### TOWARDS

# SECURE & INTERPRETABLE AI

SCALABLE METHODS, INTERACTIVE VISUALIZATIONS, PRACTICAL TOOLS



### Polo Chau

Associate Professor Associate Director, MS Analytics Georgia Tech

poloclub.github.io

#### **Polo Club of Data Science** poloclub.github.io

### **Human-Centered Al**



#### **ActiVis**

Visual Exploration of Facebook Deep Neural Network Models







#### GAN Lab

Playing with Generative Adversarial Networks in Browser



### **Cyber Security**



### Cyber MoneyBall

Predicting Cyber Threats with Vi Security Products





#### MARCO

Fake Review Detection

SDM'14 Best Student Paper

### Adversarial ML



### SHIELD

Fast, practical defense for deep learning





#### ShapeShifter

1st Targeted Physical Attack on Faster R-CNN Object Detector

### **Large Graph Mining & Visualization**



### MMap

Easy billion-scale graph computation on a PC using virtual memory

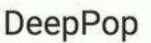


### Apolo

Explore million-node graphs in real time

### **Social Good & Health**





Deep Learning on Satellite Imagery for Population Estimation

Microsoft Al for Earth



### Firebird

Predicting Fire Risk in Atlanta

TKDD'16 Best Student Paper, runner-up



Atlanta Fire Rescue Department

### Polo Club of Data Science





Scalable interactive tools to make sense of complex large-scale datasets and models

































#### **Polo Club of Data Science** poloclub.github.io

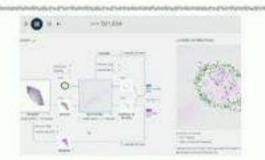
### **Human-Centered Al**



#### **ActiVis**

Visual Exploration of Facebook Deep Neural Network Models





#### GAN Lab

Playing with Generative Adversarial Networks in Browser



### **Cyber Security**



### Cyber MoneyBall

Predicting Cyber Threats with Vi Security Products





### MARCO

Fake Review Detection

SDM'14 Best Student Paper

### Adversarial ML



### SHIELD

Fast, practical defense for deep learning

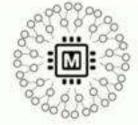




#### ShapeShifter

1st Targeted Physical Attack on Faster R-CNN Object Detector

### **Large Graph Mining & Visualization**



### **MMap**

Easy billion-scale graph computation on a PC using virtual memory



### Apolo

Explore million-node graphs in real time

### Social Good & Health





Deep Learning on Satellite Imagery for Population Estimation





### Firebird

Predicting Fire Risk in Atlanta

TKDD'16 Best Student Paper, runner-up



Atlanta Fire Rescue Department

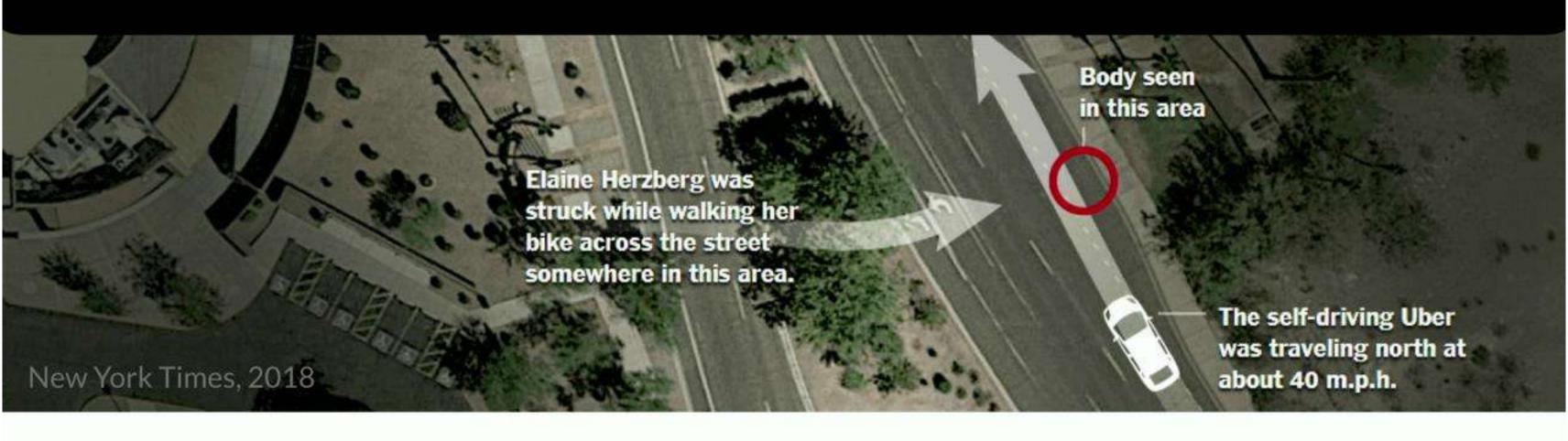
### Today's Main Topics

Secure

Interpretable

Why focus on them? How are they related?

## Al now used in safety-critical applications. Important to study threats & countermeasures.



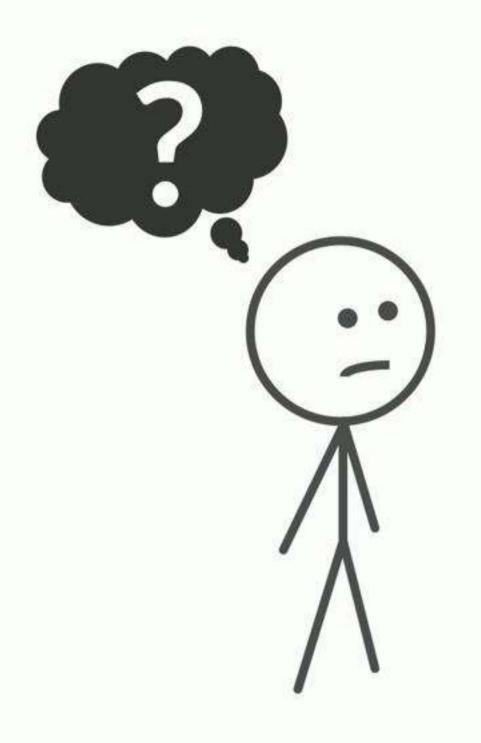
## How a Self-Driving Uber Killed a Pedestrian in Arizona

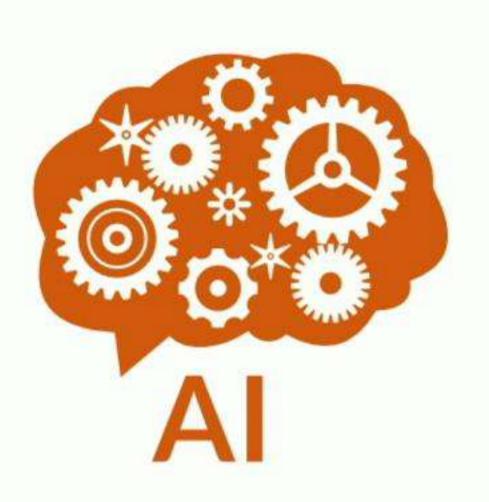
## Al now used in safety-critical applications. Important to study threats & countermeasures.



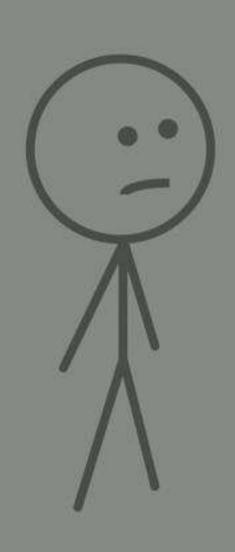
## How a Self-Driving Uber Killed a Pedestrian in Arizona

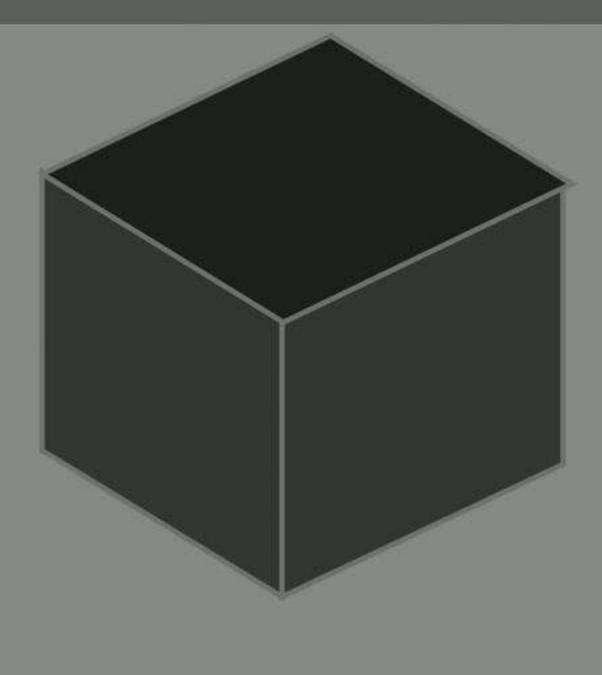
### How do we know if a defense for Al is working?

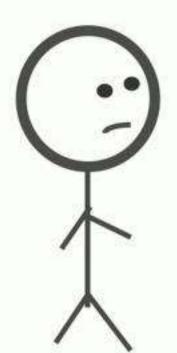


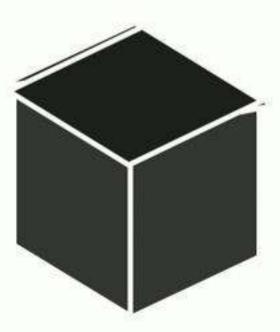


## Al models often used as black-box

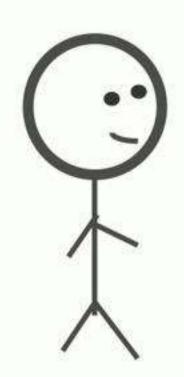








## Interpretable





## Interpretable

Via scalable, interactive, usable interfaces to help people understand complex, large-scale ML systems.



## Secure

## Interpretable

Attack & Defense (DNN)

ShapeShifter

SHIELD

Do-it-yourself Adversarial ML

ADAGIO

MLsploit

Understand Industry Models
ActiVis

Interactive Learning (Education)

GAN Lab

Research landscape

Survey, Gamut

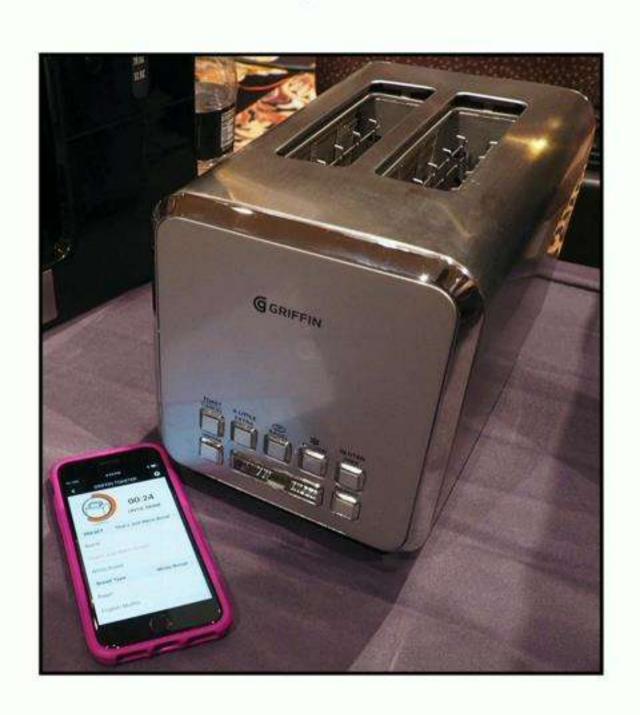


"THE TOASTER HAS BEEN HACKED INTO THINKING IT'S A BLENDER,"

### Al Security Problems Are Everywhere

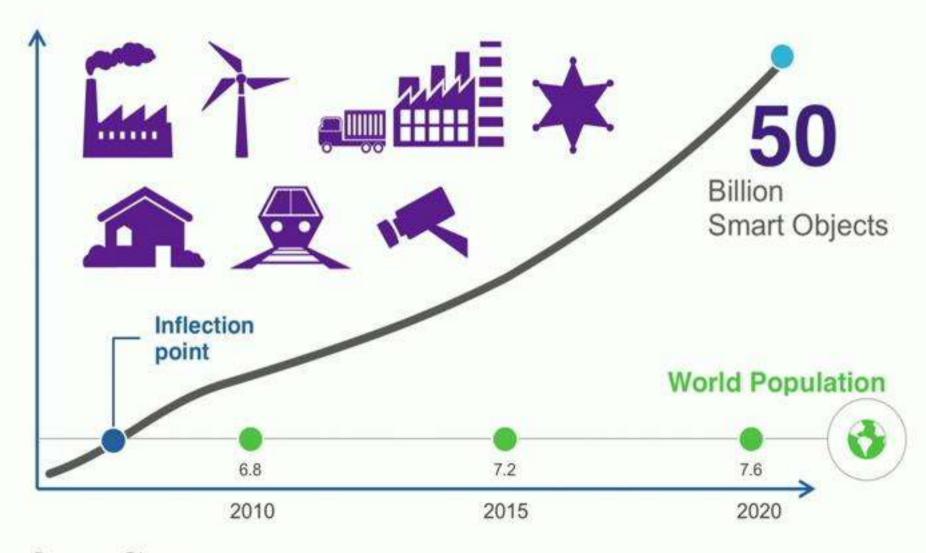


"THE TOASTER HAS BEEN HACKED INTO THINKING IT'S A BLENDER,"



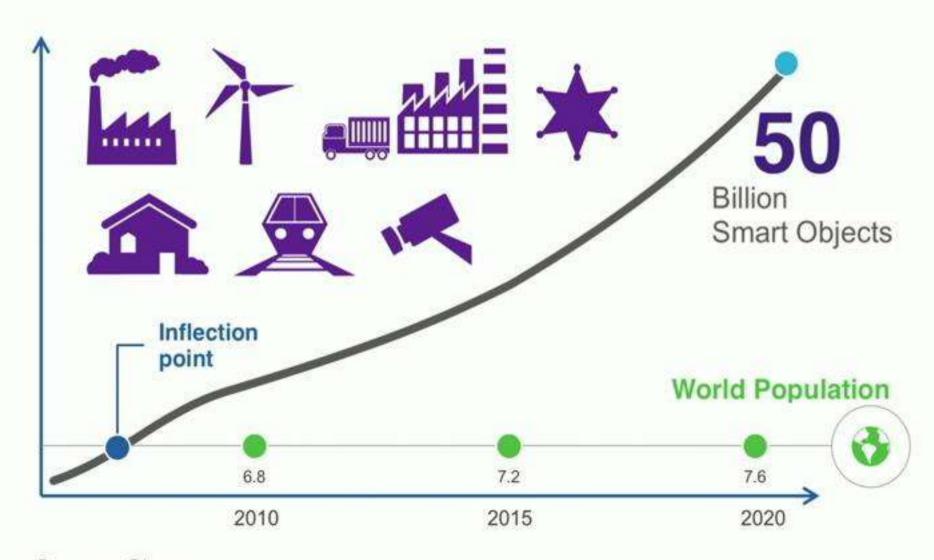
Smart toaster does exist!

## Al Security is becoming increasingly important



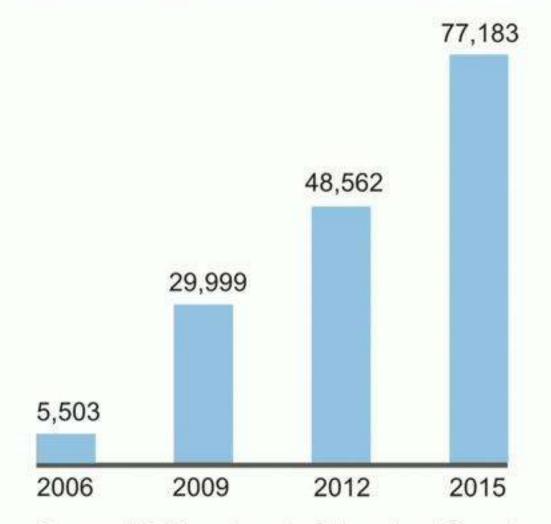
Source: Cisco

## Al Security is becoming increasingly important



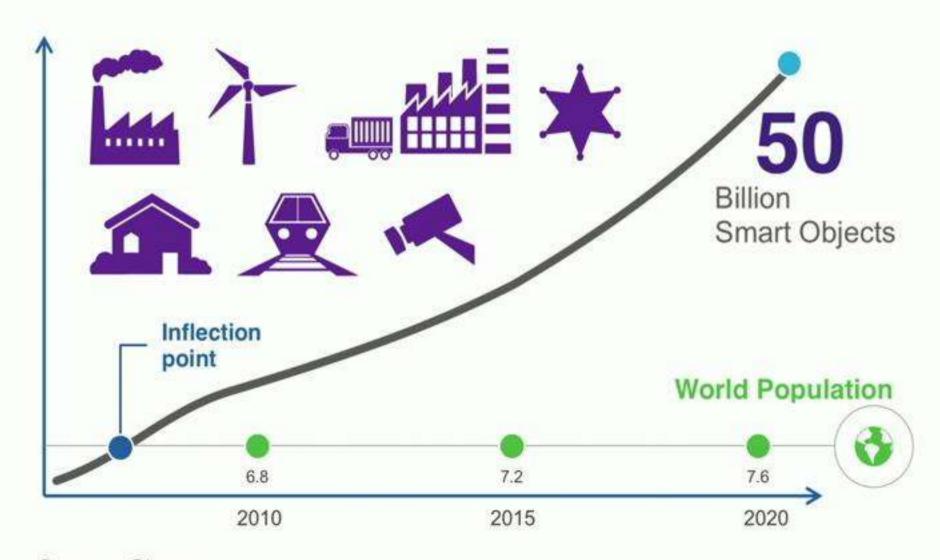
Source: Cisco

# incidents
reported by U.S. federal agencies



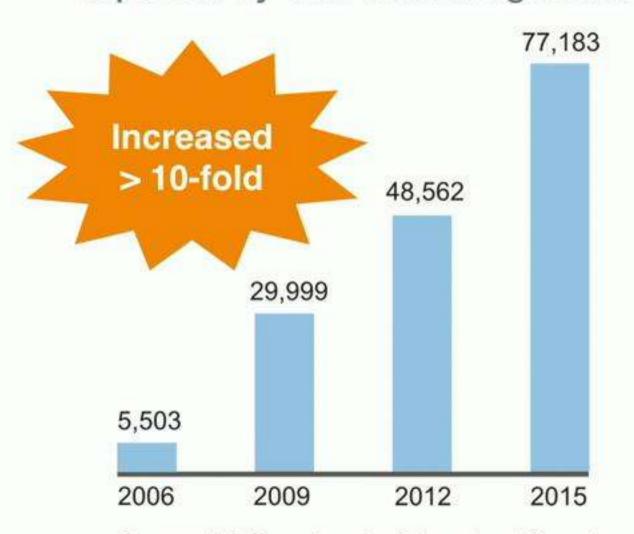
Source: US Department of Homeland Security

## Al Security is becoming increasingly important



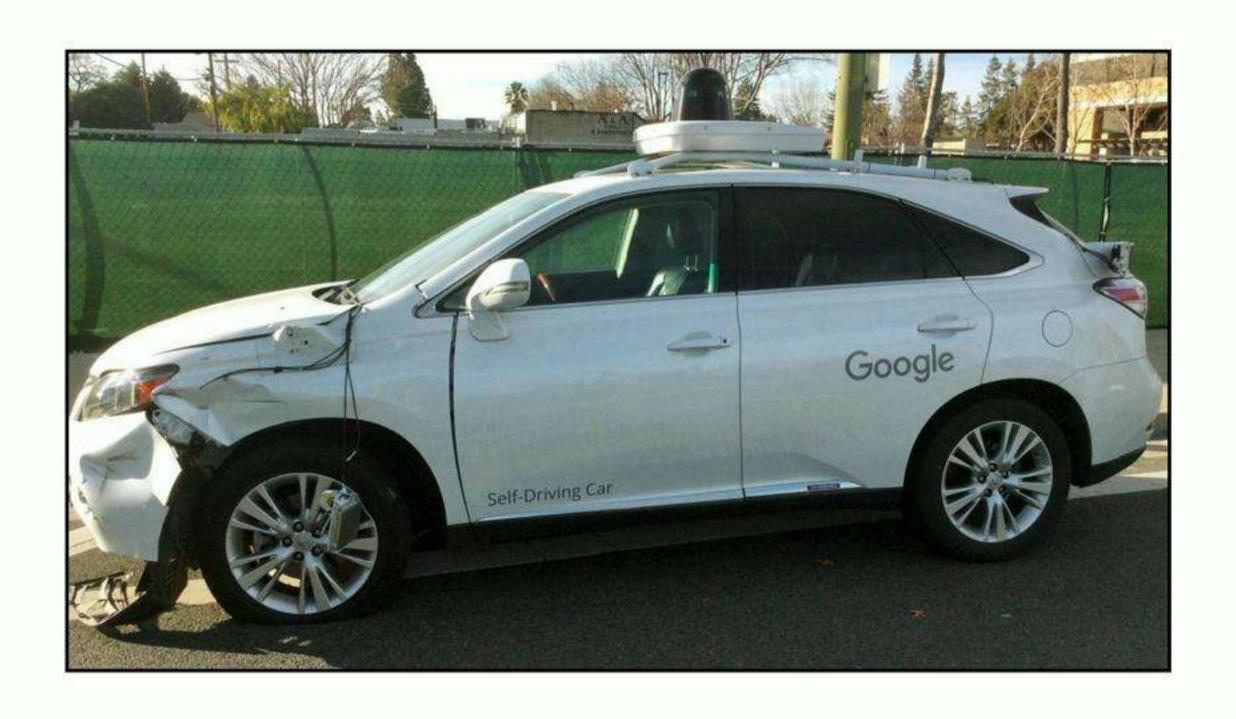
Source: Cisco

# incidents
reported by U.S. federal agencies



Source: US Department of Homeland Security

## Al in Safety-Critical Applications



## Al in Safety-Critical Applications



## Our Goal

## Study ML vulnerabilities and develop secure Al for high-stakes problems

## Secure AI

### Attack & Defense of Deep Neural Networks

ShapeShifter - Physical Adversarial Attack

SHIELD - Real-time Defense for *Images* 

### Do-it-yourself Adversarial ML

ADAGIO - Experimentation with Real-time Defense for Audio

MLsploit - Interactive Experimentation with Adversarial ML

## ShapeShifter

ECML-PKDD 2018

# First Targeted *Physical*Adversarial Attack for Object Detection



Shang-Tse Chen Georgia Tech



Cory Cornelius Intel



Jason Martin Intel

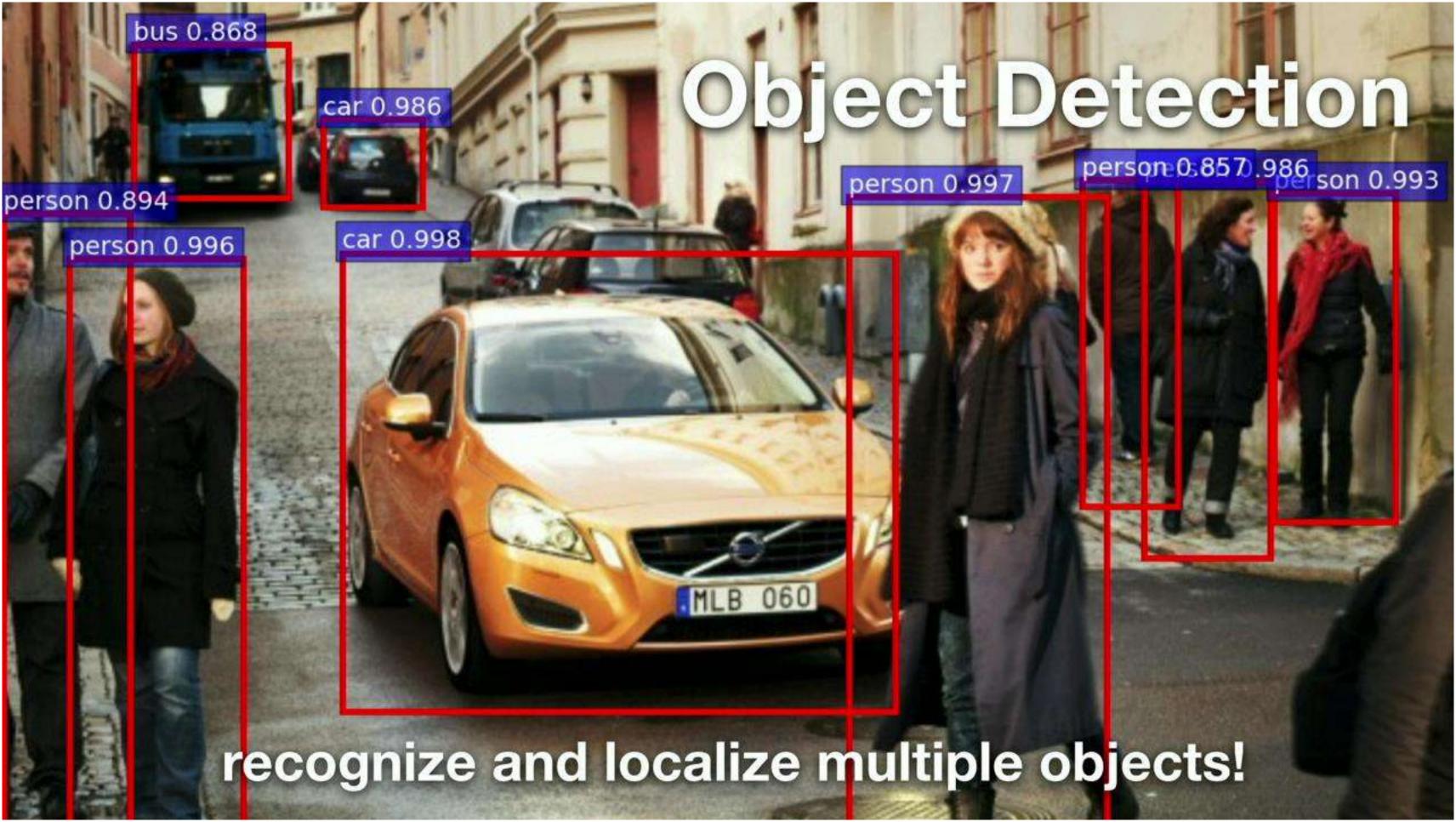


Polo Chau Georgia Tech





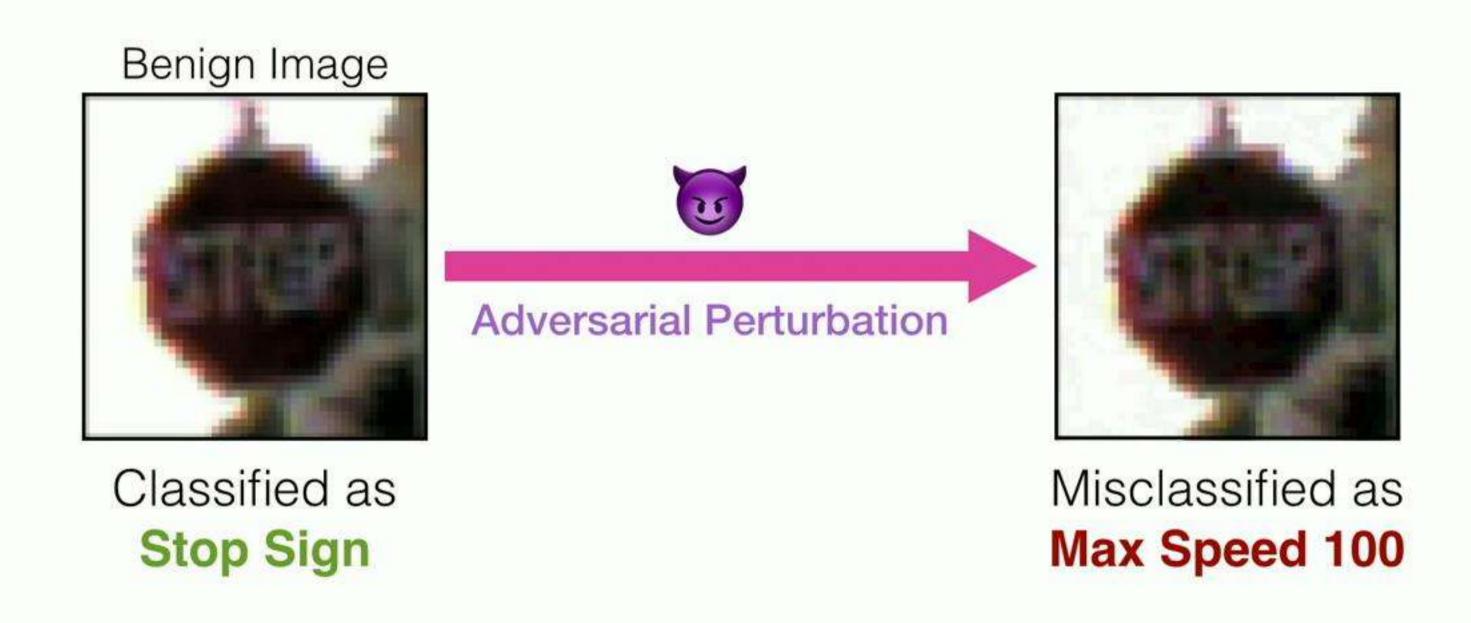








Classified as Stop Sign



But most attacks have impractical threat model



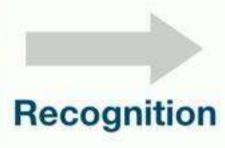


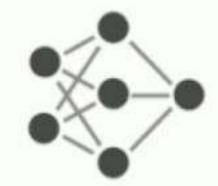














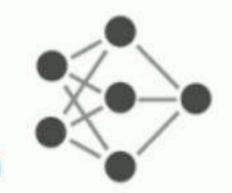












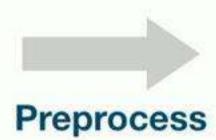






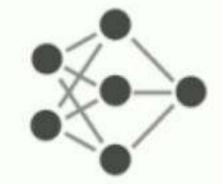
### Autonomous car system











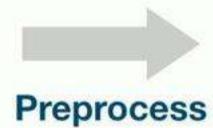






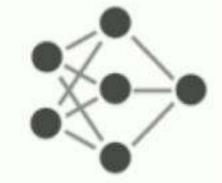
### Autonomous car system

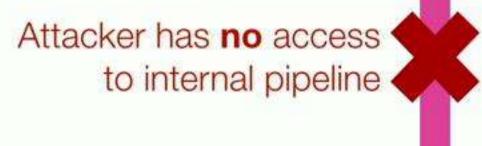






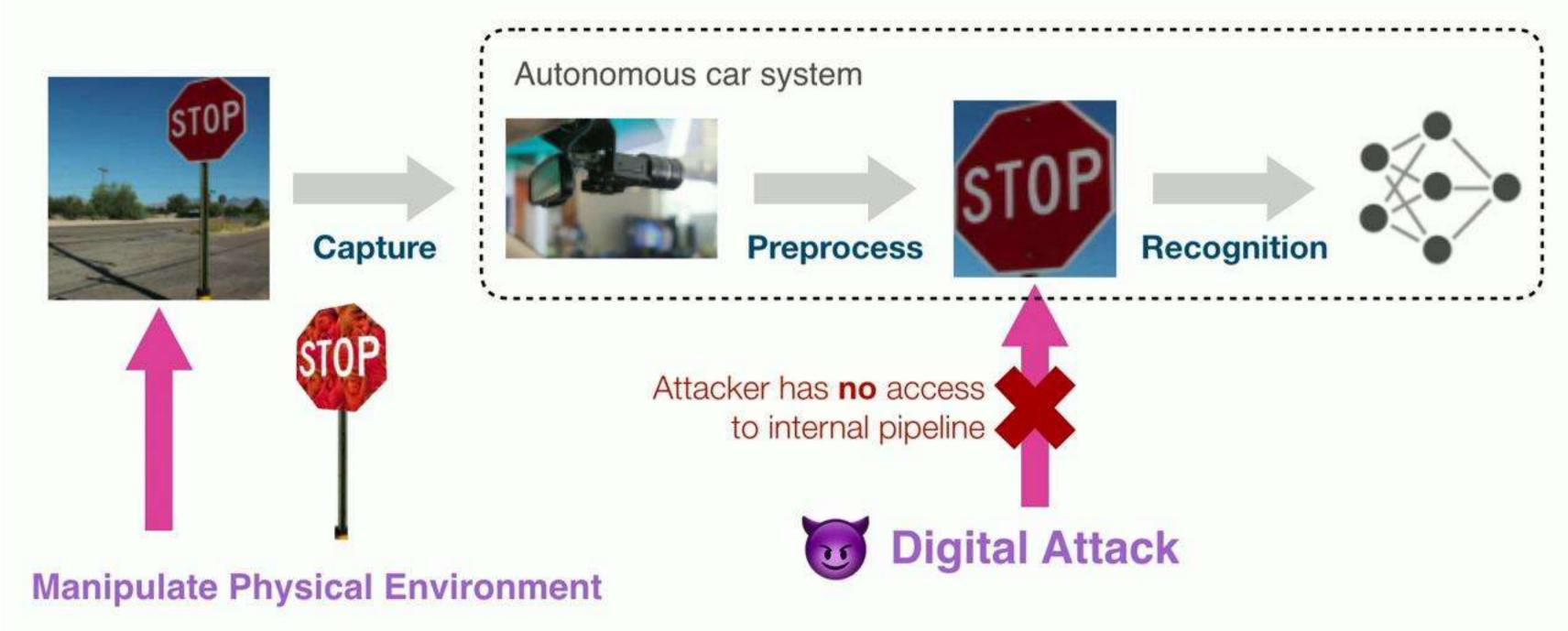
Recognition



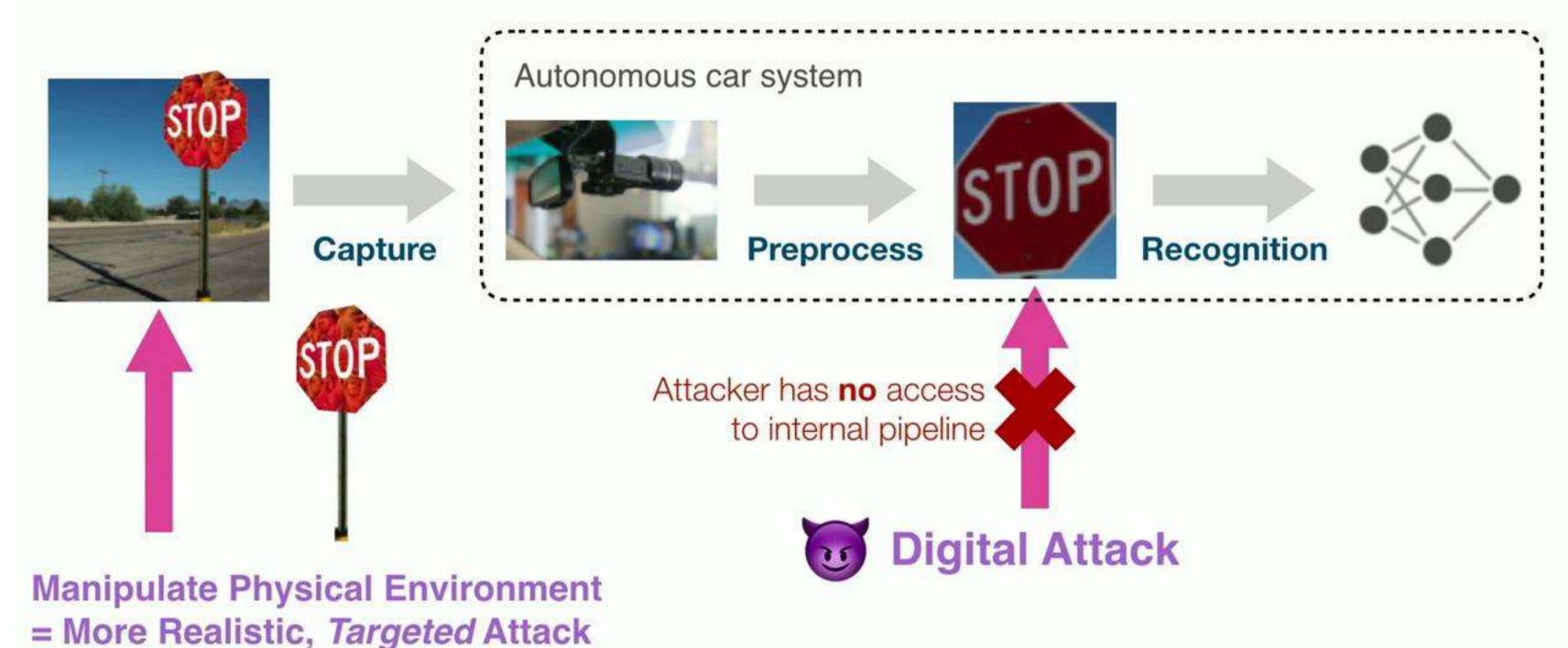


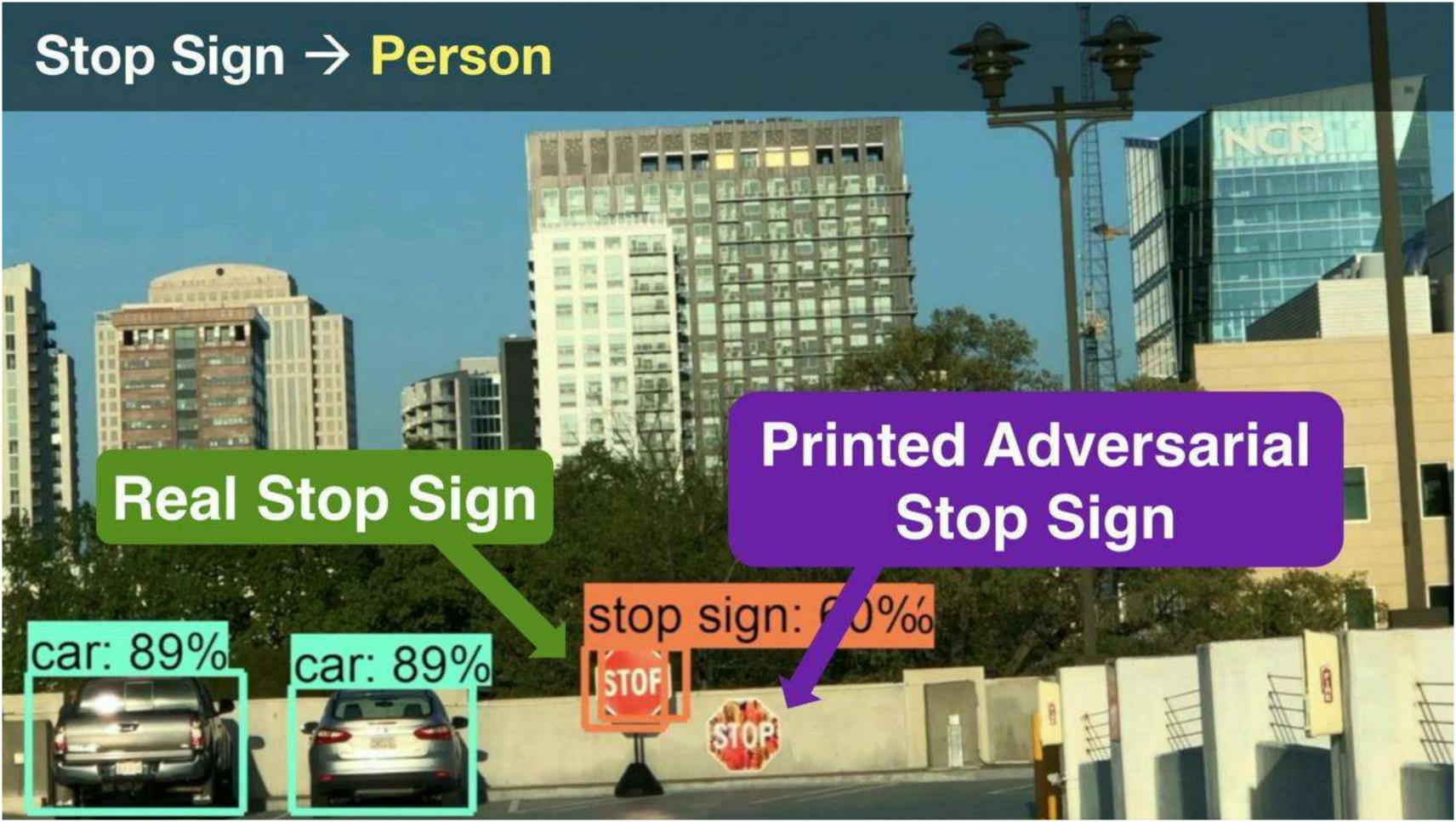


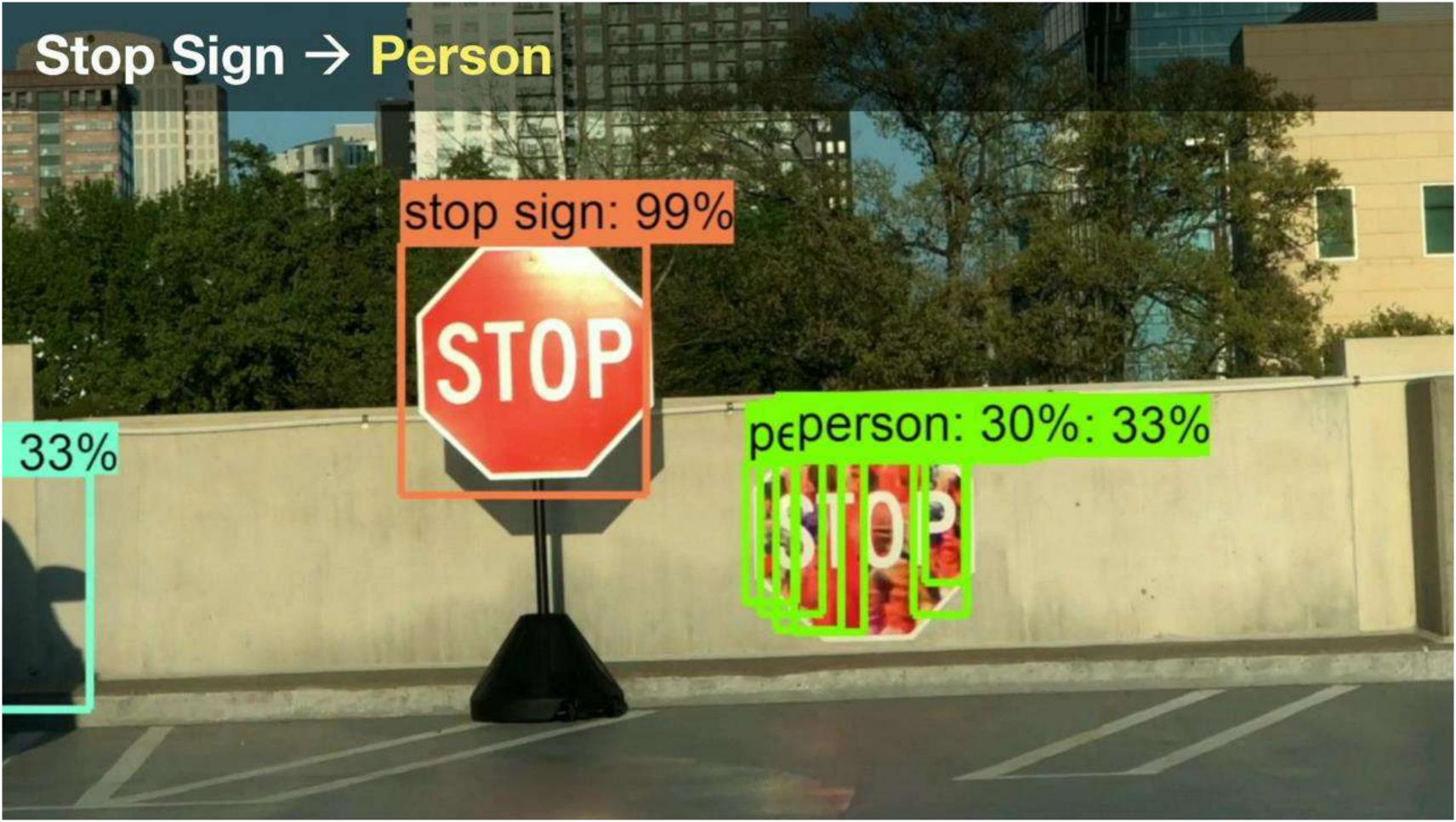
### Physically Realizable Adversarial Attack



### Physically Realizable Adversarial Attack







### **Prior Work on Physical Attacks**



Glasses that fool a face classifier [Sharif et al. CCS'16]



3D objects that fool an image classifier [Athalye et al. ICML'18]



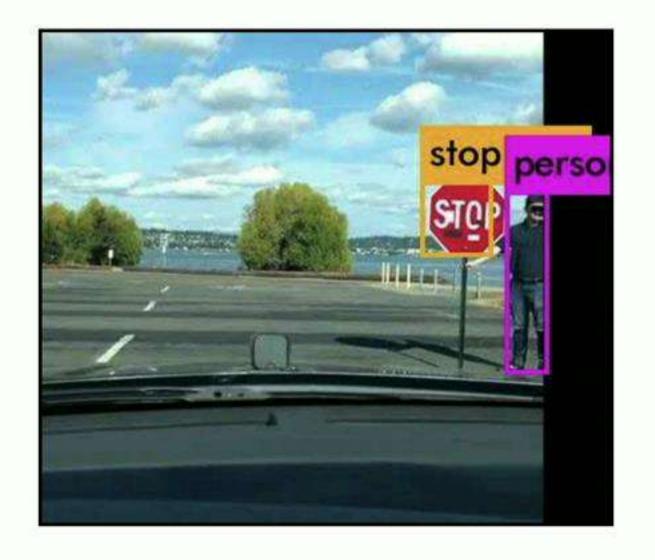
Stickers that fool a traffic sign classifier [Evtimov et al. CVPR'18]

