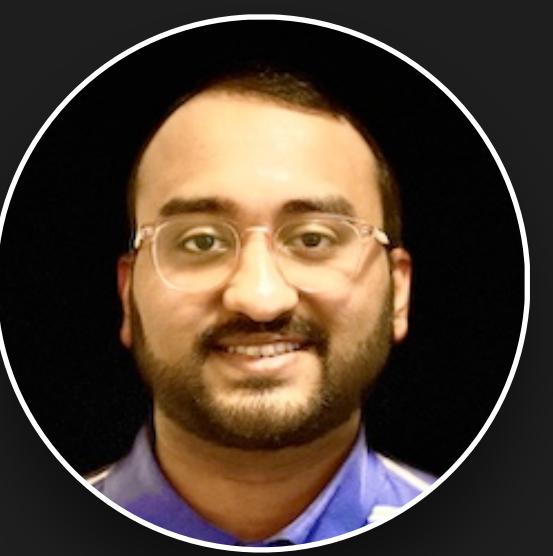


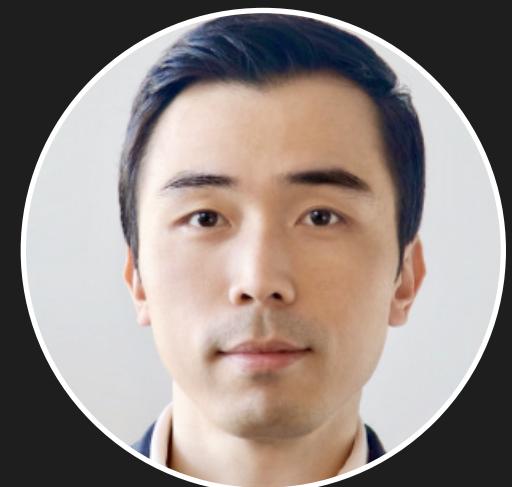
Understanding, Fortifying and Democratizing AI Security



Nilaksh Das
nilakshdas.com



Thesis Committee



PhD Advisor

Dr. Duen Horng Chau

Georgia Tech

Data Science & Adversarial ML



Dr. Srijan Kumar

Georgia Tech

Robust & Adversarial ML



Dr. Wenke Lee

Georgia Tech

Network Security & Privacy



Dr. Ponnurangam Kumaraguru ("PK")

IIT-Hyderabad

Privacy & Security in Online Social Media



Dr. Xu Chu

Georgia Tech

Data Management Systems & ML



Dr. Oliver Brdiczka

Adobe

AI & Human-Computer Interaction

AI is becoming ubiquitous

Smart voice assistants

Biometric authentication

Social media newsfeeds

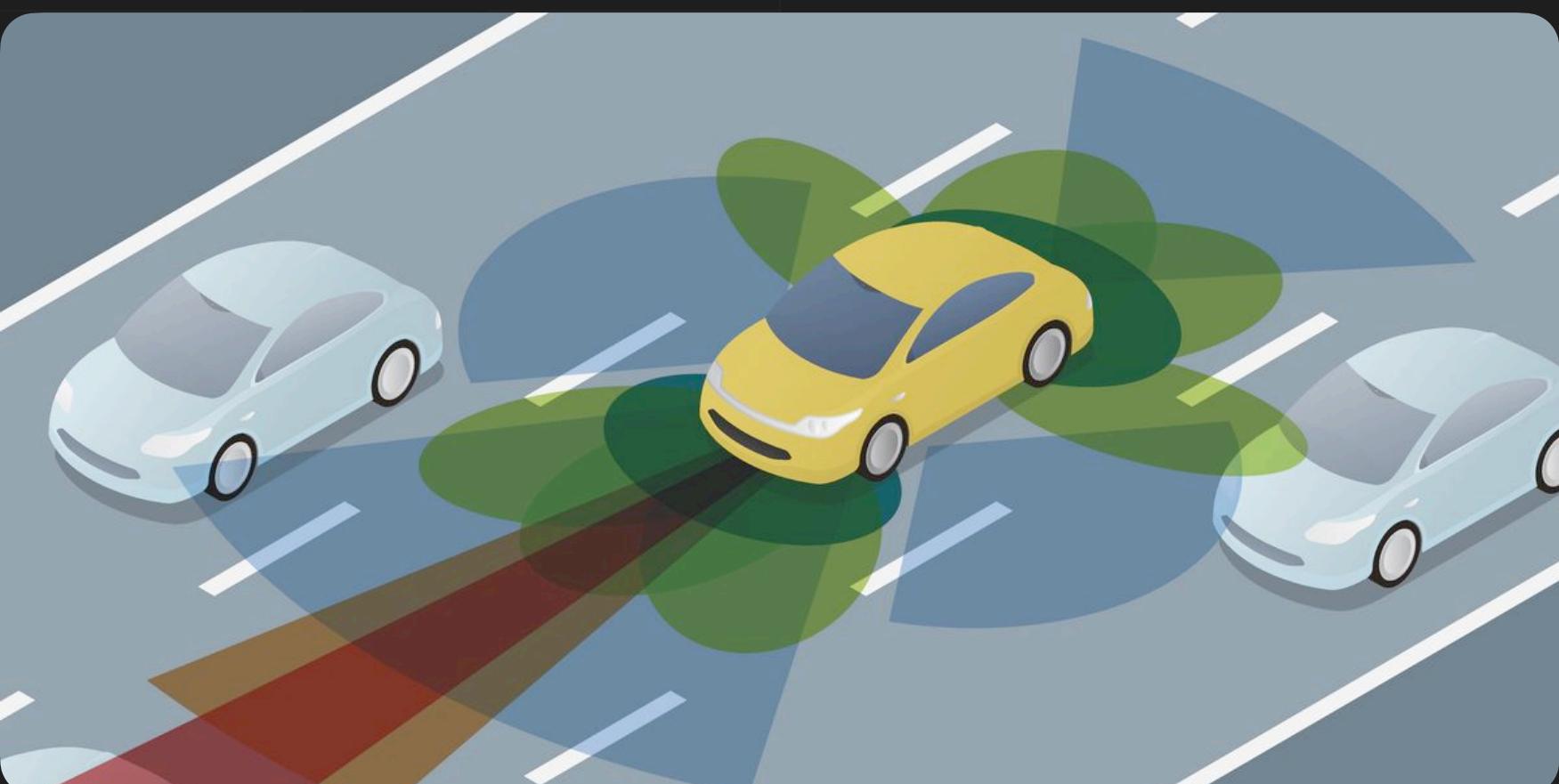
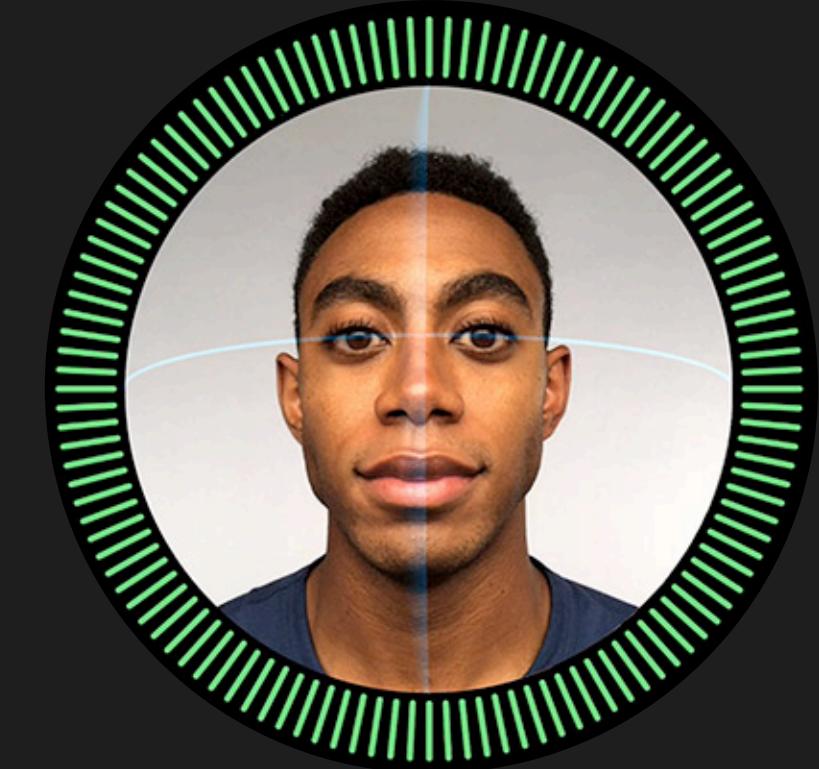
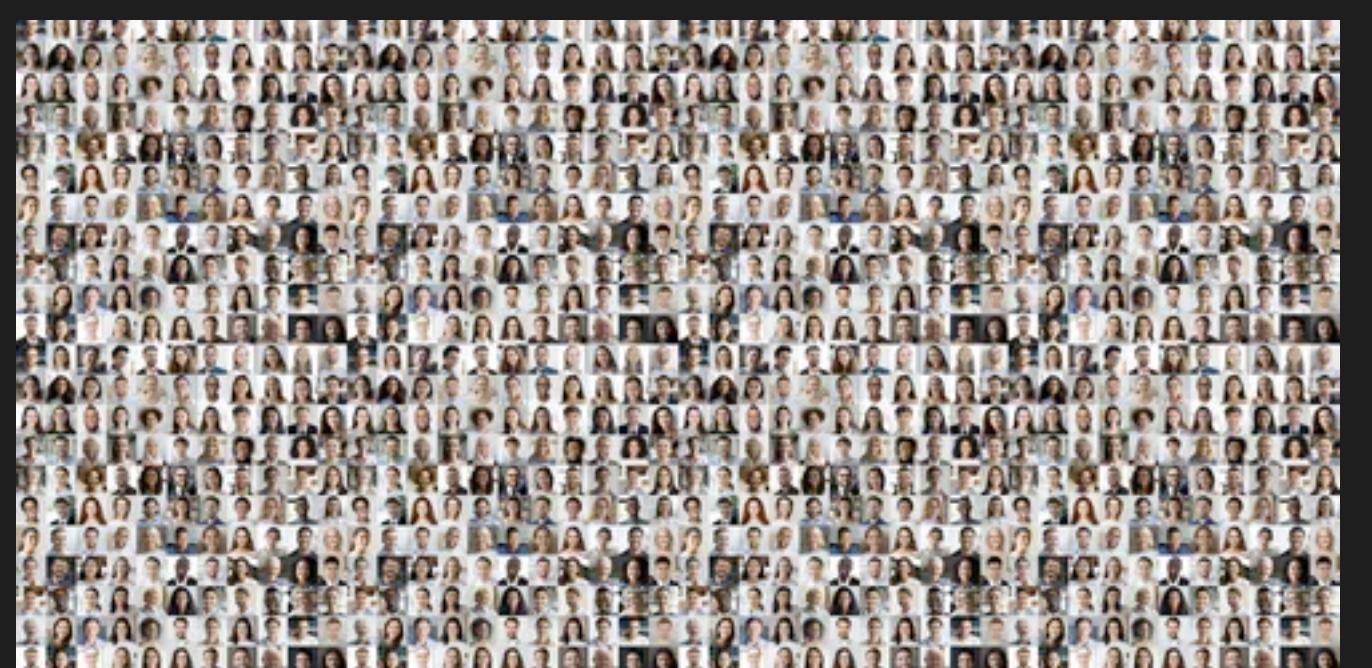
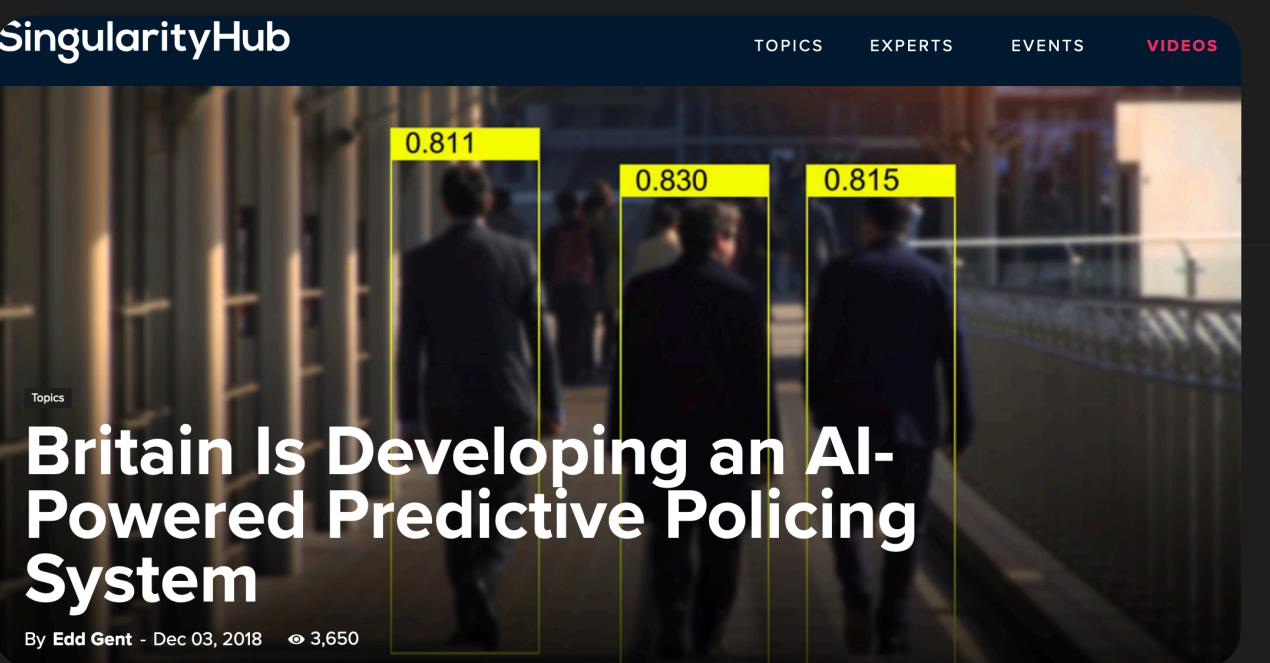
Autonomous vehicles

High-frequency trading

Healthcare

Public safety and policing

Massive data collection



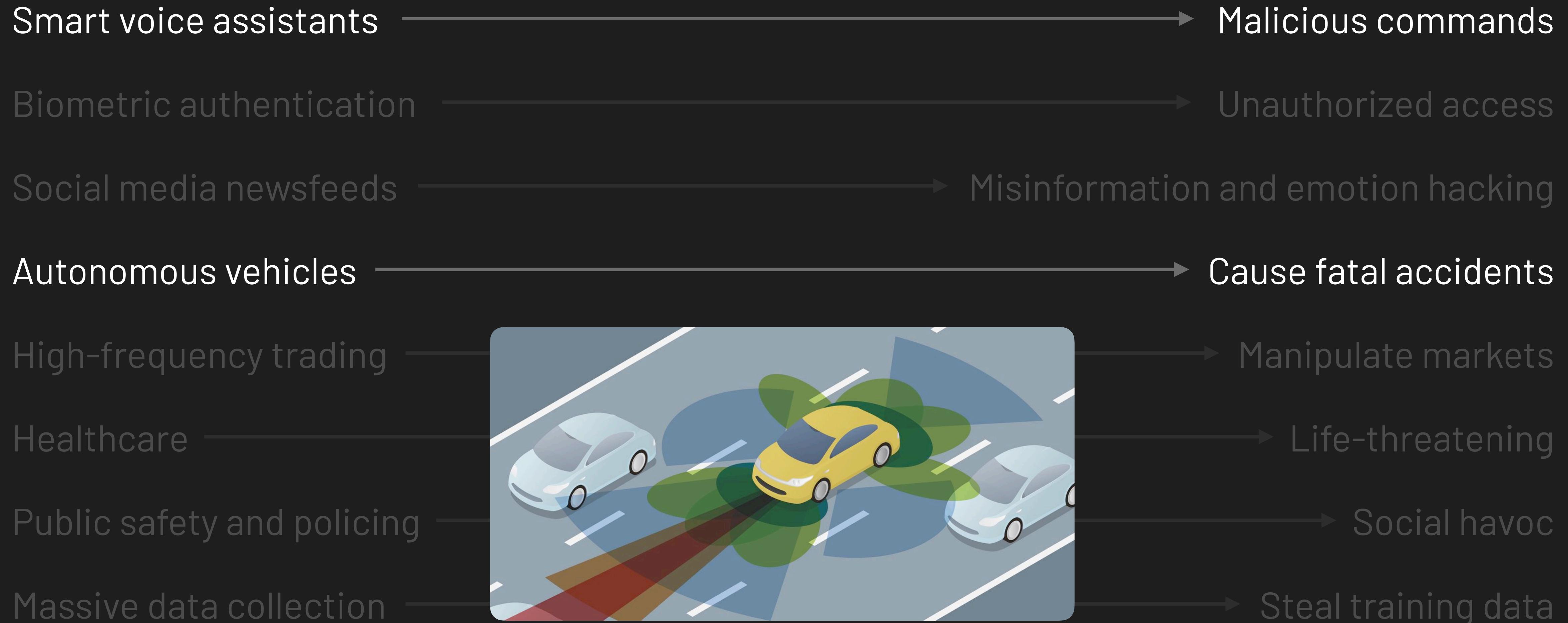
AI Security is becoming critical



AI Security is becoming critical

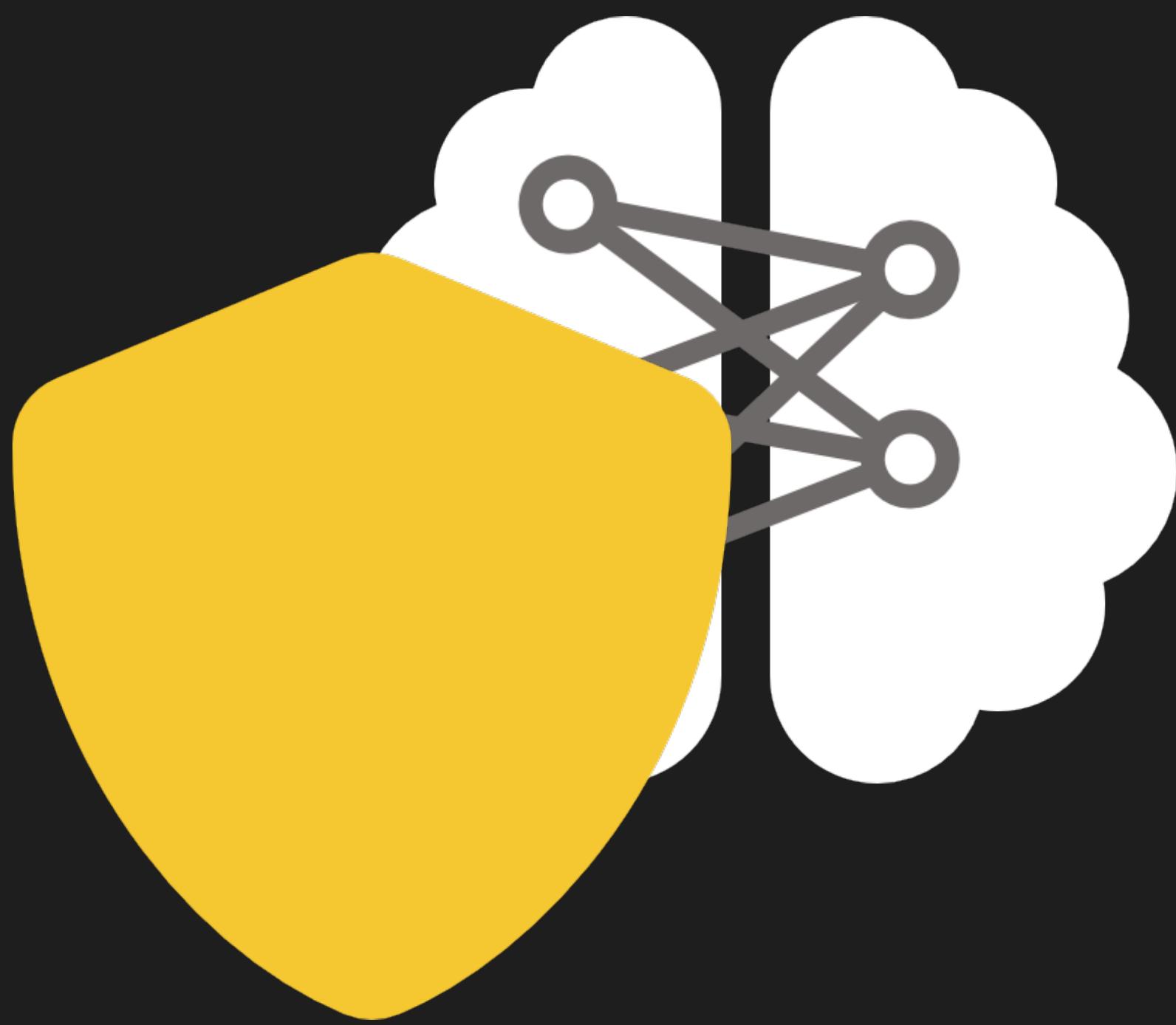


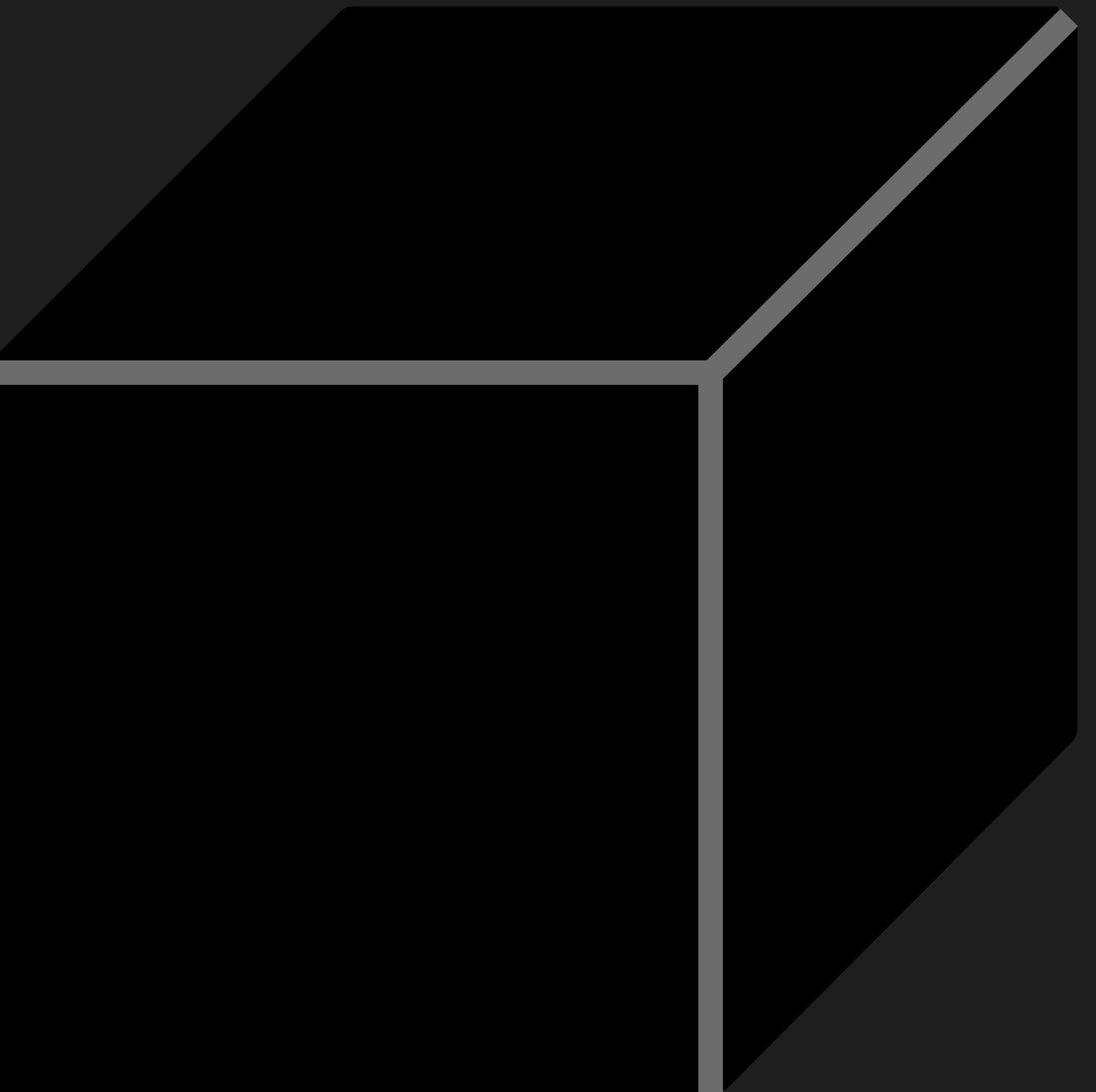
AI Security is becoming critical

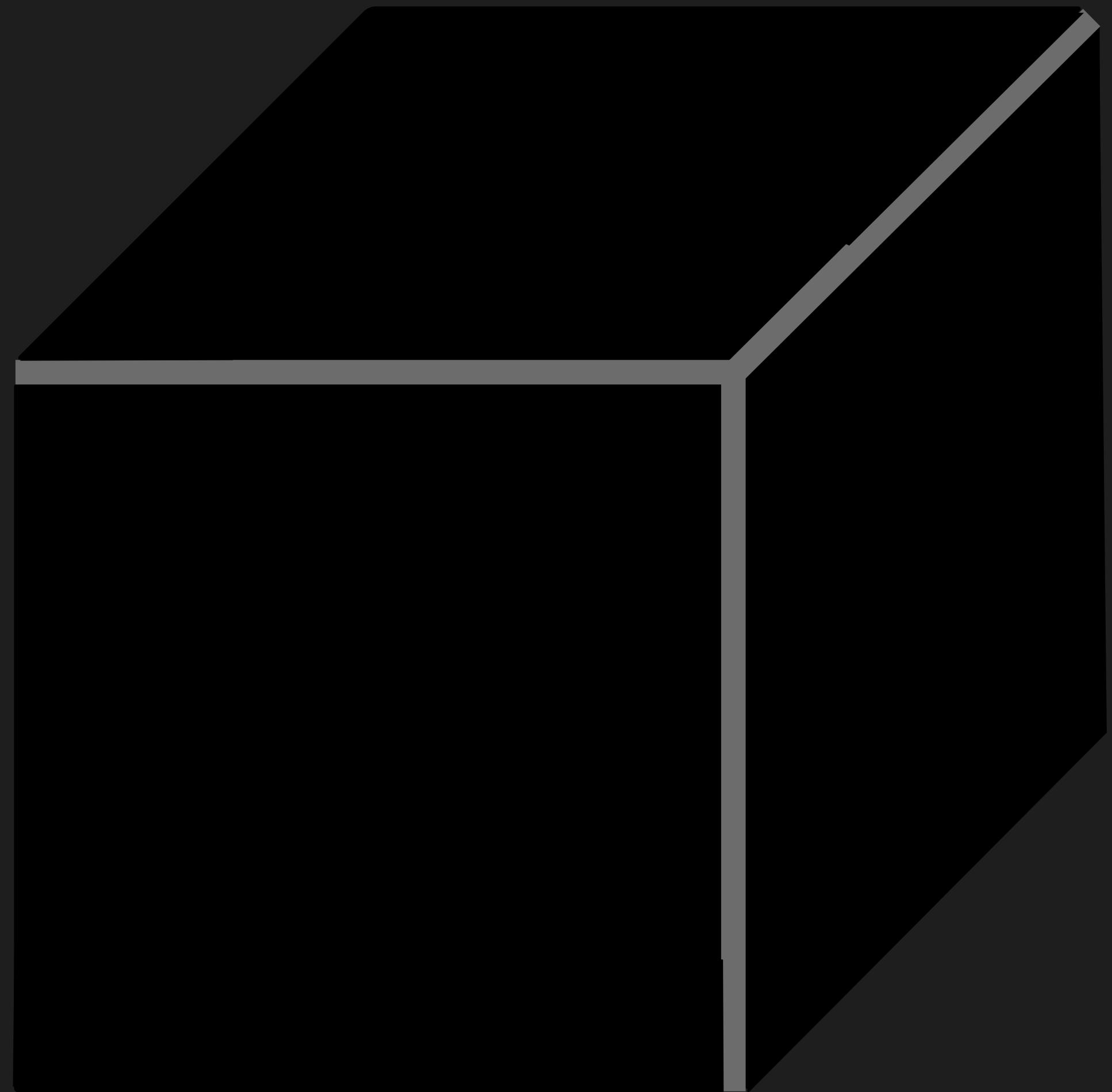


AI Security is becoming critical















Research Mission

Enable everyone to understand and fortify AI security

Research Mission

Enable everyone to understand and fortify AI security

Research Challenge I

How to **UNDERSTAND** AI vulnerabilities?

How can we intuitively interpret the core reason for AI vulnerabilities?

Research Challenge II

How to **FORTIFY** AI security?

How can we effectively mitigate the harm from adversarial examples?

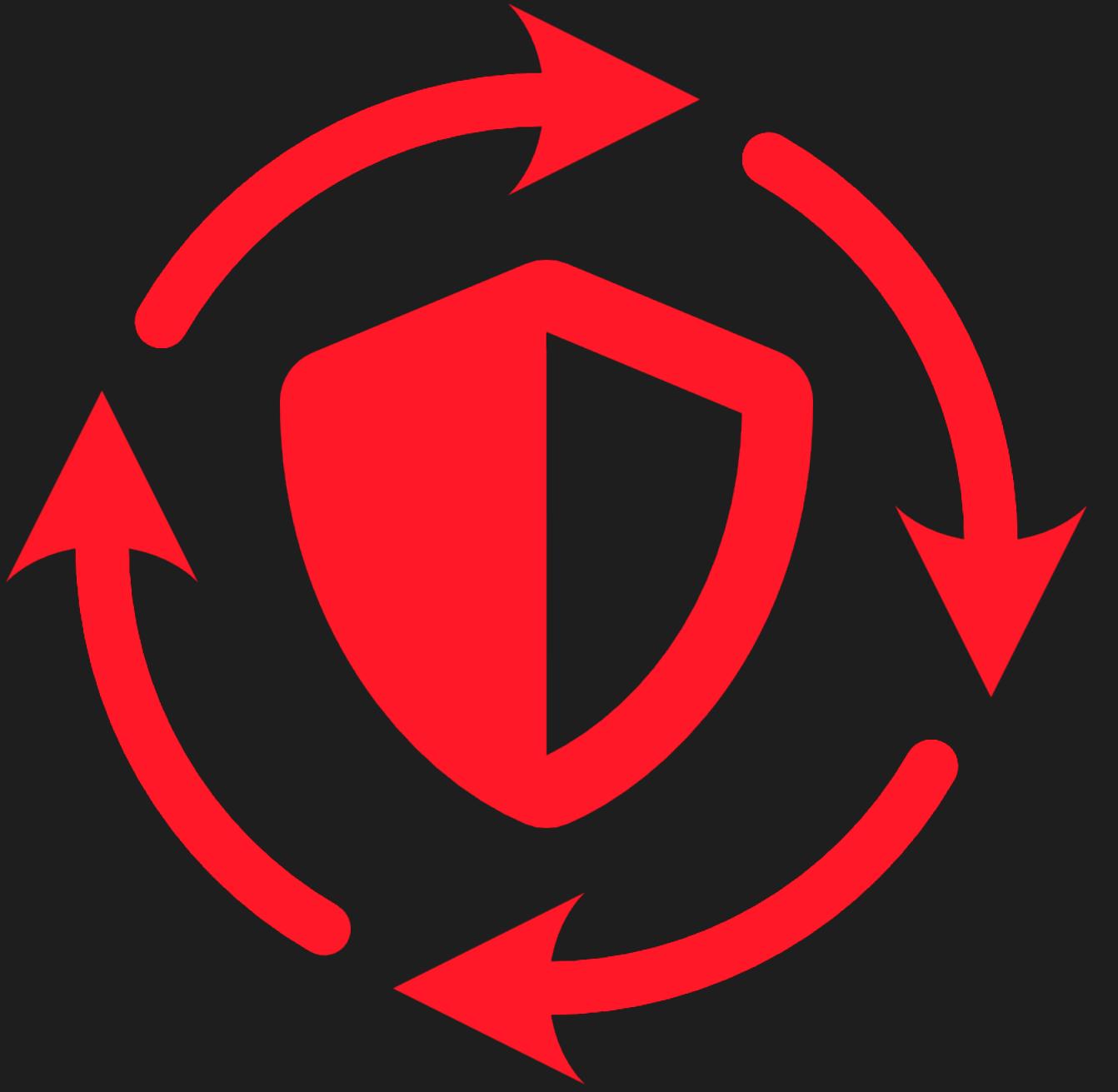
Research Challenge III

How to **ENABLE** everyone to test AI attacks & defenses?

How can we make AI security techniques more accessible?



UNDERSTAND



FORTIFY



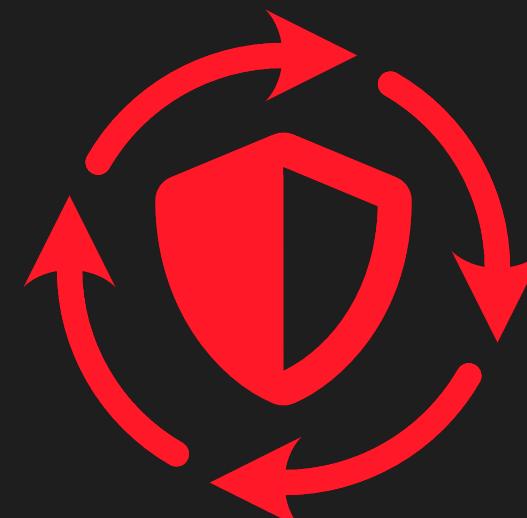
ENABLE

Research Thrust I



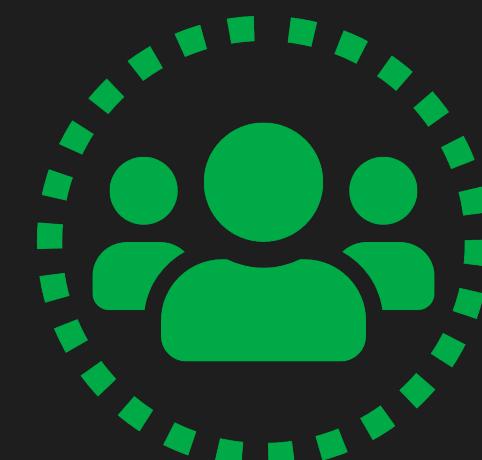
Exposing AI Vulnerabilities through Visualization & Interpretable Representations.

Research Thrust II

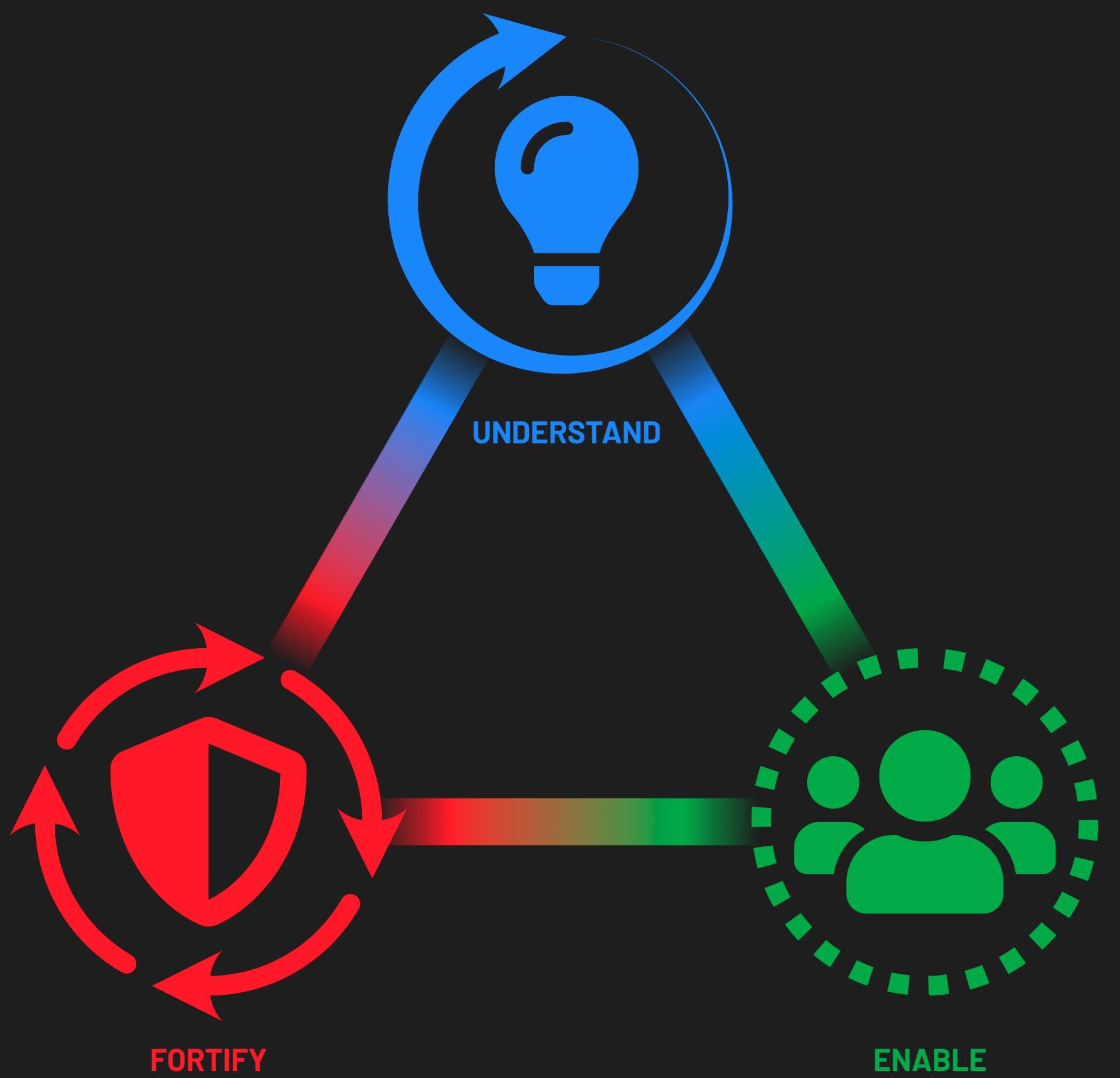


Mitigating Adversarial Examples Across Modalities & Tasks.

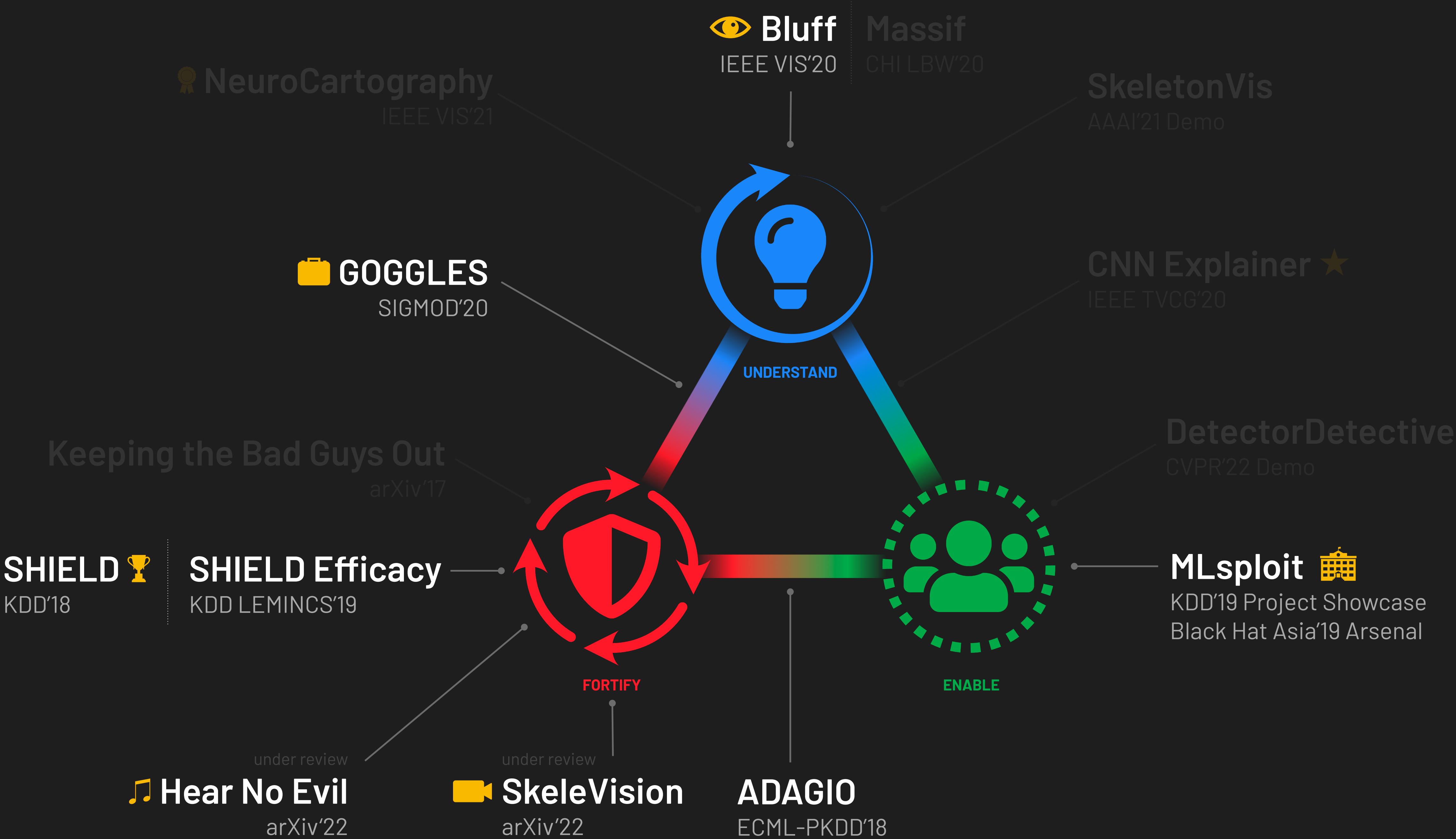
Research Thrust III



Democratizing AI Security Research & Pedagogy with Scalable Interactive Experimentation.







Thesis Statement

Overcome the threat of **adversarial attacks** by **enabling** everyone to:

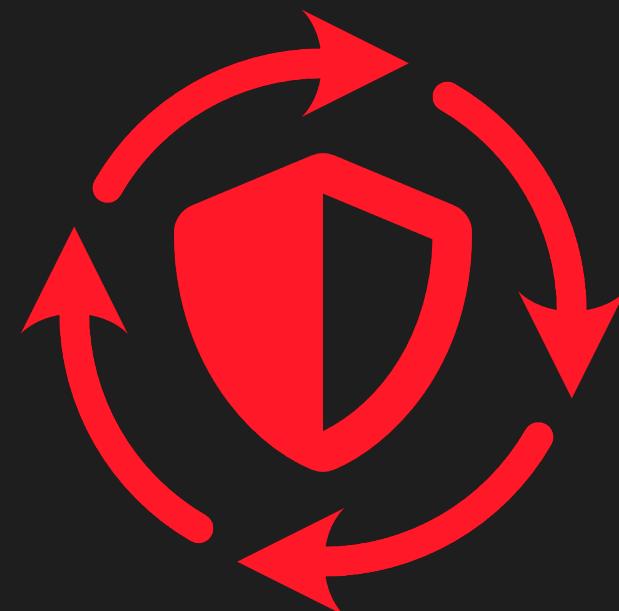
1. **understand** AI vulnerabilities through **intuitive interpretation**
2. **fortify** AI applications with **robust mitigation** of adversarial examples
3. **easily test** AI security techniques with **interactive experimentation**



Part I

Understand AI Vulnerabilities

GOGGLES SIGMOD 2020
Bluff IEEE VIS 2020



Part II

Fortify AI Security

SHIELD KDD 2018
SkeleVision arXiv 2022 (under review)
Hear No Evil arXiv 2022 (under review)



Part III

Enable Use of AI Security

ADAGIO ECML-PKDD 2018
MLsploit KDD Showcase 2019



Part I

Understand AI Vulnerabilities

GOGGLES SIGMOD 2020
Bluff IEEE VIS 2020



Part II

Fortify AI Security

SHIELD KDD 2018
SkeleVision arXiv 2022 (under review)
Hear No Evil arXiv 2022 (under review)

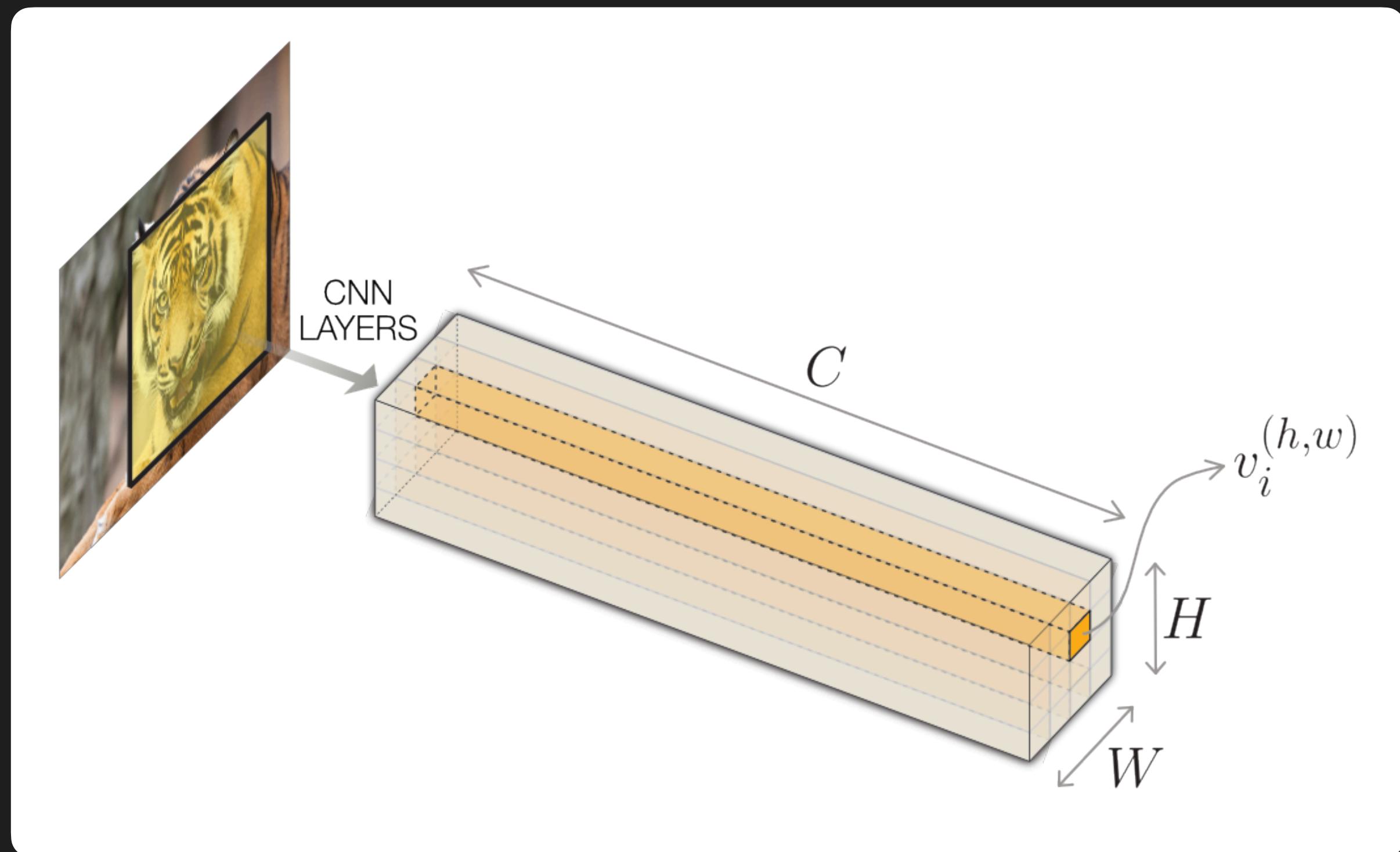


Part III

Enable Use of AI Security

ADAGIO ECML-PKDD 2018
MLsploit KDD Showcase 2019

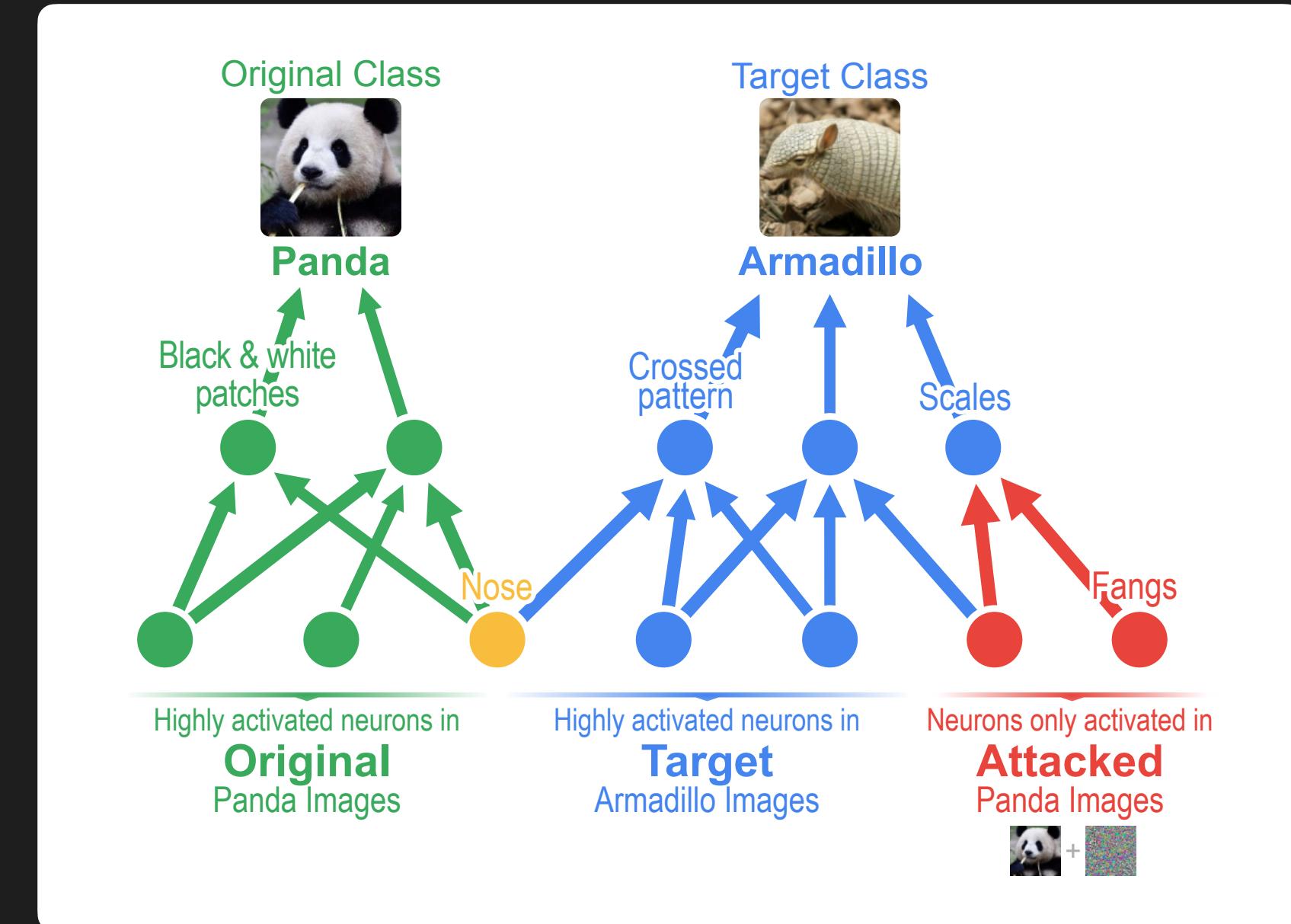
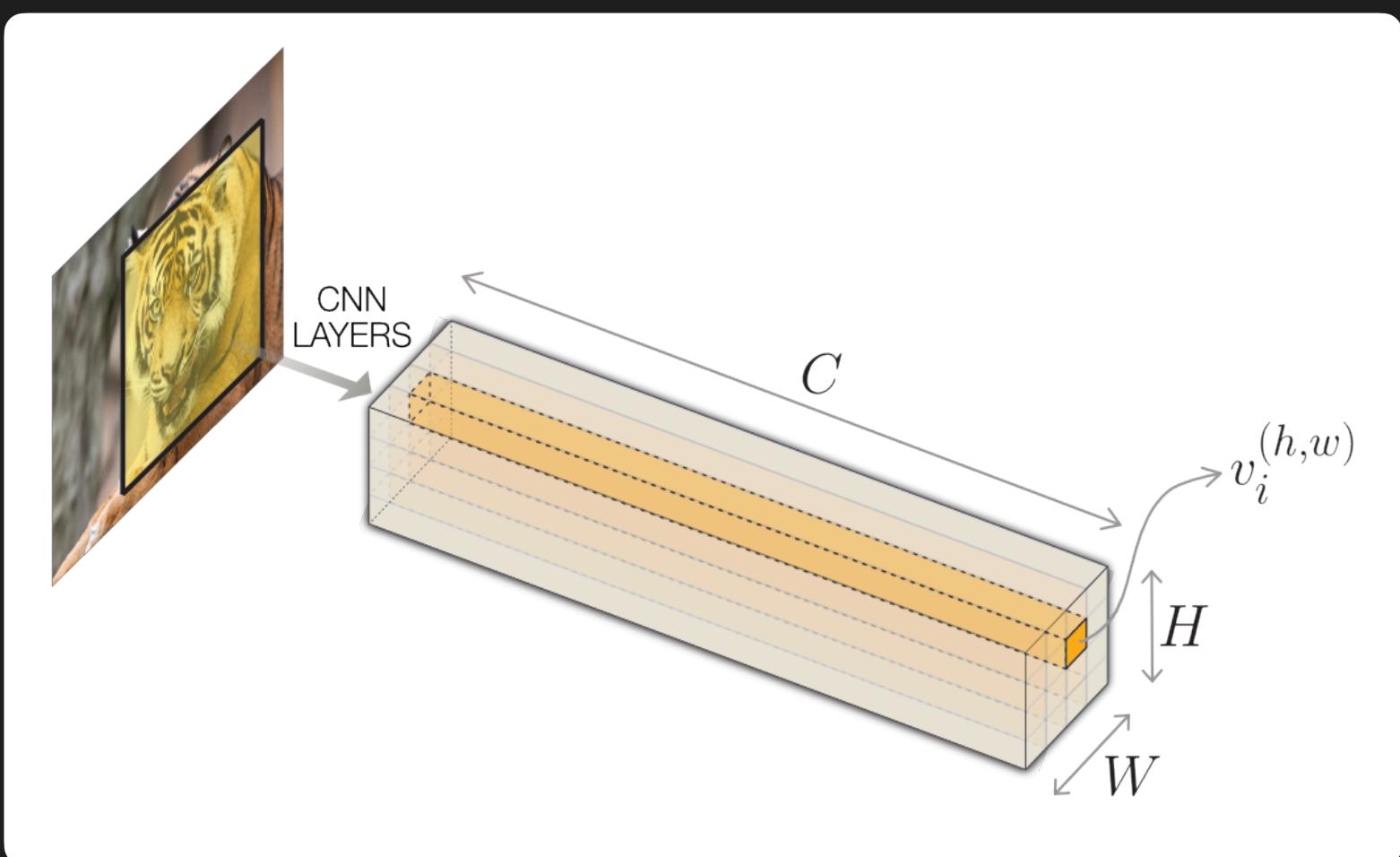
Which region of an image is the most predictive?



