**Distinguished Paper Award, Internet Defense 2nd Prize!**

# Online Website Fingerprinting: Evaluating Website Fingerprinting Attacks on Tor in the Real World

Giovanni Cherubin, ~~Alan Turing Institute~~ Microsoft Research
Rob Jansen, U.S. Naval Research Laboratory
Carmela Troncoso, EPFL SPRING Lab

**Rob Jansen, Ph.D.**
Computer Security Research Scientist
Center for High Assurance Computer Systems
U.S. Naval Research Laboratory

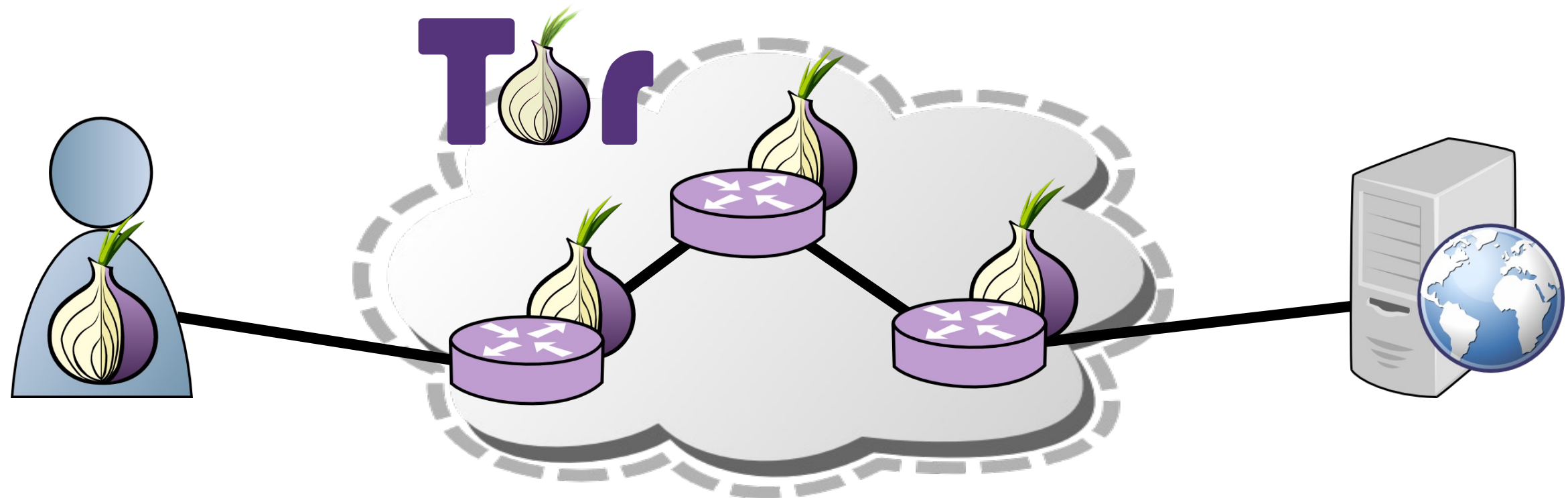31st USENIX Security Symposium
Boston, MA, USA
August 10th, 2022

## Anonymous Communication and Tor

- Separates identification from routing
- Provides unlinkable communication
- Promotes user safety and privacy online
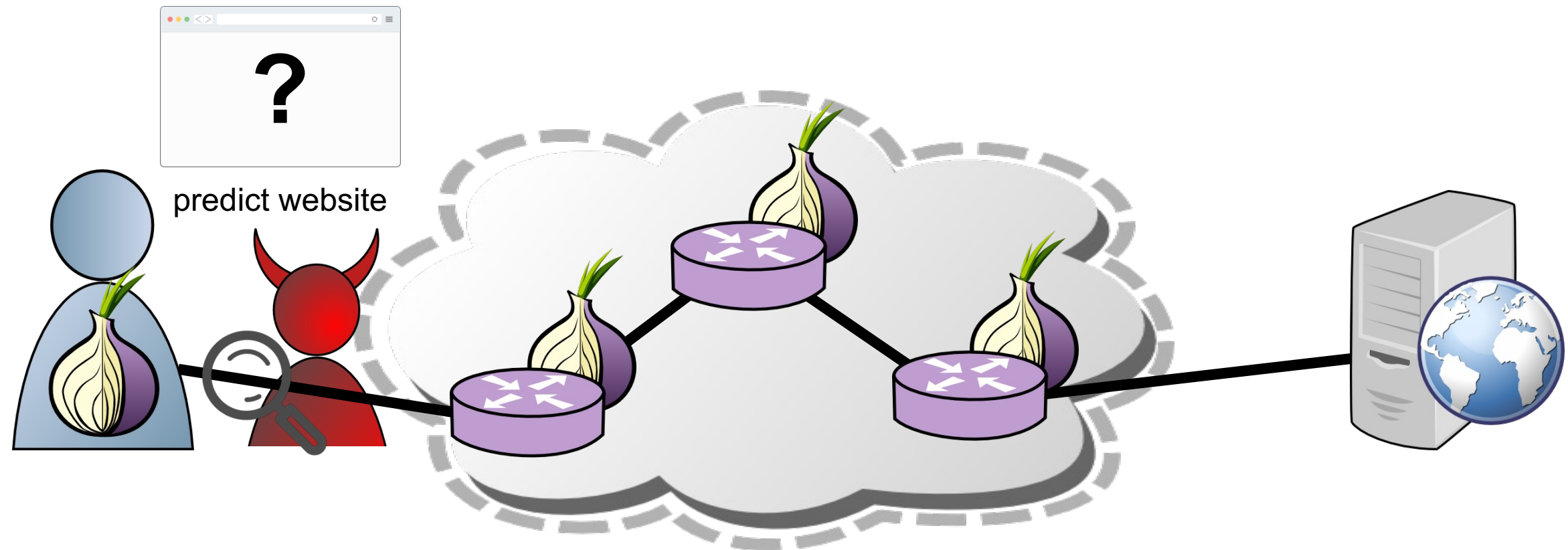
**Tor** Browse Privately.
Explore Freely.

Defend yourself against tracking and surveillance. Circumvent censorship.

