

Adversarial Domain Adaptation and Adversarial Robustness

Judy Hoffman

facebook

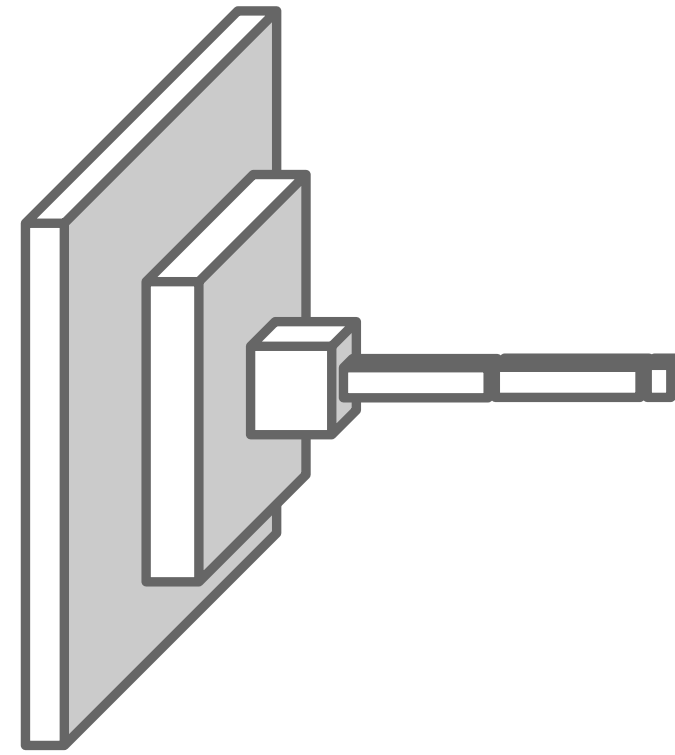
Artificial Intelligence Research





Big
data

+



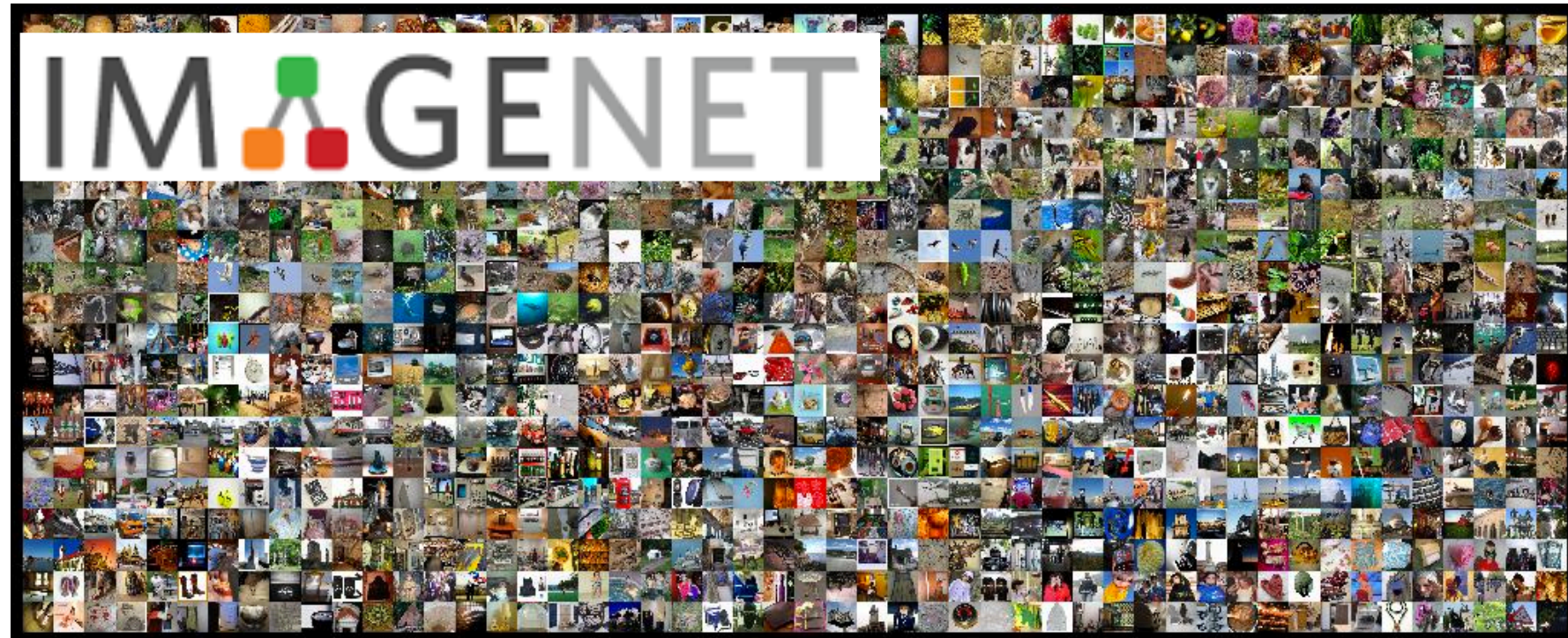
Deep
learning

=



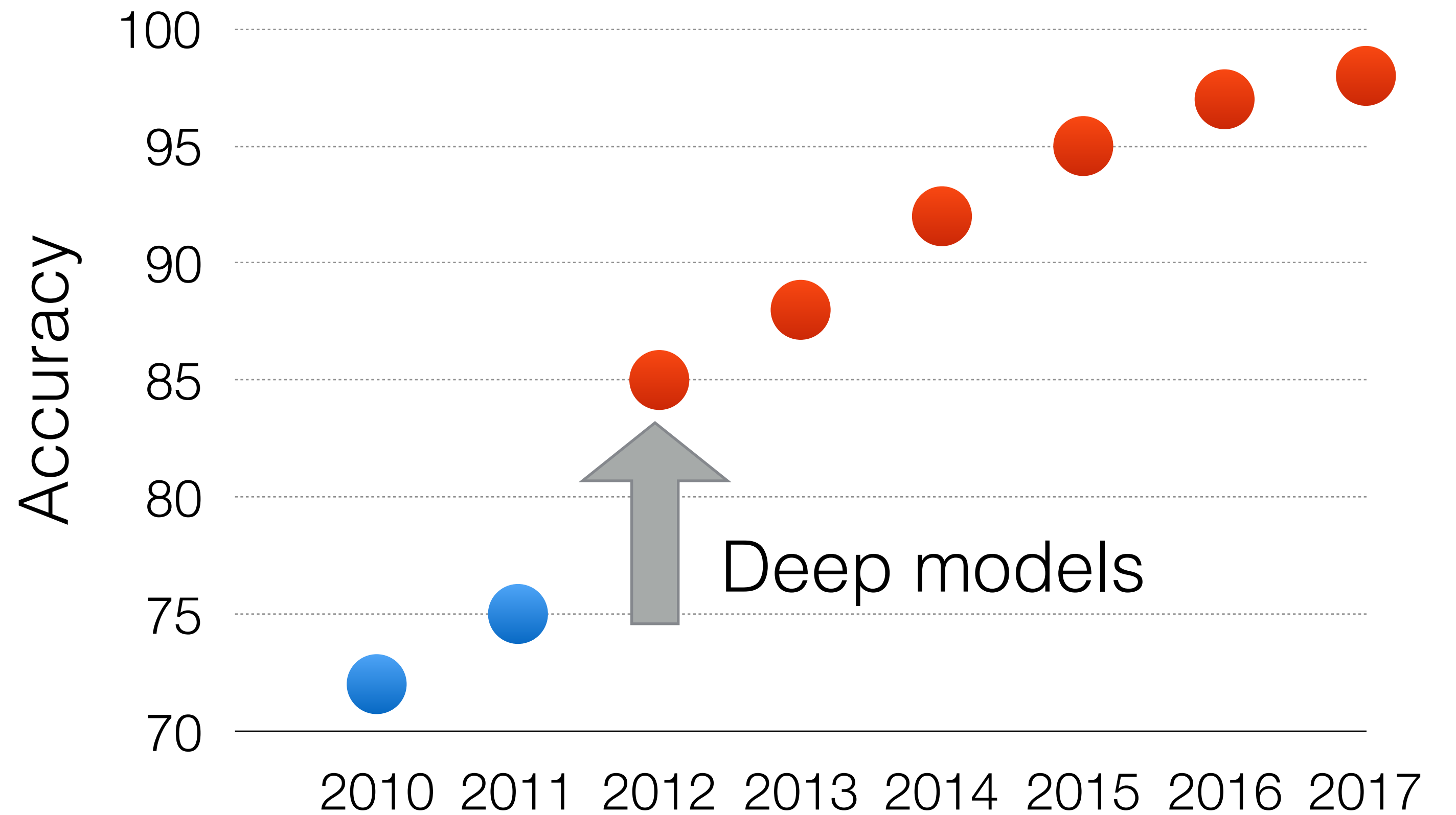
success

Benchmark Performance



Millions of Images

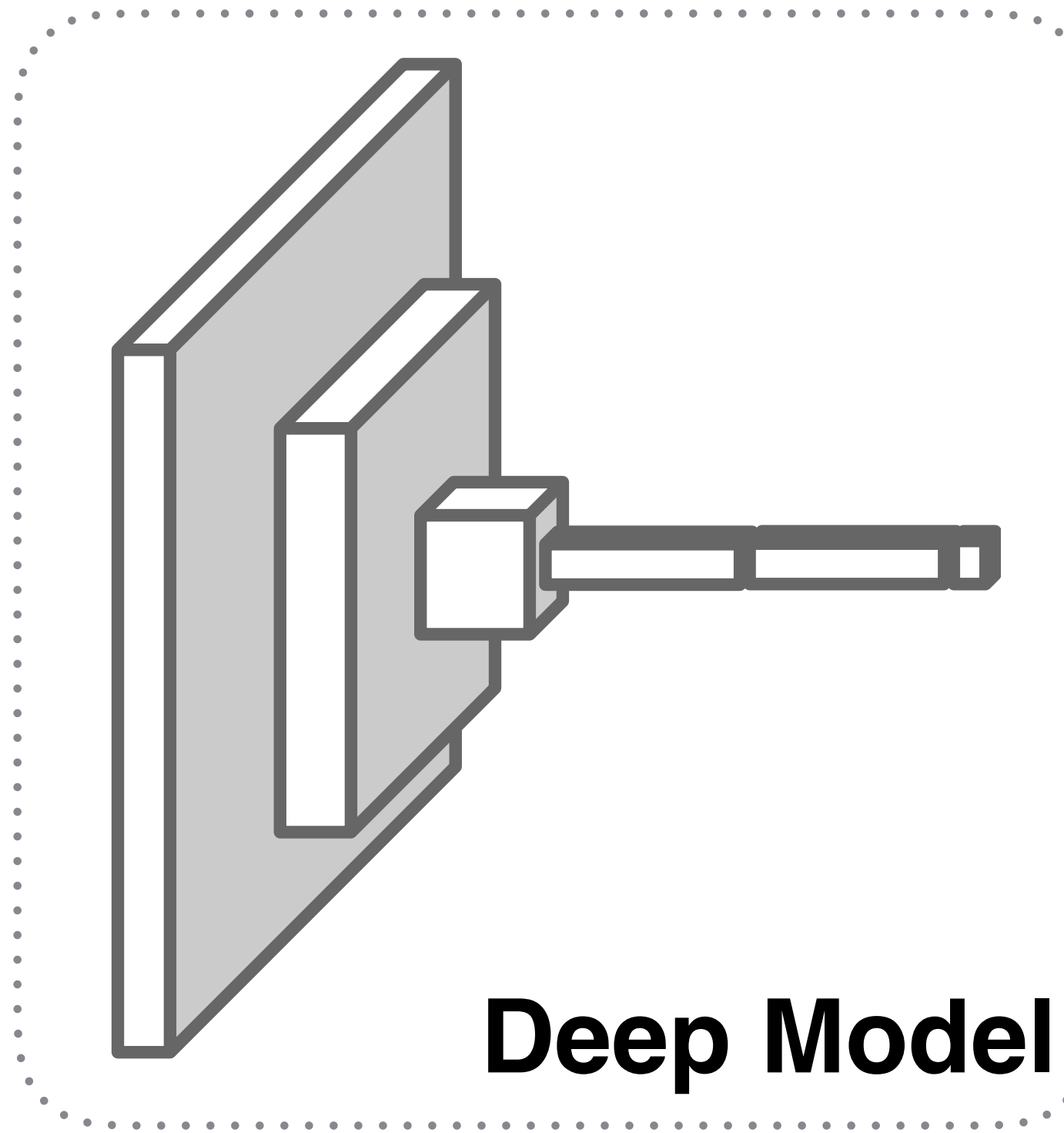
**Challenge to recognize
1000 categories**



