queried and indexed in a MapReduce fashion using JavaScript. ... also offers incremental replication with bi-directional conflict detection and resolution.

- ... provides a RESTful JSON API than can be accessed from any environment that allows HTTP requests.

So (almost) goes the introduction from http://couchdb.apache.org/
Let's skip the text and have a look at CouchDB in action

how to use it

```
CouchDB demo
$ curl -X DELETE http://couchdb:5984/users/enda?rev=6-4351162c2f4dc640708c0275
587811e8
{"ok":true,"id":"enda","rev":"7-c0ae53badc76af56ffaf9bdb9e06f68a"}
$
```

CouchDB uses standard HTTP RESTful commands - GET, PUT, POST and DELETE to access data. It uses a JSON format. Updating an existing document _requires_ having the current revision of that document which stops accidental over-writing of data by clients.

Compacting databases removes from disk the old, over-written versions of documents. In our setup, we (a) don't often care about old versions and (b) we like saving space. This space saving can be significant depending on how many updates are done to documents.
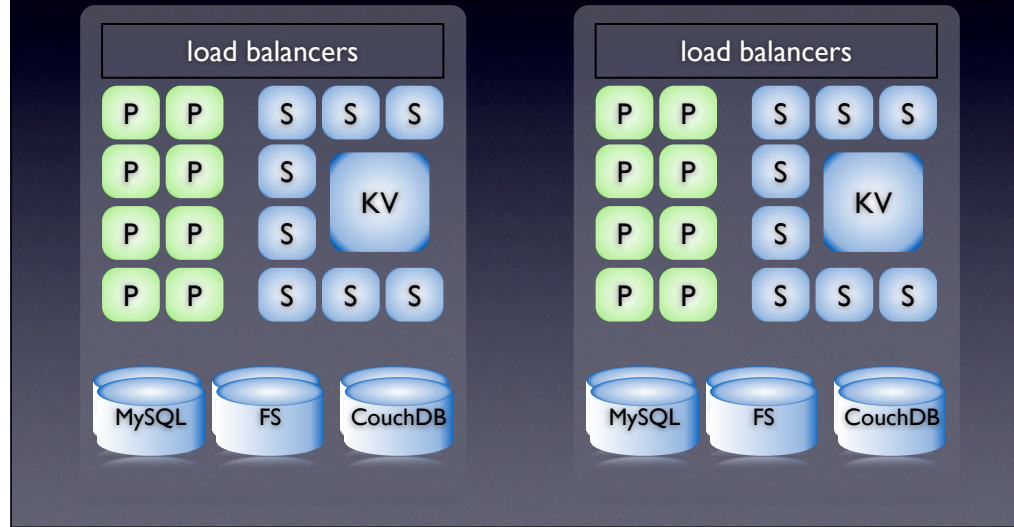
This is the old "trigger" replication which has been improved on in 0.10. Notice that even through CouchDB has an admin UI - _all_ commands to the service - like this "go replicate these" - are RESTful HTTP calls.

# what it is like at scale

- Context - one service on a new platform

- Operations

- Replication and compaction

- Some statistics

- How we use it, how we don't

(mutually authenticating) secure services (with a small "s") oriented architecture. It's not the "XML, SOAP, WSDL, UDDI" version of SOA - it is lighter, easier to code to, quicker, easier to scale and easier to manage. "P" are PHP applications assembling data. "S" are JSON/XML service providers.

To make CouchDB "fit" into our platform, we put a wrapper API above it, and to make operations simple, we put a "replication daemon" underneith.

# what it is like at scale

- Context - one service on a new platform ...

- Operations

- Replication and Compaction

- Some statistics

- How we use it, how we don't

# Operations

- Installation and running
- Instances and system utilisation
- Scalability

# Operations

- Ops folk are busy and have thankless tasks

  ```
  yum install couchdb-config

  service couchdb start|stop|restart

  service couchdb-replicatr start|stop
  ```

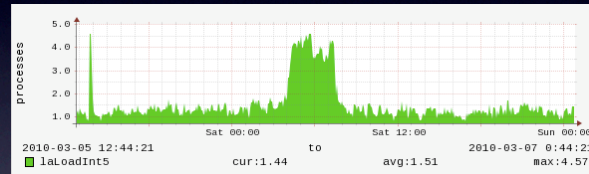We did a little work in packing RPMs and made CouchDB look act and "smell" like any other service on the platfrom

Apart from specifying IP bindings, database directories etc, the only "customisation" we have is to spin up (and down) 4 nodes per physical machine

