# PRIVACY OF THE DOMAIN NAME SYSTEM

Isn't HTTP**S** enough?

- Since Snowden: privacy is of utmost importance
  - \>90% of Web traffic is **HTTPS**
- **Every (website) visit is preceded with a bunch of DNS queries**
- DNS in a nutshell
  - Phonebook of the Internet
  - Translate hostnames to IP address
  - (Used to be) plain-text
    - *"I might not see the content you consume, but I CAN see where it comes from"*
- Main three reasons of being plain-text
  1. Historically, less focus on privacy and security
  2. DNS is an overhead → simplest → fastest
  3. Services heavily rely on DNS data

1) I want to visit example.com

2) where is example.com??

3) Go to: 93.184.216.34  DNS

4) Visit example.com at 93.184.216.34

# DNS IS A DOUBLE-EDGE SWORD
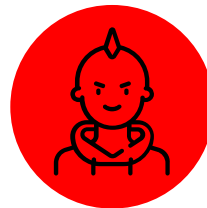
## Services based on plain-text DNS

**Your Internet Service Provider (ISP)**
- Firewall
- Parental-control
- Pay-as-you-go-models (e.g., at hotels)
- Content caching / Proxy
- Broadband router configuration
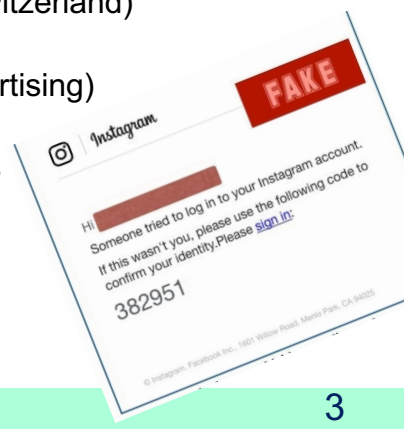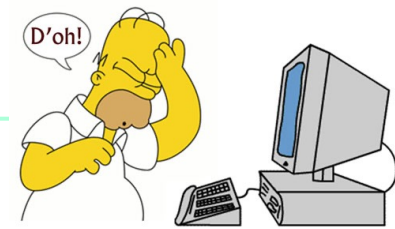- Blocking Ads
- Law-enforcements

**Malicious actors and authoritarian regimes**
- Blocking websites and contents
  - Political speech (e.g., in China)
  - Foreign gambling (e.g., in Switzerland)
- Spy on the users
  - Monetize DNS data (for advertising)
- Tampering and redirection
  - Malicious (phishing) websites

This site can't be reached

dhfhd.com's server IP address could not be found.

Try running Windows Network Diagnostics.

DNS_PROBE_FINISHED_NXDOMAIN

Instagram  FAKE

Hi

Someone tried to log in to your Instagram account.
If this wasn't you, please use the following code to confirm your identity. Please sign in:

382951

© Instagram, Facebook Inc., 1601 Willow Road, Menlo Park, CA 94025

# DoH! Headache for ISPs...
## ...privacy heaven for users and malicious actors?

**Do53 (plain-text DNS)**

spy

block

tamper

**DNS-over-TLS (RFC 7858)**

Block port 853

Fall back (RFC 8310)

ISP's DNS resolver

**ISP**

Remote Do53 Recursive Resolver (RR)

Local DNS query
Remote Do53 query
Remote DoT query
Remote DoH query

Law enforcement
Local content
Parental Control
Security

**ISP network**

**INTERNET**

Remote DoH Recursive Resolver (RR)

**DNS-over-HTTPS (RFC 8484)**

www.

**DNS-over-HTTPS: circumventing all ISP measures**
- Inherently blends into regular encrypted HTTPS traffic
  - Cannot be filtered, cannot be differentiated, cannot be blocked

- All ISP services break
  - No firewall, no parental control, no cache, no malware detection, etc.

