# Running Wikipedia.org
## Varnishcon 2016 Amsterdam

## Emanuele Rocca
Wikimedia Foundation

June 17th 2016

# 1,000,000 HTTP Requests

# Outline

- Wikimedia Foundation
- Traffic Engineering
- Upgrading to Varnish 4
- Future directions

WIKIMEDIA
FOUNDATION

# Wikimedia Foundation

- ▶ Non-profit organization focusing on free, open-content, wiki-based Internet projects
- ▶ No ads, no VC money
- ▶ Entirely funded by small donors
- ▶ 280 employees (67 SWE, 17 Ops)

**WIKIMEDIA**
**FOUNDATION**

# Alexa Top Websites

| Company | Revenue | Employees | Server count |
|---|---|---:|---:|
| Google | $75 billion | 57,100 | 2,000,000+ |
| Facebook | $18 billion | 12,691 | 180,000+ |
| Baidu | $66 billion | 46,391 | 100,000+ |
| Yahoo | $5 billion | 12,500 | 100,000+ |
| Wikimedia | $75 million | 280 | 1,000+ |

WIKIMEDIA
FOUNDATION

# Traffic Volume

- Average: ~100k/s, peaks: ~140k/s
- Can handle more for huge-scale DDoS attacks

WIKIMEDIA
FOUNDATION

# DDoS Example



Source: jimieye from flickr.com (CC BY 2.0)

# The Wikimedia Family

# Values

- ▶ Deeply rooted in the free culture and free software movements
- ▶ Infrastructure built exclusively with free and open-source components
- ▶ Design and build in the open, together with volunteers

WIKIMEDIA
**FOUNDATION**

# Build In The Open

- github.com/wikimedia
- gerrit.wikimedia.org
- phabricator.wikimedia.org
- grafana.wikimedia.org

WIKIMEDIA
FOUNDATION

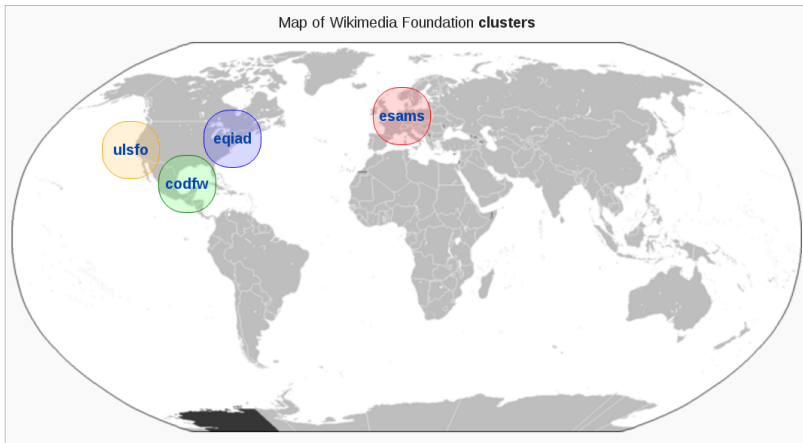# Traffic Engineering

# Traffic Engineering

- Geographic DNS routing
- Remote PoPs
- TLS termination
- Content caching
- Request routing

WIKIMEDIA
**FOUNDATION**

# Component-level Overview

- DNS resolution (gdnsd)
- Load balancing (LVS)
- TLS termination (Nginx)
- In-memory cache (Varnish)
- On-disk cache (Varnish)

# Cluster Map



Map of Wikimedia Foundation **clusters**

eqiad: Ashburn, Virginia - cp10xx
codfw: Dallas, Texas - cp20xx
esams: Amsterdam, Netherlands - cp30xx
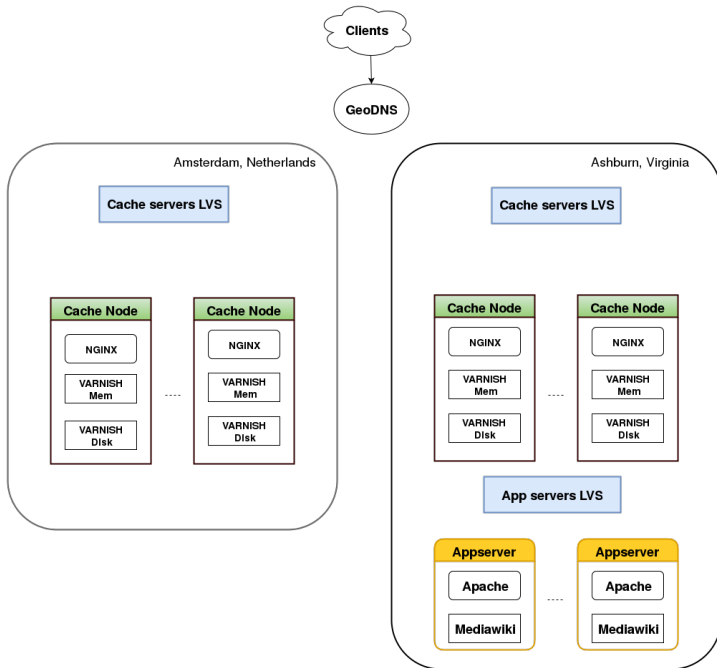ulsfo: San Francisco, California - cp40xx

