# The relationship between
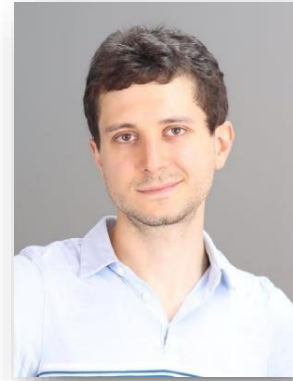# **Machine Learning & Cybersecurity**

Giovanni Apruzzese, PhD

TU Delft – May 3rd, 2022

UNIVERSITÄT
LIECHTENSTEIN

Giovanni Apruzzese, *PhD*
*giovanni.apruzzese@uni.li*

# whoami: Dr. **Giovanni Apruzzese** 🇮🇹

o **Background:**

- Did my academic studies (BSc, MSc, PhD) at University of Modena, Italy.
  – Supervisor: Prof. Michele Colajanni
- In 2019, spent 6 months at Dartmouth College, USA.
  – Supervisor: Prov. VS Subrahmanian
- Joined the University of Liechtenstein in July 2020 as a PostDoc Researcher.
  – Supervisor: Prof. Pavel Laskov
- Met Prof. Mauro Conti in 2019, with whom I have been collaborating since 2020.

o **Interests:**

- Cybersecurity, machine learning, and any network-related topic (+🎮)
- I like talking, researching and teaching – in a "pragmatic" way ☺

o **Contact information**:

- Work Email: giovanni.apruzzese@uni.li
- Feel free to contact me if you have any questions.
  – I reply fast, and will happily do so!

UNIVERSITÄT
LIECHTENSTEIN

# What I do

<div style="text-align:center">

## Machine Learning + Cybersecurity

</div>

o Applying ML to *provide security* of a given information system
  - E.g.: using ML to detect network intrusions

o *Attacking / Defending* ML applications
  - E.g.: evading a ML model that detects phishing websites

o Using machine learning *offensively* against another system
  - E.g.: artificially generating "fake" images

BONUS

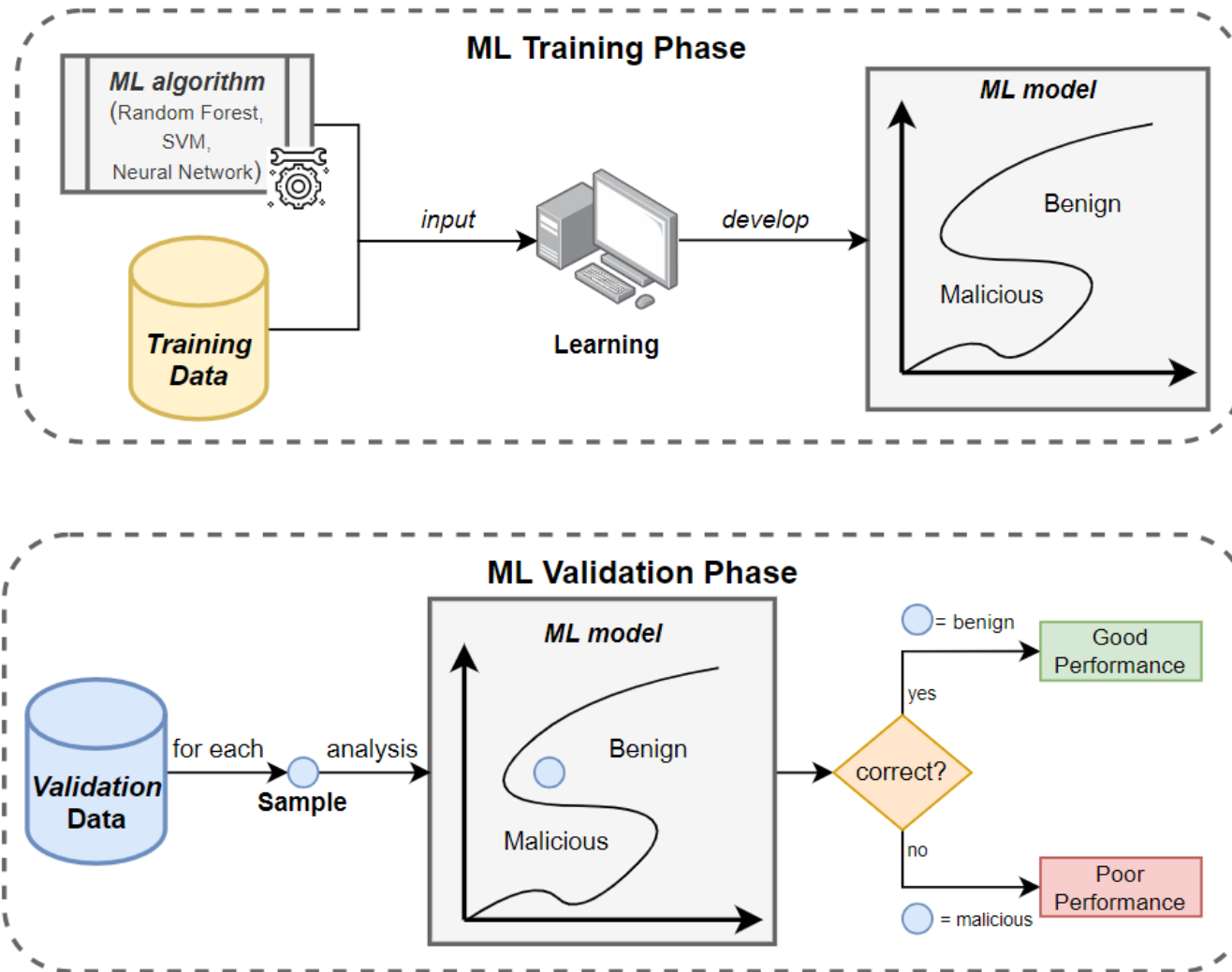o Using ML to attack a security system based on ML

UNIVERSITÄT
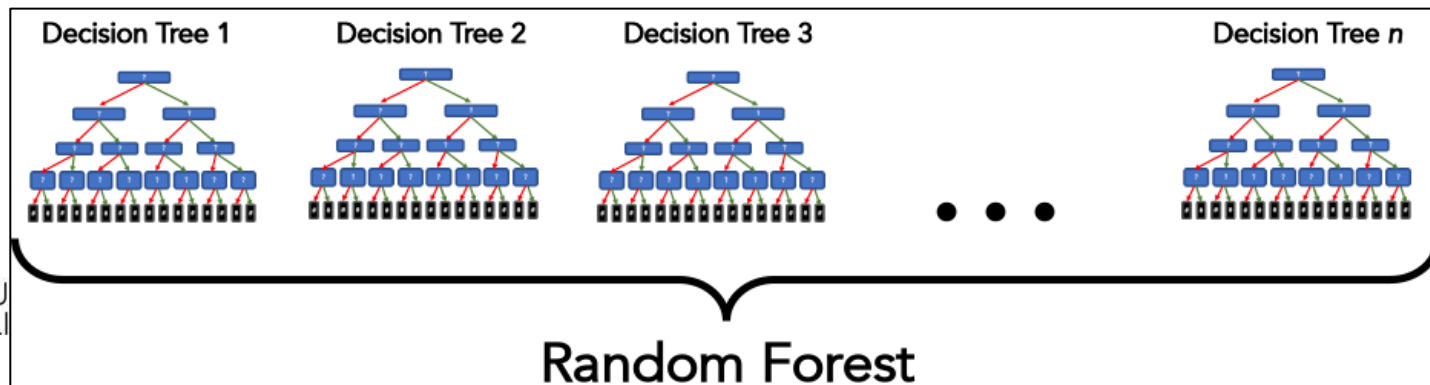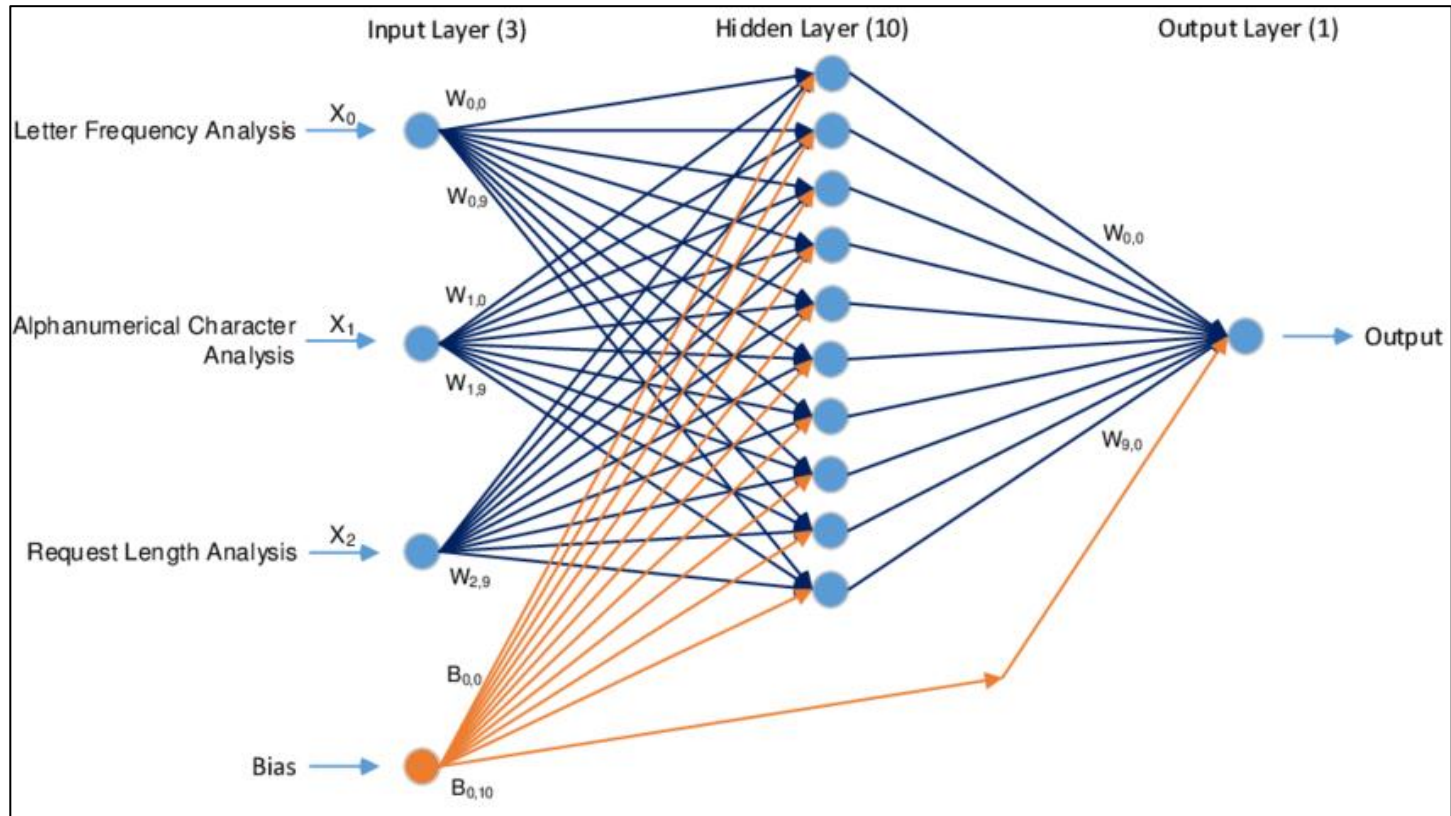LIECHTENSTEIN

# Outline of Today

o **Fundamentals of Machine Learning and Cybersecurity**

o **Using unlabelled data for Machine Learning in Cyberthreat Detection**
  - Ref: Giovanni Apruzzese, Luca Pajola, and Mauro Conti. "The Cross-evaluation of Machine Learning-based Network Intrusion Detection Systems." IEEE Transactions on Network and Service Management (2022).

o **Improving Machine Learning in Network Intrusion Detection**
  - Ref: Giovanni Apruzzese, Pavel Laskov, Aliya Tastemirova. "SoK: The Impact of Unlabelled Data for Cyberthreat Detection." IEEE European Symposium on Security and Privacy (2022).

o **The security of Machine Learning-based Phishing Website Detectors**
  - Ref: Giovanni Apruzzese, Mauro Conti, Ying Yuan. "SpacePhish: The Evasion-space of Adversarial Attacks against Phishing Website Detectors". TBD

o **Adversarial Attacks against Humans and Machine Learning**
  - Ref: Johannes Schneider, Giovanni Apruzzese. "Concept-based Adversarial Attacks: Tricking Humans and Classifiers alike." IEEE Symposium on Security and Privacy – Deep Learning and Security Workshop (2022)

UNIVERSITÄT
LIECHTENSTEIN

# Fundamentals of Machine Learning and Cybersecurity

# Machine Learning workflow: Training and Testing

Giovanni Apruzzese, *PhD*
*giovanni.apruzzese@uni.li*

# Question: do you think that training ML models is difficult?

# Question: do you think that <u>training</u> ML models is difficult?

```python
#train the classifier (rf_clf) using the training_data (train[features]) with corresponding labels (y)
print("Training...")
rf_clf.fit(train[features],y)
print("Done")
```

UNIVERSITÄT
LIECHTENSTEIN

Giovanni Apruzzese, *PhD*
*giovanni.apruzzese@uni.li*

# Question: do you think that <u>training</u> ML models is difficult?

**PROBLEMS (data)**

```
#train the classifier (rf_clf) using the training_data (train[features]) with corresponding labels (y)
print("Training...")
rf_clf.fit(train[features],y)
print("Done")
```

**PROBLEMS (tuning)**

UNIVERSITÄT
LIECHTENSTEIN

Giovanni Apruzzese, *PhD*
*giovanni.apruzzese@uni.li*

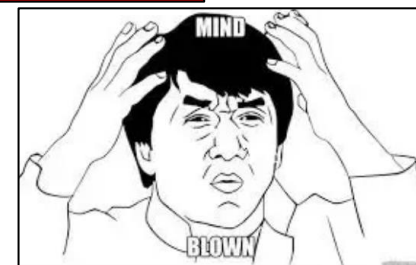Question: do you think that <u>training</u> ML models is difficult?

