

Superviz Plenary Meeting

WP4-PO4.1

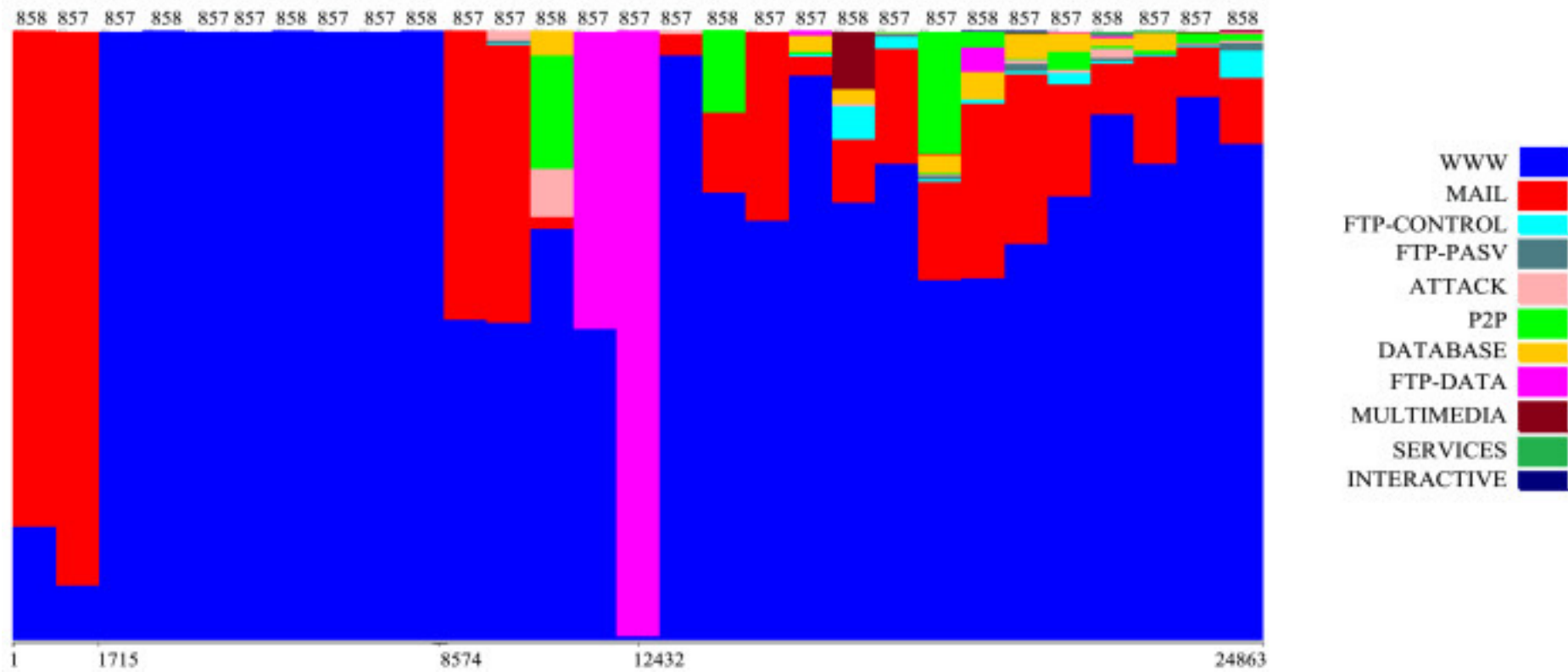
Yufei Han @ INRIA PIRAT

March 11, 2025

@Campus Cyber

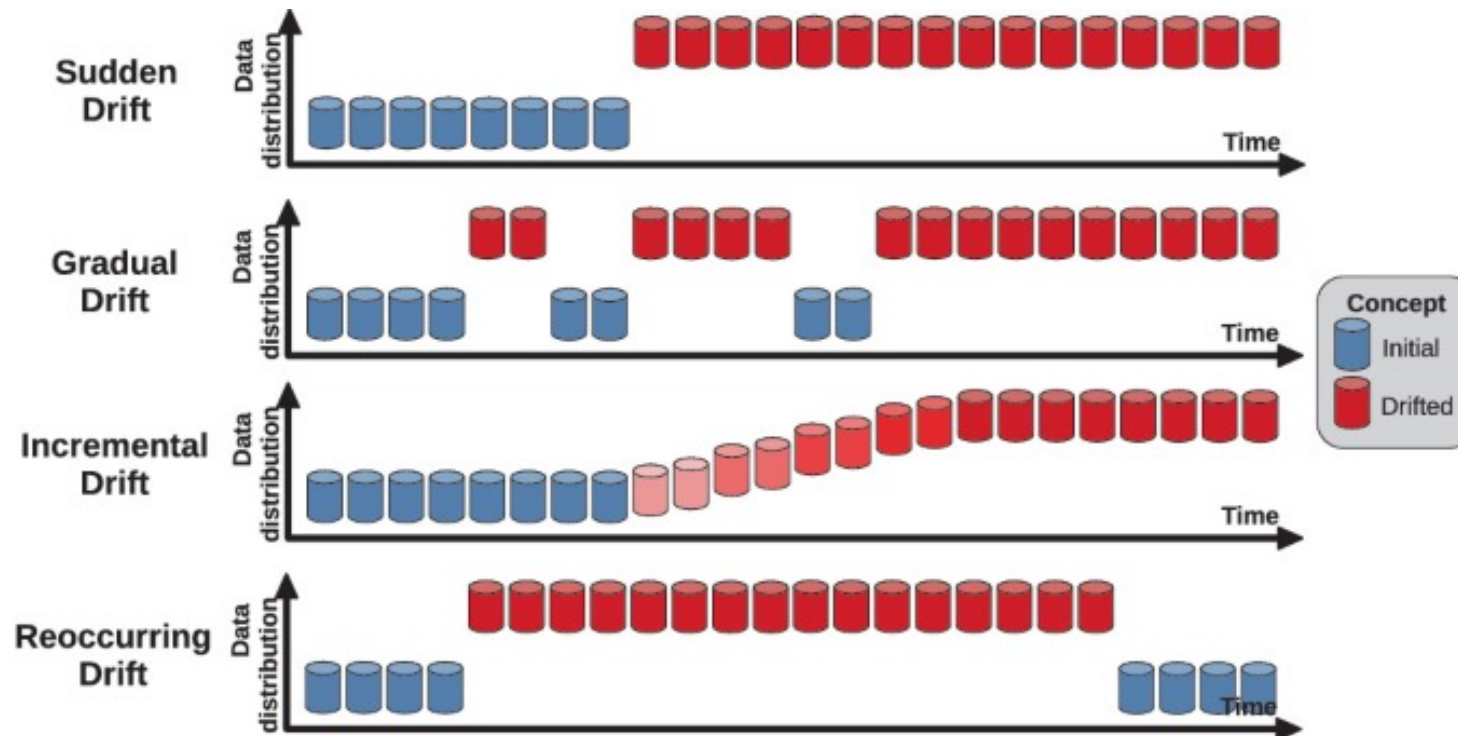
Robust and Transferable Learning for IDS

- What is the bottleneck for ML-based IDS
 - Concept drift of attack behaviours: attackers may change attack techniques to evade detection or exploit new attack surfaces



Robust and Transferable Learning for IDS

- What is the bottleneck for ML-based IDS
 - Concept drift of attack behaviours: attackers may change attack techniques to evade detection or exploit new attack surfaces




Robust and Transferable Learning for IDS

- Theoretically, every ML model has a significantly deteriorated prediction accuracy over concept drifted inputs

Theorem 1 *Let $\theta \in \mathcal{H}$ be a hypothesis, $\epsilon_s(\theta)$ and $\epsilon_t(\theta)$ be the expected risks of source and target respectively, then*

$$\epsilon_t(\theta) \leq \epsilon_s(\theta) + \boxed{2d_k(p, q)} + C, \quad (9)$$

