



European Symposium on Security and Artificial Intelligence
EU Cyber Week – November 21st, 2024

“Real Attackers Don’t Compute Gradients”: Bridging the Gap between Adversarial ML Research and Practice

Giovanni Apruzzese

(based on a joint work with: Hyrum S. Anderson, Savino Dambra, David Freeman, Fabio Pierazzi, Kevin Roundy)



**ROBUST
INTELLIGENCE**



Backstory (Dagstuhl – July 10-15th, 2022)



Giovanni Apruzzese, PhD
giovanni.apruzzese@uni.li

- Research seminar on the “Security of Machine Learning”

Backstory (Dagstuhl – July 10-15th, 2022)



SCHLOSS DAGSTUHL
Leibniz-Zentrum für Informatik

Giovanni Apruzzese, PhD
giovanni.apruzzese@uni.li

- Research seminar on the “Security of Machine Learning”
- The seminar opened with a talk by K. Grosse, showcasing the results of an extensive survey with ML practitioners about the security of ML [5]:

“Why do so?”

Backstory (Dagstuhl – July 10-15th, 2022)

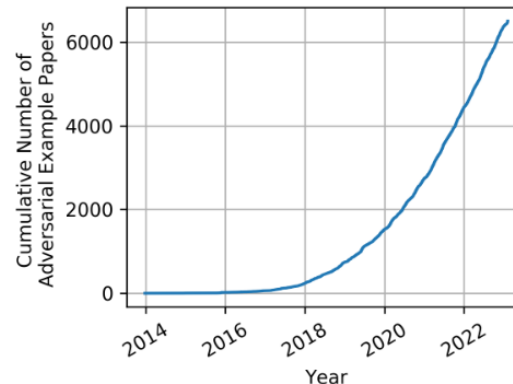


- Research seminar on the “Security of Machine Learning”
- The seminar opened with a talk by K. Grosse, showcasing the results of an extensive survey with ML practitioners about the security of ML [5]:

“Why do so?”

- Many discussions revolved around the impact of our research to the real world.

Apparently, the overwhelming number of works on adversarial ML research were not seen as problematic by practitioners!



Backstory (Dagstuhl – July 10-15th, 2022)

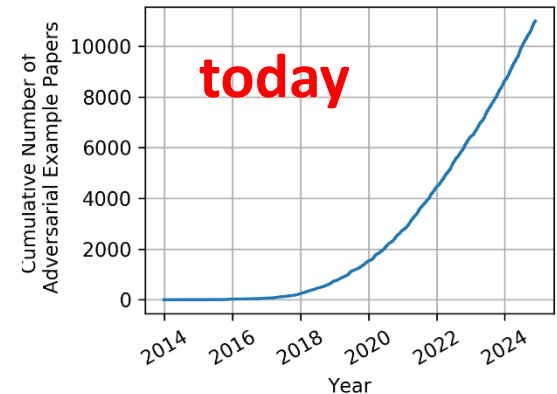
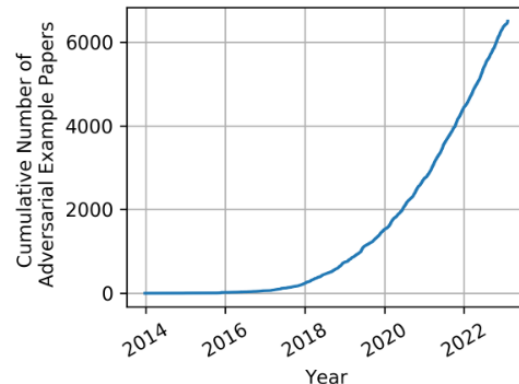


- Research seminar on the “Security of Machine Learning”
- The seminar opened with a talk by K. Grosse, showcasing the results of an extensive survey with ML practitioners about the security of ML [5]:

“Why do so?”

- Many discussions revolved around the impact of our research to the real world.

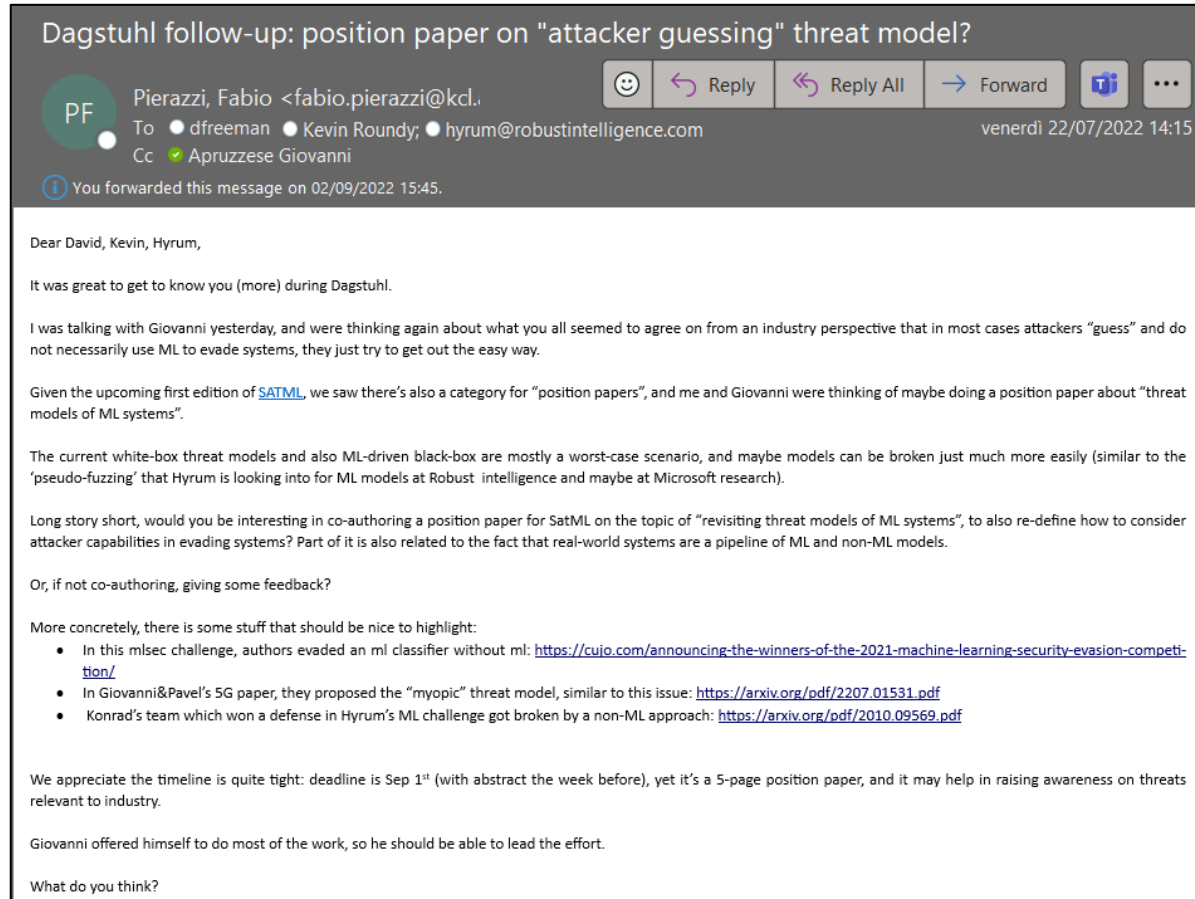
Apparently, the overwhelming number of works on adversarial ML research were not seen as problematic by practitioners!



- A recurring observation by some of the seminar’s attendees from industry was that:

Backstory (Earth – July 22nd, 2022)

- One week later, I was having a (remote) call with Fabio Pierazzi, and...



We appreciate the timeline is quite tight: deadline is Sep 1st (with abstract the week before), yet it's a 5-page position paper, and it may help in raising awareness on threats relevant to industry.

Our paper
has 26 pages!

Do real attackers compute gradients?