**PROBLEMS (data)**

```
#train the classifier (rf_clf) using the training_data (train[features]) with corresponding labels (y)
print("Training...")
rf_clf.fit(train[features],y)
print("Done")
```

**PROBLEMS (tuning)**

Of course, you're always free to go, learn and improve the *fit* function:
https://github.com/scikit-learn/scikit-learn/blob/baf828ca1/sklearn/ensemble/_forest.py#L297

# Common issues of ML in Cybersecurity

o  Applying Machine Learning requires *data* to train a ML model

o  Depending on the "problem" solved by such model, the data may require *labels*

o  **Obtaining (any) data has a <u>cost</u>, and labelled data is (very) *expensive***

o  Machine Learning models are ultimately just a component within a system

o  **Such ML models *can* be targeted by "Adversarial Attacks"**

o  Such strategies ultimately aim to compromise the functionality of the ML model.

o  The cybersecurity domain implicitly assumes the presence of attackers.

o  Attackers are *human beings,* and hence operate with a *cost/benefit* mindset

o  **Such considerations must be made when analyzing the security of (any) IT system**

UNIVERSITÄT
LIECHTENSTEIN

11

# Unlabelled data for Machine Learning in Cyberthreat Detection
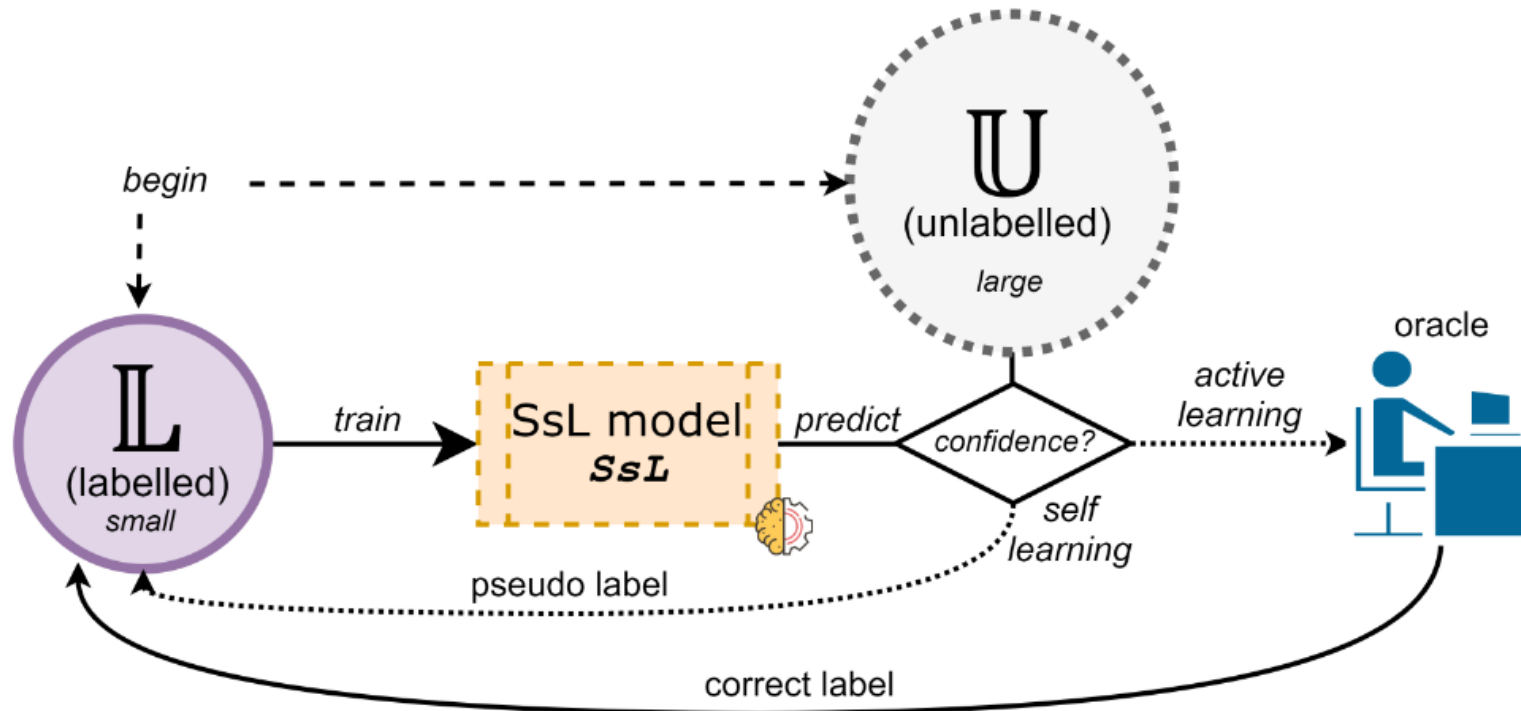
# Semisupervised Learning

o    Labelled data is expensive, but *unlabelled* data is cheap(er).

→ Why not using unlabelled data to improve the proficiency of ML models?

Mixing *labelled* with *unlabelled* data is a ML approach denoted as
**"Semisupervised Learning" (SsL)**

# Semisupervised Learning

o   Labelled data is expensive, but *unlabelled* data is cheap(er).

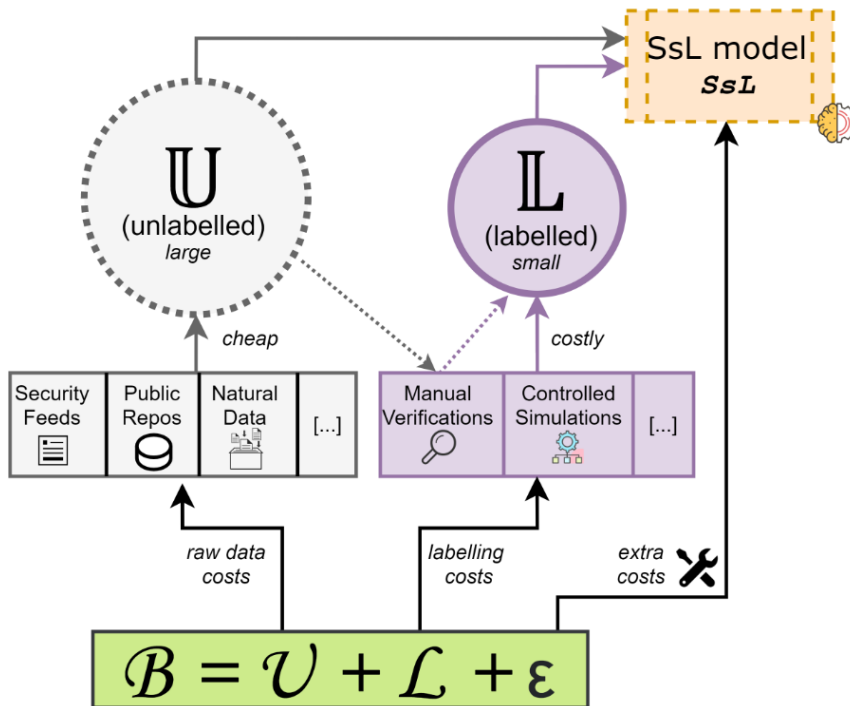→ Why not using unlabelled data to improve the proficiency of ML models?

> Mixing *labelled* with *unlabelled* data is a ML approach denoted as
> **"Semisupervised Learning" (SsL)**



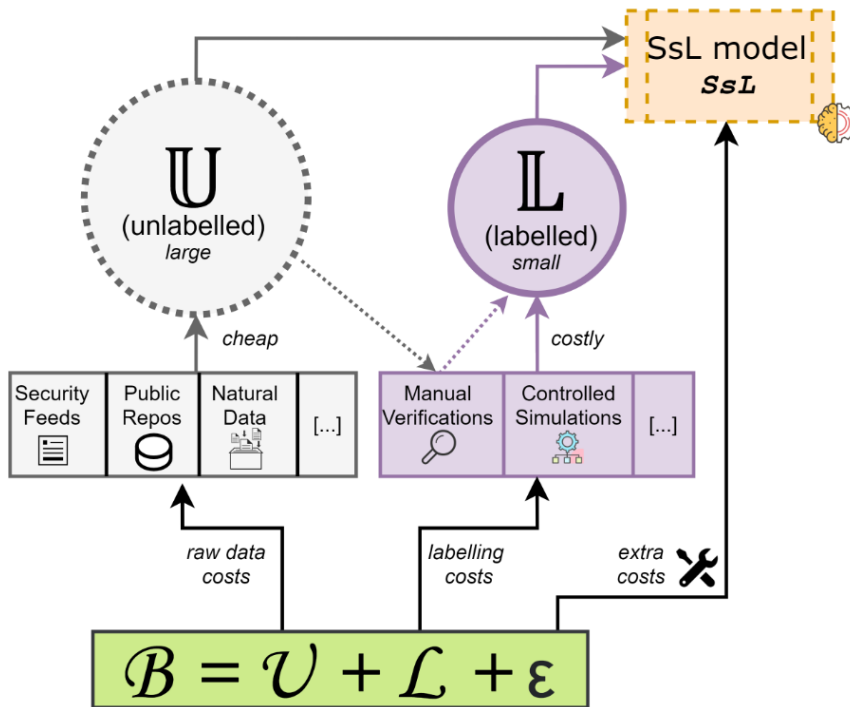Examples of SsL: *active learning* and *self learning* (e.g., *pseudo labelling*)

