

# DNS Traffic Analysis

## Opportunities, Risks, and (Self-)Defenses

Utility for forensic investigations

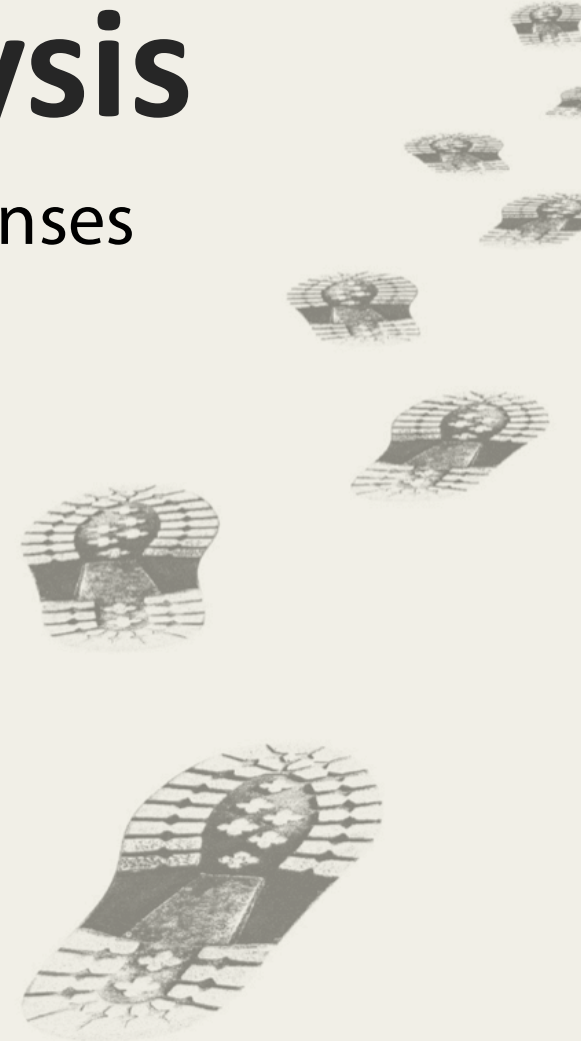
Potential threats to privacy

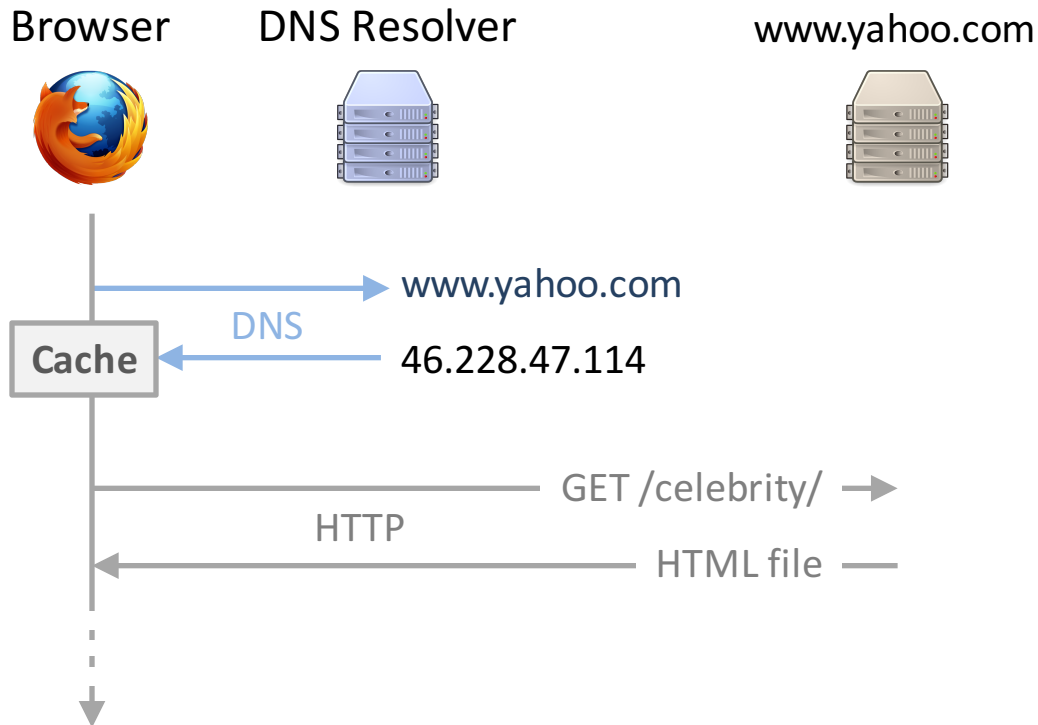
New ideas for protection

**Dr. Dominik Herrmann**

Slides available at

<https://dhgo.to/dns-traffic-analysis>





YAHOO!

Search Web

Sign in

Mail

Mail

News

Sports

Finance

Autos

Celebrity

Shopping

Movies

Politics

Beauty

More

Chris Rock lives up to hype as Oscars host

Tracy Morgan's surprise cameo was a highlight - but Stacey Dash's ultra-awkward appearance confused viewers. **Who was snubbed »**

Best and worst dressed at the Oscars

'Spotlight' wins Best Picture in Oscars upset

Rubio has momentum but still trails in polls

Scientific reason behind leap year

Trending Now

1. The Spotlight movie
2. Boston Celtics
3. Abraham Attah
4. Jared Leto
5. Evening dresses
6. Migraine headaches
7. Kim Kardashian
8. Jacob Tremblay
9. Toyota Corolla
10. Gigi Hadid

Never Old. Never New. Get yours >

foreverspin™

MADE IN CANADA

## Motivation of monitoring DNS

- block known malicious domains (e.g. phishing)
- retain log of all DNS queries for later analysis

OpenDNS

## Why is DNS monitoring interesting for forensics?

analyzing hard disk not sufficient any more  
(cloud, private browsing, disk encryption)

tail

ead><title>Yahoo</title>...

w.yahoo.com/celebrity/

*What can we infer from DNS query logs?*

oo.com

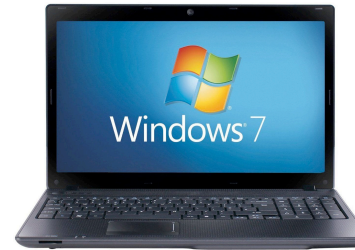
**low storage needs**

## DNS log contains essential metadata:

2016-03-05	11:14:05.124	2.240.3.12	www.yahoo.com	A
↓	↓	↓	↓	↓
date and time		user's address	domain	type

## Example 1: confirm source of traffic

Did incriminating traffic originate from **Bob's** laptop?