GOGGLES

Adversarial ML Interpretation

Data Programming



GOGGLES

Bluff

GOGGLES

SIGMOD 2020



Nilaksh Das
Georgia Tech

Automatic Image Labeling with *Affinity Coding*



Open-sourced at github.com/chu-data-lab/GOGGLES



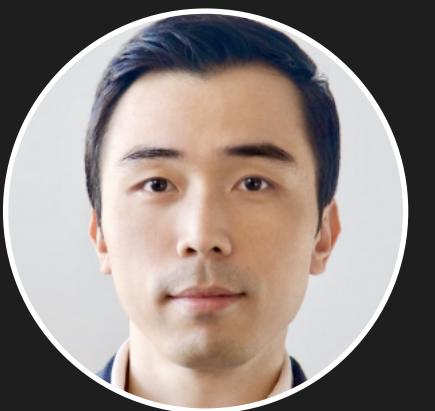
Sanya Chaba
Georgia Tech



Renzhi Wu
Georgia Tech



Sakshi Gandhi
Georgia Tech



Polo Chau
Georgia Tech



Xu Chu
Georgia Tech

Performance

Amount of Data

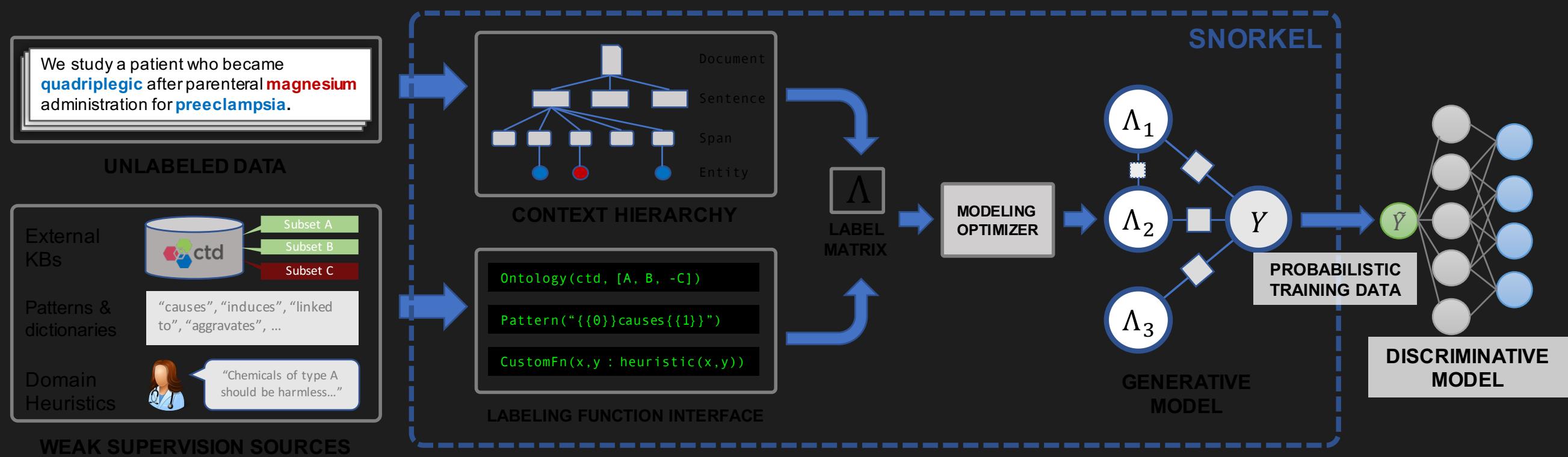
Deep Learning

**Traditional ML
Techniques**

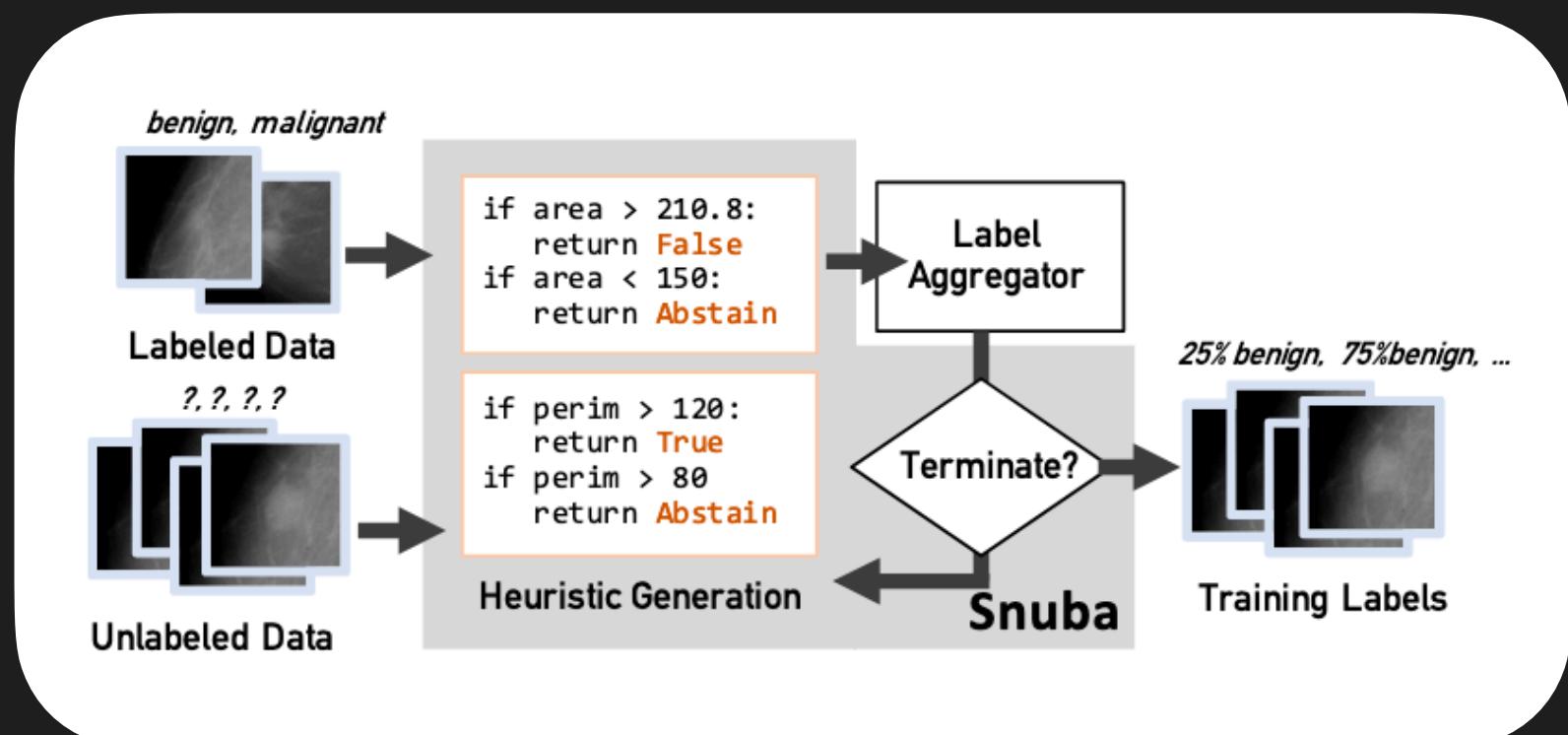
Collecting massive labeled datasets is
COSTLY

Data Programming

Snorkel



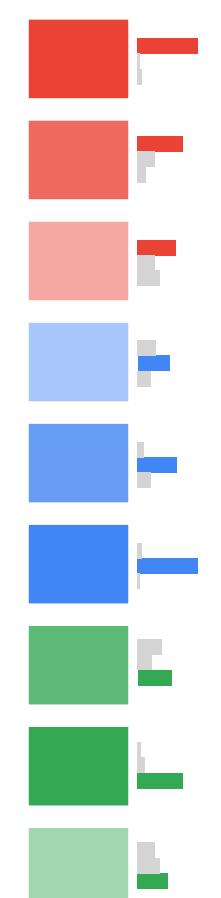
Snuba



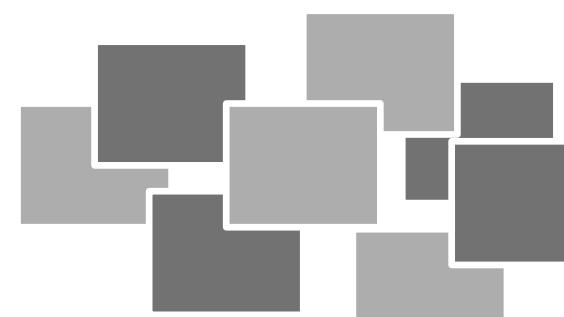
Affinity Coding

Affinity Coding with GOGGLES

Probabilistic Labels



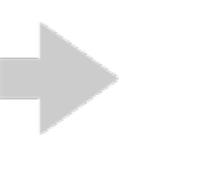
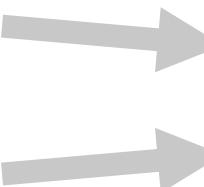
GOGGLES

All Data Instances
(without labels)

Images

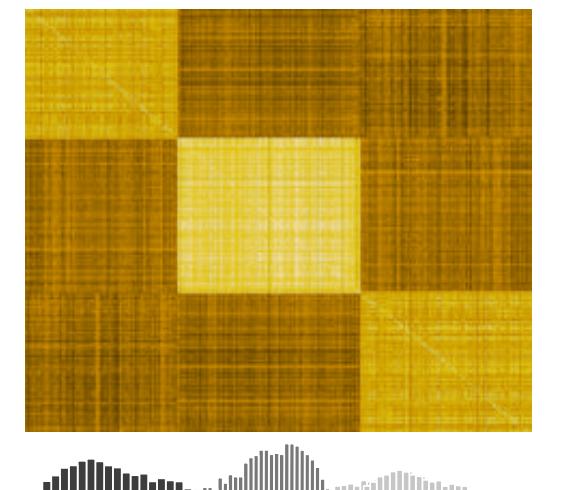
Step 1 Affinity Matrix Construction

x_i

x_j

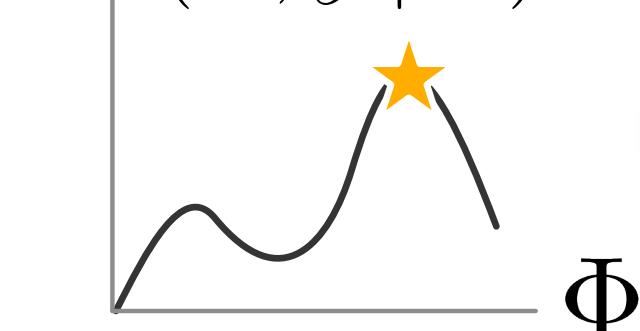


Affinity Matrix \mathcal{A}



Step 2 Class Inference

$$\Pr(\mathcal{A}, y \mid \Phi)$$



Small, Labeled Development Set

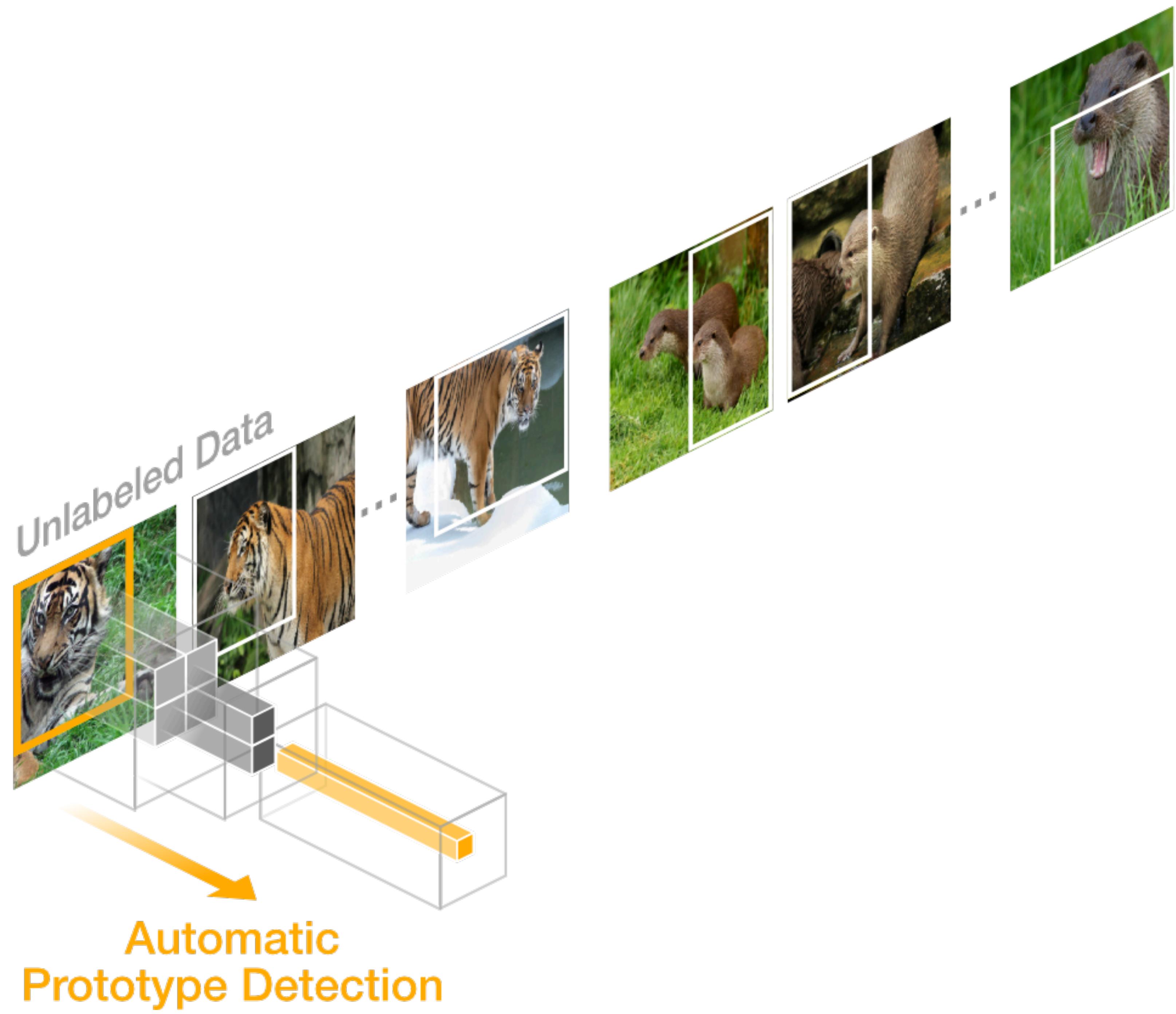


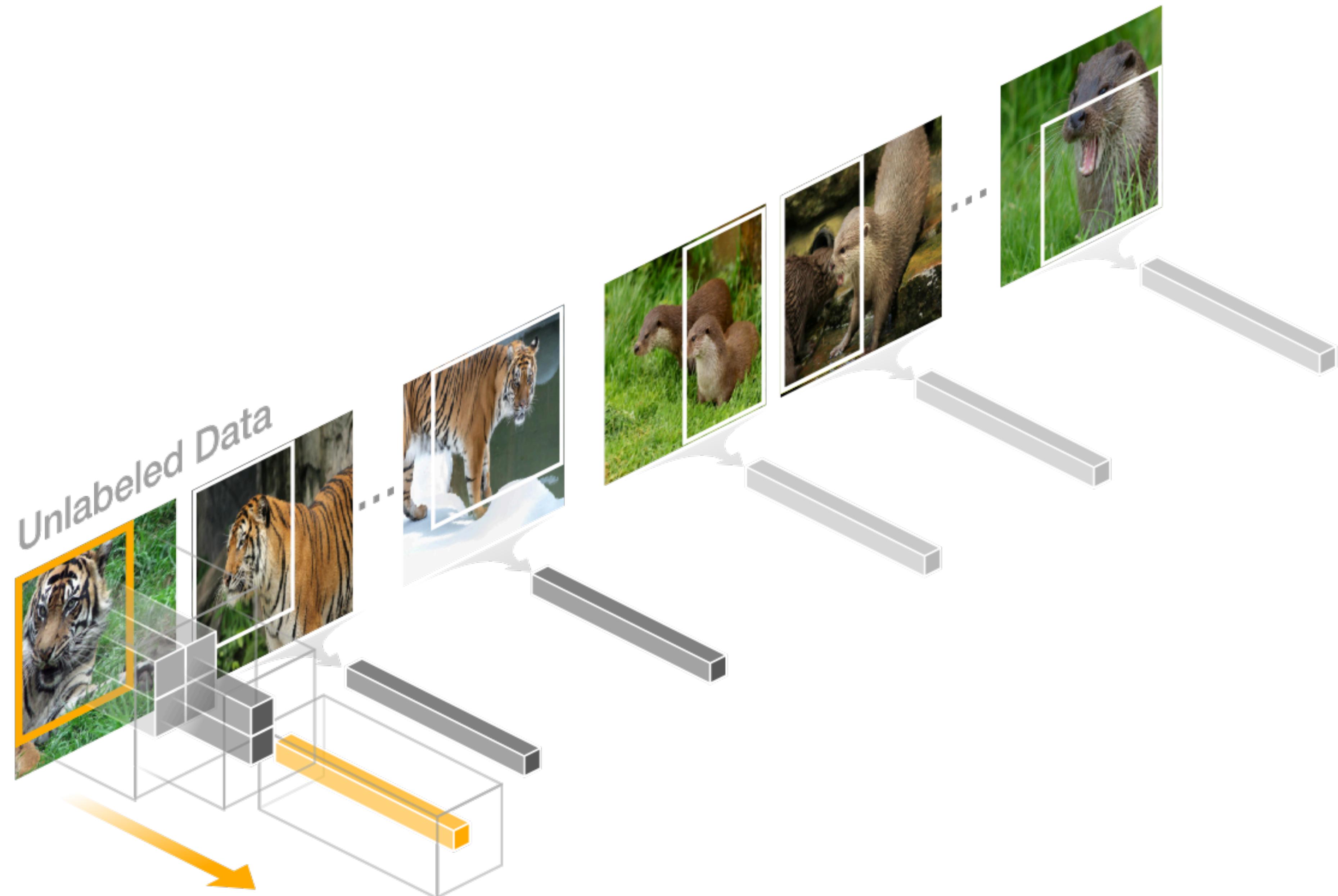
Library of Affinity Functions



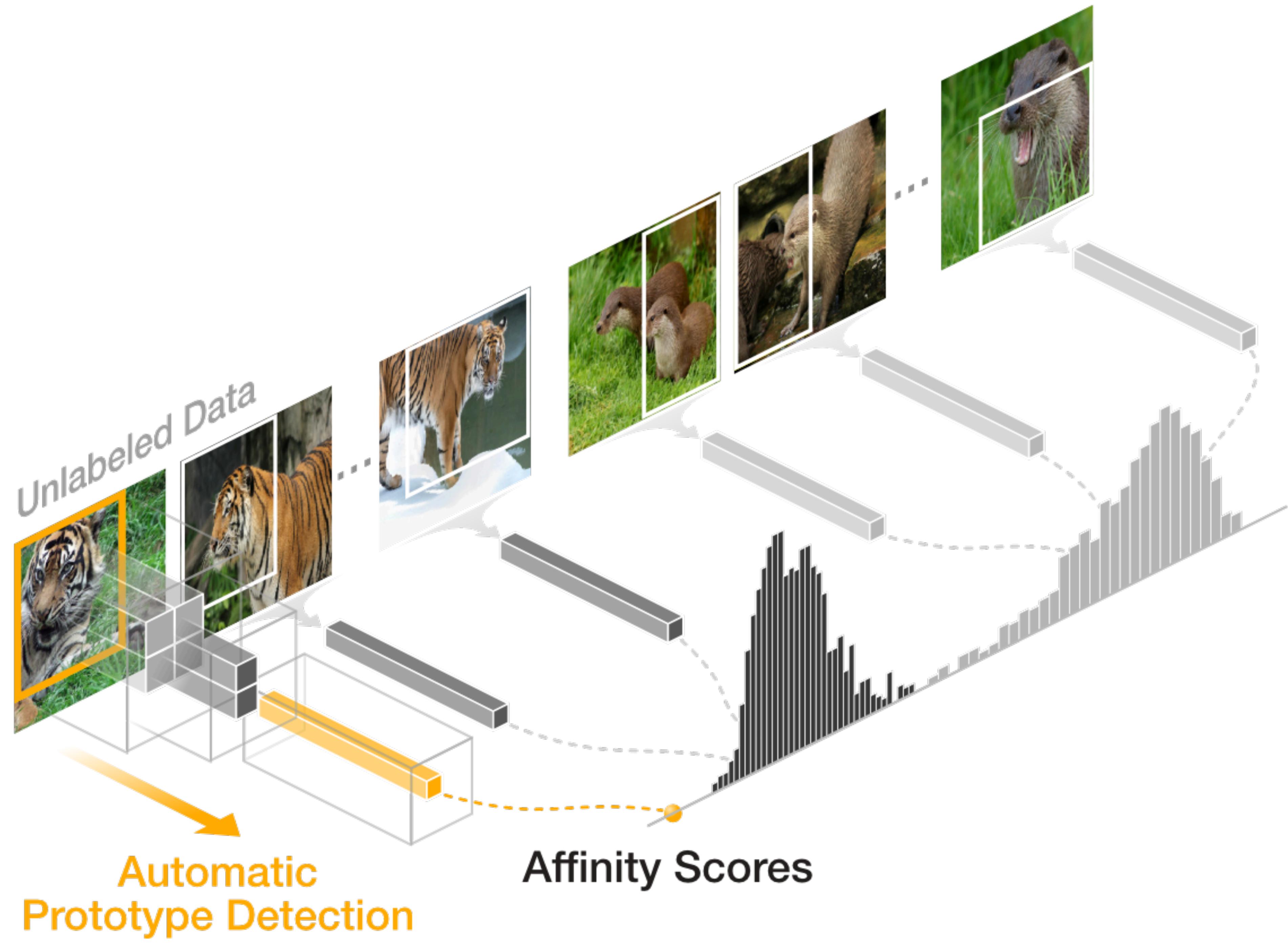


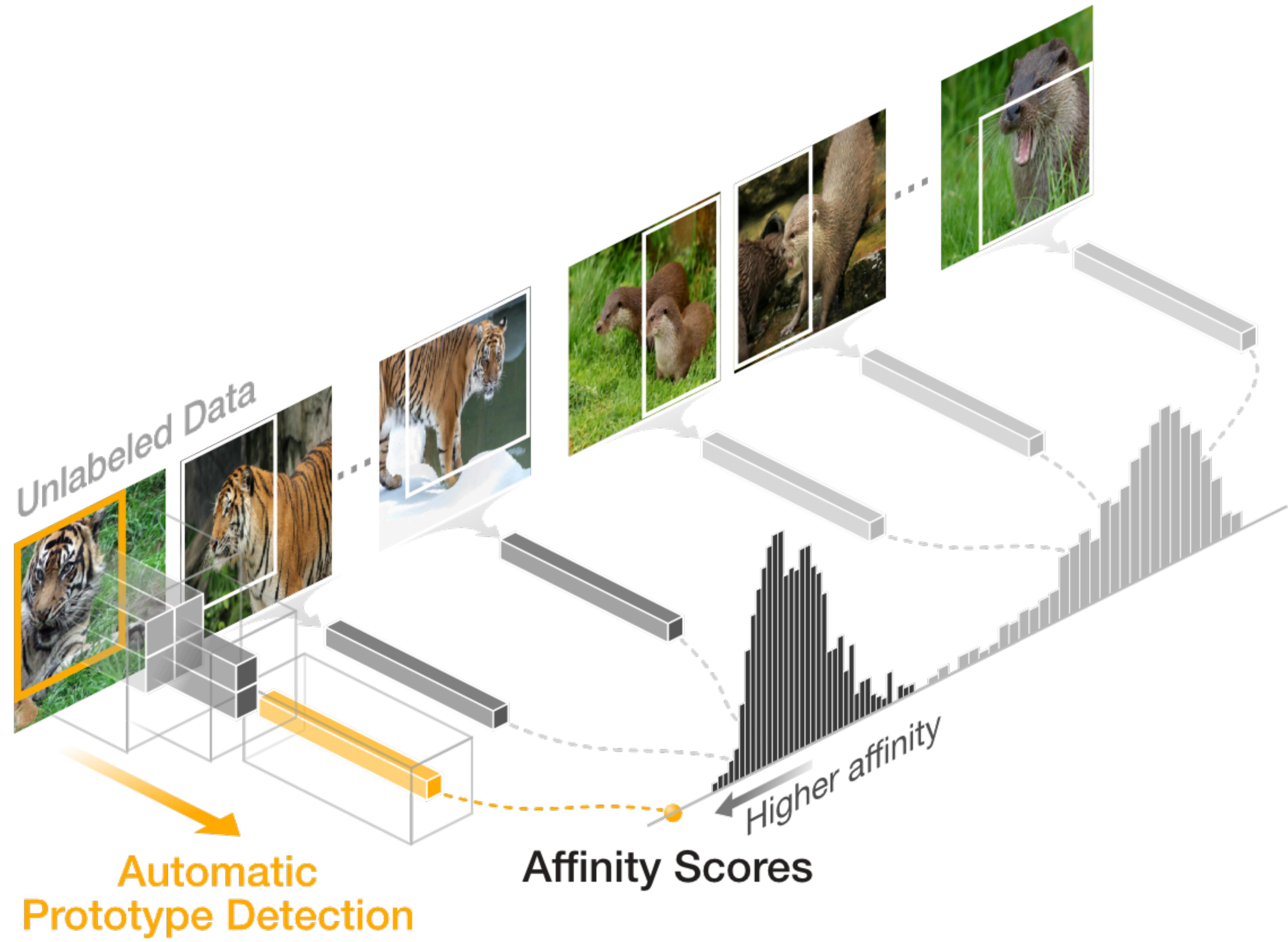


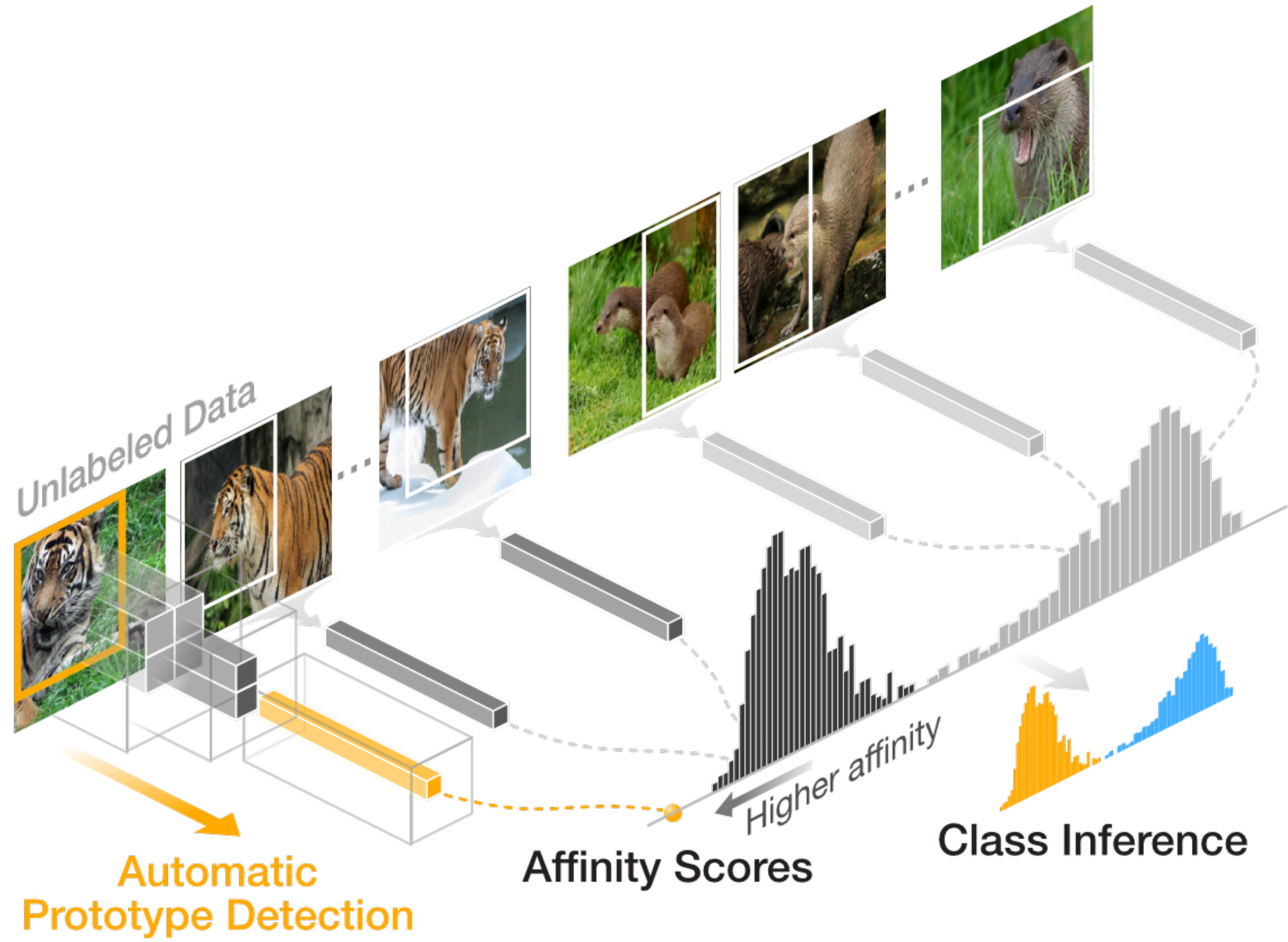


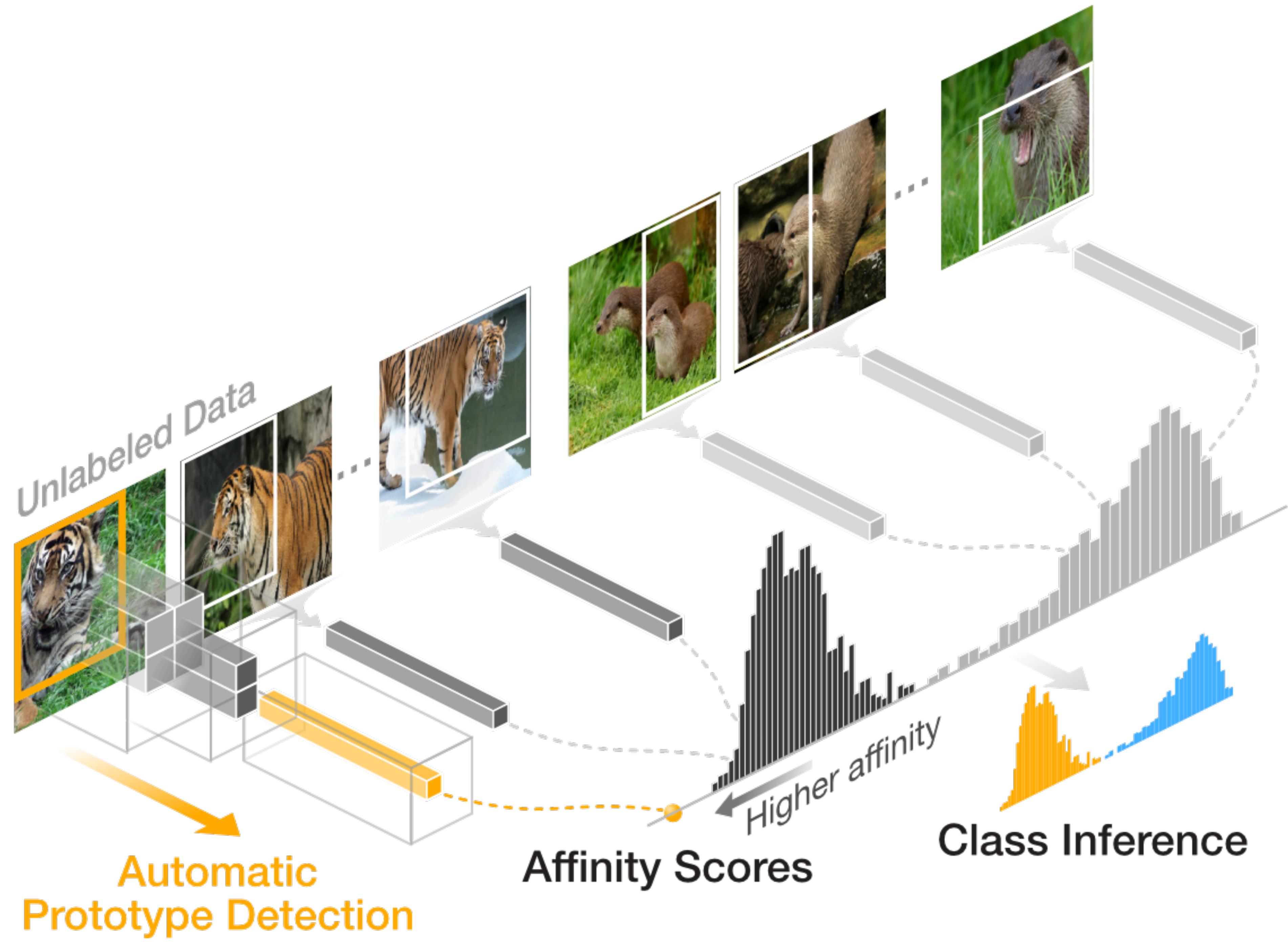


**Automatic
Prototype Detection**









VGG16

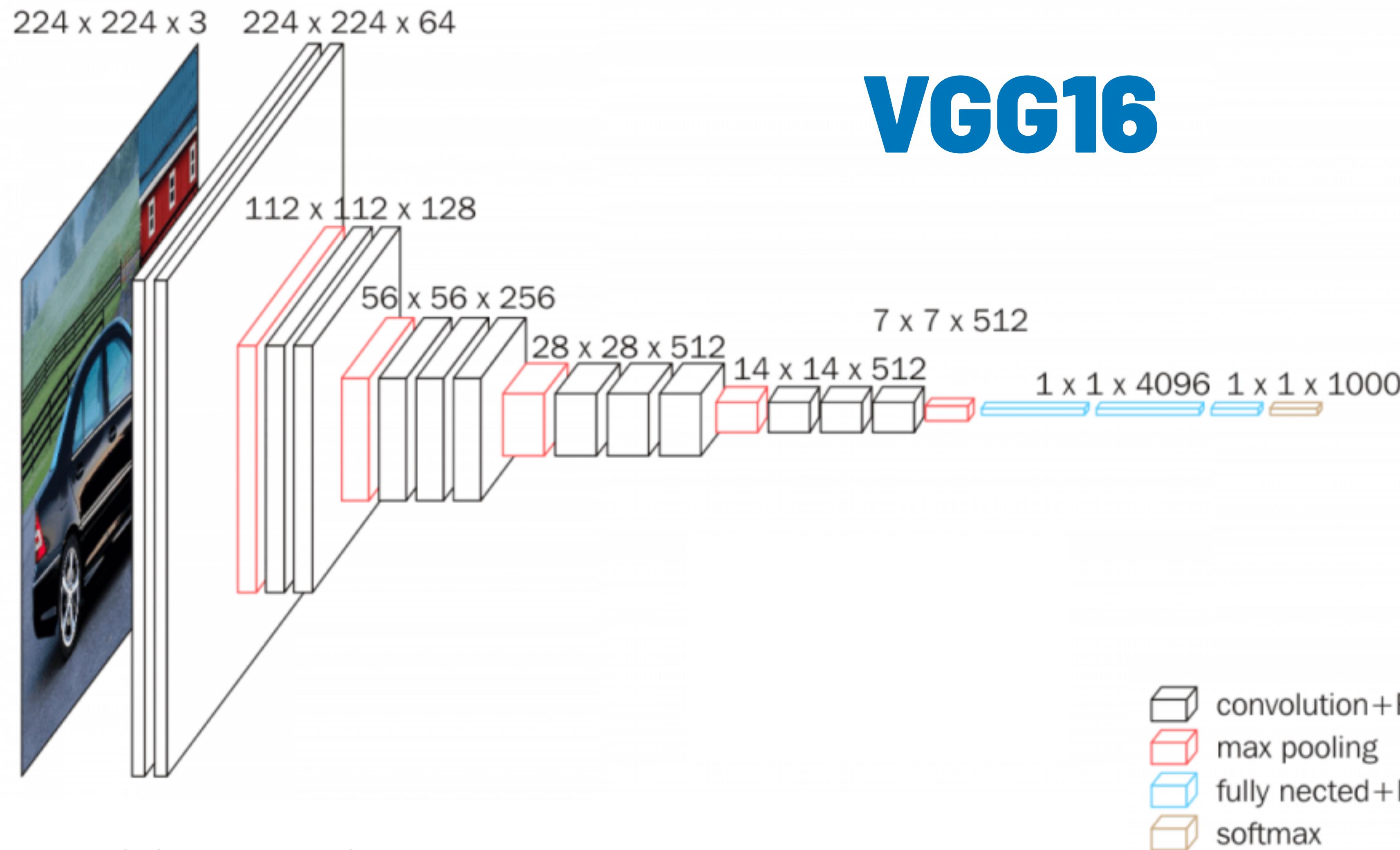
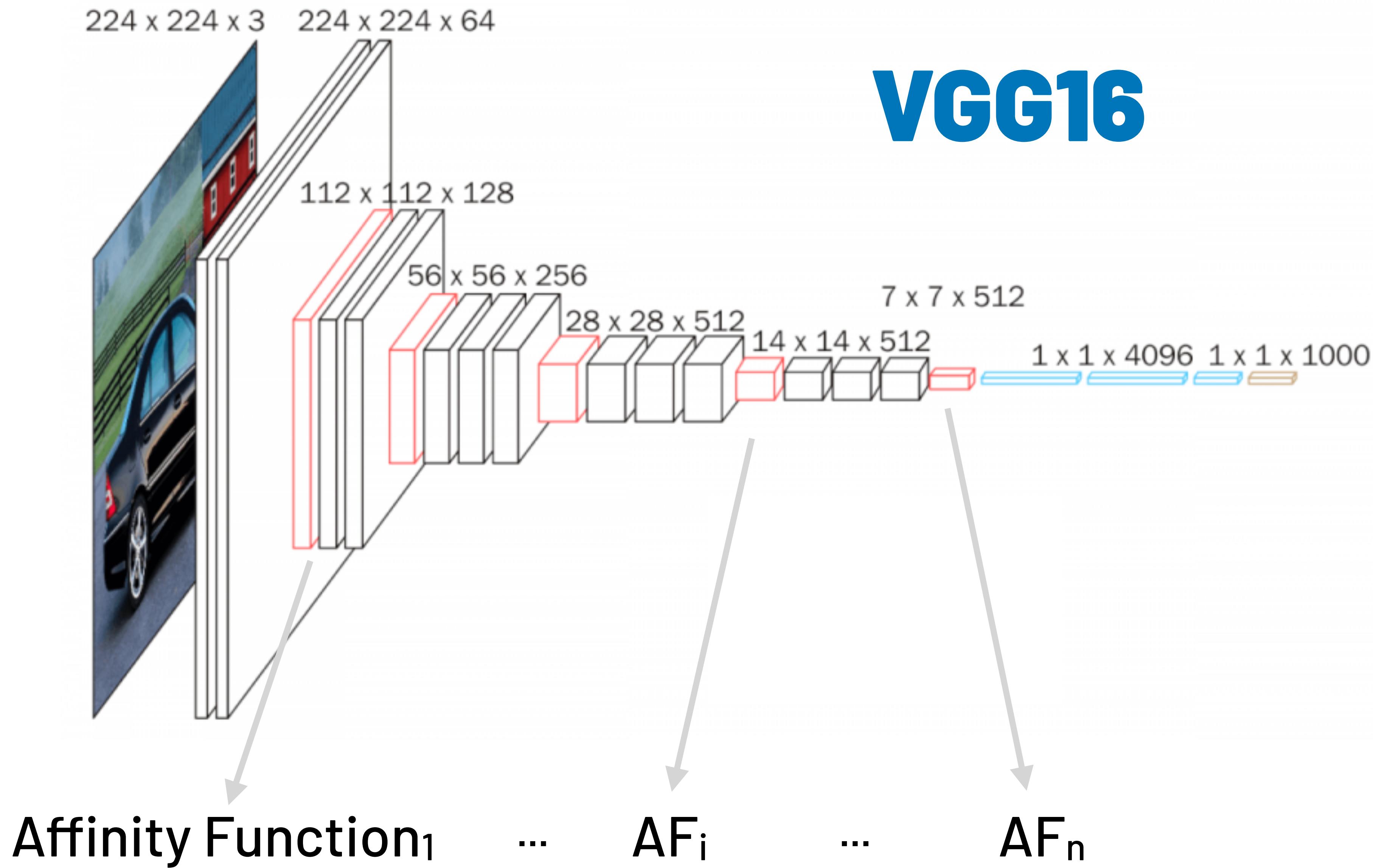
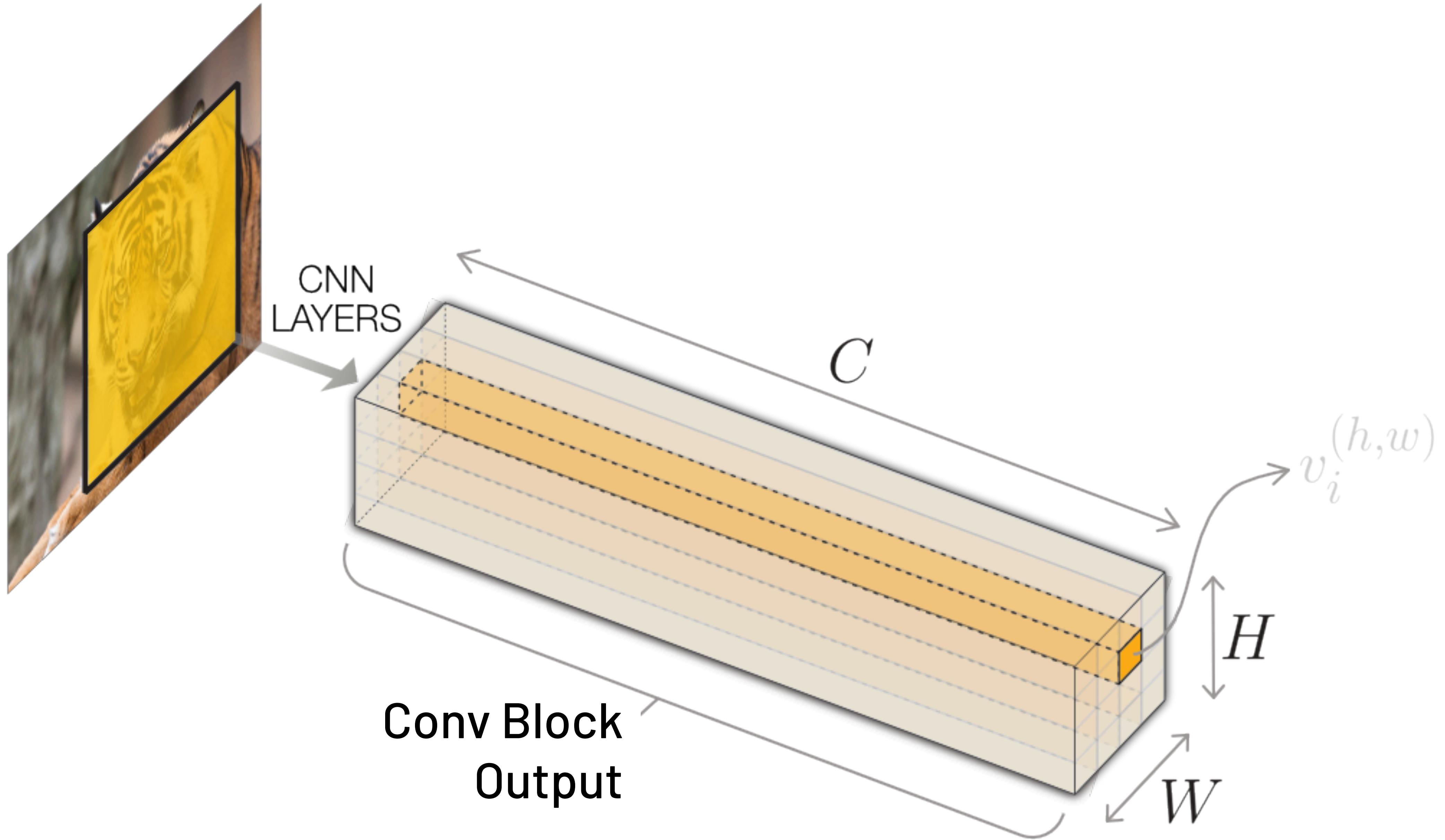


Image: neurohive.io/en/popular-networks/vgg16

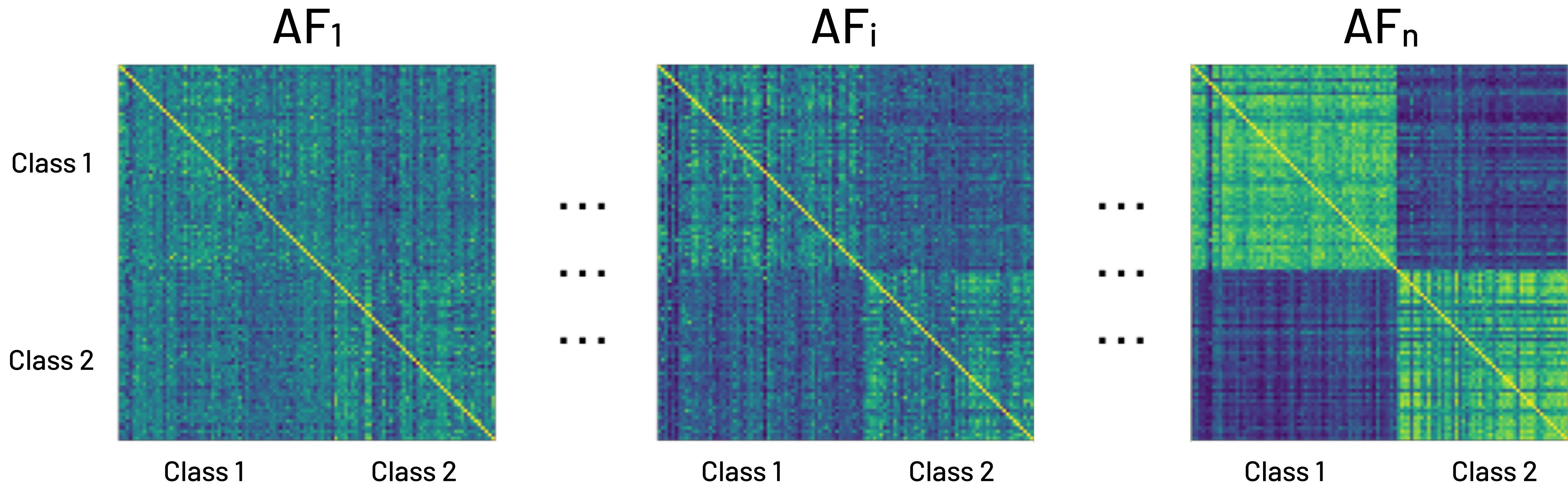
VG16





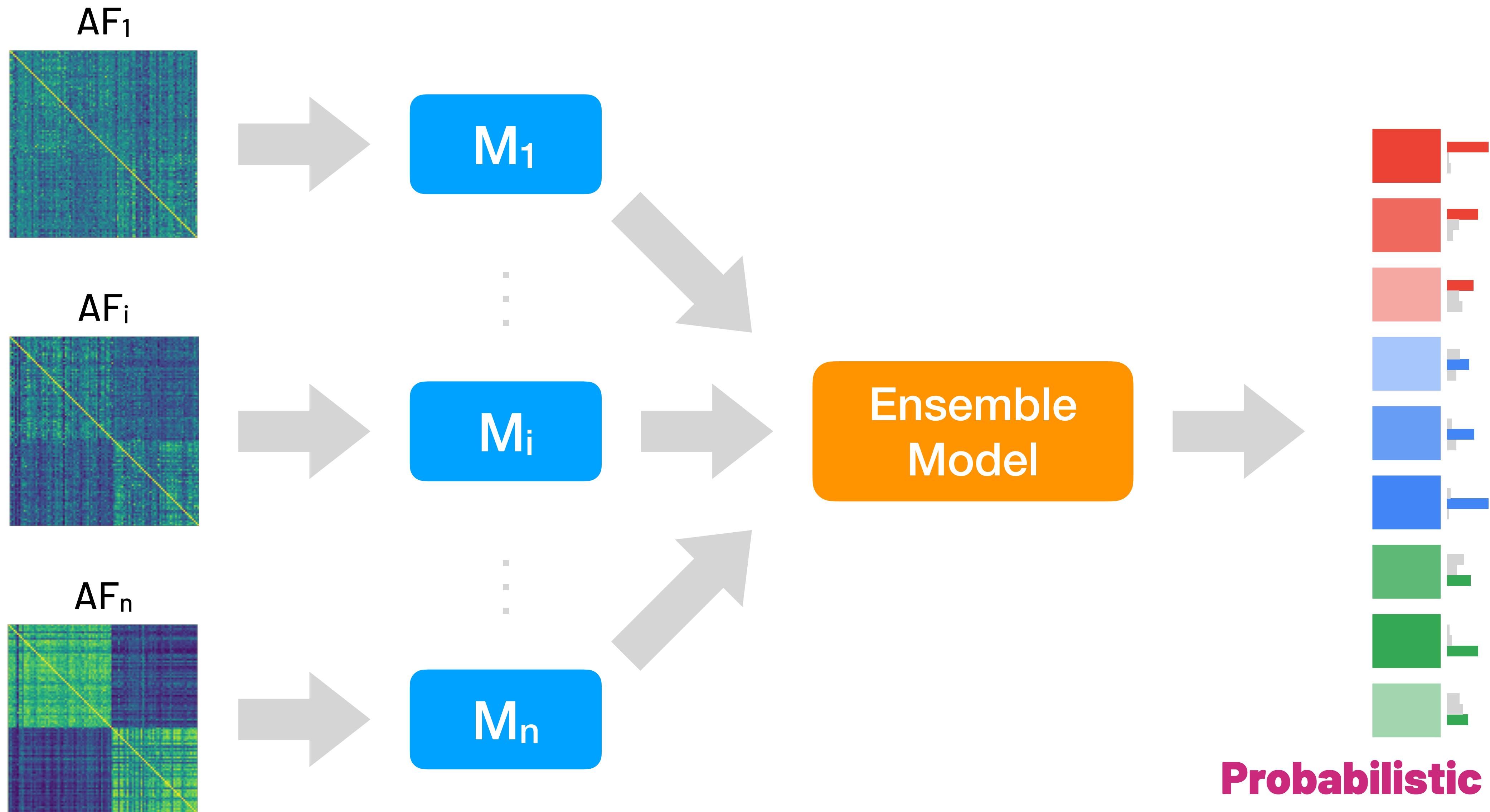
Affinity Matrices

from different **Affinity Functions**

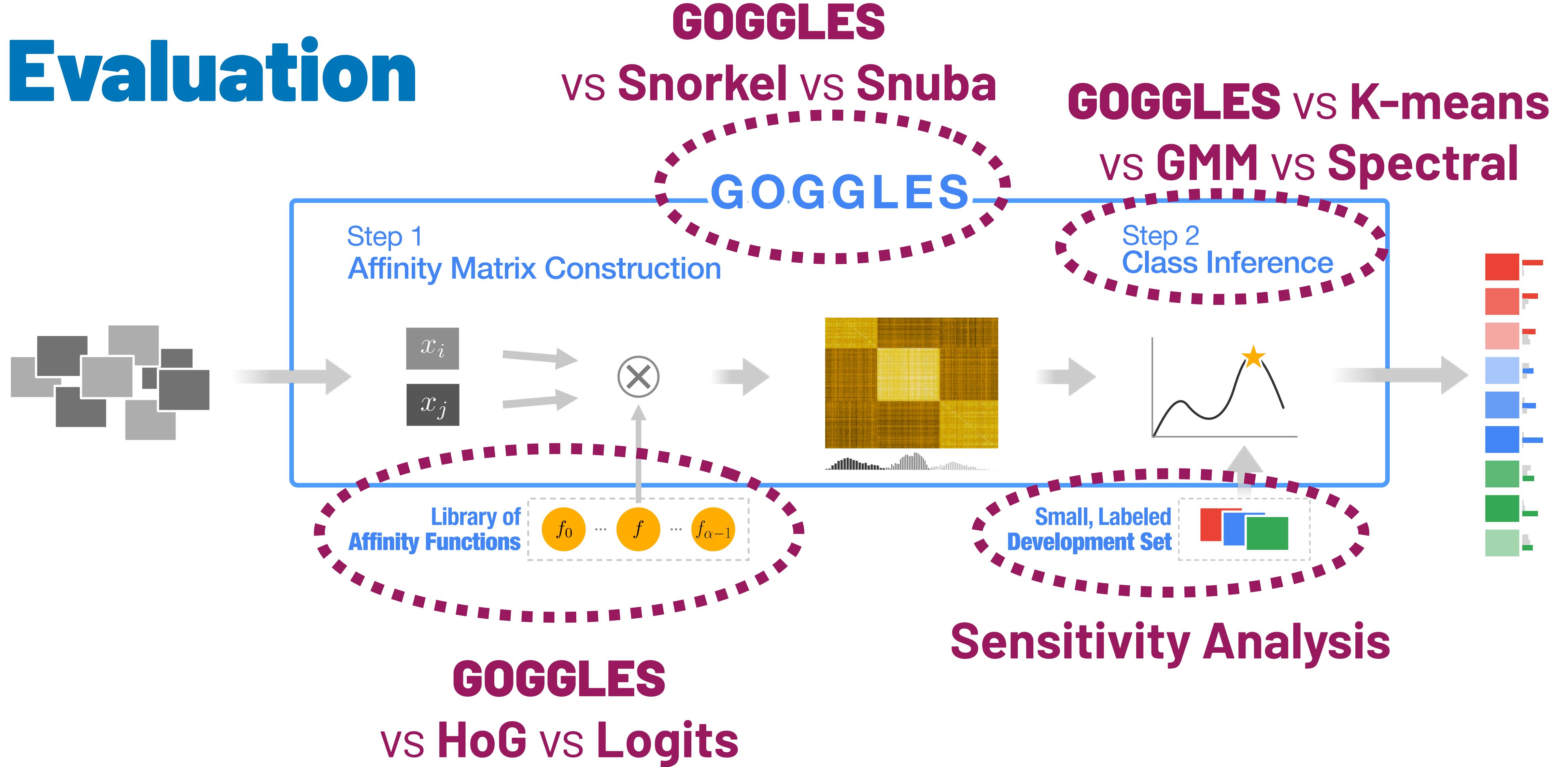


Class Inference

with **Hierarchical Generative Model**



Exhaustive Evaluation



Diverse Datasets



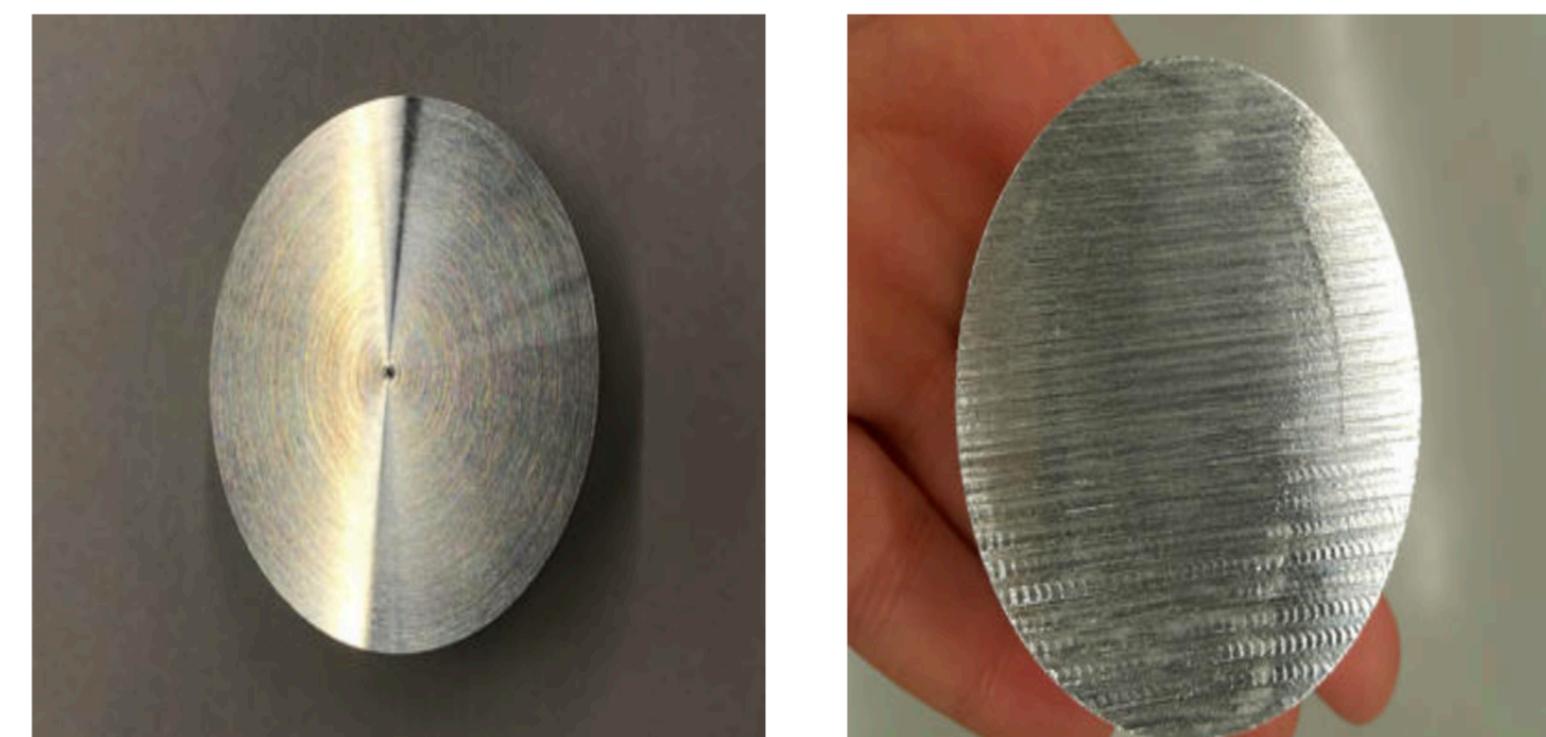
CUB (Birds)



GTSRB (Traffic Signs)



TB-Xray (Tuberculosis)
PN-Xray (Pneumonia)