Dataset Bias



Test Image



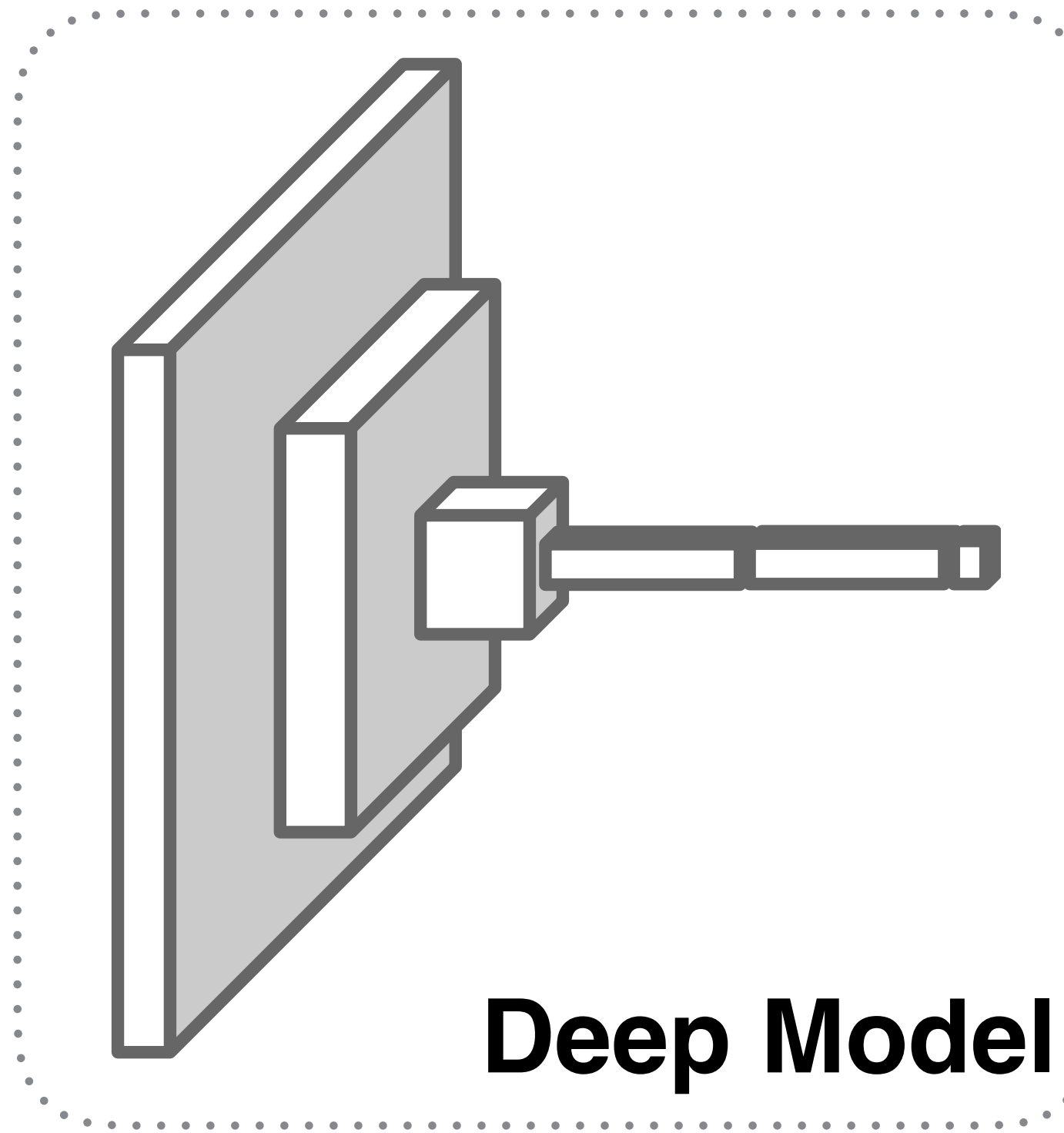
Deep Model



Dataset Bias



Test Image



Deep Model



Dataset Bias



Test Image

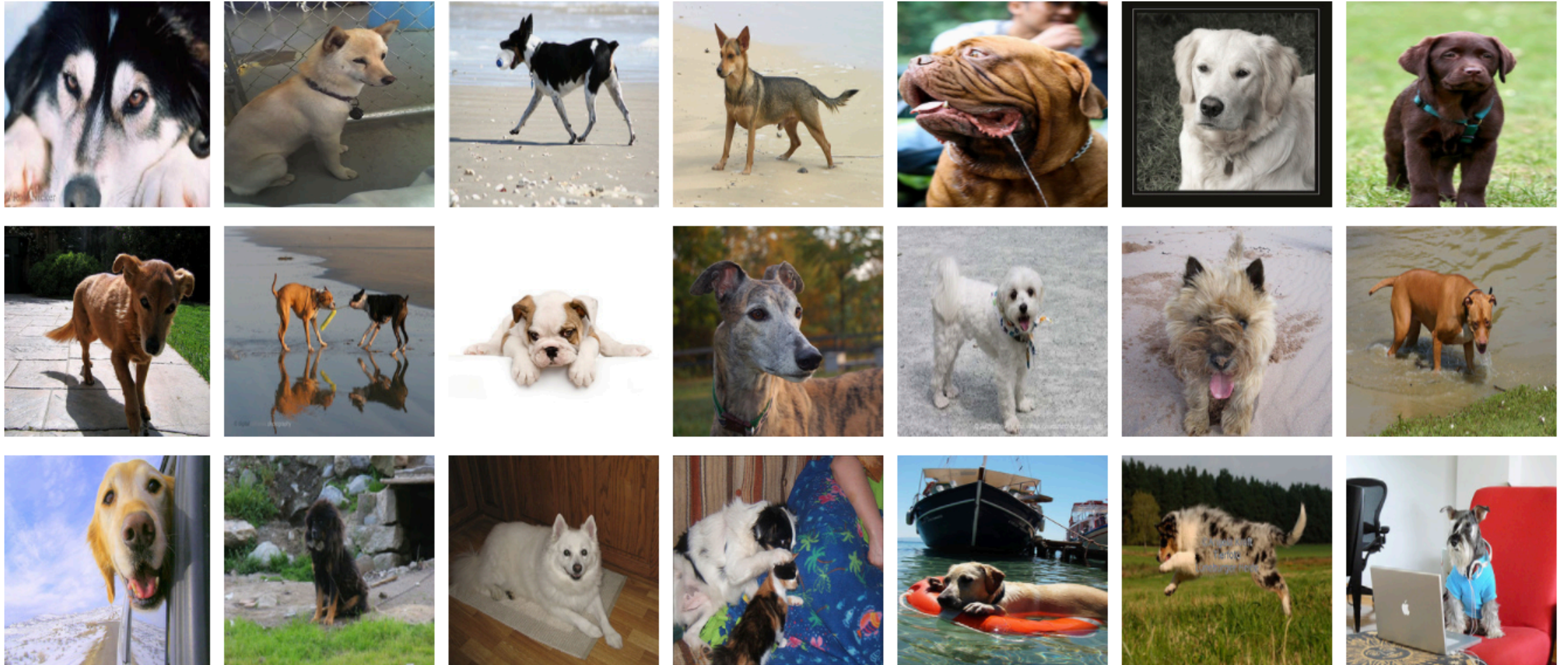
Dog is not recognized



Deep Model

?

Dataset Bias



Dataset Bias



Low resolution

Dataset Bias



Low resolution



Motion Blur

Dataset Bias



Low resolution



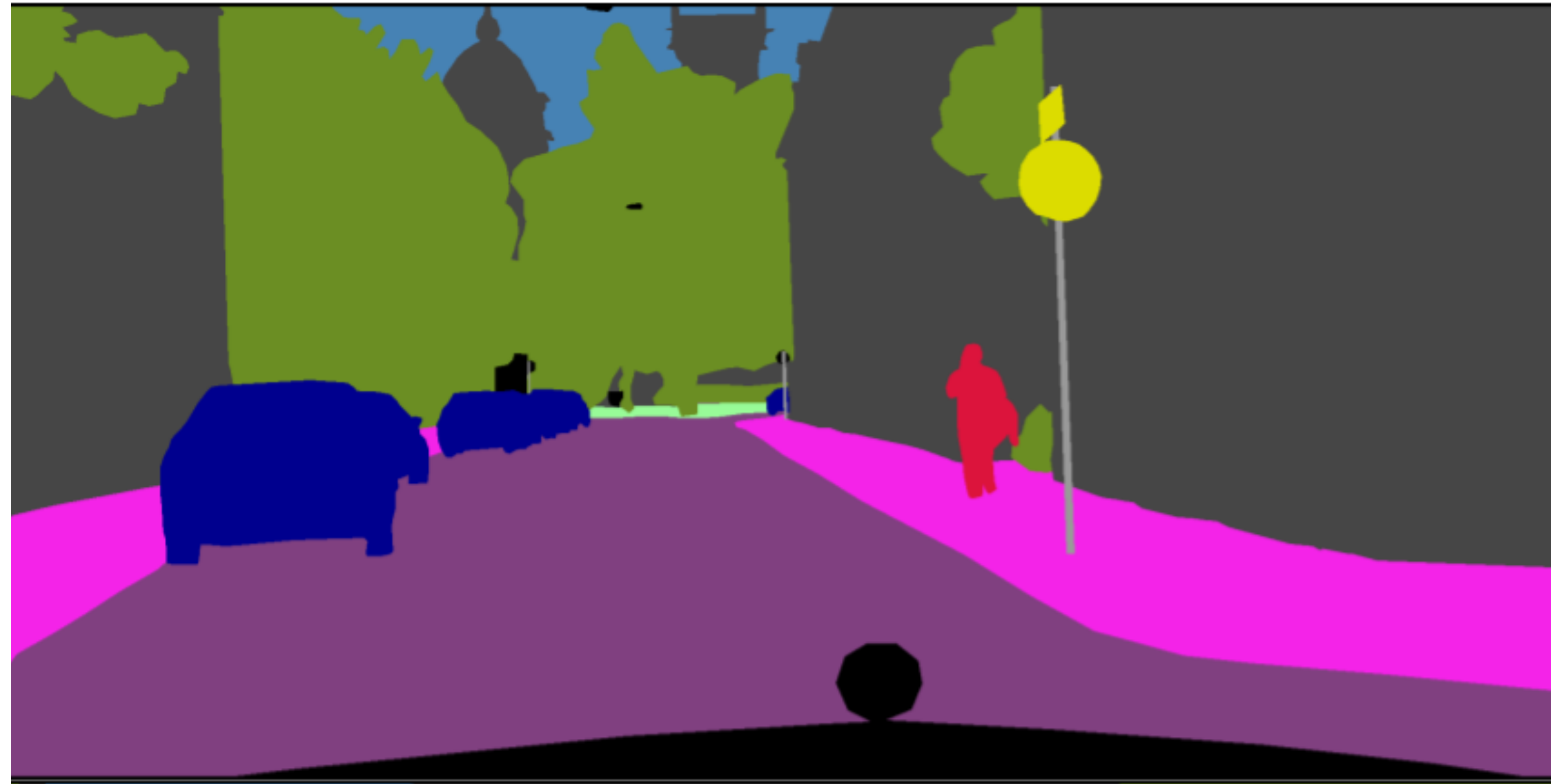
Motion Blur



Pose Variety

Why not collect new annotations?

Why not collect new annotations?



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

Why not collect new annotations?



Expensive
(\$10-12 per
image)

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

Why not collect new annotations?



Expensive
(\$10-12 per
image)

Large Potential for Change
Different: Weather, City, Car

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

Why not collect new annotations?

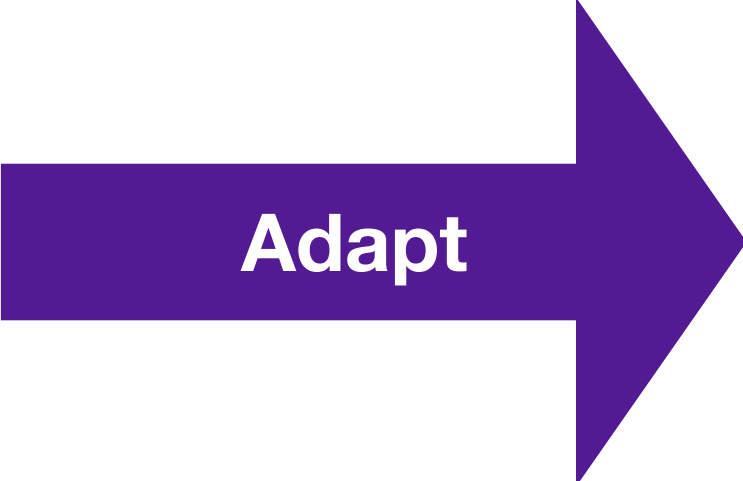
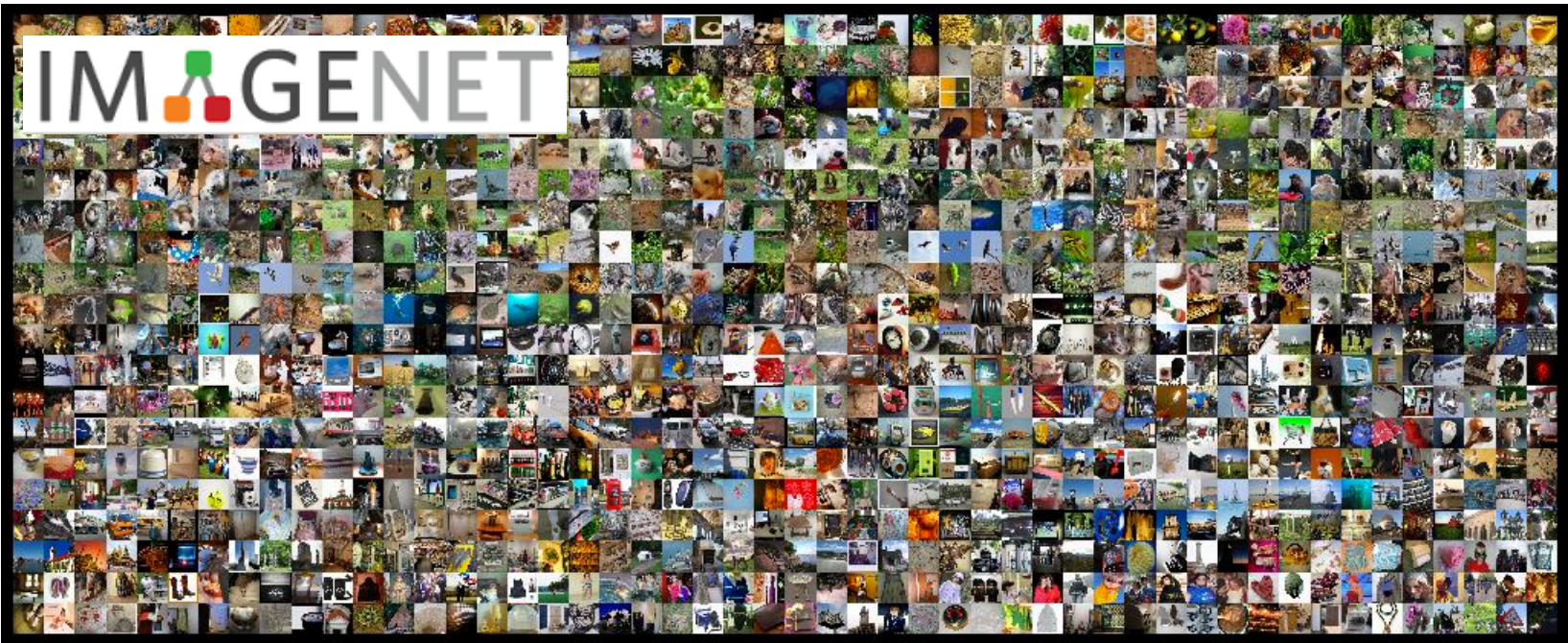


Proprietary



Private

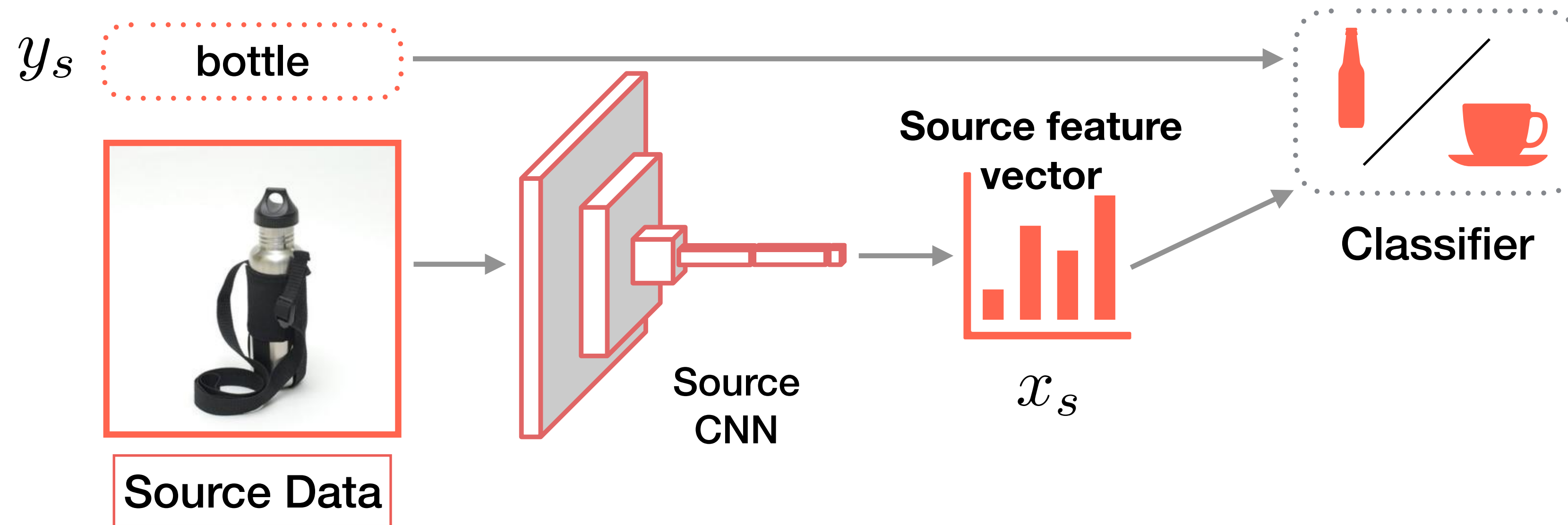
Domain Adaptation: Train on Source Test on Target



