



UCL

Trustworthy AI... **for** Systems Security

Lorenzo Cavallaro <l.cavallaro@ucl.ac.uk>

@lcavallaro — <https://s2lab.cs.ucl.ac.uk>

Keynote Talk

NSS-SocialSec 2023

Aug 16, 2023

University of Kent, Canterbury, UK



Machine Learning Revolution

Image Classification

Facial Recognition

Machine Translation

Speech Recognition

Android Malware

Malicious Javascript

Windows Malware

PDF Malware

Machine Learning Revolution

Image Classification

Facial Recognition

Machine Translation

Speech Recognition

Android Malware

Malicious Javascript

Windows Malware

PDF Malware



**A Photorealistic Medieval Paparazzi of the Year 1100 AD
(Midjourney)**



Hallucination

A Photorealistic Medieval Paparazzi of the Year 1100 AD
(Midjourney)



**A Photorealistic Medieval Paparazzi of the Year 1100 AD
(Midjourney)**

Machine Learning Revolution

Image Classification

Facial Recognition

Machine Translation

Speech Recognition

Android Malware

Malicious Javascript

Windows Malware

PDF Malware

Machine Learning Revolution

Image Classification

Facial Recognition

Machine Translation

Speech Recognition

Android Malware

Malicious Javascript

Windows Malware

PDF Malware

Machine Learning Revolution

Image Classification

Facial Recognition

Machine Translation

Speech Recognition

Android Malware

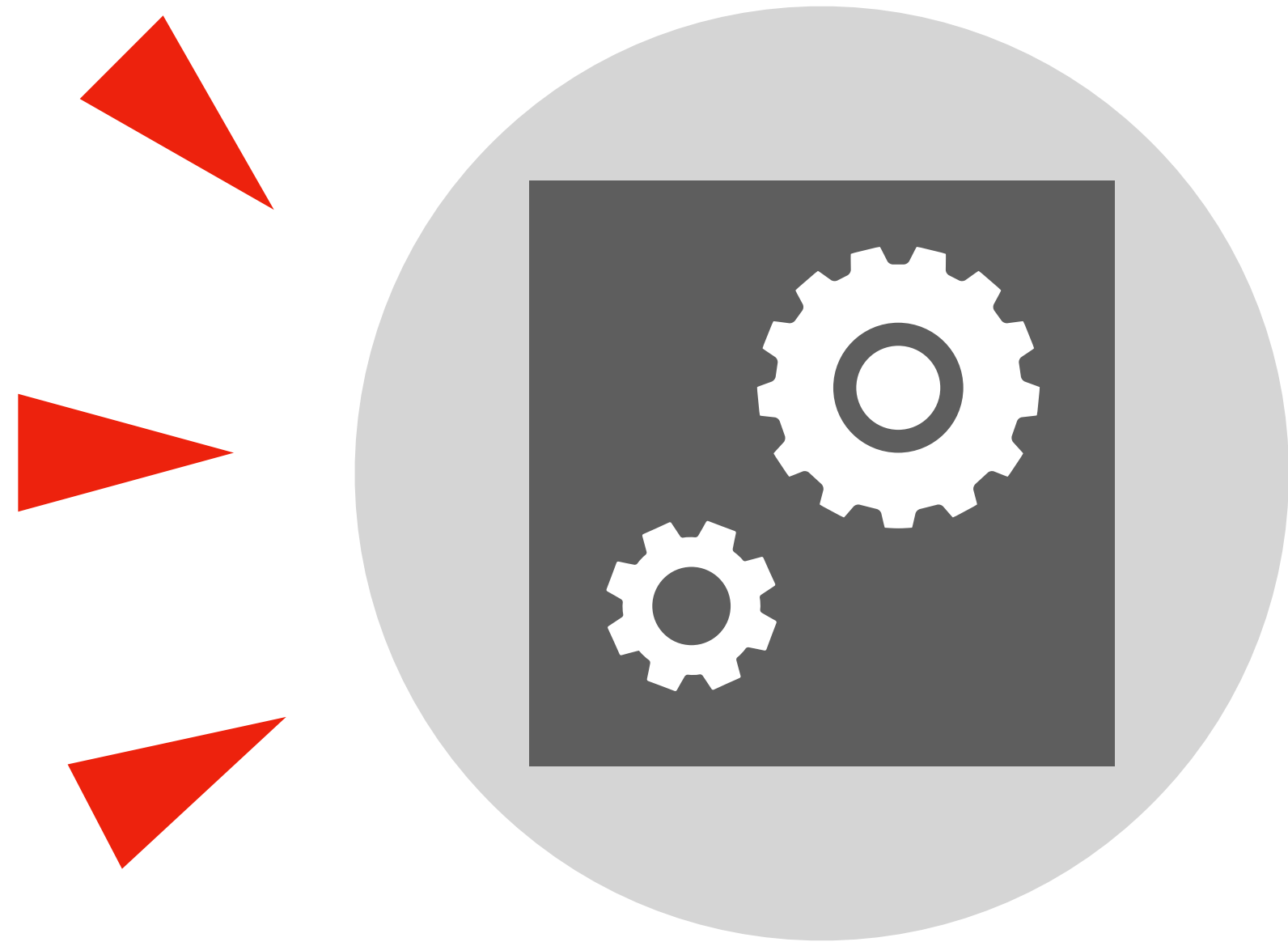
Malicious Javascript

Windows Malware

PDF Malware

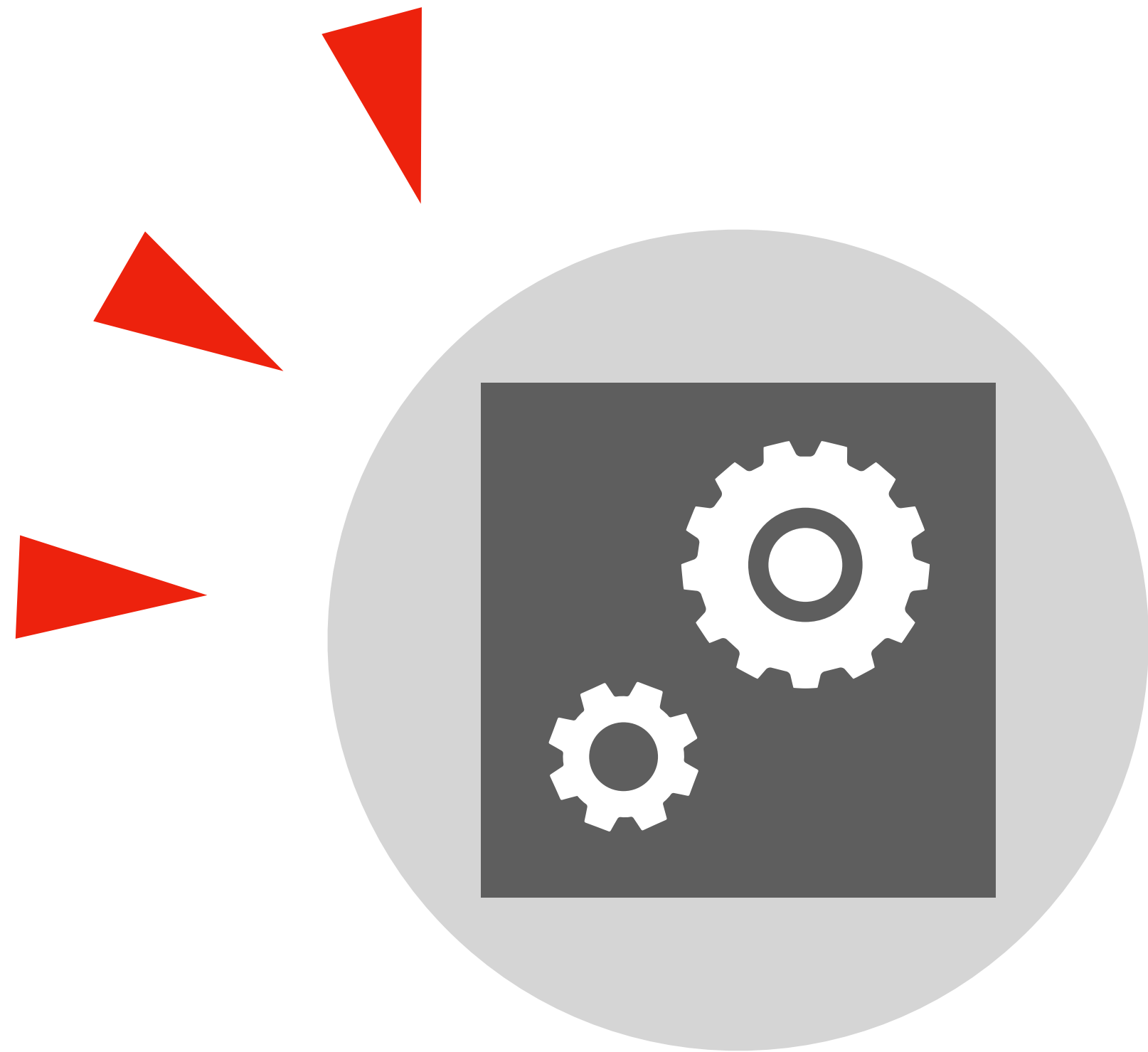
Maybe the conditions aren't ready yet?

Security *is* Adversarial



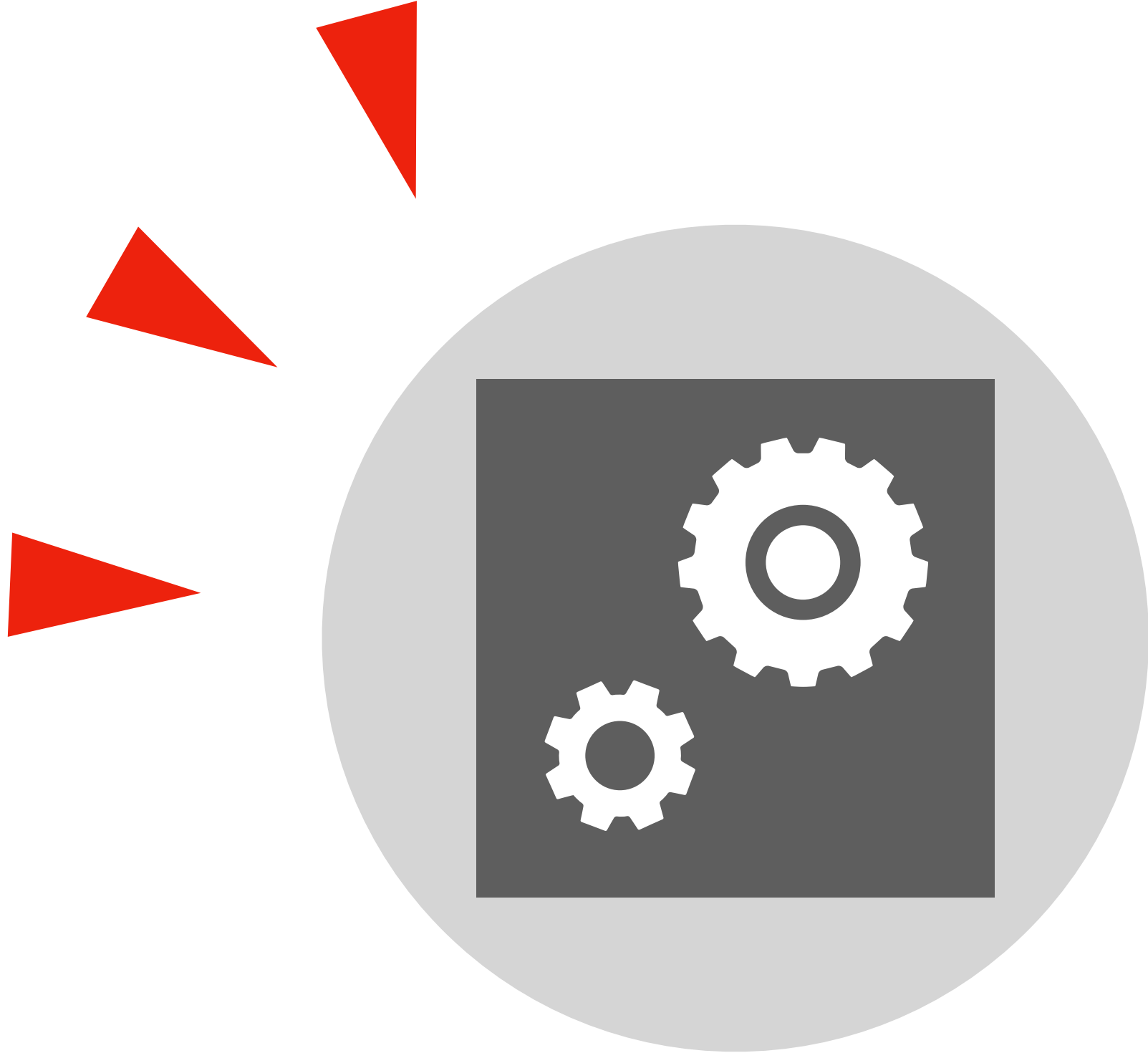
New detection systems trigger
an immediate response...

Security *is* Adversarial

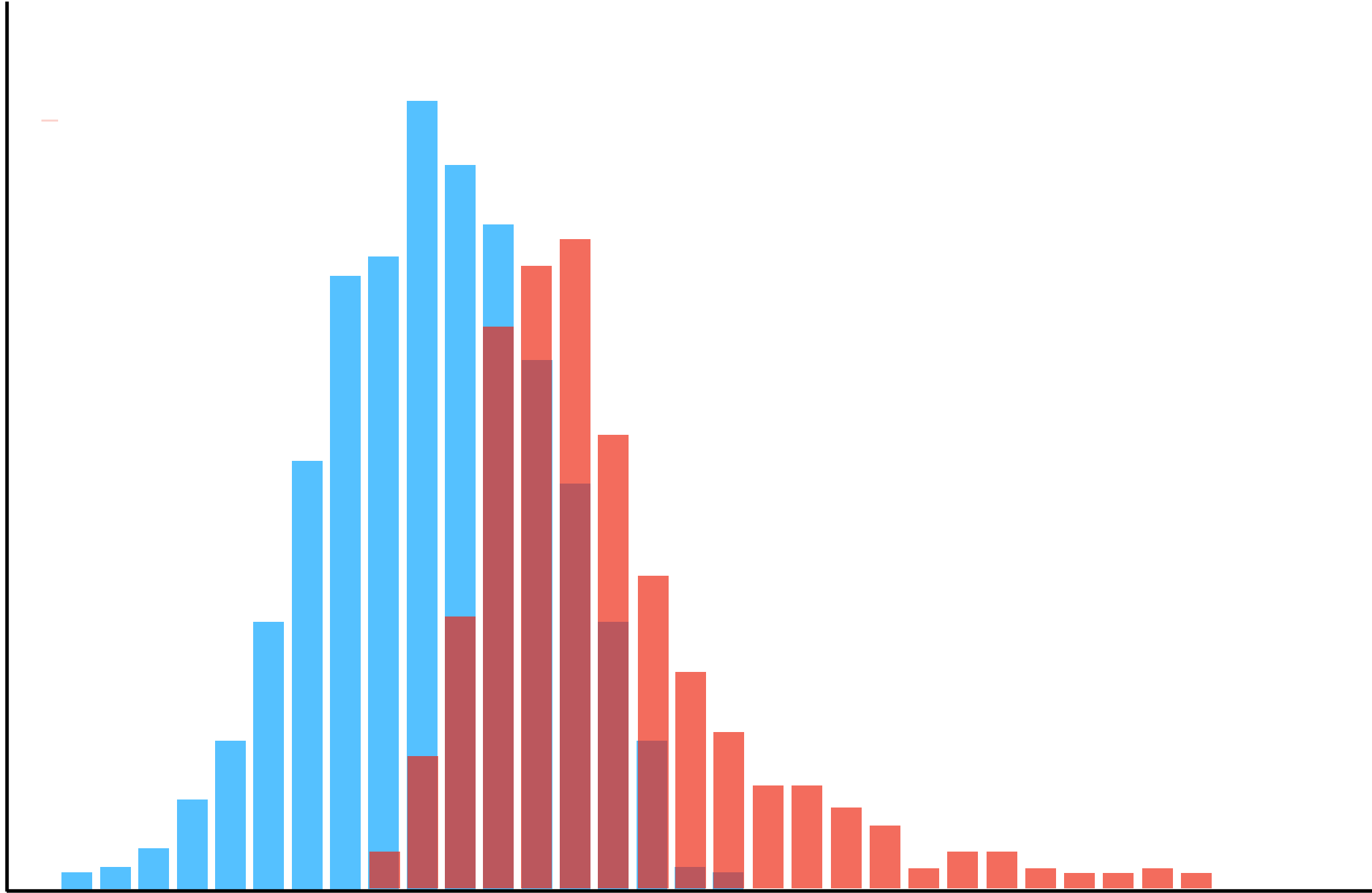


New detection systems trigger
an immediate response...

Security *is* Adversarial

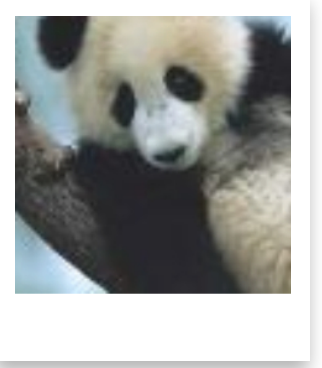
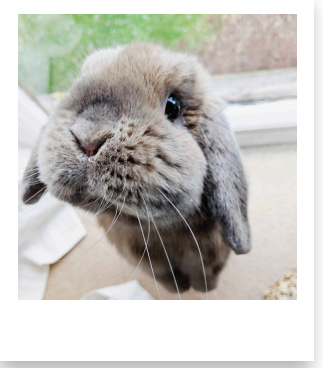
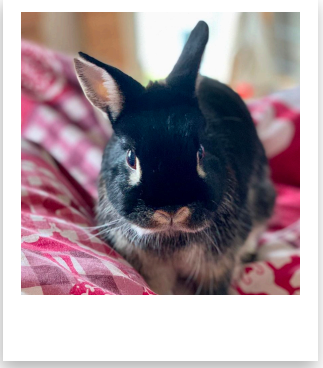
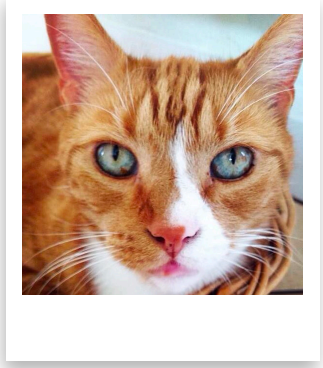
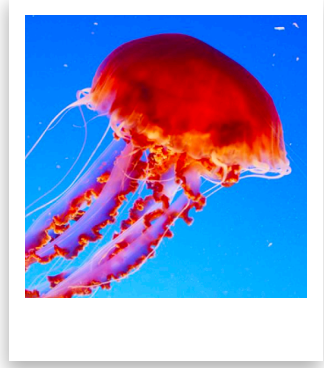
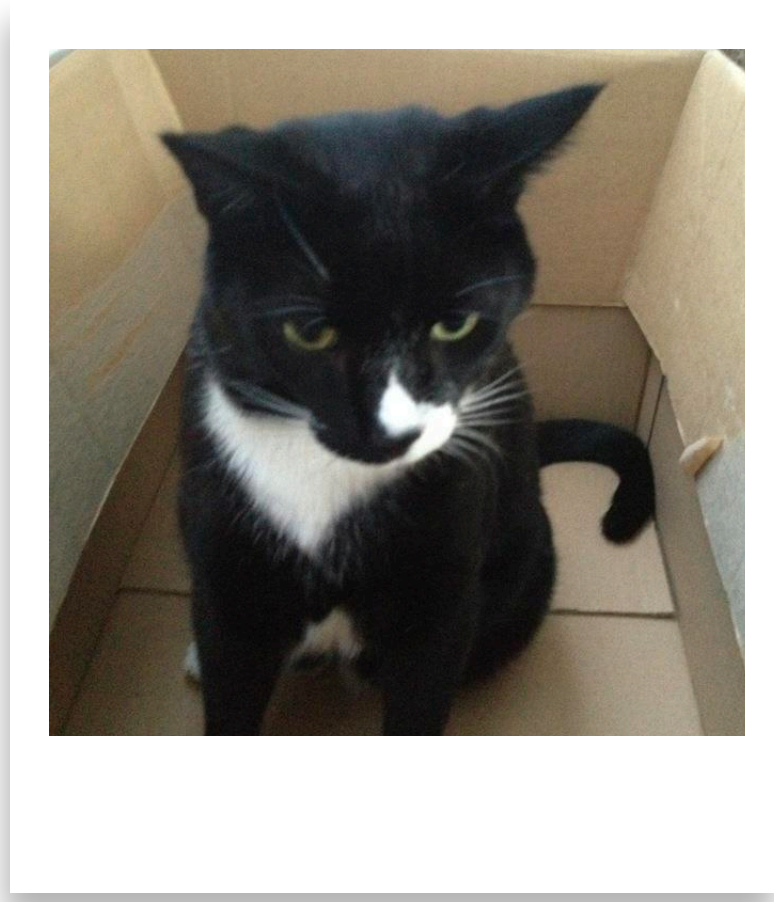


New detection systems trigger an immediate response...

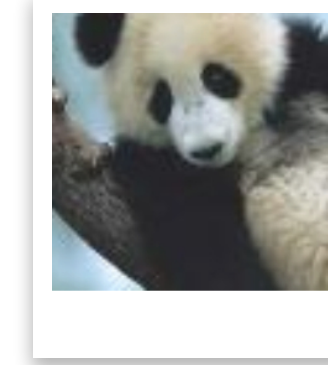
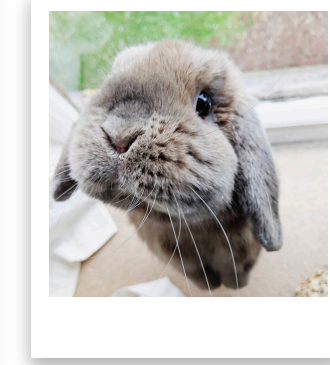
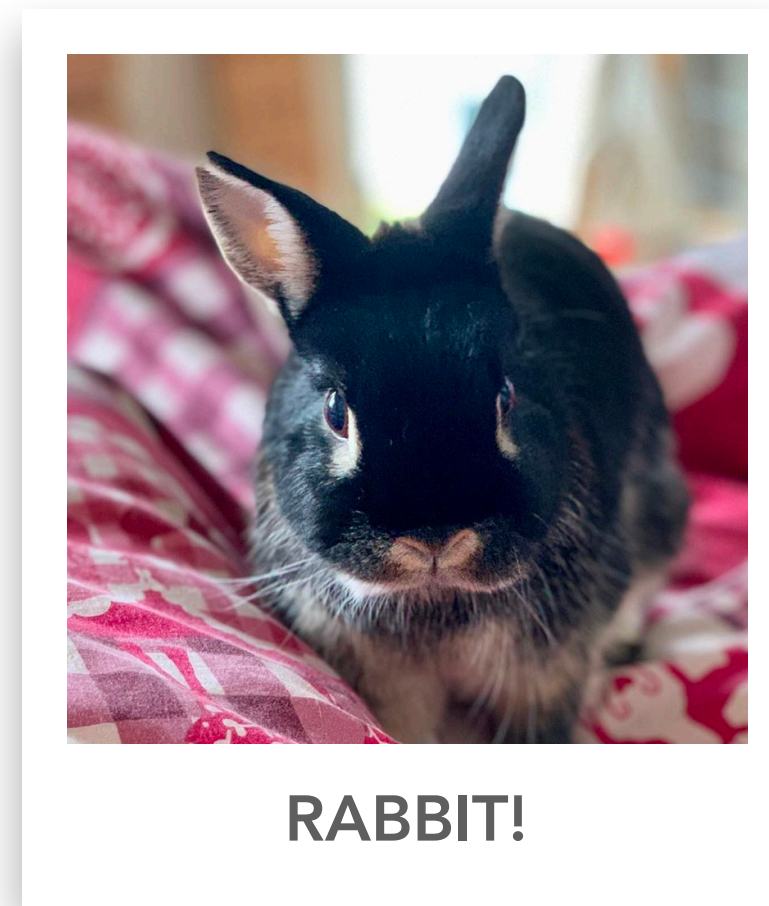
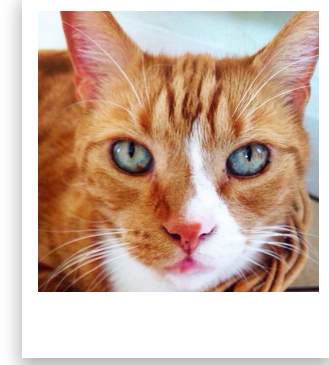
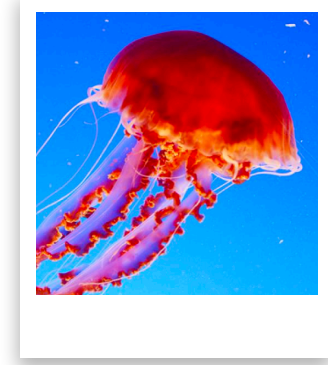
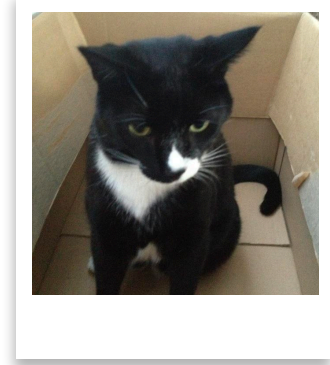


...which causes dataset shifts, often violating the i.i.d. assumption

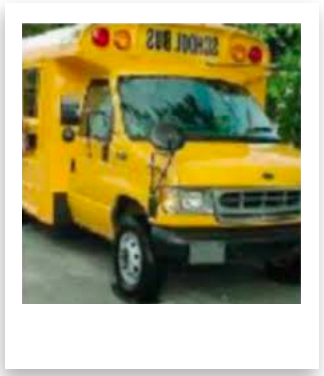
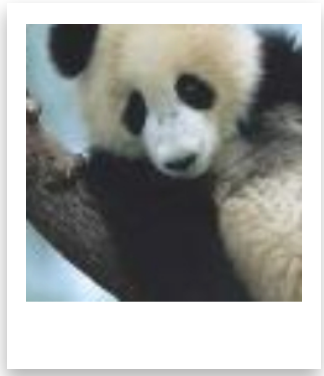
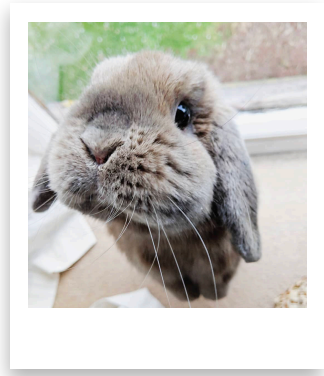
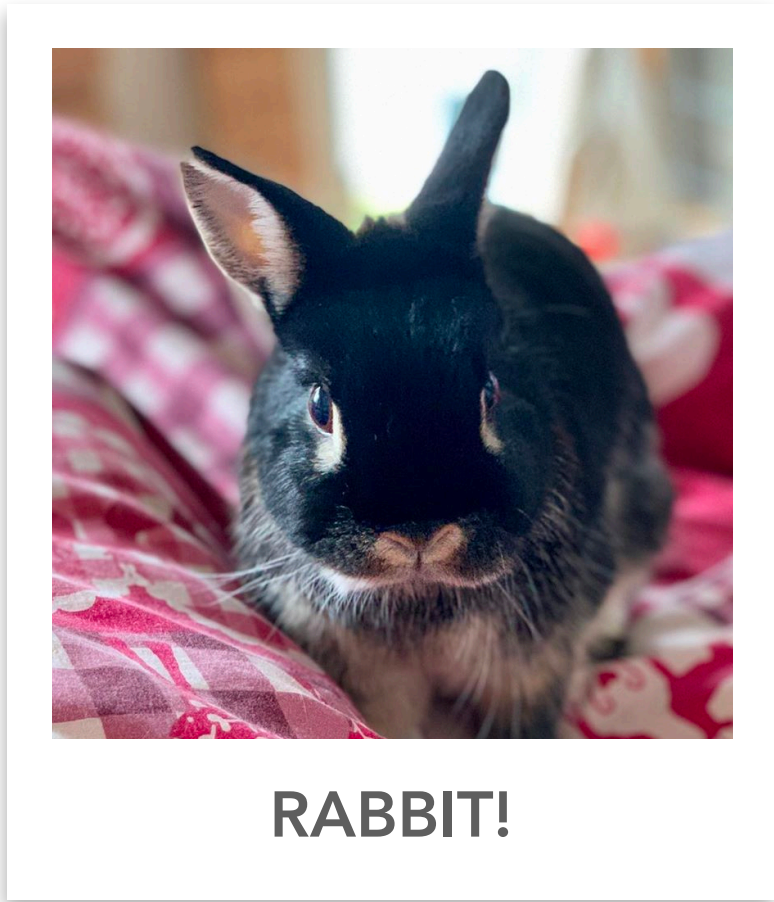
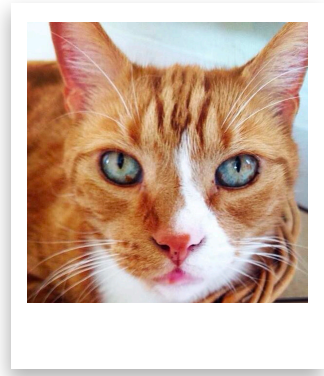
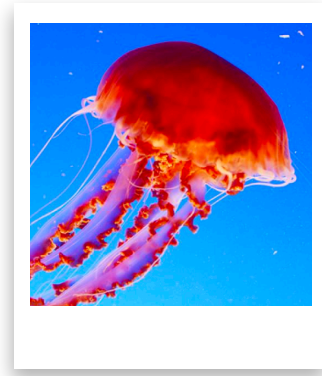
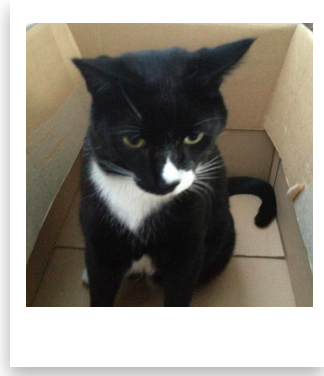
Security *is* Adversarial



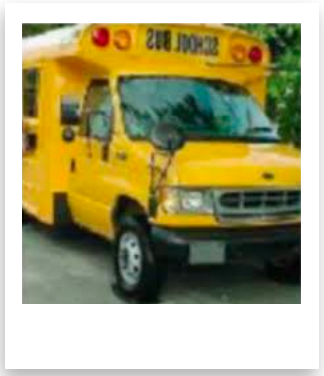
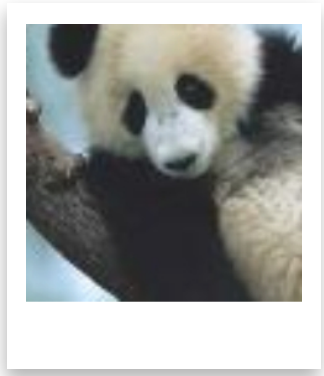
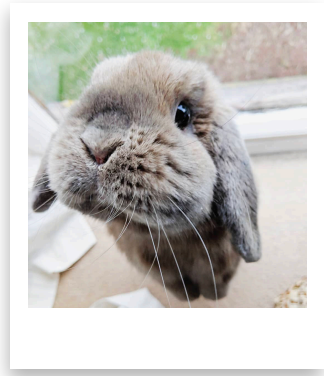
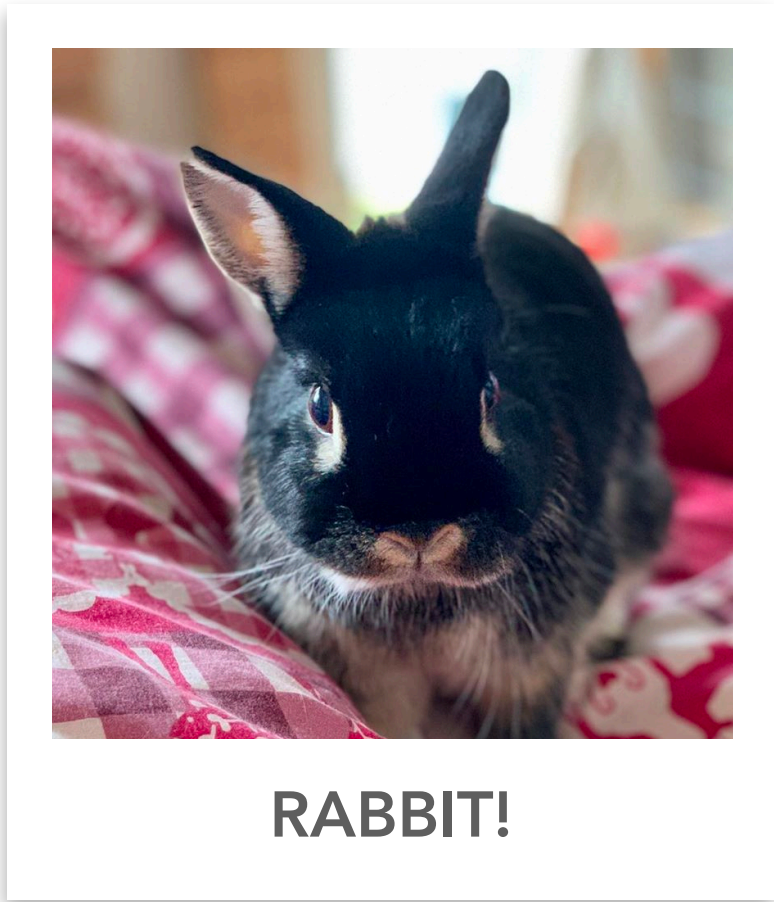
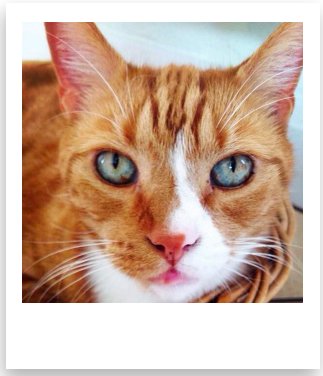
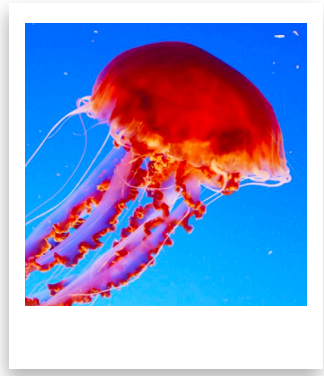
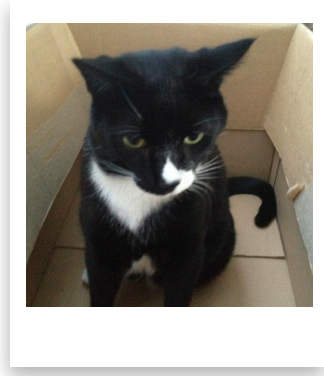
Security *is* Adversarial



Security *is* Adversarial

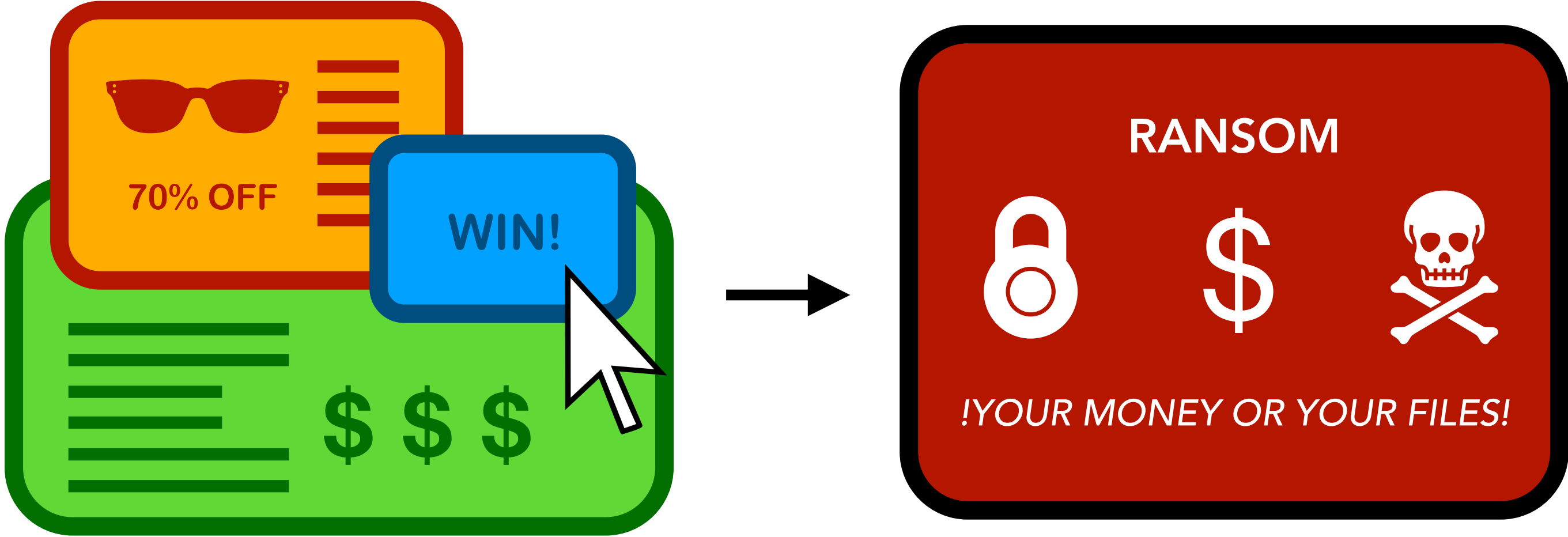
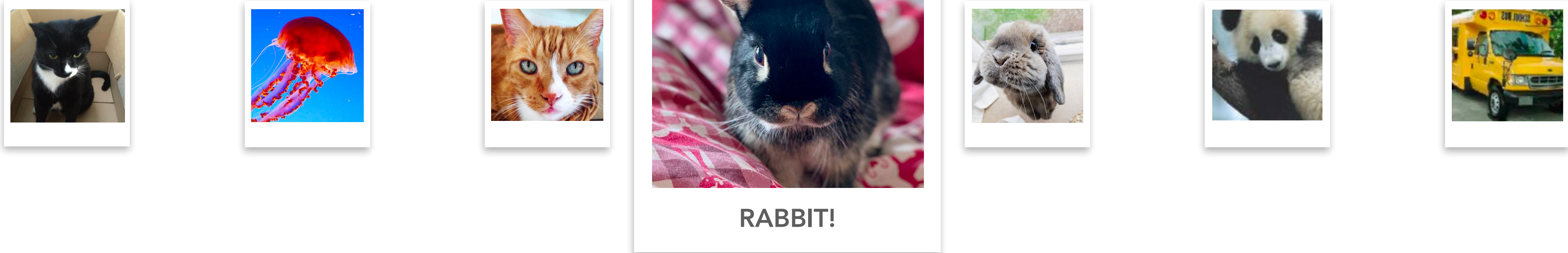


Security is Adversarial



Training on ad fraud...

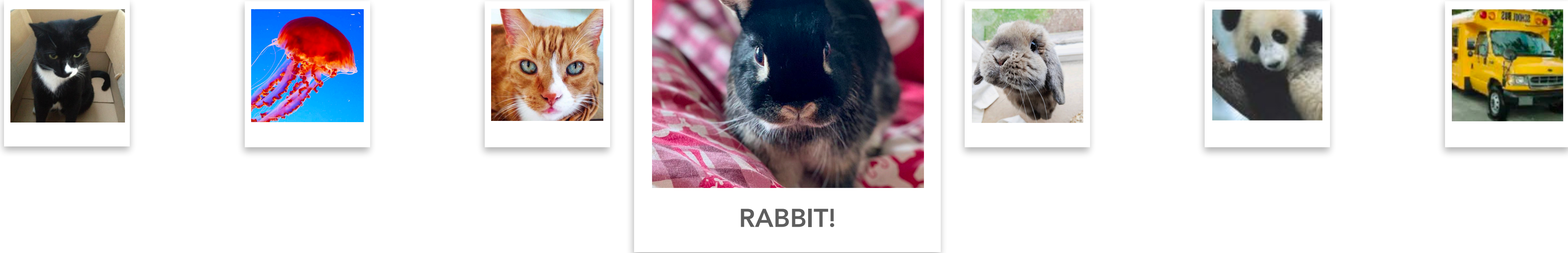
Security is Adversarial



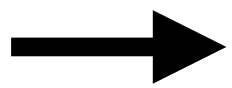
Training on ad fraud...

...attacks evolve at test time...

Security is Adversarial



Training on ad fraud...



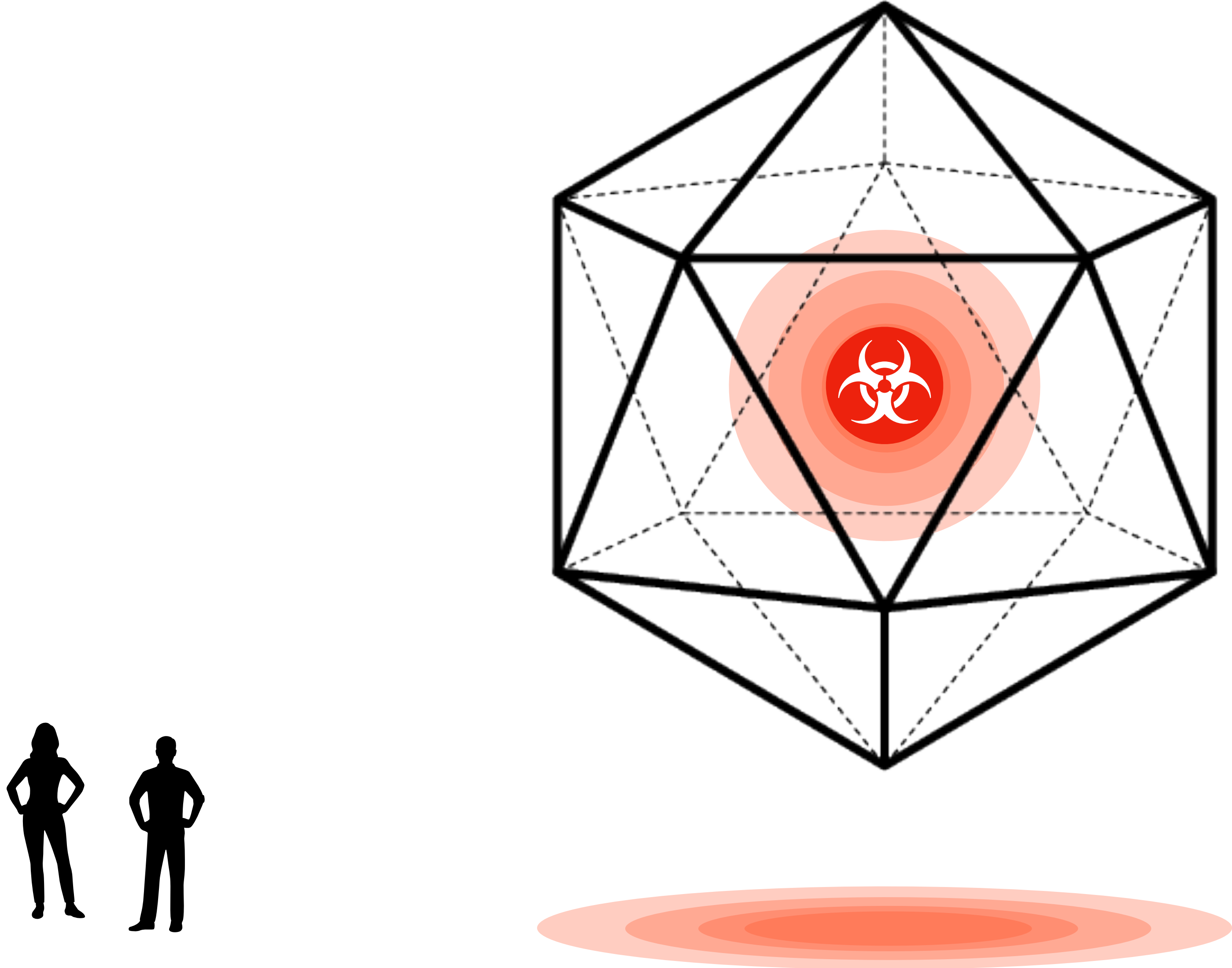
...attacks evolve at test time...



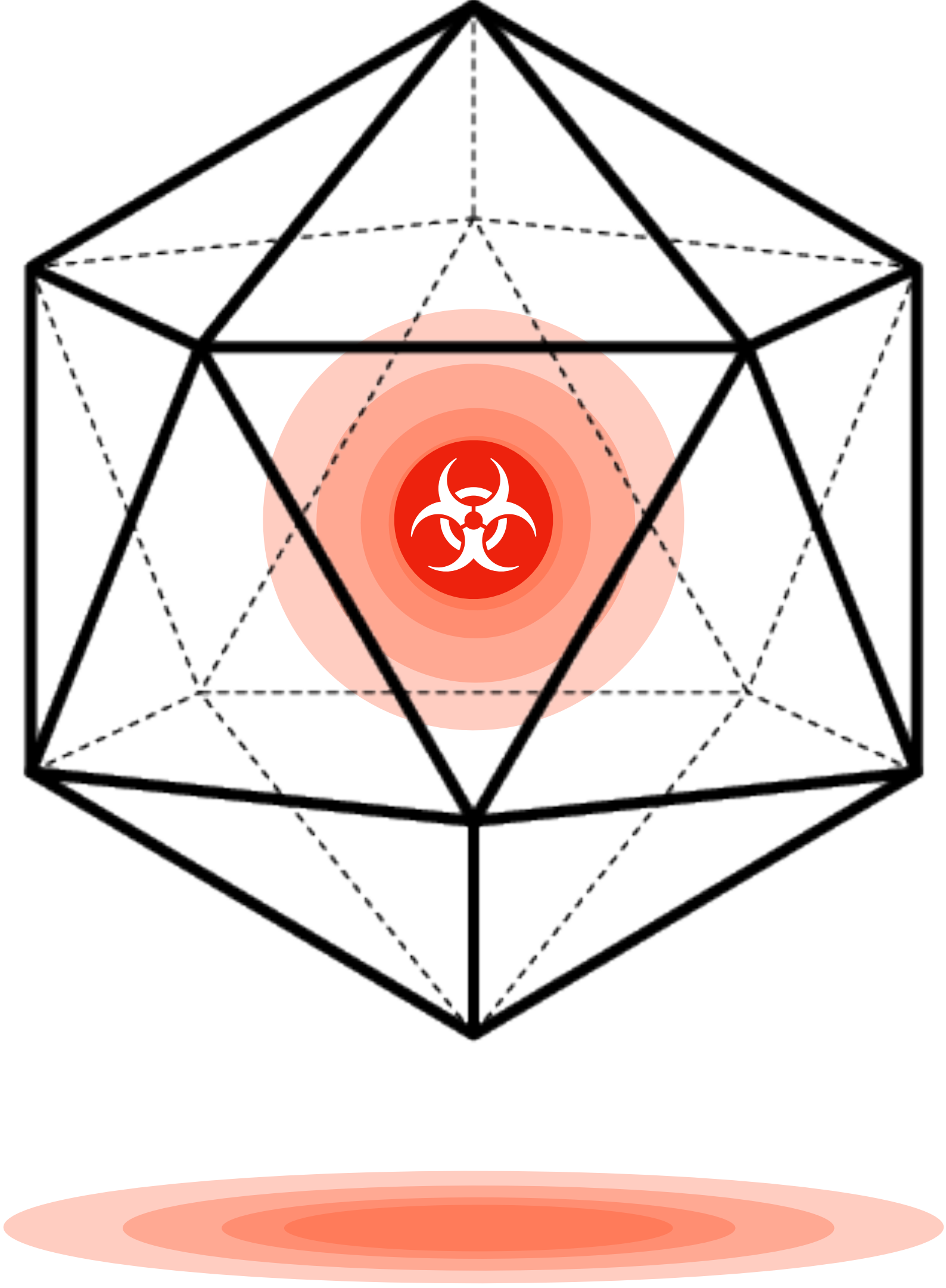
...prepare for the unknown



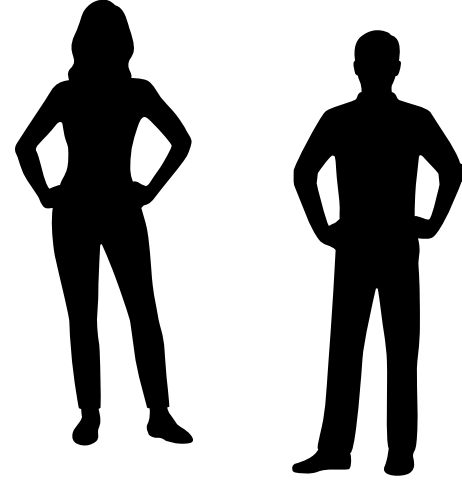
Security *is* Adversarial



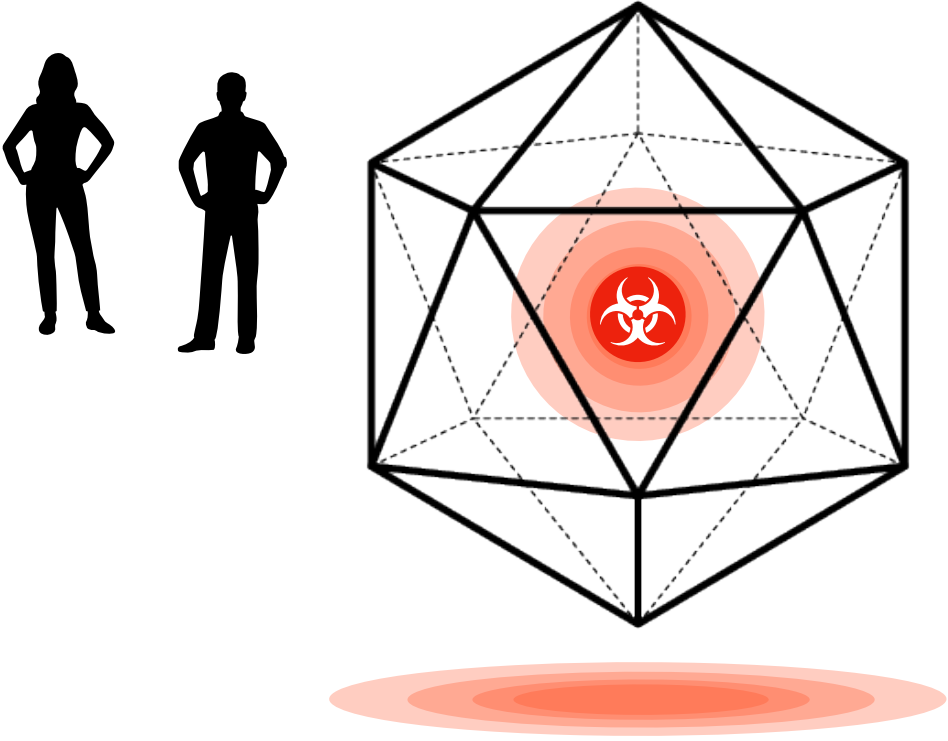
Security *is* Adversarial



Oh yeah, that's malware alright

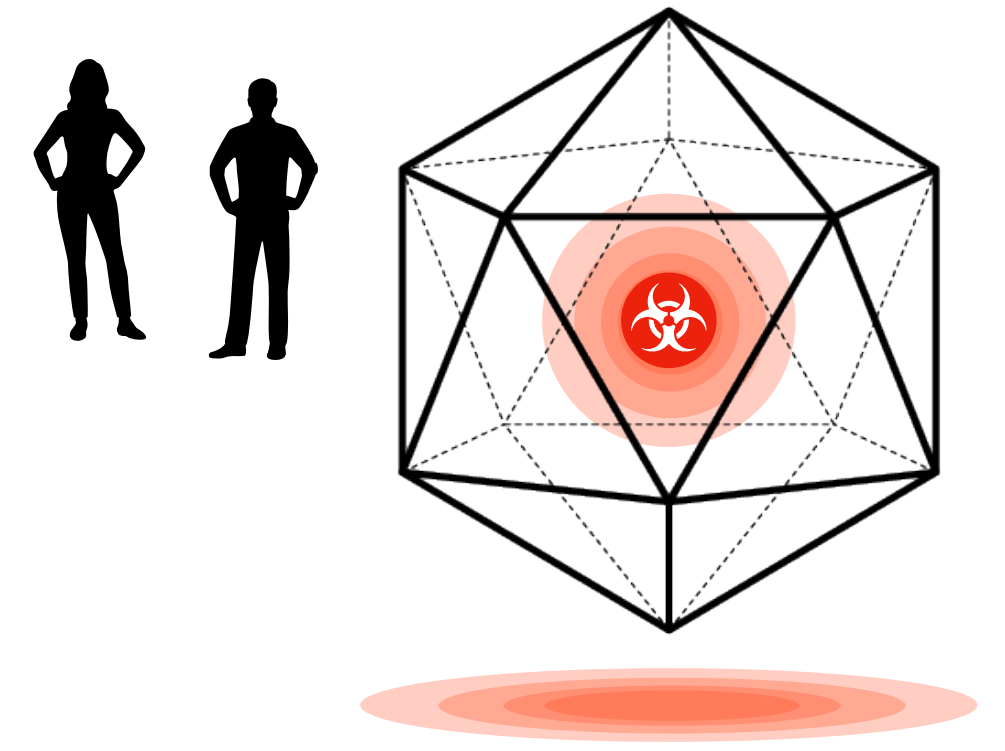


We need...



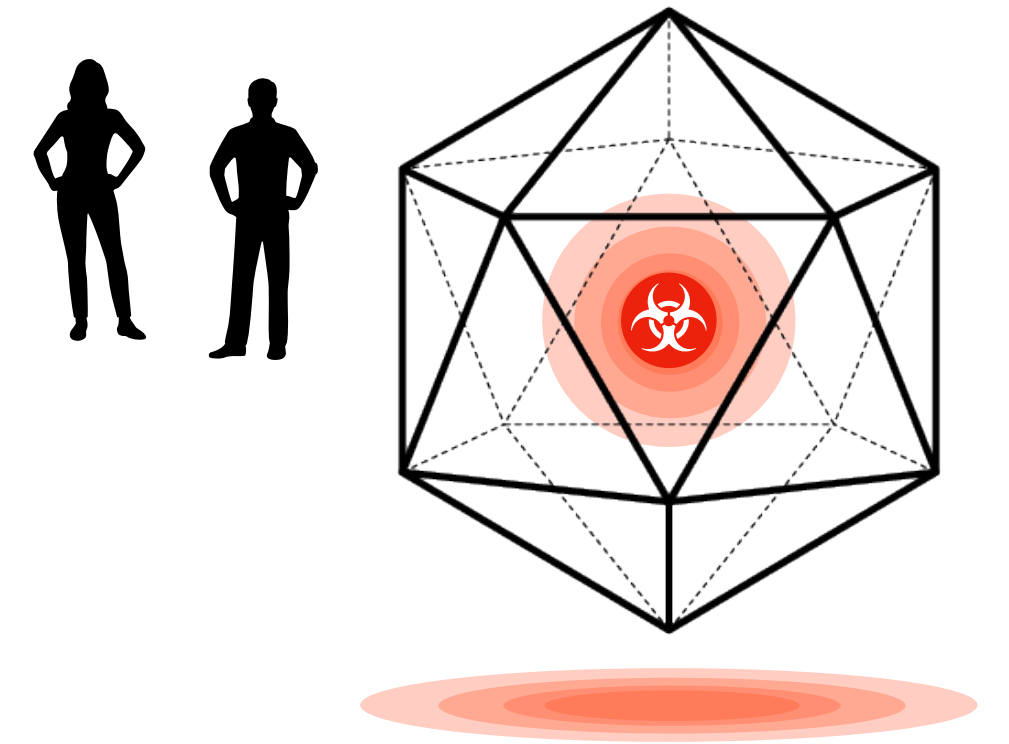
We need...

To understand and improve the effectiveness of machine learning methods for systems security in the presence of adversaries



We need...

To understand and improve the effectiveness of machine learning methods for systems security in the presence of adversaries

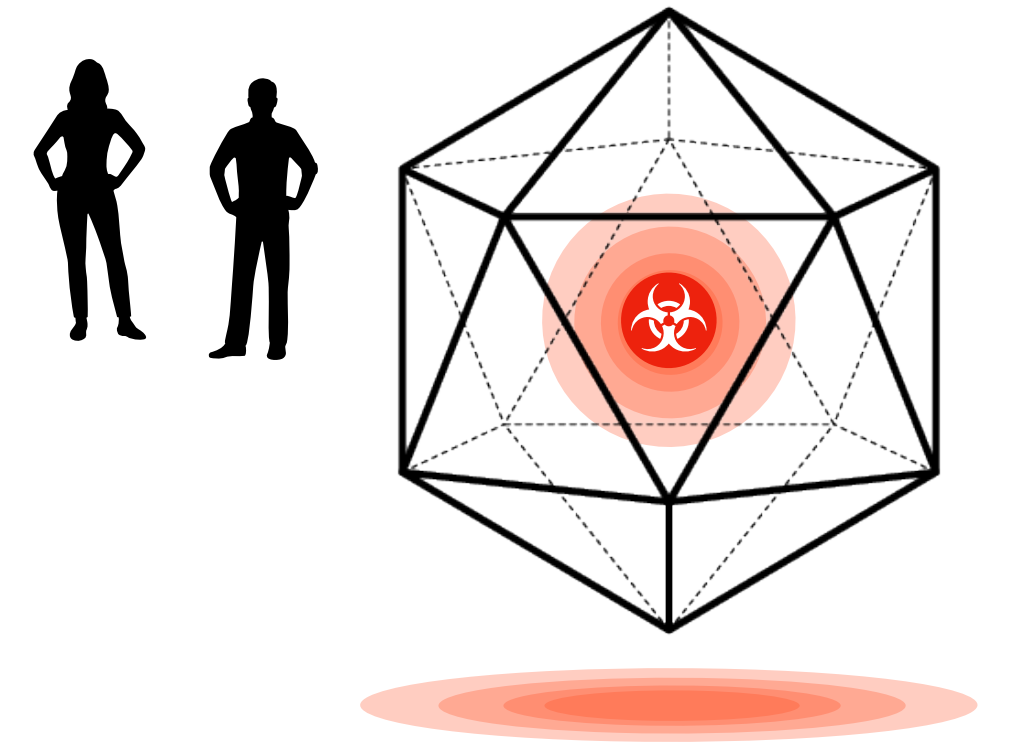


Representation of problem space objects (e.g., programs) results in a **semantic gap**

- It makes designing attacks and defenses more challenging
- It leaves room for adversarial manipulation
- It challenges the identification of causal vs non-causal (spurious) features

We need...

To understand and improve the effectiveness of machine learning methods for systems security in the presence of adversaries



Representation of problem space objects (e.g., programs) results in a **semantic gap**

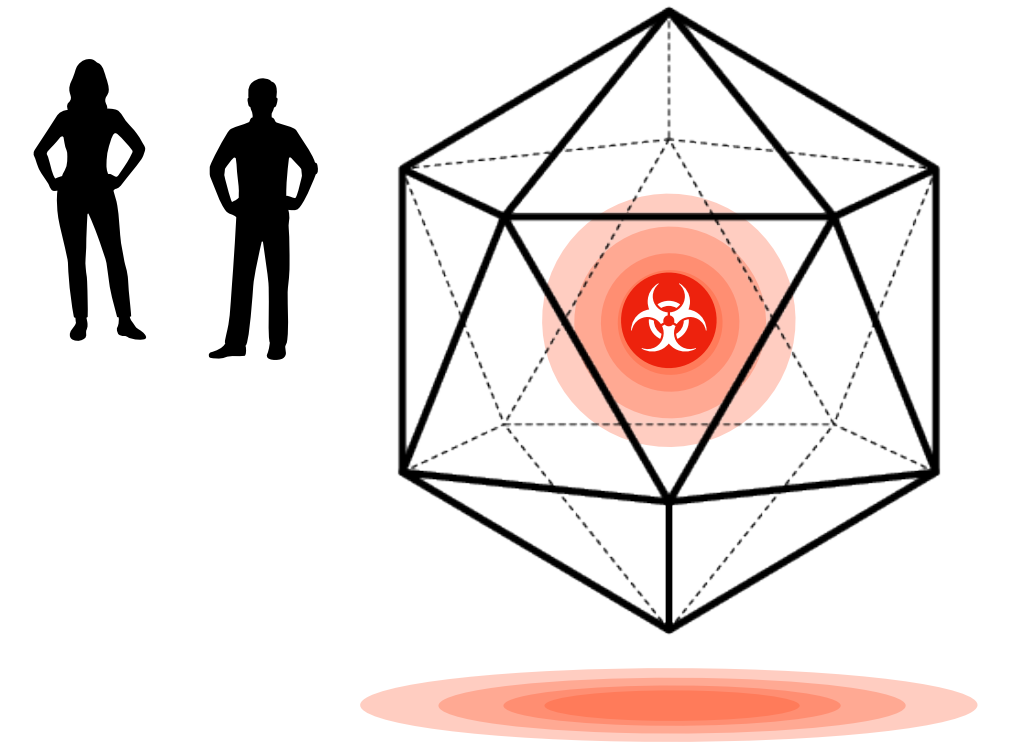
- It makes designing attacks and defenses more challenging
- It leaves room for adversarial manipulation
- It challenges the identification of causal vs non-causal (spurious) features

Effectiveness of ML for systems security is intertwined with the **underlying abstractions**, e.g., program analyses, to represent objects

- This affects robustness to adversarial drift, explainability, costs, and performance

We need...

To understand and improve the effectiveness of machine learning methods for systems security in the presence of adversaries



Representation of problem space objects (e.g., programs) results in a **semantic gap**

- It makes designing attacks and defenses more challenging
- It leaves room for adversarial manipulation
- It challenges the identification of causal vs non-causal (spurious) features

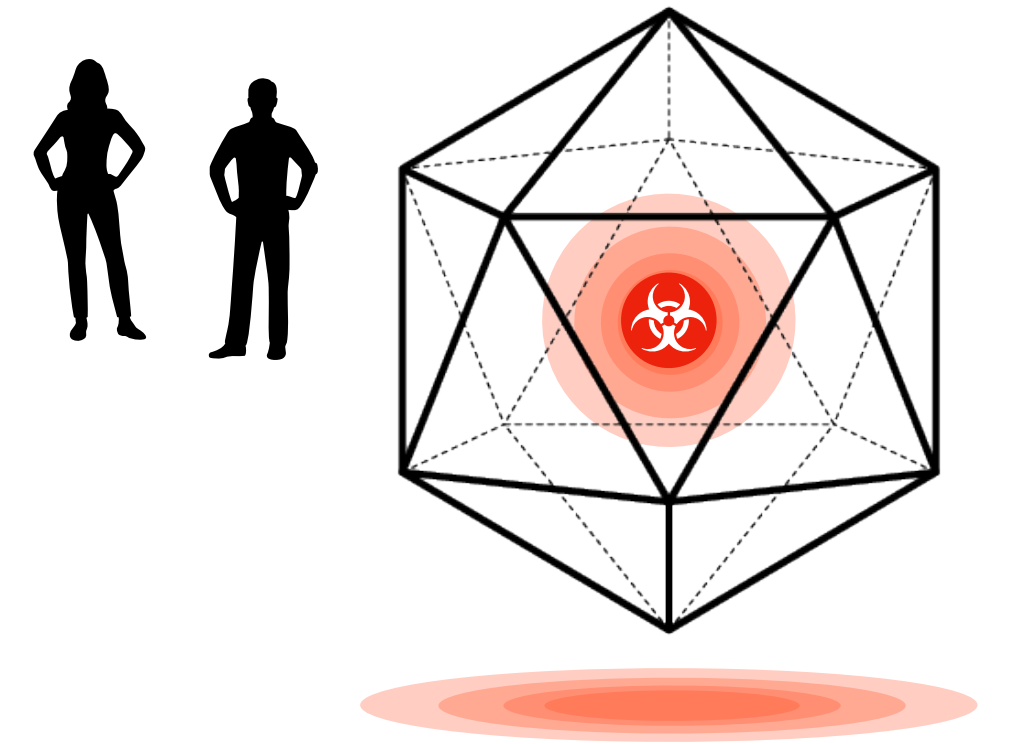
Effectiveness of ML for systems security is intertwined with the **underlying abstractions**, e.g., program analyses, to represent objects

- This affects robustness to adversarial drift, explainability, costs, and performance

Is Trustworthy AI for systems security possible?

We need...

To understand and improve the effectiveness of machine learning methods for systems security in the presence of adversaries



Representation of problem space objects (e.g., programs) results in a **semantic gap**

- It makes designing attacks and defenses more challenging
- It leaves room for adversarial manipulation
- It challenges the identification of causal vs non-causal (spurious) features

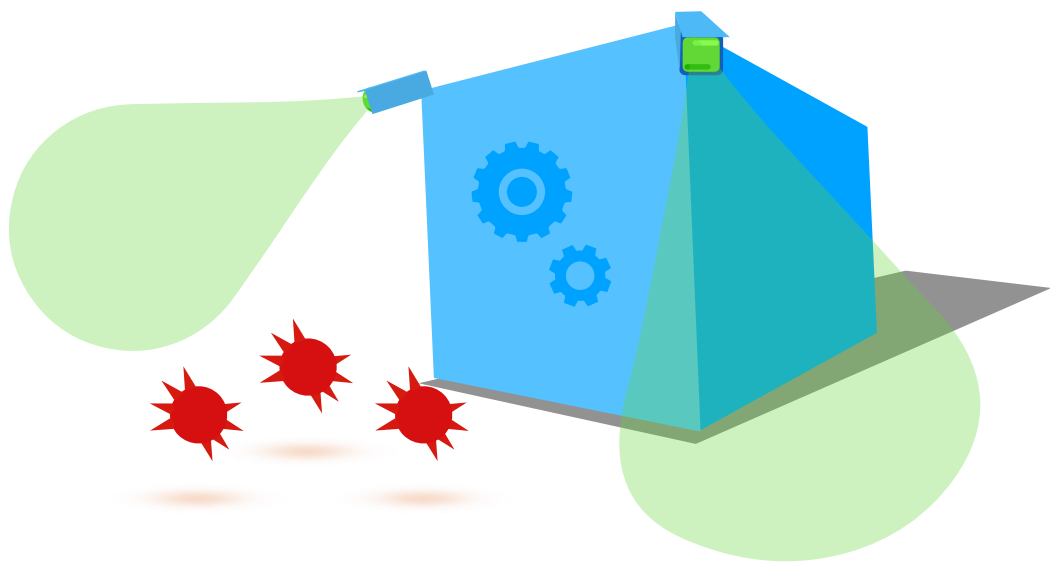
Effectiveness of ML for systems security is intertwined with the **underlying abstractions**, e.g., program analyses, to represent objects

- This affects robustness to adversarial drift, explainability, costs, and performance

Is Trustworthy AI for systems security possible? ←

*We already know the answer!
But wait, there's a plan...*

Outline

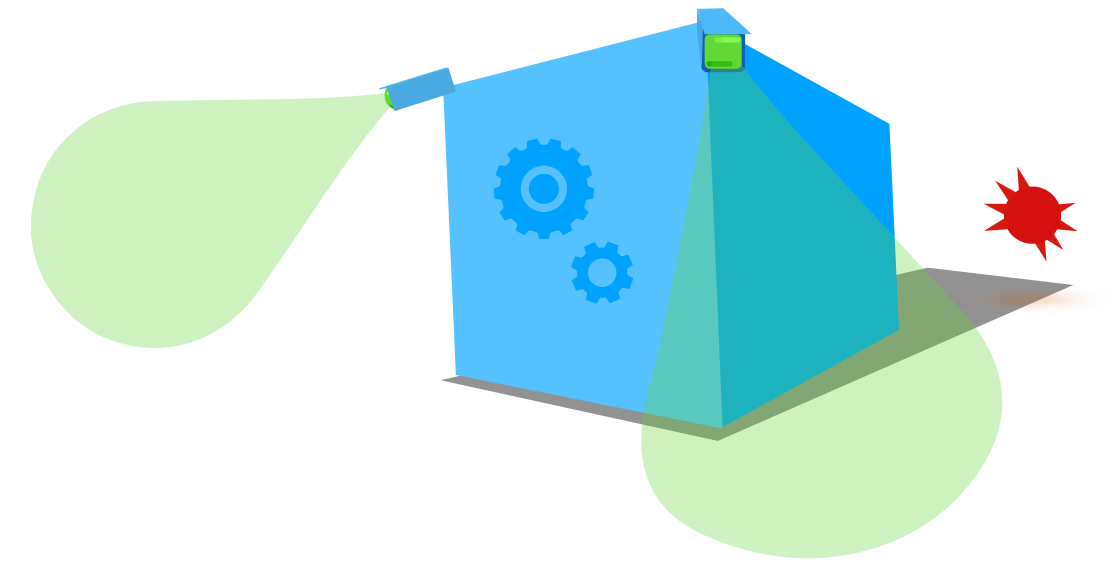


Outline

Focus

Adversarial ML evasion attacks against malware classifiers

- Classic formulation of evasion attacks is ill-suited for reasoning about realizable evasive malware
- By reformulating, we can propose stronger attacks and easily compare against alternatives

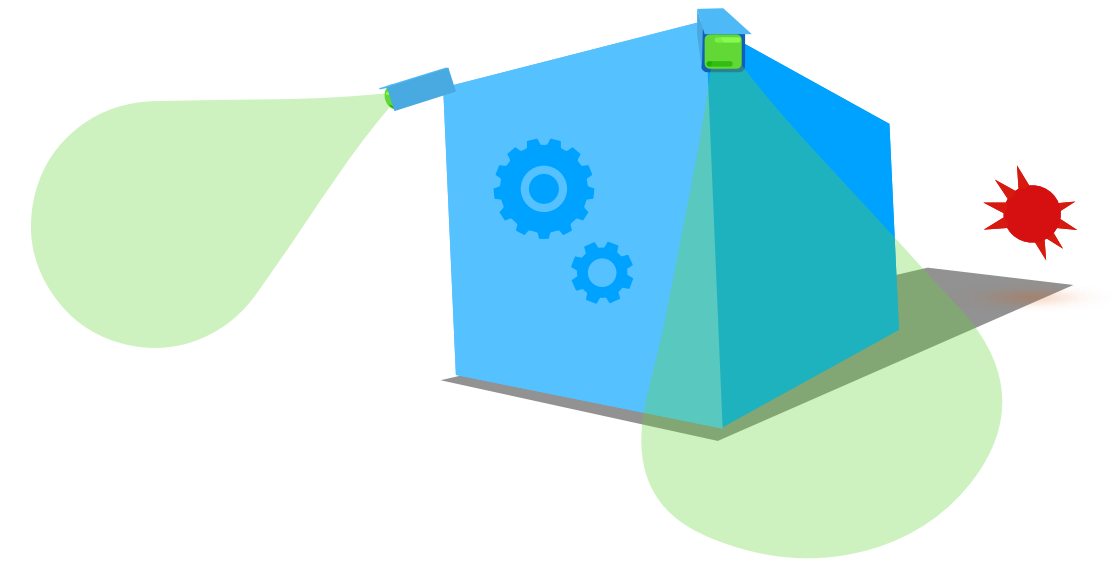


Outline

Focus

Adversarial ML evasion attacks against malware classifiers

- Classic formulation of evasion attacks is ill-suited for reasoning about realizable evasive malware
- By reformulating, we can propose stronger attacks and easily compare against alternatives

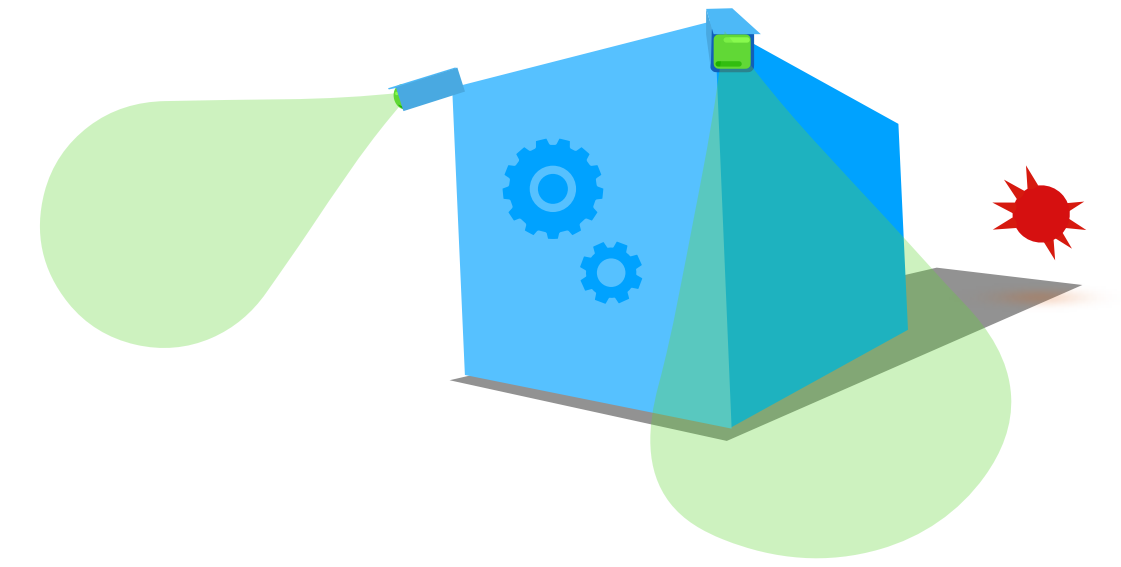


Outline

Focus

Adversarial ML evasion attacks against malware classifiers

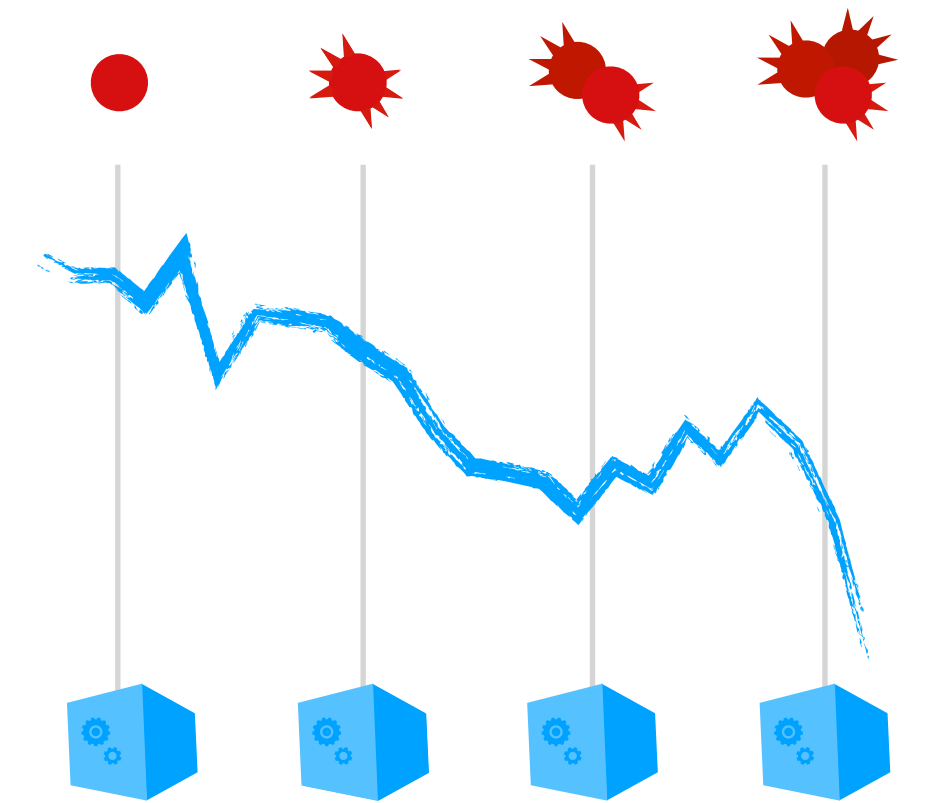
- Classic formulation of evasion attacks is ill-suited for reasoning about realizable evasive malware
- By reformulating, we can propose stronger attacks and easily compare against alternatives



Bigger Picture

Drifting scenarios caused by threats evolving over time

- How concept drift affects machine learning-based detectors in security settings
- The need for time-aware evaluations and metrics

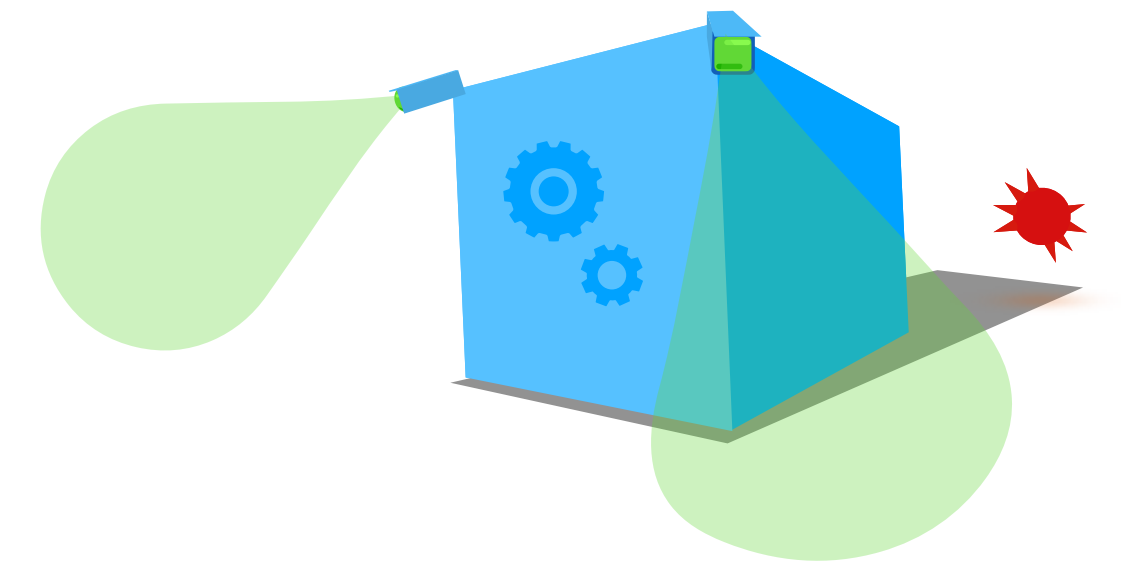


Outline

Focus

Adversarial ML evasion attacks against malware classifiers

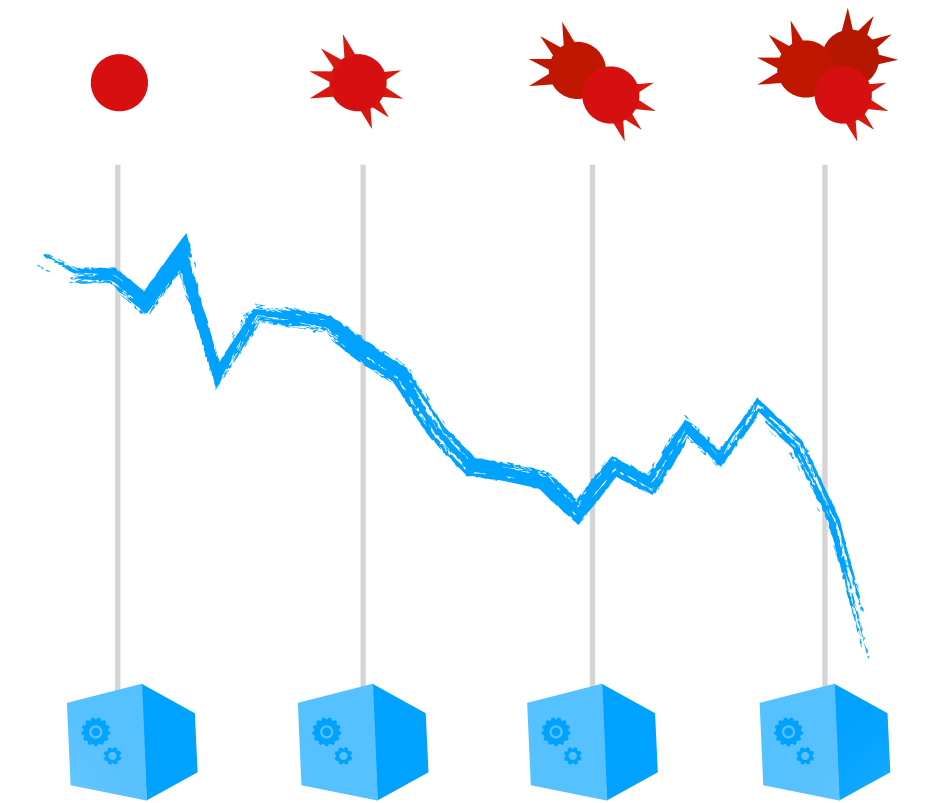
- Classic formulation of evasion attacks is ill-suited for reasoning about realizable evasive malware
- By reformulating, we can propose stronger attacks and easily compare against alternatives



Bigger Picture

Drifting scenarios caused by threats evolving over time

- How concept drift affects machine learning-based detectors in security settings
- The need for time-aware evaluations and metrics

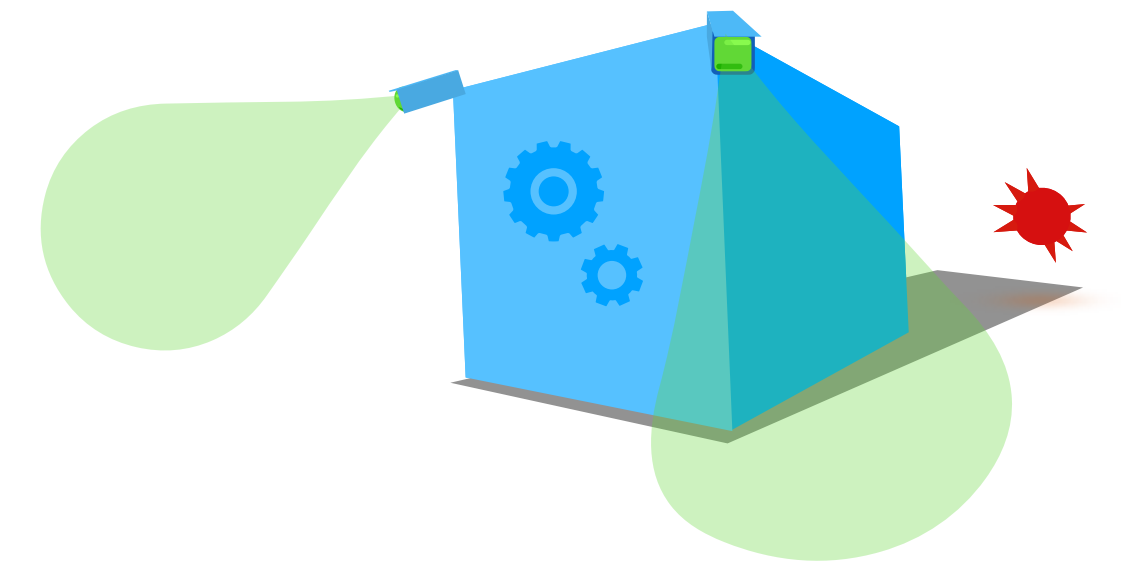


Outline

Focus

Adversarial ML evasion attacks against malware classifiers

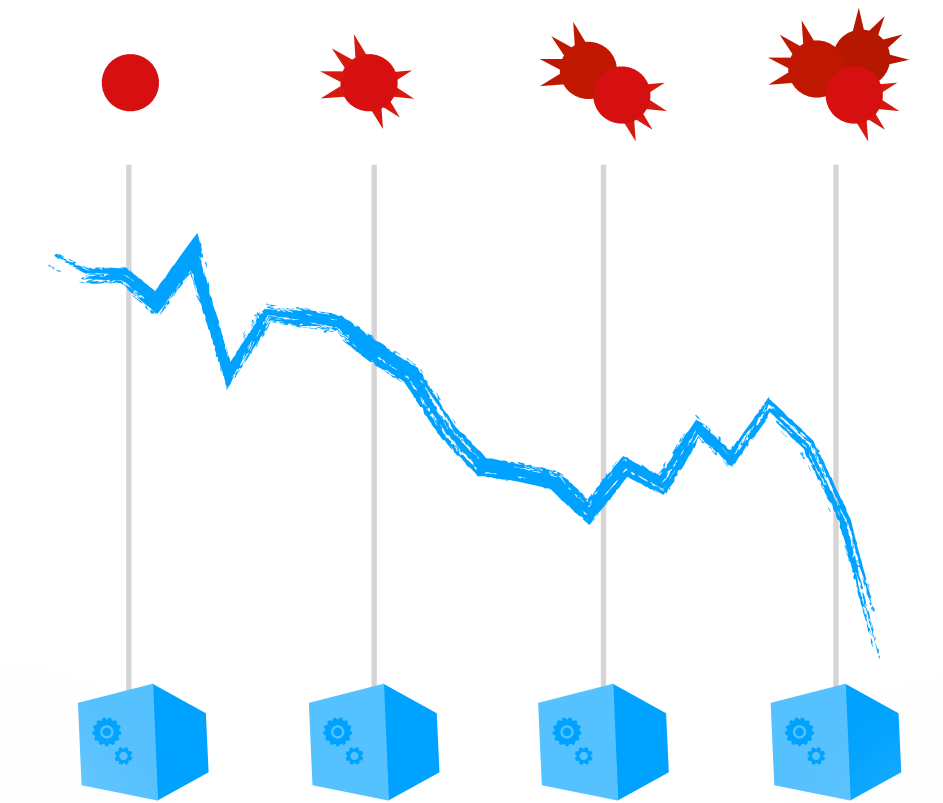
- Classic formulation of evasion attacks is ill-suited for reasoning about realizable evasive malware
- By reformulating, we can propose stronger attacks and easily compare against alternatives



Bigger Picture

Drifting scenarios caused by threats evolving over time

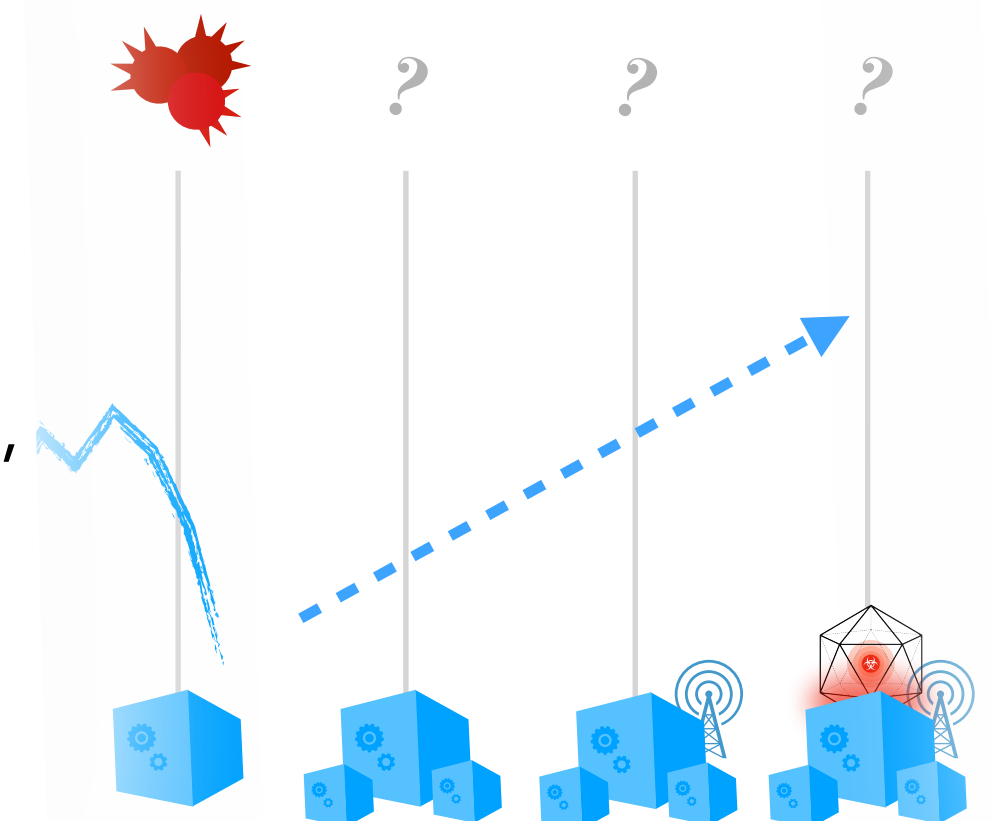
- How concept drift affects machine learning-based detectors in security settings
- The need for time-aware evaluations and metrics



Looking Ahead

Quo vadis?

- Discussion of the future of trustworthy ML for system security
- Robust feature development, universal adversarial perturbations, realizable backdoors, drift forecasting, and the role of abstractions towards the Platonic ideal of interesting behaviors

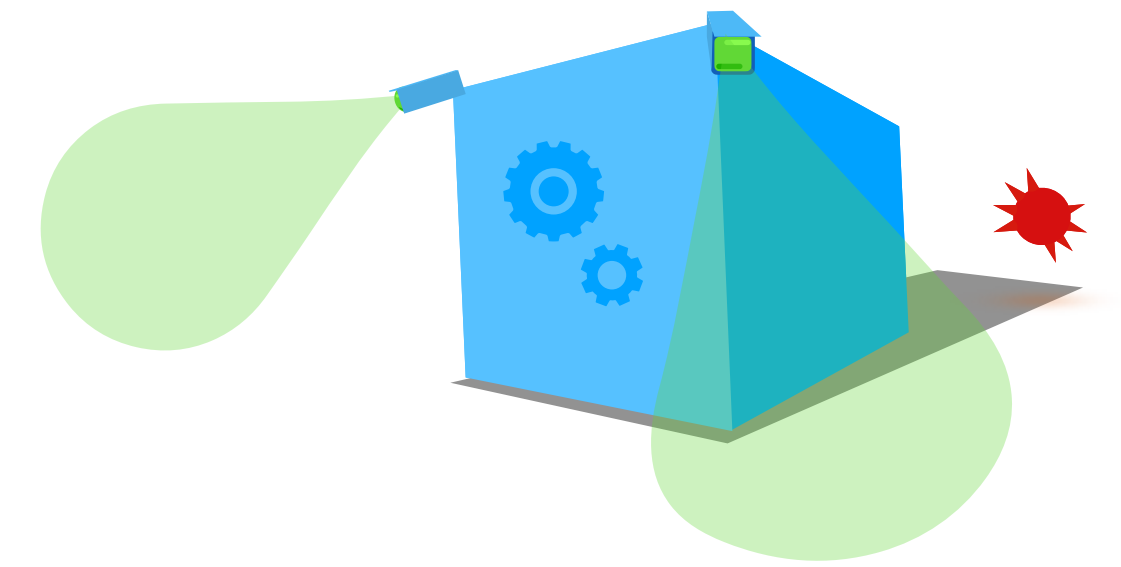


Outline

Focus

Adversarial ML evasion attacks against malware classifiers

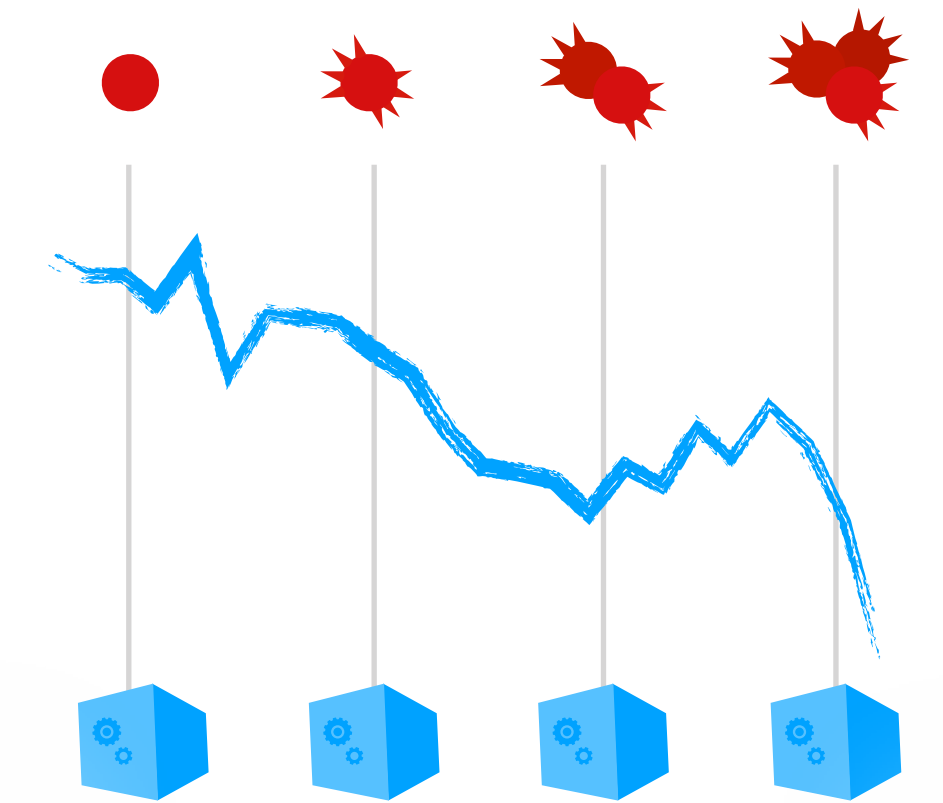
- Classic formulation of evasion attacks is ill-suited for reasoning about realizable evasive malware
- By reformulating, we can propose stronger attacks and easily compare against alternatives



Bigger Picture

Drifting scenarios caused by threats evolving over time

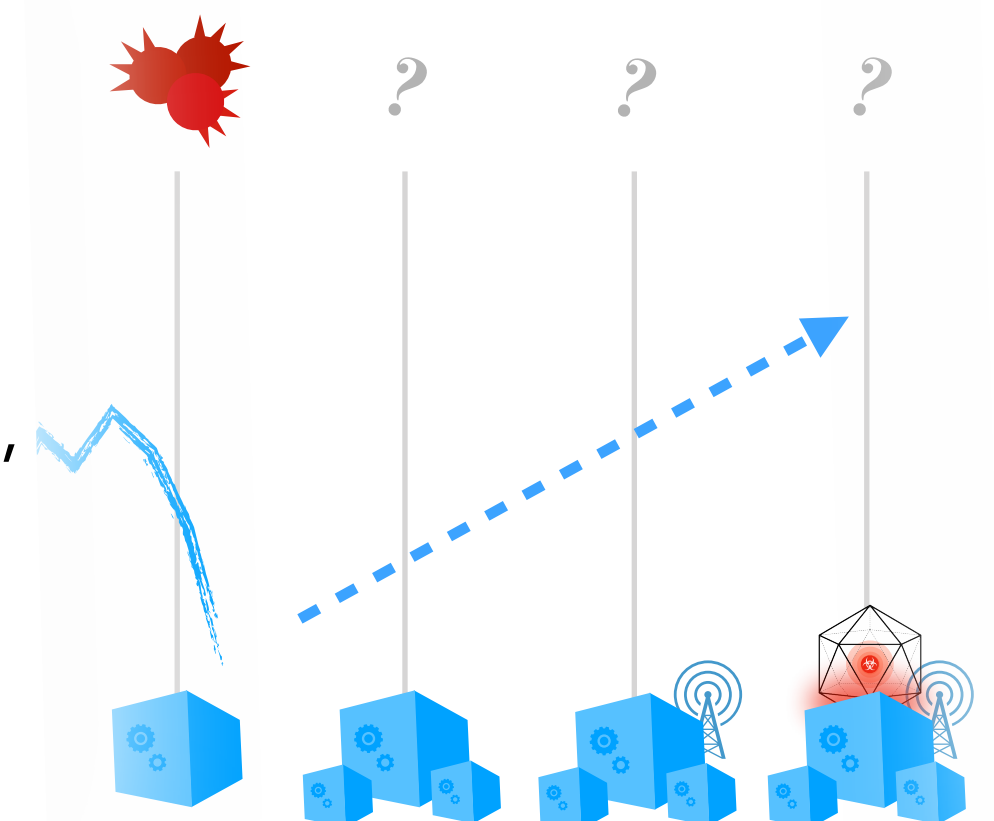
- How concept drift affects machine learning-based detectors in security settings
- The need for time-aware evaluations and metrics



Looking Ahead

Quo vadis?

- Discussion of the future of trustworthy ML for system security
- Robust feature development, universal adversarial perturbations, realizable backdoors, drift forecasting, and the role of abstractions towards the Platonic ideal of interesting behaviors

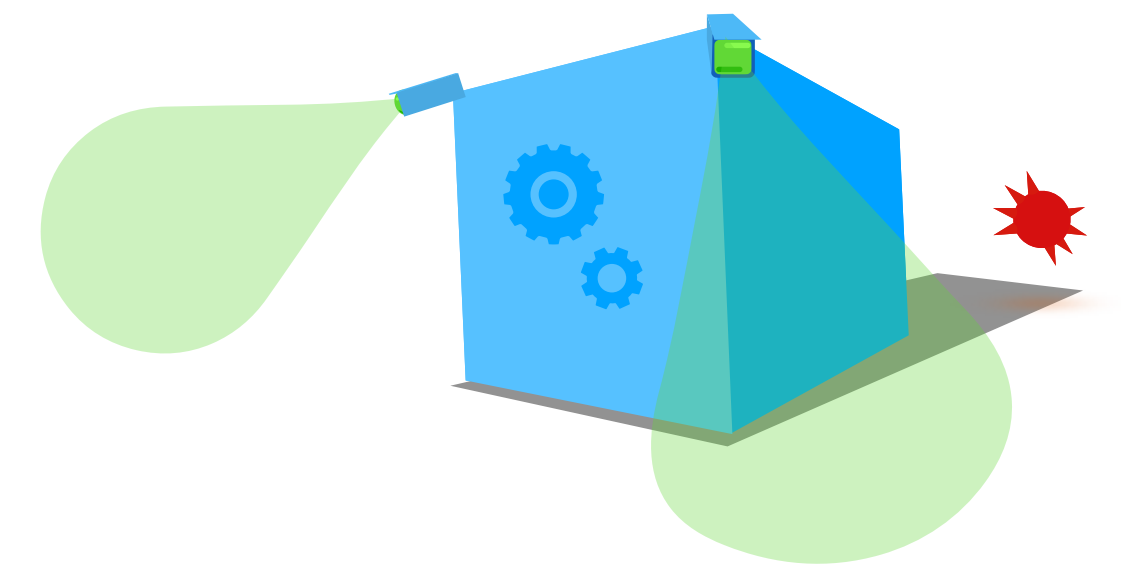


Outline

Focus

Adversarial ML evasion attacks against malware classifiers

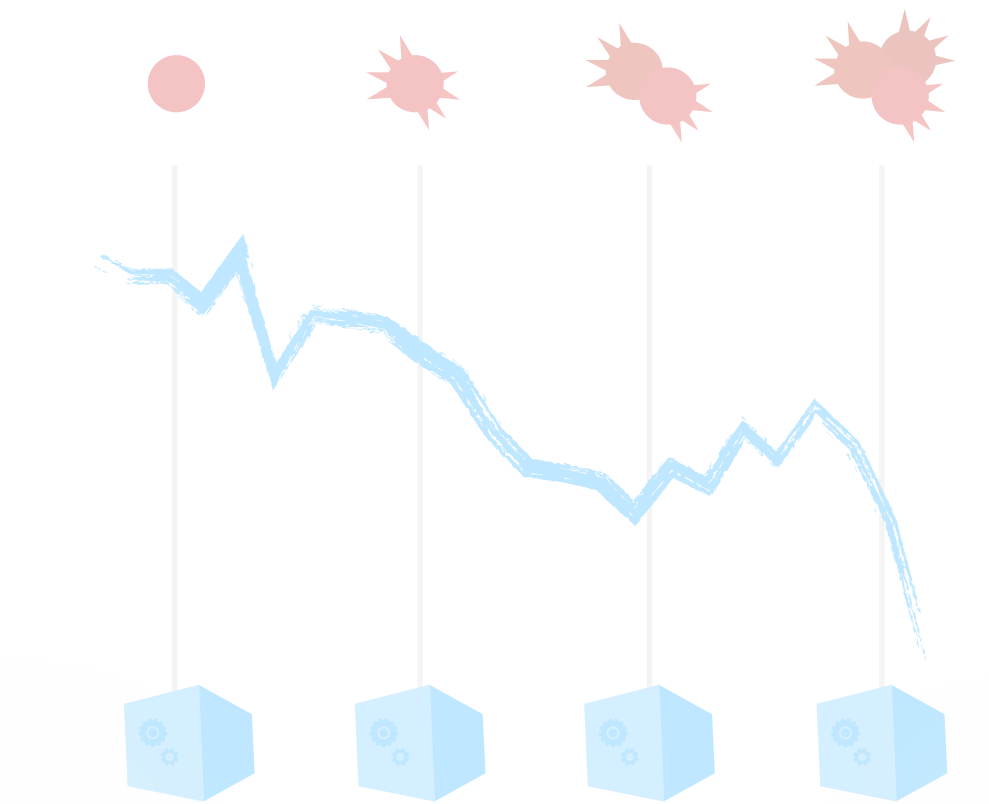
- Classic formulation of evasion attacks is ill-suited for reasoning about realizable evasive malware
- By reformulating, we can propose stronger attacks and easily compare against alternatives



Bigger Picture

Drifting scenarios caused by threats evolving over time

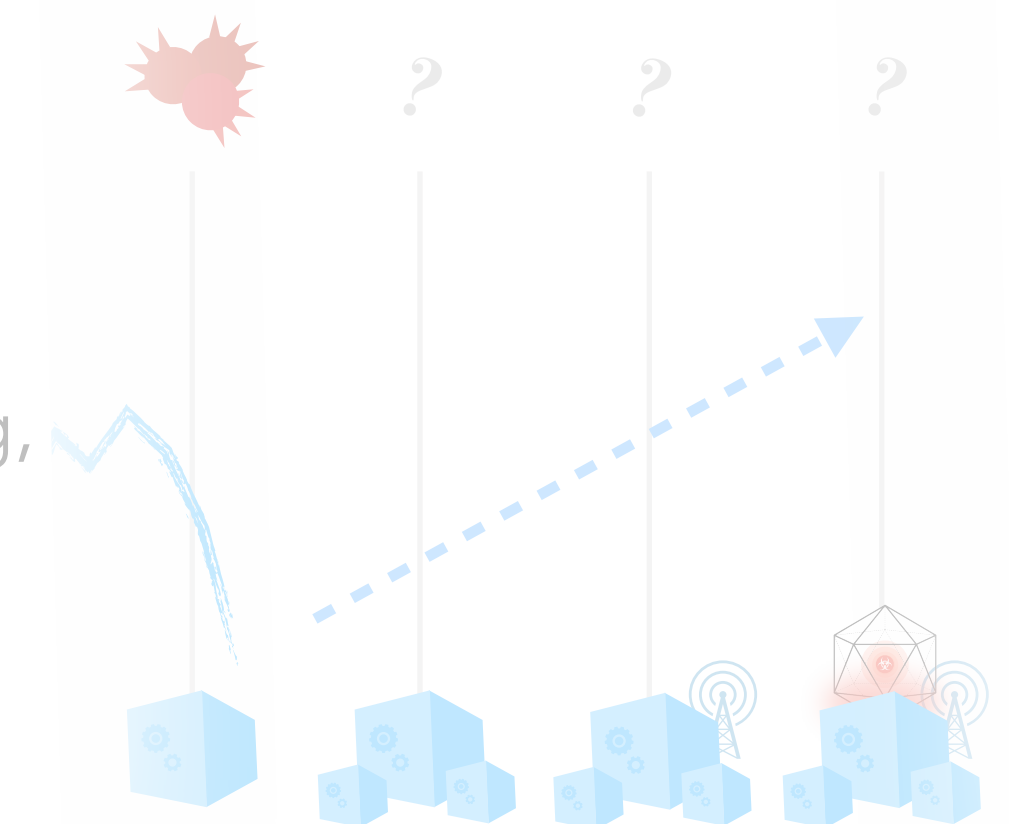
- How concept drift affects machine learning-based detectors in security settings
- The need for time-aware evaluations and metrics



Looking Ahead

Quo vadis?