They all focus on attacking image classifiers

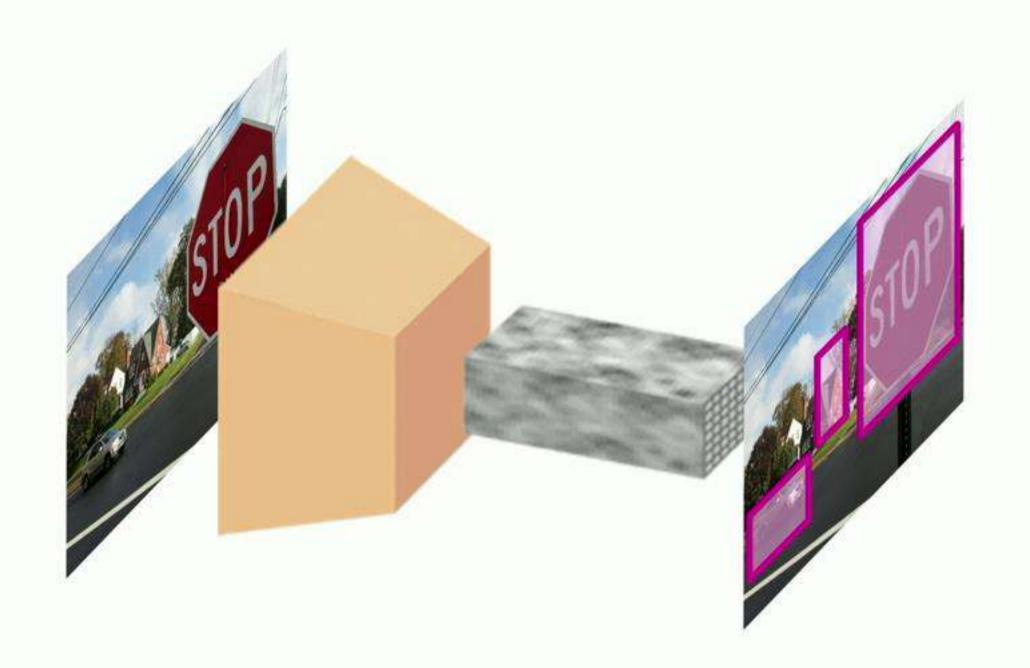
### Attack Object Detectors: Naïve Approach

Lu et al. [1] show the current technique cannot fool state-of-the-art object detectors like Faster R-CNN and YOLO



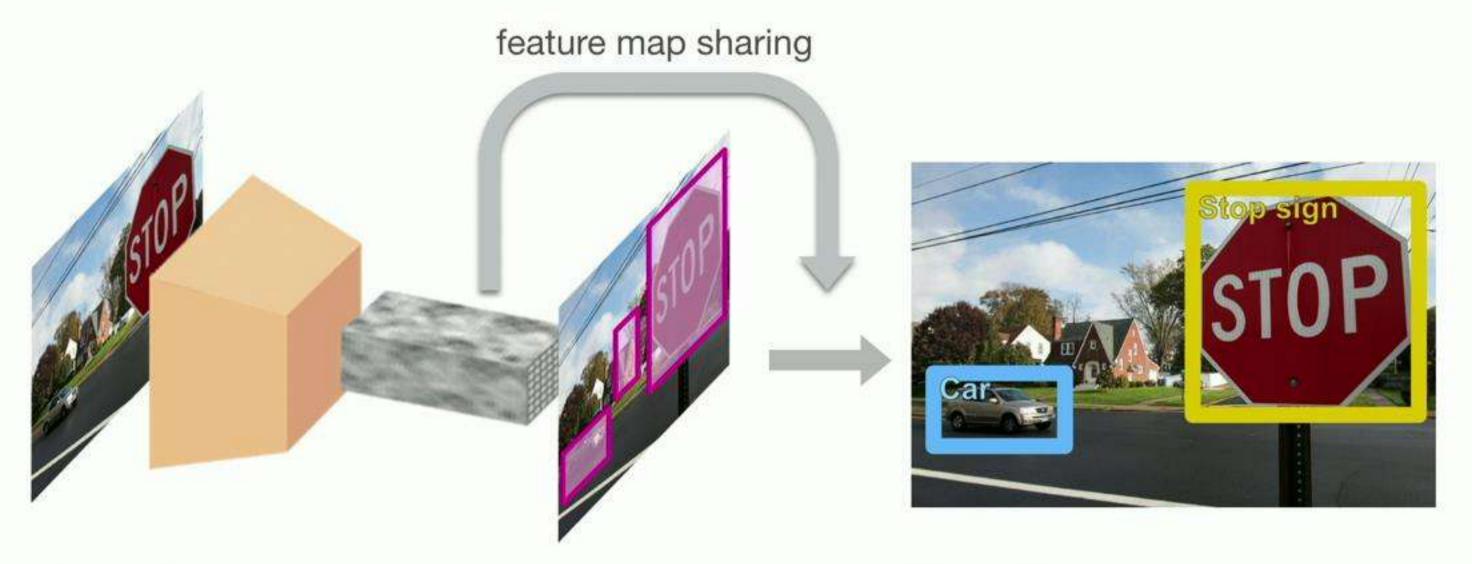


### **Brief Overview of Faster R-CNN**



A state-of-the-art Object Detector Model

### **Brief Overview of Faster R-CNN**

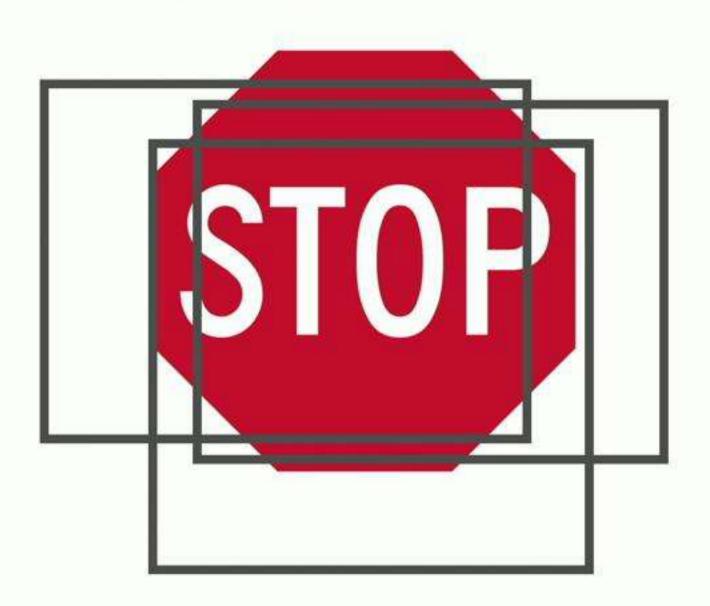


**Stage 1**: Generate region proposals

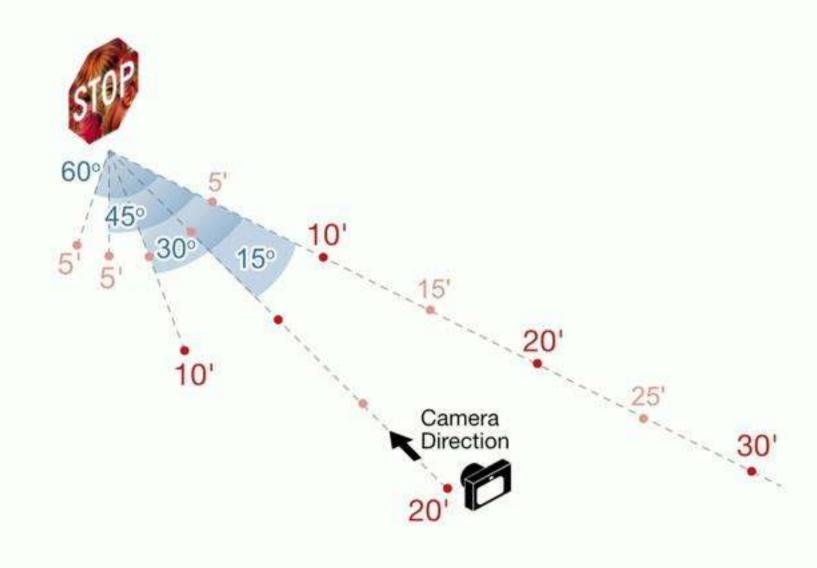
Stage 2: Refined localization and classification

### Challenges of Physically Attacking Faster R-CNN

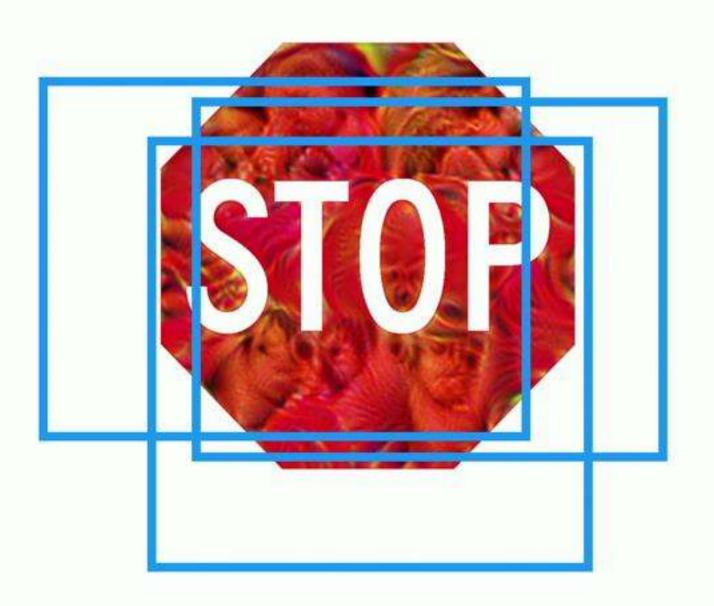
1. Multiple region proposals



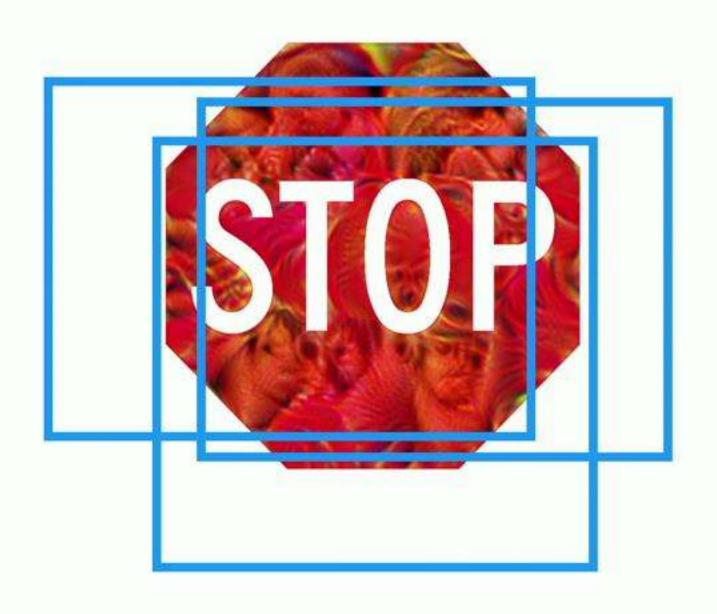
2. Distances, angles, lightings



Minimize: sum of classification losses



Minimize: sum of classification losses + deviation loss

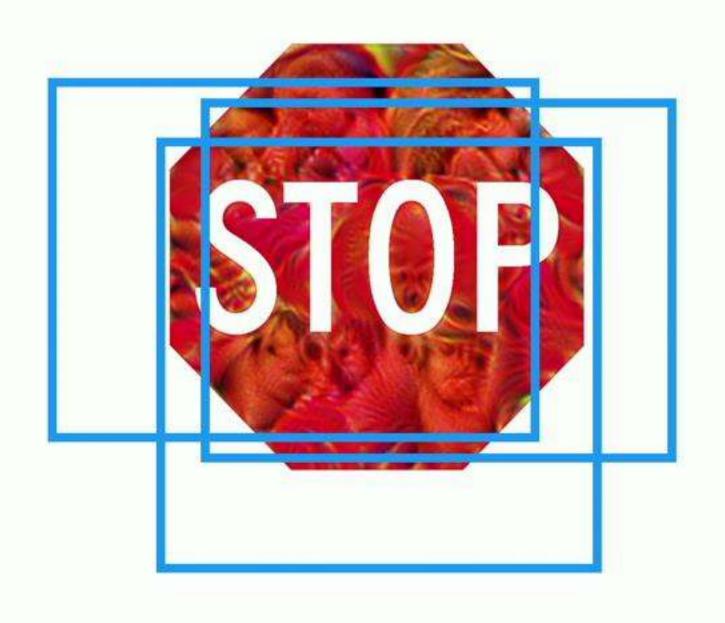


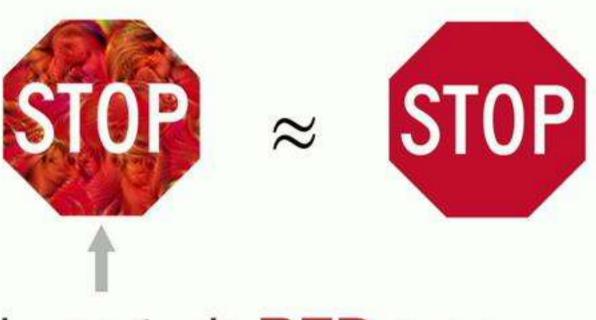






Minimize: sum of classification losses + deviation loss





Only perturb RED area

Human eye is less sensitive
to changes in darker color

### Solution 2: Robust to Real-World Distortions

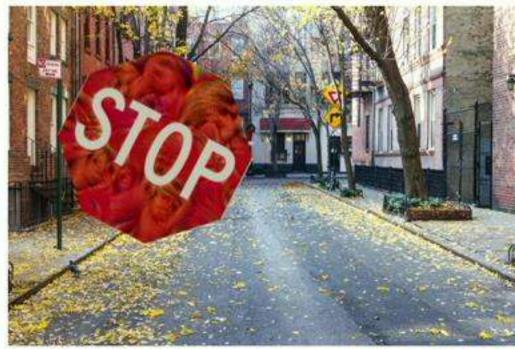
### Solution 2: Robust to Real-World Distortions

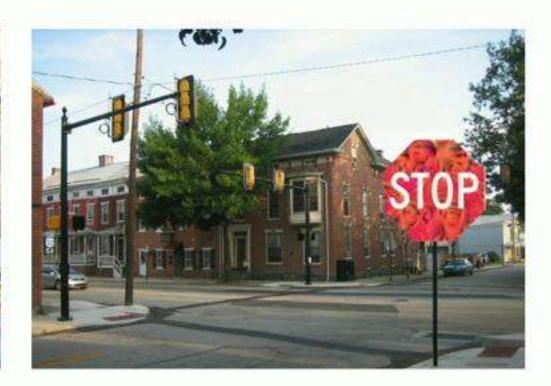
Adapt Expectation over Transformation [Athalye et al, ICML'18]

### Solution 2: Robust to Real-World Distortions

#### Adapt Expectation over Transformation [Athalye et al, ICML'18]

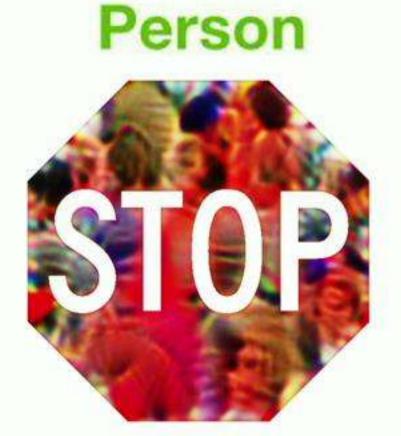


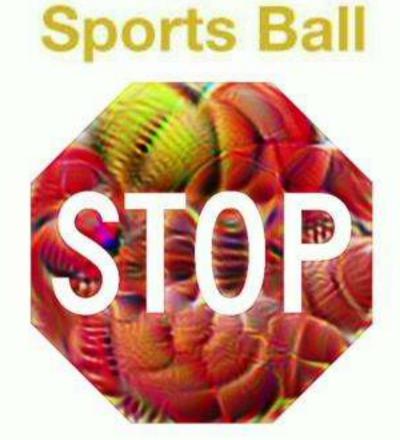




Optimize over different backgrounds, scales, rotations, lightings

### Perturbed Stop Signs



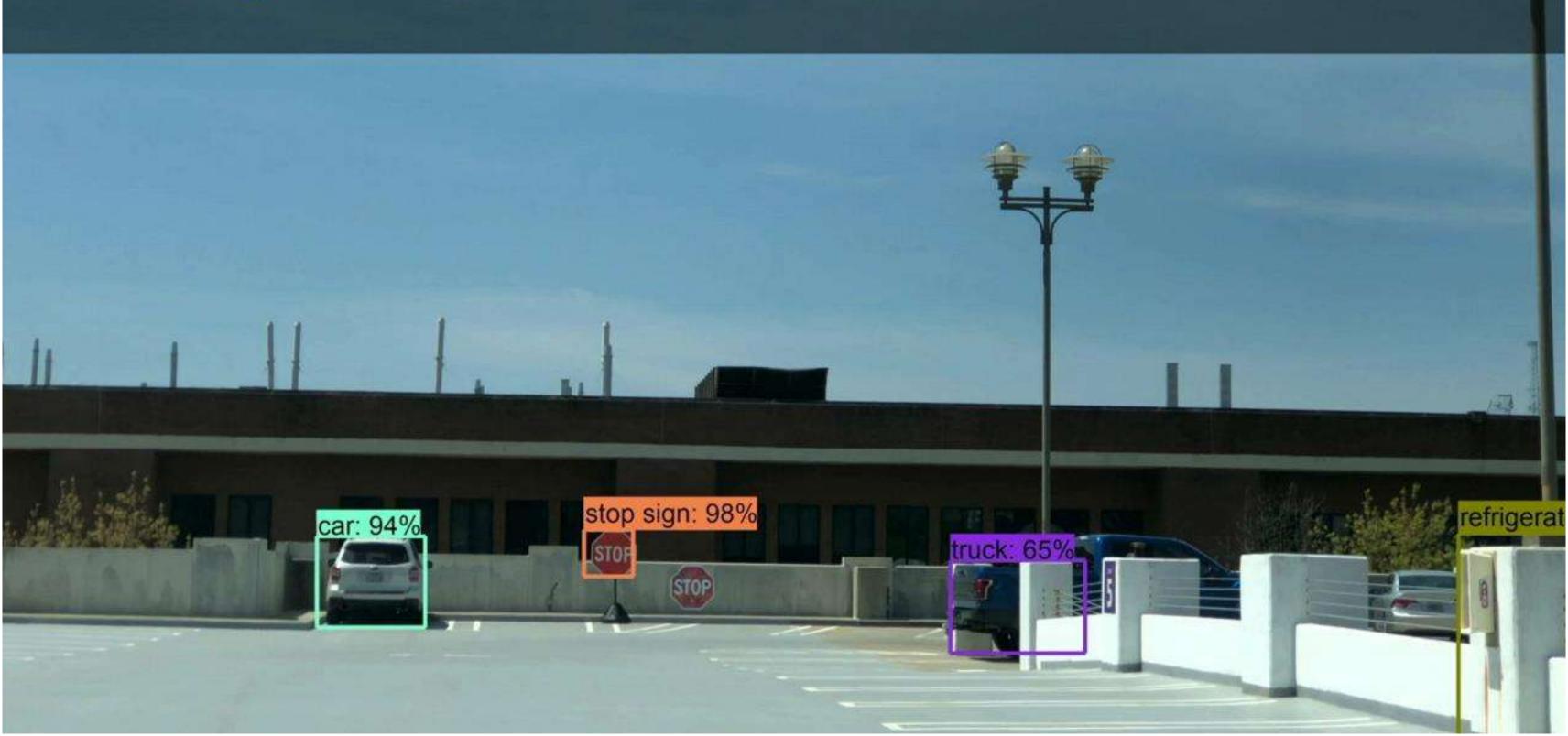






perturbations more conspicuous in order to survive viewing distances, angles, lighting conditions, camera limitations

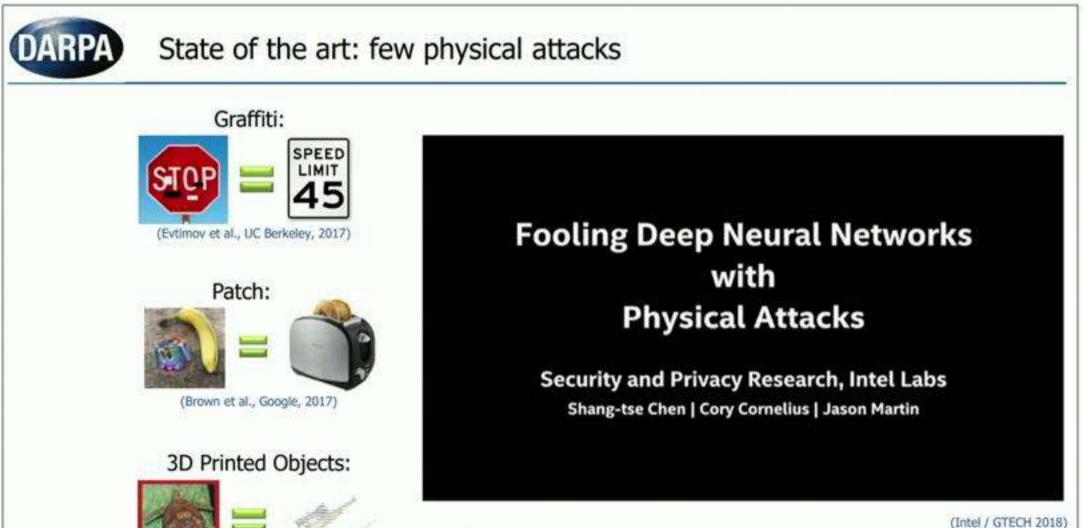
# Untargeted Attack



# Untargeted Attack



# ShapeShifter Motivates New DARPA Program GARD: Defense for Al



Highlights ShapeShifter as the state-of-the-art physical attack

(Intel / GTECH 2018)

- All physical attacks to date are White Box
- · No current consideration of resource constraints

https://www.darpa.mil/attachments/GARD\_ProposersDay.pdf

(Athalye et al., MIT, 2017)

# SHIELD Fast, Practical Defense for Image Classification

X KDD'18 Audience Appreciation Award (runner-up) KDD'19 LEMINCS

[Open-sourced]



Nilaksh Das



Madhuri Shangbogue



Shang-Tse Chen



Fred Hohman





Cory Cornelius



**Li** Chen



Michael Kounavis



Polo Chau





### Adversarial Machine Learning Landscape



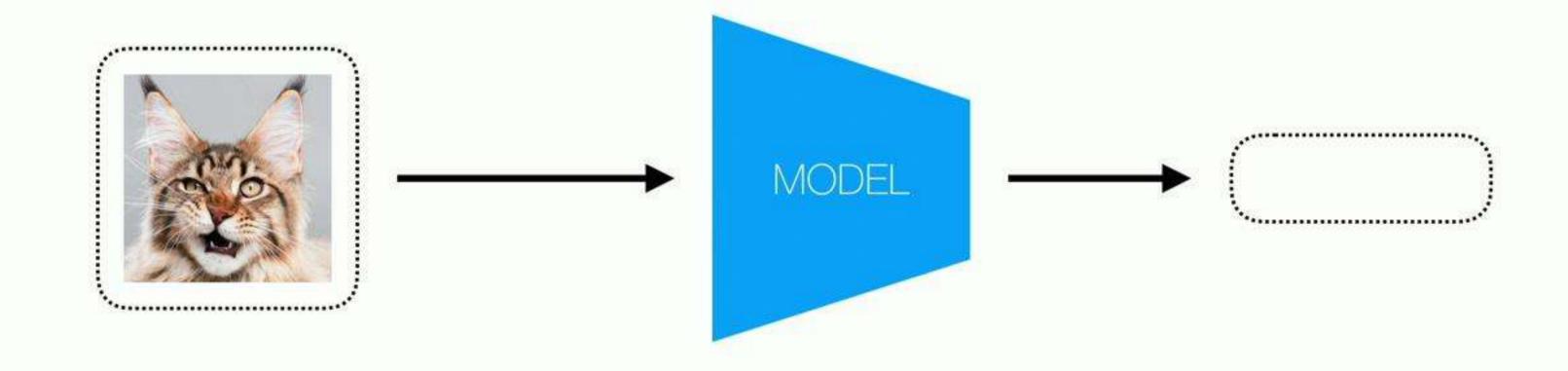


### Adversarial Machine Learning Landscape

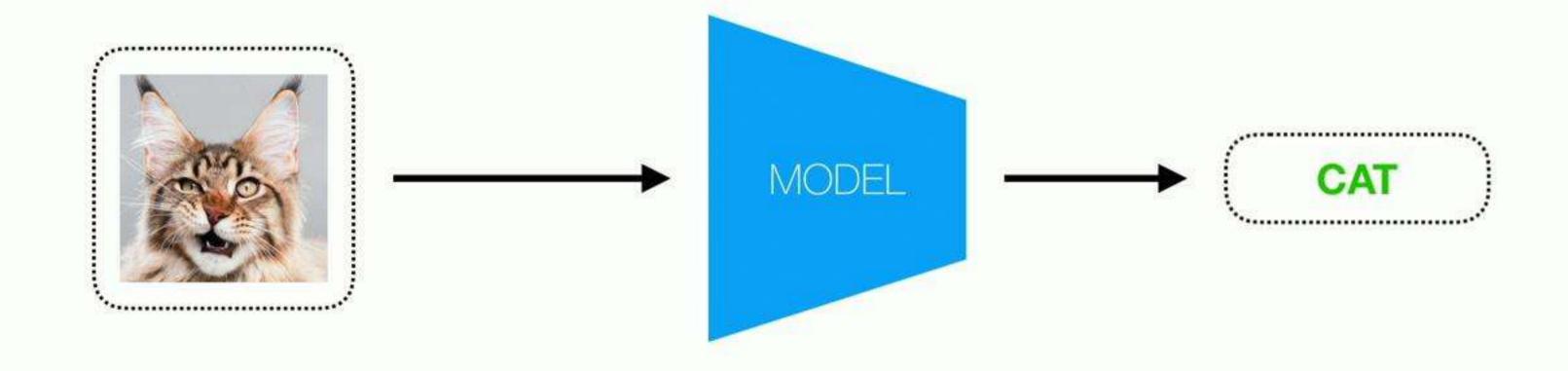




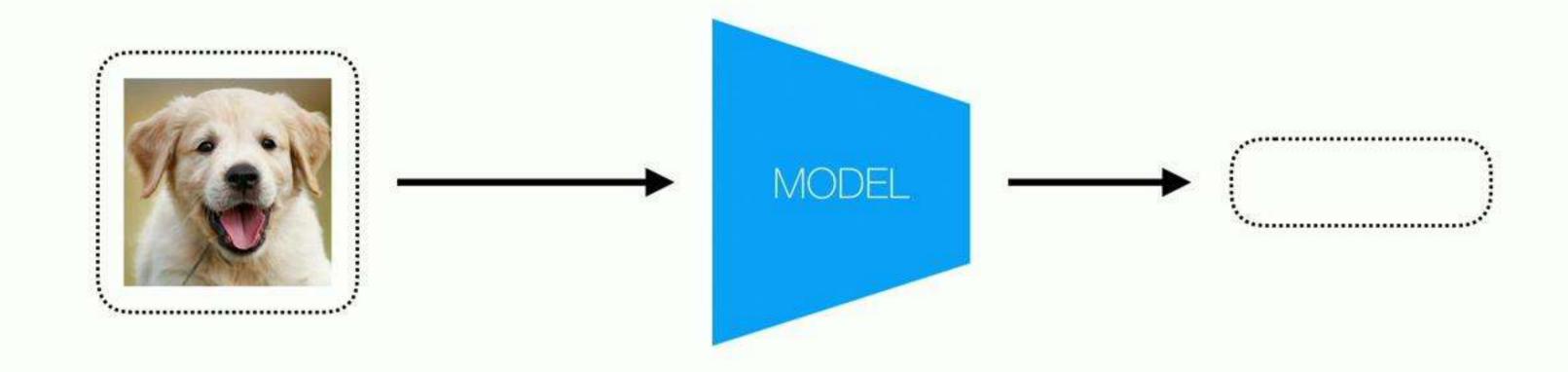
# Deep Learning for Image Classification

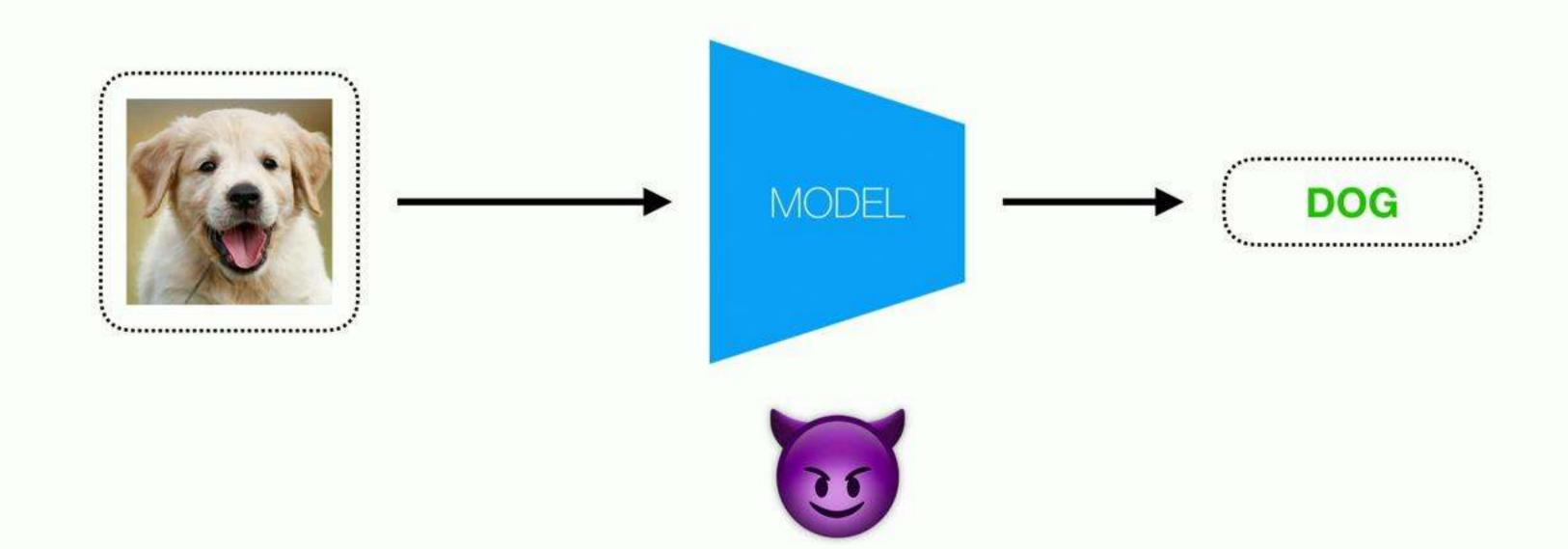


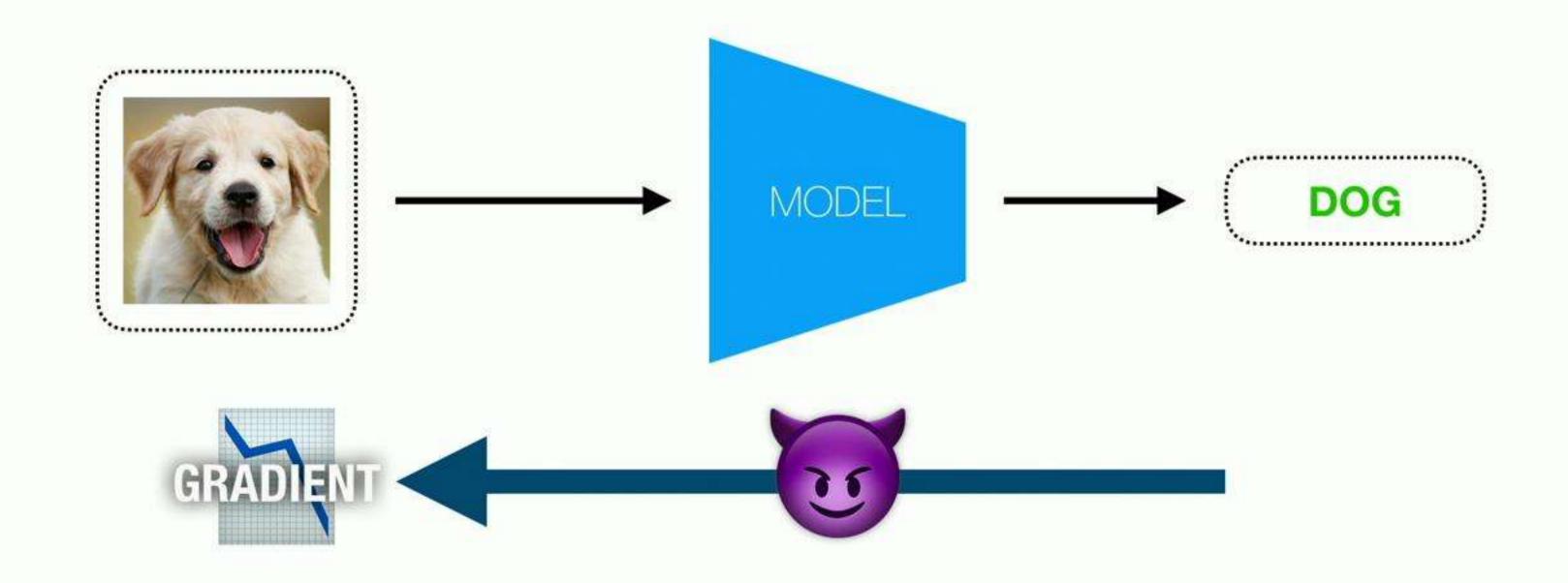
# Deep Learning for Image Classification

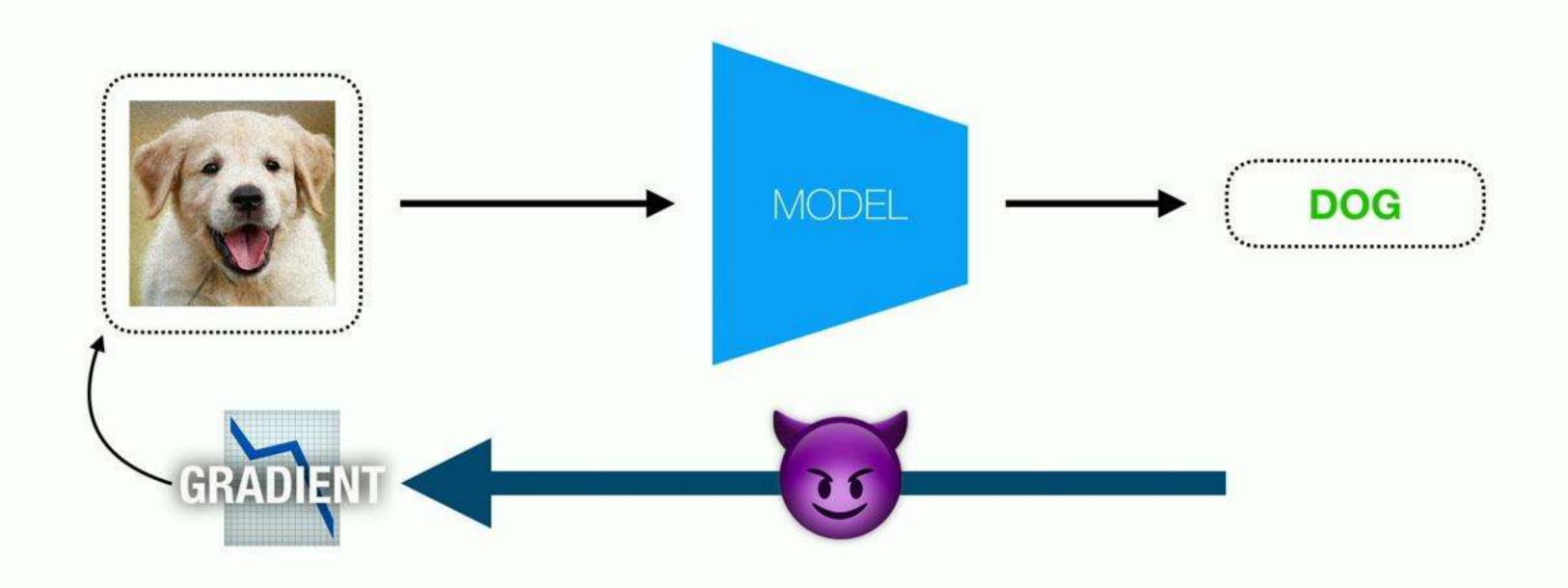


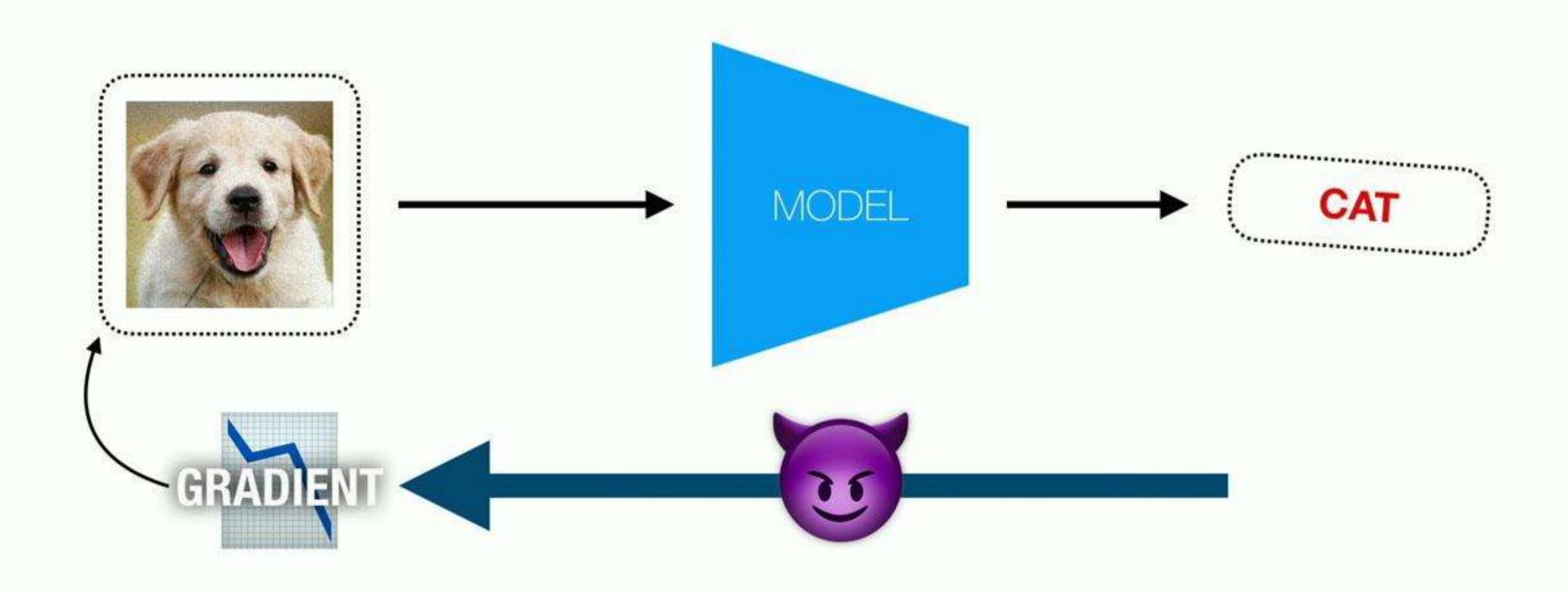
# Deep Learning for Image Classification



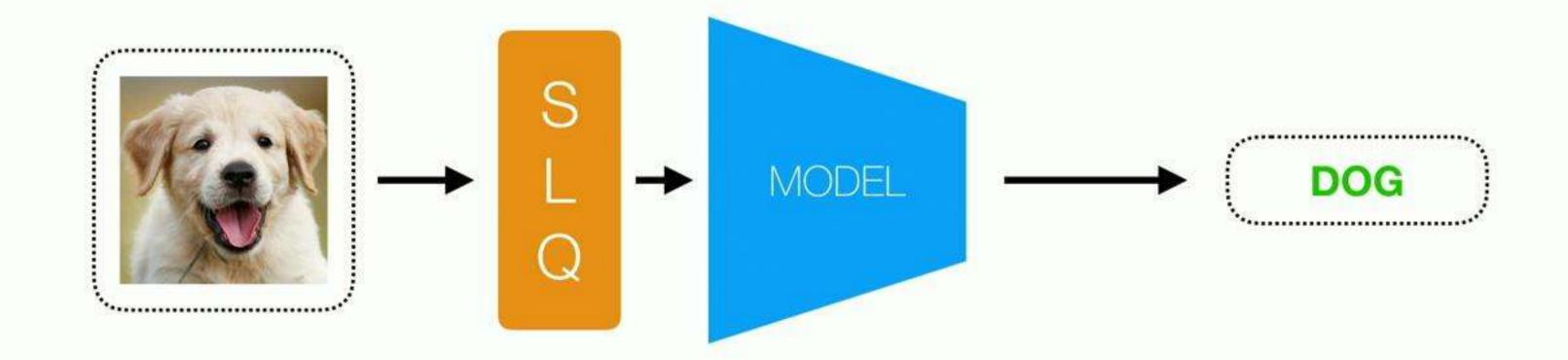




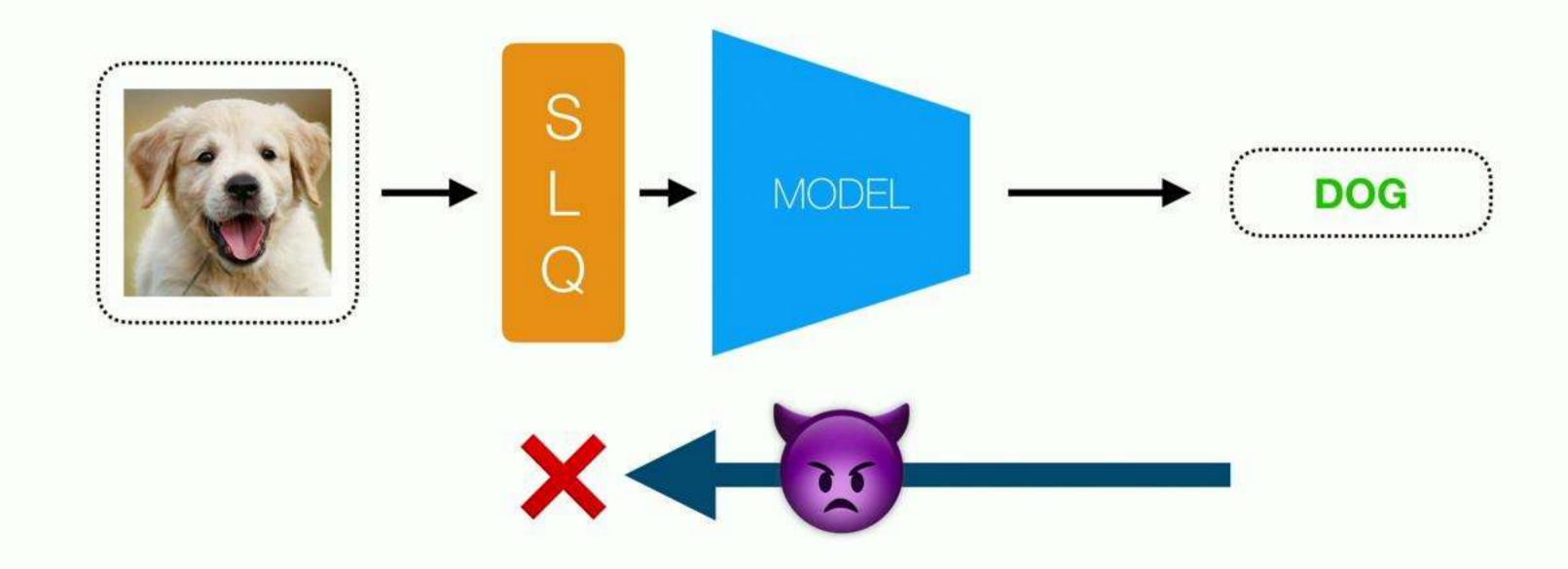




# Stochastic Local Quantization (SLQ)



### Stochastic Local Quantization (SLQ)



# **SLQ** leverages JPEG compression

### **SLQ** leverages JPEG compression

JPEQ Quality 80



JPEQ Quality 60



JPEQ Quality 40



JPEQ Quality 20





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

<sup>\*</sup> larger blocks shown for presentation

### **SLQ** leverages JPEG compression

JPEQ Quality 80



JPEQ Quality 60



JPEQ Quality 40



JPEQ Quality 20





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

<sup>\*</sup> larger blocks shown for presentation



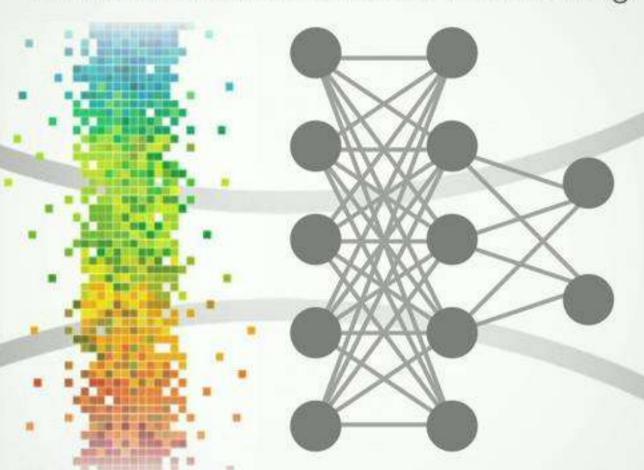


(Attacked)