*A real
attacker*

Do real attackers compute gradients? (Case Study)

- We tried answering this question by looking at the AI Incident Database [78]...
- ...but **we could not find any evidence** of real incidents stemming from “adversarial examples” (or which leverage gradient computations)

Do real attackers compute gradients? (Case Study)

- We tried answering this question by looking at the AI Incident Database [78]...
- ...but **we could not find any evidence** of real incidents stemming from “adversarial examples” (or which leverage gradient computations)
- So, we asked a well-known **cybersecurity company** to provide us with data from their (operational!) phishing website detector, empowered by *deep learning*
- Just in July 2022, there were **9K samples** for which the ML detector was “uncertain”
 - They were “close to the decision boundary”, and required manual triage by experts
- We **manually analyzed** these (phishing) samples, trying to understand the root-causes of these “adversarial webpages”

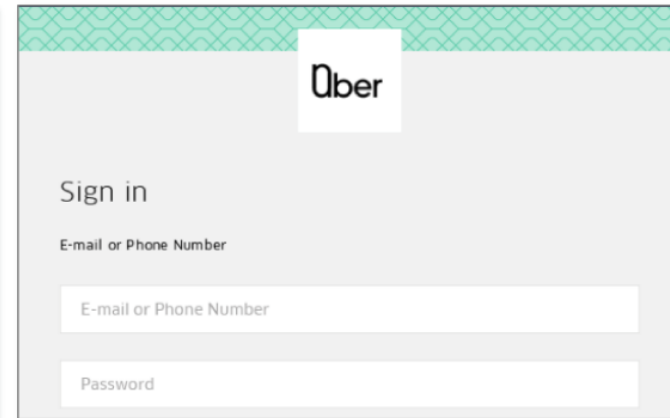
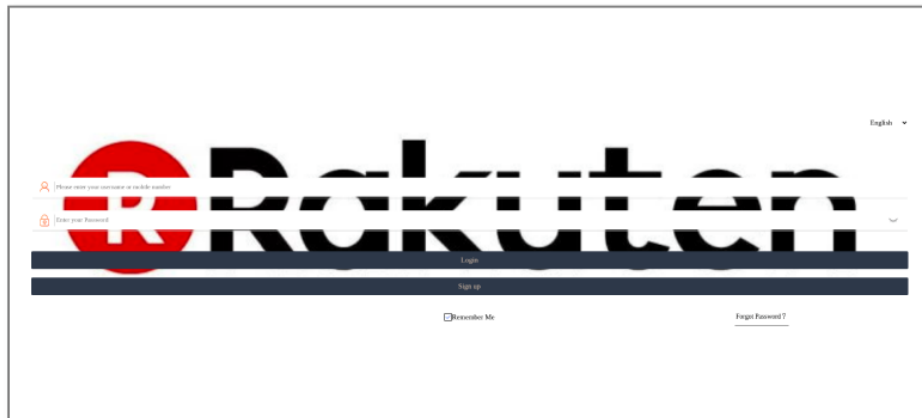
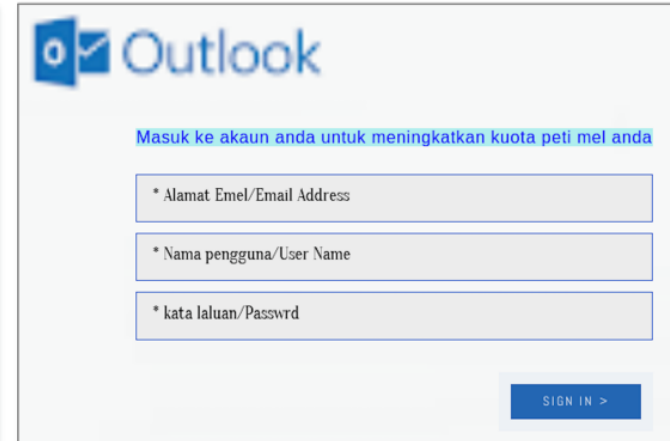
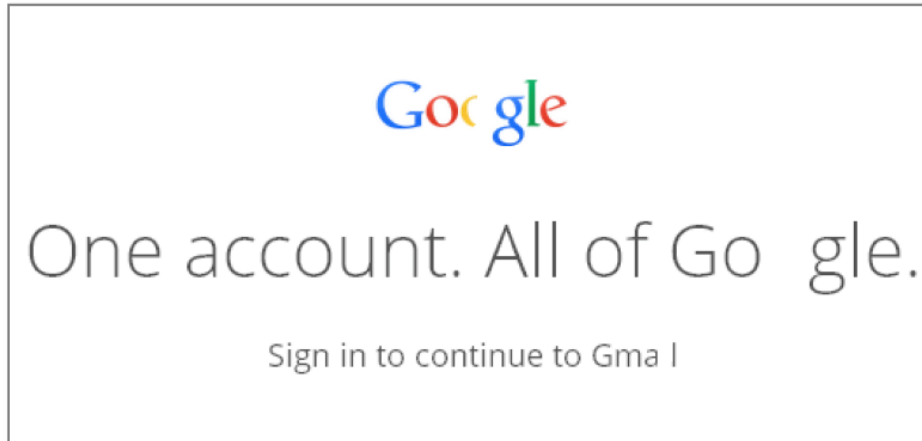
What did we find?

Do real attackers compute gradients? (Case Study) [cont'd]

- The **vast majority** of these webpages were “out of distribution”
 - They were different from any sample in the training set
- We then looked at a small subset of the remaining ones...

Do real attackers compute gradients? (Case Study) [cont'd]

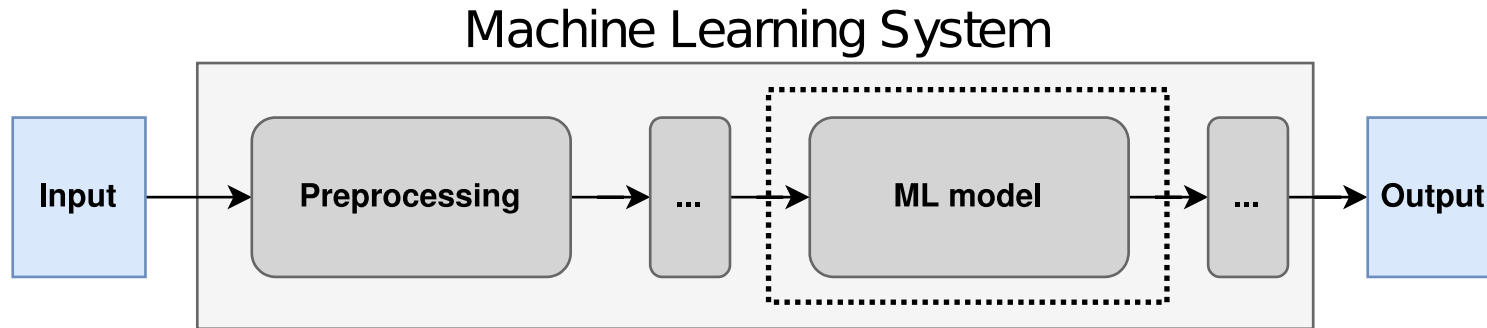
- The **vast majority** of these webpages were “out of distribution”
 - They were different from any sample in the training set
- We then looked at a small subset of the remaining ones...



Machine Learning Systems

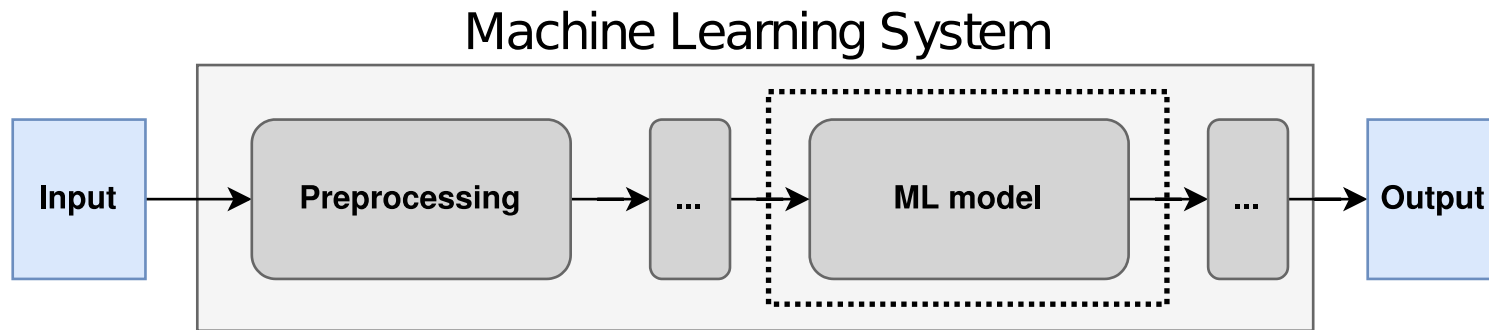
Machine Learning Systems

- In reality, ML models are a single component of a complex ML system
 - Real ML systems (are likely to) have also elements *that have nothing to do with ML*