- Discussion of the future of trustworthy ML for system security
- Robust feature development, universal adversarial perturbations, realizable backdoors, drift forecasting, and the role of abstractions towards the Platonic ideal of interesting behaviors



A Dystopian Future...

Pandas are forbidden!
Guilty of being too cute!



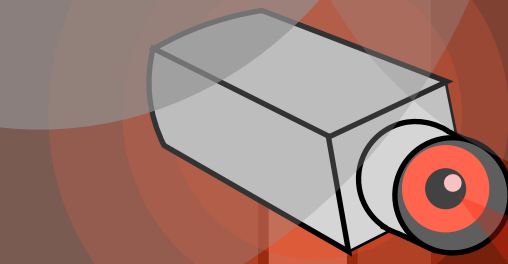
A Dystopian Future...

Pandas are forbidden!
Guilty of being too cute!



A Dystopian Future...

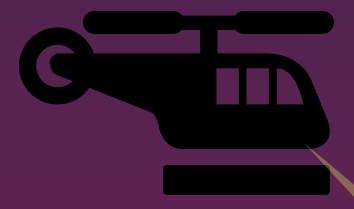
Pandas are forbidden!
Guilty of being too cute!



Luckily, pandas are fluent in math...



Luckily, pandas are fluent in math...



Luckily, pandas are fluent in math...

Intriguing properties of neural networks

Christian Szegedy
Google Inc.

Wojciech Zaremba
New York University

Ilya Sutskever
Google Inc.

Joan Bruna
New York University

Dumitru Erhan
Google Inc.

Ian Goodfellow
University of Montreal

Rob Fergus
New York University
Facebook Inc.

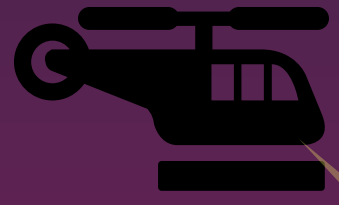
Abstract

Deep neural networks are highly expressive models that have recently achieved state of the art performance on speech and visual recognition tasks. While their expressiveness is the reason they succeed, it also causes them to learn uninterpretable solutions that could have counter-intuitive properties. In this paper we report two such properties.

First, we find that there is no distinction between individual high level units and random linear combinations of high level units, according to various methods of unit analysis. It suggests that it is the space, rather than the individual units, that contains the semantic information in the high layers of neural networks.

Second, we find that deep neural networks learn input-output mappings that are insensitive to a significant extent. We can cause the network to misclassify an input in a way that is hardly perceptible. In addition, the specific nature of the perturbation can

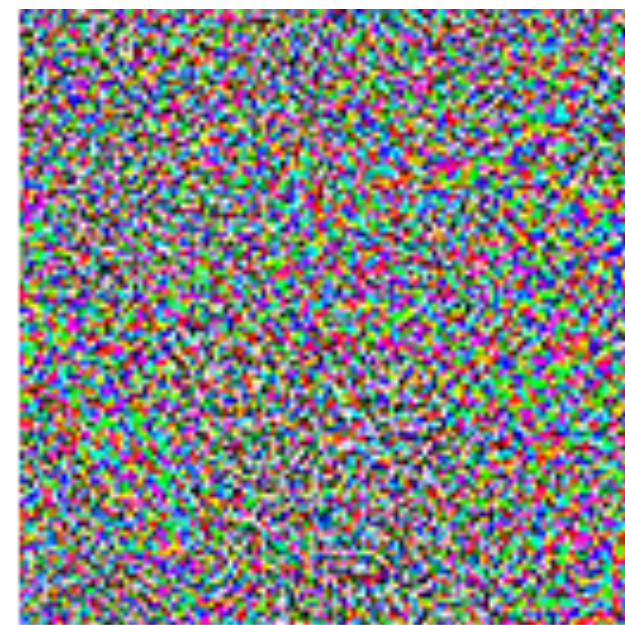
4 [cs.CV] 19 Feb 2014



Luckily, pandas are fluent in math...



+



=



"panda"
57.7% confidence

"gibbon"
99.3% confidence

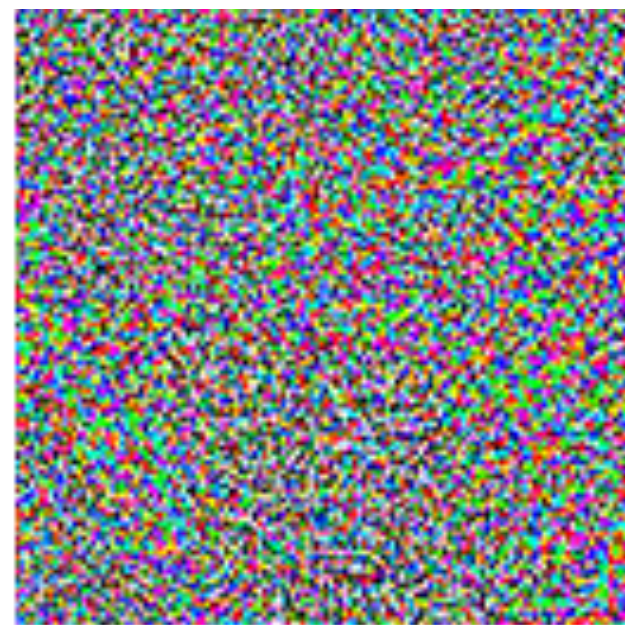


Luckily, pandas are fluent in math...



"panda"
57.7% confidence

+



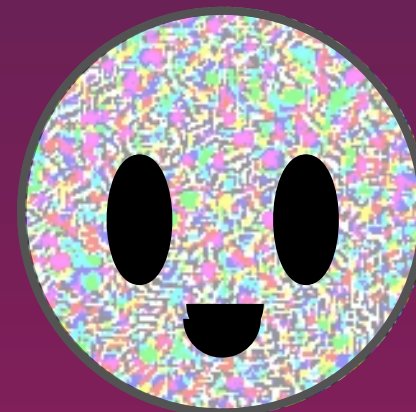
=



"gibbon"
99.3% confidence

~

Feature-space
noise mask



What happens in the **problem space**, i.e., the real world?



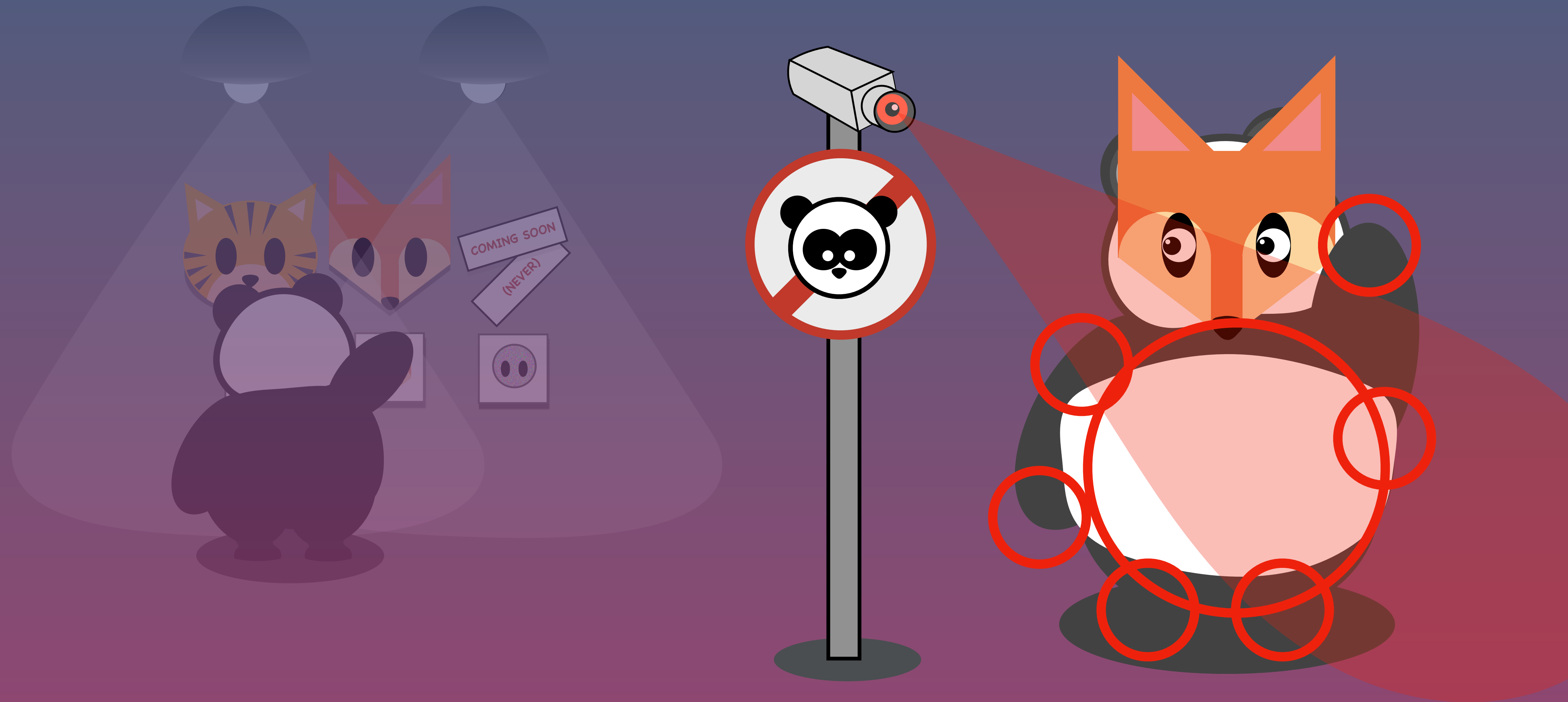
What happens in the **problem space**, i.e., the real world?



What happens in the **problem space**, i.e., the real world?



What happens in the **problem space**, i.e., the real world?



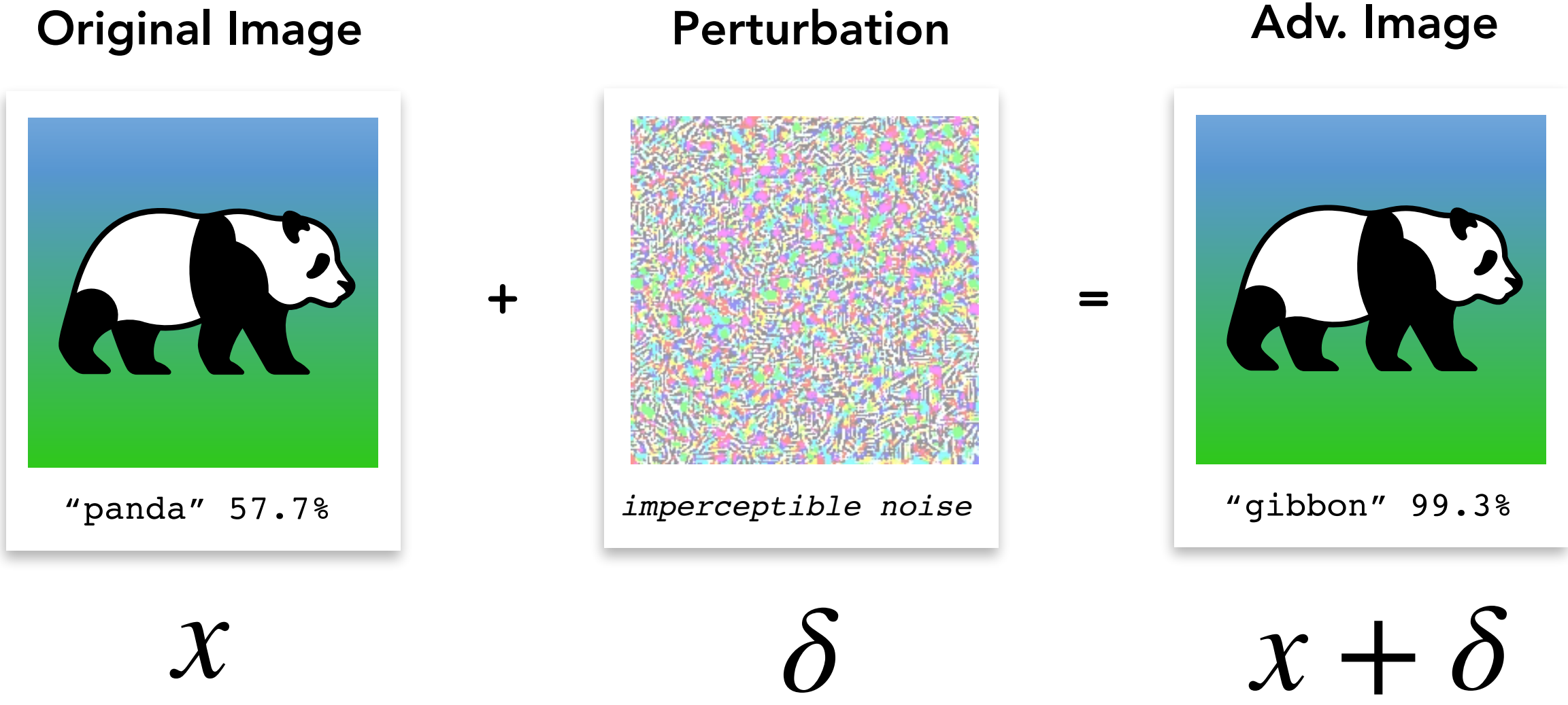
What happens in the **problem space**, i.e., the real world?



Let's Analyze What Happened

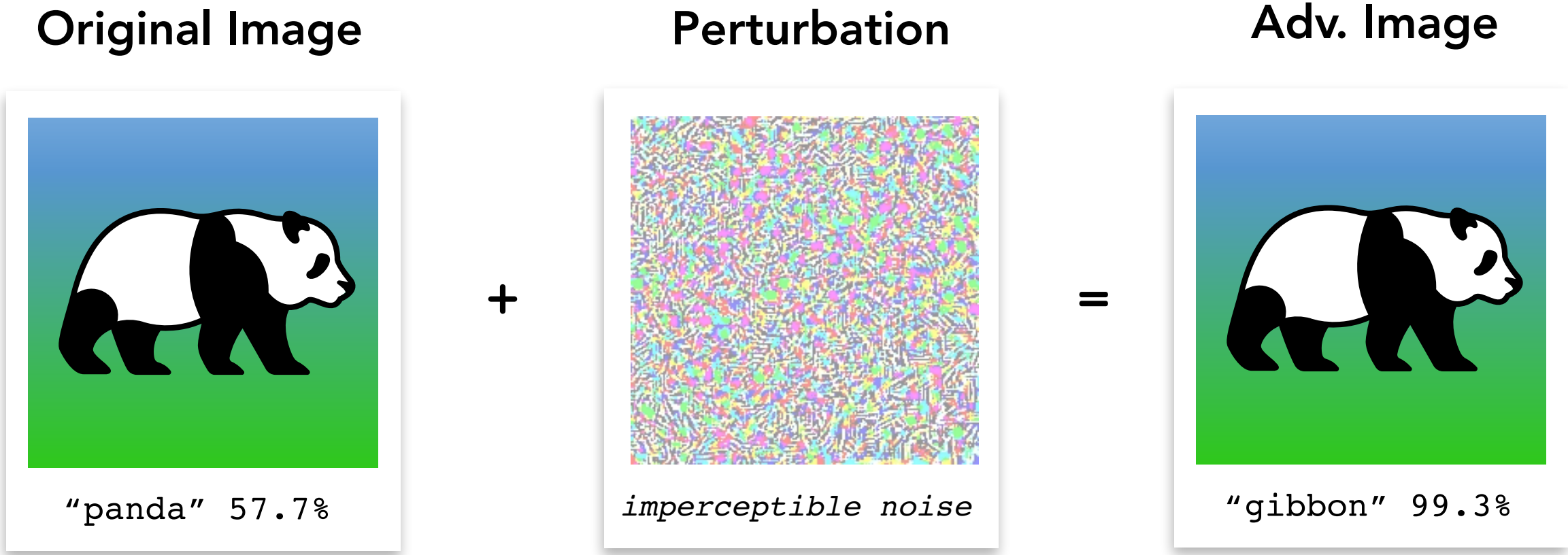
Let's Analyze What Happened

Feature-Space Attacks



Let's Analyze What Happened

Feature-Space Attacks



x

δ

$x + \delta$

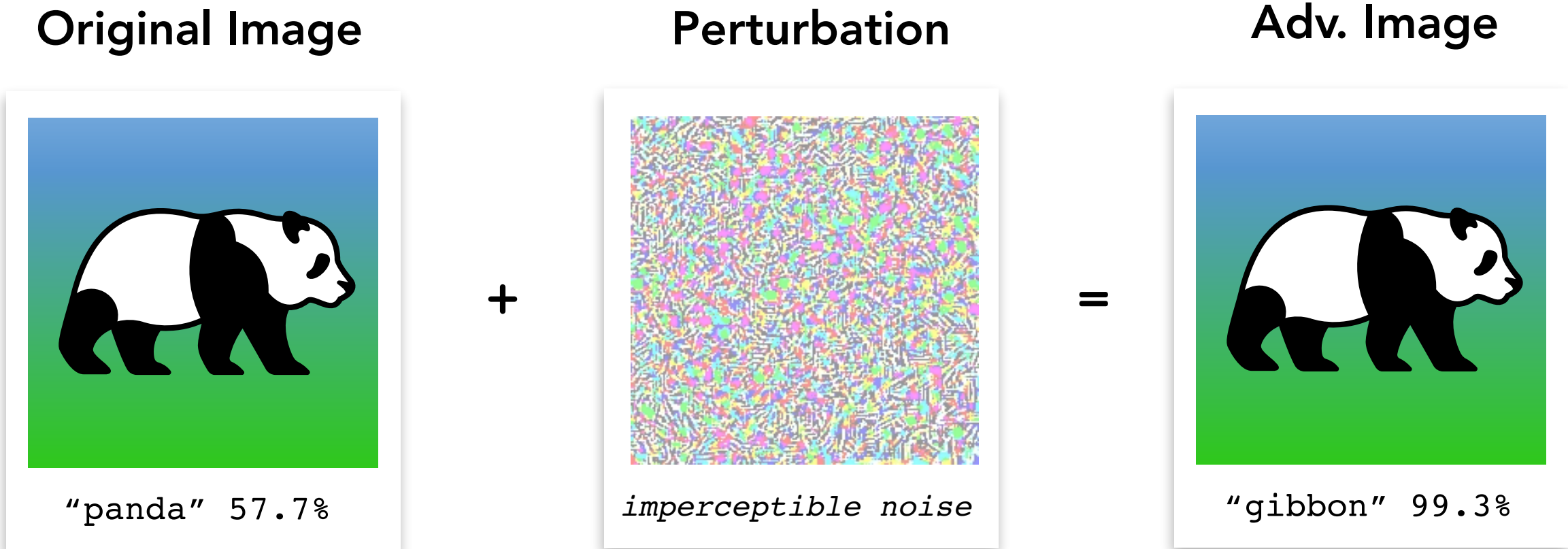


Optimization

minimize $_{\delta} ||\delta||_p + c \cdot f(x + \delta)$

Let's Analyze What Happened

Feature-Space Attacks



x

δ

$x + \delta$

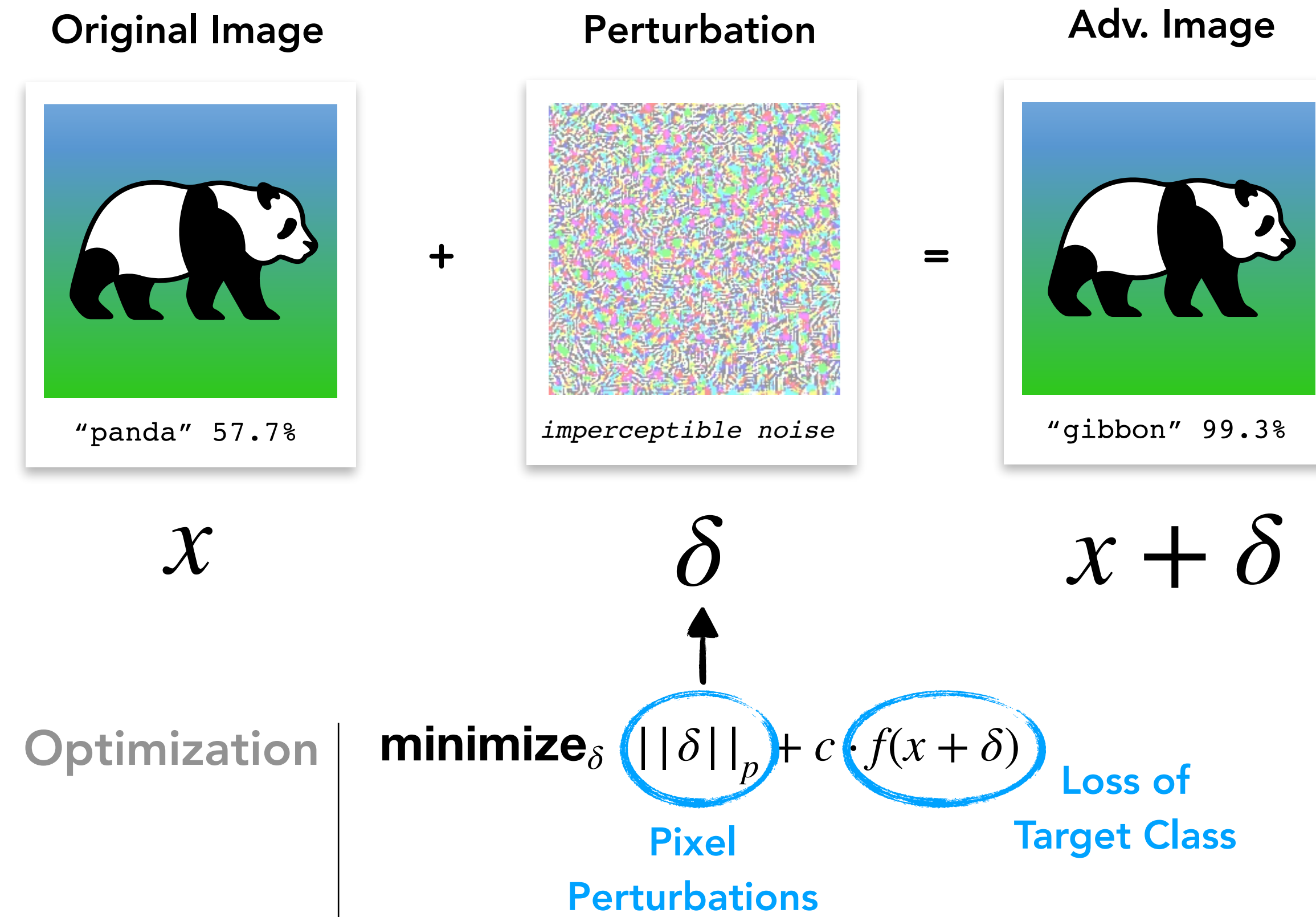
Optimization

minimize δ $\|\delta\|_p + c \cdot f(x + \delta)$

Pixel Perturbations

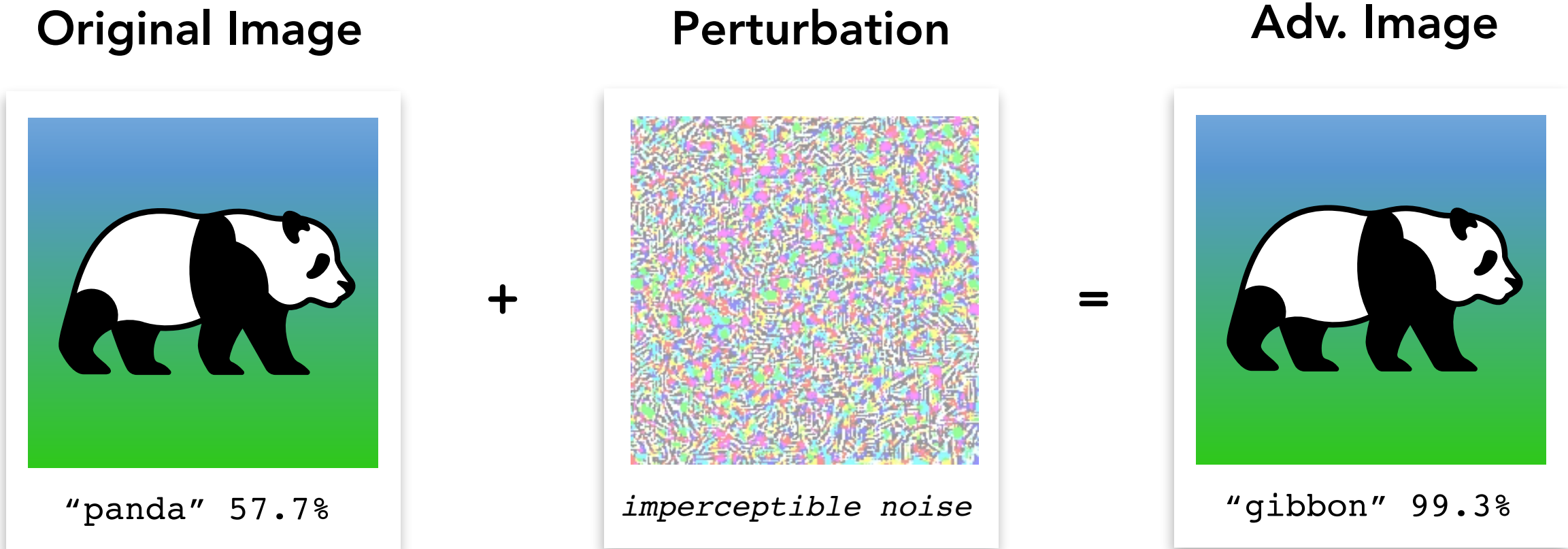
Let's Analyze What Happened

Feature-Space Attacks



Let's Analyze What Happened

Feature-Space Attacks



x

δ

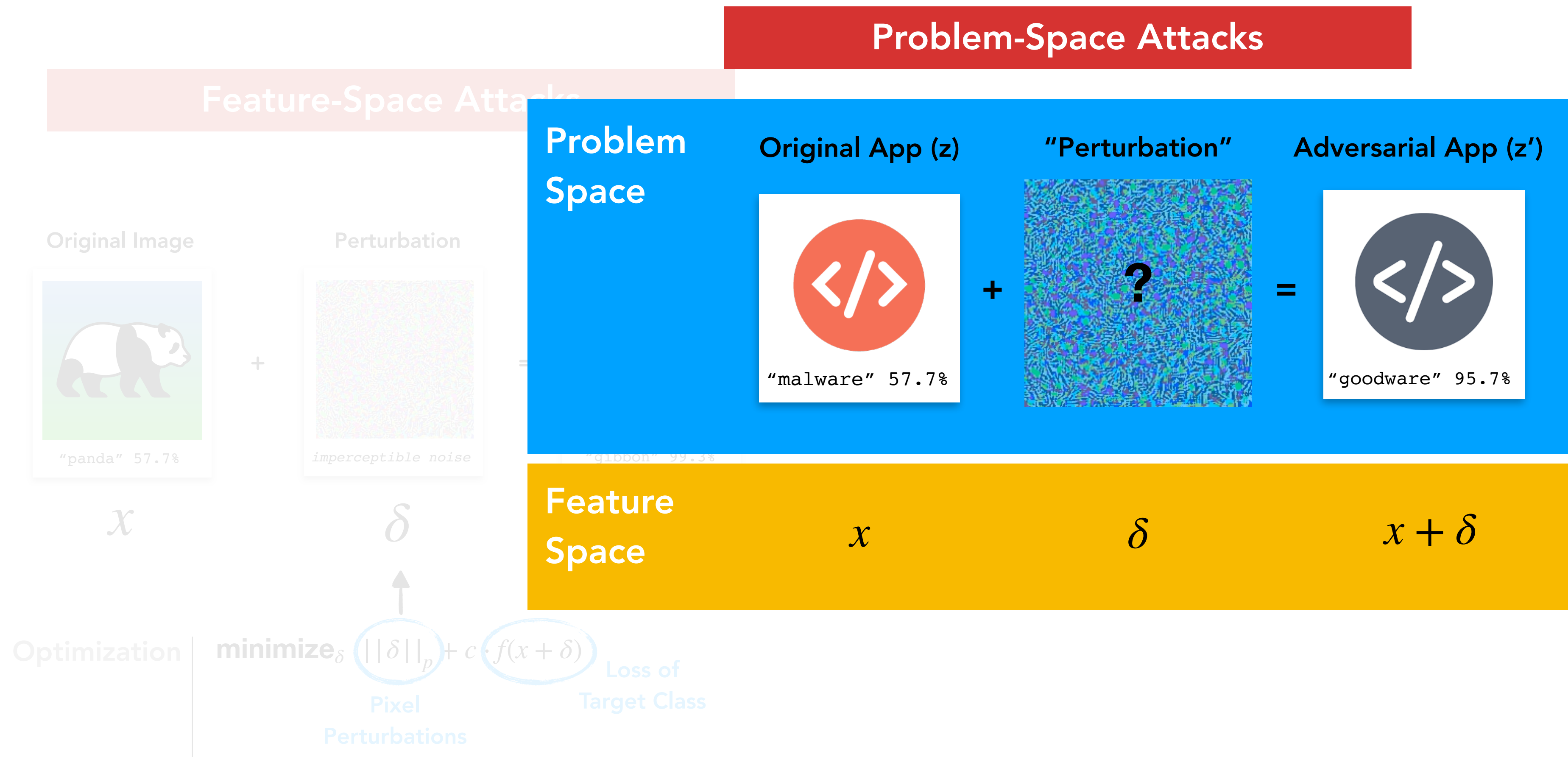
$x + \delta$

Optimization

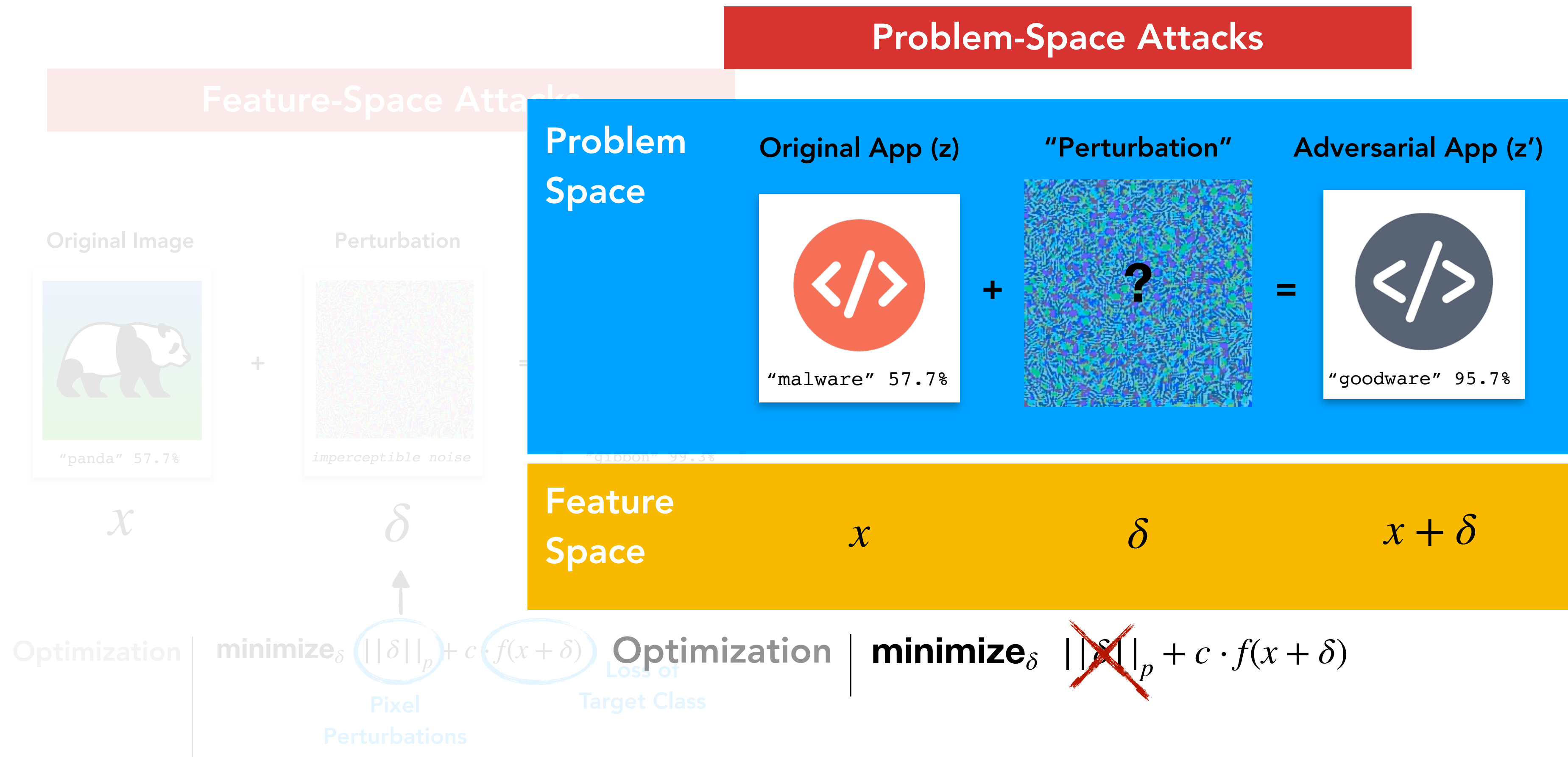
minimize δ $\|\delta\|_p + c \cdot f(x + \delta)$

Pixel Perturbations Loss of Target Class

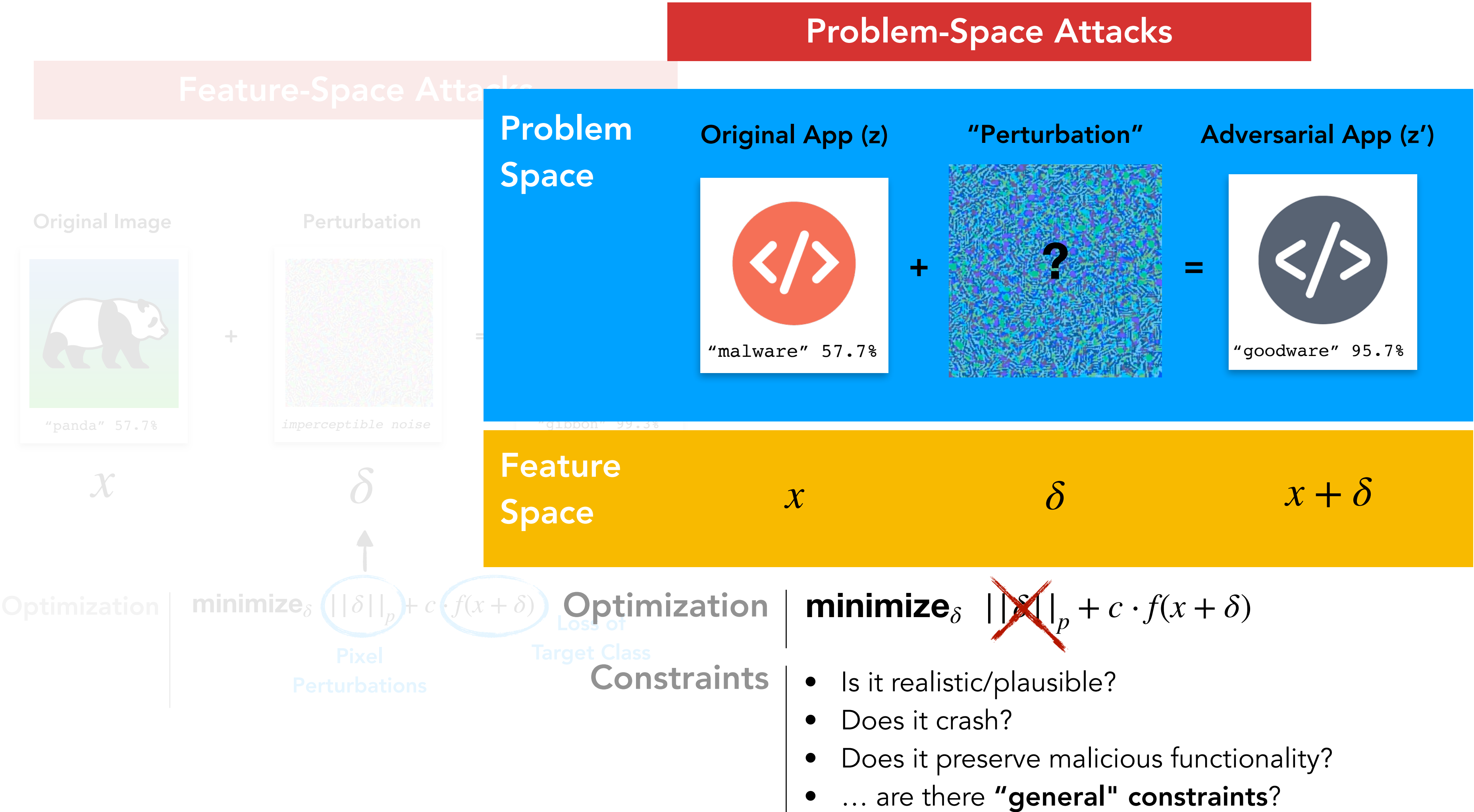
Let's Analyze What Happened



Let's Analyze What Happened



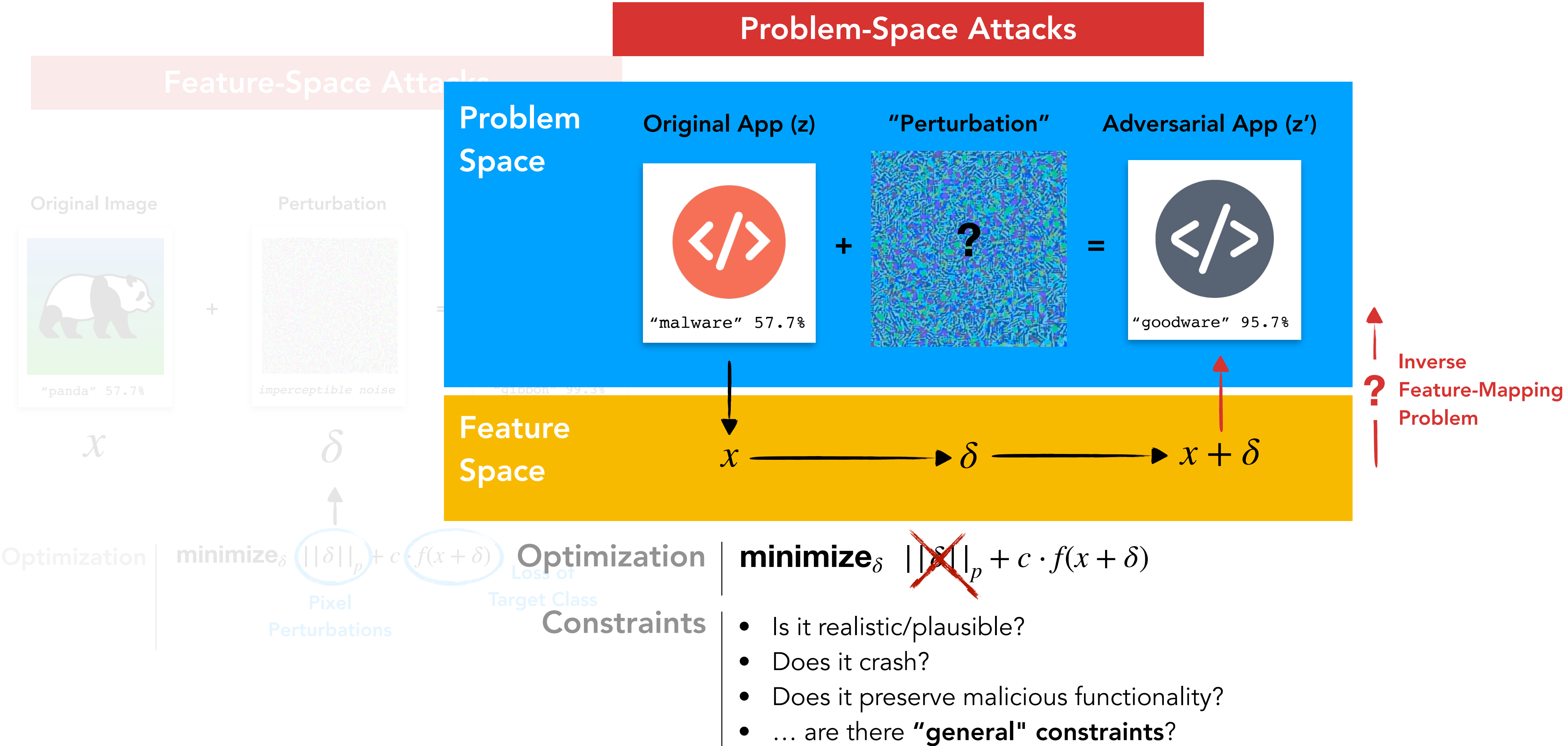
Let's Analyze What Happened



[IEEE S&P 2020] Intriguing Properties of Adversarial ML Attacks in the Problem Space

<https://s2lab.cs.ucl.ac.uk/projects/intriguing>

Let's Analyze What Happened



[IEEE S&P 2020] Intriguing Properties of Adversarial ML Attacks in the Problem Space

<https://s2lab.cs.ucl.ac.uk/projects/intriguing>

Formalization

Problem-Space Constraints

Problem-Space Constraints

Available Transformations

How can you alter problem-space objects?

Problem-Space Constraints

 Available Transformations

How can you alter problem-space objects?

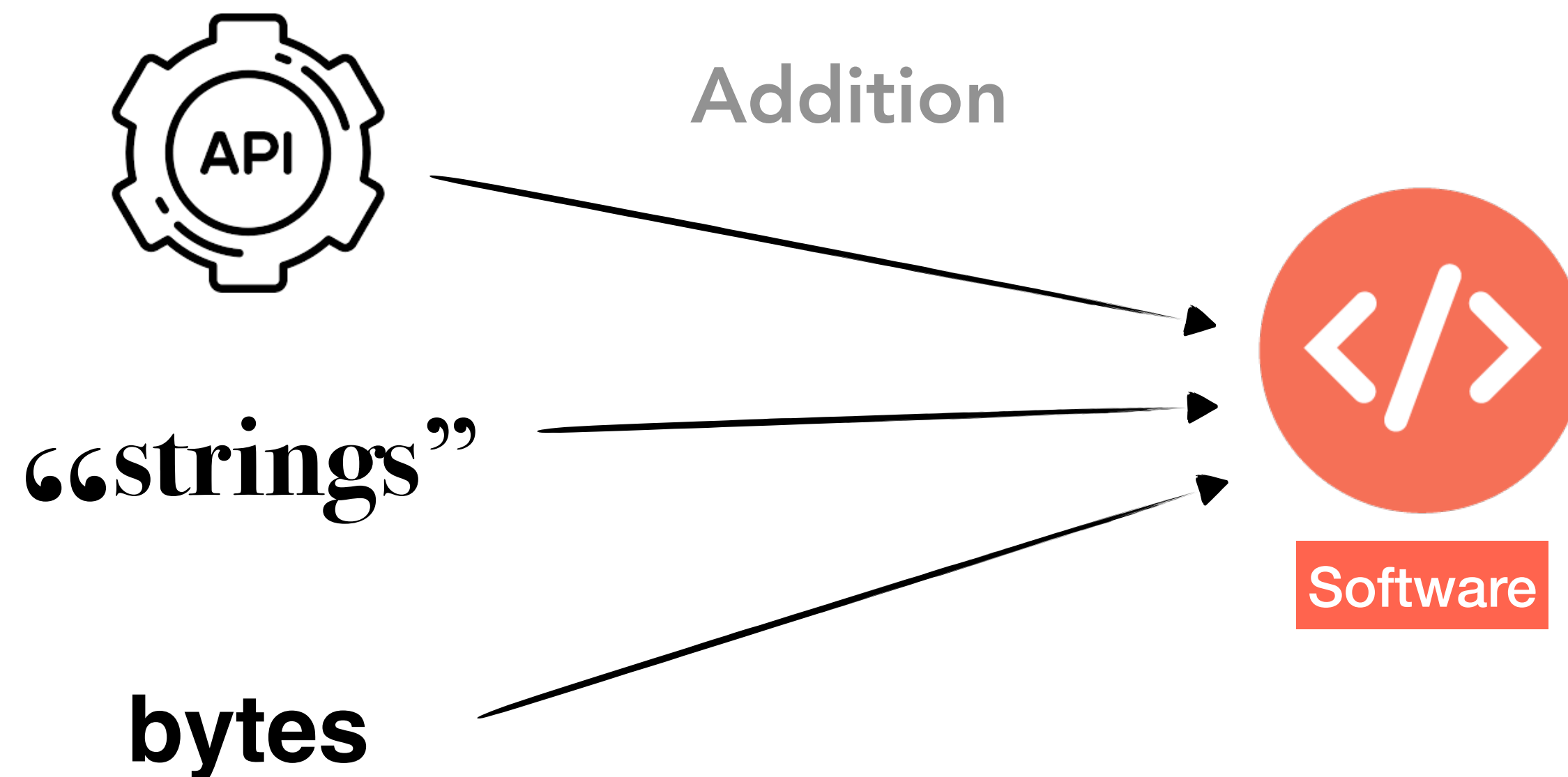


Software

Problem-Space Constraints

Available Transformations

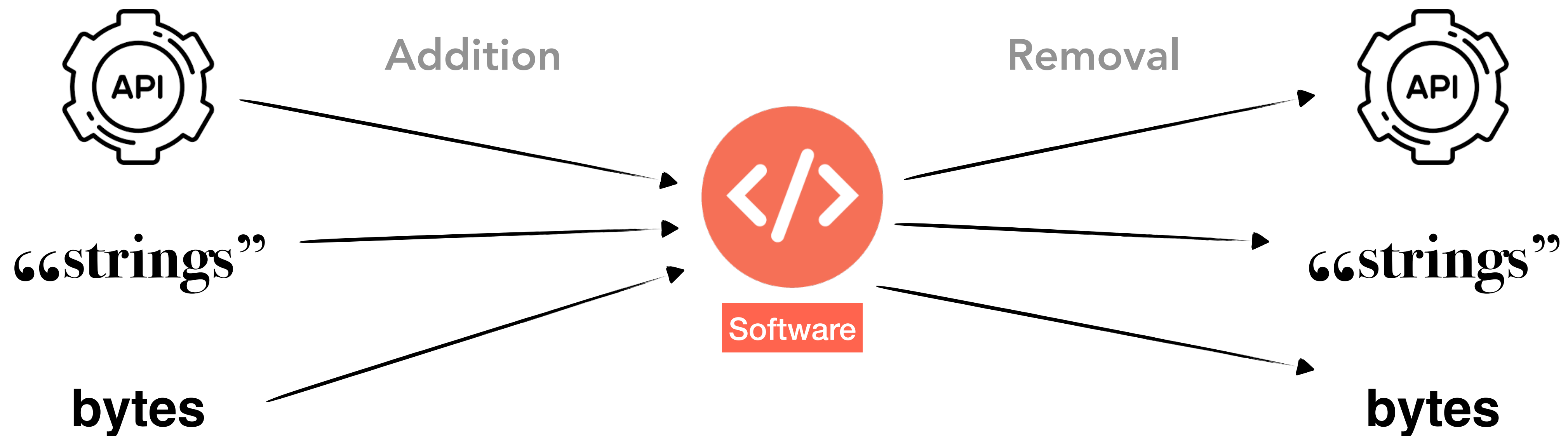
How can you alter problem-space objects?



Problem-Space Constraints

Available Transformations

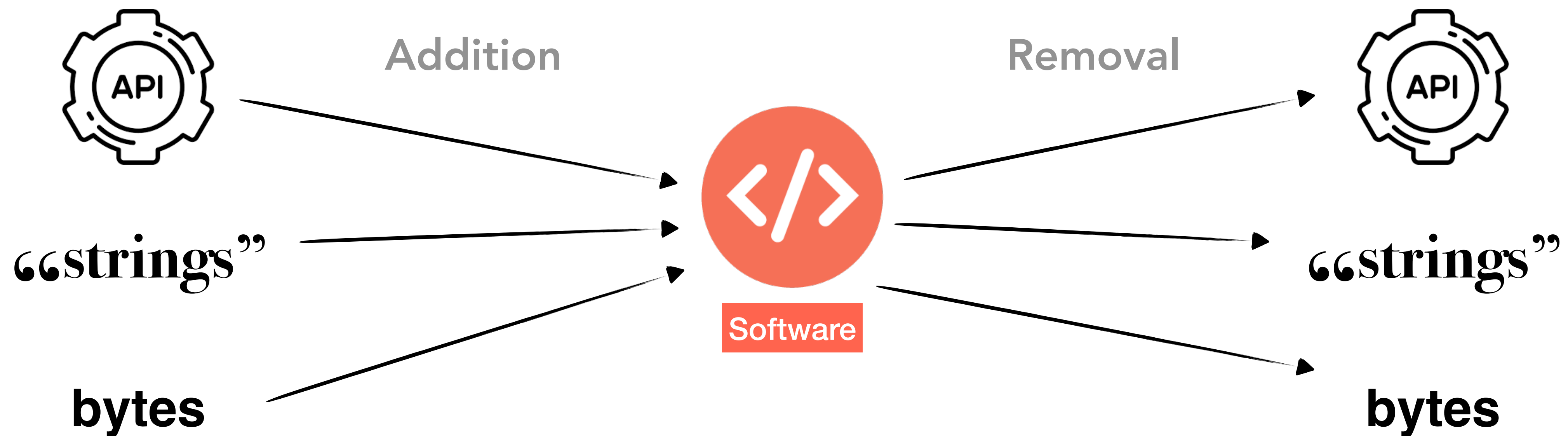
How can you alter problem-space objects?



Problem-Space Constraints

Available Transformations

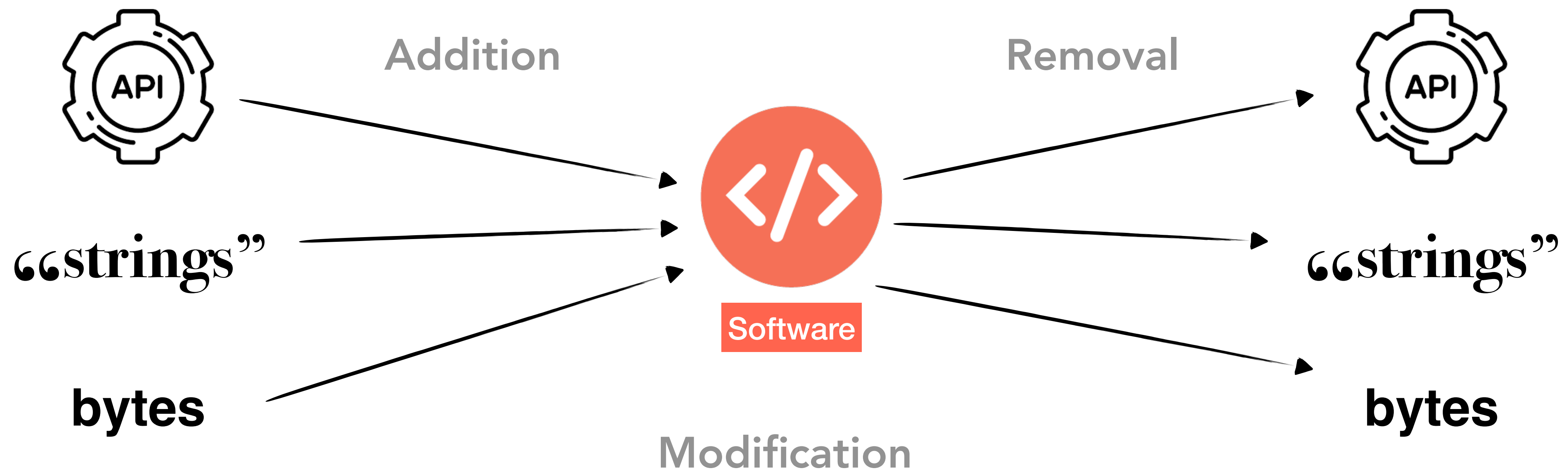
How can you alter problem-space objects?



Problem-Space Constraints

Available Transformations

How can you alter problem-space objects?



Problem-Space Constraints

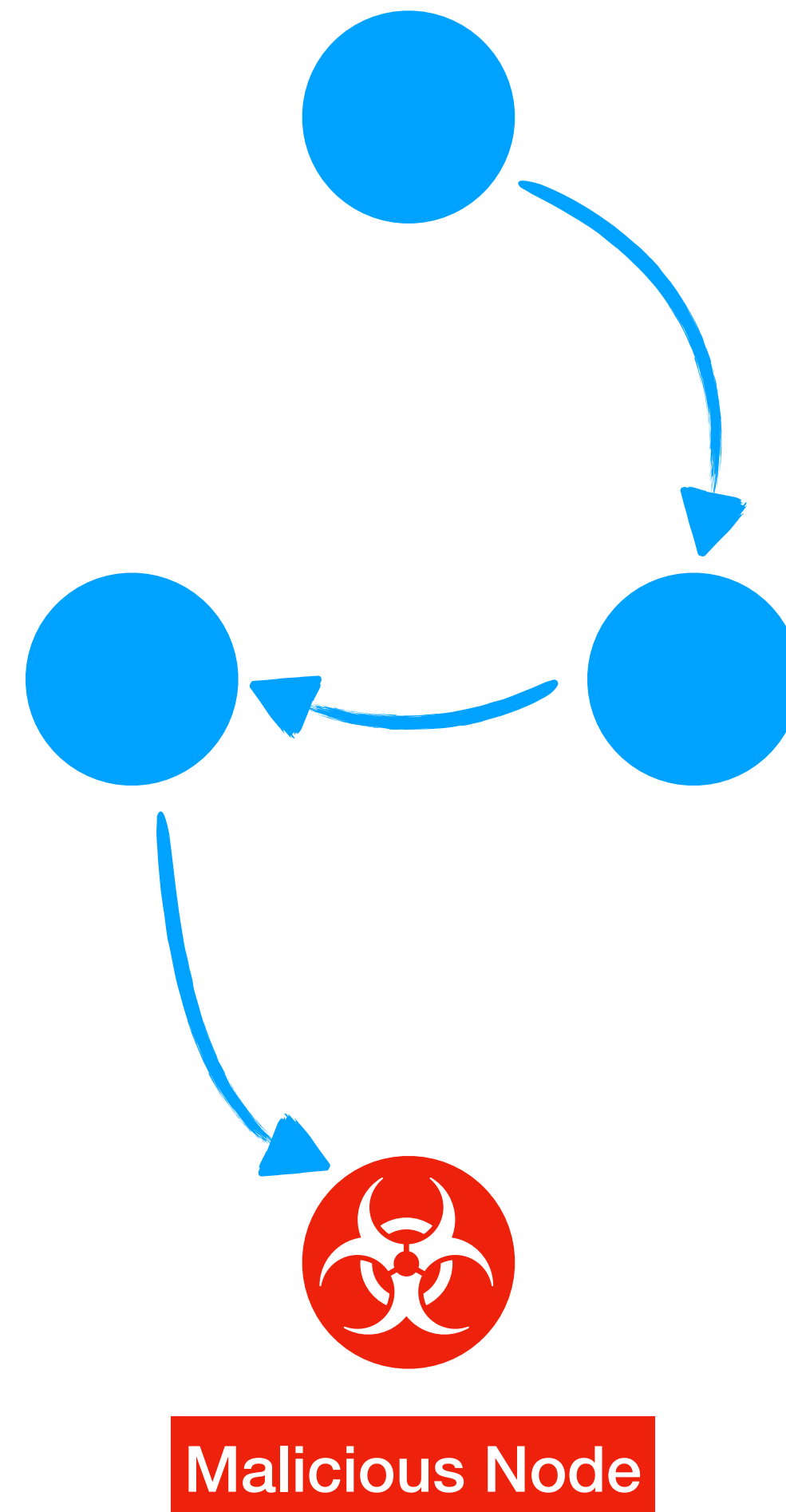
 Available Transformations

Problem-Space Constraints

 **Preserved Semantics**

 **Available Transformations**

Which semantics do you preserve? How?
Which automatic tests can verify it?

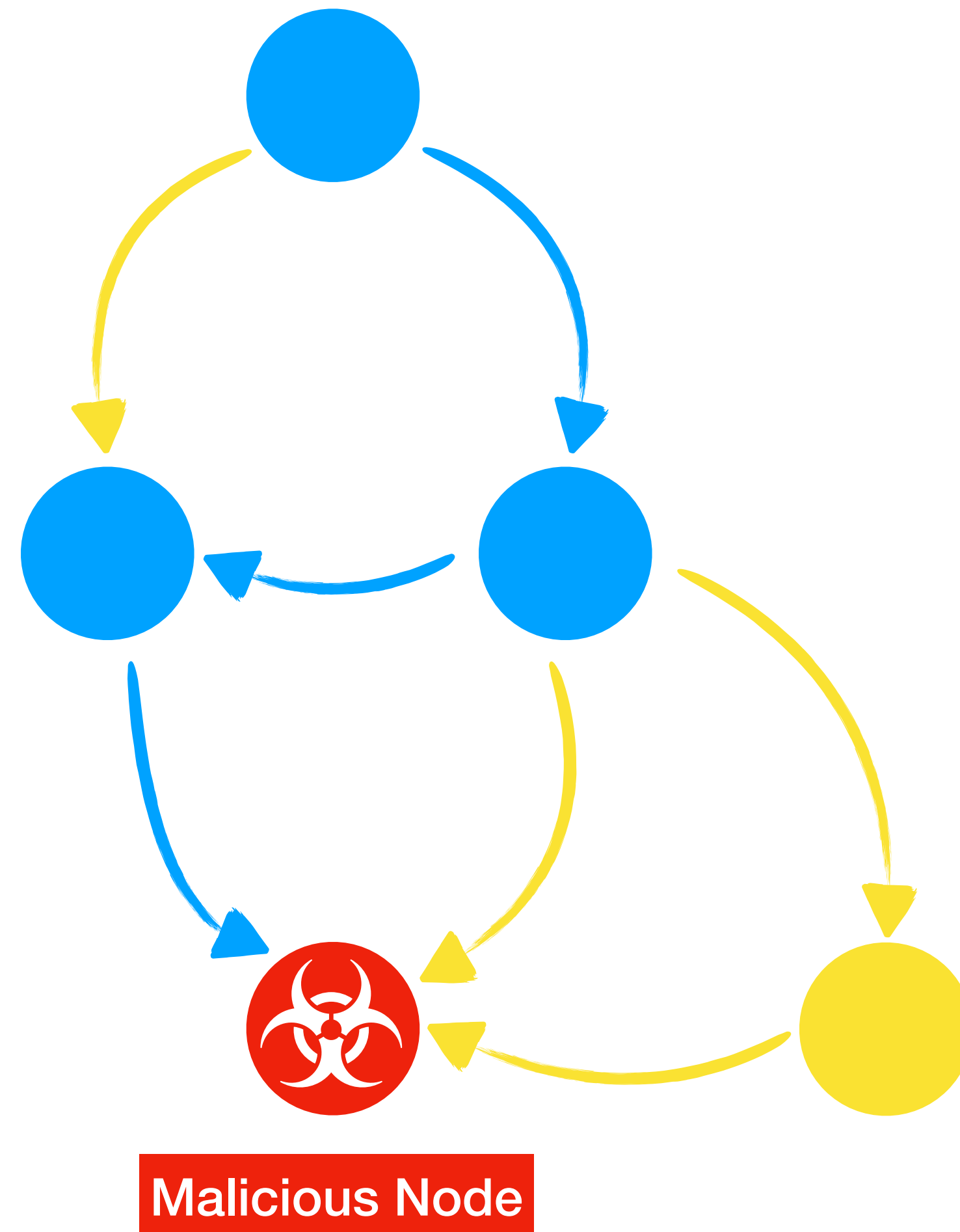


Problem-Space Constraints

 **Preserved Semantics**

 **Available Transformations**

Which semantics do you preserve? How?
Which automatic tests can verify it?



Problem-Space Constraints

 **Preserved Semantics**

 **Available Transformations**

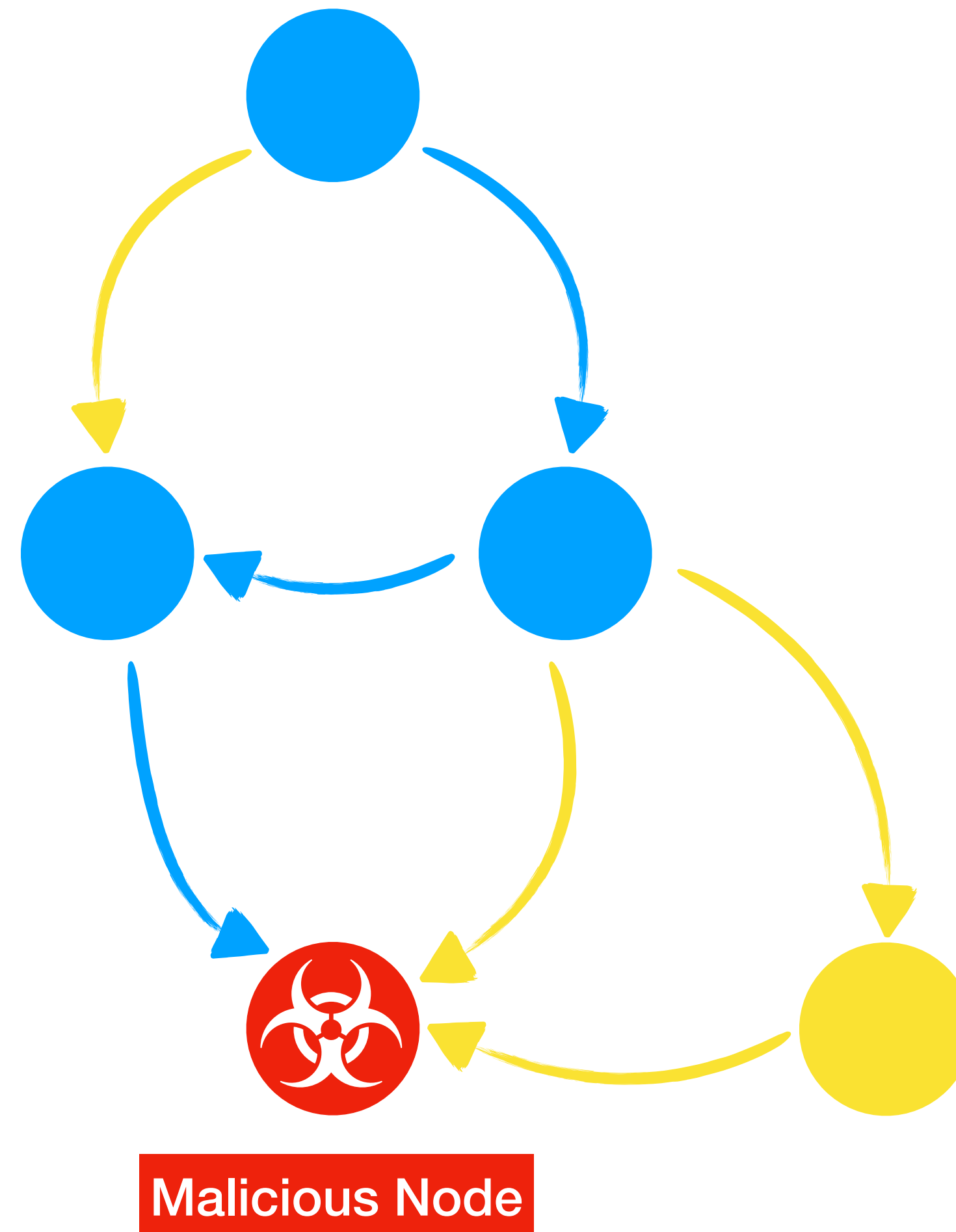
Test Suite

- Does it crash?
- Does it still communicate with CnC?
- Does it still encrypt the /home/ folder?

By Construction

- Add no-op operations
- Ensure it is not executed at runtime

Which semantics do you preserve? How?
Which automatic tests can verify it?



Problem-Space Constraints




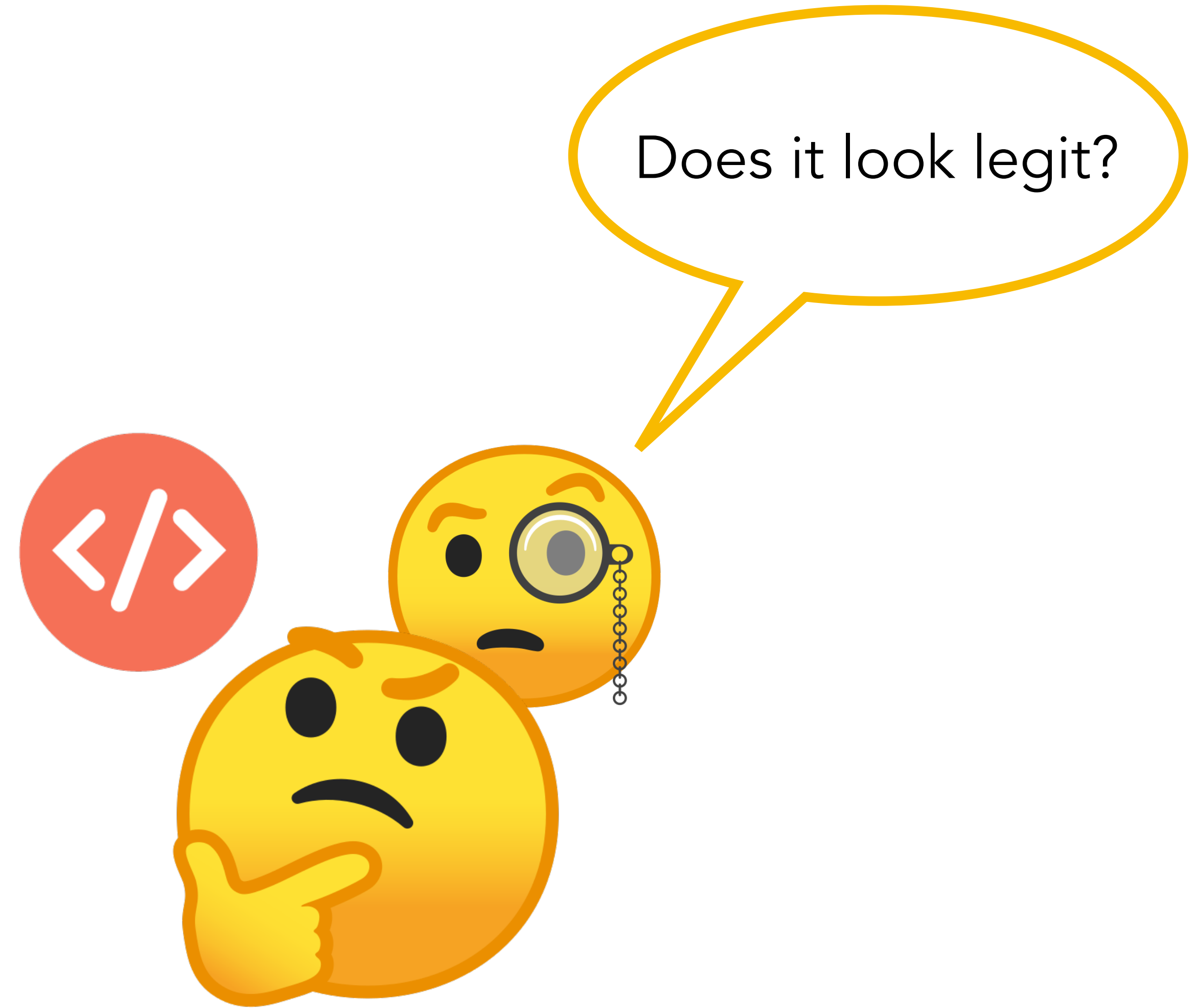
Preserved Semantics



Available Transformations

Problem-Space Constraints

-  **Plausibility**
-  Preserved Semantics
-  Available Transformations



Problem-Space Constraints

 **Plausibility**

 Preserved Semantics

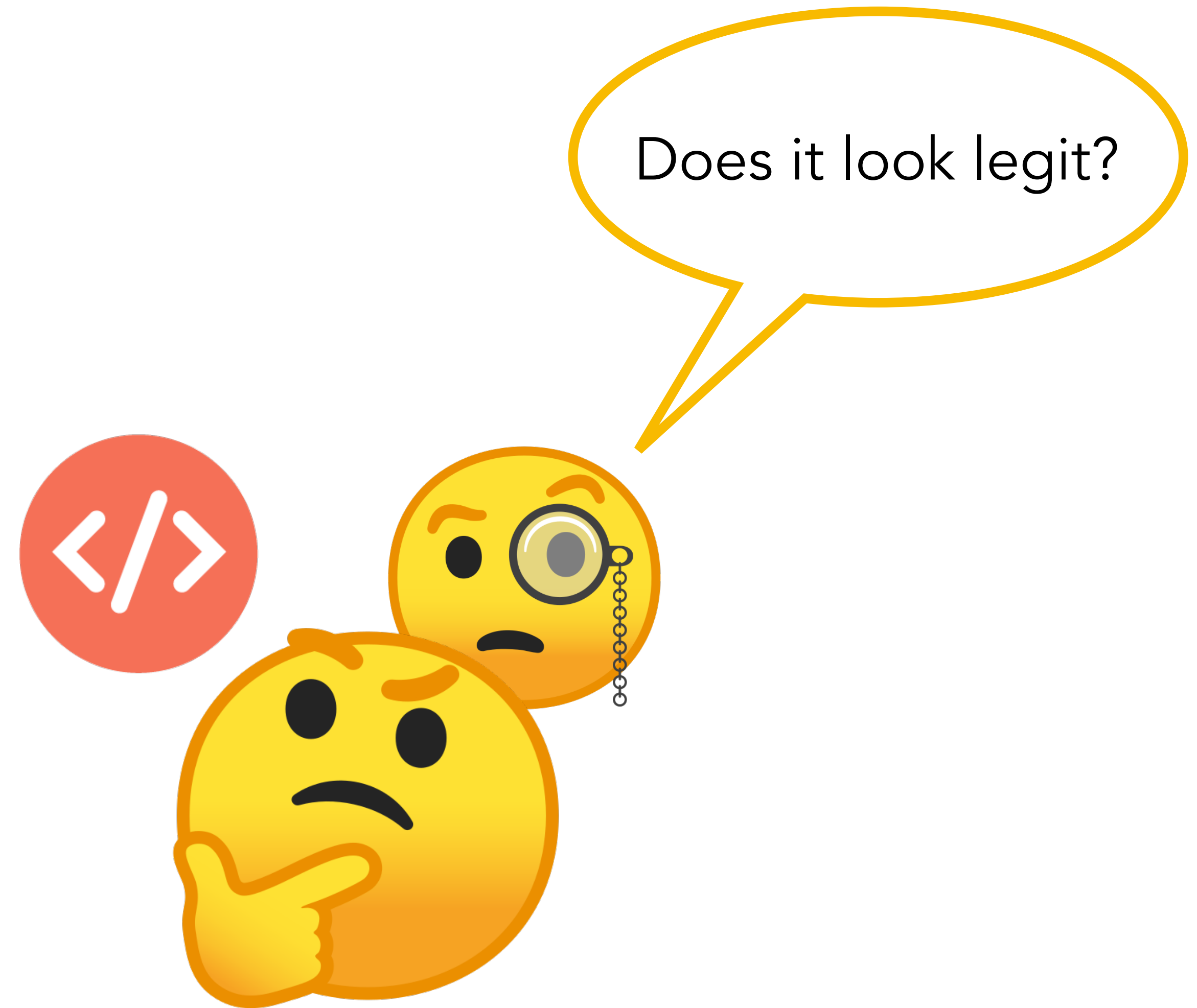
 Available Transformations

Test Suite

- User studies
- Automated heuristics

By Construction

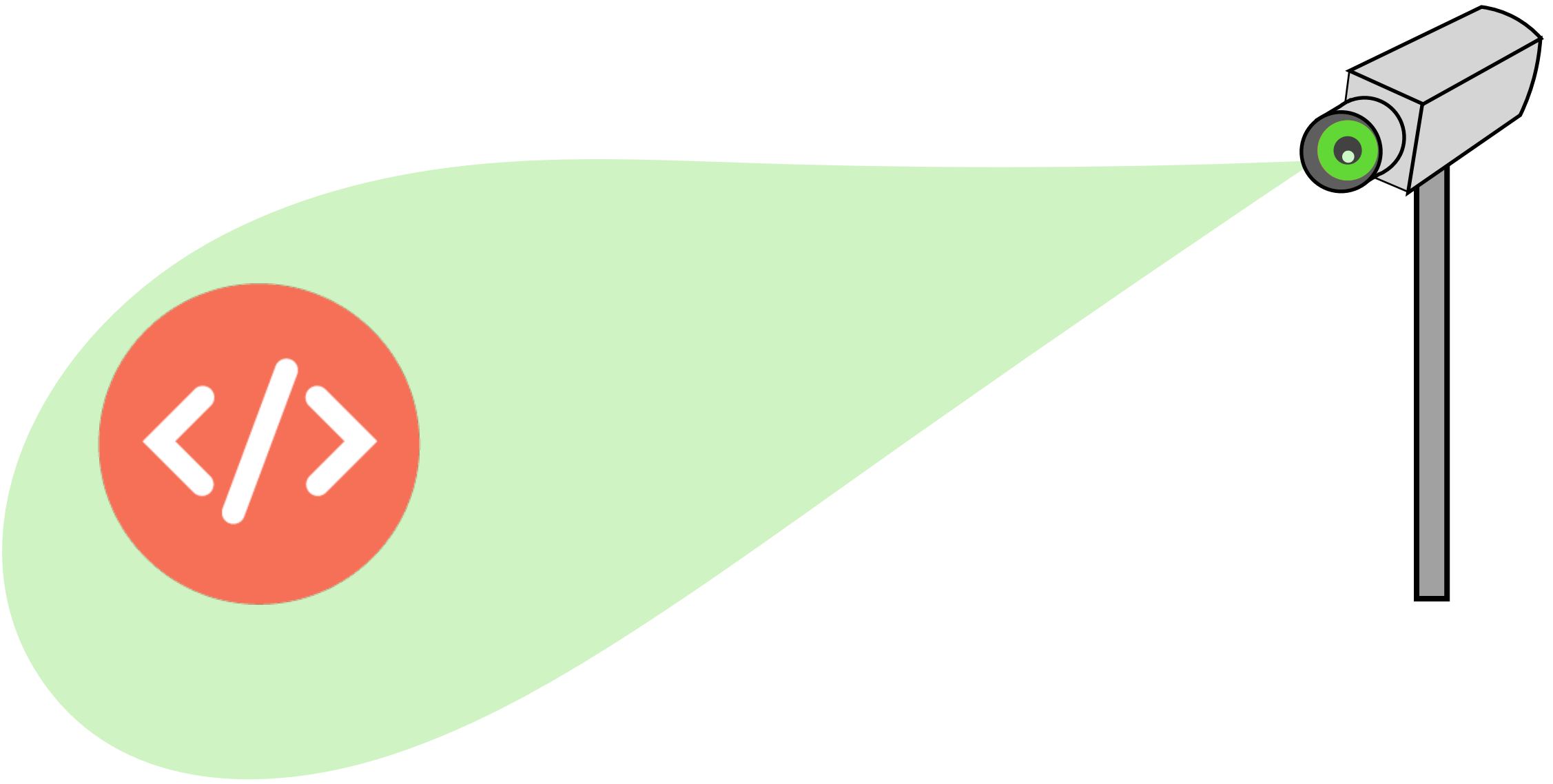
- Taking precautions during mutation







Problem-Space Constraints

-  **Plausibility**
-  Preserved Semantics
-  Available Transformations

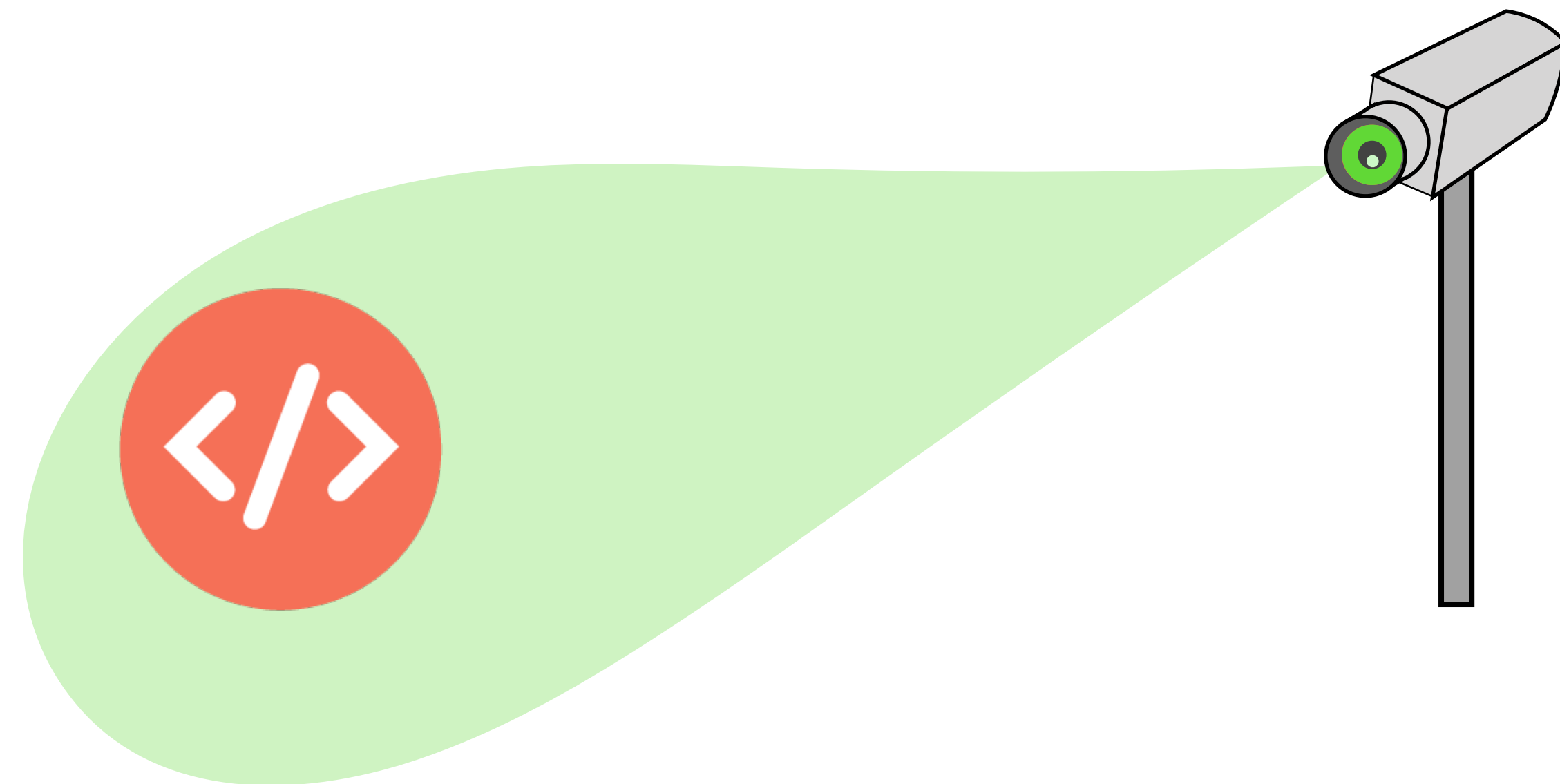
Which preprocessing are you considering?







Problem-Space Constraints

-  **Robustness to Preprocessing**
-  Plausibility
-  Preserved Semantics
-  Available Transformations

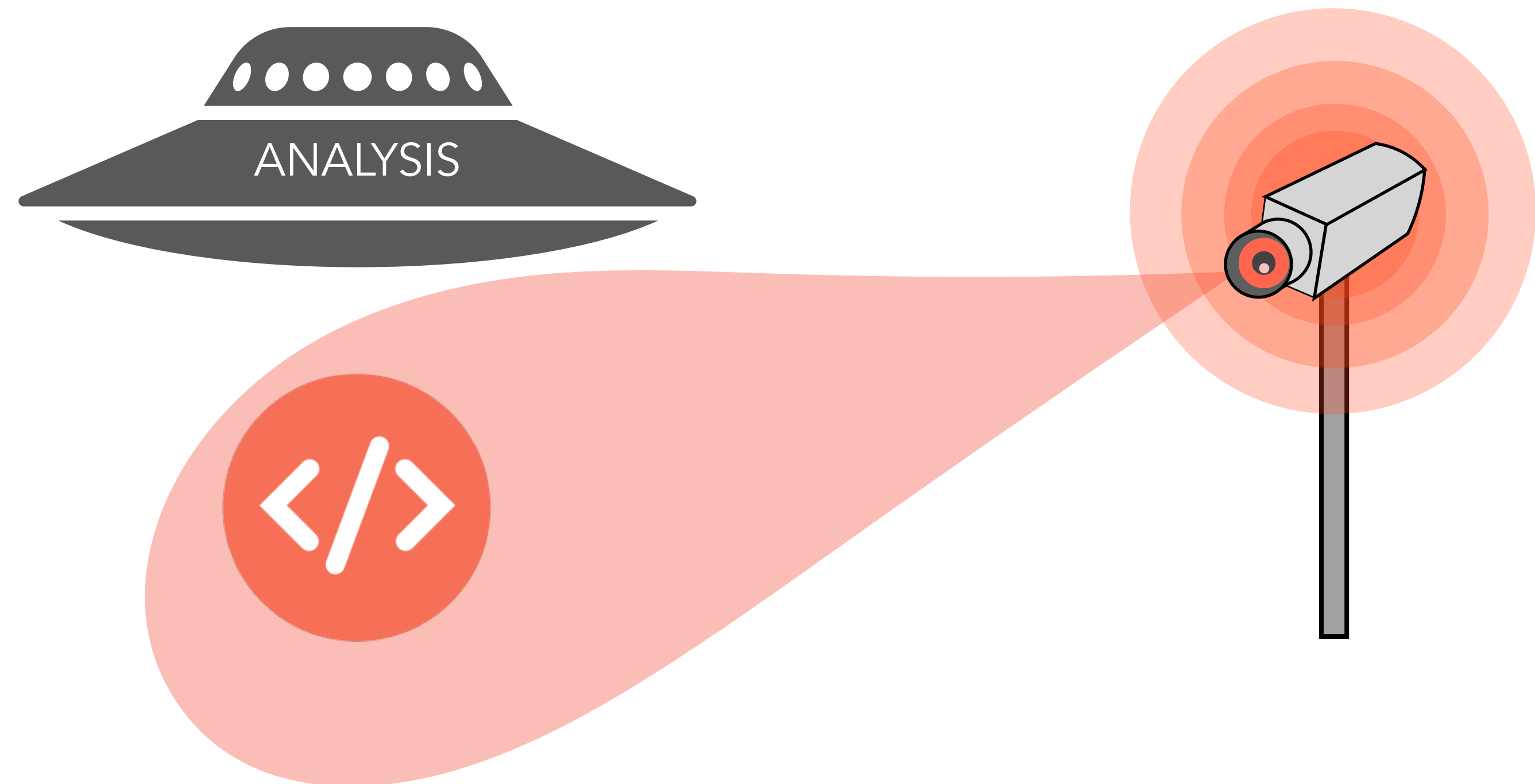
Which preprocessing are you considering?



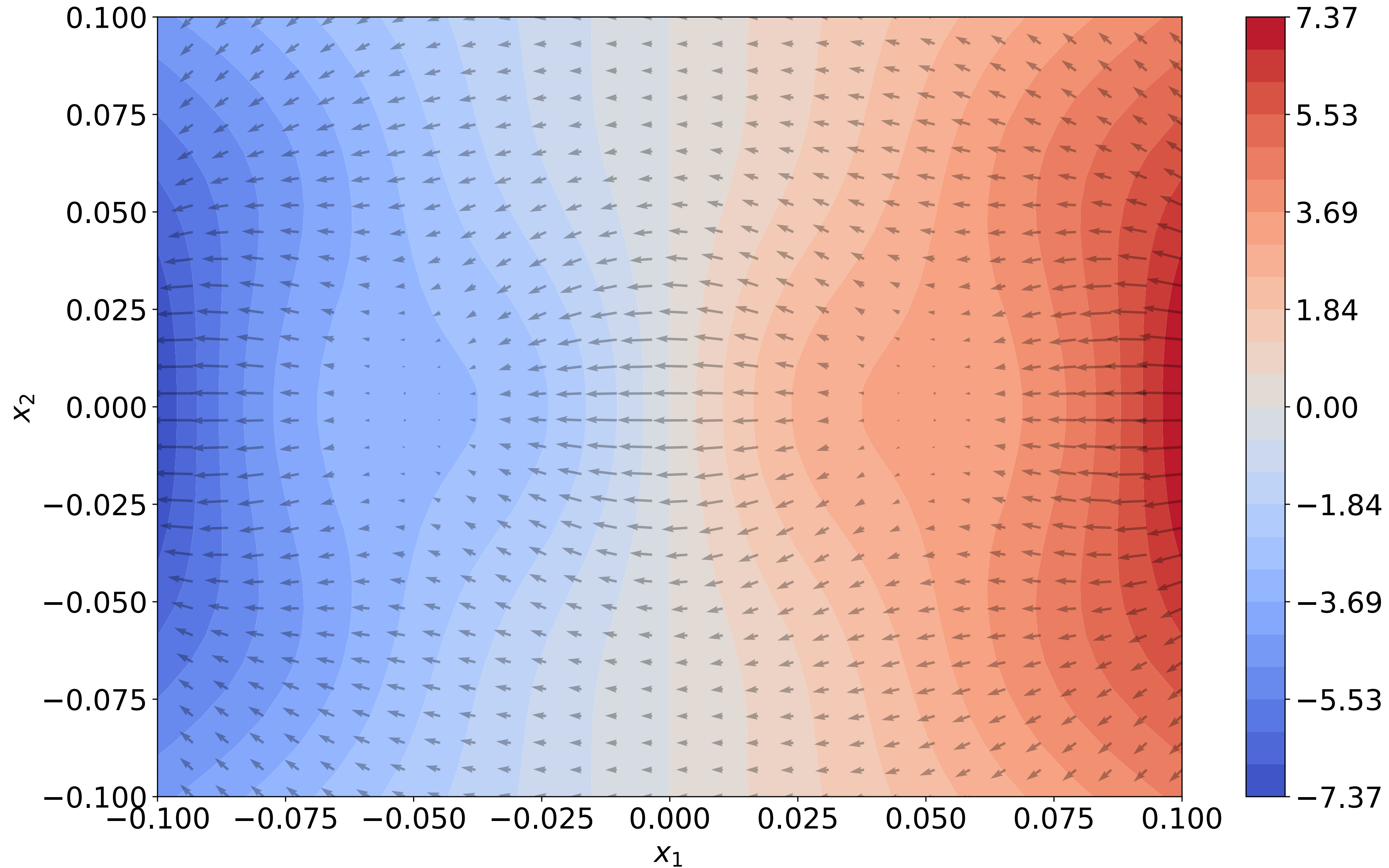
Problem-Space Constraints

-  **Robustness to Preprocessing**
-  Plausibility
-  Preserved Semantics
-  Available Transformations

Which preprocessing are you considering?



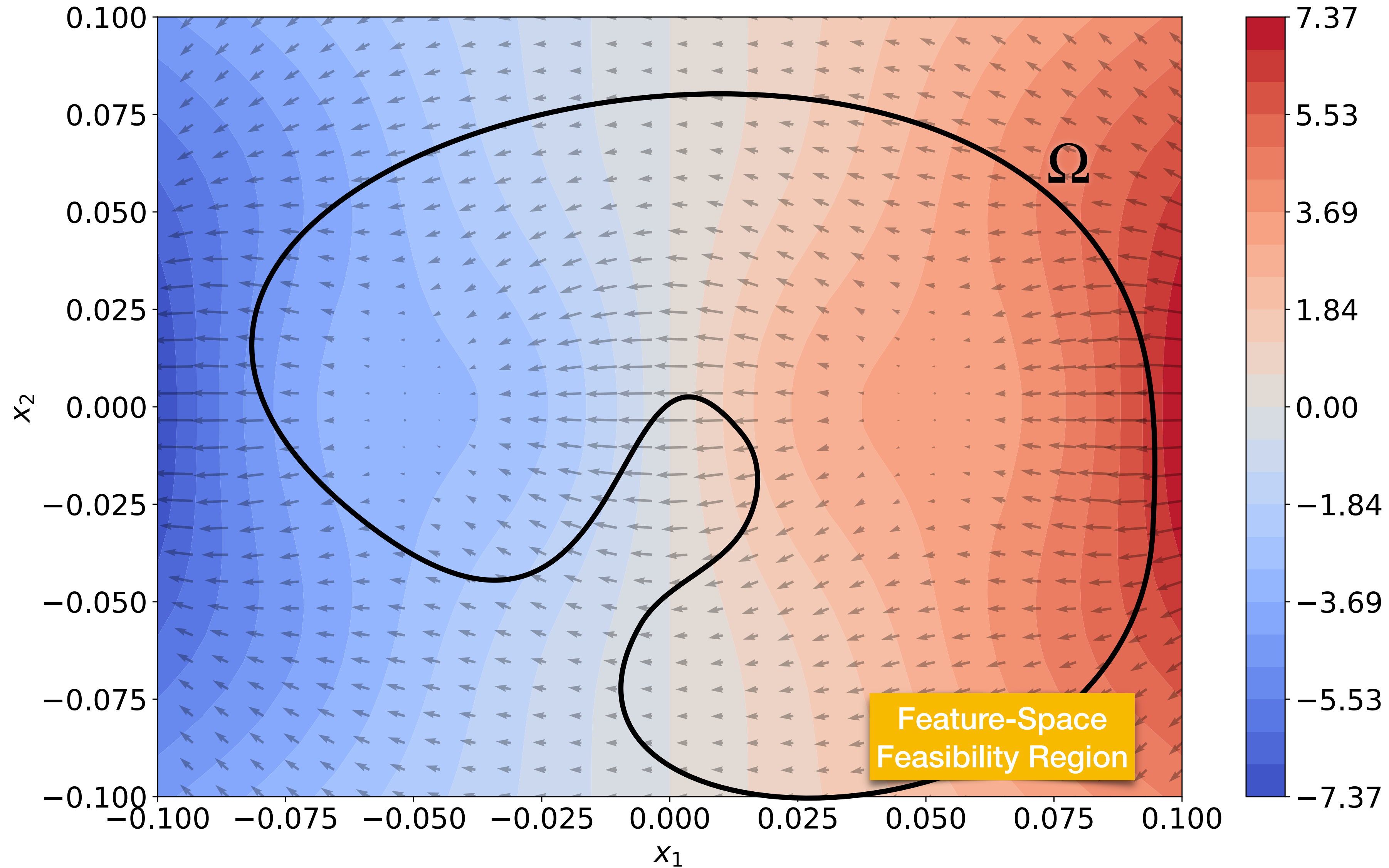
Side-effect Features: the ONE Slide to Remember



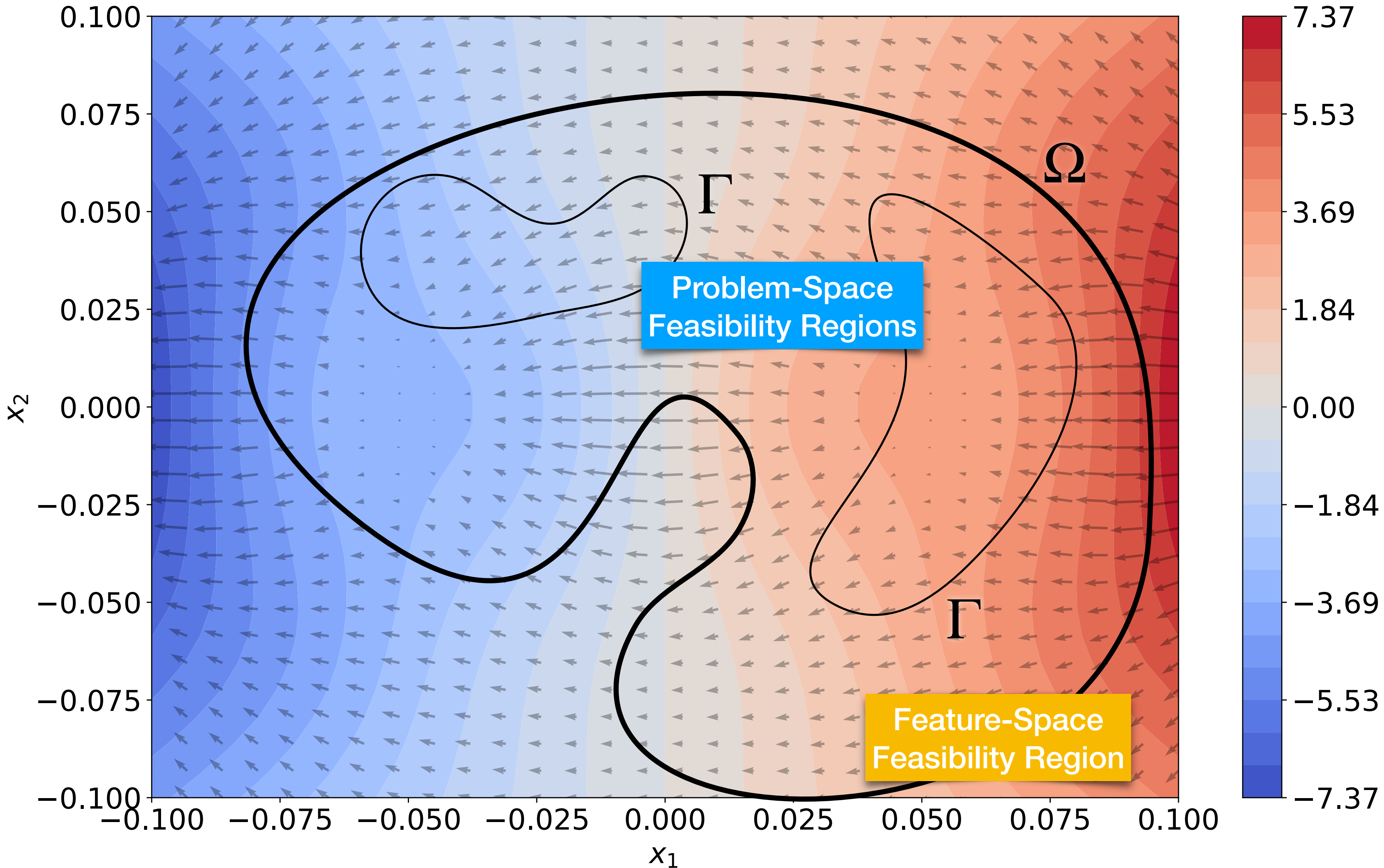
[IEEE S&P 2020] Intriguing Properties of Adversarial ML Attacks in the Problem Space

<https://s2lab.cs.ucl.ac.uk/projects/intriguing>

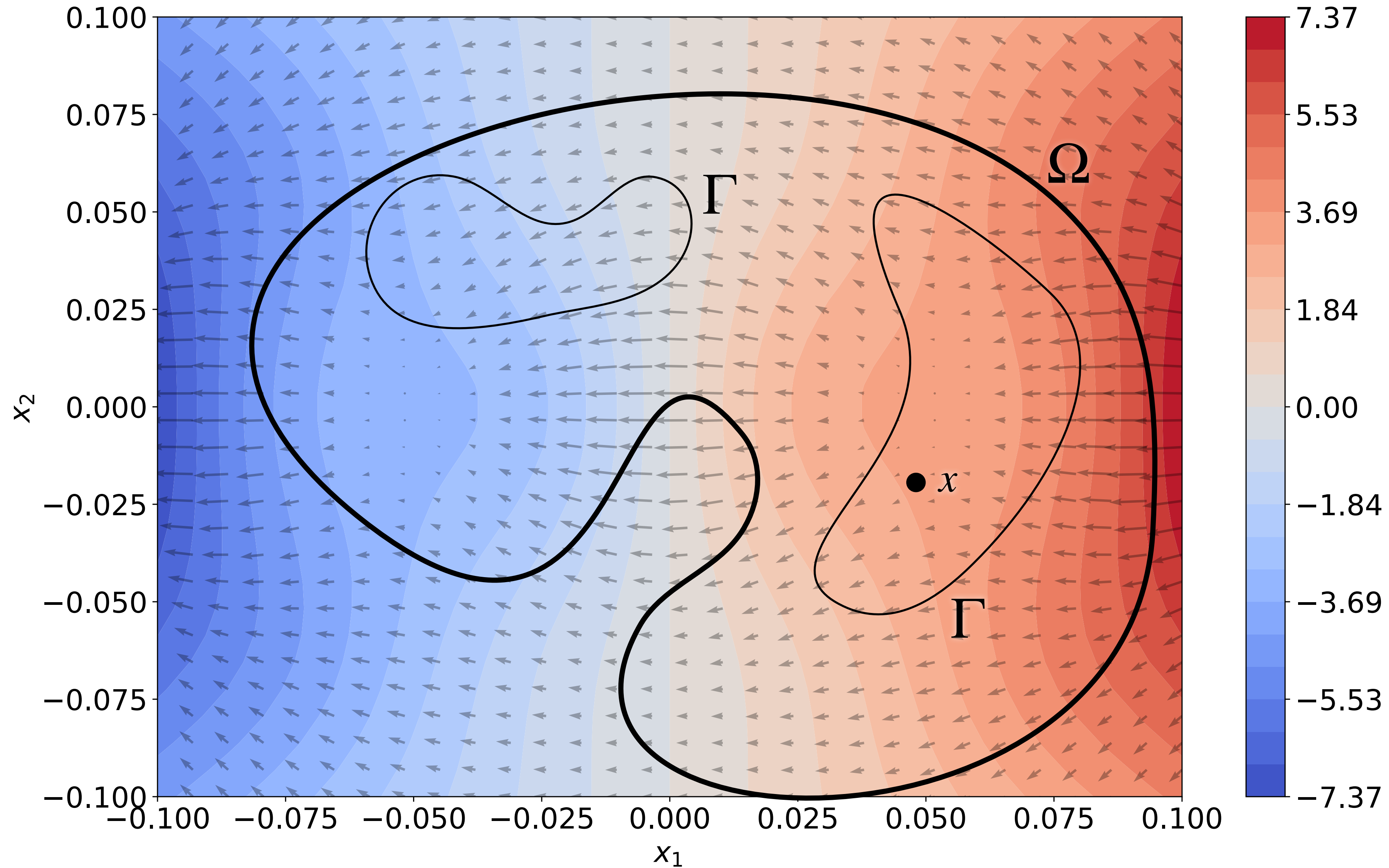
Side-effect Features: the ONE Slide to Remember



Side-effect Features: the ONE Slide to Remember



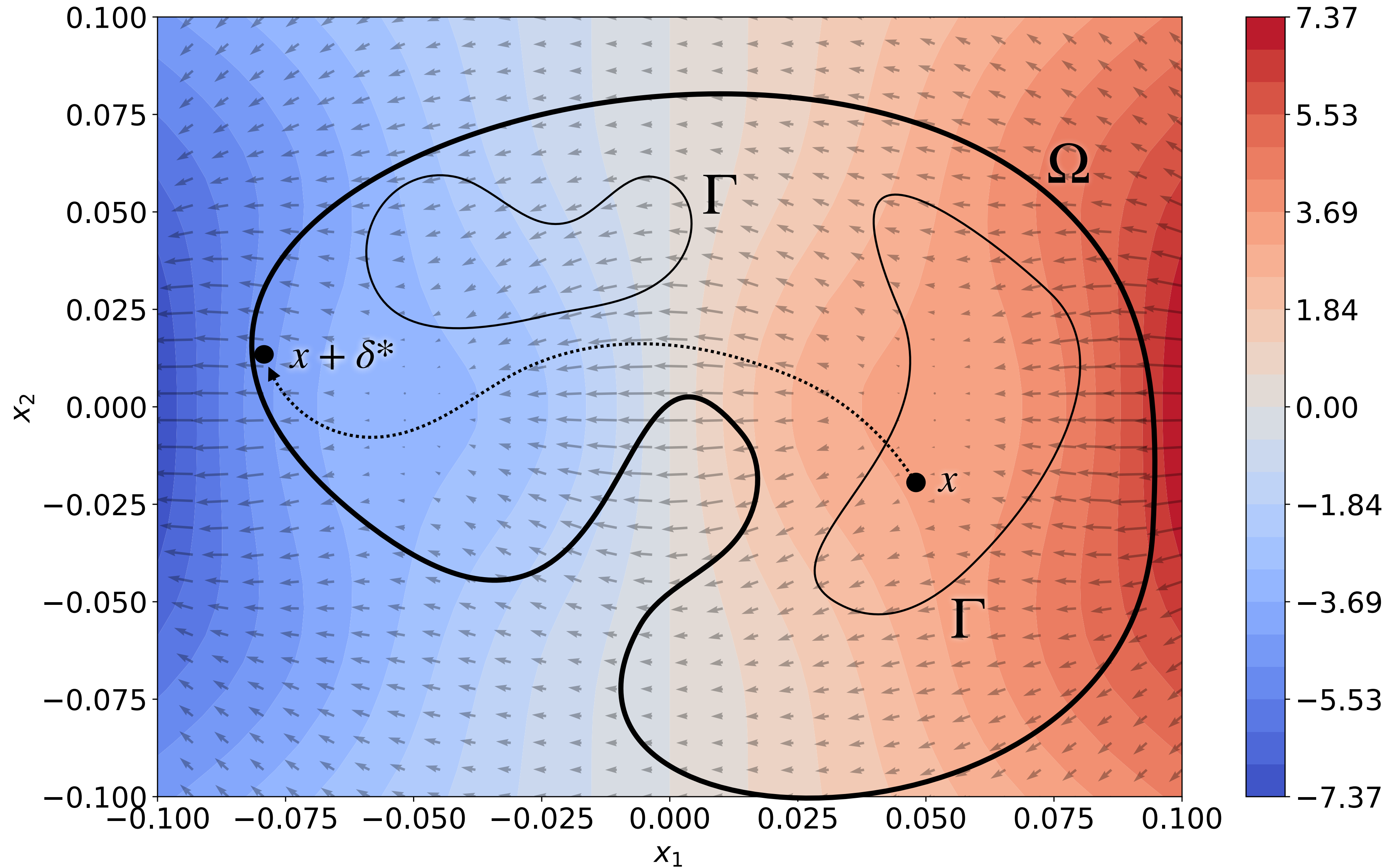
Side-effect Features: the ONE Slide to Remember



[IEEE S&P 2020] Intriguing Properties of Adversarial ML Attacks in the Problem Space

