

Responsible CV:

How do models fail and what can we do about it?

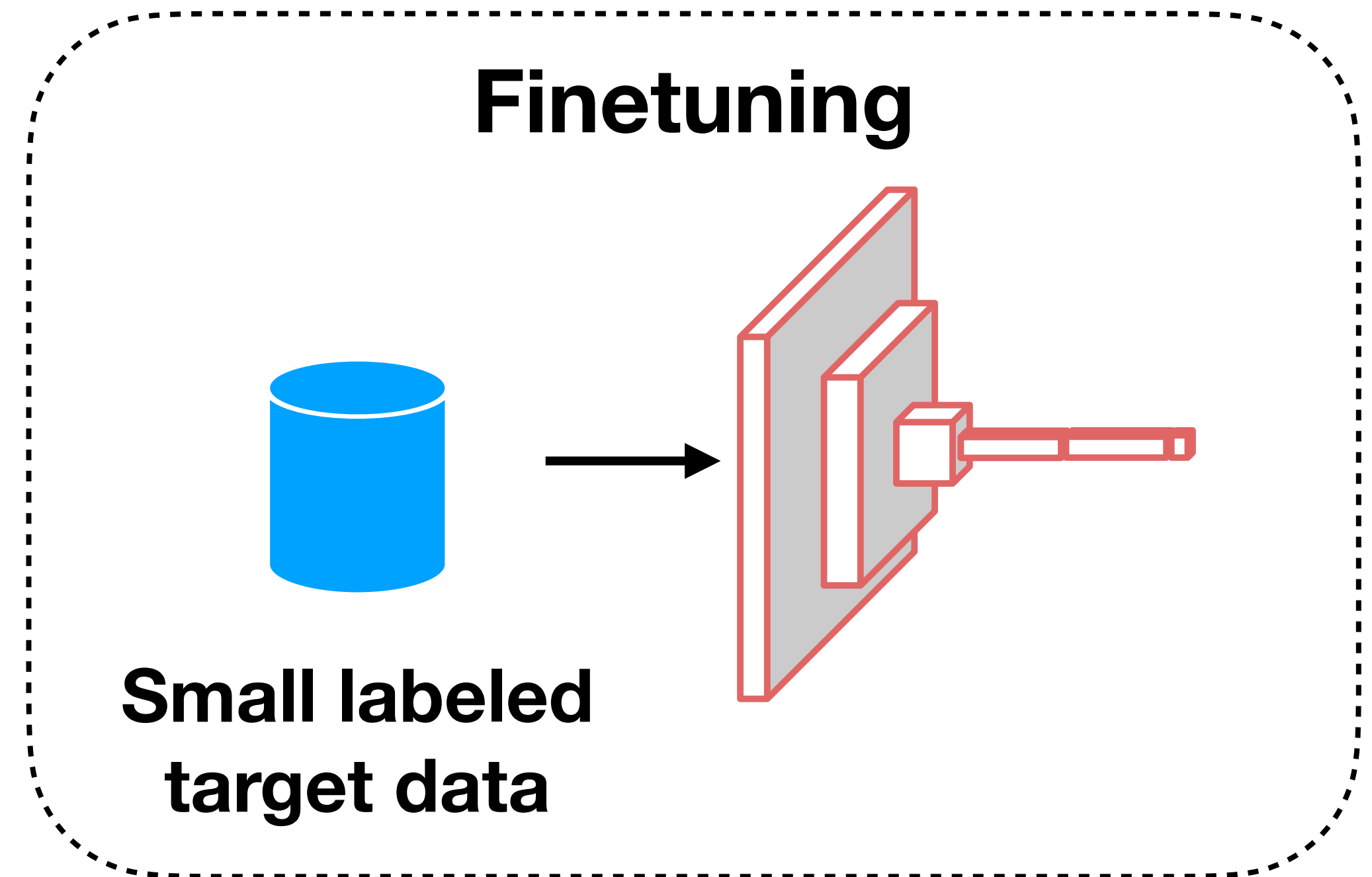
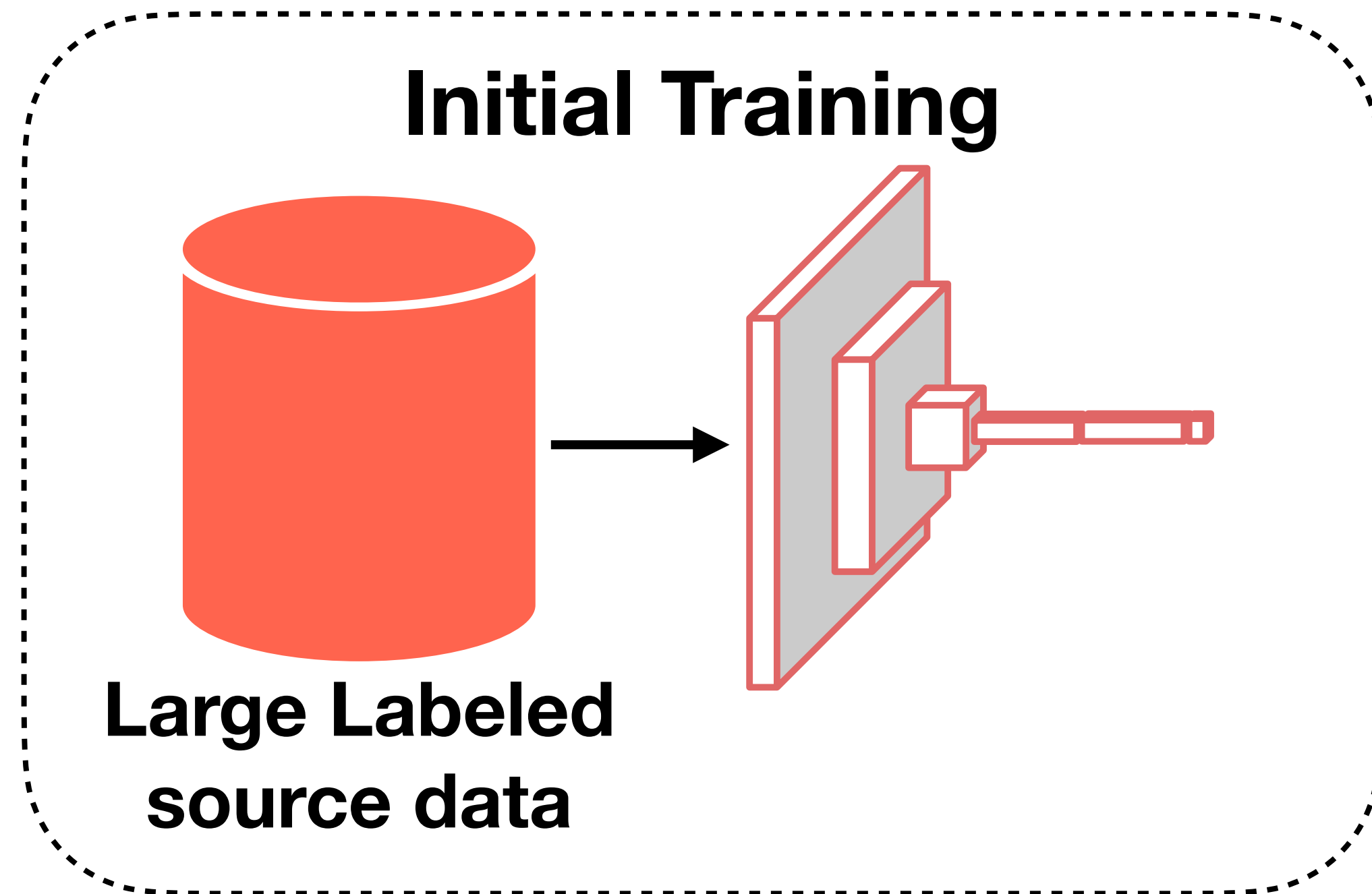
Judy Hoffman and Viraj Prabhu

Human-centered AI Tutorial @CVPR

June 20, 2022

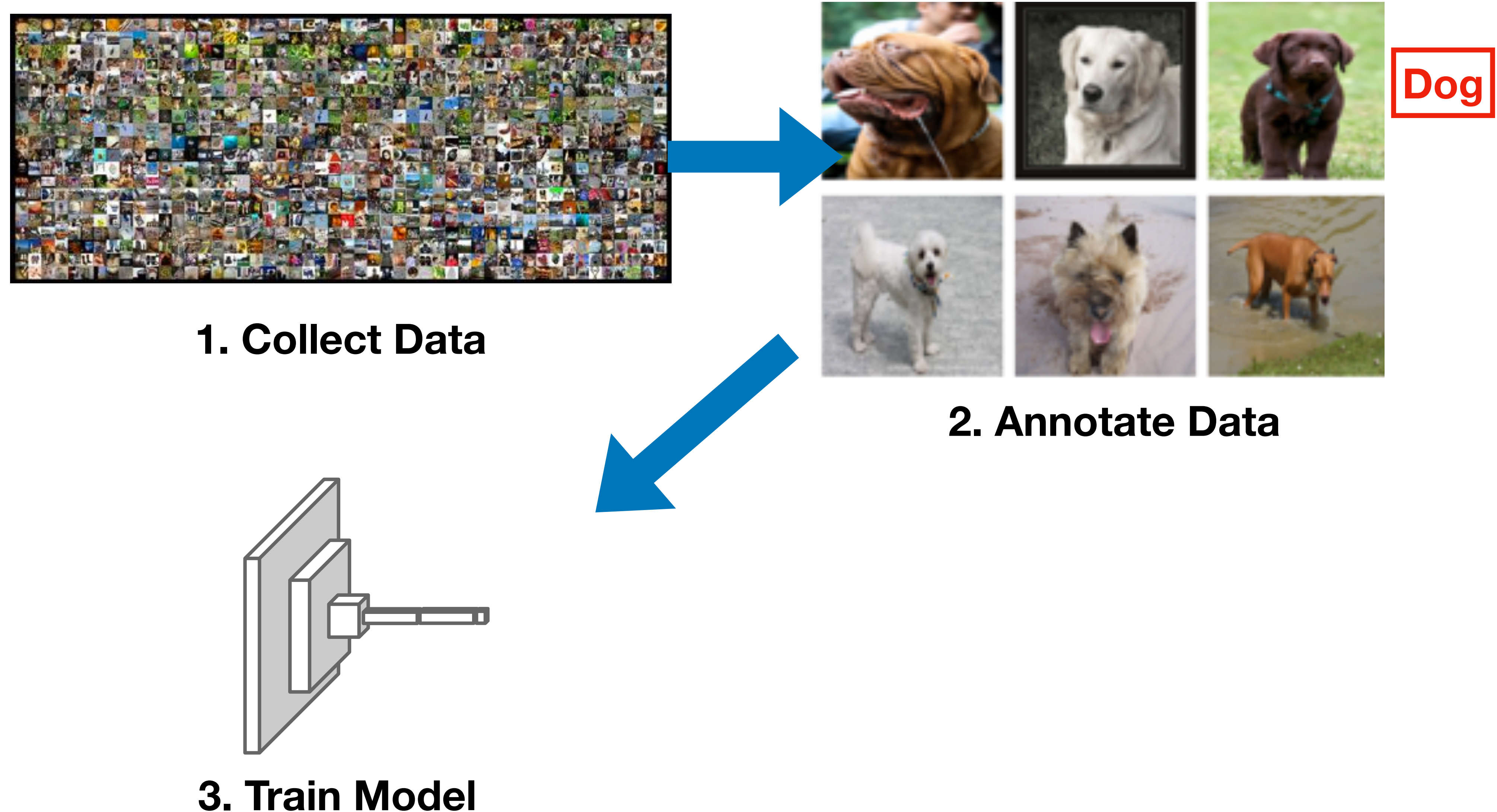


Practical Transfer Learning



Frequently select model that performs best on ImageNet

Standard Visual Recognition Pipeline



Visual Recognition Benchmarks

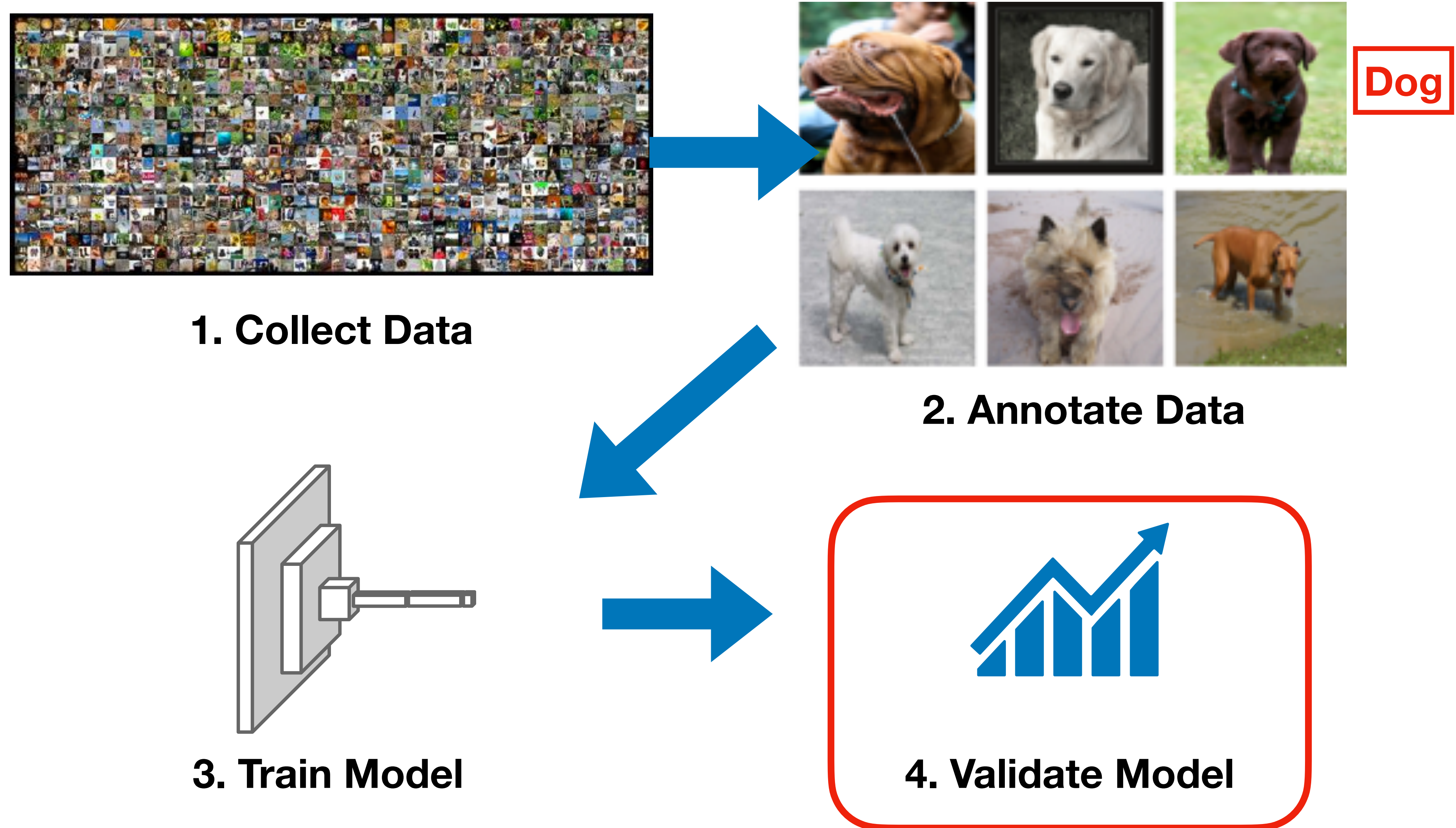


Classification



Detection / Segmentation

Standard Visual Recognition Pipeline

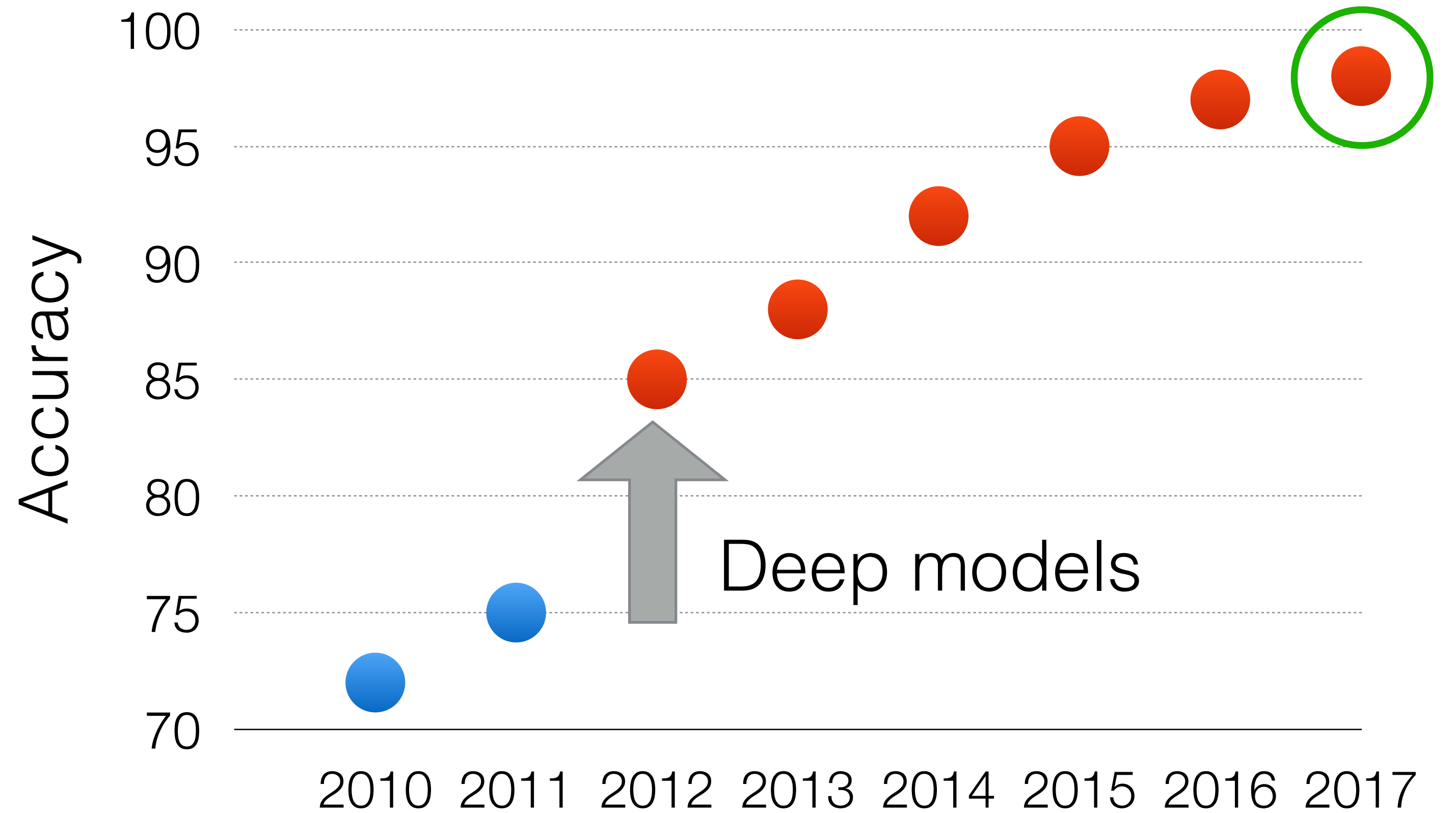


Benchmark Performance

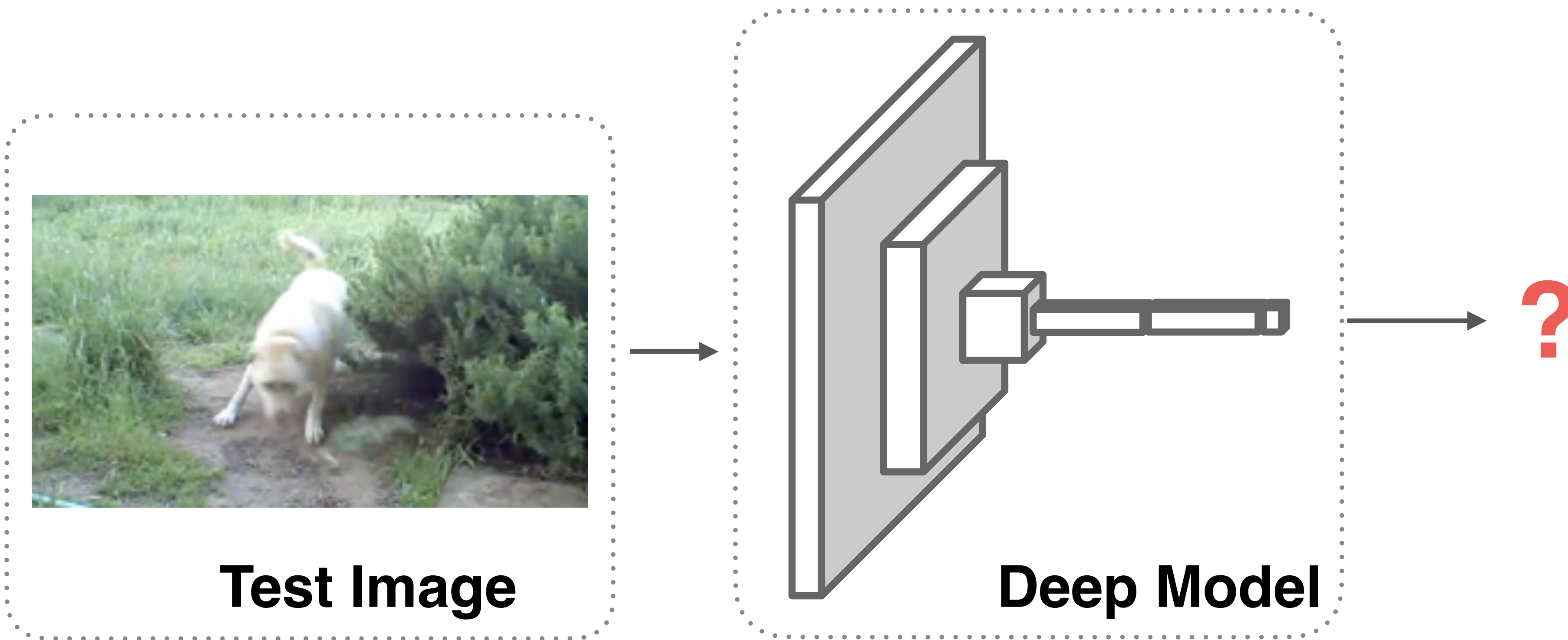


Millions of Images

**Challenge to recognize
1000 categories**



Dataset Bias



Dataset Bias



Test Image

Dog is not recognized



Deep Model



Dataset Bias



Dataset Bias



Low resolution

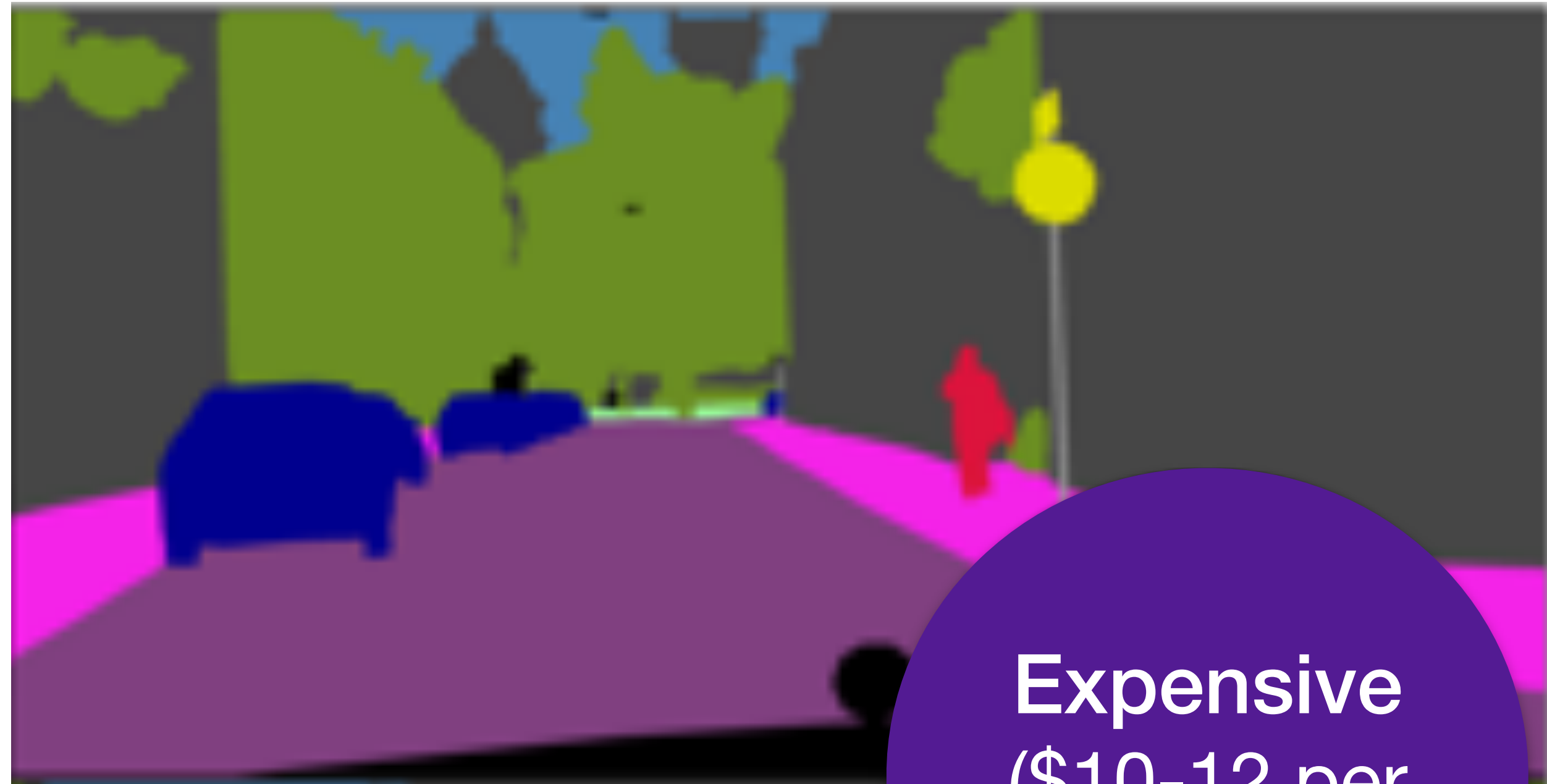


Motion Blur



Pose Variety

The world has high natural variation



Expensive
(\$10-12 per
image)

Large Potential for Change
Different: Weather, City, Car

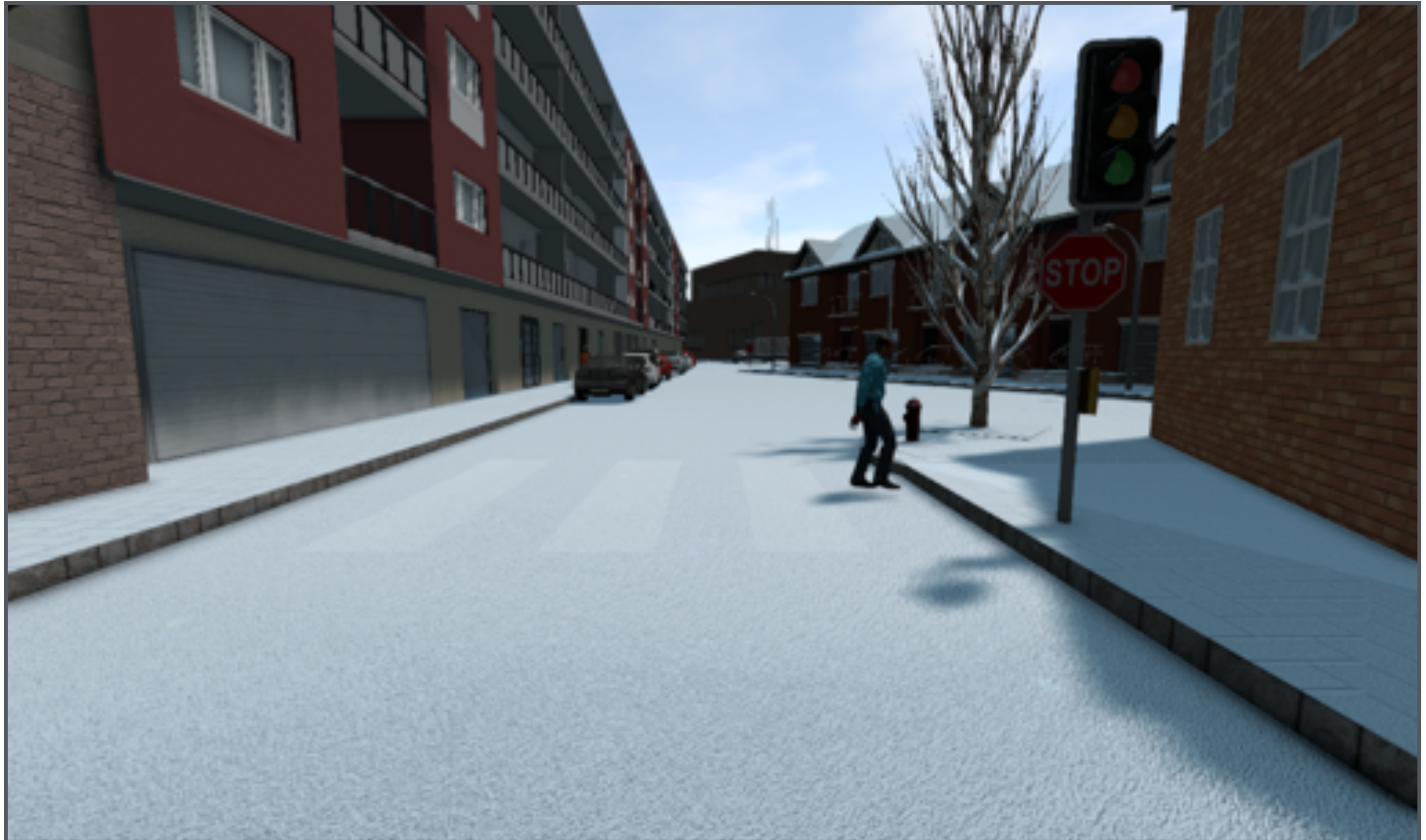
- | | |
|--|---|
|  Car |  Sky |
|  Road |  Vegetation |
|  Sidewalk |  Street Sign |
|  Person |  Building |

Train in Sunny Weather



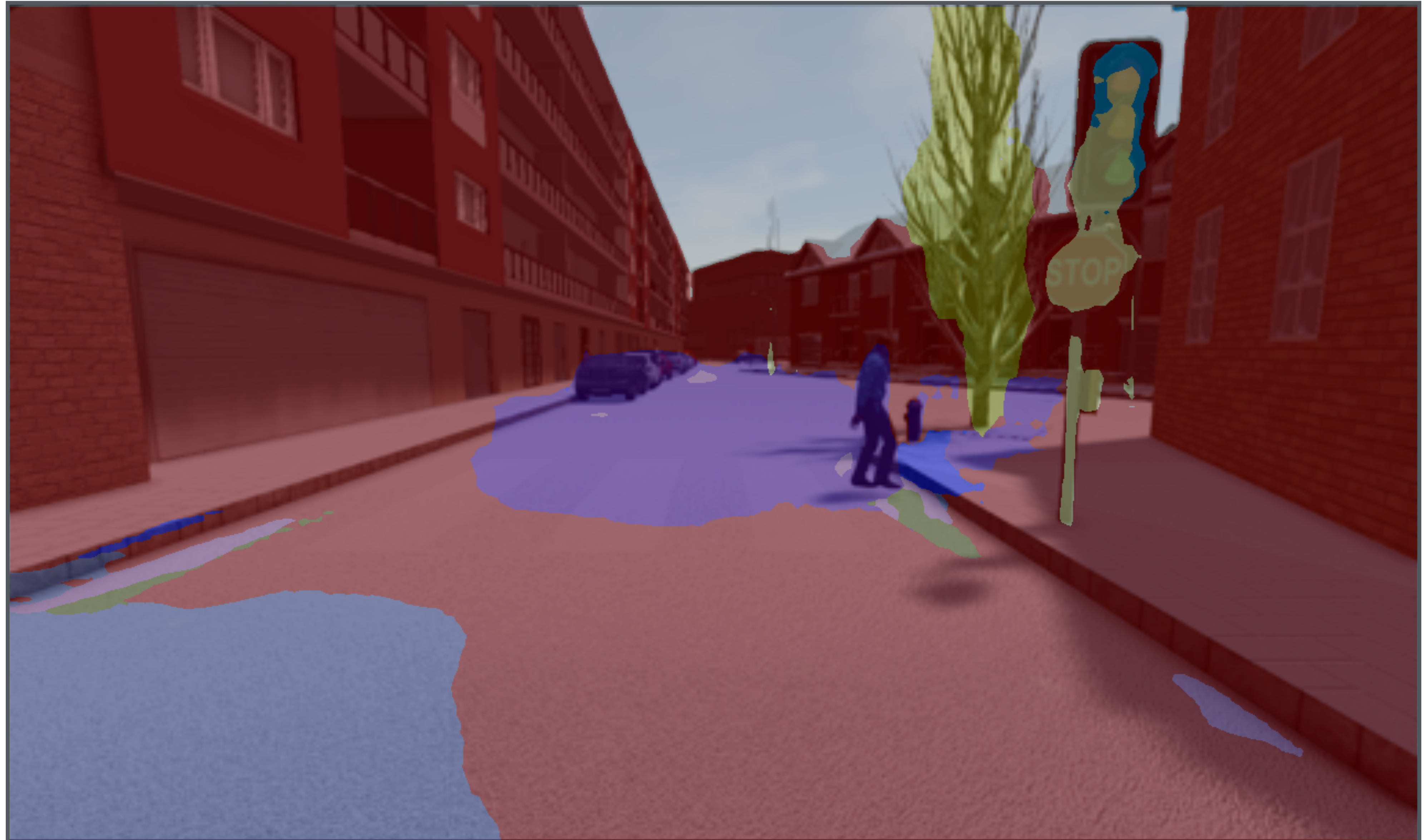
Robust to Weather Changes?

- Car
- Road
- Sidewalk
- Person
- Sky
- Vegetation
- Street Sign
- Building
- Traffic Light



Robust to Weather Changes?

