

# Detecting Reliable Instances for Learning

Judy Hoffman

Adversarial Machine Learning in Computer Vision  
CVPR Workshop 2021



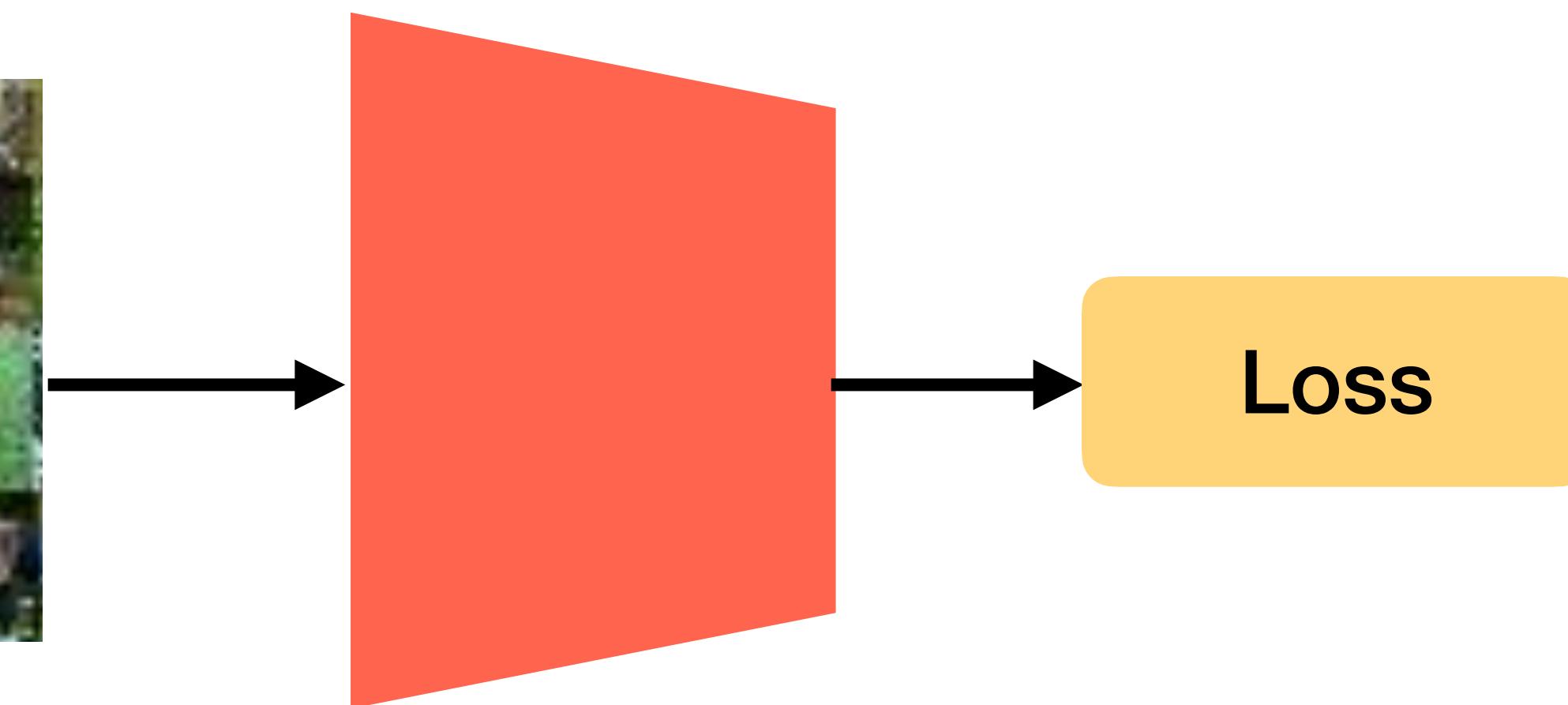
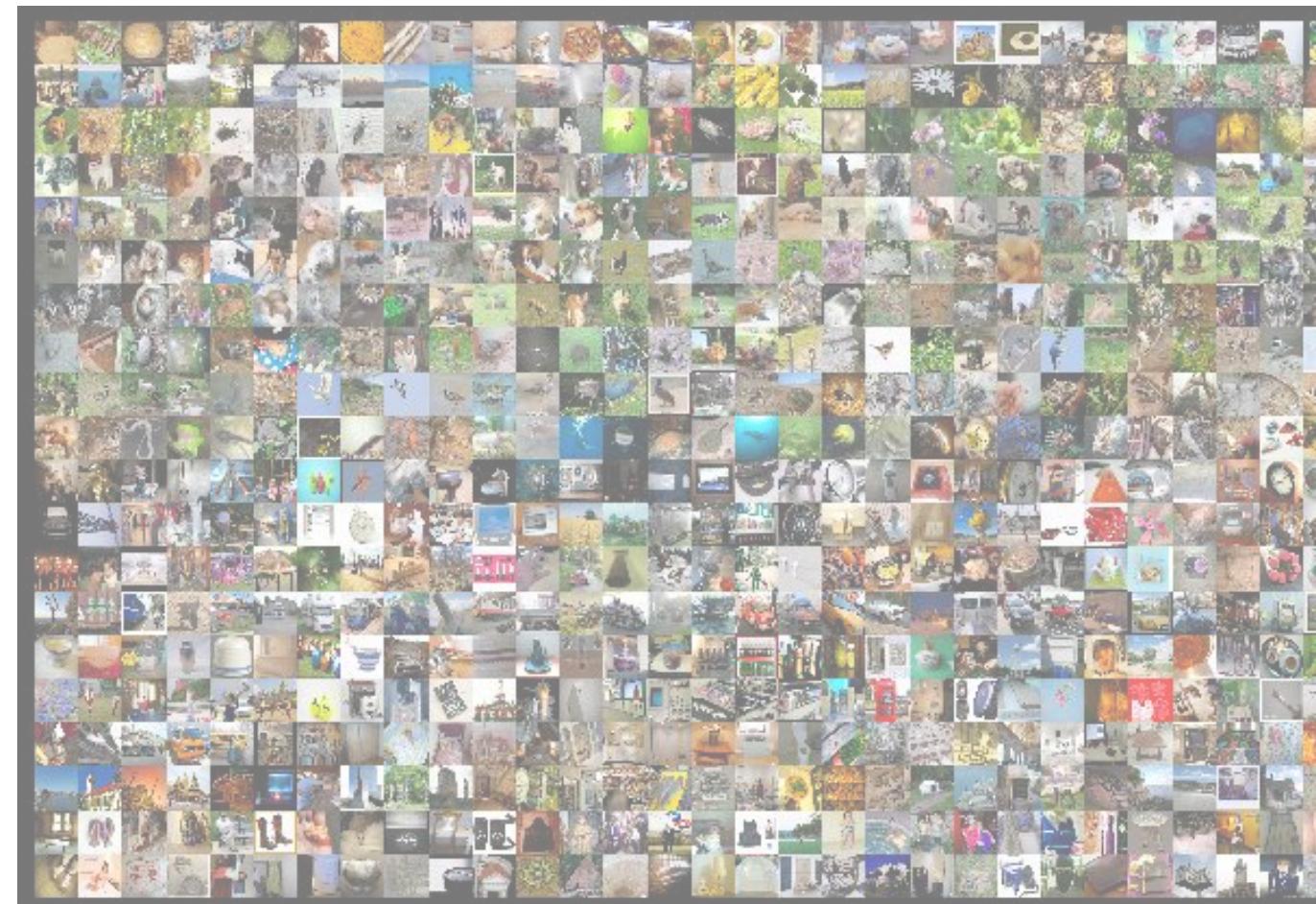
# Standard Supervised Learning

---



# Standard Supervised Learning

---



**Random Sampling**

# Potential Data Pitfalls

Incorrect Label



**Label: “Cat”**

**Prediction: “Cat”**

Adversarial Manipulations



**Label: “Dog”**

**Prediction: “Cat”**

Variable Difficulty



**Label: “Dog”**

**Prediction: “Cat”**

# Potential Data Pitfalls

Incorrect Label



**Learning with  
Noise**

**Label: “Cat”**

**Prediction: “Cat”**

Adversarial Manipulations

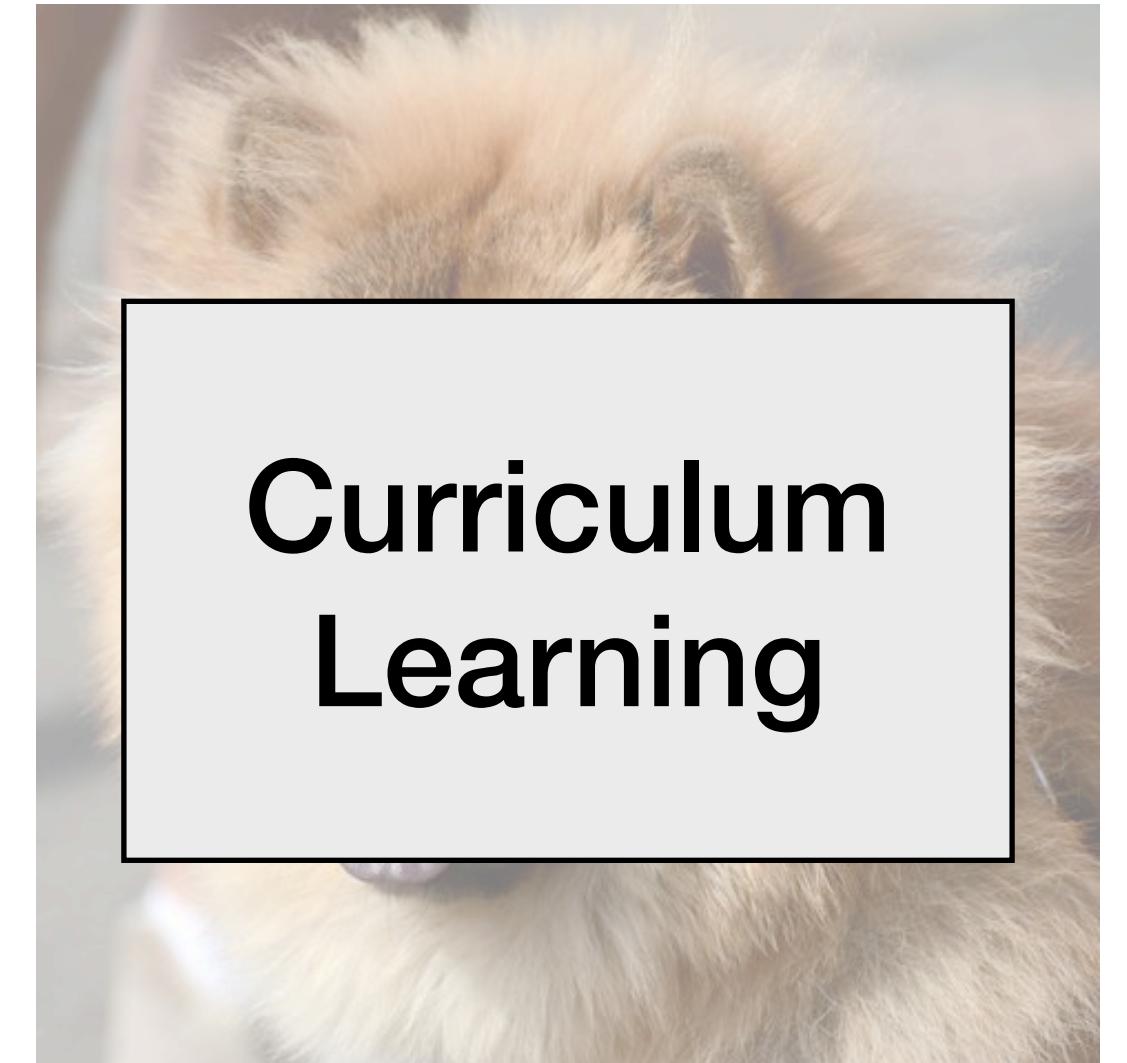


**Adversarial  
Robustness**

**Label: “Dog”**

**Prediction: “Cat”**

Variable Difficulty



**Curriculum  
Learning**

**Label: “Dog”**

**Prediction: “Cat”**

# Enforcing Reliability

# Adversarial Examples



$x$   
“panda”  
57.7% confidence

+ .007 ×



$\text{sign}(\nabla_x J(\theta, x, y))$   
“nematode”  
8.2% confidence

=



$x +$   
 $\epsilon \text{sign}(\nabla_x J(\theta, x, y))$   
“gibbon”  
99.3 % confidence