Labrador Retriever



Secure Heterogeneous Image Ensemble with Localized Denoising



Real-time Compression Preprocessing

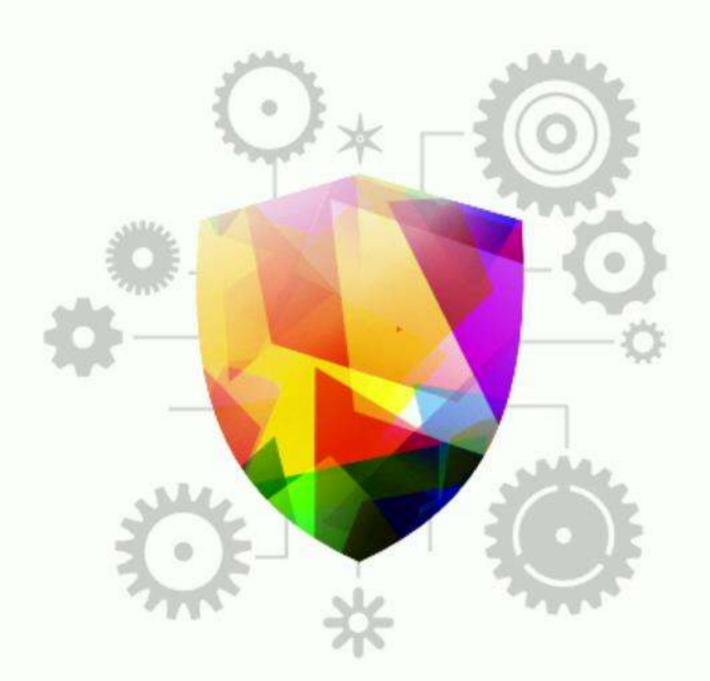
Vaccinated
Deep Neural
Network Ensemble



Correctly



Correctly
Classified

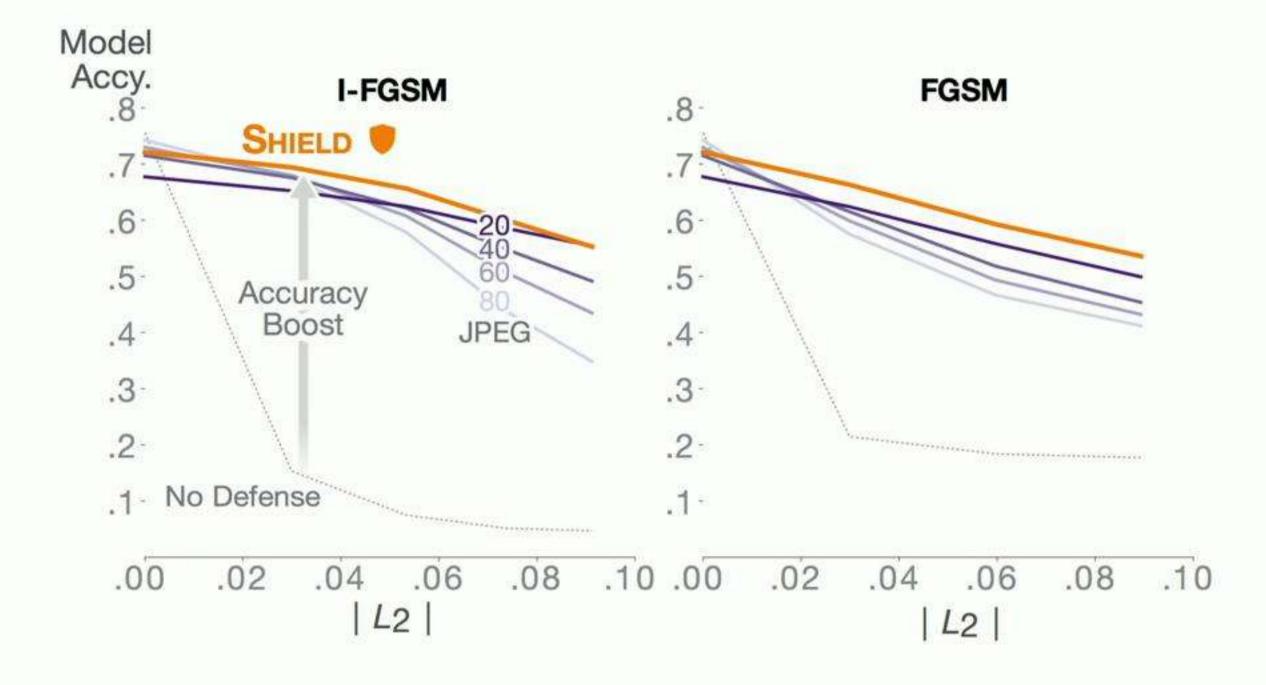


# SHIELD is multi-pronged approach that incorporates

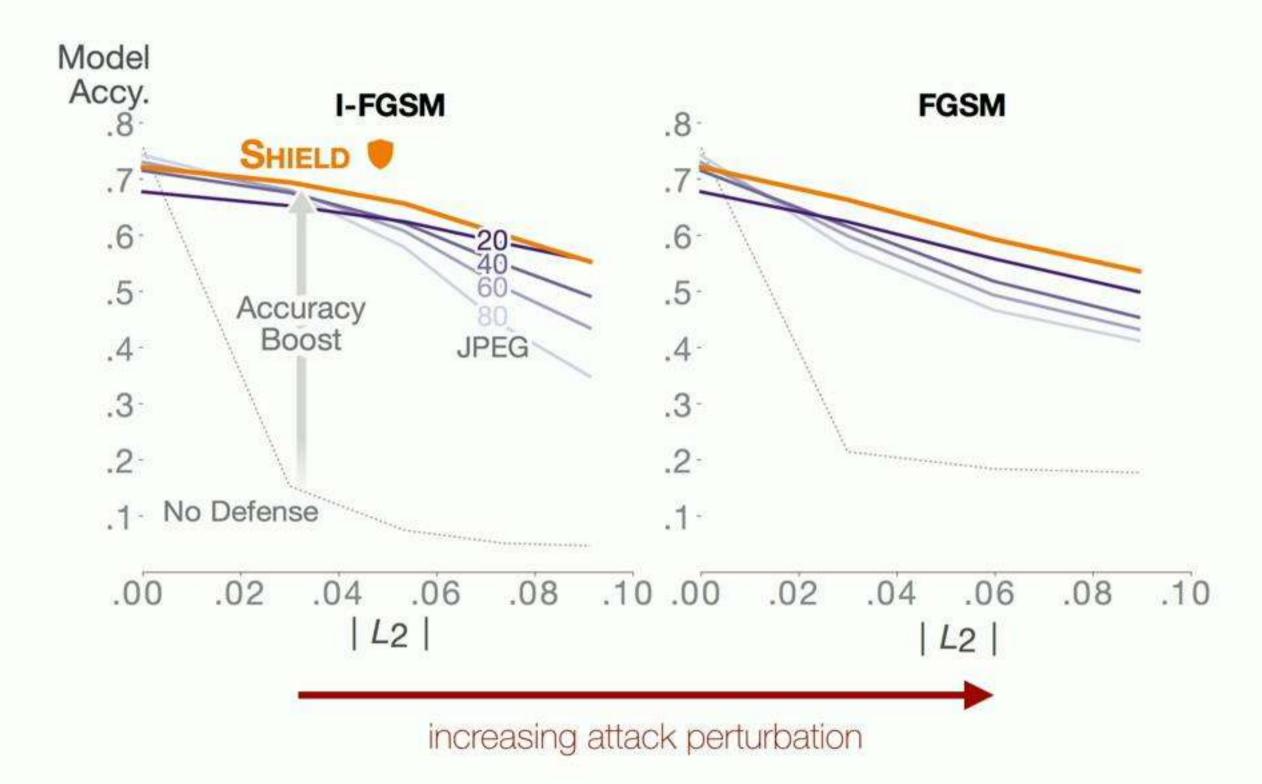
- Stochastic Local Quantization
- Model Vaccination (re-training)
- Ensembling

to mitigate adversarial attacks

#### Results with ResNet-50 v2 (on ImageNet validation set)



#### Results with ResNet-50 v2 (on ImageNet validation set)

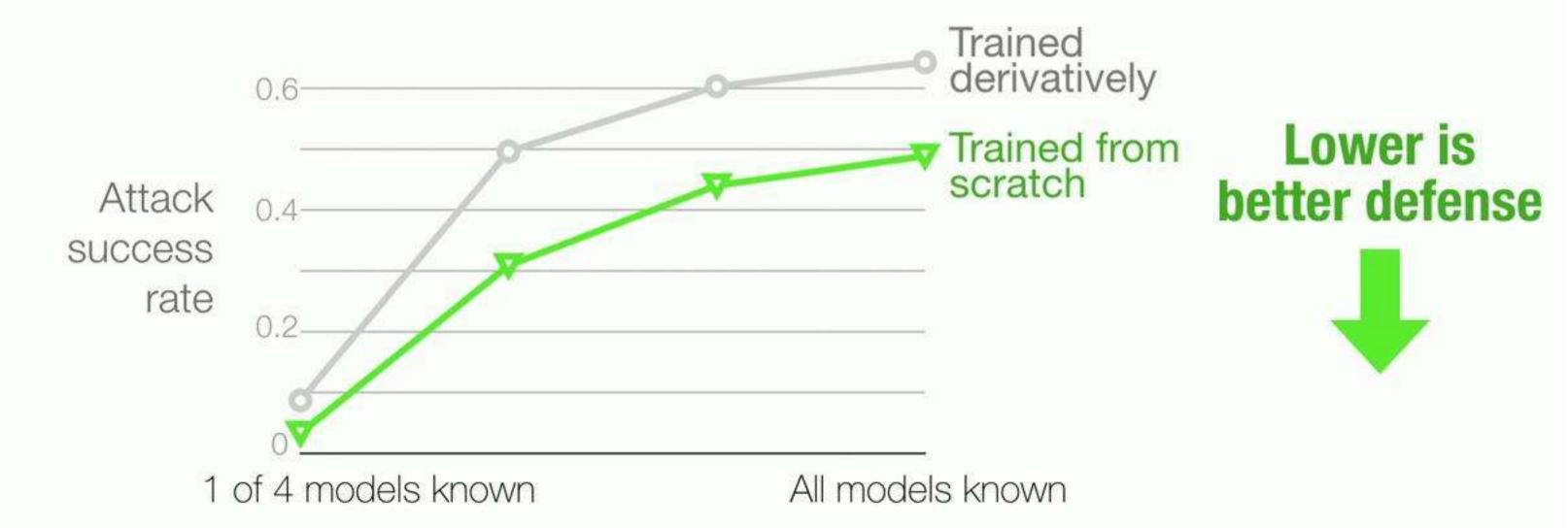


#### Defense Runtime Comparison (in seconds; shorter is better)



#### **Extending to Adaptive Attacks**

SHIELD ensemble with **less correlated** model weights are **more robust** to *targeted* adaptive attacks [KDD'19 LEMINCS]



The Efficacy of SHIELD under Different Threat Models. Cory Cornelius, Nilaksh Das, Shang-Tse Chen, Li Chen, Michael E. Kounavis, Duen Horng (Polo) Chau. KDD 2019 Workshop on Learning and Mining for Cybersecurity (LEMINCS).

## ADAGIO ECML-PKDD 2018

# Experimentation with Real-time Defense for Speech to Text



Nilaksh Das

GT



Madhuri Shangbogue GT



Shang-Tse Chen



Li Chen Intel



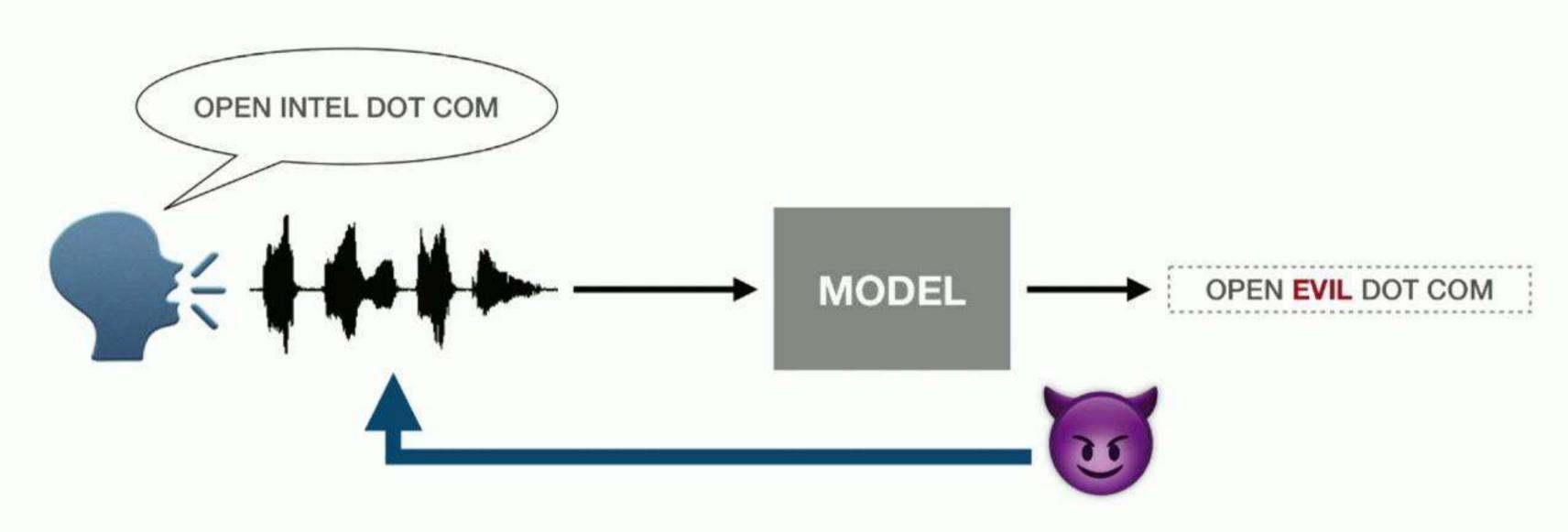
Michael Kounavis Intel



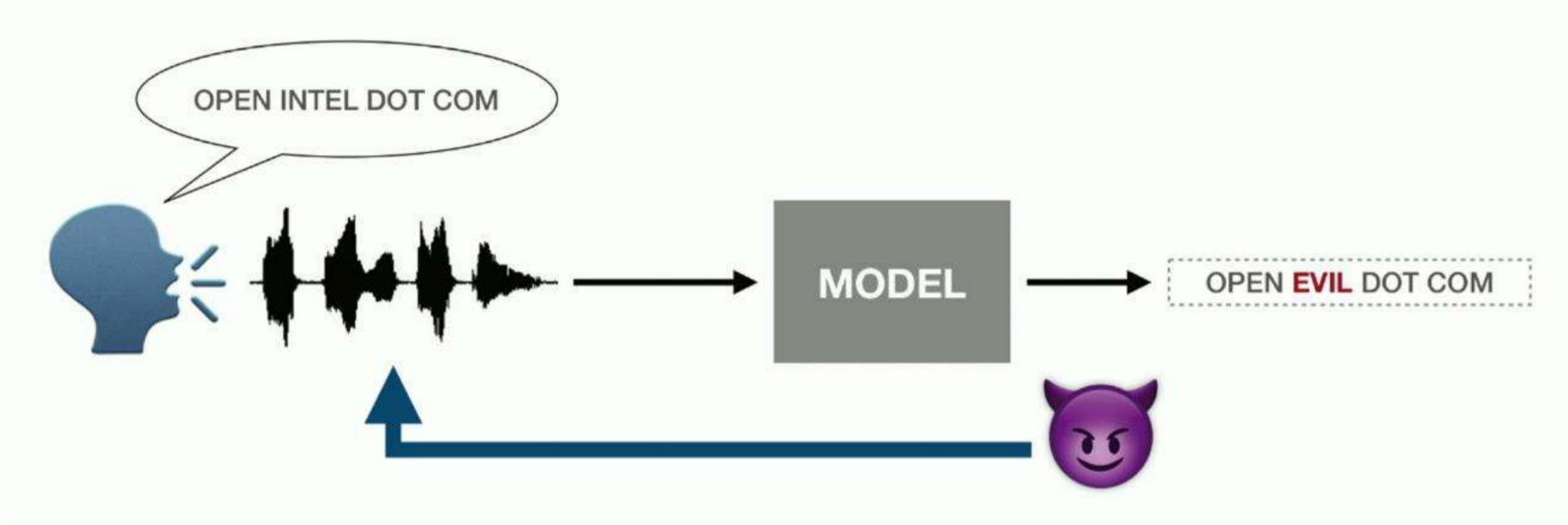
Polo Chau



#### Adversarial Attack on Speech-to-text

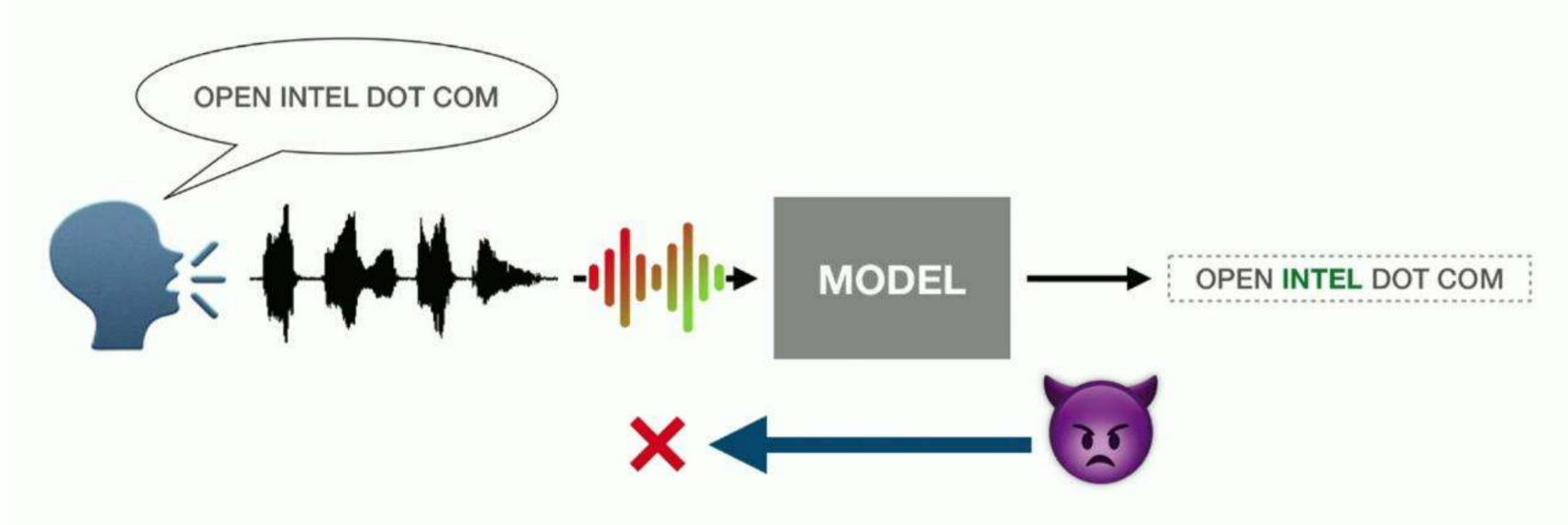


#### Adversarial Attack on Speech-to-text



An adversary uses backpropagation to attack the model.

#### Adversarial Attack on Speech-to-Text

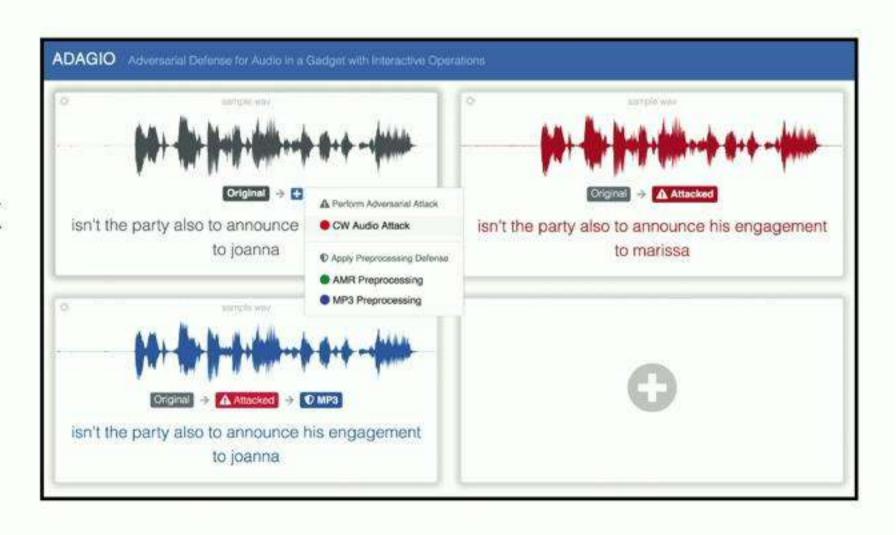


**ADAGIO** incorporates compression as defense, which blocks the gradient to the attacker.

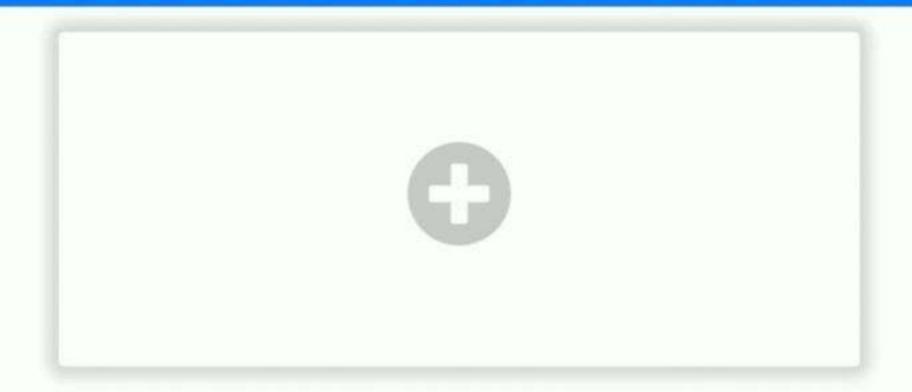


## Interactive Experimentation with Adversarial Attack & Defense for Audio

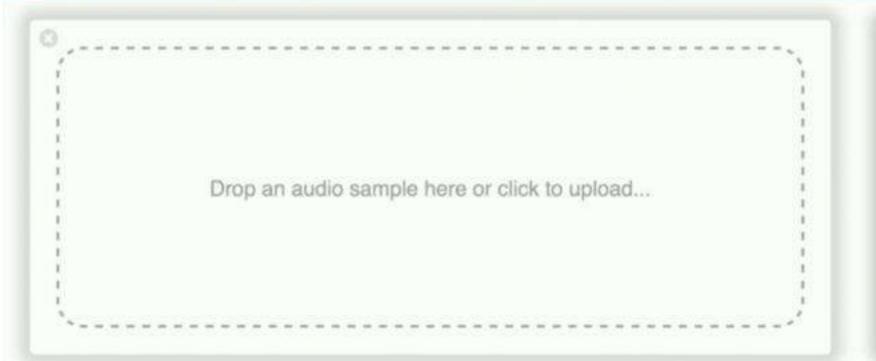
- Upload your own audio sample
- X Perform audio adversarial attack
- Apply compression to defend
- Play audio, listen for differences

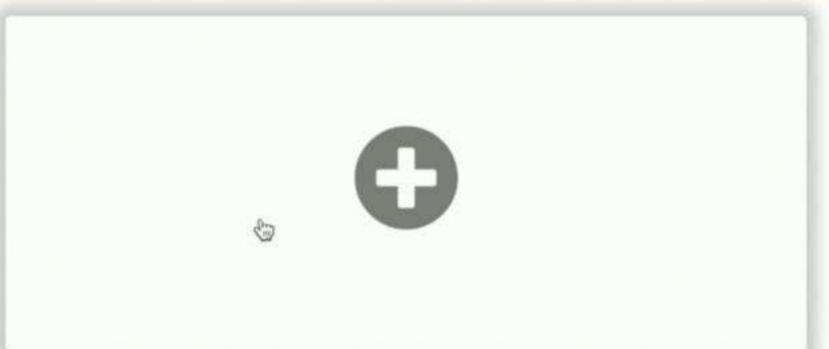


ADAGIO = Attack & Defense for Audio in a Gadget with Interactive Operations

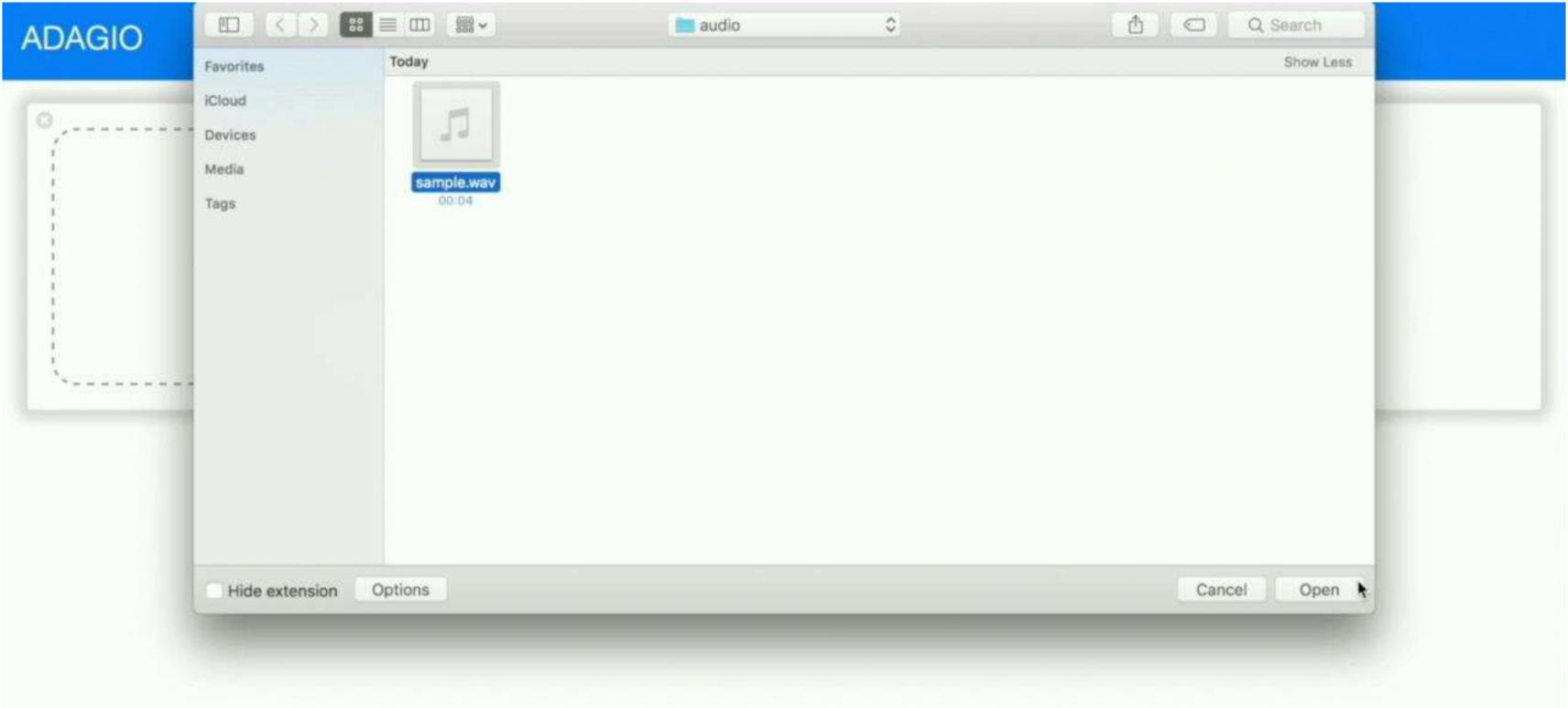


#### **ADAGIO**

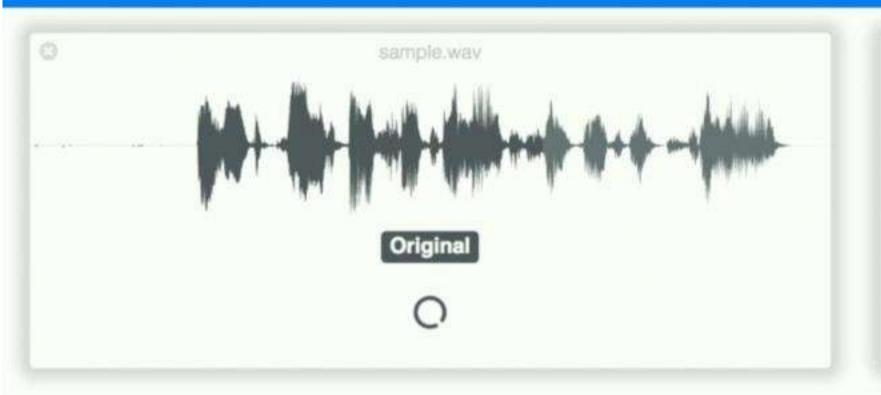




Upload an audio sample



#### **ADAGIO**

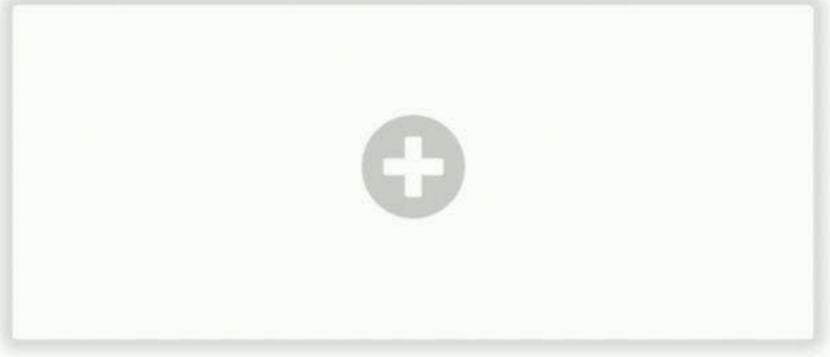




ADAGIO shows transcription from DeepSpeech

#### **ADAGIO**





ADAGIO shows transcription from DeepSpeech

## Secure AI

#### Attack & Defense of Deep Neural Networks

ShapeShifter - Physical Adversarial Attack

SHIELD - Real-time Defense for *Images* 

#### Do-it-yourself Adversarial ML

ADAGIO - Experimentation with Real-time Defense for Audio

MLsploit - Interactive Experimentation with Adversarial ML







github.com/mlsploit

## A Framework for Interactive Experimentation with Adversarial Machine Learning Research

Contributors from Intel Science and Technology Center for Adversary-Resilient Security Analytics:
Nilaksh Das, 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)

[BlackHat Asia '19, KDD'19 Showcase]

## 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

## MLsploit

- \* AVPass (leaking and bypassing Android malware detection systems)
- ★ **ELF** (bypassing Linux malware detection with API perturbation)
- **★ PE** (create and attack ML models for detecting Windows PE malware)
- Intel®-Software Guard Extensions
  - (privacy preserving adversarial ML as a service)
- \* SHIELD (attack and defend state-of-the-art image classification models)
  - \* Attacks: FGSM, DeepFool, Carlini-Wagner
  - \* Defenses: SLQ, JPEG, Median Filter, TV-Bregman



