# Trustworthy ML… for Systems Security

Lorenzo Cavallaro <l.cavallaro@ucl.ac.uk>
@lcavallaro — https://s2lab.cs.ucl.ac.uk

International Winter School on Microarchitectural Security 2022
FIAP, Paris

Dec 5, 2022

SYSTEMS
SECURITY
RESEARCH LAB

# 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
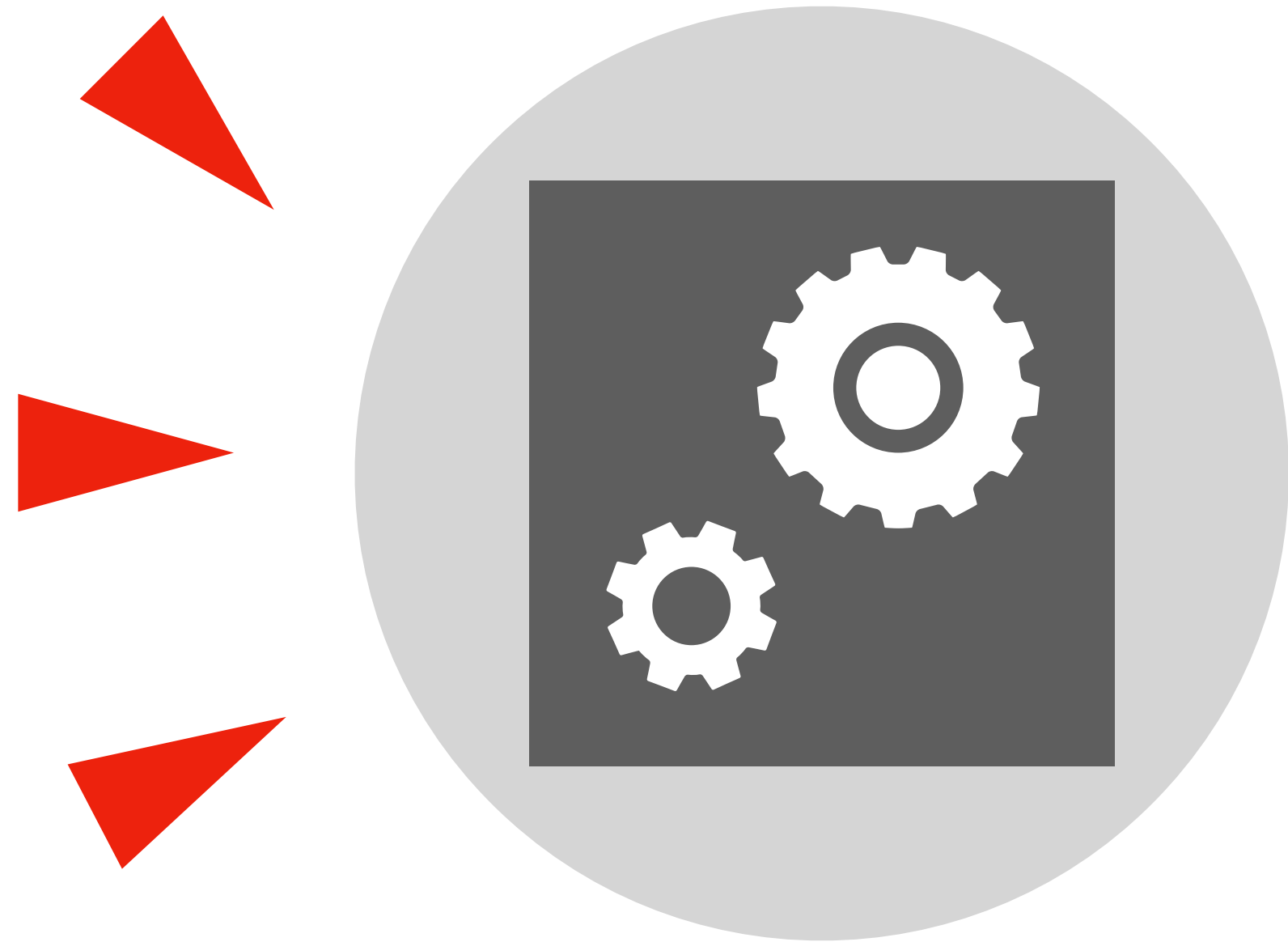
Android Malware

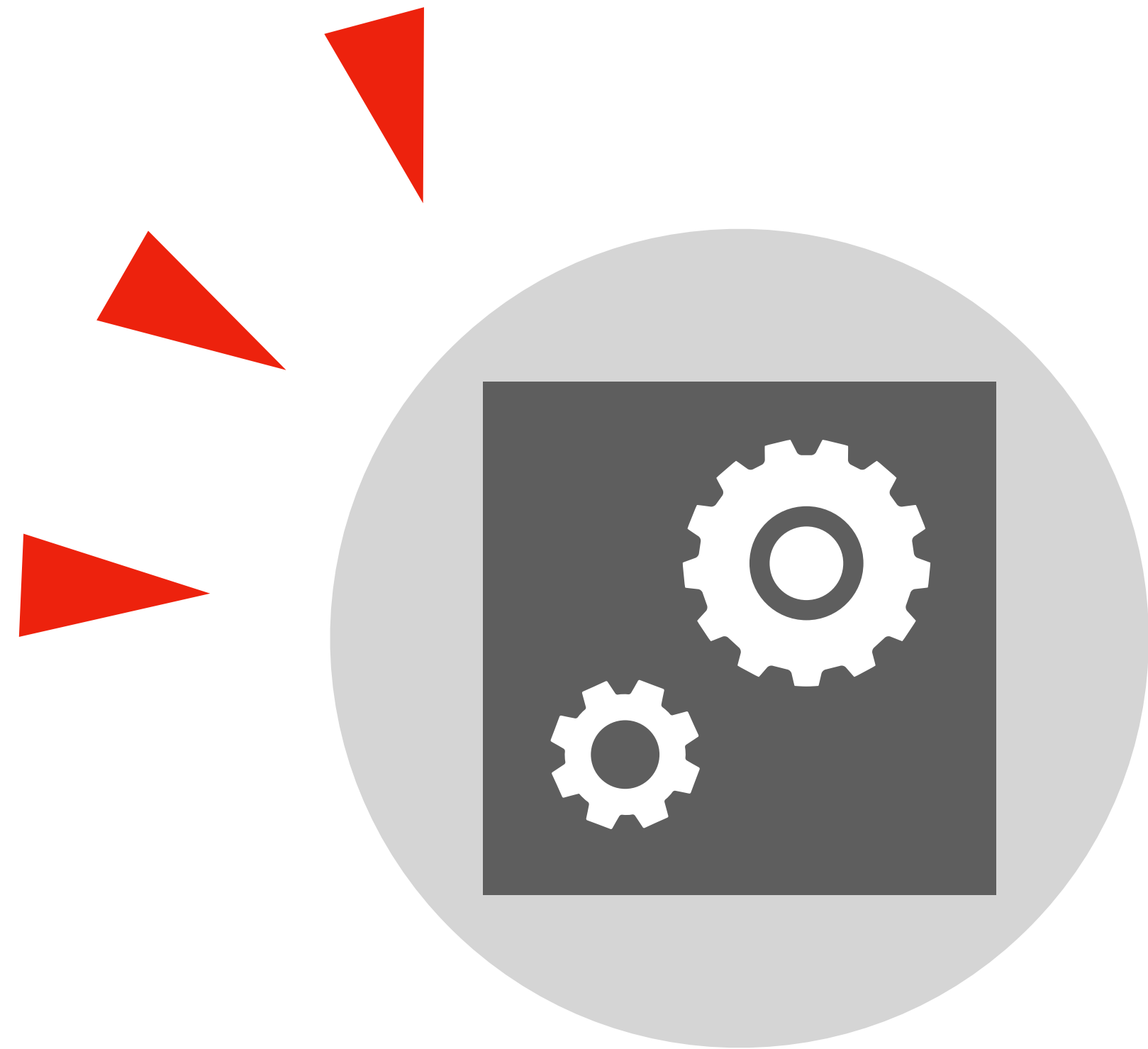Malicious Javascript

Windows Malware

PDF Malware

# Machine Learning Revolution



Image Classification

Facial Recognition

Machine Translation

Speech Recognition

Maybe the conditions aren't ready yet?

Android Malware

Malicious Javascript

Windows Malware

PDF Malware

2
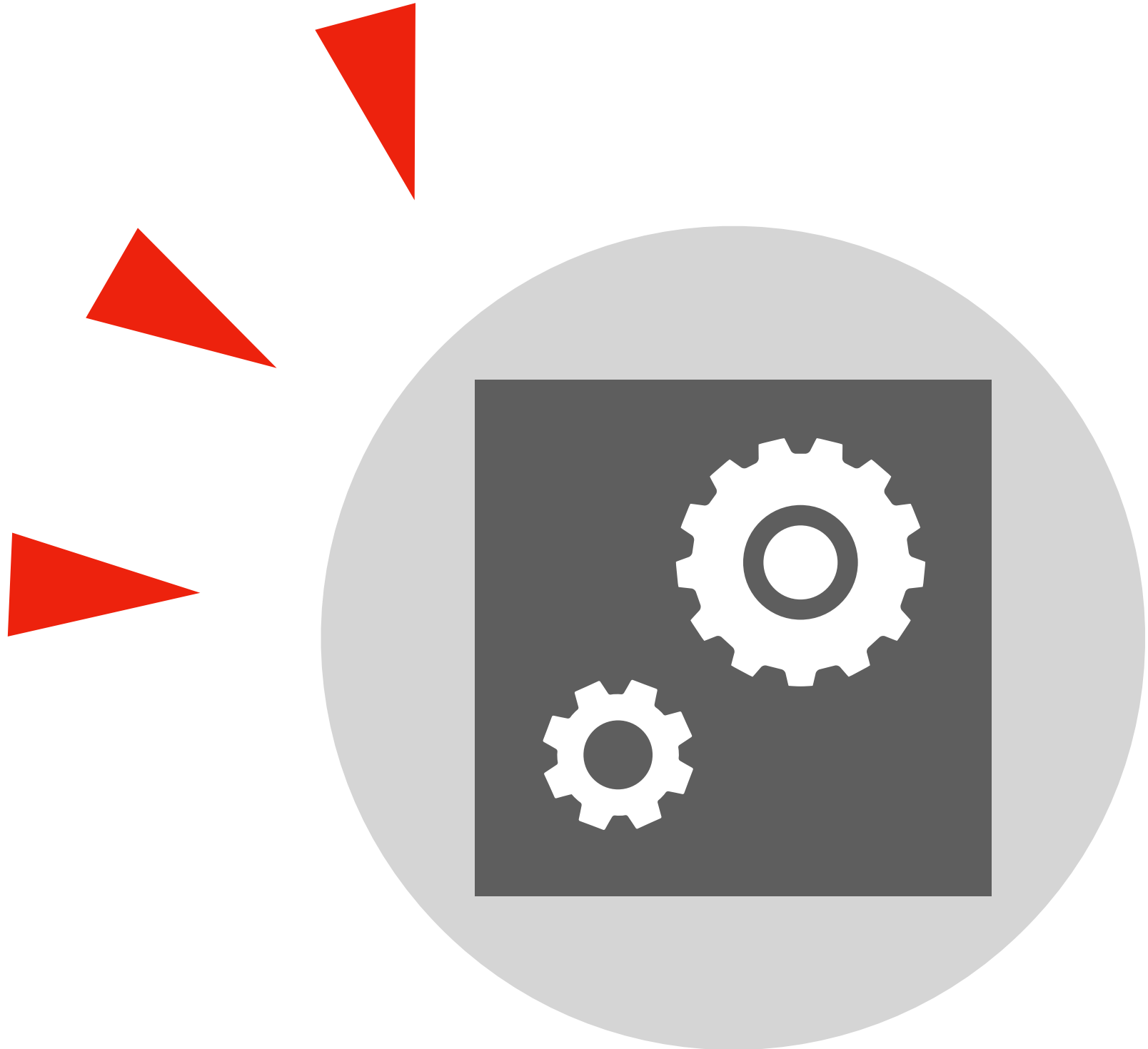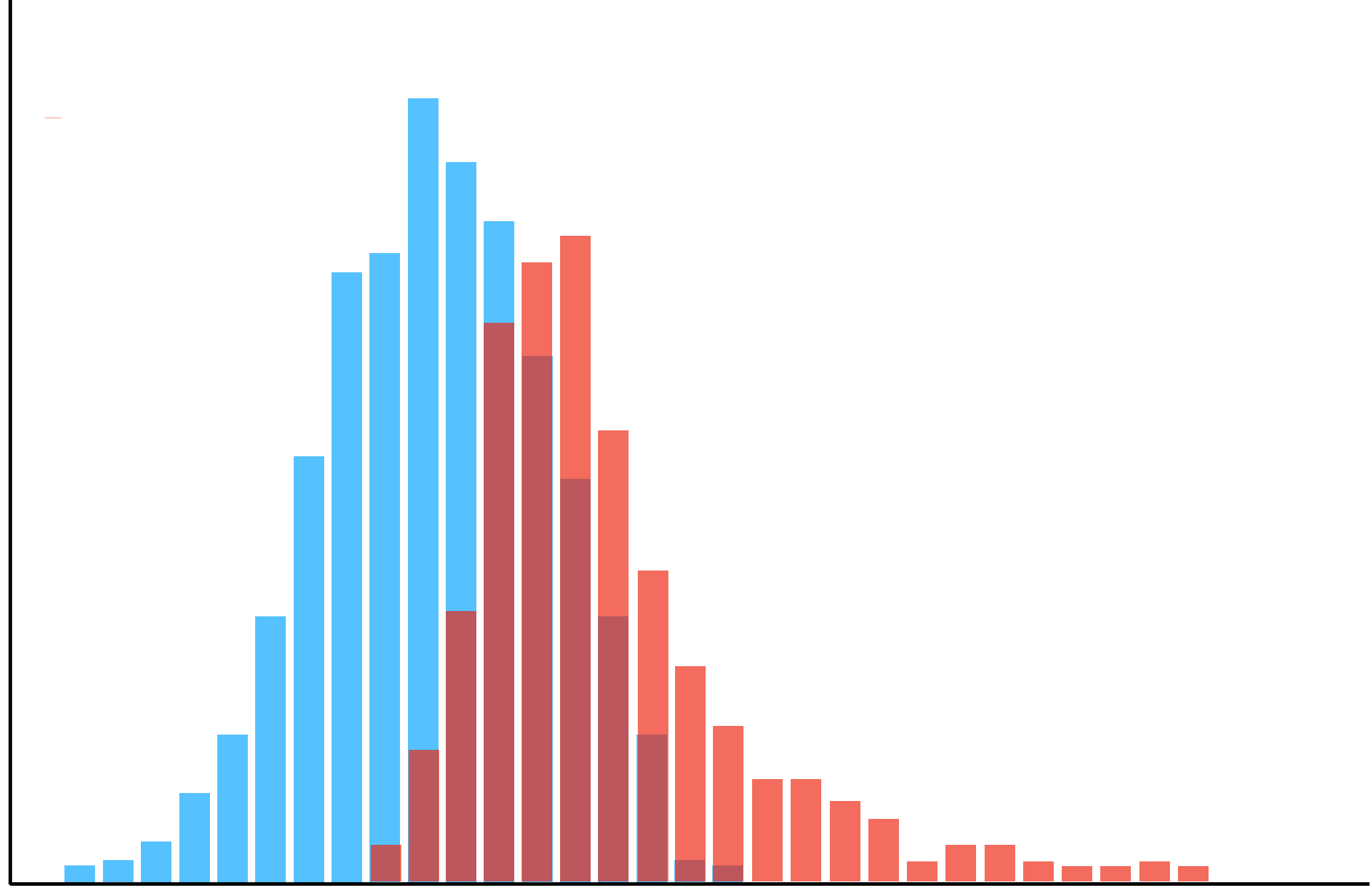
# Security *is* Adversarial