# Operations



```
                                farree02@kv001:~
top - 14:39:52 up 293 days,  2:58,  1 user,  load average: 1.68, 1.53, 1.48
Tasks: 187 total,   1 running, 185 sleeping,   0 stopped,   1 zombie
Cpu(s):  9.9%us,  1.0%sy,  0.0%ni, 88.8%id,  0.0%wa,  0.0%hi,  0.2%si,  0.0%s
Mem:  18482752k total, 18288276k used,   194476k free,   291400k buffers
Swap: 16771820k total,      924k used, 16770896k free, 16918700k cached

  PID USER      PR  NI  VIRT  RES  SHR S %CPU %MEM    TIME+  COMMAND
16012 couchdb   25   0  531m 105m 4072 S 27.3  0.6  5244:15 beam.smp
14349 root      25   0 1896m 204m 9000 S 23.6  1.1  1222:32 replicatr
15834 couchdb   25   0  627m 185m 4072 S 15.0  1.0  5909:13 beam.smp
15954 couchdb   25   0  558m 111m 4072 S 11.0  0.6  6464:57 beam.smp
15891 couchdb   25   0  506m  78m 4076 S  6.3  0.4  3998:13 beam.smp
11736 root      15   0 83520  10m 1880 S  2.7  0.1 184:31.85 snmp-collector
30501 root      15   0  106m 6392 2704 S  1.7  0.0  21:39.29 python
```

8 cores, 16 GB RAM. CouchDB is mostly kind on CPU, and if you do not run views, has a v consistent memory footprint.
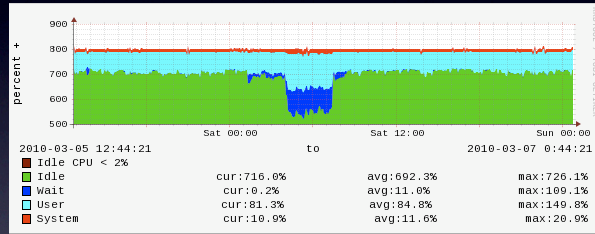
Operations

Low load average

Look - doing backups - which by the way are as simple as "copy the files in these directories" - has a big load effect. Sat/Sun are not "quiet" on the platform - this is essentially the same 7 days a week
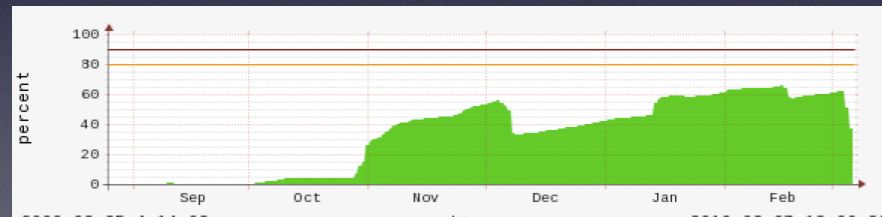
# Operations



Kind to CPU

The green is idle time on these graphs.

# Operations

- Robust - very robust

- restarts < 1 sec

- no "fix-up" if it crashes - append only B-tree

- No "scheduled downtime" needed to restart

Op ☞ scalability

- Still in our early stages
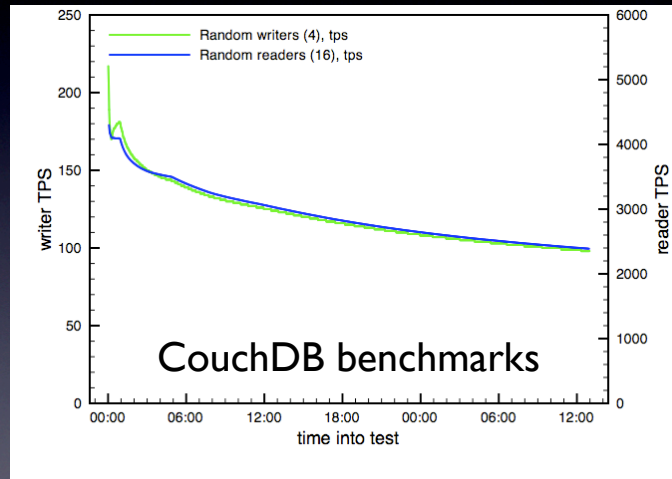- We can double and double again our infra with only small rc.d script & DNS changes

This somewhat shows how we are still "beginning" our scalability journey.

# Scalability

- What do "you" need in the next 12 months?

- If you don't know, what attributes do you rely on to deal with this?