# Goal of Semisupervised Learning

o Developing SsL models is cheaper than "supervised learning" (SL) models, **but it is not free.**
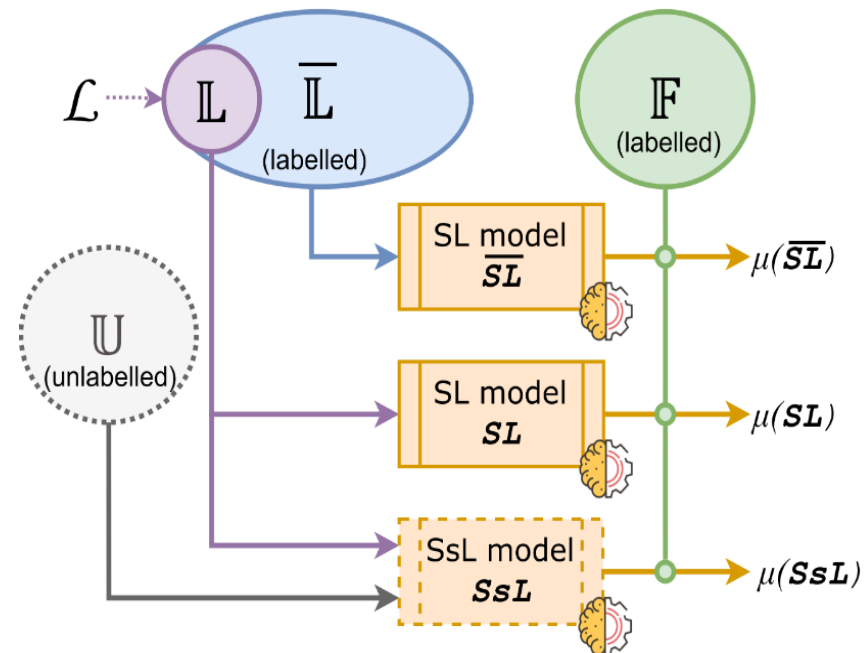
# Goal of Semisupervised Learning

o Developing SsL models is cheaper than "supervised learning" (SL) models, **but it is not free.**

o A SsL model should achieve a *performance superior* than a SL model that uses the *same labelling budget*
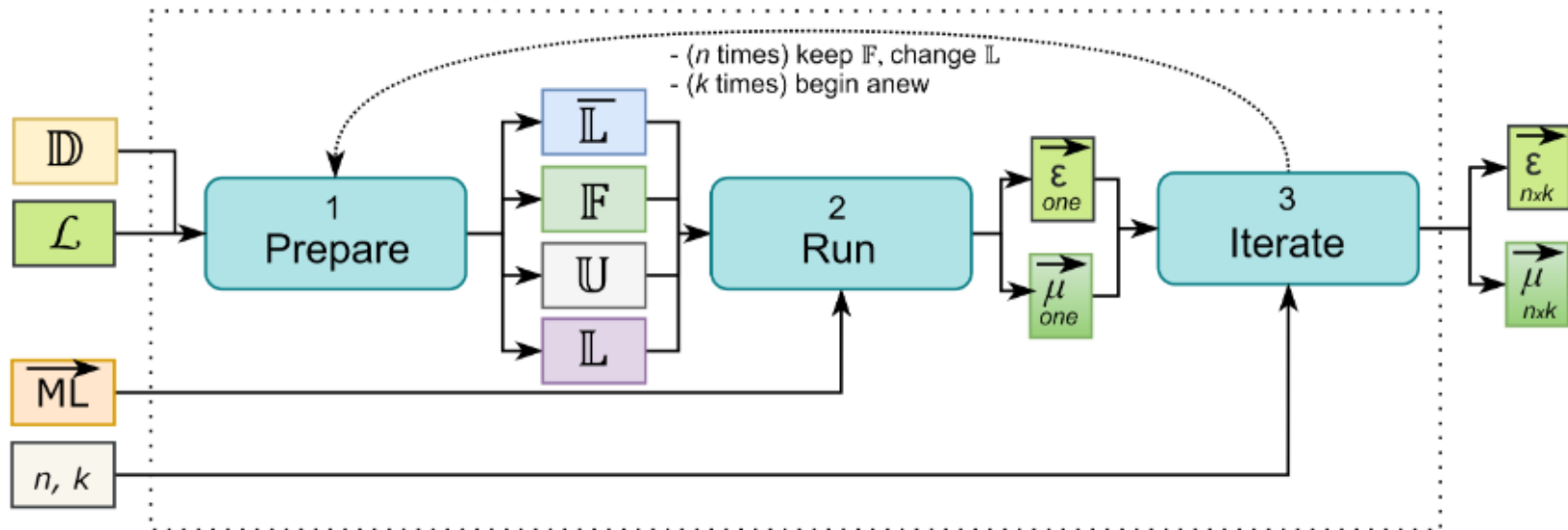
# Problem: nobody cares

The current state-of-the-art does not allow to determine whether SsL methods applied in Cyberthreat Detection are <u>truly</u> beneficial

| Task | Paper (1st Author) | Year | Lower Bound | Ablation Study | Upper Bound | Stat. Sign. | Transparency | | Repr. | Dataset |
|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | *Labels* | *Balance* | | |
| Network Intrusion Detection | Li [93] | 2007 | ✓ | ✓ | ✗ | ✗ | ✓ | ✓ | ◐ | NSL-KDD |
| | Long [94] | 2008 | ✓ | ✓ | ✗ | ◐ | ✓ | ✗ | ◐ | NSL-KDD |
| | Görnitz [95] | 2009 | ✓ | ✓ | ◐ | ✗ | ✓ | ✓ | ✗ | Private |
| | Seliya [96] | 2010 | ✓ | ✓ | ✗ | ✗ | ✓ | ✓ | ◐ | NSL-KDD |
| | Symons [97] | 2012 | ✗ | ✓ | ✓ | ◐ | ✓ | ✗ | ✗ | Kyoto2006 |
| | Wagh [98] | 2014 | ✗ | ✗ | ✗ | ✗ | ✓ | ✓ | ◐ | NSL-KDD |
| | Noorbehbahani [35] | 2015 | ✗ | ◐ | ✓ | ✗ | ✓ | ✓ | ◐ | NSL-KDD, Custom |
| | Ashfaq [99] | 2017 | ✗ | ◐ | ✗ | ✗ | ✓ | ✗ | ◐ | NSL-KDD |
| | Qiu [67] | 2017 | ✗ | ◐ | ✓ | ✗ | ✓ | ✓ | ✗ | Custom |
| | McElwee [100] | 2017 | ✗ | ◐ | ✓ | ✗ | ✓ | ✗ | ◐ | NSL-KDD |
| | Kumari [68] | 2017 | ✓ | ◐ | ✗ | ✗ | ✓ | ✗ | ◐ | NSL-KDD |
| | Yang [101] | 2018 | ◐ | ✓ | ✓ | ✗ | ✓ | ✗ | ✗ | NSL-KDD, AWID |
| | Gao [102] | 2018 | ✓ | ◐ | ✗ | ✗ | ✓ | ✗ | ✗ | NSL-KDD |
| | Shi [103] | 2018 | ◐ | ◐ | ✗ | ✗ | ✓ | ✗ | ✗ | NSL-KDD |
| | Yao [36] | 2019 | ◐ | ◐ | ✓ | ✗ | ✓ | ✓ | ◐ | NSL-KDD |
| | Yuan [104] | 2019 | ✗ | ◐ | ✗ | ◐ | ✓ | ✓ | ◐ | NSL-KDD |
| | Zhang [65] | 2020 | ◐ | ✗ | ✓ | ◐ | ✓ | ✗ | ◐ | NSL-KDD |
| | Hara [105] | 2020 | ✗ | ◐ | ✓ | ✗ | ✗ | ✗ | ✗ | NSL-KDD |
| | Ravi [106] | 2020 | ✓ | ✗ | ✗ | ✗ | ✓ | ✗ | ✗ | NSL-KDD |
| | Gao [107] | 2020 | ✗ | ✓ | ✓ | ✓ | ✓ | ✓ | ✗ | NSL-KDD |
| | Li [108] | 2020 | ✗ | ◐ | ✓ | ✓ | ✓ | ✗ | ◐ | NSL-KDD, Private |
| | Zhang [70] | 2021 | ◐ | ◐ | ✗ | ◐ | ✗ | ✓ | ◐ | CICIDS2017, CTU13 |
| | Liang [109] | 2021 | ✓ | ◐ | ✓ | ◐ | ✓ | ✓ | ◐ | NSL-KDD |
| Phishing Detection | Gyawali [110] | 2011 | ✗ | ✓ | ✓ | ✗ | ✓ | ✓ | ◐ | Private |
| | Zhao [111] | 2013 | ✓ | ✓ | ✓ | ✓ | ✗ | ✓ | ✓* | DetMalURL |
| | Gabriel [15] | 2017 | ◐ | ◐ | ✗ | ✗ | ✗ | ✗ | ◐ | Private |
| | Yang [112] | 2017 | ✓ | ◐ | ✗ | ✗ | ✓ | ✓ | ◐ | Private |
| | Bhattacharjee [113] | 2017 | ✗ | ✓ | ✗ | ◐ | ✗ | ✗ | ◐ | Private |
| | Li [55] | 2017 | ✓ | ✓ | ✓ | ◐ | ✓ | ✓ | ✗ | Custom |
| Malware Detection | Moskovitch [114] | 2008 | ✗ | ✓ | ✗ | ◐ | ✓ | ✓ | ✗ | Custom |
| | Santos [115] | 2011 | ✗ | ✗ | ✓ | ✗ | ✓ | ✓ | ◐ | Custom |
| | Nissim [116] | 2012 | ✗ | ◐ | ✓ | ◐ | ✗ | ✗ | ✗ | Private |
| | Zhao [117] | 2012 | ✗ | ✗ | ✗ | ✗ | ✓ | ✓ | ◐ | Private |
| | Nissim [118] | 2014 | ✓ | ✓ | ✗ | ◐ | ✓ | ✓ | ✗ | Custom |
| | Zhang [119] | 2015 | ◐ | ◐ | ✗ | ✗ | ✓ | ✓ | ✗ | Private |
| | Nissim [120] | 2016 | ✗ | ✓ | ✗ | ◐ | ✓ | ✓ | ◐ | Custom |
| | Ni [121] | 2016 | ✓ | ✓ | ✗ | ◐ | ✓ | ✓ | ◐ | Private |
| | Chen [122] | 2017 | ✓ | ✓ | ✗ | ◐ | ✗ | ✗ | ◐ | Private |
| | Rashidi [66] | 2017 | ✗ | ✓ | ✓ | ◐ | ✓ | ✓ | ✗ | Drebin |
| | Fu [123] | 2019 | ✓ | ✓ | ✗ | ✗ | ✓ | ✗ | ◐ | Private |
| | Irofti [124] | 2019 | ◐ | ◐ | ✗ | ◐ | ✗ | ✗ | ✓ | DREBIN, EMBER |
| | Pendlebury [86] | 2019 | ✗ | ✗ | ✓ | ◐ | ✓ | ✓ | ✓ | AndroZoo |
| | Sharmeen [125] | 2020 | ✓ | ◐ | ✗ | ◐ | ✓ | ✓ | ◐ | Drebin, AndroZoo |
| | Chen [126] | 2020 | ◐ | ◐ | ✓ | ✗ | ✓ | ✓ | ◐ | MCC |
| | Koza [11] | 2020 | ✓ | ◐ | ✓ | ◐ | ✓ | ✗ | ✓ | Private |
| | Noorbehbahani [13] | 2020 | ✓ | ✗ | ✗ | ◐ | ✓ | ✓ | ✗ | AndMal17 |
| | Li [127] | 2021 | ✗ | ◐ | ✗ | ◐ | ✓ | ✗ | ◐ | FalDroid, DREBIN, Genome |
| | Liang [109] | 2021 | ✓ | ◐ | ✓ | ◐ | ✓ | ✓ | ◐ | Custom |

# Solution: CEF-SsL

o   SsL is intriguing, but its "pragmatic" benefits are still unknown

o   Identifying (and quantifying) such benefits requires adopting a rigorous workflow

→ CEF-SsL: Cybersecurity Evaluation Framework for Semisupervised Learning

# (re)Evaluation

o Massive evaluation on 9 existing datasets for 3 cyberthreat detection tasks:

- Network Intrusion Detection (NID)

- Phishing Website Detection (PWD)

- Malware Detection (MD)

Labels range between 100 and 2400

Results
(F1-score)

| CTD | NID | | | PWD | | | MD | | |
|---|---|---|---|---|---|---|---|---|---|
| Method | CTU13 | UNB15 | IDS17 | Mend | UCI | $\delta$Phish | DREBIN | Ember | AndMal |
| | | | | | | | | | |
| $\pi SsL$ | 0.588 | 0.437 | 0.820 | 0.850 | 0.884 | 0.778 | 0.474 | 0.647 | 0.900 |
| $\widehat{\pi} SsL$ | 0.584 | 0.435 | 0.818 | 0.849 | 0.883 | 0.777 | 0.470 | 0.641 | 0.890 |
| $\alpha SsL_l$ | **0.693** | 0.582 | **0.897** | **0.863** | **0.903** | 0.770 | **0.546** | **0.687** | **0.924** |
| $\alpha SsL_o$ | 0.637 | 0.577 | 0.874 | 0.855 | 0.891 | 0.745 | 0.497 | 0.673 | 0.916 |
| $\alpha SsL_h$ | 0.510 | 0.436 | 0.786 | 0.834 | 0.851 | 0.714 | 0.423 | 0.598 | 0.892 |
| $\alpha^{\pi} SsL_l$ | 0.664 | 0.533 | 0.853 | 0.861 | 0.901 | 0.767 | 0.529 | 0.654 | 0.901 |
| $\alpha^{\pi} SsL_o$ | 0.633 | **0.595** | 0.857 | 0.854 | 0.890 | 0.745 | 0.489 | 0.647 | 0.895 |
| $\alpha^{\pi} SsL_h$ | 0.486 | 0.427 | 0.744 | 0.833 | 0.851 | 0.711 | 0.410 | 0.579 | 0.865 |