New detection systems trigger
an immediate response…

# Security *is* Adversarial
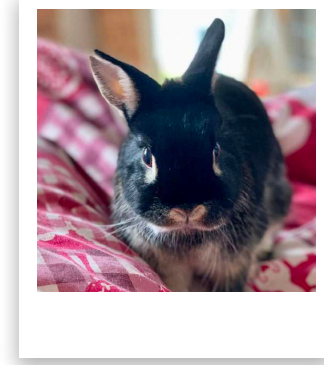
New detection systems trigger
an immediate response…

# Security *is* Adversarial
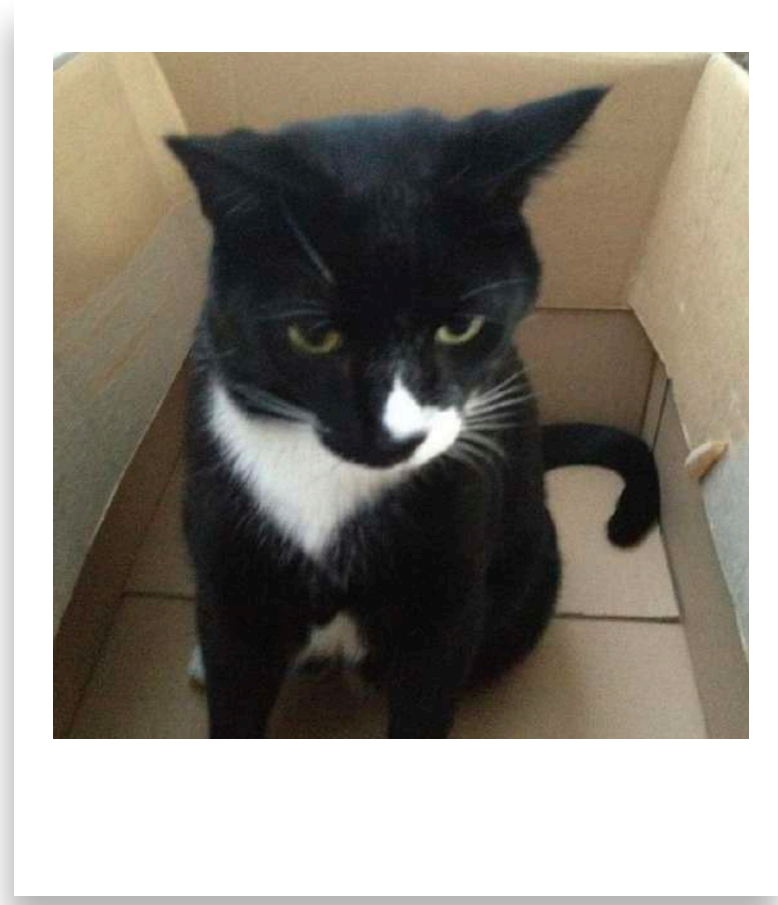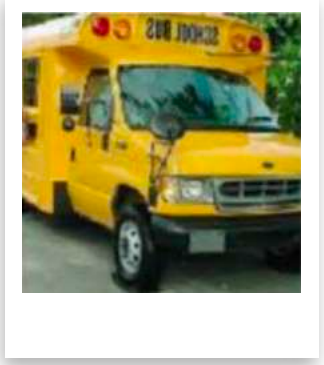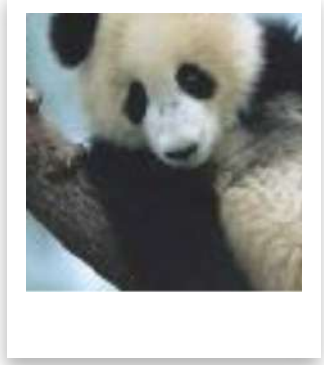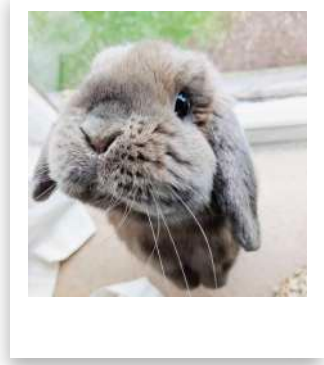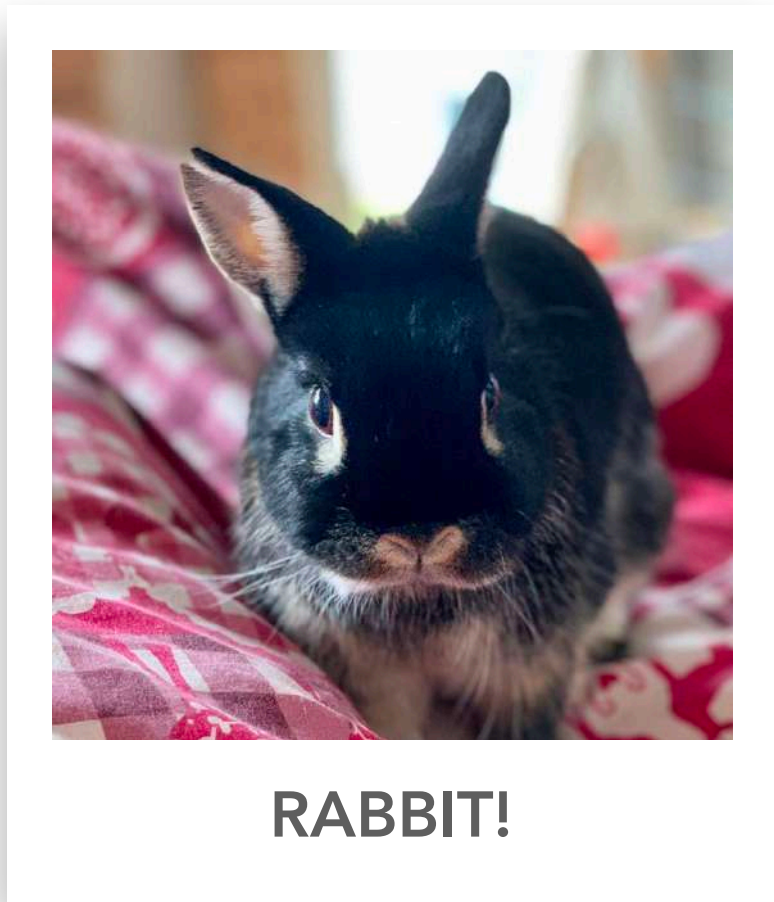


New detection systems trigger an immediate response…

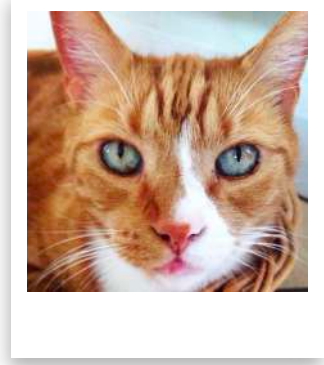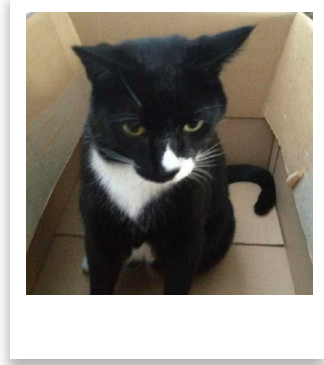…which causes dataset shifts, often violating the i.i.d. assumption

# Security *is* Adversarial

# Security *is* Adversarial



RABBIT!

# Security *is* Adversarial



RABBIT!

# Security *is* Adversarial



RABBIT!



Training on ad fraud…

# Security *is* Adversarial



RABBIT!

Training on ad fraud…            …attacks evolve at test time…

RANSOM

!YOUR MONEY OR YOUR FILES!

# Security *is* Adversarial

RABBIT!

70% OFF

WIN!

$ $ $

RANSOM

!YOUR MONEY OR YOUR FILES!

?

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...

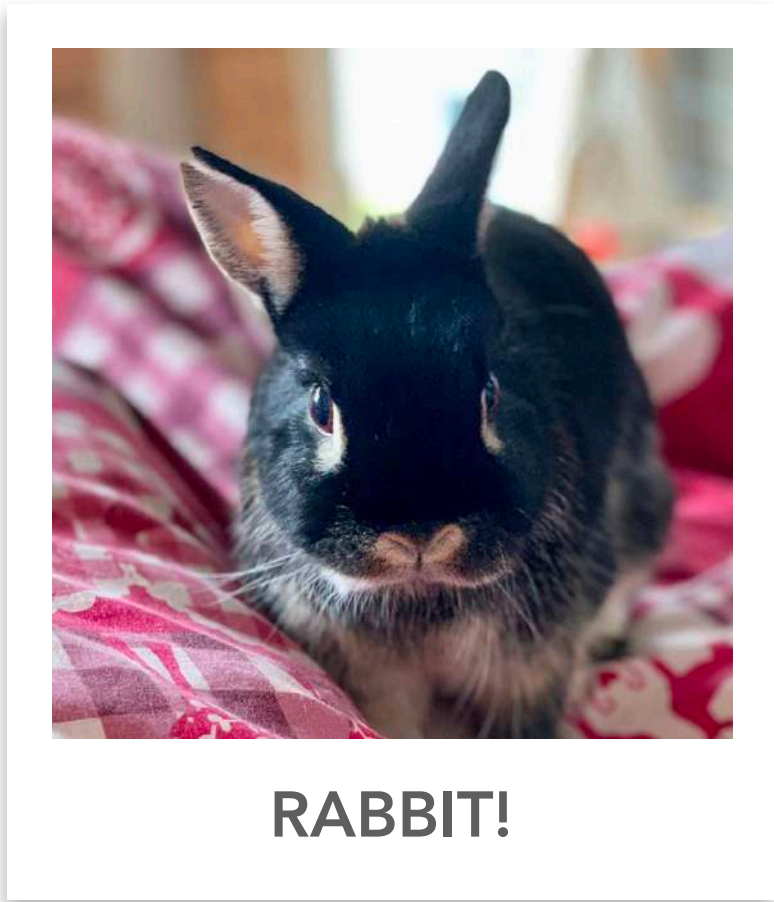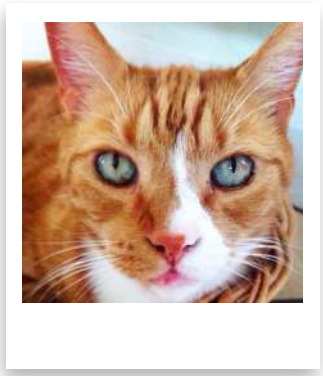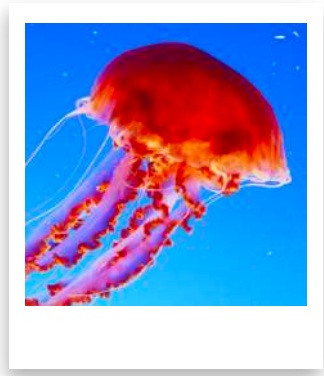To understand and improve the effectiveness of machine learning
methods for systems security in the presence of adversaries

**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 ML for systems security possible?**

# Outline

# Outline

## 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

- Practical end-to-end automatic adversarial malware as a service — how about defenses?

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

# Outline

## 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

- Practical end-to-end automatic adversarial malware as a service — how about defenses?

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

# Outline

# Outline

# Outline

## 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
- Practical end-to-end automatic adversarial malware as a service — how about defenses?

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

## Drifting scenarios caused by threats evolving over time

- How dataset drift affects machine learning-based detectors in security settings
- The need for time-aware evaluations and metrics
- Detecting shifts with abstaining classifiers and classification with rejection

**[USENIX Sec 2017 & IEEE S&P 2022]** Transcend: Detecting Concept Drift in Malware Classification Models & Transcending Transcend: Revisiting Malware Classification in the Presence of Concept Drift

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

## Quo vadis?

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

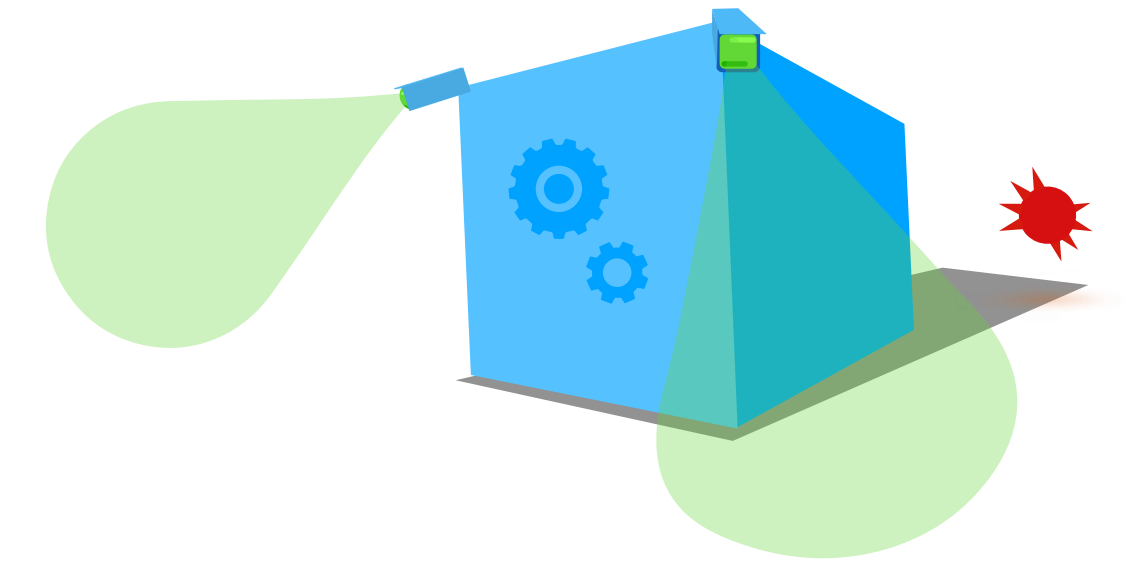**[USENIX Sec 2022]** Dos and Don'ts of Machine Learning in Computer Security

