

Towards Context-aware Intrusion Detection in Individual-oriented Information Systems: **An Empirical Study on Android Malware**

Paper at the 21st International Conference on Network and Service Management (CNSM2025), Bologna, Italy

PhD Student:

Tien NGUYEN

Supervisors:

Guillaume DOYEN, IRISA / IMT Atlantique

Daniela DRAGOMIRESCU, LAAS-CNRS / INSA Toulouse

Renzo E. NAVAS, IRISA / IMT Atlantique

Eric ALATA, LAAS-CNRS / INSA Toulouse

ICO Annual Scientific Day, 05/12/2025

Outline

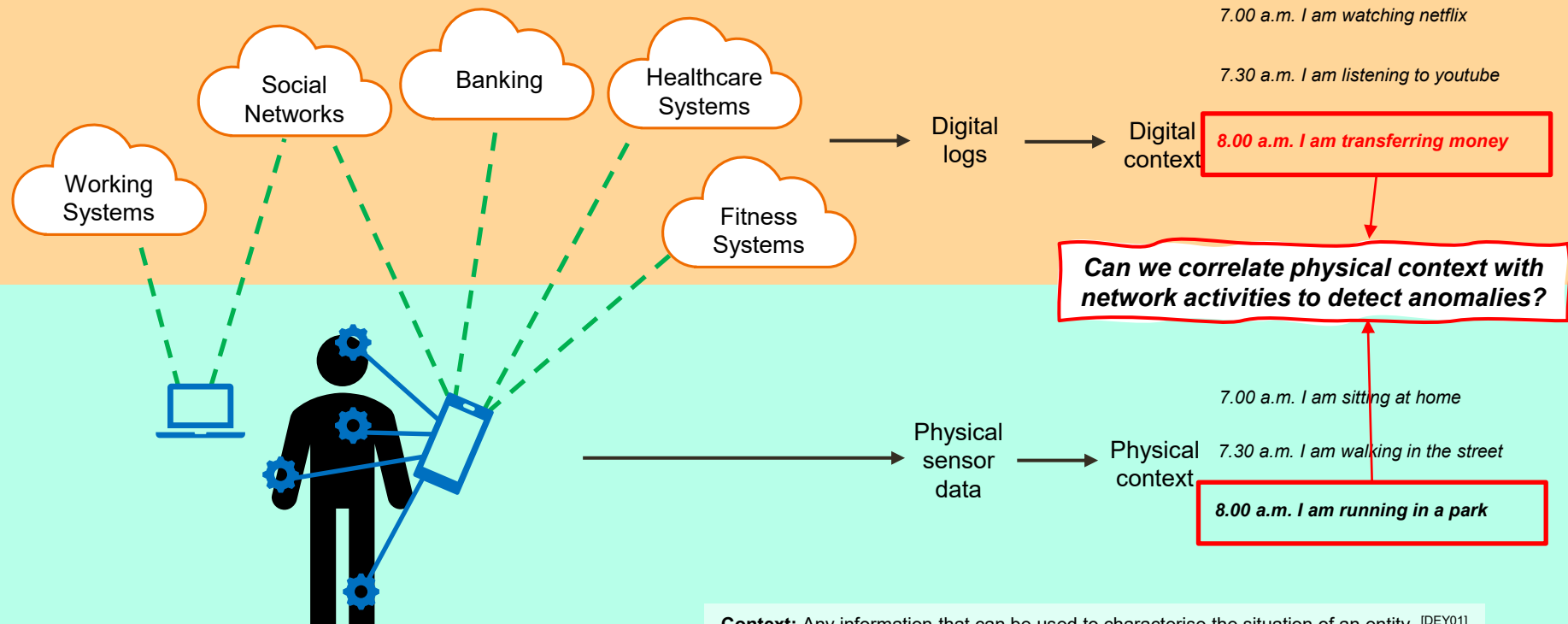
- I. Background
- II. State of The Art
- III. Proposed framework
- IV. Evaluation
- V. Conclusion

Outline

- I. Background
- II. State of The Art
- III. Proposed framework
- IV. Evaluation
- V. Conclusion