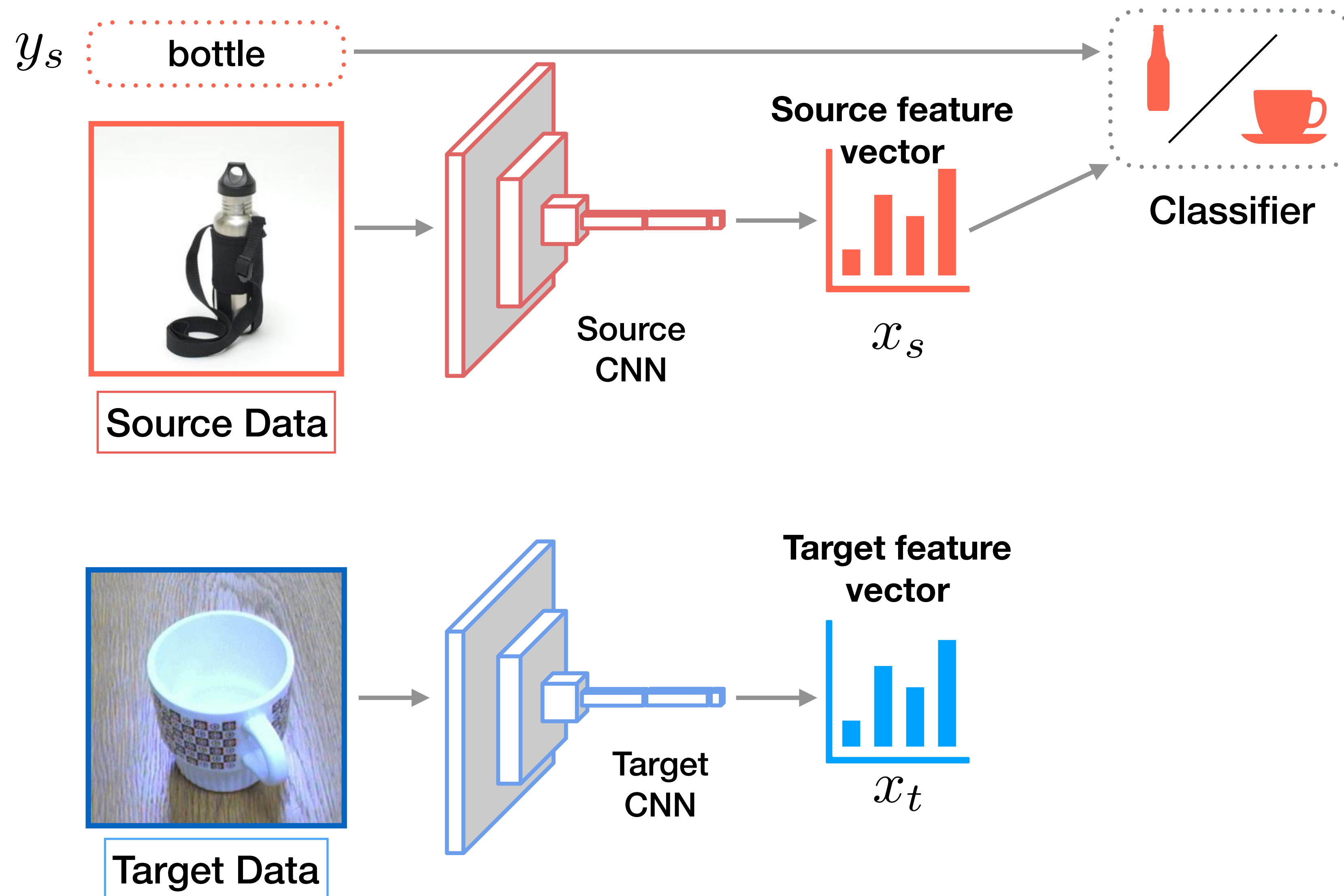
Source Domain $\sim P_S(X_S, Y_S)$
lots of **labeled** data

