

TESSERACT

Eliminating Experimental Bias in Malware Classification across Space and Time

Feargus Pendlebury*, Fabio Pierazzi*, Roberto Jordaney, Johannes Kinder, and Lorenzo Cavallaro

<https://s2lab.kcl.ac.uk/projects/tesseract>

Fabio Pierazzi

<fabio.pierazzi@kcl.ac.uk> – @fbpierazzi

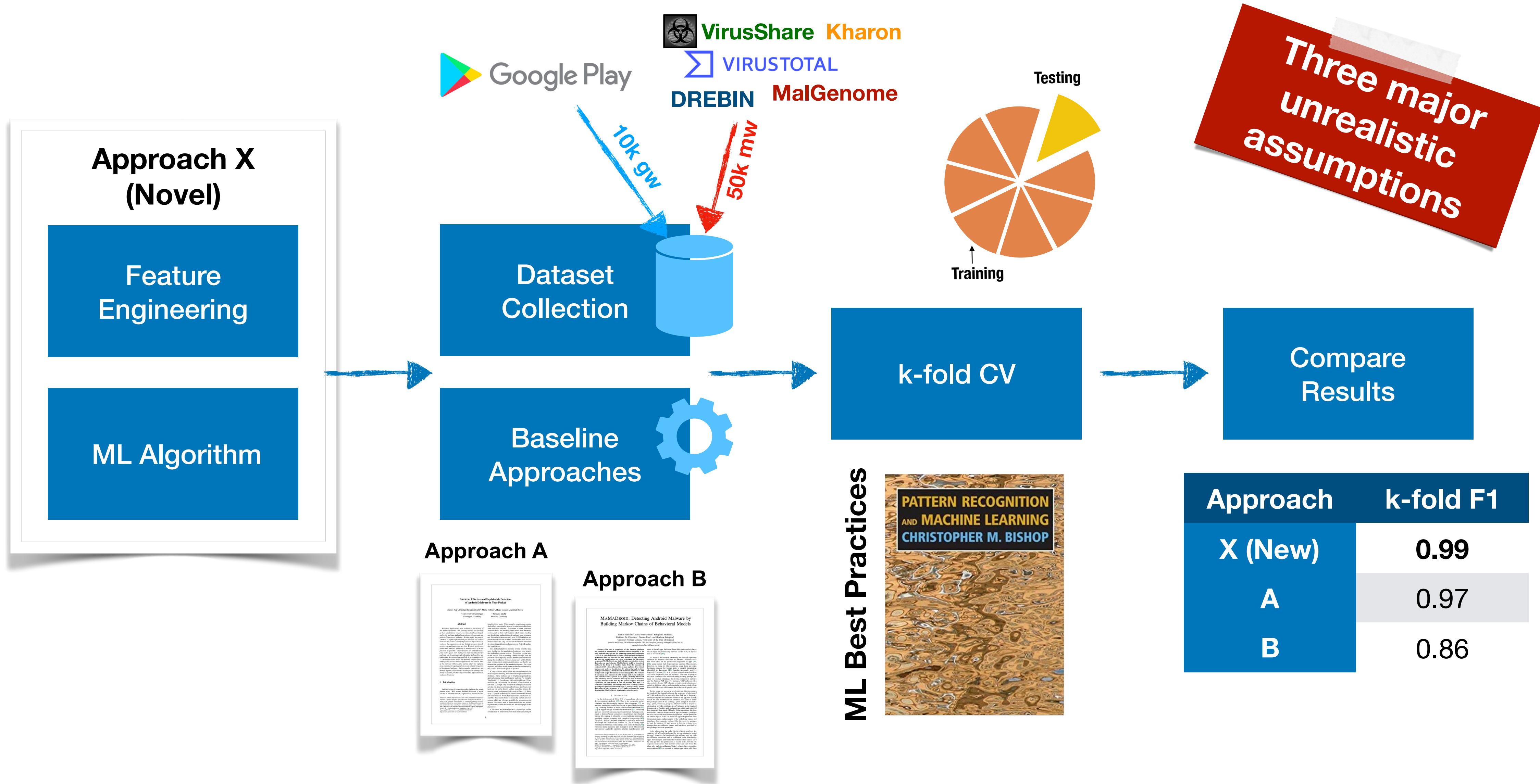
<https://fabio.pierazzi.com>

Aug 14-16, 2019

USENIX Security Symposium

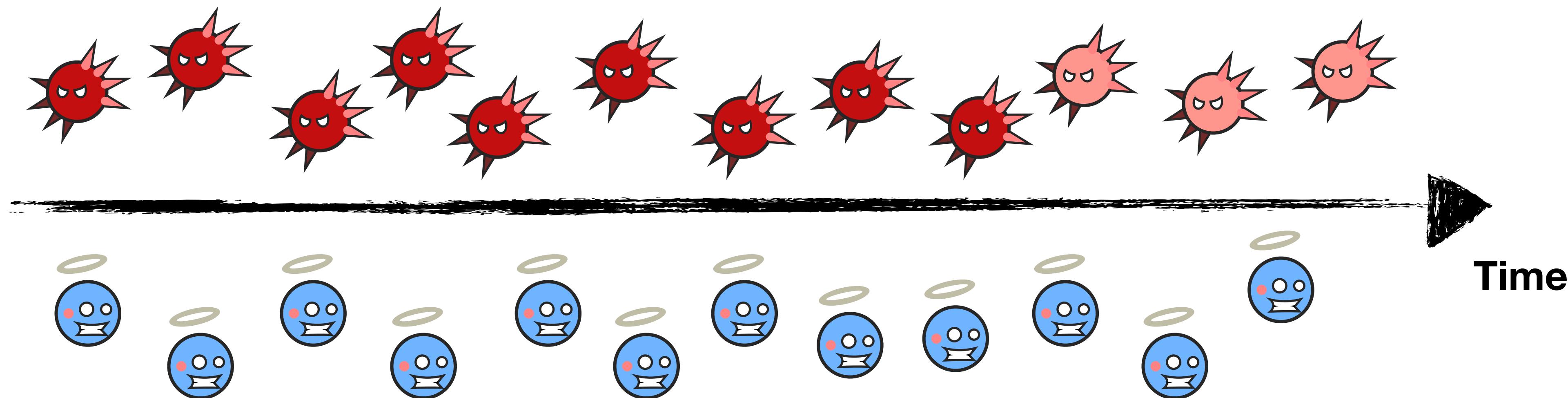
* Equal contribution

ML for Malware Detection



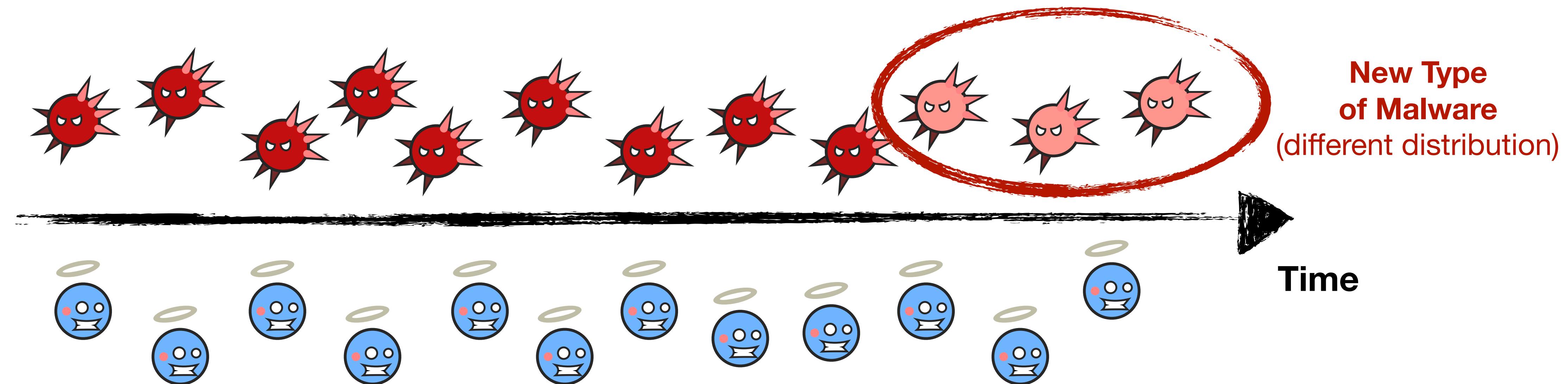
Sources of Experimental Bias (1/3)

Temporal Inconsistency in Train/Test Sets



Sources of Experimental Bias (1/3)