Background: Physical vs Digital context

Digital world



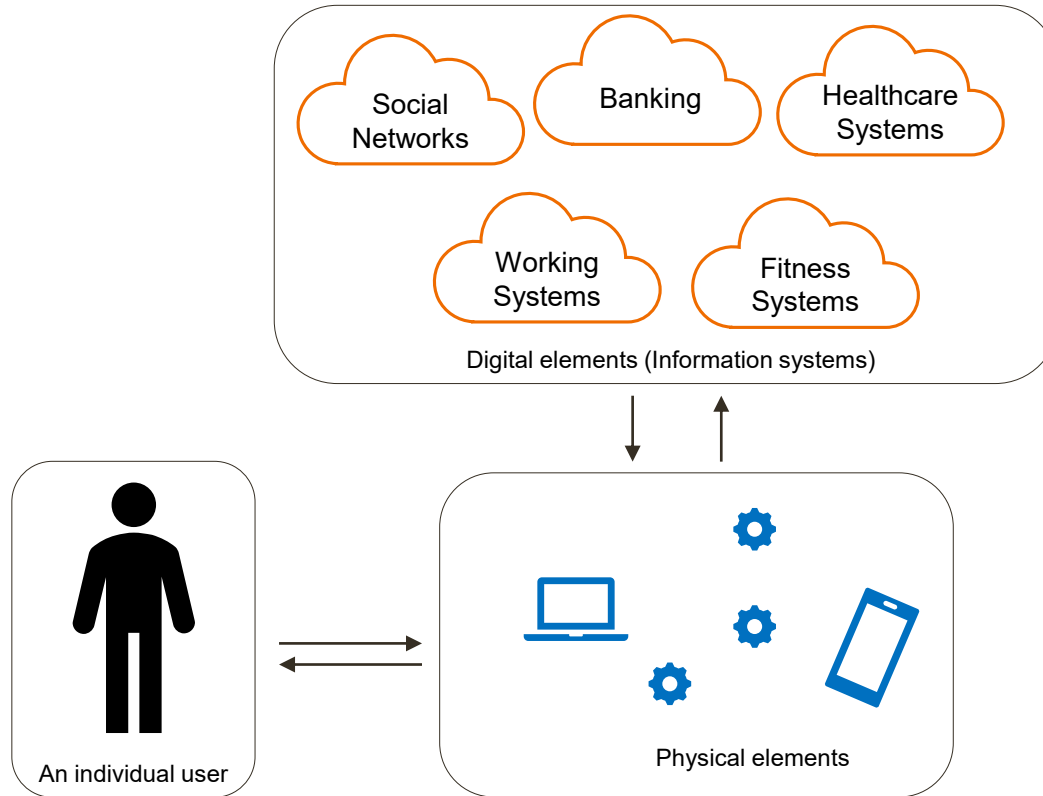
Physical world

Context: Any information that can be used to characterise the situation of an entity. [DEY01]
[DEY01] Dey, Anind K. "Understanding and using context." *Personal and ubiquitous computing* 5 (2001): 4-7.

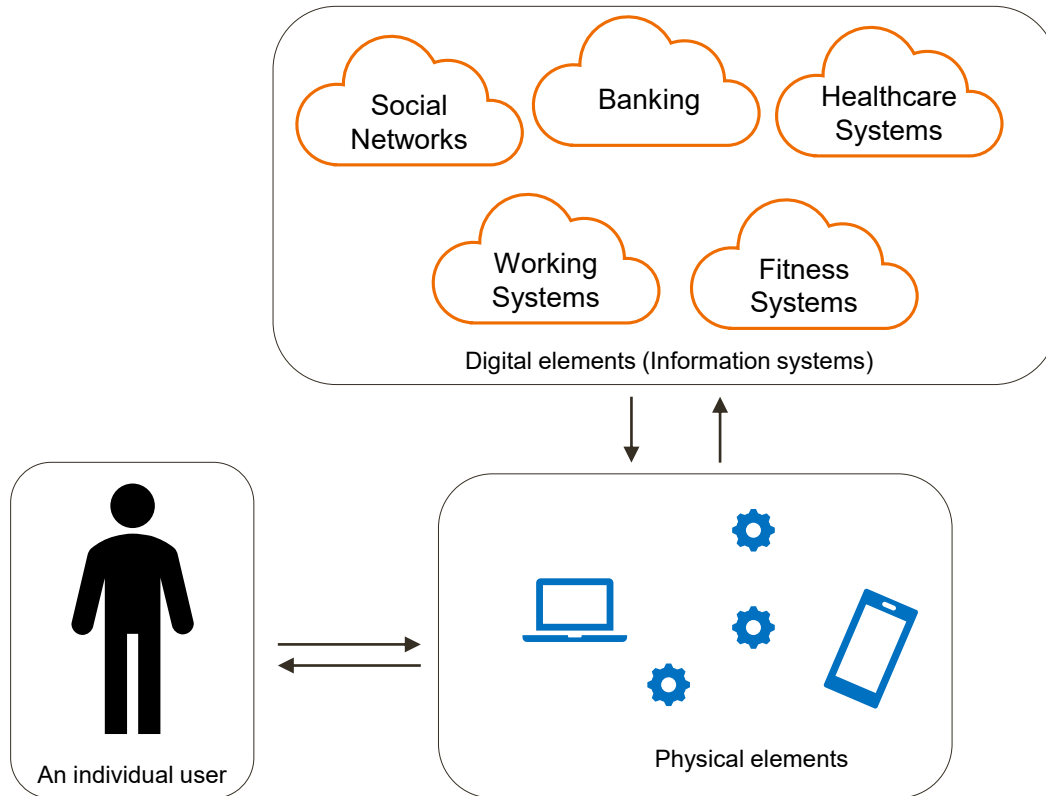
Outline

- I. Background
- II. State of The Art
- III. Proposed framework
- IV. Evaluation
- V. Conclusion

Individual-oriented Information System (IIS)



Individual-oriented Information System (IIS)



Limitations of security solutions in the literature

- System-specific
 - Require specific integration between client and server components
- Device-specific
 - Intrusive control and management
- **Fail to protect** against attacks where the adversary has sufficient knowledge to bypass these isolated solutions