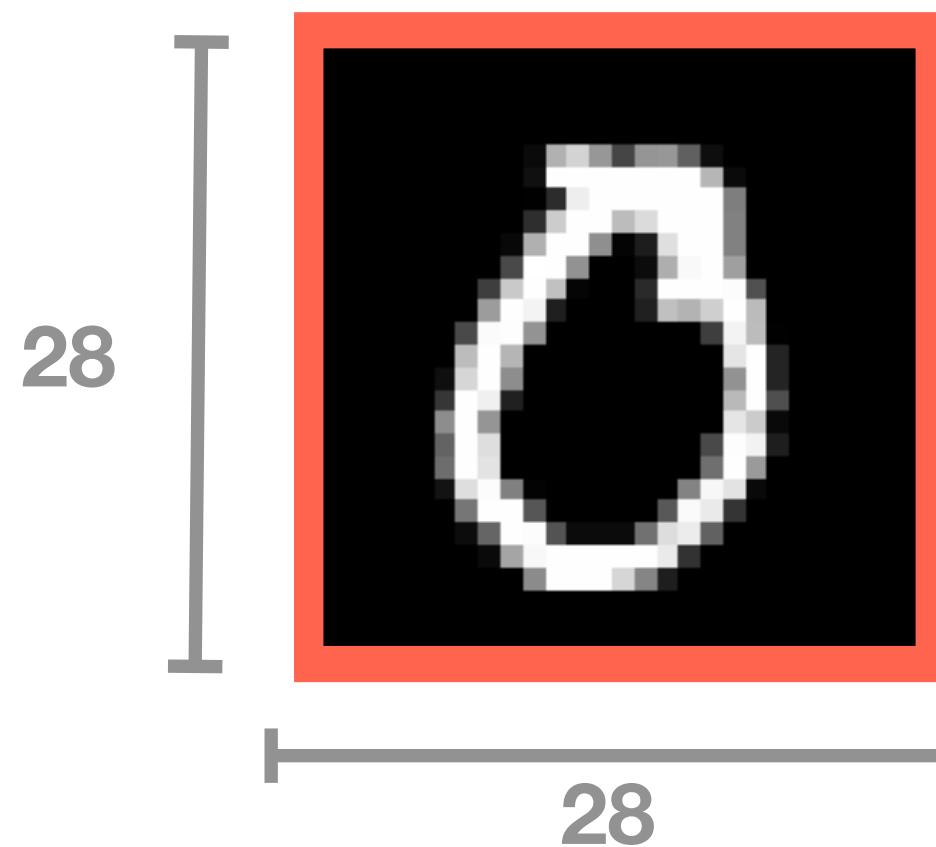
# Visualize Perturbation Space

---

# Visualize Perturbation Space

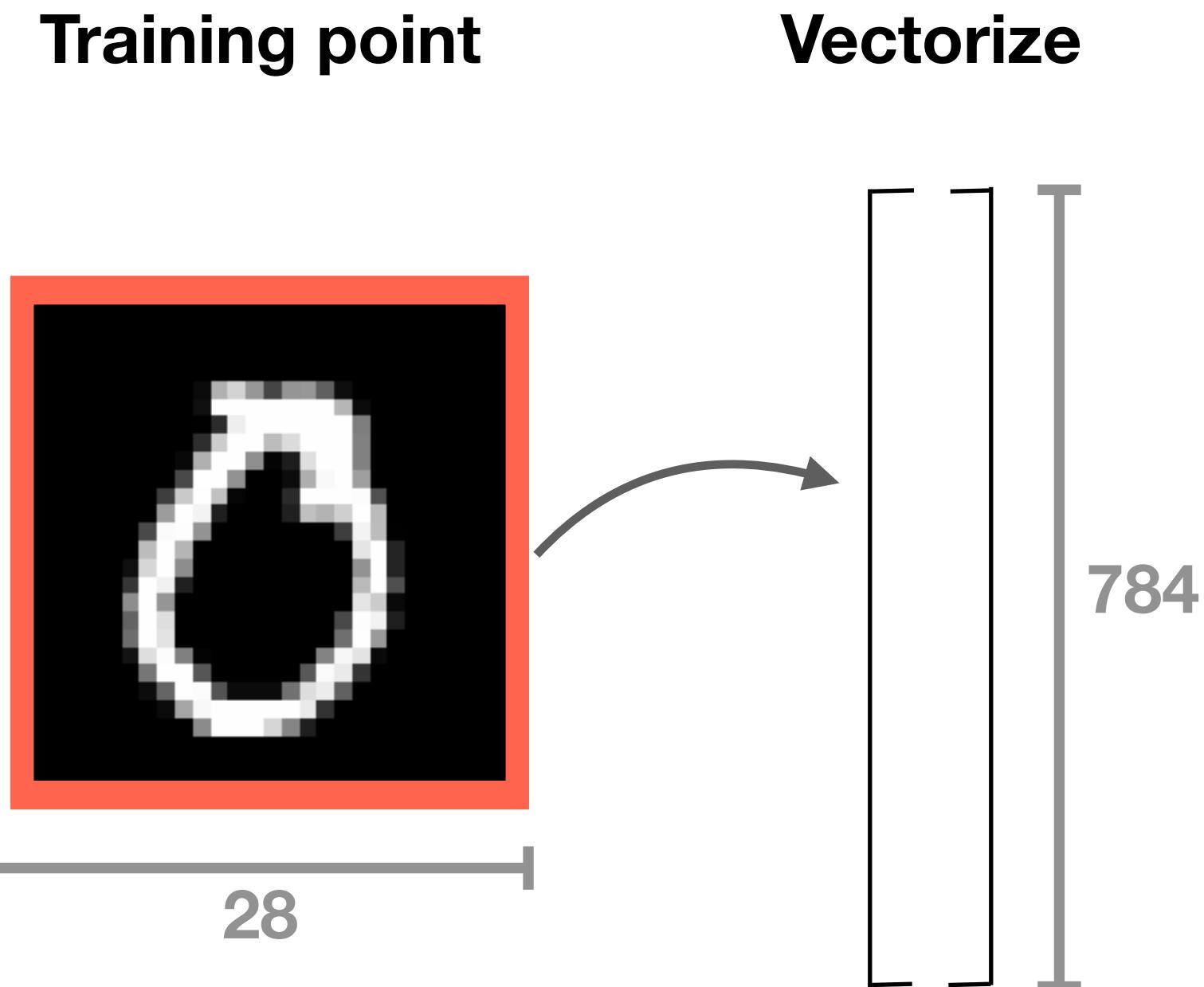
---

**Training point**

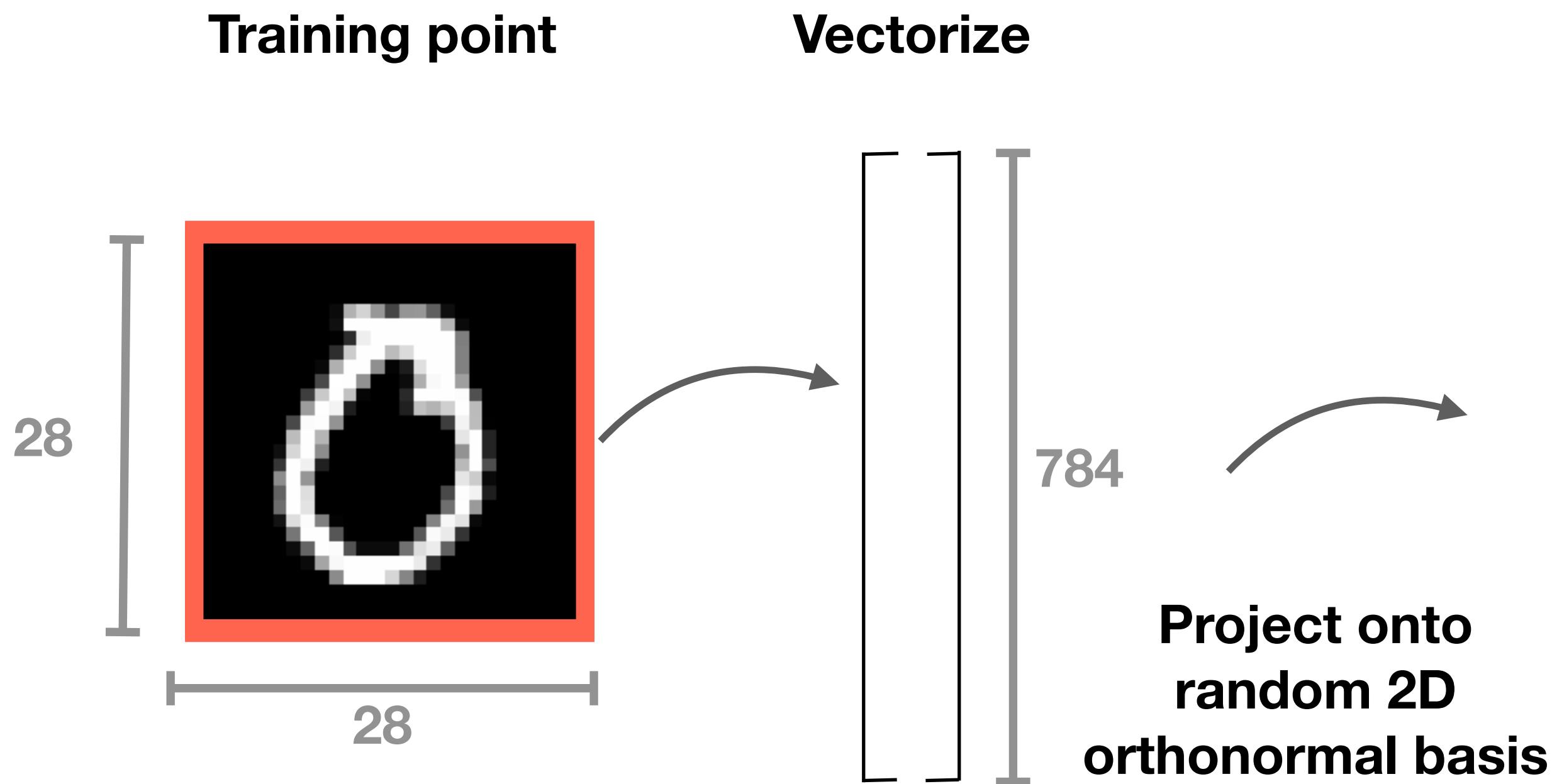


# Visualize Perturbation Space

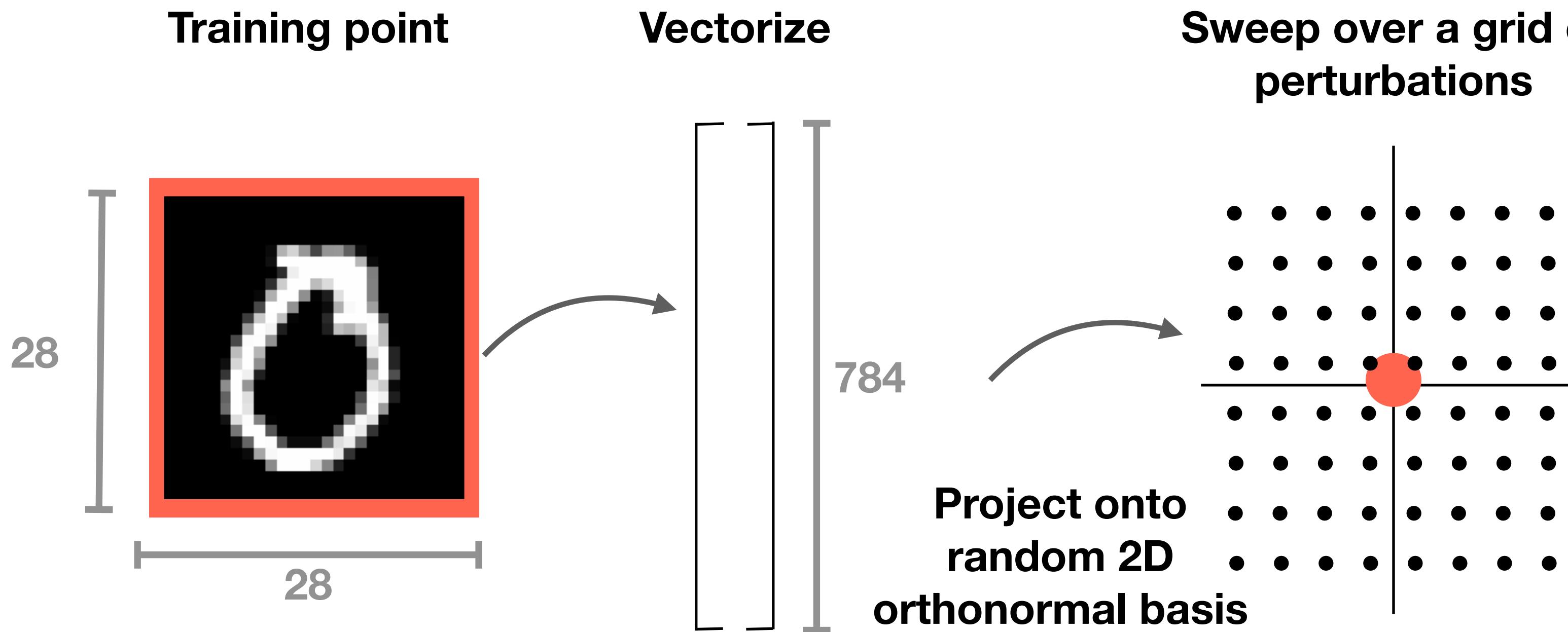
---



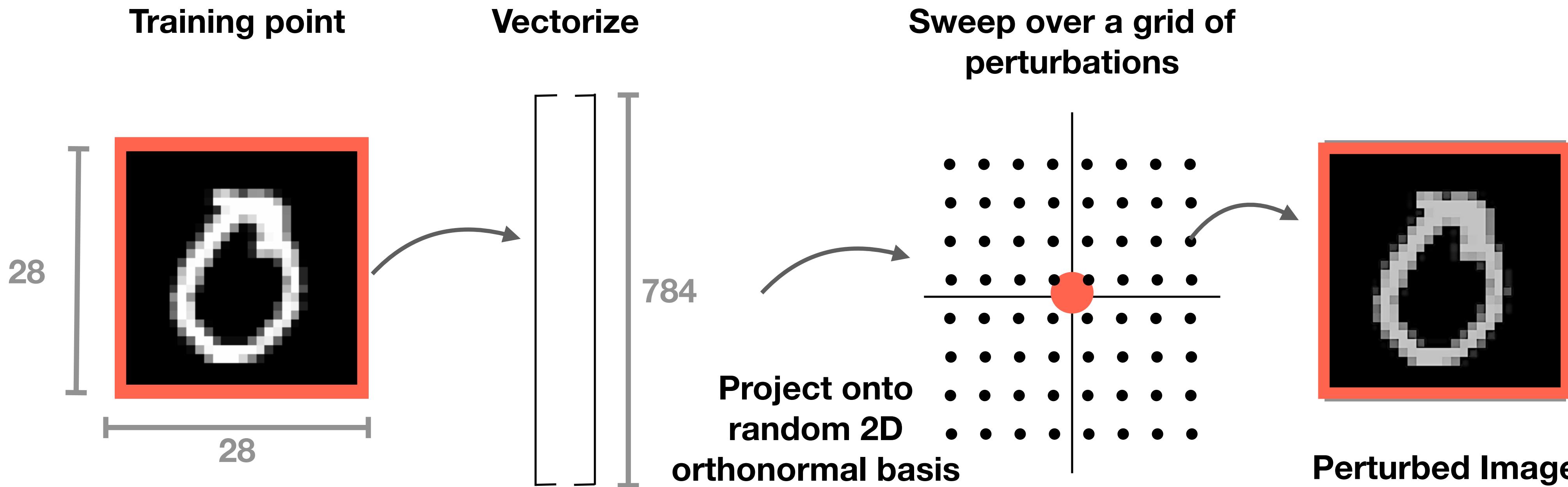
# Visualize Perturbation Space



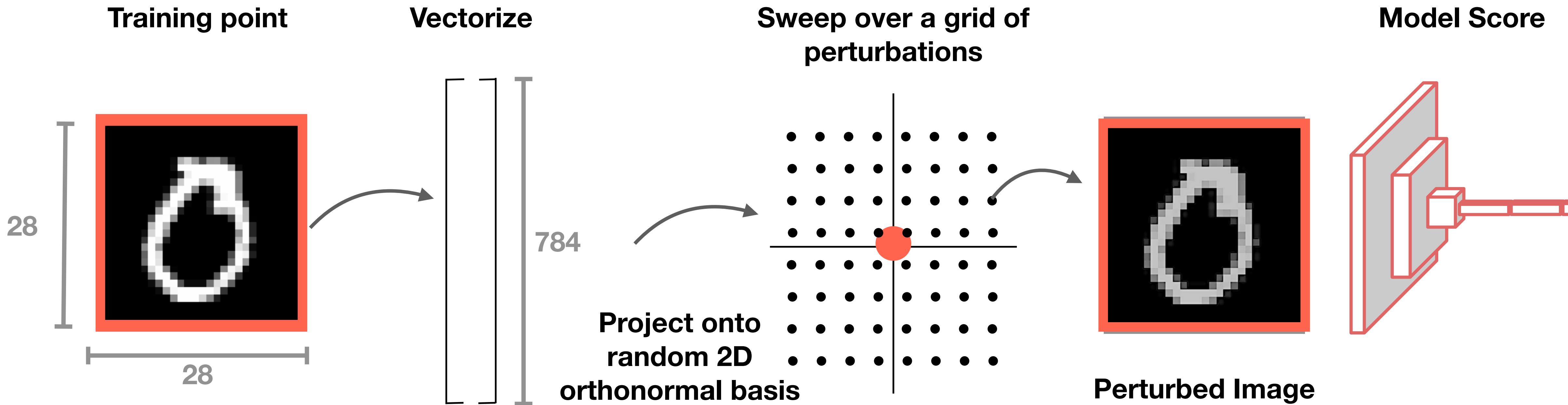
# Visualize Perturbation Space



# Visualize Perturbation Space



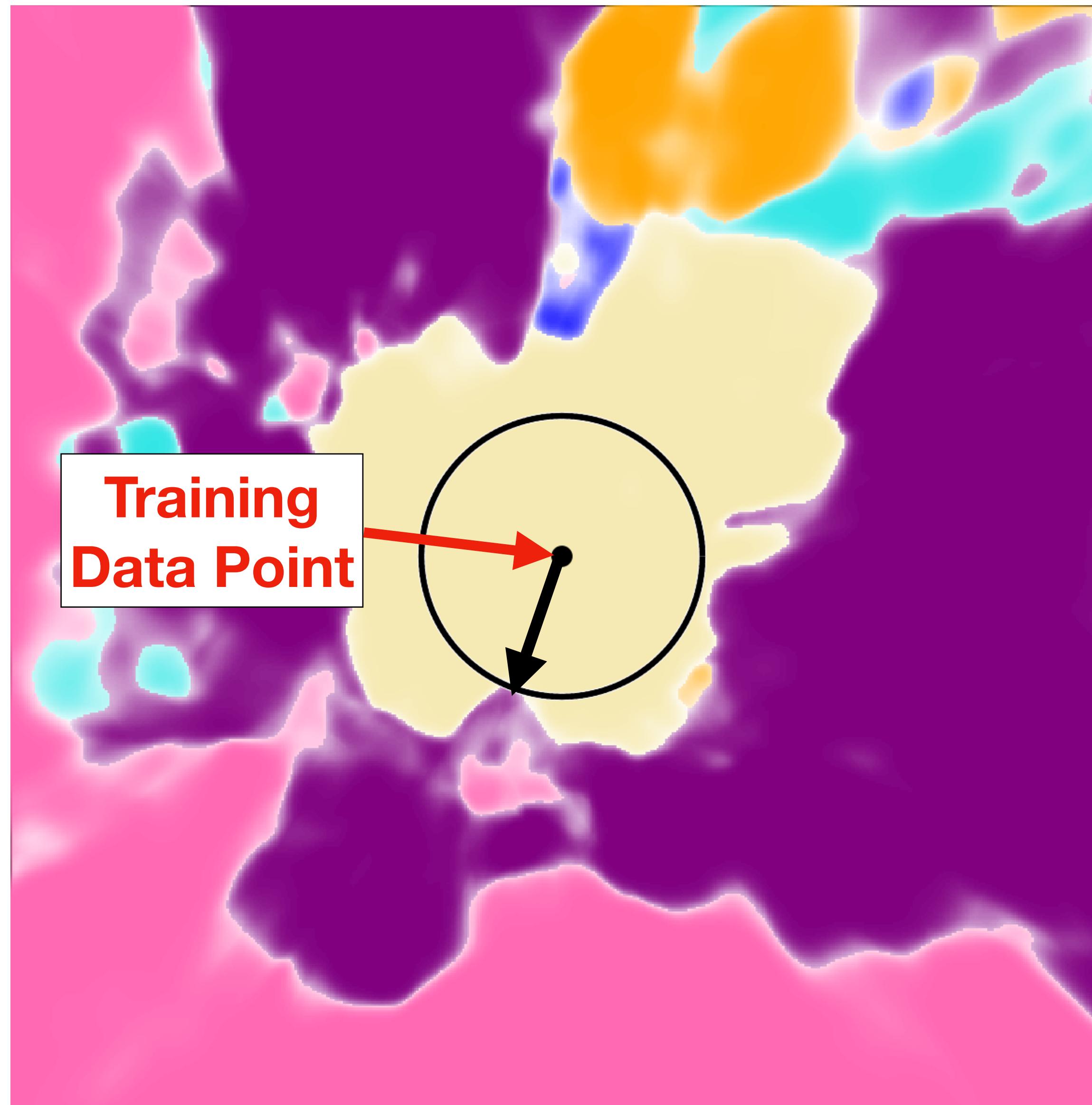
# Visualize Perturbation Space



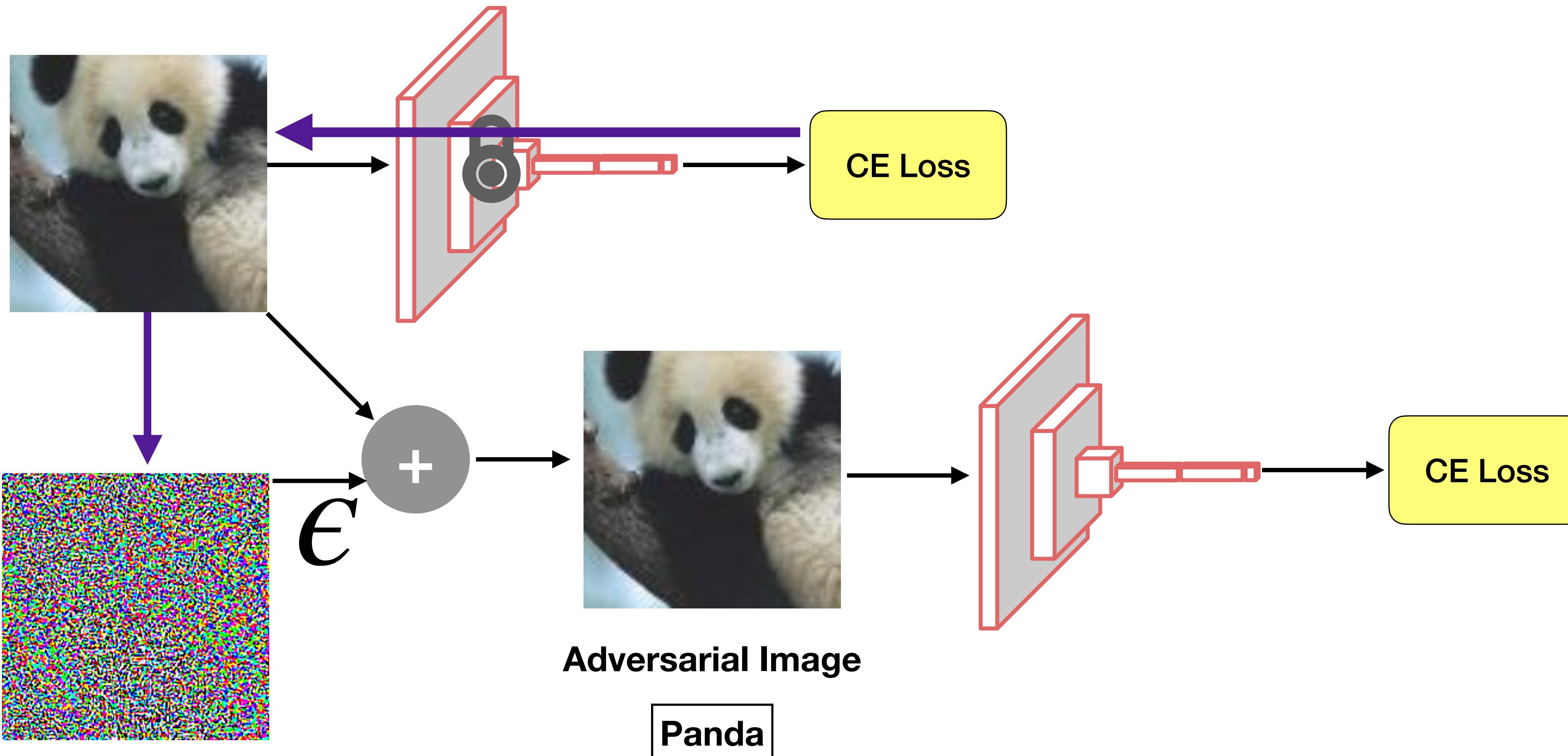
# MNIST LeNet Decisions Around Training Point

Non-smooth  
Decision Boundary

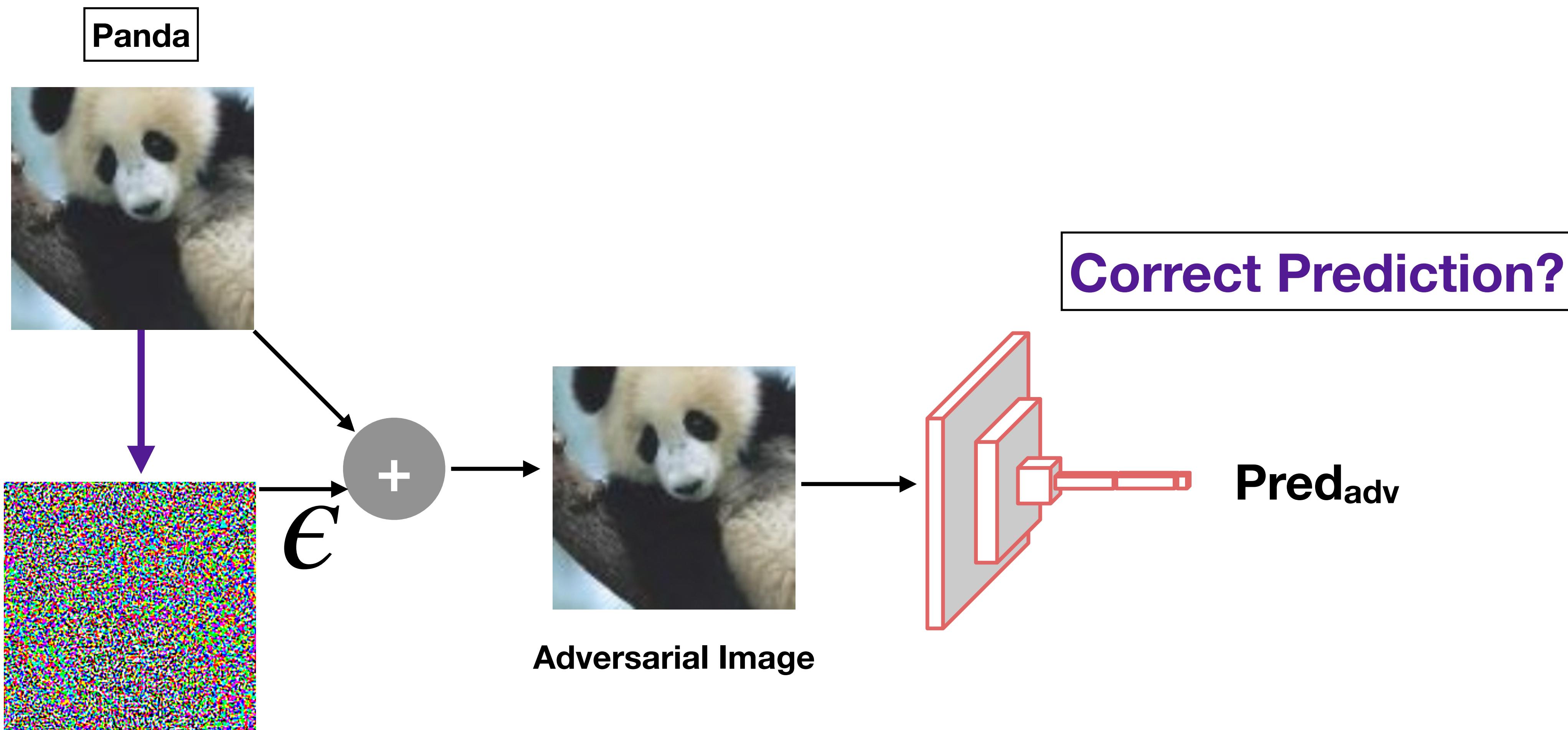
Small perturbations  
lead to new outputs



