# Systems Engineering Update
# 2021 AGM, The Ether

Matthew Pounsett

2021/10/08

## Contents

# 1   Introduction

While I have made progress on several projects in the last half year, my focus continues to be on the challenges of day to day operations and of handling unexpected issues. I'm continuing to make incremental changes to general operations that increase stability, and simplify day to day operations, which has the effect of gradually making more time available for more changes. However, we're still highly constrained in what we're able to accomplish due to the significant mismatch between what we need—or would like—to be doing, and the time available to do it in.

Although we anticipated that travel might be back to normal by now, I'm still finding work on hardware in Fremont and Ottawa to be challenging. Since the last contractor we had completed his work, we're back to relying on datacenter remote hands for anything we need done on-site, which tends to require more advance planning, and comes with more uncertainty about results.

We're still hunting for people to fill short and medium term contract vacancies that would help with some of these issues, and we're just beginning the hunt for an additional engineer for a long term position, to help cover more of the general operations work and allow me to concentrate more on larger projects intended to carry OARC forward.

# 2   OARC Services Overview

## 2.1   Data Archiving

OARC maintains a large store of multiple data sets.

**Day in the Life** OARC coordinates annual, and occasionally ad-hoc, Day in the Life (DITL) DNS traffic capture events. These involve many operators of significant DNS infrastructures—including root server operators, TLDs, and recursive operators—running packet captures of their traffic over the same 48 hour period. The data are uploaded to OARC where it is organized for use in research.

The DITL collections go back to 2009.

**RSSAC 002 Statistics** The Root Server System Advisory Committee's publication RSSAC 002 is the Advisory on Measurement of the Root Server System. It defines an initial set of statistics to be collected by root server operators from their systems. OARC collects the output of this reporting from each root server operator, daily, and maintains a history of these statistics available for analysis or review.

**Zone File Repository** OARC maintains an historical archive of zone files which includes daily updates of the root zone going back to 1993, and weekly updates of several TLDs beginning at various times between 2009 and 2018.

**Other Data** OARC also periodically accepts submissions of other data that may be relevant to researchers interested in the DNS:

- derivative data from research done on OARC's other datasets
- data collected from OARC testing tools, such as the DNS Entropy Tester
- DITL-like collections from outside regular DITL windows, such as occasional contributions from AS112 server operators
- packet captures from OARC's Open DNSSEC Validating Resolver (ODVR) which includes forwarded queries from the DNS Privacy Testbed
- Case Western Reserve University's "Case Connection Zone" FTTH data
- other ad-hoc contributions of relevant data

## 2.2  File Servers and Storage

OARC's datasets are stored on several discrete file servers. Five of the file servers, located in Fremont, California, have 343TB used of their 505TB of capacity. Two of these have multiple filesystems, marked as A and B in the chart below. The sixth file server, located in Ottawa, Ontario, is an off-site copy of a selection of datasets from the first five servers.

I am in the process of building a new Ceph storage cluster, discussed in detail in previous reports, and covered again below in section 3.5, which will replace the old discrete NFS servers. Due to their age, the old servers have been having increasing trouble with stability in the last few years. We are currently "borrowing" one of the new storage servers as a temporary addition to the older servers in order to cover for lack of available storage on the couple of stable servers we still have.

There is one file server, not listed below, which lost its filesystem due to multiple disk failures in too quick succession for us to recover the RAID. Another, FS5, listed as temporarily offline, experienced what we believe to be either a CPU or memory failure which is preventing it from booting at the moment. We will cannibalize other servers for parts in order to boot it when it comes time to migrate its data to the new cluster.

| Server/Volume | Used | Capacity | Notes |
|---|---|---|---|
| FS2a | 36TB | 42TB | |
| FS2b | 75TB | 125TB | |
| FS3 | 34TB | 42TB | |
| FS4 | 72TB | 84TB | |
| FS5a | 69TB | 84TB | Temporarily offline |
| FS5b | 33TB | 42TB | Temporarily offline |
| Stor09 | 24TB | 86TB | borrowed from new cluster |
| FS6 | 117TB | 121TB | Located in Ottawa, Canada |

Each file server uses either ZFS (RaidZ2) or XFS over software RAID for its filesystem to provide redundancy within the file server. Each dataset is intended to be stored on more than one file server in order to create cross-chassis redundancy of data; and, due to the server in Ottawa, some datasets currently have copies on three systems.