Source of discrepancy?  
Rogue hardware?

DNS queries from Bob's IP



2016-03-05 06:46:01.383 aus5.mozilla.org

2016-03-05 09:41:37.263 c1.adform.net

**2016-03-05 09:46:01.455 www.exploit-db.com**

2016-03-05 09:41:37.708 l.betrad.com

?



2016-03-05 10:22:01.814 time.apple.com

2016-03-05 09:41:38.262 lotame.nexac.com



2016-03-05 10:22:01.950 b.config.skype.com

2016-03-05 09:41:38.686 ping.chartbeat.net



2016-03-05 14:17:09.663 notify5.dropbox.com

2016-03-05 09:41:41.627 lib\_dns\_sd\_udp.fritz.box



2016-03-05 14:17:10.411 ols.officeapps.live.com

2016-03-05 09:41:48.917 college.usatoday.com

?



2016-03-05 15:29:22.510 api.textmate.org

## Example 2: reconstruct visited websites

- What websites did *Eve* visit before we fired her?
- Which users surfed to *www.yahoo.com* last week?



## Searching for www.yahoo.com ...

09:41:20.242	ad4.adition.com	
09:41:21.770	ads.nuggad.com	
09:41:40.152	skypedata.akadns.net	
09:42:41.985	dl-debug.dropbox.com	
<b>09:45:11.201</b>	<b>google.com</b>	<b>visited</b>
<b>09:46:00.033</b>	<b>www.heise.de</b>	<b>visited</b>
<b>09:46:00.133</b>	<b>dealbook.nytimes.com</b>	
<b>09:46:00.134</b>	<b>pressroom.yahoo.net</b>	<b>DNS prefetching</b>
<b>09:46:00.169</b>	<b>www.yahoo.com</b>	
<b>09:46:00.783</b>	<b>imagesrv.adition.com</b>	
<b>09:46:00.989</b>	<b>ad.atdmt.com</b>	<b>advertisements &amp;</b>
<b>09:46:00.989</b>	<b>ad.doubleclick.net</b>	<b>user tracking</b>
<b>09:46:00.991</b>	<b>imagerv2.adition.com</b>	
<b>09:46:01.017</b>	<b>jobs.heise.de</b>	<b>embedded image</b>

Simple heuristics look promising ...  
... but are not always accurate.

Heuristic search:  
 $\Delta t > 5 \text{ sec}$

09:41:20.242	ad4.adition.com	
09:41:21.770	ads.nuggad.com	
09:41:40.152	skypedata.akadns.net	
09:42:41.985	dl-debug.dropbox.com	
<b>09:45:11.201</b>	<b>google.com</b>	<b>true positive</b>
<b>09:46:00.033</b>	<b>www.heise.de</b>	<b>true positive</b>
09:46:00.133	dealbook.nytimes.com	
09:46:00.134	pressroom.yahoo.net	
09:46:00.169	www.yahoo.com	<b>true negative</b>
<b>09:46:30.812</b>	<b>[visit Yahoo website]</b>	<b>false negative</b>

**www.yahoo.com**  
cached for 1–5 min



Browser



DNS Resolver



## 51 domains resolved when Yahoo's home page is visited

**www.yahoo.com**

bs.serving-sys.com

pclick.yahoo.com

s.yimg.com

sb.scorecardresearch...

crl-ds.ws.symantec.co...

y.analytics.yahoo.com

geo.query.yahoo.com

csc.beap.bc.yahoo.com

geo.yahoo.com

comet.yahoo.com

answers.yahoo.com

everything.yahoo.com

groups.yahoo.com

login.yahoo.com

mail.yahoo.com

mobile.yahoo.com

shopping.yahoo.com

**www.flickr.com**

**www.tumblr.com**

beap.gemini.yahoo.com

finance.yahoo.com

ftw.usatoday.com

geo-um.btrll.com

googleads.g.doublecli...

match.adsrvr.org

pagead2.google syndic...

help.yahoo.com

info.yahoo.com

news.yahoo.com

na.ads.yahoo.com

pr-bh.ybp.yahoo.com

r.turn.com

rmx.pxl.ace.advertisin...

search.yahoo.com

sports.yahoo.com

**thinkprogress.org**

sync.adap.tv

sync.adaptv.advertisin...

**www.cbsnews.com**

ads.yahoo.com

**www.chicagotribune....**

**www.foxnews.com**

**www.latimes.com**

fonts.googleapis.com

tpc.google syndication...

cm.g.doubleclick.net

**www.npr.org**

**www.politico.com**

**www.sbnation.com**

**www.upi.com**

*Can we use the **set of domains** to verify  
whether a website was visited?*

**Experimental approach:**

1. Download websites with a browser
2. Record resolved hostnames
3. Determine  $k$ -identifiability of websites

**Measurements indicate:**

many websites have a unique DNS pattern

visited  
home page



inference of  
**whole (!) URL**



**Interesting problems:**

- robustness
- threshold for match
- influence of cache

$k = 1$       99 %

63 %

$k \leq 5$       99 %

76 %



Browser



DNS Resolver



**DNS log might not be available**  
(due to data protection obligations)



**www.yahoo.com**  
bs.serving-sys.com  
pclick.yahoo.com  
s.yimg.com  
...

only packet sizes are logged  
(no domain names)

however: DNS packet size correlated  
with domain name length

**logging of flow records**

(common practice)