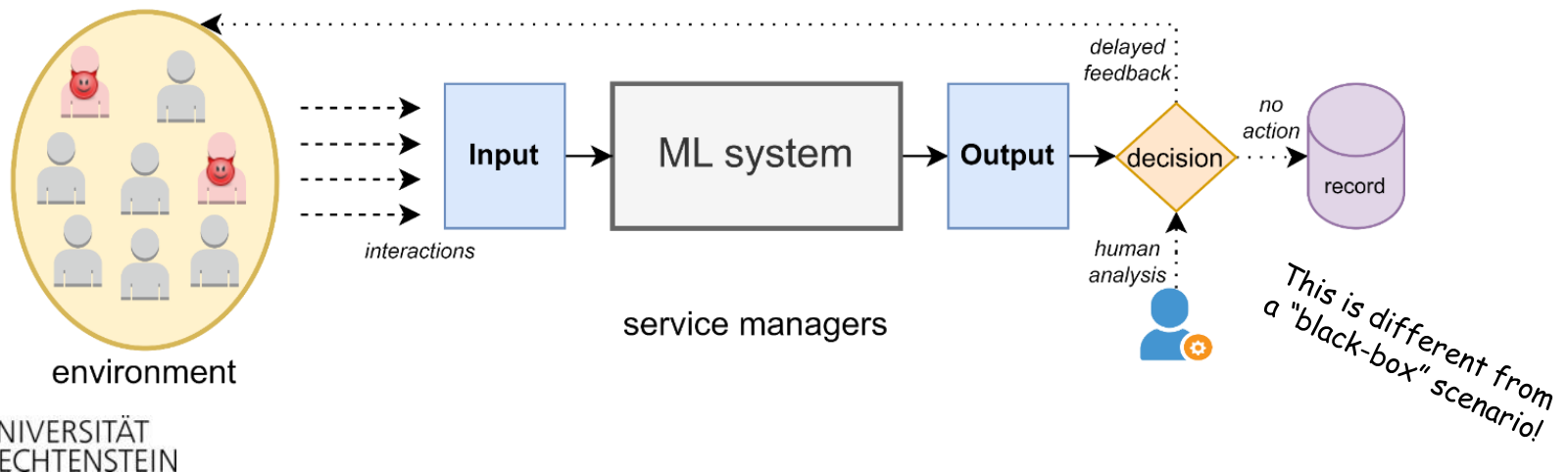


Machine Learning Systems

- In reality, ML models are a single component of a complex ML system
 - Real ML systems (are likely to) have also elements *that have nothing to do with ML*

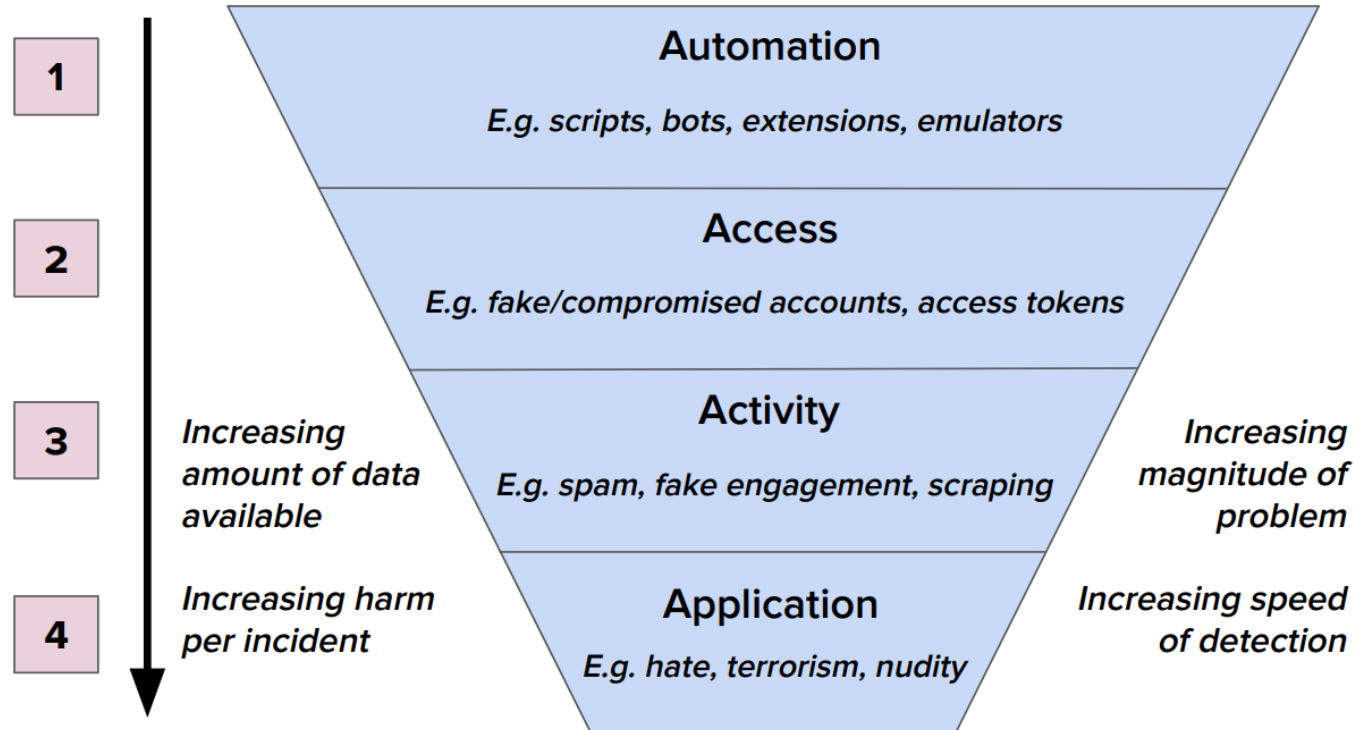


- Some ML systems are “invisible” to their users (and, hence, to real attackers)



Machine Learning Systems (Case Study)

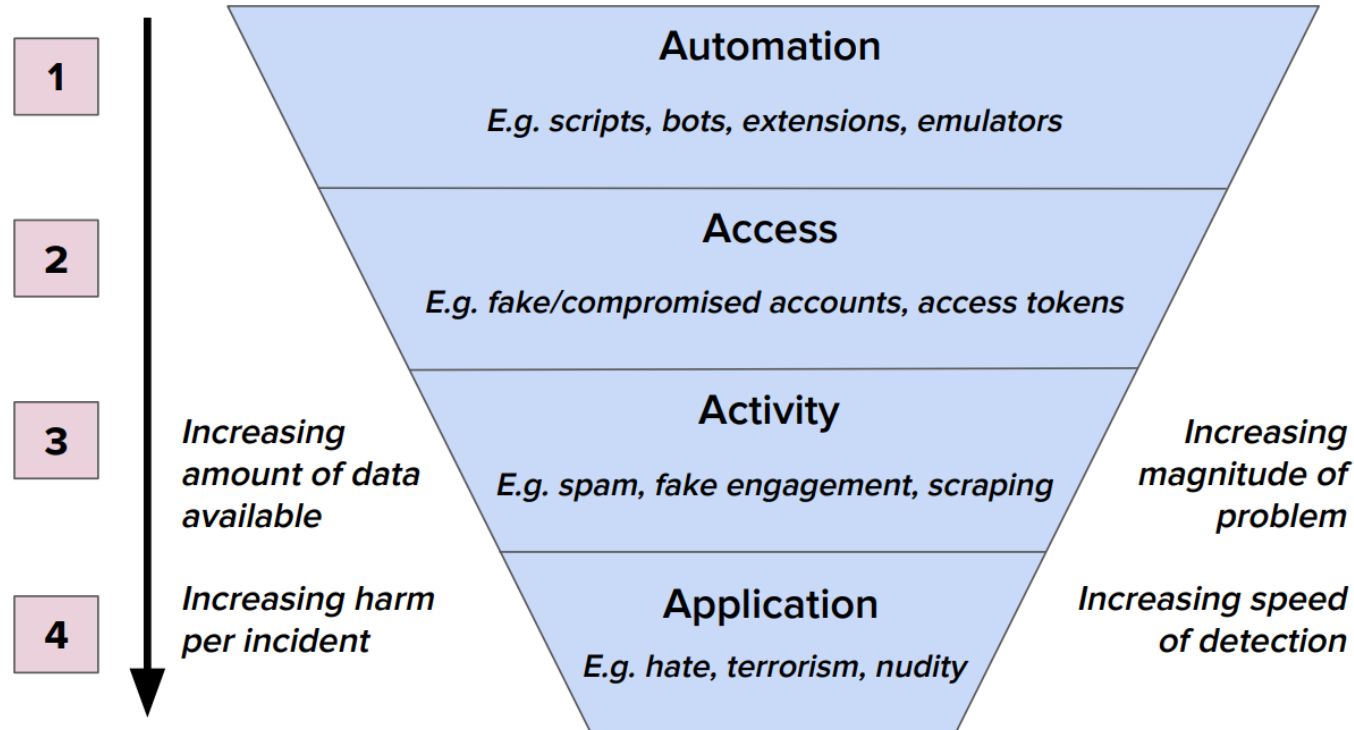
- This is the architecture of the ML-based spam detection system at **Facebook**



Machine Learning Systems (Case Study)



- This is the architecture of the ML-based spam detection system at **Facebook**



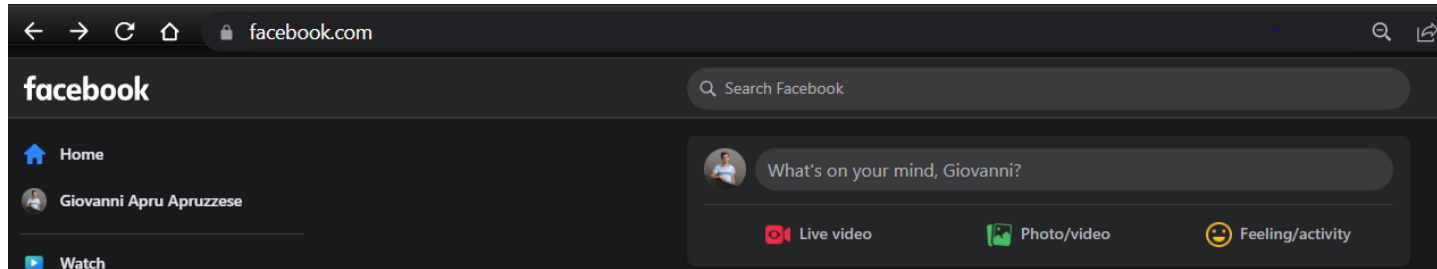
- The first layers are meant to block attacks *at scale* (e.g., query-based strategies)
- All layers use a mix of ML and non-ML techniques (not necessarily deep learning)
- Deep learning really shines at the bottom layer (few events reach this layer, though)
- The output accounts for diverse layers and is not instantaneous (an *invisible* ML system)

Real attackers have to bypass all layers to be successful.