# CDN

- No third-party CDN / cloud provider
- Own IP network: AS14907 (US), AS43821 (NL)
- Two "primary" data centers
  - Ashburn (VA)
  - Dallas (TX)
- Two caching-only PoPs
  - Amsterdam
  - San Francisco

WIKIMEDIA
FOUNDATION

# CDN

- ► Autonomy
- ► Privacy
- ► Risk of censorship

# CDN

- ▸ Full control over caching/purging policy
- ▸ Lots of functional and performance optimizations
- ▸ Custom analytics
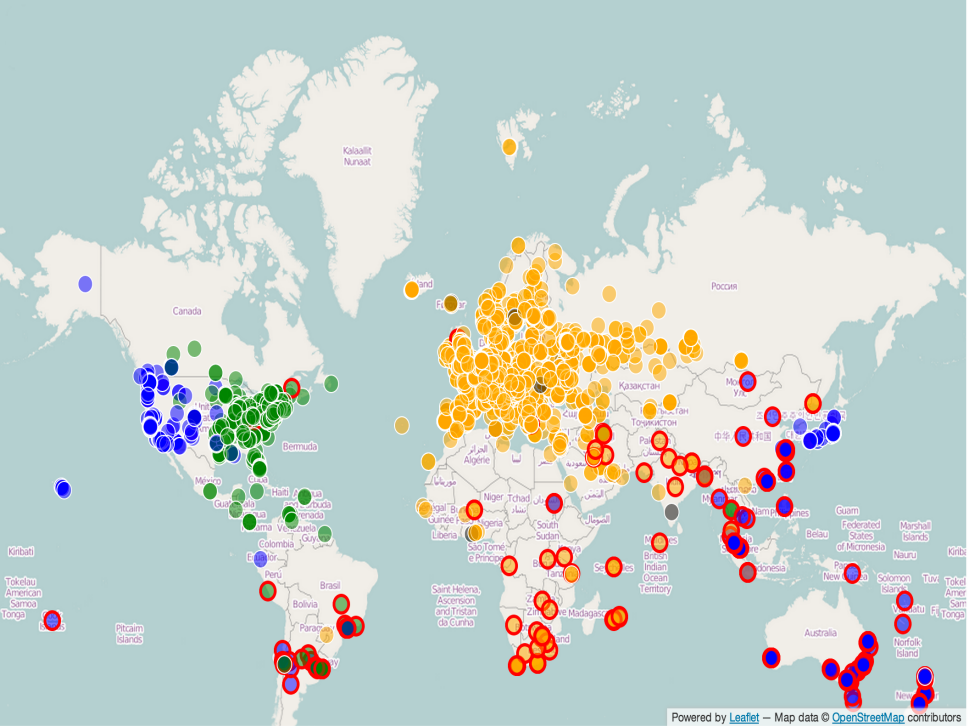- ▸ Quick VCL hacks in DoS scenarios