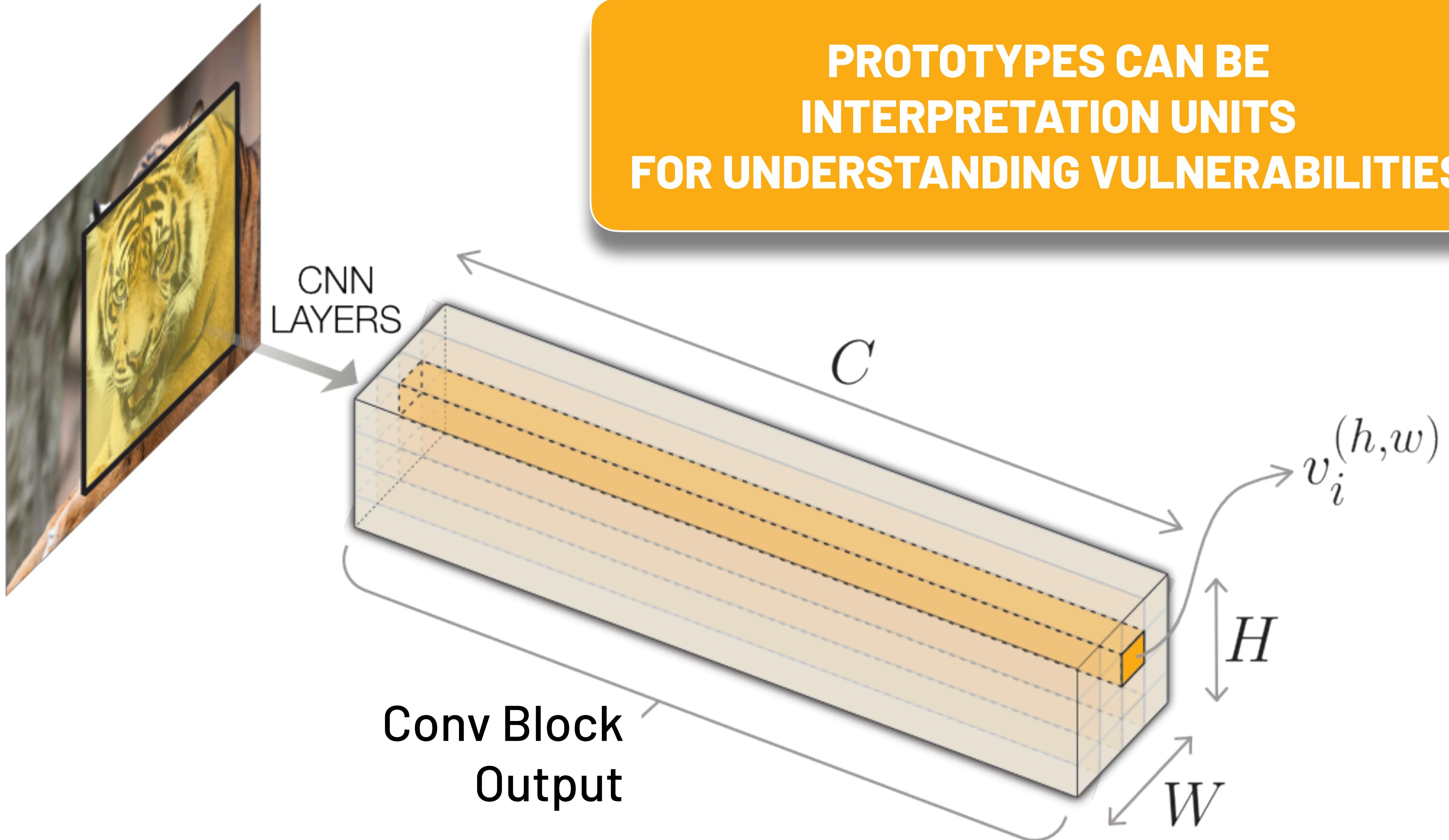
Surface (Industrial Metal Finish)

Evaluating **Labeling Accuracy** on **Training Set**

Dataset	GOGGLES (our results)	Data Programming		Representation		Class Inference Baselines		
		Snorkel	Snuba	HoG	Logits	K-Means	GMM	Spectral
CUB	97.83	89.17	58.83	62.93	96.35	98.67	97.62	72.08
GTSRB	70.51	-	62.74	75.48	64.77	70.74	69.64	62.40
Surface	89.18	-	57.86	85.82	54.08	69.08	69.14	60.82
TB-Xray	76.89	-	59.47	69.13	67.16	76.33	76.70	75.00
PN-Xray	74.39	-	55.50	53.11	71.18	50.66	68.66	75.90
Average	81.76	-	58.88	69.30	70.71	73.09	76.35	69.24

↔
~23%

PROTOTYPES CAN BE
INTERPRETATION UNITS
FOR UNDERSTANDING VULNERABILITIES



Bluff

IEEE VIS 2020

Interactively Deciphering Adversarial Attacks On Deep Neural Networks



Live Demo at poloclub.github.io/bluff



Nilaksh Das*
Georgia Tech

*equal contribution



Haekyu Park*
Georgia Tech



Zijie J. Wang
Georgia Tech



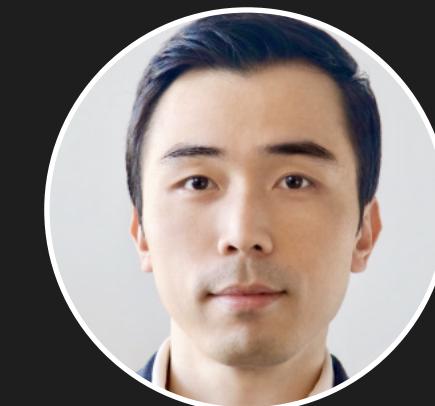
Fred Hohman
Georgia Tech



Robert Firstman
Georgia Tech

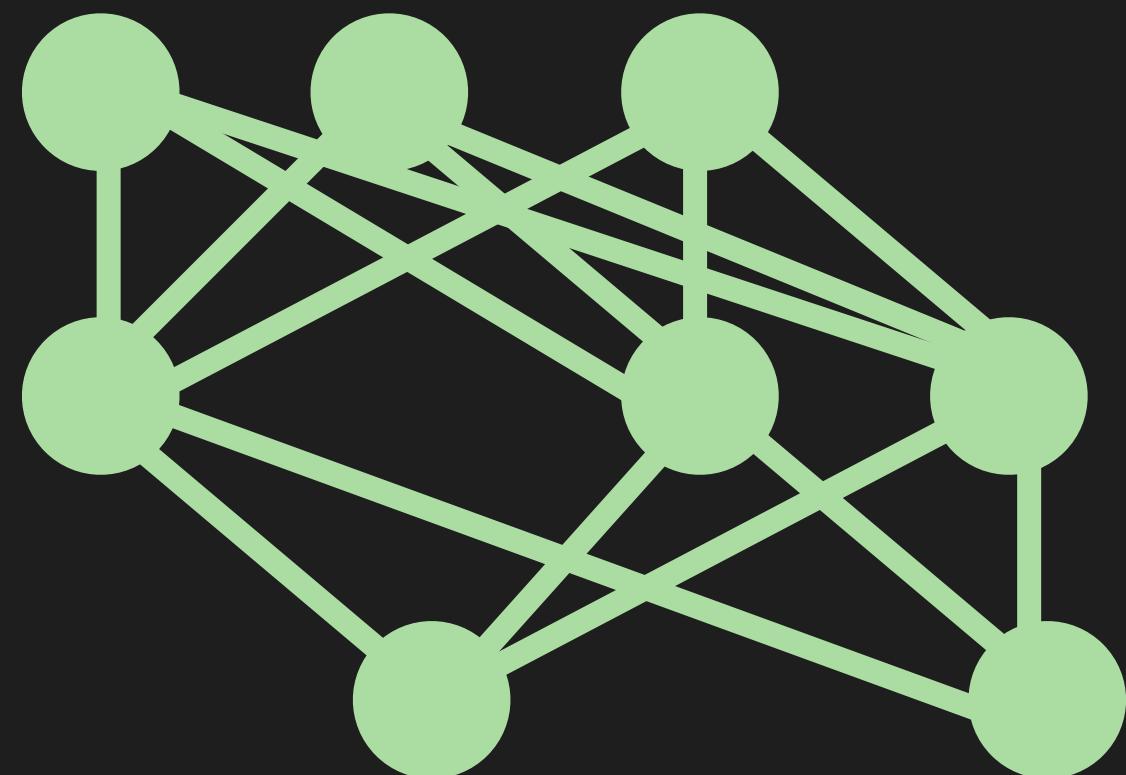


Emily Rogers
GT Research Institute

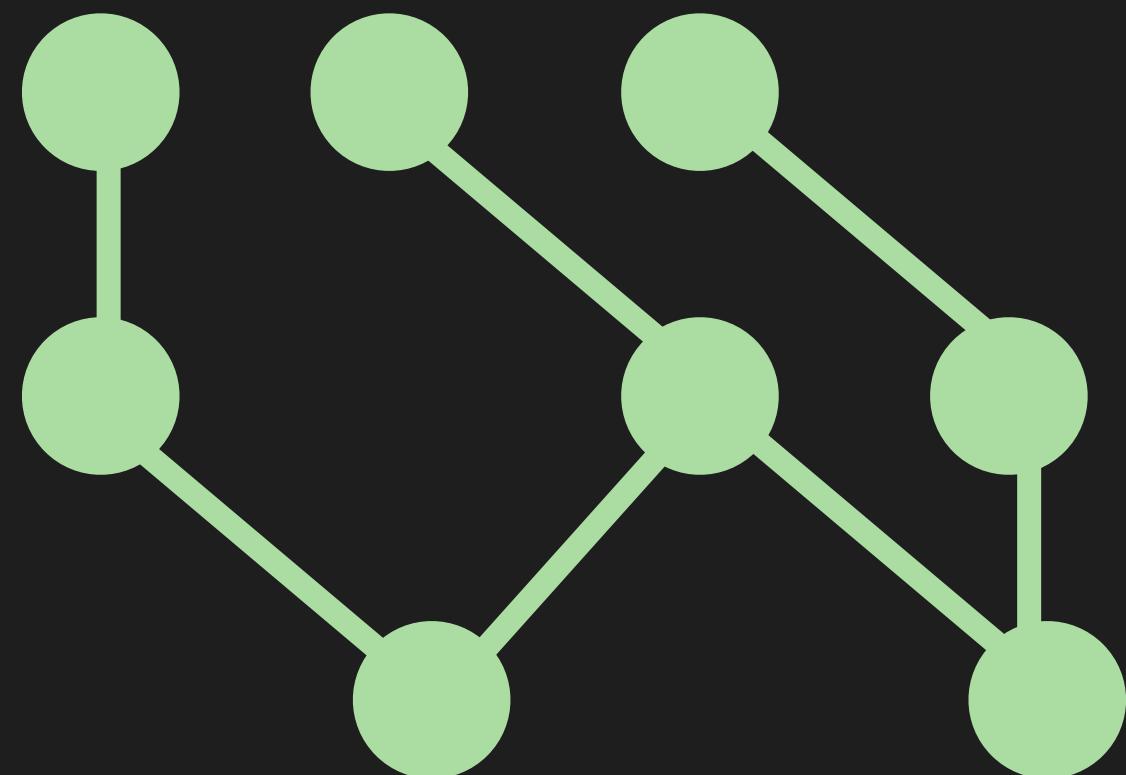


Polo Chau
Georgia Tech

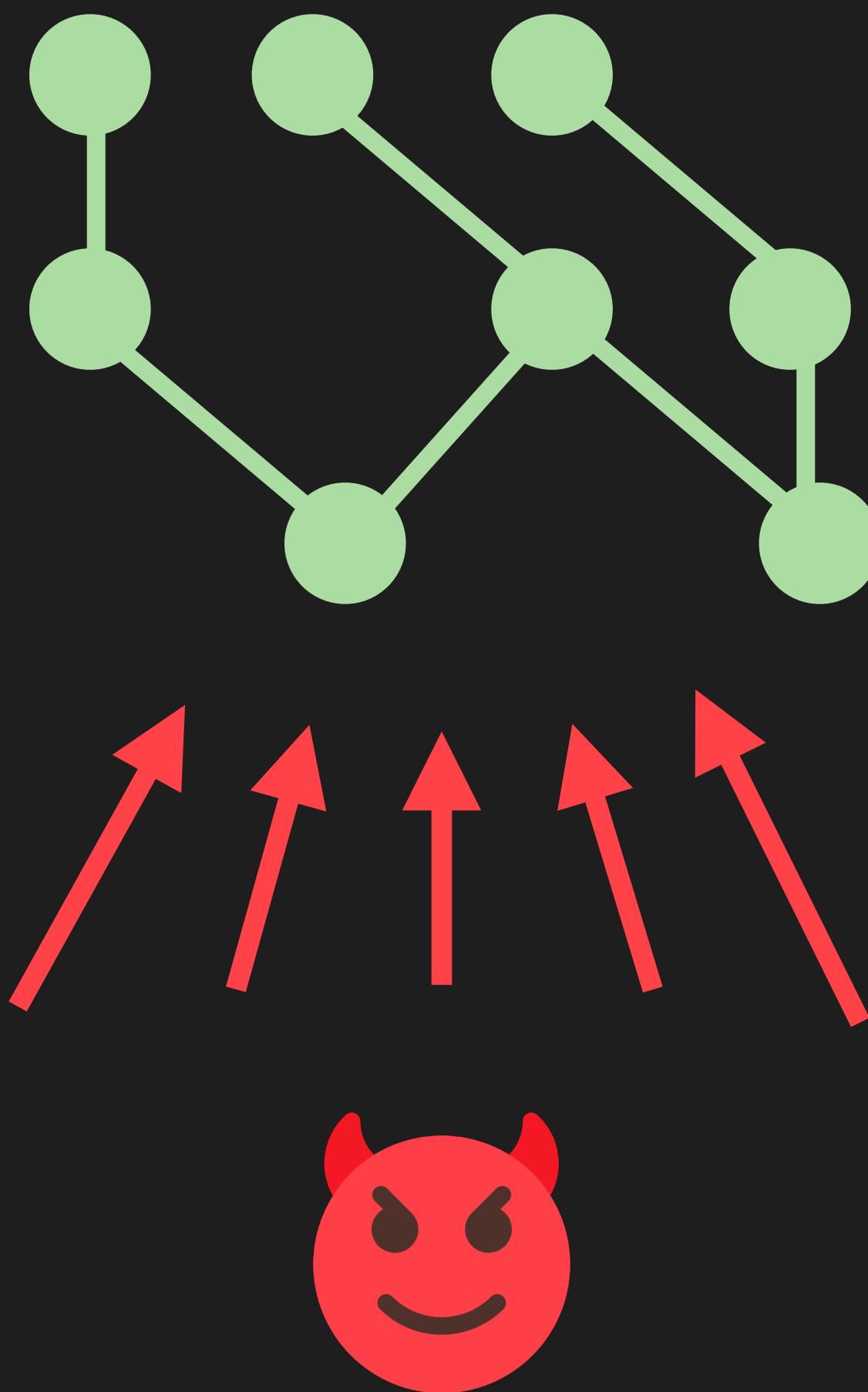
Understanding Adversarial Attacks



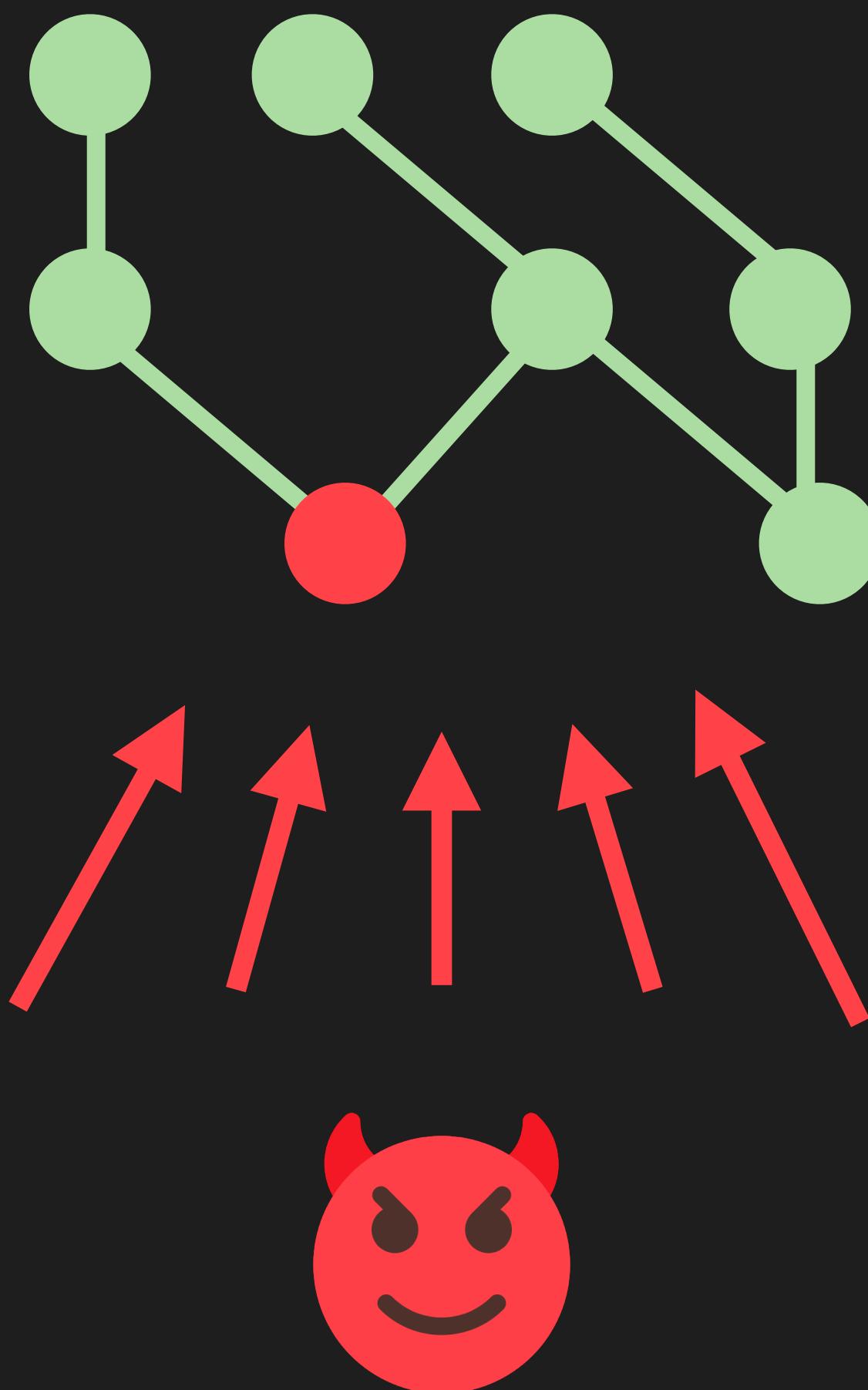
Understanding Adversarial Attacks



Understanding Adversarial Attacks

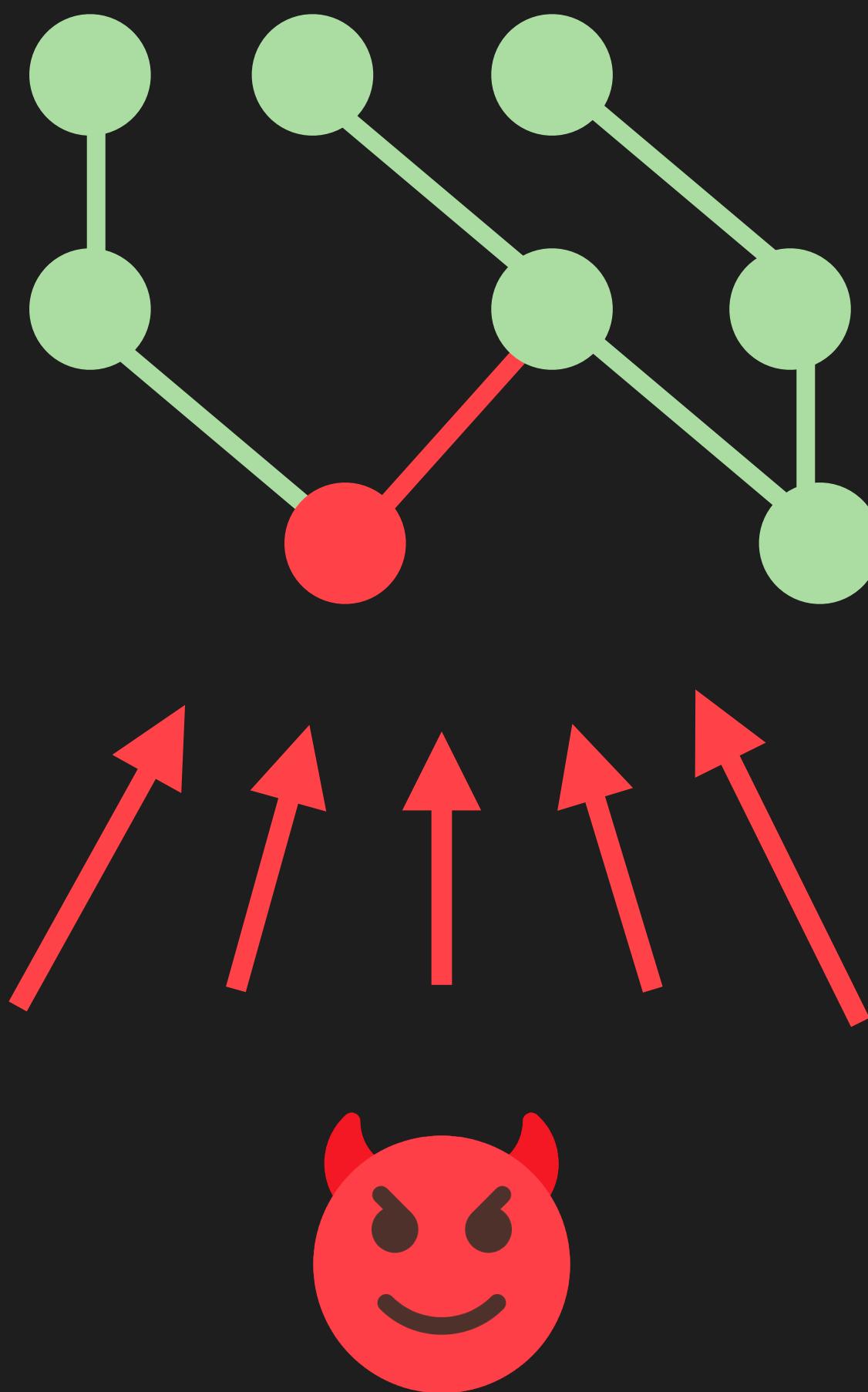


Understanding Adversarial Attacks



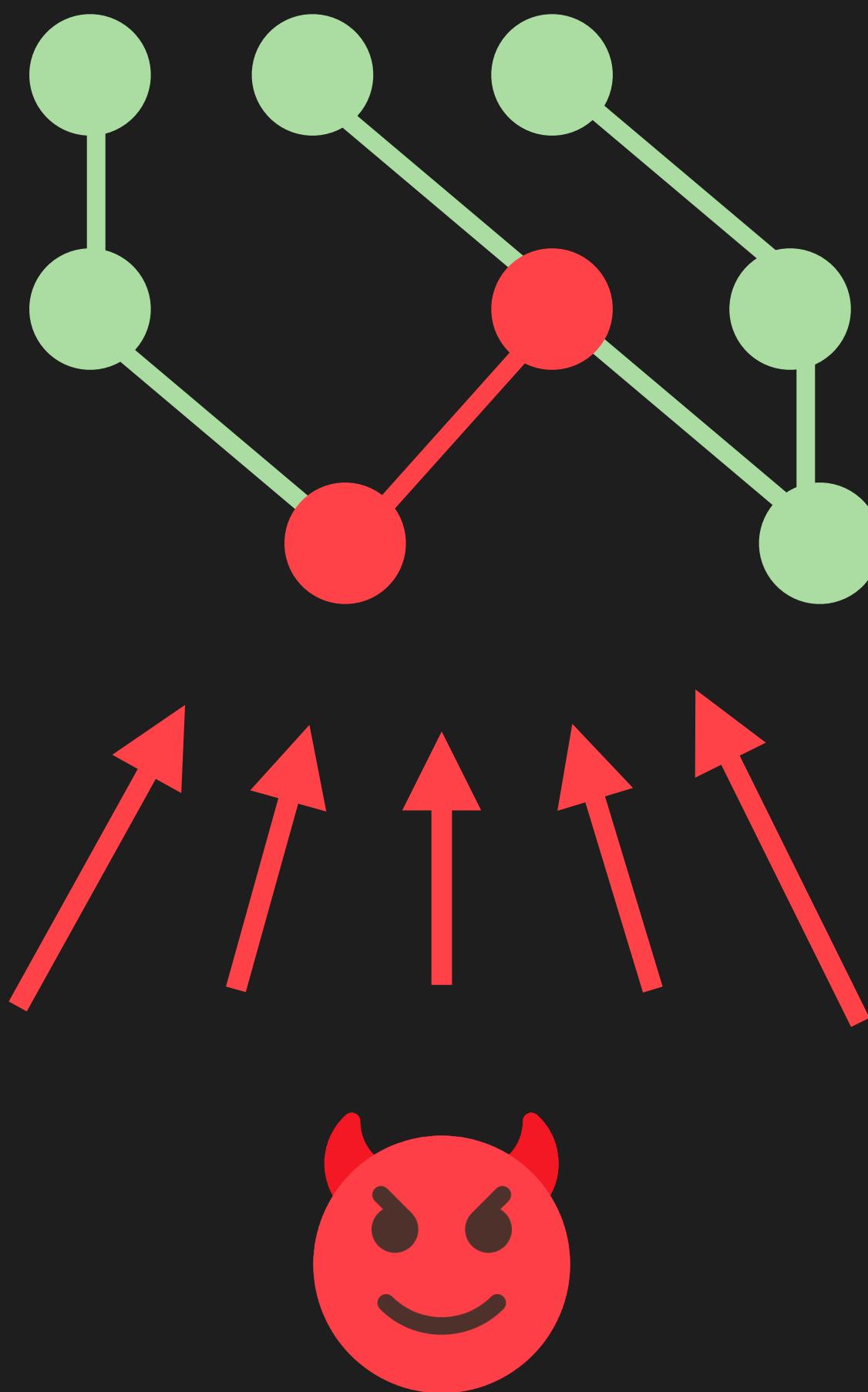
How adversarial attacks
permeate a model?

Understanding Adversarial Attacks



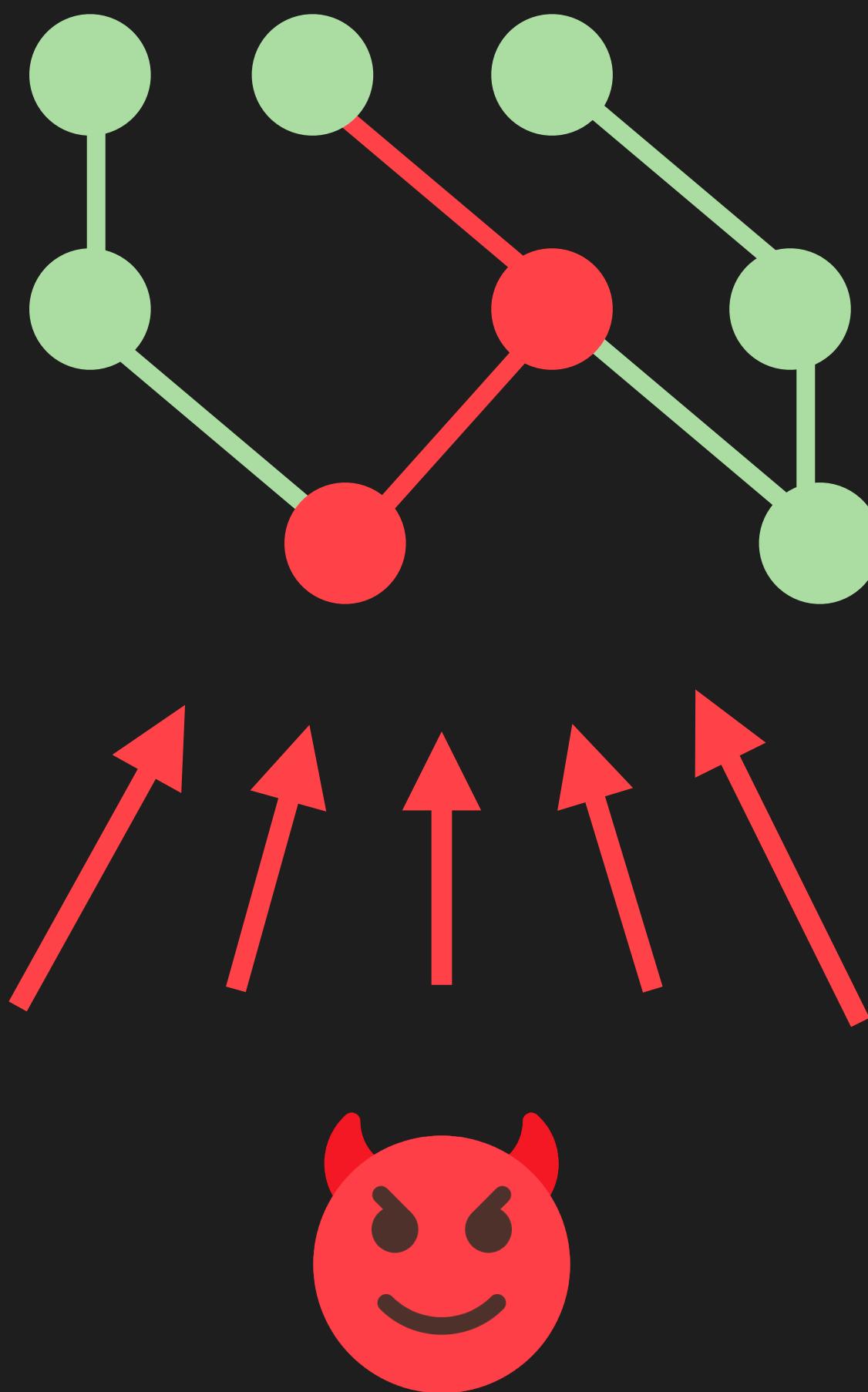
How adversarial attacks
permeate a model?

Understanding Adversarial Attacks



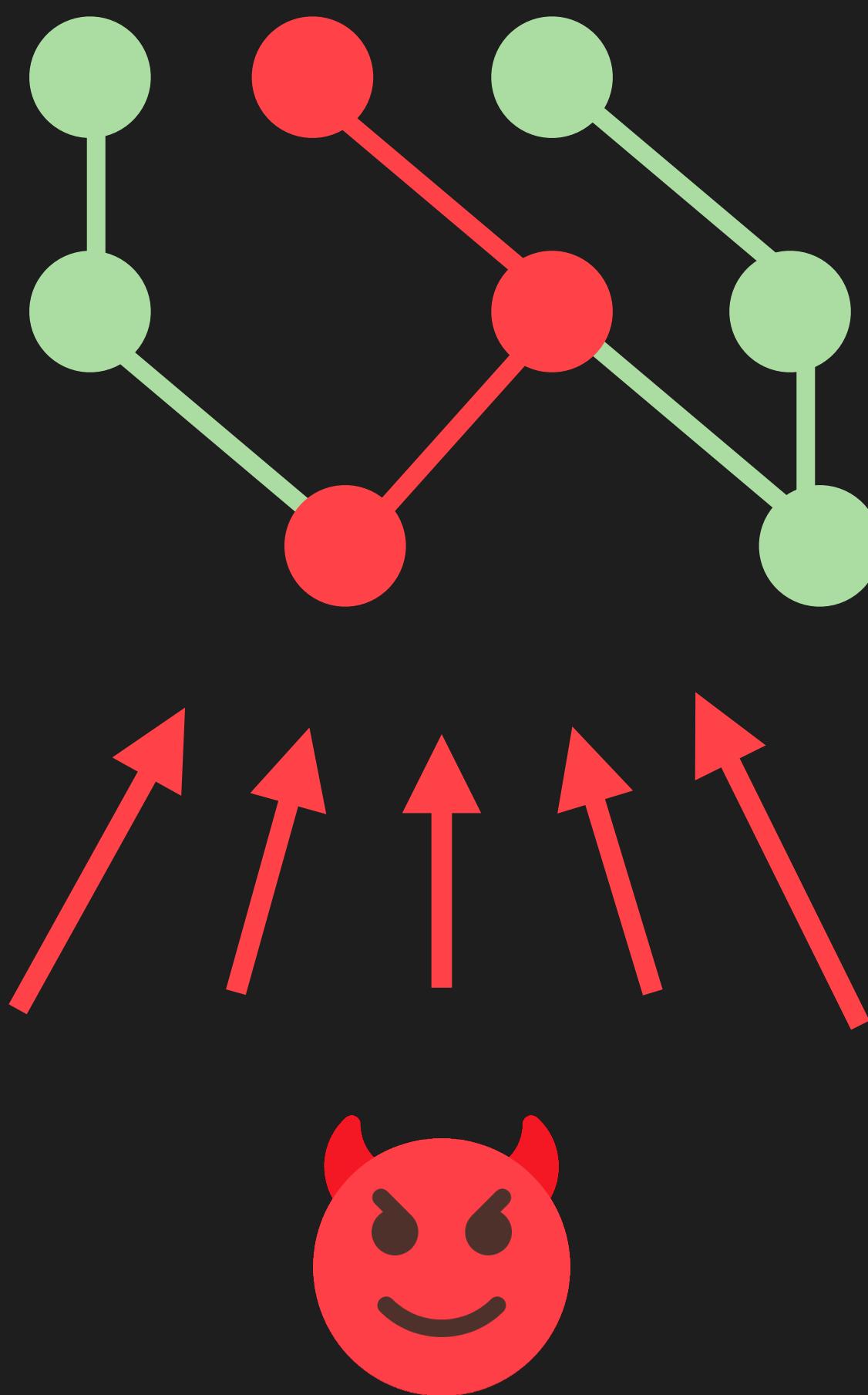
How adversarial attacks
permeate a model?

Understanding Adversarial Attacks



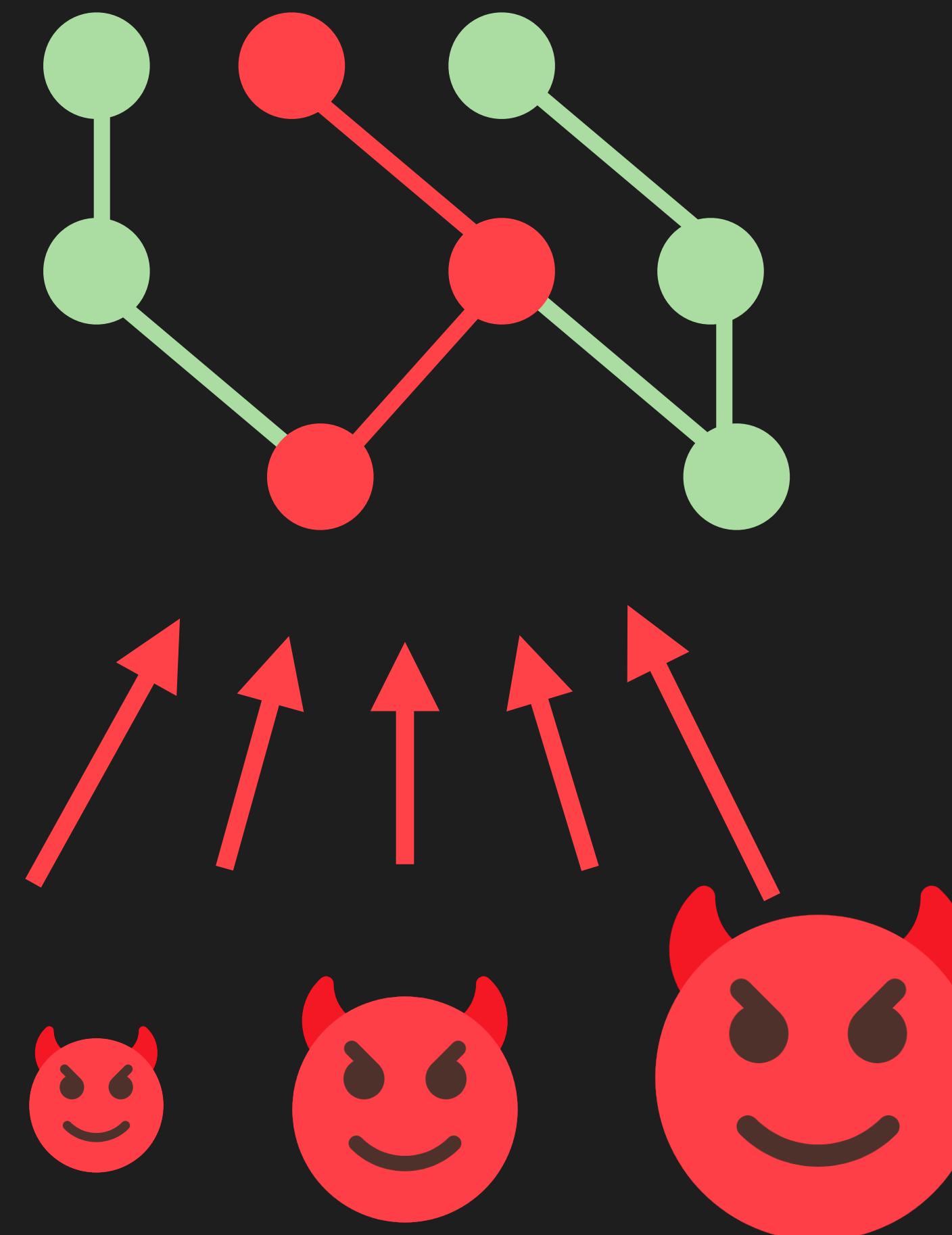
How adversarial attacks
permeate a model?

Understanding Adversarial Attacks



How adversarial attacks
permeate a model?

Understanding Adversarial Attacks



How adversarial attacks
permeate a model?

Variation in attack patterns
with varying
attack strengths?

Bluff

Understand how neural networks misclassify GIANT PANDA ▾ into ARMADILLO ▾ when attacked

A Control Sidebar

ADVERSARIAL ATTACK

PGD

Strength: 0.05



FILTER GRAPH

Show full graph

Show pinned only

Show highly activated

HIGHLIGHT PATHWAYS

Highlight pathways most excited ▾ by attack.

Neurons: top 45 % in each layer



Connections: top 50 %



Bluff

discovers and deciphers

how adversarial attacks fool DNNs

B Graph Summary View

GIANT PANDA

BOTH

ARMADILLO

EXPLOITED BY ATTACK

mixed5b

473 525 625 798

mixed5a

253 223

mixed4e

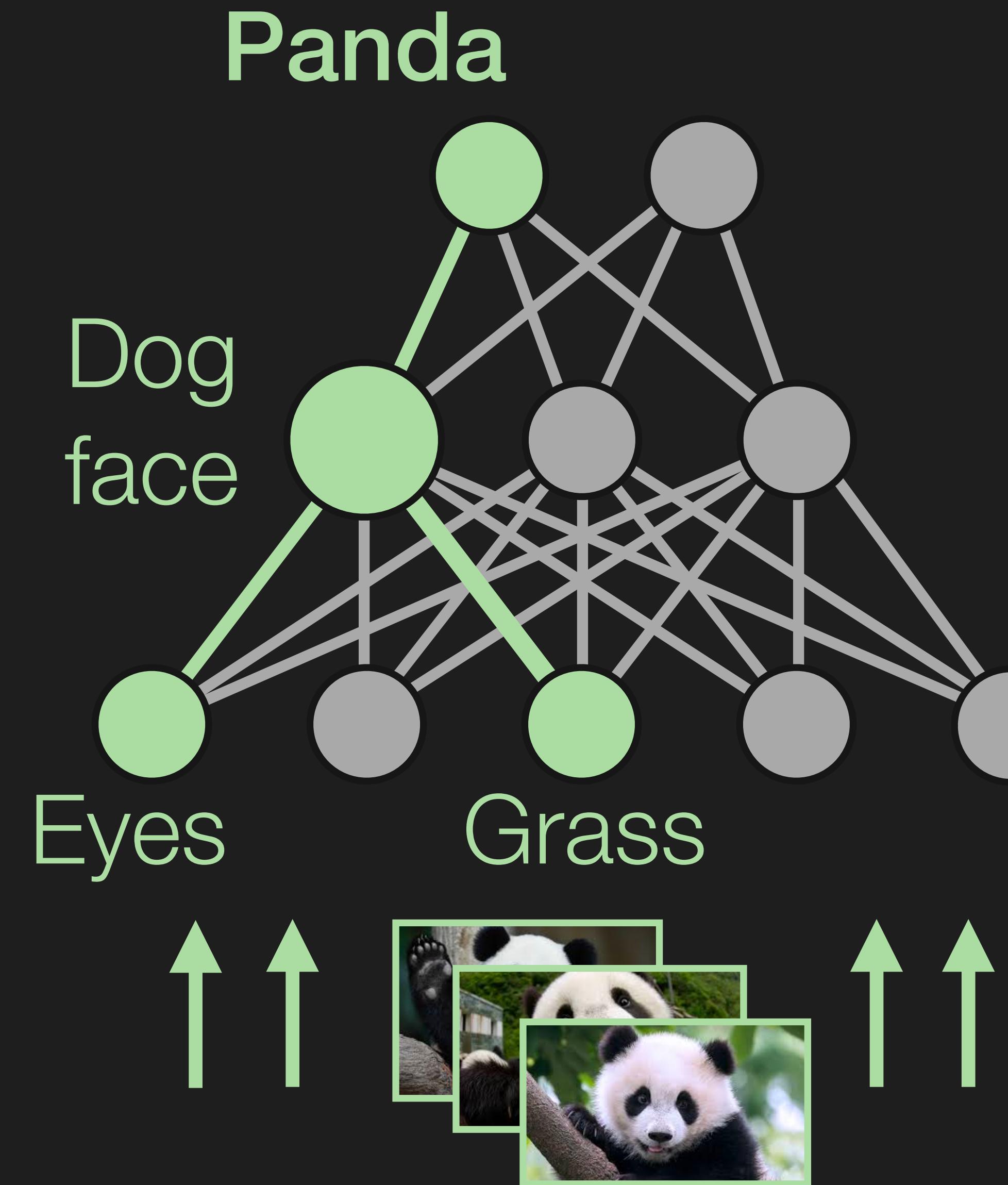
360 557 516

mixed4d

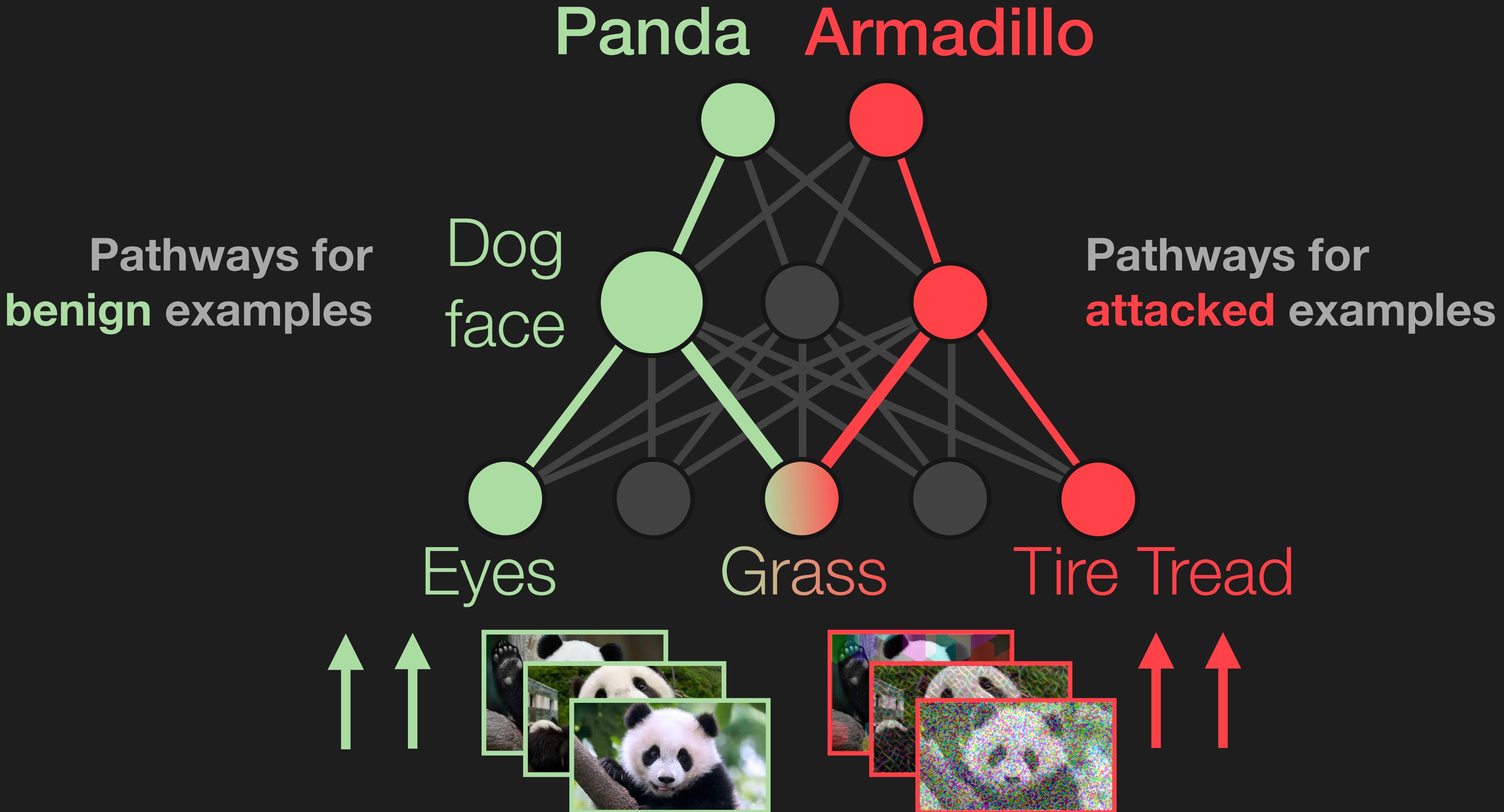
261 384 46



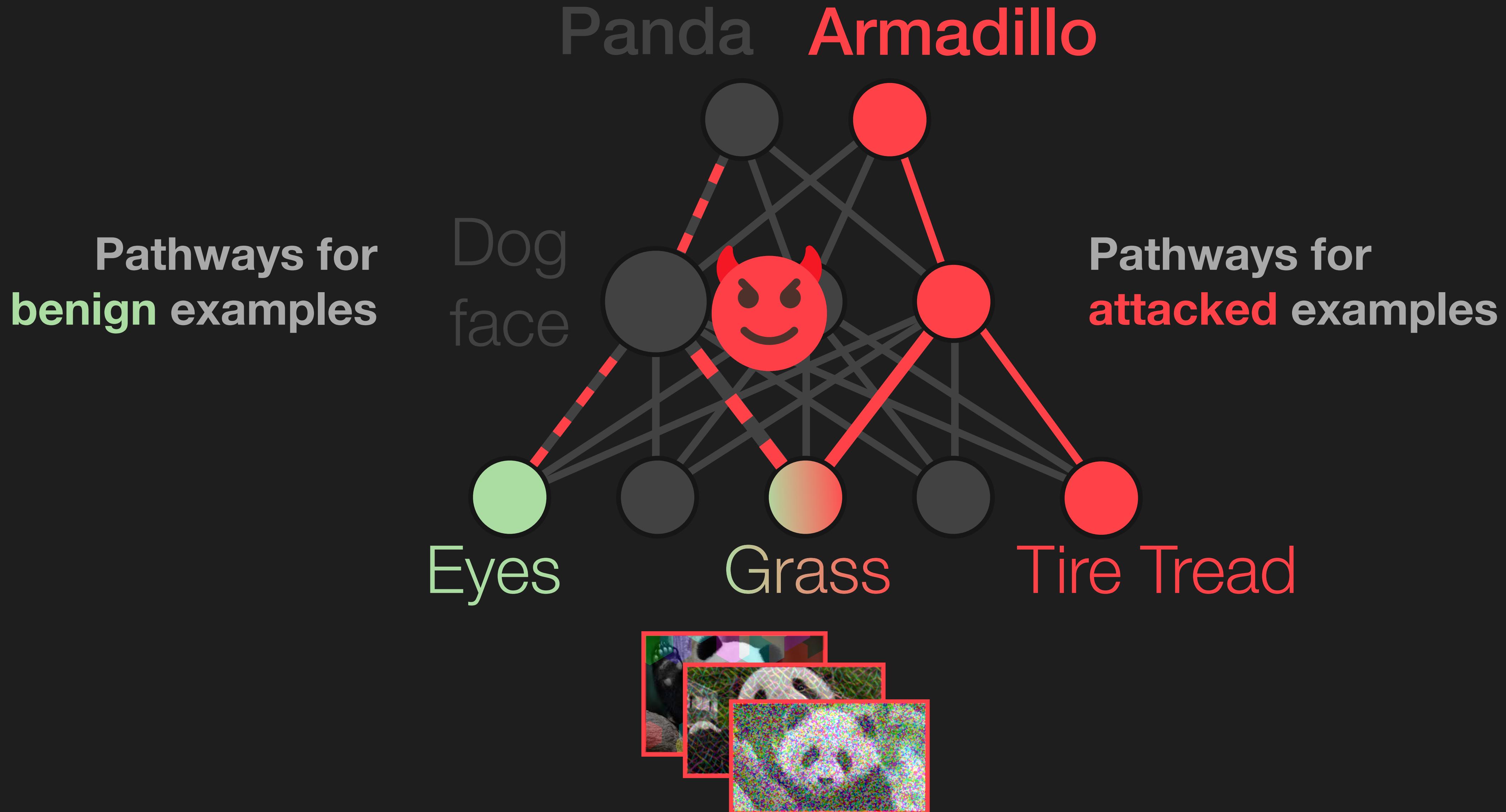
Activation Pathways — most **activated** pathways



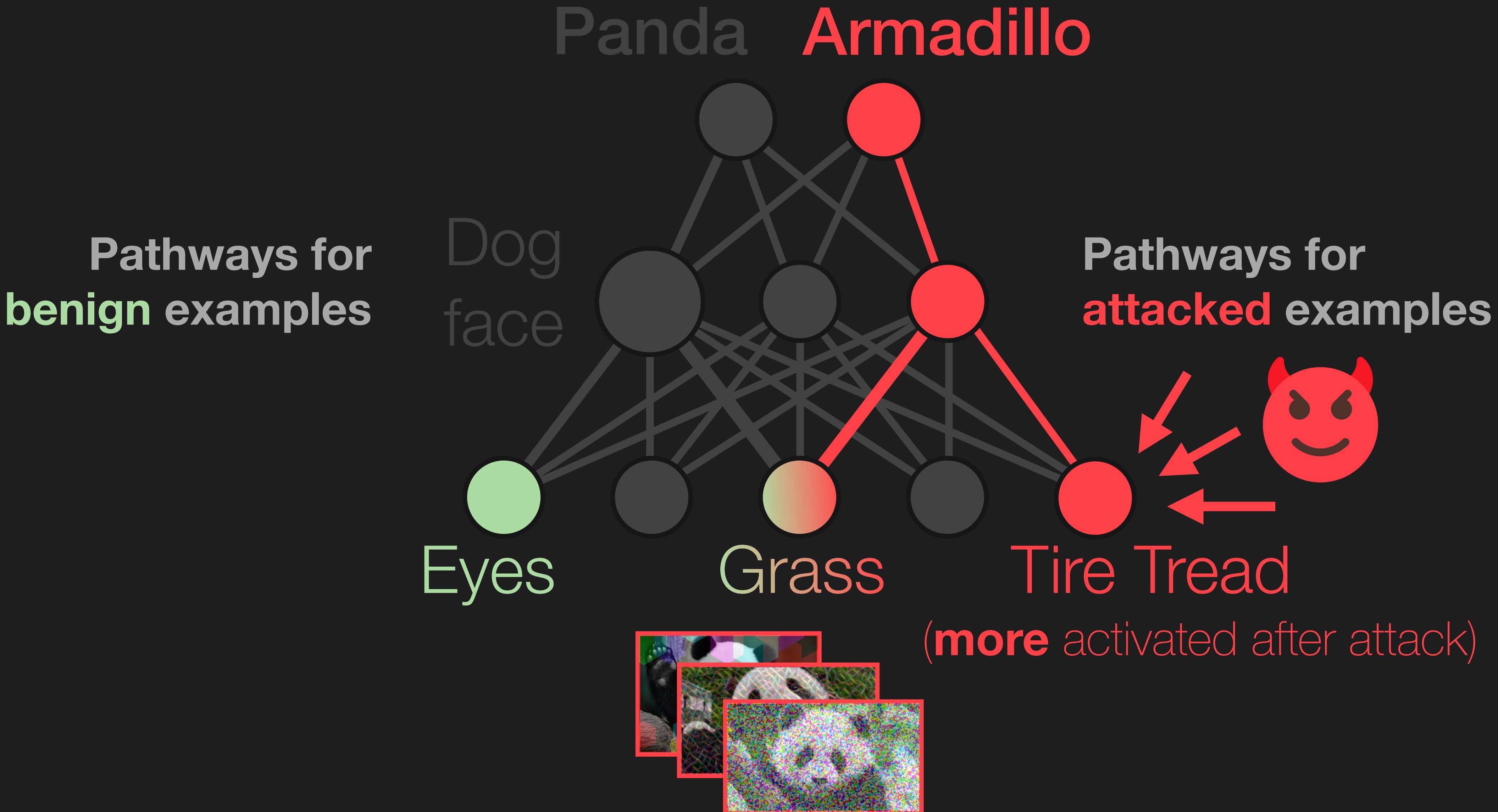
Activation Pathways — most **activated** pathways



Activation Pathways — most inhibited pathways



Activation Pathways — most excited pathways





ADVERSARIAL ATTACK

PGD

Strength: 0.05

FILTER GRAPH

- Show full graph
- Show pinned only
- Show highlighted only

HIGHLIGHT PATHWAYS

Highlight pathways most excited ▾ by attack.

Neurons: top 45 % in each layer

Connections: top 50 %

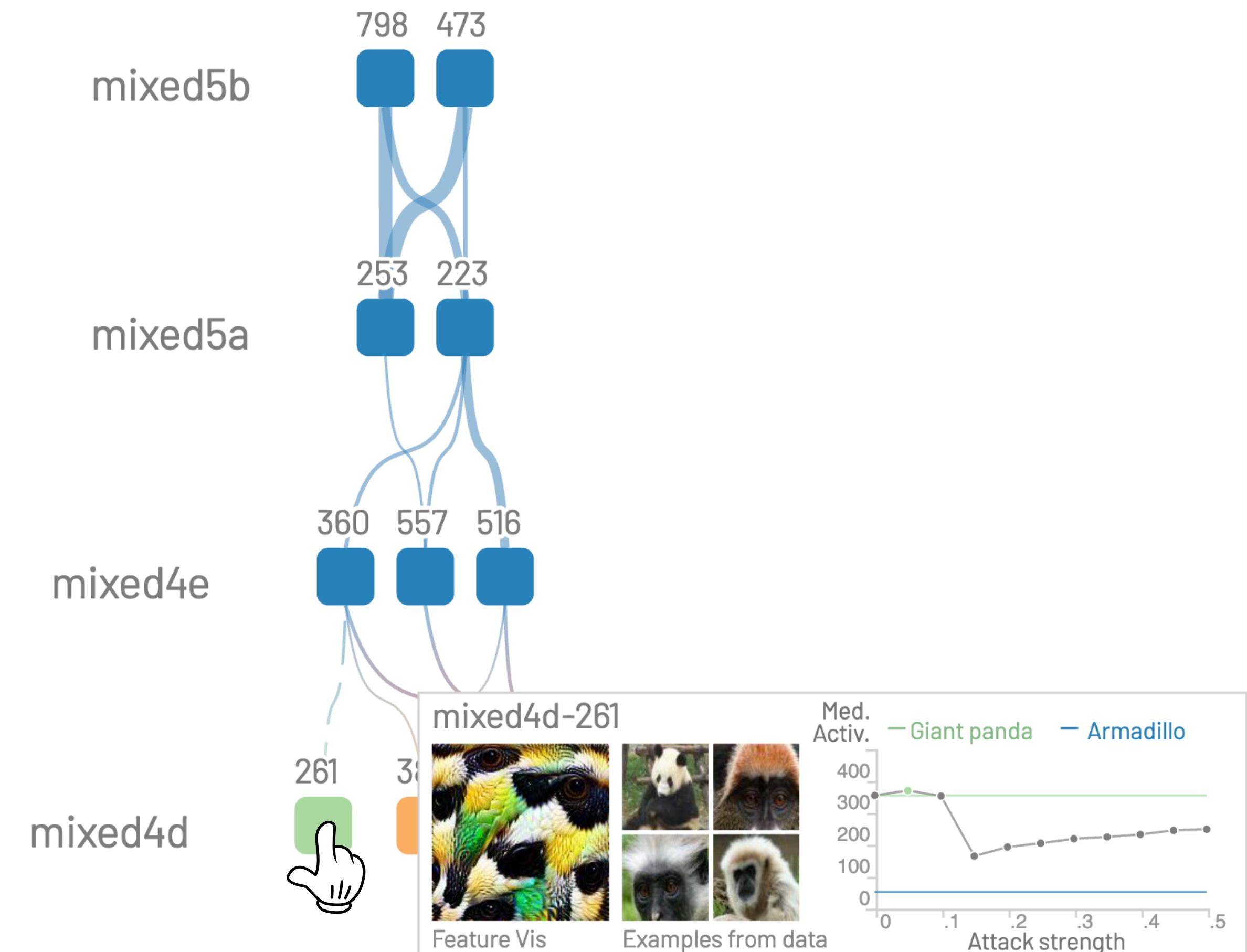
COMPARE ATTACKS

Stronger (outer): 0.45

Weaker (inner): 0.05

Only show edges most excited by weaker ▾ attack

- | | |
|---------------------------------|---|
| <input type="checkbox"/> N/A | <input type="checkbox"/> Stronger |
| <input type="checkbox"/> Weaker | <input checked="" type="checkbox"/> Weaker + Stronger |

GIANT PANDA BOTH ARMADILLO EXPLOITED BY ATTACK



ADVERSARIAL ATTACK

PGD

Strength: 0.05



FILTER GRAPH

 Show full graph Show pinned only Show highlighted only

HIGHLIGHT PATHWAYS

Highlight pathways most excited ▾ by attack.