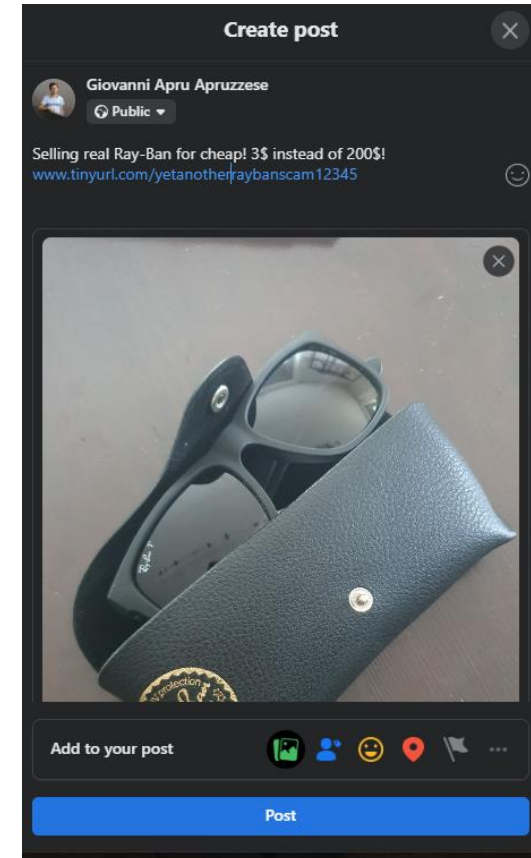
This does not mean
that this ML system
is omnipotent!

“Attacking” an *invisible* ML system

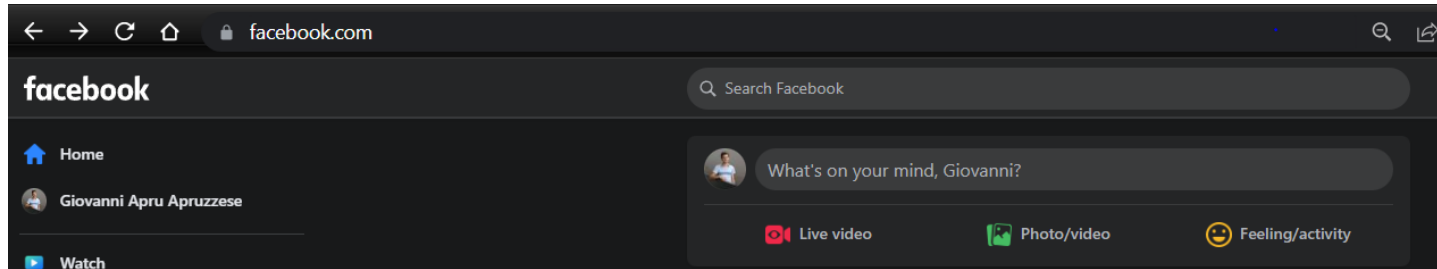
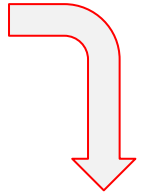
- If I go on Facebook and want to spread “spammy” content...



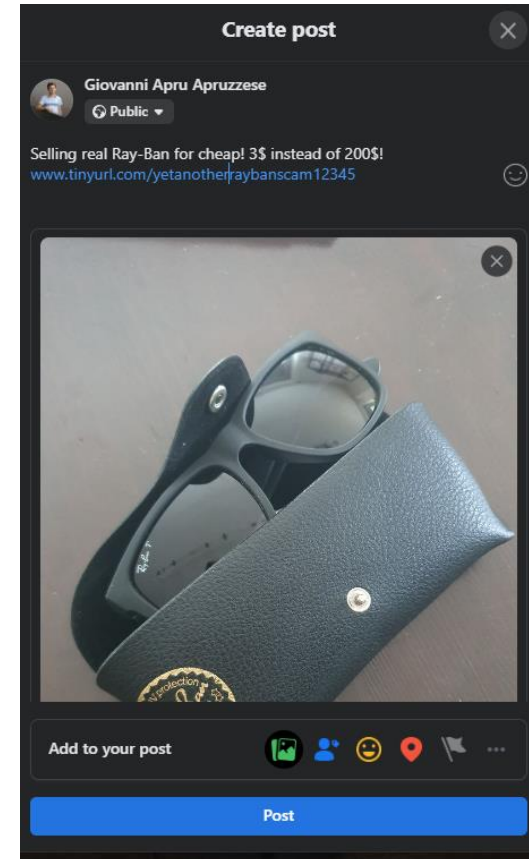
- ...the only thing I will see after “posting” it is the post itself.



“Attacking” an *invisible* ML system (cont’d)



- If I go on Facebook and want to spread “spammy” content...
- ...the only thing I will see after “posting” it is the post itself.
- I would not be able to see:
 - The architecture of Facebook’s spam detector
 - The fact that it uses ML
 - The fact that my specific post was (or not) analyzed by ML
 - The output of the system to my specific post
- If the post “appears”, does it mean that the system was evaded?
 - What if the post gets removed after 1 hour? Or 1 day?
 - What if my account is blocked after 1 week?



Machine Learning Systems (state-of-research)

- We analyzed all related papers accepted at top-4 cybersecurity conferences (NDSS, S&P, CCS, USENIX Sec) from 2019-2021.
 - Out of 1549 papers, 88 fell into the “adversarial ML” category.
 - Out of these, 78 consider *only* deep learning methods

Machine Learning Systems (state-of-research)

- We analyzed all related papers accepted at top-4 cybersecurity conferences (NDSS, S&P, CCS, USENIX Sec) from 2019-2021.
 - Out of 1549 papers, 88 fell into the “adversarial ML” category.
 - Out of these, 78 consider *only* deep learning methods

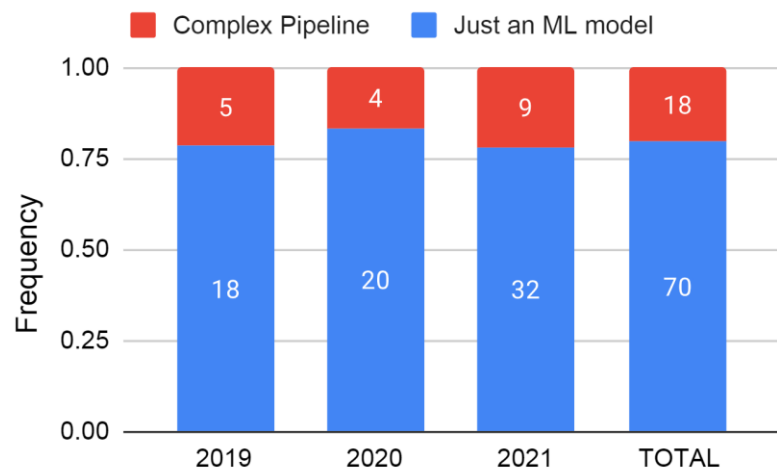


Fig. 12: Has a complex *pipeline* been reproduced in the evaluation?

Building a pipeline that resembles a (realistic) ML system is difficult.

These assets are not publicly available!

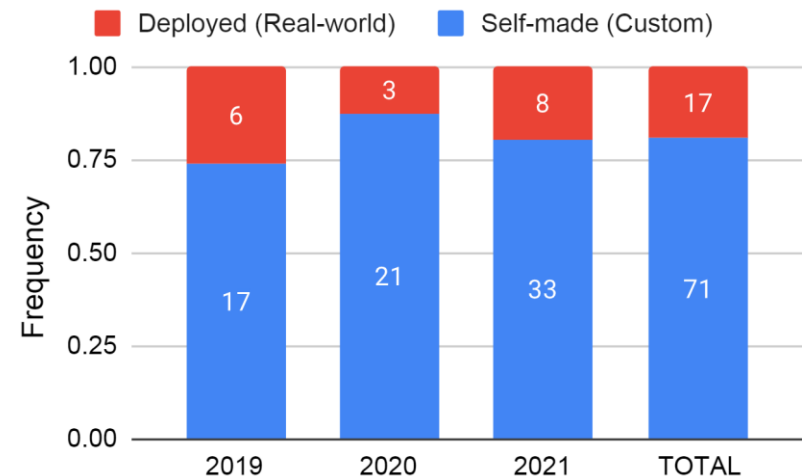


Fig. 13: Does the paper consider an ML model *deployed* in the real world?

