



Raleigh, NC, USA – February 8th, 2023

1st IEEE Conference on Secure and Trustworthy Machine Learning

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

Giovanni Apruzzese, Hyrum S. Anderson, Savino Dambra, David Freeman,
Fabio Pierazzi, Kevin Roundy



ROBUST
INTELLIGENCE



NortonLifeLock[™]



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



SCHLOSS DAGSTUHL
Leibniz-Zentrum für Informatik

- 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)



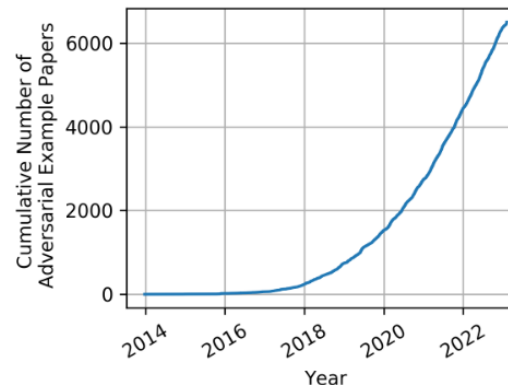
SCHLOSS DAGSTUHL
Leibniz-Zentrum für Informatik

- 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:

“Real attackers *guess*”

Backstory (Earth – July 22nd, 2022)

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

Dagstuhl follow-up: position paper on "attacker guessing" threat model?

Pierazzi, Fabio <fabio.pierazzi@kcl.>
To: dfreeman, Kevin Roundy, hyrum@robustintelligence.com
Cc: Apruzzese Giovanni
venerdì 22/07/2022 14:15

You forwarded this message on 02/09/2022 15:45.

Dear David, Kevin, Hyrum,

It was great to get to know you (more) during Dagstuhl.

I was talking with Giovanni yesterday, and were thinking again about what you all seemed to agree on from an industry perspective that in most cases attackers "guess" and do not necessarily use ML to evade systems, they just try to get out the easy way.

Given the upcoming first edition of [SATML](#), we saw there's also a category for "position papers", and me and Giovanni were thinking of maybe doing a position paper about "threat models of ML systems".

The current white-box threat models and also ML-driven black-box are mostly a worst-case scenario, and maybe models can be broken just much more easily (similar to the "pseudo-fuzzing" that Hyrum is looking into for ML models at Robust intelligence and maybe at Microsoft research).

Long story short, would you be interesting in co-authoring a position paper for SatML on the topic of "revisiting threat models of ML systems", to also re-define how to consider attacker capabilities in evading systems? Part of it is also related to the fact that real-world systems are a pipeline of ML and non-ML models.

Or, if not co-authoring, giving some feedback?

More concretely, there is some stuff that should be nice to highlight:

- In this mlsec challenge, authors evaded an ml classifier without ml: <https://cujo.com/announcing-the-winners-of-the-2021-machine-learning-security-evasion-competition/>
- In Giovanni&Pavel's 5G paper, they proposed the "myopic" threat model, similar to this issue: <https://arxiv.org/pdf/2207.01531.pdf>
- Konrad's team which won a defense in Hyrum's ML challenge got broken by a non-ML approach: <https://arxiv.org/pdf/2010.09569.pdf>

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.

Giovanni offered himself to do most of the work, so he should be able to lead the effort.

What do you think?

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!
4

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”