The IIS considers the global context (digital + physical) of a user

SoTA: User physical contextual data in security solutions

1. IDS in mobile devices

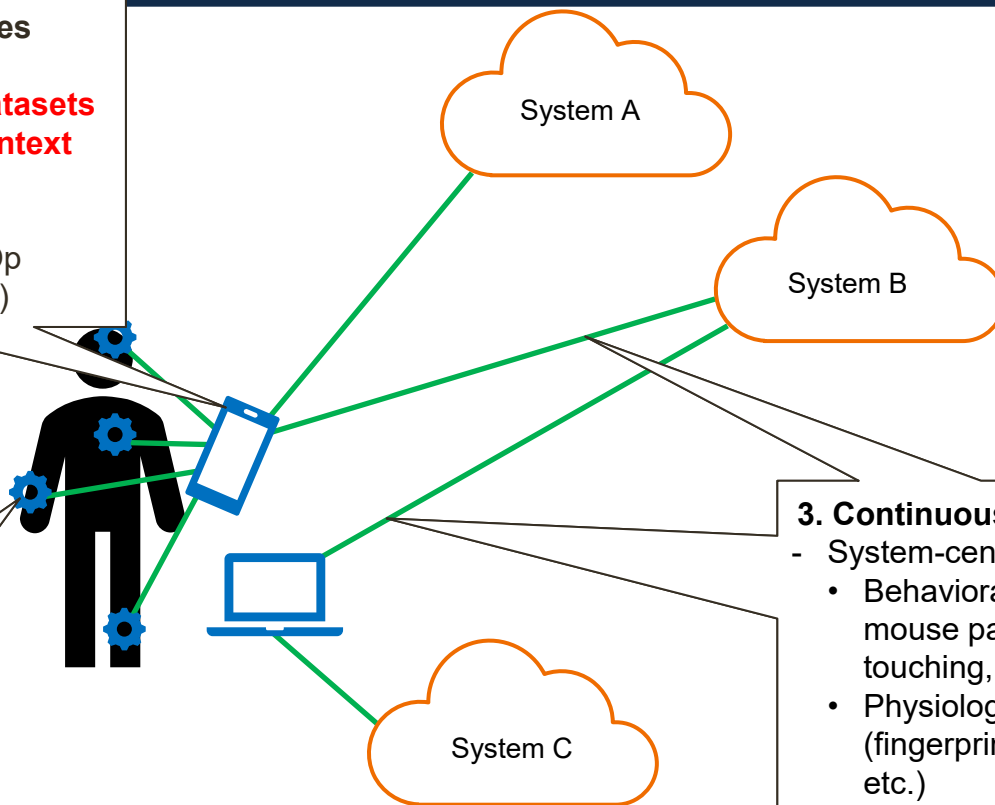
- Network-based IDS
 - **Neither studies nor datasets available with user context**
- Host-based IDS
 - Device-centric
 - Device context (CPU, Op codes, permissions, etc)

2. WBAN Security

- User physical data (biometrics) are used
 - Crypto Key management
 - Authentication
- ⇒ To protect WBAN systems

3. Continuous Authentication

- System-centric, leveraging
 - Behavioral data (key stroke, mouse pattern, gait, touching, etc.)
 - Physiological data (fingerprint, iris, blood, etc.)

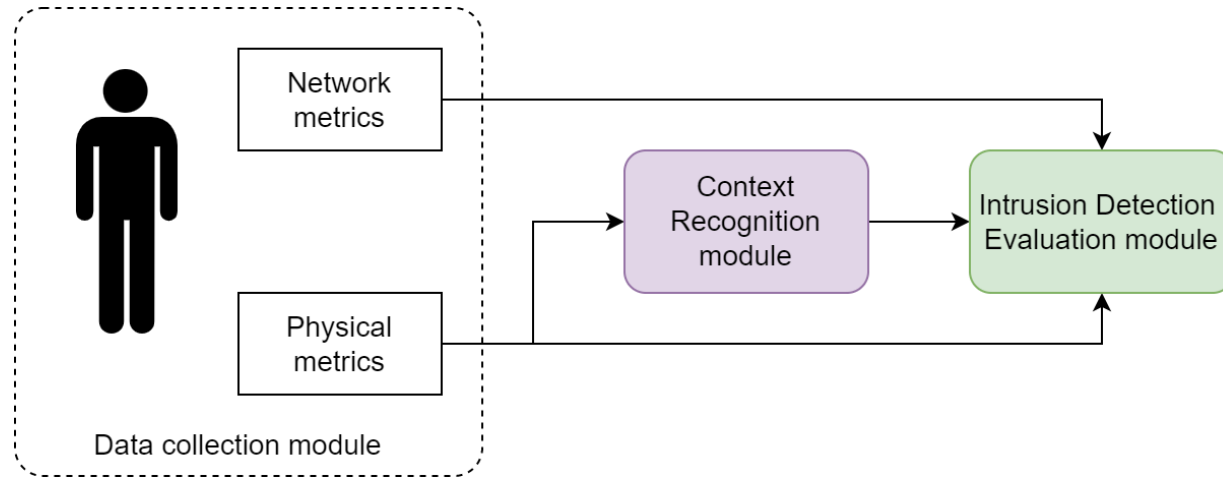


*Can a user's **physical contextual data** enhance the performance of **network intrusion detection**?*

