# AI/ML for Network Security: The Emperor has no Clothes

**Arthur Selle Jacobs[1]**

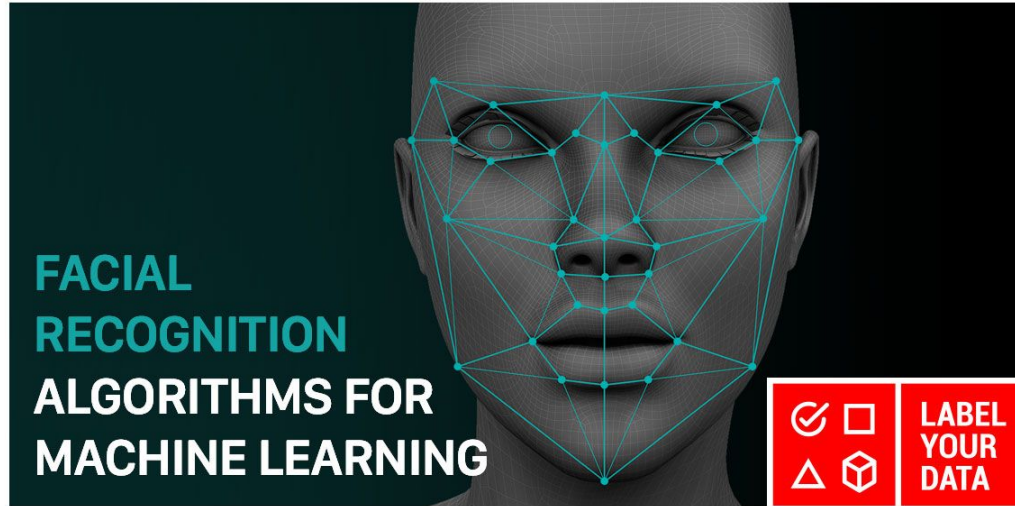Roman Beltiukov[4]          Walter Willinger[2]          Ronaldo A. Ferreira[3]

Arpit Gupta[4]          Lisandro Z. Granville[1]

*November 11th, 2022*

# The Rise of AI

# The Rise of AI

FACIAL
RECOGNITION
ALGORITHMS
MACHINE LEA

# The Rise of AI

**AI & MACHINE LEARNING**

## How Kaggle solved a spam problem in 8 days using AutoML

**Will Cukierski**
Staff Developer Advocate and
Head of Competitions, Kaggle

May 27, 2020

**Try Google Cloud**

Start building on Google Cloud
with $300 in free credits and
20+ always free products.

**FREE TRIAL**

Kaggle is a data science community of nearly 5 million users. In September of 2019, we found ourselves under a sudden siege of spam traffic that threatened to overwhelm visitors to our site. We had to come up with an effective solution, fast. Using AutoML Natural Language on Google Cloud, Kaggle was able to train, test, and deploy a spam detection model to production in just eight days. In this post, we'll detail our success story about using machine learning to rapidly solve an urgent business dilemma.

### A spam dilemma

Malicious users were suddenly creating large numbers of Kaggle accounts in order to leave spammy search engine optimization (SEO) content in the user bio section. Search engines were indexing these bios, and our existing spam detection heuristics were failing to flag them. In short, we faced a growing and embarrassing predicament.

Our problem was context. Kaggle is a community focused on data science and machine learning. As a result of our topical data-science focus, a user bio that seems harmless in isolation may be the work of a spammer. Here is a real example of one such bio:

FACIAL
RECOGNITION
ALGORITHMS F
MACHINE LEAR

# How does it work?

Traditional AI/ML Development Pipeline

Collect Data

Select Model

# How does it work?

Traditional AI/ML Development Pipeline
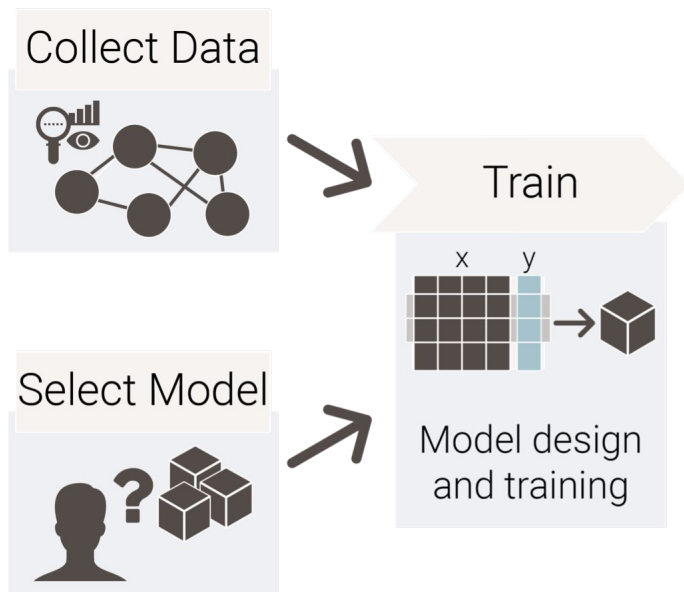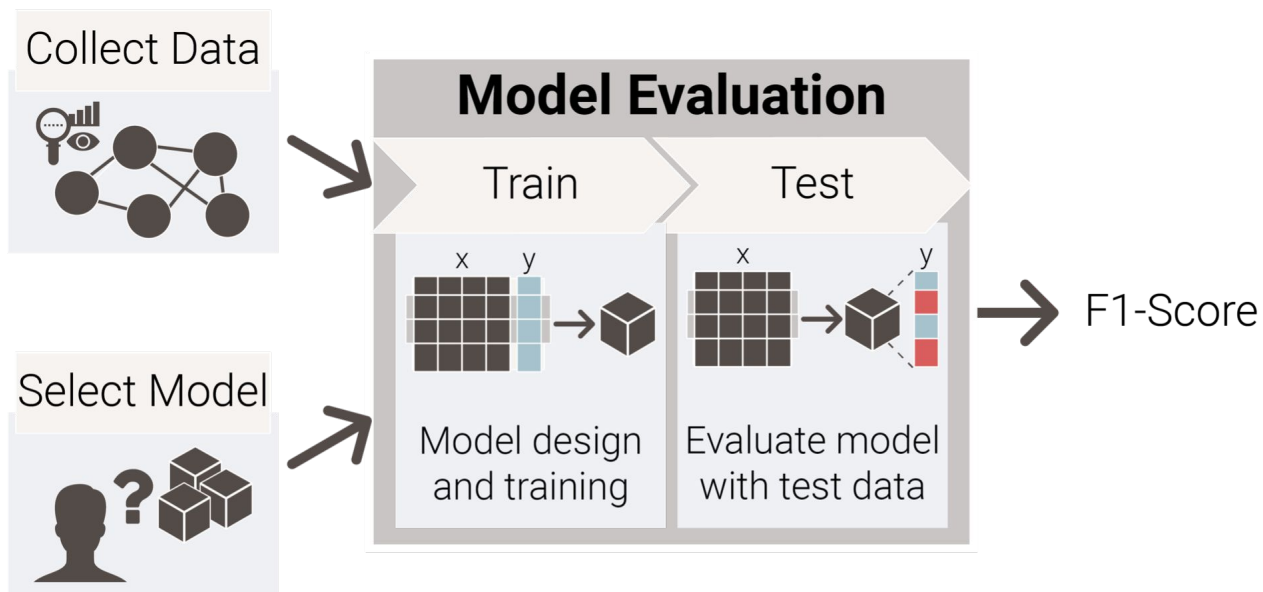
Collect Data
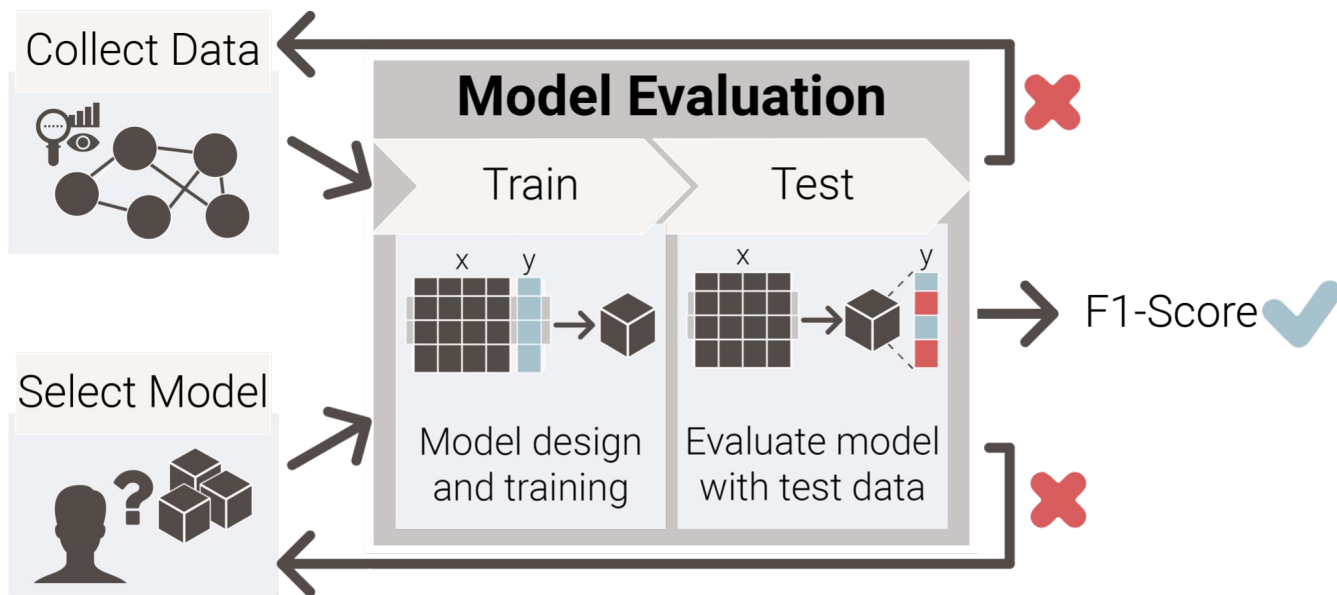
Train

