

Bangalore, June 4<sup>th</sup> 2026

21<sup>st</sup> ACM ASIA Conference on Computer and Communications Security

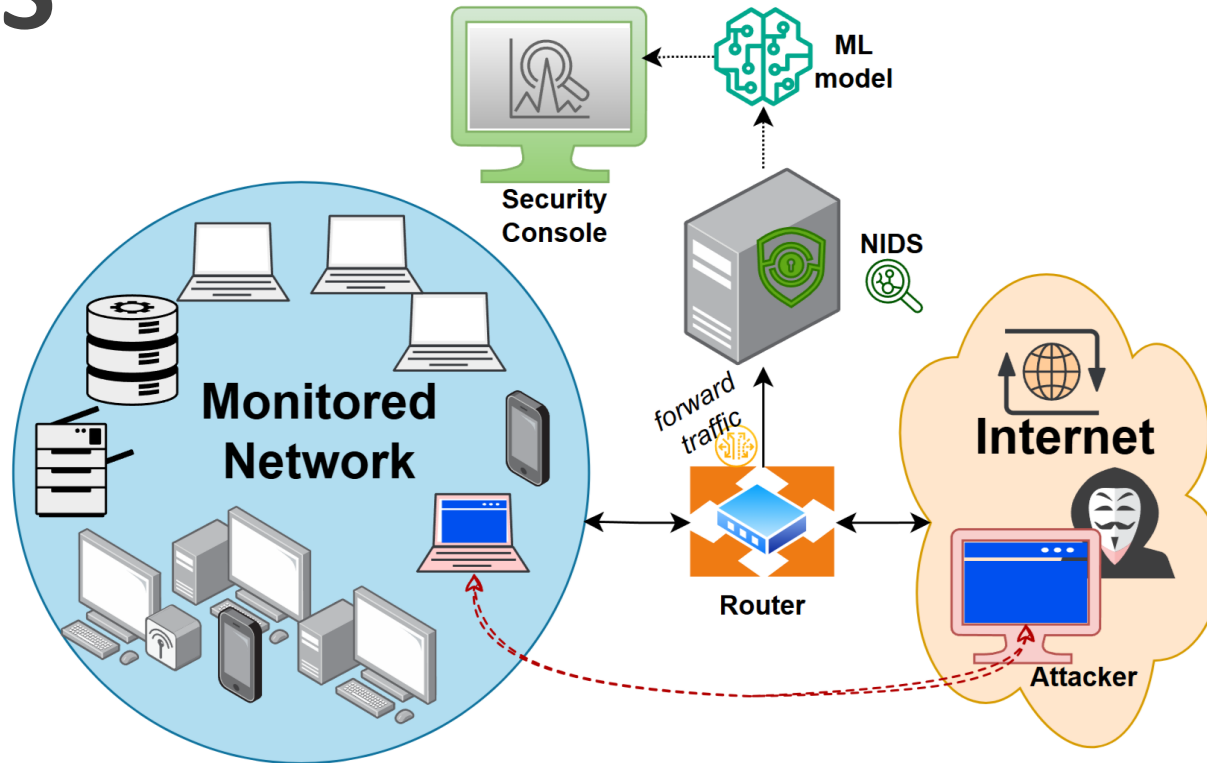
# “What is the Problem Space?”

## Defining Host-space Adversarial Perturbations against Network Intrusion Detection Systems

Miel Verkerken, Laurens d’Hooge, Bruno Volckaert, Filip de Turck, Giovanni Apruzzese