- Car
- Road
- Sidewalk
- Person
- Sky
- Vegetation
- Street Sign
- Building
- Traffic Light



Impact of Input Corruptions on Recognition

CiFAR-10, ResNet-18

Clean Acc = 94.2



Corrupt Acc = 72.7



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

Benchmark Challenge Adversarial

The Art of Robustness: Devil and Angel in Adversarial Machine Learning

Workshop at IEEE Conference on Computer Vision and Pattern Recognition 2022

RobustNav

Towards Benchmarking Robustness in Embodied Navigation

ICCV 2021



Prithvijit
Chattopadhyay¹



Judy
Hoffman¹



Roozbeh
Mottaghi²



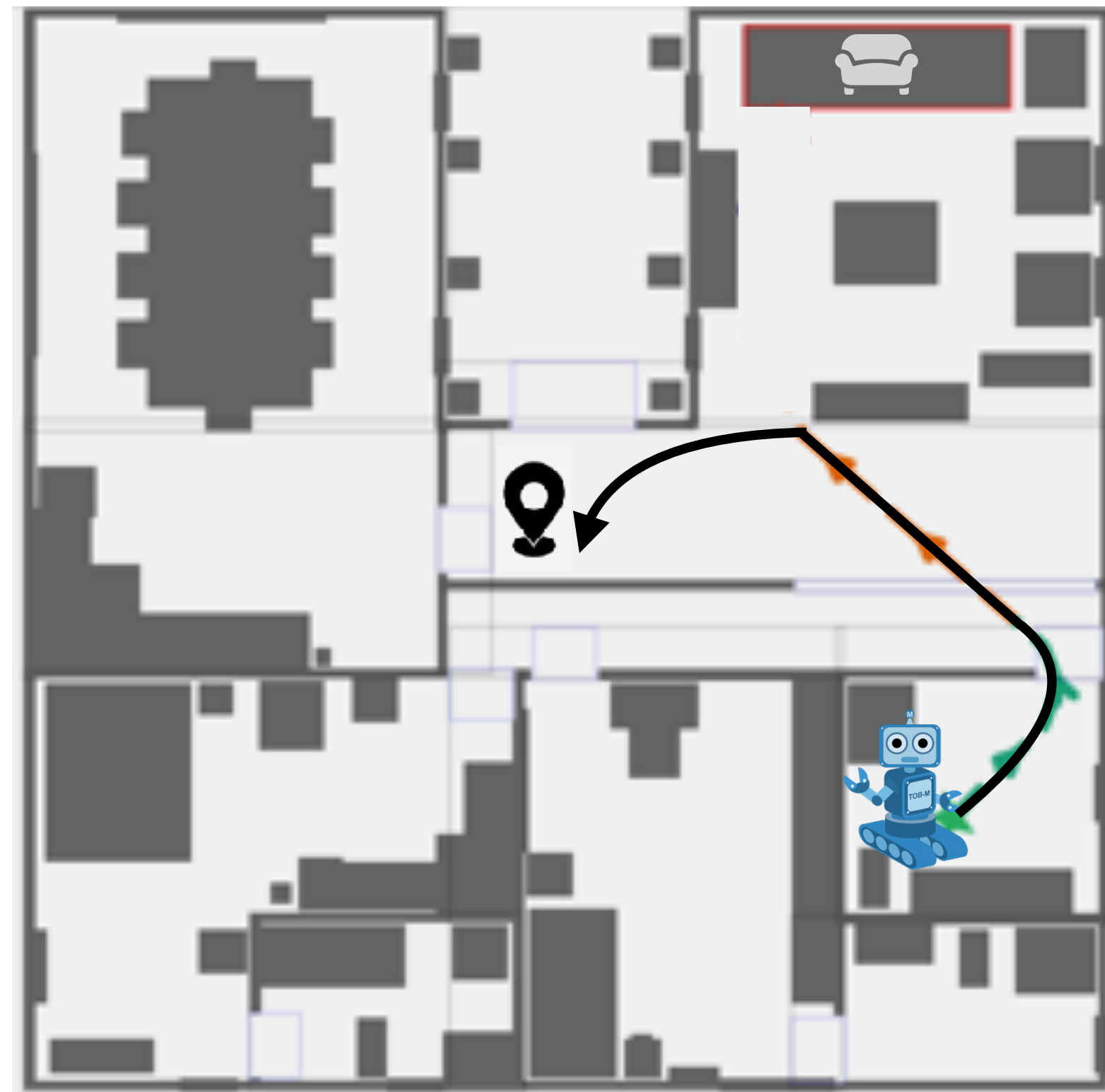
Ani
Kembhavi²



Visual Navigation (RGB+Depth)

POINTNAV

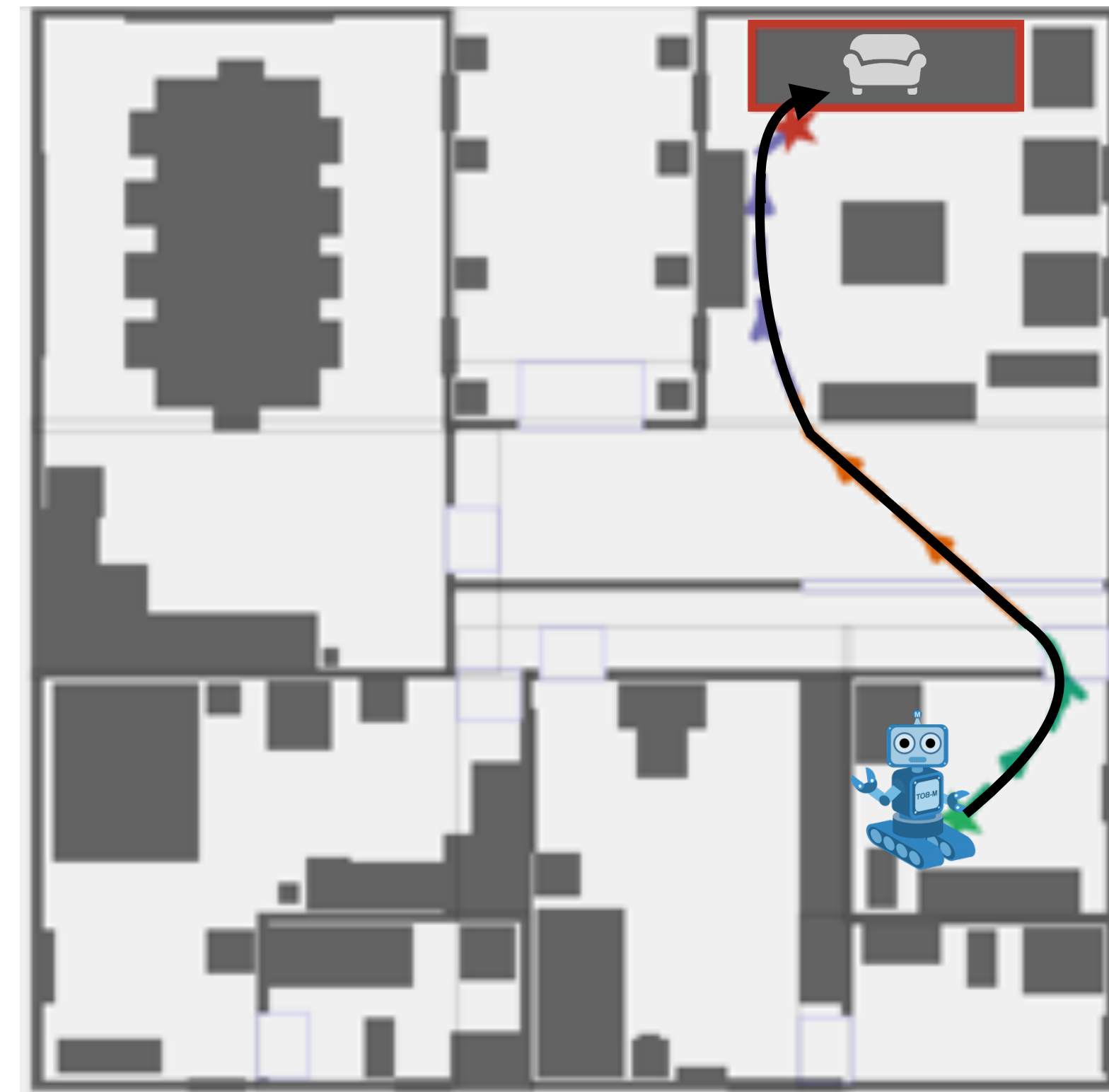
GPS + Compass enabled navigation



Task: Go to (r, θ) location

OBJECTNAV

Semantic, Target-driven Navigation



Task: Go to a "sofa"

Visual Navigation

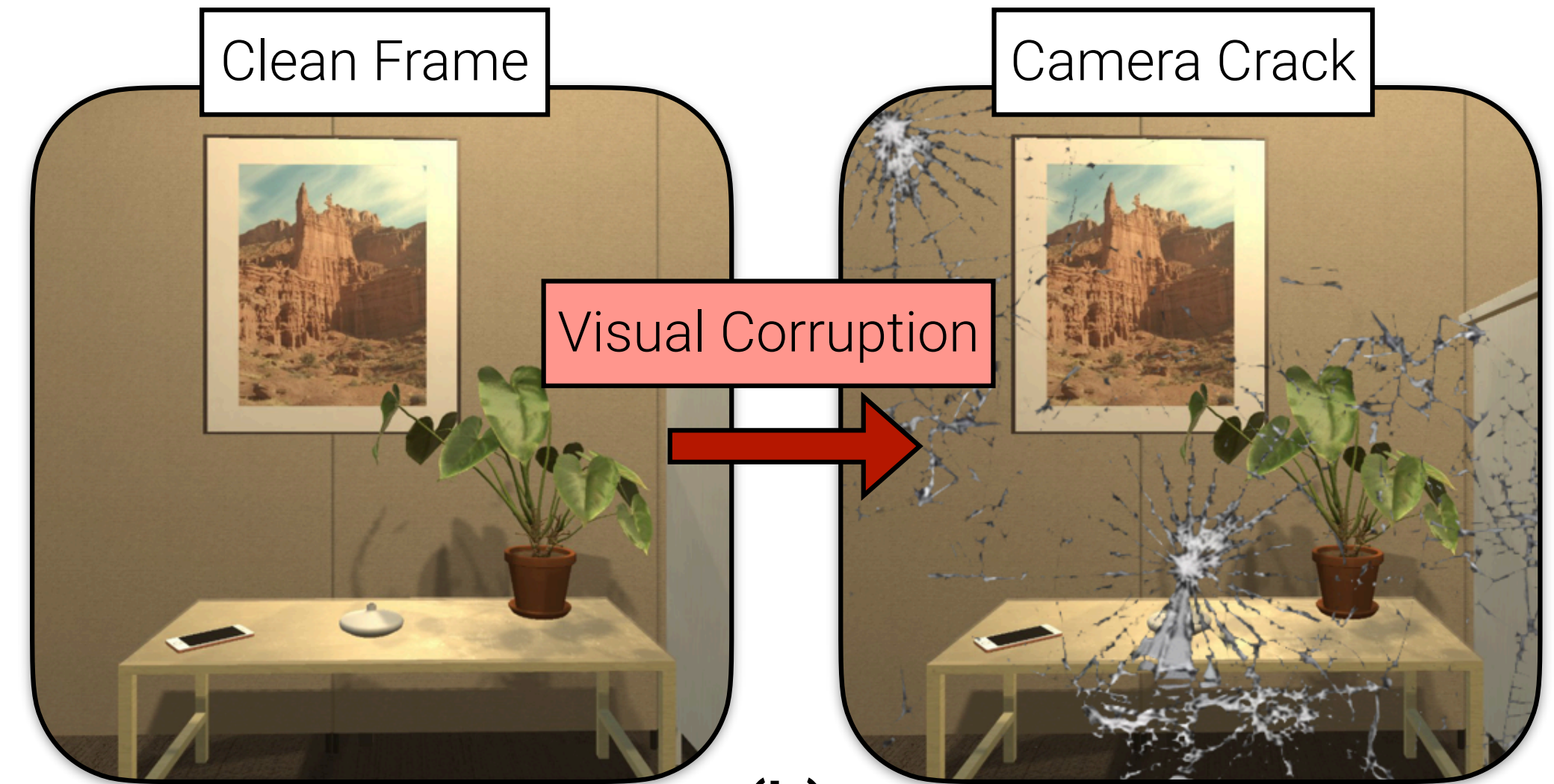
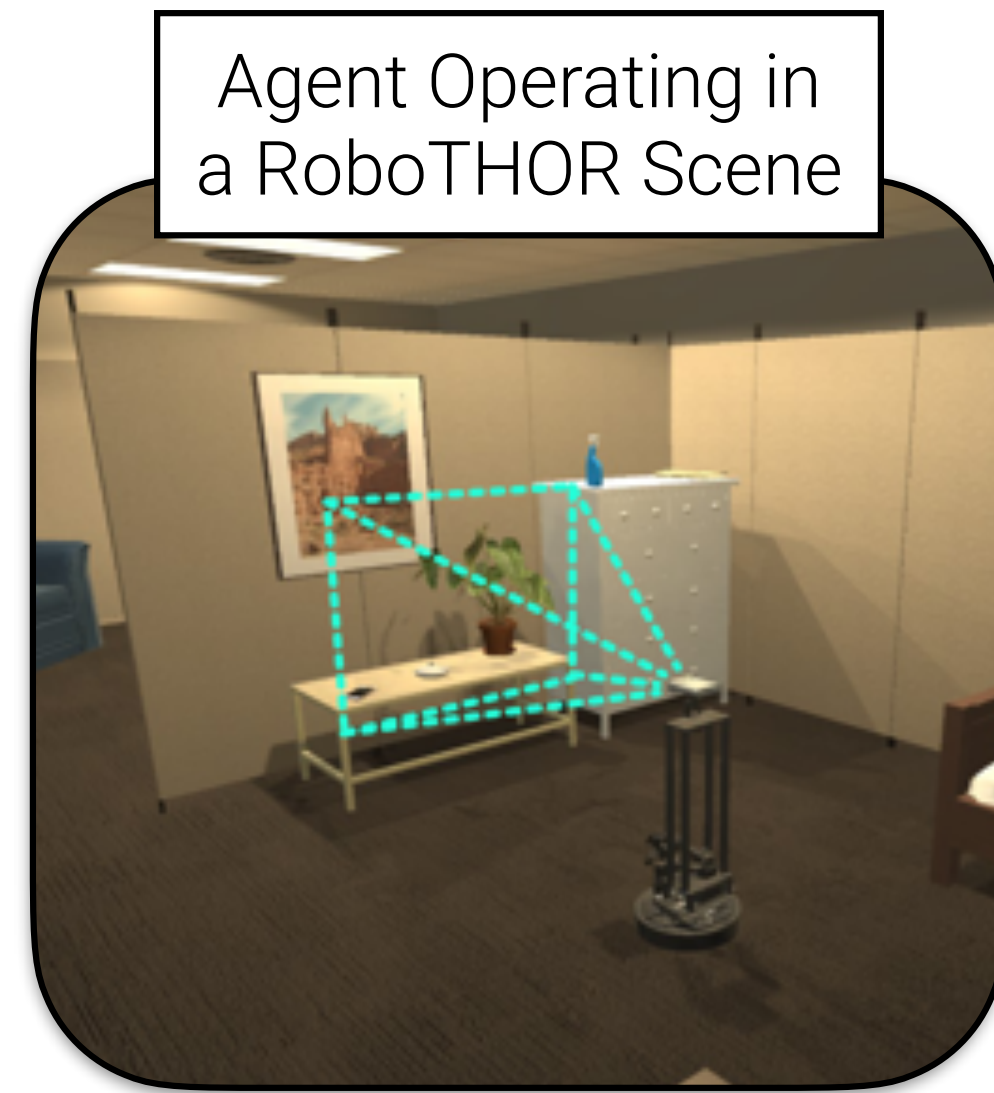
Agents don't have access to a "map", and must navigate based solely on sensory inputs



RobustNav

7 visual corruptions
at **5** levels of severity

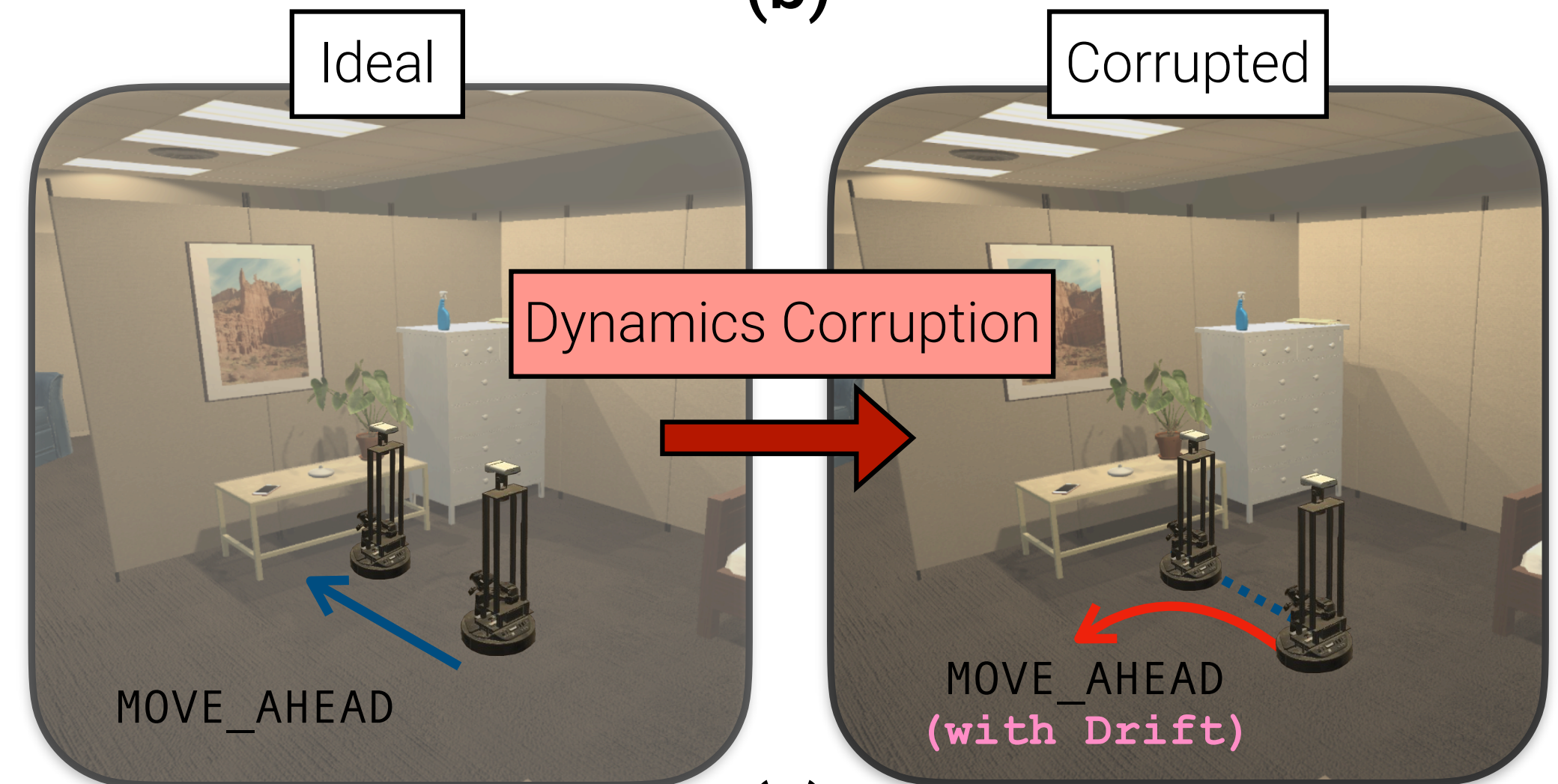
4 dynamics corruptions
Corruptions can be due
to sensor or
environment variations



Agent – LoCoBot

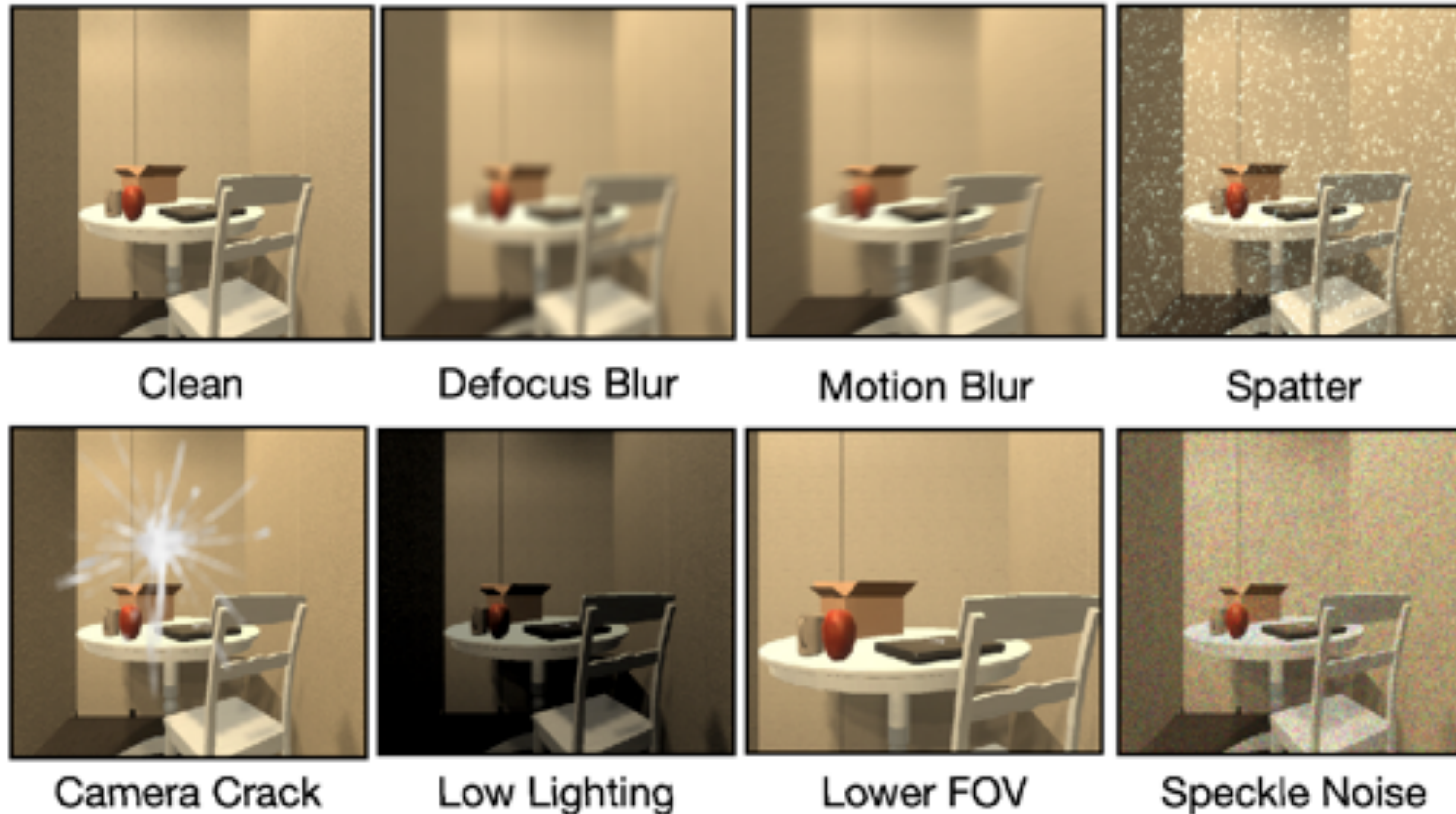


(a)



(c)

RobustNav Visual Corruptions



Severity 1

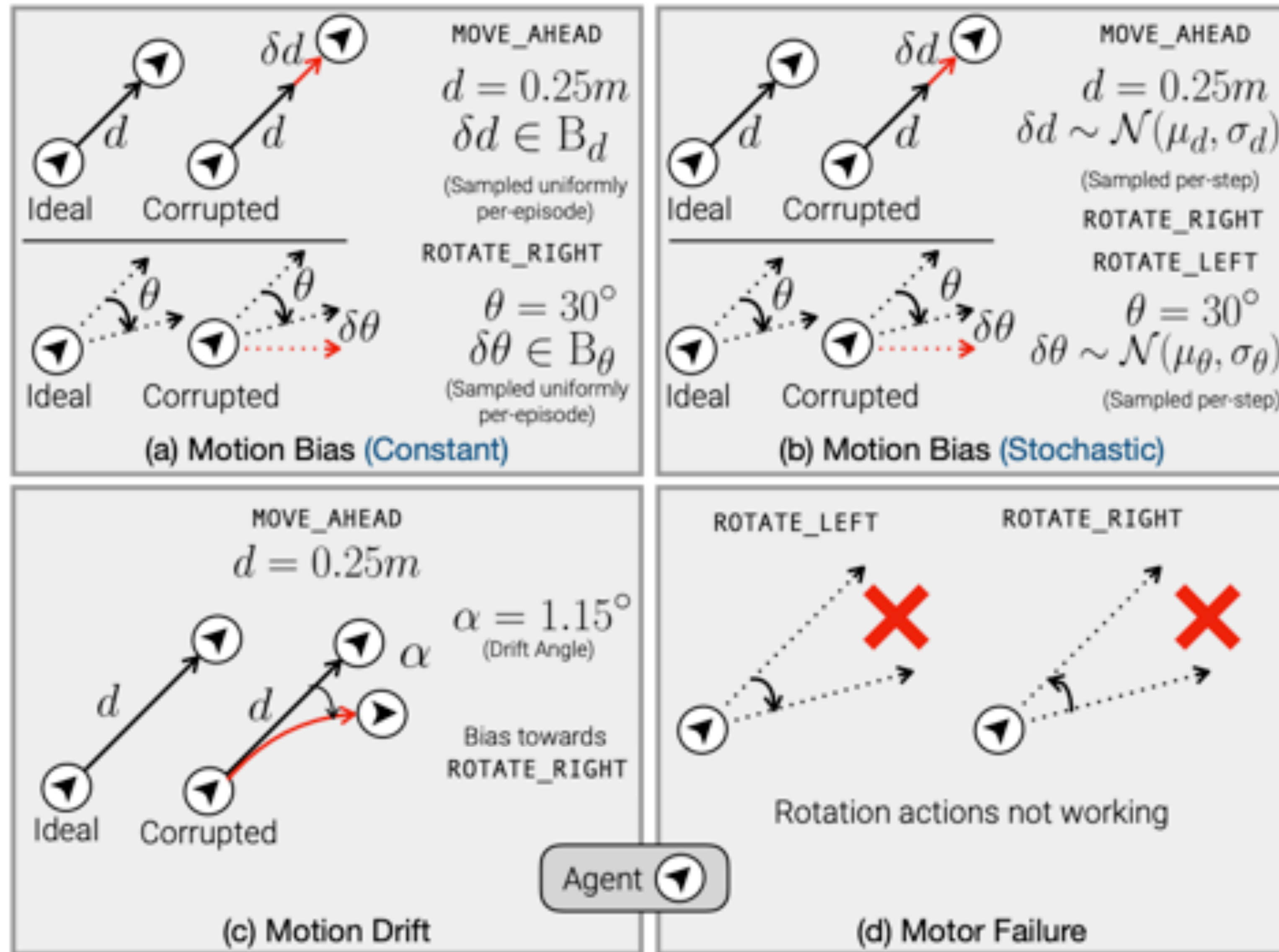
Visual Corruptions at 5 levels of severity

Severity 5

Low

High

RobustNav Dynamics Corruptions

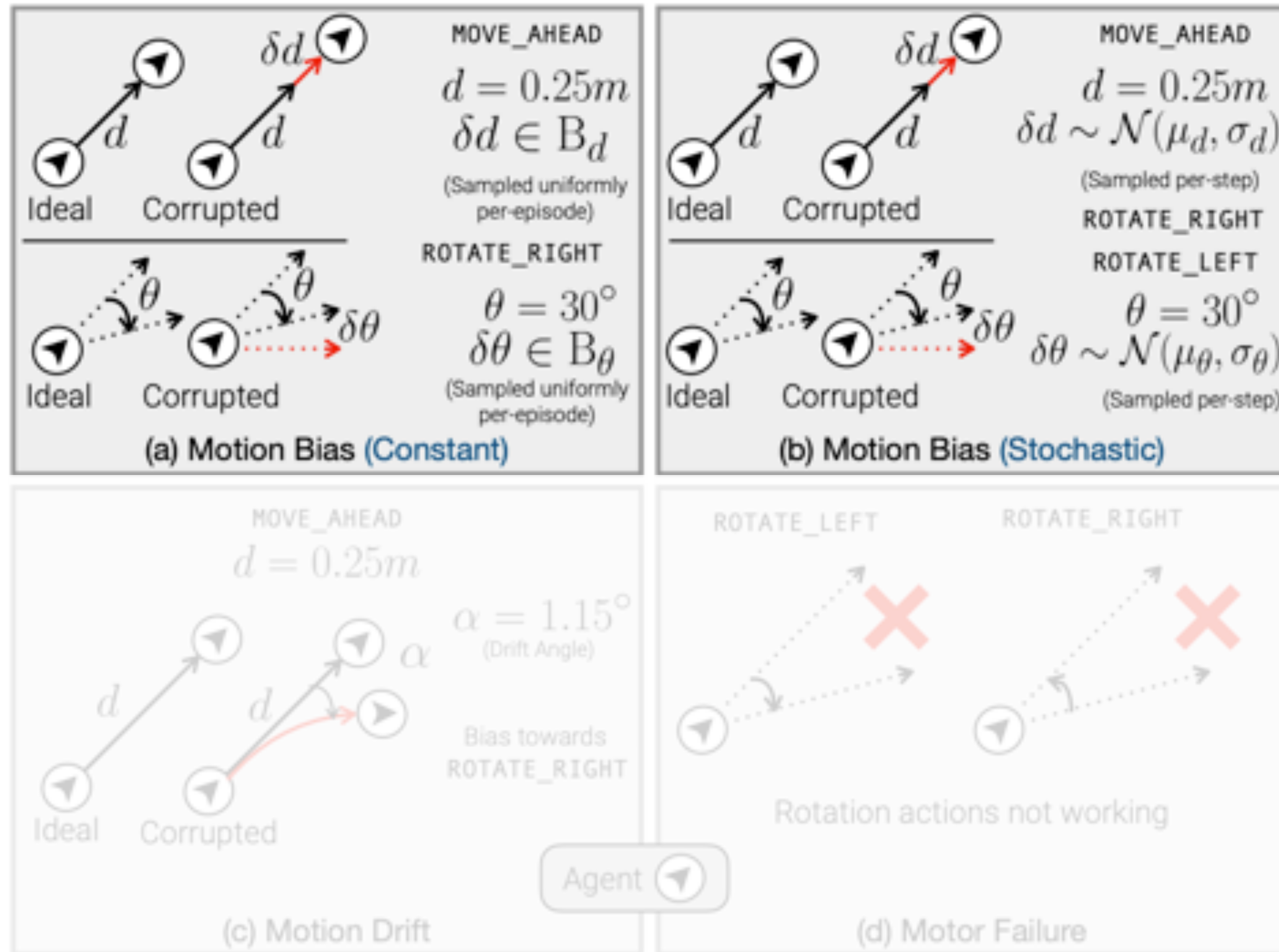


RobustNav Dynamics Corruptions

Due to Environment

Scene-level friction

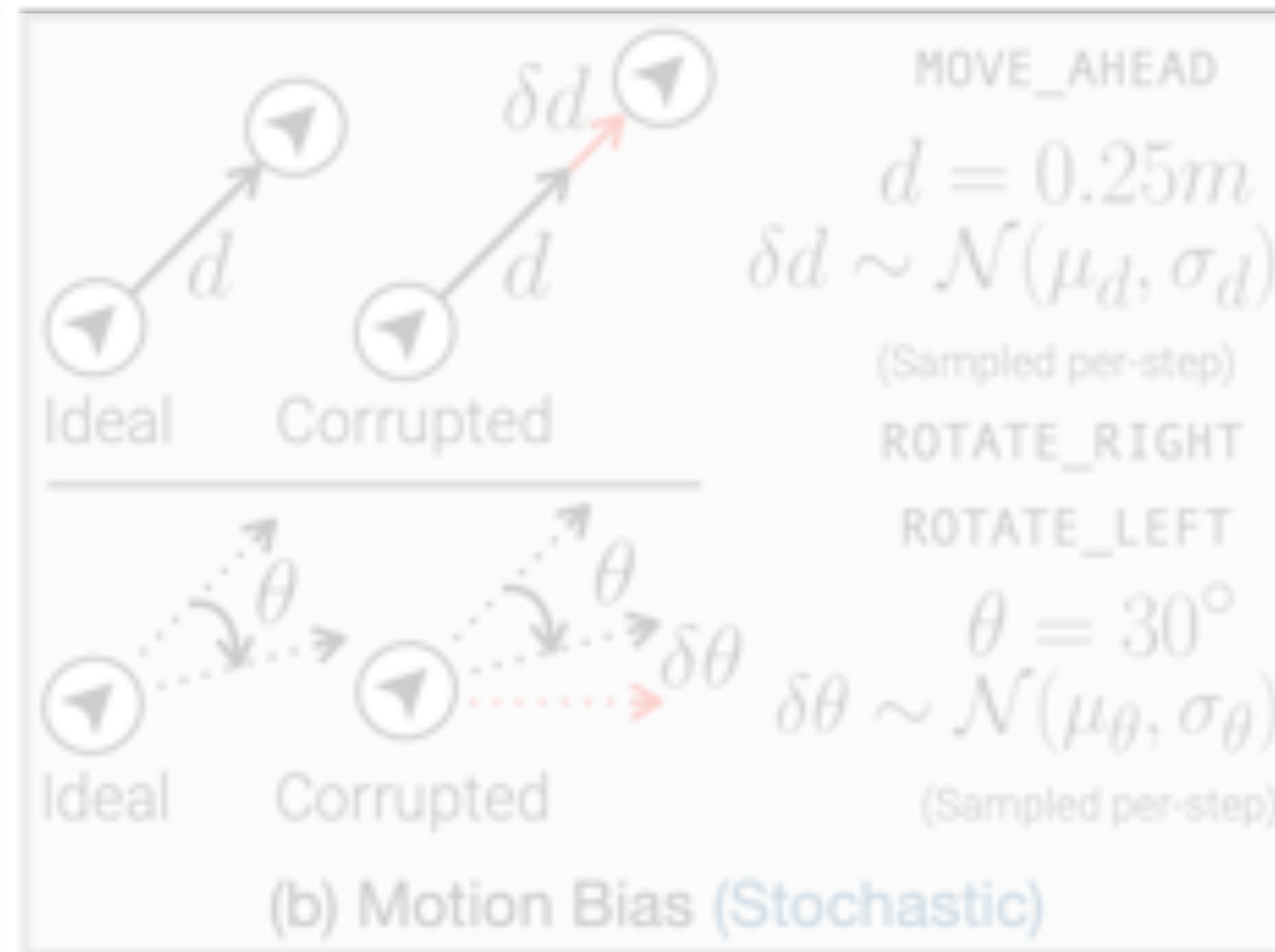
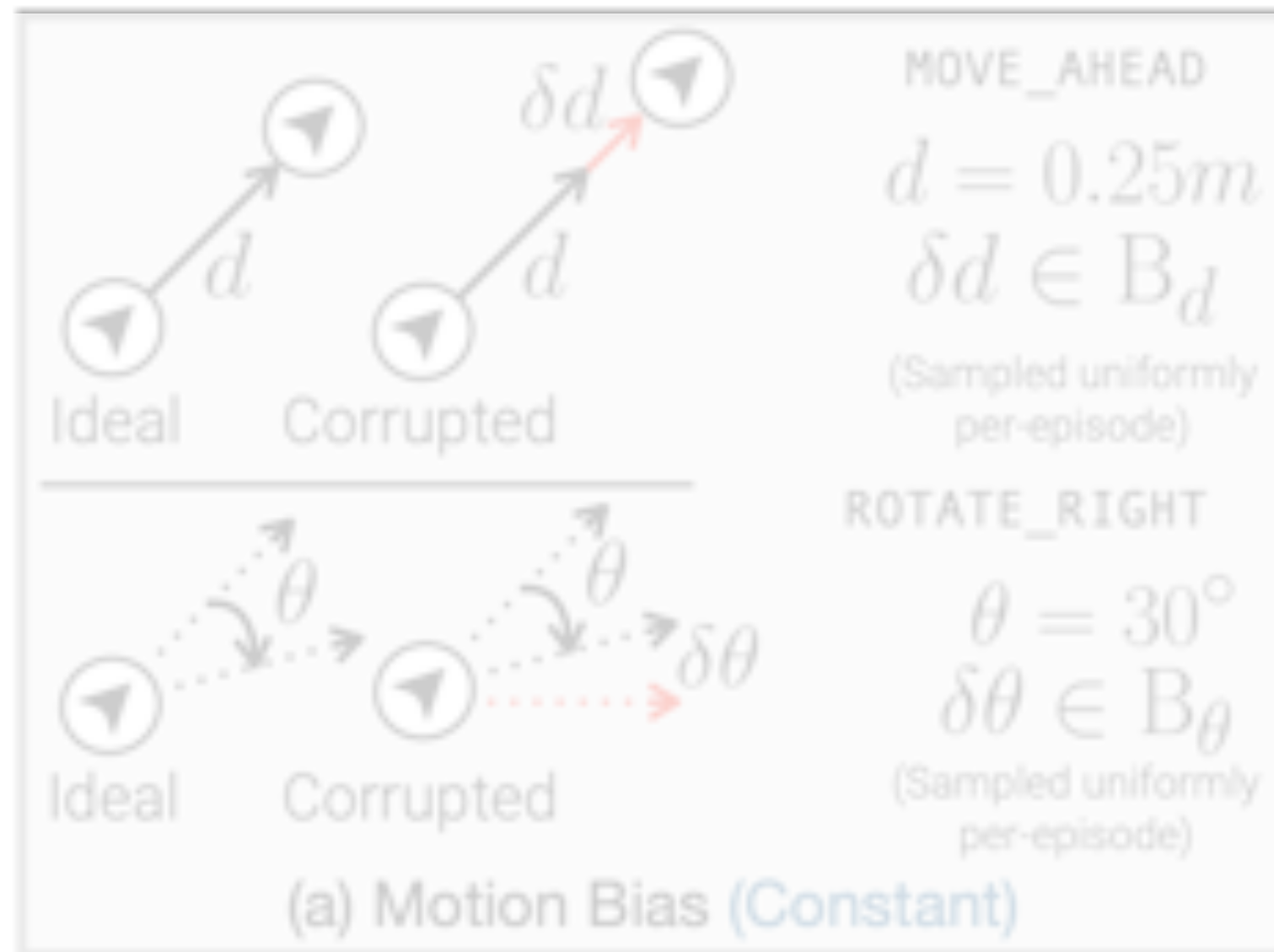
High and low friction zones



