

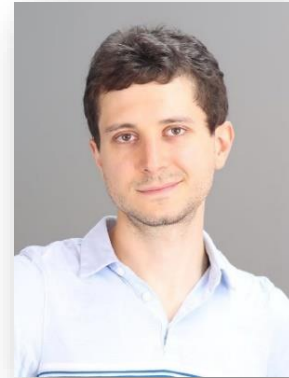


# Some Pragmatic Relationships between **Machine Learning & Cybersecurity**

Giovanni Apruzzese, PhD

May 17th, 2022

# whoami: Dr. Giovanni Apruzzese



## ○ 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: Prof. 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.

## ○ Interests:

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

## ○ Contact information:

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

## What I do

# Machine Learning + Cybersecurity

- Applying ML to *provide security* of a given information system
  - E.g.: using ML to detect network intrusions
- *Attacking / Defending* ML applications
  - E.g.: evading a ML model that detects phishing websites
- Using ML *offensively* against any target
  - E.g.: artificially generating “fake” images

## BONUS

- Using ML offensively to attack a ML-based security system



# Outline

## ○ Using **unlabelled data** for Machine Learning in Cyberthreat Detection

- Ref: Giovanni Apruzzese, Pavel Laskov, Aliya Tastemirova. "SoK: The Impact of Unlabelled Data in Cyberthreat Detection" *IEEE European Symposium on Security and Privacy*. June 2022

## ○ 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*. May 2022

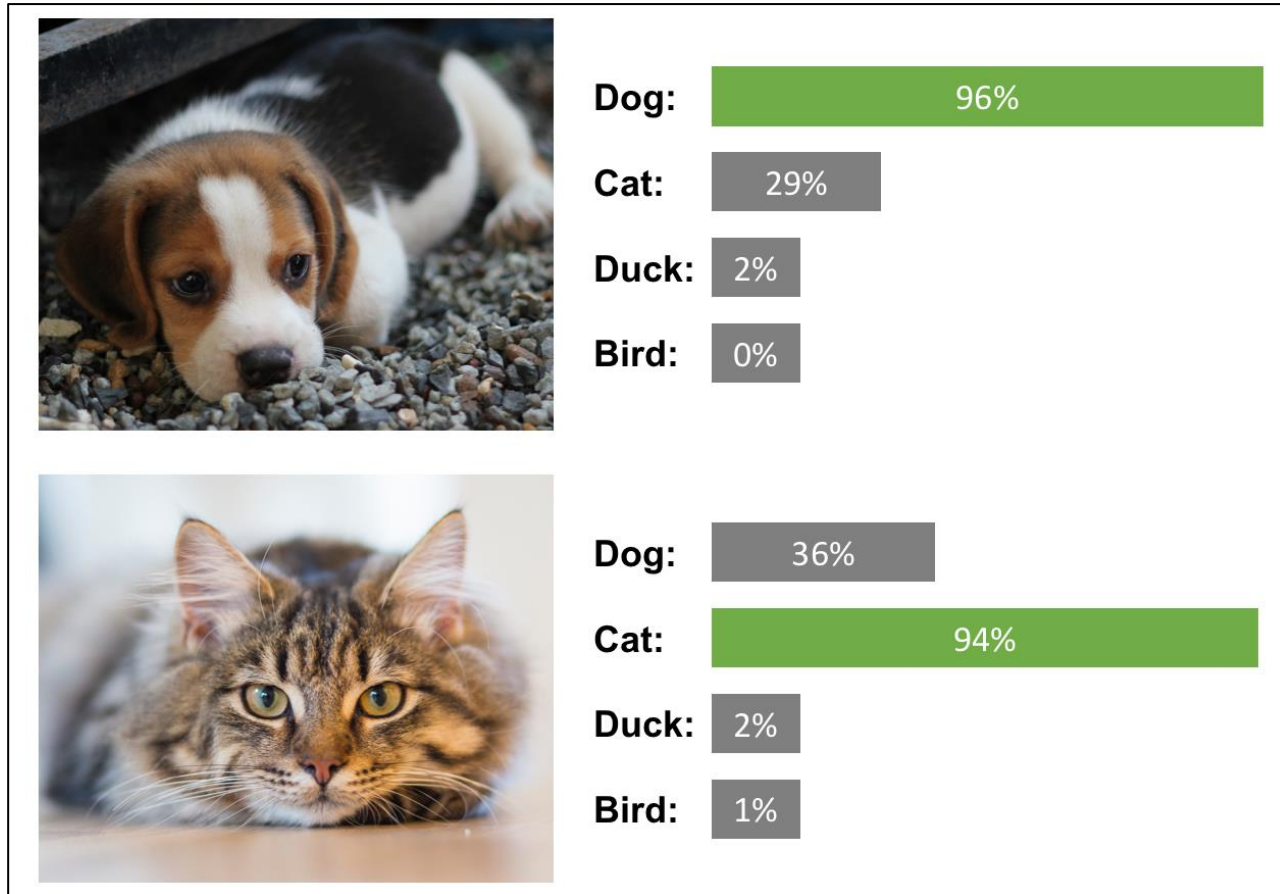
# Unlabelled data for Machine Learning in Cyberthreat Detection

## Labelling in Cyberthreat Detection (CTD)

- Using Machine Learning (ML) to *detect* cyber-threats requires *labelled data*.
- Obtaining *plenty* and *accurate* labels is expensive:
  - Human in the loop (but this is true also for Computer Vision...)

## Labelling in Cyberthreat Detection (CTD)

- Using Machine Learning (ML) to *detect* cyber-threats requires *labelled data*.
- Obtaining *plenty* and *accurate* labels is expensive:
  - Human in the loop (but this is true also for Computer Vision...)





# Labelling in Cyberthreat Detection (CTD)

- Using Machine Learning (ML) to *detect* cyber-threats requires *labelled data*.
- Obtaining *plenty* and *accurate* labels is expensive:
  - Human in the loop (but this is true also for Computer Vision...)
  - Instead, in CTD...