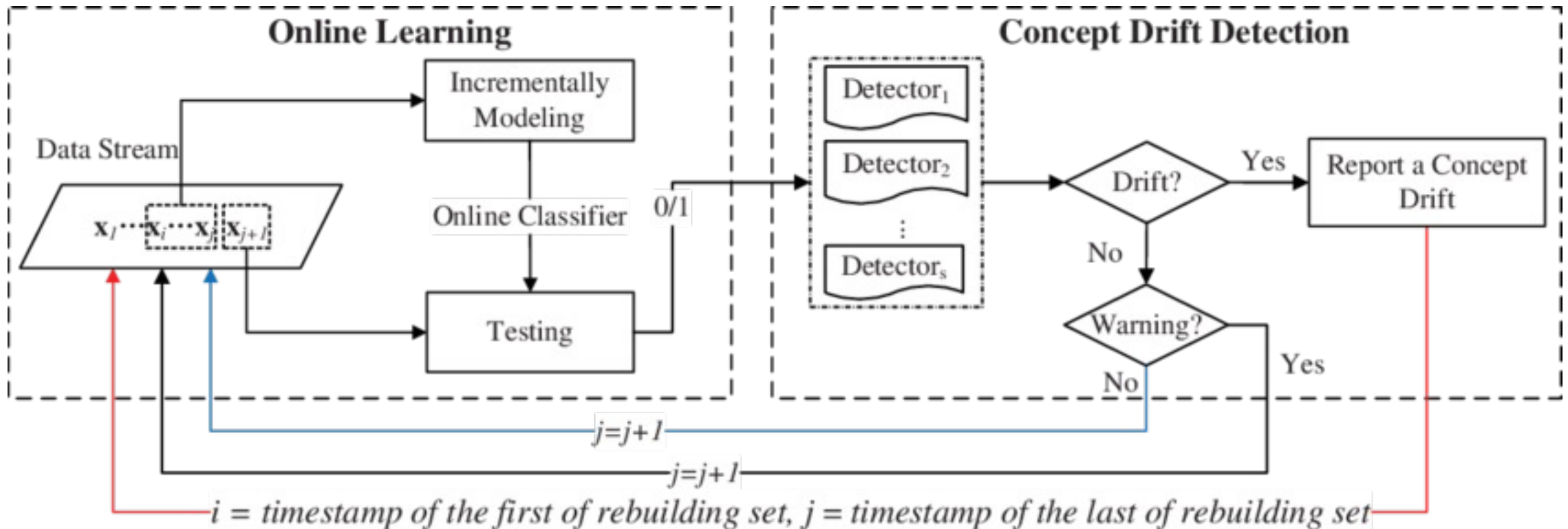
where C is a constant for the complexity of hypothesis space and the risk of an ideal hypothesis for both domains.