RobustNav Dynamics Corruptions

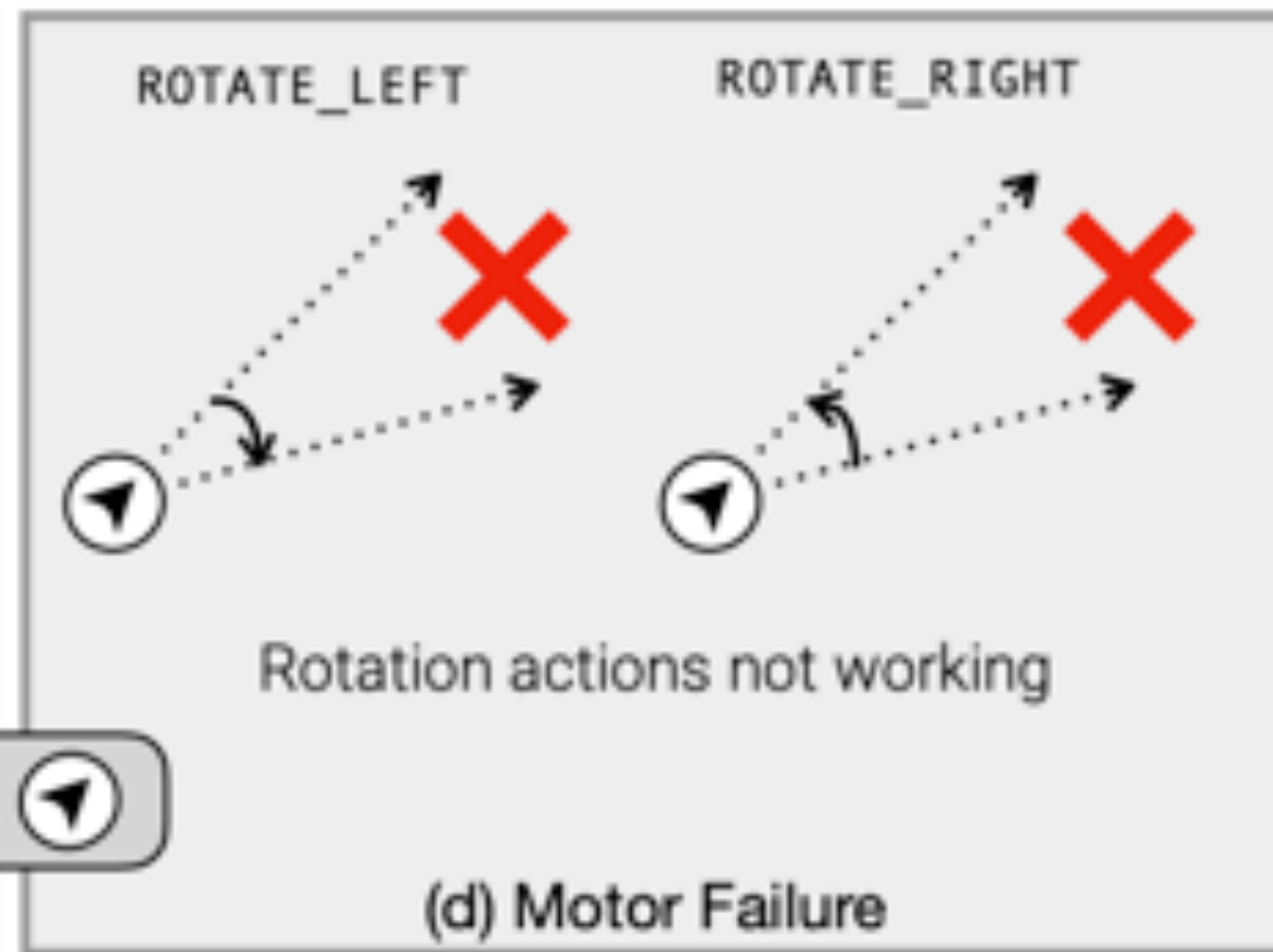
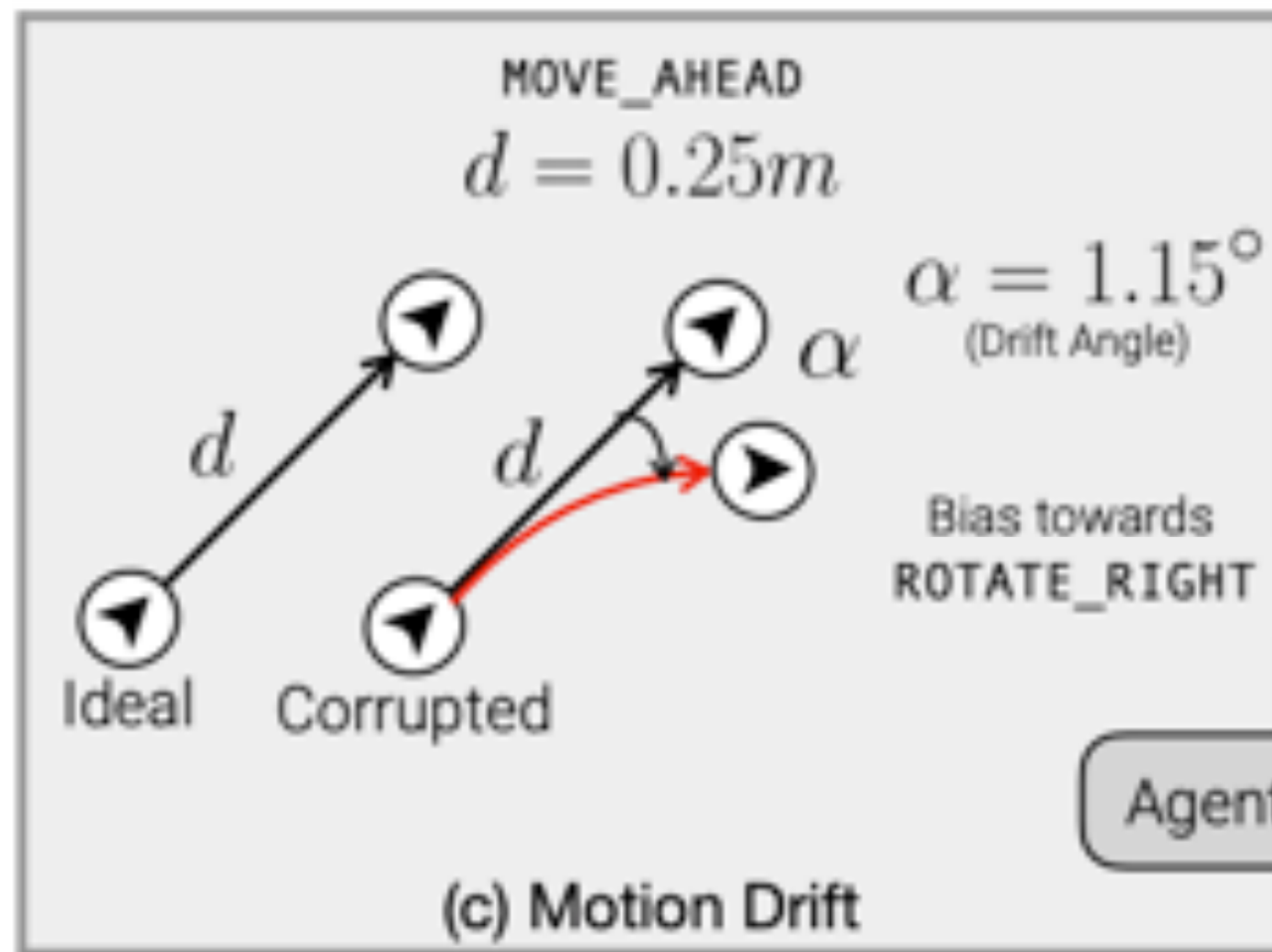
Due to Environment

Scene-level friction



High and low friction zones

Due to Faulty Movements



Malfunctioning components

ObjectNav RGBD — Target Object in “Blue”

Clean Conditions
(Success = True)

RGB

Depth

Top-Down



Synthetic to Real Pixel Adaptation

Train



GTA (synthetic)

Test



CityScapes (Germany)

Domain Adaptation: Train on Source Test on Target

Domain Adaptation: Train on Source Test on Target

Source Domain $\sim P_S(X_S, Y_S)$

lots of **labeled** data

Domain Adaptation: Train on Source Test on Target



Source Domain $\sim P_S(X_S, Y_S)$

lots of **labeled** data

Domain Adaptation: Train on Source Test on Target

backpack



.....

chair



bike



Source Domain $\sim P_S(X_S, Y_S)$
lots of **labeled** data

Domain Adaptation: Train on Source Test on Target



Source Domain $\sim P_S(X_S, Y_S)$

lots of **labeled** data

Domain Adaptation: Train on Source Test on Target

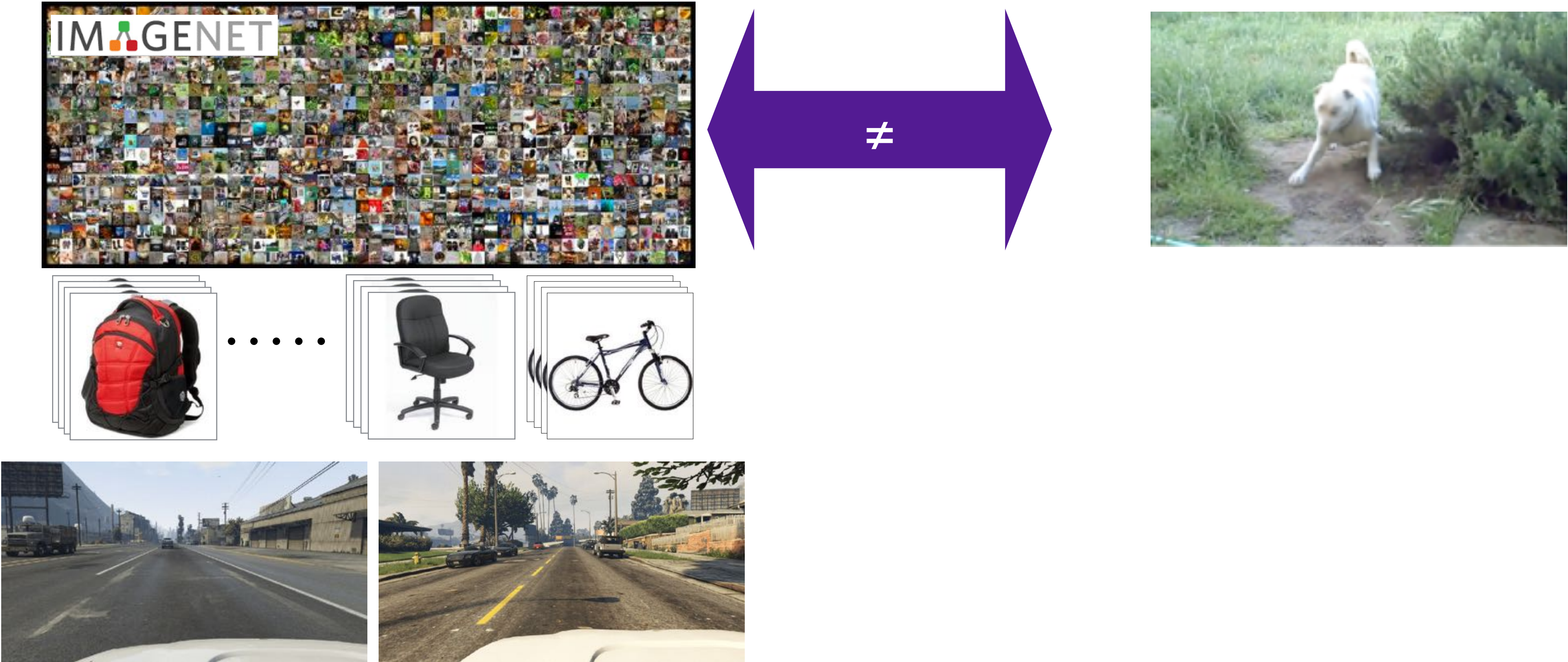


Source Domain $\sim P_S(X_S, Y_S)$
lots of **labeled** data

\neq

Target Domain $\sim P_T(X_T, Y_T)$
unlabeled or limited labels

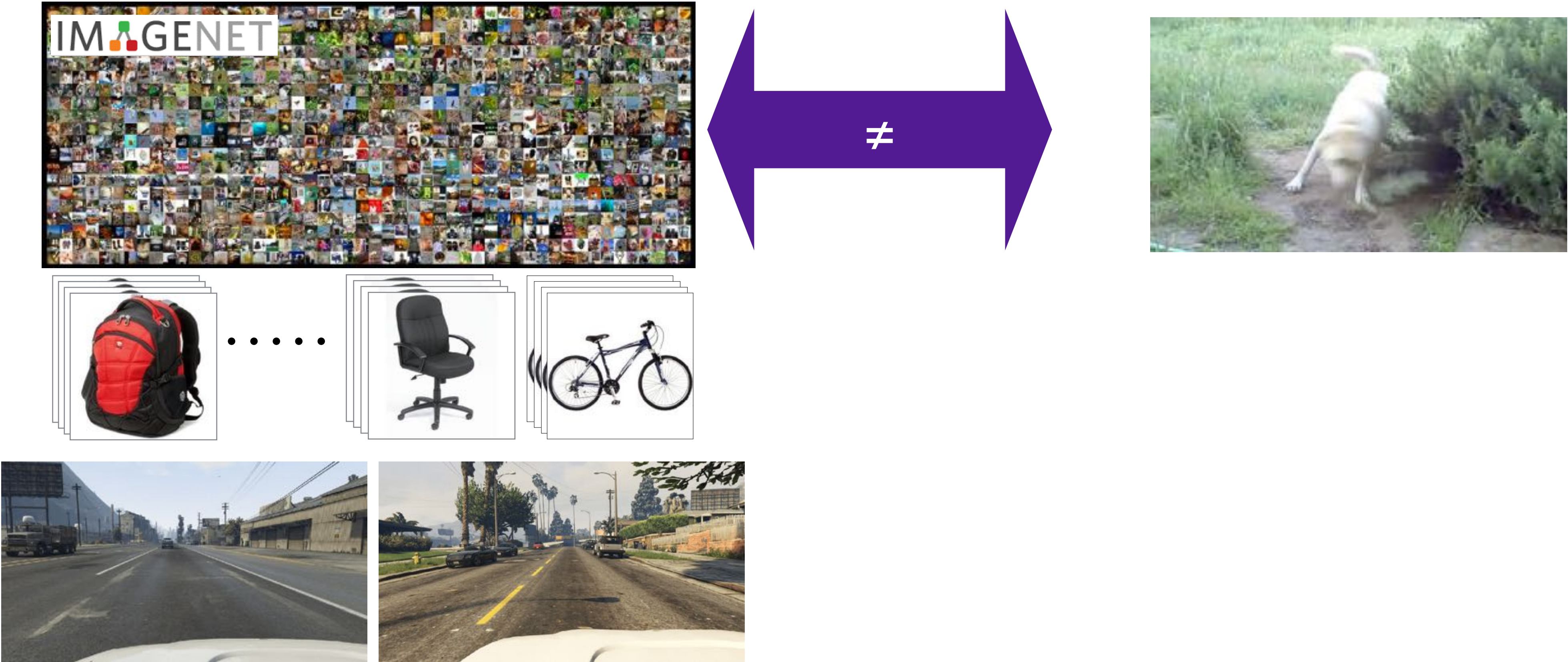
Domain Adaptation: Train on Source Test on Target



Source Domain $\sim P_S(X_S, Y_S)$
lots of **labeled** data

Target Domain $\sim P_T(X_T, Y_T)$
unlabeled or limited labels

Domain Adaptation: Train on Source Test on Target



Source Domain $\sim P_S(X_S, Y_S)$
lots of **labeled** data

Target Domain $\sim P_T(X_T, Y_T)$
unlabeled or limited labels

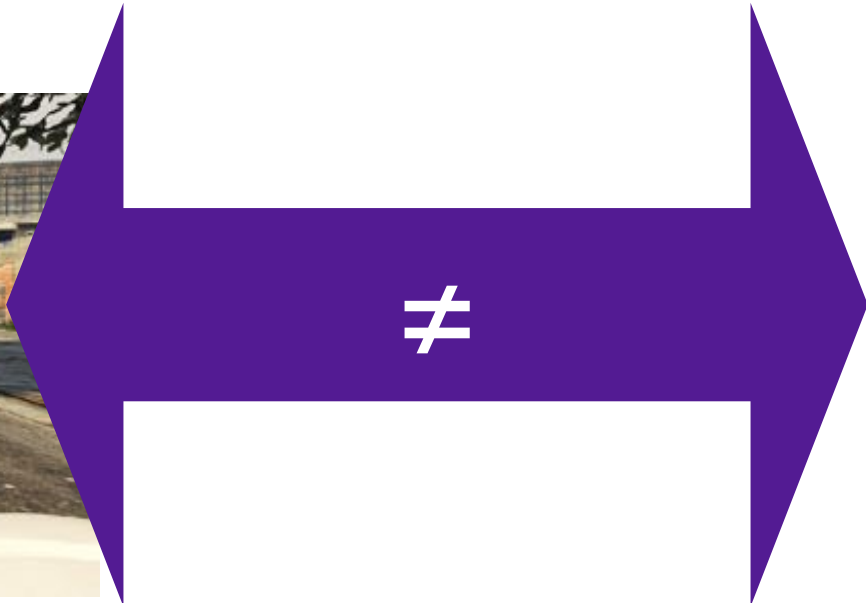
Domain Adaptation: Train on Source Test on Target



Source Domain $\sim P_S(X_S, Y_S)$
lots of **labeled** data

Target Domain $\sim P_T(X_T, Y_T)$
unlabeled or limited labels

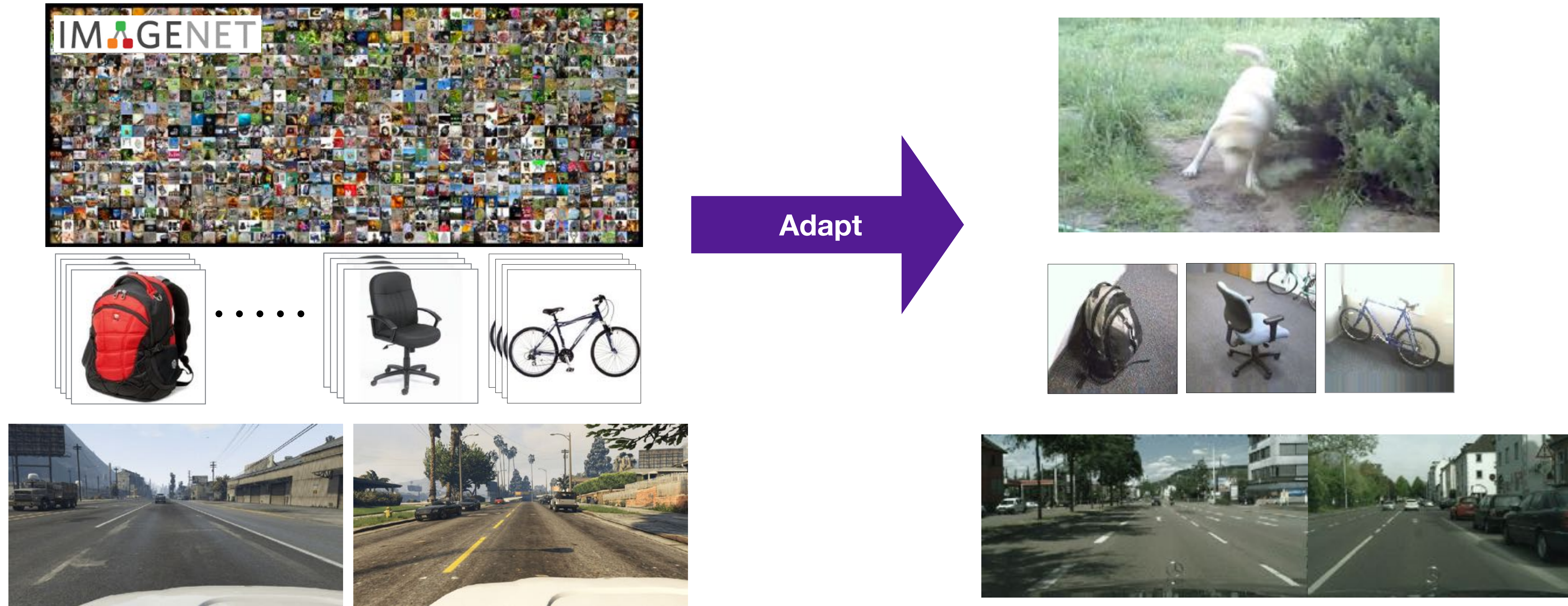
Domain Adaptation: Train on Source Test on Target



Source Domain $\sim P_S(X_S, Y_S)$
lots of **labeled** data

Target Domain $\sim P_T(X_T, Y_T)$
unlabeled or limited labels

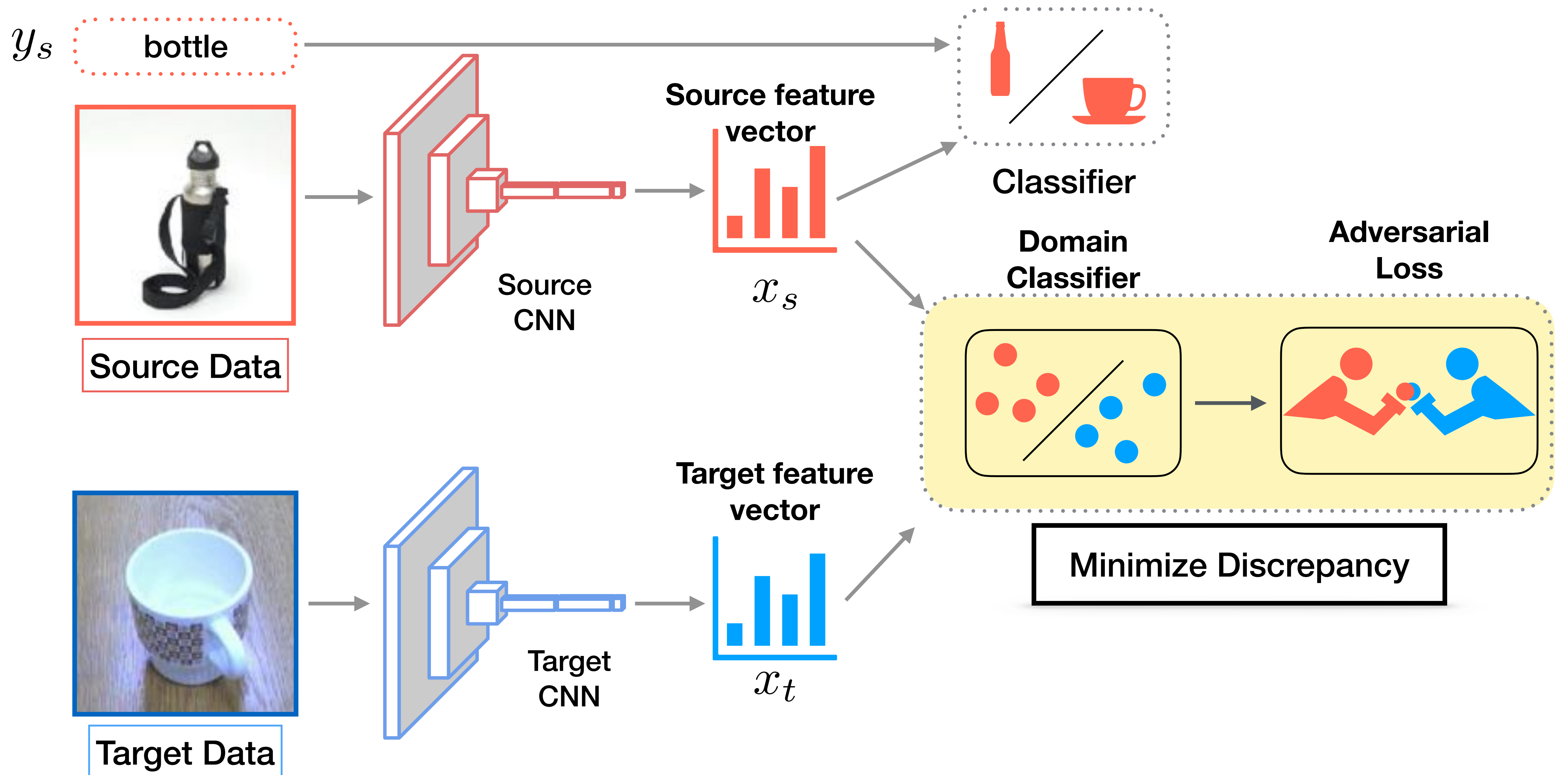
Domain Adaptation: Train on Source Test on Target



Source Domain $\sim P_S(X_S, Y_S)$
lots of **labeled** data

Target Domain $\sim P_T(X_T, Y_T)$
unlabeled or limited labels

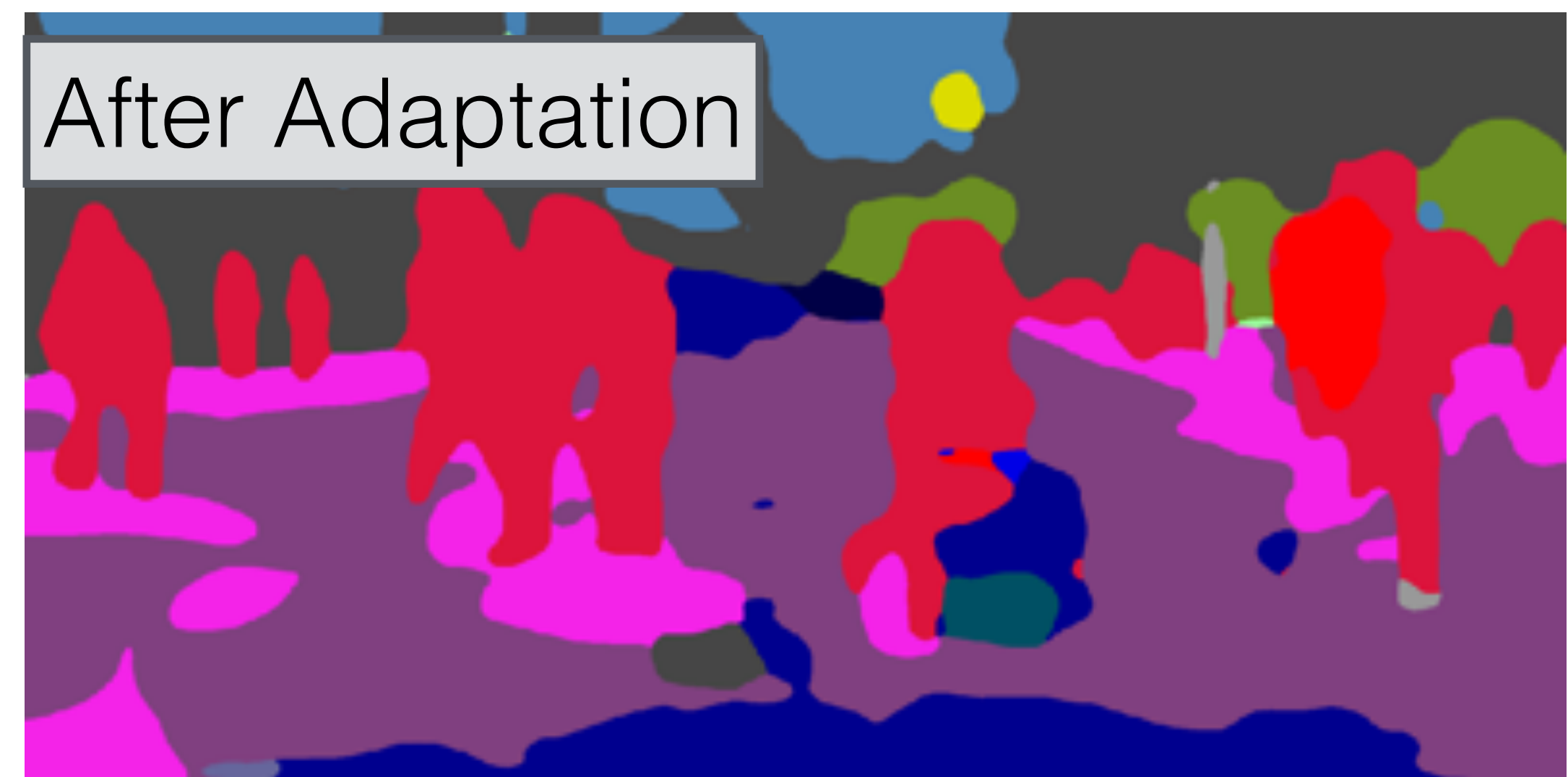
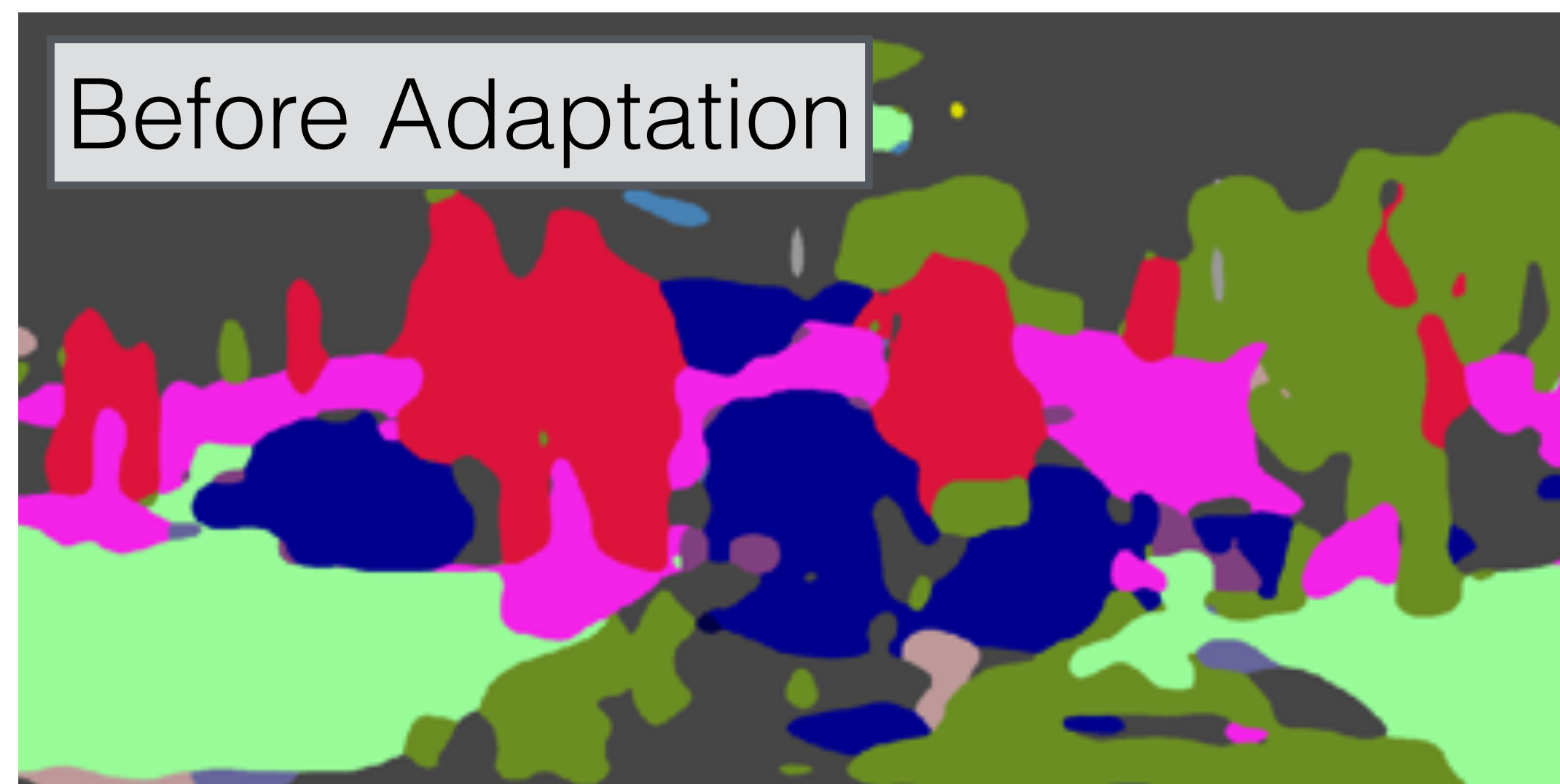
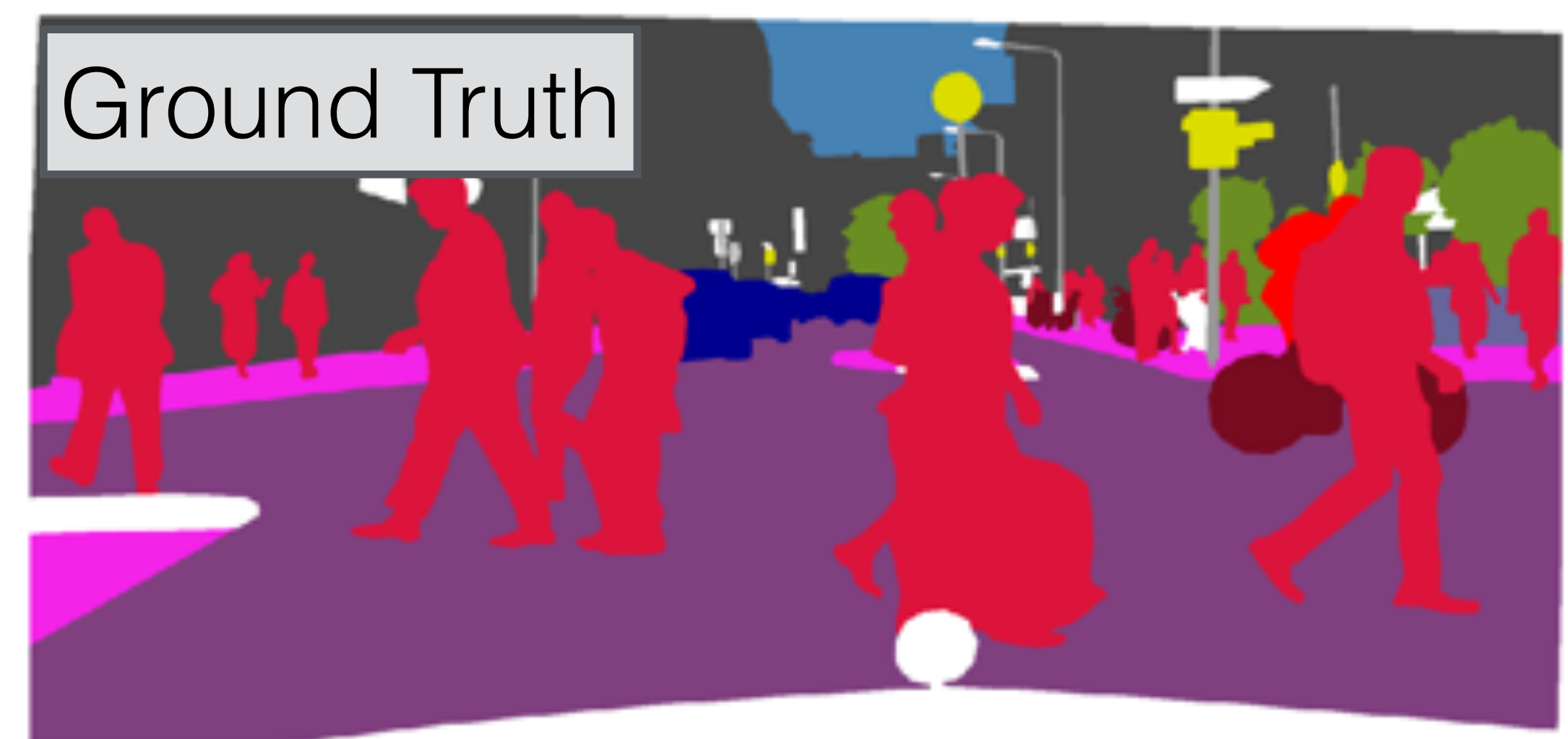
Domain Adversarial Adaptation