2 APPLY RESEARCH FUNCTIONS

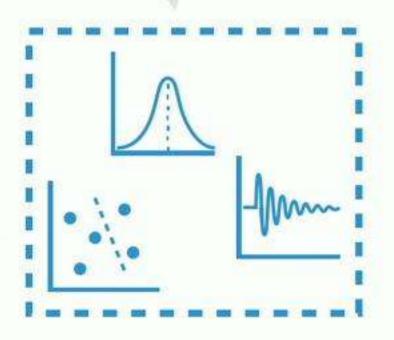


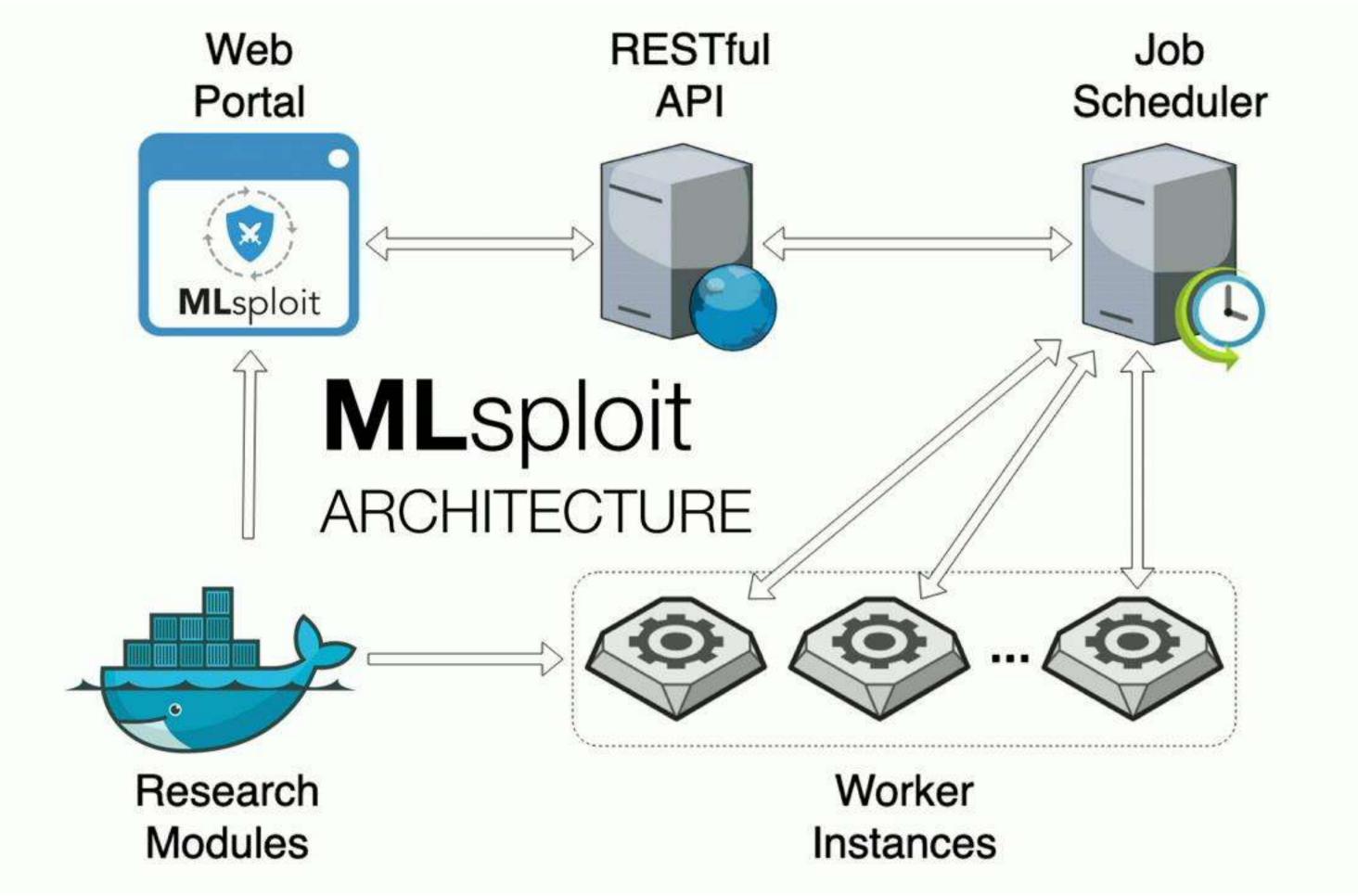
**UPLOAD** 

SAMPLES

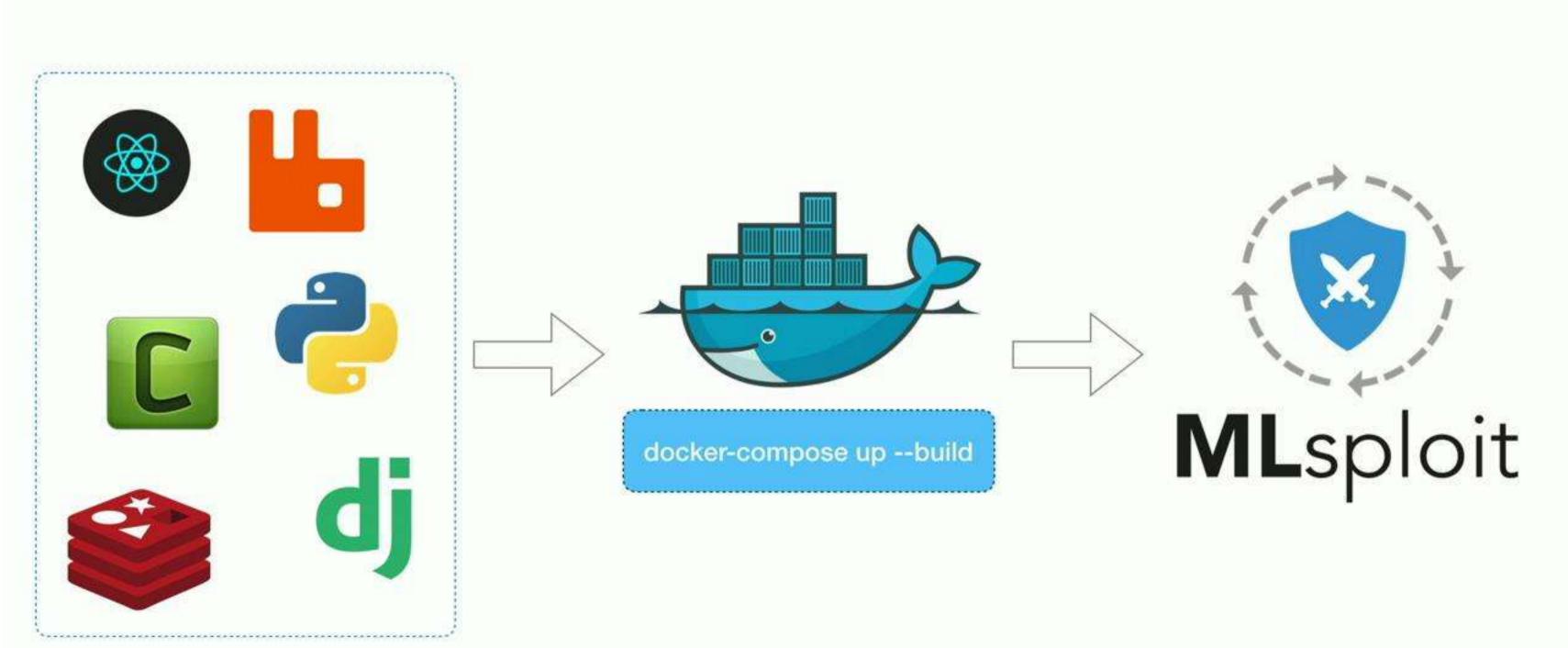
**COMPARE RESULTS** 

3

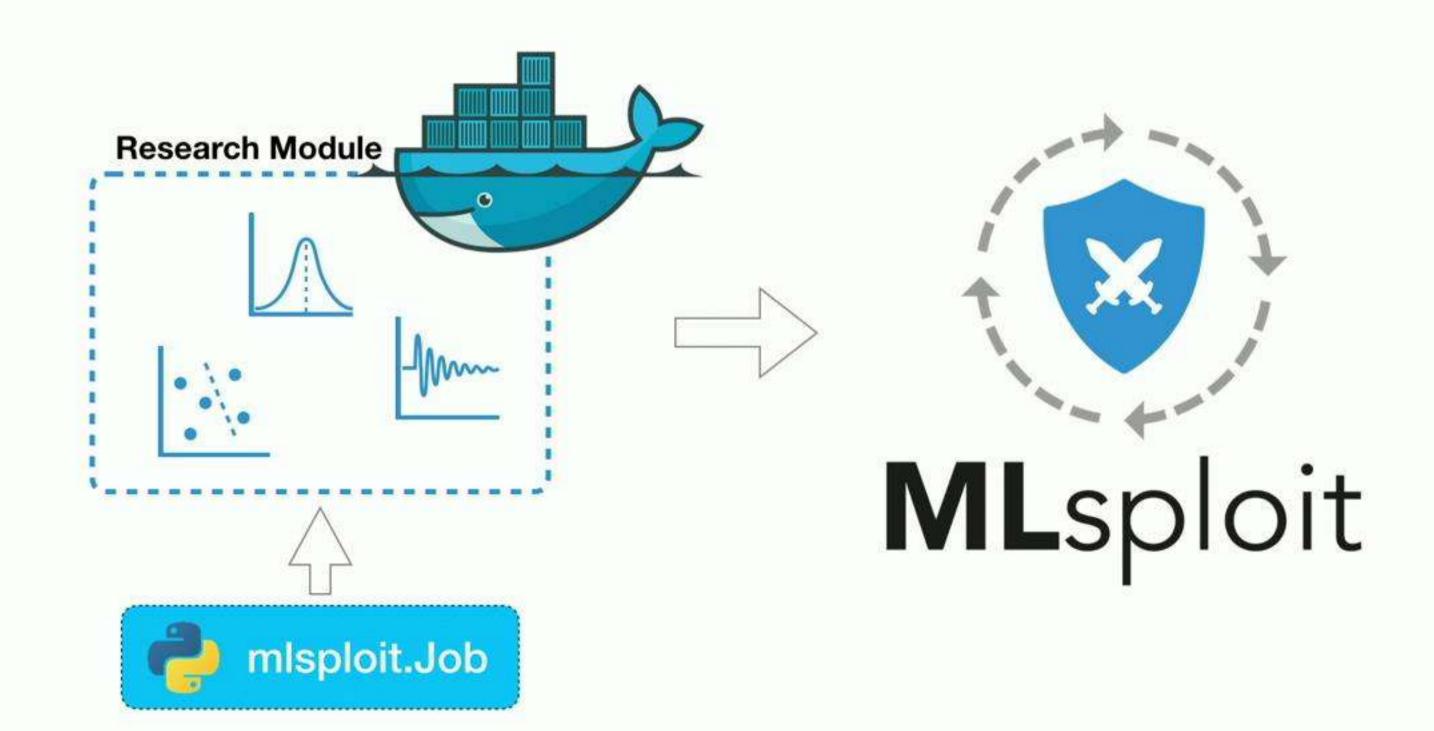


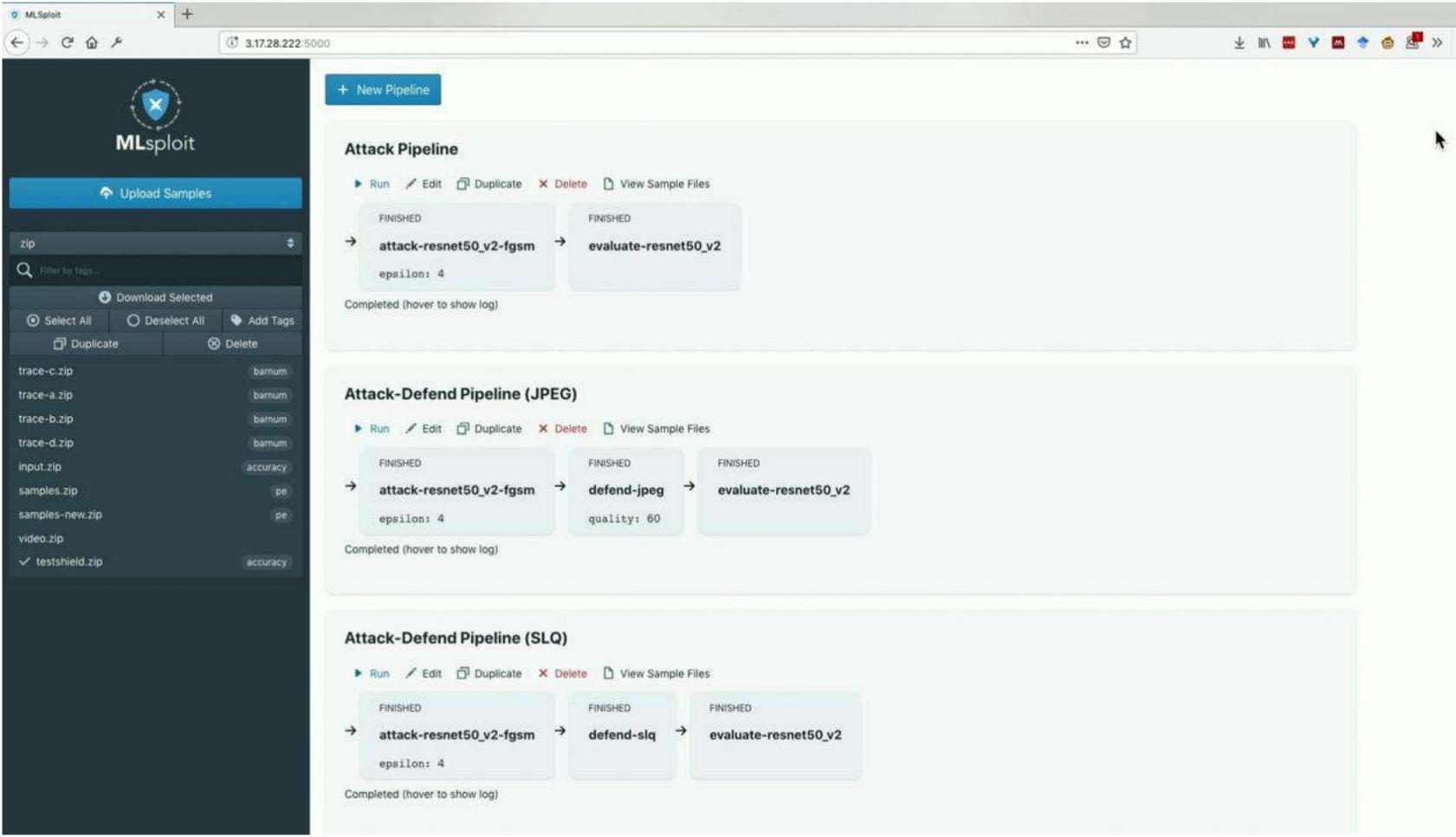


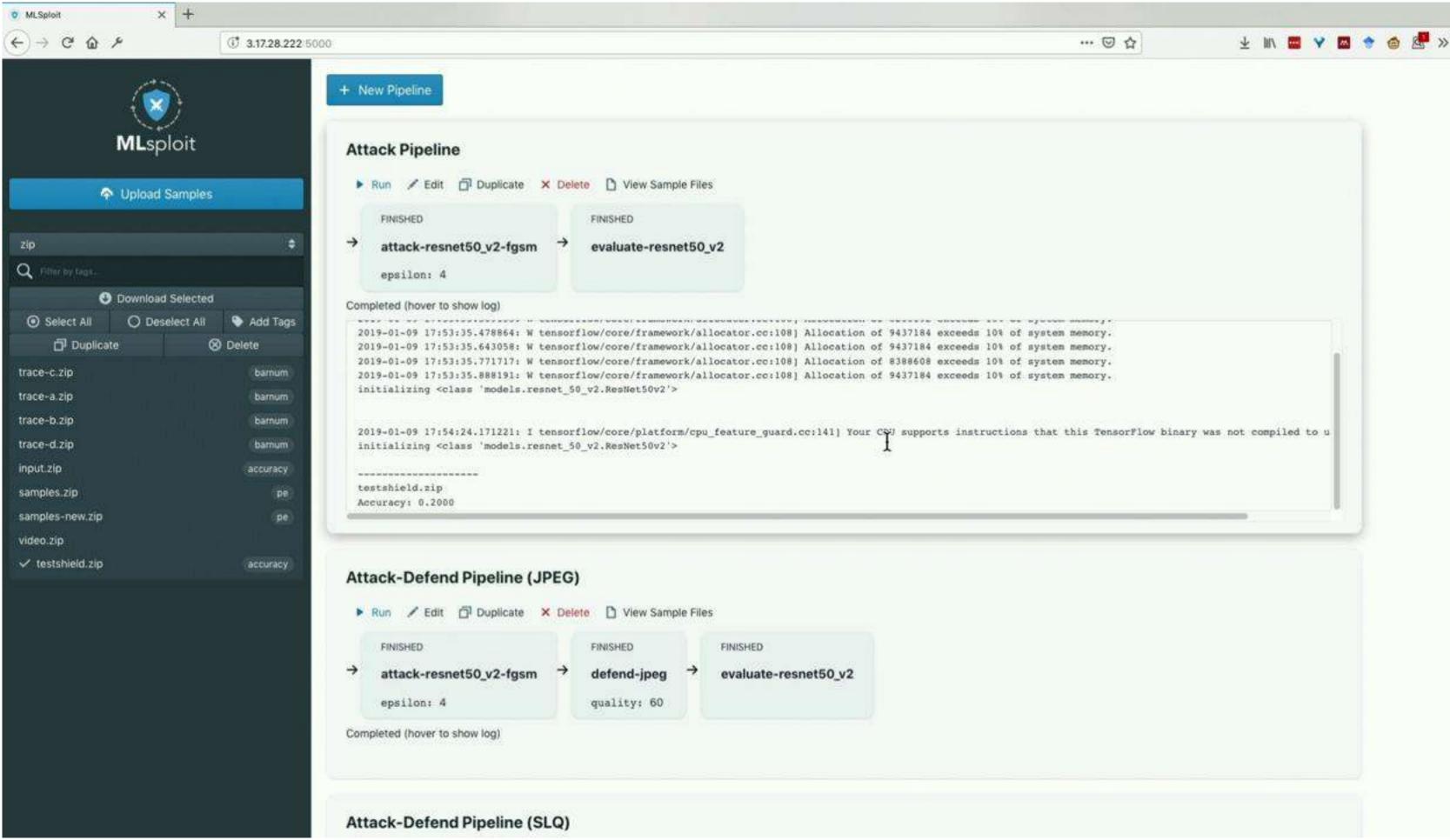
#### **ONE-STEP INSTALLATION**

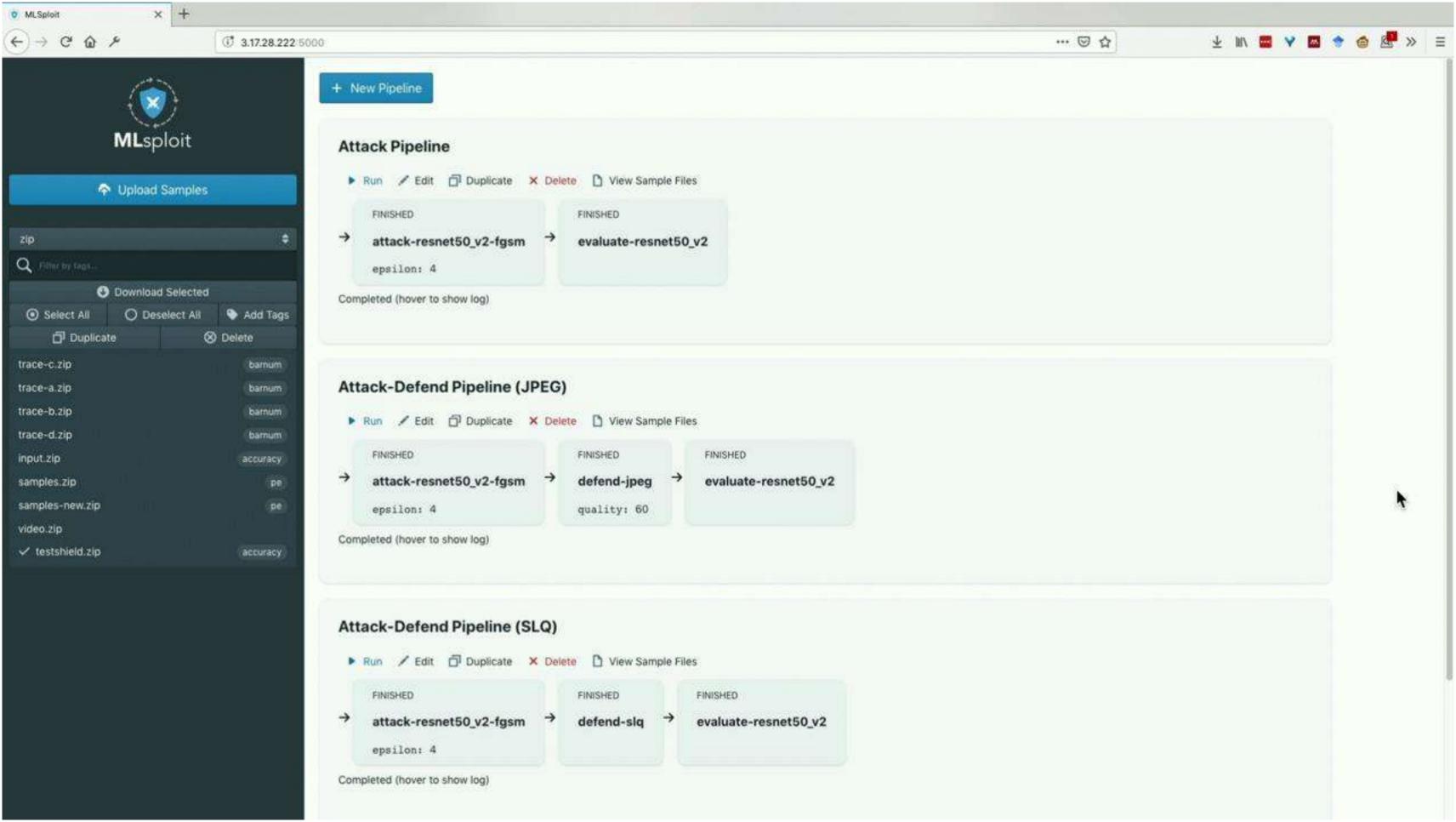


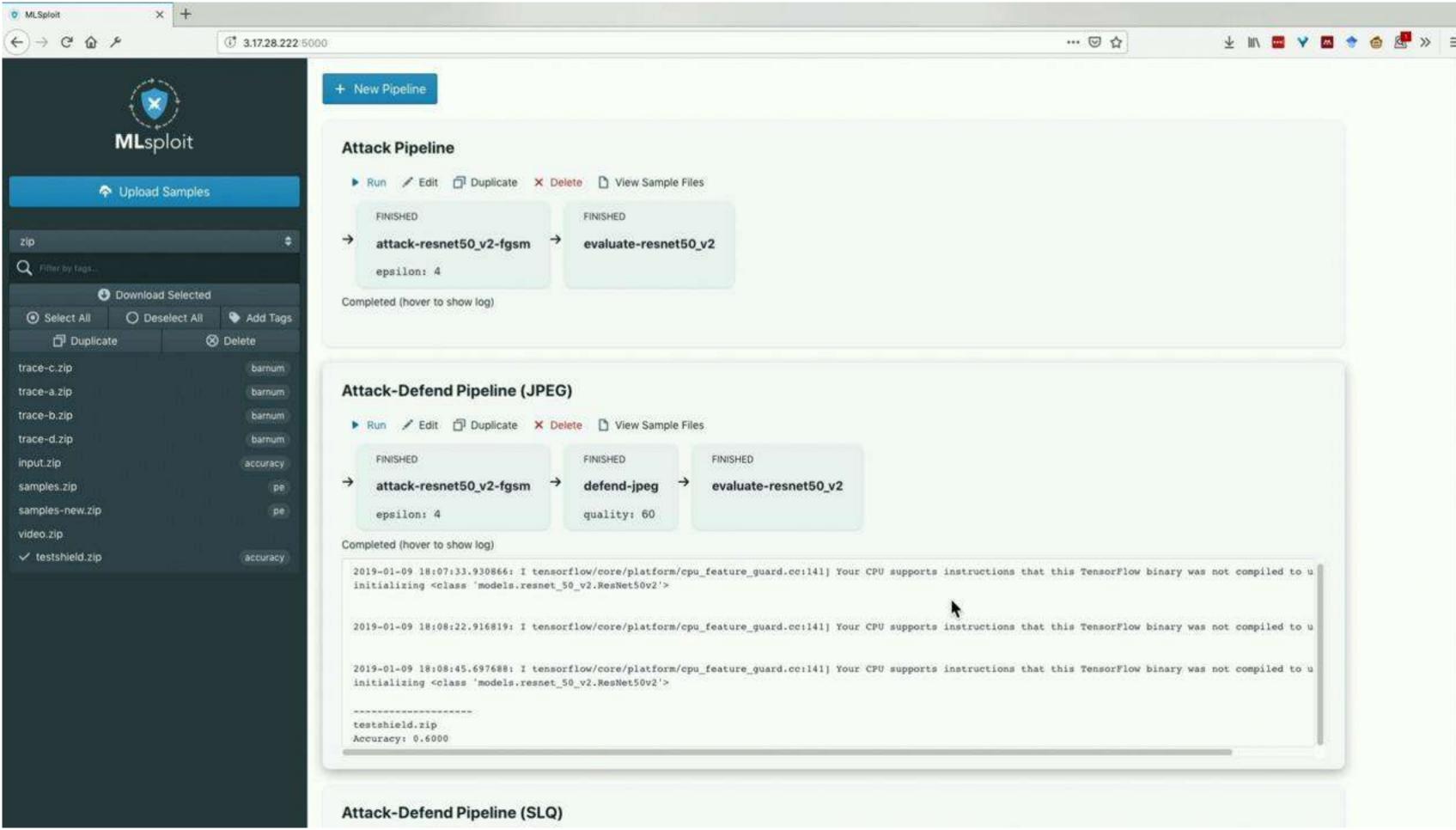
#### **EASY INTEGRATION OF RESEARCH**





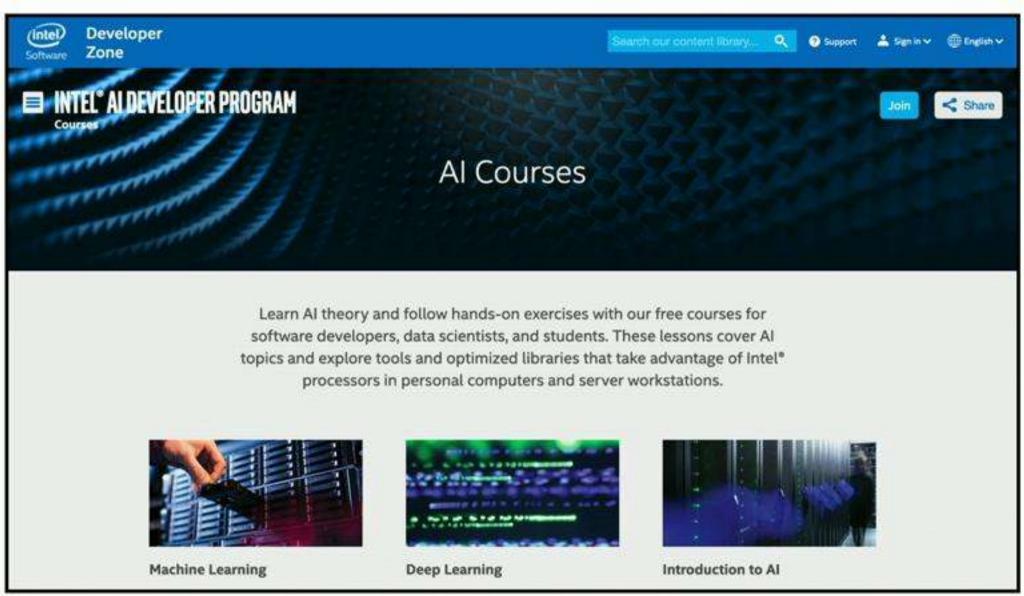






#### Intel® AI Courses





software.intel.com/en-us/ai/courses

## Secure

## Interpretable

Attack & Defense (DNN)

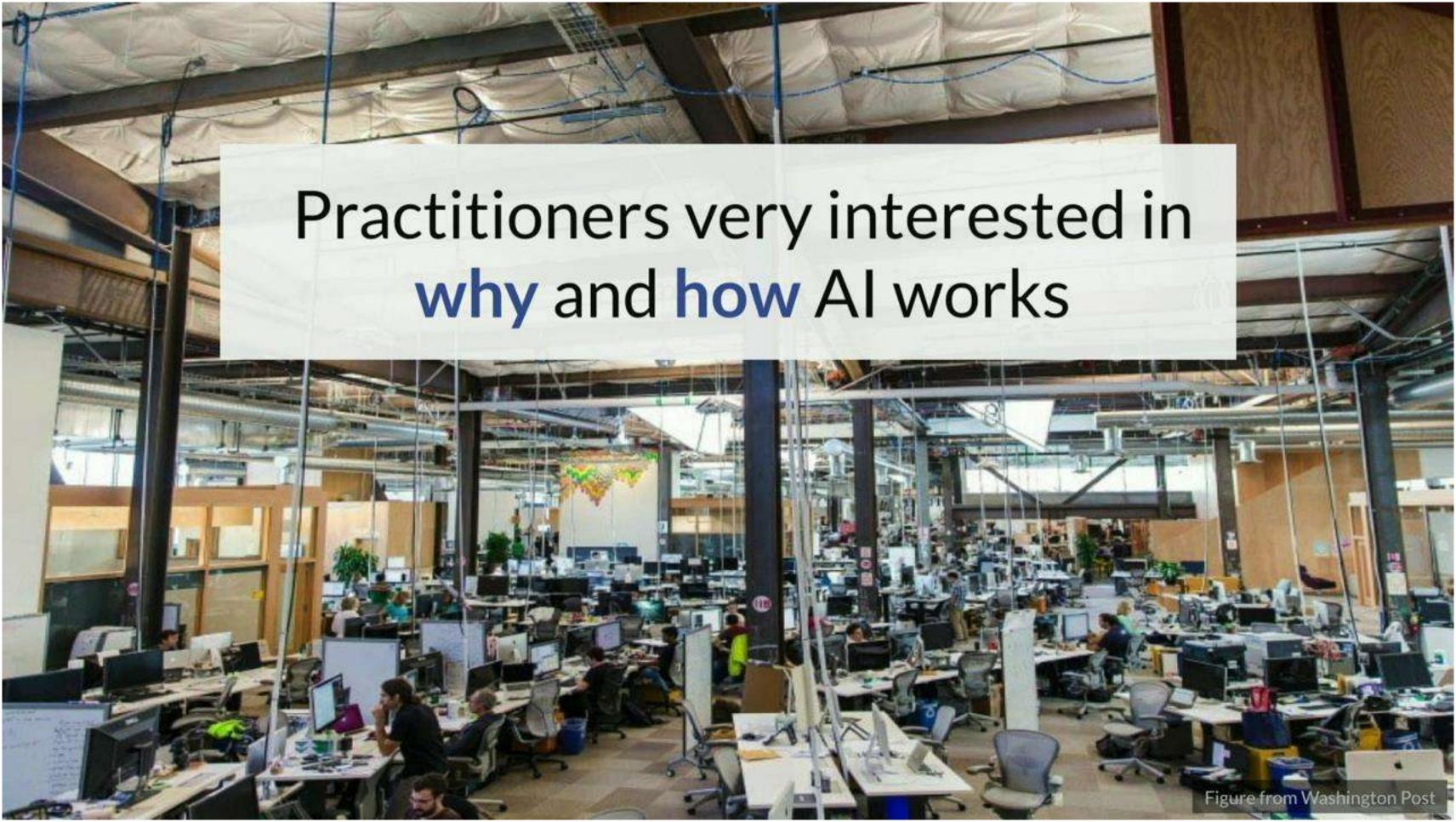
ShapeShifter Shield

Do-it-yourself Adversarial ML

ADAGIO MLsploit Understand Industry Models
ActiVis

Interactive Learning (Education)
GAN Lab

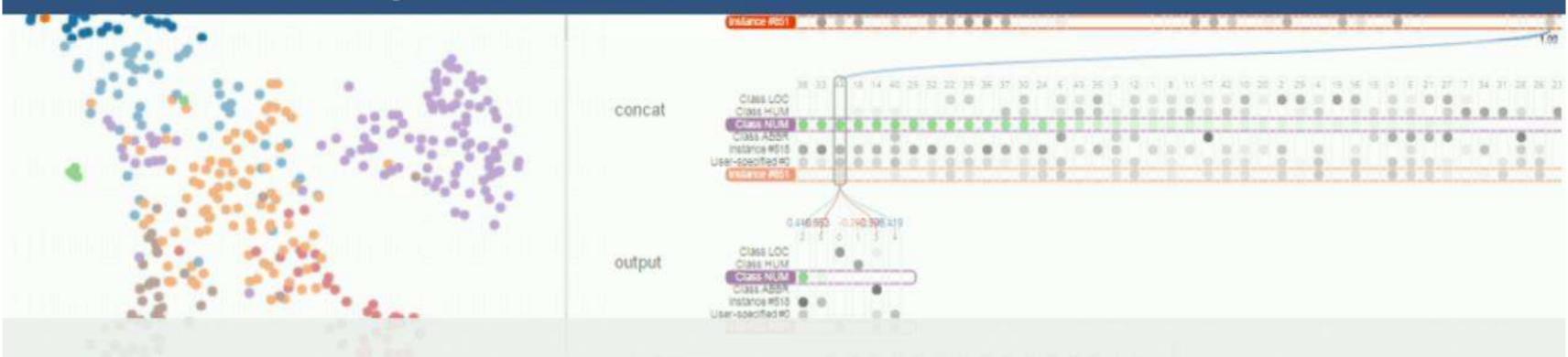
Research landscape Survey, Gamut



#### **OUR KEY IDEA**

# Scalable Interactive visualization as a medium for connecting users with ML models

#### Why interactive visualization?

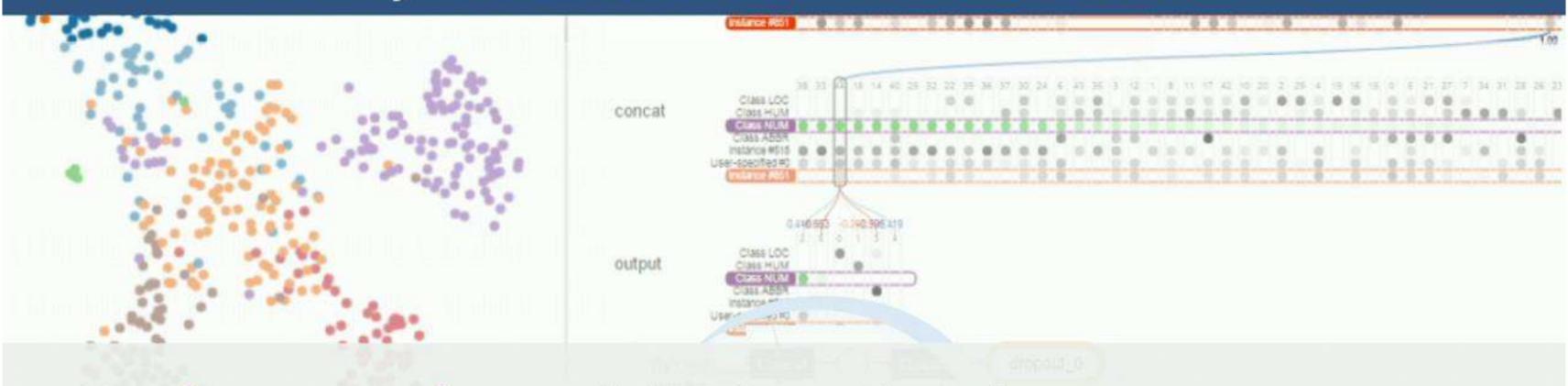


Machine learning aims to find patterns from data.

Visualization amplifies human cognition to find patterns.



#### Why interactive visualization?



By interacting with visualization, users can incrementally make sense of AI models.



## nterpretable AI via Visual Analytics

**Understanding Industry-Scale Models** 

ActiVis - Activation analysis by subsets

Interactive Learning of Complex Models

GAN Lab - Experimentation with GANs

Research Landscape

Survey, Gamut

## ACTIVIS Scaling Visualization to Industry-scale Models & Data

IEEE VIS 2017

#### Deployed by facebook



Minsuk Kahng Georgia Tech



Pierre Andrews Facebook



Aditya Kalro Facebook



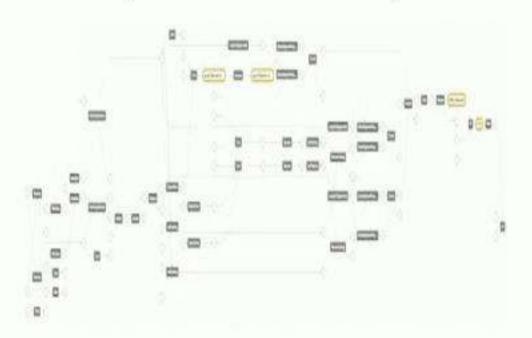
Polo Chau Georgia Tech



### Practical Design Challenges

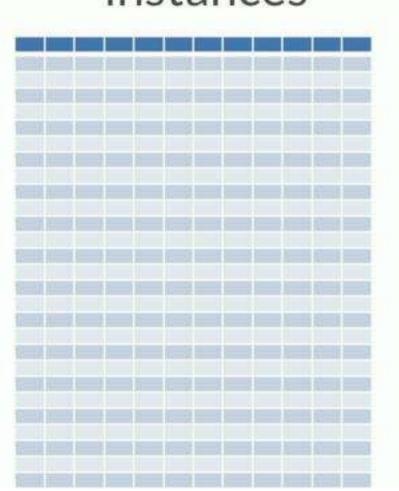
#### **DEEP/WIDE MODELS**

1,000+ operations/layers



#### **LARGE DATASETS**

1 billion+ instances

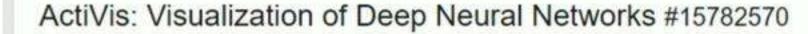


#### **DIVERSE FEATURES**

image, text, numerical, categorical, ...

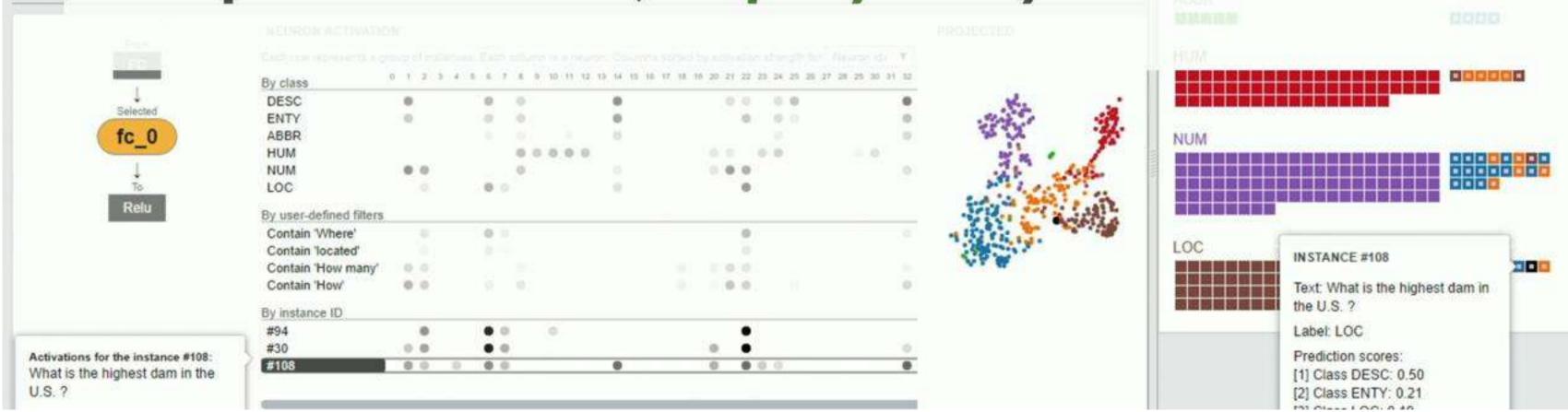








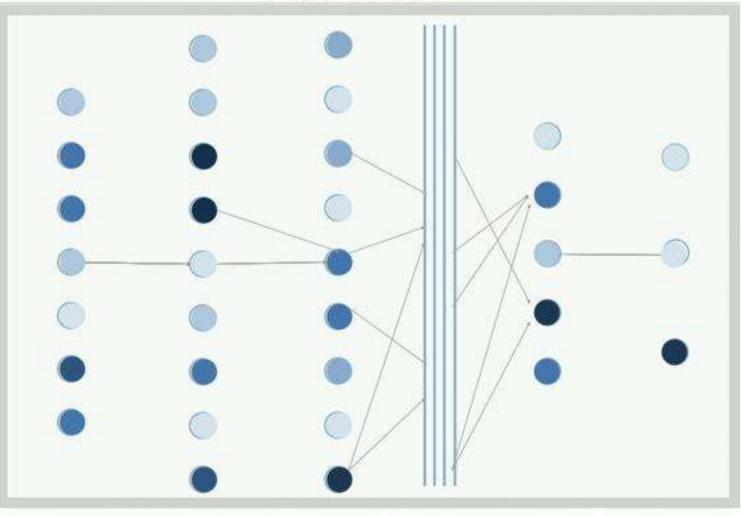
### Visualizing activation of industry-scale deep neural nets, deployed by Facebook



### How to visualize many model parameters?

**INPUT** 

Where is Mercedes-Benz Stadium located? MODEL



OUTPUT

Number 11%

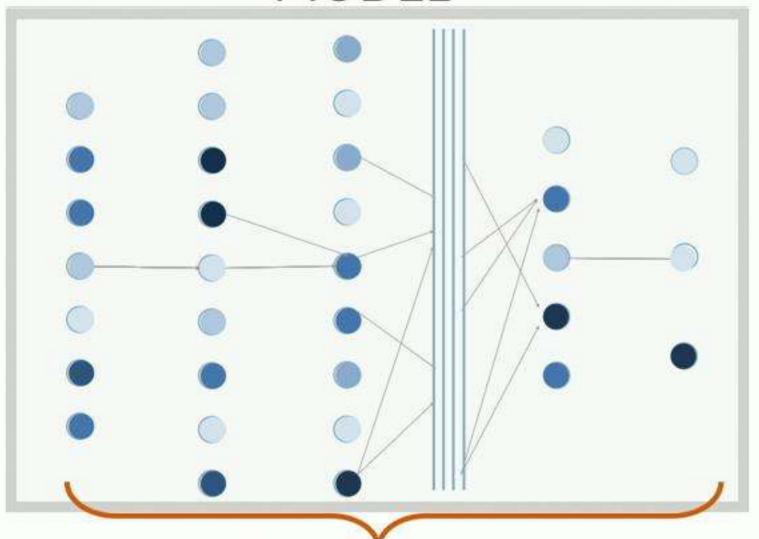
Person 8%

**Location 81%** 

### How to visualize many model parameters?

**INPUT** 

Where is Mercedes-Benz Stadium located? MODEL



OUTPUT

Number 11%

Person 8%

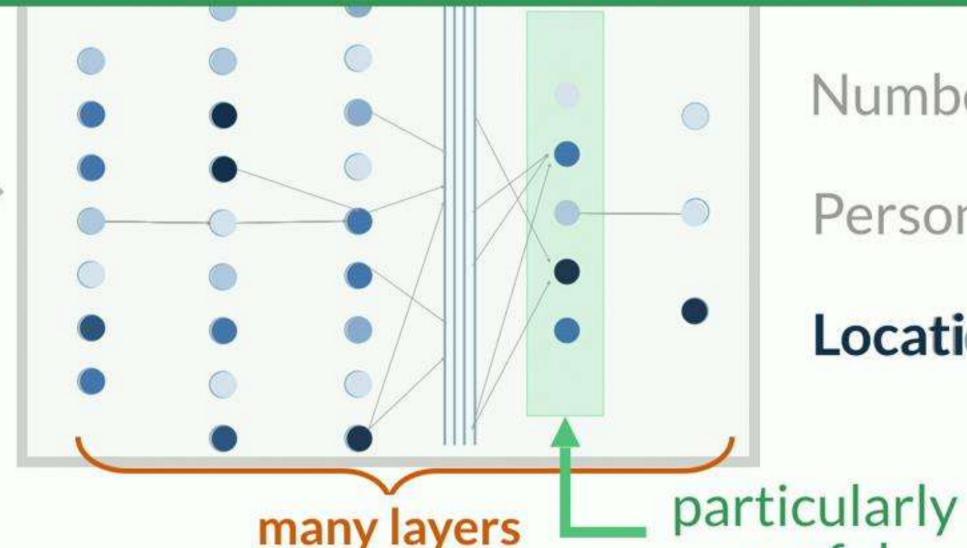
**Location 81%** 

many layers

### How to visualize many model parameters?

Observation: No need to show everything

Where is Mercedes-Benz Stadium located?



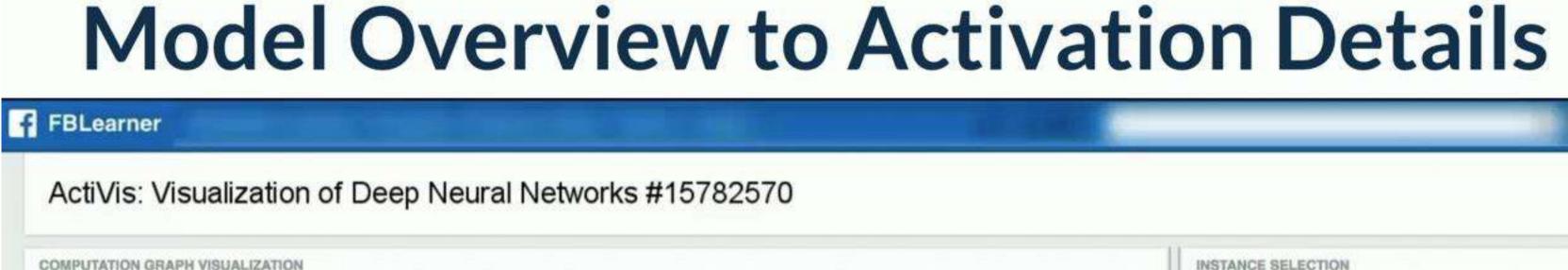
Number 11%

Person 8%

**Location 81%** 

useful

#### ActiVis Key Ideas #1

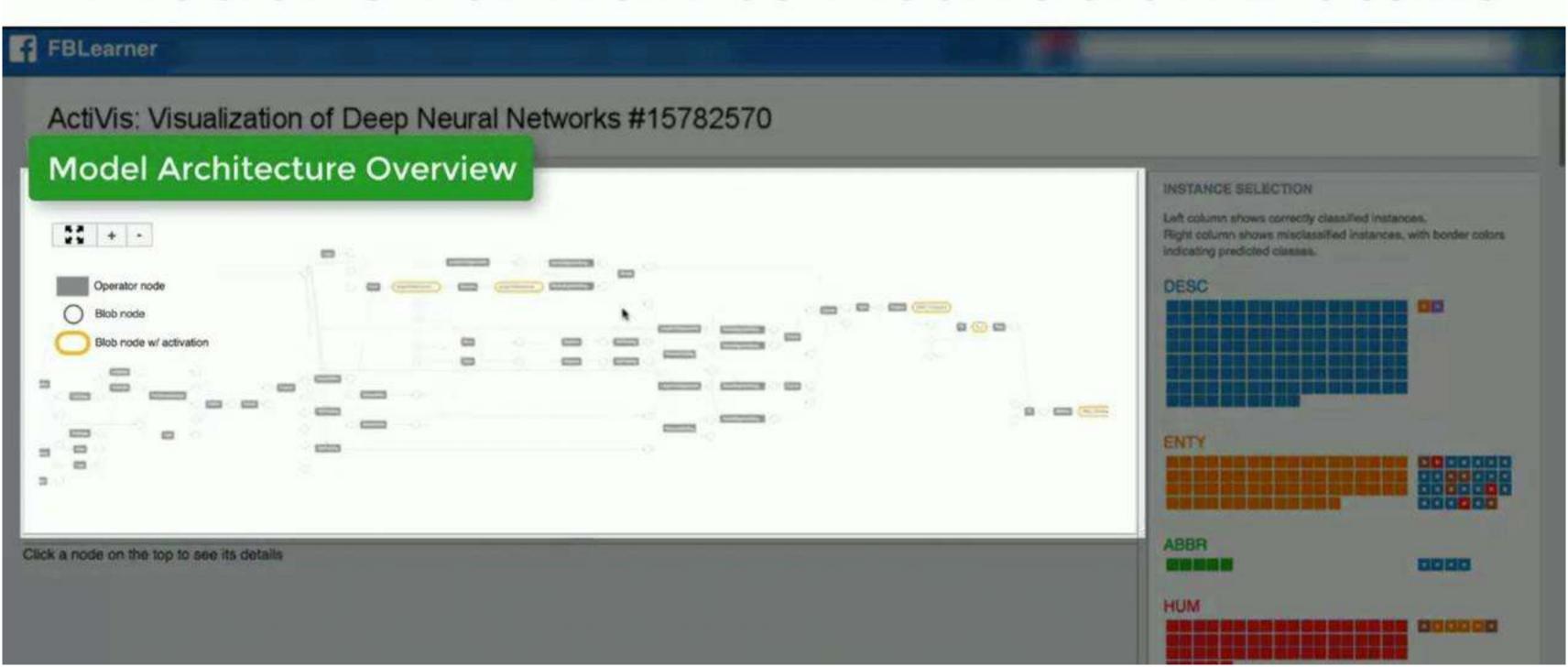




Left column shows correctly classified instances. Right column shows misclassified instances, with border colors indicating predicted classes. DESC ENTY ABBR HUM

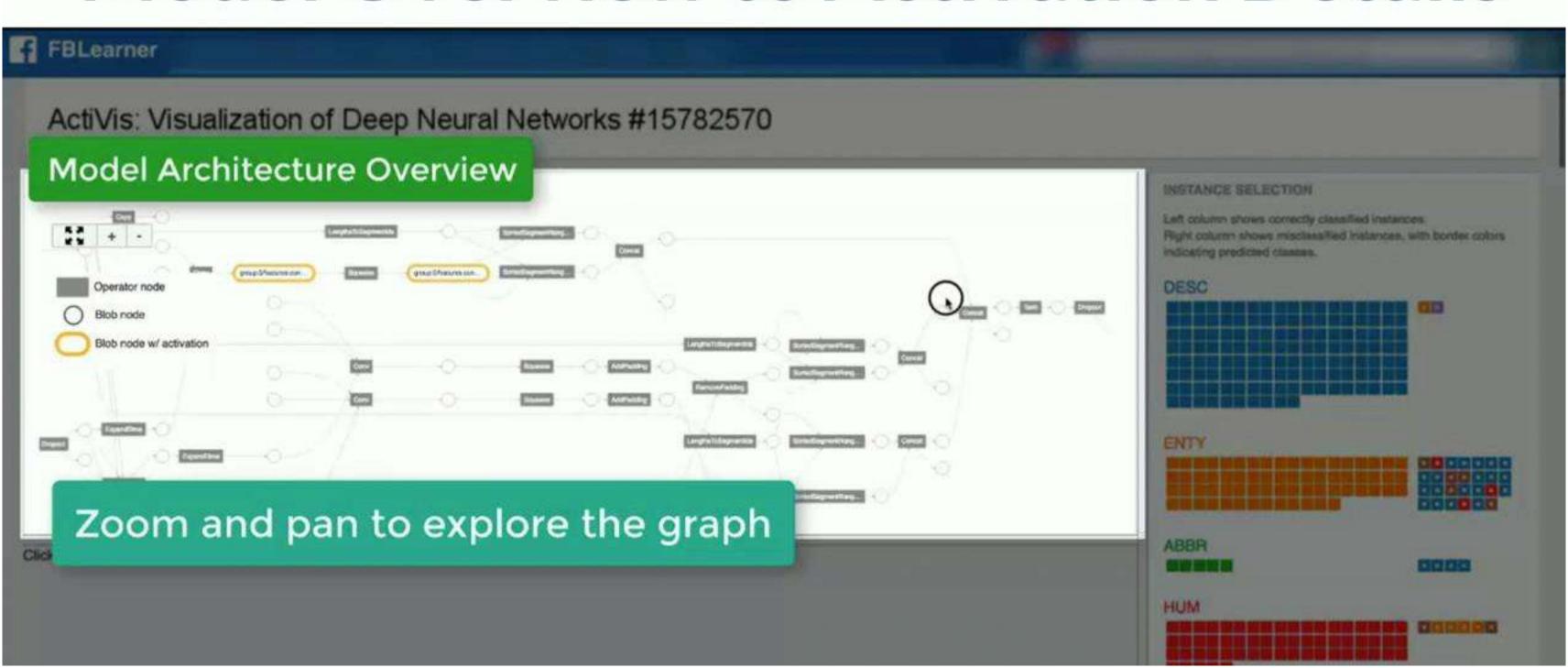
#### ActiVis Key Ideas #1

### Model Overview to Activation Details



#### ActiVis Key Ideas #1

### Model Overview to Activation Details



### How to analyze many data instances?

### Observation: Two Analytics Patterns

Complementary

INSTANCE-LEVEL



SUBSET-LEVEL

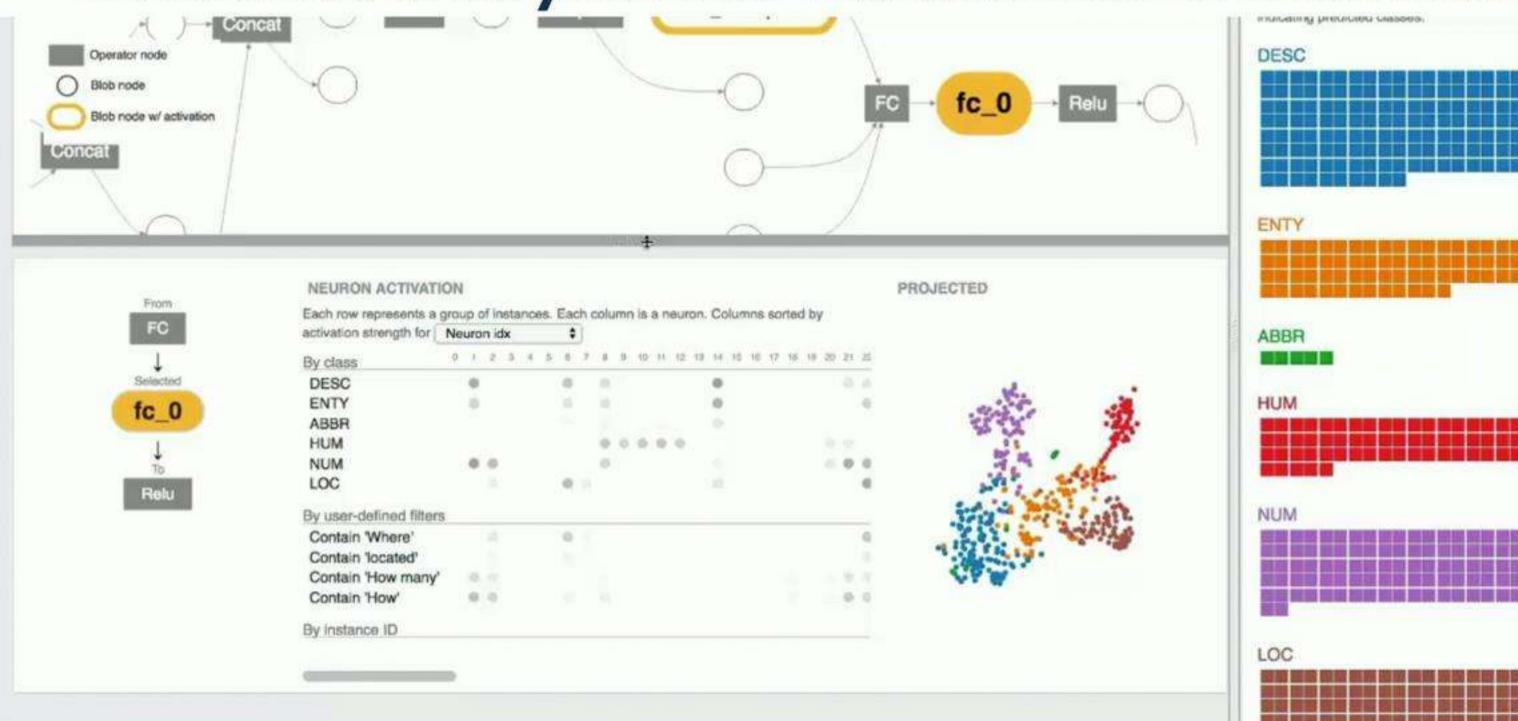
How model responds to individual instances?

How model behaves at higher-level categorization (e.g., by topic)?