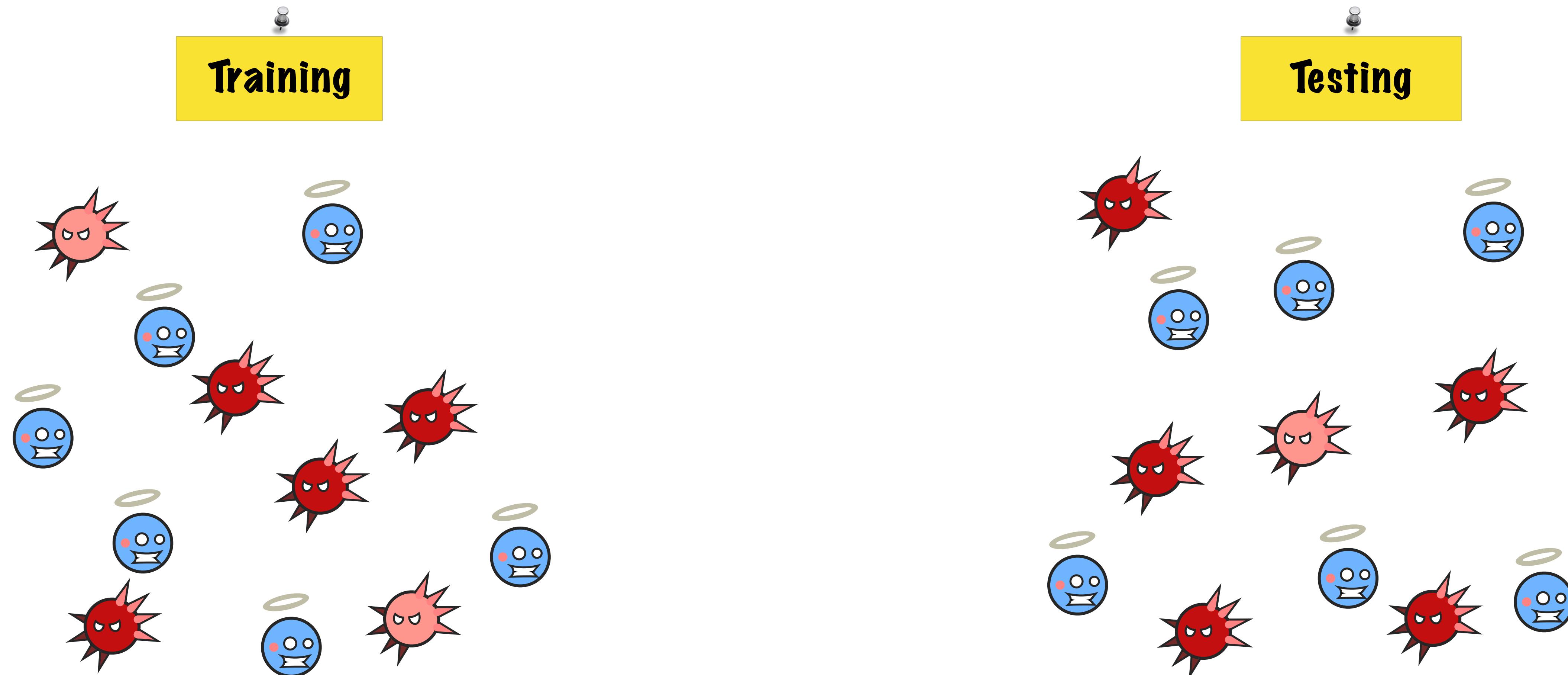
Temporal Inconsistency in Train/Test Sets



Sources of Experimental Bias (1/3)

Temporal Inconsistency in Train/Test Sets

Violations use future knowledge in training

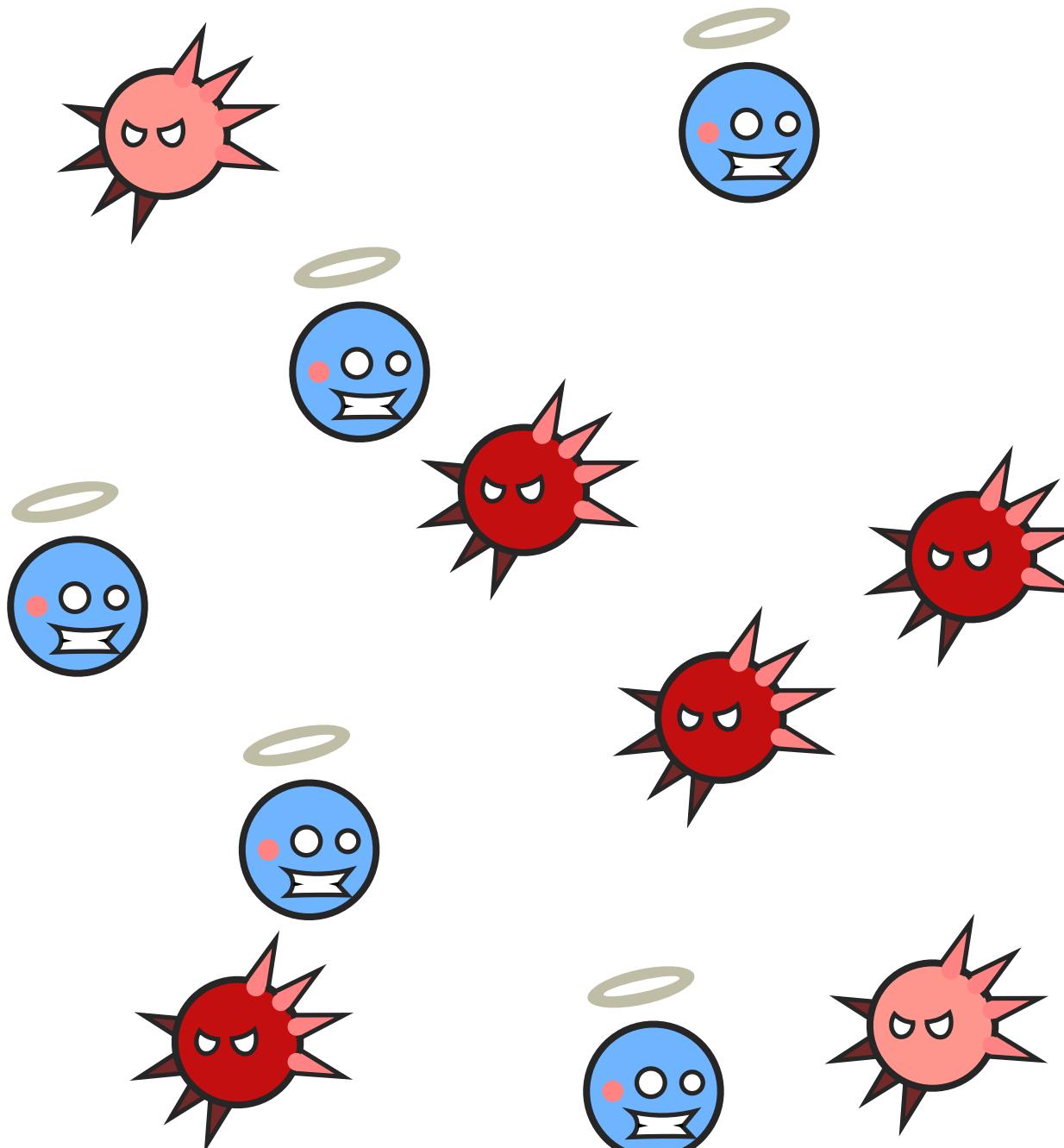


Sources of Experimental Bias (1/3)