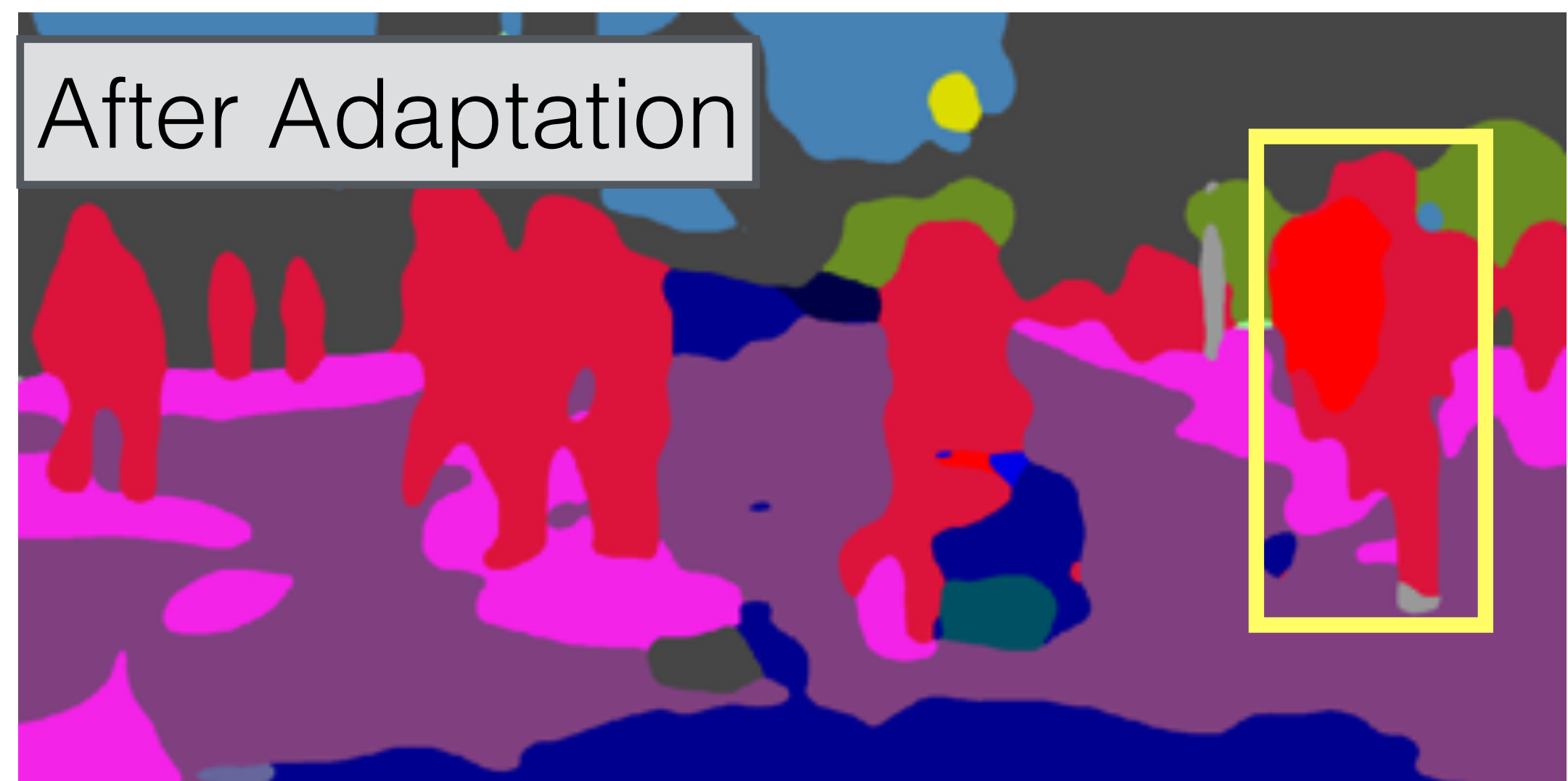
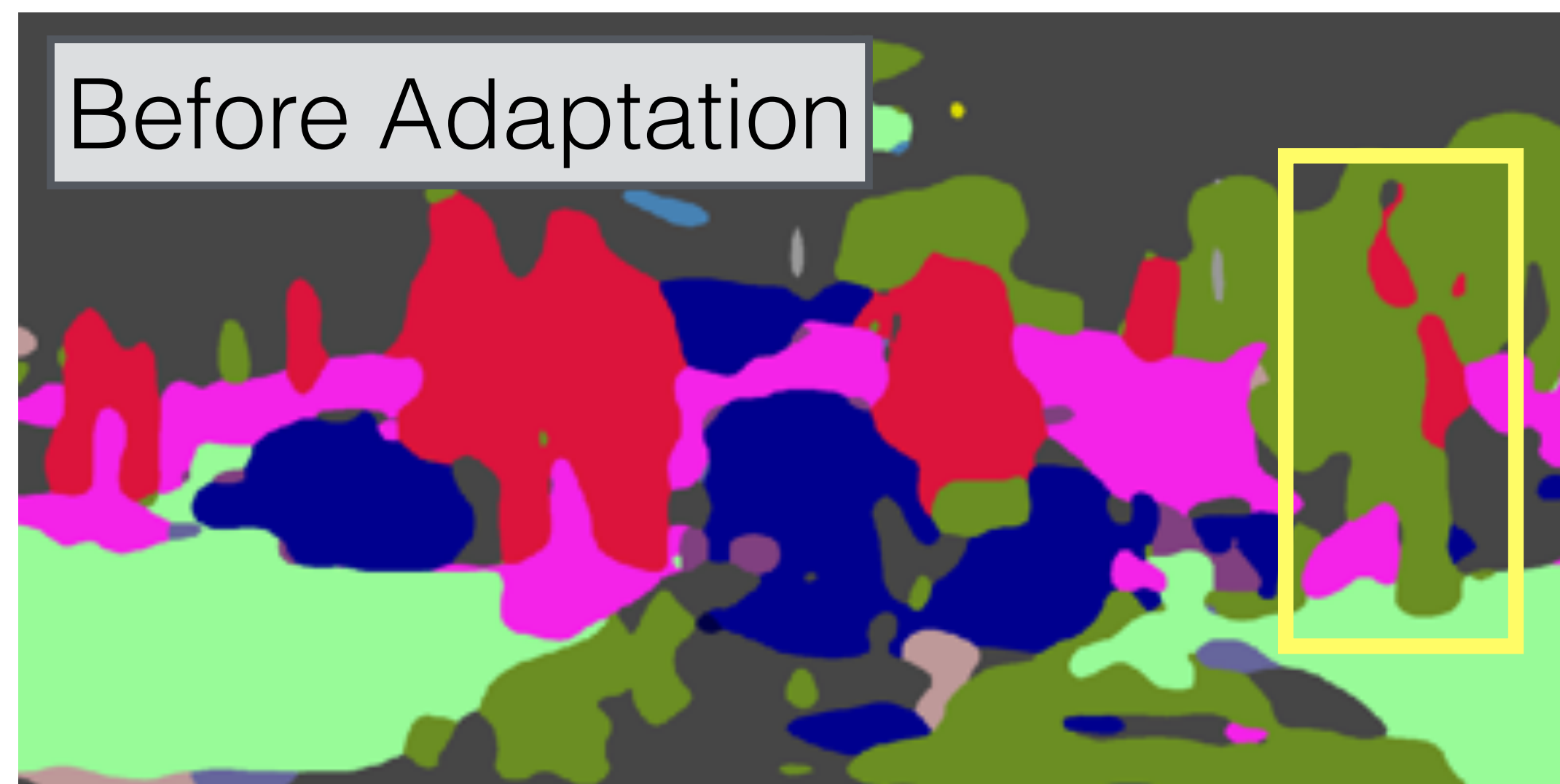
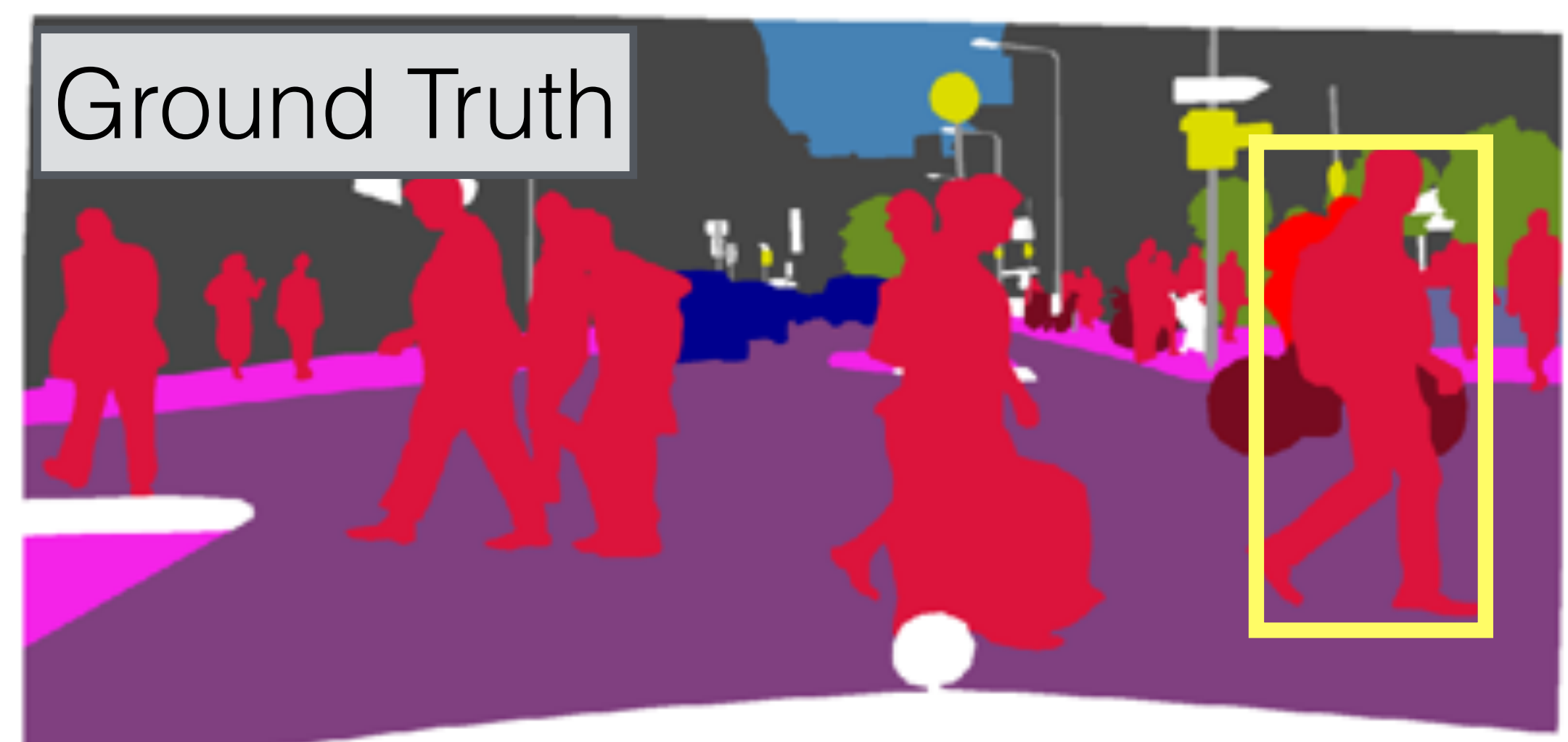
Synthetic to Real Pixel Adaptation



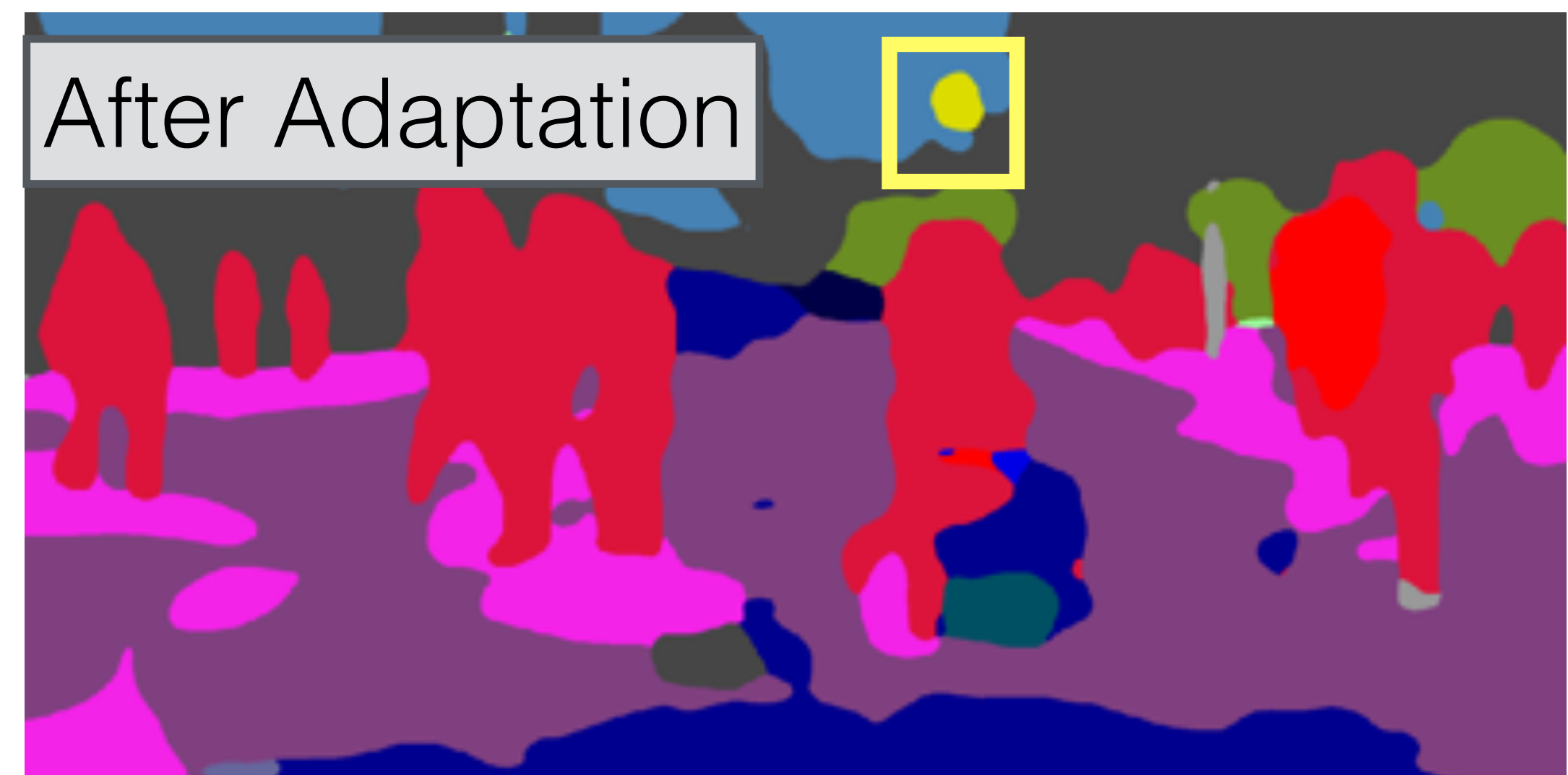
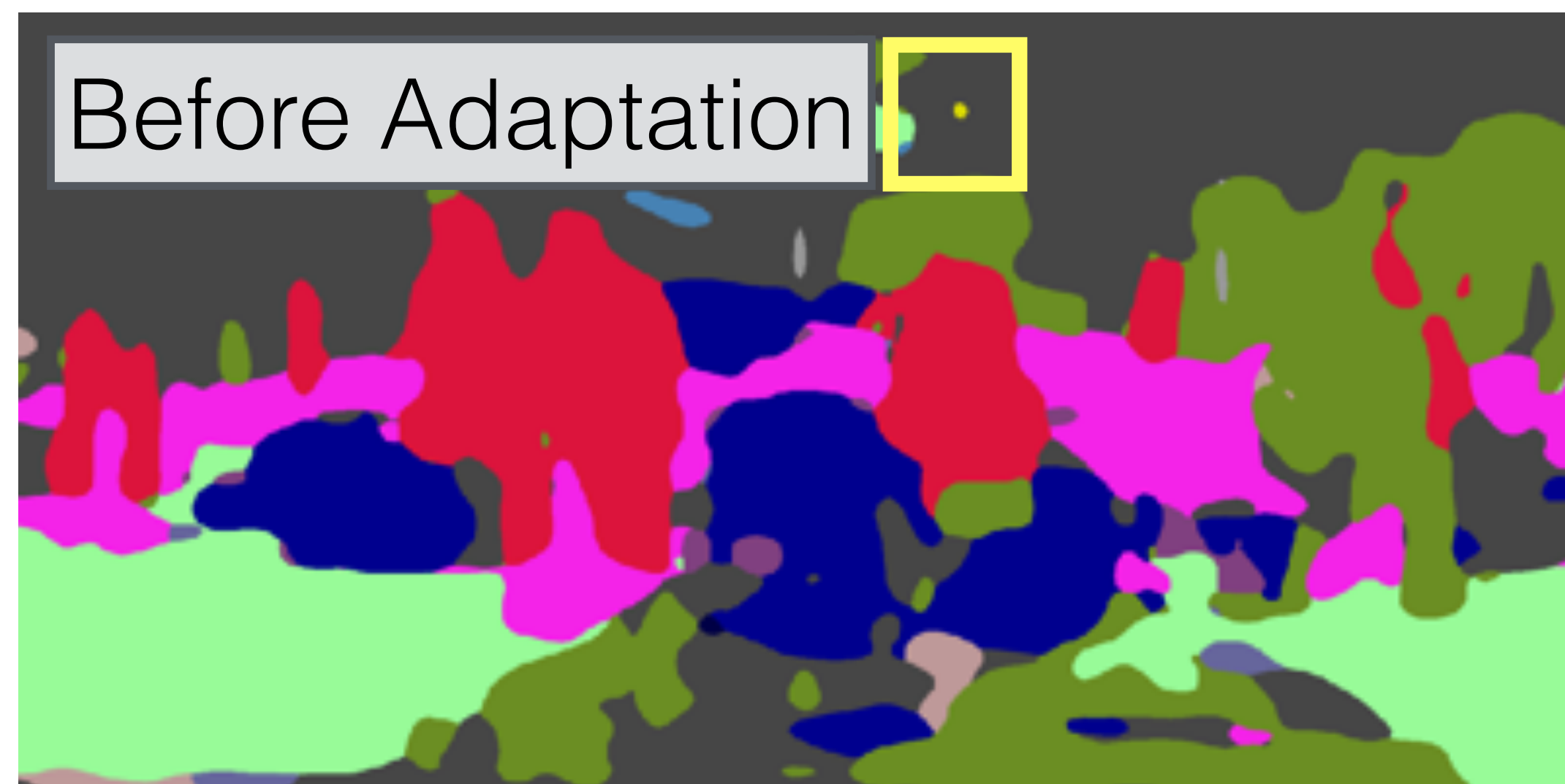
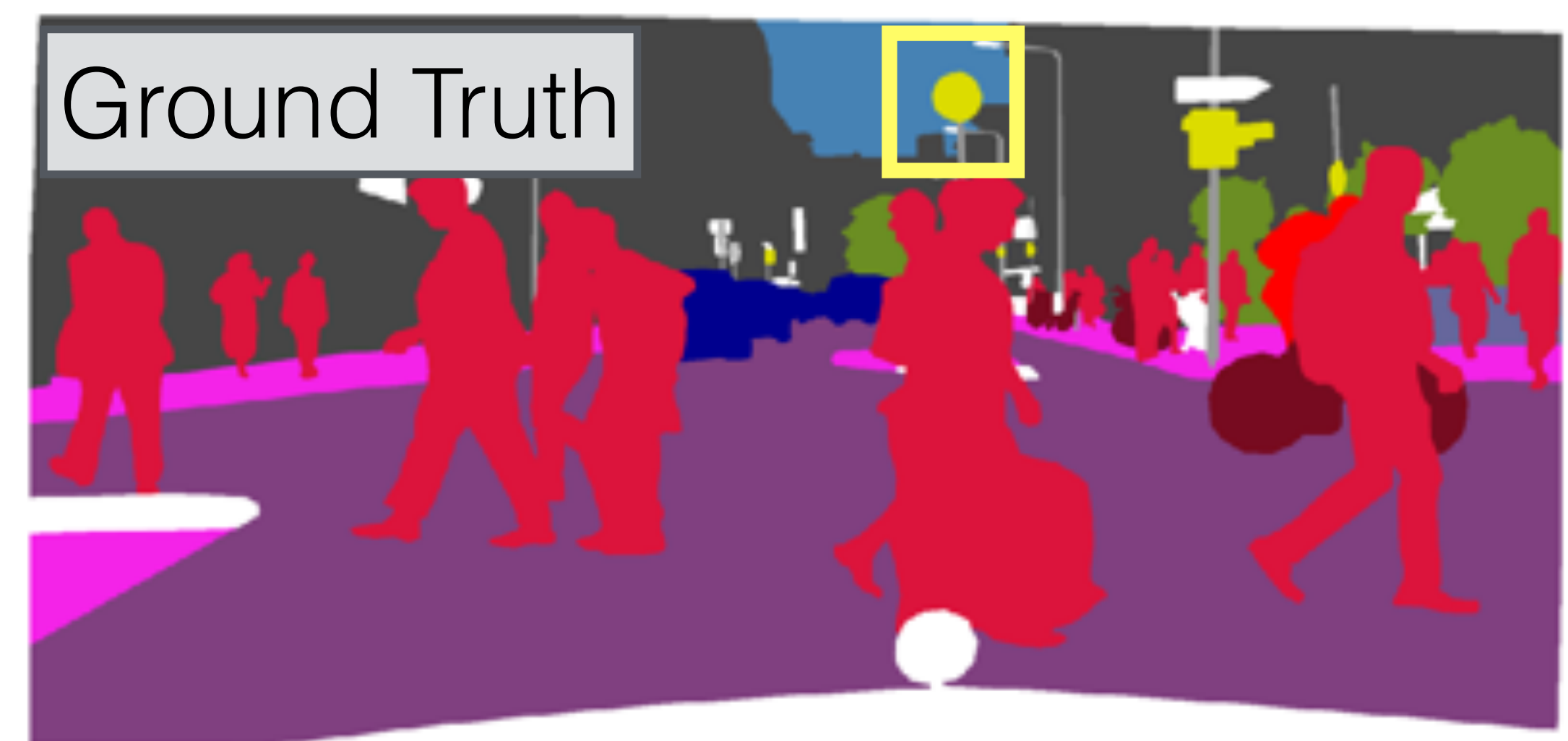
CyCADA Results: CityScapes Evaluation



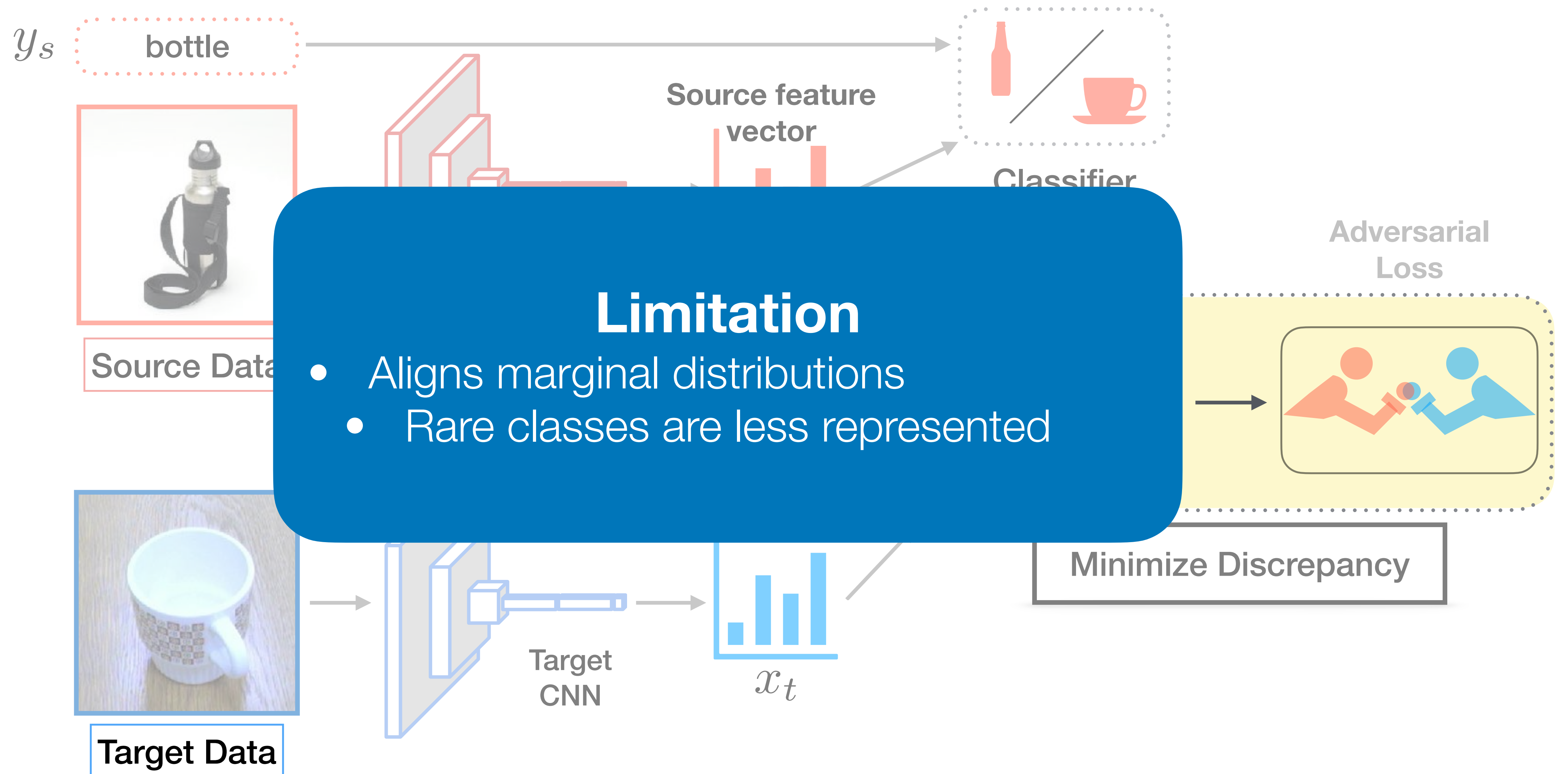
CyCADA Results: CityScapes Evaluation



CyCADA Results: CityScapes Evaluation



Domain Adversarial Adaptation



Adapting to Imbalanced Data

Source data may be curated to be balanced

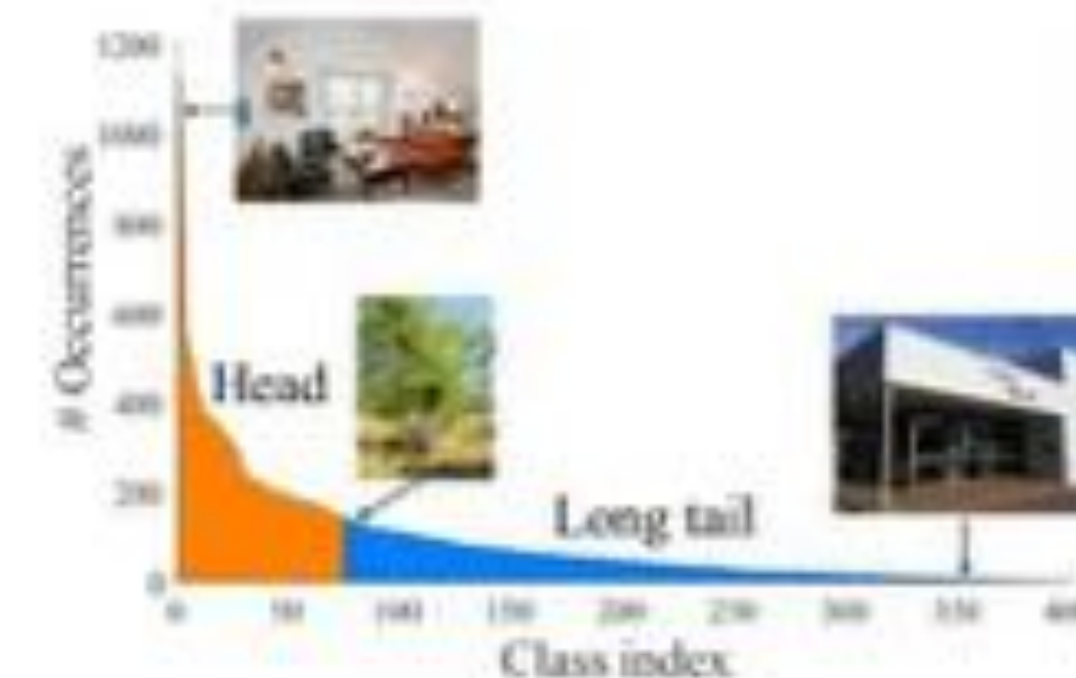
We have no control over target datasets!



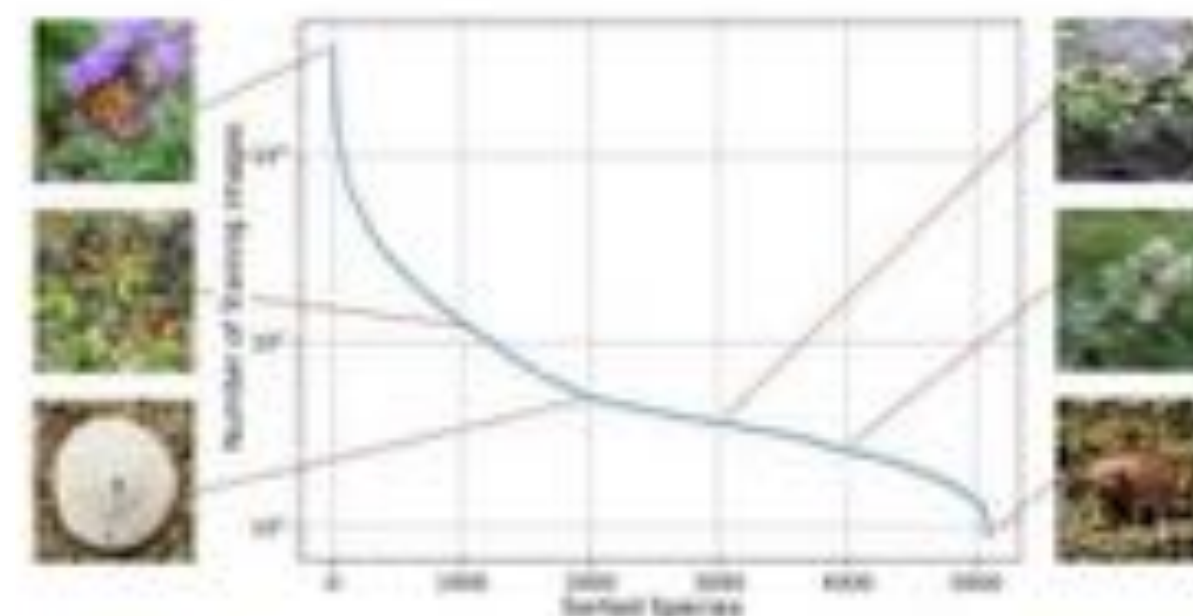
Goal: Adapt under both data and label distribution shift



Faces [Zhang et al. 2017]



Places [Wang et al. 2017]



Species [Van Horn et al. 2019]



Actions [Zhang et al. 2019]

Adapting to Imbalanced Data

- **Challenge:** Existing DA methods (eg. domain adversarial) struggle in this setting!
- Implicitly assume^{1,2} similar label distributions

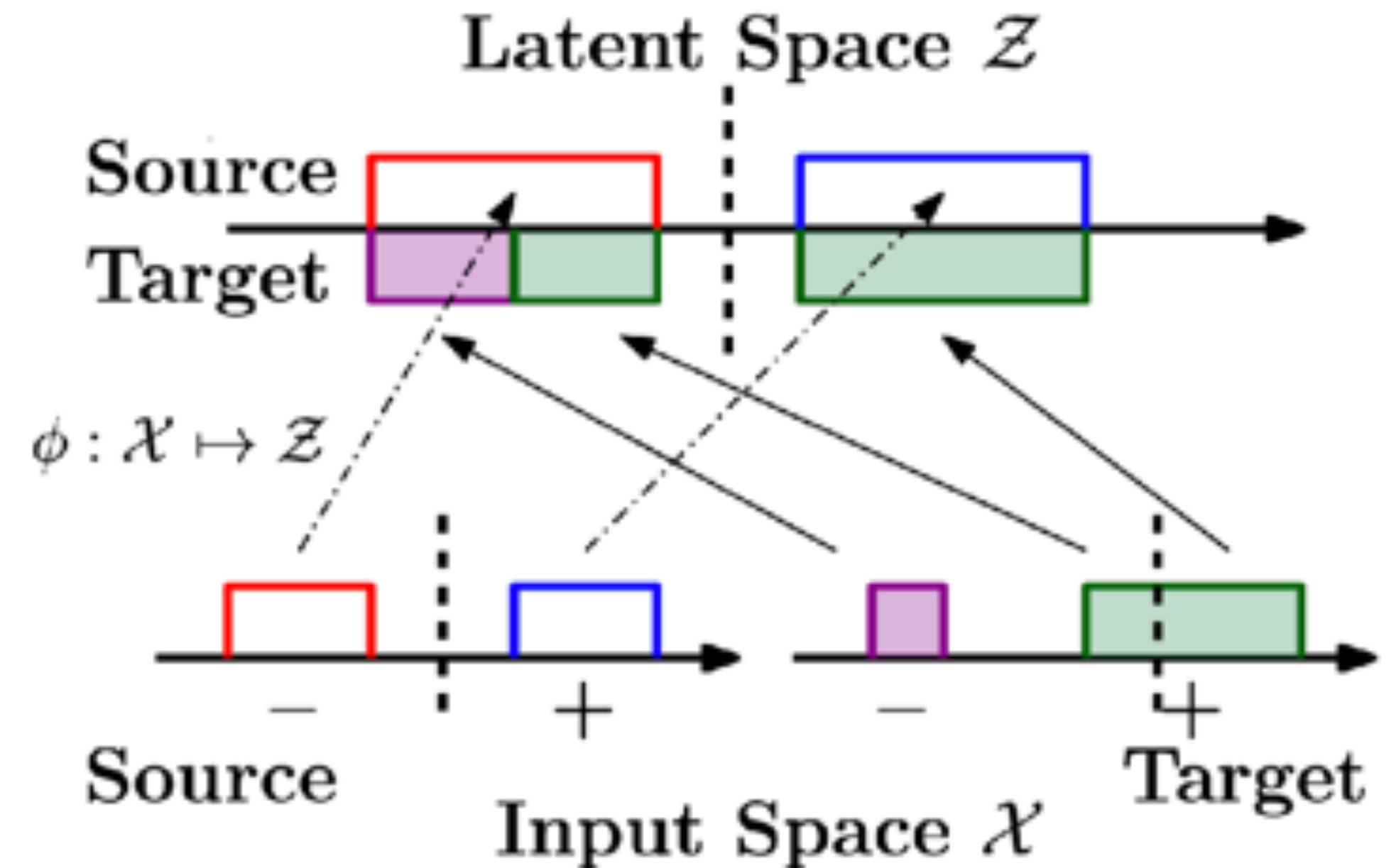
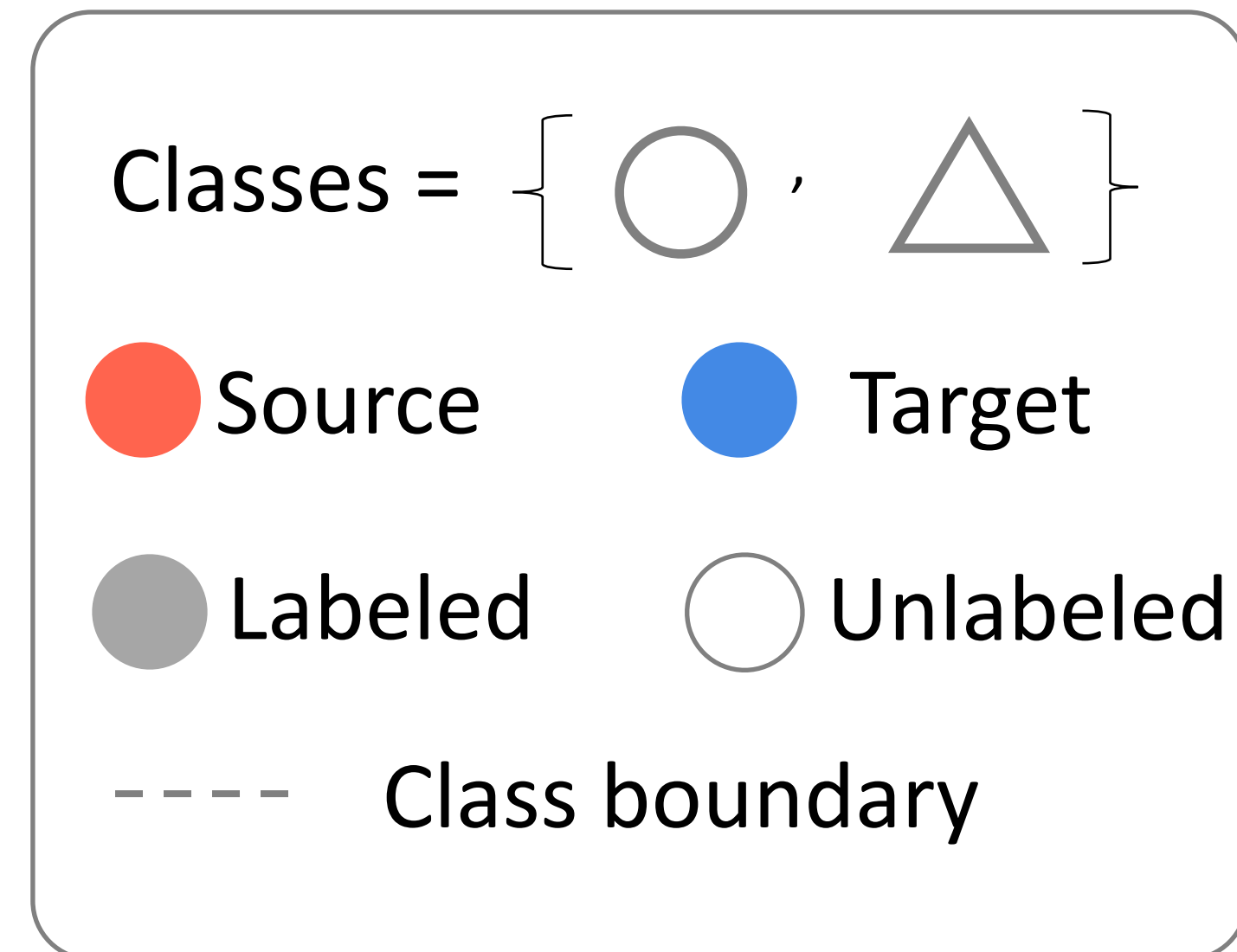
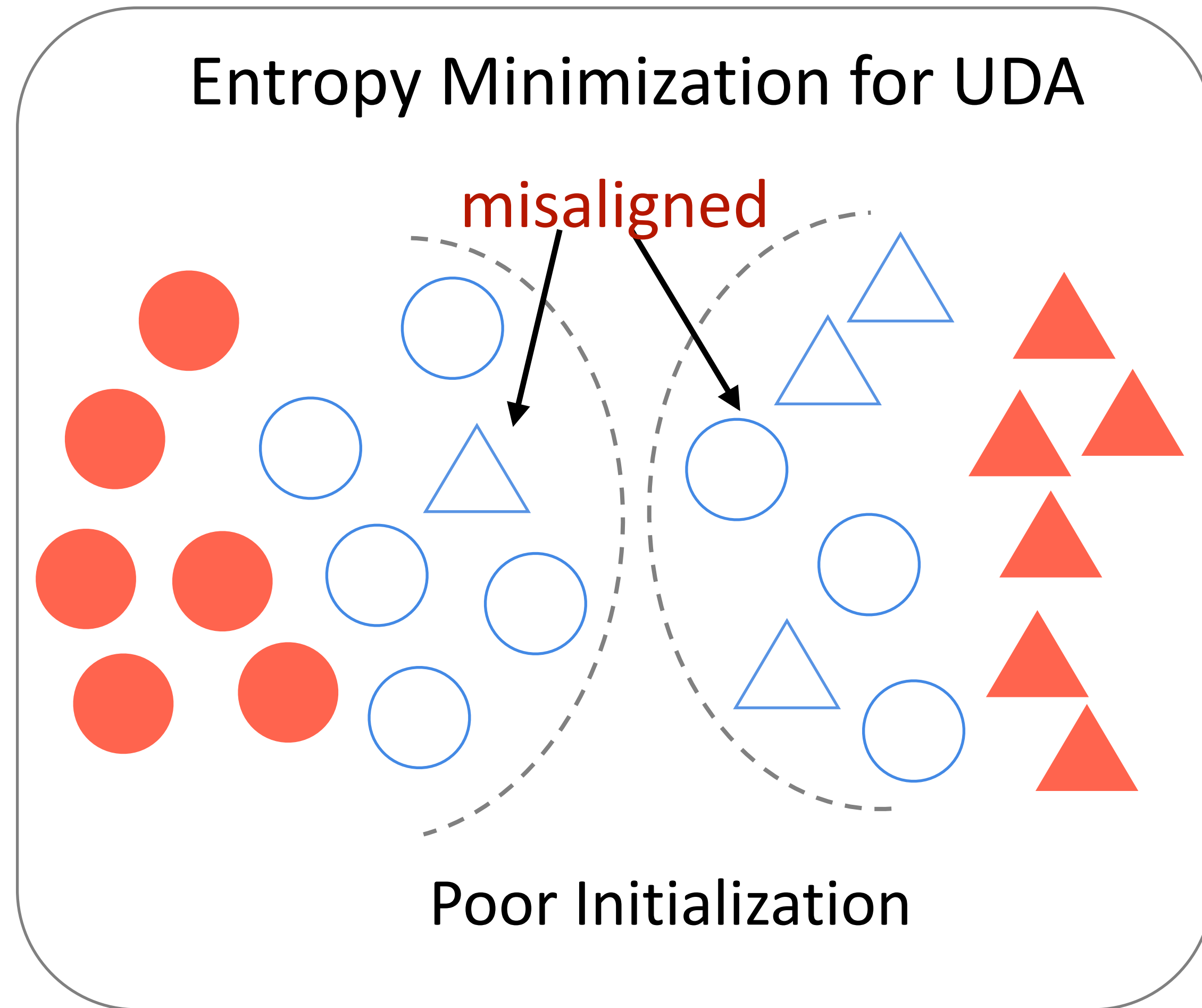