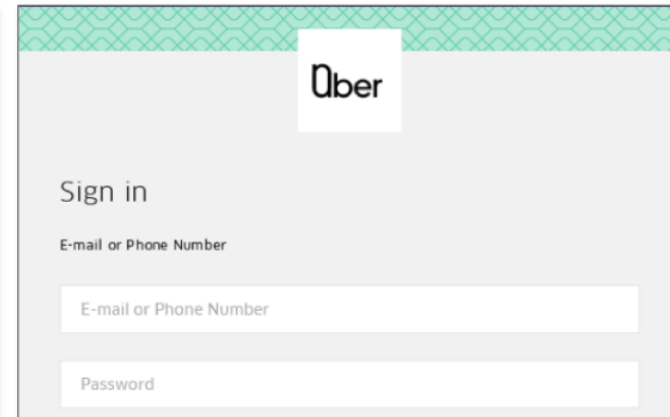
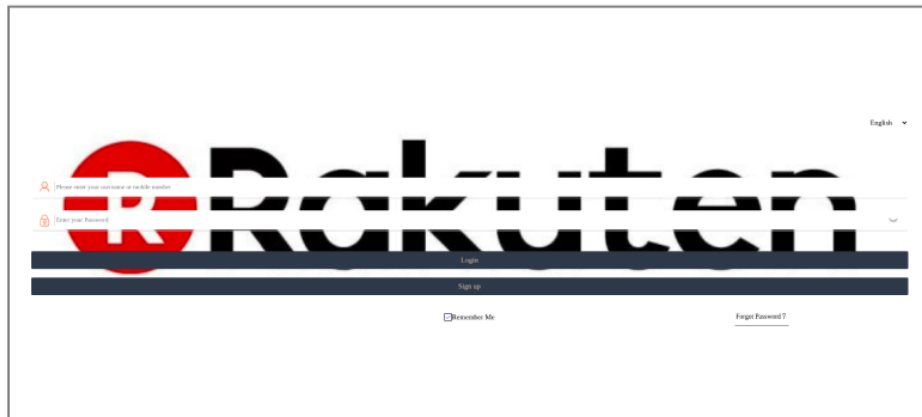
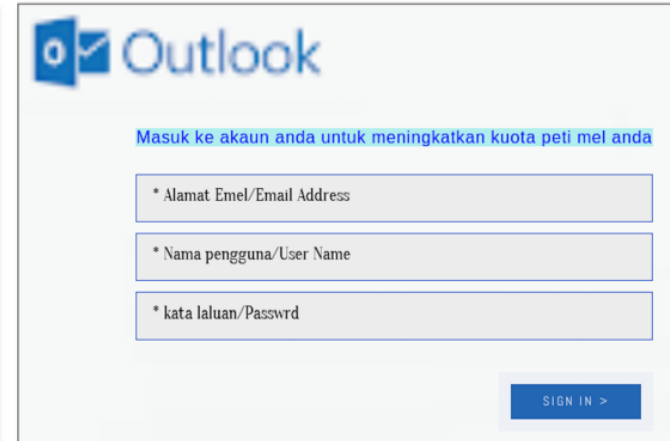
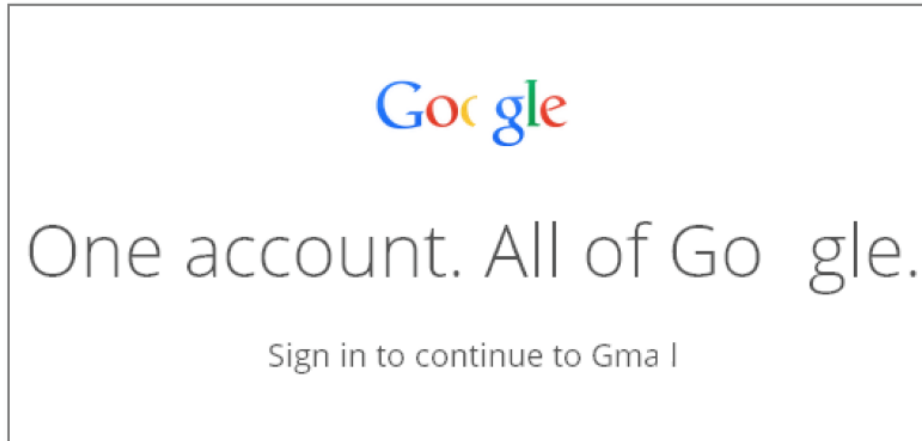
What did we find?

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]

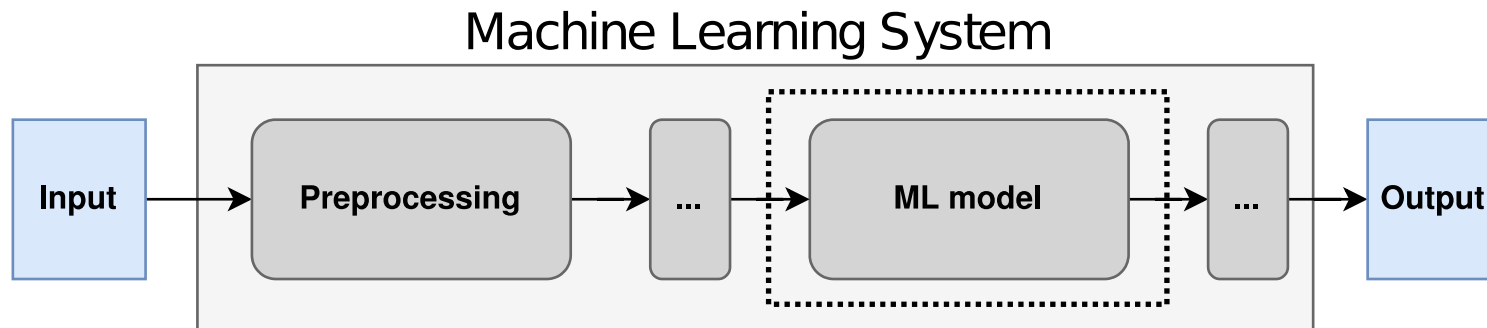
- The **vast majority** of these webpages were “out of distribution”
 - They were different from any sample in the training set
- We then looked at a small subset of the remaining ones...