Temporal Inconsistency in Train/Test Sets

Violations use future knowledge in training

Training



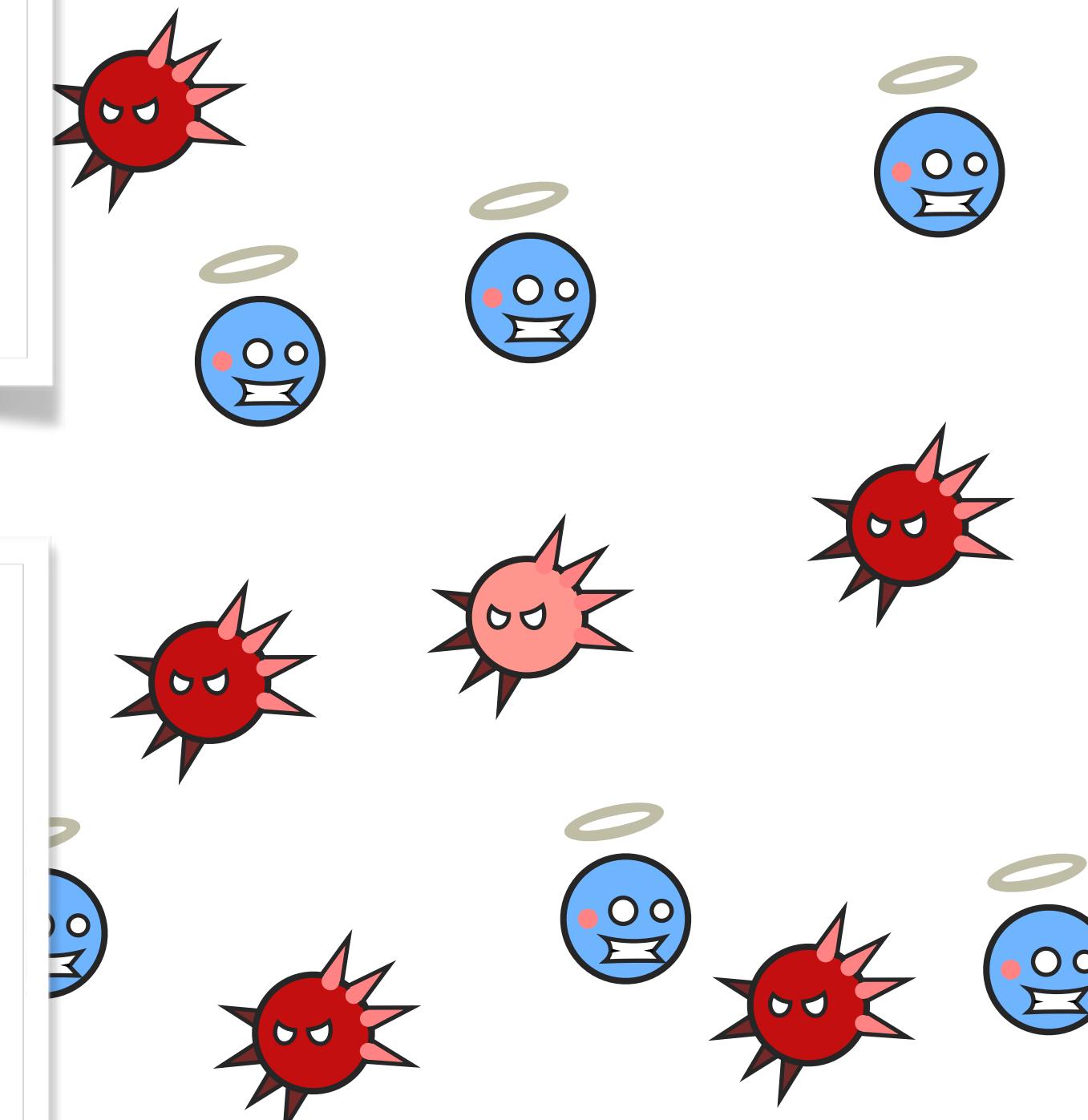
Kevin Allix et al. [ESSoS 2016]

Are Your Training Datasets Yet Relevant?
An Investigation into the Importance of Timeline in
Machine Learning-based Malware Detection

Kevin Allix, Tegawendé F. Bissyandé, Jacques Klein, and Yves Le Traon

SnT - University of Luxembourg

Testing



Brad Miller et al. [DIMVA 2016]

Reviewer Integration and Performance
Measurement for Malware Detection

Brad Miller^{1†}, Alex Kantchelian², Michael Carl Tschantz³, Sadia Afroz³,
Rekha Bachwani^{4‡}, Riyaz Faizullaboy², Ling Huang⁵, Vaishaal Shankar²,
Tony Wu², George Yiu^{6†}, Anthony D. Joseph², and J. D. Tygar²

¹ Google Inc. bradmiller@google.com

² UC Berkeley {akant, riyazdf, vaishaal, tony.wu, adj, tygar}@cs.berkeley.edu

³ International Computer Science Institute {mct, sadia}@icsi.berkeley.edu

⁴ Netflix rbachwani@netflix.com

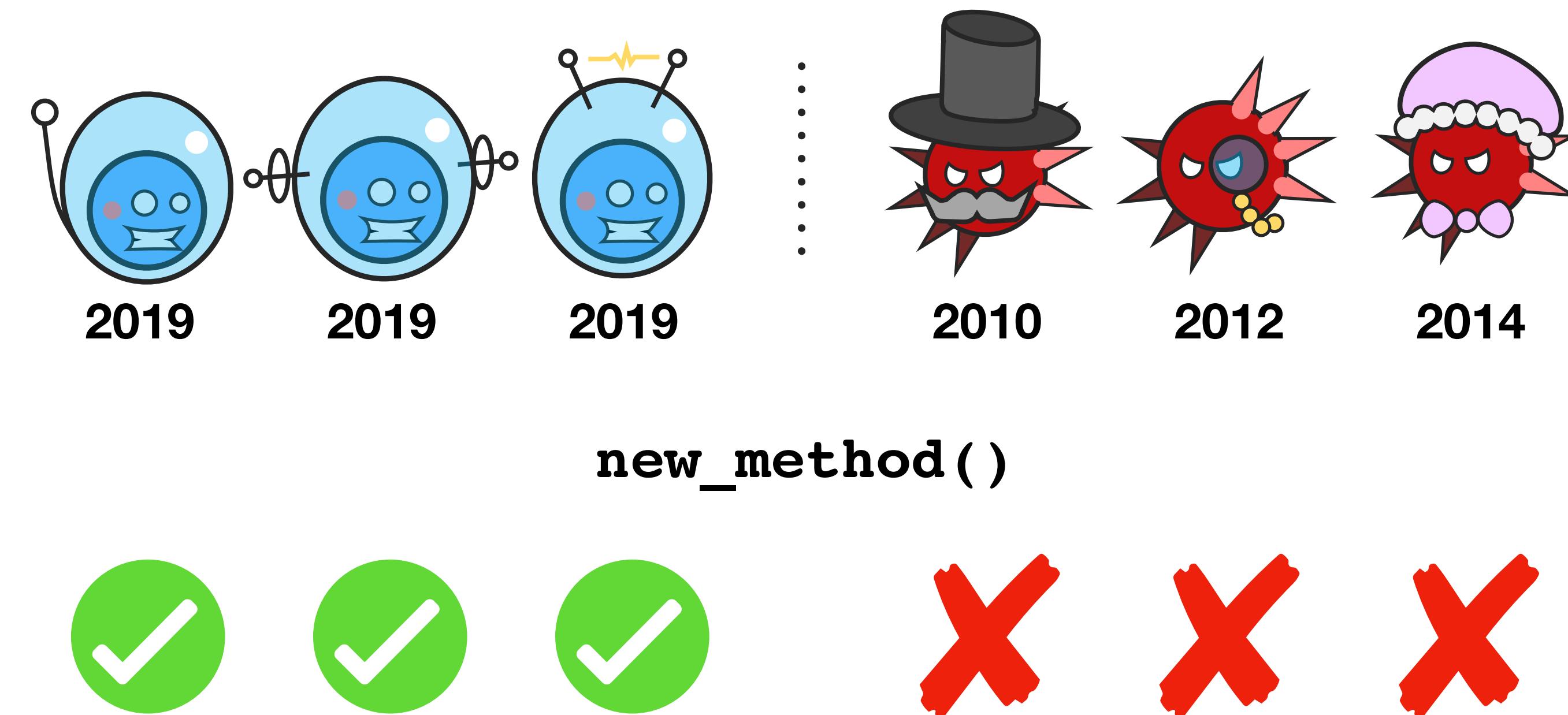
⁵ DataVisor ling.huang@datavisor.com

⁶ Pinterest george@pinterest.com

Sources of Experimental Bias (2/3)