Figure credit: Wu et al., ICML 2019

- We turn to simpler DA approaches based on **self-training**^{3,4}
 - **Algorithm:** Training on model predictions on unlabeled target
 - **No requirement of similar source/target label distributions**

Adaptation with Self-Training

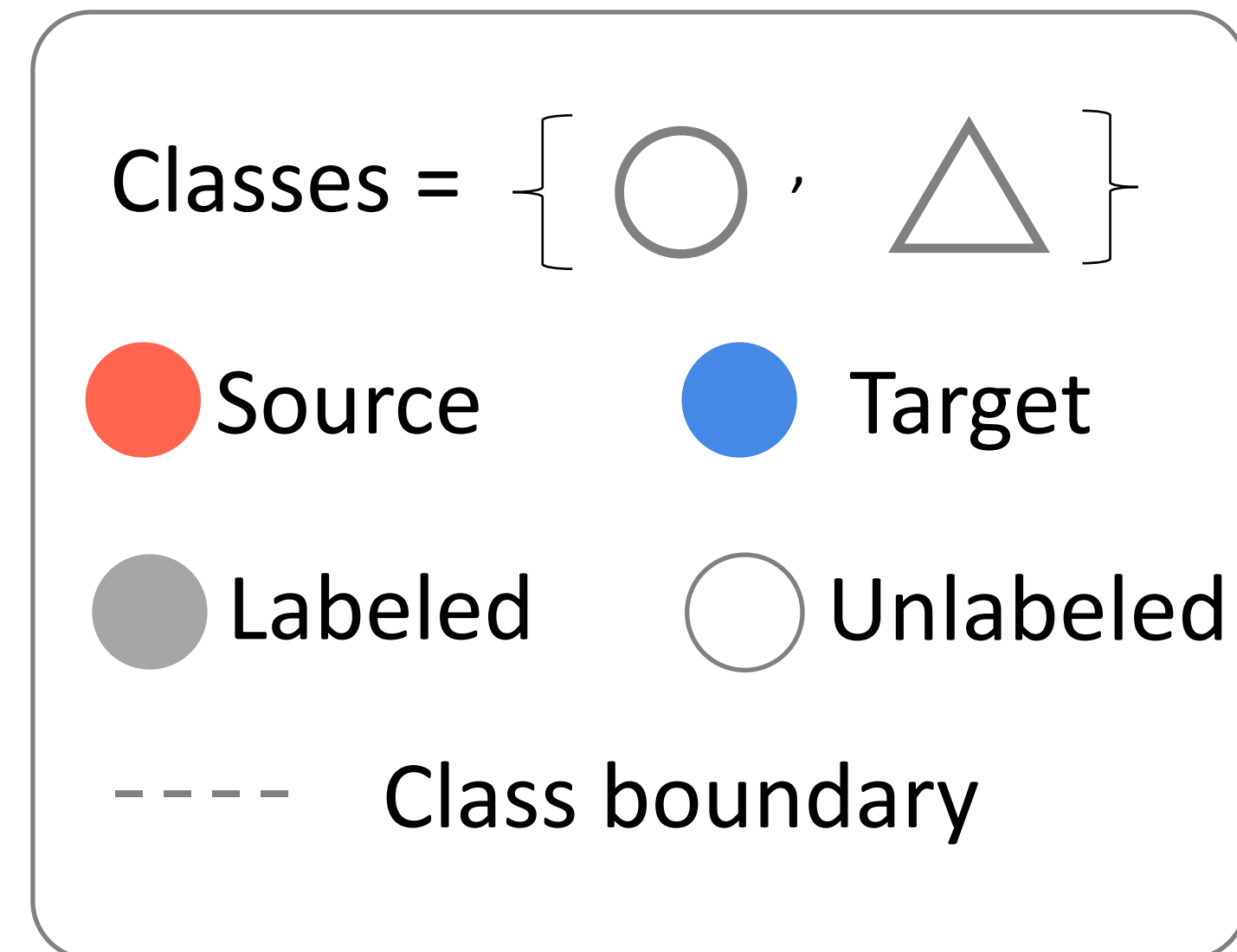
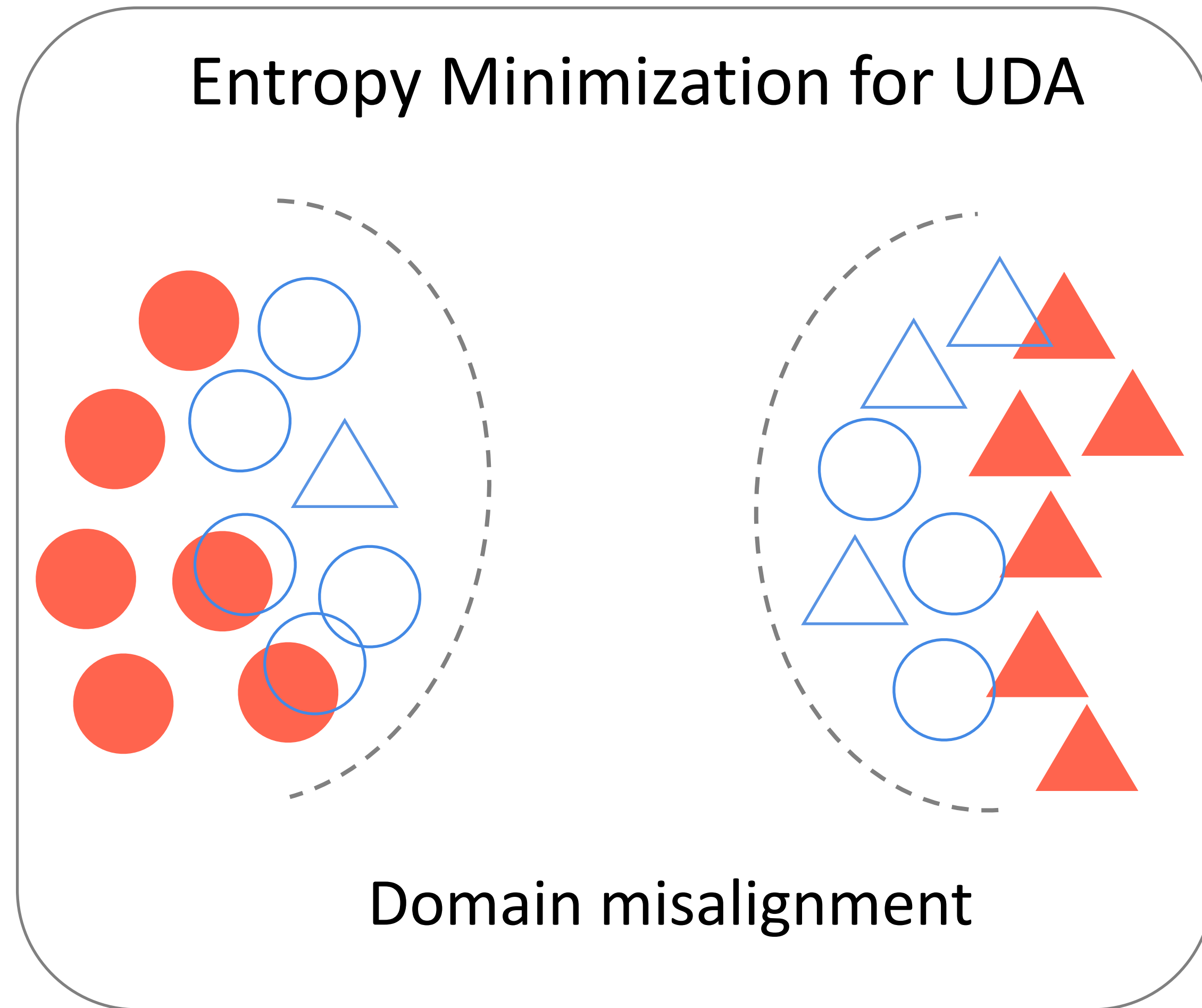
- **Domain Shift:**
Target data is misaligned
- Entropy minimization can **reinforce errors**



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

Adaptation with Self-Training

- **Domain Shift:**
Target data is misaligned
- Entropy minimization can **reinforce errors**



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

Adaptation with Self-Training

- **Domain Shift:** Target data is misaligned

- Entropy minimization can **reinforce errors**

Entropy Minimization for UDA

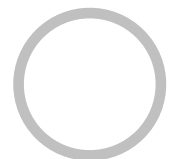
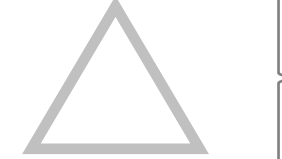
Limitation

- Adapts in response to **all** observations

Goal

- Adapt only on **reliable** observations

Domain misalignment

Classes = { ,  }

 Source

 Target

 Labeled

 Unlabeled

 Class boundary

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

SENTRY

Selective Entropy Optimization via Committee Consistency for Unsupervised Domain Adaptation



Viraj Prabhu



Shivam Khare



Deeksha Karthik

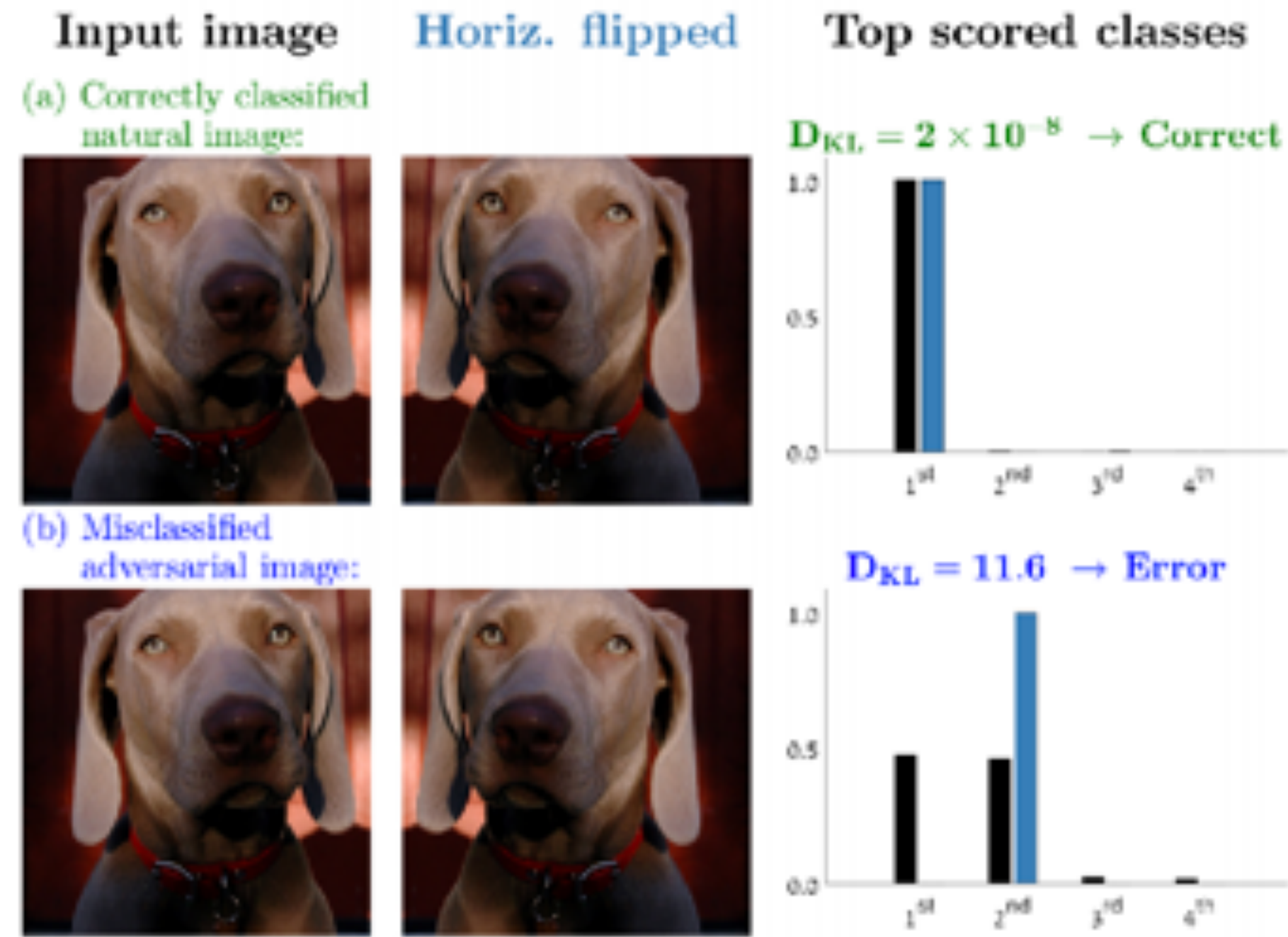


Judy Hoffman



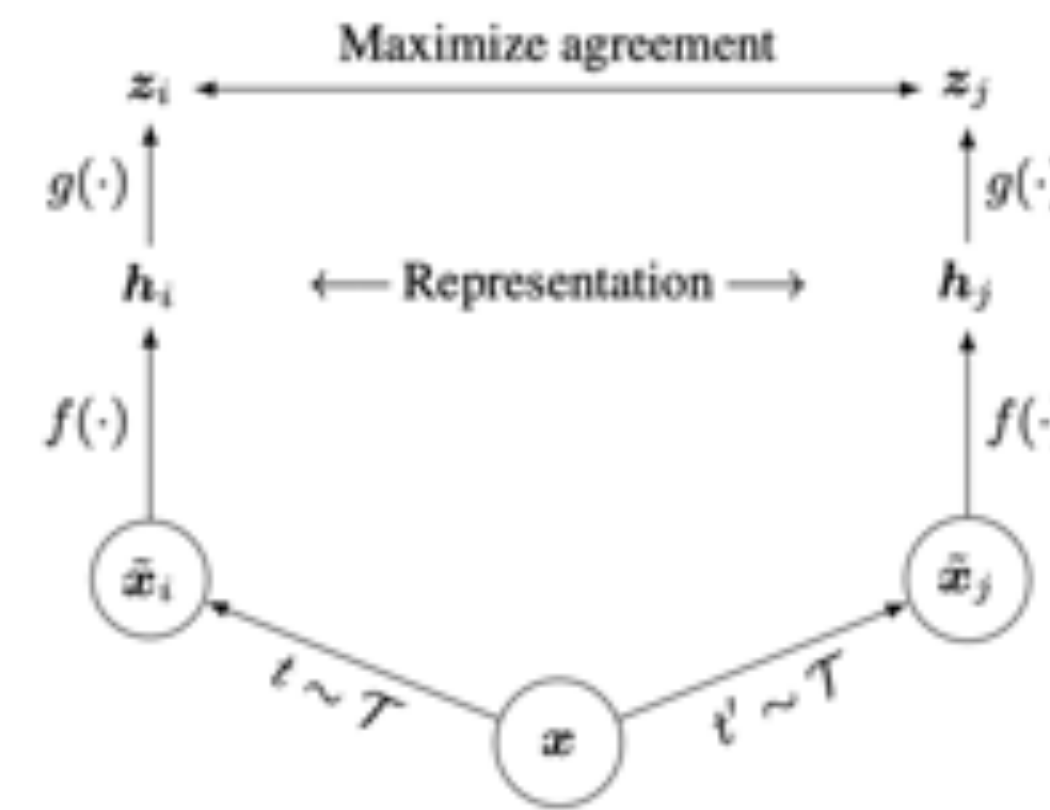
ICCV 2021

Prior Work: Predictive Consistency across Aug

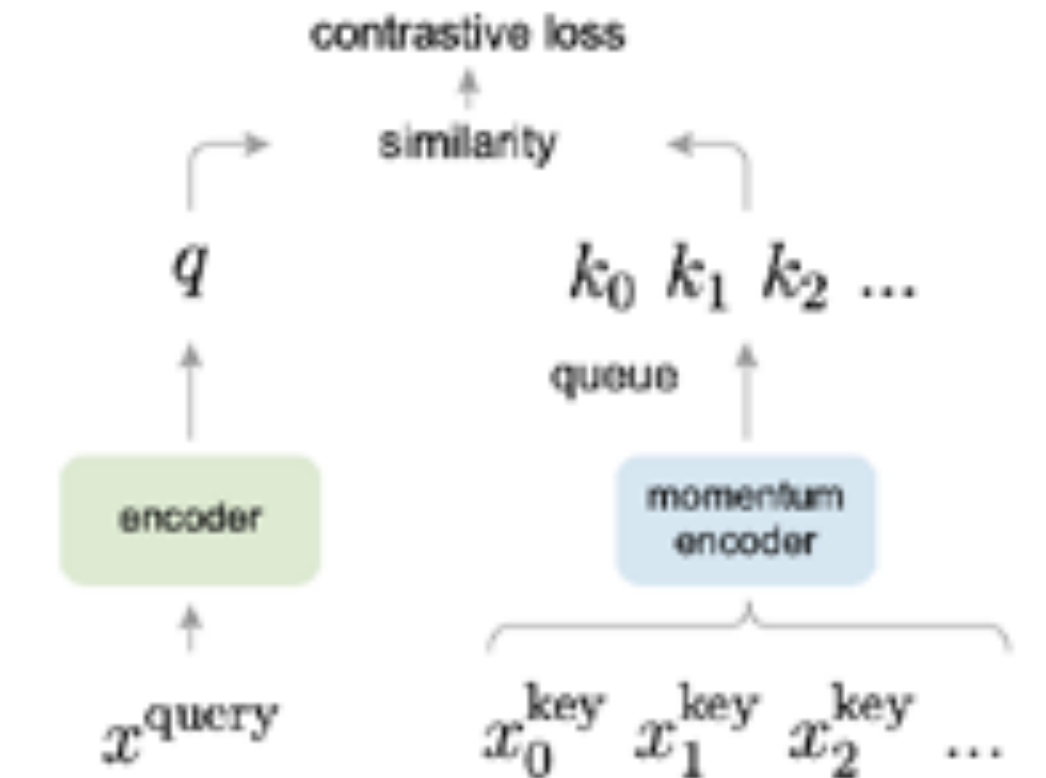


Natural and Adversarial Error Detection using Invariance to Image Transformations.
Irani *et al.*, arXiv 2019

Detecting Errors



SimCLR, Chen *et al.*
ICML 2020



MoCo, He *et al.*
CVPR 2020

Learned Invariance (Contrastive Learning)

SENTRY: Selective Entropy Optimization

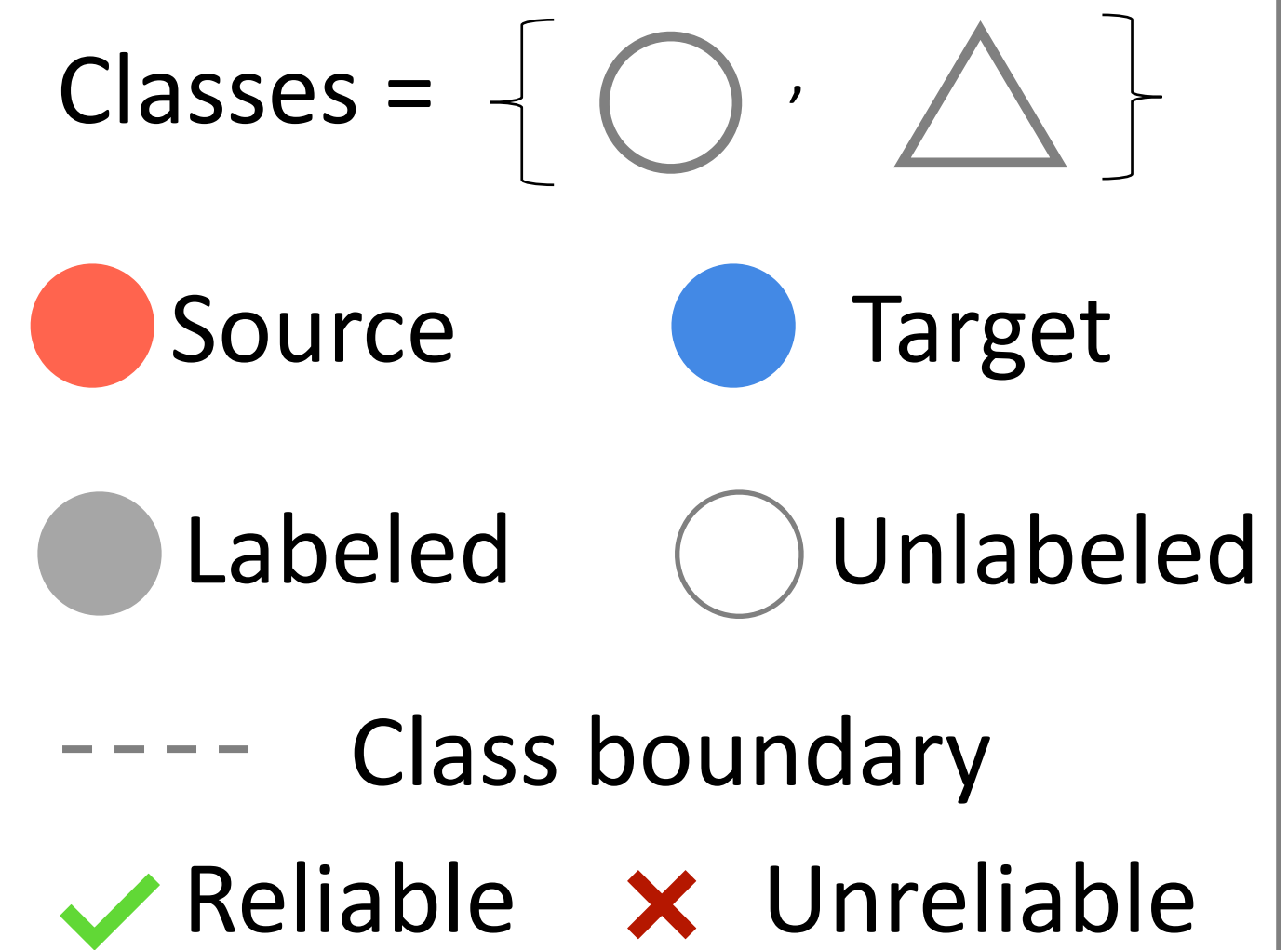
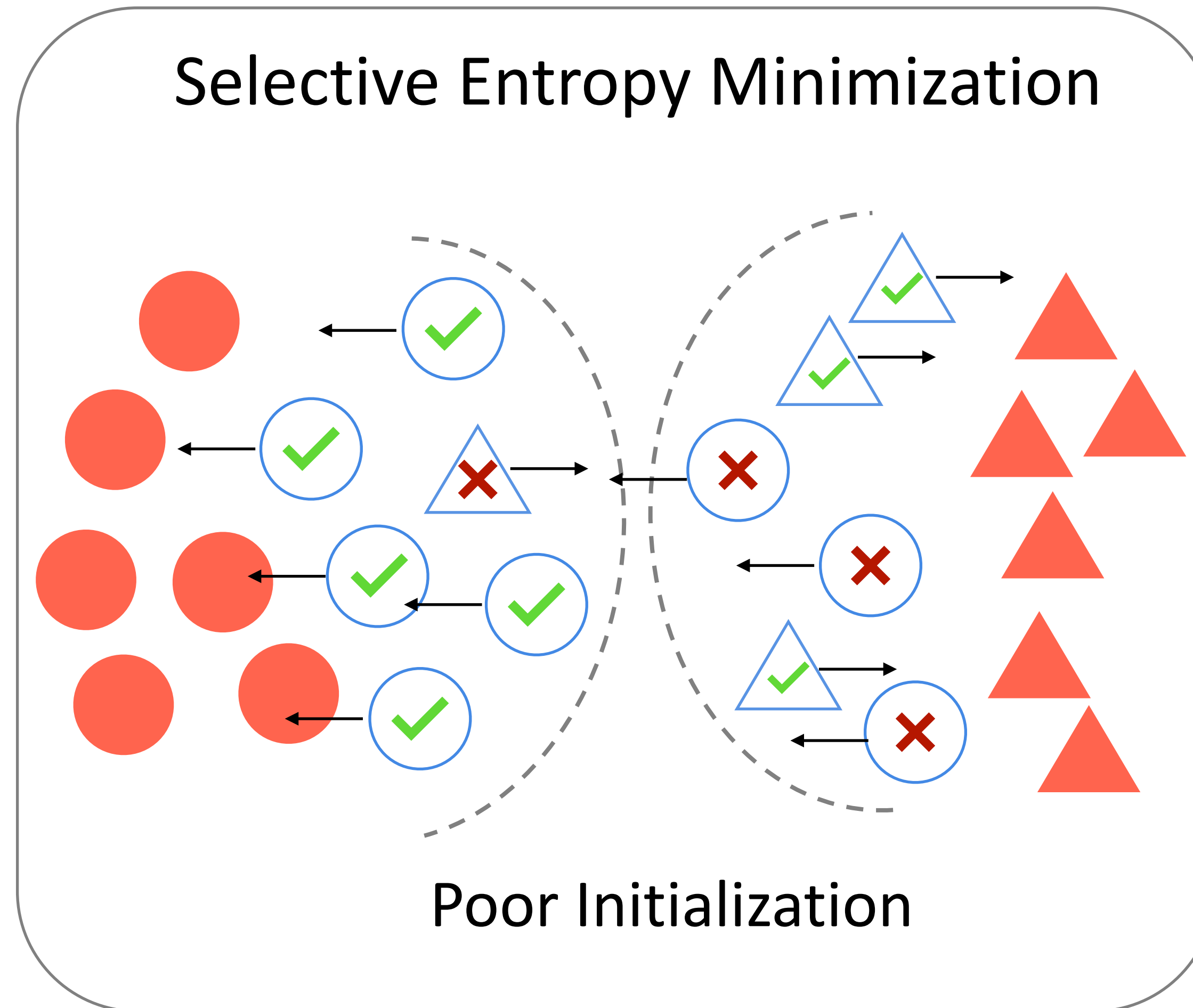
Key Idea

Identify reliable target instances via

~~model confidence~~

Predictive consistency^{1,2,3}

Increase confidence on consistent instances



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

SENTRY: Selective Entropy Optimization

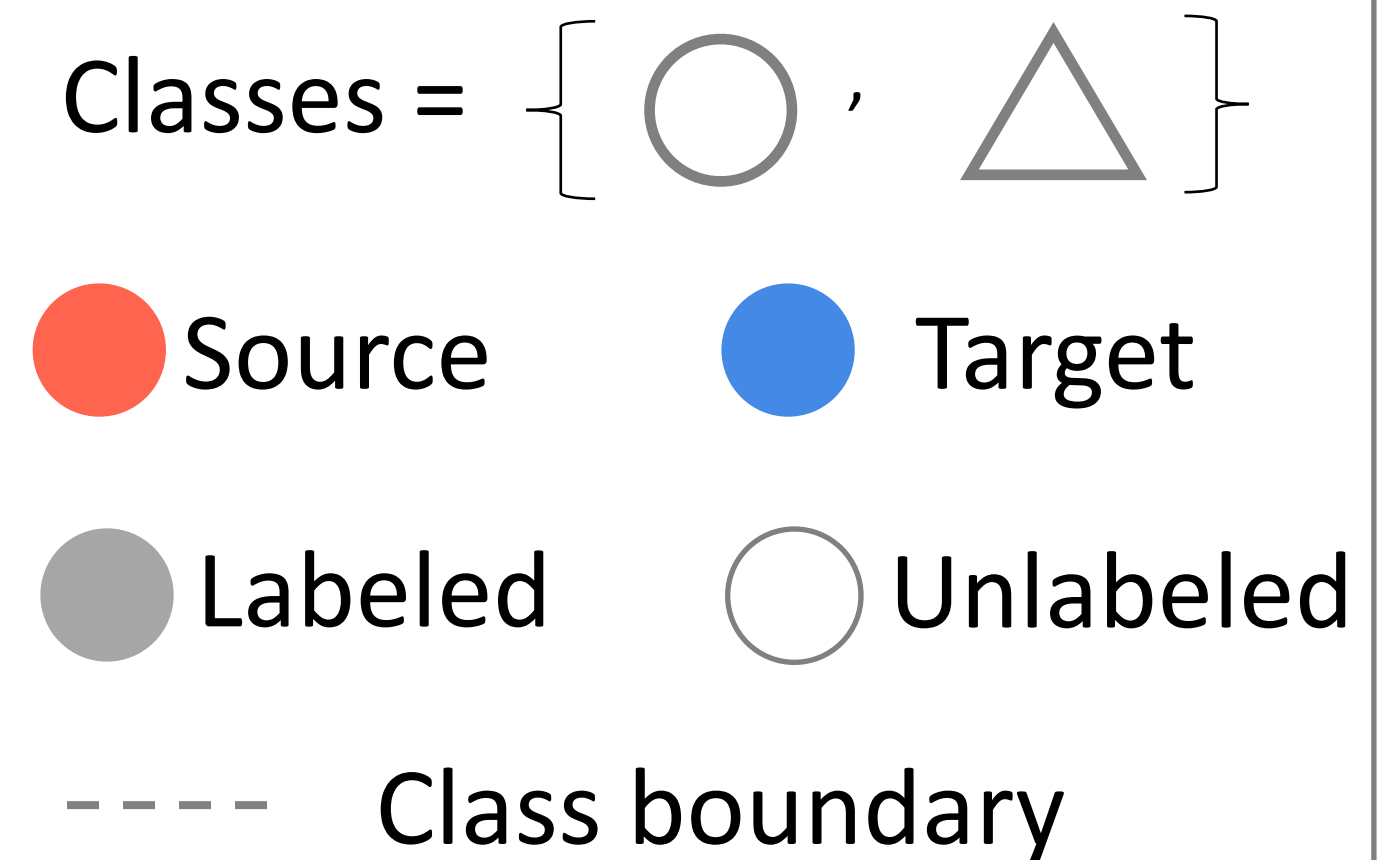
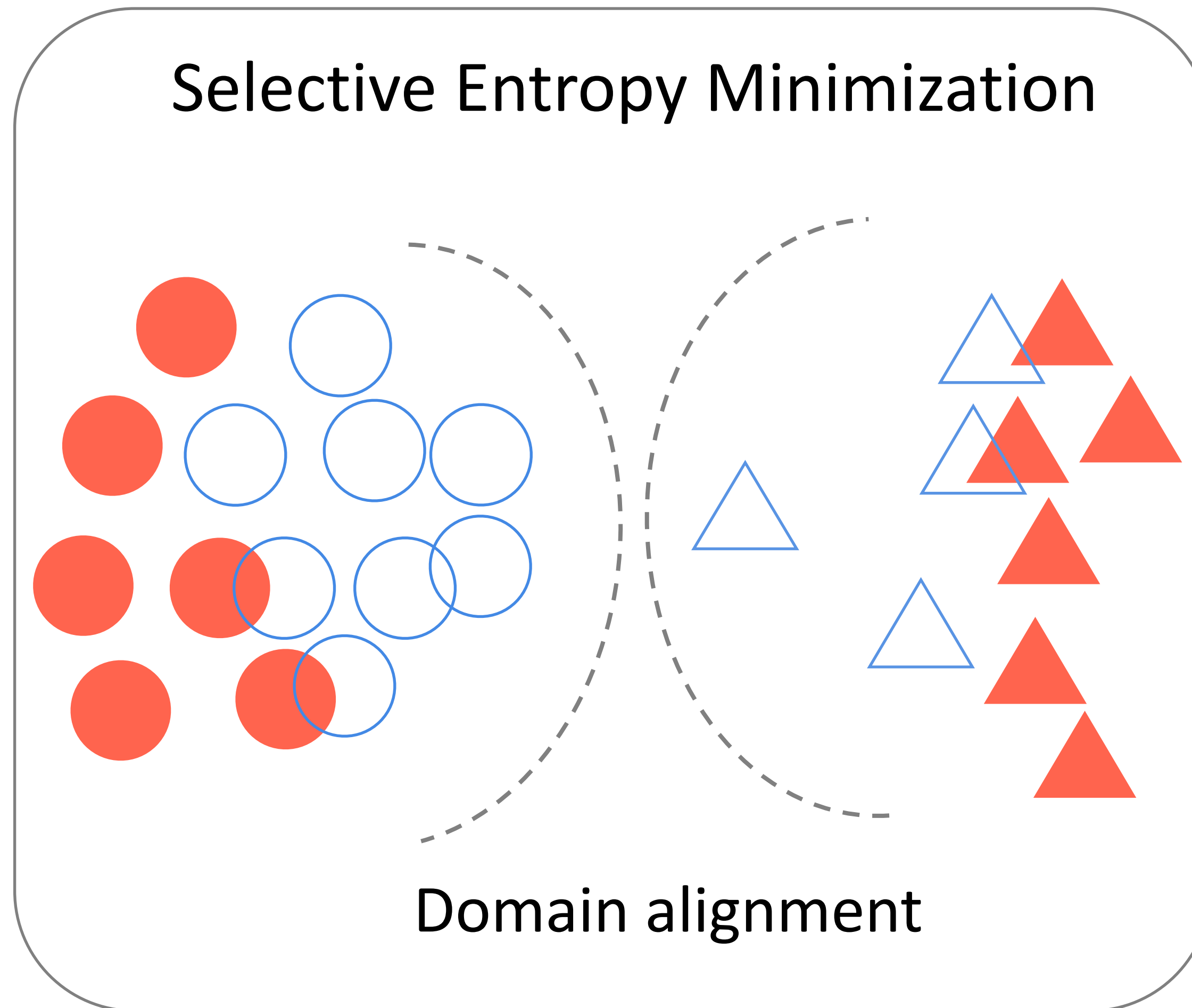
Key Idea