x    y

Select Model

Model design
and training

# How does it work?

Traditional AI/ML Development Pipeline

# How does it work?

Traditional AI/ML Development Pipeline

# What about high-stakes decision making?

**Why (and how) does the model work?**



Self-driving Cars

**When does the model not work?**



Network Security

# Underspecification issues!

**Shortcut Learning**

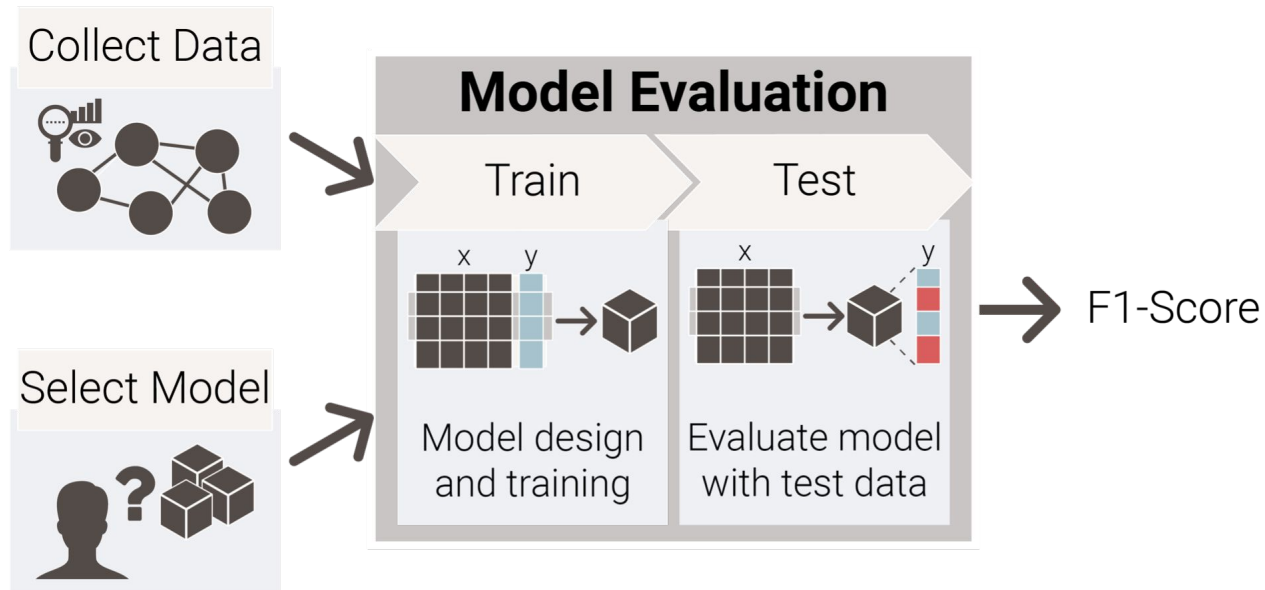Model takes shortcuts to classify data!
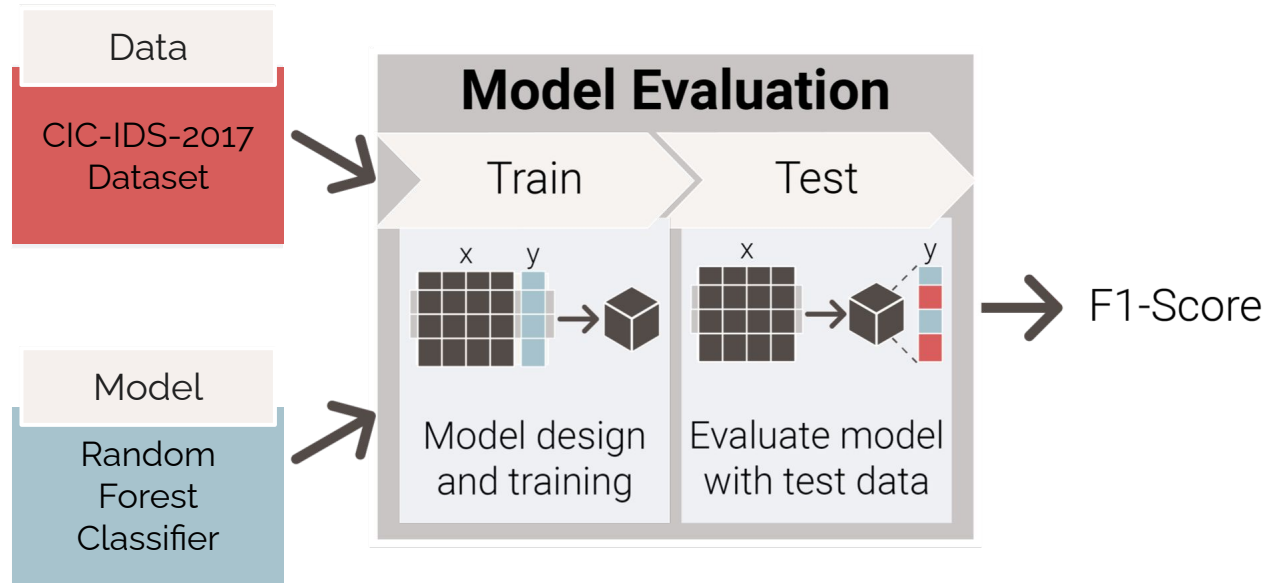
**O.O.D. Samples**

Model does not generalize!