Website fingerprinting attack
- Predict website visited by user
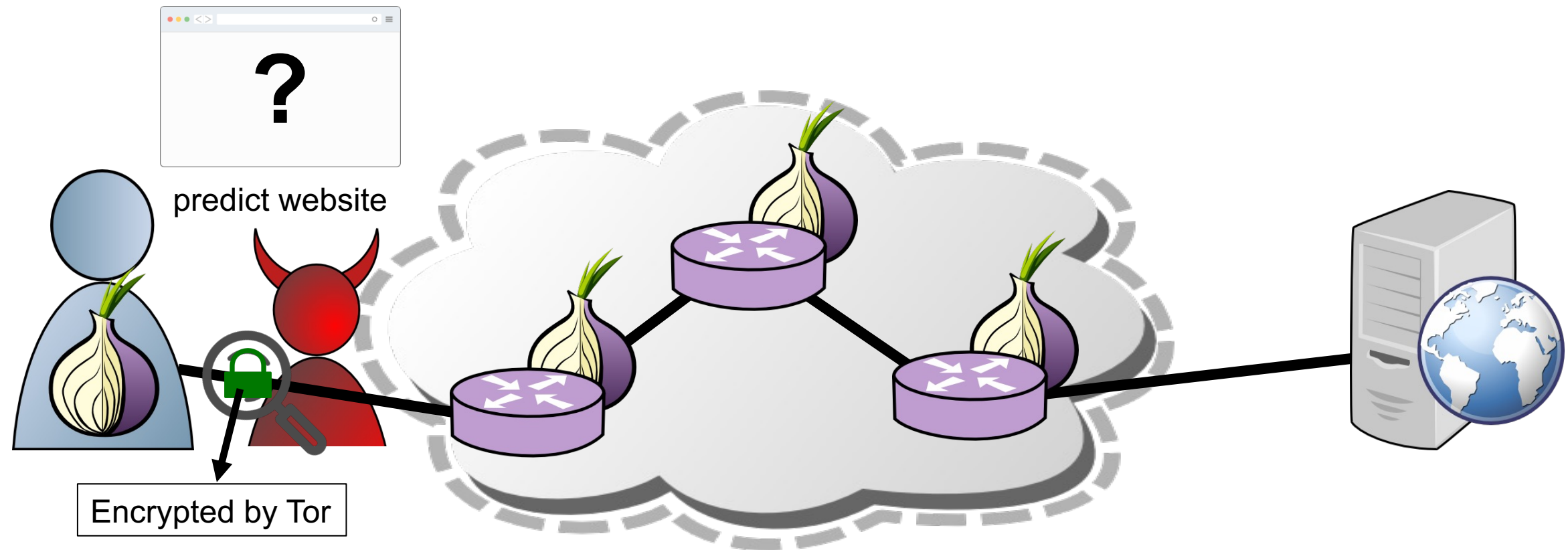- Requires access to entry side only



predict website

# Deanonymizing Tor Users

Website fingerprinting attack
- Predict website visited by user
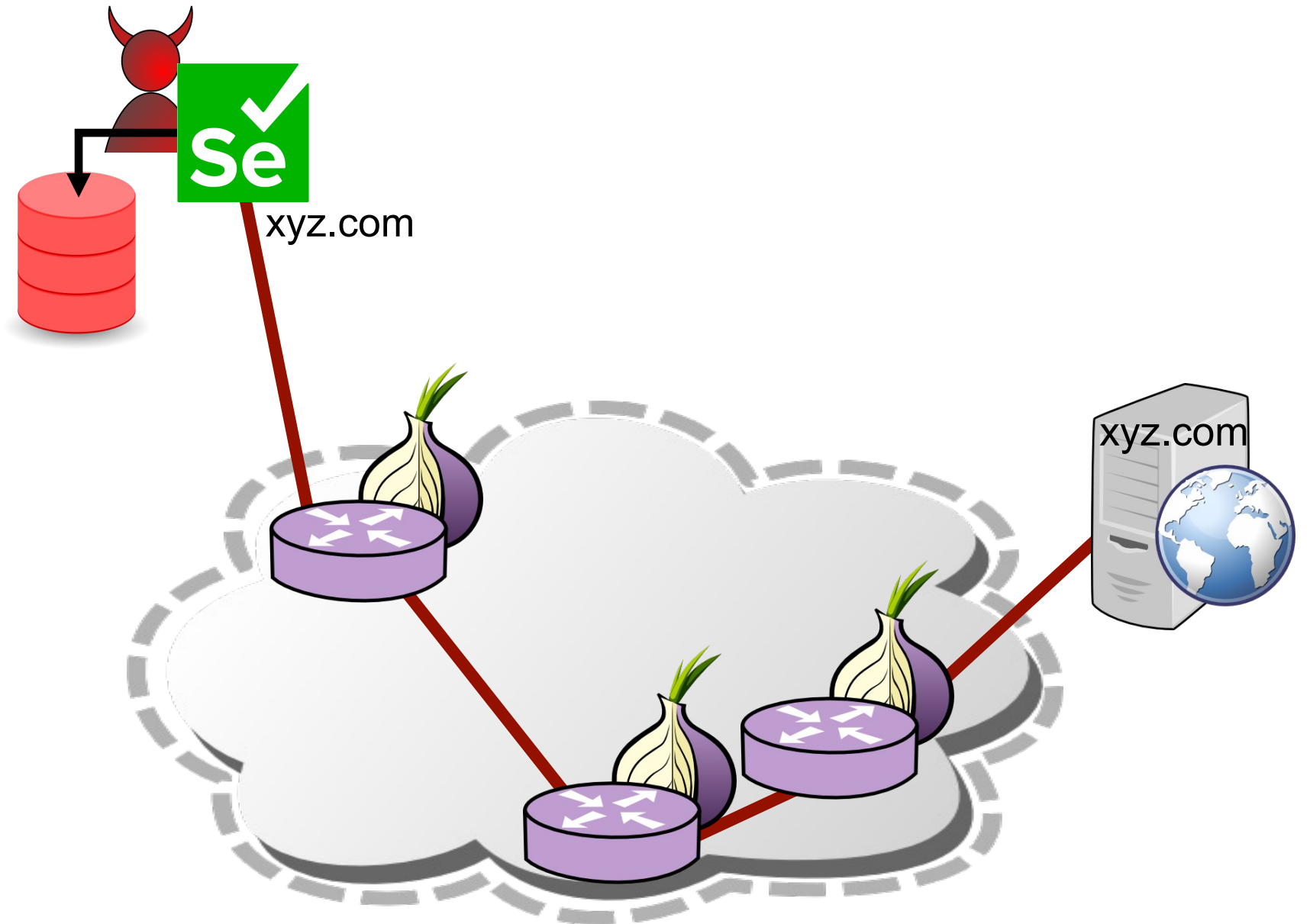- Requires access to underline{entry side only}

Problem:
- Need underline{labels} to train ML classifiers for website prediction
- Genuine labels are underline{encrypted}

? predict website

Encrypted by Tor

**Step 1: gather data & labels**
- Use automated browser (selenium) to crawl websites

xyz.com

xyz.com

Step 1: gather data & labels
- Use automated browser (selenium) to crawl websites