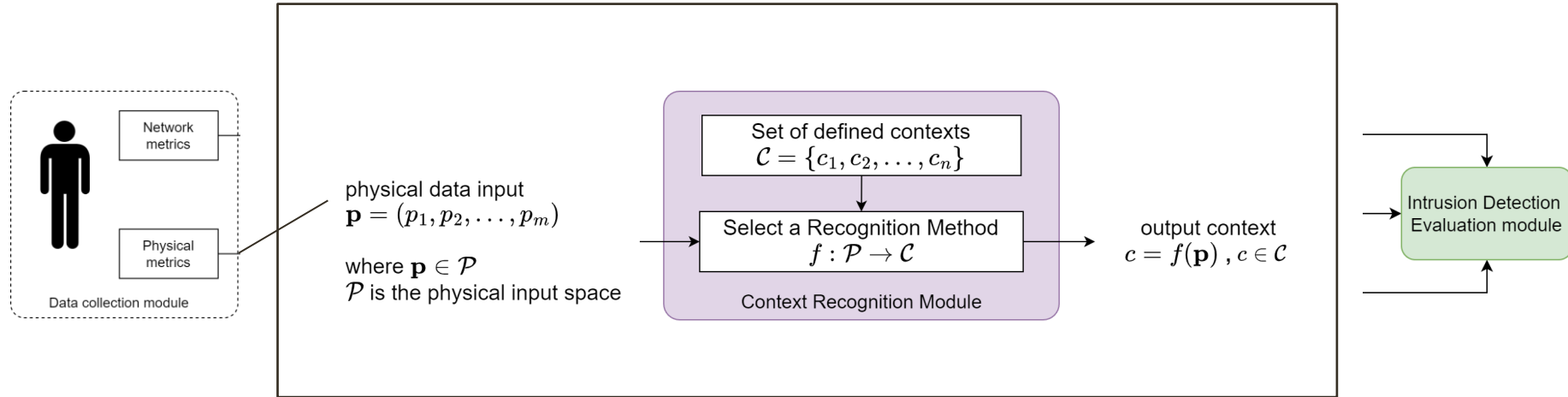
Outline

- I. Background
- II. State of The Art
- III. Proposed framework
- IV. Evaluation
- V. Conclusion