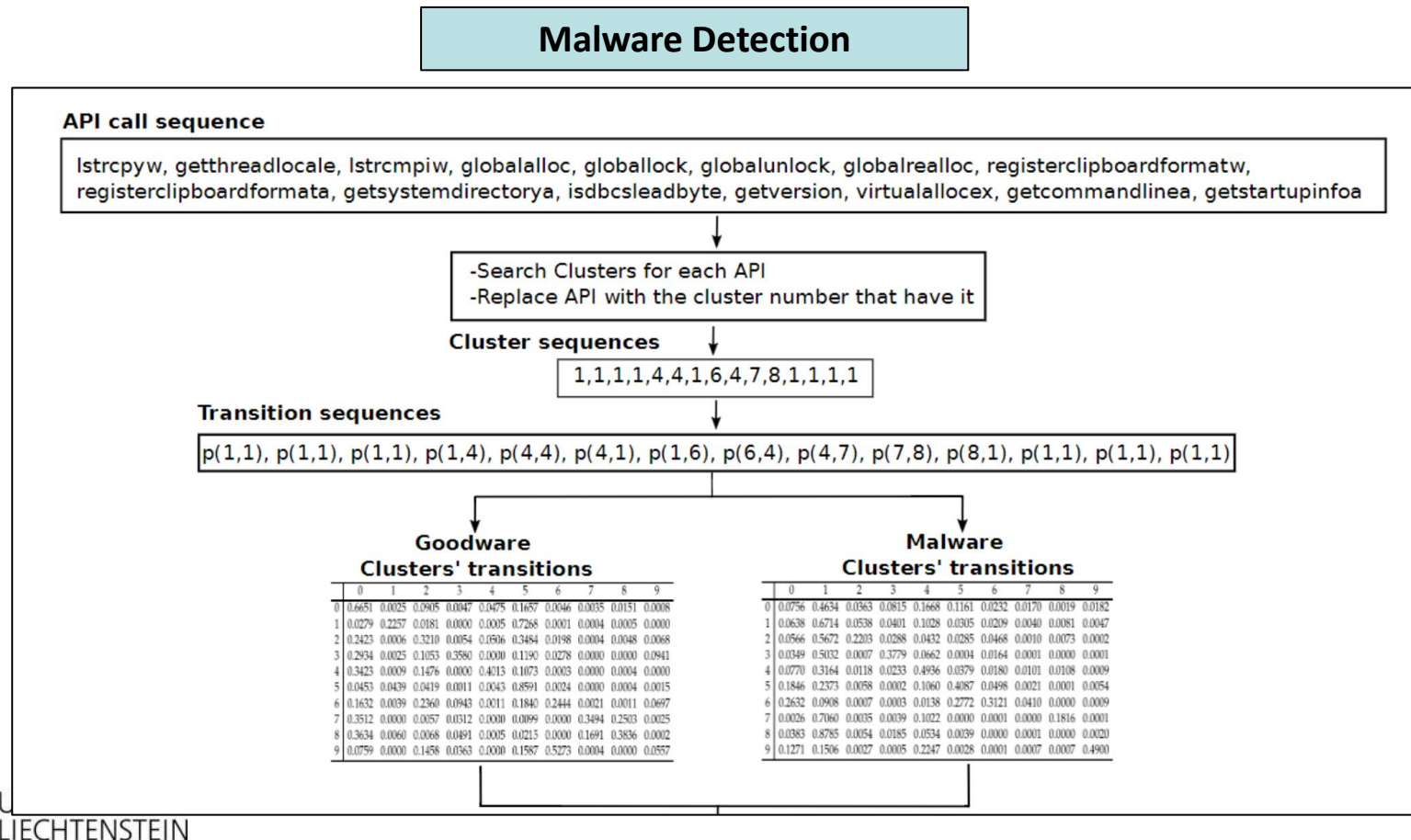
## Network Intrusion Detection

	StartTime	Dur	Proto	SrcAddr	Sport	Dir	DstAddr	Dport	State	sTos	dTos	TotPkts	TotBytes	SrcBytes
0	2011/08/10 09:46:53.047277	3550.182373	udp	212.50.71.179	39678	<->	147.32.84.229	13363	CON	0.0	0.0	12	875	413
1	2011/08/10 09:46:53.048843	0.000883	udp	84.13.246.132	28431	<->	147.32.84.229	13363	CON	0.0	0.0	2	135	75
2	2011/08/10 09:46:53.049895	0.000326	tcp	217.163.21.35	80	<?>	147.32.86.194	2063	FA_A	0.0	0.0	2	120	60
3	2011/08/10 09:46:53.053771	0.056966	tcp	83.3.77.74	32882	<?>	147.32.85.5	21857	FA_FA	0.0	0.0	3	180	120
4	2011/08/10 09:46:53.053937	3427.768066	udp	74.89.223.204	21278	<->	147.32.84.229	13363	CON	0.0	0.0	42	2856	1596
...	...	...	...	...	...	...	...	...	...	...	...	...	...	...
675872	2011/08/10 11:04:27.118993	0.020525	udp	147.32.84.165	1025	<->	147.32.80.9	53	CON	0.0	0.0	2	590	87
675873	2011/08/10 11:04:27.119042	2.919582	tcp	147.32.84.194	3112	->	147.32.80.13	80	FSPA_FSPA	0.0	0.0	172	160606	7360
675874	2011/08/10 11:04:27.124802	0.000281	udp	147.32.84.59	57993	<->	147.32.80.9	53	CON	0.0	0.0	2	235	76
675875	2011/08/10 11:04:27.125921	0.000195	udp	147.32.84.59	60616	<->	147.32.80.9	53	CON	0.0	0.0	2	218	76
675876	2011/08/10 11:04:27.129857	0.011865	tcp	147.32.84.59	52776	->	217.31.58.184	80	FSPA_FSPA	0.0	0.0	10	1615	943