# 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

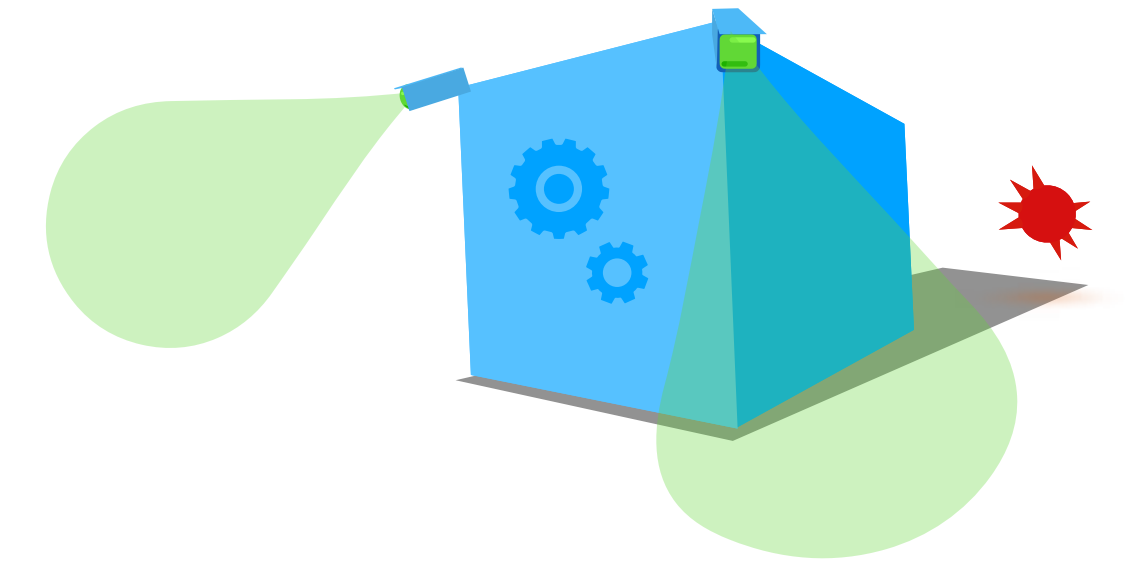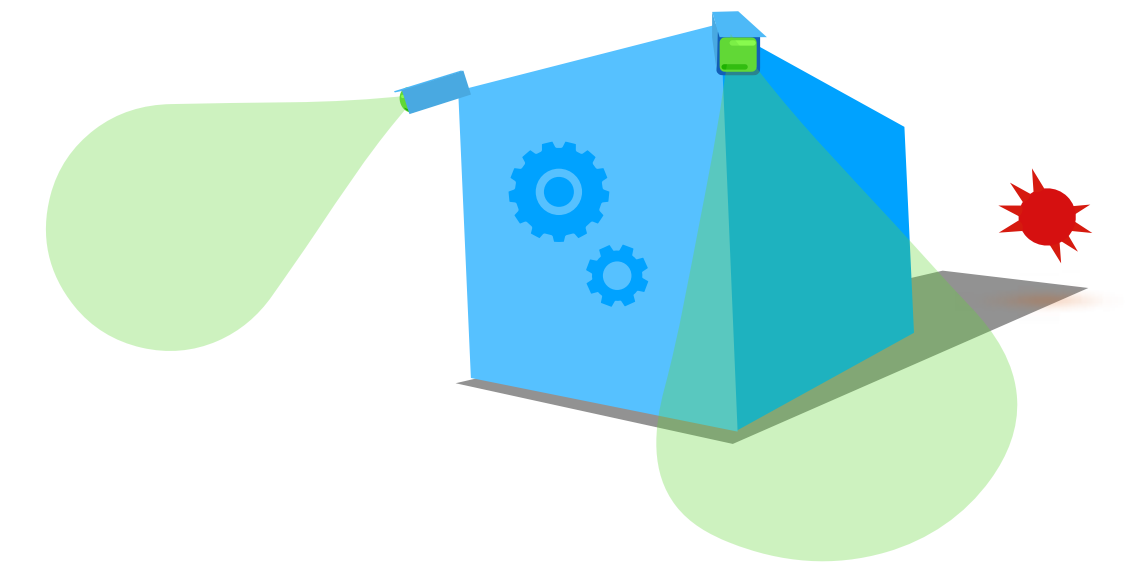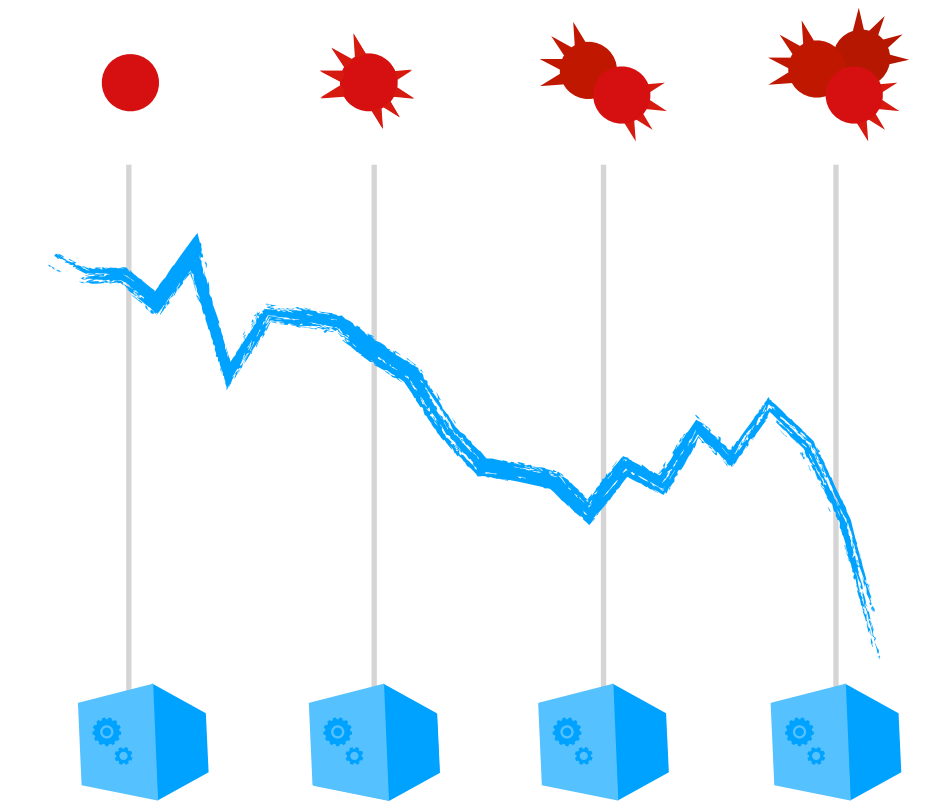- Practical end-to-end automatic adversarial malware as a service — how about defenses?

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

**Bigger Picture**

## Drifting scenarios caused by threats evolving over time

- How dataset shift affects machine learning-based detectors in security settings

- The need for time-aware evaluations and metrics

- Detecting shifts with abstaining classifiers and classification with rejection

**[USENIX Sec 2017 & IEEE S&P 2022]** Transcend: Detecting Concept Drift in Malware Classification Models & Transcending Transcend: Revisiting Malware Classification in the Presence of Concept Drift
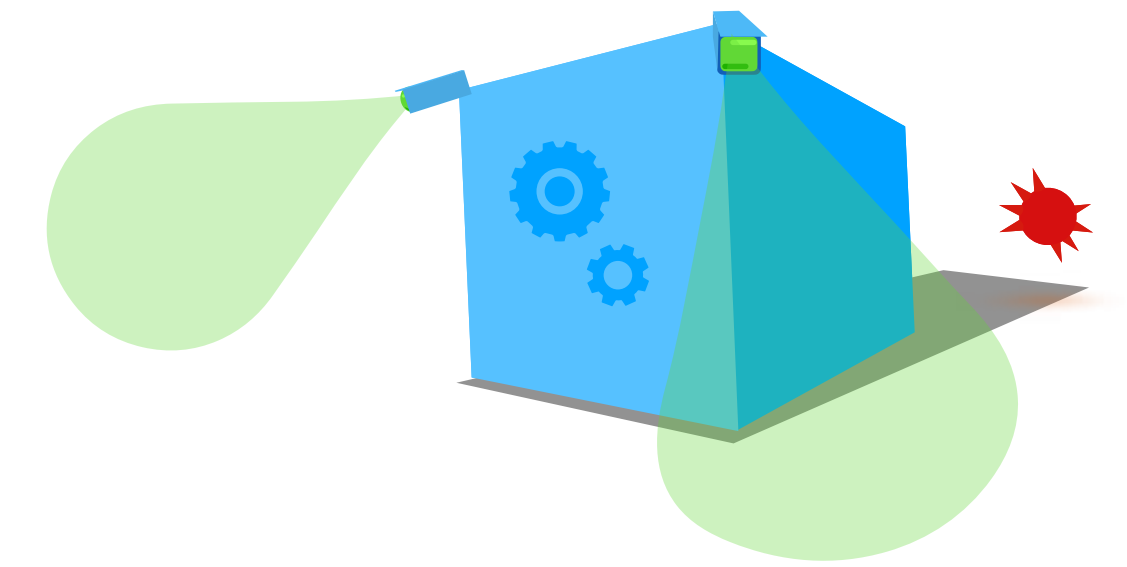
**[USENIX Sec 2019]** TESSERACT: Eliminating Experimental Bias in Malware Classification across Space and Time

**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

**[USENIX Sec 2022]** Dos and Don'ts of Machine Learning in Com

# Outline

## 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

- Practical end-to-end automatic adversarial malware as a service — how about defenses?
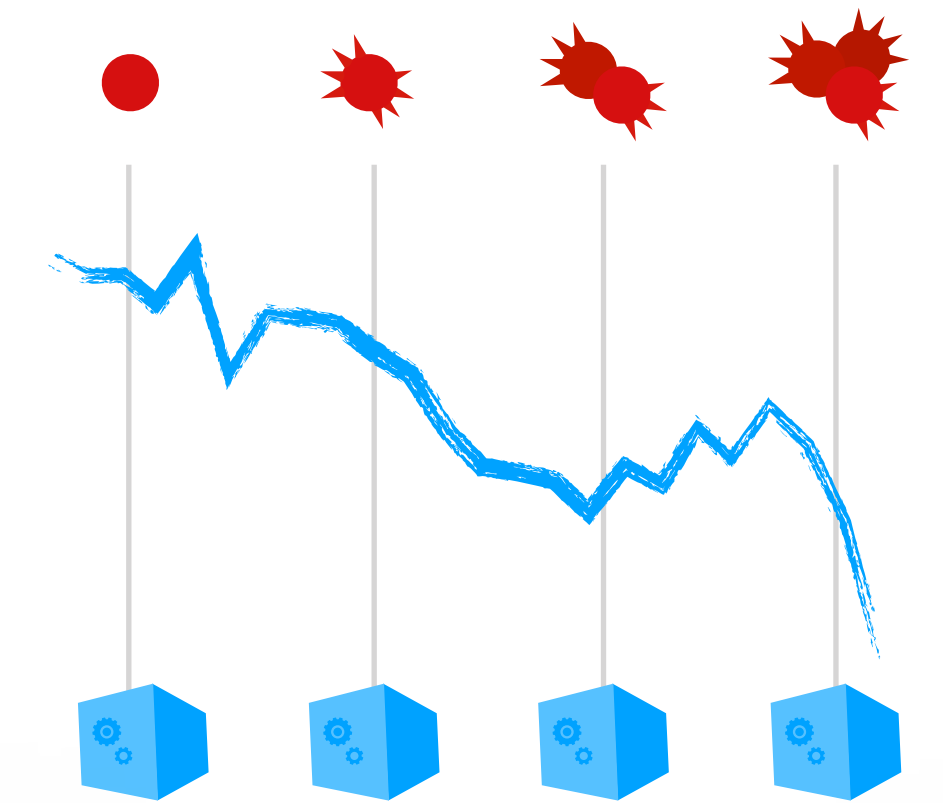
[IEEE S&P 2020] Intriguing Properties of Adversarial ML Attacks in the Problem Space

## Drifting scenarios caused by threats evolving over time

- How dataset shift affects machine learning-based detectors in security settings

- The need for time-aware evaluations and metrics

- Detecting shifts with abstaining classifiers and classification with rejection

[USENIX Sec 2017 & IEEE S&P 2022] Transcend: Detecting Concept Drift in Malware Classification Models & Transcending Transcend: Revisiting Malware Classification in the Presence of Concept Drift
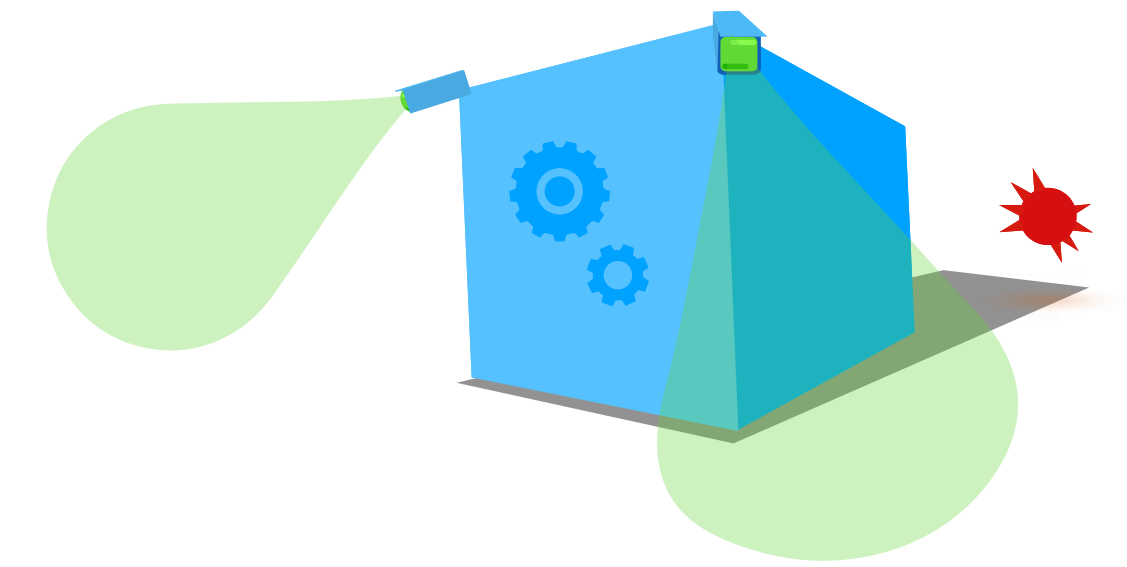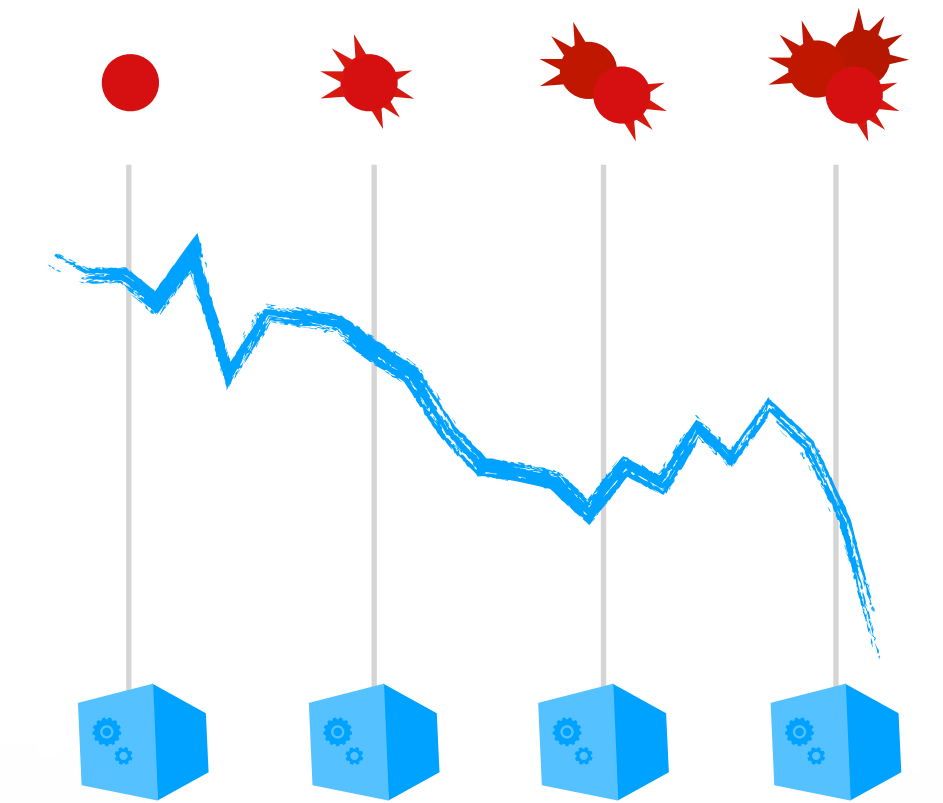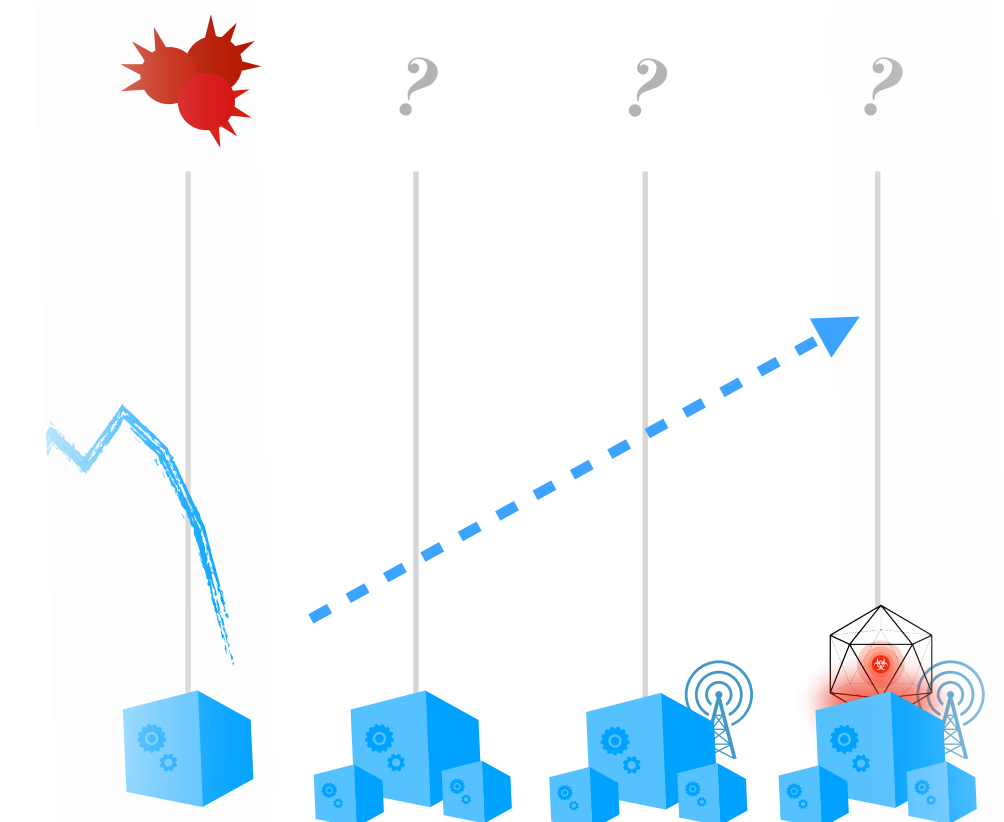
[USENIX Sec 2019] TESSERACT: Eliminating Experimental Bias in Malware Classification across Space and Time