Neurons: top 45 % in each layer



Connections: top 50 %



COMPARE ATTACKS

Stronger (outer): 0.45



Weaker (inner): 0.05

Only show edges most excited by weaker ▾ attack

- N/A
- Stronger
- Weaker
- Weaker + Stronger

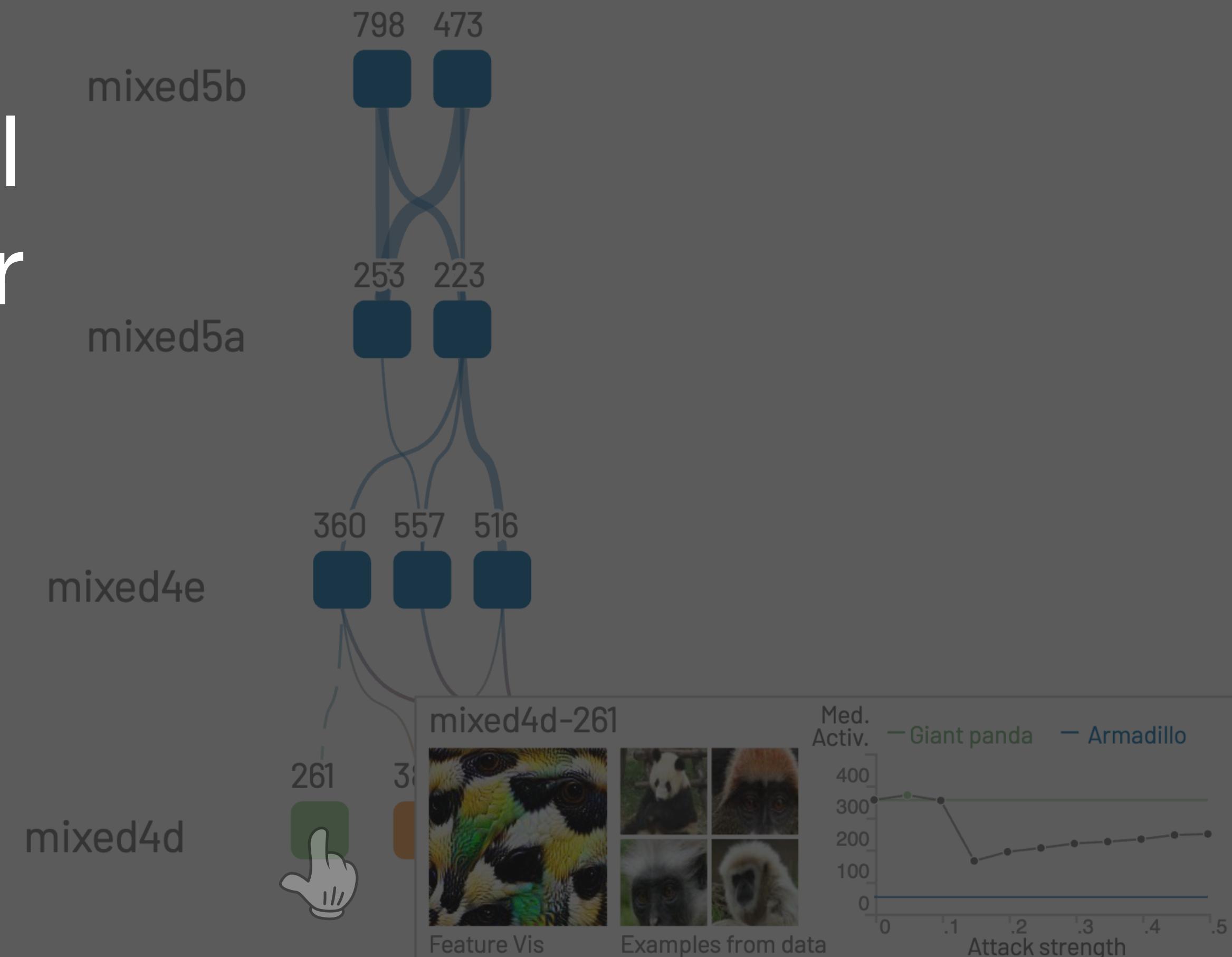
GIANT PANDA

BOTH

ARMADILLO

EXPLOITED BY ATTACK

Control sidebar





ADVERSARIAL ATTACK

PGD

Strength: 0.05

FILTER GRAPH

- Show full graph
- Show pinned only
- Show highlighted only

HIGHLIGHT PATHWAYS

Highlight pathways most excited ▾ by attack.

Neurons: top 45 % in each layer

Connections: top 50 %

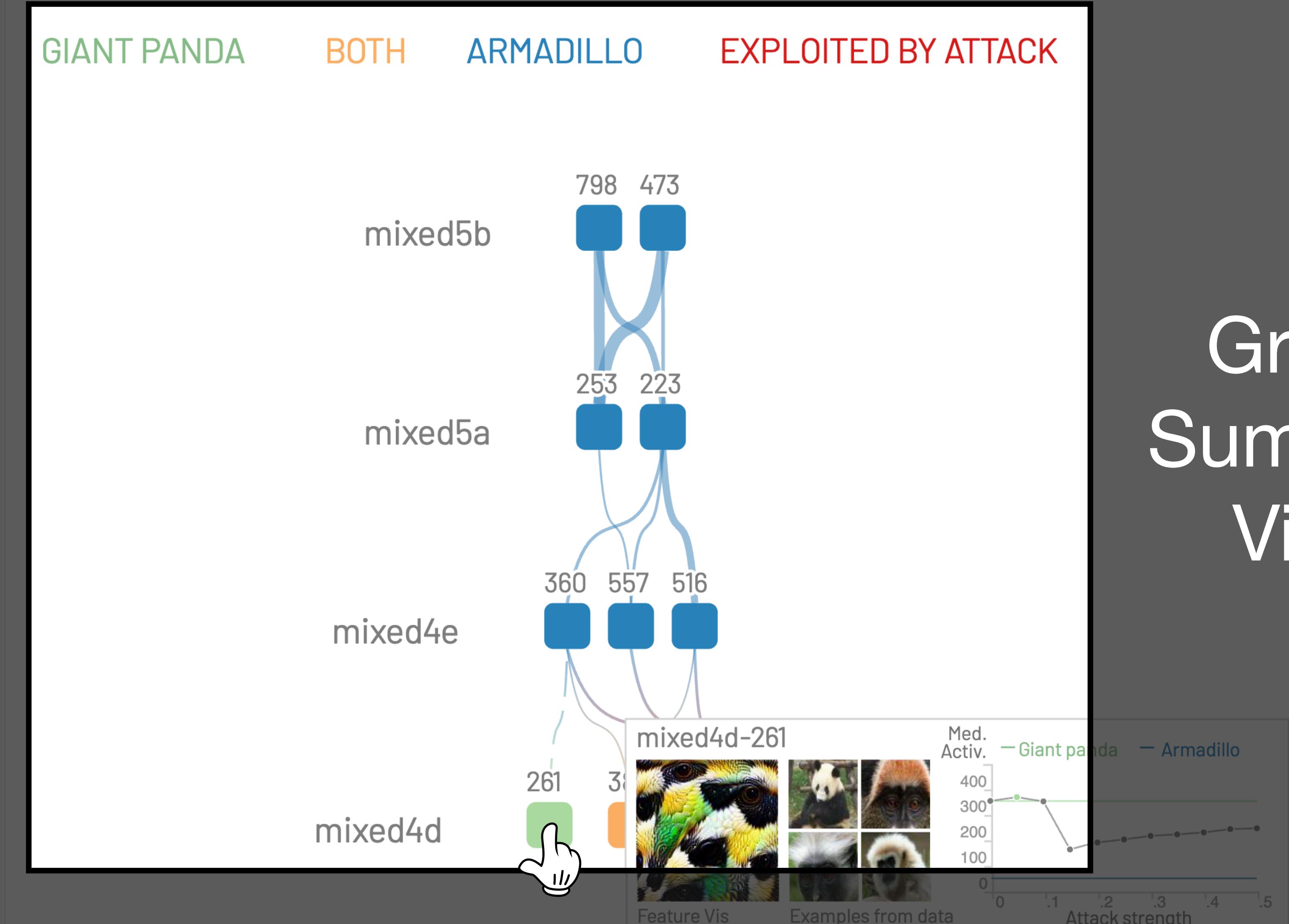
COMPARE ATTACKS

Stronger (outer): 0.45

Weaker (inner): 0.05

Only show edges most excited by weaker ▾ attack

- N/A
- Stronger
- Weaker
- Weaker + Stronger



Graph Summary View



ADVERSARIAL ATTACK

PGD

Strength: 0.05

FILTER GRAPH

- Show full graph
- Show pinned only
- Show highlighted only

HIGHLIGHT PATHWAYS

Highlight pathways most excited ▾ by attack.

Neurons: top 45 % in each layer

Connections: top 50 %

COMPARE ATTACKS

Stronger (outer): 0.45

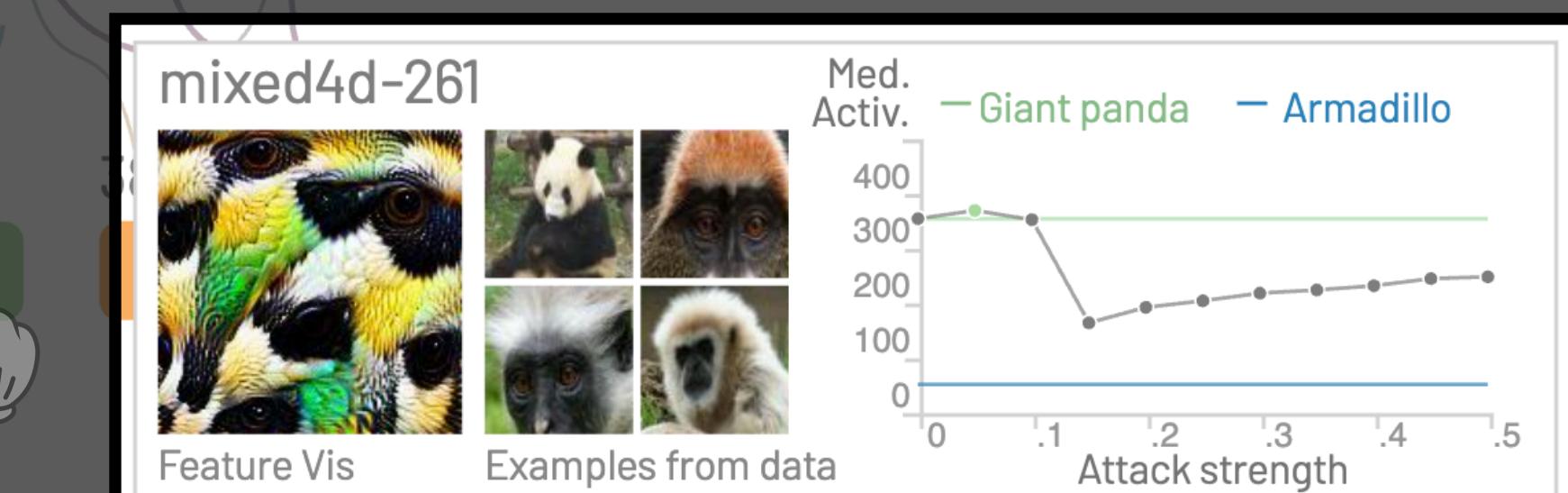
Weaker (inner): 0.05

Only show edges most excited by weaker ▾ attack

- | | |
|--------|-------------------|
| N/A | Stronger |
| Weaker | Weaker + Stronger |

GIANT PANDA BOTH ARMADILLO EXPLOITED BY ATTACK

Detail View





ADVERSARIAL ATTACK

PGD

PIECEWISE GRAPHS

- Show full graph
- Show pinned only
- Show highlighted only

HIGHLIGHT PATHWAYS

Highlight pathways most excited by attack.

Neurons: top 45 % in each layer

Connections: top 50 %

COMPARE ATTACKS

Stronger (outer): 0.45

Weaker (inner): 0.05

Only show edges most excited by weaker attack

- N/A
- Stronger
- Weaker
- Weaker + Stronger

misclassify GIANT PANDA ▾ into ARMADILLO ▾ when attacked

NDA BOTH ARMADILLO

mixed5b

473 798



mixed5a

253 223

Panda

mixed4e

360 557 516

mixed4d

261 384 46



PGD
Attack



Armadillo



ADVERSARIAL ATTACK

PGD

Strength: 0.05

FILTER GRAPH

- Show full graph
- Show pinned only
- Show highlighted only

HIGHLIGHT PATHWAYS

Highlight pathways most excited ▾ by attack.

Neurons: top 45 % in each layer

Connections: top 50 %

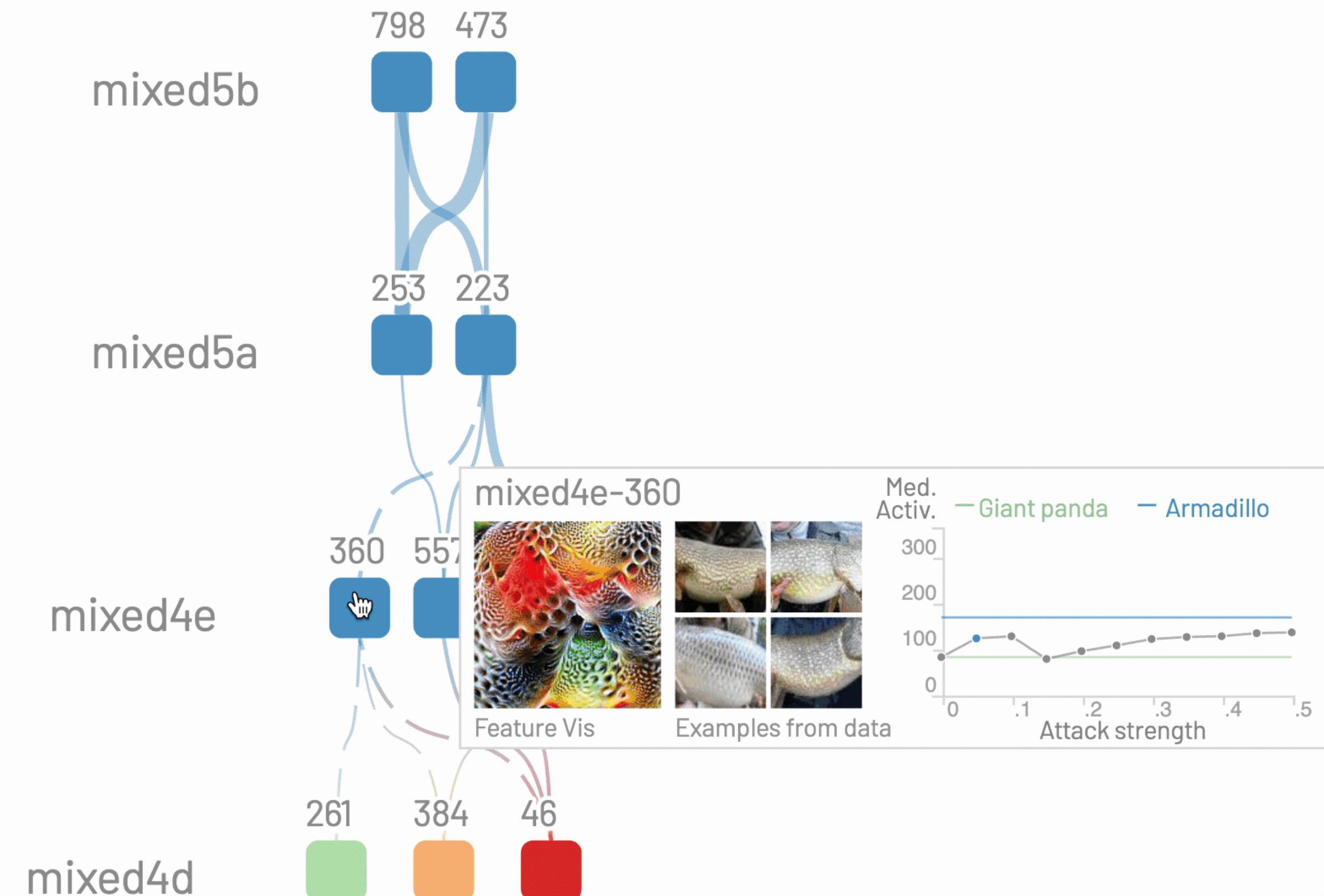
COMPARE ATTACKS

Stronger (outer): 0.45

Weaker (inner): 0.05

Only show edges most excited by weaker ▾ attack

- N/A
- Stronger
- Weaker
- Weaker + Stronger

GIANT PANDA BOTH ARMADILLO EXPLOITED BY ATTACK



ADVERSARIAL ATTACK

PGD

Strength: 0.05



FILTER GRAPH

 Show full graph Show pinned only Show highlighted only

HIGHLIGHT PATHWAYS

Highlight pathways most excited [▼](#) by attack.

Neurons: top 45 % in each layer



Connections: top 50 %

COMPARE ATTACKS [○](#)

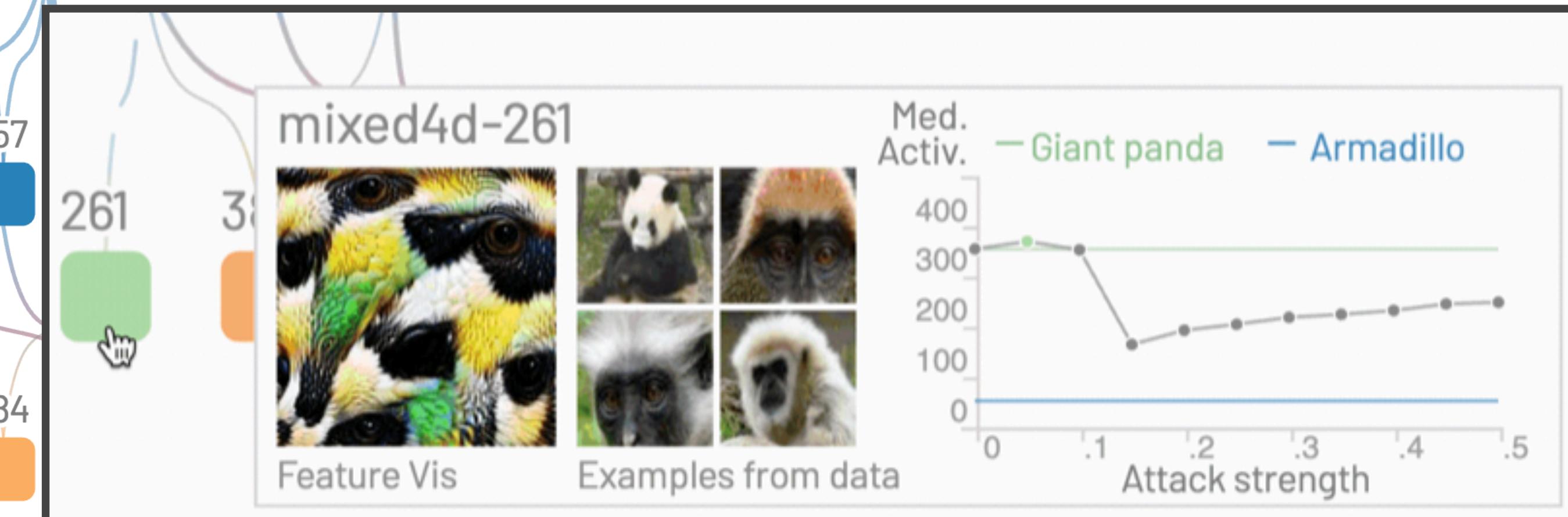
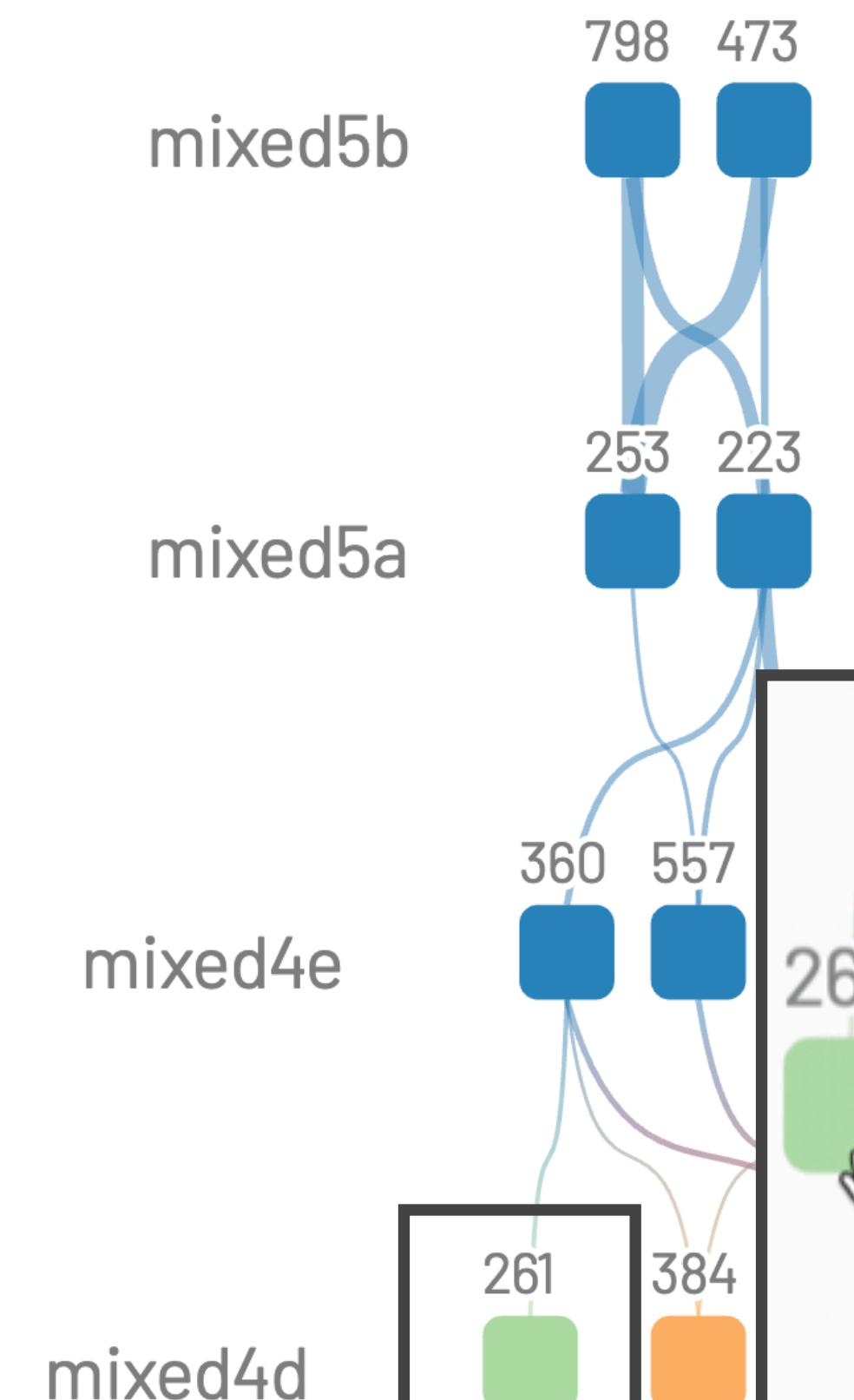
Stronger (outer): 0.45



Weaker (inner): 0.05

Only show edges most excited by weaker [▼](#) attack

- N/A
- Stronger
- Weaker
- Weaker + Stronger

GIANT PANDA**BOTH****ARMADILLO****EXPLOITED BY ATTACK**



ADVERSARIAL ATTACK

PGD

Strength: 0.05



FILTER GRAPH

 Show full graph Show pinned only Show highlighted only

HIGHLIGHT PATHWAYS

Highlight pathways most excited ▾ by attack.

Neurons: top 45 % in each layer



Connections: top 50 %



COMPARE ATTACKS

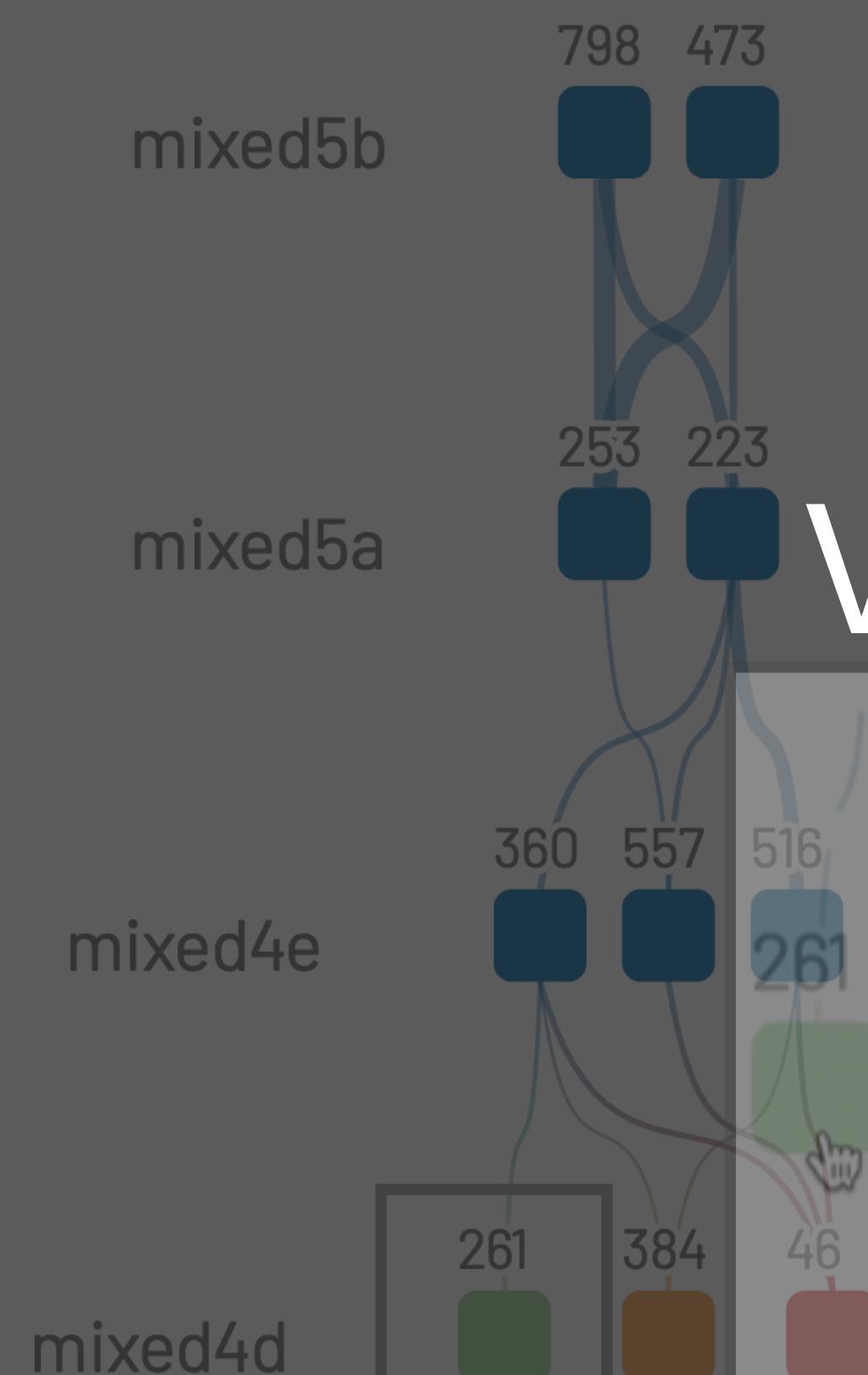
Stronger (outer): 0.45



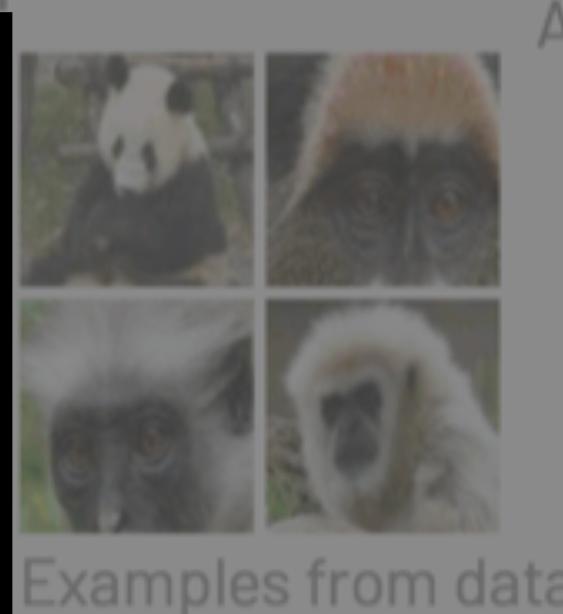
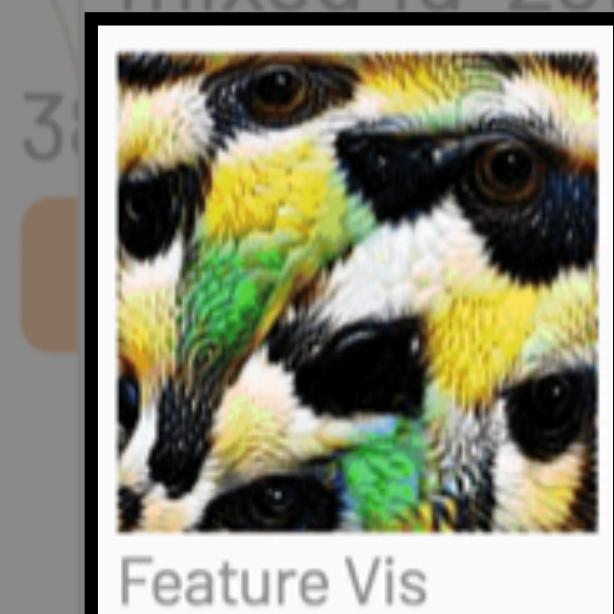
Weaker (inner): 0.05

Only show edges most excited by weaker ▾ attack

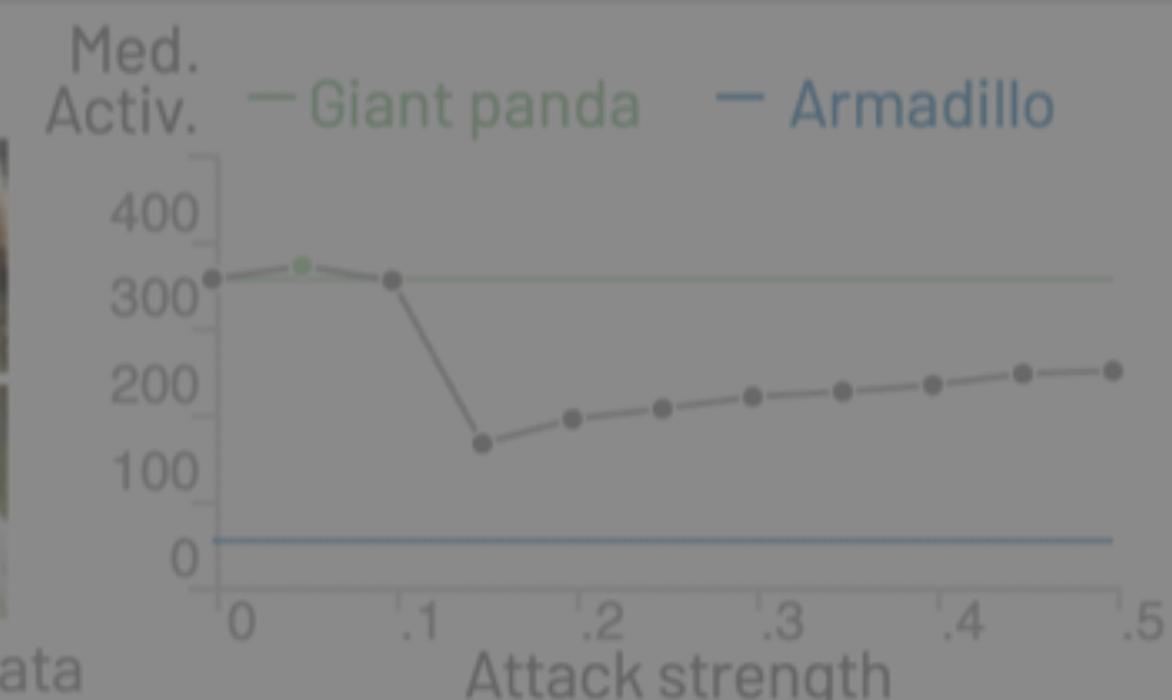
- N/A
- Stronger
- Weaker
- Weaker + Stronger

GIANT PANDA**BOTH****ARMADILLO****EXPLOITED BY ATTACK**

Feature Visualization



Examples from data





ADVERSARIAL ATTACK

PGD

Strength: 0.05



FILTER GRAPH

 Show full graph Show pinned only Show highlighted only

HIGHLIGHT PATHWAYS

Highlight pathways most excited ▾ by attack.

Neurons: top 45 % in each layer



Connections: top 50 %



COMPARE ATTACKS

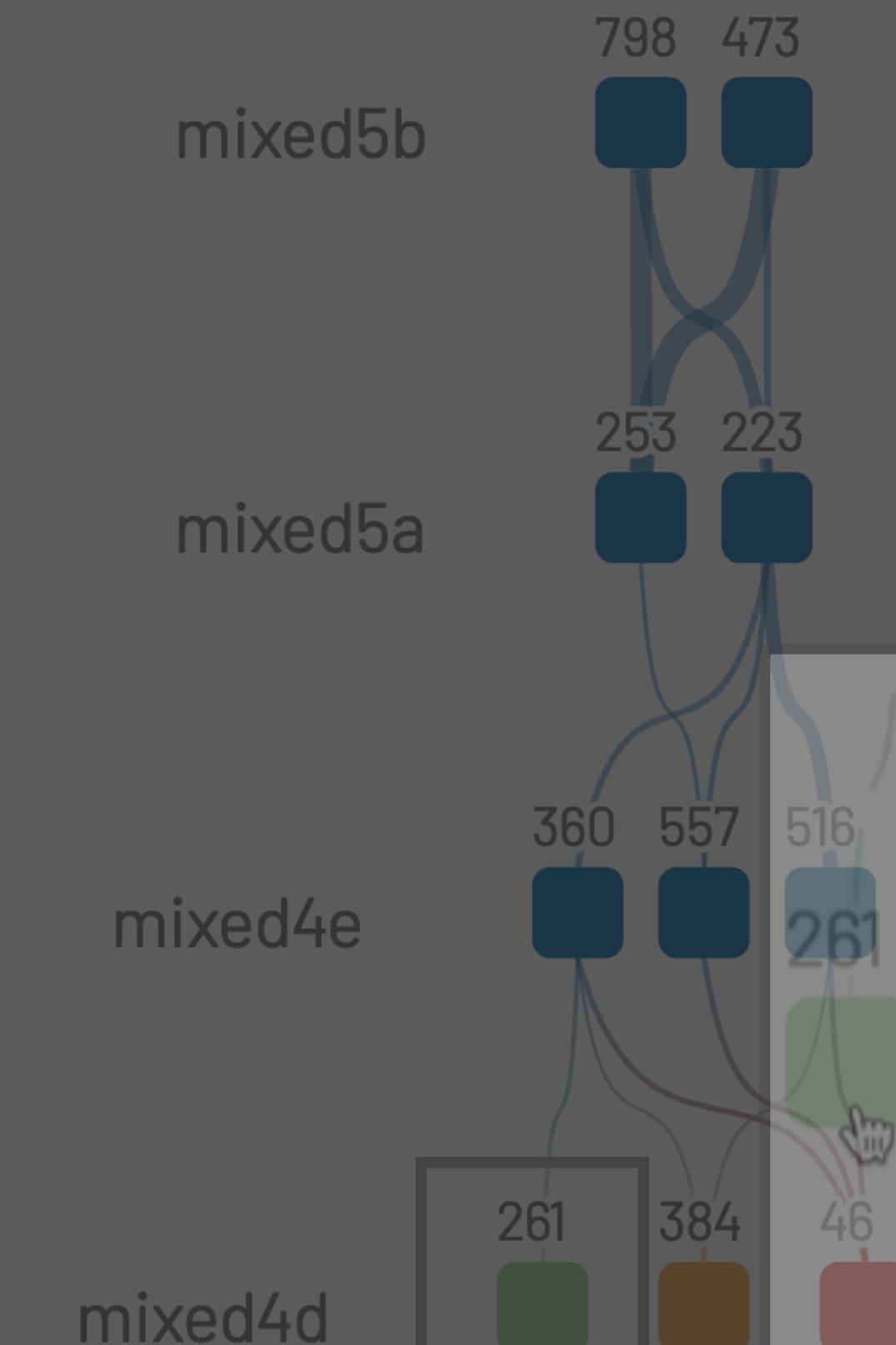
Stronger (outer): 0.45



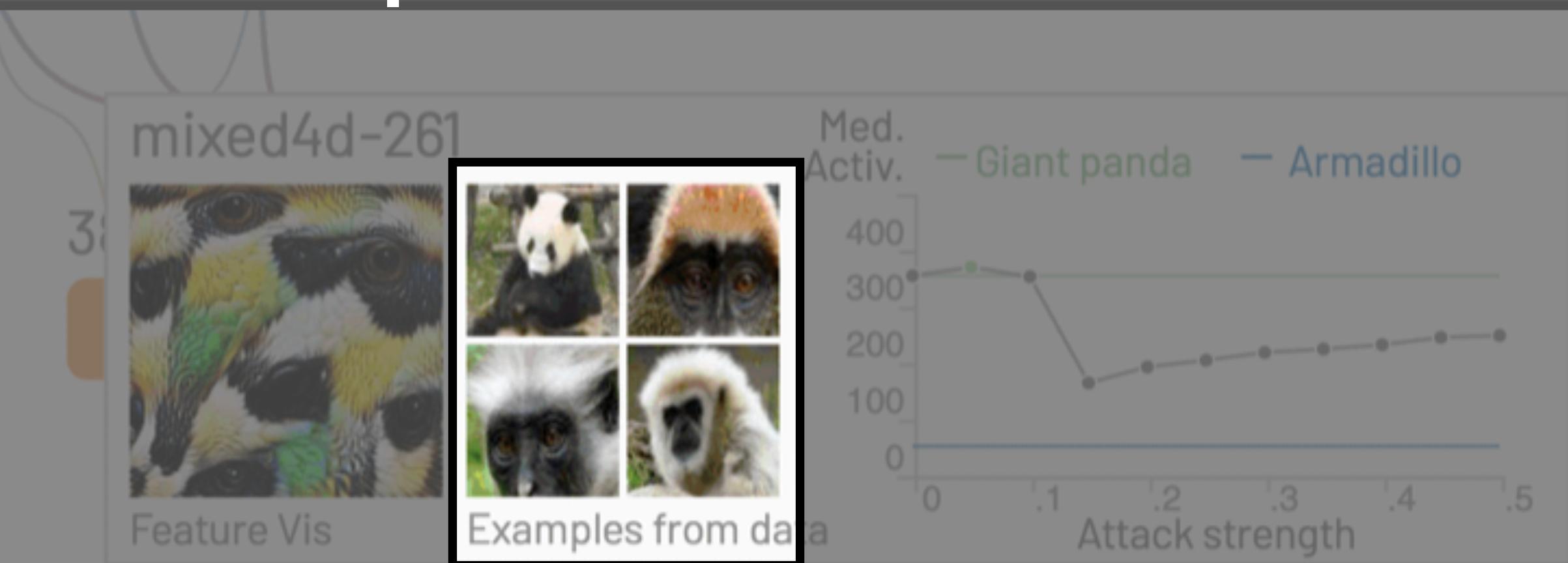
Weaker (inner): 0.05

Only show edges most excited by weaker ▾ attack

- N/A
- Stronger
- Weaker
- Weaker + Stronger

GIANT PANDA**BOTH****ARMADILLO****EXPLOITED BY ATTACK**

Example patches





ADVERSARIAL ATTACK

PGD

Strength: 0.05



FILTER GRAPH

 Show full graph Show pinned only Show highlighted only

HIGHLIGHT PATHWAYS

Highlight pathways most excited ▾ by attack.

Neurons: top 45 % in each layer



Connections: top 50 %



COMPARE ATTACKS

Stronger (outer): 0.45



Weaker (inner): 0.05



Only show edges most excited by weaker ▾ attack

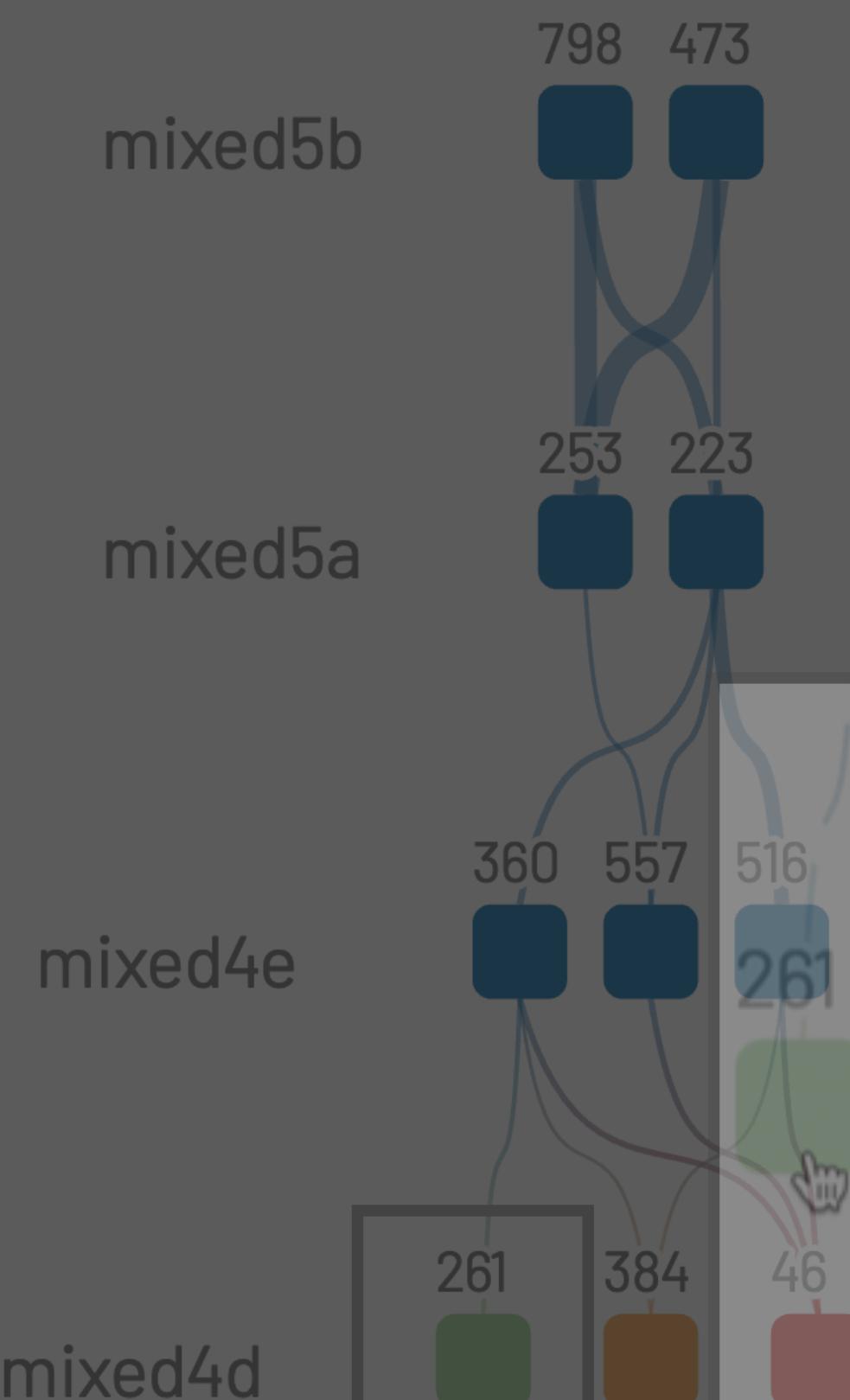


N/A

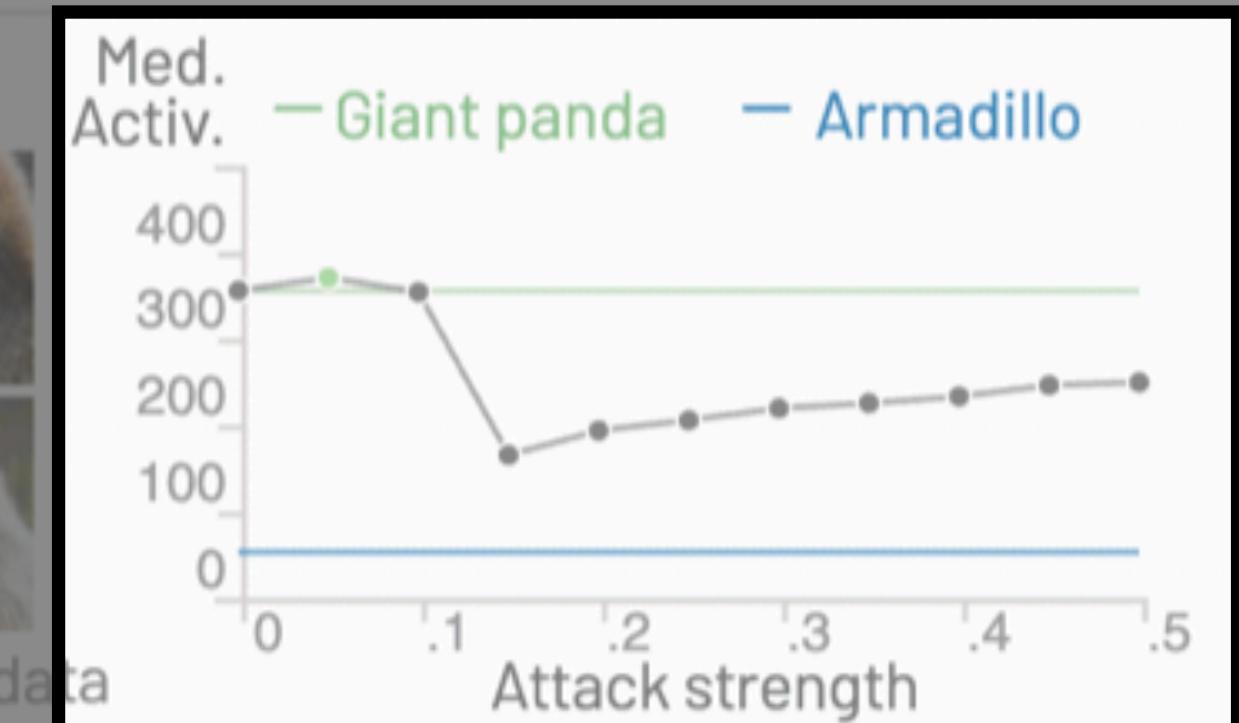
Stronger

Weaker

Weaker + Stronger

GIANT PANDA**BOTH****ARMADILLO****EXPLOITED BY ATTACK**

Median Activation across attack strength





ADVERSARIAL ATTACK

PGD

Strength: 0.05

FILTER GRAPH

- Show full graph
- Show pinned only
- Show highlighted only

HIGHLIGHT PATHWAYS

Highlight pathways most excited by attack.

Neurons: top 45 % in each layer

Connections: top 50 %

COMPARE ATTACKS

Stronger (outer): 0.45

Weaker (inner): 0.05

Only show edges most excited by weaker attack

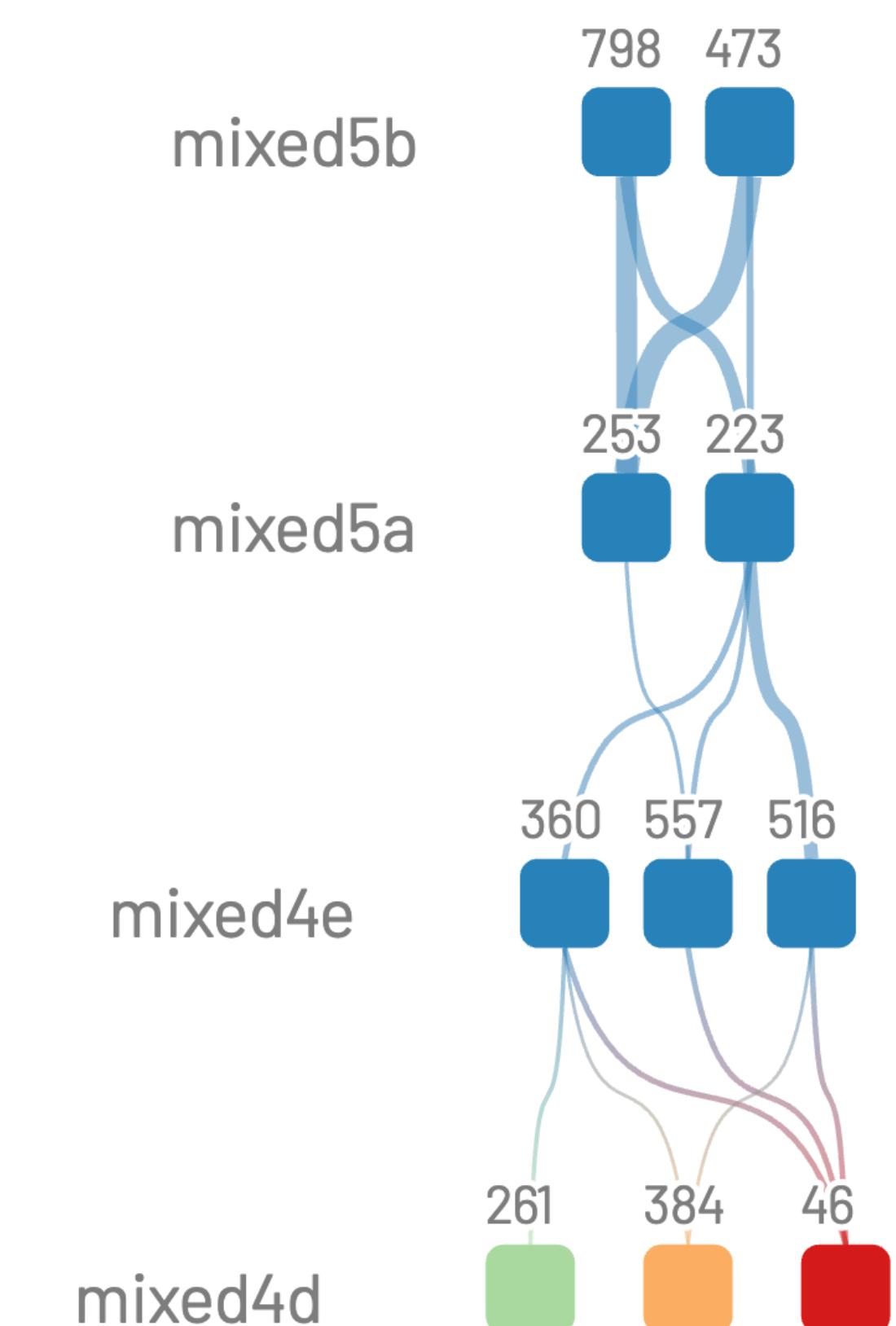
- | | |
|--|--|
| <input type="checkbox"/> N/A | <input type="checkbox"/> Stronger |
| <input checked="" type="checkbox"/> Weaker | <input type="checkbox"/> Weaker + Stronger |

GIANT PANDA

BOTH

ARMADILLO

EXPLOITED BY ATTACK





ADVERSARIAL ATTACK

PGD

Strength: 0.05



FILTER GRAPH

 Show full graph Show pinned only Show highlighted only

HIGHLIGHT PATHWAYS

Highlight pathways most excited ▾ by attack.

Neurons: top 45 % in each layer



Connections: top 50 %



COMPARE ATTACKS

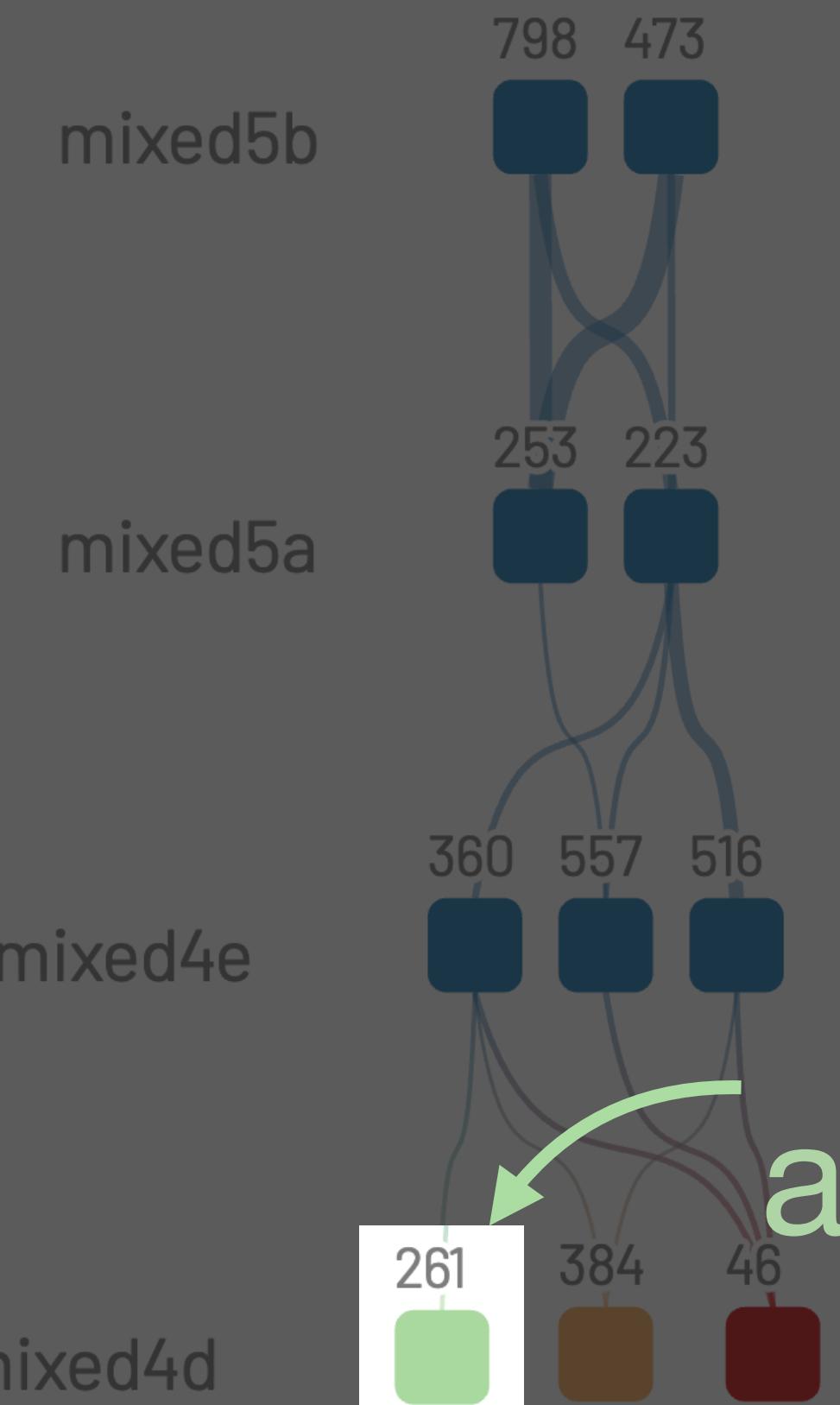
Stronger (outer): 0.45



Weaker (inner): 0.05

Only show edges most excited by weaker ▾ attack

- | | |
|---------------------------------|---|
| <input type="checkbox"/> N/A | <input type="checkbox"/> Stronger |
| <input type="checkbox"/> Weaker | <input checked="" type="checkbox"/> Weaker + Stronger |

GIANT PANDA**BOTH****ARMADILLO****EXPLOITED BY ATTACK****(Original Class)**

highly
activated by

Benign
Giant Panda
Images



ADVERSARIAL ATTACK

PGD

Strength: 0.05



FILTER GRAPH

 Show full graph Show pinned only Show highlighted only

HIGHLIGHT PATHWAYS

Highlight pathways most excited ▾ by attack.