**Spurious Correlations**
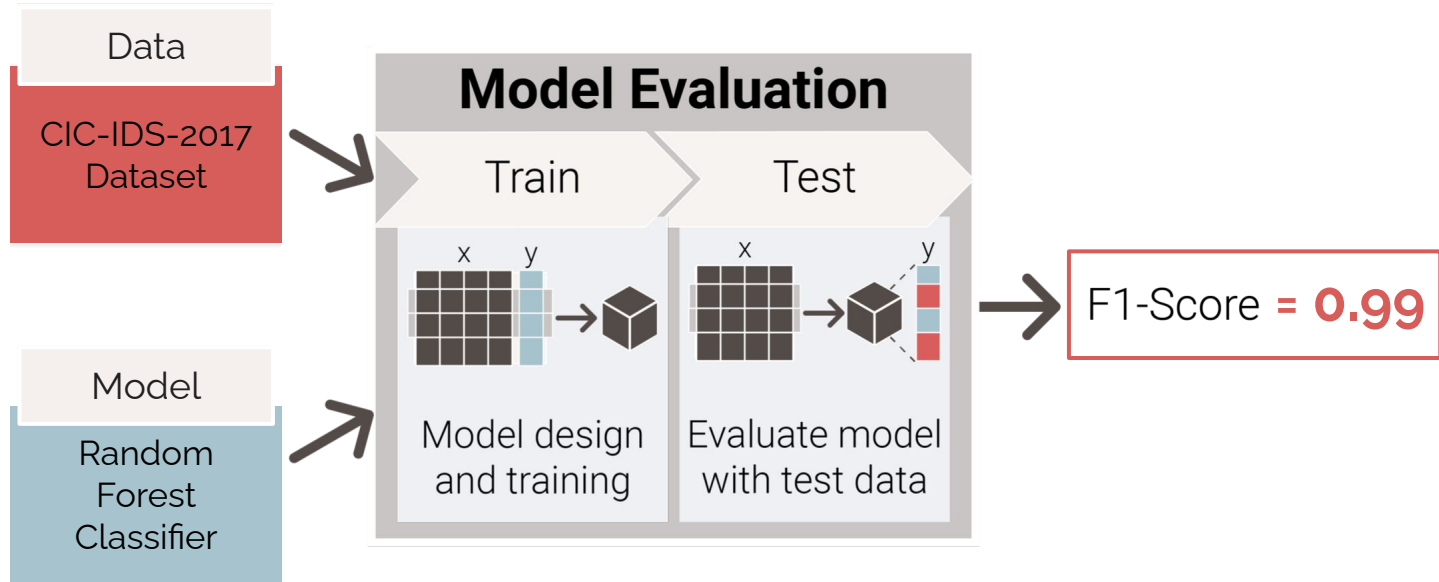
Model picks up wrong correlations in the data!

# Consider this example…

# Consider this example...

# Consider this example...



**Data**
CIC-IDS-2017 Dataset

**Model**
Random Forest Classifier

**Model Evaluation**

Train — Model design and training

Test — Evaluate model with test data

F1-Score = **0.99**

# Can you answer these questions?

**Why (and how) does the model work?**   **When does the model not work?**

# Can you answer these questions?

**Why (and how) does the model work?**
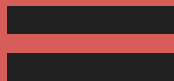
**When does the model not work?**
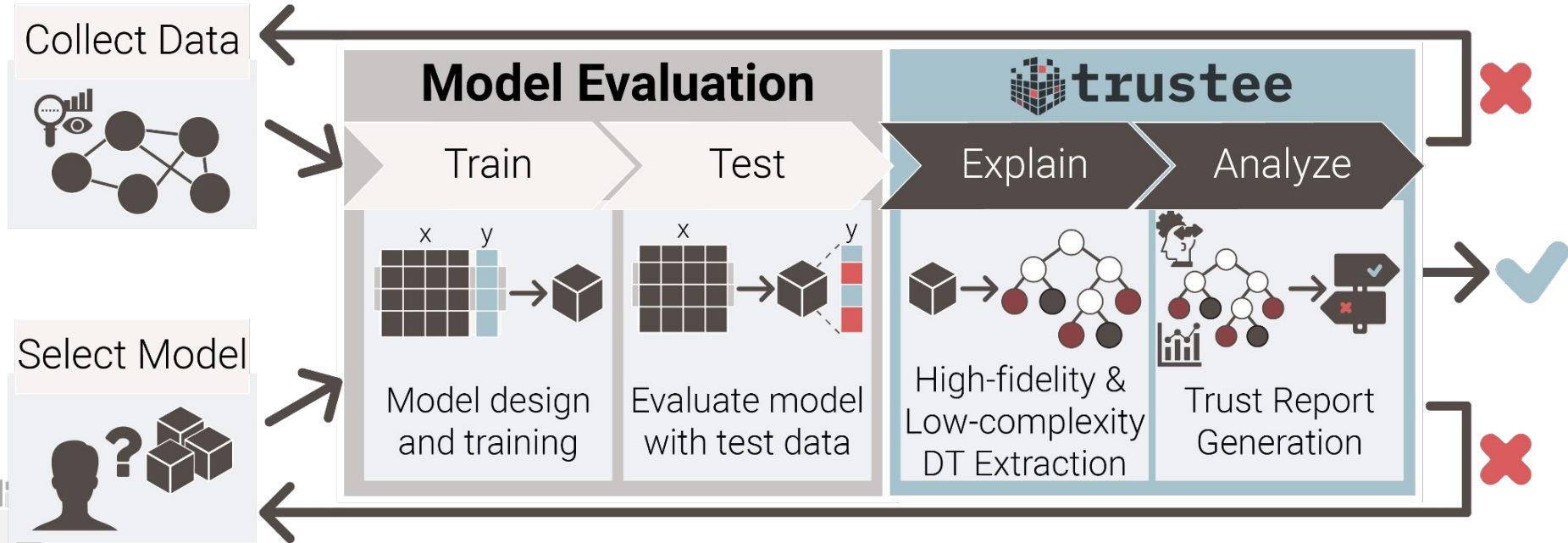
Can you trust this model?
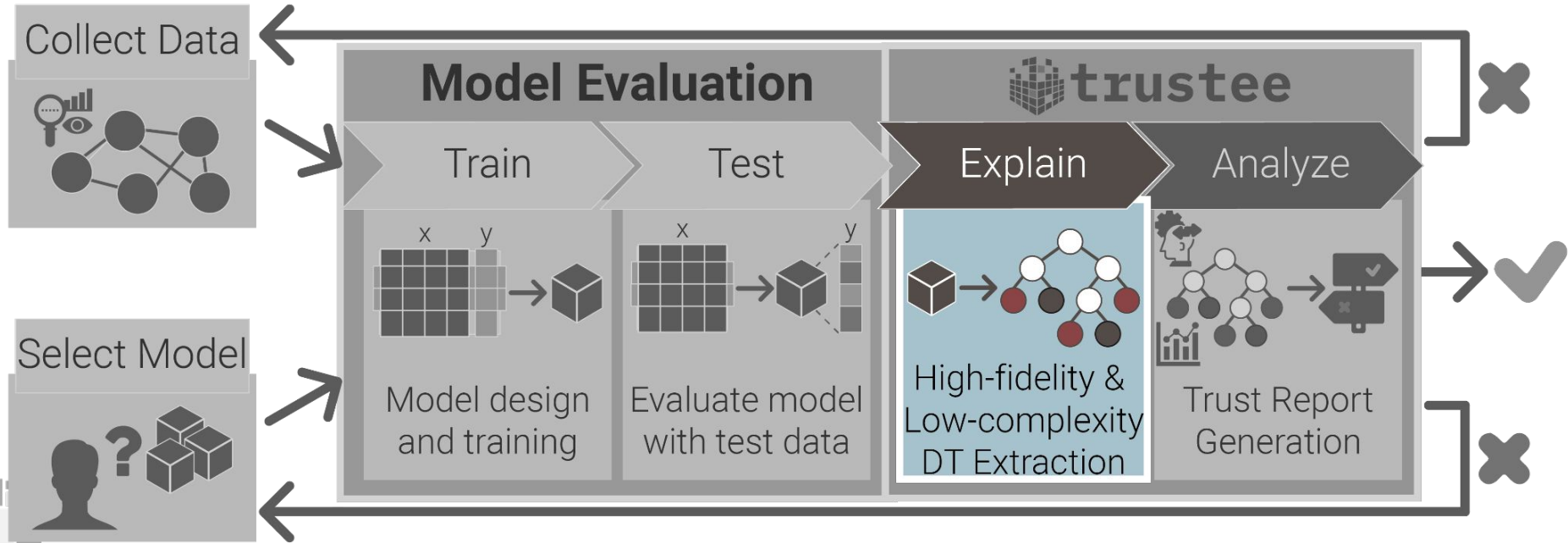
# Can you trust this model?