Temporal {good|mal}ware inconsistency

Violations may learn artifacts

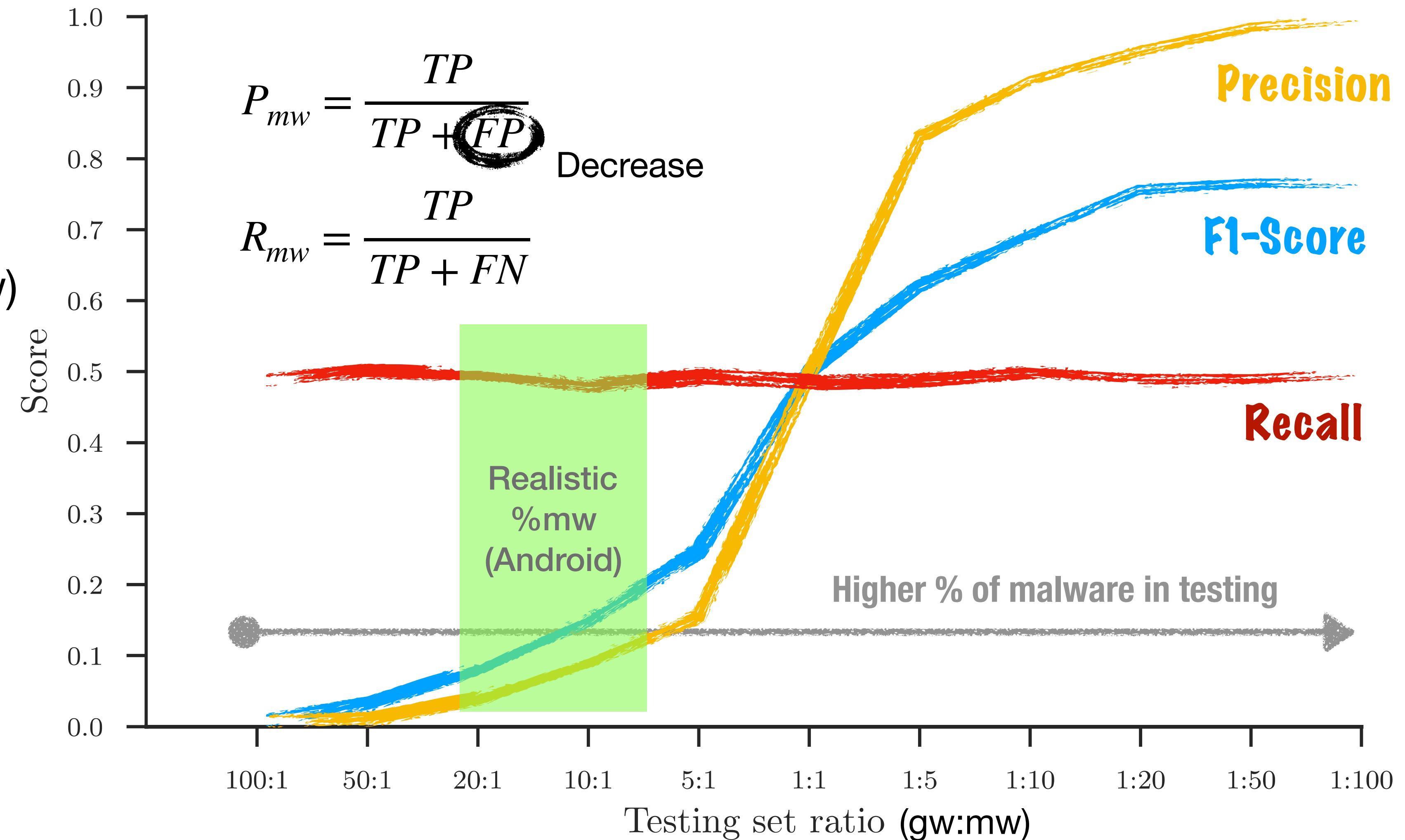


Sources of Experimental Bias (3/3)

Unrealistic Test Class Ratio

- Training set: Fixed
- Testing set: Varying % of mw (by downsampling gw)

Violations produce unrealistic results

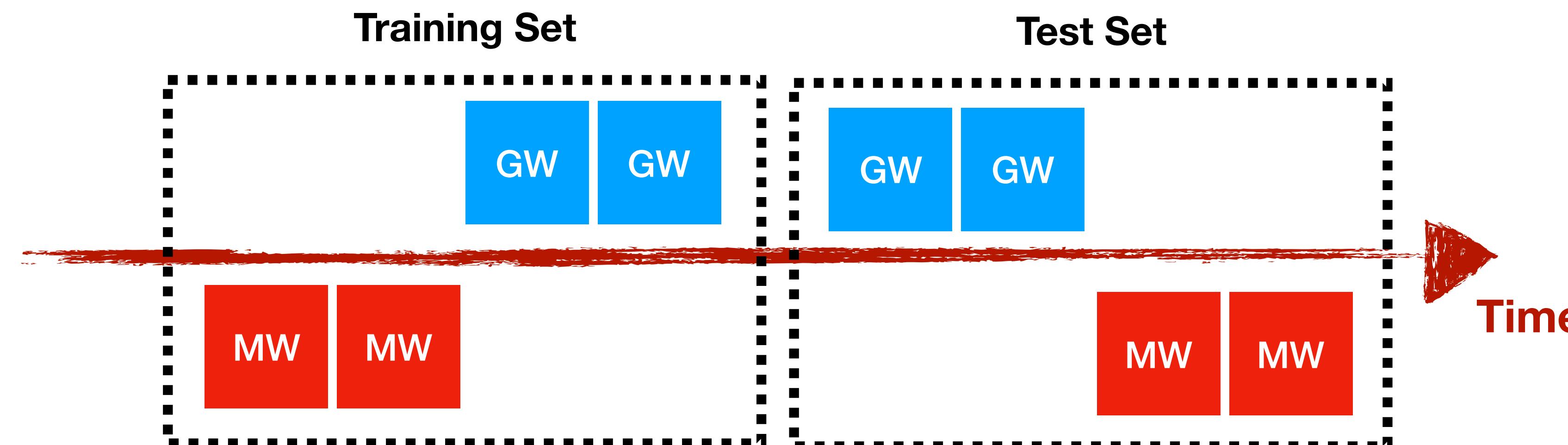


TESSERACT Framework

Experimental Constraints

C1 Temporal training consistency

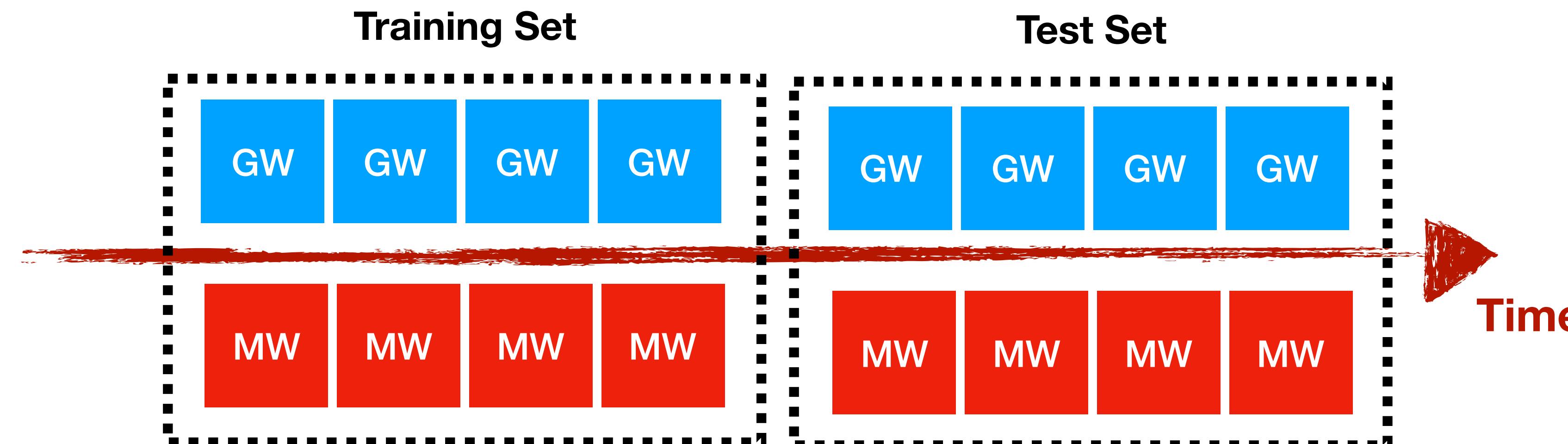
→ time(training) < time(testing)