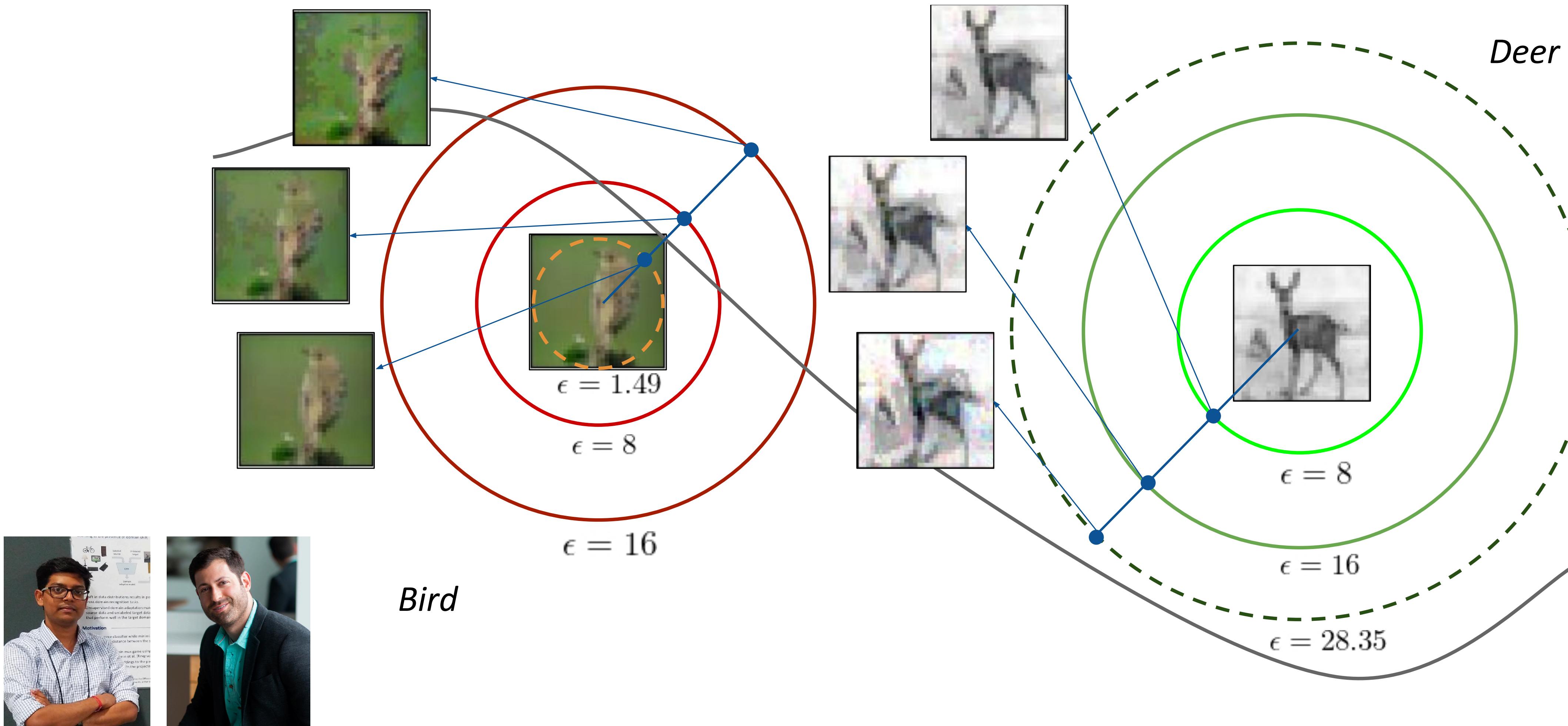
# Adversarial Training



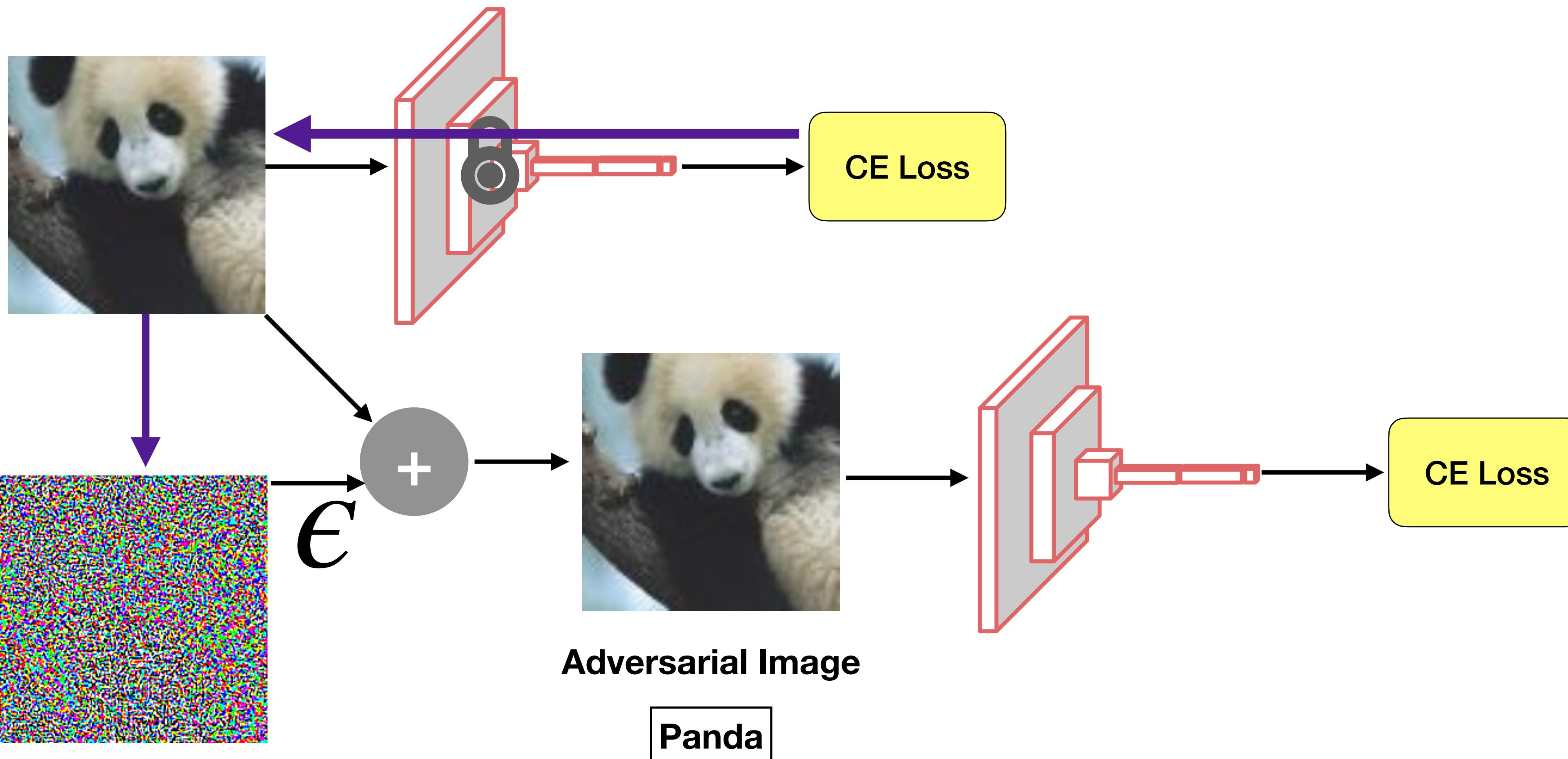
# Adversarial Stability



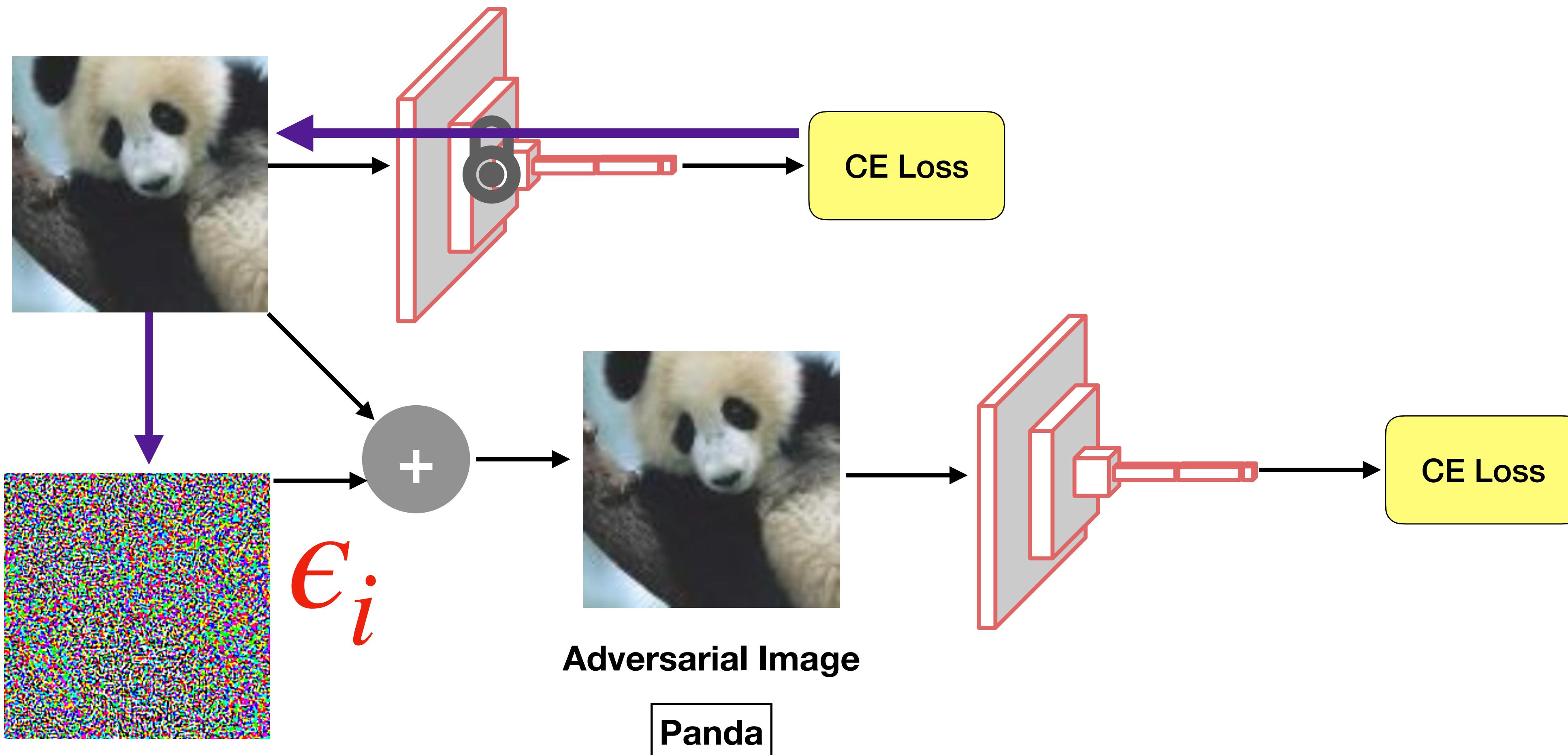
# Instance Adaptive Adversarial Training



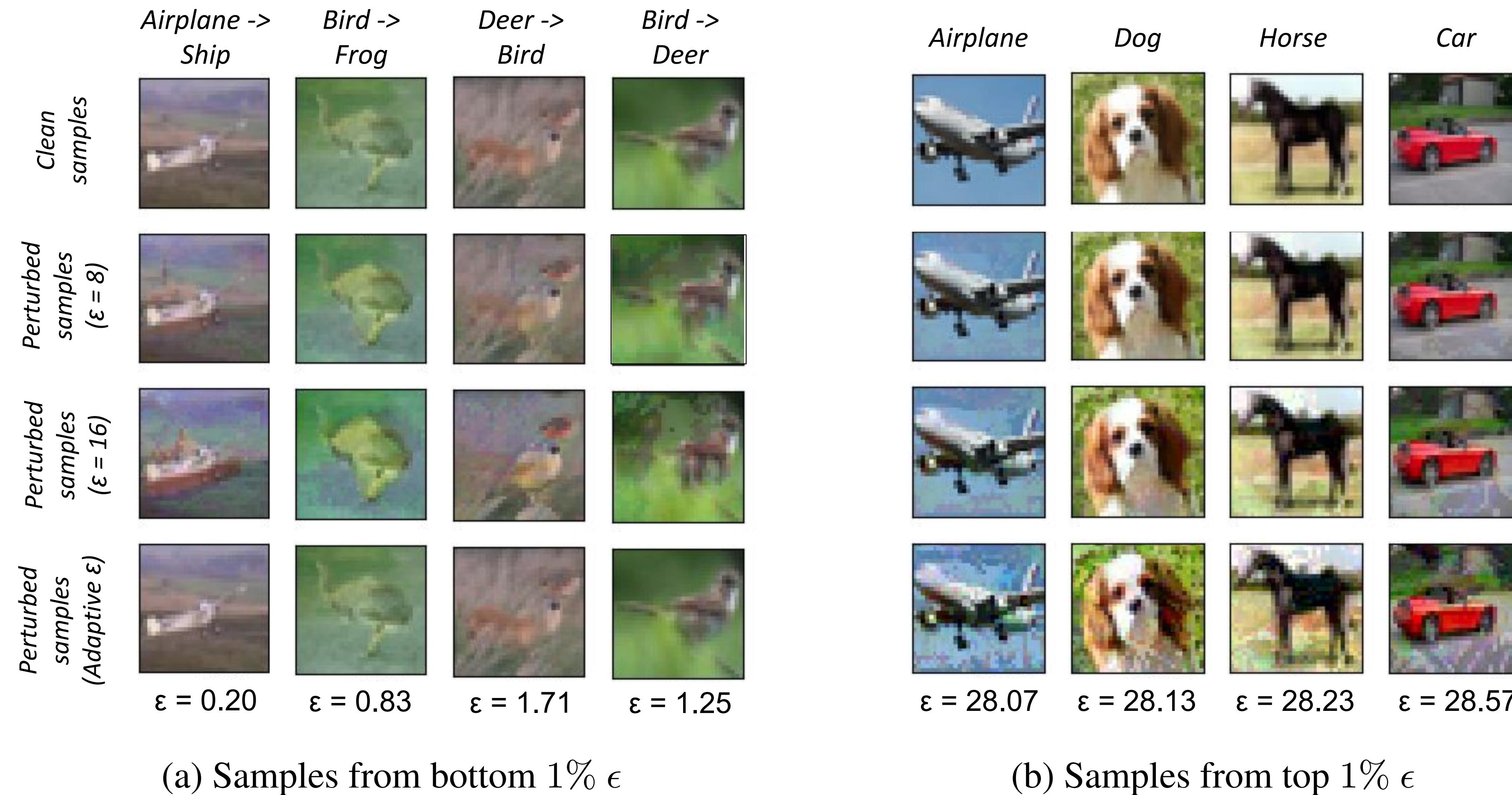
# Adversarial Training



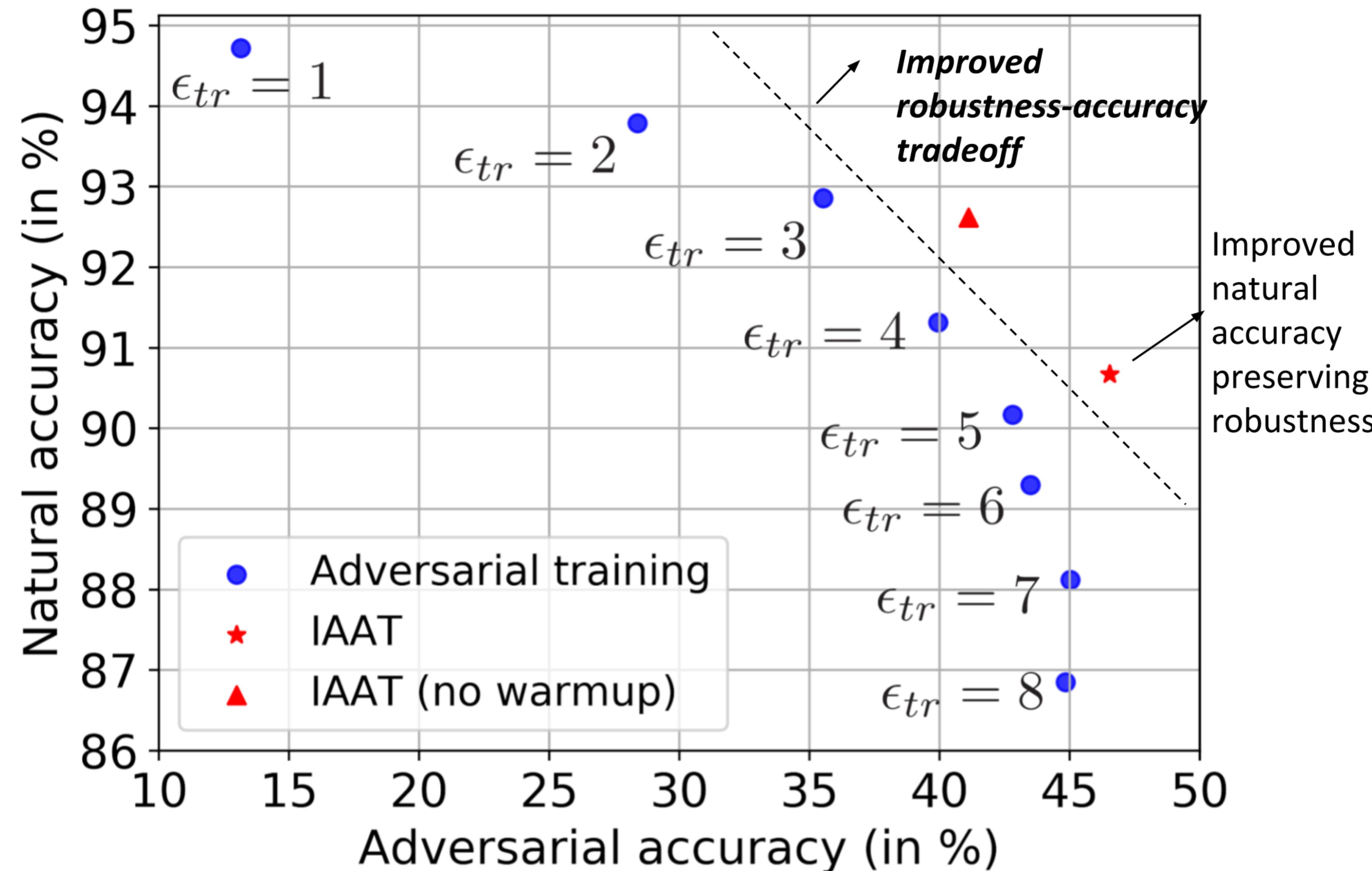
# Instance Adaptive Adversarial Training