Trust in AI/ML model
=
Hand over control to the AI/ML model
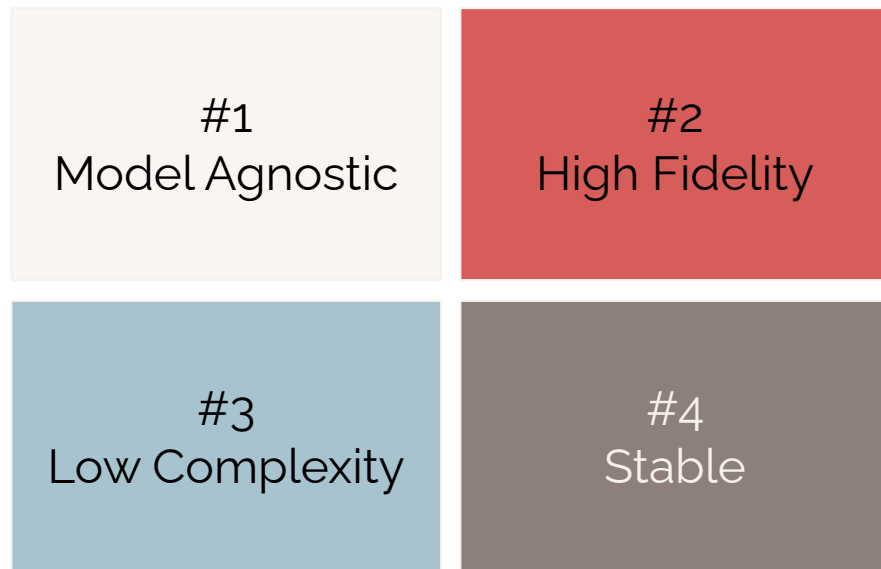
# Augmented
## AI/ML Development Pipeline



Collect Data

Select Model

**Model Evaluation**

**trustee**

Train — Test — Explain — Analyze

Model design and training

Evaluate model with test data

High-fidelity & Low-complexity DT Extraction

Trust Report Generation

# Augmented
## AI/ML Development Pipeline

# trustee

## Explanation Requirements

| | |
|---|---|
| **#1** <br> **Model Agnostic** | **#2** <br> **High Fidelity** |
| **#3** <br> **Low Complexity** | **#4** <br> **Stable** |

# trustee

Dataset

Black-box
Model

trustee

70% Train Dataset

30% Test Dataset

Black-box Model

# trustee



70%

Train Dataset

30%

Test Dataset

Black-box
Model

Expected Output

70%

Train Dataset

M samples

30%

Test Dataset

Expected Output

# trustee

Inner Loop #1...N

Train Dataset

M samples

70%

30%

Test Dataset

Expected Output

Iteration #N

70% Train Dataset 1

30% Test Dataset 1

**trustee**

Inner Loop
#1...N

DT with Best
Fidelity

# trustee

Inner Loop
#1...N

⇨

DT with Best
Fidelity

trustee

Inner Loop
#1...N

Size matters!

DT with Best
Fidelity

# trustee

Inner Loop
#1...N

$\Rightarrow$

Top-k Pruning



DT with Best
Fidelity

# Top-k Pruning

Fidelity

Samples

# Top-k Pruning

Fidelity

Samples



Diminishing returns!

# Top-k Pruning

## Fidelity



## Samples

trustee

Outer Loop
#1...S

Inner Loop
#1...N

Top-k Pruning

DT with Best
Fidelity

DT with Highest
Agreement

#4
Stable

trustee

Outer Loop
#1...S

Inner Loop
#1...N

Top-k Pruning

DT with Best
Fidelity

DT with Highest
Agreement

Dataset

Black-box Model

**trustee**

DT with Highest Agreement

# Augmented
## AI/ML Development Pipeline

# Generating Trust Reports

Underspecification issues!

(revisited)

| Shortcut Learning | O.O.D. Samples | Spurious Correlations |
|---|---|---|
| Model takes shortcuts to classify data! | Model does not generalize! | Model makes the picks up wrong correlations in the data! |

# Generating Trust Reports

# Generating Trust Reports



48

# Generating Trust Reports

# Generating Trust Reports

# Generating Trust Reports

# Use Case #1: Detecting VPN vs Non-VPN Traffic

## Problem Setup

- **Selected publication**:
  - *"End-to-end encrypted traffic classification with one-dimensional convolution neural networks"* — Wang et al., 2017
- **Proposal**:
  - **Model**: 1D-CNN to classify traffic between encrypted VPN traffic and non-encrypted traffic (i.e. VPN vs Non-VPN)
  - **Features**: first 784 raw bytes of each PCAP file
  - **Dataset**: ISCX VPN-nonVPN 2016 [https://www.unb.ca/cic/datasets/vpn.html]
- **Results**:
  - Reported F1-score: 0.99
  - Reproduced F1-score: 0.959

# Use Case #1: Detecting VPN vs Non-VPN Traffic

**Explanation**

Fidelity: 1.000
No pruning
7 nodes

$B_{49} \leq 17$

True ⟶ 33%

False ⟶ 67%

$B_{43} \leq 1$

$B_{47} \leq 251$

1% — Non VPN

32% — VPN

66% — Non VPN

1% — VPN

# Use Case #1: Detecting VPN vs Non-VPN Traffic

**Explanation**

## Non VPN

| | 0 | | | | | | | | 9 | 10 | | | | | | | | | 19 | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **Pcap** | 161 | 178 | 195 | 212 | 0 | 2 | 0 | 4 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 255 | 255 |
| **Meta** | 0 | 0 | 0 | 1 | 85 | 65 | 10 | 69 | 0 | 5 | 80 | 24 | 0 | 0 | 0 | 64 | 0 | 0 | 0 | 64 |
| **Eth** | 1 | 0 | 94 | 0 | 0 | 252 | 184 | 172 | 111 | 54 | 28 | 162 | 8 | 0 | 69 | 0 | 0 | 50 | 65 | 228 |
| **IPv4** | 0 | 0 | 1 | 17 | 34 | 185 | 131 | 202 | 240 | 87 | 224 | 0 | 0 | 252 | 201 | 86 | 20 | 235 | 0 | ... |

*(Eth row labels: Destination MAC Address, Source MAC Address)*

## VPN

| | 0 | | | | | | | | 9 | 10 | | | | | | | | | 19 | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **Pcap** | 161 | 178 | 195 | 212 | 0 | 2 | 0 | 4 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 255 | 255 |
| **Meta** | 0 | 0 | 0 | 101 | 85 | 45 | 101 | 91 | 0 | 0 | 111 | 11 | 0 | 0 | 0 | 56 | 0 | 0 | 0 | 56 |
| **IPv4** | 69 | 0 | 0 | 56 | 99 | 213 | 64 | 0 | 64 | 5 | 254 | 10 | 8 | 0 | 10 | 69 | 171 | 255 | 36 | |
| **UDP** | 146 | 214 | 13 | 150 | 0 | 36 | 120 | 43 | 0 | 1 | 0 | 8 | 33 | 18 | 164 | 66 | 52 | 167 | 9 | ... |

*(IPv4 row labels: Total Length, Frag.Off., Protocol)*

# Use Case #1: Detecting VPN vs Non-VPN Traffic

**Explanation**

VPN (without Ethernet):
IPv4 Protocol (6 or 17)

No-VPN (with Ethernet):
Source Mac Address
(Random)

True $B_{49} \leq 17$ False

33% 67%

$B_{43} \leq 1$ $B_{47} \leq 251$

1% 32% 66% 1%

Non VPN VPN Non VPN VPN

# Use Case #1: Detecting VPN vs Non-VPN Traffic

**Explanation**



VPN (without Ethernet): Fragment Offset

No-VPN (with Ethernet): Source Mac Address (Random)

True $B_{49} \leq 17$ False

33%   67%

$B_{43} \leq 1$   $B_{47} \leq 251$

1%   32%   66%   1%

Non VPN   VPN   Non VPN   VPN

# Use Case #1: Detecting VPN vs Non-VPN Traffic

**Explanation**

VPN (without Ethernet):
IP Total Length

No-VPN (with Ethernet):
Destination Mac Address
(Always 0)

True    $B_{49} \leq 17$    False

33%    67%

$B_{43} \leq 1$    $B_{47} \leq 251$

1%    32%    66%    1%

Non VPN    VPN    Non VPN    VPN

# Use Case #1: Detecting VPN vs Non-VPN Traffic

## Validation

- Validation dataset:
  - Tampering with packet headers from original PCAPs

| Validation Dataset | Avg. Precision | Avg. Recall | Avg. F1 |
|---|---|---|---|
| *Untampered* | 0.959 | 0.956 | 0.955 |
| *Tampered-43-47-49* | 0.959 | 0.956 | 0.955 |