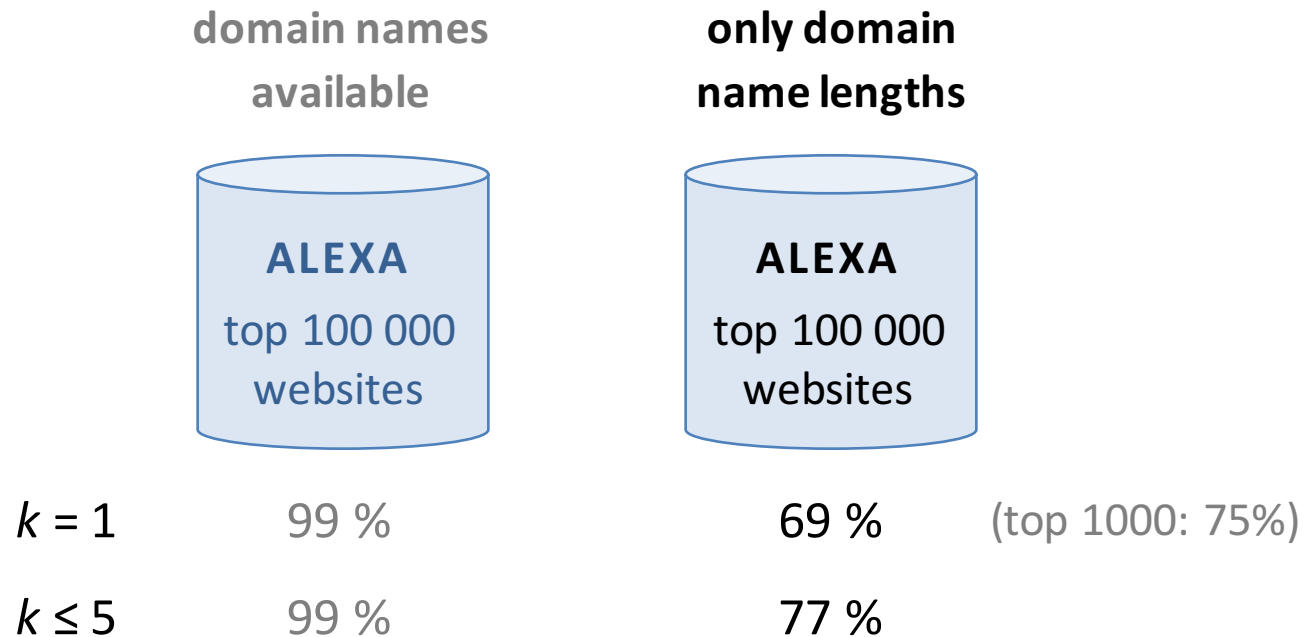
**13 18 16 10 24 34 21 19 21 13**  
**15 17 20 16 15 14 16 18 14 14**  
**21 17 16 16 27 16 29 14 14 14**  
**16 19 10 27 16 16 17 12 27 15**  
**13 22 15 15 20 25 20 11 16 16**  
**11**

*Is DNS-based visited website  
verification still possible?*

**Yahoo's DNS flow record fingerprint**  
(multiset of 51 domain name lengths)

## Measurements indicate:

domain lengths multiset is characteristic



## drawing inferences from DNS logs and flow records

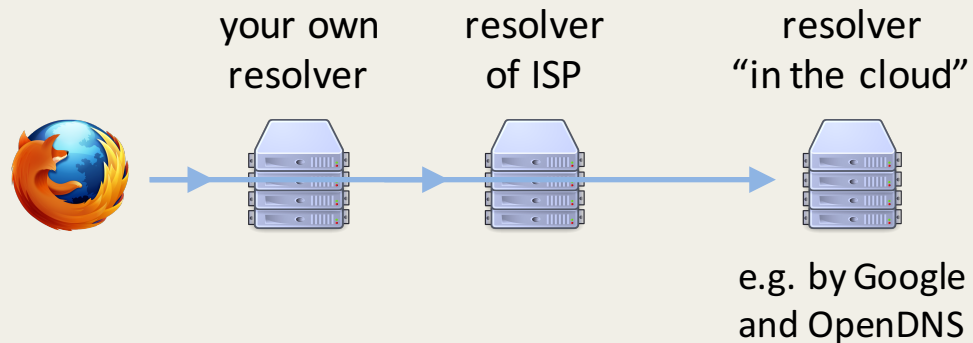
useful for forensics

*real-world  
accuracy?*

*utility for law  
enforcement?*

*probative value  
of evidence?*

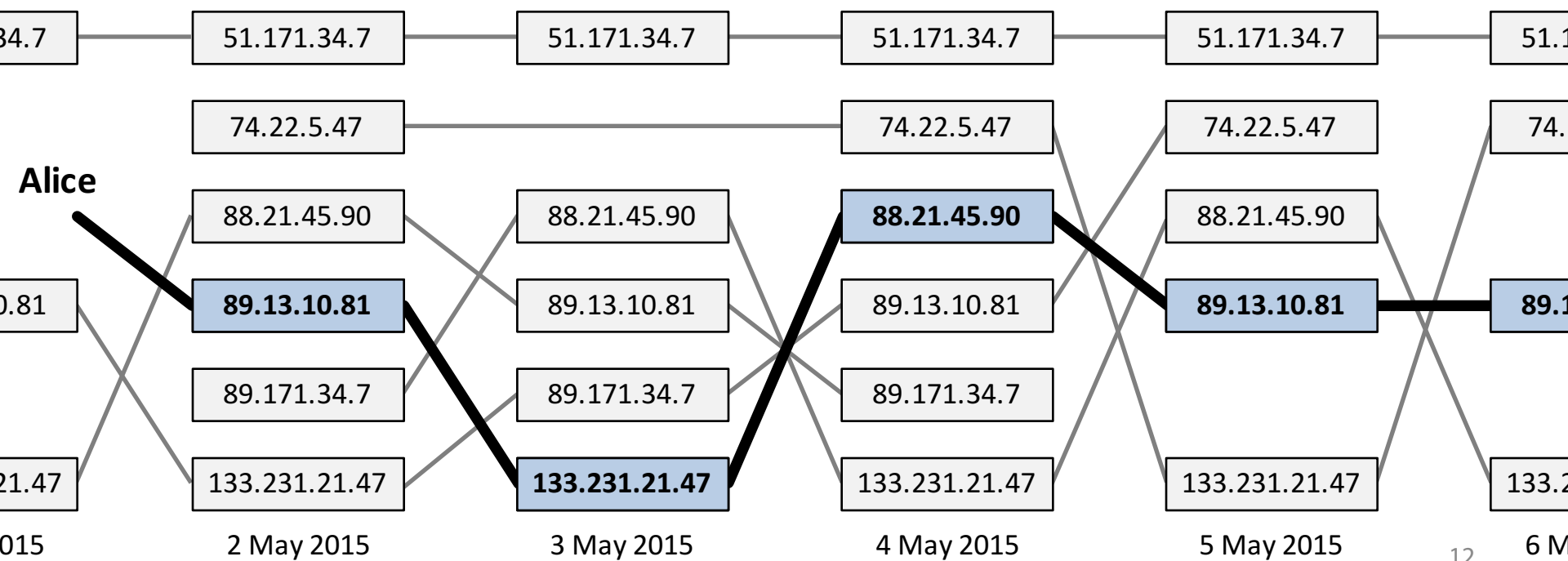
privacy concerns



*But third-party DNS resolvers  
cannot track their users – **or can they?***

### Challenge:

IP address changes frequently (daily)