# Instance Adaptive Adversarial Training



# Adaptive Adversarial Training: CIFAR-10

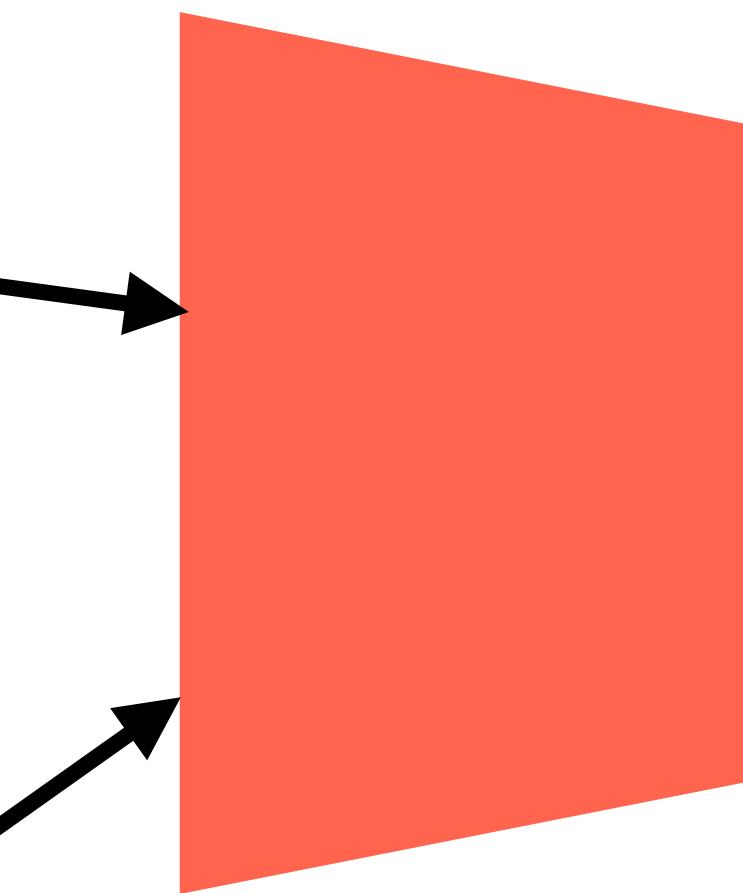


# Semi-supervised Learning

Labeled Data



Unlabeled Data

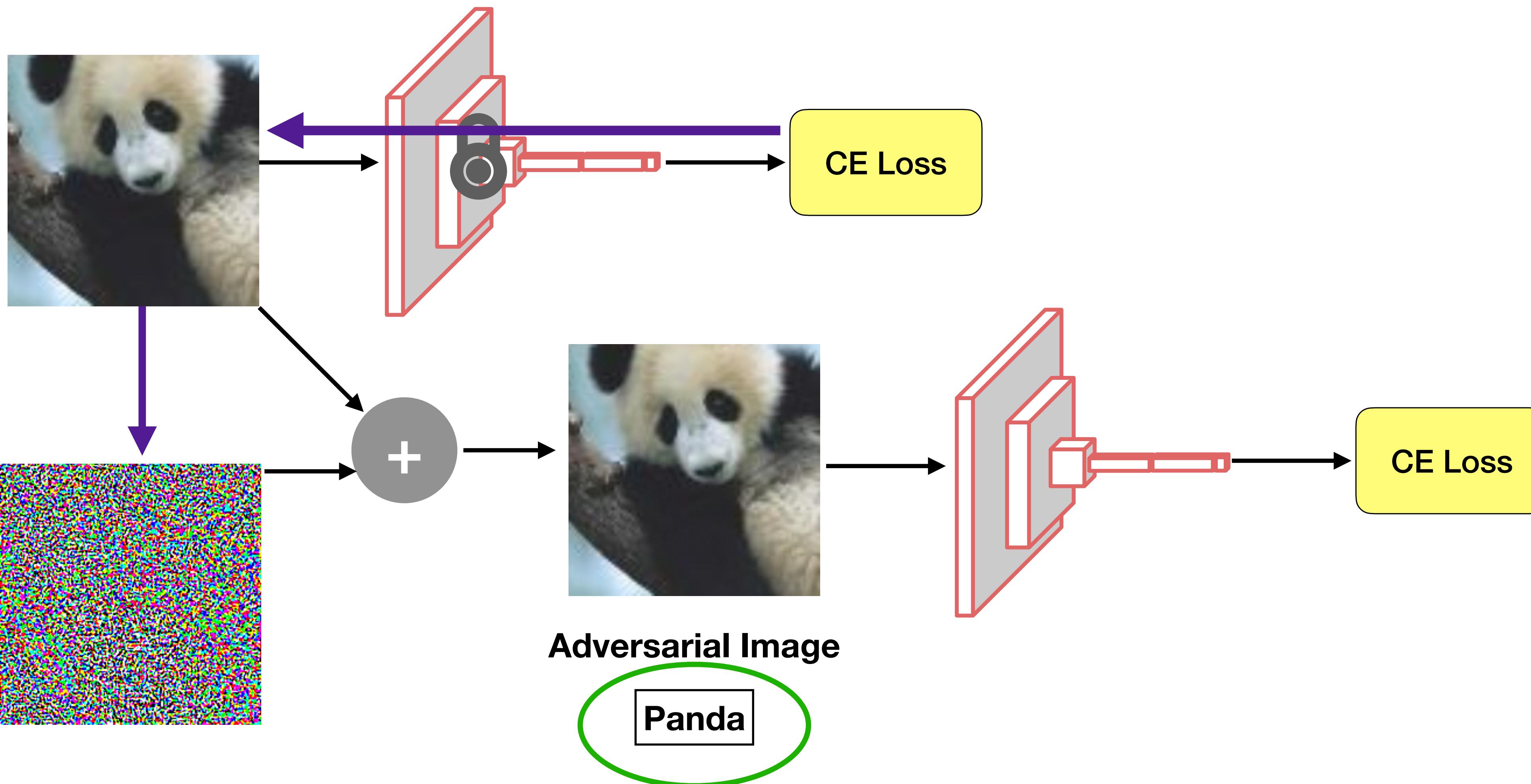


Supervised Loss

Unsupervised Loss

Leverages labeled data to  
“pseudo-label” unlabeled data

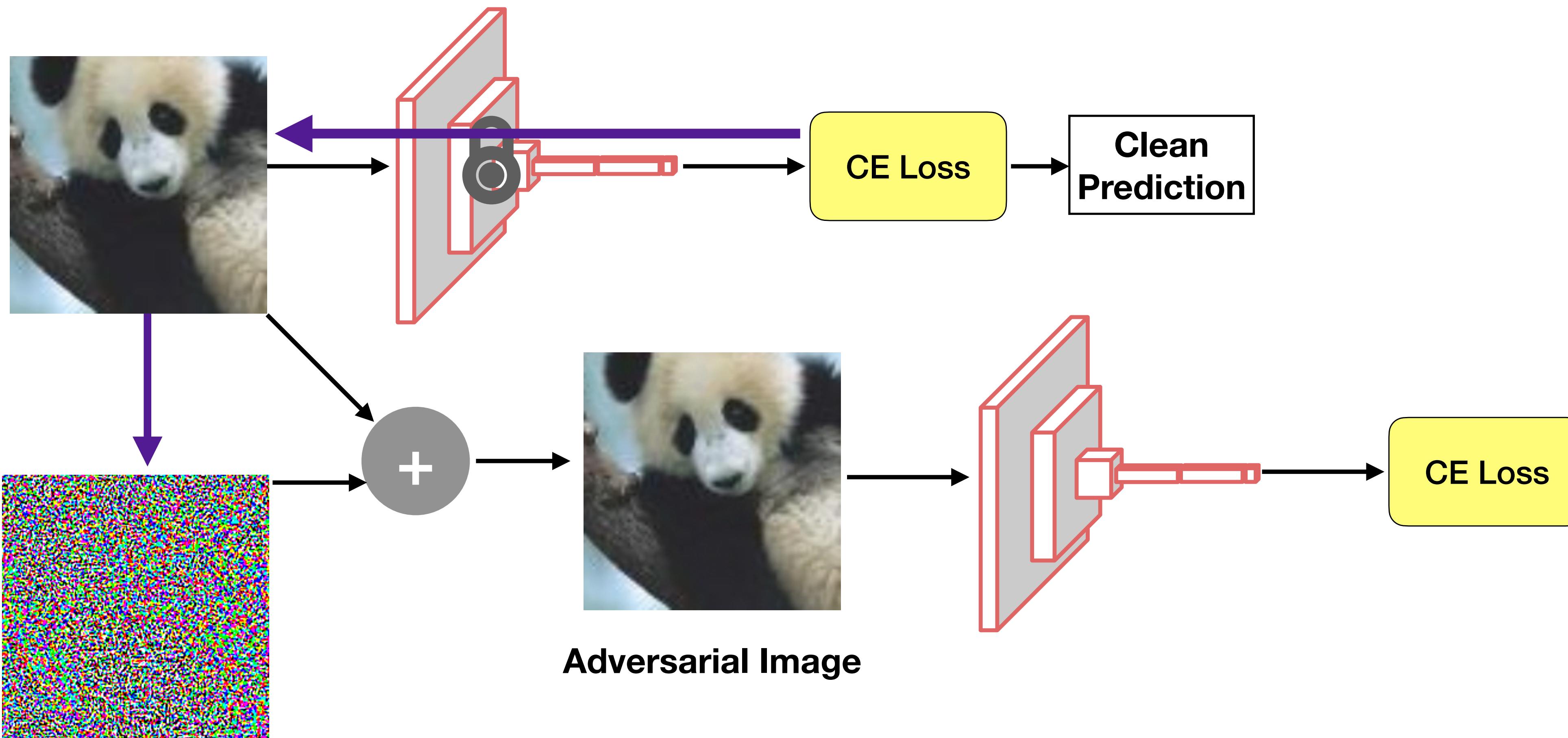
# Adversarial Training



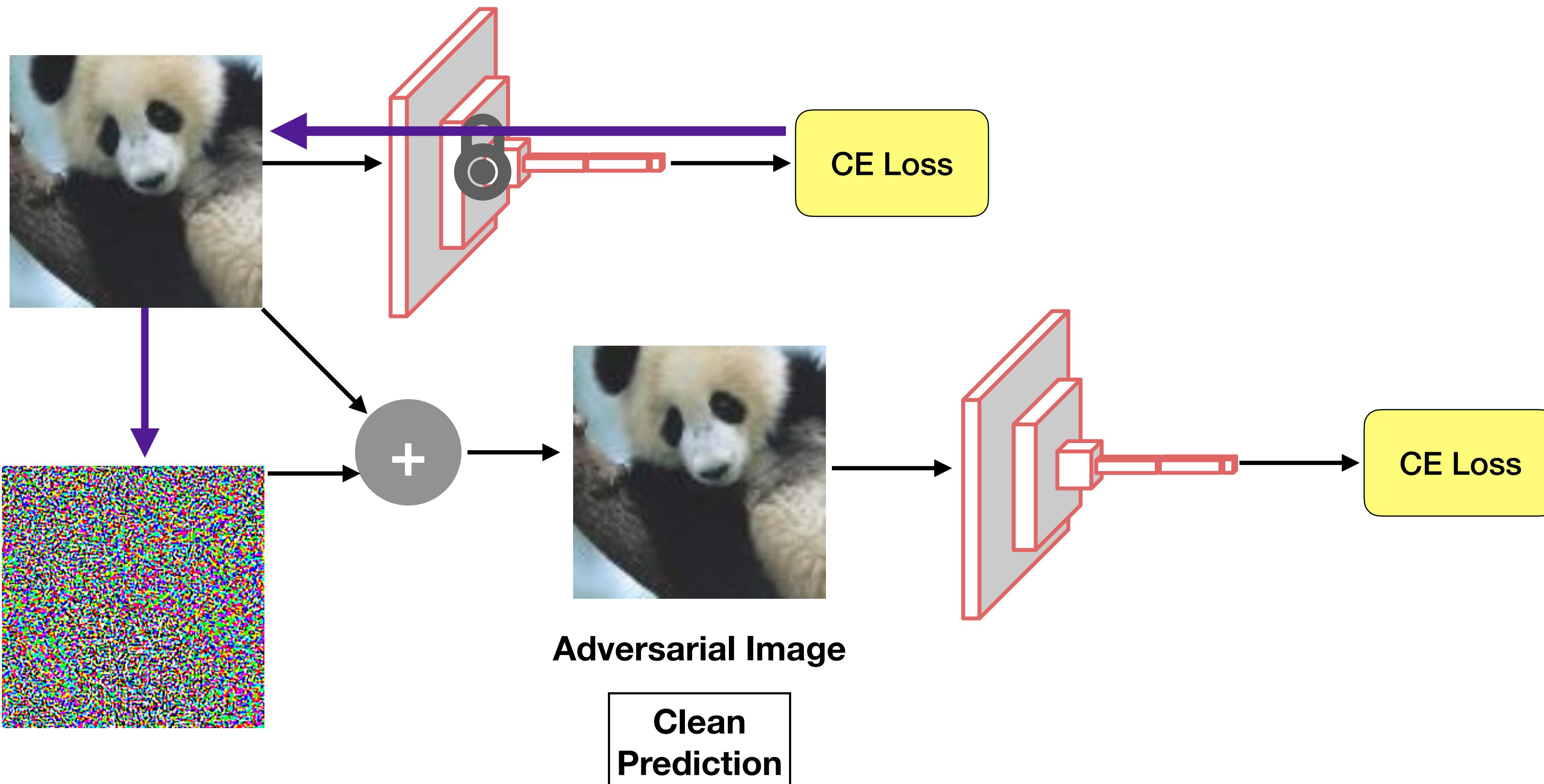
What if the label is unknown?

[Madry et. al. ICLR 18]

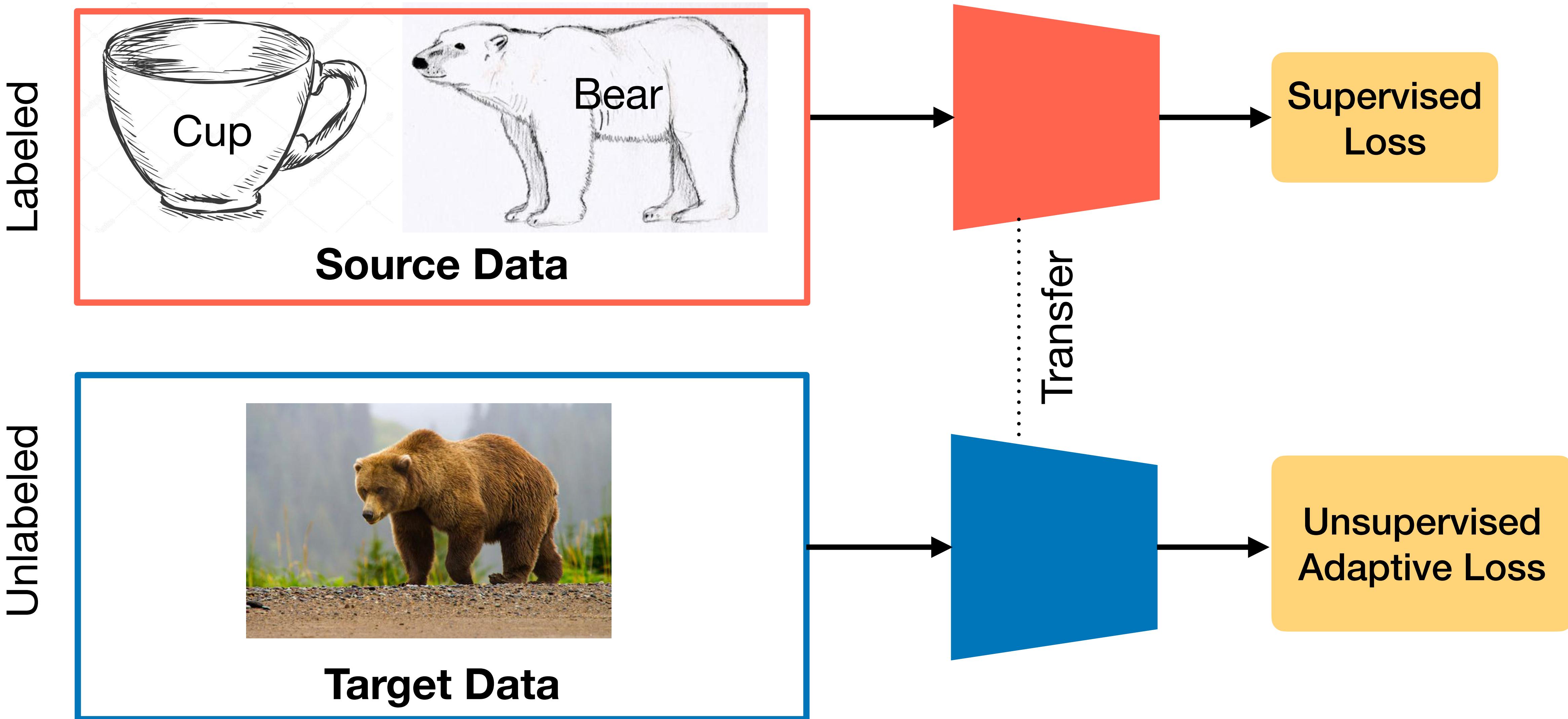
# Virtual Adversarial Training



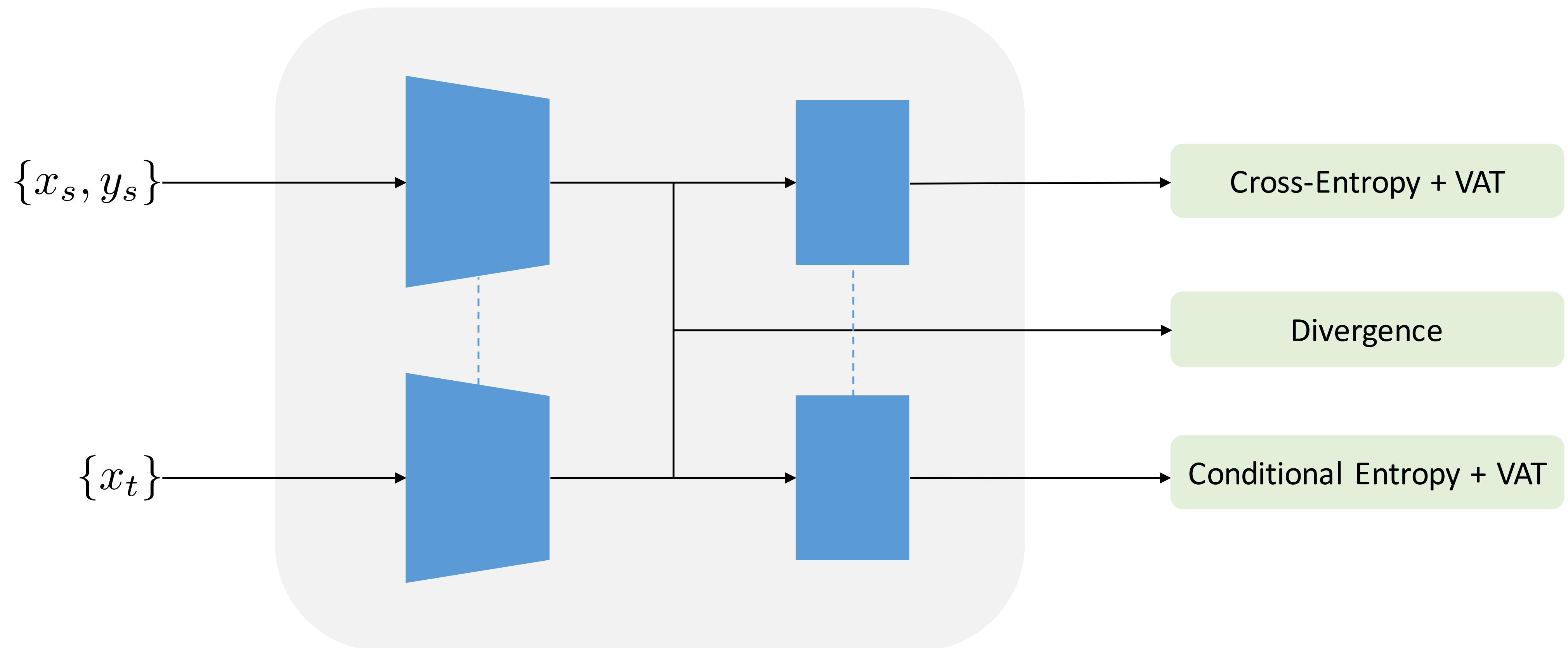
# Virtual Adversarial Training



# Unsupervised Domain Adaptation

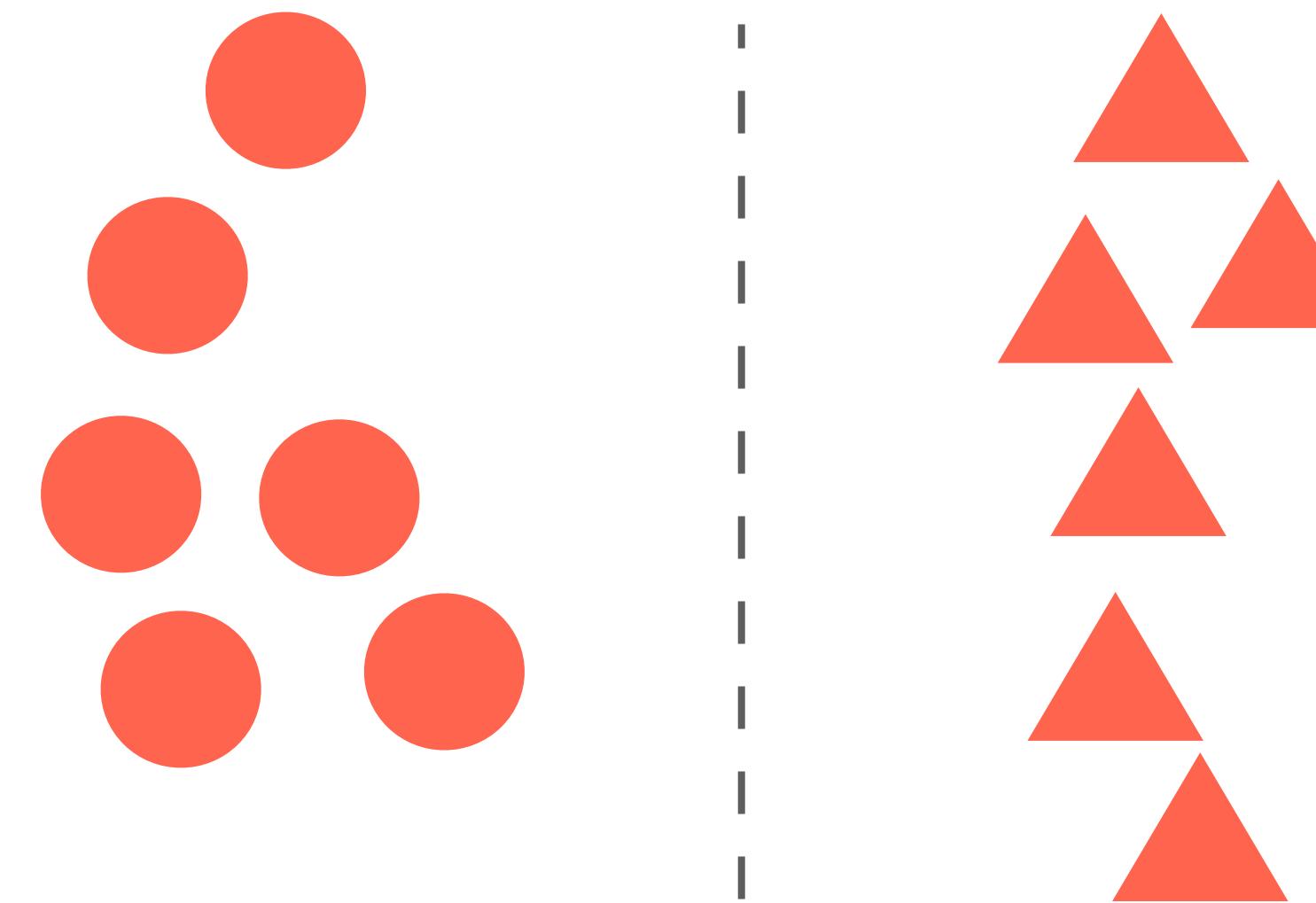


# DIRT-T / VADA



# Example: Entropy Minimization

Supervised Decision Boundary



Classes = { , }

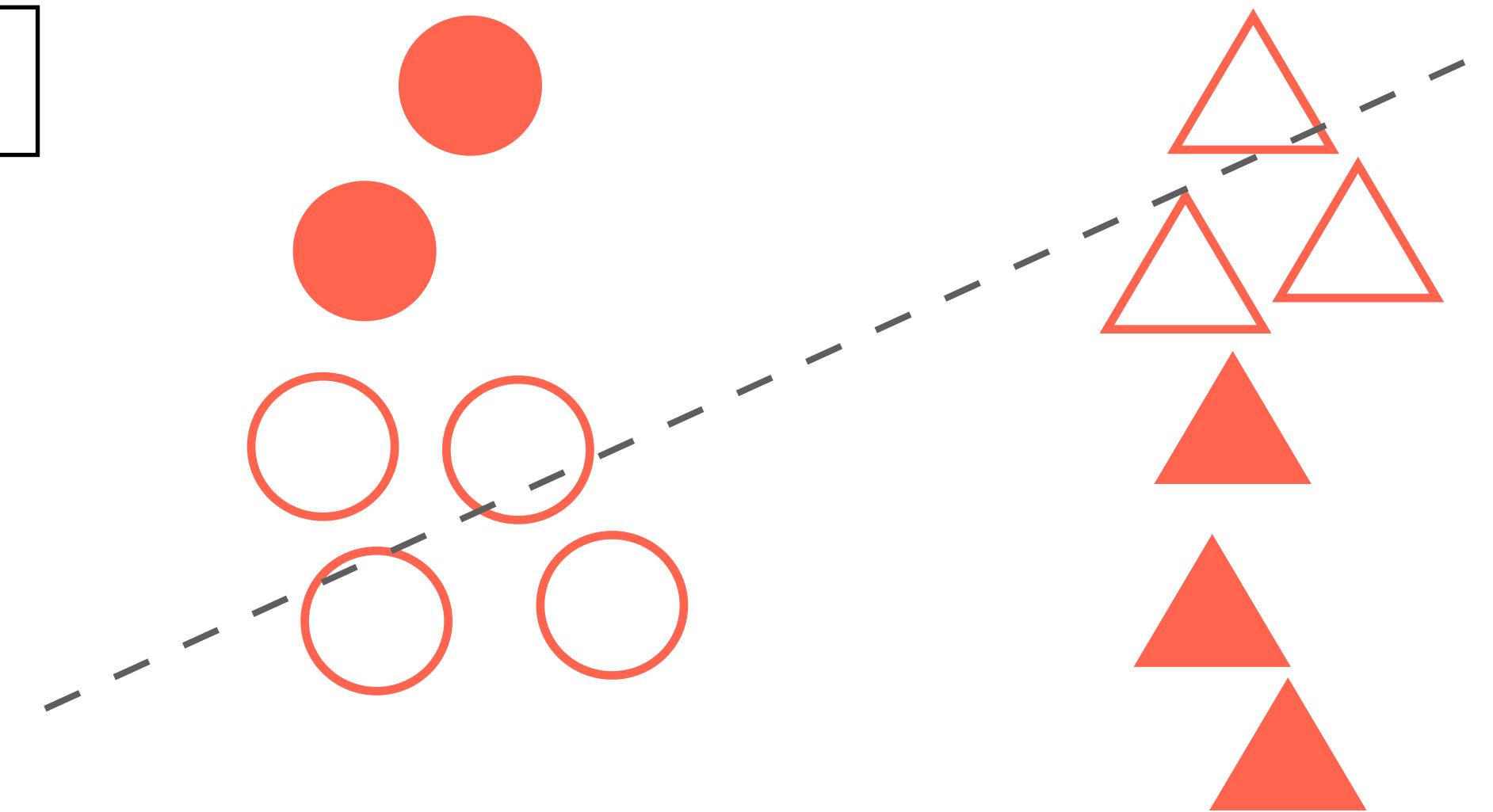
Labeled

Unlabeled

----- Class  
boundary

# Example: Entropy Minimization

## Error Accumulation



Classes = { , }

Labeled

Unlabeled

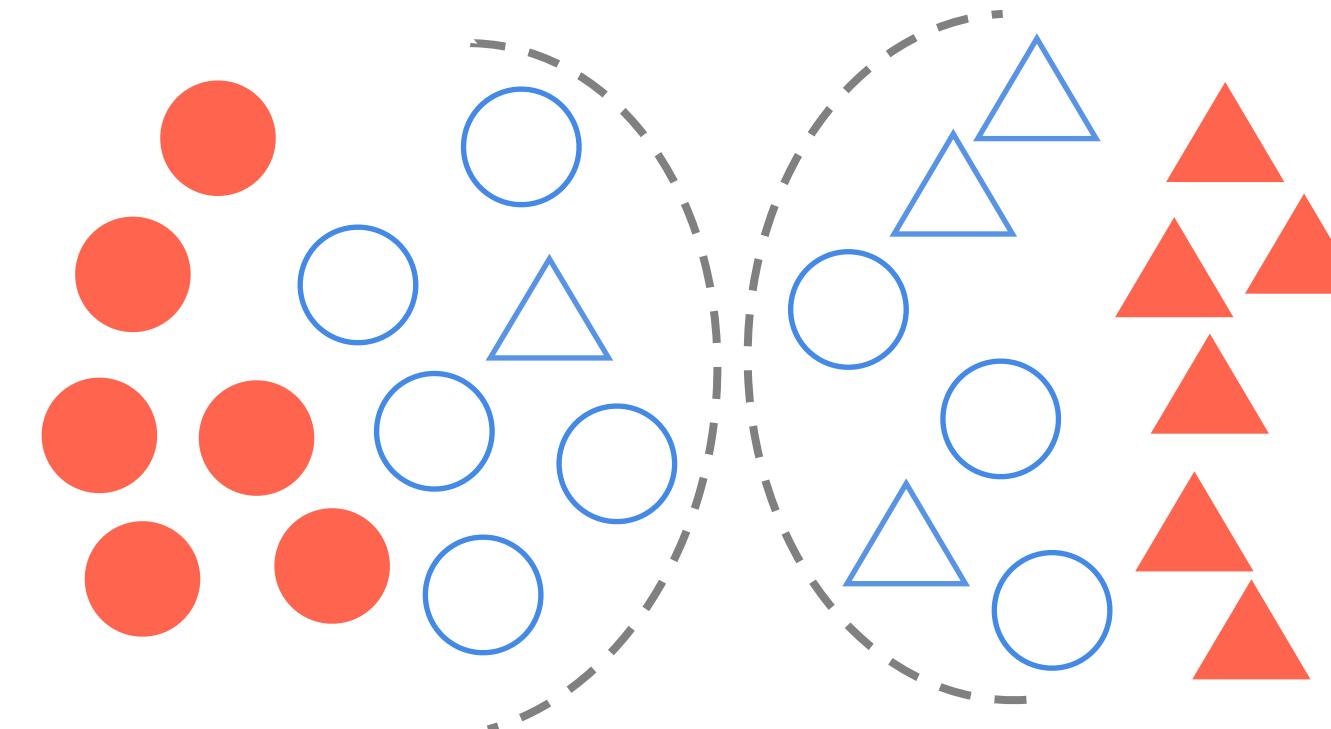
----- Class  
boundary

# Self-Training with Unreliable Instances

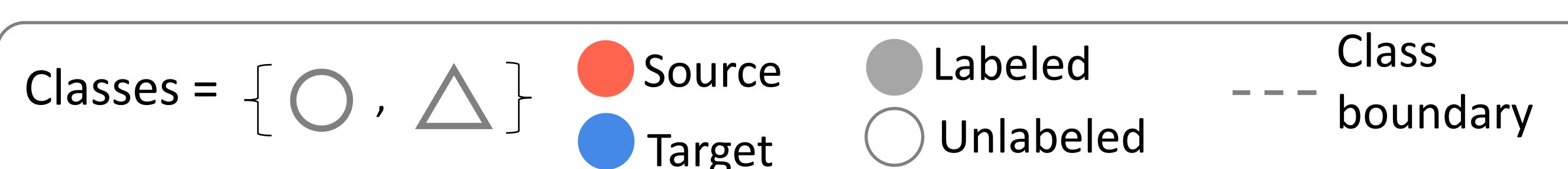
- Under a domain shift, some target categories may be misaligned
- Entropy minimization on such instances would increase model confidence, **reinforcing errors**

$$\begin{aligned}\mathcal{L}_{CEM} &= \mathbb{E}_{\mathbf{x}_T \sim \mathcal{P}_T} [\mathcal{H}_\Theta(y \mid \mathbf{x}_T)] \\ &= \mathbb{E}_{\mathbf{x}_T \sim \mathcal{P}_T} \left[ \sum_{c=1}^C -p_\Theta(y = c \mid \mathbf{x}_T) \log p_\Theta(y = c \mid \mathbf{x}_T) \right]\end{aligned}$$

