# Towards Web Services Classification

## Problem

Users

73.45.5.15

Encrypted

Passive Monitoring

Where is this flow directed?

liverail.com

instagram.com

adnxs.com

facebook.com

akamaihd.net

73.45.5.15

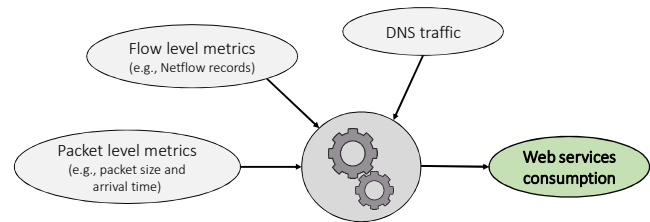A flow to 73.45.5.15 is to Facebook, Instagram or … ?
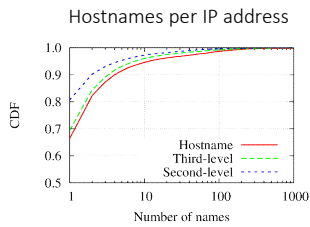
- ❑ Traffic classification difficult in modern web, traffic more and more **encrypted, CDNs and Cloud Providers** complicate the scenario
- ❑ A classifier can still rely on **flow level metrics** ( flow records)
- ❑ **Server IP address** doesn't give much information
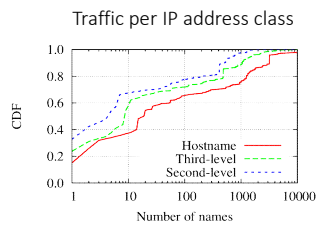
## Goal

- ❑ **Account** traffic to web service generating it
- ❑ Evaluate which **features** of traffic are useful for classification
- ❑ Leverage **machine learning** techniques
- ❑ **Self-learning** approach when possible

Flow level metrics (e.g., Netflow records)

DNS traffic

Packet level metrics (e.g., packet size and arrival time)

Web services consumption

## How many names have IP addresses?

Hostnames per IP address

Traffic per IP address class



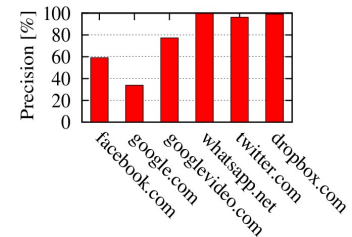CDF — Hostname, Third-level, Second-level; Number of names

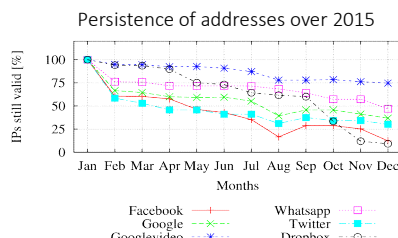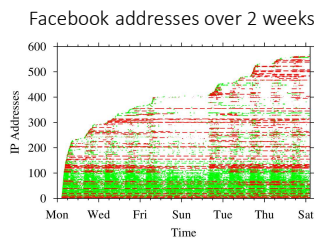The majority of addresses associated to 1 hostname.

Those addresses carry a little share of the traffic (20%).

## Bags of IP addresses

- ❑ Enumerate all the IP addresses of some popular services.
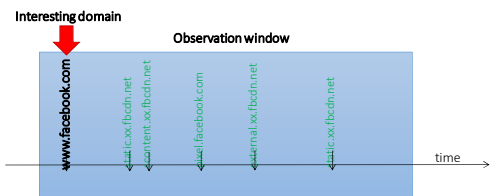- ❑ Consider all traffic going to those addresses as belonging to the respective service.



Precision [%]: facebook.com, google.com, googlevideo.com, whatsapp.net, twitter.com, dropbox.com

## How stable are IP addresses in the time?

Facebook addresses over 2 weeks

Persistence of addresses over 2015



IP Addresses; Time — Mon Wed Fri Sun Tue Thu Sat

IPs still valid [%]; Months — Jan Feb Mar Apr May Jun Jul Aug Sep Oct Nov Dec

Facebook, Whatsapp, Google, Twitter, Googlevideo, Dropbox

## Temporal correlation

Characterize interesting domains neighborhood

Interesting domain

Observation window

www.facebook.com

time

## Future work

- ❑ A classifier relying only on server IP address doesn't achieve high performance ( about 20-30% of coverage)
- ❑ Other traffic characteristic can be exploited for classification:
  - ❖ DNS requests and replies
  - ❖ Packet level features (e.g., packets size, arrival time…)
  - ❖ Temporal and spatial correlation among flows
- ❑ Combine such approaches in a unique system

Traffic share



- Google 16%
- Youtube 27%
- Facebook 13%
- Storage 7%
- News 8%
- Other 29%

**Martino Trevisan**
Marco Mellia

Politecnico di Torino

martino.trevisan@polito.it