TESSERACT Framework

Experimental Constraints

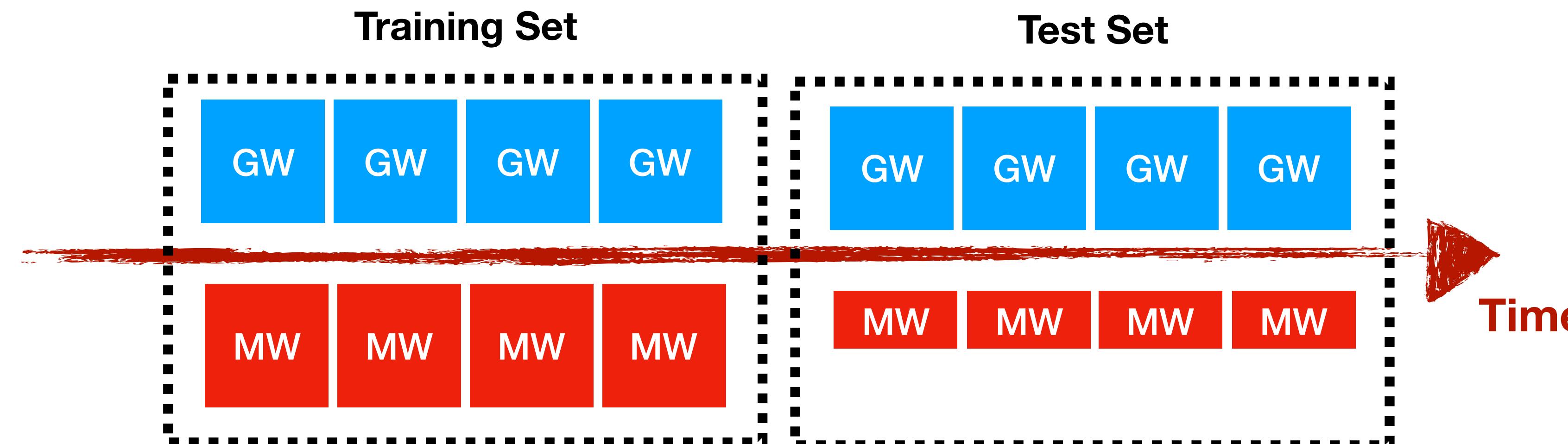
- C1** Temporal training consistency → time(training) < time(testing)
- C2** {good|mal}ware temporal consistency → time(gw) = time(mw)



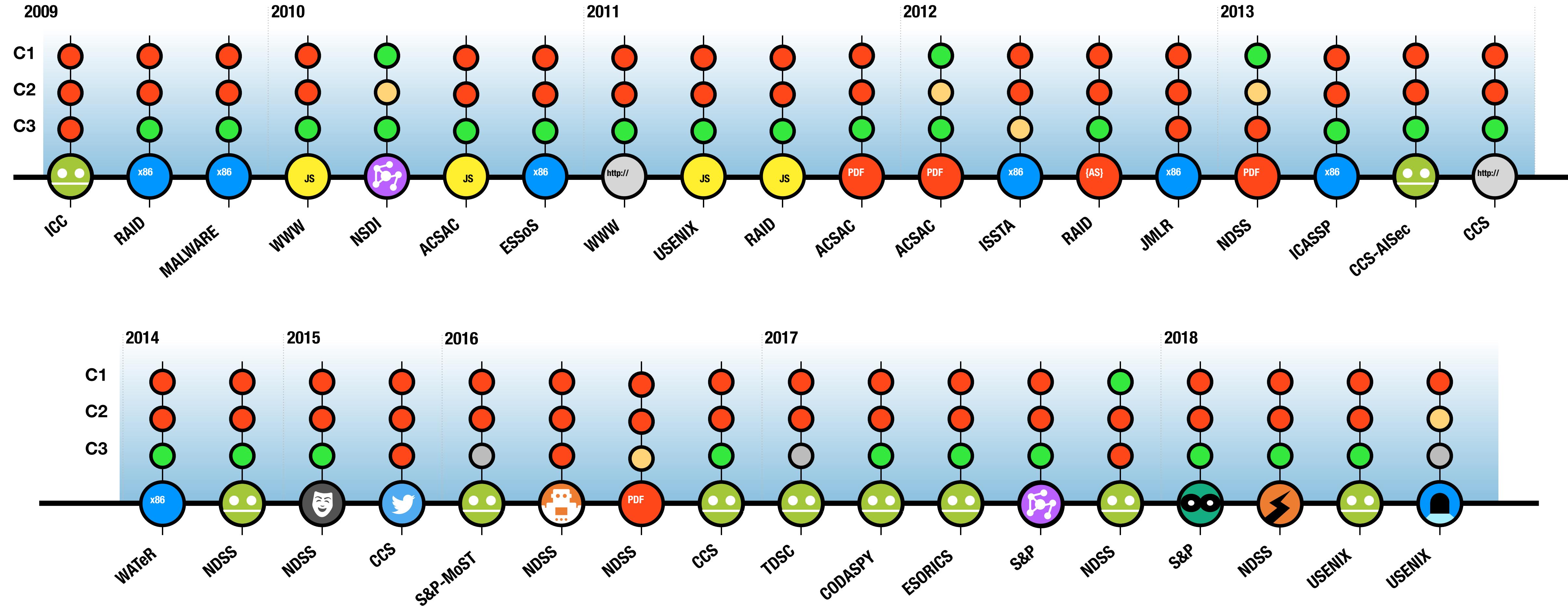
TESSERACT Framework

Experimental Constraints

C1	Temporal training consistency	→ time(training) < time(testing)
C2	{good mal}ware temporal consistency	→ time(gw) = time(mw)
C3	Realistic testing classes ratio	→ realistic %mw in test

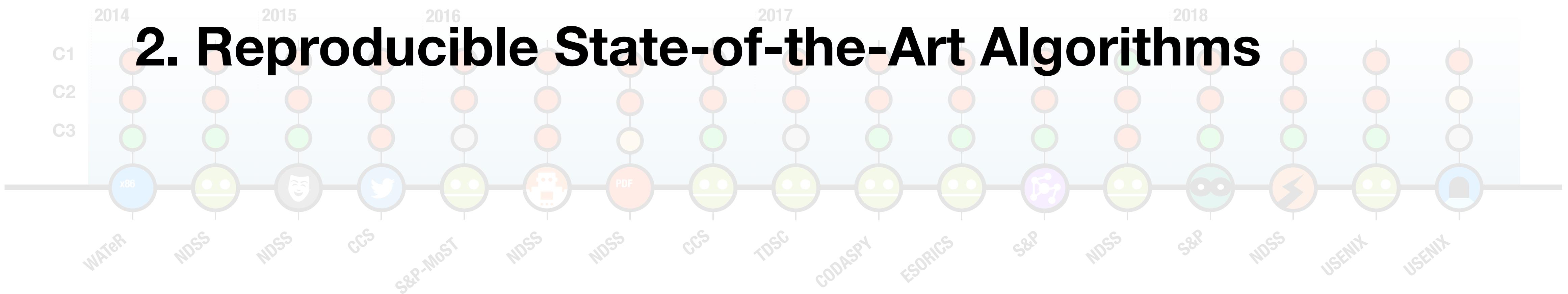
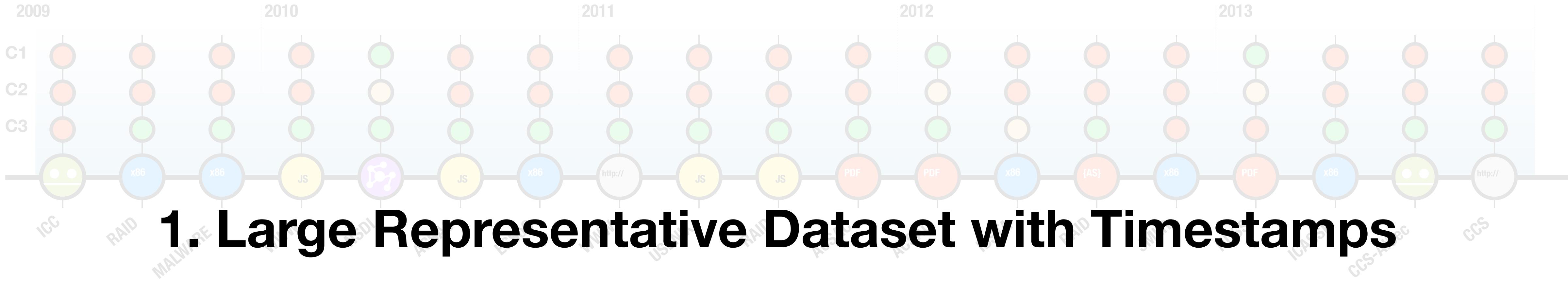


Endemic Problem



Details: <https://s2lab.kcl.ac.uk/projects/tesseract/poster-references.pdf>

Endemic Problem



Details: <https://s2lab.kcl.ac.uk/projects/tesseract/poster-references.pdf>

Dataset

- **129,729** Android applications from **AndroZoo** [1]
- **10%** malware
- Covering **3 years** (2014 to 2016)

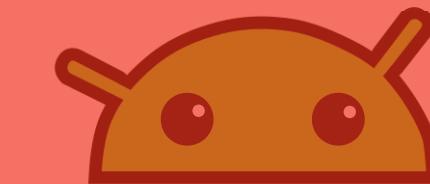
[1] K. Allix, T. F. Bissyandé, J. Klein, and Y. Le Traon.
AndroZoo: Collecting Millions of Android Apps for the Research Community
In *Mining Software Repositories (MSR)*, 2016

Benchmark Algorithms

Representative

Reproducible

Algorithm 1: DREBIN [NDSS14]



- Statically extracted features are bit vectors, present (1) or not (0)
- Linear Support Vector Machine (**SVM**) classifier

Algorithm 2: MaMaDroid [NDSS17]



- Markov chains from static call graphs through flow analysis
- Features are transition matrices of the Markov chains
- Random Forest (**RF**) classifier

Algorithm 3: Deep Learning [ESORICS17]



- DREBIN features and dataset
- **Deep Feed-Forward Neural Network**

[NDSS14] Arp et al., *DREBIN: Effective and Explainable Detection of Android Malware in Your Pocket*.