Neurons: top 45 % in each layer



Connections: top 50 %



COMPARE ATTACKS

Stronger (outer): 0.45



Weaker (inner): 0.05

Only show edges most excited by weaker ▾ attack

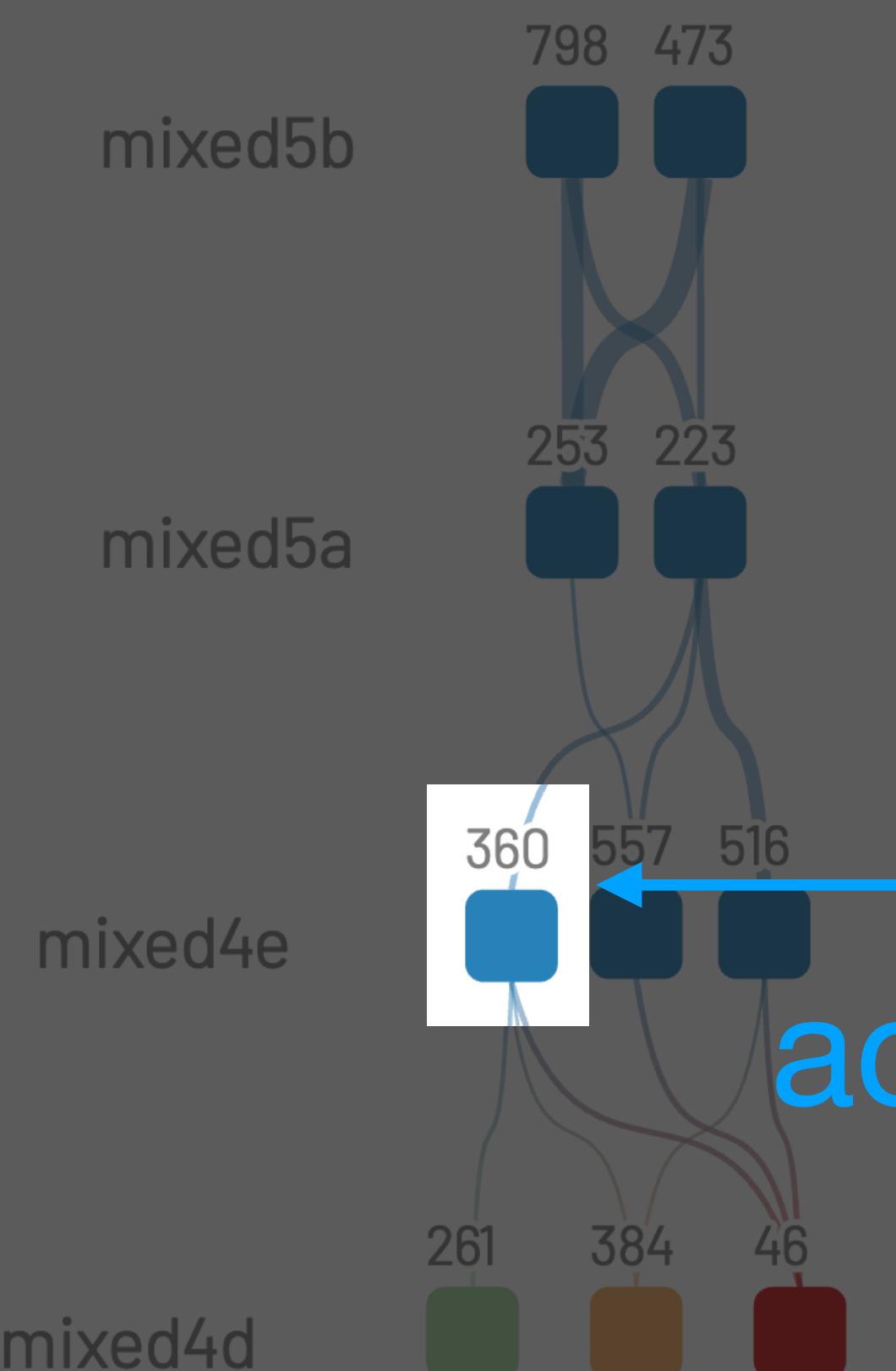
- N/A
- Stronger
- Weaker
- Weaker + Stronger

GIANT PANDA

BOTH

ARMADILLO

EXPLOITED BY ATTACK



highly
activated by

(Target Class)



Benign
Armadillo
Images



ADVERSARIAL ATTACK

PGD

Strength: 0.05

FILTER GRAPH

- Show full graph
- Show pinned only
- Show highlighted only

HIGHLIGHT PATHWAYS

Highlight pathways most excited by attack.

Neurons: top 45 % in each layer

Connections: top 50 %

COMPARE ATTACKS

Stronger (outer): 0.45

Weaker (inner): 0.05

Only show edges most excited by weaker attack

- N/A
- Stronger
- Weaker
- Weaker + Stronger

GIANT PANDA

BOTH

ARMADILLO

ATTACK

mixed5b

798 473



mixed5a

253 223

Benign
Giant Panda
Images

mixed4e

360 557 516

mixed4d

261 384 46

highly activated by Both class
images



Benign
Armadillo
Images



ADVERSARIAL ATTACK

PGD

Strength: 0.05

FILTER GRAPH

- Show full graph
- Show pinned only
- Show highlighted only

HIGHLIGHT PATHWAYS

Highlight pathways most excited ▾ by attack.

Neurons: top 45 % in each layer

Connections: top 50 %

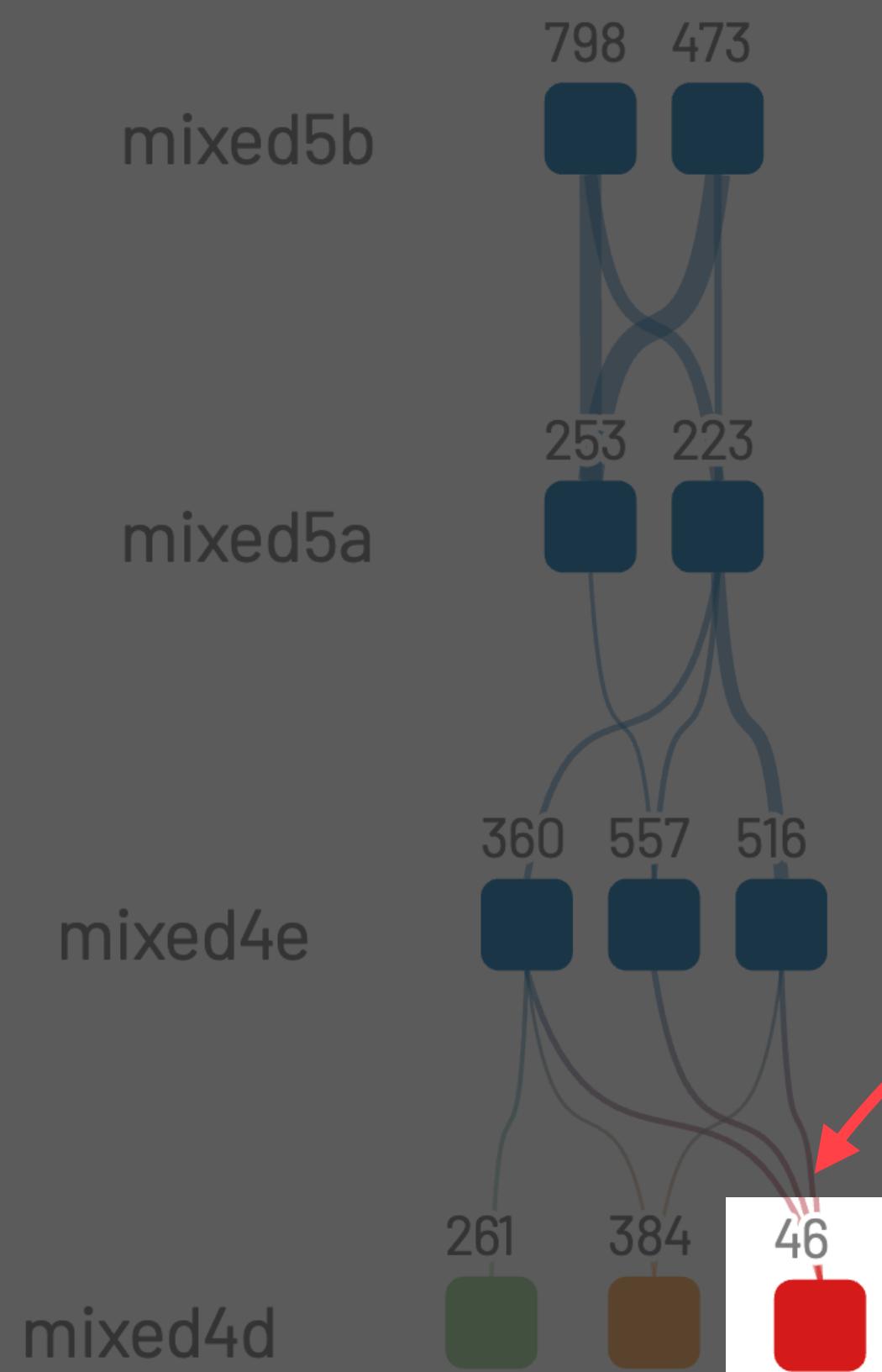
COMPARE ATTACKS

Stronger (outer): 0.45

Weaker (inner): 0.05

Only show edges most excited by weaker ▾ attack

- | | |
|---------------------------------|---|
| <input type="checkbox"/> N/A | <input type="checkbox"/> Stronger |
| <input type="checkbox"/> Weaker | <input checked="" type="checkbox"/> Weaker + Stronger |

GIANT PANDA **BOTH** **ARMADILLO****EXPLOITED BY ATTACK****(Attacked images)**

highly activated by **Adversarial Giant Panda Images**



ADVERSARIAL ATTACK

PGD

Strength: 0.05

FILTER GRAPH

- Show full graph
- Show pinned only
- Show highlighted only

HIGHLIGHT PATHWAYS

Highlight pathways most excited ▾ by attack.

Neurons: top 45 % in each layer

Connections: top 50 %

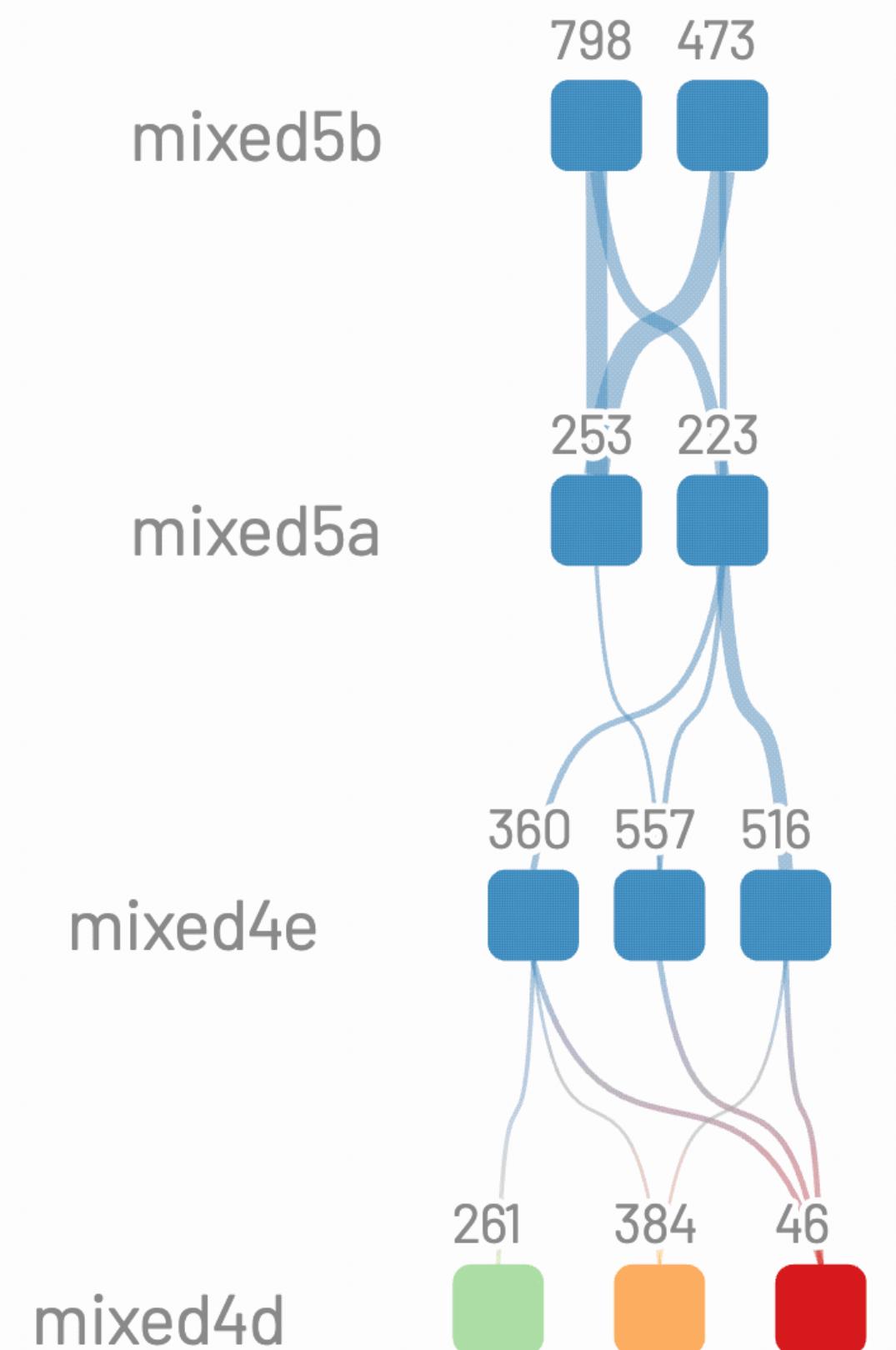
COMPARE ATTACKS

Stronger (outer): 0.45

Weaker (inner): 0.05

Only show edges most excited by weaker ▾ attack

- N/A
- Stronger
- Weaker
- Weaker + Stronger

GIANT PANDA BOTH ARMADILLO EXPLOITED BY ATTACK



ADVERSARIAL ATTACK

PGD

Strength: 0.05

FILTER GRAPH

 Show full graph Show pinned only Show highlighted only

HIGHLIGHT PATHWAYS

Highlight pathways most excited ▾ by attack.

Neurons: top 45 % in each layer

Connections: top 50 %

COMPARE ATTACKS

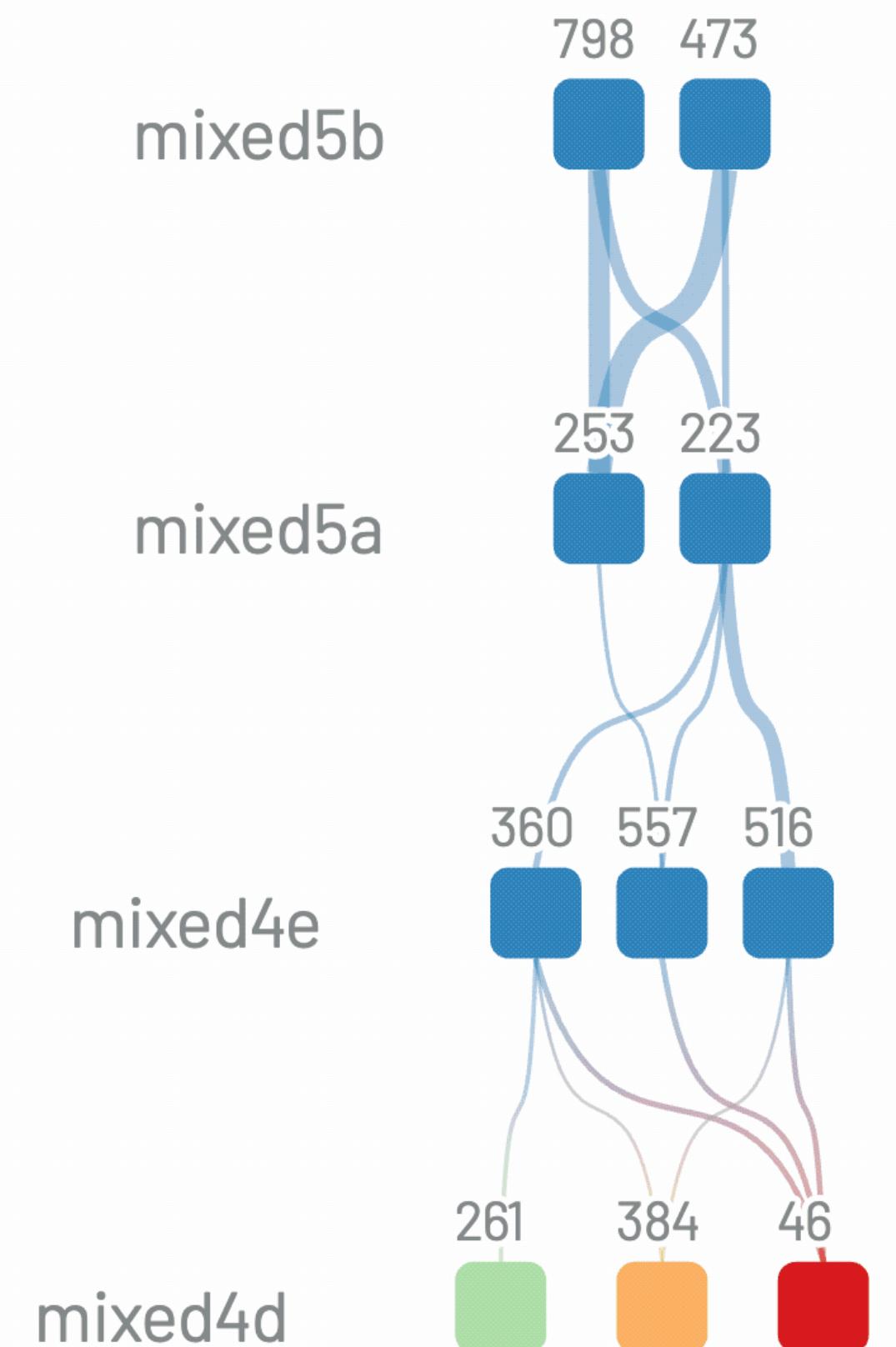
Stronger (outer): 0.45



Weaker (inner): 0.05

Only show edges most excited by weaker ▾ attack

- | | |
|---------------------------------|--|
| <input type="checkbox"/> N/A | <input type="checkbox"/> Stronger |
| <input type="checkbox"/> Weaker | <input type="checkbox"/> Weaker + Stronger |

GIANT PANDA BOTH ARMADILLO EXPLOITED BY ATTACK

Bluff

Understand how neural networks misclassify GIANT PANDA into ARMADILLO when attacked

A Control Sidebar

ADVERSARIAL ATTACK

PGD

Strength: 0.05

FILTER GRAPH

Show full graph

Show pinned only

Show highlighted only

HIGHLIGHT PATHWAYS

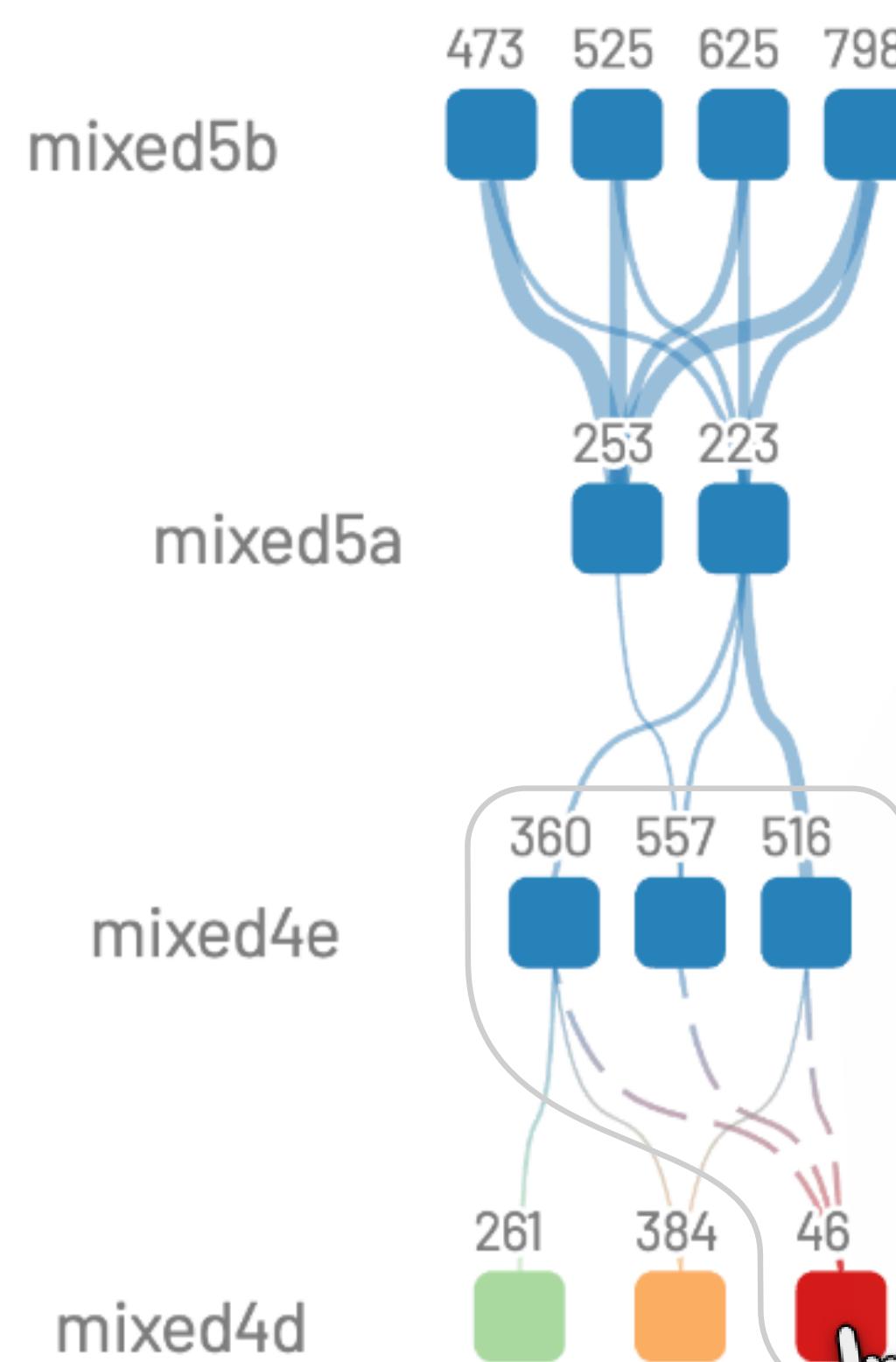
Highlight pathways most excited by attack.

Neurons: top 45 % in each layer

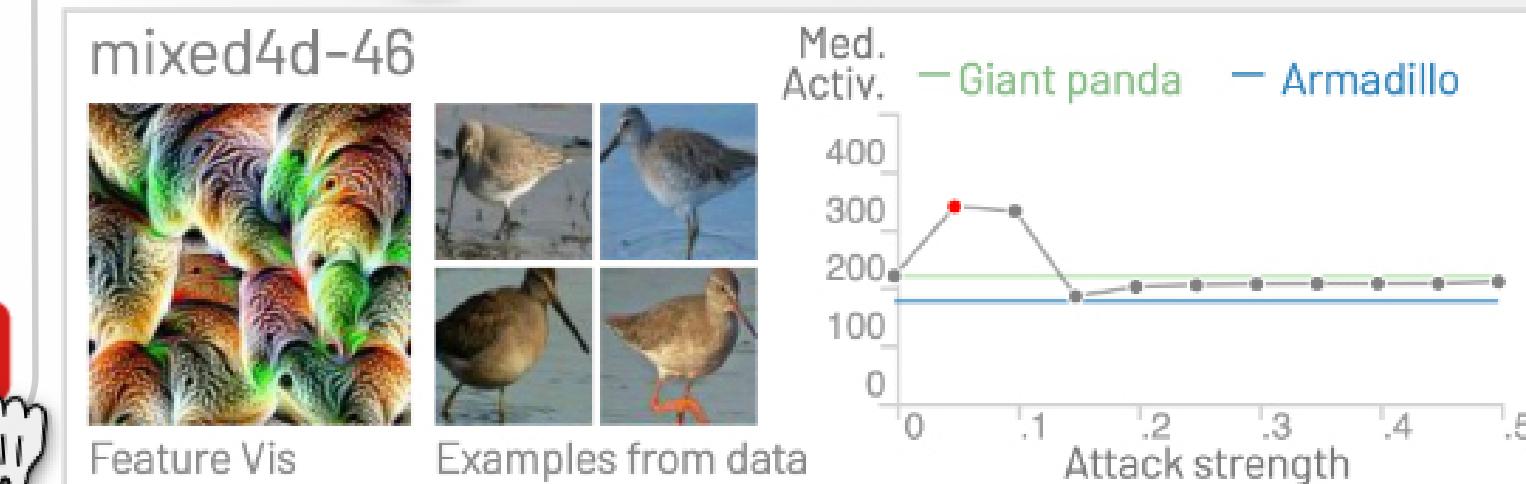
Connections: top 50 %

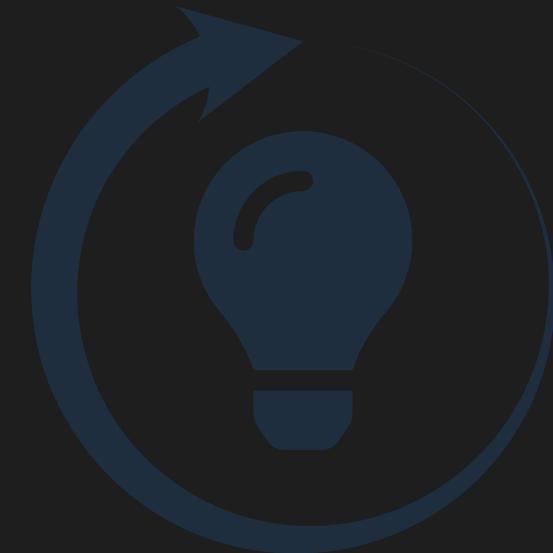
B Graph Summary View

GIANT PANDA BOTH ARMADILLO EXPLOITED BY ATTACK



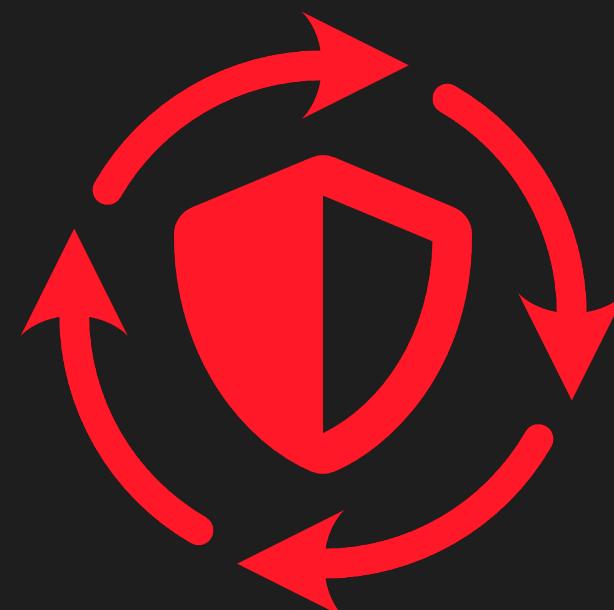
C Detail View





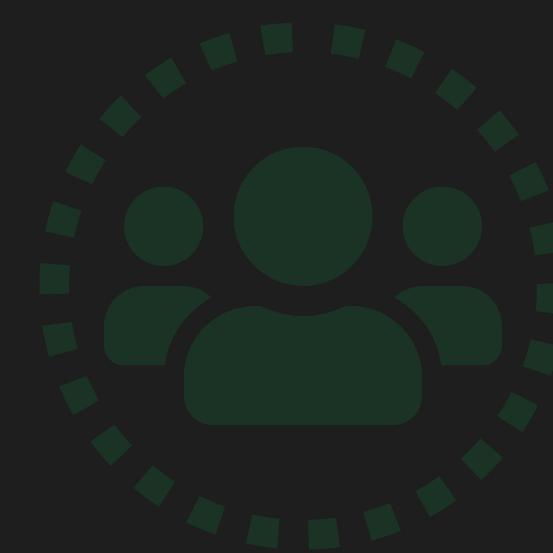
Part I
Understand
AI Vulnerabilities

GOGGLES SIGMOD 2020
Bluff IEEE VIS 2020



Part II
Fortify
AI Security

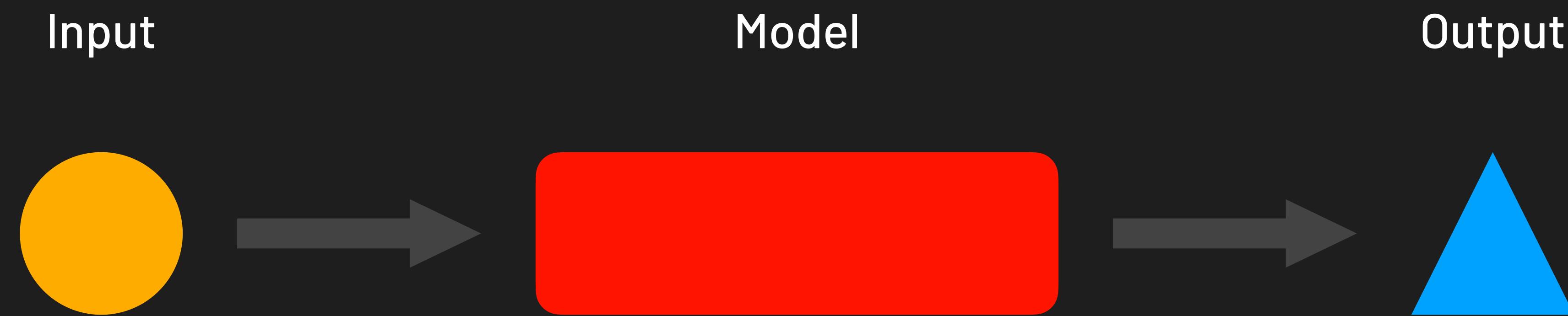
SHIELD KDD 2018
SkeleVision arXiv 2022 (under review)
Hear No Evil arXiv 2022 (under review)



Part III
Enable
Use of AI Security

ADAGIO ECML-PKDD 2018
MLsploit KDD Showcase 2019

Machine Learning Pipeline



Automatic Speech Recognition

Machine Learning Pipeline

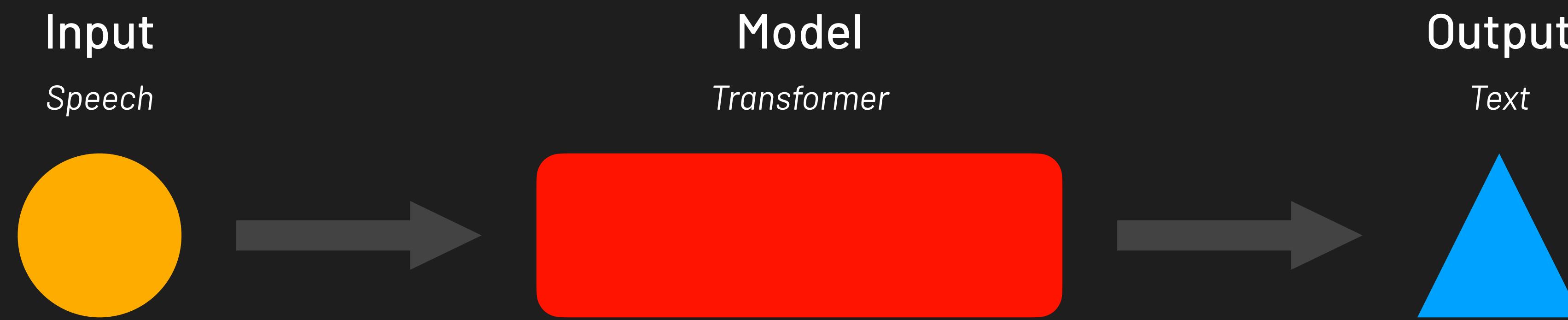
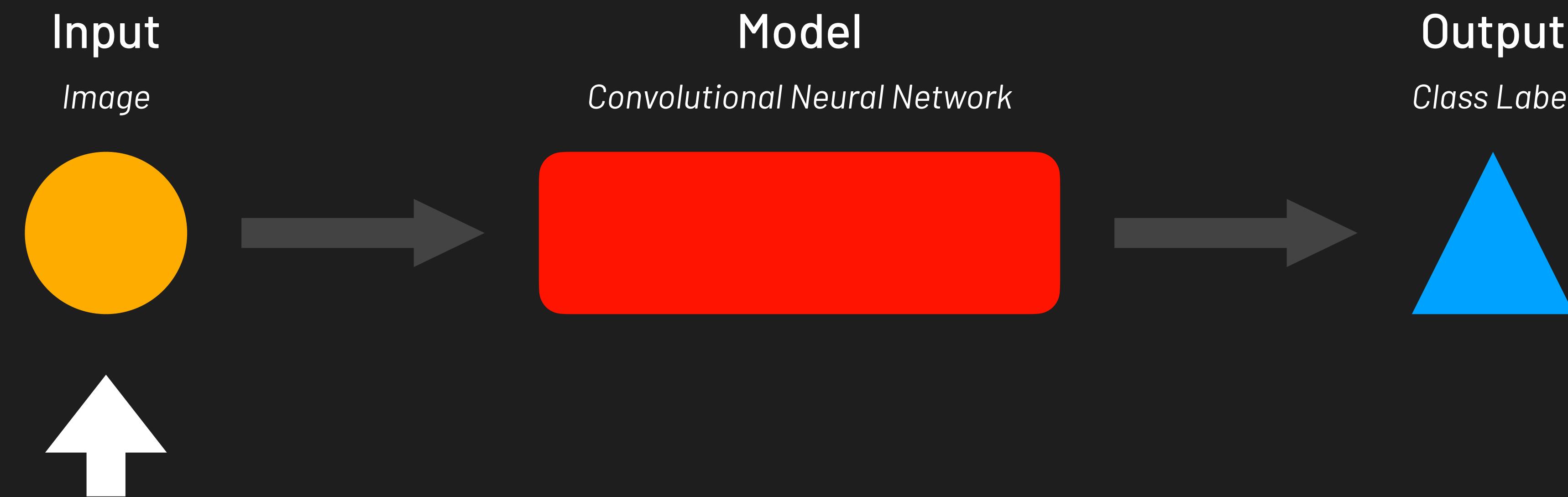


Image Classification

Machine Learning Pipeline



SHIELD

KDD 2018

KDD 2019 LEMINCS

Fast, Practical Defense for Image Classification



KDD'18 Audience Appreciation Award (runner-up)



Tech-transferred to Intel Labs



Open-sourced at github.com/poloclub/jpeg-defense



Nilaksh Das
Georgia Tech



**Madhuri
Shanbhogue**
Georgia Tech



**Shang-Tse
Chen**
Georgia Tech



**Fred
Hohman**
Georgia Tech



**Siwei
Li**
Georgia Tech



**Cory
Cornelius**
Intel Labs



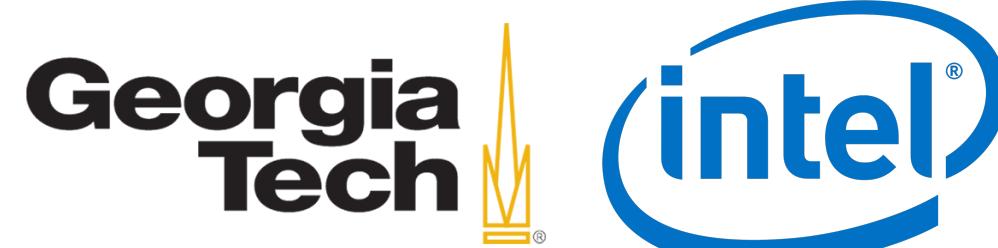
**Li
Chen**
Intel Labs



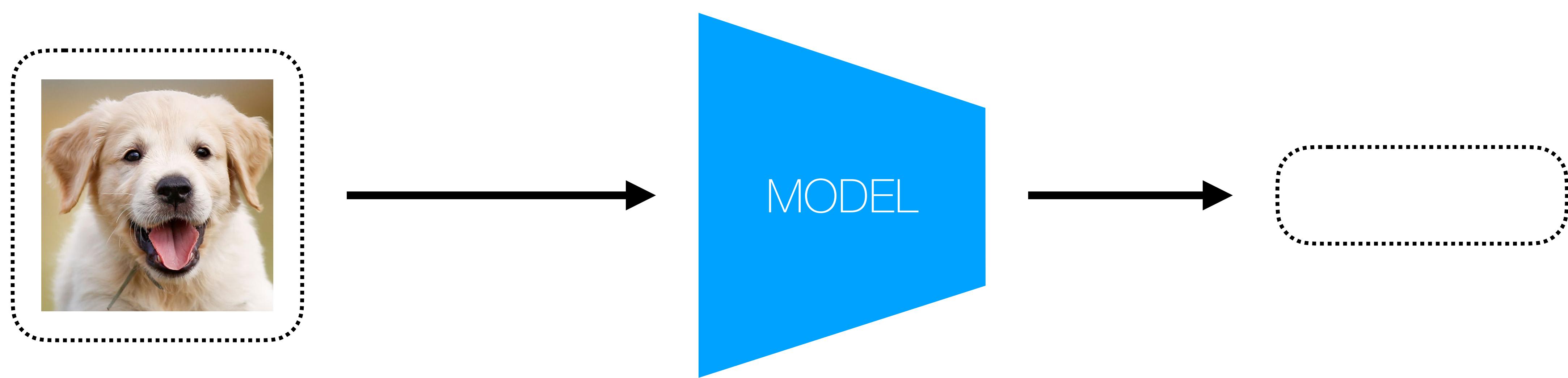
**Michael
Kounavis**
Intel Labs



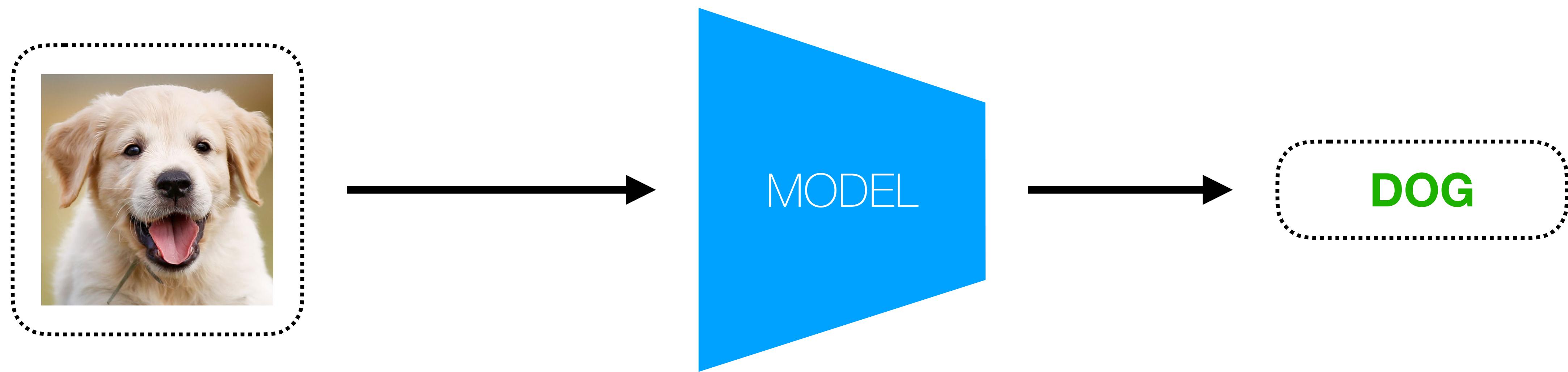
**Polo
Chau**
Georgia Tech



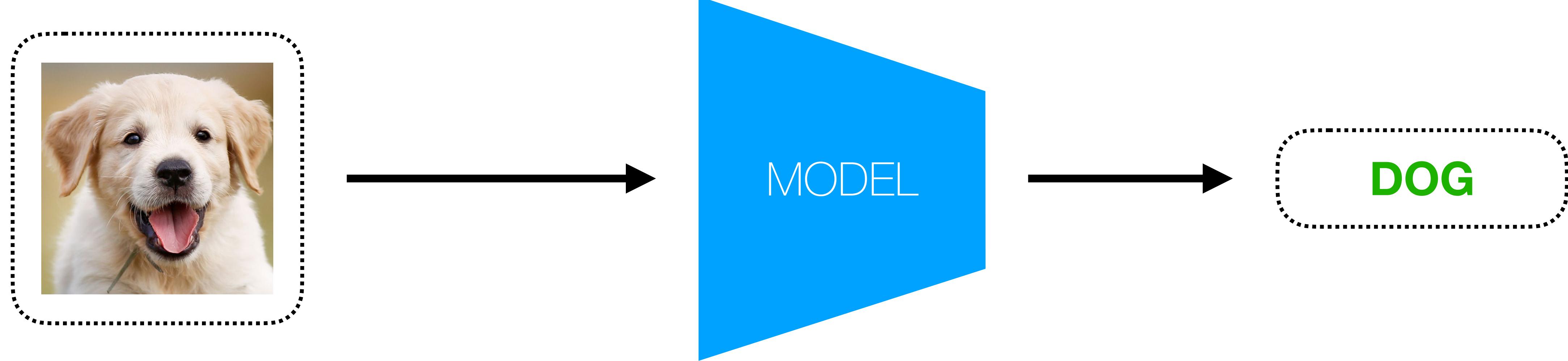
Deep Learning for Image Classification



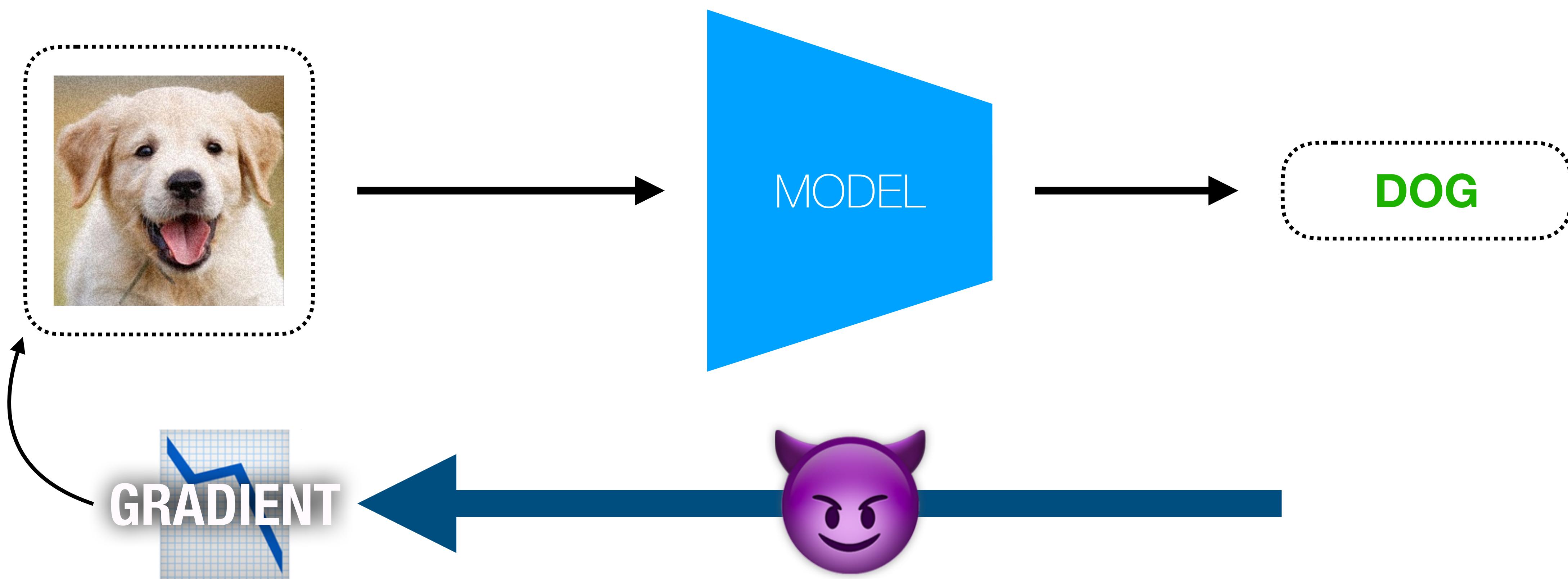
Deep Learning for Image Classification



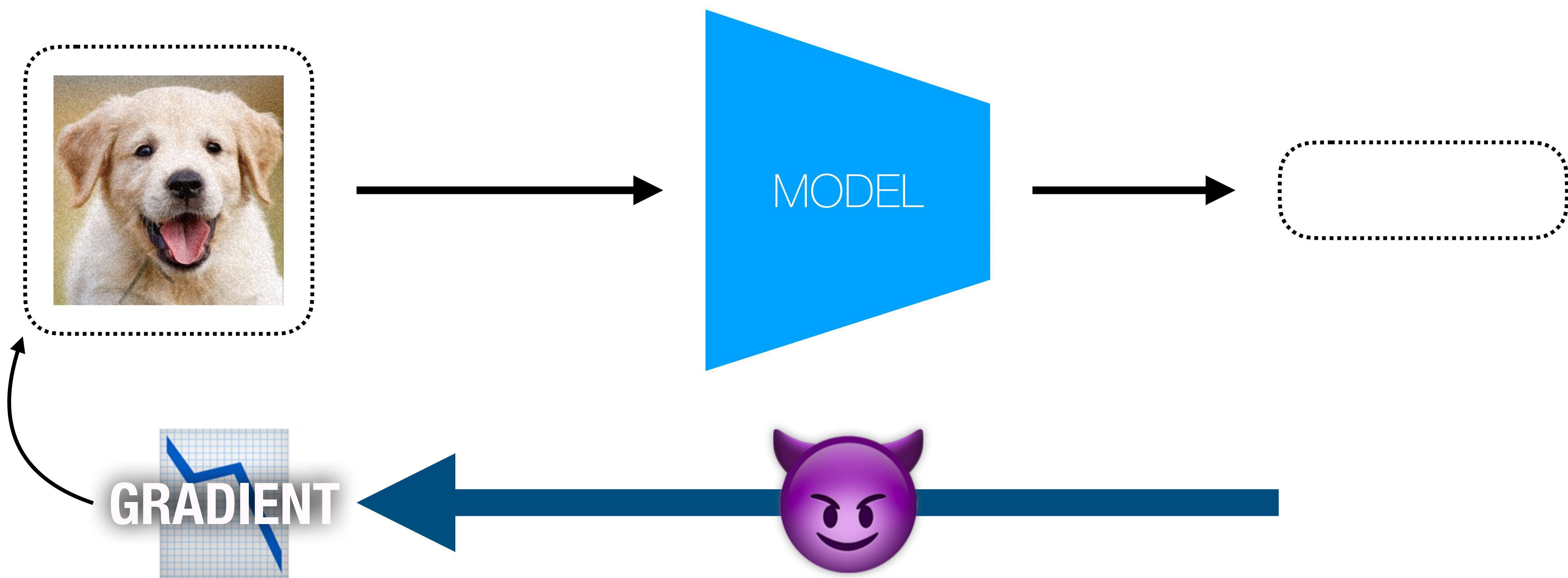
Adversarial Attack on Deep Learning



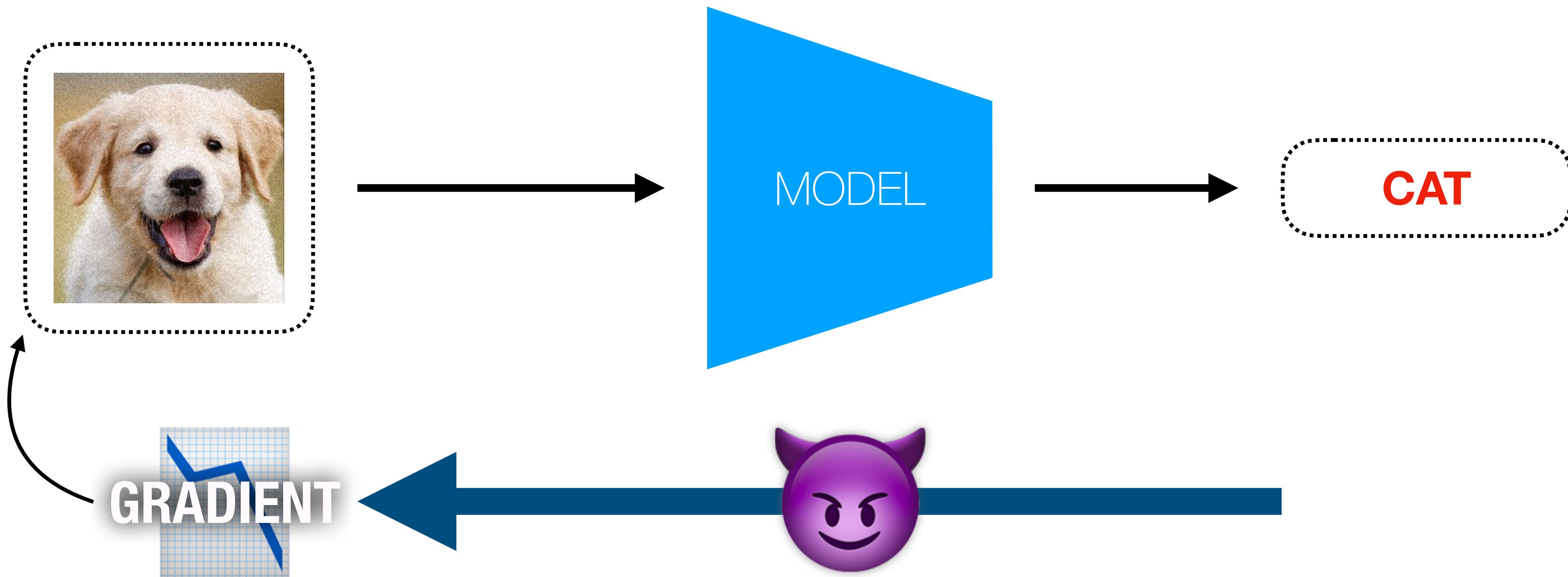
Adversarial Attack on Deep Learning



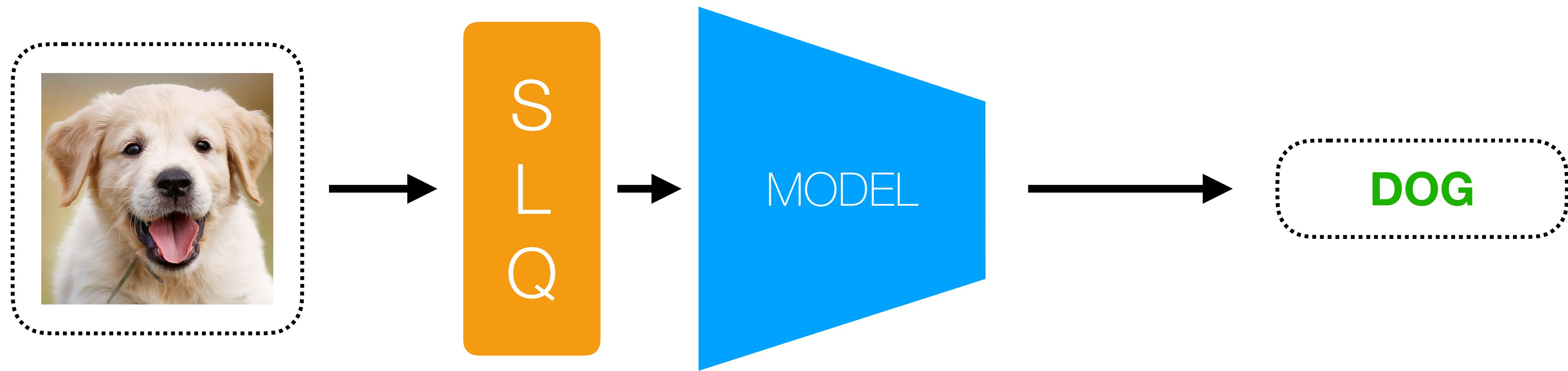
Adversarial Attack on Deep Learning



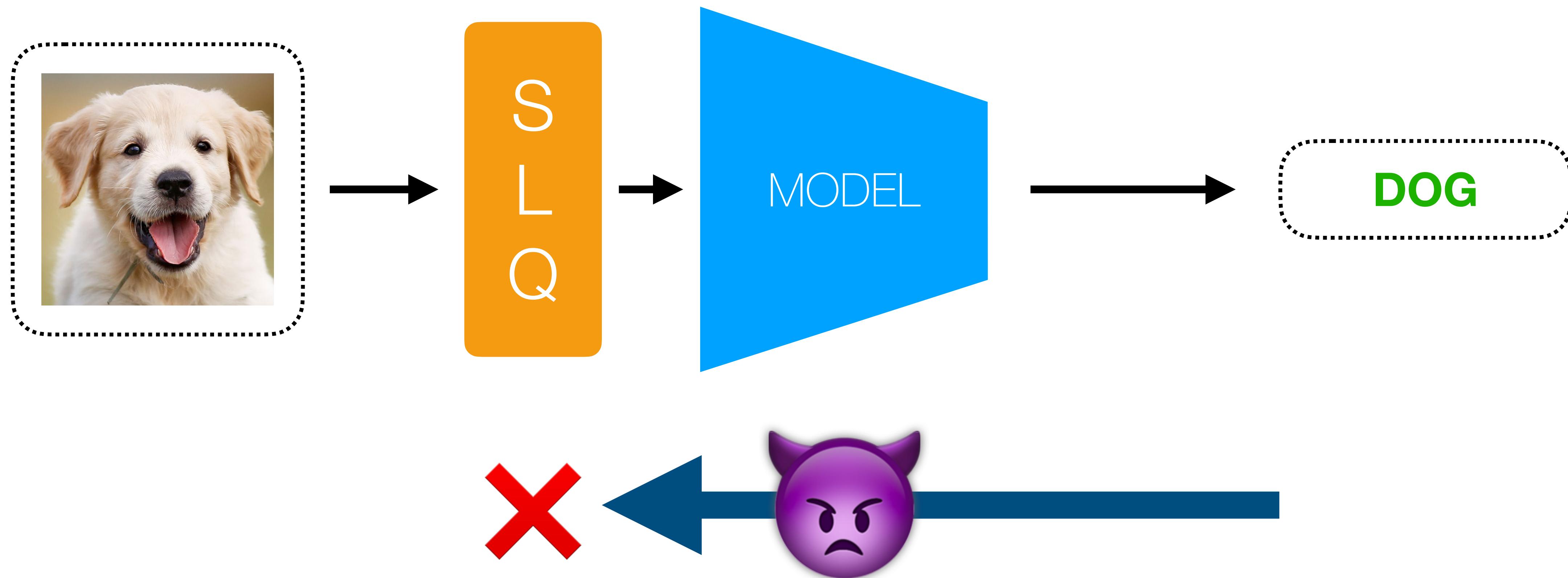
Adversarial Attack on Deep Learning



Stochastic Local Quantization (SLQ)



Stochastic Local Quantization (SLQ)