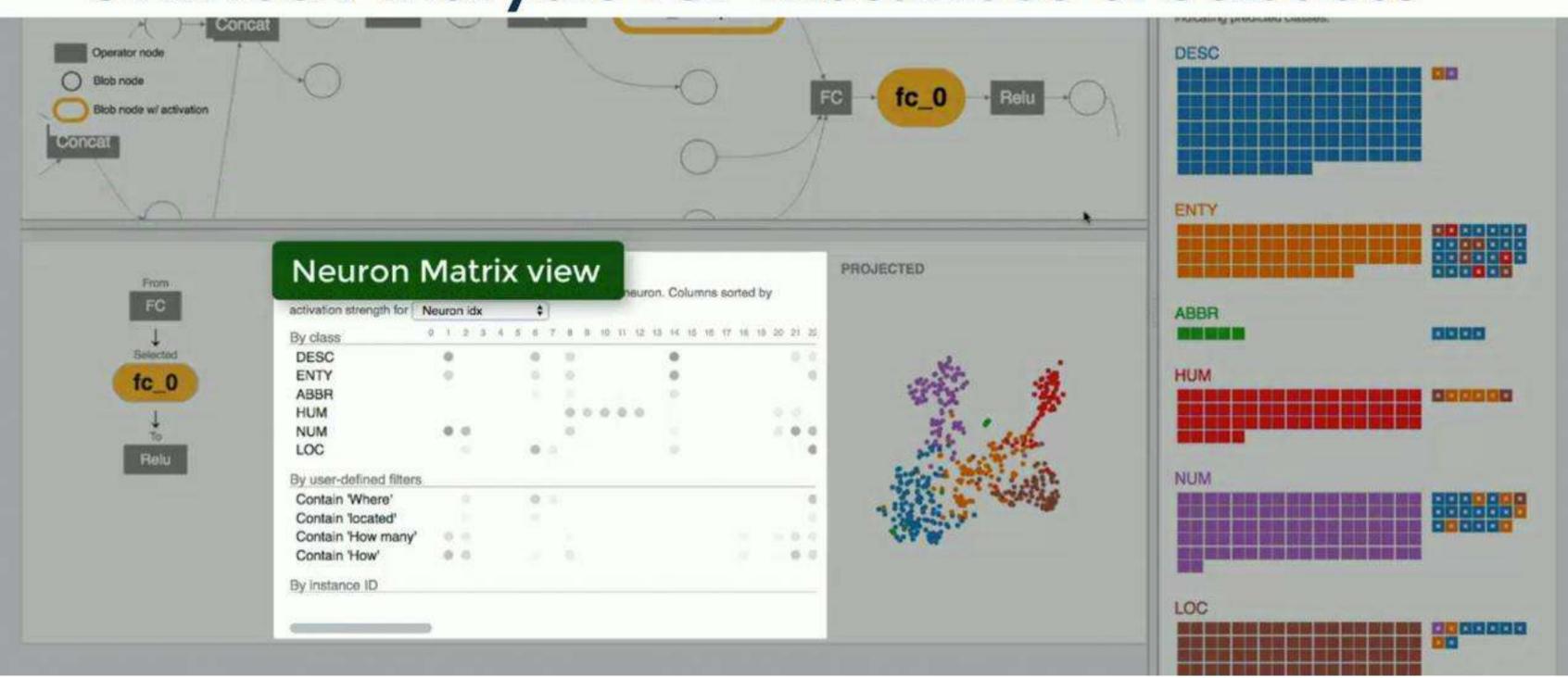
Useful for debugging

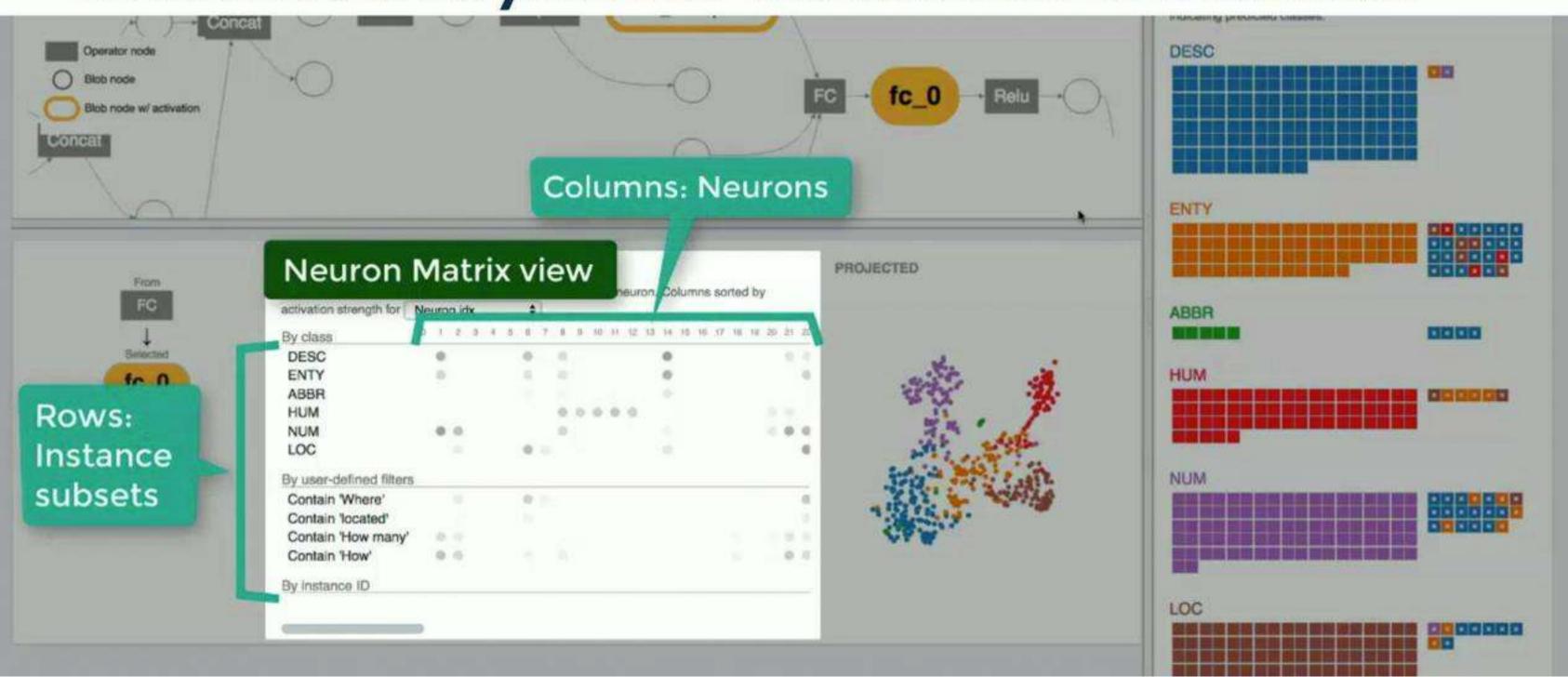
Useful for large datasets

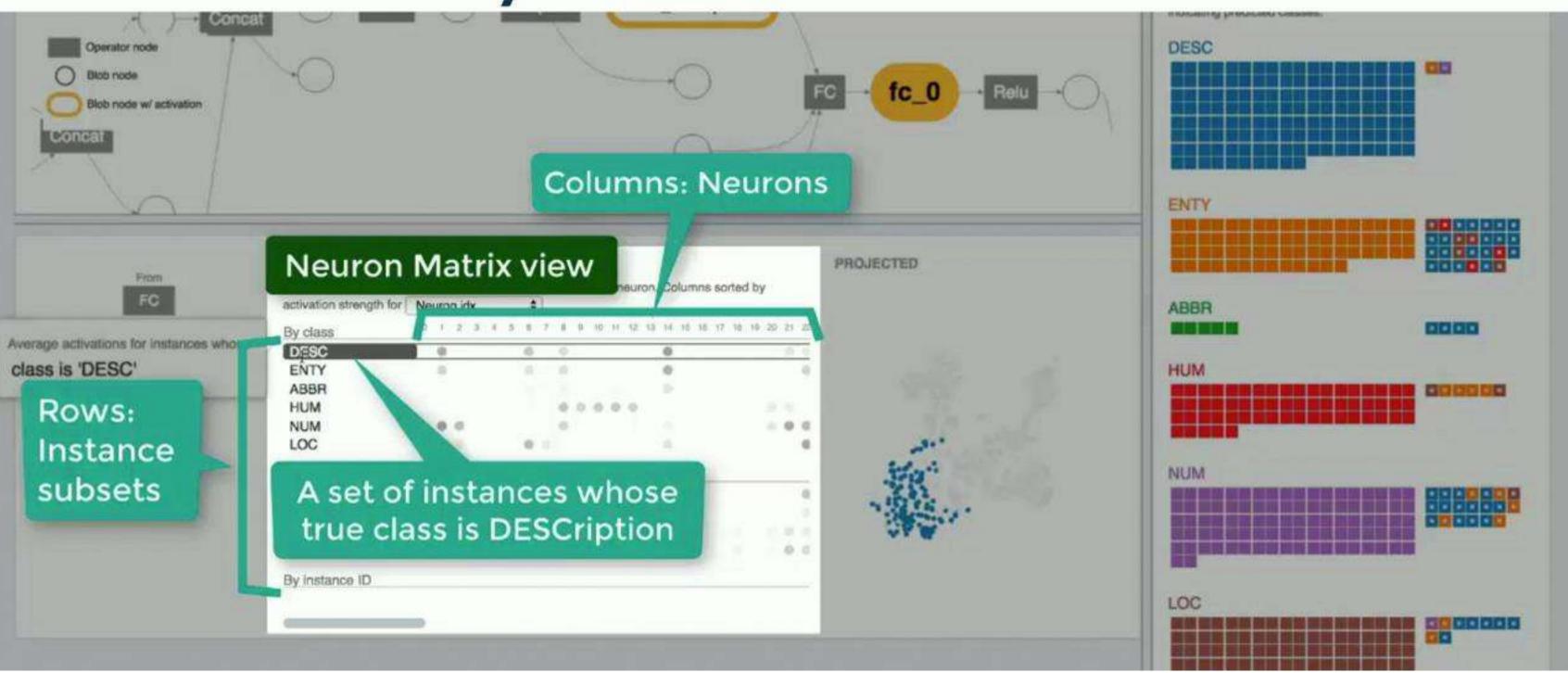
Unified Analysis for Instances & Subsets

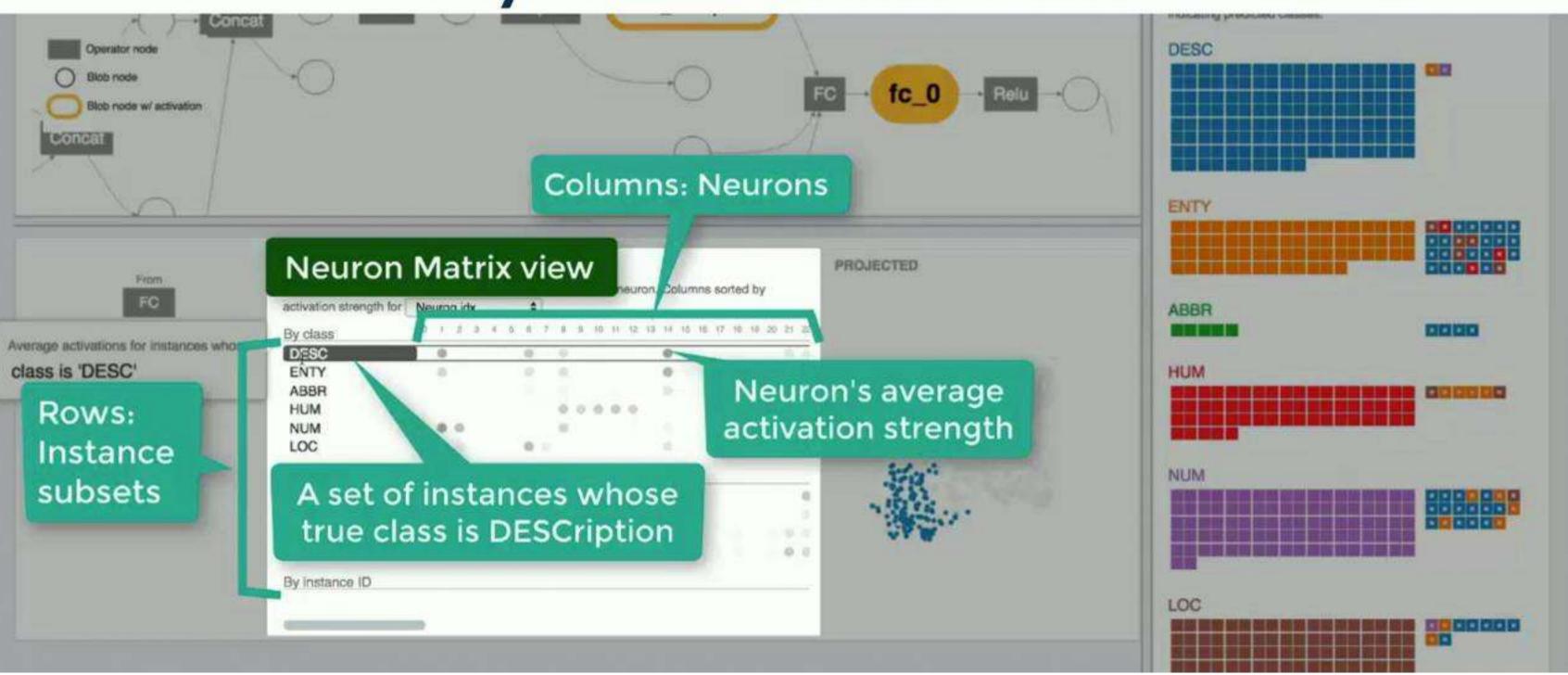


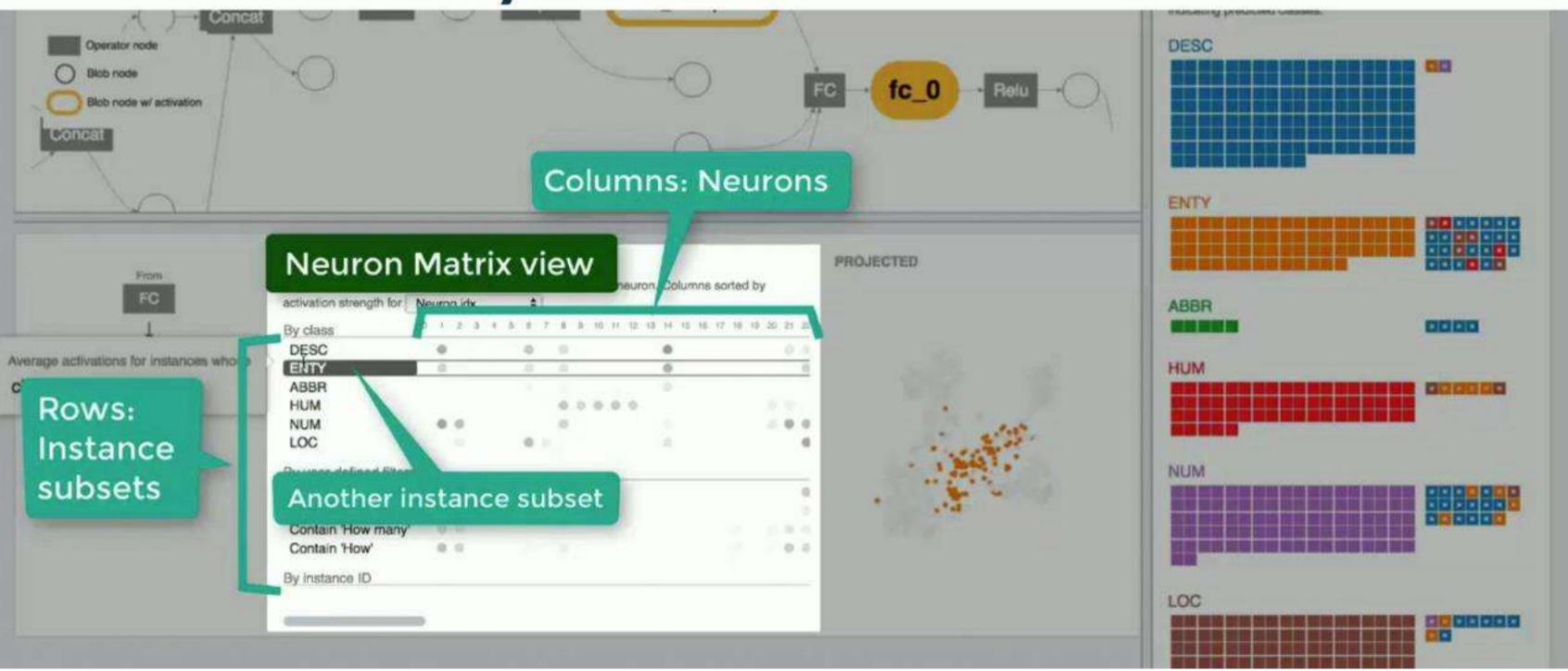












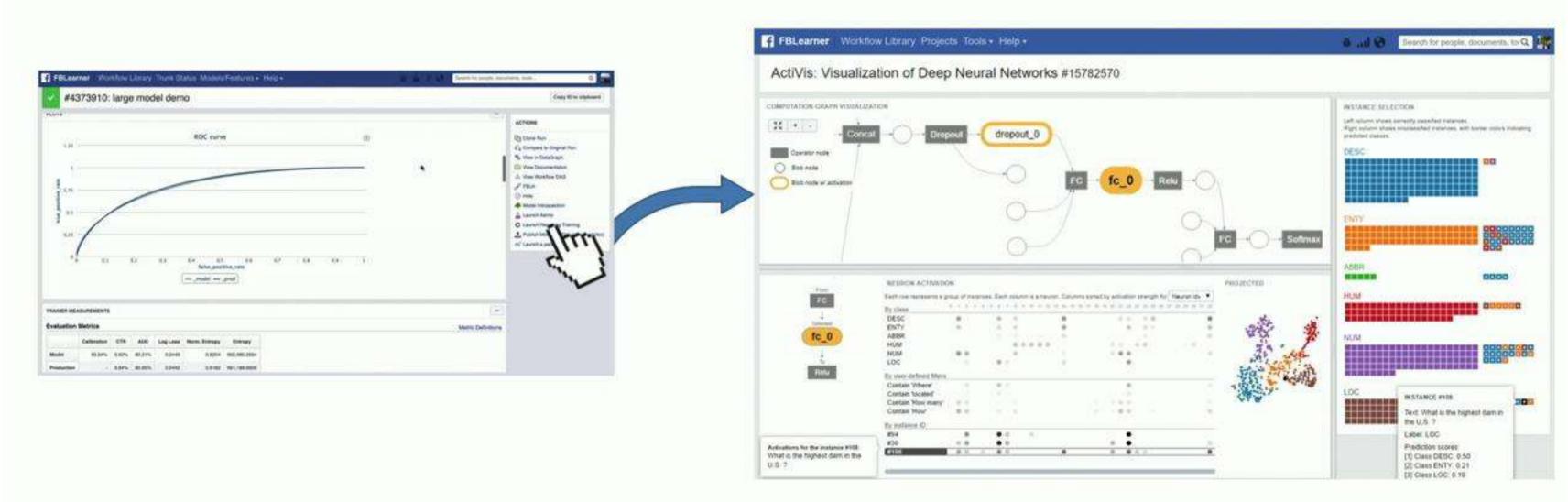
### Scaling Up ActiVis for Facebook

- 1. User-guided Instance Sampling
- 2. Selective Pre-computation of Layers
- 3. Matrix Computation for Billion-Scale Instances

### Deployed on FBLearner

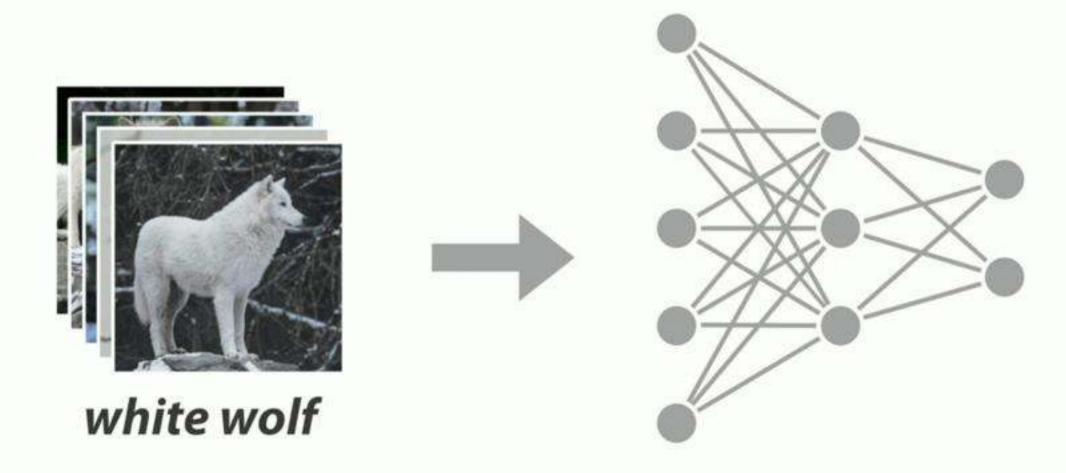
### Facebook's ML platform

used by >25% of engineering team



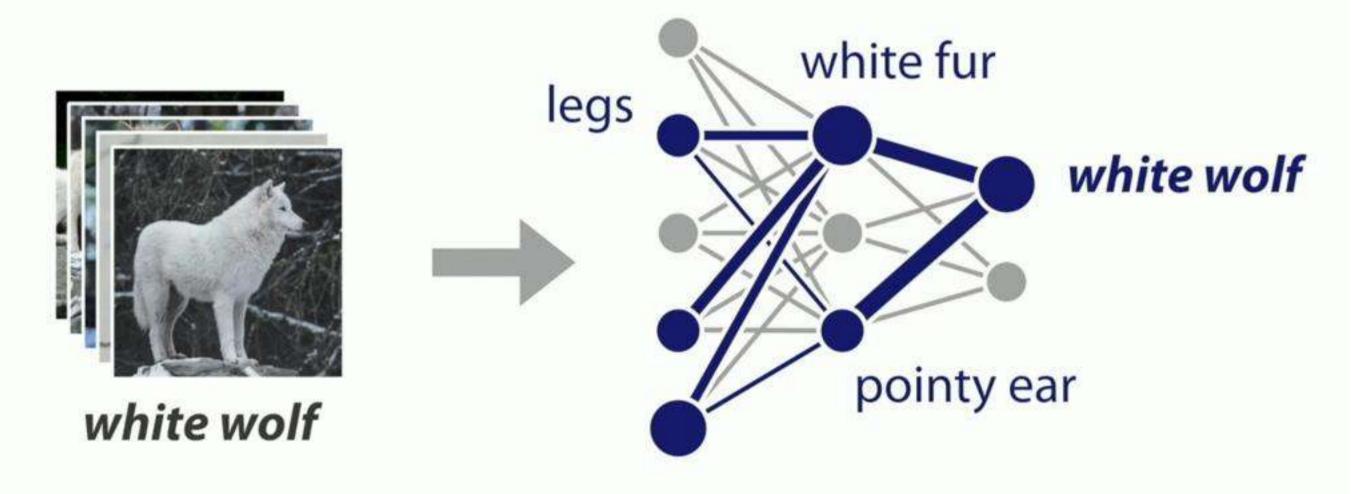
### SUMMIT

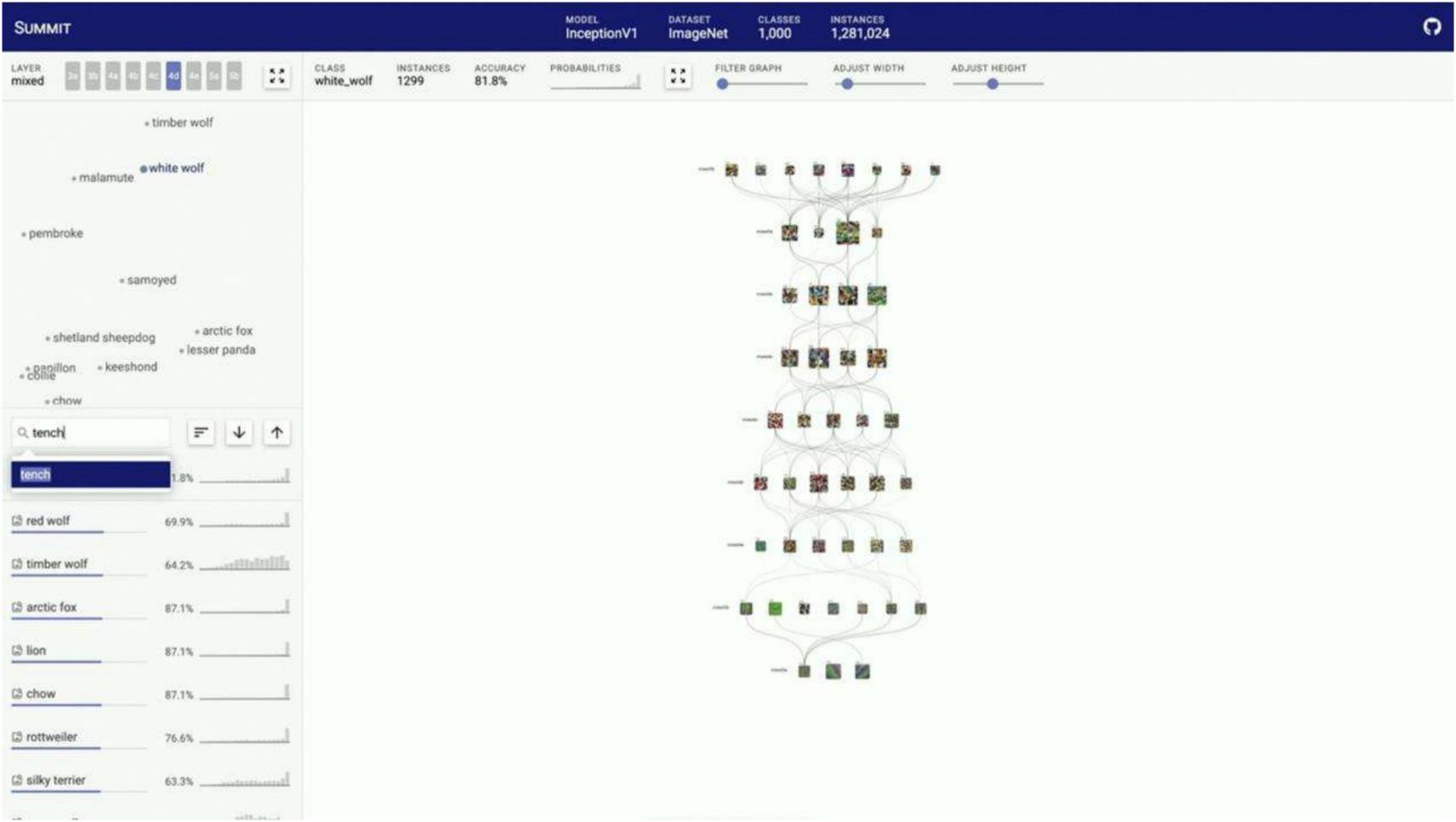
Scalably summarize and interactively visualize neural network feature representations for millions of images

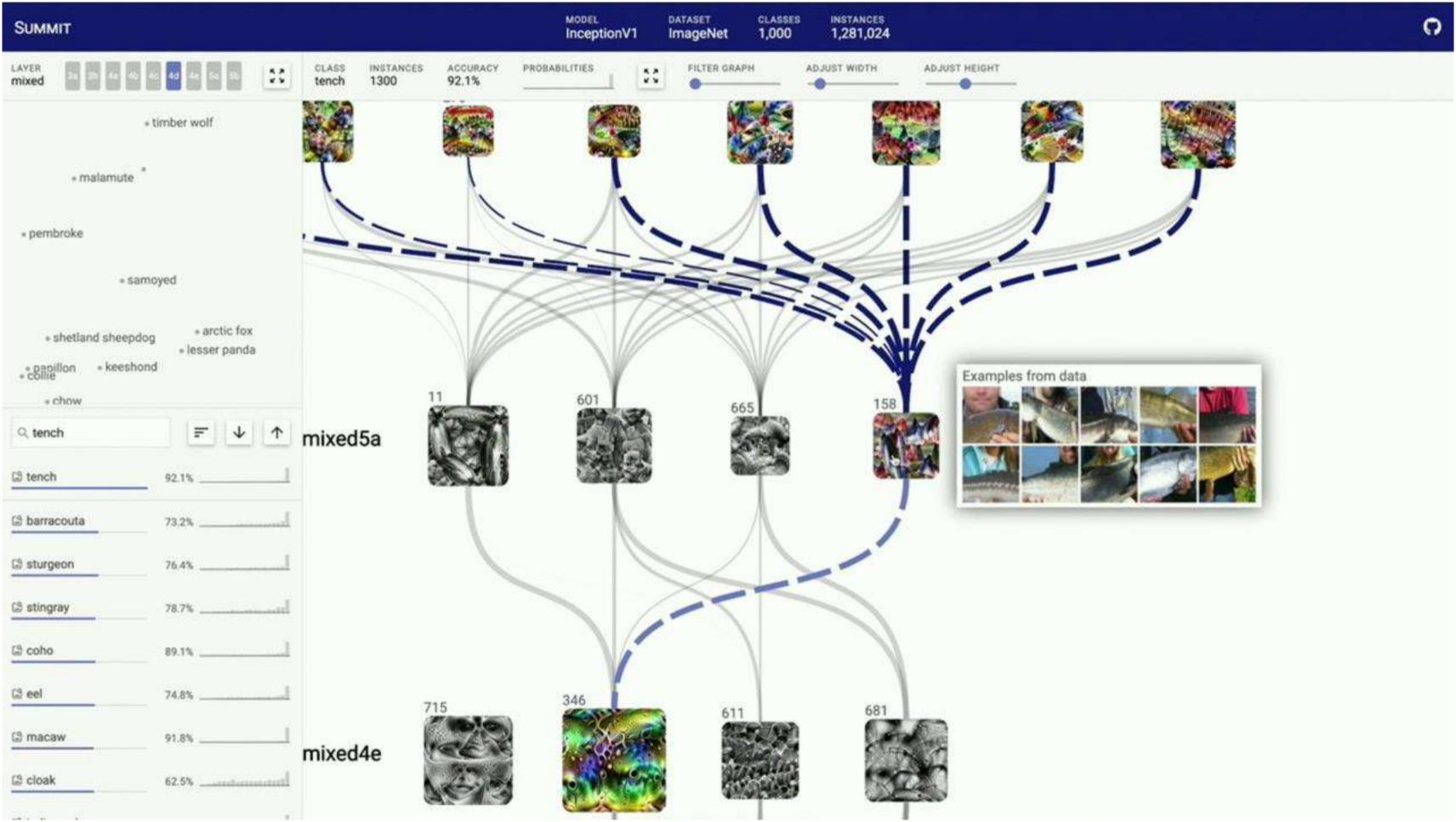


### SUMMIT

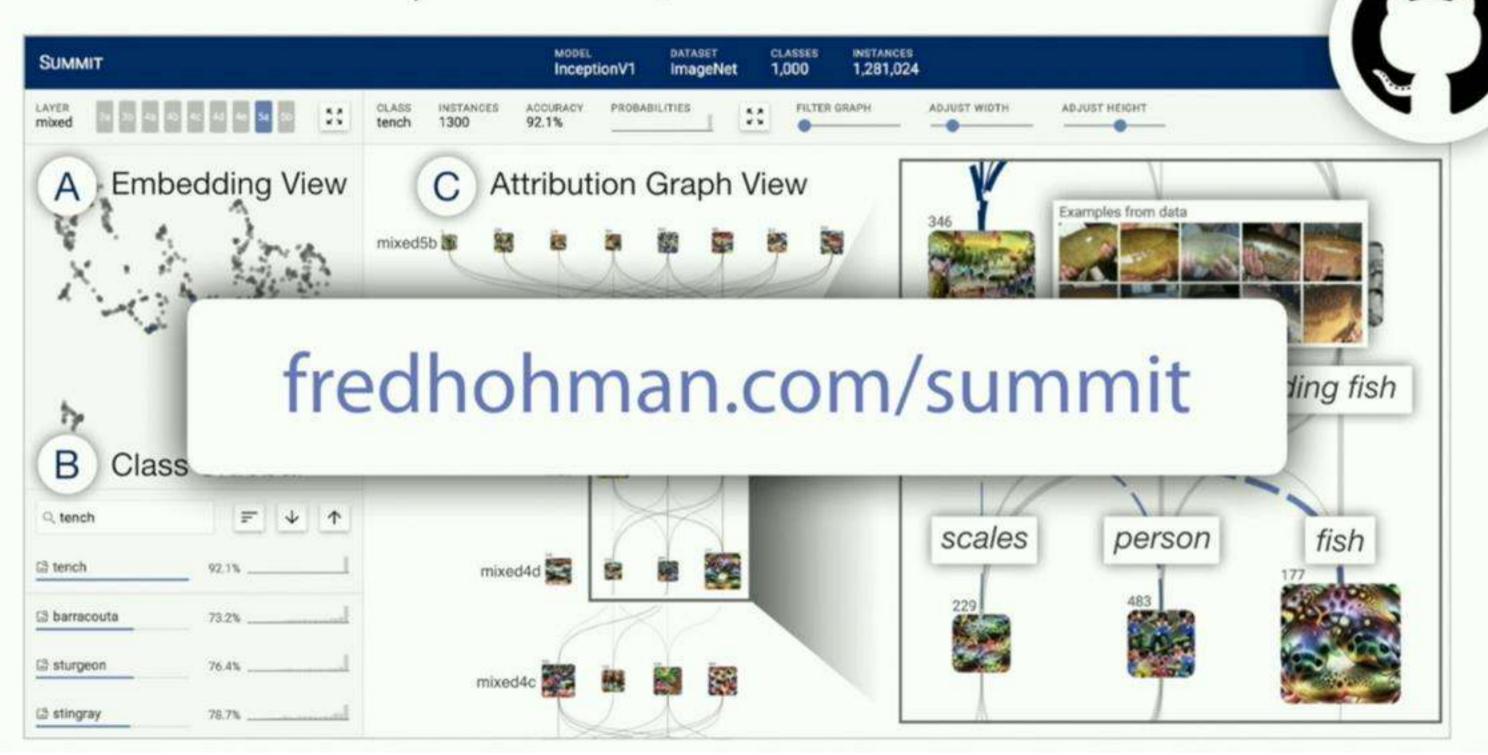
Scalably summarize and interactively visualize neural network feature representations for millions of images







See more in the paper, and try our open-source demo!



### SUMMIT

Scaling Deep Learning Interpretability by Visualizing Activation and Attribution Summarizations

Fred Hohman, Haekyu Park, Caleb Robinson, Duen Horng (Polo) Chau

IEEE VIS 2019 Vancouver, Canada



# nterpretable Al via Visual Analytics

**Understanding Industry-Scale Models** 

ActiVis - Activation analysis by subsets

Interactive Learning of Complex Models

GAN Lab - Experimentation with GANs

Research Landscape

Survey, Gamut

## GAN Lab

Understanding Complex Deep Generative Models using Interactive Visual Experimentation



Minsuk Kahng Georgia Tech



Nikhil Thorat Google



Polo Chau Georgia Tech



Fernanda Viégas Google



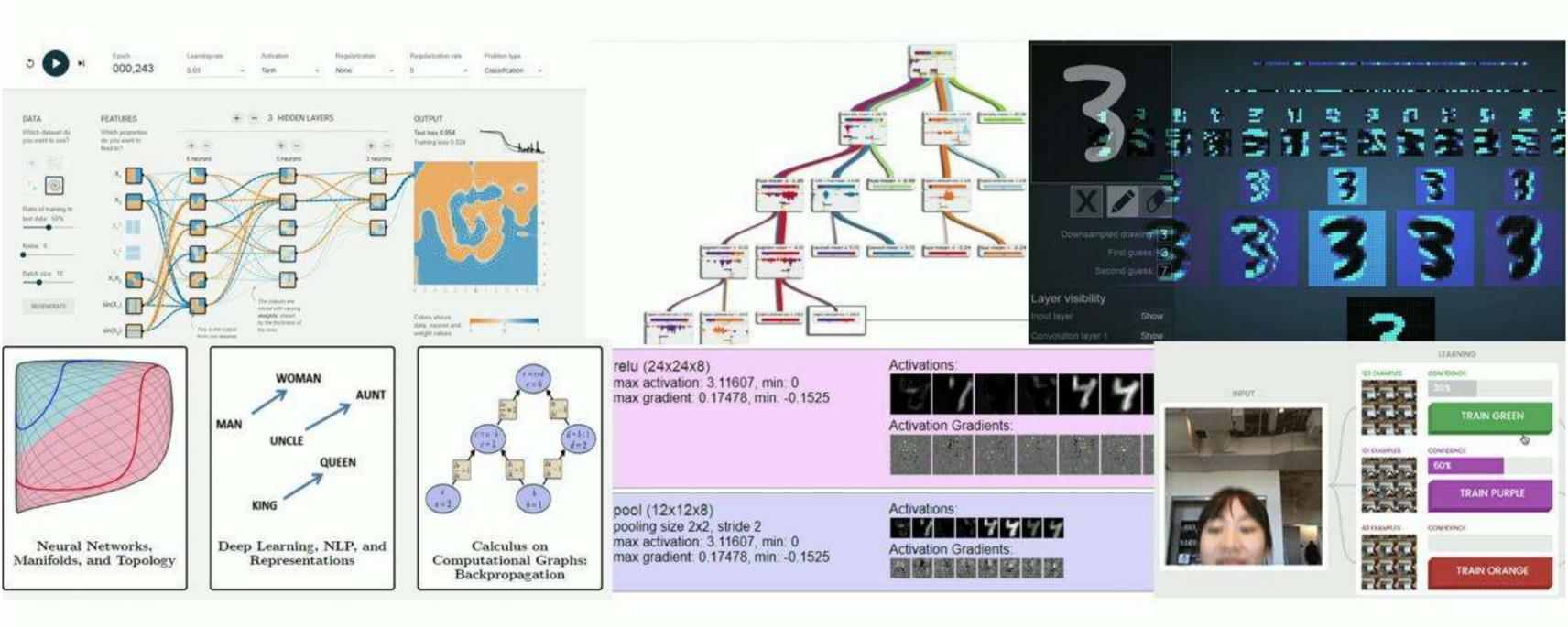
Martin Wattenberg Google



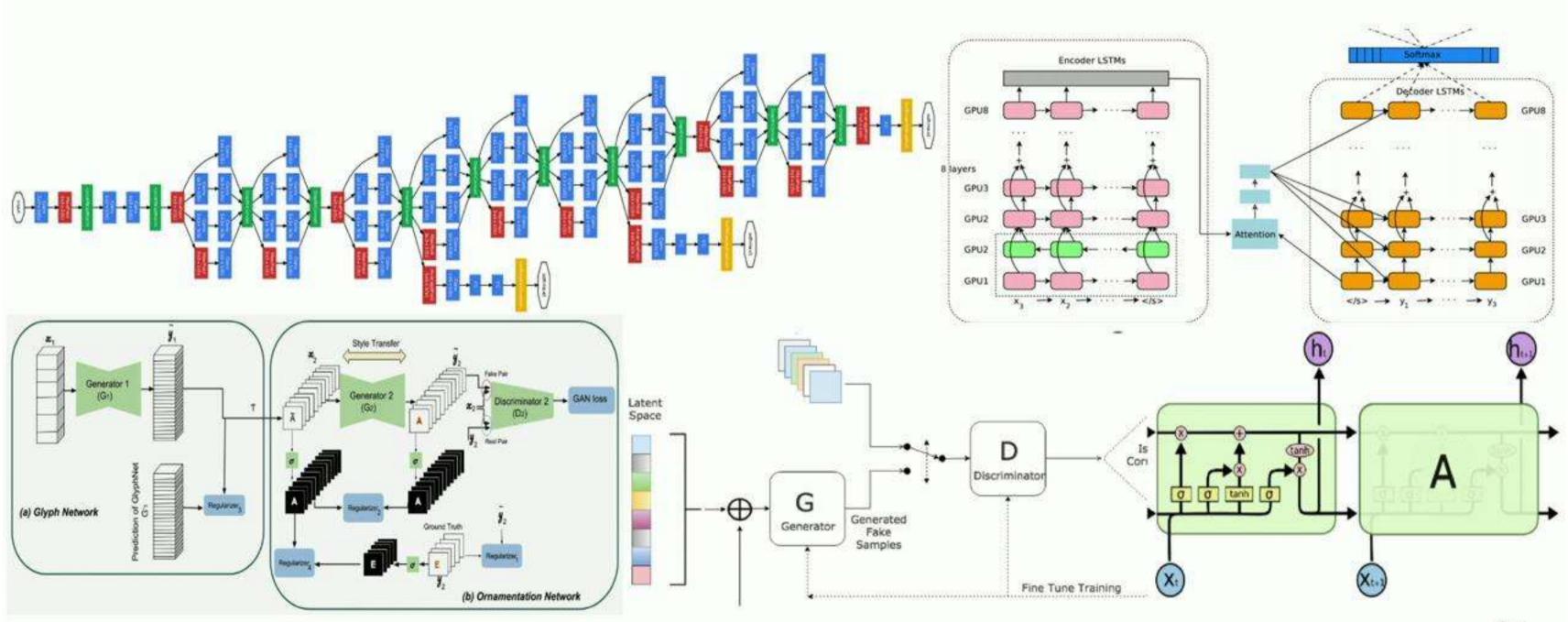


PAIR | People + AI Research Initiative

### Visualization for ML Education



### Modern deep models are complex



### Generative Adversarial Networks (GANs)

"the most interesting idea in the last 10 years in ML" - Yann LeCun



Face images generated by BEGAN

### Generative Adversarial Networks (GANs)

### Hard to understand and train even for experts

$$\min_{G} \max_{D} V(D, G) = \mathbb{E}_{\boldsymbol{x} \sim p_{\text{data}}(\boldsymbol{x})}[\log D(\boldsymbol{x})] + \mathbb{E}_{\boldsymbol{z} \sim p_{\boldsymbol{z}}(\boldsymbol{z})}[\log(1 - D(G(\boldsymbol{z})))].$$

#### Discriminator

$$\nabla_{\theta_d} \frac{1}{m} \sum_{i=1}^m \left[ \log D\left(\boldsymbol{x}^{(i)}\right) + \log\left(1 - D\left(G\left(\boldsymbol{z}^{(i)}\right)\right)\right) \right].$$

#### Generator

$$\nabla_{\theta_g} \frac{1}{m} \sum_{i=1}^m \log \left( 1 - D\left( G\left(\boldsymbol{z}^{(i)}\right) \right) \right).$$

A GAN uses two competing neural networks

A GAN uses two competing neural networks

Generator synthesizes outputs

Discriminator spots fake

A GAN uses two competing neural networks

Generator synthesizes outputs



Counterfeiter makes fake bills

Discriminator spots fake

A GAN uses two competing neural networks

Generator synthesizes outputs



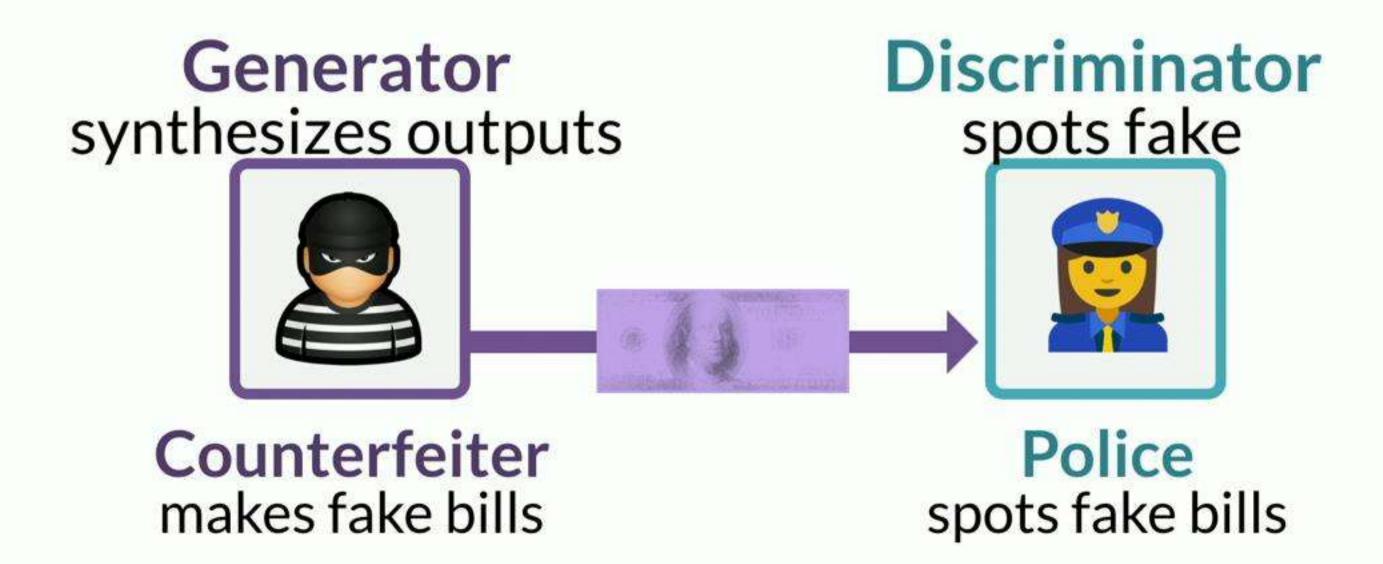
Counterfeiter makes fake bills

Discriminator spots fake

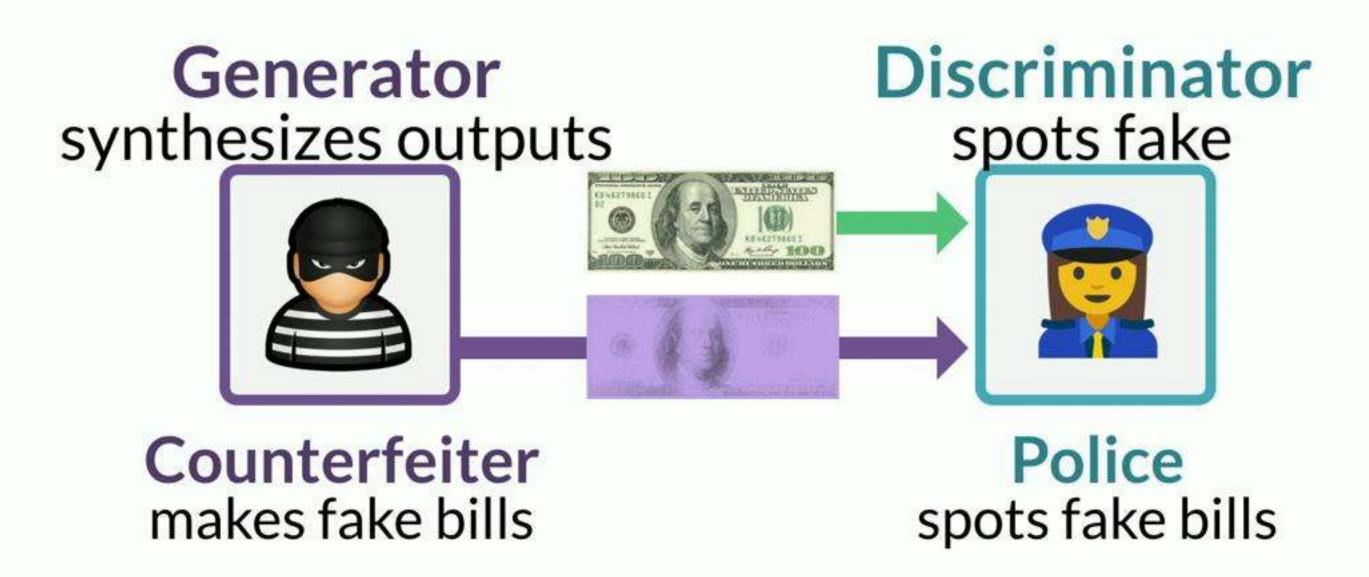


Police spots fake bills

A GAN uses two competing neural networks

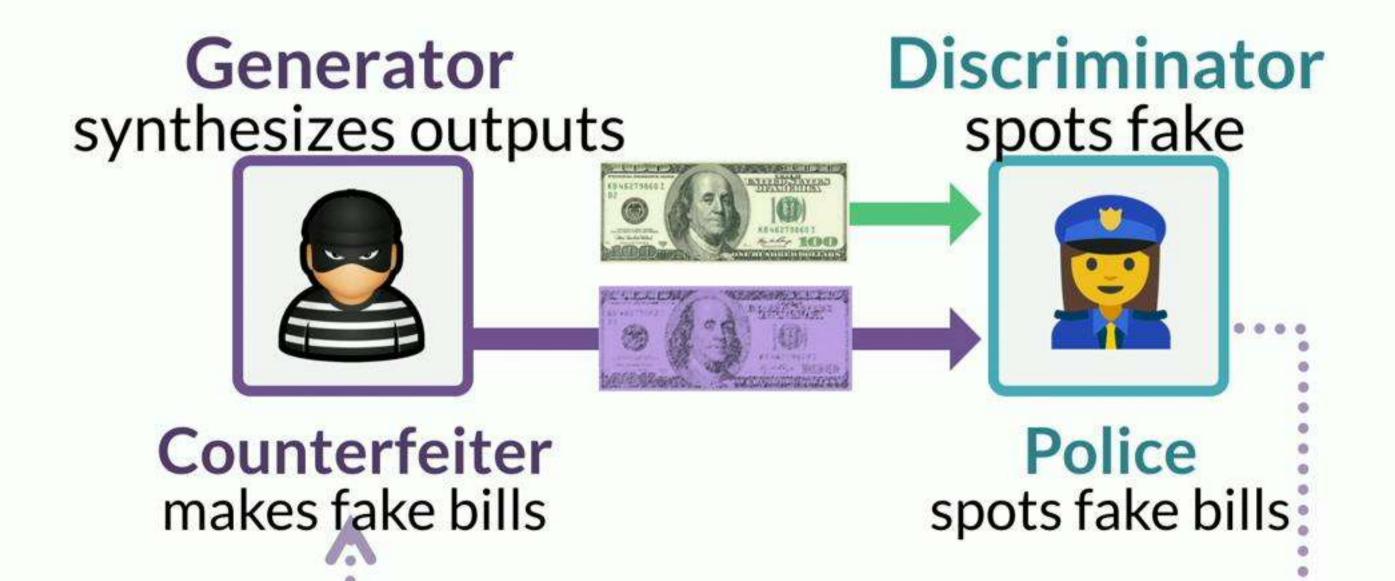


A GAN uses two competing neural networks



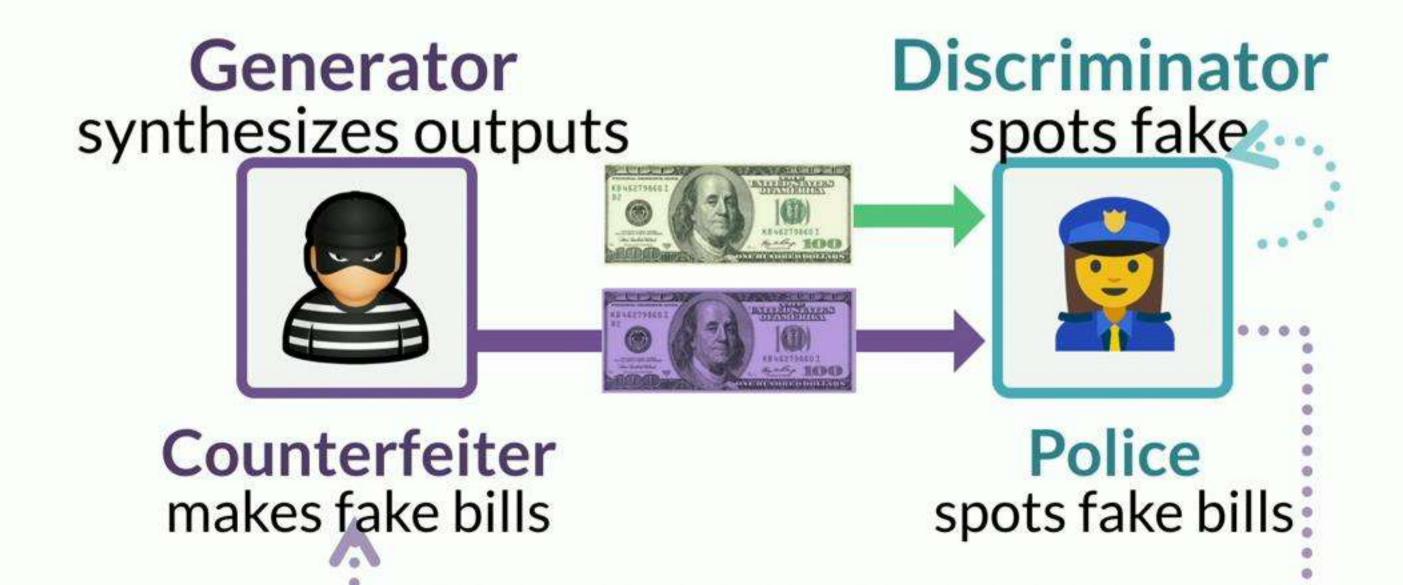
### Why GANs are hard?