[NDSS17] Mariconti et al., *MaMaDroid: Detecting Android Malware by Building Markov Chains of Behavioural Models*.

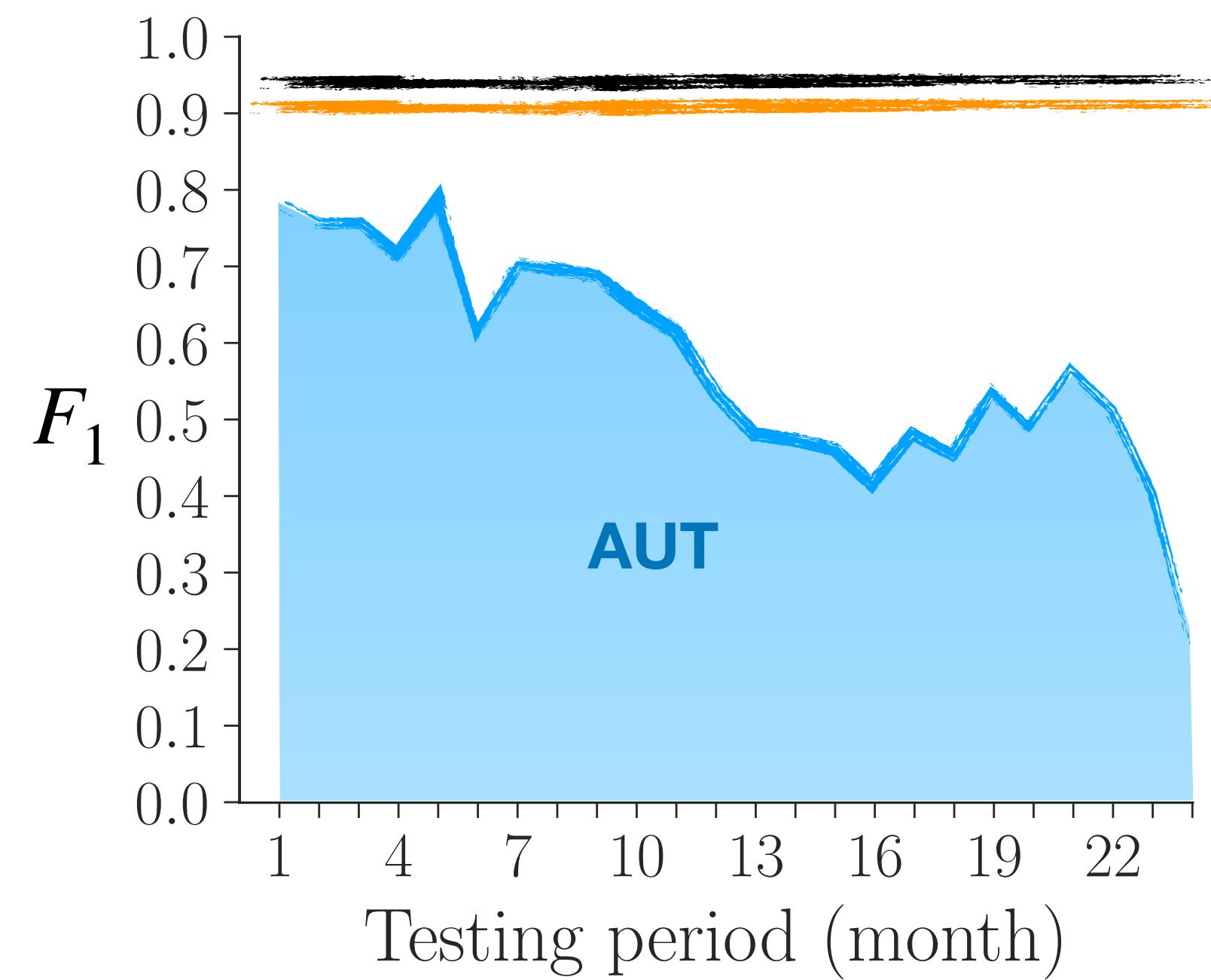
[ESORICS17] Grosse et al., *Adversarial Examples for Malware Detection*.

TESSERACT Evaluations

Experimental Constraints

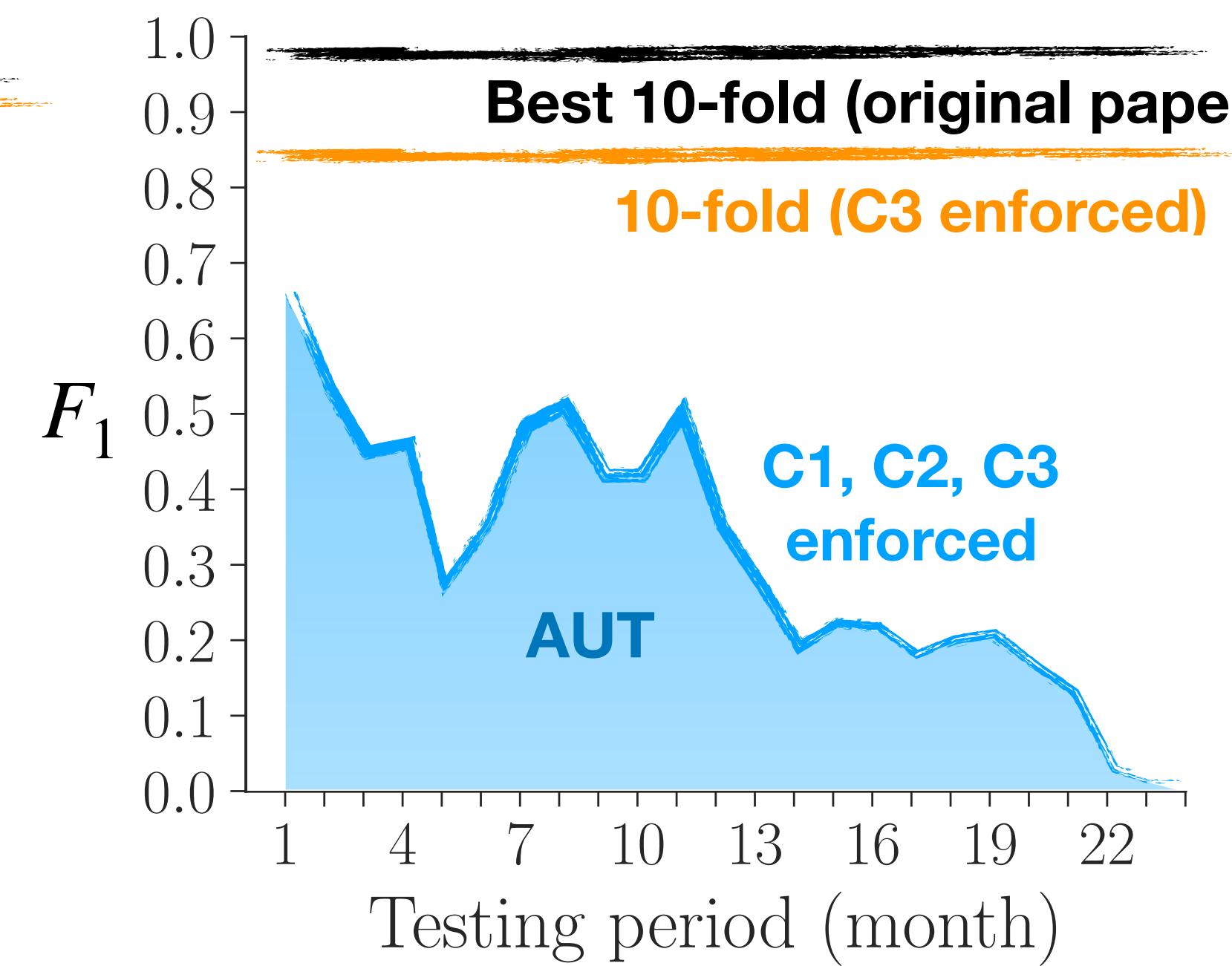
- C1** Temporal training consistency
- C2** {good|mal}ware temporal consistency
- C3** Realistic testing classes ratio