# Labelling in Cyberthreat Detection (CTD)

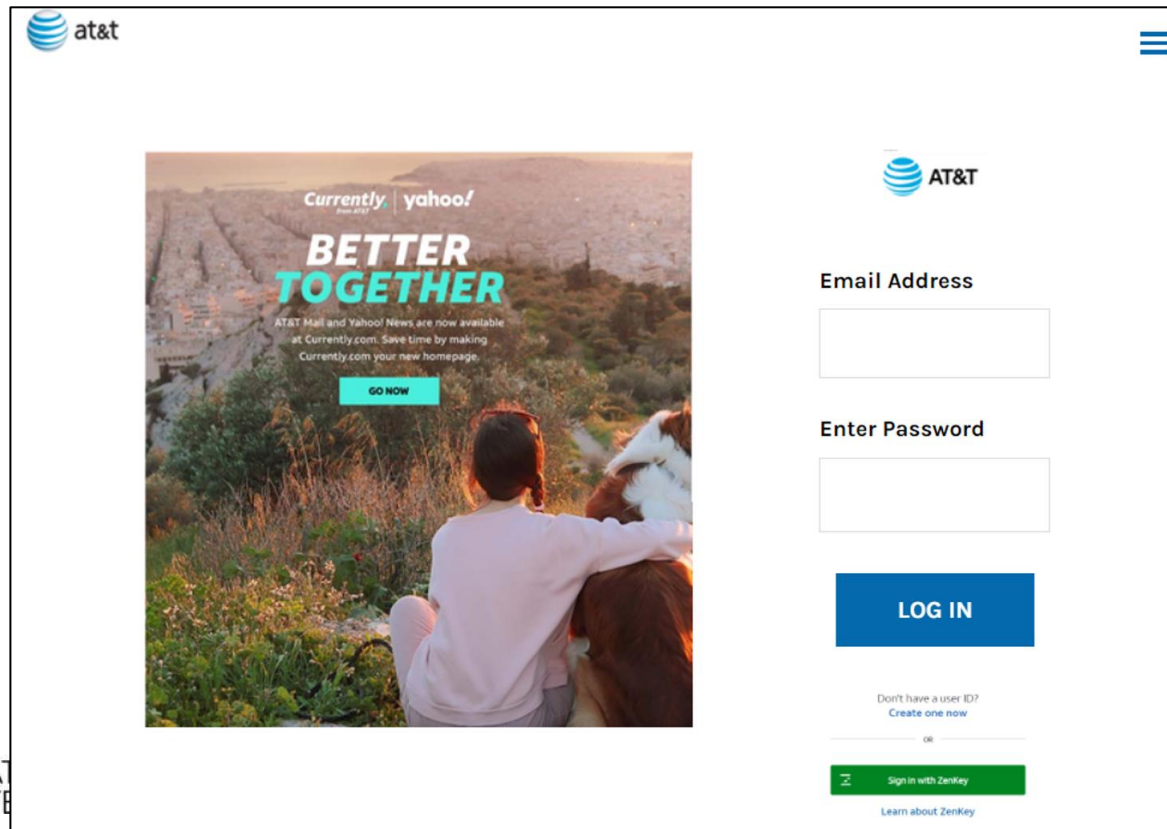
- Using Machine Learning (ML) to *detect* cyber-threats requires *labelled data*.
- Obtaining *plenty* and *accurate* labels is expensive:
  - Human in the loop (but this is true also for Computer Vision...)
  - Instead, in CTD...



# Labelling in Cyberthreat Detection (CTD)

- Using Machine Learning (ML) to *detect* cyber-threats requires *labelled data*.
- Obtaining *plenty* and *accurate* labels is expensive:
  - Human in the loop (but this is true also for Computer Vision...)
  - Instead, in CTD...

## Phishing Website Detection



## Labelling in Cyberthreat Detection (CTD)

- Using Machine Learning (ML) to *detect* cyber-threats requires *labelled data*.
- Obtaining *plenty* and *accurate* labels is expensive:
  - Human in the loop (but this is true also for Computer Vision...)
  - Instead, in CTD...

**For CTD, labelling requires *expert knowledge* and is an *error prone* task.**

And the “concept drift” further aggravates this issue...

# Semisupervised Learning

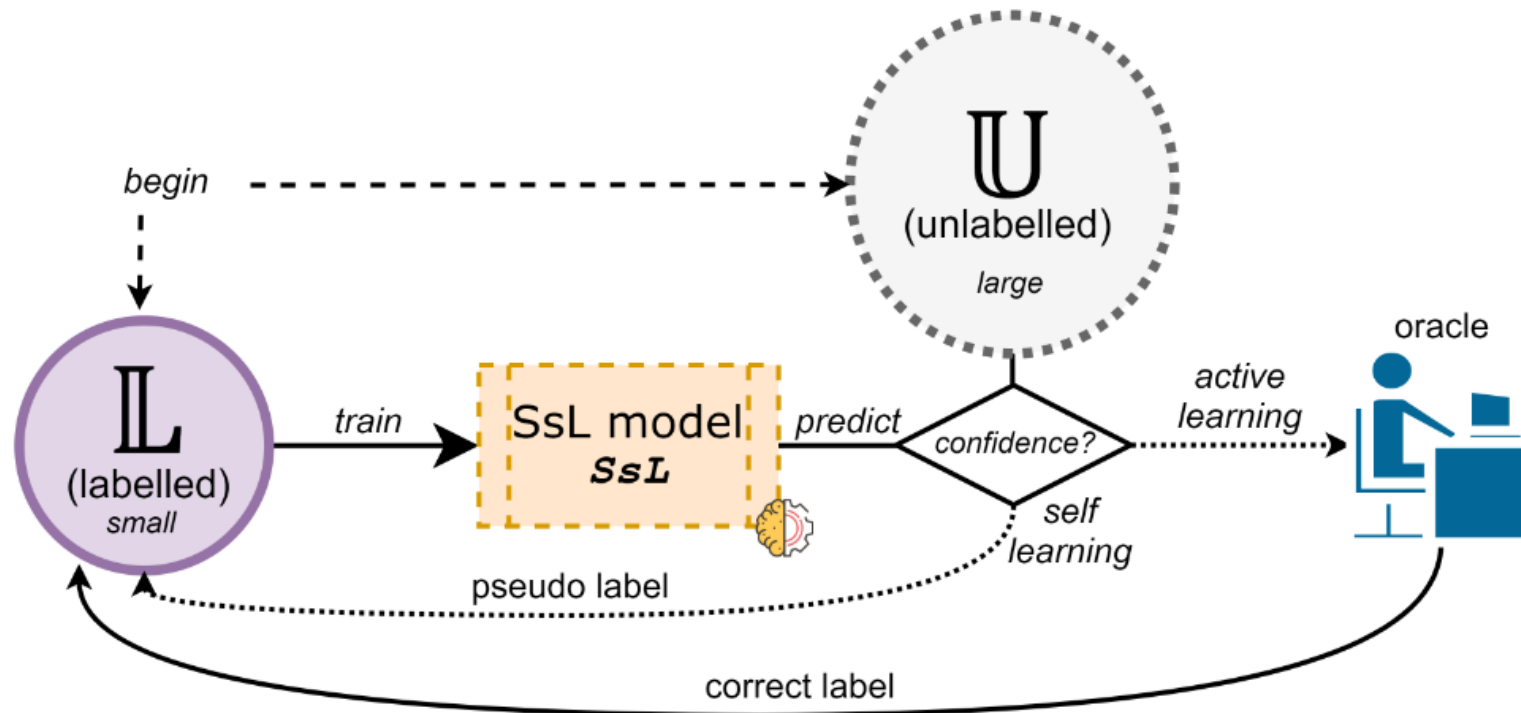
- Labelled data is expensive, but *unlabelled* data is cheap(er).
- Why not using unlabelled data to improve the proficiency of ML models?