Step 2: train ML classifier
- Use collected data & labels
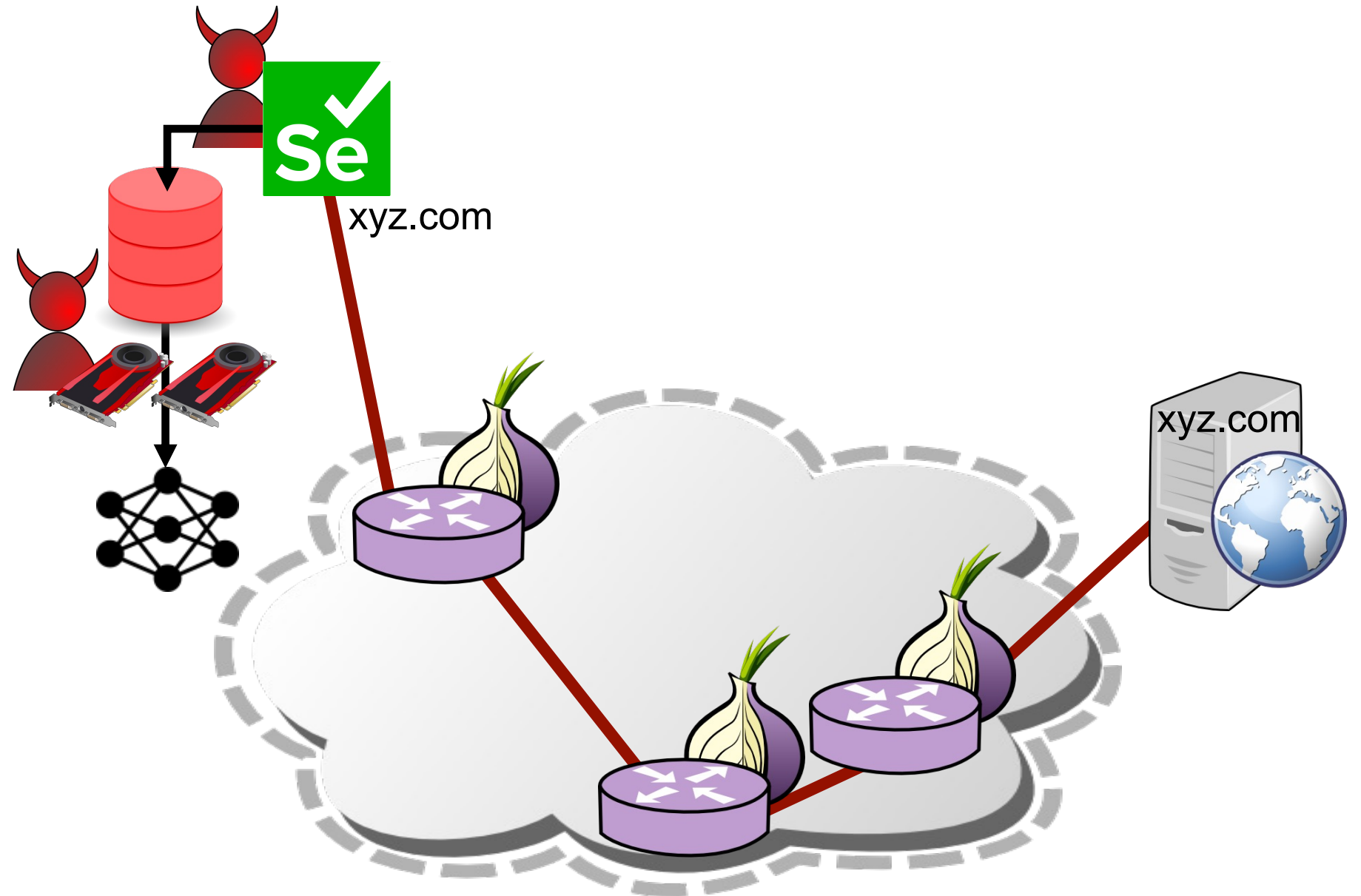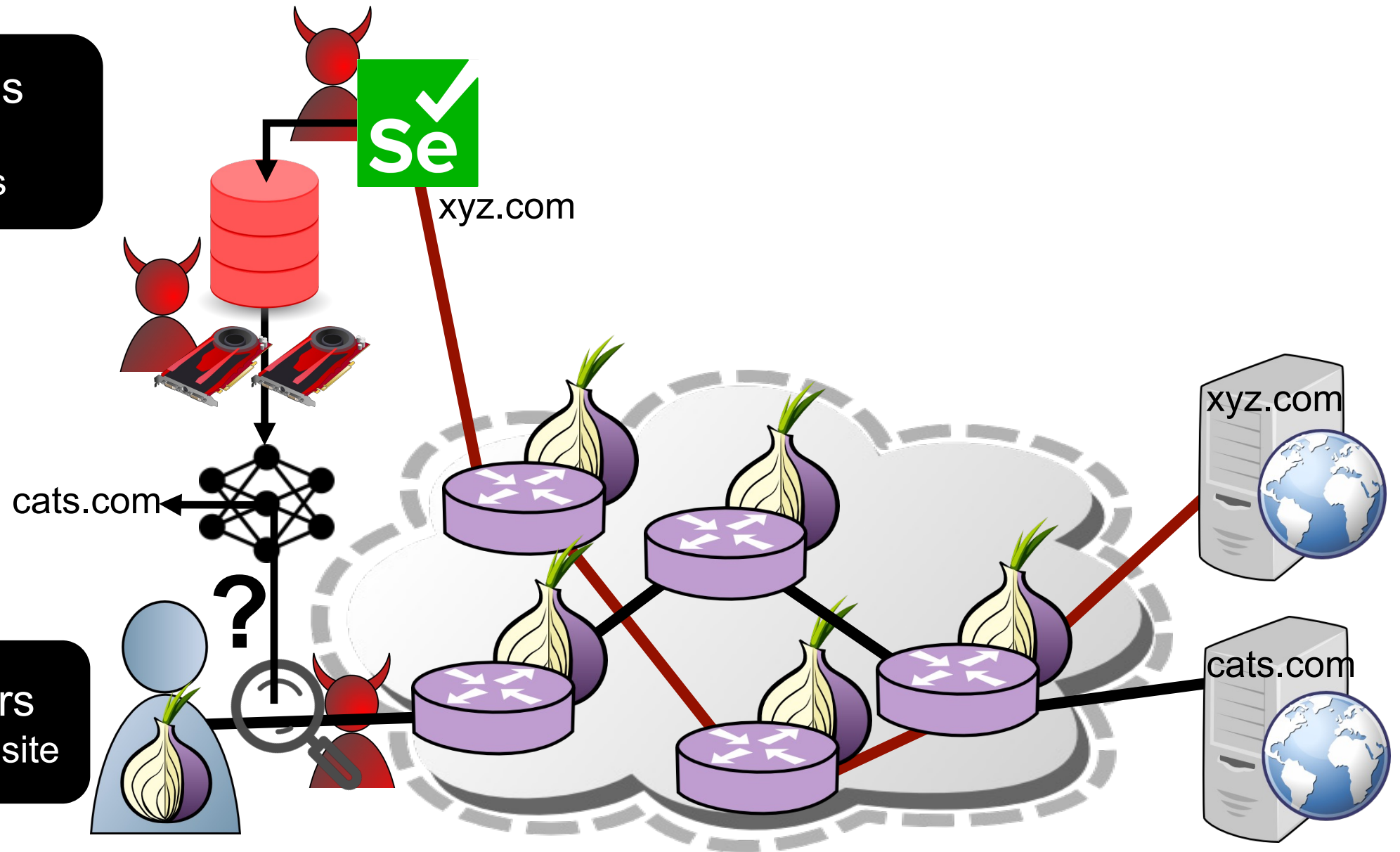
**Step 1: gather data & labels**
- Use automated browser (selenium) to crawl websites

**Step 2: train ML classifier**
- Use collected data & labels

**Step 3: deploy against users**
- Use ML model to predict website

# What is the threat of WF attacks in the *real* world?

**Synthetic model**
- Overly simple and unrealistic
- High ML accuracy in simple model

Stop using!!