SLQ leverages JPEG compression

SLQ leverages JPEG compression



JPEG Quality 80



JPEG Quality 60



JPEG Quality 40



JPEG Quality 20

SLQ leverages JPEG compression



JPEG Quality 80



JPEG Quality 60

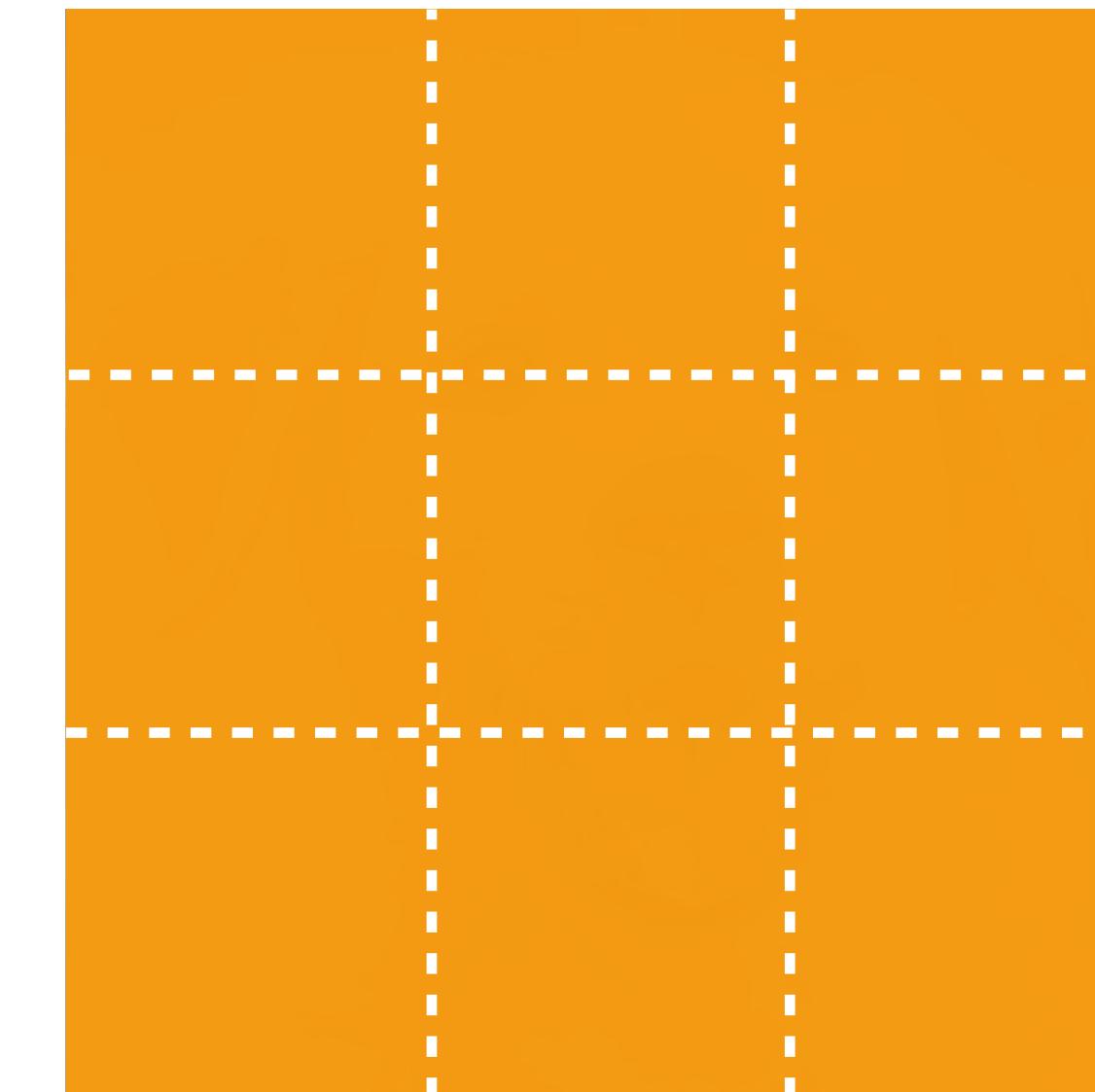
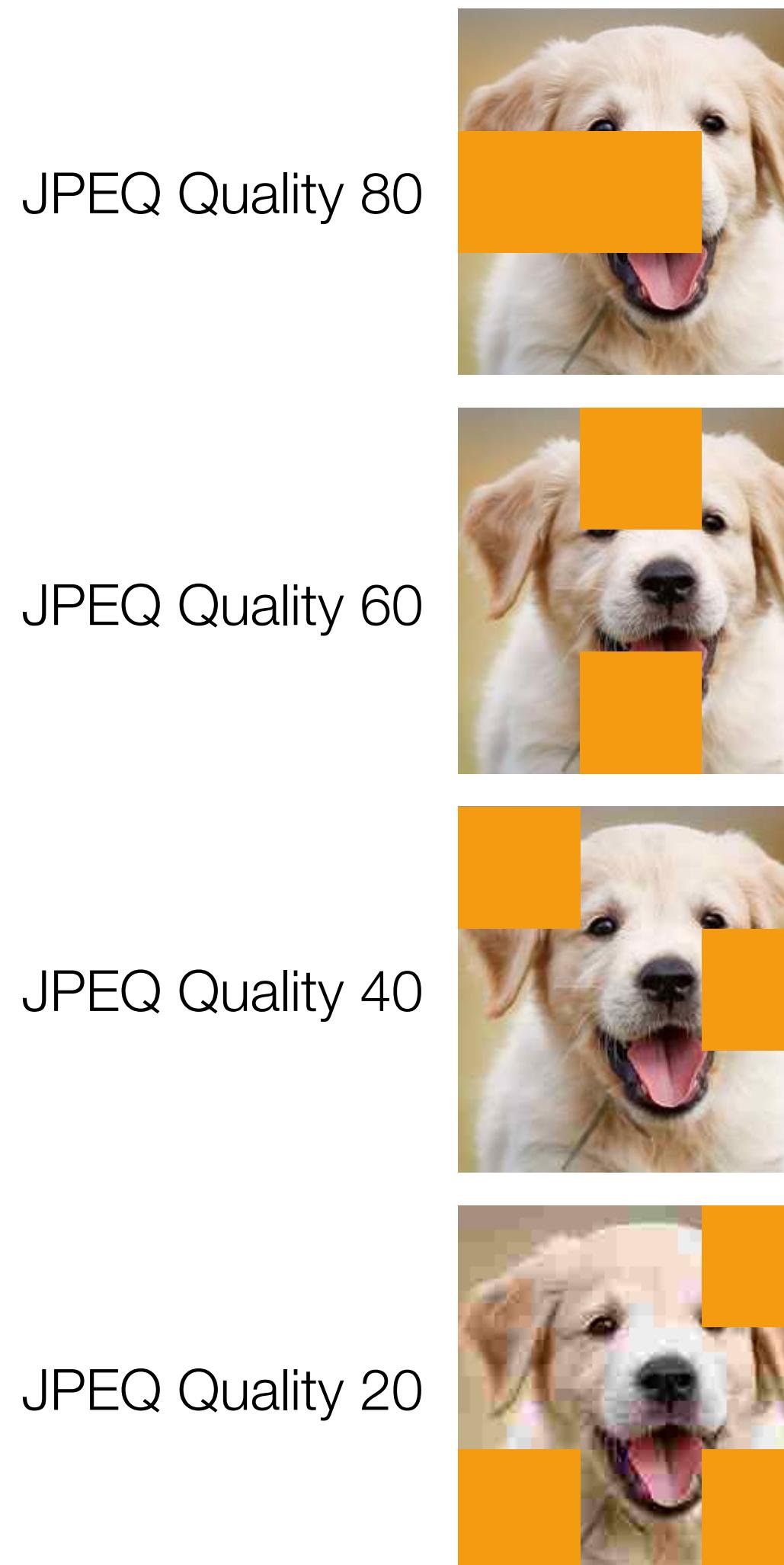


JPEG Quality 40



JPEG Quality 20

SLQ leverages JPEG compression with randomization



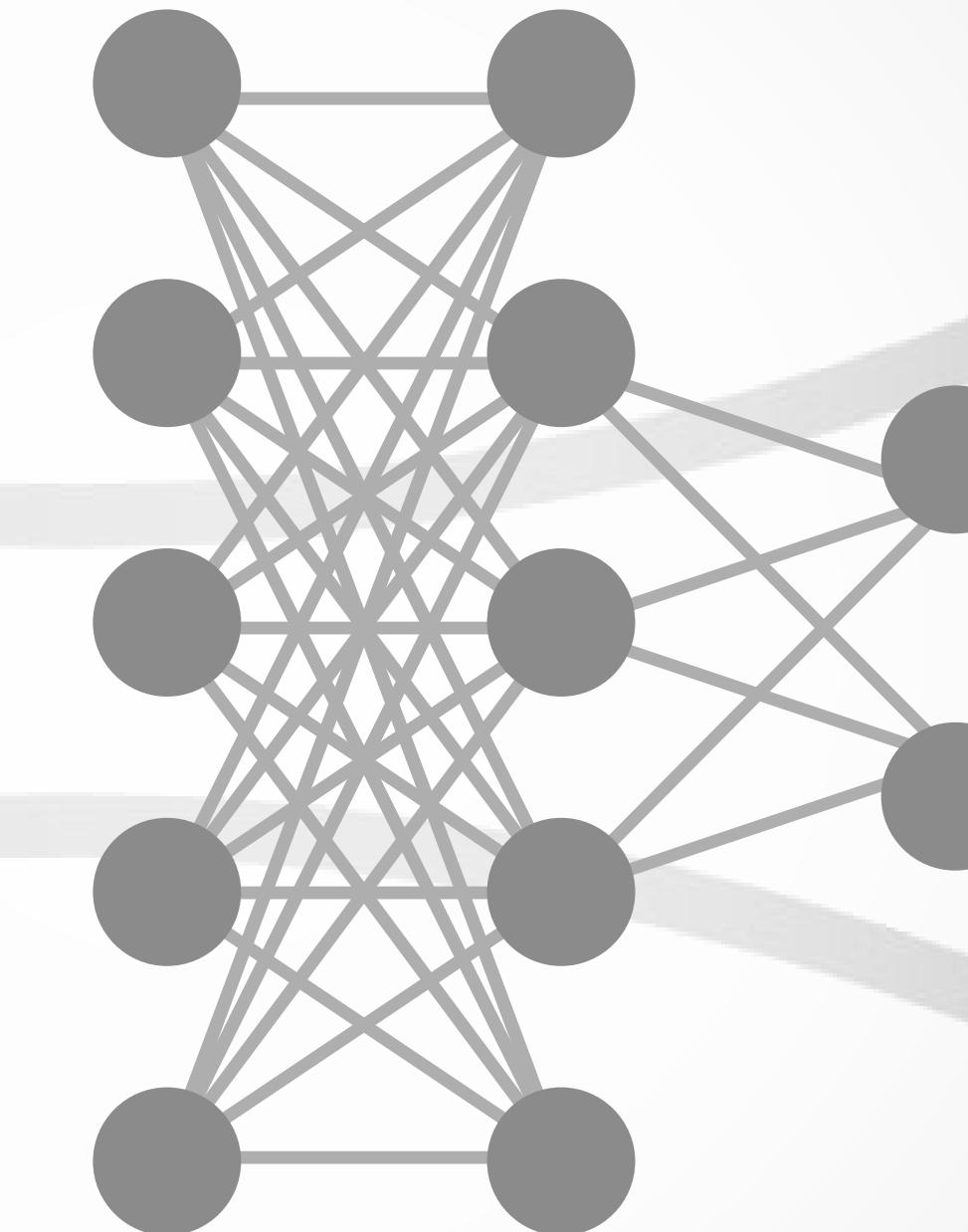
SLQ applies JPEG compression
of a random quality to each
8 x 8 block of the image

* larger blocks shown for presentation



SHIELD

Secure Heterogeneous Image Ensemble with Localized Denoising



Real-time
Compression
Preprocessing

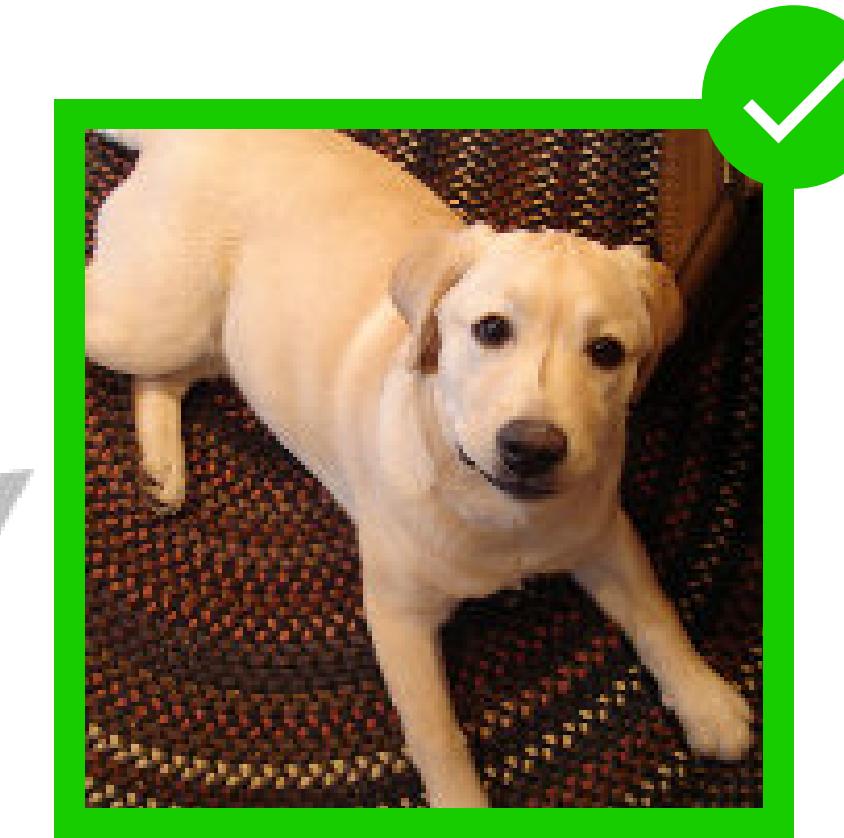
Vaccinated
Deep Neural
Network Ensemble



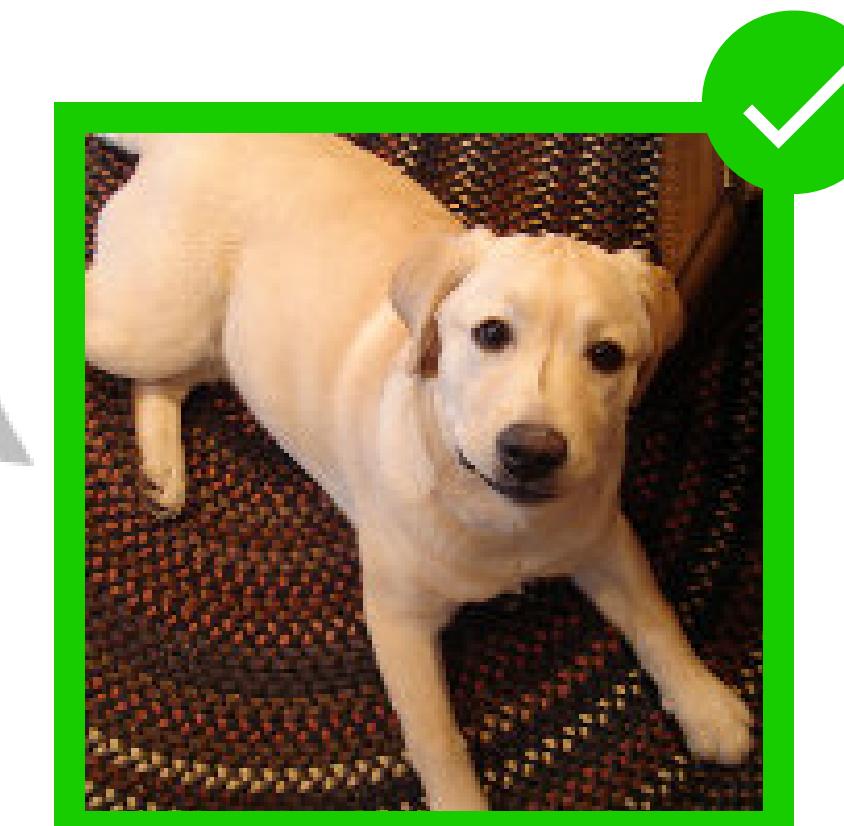
"Chain Mail"
(Attacked)



Labrador
Retriever



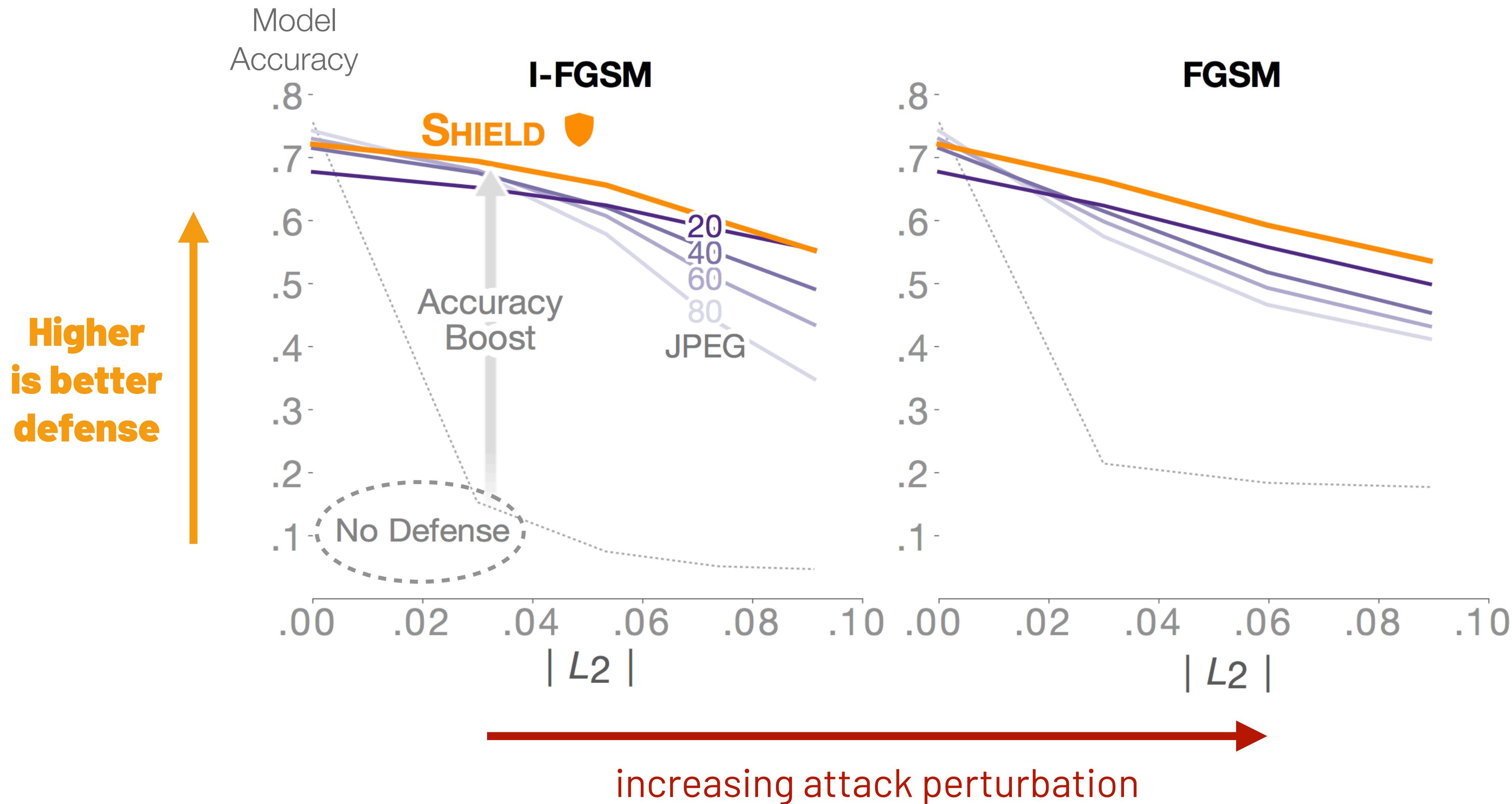
Correctly
Classified



Correctly
Classified

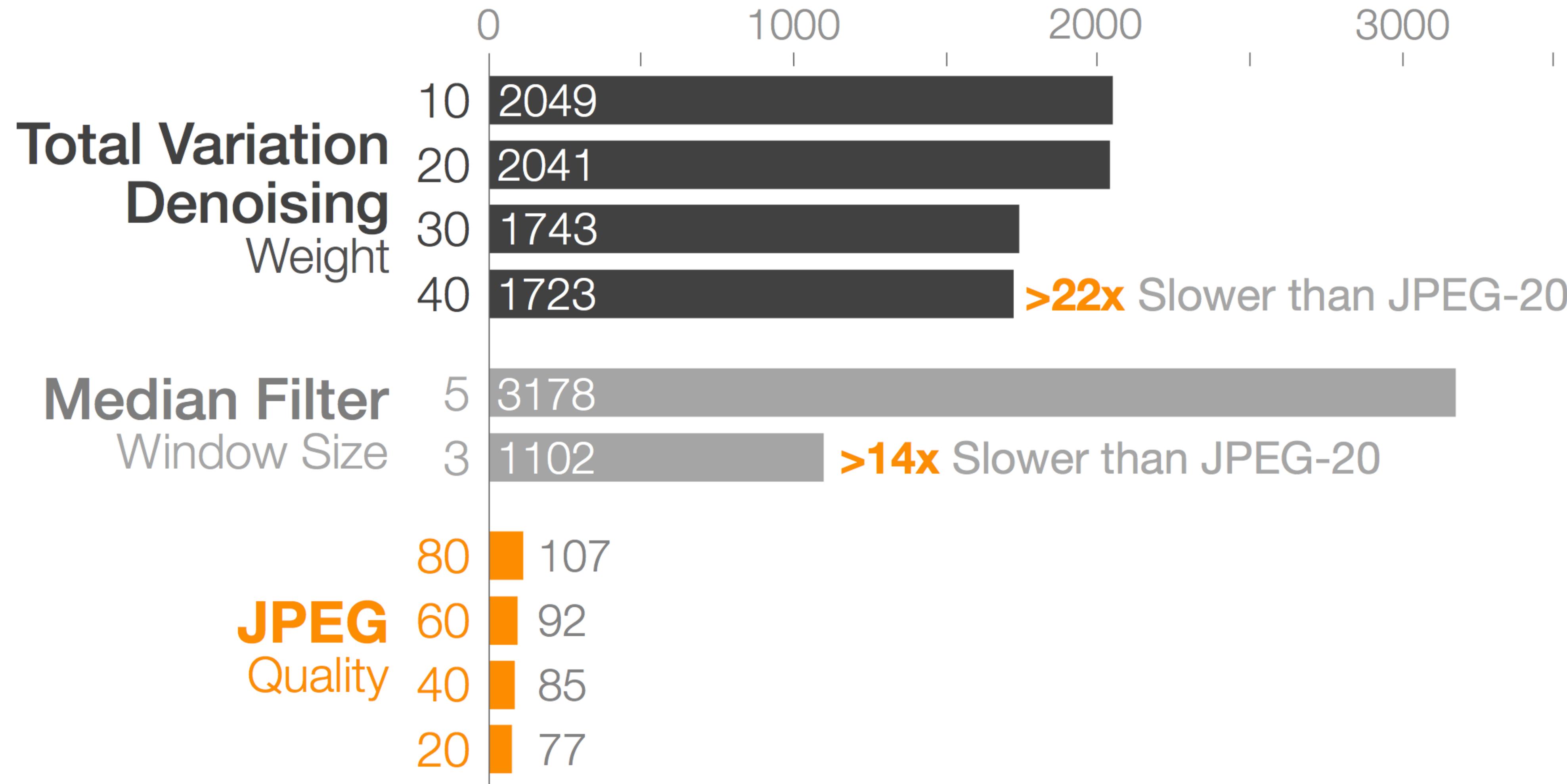
Results with ResNet-50 v2

(on ImageNet validation set)



Defense Runtime Comparison

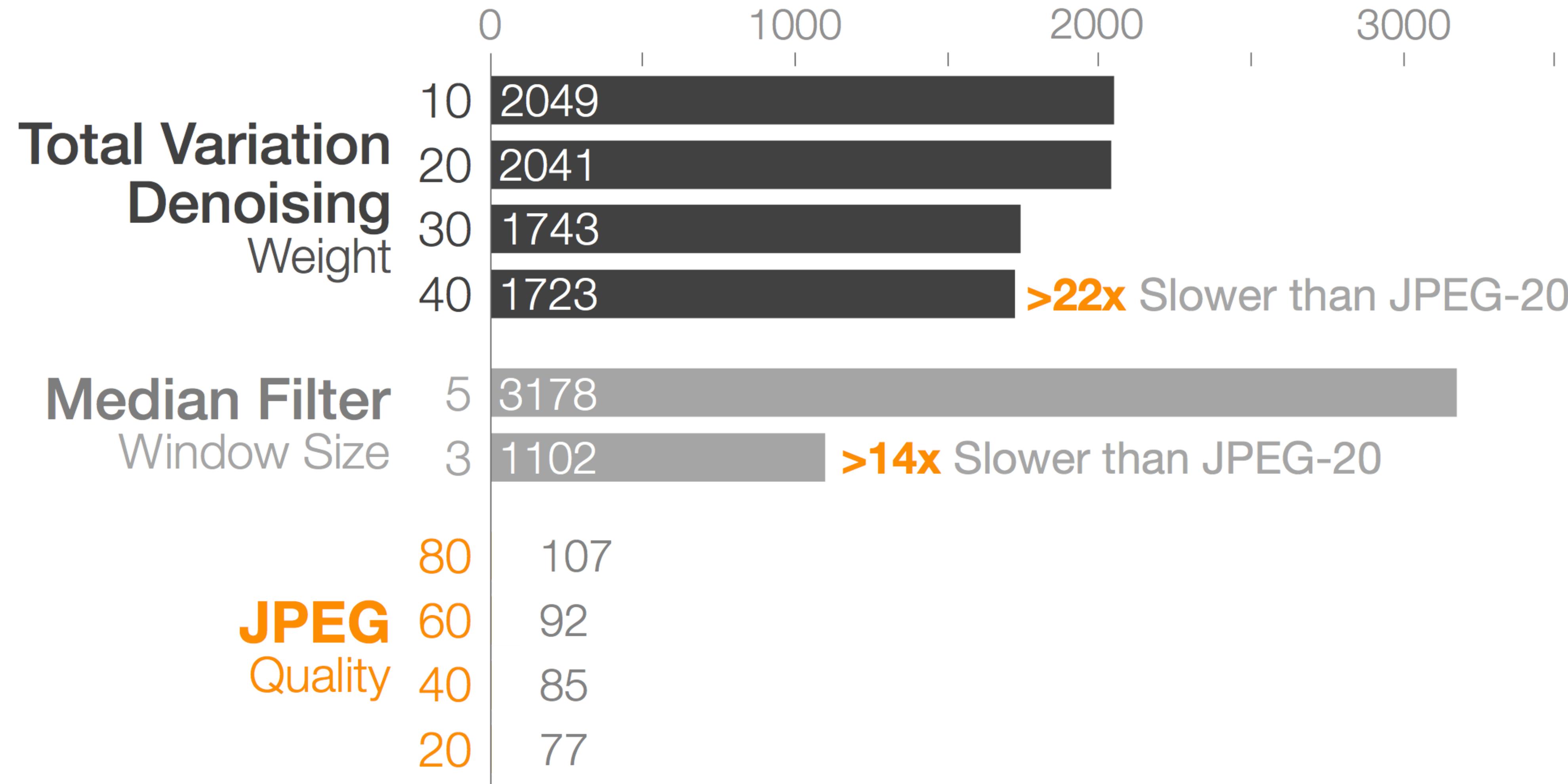
(in seconds; shorter is better)



tested on 50,000 images from the ImageNet validation set

Defense Runtime Comparison

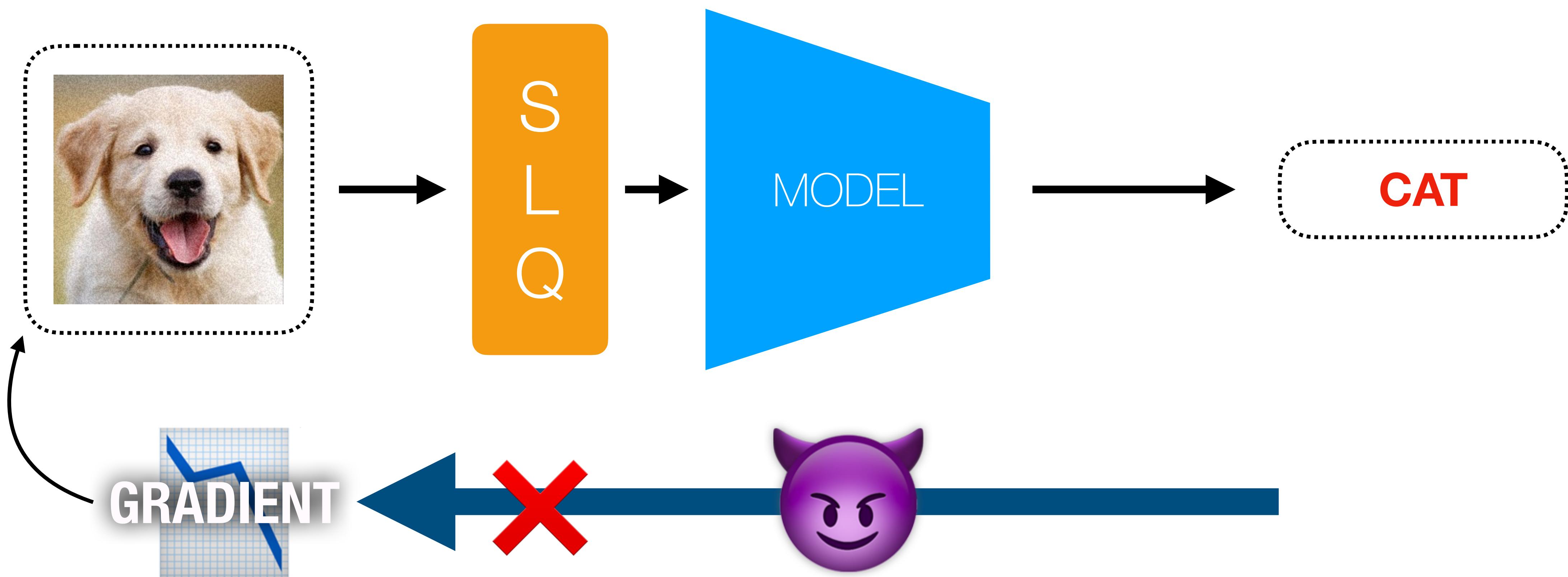
(in seconds; shorter is better)



tested on 50,000 images from the ImageNet validation set

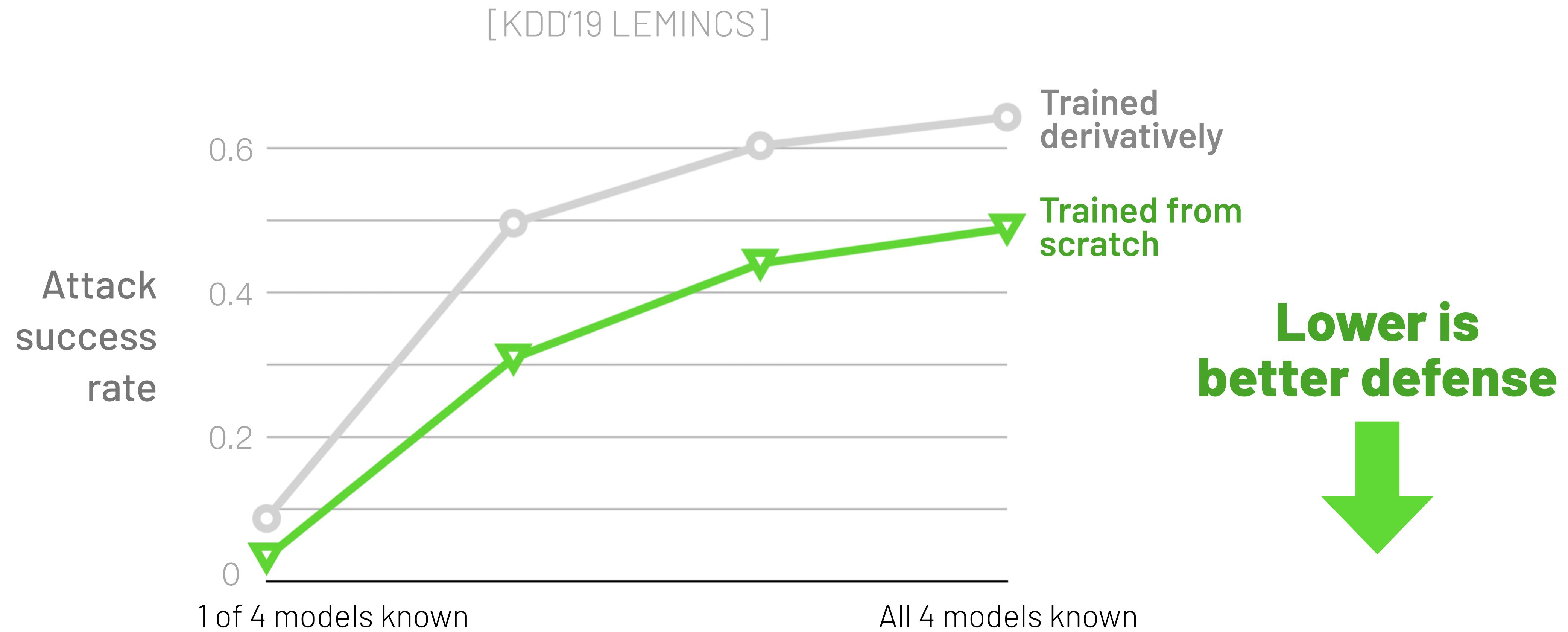
Adaptive Attacks

- Backward Pass Differentiable Approximation (BPDA)
- Expectation over Transformations (EoT)

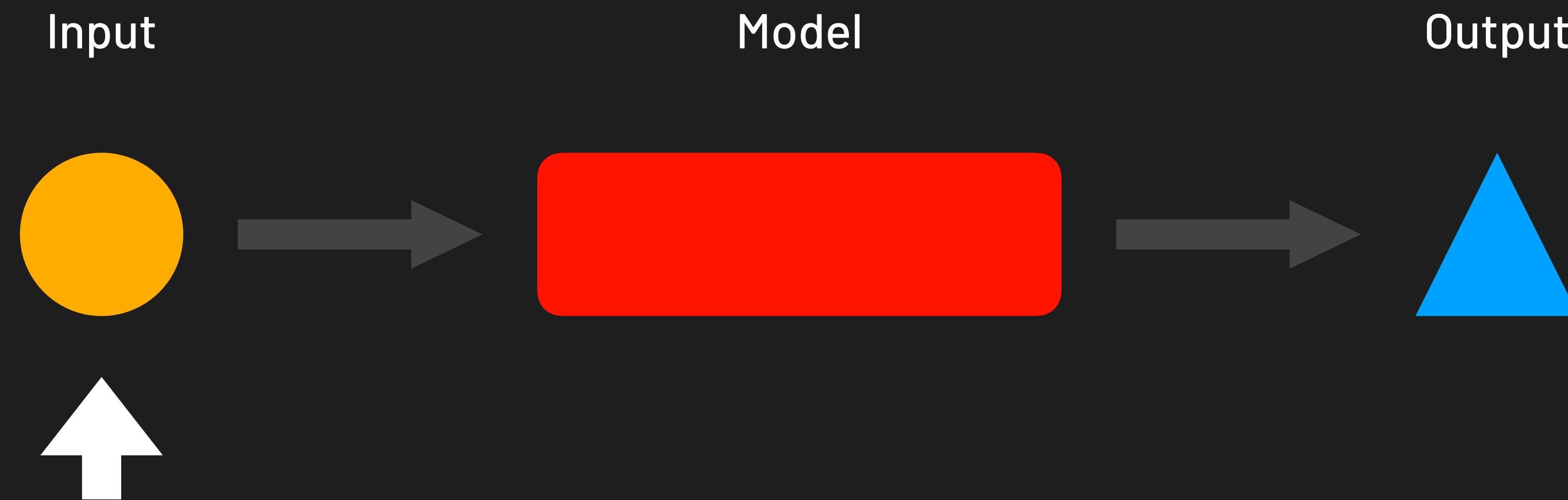


Extending to Adaptive Attacks

SHIELD ensemble with *less correlated* model weights
are *more robust* to targeted **adaptive attacks**



Machine Learning Pipeline



SkeleVision

arXiv 2022
(under review)

Adversarial Resiliency of Person Tracking
with Multi-Task Learning

- ▶ Robust inference in **video** domain
- 🌐 Open-sourced at *github.com/nilakshdas/SkeleVision*



Nilaksh Das

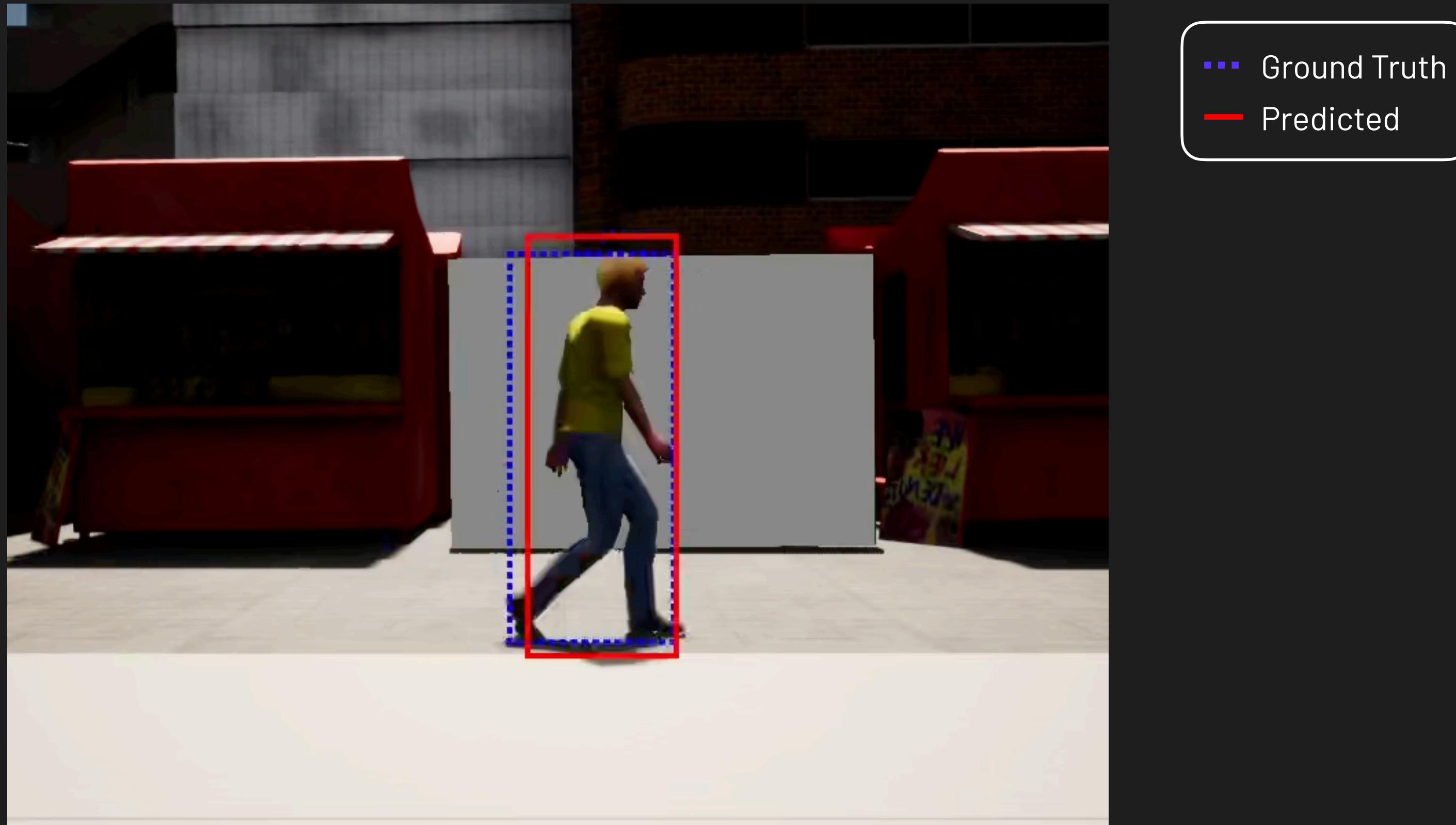


Sheng-Yun Peng

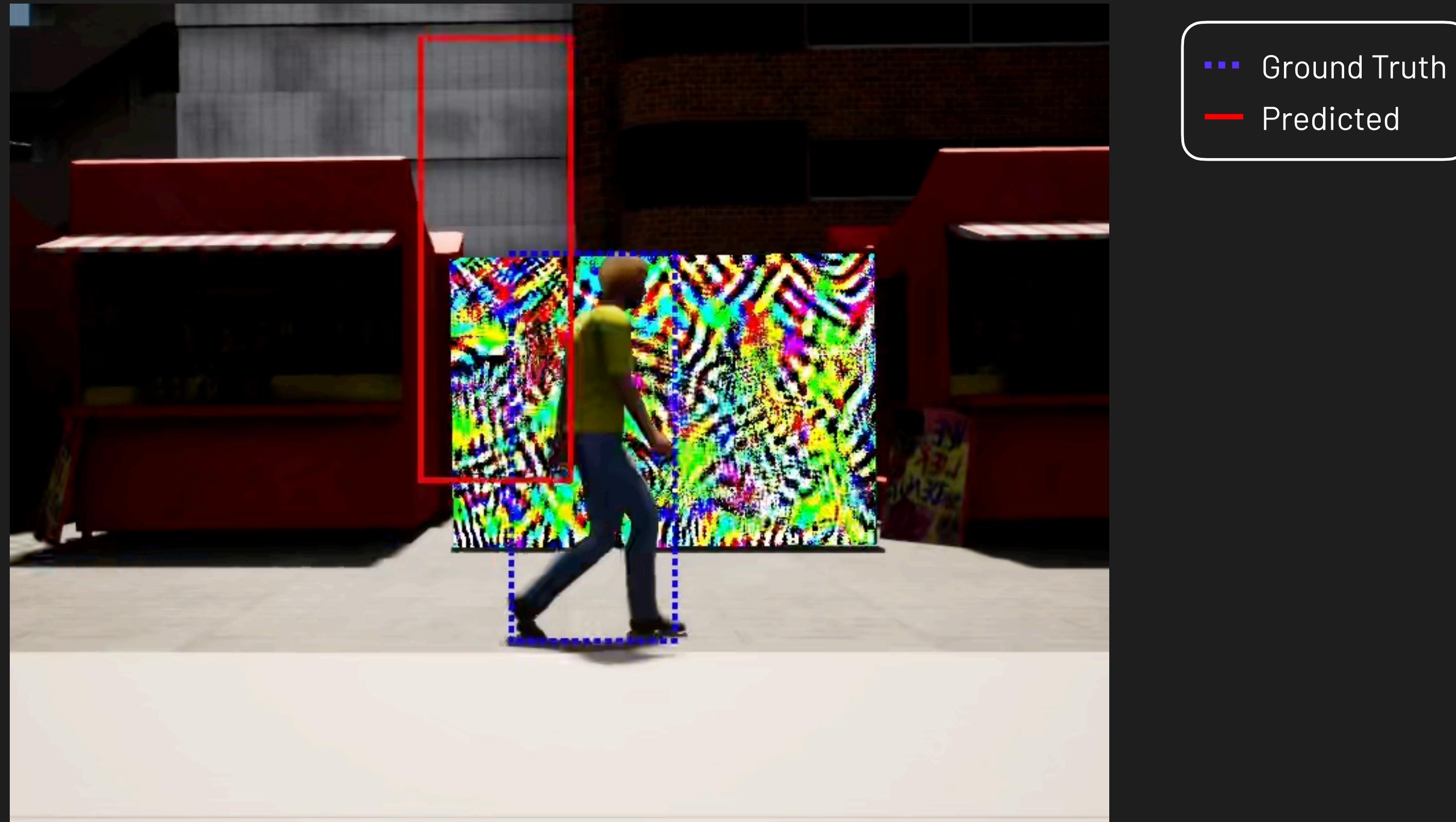


Polo Chau

Person Tracking



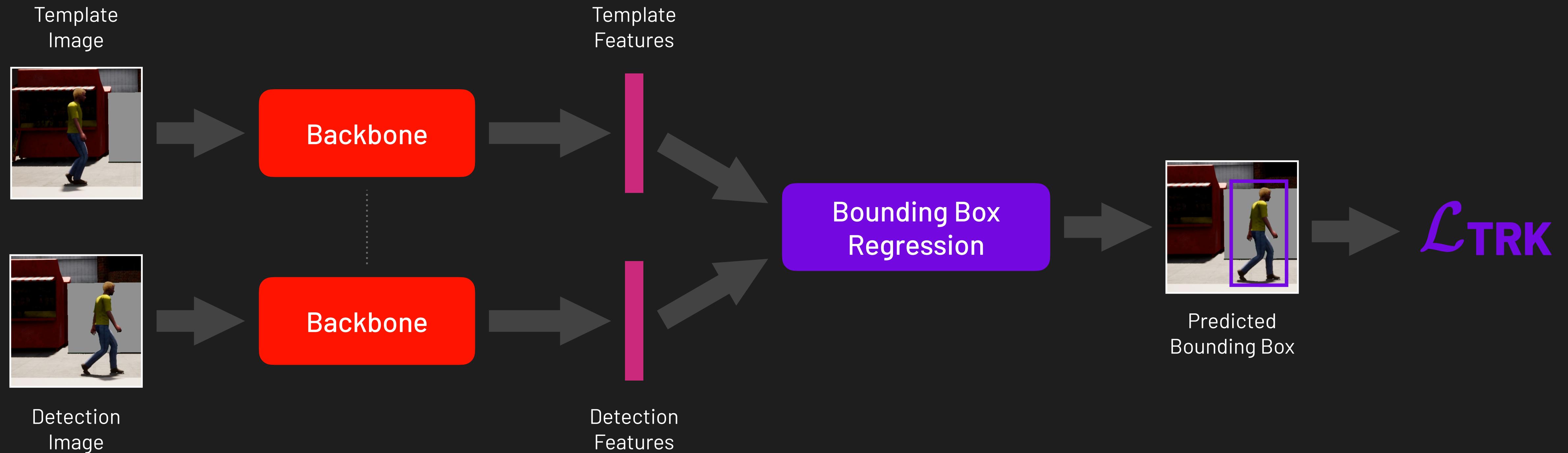
Adversarial Attack on Person Tracking



Physically Realizable

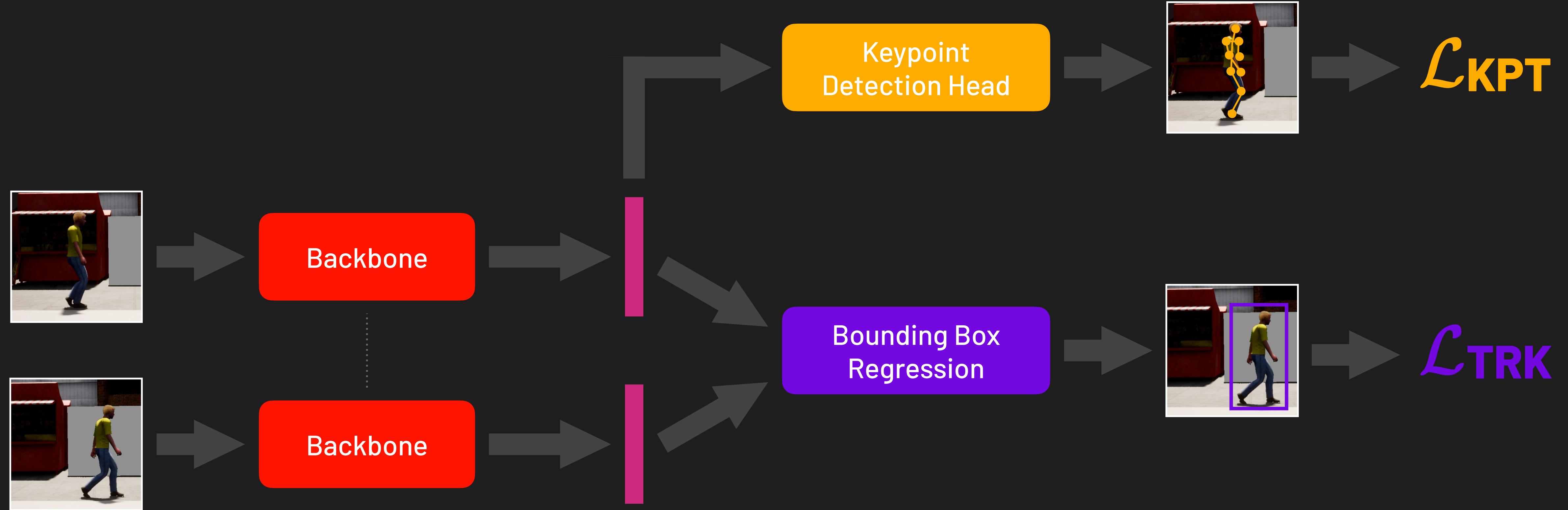


Perceptually Unbounded

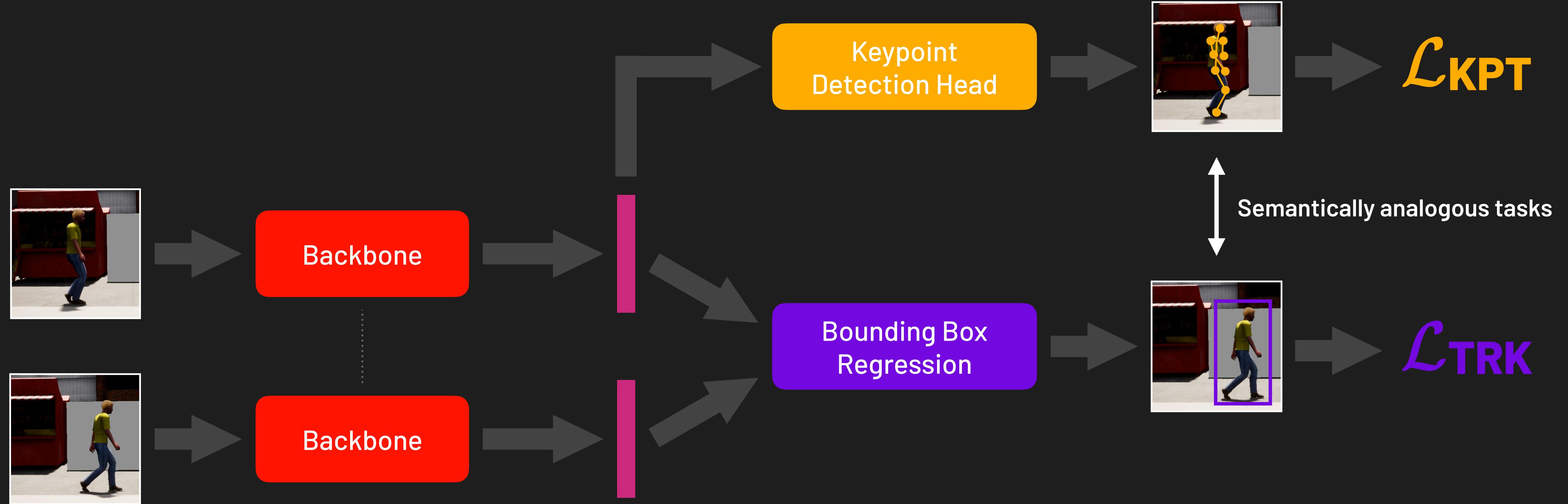


SiamRPN Model

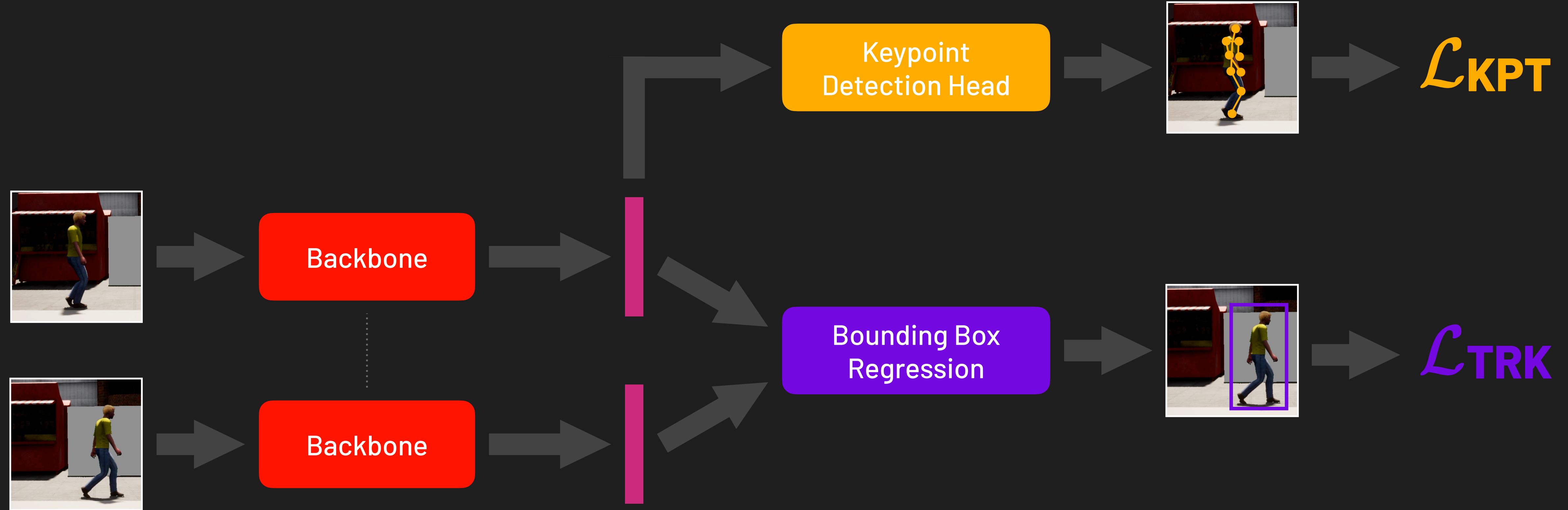
High Performance Visual Tracking with Siamese Region Proposal Network
B. Li, J. Yan, W. Wu, Z. Zhu and X. Hu
CVPR 2018



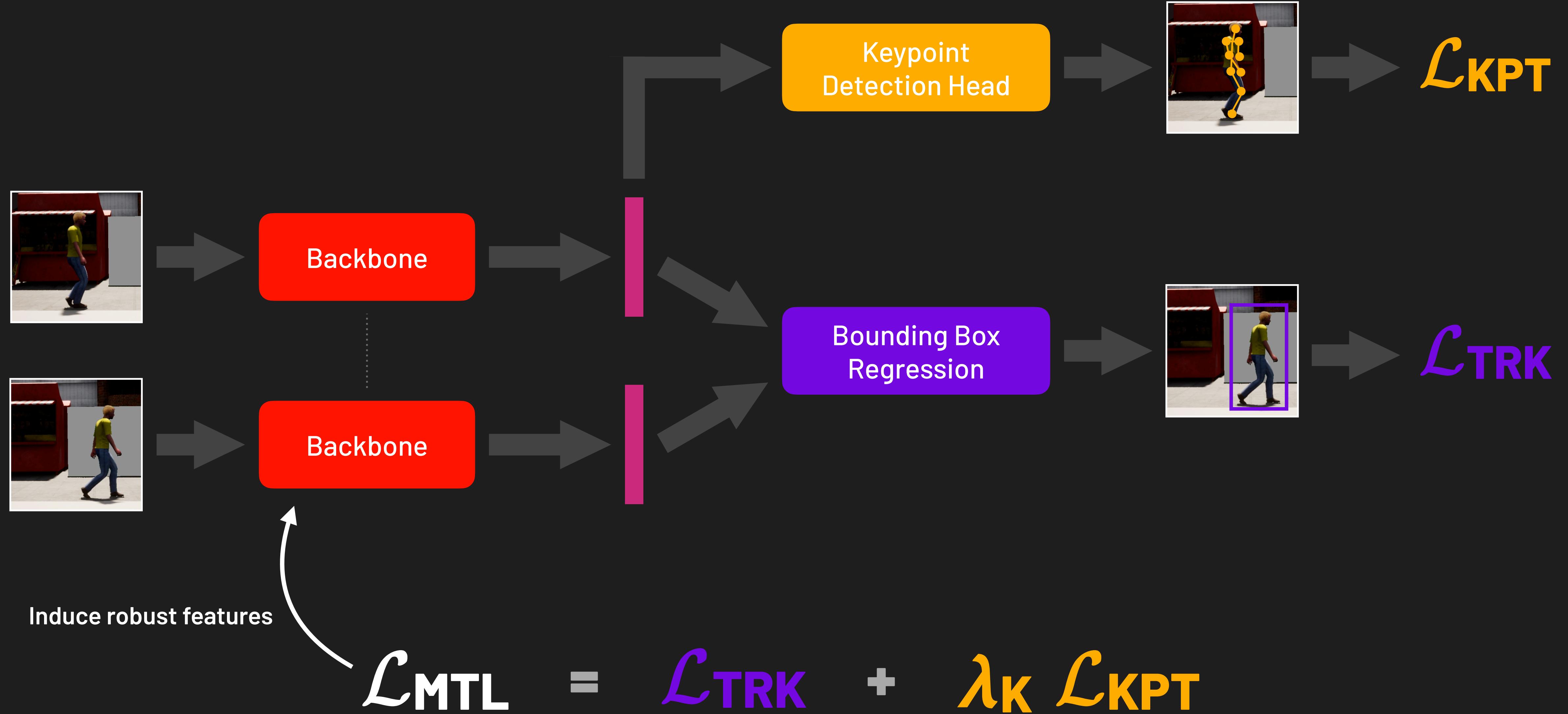
Multi-Task Learning with Tracking + Keypoint Detection