**3 May 2015**

spiegel.de 4 x  
google.de 15 x  
apple.com 1 x  
**airbus.com 3 x**  
**mpg.de 2 x**

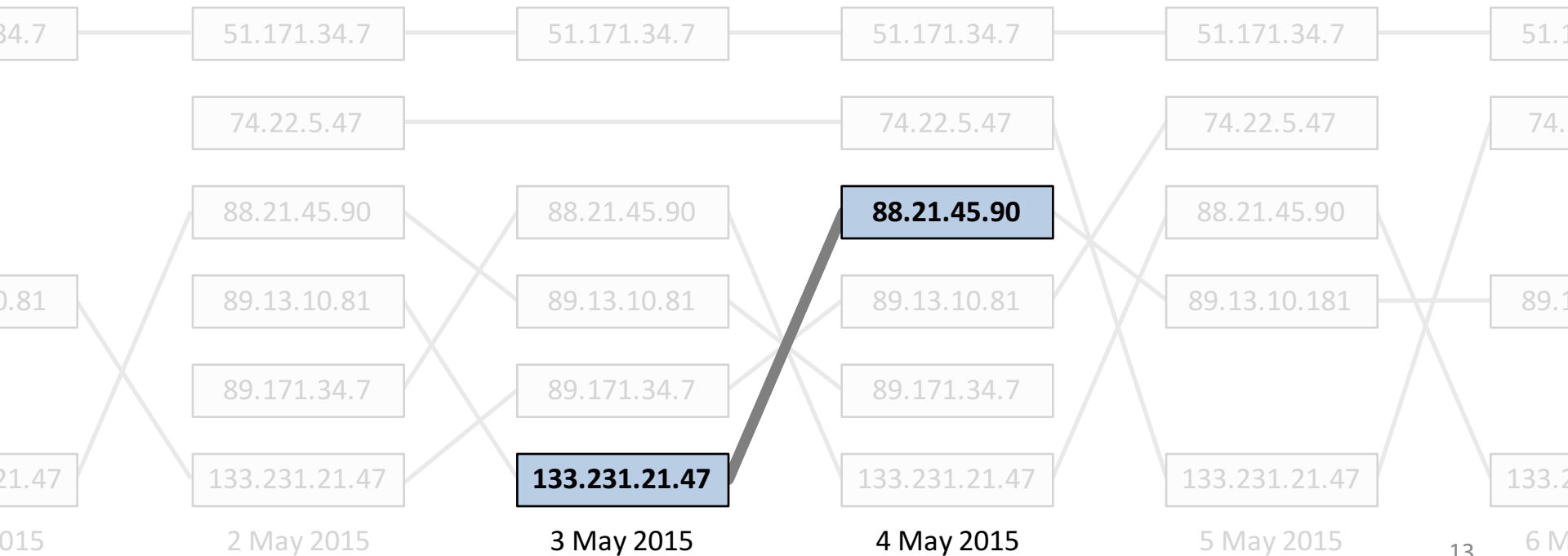


**re-identification via  
resolved domains**

**4 May 2015**

1 x spiegel.de  
9 x google.de  
0 x apple.com  
**6 x airbus.com**  
**3 x mpg.de**

*Do users have  
distinct habits?*

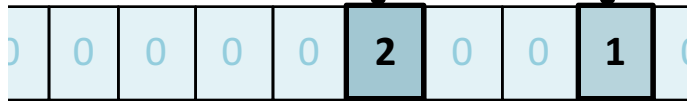


**Sessions** are modelled as **vectors** that are compared with **cosine similarity**  
(nearest-neighbor classifier)

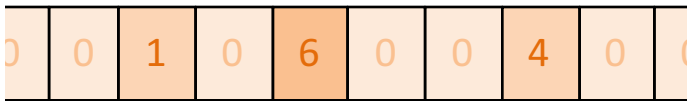
airbus.com

bahn.de

1



2



3

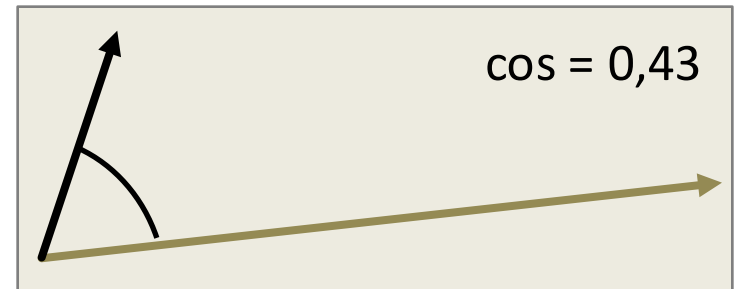
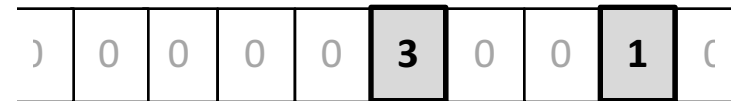
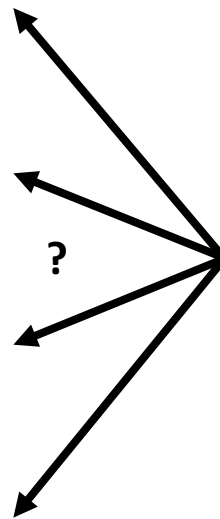