  - consistency - linear or O(log n) graphs

  - reliable empirical stats
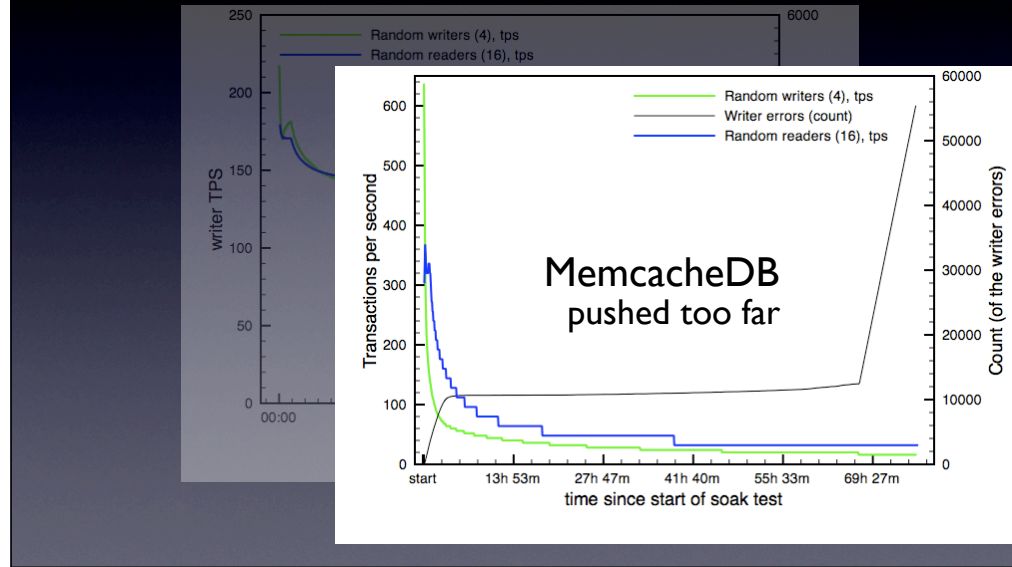
  - known break points - stress tests

Order log n decay of performance with data sizes - watch the blip as we break through the machine's working set.
We ran out of disk before we hit a break-point of these tests.
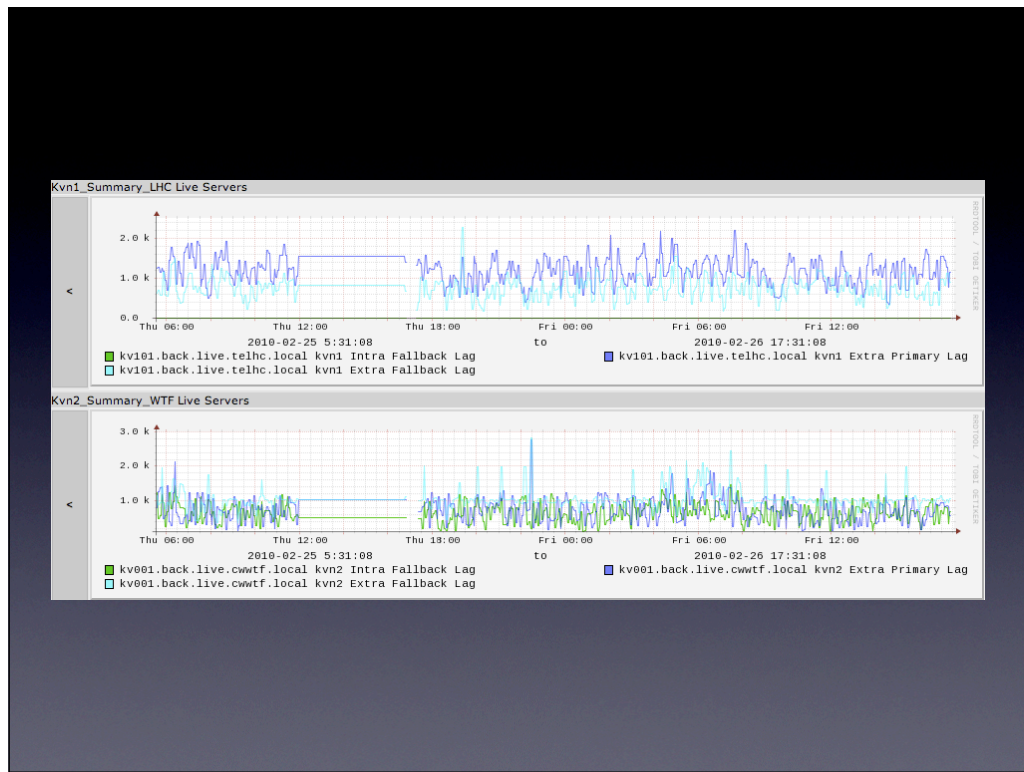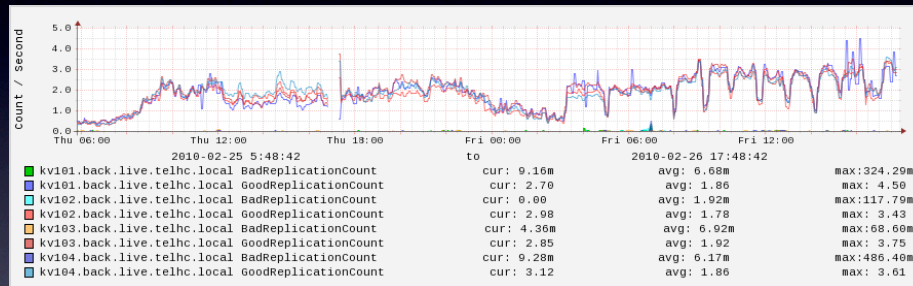Writers finished at 100 tps, readers at 2400 in this test

When you push a system too far - like an in-memory DB beyond the working set - you see this sort of graph. Exponential decay, order of magnitude drops beyond the working set, a findable break point beyond which you cannot scale. Writers finished at 40 tps, readers at 60 in this test - though started much better.