Entropy Minimization for UDA



Poor Initialization

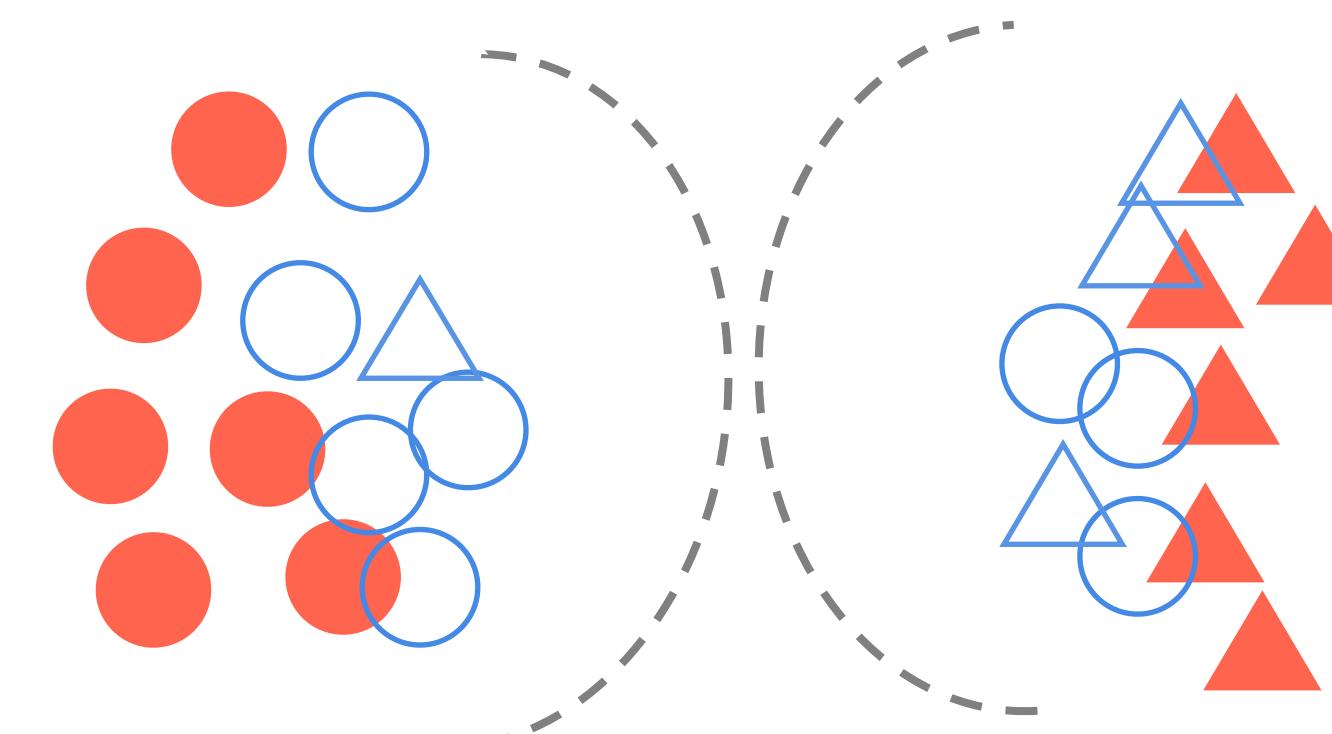


# Self-Training with Unreliable Instances

- Under a domain shift, some target categories may be misaligned
- Entropy minimization on such instances would increase model confidence, **reinforcing errors**

$$\begin{aligned}\mathcal{L}_{CEM} &= \mathbb{E}_{\mathbf{x}_T \sim \mathcal{P}_T} [\mathcal{H}_\Theta(y \mid \mathbf{x}_T)] \\ &= \mathbb{E}_{\mathbf{x}_T \sim \mathcal{P}_T} \left[ \sum_{c=1}^C -p_\Theta(y = c \mid \mathbf{x}_T) \log p_\Theta(y = c \mid \mathbf{x}_T) \right]\end{aligned}$$

Entropy Minimization for UDA

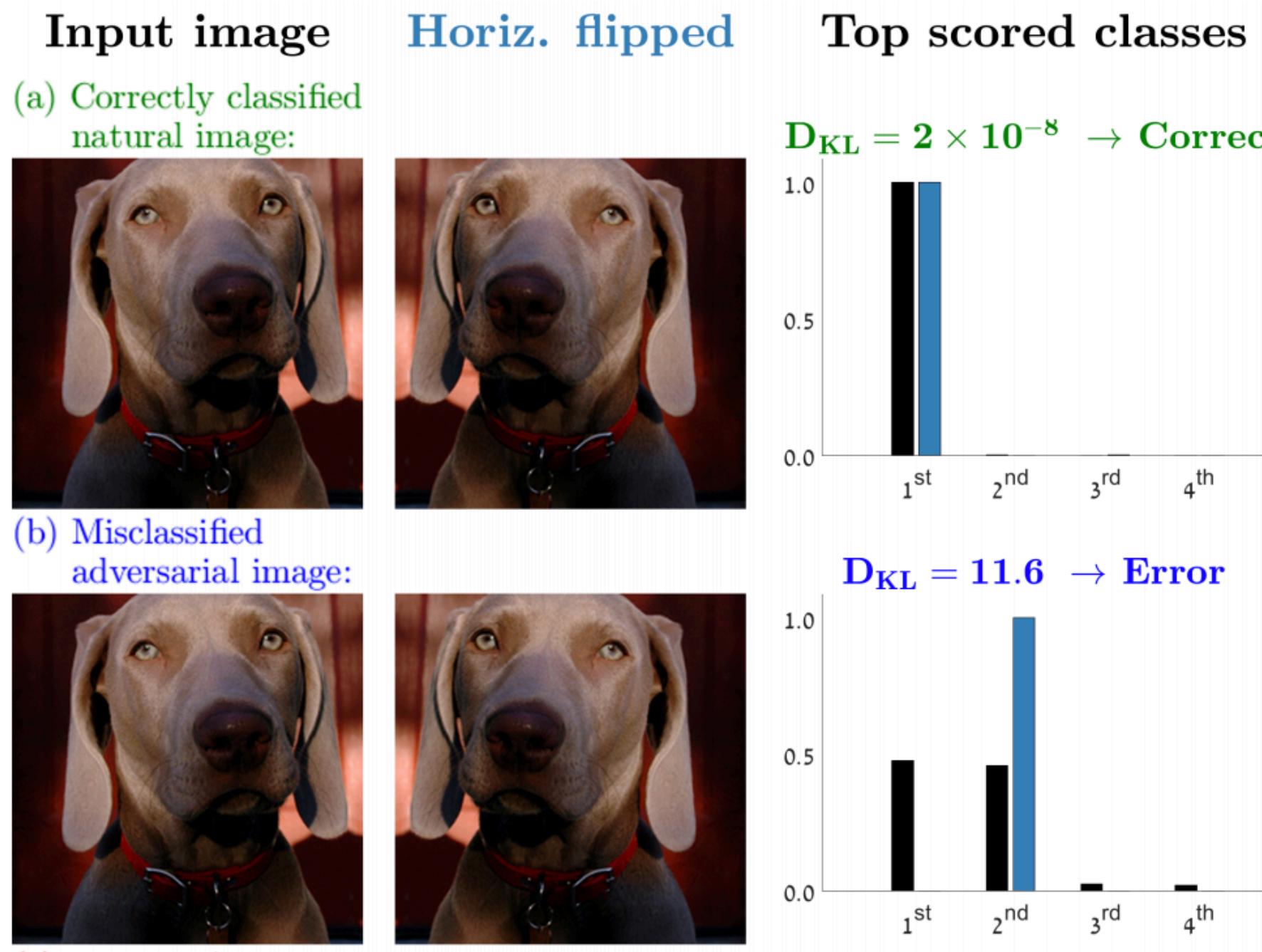


Classes = {○, △}

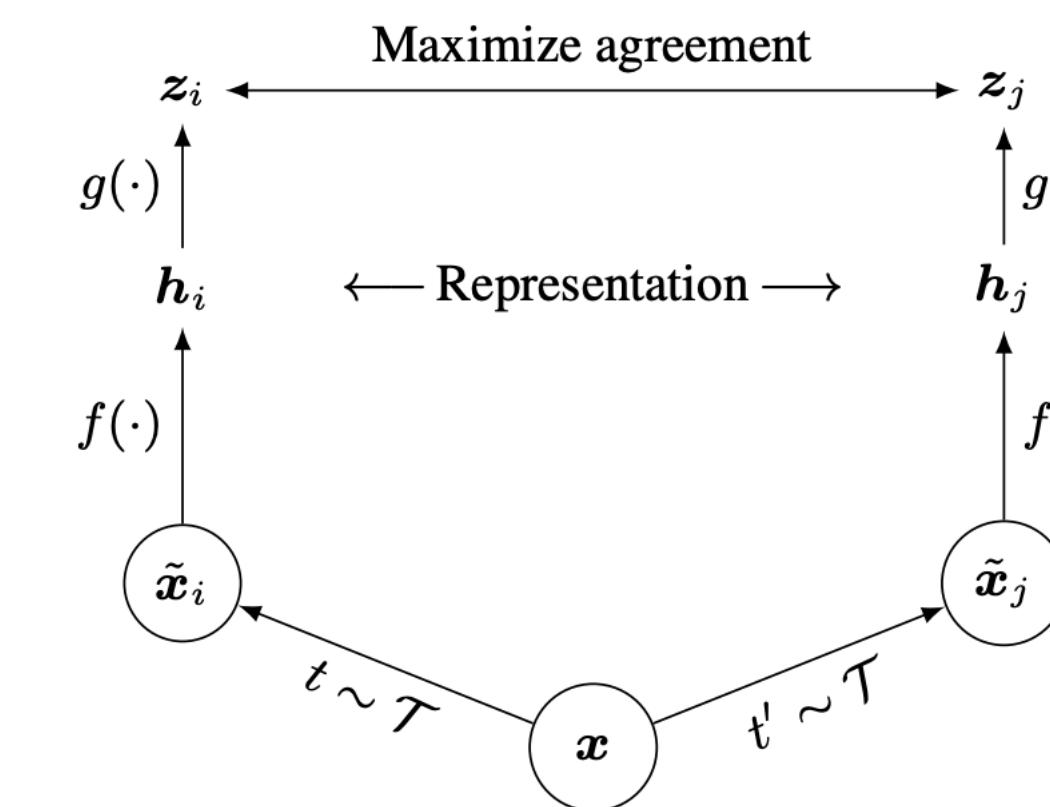
Source	●
Target	●
Labeled	●
Unlabeled	○

Class boundary

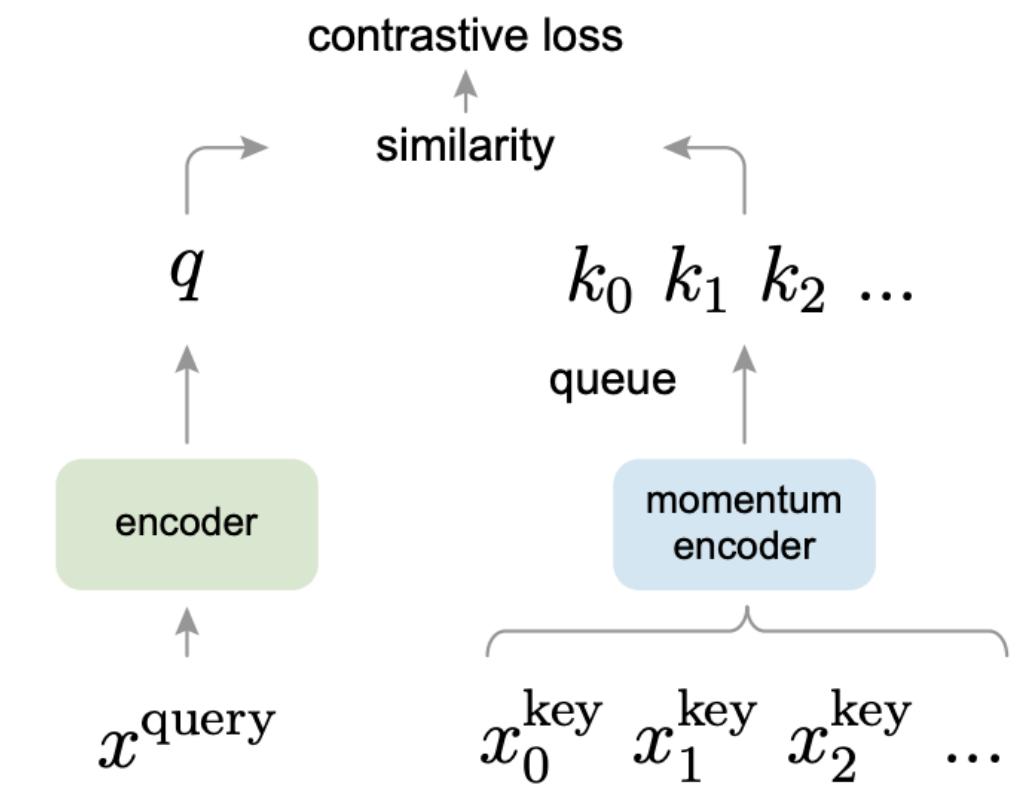
# Prior Work: Measure Image Aug Differences



Natural and Adversarial Error Detection using Invariance to Image Transformations. Bahat, Irani, Shakhnarovich, arXiv 2019



SimCLR, Chen et al.  
ICML 2020



MoCo, He et al.  
CVPR 2020

## Detecting Errors

## Learned Invariance (Contrastive Learning)

# **SENTRY**

## Selective Entropy Optimization via Committee Consistency for Unsupervised Domain Adaptation



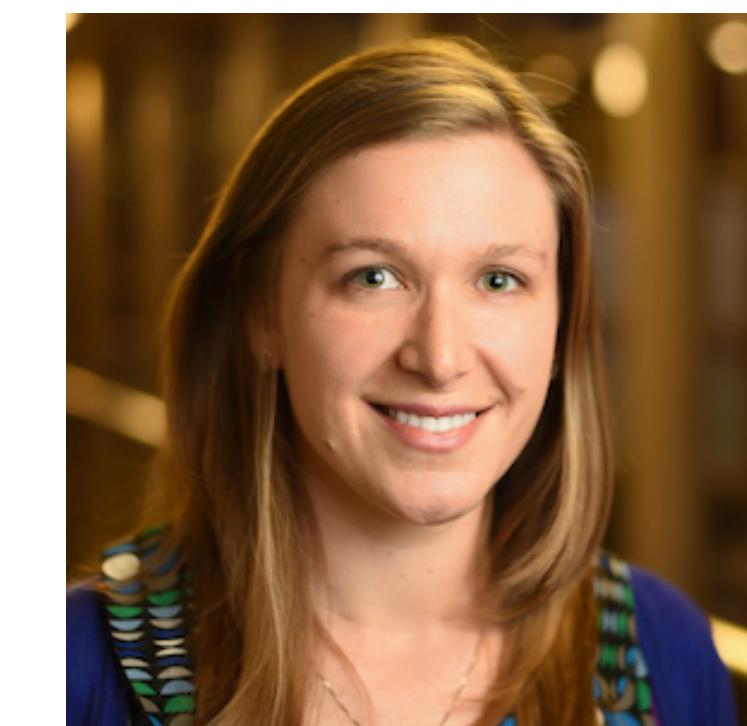
Viraj Prabhu



Shivam Khare



Deeksha Karthik



Judy Hoffman

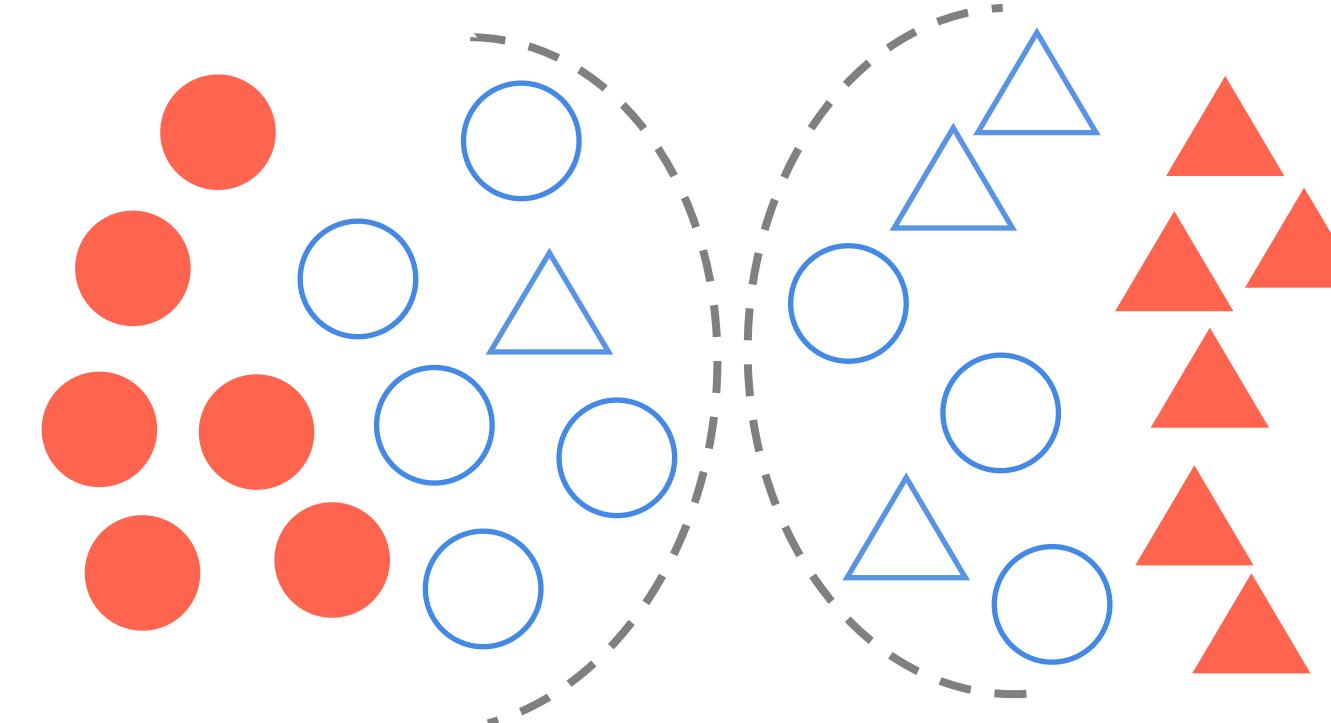


# SENTRY: Selective Entropy Optimization via Committee Consistency

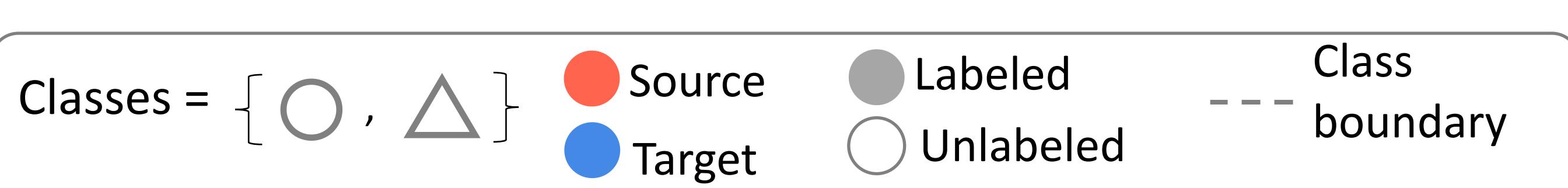
## Key Idea

1. Identify reliable instances using predictive consistency<sup>1,2,3</sup>
  - Model confidence is known to be uncalibrated under distribution shift [Ovadia NeurIPS 2019]
2. Increase model confidence on highly consistent target instances, reduce on inconsistent

### Selective Entropy Minimization



Poor Initialization



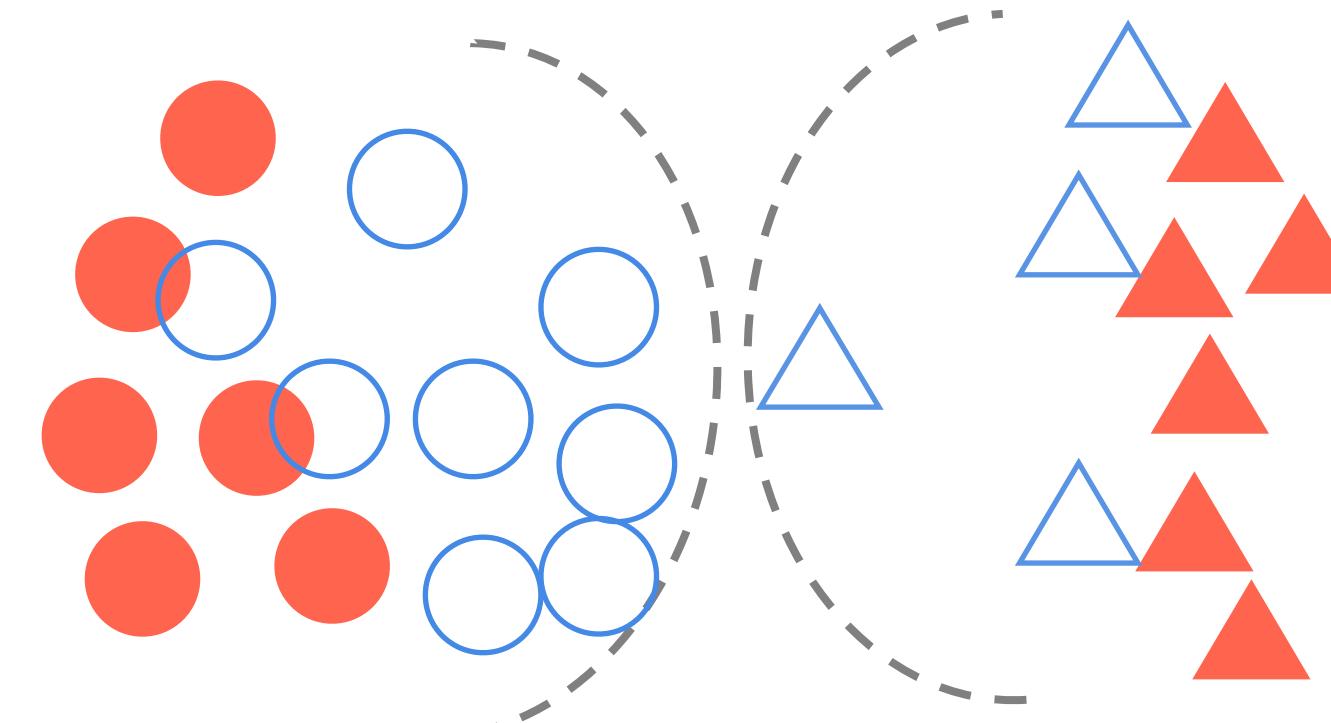
1. Bahat et al. arXiv 2019.
2. Chen et al. ICML 2020.
3. Sohn et al., NeurIPS 2020.