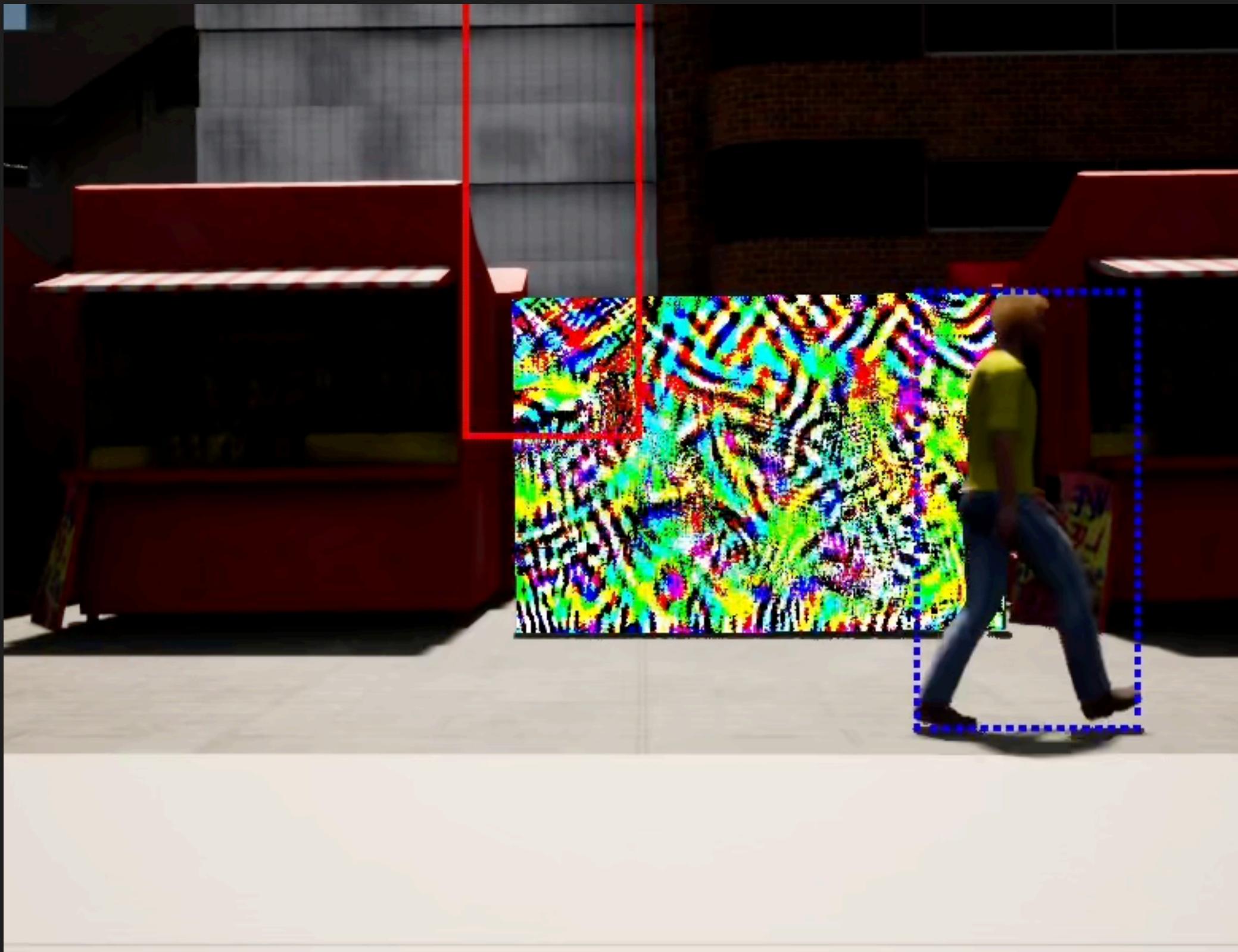
Multi-Task Learning with Tracking + Keypoint Detection



$$\mathcal{L}_{MTL} = \mathcal{L}_{TRK} + \lambda_K \mathcal{L}_{KPT}$$



Single-Task Learning (STL) Tracker



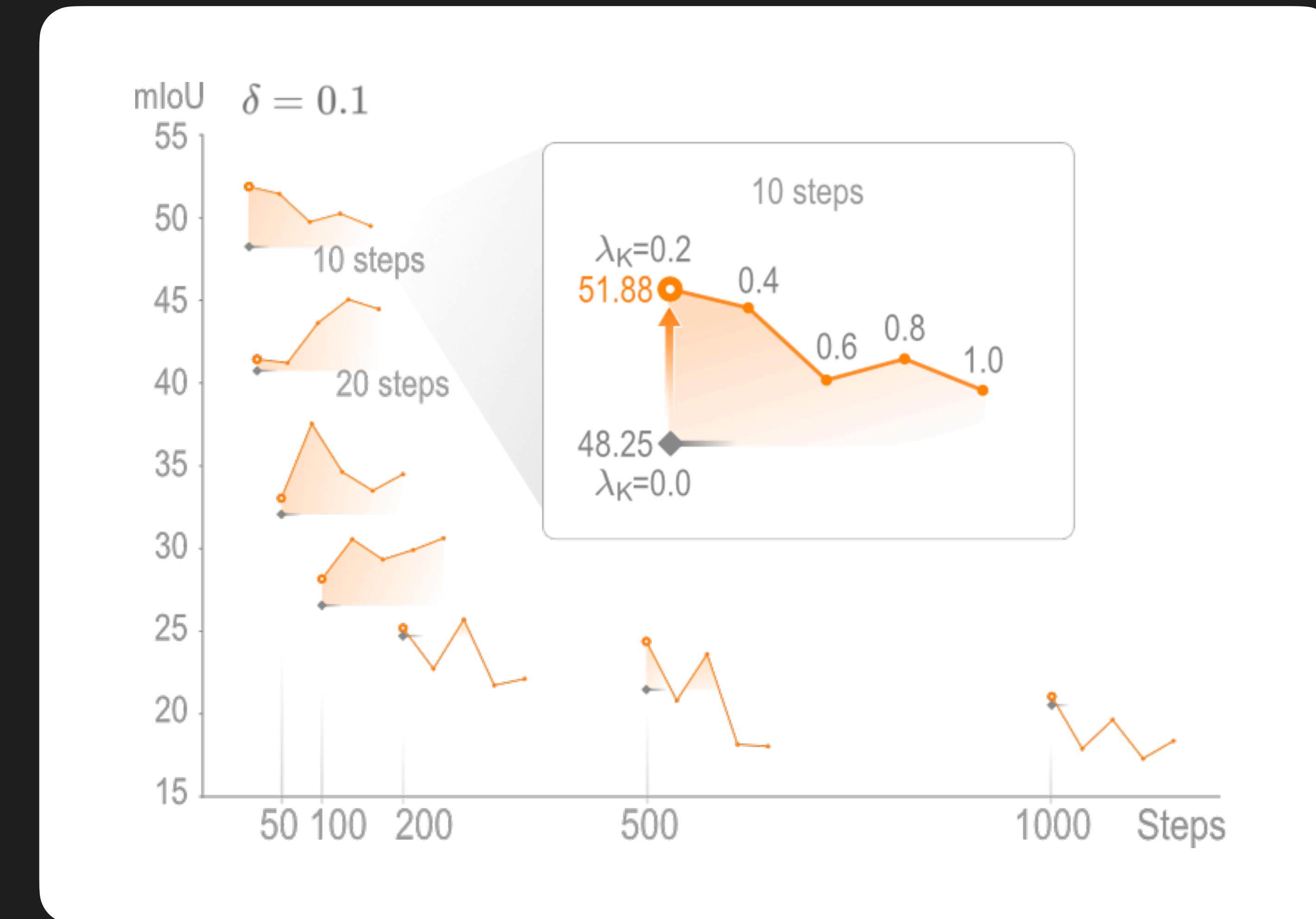
STL Tracker loses target

Multi-Task Learning (MTL) Tracker

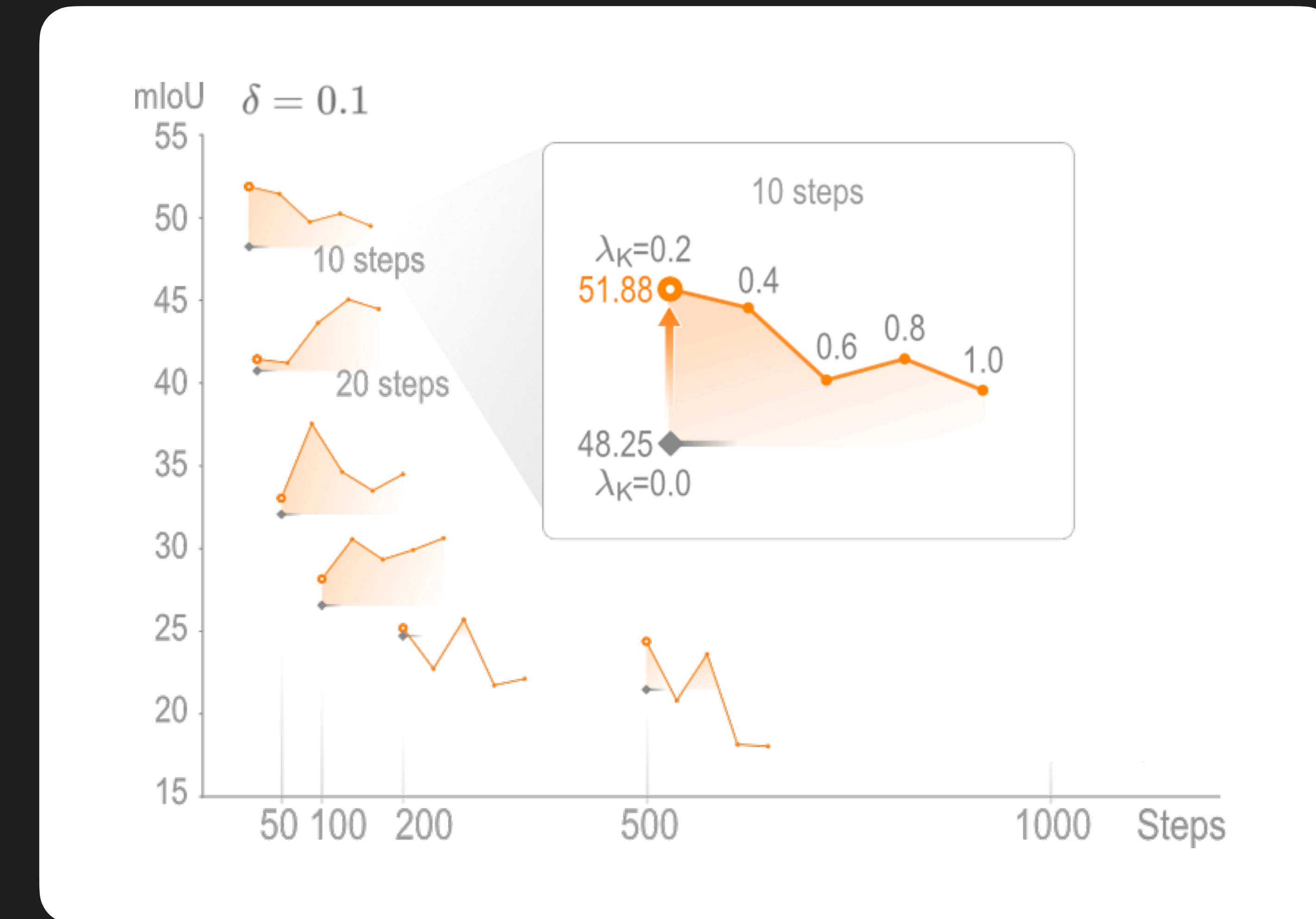


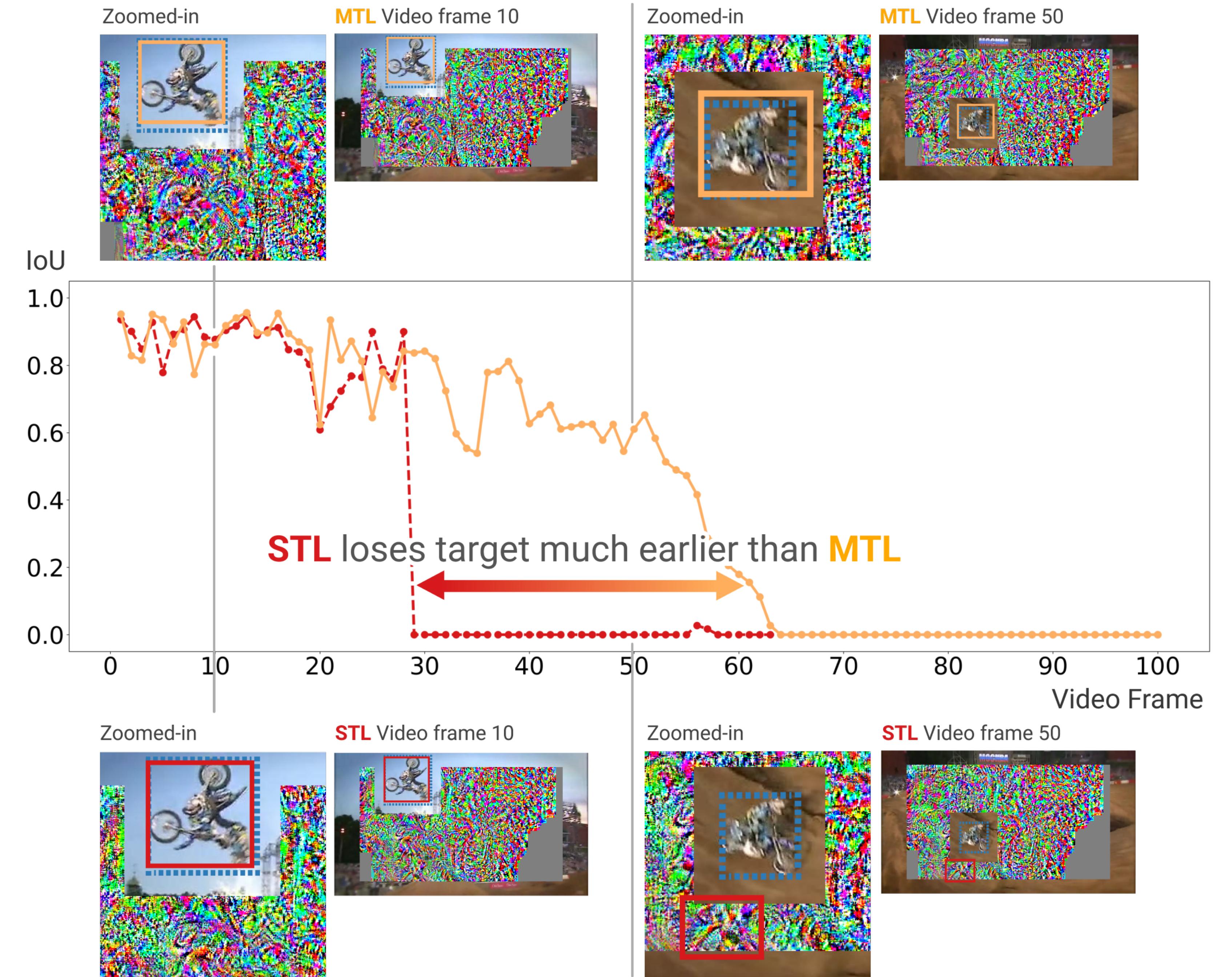
MTL Tracker recovers target

Mean Intersection-over-Union for ARMORY-CARLA dataset



Mean Intersection-over-Union for ARMORY-CARLA dataset





Hear No Evil

Towards Adversarial Robustness
of Automatic Speech Recognition
via Multi-Task Learning

arXiv 2022
(under review)



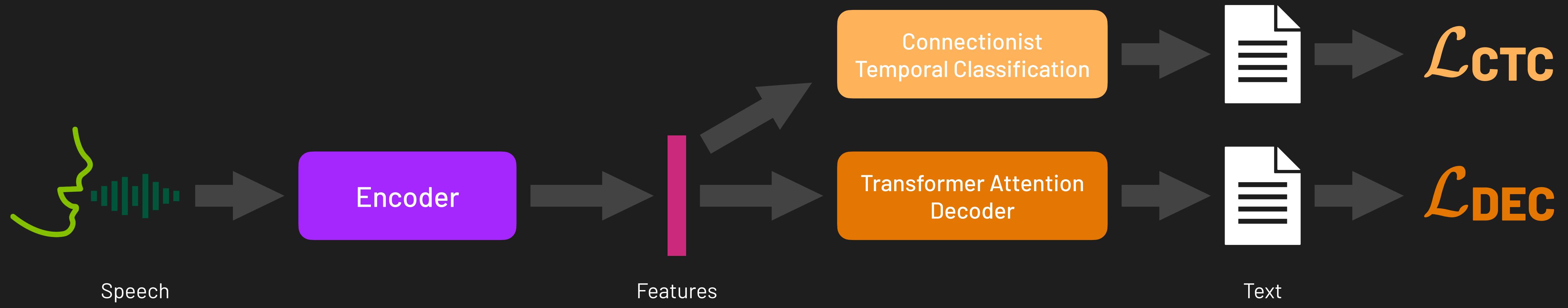
Robust inference in **audio** domain



Nilaksh Das



Polo Chau

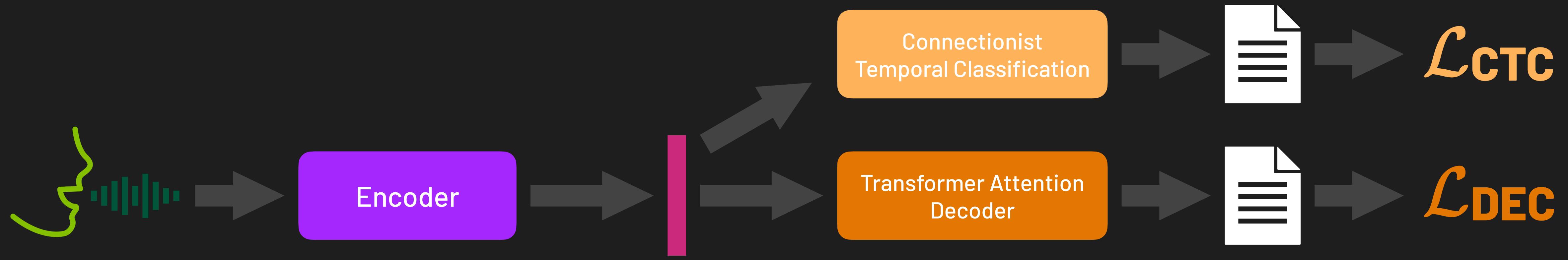


Hybrid ASR Model

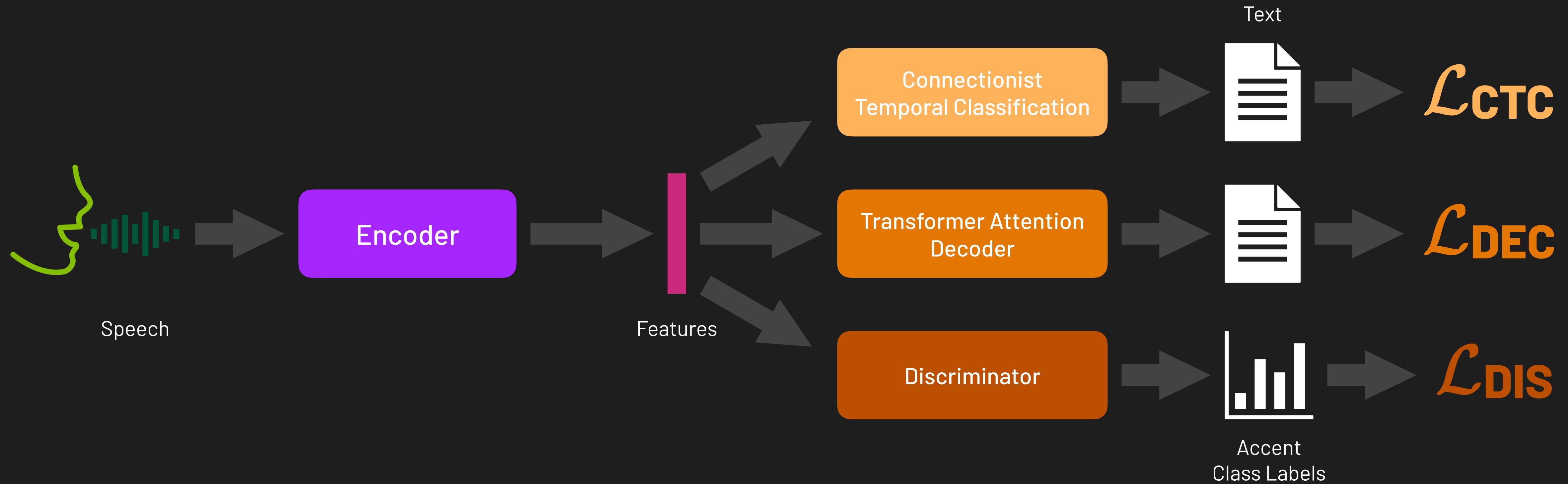
Hybrid CTC/Attention Architecture for End-to-End Speech Recognition

S. Watanabe, T. Hori, S. Kim, J. R. Hershey and T. Hayashi

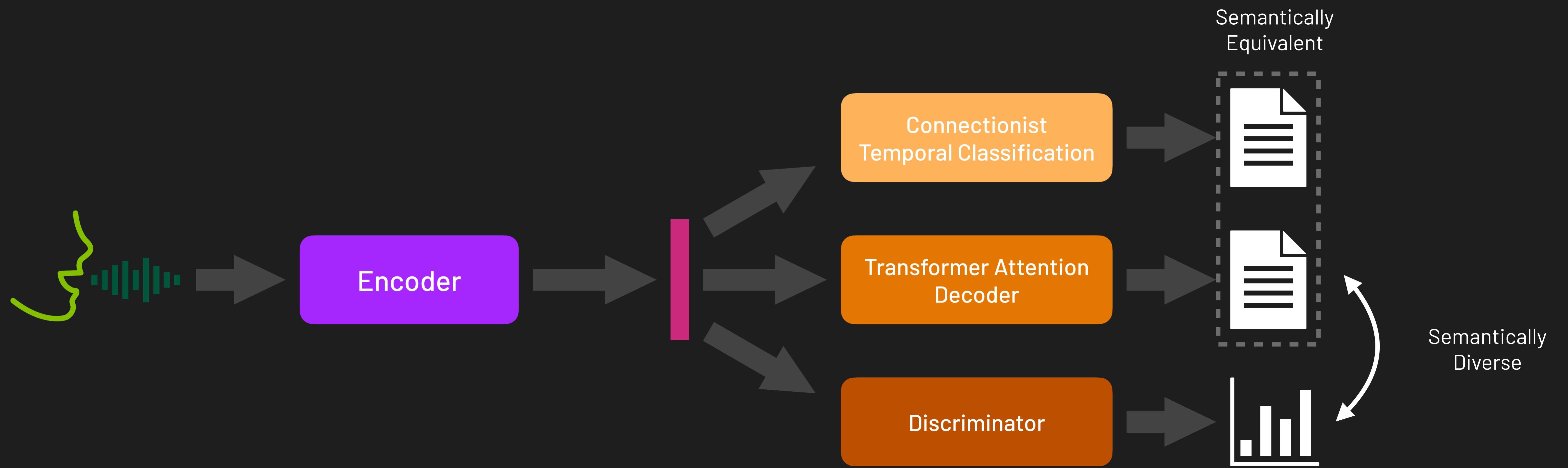
IEEE Journal of Selected Topics in Signal Processing, vol. 11

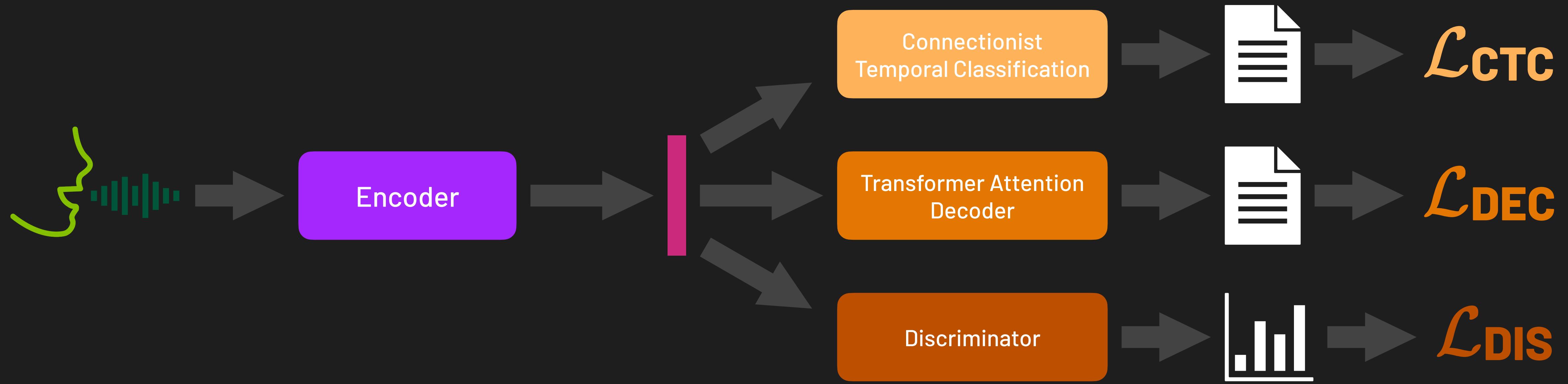


$$\mathcal{L}_{\text{ASR}} = \lambda_c \mathcal{L}_{\text{CTC}} + (1-\lambda_c) \mathcal{L}_{\text{DEC}}$$

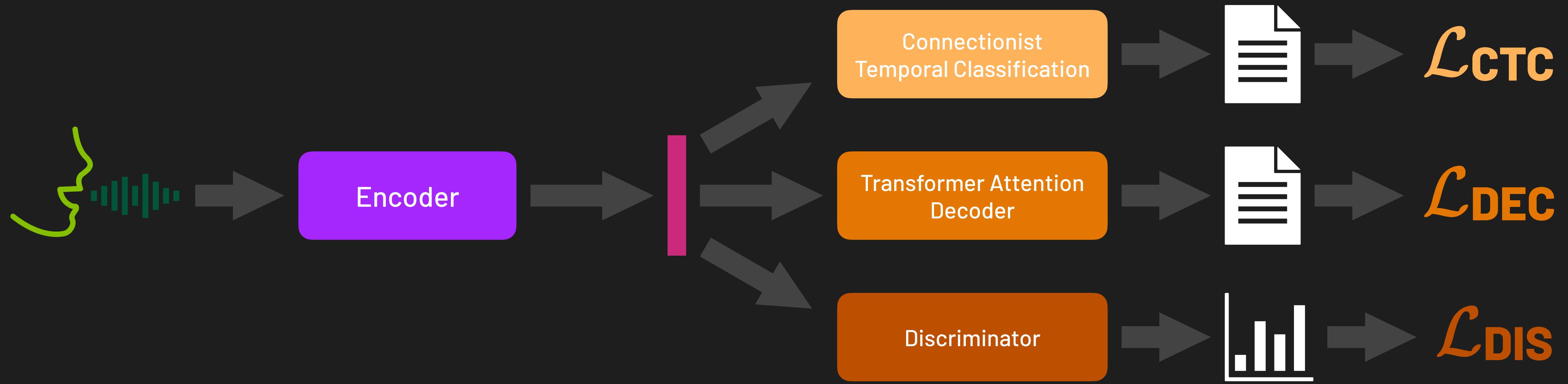


Multi-Task Learning with ASR + Accent Classification



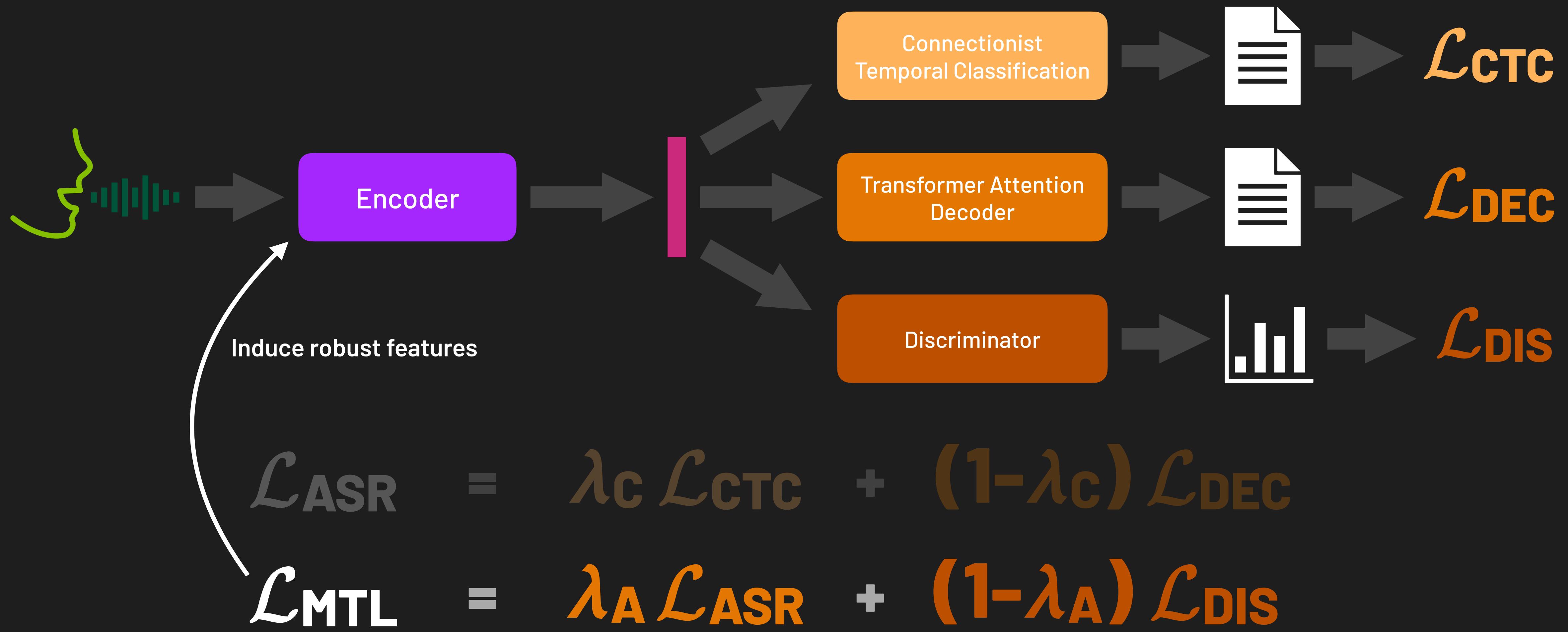


$$\mathcal{L}_{ASR} = \lambda_c \mathcal{L}_{CTC} + (1-\lambda_c) \mathcal{L}_{DEC}$$

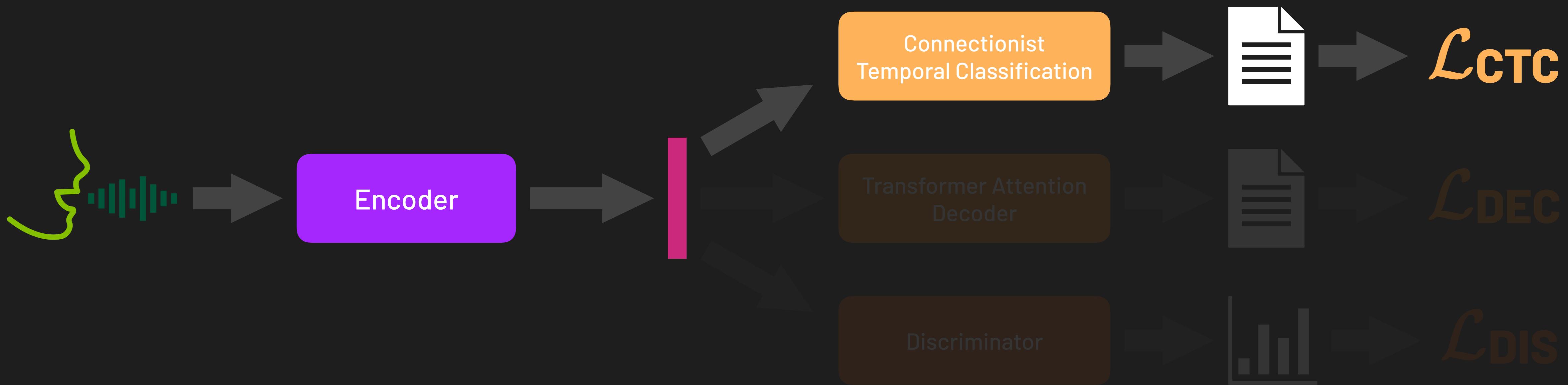


$$\mathcal{L}_{\text{ASR}} = \lambda_c \mathcal{L}_{\text{CTC}} + (1-\lambda_c) \mathcal{L}_{\text{DEC}}$$

$$\mathcal{L}_{\text{MTL}} = \lambda_A \mathcal{L}_{\text{ASR}} + (1-\lambda_A) \mathcal{L}_{\text{DIS}}$$



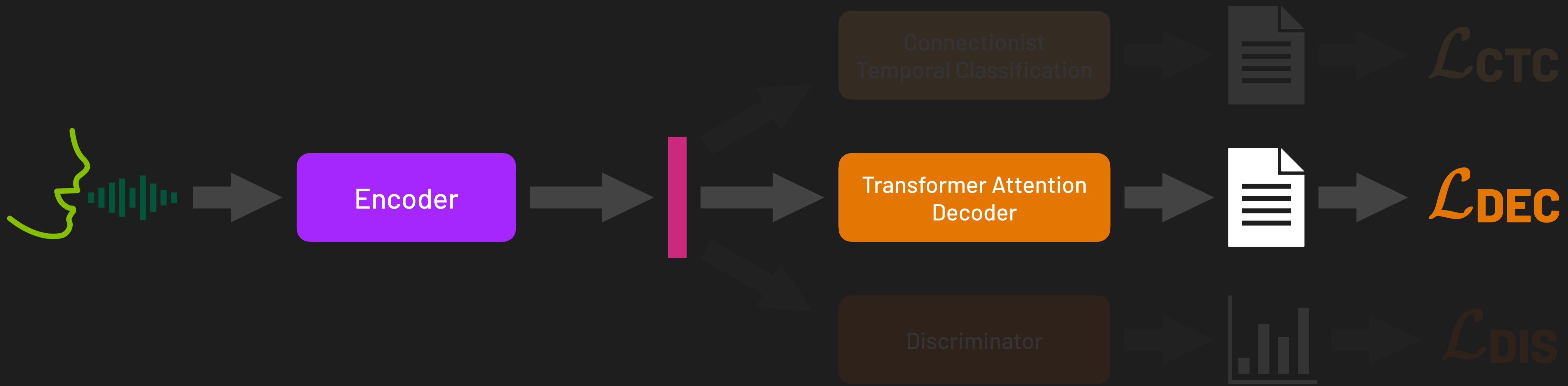
Single-Task Learning with CTC Loss



$$\mathcal{L}_{ASR} = \lambda_c \mathcal{L}_{CTC} + (1-\lambda_c) \mathcal{L}_{DEC} \quad \lambda_c = 1.0$$

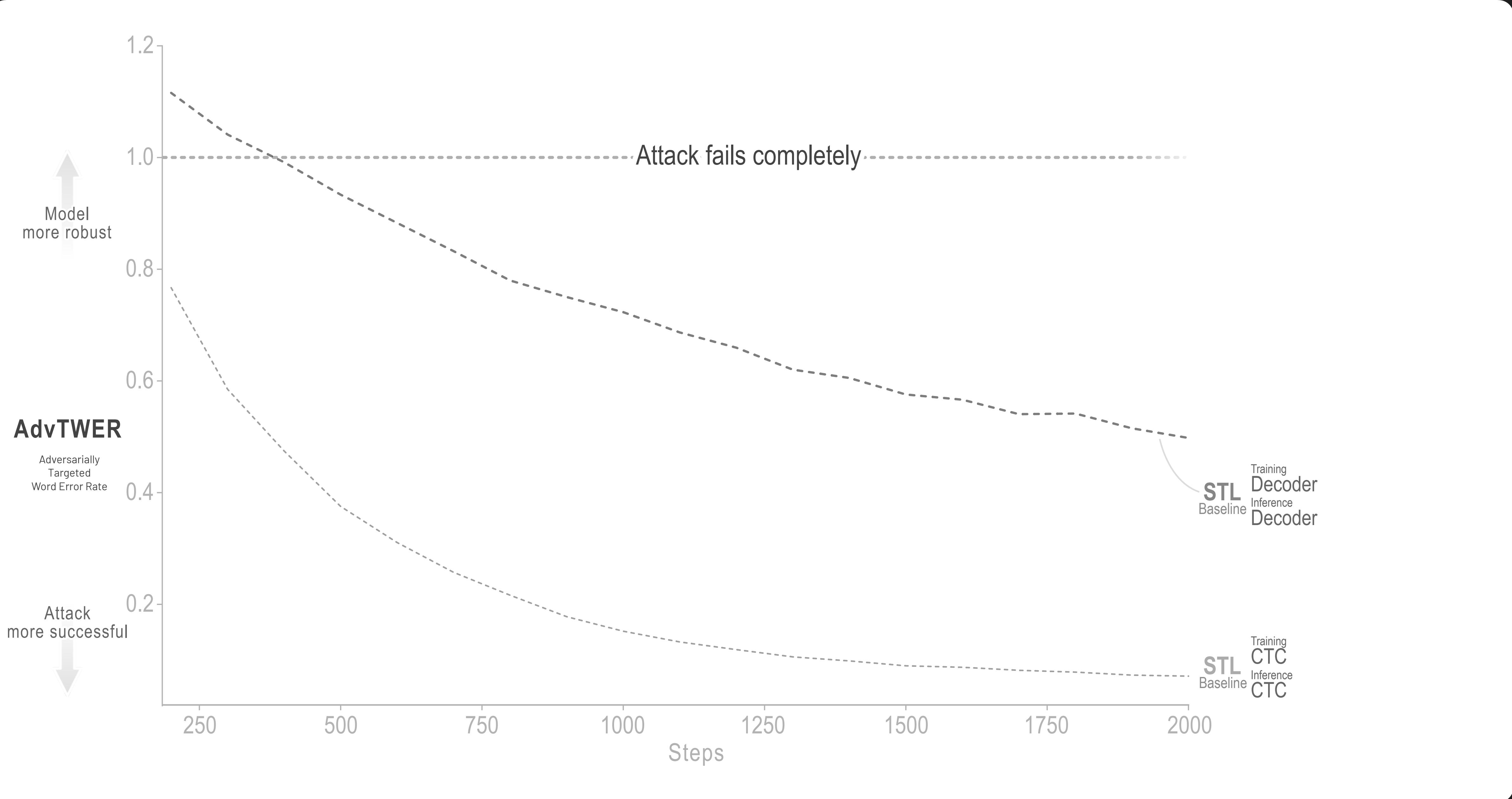
$$\mathcal{L}_{MTL} = \lambda_A \mathcal{L}_{ASR} + (1-\lambda_A) \mathcal{L}_{DIS} \quad \lambda_A = 0.0$$

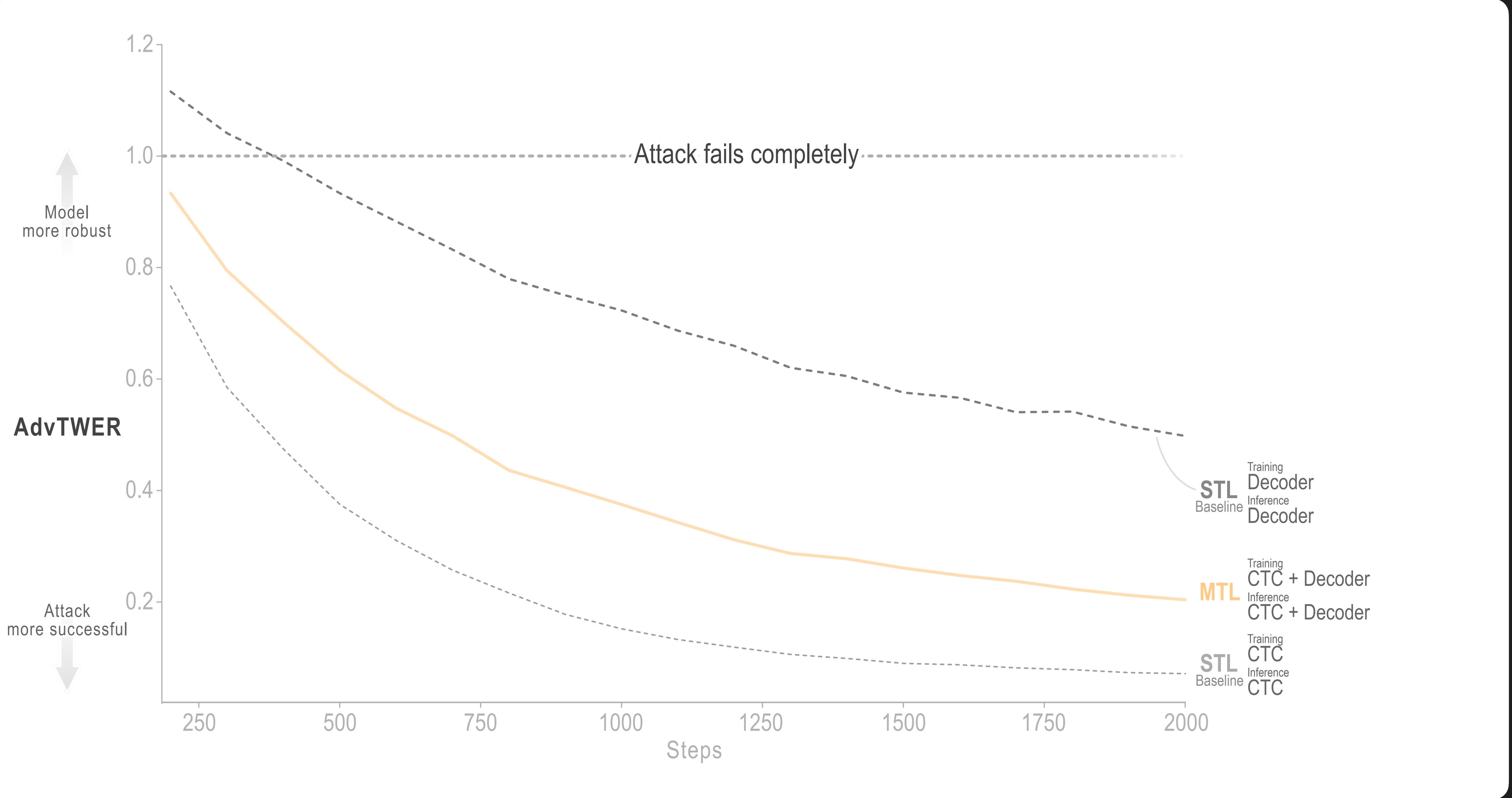
Single-Task Learning with Decoder Loss

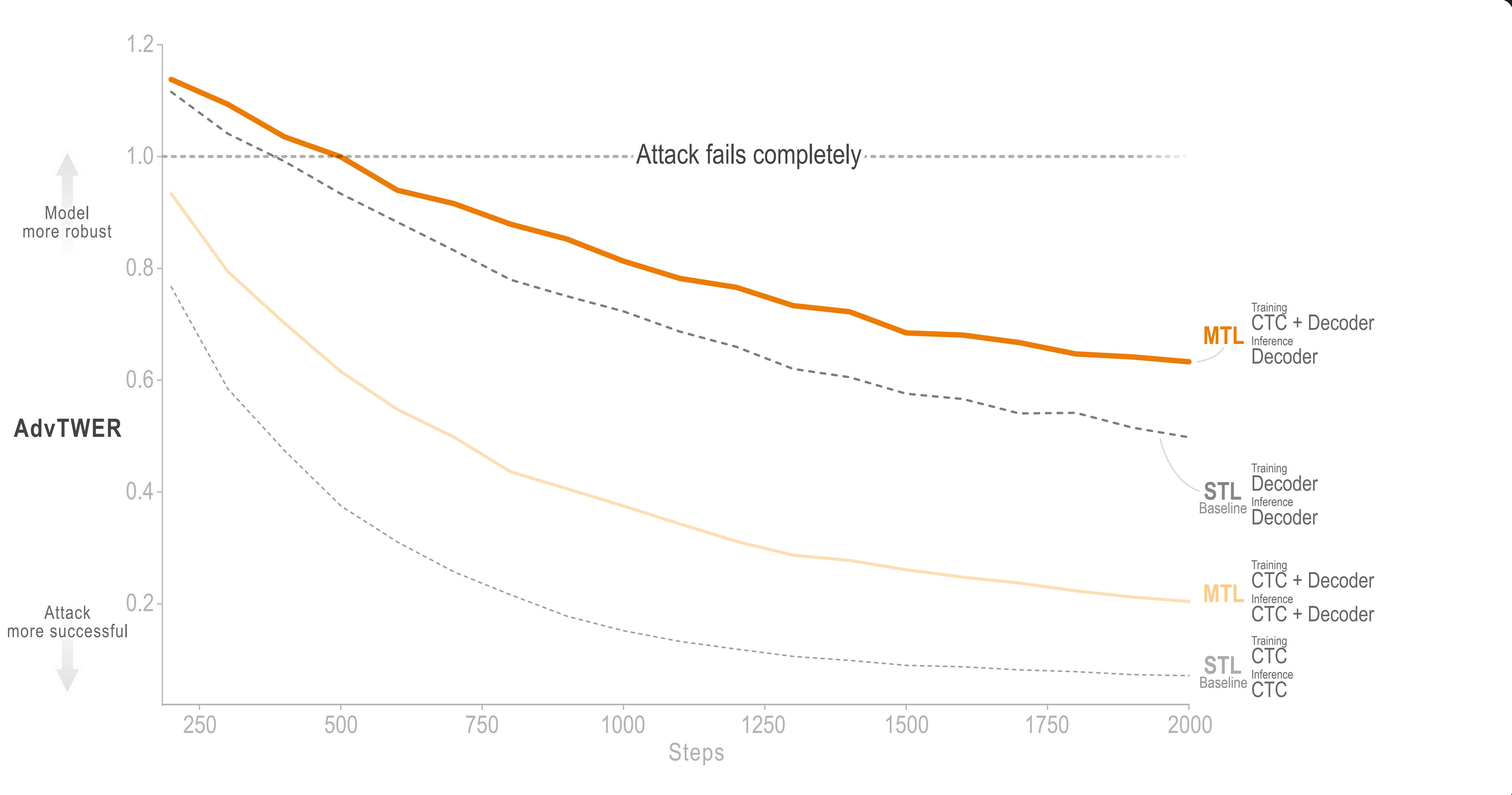


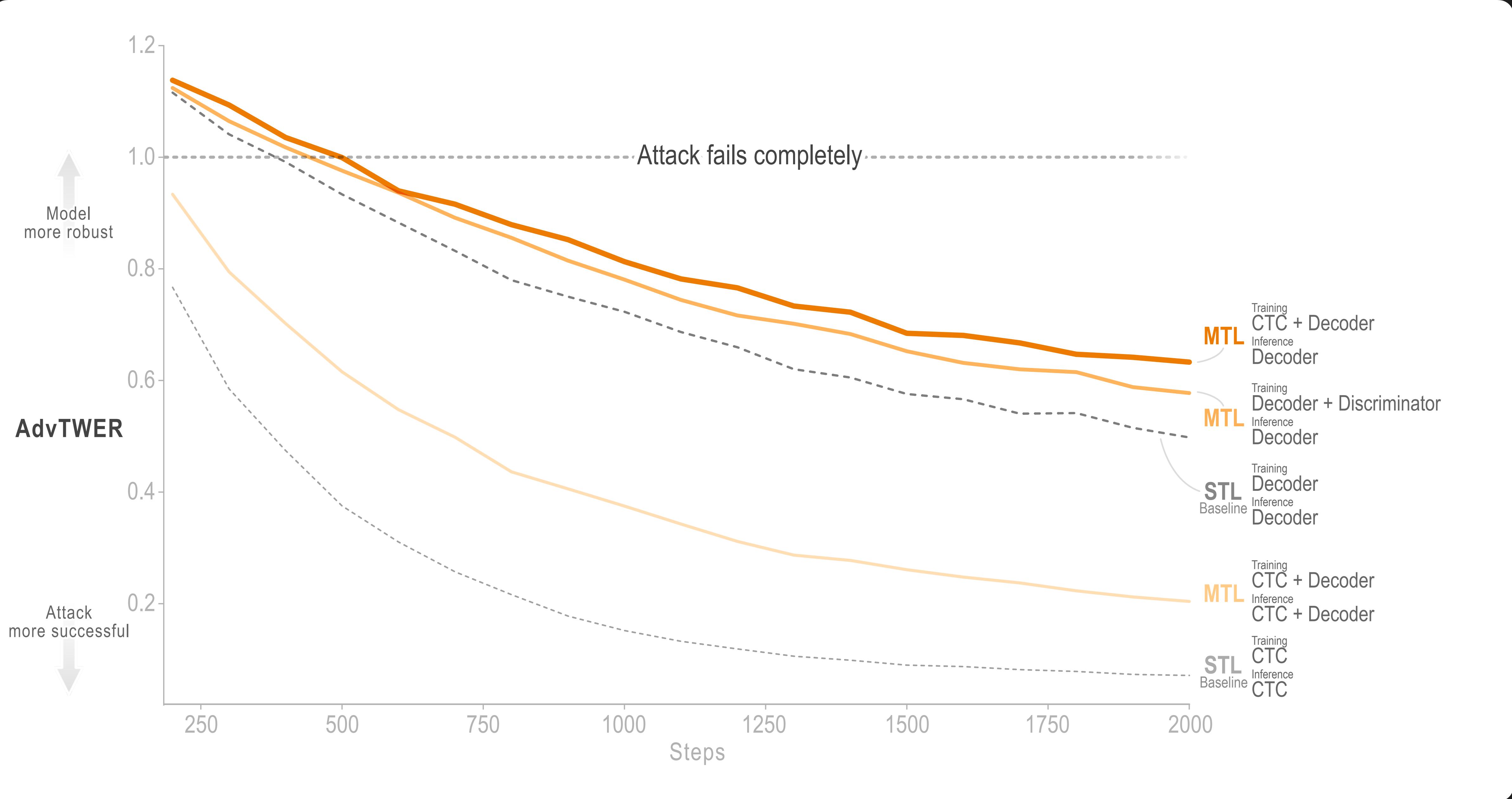
$$\mathcal{L}_{ASR} = \lambda_c \mathcal{L}_{CTC} + (1-\lambda_c) \mathcal{L}_{DEC} \quad \lambda_c = 0.0$$

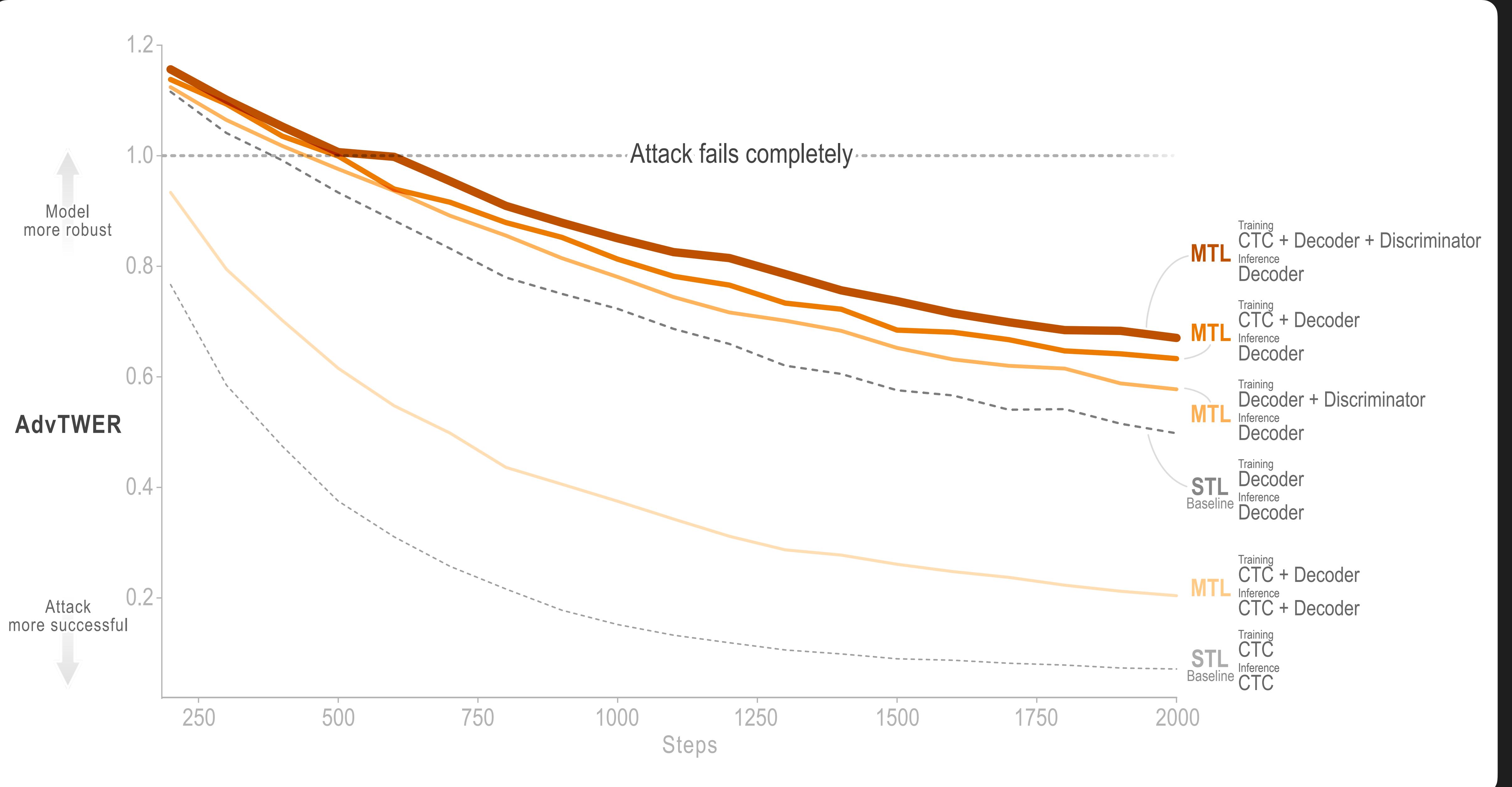
$$\mathcal{L}_{MTL} = \lambda_A \mathcal{L}_{ASR} + (1-\lambda_A) \mathcal{L}_{DIS} \quad \lambda_A = 0.0$$

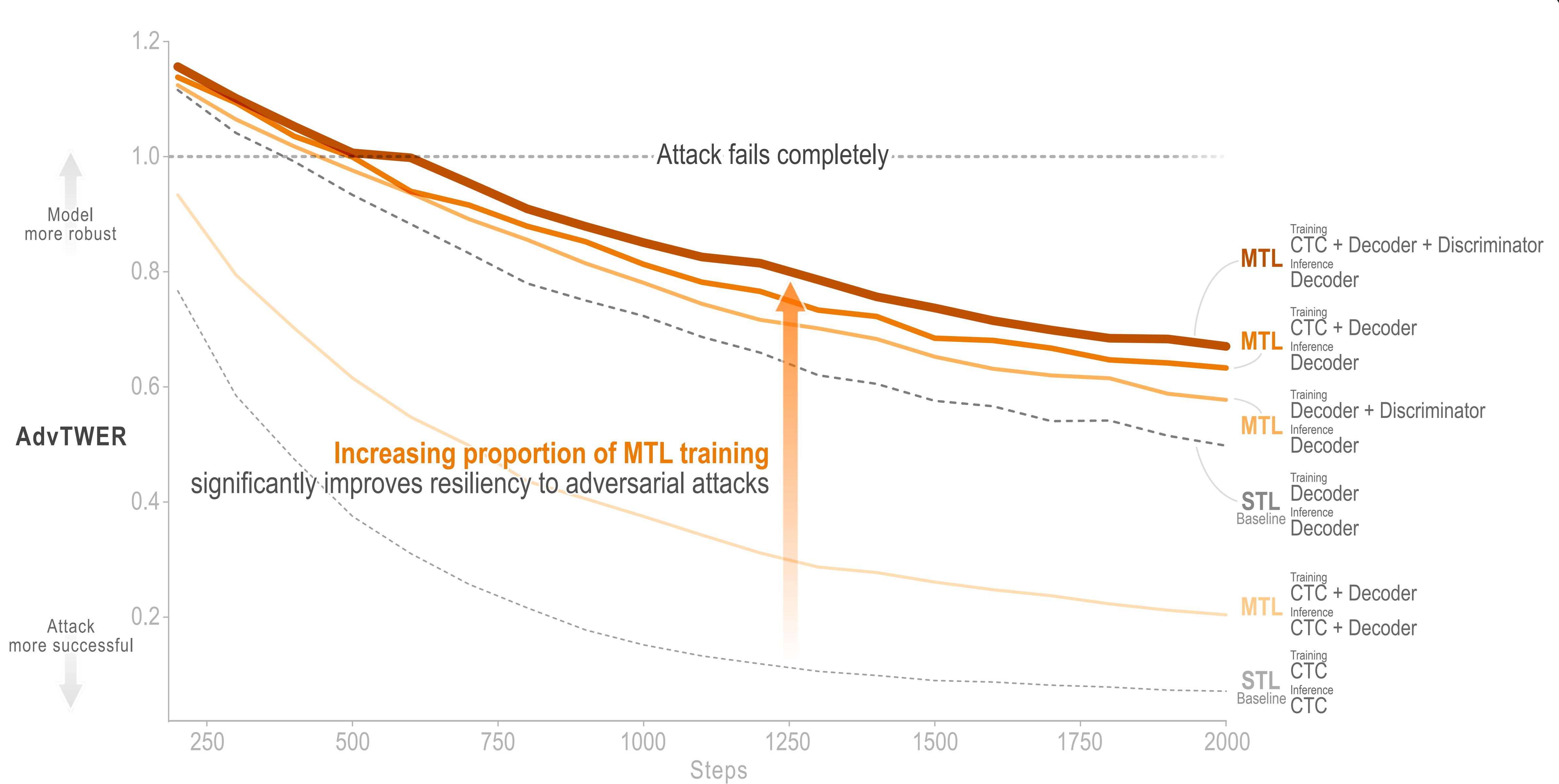






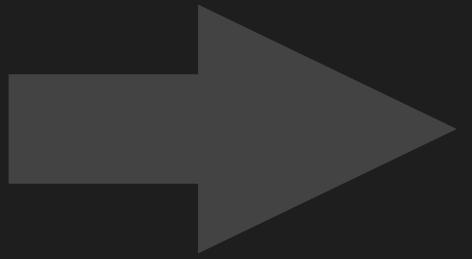
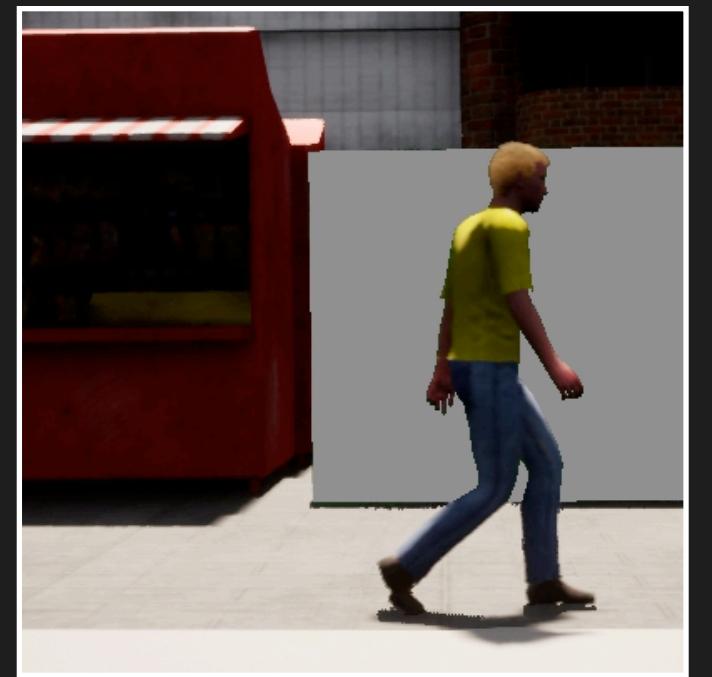






