

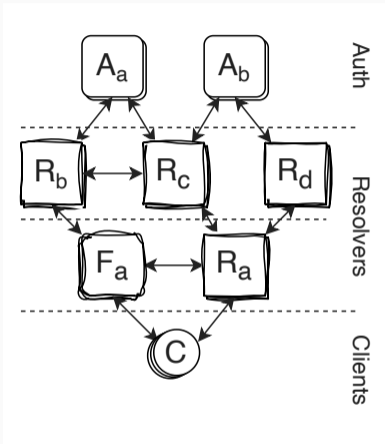
Challenges on Working with DNS Data

Alfred Arouna¹² • Mattijs Jonker³ • Ioana Livadariu¹
alfred@simula.no

¹Simula Metropolitan • ²Oslo Metropolitan University • ³University of Twente

February 22, 2023 • DNS and Internet Naming Research (DINR) 2023 • Virtual Workshop

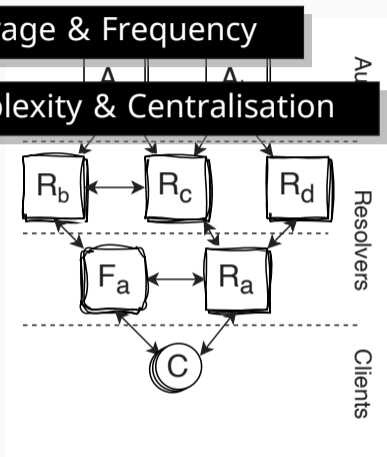
DNS Infrastructure: Highly Distributed and Verbose



(a) Distributed

Simplified Overview of the DNS Infrastructure

DNS Infrastructure: Highly Distributed and Verbose



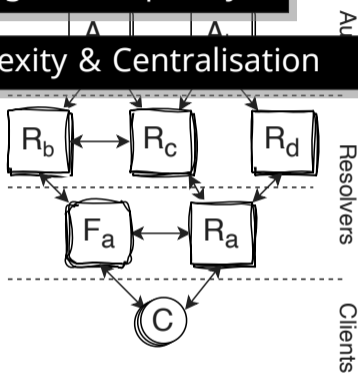
(a) Distributed

Simplified Overview of the DNS Infrastructure

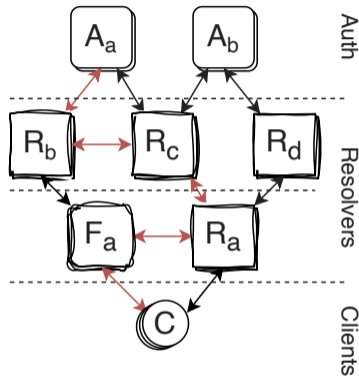
DNS Infrastructure: Highly Distributed and Verbose