# Scalability - reliable stats

- Throughout the platform we use SNMP to collect, organise, store and present the data

- We can scale by looking at where we need to - proactively

CouchDB access speeds, num accesses, replication lag, counts of http actions, KV access speeds, KV namespace stats, replication stats

Summary charts for replication statistics

CouchDB users will be familiar with the white background - "Futon" a relaxing admin UI (which does NOT have any "special" hooks - it just uses the same API calls). The panel on the left is an addition of ours - showing the shards across different DCs for different environments (live, stage, test, int). Every few seconds, some funky AJAX goes and checks each - giving it a set of colours if not.

## Scalability - stress tests

- Everything breaks

- The question is - "where?"

- No - the question is "why?"

- No - the question is "when?"

- Aaagh!

CPU on firewalls, network interrupts on NICs, high churn data evicts memcache and > 10% f/e calls go back to service, bandwidth of traffic managers - all platforms break. Code sometimes breaks too ;-)

## Scalability - stress tests

- Known break points:
- RAID controller throughput to disk
- Inter-DC VPN drops packets, bad HTTP
- Poor JavaScript breaking views
- Early adopter CouchDB bugs - all now fixed
- Network devices caching on URLs

1 - Our RAID controllers are a bottle neck - if we try to push MORE than they can handle, the OS on the box starts to back up and that causes problems. Not a CouchDB issue.
2 - Can cause sessions to hang as ACKs are not reliably delivered. If the session is a replication, it makes it look like its hung. Can't really blame CouchDB for that!
3 - Traffic manager CPU - (platform wide, but as one of the most shared network resources, seen on the KV service) - hit that and requests back up
4 - Poor Javascript in views - can completely kill the use of that database on that node - slow response times leading to repeated requests when timeouts occur, leading to a snowball of higher and higher load
5 - Compaction, replication 404  ===  6 - Too clever for its own good - poor corporate networks

# what it is like at scale

- Context - one service on a new platform ...
- Operations
- Replication and Compaction
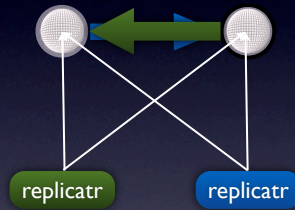- Some statistics
- How we use it, how we don't