Due to the loss of the `FS1` server, not all datasets have their expected number of copies available, and with `FS5` offline, some datasets have either no copy accessible or are only accessible on the Ottawa server.

All capacity numbers above are the filesystem capacity, rather than the raw size of the disks in service.

## 2.3 Data Analysis Servers

OARC maintains four Linux shell servers with access to the above data sets. Three in Fremont, CA (an1, an2, an4) and one in Ottawa, ON (an3). Members and Supporters who have signed a Data Sharing Agreement and request access are given accounts on these analysis servers, which they can use to do research into the DNS using any of OARC's datasets.

> **Note Well**: No data, even derived data, may leave OARC analysis systems without express written authorization, in compliance with the Data Sharing Agreement. Contact admin@dns-oarc.net first, *always*.

## 3 System and Service Status

### 3.1 General Condition

The list of high-maintenance services continues to shrink, as we either retire services that the community has lost interest in, or do cleanup and reimplementation to make maintenance easier. The list of remaining work is still very long, but I believe I'm still managing to remove things faster than they are added.

I am still unable to travel to Fremont, and unable to conveniently travel to Ottawa, so hardware maintenance still relies on remote hands at the local datacenters. This arrangement is more reliable in Fremont than it is in Ottawa, which ironically means that the more distant site (from my perspective) is getting more care and attention.

We are still on the lookout for short term contractors to handle some specific tasks, including doing the leg-work to specify and deploy new edge routers for the OARC sites, as well as an on-call contract for emergency response in Fremont. We are also just beginning a search for a second long-term Systems/Network engineer to take up some of the regular maintenance load, and

free my time up for more long term project work. We expect to be posting job requirements for that soon.

Please see the OARC careers page for more detail.

## 3.2  Authoritative DNS Improvements

In March, I deployed a new hidden primary name server which is allowing us to more easily manage our DNSSEC-signed zones, and to begin to do regular key rolls for our signed zones. The first zone on this server, mentioned in the previous report, was `dnsviz.net`. In July, I moved `dns-oarc.net` to this server, and almost immediately encountered an issue with the zone going BOGUS. This turned out to be a bug in the KSK roll logic of the software, which removed the old KSK prematurely after finding the new DS record present in the parent. The Knot developers quickly resolved the bug, and we have been humming along with the new version ever since.

I allowed a fairly long wait period after upgrading to the patched version before planning new zone migrations, but by that point we were into summer vacations, and preparing for workshops and the AGM, so those have continued to be put off. I'm intending to resume migrating zones next month (November), with the intention of turning off the old signer at least a couple of weeks before winter vacations begin in December.

Beginning in January I will begin changing the signing profiles of zones, triggering algorithm rolls from RSASHA256 to ECDSAP256SHA256.

## 3.3  Discontinuing the Shared DSC Platform

OARC's shared DNS Statistics Collector (DSC) platform has been discontinued, and the last member contributing data has turned off their feed.

Early in DNS-OARC's existence, members configured their DSC instances to feed data to a shared setup on OARC's servers, allowing members to compare traffic levels and other statistics among themselves. Some found this useful for seeing trends across the industry, such as whether a traffic spike was seen by multiple root servers, or TLD operators, or not.

In the two years since OARC replaced its member portal, there has been no way to view this shared repository of data. Beginning earlier than that, the number of contributors to the dataset had begun to drop off.

This time last year, there was only one member other than OARC itself contributing data, and no way to view it. After asking the community for feedback on whether OARC should continue this data collection activity, and hearing none, we decided it was time to retire the service. It was shut down this summer, and its data archived on OARC's file server infrastructure with other historical DNS data.

## 3.4   Networks, Routing, and Routers

We have been unable to find a contractor to take on the task of specifying and deploying new routing equipment for OARC's sites. This means that we are now nearly a year behind the desired schedule for getting our routers replaced, and deploying our first redundant transit connection. We now have a couple of recommendations and expressions of interest, but unfortunately the workload that Keith and I have had in recent weeks has made it impossible to progress discussions. We are hoping to be able to resume that search within the next few months.

At the same time, we will soon be opening head-count for an additional Network/Systems engineer which may obviate the need for a short term contractor.

Details are available on OARC's web site on our careers page.

## 3.5   File Servers

### The Old Stuff

Little has changed with the old storage servers since the last member report. FS5 continues to be offline (but expected to be recoverable when necessary), and our other servers continue to be full or nearly full.

