# The Internet of Names: A DNS Big Dataset

## Actively Measuring 50% of the Entire DNS Name Space, Every Day

Roland van Rijswijk-Deij[1,2], Mattijs Jonker[1], Anna Sperotto[1], Aiko Pras[1]
[1]University of Twente, [2]SURFnet bv
{r.m.vanrijswijk, m.jonker, a.sperotto, a.pras}@utwente.nl

## ABSTRACT

The Domain Name System (DNS) is part of the core infrastructure of the Internet. Tracking changes in the DNS over time provides valuable information about the evolution of the Internet's infrastructure. Until now, only one large-scale approach to perform these kinds of measurements existed, passive DNS (pDNS). While pDNS is useful for applications like tracing security incidents, it does not provide sufficient information to reliably track DNS changes over time. We use a complementary approach based on active measurements, which provides a unique, comprehensive dataset on the evolution of DNS over time. Our high-performance infrastructure performs Internet-scale active measurements, currently querying over 50% of the DNS name space on a daily basis. Our infrastructure is designed from the ground up to enable big data analysis approaches on, e.g., a Hadoop cluster. With this novel approach we aim for a quantum leap in DNS-based measurement and analysis of the Internet.

## Categories and Subject Descriptors

C.4 [**Performance of Systems**]: Measurement Techniques

## Keywords

DNS; active measurements; big data; Internet evolution

## 1. INTRODUCTION

Next to IP, DNS is arguably the most important infrastructure on the Internet. DNS is pervasive as almost all networked applications and services rely on DNS to map names to IP addresses. Consequently, measuring what is in the DNS can teach us a lot about the state of the Internet. If performed systematically over time, such measurements allow us to observe the evolution of the Internet.

The applications of measuring DNS over time are myriad. An important area of application is network security. Knowledge of what names an IP address mapped to in the past, for example, can be a valuable tool to track malicious
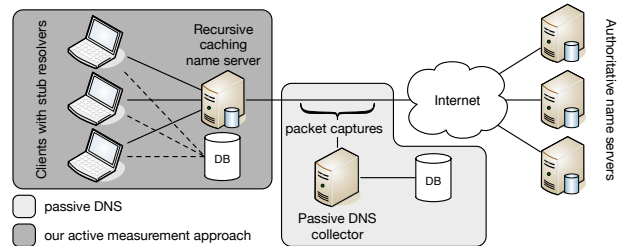
**Figure 1: pDNS compared to our approach**

activity. There are also many applications in network research. Knowledge of DNS content over time provides empirical data about operational practices and deployment of new protocols. If, for example, we wanted to answer the question how the use of cloud e-mail providers develops over time, knowledge of the DNS is vital (as who handles mail for a domain is configured in DNS through the `MX` record type).

### Passive DNS

Development of the only existing large-scale approach to DNS measurements, passive DNS (pDNS) [1], was driven by security benefits. Fig. 1 shows an abstract view of the DNS. pDNS typically collects data on the link between recursive caching name servers ("resolvers") and authoritative name servers (light grey area). Data from pDNS setups[1] can, e.g., be used to track names associated with IP addresses that exhibit malicious activity. From a network research perspective, pDNS is also of interest but it suffers from one problem: it does not provide reliable data over time. This is because 1) pDNS will only record data for domains in which clients behind the resolvers where pDNS data is collected are interested and 2) pDNS has no influence over temporal spacing of queries.

### Our approach

Driven by a research need for reliable data from the DNS over time, we have developed a complementary approach to pDNS, based on active measurements. Given the DNS zone files from top-level domains (TLDs) as input, we send a fixed selection of queries for each domain in a TLD once per 24 hours. Effectively, if we compare our approach to pDNS, we control the behaviour of the clients performing queries (shown in dark grey in Fig. 1). It is highly challenging to make such an approach scale. For example, the `.com` domain alone (the largest TLD on the Internet), already contains >116M names. In the remainder of this poster abstract, we provide a brief outline of our approach and we highlight the potential of the resulting dataset with a case study.