Identify reliable target instances via

~~model confidence~~

Predictive consistency^{1,2,3}

Increase confidence on consistent instances



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

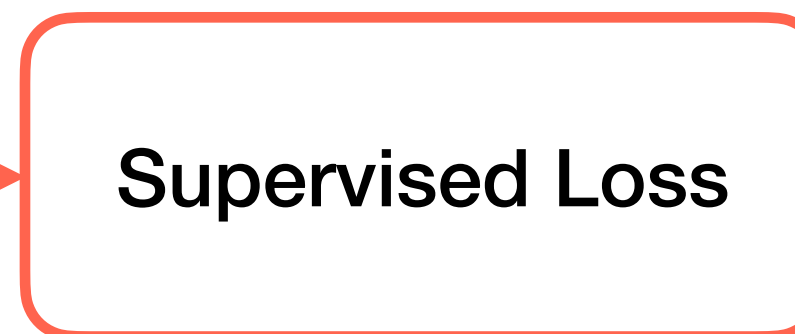
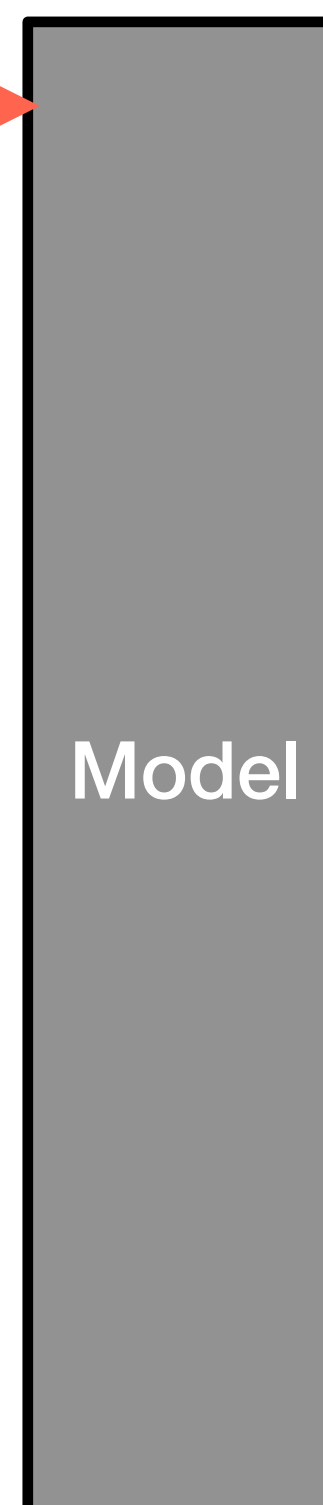
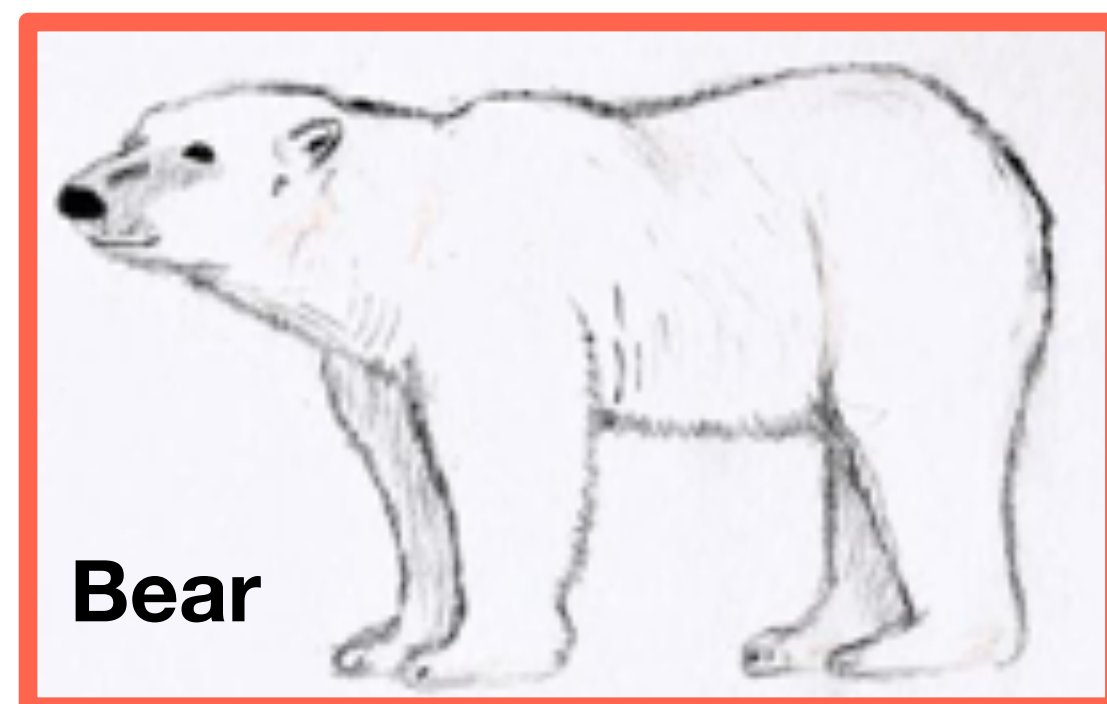
SENTRY: Selective Entropy Optimization via Committee Consistency

Sampled w/
Class Balancing

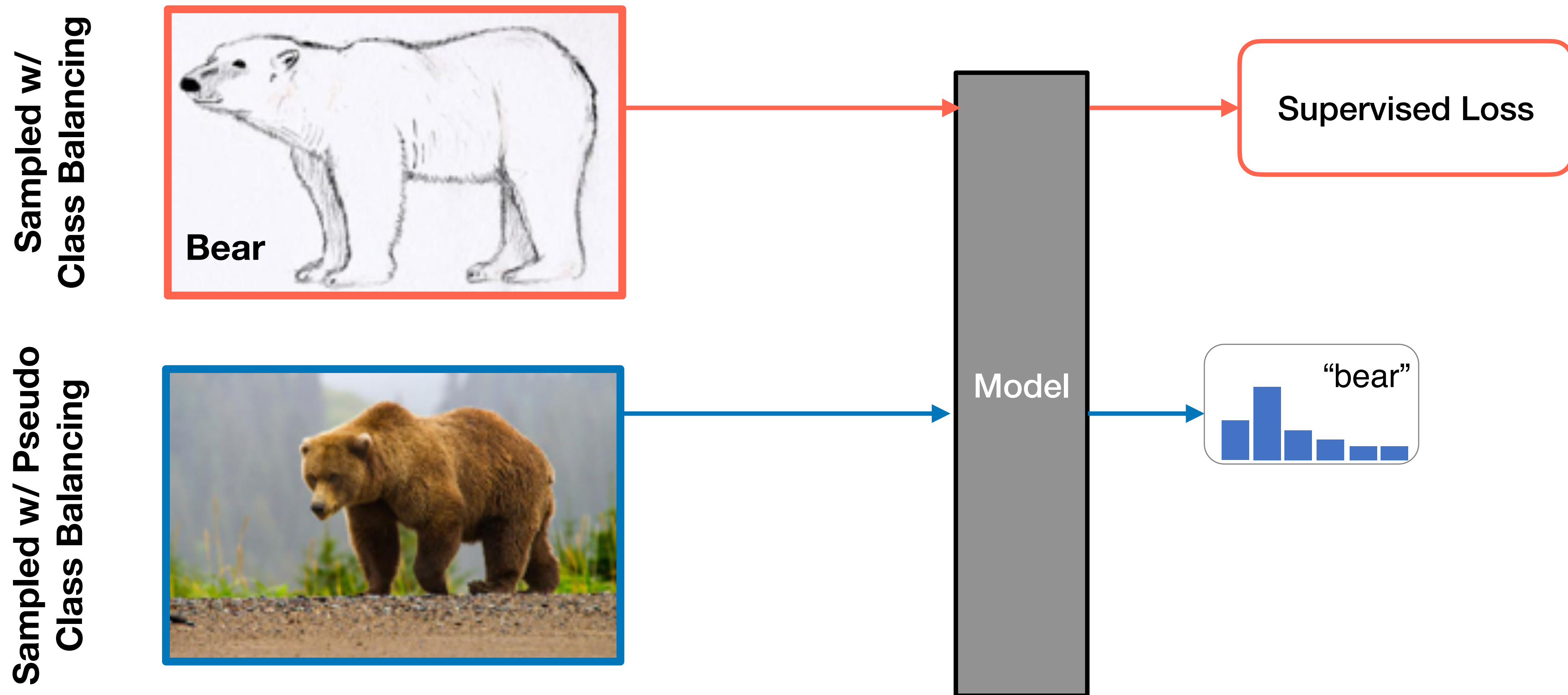


SENTRY: Selective Entropy Optimization via Committee Consistency

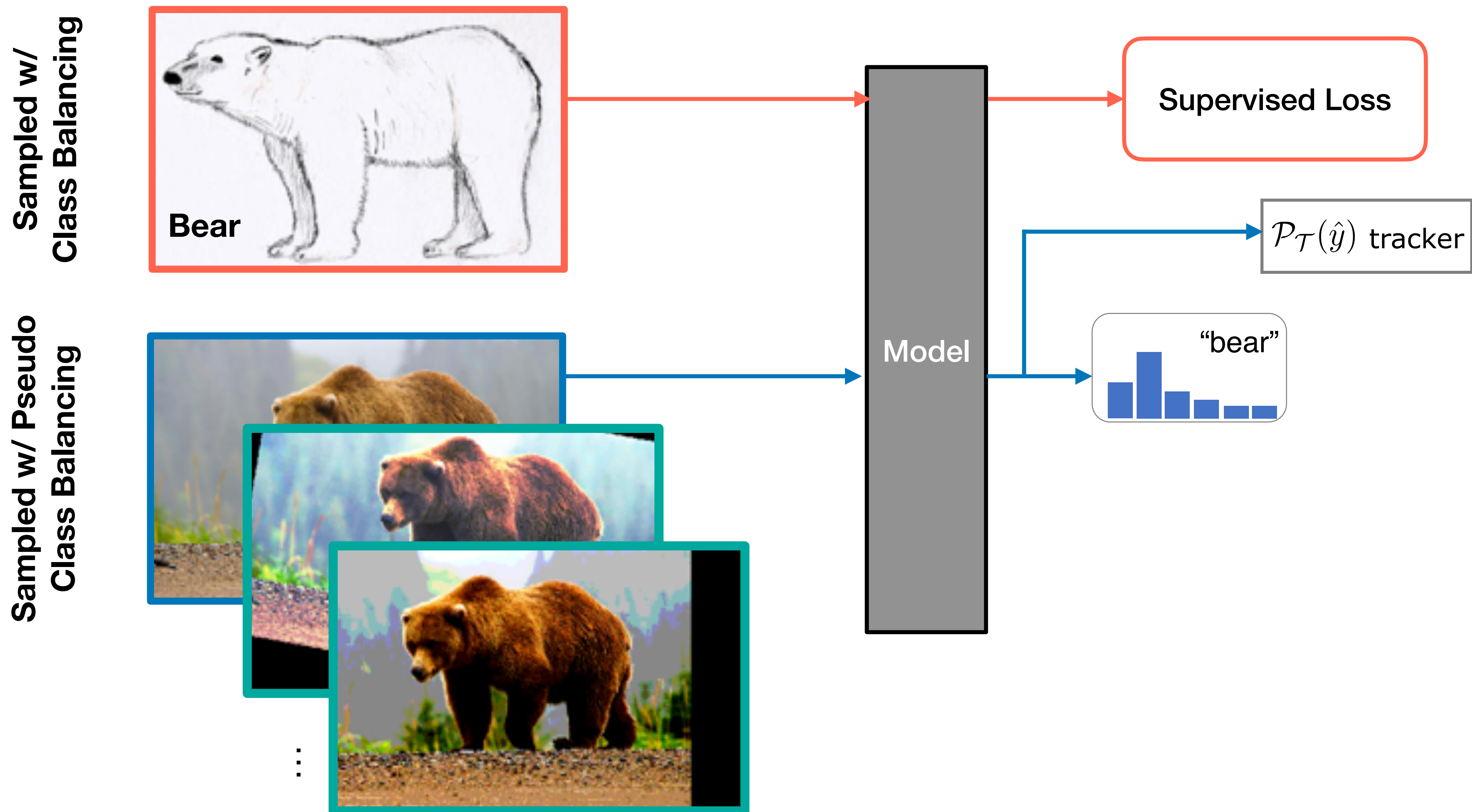
Sampled w/
Class Balancing



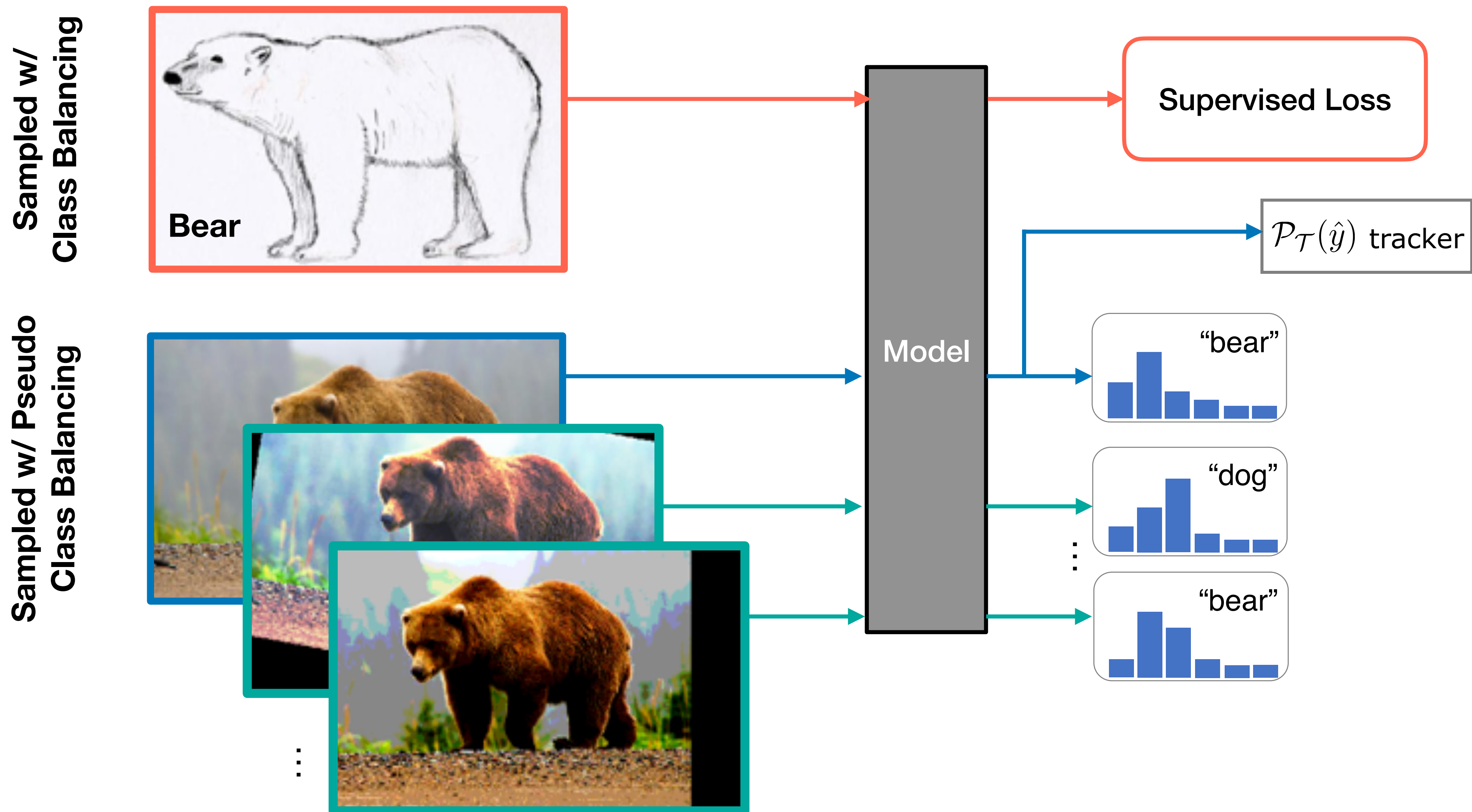
SENTRY: Selective Entropy Optimization via Committee Consistency



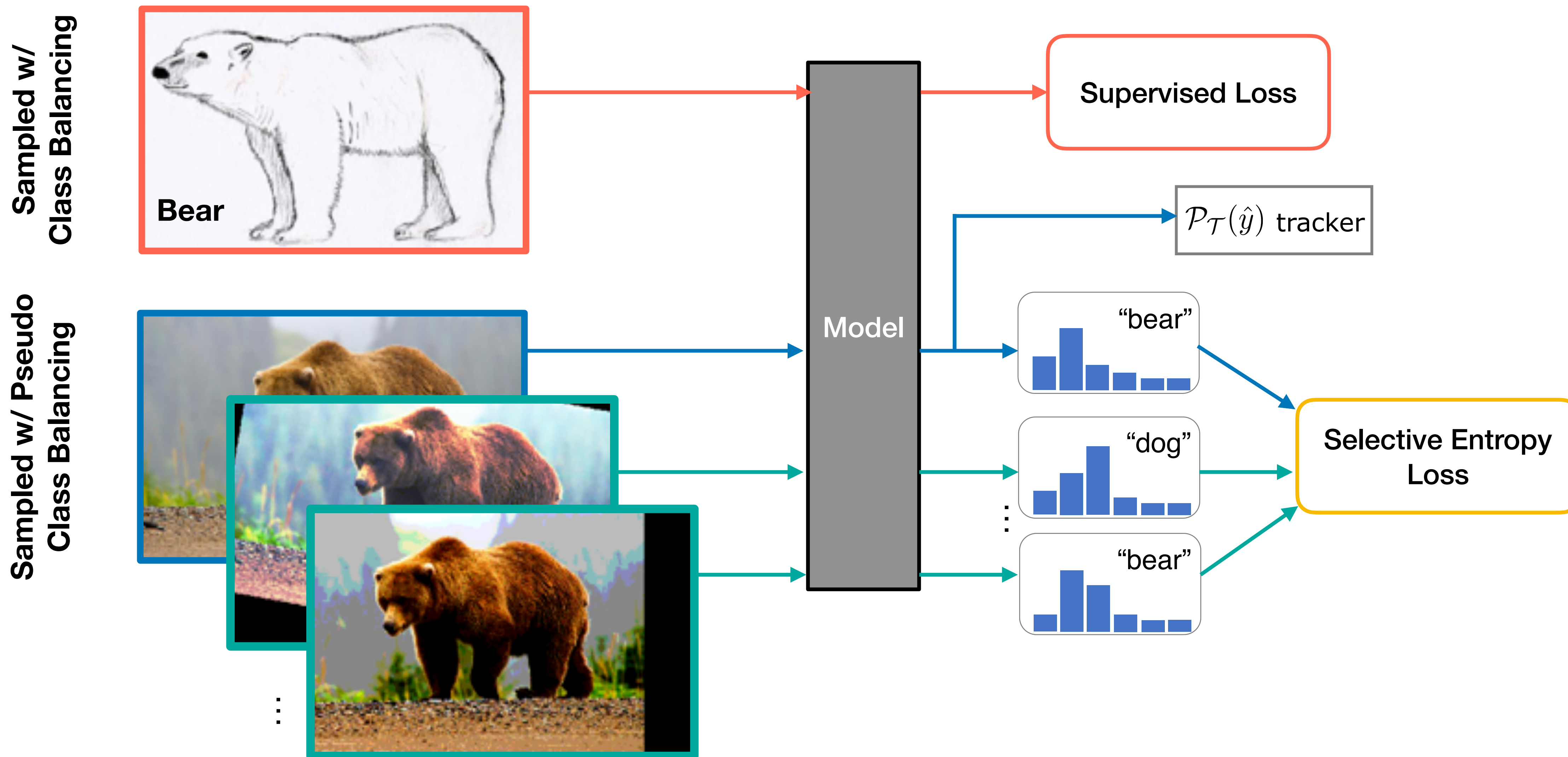
SENTRY: Selective Entropy Optimization via Committee Consistency



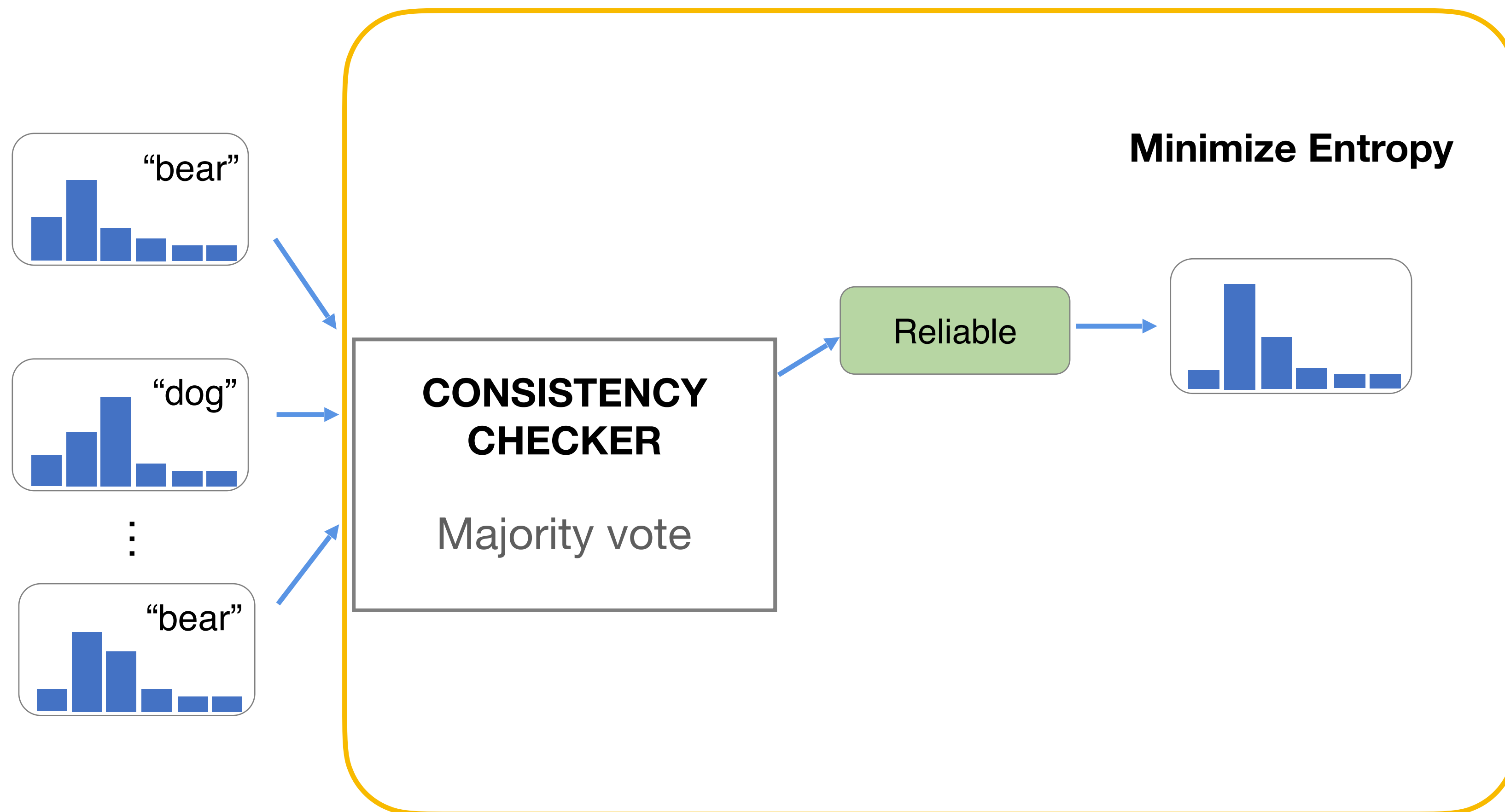
SENTRY: Selective Entropy Optimization via Committee Consistency



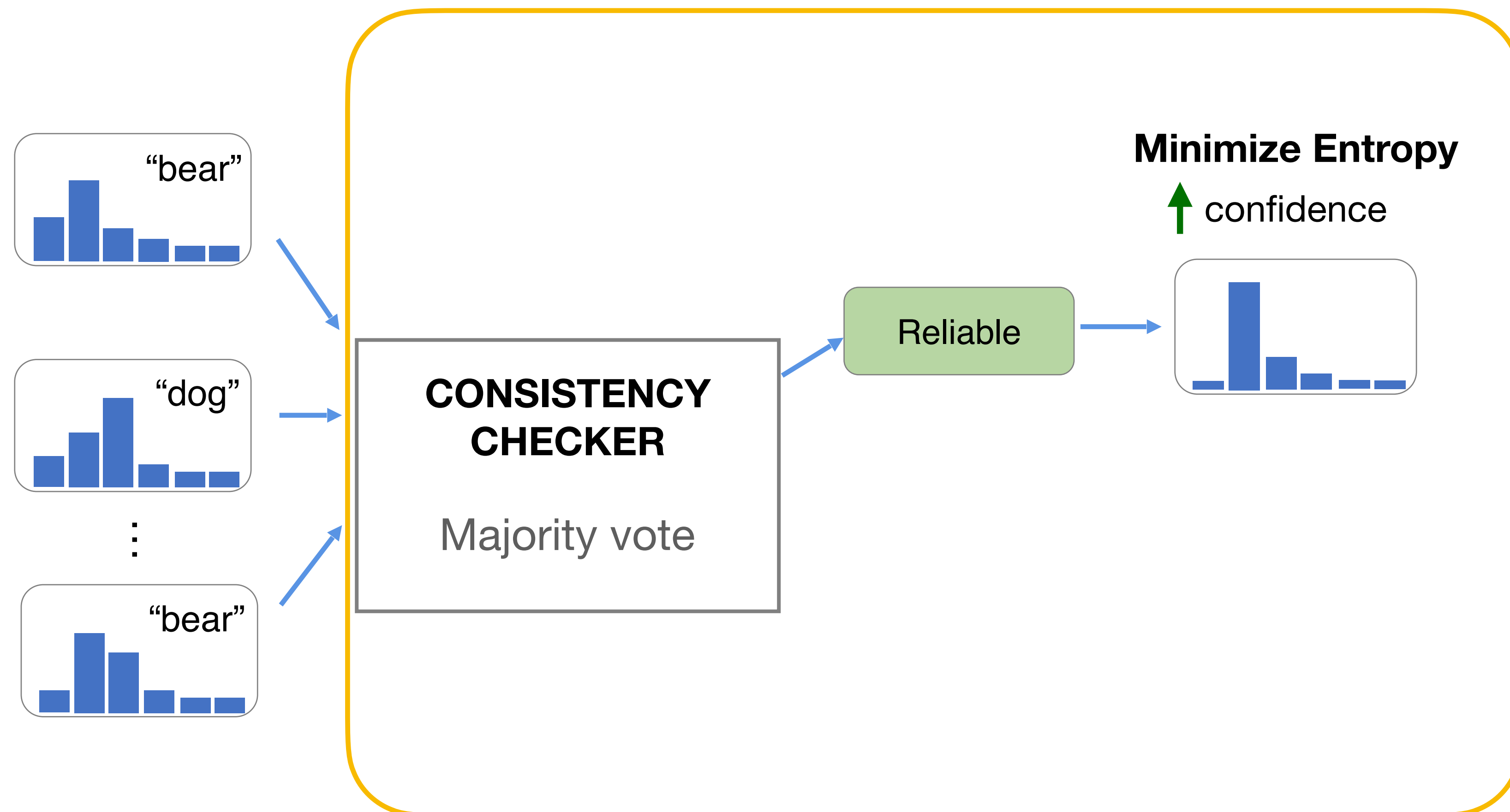
SENTRY: Selective Entropy Optimization via Committee Consistency