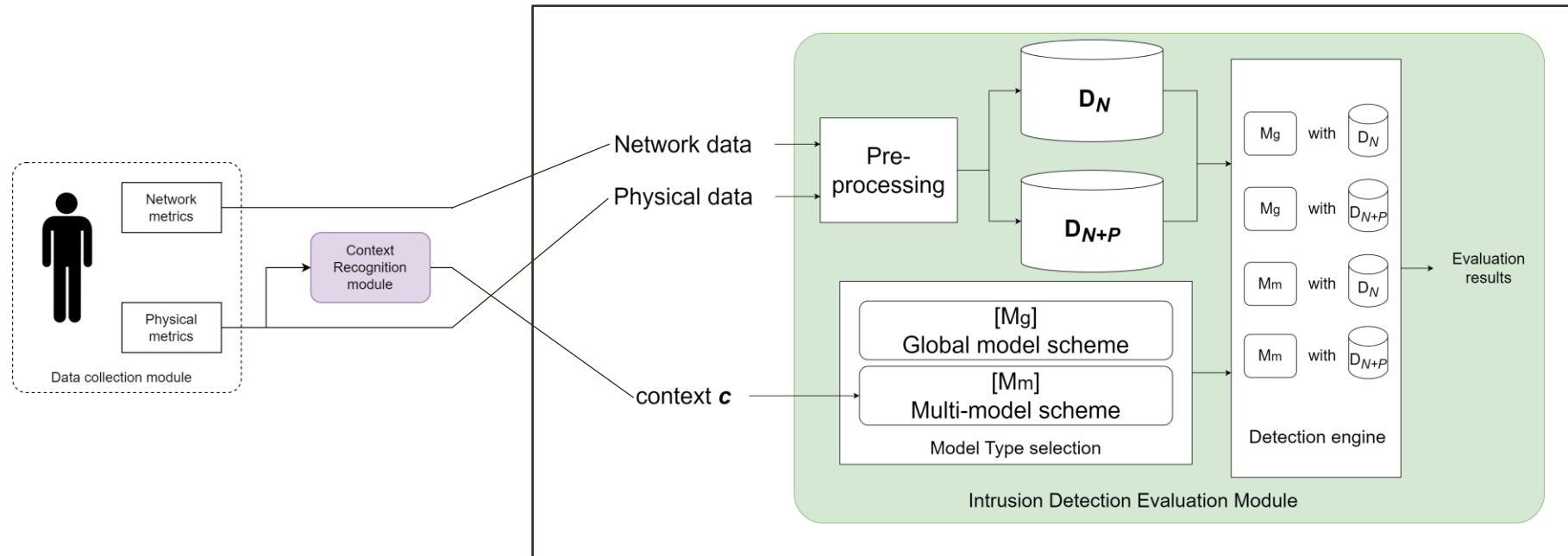
Overview



Context Recognition Module



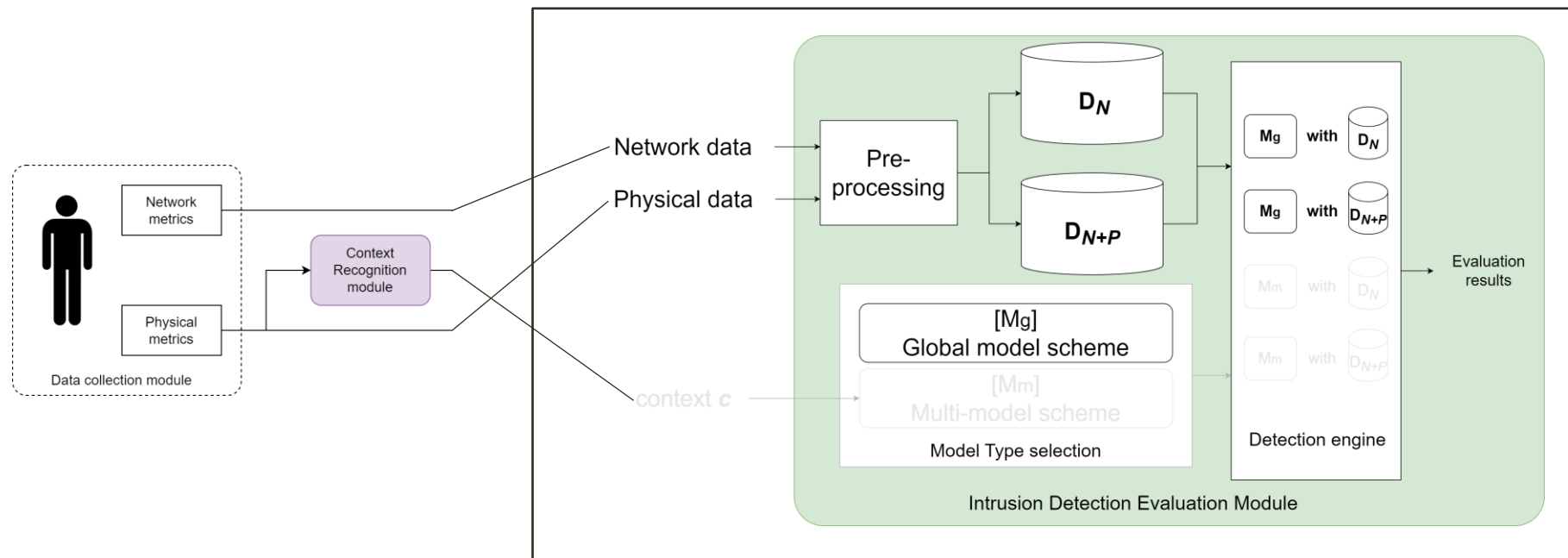
Intrusion Detection Evaluation Module



D_N : Dataset with network-only features (netflow-based features)

D_{N+P} : Dataset with network + physical features (accelerometer, ambient light, user speed)

Evaluation Scope



D_N : Dataset with network-only features (netflow-based features)

D_{N+P} : Dataset with network + physical features (accelerometer, ambient light, user speed)

Outline

- I. Background: Individual-oriented Information System (IIS)
- II. State of The Art
- III. Proposed framework
- IV. Evaluation
- V. Conclusion

Dataset

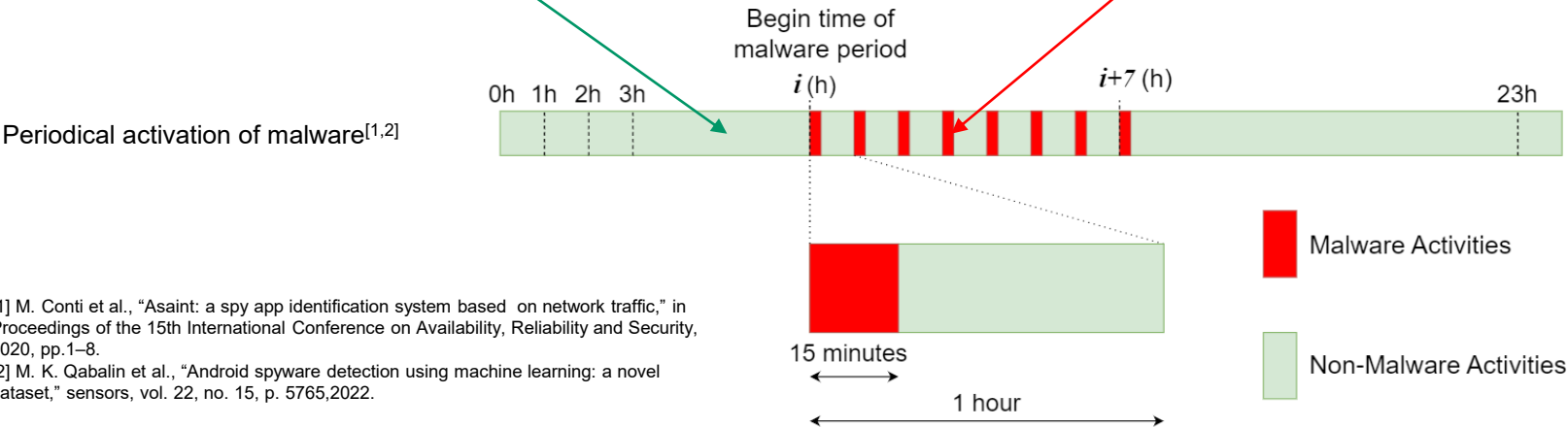
Benign data

Our in-situ data collection experiment
User's daily activities (2 x 24h),
From a personal smartphone

Malware data

Extracted from CIC-AndMal2017 dataset

- 5 Adware families
- 8 captures / family
- 15 minutes / capture



[1] M. Conti et al., "Asaint: a spy app identification system based on network traffic," in Proceedings of the 15th International Conference on Availability, Reliability and Security, 2020, pp. 1–8.

[2] M. K. Qabalin et al., "Android spyware detection using machine learning: a novel dataset," sensors, vol. 22, no. 15, p. 5765, 2022.

Injection and Evaluation Pipeline

With each malware in [Ewind, Feiwo, Gooligan, Kemoge, Youmi]:

Inject 8 captures starting from hour i ($i = 0, 1, \dots, 23$)

With each injection: Do the detection evaluation

Repeat 10 times:

Randomly split the dataset into train/validate/test set

Create a machine learning model (**XGBoost**)