# Use Case #1: Detecting VPN vs Non-VPN Traffic

**No VPN**

Byte 23: PCAP Link Type

No-VPN (With Ethernet): 1

| | 0 | | | | | | | | | | 10 | | | | | | | | 19 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Pcap | 161 | 178 | 195 | 2 | 0 | 2 | 0 | 4 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 255 | 255 |
| Meta | 0 | 0 | 0 | 1 | 85 | 65 | 10 | 69 | 0 | 5 | 80 | 24 | 0 | 0 | 0 | 64 | 0 | 0 | 0 | 64 |
| Ethernet | 1 | 0 | 94 | 0 | 0 | 252 | 184 | 172 | 111 | 54 | 28 | 162 | 8 | 0 | 69 | 0 | 0 | 50 | 65 | 228 |
| IPv4 | 0 | 0 | 1 | 17 | 34 | 185 | 131 | 202 | 240 | 87 | 224 | 0 | 0 | 252 | 201 | 86 | 20 | 235 | 0 | ... |

Destination MAC Address — Source MAC Address

Byte 23: PCAP Link Type

VPN (Without Ethernet): 101

**VPN**

| | 0 | | | | | | | | | | 10 | | | | | | | | 19 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Pcap | 161 | 178 | 195 | 2 | 0 | 2 | 0 | 4 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 255 | 255 |
| Meta | 0 | 0 | 0 | 101 | 85 | 45 | 101 | 91 | 0 | 0 | 111 | 11 | 0 | 0 | 0 | 56 | 0 | 0 | 0 | 56 |
| IPv4 | 69 | 0 | 0 | 56 | 199 | 213 | 64 | 0 | 64 | 17 | 35 | 254 | 10 | 8 | 0 | 10 | 69 | 171 | 255 | 36 |
| UDP | 146 | 214 | 13 | 150 | 0 | 36 | 120 | 43 | 0 | 1 | 0 | 8 | 33 | 18 | 164 | 66 | 52 | 167 | 9 | ... |

Total Length — Frag.Off. — Protocol

59

# Use Case #1: Detecting VPN vs Non-VPN Traffic

**Validation**

- Validation dataset:
  - Tampering with packet headers from original PCAPs

| Validation Dataset | Avg. Precision | Avg. Recall | Avg. F1 |
|---|---|---|---|
| *Untampered* | 0.959 | 0.956 | 0.955 |
| *Tampered-43-47-49* | 0.959 | 0.956 | 0.955 |
| *Tampered-32-to-63* | 0.889 | 0.867 | 0.856 |
| *Tampered-0-to-63* | 0.831 | 0.757 | 0.734 |
| *Tampered-0-to-127* | 0.753 | 0.555 | **0.398** |

# Use Case #1: Detecting VPN vs Non-VPN Traffic

## Validation

- Validation dataset:
  - Tampering with packet headers from original PCAPs

| Validation Dataset | Avg. Precision | Avg. Recall | Avg. F1 |
|---|---|---|---|
| *Untampered* | 0.959 | 0.956 | 0.955 |
| *Tampered-43-47-49* | 0.959 | 0.956 | 0.955 |
| *Tampered-32-to-63* | 0.889 | 0.867 | 0.856 |
| *Tampered-0-to-63* | 0.831 | 0.757 | 0.734 |
| *Tampered-0-to-127* | 0.753 | 0.555 | **0.398** |

**Takeaway: the model suffers from shortcut learning!**

# Use Case #2: Detecting Heartbleed Traffic

## Problem Setup

- Selected publications:
  - *Many papers that rely on the CIC-IDS-2017 dataset*
  - *"Toward Generating a New Intrusion Detection Dataset and Intrusion Traffic Characterization" — Sharafaldin et al., 2018*
- Proposal:
  - **Model**: Random Forest to classify traffic between benign traffic and 13 different attacks (e.g. PortScan, DDoS, Heartbleed)
  - **Features**: 78 pre-computed features, from flow statistics (e.g. flow duration, mean IAT)
  - **Dataset**: CIC-IDS-2017 [https://www.unb.ca/cic/datasets/ids-2017.html]
- Results:
  - Reported F1-score: 0.99
  - Reproduced F1-score: 0.99

# Use Case #2: Detecting Heartbleed Traffic



**Explanation**

Fidelity: 0.99
Top-3 pruning
6 nodes

True — Bwd Packet Length Max ≤ 12k — False

93%

7% — Heartbleed

Dest. Port ≤ 21.5

7% — FTP-Patator

86%

Dest. Port ≤ 22.5

7% — SSH-Patator

79% — …

# Use Case #2: Detecting Heartbleed Traffic

**Explanation**



True    Bwd Packet Length Max ≤ 12k    False

93%    7%

Dest. Port ≤ 21.5    Heartbleed

7%    86%

FTP-Patator    Dest. Port ≤ 22.5

7%    79%

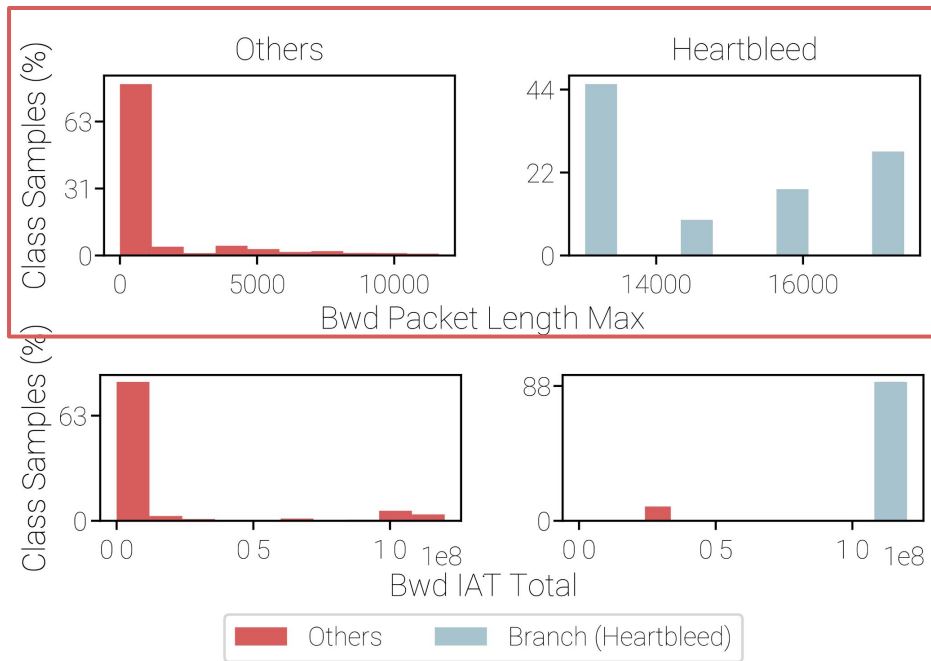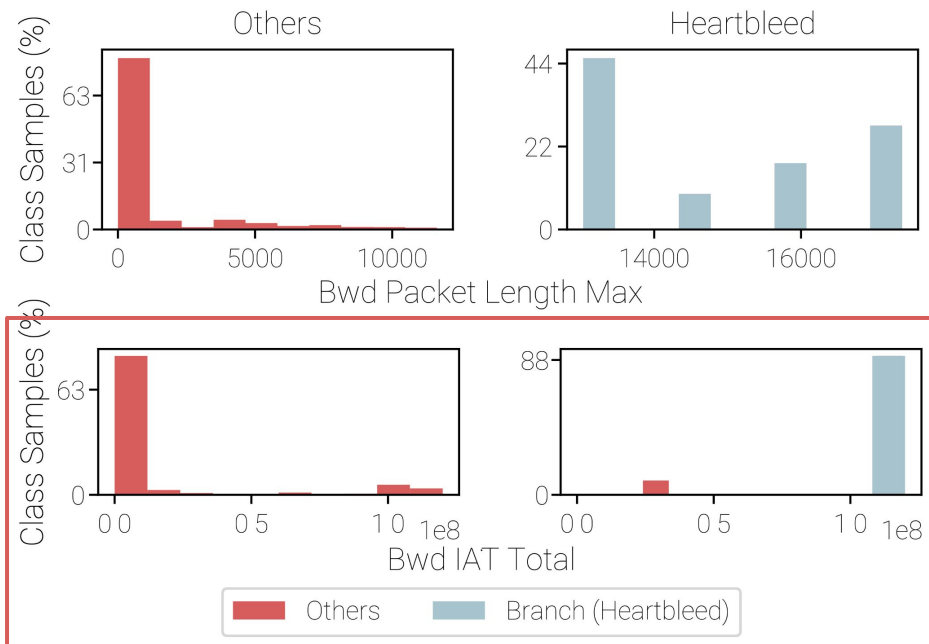SSH-Patator    ...