UNIVERSITÄT
LIECHTENSTEIN

# (re)Evaluation

o  Massive evaluation on 9 existing datasets for 3 cyberthreat detection tasks:

- Network Intrusion Detection (NID)
- Phishing Website Detection (PWD)
- Malware Detection (MD)

Labels range between 100 and 2400

Results
(F1-score)

| CTD | NID | | | PWD | | | MD | | |
|---|---|---|---|---|---|---|---|---|---|
| Method | CTU13 | UNB15 | IDS17 | Mend | UCI | $\delta$Phish | DREBIN | Ember | AndMal |
| $\overline{SL}$ | 0.979 | 0.942 | 0.989 | 0.958 | 0.974 | 0.958 | 0.907 | 0.970 | 0.986 |
| $SL$ | 0.611 | 0.447 | 0.878 | 0.852 | 0.884 | 0.780 | 0.480 | 0.667 | 0.910 |
| $\underline{SsL}$ | 0.613 | 0.447 | 0.879 | 0.852 | 0.886 | **0.778** | 0.486 | 0.662 | 0.910 |
| $\pi\,SsL$ | 0.588 | 0.437 | 0.820 | 0.850 | 0.884 | 0.778 | 0.474 | 0.647 | 0.900 |
| $\widehat{\pi}\,SsL$ | 0.584 | 0.435 | 0.818 | 0.849 | 0.883 | 0.777 | 0.470 | 0.641 | 0.890 |
| $\alpha\,SsL_l$ | **0.693** | 0.582 | **0.897** | **0.863** | **0.903** | 0.770 | **0.546** | **0.687** | **0.924** |
| $\alpha\,SsL_o$ | 0.637 | 0.577 | 0.874 | 0.855 | 0.891 | 0.745 | 0.497 | 0.673 | 0.916 |
| $\alpha\,SsL_h$ | 0.510 | 0.436 | 0.786 | 0.834 | 0.851 | 0.714 | 0.423 | 0.598 | 0.892 |
| $\alpha^{\pi}SsL_l$ | 0.664 | 0.533 | 0.853 | 0.861 | 0.901 | 0.767 | 0.529 | 0.654 | 0.901 |
| $\alpha^{\pi}SsL_o$ | 0.633 | **0.595** | 0.857 | 0.854 | 0.890 | 0.745 | 0.489 | 0.647 | 0.895 |
| $\alpha^{\pi}SsL_h$ | 0.486 | 0.427 | 0.744 | 0.833 | 0.851 | 0.711 | 0.410 | 0.579 | 0.865 |

## Is SsL truly advantageous?

UNIVERSITÄT
LIECHTENSTEIN

# (re)Evaluation

o Massive evaluation on 9 existing datasets for 3 cyberthreat detection tasks:

- Network Intrusion Detection (NID)
- Phishing Website Detection (PWD)
- Malware Detection (MD)

*Labels range between 100 and 2400*

**Results (F1-score)**

| CTD | NID | | | PWD | | | MD | | |
|---|---|---|---|---|---|---|---|---|---|
| Method | CTU13 | UNB15 | IDS17 | Mend | UCI | $\delta$Phish | DREBIN | Ember | AndMal |
| $\overline{SL}$ | 0.979 | 0.942 | 0.989 | 0.958 | 0.974 | 0.958 | 0.907 | 0.970 | 0.986 |
| $SL$ | 0.611 | 0.447 | 0.878 | 0.852 | 0.884 | 0.780 | 0.480 | 0.667 | 0.910 |
| $\underline{SsL}$ | 0.613 | 0.447 | 0.879 | 0.852 | 0.886 | **0.778** | 0.486 | 0.662 | 0.910 |
| $\pi\, SsL$ | 0.588 | 0.437 | 0.820 | 0.850 | 0.884 | 0.778 | 0.474 | 0.647 | 0.900 |
| $\widehat{\pi}\, SsL$ | 0.584 | 0.435 | 0.818 | 0.849 | 0.883 | 0.777 | 0.470 | 0.641 | 0.890 |
| $\alpha\, SsL_l$ | **0.693** | 0.582 | **0.897** | **0.863** | **0.903** | 0.770 | **0.546** | **0.687** | **0.924** |
| $\alpha\, SsL_o$ | 0.637 | 0.577 | 0.874 | 0.855 | 0.891 | 0.745 | 0.497 | 0.673 | 0.916 |
| $\alpha\, SsL_h$ | 0.510 | 0.436 | 0.786 | 0.834 | 0.851 | 0.714 | 0.423 | 0.598 | 0.892 |
| $\alpha^{\pi}SsL_l$ | 0.664 | 0.533 | 0.853 | 0.861 | 0.901 | 0.767 | 0.529 | 0.654 | 0.901 |
| $\alpha^{\pi}SsL_o$ | 0.633 | **0.595** | 0.857 | 0.854 | 0.890 | 0.745 | 0.489 | 0.647 | 0.895 |
| $\alpha^{\pi}SsL_h$ | 0.486 | 0.427 | 0.744 | 0.833 | 0.851 | 0.711 | 0.410 | 0.579 | 0.865 |

**Statistical Validation**

| Dataset | PopSize | Best 'pure' pseudo-labelling | | | Best active learning | | |
|---|---|---|---|---|---|---|---|
| | | Method | $p$-value | $z$-value | Method | $p$-value | $z$-value |
| CTU13 | 396 | $\underline{SsL}$ | **0.873** | 0.159 | $\alpha\, SsL_l$ | $< 0.001$ | 4.310 |
| UNB15 | 1104 | $\underline{SsL}$ | **0.964** | $-0.044$ | $\alpha^{\pi}SsL_o$ | $< 0.001$ | 15.98 |
| IDS17 | 540 | $\underline{SsL}$ | **0.932** | 0.085 | $\alpha\, SsL_l$ | **0.978** | $-0.027$ |
| UCI | 1200 | $\underline{SsL}$ | **0.473** | 0.717 | $\alpha\, SsL_l$ | $< 0.001$ | 7.386 |
| Mend. | 1200 | $\underline{SsL}$ | **0.713** | 0.368 | $\alpha\, SsL_l$ | $< 0.001$ | 6.757 |
| $\delta$Phish | 1200 | $\underline{SsL}$ | **0.554** | $-0.590$ | $\alpha\, SsL_l$ | 0.002 | $-3.113$ |
| Drebin | 1200 | $\underline{SsL}$ | **0.310** | 1.015 | $\alpha\, SsL_l$ | $< 0.001$ | 11.78 |
| Ember | 1200 | $\underline{SsL}$ | **0.603** | $-0.512$ | $\alpha\, SsL_l$ | $< 0.001$ | 3.407 |
| AndMal | 1200 | $\underline{SsL}$ | **0.712** | $-0.370$ | $\alpha\, SsL_l$ | $< 0.001$ | 12.01 |

UNIVERSITÄT LIECHTENSTEIN

# Improving Machine Learning in Network Intrusion Detection
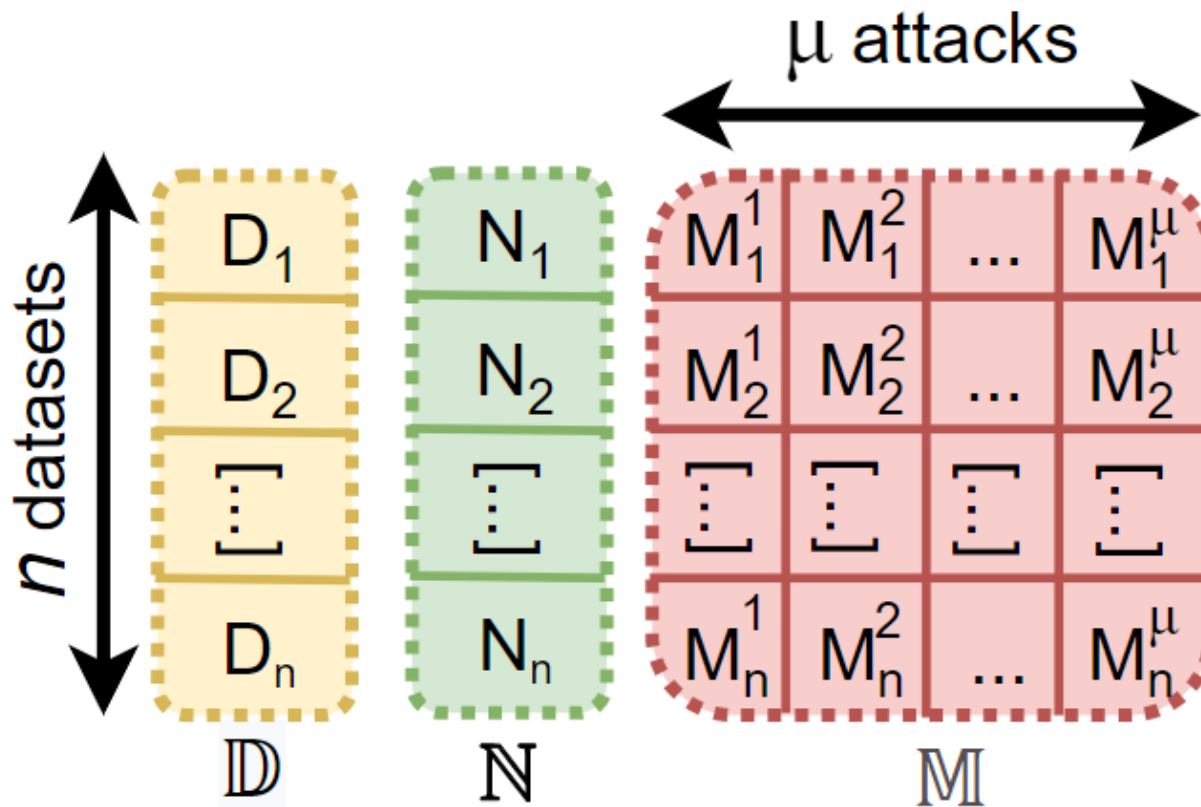
# Problem Statement

o   Most organizations adopt Network Intrusion Detection Systems (NIDS)

o   Such NIDS are starting to actively leverage Machine Learning (ML-NIDS)



o   **Problem: every network environment is *unique***

- This characteristic conflicts with the "iid" assumption, which is fundamental for ML
    - iid: independent and identically distributed random variables
- Training data must be collected <u>*from and for*</u> each network monitored by ML-NIDS

UNIVERSITÄT
LIECHTENSTEIN

# Intuition: Cross-evaluation of ML-NIDS

o  It is true that every network is unique…

o  … but (some) *malicious* events are malicious *everywhere* and *everytime*

→  Why not using <u>malicious</u> samples taken from *different* networks to "augment" the data used for training/testing my ML models?

# Intuition: Cross-evaluation of ML-NIDS

o   It is true that every network is unique…

o   … but (some) *malicious* events are malicious *everywhere* and *everytime*

→   Why not using <u>malicious</u> samples taken from *different* networks to "augment" the data used for training/testing my ML models?