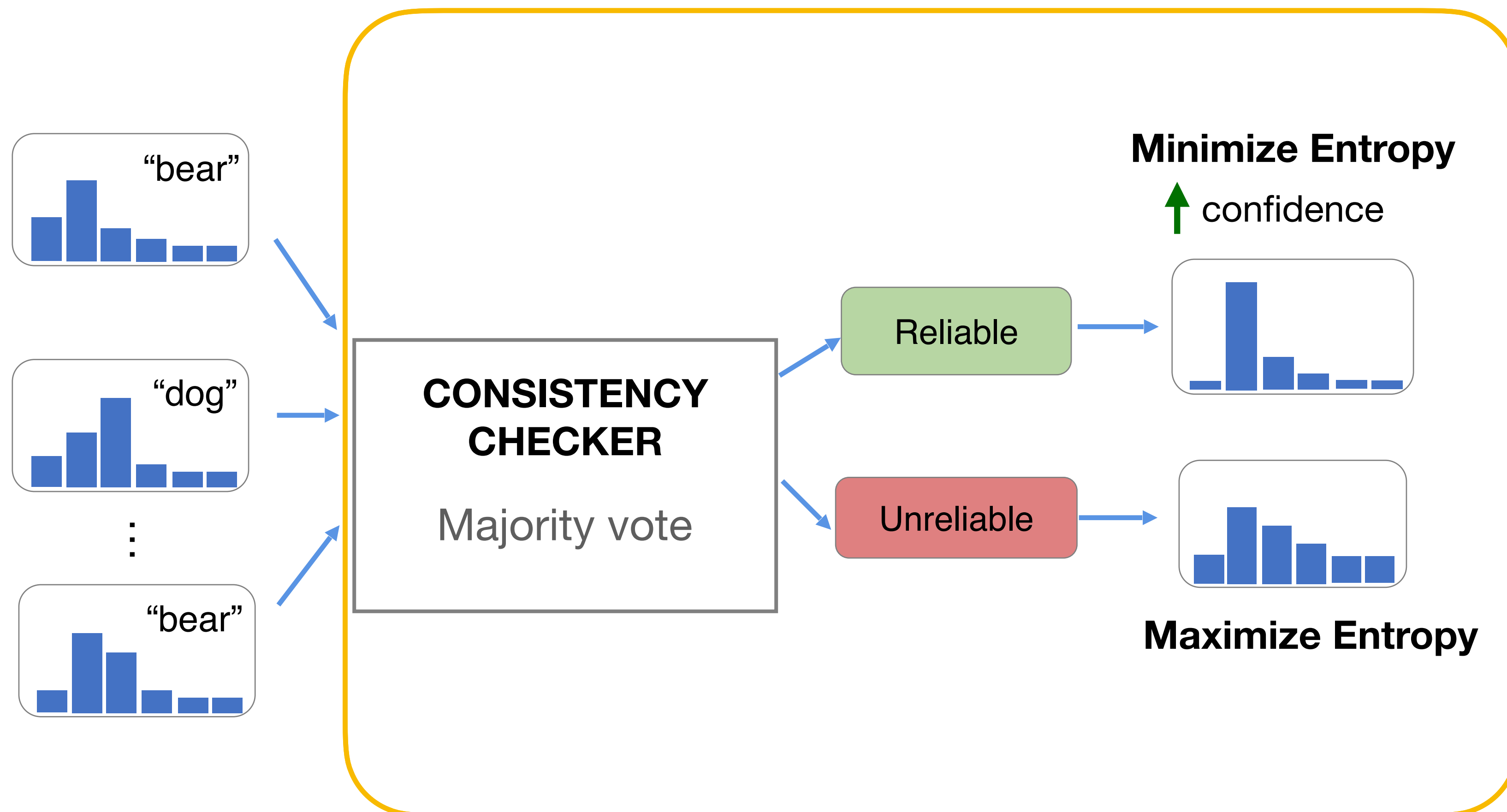
Selective Entropy Loss



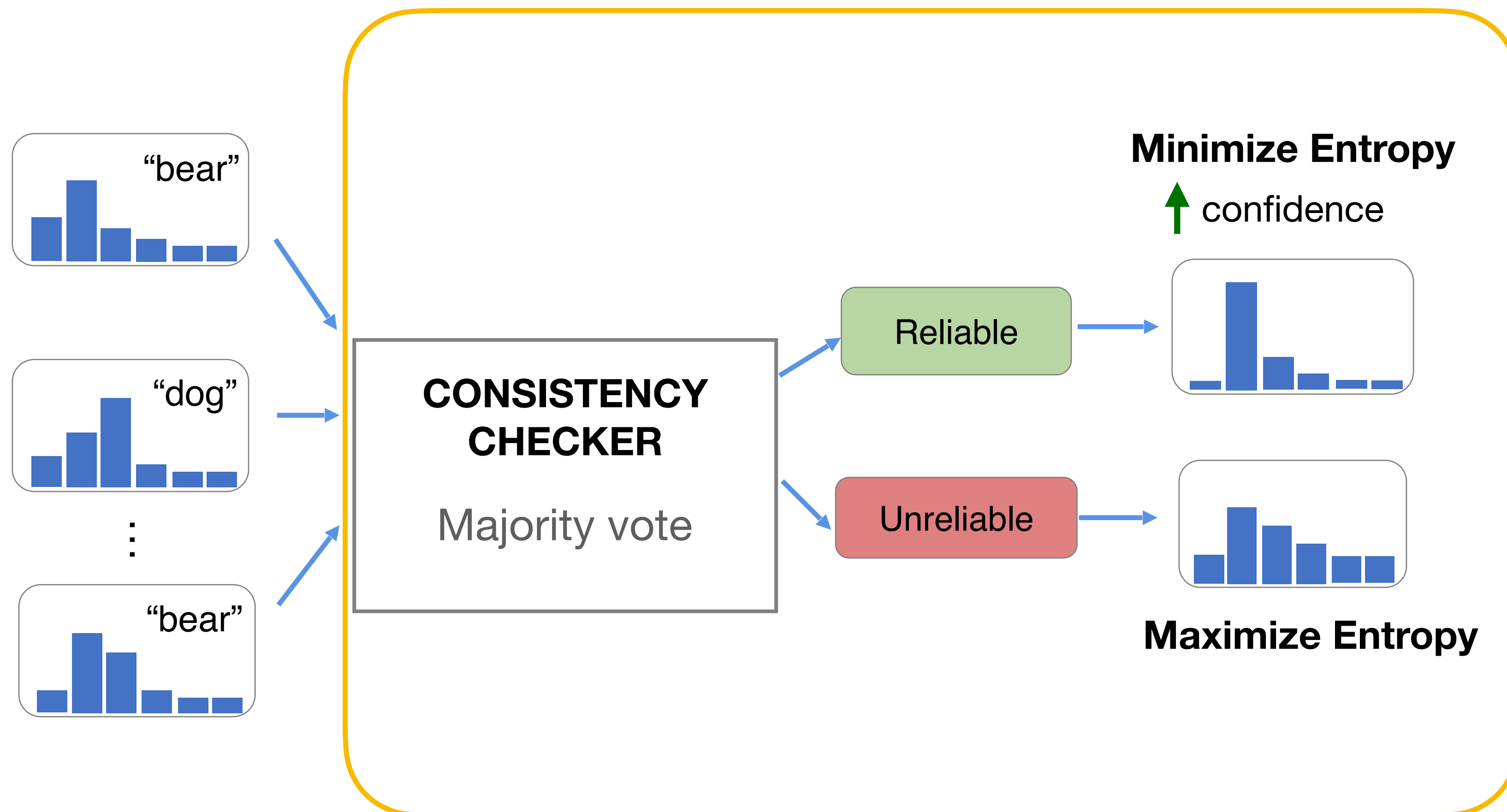
Selective Entropy Loss



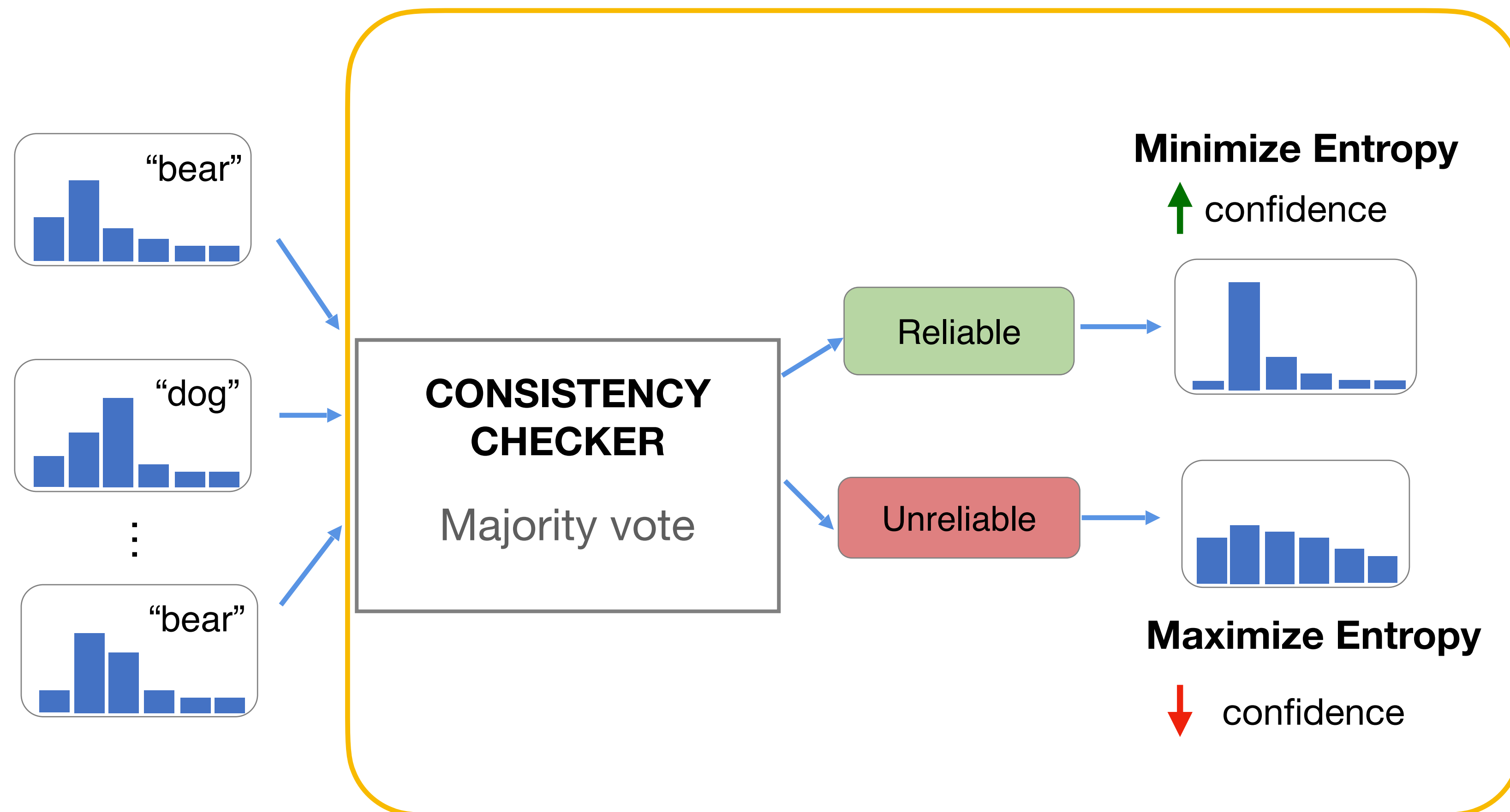
Selective Entropy Loss



Selective Entropy Loss



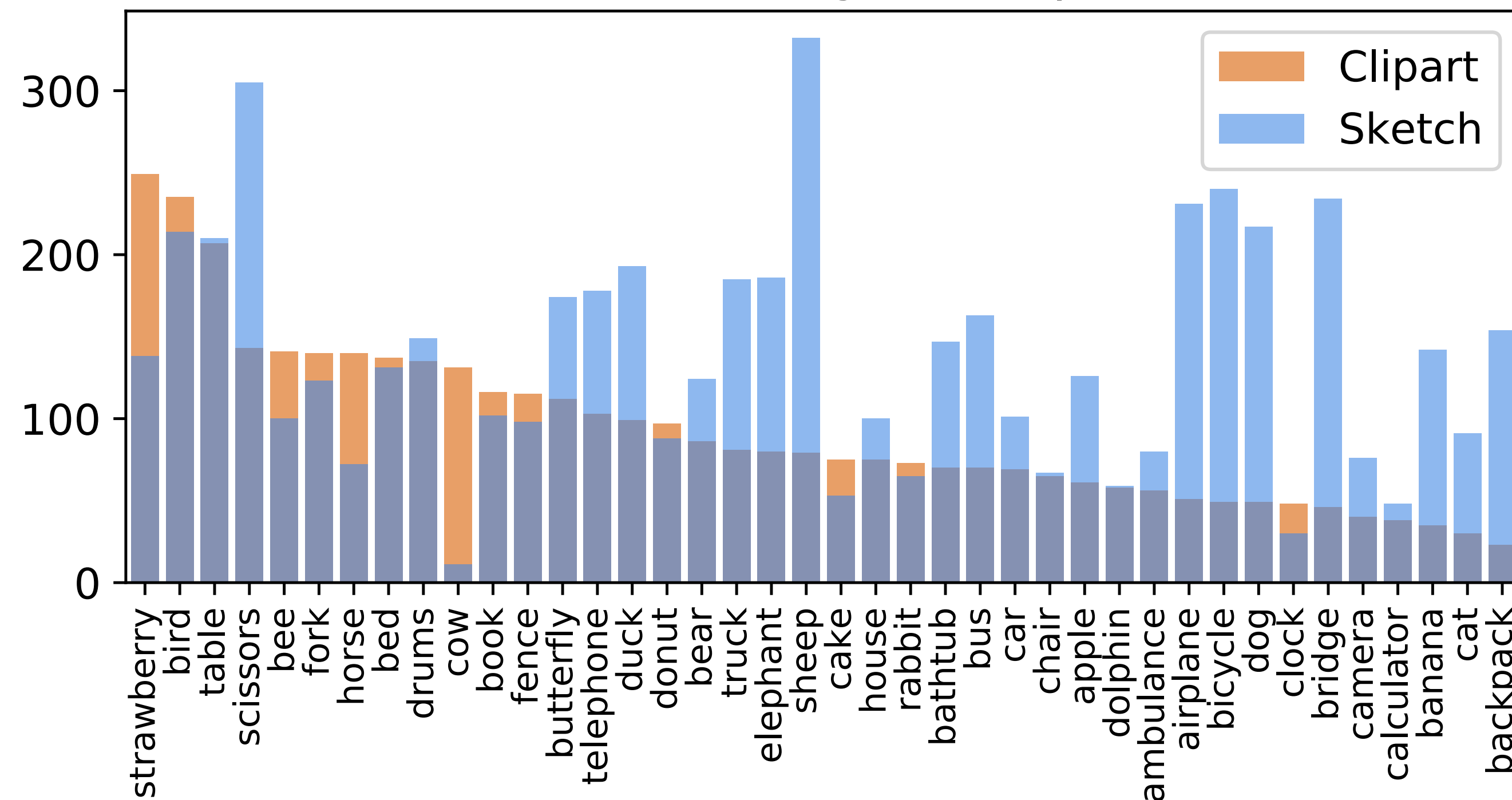
Selective Entropy Loss



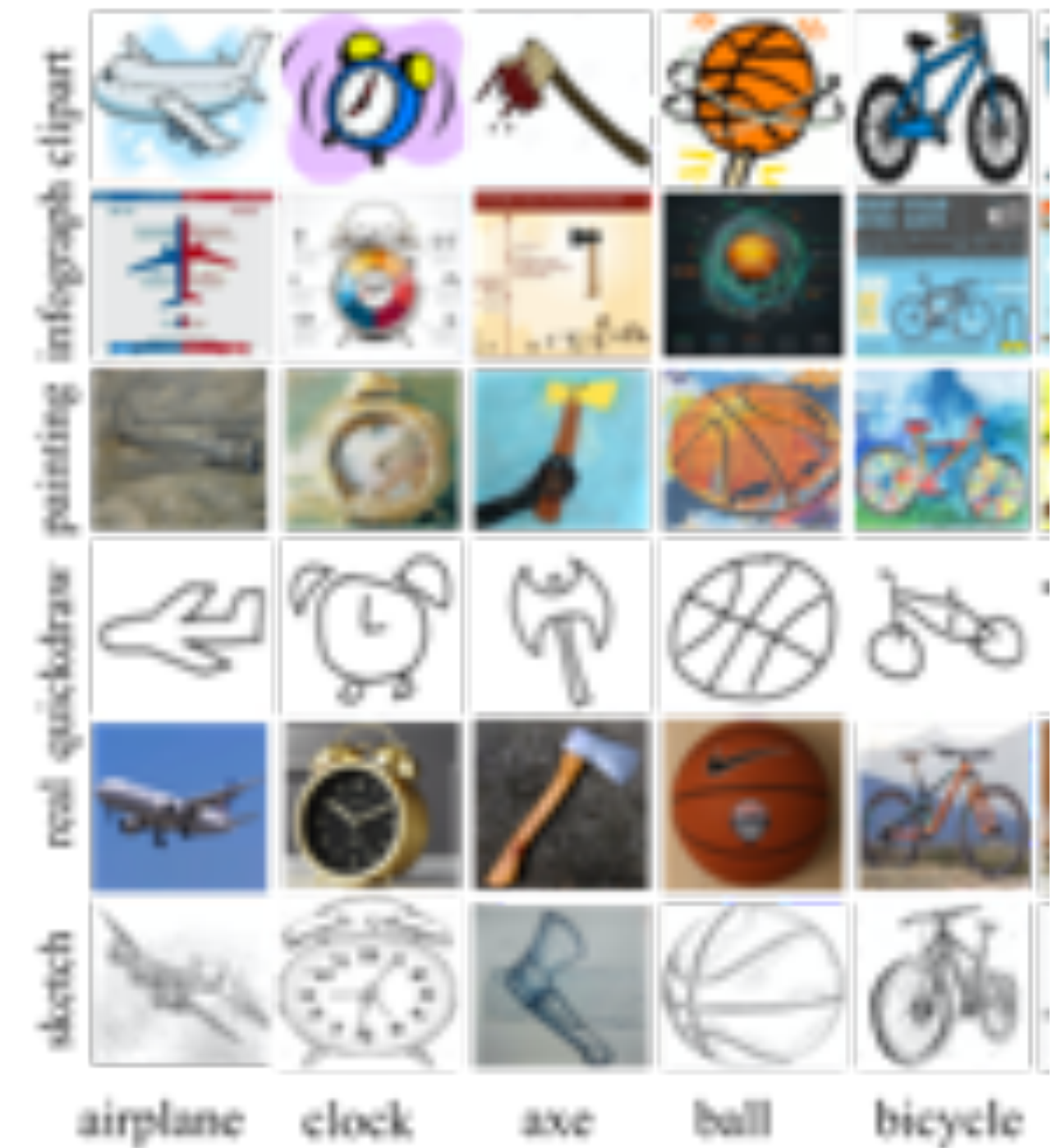
SENTRY Results: Image Classification

Natural label shifts

DomainNet Label Histogram: clipart to sketch

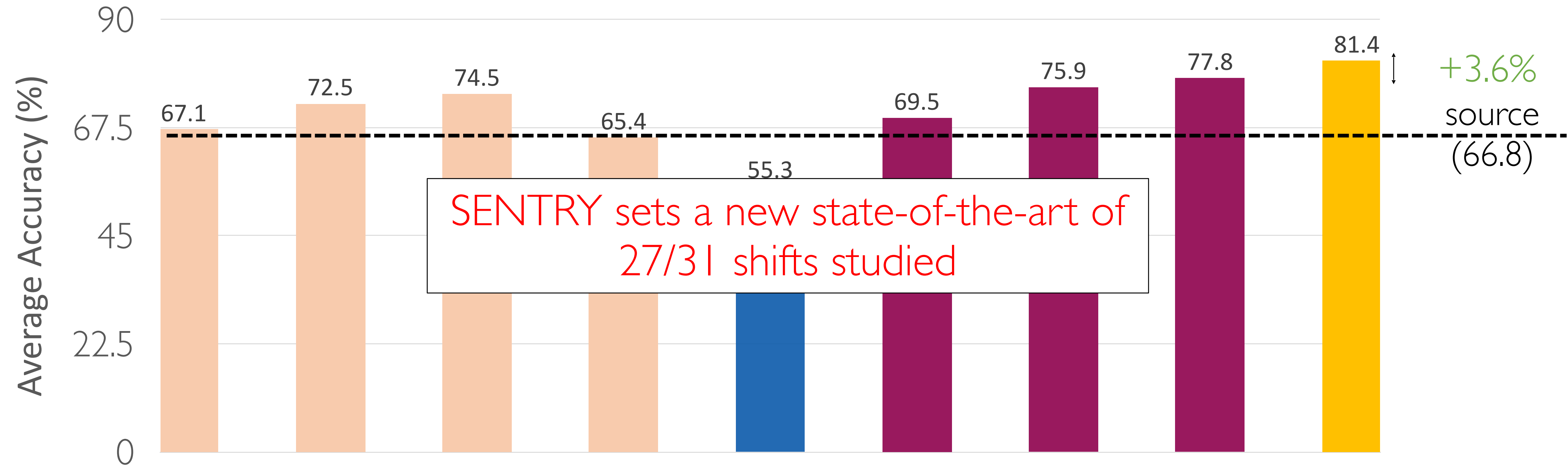


MiniDomainNet^{1,2}



SENTRY Results: MiniDomainNet

MiniDomainNet (40 classes, 12 shifts)



SENTRY sets a new state-of-the-art of 27/31 shifts studied

+3.6% source (66.8)

DAN JAN DANN MCD BBSE F-DANN COAL InstaPBM SENTRY

Distribution-matching based

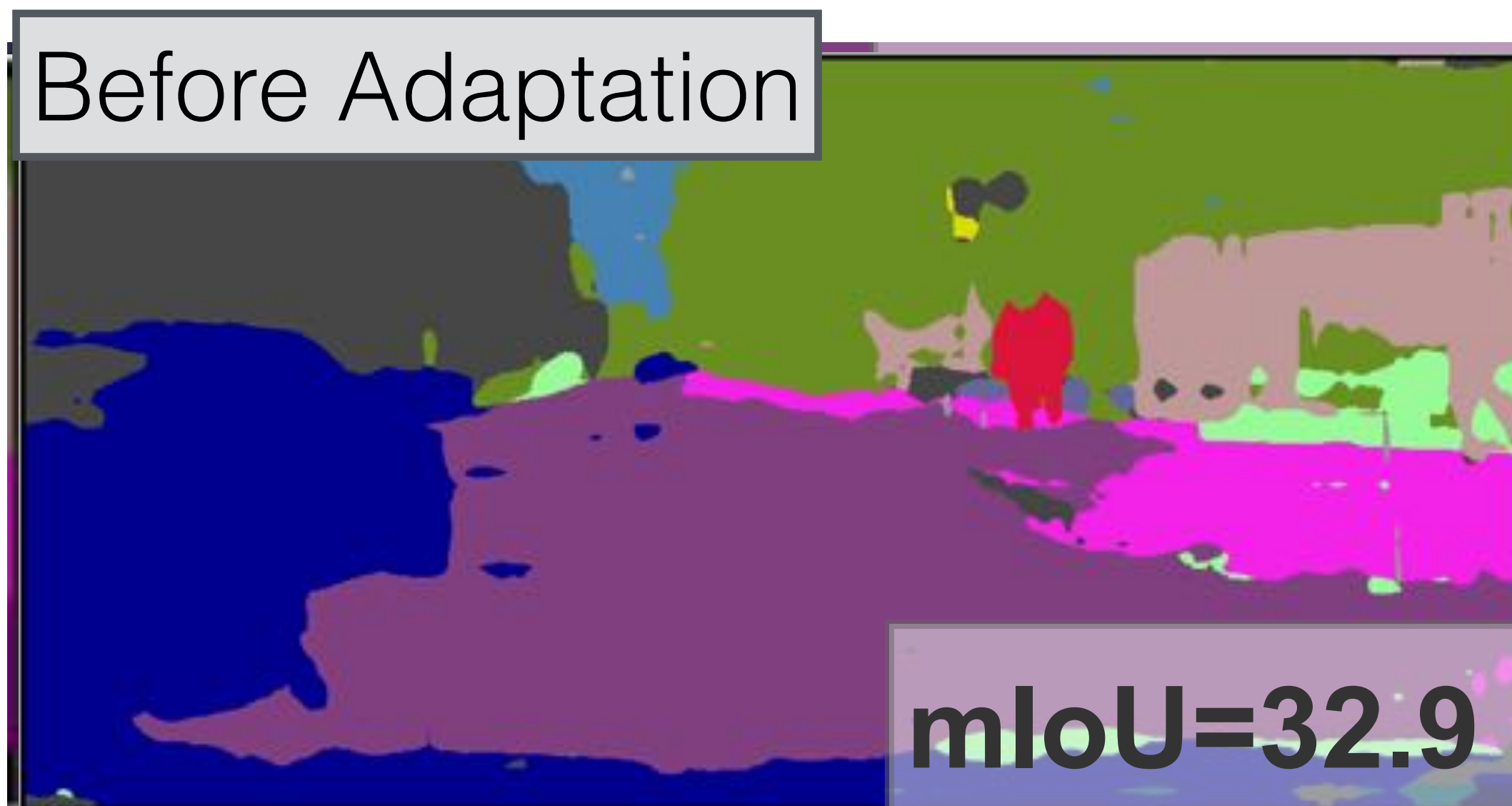
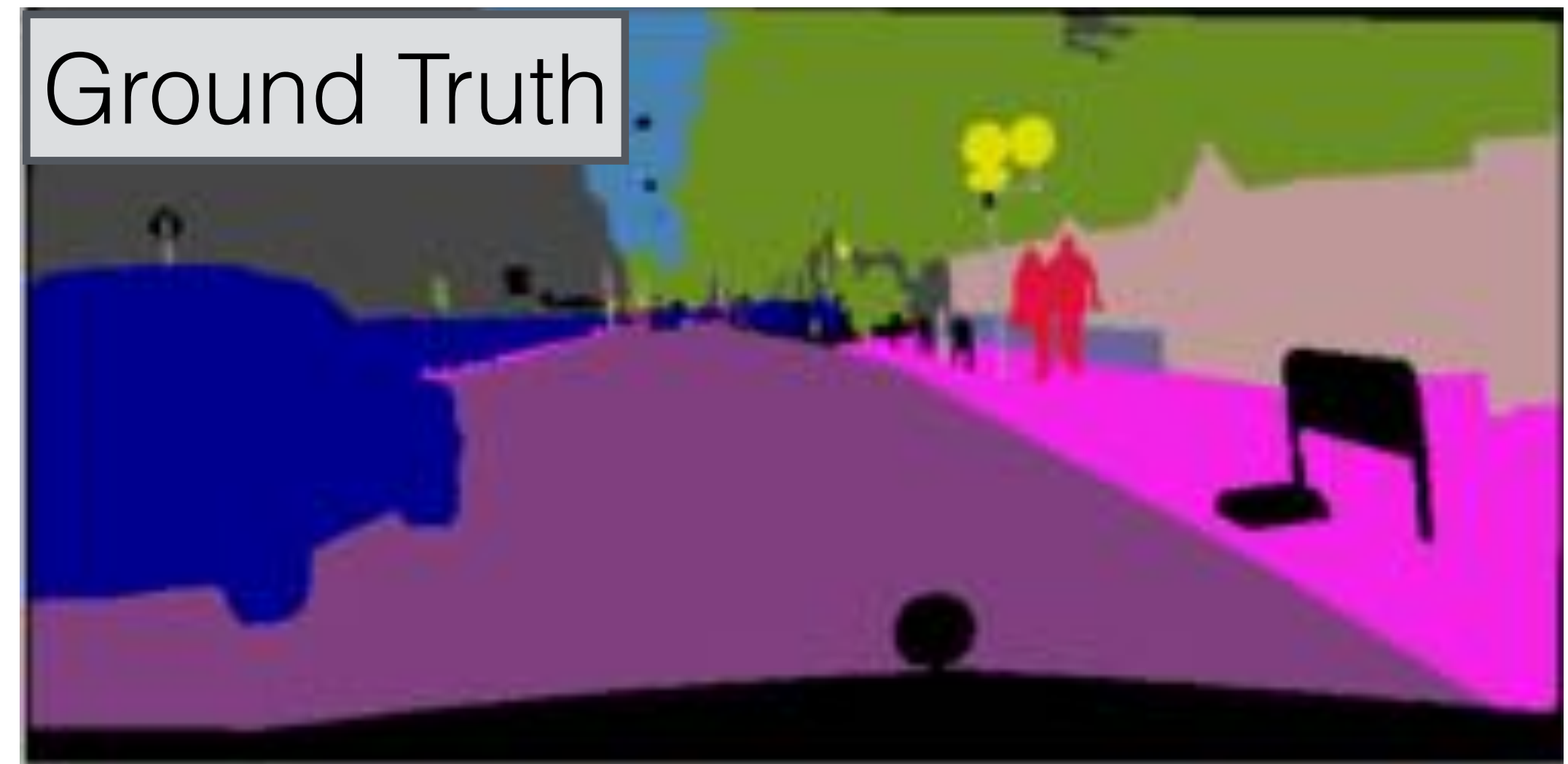
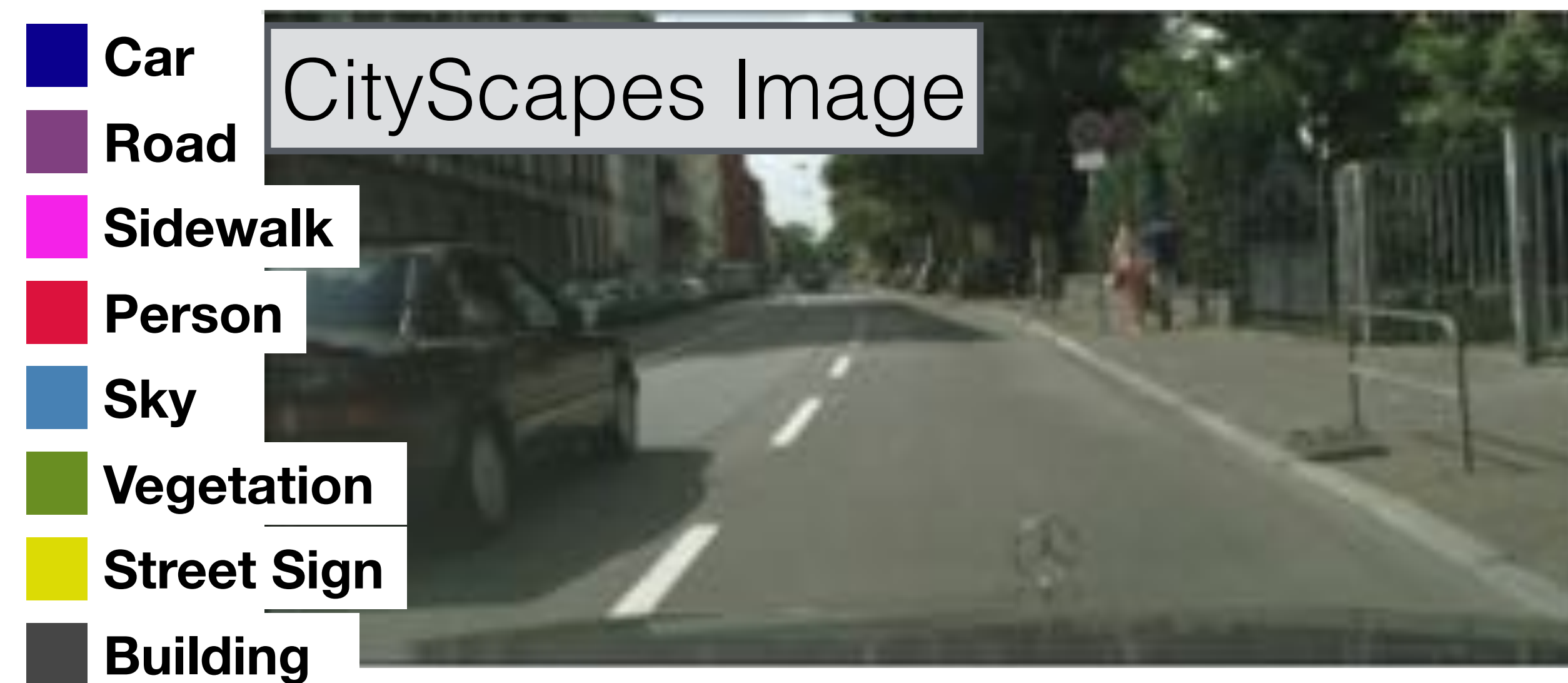
Label shift

"relax self-train on c pseudola Entropy minimization + contrastive loss + mixup loss

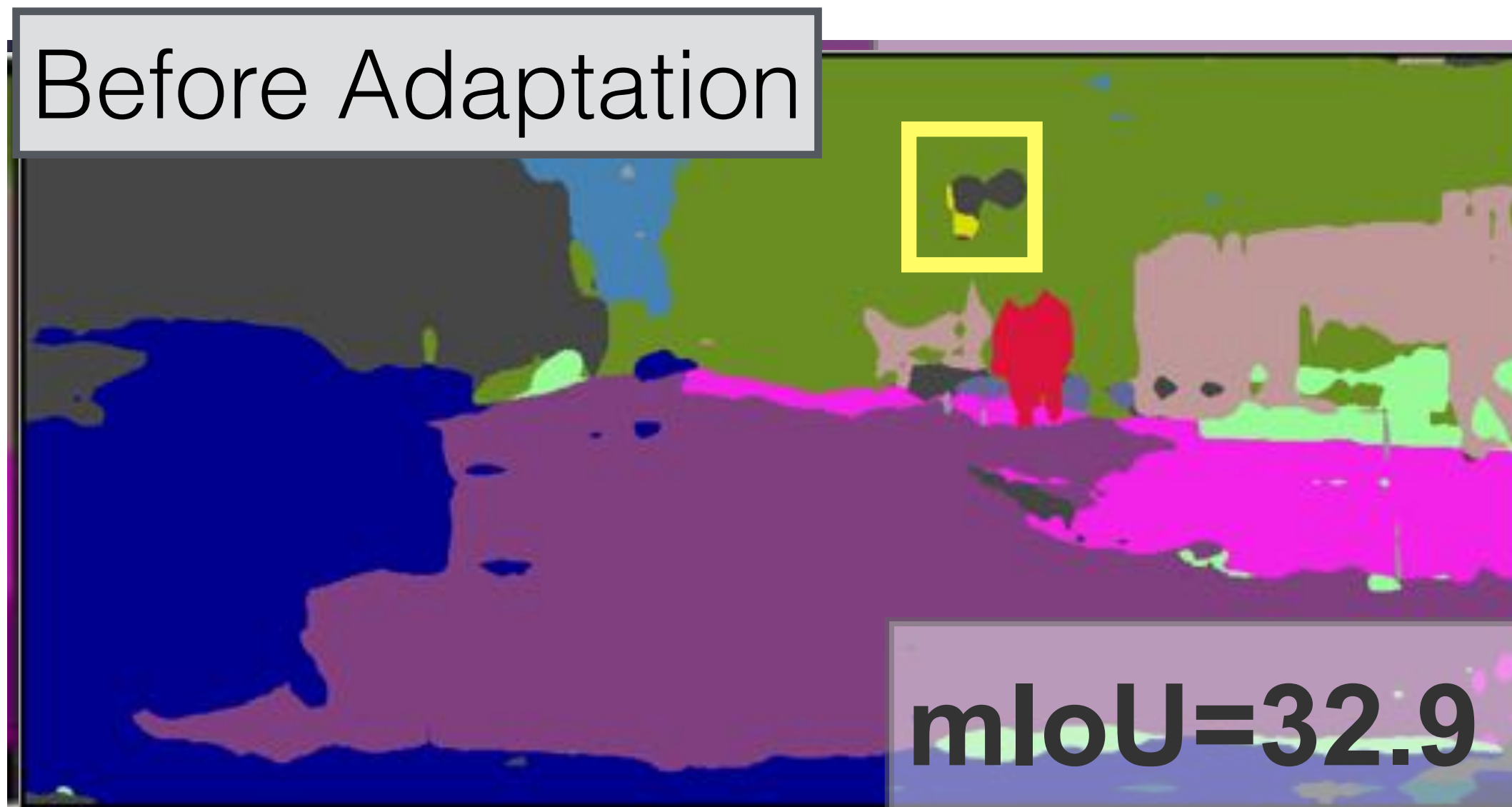
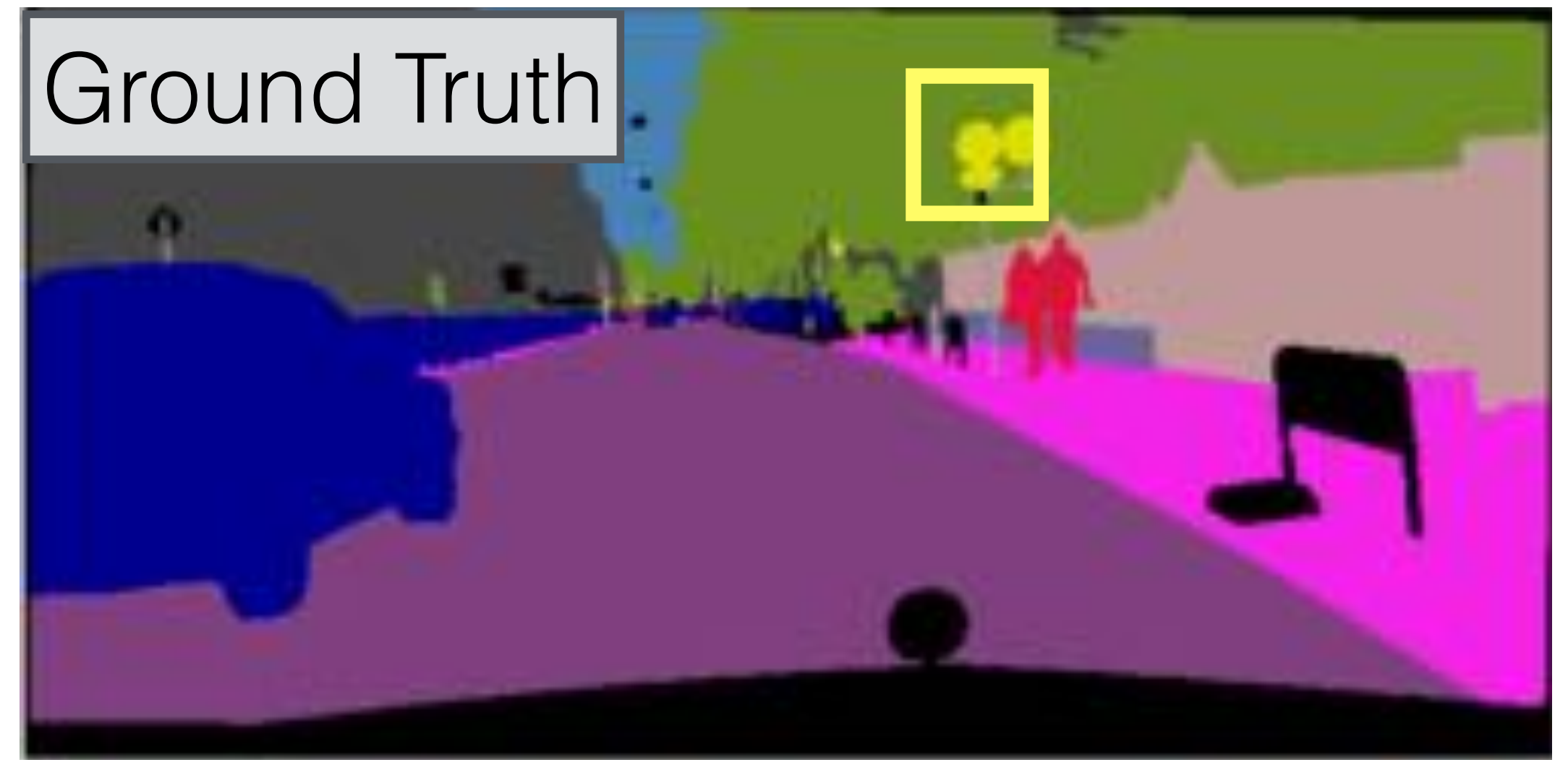
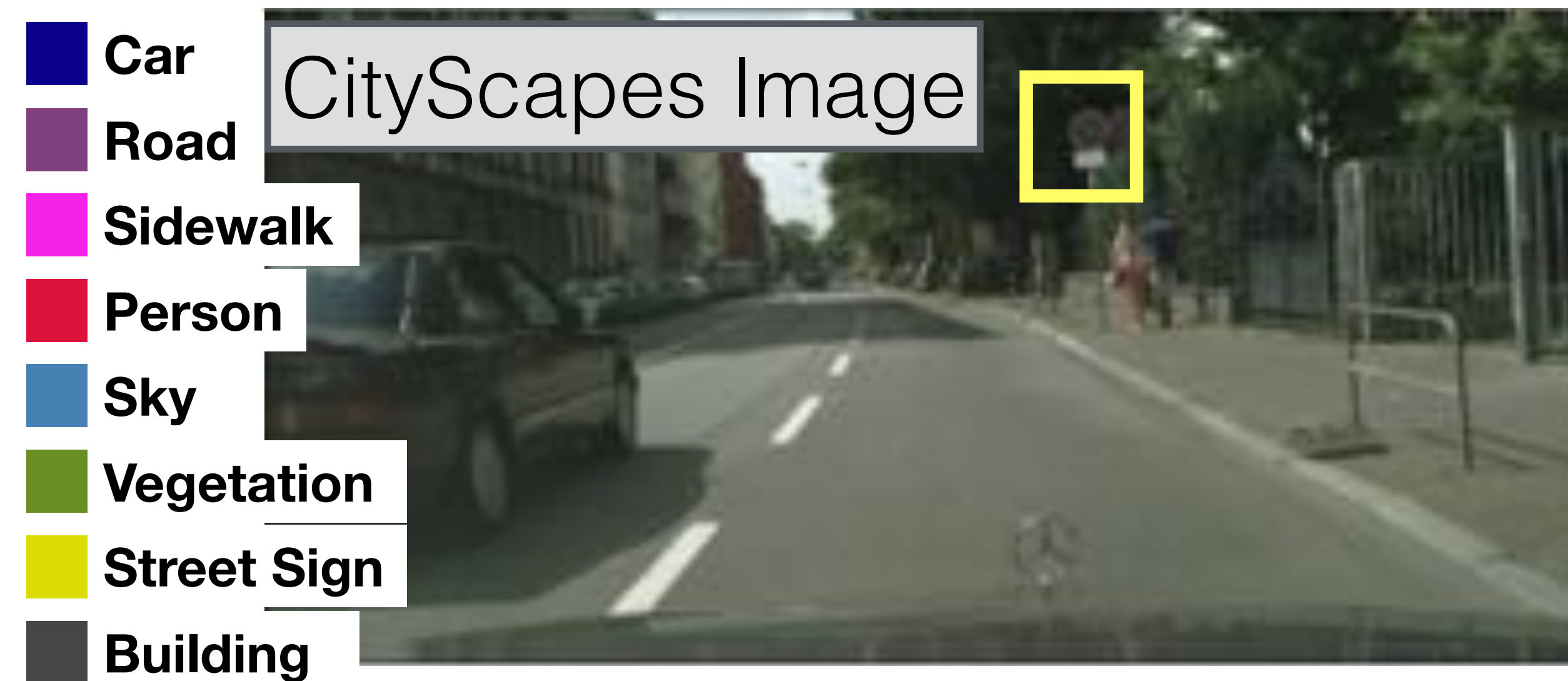
Data + label distribution shift

1. Long et al., ICML 2015.
2. Long et al., ICML 2017.
3. Ganin et al., ICML 2015.
4. Saito et al., CVPR 2018.
5. Lipton et al., ICML 2018.
6. Wu et al., ICML 2019.
7. Tan et al., ECCVW 2020.
8. Li et al., arXiv 2020.