<https://s2lab.cs.ucl.ac.uk/projects/intriguing>

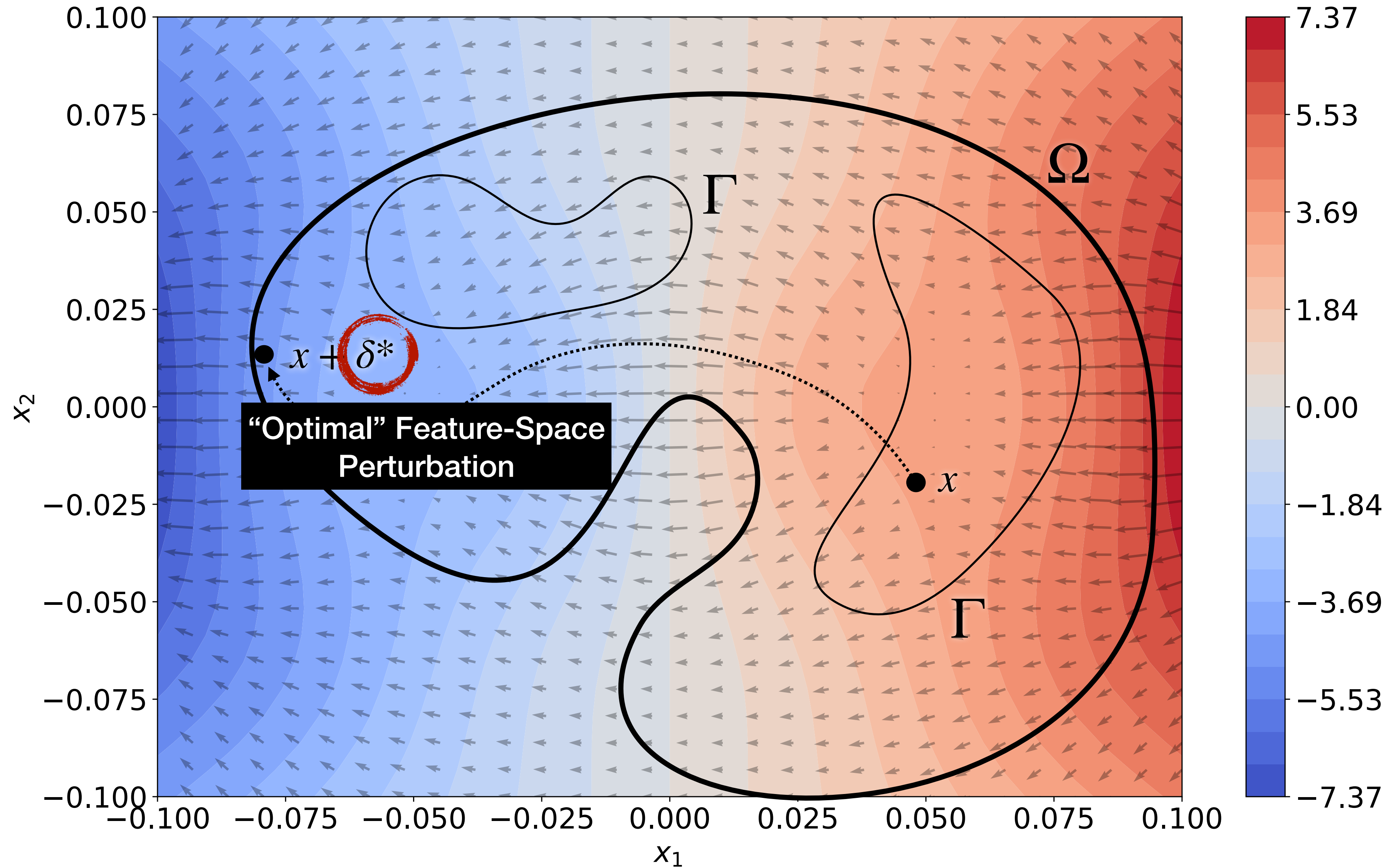
Side-effect Features: the ONE Slide to Remember



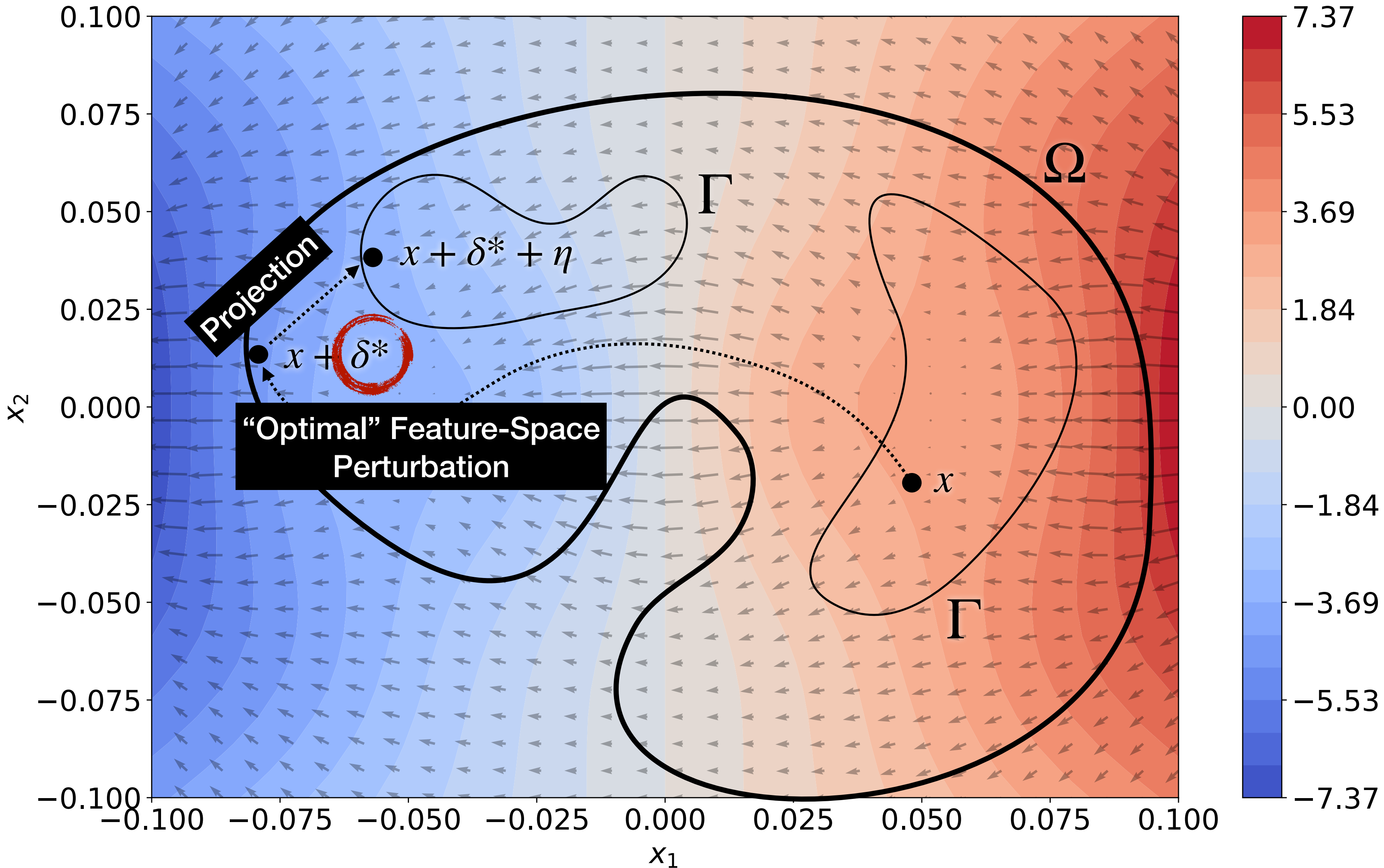
[IEEE S&P 2020] Intriguing Properties of Adversarial ML Attacks in the Problem Space

<https://s2lab.cs.ucl.ac.uk/projects/intriguing>

Side-effect Features: the ONE Slide to Remember



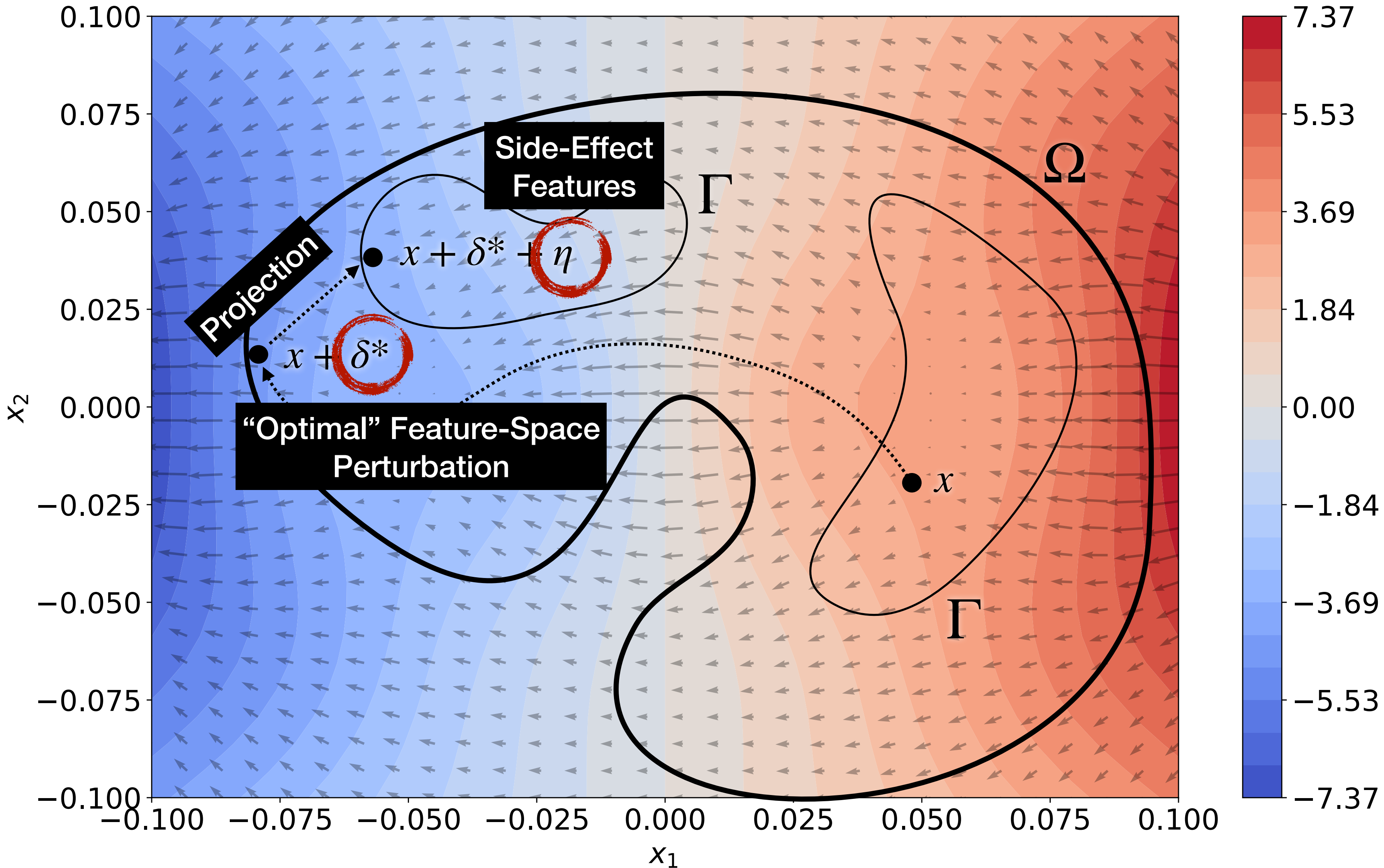
Side-effect Features: the ONE Slide to Remember



[IEEE S&P 2020] Intriguing Properties of Adversarial ML Attacks in the Problem Space

<https://s2lab.cs.ucl.ac.uk/projects/intriguing>

Side-effect Features: the ONE Slide to Remember



[IEEE S&P 2020] Intriguing Properties of Adversarial ML Attacks in the Problem Space

<https://s2lab.cs.ucl.ac.uk/projects/intriguing>

Feature Space vs. Problem Space: the OTHER Slide to Remember

$$\begin{aligned} \delta^* &= \arg \min_{\delta \in \mathbb{R}^n} f_t(\mathbf{x} + \delta) \\ \text{subject to: } &\delta \models \Omega. \end{aligned}$$

Feature-Space Constraints

- Lp perturbations
- Domain constraints for vectors

Search Strategy

- Gradient-driven

$$\begin{aligned} \operatorname{argmin}_{\mathbf{T} \in \mathcal{T}} \quad & f_t(\varphi(\mathbf{T}(z))) = f_t(\mathbf{x} + \delta^* + \eta) \\ \text{subject to: } \quad & \llbracket z \rrbracket^\tau = \llbracket \mathbf{T}(z) \rrbracket^\tau, \quad \forall \tau \in \Upsilon \\ & \pi(\mathbf{T}(z)) = 1, \quad \forall \pi \in \Pi \\ & \mathbf{A}(\mathbf{T}(z)) = \mathbf{T}(z), \quad \forall \mathbf{A} \in \Lambda \end{aligned}$$

Problem-Space Constraints

- Available Transformations
- Preserved Semantics
- Plausibility
- Robustness to Preprocessing

Search Strategy

- Gradient-driven
- Problem-driven
- Hybrid

Feature Space vs. Problem Space: the OTHER Slide to Remember

$$\delta^* = \arg \min_{\delta \in \mathbb{R}^n} f_t(\mathbf{x} + \delta)$$

subject to: $\delta \models \Omega.$

$$\operatorname{argmin}_{\mathbf{T} \in \mathcal{T}} f_t(\varphi(\mathbf{T}(z))) = f_t(\mathbf{x} + \delta^* + \eta)$$

subject to: $\llbracket z \rrbracket^\tau = \llbracket \mathbf{T}(z) \rrbracket^\tau, \quad \forall \tau \in \Upsilon$

$$\pi(\mathbf{T}(z)) = 1, \quad \forall \pi \in \Pi$$
$$\mathbf{A}(\mathbf{T}(z)) = \mathbf{T}(z), \quad \forall \mathbf{A} \in \Lambda$$





Feature-Space Constraints

- Lp perturbations
- Domain constraints for vectors

Search Strategy

- Gradient-driven

Problem-Space Constraints

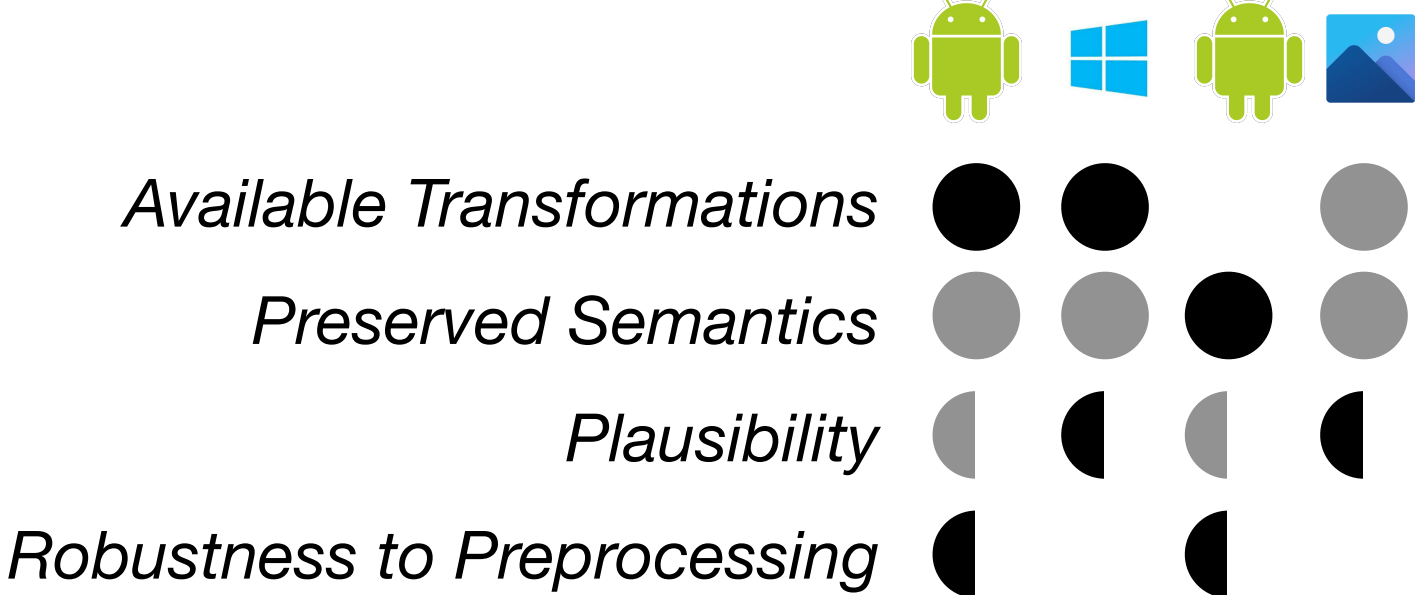
- Available Transformations 
- Preserved Semantics 
- Plausibility 
- Robustness to Preprocessing 

Search Strategy

- Gradient-driven
- Problem-driven
- Hybrid

Works for Multiple Domains

Compare existing methods & improve SoA



See Table I in §II.C

Works for Multiple Domains

Compare existing methods & improve SoA

TABLE I
PROBLEM-SPACE EVASION ATTACKS FROM PRIOR WORK ACROSS DIFFERENT SPACES, MODELED WITH OUR FORMALIZATION.

See Table I in §II.C

	Image Classification [16]	Facial Recognition [22]	Audio [17]	Text [33]	Code Attribution [31]	Java [32]	Windows [35]	Windows RNN [20]	Android Transplantation [25]	Our Android Attack (see [30])	
Knowledge θ	PK.	PK.	PK.	PK.	PK.	ZK.	PK.	ZK.	ZK.	PK.	
Feature mapping φ	Invertible: no. Differentiable: yes.	Invertible: no. Differentiable: yes.	Invertible: no. Differentiable: yes.	Invertible: no. Differentiable: yes.	Invertible: no. Differentiable: no.	Invertible: no. Differentiable: no.	Invertible: no. Differentiable: no.	Invertible: no. Differentiable: no.	Invertible: no. Differentiable: no.	Invertible: no. Differentiable: no.	
Available Transformations	None	None	None	None	None	None	None	None	None	None	
Preserved Semantics	None	None	None	None	None	None	None	None	None	None	
Plausibility	None	None	None	None	None	None	None	None	None	None	
Robustness to Preprocessing	None	None	None	None	None	None	None	None	None	None	
Problem space \mathcal{Z}	Image (pixels).	Image (pixels).	Audio (samples).	Text (words).	Code (source code).	Code (source code).	Software (binary).	Software (bytecode).	Software (bytecode).	Software (bytecode).	
Classifier g	Deep learning.	Deep learning.	Deep learning.	Classifier.	Any classifier.	SVM-RBF (Hidost [41]), RF (PDFRate [41]).	Deep learning (MalConv [52]).	RNN/LSTM variants, and transferability to traditional classifiers (e.g., RF, SVM).	kNN, DT, SVM (and VirusTotal [50]).	Linear SVM (DREBIN [5]) and its hardened version (Sec-SVM [21]).	
\mathcal{T}	(i) Addition of salt noise. $(\alpha + \delta \in [0, 1]^n)$. (ii) Pixel values must be integers from 0 to 255 (discretization problem).	(i) Addition of salt noise. $(\alpha + \delta \in [0, 1]^n)$. (ii) Pixel values must be integers from 0 to 255. (iii) Pixels are printable. (iv) Robust to 3D rotations.	(i) Addition of salt noise. values bounded (i.e., $\alpha + \delta \in [-M, +M]$). (ii) Word-level perturbations.	(i) Addition of salt noise. values bounded (i.e., $\alpha + \delta \in [-M, +M]$). (ii) No changes to the layout of the code.	(i) Addition of salt noise. values bounded (i.e., $\alpha + \delta \in [-M, +M]$). (ii) No changes to the layout of the code.	(i) Addition of salt noise. values bounded (i.e., $\alpha + \delta \in [-M, +M]$). (ii) No changes to the layout of the code.	(i) Addition of salt noise. values bounded (i.e., $\alpha + \delta \in [-M, +M]$). (ii) No changes to the layout of the code.	(i) Addition of salt noise. values bounded (i.e., $\alpha + \delta \in [-M, +M]$. (ii) No changes to the layout of the code.	(i) Addition of salt noise. values bounded (i.e., $\alpha + \delta \in [-M, +M]$. (ii) No changes to the layout of the code.	(i) Addition of salt noise. values bounded (i.e., $\alpha + \delta \in [-M, +M]$. (ii) No changes to the layout of the code.	
Preserved Semantics Υ	An image should not trivially become an image of another class, so perturbation is constrained $\ \delta\ _p \leq \delta_{max}$.	Human subjects retain their original identity and their recognizability to other humans (compared to using full face masks, disguises, etc).	Semantics of original audio preserved by constraining the perturbation $(\ \delta\ _p \leq \delta_{max})$.	Sentence meaning preserved by (i) replacing like characters (ii) using the GloVe model [36] (not syntactically similar words).	Source code semantics preserved by construction through use of semantics-preserving transformations.	Malicious semantics preserved by construction through use of AST-based transplantation.	Malicious network functionality is still present (verification with Cuckoo Sandbox).	Malicious code is unaffected by only appending redundant bytes.	API sequences and function return values are unchanged (verification with Cuckoo Monitor).	Malicious semantics preserved, tested by installing and executing each application.	Malicious semantics preserved by construction with opaque predicates (newly inserted code is not executed at runtime).
Robustness to Preprocessing Λ	None explicitly considered.	Discussed but not robust to: the use of specific illumination or distance of the camera.	Robust to: (i) Addition of pointwise random noise (ii) MP3 compression. Discussed but not robust to: Over-the-air playing.	Not explicitly considered.	Robust to: removal of layout features (i.e., use of tabs vs spaces) which are trivial to alter.	Robust to: removal of name inconsistencies of functions and variables.	Discussed but not robust to: removal of spurious features such as presence or absence of font objects (discovered post-attack).	Discussed but not robust to: removal of redundant (non-text) bytes.	Robust to: removal of redundant code, undeclared variables, unlinked resources, undefined references, name conflicts.	Not explicitly considered.	Robust to: removal of redundant code, undeclared variables, undefined references, name conflicts, no-op instructions.
Plausibility Π	Perturbation constrained $(\ \delta\ _p \leq \delta_{max})$ to ensure the changes are imperceptible to a human.	(i) Perturbation constrained $(\ \delta\ _p \leq \delta_{max})$. (ii) Smooth pixel transitions so the eyeglass frames look legitimate with plausible deniability.	Perturbation constrained $(\ \delta\ _p \leq \delta_{max})$, so that added noise resembles white background noise largely imperceptible to a human.	(i) Ensure short distance (e.g., edit distance) of modifications (ii) User study to verify plausibility.	The code does not look suspicious and seems written by a human (survey with developers).	By construction through automated AST transplantation (although plausibility is inhibited if certain objects are used, e.g., obsolete ActiveX components).	PDFs can still be parsed and opened by a reader.	None explicitly considered.	The added no-op API calls do not raise errors.	Code is realistic by construction through automated software transplantation.	(i) Code is realistic by construction through use of automated software transplantation. (ii) Mutated apps install and start on an emulator.
Search Strategy	Gradient-driven. Stochastic Gradient Descent in the feature space.	Gradient-driven. Stochastic Gradient Descent in the feature space.	Gradient-driven. Adam optimizer with learning rate 10 and 5,000 max iterations.	Hybrid (PK). Gradients used to choose 'top' words. Problem-driven (ZK). Without gradients, importance of words is estimated by scoring without each word.	Problem-driven. New Monte-Carlo Search algorithm, applied to the problem space.	Problem-driven. Search of isomorphic sub-AST graphs in benign samples that are equivalent to malicious sub-ASTs.	Problem-driven. Genetic Programming.	Gradient-driven. Although the feature mapping is not invertible and not differentiable, the authors devise an algorithm to project byte padding on to the negative gradient.	Hybrid. Greedy algorithm selects API calls in order to minimize difference between current and previous iterations w.r.t. the direction of the Jacobian.	Gradient-driven. Prioritizing mutations that affect features typical of malware evolution (e.g., phylogenetic trees) and those present in both malware and goodware.	Gradient-driven. We use an approximate inverse of the feature mapping, and then a greedy algorithm in the problem space to follow the negative gradient.
Side-effect features η	$\eta = 0$	$\eta = 0$	$\eta = 0$	$\eta = 0$	$\eta \approx 0$	$\eta \neq 0$	$\eta \approx 0$	$\eta = 0$	$\eta \approx 0$	$\eta \neq 0$	$\eta \neq 0$

Works for Multiple Domains

Compare existing methods & improve SoA

See Table I in §II.C

TABLE I
PROBLEM-SPACE EVASION ATTACKS FROM PRIOR WORK ACROSS DIFFERENT DOMAINS, MODELED WITH OUR FORMALIZATION.

	Image Classification [16]	Facial Recognition [22]	Audio [17]	Text [53]	Code Attribution [51]	Windows [55]	Windows RNN [20]	Android Transplantation [25]	Our Android Attack (see [11])			
Knowledge θ	PK.	PK.	PK.	PK.	PK.	ZK.	ZK.	ZK.	PK.			
Feature mapping ψ	Invertible: no. Differentiable: yes.	Invertible: no. Differentiable: yes.	Invertible: no. Differentiable: yes.	Invertible: no. Differentiable: yes.	Invertible: no. Differentiable: no.	Invertible: no. Differentiable: no.	Invertible: no. Differentiable: no.	Invertible: no. Differentiable: no.	Invertible: no. Differentiable: no.			
Available Transformations	None	None	None	None	Syntactic and lexical features.	Static syntactic, based on AST, PDG, CFG.	Static (metadata, object keywords and properties, structural).	Feature mapping of MalConv [57]	Dynamic API sequences, static printable strings (also in latent feature space).	Static analysis (RTLD model [25]).	Lightweight static analysis (binary features).	
Preserved Semantics	None	None	None	None	None	Software (source code).	Software (source code).	PDF.	Software (binary).	Software (bytecode).	Software (bytecode).	
Plausibility	None	None	None	None	None	Classifier.	Any classifier.	SVM-RBF (Hidost [21]), RF (PDFRaz [21]).	Deep learning (MalConv [57]).	RNN/LSTM variants, and transferability to traditional classifiers (e.g., RF, SVM).	kNN, DT, SVM (and VirusTotal [58]).	Linear SVM (DREBIN [5]) and its hardened version (Sec-SVM [21]).
Robustness to Preprocessing	(i) Additive noise. (ii) Pixel values must be integers from 0 to 255 (discretization problem).	(i) Additive noise. (ii) Pixel values must be integers from 0 to 255 (discretization problem). (iv) Robust to 3D rotations.	(i) Additive noise. (ii) Word-level perturbations.	(i) Additive noise. (ii) Word-level perturbations.	(i) Additive noise. (ii) Word-level perturbations.	(i) Additive noise. (ii) Word-level perturbations.	(i) Additive noise. (ii) Word-level perturbations.	(i) Additive noise. (ii) Word-level perturbations.	(i) Additive noise. (ii) Word-level perturbations.	(i) Additive noise. (ii) Word-level perturbations.	(i) Additive noise. (ii) Word-level perturbations.	(i) Additive noise. (ii) Word-level perturbations.
Preserved Semantics Υ	An image should not trivially become an image of another class, so perturbation is constrained $\ \delta \ _p \leq \delta_{max}$.	Human subjects retain their original identity and their recognizability to other humans (compared to using full face masks, disguises, etc.).	Semantics of original audio preserved by constraining the perturbation $\ \delta \ _p \leq \delta_{max}$.	Sentence meaning preserved by replacing like words (e.g., using a thesaurus).	Source code semantics preserved by construction through use of semantics-preserving transformations.	Malicious semantics preserved by construction through use of AST-based transplantation.	Malicious network functionality is still present (verification with Cuckoo).	Malicious code is unaffected by only appending redundant bytes.	API sequences and function return values are unchanged (verification with Cuckoo Monitor).	Malicious semantics preserved, tested by installing and executing on application.	Malicious semantics preserved by construction with Cuckoo.	
Robustness	None explicitly considered.	Discussed but not robust to: the use of specific illumination or distance of the camera.	Robust to: removal of layout features (i.e., use of tabs vs spaces) which are trivial to alter.	Robust to: removal of name inconsistencies of functions and variables.	Robust to: removal of name inconsistencies of functions and variables.	Robust to: removal of name inconsistencies of functions and variables.	Robust to: removal of name inconsistencies of functions and variables.	Robust to: removal of name inconsistencies of functions and variables.	Robust to: removal of name inconsistencies of functions and variables.	Robust to: removal of name inconsistencies of functions and variables.	Robust to: removal of name inconsistencies of functions and variables.	
Search Strategy	Gradient-driven. Stochastic Gradient Descent (SGD).	Gradient-driven. Stochastic Gradient Descent (SGD).	Gradient-driven. Adam optimizer with learning rate 10 and 5,000 max iterations.	Hybrid (PK). Gradients used to choose 'top' words. Problem-driven (ZK). Without gradients, importance of words is estimated by scoring without each word.	Problem-driven. Genetic Programming.	Gradient-driven. Although the feature mapping is not invertible and not differentiable, the authors devise an algorithm to project byte padding on to the negative gradient.	Gradient-driven. We use an approximate inverse of the feature mapping, and then a greedy algorithm in the problem space to follow the negative gradient.					
Side-effect features η	$\eta = 0$	$\eta = 0$	$\eta = 0$	$\eta = 0$	$\eta = 0$	$\eta = 0$	$\eta = 0$	$\eta = 0$	$\eta = 0$	$\eta = 0$	$\eta = 0$	



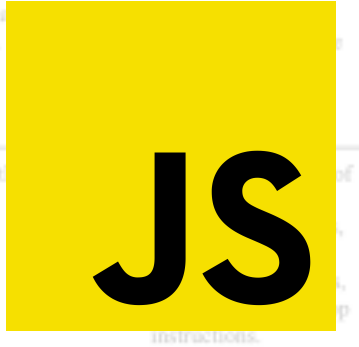
Facial Recognition



Speech Recognition



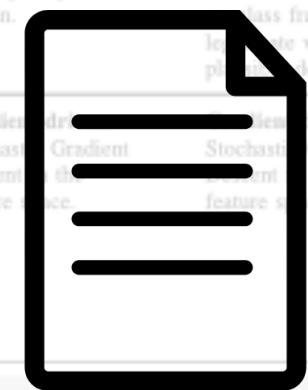
Android Malware



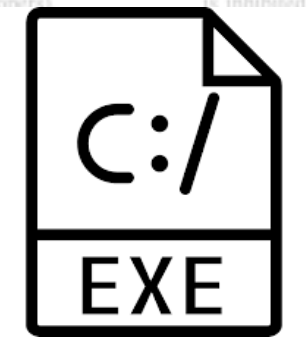
Malicious Javascript



Image Classification



Text Classification



Windows Malware



Code Attribution



PDF Malware



How hard could it be?

- Despite hype on adversarial learning: No suitable work for us 😞
- Two tricky challenges
 - No inverse map from topic space back to problem space
 - Unobtrusive changes lead to side effects in the feature space



Slide 13

Machine Learning
and Security



19:31 / 51:04



Konrad Rieck: When Papers Choose their Reviewers: Adversarial ML in Conference Management Systems.



ViSP

34 subscribers

Subscribe



Like



Share



Download



Clip



How hard could it be?

- Despite hype on adversarial learning: No suitable work for us 😞
- Two tricky challenges
 - No inverse map from topic space back to problem space
 - Unobtrusive changes lead to side effects in the feature space



[USENIX Sec 2023] Eisenhofer et al. No more
Reviewer #2: Subverting Automatic Paper-
Reviewer Assignment using Adversarial
Learning

Slide 13

Machine Learning
and Security

Konrad Rieck: When Papers Choose their Reviewers: Adversarial ML in Conference Management Systems.



ViSP

34 subscribers

Subscribe

Like



Share

Download

Clip



How hard could it be?

- Despite hype on adversarial learning: No suitable work for us 😞
- Two tricky challenges
 - No inverse map from topic space back to problem space
 - Unobtrusive changes lead to side effects in the feature space



Slide 13

Machine Learning
and Security

[USENIX Sec 2023] Eisenhofer et al. No more
Reviewer #2: Subverting Automatic Paper-
Reviewer Assignment using Adversarial
Learning

Problem space reasoning applies and affects (backdoor) poisoning attacks too!

[IEEE S&P 2023] Yang et al. Jigsaw Puzzle: Selective Backdoor Attack to Subvert Malware Classifiers

Android Attack: Experiments

Results: What is the impact on the Feature Space?

Results: What is the impact on the Feature Space?

- **Dataset:** ~170K Android apps (10% malware) from Jan 2017 to Dec 2018

Results: What is the impact on the Feature Space?

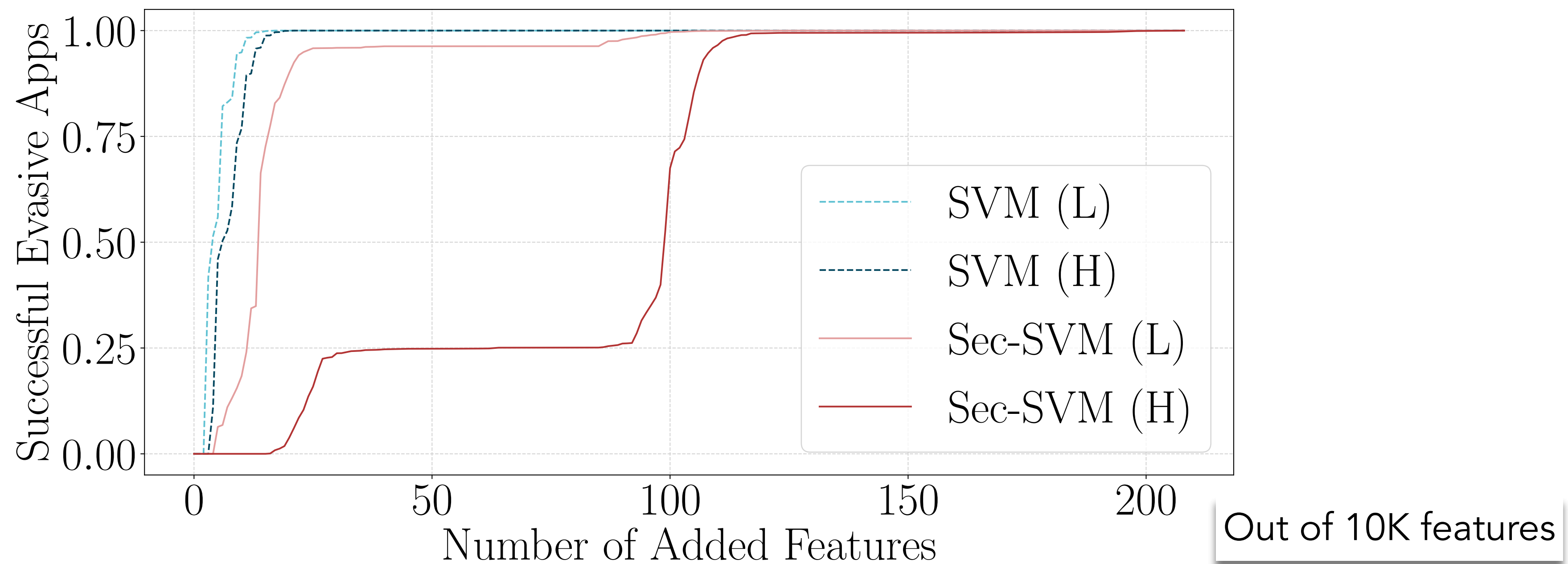
- **Dataset:** ~170K Android apps (10% malware) from Jan 2017 to Dec 2018
- **DREBIN** [NDSS'14]: Linear SVM, binary feature space
- **Sec-SVM** [TDSC'17]: Feature-space defense for DREBIN (evenly distributes weights)

Results: What is the impact on the Feature Space?

- **Dataset:** ~170K Android apps (10% malware) from Jan 2017 to Dec 2018
- **DREBIN** [NDSS'14]: Linear SVM, binary feature space
- **Sec-SVM** [TDSC'17]: Feature-space defense for DREBIN (evenly distributes weights)
- **Low** (L) vs **High** (H) **confidence**: cross decision boundary or cross into Q1 of benign

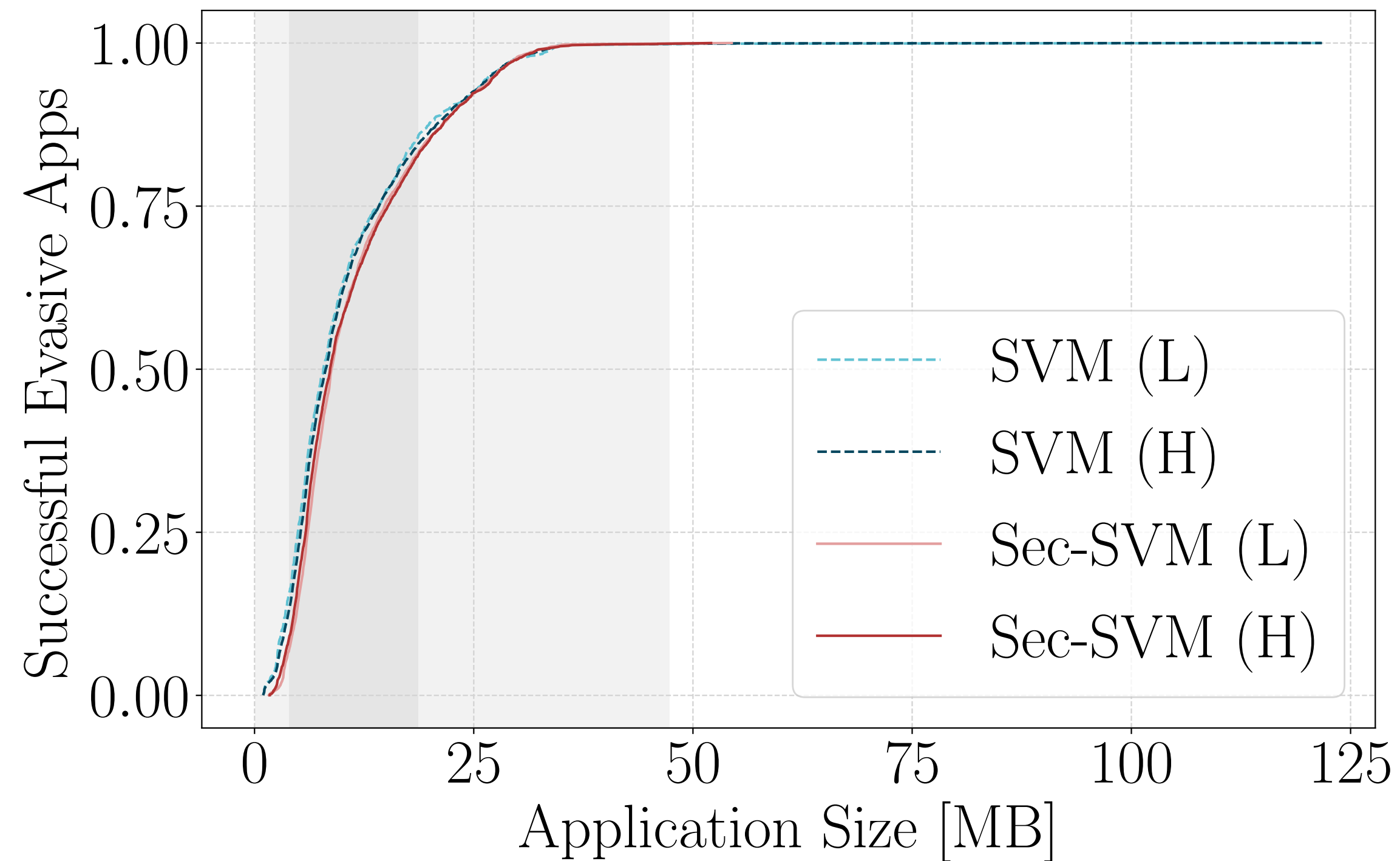
Results: What is the impact on the Feature Space?

- **Dataset:** ~170K Android apps (10% malware) from Jan 2017 to Dec 2018
- **DREBIN** [NDSS'14]: Linear SVM, binary feature space
- **Sec-SVM** [TDSC'17]: Feature-space defense for DREBIN (evenly distributes weights)
- **Low (L) vs High (H) confidence:** cross decision boundary or cross into Q1 of benign



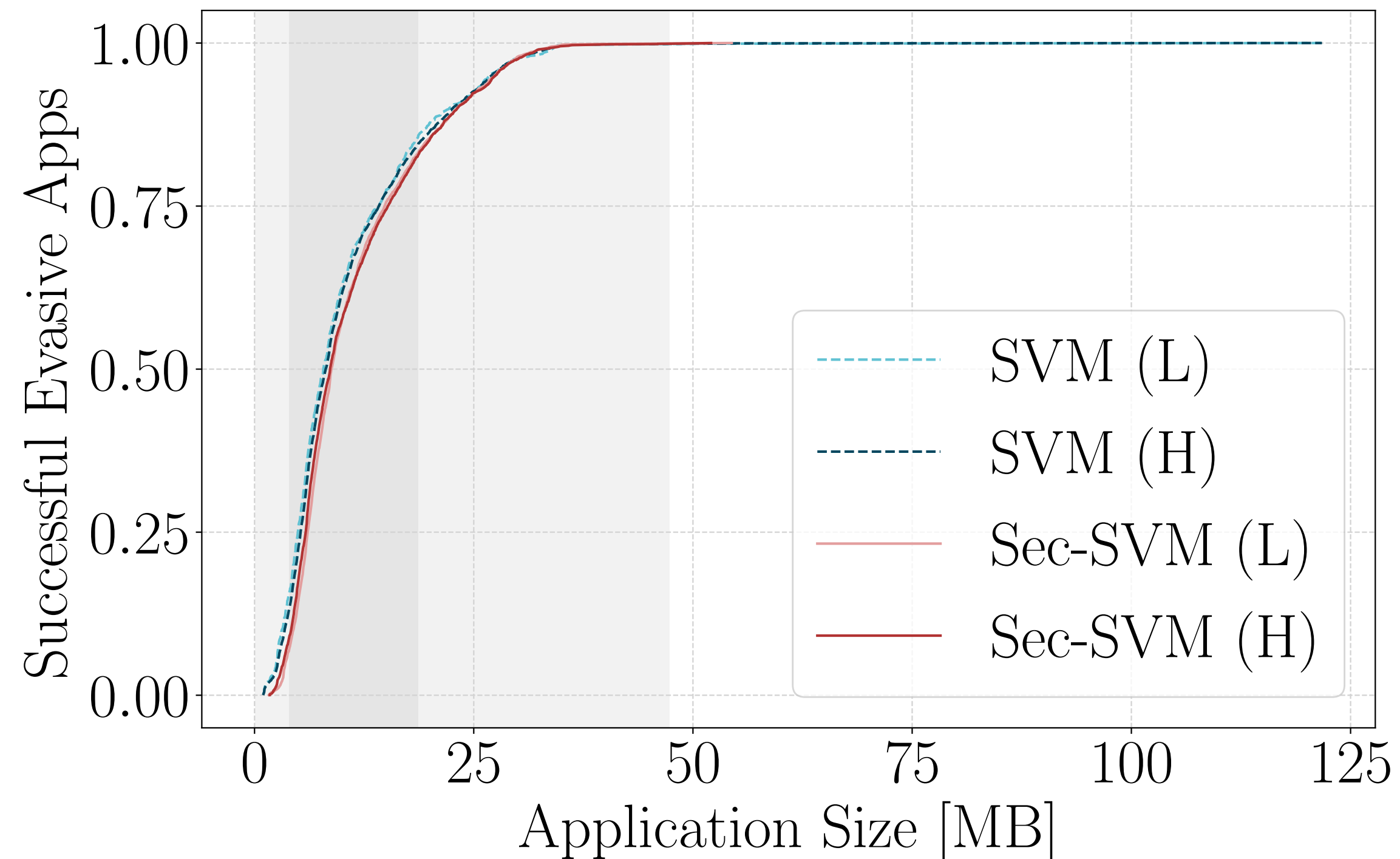
Results: What is the impact on the Problem Space?

Results: What is the impact on the Problem Space?



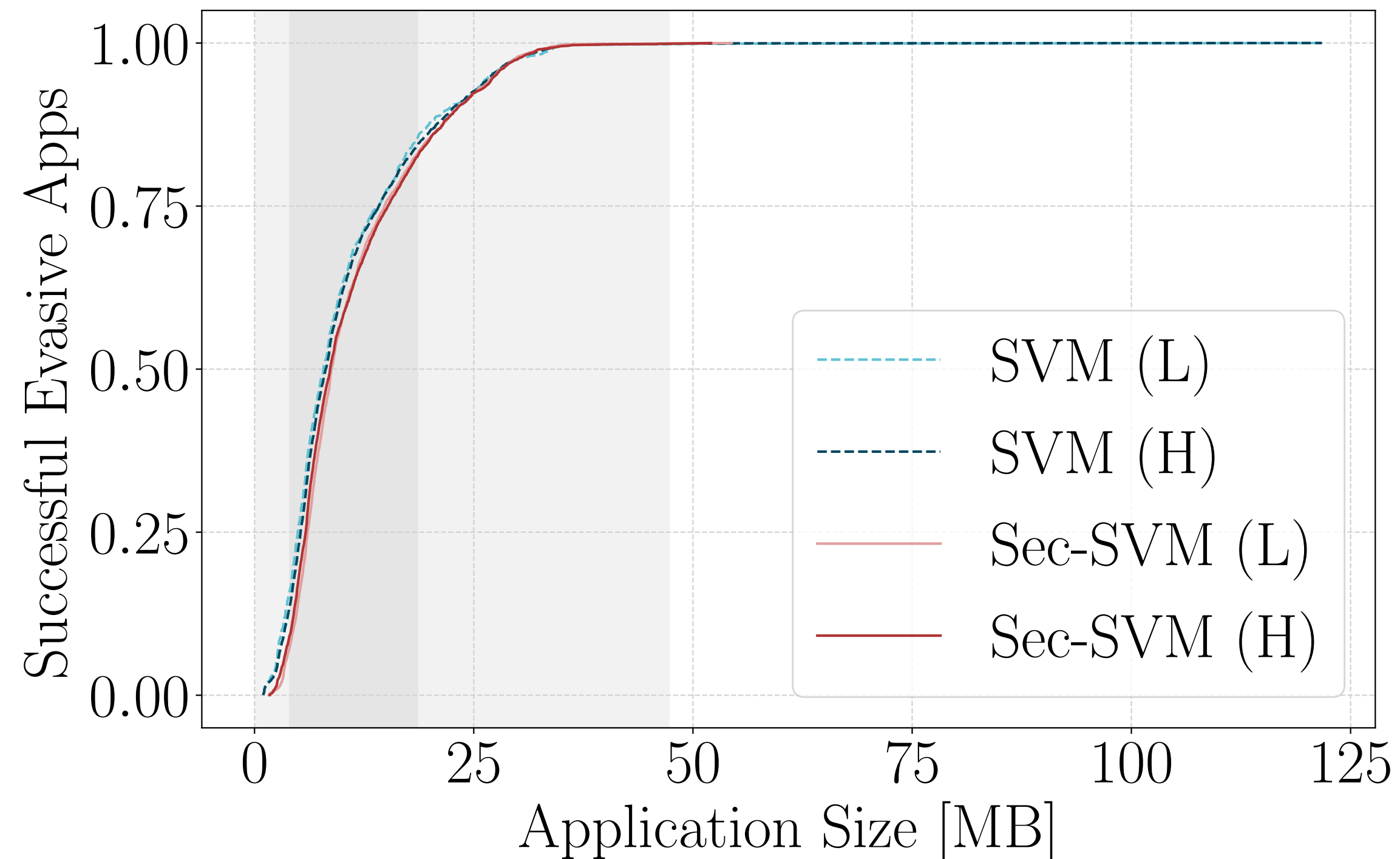
Results: What is the impact on the Problem Space?

- Adversarial generation < 2 minutes per app



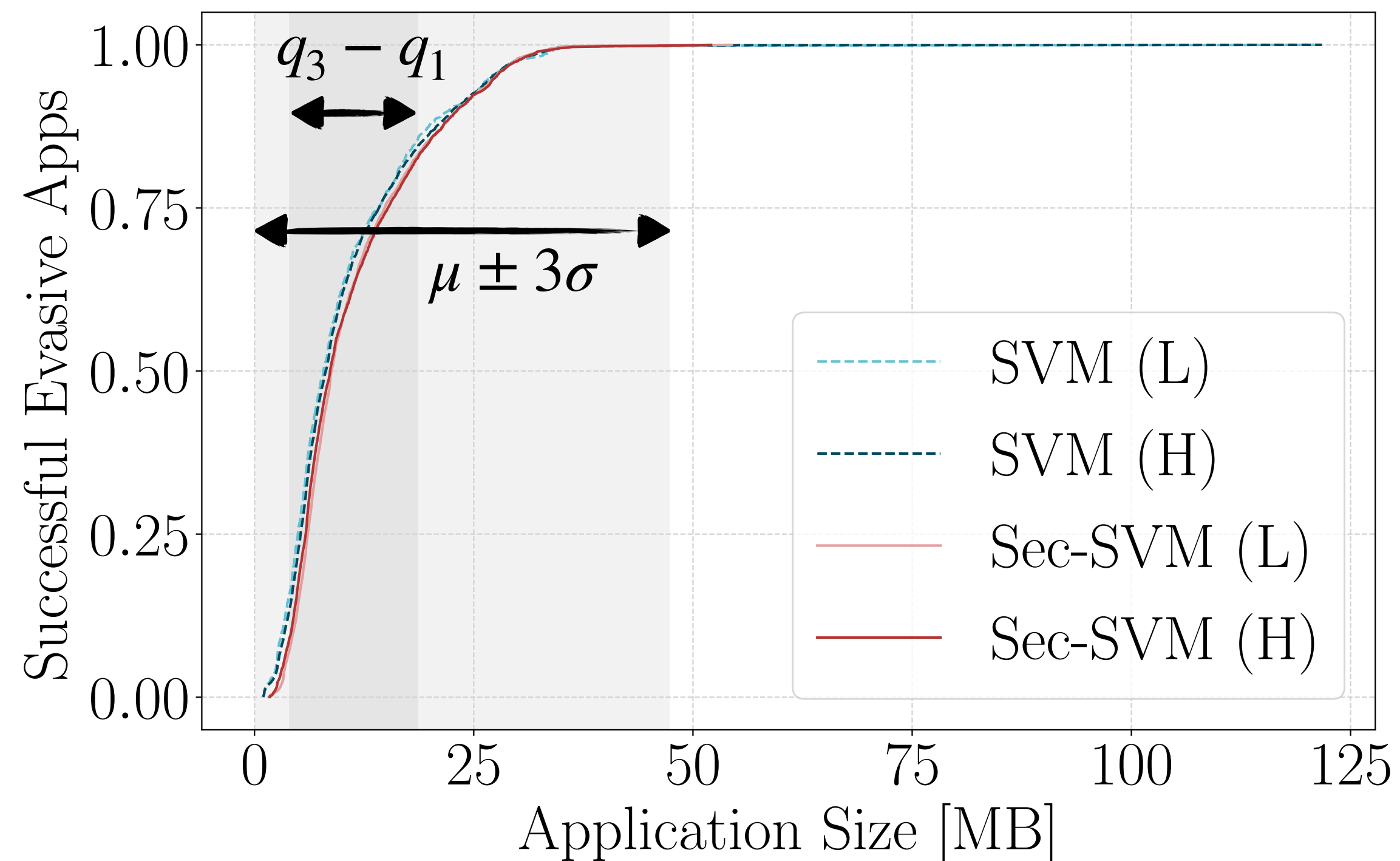
Results: What is the impact on the Problem Space?

- Adversarial generation < 2 minutes per app
- Restricting feature-space perturbations δ does not hinder problem-space attack



Results: What is the impact on the Problem Space?

- Adversarial generation < 2 minutes per app
- Restricting feature-space perturbations δ does not hinder problem-space attack
- App statistics (e.g., size) do not become anomalous after injection (other stats in the paper)



Problem Space Adversarial Training

Problem vs Feature Space Adversarial Training

[IEEE S&P 2020] Intriguing Properties of Adversarial ML Attacks in the Problem Space

<https://s2lab.cs.ucl.ac.uk/projects/intriguing>

[IEEE S&P 2023] Dyrmishi et al. On The Empirical Effectiveness of Unrealistic Adversarial Hardening Against Realistic Adversarial Attacks

Problem vs Feature Space Adversarial Training

- Exciting work [IEEE S&P 2023] on Text, Botnet Traffic, Windows Malware Classification Tasks
 - Text: Problem Space AT 16.94% more effective than Feature Space AT
 - Botnet Traffic: Problem Space AT robustness \approx Feature Space AT
 - Windows Malware: Problems Space AT outperforms Feature Space AT robustness

[IEEE S&P 2020] **Intriguing Properties of Adversarial ML Attacks in the Problem Space**

<https://s2lab.cs.ucl.ac.uk/projects/intriguing>

[IEEE S&P 2023] Dyrmishi et al. **On The Empirical Effectiveness of Unrealistic Adversarial Hardening Against Realistic Adversarial Attacks**

Problem vs Feature Space Adversarial Training

- Exciting work [IEEE S&P 2023] on Text, Botnet Traffic, Windows Malware Classification Tasks
 - **(Marginal)** Text: Problem Space AT 16.94% more effective than Feature Space AT
 - **(Not Required)** Botnet Traffic: Problem Space AT robustness \sim Feature Space AT
 - **(Required)** Windows Malware: Problem Space AT $>$ Feature Space AT robustness
- It may seem a task-dependent result... let's fix some variables

[IEEE S&P 2020] **Intriguing Properties of Adversarial ML Attacks in the Problem Space**

<https://s2lab.cs.ucl.ac.uk/projects/intriguing>

[IEEE S&P 2023] Dyrmishi et al. **On The Empirical Effectiveness of Unrealistic Adversarial Hardening Against Realistic Adversarial Attacks**

Problem vs Feature Space Adversarial Training

- Exciting work [IEEE S&P 2023] on Text, Botnet Traffic, Windows Malware Classification Tasks
 - **(Marginal)** Text: Problem Space AT 16.94% more effective than Feature Space AT
 - **(Not Required)** Botnet Traffic: Problem Space AT robustness \approx Feature Space AT
 - **(Required)** Windows Malware: Problem Space AT $>$ Feature Space AT robustness
- It may seem a task-dependent result... let's fix some variables
- **Dataset:** ~170K Android apps (10% malware) from Jan 2017 to Dec 2018
- **DREBIN** [NDSS'14]: Linear SVM, binary feature space
- **Sec-SVM** [TDSC'17]: Feature-space defense for DREBIN (evenly distributes weights)
- Increasing level of perturbation budget and [IEEE S&P 2020] Android attack

[IEEE S&P 2020] **Intriguing Properties of Adversarial ML Attacks in the Problem Space**

<https://s2lab.cs.ucl.ac.uk/projects/intriguing>

[IEEE S&P 2023] Dyrmishi et al. **On The Empirical Effectiveness of Unrealistic Adversarial Hardening Against Realistic Adversarial Attacks**

Problem vs Feature Space Adversarial Training (SVM)

- **Dataset:** ~170K Android apps (10% malware) from Jan 2017 to Dec 2018
- **DREBIN** [NDSS'14]: Linear SVM, binary feature space
- **Sec-SVM** [TDSC'17]: Feature-space defense for DREBIN (evenly distributes weights)
- Increasing level of perturbation budget and [IEEE S&P 2020] Android attack

SVM Adversarial Training

Budget (%)	FS ASR	PS ASR	FS F1 Clean	PS F1 Clean
10	0.22	0.02	0.9	0.9
20	0.14	0.01	0.9	0.9
30	0.12	0	0.9	0.9
40	0.1	0	0.9	0.9

Problem vs Feature Space Adversarial Training (SecSVM)

- **Dataset:** ~170K Android apps (10% malware) from Jan 2017 to Dec 2018
- **DREBIN** [NDSS'14]: Linear SVM, binary feature space
- **Sec-SVM** [TDSC'17]: Feature-space defense for DREBIN (evenly distributes weights)
- Increasing level of perturbation budget and [IEEE S&P 2020] Android attack

SecSVM Adversarial Training

Budget (%)	FS ASR	PS ASR	FS F1 Clean	PS F1 Clean
10	0.95	0.17	0.81	0.81
20	0.91	0.12	0.81	0.81
30	0.92	0.1	0.81	0.81
40	0.89	0.07	0.81	0.81

Problem vs Feature Space Adversarial Training (SecSVM)

SecSVM Adversarial Training

Budget (%)	FS ASR	PS ASR	FS F1 Clean	PS F1 Clean
10	0.95	0.17	0.81	0.81
20	0.91	0.12	0.81	0.81
30	0.92	0.1	0.81	0.81
40	0.89	0.07	0.81	0.81

SVM Adversarial Training

Budget (%)	FS ASR	PS ASR	FS F1 Clean	PS F1 Clean
10	0.22	0.02	0.9	0.9
20	0.14	0.01	0.9	0.9
30	0.12	0	0.9	0.9
40	0.1	0	0.9	0.9

Problem vs Feature Space Adversarial Training

SecSVM Adversarial Training

Budget (%)	FS ASR	PS ASR	FS F1 Clean	PS F1 Clean
10	0.95	0.17	0.81	0.81
20	0.91	0.12	0.81	0.81
30	0.89	0.07	0.81	0.81
40	0.89	0.07	0.81	0.81

Perhaps not task-dependent but affected by

- Program abstractions
- Feature representations
- ML models

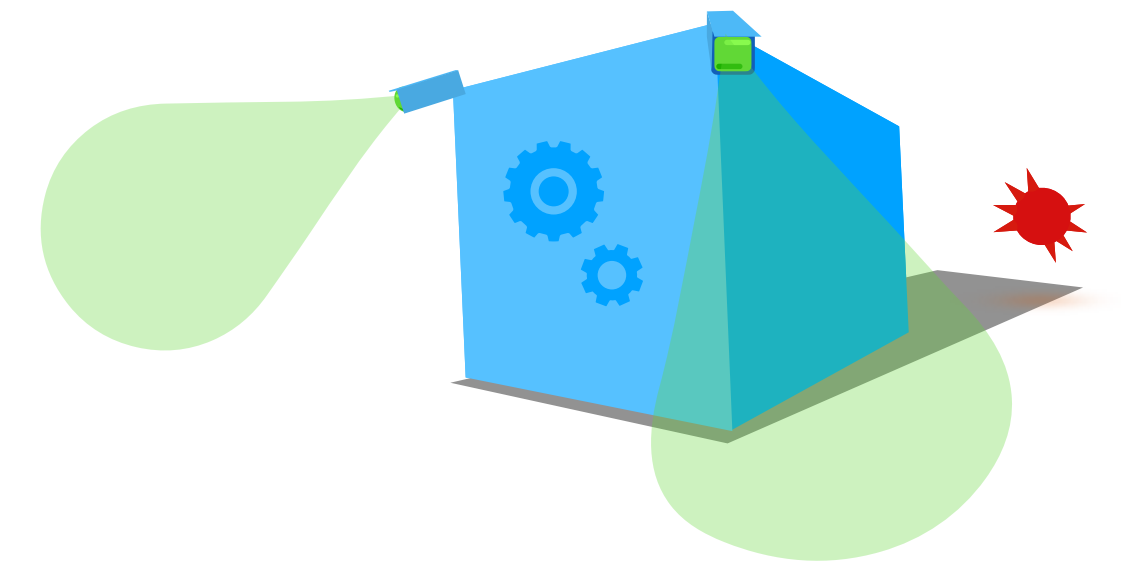
Budget (%)	FS ASR	PS ASR	FS F1 Clean	PS F1 Clean
5	0.22	0.02	0.9	0.9
20	0.14	0.01	0.9	0.9
30	0.12	0	0.9	0.9
40	0.1	0	0.9	0.9

Outline

Focus

Adversarial ML evasion attacks against malware classifiers

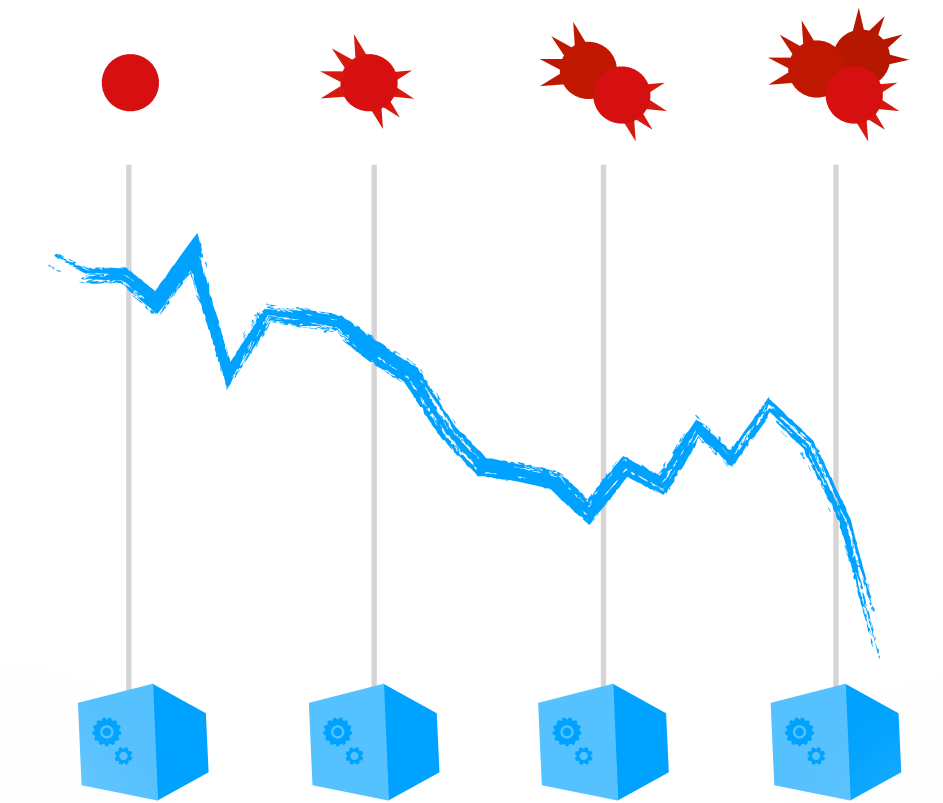
- Classic formulation of evasion attacks is ill-suited for reasoning about realizable evasive malware
- By reformulating, we can propose stronger attacks and easily compare against alternatives



Bigger Picture

Drifting scenarios caused by threats evolving over time

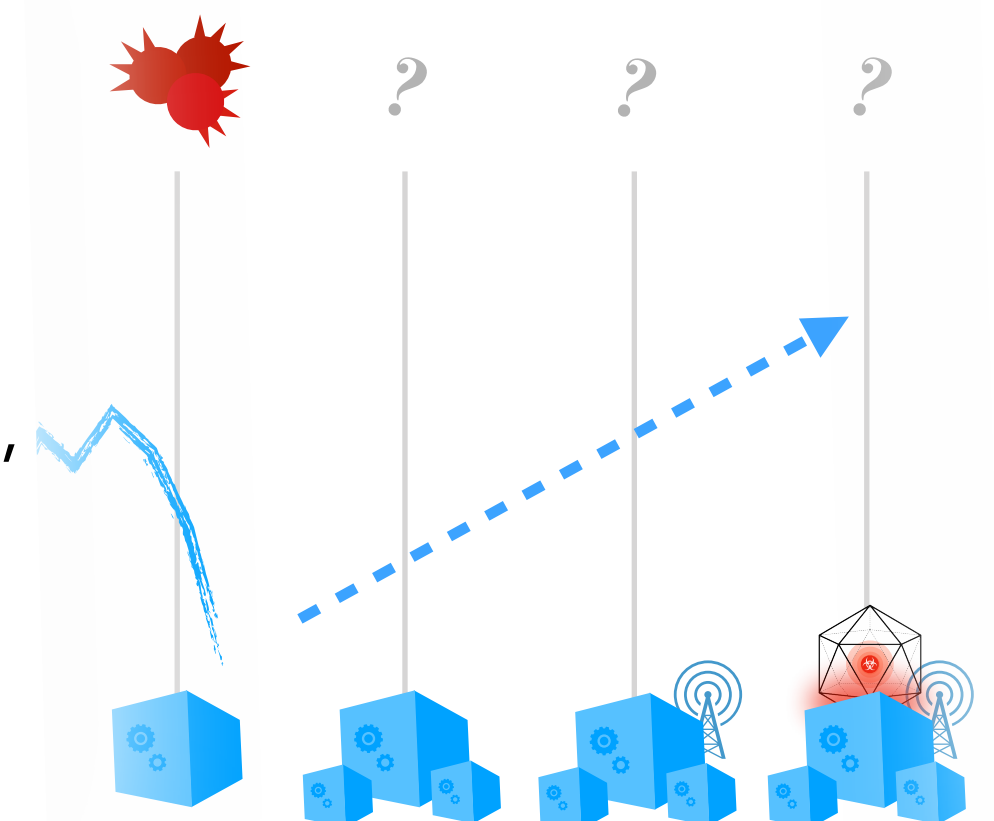
- How concept drift affects machine learning-based detectors in security settings
- The need for time-aware evaluations and metrics



Looking Ahead

Quo vadis?

- Discussion of the future of trustworthy ML for system security
- Robust feature development, universal adversarial perturbations, realizable backdoors, drift forecasting, and the role of abstractions towards the Platonic ideal of interesting behaviors

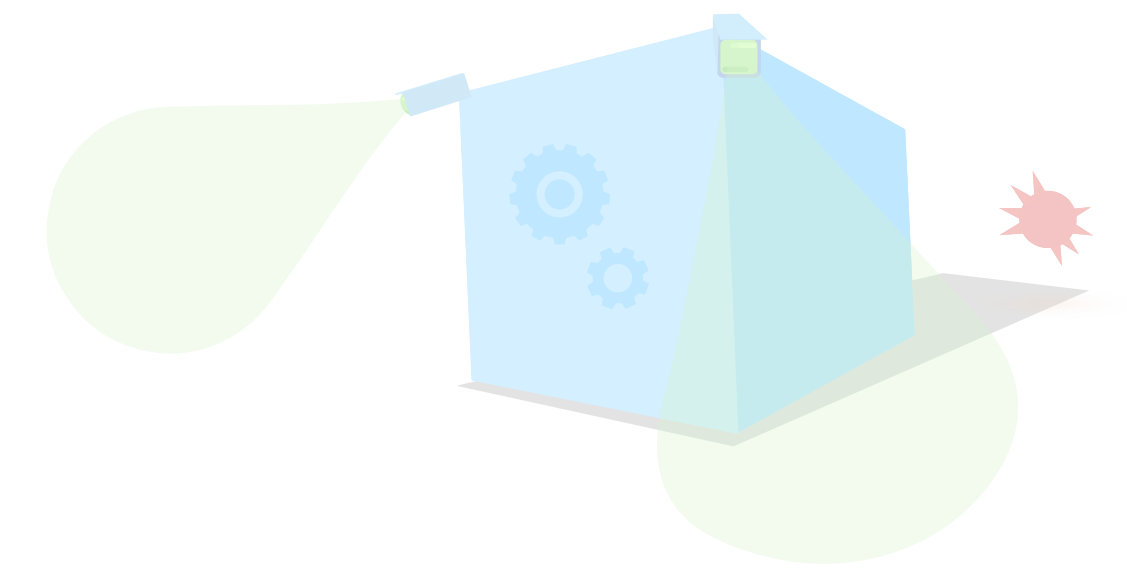


Outline

Focus

Adversarial ML evasion attacks against malware classifiers

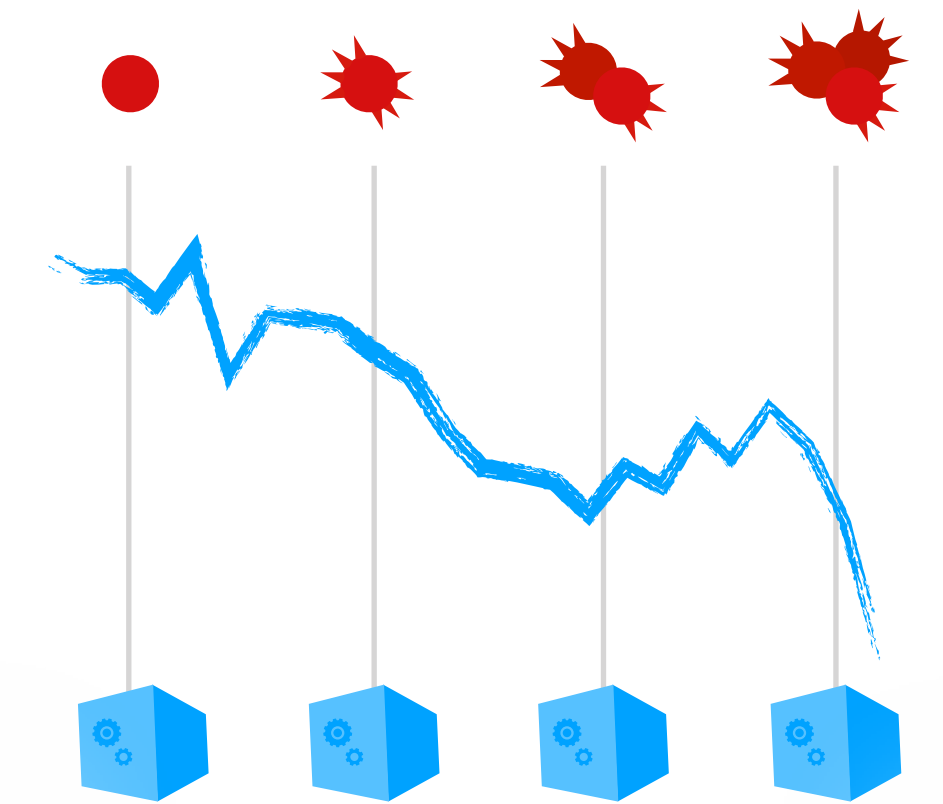
- Classic formulation of evasion attacks is ill-suited for reasoning about realizable evasive malware
- By reformulating, we can propose stronger attacks and easily compare against alternatives



Bigger Picture

Drifting scenarios caused by threats evolving over time

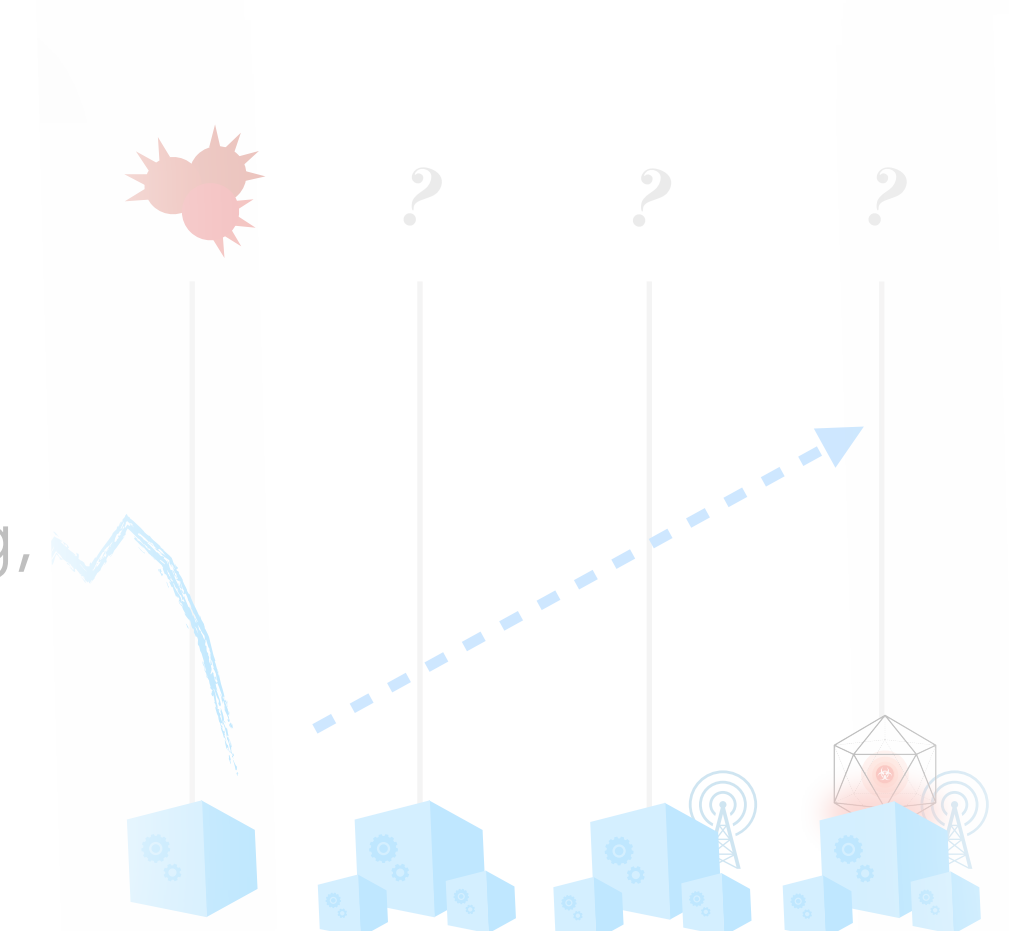
- How concept drift affects machine learning-based detectors in security settings
- The need for time-aware evaluations and metrics



Looking Ahead

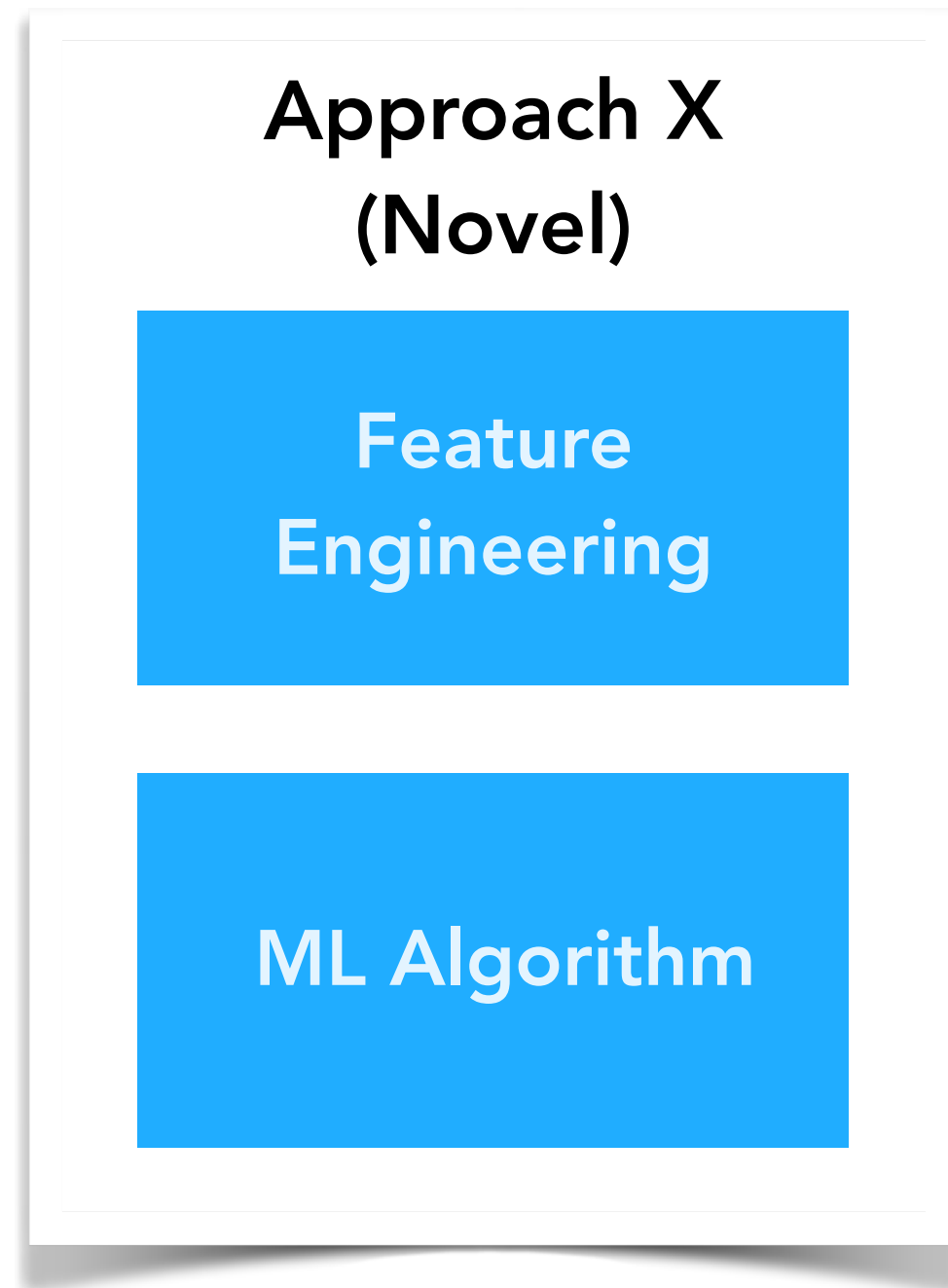
Quo vadis?

- Discussion of the future of trustworthy ML for system security
- Robust feature development, universal adversarial perturbations, realizable backdoors, drift forecasting, and the role of abstractions towards the Platonic ideal of interesting behaviors

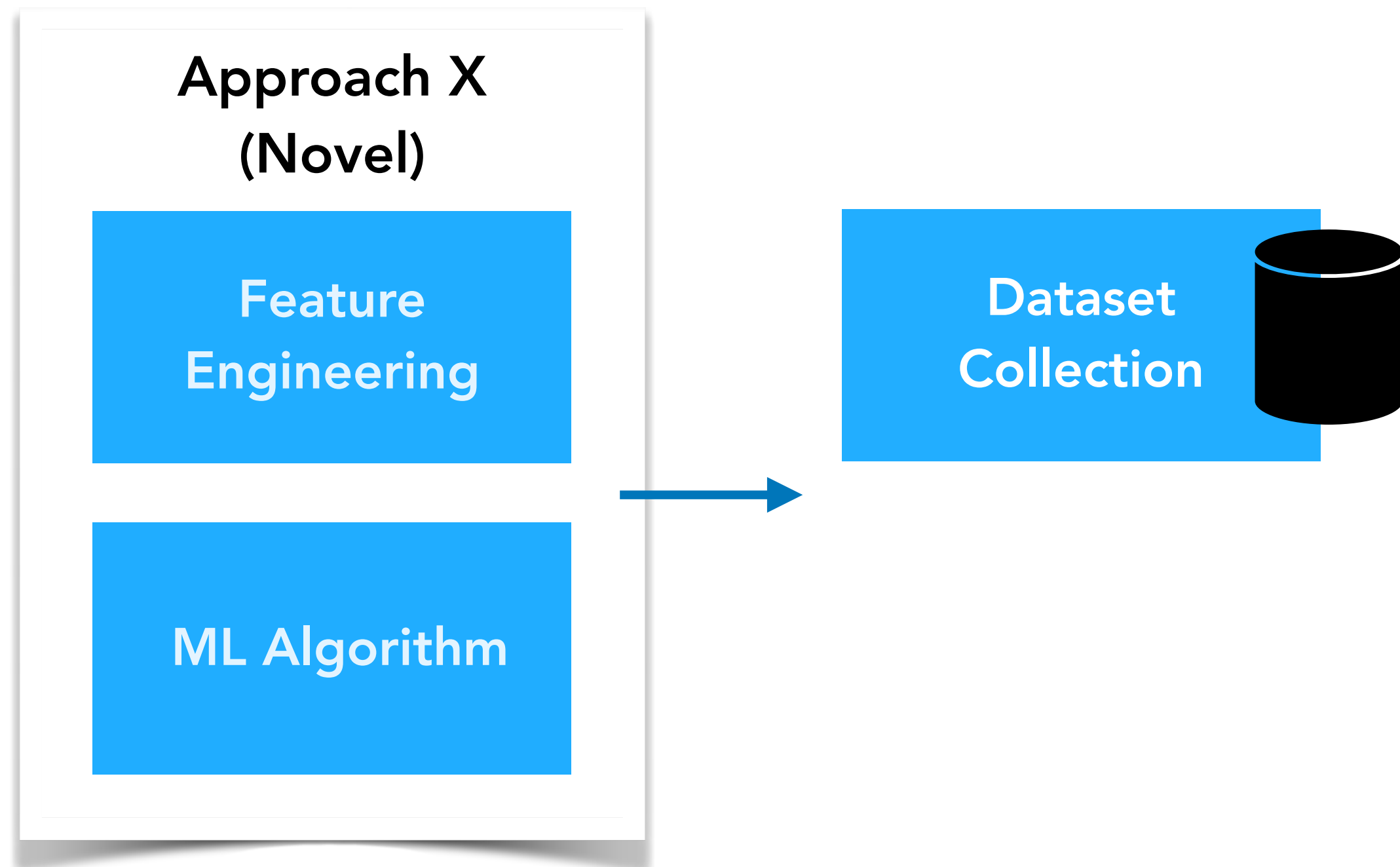


ML for Malware Detection 101


ML for Malware Detection 101



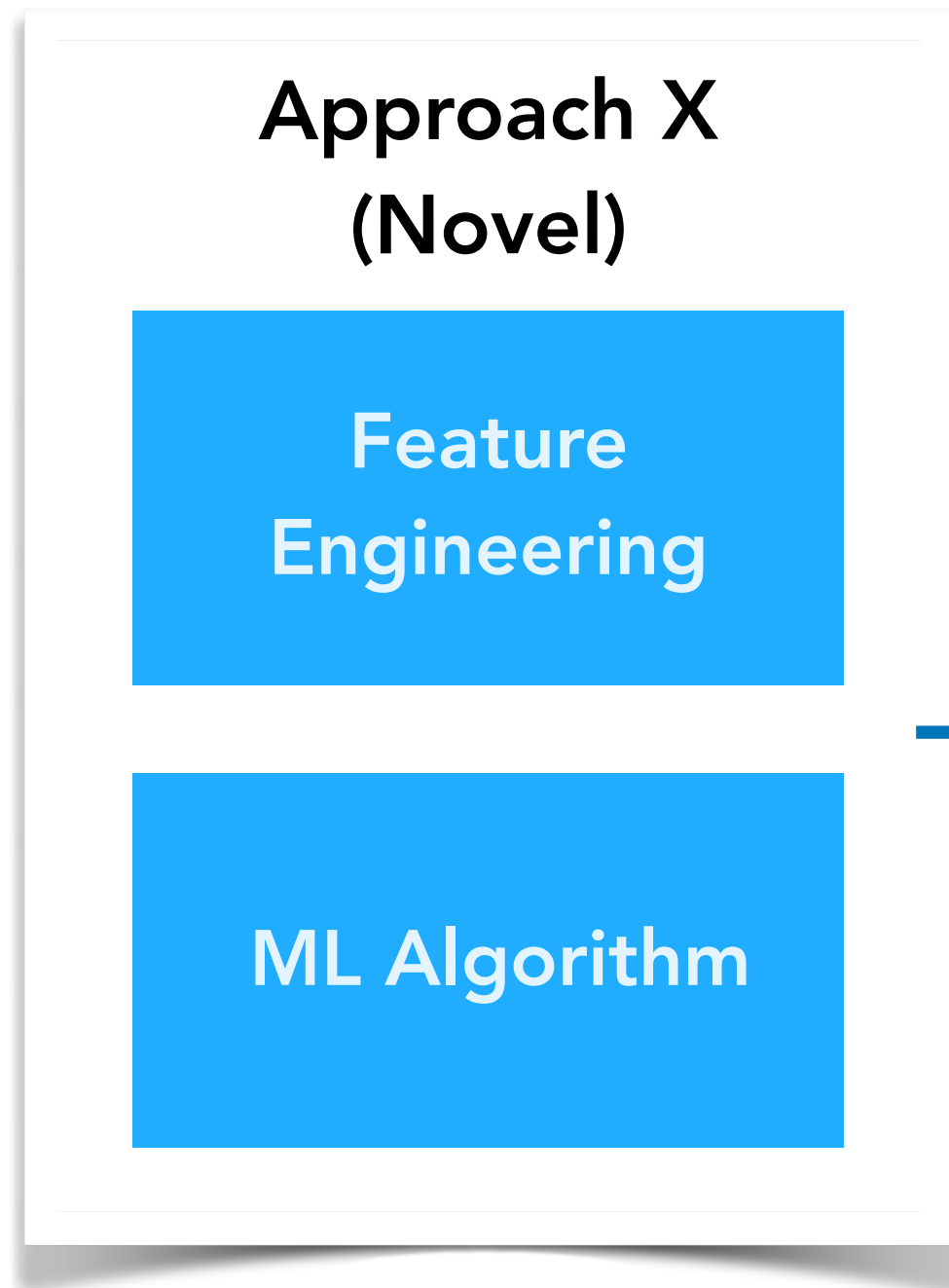
ML for Malware Detection 101



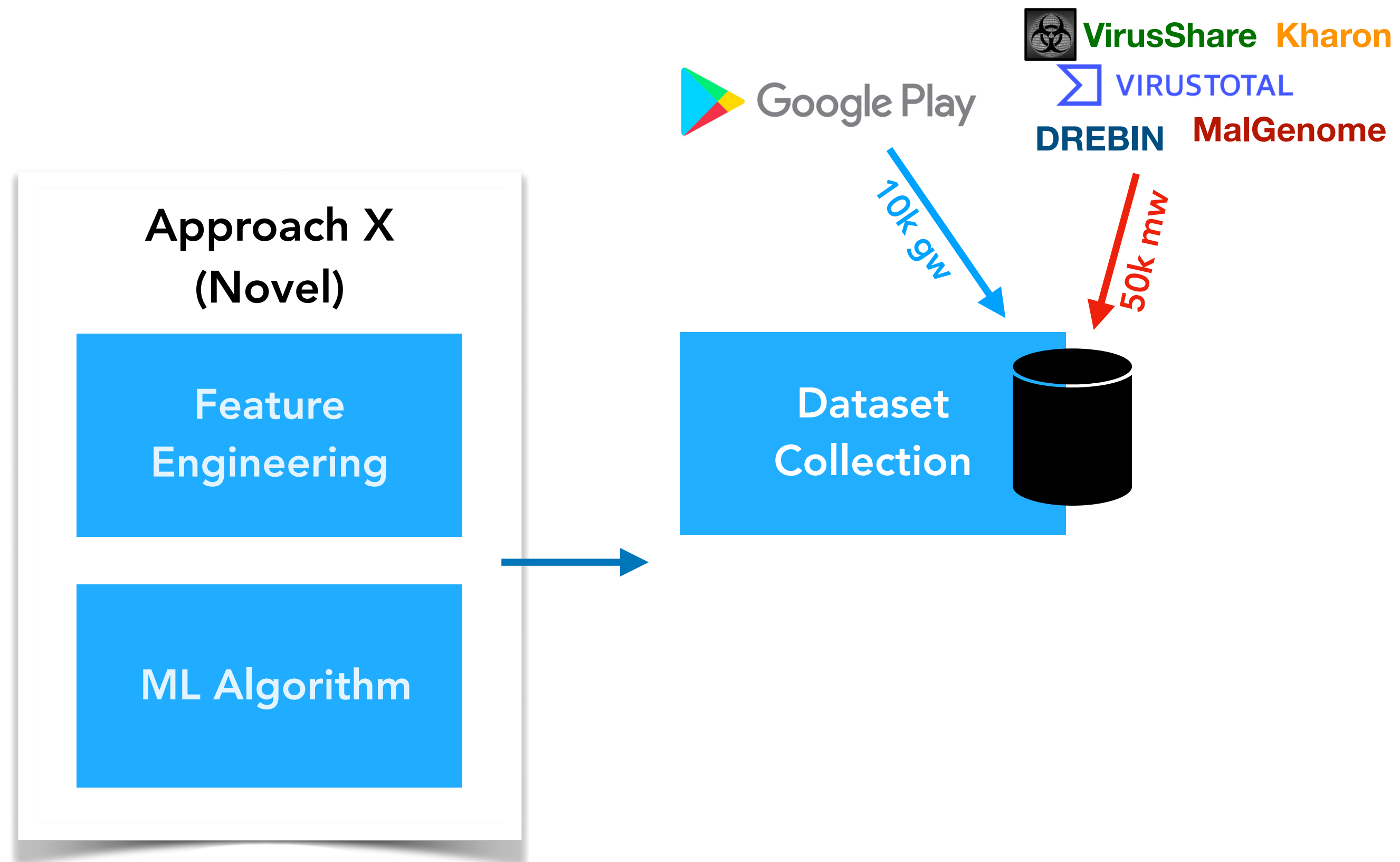
ML for Malware Detection 101

 **VirusShare** **Kharon**
 **VIRUSTOTAL**
DREBIN **MalGenome**

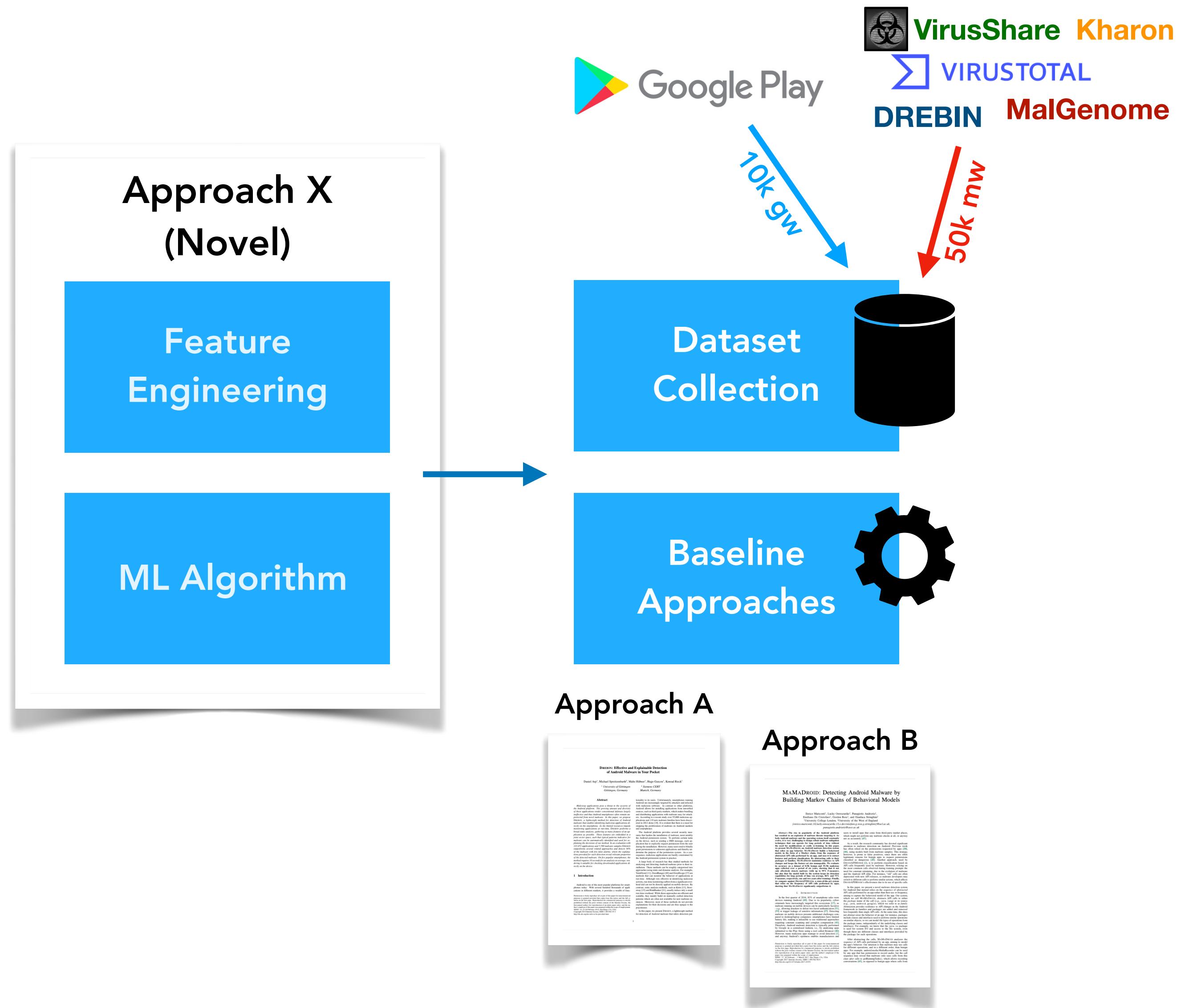
50k mw



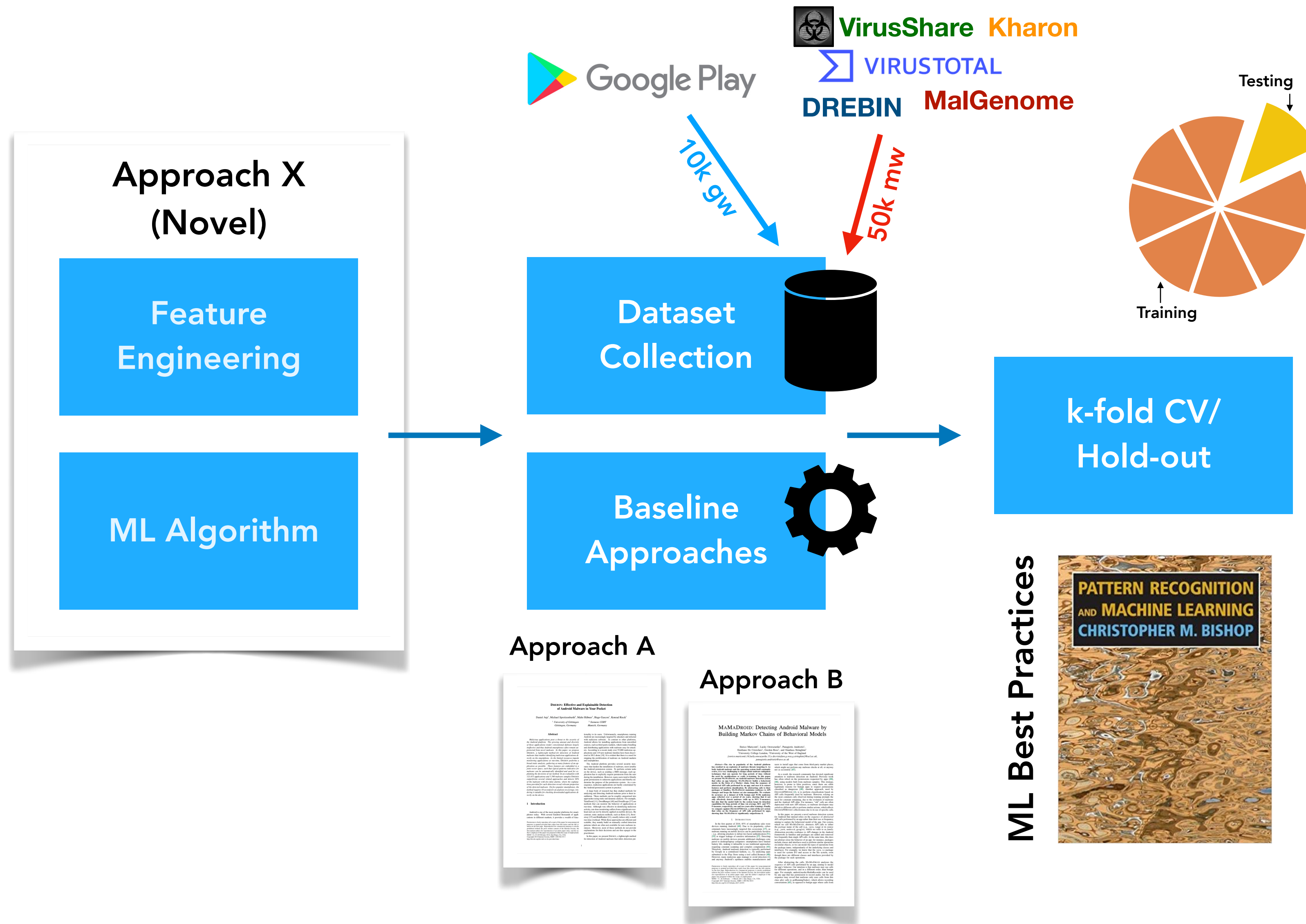
ML for Malware Detection 101



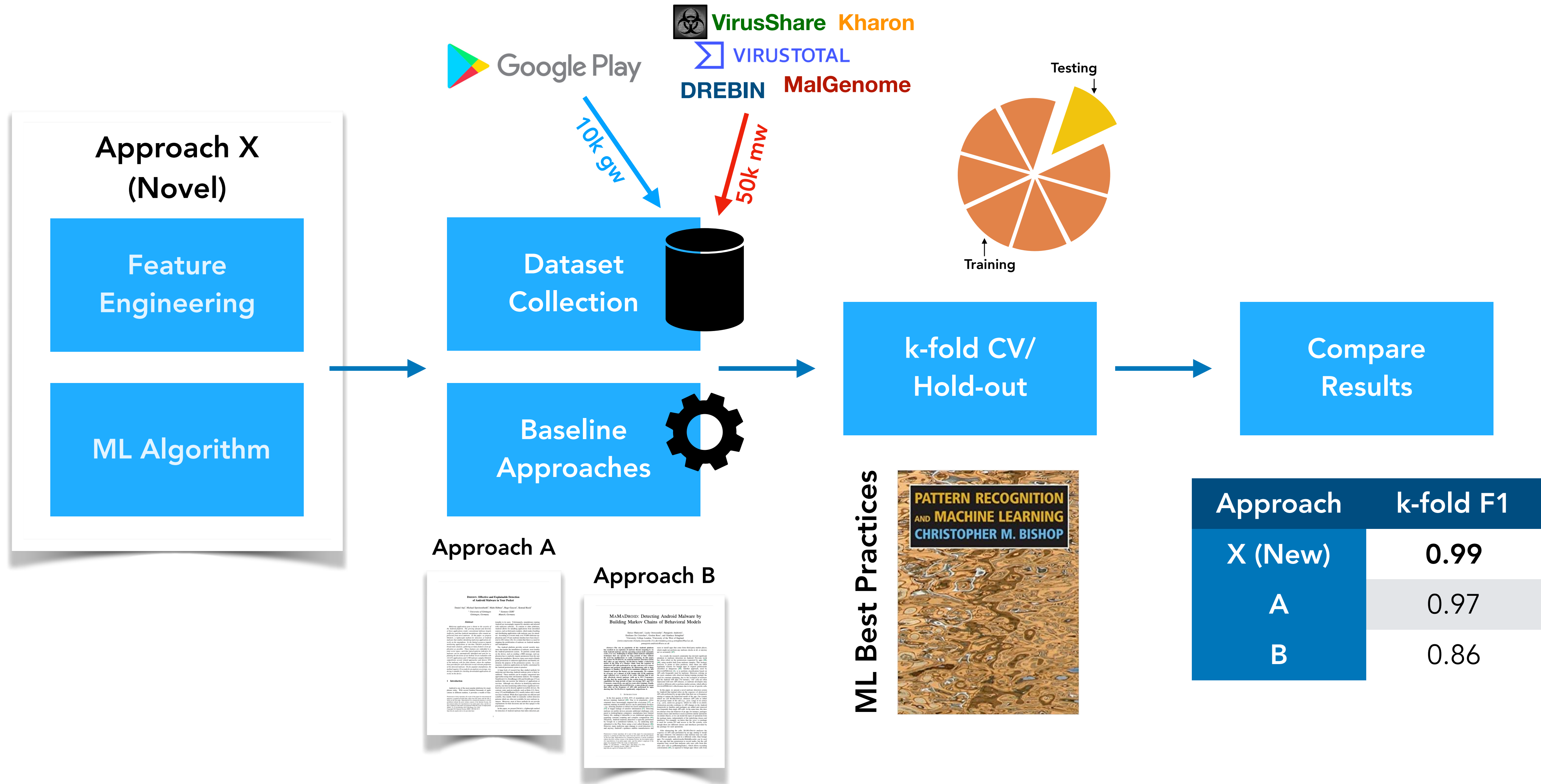
ML for Malware Detection 101



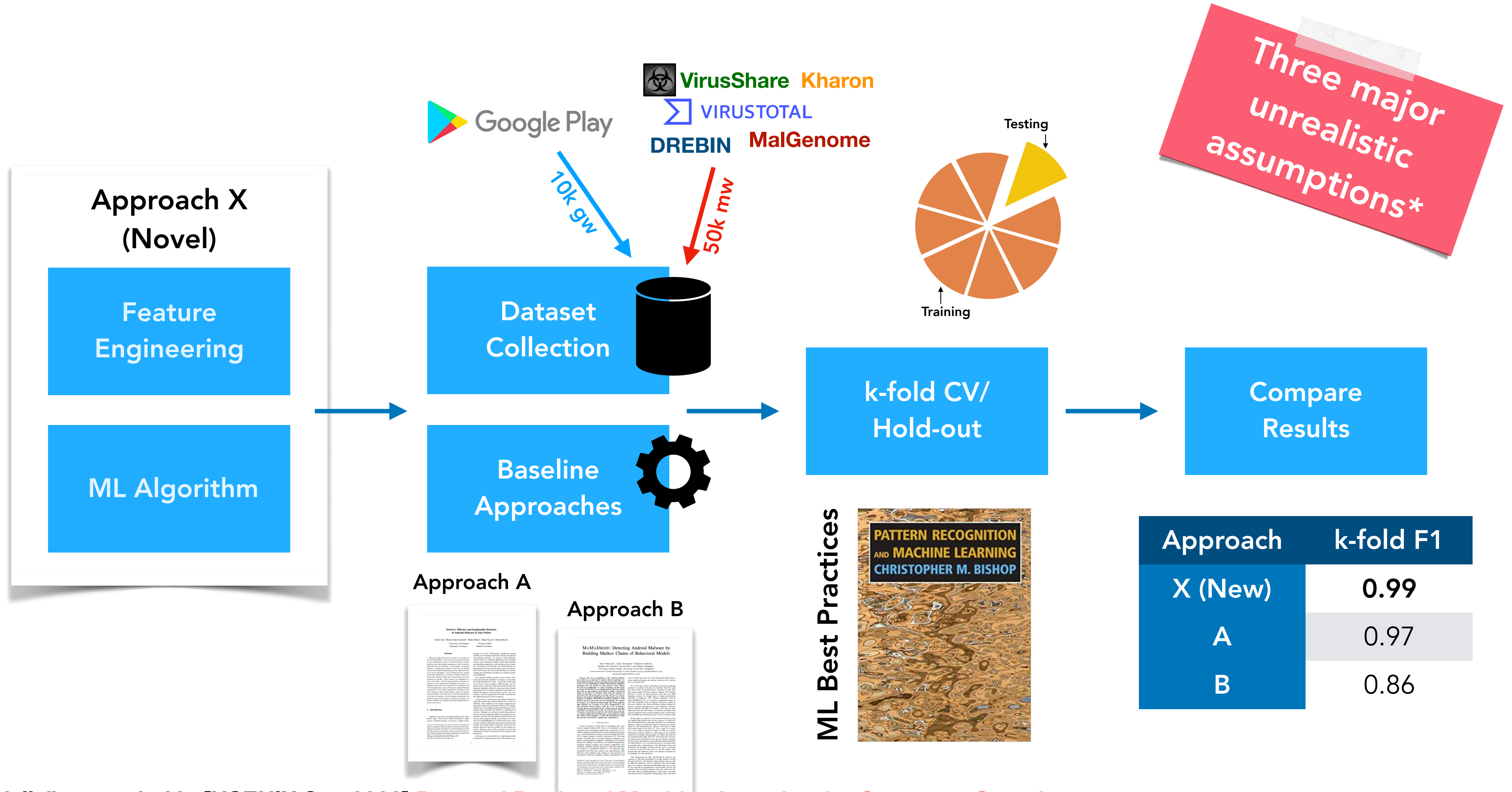
ML for Malware Detection 101



ML for Malware Detection 101



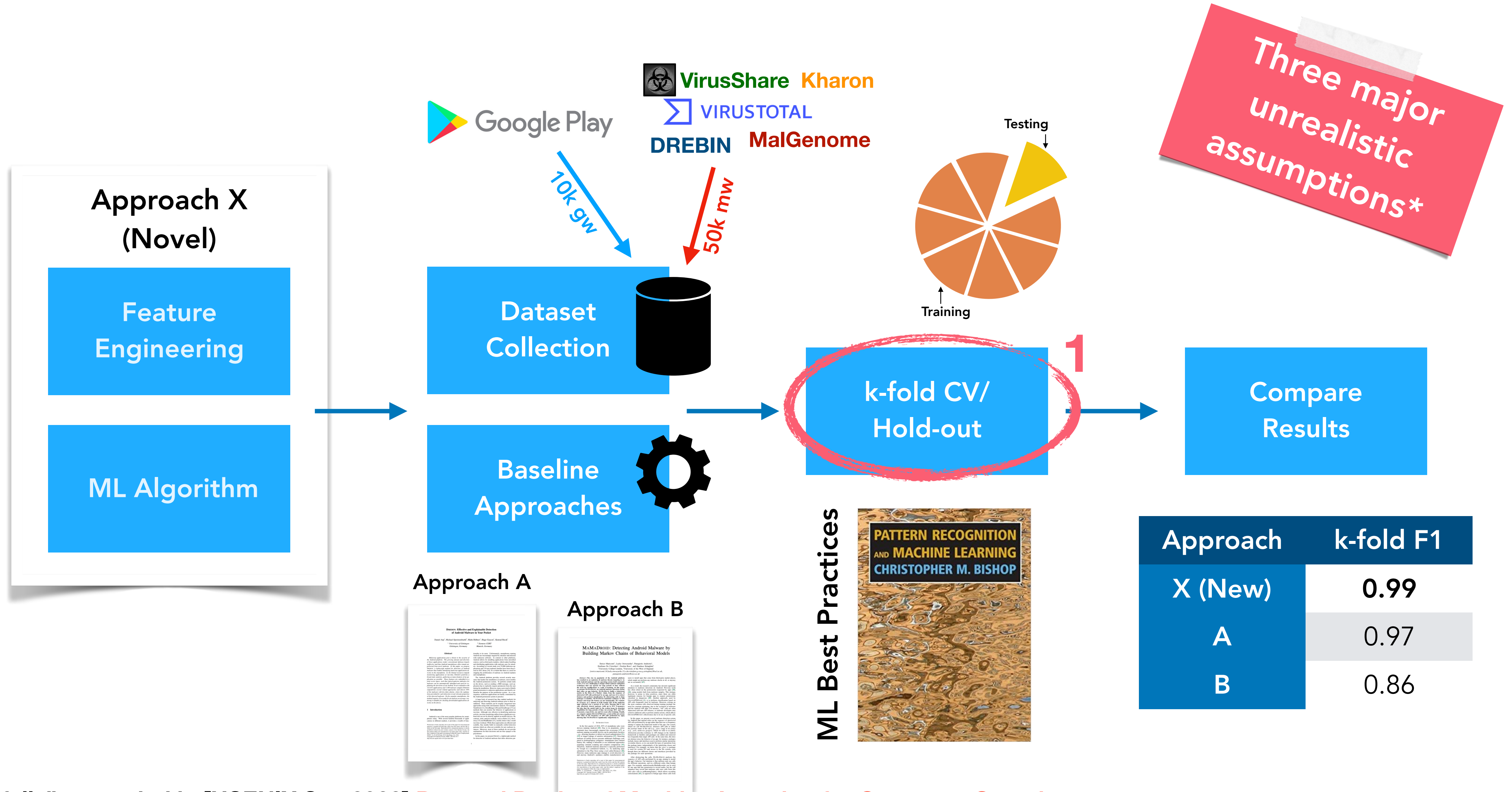
ML for Malware Detection 101



* "Pitfalls" expanded in [USENIX Sec 2022] **Dos and Don'ts of Machine Learning for Computer Security**