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

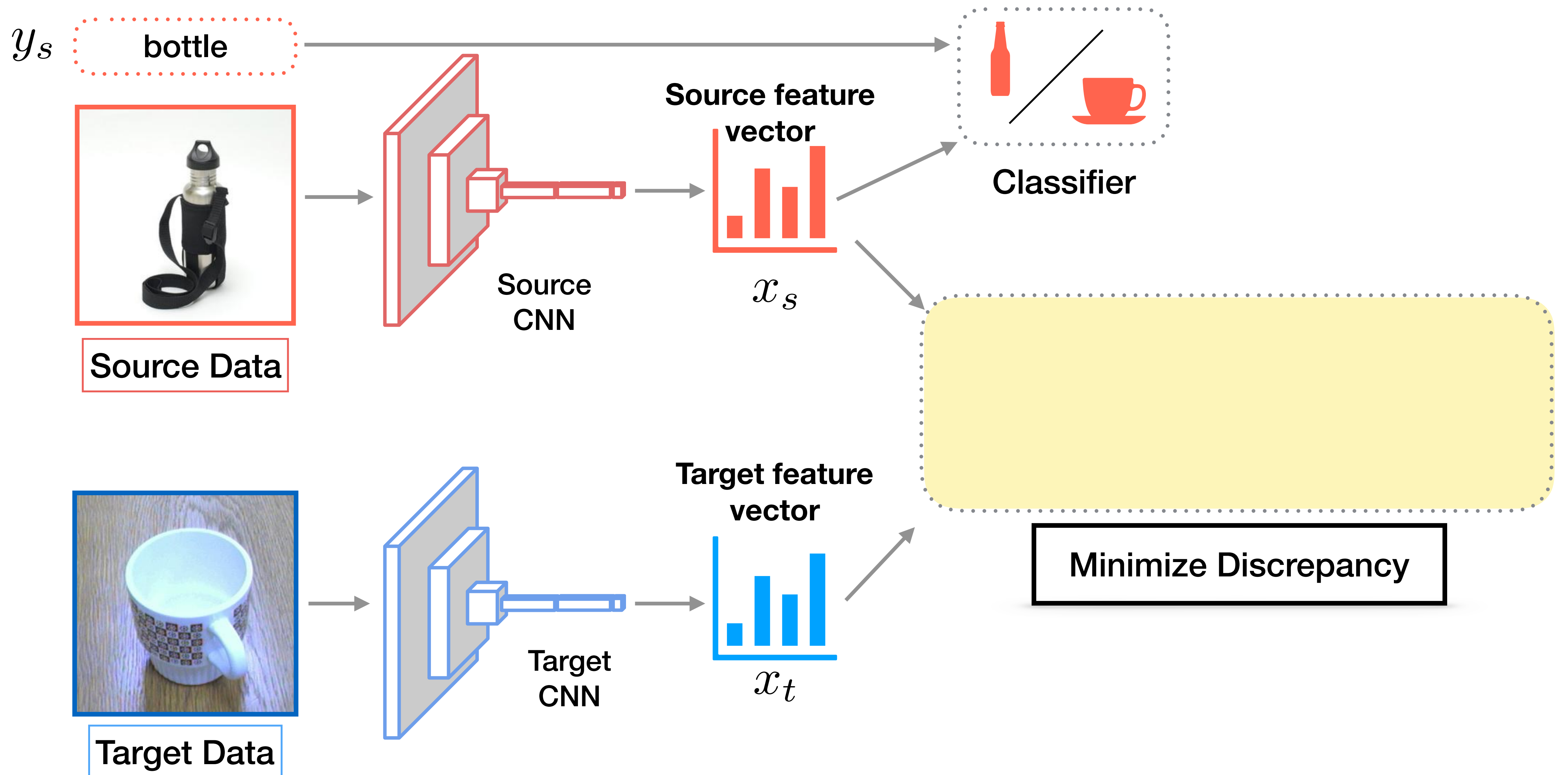
Adversarial Domain Adaptation



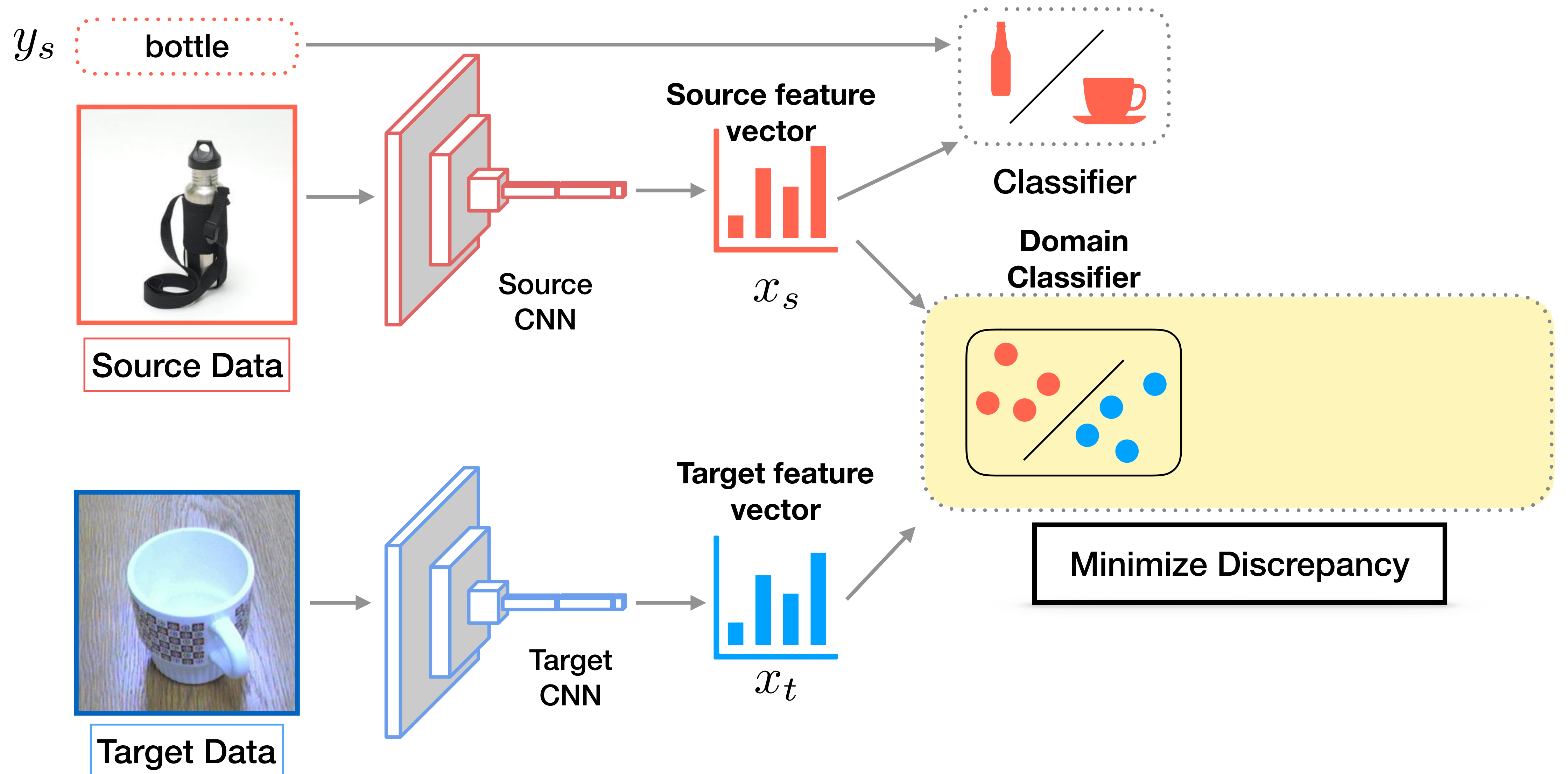
Adversarial Domain Adaptation



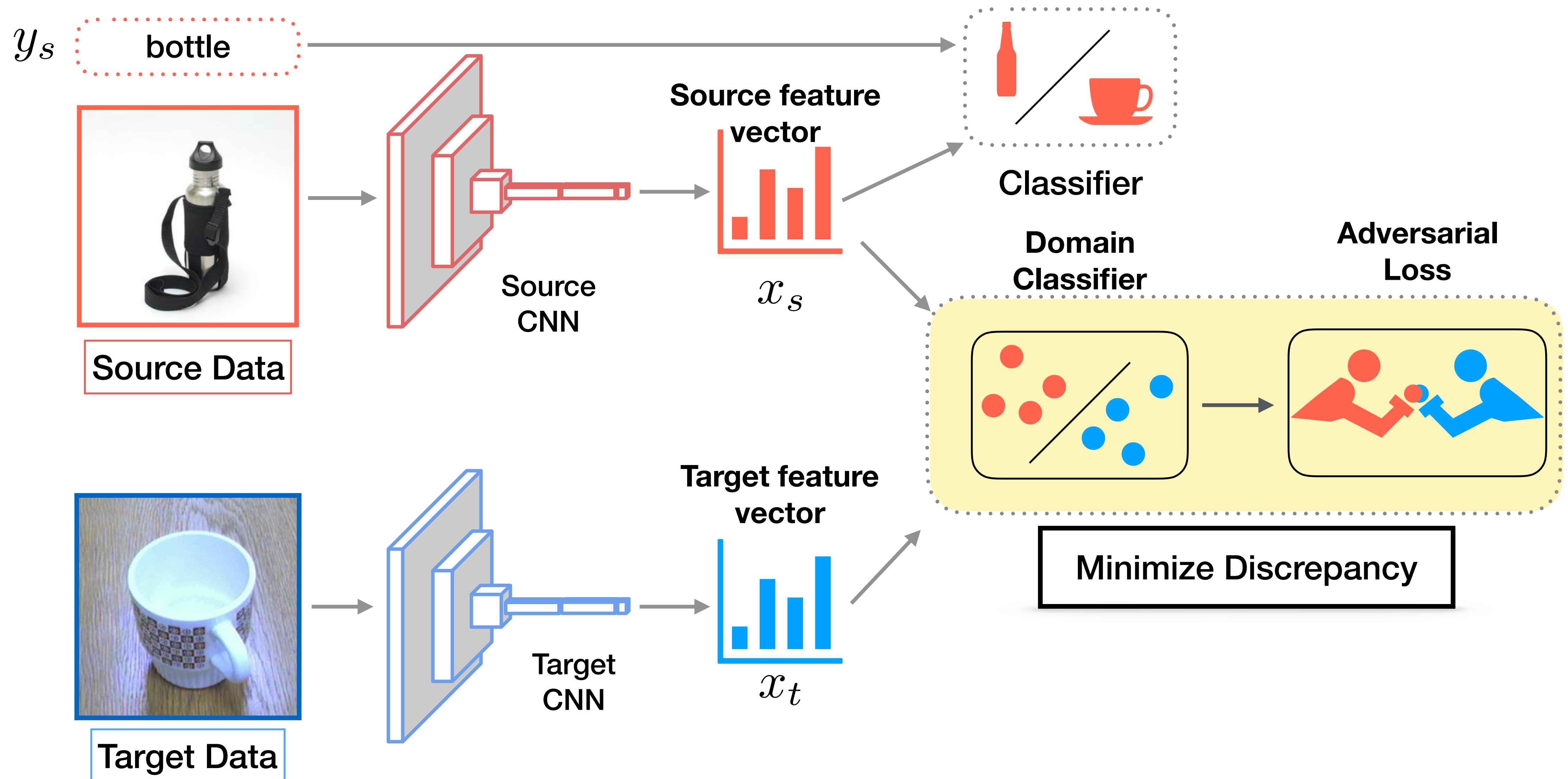
Adversarial Domain Adaptation



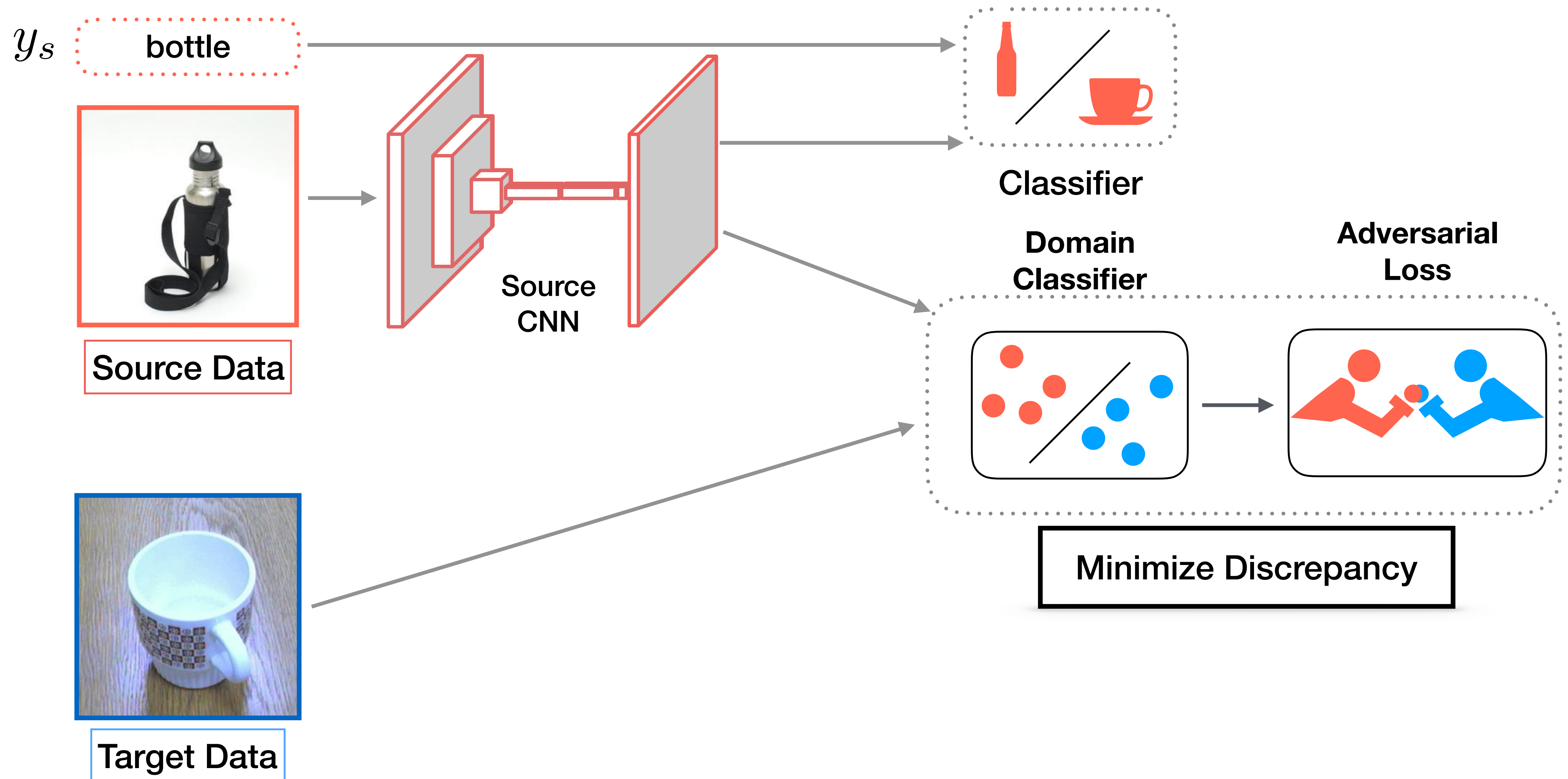
Adversarial Domain Adaptation



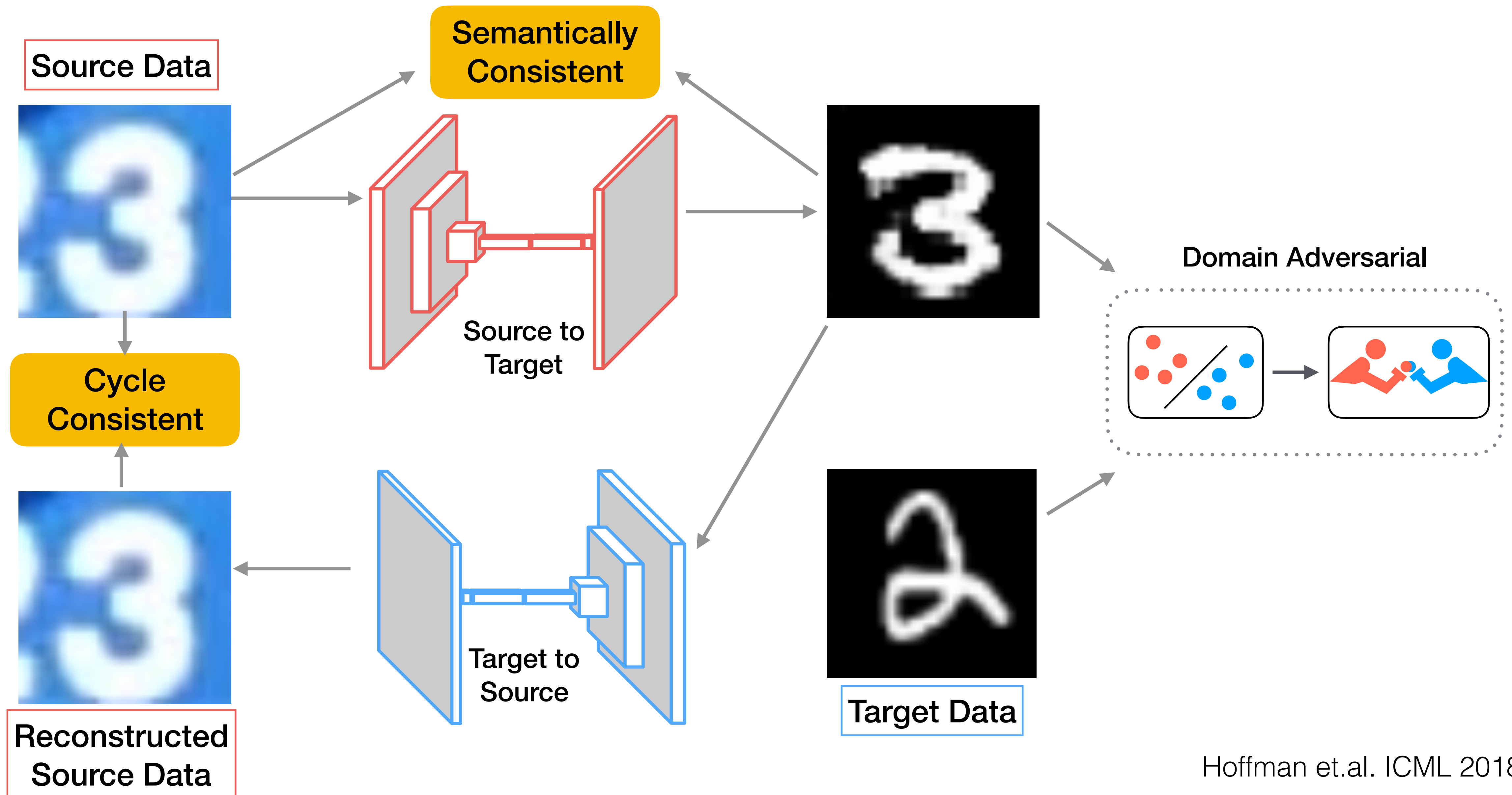
Adversarial Domain Adaptation



Adversarial Domain Adaptation



CyCADA: Cycle Consistent Adversarial DA



Synthetic to Real Pixel Adaptation

Train



GTA (synthetic)

Test



CityScapes (Germany)

Synthetic to Real Pixel Adaptation



Synthetic to Real Pixel Adaptation



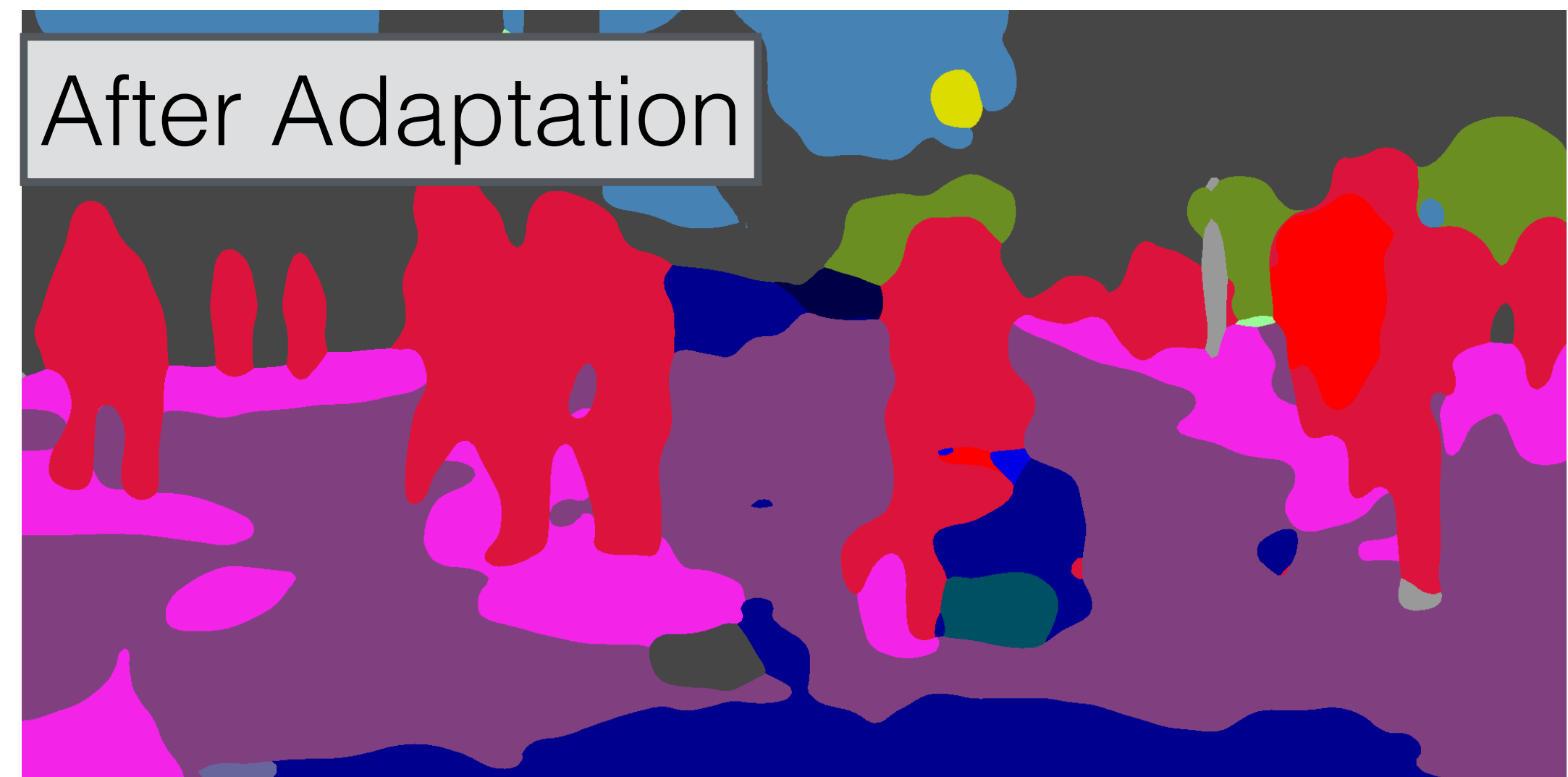
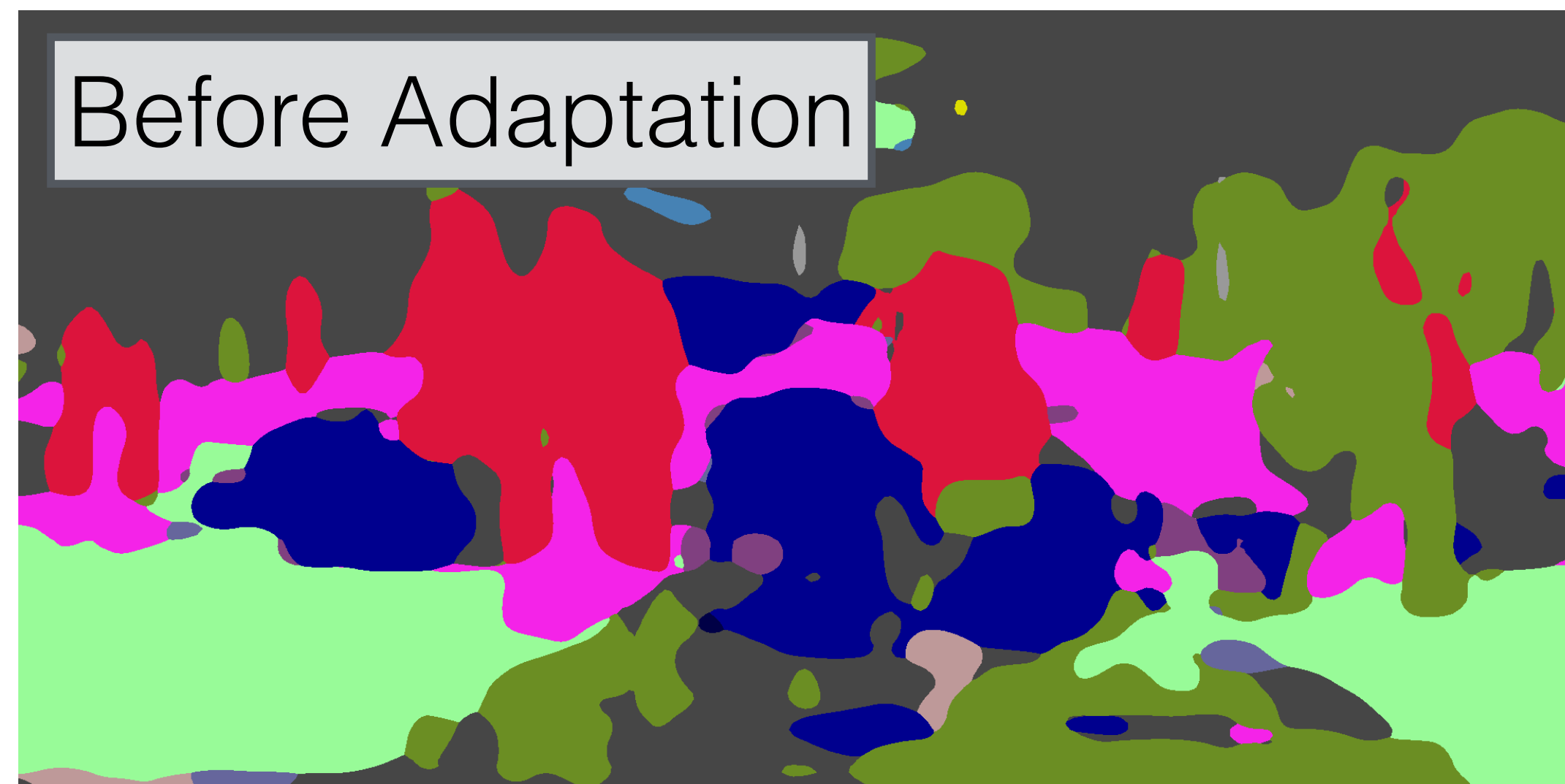
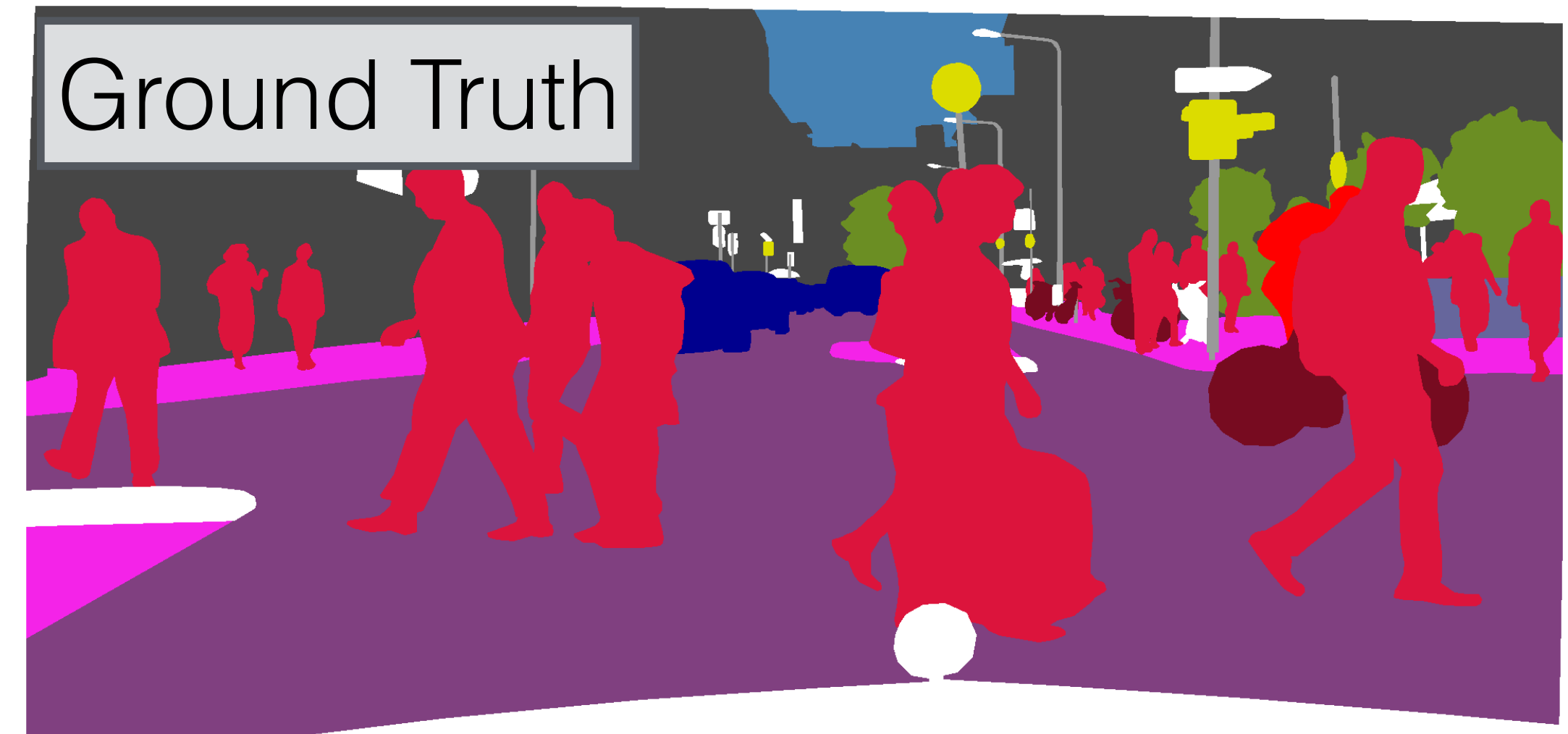
Synthetic to Real Pixel Adaptation