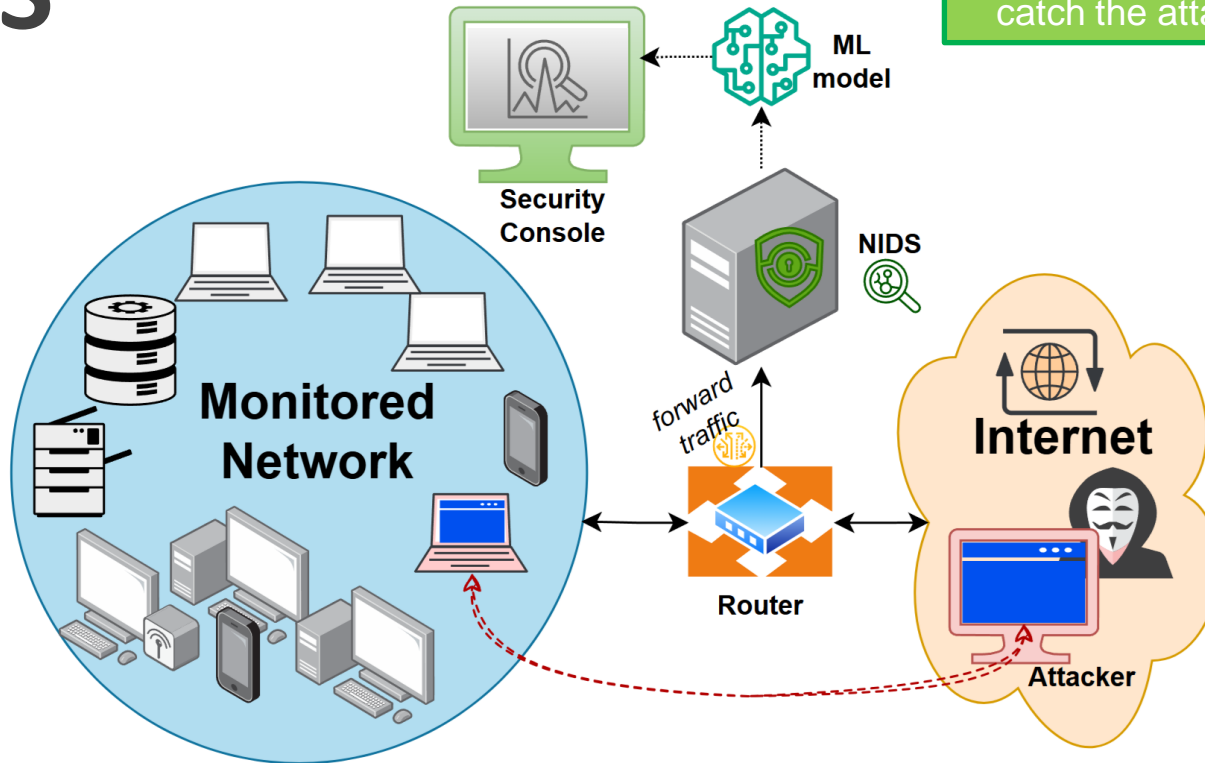


# NIDS



# NIDS

The NIDS wants to catch the attacker

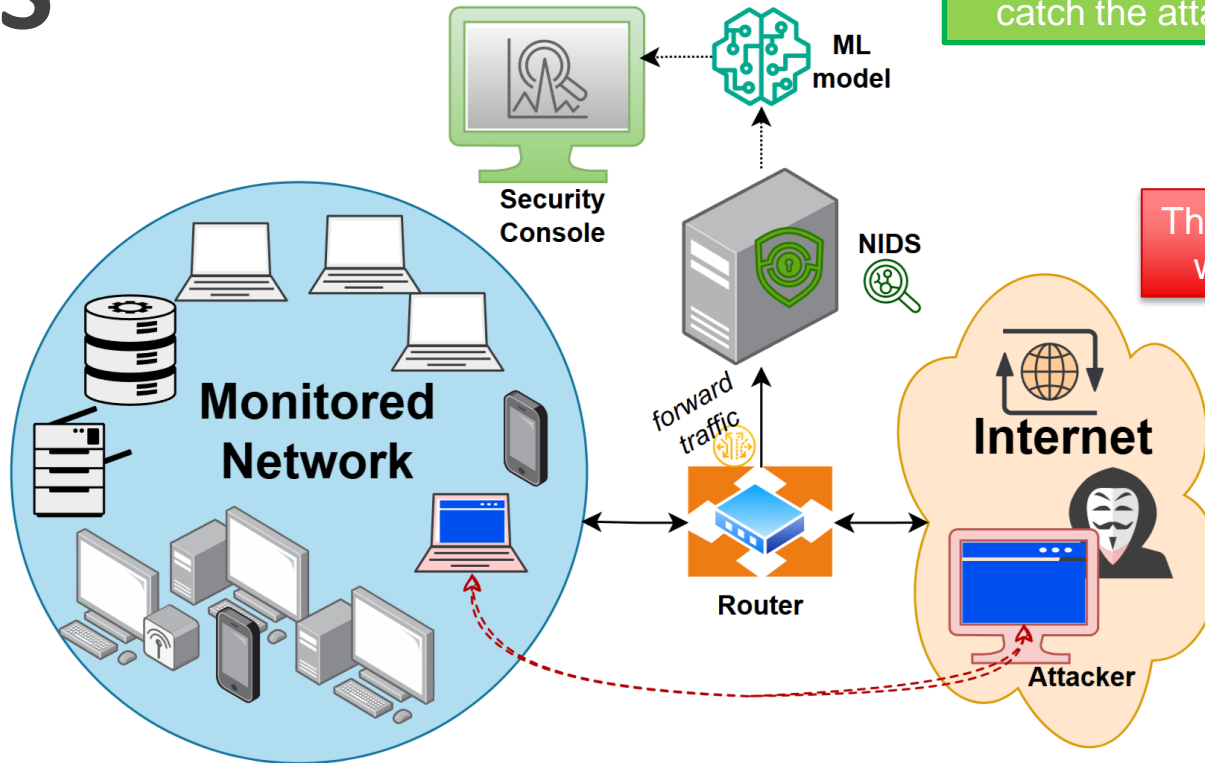


# NIDS



The NIDS wants to catch the attacker

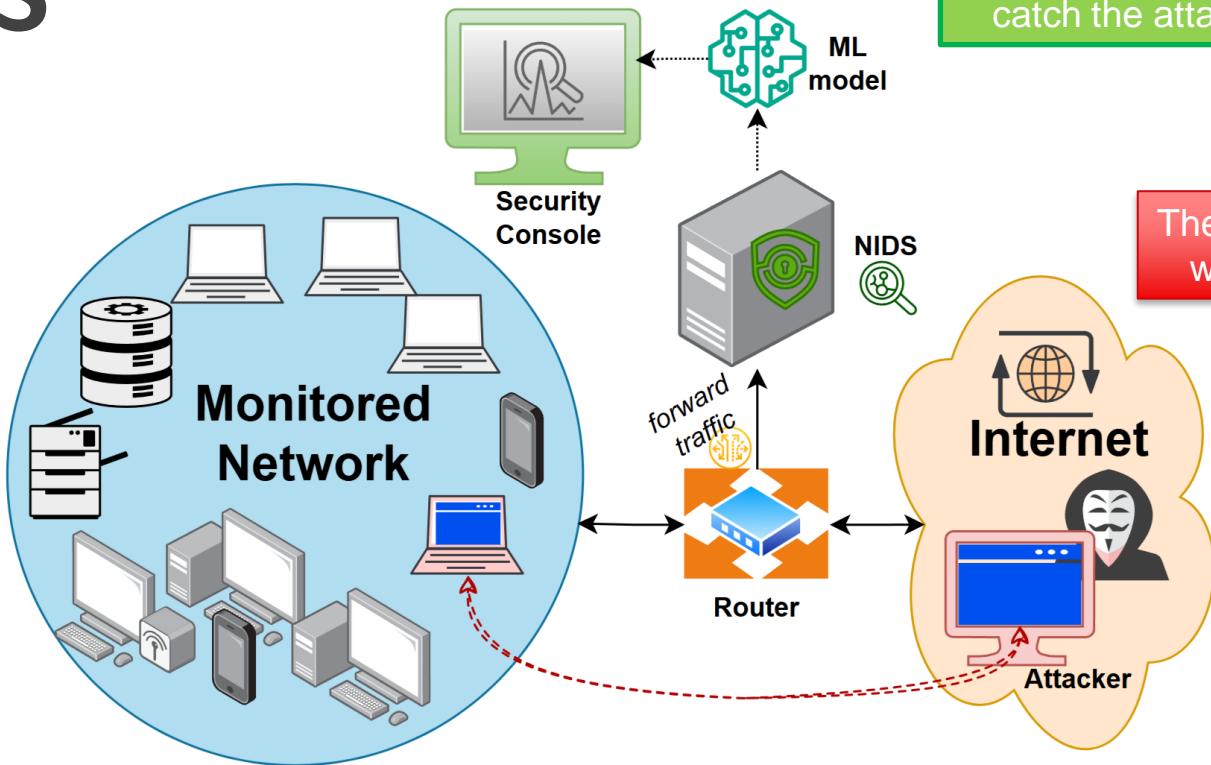
The attacker does not want to be caught



# NIDS



The NIDS wants to catch the attacker



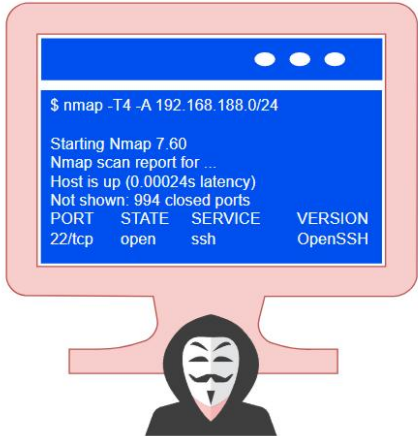
The attacker does not want to be caught

What is the attacker's problem space?

# Example: PortScan



## Host Space (Commands)





# Example: PortScan



## Host Space (Commands)

```
$ nmap -T4 -A 192.168.188.0/24
```

```
Starting Nmap 7.60
Nmap scan report for ...
Host is up (0.00024s latency)
Not shown: 994 closed ports
PORT      STATE SERVICE
22/tcp    open  ssh
VERSION  OpenSSH
```



## Traffic Space (Network Packets)



The screenshot shows a Wireshark capture of network traffic. The main pane displays a list of packets, with several TCP connections to port 9900 highlighted in red. The packet details pane shows the structure of a TCP segment, including the source and destination ports, sequence number, and window size.

## NetFlow Space (traffic metadata)



(srcIP, dstIP, srcPort, dstPort, protocol)	srcBytes	dstBytes	srcPkts	dstPkts	[...]
192.168.188.11, 192.168.188.65, 60251, 9900, TCP	650	513	8	7	[...]
192.168.188.11, 192.168.188.72, 45088, 22, TCP	321	451	6	9	[...]
[...]	[...]	[...]	[...]	[...]	[...]