## 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
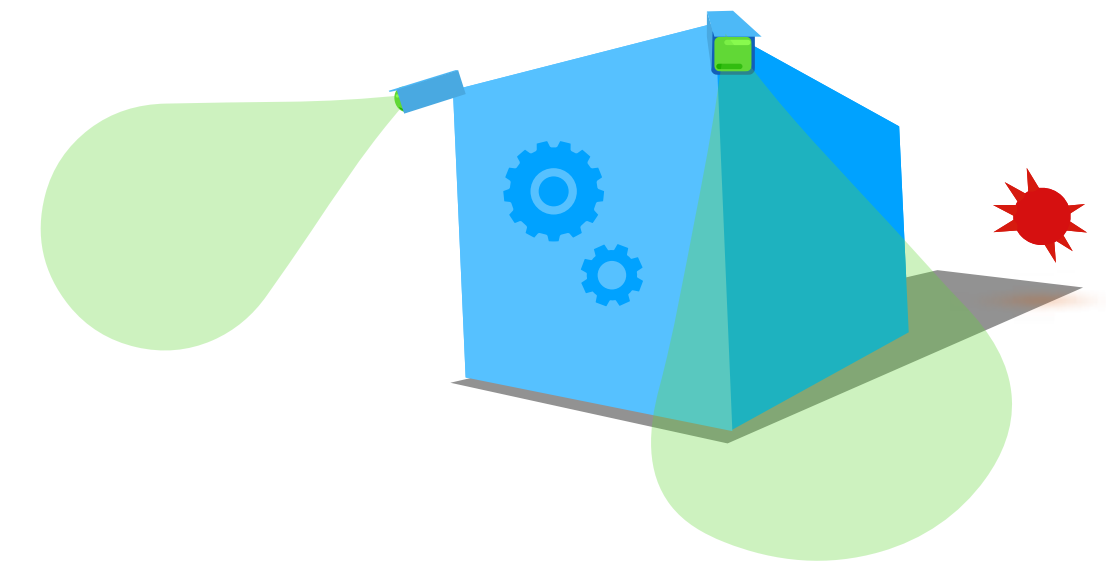
[USENIX Sec 2022] Dos and Don'ts of Machine Learning in Com

# 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 anal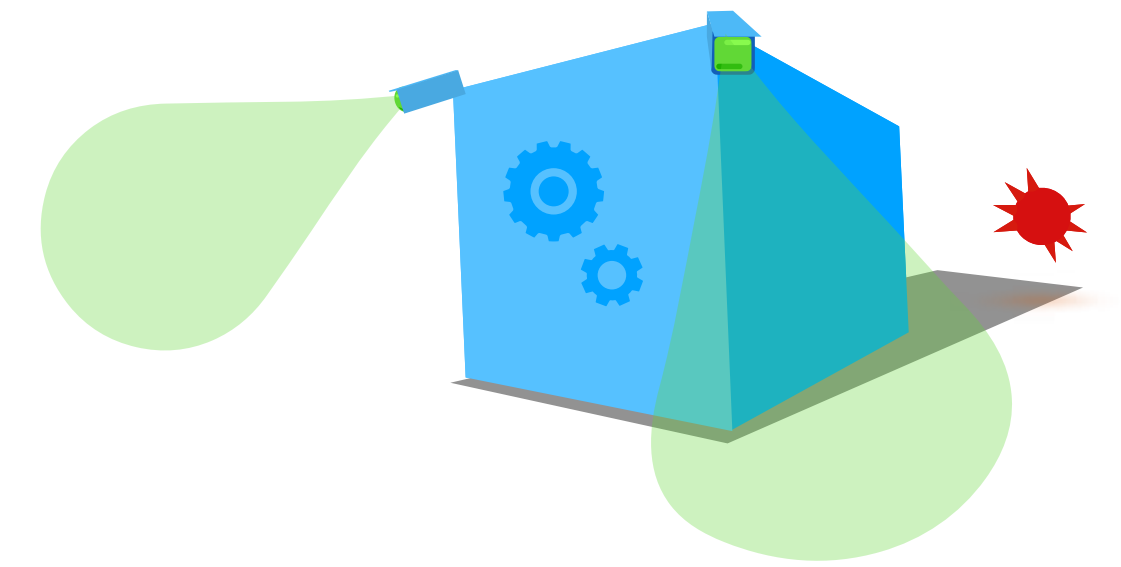ysis. It suggests that it is the space, rather than the individual units, that contains the semantic information in the high layers of neural networks.

Second, we find that deep neural networks learn input-output mappings that are ... continuous to a significant extent. We can cause the network to misclas... ... hardly perceptible perturbation, which is found ... In addition, the specific nature of ... turbation can

19 Feb 2014

[cs.CV]

13

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

Original Image

"panda" 57.7%

$x$

Perturbation

*imperceptible noise*

$\delta$

Adv. Image

"gibbon" 99.3%

$x + \delta$

# Let's Analyze What Happened



Feature-Space Attacks

| Original Image | Perturbation | Adv. Image |
|---|---|---|

"panda" 57.7%     *imperceptible noise*    "gibbon" 99.3%

$$x \qquad \delta \qquad x + \delta$$

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

# Let's Analyze What Happened



Feature-Space Attacks

Original Image

Perturbation

Adv. Image

"panda" 57.7%

*imperceptible noise*

"gibbon" 99.3%

$$x$$

$$\delta$$

$$x + \delta$$

Optimization

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

Pixel
Perturbations

# Let's Analyze What Happened

**Feature-Space Attacks**

| Original Image | Perturbation | Adv. Image |
|---|---|---|
|  |  |  |
| "panda" 57.7% | *imperceptible noise* | "gibbon" 99.3% |

$$x \qquad + \qquad \delta \qquad = \qquad x + \delta$$

Optimization  |  $\text{minimize}_\delta \; \left( ||\delta||_p \right) + c \cdot \left( f(x + \delta) \right)$

Pixel Perturbations

Loss of Target Class

# Let's Analyze What Happened

**Feature-Space Attacks**

| Original Image | Perturbation | Adv. Image |
|---|---|---|



"panda" 57.7%    +    *imperceptible noise*    =    "gibbon" 99.3%

$$x \qquad \delta \qquad x + \delta$$

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

Pixel Perturbations

Loss of Target Class

**[IEEE S&P 2020]** Intriguing Properties of Adversarial ML Attacks in the Problem Space
https://s2lab.cs.ucl.ac.uk/projects/intriguing

16

# Let's Analyze What Happened

## Problem-Space Attacks

### Feature-Space Attacks

**Problem Space**

| Original App (z) | "Perturbation" | Adversarial App (z') |
|---|---|---|
|  |  |  |
| "malware" 57.7% | ? | "goodware" 95.7% |

Original Image

Perturbation

"panda" 57.7%

*imperceptible noise*

$x$

$\delta$

**Feature Space**

$$x \qquad \delta \qquad x + \delta$$

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

Pixel
Perturbations

Loss of
Target Class

**[IEEE S&P 2020]** Intriguing Properties of Adversarial ML Attacks in the Problem Space
https://s2lab.cs.ucl.ac.uk/projects/intriguing

16

# Let's Analyze What Happened

## Problem-Space Attacks

### Feature-Space Attacks

**Problem Space**

| Original App (z) | "Perturbation" | Adversarial App (z') |
|---|---|---|
| "malware" 57.7% | ? | "goodware" 95.7% |

$+$ $=$

**Feature Space**

$$x \qquad \delta \qquad x + \delta$$

Original Image

Perturbation

"panda" 57.7%

*imperceptible noise*

$x$ $\delta$

Optimization $\quad$ minimize$_\delta \left( ||\delta||_p + c \cdot f(x + \delta) \right)$

Pixel Perturbations

Loss of Target Class

Optimization $\quad\Big|\quad$ **minimize**$_\delta \quad ||\delta||_p + c \cdot f(x + \delta)$

# Let's Analyze What Happened

## Problem-Space Attacks

**Problem Space**

| Original App (z) | "Perturbation" | Adversarial App (z') |
|---|---|---|
|  |  |  |
| "malware" 57.7% | **?** | "goodware" 95.7% |

+ =

**Feature Space**

$$x \qquad \delta \qquad x + \delta$$

Optimization

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

Constraints
- Is it realistic/plausible?
- Does it crash?
- Does it preserve malicious functionality?
- … are there **"general" constraints**?

16

# Let's Analyze What Happened

**Problem-Space Attacks**

Feature-Space Attacks

**Problem Space**

Original App (z)   "Perturbation"   Adversarial App (z')



"malware" 57.7%   **?**   "goodware" 95.7%

Original Image  Perturbation

"panda" 57.7%  *imperceptible noise*

$x$    $\delta$

**Feature Space**

$x \longrightarrow \delta \longrightarrow x + \delta$

**?** Inverse Feature-Mapping Problem

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

Pixel Perturbations

**Optimization**  $\textbf{minimize}_\delta \quad ||\cancel{\delta}||_p + c \cdot f(x + \delta)$

Loss of Target Class

**Constraints**
- Is it realistic/plausible?
- Does it crash?
- Does it preserve malicious functionality?
- … are there **"general" constraints**?

# Inverse Feature-Mapping Problem

**Problem Space**

**Feature Space**

**[IEEE S&P 2020]** Intriguing Properties of Adversarial ML Attacks in the Problem Space
https://s2lab.cs.ucl.ac.uk/projects/intriguing

# Inverse Feature-Mapping Problem



The feature mapping $\varphi$ is <u>differentiable</u>
— you can backpropagate to input

**[IEEE S&P 2020]** Intriguing Properties of Adversarial ML Attacks in the Problem Space
https://s2lab.cs.ucl.ac.uk/projects/intriguing

# Inverse Feature-Mapping Problem

**Problem Space**
Code



"malware" 57.7%

**Feature Space**
Vectors



**[IEEE S&P 2020]** Intriguing Properties of Adversarial ML Attacks in the Problem Space
https://s2lab.cs.ucl.ac.uk/projects/intriguing

# Inverse Feature-Mapping Problem



**Problem Space**
Code

**Feature Space**
Vectors

$\varphi$

$x$

| 1 | 0 | 1 | 0 | ... | 1 |

"malware" 57.7%

FS attack

$x + \delta$

| 0 | 1 | 0 | 1 | ... | 0 |

???

"goodware" 95.7%

In the software domain,
the feature mapping $\varphi$ is
neither <u>invertible</u> nor <u>differentiable</u>
— how to get back to the problem space?

**[IEEE S&P 2020]** Intriguing Properties of Adversarial ML Attacks in the Problem Space
https://s2lab.cs.ucl.ac.uk/projects/intriguing

# Many Problem-Space Attack Papers

**Android Malware**
[TDSC'17, ESORICS'17, ACSAC'19]

**Windows Malware**
[RAID'18, EUSIPCO'18]

**PDF Malware**
[ECML-PKDD'13, NDSS'16]

**Network Traffic**
[NCA'18, NCA'19]

# Many Problem-Space Attack Papers

**Android Malware**

[TDSC'17, ESORICS'17, ACSAC'19]

**Windows Malware**

[RAID'18, EUSIPCO'18]

**PDF Malware**

[ECML-PKDD'13, NDSS'16]

**Network Traffic**

[NCA'18, NCA'19]

## What is the State of the Art? How to compare them?

# Outline

## Formalization

- Problem-space attacks

- Relationships

- Actionable points

## Android Problem-Space Attack

- End-to-end adversarial malware generation at scale

- Feasible to evade feature-space defenses

# Outline

**Formalization**

- Problem-space attacks

- Relationships

- Actionable points

**Android Problem-Space Attack**

- End-to-end adversarial malware generation at scale

- Feasible to evade feature-space defenses

Evasion Attacks

# Outline

## Formalization

- Problem-space attacks
- Relationships
- Actionable points

## Android Problem-Space Attack

- End-to-end adversarial malware generation at scale
- Feasible to evade feature-space defenses

Evasion Attacks

Running example: Code

**[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

**[IEEE S&P 2020]** Intriguing Properties of Adversarial ML Attacks in the Problem Space
https://s2lab.cs.ucl.ac.uk/projects/intriguing

# Problem-Space Constraints

**Available Transformations**                    How can you alter problem-space objects?

**[IEEE S&P 2020]** Intriguing Properties of Adversarial ML Attacks in the Problem Space
https://s2lab.cs.ucl.ac.uk/projects/intriguing

# 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?

Addition

"strings"

bytes

Software

**[IEEE S&P 2020]** Intriguing Properties of Adversarial ML Attacks in the Problem Space
https://s2lab.cs.ucl.ac.uk/projects/intriguing

# Problem-Space Constraints

🔧 **Available Transformations**

How can you alter problem-space objects?