# Solution: XeNIDS

o  The idea is intriguing, but *applying* it in practice is difficult

- Adversarial poisoning
- Incompatible networks
- False-sense of security
- Performance Decrease

→ XeNIDS – framework for the Cross-evaluation of Network Intrusion Detection Systems

# Solution: XeNIDS

o  The idea is intriguing, but *applying* it in practice is difficult

- Adversarial poisoning
- Incompatible networks
- False-sense of security
- Performance Decrease

→ XeNIDS – framework for the Cross-evaluation of Network Intrusion Detection Systems

# Evaluation

o  Massive evaluation of XeNIDS on 8 datasets

| Scenario | Dataset | #Samples | #Attacks | #Features | F1-score |
|---|---|---|---|---|---|
| *Heter.* | CTU13 | 20.7M | 5 | 14 | 99.1% [17] |
| | NB15 | 2.5M | 9 | 48 | 98.7% [18] |
| | IDS18 | 3.1M | 14 | 80 | 96.2% [42] |
| | DDOS19 | 70M | 18 | 80 | 99.0% [19] |
| *Uniform* | UF-BotIoT | 600K | 4 | 12 | 97.0% [21] |
| | UF-NB15 | 1.6M | 9 | 12 | 85.0% [21] |
| | UF-IDS18 | 8.3M | 14 | 12 | 83.0% [21] |
| | UF-ToNIoT | 1.4M | 9 | 12 | 100.0% [21] |

o  XeNIDS can be used for:

- Assessing how an existing ML-NIDS fares against "unknown" attacks; and
- Increasing the robustness of an existing ML-NIDS against such "unknown" attacks

UNIVERSITÄT
LIECHTENSTEIN

# Results

o   Baseline performance against unknown attacks (F1-score)

| Heterogeneous scenario | | | | Uniform scenario | | | |
|---|---|---|---|---|---|---|---|
| **Dataset** | **Botnet** | **DoS** | **Other** | **Dataset** | **Botnet** | **DoS** | **Other** |
| CTU13 | 80.0 | 38.1 | 49.7 | UF-BotIoT | 47.8 | 69.0 | 76.8 |
| NB15 | 65.8 | 40.7 | 75.2 | UF-NB15 | 72.2 | 52.3 | 64.1 |
| IDS18 | 54.9 | 49.4 | 76.1 | UF-IDS18 | 68.2 | 81.0 | 63.3 |
| DDOS19 | 54.4 | 99.5 | 83.1 | UF-ToNIoT | 82.1 | 89.3 | 85.1 |

o   Enhanced performance against unknown attacks (F1-score)

| Heterogeneous scenario | | | | Uniform scenario | | | |
|---|---|---|---|---|---|---|---|
| **Dataset** | **Botnet** | **DoS** | **Other** | **Dataset** | **Botnet** | **DoS** | **Other** |
| CTU13 | 98.8 | 99.9 | 98.9 | UF-BotIoT | 99.7 | 99.9 | 99.2 |
| NB15 | 97.1 | 99.9 | 99.1 | UF-NB15 | 88.9 | 99.2 | 98.7 |
| IDS18 | 98.5 | 99.7 | 97.7 | UF-IDS18 | 99.9 | 99.4 | 97.8 |
| DDOS19 | 99.9 | 99.9 | 98.6 | UF-ToNIoT | 99.7 | 99.9 | 99.9 |

UNIVERSITÄT
LIECHTENSTEIN

Giovanni Apruzzese, *PhD*
*anni.apruzzese@uni.li*

# Results

Advantage?
It's all <u>free</u>

o  Baseline performance against unknown attacks (F1-score)

| Heterogeneous scenario | | | | Uniform scenario | | | |
|---|---|---|---|---|---|---|---|
| **Dataset** | **Botnet** | **DoS** | **Other** | **Dataset** | **Botnet** | **DoS** | **Other** |
| CTU13 | 80.0 | 38.1 | 49.7 | UF-BotIoT | 47.8 | 69.0 | 76.8 |
| NB15 | 65.8 | 40.7 | 75.2 | UF-NB15 | 72.2 | 52.3 | 64.1 |
| IDS18 | 54.9 | 49.4 | 76.1 | UF-IDS18 | 68.2 | 81.0 | 63.3 |
| DDOS19 | 54.4 | 99.5 | 83.1 | UF-ToNIoT | 82.1 | 89.3 | 85.1 |

o  Enhanced performance against unknown attacks (F1-score)

| Heterogeneous scenario | | | | Uniform scenario | | | |
|---|---|---|---|---|---|---|---|
| **Dataset** | **Botnet** | **DoS** | **Other** | **Dataset** | **Botnet** | **DoS** | **Other** |
| CTU13 | 98.8 | 99.9 | 98.9 | UF-BotIoT | 99.7 | 99.9 | 99.2 |
| NB15 | 97.1 | 99.9 | 99.1 | UF-NB15 | 88.9 | 99.2 | 98.7 |
| IDS18 | 98.5 | 99.7 | 97.7 | UF-IDS18 | 99.9 | 99.4 | 97.8 |
| DDOS19 | 99.9 | 99.9 | 98.6 | UF-ToNIoT | 99.7 | 99.9 | 99.9 |

UNIVERSITÄT
LIECHTENSTE

**CAUTION!**
Always analyze the results!

# The security of Machine Learning-based Phishing Website Detectors

UNIVERSITÄT
LIECHTENSTEIN

Giovanni Apruzzese, *PhD*
giovanni.apruzzese@uni.li

# Current Landscape of Phishing

o   Phishing attacks are continuously increasing

o   Current detection methods still rely on *blacklists* of malicious URLs

  • These detection techniques can be <u>evaded easily</u> by "squatting" phishing websites!



Image source: https://www.tessian.com/blog/phishing-statistics-2020/

# Current Landscape of Phishing – Countermeasures

o Countering such simple (but effective) strategies can be done via *data-driven* methods

**Website**   **Phishing Website Detector**

Preprocessing → Analysis → output → Benign / Phishing

UNIVERSITÄT LIECHTENSTEIN

# Current Landscape of Phishing – Countermeasures (ML)

o  Countering such simple (but effective) strategies can be done via *data-driven* methods



o  Such methods (obviously ☺) include (also) Machine Learning techniques:



o  Machine Learning-based Phishing Website Detectors (ML-PWD) are very effective! [1]

   •  Even popular products and web-browsers (e.g., Google Chrome) use them! [2]

UNIVERSITÄT
LIECHTENSTEIN

[1]: Tian, Ke, et al. "Needle in a haystack: Tracking down elite phishing domains in the wild." Internet Measurement Conference 2018.
[2]: El Kouari, Oumaima, Hafssa Benaboud, and Saiida Lazaar. "Using machine learning to deal with Phishing and Spam Detection: An overview." Proceedings of the 3rd International Conference on Networking, Information Systems & Security. 2020.

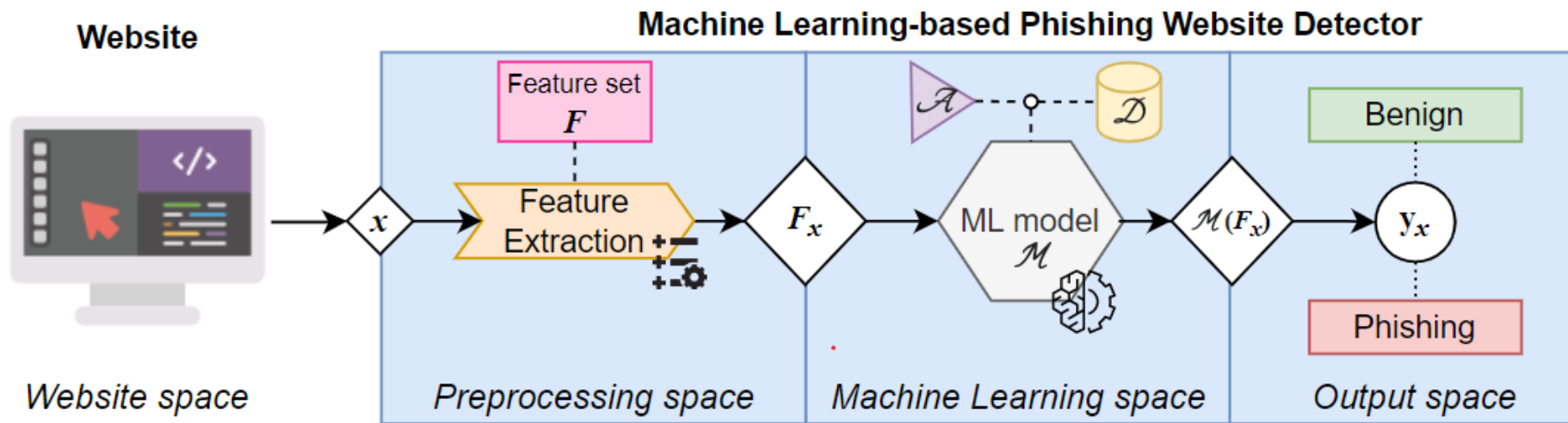# Problem Statement: Adversarial Attacks against ML

o    ML-PWD are good but…

o    …the detection of ML methods *can* be bypassed via (adversarial) *evasion* attacks!

o    Adversarial Attacks exploit a perturbation, $\varepsilon$, that induces a ML model, $\mathcal{M}$, to misclassify a given input, $x$, by producing an incorrect output ($y_x^{\varepsilon}$ instead of $y_x$)

$$\text{find } \varepsilon \text{ s.t. } \mathcal{M}(F_x) = y_x^{\varepsilon} \neq y_x$$

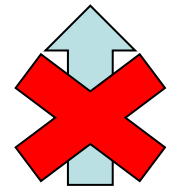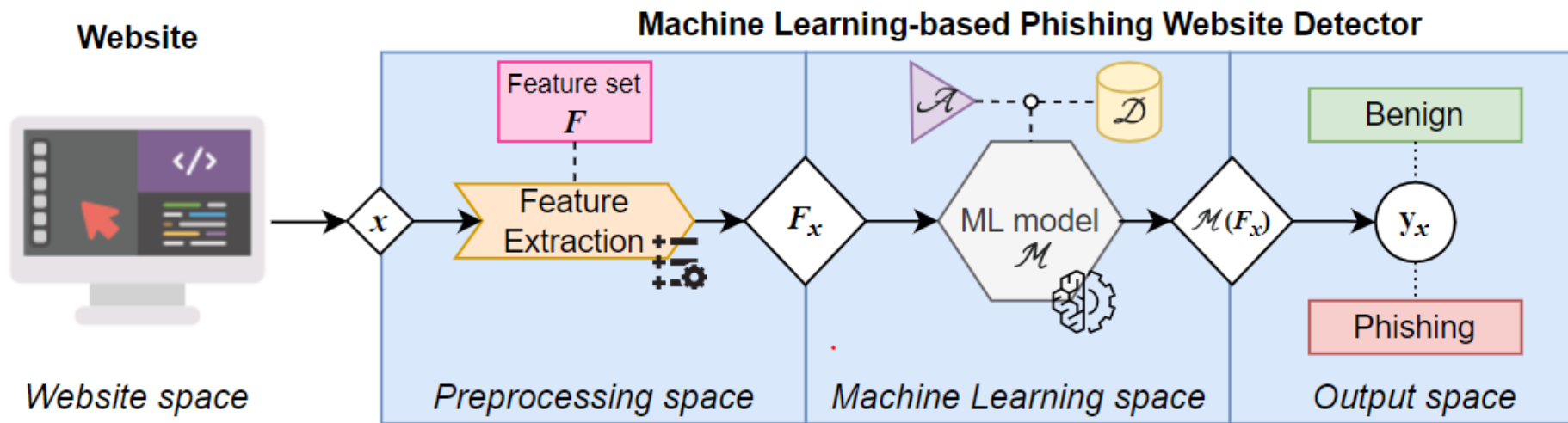# Problem Statement: Adversarial Attacks against ML-PWD

o ML-PWD are good but…

o …the detection of ML methods *can* be bypassed via (adversarial) *evasion* attacks!

o Adversarial Attacks exploit a perturbation, $\varepsilon$, that induces a ML model, $\mathcal{M}$, to misclassify a given input, $x$, by producing an incorrect output ($y_x^\varepsilon$ instead of $y_x$)

$$\text{find } \varepsilon \text{ s.t. } \mathcal{M}(F_x) = y_x^\varepsilon \neq y_x$$

o In the context of a ML-PWD, such $\varepsilon$ can be introduced in three 'spaces':
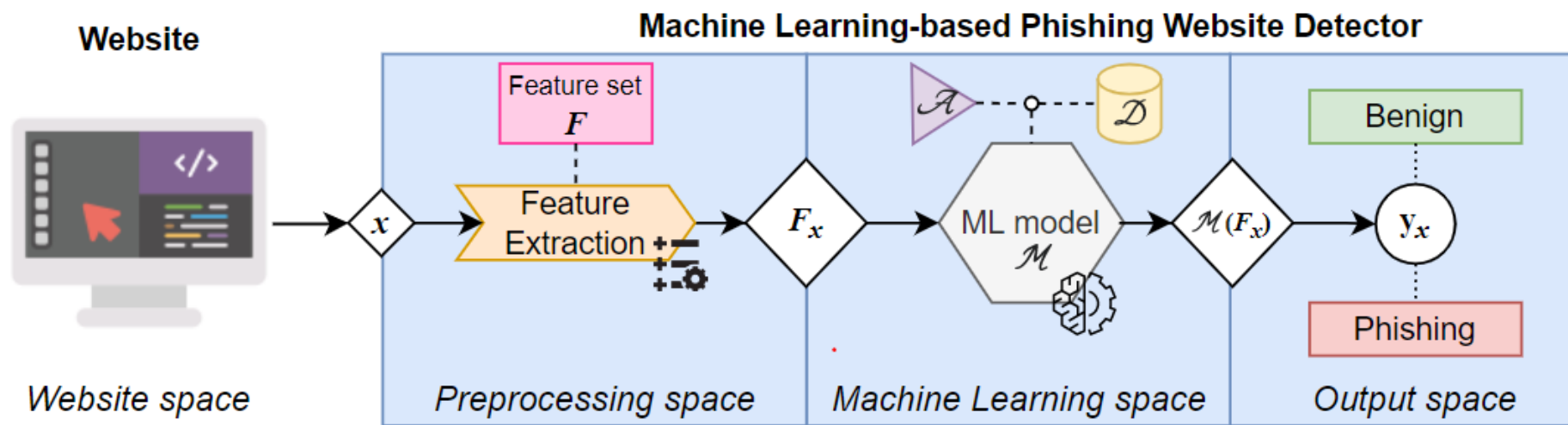