WIKIMEDIA
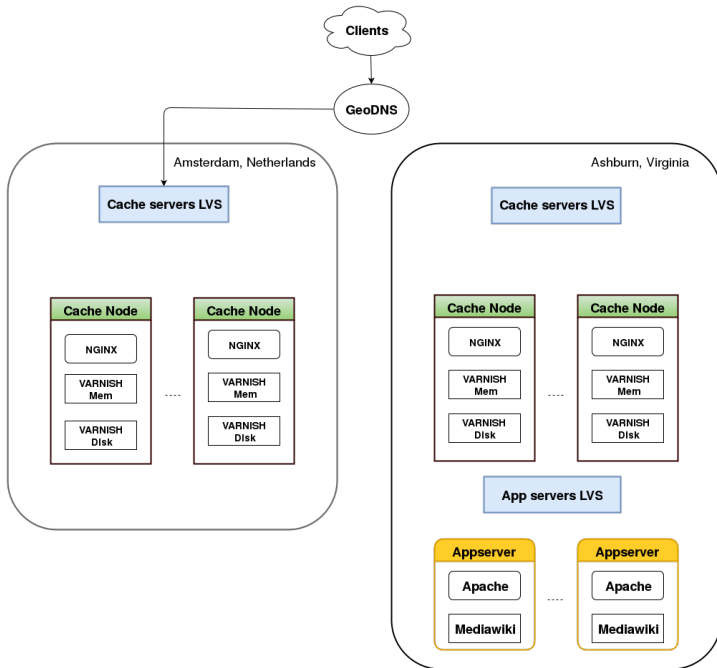**FOUNDATION**

17

# GeoDNS

- 3 authoritative DNS servers running gdnsd + geoip plugin
- GeoIP resolution, users get routed to the "best" DC
- edns-client-subnet
- DCs can be disabled through DNS configuration updates

# config-geo

```
FR => [esams, eqiad, codfw, ulsfo], # France
JP => [ulsfo, codfw, eqiad, esams], # Japan
```
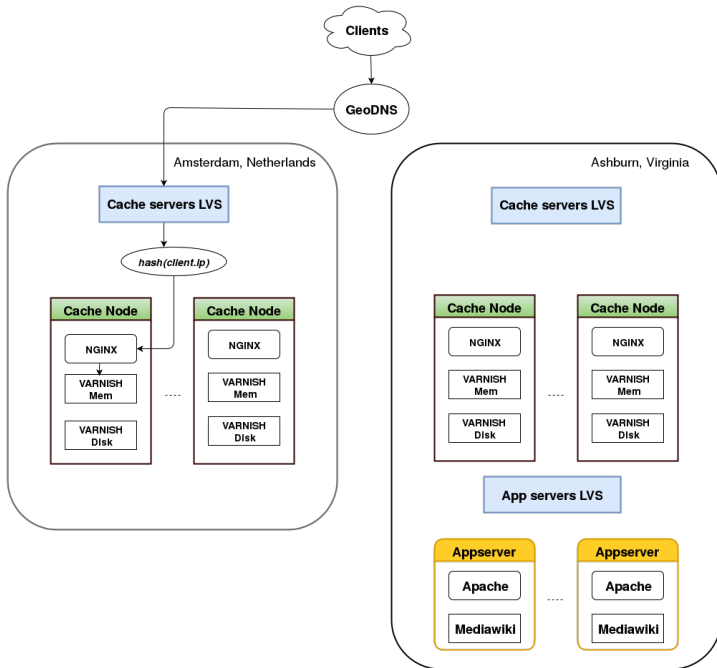
https://github.com/wikimedia/operations-dns/

Clients

GeoDNS

Amsterdam, Netherlands

Cache servers LVS

Cache Node
- NGINX
- VARNISH Mem
- VARNISH Disk

....

Cache Node
- NGINX
- VARNISH Mem
- VARNISH Disk

Ashburn, Virginia

Cache servers LVS

Cache Node
- NGINX
- VARNISH Mem
- VARNISH Disk

....

Cache Node
- NGINX
- VARNISH Mem
- VARNISH Disk

App servers LVS

Appserver
- Apache
- Mediawiki

....

Appserver
- Apache
- Mediawiki

21

# LVS

- ▶ Nginx servers behind LVS
- ▶ LVS servers active-passive
- ▶ Load-balancing hashing on client IP (TLS session persistence)
- ▶ Direct Routing