Addition

Removal

"strings"

bytes

Software

"strings"

bytes

**[IEEE S&P 2020]** Intriguing Properties of Adversarial ML Attacks in the Problem Space
https://s2lab.cs.ucl.ac.uk/projects/intriguing

22

# Problem-Space Constraints

**Available Transformations**    How can you alter problem-space objects?



Addition

Removal

"strings"

bytes

Software

"strings"

bytes

**[IEEE S&P 2020]** Intriguing Properties of Adversarial ML Attacks in the Problem Space
https://s2lab.cs.ucl.ac.uk/projects/intriguing

# Problem-Space Constraints

**Available Transformations**                How can you alter problem-space objects?

Addition                    Removal

API                                                    API

"strings"                                                    "strings"

Software

bytes                                                    bytes

Modification

**[IEEE S&P 2020]** Intriguing Properties of Adversarial ML Attacks in the Problem Space
https://s2lab.cs.ucl.ac.uk/projects/intriguing

# Problem-Space Constraints

 **Available Transformations**

# Problem-Space Constraints

🔴 **Preserved Semantics**

🟢 **Available Transformations**

Which semantics do you preserve? How?

Which automatic tests can verify it?



Malicious Node

# Problem-Space Constraints

**Preserved Semantics**

Available Transformations

Which semantics do you preserve? How?

Which automatic tests can verify it?



Malicious Node

# Problem-Space Constraints

**Preserved Semantics**

**Available Transformations**

Which semantics do you preserve? How?

Which automatic tests can verify it?

**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

Malicious Node

# Problem-Space Constraints

 **Preserved Semantics**

 **Available Transformations**

# Problem-Space Constraints

👁 **Plausibility**

⚛ Preserved Semantics

🔧 Available Transformations

Does it look legit?

**[IEEE S&P 2020]** Intriguing Properties of Adversarial ML Attacks in the Problem Space
https://s2lab.cs.ucl.ac.uk/projects/intriguing

# Problem-Space Constraints

👁 **Plausibility**

⚛ Preserved Semantics

🔧 Available Transformations

**Does it look legit?**

**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?



**[IEEE S&P 2020]** Intriguing Properties of Adversarial ML Attacks in the Problem Space
https://s2lab.cs.ucl.ac.uk/projects/intriguing

# Problem-Space Constraints

⚙️ **Robustness to Preprocessing**

👁️ Plausibility

⚛️ Preserved Semantics

🔧 Available Transformations

Which preprocessing are you considering?

**[IEEE S&P 2020]** Intriguing Properties of Adversarial ML Attacks in the Problem Space
https://s2lab.cs.ucl.ac.uk/projects/intriguing

# Problem-Space Constraints

⚙️ **Robustness to Preprocessing**          Which preprocessing are you considering?

👁️ Plausibility

⚛️ Preserved Semantics

🔧 Available Transformations



ANALYSIS

25

# Side-effect Features

26

# Side-effect Features

# Side-effect Features

# Side-effect Features

# Side-effect Features

# Side-effect Features

# Side-effect Features

# Side-effect Features

# Actionable Points

**[IEEE S&P 2020]** Intriguing Properties of Adversarial ML Attacks in the Problem Space
https://s2lab.cs.ucl.ac.uk/projects/intriguing

# Actionable Points

Verify existence of feature-space attack

**Necessary Condition for problem-space attacks**

$\exists$ **problem-space attack** $\implies$ $\exists$ **feature-space attack**

Proof 1
in paper

**[IEEE S&P 2020]** Intriguing Properties of Adversarial ML Attacks in the Problem Space
https://s2lab.cs.ucl.ac.uk/projects/intriguing

# Our Android Attack

32

# 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.

**[IEEE S&P 2020]** Intriguing Properties of Adversarial ML Attacks in the Problem Space
https://s2lab.cs.ucl.ac.uk/projects/intriguing

# Organic Harvesting

34

# Organ Harvesting

DEX

① Identify feature entry point

Identify activity in dex

**[IEEE S&P 2020]** Intriguing Properties of Adversarial ML Attacks in the Problem Space
https://s2lab.cs.ucl.ac.uk/projects/intriguing

# Organ Harvesting

**①** Identify feature entry point

DEX

**Identify activity in dex**

**[IEEE S&P 2020]** Intriguing Properties of Adversarial ML Attacks in the Problem Space
https://s2lab.cs.ucl.ac.uk/projects/intriguing

# Organn Harvesting



② Choose any vein (backward slice)

**Extract intent creation and startActivity()**

36

# Organ Harvesting

**DEX**

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**

DEX

# Organ Harvesting



④ Include transitive dependencies

**Recursively collect dependencies**

**[IEEE S&P 2020]** Intriguing Properties of Adversarial ML Attacks in the Problem Space
https://s2lab.cs.ucl.ac.uk/projects/intriguing

# Organ Harvesting

**④** Include transitive dependencies

DEX

**Recursively collect dependencies**

# Organ Harvesting

⑤ Store organ in an "ice box"

DEX

**Save gadget to a database ready for the attack**

**[IEEE S&P 2020]** Intriguing Properties of Adversarial ML Attacks in the Problem Space
https://s2lab.cs.ucl.ac.uk/projects/intriguing

# Organ Harvesting

⑤ Store organ in an "ice box"

Save gadget to a database ready for the attack

# Attack Overview

40

# Attack Overview

Given a trained target model

# Attack Overview

Given a trained target model
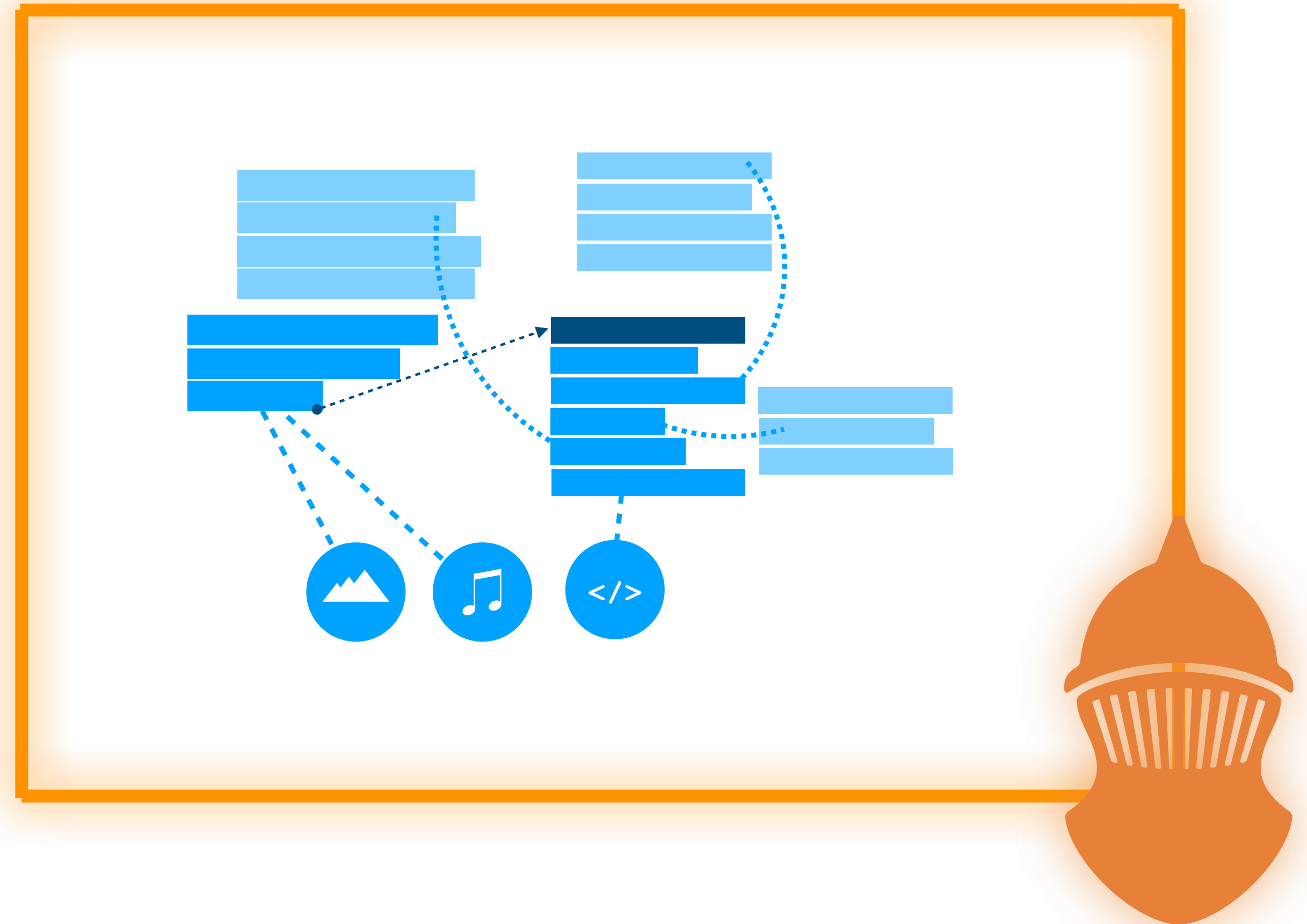
First pick feature with greatest 'benign' weight

# Attack Overview



Given a trained target model

First pick feature with greatest 'benign' weight

Find a corresponding organ from the ice box
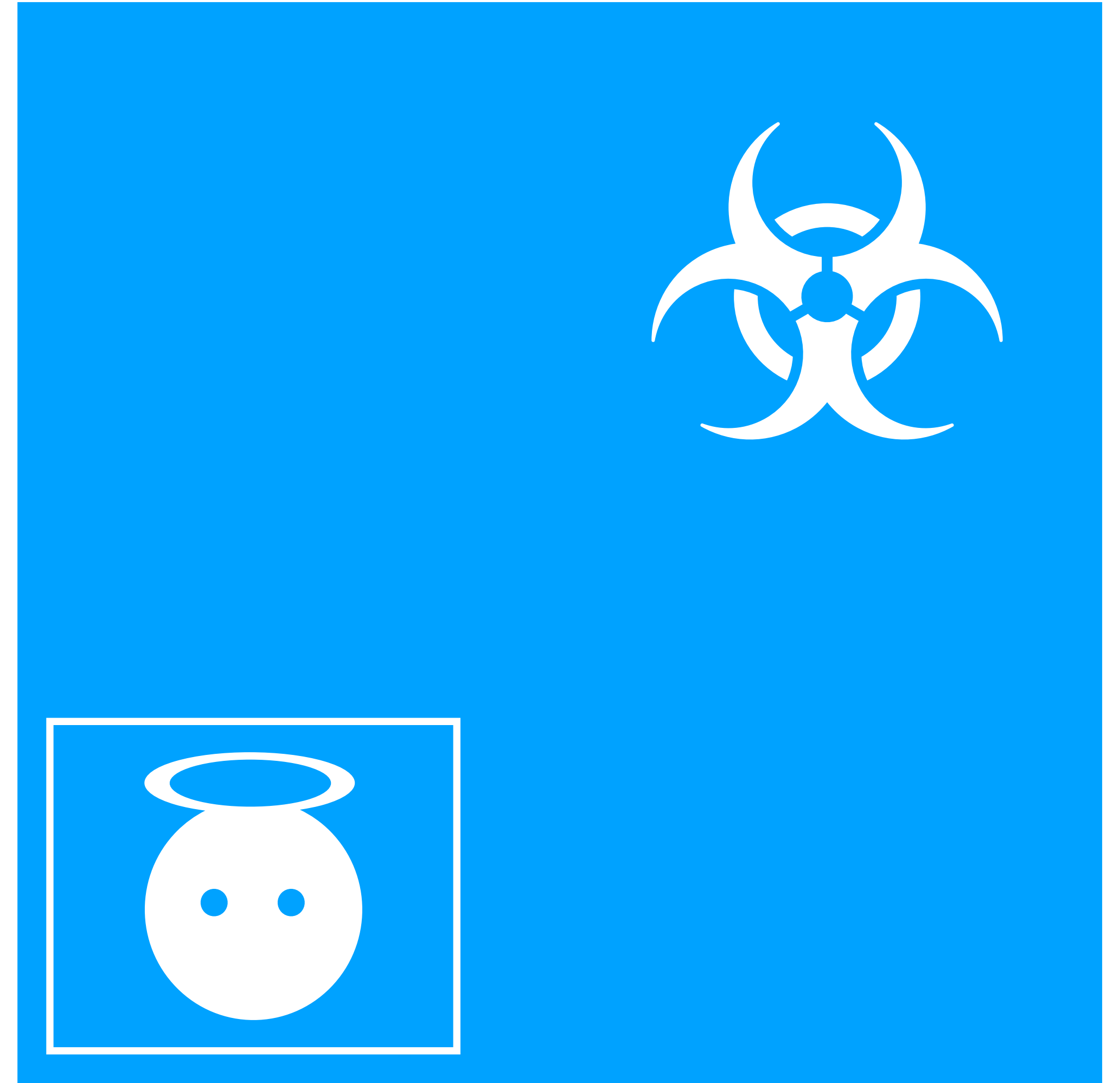
# Attack Overview