UNIVERSITÄT
LIECHTENSTEIN

36

# Problem Statement: Adversarial Attacks against ML-PWD

o   ML-PWD are good but…

o   …the detection of ML methods *can* be bypassed via (adversarial) *evasion* attacks!

o   Adversarial Attacks exploit a perturbation, $\varepsilon$, that induces a ML model, $\mathcal{M}$, to misclassify a given input, $x$, by producing an incorrect output ($y_x^\varepsilon$ instead of $y_x$)

$$\text{find } \varepsilon \text{ s.t. } \mathcal{M}(F_x) = y_x^\varepsilon \neq y_x$$

o   In the context of a ML-PWD, such $\varepsilon$ can be introduced in three 'spaces':

# Problem Statement: Adversarial Attacks against ML-PWD

o   ML-PWD are good but…

o   …the detection of ML methods *can* be bypassed via (adversarial) *evasion* attacks!

o   Adversarial Attacks exploit a perturbation, $\varepsilon$, that induces a ML model, $\mathcal{M}$, to misclassify a given input, $x$, by producing an incorrect output ($y_x^{\varepsilon}$ instead of $y_x$)

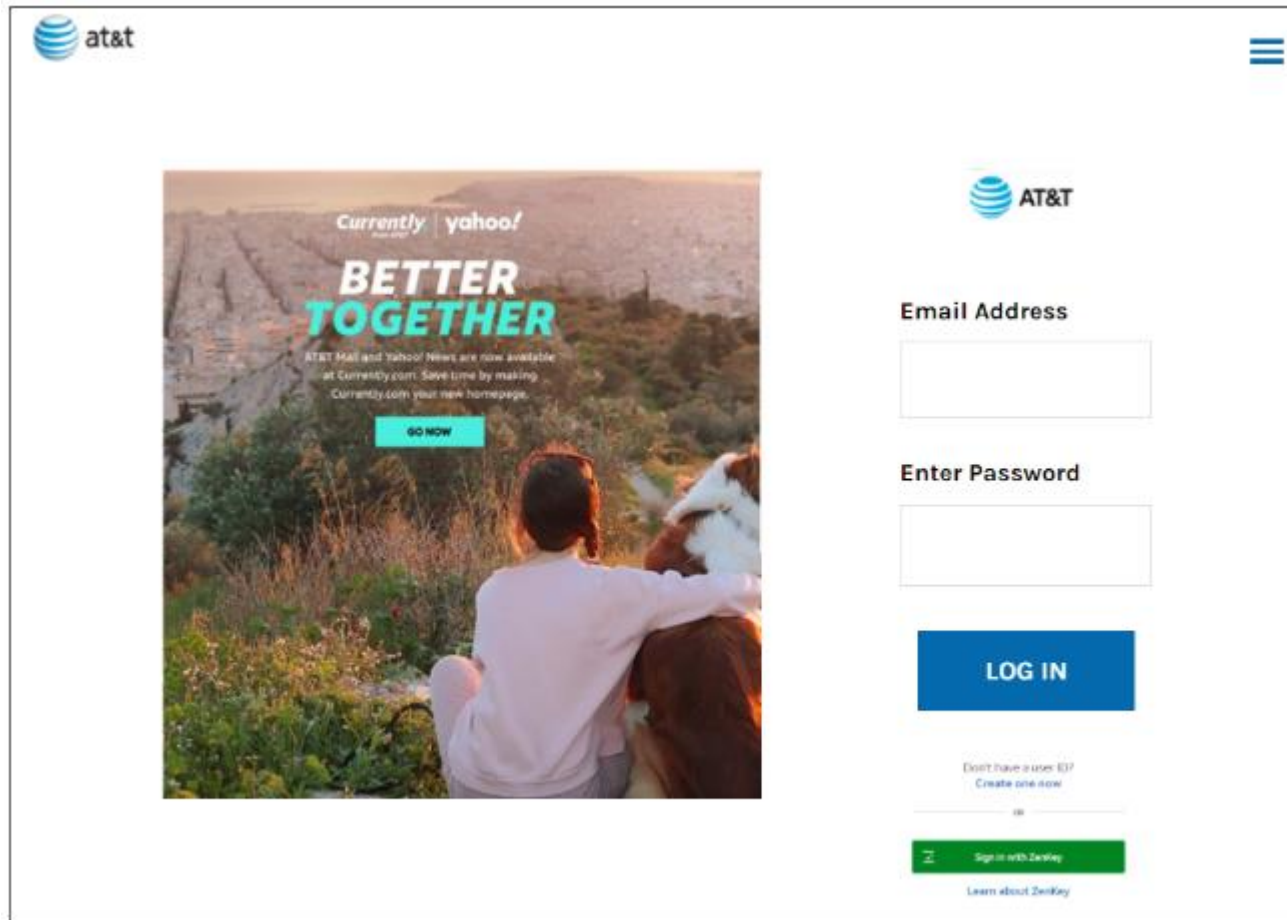$$\text{find } \varepsilon \text{ s.t. } \mathcal{M}(F_x) = y_x^{\varepsilon} \neq y_x$$

o   In the context of a ML-PWD, such $\varepsilon$ can be introduced in three 'spaces':



UNIVERSITÄT
LIECHTENSTEIN

38

Giovanni Apruzzese, *PhD*
giovanni.apruzzese@uni.li

# Problem Statement: Adversarial Attacks against ML-PWD

o   ML-PWD are good but…

o   …the detection of ML methods *can* be bypassed via (adversarial) *evasion* attacks!

o   Adversarial Attacks exploit a perturbation, $\varepsilon$, that induces a ML model, $\mathcal{M}$, to misclassify a given input, $x$, by producing an incorrect output ($y_x^{\varepsilon}$ instead of $y_x$)

$$\text{find } \varepsilon \text{ s.t. } \mathcal{M}(F_x) = y_x^{\varepsilon} \neq y_x$$

o   In the context of a ML-PWD, such $\varepsilon$ can be introduced in three 'spaces':



UNIVERSITÄT
LIECHTENSTEIN

39

# Problem Statement: Adversarial Attacks against ML-PWD

o   ML-PWD are good but…

o   …the detection of ML methods *can* be bypassed via (adversarial) *evasion* attacks!

o   Adversarial Attacks exploit a perturbation, $\varepsilon$, that induces a ML model, $\mathcal{M}$, to misclassify a given input, $x$, by producing an incorrect output ($y_x^\varepsilon$ instead of $y_x$)

$$\text{find } \varepsilon \text{ s.t. } \mathcal{M}(F_x) = y_x^\varepsilon \neq y_x$$

o   In the context of a ML-PWD, such $\varepsilon$ can be introduced in three 'spaces':

# Problem Statement: Adversarial Attacks against ML-PWD

o  ML-PWD are good but…

o  …the detection of ML methods *can* be bypassed via (adversarial) *evasion* attacks!

o  Adversarial Attacks exploit a perturbation, $\varepsilon$, that induces a ML model, $\mathcal{M}$, to misclassify a given input, $x$, by producing an incorrect output ($y_x^{\varepsilon}$ instead of $y_x$)

$$\text{find } \varepsilon \text{ s.t. } \mathcal{M}(F_x) = y_x^{\varepsilon} \neq y_x$$

o  In the context of a ML-PWD, such $\varepsilon$ can be introduced in three "spaces":



Question: Which "space" do you think an *attacker* is **most likely** to use?

# Website-space Perturbations – In practice (original example)

**Figure 4: An exemplary (and true) Phishing website, whose URL is https://www.63y3hfh-fj39f30-f30if0f-f392.weebly.com/.**

# Website-space Perturbations – In practice (changing the URL)

https://www.63y3hfh-fj39f30-f30if0f-f392.weebly.com/  ➡  https://bit.ly/3MZHjt7"

UNIVERSITÄT
LIECHTENSTEIN

43

# Website-space Perturbations – In practice (changing the HTML)

Giovanni Apruzzese, *PhD*
*giovanni.apruzzese@uni.li*

# Website-space Perturbations – In practice (change URL + HTML)



https://www.63y3hfh-fj39f30-f30if0f-f392.weebly.com/

https://bit.ly/3MZHjt7"

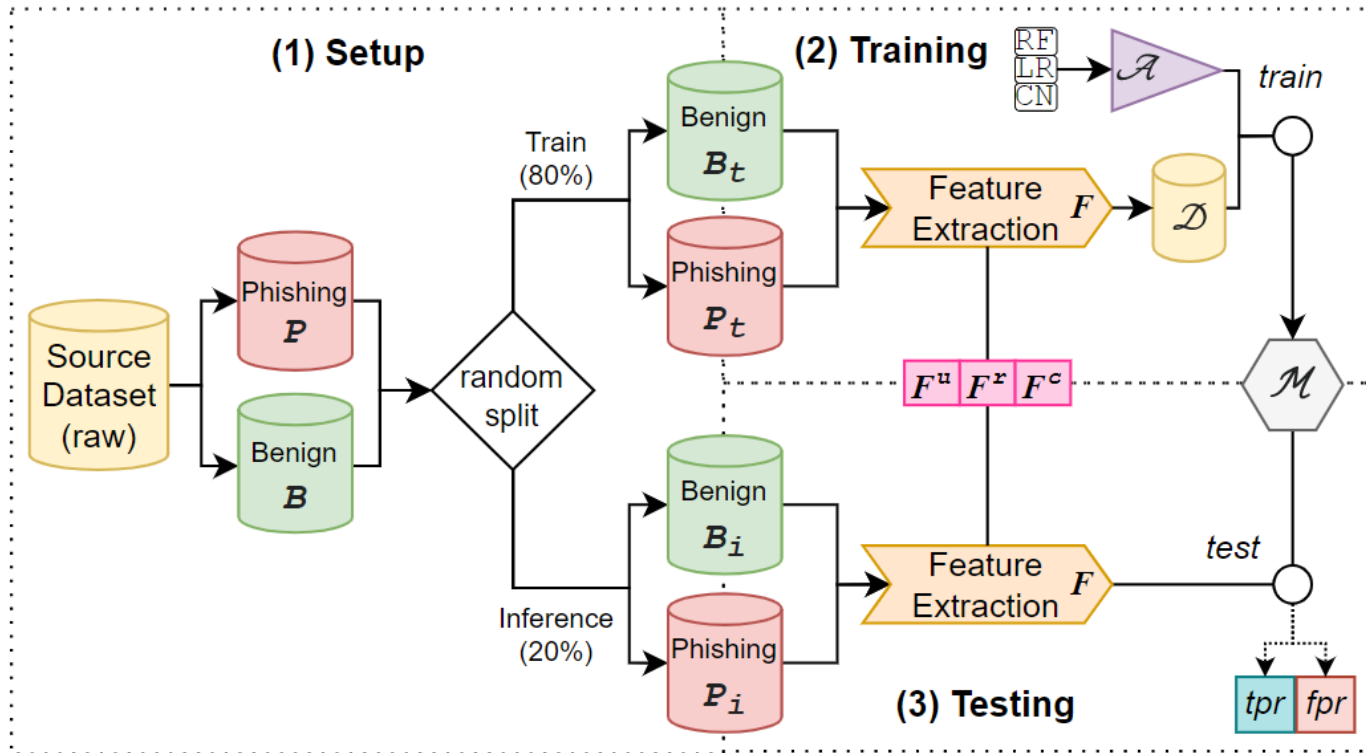

```
1   <div>
2       <form enctype="multipart/form-data" action=
        "//www.weebly.com/weebly/apps/formSubmit.php" method="POST" id=
        "form-723155629711391878">
3           <div id="723155629711391878-form-parent" class="wsite-form-container"
4               style="margin-top:10px;">
5           <ul class="formlist" id="723155629711391878-form-list">
6               <div><div class="wsite-form-field" style="margin:5px 0px 5px 0px;">
7               <label class="wsite-form-label" for="input-2279820181179653776">Email
                Address <span class="form-not-required">*</span></label>
8               <div class="wsite-form-input-container">
9                   <input id="input-2279820181179653776" class="wsite-form-input
                    wsite-input wsite-input-width-370px" type="text" name=
                    "_u2279820181179653776" />
10              </div>
11              <div id="instructions-2279820181179653776" class="wsite-form-instructions"
                style="display:none;"></div>
12          </div></div>
13
14          <a href="../fake-link-to-nonexisting-resource">
15          <font style="visibility:hidden">Resource</font></a>
16
17   <div><div class="wsite-form-field" style="margin:5px 0px 5px 0px;">
18              <label class="wsite-form-label" for="input-435728988405554593">Enter
                Password <span class="form-not-required">*</span></label>
19              <div class="wsite-form-input-container">
20                  <textarea id="input-435728988405554593" class="wsite-form-input
                    wsite-input wsite-input-width-370px" name="_u435728988405554593" style
                    ="height: 50px"></textarea>
21          </div>
```
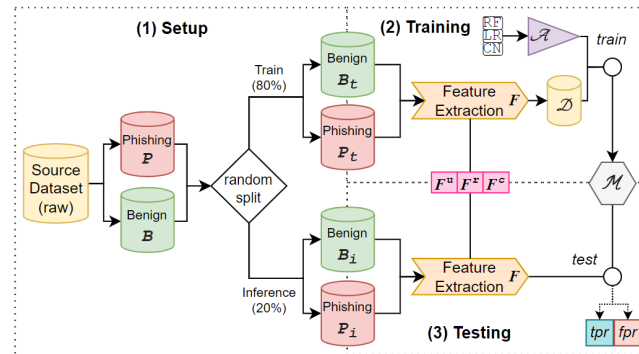
⇐ ε (WsP)