Hyper-parameter turning

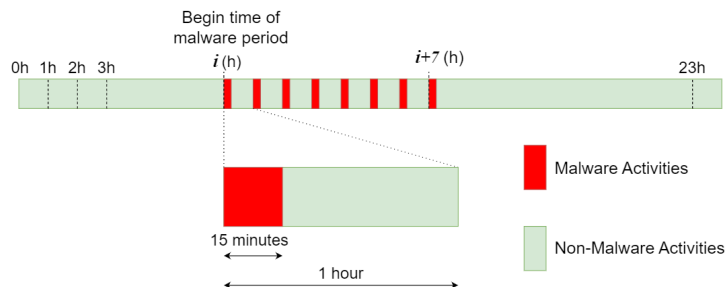
Cross Validation

Record the evaluation metric (**PR AUC**) on testing set

Aggregate the metrics over 10 trials

Aggregate the results over 24 injection of current malware

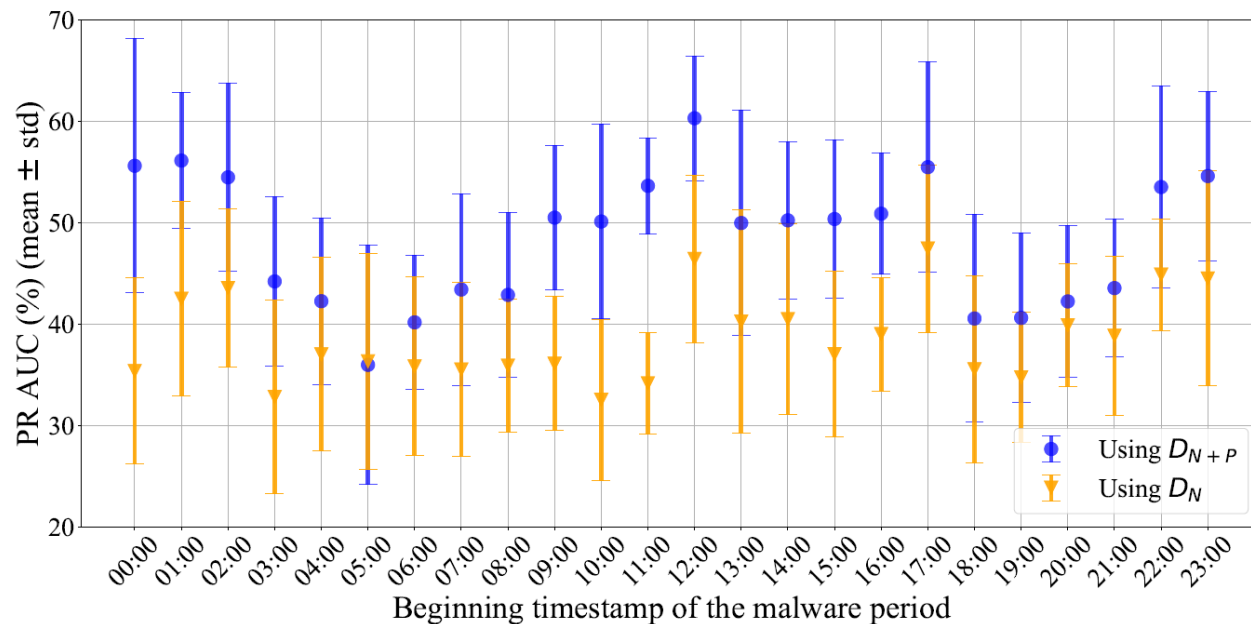
Aggregate the results across all malware families



Imbalanced dataset

Results (1/4)

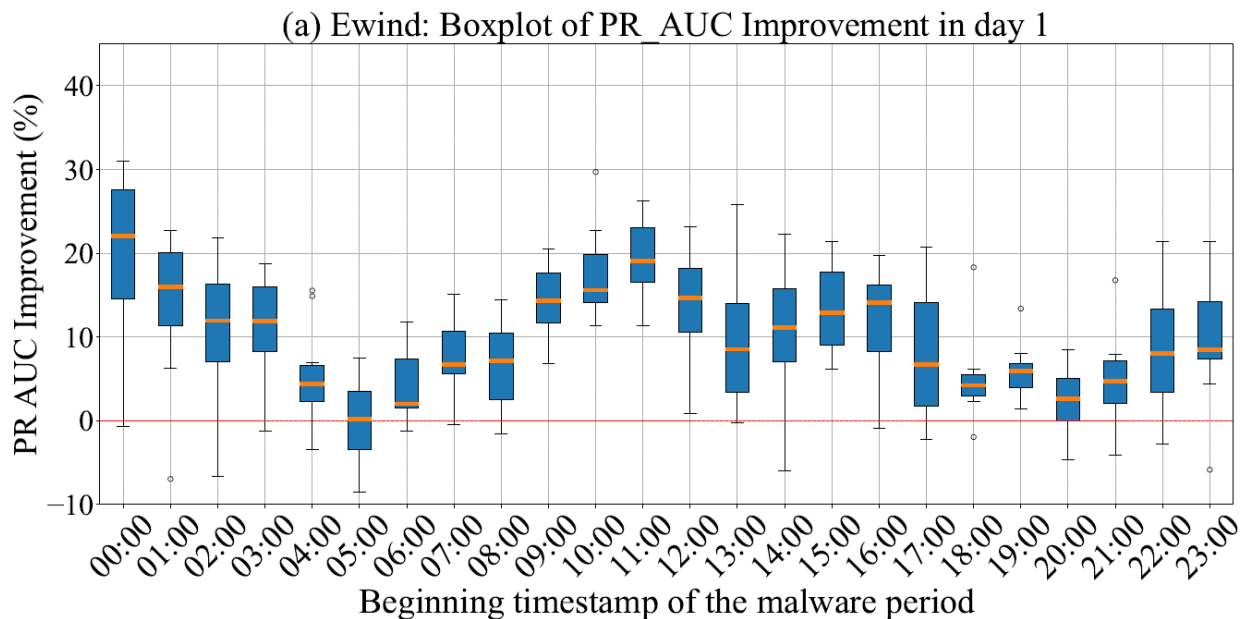
(a) Comparison of PR AUC between using D_{N+P} and D_N