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

# Semisupervised Learning

- Labelled data is expensive, but *unlabelled* data is cheap(er).
- Why not using unlabelled data to improve the proficiency of ML models?

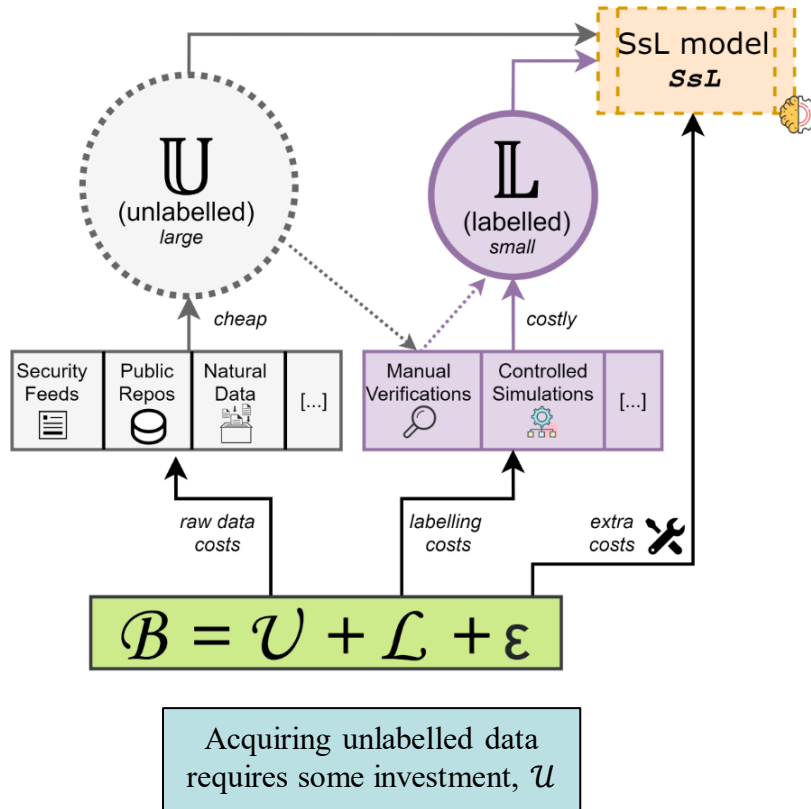
Mixing *labelled* with *unlabelled* data is a ML approach denoted as  
**“Semisupervised Learning” (SsL)**



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

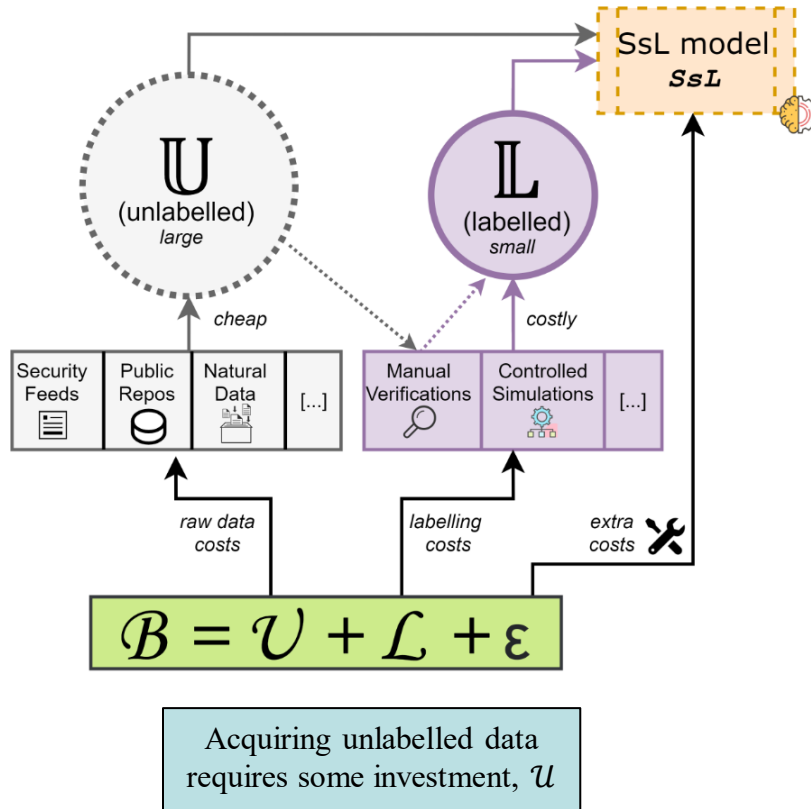
# Goal of Semisupervised Learning

- Developing SsL models is cheaper than “supervised learning” (SL) models, **but it is not free.**