Given a trained target model

First pick feature with greatest 'benign' weight

Find a corresponding organ from the ice box

Wrap the organ in an opaque predicate

# Attack Overview

Given a trained target model

First pick feature with greatest 'benign' weight

Find a corresponding organ from the ice box

Wrap the organ in an opaque predicate

# Attack Overview

Given a trained target model

First pick feature with greatest 'benign' weight

Find a corresponding organ from the ice box

Wrap the organ in an opaque predicate
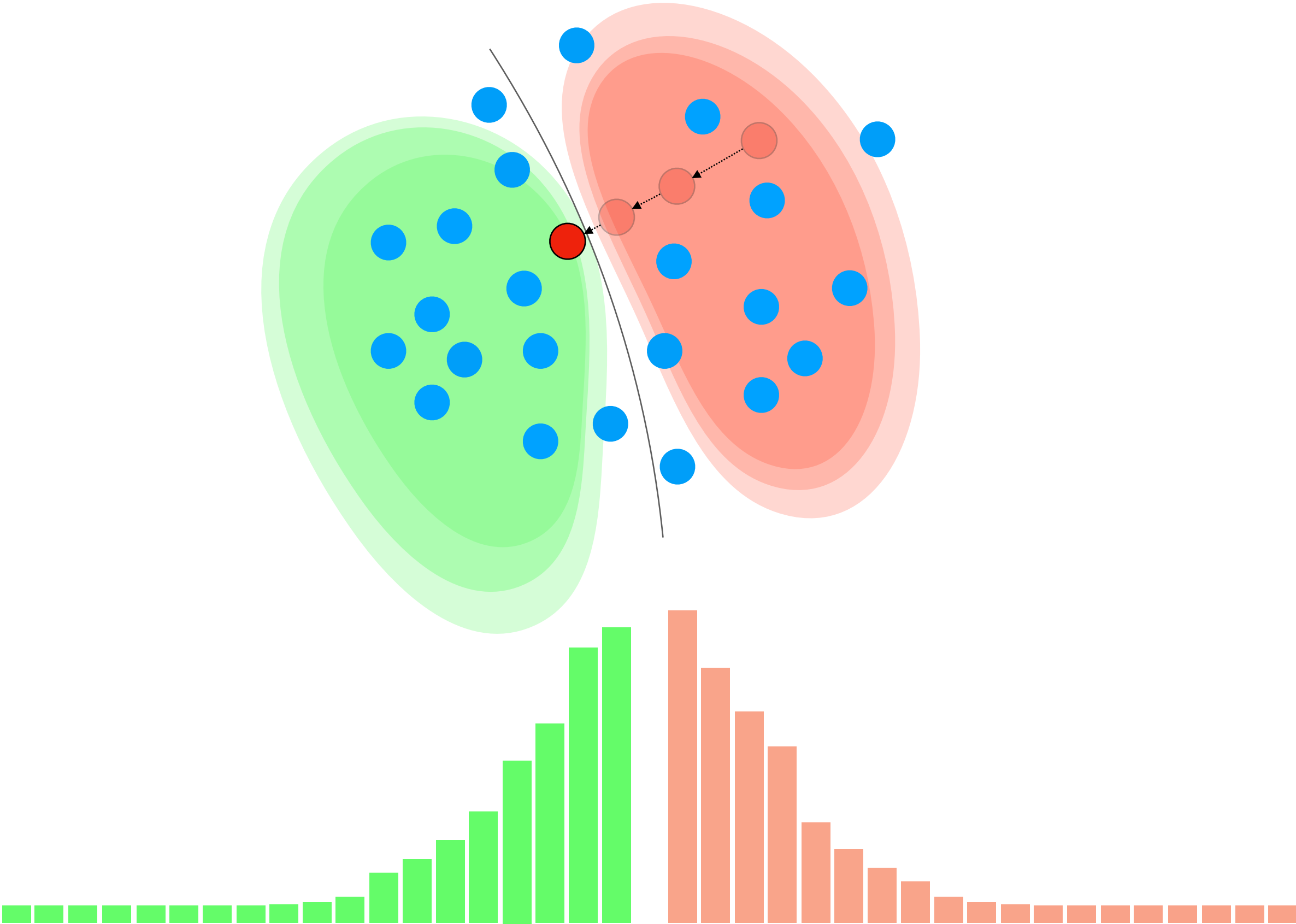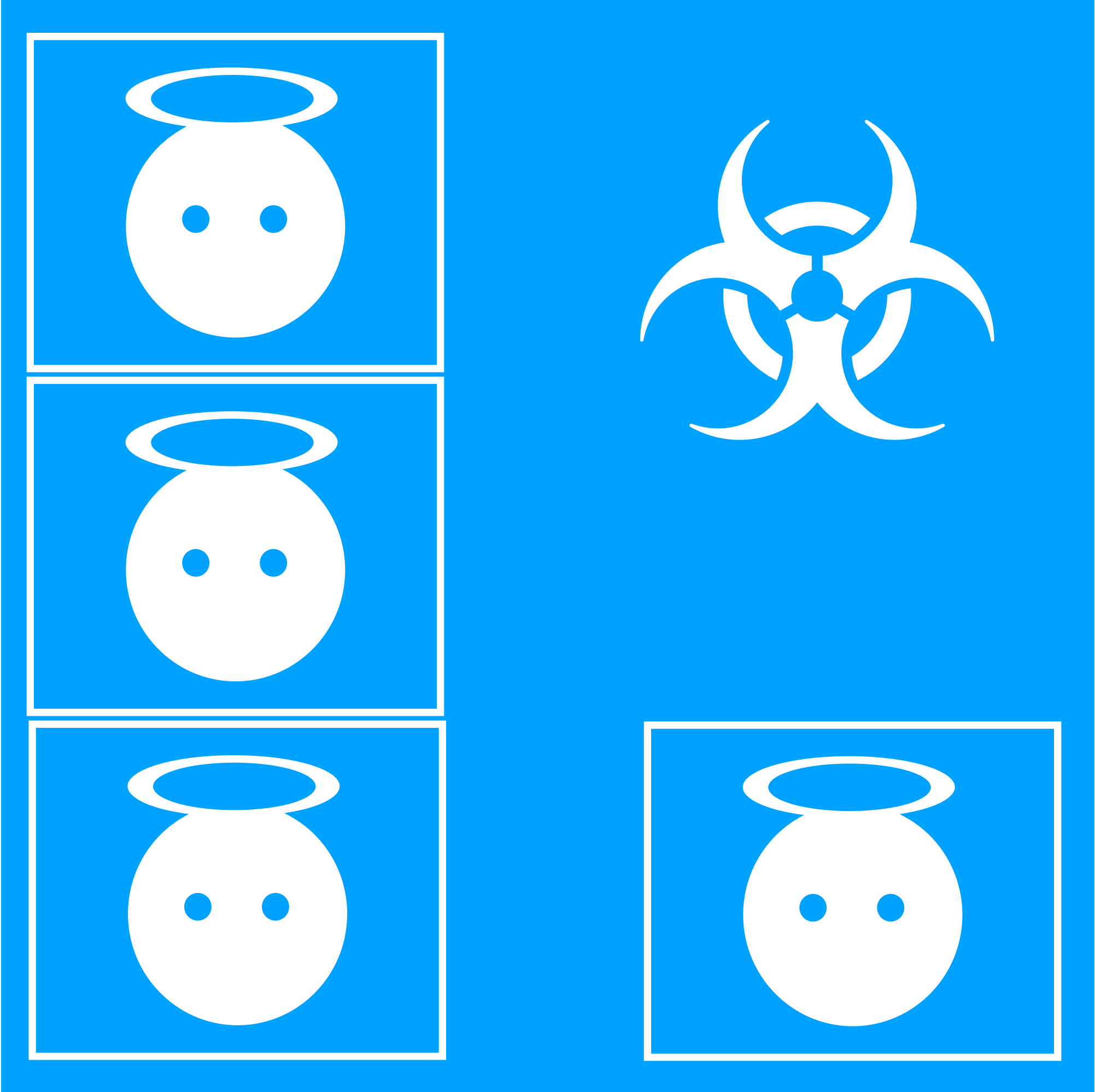
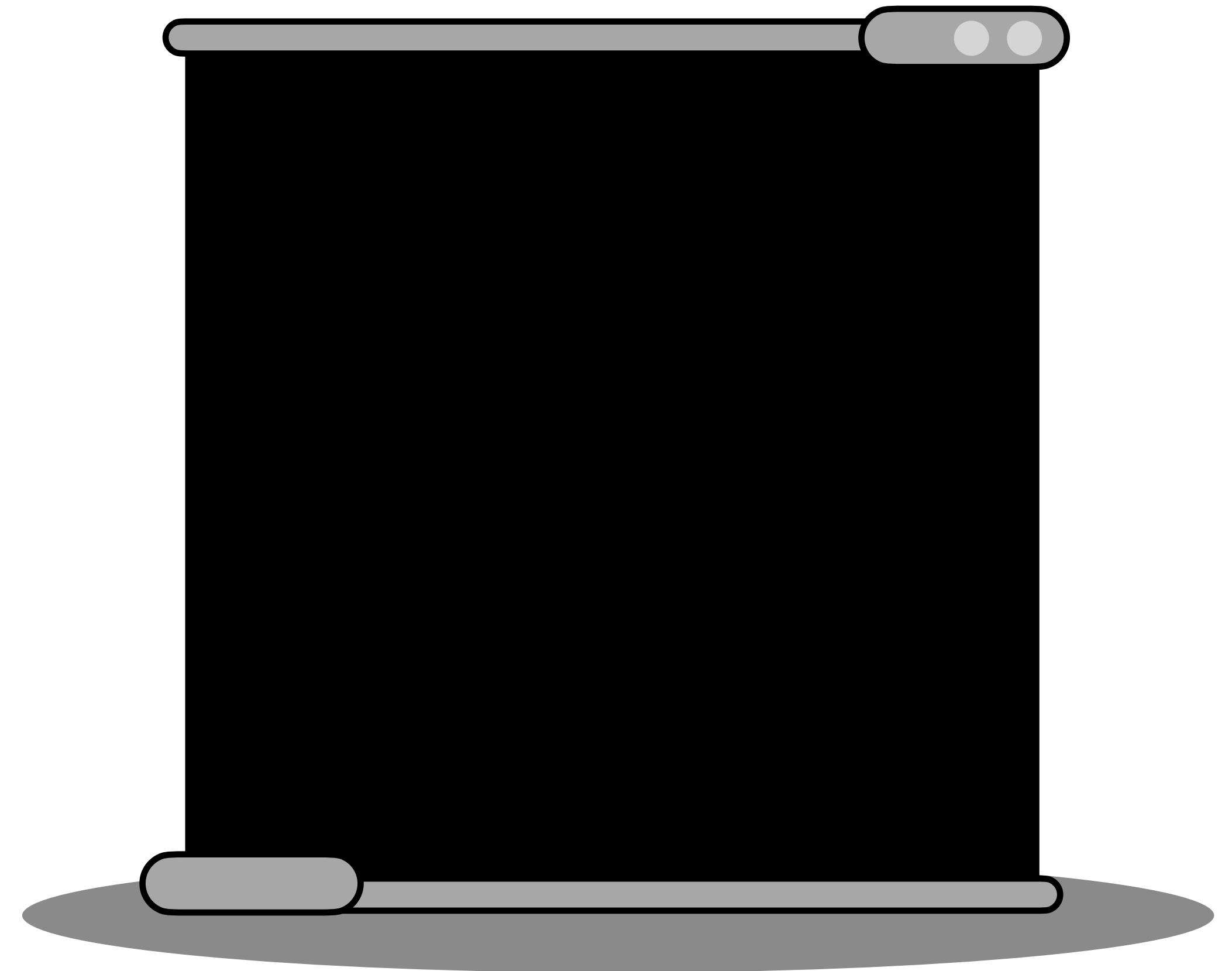Inject the new benign code and repackage

# Attack Overview

Given a trained target model

First pick feature with greatest 'benign' weight

Find a corresponding organ from the ice box

Wrap the organ in an opaque predicate

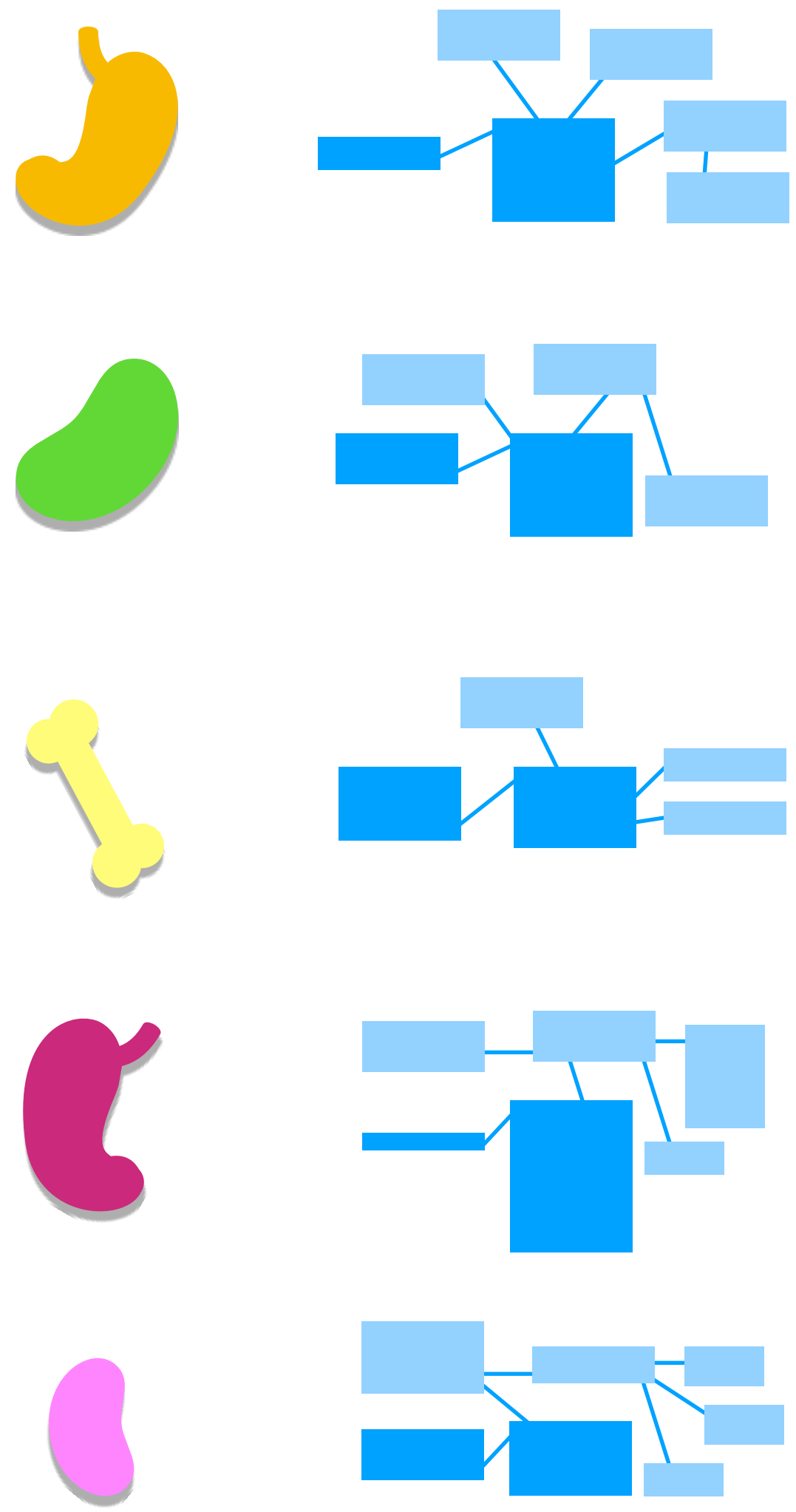Inject the new benign code and repackage

# Attack Overview

# Attack Overview

Continue choosing benign features until the app is misclassified

**[IEEE S&P 2020]** Intriguing Properties of Adversarial ML Attacks in the Problem Space
https://s2lab.cs.ucl.ac.uk/projects/intriguing