- **Enhanced Privacy vs. Weakened Protection**

# IN THIS PAPER
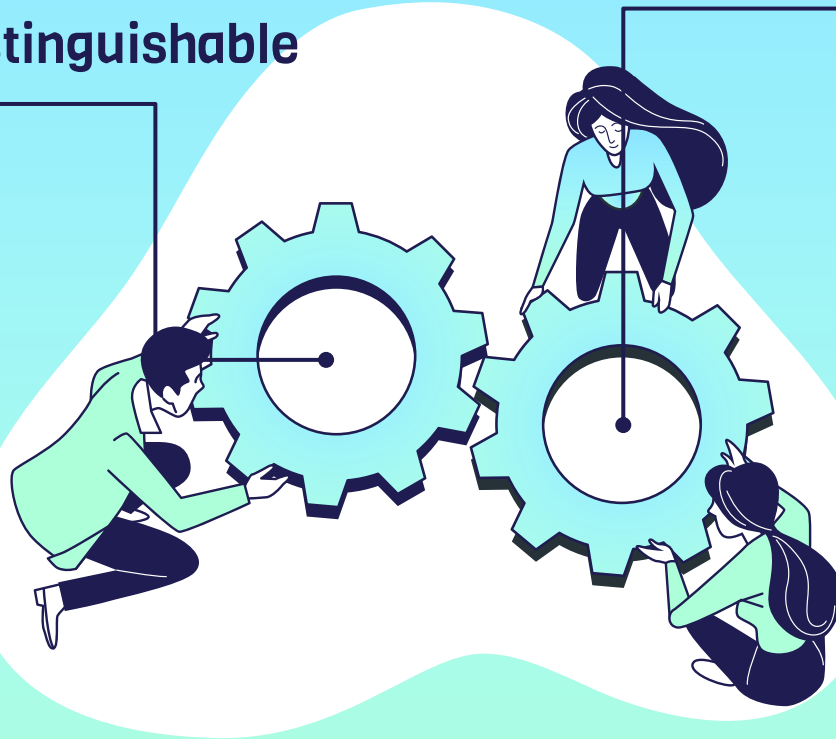
## Is DoH indeed indistinguishable from Web traffic?

**Countermeasures?**

### NO

### YES

We build a **Machine Learning model to identify** each encrypted packet (on port 443) as **Doh** or **Web**

We study a wide set of **padding techniques**

To disguise the DoH identification model trained on the padded data

| 97.4 % Closed-world | 90 % Open-world |

At extremely low false-positive rate (FPR=$10^{-4}$)

ISPs can identify and block it → transparently fall back to Do53

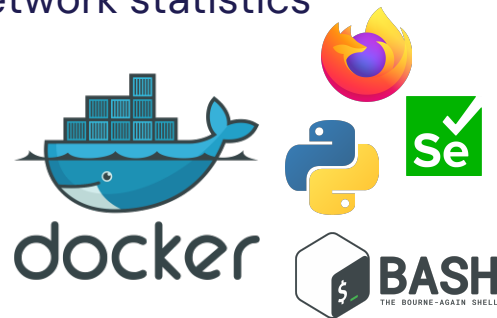The anti–identification model **significantly reduces** the classification accuracy

| 53 % Closed-world | 0 % Open-world |

5

# DATASET COLLECTION

Multiple cities across multiple continents to capture diverse network statistics

- Easy to deploy Docker containers
- Alexa top-1M domains
  - Visit the first 20K one-by-one (to flush DNS cache)
- Using 25 DoH resolvers[1]
  - Well-known: Google, Cloudflare, Quad9, CleanBrowsing
- Containers deployed "world-wide"
  - *LocA* – South America: University of Campinas, Brazil (`x86`)
  - *LocB* – North America: Multiple Cloudlab sites (`x86, arm64`)
  - *LocC* – Asia: National University of Singapore (`x86`)