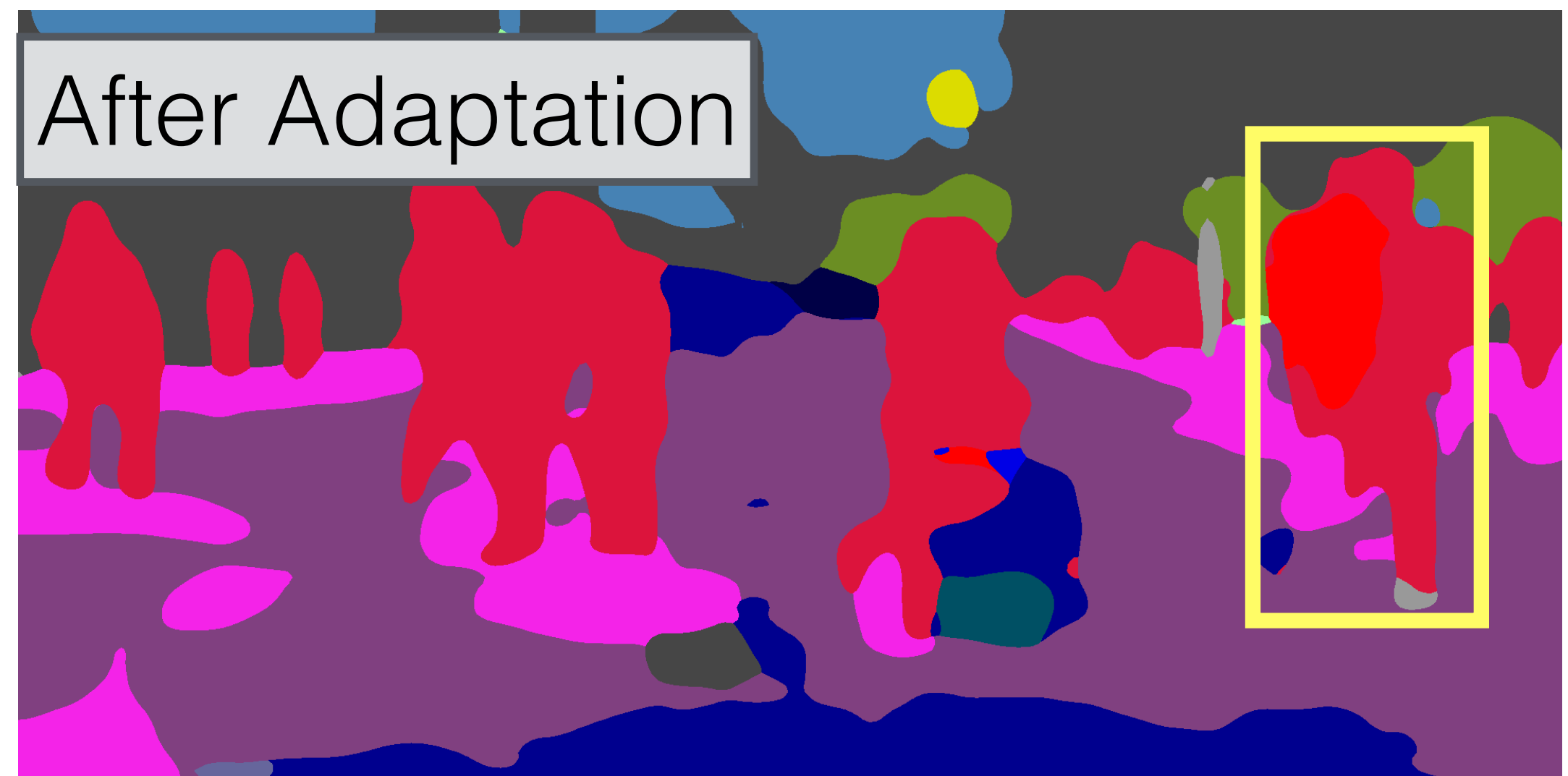
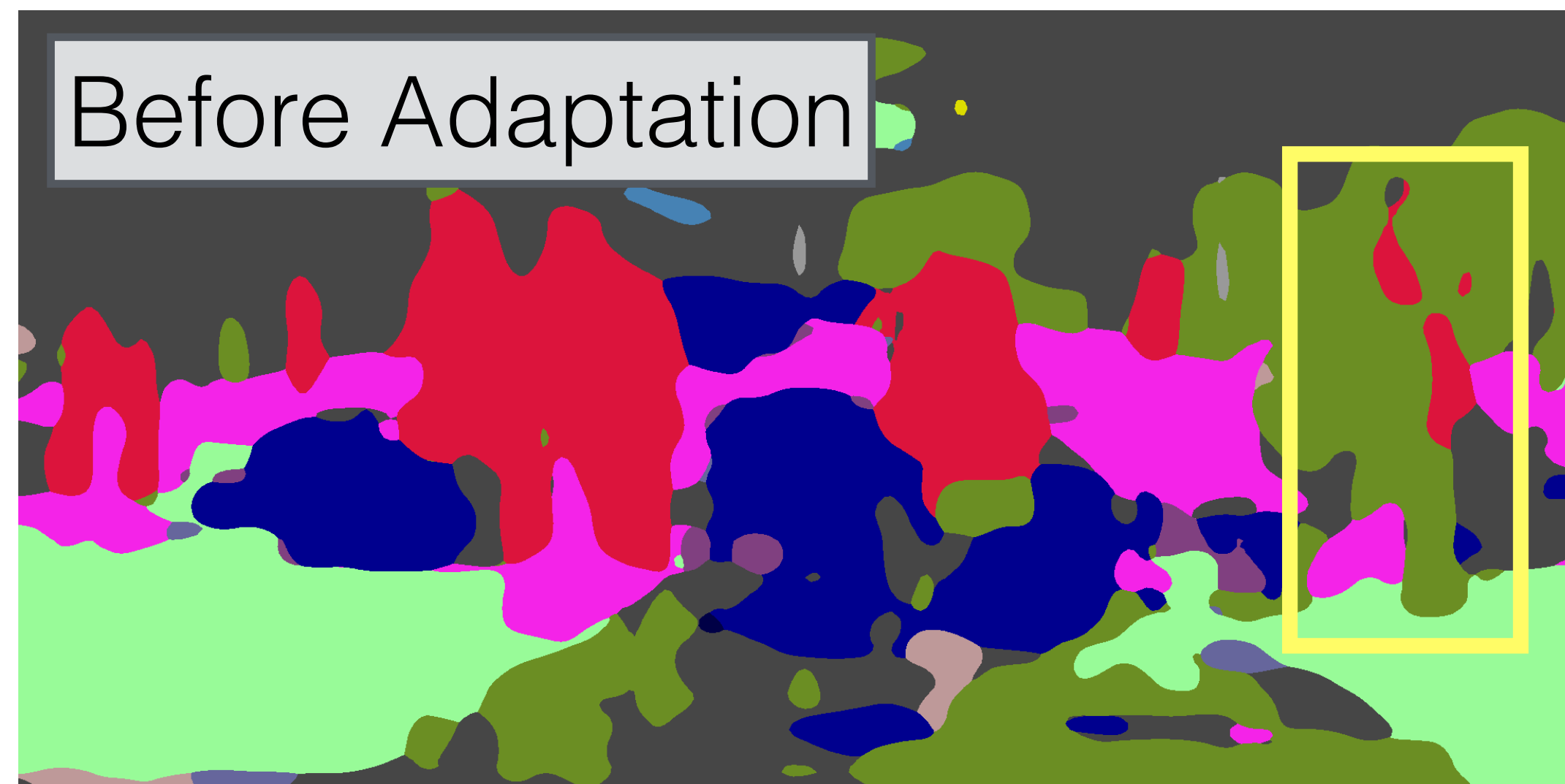
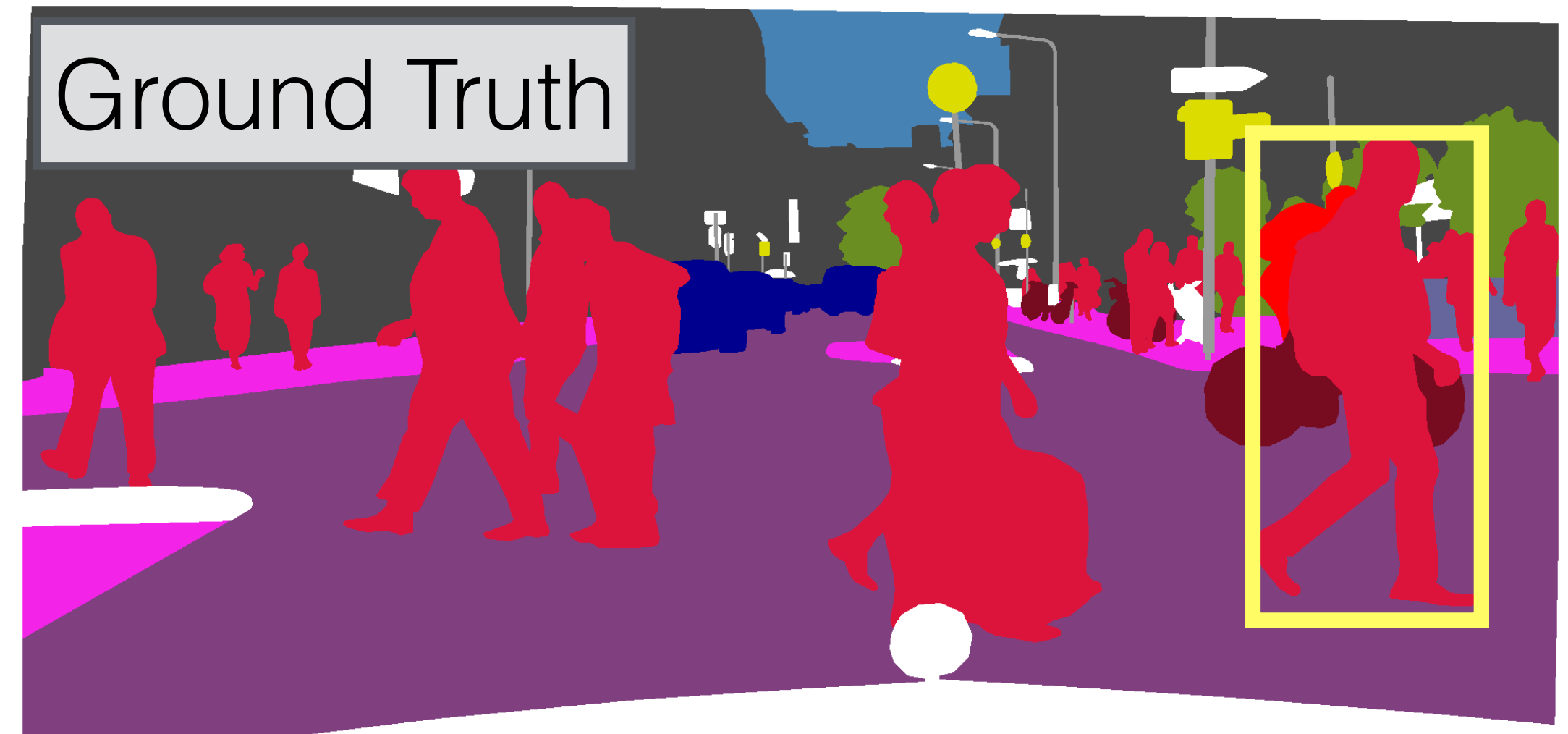
Synthetic to Real Pixel Adaptation



