# Swift and Media Storage at Wikimedia

Ben Hartshorne
Operations Engineer
<bhartshorne@wikimedia.org>

# What is Media Storage?

- All images, sounds, and videos on all wikis
- All scaled versions of all those images
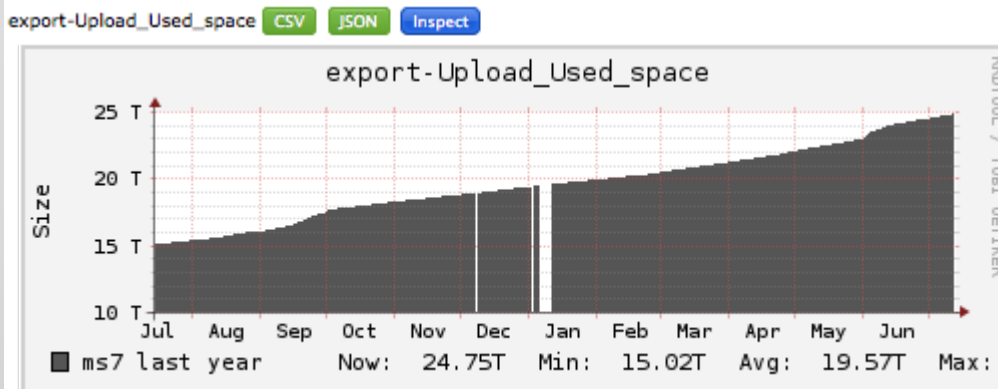- It just keeps growing..

## Commons:MIME type statistics

This page is updated weekly by MIMEStatBot. Any other edits made to

Files on Commons by MIME type as of 2012-07-08 06:01:27 (UTC)

See also: Commons:Project scope/Allowable file types

| MIME type ⇕ | Media type ⇕ | Files ⇕ | Bytes ⇕ |
|---|---|---|---|
| application/ogg | AUDIO | 159,119 | 137,740,897,523 |
| application/ogg | VIDEO | 15,135 | 416,045,941,908 |
| application/pdf | OFFICE | 23,970 | 120,378,314,149 |
| audio/midi | AUDIO | 2,451 | 13,298,263 |
| image/gif | BITMAP | 130,459 | 23,933,529,446 |
| image/jpeg | BITMAP | 11,250,437 | 14,747,665,370,021 |
| image/png | BITMAP | 982,320 | 543,671,948,791 |
| image/svg+xml | DRAWING | 593,115 | 137,535,894,060 |
| image/tiff | BITMAP | 84,722 | 842,335,343,318 |
| image/vnd.djvu | BITMAP | 22,295 | 274,908,628,366 |
| image/x-xcf | BITMAP | 312 | 955,816,108 |
| video/3gpp | UNKNOWN | 1 | 1,015,808 |
| video/mp4 | MULTIMEDIA | 1 | 11,900,528 |
| **Total** | | 13,264,337 | 17,245,197,898,289 |

export-Upload_Used_space  CSV  JSON  Inspect

export-Upload_Used_space

Size

25 T
20 T
15 T
10 T

Jul  Aug  Sep  Oct  Nov  Dec  Jan  Feb  Mar  Apr  May  Jun

RRDTOOL / TOBI OETIKER

■ ms7 last year    Now: 24.75T  Min: 15.02T  Avg: 19.57T  Max:

# Original Media

# Thumbnails (Scaled Media)

# What do we need from a Media Store?

- Large Capacity
    - currently 25TB; minimum for growth: 50-100TB

- Fault Tolerence
    - any component must be able to fail without impact

- Medium Throughput
    - rate of image requests, additions is about 100/s

- Medium Latency
    - most end-user actions are served from cache

# What is Swift?

OpenStack Storage (http://openstack.org/software/openstack-storage/)

<buzzword>scalable fault tolerant object store</buzzword>

- Scalable: Increase cluster size (and throughput) by adding additional hardware; latency shouldn't increase with cluster size
- Fault Tolerant: No single point of failure
- Object Store: Not a filesystem - stores whole objects

# Swift Architecture

Four main
server processes

Traffic

Proxy Server

Container Server

Object Server

Account Server

Many ancillary processes for background jobs: synchronization, auditing,
replication, etc.