4



⋮

yesterday



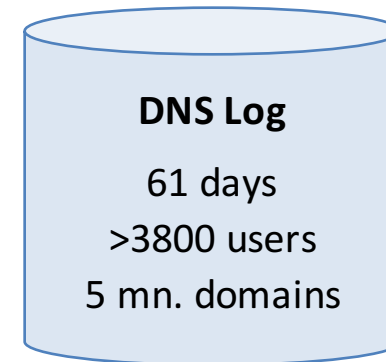
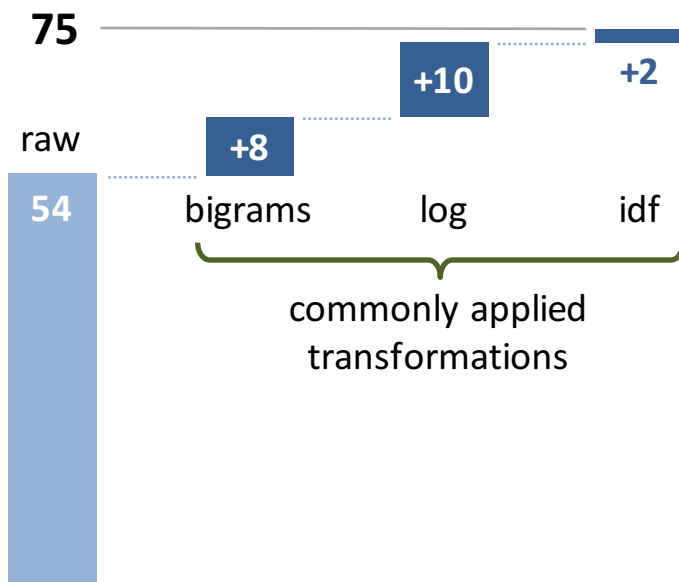
today

*How accurate is behavior-based tracking in practice?*

### Experimental approach:

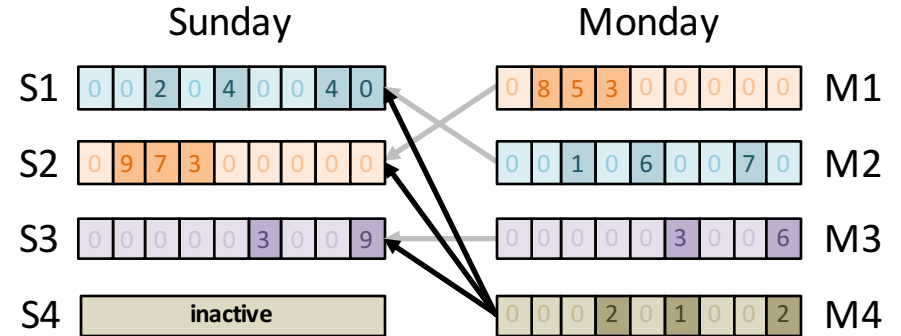
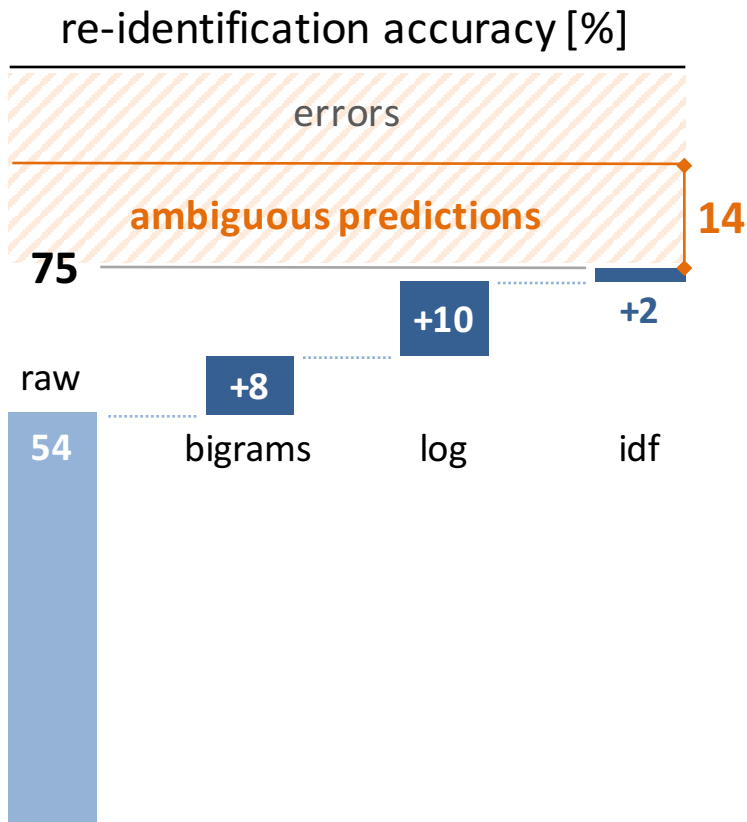
1. Obtain DNS log with realistic traffic
2. Track users day to day (24h sessions)
3. Determine overall accuracy

re-identification accuracy [%]

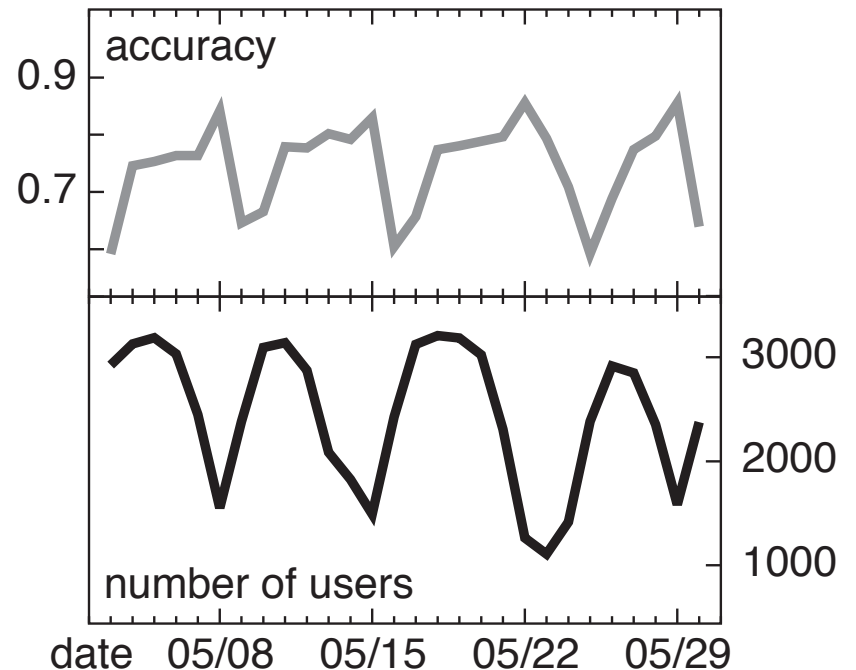


with »ground truth«  
(pseudonymized)