Note: $\text{PR AUC}_{\text{Random guesser}} = \# \text{ positive samples} / \# \text{ all samples} (= 8.3\% \text{ in this work})$

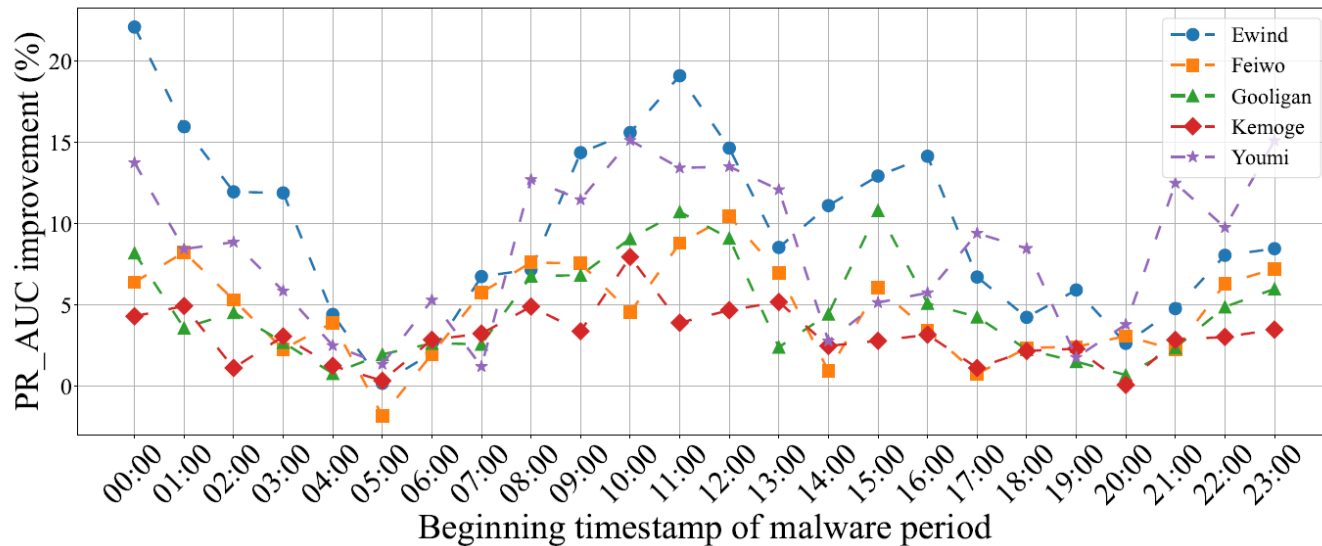
Results (2/4)

(b) PR AUC Improvement (= PR AUC using D_{N+P} - PR AUC using D_N)



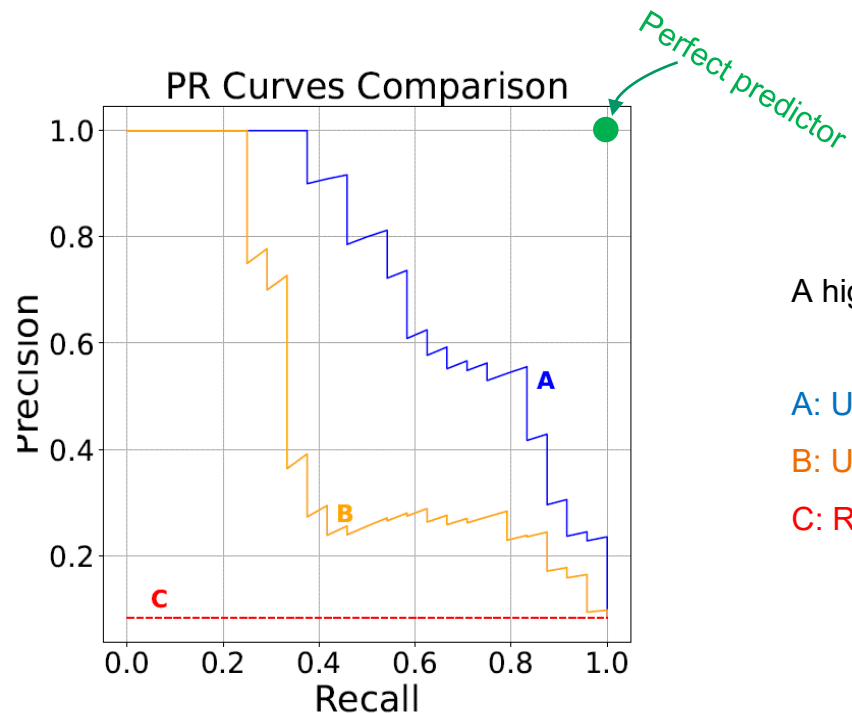
Results (3/4)

(c) PR AUC Improvement (mean values) across 5 malware families



Results (4/4)

(d) PR curve comparison at a specific timestamp (00h00, day 1)



A higher PR AUC indicates better performance.