# Swift Architecture

Grouped onto
two machine types

Traffic

Proxy Server

Container Server     Object Server

Account Server

# Machine Types

- Frontend Proxy Server
  - dual 6-core CPU
  - 16GB RAM
  - two 250GB SATA disks RAID1

- Backend Storage Server
  - dual 6-core CPU
  - 48GB RAM
  - two 160GB SSDs
  - twelve 2TB SATA disks (no RAID)

# Swift Architecture

Grouped onto
two machine types

Traffic

Proxy Server

Container Server

Account Server

Object Server

SSD

SATA

# Swift Architecture

Grouped onto
two machine types

Traffic

Proxy Server

**5 Proxies**

Container Server

Object Server

**12 Storage**

Account Server

**SATA**

**SSD**

# How is Swift used in MW?



- Thumbnails

- Originals
  - Mediawiki FileBackend class has multiple modules; calls Swift using CloudFiles

# Thanks!

Ben Hartshorne
Operations Engineer
<bhartshorne@wikimedia.org>

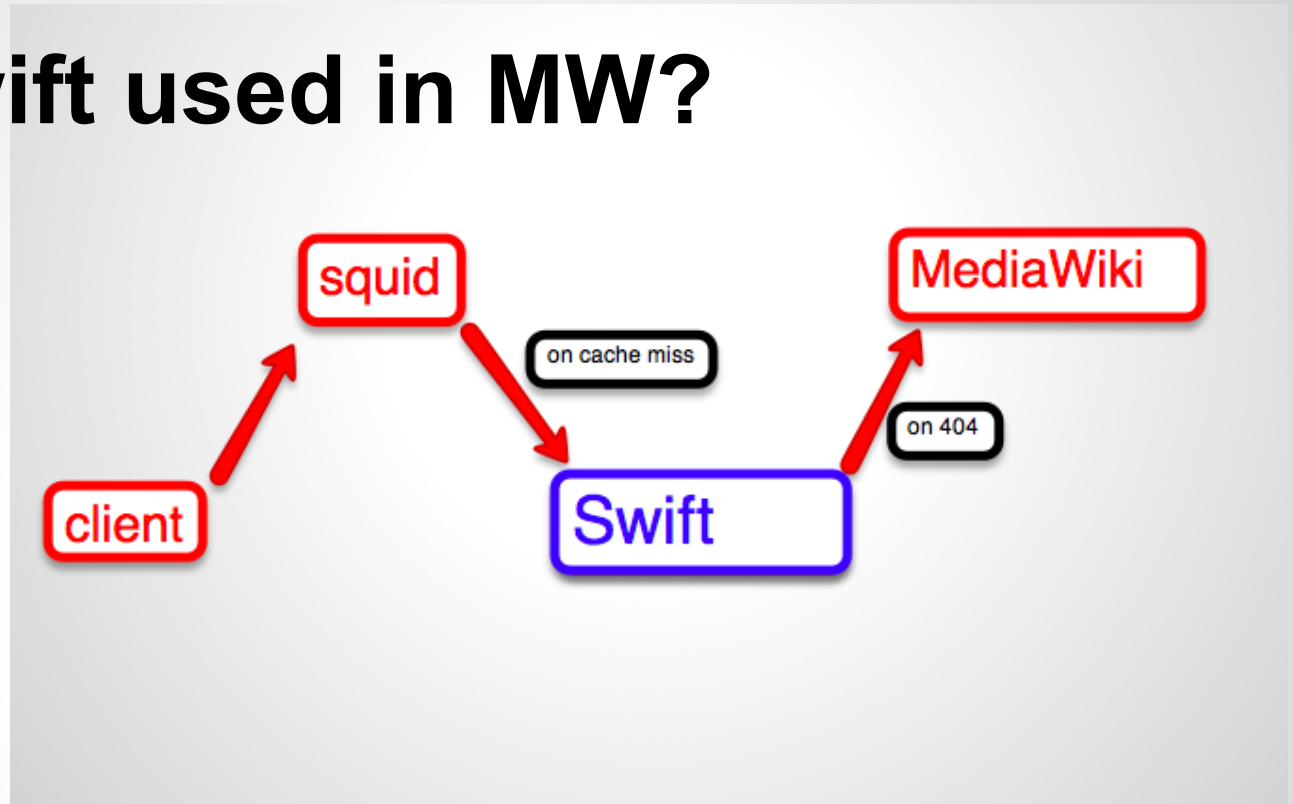http://wikitech.wikimedia.org/view/Swift

# end of presentation

some optional slides follow - they might be used if people ask specific questions