*How accurate is behavior-based tracking in practice?*

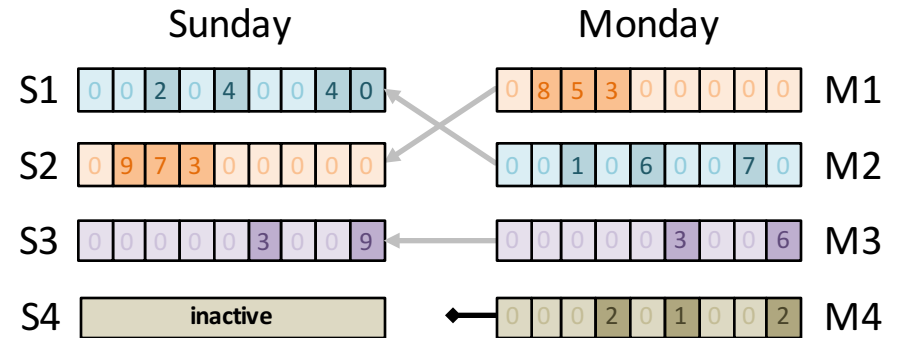


**ambiguous prediction  
... can be resolved**

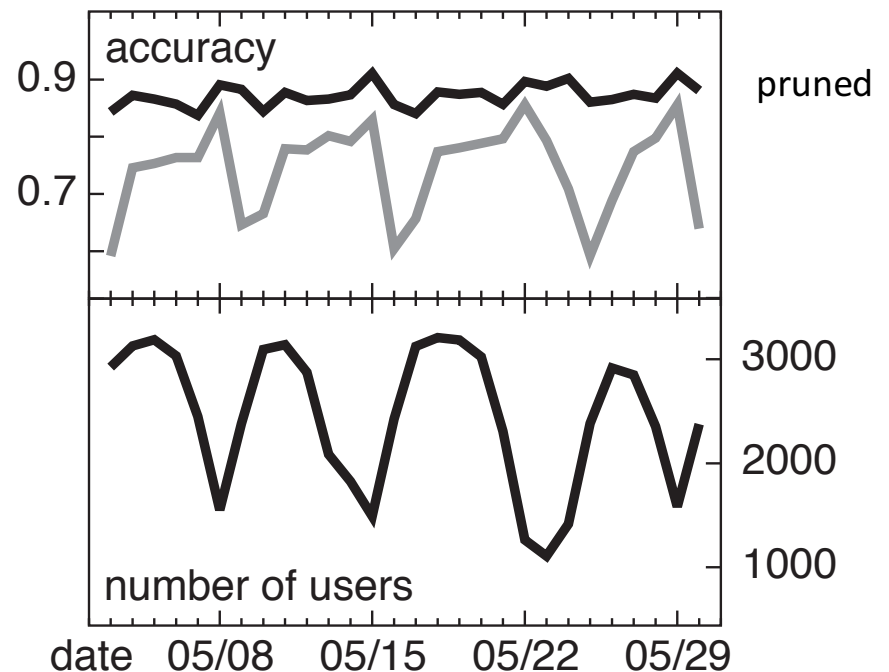
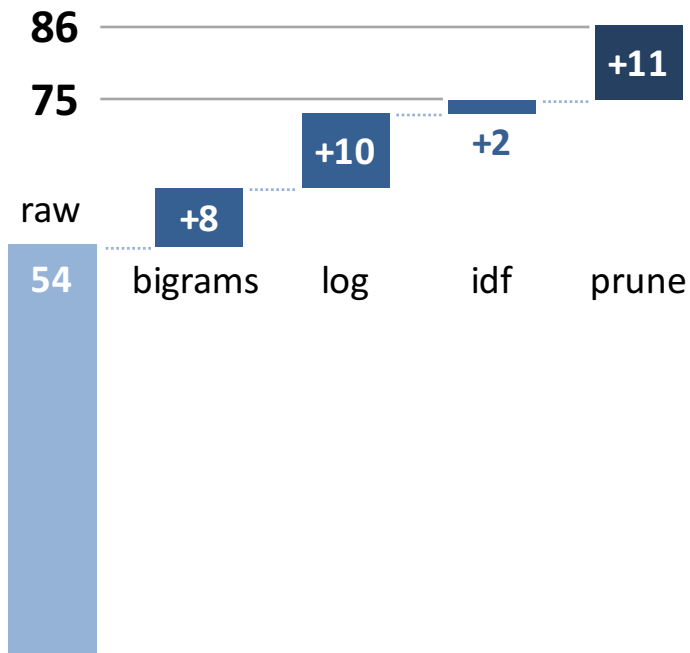




*How accurate is behavior-based tracking in practice?*



re-identification accuracy [%]

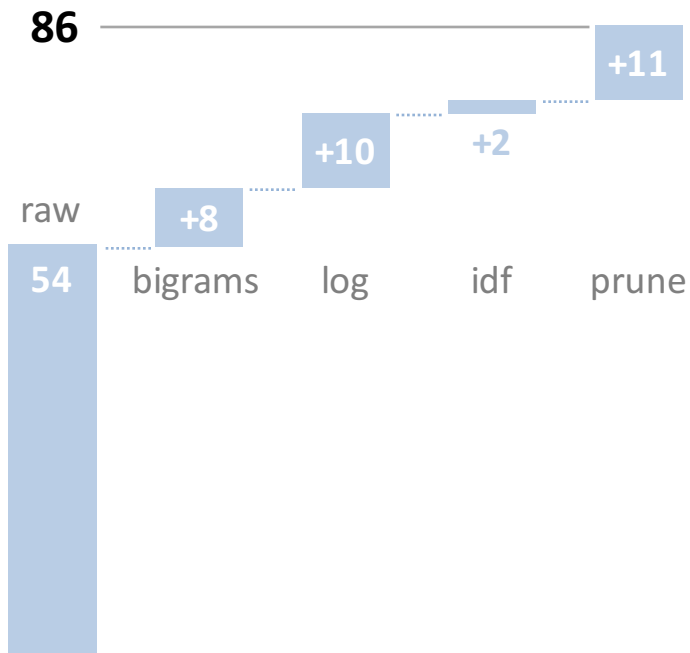


## Application to network forensics:

*How accurate is user re-identification with **flow records only**?*

### domain names

re-identification accuracy [%]