A GAN uses two competing neural networks



# Why GANs are hard?

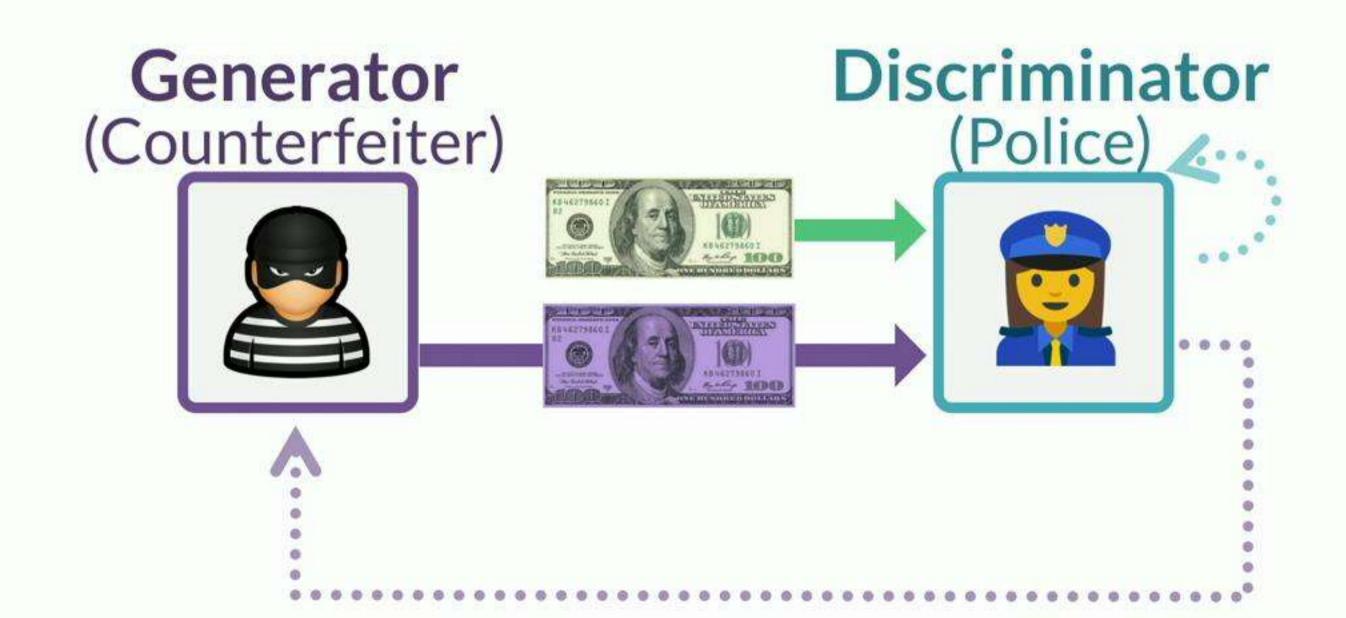
A GAN uses two competing neural networks



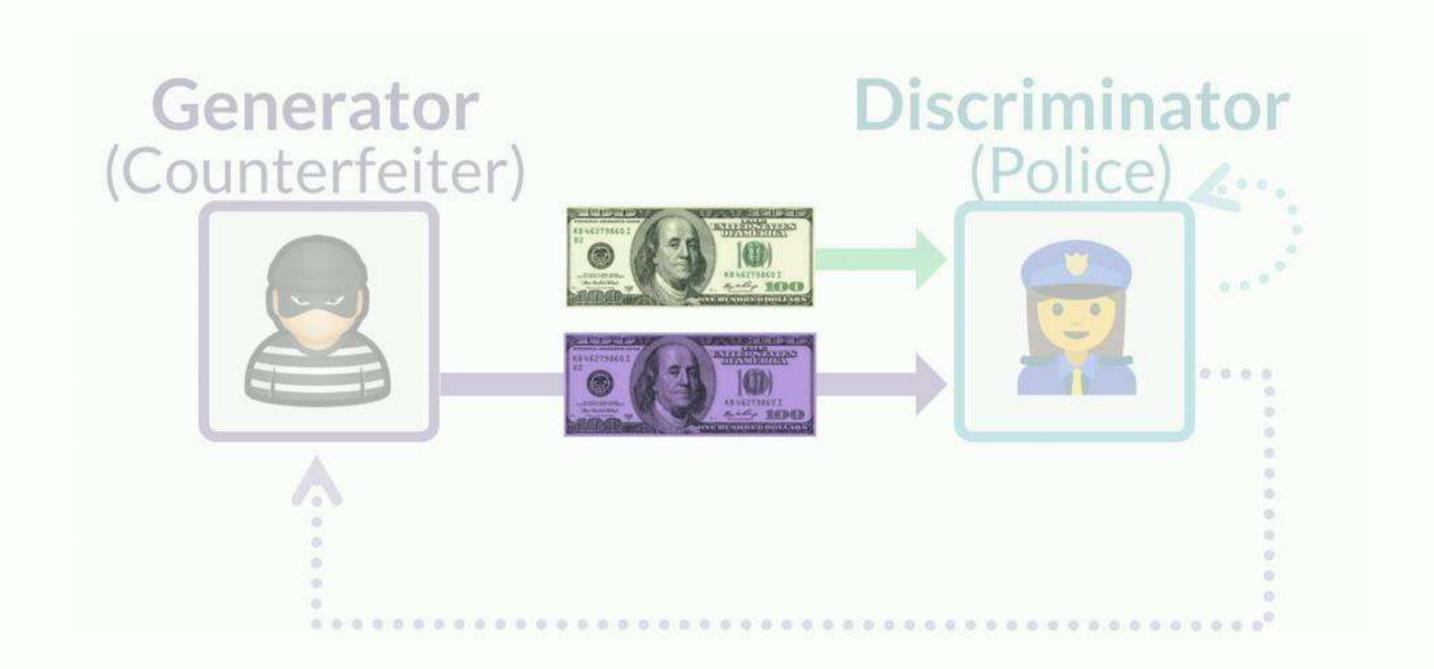
#### **GAN Lab Research Challenges**

#### Can we design an interactive tool for GANs?

- 1. Conceptual understanding of GANs
- 2. Interactive model training
- 3. Easily accessible by anyone

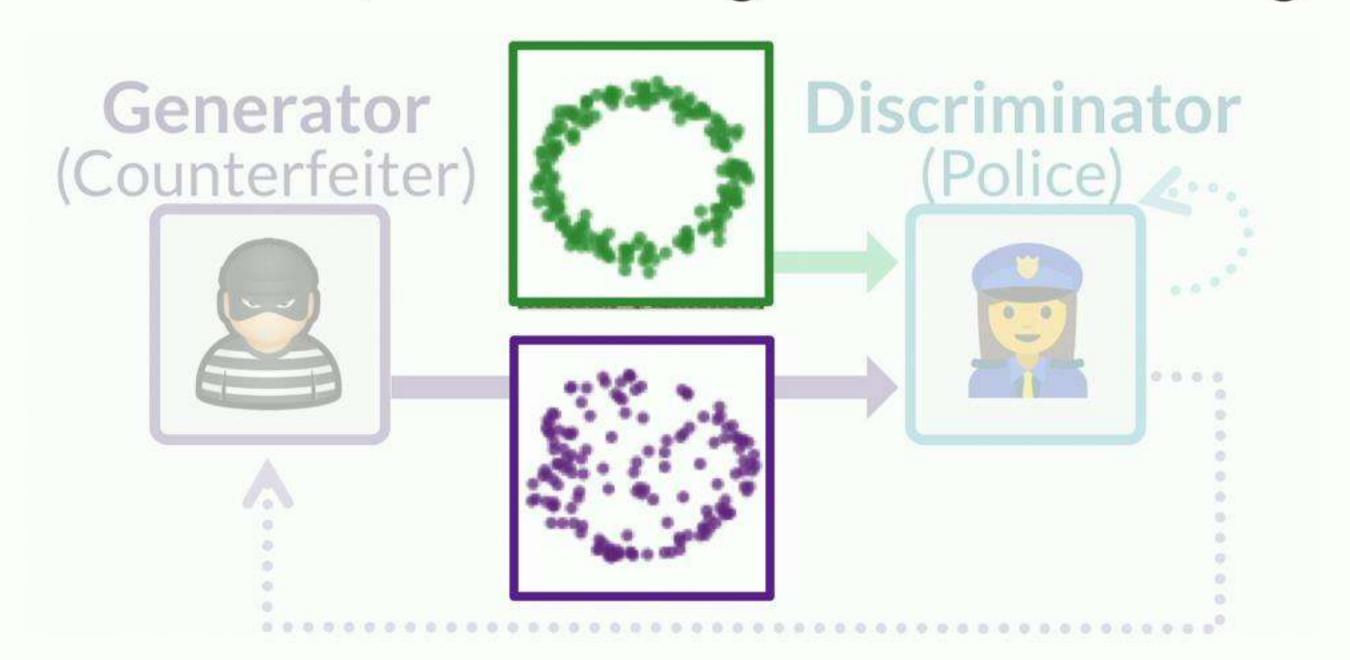


#### What type of data to visualize?



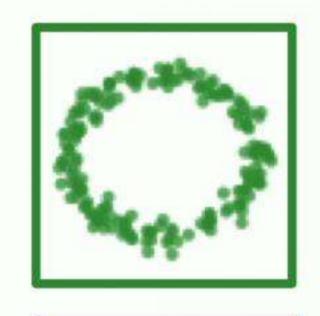
### What type of data to visualize?

2D distribution, instead of high-dimensional images



#### What type of data to visualize?

2D distribution, instead of high-dimensional images



#### Why 2D data points?

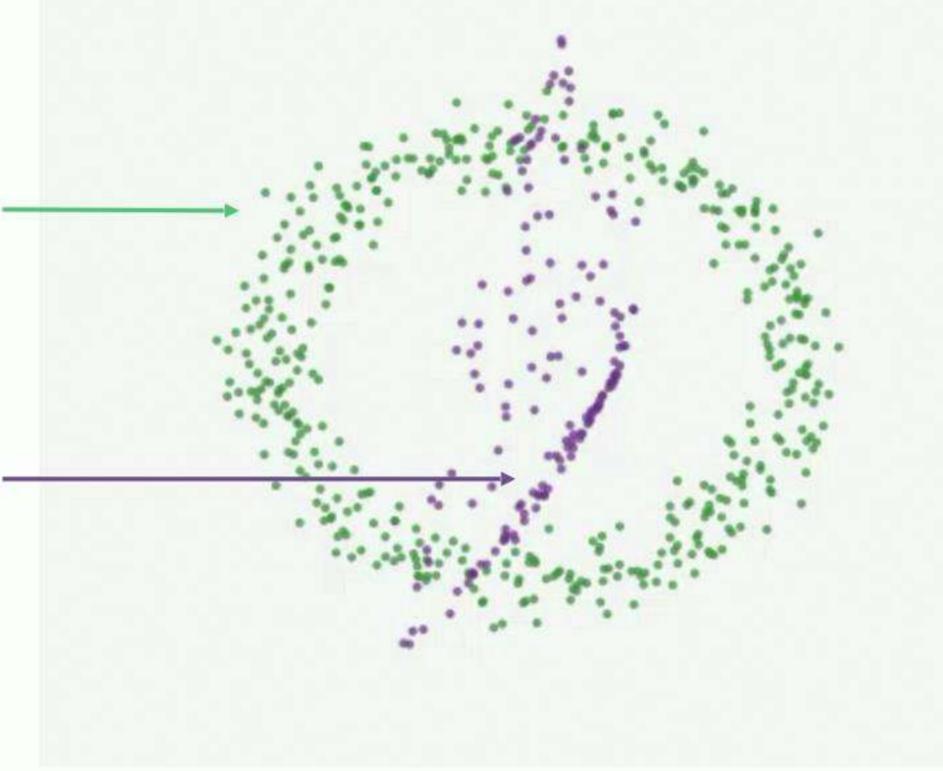
- 1. To focus on GAN's main concepts
- 2. To easily visualize data distribution

# 5 PI 000,000

#### VER. 0.1

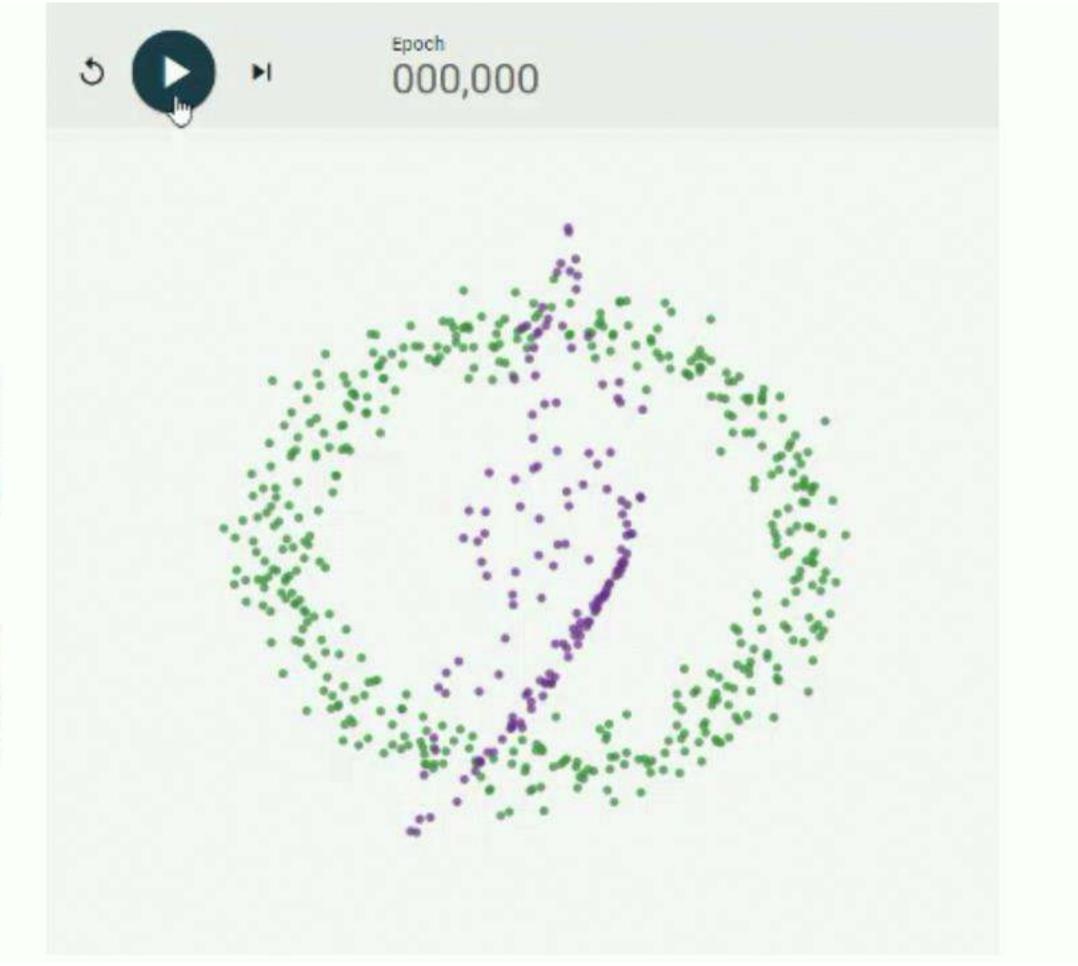


Generated (purple)



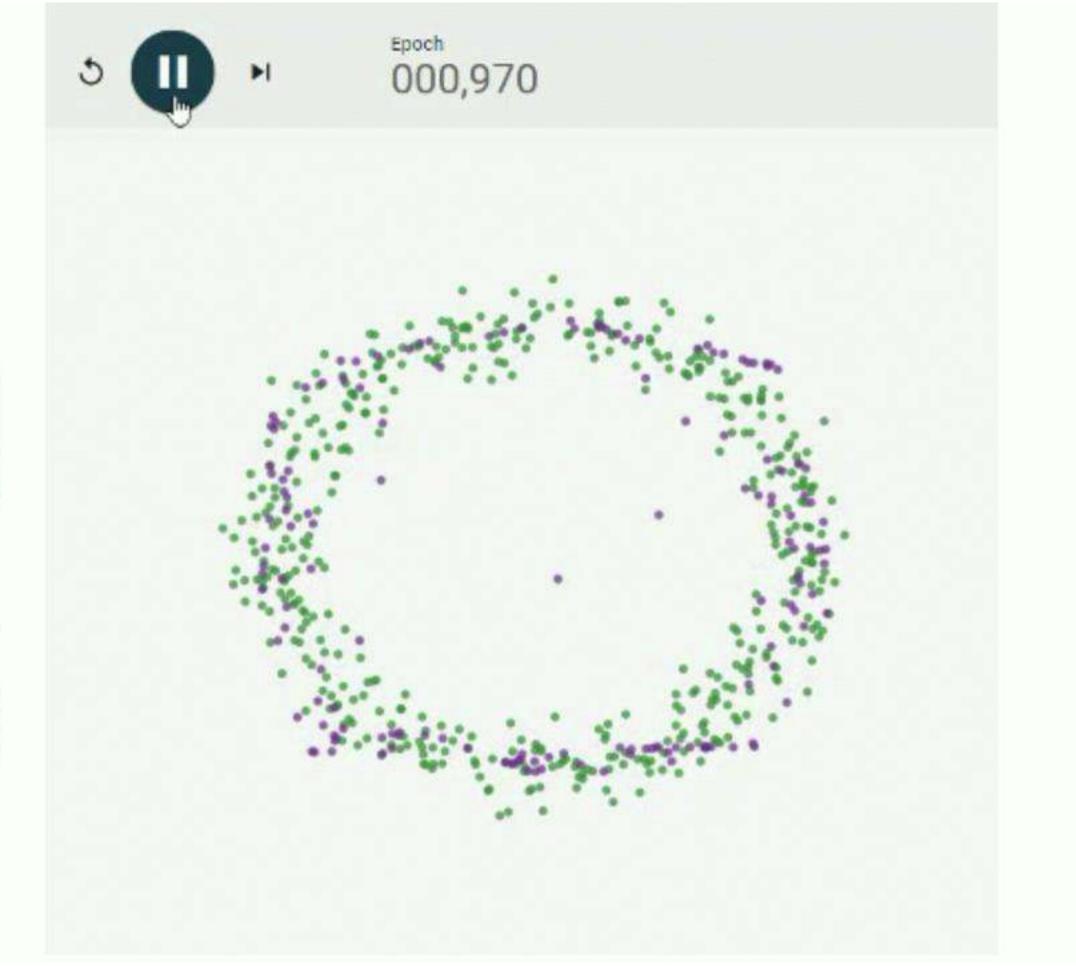
Real (green)

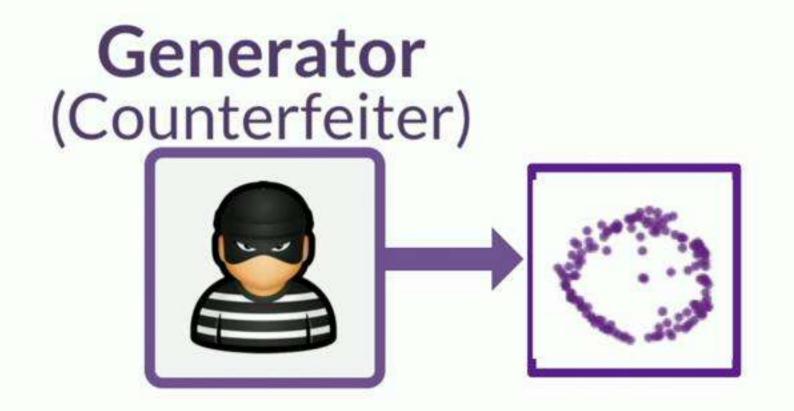
Generated (purple)

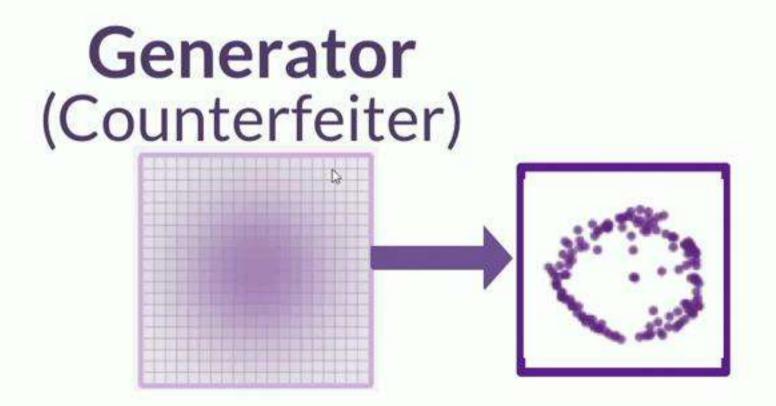


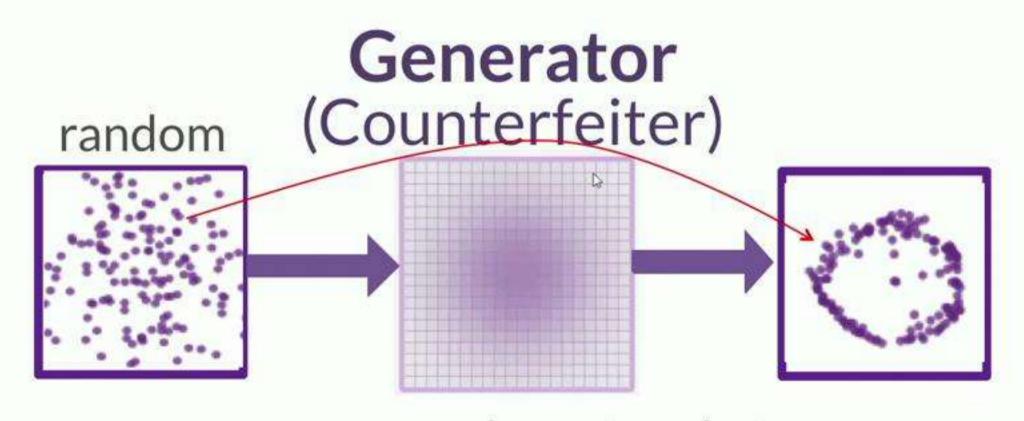
Real (green)

Generated (purple)

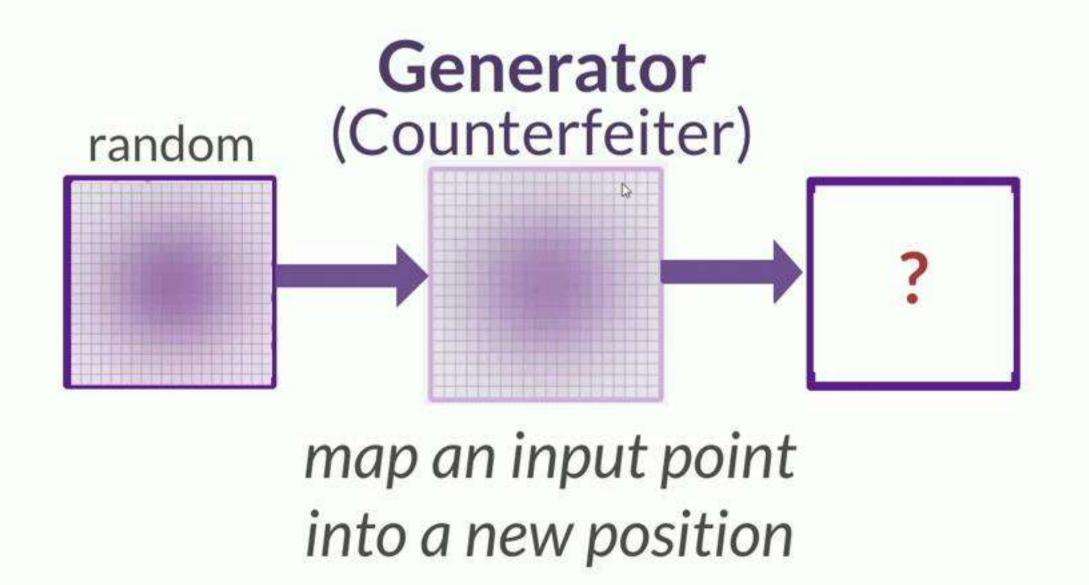


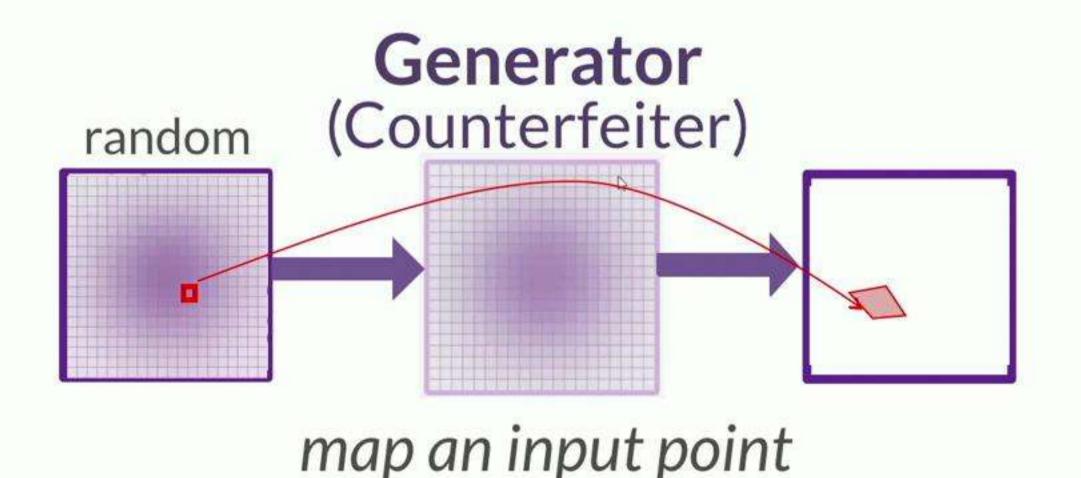




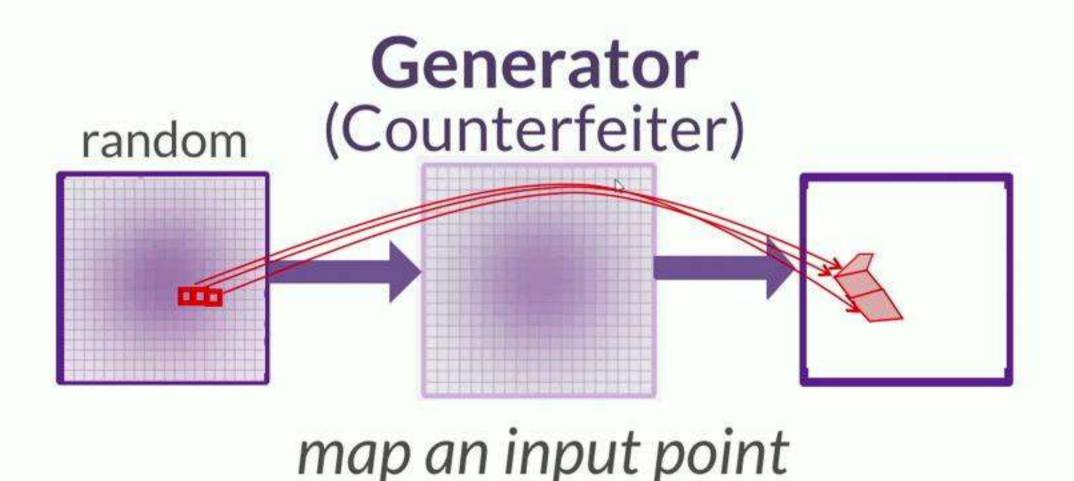


map an input point into a new position

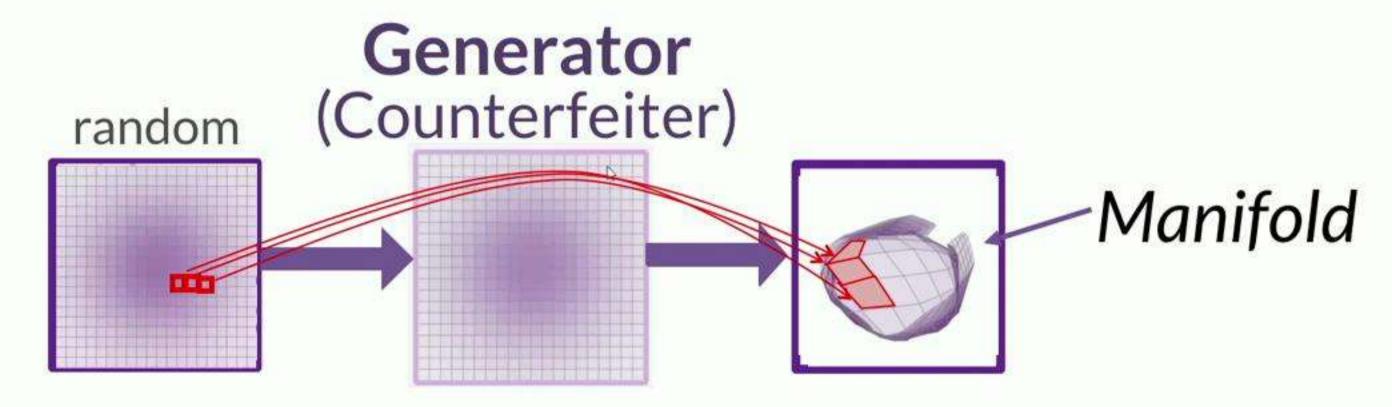




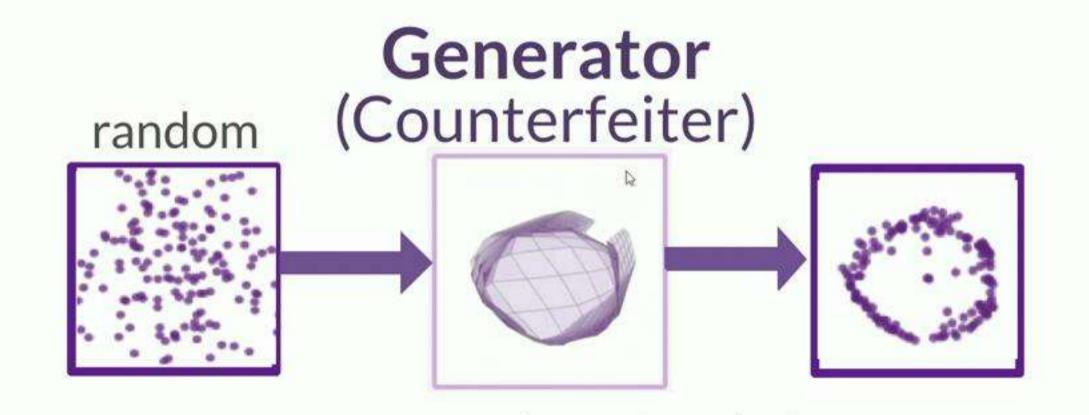
into a new position



into a new position



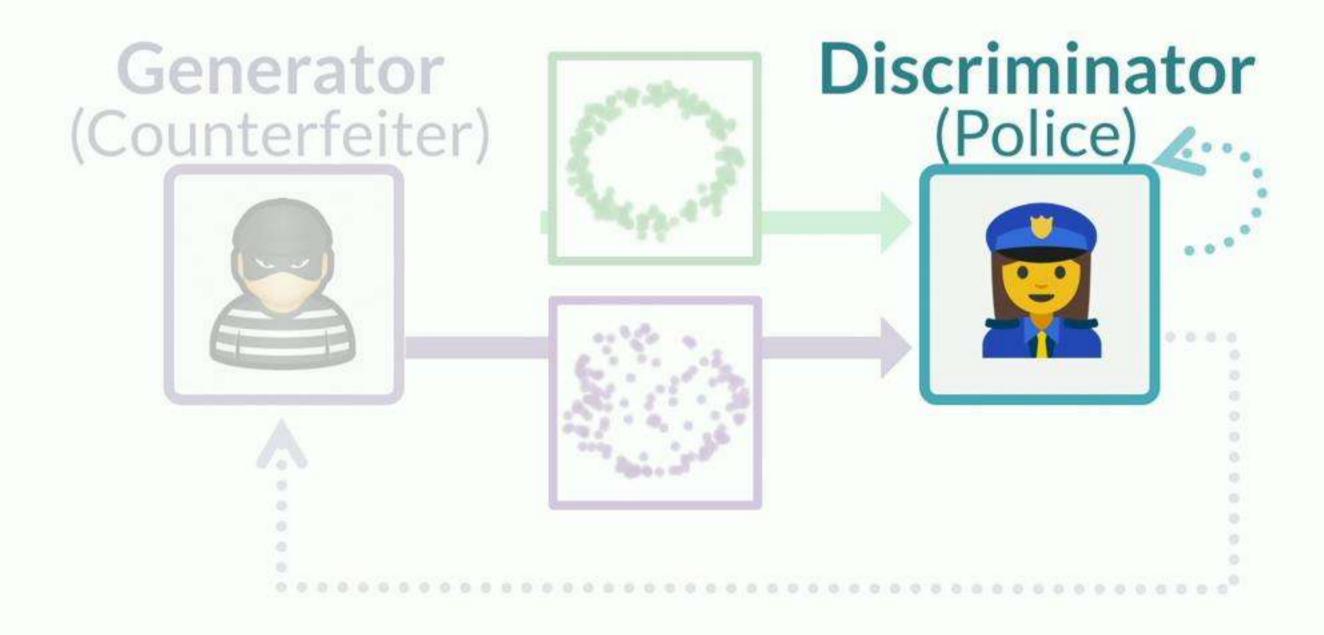
map an input point into a new position



map an input point

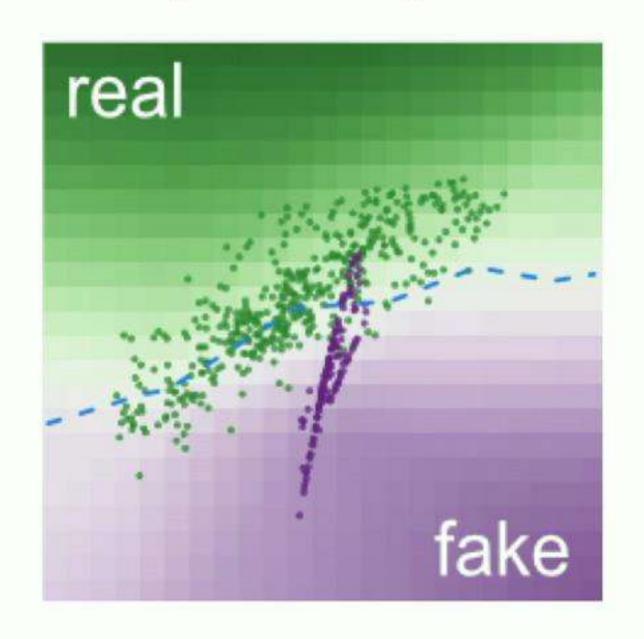
into a new position

#### How to visualize the discriminator?



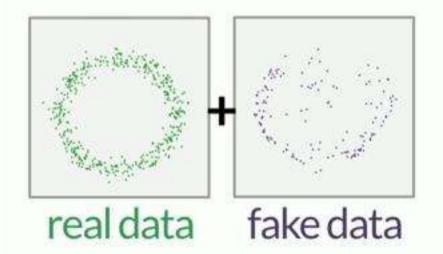
#### How to visualize the discriminator?

2D heatmap, to represent binary classification

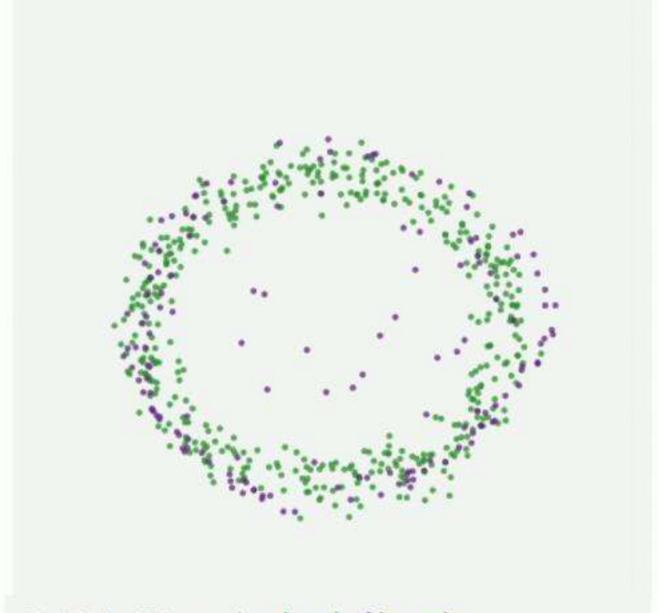


Data points in this region are likely real.

Data points are likely fake.



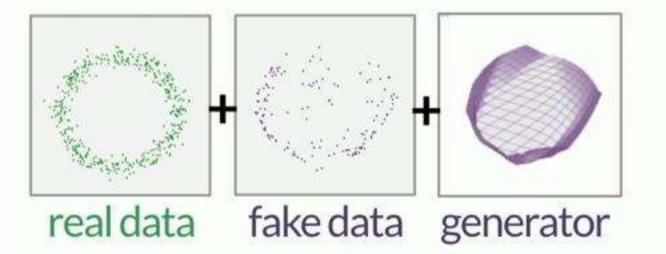




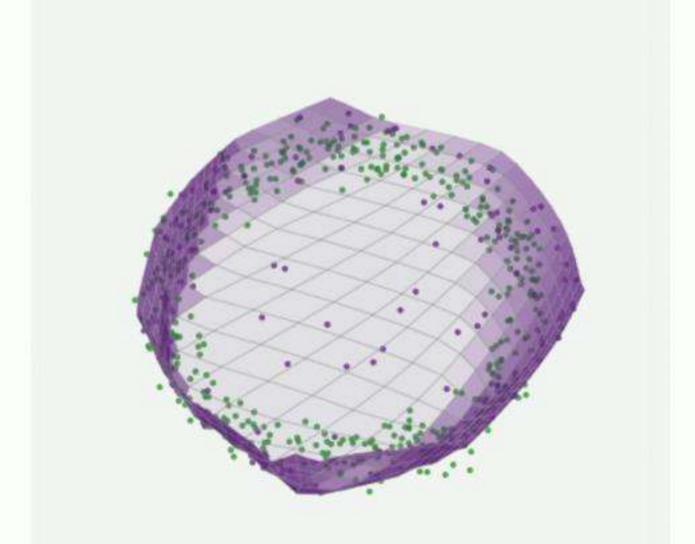
Each dot is a 2D data sample: real samples: fake samples.

Background colors of grid cells represent <u>discriminator</u>'s classifications.

Samples in green regions are likely to be real; those in purple regions likely fake.



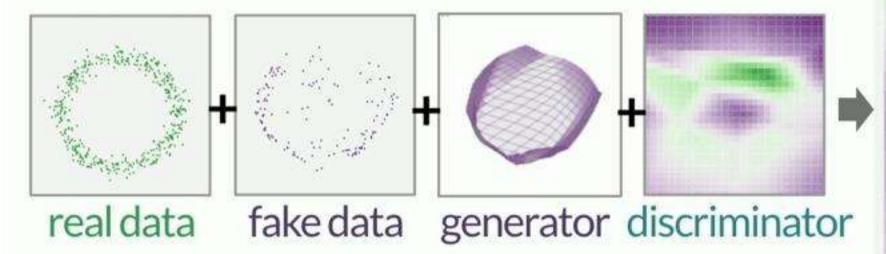




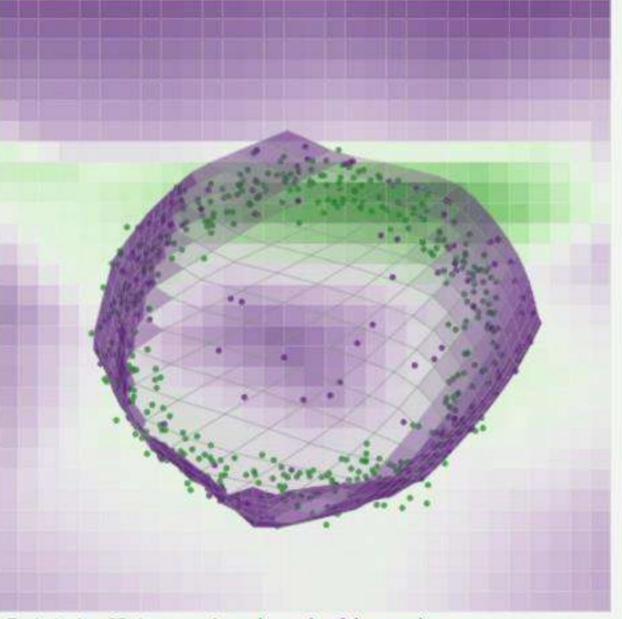
Each dot is a 2D data sample: real samples: fake samples.

Background colors of grid cells represent <u>discriminator</u>'s classifications.

Samples in green regions are likely to be real; those in purple regions likely fake.







Each dot is a 2D data sample: real samples; fake samples.

Background colors of grid cells represent <u>discriminator</u>'s classifications.

Samples in green regions are likely to be real; those in purple regions likely fake.