WIKIMEDIA
**FOUNDATION**

# Pybal

- Real servers are monitored by a software called Pybal
- Health checks to determine which servers can be used
- Pool/depool decisions
- Speaks BGP with the routers
  - Announces service IPs
  - Fast failover to backup LVS machine
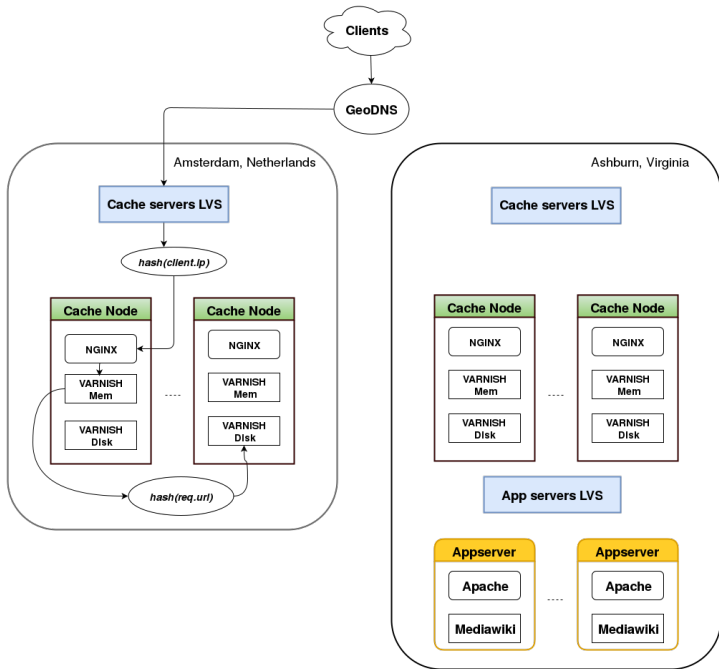
WIKIMEDIA
FOUNDATION

# Pybal + etcd

- ▸ Nodes pool/weight status defined in etcd
- ▸ confctl: CLI tool to update the state of nodes
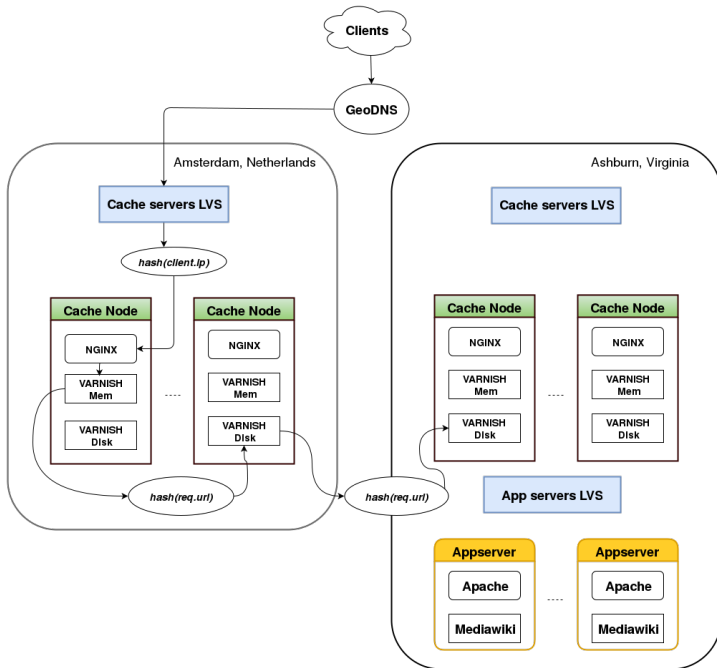- ▸ Pybal consuming from etcd with HTTP Long Polling

WIKIMEDIA
FOUNDATION

# Nginx + Varnish

- 2x varnishd running on all cache nodes
  - :80 -smalloc
  - :3128 -spersistent
- Nginx running on all cache nodes for TLS termination
- Requests sent to in-memory varnishd on the same node

WIKIMEDIA
FOUNDATION

# Persistent Varnish

- ► Much larger than in-memory cache
- ► Survives restarts
- ► Effective in-memory cache size: ~avg(mem size)
- ► Effective disk cache size: ~sum(disk size)

WIKIMEDIA
FOUNDATION

Clients

GeoDNS

**Amsterdam, Netherlands**

Cache servers LVS

hash(client.Ip)

Cache Node
- NGINX
- VARNISH Mem
- VARNISH Disk

Cache Node
- NGINX
- VARNISH Mem
- VARNISH Disk

....

hash(req.url)

hash(req.url)

**Ashburn, Virginia**

Cache servers LVS

Cache Node
- NGINX
- VARNISH Mem
- VARNISH Disk

Cache Node
- NGINX
- VARNISH Mem
- VARNISH Disk

....

App servers LVS

Appserver
- Apache
- Mediawiki

Appserver
- Apache
- Mediawiki

....