# Use Case #2: Detecting Heartbleed Traffic



**Explanation**

# Use Case #2: Detecting Heartbleed Traffic



**Explanation**

# Use Case #2: Detecting Heartbleed Traffic

- Heartbleed attack:
  - An attacker sends an HTTPS heartbeat message with a value in the size field bigger than the message
    - e.g., 16k bytes packet with 64k bytes size value
  - A vulnerable server responds with a message with the size equal to the value specified in the size field and reveals information stored locally in its memory
    - e.g. server returns 64k bytes (16k from packet and 48k from memory)

- In the CIC-IDS-2017 dataset:
  - HTTPS connection was never closed during the duration of the attack
    - Huge number of backward bytes and very high IAT in the flow!

# Use Case #2: Detecting Heartbleed Traffic

## Validation

- Validation dataset:
  - 1000 new heartbleed flows closing connection after every heartbeat
  - Backward bytes and IAT similar to benign traffic

| Class | Precision | Recall | F1 |
|---|---|---|---|
| *Heartbleed (i.i.d.)* | 1.000 | 1.000 | 1.000 |
| *Heartbleed (o.o.d)* | 0.000 | 0.000 | 0.000 |

# Use Case #2: Detecting Heartbleed Traffic

## Validation

- Validation dataset:
  - 1000 new heartbleed flows closing connection after every heartbeat
  - Backward bytes and IAT similar to benign traffic

| Class | Precision | Recall | F1 |
|---|---|---|---|
| *Heartbleed (i.i.d.)* | 1.000 | 1.000 | 1.000 |
| *Heartbleed (o.o.d)* | 0.000 | 0.000 | 0.000 |

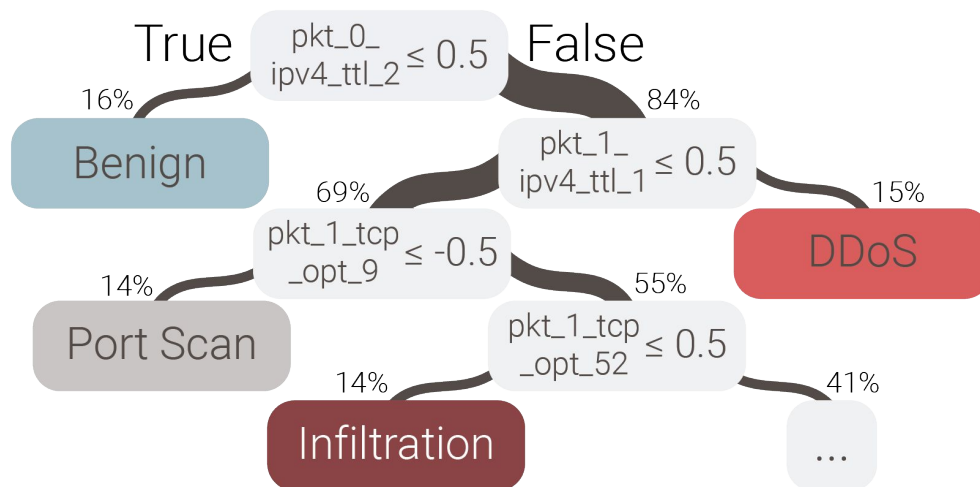**Takeaway: the model is overfitted to training data and fails to identify o.o.d. samples!**

# Use Case #3: Inferring Malicious Traffic for IDS

## Problem Setup

- Selected publications:
  - *"New Directions in Automated Traffic Analysis" — Holland et al., 2020*
- Proposal:
  - **Model**: nPrintML, an AutoML model for an Intrusion Detection System (IDS)
  - **Features**: 4,480 features with values -1, 0, or 1, each feature represents a bit of a set of pre-established protocol headers.
  - **Dataset**: CIC-IDS-2017 [https://www.unb.ca/cic/datasets/ids-2017.html]
- Results:
  - Reported F1-score: 0.99
  - Reproduced F1-score: 0.99
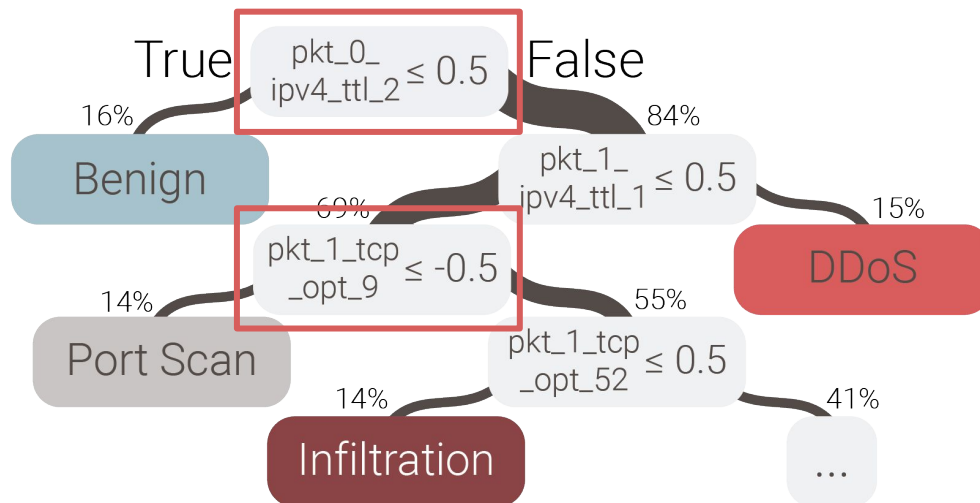
# Use Case #3: Inferring Malicious Traffic for IDS

**Explanation**



Fidelity: 0.99
Top-4 pruning
8 nodes

# Use Case #3: Inferring Malicious Traffic for IDS

**Explanation**

# Use Case #3: Inferring Malicious Traffic for IDS

## Validation

- Validation dataset:
  - Curated balanced dataset with 4,047 flows from real-world traffic in UCSB network
  - Used Suricata-IDS to generate flow labels

| Class | Precision | Recall | F1 |
|-------|-----------|--------|------|
| *Benign* | 0.653 | 0.806 | 0.722 |
| *DoS* | 0.000 | 0.000 | 0.000 |
| *Port Scan* | 0.120 | 0.143 | 0.130 |
| Average | 0.256 | 0.315 | 0.282 |

**Takeaway: the model suffers from spurious correlations in the training data!**
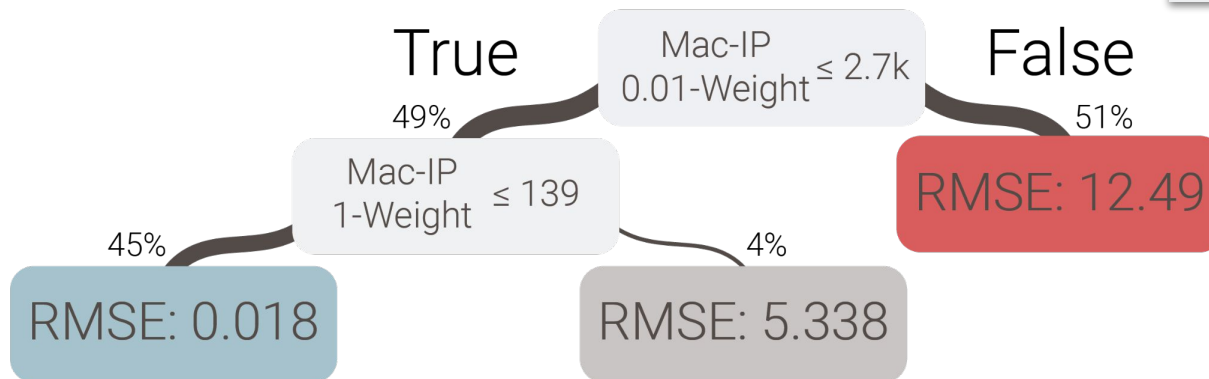
# Use Case #4: Anomaly Detection for Mirai Attacks