<https://s2lab.cs.ucl.ac.uk/downloads/dodo.pdf>

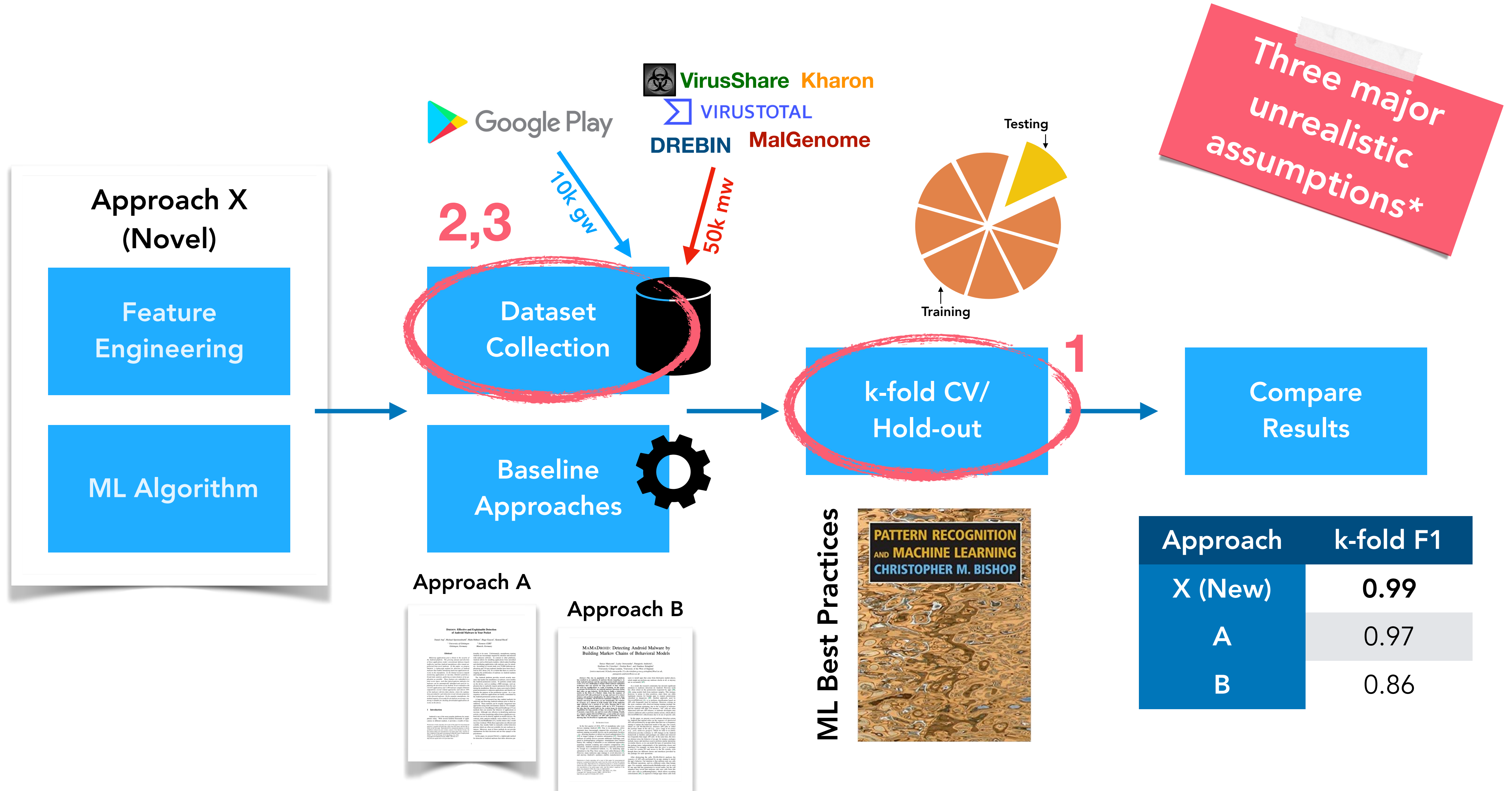
ML for Malware Detection 101



* "Pitfalls" expanded in [USENIX Sec 2022] **Dos and Don'ts of Machine Learning for Computer Security**

<https://s2lab.cs.ucl.ac.uk/downloads/dodo.pdf>

ML for Malware Detection 101



* "Pitfalls" expanded in [USENIX Sec 2022] **Dos and Don'ts of Machine Learning for Computer Security**

<https://s2lab.cs.ucl.ac.uk/downloads/dodo.pdf>

TESSERACT Framework

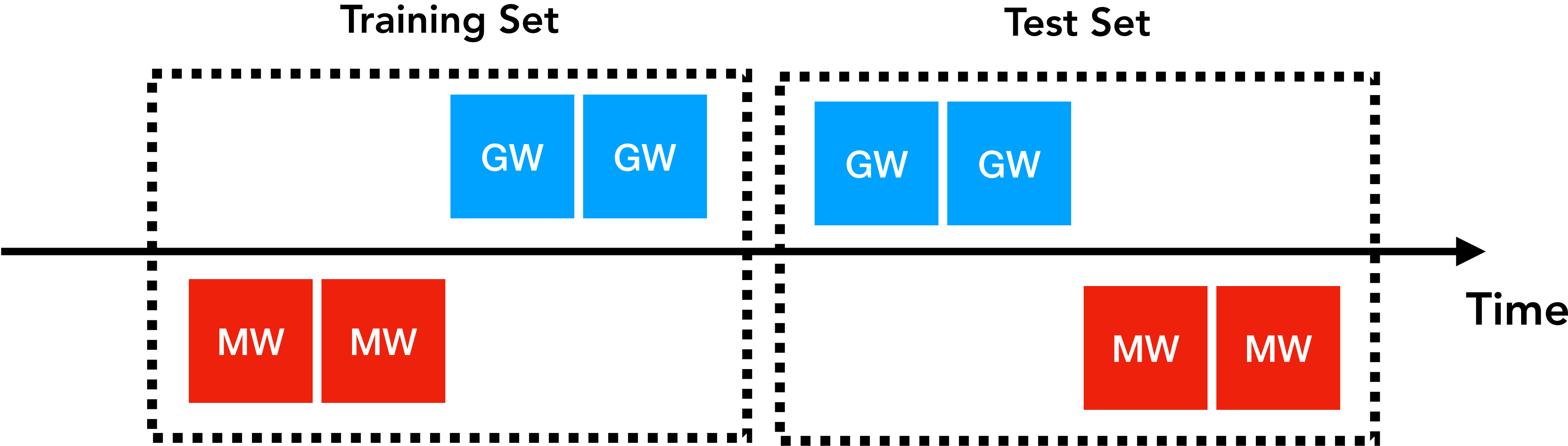
TESSERACT Framework

Experimental
Constraints



TESSERACT Framework

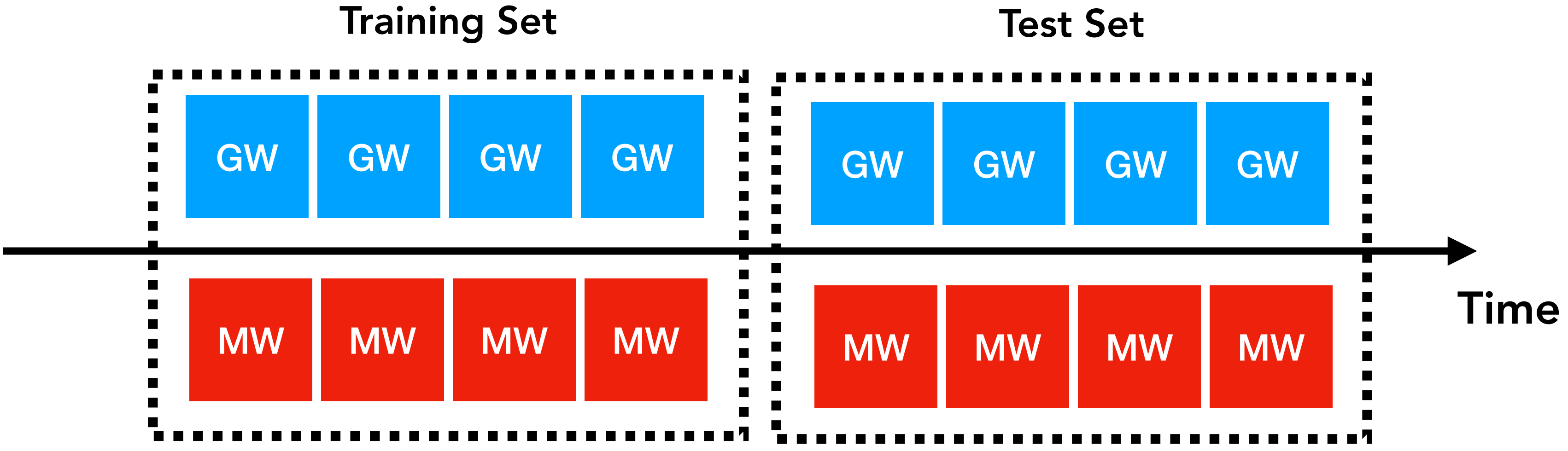
Experimental Constraints **C1** Temporal training consistency \rightarrow $\text{time}(\text{training}) < \text{time}(\text{testing})$



TESSERACT Framework

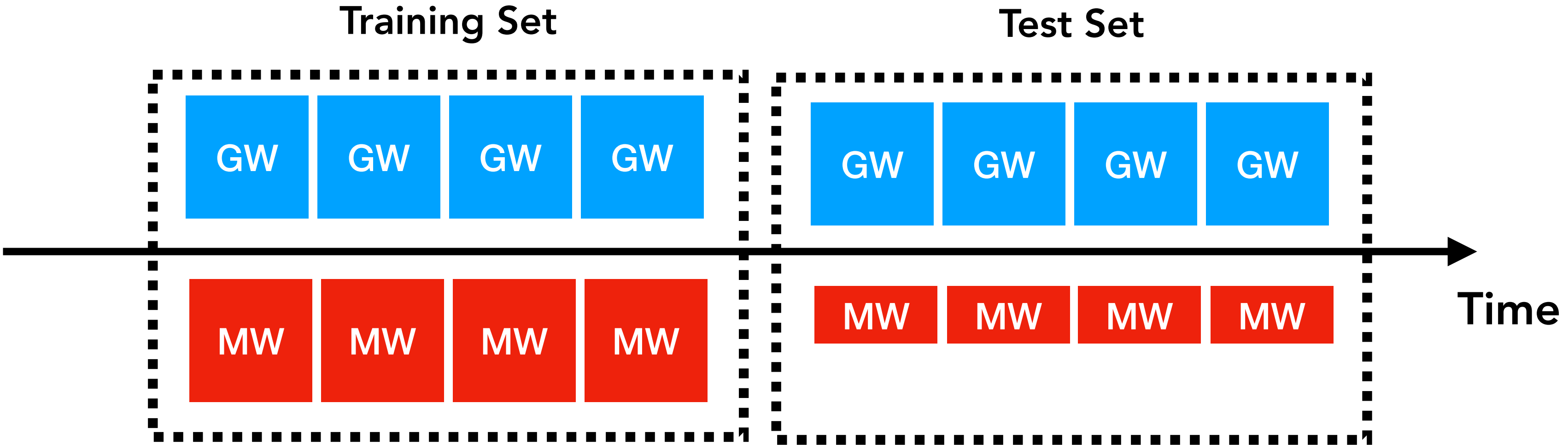
Experimental Constraints

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

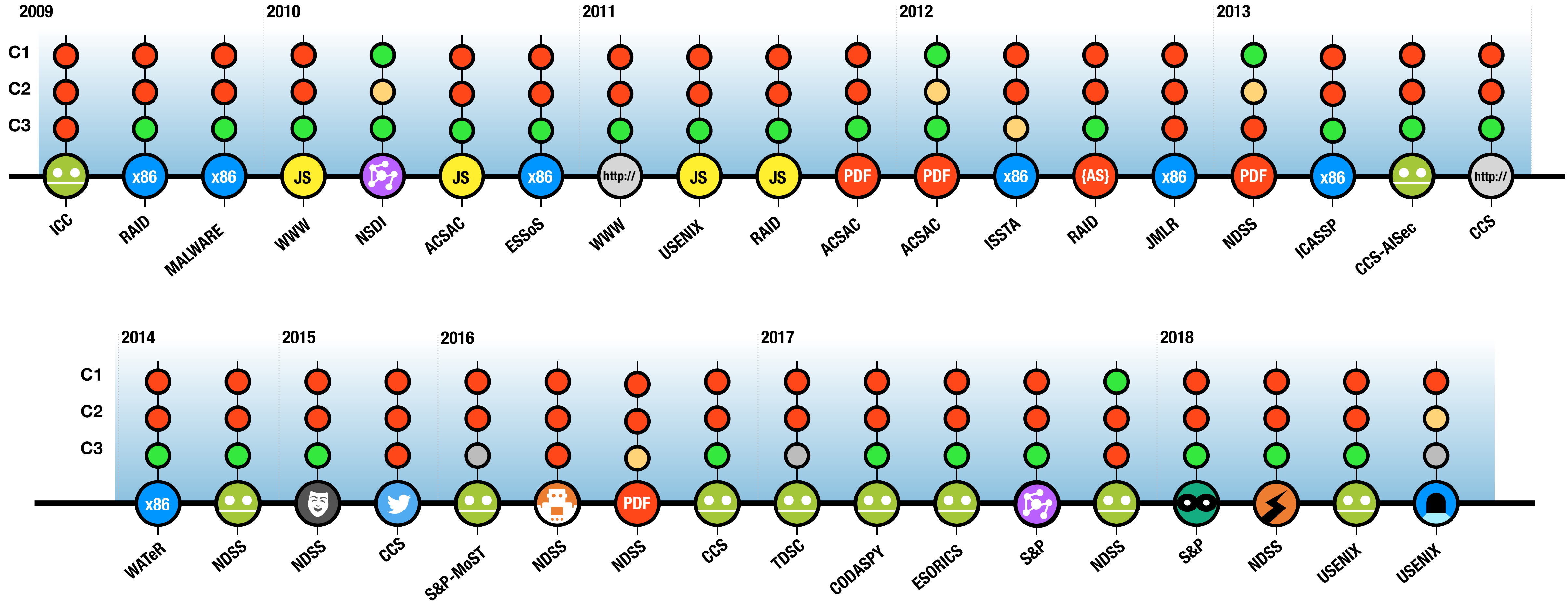


TESSERACT Framework

- Experimental Constraints
- C1** Temporal training consistency → $\text{time}(\text{training}) < \text{time}(\text{testing})$
 - C2** {good|mal}ware temporal consistency → $\text{time}(\text{gw}) = \text{time}(\text{mw})$
 - C3** Realistic testing classes ratio → realistic %mw in test



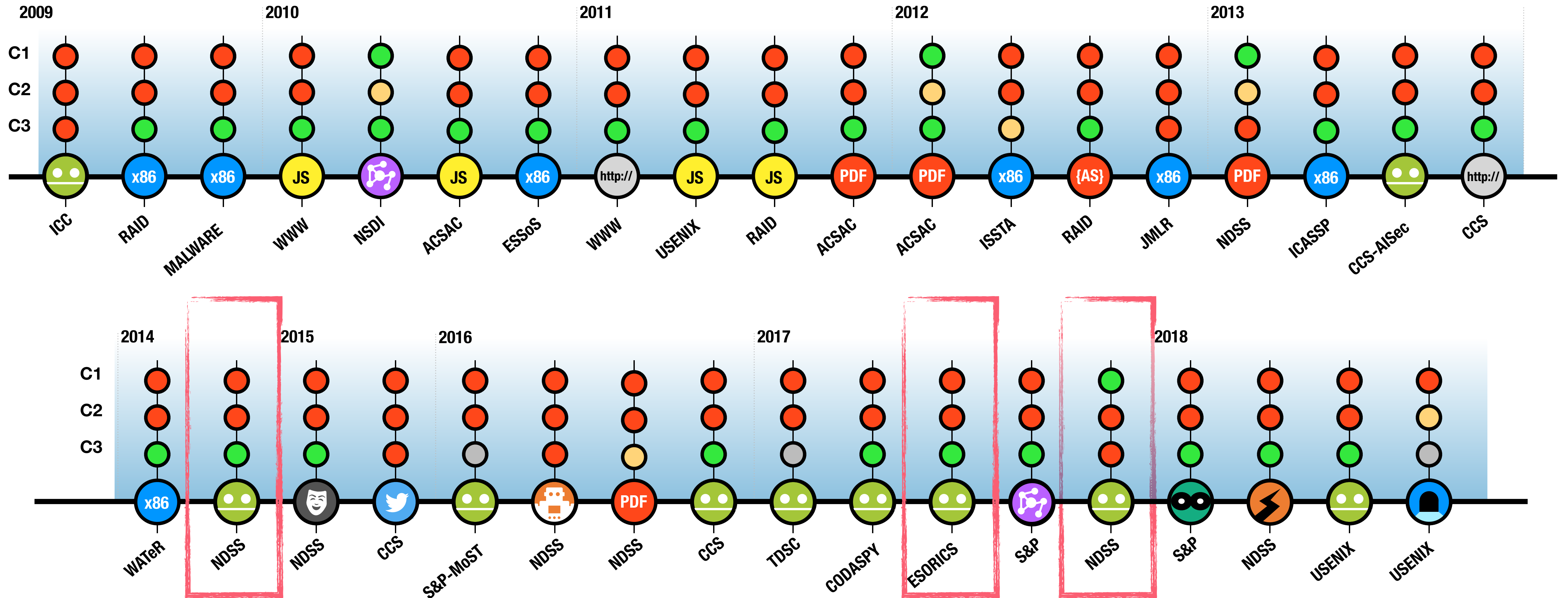
Endemic Problem



[USENIX Sec 2019] **TESSERACT: Eliminating Experimental Bias in Malware Classification across Space and Time**

<https://s2lab.cs.ucl.ac.uk/projects/tesseract>

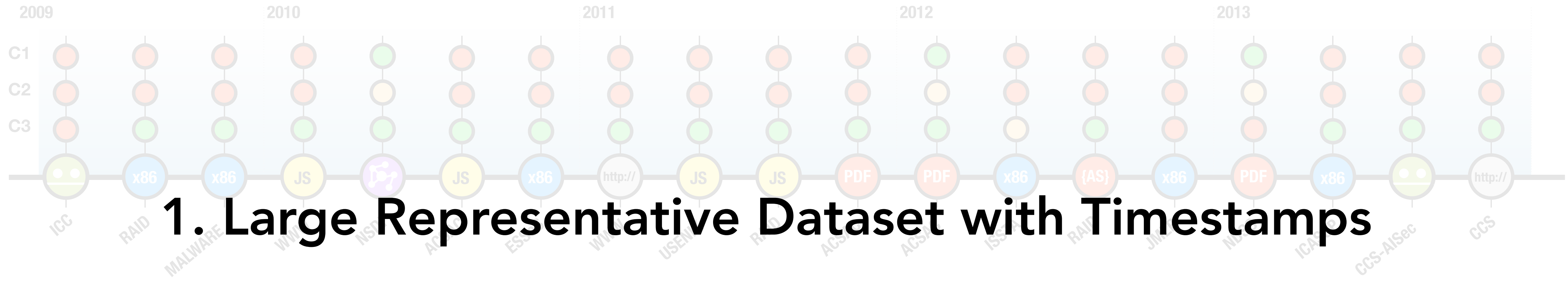
Endemic Problem



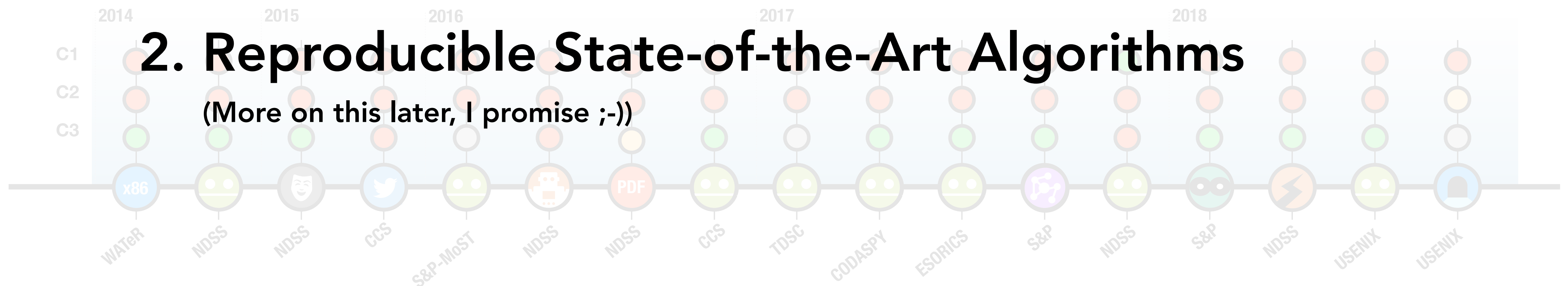
[USENIX Sec 2019] **TESSERACT: Eliminating Experimental Bias in Malware Classification across Space and Time**

<https://s2lab.cs.ucl.ac.uk/projects/tesseract>

Endemic Problem



1. Large Representative Dataset with Timestamps



2. Reproducible State-of-the-Art Algorithms

(More on this later, I promise ;-))

Dataset

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

TESSERACT Evaluations

TESSERACT Evaluations

Experimental
Constraints

C1 Temporal training consistency

C2 {good|mal}ware temporal consistency

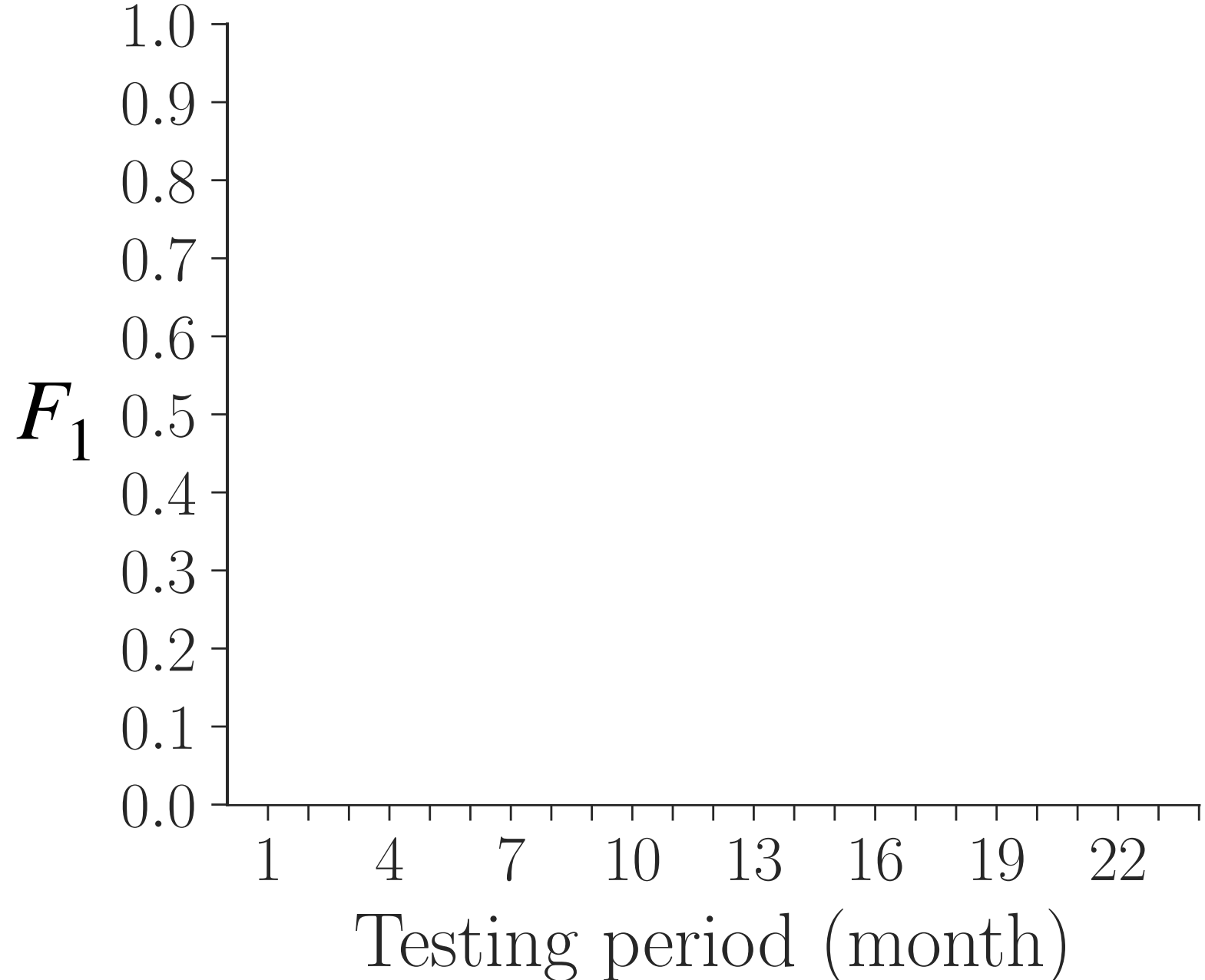
C3 Realistic testing classes ratio

TESSERACT Evaluations

Experimental Constraints

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

NDSS14

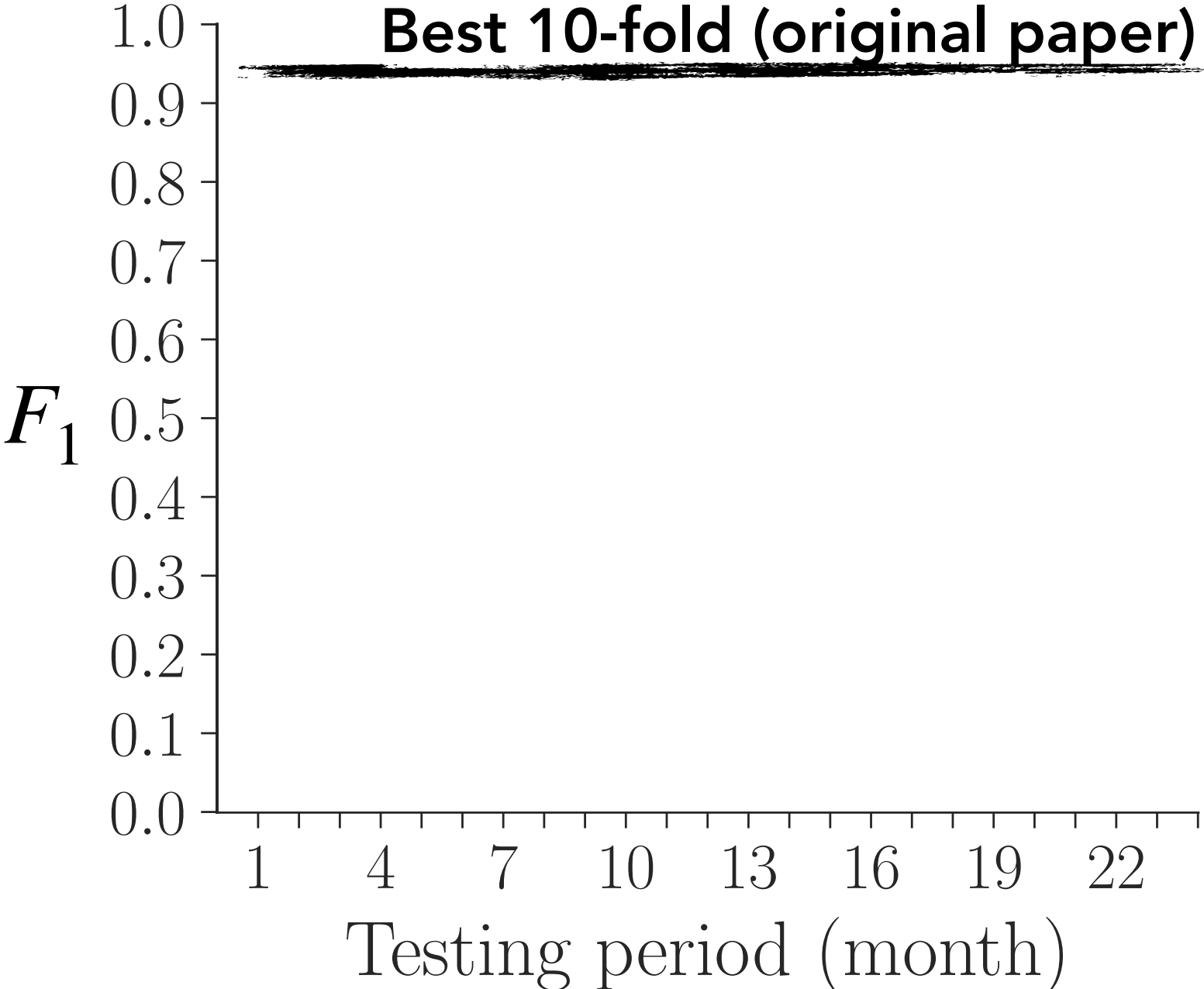


TESSERACT Evaluations

Experimental Constraints

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

NDSS14

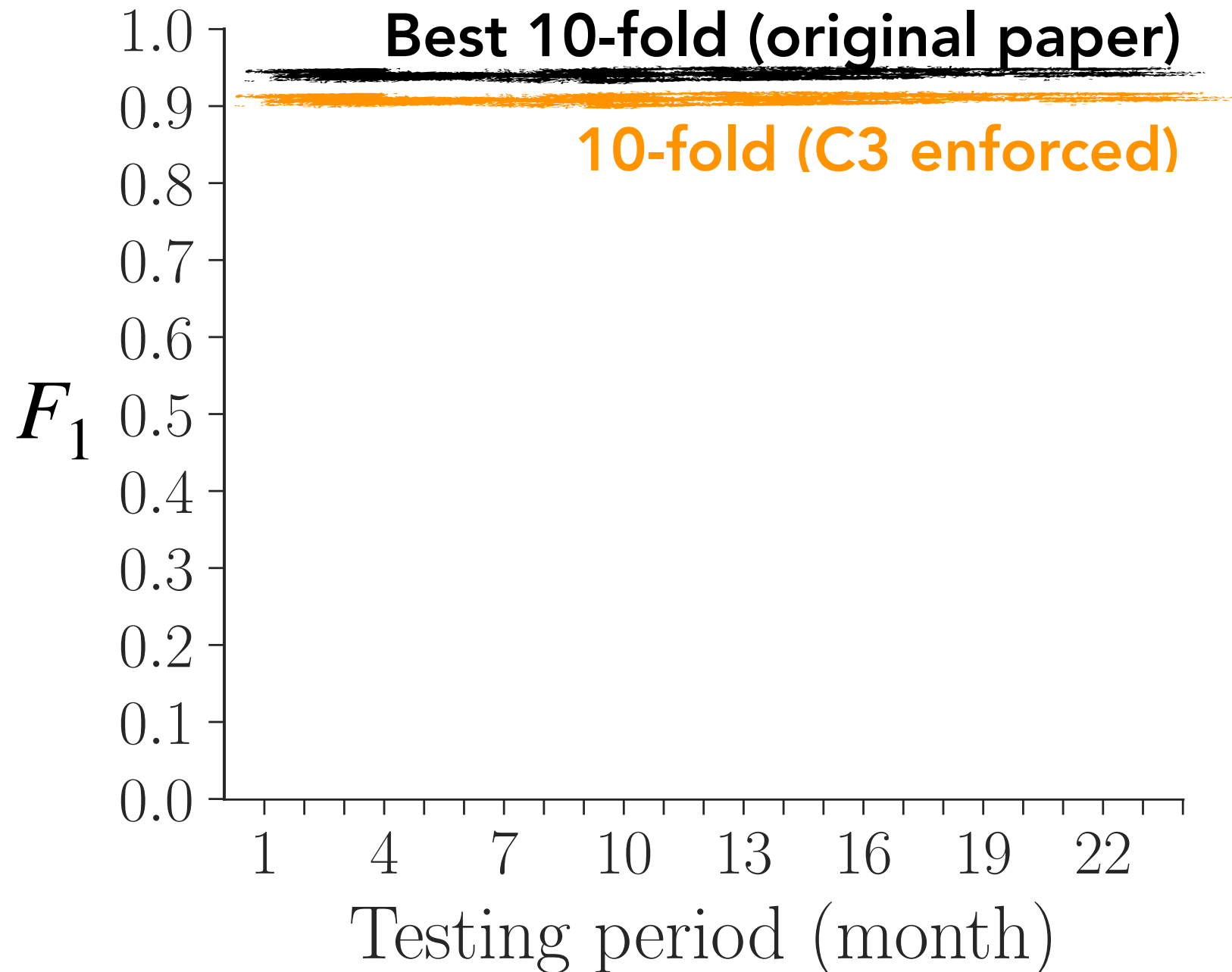


TESSERACT Evaluations

Experimental Constraints

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

NDSS14

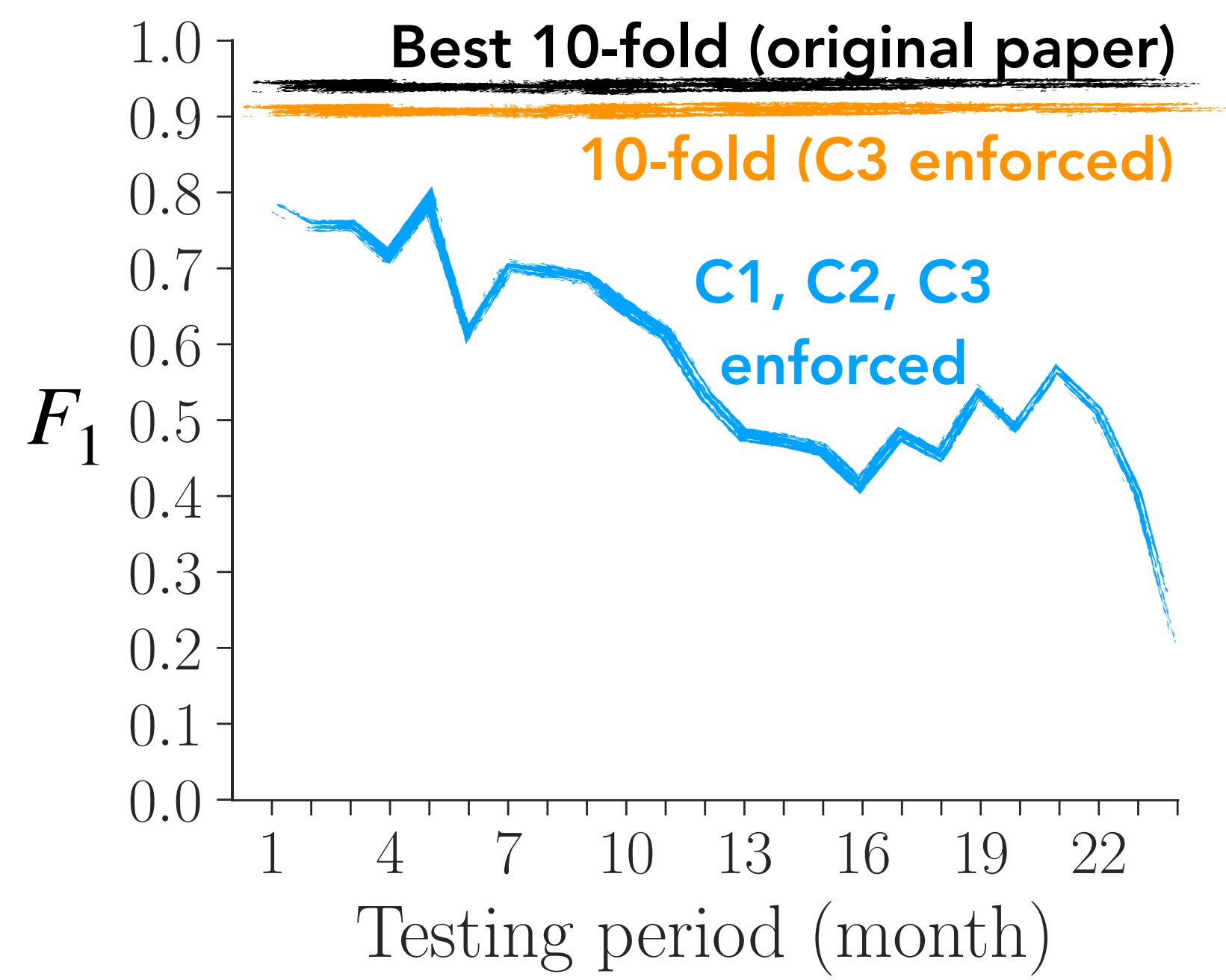


TESSERACT Evaluations

Experimental Constraints

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

NDSS14

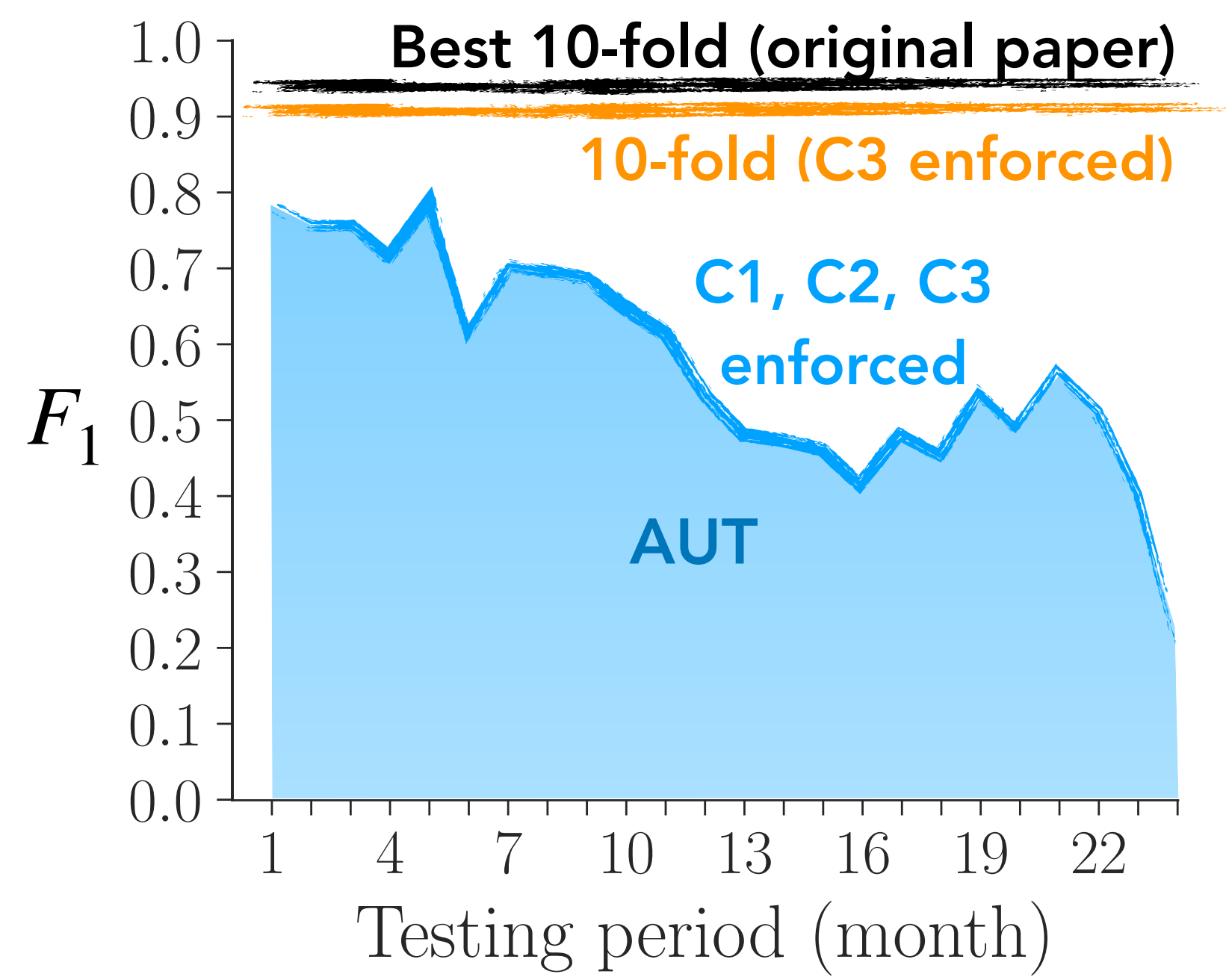


TESSERACT Evaluations

Experimental Constraints

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

NDSS14



Area Under Time
AUT(Metric, Period)

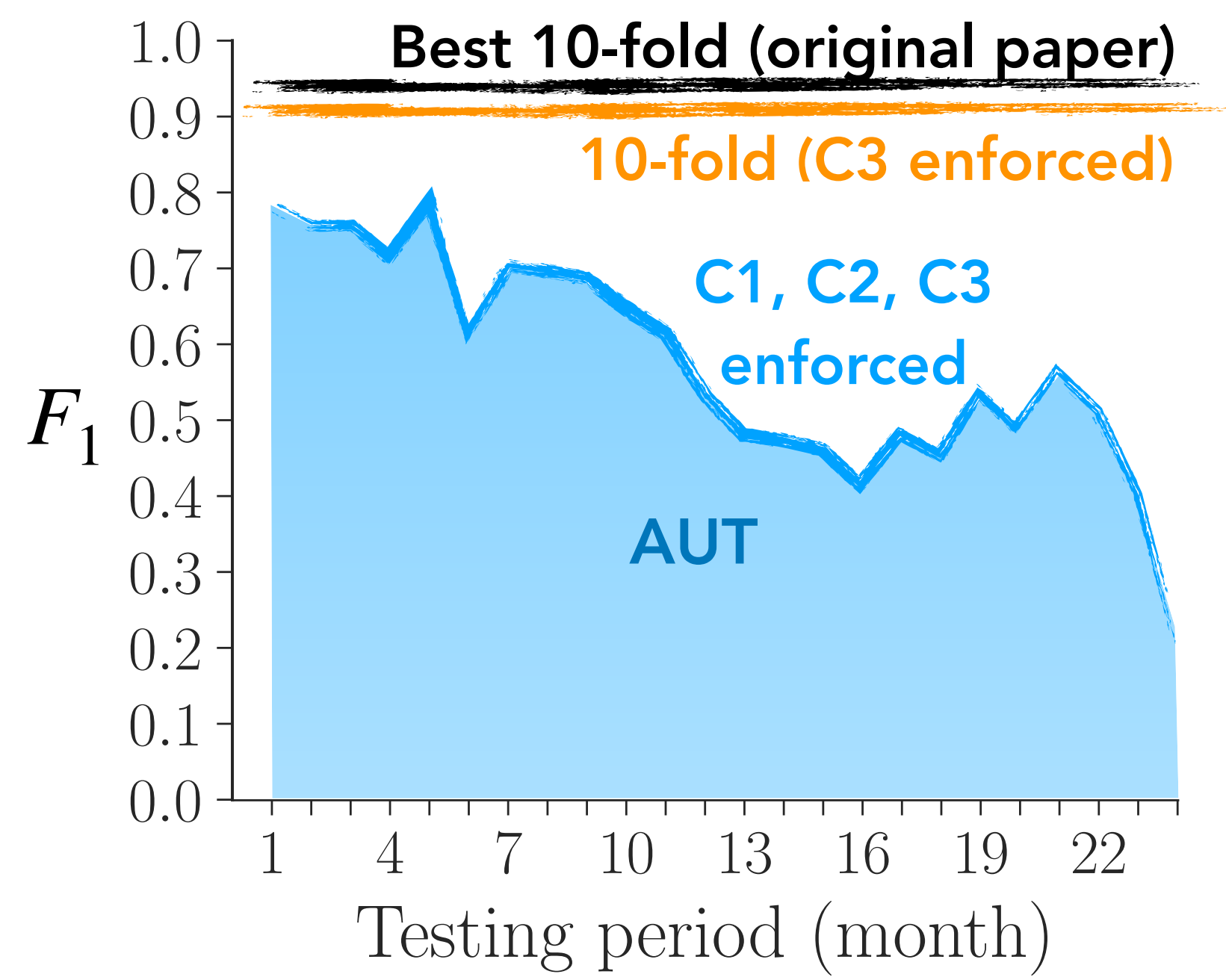
$$AUT(F_1, 24m) = 0.58$$

TESSERACT Evaluations

Experimental Constraints

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

NDSS14



Area Under Time
AUT(Metric, Period)

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

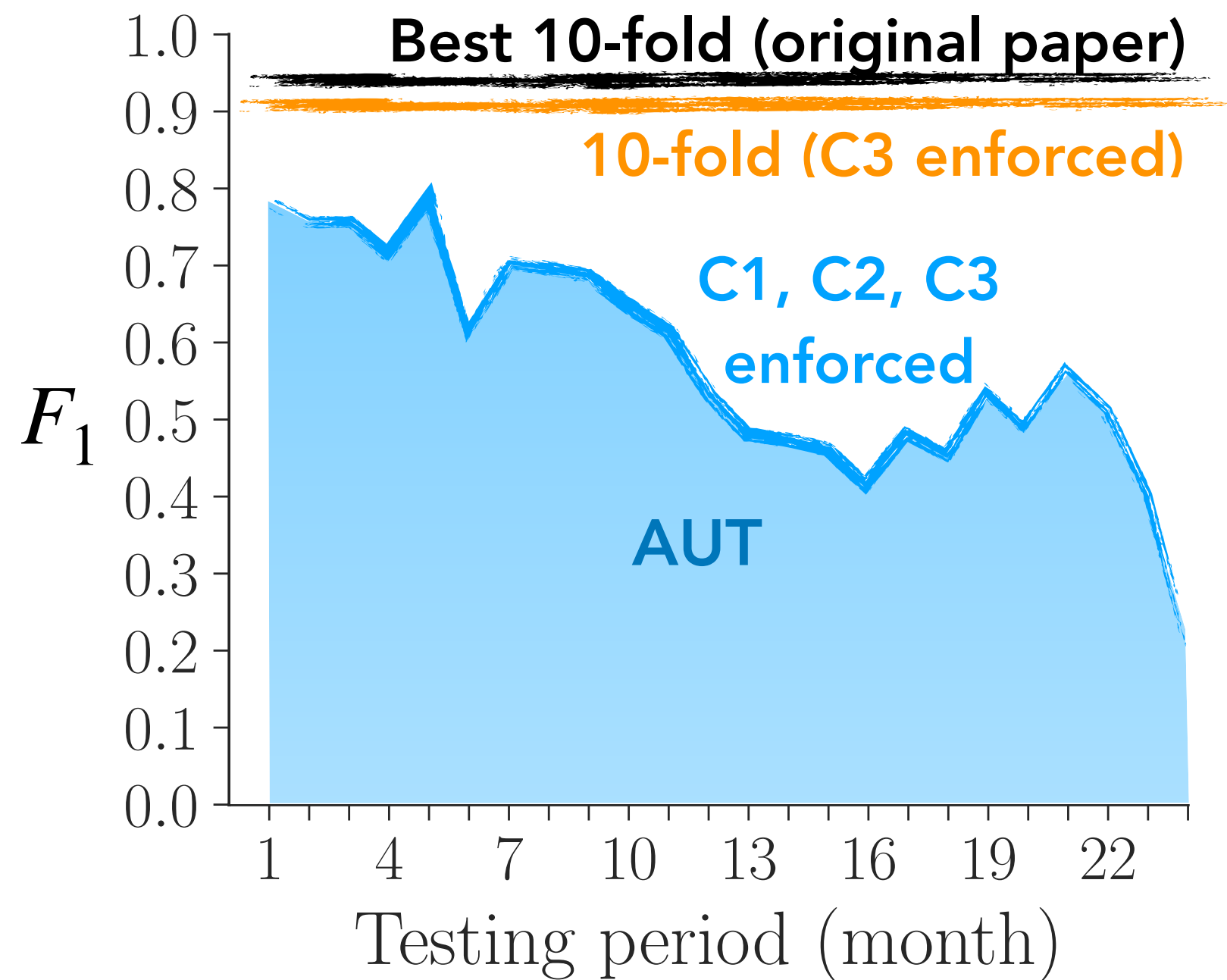
TESSERACT Evaluations

Experimental Constraints

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

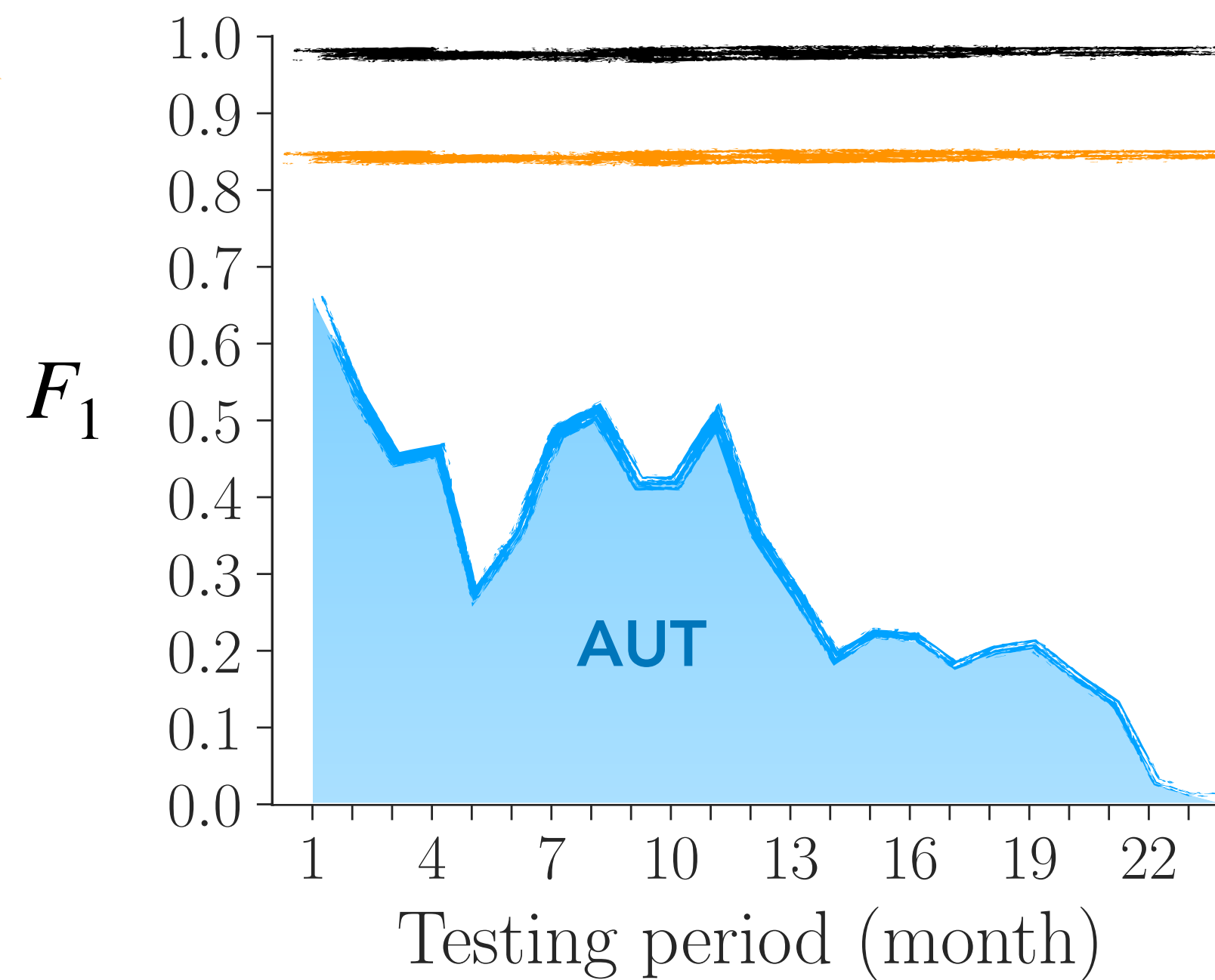
Area Under Time
AUT(Metric, Period)

NDSS14



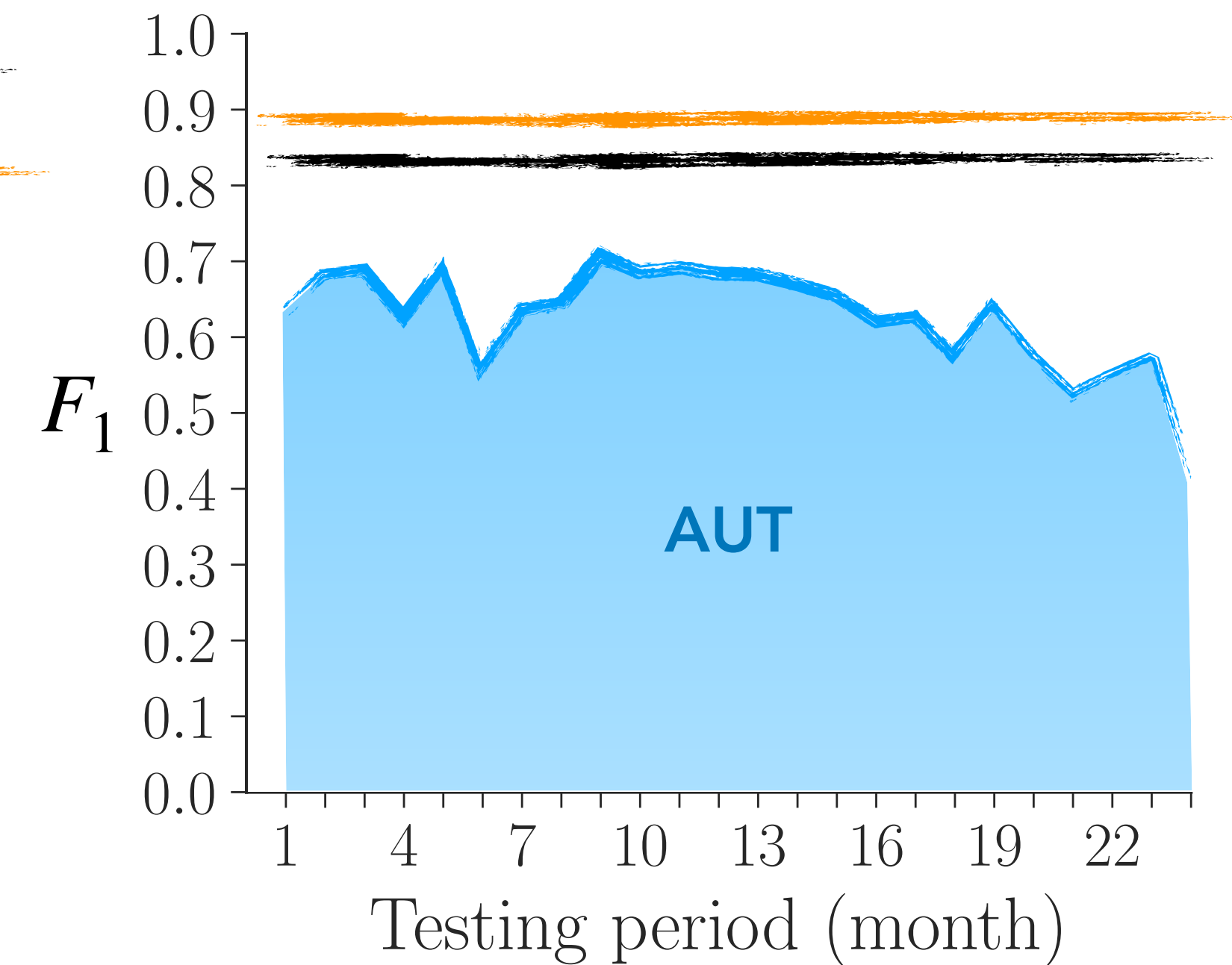
$$AUT(F_1, 24m) = 0.58$$

NDSS17



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

ESORICS17



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

TESSERACT: Actionable Points

TESSERACT: Actionable Points

Realistic Evaluations

- Reveals performance in more realistic setting
- Removes space-time experimental bias
- **Practitioners:** Choose Best Solution
- **Researchers:** Evaluate New Solutions

TESSERACT: Actionable Points

Realistic Evaluations

- Reveals performance in more realistic setting
- Removes space-time experimental bias
- **Practitioners:** Choose Best Solution
- **Researchers:** Evaluate New Solutions

Performance-Cost Trade Offs

- **Detection Performance** (e.g., AUT F_1)
- **Labeling Cost** for retraining (e.g., manpower)
- **Quarantine Cost** for rejection (e.g., low-confidence decisions)

Incremental Retraining

Active Learning

Online Learning

Rejection*

TESSERACT: Actionable Points

Realistic Evaluations

- Reveals performance in more realistic setting
- Removes space-time experimental bias
- **Practitioners:** Choose Best Solution
- **Researchers:** Evaluate New Solutions

Performance-Cost Trade Offs

- **Detection Performance** (e.g., AUT F_1)
- **Labeling Cost** for retraining (e.g., manpower)
- **Quarantine Cost** for rejection (e.g., low-confidence decisions)

Incremental Retraining

Active Learning

Online Learning

Rejection*

- Incremental retraining costly?*
- Active learning and query strategies* **
- Online learning and the risk of self-poisoning+

*[USENIX Sec 2019] **TESSERACT: Eliminating Experimental Bias in Malware Classification across Space and Time**

[arXiv 2023] Chen et al. **Continuous Learning for Android Malware Detection

+ [AISec 2021] **Investigating Labelless Drift Adaptation for Malware Detection** — <https://s2lab.cs.ucl.ac.uk/downloads/aisec47-kanA.pdf>

TESSERACT: Actionable Points

Realistic Evaluations

- Reveals performance in more realistic setting
- Removes space-time experimental bias
- **Practitioners:** Choose Best Solution
- **Researchers:** Evaluate New Solutions

Performance-Cost Trade Offs

- **Detection Performance** (e.g., AUT F_1)
- **Labeling Cost** for retraining (e.g., manpower)
- **Quarantine Cost** for rejection (e.g., low-confidence decisions)

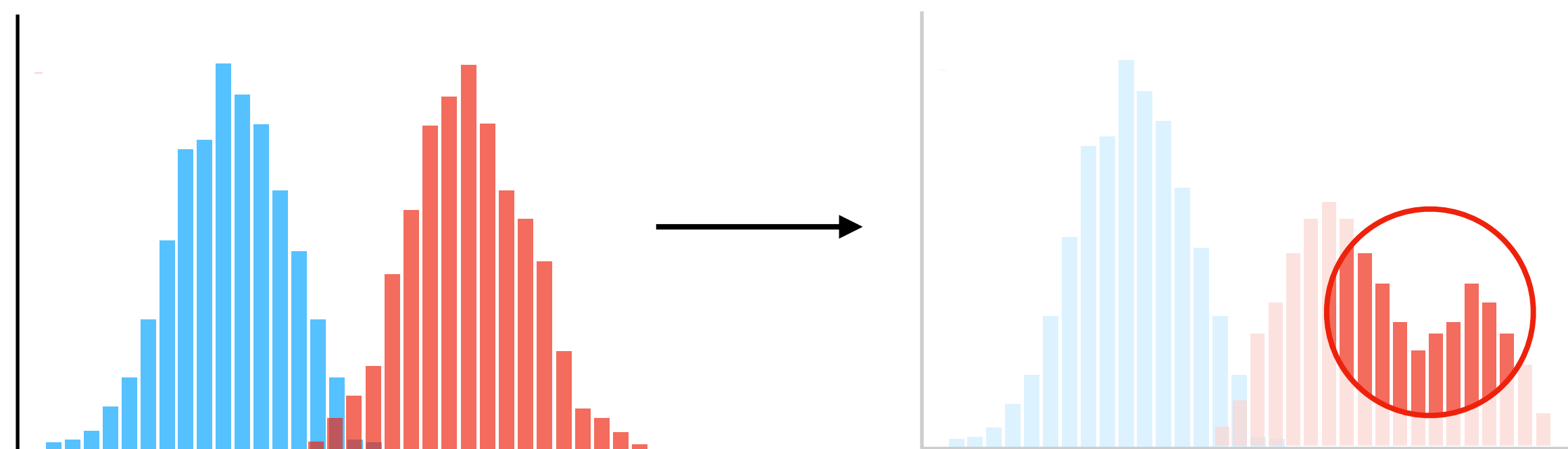
Incremental Retraining

Active Learning

Online Learning

Rejection*

As well as measuring the overall effect of drift we can **identify** specific aspects of the drift and **reject** objects that are likely to be misclassified.



* [USENIX Sec 2017] **Transcend: Detecting Concept Drift in Malware Classification Models**

* [IEEE S&P 2022] **Transcending Transcend: Revisiting Malware Classification in the Presence of Concept Drift**

Revisiting Classification in the Presence of Concept Drift

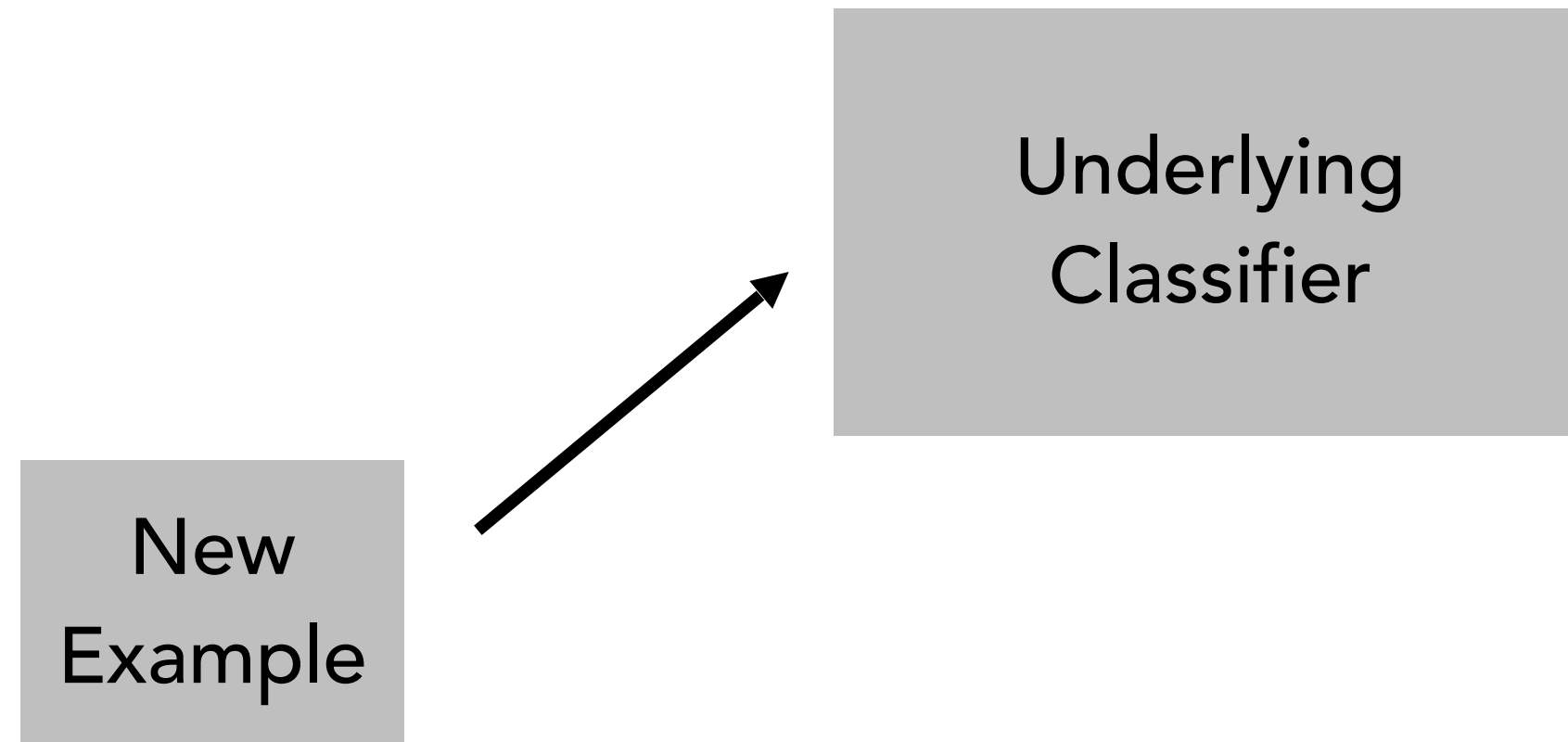
Revisiting Classification in the Presence of Concept Drift

Covariate Shift: Change in feature distribution $P(x \in X)$

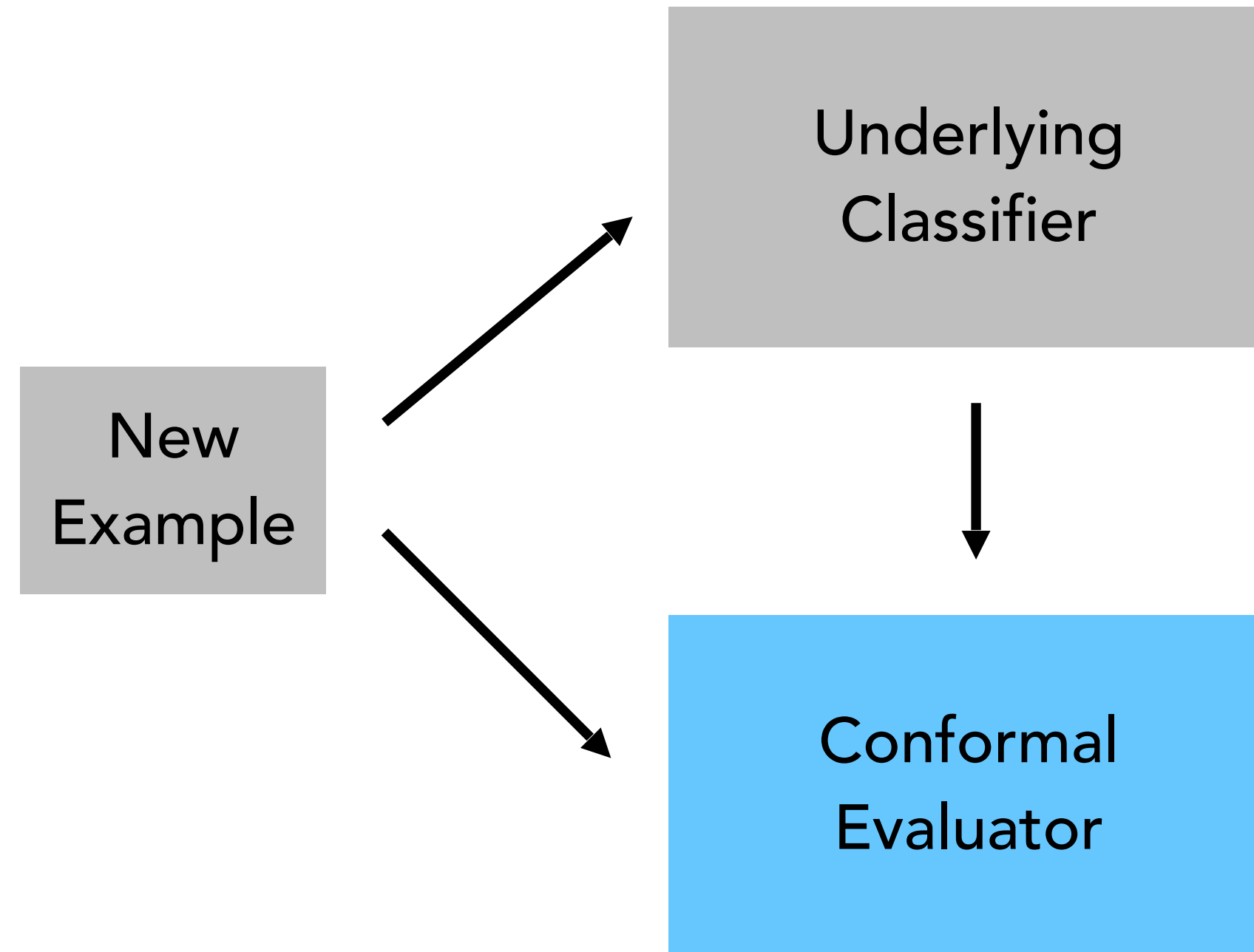
Prior-probability Shift: Change in class base rate $P(y \in Y)$

Concept Drift: Change in ground truth definition $P(y \in Y|x \in X)$

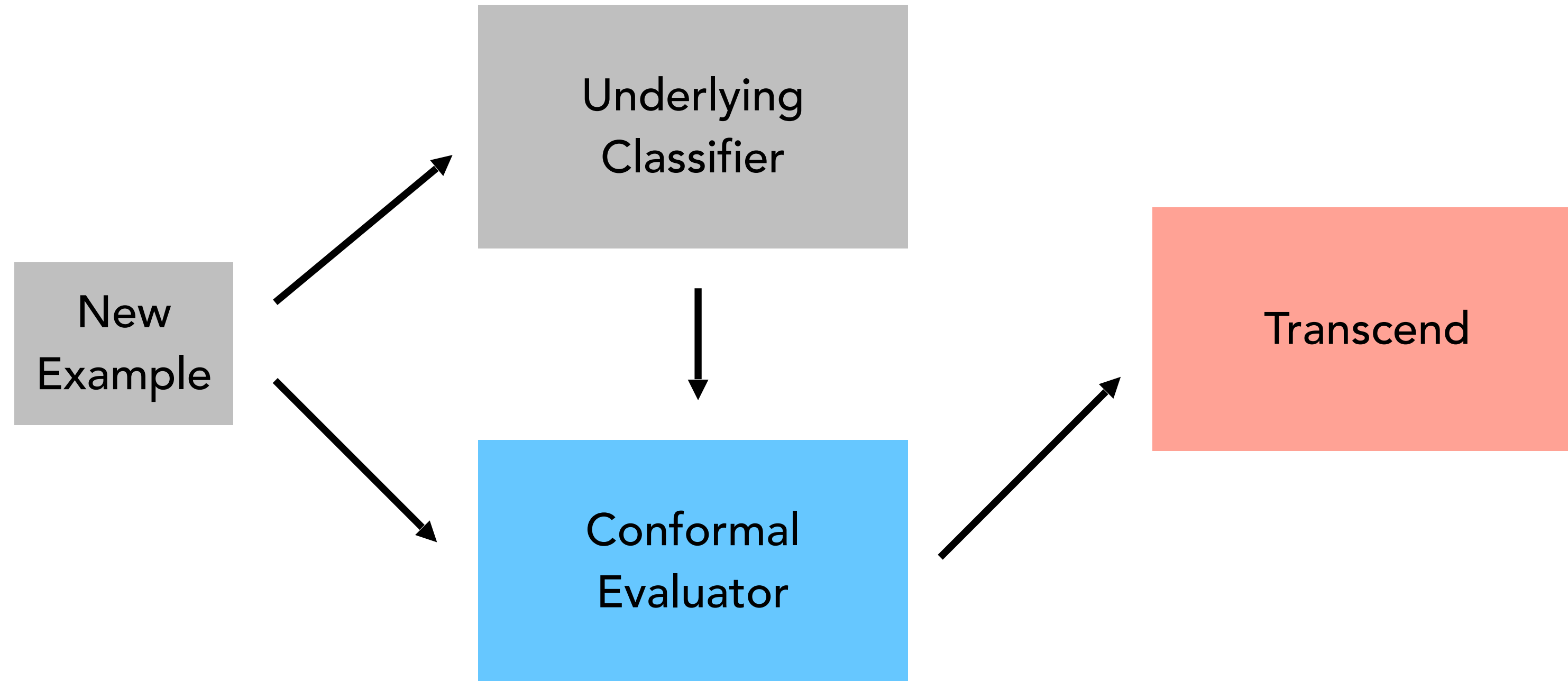
Transcend at Test Time



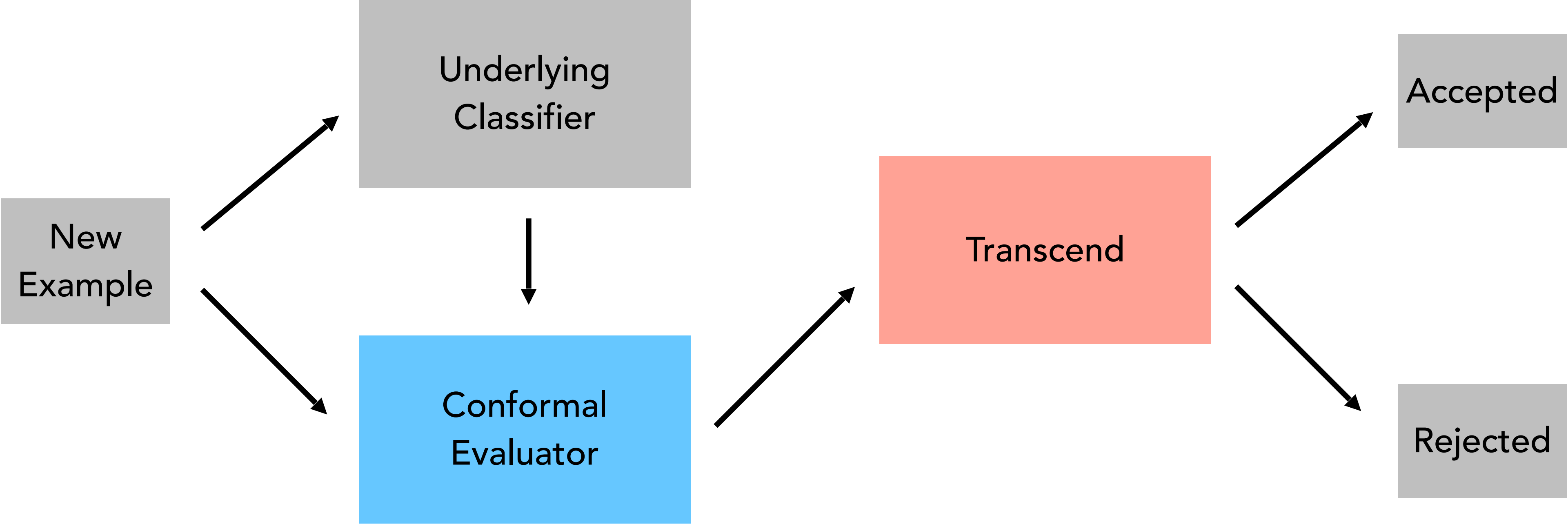
Transcend at Test Time



Transcend at Test Time



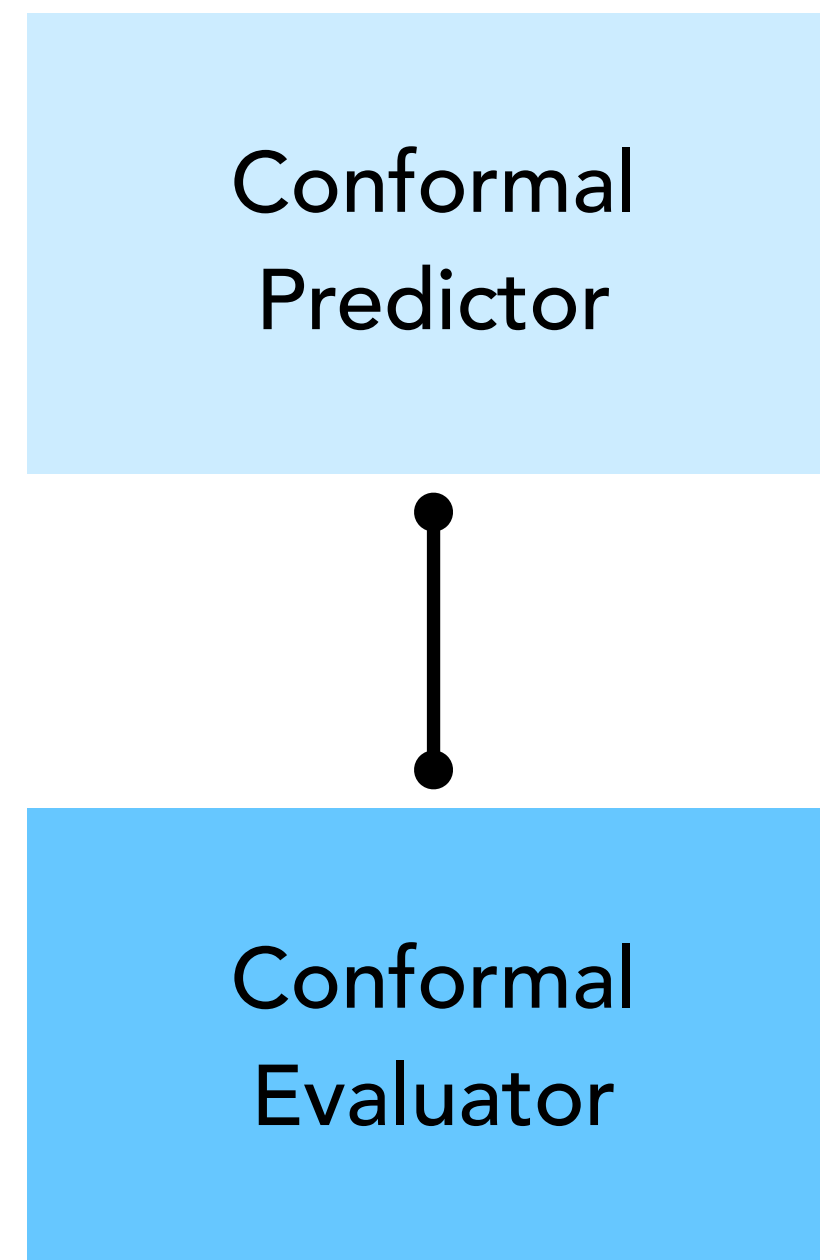
Transcend at Test Time



[USENIX Sec 2017] **Transcend: Detecting Concept Drift in Malware Classification Models**

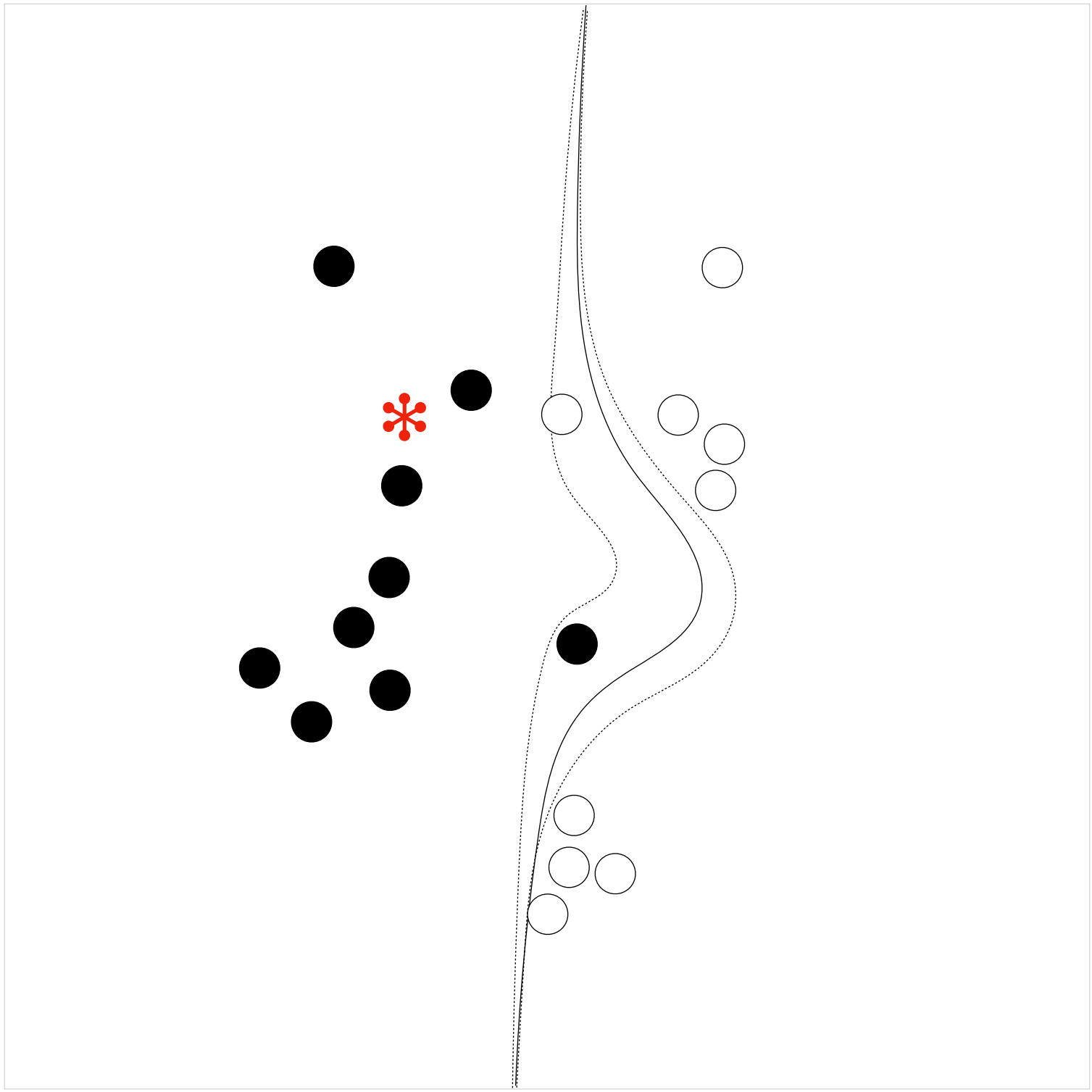
<https://s2lab.cs.ucl.ac.uk/projects/transcend/>

Conformal Prediction and Evaluation



- CP theory lays foundation for CE
- CPs outputs prediction sets with guaranteed confidence $1 - \epsilon$
- CPs rely on two assumptions:
 - **Exchangeability**: Generalization of i.i.d.
 - **Closed-world**: Fixed label space

Conformal Prediction and Non-Conformity Measure (NCM)

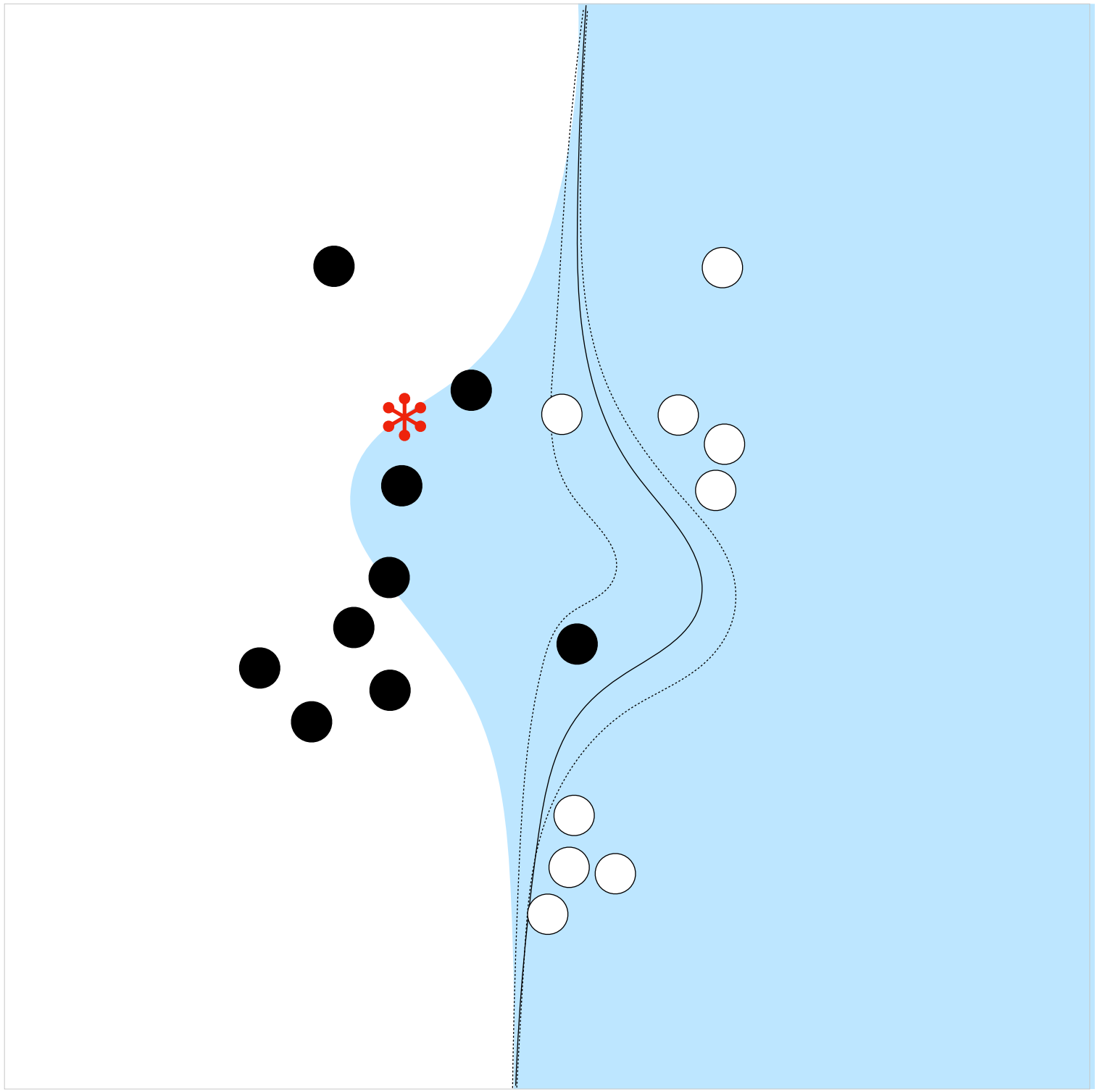


SVM Polynomial

[IEEE S&P 2022] **Transcending TRANSCEND: Revisiting Malware Classification in the Presence of Concept Drift**

<https://s2lab.cs.ucl.ac.uk/projects/transcend/>

Conformal Prediction and Non-Conformity Measure (NCM)

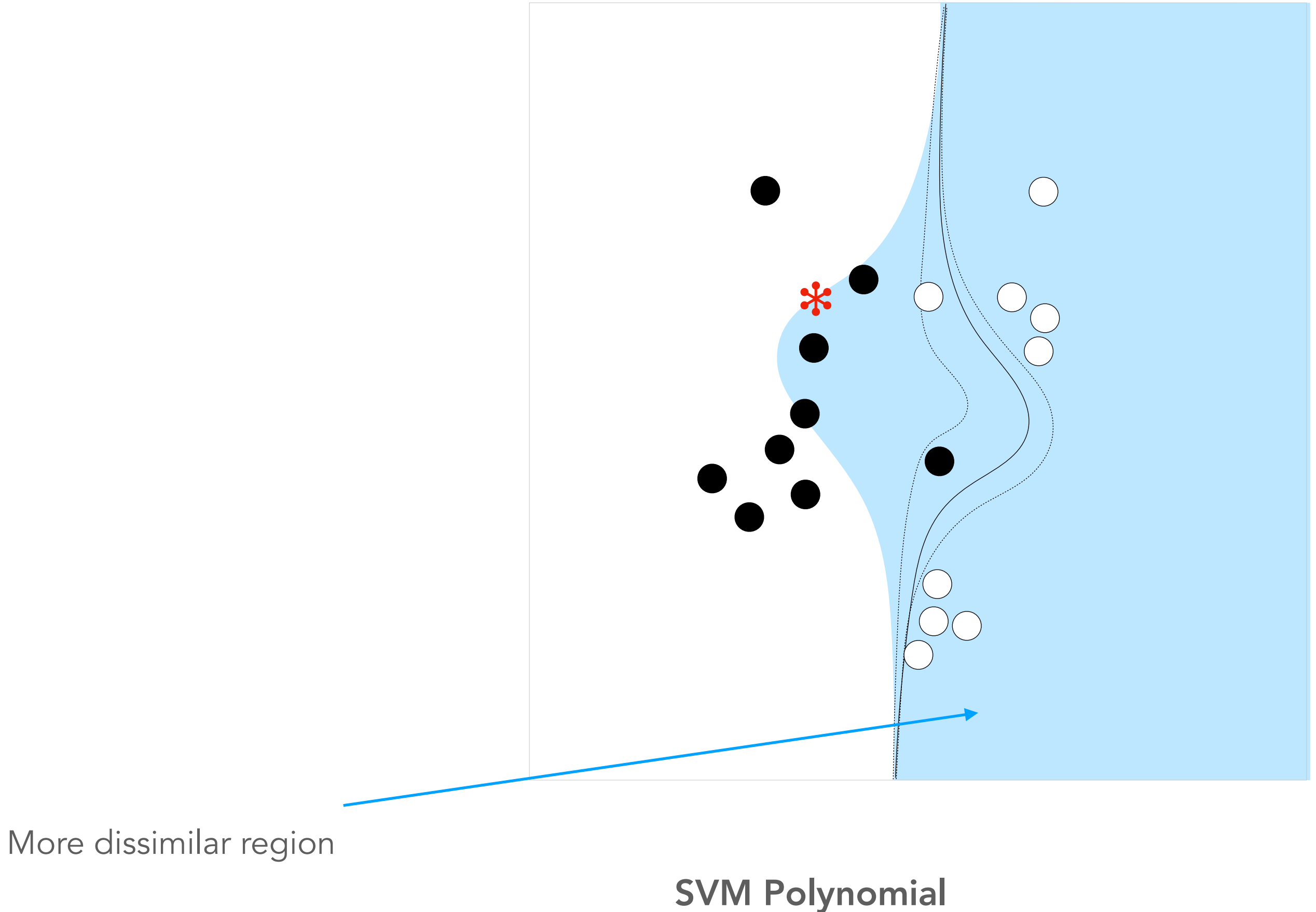


SVM Polynomial

[IEEE S&P 2022] **Transcending TRANSCEND: Revisiting Malware Classification in the Presence of Concept Drift**

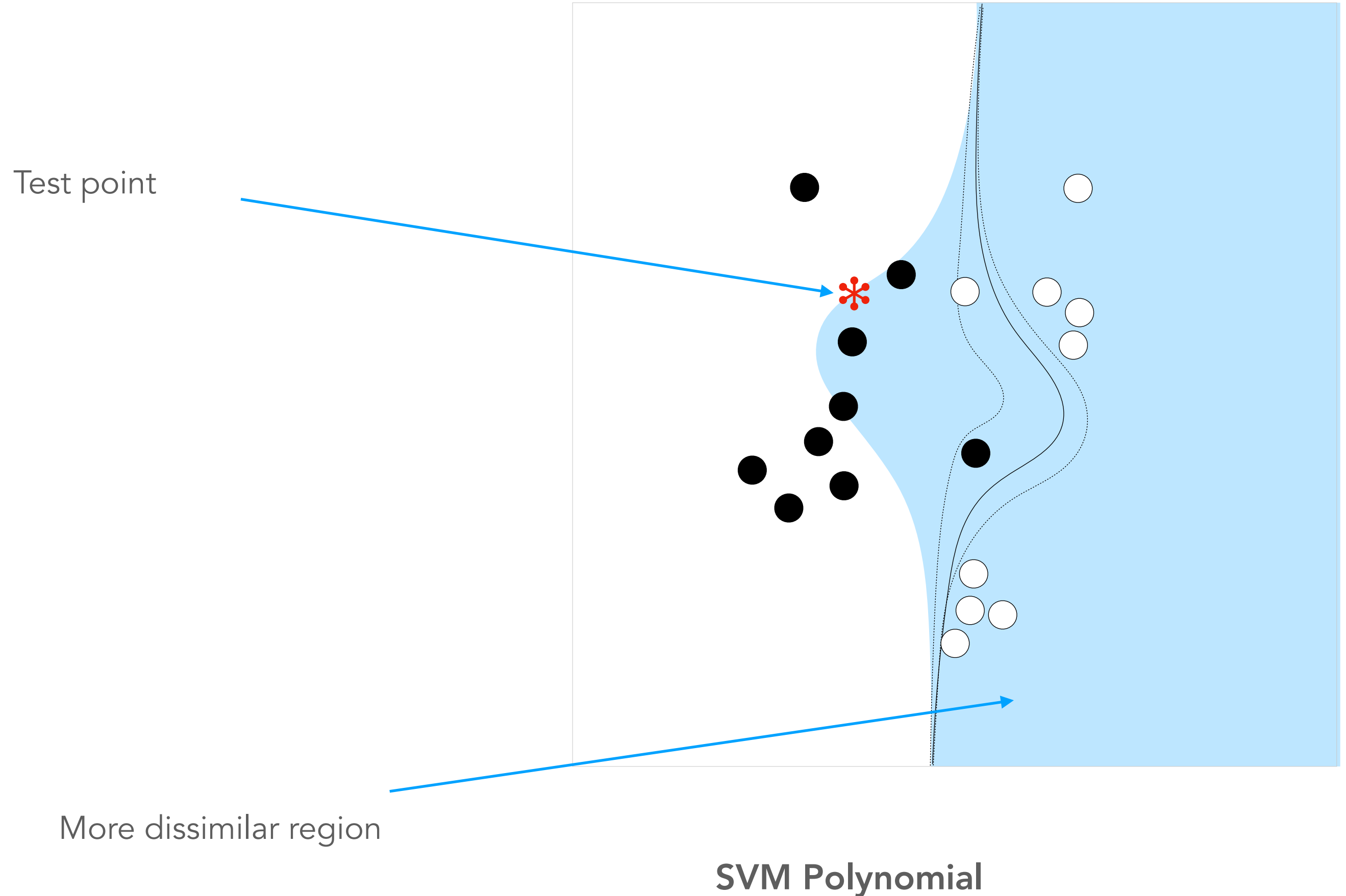
<https://s2lab.cs.ucl.ac.uk/projects/transcend/>

Conformal Prediction and Non-Conformity Measure (NCM)



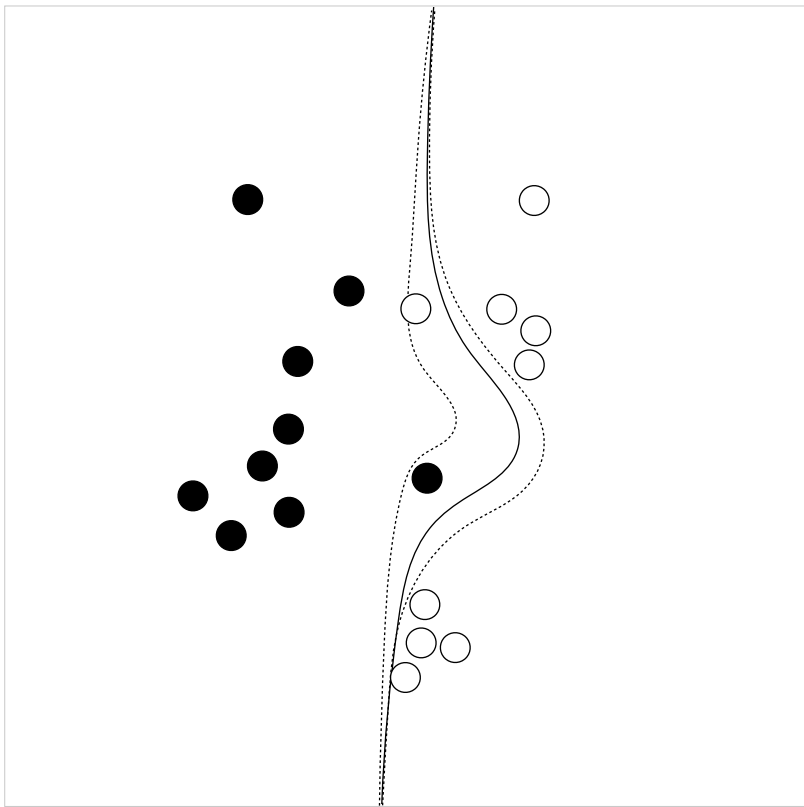
[IEEE S&P 2022] **Transcending TRANSCEND: Revisiting Malware Classification in the Presence of Concept Drift**
<https://s2lab.cs.ucl.ac.uk/projects/transcend/>

Conformal Prediction and Non-Conformity Measure (NCM)

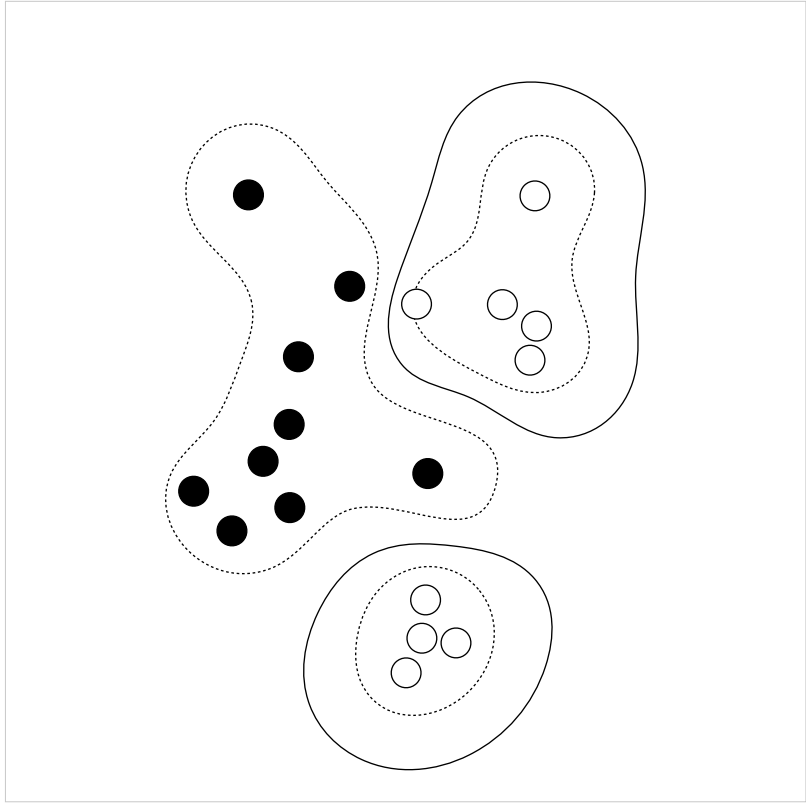


[IEEE S&P 2022] **Transcending TRANSCEND: Revisiting Malware Classification in the Presence of Concept Drift**
<https://s2lab.cs.ucl.ac.uk/projects/transcend/>

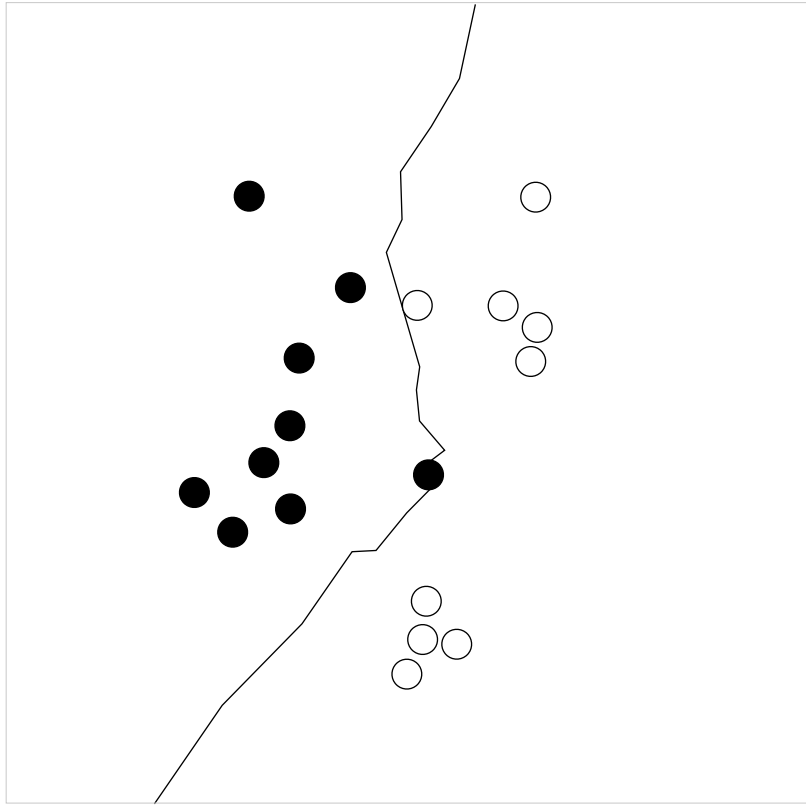
Conformal Prediction and Non-Conformity Measure (NCM)