NDSS14



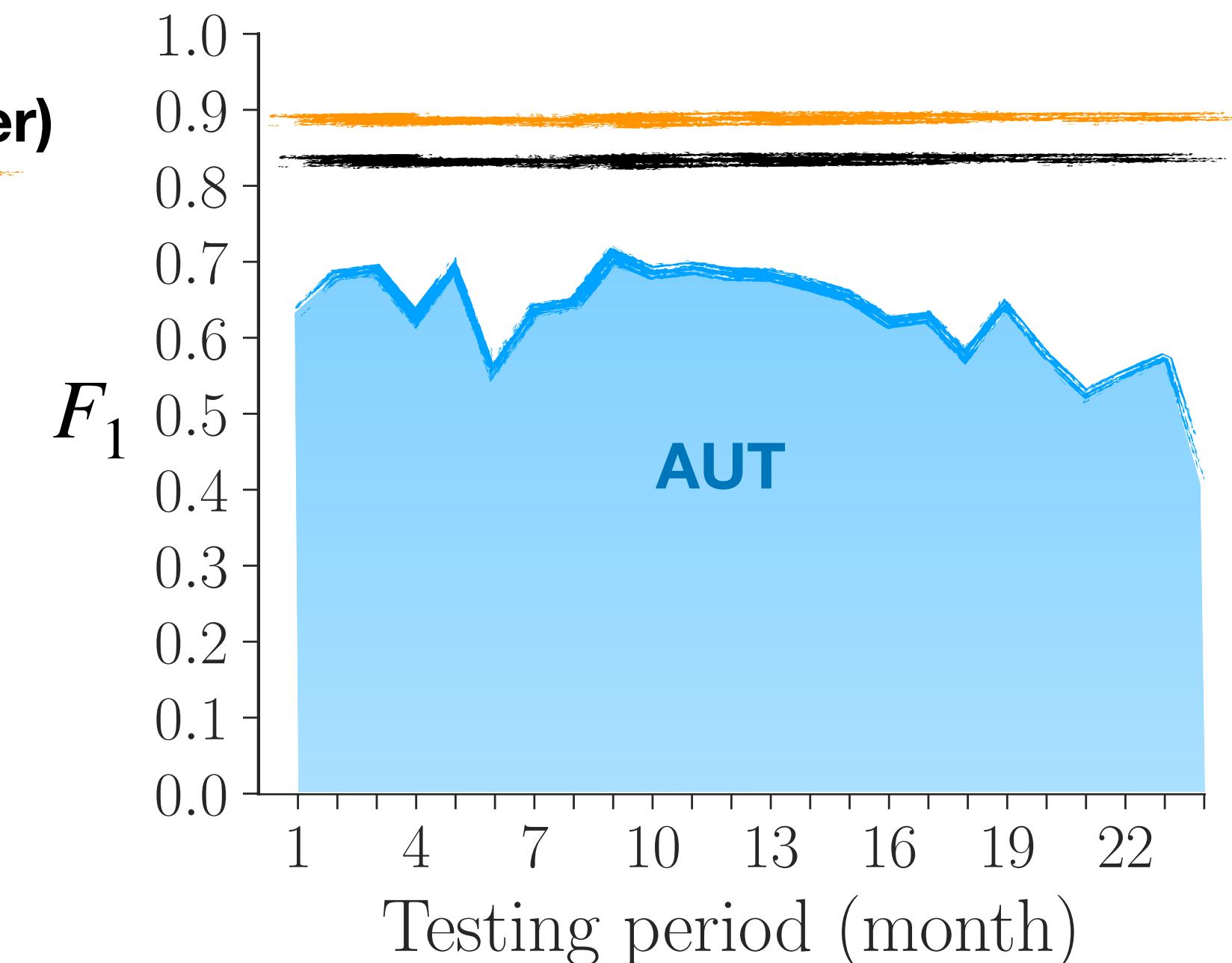
$$AUT(F_1, 24m) = 0.58$$

NDSS17



$$AUT(F_1, 24m) = 0.32$$

ESORICS17



$$AUT(F_1, 24m) = 0.64$$

TESSERACT: Actionable Points

Realistic Evaluations

- Unveils performance in realistic deployment
- Removes space-time experimental bias
- **Practitioners:** Choose Best Solution
- **Researchers:** Evaluate New Solutions

Tunable
Parameters
(e.g., time window,
granularity)

Performance-Cost Trade Offs

- **Detection Performance** (e.g., AUT F1)
- **Labeling Cost** for retraining (e.g., manpower)
- **Quarantine Cost** for rejection (e.g., low-confidence decisions)