Machine Learning Systems

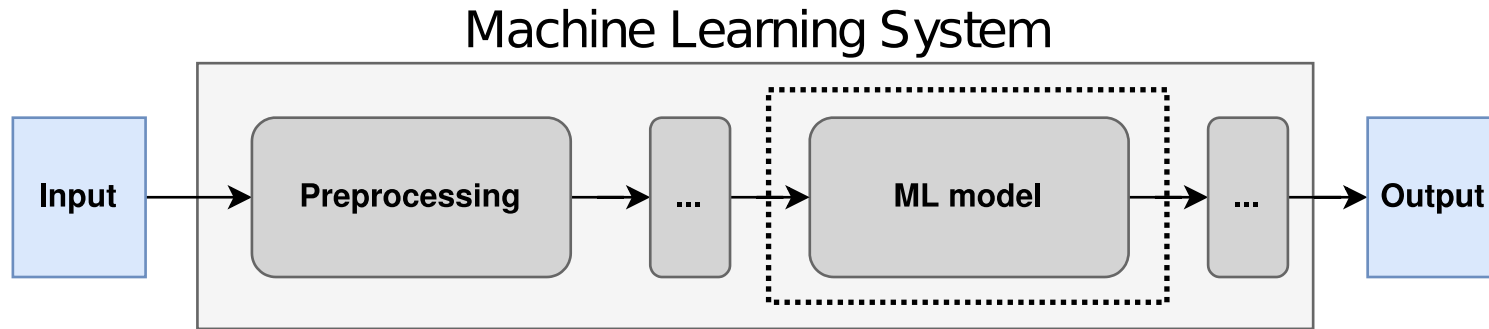
Machine Learning Systems

- 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*

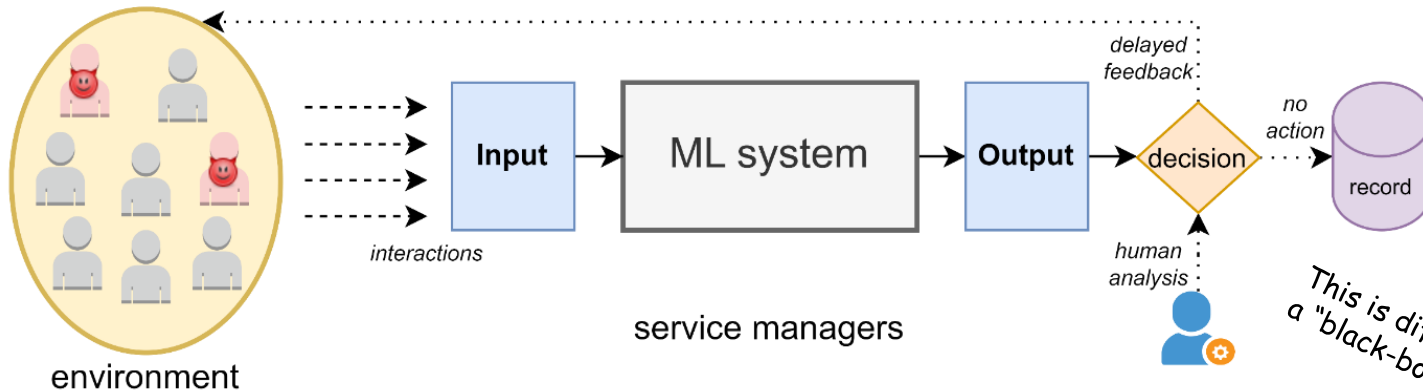


Machine Learning Systems