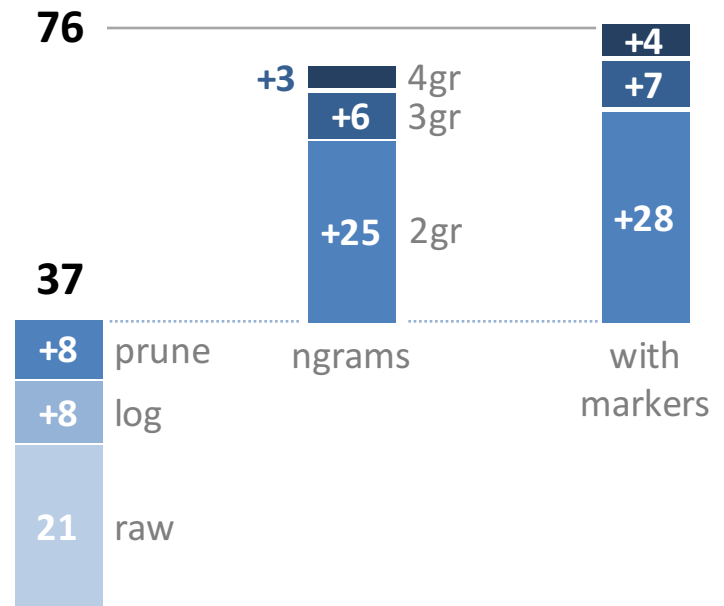
## Idea of ngram markers:

observed: 15 30 [ pause  $\geq 5$  s ] 18

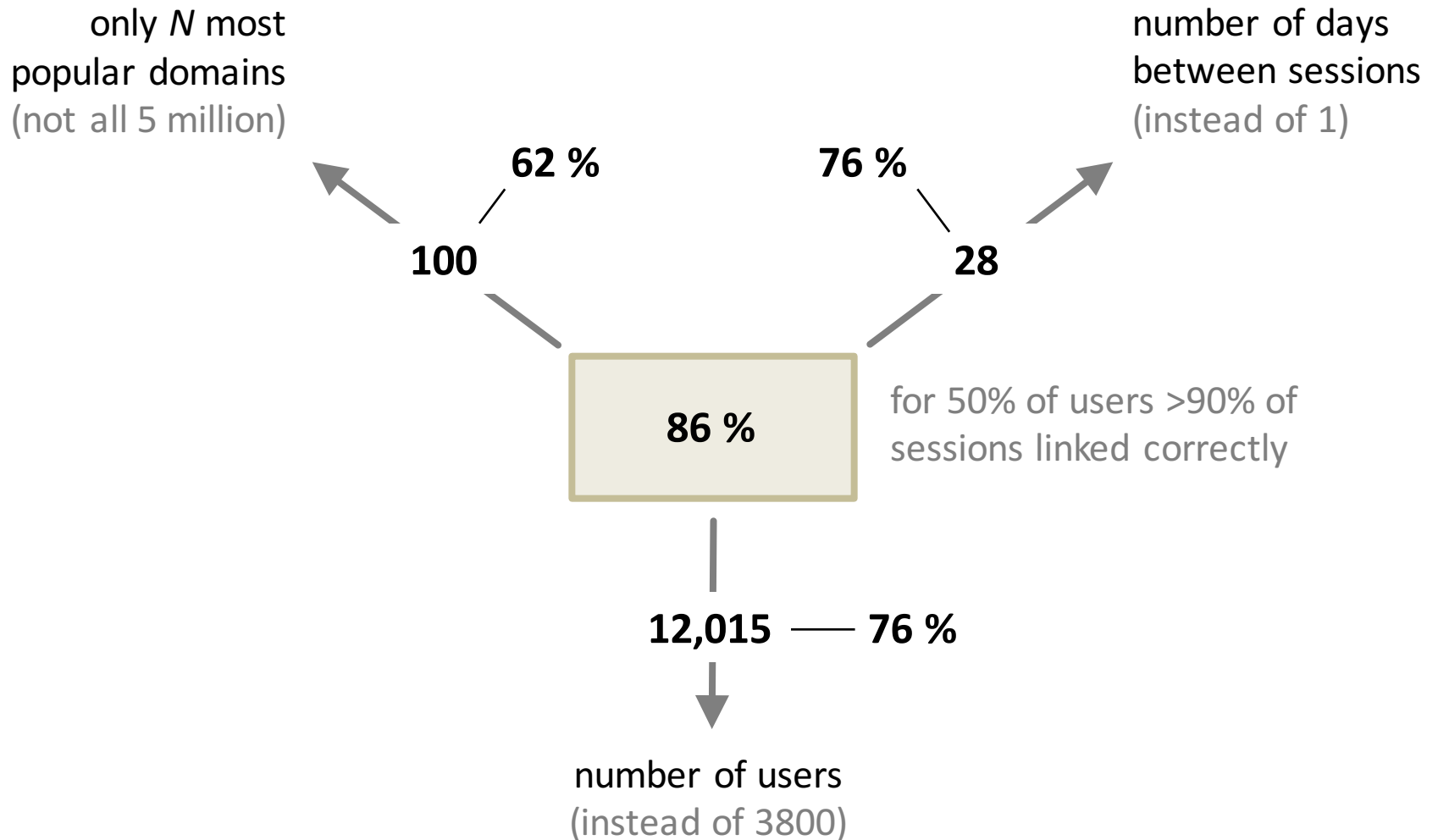
bigrams: 15-30 30-P P-18

### domain name lengths

re-identification accuracy [%]