# Rewrite middleware

- New thumbnails are scaled on demand
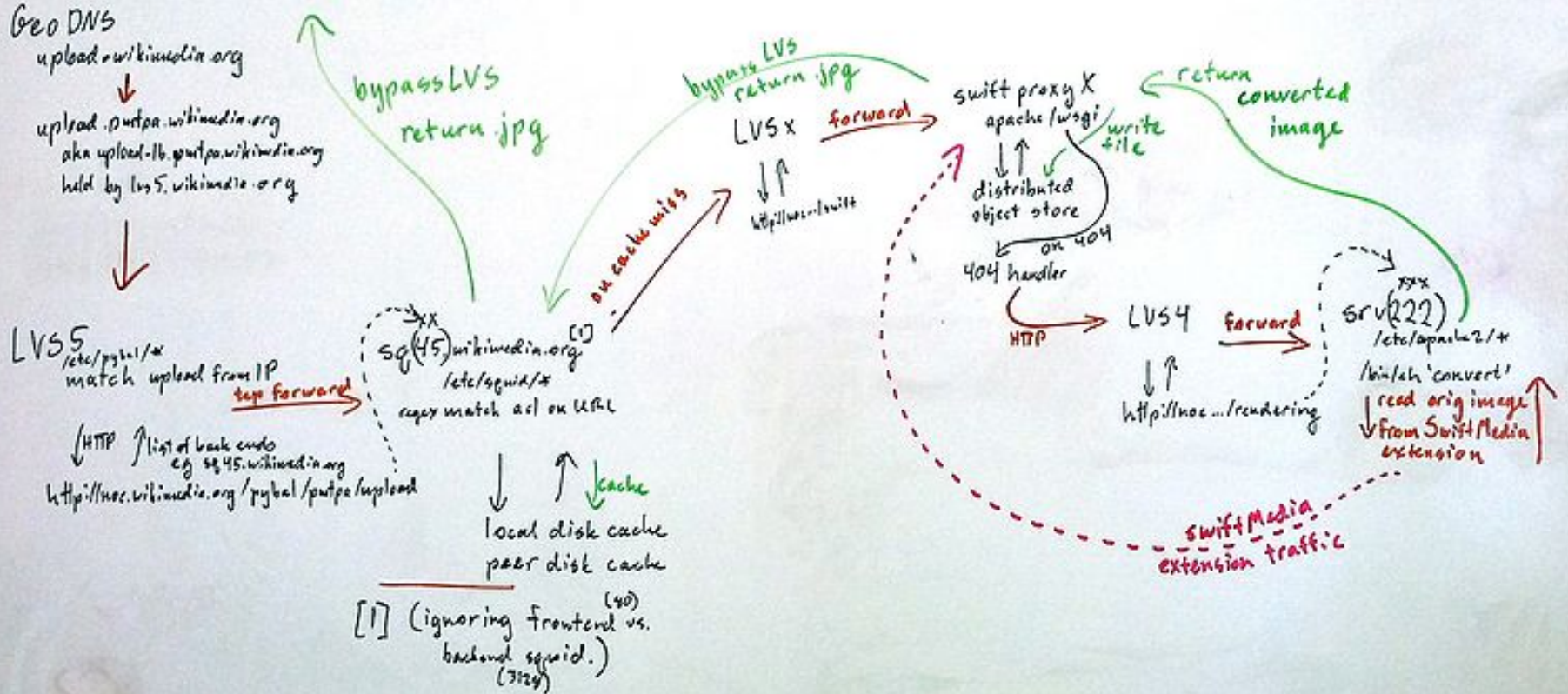- 404 handler tries to scale images
    that don't exist
- swift-proxy is built for this
  - in /etc/swift/proxy-server.conf:
    ```
    [pipeline:main]
    pipeline = rewrite healthcheck cache swauth proxy-server
    ```
- rewrite does two things
  - call back to get the scaled version of the image
  - write that scaled version into swift

# Query flow: client to scaled image

# What about that 404 handler?

Perfect for middleware in the proxy pipeline

```
[pipeline:main]
pipeline = rewrite healthcheck cache swauth proxy-server
```

Rewrite does two things:

- Handle 404s
  - if the object doesn't exist in swift
  - call back to mediawiki to generate the image
  - optionally write the generated image into swift