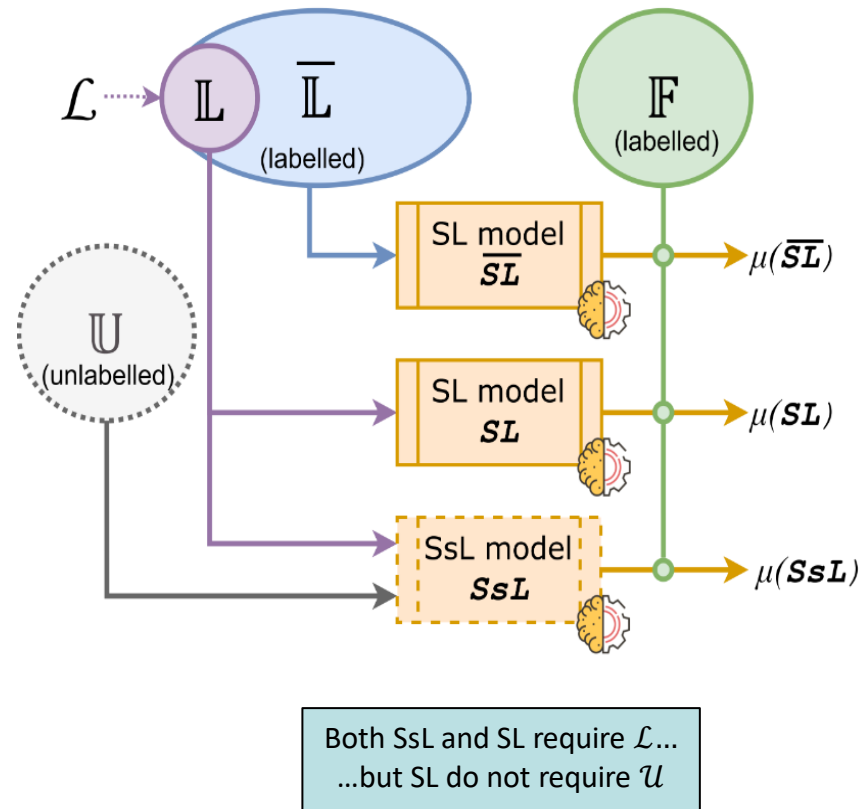


# Goal of Semisupervised Learning

- Developing SsL models is cheaper than “supervised learning” (SL) models, **but it is not free.**



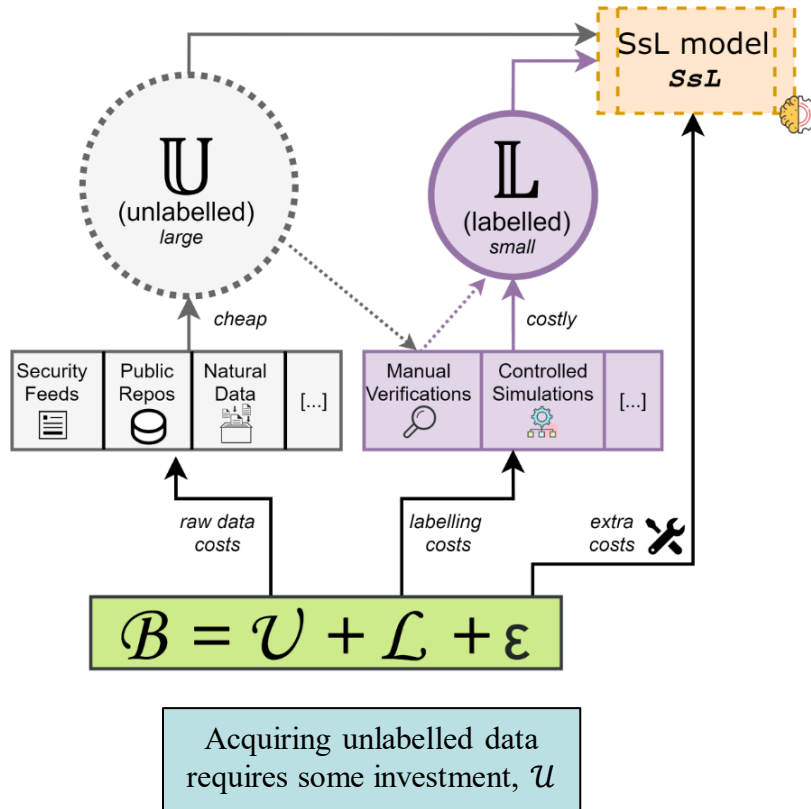
- A SsL model should achieve a *superior performance* than a SL model that uses the *same labelling budget,  $L$*



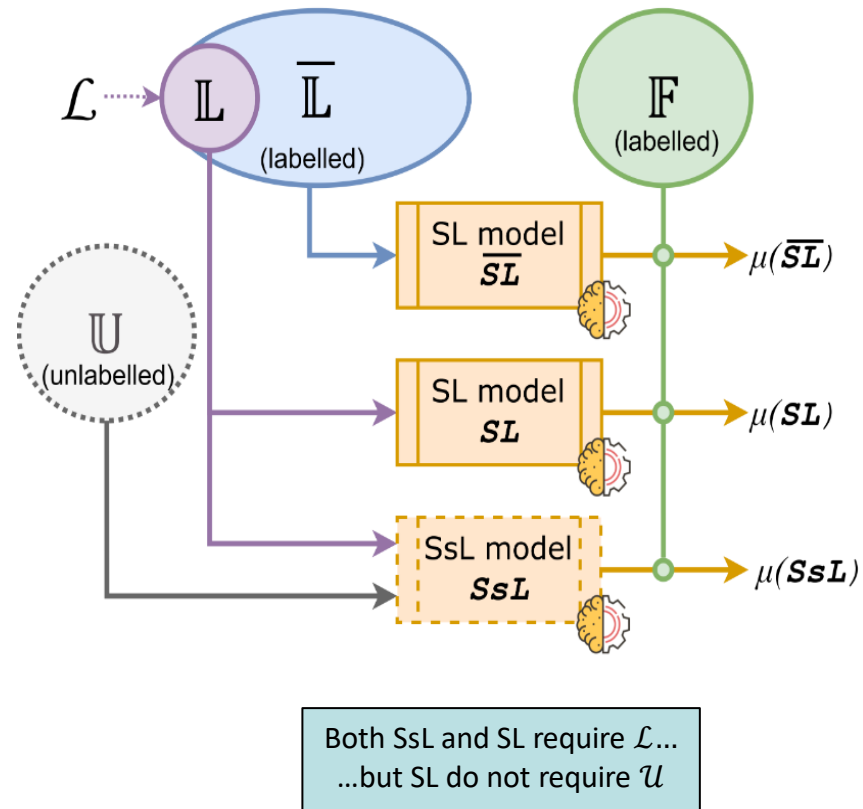


# Goal of Semisupervised Learning

- Developing SsL models is cheaper than “supervised learning” (SL) models, **but it is not free.**



- A SsL model should achieve a *superior performance* than a SL model that uses the *same labelling budget*,  $L$



**Definition 1.** The *goal* of a Semisupervised Learning (SsL) method is using  $U$  alongside any  $L$  obtained with  $L$  to devise a model  $SsL$ . After deployment, such  $SsL$  should predict the ground truth of the samples in  $F$  by achieving a performance  $\mu(SsL)$  that is:  $\mu(SL) < \mu(SsL) \leq \mu(\bar{SL})$ .