Finding a ML system that is openly available for research-focused (security) assessments is hard.

Getting in touch with companies is tough!

Disclaimer: the findings of all these papers are still significant!

Cybersecurity is rooted in *economics*

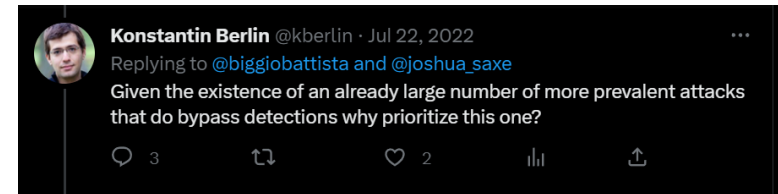
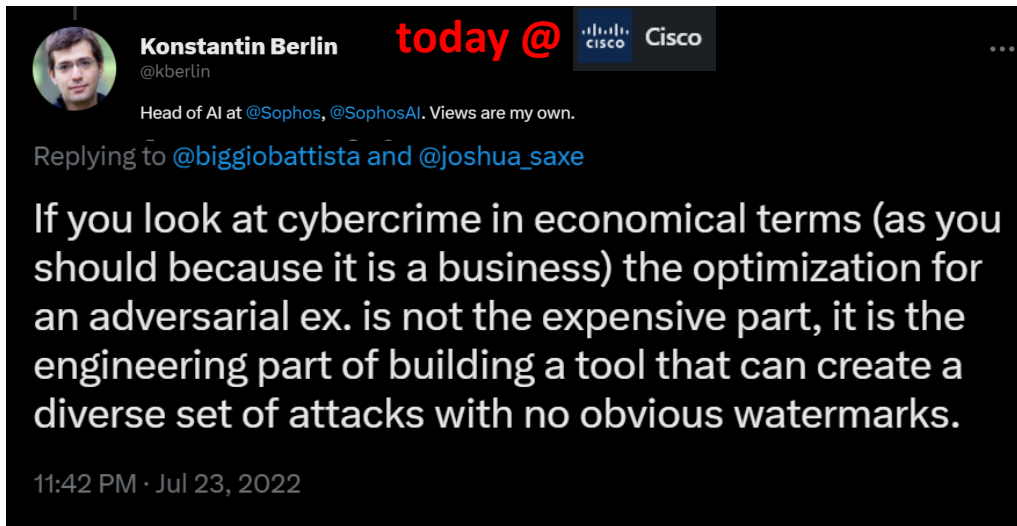
Cybersecurity ↔ Economics

- Given enough resources, any attack will be successful
- The goal of a defense is to “raise the bar” for the attacker

“There is no such a thing
as a foolproof system.”

→ A real attacker will opt for the **cheaper** strategy to reach their objective

→ A real defender will prioritize the **most likely** threats.



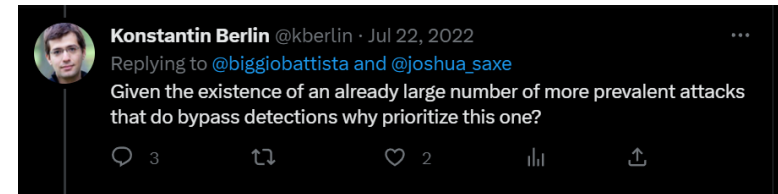
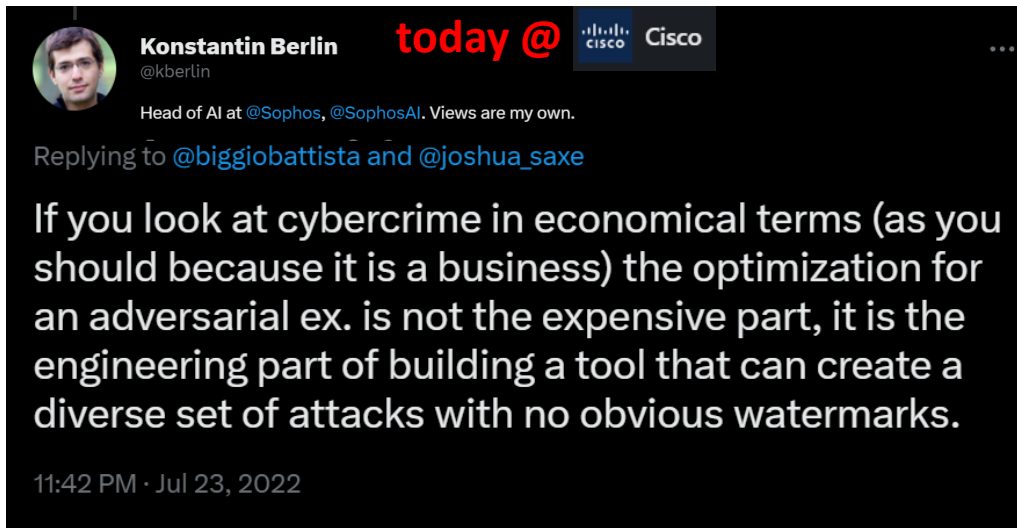
Cybersecurity ↔ Economics

"There is no such a thing
as a foolproof system."

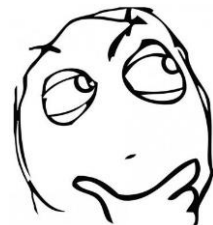
- Given enough resources, any attack will be successful
- The goal of a defense is to “raise the bar” for the attacker

→ A real attacker will opt for the **cheaper** strategy to reach their objective

→ A real defender will prioritize the **most likely** threats.



- In our domain, the **cost** of an attack is typically measured by means of “queries”
 - More queries → higher cost → “less effective” attack

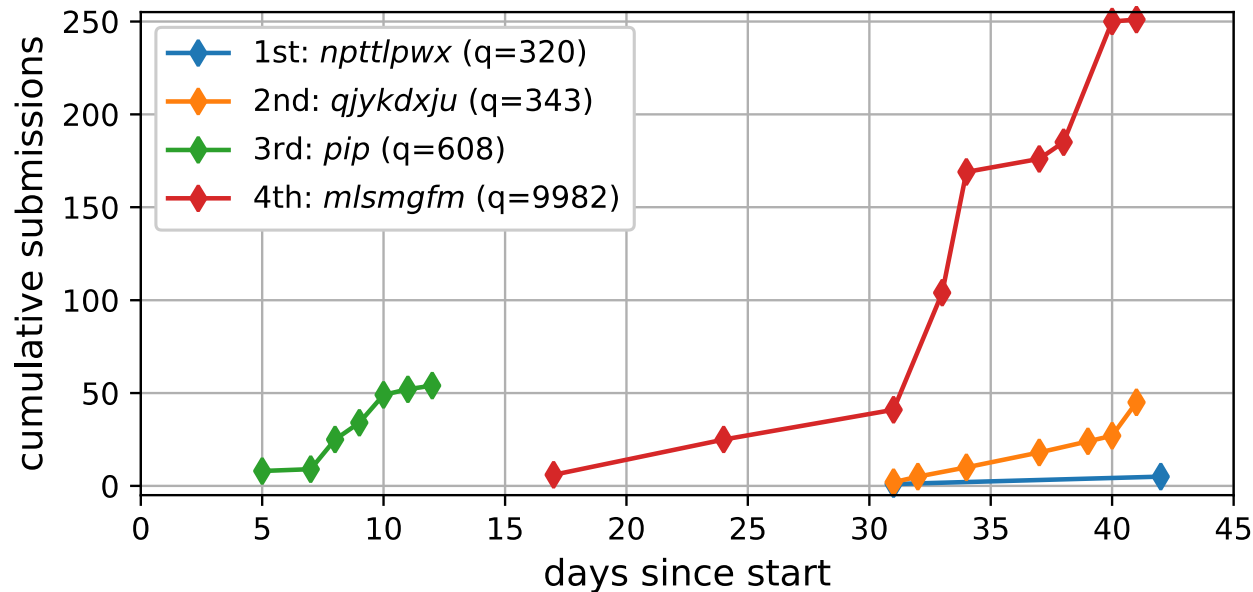


Cybersecurity \Leftrightarrow Economics (Case Study)

- We performed an in-depth look at the MLSEC anti-phishing challenge of 2021
 - Participants had to “evade the black-box detector” with as few queries as possible

Cybersecurity \leftrightarrow Economics (Case Study)

- We performed an in-depth look at the MLSEC anti-phishing challenge of 2021
 - Participants had to “evade the black-box detector” with as few queries as possible