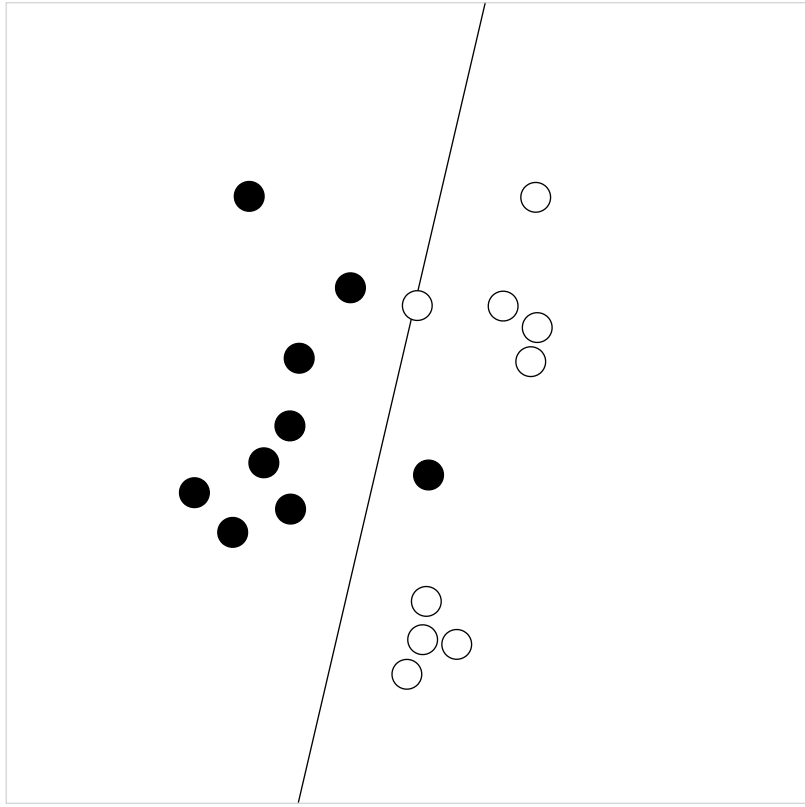
SVM Polynomial



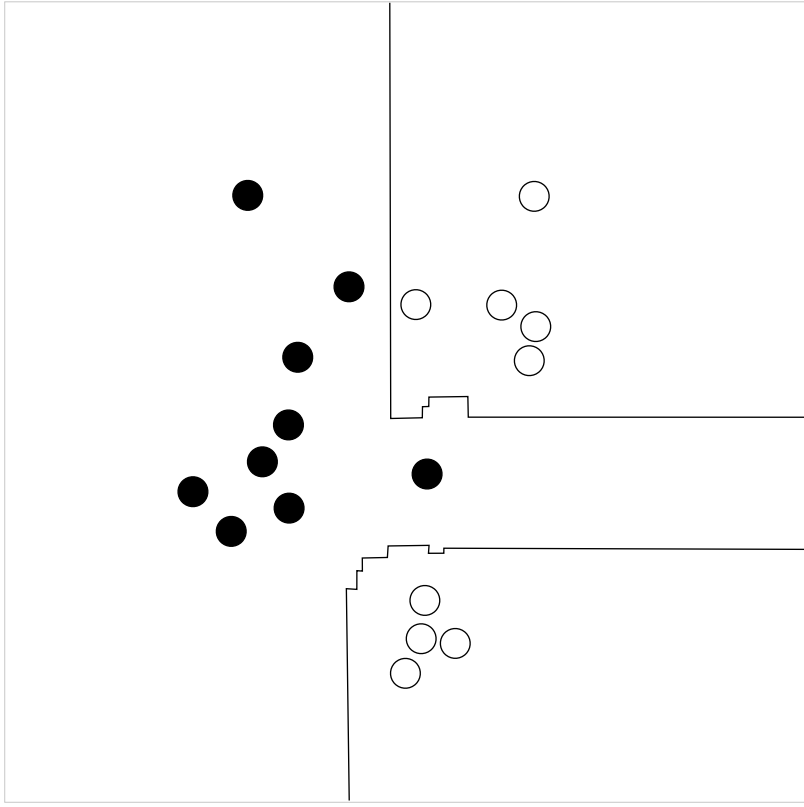
SVM RBF



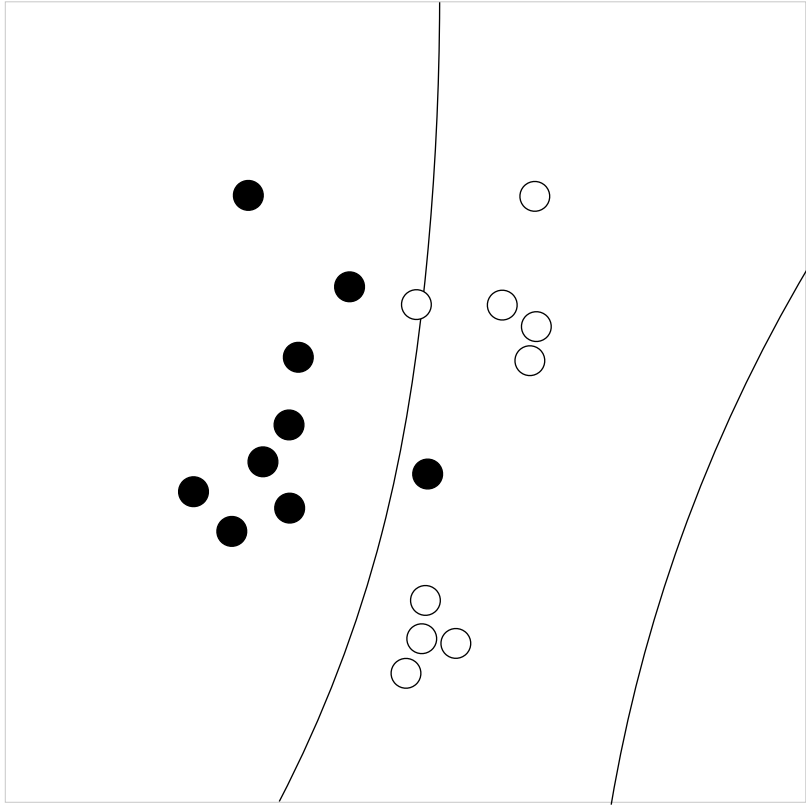
3NN



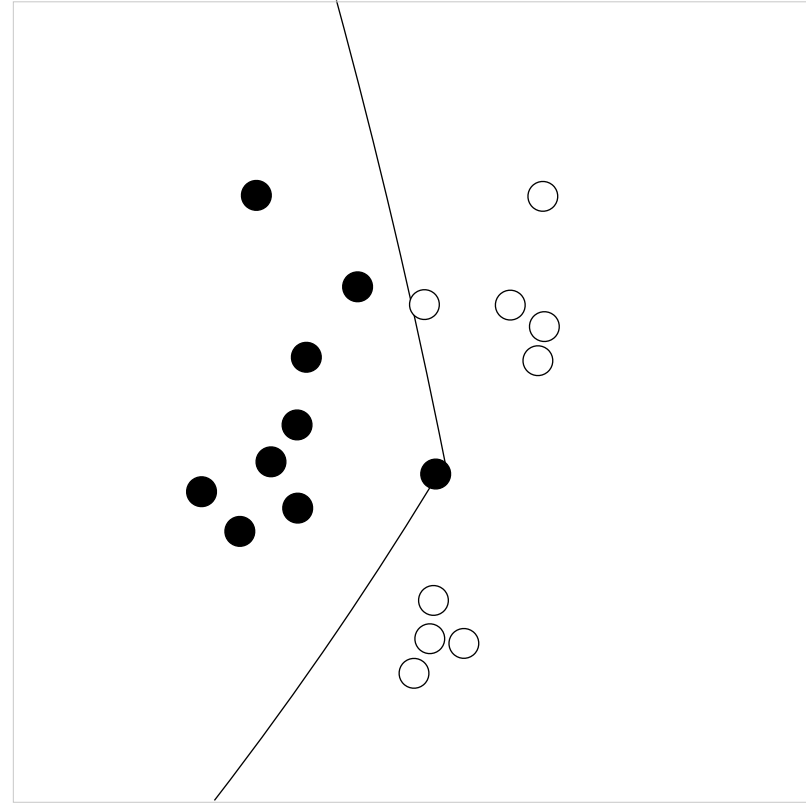
Nearest Centroid



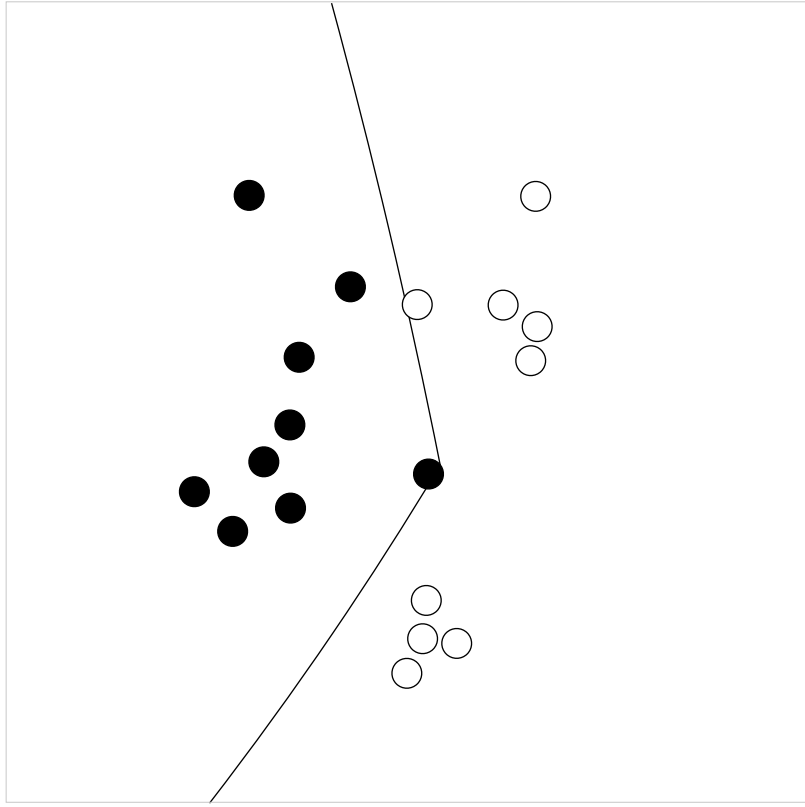
Random Forests



QDA

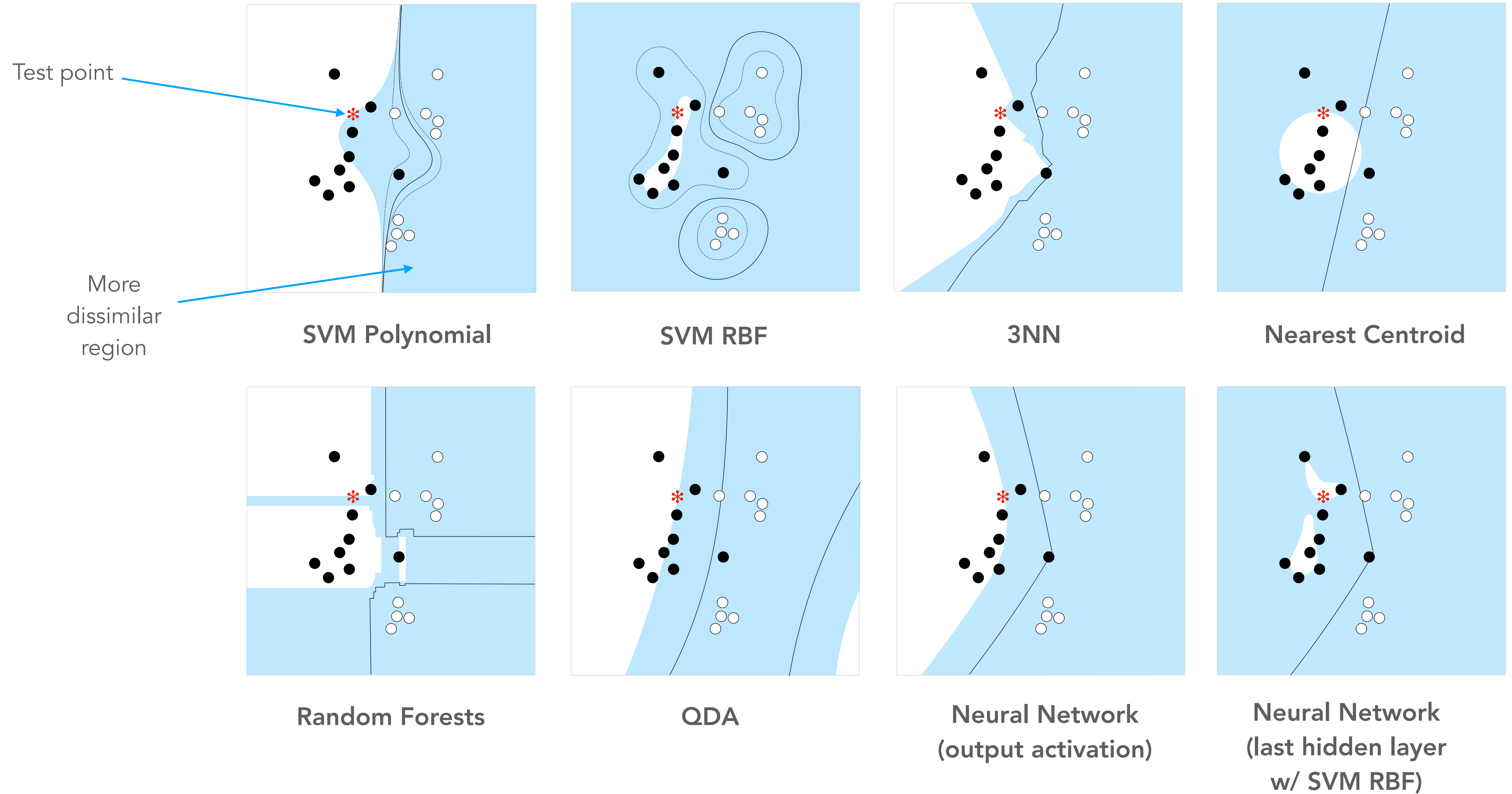


Neural Network
(output activation)



Neural Network
(last hidden layer
w/ SVM RBF)

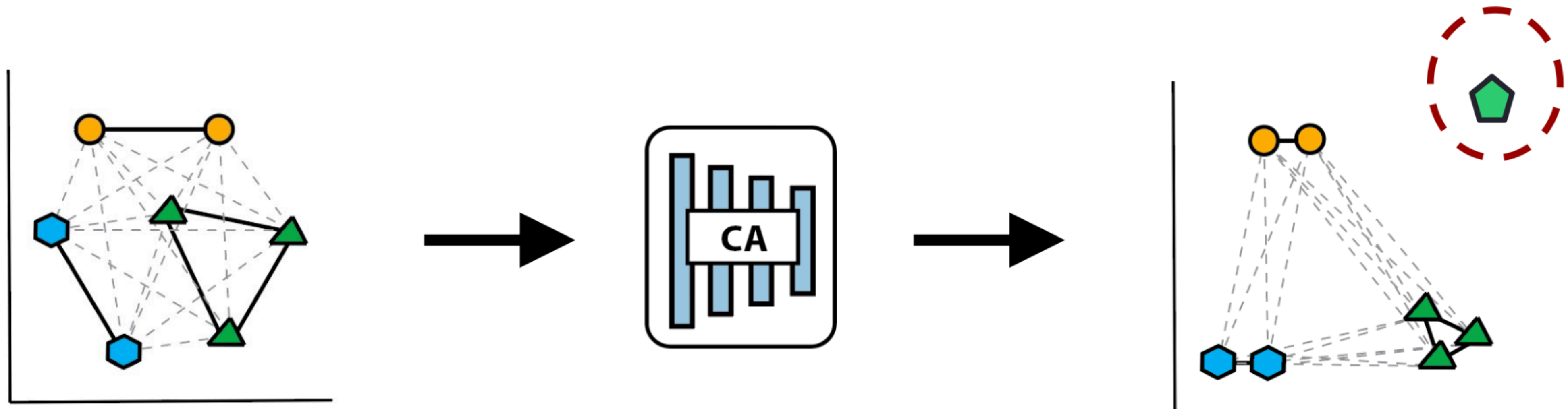
Conformal Prediction and Non-Conformity Measure (NCM)



Contrastive Learning Representation and Transcendent's NCM

Aside

- Use contrastive learning to learn a compressed representation of the training data by **contrasting with existing samples**

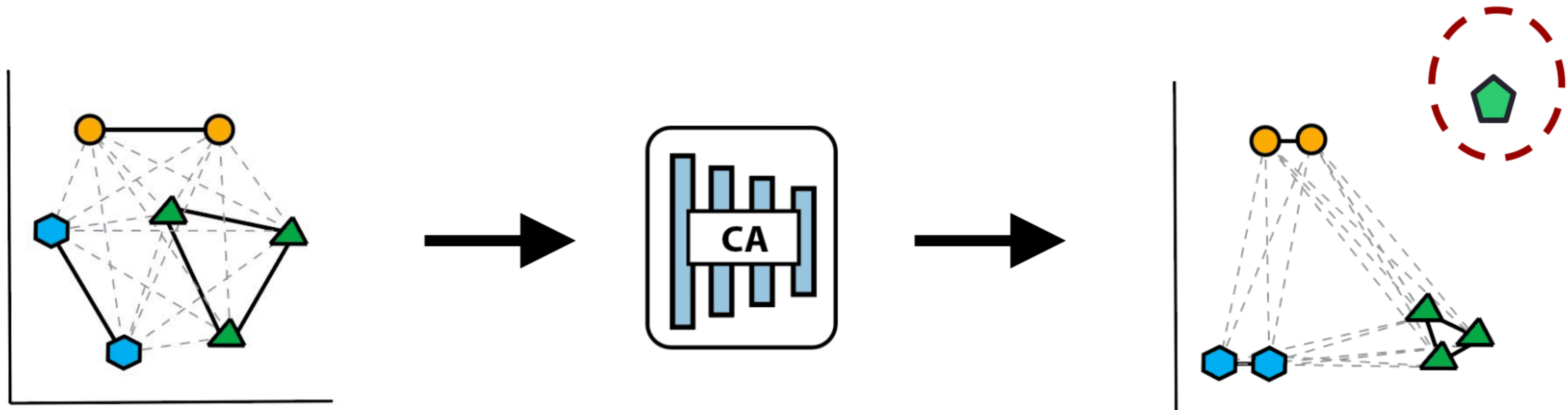


Yang et al. CADE: Detecting and Explaining Concept Drift Samples for Security Applications. USENIX Sec 2021
(Slide courtesy of Gang Wang)

Contrastive Learning Representation and Transcendent's NCM

Aside

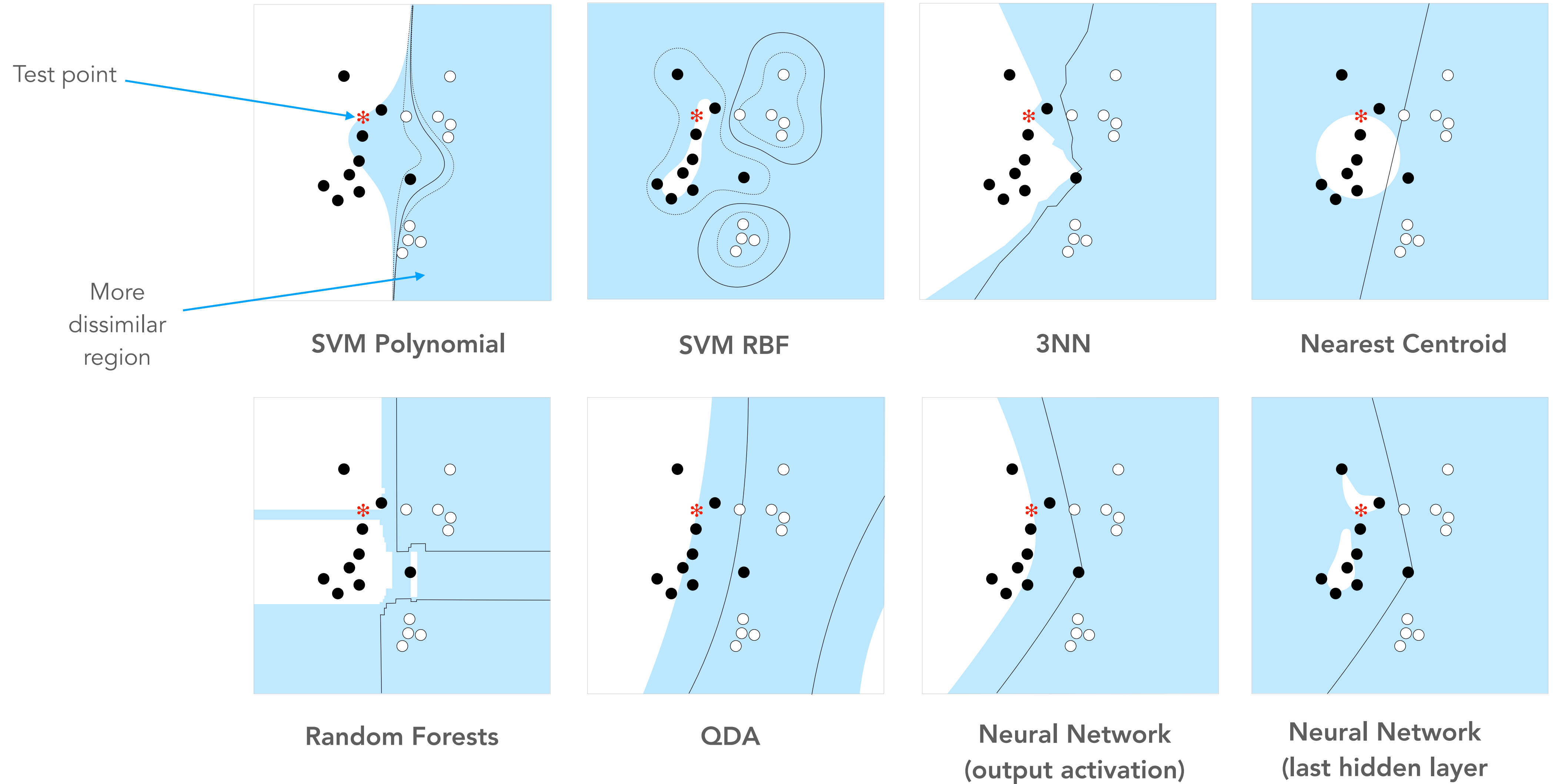
- Use contrastive learning to learn a compressed representation of the training data by **contrasting with existing samples**



Yang et al. CADE: Detecting and Explaining Concept Drift Samples for Security Applications. USENIX Sec 2021
(Slide courtesy of Gang Wang)

- CL seems ideal to be used as a non-conformity measure (!)

Conformal Prediction and Non-Conformity Measure (NCM)

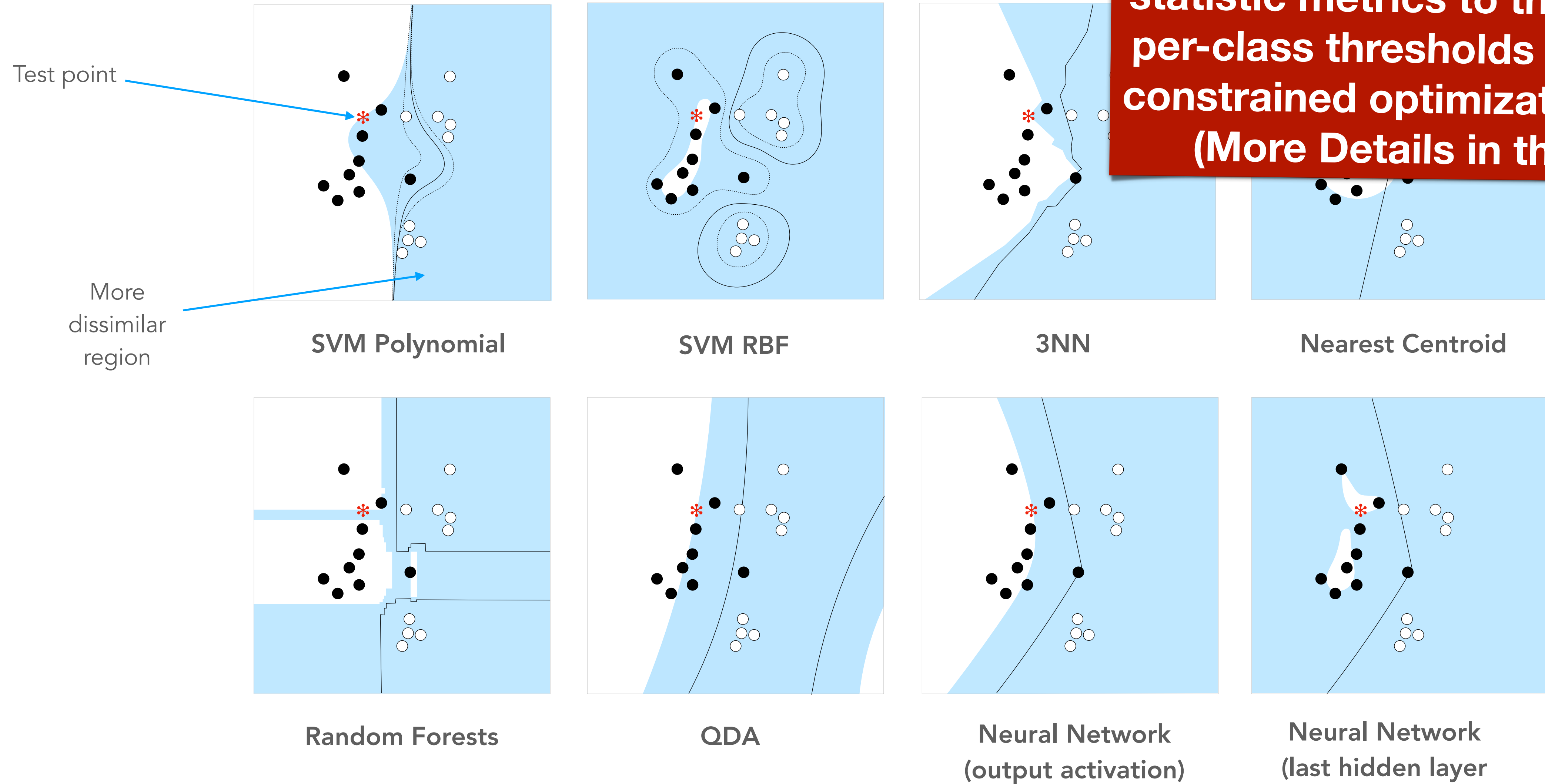


[IEEE S&P 2022] **Transcending TRANSCEND: Revisiting Malware Classification in the Presence of Concept Drift**

<https://s2lab.cs.ucl.ac.uk/projects/transcend/>

Conformal Prediction and Non-Conformity Measure (NCM)

From NCMs to p-values and statistic metrics to then compute per-class thresholds by solving a constrained optimization problem (More Details in the Paper)

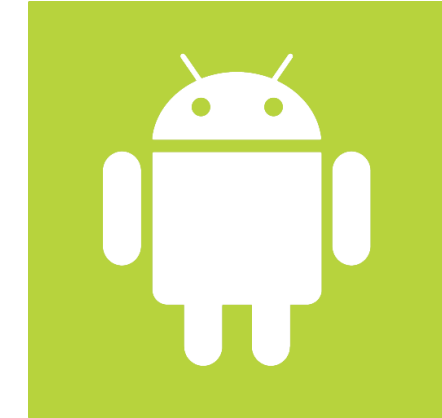


Experimental Setup

Experimental Setup

Android

- DREBIN w/ ~260K apps (Jan 2014 - Dec 2018)
- Linear SVM, binary feature space



Experimental Setup

Android

- DREBIN w/ ~260K apps (Jan 2014 - Dec 2018)
- Linear SVM, binary feature space



Windows PE

- EMBER v2 w/ ~117K apps (Jan 2017 - Dec 2017)
- Gradient Boosted Decision Tree (GBDT)



Experimental Setup

Android

- DREBIN w/ ~260K apps (Jan 2014 - Dec 2018)
- Linear SVM, binary feature space



Windows PE

- EMBER v2 w/ ~117K apps (Jan 2017 - Dec 2017)
- Gradient Boosted Decision Tree (GBDT)



PDF

- Hidost w/ ~189k apps (Aug 2017 - Sep 2017)
- Random Forest, features robust to drift



Experimental Setup

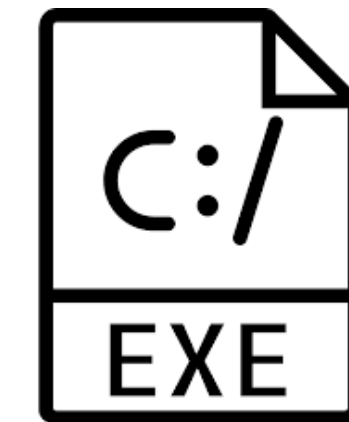
Android

- DREBIN w/ ~260K apps (Jan 2014 - Dec 2018)
- Linear SVM, binary feature space



Windows PE

- EMBER v2 w/ ~117K apps (Jan 2017 - Dec 2017)
- Gradient Boosted Decision Tree (GBDT)



PDF

- Hidost w/ ~189k apps (Aug 2017 - Sep 2017)
- Random Forest, features robust to drift



Thresholding Optimization

- Constraints: minimum F1 of 0.9 for kept elements @ rejection rate < 15%

Experimental Setup

Android

- DREBIN w/ ~260K apps (Jan 2014 - Dec 2018)
- Linear SVM, binary feature space



Windows PE

- EMBER v2 w/ ~117K apps (Jan 2017 - Dec 2017)
- Gradient Boosted Decision Tree (GBDT)



PDF

- Hidost w/ ~189k apps (Aug 2017 - Sep 2017)
- Random Forest, features robust to drift



Thresholding Optimization

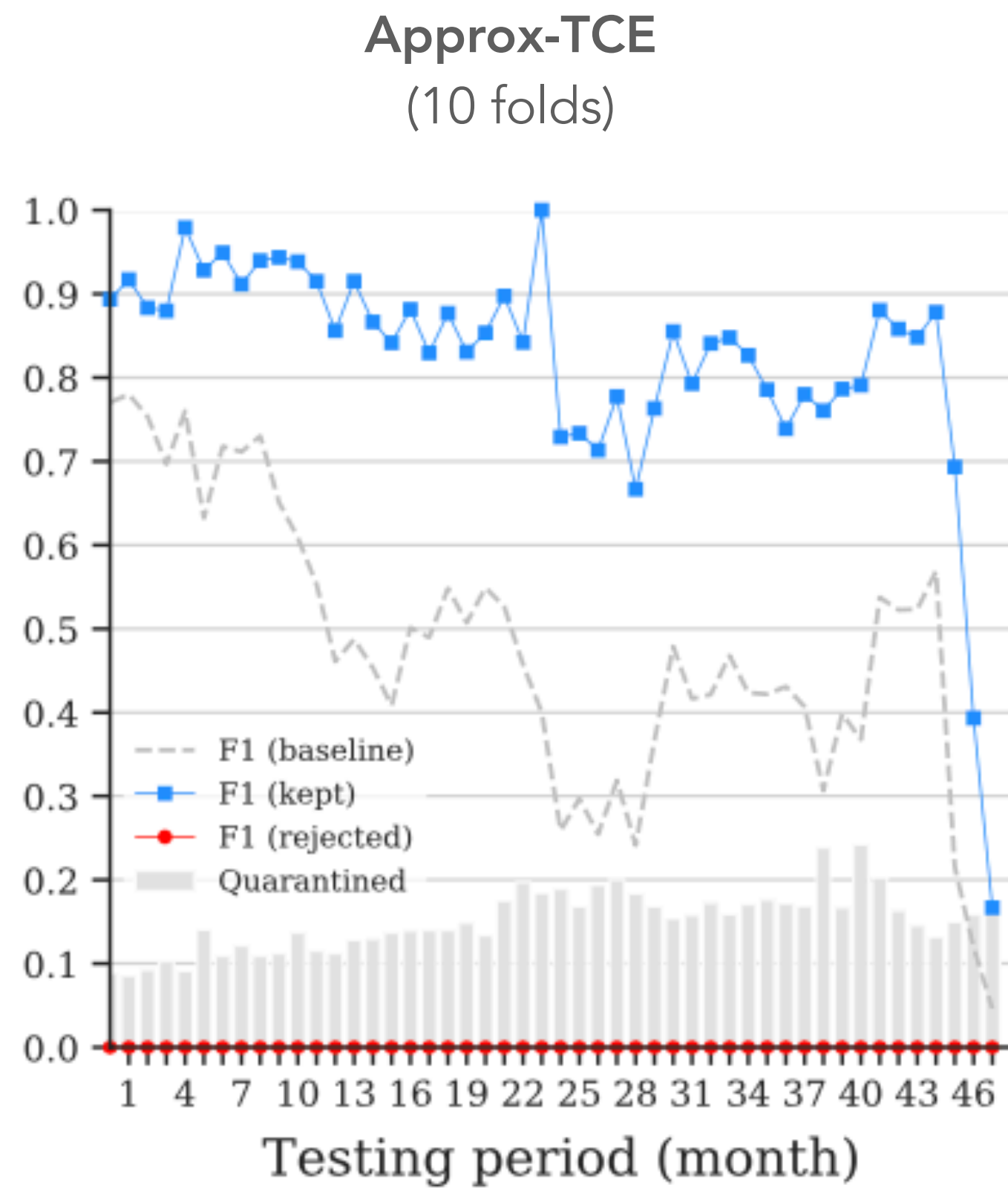
- Constraints: minimum F1 of 0.9 for kept elements @ rejection rate < 15%

Results: Rejection Performance

[IEEE S&P 2022] **Transcending TRANSCEND: Revisiting Malware Classification in the Presence of Concept Drift**

<https://s2lab.cs.ucl.ac.uk/projects/transcend/>

Results: Rejection Performance



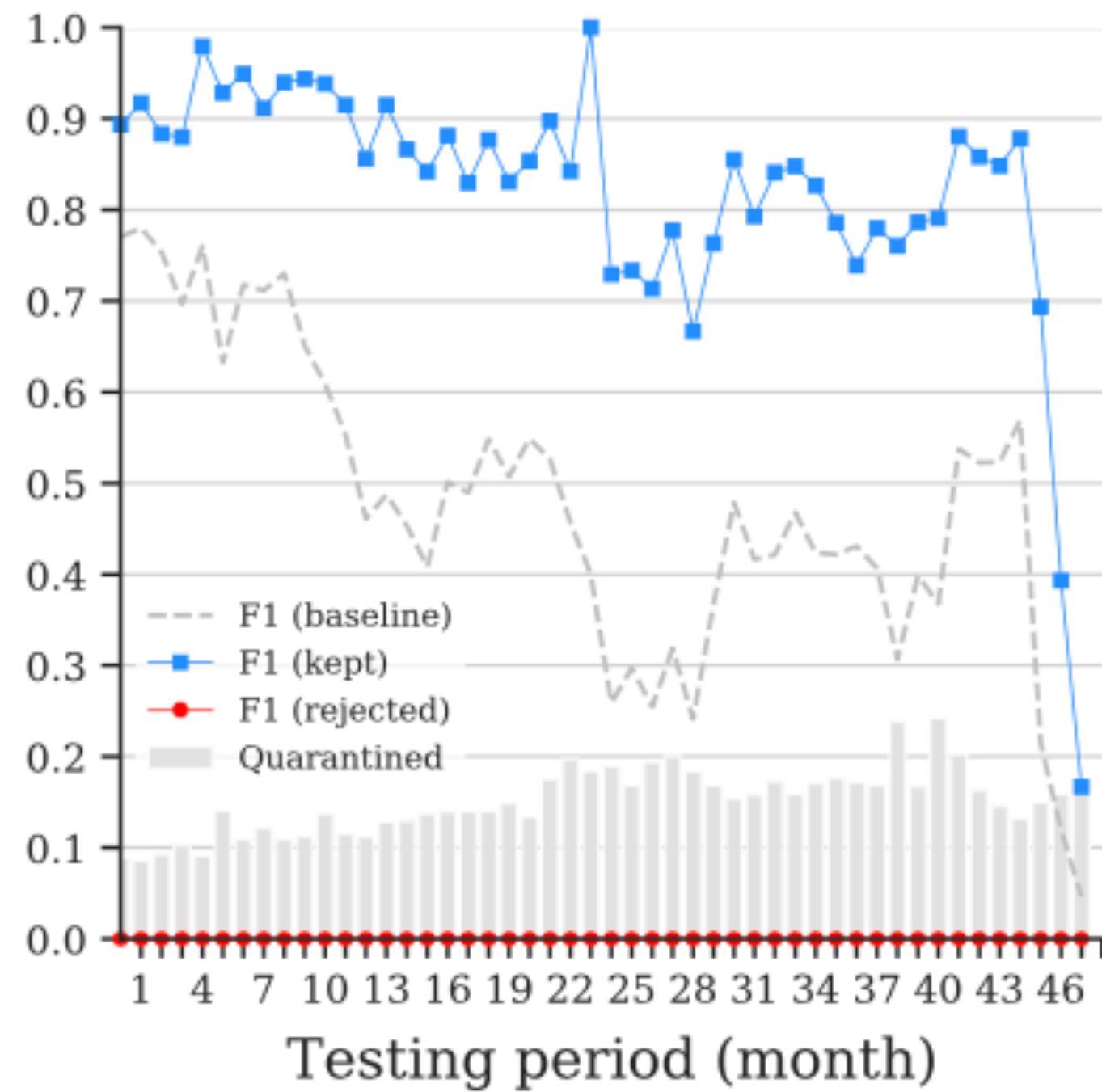
AUT(F1, kept):	0.82
AUT(F1, rejected):	0.00
CPU hours:	46.1

[IEEE S&P 2022] **Transcending TRANSCEND: Revisiting Malware Classification in the Presence of Concept Drift**

<https://s2lab.cs.ucl.ac.uk/projects/transcend/>

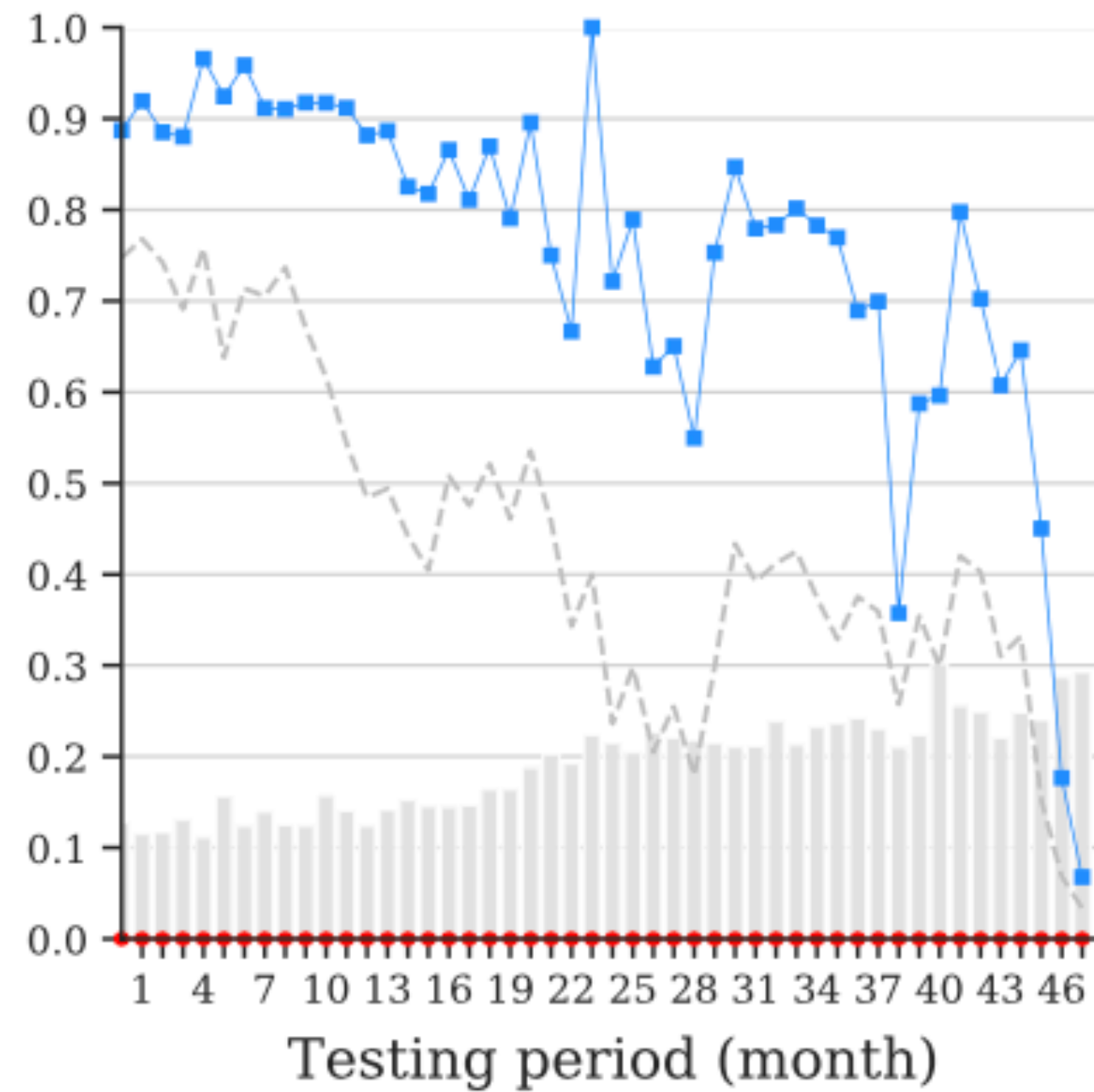
Results: Rejection Performance

Approx-TCE
(10 folds)



AUT(F1, kept):	0.82
AUT(F1, rejected):	0.00
CPU hours:	46.1

ICE
(0.33 calibration split)



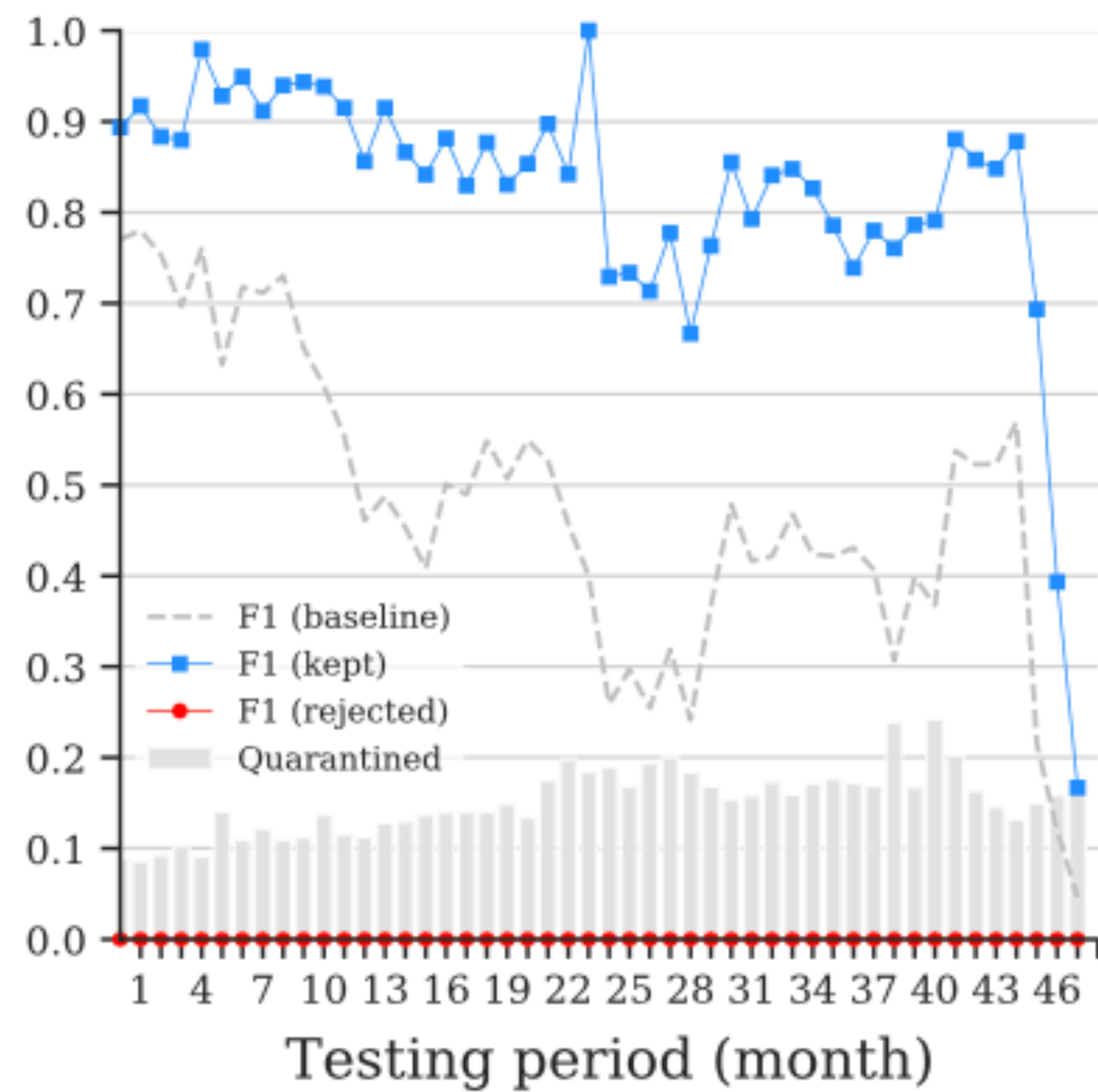
AUT(f1, kept):	0.76
AUT(f1, rejected):	0.00
CPU hours:	11.5

[IEEE S&P 2022] **Transcending TRANSCEND: Revisiting Malware Classification in the Presence of Concept Drift**

<https://s2lab.cs.ucl.ac.uk/projects/transcend/>

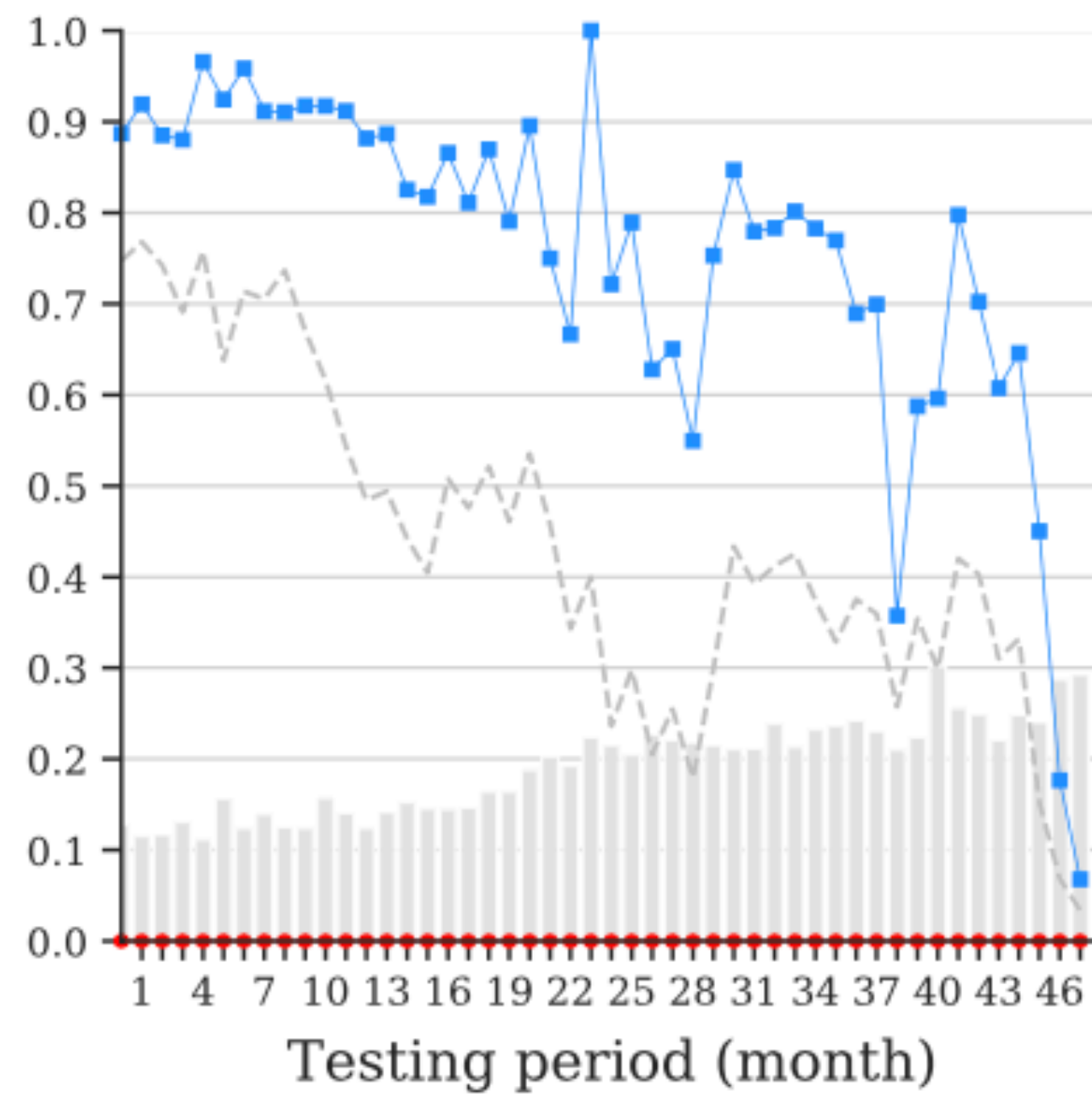
Results: Rejection Performance

Approx-TCE
(10 folds)



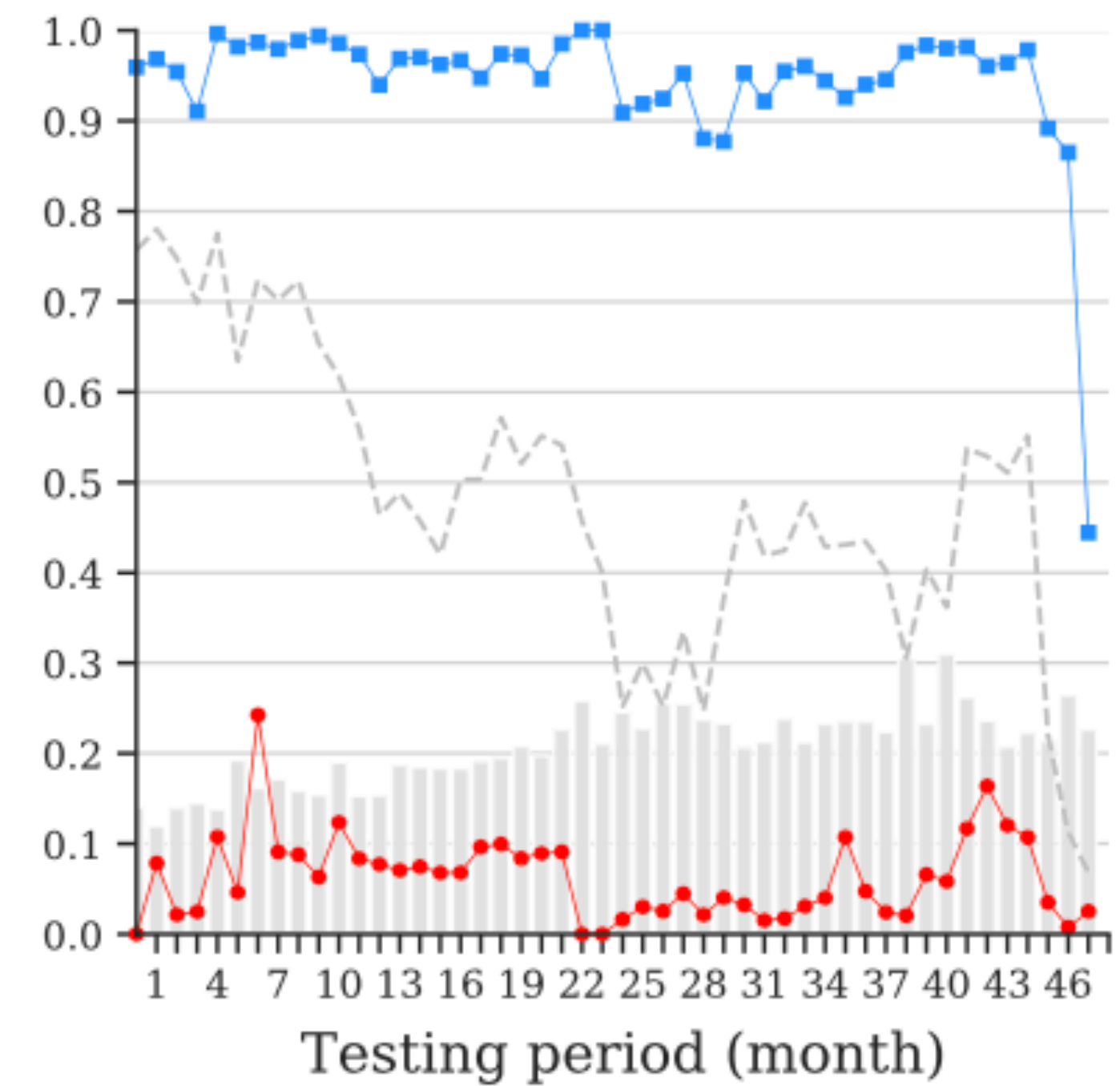
AUT(F1, kept):	0.82
AUT(F1, rejected):	0.00
CPU hours:	46.1

ICE
(0.33 calibration split)



AUT(f1, kept):	0.76
AUT(f1, rejected):	0.00
CPU hours:	11.5

CCE
(10 folds)



AUT(f1, kept):	0.94
AUT(f1, rejected):	0.06
CPU hours:	35.6

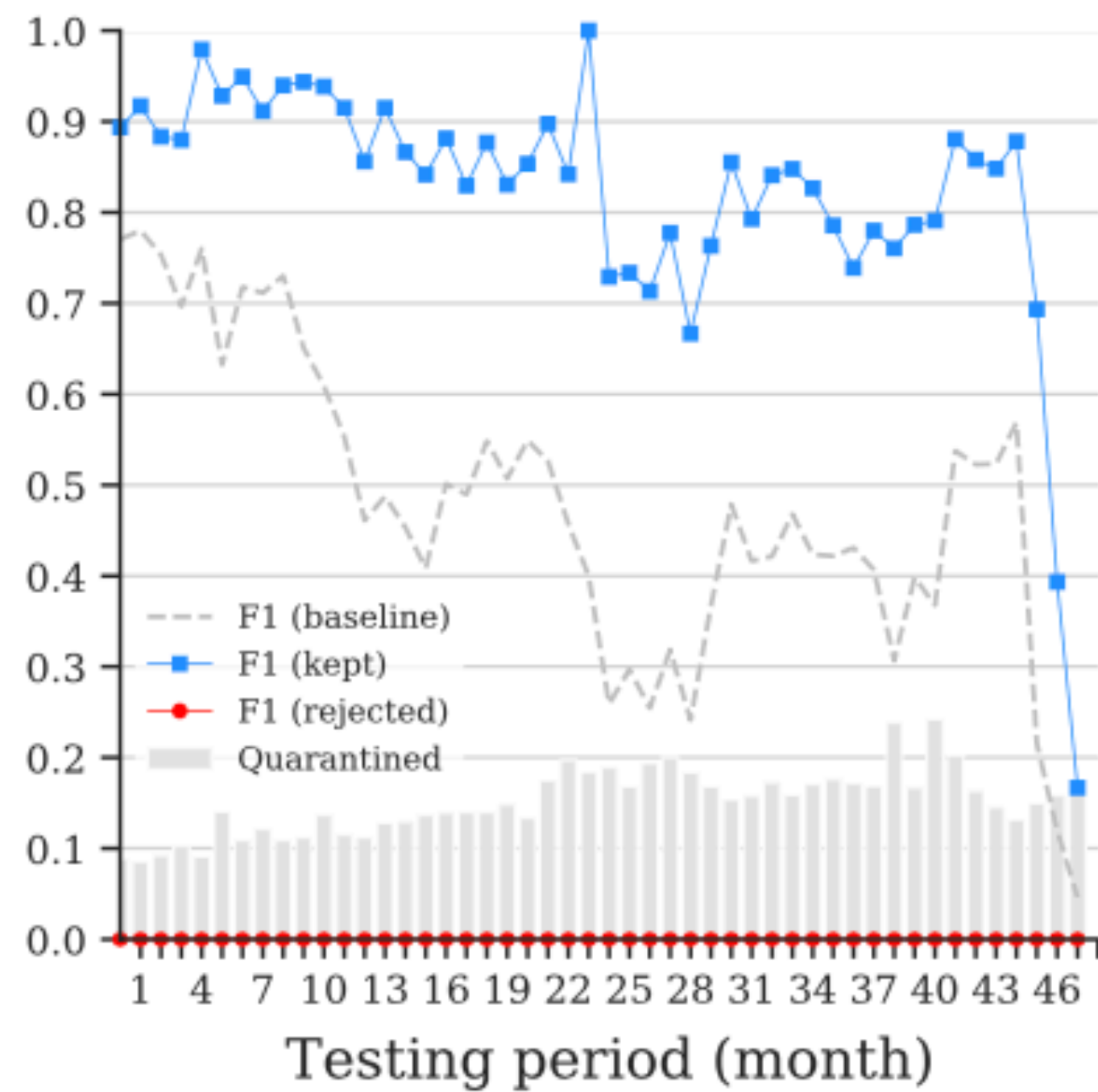
[IEEE S&P 2022] **Transcending TRANSCEND: Revisiting Malware Classification in the Presence of Concept Drift**

<https://s2lab.cs.ucl.ac.uk/projects/transcend/>

Results: Rejection Performance

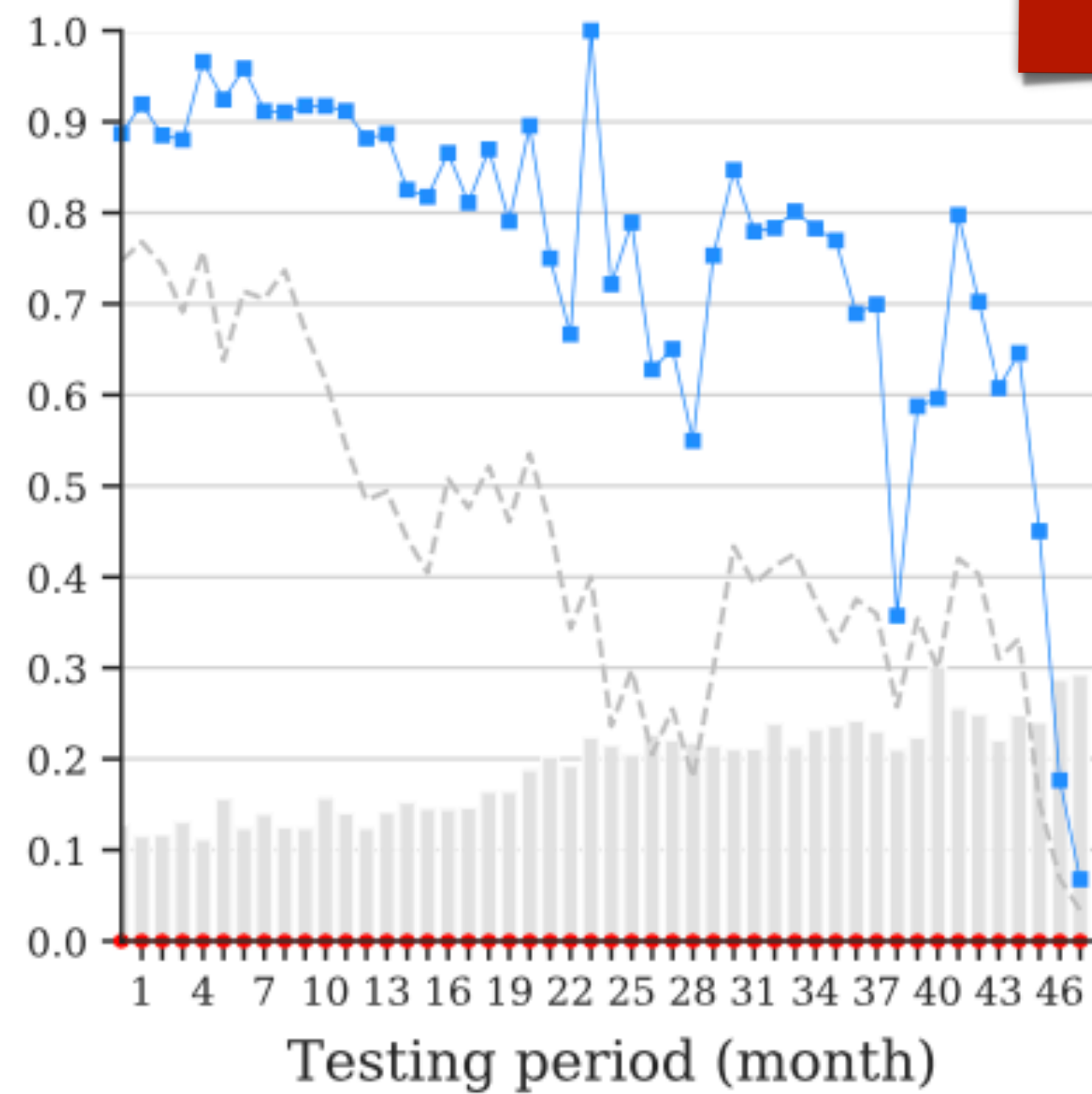
More experiments on different datasets, SOTA, and optimization in the paper

Approx-TCE
(10 folds)

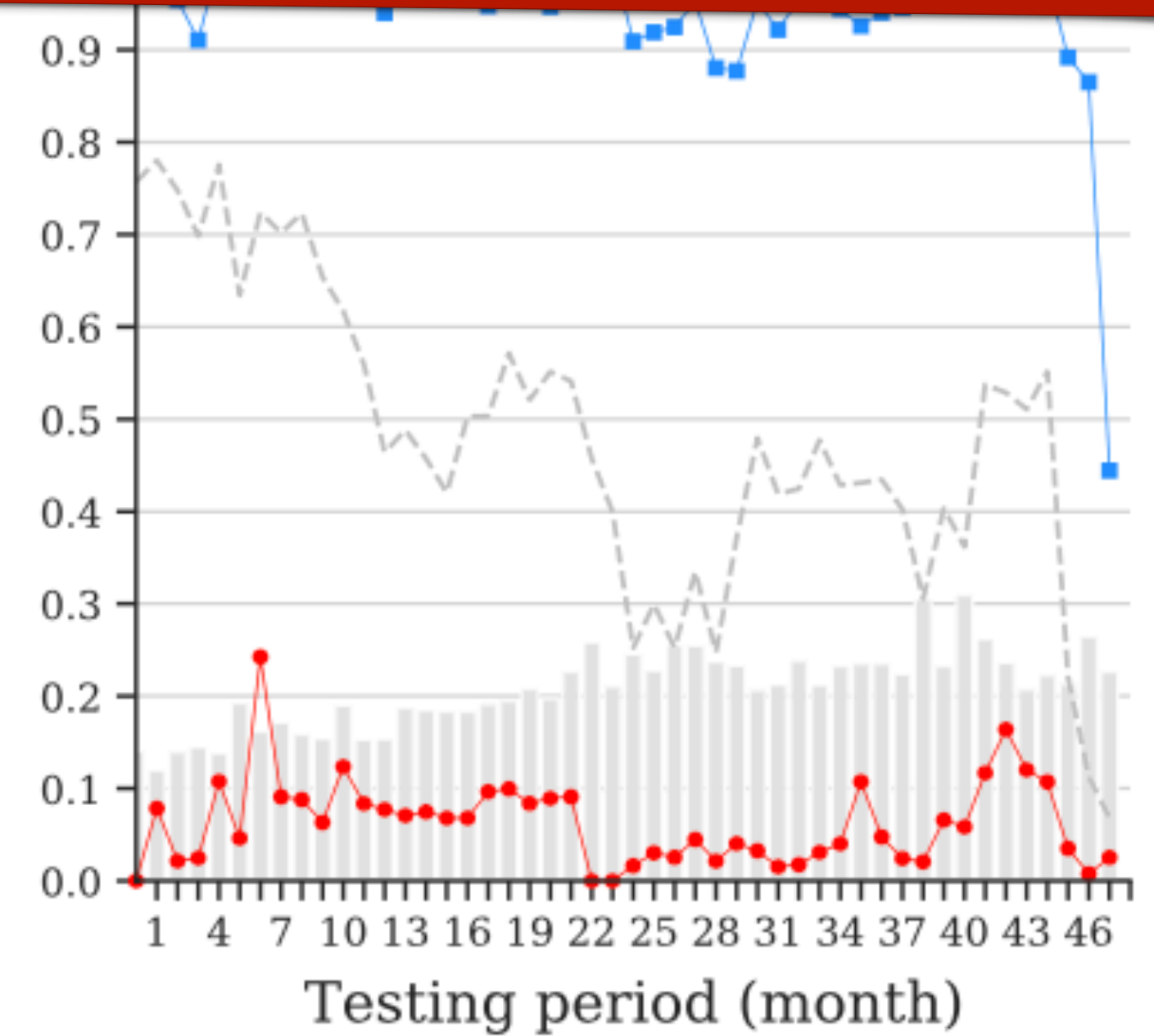


AUT(F1, kept):	0.82
AUT(F1, rejected):	0.00
CPU hours:	46.1

ICE
(0.33 calibration split)



AUT(f1, kept):	0.76
AUT(f1, rejected):	0.00
CPU hours:	11.5



AUT(f1, kept):	0.94
AUT(f1, rejected):	0.06
CPU hours:	35.6

Reproducibility Crisis

NEWS

Home | UK | World | Business | Politics | Tech | Science | Health | Family & Education | Entertainment & Arts | Stories | More ▾

Science & Environment

February 2019

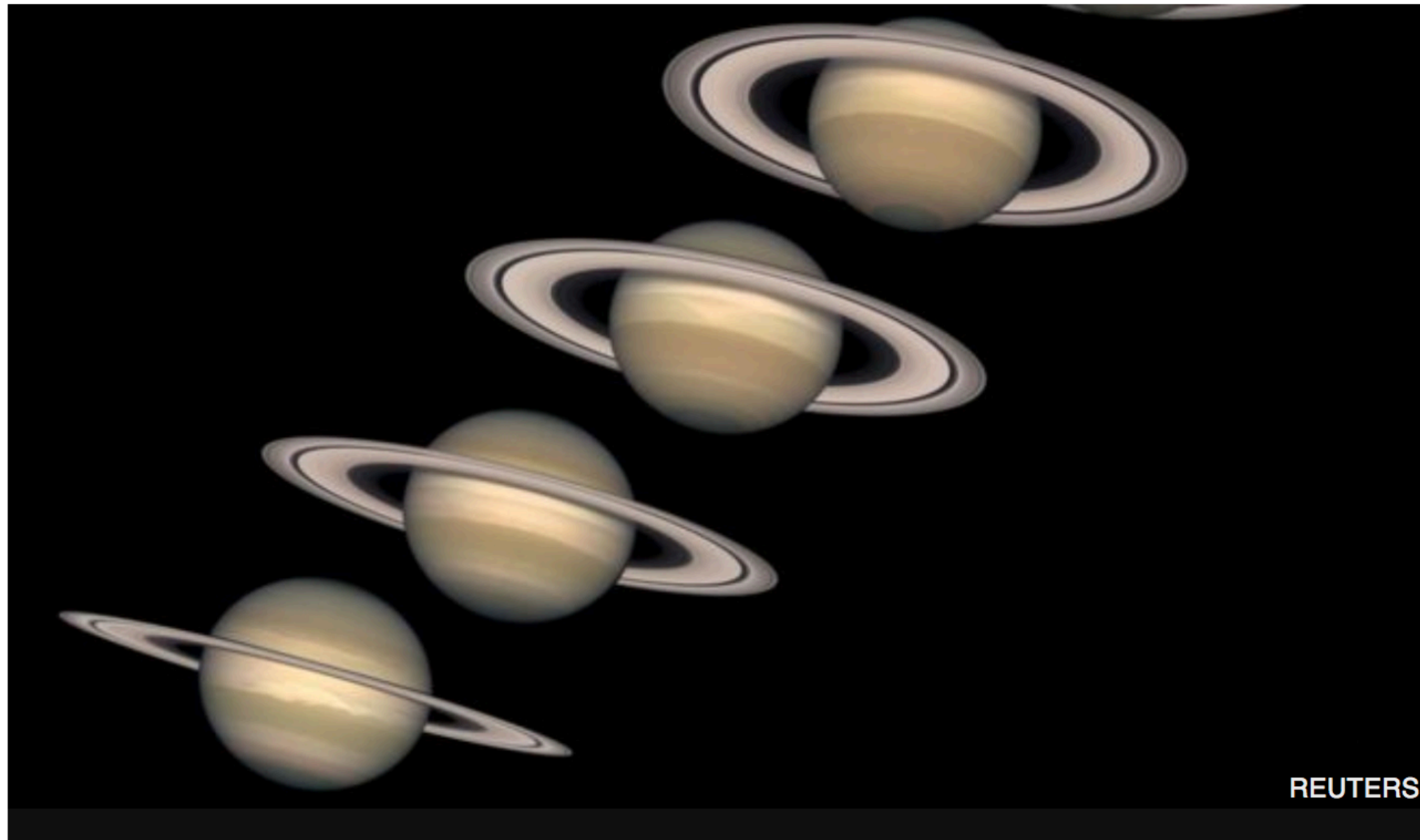
AAAS: Machine learning 'causing science crisis'

By Pallab Ghosh
Science correspondent, BBC News, Washington

🕒 16 February 2019

[f](#) [m](#) [t](#) [e](#) [Share](#)

AAAS meeting



Top Stories

MPs demand 'urgent' Facebook regulation

The House of Commons publishes its report into fake news with some strong criticism of Facebook.

🕒 7 hours ago

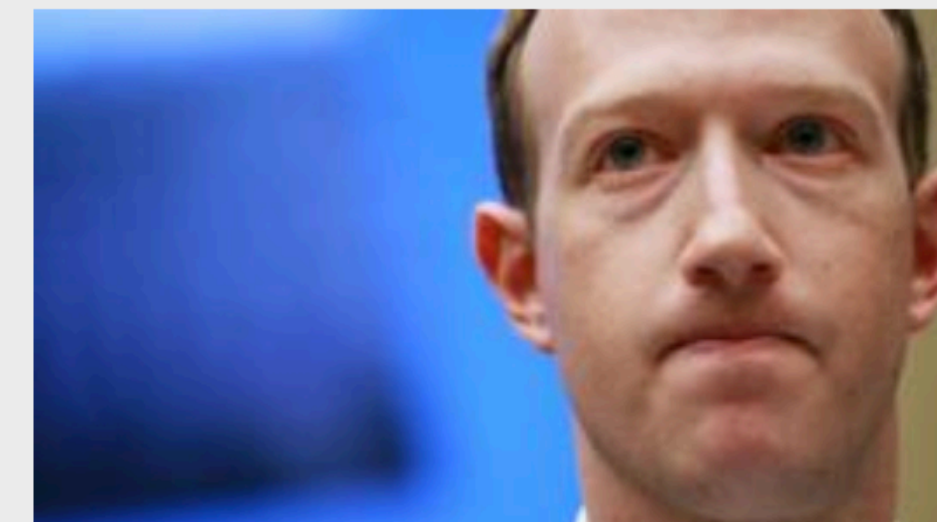
UK cyber-security: Huawei risk manageable

🕒 58 minutes ago

Several Labour MPs 'set to quit party'

🕒 2 hours ago

Features



In a hurry? Here's what you need

NEWS

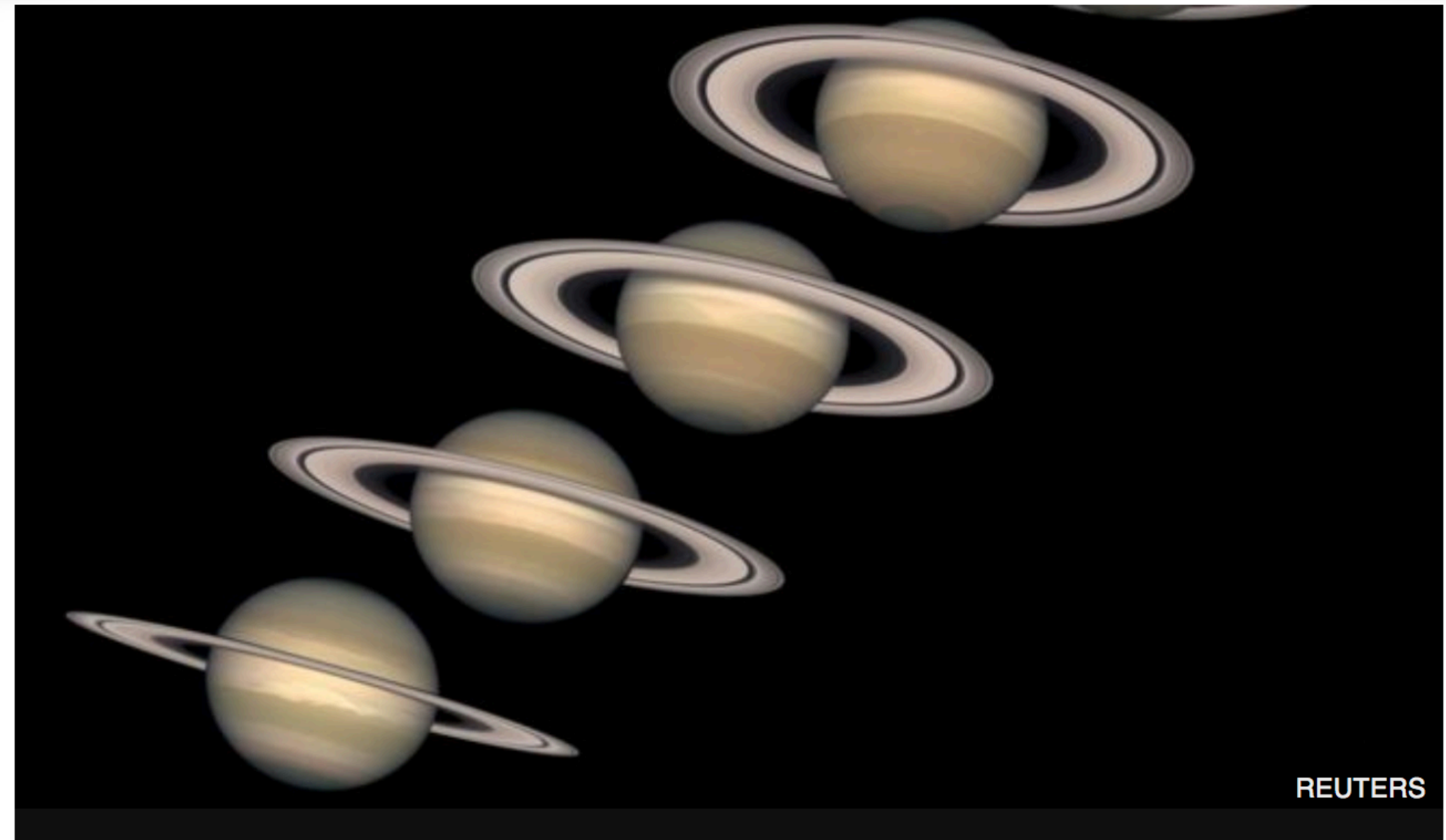
Home | UK | World | Business | Politics | Tech | **Science** | Health | Family & Education | Entertainment & Arts | Stories | More ▾

Science & Environment

February 2019

AAAS: Machine learning 'causing science crisis'

Machine learning techniques used by thousands of scientists to analyse data are producing results that are misleading and often completely wrong



Top Stories

MPs demand 'urgent' Facebook regulation

Commons publishes its news with some strong facebook.

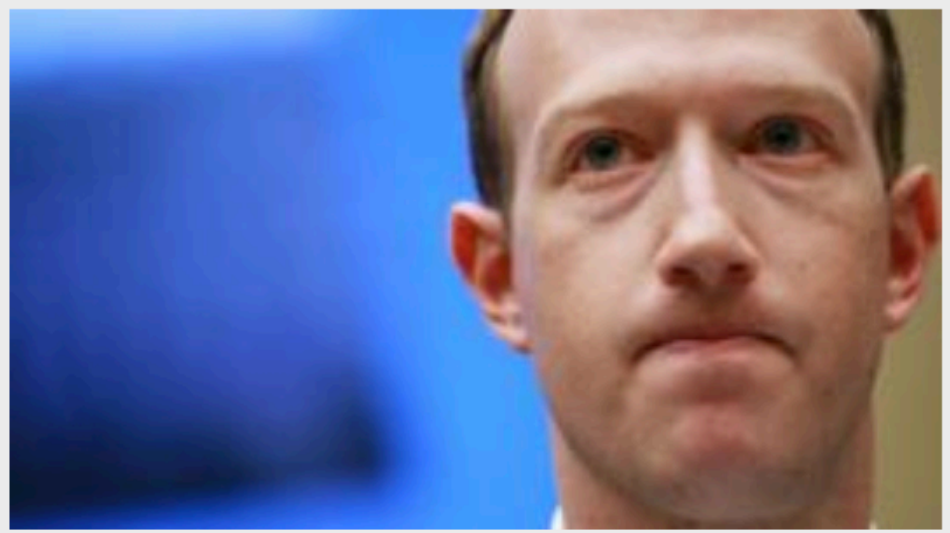
Security: Huawei risk

58 minutes ago

Several Labour MPs 'set to quit party'

2 hours ago

Features



In a hurry? Here's what you need

Leakage and the Reproducibility Crisis in ML-based Science



Draft paper

July '22 online workshop

We argue that there is a reproducibility crisis in ML-based science. We compile evidence of this crisis across fields, identify data leakage as a pervasive cause of reproducibility failures, conduct our own reproducibility investigations using in-depth code-review, and propose a solution.

[Top](#)

[List of failures](#)

[Taxonomy](#)

[Model info sheets](#)

[Case study](#)

[Terminology](#)

[Citation](#)

[About us](#)

Context

Many quantitative science fields are [adopting](#) the paradigm of predictive modeling using machine learning. We welcome this development. At the same time, as researchers whose interests include the strengths and limits of machine learning, we have concerns about reproducibility and overoptimism.

There are many reasons for caution:

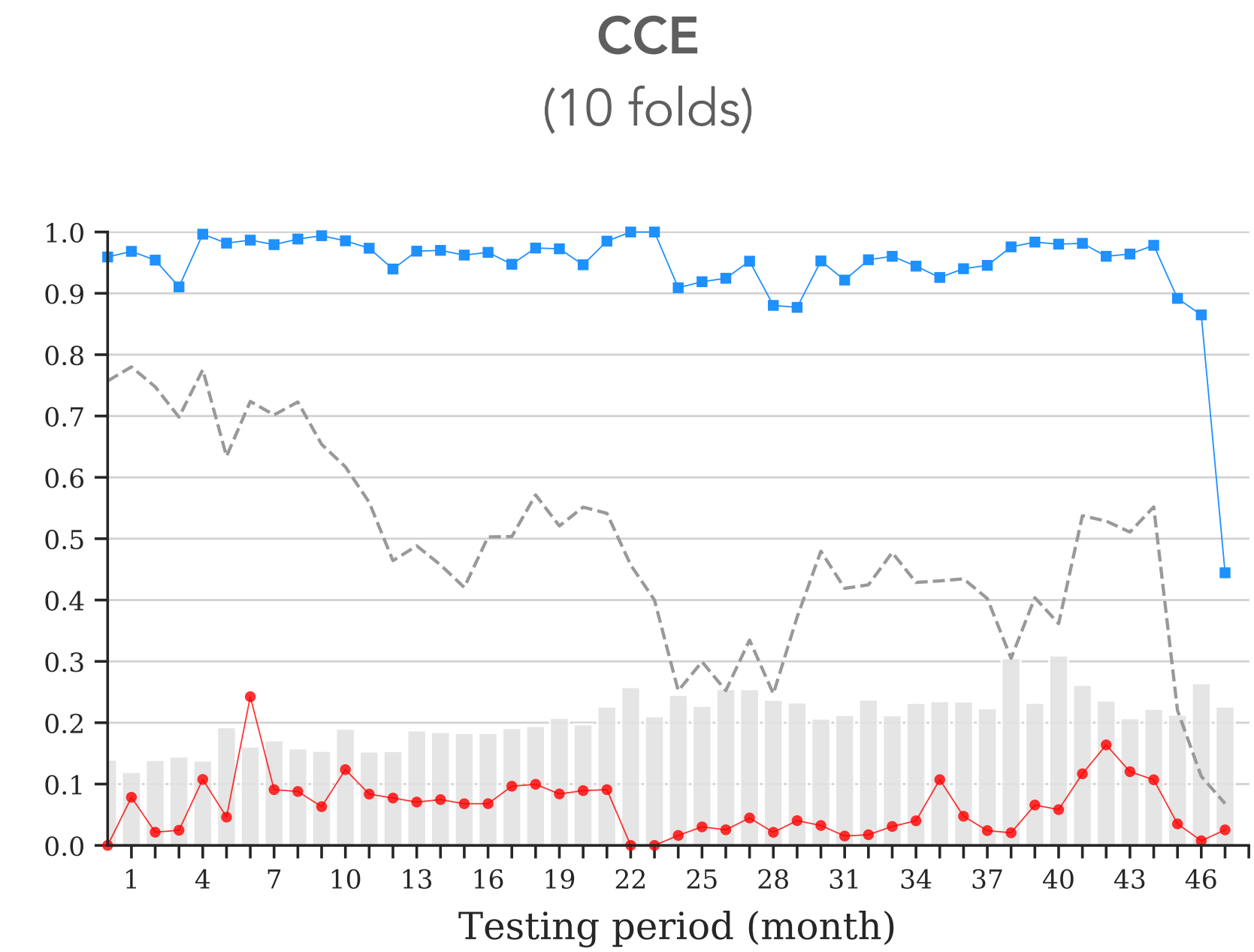
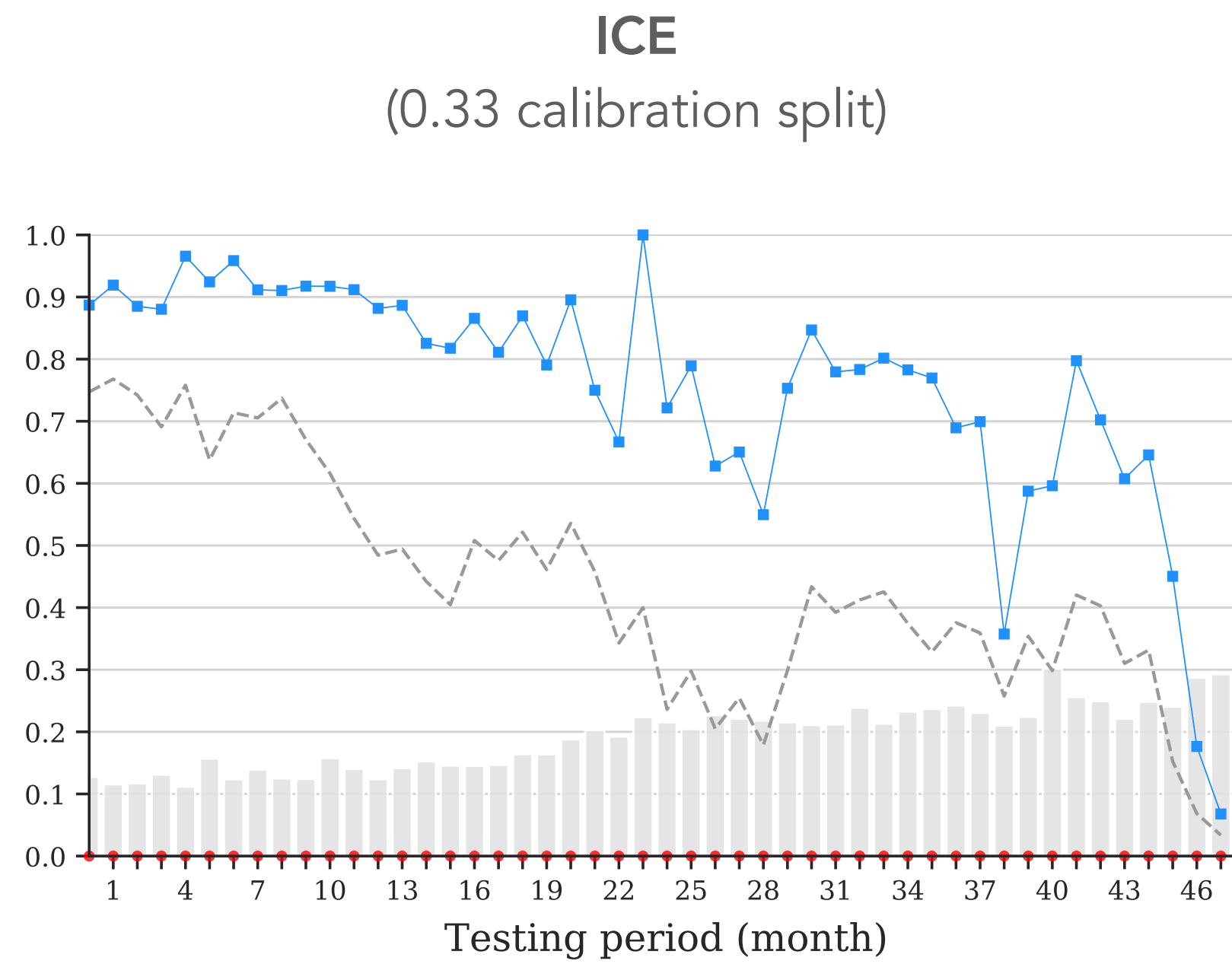
- Performance evaluation is notoriously tricky in machine learning.
- ML code tends to be complex and as yet lacks standardization.
- Subtle pitfalls arise from the differences between explanatory and predictive modeling.
- The hype and overoptimism about commercial AI may spill over into ML-based scientific research.
- Pressures and publication biases that have led to past reproducibility crises are also present in ML-based science.

Given these reasons, we view reproducibility difficulties as the expected state of affairs until best practices become better established and understood. The spate of reproducibility failures (that we have compiled below) highlights the immaturity of ML-based science, the critical need for ongoing work on methods and best practices, and the importance of treating the results from this body of work with caution.

[Sayash Kapoor](#) and [Arvind Narayanan](#) — <https://reproducible.cs.princeton.edu/>

Unintentional Data Snooping

Conformal
Evaluation



Unintentional Data Snooping

Conformal Evaluation

BUGFIX: fixed data snooping using ground truth labels when computing ...
...SVM scores for pvals

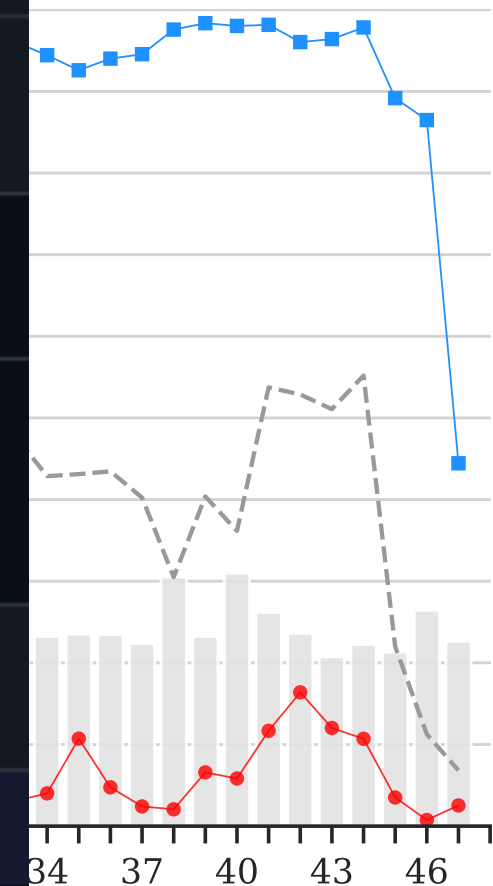
main

Icavallaro committed on Dec 8, 2022

Showing 1 changed file with 1 addition and 1 deletion.

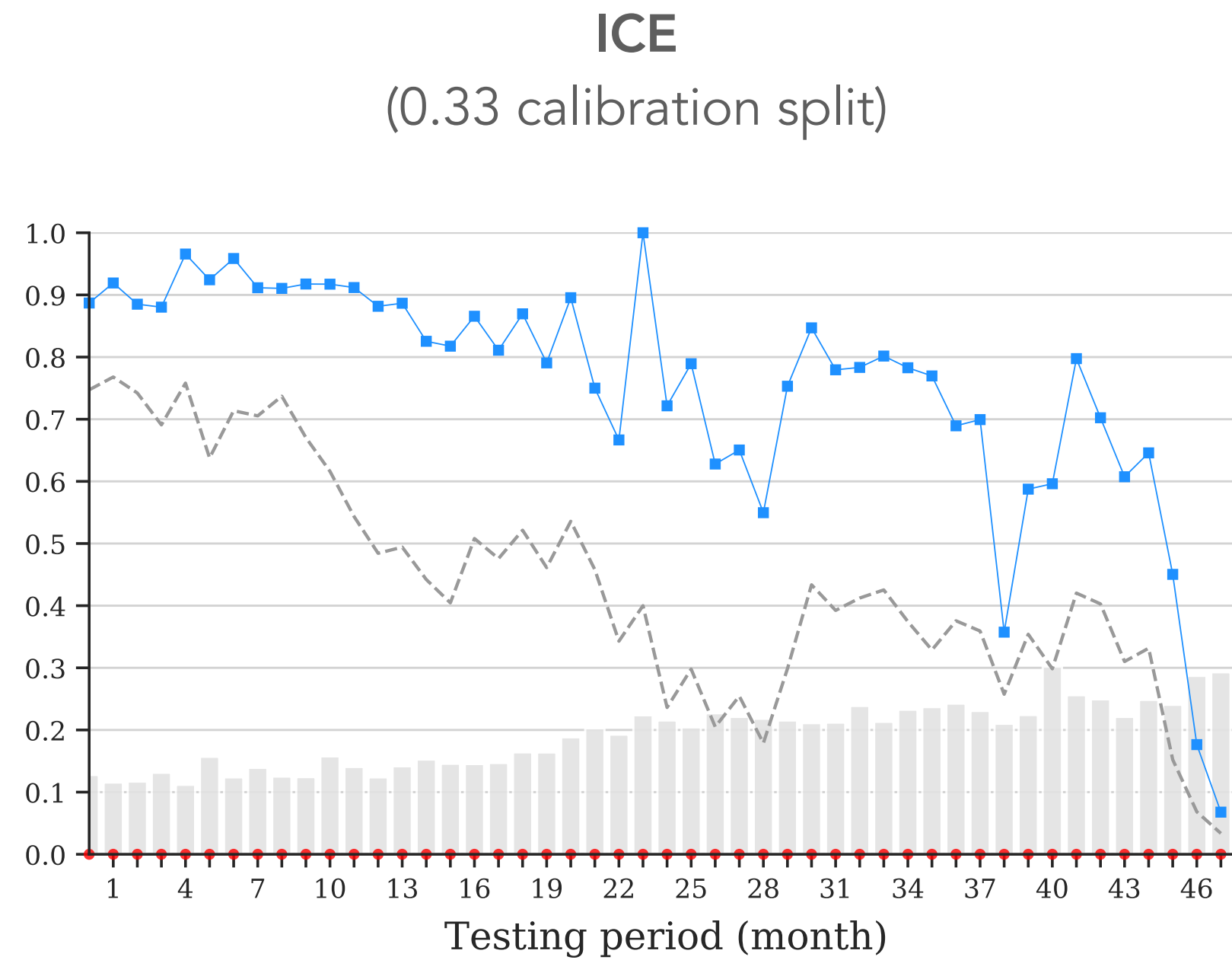
ice.py

```
@@ -282,7 +282,7 @@ def main():
    args.pval_consider))
282 282
283 283
284 284     logging.info('Getting NCMs for test ({})...'.format(argstest))
285 -     ncms_full_test = scores.get_svm_ncms(svm, X_test, y_test)
285 +     ncms_full_test = scores.get_svm_ncms(svm, X_test, pred_test)
286 286
287 287     p_val_test_dict = scores.compute_p_values_cred_and_conf(
288 288         train_ncms=ncms,
```

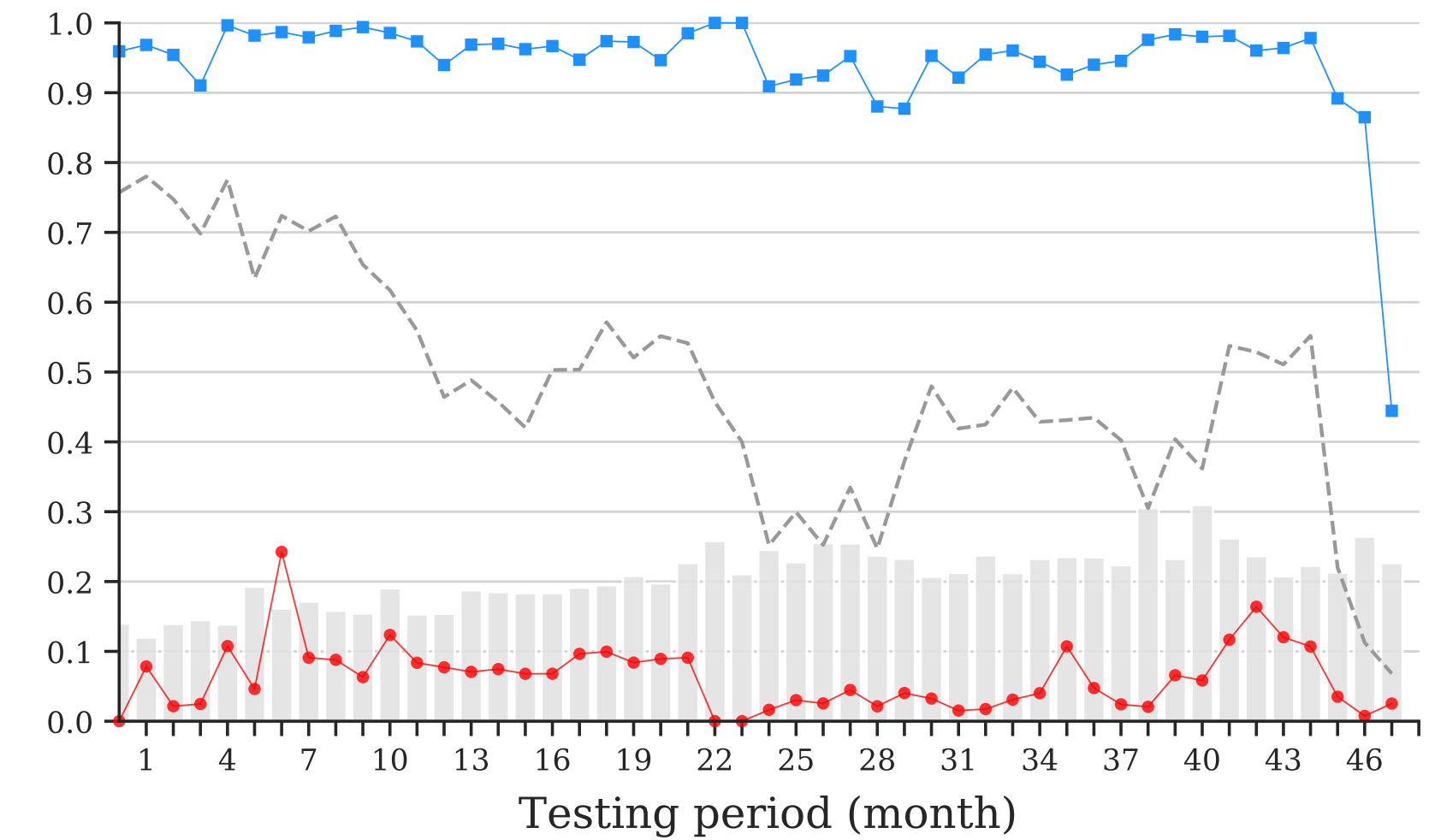


Unintentional Data Snooping

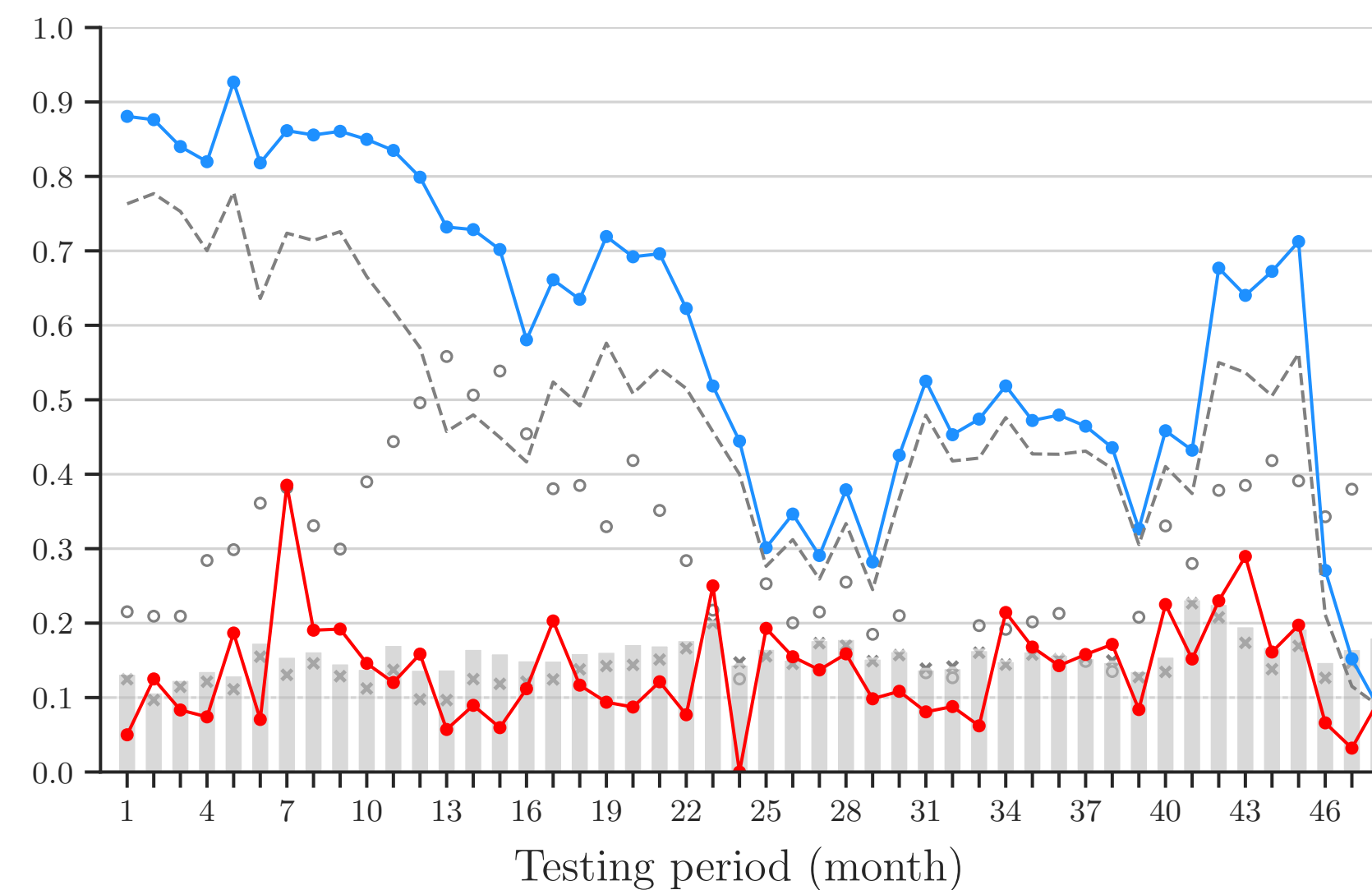
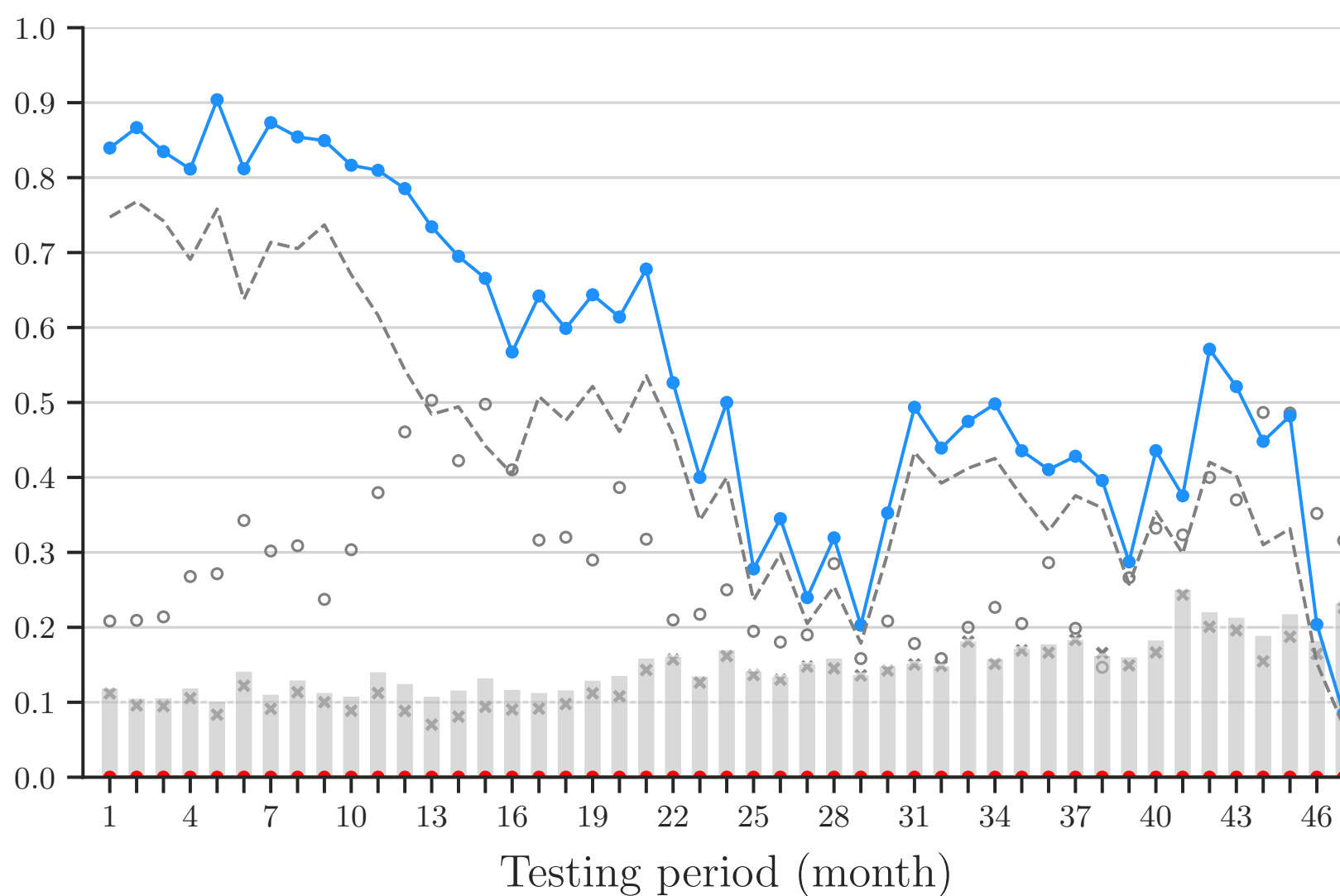
Conformal Evaluation



CCE
(10 folds)



Conformal Evaluation
(Fixed bug)



Unintentional Data Snooping^(*)

ICE
(0.33 calibration split)

CCE
(10 folds)