## Problem Setup

- Selected publications:
  - *"Kitsune: An Ensemble of Autoencoders for Online Network Intrusion Detection" — Mirsky et al., 2018*
- Proposal:
  - **Model**: Kitsune, an ensemble of neural networks, trained with unsupervised learning, for anomaly detection
  - **Features**: 110 features based on traffic statistics (*e.g.,* number of packets per **time window**).
  - **Dataset**: synthetic Mirai attack trace.
- Results:
  - Reported R-squared: 0.99
  - Reproduced R-squared: 0.99

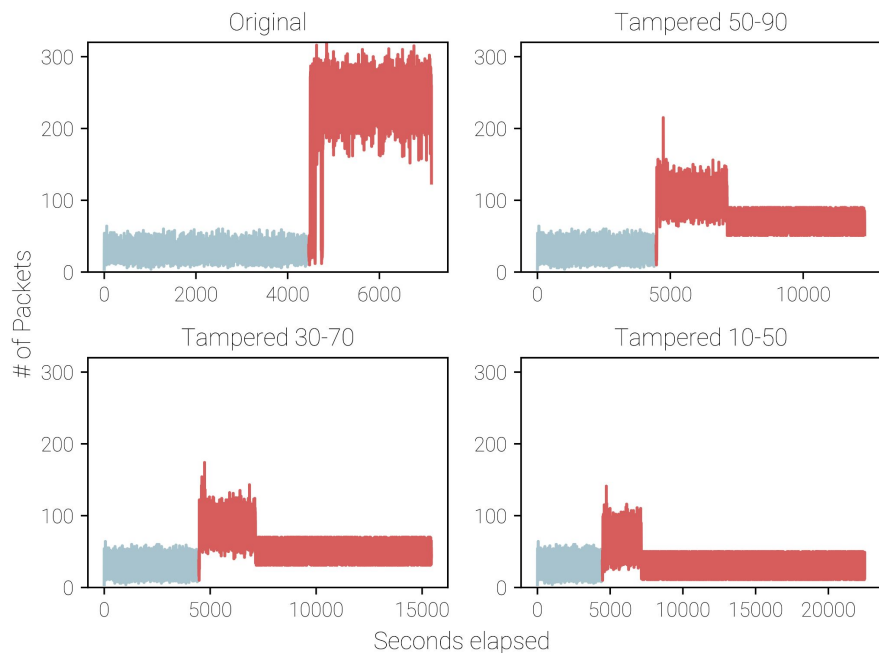# Use Case #4: Anomaly Detection for Mirai Attacks

**Explanation**

Fidelity: 0.99
Top-3 pruning
5 nodes

True

Mac-IP
0.01-Weight ≤ 2.7k

False

49%

51%

Mac-IP
1-Weight ≤ 139

RMSE: 12.49

45%

4%

RMSE: 0.018

RMSE: 5.338

# Use Case #4: Anomaly Detection for Mirai Attacks
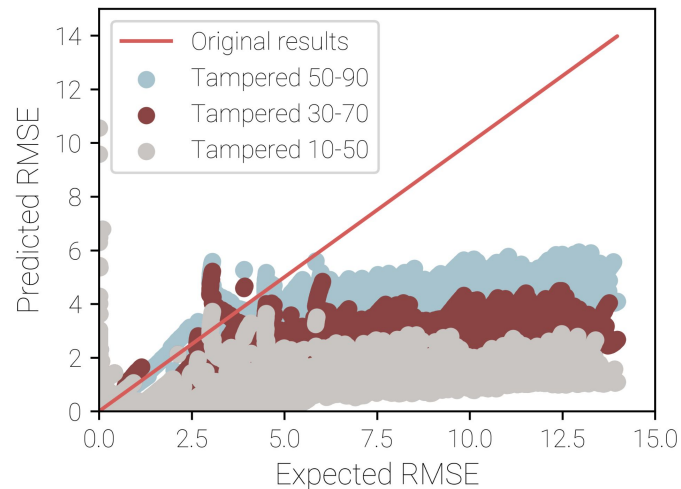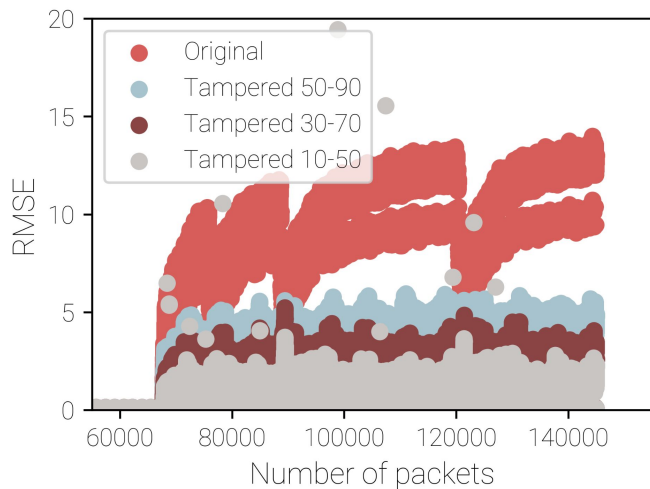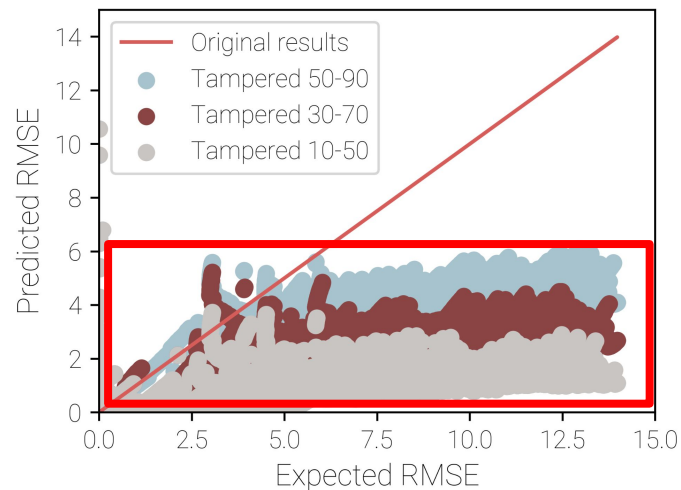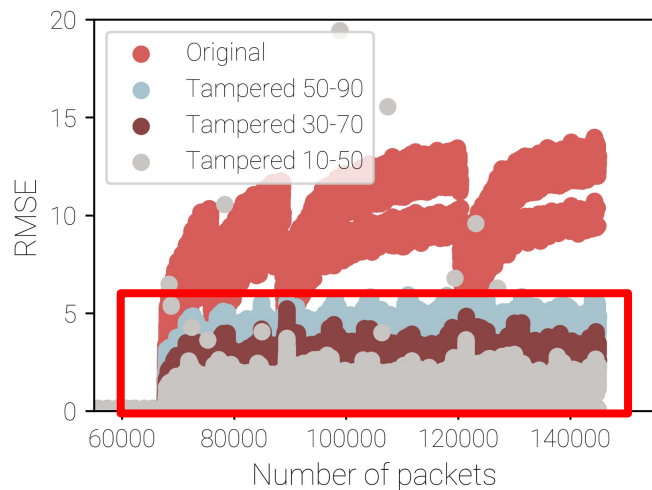
- Validation datasets:

# Use Case #4: Anomaly Detection for Mirai Attacks

**Validation**

# Use Case #4: Anomaly Detection for Mirai Attacks

**Validation**



**Takeaway: the model is overfitted to training data and fails to identify o.o.d. samples!**

# Other Use Cases

| Problem | Model(s) | Dataset(s) | Trustee Fidelity | Inductive Bias |
|---|---|---|---|---|
| Detect VPN traffic (Wang *et al.*, ISI'17) | 1-D CNN | ISCX VPN-nonVPN | 1.00 | Shortcut learning |
| Detect Heartbleed traffic (Sharafaldin *et al.*, ICISSP'18) | RFC | CIC-IDS-2017 | 0.99 | O.O.D. |
| Detect Malicious traffic (IDS) (Holland *et al.*, CCS'21) | nPrintML | CIC-IDS-2017 | 0.99 | Spurious Correlation |
| Anomaly Detection (Mirsky *et al.*, NDSS'18) | Kitsune | Mirai dataset | 0.99 | O.O.D |
| OS Fingerprinting (Holland *et al.*, CCS'21) | nPrintML | CIC-IDS-2017 | 0.99 | O.O.D |
| IoT Device Fingerprinting (Xiong *et al.*, HotNets'19) | Iisy | UNSW-IoT | 0.99 | Shortcut learning |
| Adaptive Bit-rate (Mao *et al.*, SIGCOMM'17) | Pensieve | HSDPA Norway | 0.99 | O.O.D |