# Takeaway



Real attackers perform *actions* leading to changes in the network data generated by their controlled hosts.



# Takeaway

Real attackers perform *actions* leading to changes in the network data generated by their controlled hosts.

From an adversarial ML viewpoint, this can be simulated through “host-space perturbations” (HsP), which require operating at the host level—and not on the data that, despite having been generated by such a host, was captured by an appliance outside the attacker’s control.



# Threat Model

Given an attacker having a goal  $\mathcal{G}$ , which they pursue by using their knowledge  $\mathcal{K}$  and capabilities  $\mathcal{C}$  on the targeted system to devise a certain strategy  $\mathcal{S}$ .

To realise a valid HsP, the following must be met:



# Threat Model – Goal

The overarching goal is the same as in the original threat model (i.e.,  $\mathcal{G}$ ), but it must also include evasion of an (ML-based) NIDS.

*Otherwise, the attacker would not even attempt any evasion!*

# Threat Model – Knowledge



The knowledge is identical to that of the original threat model (i.e.,  $\mathcal{K}$ ), but it must be assumed that the attacker expects that the ML-NIDS would detect the strategy  $\mathcal{S}$  the attacker would follow if they did not apply any HsP

*Otherwise, why would the attacker even attempt using an HsP?*

# Threat Model – Capabilities



The capabilities must not change w.r.t.  $\mathcal{C}$ : the attacker cannot be assumed to, e.g., control more hosts, or have more privileges on their already controlled hosts.

*Otherwise, there would be a logical fallacy in the threat model*

# Threat Model – Strategy



The strategy *must* change w.r.t.  $\mathcal{S}$ , and is still dictated by the assumed  $\mathcal{K}$  and  $\mathcal{C}$ , and must only include operations that involve the attacker's controlled hosts.

E.g., an attacker cannot "bribe" an employee to obtain a password; or physically go to a host they wish to portscan in order to check which ports are open locally.

# Research Question 1



“Has prior work on the robustness of ML-NIDS considered perturbations involving launching different commands on the attacker-controlled host (i.e., HsP)?”

# Research Question 1



“Has prior work on the robustness of ML-NIDS considered perturbations involving launching different commands on the attacker-controlled host (i.e., HsP)?”

- Method: systematic literature review (n=316)
- Findings: no paper allows us to answer “yes” to RQ1
- Notable mentions: [32] and [18]

[32]: Catillo, Marta, et al. "Towards realistic problem-space adversarial attacks against machine learning in network intrusion detection." Proceedings of the 19th International Conference on Availability, Reliability and Security. 2024.

[18]: Apruzzese, Giovanni, Aurore Fass, and Fabio Pierazzi. "When adversarial perturbations meet concept drift: an exploratory analysis on ml-nids." Proceedings of the 2024 Workshop on Artificial Intelligence and Security. 2024.

# Research Question 1



“Has prior work on the robustness of ML-NIDS considered perturbations involving launching different commands on the attacker-controlled host (i.e., HsP)?”

- Method: systematic literature review (n=316)
- Findings: no paper allows us to answer “yes” to RQ1
- Notable mentions: [32] and [18]

*Studying HsP requires performing actions on physical hosts*

[32]: Catillo, Marta, et al. "Towards realistic problem-space adversarial attacks against machine learning in network intrusion detection." Proceedings of the 19th International Conference on Availability, Reliability and Security. 2024.

[18]: Apruzzese, Giovanni, Aurore Fass, and Fabio Pierazzi. "When adversarial perturbations meet concept drift: an exploratory analysis on ml-nids." Proceedings of the 2024 Workshop on Artificial Intelligence and Security. 2024.

# Research Question 2



“What are some effects that HsP can have on existing ML-NIDS?”

# Research Question 2



“What are some effects that HsP can have on existing ML-NIDS?”

- Method: experiments on benchmark datasets, and on a real-world network
- Attacker’s Goal: SSH-bruteforcing, particularly by using *ssh-patator* [5]
- ML-NIDS: trained on *ssh-patator --persistent=1* (or *--persistent=0*)
- Attacker’s Strategy (HsP): *ssh-patator --persistent=0* (or *--persistent=1*)

# Research Question 2



“What are some effects that HsP can have on existing ML-NIDS?”

- Method: experiments on benchmark datasets, and on a real-world network
- Attacker’s Goal: SSH-bruteforcing, particularly by using *ssh-patator* [5]
- ML-NIDS: trained on *ssh-patator --persistent=1* (or *--persistent=0*)
- Attacker’s Strategy (HsP): *ssh-patator --persistent=0* (or *--persistent=1*)

Dataset	CICIDS17		CICIDS18
Train Set	Benign + P=1	Benign + all malicious	Benign + P=1
Test Set			
DT			
RF			
XGB			
SVM			
DNN			

Benchmarks

Train Set	Benign + P=1	Benign + P=0
Test Set		
DT		
RF		
XGB		
SVM		
DNN		

Real World

# Research Question 2



“What are some effects that HsP can have on existing ML-NIDS?”