UNIVERSITÄT
LIECHTENSTEIN

# Evaluation – Workflow

o   Such attacks appear cheap, but are they effective? Let's assess their impact!
o   First step: develop proficient ML-PWD (high *tpr*, low *fpr*)
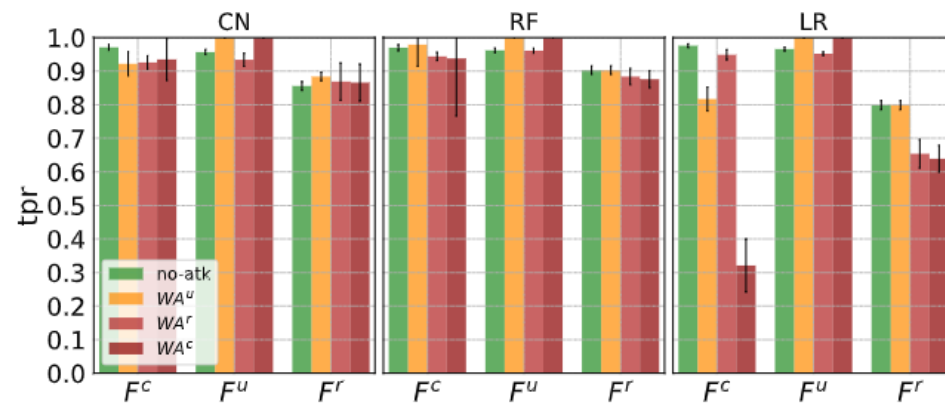
# Evaluation – Baseline

o   Such attacks appear cheap, but are they effective? Let's assess their impact!

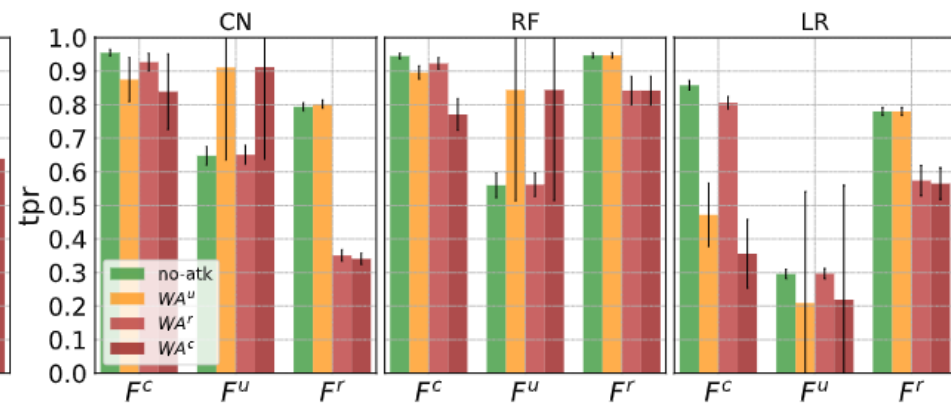o   First step: develop proficient ML-PWD (high *tpr*, low *fpr*)



o   Results comparable to the state-of-the-art ☺

o   Let's attack such ML-PWD
  • The *tpr* will decrease!

| $\mathcal{A}$ | $F$ | Zenodo | | $\delta$phish | |
|---|---|---|---|---|---|
| | | *tpr* | *fpr* | *tpr* | *fpr* |
| CN | $F^u$ | $0.96_{\pm 0.008}$ | $0.021_{\pm 0.0077}$ | $0.55_{\pm 0.030}$ | $0.037_{\pm 0.0076}$ |
| | $F^r$ | $0.88_{\pm 0.018}$ | $0.155_{\pm 0.0165}$ | $0.81_{\pm 0.019}$ | $0.008_{\pm 0.0020}$ |
| | $F^c$ | $0.97_{\pm 0.006}$ | $0.018_{\pm 0.0088}$ | $0.93_{\pm 0.013}$ | $0.005_{\pm 0.0025}$ |
| RF | $F^u$ | $0.98_{\pm 0.004}$ | $0.007_{\pm 0.0055}$ | $0.75_{\pm 0.022}$ | $0.003_{\pm 0.0014}$ |
| | $F^r$ | $0.93_{\pm 0.013}$ | $0.025_{\pm 0.0118}$ | $0.94_{\pm 0.016}$ | $0.006_{\pm 0.0025}$ |
| | $F^c$ | $\mathbf{0.98}_{\pm 0.006}$ | $\mathbf{0.007}_{\pm 0.0046}$ | $\mathbf{0.97}_{\pm 0.007}$ | $\mathbf{0.001}_{\pm 0.0011}$ |
| LR | $F^u$ | $0.95_{\pm 0.009}$ | $0.037_{\pm 0.0100}$ | $0.24_{\pm 0.017}$ | $0.011_{\pm 0.0026}$ |
| | $F^r$ | $0.82_{\pm 0.017}$ | $0.144_{\pm 0.0171}$ | $0.74_{\pm 0.025}$ | $0.018_{\pm 0.0036}$ |
| | $F^c$ | $0.96_{\pm 0.007}$ | $0.025_{\pm 0.0077}$ | $0.81_{\pm 0.020}$ | $0.013_{\pm 0.0037}$ |

UNIVERSITÄT LIECHTENSTEIN

47

# Results – Are WsP effective?



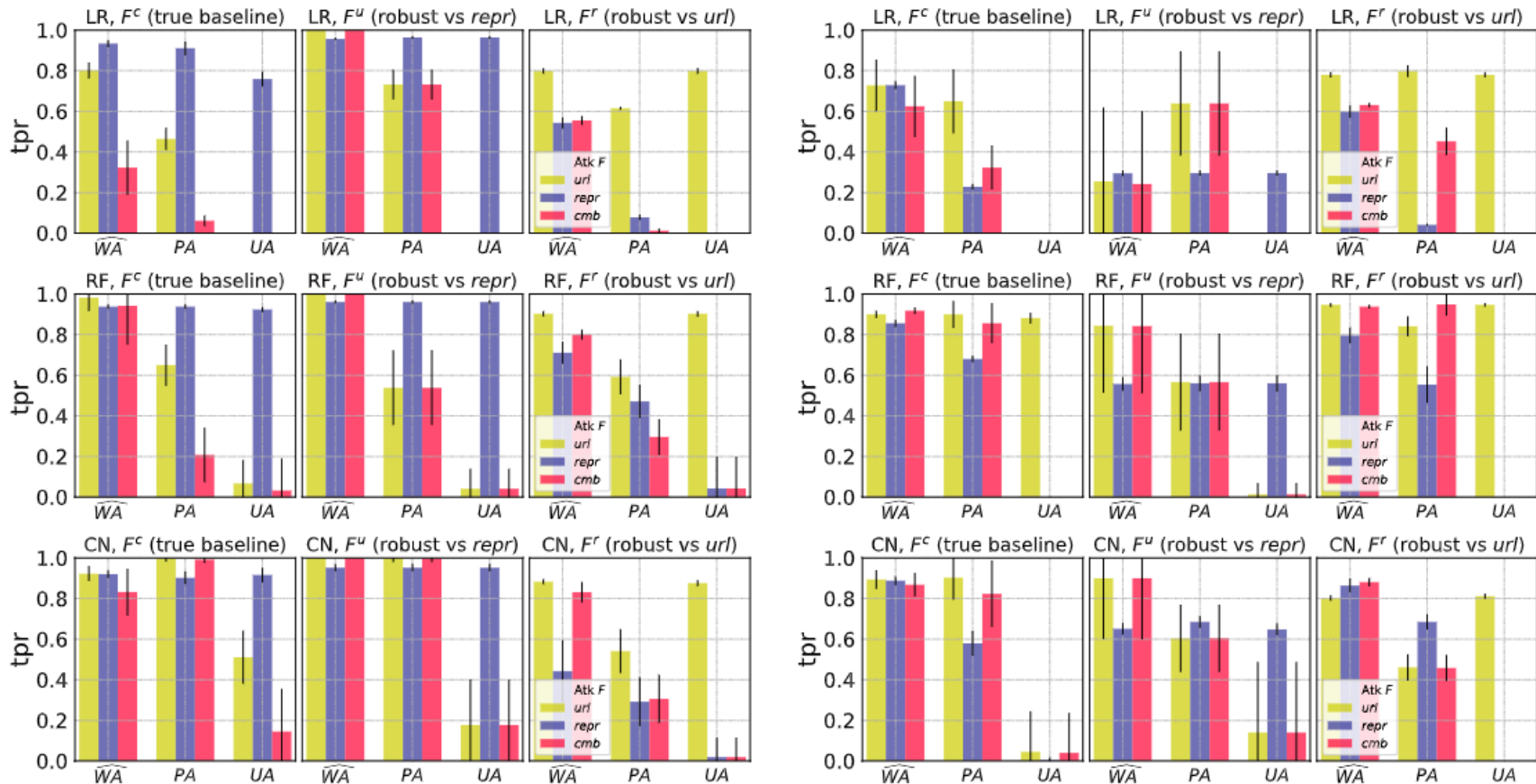(a) Impact of WA on the ML-PWD trained on Zenodo.

(b) Impact of WA on the ML-PWD trained on $\delta$Phish.

- In some cases, NO
  - This is *significant* because most past studies show ML-PWD being bypassed very easily!
- In some cases, VERY LITTLE
  - This is also significant, because even a 1% decrease in detection rate can be problematic when dealing with *millions of samples*!
- In other cases, YES
  - This is very significant, because WsP are cheap and are likely to be exploited by attackers!

UNIVERSITÄT
LIECHTENSTEIN

Bottom line: no free lunch!

# Results – What about attacks in the other spaces?

In general, attacks in the other spaces (PA and UA) are more disruptive…



However, such attacks also have a *higher cost*!
Will real attackers truly use them *just to evade* a ML-PWD?