Distribution gap between training and testing data

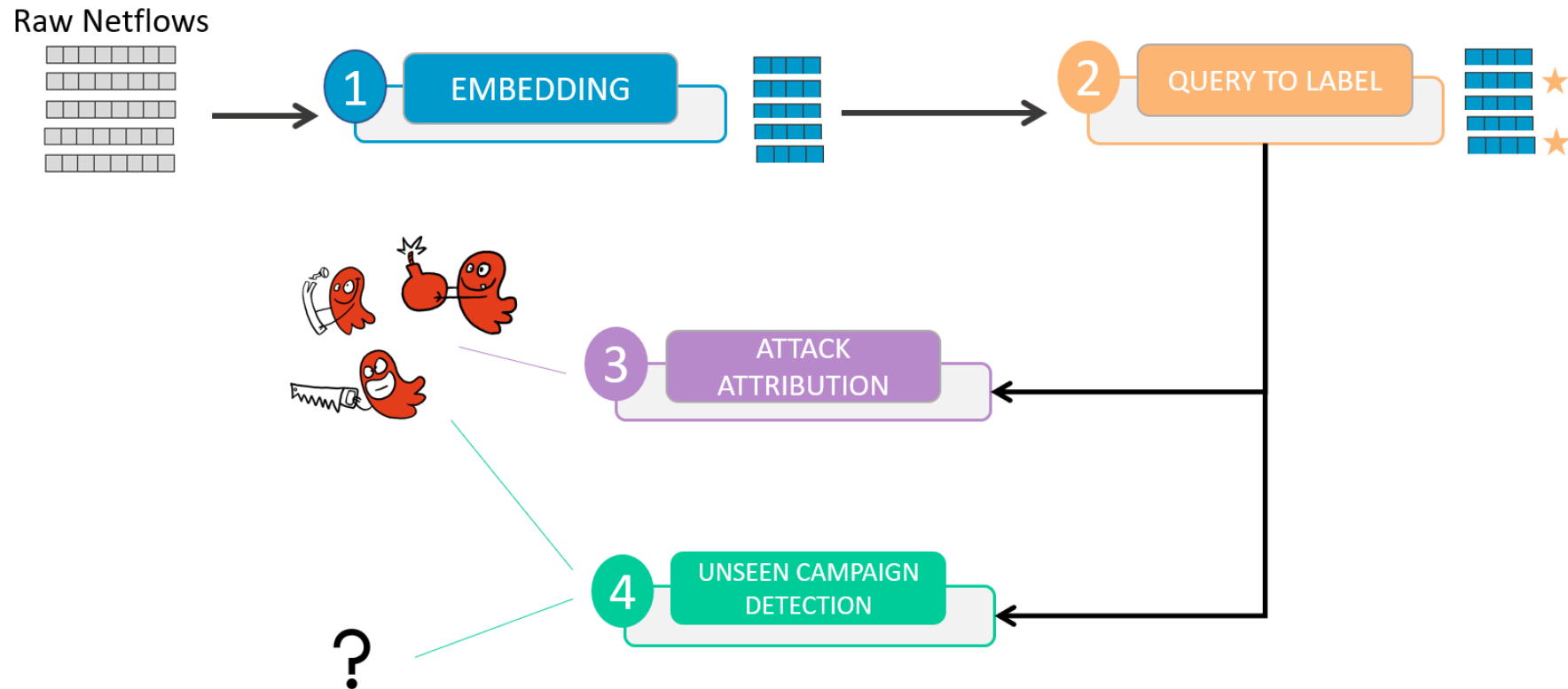
Robust and Transferable Learning for IDS

- How to reach this goal ?



Robust and Transferable Learning for IDS

- Active Learning as a protocol for incrementally update the knowledge for network attack classification



Helene Orsini and Yufei Han, DYNAMO: Towards Network Attack Campaign Attribution via Density-Aware Active Learning, <https://hal-emse.ccsd.cnrs.fr/UNIV-UBS/hal-04877620v1>