29

# Inter-DC traffic routing

```
cache::route_table:
    eqiad:  'direct'
    codfw:  'eqiad'
    ulsfo:  'codfw'
    esams:  'eqiad'
```

# Inter-DC traffic routing

- Varnish backends from etcd:
  directors.vcl.tpl.erb
  - puppet template -> golang template -> VCL file
- IPSec between DCs

WIKIMEDIA
FOUNDATION

# X-Cache

Cache miss:

```
$ curl —v https://en.wikipedia.org?test=$RANDOM 2>&1 | grep X—Cache
X—Cache: cp1068 miss, cp3040 miss, cp3042 miss
```

# X-Cache

Cache miss:

```
$ curl -v https://en.wikipedia.org?test=$RANDOM 2>&1 | grep X-Cache
X-Cache: cp1068 miss, cp3040 miss, cp3042 miss
```

Cache hit:

```
$ curl -v https://en.wikipedia.org | grep X-Cache
X-Cache: cp1066 hit/3, cp3043 hit/5, cp3042 hit/21381
```
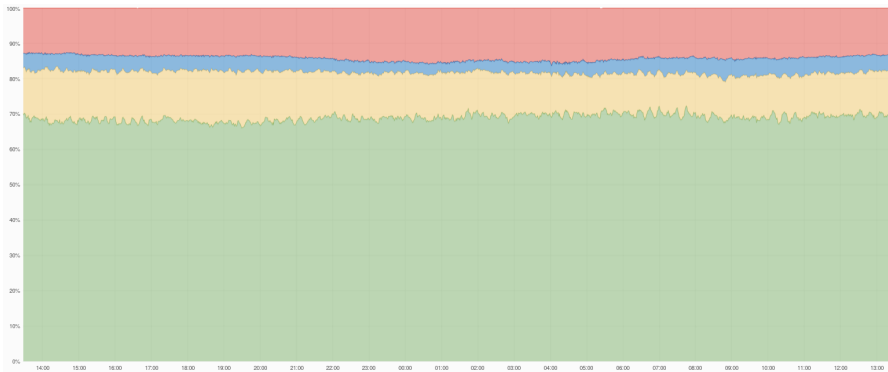
# X-Cache

### Cache miss:

```
$ curl −v https://en.wikipedia.org?test=$RANDOM 2>&1 | grep X−Cache
X−Cache: cp1068 miss, cp3040 miss, cp3042 miss
```

### Cache hit:

```
$ curl −v https://en.wikipedia.org | grep X−Cache
X−Cache: cp1066 hit/3, cp3043 hit/5, cp3042 hit/21381
```

### Forcing a specific DC:

```
$ curl −v https://en.wikipedia.org?test=$RANDOM \
    −−resolve en.wikipedia.org:443:208.80.153.224 2>&1 | grep X−Cache
X−Cache: cp1066 miss, cp2016 miss, cp2019 miss
```

# Cache clusters

- ► Text: primary wiki traffic
- ► Upload: multimedia traffic (OpenStack Swift)
- ► Misc: other services (phabricator, gerrit, …)
- ► Maps: maps.wikimedia.org

WIKIMEDIA
FOUNDATION

# Terminating layer - text cluster



- ▶ Memory cache: 69%
- ▶ Local disk cache: 13%
- ▶ Remote disk cache: 4%
- ▶ Applayer: 14%

WIKIMEDIA
**FOUNDATION**

# Terminating layer - upload cluster



- ▶ Memory cache: 68%
- ▶ Local disk cache: 29%
- ▶ Remote disk cache: 1%
- ▶ Applayer: 2%

WIKIMEDIA
FOUNDATION

# Upgrading to Varnish 4

# Varnish VCL

- Puppet ERB templating on top of VCL
- 22 files, 2605 lines
- Shared across:
    - clusters (text, upload, ...)
    - layers (in-mem, on-disk)
    - tiers (primary, secondary)
- 21 VTC test cases, 715 lines

WIKIMEDIA
FOUNDATION

# Varnish 3

- 3.0.6-plus with WMF patches
  - consistent hashing
  - VMODs (in-tree!)
  - bugfixes
- V3 still running on two clusters: text and upload

WIKIMEDIA
**FOUNDATION**

# Varnish 4 upgrade

- Bunch of patches forward ported
- VMODs now built out-of-tree
- VCL code upgrades
- Custom python modules reading VSM files forward ported
- Varnishkafka

V4 running on two clusters: misc and maps

# V4 packages

- Official Debian packaging:
  `git://anonscm.debian.org/pkg-varnish/pkg-varnish.git`

- WMF patches:

  `https://github.com/wikimedia/operations-debs-varnish4/`

  `tree/debian-wmf`

- Need to co-exist with v3 packages (main vs. experimental)

- APT pinning

WIKIMEDIA
FOUNDATION

# VMODs

- ▶ vmod-vslp replacing our own chash VMOD
- ▶ vmod-netmapper forward-ported
- ▶ Packaged vmod-tbf and vmod-header