- 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*



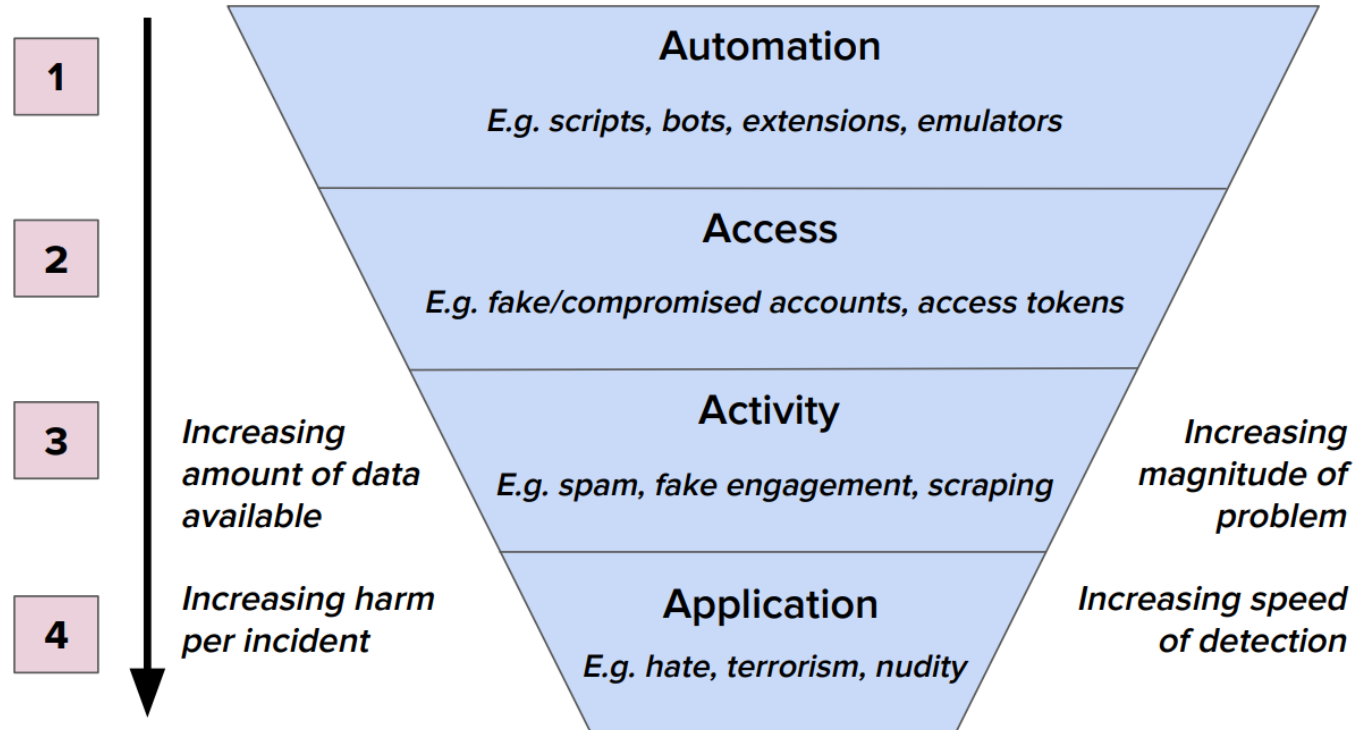
- Some ML systems are “invisible” to their users (and, hence, to real attackers)



This is different from a “black-box” scenario!

Machine Learning Systems (Case Study)