Coverage & Frequency

Complexity & Centralisation



(a) Distributed



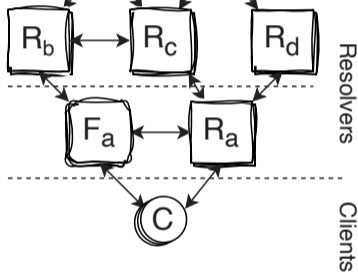
(b) Verbose

Simplified Overview of the DNS Infrastructure

DNS Infrastructure: Highly Distributed and Verbose

Coverage & Frequency

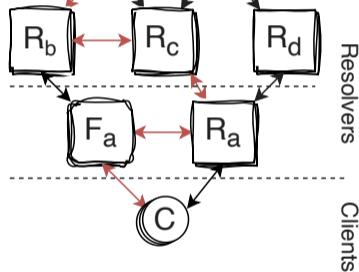
Complexity & Centralisation



(a) Distributed

Privacy & Confidentiality

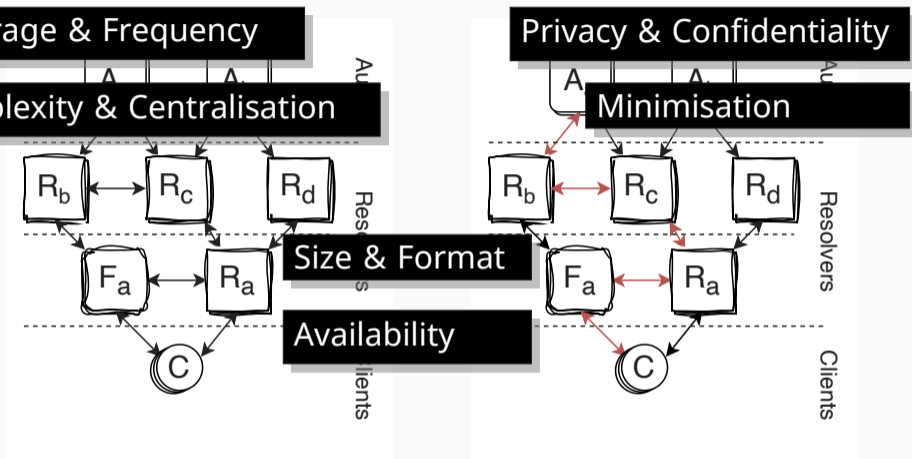
Minimisation



(b) Verbose

Simplified Overview of the DNS Infrastructure

DNS Infrastructure: Highly Distributed and Verbose

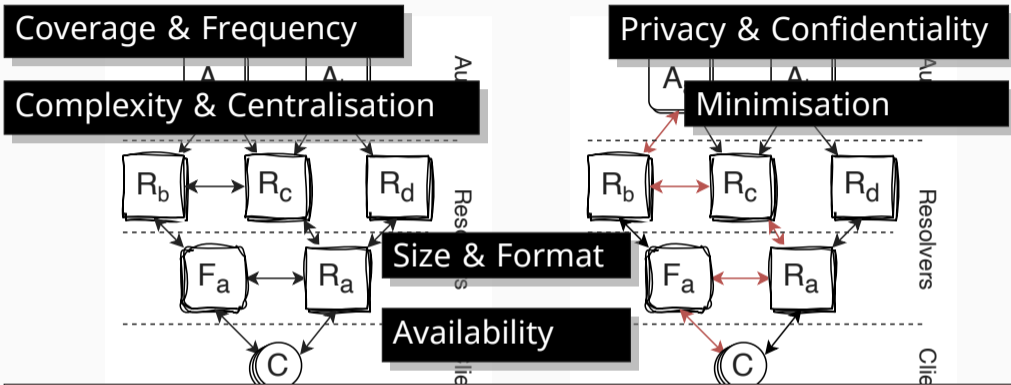


(a) Distributed

(b) Verbose

Simplified Overview of the DNS Infrastructure

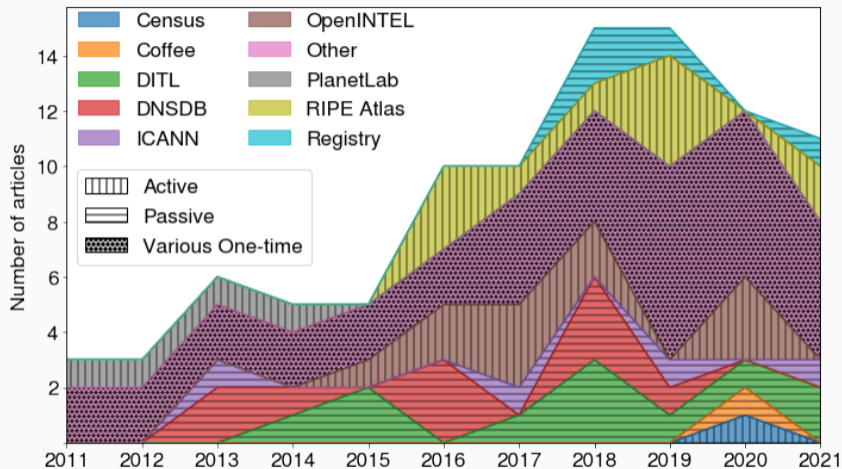
DNS Infrastructure: Highly Distributed and Verbose



How has the DNS research community so far managed to address these data collection challenges?
Evaluate existing datasets usage from most impacting publications on DNS.

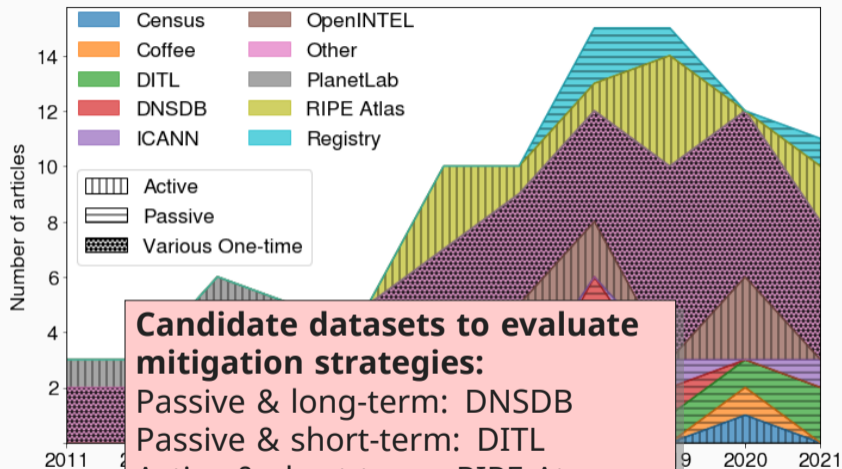
Simplified Overview of the DNS Infrastructure

