# Poster: Data Collection for ML Classification of Encrypted Messaging Applications

Jason Hussey, *Colorado School of Mines (presenting)*

Ethan Taylor, *Colorado School of Mines*

Kerri Stone, *iCR, inc.*

Tracy Camp , *Colorado School of Mines*

CS@Mines

**WHATSAPP RIVAL SIGNAL GETS 'MILLIONS' OF NEW USERS IN THE WAKE OF FACEBOOK'S DRAMATIC SIX-HOUR OUTAGE**

https://www.independent.co.uk/life-style/gadgets-and-tech/facebook-outage-instagram-whatsapp-signal-down-b1932505.html

"Signal is regularly used by journalists and investigators to protect sources identity"

<u>Users in 2020:</u>
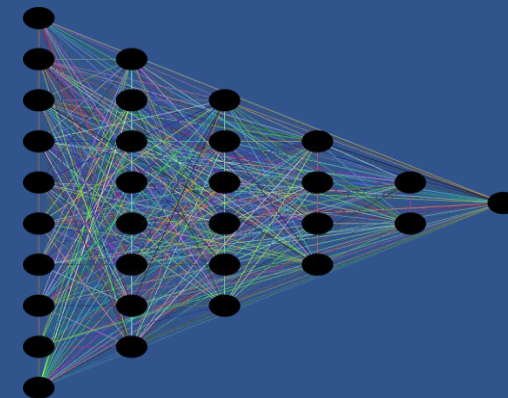WhatsApp, 2 billion
Telegram, 400 million
Signal, 20 million

https://www.businessofapps.com/data/signal-statistics/
https://www.businessofapps.com/data/telegram-statistics/
https://www.businessofapps.com/data/whatsapp-statistics/

## Signal: The Pros and Cons of a Truly Private Chat App

Signal, the encrypted messaging app, is seeing record numbers of downloads amid the pandemic and nationwide protests. It might make sense for you, too.

https://www.wsj.com/articles/signal-the-pros-and-cons-of-a-truly-private-chat-app-11592127002

Poster: Data Collection for ML Classification of Encrypted Messaging Applications

# Research Summary

- **Network traffic classification** is used to identify the nature of traffic on a network.

- Entities capable of monitoring network traffic use classification for all manner of reasons, including **identification of mobile applications being used on the network**.

- It is possible that the usage of encrypted messaging applications by users on these networks can be detected, **betraying elements of their privacy.**

- We describe a system that:
  - leverages campus network resources to generate real-world data
  - alongside a more curated dataset captured from Android application traffic.

- We also explore the ability of machine learning (ML) models to accurately classify traffic from these encrypted messaging applications.

# Methodology – Data Collection

## WiFi Data Collection

- Partner with the ITS office to collect anonymous WiFi packet headers
- Leverage ntop's n2disk utility
  - Zero copy drivers
- Extract just the IP and TCP/UDP headers and pre-process with tshark
- Multiprocess the tshark output into mongodb