# Other Use Cases

| Problem | Model(s) | Dataset(s) | Trustee Fidelity | Inductive Bias |
| --- | --- | --- | --- | --- |
| Detect VPN traffic (Wang et al., ISI'17) | 1-D CNN | ISCX VPN-nonVPN | 1.00 | Shortcut learning |
| Detect Heartbleed traffic (Sharafaldin et al., ICISSP'18) | RFC | CIC-IDS-2017 | 0.99 | O.O.D. |
| Detect Malicious traffic (IDS) (Holland et al., CCS'21) | nPrintML | CIC-IDS-2017 | 0.99 | Spurious Correlation |
| Anomaly Detection (Mirsky et al., NDSS'18) | Kitsune | Mirai dataset | 0.99 | O.O.D |
| OS Fingerprinting (Holland et al., CCS'21) | nPrintML | CIC-IDS-2017 | 0.99 | O.O.D |
| IoT Device Fingerprinting (Xiong et al., HotNets'19) | Iisy | UNSW-IoT | 0.99 | Shortcut learning |
| Adaptive Bit-rate (Mao et al., SIGCOMM'17) | Pensieve | HSDPA Norway | 0.99 | O.O.D |

# Trustee Python package



> 500 Downloads

# Conclusions

1. ML in high-stakes requires trust

2. Trustee improves trust!

3. Trustee can be used with any existing model

4. Trustee is ready to be used!
   - Just download our Python package

# Thank you!

**Arthur Jacobs**
*asjacobs@inf.ufrgs.br*

## trustee

*https://trusteeml.github.io*

Trustee Python package

- *https://pypi.org/project/trustee/*

Trustee Repository

- *https://github.com/TrusteeML/trustee*

Use Cases Repository

- *https://github.com/TrusteeML/emperor*

.Inf
INSTITUTO
DE INFORMÁTICA
UFRGS

NiKSUN

UFMS

# Backup

# But Network Practitioners remain skeptical…

## COMMUNICATIONS OF THE ACM

HOME | CURRENT ISSUE | NEWS | BLOGS | OPINION | RESEARCH | PRACTICE

Home / Magazine Archive / November 2021 (Vol. 64, No. 11) / November 2021 (Vol. 64, No. 11) / Full Text

CONTRIBUTED ARTICLES

## There Is No AI Without Data

By Christoph Gröger

Comments

VIEW AS: | | | | | SHARE: | | | | | | |

# But Network Practitioners remain skeptical…



COMMUNICATIONS
OF THE
ACM
HOME | CURRENT ISSUE | NEWS |

Home / Magazine Archive / November 2021 (Vol. 64, No. 11) / November 202

CONTRIBUTED ARTICLES

There Is No AI Without Data

By Christoph Gröger
Communications of the ACM, November 2021, Vol. 64 No. 11, Pages 98-108
10.1145/3448247
Comments

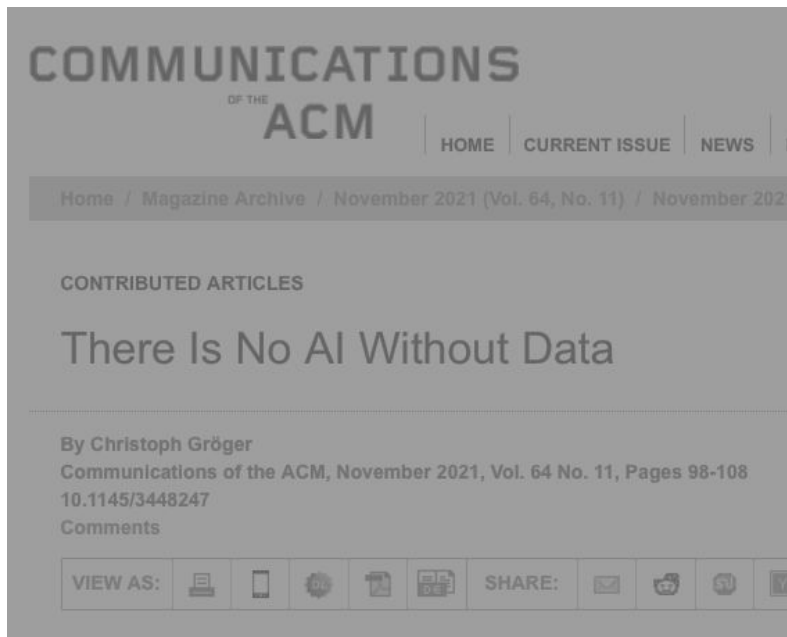VIEW AS:    SHARE:



AI & ML IN CYBERSECURITY – Why Algorithms
Are Dangerous

Category: Artificial Intelligence, Security Intelligence — Raffael Marty @ 10:28 am

WARNING
ALGORITHMS ARE DANGEROUS

# But Network Practitioners remain skeptical...

# Black-box nature of ML Models



Data → ML → Prediction

This issue is not unique to network security:

**eXplainable Artificial Intelligence (XAI)**

# Black-box nature of ML Models

# Black-box nature of ML Models

| | 0 | | | | | | | | | 9 | 10 | | | | | | | | | 19 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **Pcap Meta** | 161 | 178 | 195 | 212 | 0 | 2 | 0 | 4 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 255 | 255 |
| | 0 | 0 | 0 | 1 | 85 | 65 | 10 | 69 | 0 | 5 | 80 | 24 | 0 | 0 | 0 | 64 | 0 | 0 | 0 | 64 |
| **Eth** | 1 | 0 | 94 | 0 | 0 | 252 | 184 | 172 | 111 | 54 | 28 | 162 | 8 | 0 | 69 | 0 | 0 | 50 | 65 | 228 |
| **IPv4** | 0 | 0 | 1 | 17 | 34 | 185 | 131 | 202 | 240 | 87 | 224 | 0 | 0 | 252 | 201 | 86 | 20 | 235 | 0 | ... |

DNN

No VPN

VPN

# Existing approaches

| Method | Model Agnostic | High Fidelity | Domain-specific Pruning |
|---|---|---|---|
| Trepan | ✓ | — | — |
| *dtextract* | ✓ | — | — |
| VIPER | — | — | — |
| Metis | — | — | — |
| **trustee** | ✓ | ✓ | ✓ |

# Underspecification issues!

**Shortcut Learning**

Model 'learns' to classify based on feature values unrelated to classification problem.

**O.O.D. Samples**

Model overfits to training dataset distribution, and fails when faced with out of distribution (o.o.d) samples.

**Spurious Correlations**

Model relies on spurious correlations between features to achieve perfect accuracy.

# Underspecification issues!

**Shortcut Learning**

Model 'learns' to classify based on feature values unrelated to classification problem.

**O.O.D. Samples**

Model overfits to training dataset distribution, and fails when faced with out of distribution (o.o.d) samples.

**Spurious Correlations**

Model relies on spurious correlations between features to achieve perfect accuracy.

⚠ SPOILER ALERT! ⚠

These issues usually come from the same underlying problem: **bad data**.

# Conclusions

1. Do not blindly trust AI/ML!

2. Make reproducibility artifacts available!

3. Collect your own data!
   - Ask for your university IT staff for help.

# Thank you!

**Arthur Jacobs**
*asjacobs@inf.ufrgs.br*

## trustee

*https://trusteeml.github.io*

Trustee Python package

- *https://pypi.org/project/trustee/*

Trustee Repository

- *https://github.com/TrusteeML/trustee*

Use Cases Repository

- *https://github.com/TrusteeML/emperor*

.inf
INSTITUTO
DE INFORMÁTICA
UFRGS

NIKSUN

UFMS