# Side-Effects

# Side-Effects

**[IEEE S&P 2020]** Intriguing Properties of Adversarial ML Attacks in the Problem Space
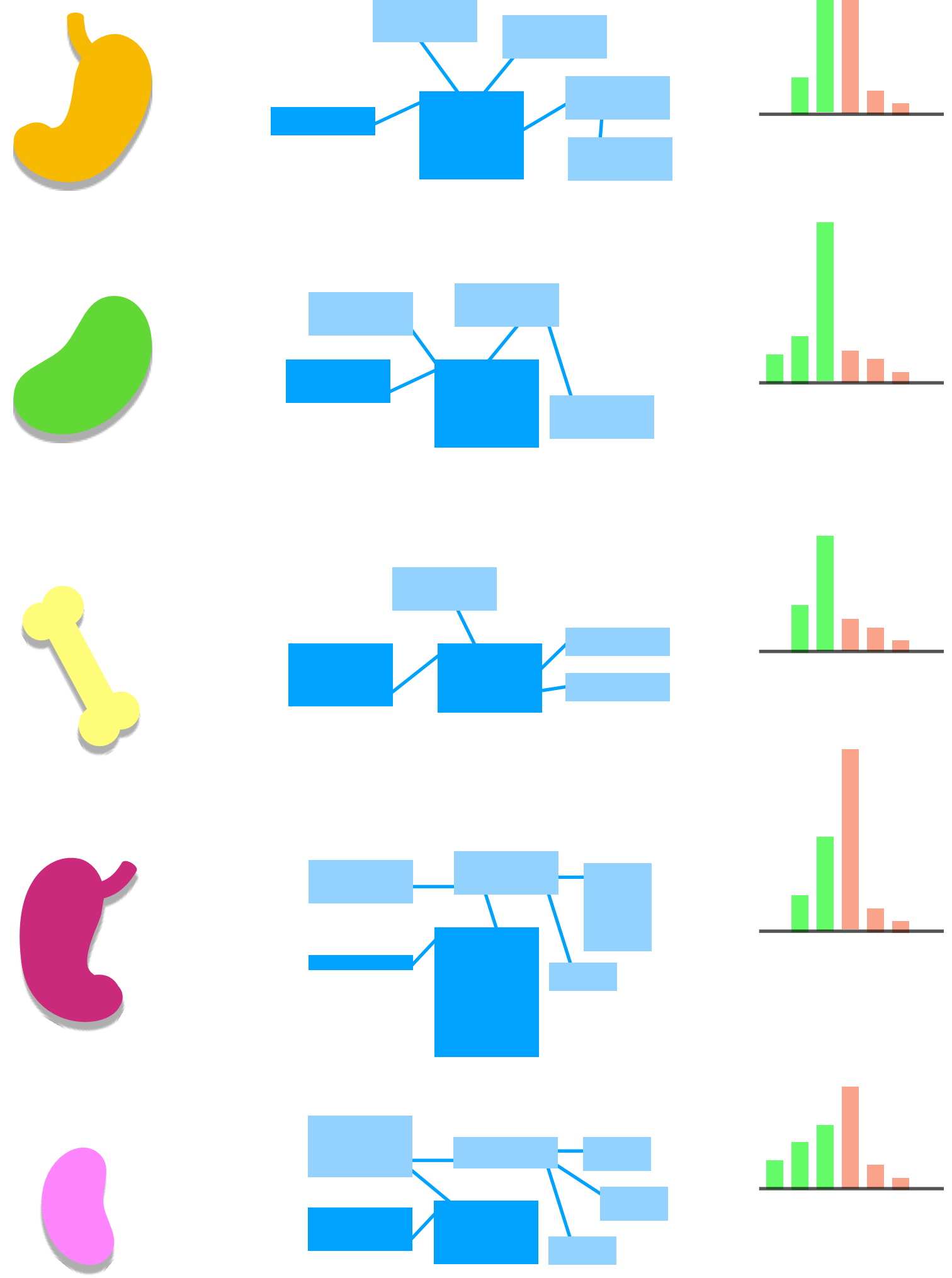https://s2lab.cs.ucl.ac.uk/projects/intriguing

# Side-Effects

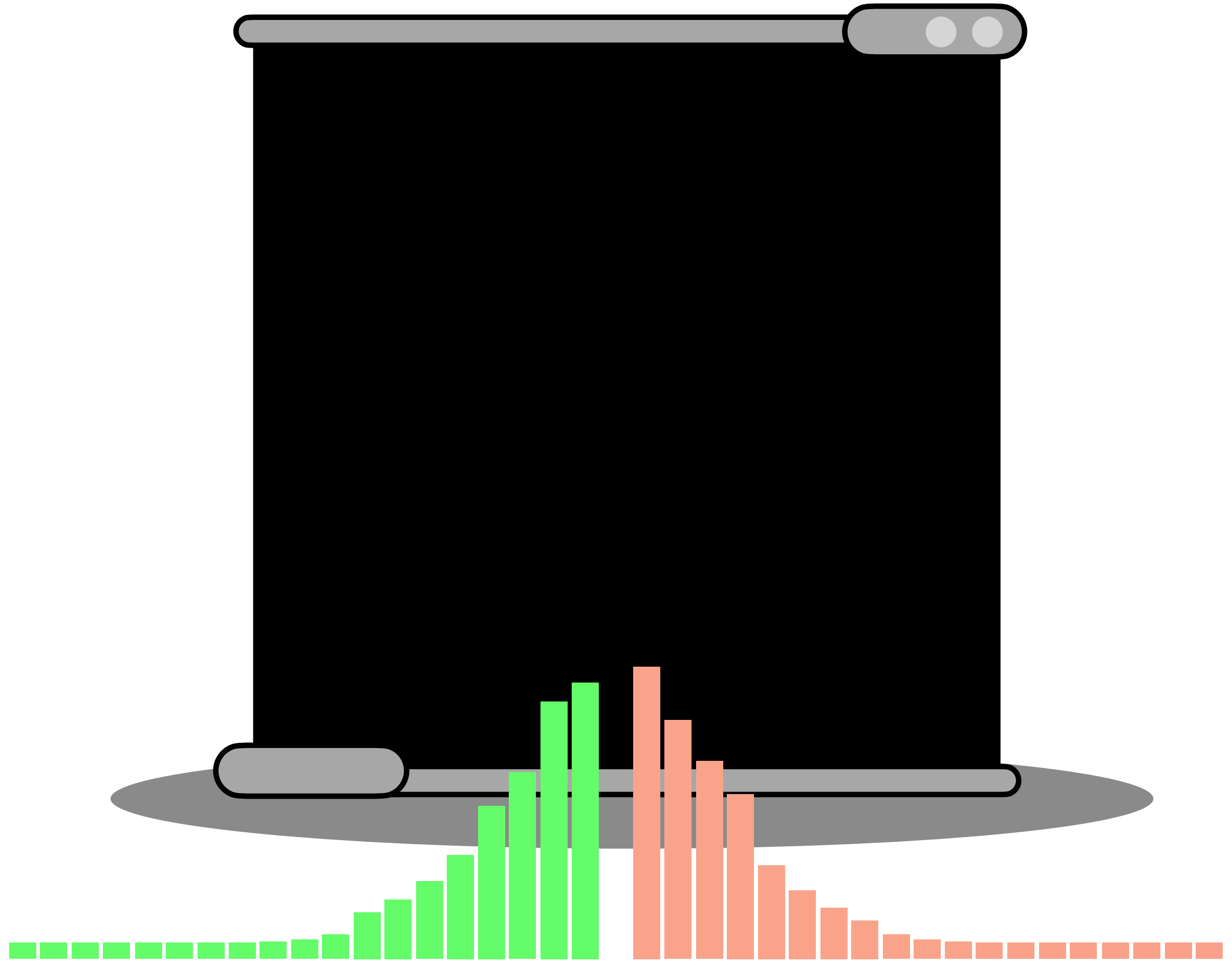Each organ contains side-effect features.
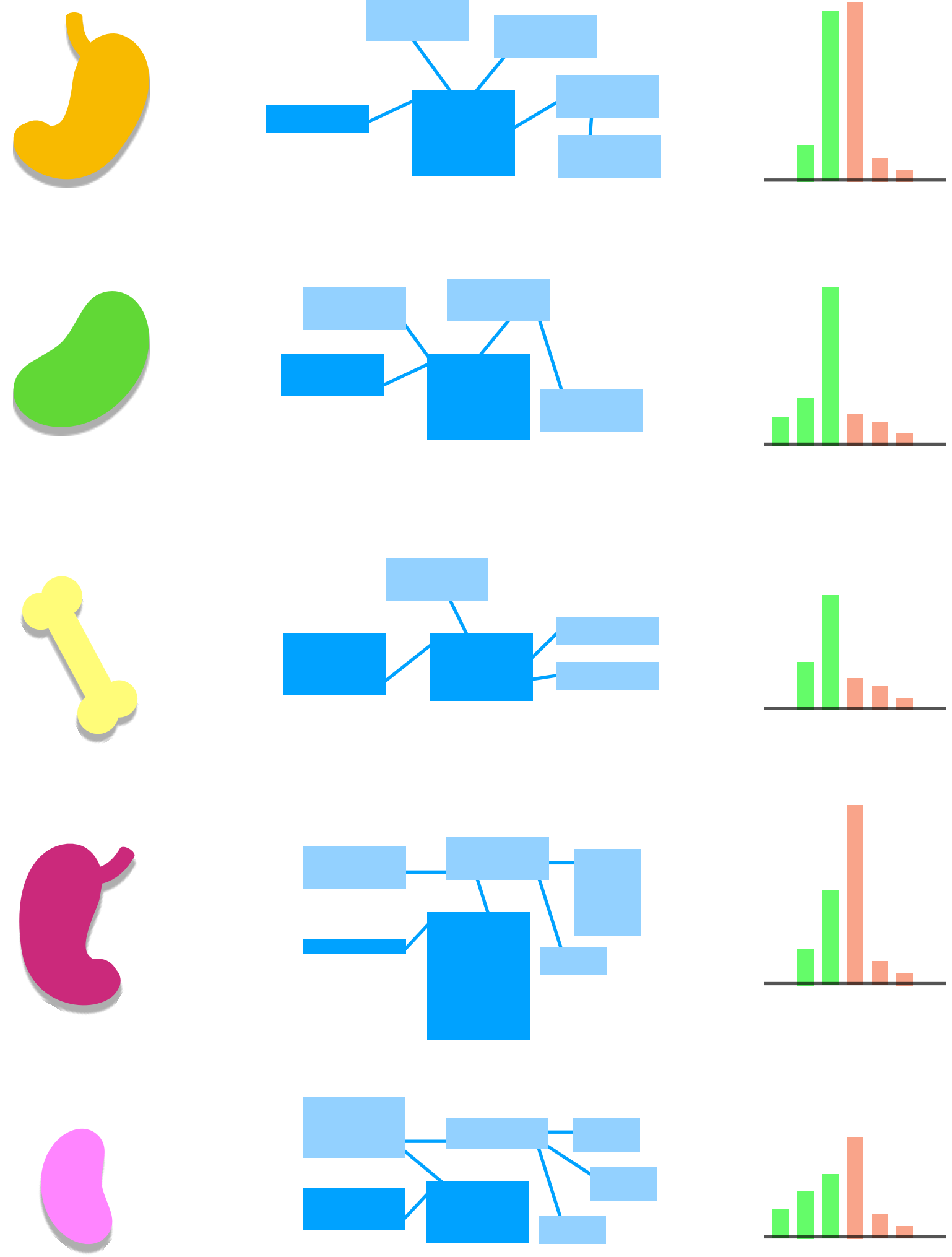
# Side-Effects

Each organ contains side-effect features.

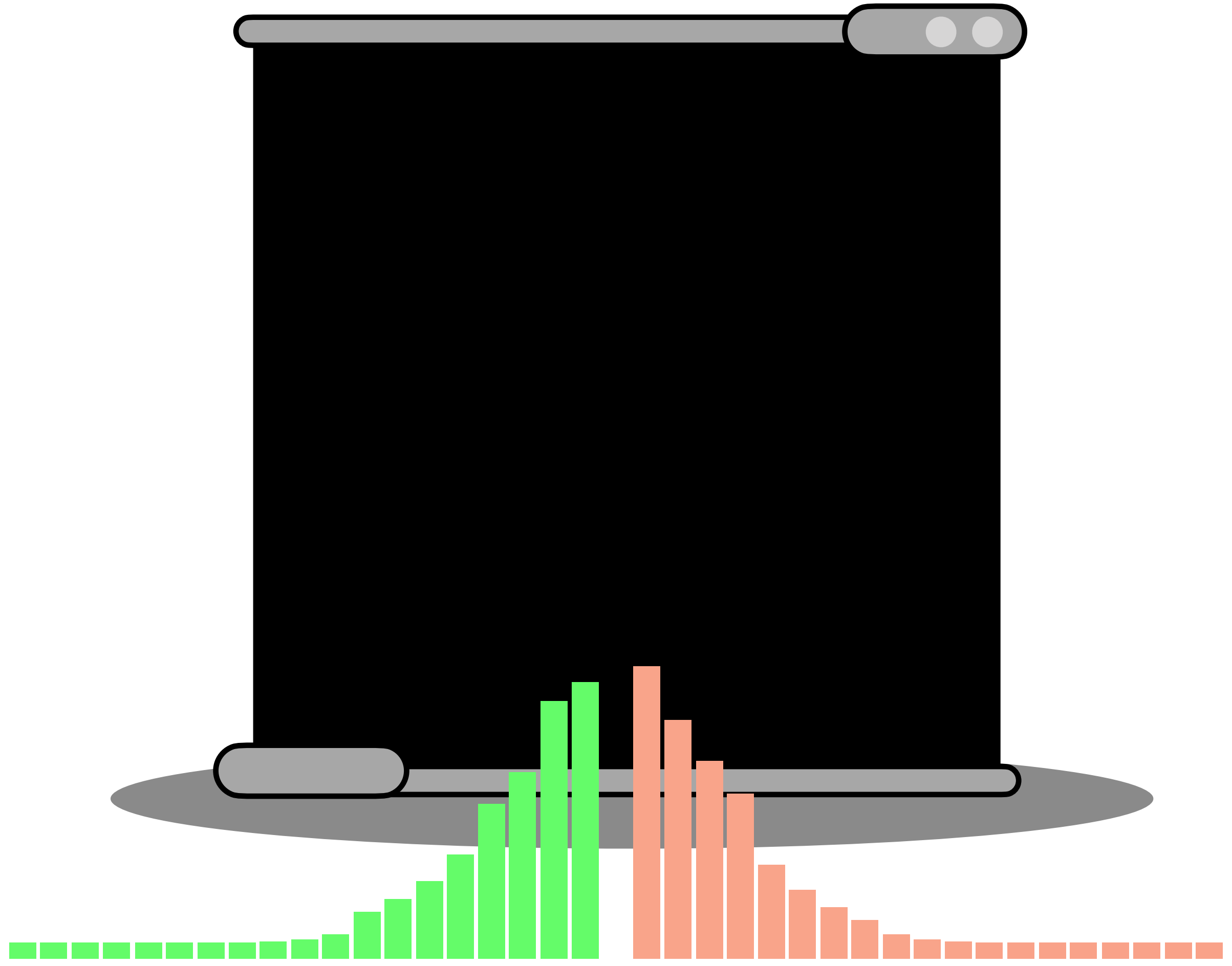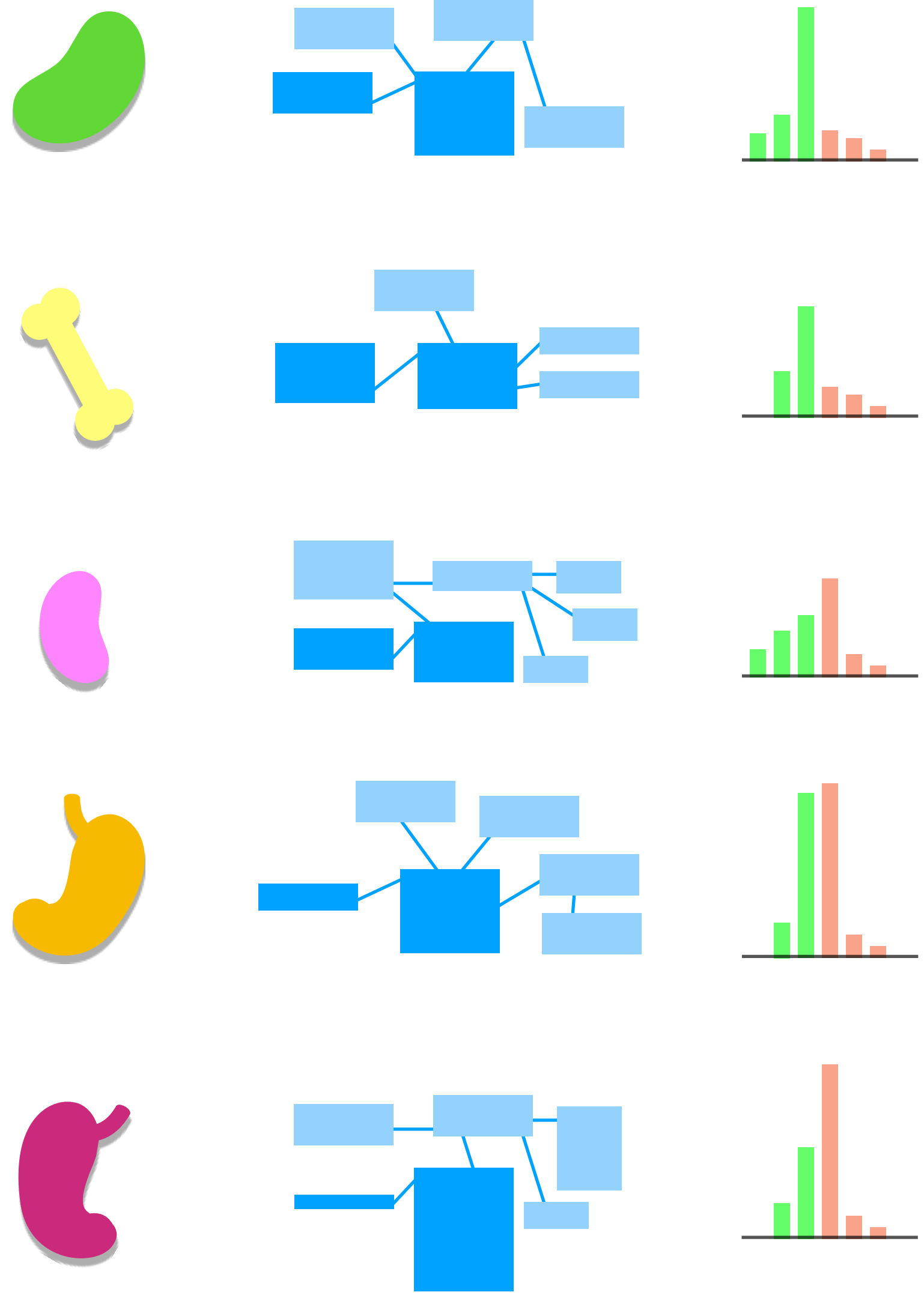We can sum target features, positive, and negative side effects