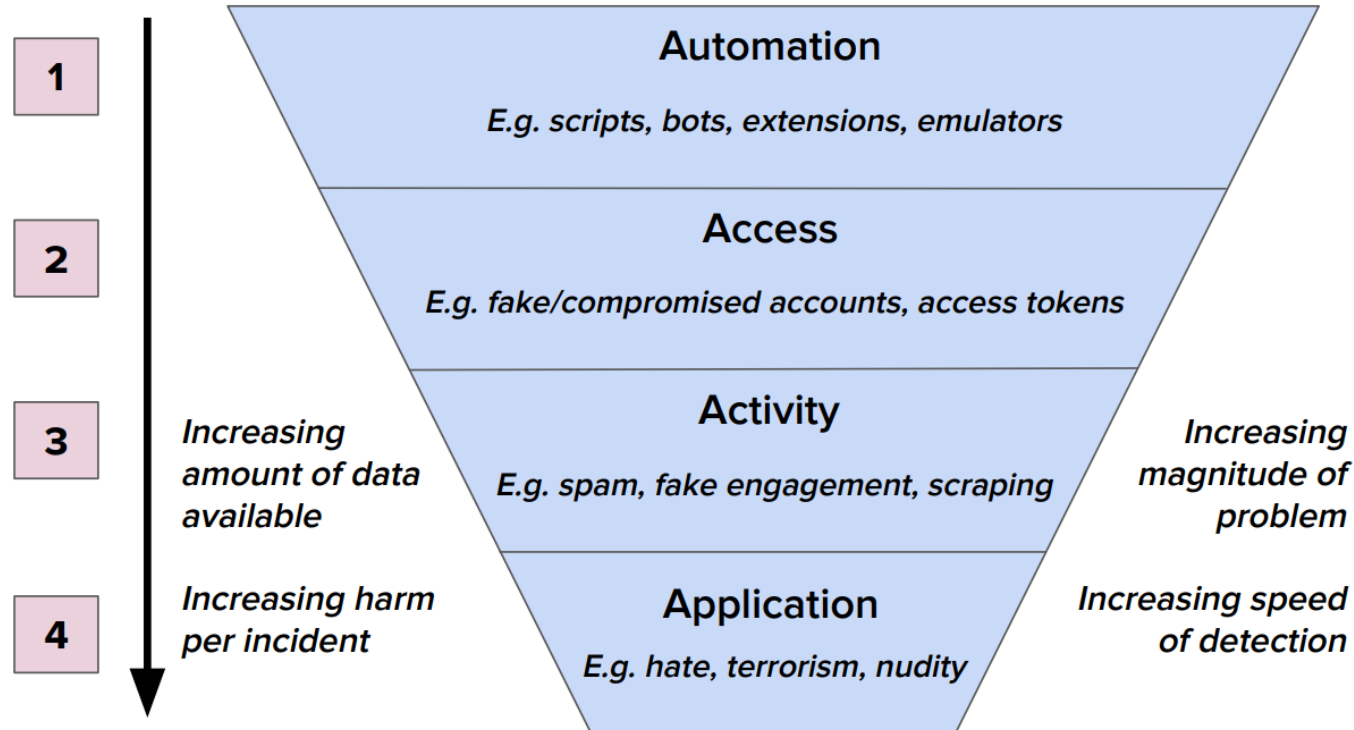
- This is the architecture of the ML-based spam detection system at **Facebook**



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!

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

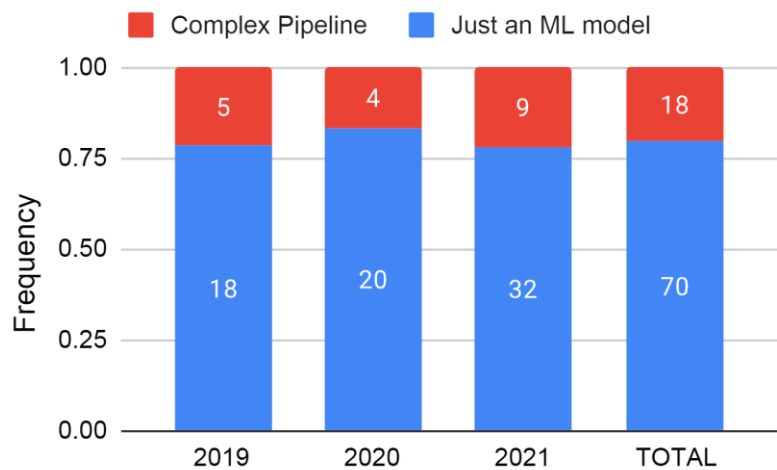


Fig. 12: Has a complex *pipeline* been reproduced in the evaluation?

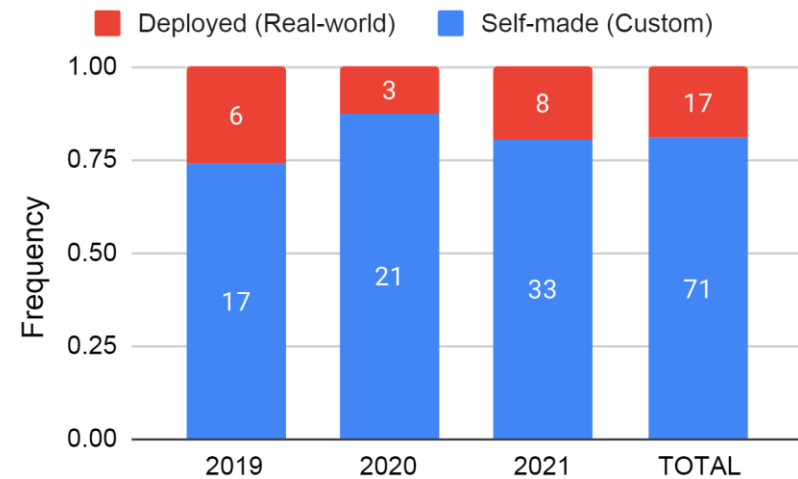


Fig. 13: Does the paper consider an ML model *deployed* in the real world?