# SENTRY: Selective Entropy Optimization via Committee Consistency

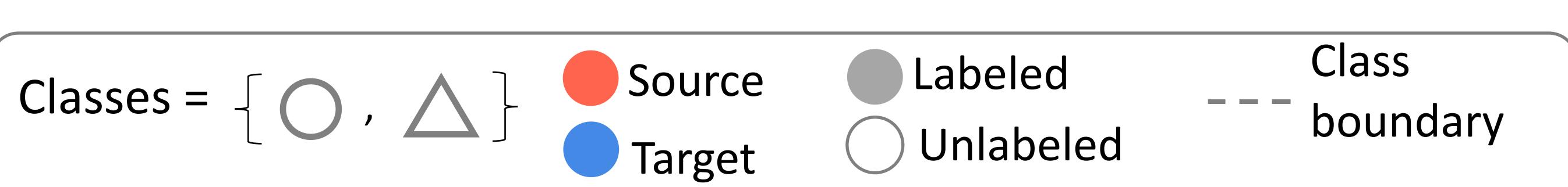
## Key Idea

1. Identify reliable instances using predictive consistency<sup>1,2,3</sup>
  - Model confidence is known to be uncalibrated under distribution shift [Ovadia NeurIPS 2019]
2. Increase model confidence on highly consistent target instances, reduce on inconsistent

### Selective Entropy Minimization



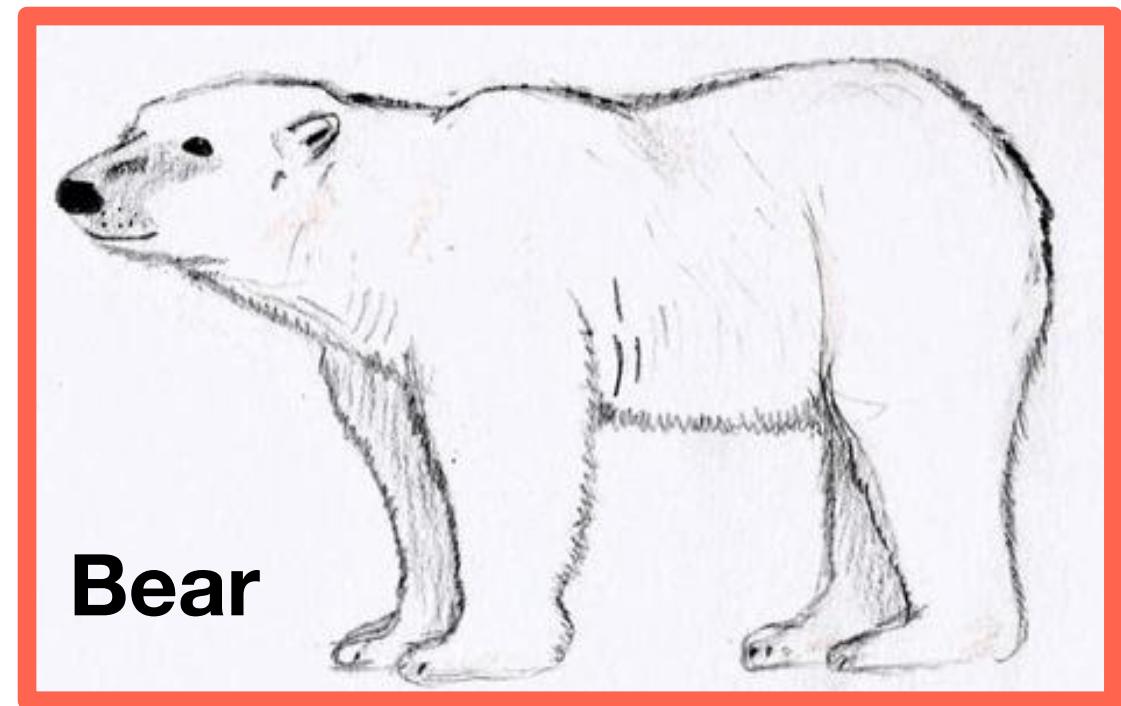
Poor Initialization



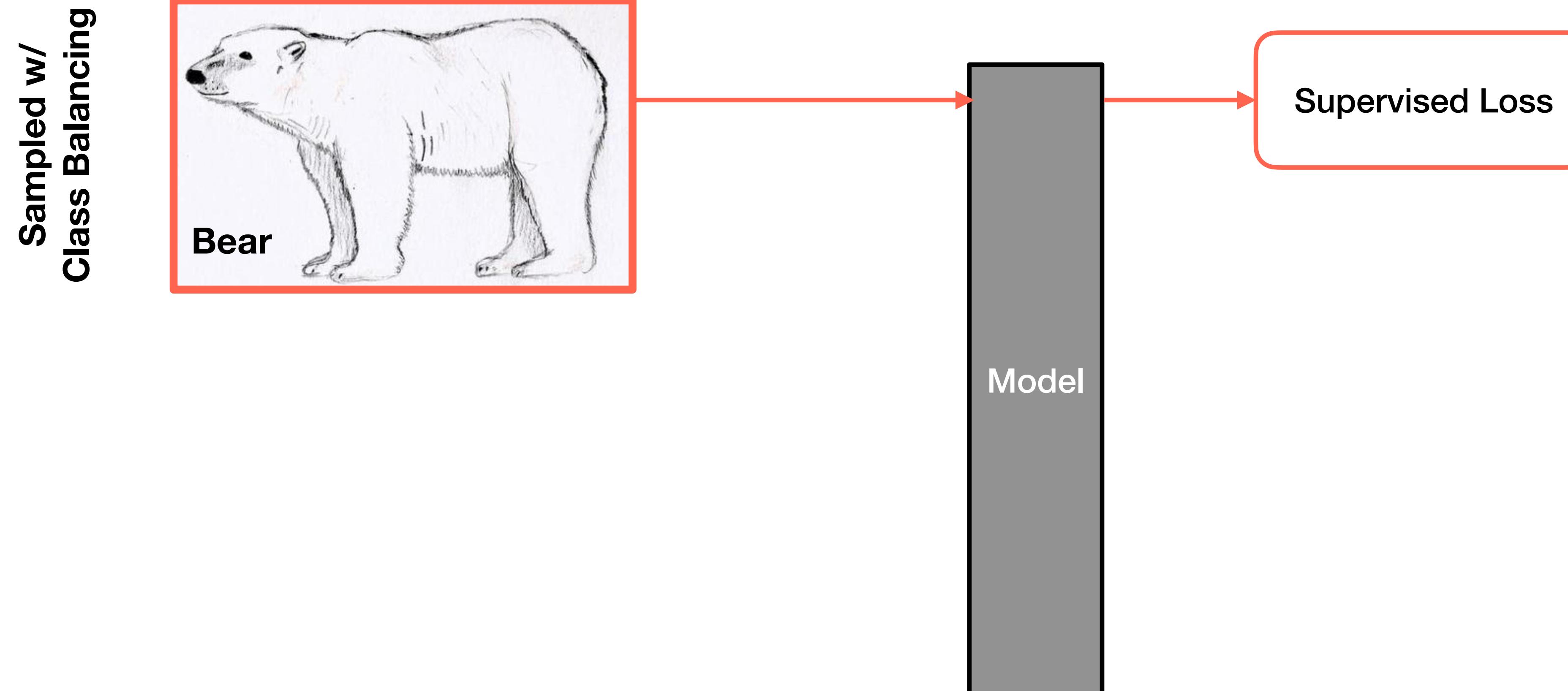
1. Bahat et al. arXiv 2019.
2. Chen et al. ICML 2020.
3. Sohn et al., NeurIPS 2020.

# SENTRY: Selective Entropy Optimization via Committee Consistency

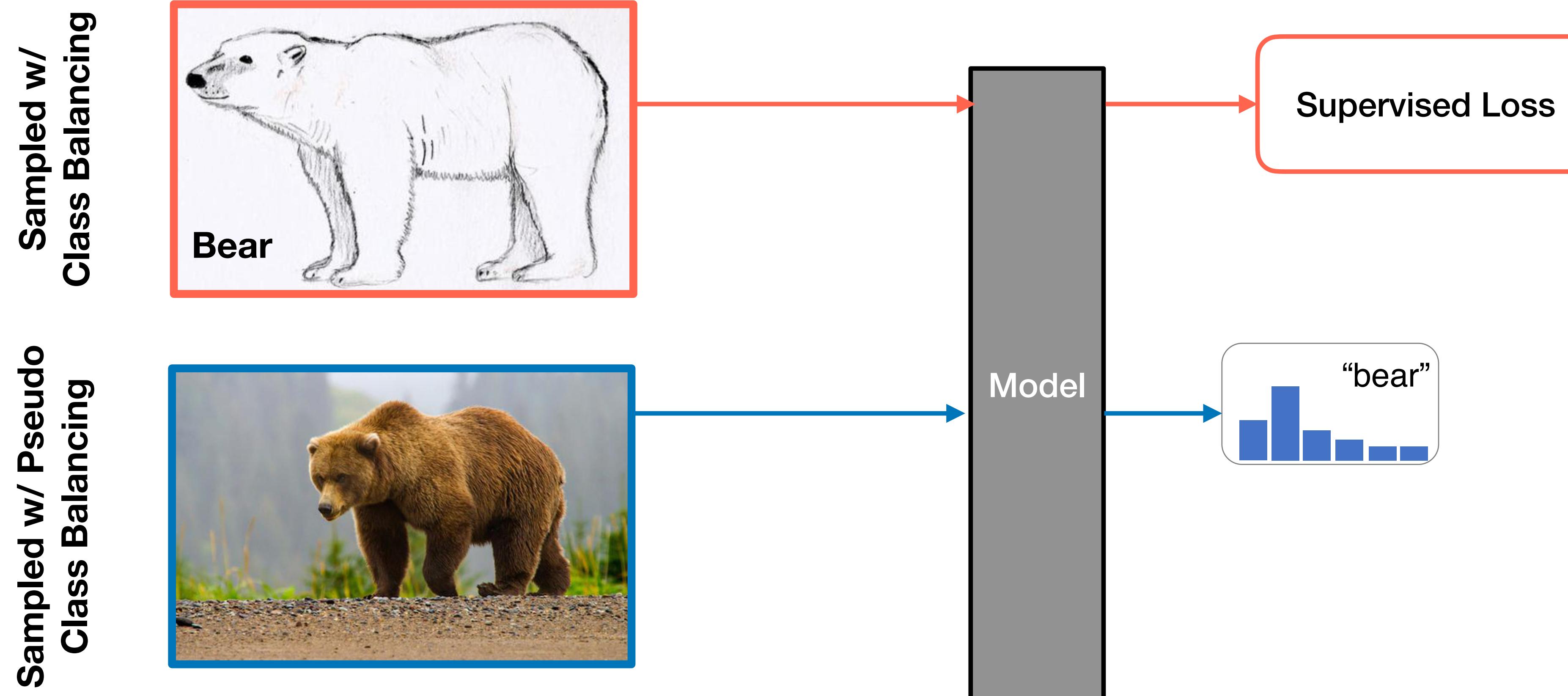
Sampled w/  
Class Balancing



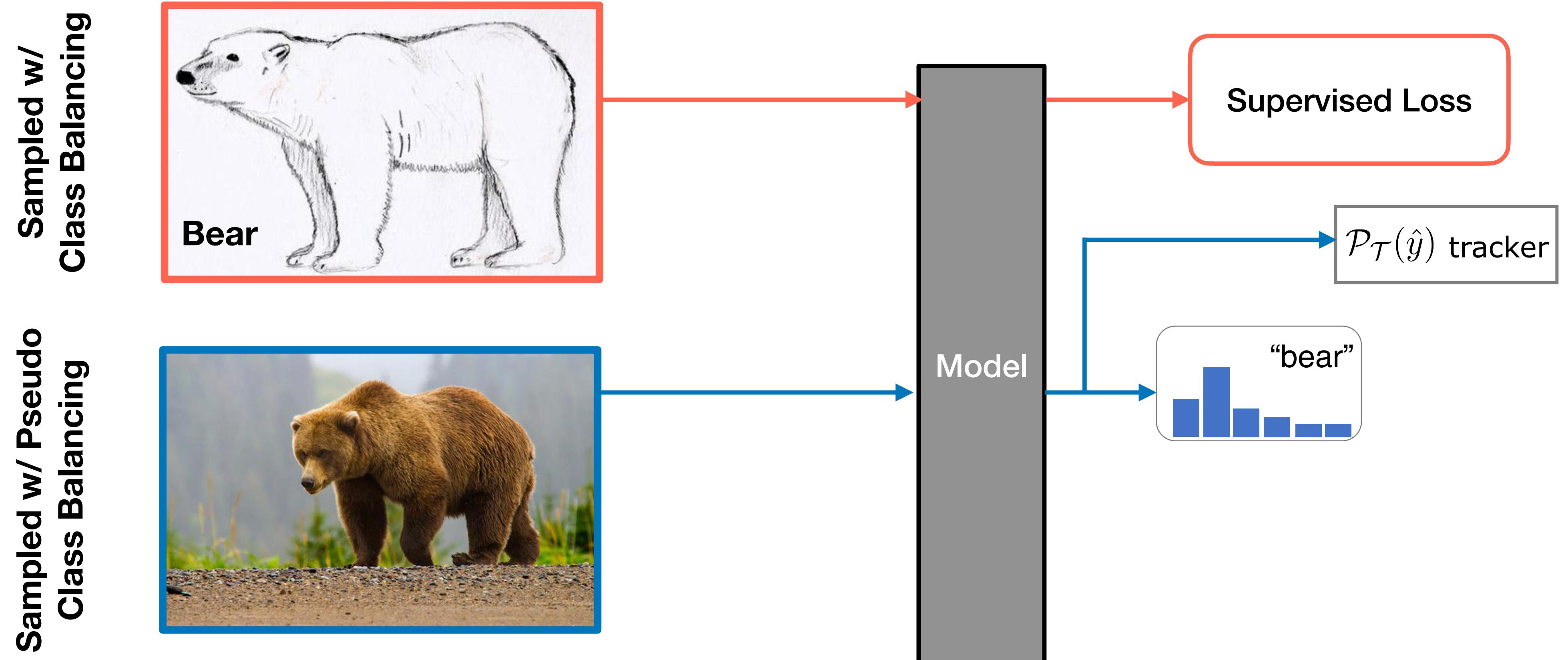
# SENTRY: Selective Entropy Optimization via Committee Consistency



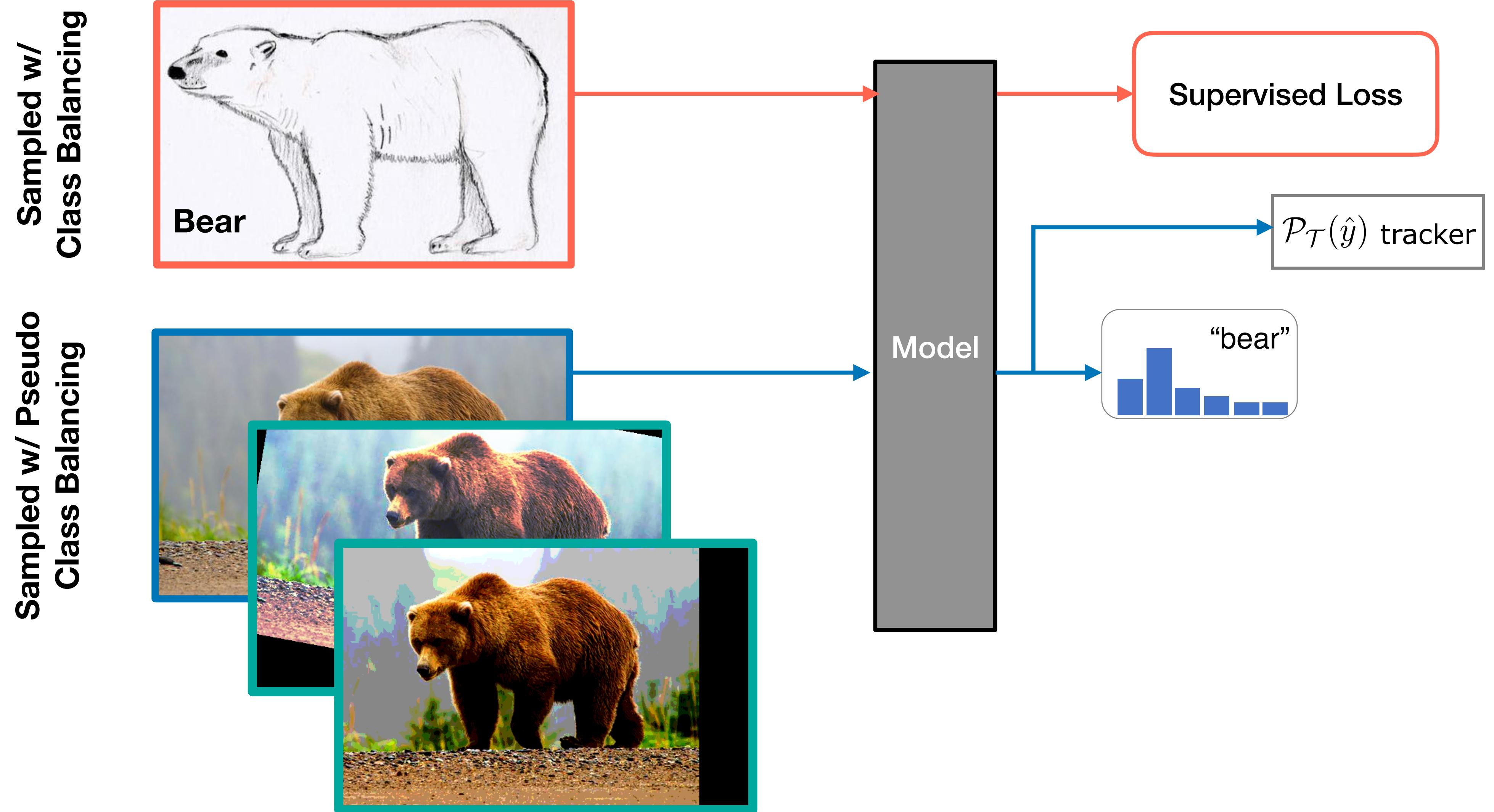
# SENTRY: Selective Entropy Optimization via Committee Consistency



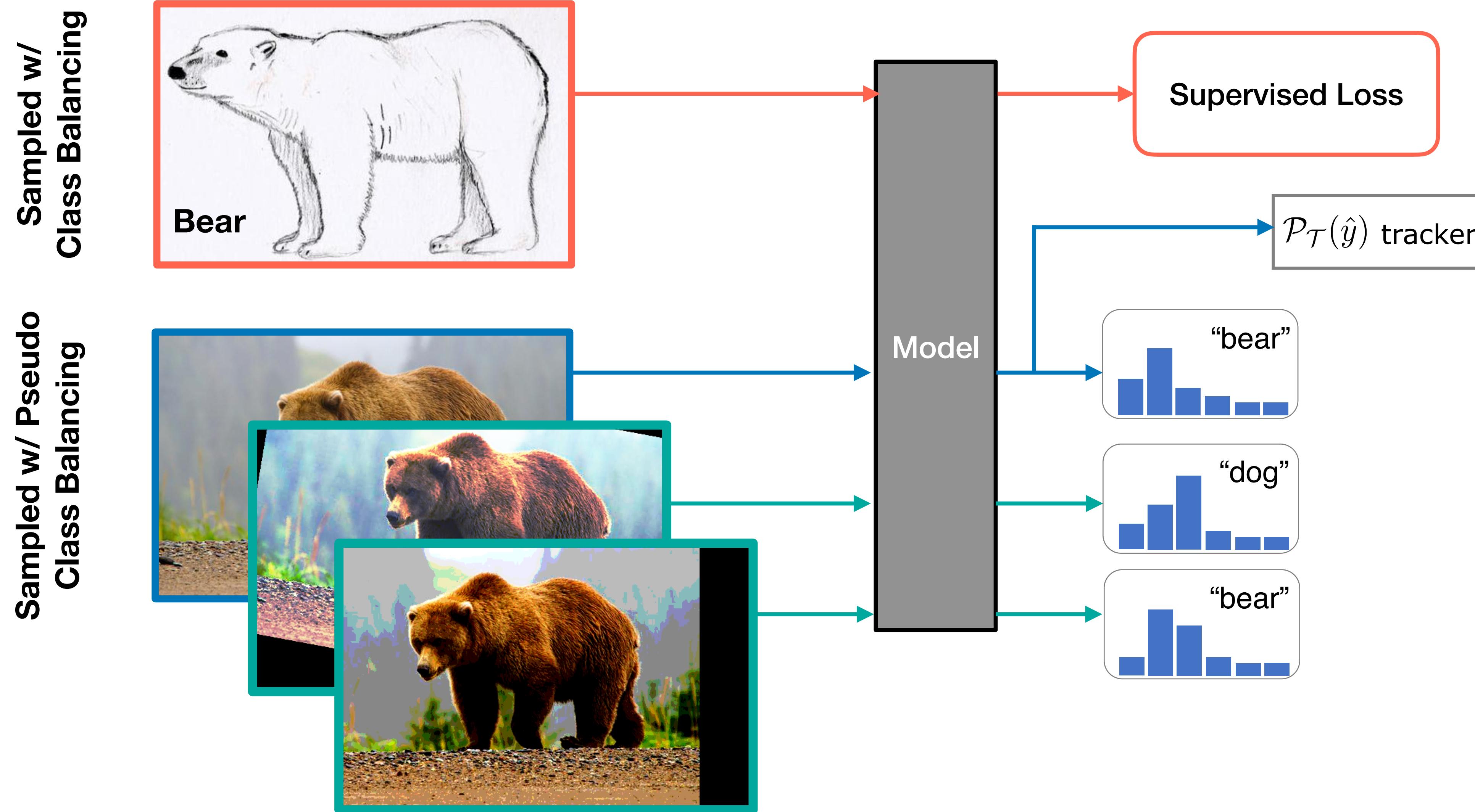
# SENTRY: Selective Entropy Optimization via Committee Consistency