# Side-Effects

Each organ contains side-effect features.

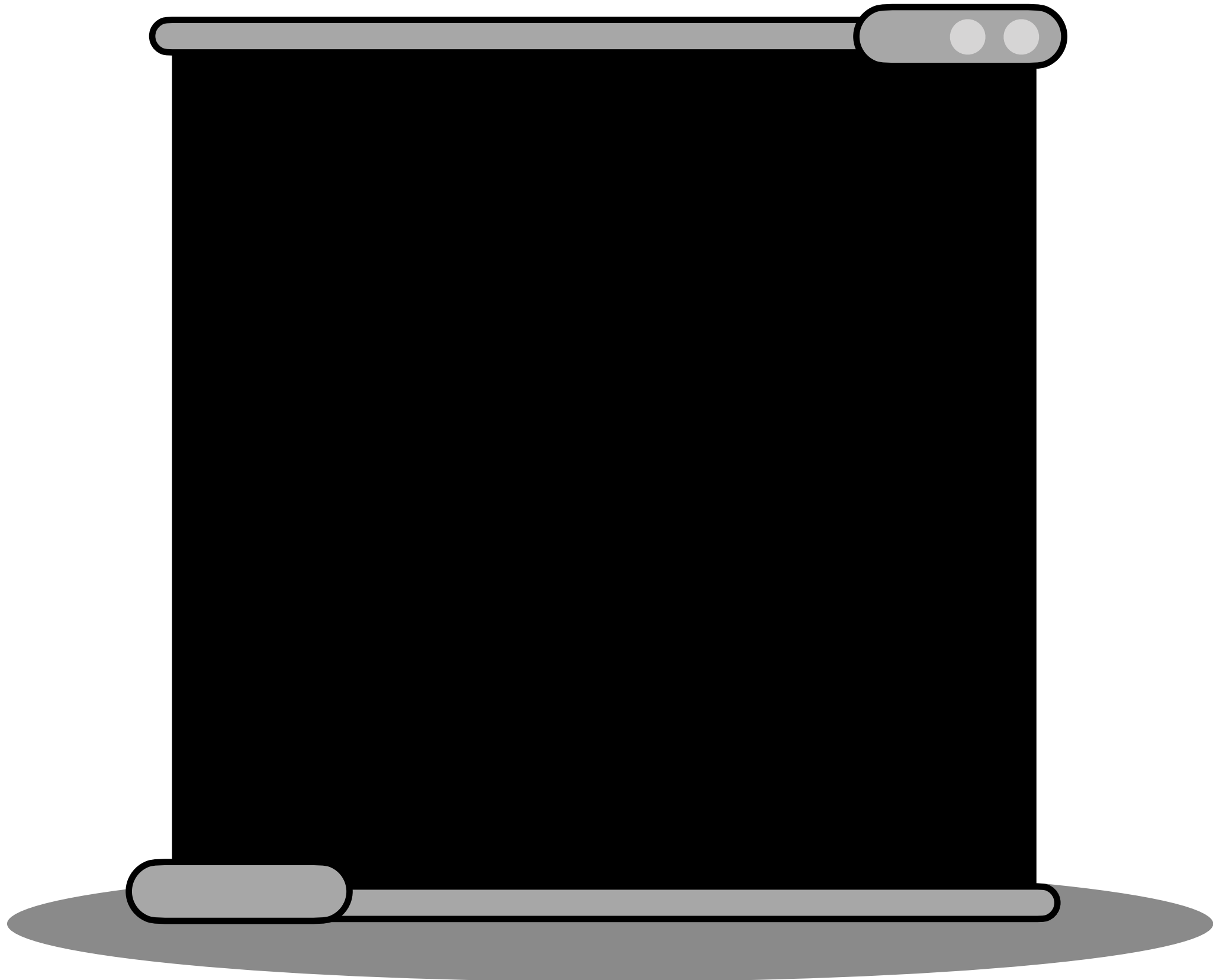We can sum target features, positive, and negative side effects

# Side-Effects

Each organ contains side-effect features.

We can sum target features, positive, and negative side effects

to choose organs in order of their overall benign weight

**[IEEE S&P 2020]** Intriguing Properties of Adversarial ML Attacks in the Problem Space
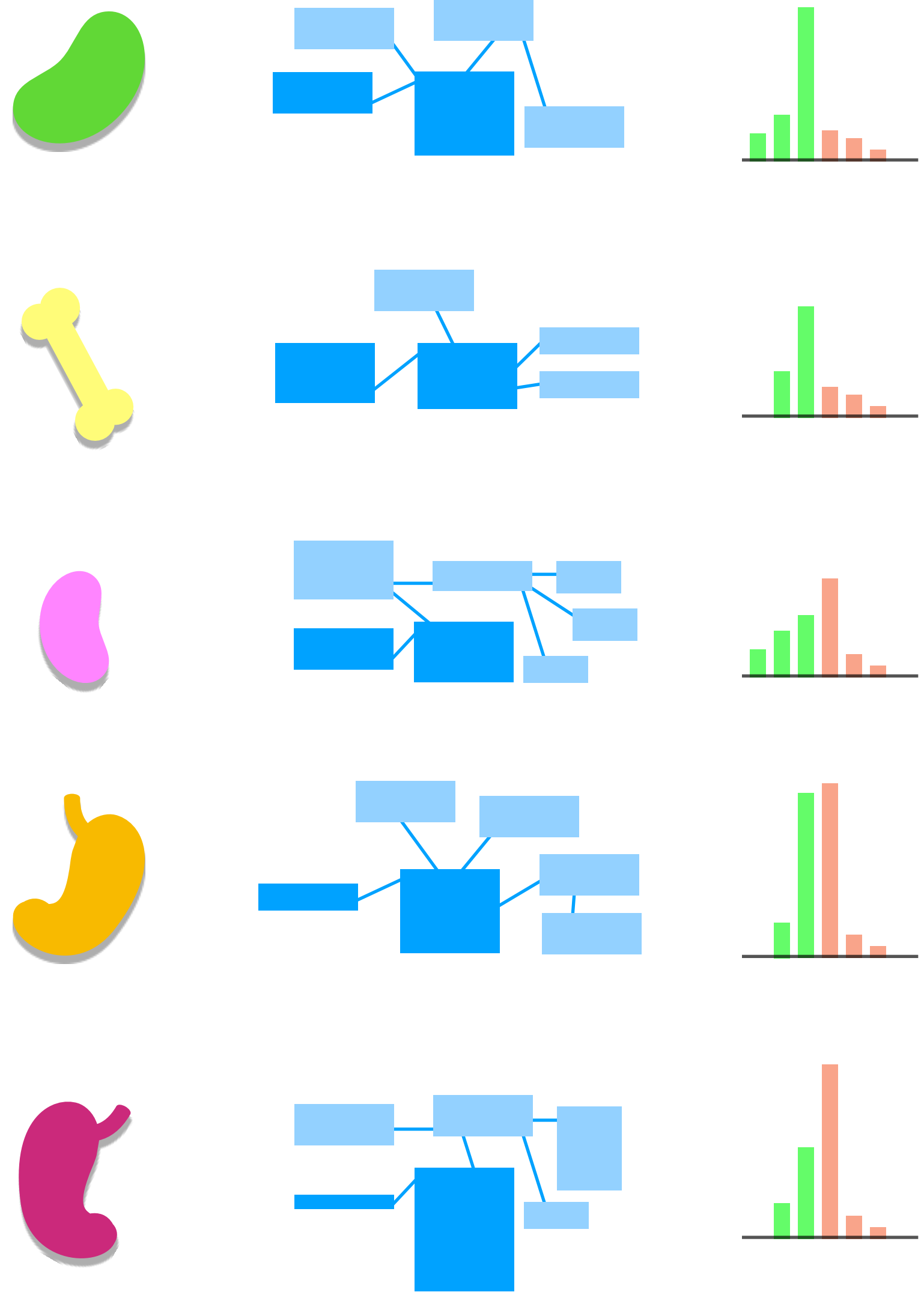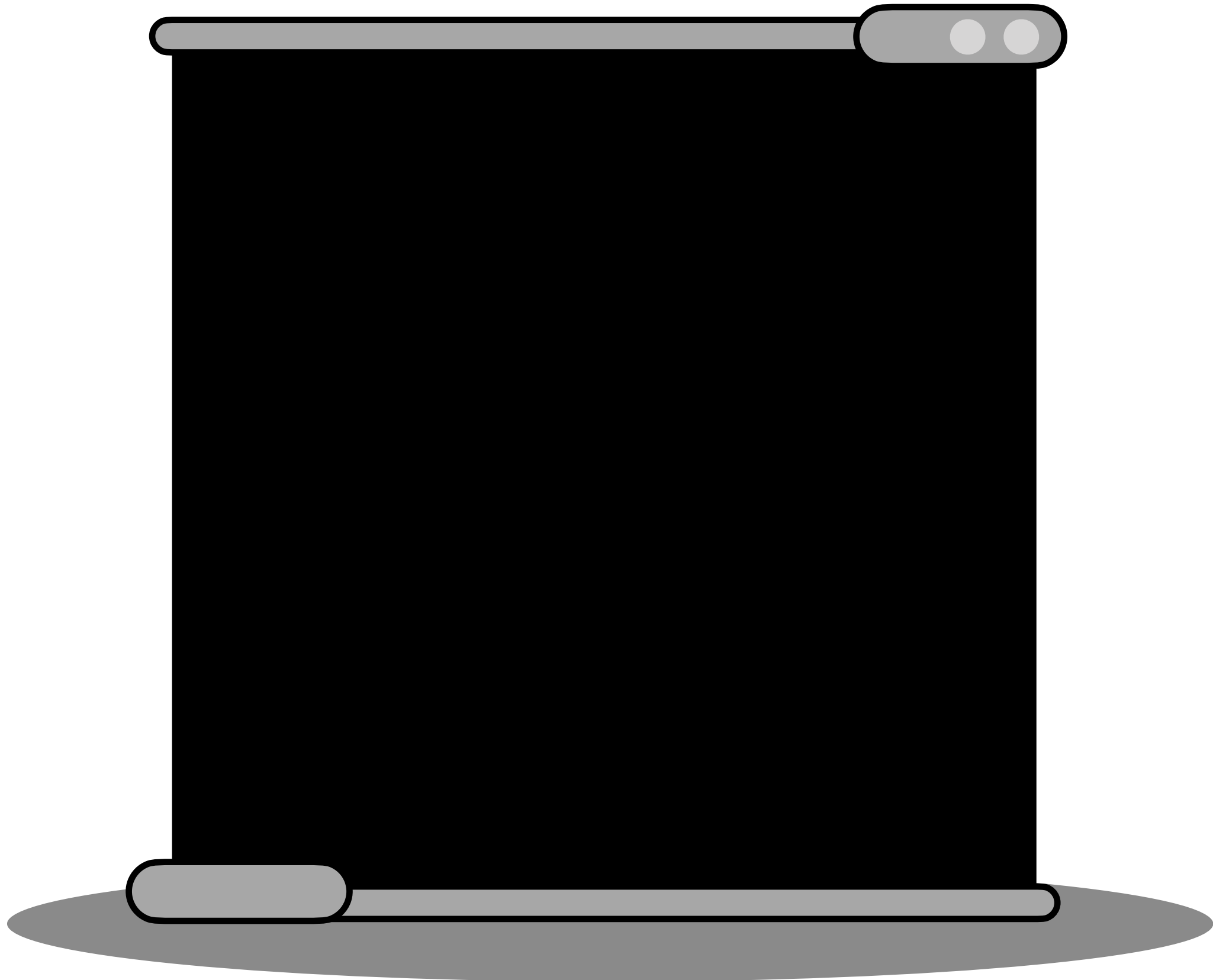https://s2lab.cs.ucl.ac.uk/projects/intriguing

# Side-Effects



Each organ contains side-effect features.

We can sum target features, positive, and negative side effects

to choose organs in order of their overall benign weight

# Side-Effects

Each organ contains side-effect features.

We can sum target features, positive, and negative side effects

to choose organs in order of their overall benign weight