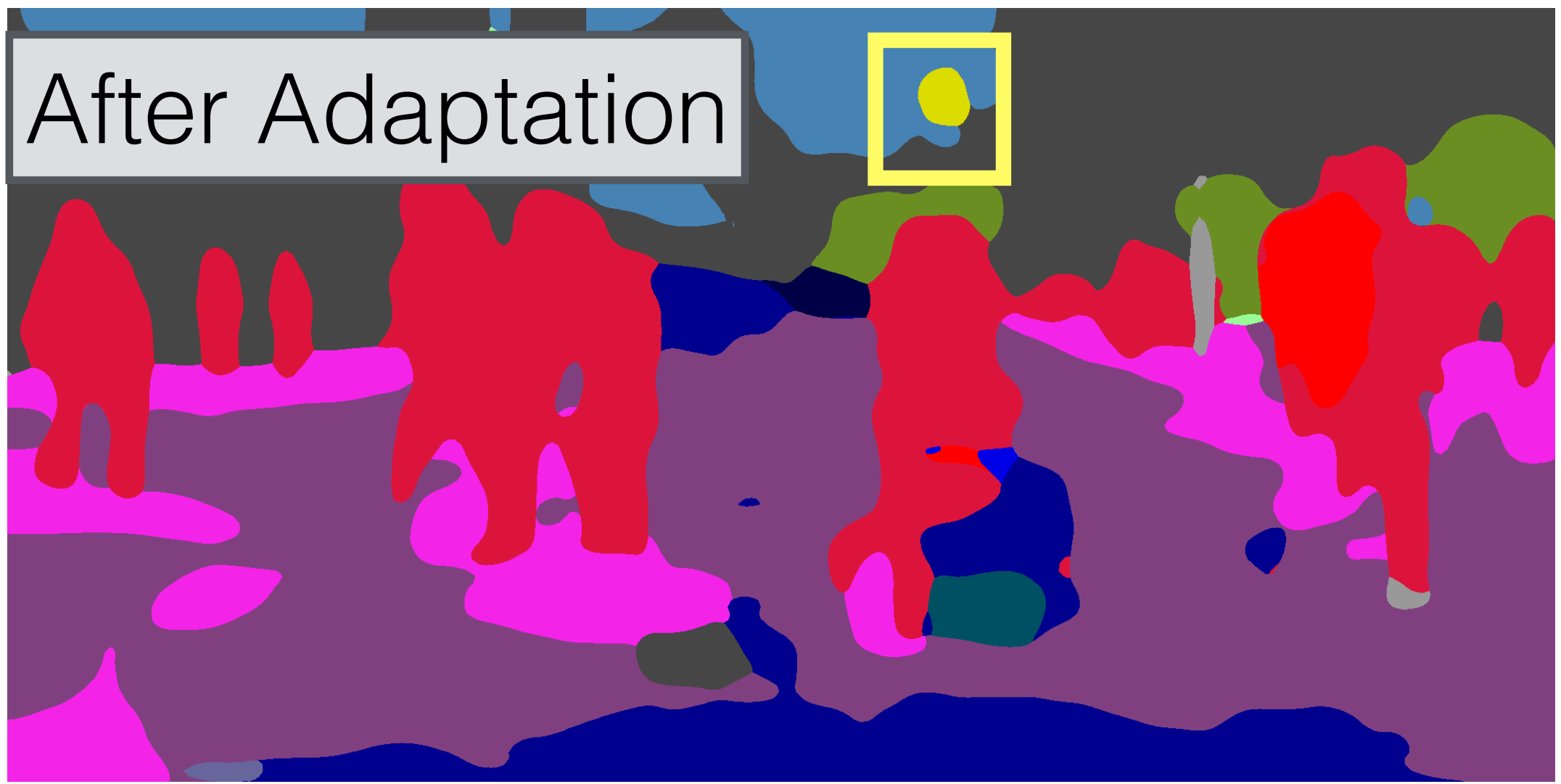
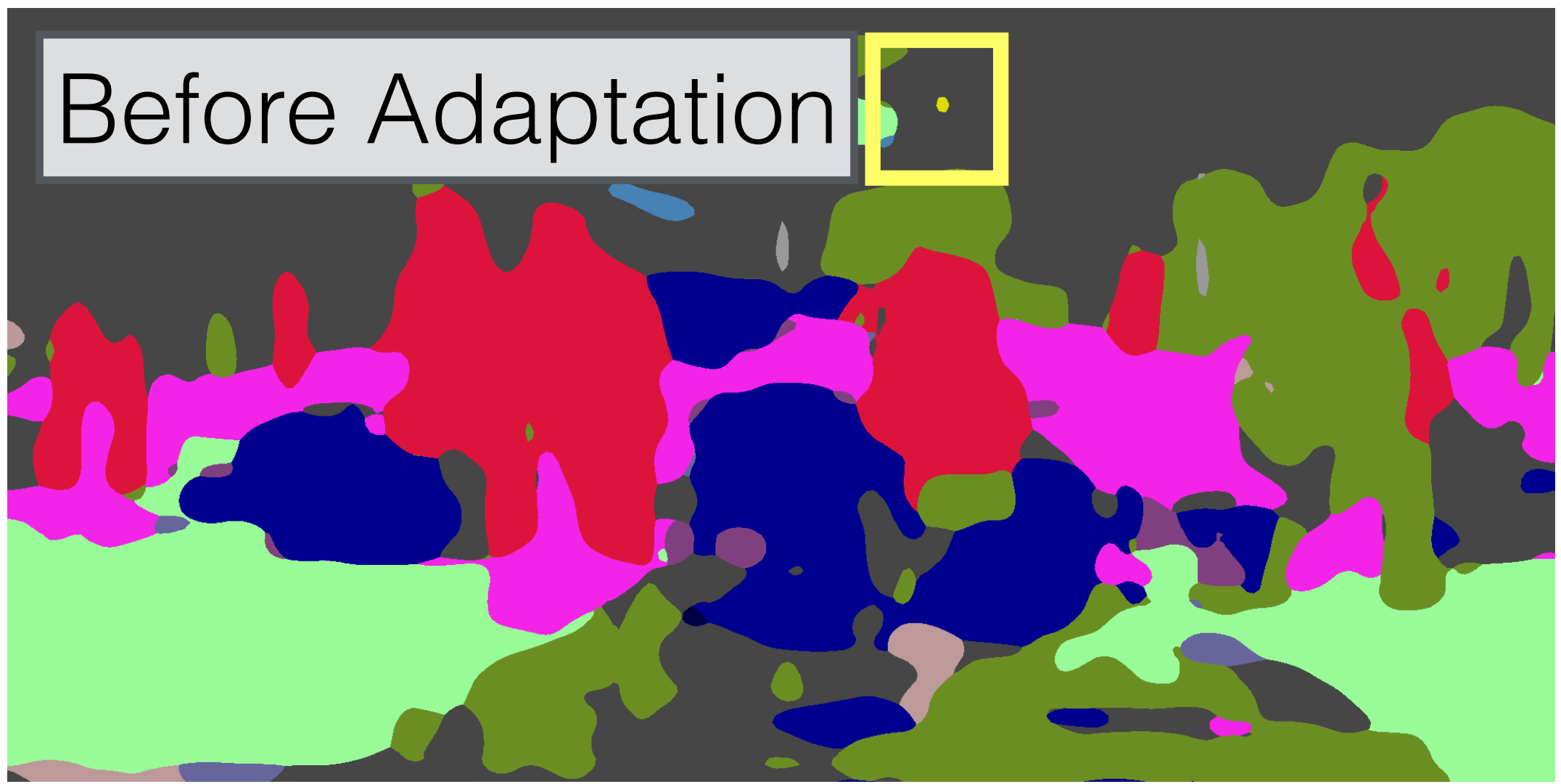
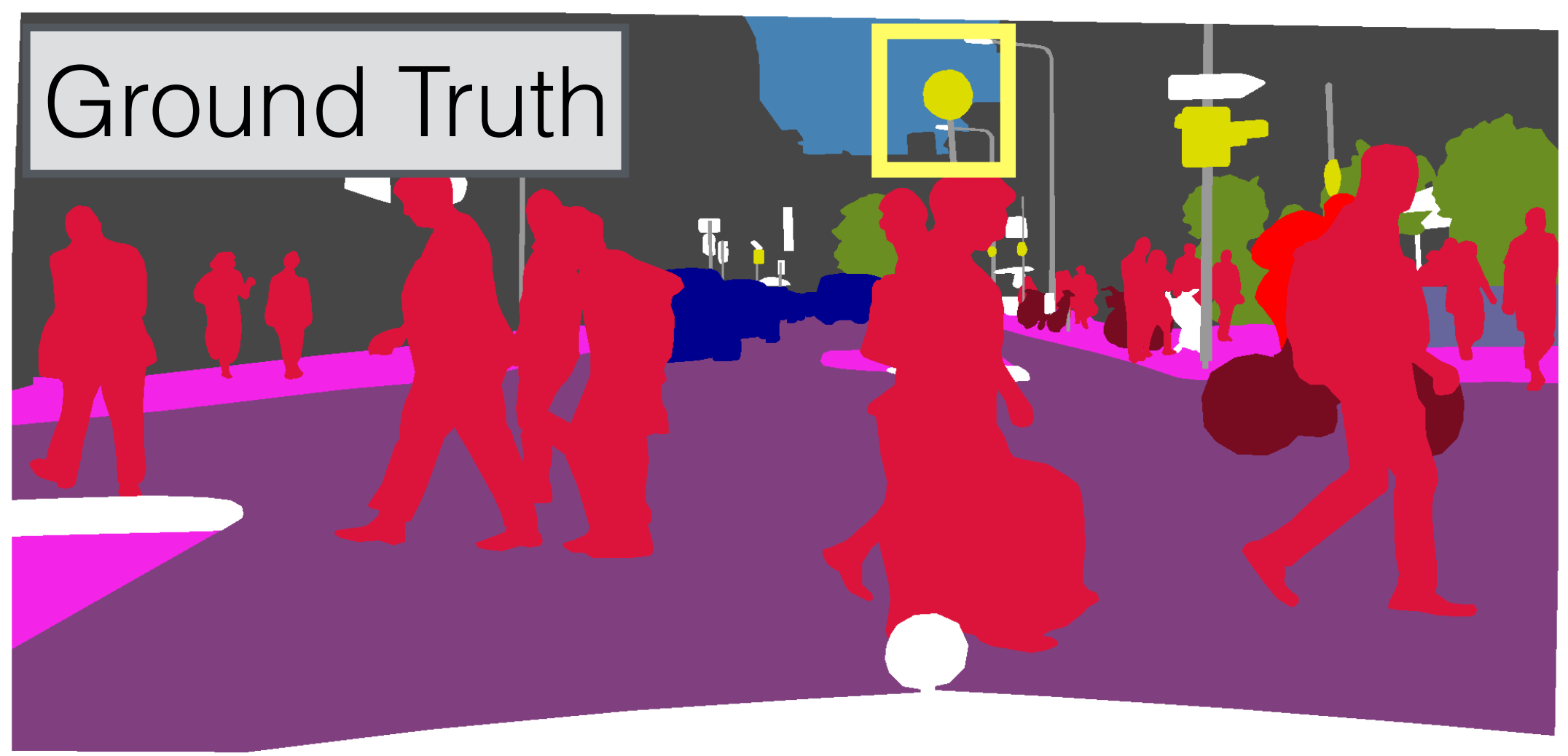
CyCADA Results: CityScapes Evaluation



CyCADA Results: CityScapes Evaluation



CyCADA Results: CityScapes Evaluation



So Far: Adapting to Natural Shifts



Adapt



So Far: Adapting to Natural Shifts



Adapt



What about
adversarial shifts?