[1] https://github.com/curl/curl/wiki/DNS-over-HTTPS

## TRAIN A SUPERVISED MACHINE LEARNING MODEL

**01**

- Chosen ML model: Random forest
  - out of six models evaluated
- Train–test ratio: 90–10%
  - Found similar results with 80–20%, 70–30%

## IMPORTANT METRICS

**03**

- Precision, Recall and $F_1$–score
- False-positive rate (FPR)
- Recall at low FPR
  - Not deployable if Web packets are blocked due to misclassification

## FEATURES

**02**

- IP length of current packet
- IP length of the previous packet
- Inter-packet arrival, i.e., time lag, of current packet
- Time lag of previous packet

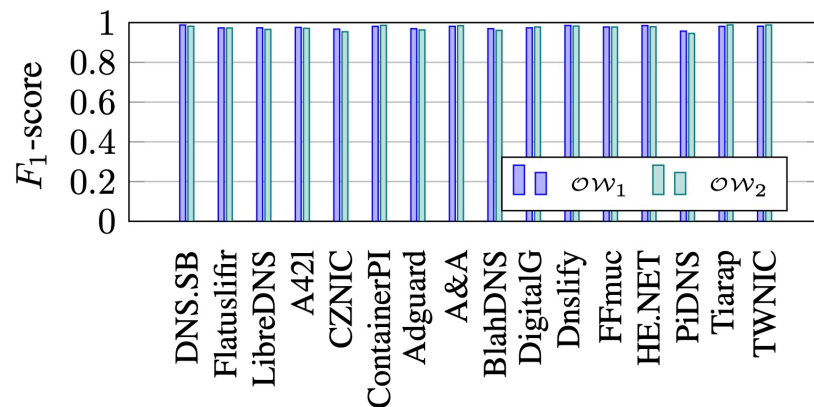## SETTINGS FOR EVALUATION

**04**

- Closed–world
  - Same resolvers, same domains visited
- Open–world 1
  - Different resolvers, same domains visited
- Open–world 2
  - Different resolvers, different domains visited

Closed- and Open-world results for the data gathered in $LocC$ (Asia)

- Best-case:
  - Resolvers used for training
    - Prominent: `Google, Cloudflare, Quad9, CleanBrowsing`
    - + worst-performing: `Comcast, OpenDNS, Doh.li`
  - Closed-world: $F_1$-score **>0.99** (FPR=**0.009**)
  - Open-world: $F_1$-score = **~0.975** (FPR=**0.0055**)
- Best-case at **low FPR < 0.0001**:
  - Closed-world   Recall = **0.974**
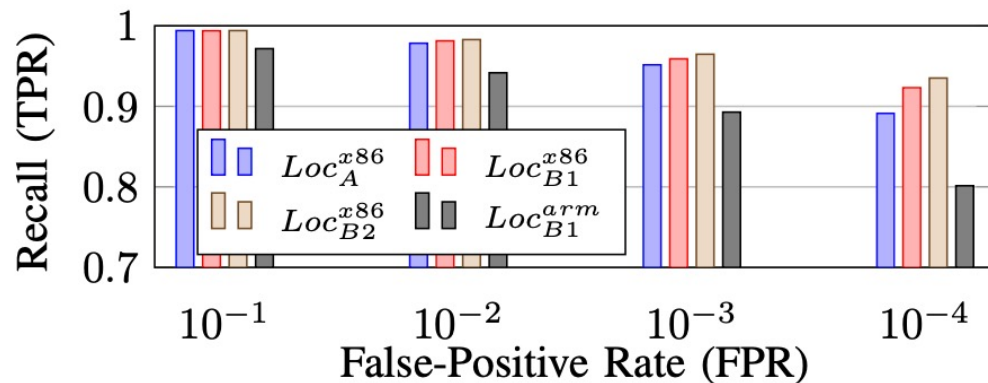  - Open-world:    Recall = **0.9**