I have temporarily loaned one of the storage servers from the new cluster to the old NFS architecture in order to deal with new, incoming datasets. This server's data (currently, the 2021 DITL collection) will be among the first data migrated to the new Ceph clustered storage when it is in production, and the server will then rejoin the cluster.

The problem of having no off-site replication of the data housed on our file servers remains an issue, and the double failures of FS1 and FS5 highlight the danger. We are constrained by both the Data Sharing Agreement, which prevents us from copying data into long term cloud storage, and by available budget, which prevents us from duplicating our existing storage architecture at a new site.

We are hoping that the Board's Privacy Committee, which has seen renewed activity, will allow us some options.

### The New Stuff

The pace of implementation of the new Ceph cluster has slowed a bit during this past summer. We have run into a few delays, but the main problem has been related to time limitations and conflicts with tasks that cannot be deprioritized. With my attention pulled in many different directions, it has been difficult to focus on working thorough problems. The long schedule on this project basically comes down to my failure to make two predictions accurately: the time it would take to deploy, and the time that would available to work on deployment.

We temporarily solved the problem of Dell misreporting the power requirements of the hardware by ordering in extra circuits. We're currently running two 20A circuits per cabinet in order to bring the load down under 80% per circuit. Once easy travel to Fremont is possible, I'll re-deploy the two cabinets of hardware into three, and we'll go back down to one circuit per cabinet.

Dell contacted us this summer to let us know that the hardware we received contained drives with a fatal defect in their firmware that could result in catastrophic data loss. This prompted a BIOS/firmware patching effort that ate up a fair amount of time, as many of the servers were 11 (or more) BIOS versions behind, along with additional firmware updates that needed to be applied to other components (such as the disks). Following Dell's recommended procedure, the first couple of servers took approximately four hours each to update, which would have led to a couple of weeks of server patching. I'd like to thank Dan Mahoney of ISC for coming to the rescue with pointers to a Dell tool that was actually installable on Debian hosts, and which we could use to semi-automate the firmware updates. In addition to bringing the process down to 30-40 minutes per server, the tool found firmware updates that I'd been unable to locate in Dell's catalogue.

The remainder of the summer has been spent dealing with the sort of unpredictable surprises that are common when setting up a complex system for the first time. We haven't quite reached production capability yet, but I'm expecting it to be not too much longer before we get to that point. Unfortunately, predicting how much longer is required to complete the project is going to be like predicting the length of a piece of string; I won't know what problems are still waiting to be found, or how long they will take to resolve until I get there.

## 3.6   Day in the Life Dataset

Over the last few years we have seen a growing delay between the end of a DITL collection event and the release of the data to researchers. This has been due to a number of reasons. During the last couple of years we were having issues with hardware stability which meant we were unable to run the regular post-processing. This year, we "borrowed" hardware from the new Ceph cluster to store the 2021 collection, but were initially prevented from running the post-processing due to a power shortage in the server's cabinet; running the processing tools would have created enough load to trip the breaker.

Once we dealt with the power problem, we immediately started the data cleaning process, but ran into some anomalies which we felt required correction or explanation before releasing the data to researchers. These sorts of anomalies are becoming more common, and to prevent them from causing further delays in the future we're making a few changes in our procedures around DITL data.

First, we are going to begin immediately releasing the raw, un-cleaned dataset to researchers as soon as we are sure that all contributors have completed their data upload. Due to the fact that some contributors run infras-

tructure in distant, poorly connected parts of the world, this is sometimes a week or two after the end of the DITL collection window.

Second, even when there are anomalies in the cleaning process, we will release cleaned data as soon as we have about three quarters of the raw dataset successfully cleaned. We will withhold any components of the dataset that are exhibiting anomalies we cannot explain until they are either corrected, or explained sufficiently clearly that we can attach a README document to that part of the data set.

In addition to reaching out to individual data contributors to help explain why we see odd interactions between our cleaning process and their data, we will also be enlisting the help of Supporters—looking for in-kind contributions they can make in exchange for their membership in OARC—to investigate these anomalies.

An email was sent to the Members' mailing list on October 6th explaining these changes, our cleaning process, and related circumstances in much more detail. I recommend interested parties refer to that email for more detail, or contact me by email at matt@dns-oarc.net or on OARC's Mattermost as `@matt` with questions.