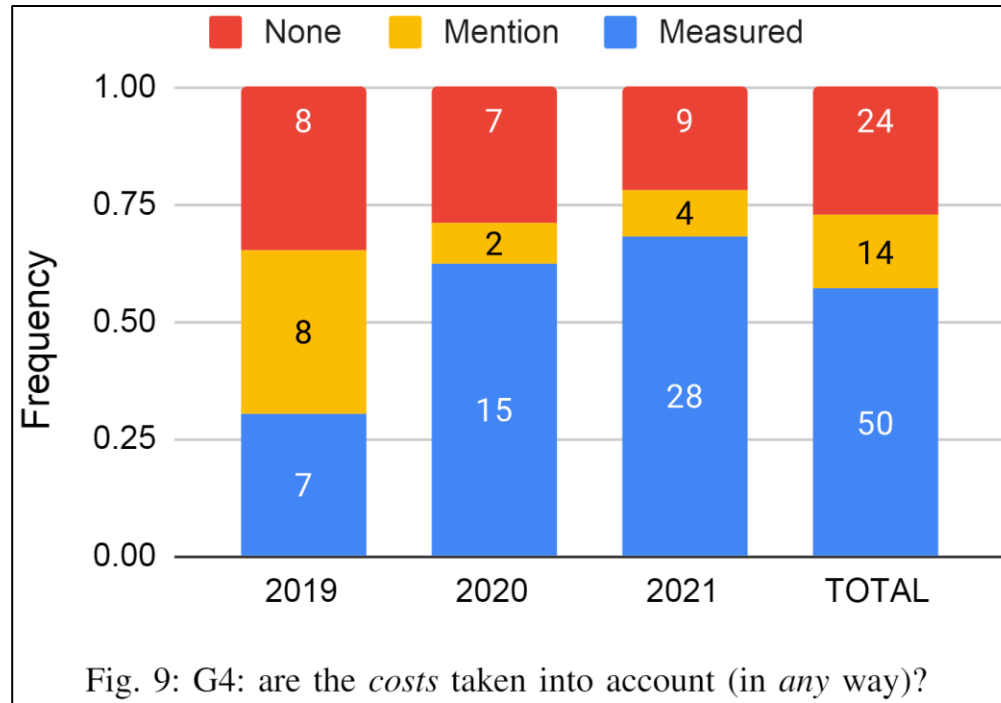
- The team arriving first (320 queries)... was **the last** to submit their solution
- The team arriving third (608 queries)... was **the first** to submit their solution
- Both of these teams only relied on their **domain expertise**

Queries do not tell
the whole story!

No gradient was
computed here!

Cybersecurity \leftrightarrow Economics (state-of-research)

- Do research papers on adversarial ML take economics into account?



Positive trend!

- Only 3 papers provided an *actual cost* in \$\$ (but only for “expenses”)
- The measurements never considered the *human factor*
 - Attack papers measured “queries”, defense papers measured “performance degradation”

At least in the adversarial ML domain, economics appears to be overlooked.

Objectively measuring
the human factor is hard!

A few words on the state-of-research

Data and Reproducibility (state-of-research)

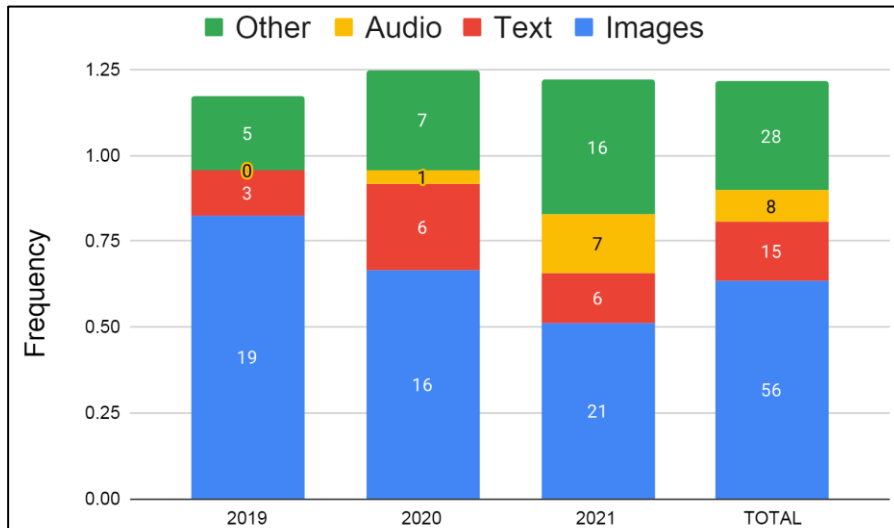


Fig. 10: What are the *data-types* considered in the evaluation?

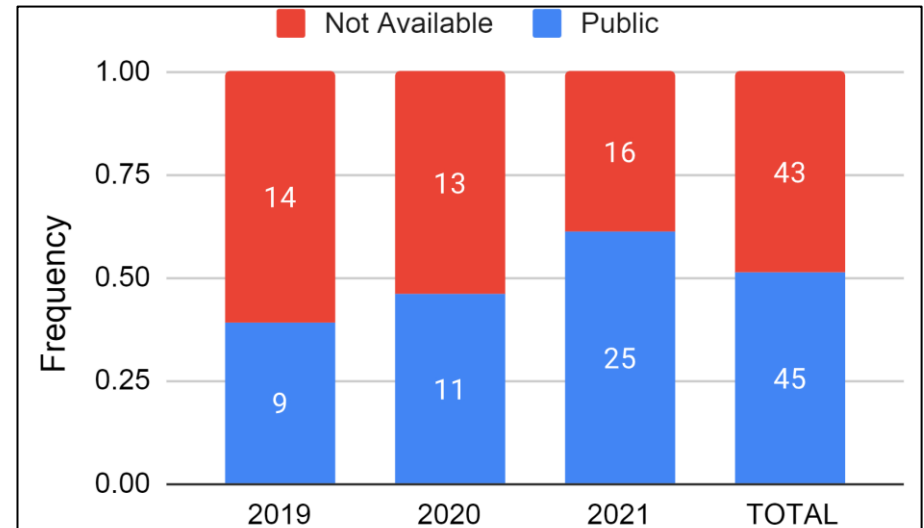


Fig. 11: Has the *source-code* been publicly released?

- Over 50% of the papers focus on image data (decreasing trend)
 - Only 12 papers (out of 88) focus on ML applications for cybersecurity (e.g., phishing, malware)

Some ML application domains (e.g., finance) are rarely discussed in adversarial ML literature.

- Only 50% of the papers release their implementations publicly (increasing trend)

In cybersecurity conferences!

Inconsistent Terminology (“What does the attacker know?”)

- The terms “white-box” and “black-box” are widespread, but often denote different degrees of attacker’s knowledge. Here are some examples, taken verbatim.

Co et al. [101]: “In **white-box** settings, the adversary has complete knowledge of the model architecture, parameters, and training data.[...] In a **black-box** setting, the adversary has no knowledge of the target model and no access to surrogate datasets.”

Aligns with Srndic and Laskov [43]

Inconsistent Terminology (“What does the attacker know?”)

- The terms “white-box” and “black-box” are widespread, but often denote different degrees of attacker’s knowledge. Here are some examples, taken verbatim.

Co et al. [101]: “In **white-box** settings, the adversary has complete knowledge of the model architecture, parameters, and training data.[...] In a **black-box** setting, the adversary has no knowledge of the target model and no access to surrogate datasets.”

Aligns with Srndic and Laskov [43]

Shan et al. [102]: “We assume a basic **white box** threat model, where adversaries have direct access to the the ML model, its architecture, and its internal parameter values [...] but do not have access to the training data.”

...this is different from [101] (“white-box”)!

Inconsistent Terminology (“What does the attacker know?”)

- The terms “white-box” and “black-box” are widespread, but often denote different degrees of attacker’s knowledge. Here are some examples, taken verbatim.

Co et al. [101]: “In **white-box** settings, the adversary has complete knowledge of the model architecture, parameters, and training data.[...] In a **black-box** setting, the adversary has no knowledge of the target model and no access to surrogate datasets.”

Aligns with Srndic and Laskov [43]

Xiao et al. [22]: “In this paper, we focus on the **white-box** adversarial attack, which means we need to access the target model (including its structure and parameters).”

...what about the training data?

Inconsistent Terminology (“What does the attacker know?”)

- The terms “white-box” and “black-box” are widespread, but often denote different degrees of attacker’s knowledge. Here are some examples, taken verbatim.

Co et al. [101]: “In **white-box** settings, the adversary has complete knowledge of the model architecture, parameters, and training data.[...] In a **black-box** setting, the adversary has no knowledge of the target model and no access to surrogate datasets.”

Aligns with Srndic and Laskov [43]

Suya et al. [103] assume a “**black-box**” attacker that “does not have direct access to the target model or knowledge of its parameters,” but that “has access to pre-trained local models for the same task as the target model” which could be “directly available or produced from access to similar training data.”

...this is different from [101] (“black-box”)!

Inconsistent Terminology (“What does the attacker know?”)

- The terms “white-box” and “black-box” are widespread, but often denote different degrees of attacker’s knowledge. Here are some examples, taken verbatim.

Shan et al. [102]: “We assume a basic **white box** threat model, where adversaries have direct access to the the ML model, its architecture, and its internal parameter values [...] but do not have access to the training data.”