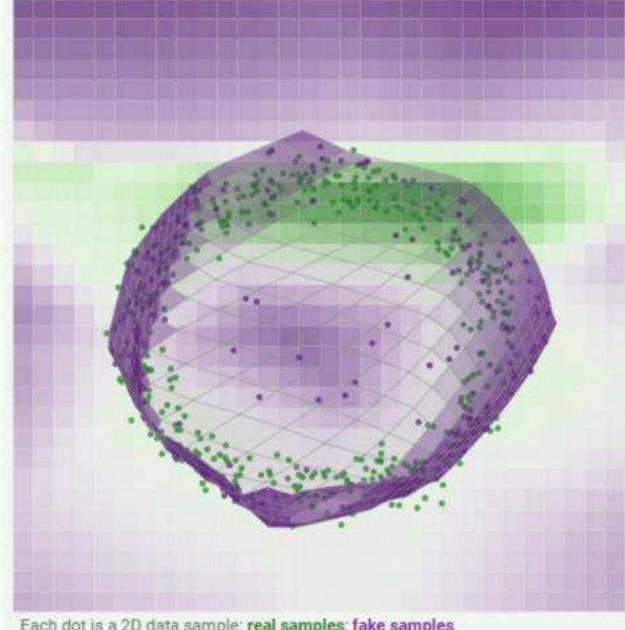


GENERATOR: # Hidden Layers: 2 # Neurons in Hidden Layer: 10

DISCRIMINATOR: # Hidden Layers: 2

# Neurons in Hidden Layer: 10

HYPERPARAMETERS: More Steps for Discriminator: 2



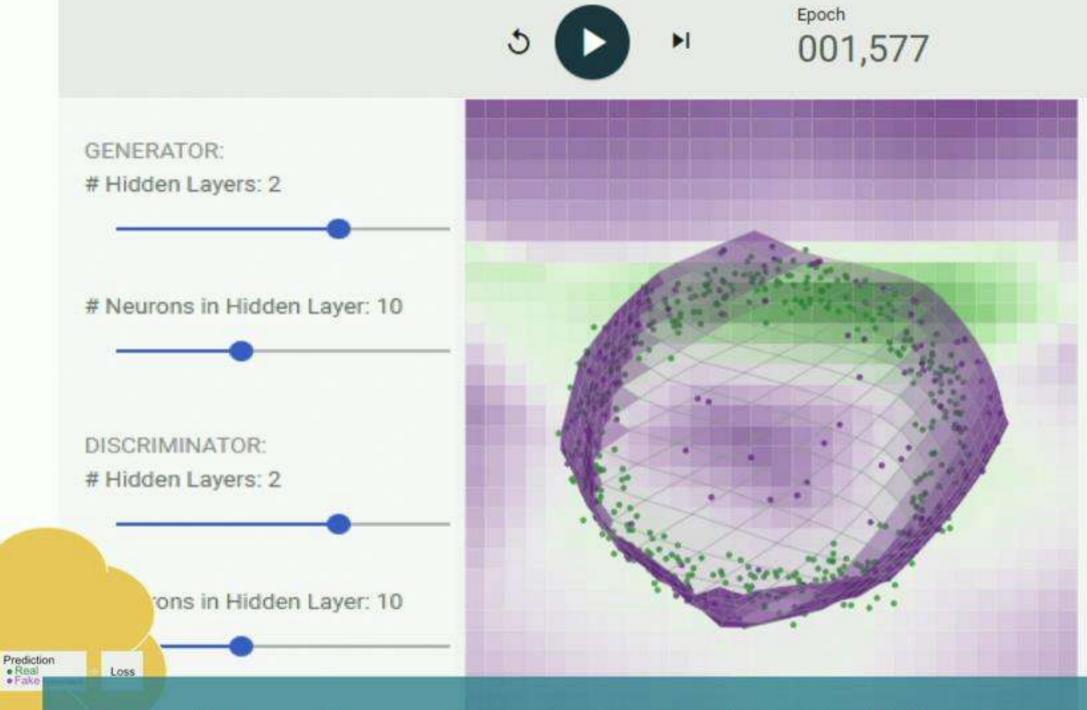
Each dot is a 2D data sample: real samples: fake samples.

Background colors of grid cells represent discriminator's classifications. Samples in green regions are likely to be real; those in purple regions likely fake.

Sample

Sample

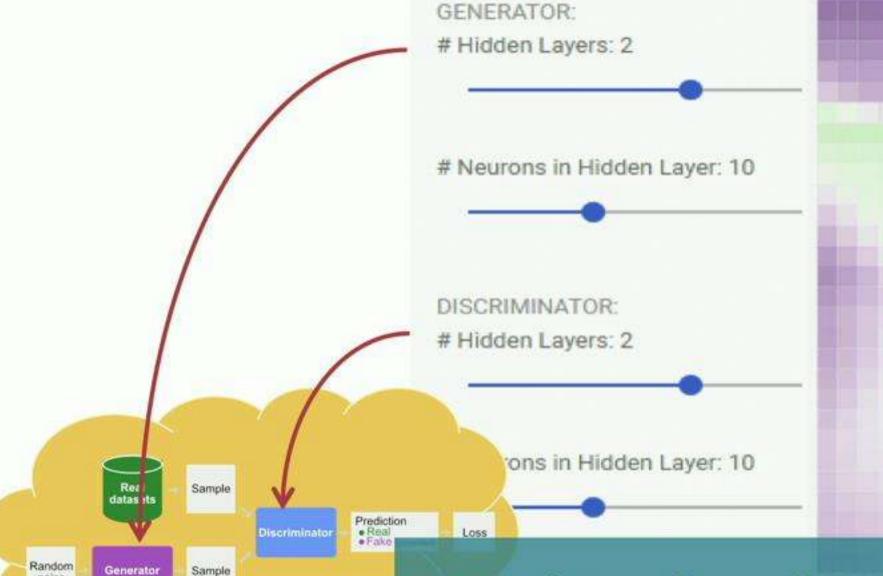
Random

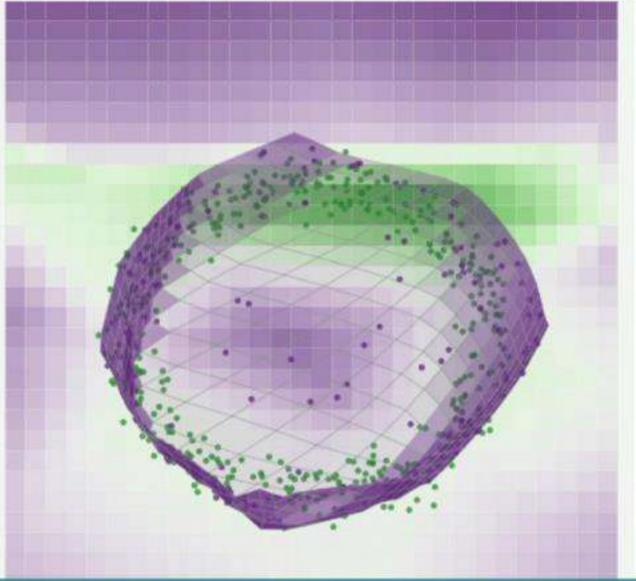


#### Hard to develop mental models for GANs



001,577





Hard to develop mental models for GANs

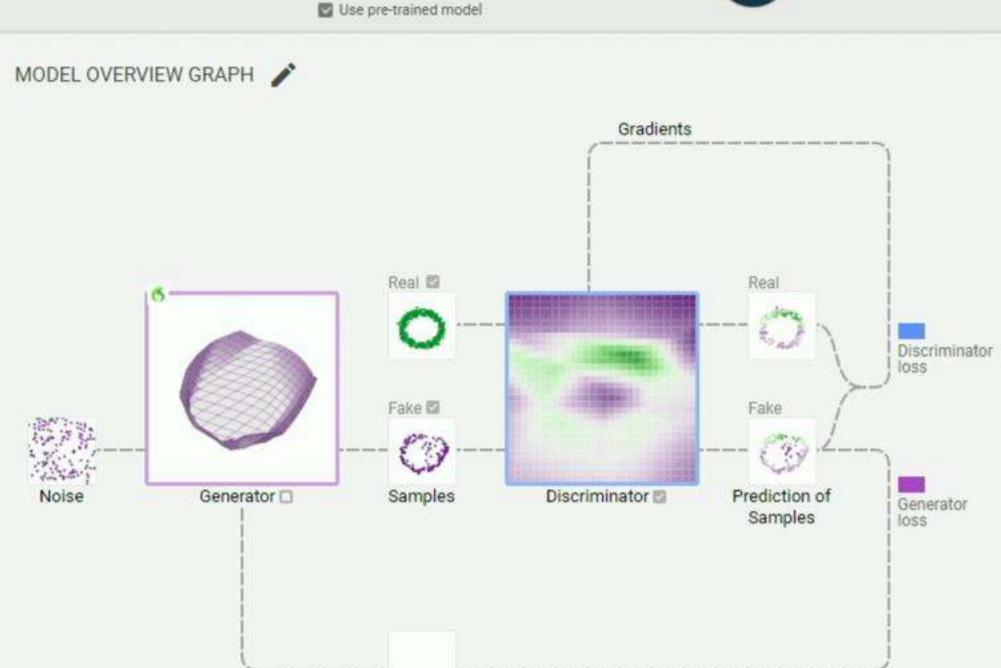






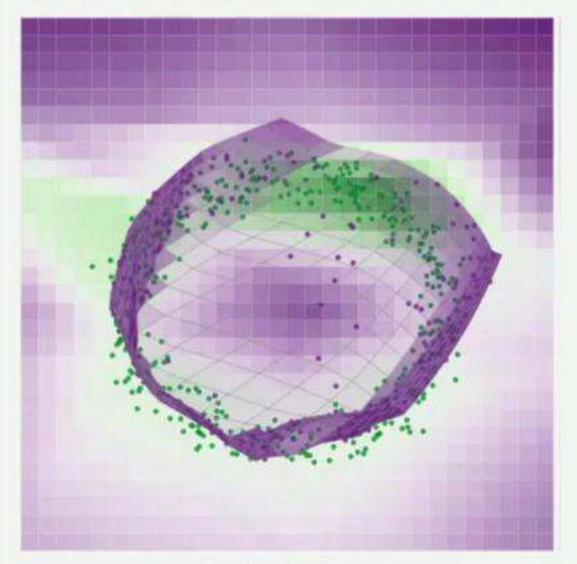


001,577



Gradients

#### LAYERED DISTRIBUTIONS



Each dot is a 2D data sample: real samples; fake samples.

Background colors of grid cells represent <u>discriminator</u>'s classifications.

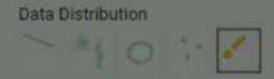
Samples in green regions are likely to be real, those in purple regions likely fake.

Manifold represents generator's transformation results from noise space.

Opacity encodes density: darker purple means more samples in smaller area.

Pink lines from fake samples represent gradients for generator.

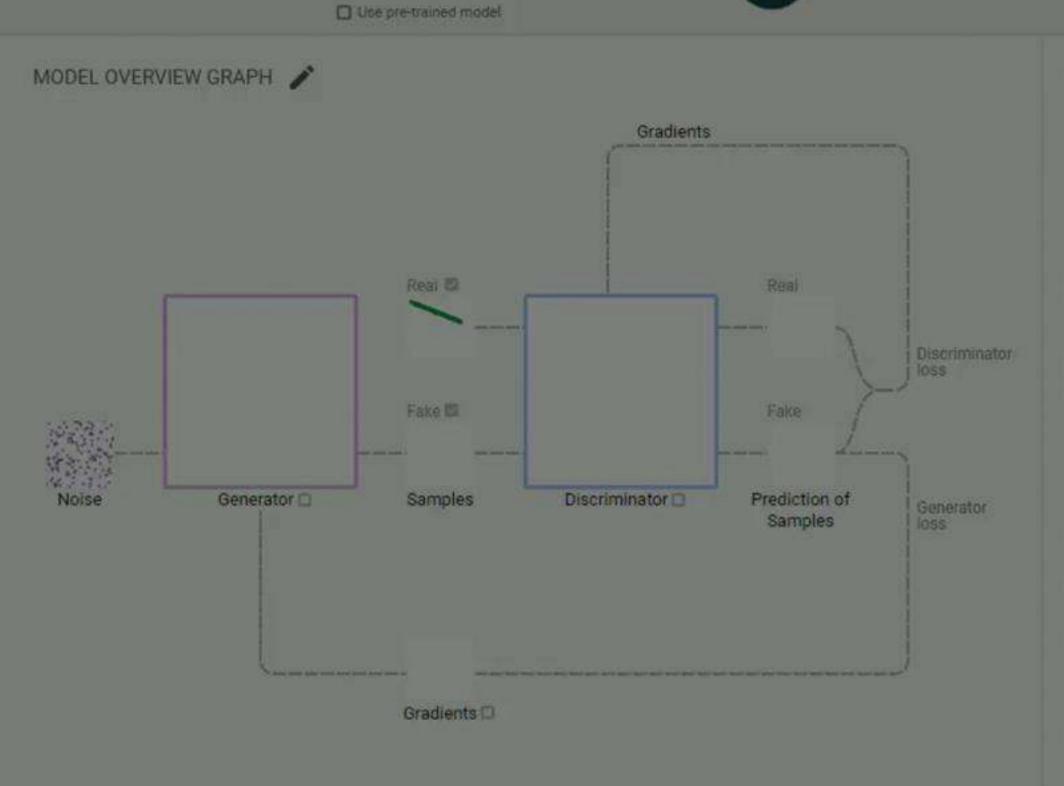
✓ This sample needs to move upper right to decrease generator's loss.







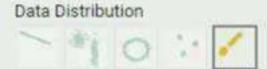
Epoch 000,000



#### LAYERED DISTRIBUTIONS

Draw a distribution above, then click the apply button.

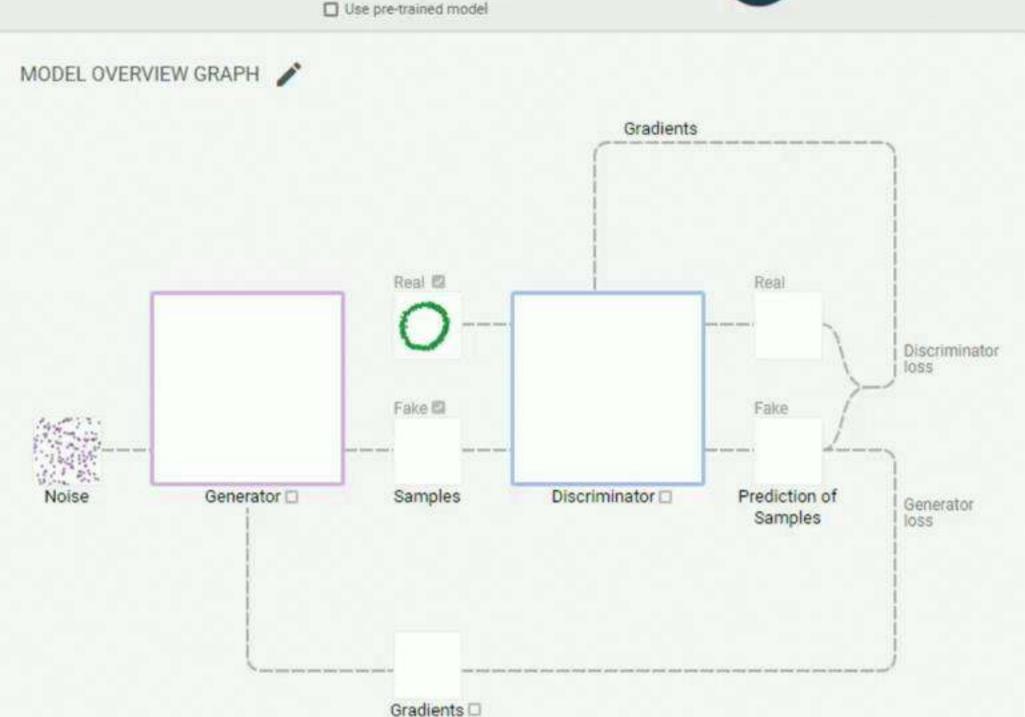
APPLY







000,000



#### LAYERED DISTRIBUTIONS



Each dot is a 2D data sample real samples, fake samples.

Background holors of grid cells represent discriminator's classifications.

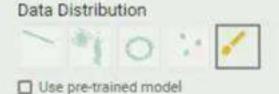
Samples in seen regions are likely to be real; those in purple regions likely fake.

Manifold represents generator's transformation results from noise space.

Opacity encodes density: darker purple means more samples in smaller area.

Pink lines from fake samples represent gradients for generator.

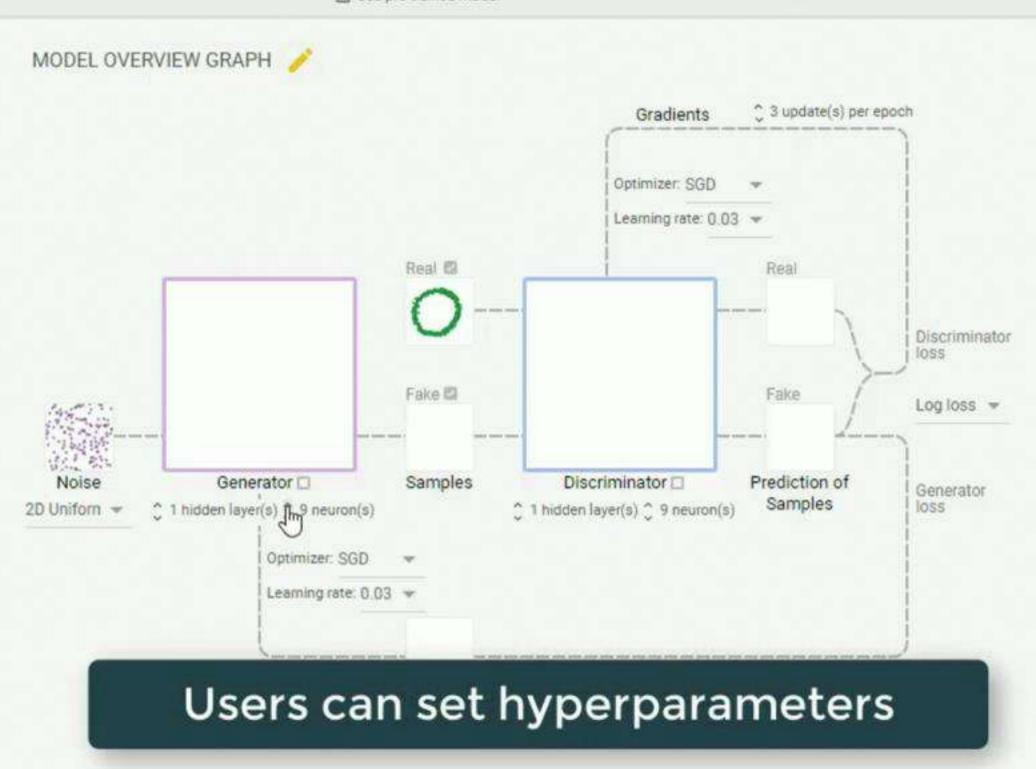
This sample needs to move upper right to decrease generator's loss.



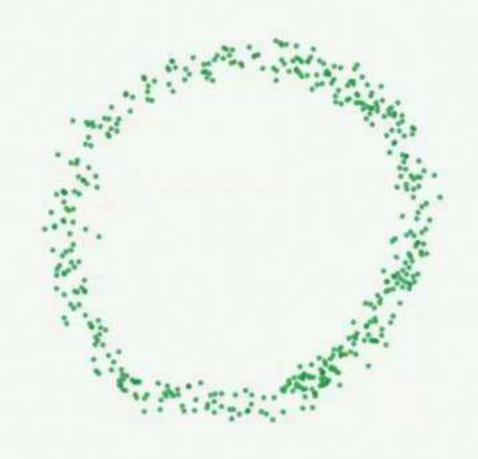




Epoch 000,000



#### LAYERED DISTRIBUTIONS



Each dot is a 2D data sample: real samples, fake samples.

Background colors of grid cells represent **discriminator**'s classifications.

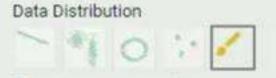
Samples in green regions are likely to be real; those in purple regions likely fake.

Manifold represents generator's transformation results from noise space.

Opacity encodes density: darker purple means more samples in smaller area.

Pink lines from fake samples represent **gradients** for generator.

This sample needs to move upper right to decrease generator's loss.

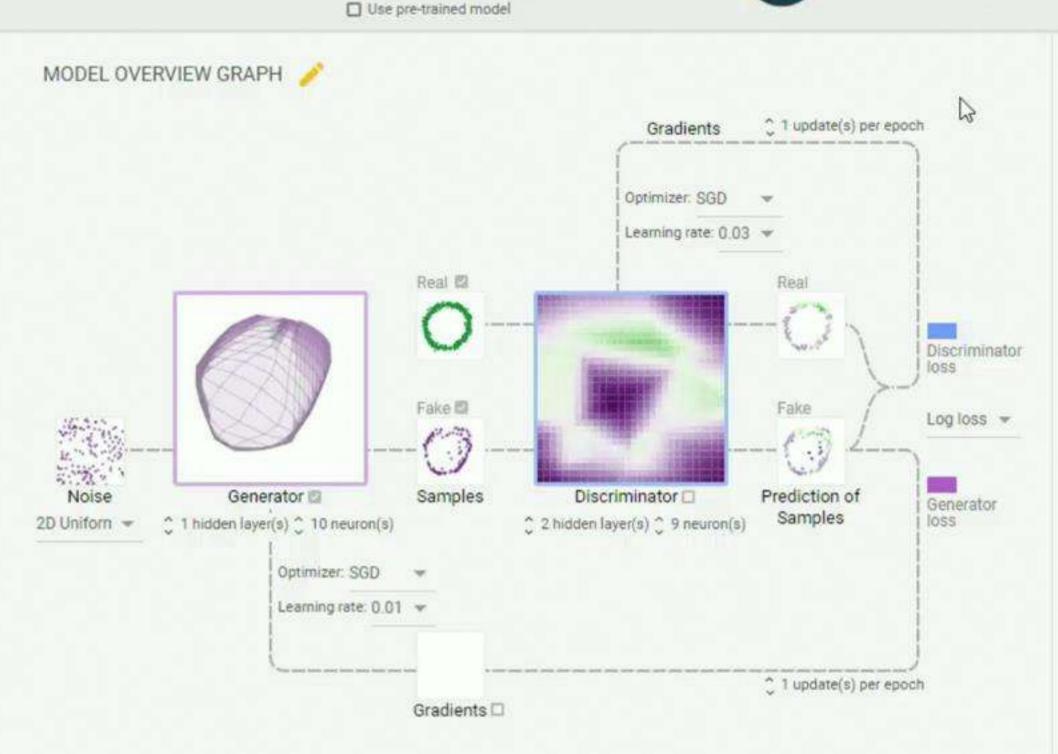




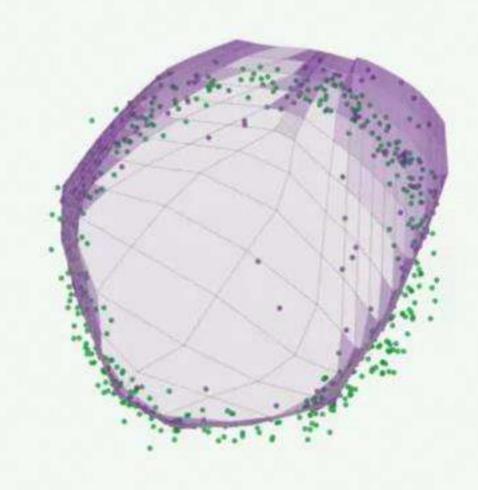




Epoch 002,824



#### LAYERED DISTRIBUTIONS



Each dot is a 2D data sample: real samples fake samples

Background colors of grid cells represent **discriminator**'s classifications.

Samples in green regions are likely to be real; those in purple regions likely fake.

Manifold represents generator's transformation results from noise space.

Opacity encodes density: darker purple means more samples in smaller area.

Pink lines from fake samples represent **gradients** for generator.

This sample needs to move upper right to decrease generator's loss.

Playing at 15x speed

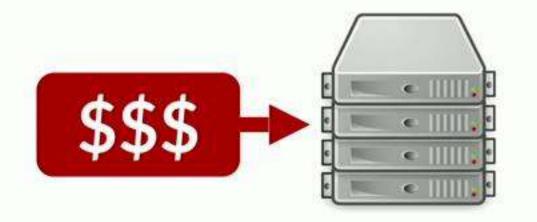
#### GAN Lab broadens education access

Conventional Deep Learning Visualization



Visualization in JavaScript





Model Training in Python with GPU

#### GAN Lab broadens education access

Everything done in browser, powered by TensorFlow.js



Visualization in JavaScript

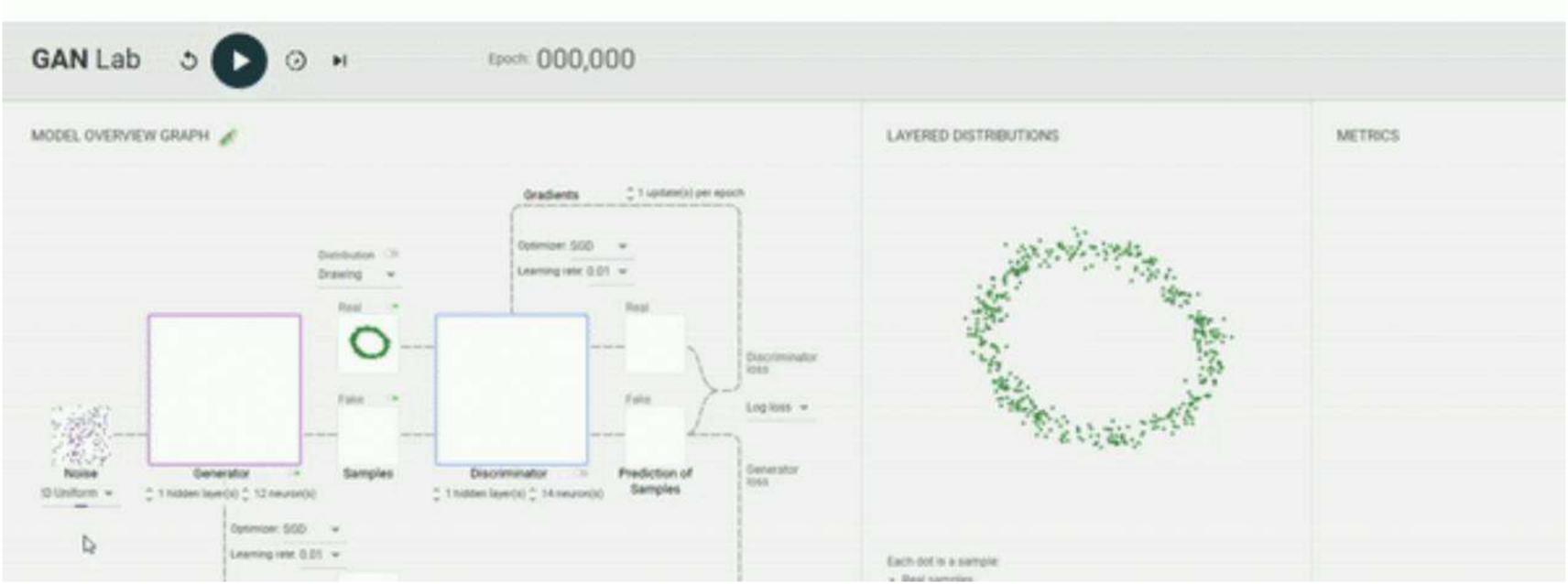
Model Training also in JavaScript

Accelerated by WebGL

#### GAN Lab is Live! Try at bit.ly/gan-lab

30K visitors, 135 countries

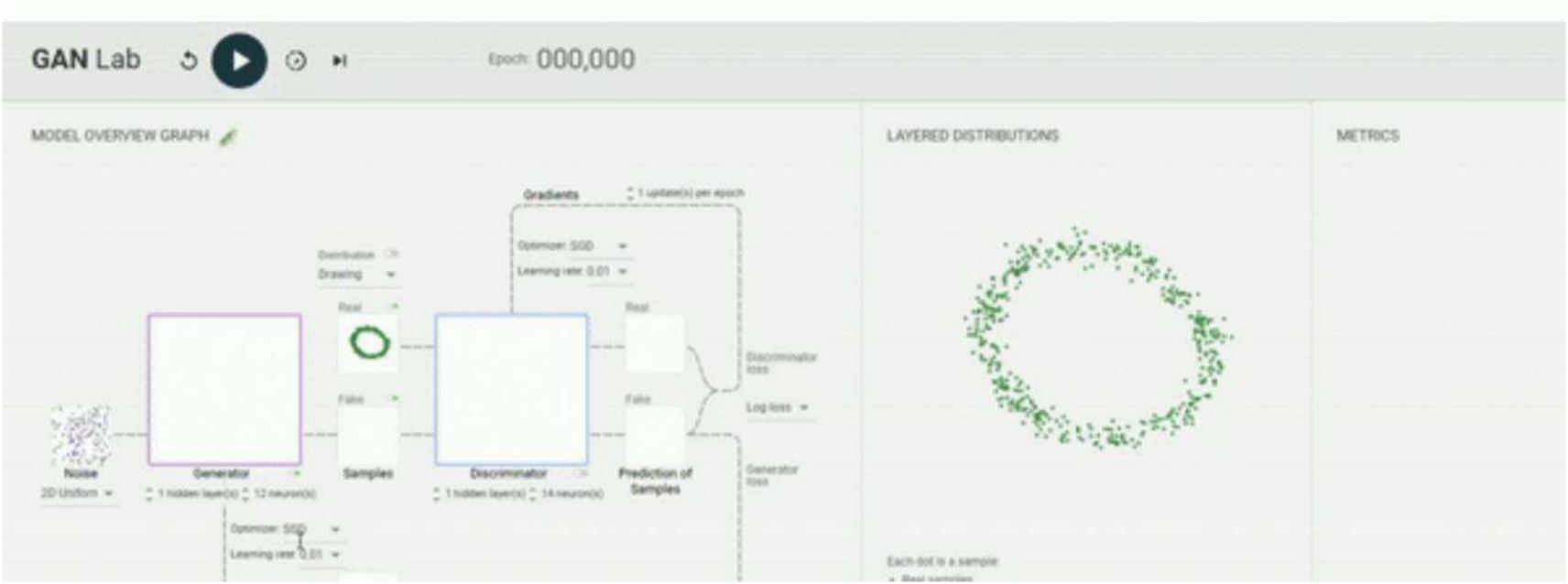
♥ 1.9K Likes 13800+ Retweets



#### GAN Lab is Live! Try at bit.ly/gan-lab

30K visitors, 135 countries

♥ 1.9K Likes 13800+ Retweets



# nterpretable AI via Visual Analytics

Understanding Industry-Scale Models

ActiVis - Activation analysis by subsets

Interactive Learning of Complex Models

GAN Lab - Experimentation with GANs

Research Landscape

Survey, Gamut

# Visual Analytics in Deep Learning An Interrogative Survey for the Next Frontiers



Fred Hohman Georgia Tech



Minsuk Kahng Georgia Tech



Robert Pienta Symantec



Polo Chau Georgia Tech





### Visual Analytics in Deep Learning

#### Interrogative Survey Overview



Why would one want to use visualization in deep learning?

Interpretability & Explainability
Debugging & Improving Models
Comparing & Selecting Models
Teaching Deep Learning Concepts



What data, features, and relationships in deep learning can be visualized?

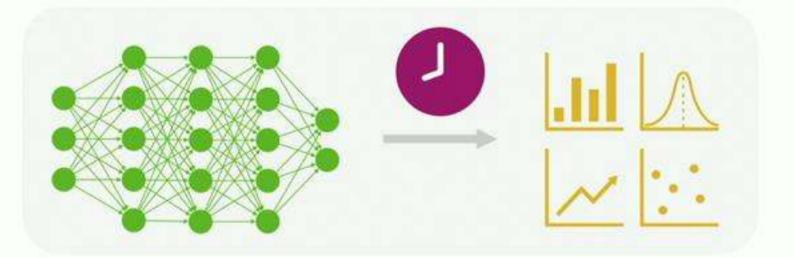
Computational Graph & Network Architecture Learned Model Parameters Individual Computational Units Neurons In High-dimensional Space Aggregated Information



When in the deep learning process is visualization used?

During Training After Training









Who would use and benefit from visualizing deep learning?

Model Developers & Builders Model Users Non-experts

#### §7 HOW

How can we visualize deep learning data, features, and relationships?

Node-link Diagrams for Network Architecture Dimensionality Reduction & Scatter Plots Line Charts for Temporal Metrics Instance-based Analysis & Exploration Interactive Experimentation



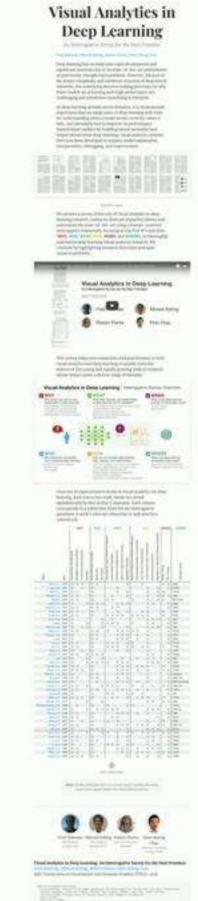
Where has deep learning visualization been used?

Application Domains & Models A Vibrant Research Community

## **Key Takeaways**

## bit.ly/va-dl-survey

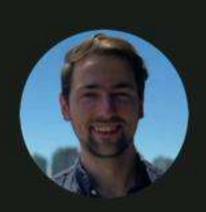
- 1. Most tools aimed at expert users
- 2. Instance-based analysis
- 3. Inherently interdisciplinary
- 4. Lacks actionability
- 5. Evaluation is hard
- 6. State-of-the-art models not robust



# Gamut

A Design Probe to Understand How Data Scientists Understand Machine Learning Models

CHI 2019



Fred Hohman

@fredhohman

Georgia Tech



Andrew Head UC Berkeley



Rob DeLine Microsoft Research



Rich Caruana Microsoft Research



Steven Drucker Microsoft Research

## What is interpretability?

internals e.g., components [Gilpin, 2018] operations Human understanding e.g., math [Biran, 2017] of a system's... data mapping e.g., input to output [Montavon, 2017] representation in an explanation [Ribeiro, 2016]

## What is interpretability?

**Human understanding** of a system's...

No formal, agreed upon definition [Lipton, 2016]

internals e.g., components [Gilpin, 2018]

operations e.g., math [Biran, 2017]

data mapping e.g., input to output [Montavon, 2017]

representation in an explanation [Ribeiro, 2016]

## **Gamut Contributions**

- 1. Capabilities of interpretability
- 2. Design Probe embodying capabilities
- 3. Evaluation & Investigation of probe & emerging practice of interpretability w/real users



#### From formative research

## Explainable ML Interface Questions

Why does this house cost that much?

What is the difference between these two?

What if I added...

What are similar homes?

Where is it wrong?

What is most important?

## Explainable ML Interface Capabilities

Why does this house cost that much?

What is the difference between these two?

What if I added...

What are similar homes?

Where is it wrong?

What is most important?

#### From formative research

## Explainable ML Interface Capabilities

- Why does this house cost that much?

  Local instance explanations
- What is the difference between these two?
  Instance explanation comparisons
- What if I added...
  Counterfactuals
- What are similar homes?
  Nearest neighbors
- Where is it wrong?
  Regions of error
- What is most important?
  Feature importance

#### From formative research

## Explainable ML Interface Capabilities



Why does this house cost that much?

Local instance explanations



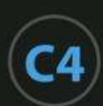
What is the difference between these two?

Instance explanation comparisons



What if I added...

Counterfactuals



What are similar homes?

Nearest neighbors



Where is it wrong?

Regions of error



What is most important?

Feature importance

#### GAMUT: A Design Probe to Understand How Data Scientists Understand Machine Learning Models

Atlanta, GA, USA nefleckmana/typicals eds.

#### UC Blakeley Berkeley, CA, USA anders hout 2 herbeley orbi

Microsoft Research Biolimond, WA, USA

#### Microsoft Research Redmond, WA, USA.

Without good models and the right tools to interprit them. data scientists risk making decisions based on hidden bisses. no correlations, and false presentinations. This has led to a rallying cry for model interpretability. Yet the concept of interpretability remains rebuleus, each that rewarders and total dealgrants lack actionable qualifilines for how to inconstrainty interpretability into models and accompanying tools. Through an iterative design process with expert me risual analytics restore. Casery, to explore how interactive nertices could better support model interpretation. Using out't us a peobe, we investigated why and how profesocial data retreation interpret models, and how introduce alfunction that support data is livelists in assessing specificon

#### Birdmond, WA, USA

Microsoft Research

· Human centered computing -- Empirical studies in visualization. Visualization systems and tools - Compating methodologies -- Machine learning

Machine learning interpretability, design peole, visual ana bytics, data examination, interactive interfaces

Total Histories, Andrew Hond, Ruth Campana, Robert Dell, sec. und Steve N. Drucker, 2017 Green: A Drugo Parks to Understand New Data Incontast Understand Hardon Learning Models. In CHI Conference or Names Factors in Computing Systems Proceedings (CAF 3019). May 8-9, 2018. Glasgow, Southead LW, ACM, New York, NY, USA,

Definitions + examples in the paper!

## Takeaways

- Consider interpretability capabilities for your interfaces Interpretability is not a singular, rigid concept
- Tailor explanations for specific audiences

  Balance simplicity and completeness
- Design and integrate effective interaction

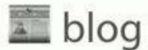
Interaction key to realizing interpretability & solidify model understanding [Weld & Bansal, 2018]

Gamut A Design Probe to
Understand How Data Scientists
Understand Machine Learning Models

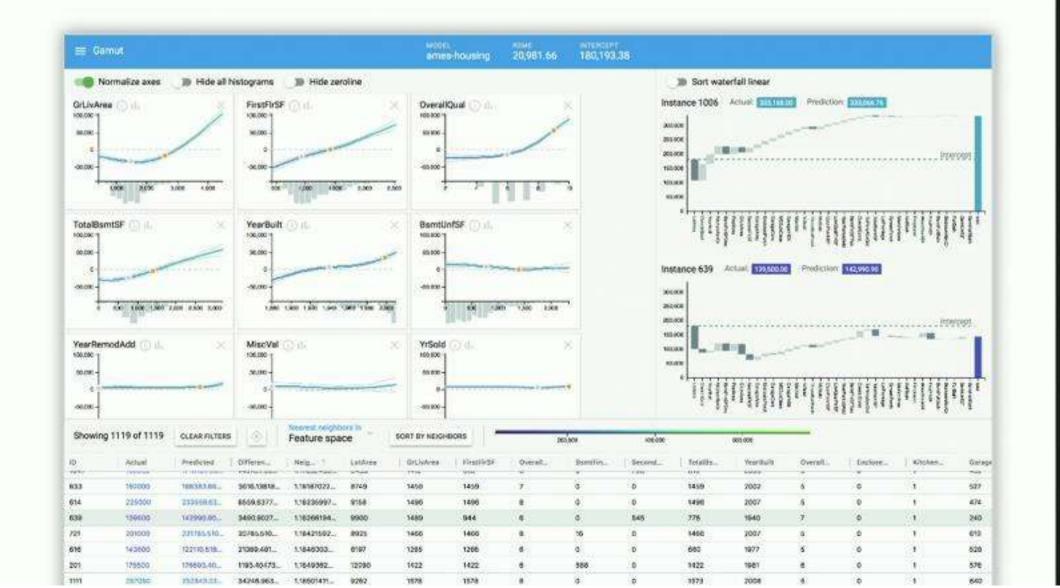
## bit.ly/gamut-chi













**Fred Hohman** @fredhohman Georgia Tech



**Andrew Head UC** Berkeley



**Rich Caruana** Microsoft Research



**Rob DeLine** Microsoft Research



Steven Drucker Microsoft Research





Berkeley Research

## Secure

# Interpretable

Attack & Defense (DNN)

ShapeShifter Shield

Do-it-yourself Adversarial ML

ADAGIO MLsploit Understand Industry Models
ActiVis

**ML** Education

**GAN** Lab

Research landscape

Survey, Gamut

#### TOWARDS



# SECURE & INTERPRETABLE AI

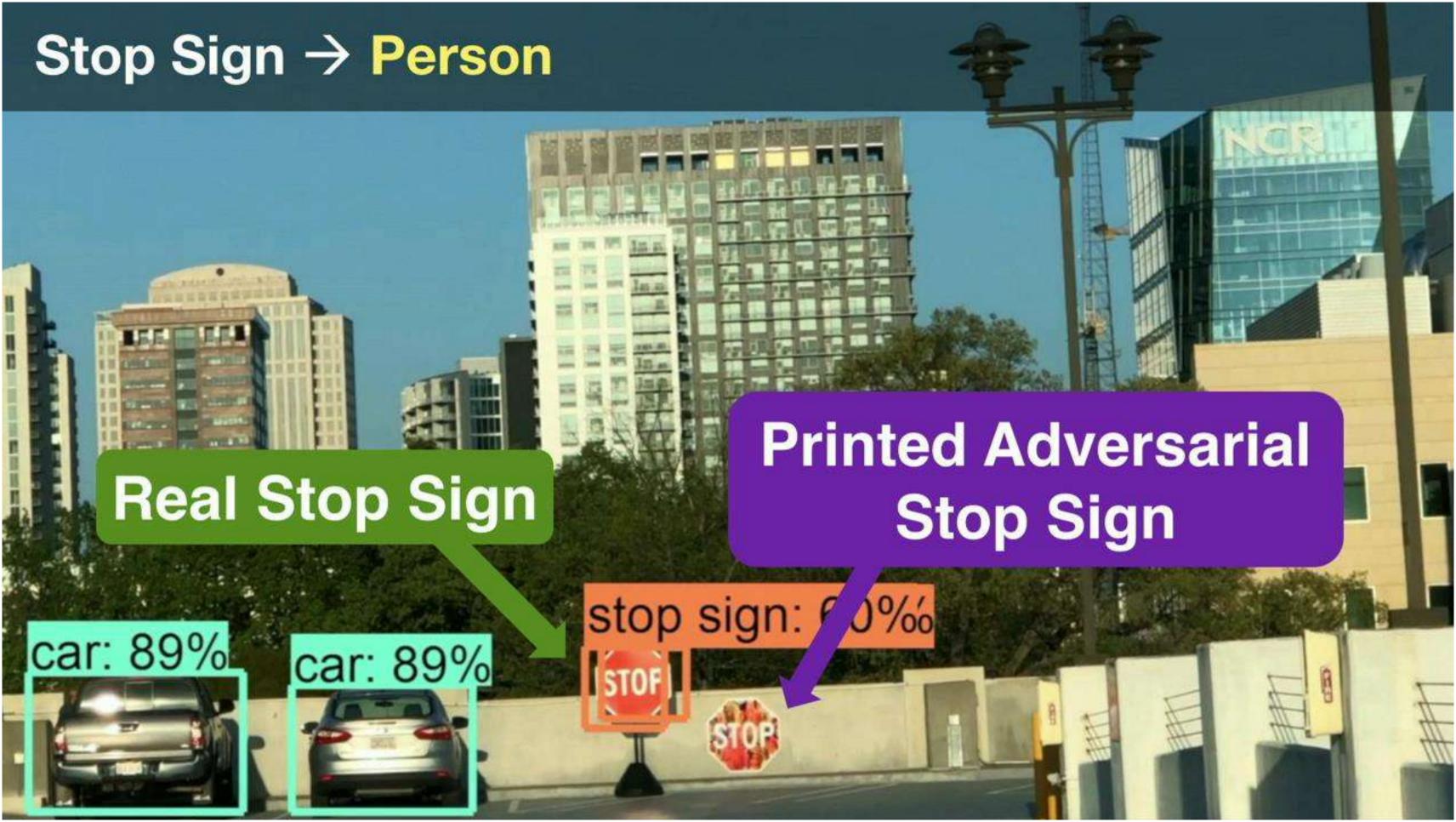
SCALABLE METHODS, INTERACTIVE VISUALIZATIONS, PRACTICAL TOOLS

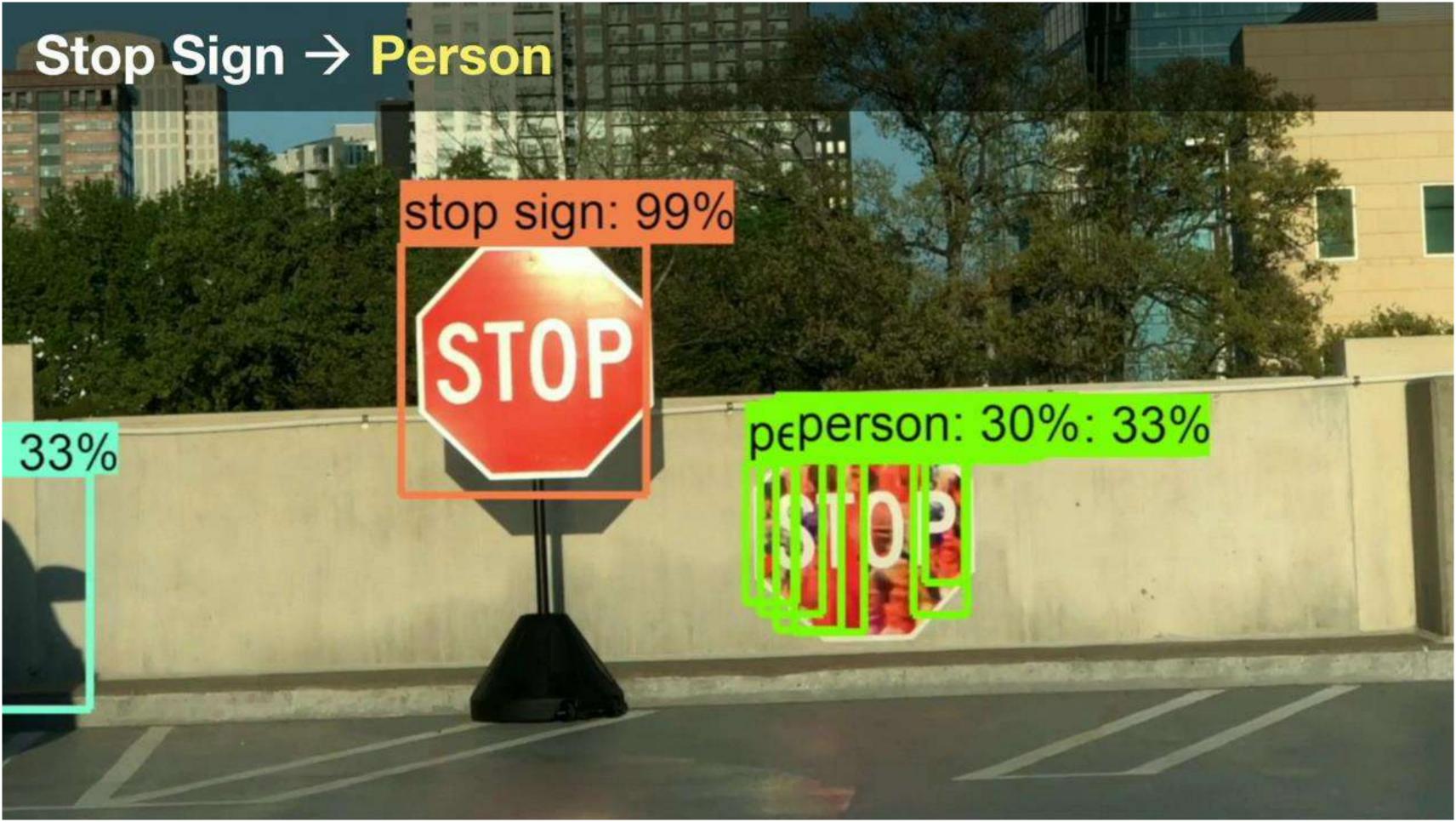


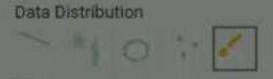
### Polo Chau

Associate Professor Associate Director, MS Analytics Georgia Tech

poloclub.github.io



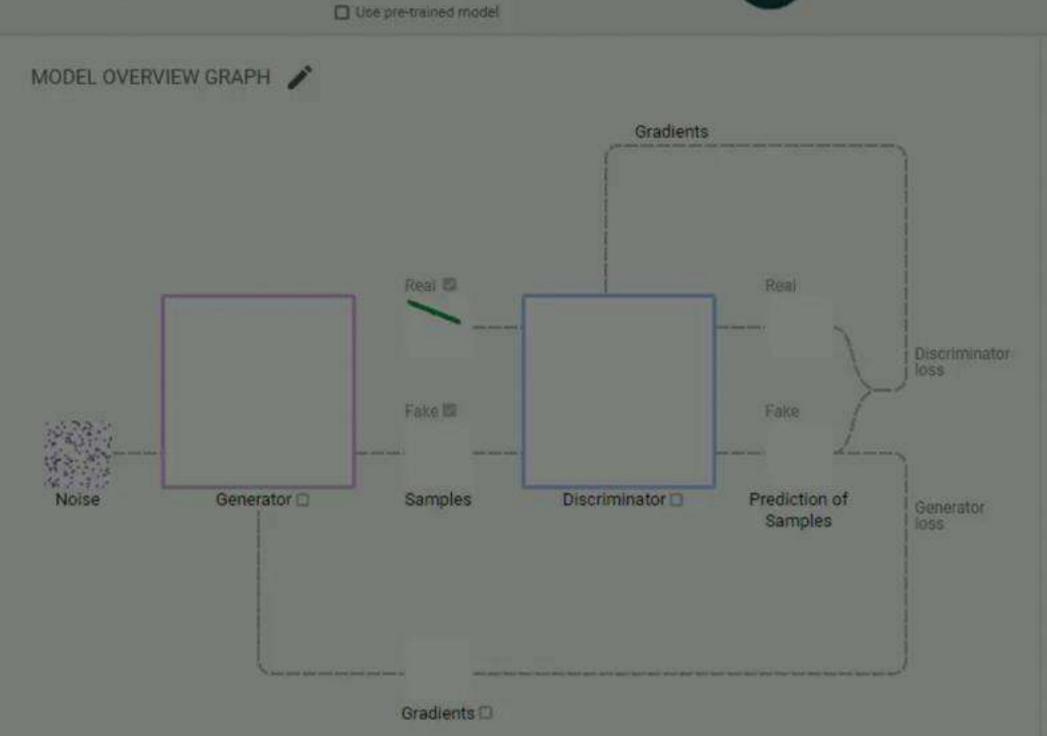








Epoch 000,000

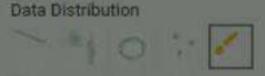


#### LAYERED DISTRIBUTIONS

[hm]

Draw a distribution above, then click the apply button.