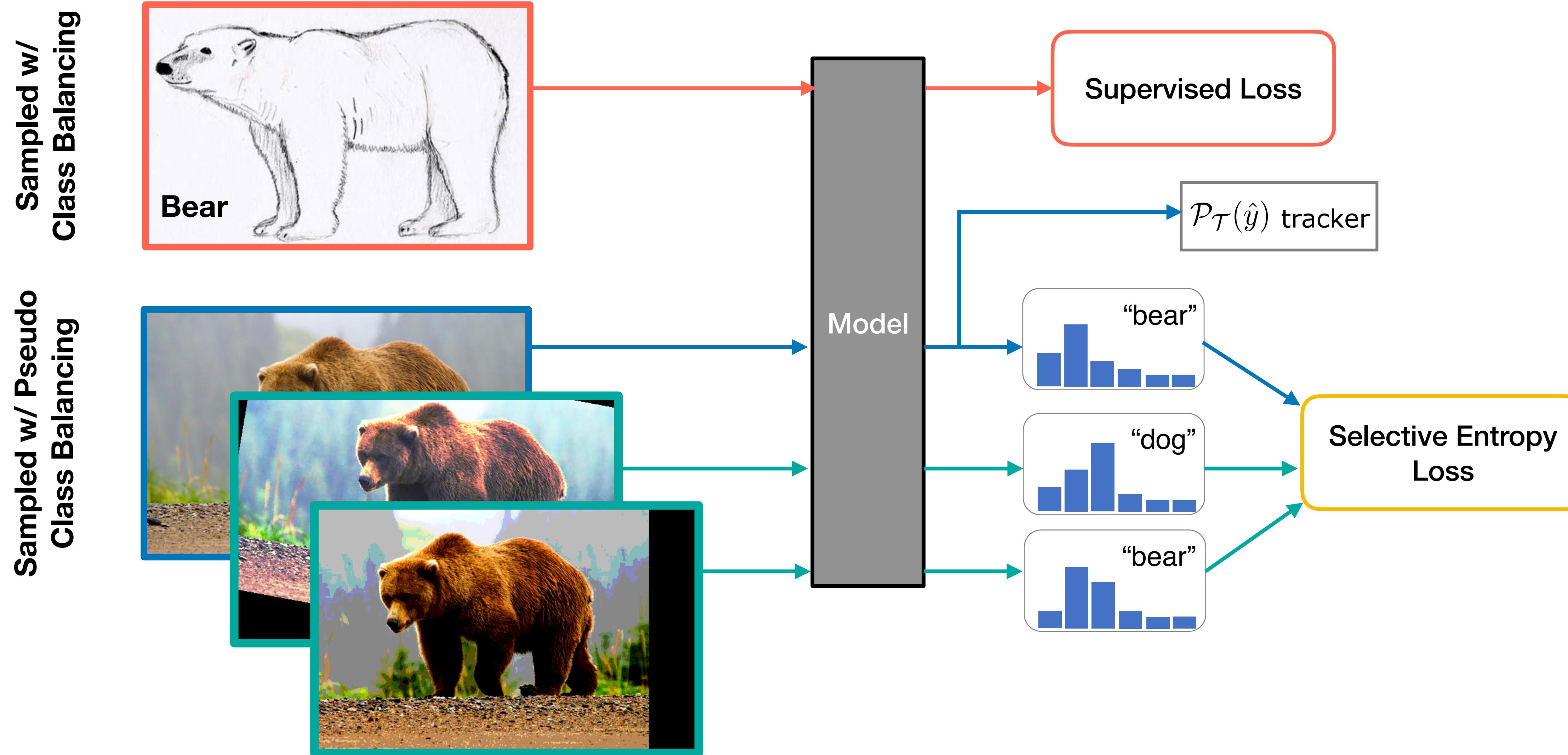
# SENTRY: Selective Entropy Optimization via Committee Consistency



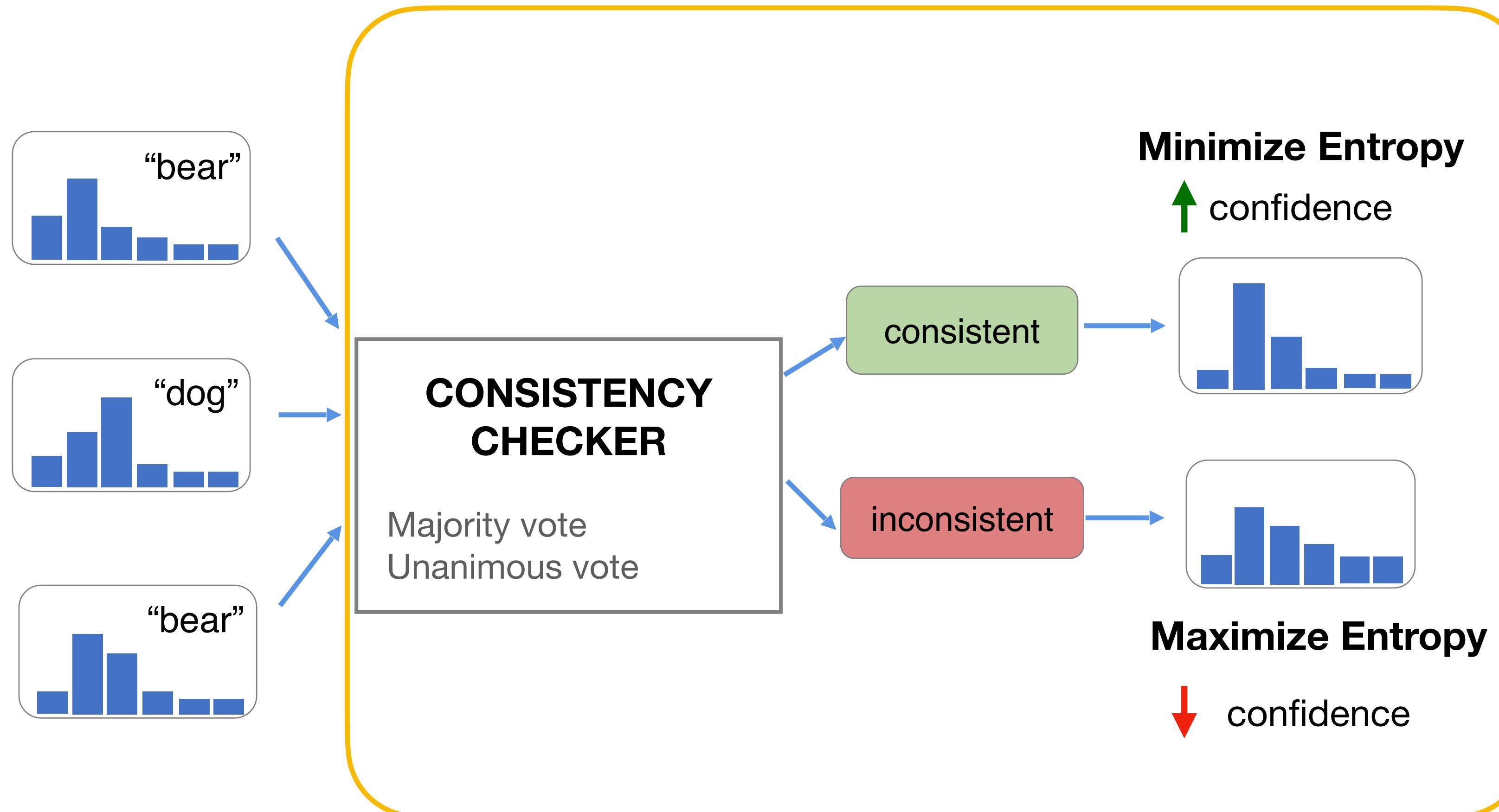
# SENTRY: Selective Entropy Optimization via Committee Consistency



# SENTRY: Selective Entropy Optimization via Committee Consistency



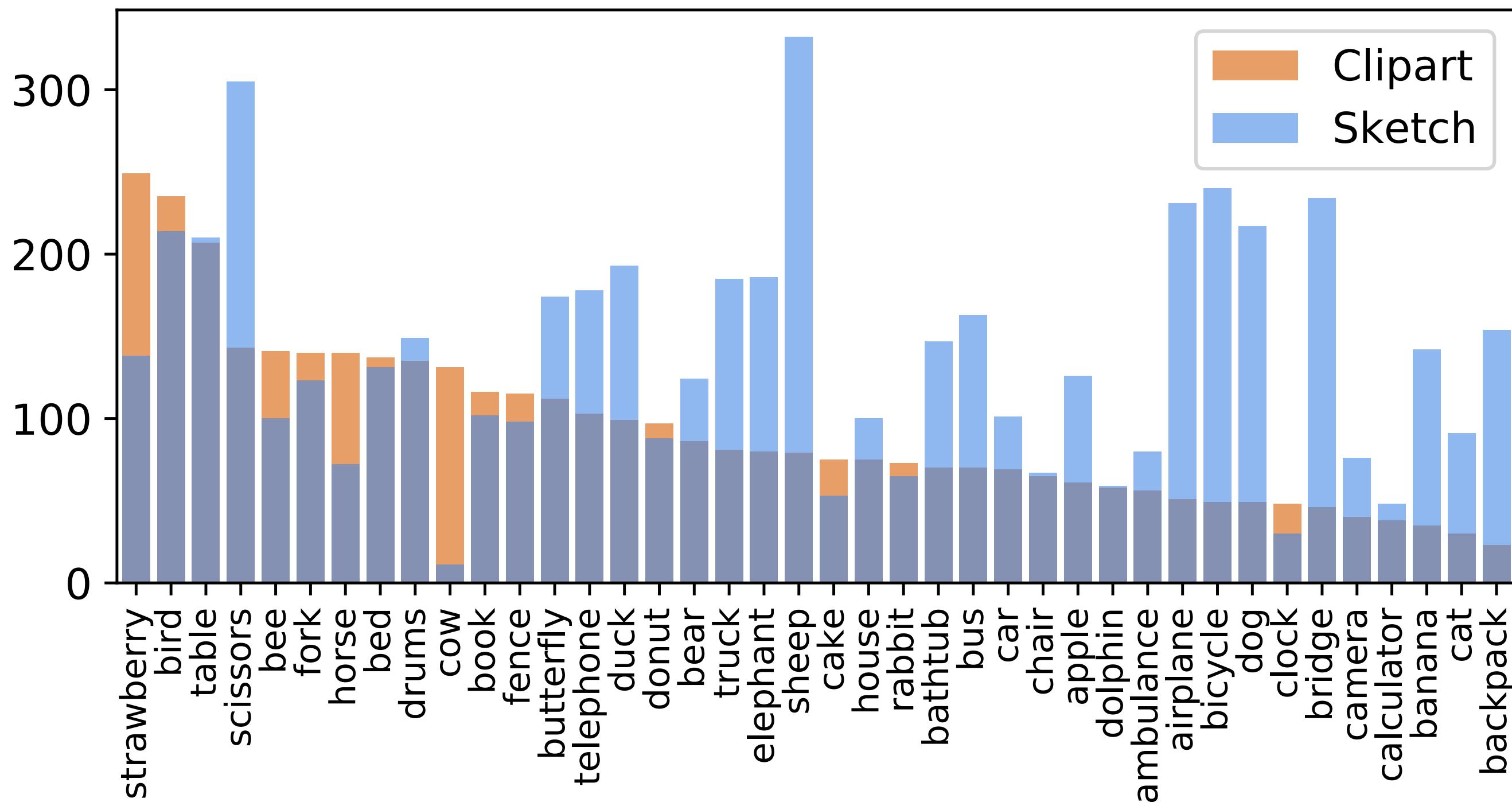
# Selective Entropy Loss



# SENTRY Results: Image Classification

Natural label shifts

DomainNet Label Histogram: clipart to sketch

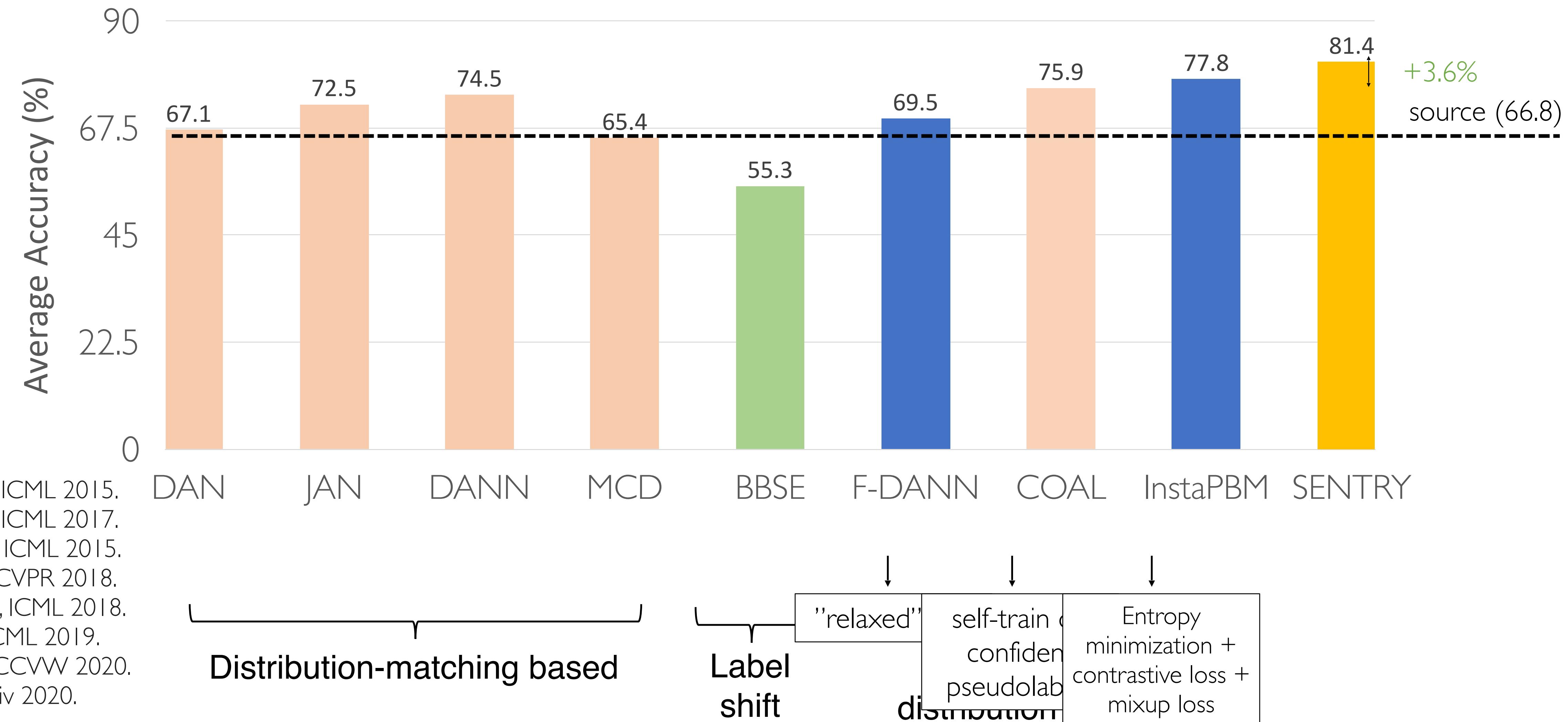


MiniDomainNet<sup>1,2</sup>



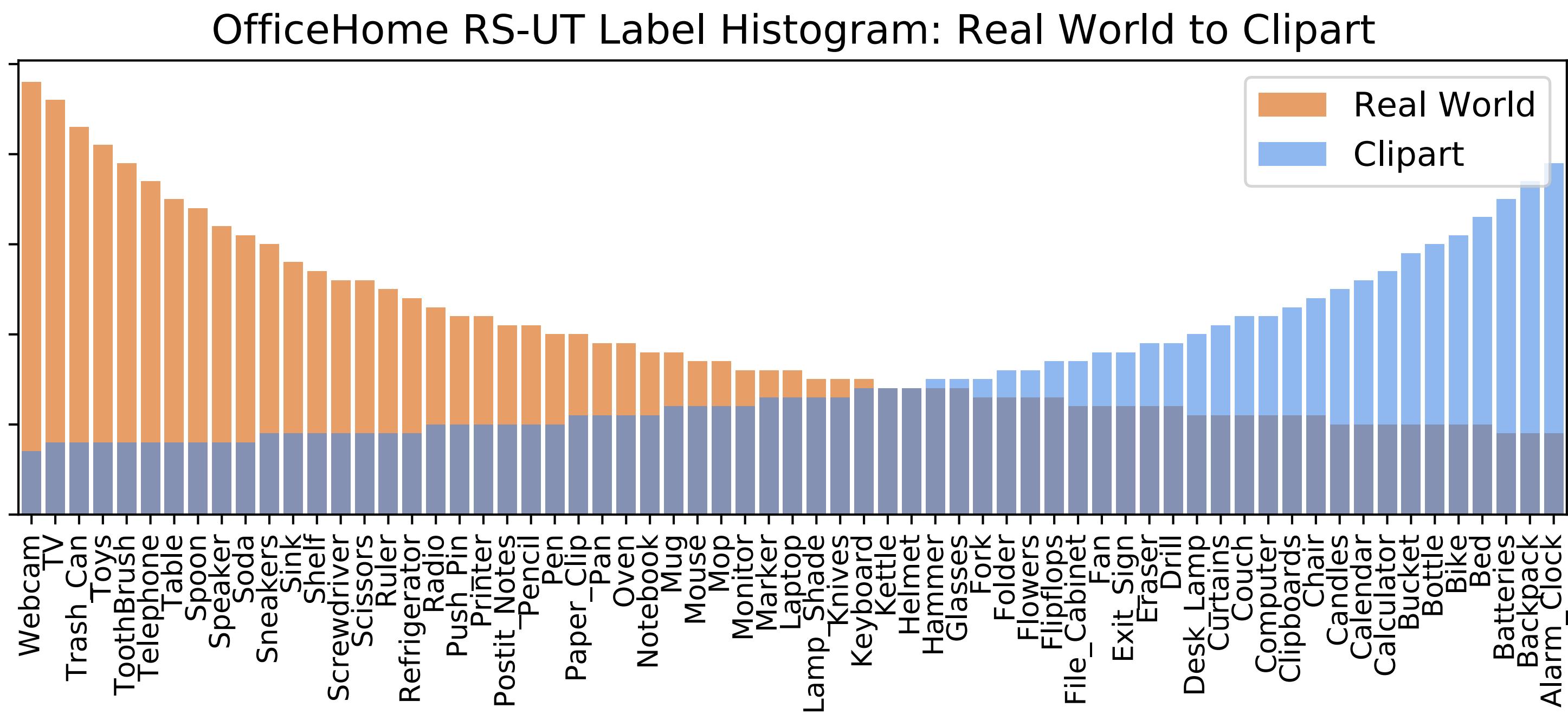
# SENTRY Results: MiniDomainNet

MiniDomainNet (40 classes, 12 shifts)