More
in Paper

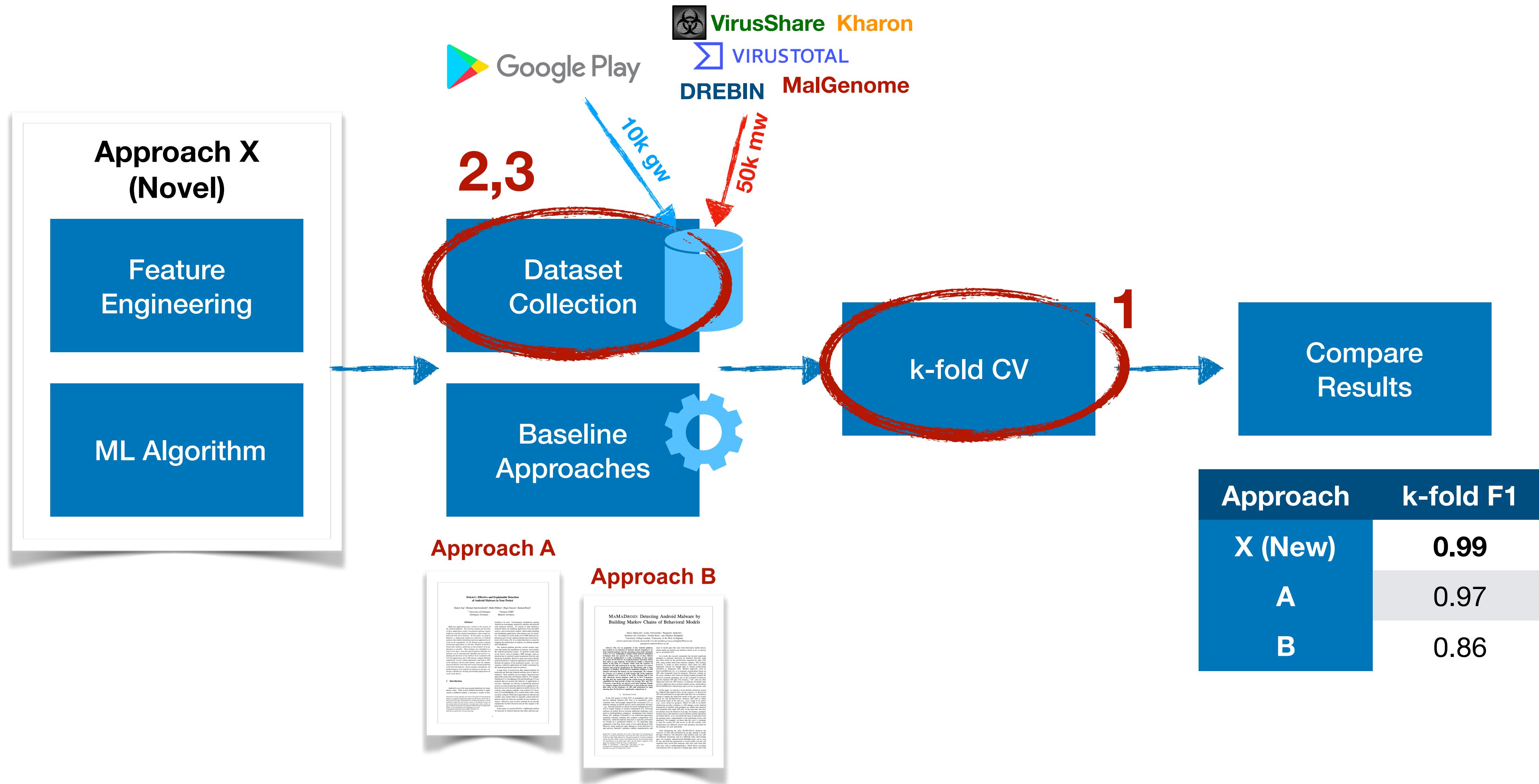
Incremental Retraining

Active Learning

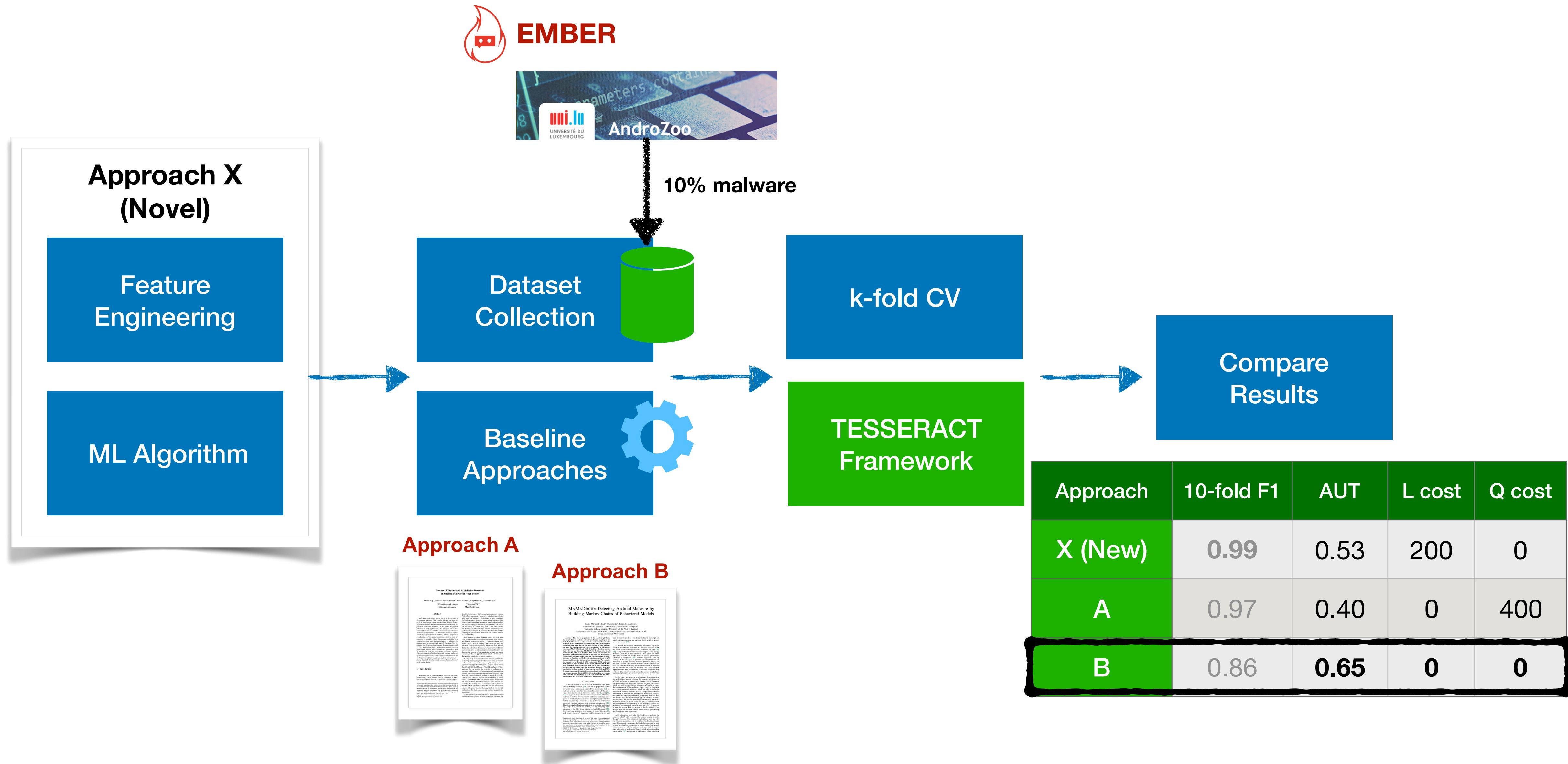
Rejection*

* R. Jordaney, K. Sharad, S. Dash, Z. Wang, D. Papini, I. Nouretdinov, and L. Cavallaro. **Transcend: Detecting Concept Drift in Malware Classification Models**. In USENIX Security Symposium, 2017

ML for Malware Detection



ML for Malware Detection



TESSERACT open-source library

<https://s2lab.kcl.ac.uk/projects/tesseract/>



UNIVERSITY OF
TORONTO



UNIVERSITY
OF CAGLIARI

IDC
HERZLIYA



MITRE



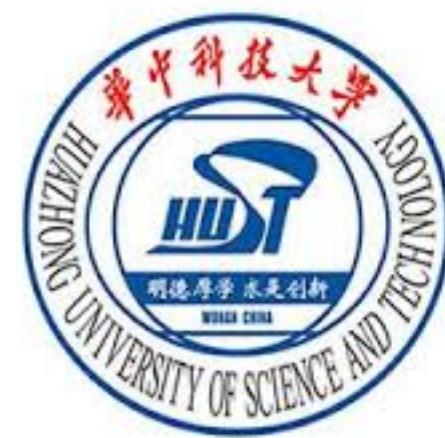
WASHINGTON STATE
UNIVERSITY



Duke
UNIVERSITY



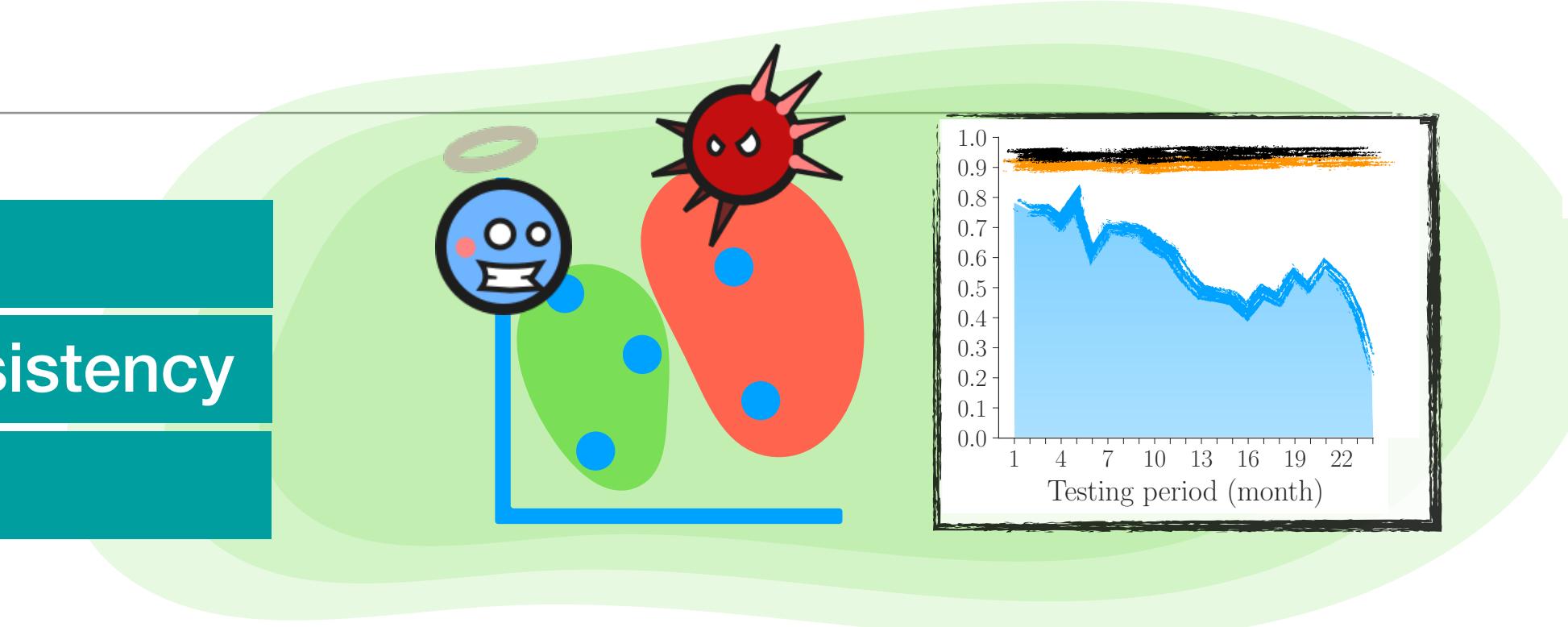
Universidad
Carlos III de Madrid



Conclusions

Experimental Constraints

- C1** Temporal training consistency
- C2** {good|mal}ware temporal consistency
- C3** Realistic testing classes ratio



- (Android) Malware Detection is still an **open problem**
- We propose TESSERACT framework
 - › **Sound time-aware evaluations** to compare different solutions
 - › **Performance-cost trade-offs** for retraining/rejection strategies

- **Open-source** code, dataset, features

<https://s2lab.kcl.ac.uk/projects/tesseract/>