## Android Application Collection

- Rooted Android phones (Samsung and Xiaomi)
- X-compiled strace attached to Signal messaging app process
- netstat polling for verification
- tcpdump on a Ubuntu station serving as AP
- Filter the PCAP file to only those flows identified by socket calls in trace

```
23903 getsockopt(84, SOL_SOCKET, SO_DOMAIN, [10], [4]) = 0
23903 socket(AF_UNIX, SOCK_STREAM|SOCK_CLOEXEC, 0) = 93
23903 connect(93, {sa_family=AF_UNIX, sun_path="/dev/socket/fwmarkd"}, 110) = 0
23903 sendmsg(93, {msg_name=NULL, msg_namelen=0, msg_iov=[{iov_base="\1\0\0\0\0\0\0\0\0\0\0\0\0\0\0\0", iov_len=16}, {iov_base=NULL, iov_len=0}], msg_iovlen=2,
msg_control=[{cmsg_len=20, cmsg_level=SOL_SOCKET, cmsg_type=SCM_RIGHTS, cmsg_data=[84]}], msg_controllen=24, msg_flags=0}, 0) = 16
23903 recvfrom(93, <unfinished ...>
23904 socket(AF_INET6, SOCK_STREAM, IPPROTO_IP <unfinished ...>
23903 <... recvfrom resumed> "\0\0\0\0", 4, 0, NULL, NULL) = 4
23904 <... socket resumed>)       = 98
23903 connect(84, {sa_family=AF_INET6, sin6_port=htons(443), inet_pton(AF_INET6, "::ffff:76.223.92.165", &sin6_addr), sin6_flowinfo=htonl(0), sin6_scope_id=0}, 28) = -1
EINPROGRESS (Operation now in progress)
23903 socket(AF_UNIX, SOCK_STREAM|SOCK_CLOEXEC, 0) = 93
23903 connect(93, {sa_family=AF_UNIX, sun_path="/dev/socket/fwmarkd"}, 110) = 0
23903 sendmsg(93, {msg_name=NULL, msg_namelen=0, msg_iov=[{iov_base="\6\0\0\0\0\0\0\0\0\0\0\0\0\0\0\0", iov_len=16},
```
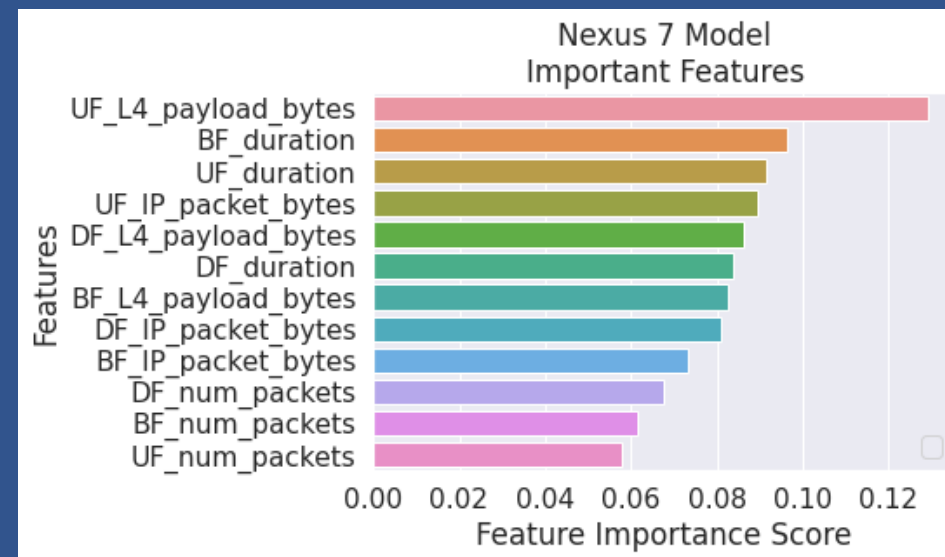
# Methodology – Data Analysis

- Traffic object we examine is the bi-directional flow
  - Uniquely identified by the 5-tuple of source IP, source port, destination IP, destination port, and which protocol (TCP or UDP)
  - These are not features, just unique identifiers
- Direction, timing, and size are preserved as a 'feature'

- Many other statistical features can then be created to describe these flows
  - E.g., total bytes sent, momentum of the conversation, in addition to the mean, max, min, variance, etc.

Poster: Data Collection for ML Classification of Encrypted Messaging Applications

# ML applications



- **Some initial proof-of-concept multi-class classification**

- **Off the shelf classifiers; in our experiments Random Forests worked very well.**

- **Trained a classifier on MIRAGE data's Nexus 7 flows to classify apps from a different phone's flows**

- **In this particular case, the upstream L4 payload was of high importance.**

  - **This intuitively suggests that the client side behavior is an important discriminator**

# Future Work

- Describe the system and considerations in greater detail to assist researchers
  - Emphasizing the partnership opportunities with host institutions
  - Allow other researchers to similarly extend the MIRAGE dataset

- ML applications
  - Extending the MIRAGE dataset with our own custom applications in the same format
  - Applying classifiers to 'real world' WiFi dataset from Mines
  - Expanding the 'positive class' from just a single application to the genre of Encrypted Messaging Applications