- Method: experiments on benchmark datasets, and on a real-world network
- Attacker’s Goal: SSH-bruteforcing, particularly by using *ssh-patator* [5]
- ML-NIDS: trained on *ssh-patator --persistent=1* (or *--persistent=0*)
- Attacker’s Strategy (HsP): *ssh-patator --persistent=0* (or *--persistent=1*)

Dataset	CICIDS17						CICIDS18		
	Benign + P=1		Benign + all malicious				Benign + P=1		
Train Set	Benign	P=1		Benign	(all)		Benign	P=1	
DT	<0.001	0.997		<0.001	>0.999		0.000	1.000	
RF	0.000	0.998		<0.001	>0.999		0.000	1.000	
XGB	<0.001	0.998		<0.001	>0.999		0.000	1.000	
SVM	<0.001	0.993		<0.001	0.997		0.000	1.000	
DNN	0.000	0.998		<0.001	>0.999		0.000	1.000	

Benchmarks

Train Set	Benign + P=1			Benign + P=0		
	Benign	P=1		Benign	P=0	
DT	<0.001	1.000		<0.001	1.000	
RF	0.000	1.000		0.000	1.000	
XGB	0.000	>0.999		<0.001	1.000	
SVM	0.000	1.000		0.000	1.000	
DNN	<0.001	1.000		<0.001	1.000	

Real World

# Research Question 2



“What are some effects that HsP can have on existing ML-NIDS?”

- Method: experiments on benchmark datasets, and on a real-world network
- Attacker’s Goal: SSH-bruteforcing, particularly by using *ssh-patator* [5]
- ML-NIDS: trained on *ssh-patator --persistent=1* (or *--persistent=0*)
- Attacker’s Strategy (HsP): *ssh-patator --persistent=0* (or *--persistent=1*)

Dataset	CICIDS17						CICIDS18		
	Benign + P=1			Benign + all malicious			Benign + P=1		
Train Set	Benign	P=1	P=0	Benign	(all)	P=0	Benign	P=1	P=0
DT	<0.001	0.997	0.224	<0.001	>0.999	0.214	0.000	1.000	0.000
RF	0.000	0.998	0.490	<0.001	>0.999	0.073	0.000	1.000	0.000
XGB	<0.001	0.998	0.986	<0.001	>0.999	<0.001	0.000	1.000	0.000
SVM	<0.001	0.993	0.000	<0.001	0.997	<0.001	0.000	1.000	0.000
DNN	0.000	0.998	0.000	<0.001	>0.999	<0.001	0.000	1.000	0.000

Benchmarks

Train Set	Benign + P=1		Benign + P=0	
	Benign	P=1	Benign	P=0
DT	<0.001	1.000	<0.001	1.000
RF	0.000	1.000	0.000	1.000
XGB	0.000	>0.999	<0.001	1.000
SVM	0.000	1.000	0.000	1.000
DNN	<0.001	1.000	<0.001	1.000

Real World

# Research Question 2



“What are some effects that HsP can have on existing ML-NIDS?”

- Method: experiments on benchmark datasets, and on a real-world network
- Attacker’s Goal: SSH-bruteforcing, particularly by using *ssh-patator* [5]
- ML-NIDS: trained on *ssh-patator --persistent=1* (or *--persistent=0*)
- Attacker’s Strategy (HsP): *ssh-patator --persistent=0* (or *--persistent=1*)

Dataset	CICIDS17						CICIDS18		
	Benign + P=1			Benign + all malicious			Benign + P=1		
Train Set	Benign	P=1	P=0	Benign	(all)	P=0	Benign	P=1	P=0
DT	<0.001	0.997	0.224	<0.001	>0.999	0.214	0.000	1.000	0.000
RF	0.000	0.998	0.490	<0.001	>0.999	0.073	0.000	1.000	0.000
XGB	<0.001	0.998	0.986	<0.001	>0.999	<0.001	0.000	1.000	0.000
SVM	<0.001	0.993	0.000	<0.001	0.997	<0.001	0.000	1.000	0.000
DNN	0.000	0.998	0.000	<0.001	>0.999	<0.001	0.000	1.000	0.000

Benchmarks

Train Set	Benign + P=1			Benign + P=0		
	Benign	P=1	P=0	Benign	P=0	P=1
DT	<0.001	1.000	0.000	<0.001	1.000	0.000
RF	0.000	1.000	0.000	0.000	1.000	0.000
XGB	0.000	>0.999	0.000	<0.001	1.000	0.000
SVM	0.000	1.000	0.490	0.000	1.000	1.000
DNN	<0.001	1.000	0.000	<0.001	1.000	1.000

Real World

# Research Question 2



“What are some effects that HsP can have on existing ML-NIDS?”