# What about that 404 handler?

Perfect for middleware in the proxy pipeline

```
[pipeline:main]
pipeline = rewrite healthcheck cache swauth proxy-server
```

Rewrite does two things:

● Change URL into Container / Object

# Integration with Mediawiki

- MW storage mechanisms abstracted to a FileBackend class with multiple subclasses
  - local filesystem, swift, azure, S3, etc.
- All interactions with the FileBackend implemented as appropriate for each backend storage module
- Swift storage implemented using CloudFiles
  - https://github.com/rackspace/php-cloudfiles
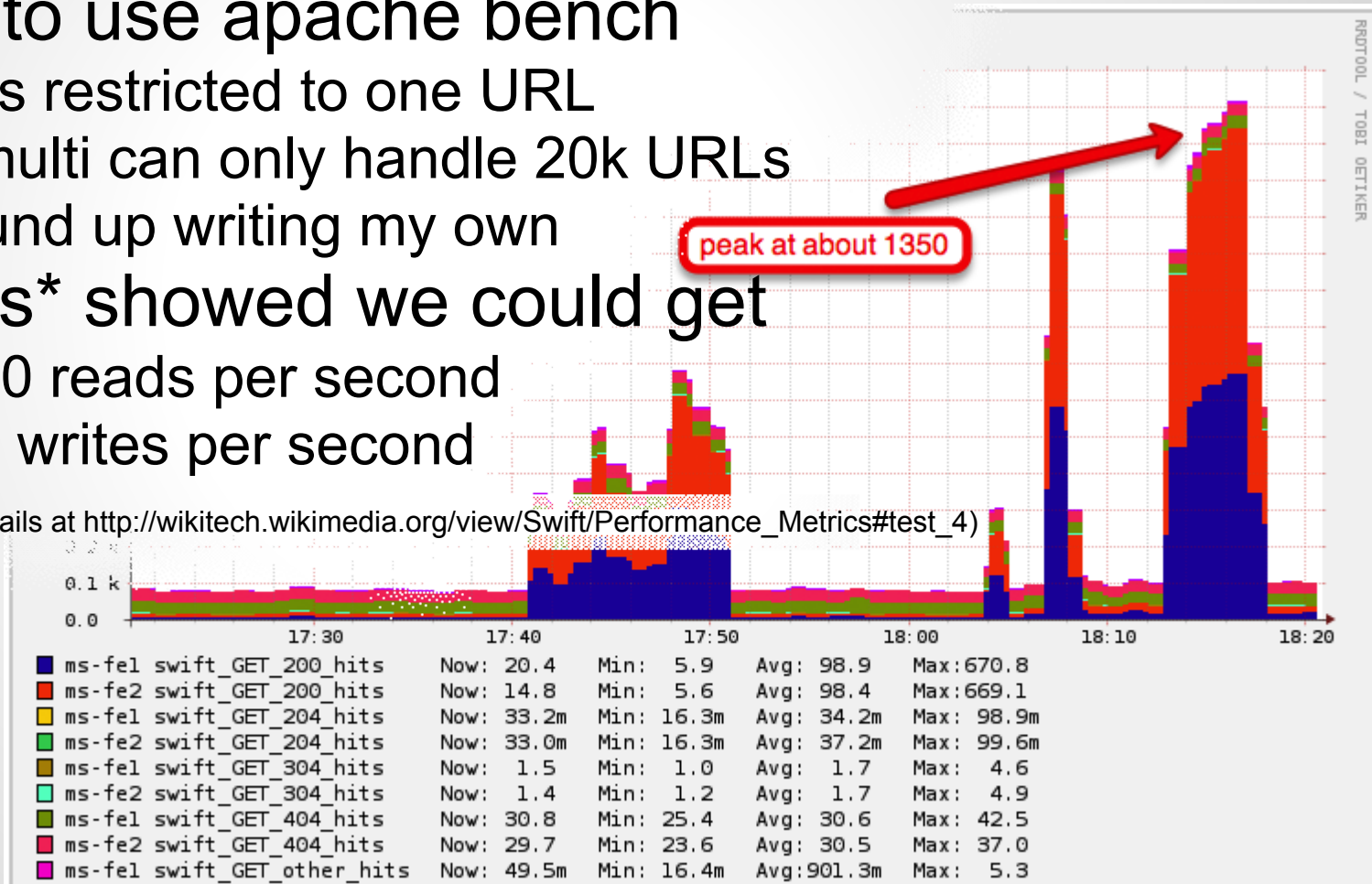- More detail on this part: Aaron Schulz