WIKIMEDIA
FOUNDATION

# V4 VMOD porting

# V4 VMOD packaging

- Modifications to vmod-tbf to build out-of-tree
  - Header files path
  - Autotools

- vmod-header was done already, minor packaging changes

WIKIMEDIA
FOUNDATION

# VCL code upgrades

- Need to support both v3 and v4 syntax (shared code)
- Hiera attribute to distinguish between the two
- ERB variables for straightforward replacements
    - $req_method $\to$ req.method vs. req.request
    - $resp_obj $\to$ resp vs. obj
    - ...
- 42 if @varnish_version4

WIKIMEDIA
FOUNDATION

# varnishlog.py

- Python callbacks on VSL entries matching certain filters

- Ported to new VSL API using python-varnishapi: https://github.com/xcir/python-varnishapi

- Scripts depending on it also ported
    - TxRequest → BereqMethod
    - RxRequest → ReqMethod
    - RxStatus → BereqStatus
    - TxStatus → RespStatus

WIKIMEDIA
FOUNDATION

# varnishkafka

- Analytics
- C program reading VSM files and sending data to kafka
- https://github.com/wikimedia/varnishkafka
- Lots of changes:
- 6 files changed, 612 insertions(+), 847 deletions(-)

WIKIMEDIA
FOUNDATION

# varnishtest

- Started using it after Varnish Summit Berlin
- See ./modules/varnish/files/tests/
- Mocked backend (vtc_backend)
- Include test version of VCL files
- VCL code depends heavily on the specific server

WIKIMEDIA
FOUNDATION

```
[...]

varnish v1 -arg "-p vcc_err_unref=false" -vcl+backend {
    backend vtc_backend {
        .host = "${s1_addr}"; .port = "${s1_port}";
    }

    include "/usr/share/varnish/tests/wikimedia_misc-frontend.vcl";
} -start

client c1 {
    txreq -hdr "Host: git.wikimedia.org" -hdr "X-Forwarded-Proto: https"
    rxresp
    expect resp.status == 200
    expect resp.http.X-Client-IP == "127.0.0.1"

    txreq -hdr "Host: git.wikimedia.org"
    rxresp
    # http -> https redirect through _synth, we should still get X-Client-IP
    # (same as in _deliver)
    expect resp.status == 301
    expect resp.http.X-Client-IP == "127.0.0.1"
} -run
```

# Future plans

# Future plans - TLS

- ▸ Outbound TLS
- ▸ Add support for listening on unix domain socket

WIKIMEDIA
**FOUNDATION**

# Future plans - backends

- ▸ Make backend routing more dynamic: eg, bypass layers on pass at the frontend
- ▸ etcd-backed director to dynamically depool/repool/re-weight

WIKIMEDIA
**FOUNDATION**

# Future plans - caching strategies

- Only-If-Cached to probe other cache datacenters for objects before requesting from the applayer
- XKey integration to "tag" different versions of the same content and purge them all at once (eg: desktop vs. mobile)

WIKIMEDIA
FOUNDATION

# Future plans - bloom filters

Very fast and space-efficient way to find out if something is definitely not in the set

- ▶ cache-on-second-fetch: avoid caching "rare" items
- ▶ 404 filter with the bloom set representing all legal URLs to help against randomized URL paths from botnets

WIKIMEDIA
FOUNDATION

# Conclusions

- One of the most popular CDNs in the world is built in the open using FOSS
- Multi-layered Varnish setup
- Currently upgrading to Varnish 4
- Big plans for the future!

# Cache servers

101 bare-metal servers
- ▸ 28 Amsterdam
- ▸ 27 Virginia
- ▸ 26 Texas
- ▸ 20 California

WIKIMEDIA
**FOUNDATION**

# edns-client-subnet

```python
import dns
import clientsubnetoption

def resolve(client_ip):
    cso = clientsubnetoption.ClientSubnetOption(client_ip)
    message = dns.message.make_query('en.wikipedia.org', 'A')
    message.use_edns(options=[cso])

    # ns0.wikimedia.org
    r = dns.query.udp(message, '208.80.154.238')

    for a in r.answer:
        print a

print "United States"
resolve('199.217.118.41')

print "Italy"
resolve('151.1.1.1')
```

WIKIMEDIA
FOUNDATION

# edns-client-subnet

```
$ python resolve.py
United States
en.wikipedia.org. 600 IN A 208.80.153.224
 Italy
en.wikipedia.org. 600 IN A 91.198.174.192
```