Hui et al. [104] envision a “**gray-box**” setting which “gives full knowledge to the adversary in terms of the model details. Specifically, except for the training data, the adversary knows almost everything about the model, such as the architecture and the hyper-parameters used for training.”

This is the exact same as [102]... which describes a “white-box” setting!

Inconsistent Terminology (“What does the attacker know?”)

- The terms “white-box” and “black-box” are widespread, but often denote different degrees of attacker’s knowledge. Here are some examples, taken verbatim.

Co et al. [101]: “In **white-box** settings, the adversary has complete knowledge of the model architecture, parameters, and training data.[...] In a **black-box** setting, the adversary has no knowledge of the target model and no access to surrogate datasets.”

Aligns with Srndic and Laskov [43]

Shan et al. [102]: “We assume a basic **white box** threat model, where adversaries have direct access to the the ML model, its architecture, and its internal parameter values [...] but do not have access to the training data.”

...this is different from [101] (“white-box”)!

Xiao et al. [22]: “In this paper, we focus on the **white-box** adversarial attack, which means we need to access the target model (including its structure and parameters).”

...what about the training data?

Suya et al. [103] assume a “**black-box**” attacker that “does not have direct access to the target model or knowledge of its parameters,” but that “has access to pre-trained local models for the same task as the target model” which could be “directly available or produced from access to similar training data.”

...this is different from [101] (“black-box”)!

Hui et al. [104] envision a “**gray-box**” setting which “gives full knowledge to the adversary in terms of the model details. Specifically, except for the training data, the adversary knows almost everything about the model, such as the architecture and the hyper-parameters used for training.”

This is the exact same as [102]... which describes a “white-box” setting!

Taken individually, all past work are correct. The problems arise when analyzing the situation **as a whole!**

Our four Positions

P1: Adapt threat models to ML systems

Attacker's **Goal, Knowledge, Capabilities** and **Strategy** should reflect the ML system (and not just the ML model!)

→ Real attackers have **broader objectives** and do not want just to “evade the ML model.”

Each of those elements should be **precisely defined**.

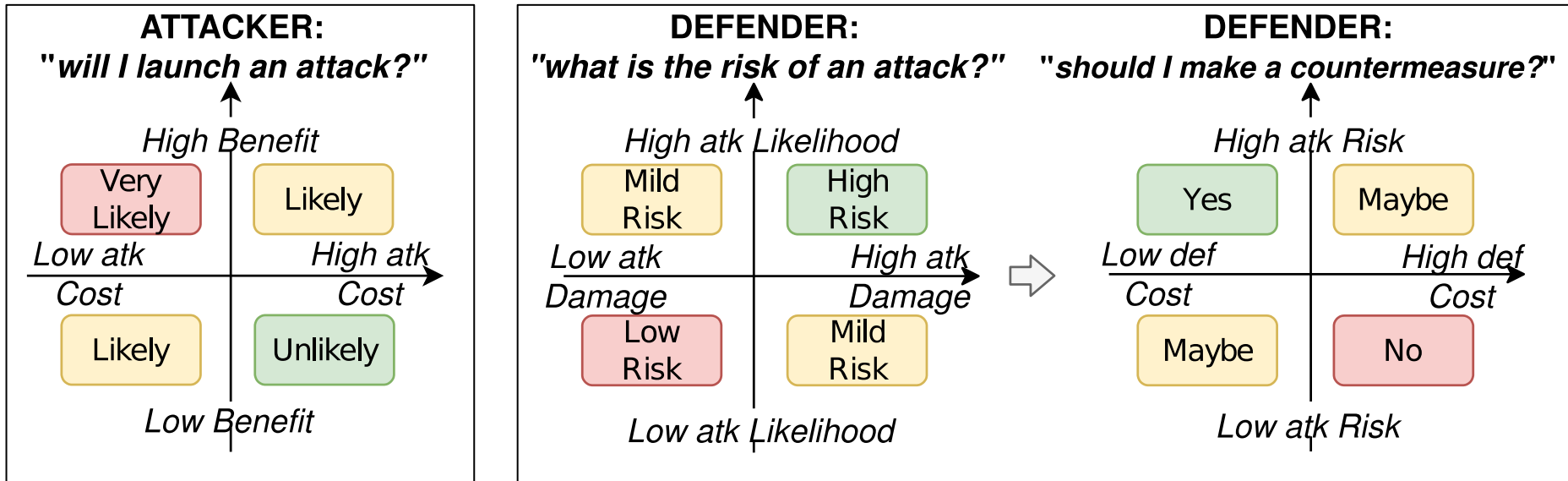
→ Existing **terminology** is often used inconsistently.

Problematic Terms:

- “Box-based” terminology
- “Access”
- “Adversarial”
- “Evasion”

*Solutions and
recommendations in
the paper!*

P2: Cost-based threat modeling



Both attacks and defenses have a **cost**. Real attackers do not launch an attack if it is *too expensive*; and real developers will not develop a countermeasure if the attack is *unlikely to occur in reality*.

→ Cost measurements should account for the **human factor** (queries / computation are not enough)

More on this in the paper!

→ There is value also in defenses that work "only" against attackers with **limited knowledge** (they are more common in reality).

P3: Collaborations between *industry* and *academia*

Practitioners should be **more willing** to cooperate with researchers: both have the same goal!

- 💡 Streamline research collaboration process
- 💡 Bug Bounties
- 💡 Releasing Schematics

P4: *Source-code* disclosure with “just culture”

Just Culture: assumes that mistakes are bound to occur and derive from organizational issues. Mistakes are avoided by understanding their root causes and using them as constructive learning experiences.

Embracing a just culture naturally promotes the **gradual improvement** at the base of research efforts.

→ The fast pace of research in ML can lead to errors in experiments (not always spotted during the peer-review)

→ By releasing the source code, future works can correct such mistakes, potentially systematizing them, and hence **turning “negative results” into positive outcomes** for our community.

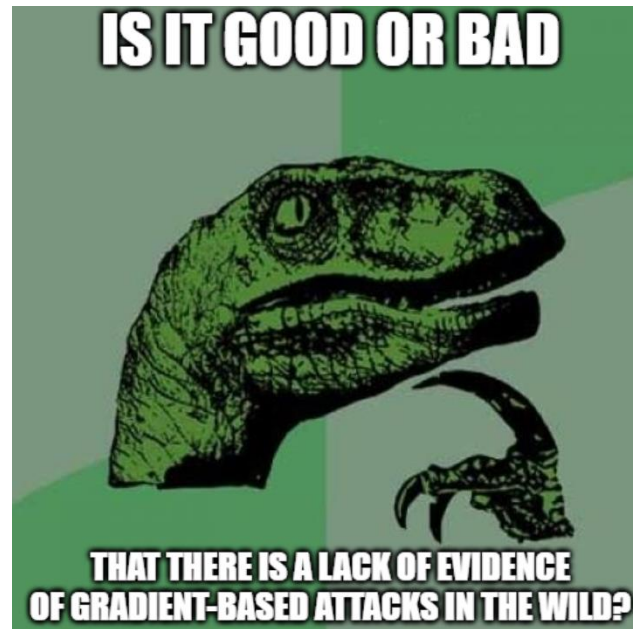
State-of-research [bonus]

TABLE IV: The 88 papers considered in our analysis. Each column reports the answer to one of the 12 research questions we used during our survey available, the G6 column provides the *hyperlink* to the websites hosting the source-code of a given paper. Explanations are in Appendix B-I.