# Throughput and Latency Performance

# Initial tests

- Tried to use apache bench
  - ab is restricted to one URL
  - abmulti can only handle 20k URLs
  - wound up writing my own
- geturls* showed we could get
  - 1300 reads per second
  - 120 writes per second
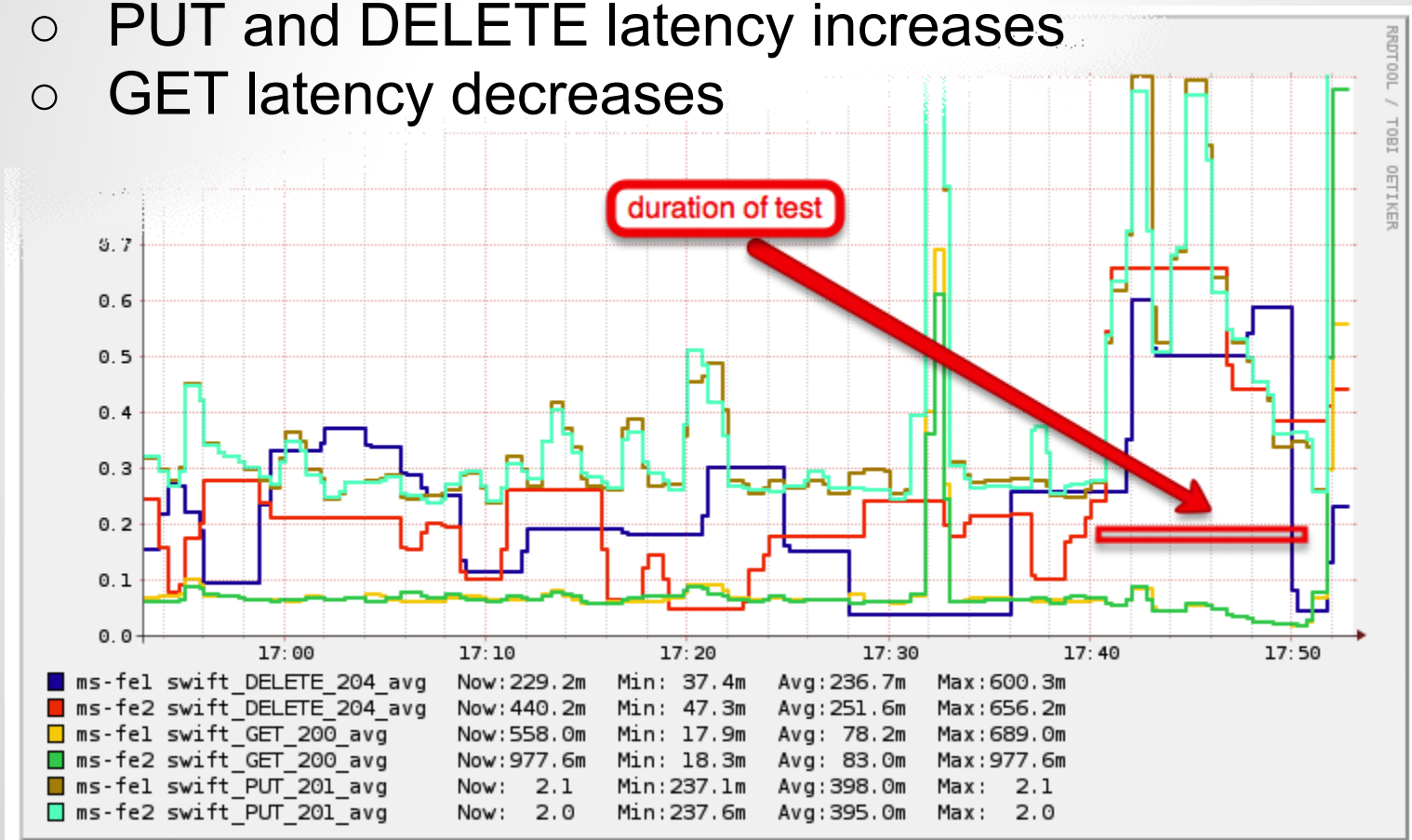  - (full details at http://wikitech.wikimedia.org/view/Swift/Performance_Metrics#test_4)