**Genuine model**
- Consider genuine data & labels from a Tor exit relay

Our new approach

# Key Insight: Exits Observe Genuine Data & Labels

Step 1: gather data & labels
- Run a Tor exit relay and use to to collect genuine Tor traffic

Exit can observe:
1. New circuit
2. DNS lookup
3. Website load

Genuine labels: resolved domains

Genuine data: circuit traffic patterns

cats.com?

198.71.232.3

DNS

cats.com
198.71.232.3

# Key Insight: Exits Observe Genuine Data & Labels

**Step 1: gather data & labels**
- Run a Tor exit relay and use to to collect genuine Tor traffic

**Step 2: train ML classifier**
- Use collected data & labels

Genuine labels: resolved domains

Genuine data: circuit traffic patterns

**Exit can observe:**
1. New circuit
2. DNS lookup
3. Website load

cats.com?

DNS

198.71.232.3

cats.com
198.71.232.3

Step 1: gather data & labels
- Run a Tor exit relay and use to to collect genuine Tor traffic

Step 2: train ML classifier
- Use collected data & labels

Step 3: deploy against users
- Use ML model to predict website

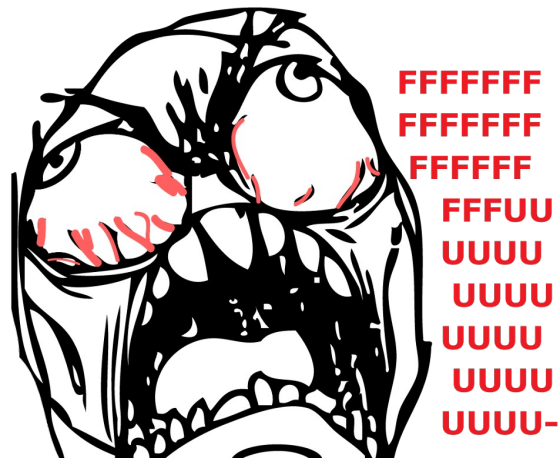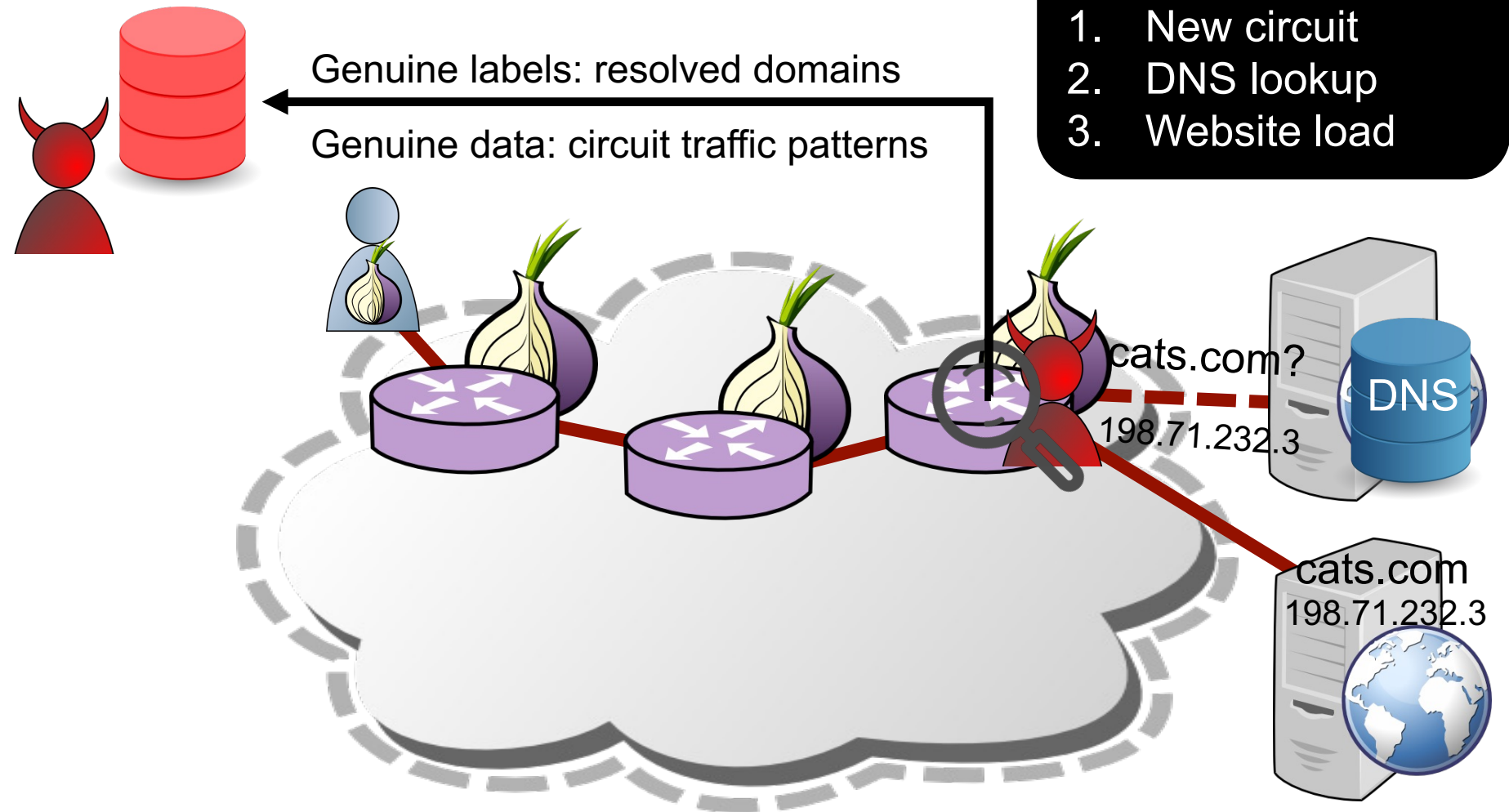Genuine labels: resolved domains

Genuine data: circuit traffic patterns

Exit can observe:
1. New circuit
2. DNS lookup
3. Website load

cats.com

cats.com?

DNS

198.71.232.3

cats.com
198.71.232.3

# Key Insight: Exits Observe Genuine Data & Labels

**Benefits**
- Captures real world diversity of browsers, behavior, world size, choice of pages
- Can stop trying to fix the synthetic model

**Exit can observe:**
1. New circuit
2. DNS lookup
3. Website load

Genuine labels: resolved domains

Genuine data: circuit traffic patterns

cats.com

cats.com?

198.71.232.3

DNS

cats.com
198.71.232.3

### Benefits
- Captures real world diversity of browsers, behavior, world size, choice of pages
- Can stop trying to fix the synthetic model

### Caveats
- Train at exit, deploy at entry → noise
- Domain, not page label
- Need safe eval methods

### Exit can observe:
1. New circuit
2. DNS lookup
3. Website load

Genuine labels: resolved domains

Genuine data: circuit traffic patterns

cats.com

?

cats.com?