A: Using D_{N+P}

PR AUC = 0.7393

B: Using D_N

PR AUC = 0.4767

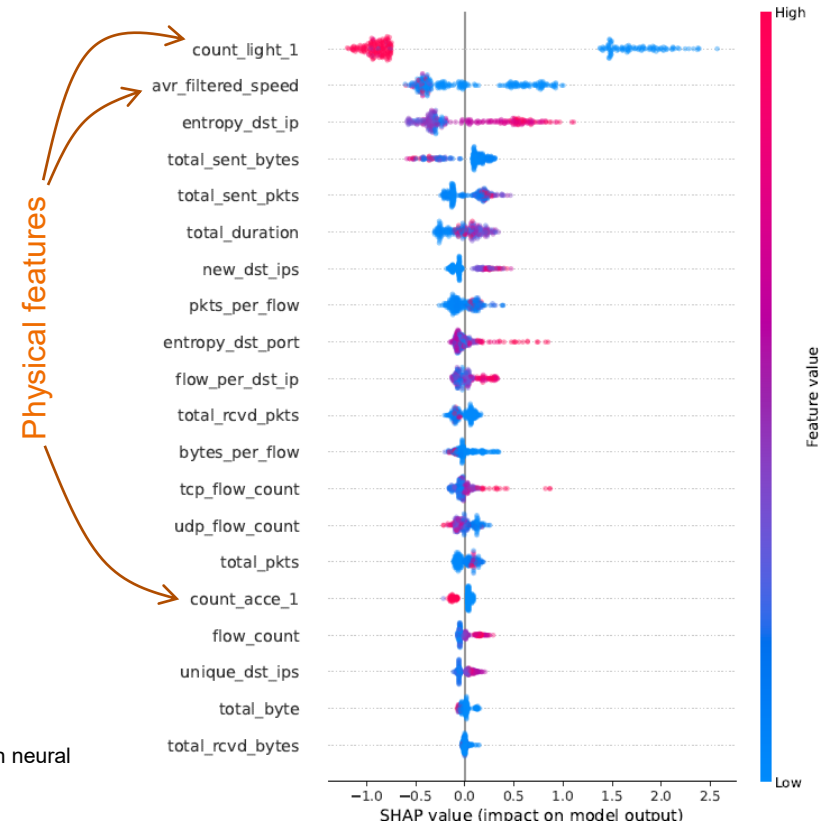
C: Random guesser

PR AUC = 0.0830

Explanation by SHAP

- The **SHAP** (SHapley Additive exPlanations) methodology^[3]
 - A game-theoretic approach
 - Explain the output of machine learning models
 - Assign importance values to individual variables (features)
- ⇒ Specific physical signals are critical (ambient light, user speed)
- ⇒ The relevance of physical data in the decision-making of the IDS algorithm

[3] S. M. Lundberg and S.-I. Lee, "A unified approach to interpreting model predictions," Advances in neural information processing systems, vol. 30, 2017.



Outline

- I. Background
- II. State of The Art
- III. Proposed framework
- IV. Evaluation
- V. Conclusion

Conclusion

Our contributions

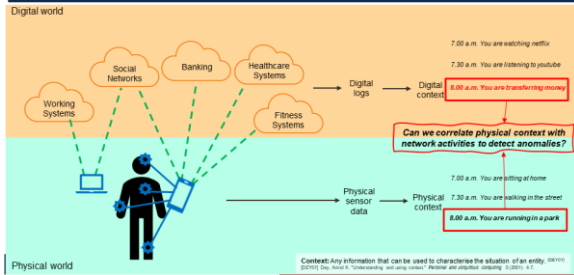
- **A new dataset** combining network traffic and physical sensor data collected from a real person's daily activities
- **A framework** leveraging user physical context data in network intrusion detection systems
- **An experimental validation** of the hypothesis that physical contextual information enhances NIDS performance

Ongoing work

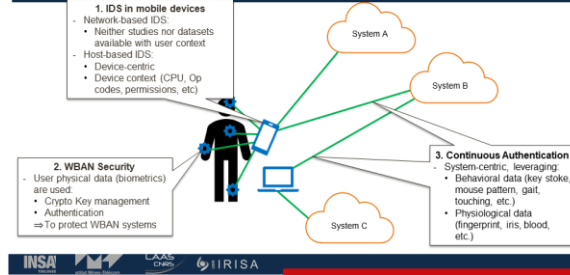
- A **large-scale dataset collection** campaign at IMT Atlantique
- Unsupervised Learning methods; Multi-model Validation

Thank you!

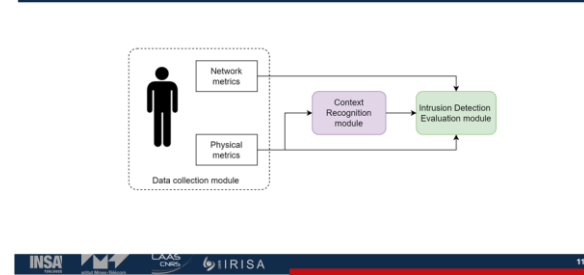
Background



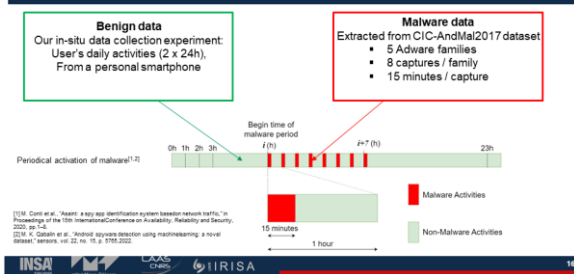
SoTA: User physical contextual data in security solutions



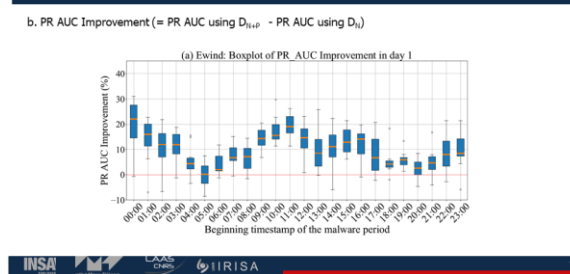
Overview



Dataset



Results (2/4)



Explanation by SHAP