**ANSWER TO RQ2:** Some HsP can drop the *tpr* of an ML-NIDS from 1 to 0 by changing one character in the command launched by the attacker’s controlled host. Even HsP leveraging a different tool but for the same malicious goal can lead to complete bypass. HsP requiring higher privileges are not necessarily more evasive. The effectiveness of HsP against specific ML models varies.

# Explanation

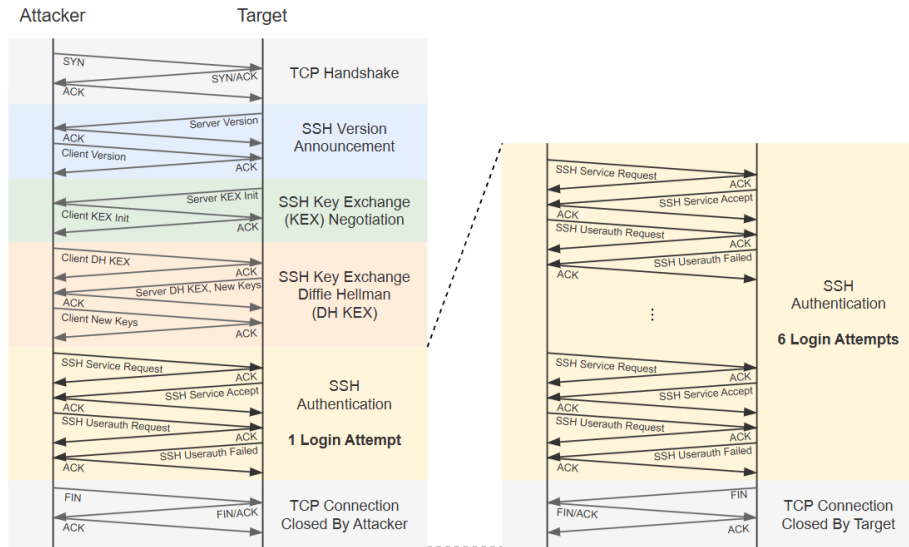


“Why is it that, for *ssh-patator*, changing a single character can lead to a complete bypass of the ML-NIDS?”

# Explanation



“Why is it that, for *ssh-patator*, changing a single character can lead to a complete bypass of the ML-NIDS?”



Non-Persistent Brute Force Attack

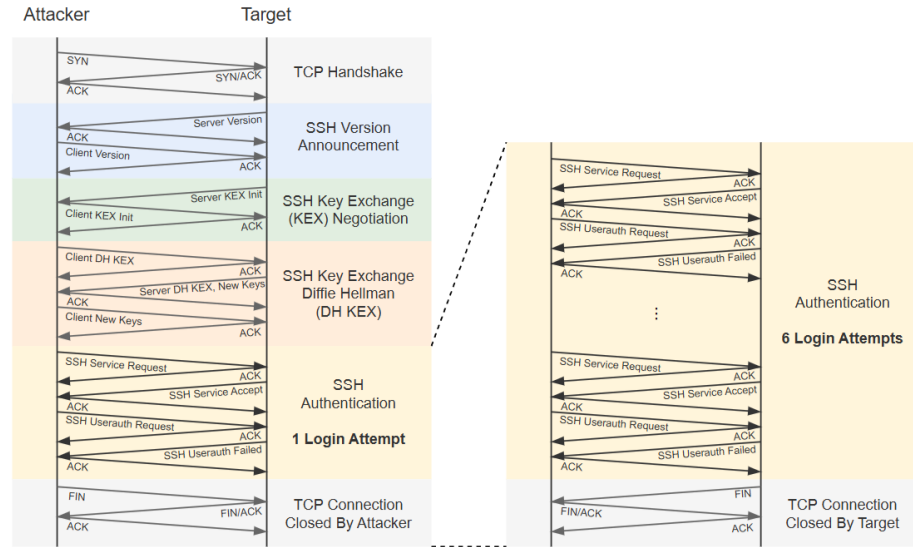
Persistent Brute Force Attack

(a) **Packet-level.** By launching *patator* -P=0 (left), only one login attempt is made. However, launching *patator* -P=1 (right) leads to 6x more attempts. This difference is the “packet-level perturbation” induced by an HsP that changes -P=0 to *patator* -P=1 (i.e., a one-character change in the command launched by the attacker).