Year (subs)	Venue (subs)	Paper (1st author)	G1	G2	G3	G4	G5 (Evaluation Data)				G6	G7	G8	T1	T2	T3	T4
			Focus	Attack	Paradigm	Cost	Img	Text	Audio	Other	Code	Pipeline	Type	Param.	Sem.	Output	Training
2019 (20/00)	NDSS (000)	Salem [158]	atk	Member.	DL	●	✓	✓			✓		CLOSED	✓	✓	p	R
		Li [138]	def	Evasion	DL	●	✓	✓			✓		CLOSED	✓	✓	p	R
		Ma [163]	def	Evasion	DL+SL	●	✓	✓			✓	✓	CLOSED	✓	✓	p	D
	SP (004)	Ling [152]	def	Poison.	DL	✓	✓	✓			✓			✓	✓	p	R
		Wang [164]	def	Poison.	DL	✓	✓	✓		Finance	✓			✓	✓	p	D
	SEC (011)	Nasr [100]	atk	Member.	DL	✓	✓	✓		Malware	✓			✓	✓	p	R
		Tong [114]	def	Evasion	DL+SL	✓	✓	✓		Malware	✓			✓	✓	p	S
		Demonitis [44]	atk	Evasion	DL+SL	✓	✓	✓		Malware	✓	✓	CLOSED	✓	✓	p	R
		Xiao [22]	atk	Evasion	DL	✓	✓	✓			✓			✓	✓	p	R
		Quiring [165]	atk	Evasion	DL+SL	✓	✓	✓			✓			✓	✓	p	R
2020 (24/00)	NDSS (000)	Hong [110]	atk	Evasion	DL	✓	✓	✓			✓			✓	✓	p	R
		Batima [111]	atk	Stealing	DL	●	✓	✓			✓			✓	✓	p	R
		Song [166]	atk	Member.	DL	✓	✓	✓		+	✓			✓	✓	p	R
		Jia [167]	def	Member.	DL	✓	✓	✓			✓			✓	✓	p	R
		Co [101]	atk	Evasion	DL	✓	✓	✓			✓			✓	✓	p	R
	CCS (00149)	Liu [168]	def	Poison.	DL	✓	✓	✓			✓			✓	✓	p	✓
		Baluta [169]	def	Poison.	DL	●	✓	✓			✓			✓	✓	p	✓
		Zhao [105]	atk	Evasion	DL	✓	✓	✓			✓			✓	✓	p	S
		Tramer [12]	atk	Evasion	DL	✓	✓	✓			✓	✓	CLOSED	✓	✓	p	R
		Wang [170]	atk	Evasion	DL	✓	✓	✓		Graphs	✓			✓	✓	p	C
		Yao [52]	atk	Poison.	DL	●	✓	✓			✓			✓	✓	p	C
		Yang [171]	atk	Stealing	DL	●	✓	✓			✓		CLOSED	✓	✓	p	D
2021 (41/044)	NDSS (000)	Aghakhani [142]	atk	Evasion*	DL+SL	✓	✓	✓		Malware	✓		CLOSED	✓	✓	p	D
		Yu [34]	atk	Stealing	DL	✓	✓	✓			✓			✓	✓	p	C
		Schuster [172]	atk	Poison.	DL	✓	✓	✓			✓			✓	✓	p	R
	SP (0006)	Pierazza [49]	atk	Evasion	SL	✓	✓	✓		Malware	✓	✓		✓	✓	p	R
		Chen [88]	def	Evasion	DL	✓	✓	✓			✓			✓	✓	p	C
	SEC (0127)	Jan [147]	atk	Evasion*	DL	✓	✓	✓		Network	✓			✓	✓	p	C
		Salem [146]	atk	Member.	DL	✓	✓	✓		Location	✓			✓	✓	p	S
		Chandrasekaran [148]	atk	Stealing	DL+SL	✓	✓	✓		+	✓			✓	✓	p	R
		Suya [103]	atk	Evasion	DL	✓	✓	✓			✓			✓	✓	p	S
		Jagielski [107]	atk	Stealing	DL	●	✓	✓			✓			✓	✓	p	D
2021 (41/044)	NDSS (007)	Quiring [173]	atk	Evasion	DL	✓	✓	✓			✓	✓		✓	✓	p	R
		Li [151]	def	Evasion	DL	✓	✓	✓			✓			✓	✓	p	D
		Leino [174]	atk	Member.	DL	✓	✓	✓		+	✓			✓	✓	p	R
		Zhang [175]	atk	Evasion	DL	✓	✓	✓			✓		CLOSED	✓	✓	p	R
		Nasr [25]	atk	Evasion	DL	✓	✓	✓			✓	✓		✓	✓	p	S
	CCS (00121)	Li [176]	atk	Evasion	DL	✓	✓	✓			✓			✓	✓	p	C
		Shan [102]	def	Evasion	DL	●	✓	✓			✓			✓	✓	p	R
		Pang [162]	atk	Poison.	DL	✓	✓	✓			✓			✓	✓	p	R
		Abdelnabi [177]	atk	Evasion*	DL	✓	✓	✓		Phishing	✓	✓		✓	✓	p	R
		Li [178]	atk	Evasion	DL	✓	✓	✓			✓			✓	✓	p	R
		Lin [179]	atk	Poison.	DL	✓	✓	✓			✓			✓	✓	p	R
		Chen [180]	atk	Member.	DL	✓	✓	✓			✓			✓	✓	p	D
2021 (41/044)	NDSS (007)	Zanella [24]	atk	Member.	DL	✓	✓	✓			✓			✓	✓	p	R
		Song [23]	atk	Member.	DL	✓	✓	✓			✓			✓	✓	p	S
		Hui [104]	atk	Member.	DL	✓	✓	✓		+	✓			✓	✓	p	R
	SP (011)	Huang [181]	atk	Poison.	DL	✓	✓	✓		Ratings	✓			✓	✓	p	C
		Barradas [90]	atk	Evasion*	SL	✓	✓	✓		Network	✓			✓	✓	p	R
		Xu [170]	def	Poison.	DL	✓	✓	✓			✓			✓	✓	p	R
		Abdelnabi [182]	atk	Evasion	DL	●	✓	✓			✓	✓		✓	✓	p	R
		Chen [91]	atk	Evasion	DL	✓	✓	✓			✓	✓	(both)	✓	✓	p	S
2021 (41/044)	SEC (01246)	Abdullah [145]	atk	Evasion	DL	✓	✓	✓			✓	✓	CLOSED	✓	✓	p	R
		Nasr [35]	def	Evasion	DL	✓	✓	✓		Finance	✓			✓	✓	p	R
		Sato [28]	atk	Evasion	DL	✓	✓	✓			✓	✓	OPEN	✓	✓	p	R
		Nasr [89]	atk	Evasion	DL	✓	✓	✓		Network	✓	✓		✓	✓	p	D
		He [183]	atk	Member.	DL	●	✓	✓		Graph	✓			✓	✓	p	R
	CCS (0196)	Severi [184]	atk	Poison.	DL+SL	✓	✓	✓		Malware	✓			✓	✓	p	✓
		Bagdasaryan [95]	atk	Poison.	DL	✓	✓	✓			✓			✓	✓	p	R
		Xi [155]	atk	Poison.	DL	✓	✓	✓		Graph	✓			✓	✓	p	S
		Tang [96]	def	Poison.	DL	✓	✓	✓			✓			✓	✓	p	C
		Schuster [185]	atk	Poison.	DL	✓	✓	✓			✓		OPEN	✓	✓	p	C
		Carlini [54]	atk	Poison.	DL	✓	✓	✓			✓			✓	✓	p	C
		Vicarte [186]	atk	Poison.	DL	✓	✓	✓			✓			✓	✓	p	R
2021 (41/044)	SEC (01246)	Lovisotto [150]	atk	Evasion	DL	✓	✓	✓			✓		OPEN	✓	✓	p	C
		Carlini [30]	atk	Member.	DL	✓	✓	✓			✓	✓	OPEN	✓	✓	p	R
		Han [97]	atk	Evasion*	DL	✓	✓	✓		Graph	✓			✓	✓	p	R
		Eisenhofer [153]	def	Evasion	DL	✓	✓	✓			✓	✓	OPEN	✓	✓	p	R
		Wu [156]	atk	Poison.	DL	✓	✓	✓		Games	✓			✓	✓	p	R
	CCS (0196)	He [187]	atk	Stealing	DL	✓	✓	✓			✓			✓	✓	p	S
		Rakin [112]	atk	Evasion	DL	✓	✓	✓			✓			✓	✓	p	R
		Jia [188]	def	Stealing	DL	✓	✓	✓			✓			✓	✓	p	C
		Zhu [189]	def	Stealing	DL	✓	✓	✓			✓			✓	✓	p	R
		Xiang [190]	def	Evasion	DL	✓	✓	✓			✓			✓	✓	p	R
		Lin [191]	atk	Evasion *	DL	✓	✓	✓		Phishing	✓	✓		✓	✓	p	R
		Aziz [192]	def	Poison.	DL	✓	✓	✓			✓			✓	✓	p	R
2021 (41/044)	NDSS (007)	Hussain [93]	def	Evasion	DL	✓	✓	✓			✓			✓	✓	p	R
		Song [99]	def	Member.	DL	✓	✓	✓		+	✓			✓	✓	p	R
		Zheng [92]	atk	Evasion	DL	●	✓	✓			✓		CLOSED	✓	✓	p	R
		Mu [93]	atk	Evasion	DL	✓	✓	✓			✓			✓	✓	p	R
		Bahrāmali [194]	atk	Evasion	DL	✓	✓	✓		Graphs	✓			✓	✓	p	S
	CCS (0196)	Sheatsley [157]	atk	Evasion	DL	✓	✓	✓		Network	✓			✓	✓	p	R
		Du [94]	def	Evasion	DL	✓	✓	✓			✓			✓	✓	p	R
		Li [195]	def	Evasion	DL	✓	✓	✓			✓			✓	✓	p	R
		He [196]	def	Member.	DL	✓	✓	✓			✓			✓	✓	p	R
		Li [160]	atk	Member.	DL	✓	✓	✓			✓			✓	✓	p	R
		Chen [197]	def	Member.	DL+SL	✓	✓	✓		+	✓			✓	✓	p	R

Last page of
our paper

This leads to the
code repository!



Do real attackers compute gradients?

→ We cannot prove it ☹️ (yet).

Maybe they do!

“Real Attackers Don’t Compute Gradients”: Bridging the Gap between Adversarial ML Research and Practice



Please get his name right!
“Savino Dambra”