---

[1]e.g., DNSDB – https://www.dnsdb.info/

**Figure 2: High-level infrastructure overview**

| TLD | #domains | #workers | avg. time | #queries/day | data/day |
|---|---|---|---|---|---|
| .org | 10.5 mil. | 10 | 6h45m | 127 mil. | 2.4GB |
| .net | 15.0 mil. | 10 | 13h30m | 181 mil. | 3.3GB |
| .com | 116.5 mil. | 80 | 17h30m | 1427 mil. | 27.2GB |
| *total* | 142.0 mil. | 100 | - | 1735 mil. | 32.9GB |

**Table 1: Active measurement characteristics**

## 2. INFRASTRUCTURE

### 2.1 High-level overview

Fig. 2 gives a bird's eye view of our measurement setup. We divide the measurement process into three stages:

- **Stage 1:** input collection – in this stage, we collect the DNS zones for the TLDs to measure. We compute daily deltas and track both the active zone content as well as changes in the zone over time in a database.

- **Stage 2:** main measurement – this is the active measurement stage; we will explain this stage in more detail in Sec. 2.2 below.

- **Stage 3:** aggregate and prepare for analysis – in this final stage we convert the output from the measurement to the Parquet columnar storage format, which is well-suited for processing on a Hadoop cluster.

### 2.2 Main measurement

Our main active measurement runs on a cloud-based cluster. Every TLD measurement is orchestrated by a cluster management host. This host is responsible for distributing chunks of work, of 100k domains each, to a set of worker nodes. Each worker node runs custom-built software that performs a pre-defined selection of DNS queries for each domain in a chunk of work. Queries are performed against a local DNS resolver instance running on the worker node. Data collected by workers is sent to a central aggregation point for further processing and analysis. Tab. 1 shows an overview of our current setup. It shows the TLDs we measure, the number of worker nodes, the average time to complete a full measurement, the average number of queries per day, and the amount of data collected per day.

### 2.3 Analysis

As mentioned in Sec. 2.1, we store data such that it is well-suited for processing on a Hadoop cluster. We use such a cluster for analysis, defining map/reduce operations on the data, but also more advanced forms of analysis using, e.g., the Impala massively parallel query engine[2].

---

[2]http://www.cloudera.com/content/cloudera/en/products-and-services/cdh/impala.html
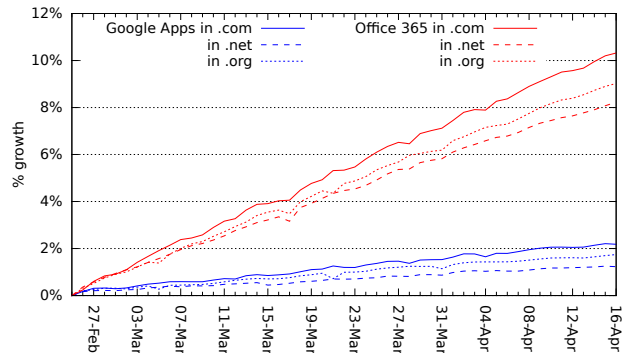


**Figure 3: Use of cloud mail platforms over time**

## 3. CASE STUDY

As case study, we have analysed the use of cloud mail platforms over a 50-day period. We focused on two particular platforms, Microsoft's Office 365 and Google Apps. Fig. 3 shows the growth in the fraction of domains per TLD that use either of these platforms. Growth is presented as a percentage relative to the start of the 50-day period. To perform the analysis, our platform processed over 84 billion query results. The full analysis was performed by a single 40-core node in about 7.5 hours. This can easily be improved by running the analysis on a larger Hadoop cluster.

This simple example showcases what can be achieved using our measurement platform and data. The growth in use of, e.g., cloud mail platforms illustrates how the Internet is evolving from every organisation managing its own services to a few large providers offering these services in bulk.

## 4. CONCLUSIONS AND FUTURE WORK

We created a unique active measurement infrastructure for the DNS. Our infrastructure actively measures over 50% of the total DNS name space on a daily basis. The resulting dataset enables reliable DNS-based analysis of the evolution of the Internet for the first time. And not only do we measure on a large scale, we have also carefully designed for optimal analysis of the collected data through the Hadoop toolchain. The simple case study included in this abstract showcases use of our dataset. It answers the simple question about cloud mail platforms that we provided as an example in the introduction.

The goal of this poster is to invite other researchers to collaborate with us, to analyse this unique new dataset. To provide insight into its potential, we plan to create a web portal with daily statistics. Furthermore, we have already started several research projects that investigate Internet phenomena and that rely on measurement data from this platform.

## 5. ACKNOWLEDGEMENTS

## 6. REFERENCES

[1] Florian Weimer. Passive DNS Replication. In *Proc. of the 17th FIRST Conference (FIRST 2005)*, 2005.