# Problem: nobody cares

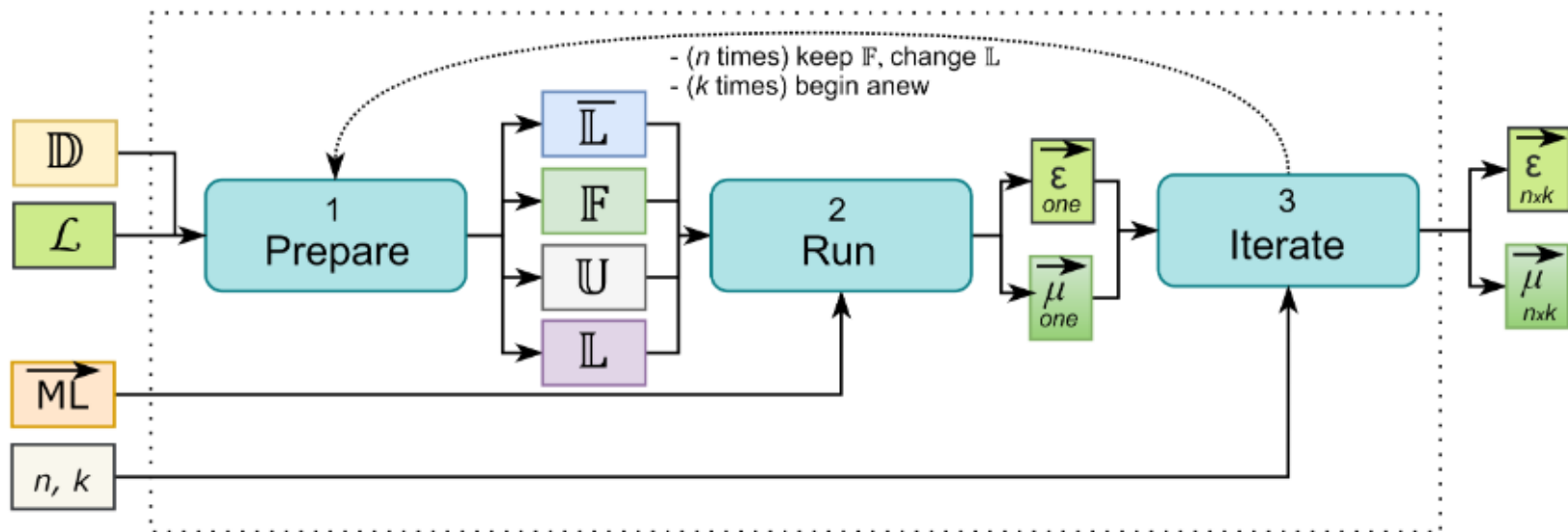
The current state-of-the-art does not allow to determine whether SsL methods applied in Cyberthreat Detection are truly beneficial

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
	Görmitz [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
Phishing Detection	Zhang [70]	2021	●	●	✗	●	✗	✓	●	CICIDS2017, CTU13
	Liang [109]	2021	✓	●	✓	●	✓	✓	●	NSL-KDD
	Gyawali [110]	2011	✗	✓	✓	✗	✓	✓	●	Private
	Zhao [111]	2013	✓	✓	✓	✓	✗	✓	✓*	DetMalURL
	Gabriel [15]	2017	●	●	✗	✗	✗	✗	●	Private
	Yang [112]	2017	✓	●	✗	✗	✓	✓	●	Private
Malware Detection	Bhattacharjee [113]	2017	✗	✓	✗	●	✗	✗	●	Private
	Li [55]	2017	✓	✓	✓	●	✓	✓	✗	Custom
	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

- SsL is intriguing, but its “pragmatic” benefits are still unknown
- Identifying (and quantifying) such benefits requires adopting a rigorous workflow

→ CEF-SsL: Cybersecurity Evaluation Framework for Semisupervised Learning



# (re)Evaluation

- 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
$\hat{\pi} SsL$	0.584	0.435	0.818	0.849	0.883	0.777	0.470	0.641	0.890
$\alpha SsL_l$	<b>0.693</b>	0.582	<b>0.897</b>	<b>0.863</b>	<b>0.903</b>	0.770	<b>0.546</b>	<b>0.687</b>	<b>0.924</b>
$\alpha SsL_o$	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	<b>0.595</b>	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

## (re)Evaluation

- 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
$SsL$	0.613	0.447	0.879	0.852	0.886	<b>0.778</b>	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
$\hat{\pi} SsL$	0.584	0.435	0.818	0.849	0.883	0.777	0.470	0.641	0.890
$\alpha SsL_l$	<b>0.693</b>	0.582	<b>0.897</b>	<b>0.863</b>	<b>0.903</b>	0.770	<b>0.546</b>	<b>0.687</b>	<b>0.924</b>
$\alpha SsL_o$	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	<b>0.595</b>	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?