UNIVERSITÄT
LIECHTENSTEIN

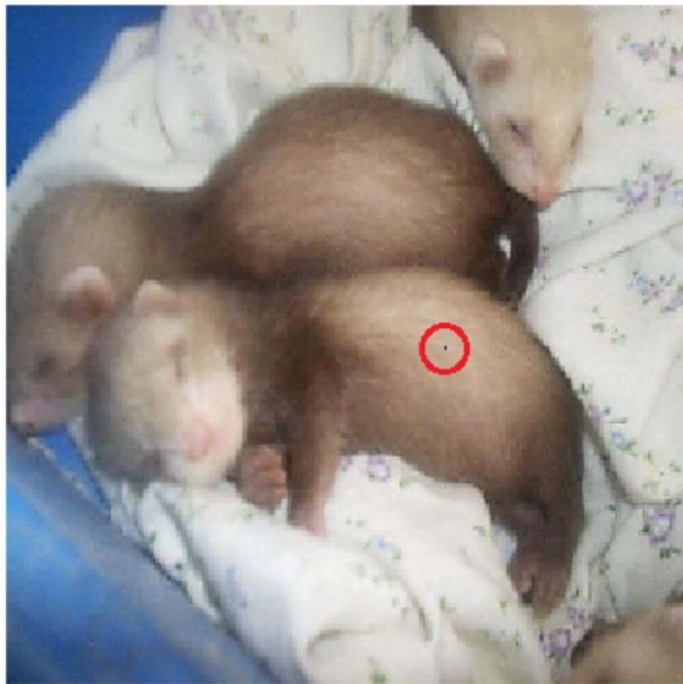# Adversarial Attacks against Humans and Machine Learning

# Scenario

o ML is used not only for cybersecurity, but for a plethora of other applications

o In some cases, the "decision making" is based on:

- The <u>output</u> of a *ML model*
- The interpretation of a *human* to such <u>output</u>

# Scenario

o ML is used not only for cybersecurity, but for a plethora of other applications

o In some cases, the "decision making" is based on:

- The <u>output</u> of a *ML model*
- The interpretation of a *human* to such <u>output</u>

o Case in point: online marketplace

- A person wants to sell an item (e.g., a car)
- This person (i.e., the seller) uploads the images of such an item on an online marketplace
- The marketplace automatically provides an estimate of the "value" of the corresponding item
    - This is done via ML
- Another person (i.e., a potential buyer) looks at the images, then looks at the "suggested" price, and determines whether to buy or not the corresponding item
    - The human uses the output of the ML model to make their decisions

UNIVERSITÄT
LIECHTENSTEIN

# Attack – what if…

o    What if the seller has malicious intentions?

→ The seller may want to induce the ML model to estimate a higher price

o    Doing this by introducing "imperceptible" perturbations may trick the ML…

o    …but not the human!

# Attack – what if…

o   What if the seller has malicious intentions?

→ The seller may want to induce the ML model to estimate a higher price

o   Doing this by introducing "imperceptible" perturbations may trick the ML…
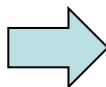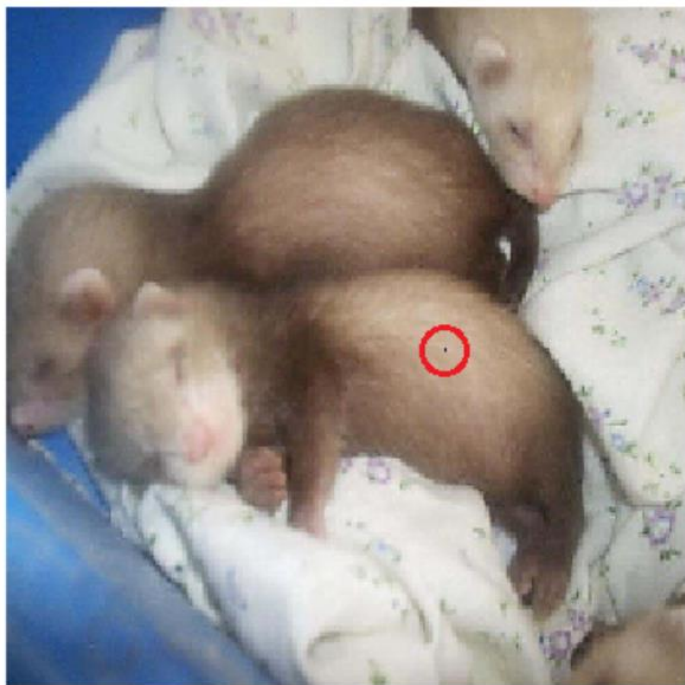
o   …but not the human!



Hamster(35.79%)
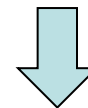Nipple(42.36%)

Reference: Su Jiawei, Danilo Vasconcellos Vargas, and Kouichi Sakurai. "One pixel attack for fooling deep neural networks." *IEEE Transactions on Evolutionary Computation* (2019)

54

Giovanni Apruzzese, *PhD*
*giovanni.apruzzese@uni.li*

# Attack – what if…

o   What if the seller has malicious intentions?

→ The seller may want to induce the ML model to estimate a higher price

o   Doing this by introducing "imperceptible" perturbations may trick the ML…

o   …but not the human!



In some cases, "imperceptible" perturbations **may not be what an attacker wants!**

This is especially true when there is a "human-in-the-loop".

Hamster(35.79%)
Nipple(42.36%)

Reference: Su Jiawei, Danilo Vasconcellos Vargas, and Kouichi Sakurai. "One pixel attack for fooling deep neural networks." *IEEE Transactions on Evolutionary Computation* (2019)

# Solution (high-level)

o   If humans are involved in the "decision making" process, then such humans <u>will react</u> to clearly incorrect outputs of ML models.

- Humans may suspect an adversarial <u>attack taking place</u>; or
- They may think that the ML model is faulty, and hence <u>not trust/believe its output</u>
- Both of the above are **detrimental** for the attacker!

UNIVERSITÄT
LIECHTENSTEIN

# Solution (high-level)

o   If humans are involved in the "decision making" process, then such humans will react to clearly incorrect outputs of ML models.

- Humans may suspect an adversarial attack taking place; or
- They may think that the ML model is faulty, and hence not trust/believe its output
- Both of the above are **detrimental** for the attacker!

## (Malicious) solution: deceive both the human *and* the ML model!

o   A ML model that thinks that a "FIAT Panda" is a "VW Polo" will output a very high price

- But if the "perturbation" only affects a single pixel, nobody will fall for it!

o   A FIAT Panda is clearly different than a VW Polo, so the perturbation (whatever it is) must be *perceived* by the human

→ The FIAT Panda must be changed in such a way that the human can be somewhat fooled

- E.g.: the human should think that "it could be a Panda… but it could also be a Polo"



- FIAT Panda MSRP: ~10k $
- VW Polo MSRP: ~20k $

# Solution (low-level)

o   How to achieve this in practice?
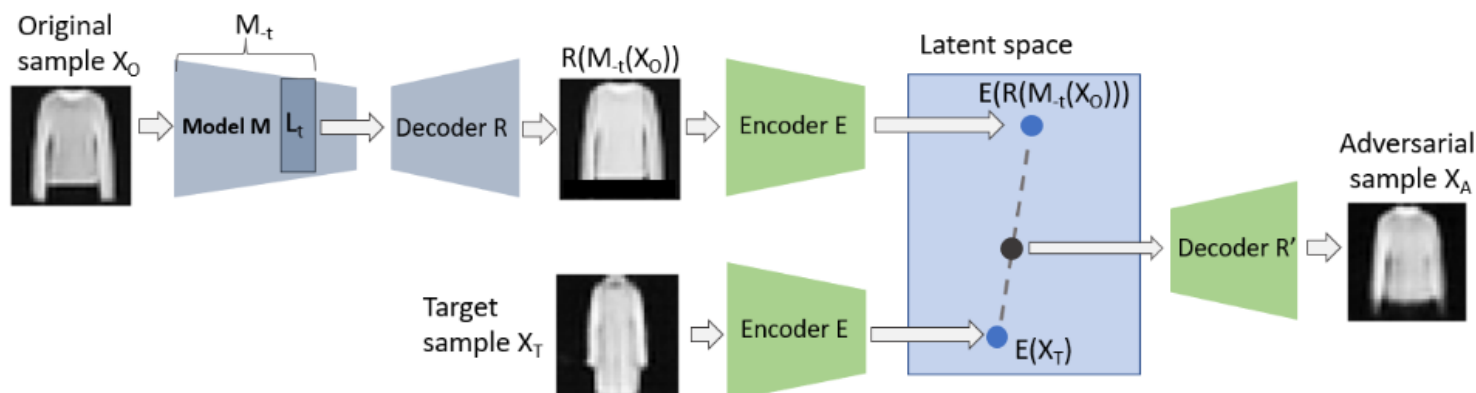
## Semantics Adversarial Attacks

o   The idea is using "explainability" techniques to create adversarial examples.

UNIVERSITÄT
LIECHTENSTEIN

58

# Solution (low-level)

o How to achieve this in practice?

## Semantics Adversarial Attacks

o The idea is using "explainability" techniques to create adversarial examples.

o Requirements:

- An "original sample" (i.e., a FIAT Panda)
- A desired "target sample" (i.e., a VW Polo)
- A given magnitude of the perturbation (neither too big nor too small)
  - If the FIAT Panda "becomes" a VW Polo, then the adversarial attack would be unfair
  - …and the "buyer" will complain ☺
- The details of a ML model (which must be based on Convolutional Neural Networks)
  - These attacks <u>can</u> be transferred!

o Output: an "adversarial example" that is a mix between the original and target sample
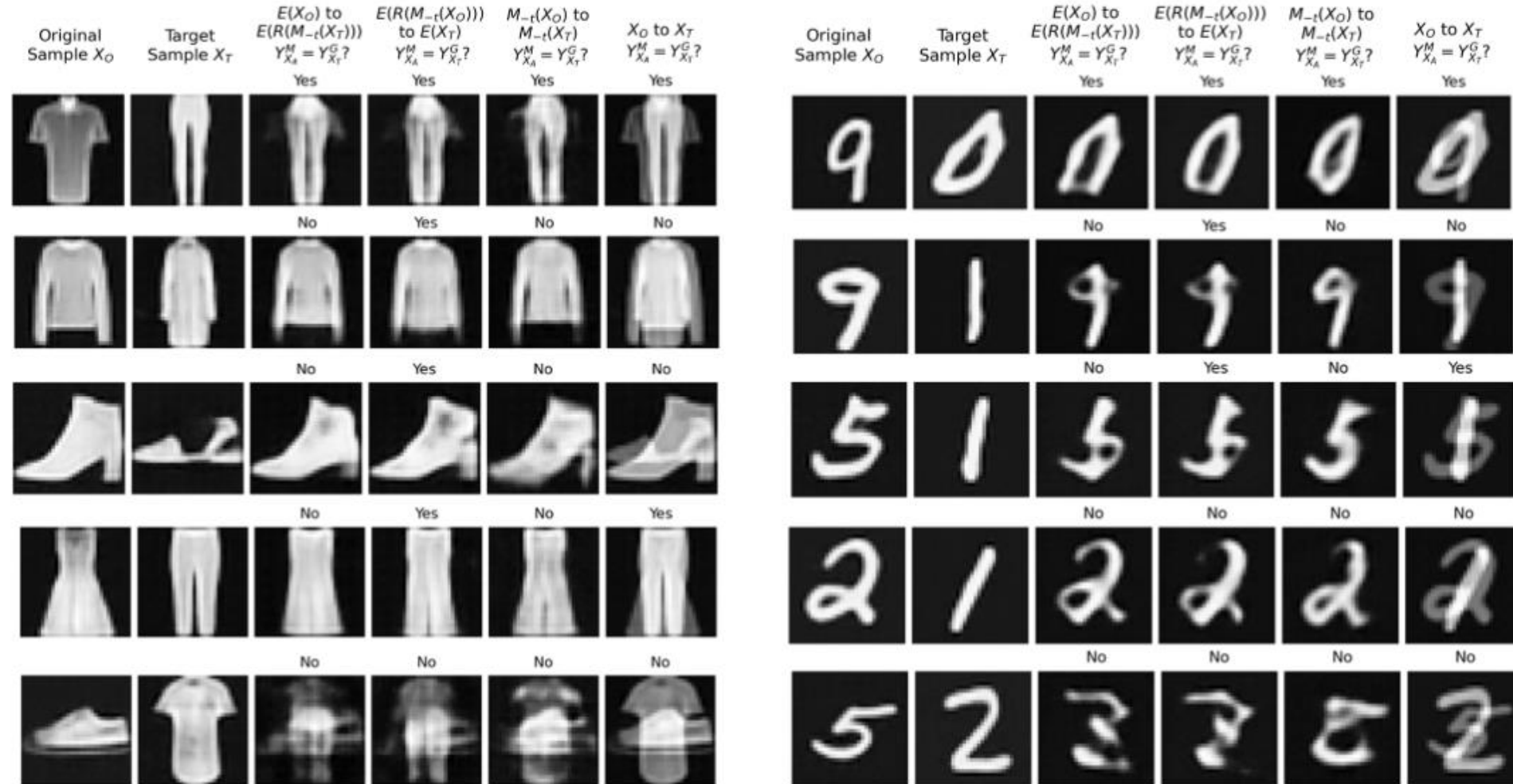
# Experiments



Fig. 2: Original, target and adversarial samples for different en-/decodings and interpolation for Fashion-MNIST(left) and MNIST(right). Yes/No indicates, whether the model got fooled by $X_A$, i.e. it outputs the class of $X_T$ for $X_A$

# The relationship between
# **Machine Learning & Cybersecurity**

Giovanni Apruzzese, PhD
TU Delft – May 3rd, 2022

UNIVERSITÄT
LIECHTENSTEIN