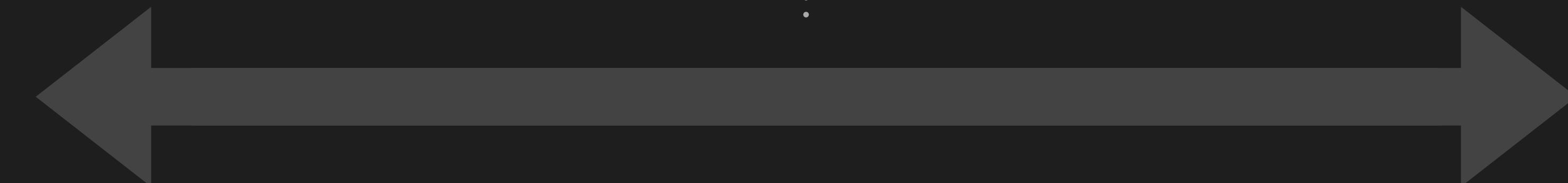
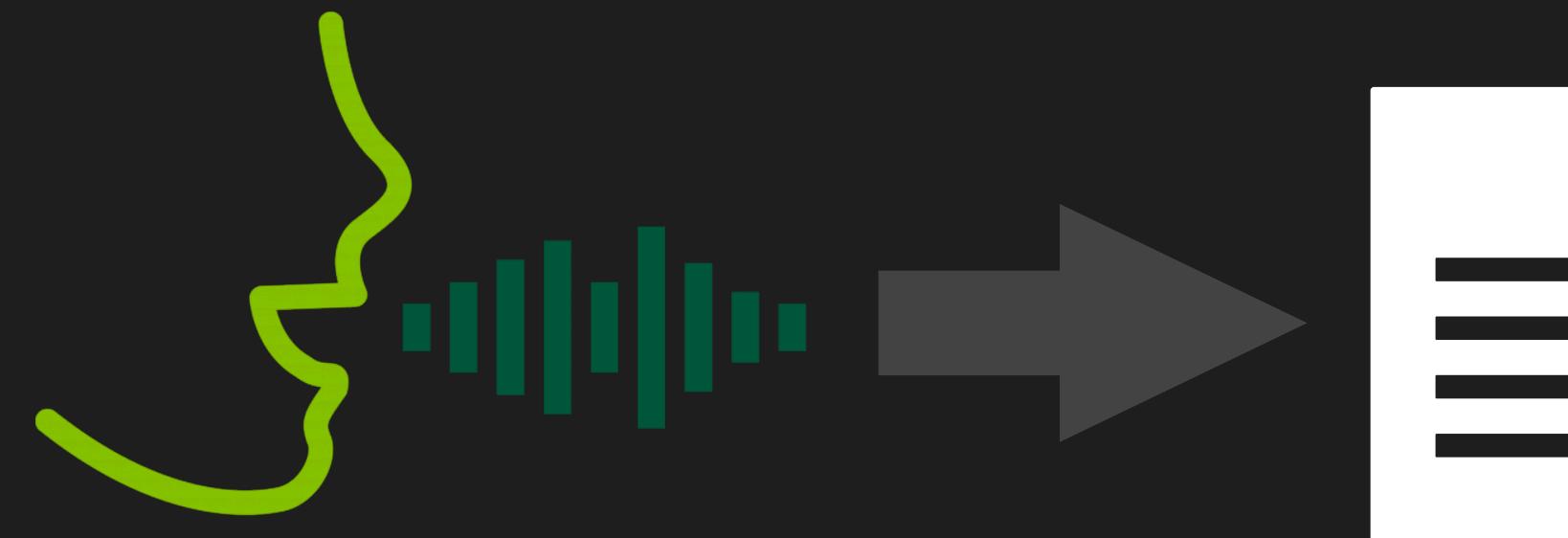
Person Tracking

(video domain)



Automatic Speech Recognition

(audio domain)



Multi-Task Learning

makes DNNs harder to attack **across AI tasks**



Part I
Understand
AI Vulnerabilities

GOGGLES SIGMOD 2020
Bluff IEEE VIS 2020



Part II
Fortify
AI Security

SHIELD KDD 2018
SkeleVision arXiv 2022 (under review)
Hear No Evil arXiv 2022 (under review)



Part III
Enable
Use of AI Security

ADAGIO ECML-PKDD 2018
MLsploit KDD Showcase 2019

ADAGIO

Enabling *Interactive Experimentation*
with Attacks and Defenses for Audio

ECML-KDD 2018



First user interface for security testing with user's examples



Nilaksh Das
Georgia Tech



Madhuri
Shanbhogue
Georgia Tech



Shang-Tse
Chen
Georgia Tech



Li
Chen
Intel Labs



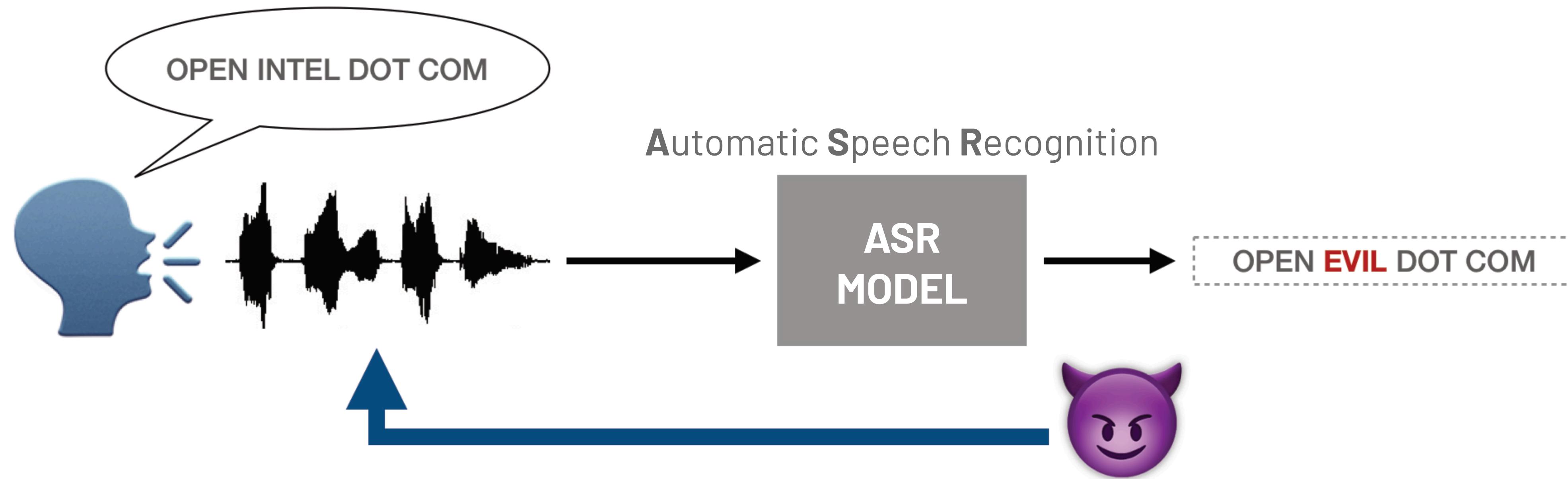
Michael
Kounavis
Intel Labs



Polo
Chau
Georgia Tech

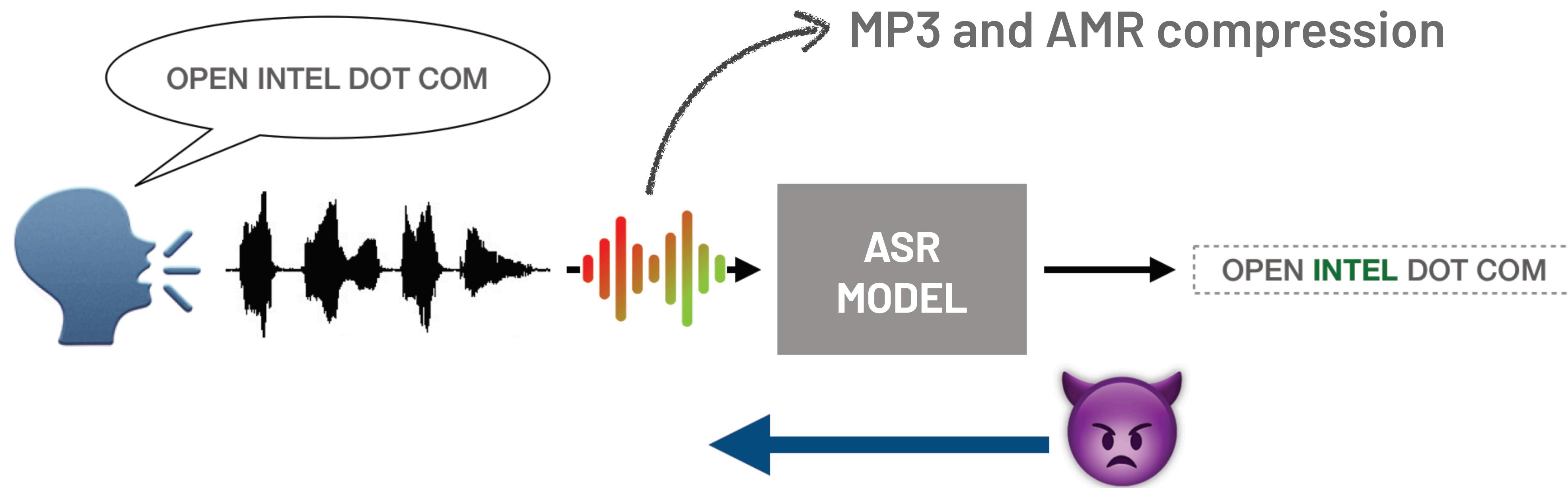


Adversarial Attack on Speech-to-Text



An adversary uses backpropagation to attack the model.

Adversarial Attack on Speech-to-Text



ADAGIO incorporates compression as defense, which blocks the gradient to the attacker.

Defense	WER (no attack)	WER (with attack)	Targeted attack success rate
None	0.369	1.287	92.45%
AMR	0.488	0.666	0.00%
MP3	0.400	0.780	0.00%

Targeted attack is eliminated

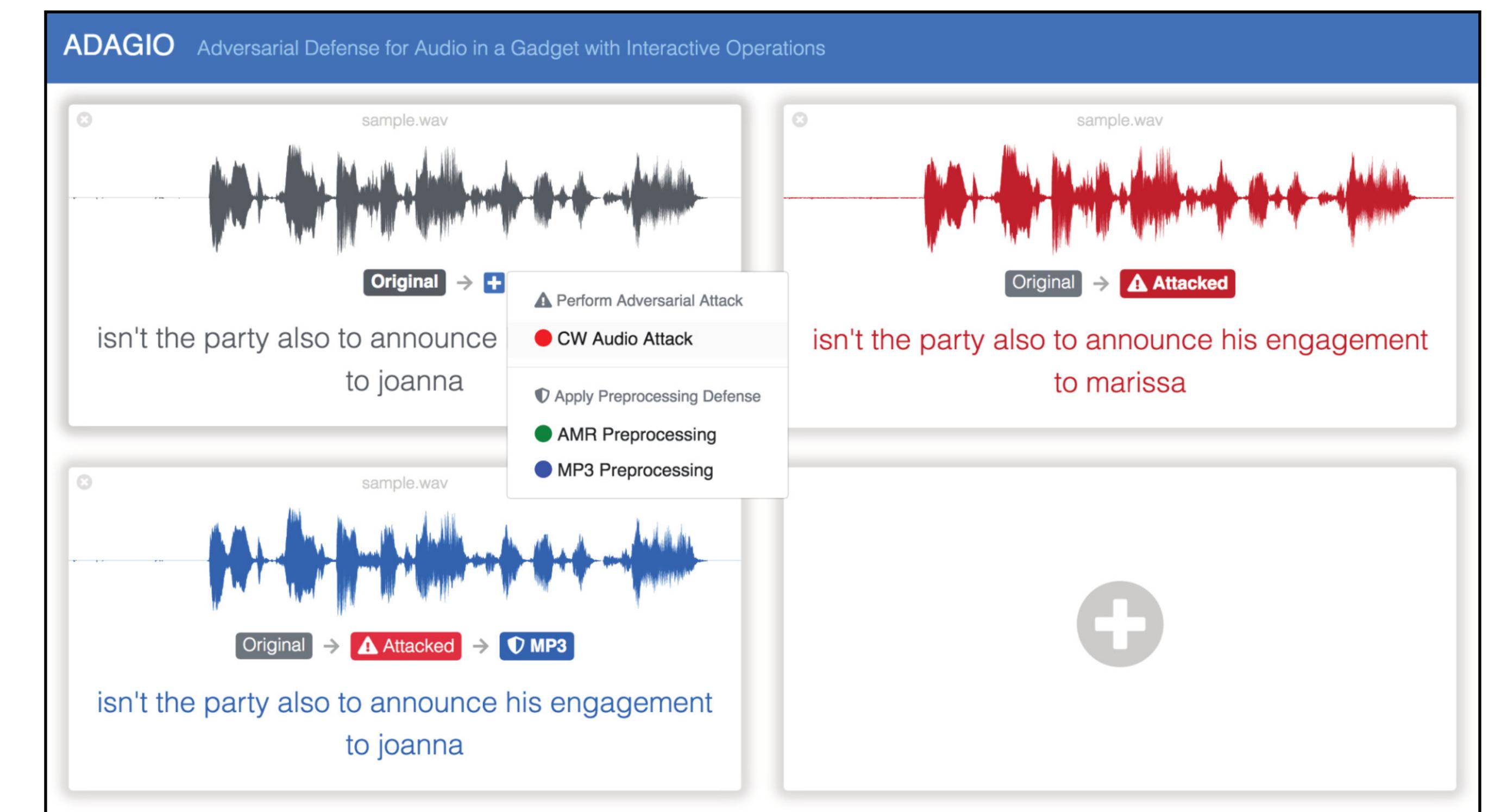
Word Error Rate (WER) and the targeted attack success rate on the DeepSpeech ASR model
(lower is better for both)



ADAGIO

Interactive Experimentation with
Adversarial Attack & Defense for Audio

- Upload your own audio sample
- Perform audio adversarial attack
- Apply compression to defend
- Play audio, listen for differences



ADAGIO = Attack & Defense for Audio in a Gadget with Interactive Operations

ADAGIO

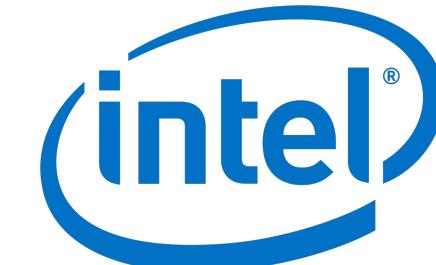


MLSploit

KDD 2019 Project Showcase
Black Hat Asia, Arsenal 2019



Nilaksh Das
Georgia Tech



Interactive Experimentation with Adversarial ML Research

- 🏢 Employed in teaching at Georgia Tech
- ↗ Tech-transferred to Intel Labs
- 🌐 Open-sourced at github.com/mlsploit

Contributors from

Intel Science and Technology Center

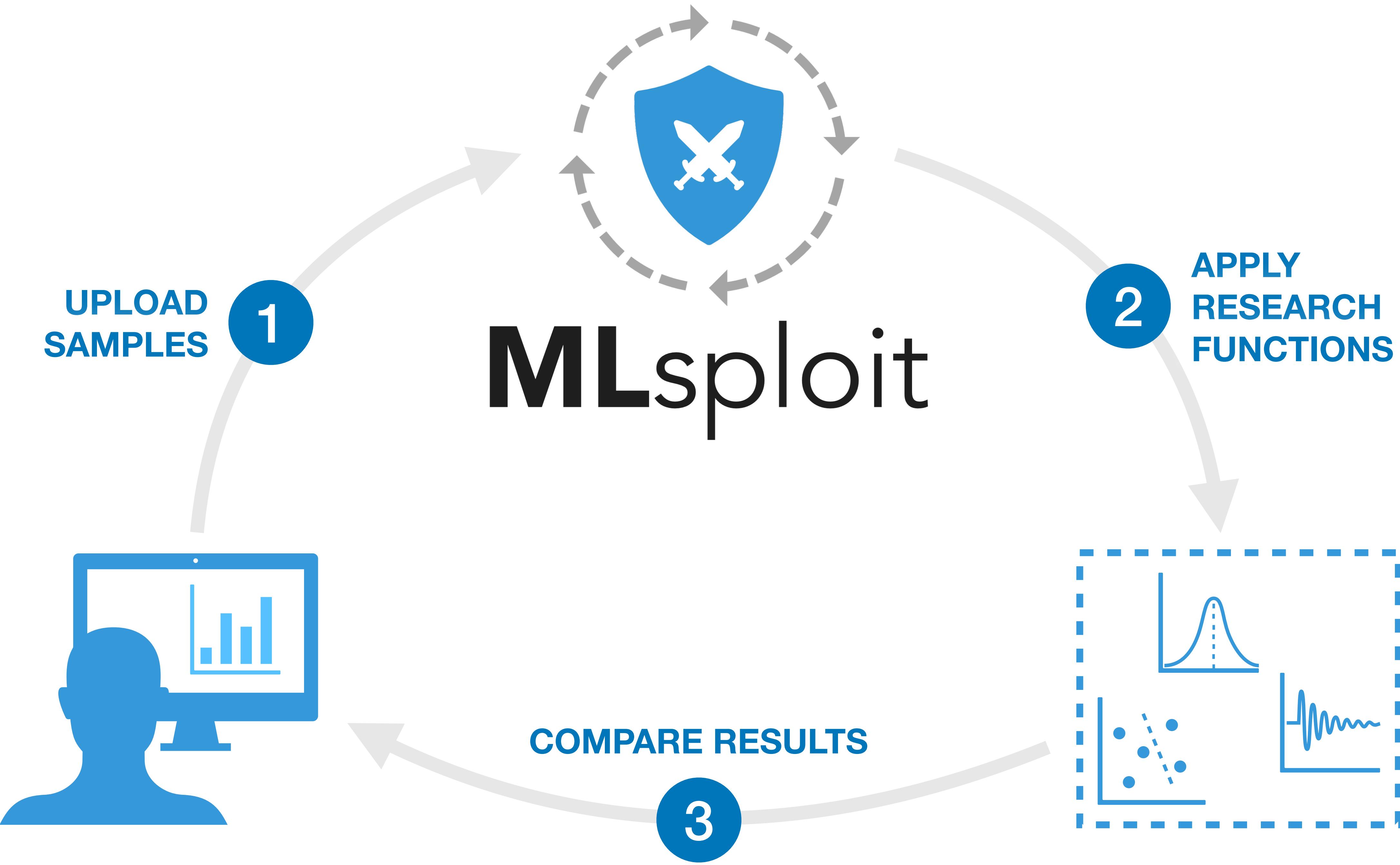
for Adversary-Resilient Security Analytics:

Siwei Li, Chanil Jeon, Jinho Jung*, Shang-Tse Chen*,
Carter Yagemann*, Evan Downing*, Haekyu Park,
Evan Yang, Li Chen, Michael Kounavis, Ravi Sahita,
David Durham, Scott Buck,
Polo Chau, Taesoo Kim, Wenke Lee

*equal contribution

MLsploit

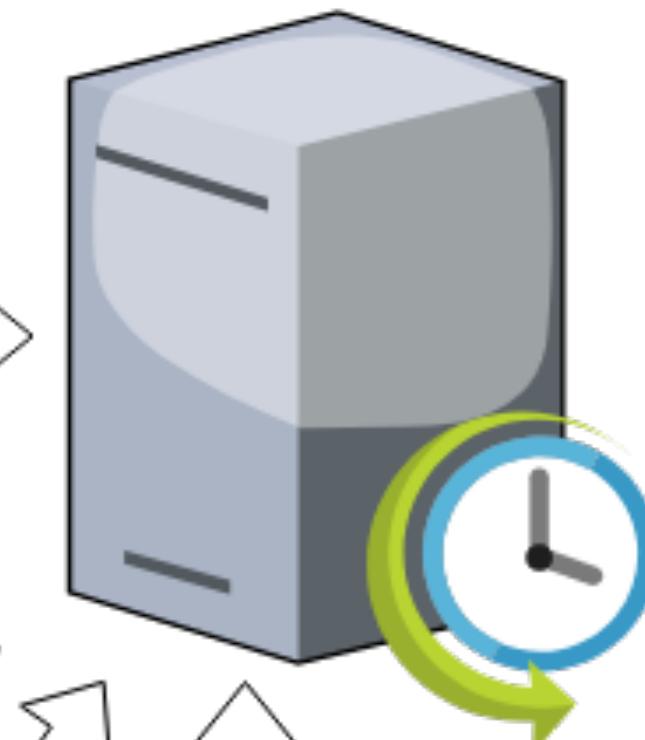
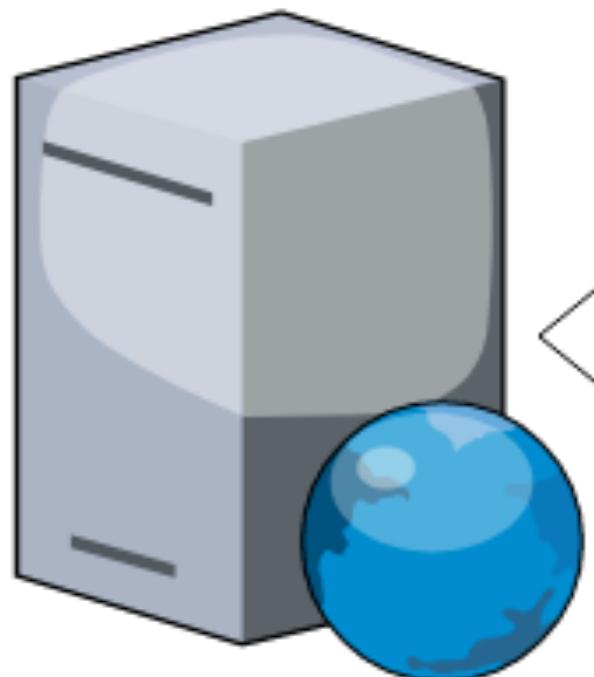
- ★ **Research modules** for adversarial ML
 - * Enables **comparison** of attacks and defenses
- ★ **Interactive experimentation** with ML research
- ★ Researchers can **easily integrate** novel research into an intuitive and seamless **user interface**



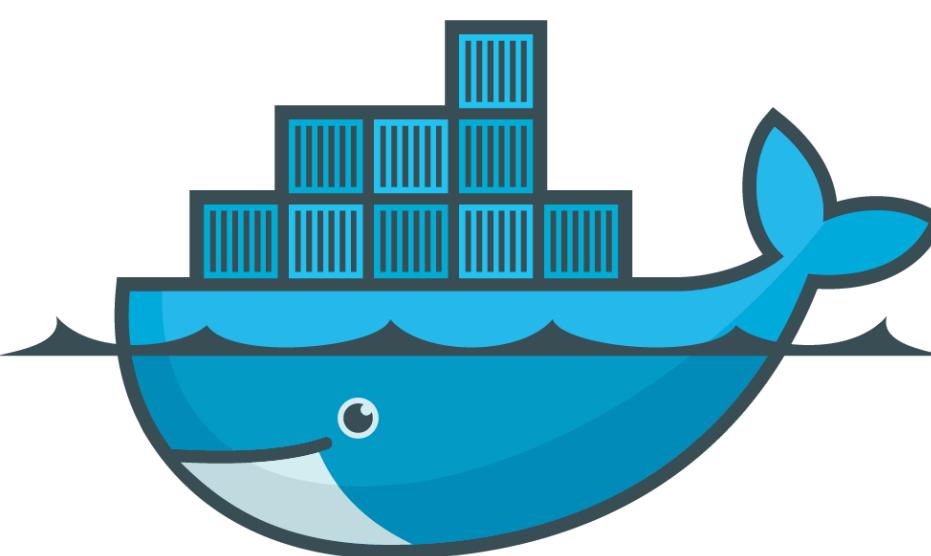
Web
Portal

RESTful
API

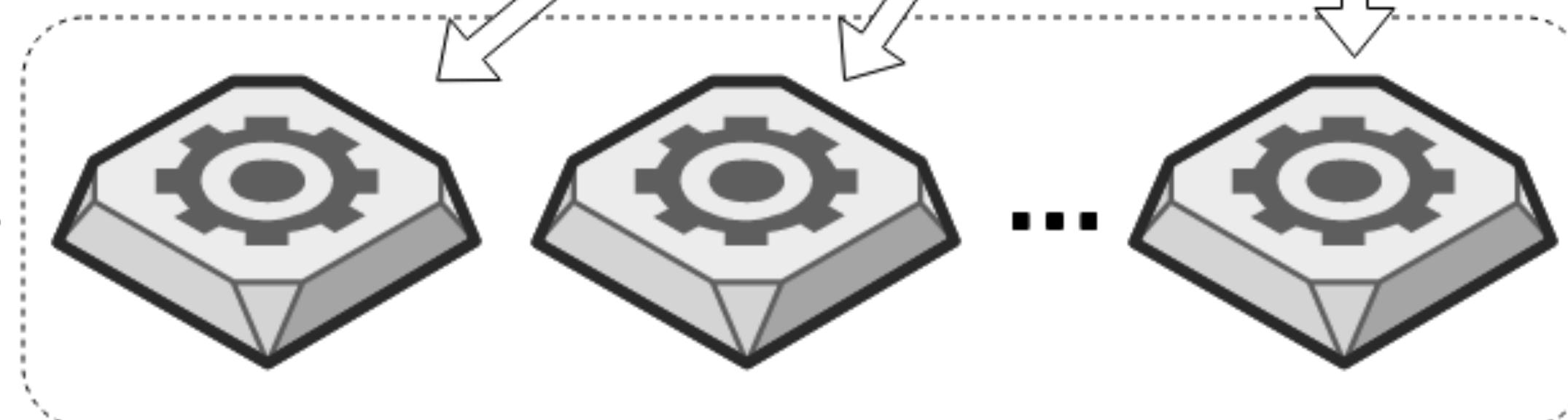
Job
Scheduler



MLsploit ARCHITECTURE

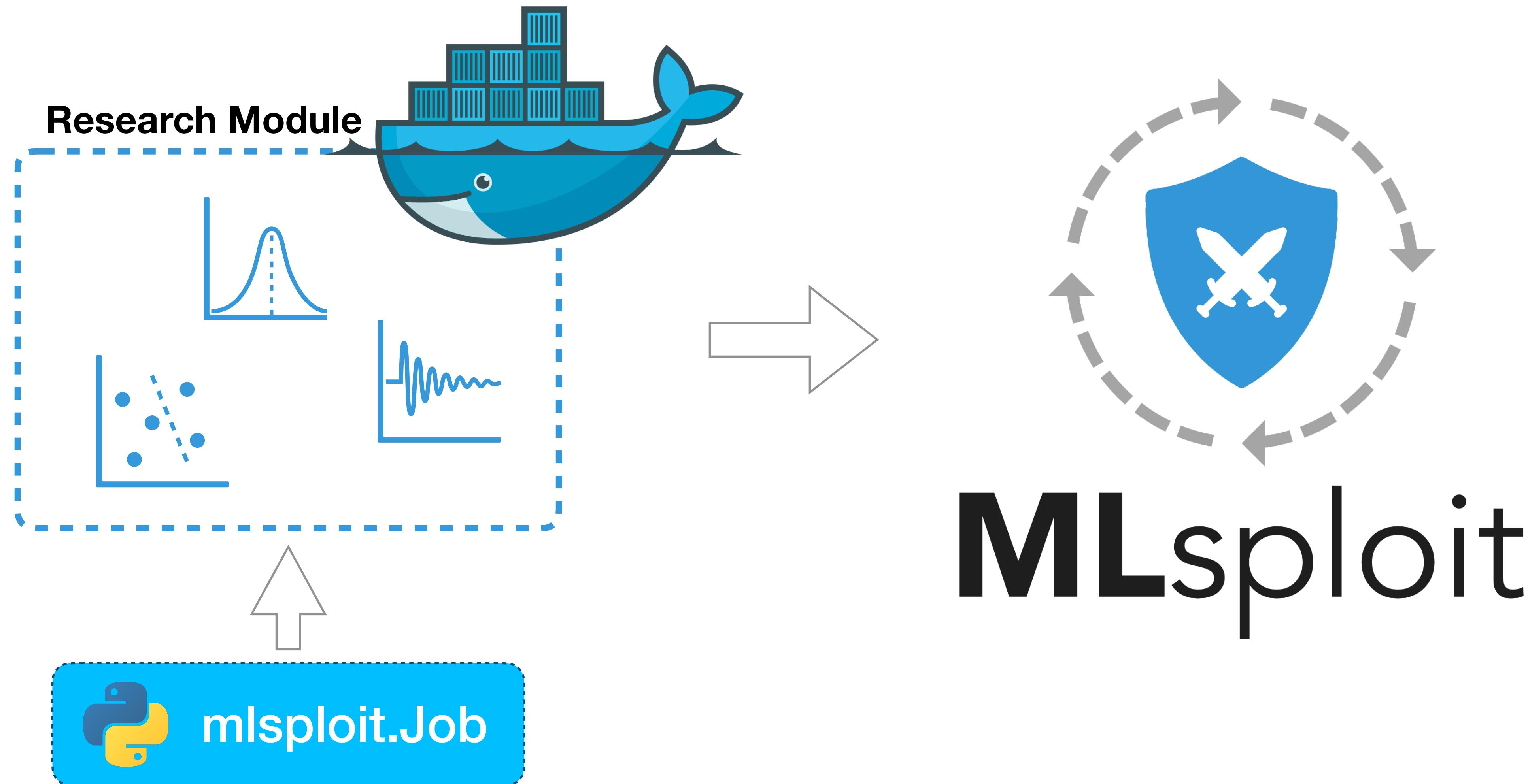


Research
Modules



Worker
Instances

EASY INTEGRATION OF RESEARCH







MLsploit

Experiment with AI Security in your browser!

Research Modules

**SHIELD**

Fast, Practical Defense for Deep Learning

**Foolbox**

Fool neural networks!

**AVPASS**

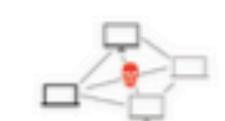
Android Malware Detection Bypass

**Barnum**

Deep Learning Software Anomaly Detection

**ELF**

ELF File Malware Detection and Bypassing

**Network**

Network Intrusion Detection and Evasion

Attack and Defend



→

FGSM

→



→

SLQ



Attack



→

FGSM

→



Classify



→

Classify

→





MLsploit

Experiment with AI Security in your browser!

Research Modules

**SHIELD**

Fast, Practical Defense for Deep Learning

**Foolbox**

Fool neural networks!

**AVPASS**

Android Malware Detection Bypass

**Barnum**

Deep Learning Software Anomaly Detection

**ELF**

ELF File Malware Detection and Bypassing

**Network**

Network Intrusion Detection and Evasion

Attack and Defend



→

FGSM

→



→

SLQ



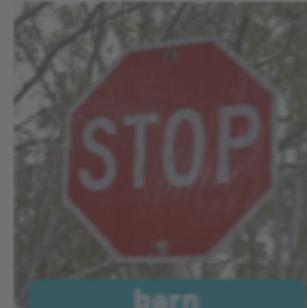
Attack



→

FGSM

→



Classify



→

Classify

→





MLsploit

Experiment with AI Security in your browser!

Research Modules



SHIELD

Fast, Practical Defense for Deep Learning



Foolbox

Fool neural networks!

users can create multiple pipelines with different research functions



Barnum

Deep Learning Software Anomaly Detection



ELF

ELF File Malware Detection and Bypassing



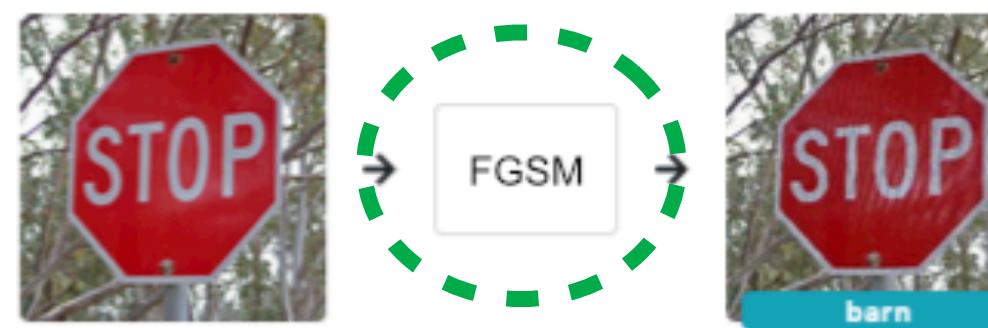
Network

Network Intrusion Detection and Evasion

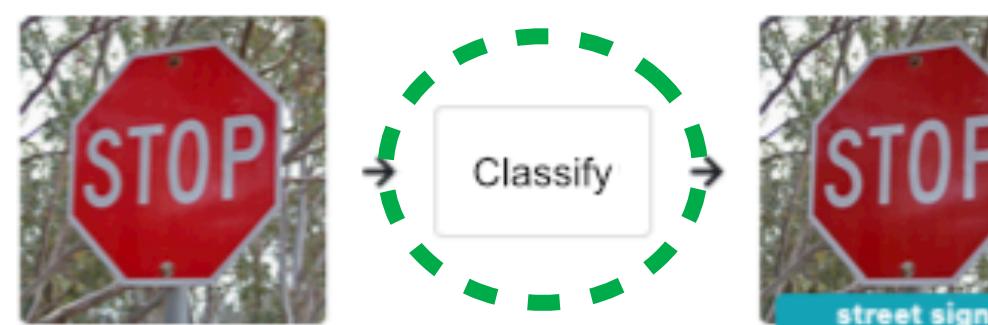
Attack and Defend



Attack



Classify





MLsploit

Experiment with AI Security in your browser!

Research Modules



SHIELD

Fast, Practical Defense for Deep Learning



Foolbox

Fool neural networks!



AVPASS

Android Malware Detection Bypass



Barnum

Deep Learning Software Anomaly Detection



ELF

ELF File Malware Detection and Bypassing

... and experiment with
their own samples!

Attack and Defend



→ FGSM



→ SLQ



Attack



→ FGSM



Classify



→ Classify





Research Modules



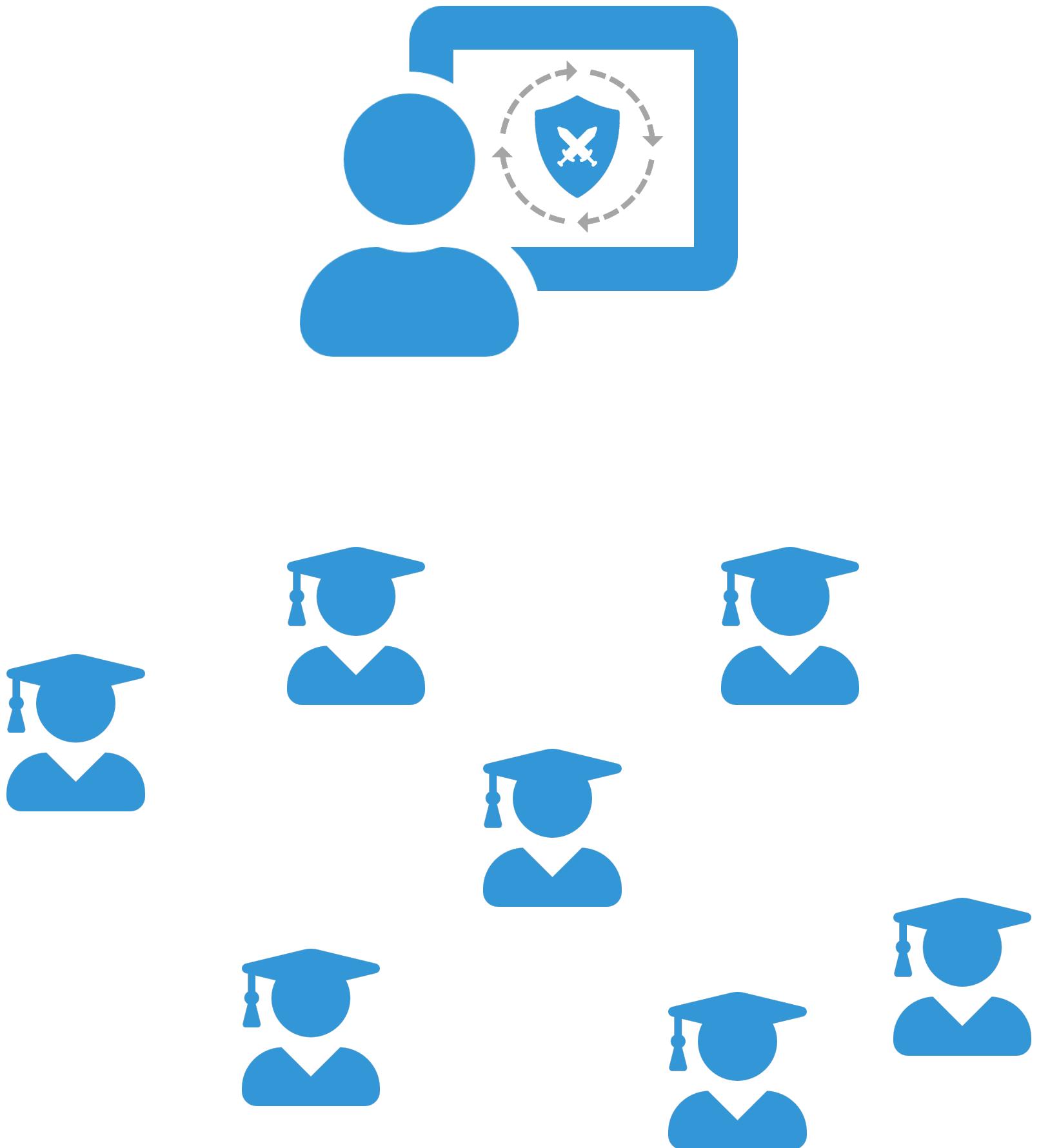
Foolbox

Fool neural networks!

Expand the “**Foolbox**” Research Module
to see various research functions

IMPACT

MLsploit employed at **GT**
to teach students how to
train, evade and defend
DNNs for malware detection

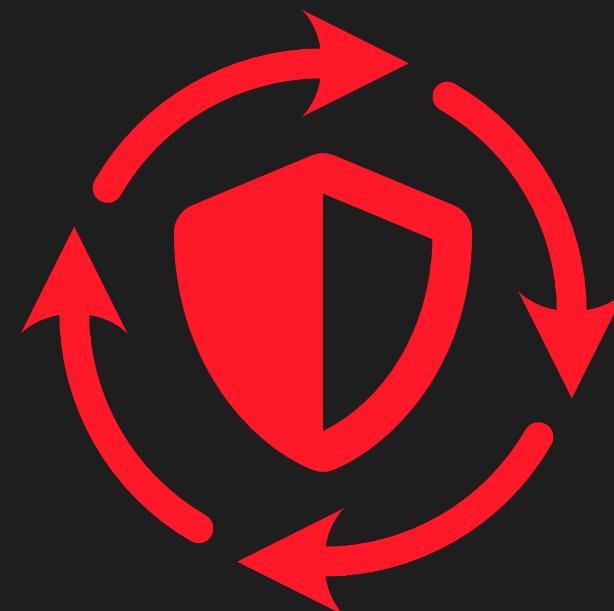




Part I

Understand AI Vulnerabilities

GOGGLES SIGMOD 2020
Bluff IEEE VIS 2020



Part II

Fortify AI Security

SHIELD KDD 2018
SkeleVision arXiv 2022 (under review)
Hear No Evil arXiv 2022 (under review)



Part III

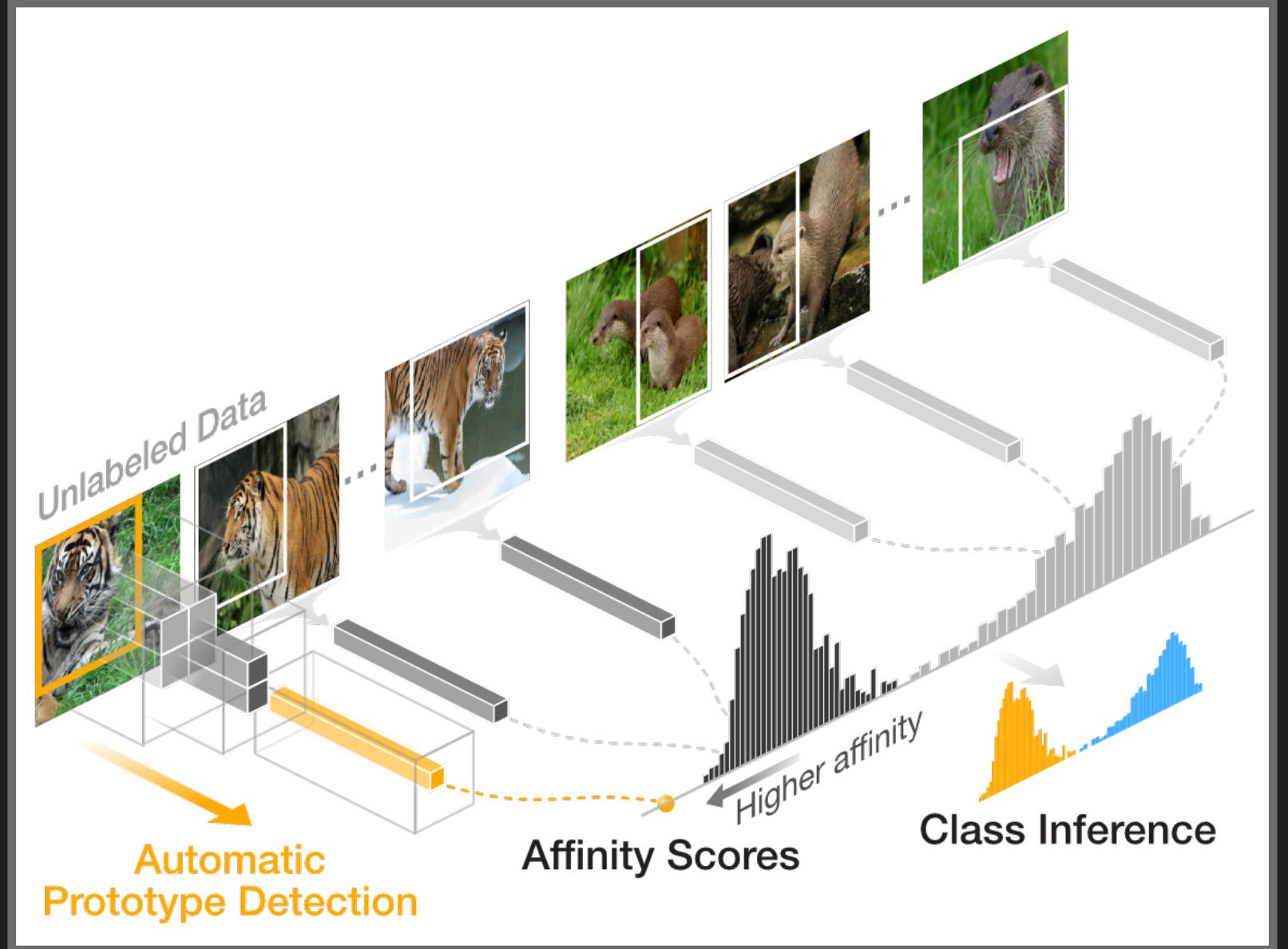
Enable Use of AI Security

ADAGIO ECML-PKDD 2018
MLsploit KDD Showcase 2019

Research Contributions

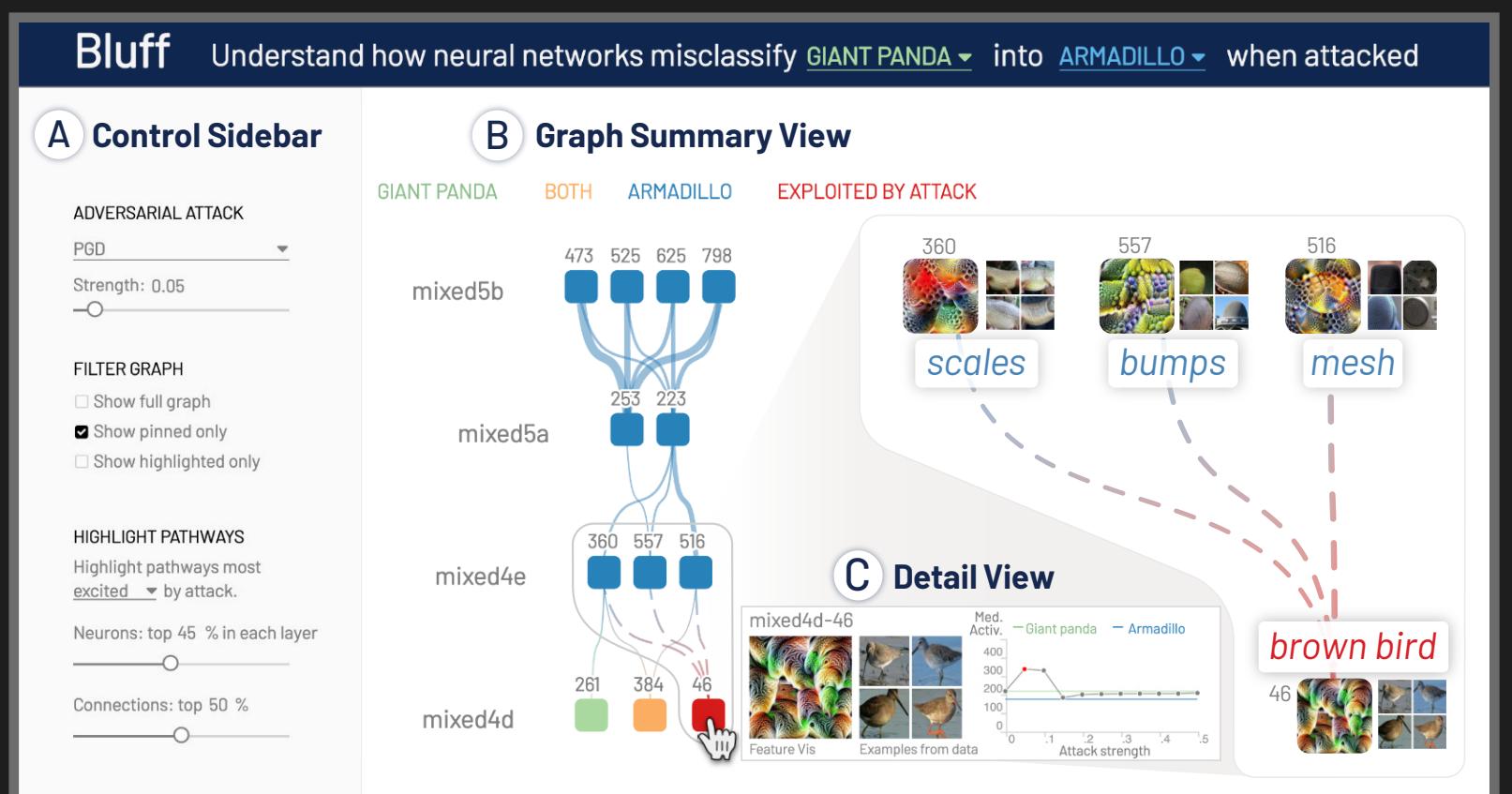
Novel principled approaches

- GOGGLES contributes novel, theoretically principled approaches to extract **interpretable prototypes**
- weakly-supervised GOGGLES framework is only **7% away** from the fully-supervised baseline



Novel visualization technique for deciphering attacks

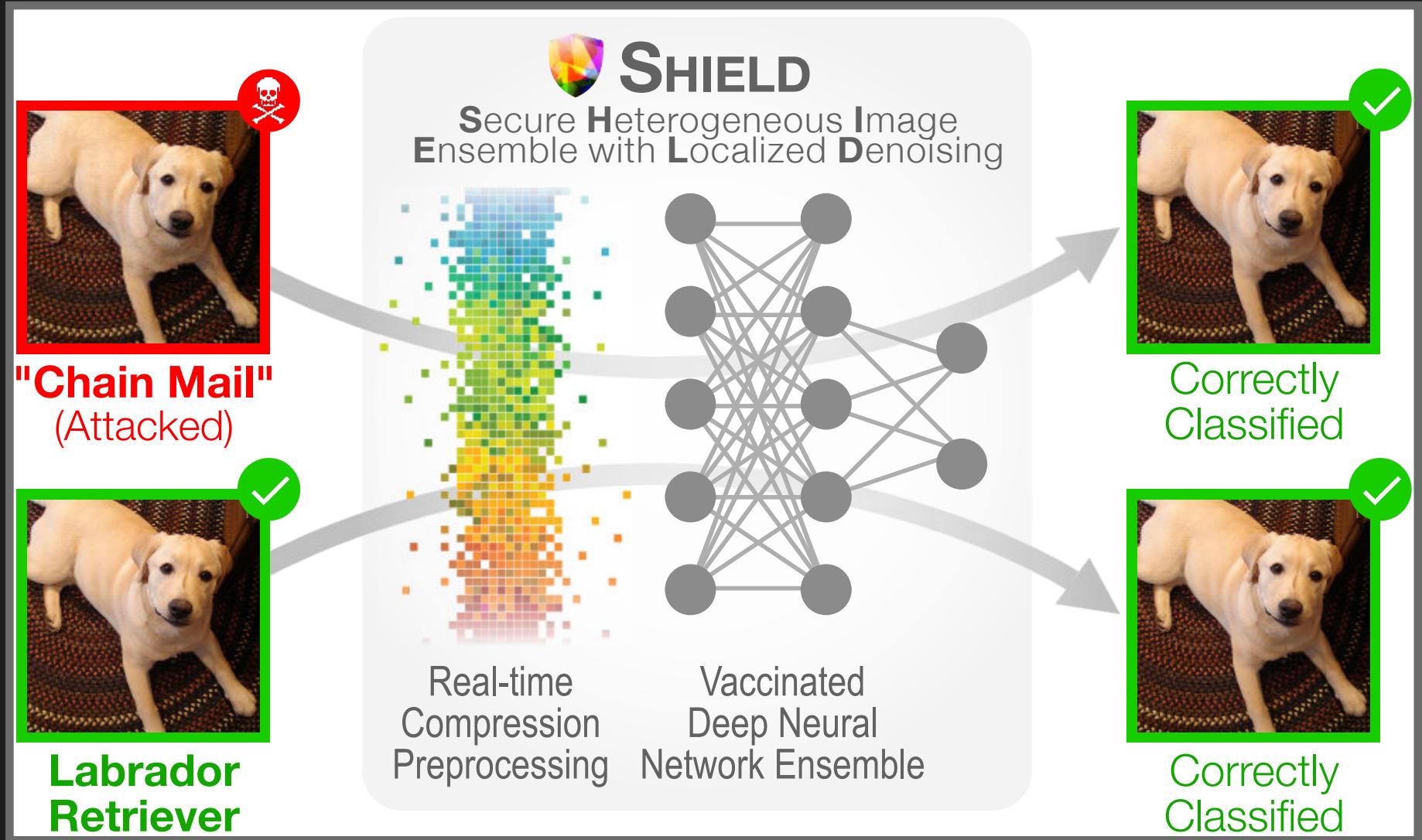
- Bluff helps in identifying the **vulnerable neurons** in a DNN that are non-performant and lead to misclassifications under attack



Research Contributions

Faster, generalizable defenses

- SHIELD defense is up to **22x faster** than other preprocessing defenses
- ADAGIO defense removes **all** non-adaptive targeted adversarial attacks



Fundamentally unifying approach for training robust models across AI tasks

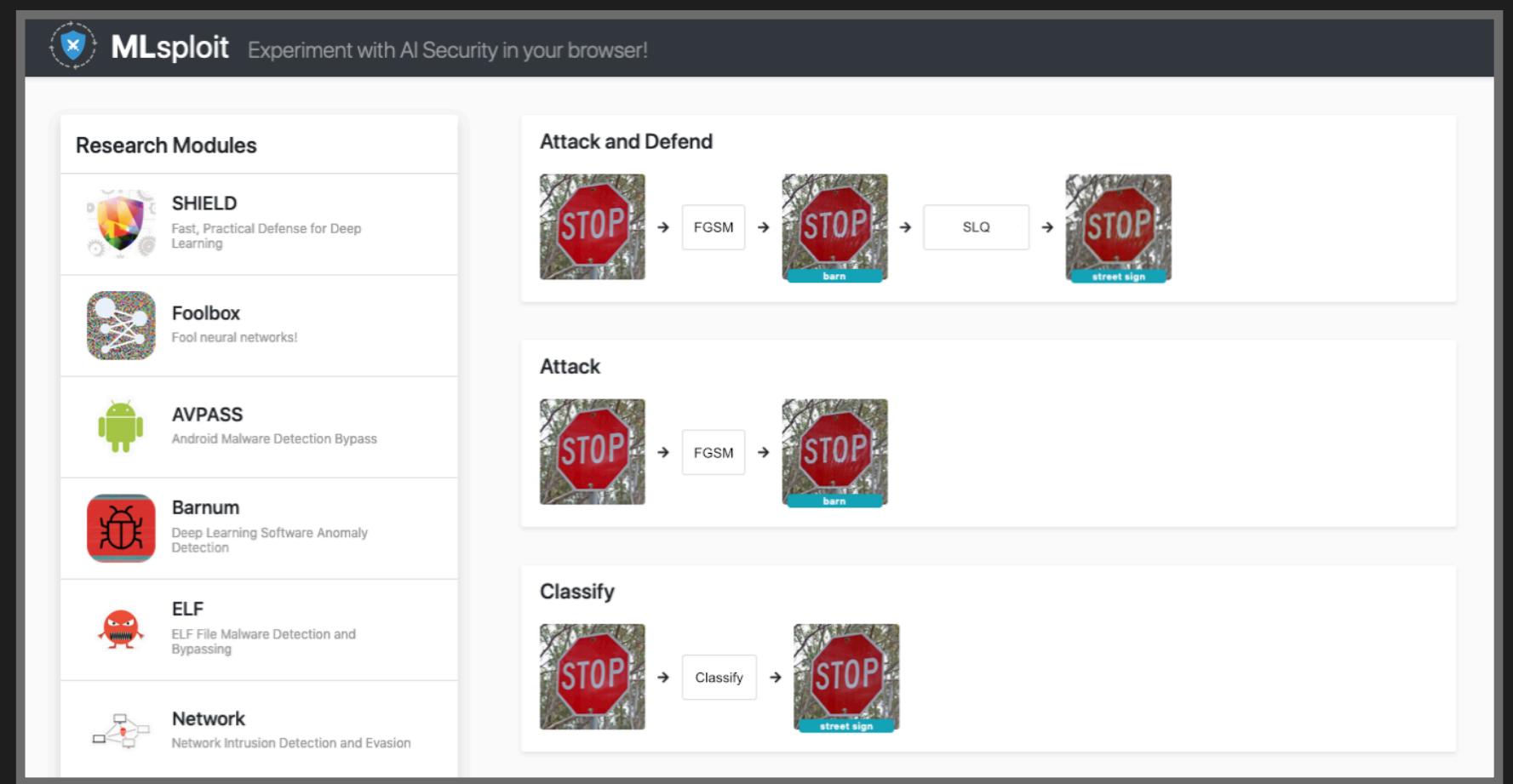
- Our in-depth study of **Multi-Task Learning** establishes it as a fundamentally unifying approach that helps learn robust features for resisting adversarial attacks **across AI tasks**



Research Contributions

First scalable system for interactive Adversarial ML research

- MLsploit is the first open-source interactive system that allows in-depth security testing of AI models, and lowers barriers to entry for everyone



Understanding, Fortifying and Democratizing AI Security



Nilaksh Das
nilakshdas.com