replication

source data on a CouchDB node

This is "trigger" replication, to be replaced with 0.10's "continuous" replication
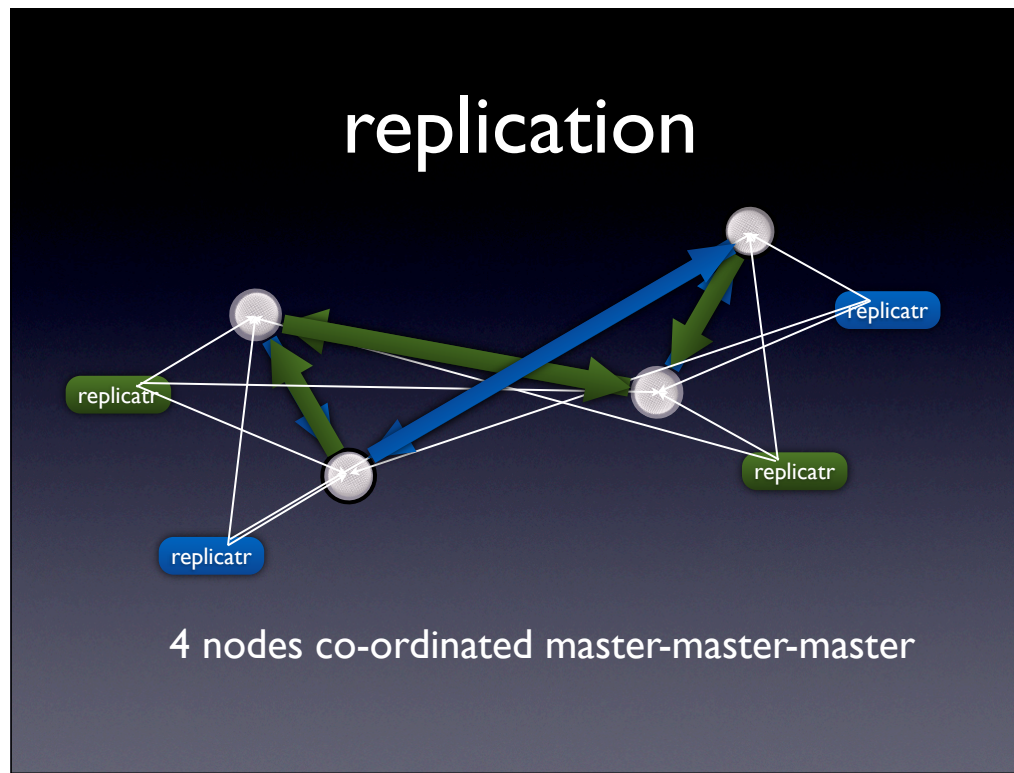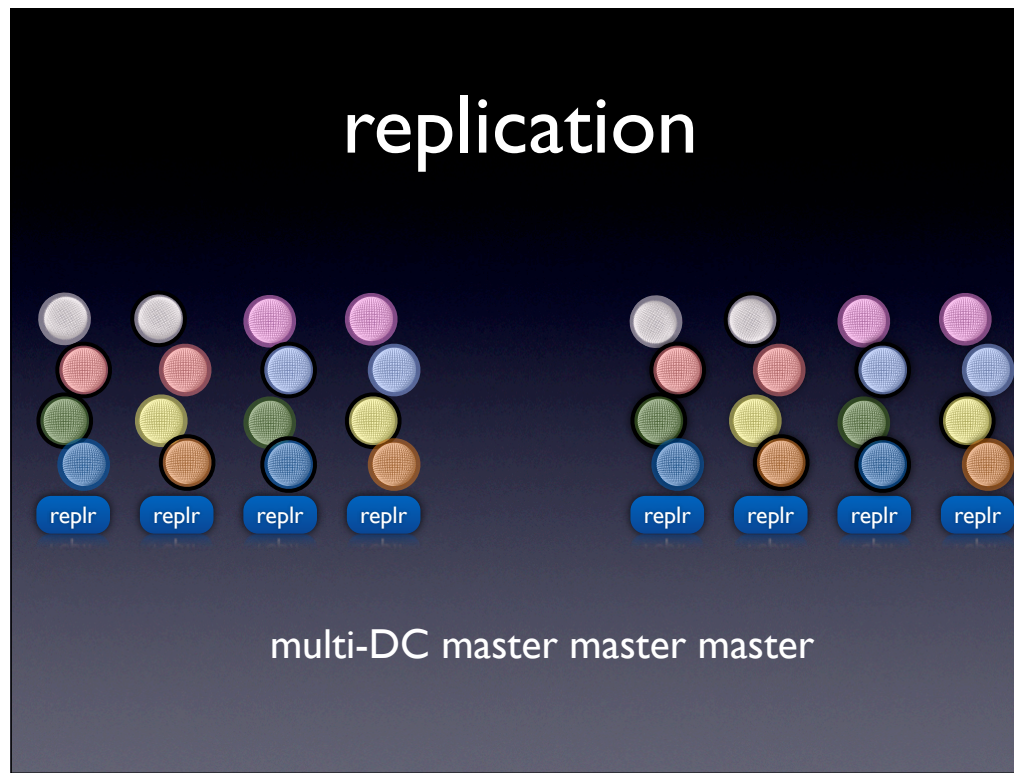
# replication



has source changed?

`POST /db/_replicate`

CouchDB replicates

replicated pair

replication

4 nodes co-ordinated master-master-master

OK - this "looks" scary - but it's quite normal on our platform, and across the web. It looks good though - helps the business understand some of the hidden complexities

It's a step up from master-slave to master master.  Another one to go to 4 node co-ordinated master-master-master. Another one when you see all such shards together. There's another step up when you remember that replication is per database - we will have 100s.

## replication

- No other data store on the platform gives master master updates

- Deploy to one, the other, both DCs

- Application code simpler - no "I can read but not write" logic that our MySQL users have

- Eventual consistency is really quite OK on our operational platform due to DC affinity

What business advantages come from this? Cool graphs - perhaps! Most importantly, other code using the KV store can be simpler, easier to understand, easier to deploy, and perhaps significantly does NOT need to know whether they are running in a DC which allows writes.
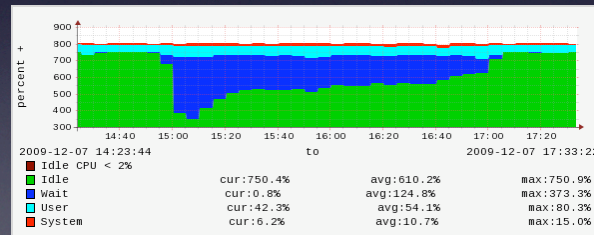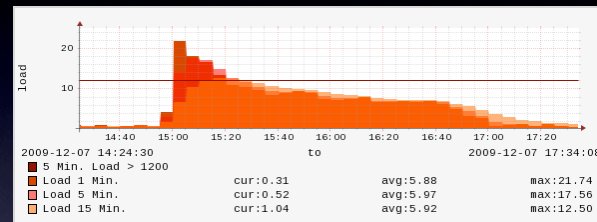
# Compaction

- Compaction removes old revisions of documents, saving space

- Be advised - compact namespaces serially!