Robust and Transferable Learning for IDS

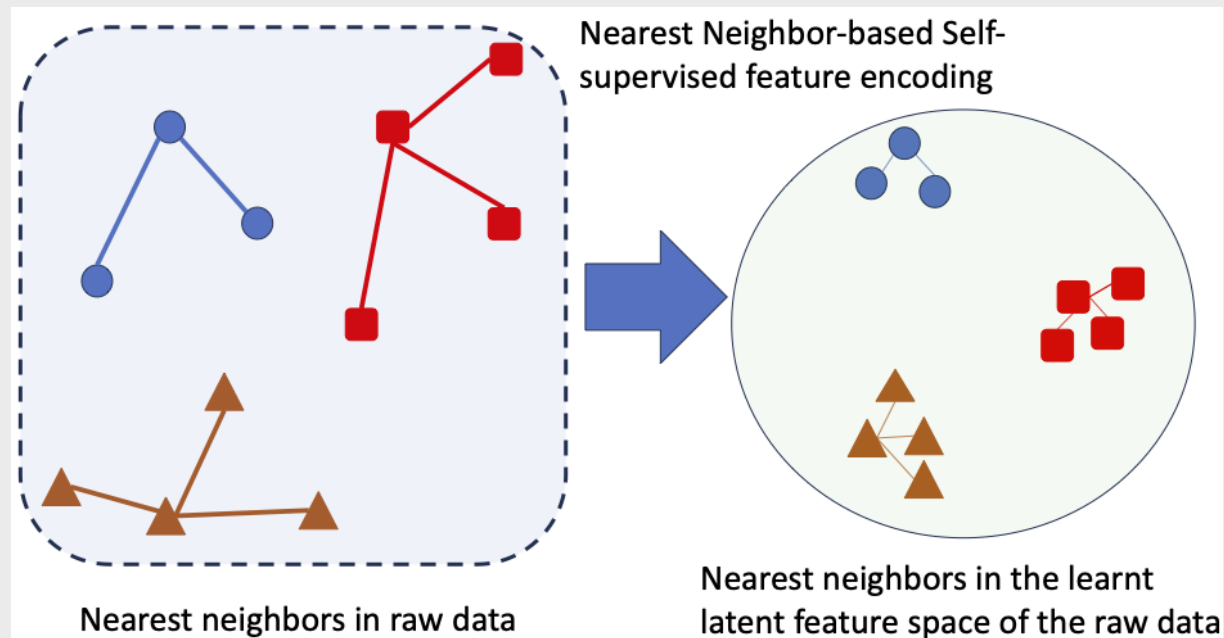
1

EMBEDDING

Raw feature vectors from Netflow : GraphSage method

$\theta^* =$

$$\arg \min_{\theta} - \frac{1}{nK} \sum_{i=1}^n [\sum_{k=1}^K \log(\sigma(h_{\theta}^T(x_i)h_{\theta}(x_i^{NN,k}))) - \lambda \sum_{j, x_v \notin KNN(x_i)} \log(\sigma(-h_{\theta}^T(x_i)h_{\theta}(x_v)))]$$



Robust and Transferable Learning for IDS

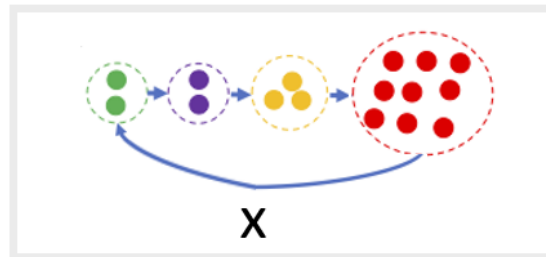
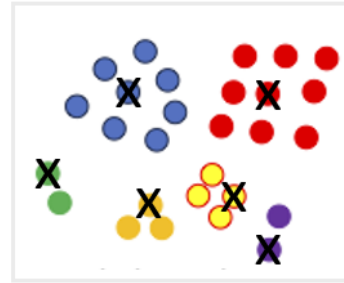
2

QUERY TO LABEL

3

ATTACK
ATTRIBUTION

Hierachical
clustering



Classifier

3

2



Robust and Transferable Learning for IDS

4

UNSEEN CAMPAIGN
DETECTION

Pu learning

Train a classifier to distinguish between positive and negative.

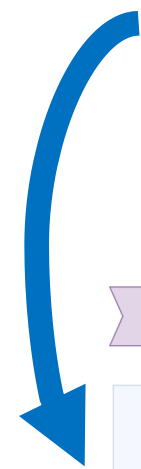
Learning phase: **Positive and Unlabelled** (*only some of the positive examples in the training data are labeled and none of the negative examples are*)

$$g_{\phi}^{pu} = \arg \min_{\phi} \frac{\pi}{n_p} \sum_{x_i \in S} [\ell(g_{\phi}^{pu}(h_{\theta}(x_i)), y_i = +1) - \ell(g_{\phi}^{pu}(h_{\theta}(x_i)), y_i = -1)] + \frac{1}{n_u} \sum_{x_i \in X_{\text{unlabeled}}} \ell(g_{\phi}^{pu}(h_{\theta}(x_i)), y_i = -1)$$