198.71.232.3

DNS

cats.com
198.71.232.3

# Our safe evaluation plan:

- Hash domain labels using keyed HMAC
    - Never learn true labels



data: (-1,+1,…)
label: HMAC(cats.com)

DNS

cats.com?

cats.com

198.71.232.3

## Our safe evaluation plan:

- ## Hash domain labels using keyed HMAC
  - Never learn true labels

- ## Use online learning
  - Adapted Triplet Fingerprinting [CCS'19]
  - Compute means in real time, discard data
  - Individual data items never stored



k-nn model

1. predict label

HMAC(cats.com)
correct?
yes or no

2. update k-nn mev

triplet feature extractor

data: (-1,+1,…)
label: HMAC(cats.com)

DNS

cats.com?

cats.com

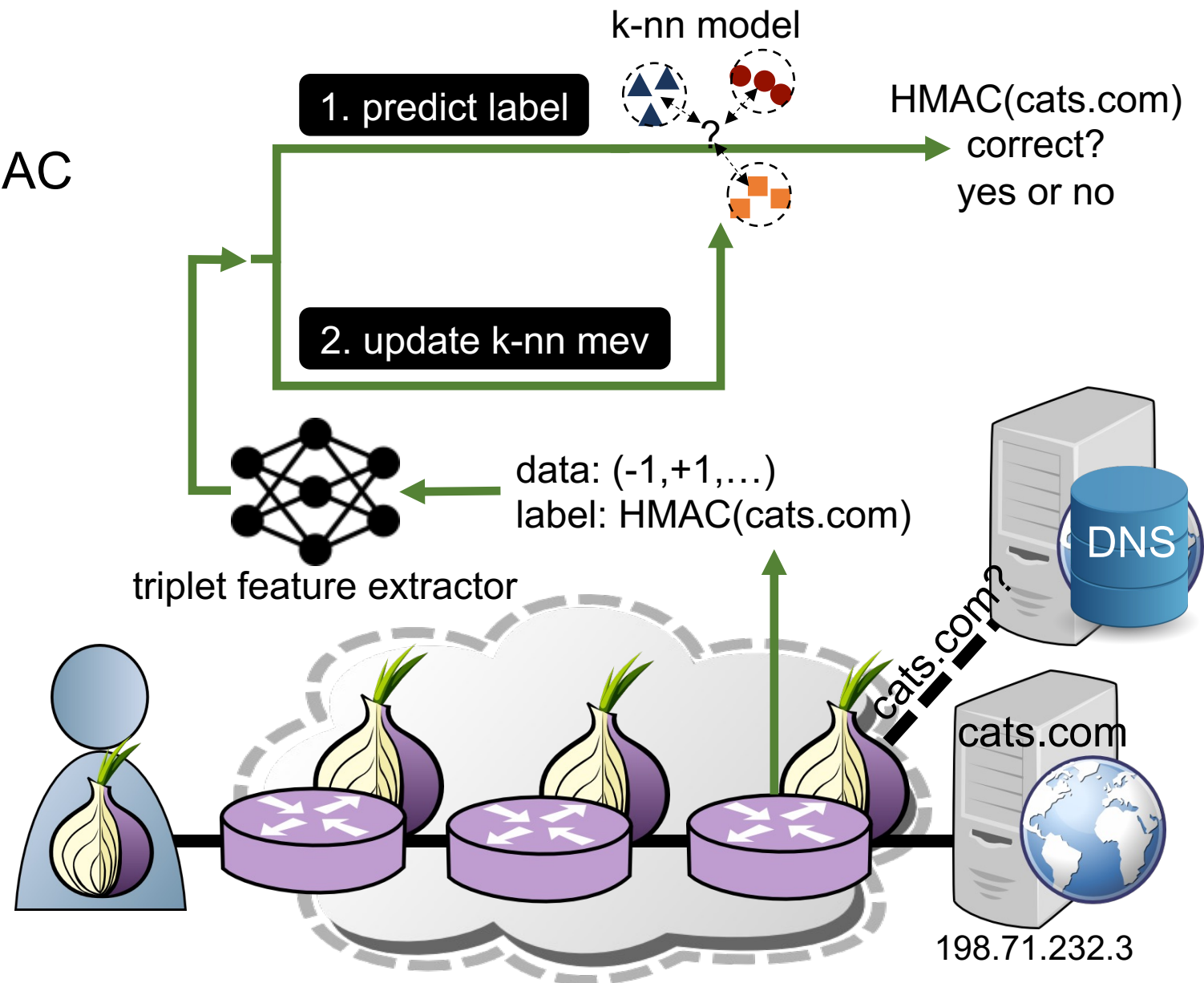198.71.232.3

## Our safe evaluation plan:

- Hash domain labels using keyed HMAC
  - Never learn true labels

- Use online learning
  - Adapted Triplet Fingerprinting [CCS'19]
  - Compute means in real time, discard data
  - Individual data items never stored

- Other safety precautions
  - Never deanonymizes Tor users
  - Destroyed models, HMAC key after eval

- Tor Safety Board reviewed plan
  - See paper for details!



k-nn model

1. predict label

HMAC(cats.com)
correct?
yes or no

2. update k-nn mev

data: (-1,+1,…)
label: HMAC(cats.com)

triplet feature extractor

cats.com?

DNS

cats.com

198.71.232.3

## Train and evaluate at exit relay

- No noise from transferring to entry
- Upper bound on attack accuracy

## Details

- 1 week evaluation
  - 3.9M data sequences, 671k unique sites

- Multi-class classification
  - predict a monitored site, or 'unmonitored'

- Performance metric
  - instant accuracy (i.e., moving average)
  - # correct / # total predictions (10k window)