## Behavior-based re-identification is quite robust.



## behavior-based linkage of browsing sessions

significant because undetectable  
threatens informational self-determination

*accuracy improvements?*

**yes**  
work in progress

*exploitable  
by ad-networks?*

*other applications?*

forensics  
authentication  
anomaly detection

*affordable protection?*

**yes**  
stay tuned

*What should a privacy-preserving  
DNS resolver look like?*

generic anonymization  
services (Tor) too slow

## Tailored solution: EncDNS

repurpose resolver of ISP as a proxy for encrypted queries

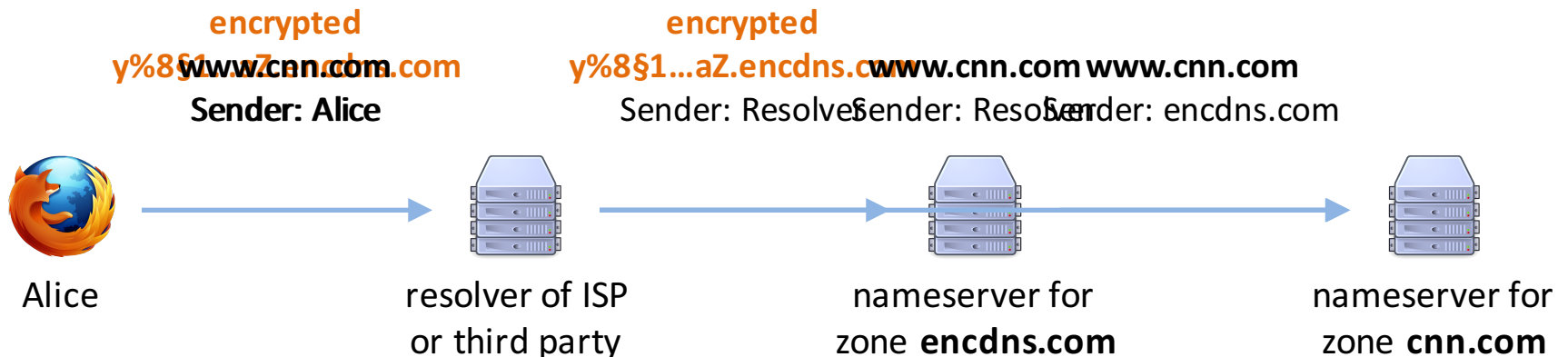
### Challenge:

limited space (255 bytes)

cryptobox of Bernstein's NaCl library  
(Curve25519)

### Measurements indicate:

fast and scalable (>6000 queries/sec)



We can exploit **peculiarities of DNS** to improve performance and privacy.

**Observation 1:**

few domains are very popular (power law)  
top 10,000 domains: 80% of all queries

**Observation 2:**

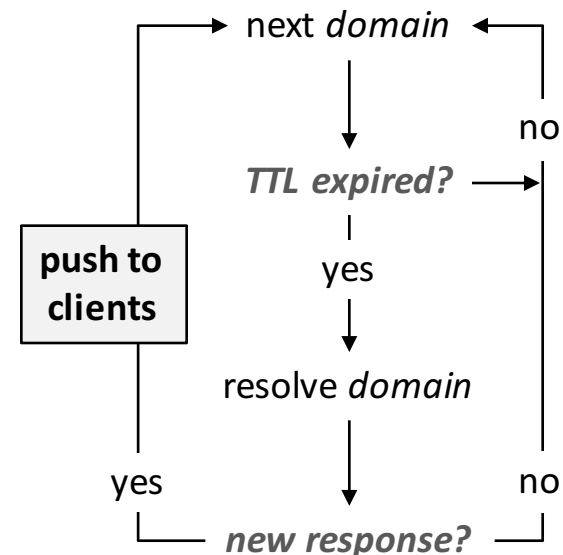
most IPs constant over long time  
for 50% of domains: TTL > 5 min

**Tailored solution: PushDNS Service**

send DNS records of most popular domains to connected clients

**Traffic requirements (10,000 domains):**

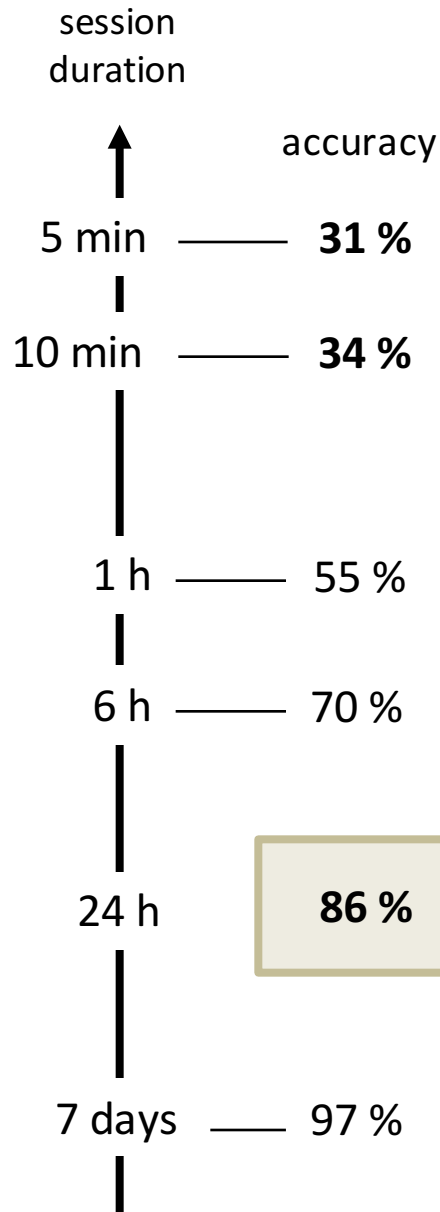
- resolving domains: 350 MB per day
- pushing updates: 0.8 KB/s per user



**Consequence:** majority of queries **unobservable** and resolved **instantaneously**

## Protection against behavior-based tracking

... can be delegated to Internet Service Provider



**Change IP address frequently!**

**Chance for ISPs**

Effortless protection with  
**IPv6 Prefix Bouquets**

**ANON-Next**

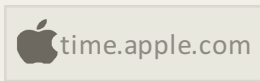
(BMBF, 2016 – 2019)

**manitu**

**opportunity  
for forensics**

# A Double-Edged Sword: Traffic Analysis in the Domain Name System

**threat to  
privacy**



DNS patterns of software and websites

13 18 16 10 24 34  
15 17 20 16 15 14

behavior-based tracking of users

0 2 0 1 0 0 2

## INFERENCE IN NETWORKED SYSTEMS

## PRIVACY ENHANCING TECHNOLOGIES

EncDNS

tailored protection tools promising

PushDNS

effortless tracking protection by delegation

IPv6 Prefix Bouquets