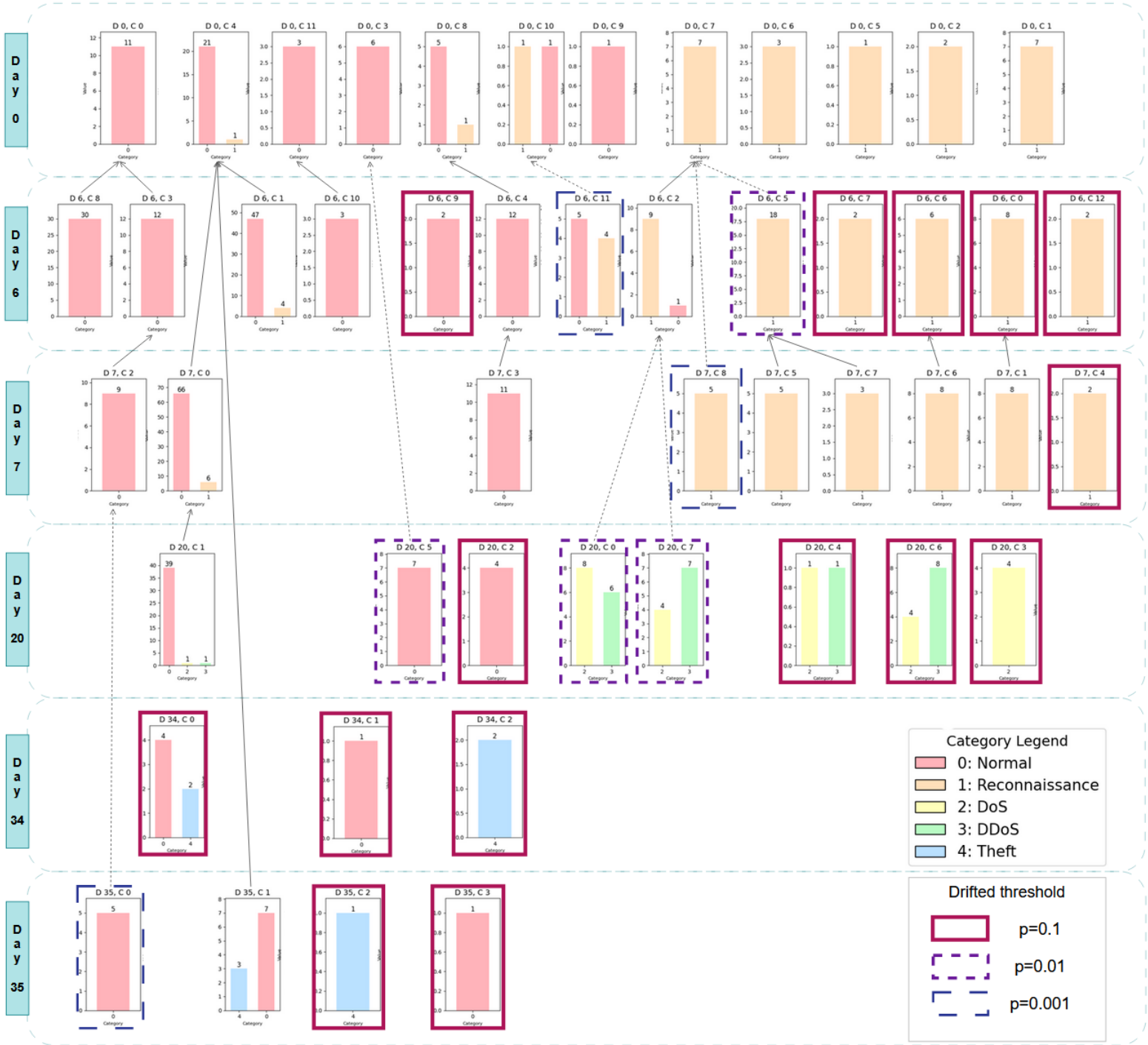
Robu

• Extenc

DS



packi
packi
packi
packi
packi
packi
packi

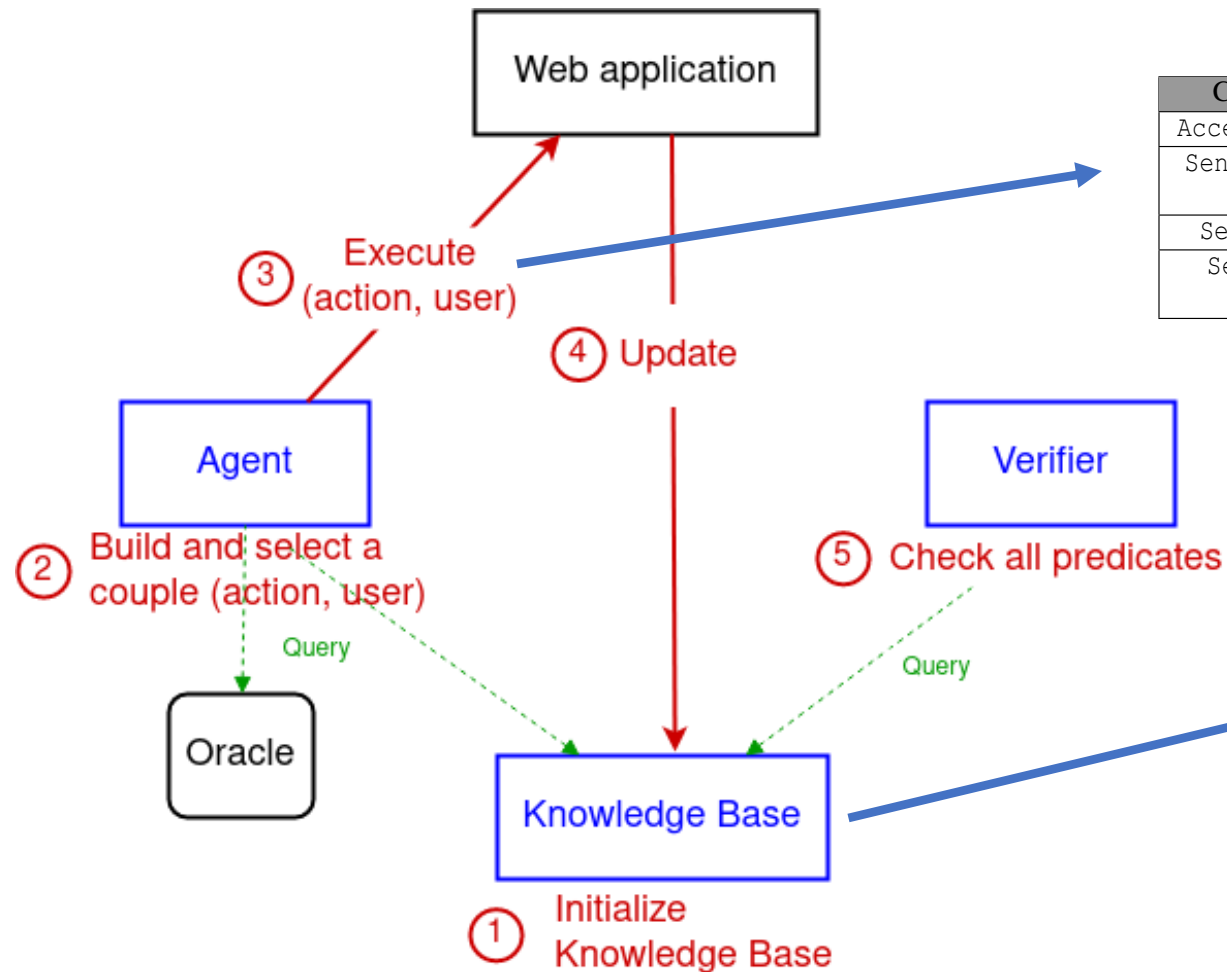


Robust and Transferable Learning for IDS

- Back to the question “How to reach a robust IDS facing concept drift ?”
- Can we predict the attack behaviour ?
- Natan Talon et al, SCWAD: Automated Pentesting of Web Applications, <https://inria.hal.science/hal-04874868v1/document>

Robust and Transferable Learning for IDS

Vulnerabilities, e.g. Broken Access Control and Reflected Cross-Site Scripting.



Action set

Command	Parameters	Expected behavior on the Web application
AccessWebPage	u: <i>url</i>	change current page for the current user
SendHttpForm	x: <i>xpath</i> , input: <i>dict</i> { <i>key</i> : <i>value</i> }	send the form identified by x and filled with input to the web application
SearchData	data: <i>string</i>	return <code>true</code> if data has been found in the current page
SetCookie	cookie: <i>dict</i> { <i>key</i> : <i>value</i> }	add cookie to user's cookies if this cookie doesn't exists and replace it otherwise