Adversarial Examples



x
“panda”
57.7% confidence

+ .007 ×



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

=

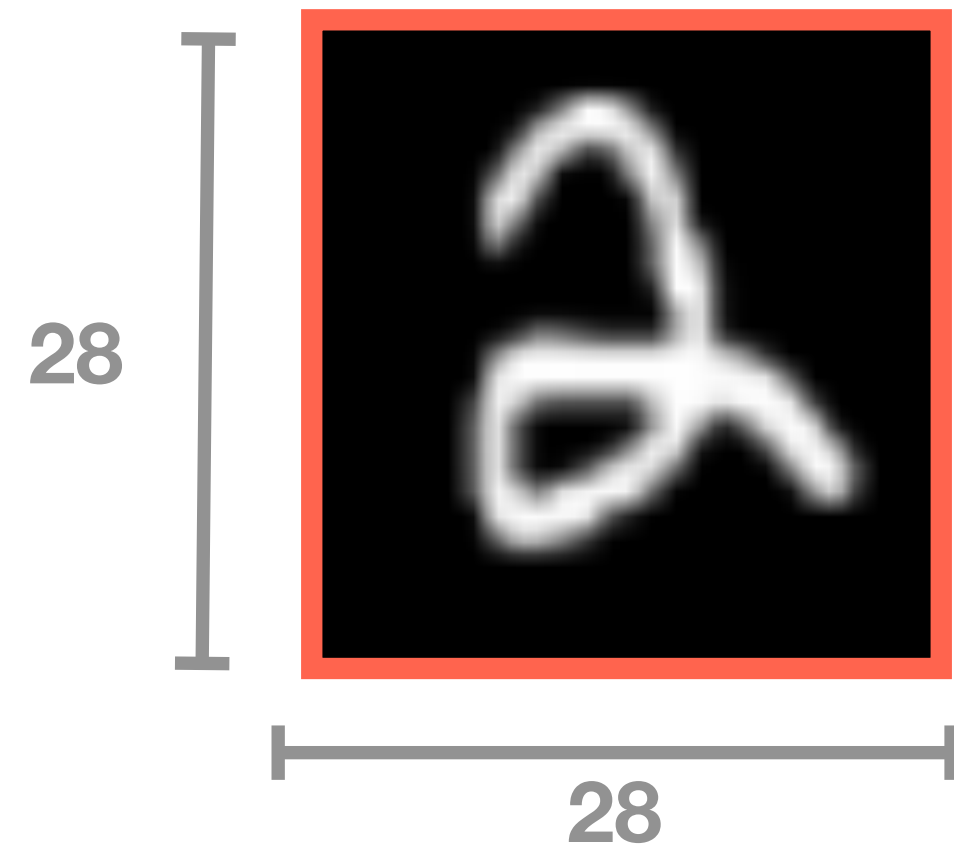


$x + \epsilon \text{sign}(\nabla_x J(\theta, x, y))$
“gibbon”
99.3 % confidence

Visualize Perturbation Space

Visualize Perturbation Space

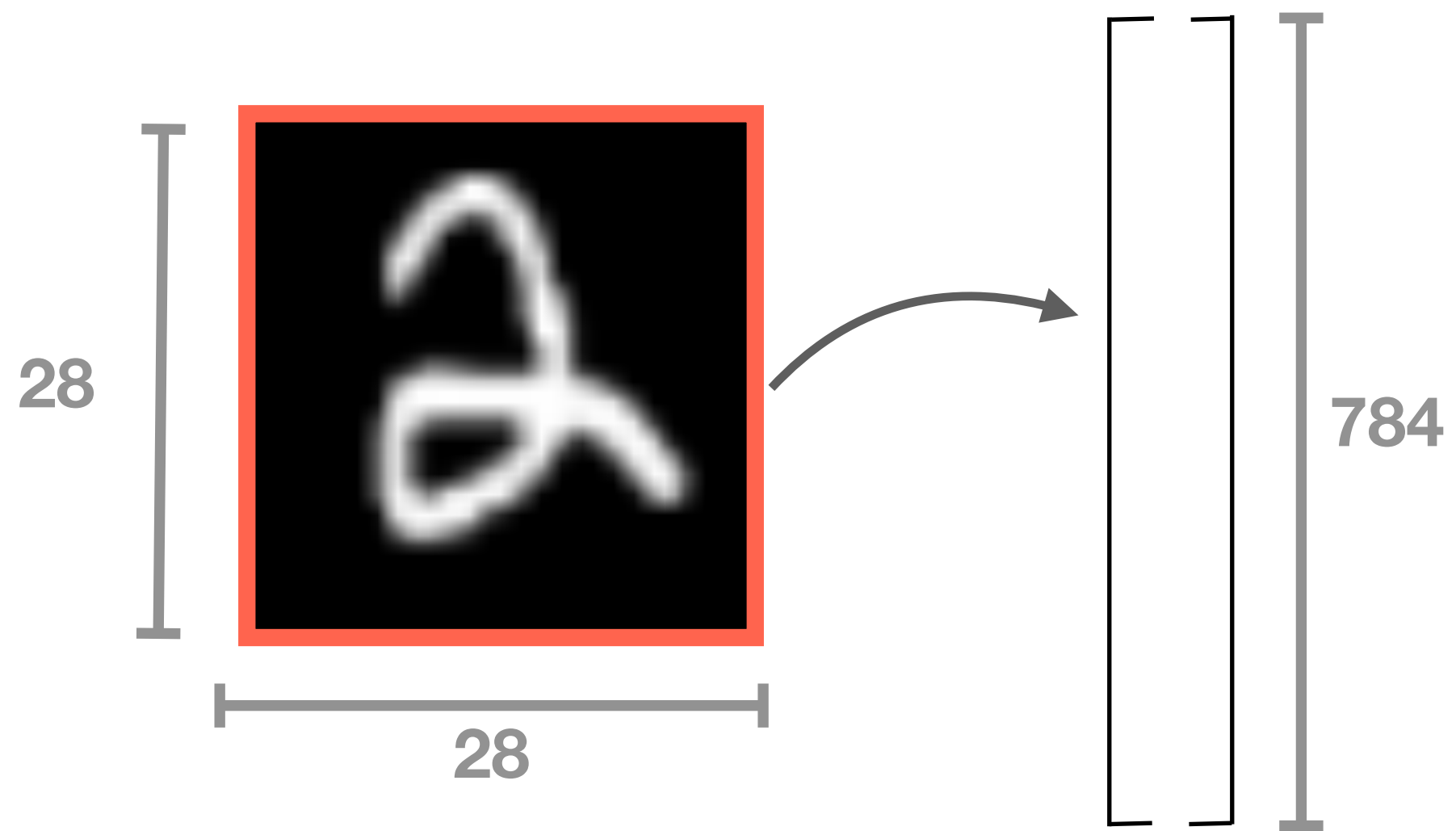
Training point



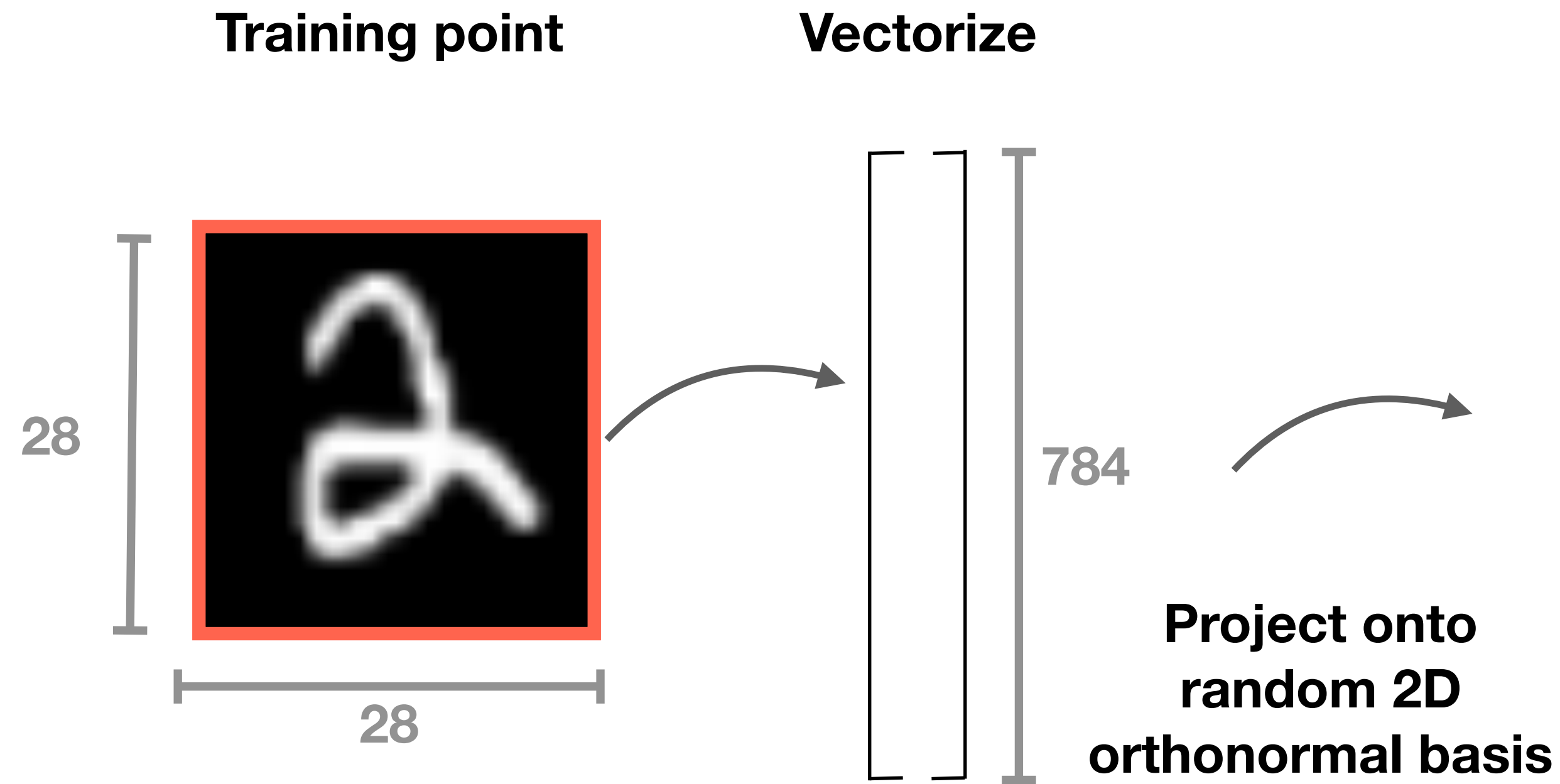
Visualize Perturbation Space

Training point

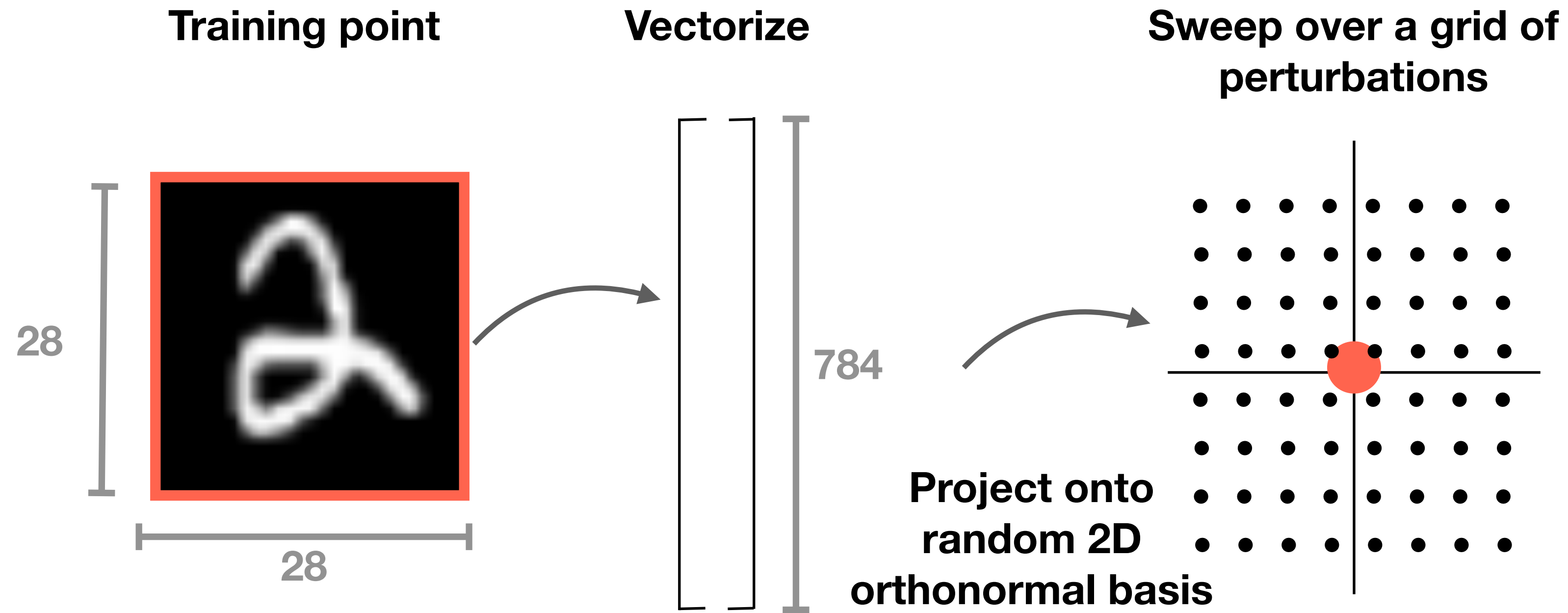
Vectorize



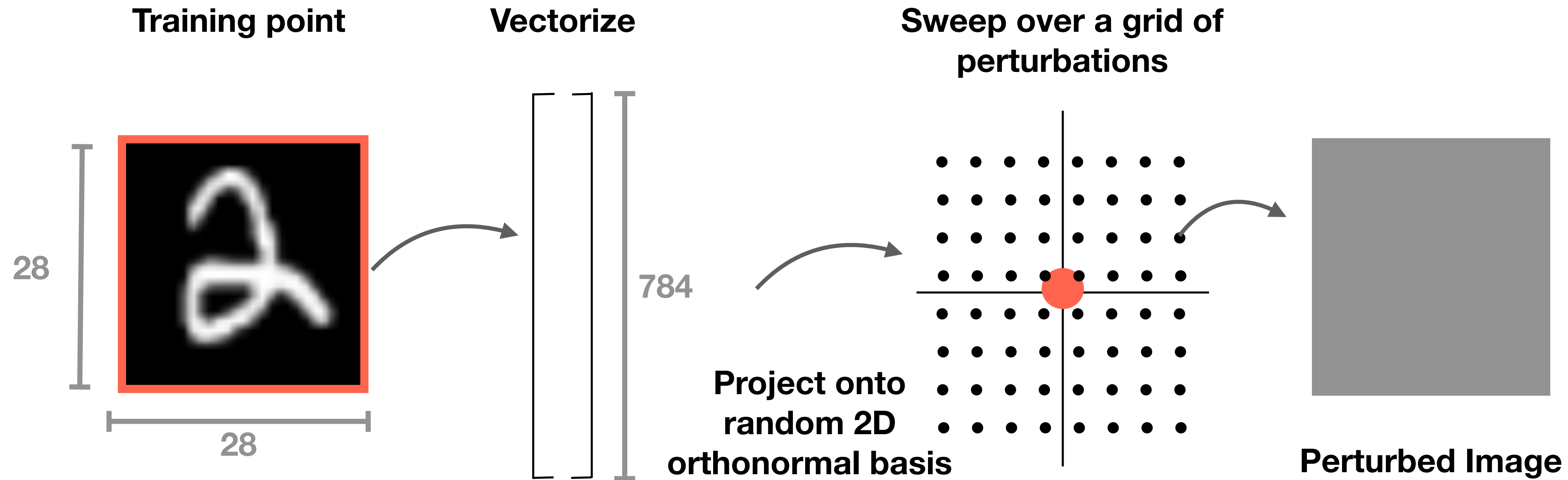
Visualize Perturbation Space



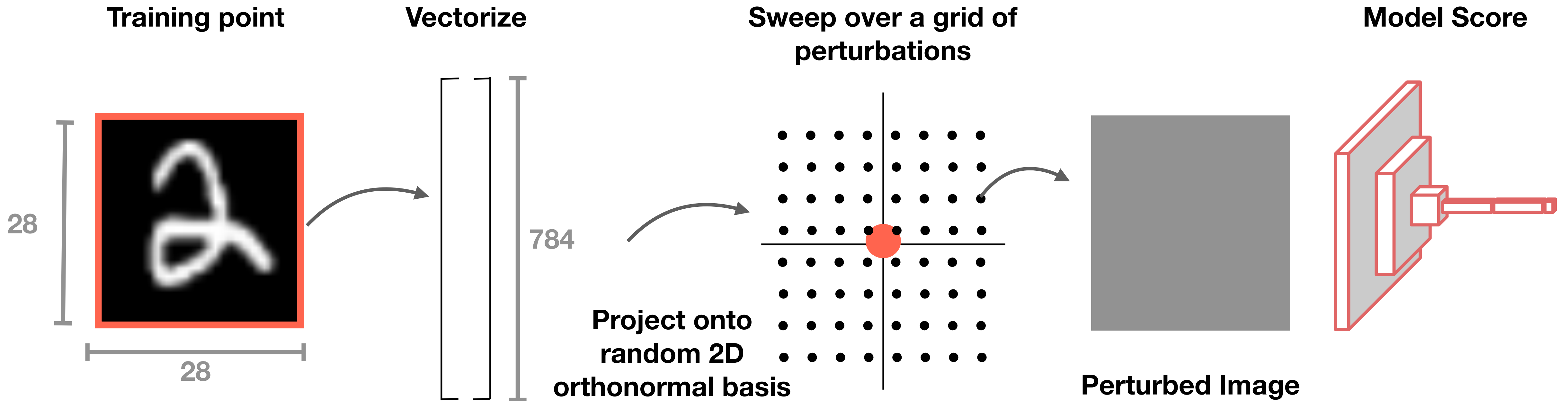