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

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

These assets are not publicly available!

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.



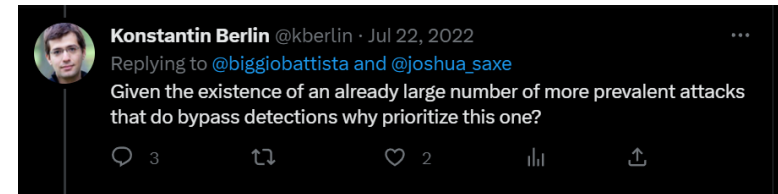
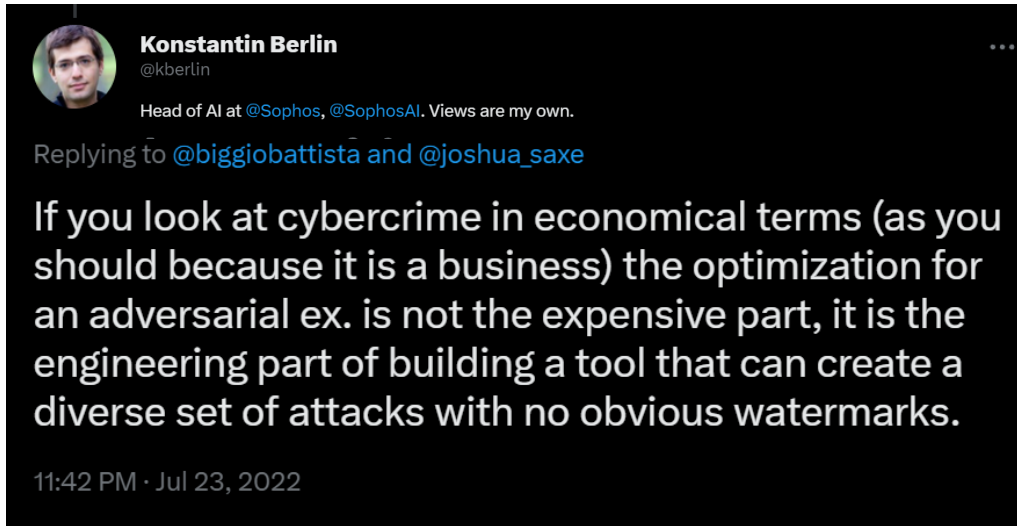
Cybersecurity ↔ Economics

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

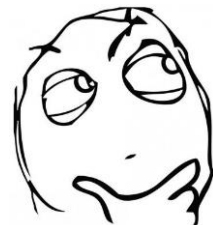
- Given enough resources, any attack will be successful
- The goal of a defense is to “raise the bar” for the attacker