Conformal Evaluation

● **Let's not panic: liaised with IEEE S&P22 Chairs**

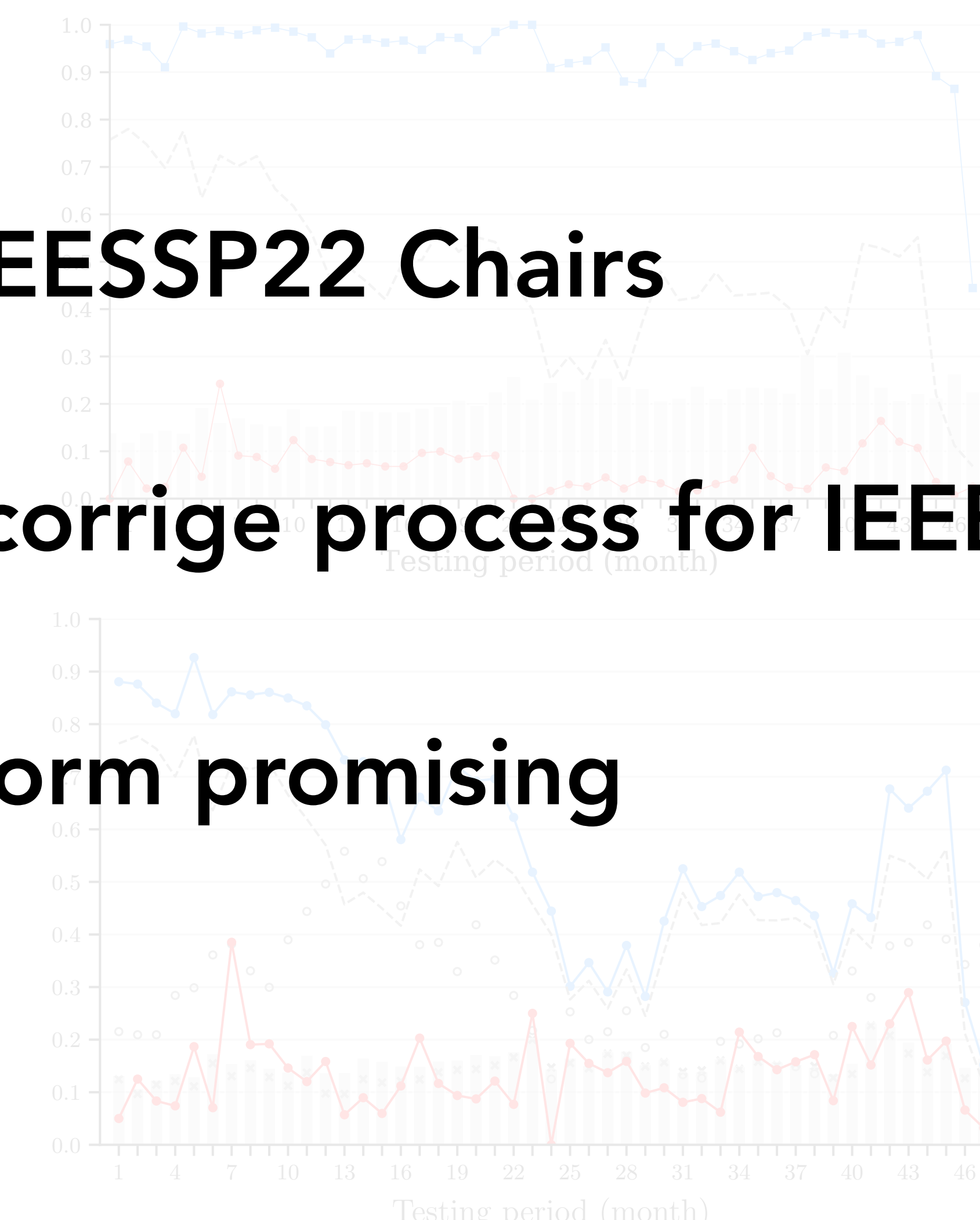
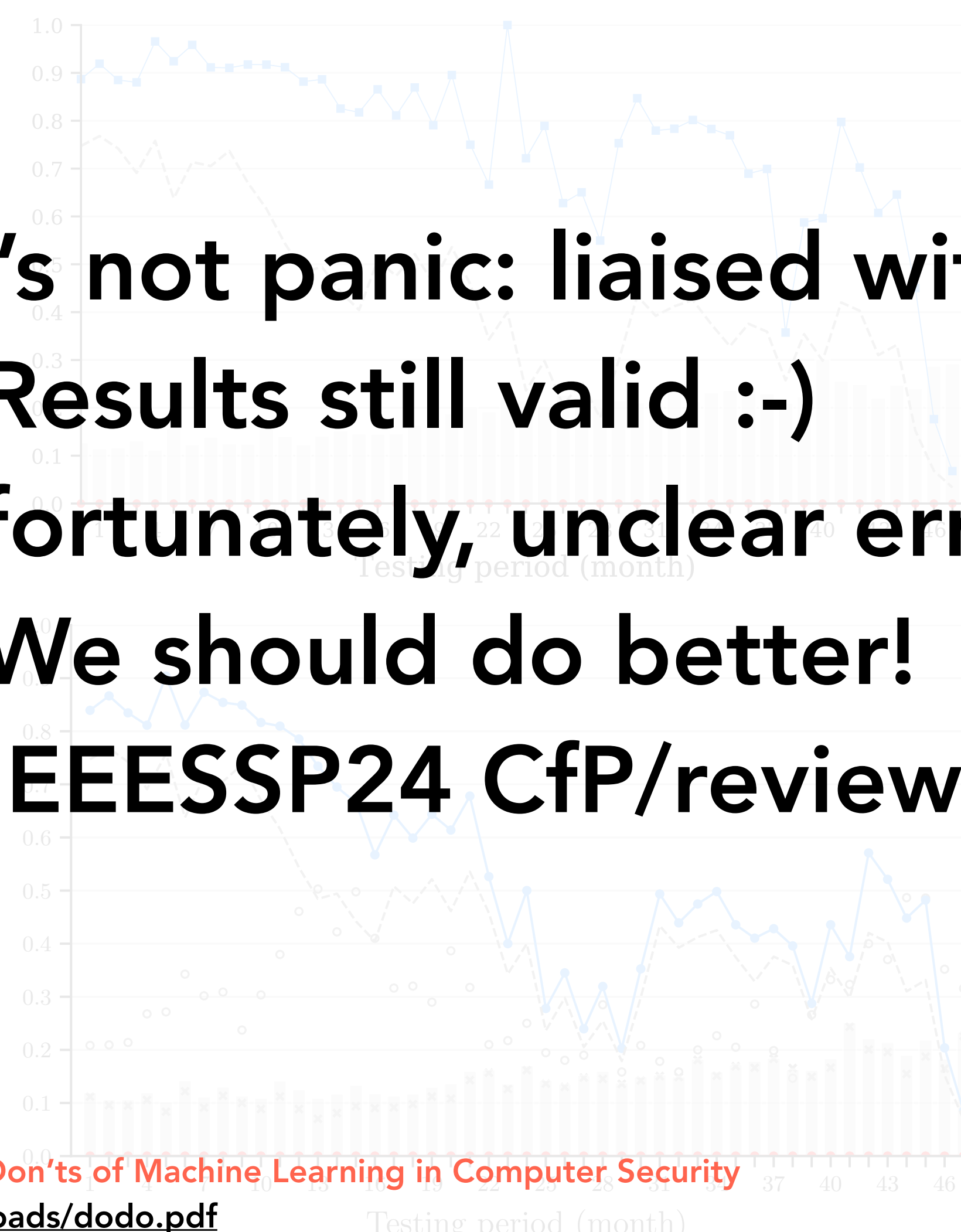
● **Results still valid :-)**

● **Unfortunately, unclear errata corrigee process for IEEE**

● **We should do better!**

● **IEEE S&P24 CfP/reviewing form promising**

Conformal Evaluation
(Fixed bug)



(*) [USENIX Sec 2022] **Dos and Don'ts of Machine Learning in Computer Security**
<https://s2lab.cs.ucl.ac.uk/downloads/dodo.pdf>

Unintentional Data Snooping(*)

IEEE S&P 2024 Review Form (Excerpt)

Scientific Contribution

Check all that apply

- 1. Independent Confirmation of Important Results with Limited Prior Research
- 2. Provides a New Data Set For Public Use
- 3. Creates a New Tool to Enable Future Science
- 4. Addresses a Long-Known Issue
- 5. Identifies an Impactful Vulnerability
- 6. Provides a Valuable Step Forward in an Established Field
- 7. Establishes a New Research Direction
- 8. Other

Scientific Contribution Comments

Conformal Evaluation

Let's

• R

• Unfo

• V

• IE

Conformal Evaluation (Fixed bug)

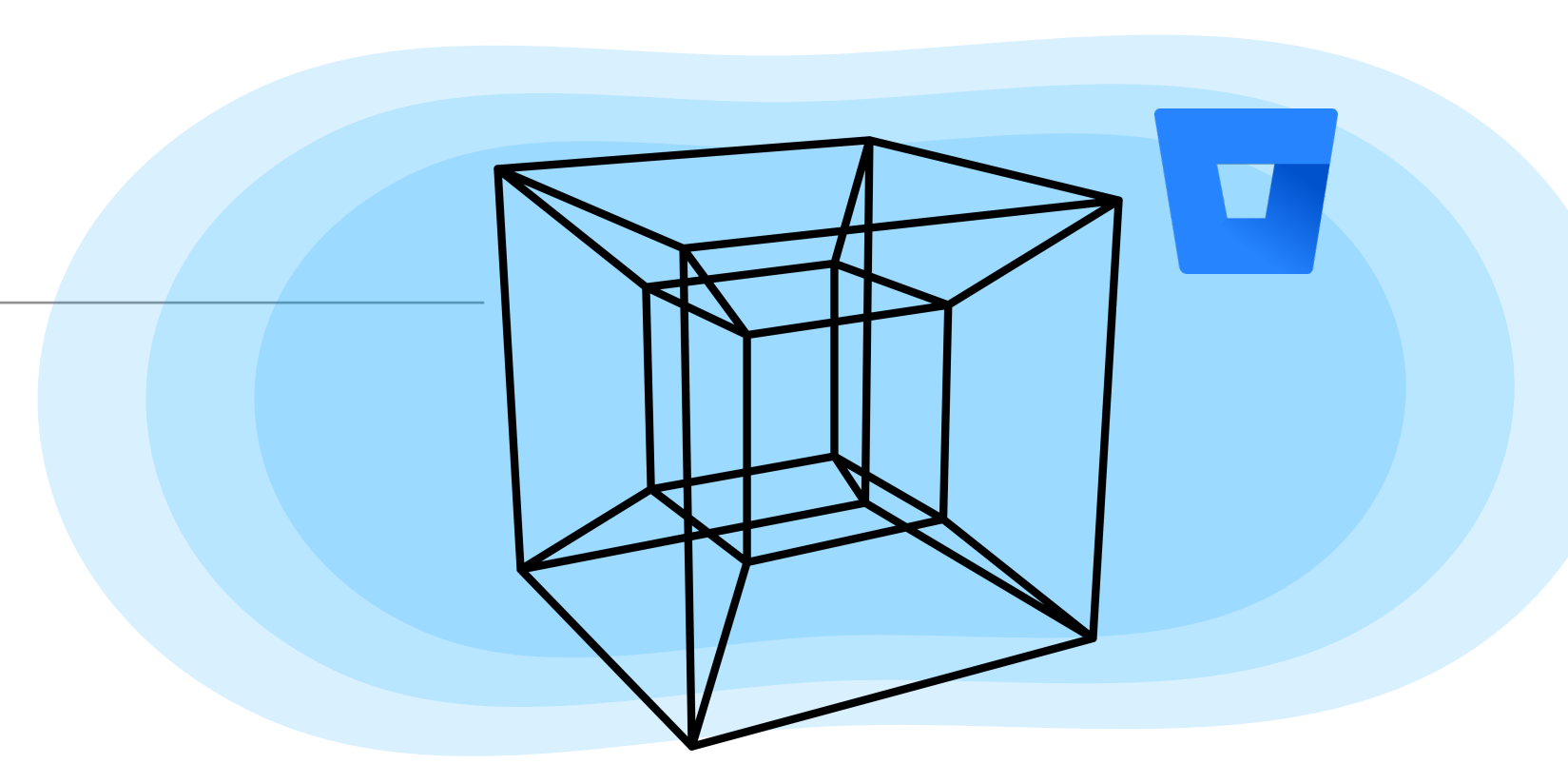
or IEEE

(*) [USENIX Sec 2022] Dos and Don'ts of Machine Learning in Computer Security

<https://s2lab.cs.ucl.ac.uk/downloads/dodo.pdf>

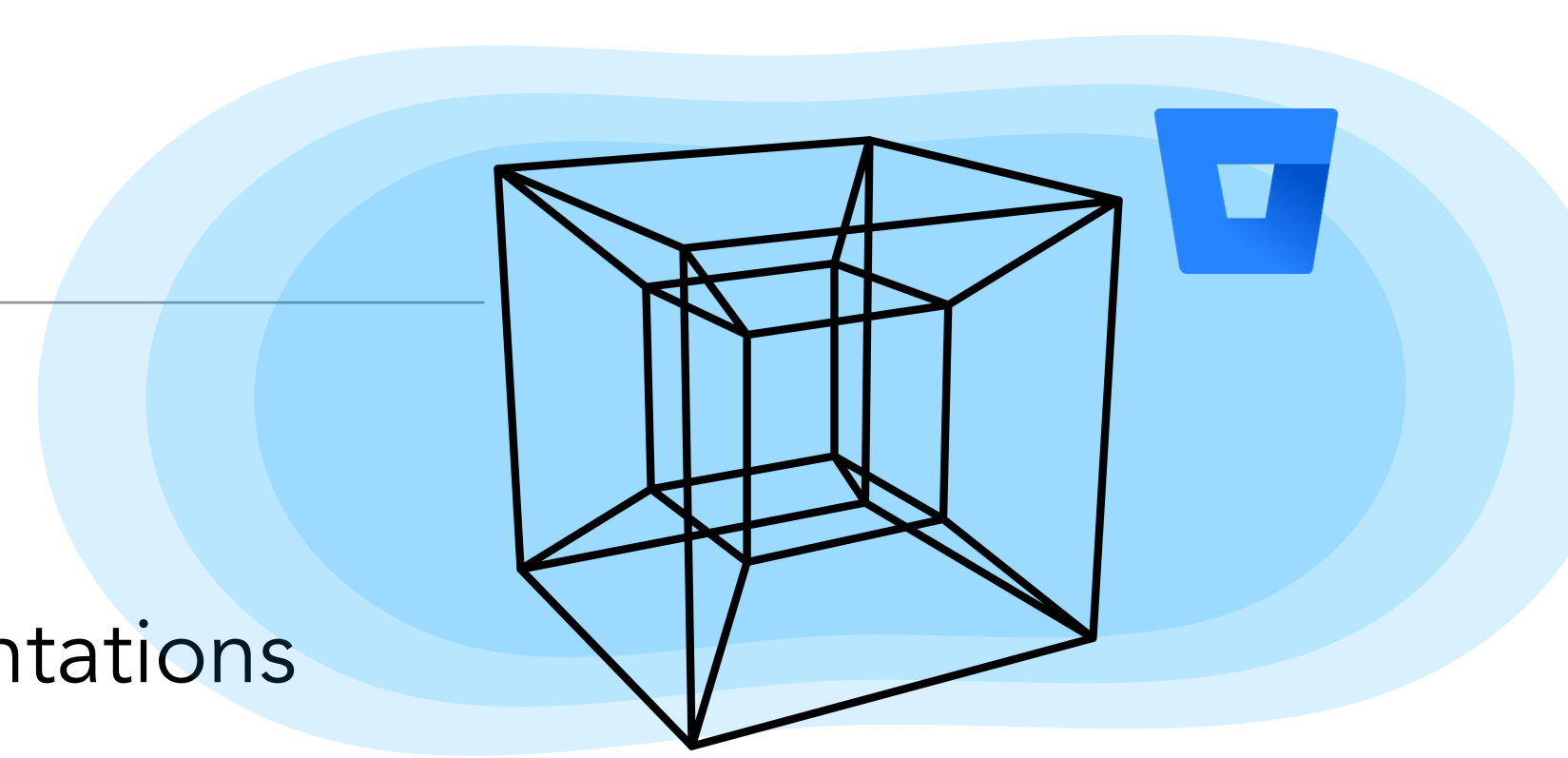
Take Aways

Take Aways



Take Aways

- Computer Security is highly **non-stationary** [1] and often **class-imbalanced**
 - › Arms-race between attackers and defenders; role of abstractions/representations
 - › Perform **time-aware** evaluations [1], and avoid pitfalls [5]
 - › Assume things go wrong: explore **rejection options** [3,4], **active learning** [1], **online learning** [6]



[1] Pendlebury et al., **TESSERACT: Eliminating experimental bias in malware classification across space and time**, USENIX Security 2019

[2] Pierazzi et al., **Intriguing Properties of Adversarial ML attacks in the problem space**, IEEE S&P 2020

[3] Jordaney et al., **Transcend: Detecting concept drift in malware classification models**, USENIX Security 2017

[4] Barbero et al., **Transcending Transcend: Revisiting malware classification in the presence of concept drift**, IEEE S&P 2022

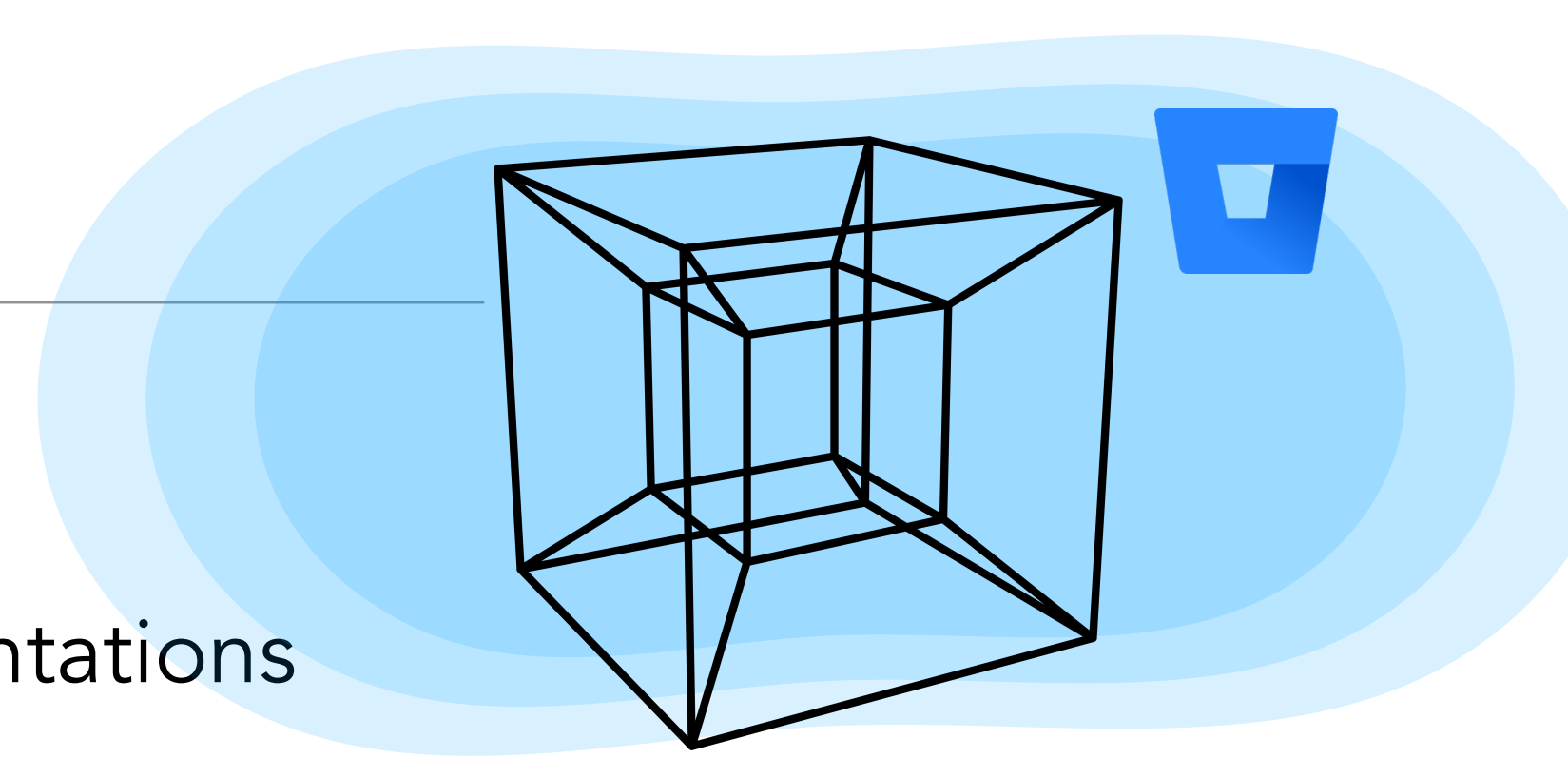
[5] Arp et al., **Dos and Dont's of Machine Learning in Security**, USENIX Security 2022

[6] Kan et al. **Investigating Labelless Drift Adaptation for Malware Detection**, AISEC 2021

[7] Yang et al. **Jigsaw Puzzle: Selective Backdoor Attack to Subvert Malware Classifiers**. IEEE S&P 2023

Take Aways

- Computer Security is highly **non-stationary** [1] and often **class-imbalanced**
 - › Arms-race between attackers and defenders; role of abstractions/representations
 - › Perform **time-aware** evaluations [1], and avoid pitfalls [5]
 - › Assume things go wrong: explore **rejection options** [3,4], **active learning** [1], **online learning** [6]
- Reason about problem space **(reliable) adversarial attacks and defenses** [2, 7]



[1] Pendlebury et al., **TESSERACT: Eliminating experimental bias in malware classification across space and time**, USENIX Security 2019

[2] Pierazzi et al., **Intriguing Properties of Adversarial ML attacks in the problem psace**, IEEE S&P 2020

[3] Jordaney et al., **Transcend: Detecting concept drift in malware classification models**, USENIX Security 2017

[4] Barbero et al., **Transcending Transcend: Revisiting malware classification in the presence of concept drift**, IEEE S&P 2022

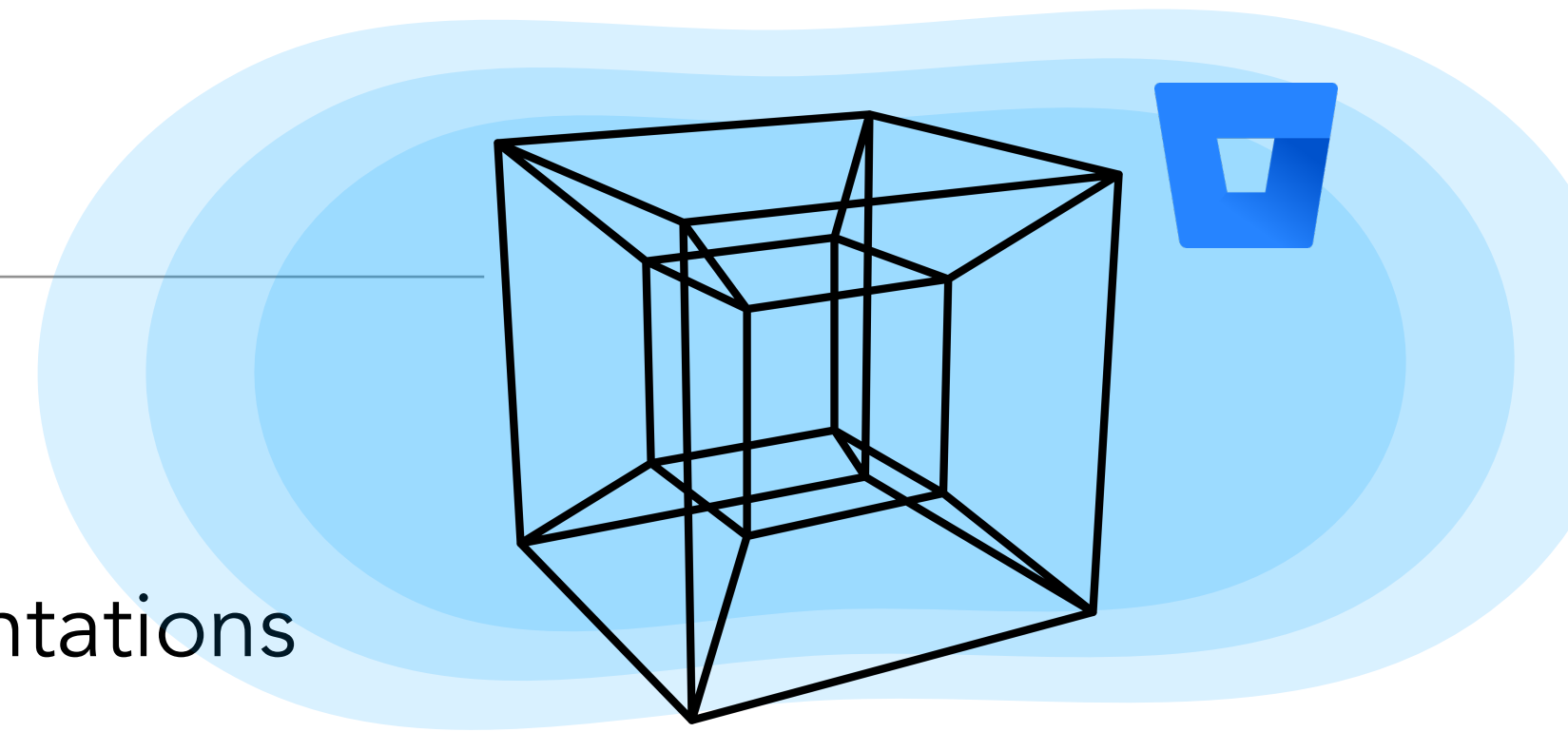
[5] Arp et al., **Dos and Dont's of Machine Learning in Security**, USENIX Security 2022

[6] Kan et al. **Investigating Labelless Drift Adaptation for Malware Detection**, AISec 2021

[7] Yang et al. **Jigsaw Puzzle: Selective Backdoor Attack to Subvert Malware Classifiers**. IEEE S&P 2023

Take Aways

- Computer Security is highly **non-stationary** [1] and often **class-imbalanced**
 - › Arms-race between attackers and defenders; role of abstractions/representations
 - › Perform **time-aware** evaluations [1], and avoid pitfalls [5]
 - › Assume things go wrong: explore **rejection options** [3,4], **active learning** [1], **online learning** [6]
- Reason about problem space **(reliable) adversarial attacks and defenses** [2, 7]
- Reason about the **relationship** between **adversarial ML** and **dataset shifts**



[1] Pendlebury et al., **TESSERACT: Eliminating experimental bias in malware classification across space and time**, USENIX Security 2019

[2] Pierazzi et al., **Intriguing Properties of Adversarial ML attacks in the problem space**, IEEE S&P 2020

[3] Jordaney et al., **Transcend: Detecting concept drift in malware classification models**, USENIX Security 2017

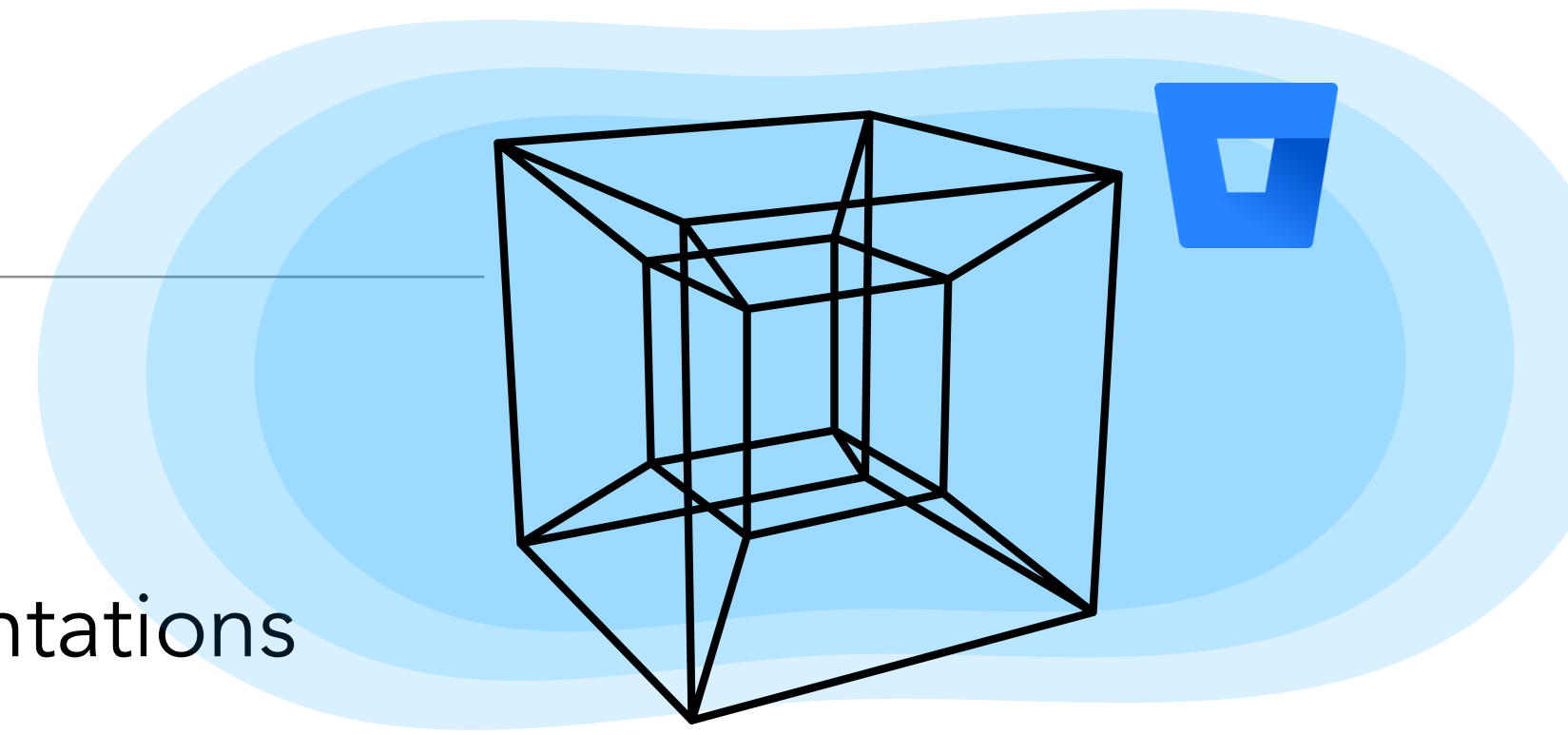
[4] Barbero et al., **Transcending Transcend: Revisiting malware classification in the presence of concept drift**, IEEE S&P 2022

[5] Arp et al., **Dos and Dont's of Machine Learning in Security**, USENIX Security 2022

[6] Kan et al. **Investigating Labelless Drift Adaptation for Malware Detection**, AISec 2021

[7] Yang et al. **Jigsaw Puzzle: Selective Backdoor Attack to Subvert Malware Classifiers**. IEEE S&P 2023

Take Aways



- Computer Security is highly **non-stationary** [1] and often **class-imbalanced**
 - › Arms-race between attackers and defenders; role of abstractions/representations
 - › Perform **time-aware** evaluations [1], and avoid pitfalls [5]
 - › Assume things go wrong: explore **rejection options** [3,4], **active learning** [1], **online learning** [6]
- Reason about problem space **(reliable) adversarial attacks and defenses** [2, 7]
- Reason about the **relationship** between **adversarial ML** and **dataset shifts**
- Reason about **abstractions** and **representations** and their effect on the entire ML pipeline

[1] Pendlebury et al., **TESSERACT: Eliminating experimental bias in malware classification across space and time**, USENIX Security 2019

[2] Pierazzi et al., **Intriguing Properties of Adversarial ML attacks in the problem space**, IEEE S&P 2020

[3] Jordaney et al., **Transcend: Detecting concept drift in malware classification models**, USENIX Security 2017

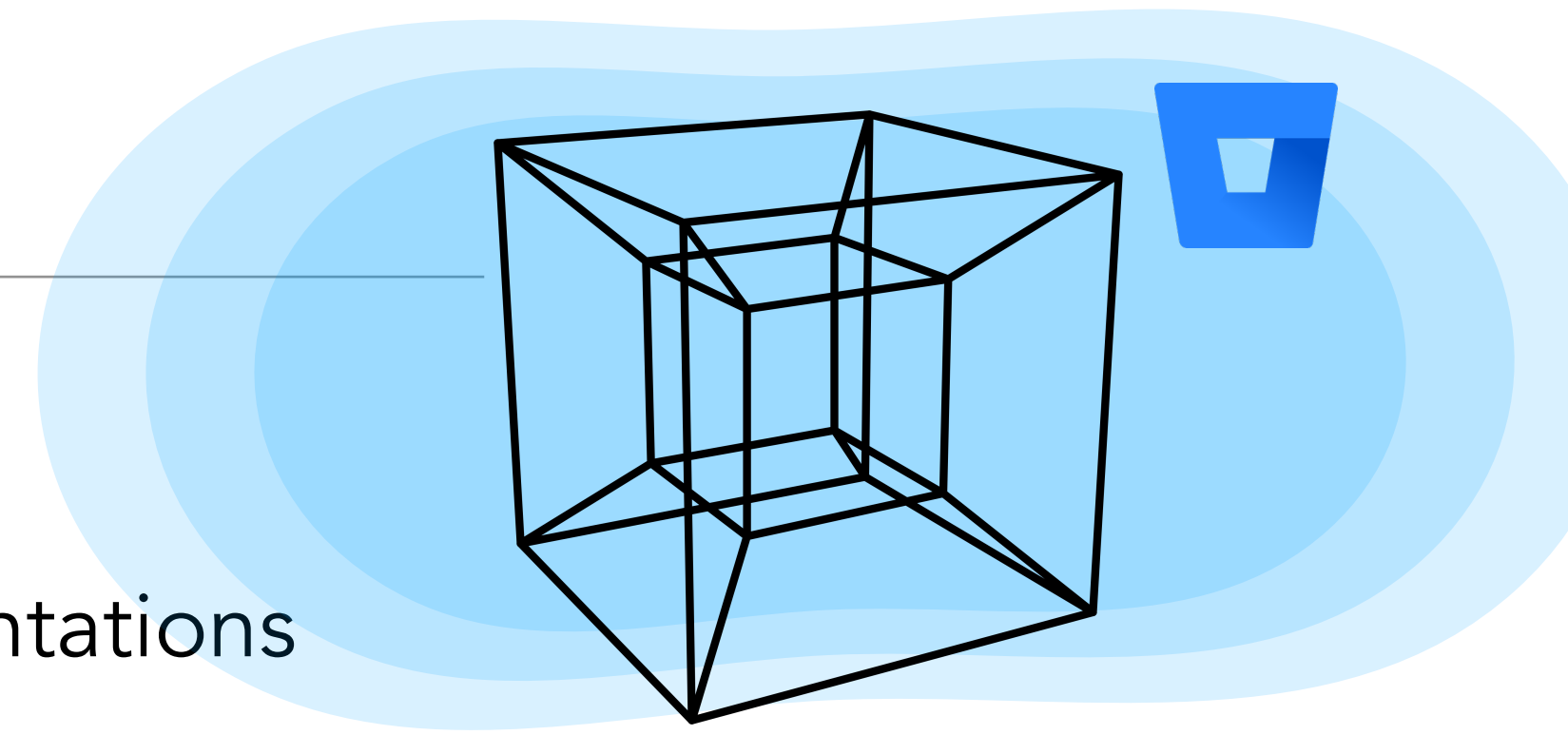
[4] Barbero et al., **Transcending Transcend: Revisiting malware classification in the presence of concept drift**, IEEE S&P 2022

[5] Arp et al., **Dos and Dont's of Machine Learning in Security**, USENIX Security 2022

[6] Kan et al. **Investigating Labelless Drift Adaptation for Malware Detection**, AISEC 2021

[7] Yang et al. **Jigsaw Puzzle: Selective Backdoor Attack to Subvert Malware Classifiers**. IEEE S&P 2023

Take Aways



- Computer Security is highly **non-stationary** [1] and often **class-imbalanced**
 - › Arms-race between attackers and defenders; role of abstractions/representations
 - › Perform **time-aware** evaluations [1], and avoid pitfalls [5]
 - › Assume things go wrong: explore **rejection options** [3,4], **active learning** [1], **online learning** [6]
- Reason about problem space **(reliable) adversarial attacks and defenses** [2, 7]
- Reason about the **relationship** between **adversarial ML** and **dataset shifts**
- Reason about **abstractions** and **representations** and their effect on the entire ML pipeline
- Bridging the academia-industry gap
 - › See <https://s2lab.cs.ucl.ac.uk> for access

[1] Pendlebury et al., **TESSERACT: Eliminating experimental bias in malware classification across space and time**, USENIX Security 2019

[2] Pierazzi et al., **Intriguing Properties of Adversarial ML attacks in the problem space**, IEEE S&P 2020

[3] Jordaney et al., **Transcend: Detecting concept drift in malware classification models**, USENIX Security 2017

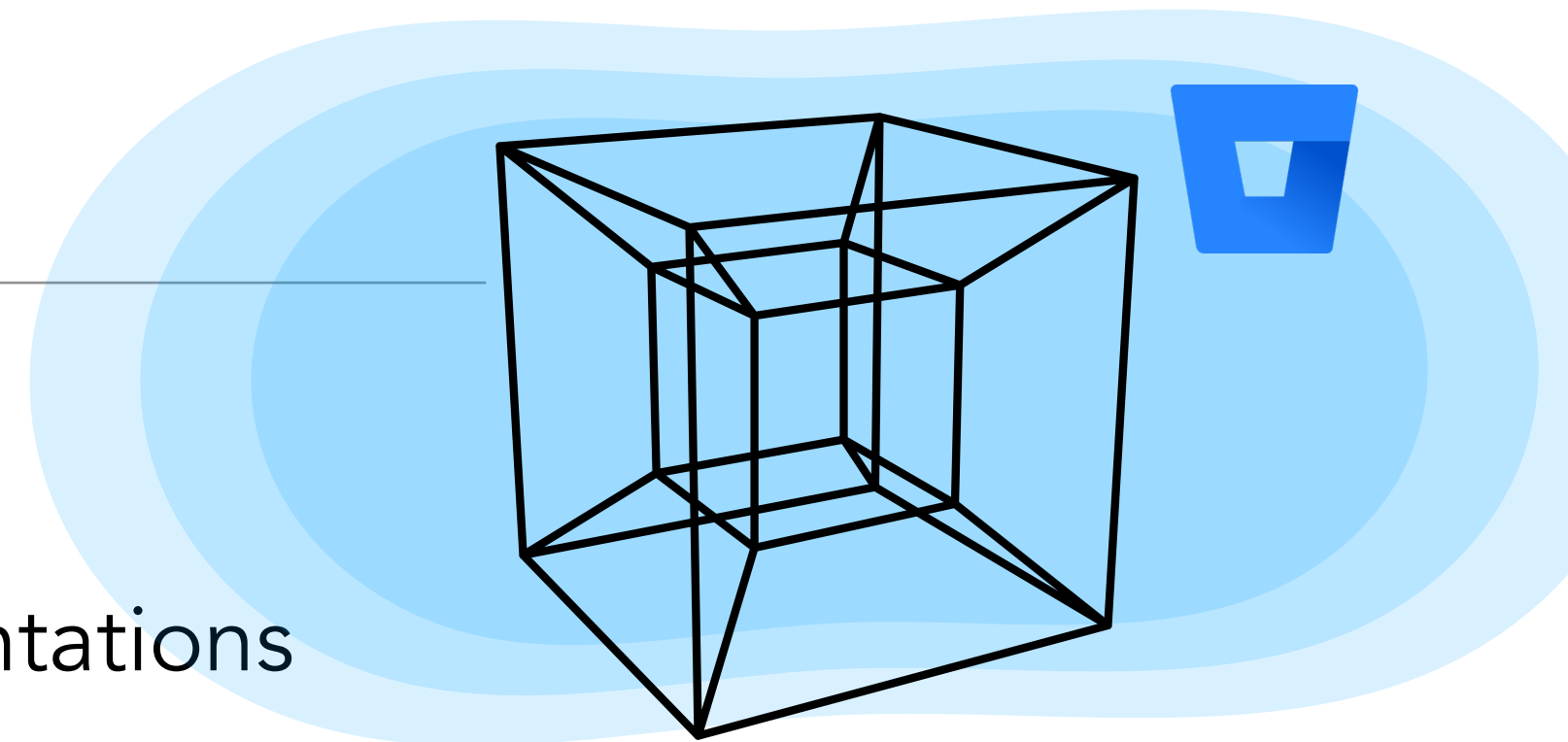
[4] Barbero et al., **Transcending Transcend: Revisiting malware classification in the presence of concept drift**, IEEE S&P 2022

[5] Arp et al., **Dos and Dont's of Machine Learning in Security**, USENIX Security 2022

[6] Kan et al. **Investigating Labelless Drift Adaptation for Malware Detection**, AISEC 2021

[7] Yang et al. **Jigsaw Puzzle: Selective Backdoor Attack to Subvert Malware Classifiers**. IEEE S&P 2023

Take Aways



- Computer Security is highly **non-stationary** [1] and often **class-imbalanced**
 - › Arms-race between attackers and defenders; role of abstractions/representations
 - › Perform **time-aware** evaluations [1], and avoid pitfalls [5]
 - › Assume things go wrong: explore **rejection options** [3,4], **active learning** [1], **online learning** [6]
- Reason about problem space **(reliable) adversarial attacks and defenses** [2, 7]
- Reason about the **relationship** between **adversarial ML** and **dataset shifts**
- Reason about **abstractions** and **representations** and their effect on the entire ML pipeline
- Bridging the academia-industry gap
 - › See <https://s2lab.cs.ucl.ac.uk> for access

https://twitter.com/joshua_saxe/status/1550545466072264704



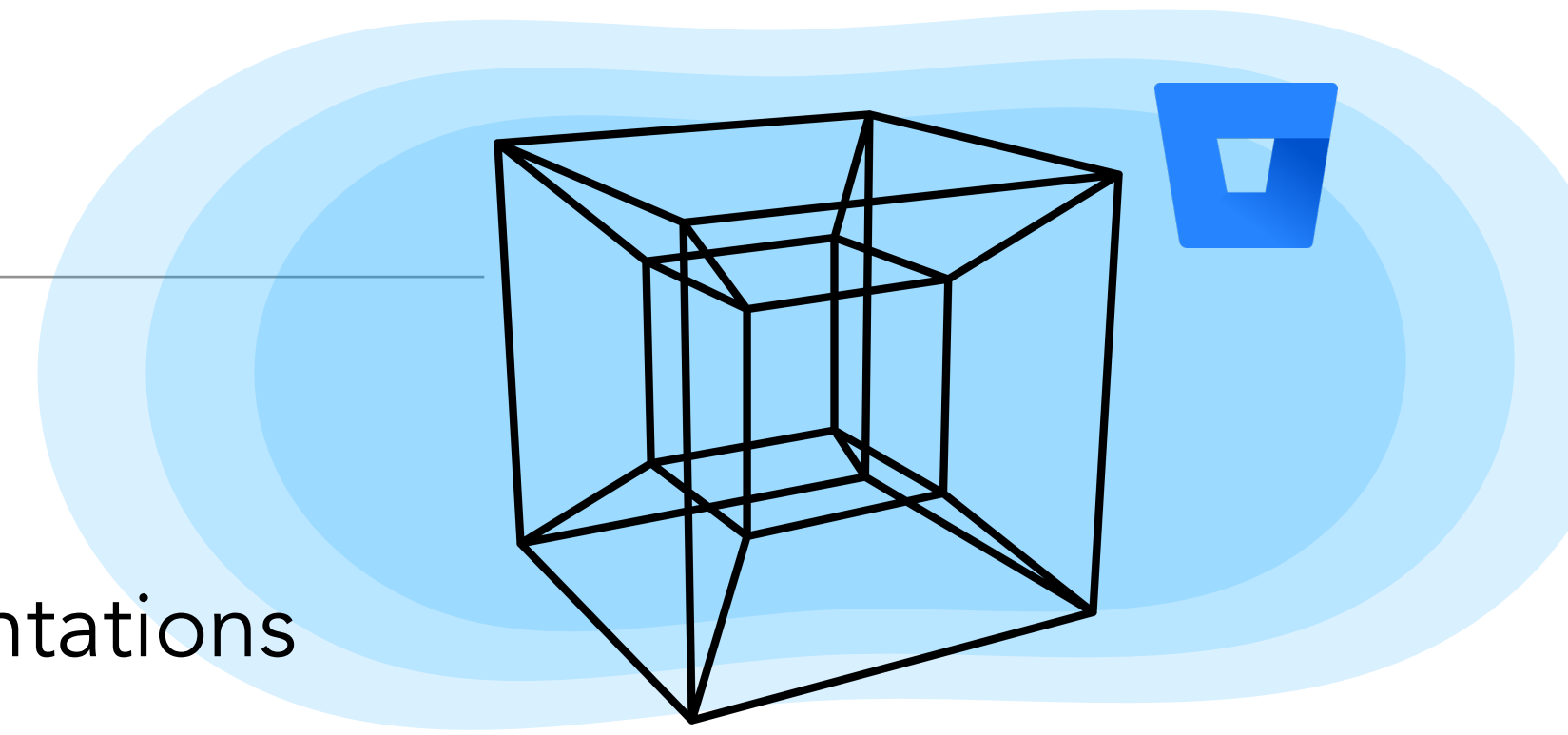
[1] Pendlebury et al., TESSERACT: Eliminating experimental bias in malware classification
[2] Pierazzi et al., Intriguing Properties of Adversarial ML attacks in the wild
[3] Jordaney et al., Transcend: Detecting concept drift in malware classification
[4] Barbero et al., Transcending Transcend: Revisiting malware classification
[5] Arp et al., Dos and Dont's of Machine Learning in Security, USENIX Security 2021
[6] Kan et al. Investigating Labelless Drift Adaptation for Malware Detection, AISEC 2021
[7] Yang et al. Jigsaw Puzzle: Selective Backdoor Attack to Subvert Malware Classifiers. IEEE S&P 2023

Our Open-Source Libraries

- Requested access by 120+ organizations, including (honorable mentions):



Take Aways



- Computer Security is highly **non-stationary** [1] and often **class-imbalanced**
 - › Arms-race between attackers and defenders; role of abstractions/representations
 - › Perform **time-aware** evaluations [1], and avoid pitfalls [5]
 - › Assume things go wrong: explore **rejection options** [3,4], **active learning** [1], **online learning** [6]
- Reason about problem space (**reliable**) **adversarial attacks and defenses** [2, 7]
- Reason about the **relationship** between **adversarial ML** and **dataset shifts**
- Reason about **abstractions** and **representations** and their effect on the entire ML pipeline
- Bridging the academia-industry gap
 - › See <https://s2lab.cs.ucl.ac.uk> for access

[1] Pendlebury et al., **TESSERACT: Eliminating experimental bias in malware classification across space and time**, USENIX Security 2019

[2] Pierazzi et al., **Intriguing Properties of Adversarial ML attacks in the problem space**, IEEE S&P 2020

[3] Jordaney et al., **Transcend: Detecting concept drift in malware classification models**, USENIX Security 2017

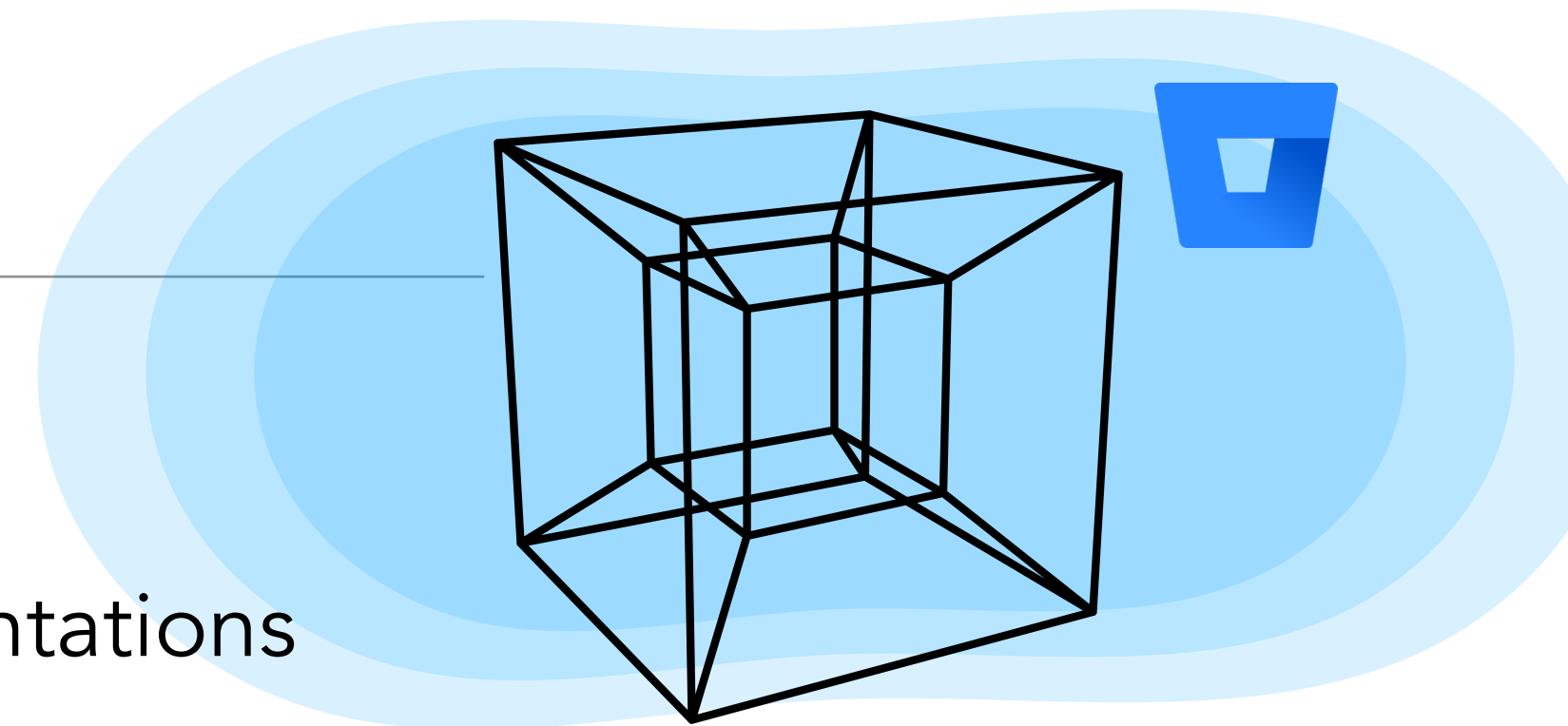
[4] Barbero et al., **Transcending Transcend: Revisiting malware classification in the presence of concept drift**, IEEE S&P 2022

[5] Arp et al., **Dos and Dont's of Machine Learning in Security**, USENIX Security 2022

[6] Kan et al. **Investigating Labelless Drift Adaptation for Malware Detection**, AISEC 2021

[7] Yang et al. **Jigsaw Puzzle: Selective Backdoor Attack to Subvert Malware Classifiers**. IEEE S&P 2023

Take Aways



- Computer Security is highly **non-stationary** [1] and often **class-imbalanced**
 - › Arms-race between attackers and defenders; role of abstractions/representations
 - › Perform **time-aware** evaluations [1], and avoid pitfalls [5]
 - › Assume things go wrong: explore **rejection options** [3,4], **active learning** [1], **online learning** [6]

- Reason about problem space **(reliable) adversarial attacks and defenses** [2, 7]

- Reason about the **relationship** between **adversarial ML** and **dataset shifts**

Exciting Recent Effort [8]

- Reason about **abstractions** and **representations** and their effect on the entire ML pipeline

- Bridging the academia-industry gap

[8] Pei et al. **Symmetry-Preserving Program Representations for Learning Code Semantics**
Collaboration with Columbia University — <https://arxiv.org/abs/2308.03312>

- › See <https://s2lab.cs.ucl.ac.uk> for access

[1] Pendlebury et al., **TESSERACT: Eliminating experimental bias in malware classification across space and time**, USENIX Security 2019

[2] Pierazzi et al., **Intriguing Properties of Adversarial ML attacks in the problem space**, IEEE S&P 2020

[3] Jordaney et al., **Transcend: Detecting concept drift in malware classification models**, USENIX Security 2017

[4] Barbero et al., **Transcending Transcend: Revisiting malware classification in the presence of concept drift**, IEEE S&P 2022

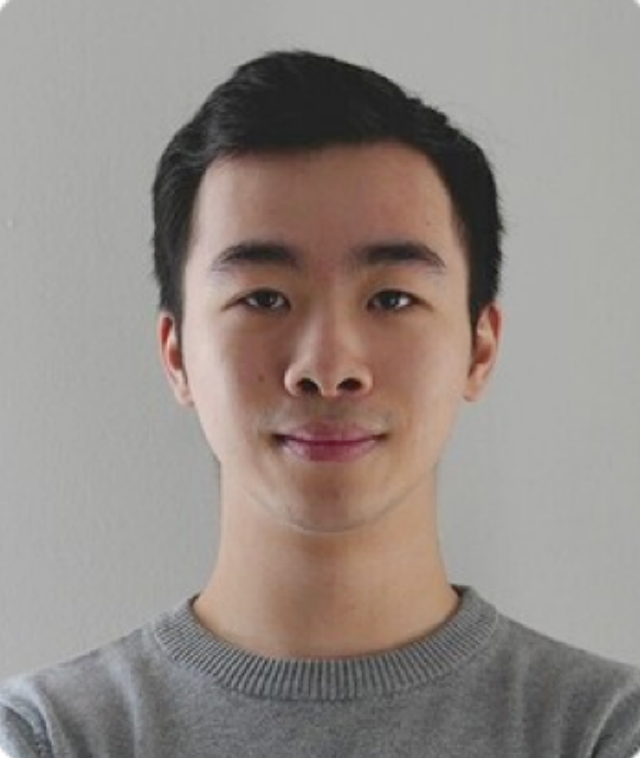
[5] Arp et al., **Dos and Dont's of Machine Learning in Security**, USENIX Security 2022

[6] Kan et al. **Investigating Labelless Drift Adaptation for Malware Detection**, AISEC 2021

[7] Yang et al. **Jigsaw Puzzle: Selective Backdoor Attack to Subvert Malware Classifiers**. IEEE S&P 2023

Core Team

Ph.D. Students



Theo



Jacopo



Mark



Giulio



Feargus



Fabio Pierazzi

Ph.D. Alumni

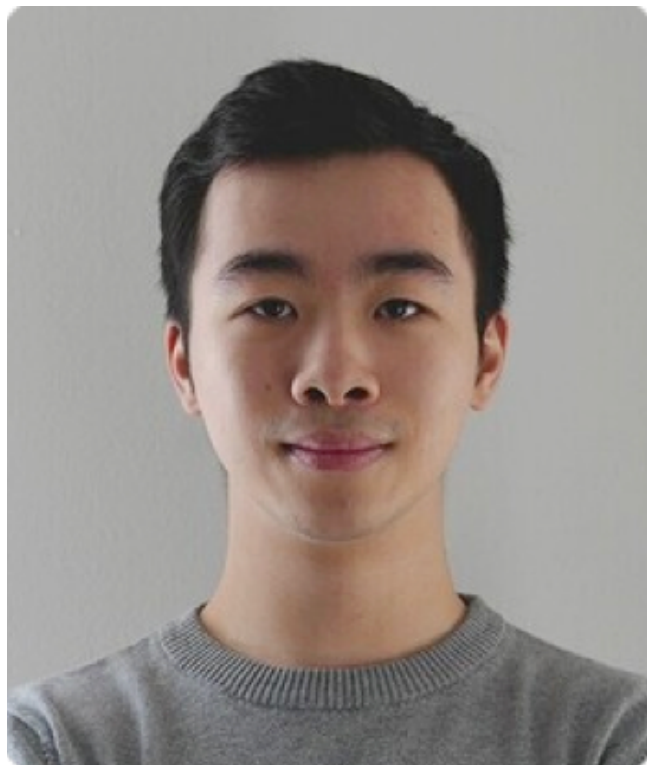
Team-ups

Current Research Collaborators



Core Team

Ph.D. Students



Theo



Jacopo



Mark



Giulio



Feargus



Fabio Pierazzi

Ph.D. Alumni

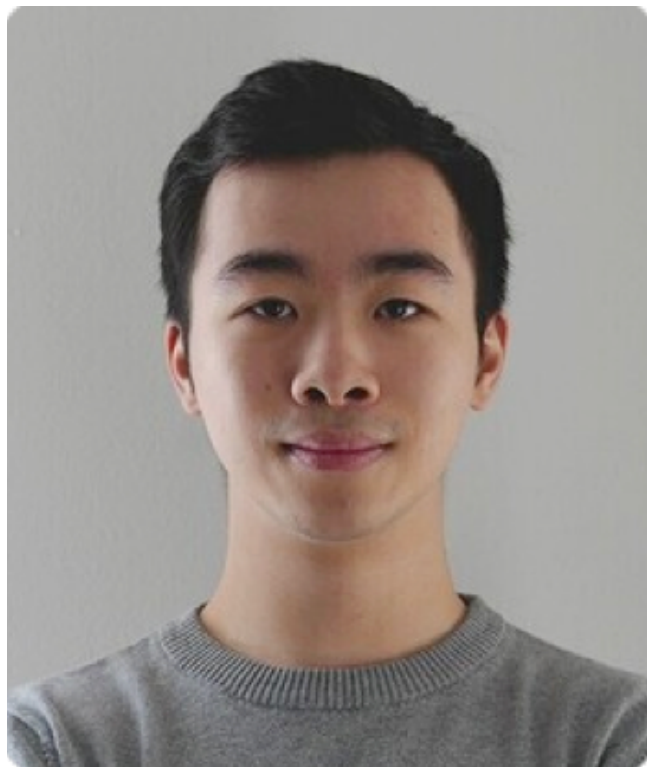
Team-ups

Current Research Collaborators



Core Team

Ph.D. Students



Theo



Jacopo



Mark



Giulio



Feargus



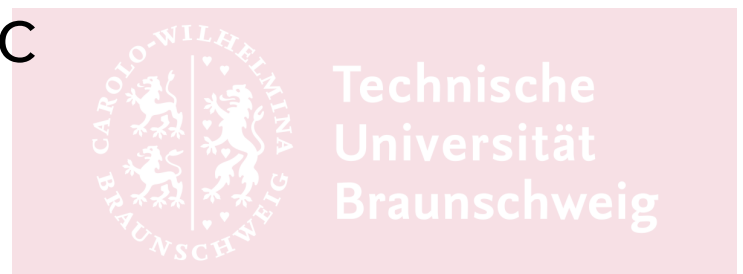
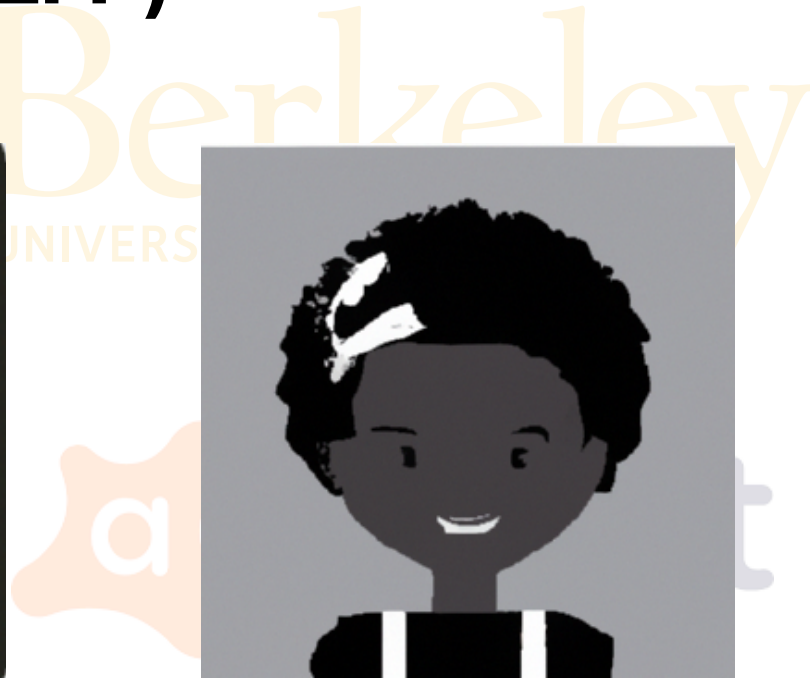
Fabio Pierazzi

Ph.D. Alumni

Team-ups

Current Research Collaborators

I am hiring at UCL! :-)



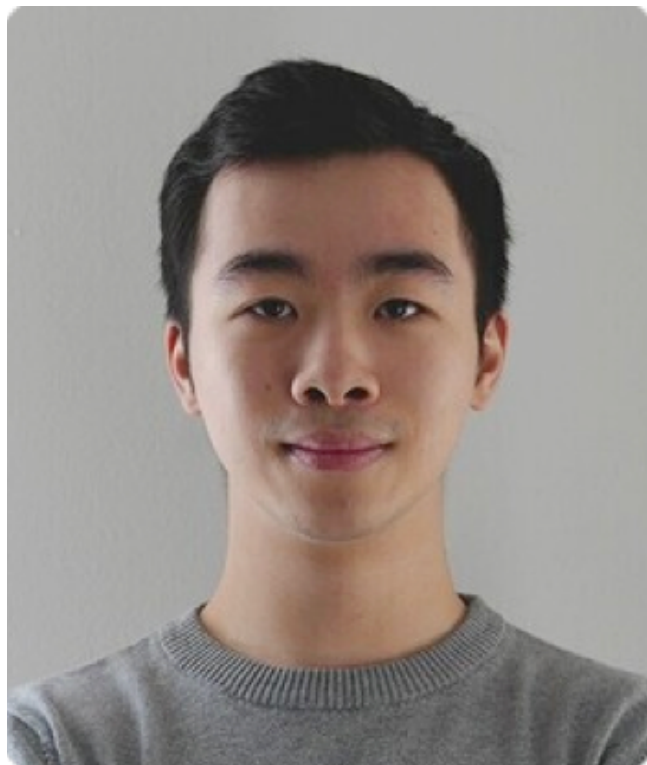
PhD1

PhD2

PostDoc

Core Team

Ph.D. Students



Theo



Jacopo



Mark



Giulio



Feargus



Fabio Pierazzi

Ph.D. Alumni

Team-ups

I am hiring at UCL! :-)

Current Research Collaborators



Xinran Zheng



Xingzhi Qiang



PostDoc



Take Aways

- Computer Security is highly **non-stationary** [1] and often **class-imbalanced**
 - › Arms-race between attackers and defenders; role of abstractions/representations
 - › Perform **time-aware** evaluations [1], and avoid pitfalls [5]
 - › Assume things go wrong: explore **rejection options** [3,4], **active learning** [1], **online learning** [6]

- Reason about problem space **(reliable) adversarial attacks and defenses** [2, 7]

- Reason about the **relationship** between **adversarial ML** and **dataset shifts**

Exciting Recent Effort [8]

- Reason about **abstractions** and **representations** and their effect on the entire ML pipeline

- Bridging the academia-industry gap

[8] Pei et al. **Symmetry-Preserving Program Representations for Learning Code Semantics**
Collaboration with Columbia University — <https://arxiv.org/abs/2308.03312>

- › See <https://s2lab.cs.ucl.ac.uk> for access

[1] Pendlebury et al., **TESSERACT: Eliminating experimental bias in malware classification across space and time**, USENIX Security 2019

[2] Pierazzi et al., **Intriguing Properties of Adversarial ML attacks in the problem space**, IEEE S&P 2020

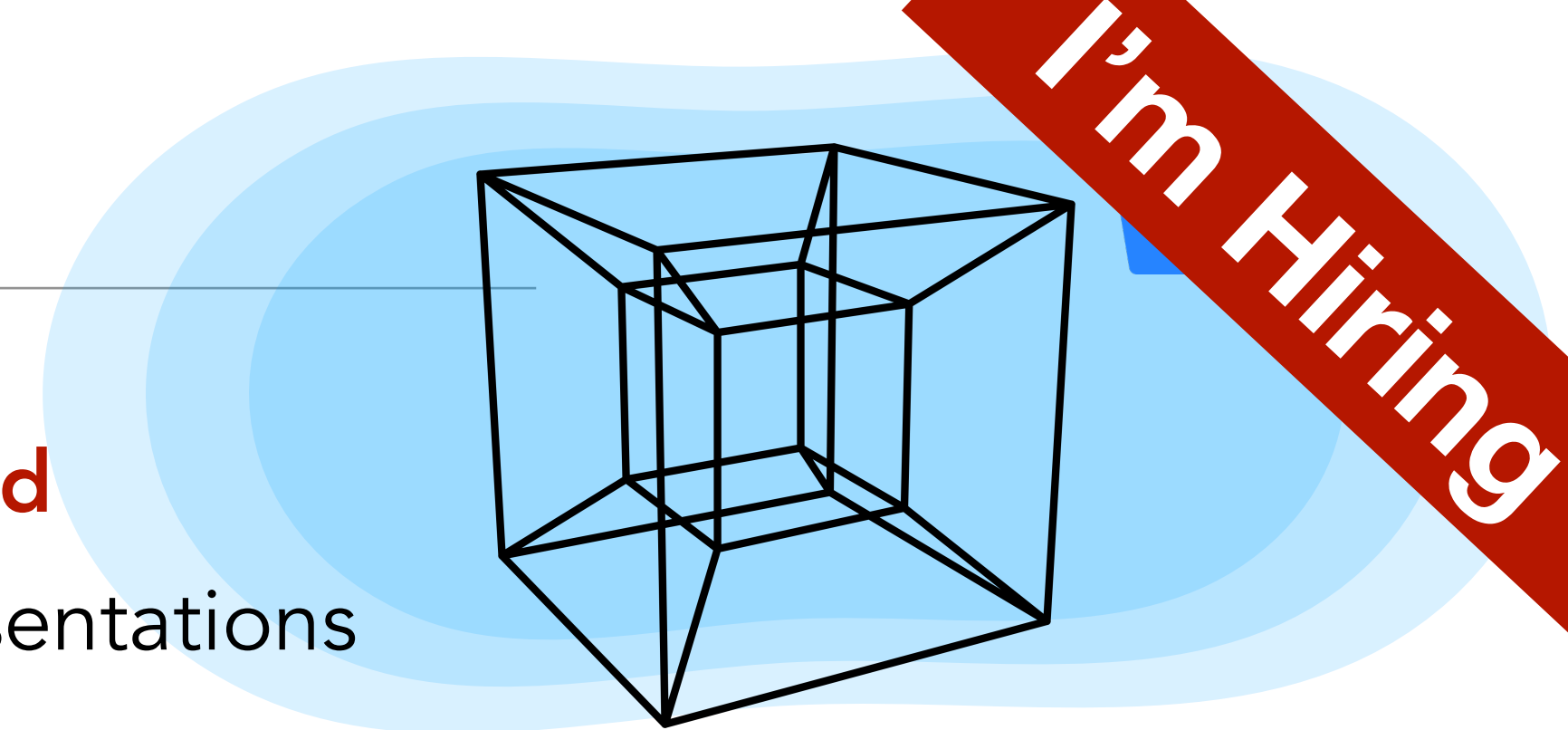
[3] Jordaney et al., **Transcend: Detecting concept drift in malware classification models**, USENIX Security 2017

[4] Barbero et al., **Transcending Transcend: Revisiting malware classification in the presence of concept drift**, IEEE S&P 2022

[5] Arp et al., **Dos and Dont's of Machine Learning in Security**, USENIX Security 2022

[6] Kan et al. **Investigating Labelless Drift Adaptation for Malware Detection**, AISEC 2021

[7] Yang et al. **Jigsaw Puzzle: Selective Backdoor Attack to Subvert Malware Classifiers**. IEEE S&P 2023



Backup Slides

Sources of Experimental Bias (1/3)

Temporal Inconsistency in Train/Test Sets

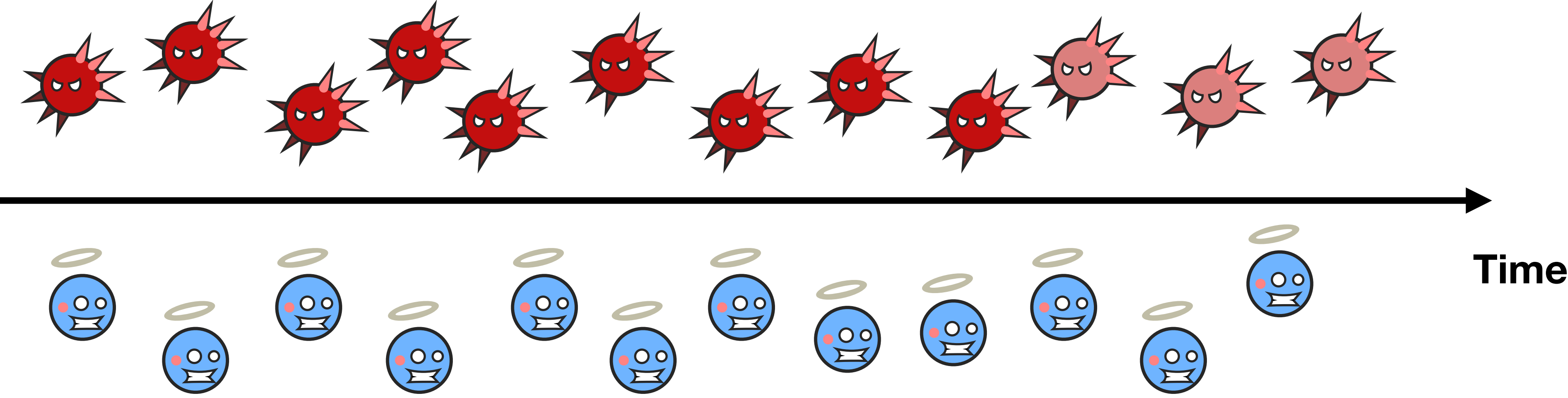
Sources of Experimental Bias (1/3)

Temporal Inconsistency in Train/Test Sets



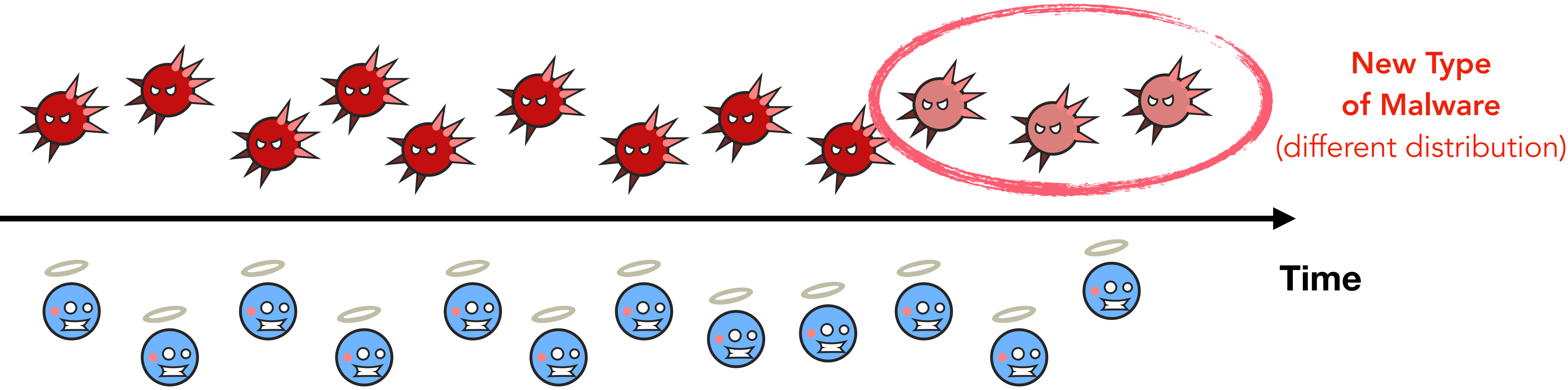
Sources of Experimental Bias (1/3)

Temporal Inconsistency in Train/Test Sets



Sources of Experimental Bias (1/3)

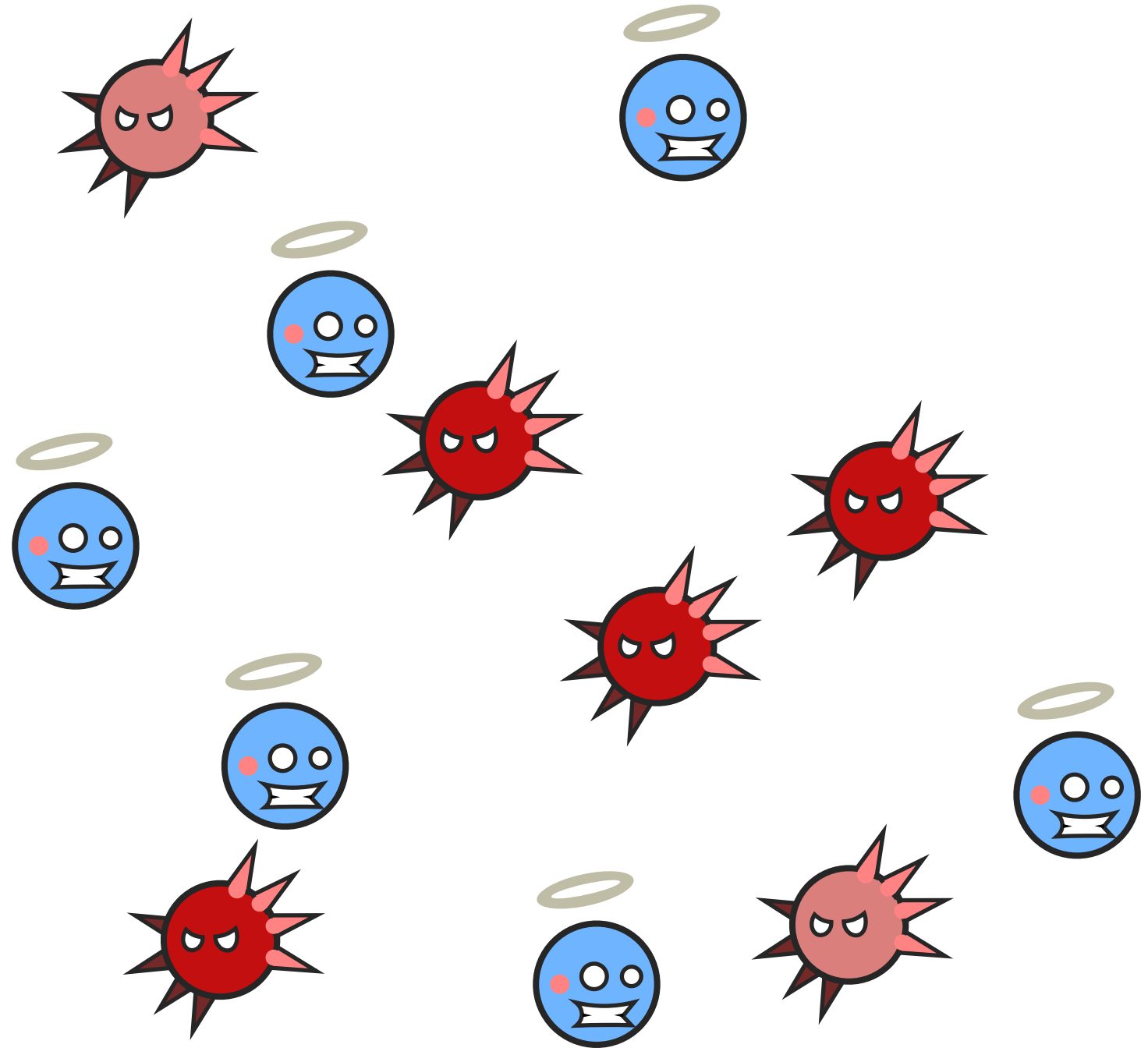
Temporal Inconsistency in Train/Test Sets



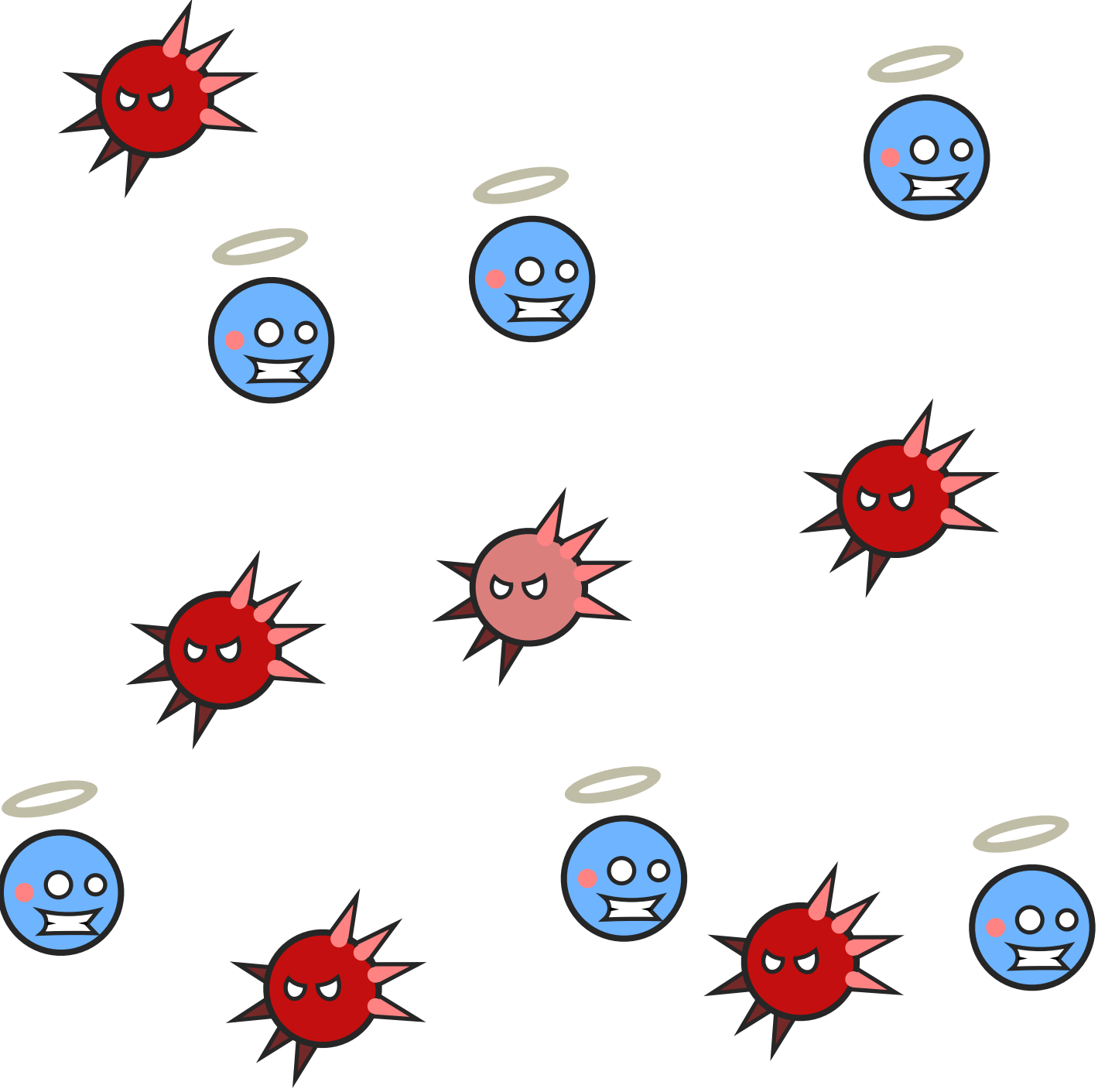
Sources of Experimental Bias (1/3)

Temporal Inconsistency in Train/Test Sets

Training



Testing

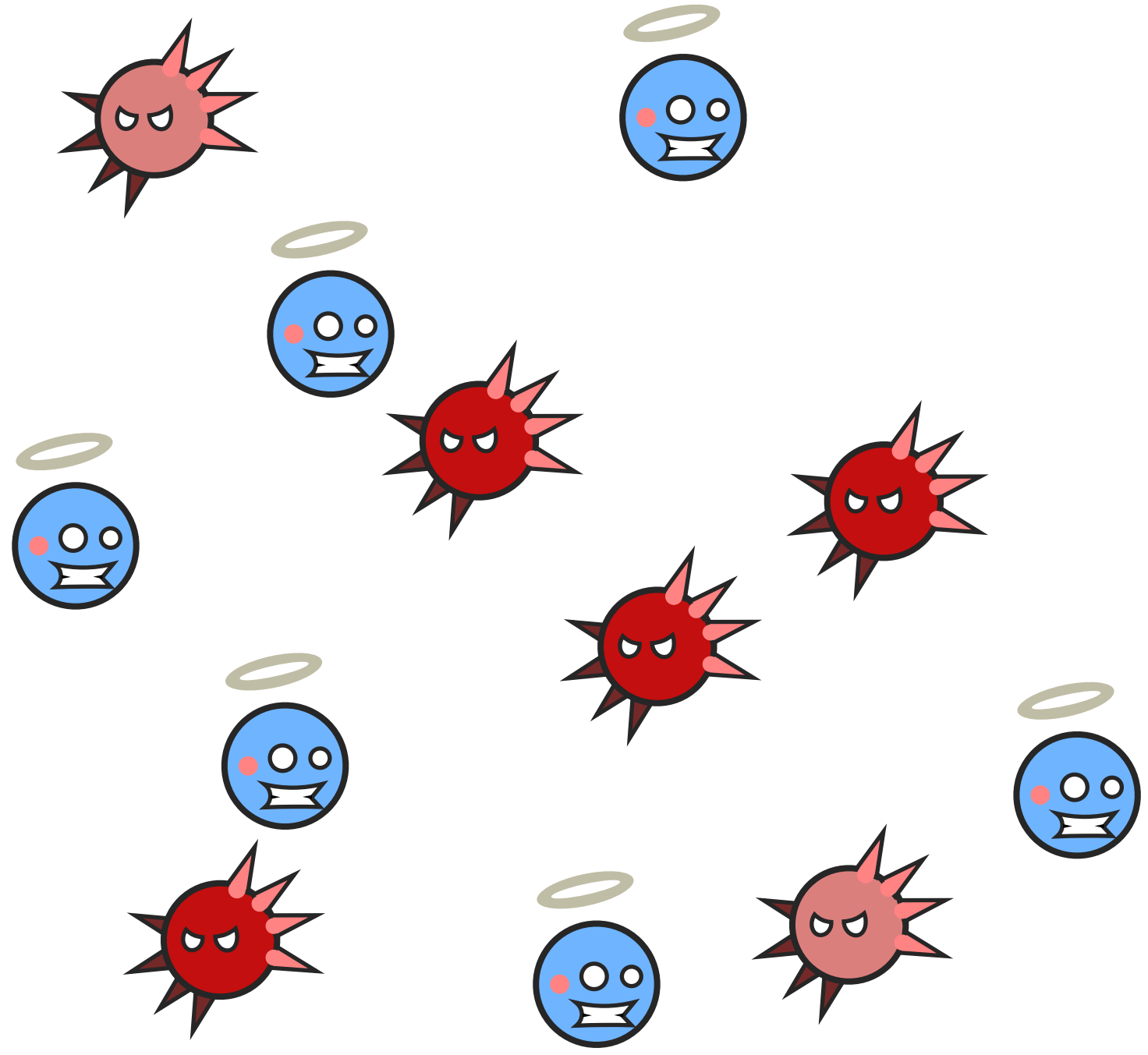


Sources of Experimental Bias (1/3)

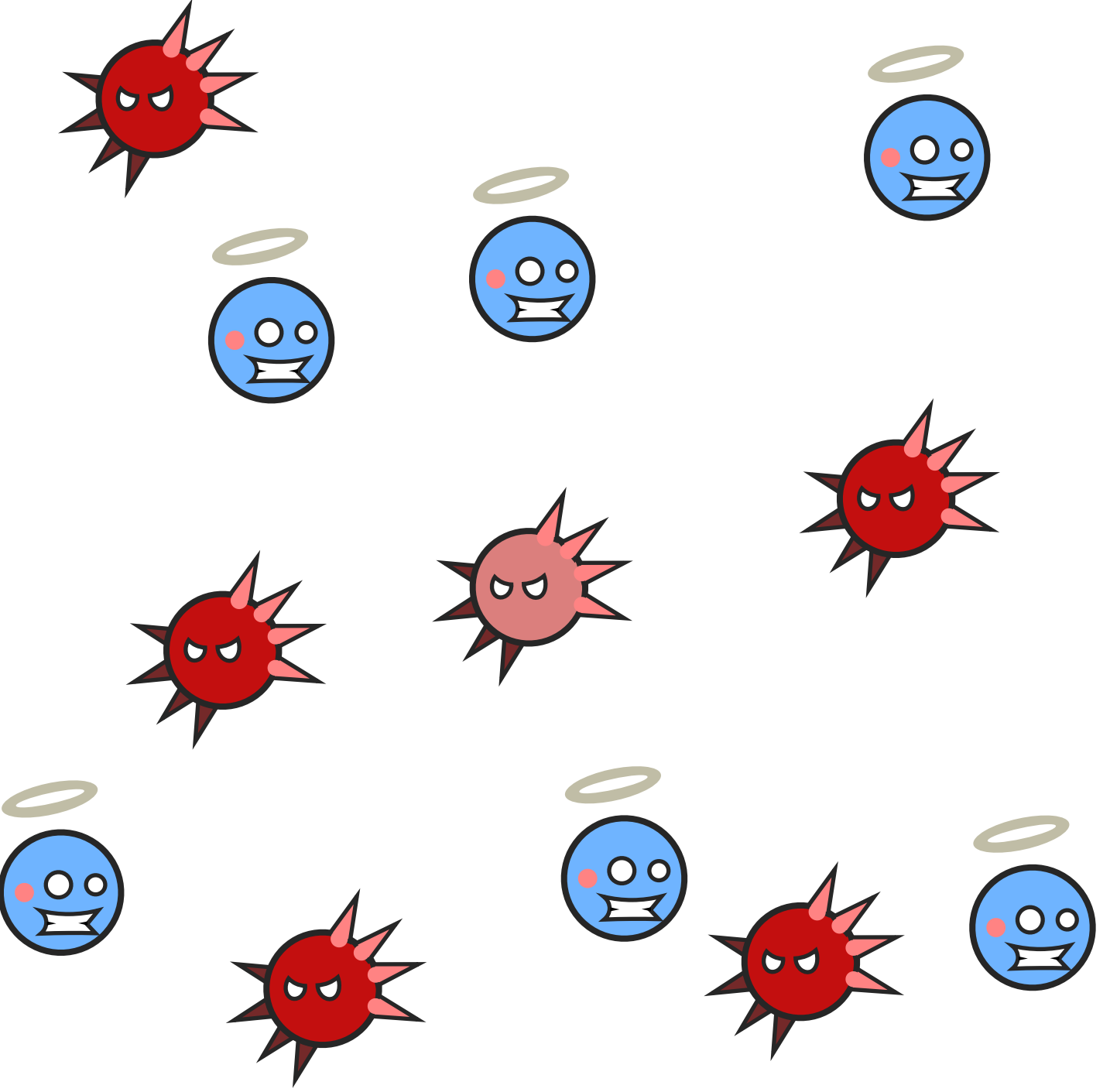
Temporal Inconsistency in Train/Test Sets

Violations use future knowledge in training

Training



Testing



Sources of Experimental Bias (1/3)

Temporal Inconsistency in Train/Test Sets

Violations use future knowledge in training

Training

Testing

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

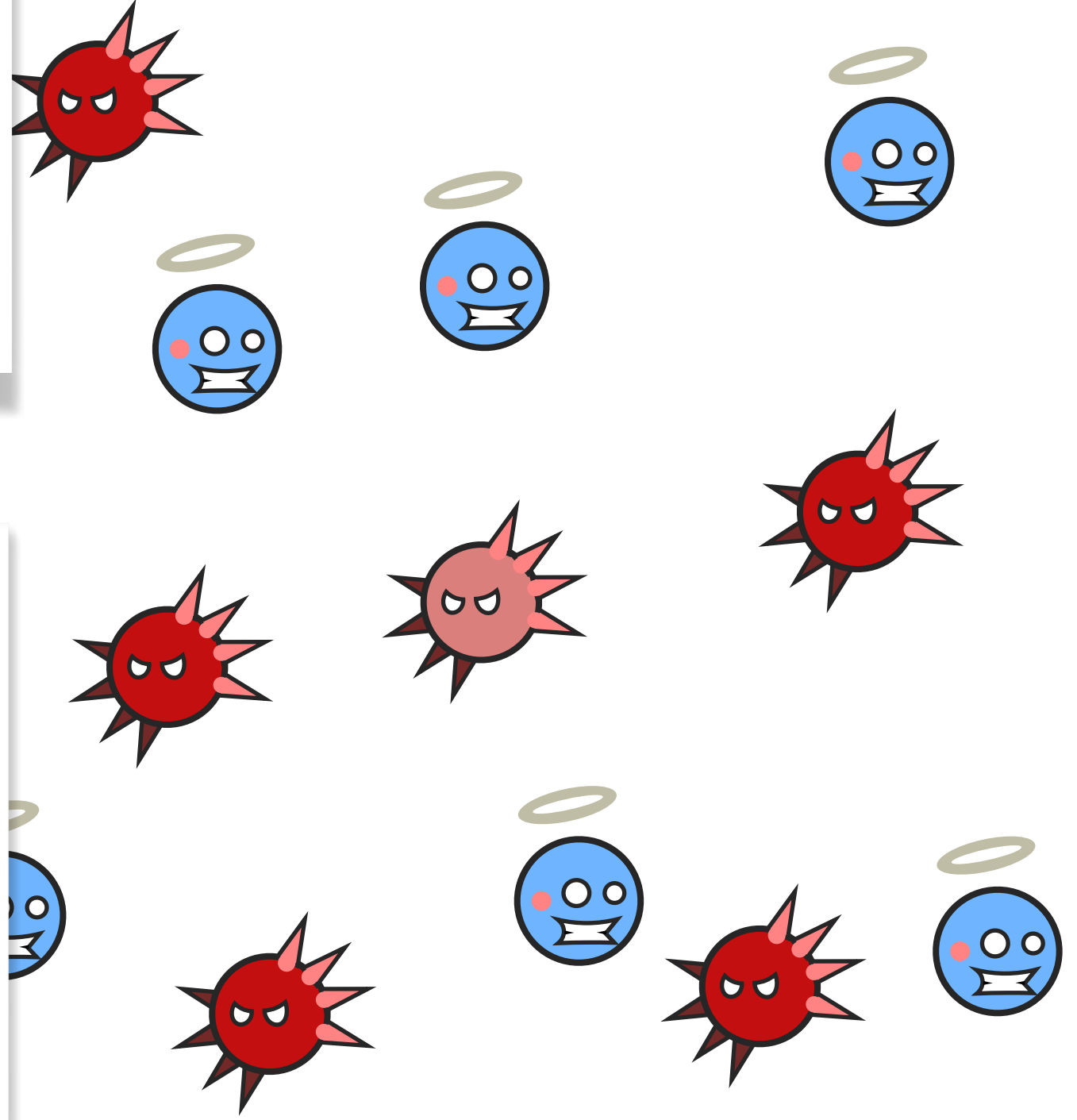
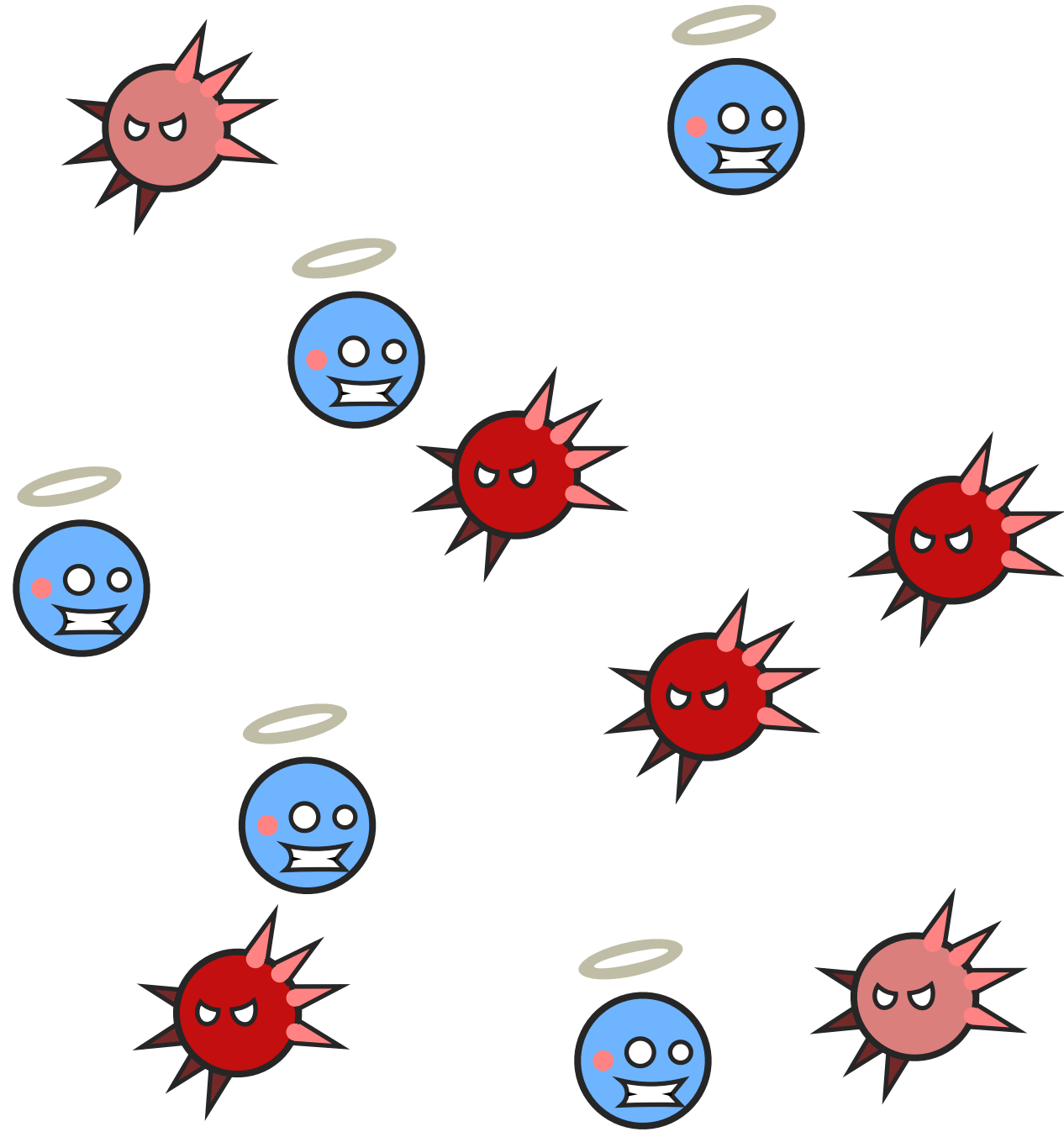
Kevin Allix, Tegawendé F. Bissyandé, Jacques Klein, and Yves Le Traon
SnT - University of Luxembourg

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 Faizullahoy², 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

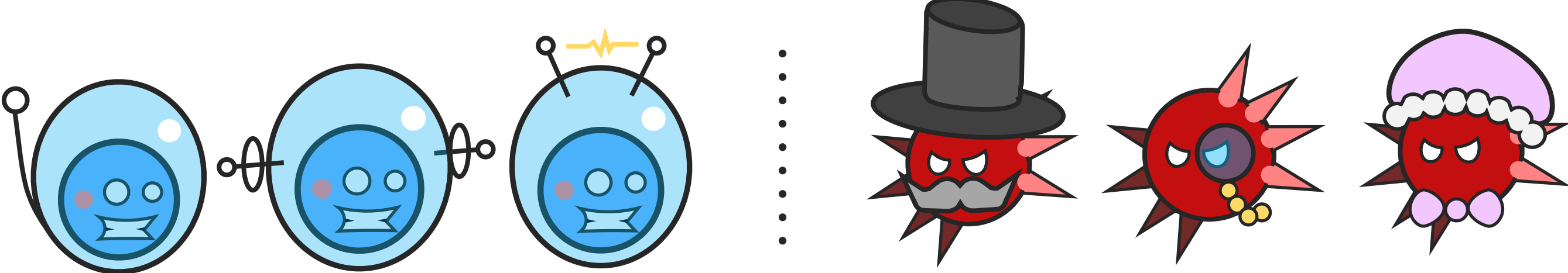


Sources of Experimental Bias (2/3)

Temporal {good|mal}ware inconsistency

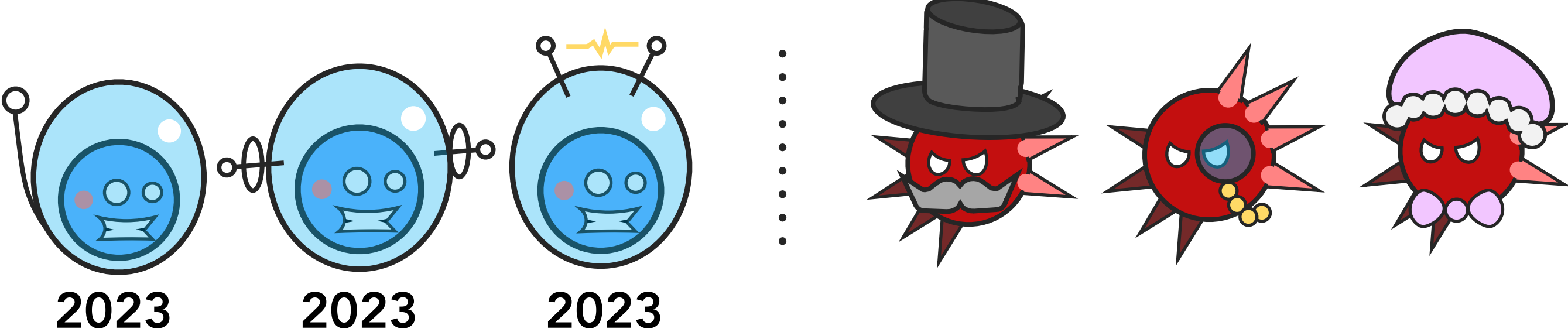
Sources of Experimental Bias (2/3)

Temporal {good|mal}ware inconsistency



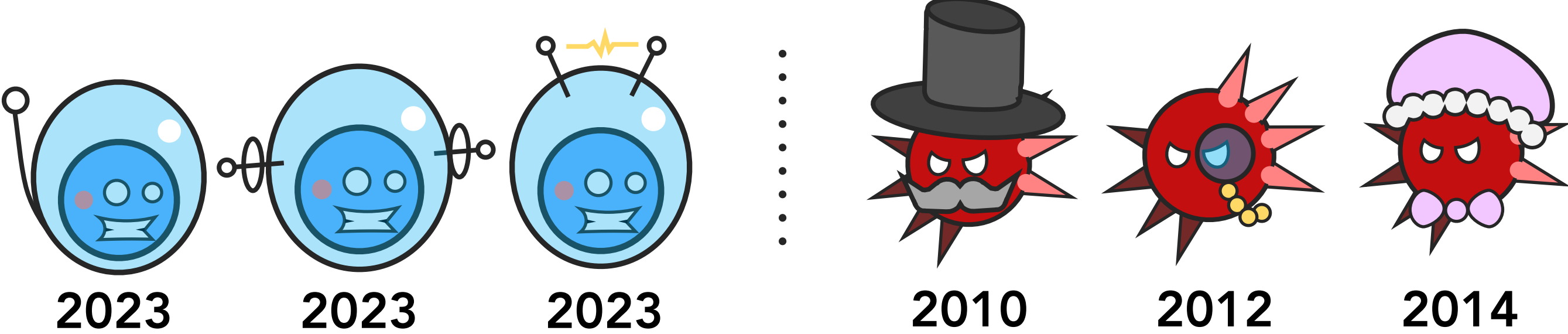
Sources of Experimental Bias (2/3)

Temporal {good|mal}ware inconsistency



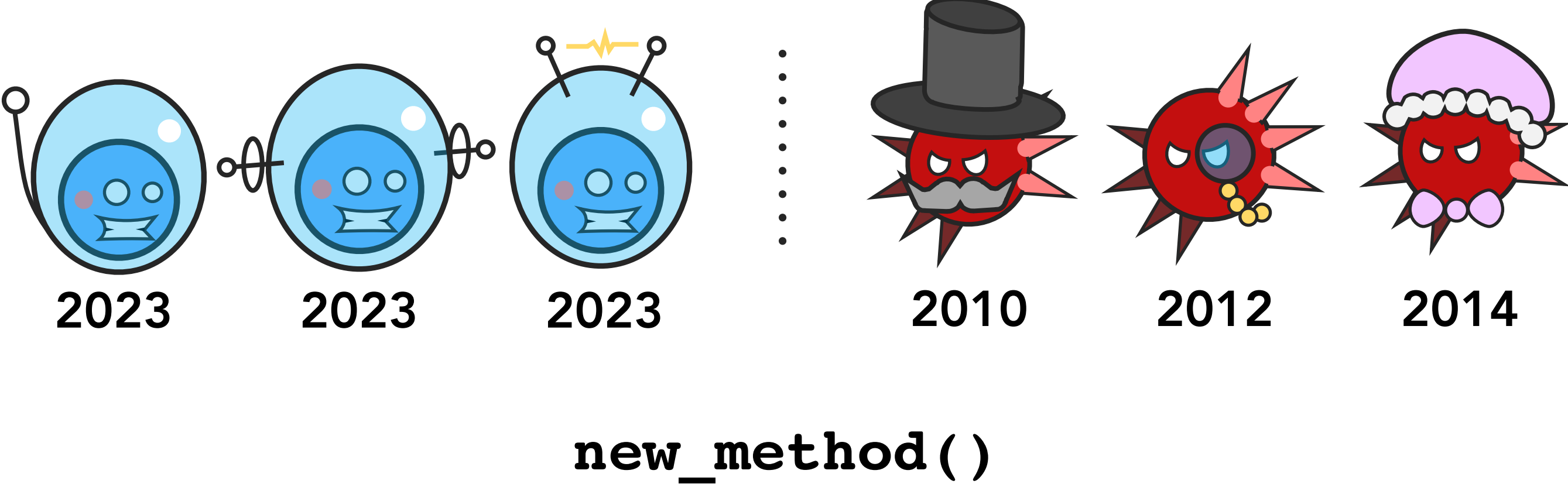
Sources of Experimental Bias (2/3)

Temporal {good|mal}ware inconsistency



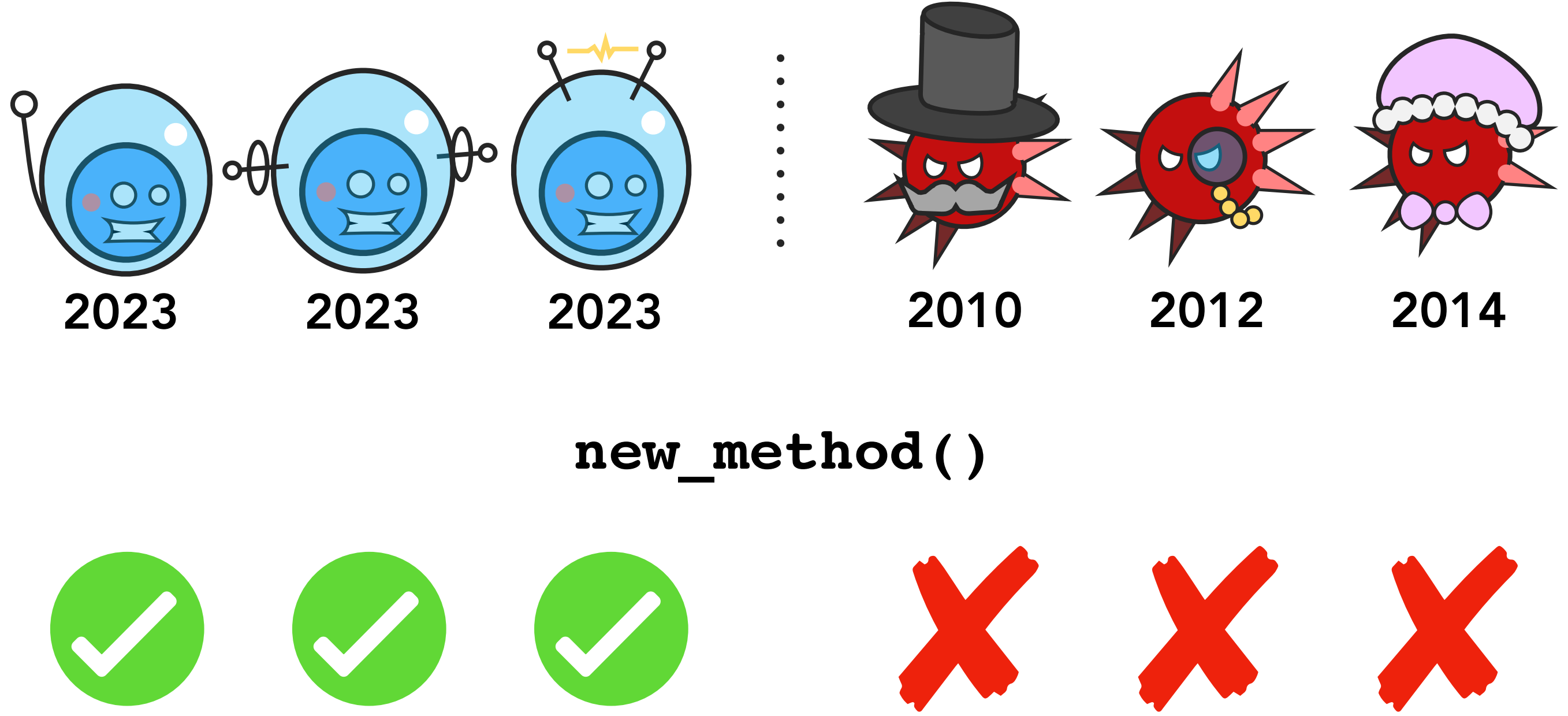
Sources of Experimental Bias (2/3)