- It's included as part of our replicatr daemon

If you do not update existing docs, there is little benefit in compacting, with the cost of slower access during the compaction process. Some logic is therefore advisable in deciding on when.
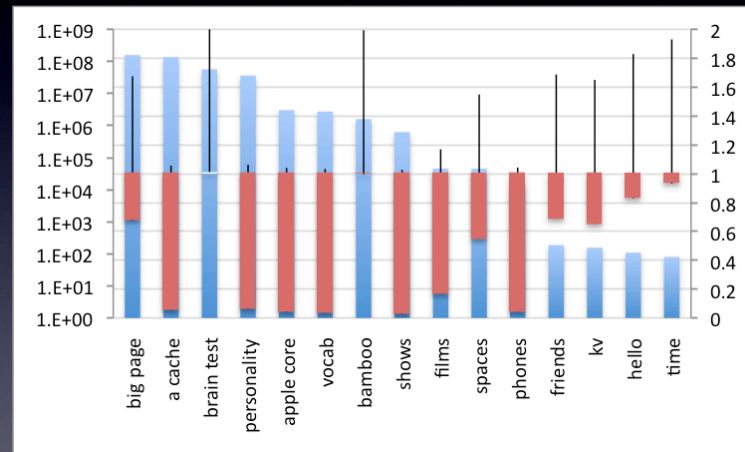On our platform, previous revisions are not important, so we save space.

In Dec we were starting to hit 60% disk, I was going to be away for almost a month and we hadn't compacted for a while. Instead of doing it serially I compacted almost everything together at once. It was more "ouch" than we had expected!

Log 10 scale on the left - the size of the databases. Red is compaction RATIO saved. "a cache" compacted fantastically, bamboo not at all and ran at a high possible cost.

# what it is like at scale

- Context - one service on a new platform ...

- Operations

- Replication and Compaction

- Some statistics

- How we use it, how we don't

hourly, we gather up lots of stats and as we eat our own food, being fans of our own service, we keep them for posterity in CouchDB ;-)
So - they were some stats.

# 1,030
GB

# 155 million

requests on an average day

5 billion

5,036,466,928 requests, since last summer

Chris Anderson "having the largest known CouchDB installation" - one of the 3 biggest installations

96% are GET, 3% are POST (for replication, compaction and new database creation), 1% are the PUTs which create and update document, we discourage DELETEs due to the highly parallel nature of our platform. One of the 3 biggest known CouchDB installations in the world.

# what it is like at scale

- Context - one service on a new platform ...

- Operations

- Replication and Compaction

- Some statistics

- How we use it, how we don't

# How we don't use it

- Views
- Attachments

**no views**

- they are cool, but on the platform we want a simple "Key Value" store

- poor javascript concerns mean we'll move slowly here

Simple - we want to use things in a way that is simpler than CouchDB CAN be used.

My engineering team does not "use" the service much, other developers at the BBC do. Given that CouchDB is schema-less, the structure of documents can change. I can't trust that every developer will take each of their own edge-cases into account here.
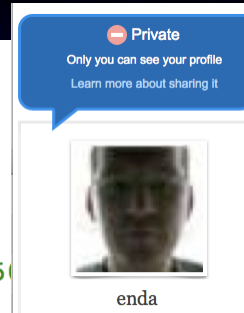
# no attachments
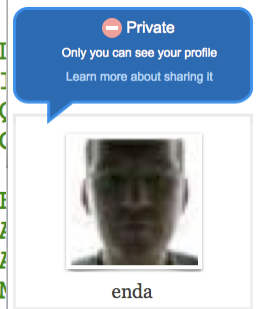
# no attachments

```
{
    _id: "u_enda",
    _rev: "1-3019978615",
    avatarChoice: "upload",
    dateJoined: 1239814800,
    avatarLastModified: 1242927467,
    TsAndCs: "1.2.0",
    guid: "e2cc1b6ed1282775aa64b2aa26150
  - hiddenNotifications: {
        profile_firstlook: "1.0.1"
    },
    activeAvatar: "2009-05-21_17-37-35_286224222.jpg",
    biog: "Hello!"
}
```

no attachments

Actually, in our environment, attachments are usually images or media assets and they really ought not be served from inside a database - so this too is a platform architecture restriction.