# (re)Evaluation

- 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
$SsL$	0.613	0.447	0.879	0.852	0.886	<b>0.778</b>	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
$\hat{\pi} SsL$	0.584	0.435	0.818	0.849	0.883	0.777	0.470	0.641	0.890
$\alpha SsL_l$	<b>0.693</b>	0.582	<b>0.897</b>	<b>0.863</b>	<b>0.903</b>	0.770	<b>0.546</b>	<b>0.687</b>	<b>0.924</b>
$\alpha SsL_o$	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	<b>0.595</b>	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	$SsL$	<b>0.873</b>	0.159	$\alpha SsL_l$	< 0.001	4.310
UNB15	1104	$SsL$	<b>0.964</b>	-0.044	$\alpha^\pi SsL_o$	< 0.001	15.98
IDS17	540	$SsL$	<b>0.932</b>	0.085	$\alpha SsL_l$	<b>0.978</b>	-0.027
UCI	1200	$SsL$	<b>0.473</b>	0.717	$\alpha SsL_l$	< 0.001	7.386
Mend.	1200	$SsL$	<b>0.713</b>	0.368	$\alpha SsL_l$	< 0.001	6.757
$\delta$ Phish	1200	$SsL$	<b>0.554</b>	-0.590	$\alpha SsL_l$	0.002	-3.113
Drebin	1200	$SsL$	<b>0.310</b>	1.015	$\alpha SsL_l$	< 0.001	11.78
Ember	1200	$SsL$	<b>0.603</b>	-0.512	$\alpha SsL_l$	< 0.001	3.407
AndMal	1200	$SsL$	<b>0.712</b>	-0.370	$\alpha SsL_l$	< 0.001	12.01

# **Adversarial Attacks against Humans and Machine Learning**



## Scenario

- ML is used not only for cybersecurity, but for a plethora of other applications
- In some cases, the “decision making” is based on:
  - The output of a *ML model*
  - The interpretation of a *human* to such output

## Scenario

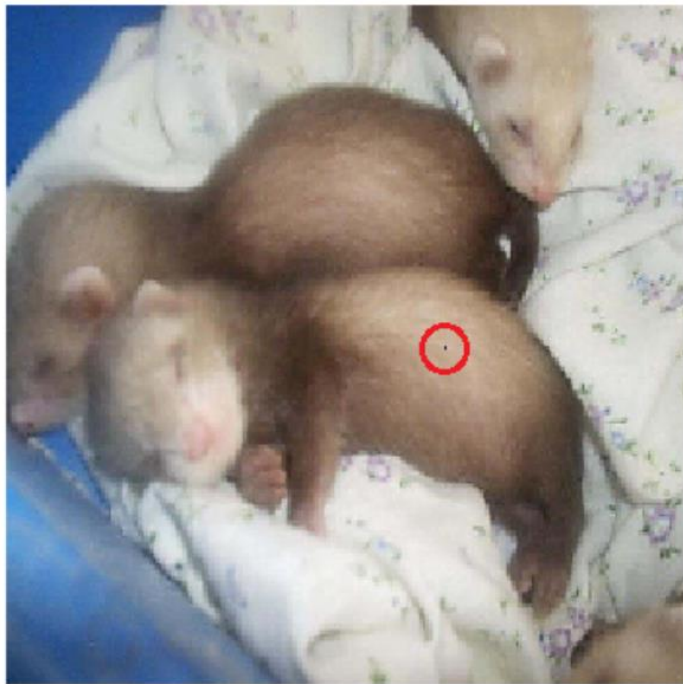
- ML is used not only for cybersecurity, but for a plethora of other applications
- In some cases, the “decision making” is based on:
  - The output of a *ML model*
  - The interpretation of a *human* to such output
- Case in point: online marketplace
  - A person wants to sell an item (e.g., a car)
  - 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

## Attack – what if...

- What if the seller has malicious intentions?
  - The seller may want to induce the ML model to estimate a higher price
- Doing this by introducing “imperceptible” perturbations may trick the ML...
- ...but not the human!

## Attack – what if...

- What if the seller has malicious intentions?  
→ The seller may want to induce the ML model to estimate a higher price
- Doing this by introducing “imperceptible” perturbations may trick the ML...
- ...but not the human!

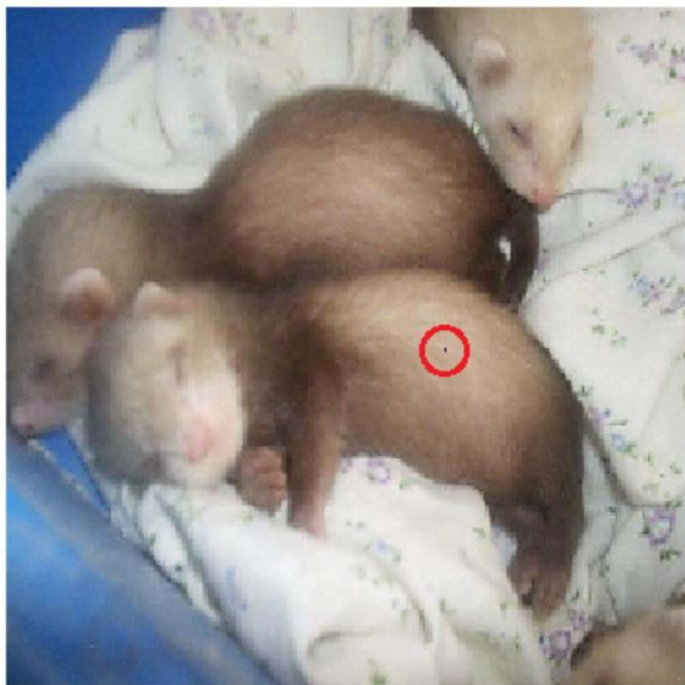


**Hamster(35.79%)**

**Nipple(42.36%)**

## Attack – what if...

- What if the seller has malicious intentions?  
→ The seller may want to induce the ML model to estimate a higher price
- Doing this by introducing “imperceptible” perturbations may trick the ML...
- ...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%)**

## Solution (high-level)

- 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!

## Solution (high-level)

- 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!

- 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!
- 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)

- How to achieve this in practice?