Temporal {good|mal}ware inconsistency



Sources of Experimental Bias (2/3)

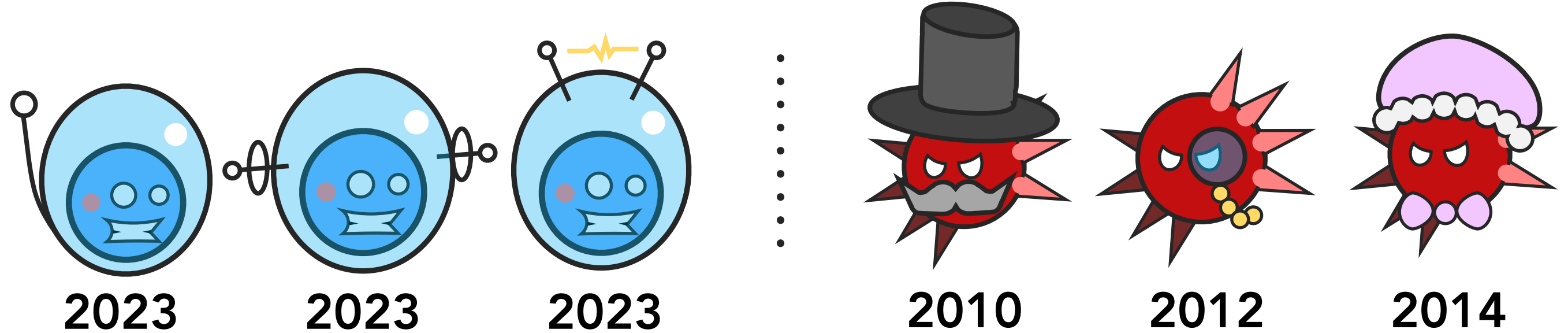
Temporal {good|mal}ware inconsistency



Sources of Experimental Bias (2/3)

Temporal {good|mal}ware inconsistency

Violations may learn artifacts



`new_method()`



Sources of Experimental Bias (3/3)

Unrealistic Test Class Ratio

Sources of Experimental Bias (3/3)

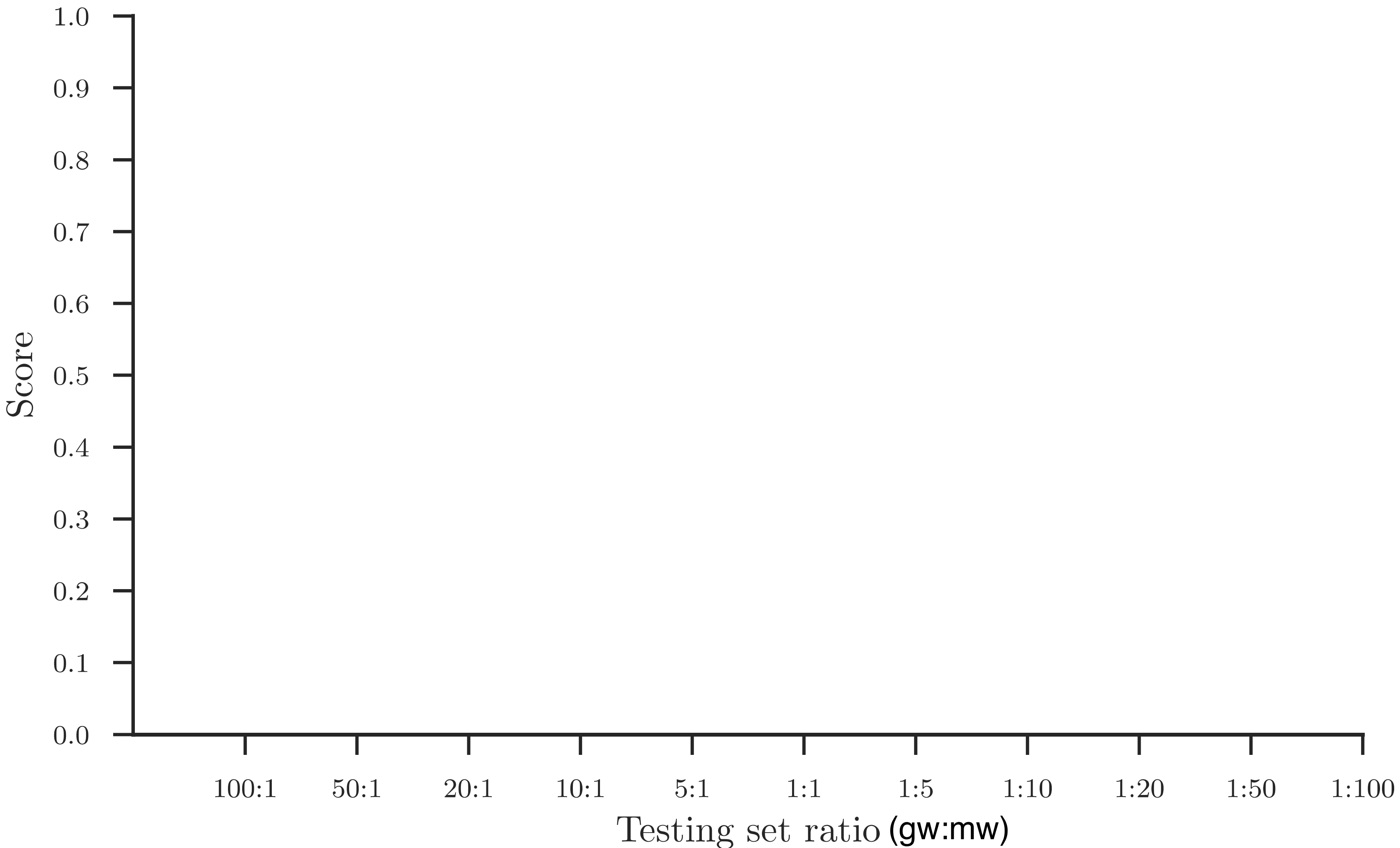
Unrealistic Test Class Ratio

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

Sources of Experimental Bias (3/3)

Unrealistic Test Class Ratio

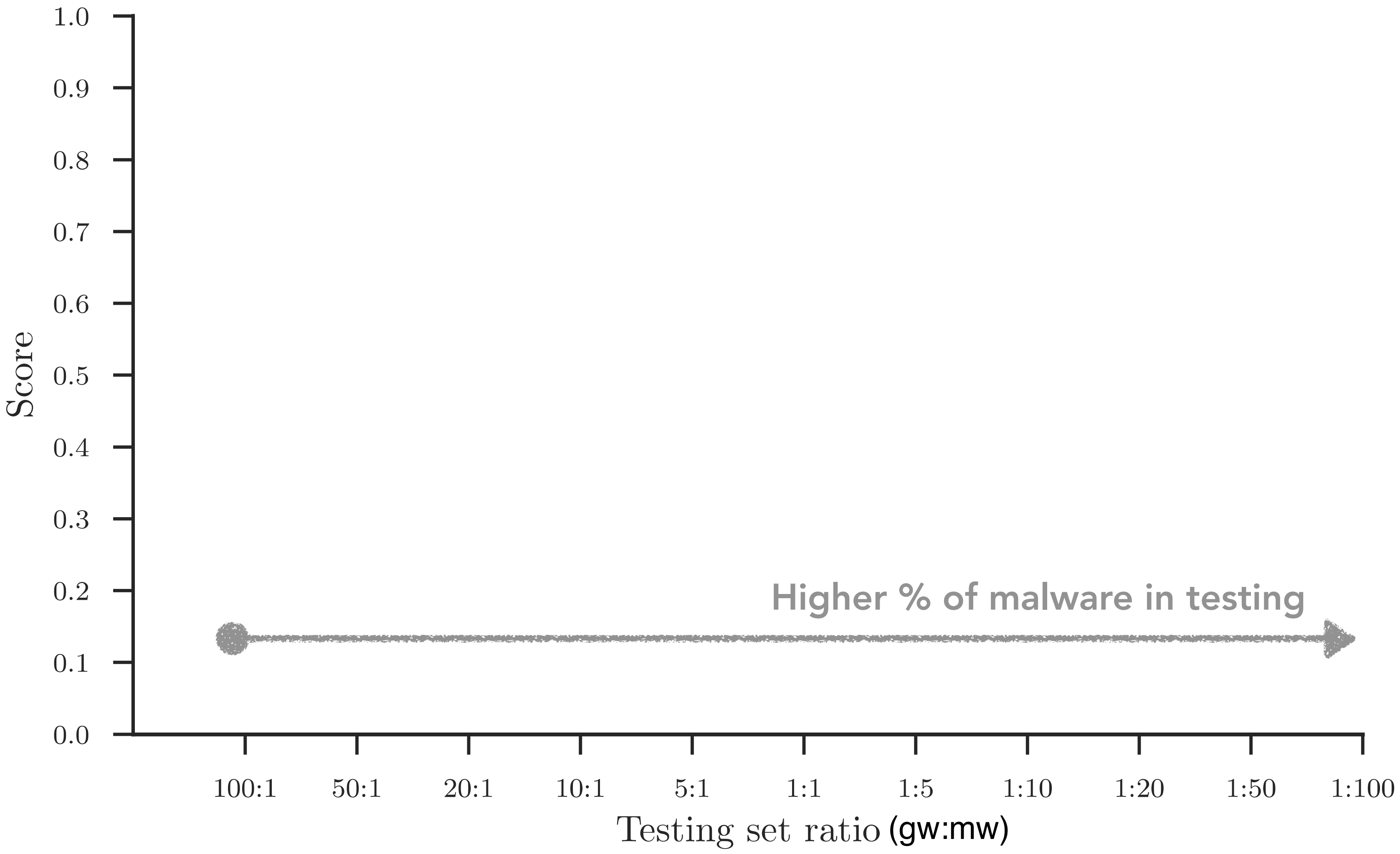
- **Training set:** Fixed
- **Testing set:** Varying % of mw (by downsampling gw)



Sources of Experimental Bias (3/3)

Unrealistic Test Class Ratio

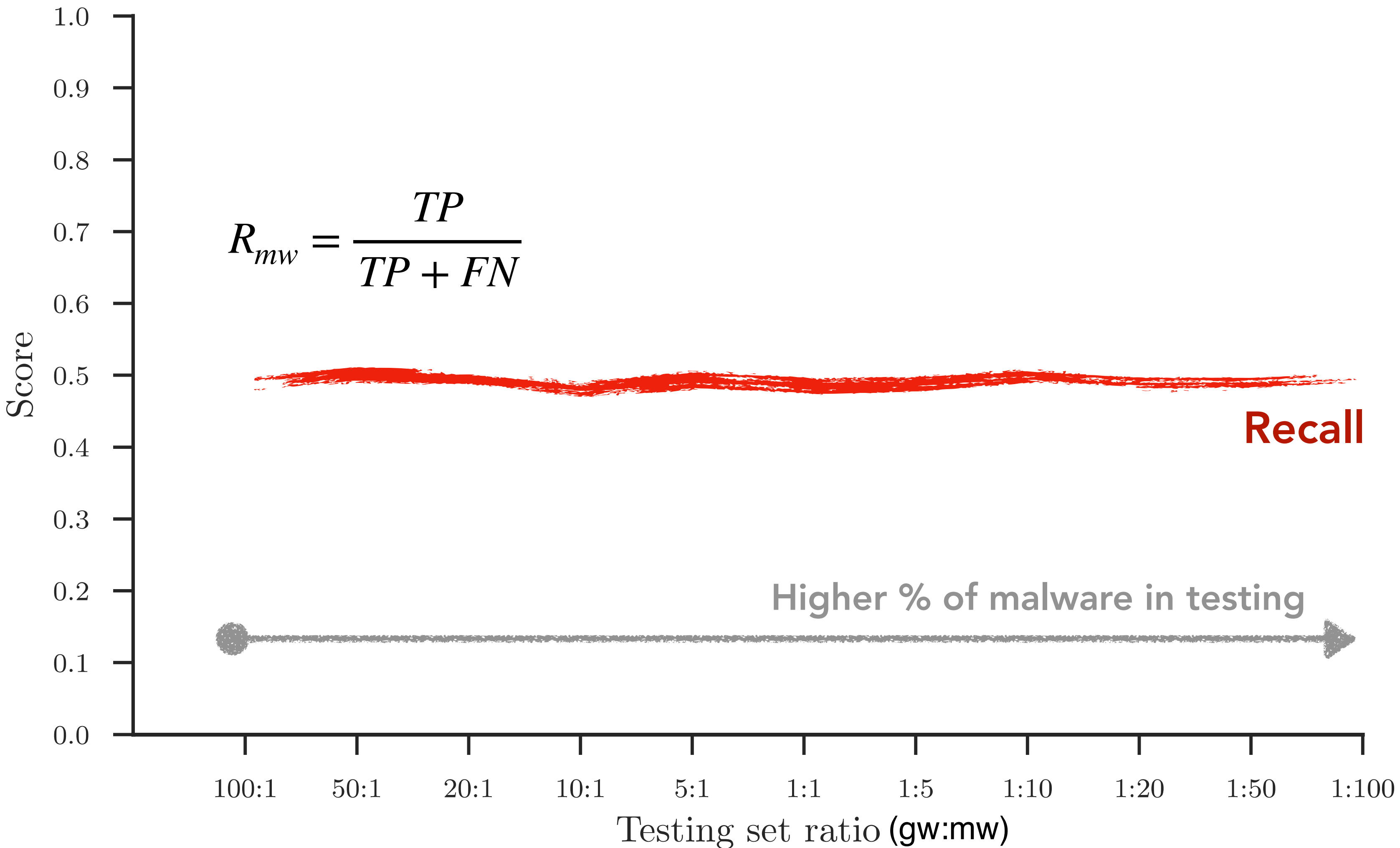
- **Training set:** Fixed
- **Testing set:** Varying % of mw (by downsampling gw)



Sources of Experimental Bias (3/3)

Unrealistic Test Class Ratio

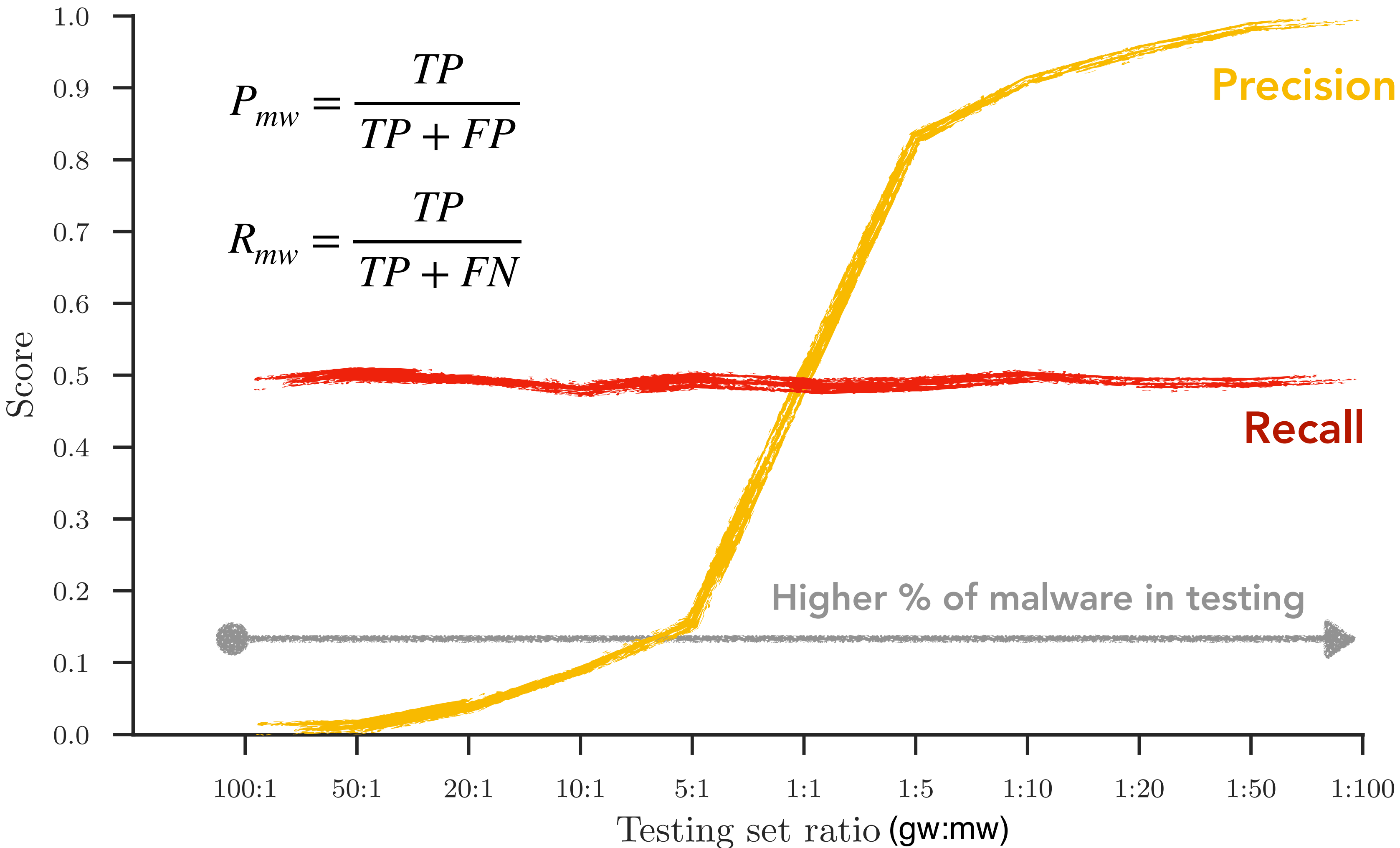
- **Training set:** Fixed
- **Testing set:** Varying % of mw (by downsampling gw)



Sources of Experimental Bias (3/3)

Unrealistic Test Class Ratio

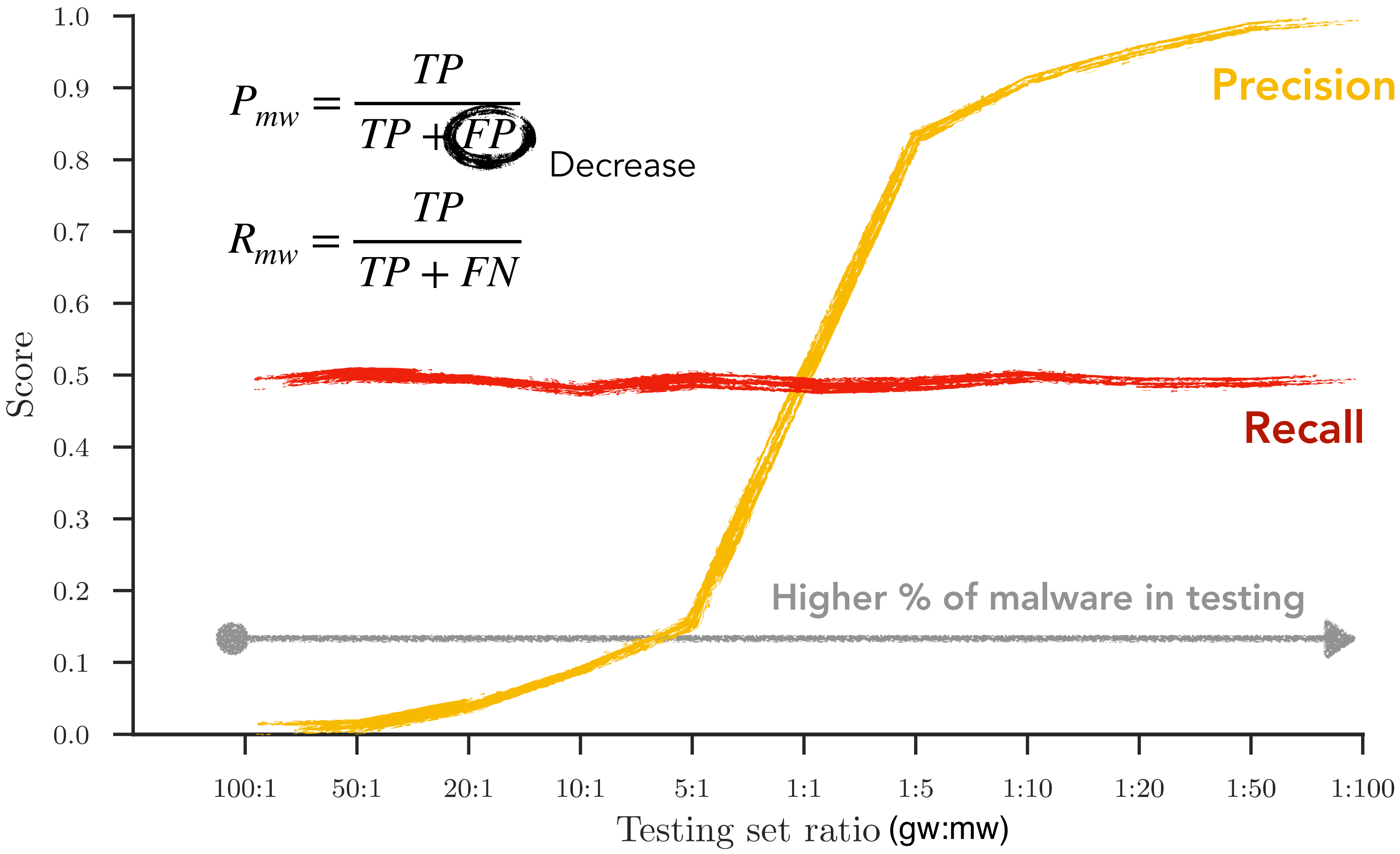
- **Training set:** Fixed
- **Testing set:** Varying % of mw (by downsampling gw)



Sources of Experimental Bias (3/3)

Unrealistic Test Class Ratio

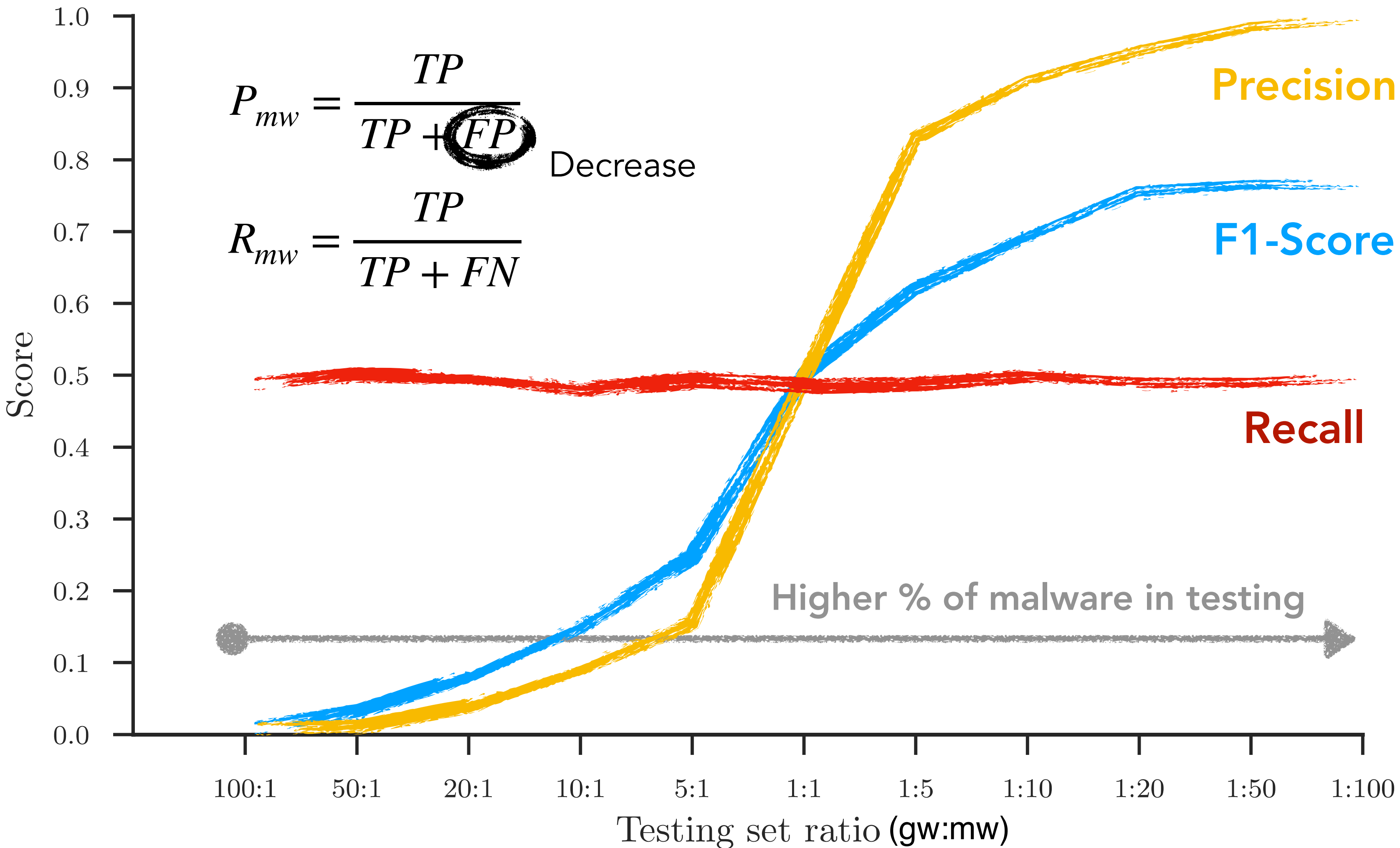
- **Training set:** Fixed
- **Testing set:** Varying % of mw (by downsampling gw)



Sources of Experimental Bias (3/3)

Unrealistic Test Class Ratio

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

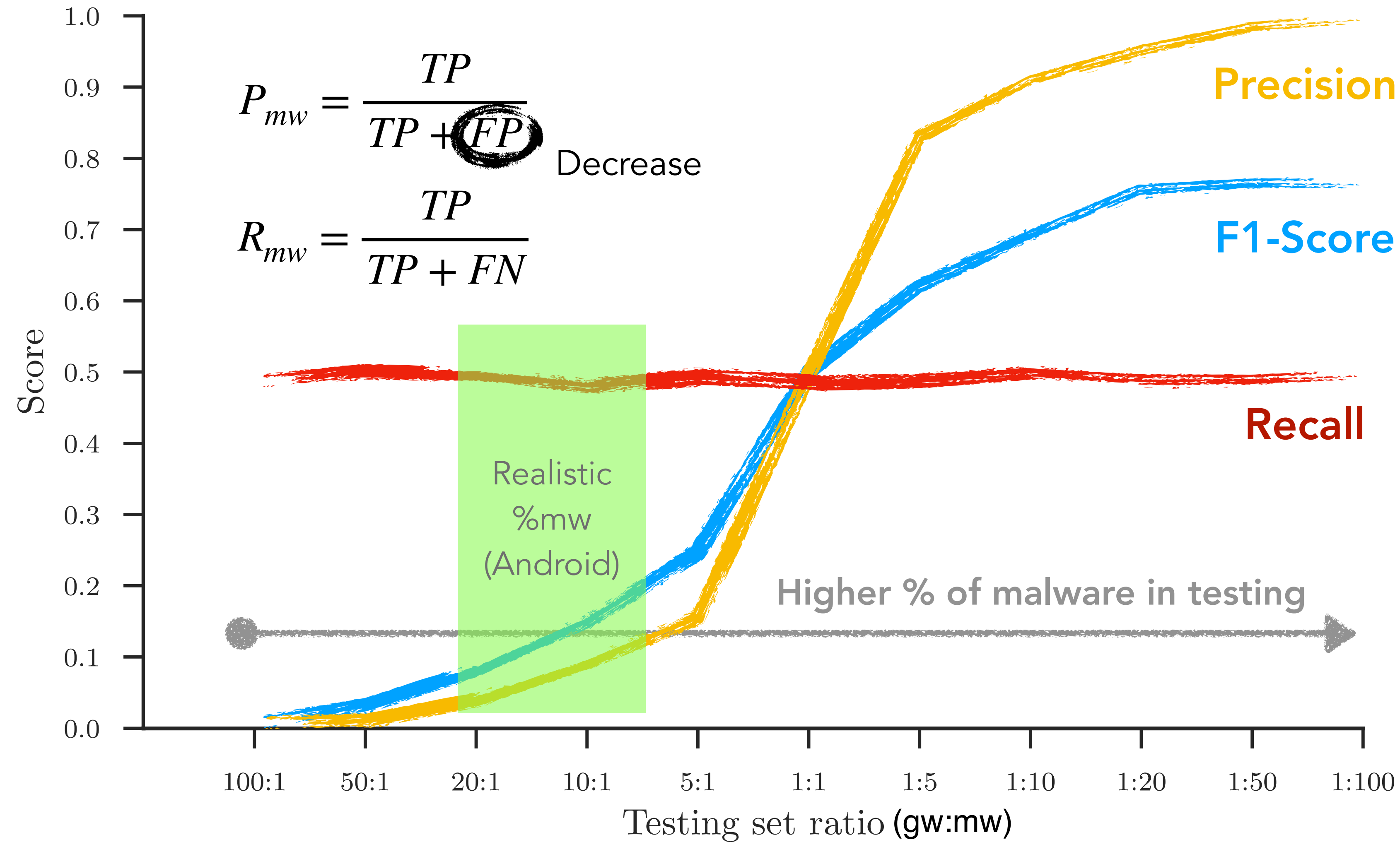


Sources of Experimental Bias (3/3)

Unrealistic Test Class Ratio

Violations produce unrealistic results

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



Prior Work on Adv. Malware

Prior Work on Adv. Malware

Prior work was fundamental to initially explore problem-space attacks.
We propose a principled approach that supports reasoning.

Available Transformations

- Limiting #features modified [ESORICS'17]

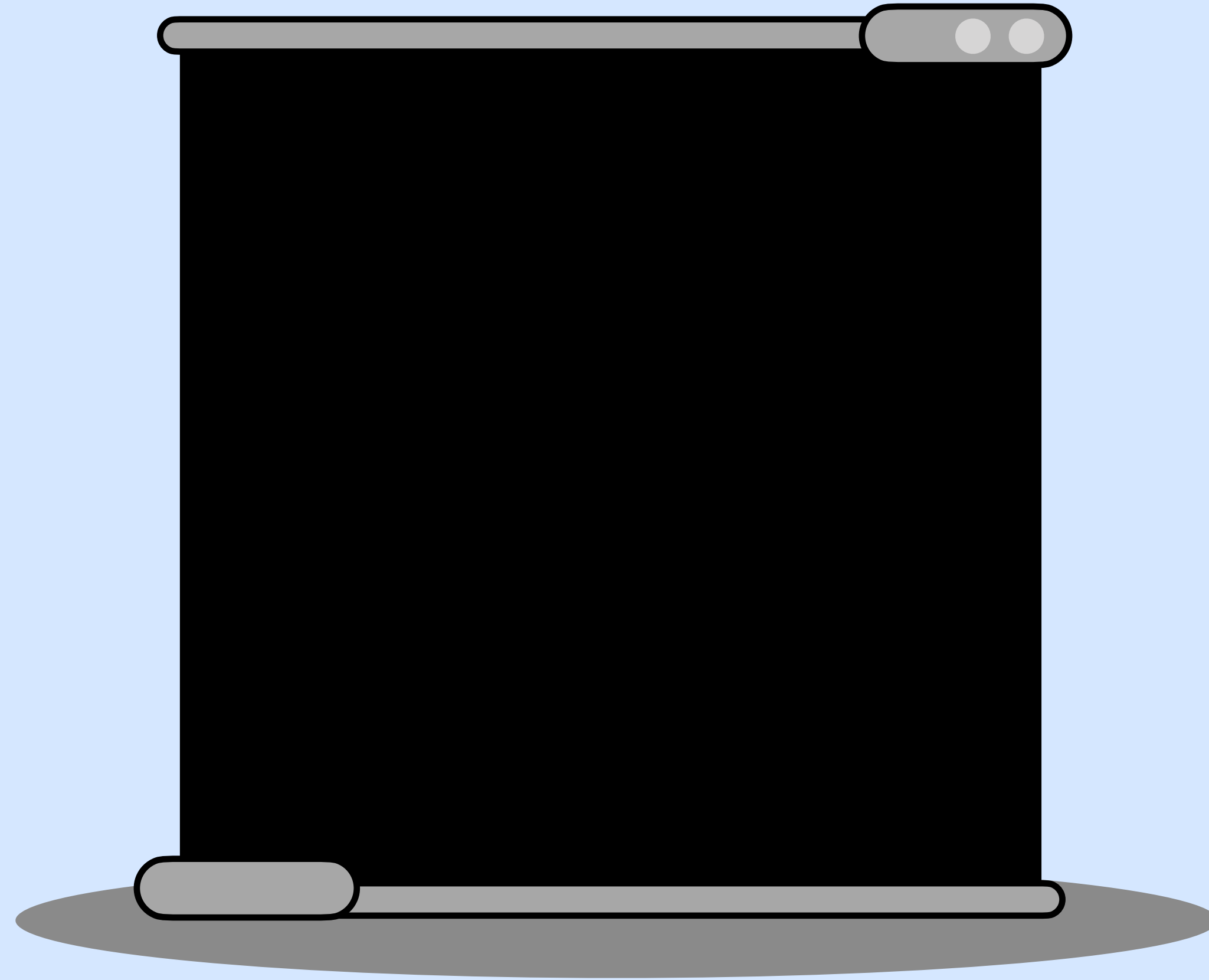
Robustness to Preprocessing

- Removable unused permissions [ESORICS'17]
- Removal code (unreachable, no-op) [EUSIPCO'18, RAID'18]
- Unclear [ACSAC'18]

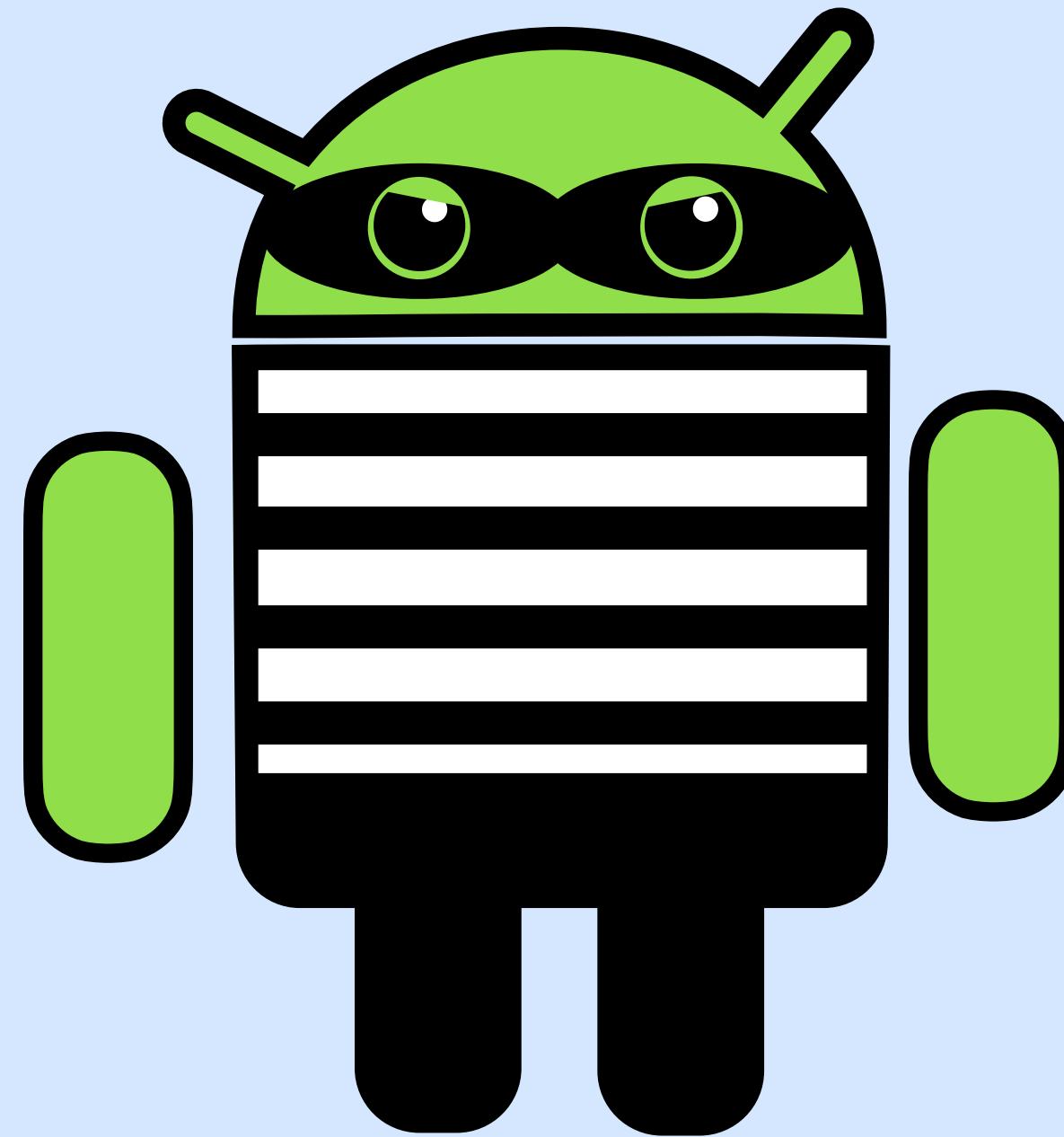
Preserved Semantics

- Highly unstable transformations [ACSAC'18]

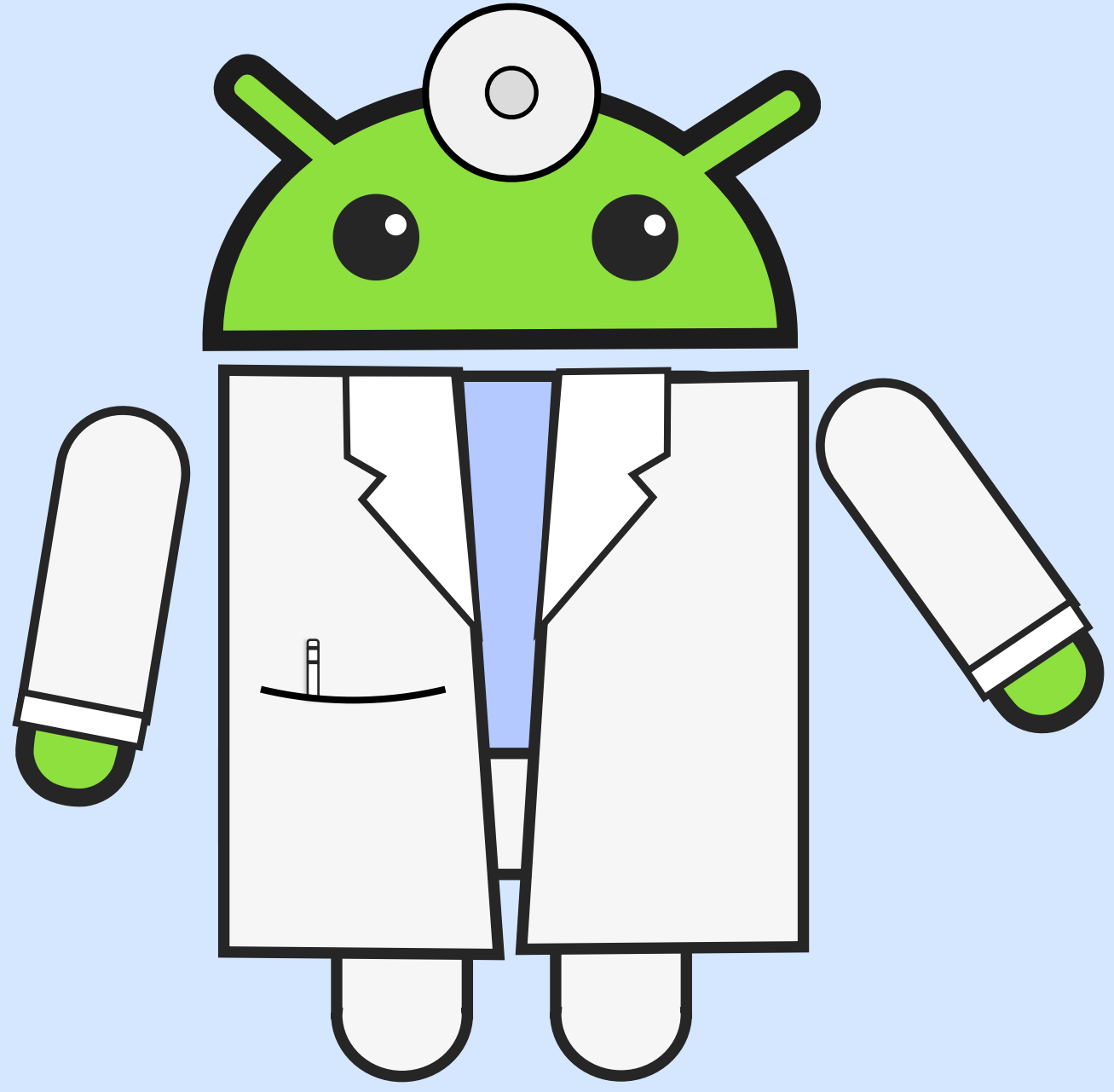
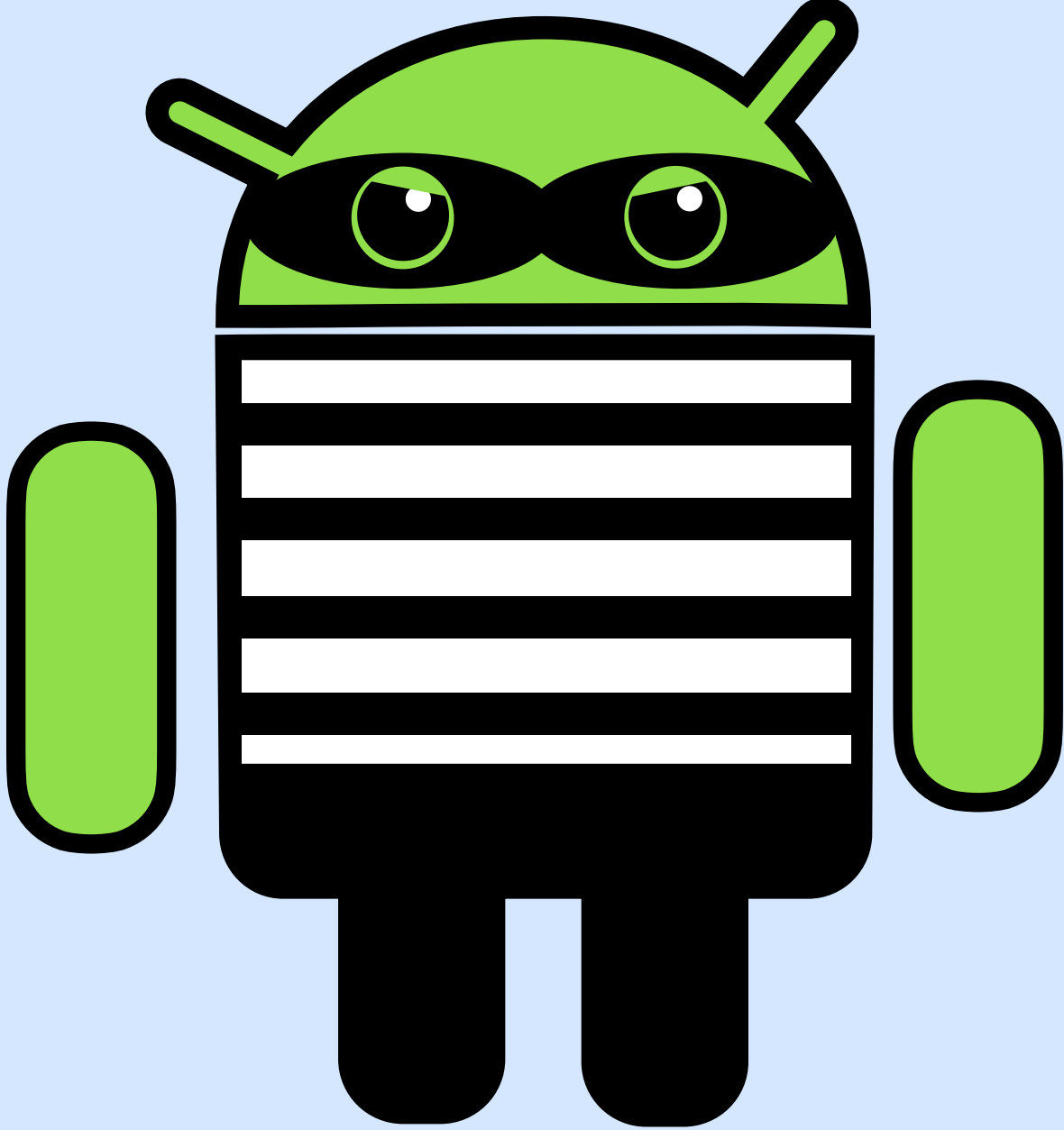
Our Android Attack




Our Android Attack

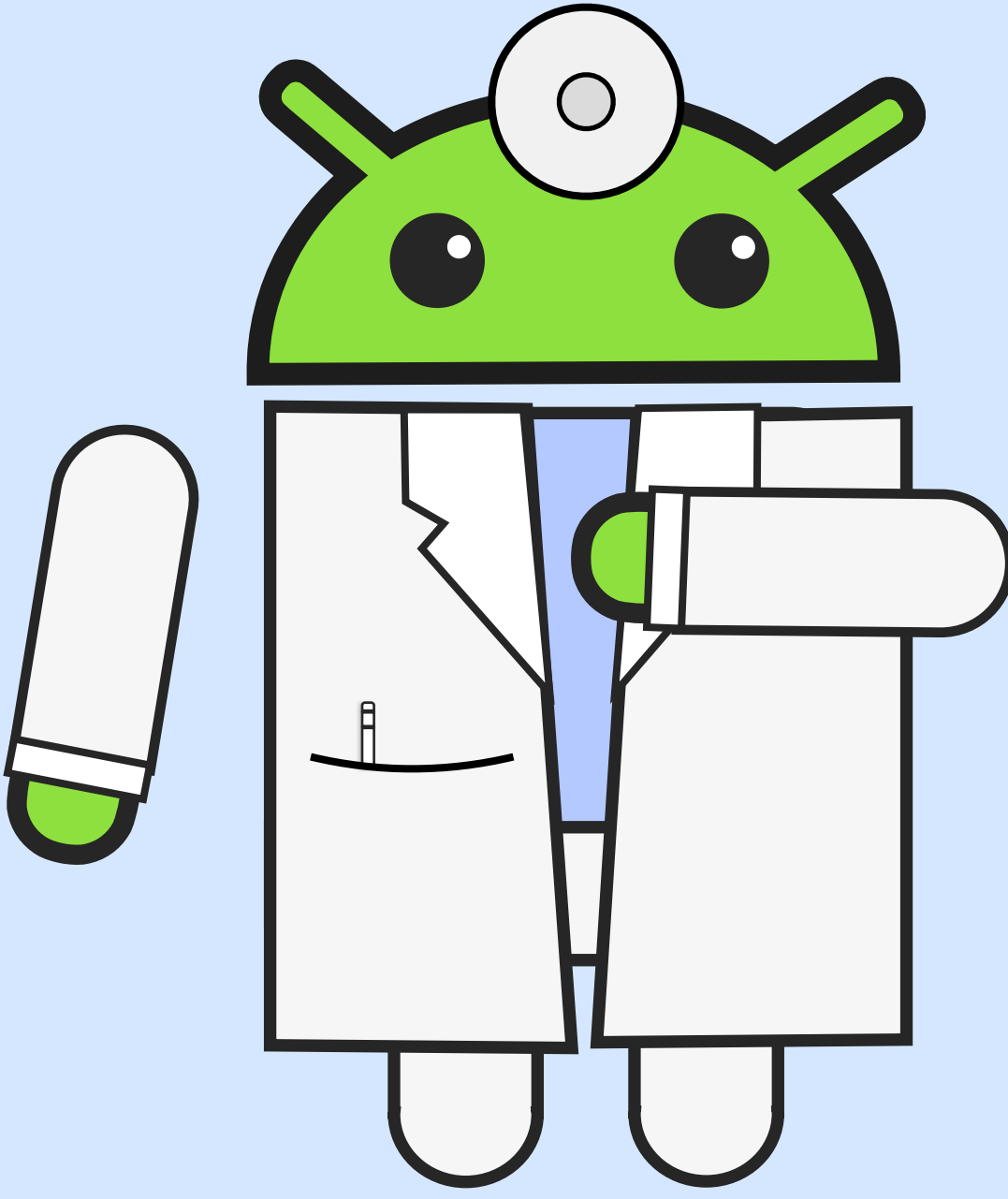
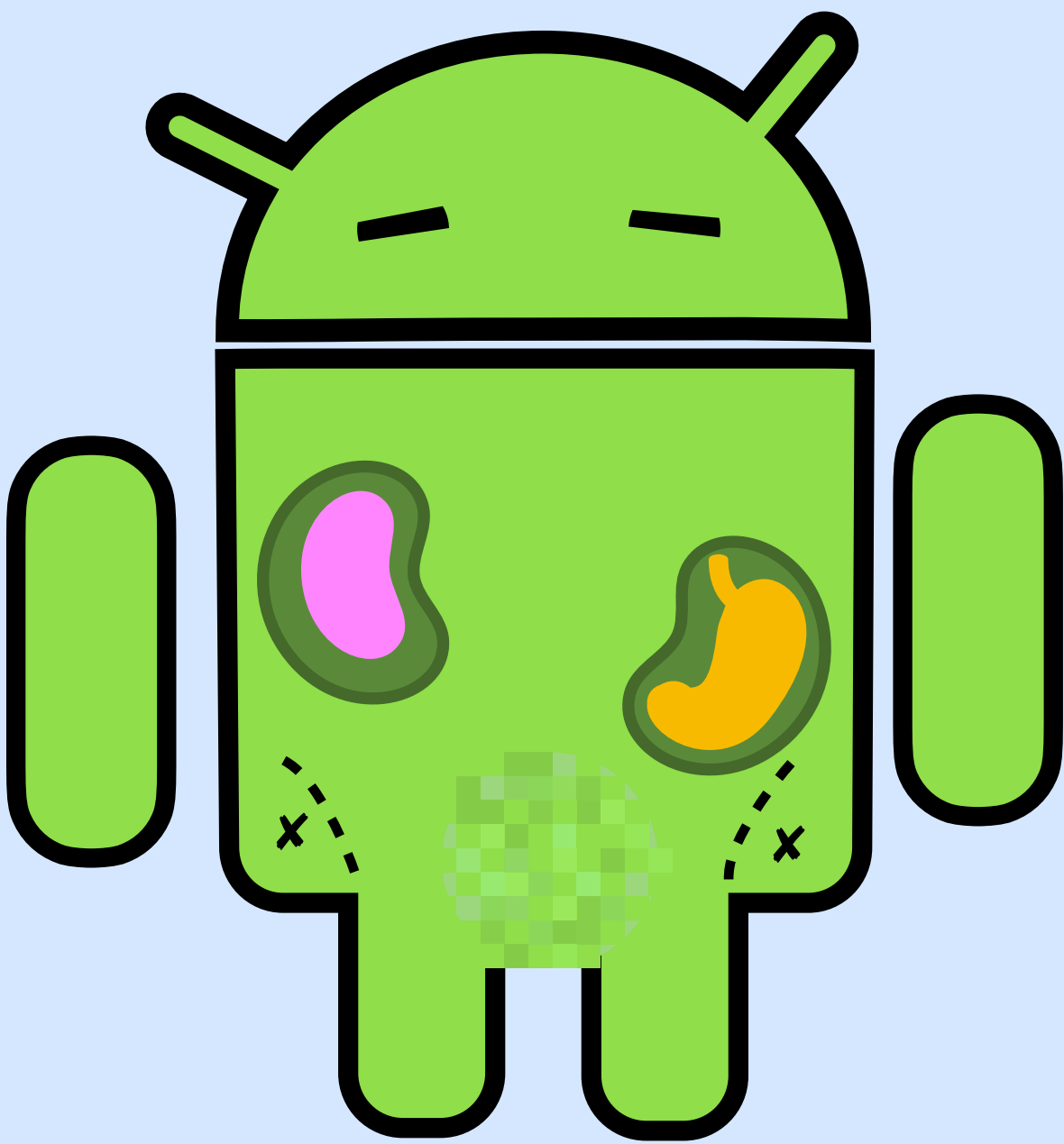


Our Android Attack





Our Android Attack

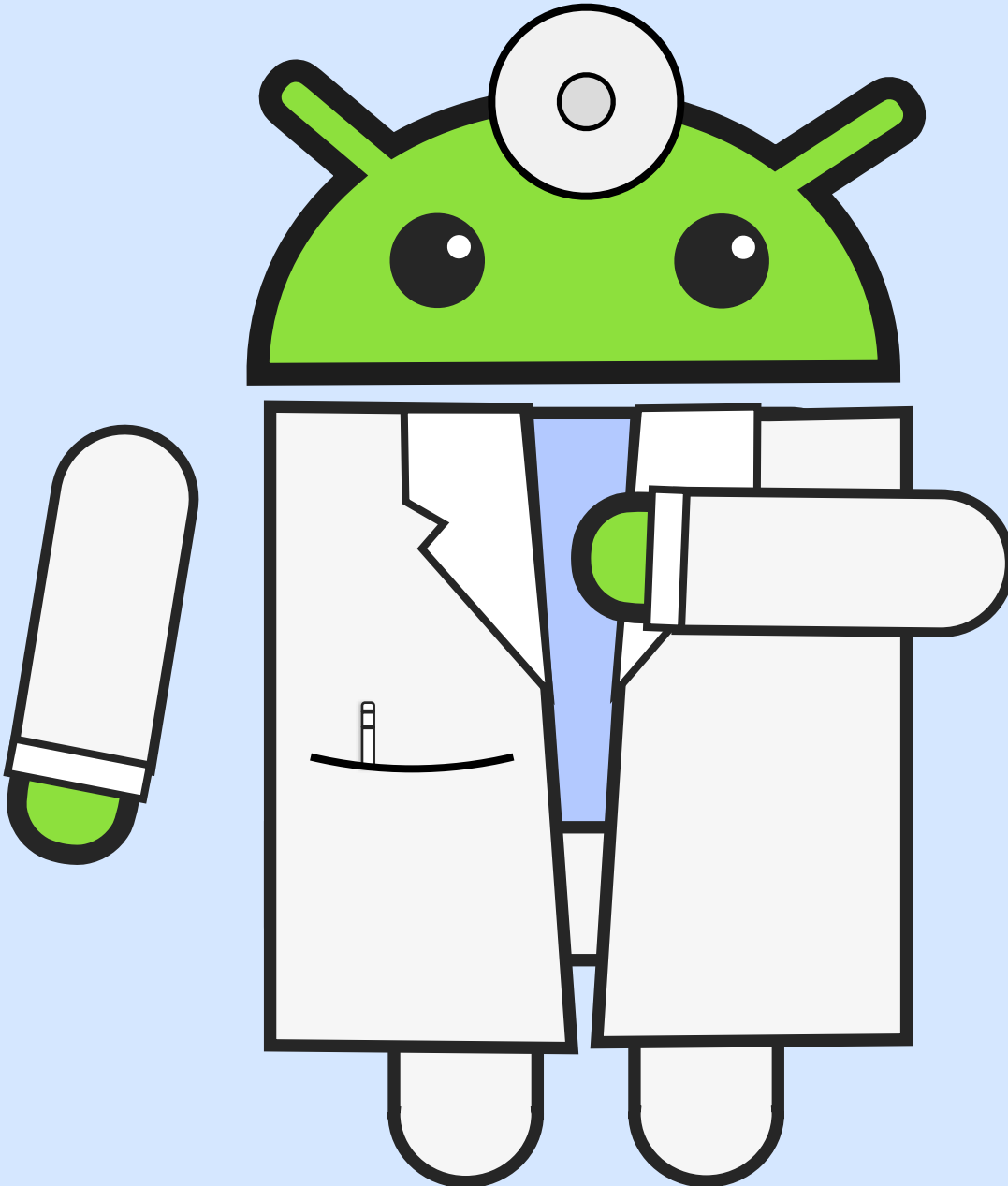
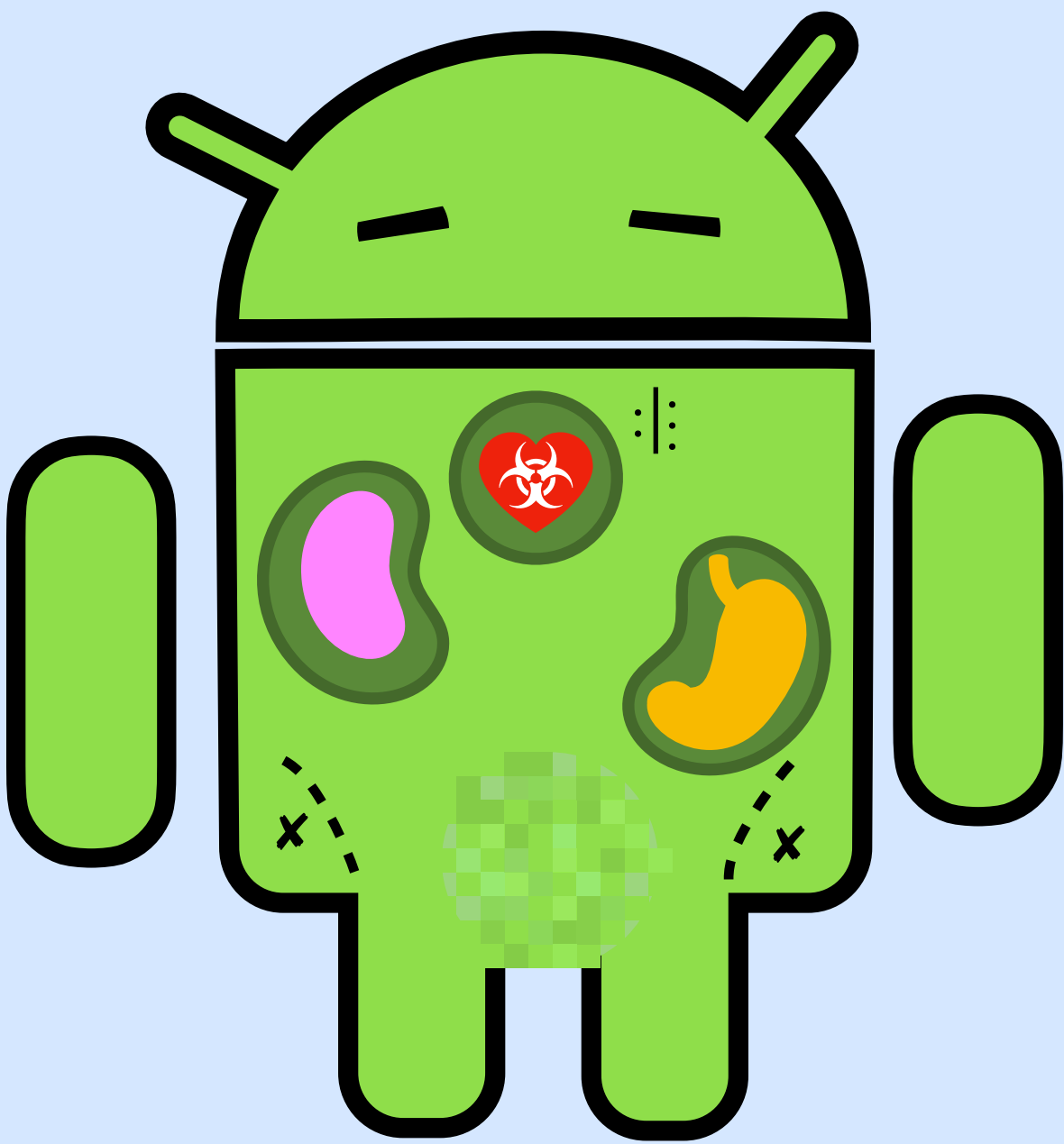
 **Available Transformations**
Code addition through automated software transplantation.



Our Android Attack

 **Available Transformations**
Code addition through automated software transplantation.

 **Preserved Semantics**
Malicious semantics preserved by construction using opaque predicates (inserted code is not executed at runtime).



Our Android Attack



Available Transformations

Code addition through automated software transplantation.



Preserved Semantics

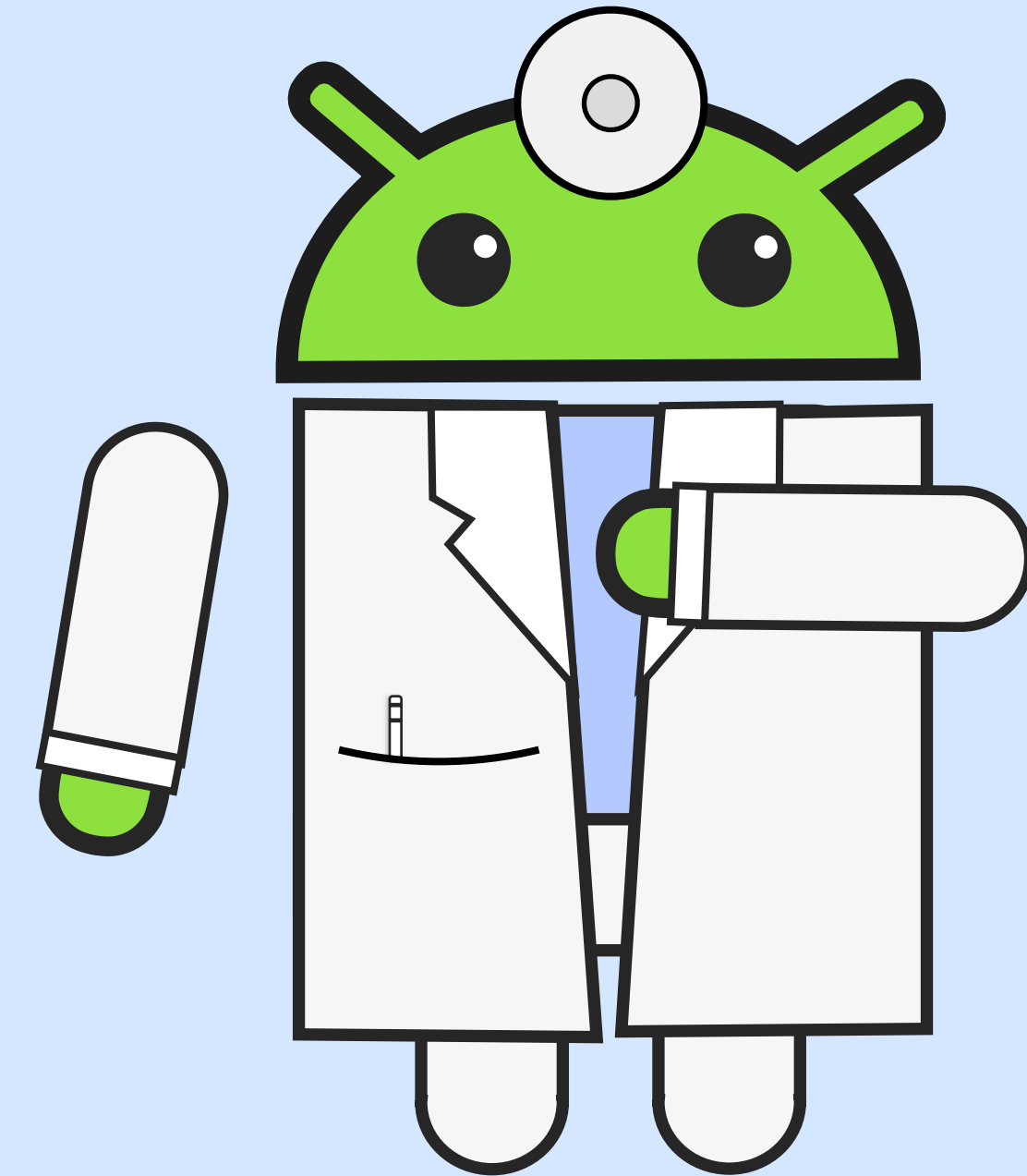
Malicious semantics preserved by construction using opaque predicates (inserted code is not executed at runtime).

Robustness to Preprocessing

We're robust to:



- removal of redundant code
- undeclared variables
- unlinked resources
- undefined references
- naming conflicts
- no-op instructions.



Our Android Attack



Available Transformations

Code addition through automated software transplantation.



Preserved Semantics

Malicious semantics preserved by construction using opaque predicates (inserted code is not executed at runtime).

Robustness to Preprocessing

We're robust to:

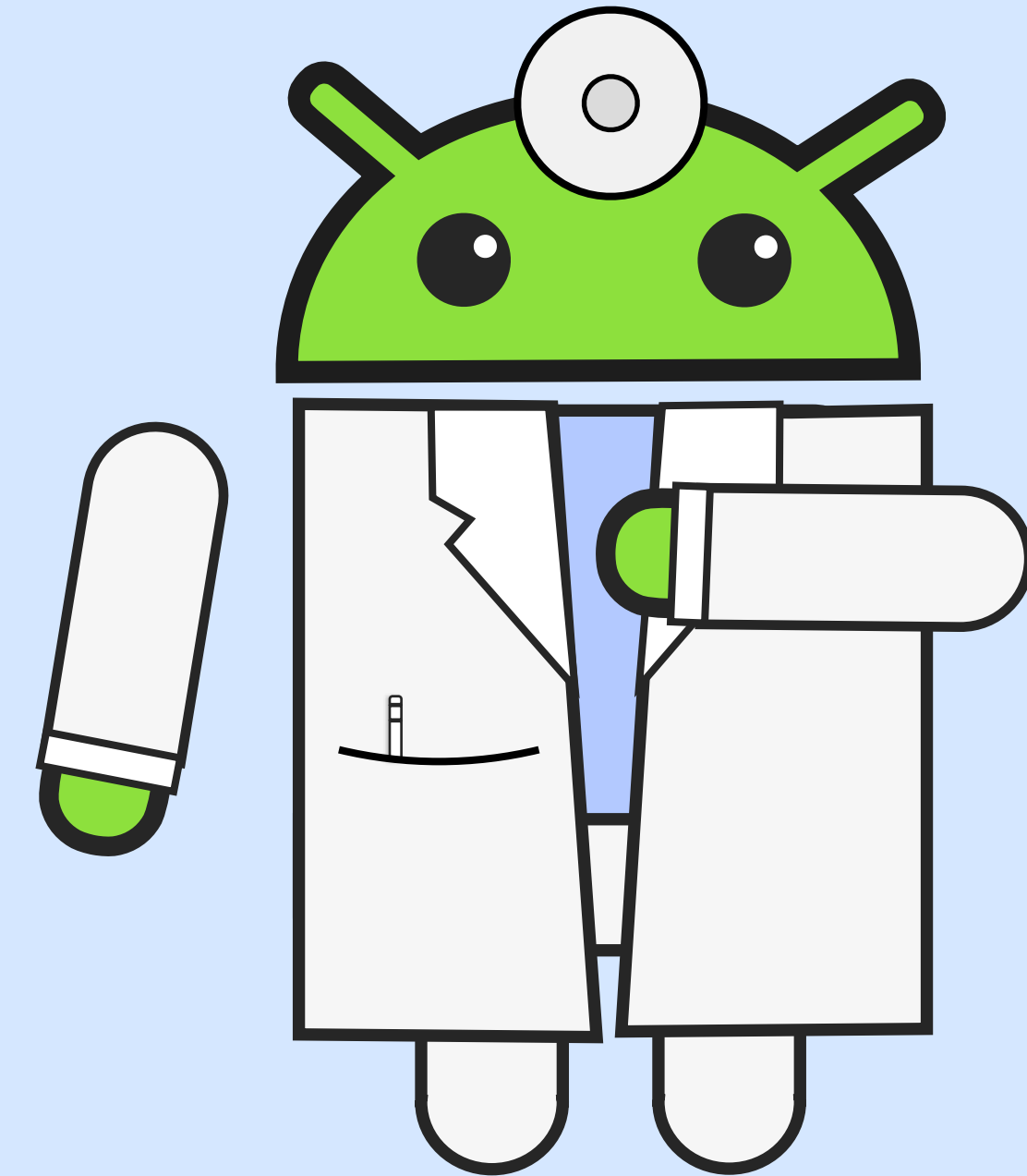
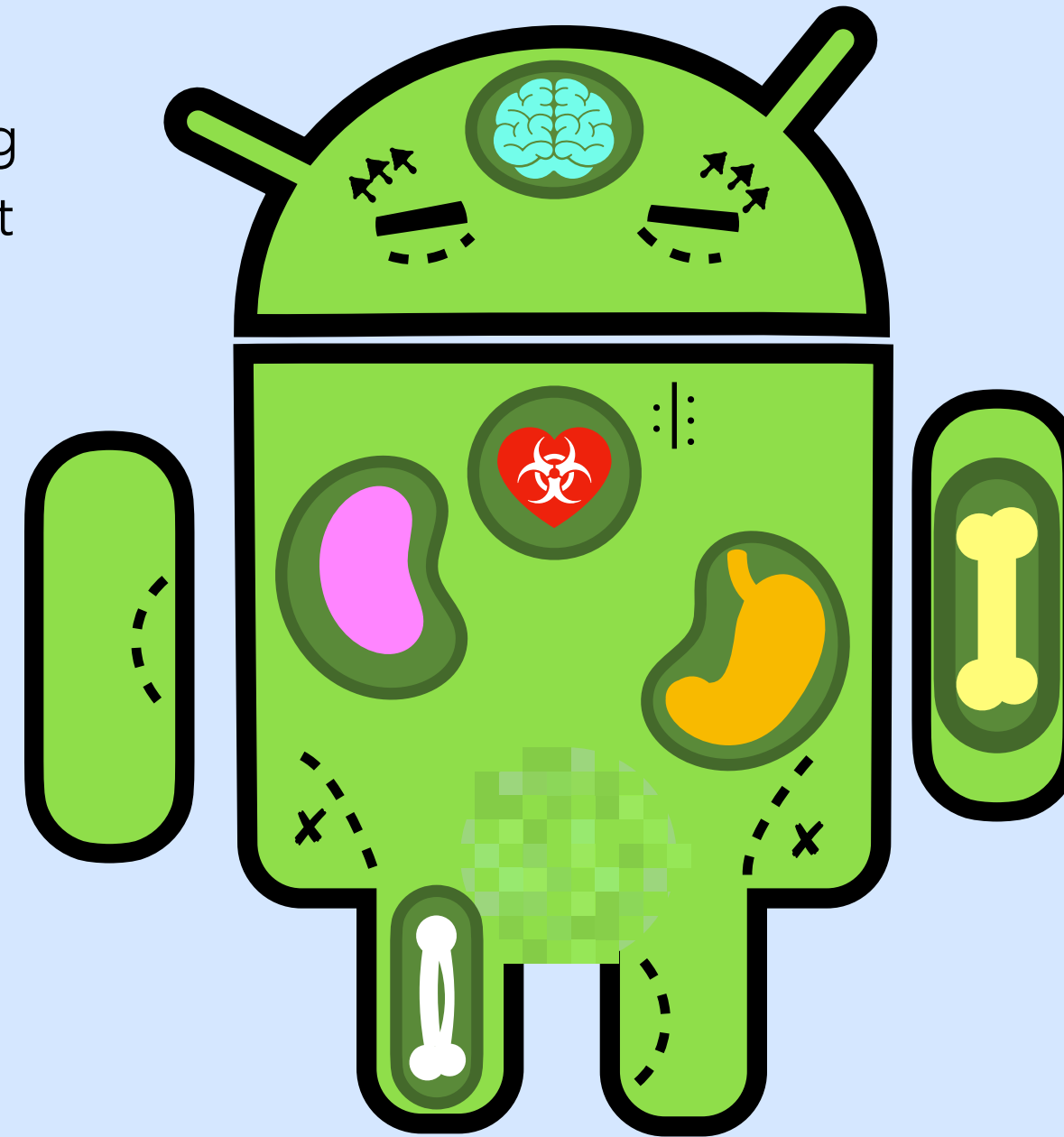


- removal of redundant code
- undeclared variables
- unlinked resources
- undefined references
- naming conflicts
- no-op instructions.

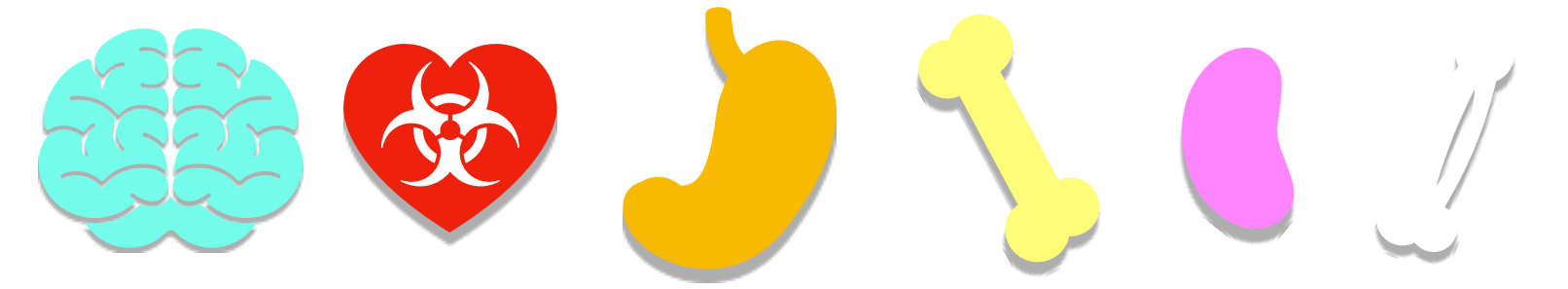


Plausibility

Only realistic code is injected (rather than orphaned urls, api calls, etc.)
Mutated apps install and start on an emulator.



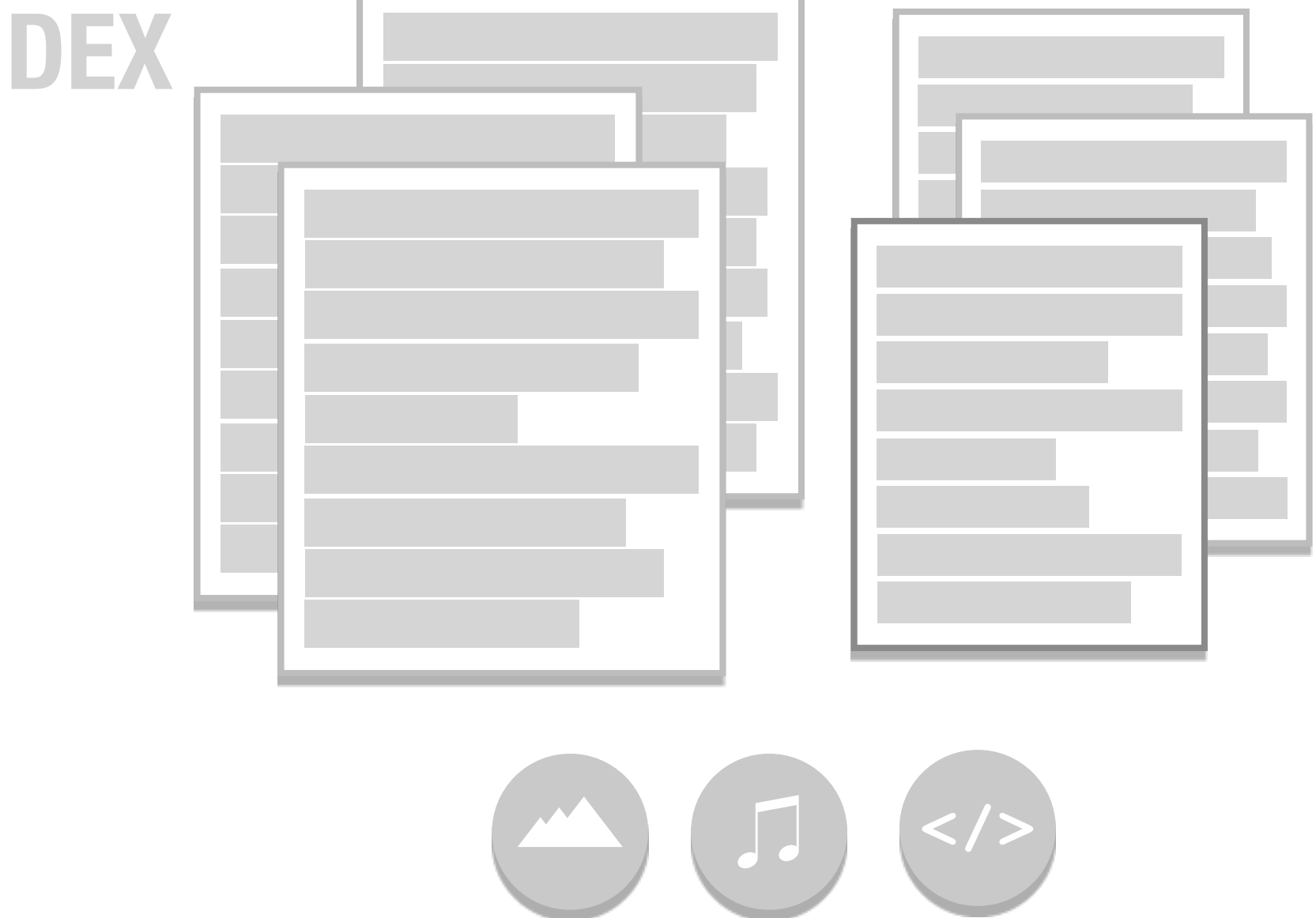
Organ Harvesting



Organ Harvesting

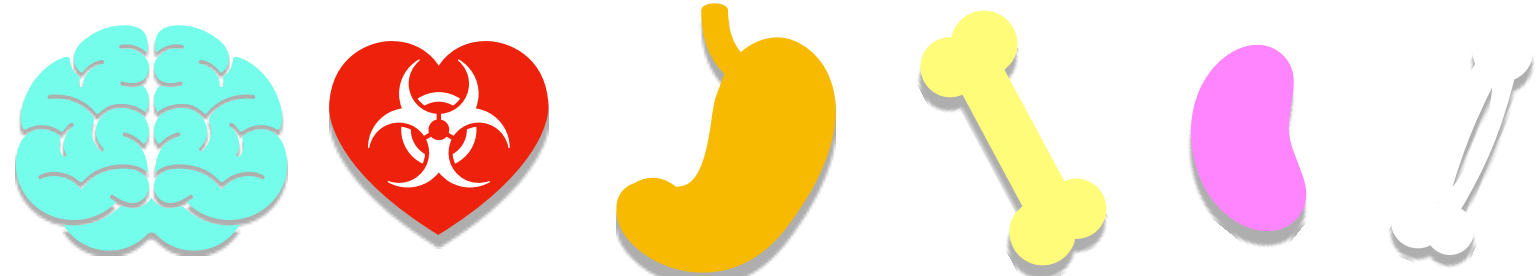


1 Identify feature entry point

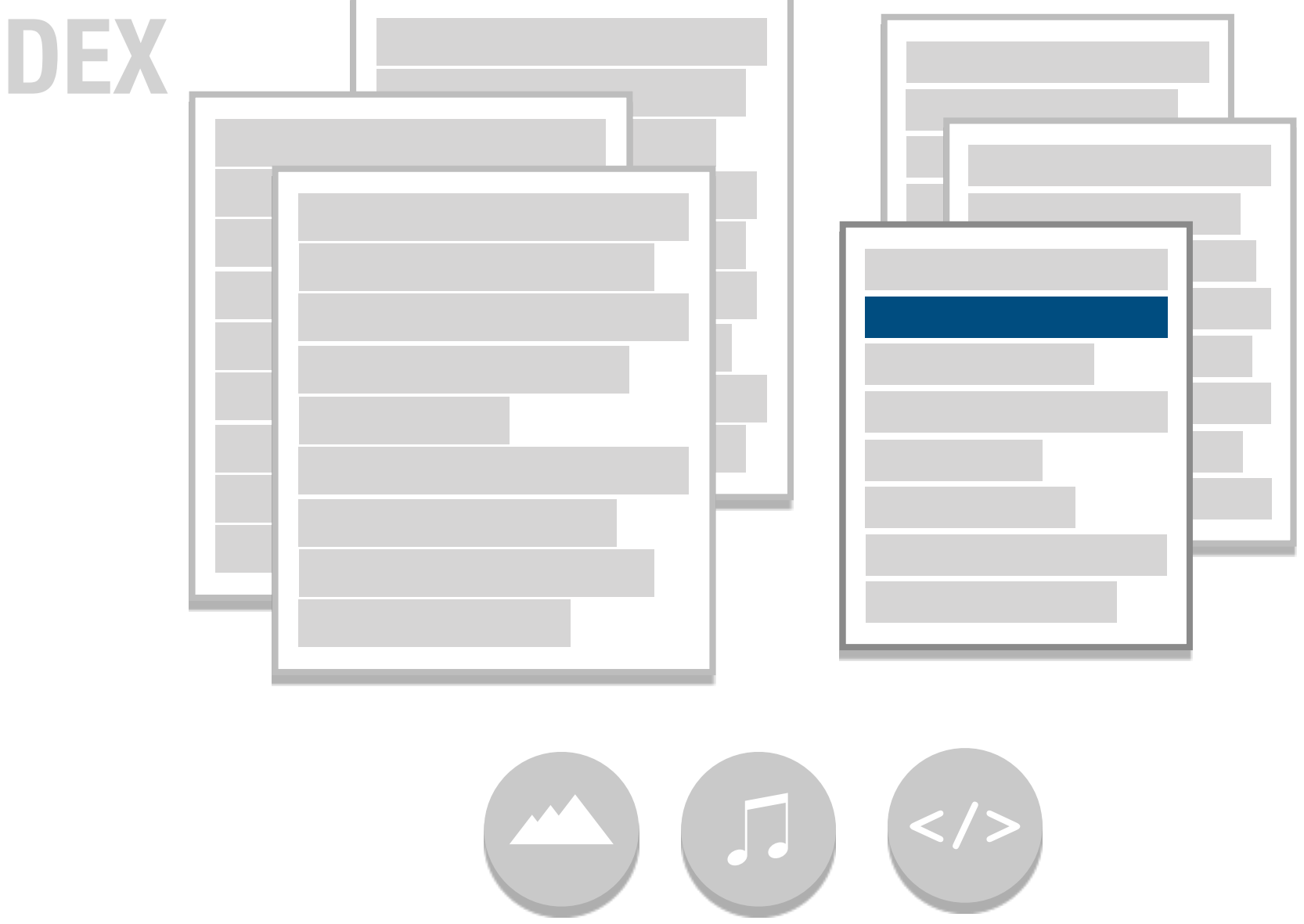


Identify activity in dex

Organ Harvesting

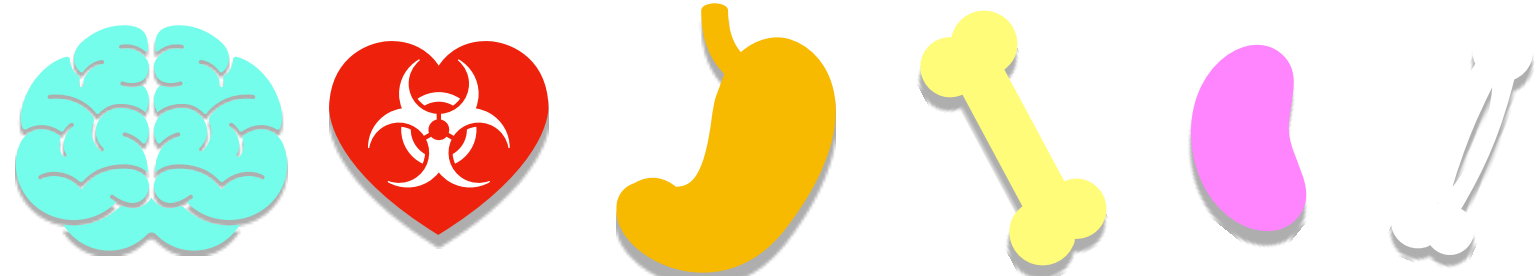


1 Identify feature entry point

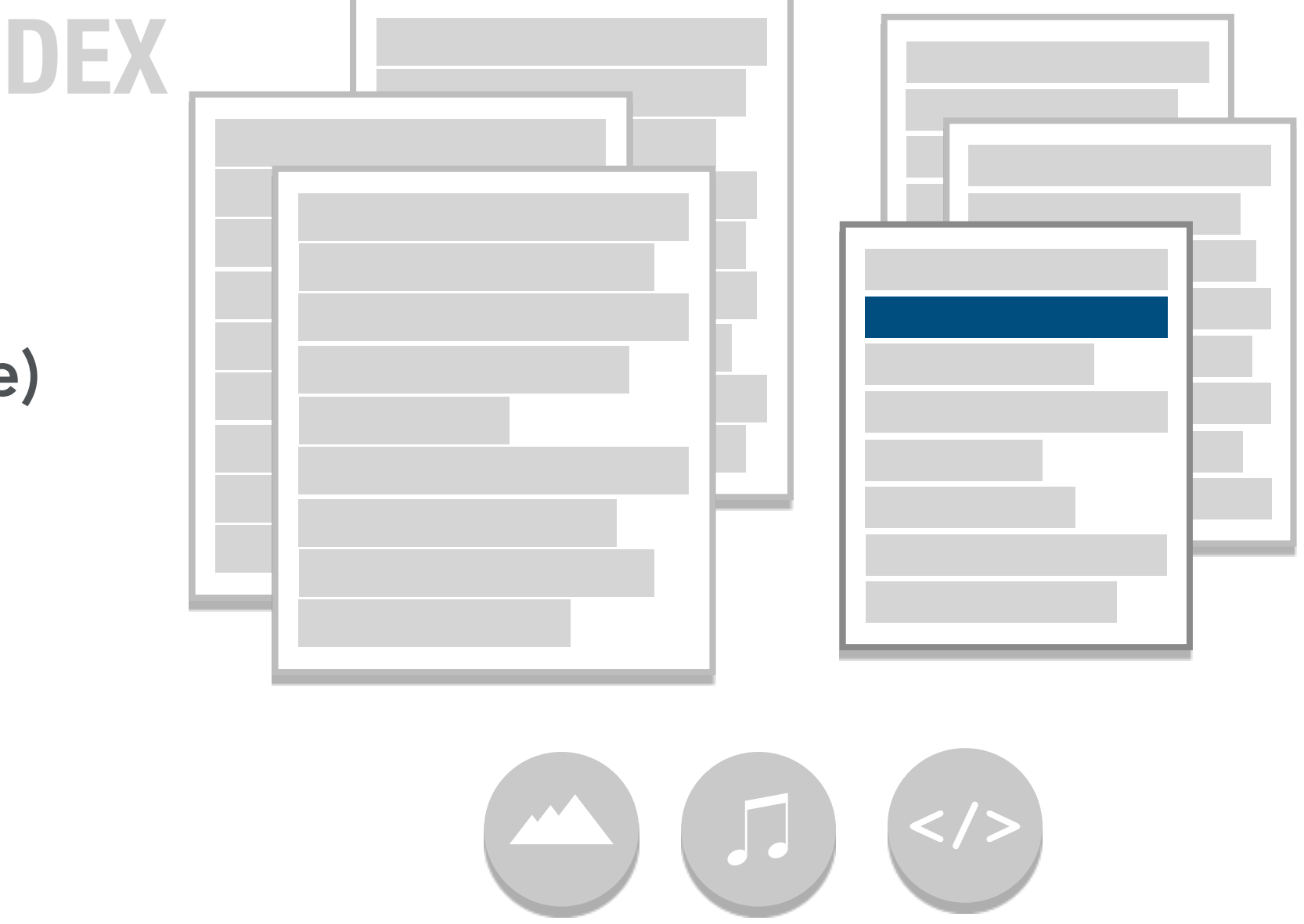


Identify activity in dex

Organ Harvesting



2 Choose any vein (backward slice)

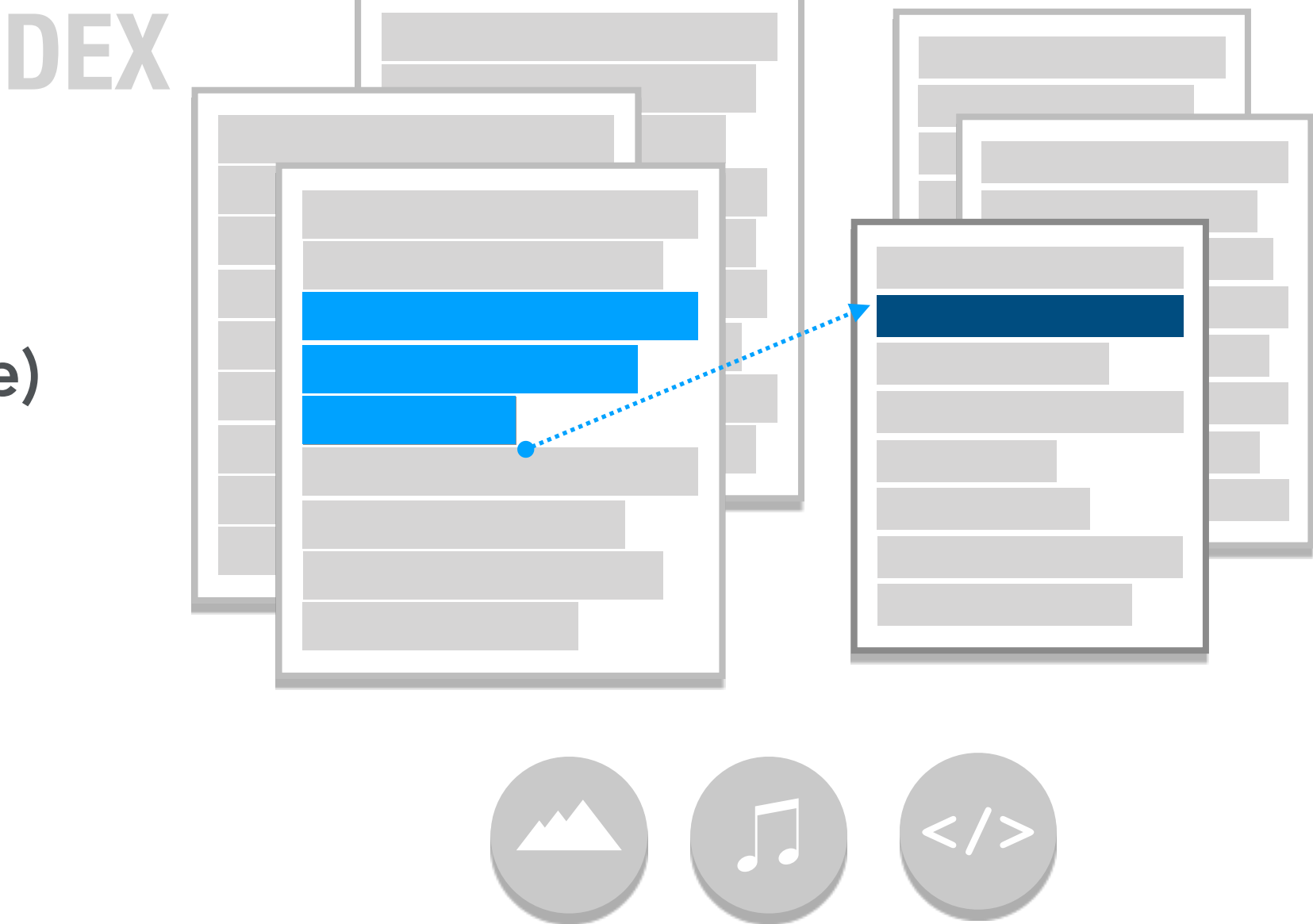


Extract intent creation and startActivity()

Organ Harvesting

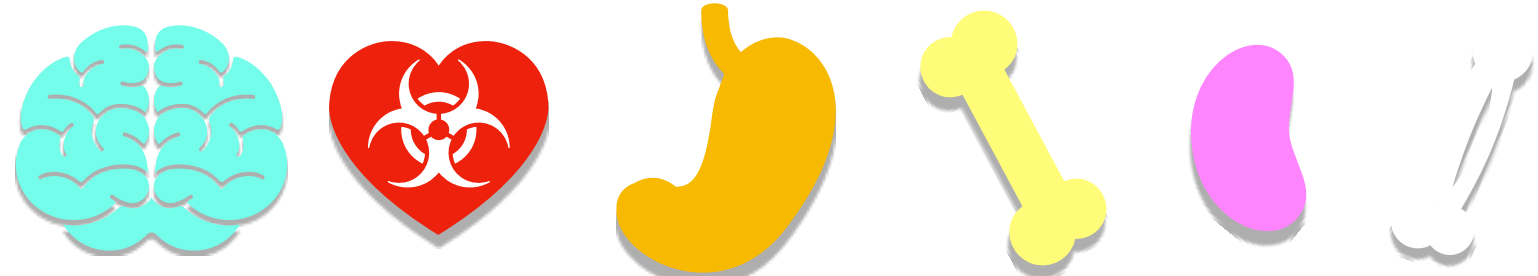


2 Choose any vein (backward slice)

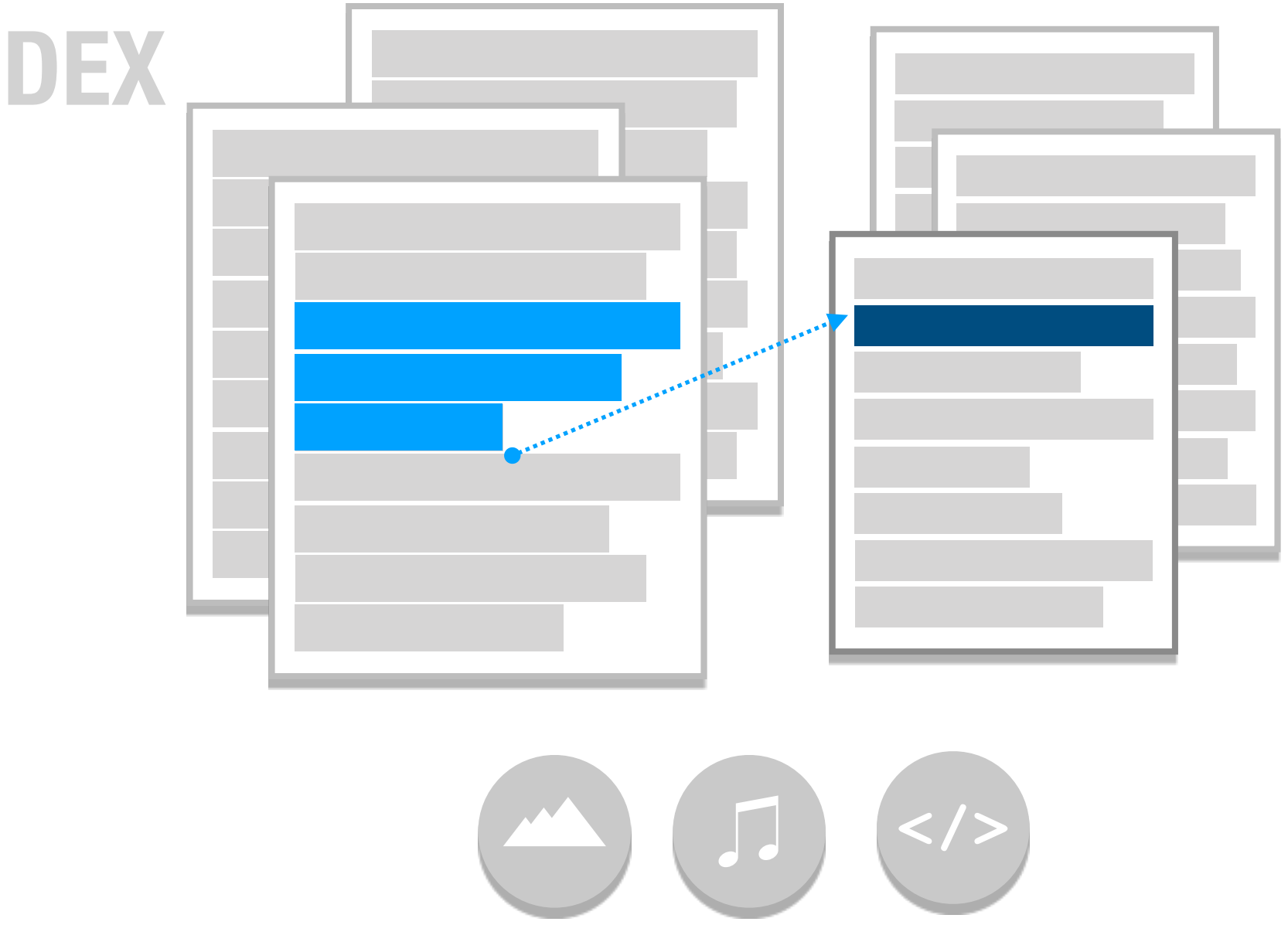


Extract intent creation and startActivity()