# Why CouchDB?
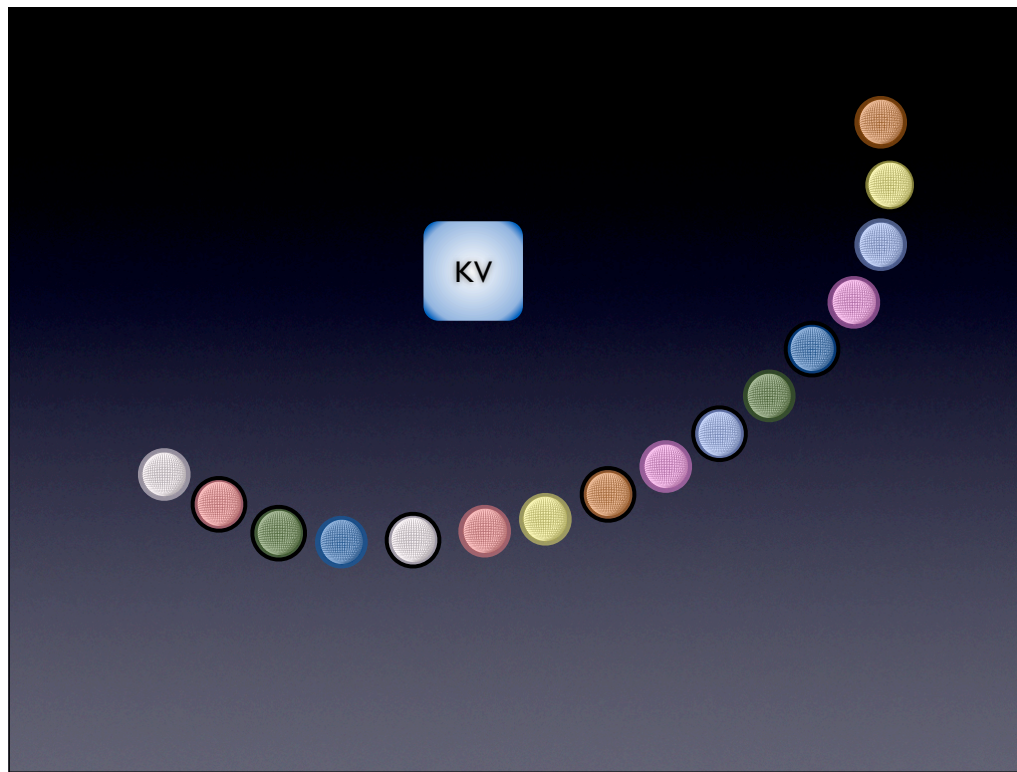
- master master master replication
- operational robustness & consistency
- ease of use

Even though we knew the code wasn't even "beta", and we knew that some high-profile sites would depend on it, it did exactly what it said on the tin, and we could wrap it to stop over-zealous developers (who don't spend enough time thinking about operational impact) using features which may cause headaches
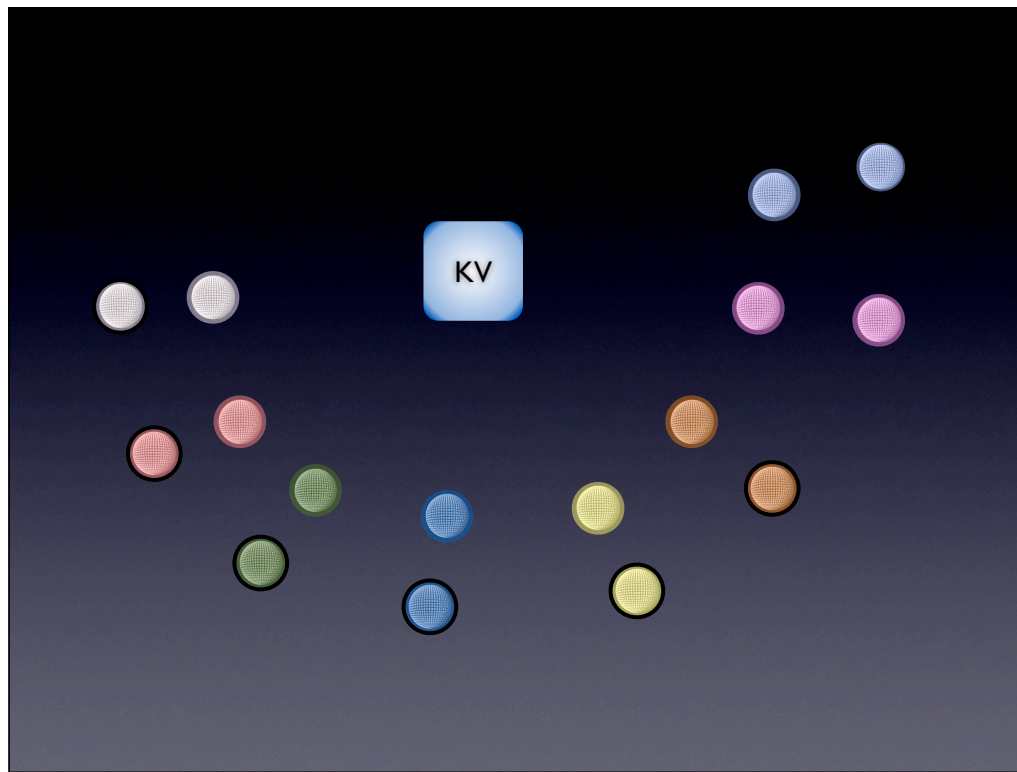
# OuchDB*