→ 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

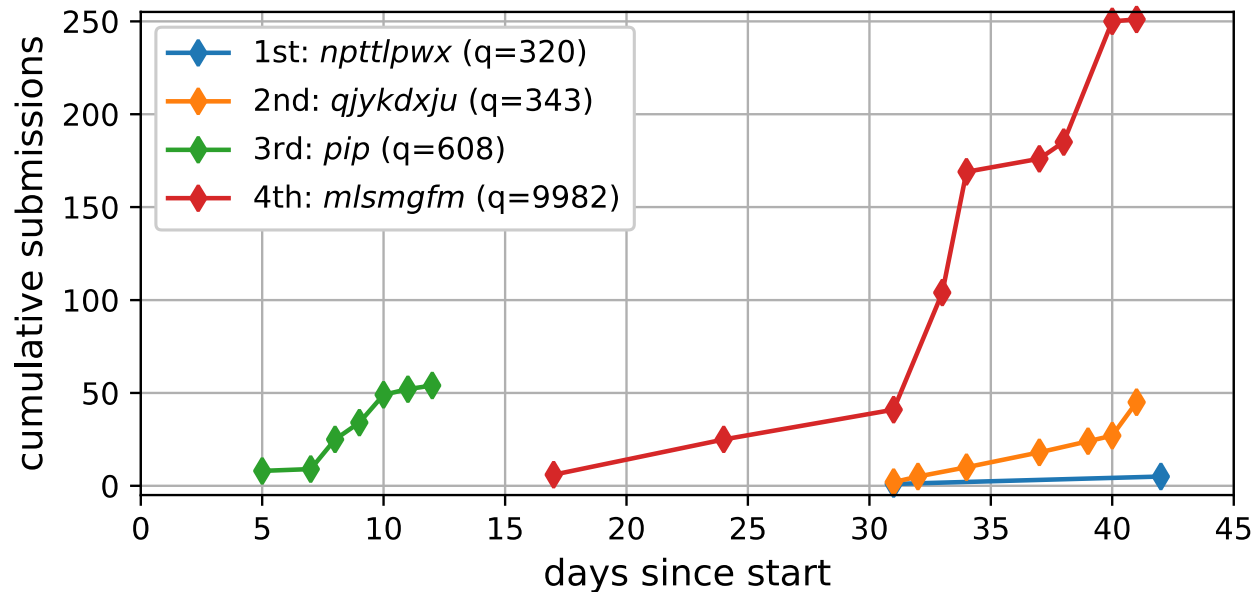


Cybersecurity ↔ Economics (Case Study)

- 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 ↔ 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



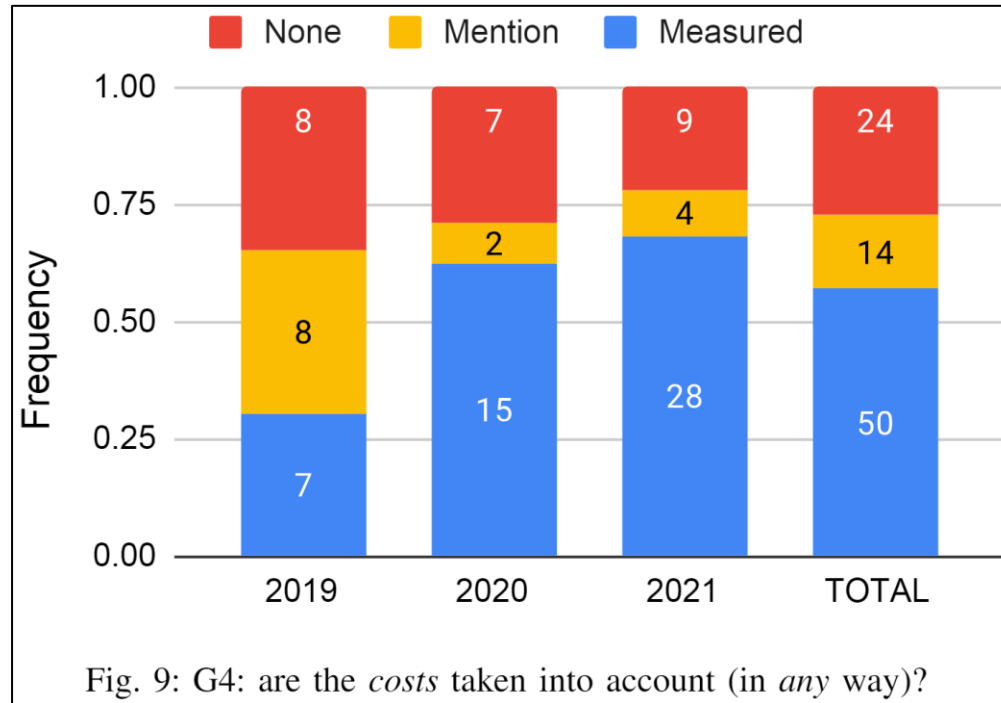
- The team arriving first (320 queries)... was **the last** to submit their solution
- The team arriving third (608 queries)... was **the first** to submit their solution
- Both of these teams only relied on their **domain expertise**

Queries do not tell the whole story!

No gradient was computed here!

Cybersecurity ↔ 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!

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!)