Organ Harvesting

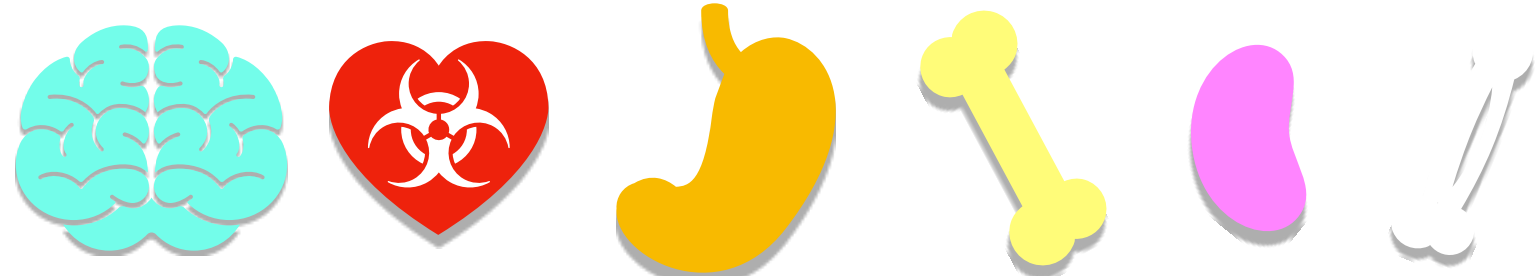


3 Collect organ (forward slice)

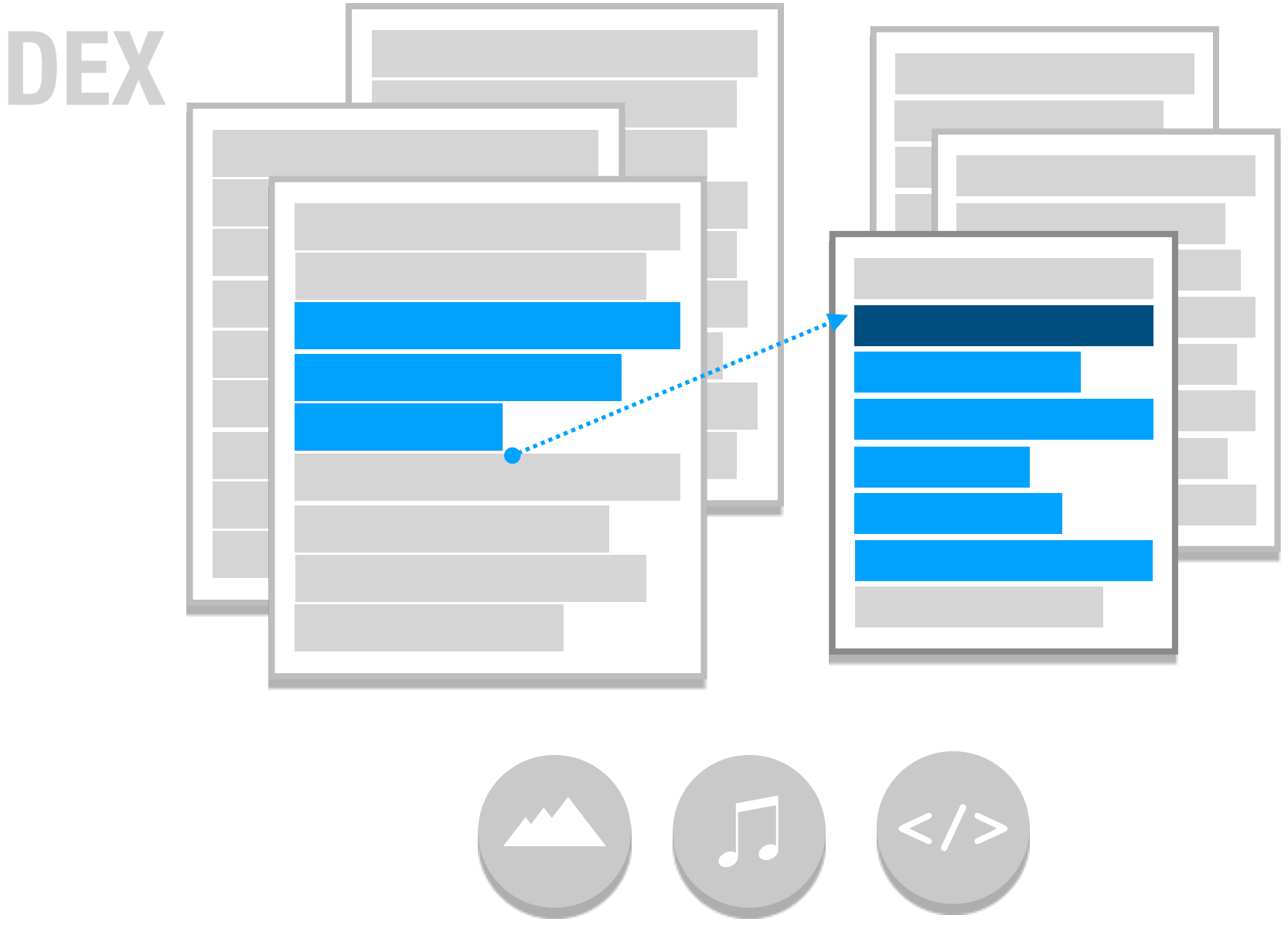


Gather activity definition

Organ Harvesting

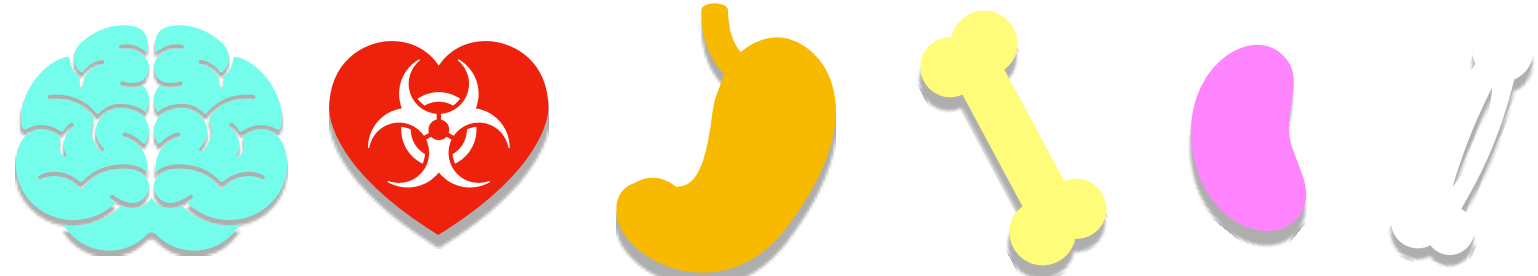


3 Collect organ (forward slice)

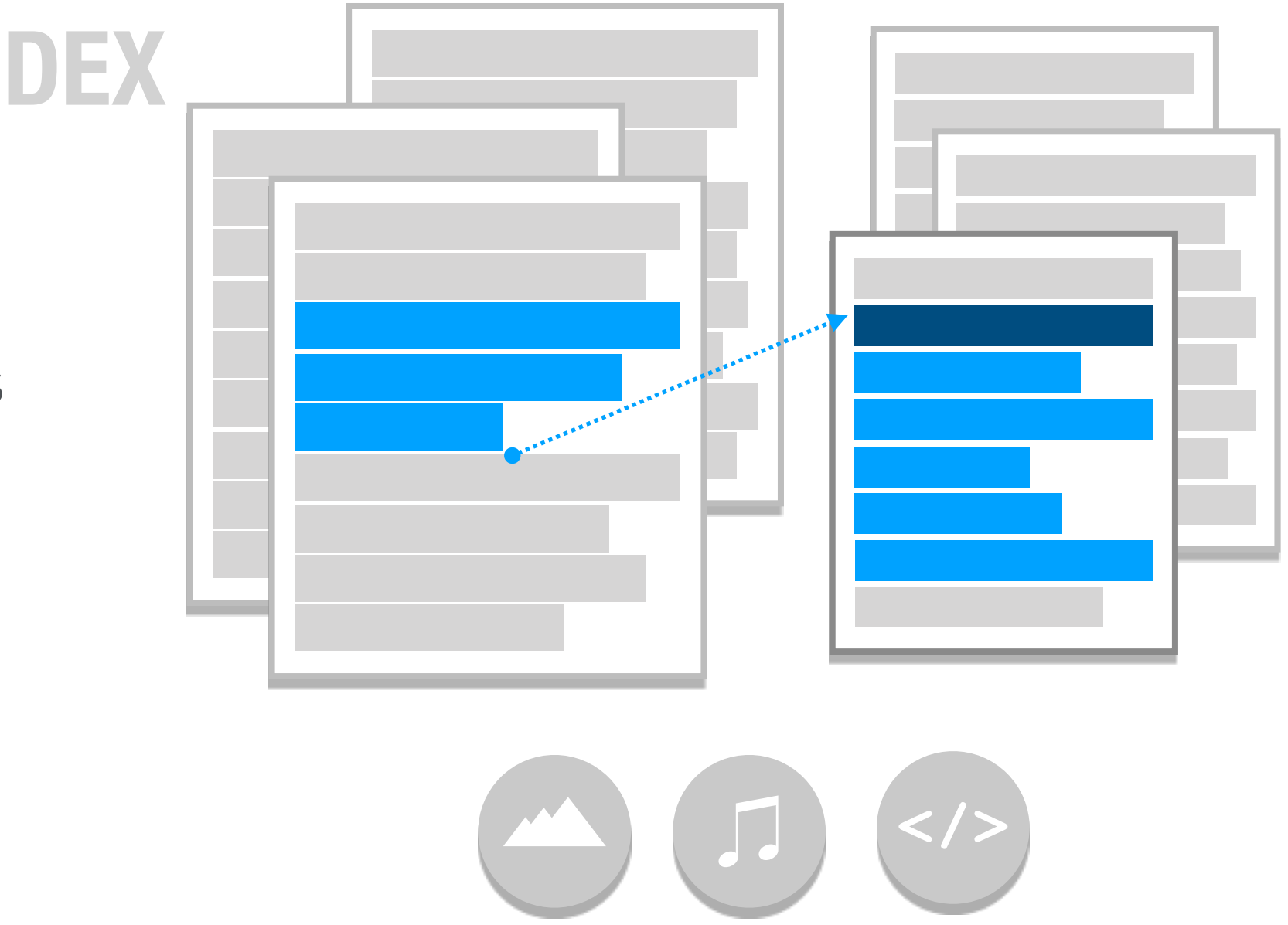


Gather activity definition

Organ Harvesting

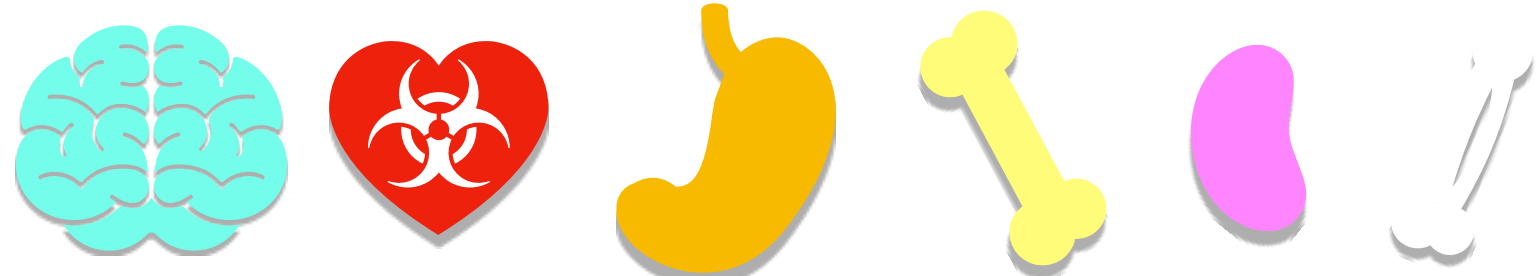


4 Include transitive dependencies

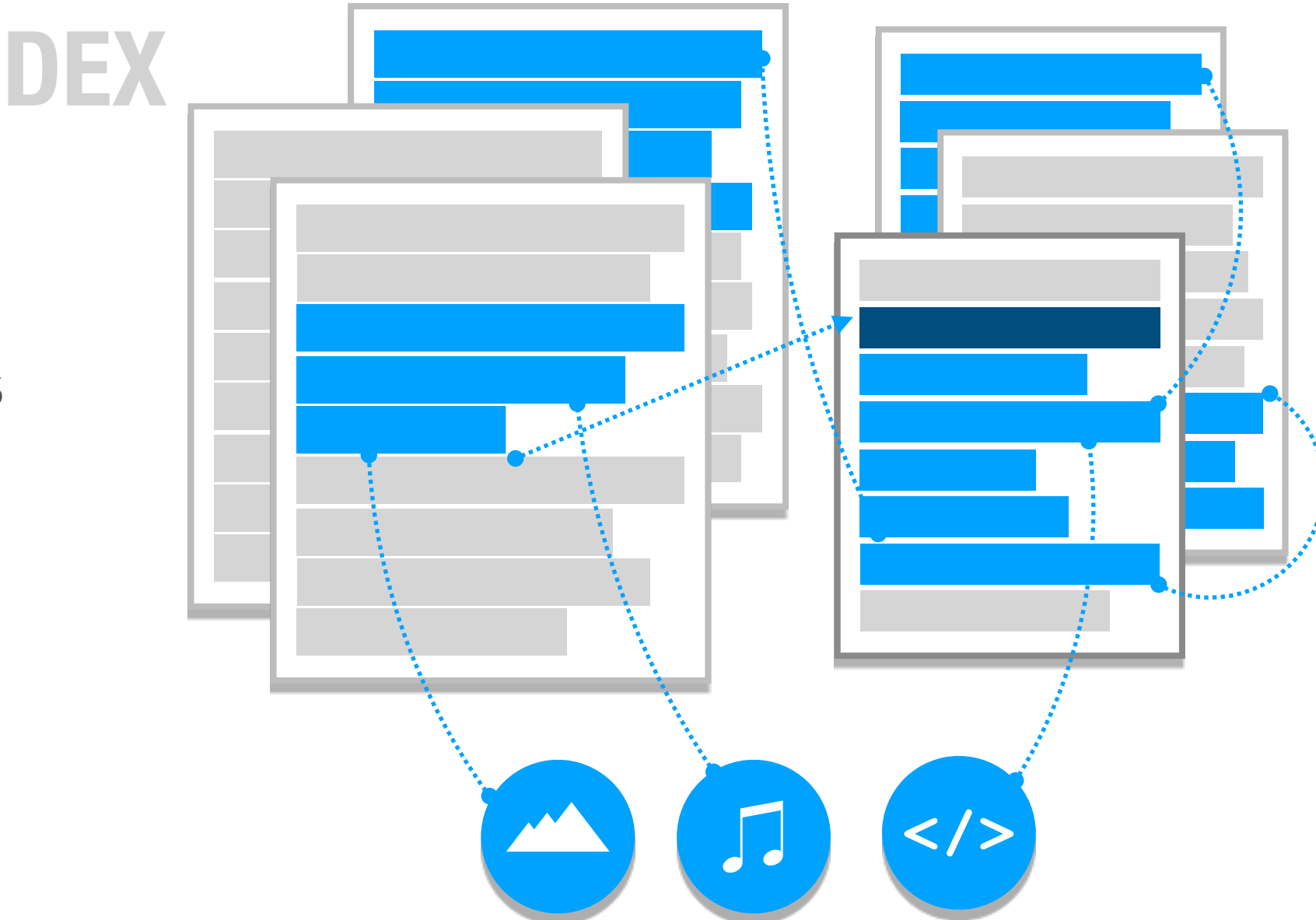


Recursively collect dependencies

Organ Harvesting



4 Include transitive dependencies

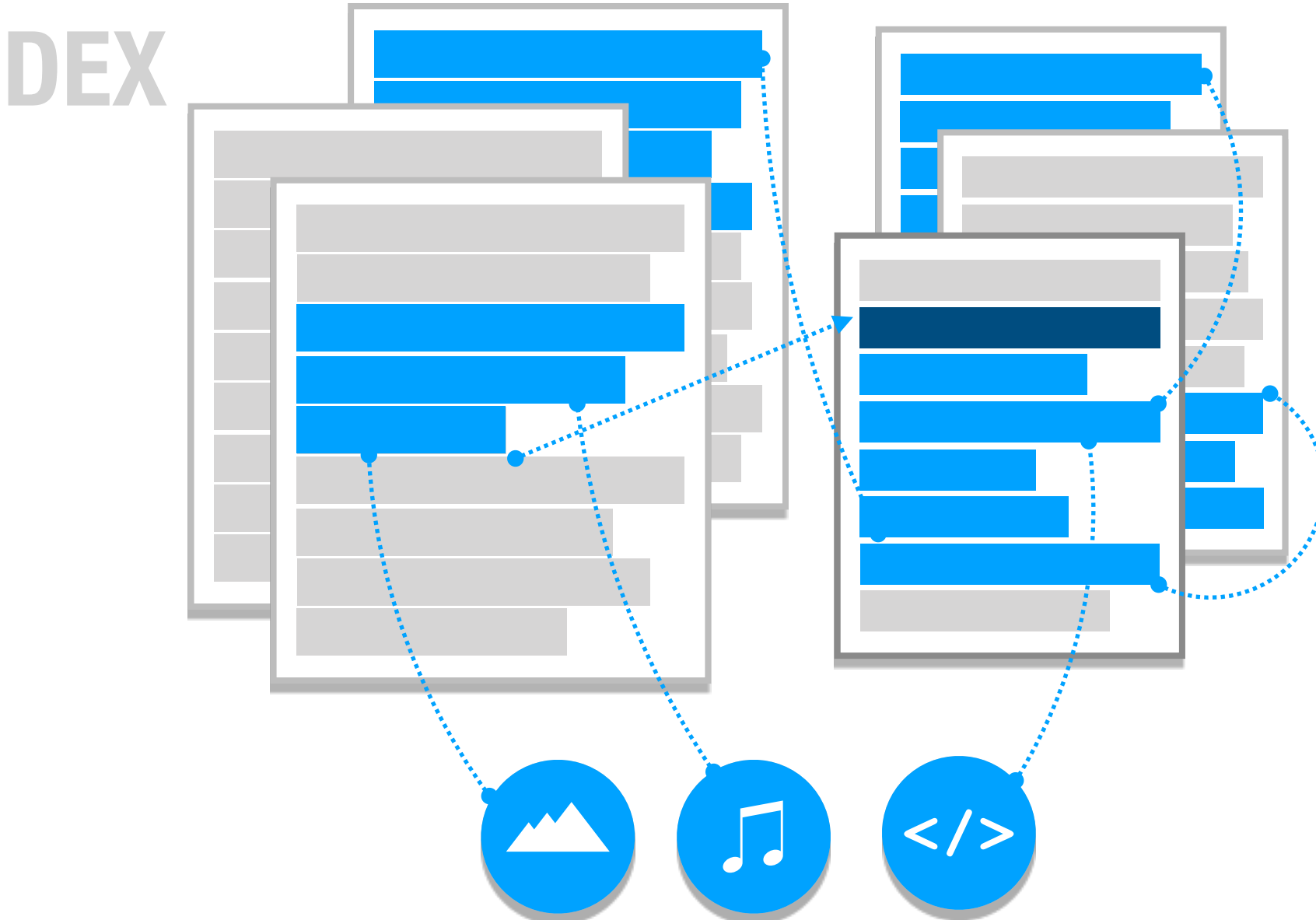


Recursively collect dependencies

Organ Harvesting



5 Store organ in an "ice box"

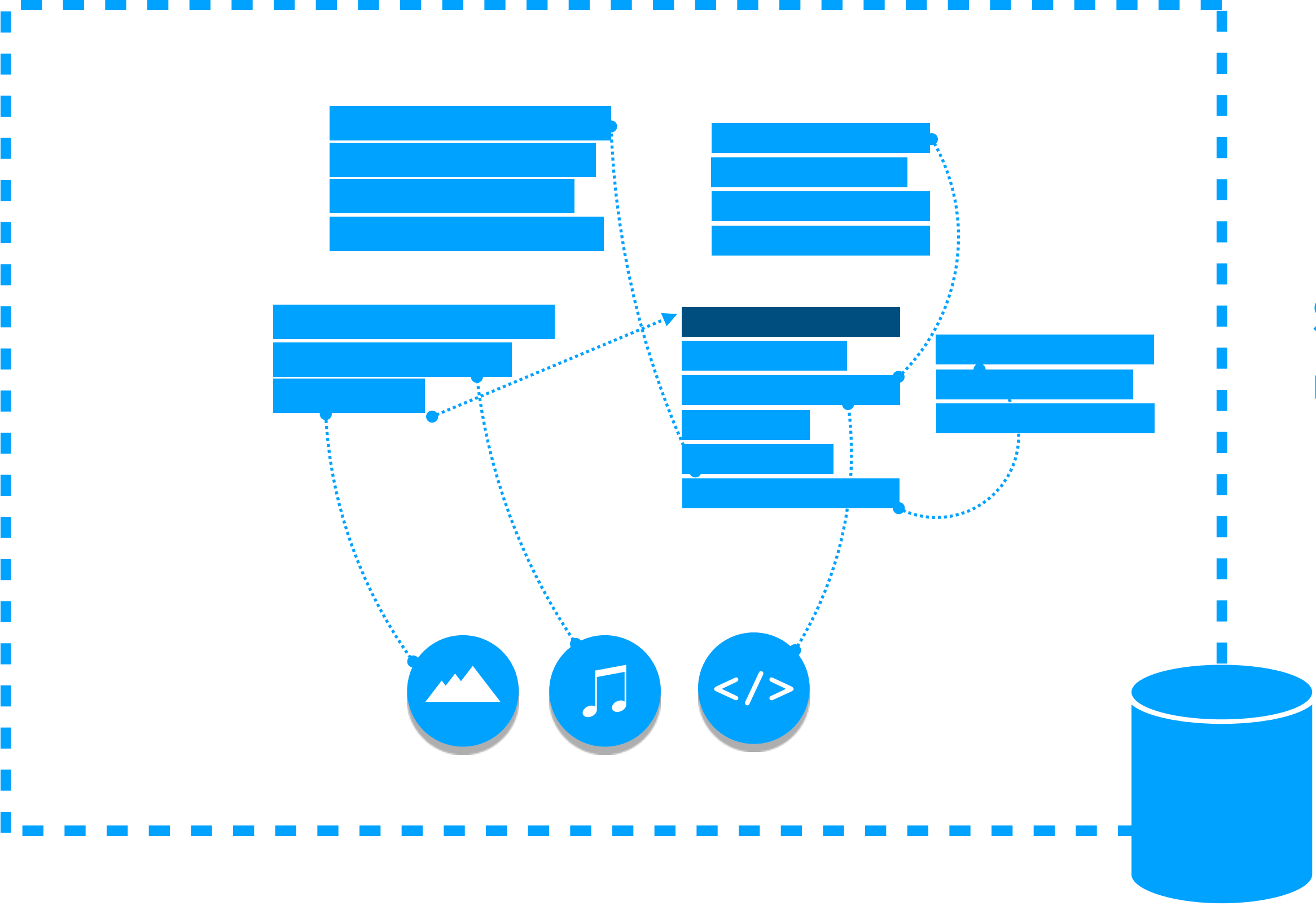


Save gadget to a database ready for the attack

Organ Harvesting

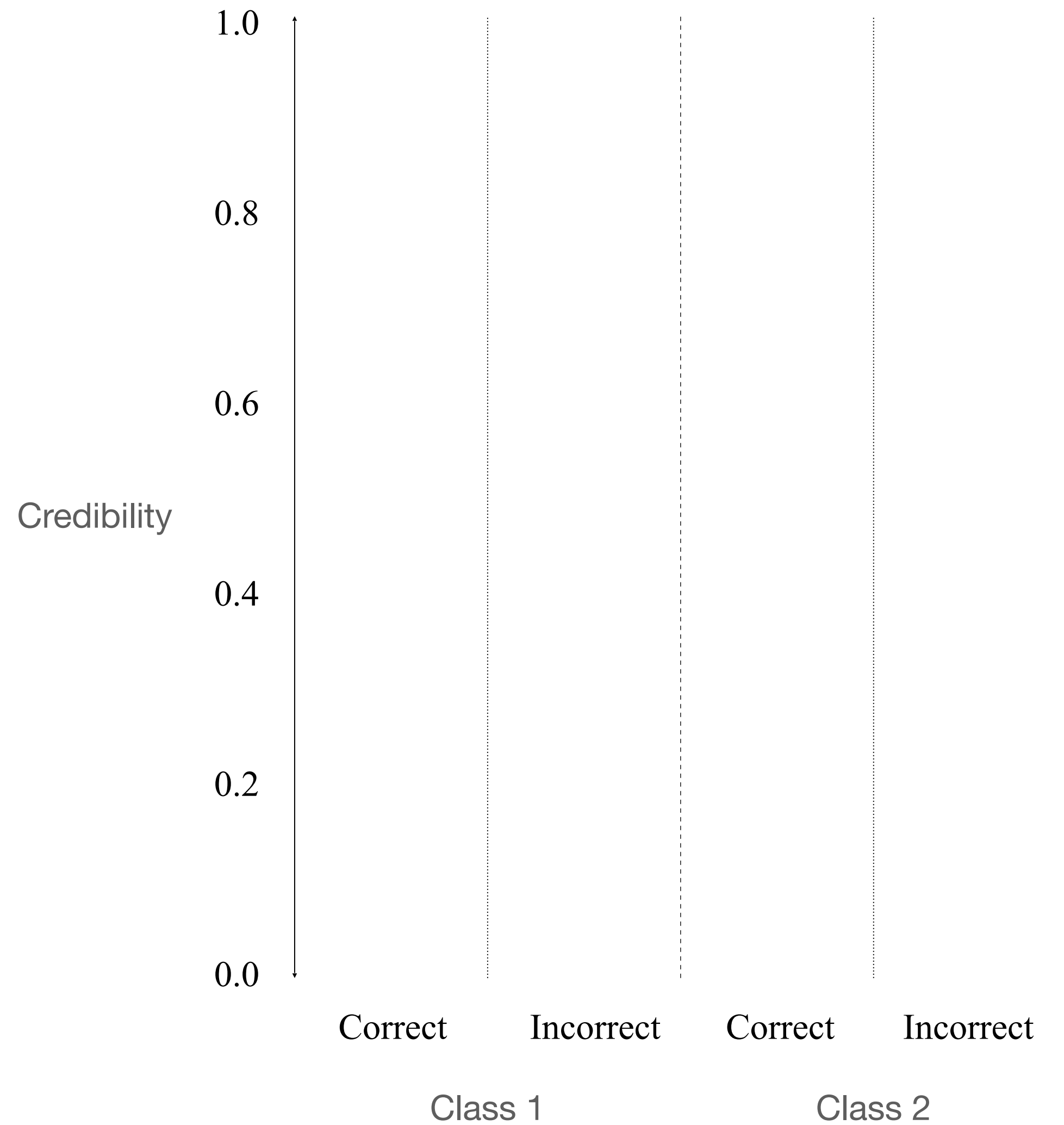


5 Store organ in an "ice box"



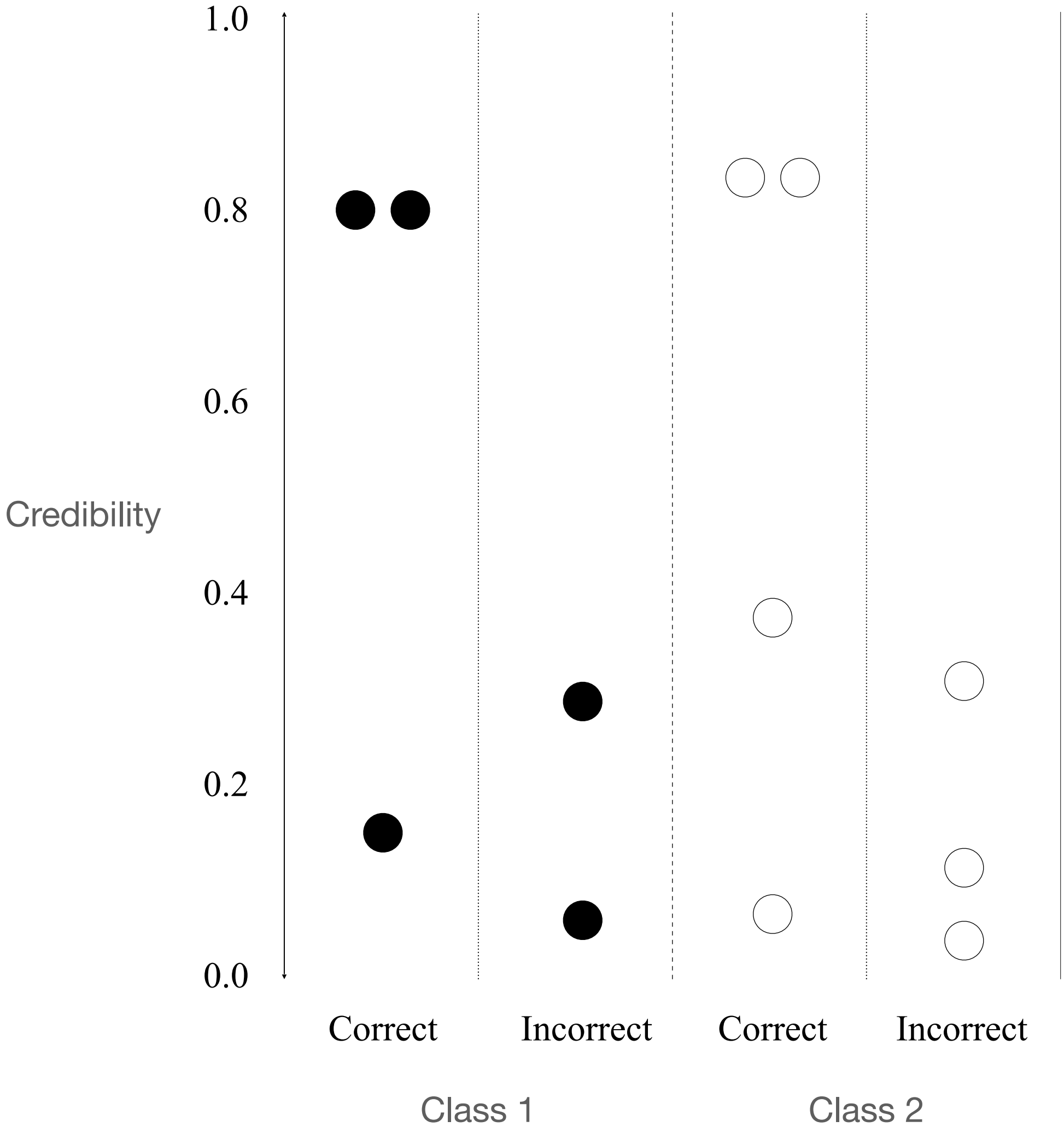
Save gadget to a database ready for the attack

Transcend Calibration



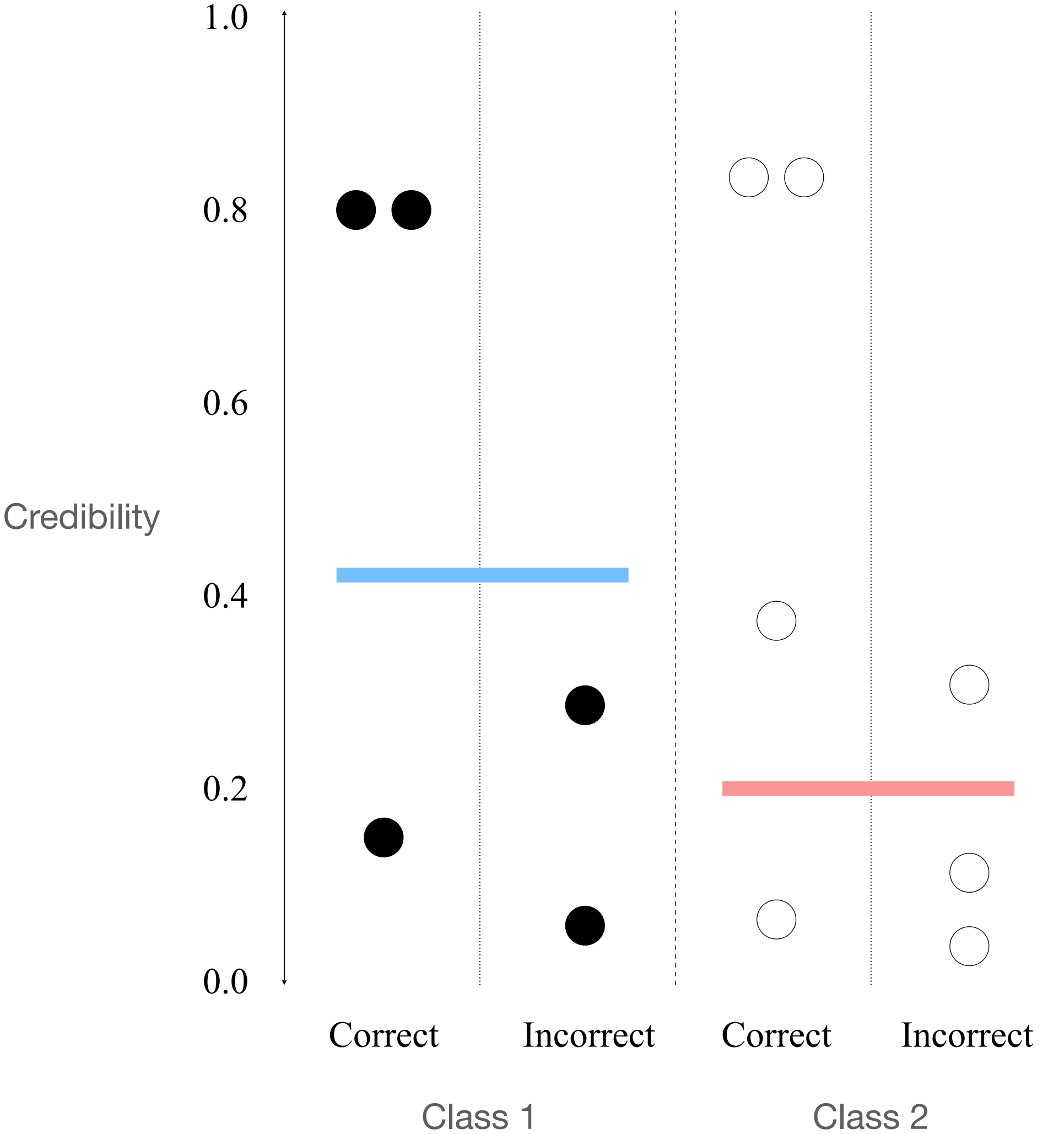
- **How much drift is too much?**
- Produce a threshold for each class
- Optimize cost vs performance on training and calibration sets
- Maximise separation between credibility of correct and incorrect decisions

Transcend Calibration



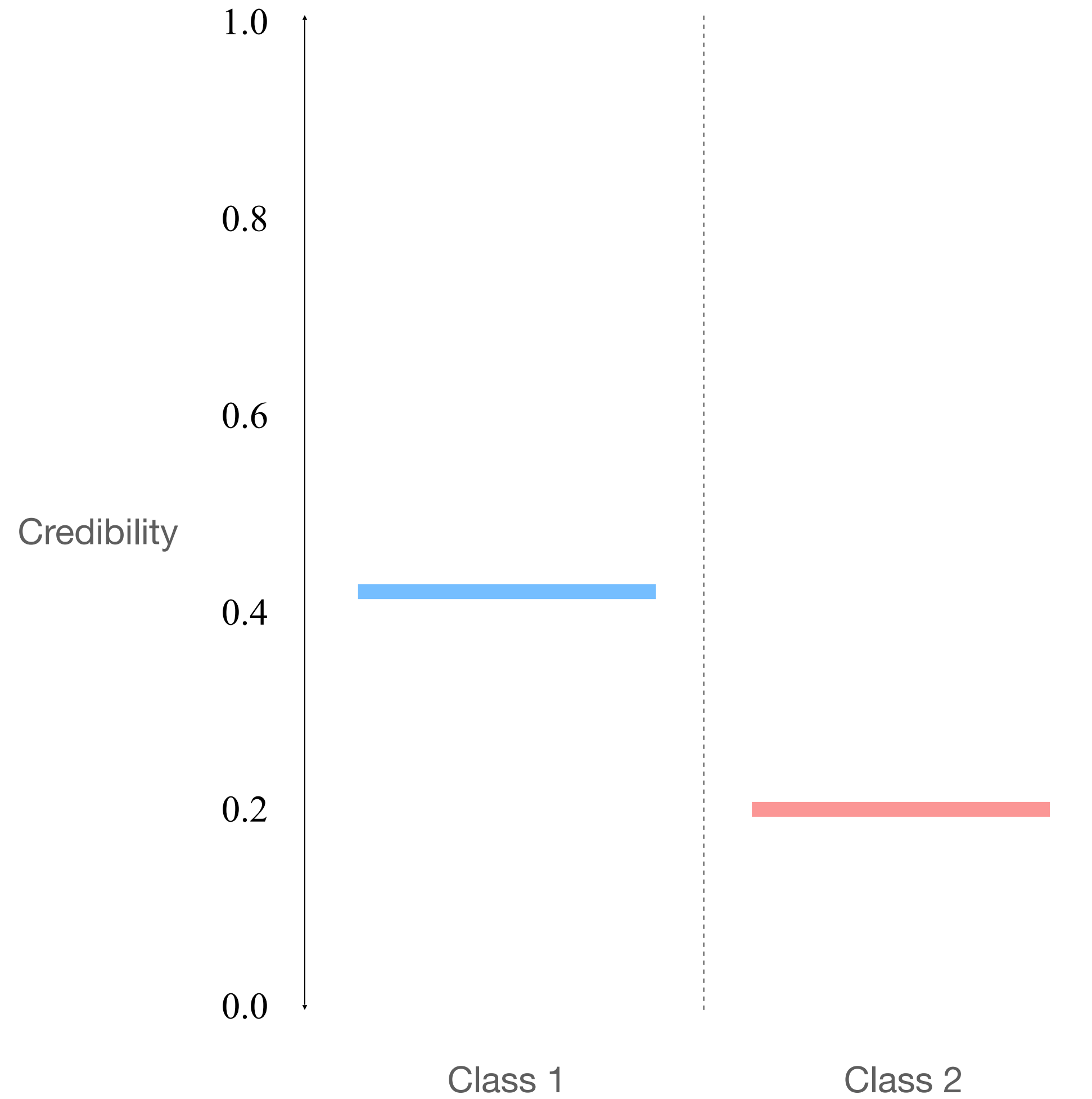
- **How much drift is too much?**
- Produce a threshold for each class
- Optimize cost vs performance on training and calibration sets
- Maximise separation between credibility of correct and incorrect decisions

Transcend Calibration



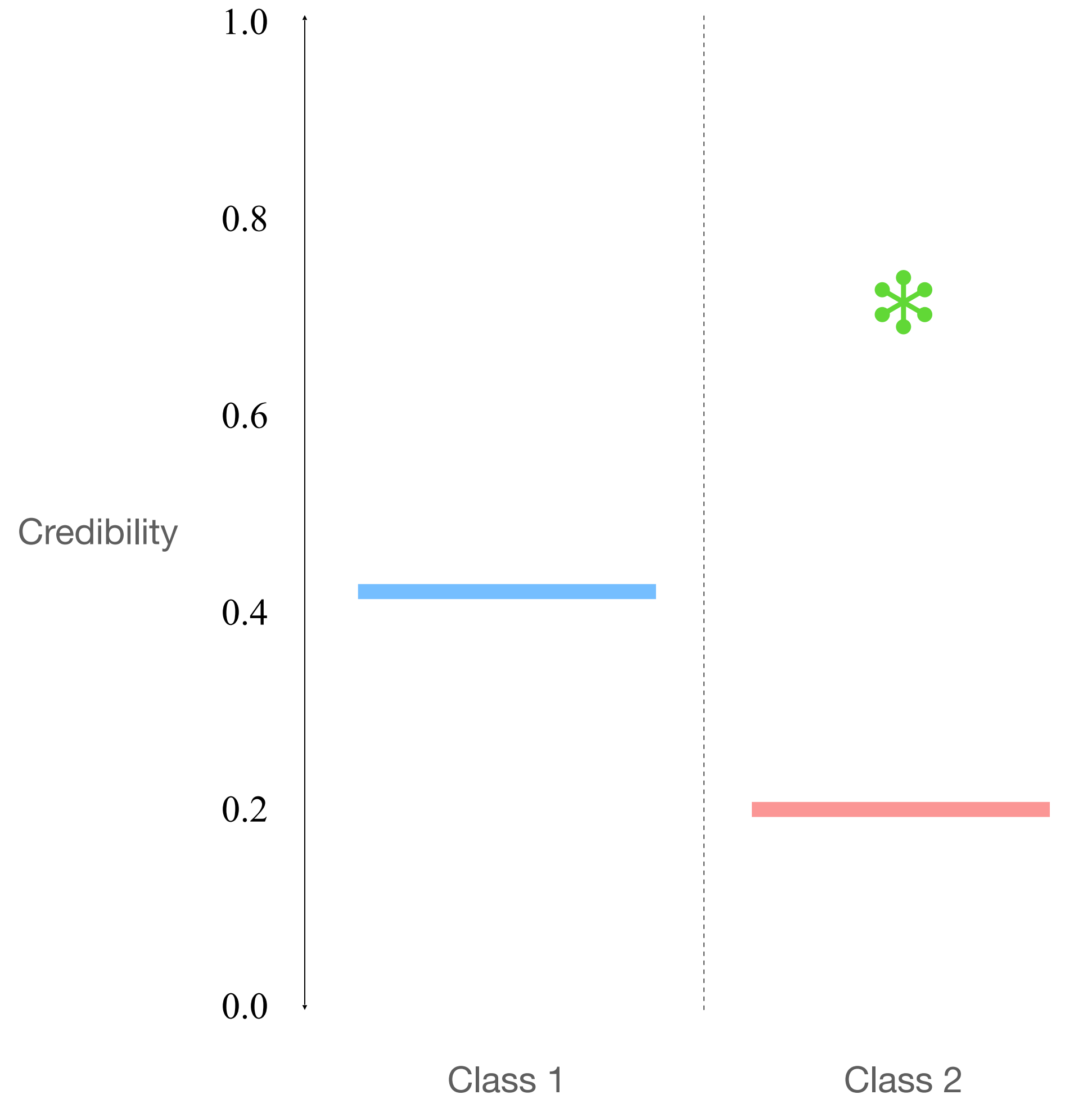
- **How much drift is too much?**
- Produce a threshold for each class
- Optimize cost vs performance on training and calibration sets
- Maximise separation between credibility of correct and incorrect decisions

Transcend at Test Time



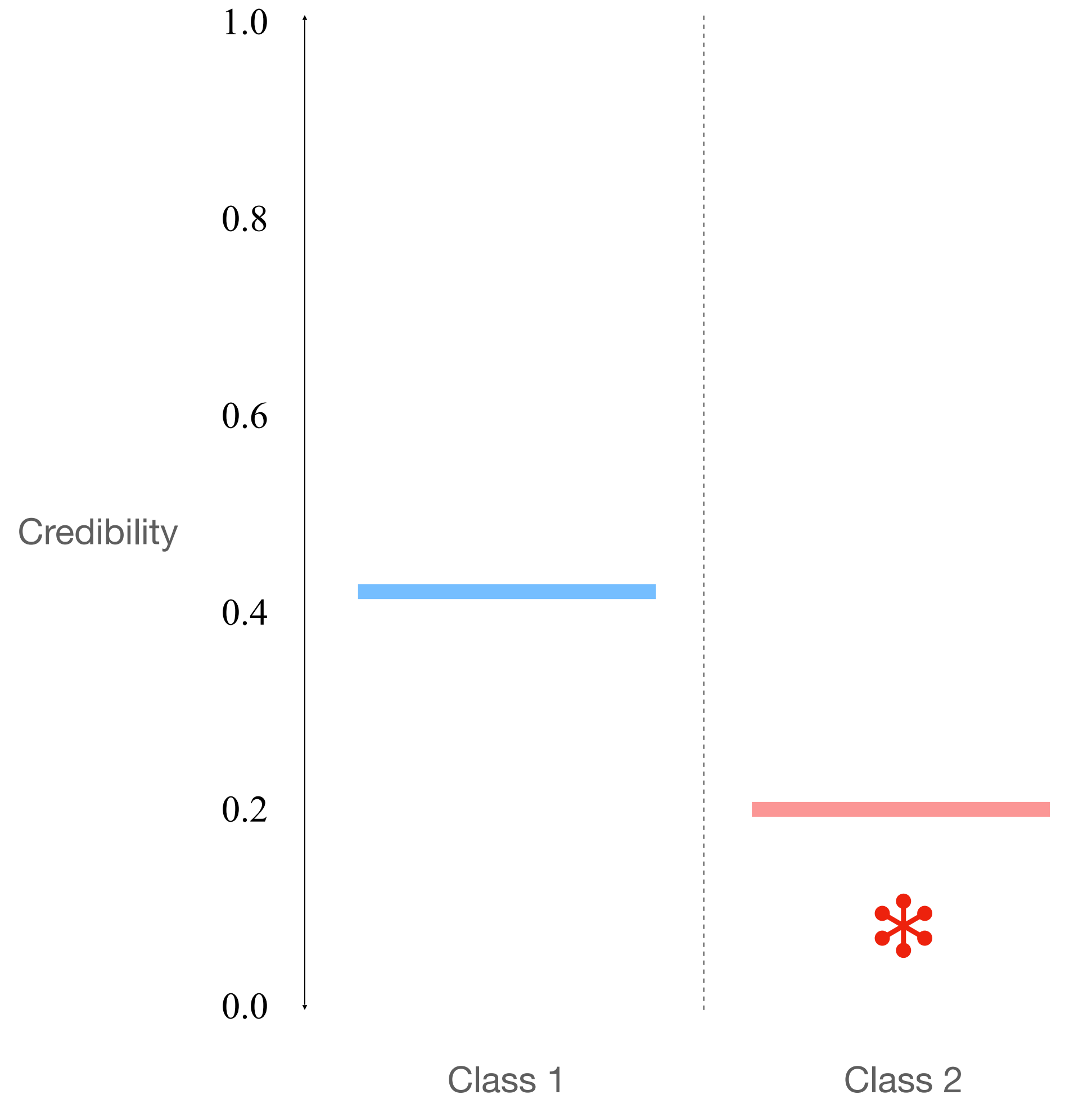
- Credibilities of new examples are compared against the threshold of their predicted class
- Above = keep the prediction
- Below = reject the prediction

Transcend at Test Time



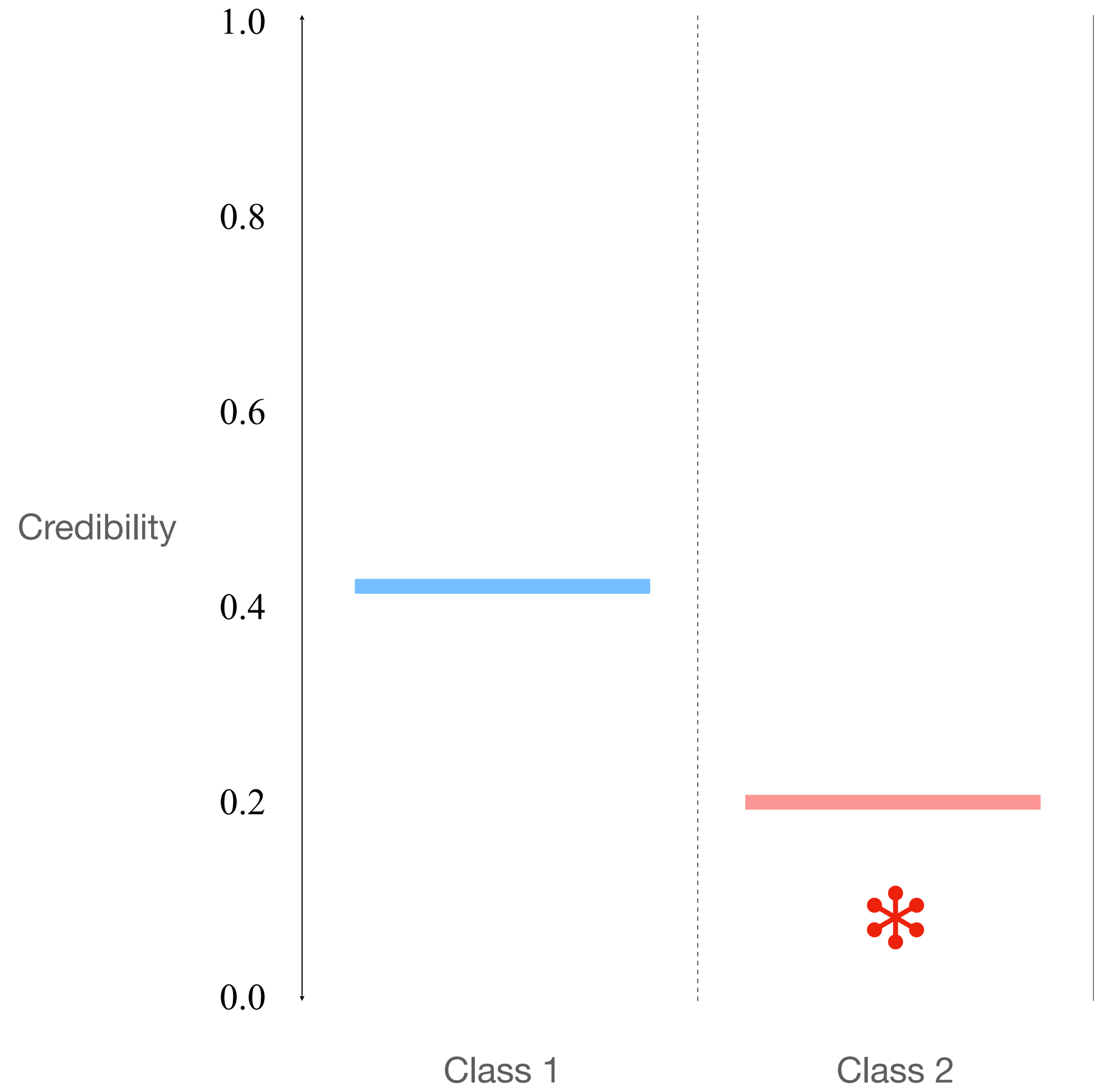
- Credibilities of new examples are compared against the threshold of their predicted class
- Above = keep the prediction
- Below = reject the prediction

Transcend at Test Time



- Credibilities of new examples are compared against the threshold of their predicted class
- Above = keep the prediction
- Below = reject the prediction

Rejection Cost



Rejection Cost

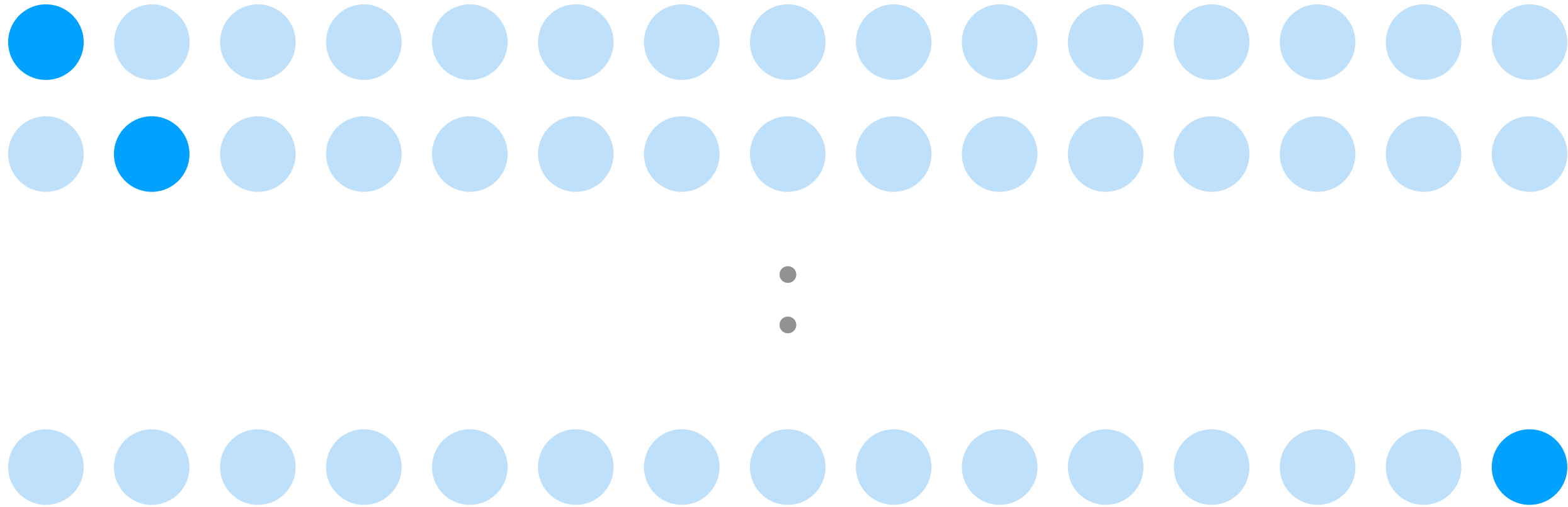


- Actions for rejected points *:
 - Manual inspection
 - Downstream analysis
 - Quarantine
 - Exemption

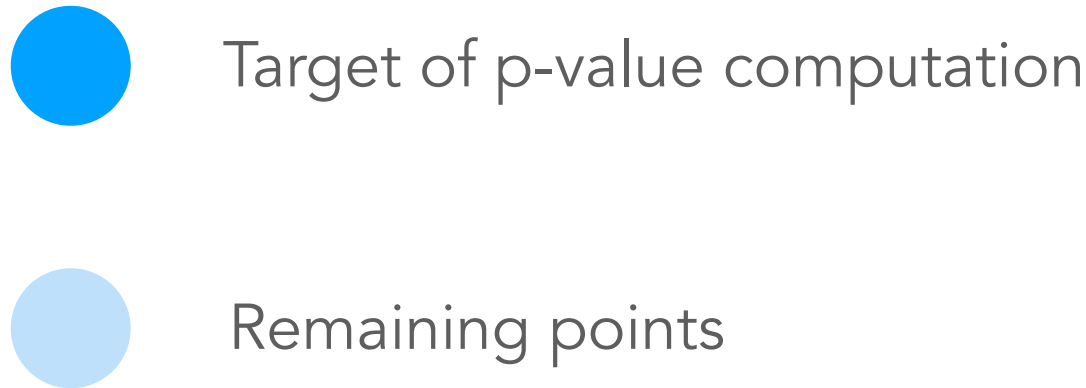
* [AISec 2021] [Investigating Labelless Drift Adaptation for Malware Detection](#)

* [AISec 2021] [INSOMNIA: Towards Concept-Drift Robustness in Network Intrusion Detection](#)

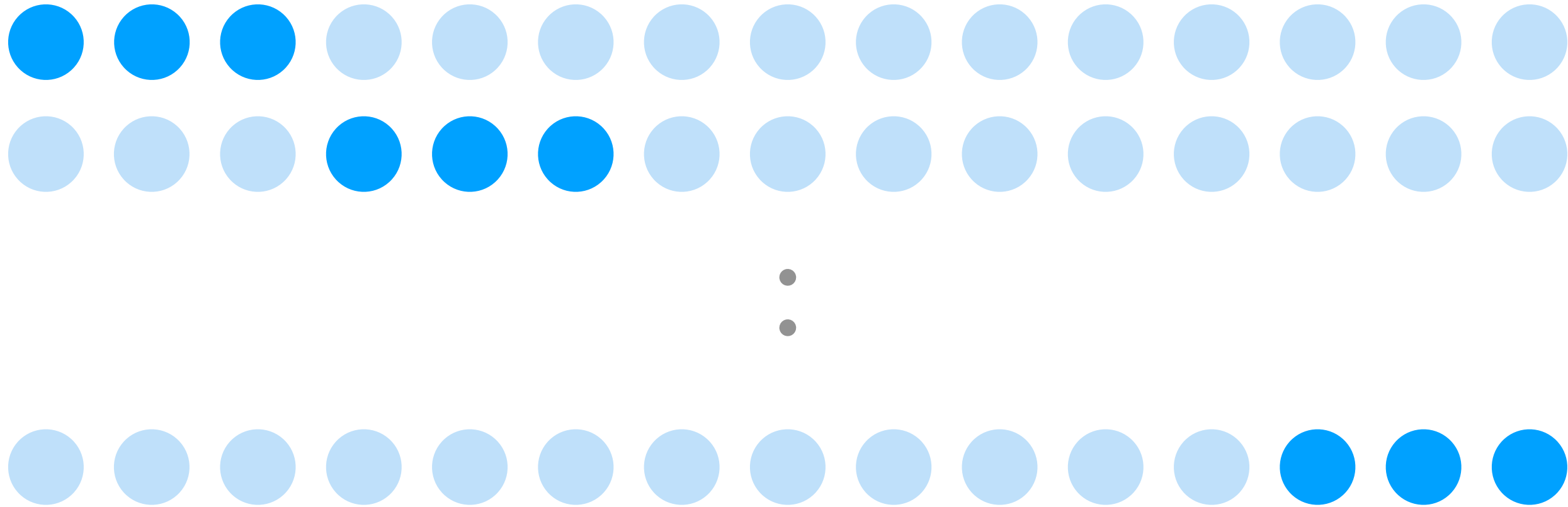
The Cost of Transductive Conformal Evaluators




- Underlying classifier retrained for every training point
- Rooted in CP theory
- Often computationally infeasible




Approximate TCE

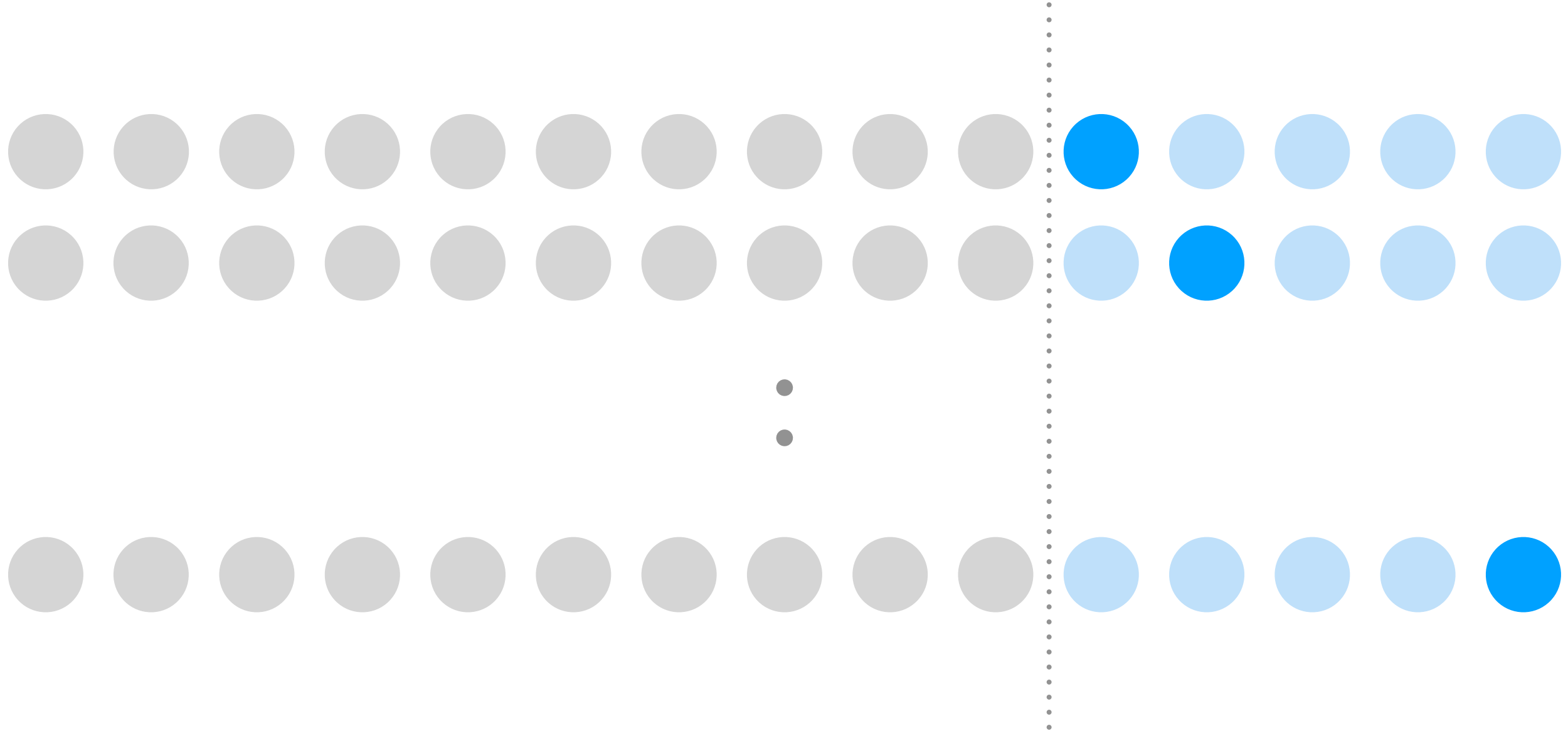





- First attempt to improve on the TCE
- P-values computed in batches
- Relies on unsound assumption

 Target of p-value computation

 Remaining points

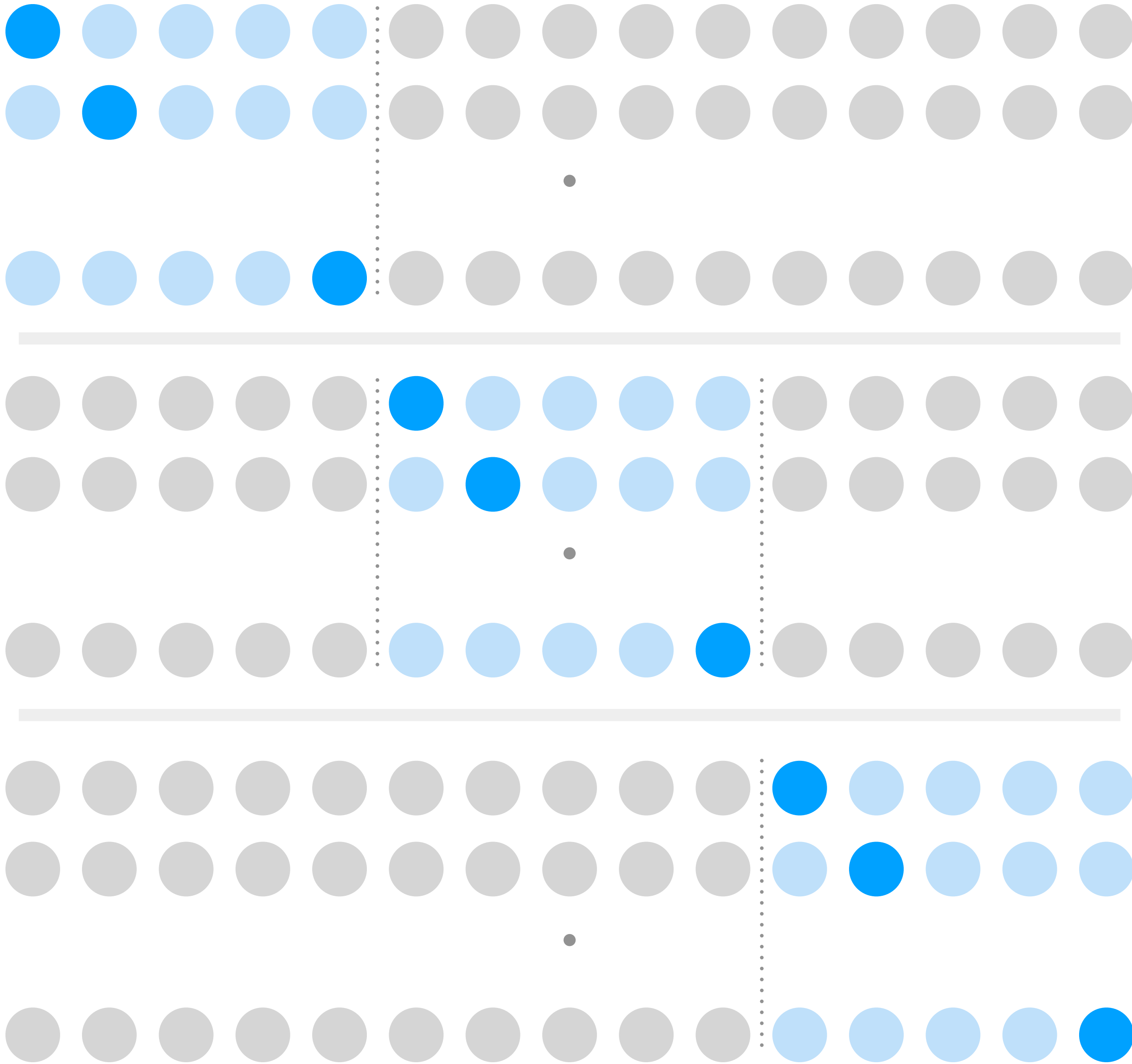
Inductive Conformal Evaluator (ICE)



-  Target of p-value computation
-  Remaining points
-  Excluded points used for prediction but not evaluation

- Increase speed by splitting into training and calibration sets
- Rooted in CP theory
- Computationally efficient
- Informationally inefficient

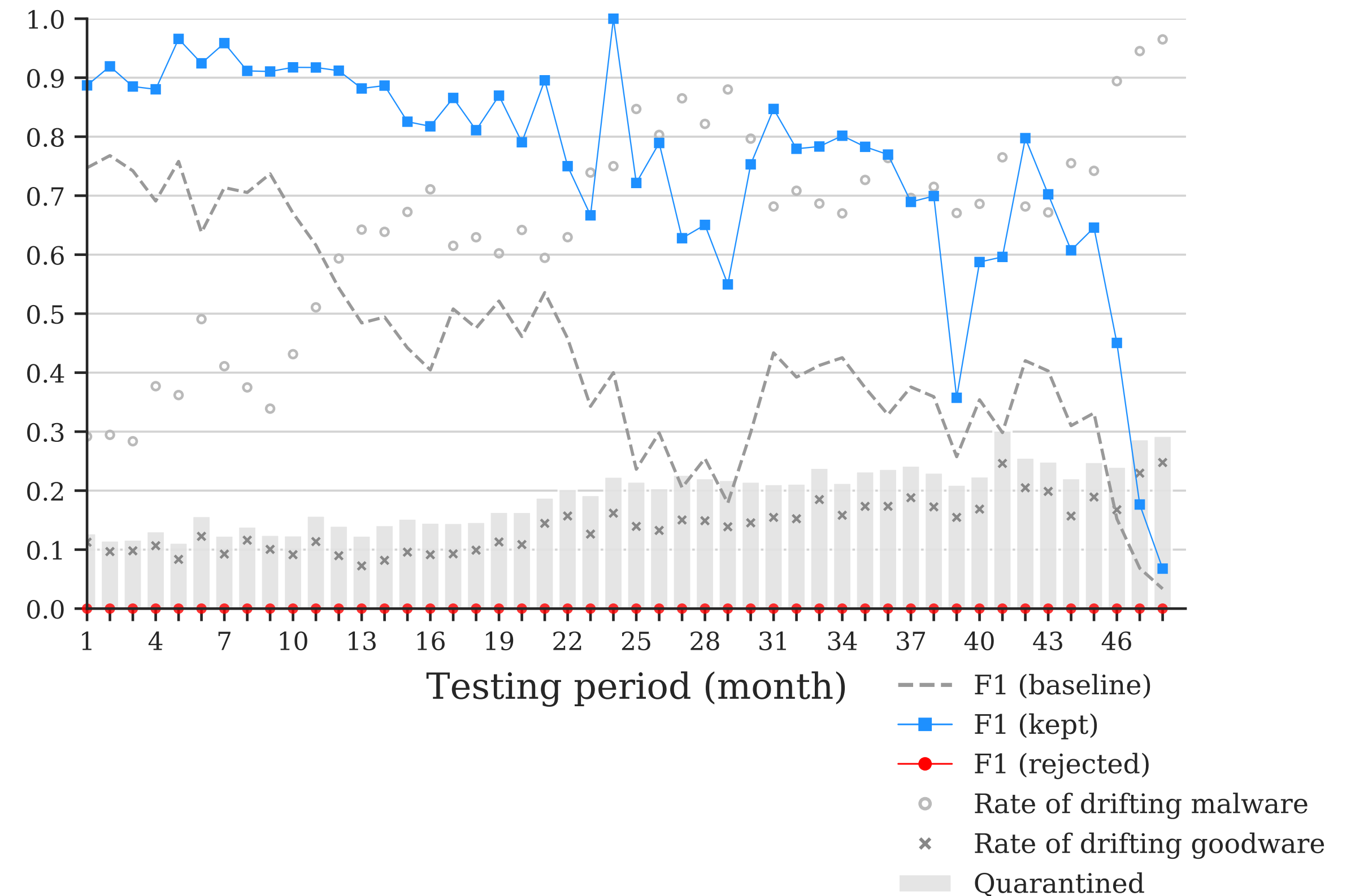
Cross-Conformal Evaluator (CCE)



- Inspired by cross validation - multiple ICEs in parallel vote on evaluation
- Rooted in CP theory
- Computationally efficient
- Informationally efficient

Results: Rejection Performance — Drift Rate

Android Malware (maximizing F1)

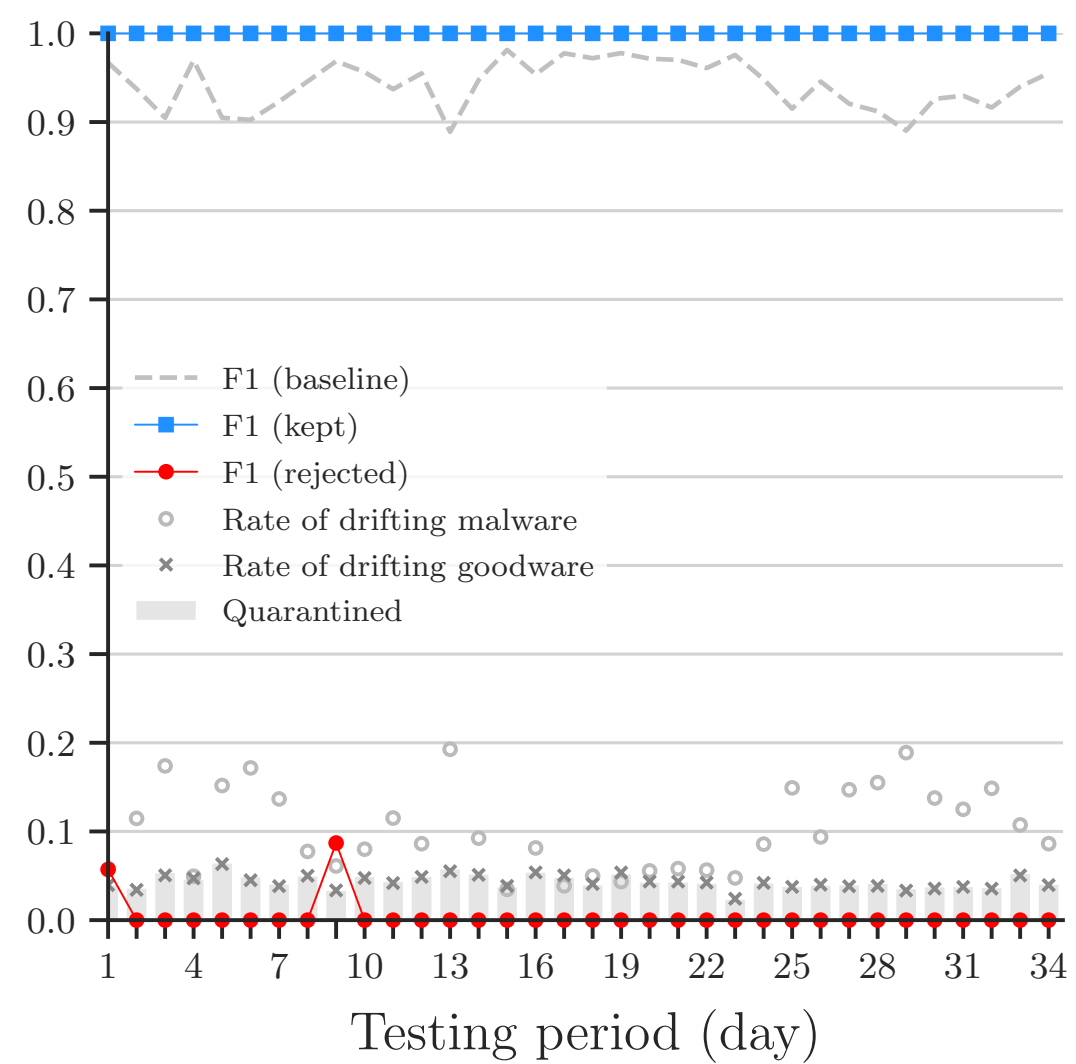


[IEEE S&P 2022] **Transcending TRANSCEND: Revisiting Malware Classification in the Presence of Concept Drift**

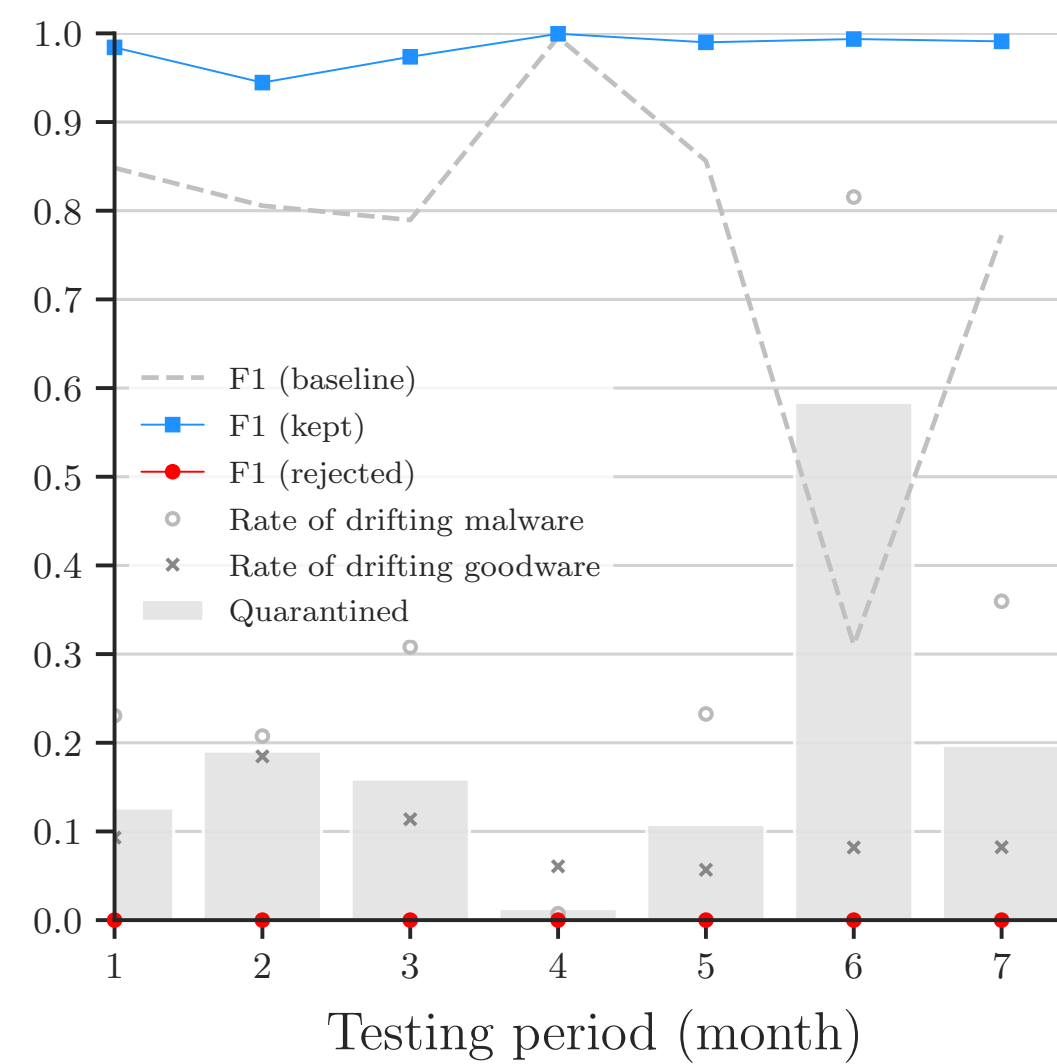
<https://s2lab.cs.ucl.ac.uk/projects/transcend/>

Results: Rejection Performance — Drift Rate

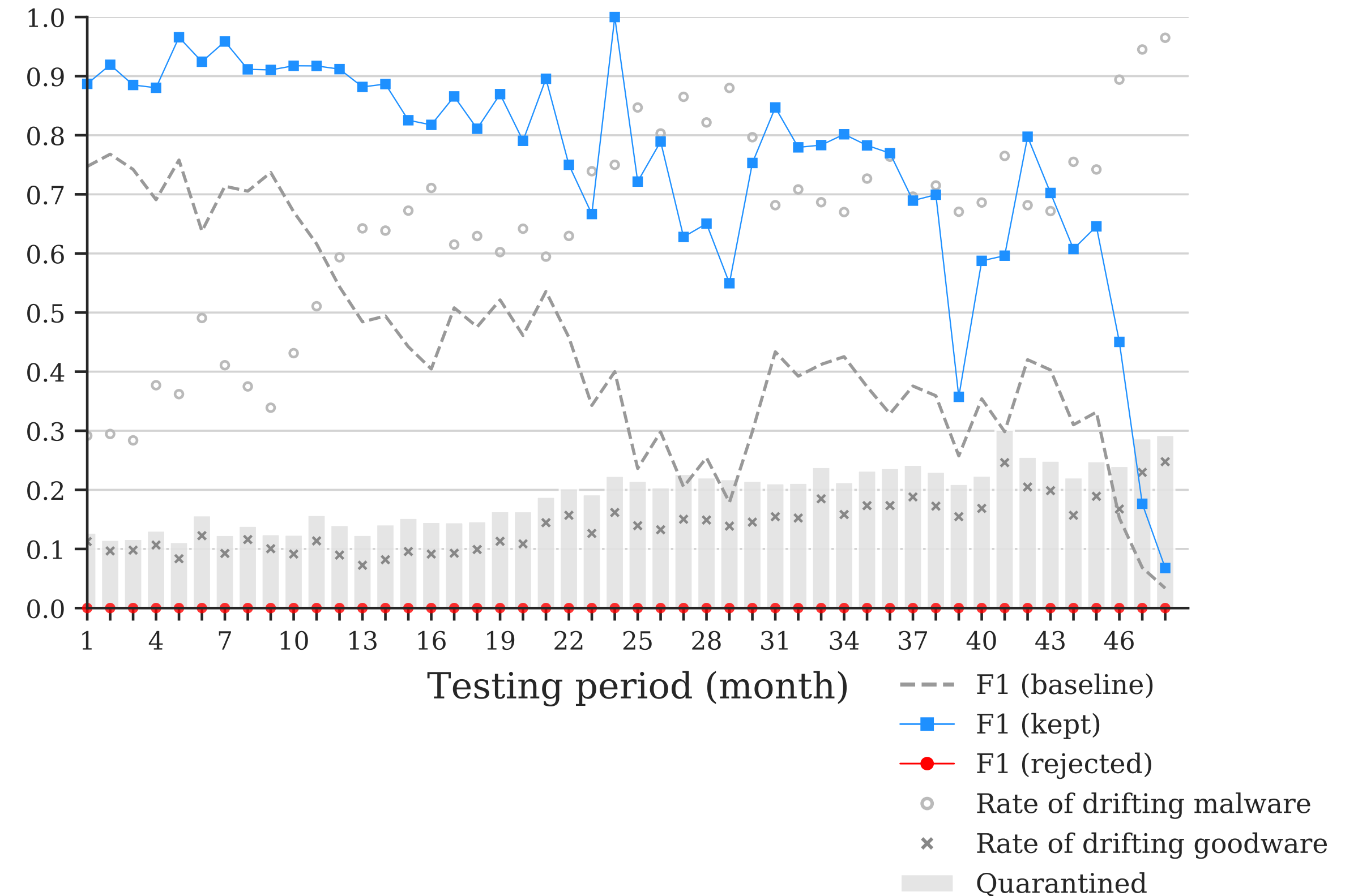
PDF Malware



Windows PE Malware



Android Malware (maximizing F1)

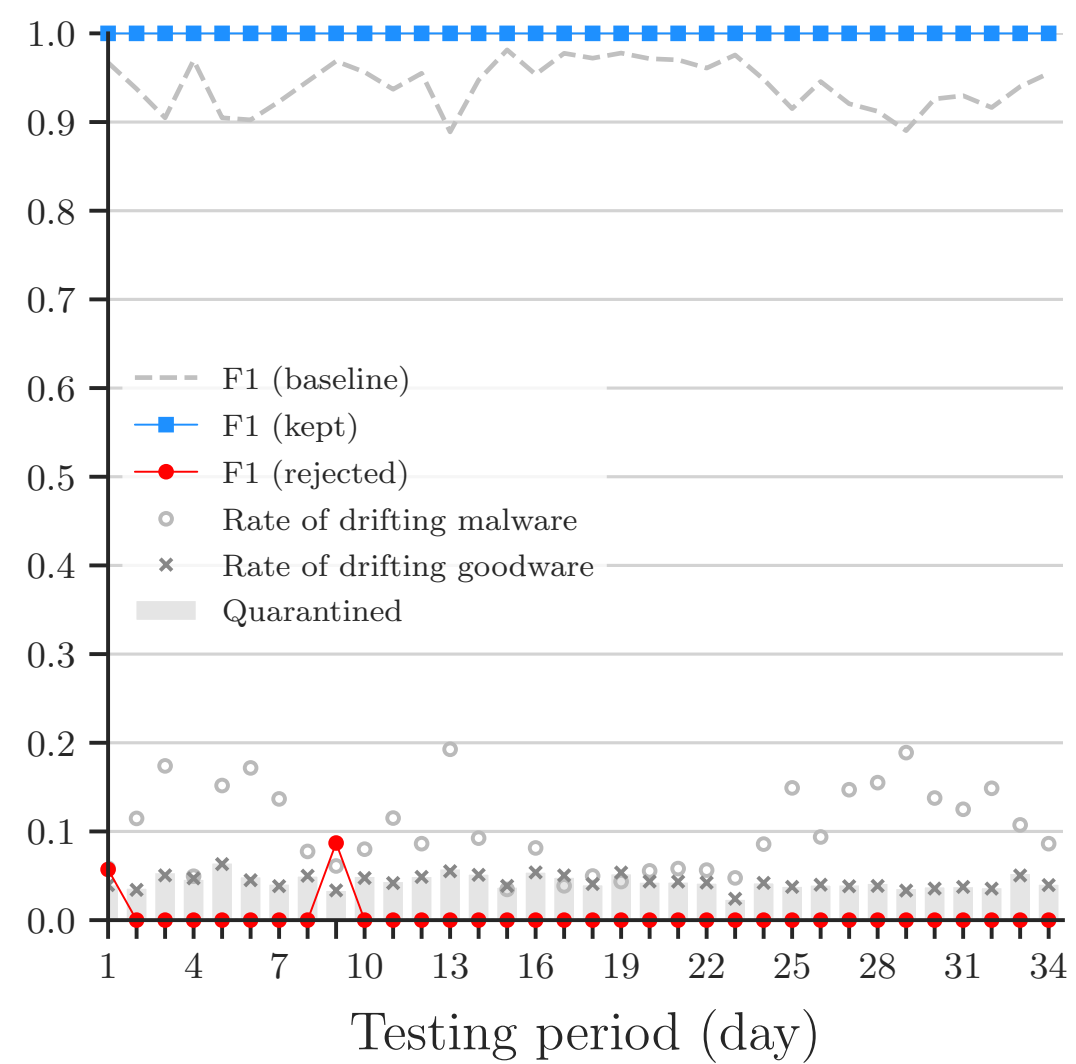


[IEEE S&P 2022] Transcending TRANSCEND: Revisiting Malware Classification in the Presence of Concept Drift

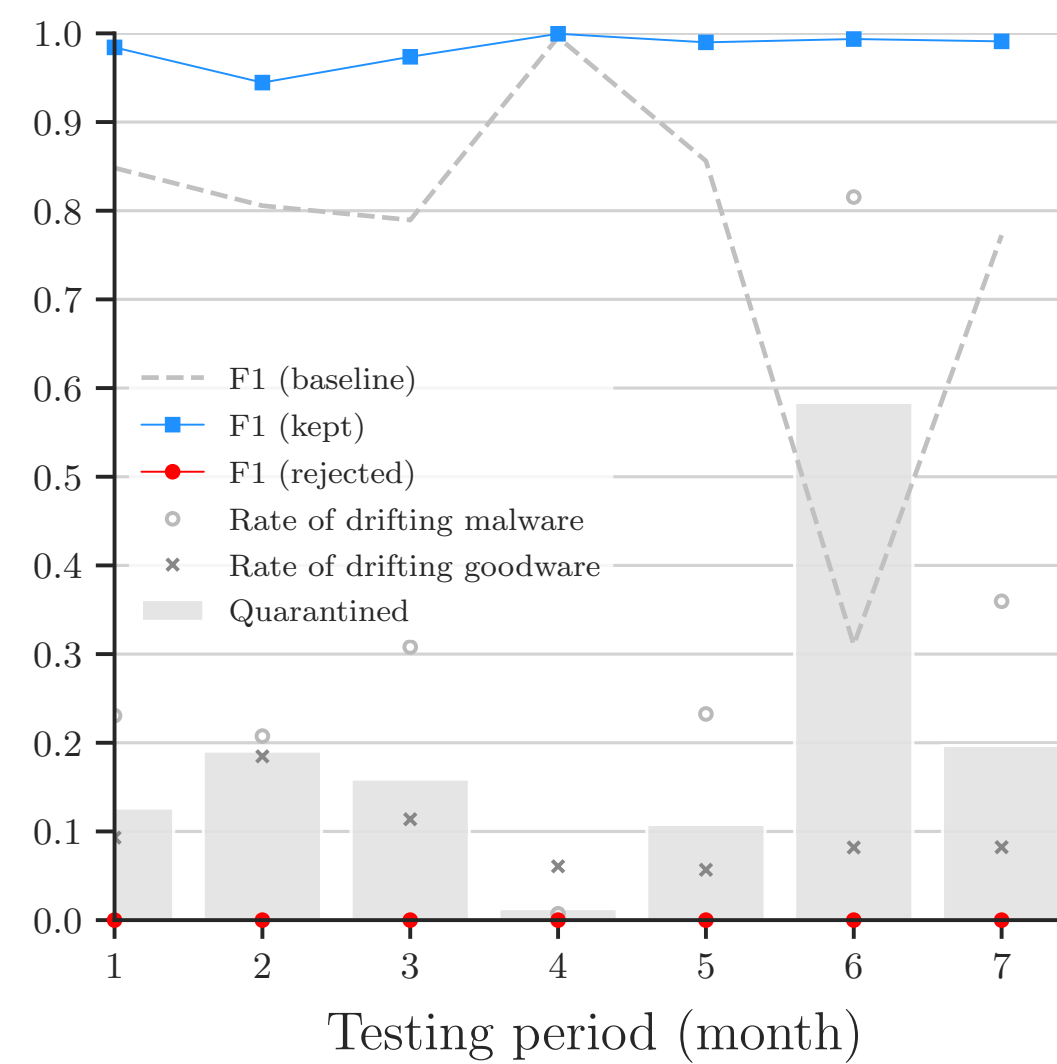
<https://s2lab.cs.ucl.ac.uk/projects/transcend/>

Results: Rejection Performance — Drift Rate

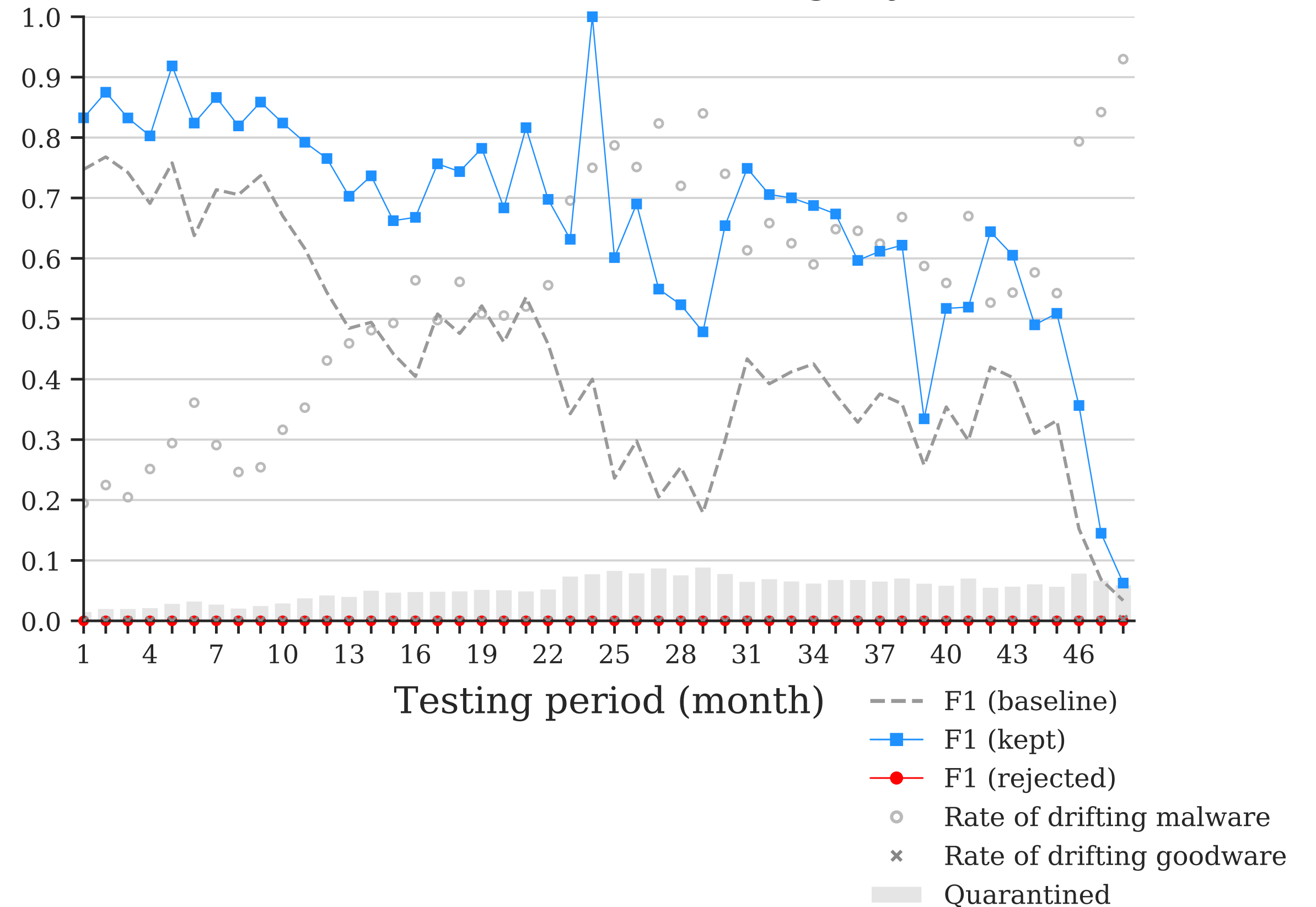
PDF Malware



Windows PE Malware



Android Malware (minimizing rejections)

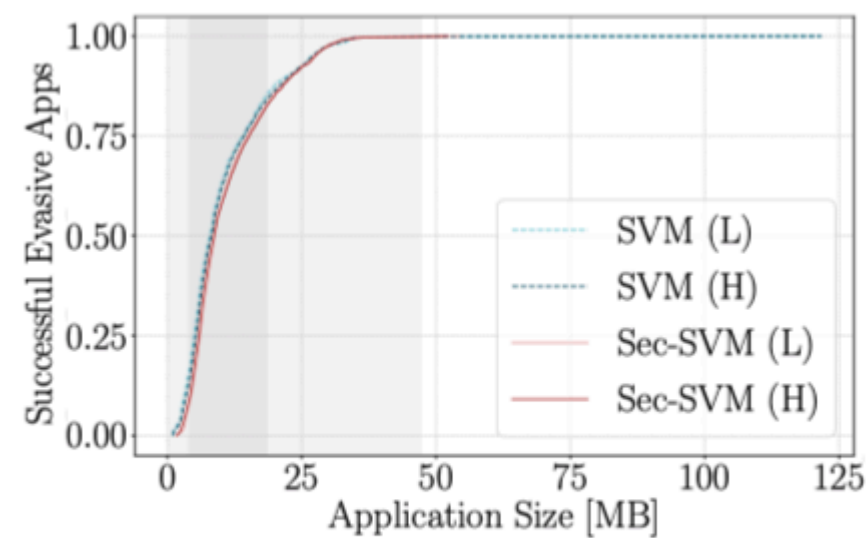


[IEEE S&P 2022] Transcending TRANSCEND: Revisiting Malware Classification in the Presence of Concept Drift

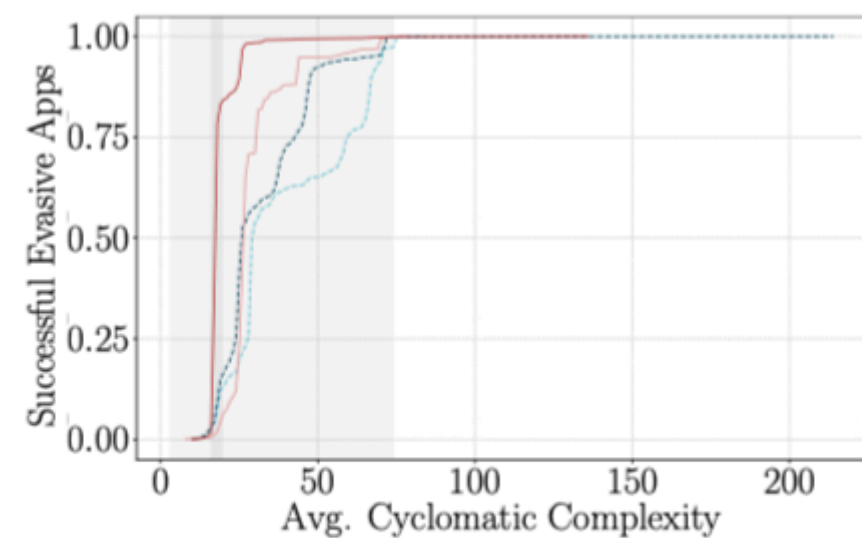
<https://s2lab.cs.ucl.ac.uk/projects/transcend/>

Results: How much are app statistics affected?

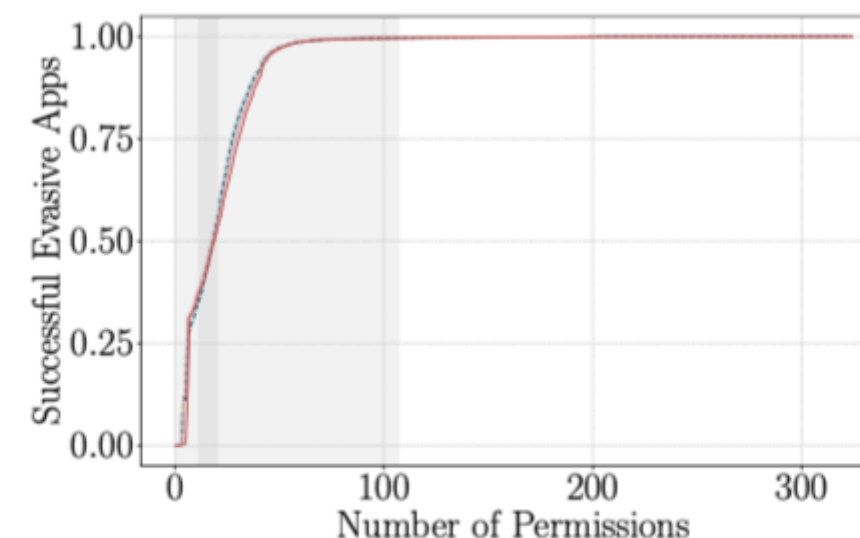
- Adding all these features (+ side-effect features), what does it do to app statistics?
 - › Limiting feature-space perturbations δ does not affect problem-space attack



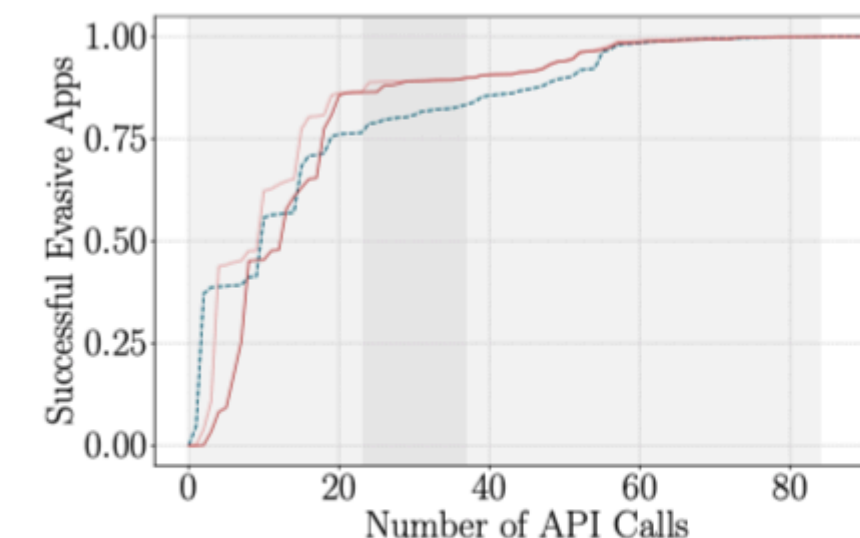
(a) Size



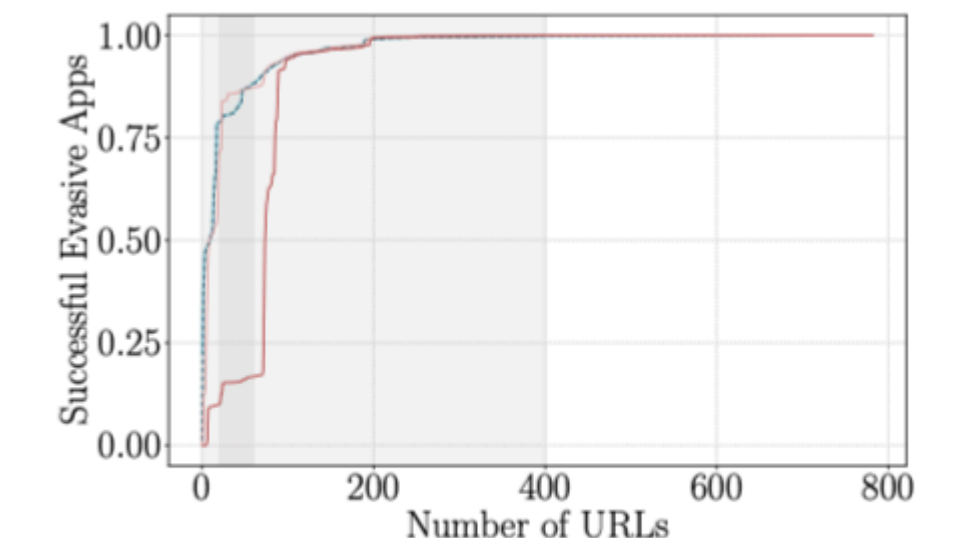
(b) Avg. CC



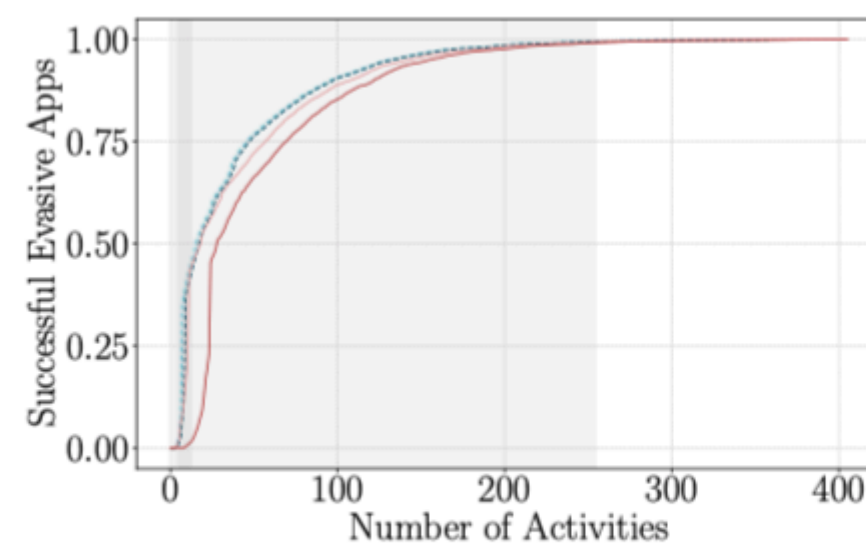
(c) Permissions



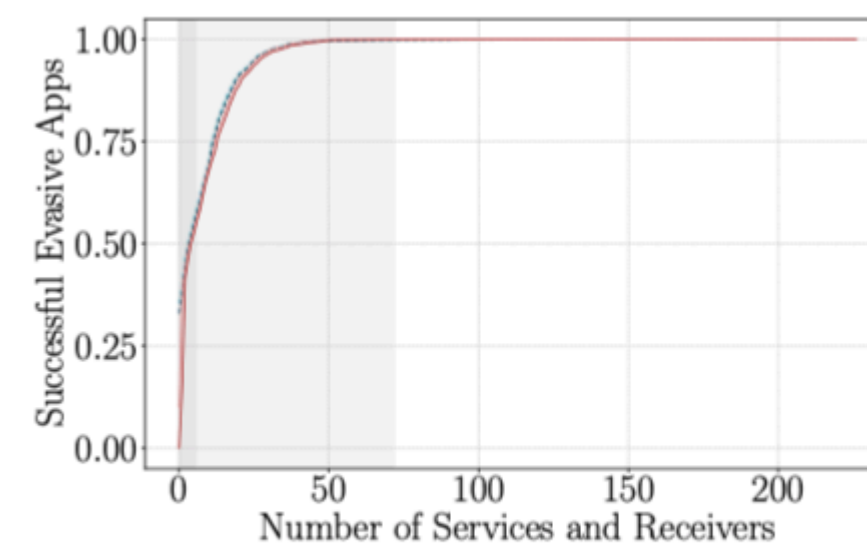
(d) API calls



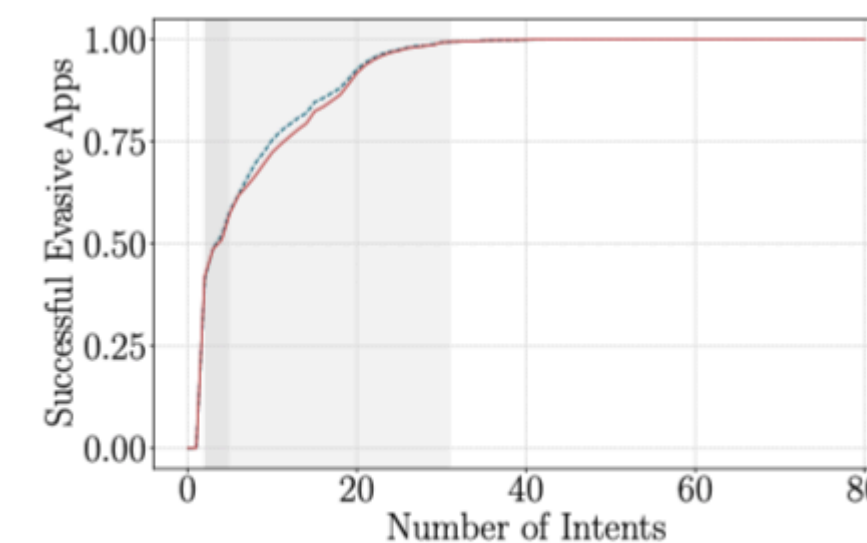
(e) URLs



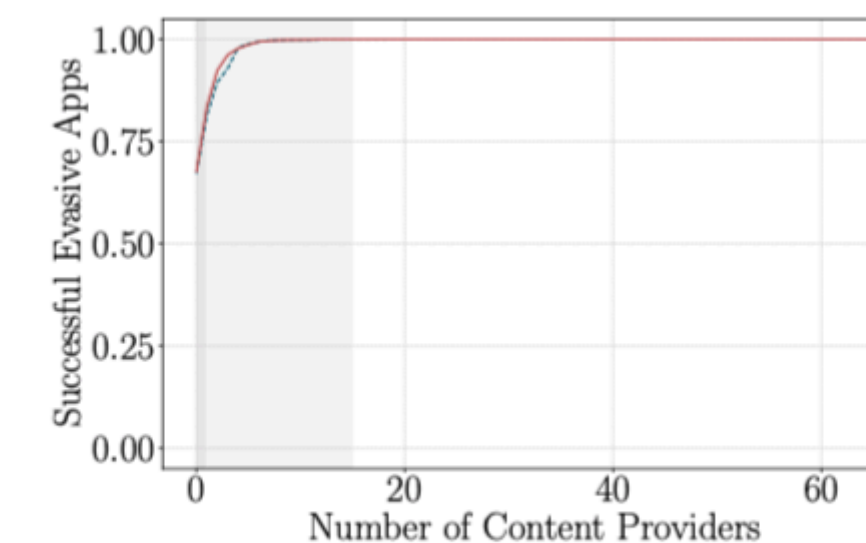
(f) Activities



(g) Services and Receivers



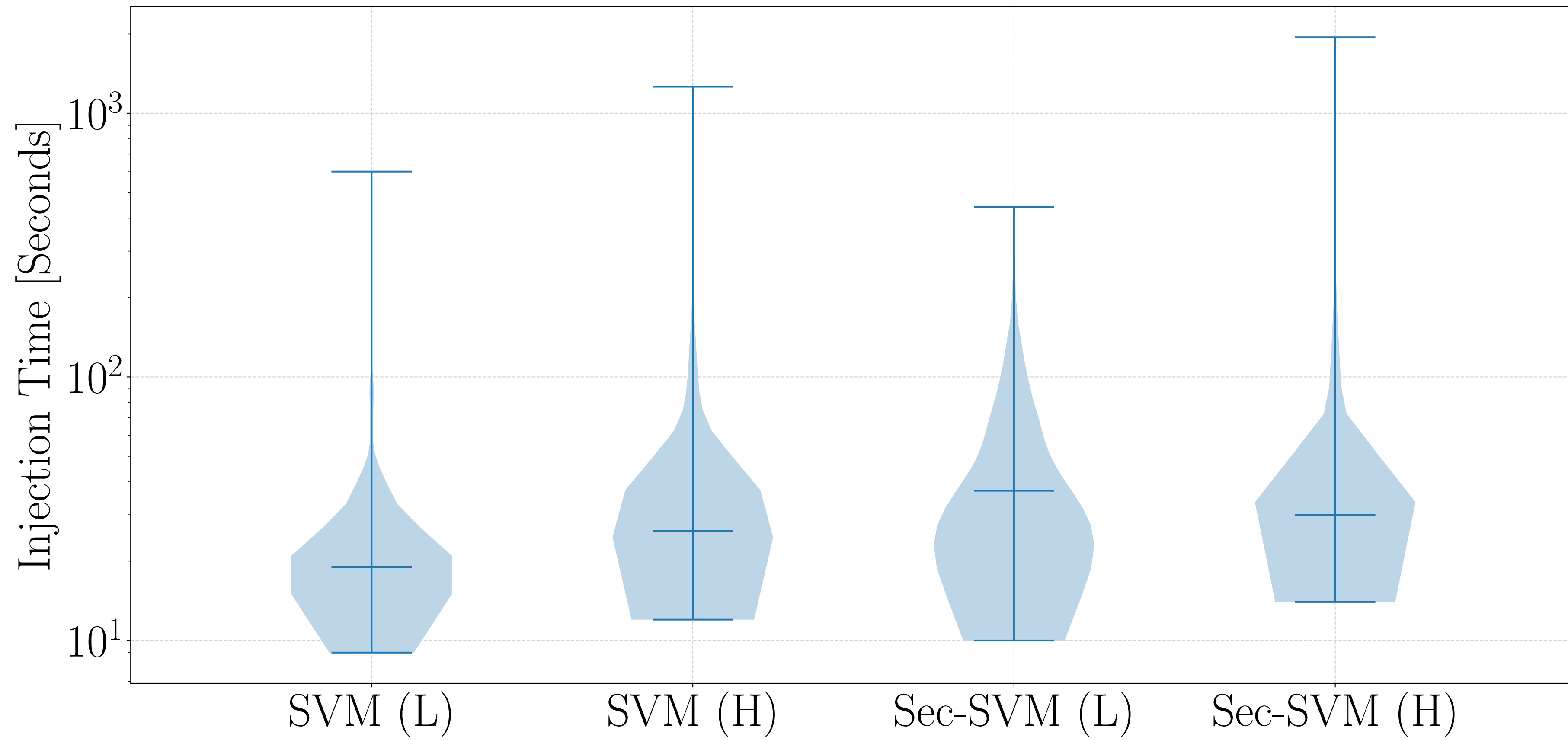
(h) Intents



(i) Content Providers

Results: How much time does an attack take?

- In most cases, **less than 2 minutes** to create an adversarial example



Results: How much time does an attack take?

- In most cases, **less than 2 minutes** to create an adversarial example