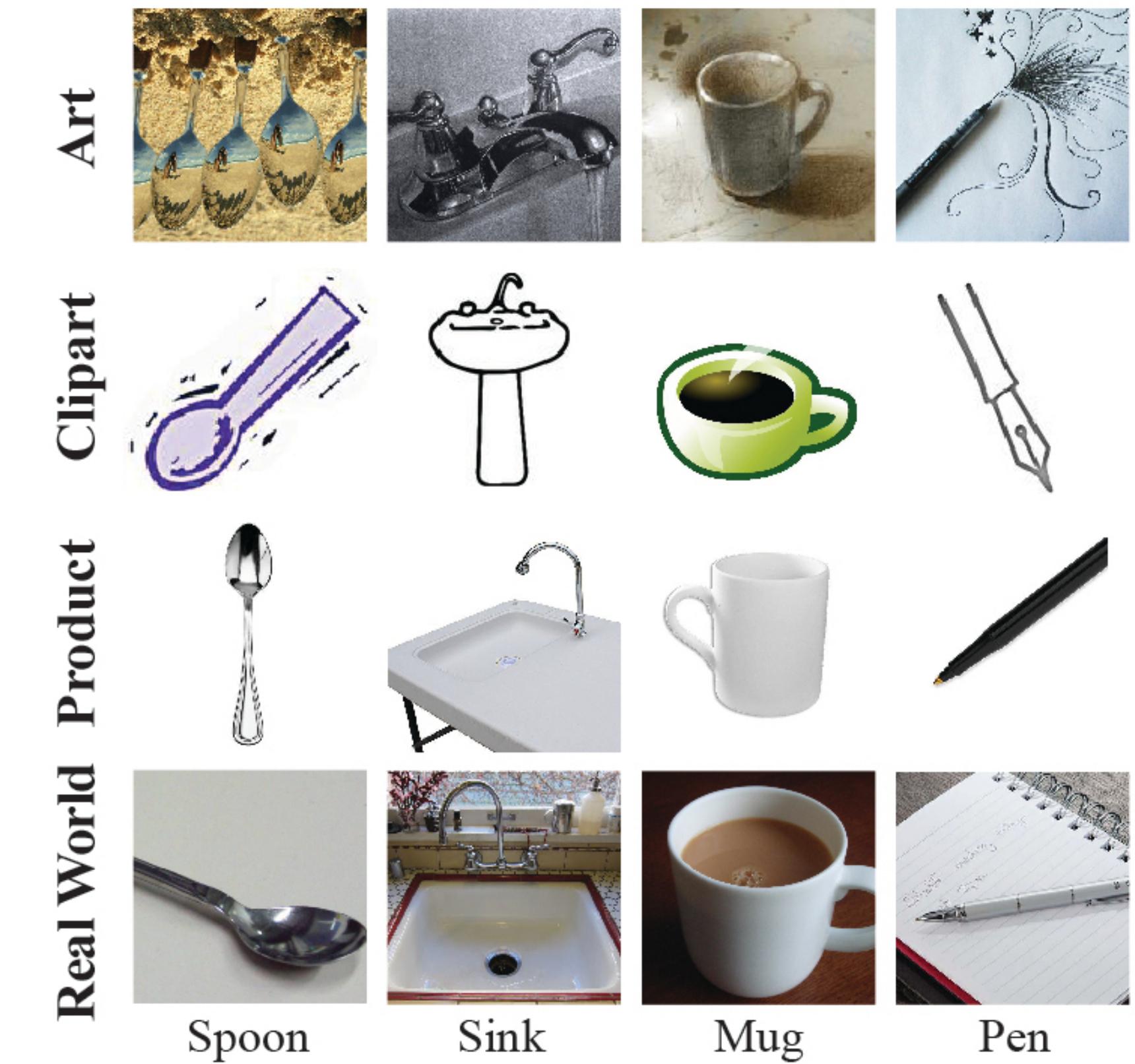


# SENTRY Results: Office Home

Custom label shifts<sup>2</sup>

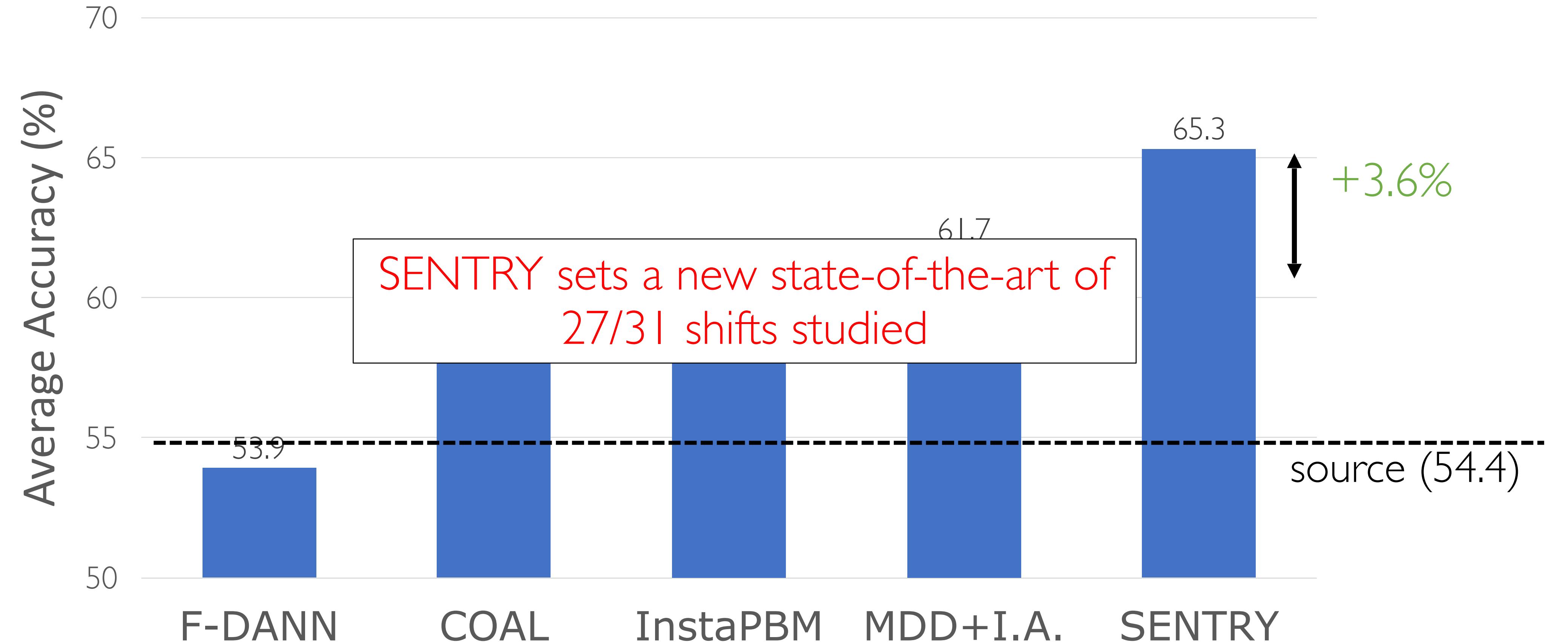


Office Home<sup>1</sup>



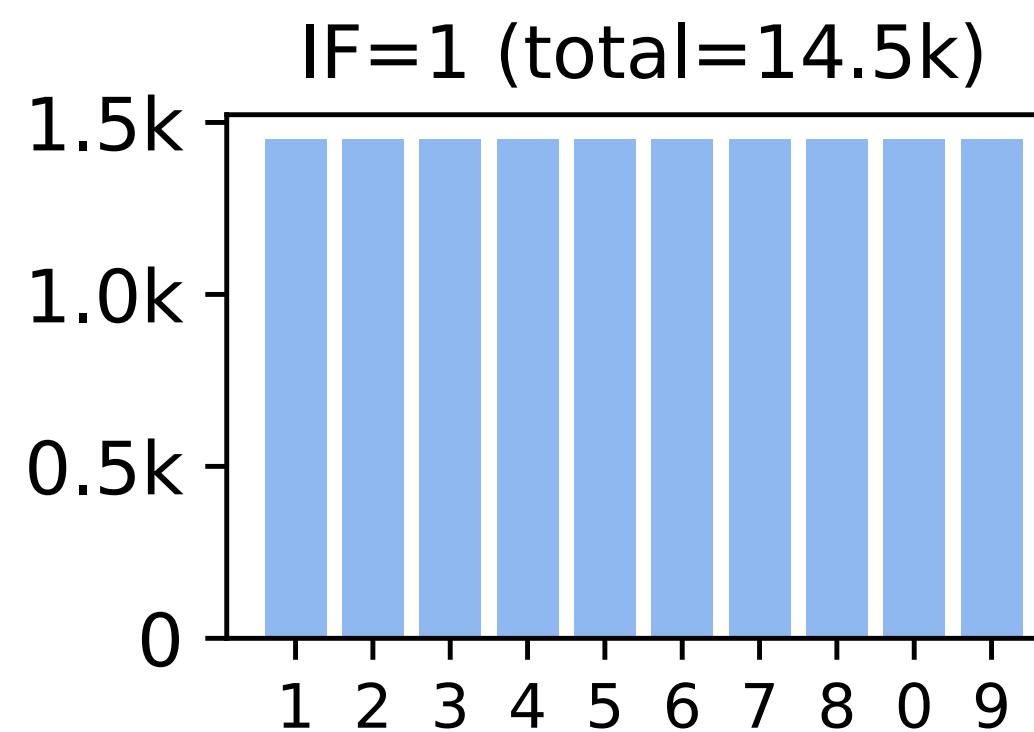
# SENTRY Results: Office Home

OfficeHome-LDS (65 classes, 6 shifts)

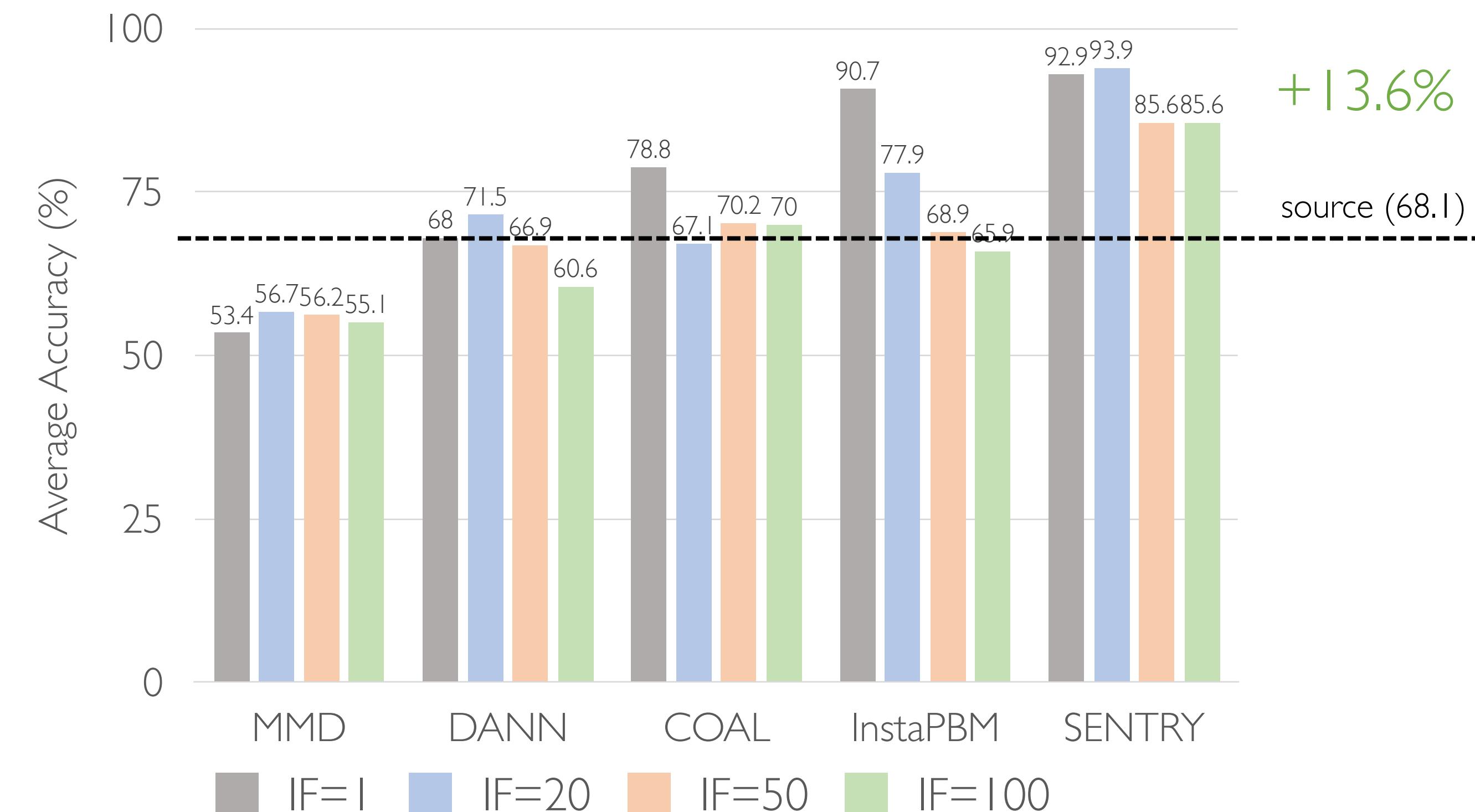


# Results: Controlling Target Distribution

MNIST-LT label histograms

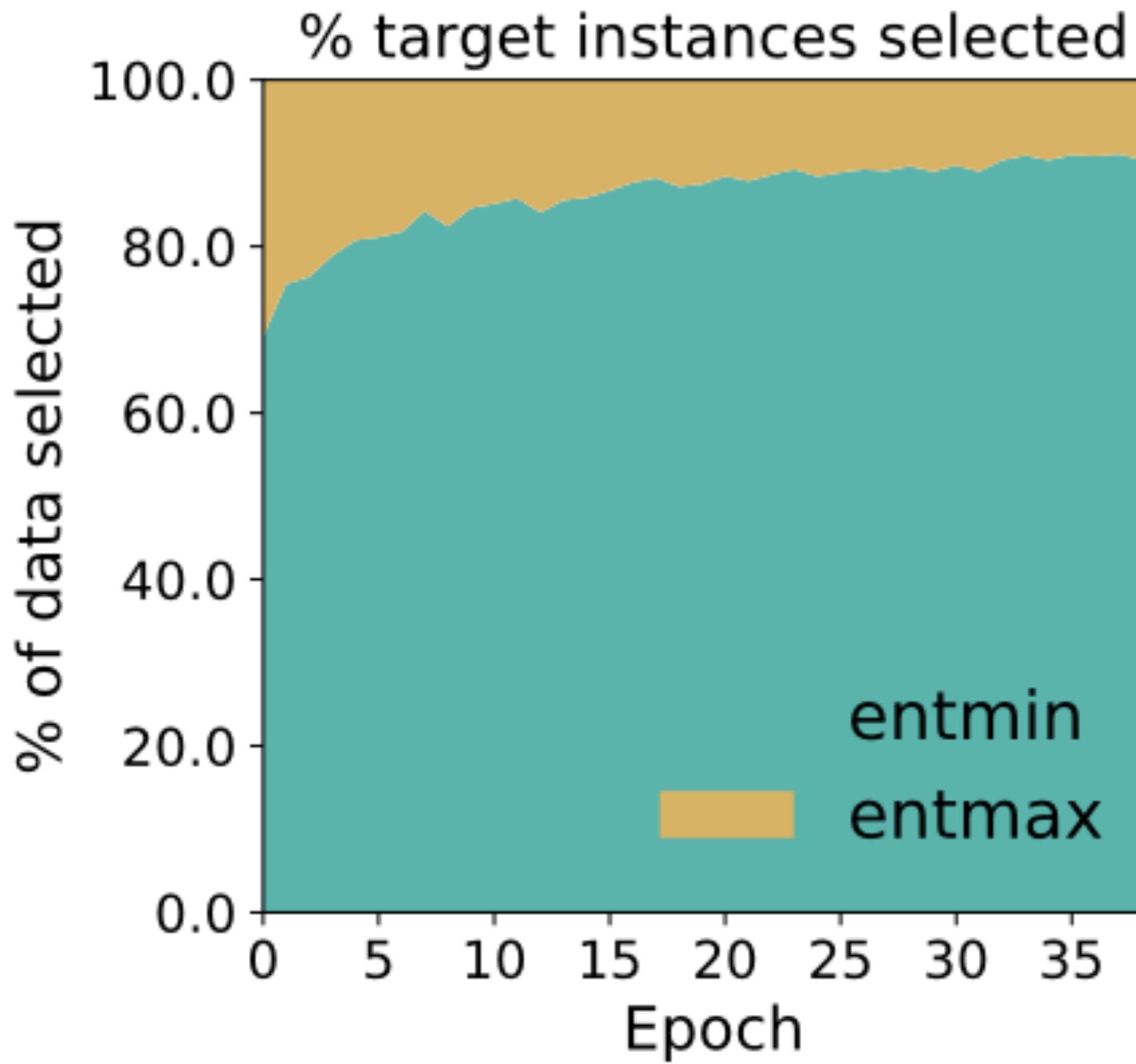


SVHN→MNIST-LT

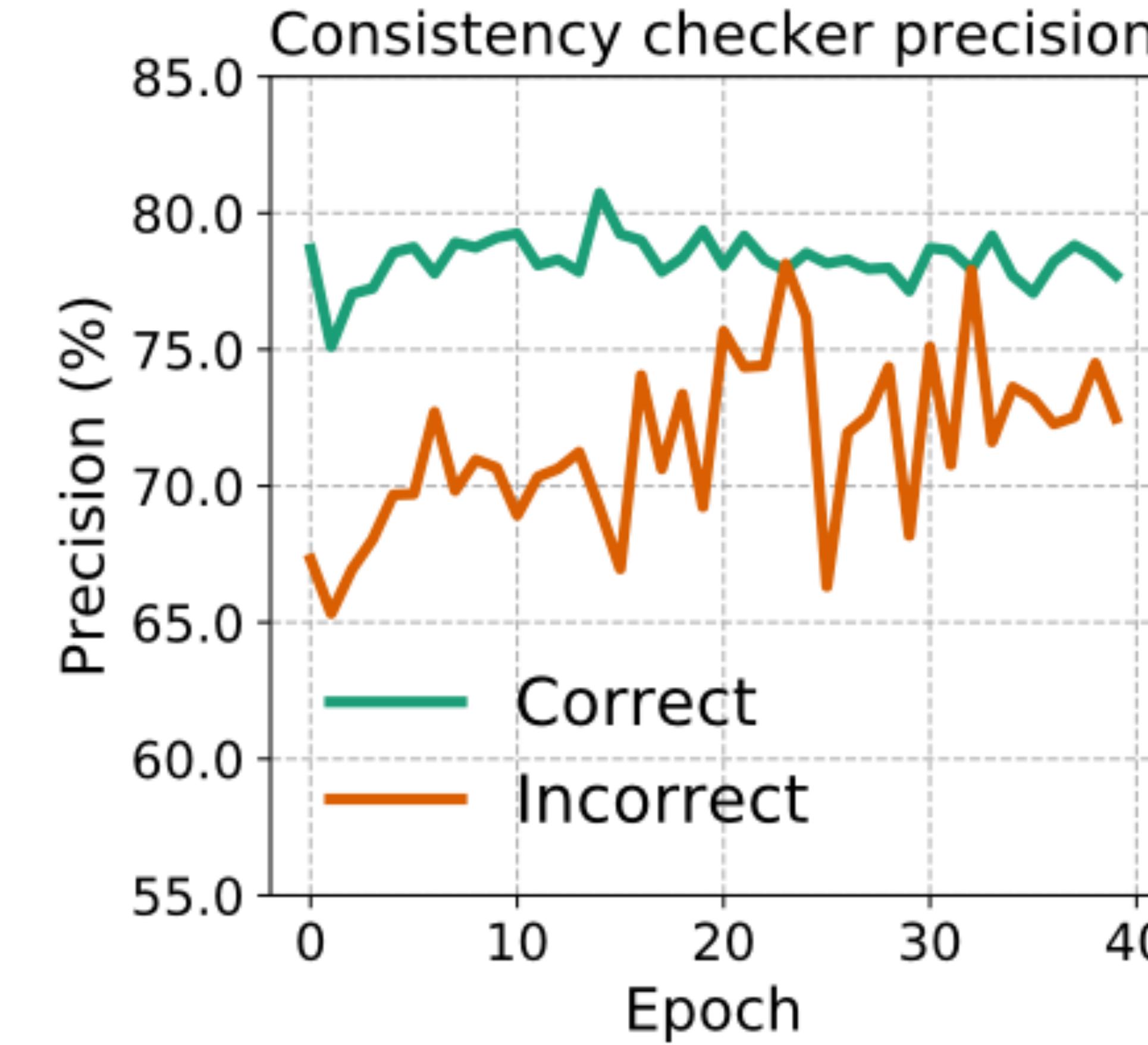


# Ablating SENTRY: Selection

Increasing % of data is aligned over time



Predictive consistency is a good reliability measure

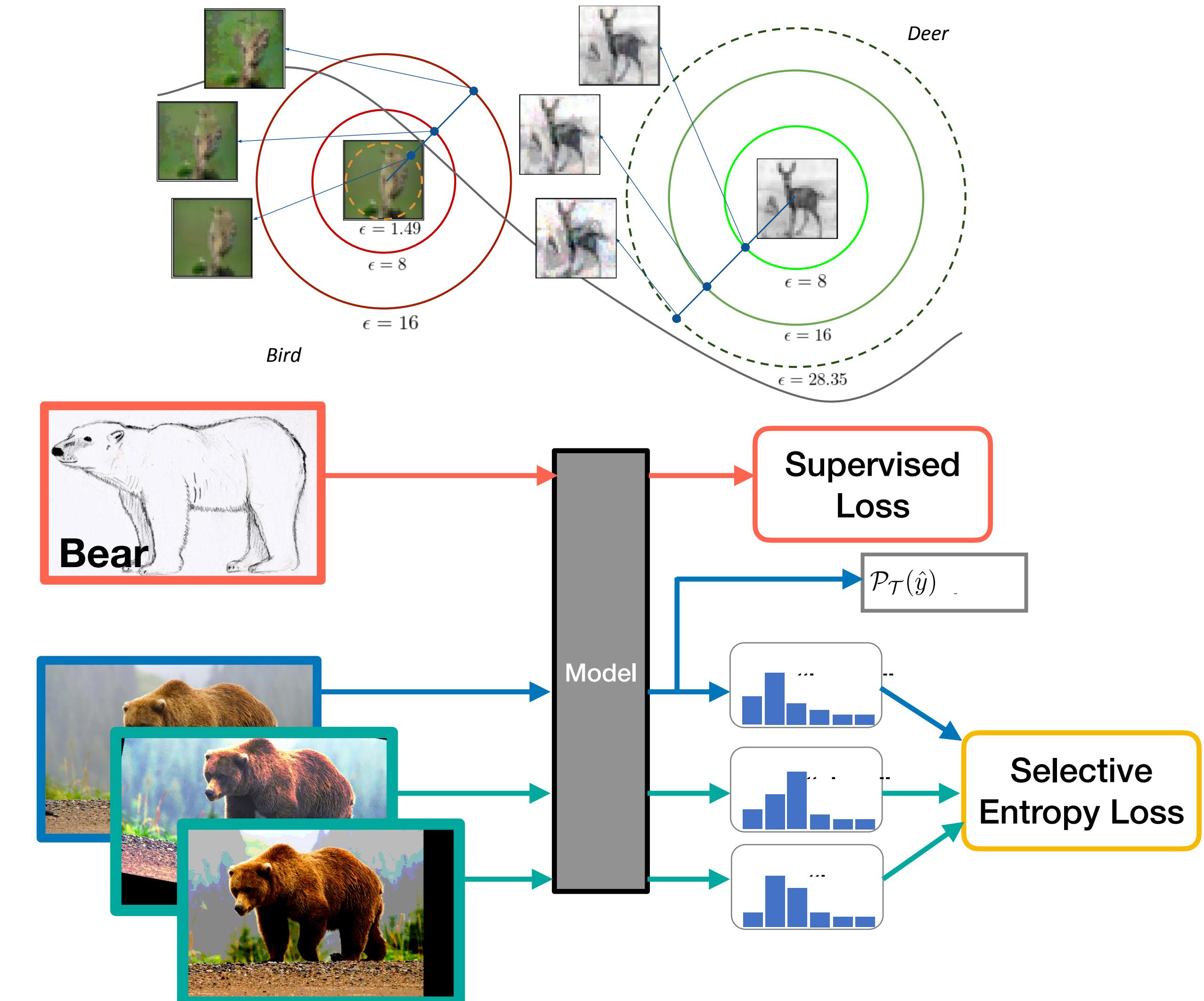


# Summary

**Key Idea:** Decide when to learn

Learning can be derailed by

- Unreliable labels
  - Label noise
  - Manipulated
  - Model misalignment
- Unreliable samples
  - Inherently ambiguous
  - Different from prior data
  - Manipulated



# Thank you



Sean Foley



Daniel Bolya



Sruthi Sudhakar



Arvind  
Krishnakumar



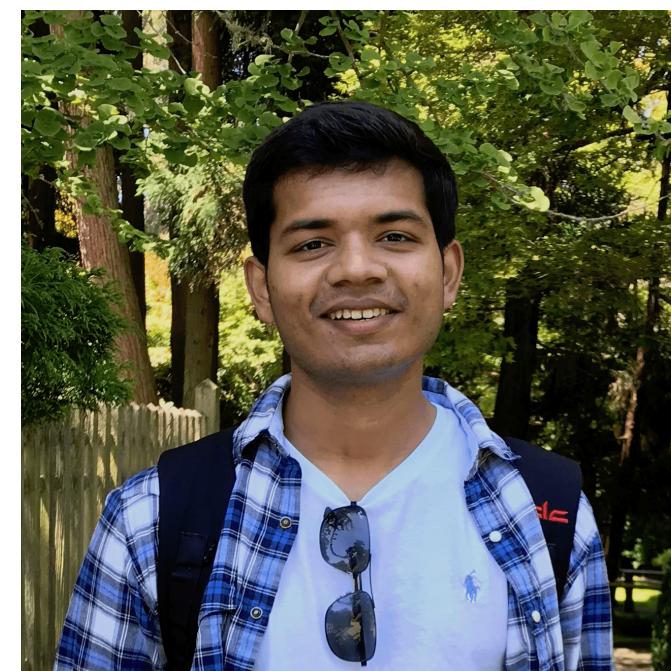
Rohit Mittapali



Kartik  
Sarangmath



Prithvijit  
Chattopadhyay



Viraj Prabhu



Shivam Khare



Deeksha Karthik



Bhavika Devnani



Luis Bermudez

# Summary

**Key Idea:** Decide when to learn

Learning can be done

- Unreliable labels
  - Label noise
  - Manipulated
  - Model misaligned
- Unreliable samples
  - Inherently ambiguous
  - Different from prior data
  - Manipulated

Questions?  
[judy@gatech.edu](mailto:judy@gatech.edu)