# Explanation

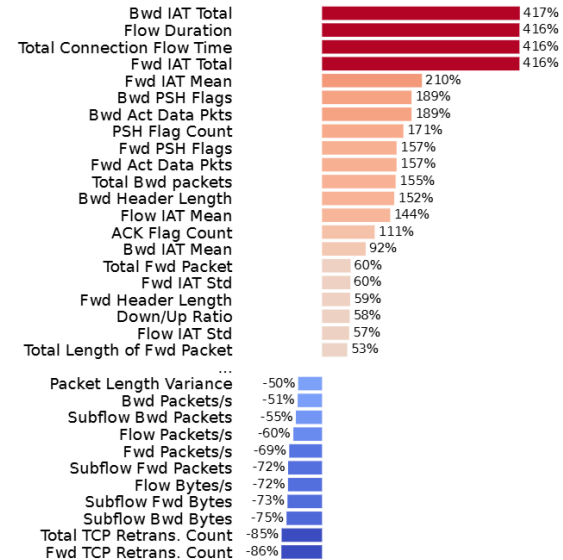
“Why is it that, for *ssh-patator*, changing a single character can lead to a complete bypass of the ML-NIDS?”



Non-Persistent Brute Force Attack

Persistent Brute Force Attack

(a) **Packet-level.** By launching patator -P=0 (left), only one login attempt is made. However, launching patator -P=1 (right) leads to 6x more attempts. This difference is the “packet-level perturbation” induced by an HsP that changes -P=0 to patator -P=1 (i.e., a one-character change in the command launched by the attacker).



(b) **NetFlow-level.** In a the NetFlow-based feature space, an HsP switching -P=0 to -P=1 leads to substantial differences (numbers in the plot are obtained by computing  $-P=0 / -P=1$ )

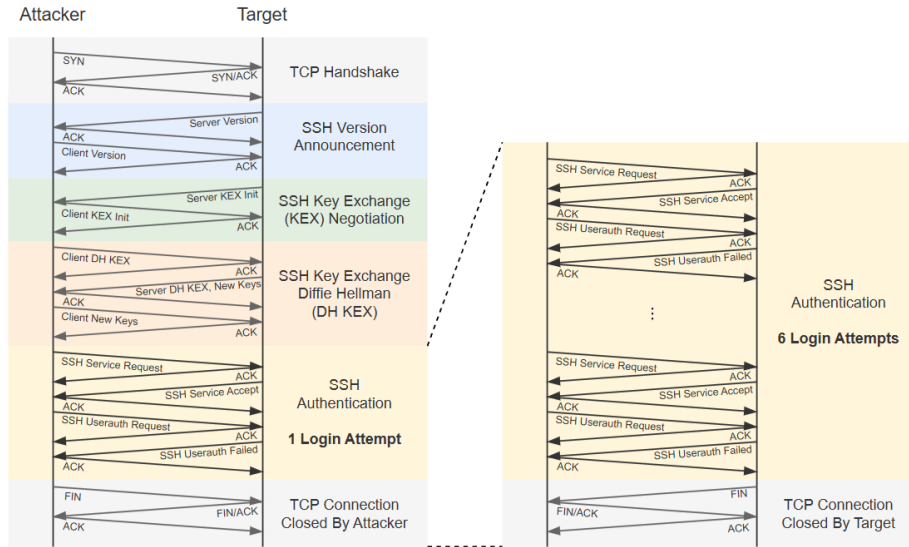


# Exp

## “Why C

**TAKEAWAY:** Changing a single option (or parameter value) of a malicious command can lead to substantially different network activities—which achieve the same goal, but can cause misclassifications due to OOD (exploitable via HsP). Attempts to reproduce, or approximate, such changes by manipulating pre-collected datapoints is challenging and may result in misleading evaluations.

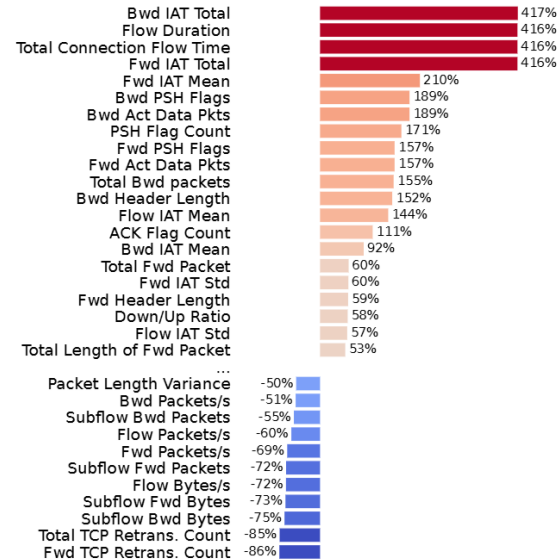
## Character ?”