# Train and evaluate at exit relay

- No noise from transferring to entry
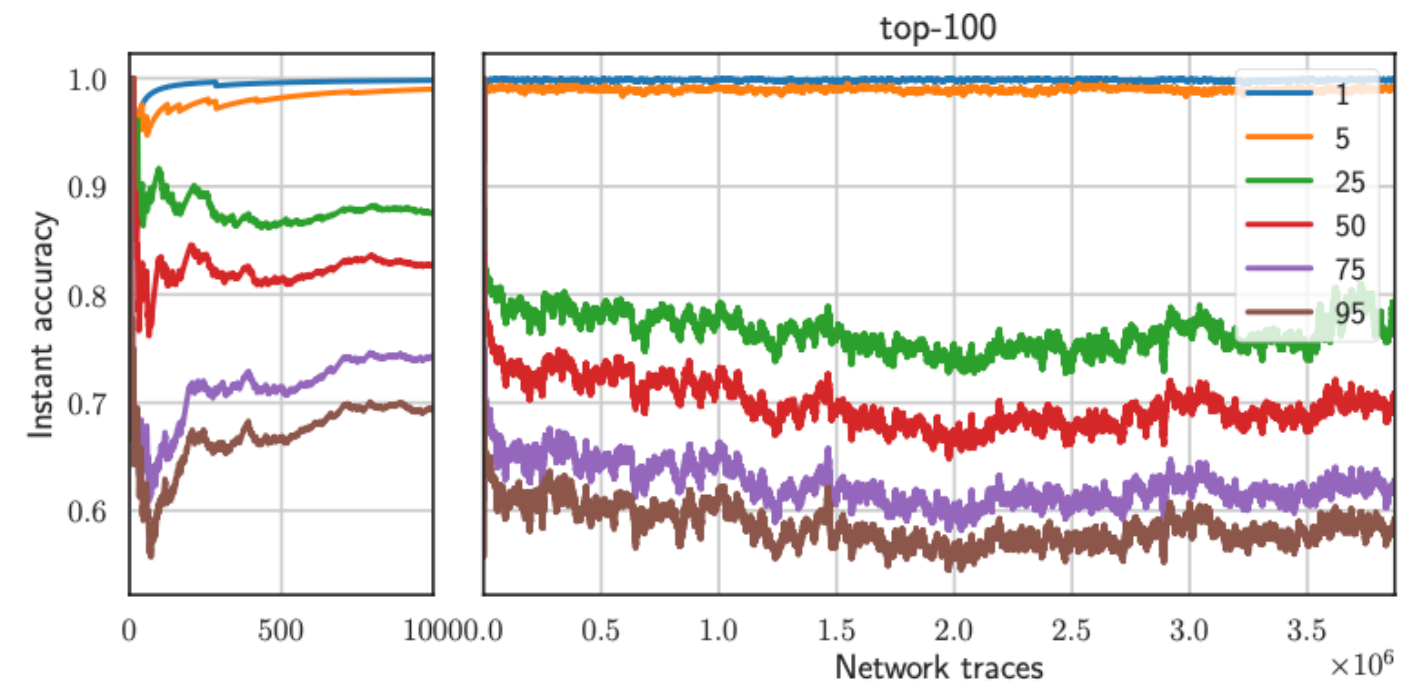- Upper bound on attack accuracy

# Details

- 1 week evaluation
  - 3.9M data sequences, 671k unique sites

- Multi-class classification
  - predict a monitored site, or 'unmonitored'

- Performance metric
  - instant accuracy (i.e., moving average)
  - # correct / # total predictions (10k window)

## Train and evaluate at exit relay

- No noise from transferring to entry
- Upper bound on attack accuracy

## Details

- 1 week evaluation
  - 3.9M data sequences, 671k unique sites

- Multi-class classification
  - predict a monitored site, or 'unmonitored'

- Performance metric
  - instant accuracy (i.e., moving average)
  - # correct / # total predictions (10k window)

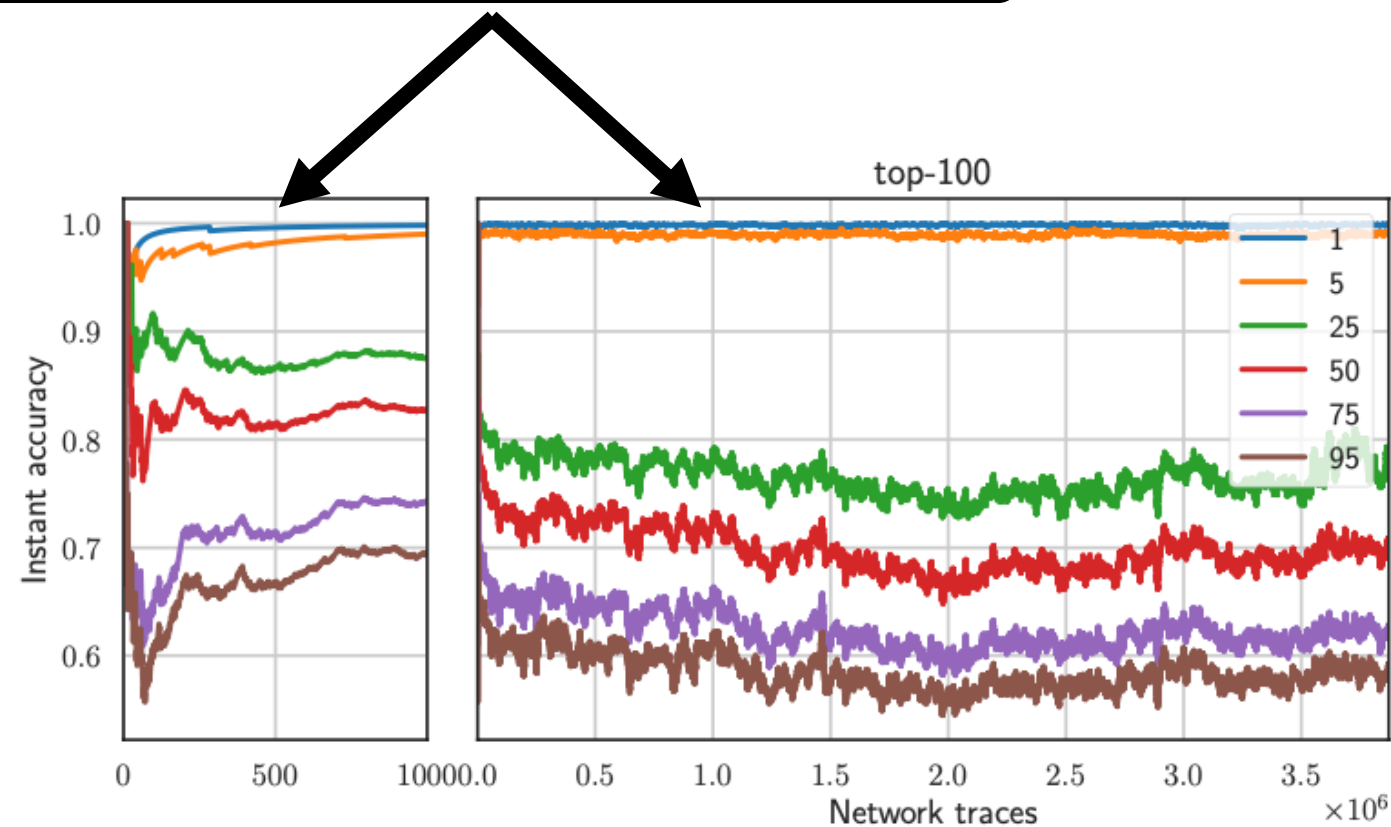**accuracy above 95% when monitoring ≤ 5 sites**

# Train and evaluate at exit relay

- No noise from transferring to entry
- Upper bound on attack accuracy

# Details

- 1 week evaluation
  - 3.9M data sequences, 671k unique sites

- Multi-class classification
  - predict a monitored site, or 'unmonitored'

- Performance metric
  - instant accuracy (i.e., moving average)
  - # correct / # total predictions (10k window)