just one horror story

Rachel's response to hearing about one of our mishaps

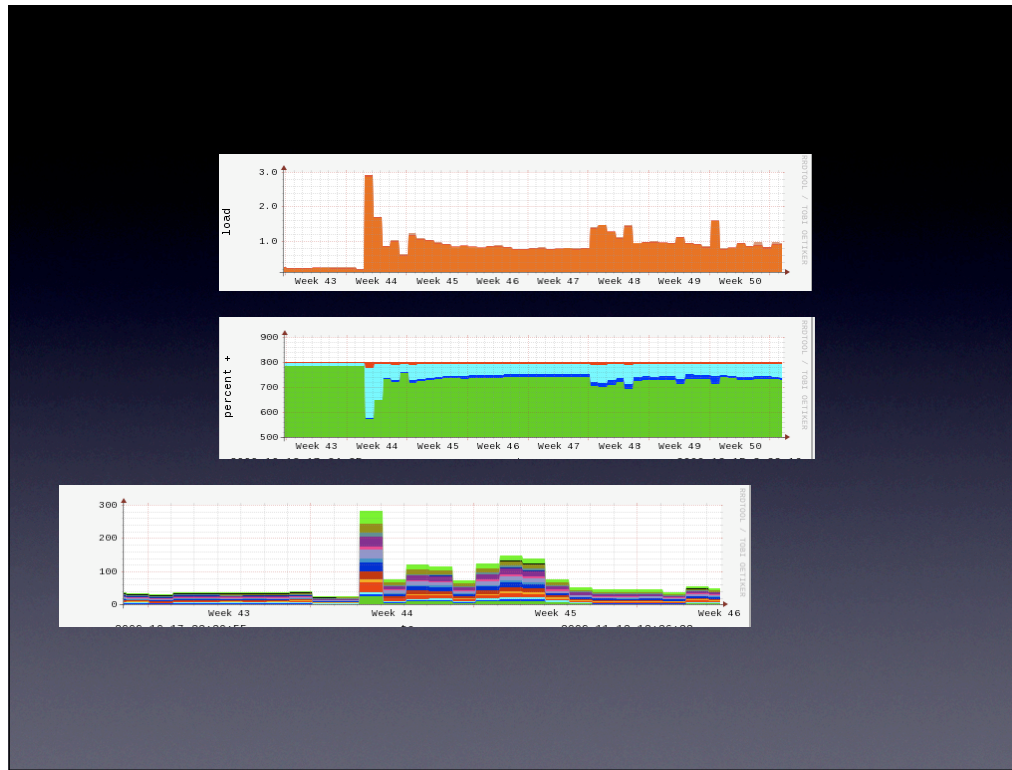At first we did not have hardware redundancy, so had 16 shards

Our intention was to reduce to 8 shards, freeing up the other 8 to be hot-failovers in the event of hardware problems. We had a config error, resulting in us not being able to find some data in one DC and not being able to find other data until replication had finished.

$\frac{2}{8}$ unavailable 50%

$\frac{1}{2}$ unavailable until replication done

The load - partially driven by 404 snowballs - was much higher than expected.

- 3.5% of our requests result in 404s
  - this we consider normal
- Some applications created new docs
- The load was not expected
- Replication ought to have taken 30 min - but with bad config ~ 7+ hours

- What did we do?
- Shut down compaction
  - Kept all revisions of all data
- Many smart folks spent long hours writing scripts to re-assemble data, using these revisions.
- Saved - no data was lost in the end

CouchDB to the rescue! If you don't compact, CouchDB's MVCC can come to the rescue

## er, testing?

- Scaling can be cruel

  traffic:    live 10  1  stage

   docs:     live  8  1  stage

- replication had finished before we noticed
- we now have **lots** of new know-how :-)

This is "glib" in comparison to what we went through, but in summary, it's what mattered.

CouchDB: we like it.

# thank you

twitter: @endafarrell
blog: endafarrell.net