APPLY



Use pre-trained model

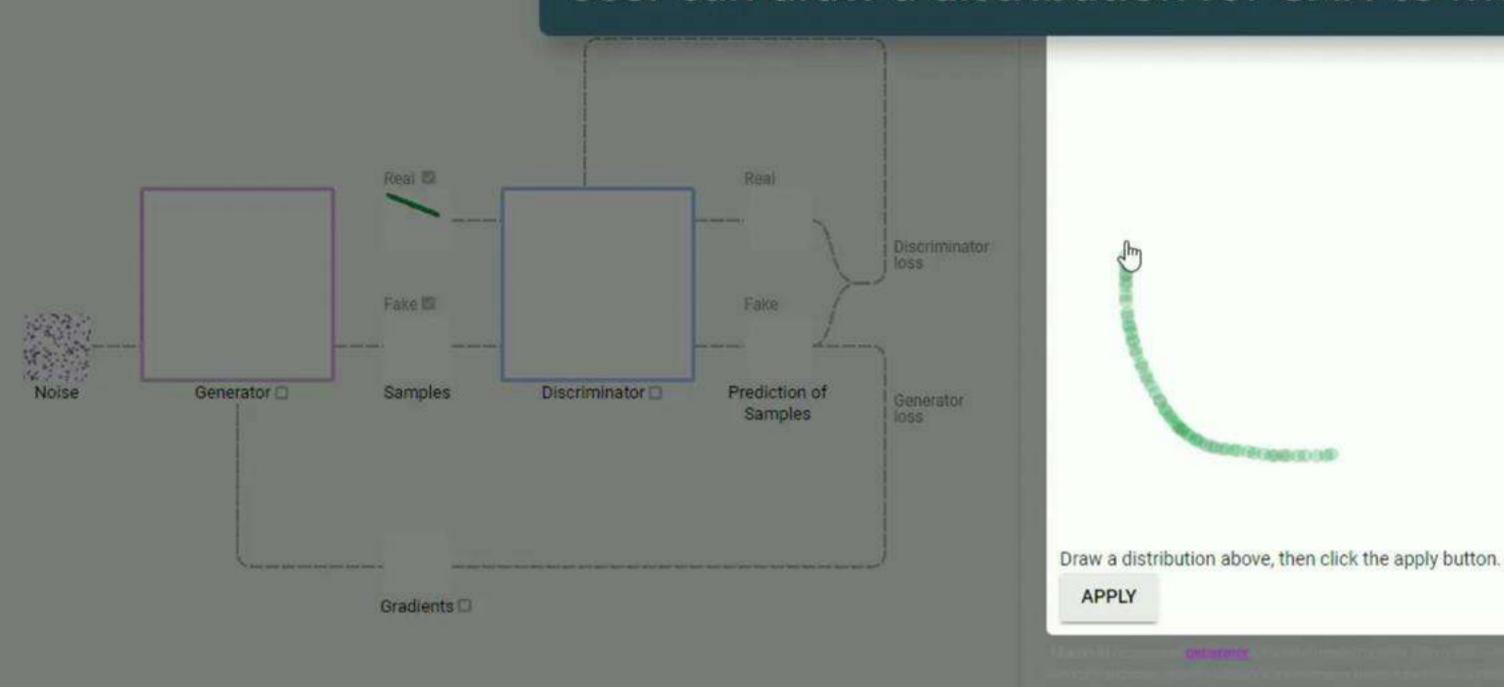


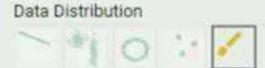


Epoch 000,000

MODEL OVERVIEW GRAPH 🧪

#### User can draw a distribution for GAN to model

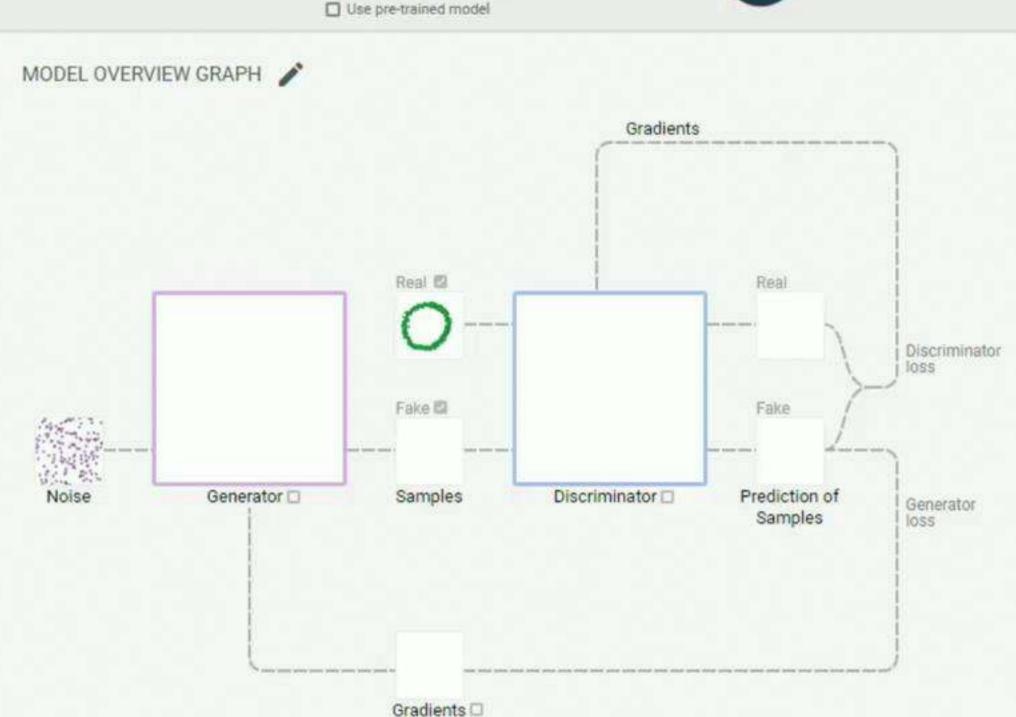








Epoch 000,000



#### LAYERED DISTRIBUTIONS



Each dot is a 2D data sample: real samples, fake samples.

Background Inlors of grid cells represent discriminator's classifications.

Samples in green regions are likely to be real; those in purple regions likely fake.

Manifold represents generator's transformation results from noise space.

Opacity encodes density: darker purple means more samples in smaller area.

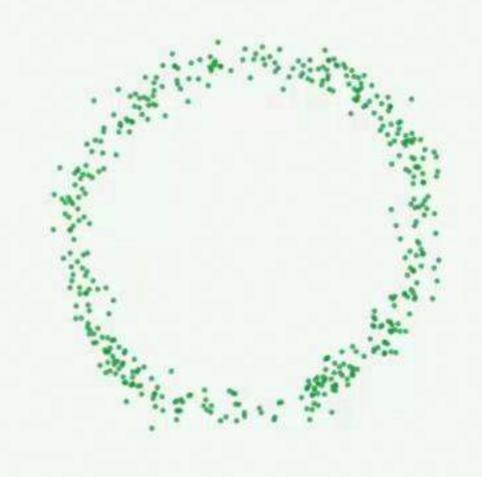
**Data Distribution GAN** Lab Click play to ☐ Use p Run/Pause training start training MODEL OVERVIEW GRAPH 3 update(s) per epoch Gradients Optimizer: SGD Learning rate: 0.1 Real 🖾 Real Discriminator Fake 🖾 Fake Log loss \* Discriminator Prediction of Noise Generator Samples Generator Samples 🗘 1 hidden layer(s) 🐧 10 neuron(s) 2D Uniforn -2 hidden layer(s) 2 neuron(s) loss Optimizer: SGD Learning rate: 0.03 \*

Gradients

1 update(s) per epoch

000,000

#### LAYERED DISTRIBUTIONS



Each dot is a 2D data sample: real samples, fake samples.

Background colors of grid cells represent **discriminator**'s classifications.

Samples in green regions are likely to be real; those in purple regions likely fake.

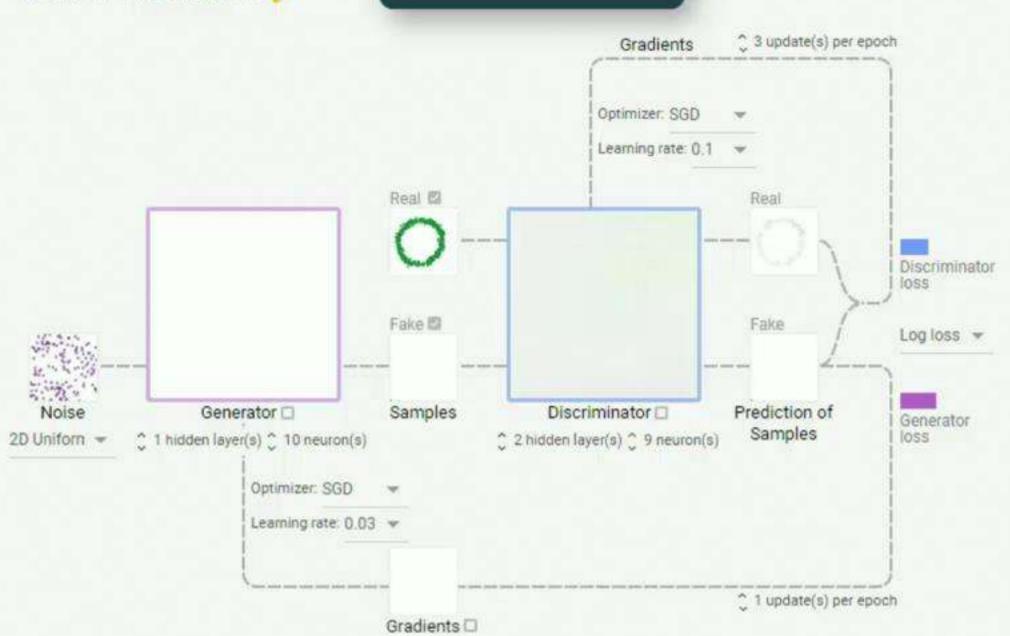
Manifold represents generator's transformation results from noise space.

Opacity encodes density: darker purple means more samples in smaller area.

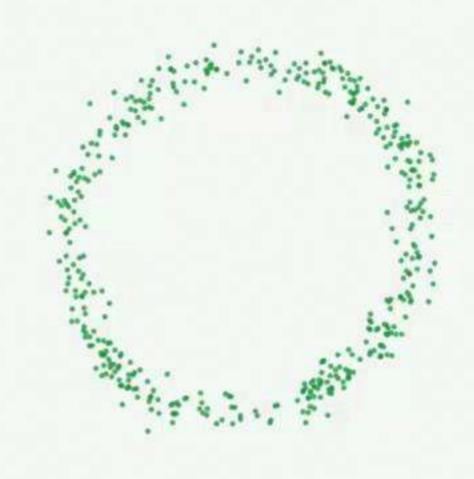
Click play to start training

000,014

MODEL OVERVIEW GRAPH 🧪



#### LAYERED DISTRIBUTIONS



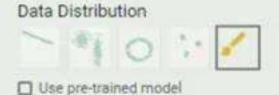
Each dot is a 2D data sample: real samples; fake samples.

Background colors of grid cells represent **discriminator**'s classifications.

Samples in green regions are likely to be real; those in purple regions likely fake.

Manifold represents generator's transformation results from noise space.

Opacity encodes density: darker purple means more samples in smaller area.



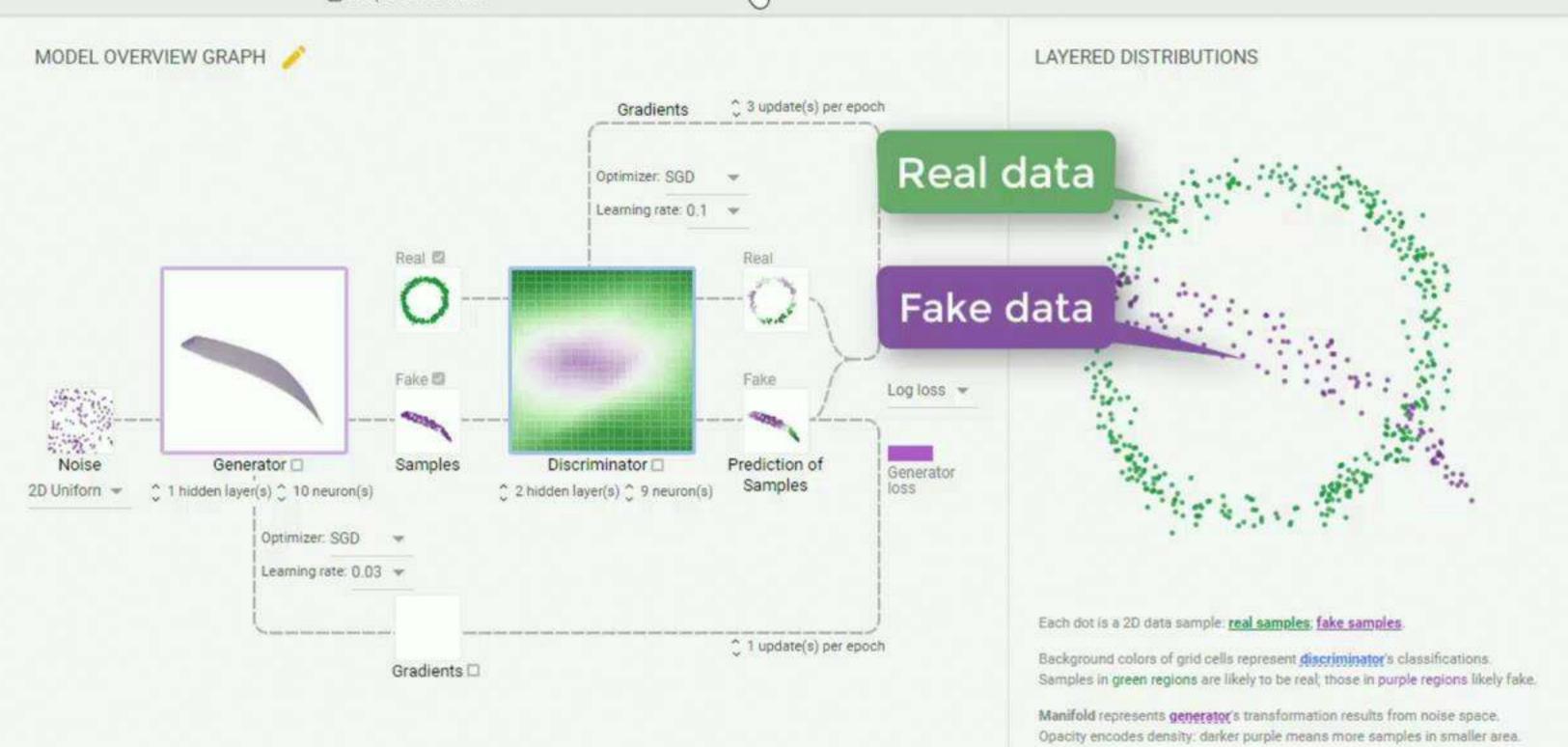




Epoch 000,092

Pink lines from fake samples represent gradients for generator.

✓ This sample needs to move upper right to decrease generator's loss.



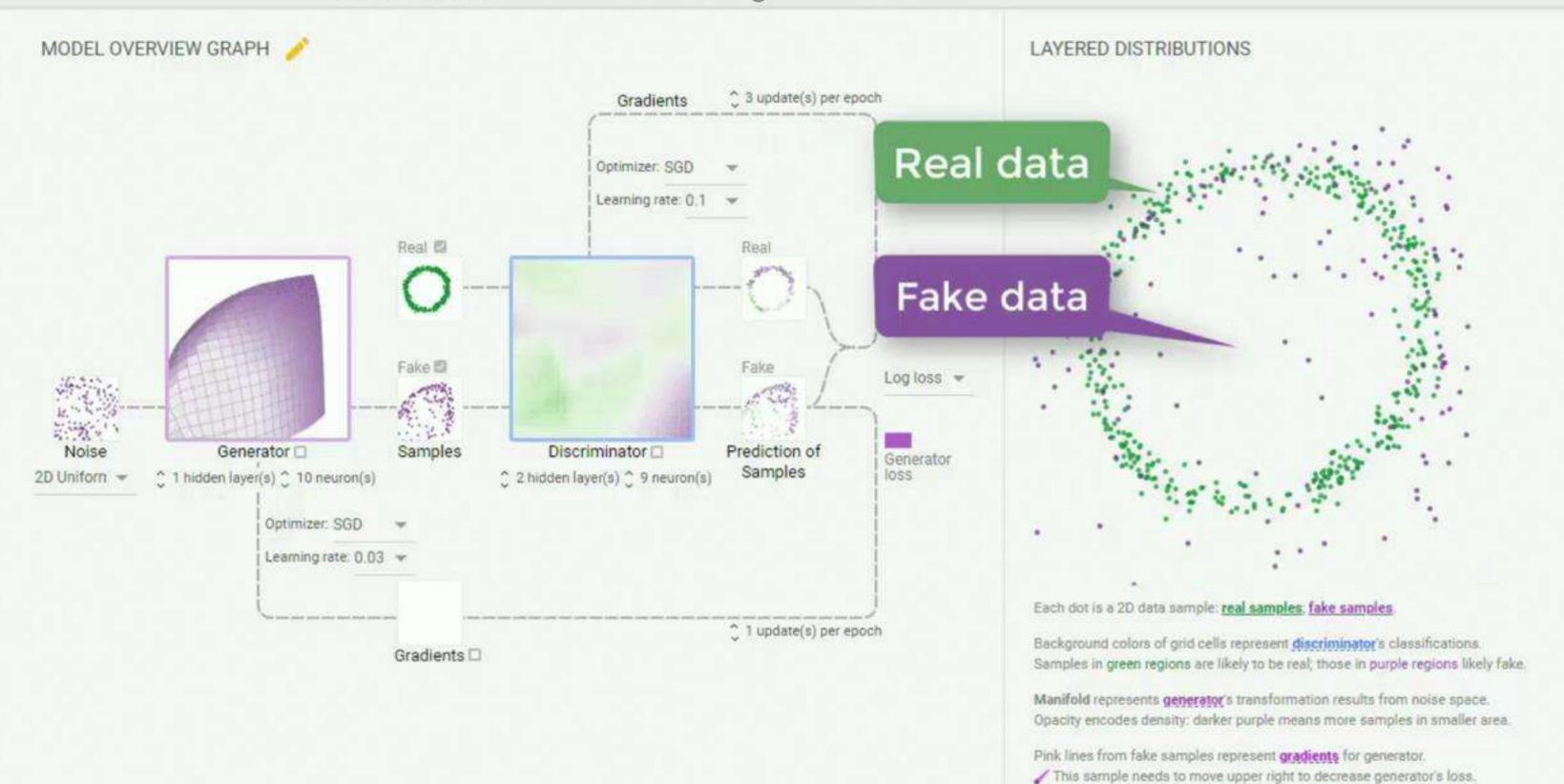


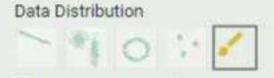




M

Epoch 000,162

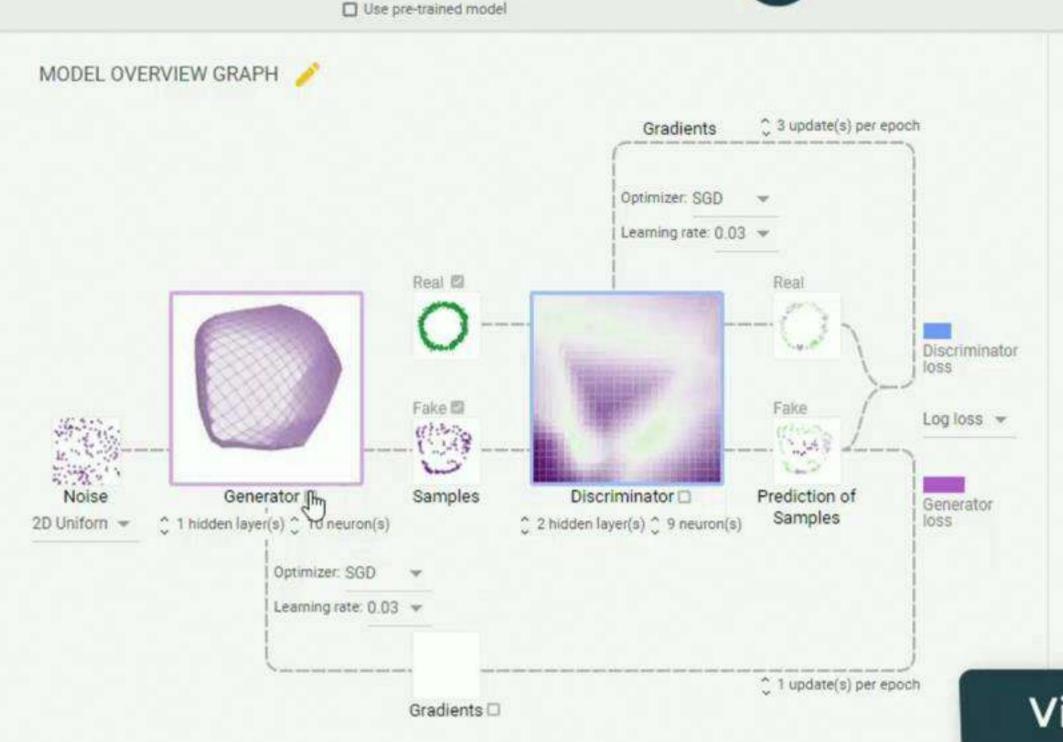


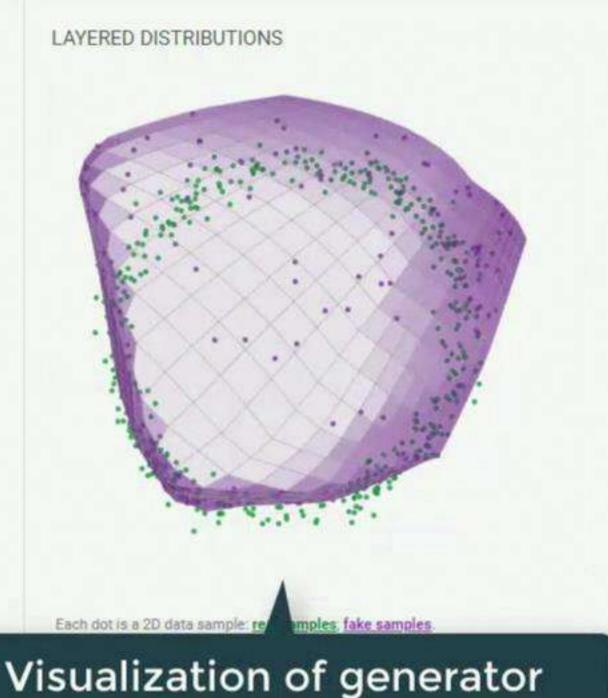






Epoch 000,574

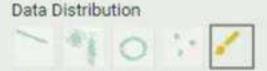




Pink lines from fake samples represent gradients for generator.

using manifold

✓ This sample needs to move upper right to decrease generator's loss.



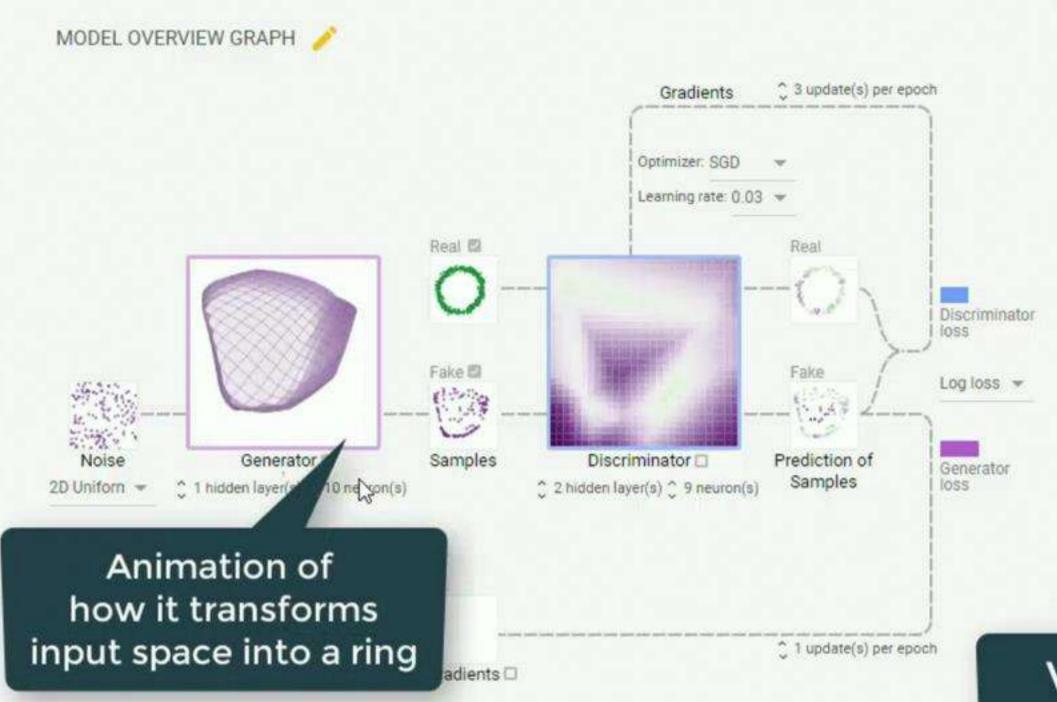
Use pre-trained model





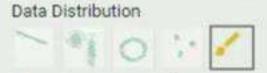
M

Epoch 000,662



# LAYERED DISTRIBUTIONS Each dot is a 20 data sample: re

Visualization of generator using manifold

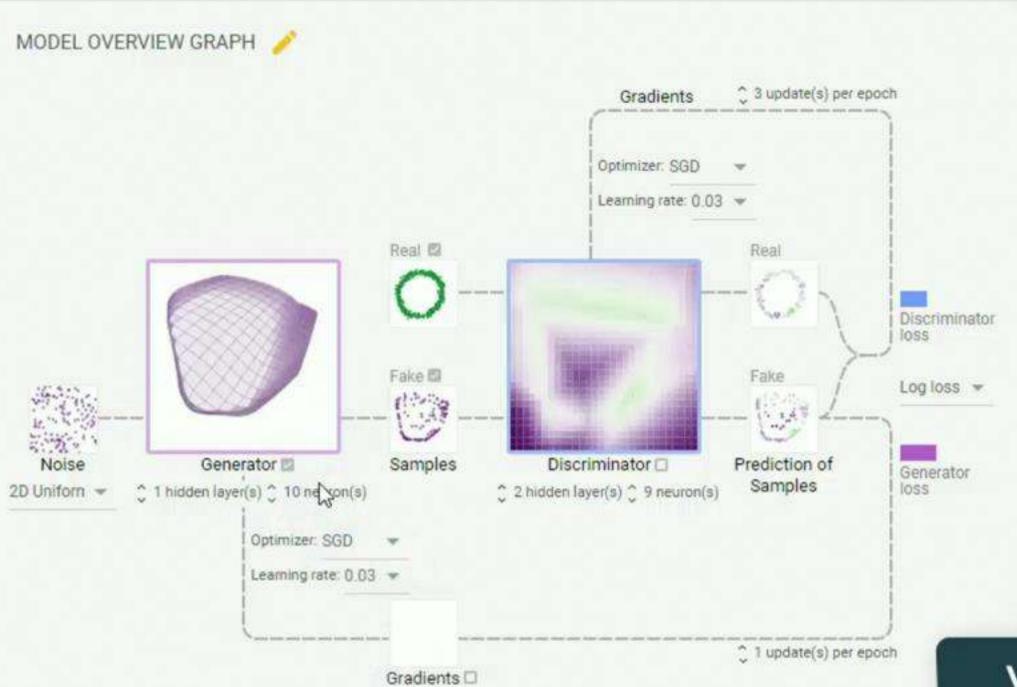


Use pre-trained model

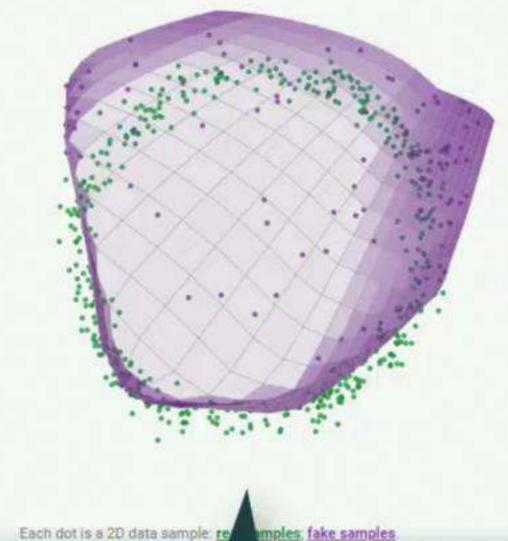




Epoch 000,744



## LAYERED DISTRIBUTIONS



Visualization of generator using manifold

Pink lines from fake samples represent gradients for generator.

✓ This sample needs to move upper right to decrease generator's loss.

Data Distribution

Use pre-trained model

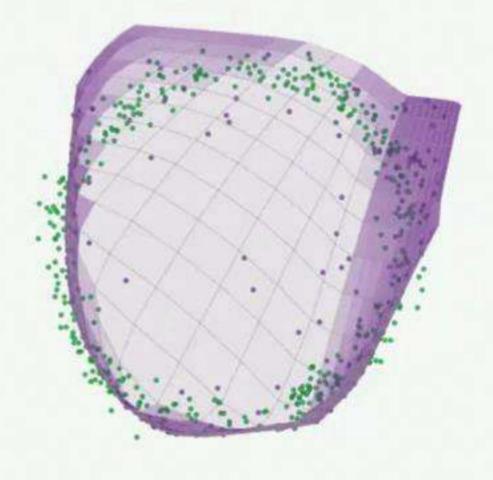
Users can dynamically adjust hyperparameters

2

Epoch

RED DISTRIBUTIONS

MODEL OVERVIEW GRAPH 🧪 3 update(s) per epoch Gradients Optimizer: SGD Learning rate: 0.03 -Real @ Discriminator Fake 🖾 Fake Log loss \* Generator 🖾 Samples Discriminator ... Prediction of Noise Generator Samples 1 hidden layer(s) 10 neuron(s) 2D Uniforn -2 hidden layer(s) 2 neuron(s) loss Optimizer: SGD Learning rate: 0.03 -1 update(s) per epoch Gradients 🗆



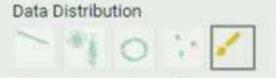
Each dot is a 2D data sample: real samples: fake samples.

Background colors of grid cells represent **discriminator**'s classifications.

Samples in green regions are likely to be real; those in purple regions likely fake.

Manifold represents generator's transformation results from noise space.

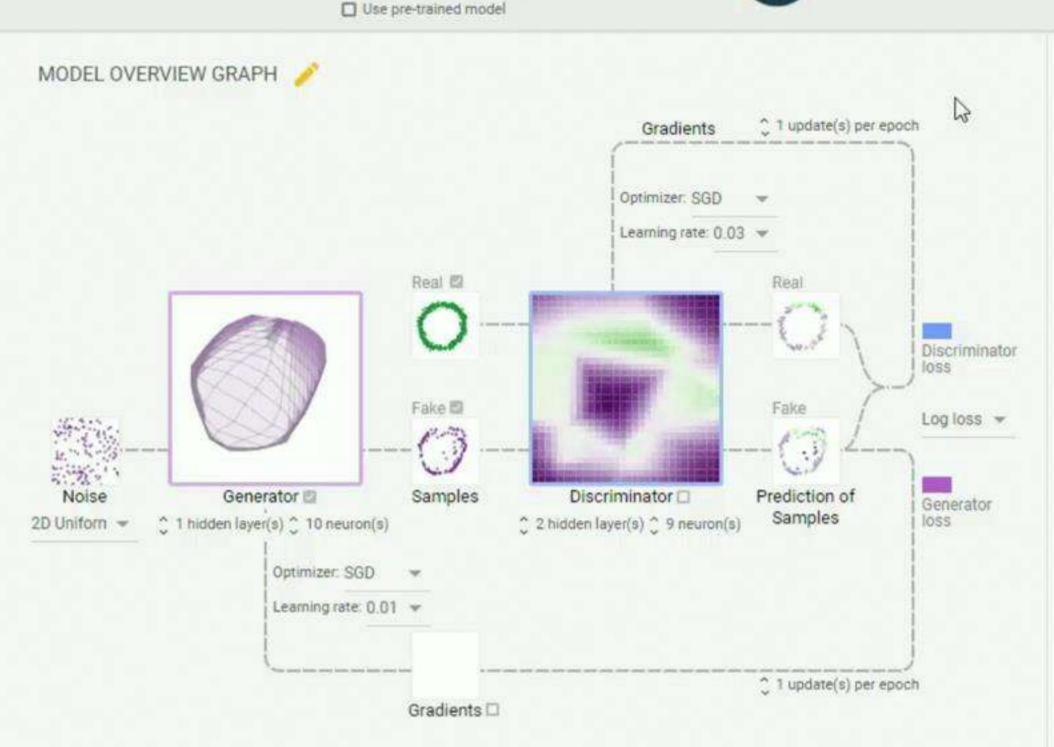
Opacity encodes density: darker purple means more samples in smaller area.



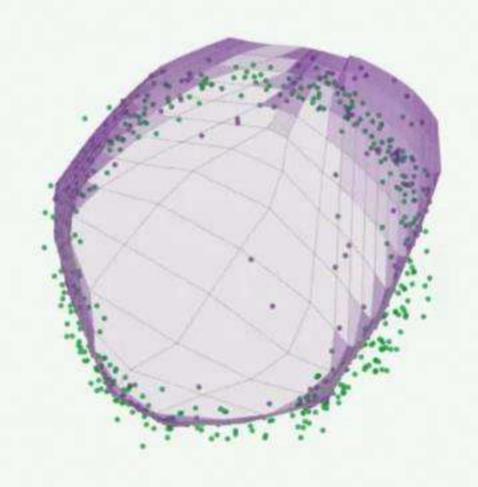




Epoch 002,824



#### LAYERED DISTRIBUTIONS



Each dot is a 20 data sample: real samples; fake samples.

Background colors of grid cells represent discriminator's classifications.

Samples in green regions are likely to be real; those in purple regions likely fake.

Manifold represents generator's transformation results from noise space.

Opacity encodes density: darker purple means more samples in smaller area.

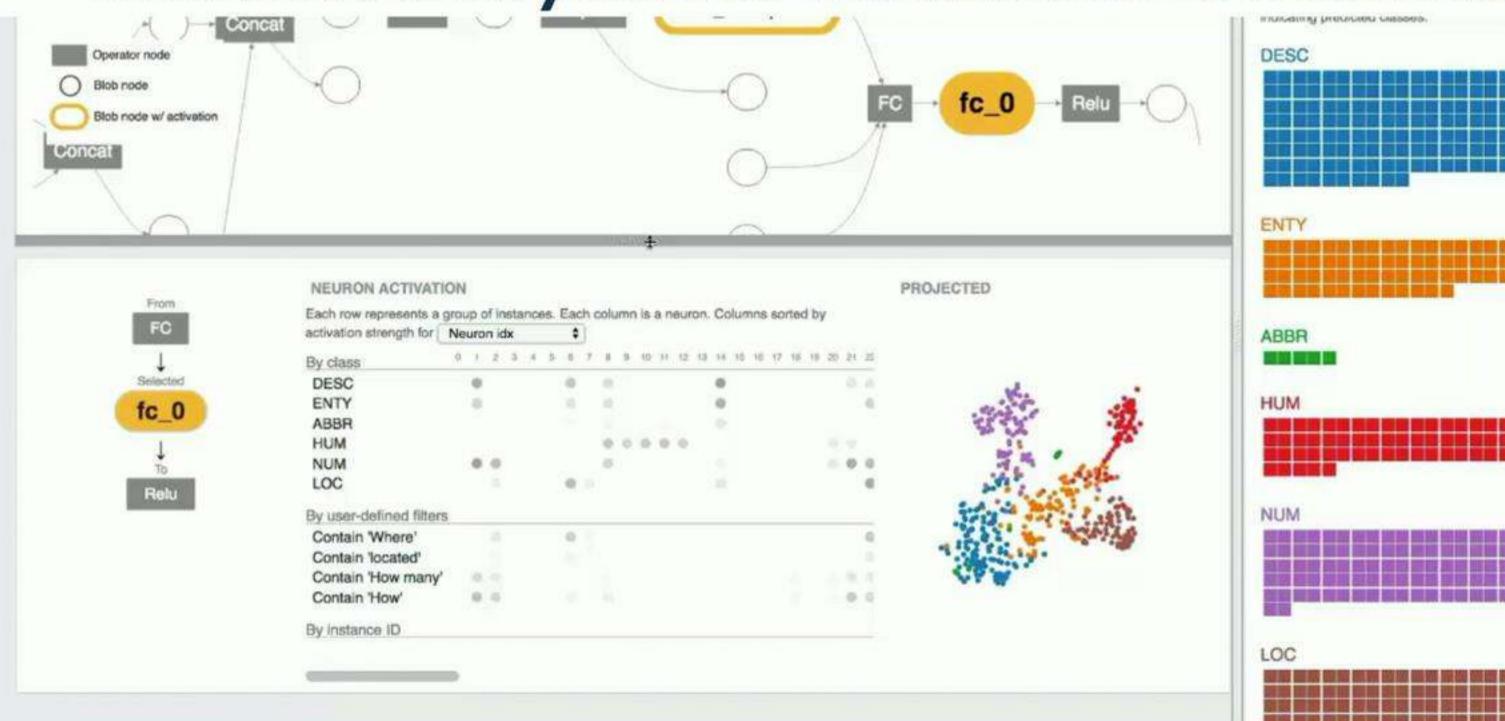
Pink lines from fake samples represent gradients for generator.

This sample needs to move upper right to decrease generator's loss.

Playing at 15x speed

## ActiVis Key Ideas (2)

Unified Analysis for Instances & Subsets

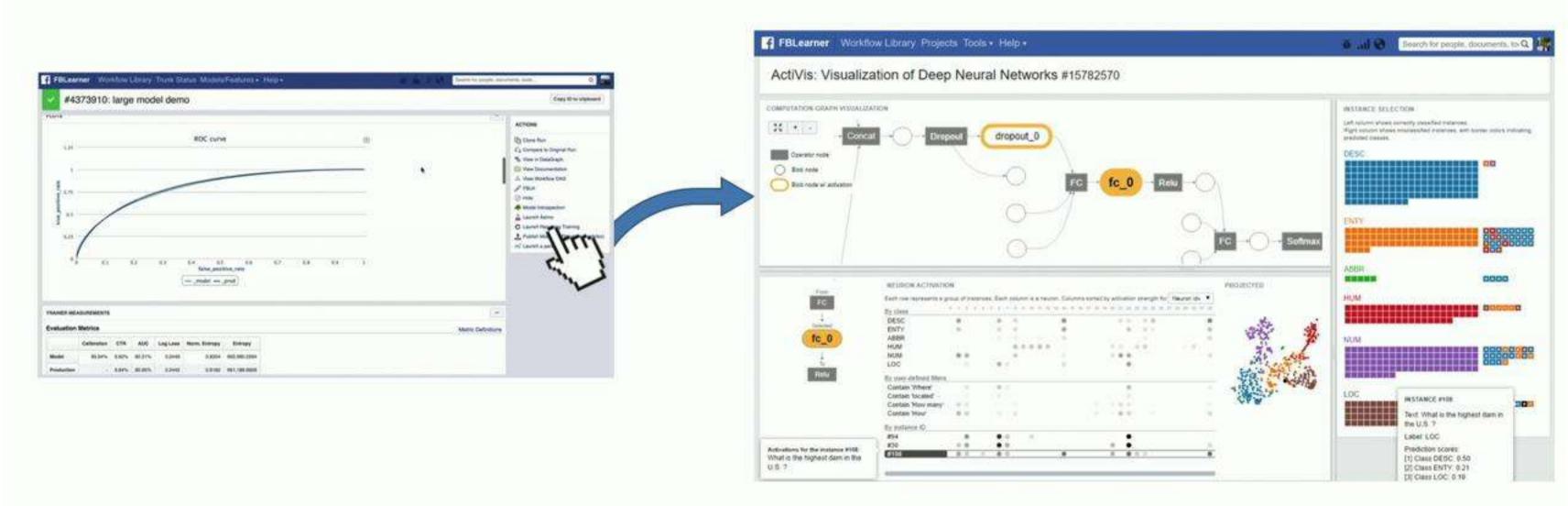


BEN RICH

# Deployed on FBLearner

## Facebook's ML platform

used by >25% of engineering team



## Secure

# Interpretable

Attack & Defense (DNN)

ShapeShifter Shield

Do-it-yourself Adversarial ML

ADAGIO MLsploit Understand Industry Models
ActiVis

**ML** Education

**GAN** Lab

Research landscape

Survey, Gamut