User login: user name

User credentials: The credentials for the application

User cookies: The set of cookies, active or not

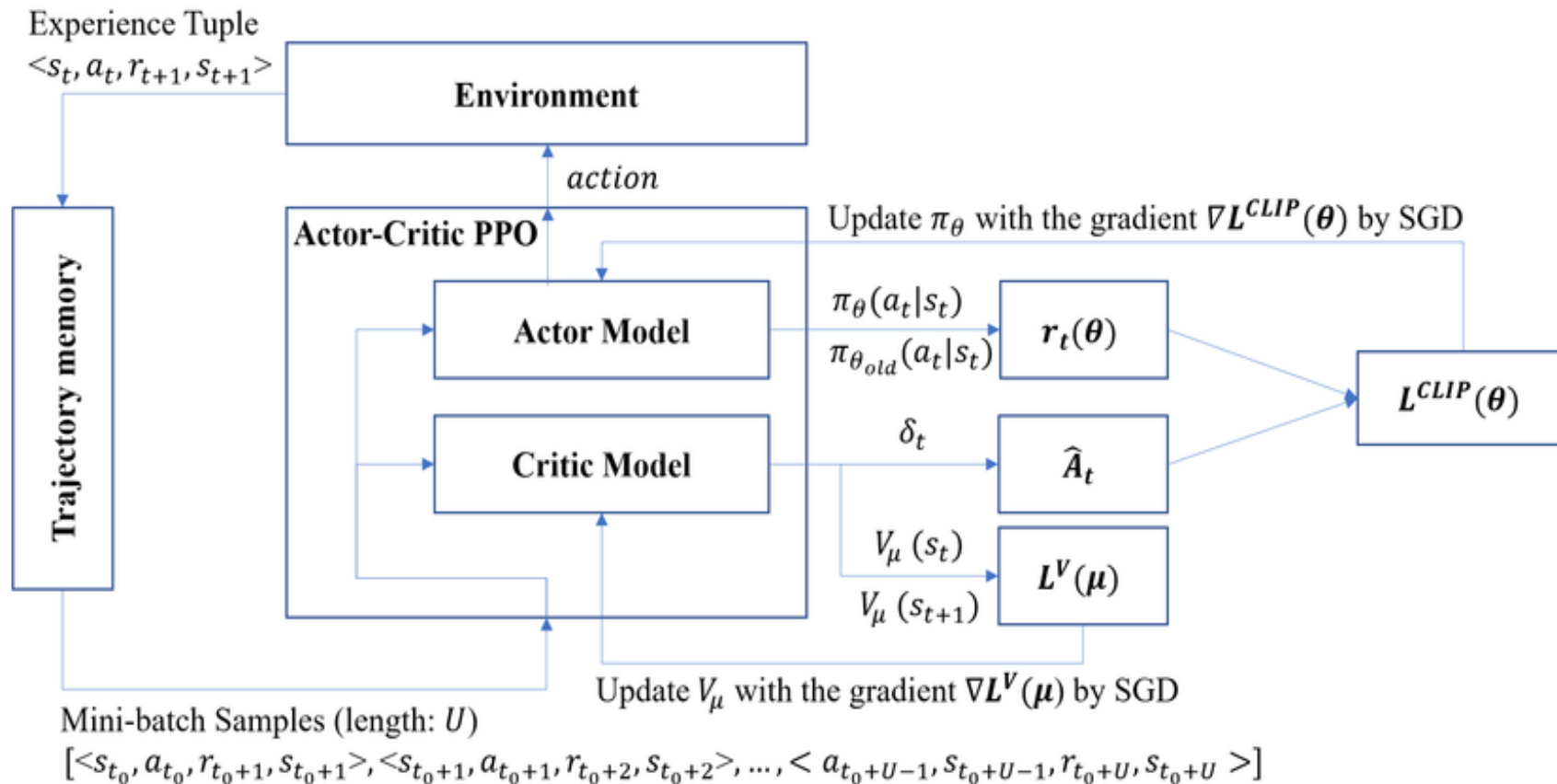
User pages: The pages visited by the user and/or those to which the user has a link

User current page: The current page visited by the user

User allowed paths: All the links to the pages reachable by the user

Robust and Transferable Learning for IDS

- Learn an attack policy with reinforcement learning



Key question to answer:

Given the current state of the target web application and the vulnerability to explore, what shall be the most likely action that an attack can take to compromise the application ?

Thanks !