### Concept-based Adversarial Attacks

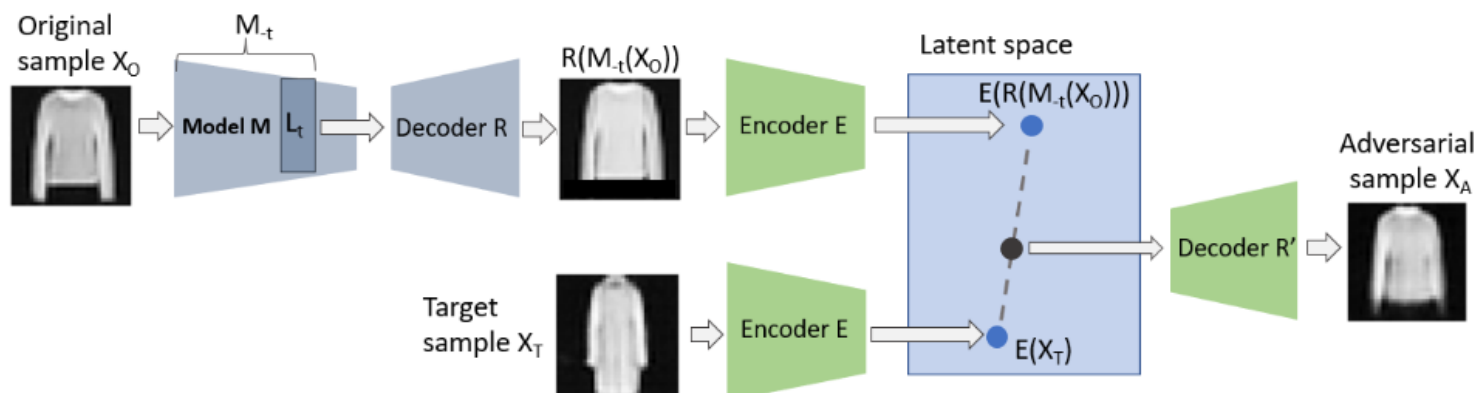
- The idea is using “explainability” techniques to create adversarial examples.

## Solution (low-level)

- How to achieve this in practice?

### Concept-based Adversarial Attacks

- The idea is using “explainability” techniques to create adversarial examples.
- 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 can be transferred!
- Output: an “adversarial example” that is a mix between the original and target sample



# Experiments

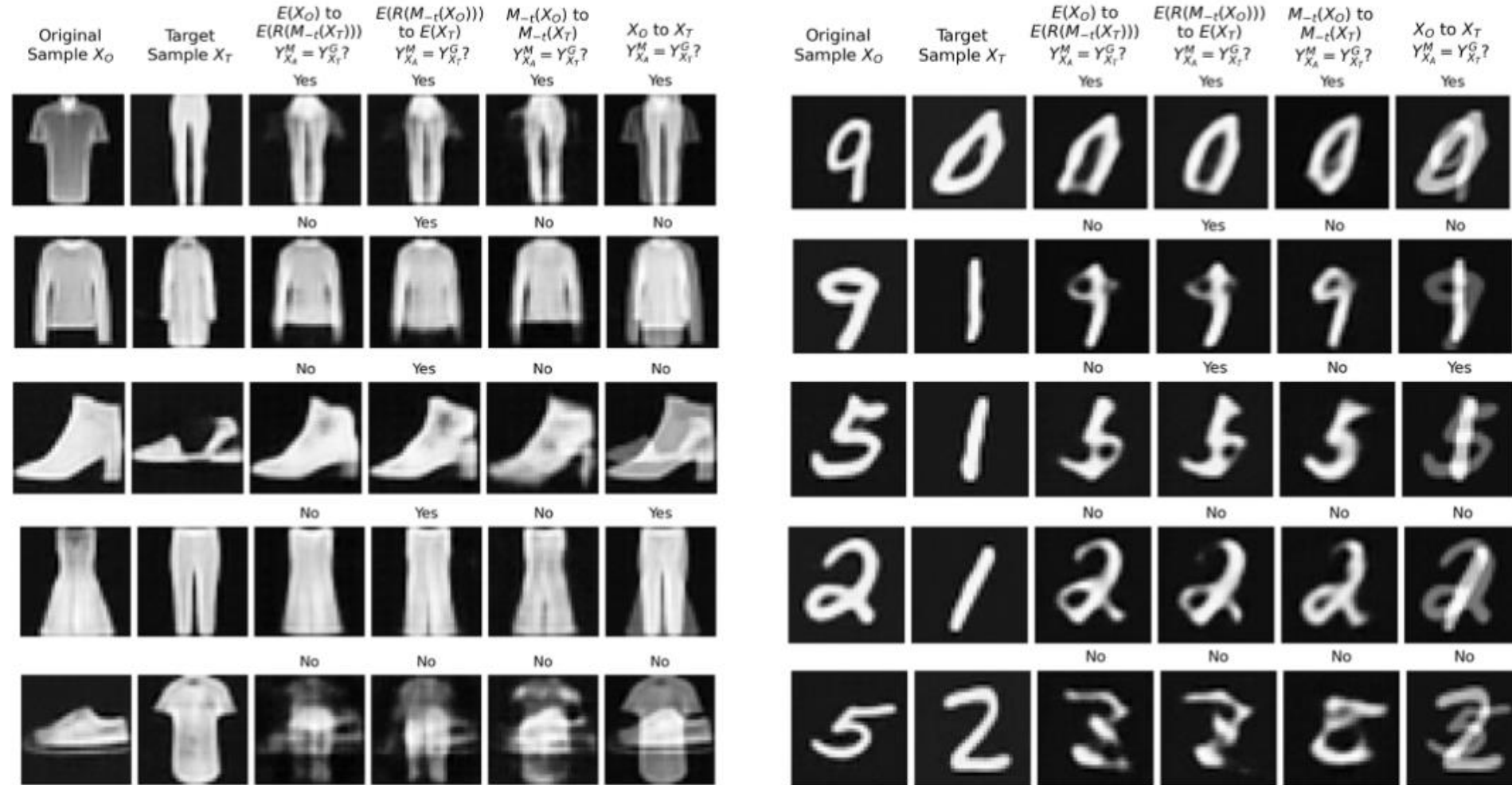


Fig. 2: Original, target and adversarial samples for different en-/decodings and interpolation for Fashion-MNIST(left) and MNIST(right). Yes/No indicates, whether the model got fooled by  $X_A$ , i.e. it outputs the class of  $X_T$  for  $X_A$



# Some Pragmatic Relationships between **Machine Learning & Cybersecurity**

Giovanni Apruzzese, PhD

May 17th, 2022