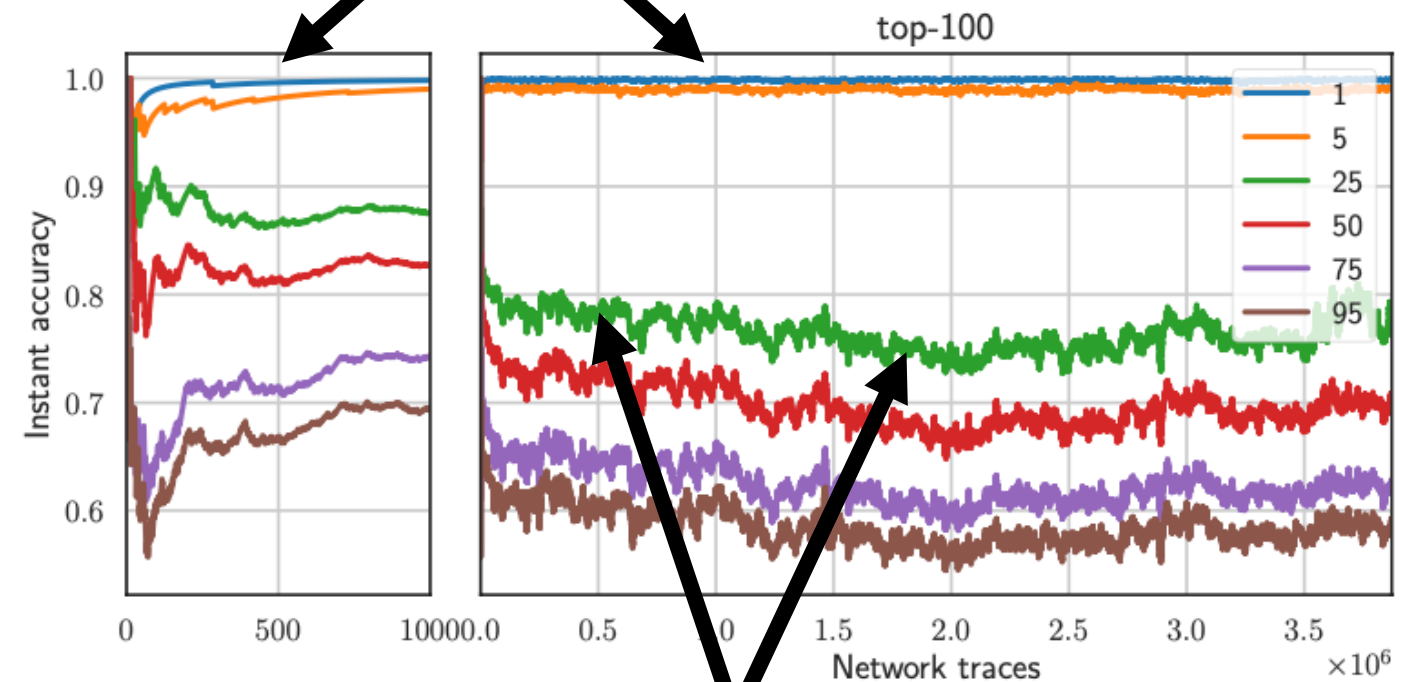accuracy above 95% when monitoring ≤ 5 sites

accuracy quickly falls below 80% when monitoring ≥ 25 sites

# Genuine vs. Synthetic Data

## Offline phase

- Crawl 'synthetic' list of domains
  - <u>Synthetic</u>: use crawl to train a classifier offline

## Online phase

- Train two classifiers online
  - <u>Hybrid</u>: update copy of synthetic classifier with genuine data
  - <u>Real</u>: train new classifier on genuine data only

- 1 week evaluation
  - 1.2M data sequences
  - observed 183 of 1,074 'synthetic' domains

- Binary classification
  - monitored set contains 5 sites
  - predict either 'monitored' or 'unmonitored'

## Offline phase

- Crawl 'synthetic' list of domains
  - <u>Synthetic</u>: use crawl to train a classifier offline

## Online phase

- Train two classifiers online
  - <u>Hybrid</u>: update copy of synthetic classifier with genuine data
  - <u>Real</u>: train new classifier on genuine data only

- 1 week evaluation
  - 1.2M data sequences
  - observed 183 of 1,074 'synthetic' domains

- Binary classification
  - monitored set contains 5 sites
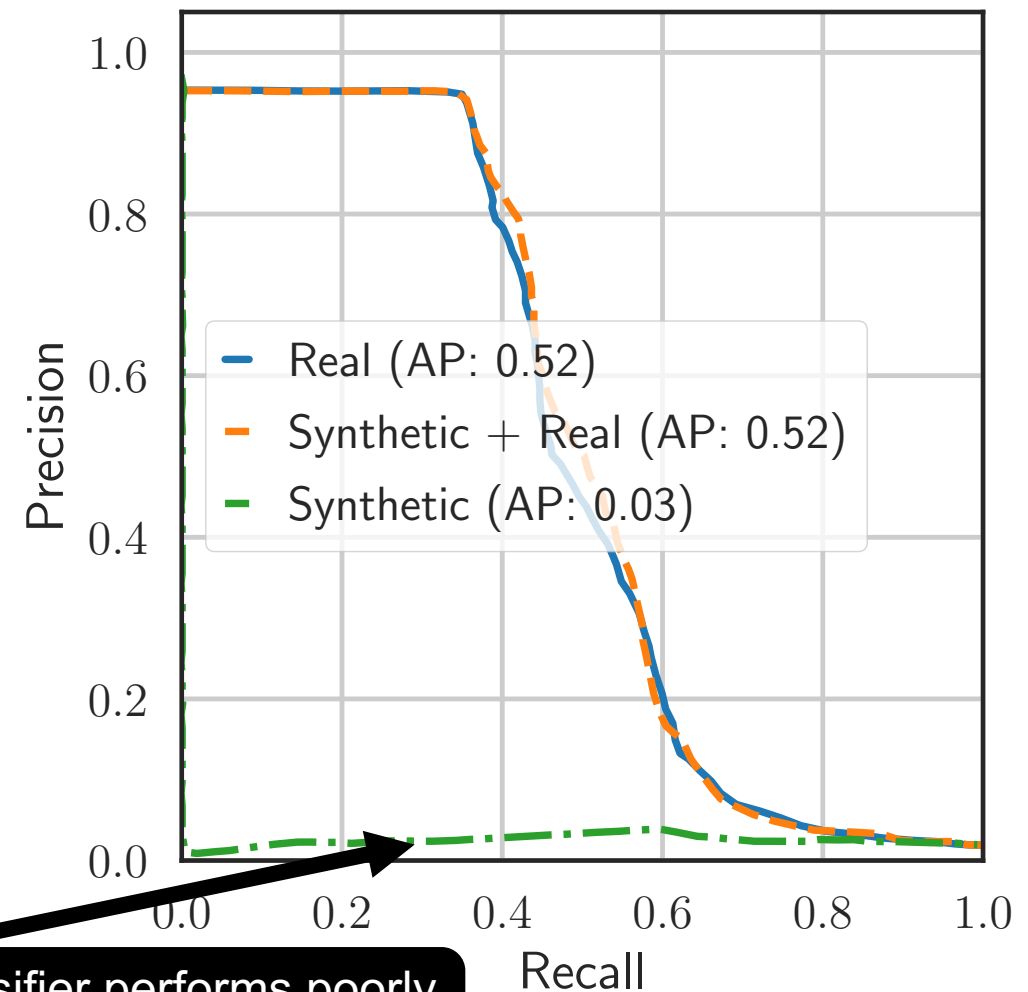  - predict either 'monitored' or 'unmonitored'