Non-Persistent Brute Force Attack

Persistent Brute Force Attack

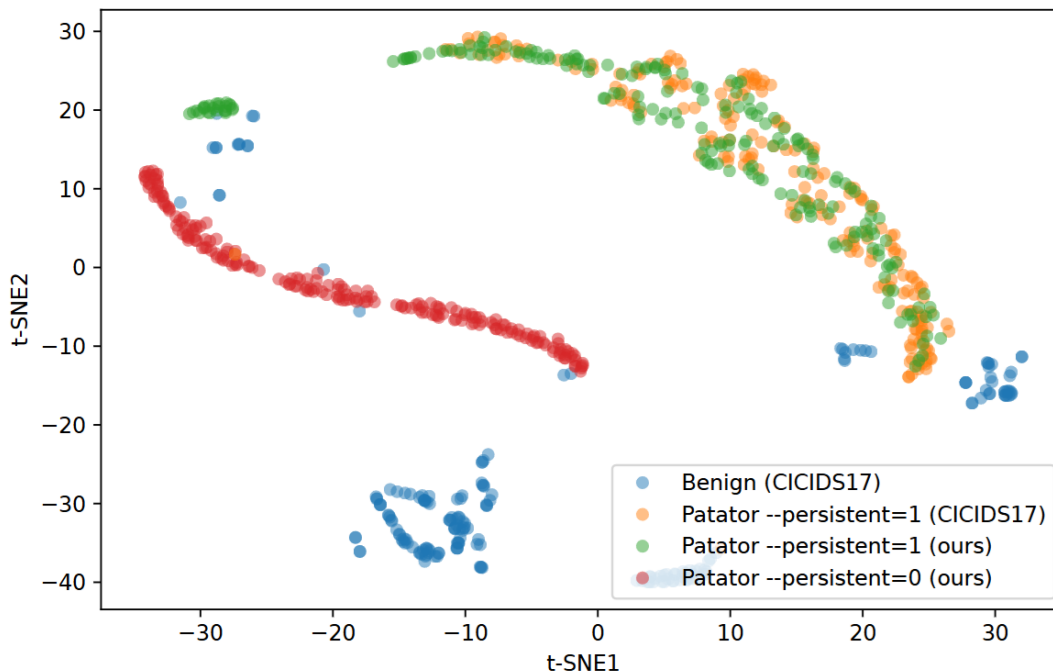
(a) **Packet-level.** By launching patator -P=0 (left), only one login attempt is made. However, launching patator -P=1 (right) leads to 6x more attempts. This difference is the “packet-level perturbation” induced by an HsP that changes -P=0 to patator -P=1 (i.e., a one-character change in the command launched by the attacker).



(b) **NetFlow-level.** In the NetFlow-based feature space, an HsP switching -P=0 to -P=1 leads to substantial differences (numbers in the plot are obtained by computing  $-P=0 / -P=1$ )

Giovanni Apruzzese  
giovannia@ru.is

# Out of Distribution?



**Fig. 4: OOD verification.** We use a tSNE plot to visualize the distribution of NetFlows generated by patator (as well as benign ones).

# But!



- Only one attack! → We show more attacks in the paper

# But!



- Only one attack! → We show more attacks in the paper
- Only one model type! → We evaluate other model types in the paper!

# But!



- Only one attack! → We show more attacks in the paper
- Only one model type! → We evaluate other model types in the paper!
- Only one benchmark dataset! → We evaluate also in other datasets

# But!



- Only one attack! → We show more attacks in the paper
- Only one model type! → We evaluate other model types in the paper!
- Only one benchmark dataset! → We evaluate also in other datasets
- HsP are not adversarial perturbations! → Let's discuss 😊

# Conclusions



- HsP describe “alternative strategies” to bypass (ML-)NIDS
- Prior work on robustness of ML-NIDS studied “hard-to-realise” adversarial perturbations, but we showed that evading ML-NIDS via HsP is simple
- Training on a single “form” of a given attack is risky

# Conclusions – so what now?



- HsP describe “alternative strategies” to bypass (ML-)NIDS
- Prior work on robustness of ML-NIDS studied “hard-to-realise” adversarial perturbations, but we showed that evading ML-NIDS via HsP is simple
- Training on a single “form” of a given attack is risky
- Future work should investigate “easy-to-apply” perturbations implemented at the host level (i.e., HsP)
- (hard problem) mapping fine-grained feature-space perturbations to HsP

Bangalore, June 4<sup>th</sup> 2026

21<sup>st</sup> ACM ASIA Conference on Computer and Communications Security

# “What is the Problem Space?” Defining Host-space Adversarial Perturbations against Network Intrusion Detection Systems

Miel Verkerken, Laurens d’Hooge, Bruno Volckaert, Filip de Turck, Giovanni Apruzzese





IEEE SaTML'26

# ConCap: Practical Network Traffic Generation for (ML- and) Flow-based Intrusion Detection Systems

Miel Verkerken<sup>†</sup>, Laurens D'hooge<sup>†</sup>, Bruno Volckaert<sup>†</sup>, Filip De Turck<sup>†</sup>, Giovanni Apruzzese<sup>¶</sup>  
<sup>†</sup>*Ghent University – imec*, <sup>¶</sup>*University of Liechtenstein*, <sup>¶</sup>*Reykjavik University*,  
{name.surname}@{<sup>†</sup>ugent.be, <sup>¶</sup>uni.li}