FPR=$10^{-4}$ → 1 out of 10,000 Web packet is misclassified as DoH

Robustness of the DoH identification model

- Worst-case:
  - Model trained in one location and tested at other locations
    - Trained at `LocC` (`x86`), tested at `LocA` and `LocB` (`arm` and `x86`)
  - Closed-world:
    - `x86`: Recall = **~0.90** (FPR=**0.0001**)
    - `arm`: Recall = **~0.80** (FPR=**0.0001**)



FPR=$10^{-4}$ → 1 out of 10,000 Web packet is misclassified as DoH

# COUNTERMEASURES?

ISPs deploying the DoH identification model can filter out DoH packets with high accuracy and low FPR
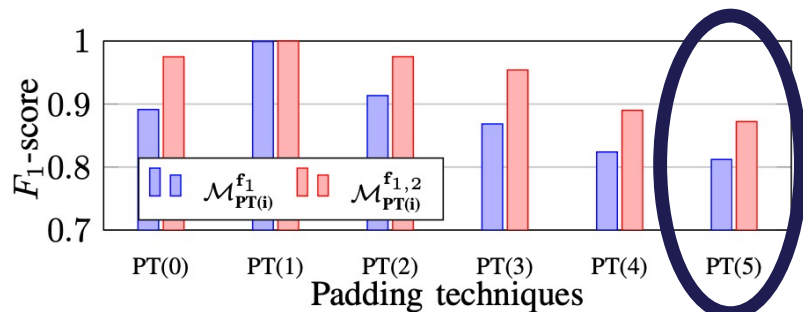
- **Two fundamental packet characteristics** to manipulate:
  - Packet length and time lag
- **Idea**: Pad the DNS packets to look more like Web packets
    1) Fix padding (RFC8467) – closest multiple of 128B
    2) Random padding
    3) Pad to the average of the Web packets
    4) Pad to a random recent Web packet
    5) Pad a sequence of DoH packets to a recent sequence of Web packets
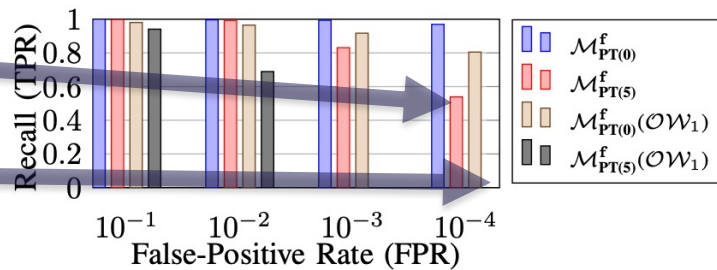
# EVALUATION OF PADDING TECHNIQUES

Proposal for a DoH anti-identification model

- Preliminary analysis
  - Padding **packet lengths only**
  - **Closed-world** setting
    - Feature $f_1$ = packet length
    - Feature $f_2$ = previous packet length
    - PT($i$) = different padding techniques
    - PT(0) = original non-padded data
- Apply PT(5) on all features as well
  - at low FPR=**0.0001**

53 %
Closed-world

0 %
Open-world

# SUMMARY

- Privacy of the DNS is important

- DoH is designed to blend DNS traffic into HTTPS Web traffic
  - Cannot be monitored, cannot be filtered, cannot be blocked

- Our main contributions
  - DoH identification model to distinguish DoH and Web packets with high accuracy
    - **97.4%** and **90%** in the **closed–** and **open–world** setting, respectively
    - With a **false–positive rate of 0.0001**
  - Develop DoH anti–identification model as a counter–measure
    - **53%** and **0%** in the **closed–** and **open–world** setting, respectively
    - With a **false–positive rate of 0.0001**

# Q&A

## Levente Csikor

Trustwave

`levente.csikor@gmail.com`

Container for data collection:
- `https://github.com/cslev/doh_docker`
- `https://hub.docker.com/r/cslev/doh_docker`

Machine Learning algorithm:
- `https://github.com/cslev/doh_ml`