

Scaling up: How we made millions of domains happier

Tom Arnfeld, DNS Engineer Pavel Odintsov, DNS Engineer

Cloudflare DNS Services

- Authoritative DNS
 - IPv4 and IPv6 support
 - Low latency propagation
 - DNSSEC on the fly signing
 - CNAME flattening
 - Load balancing
- DNS Firewall (former VDNS), DNS proxy with bundled security
- Rich Analytics





Cloudflare DNS Scale

100+ Points of Presence Thousands of servers



Cloudflare DNS Scale 5.5m+DNS Zones

.om .vivo .rent .men .prof .gmail .final .games .1WC .army .bz mt lpl .dad .de .ga .jetzt .java .shell .fm .adac wang . **now .now .now .now .now** farm .01.bn .osaka .hiv .ceb .lr .gm .1] .tmall .boats .VU .Za .jr i .bio .cf .hot .in .tube .parts .cy .tours .food





Cloudflare DNS Scale ~1M QPS Throughput







DNS Packets/s

.::

Cloudflare DNS Scale O(100M) QPS Attack Throughput









rrDNS



What is rrDNS?

- Written in Golang
- Powered by Miek Gieben's library
- Modular design for different query paths
- Detailed statistics and logging for queries and responses
- Uses a custom globally replicated key-value store



Issues with rrDNS v1

- Poor architecture (no strict modularity)
- Bloated and slow storage encoding (protobuf)
- Everything stored as separate RR's
- Poor performance
- Unpredictable behaviour in certain complex cases
- Hard to maintain



Goals for rrDNS v2 business

- Scale to serve >= 100 million zones
- Better query performance
- Easier development of new features / products
- No single switch when migrating to v2





Goals for rrDNS v2 technical

- More scalable data encoding
- Improve correctness against DNS RFC's
- Optimise based on experience and numbers
- Predictable and documented business logic
- Support for arbitrary RR types (only require changes in our UI)
- Considerably more maintainable, testable and well documented code





rrDNS v2 technical decisions

- Completely reworked design
- Totally new, optimised storage and encoding
- MessagePack (binary json encoding) as new storage encoding
- Faster paths for common query types (A, AAAA, NS)
- Predictable behaviour with >95% test coverage





rrDNS v2 achievements

- Increased max throughput by 3-4x
- Improved correctness against DNS Clear RFCs
- Mixed-mode DNS for a smooth, safe migration
- We love this code





rrDNS v2 RFC correctness

- Fixed DS query for unsigned delegation, DVE-2017-0008 https://github.com/DNS-OARC/dns-violations/blob/master/2017/DVE-2017-0008.md
- Do not return authority flag for delegations
- Do not return useless glue records





DNS Data Pipeline



What is our Data Pipeline?

Primary DB **Serialization & Business Logic** Replicated **KV** Store



CC-BY 2.0 image by Steve Jurvetson



Data Pipeline v1

- Extremely complicated and ancient design
- Written in PHP/C/C++/bash
- No documentation or test coverage
- Very poor performance, especially for larger zones
- Impossible to offer guaranteed processing time
 - Falls over with large bursts of changes



Data Pipeline v2

- More efficient encoding/storage
- Drastically improve performance
- Guarantee low propagation time for changes
 - Even for zones with millions of records
- Predictable behaviour
- Handle bursts of changes gracefully



Data Pipeline v2 - what did we get?

- Written in pure Go
- Scales linearly with number of CPU cores
- ~95% test coverage
- Extremely well documented
- Lots of new debugging tools for SREs/Support/Operations
- Guaranteed propagation time
- 10x performance increase





Monitoring Propagation before







Monitoring Propagation after







Migrating V1 \rightarrow V2





How do we know it works?

- and compare responses
- Full comparison Iterate over all zones and records to compare answers from v1 and v2, ensuring a "safe" migration
- **Tests** Cover every new discovered case with a test, *religiously*
- Smart Comparisons Wrote "smart comparators" to filter acceptable differences, because of bad behaviour in v1





• **DNS mirror** Feed a sample of production queries to both old and new

Migration Process

- Tools to perform validation and migration of each zone automatically
- Treated like a prod. service meaning metrics, logging etc.
- Many iterations fixing bugs to improve the validation pass rate







DNS Mirror

- Samples real DNS traffic from entire Cloudflare network
- Compares responses between multiple instances of rrDNS
- Tools for debugging detected differences
- Part of every rrDNS release to minimise regressions



DNS Mirror



Migrating to v2

- Total consumed time (development)
- Time consumed by migration
- Number of new test cases
- Number of Pull Requests
- Number of rrDNS releases
- Total number of working hours consumed by project
- Percent of zones with noticeable impact

~1 year < 3 months 183 250+ ~50 (about one per week) ~3 years 0.005%

Thank you. Questions?

https://blog.cloudflare.com/how-we-made-our-dns-stack-3x-faster/

DNS is hard!