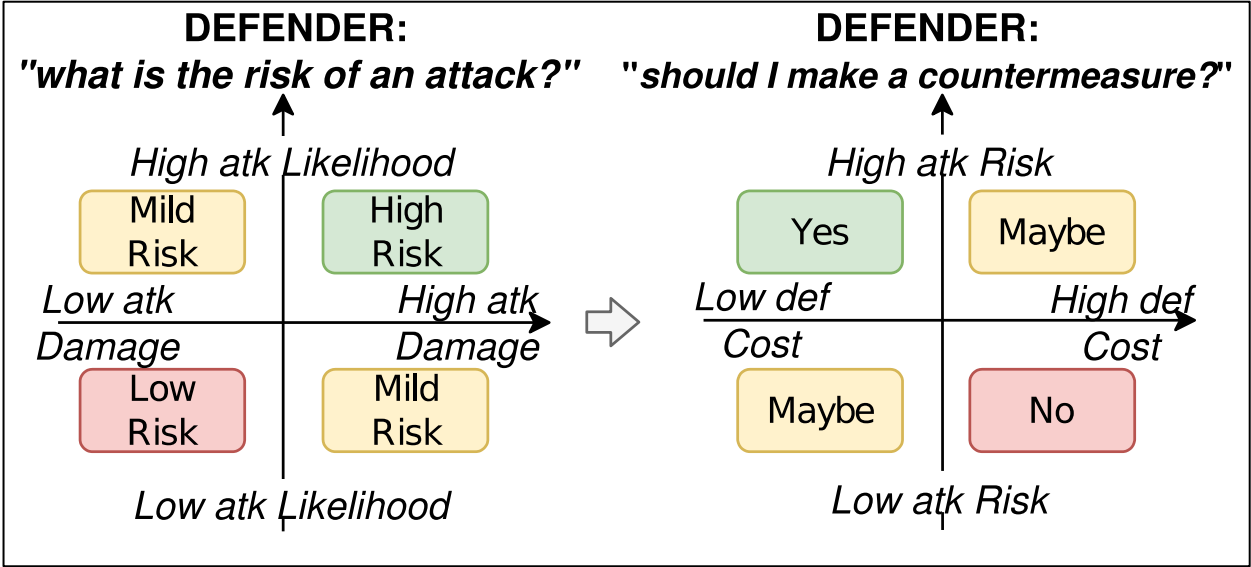
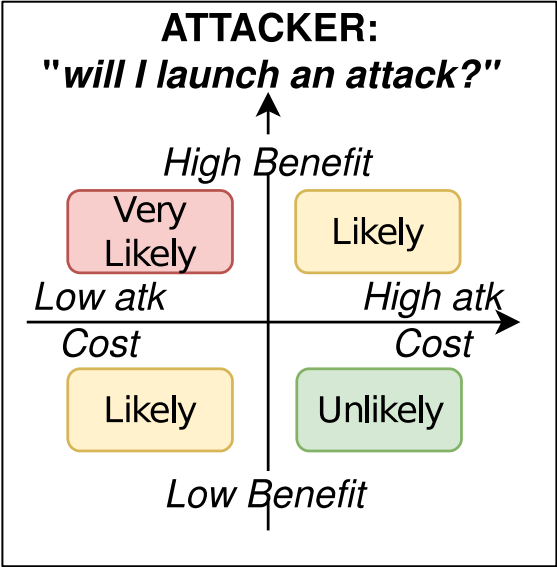
→ Real attackers have **broader objectives** and do not want just to “evade the ML model.”

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

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

*More on this in the paper
(or in the poster!)*

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

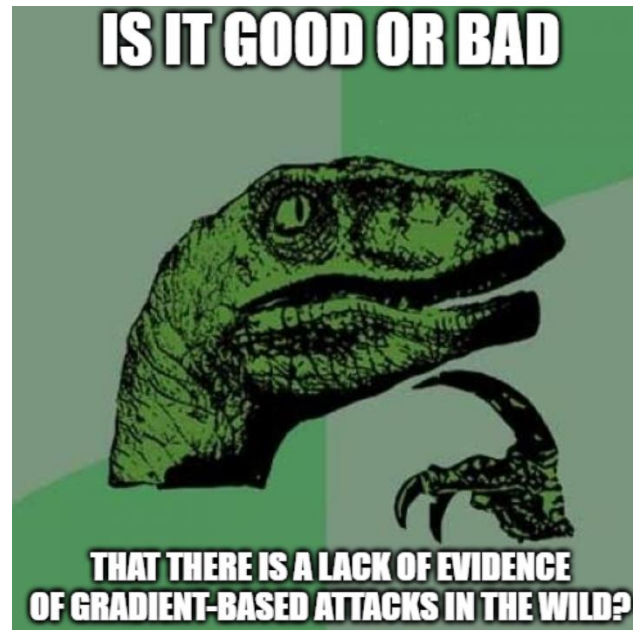
P4: *Source-code* disclosure with “just culture”

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.



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”

Meet the team