Legend:
- Real (AP: 0.52)
- Synthetic + Real (AP: 0.52)
- Synthetic (AP: 0.03)

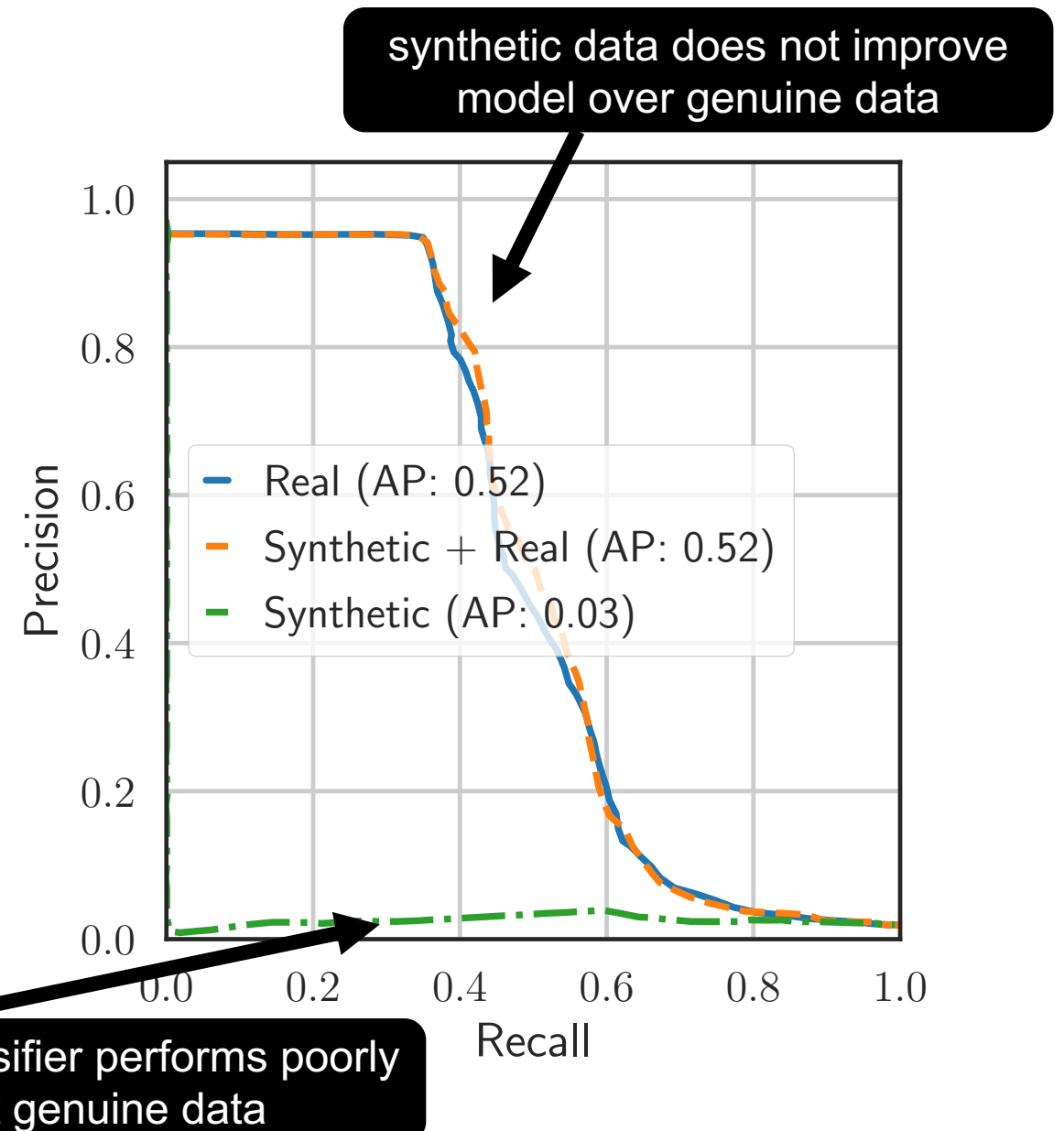synthetic classifier performs poorly against genuine data

# Genuine vs. Synthetic Data

## Offline phase

- Crawl 'synthetic' list of domains
  - <u>Synthetic</u>: use crawl to train a classifier offline

## Online phase

- Train two classifiers online
  - <u>Hybrid</u>: update copy of synthetic classifier with genuine data
  - <u>Real</u>: train new classifier on genuine data only

- 1 week evaluation
  - 1.2M data sequences
  - observed 183 of 1,074 'synthetic' domains

- Binary classification
  - monitored set contains 5 sites
  - predict either 'monitored' or 'unmonitored'

synthetic data does not improve model over genuine data

- Real (AP: 0.52)
- Synthetic + Real (AP: 0.52)
- Synthetic (AP: 0.03)

synthetic classifier performs poorly against genuine data

## Fully synthetic evaluation

- Crawled 1k URLs 10x each

- Pinned entry and exit on each circuit

- Collected data sequences in both positions on each circuit

- Closed-world batch classification
  - 50%-50% train-test split

| Monitored set size: | 5 | 50 | 750 |
|---|---|---|---|
| Train and test on <u>exit</u> | 91.2% | 76.2% | 52.2% |
| Train on <u>exit</u>, test on <u>entry</u> | 86.4% | 65.1% | 34.1% |
| **Loss in accuracy:** | **4.8%** | **11.1%** | **18.1%** |

loss in accuracy is low for feasible (i.e. small) monitored sets
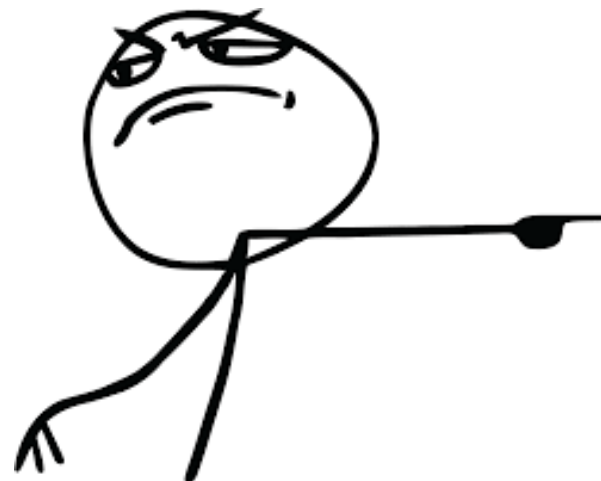
# Main Takeaways

## Insights

- WF can be feasible with genuine data and small monitored sets, online learning can mitigate concept drift

- Synthetic data is not useful when the adversary deploys in the real world

- Simple defenses may be more effective than we thought
  - Adversary has to simulate defense on top of undefended exit data

## Contact

- rob.g.jansen@nrl.navy.mil
- robgjansen.com
- @robgjansen

## Future Research Areas

- Improve accuracy when training on genuine data

- Reduce distortion when transferring models from exit to entry

- Defenses that make it harder to learn from genuine data, increase distortion

**Read the paper!**