Popular DNS Datasets



DNS datasets usage over the last 10 years. Although the increase of publications can be correlated with the rise of long-term datasets, DNS researchers relied in majority on one-time snapshot of the state of (parts of) the DNS.

Popular DNS Datasets



DNS datasets used can be correlated with the rise of long-term datasets, DNS researchers relied in majority on one-time snapshot of the state of (parts of) the DNS. increase of publications

Distributed: Addressing Challenges

Table: Mitigation approaches for challenges resulting from the distributed DNS infrastructure. Long-term datasets have large coverage based on consistent data collection frequency.

	Coverage	Frequency	Complexity	Centralisation
OpenINTEL	Large	Consistent	High: Workers	Low
RIPE Atlas	Limited	Variable	High: probes	Low
DNSDB	Large	Consistent	High: VPs	Medium
DITL	Limited	Variable	High:servers	Medium

Distributed: Addressing Challenges

Table: Mitigation approaches for challenges resulting from the distributed DNS infrastructure. Long-term datasets have large coverage based on consistent data collection frequency.

	Coverage	Frequency	Complexity	Centralisation
OpenINTEL	Large	Consistent	High: Workers	Low
RIPE Atlas	Limited	Variable	High: probes	Low

Dataset should apply a distributed systematic collection data at an established frequency with minimal burden on the DNS.

Verbose: Addressing Challenges

Table: Mitigation approaches for challenges resulting from DNS verbosity. Addressing privacy and confidentiality is still challenging. However, minimization can help to reduce privacy and confidentiality risks. Active dataset by controlling the resolver are less impacted by minimization.

	Privacy	Confidentiality	Minimization
OpenINTEL	High: zones	High: zones	Medium
RIPE Atlas	High: user	High: users	Medium
DNSDB	Medium ¹	Medium ¹	High
DITL	Medium ²	Medium ²	High

¹resolver-authoritative

²root-servers

Verbose: Addressing Challenges

Table: Mitigation approaches for challenges resulting from DNS verbosity. Addressing privacy and confidentiality is still challenging. However, minimization can help to reduce privacy and confidentiality risks. Active dataset by controlling the resolver are less impacted by minimization.

	Privacy	Confidentiality	Minimization
OpenINTEL	High: zones	High: zones	Medium
RIPE Atlas	High: user	High: users	Medium

Dataset should consider the increasing adoption of the principle of minimum disclosure to minimize privacy and confidentiality risks

¹resolver-authoritative

²root-servers

Distributed & Verbose: Addressing Challenges

Table: Mitigation approaches for challenges resulting from DNS distributed infrastructure and verbosity. The variability in data formats limits for large-scale and long-term analysis. Active dataset are publicly available while passive dataset access are on demand.

	Format	Size	Availability
OpenINTEL	Avro/Parquet	+10TB ³	Public
RIPE Atlas	JSON	+25TB ⁴	Public
DNSDB	ISC/dnsqr	n.a.	Restricted
DITL	PCAP	n.a.	Restricted

³10TB of compressed data as of Feb. 2015 (1 year data).

⁴July 2015 in Hadoop/HBase (5 years of data).

Distributed & Verbose: Addressing Challenges

Table: Mitigation approaches for challenges resulting from DNS distributed infrastructure and verbosity. The variability in data formats limits for large-scale and long-term analysis. Active dataset are publicly available while passive dataset access are on demand.

	Format	Size	Availability
OpenINTEL	Avro/Parquet	+10TB ³	Public
RIPE Atlas	JSON	+25TB ⁴	Public

The variability in data formats limits for large-scale and long-term analysis. Compacted-DNS (C-DNS) seems to be a good candidate for a common data format. However, Avro/Parquet has been proven for large-scale analysis as part of DNS big data pipeline.

³ 10TB of compressed data as of Feb. 2015 (1 year data).

⁴ July 2015 in Hadoop/HBase (5 years of data).

Comments?

**Ideas to improve the work?
Questions?**