peak at about 1350

```
                                    Now:      Min:      Avg:       Max:
ms-fe1 swift_GET_200_hits           20.4       5.9      98.9      670.8
ms-fe2 swift_GET_200_hits           14.8       5.6      98.4      669.1
ms-fe1 swift_GET_204_hits           33.2m     16.3m     34.2m      98.9m
ms-fe2 swift_GET_204_hits           33.0m     16.3m     37.2m      99.6m
ms-fe1 swift_GET_304_hits            1.5       1.0       1.7        4.6
ms-fe2 swift_GET_304_hits            1.4       1.2       1.7        4.9
ms-fe1 swift_GET_404_hits           30.8      25.4      30.6       42.5
ms-fe2 swift_GET_404_hits           29.7      23.6      30.5       37.0
ms-fe1 swift_GET_other_hits         49.5m     16.4m    901.3m       5.3
```
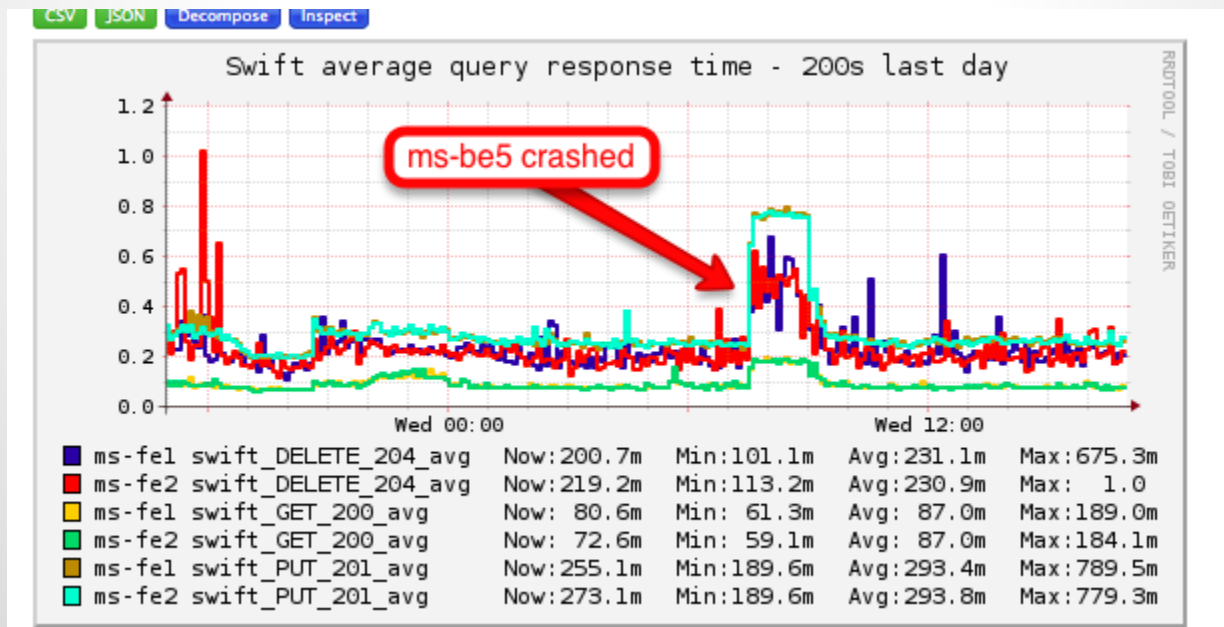
# Effect of load on performance

- Under heavy read load
  - PUT and DELETE latency increases
  - GET latency decreases

# Effect of node failure

- One (out of 5) storage nodes crashing
    - 0.5s timer on connection failures - adjustable
    - 2x read latency (from 100ms to 200ms)
    - 3x write latency (250ms to 750ms)
    - 2.5x delete latency (200ms to 500ms)
- No data (yet) on proxy nodes crashing

# Some problems encountered along the way

- Effect of one storage node crashing on performance is too large
  - solved by reducing the connection timeout from 0.5s to 0.1s
- Container listing latency is high
  - solved by moving container and account servers to SSD leaving objects on spinning media
- Consistency issues with rewrite middleware
  - ETags help
  - Still have issues sometimes (cleaner script)
  - solved by having mediawiki write to swift instead
- It's difficult diagnosing problems with rewrite
  - natural effect of asynchronous code (eventlet)
  - eg. stack trace in proxy logs