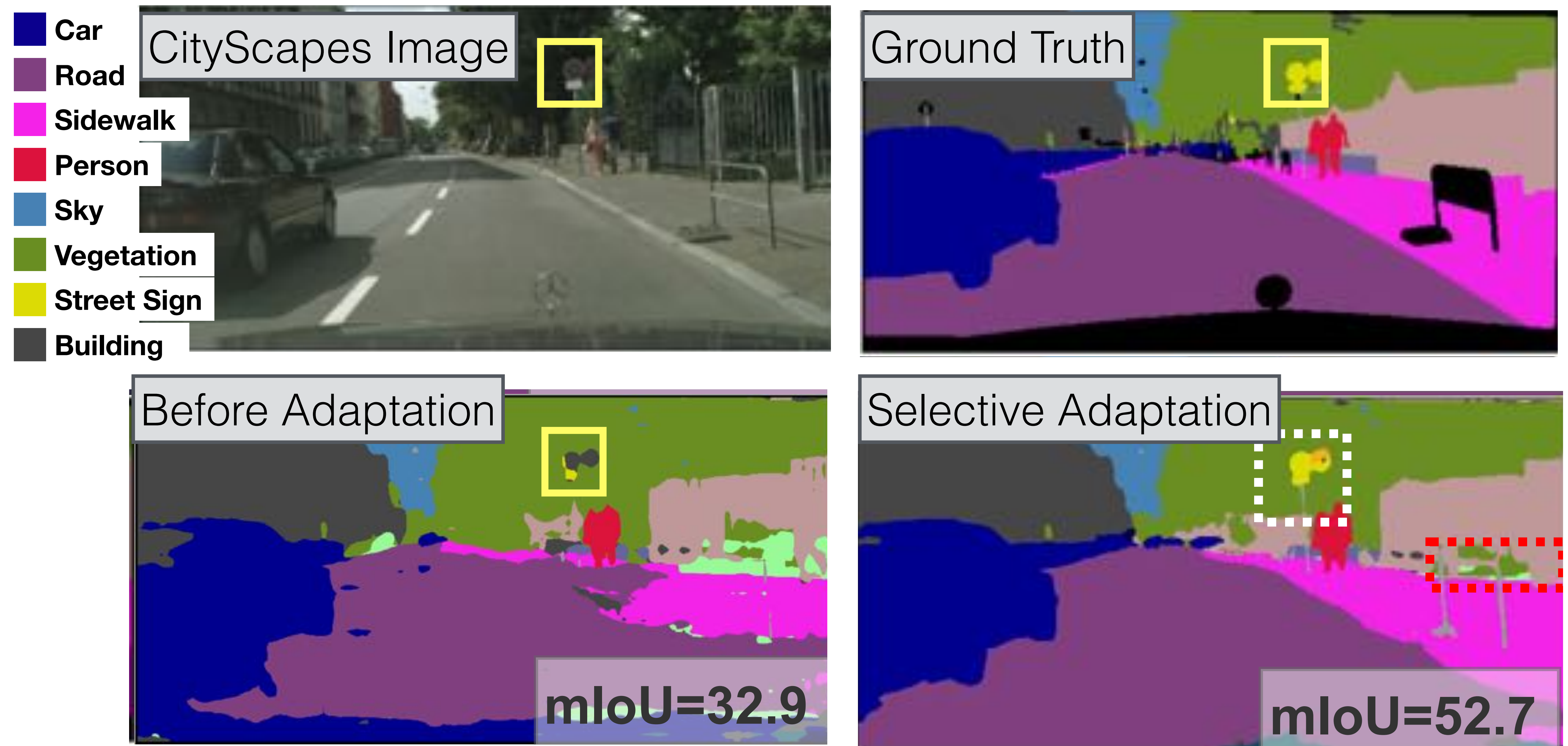
Extension to Semantic Segmentation



Extension to Semantic Segmentation



Extension to Semantic Segmentation



Consistency via attention-conditioned masking

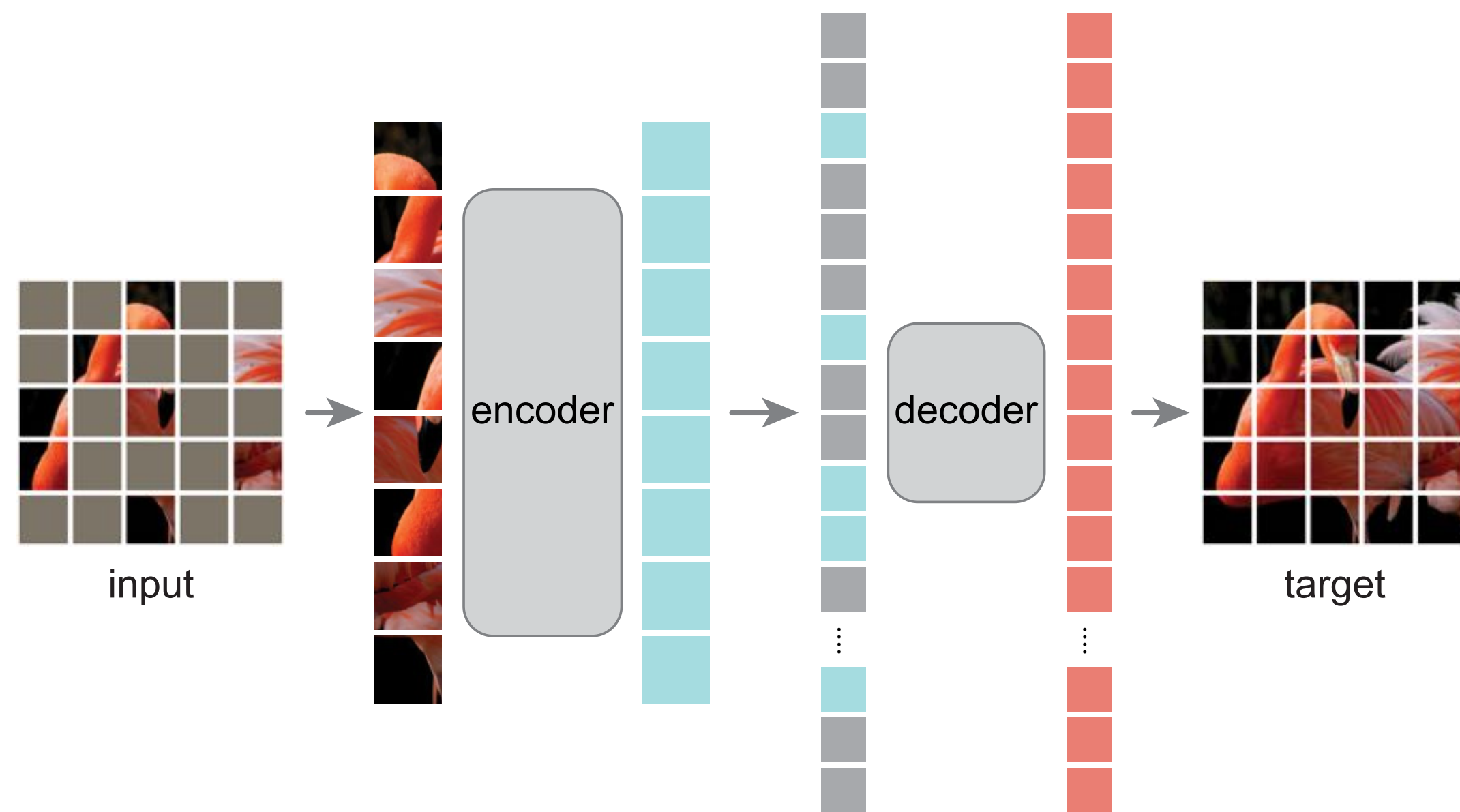
CNN \Rightarrow Vision Transformer (ViT)

Supervised \Rightarrow Self-Supervised Init

Key Idea

Measure predictive consistency under:

Random augmentations \Rightarrow **Self-supervised proxy task**



“marker”



**ATTENTION-
CONDITIONED
MASKING**

“pencil”



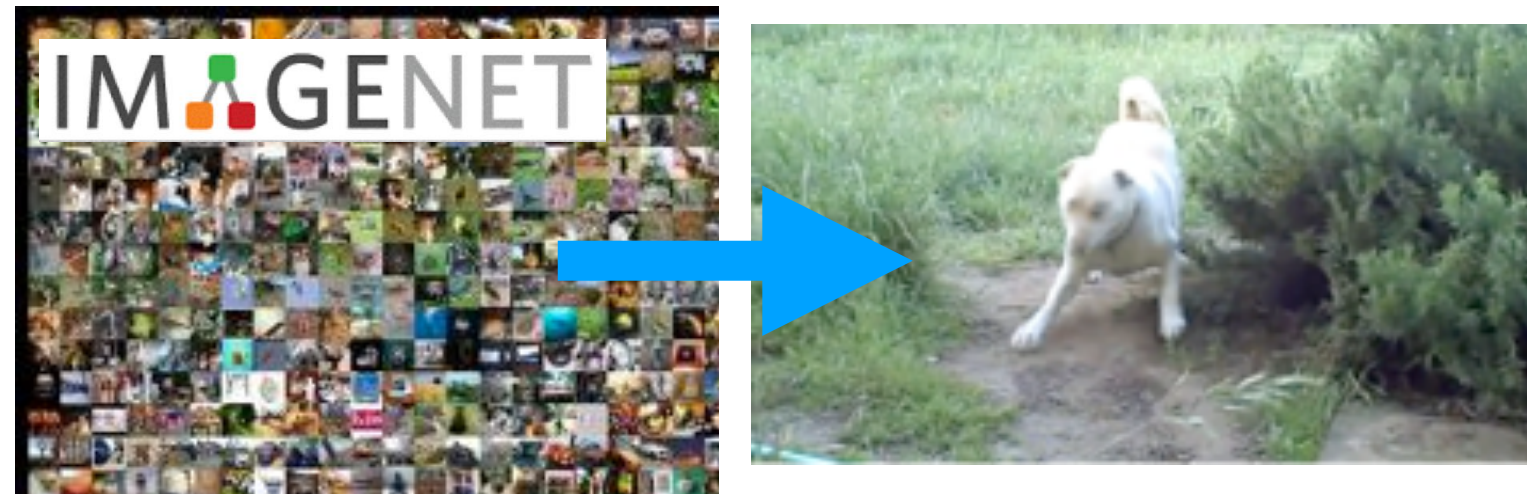
“flowers”



Masked Autoencoders, He *et al.*, CVPR 2022

Performance Degradation from Bias

Curation



Weather



Sensor

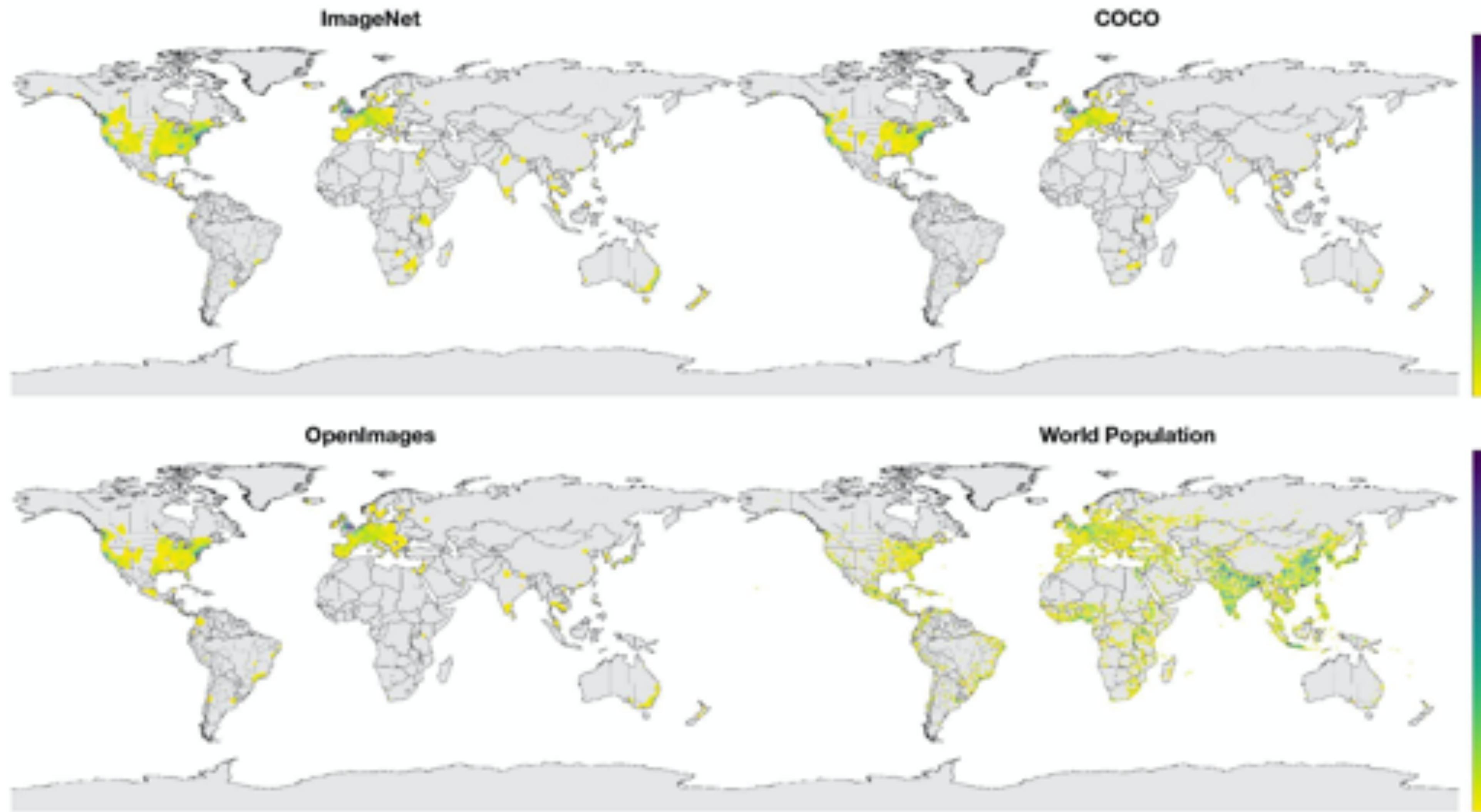


Demographic

Systems can underperform for certain subpopulations

Often caused by underrepresentation

Geographic Bias



Does object recognition work for everyone?



Ground truth: Soap

Nepal, 288 \$/month

Azure: food, cheese, bread, cake, sandwich

Clarifai: food, wood, cooking, delicious, healthy

Google: food, dish, cuisine, comfort food, spam

Amazon: food, confectionary, sweets, burger

Watson: food, food product, turmeric, seasoning

Tencent: food, dish, matter, fast food, nutriment



Ground truth: Soap

UK, 1890 \$/month

Azure: toilet, design, art, sink

Clarifai: people, faucet, healthcare, lavatory, wash closet

Google: product, liquid, water, fluid, bathroom accessory

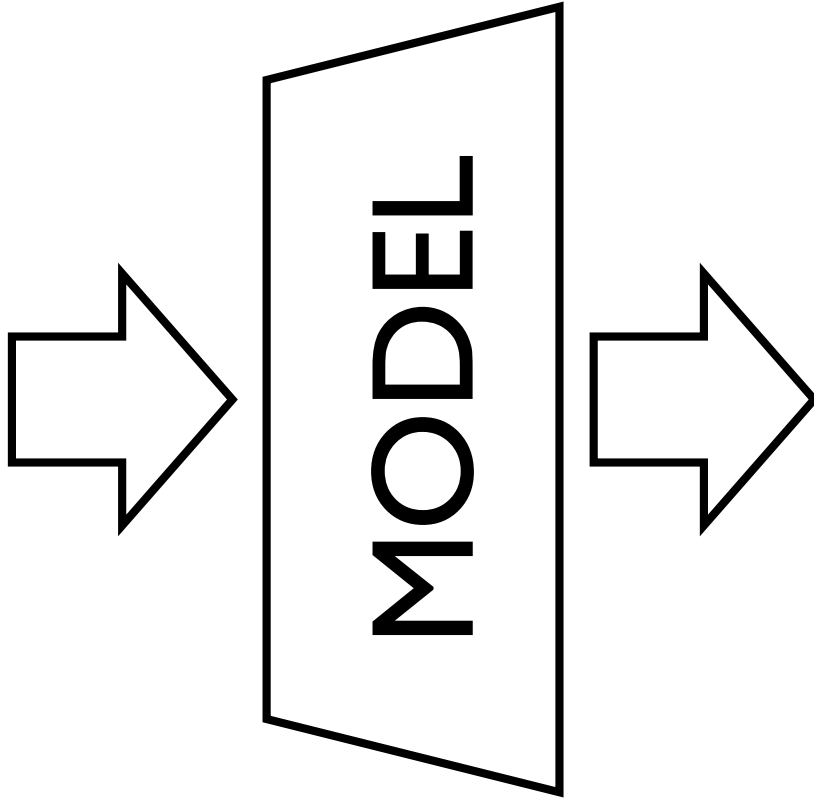
Amazon: sink, indoors, bottle, sink faucet

Watson: gas tank, storage tank, toiletry, dispenser, soap dispenser

Tencent: lotion, toiletry, soap dispenser, dispenser, after shave

Can domain adaptation make obj rec work for everyone?

Train (North America)
label = "statue"



Test (Rest of the world)



Geographically diverse data



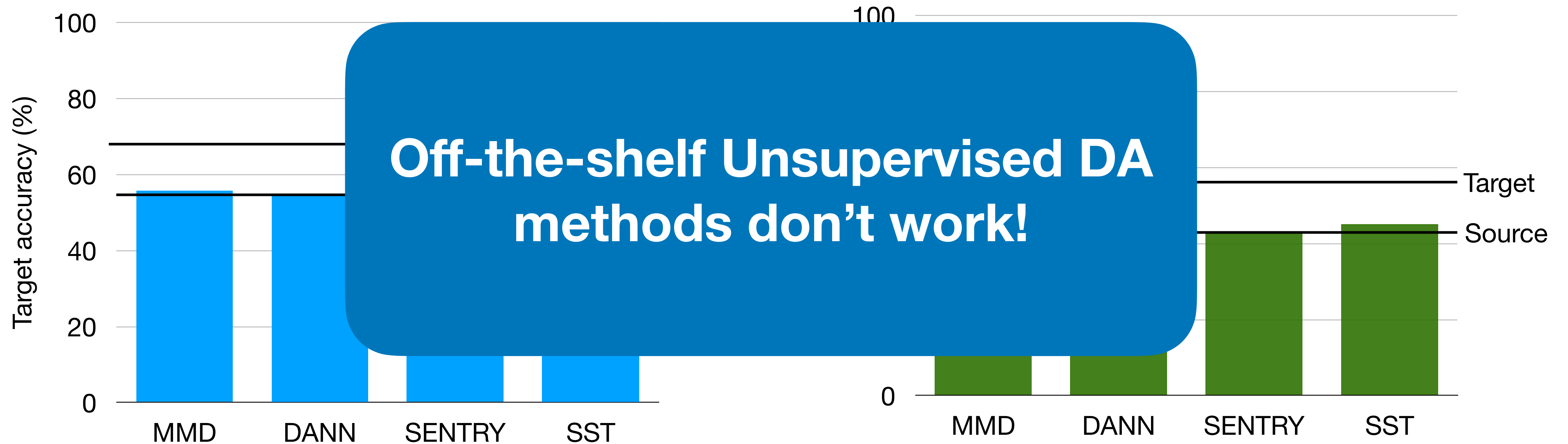
Results

Dollarstreet-DA

{N. America, Europe} → {Asia, Africa, S. America}

GeoYFCC-DA

{N. America} → {Asia, Australia, S. America}



1. Long *et al.*, ICML 2015
2. Ganin *et al.*, ICML 2015
3. Prabhu *et al.*, ICCV 2021

Additional challenges in GeoDA

Context Shift

$$P_S(c(\mathbf{x}) | y) \neq P_T(c(\mathbf{x}) | y)$$

Specialized solutions are needed for Geo DA!

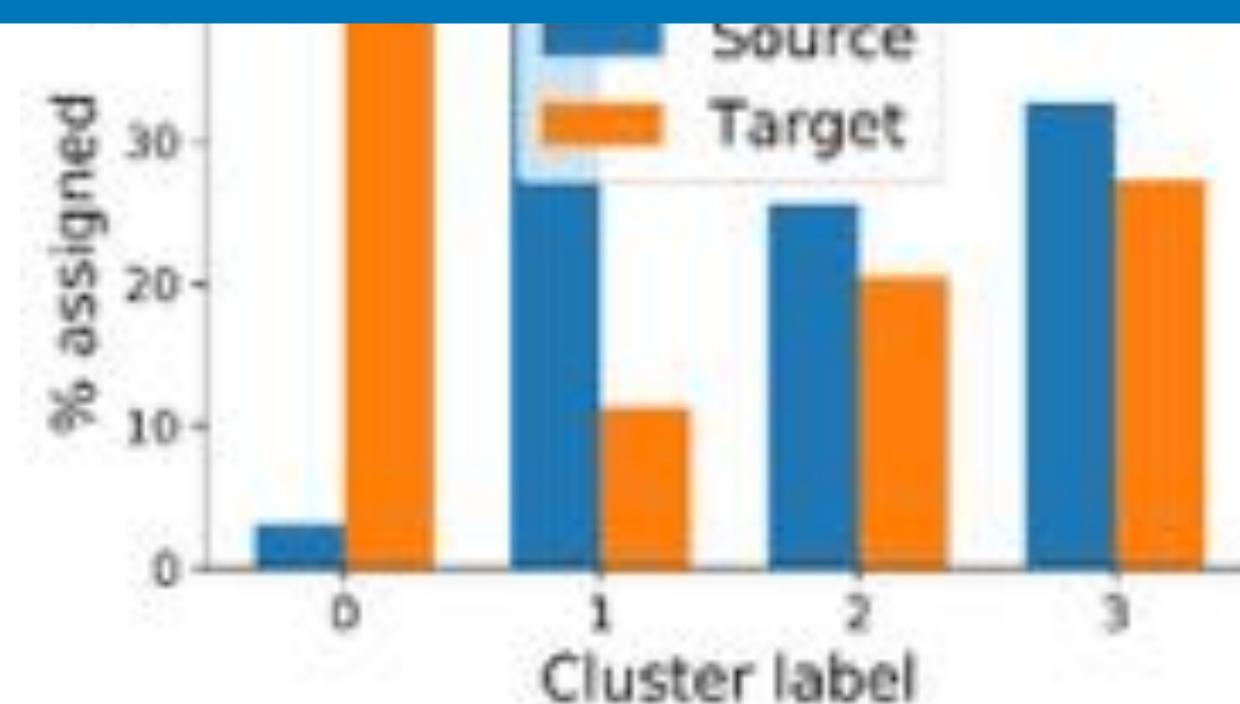
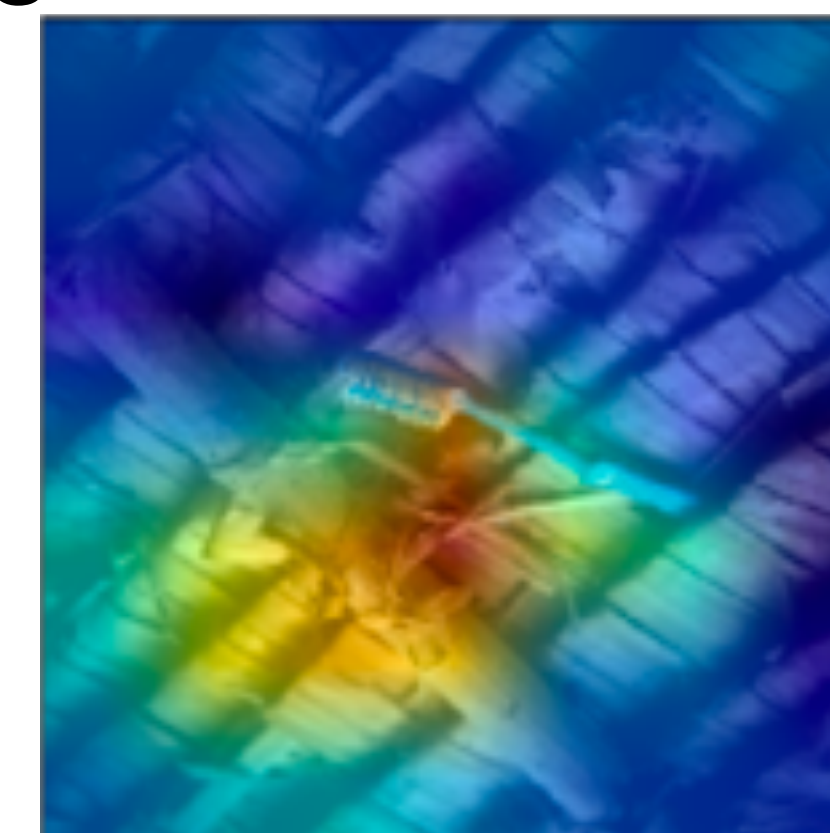
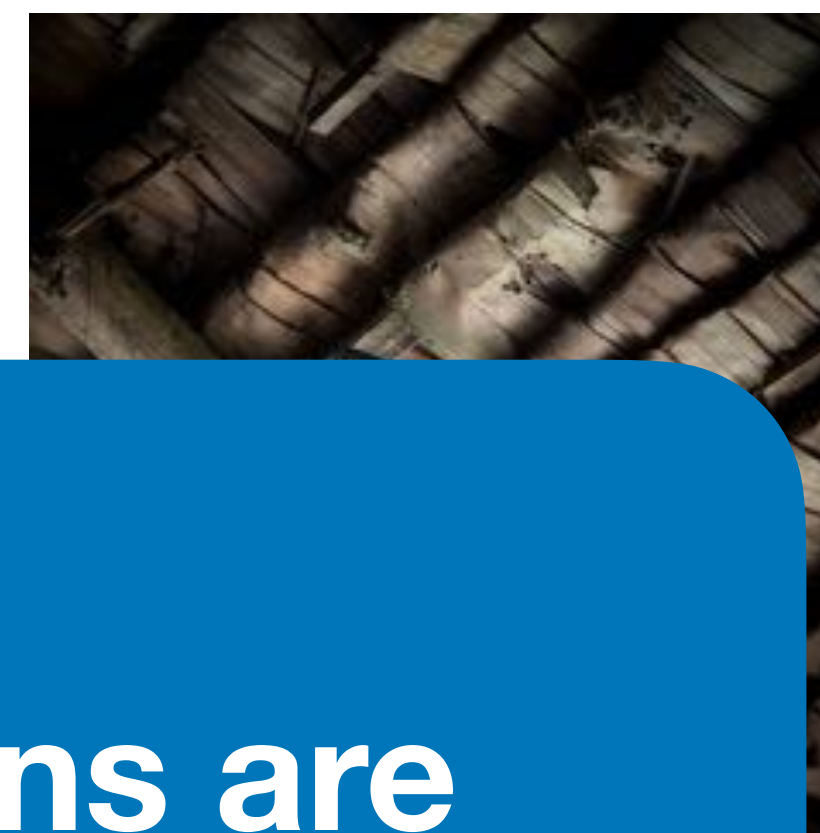
Subpopulation Shift

$$P_S(\mathbf{x} | y) \neq P_T(\mathbf{x} | y)$$

source



target



Summary: Responsible Vision

Reliability Goal: Perform vision tasks as expected at deployment time.

Benchmarks

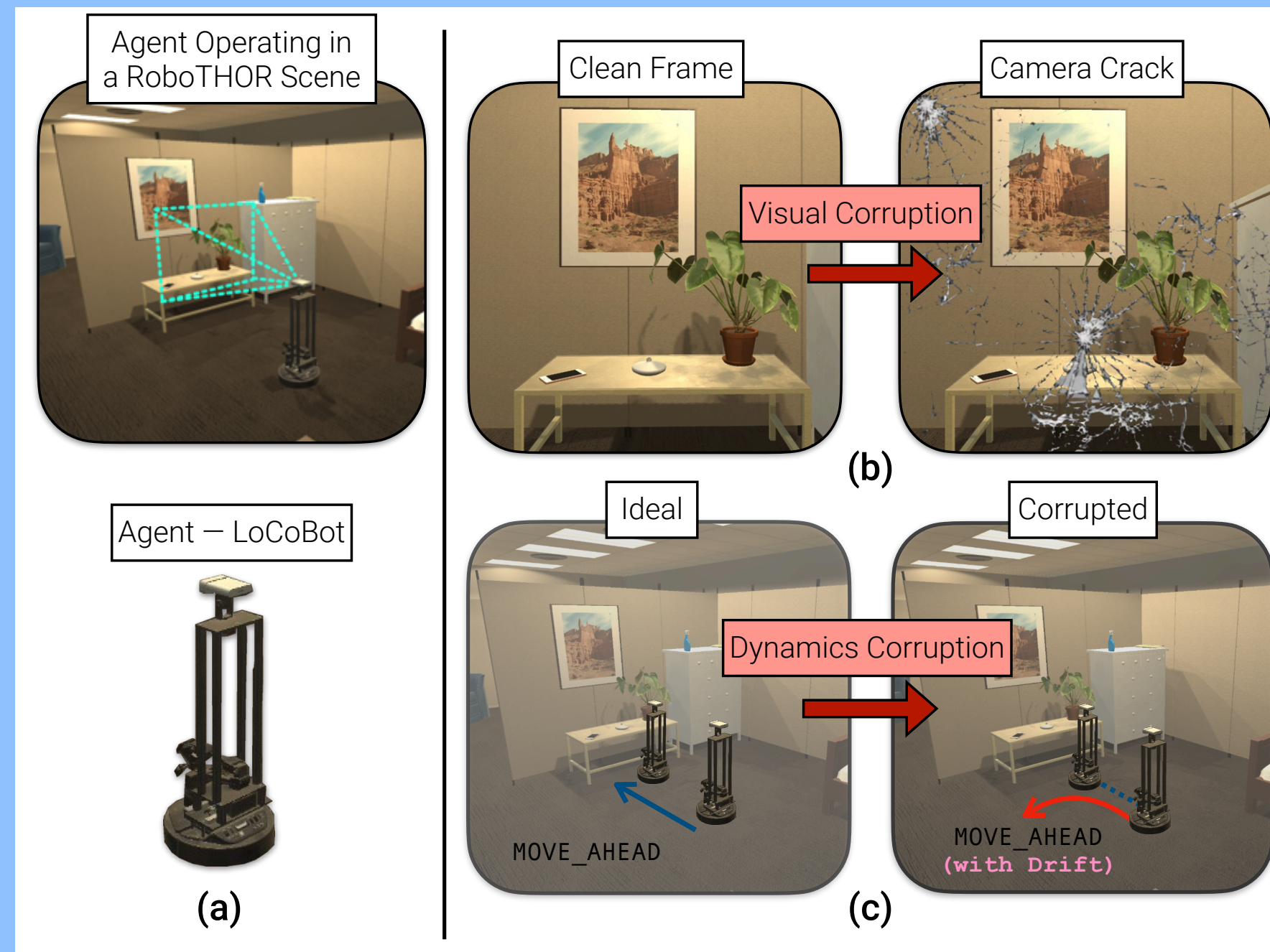
Need benchmarks to define expectations

Resilience

Withstand or adapt to a diverse set of visual conditions

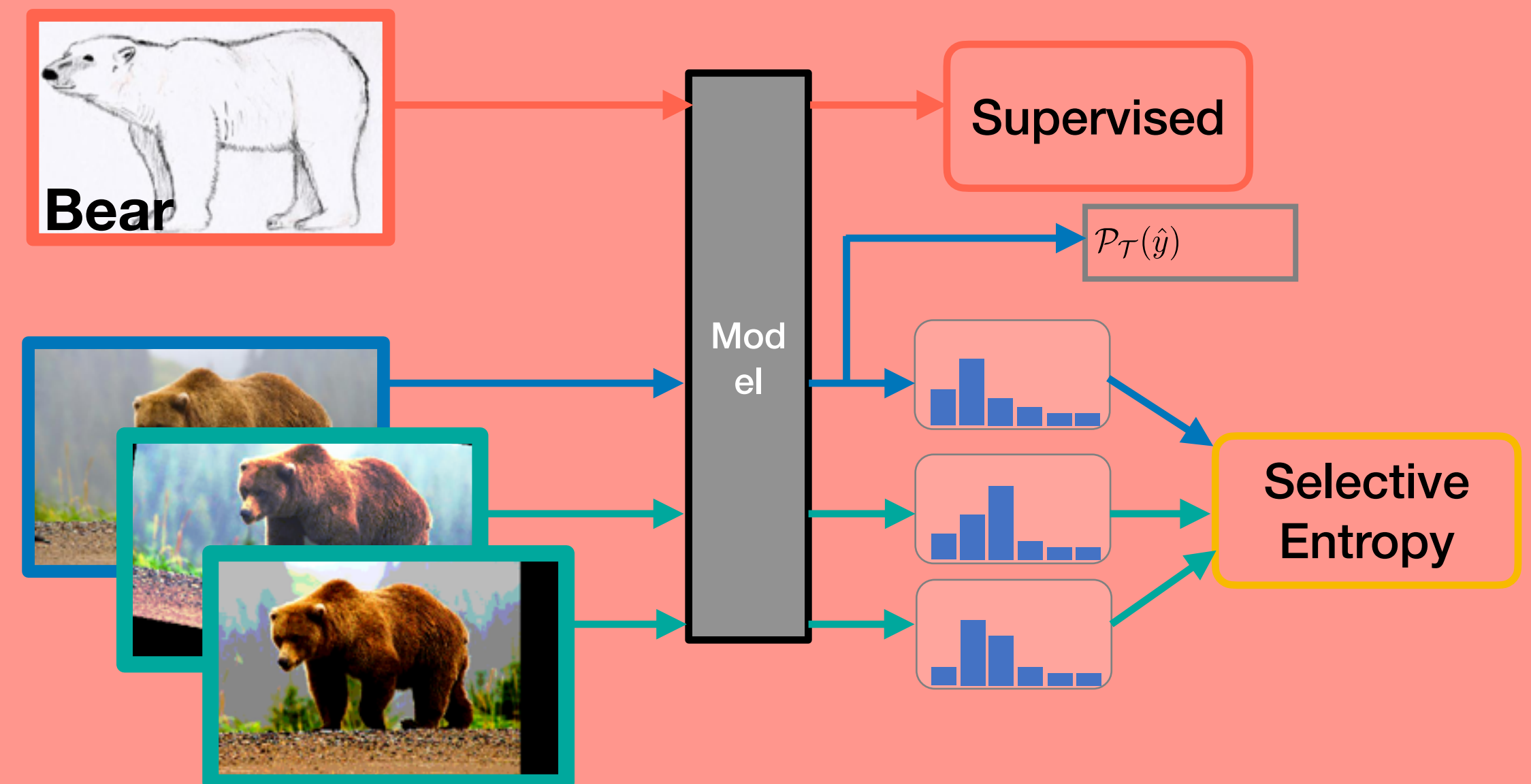
Summary: Responsible Vision

Benchmarks for Analysis



RobustNav for Embodied Nav study
Chattopadhyay et al, ICCV 2021

Domain Adaptation

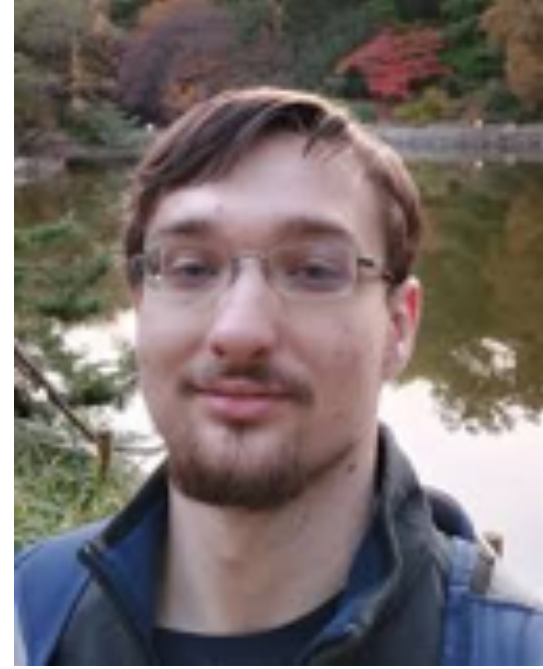


SENTRY: Selective Updates
Prabhu et al, ICCV 2021

Thank you



Sean Foley



Daniel Bolya



Sruthi Sudhakar



George Stoica



Aayushi Agarwal



Kartik
Sarangmath



Prithvijit
Chattopadhyay



Viraj Prabhu



Shivam Khare



Deeksha Karthik

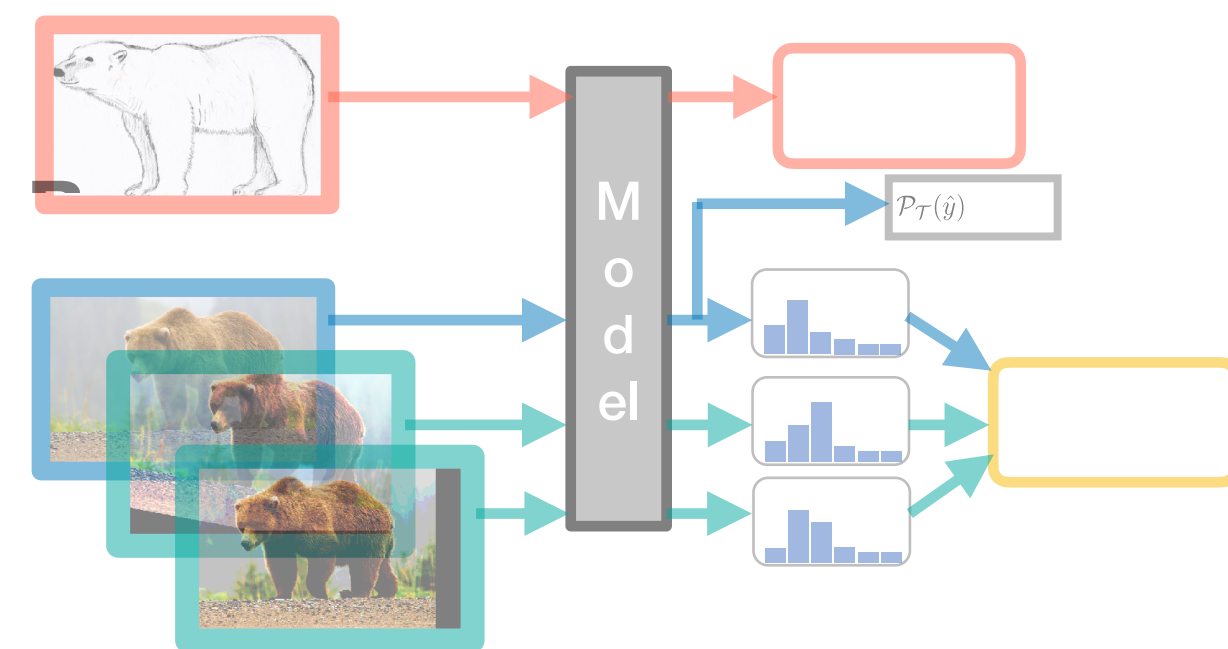
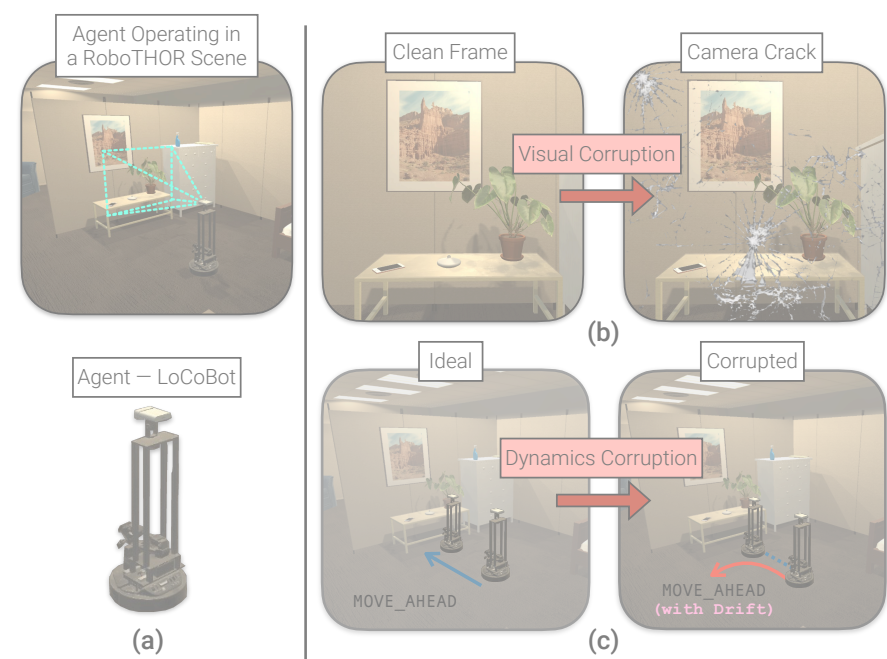


Bhavika Devnani



Deepanshi
Deepanshi

Summary: Responsible Vision



Domain Adaptation

SENTRY: Selective Updates
Prabhu et al, ICCV 2021

Benchmarks for Analysis

RobustNav for Embodied Nav study
Chattopadhyay et al, ICCV 2021

Thank you!
Questions?
{judy,virajp}@gatech