Visualize Perturbation Space



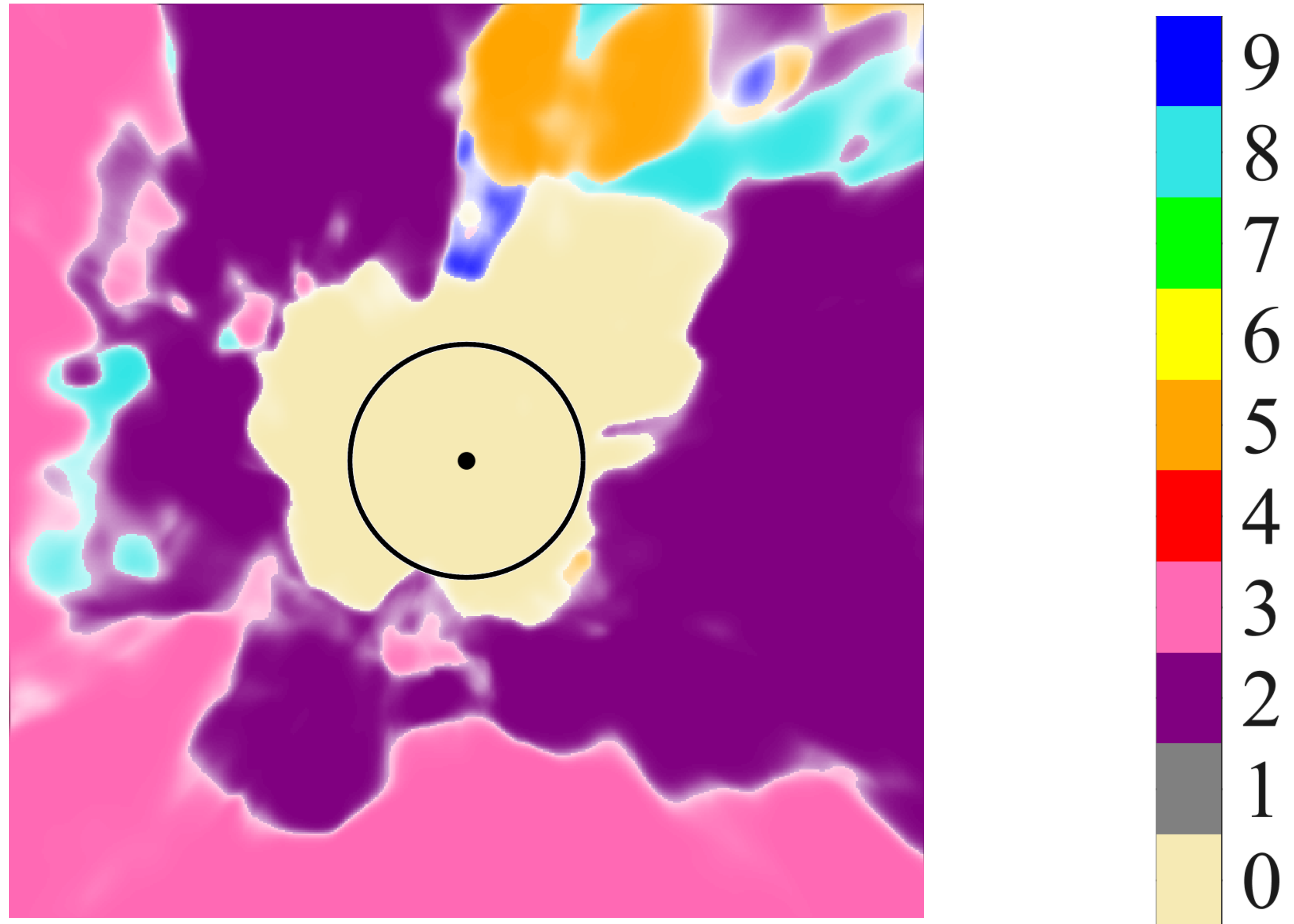
Visualize Perturbation Space



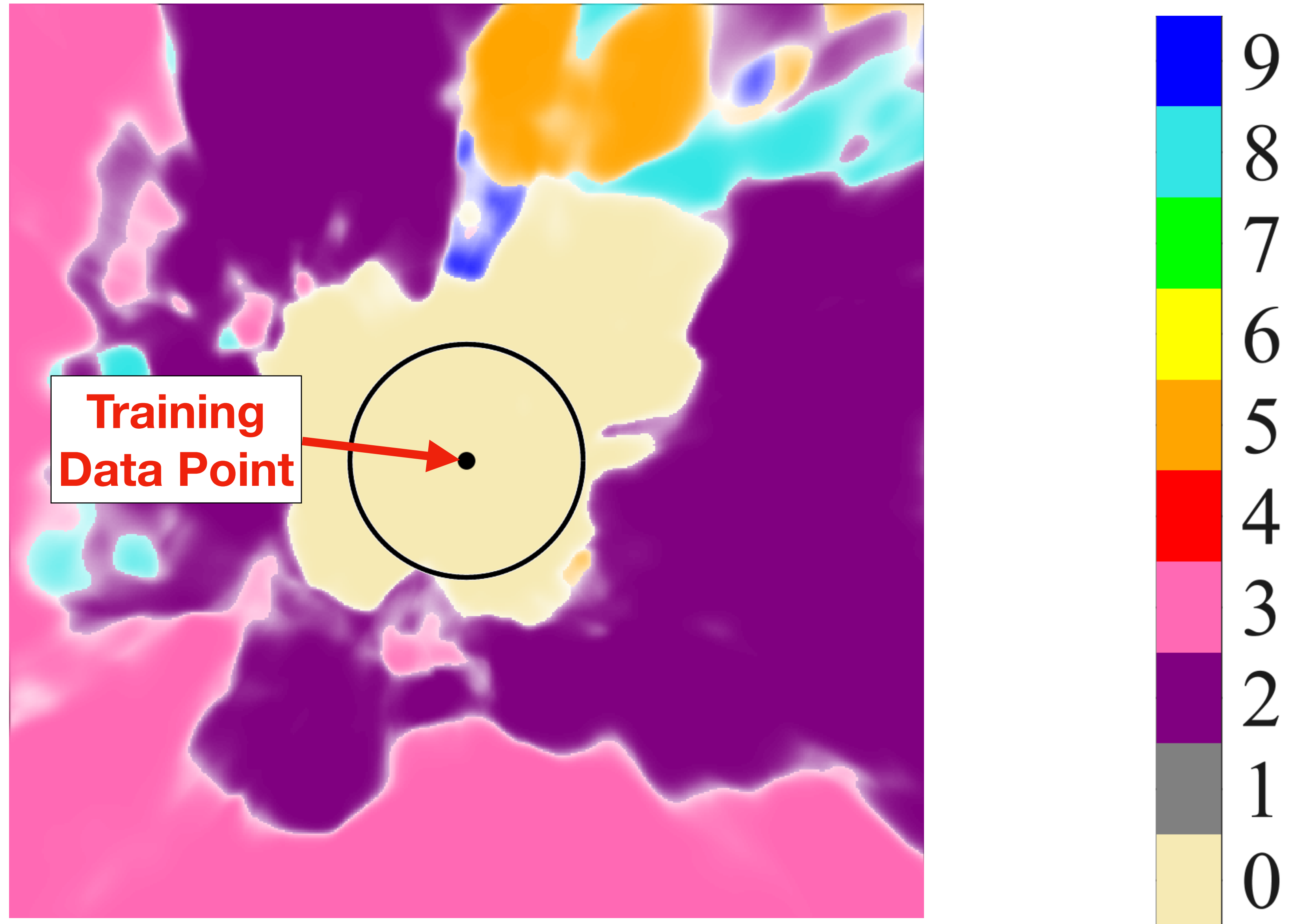
Visualize Perturbation Space



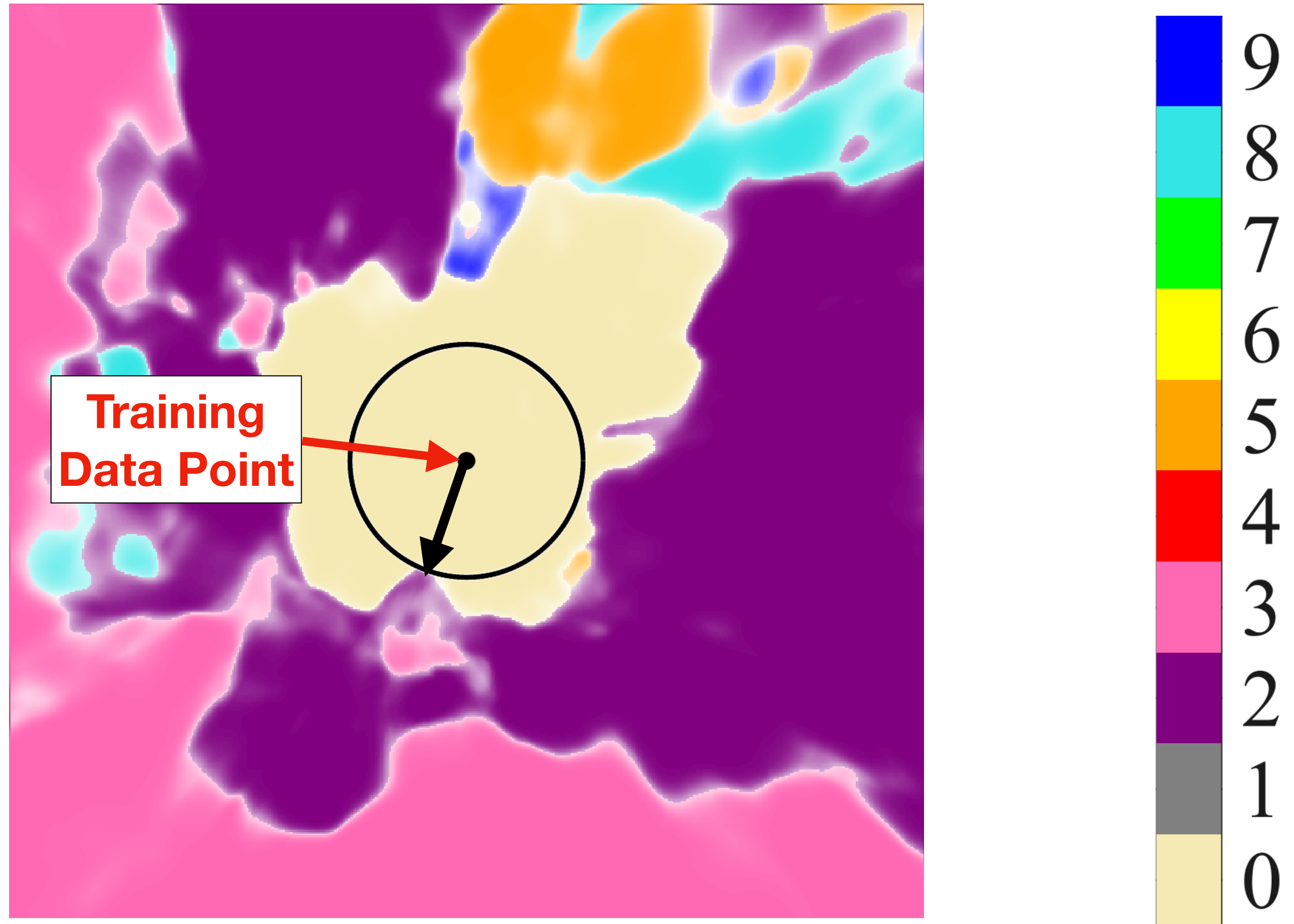
MNIST LeNet Decisions Around Training Point



MNIST LeNet Decisions Around Training Point

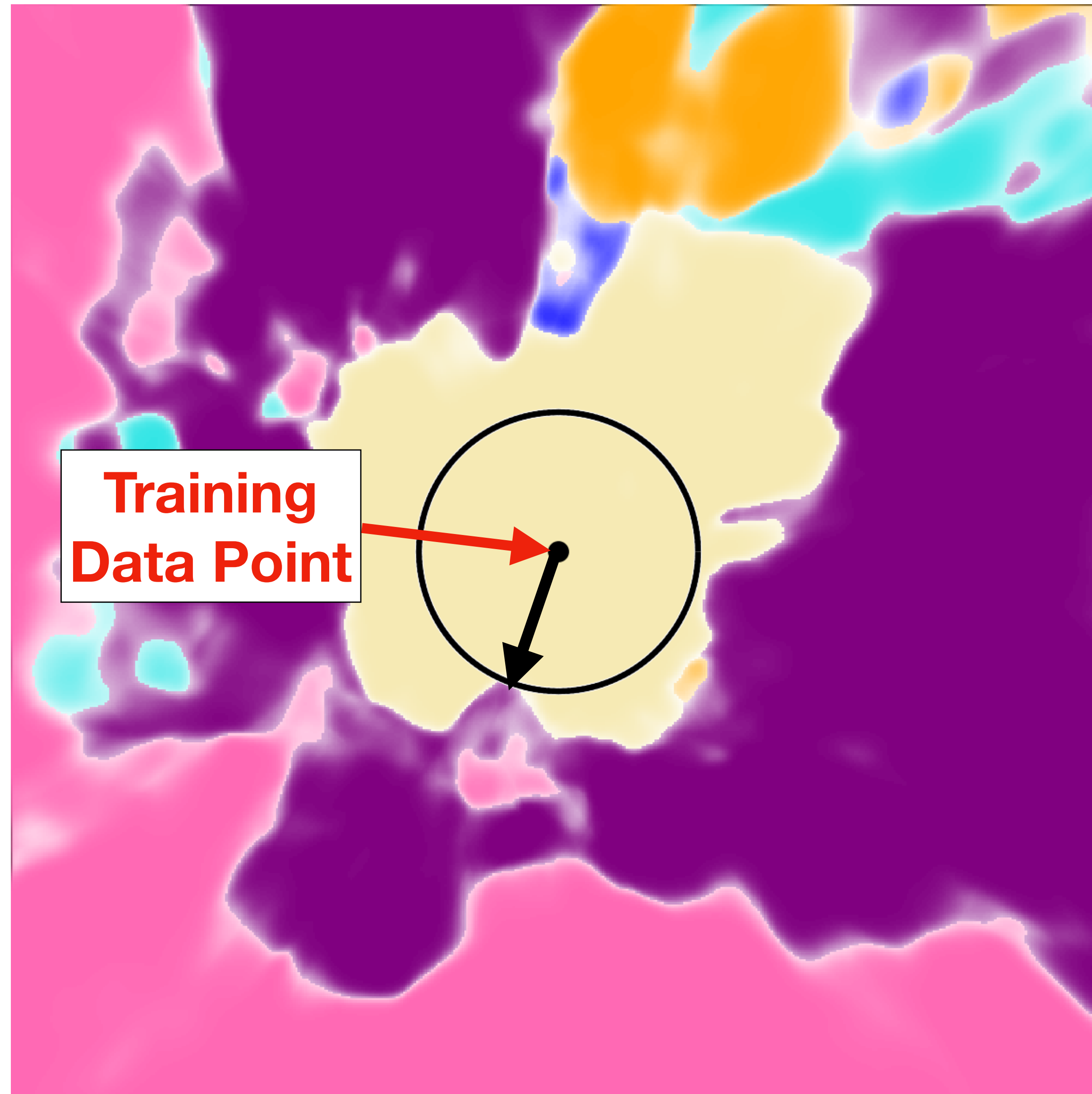


MNIST LeNet Decisions Around Training Point



MNIST LeNet Decisions Around Training Point

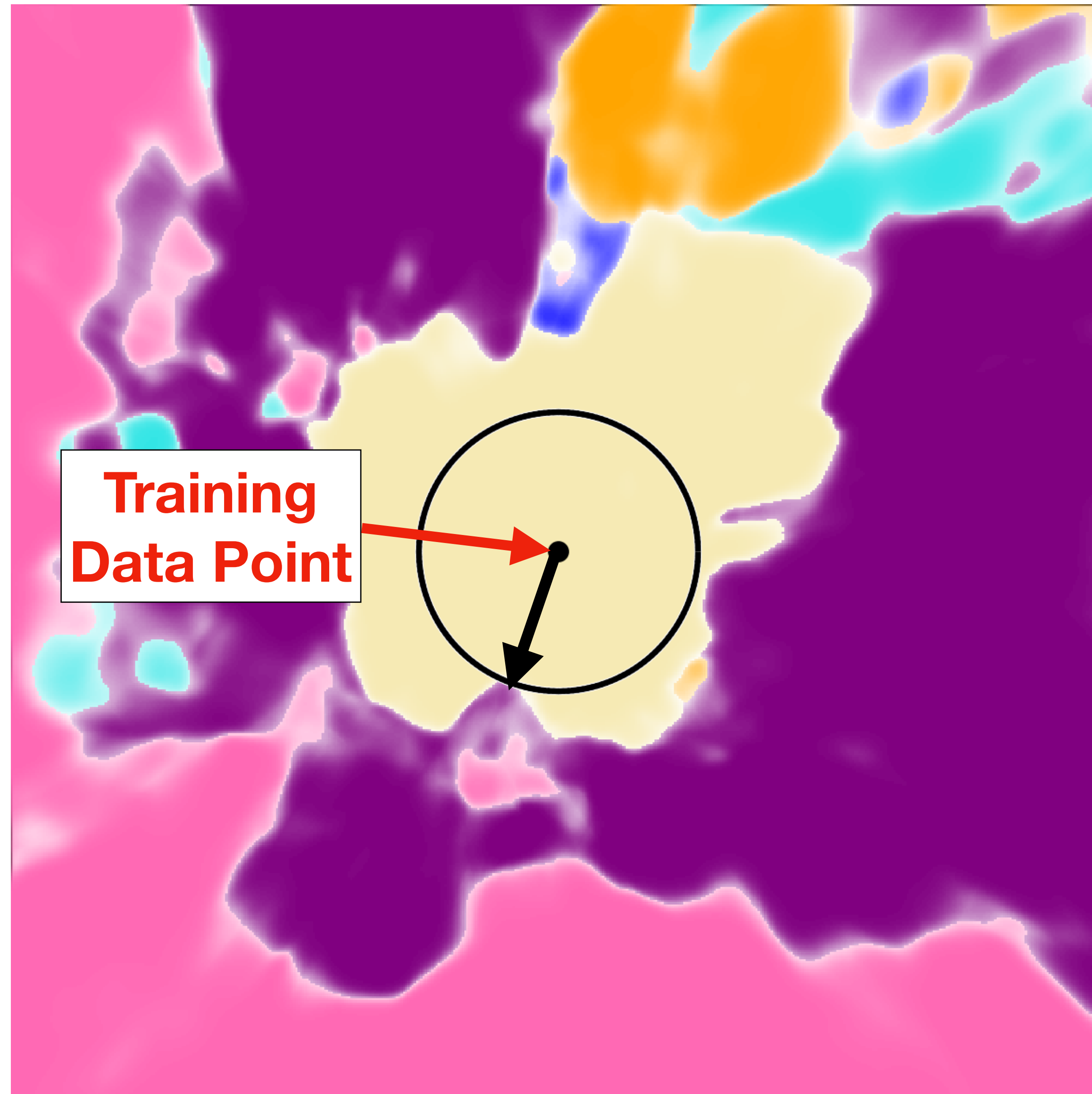
**Non-smooth
Decision Boundary**



MNIST LeNet Decisions Around Training Point

**Non-smooth
Decision Boundary**

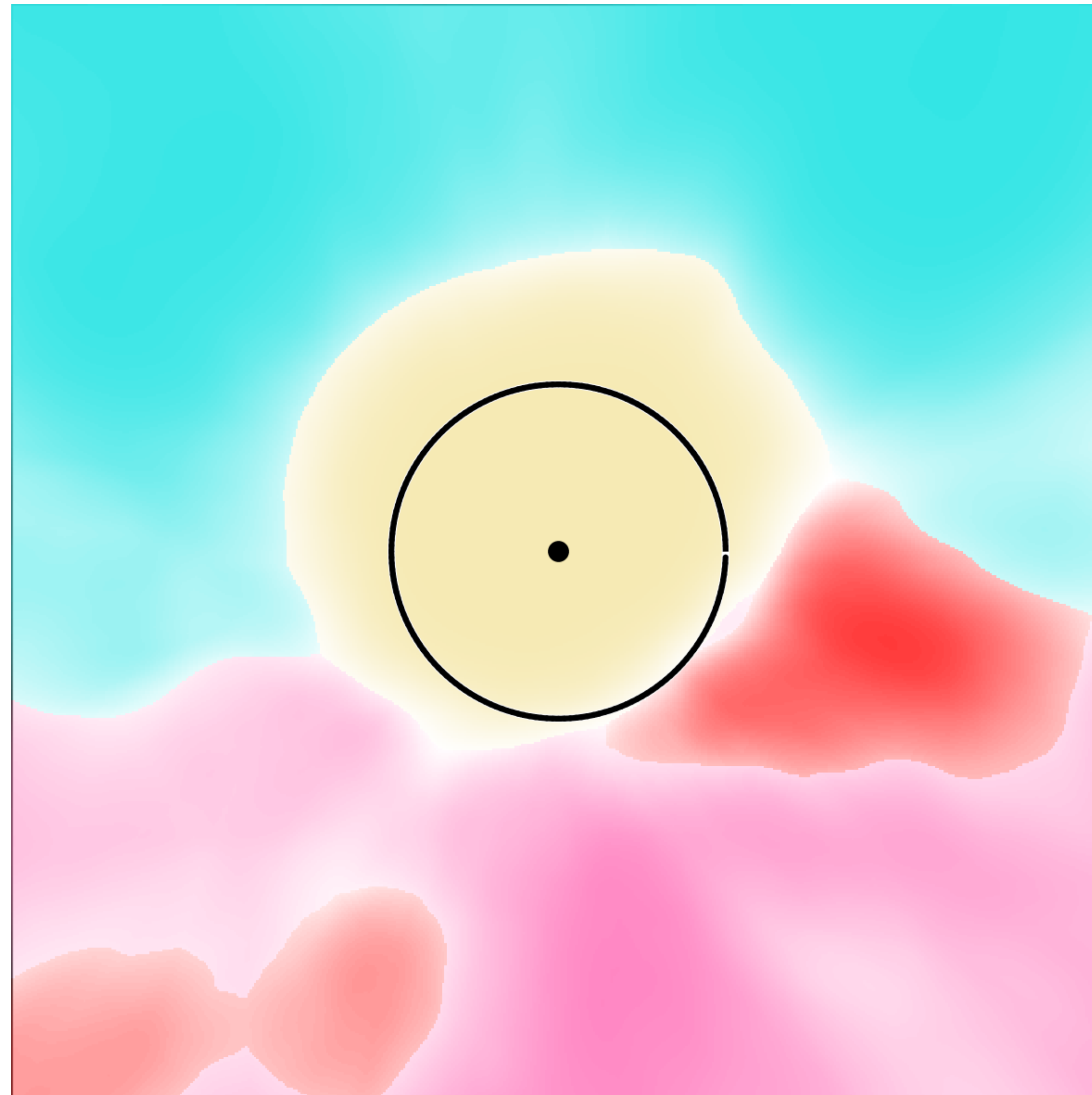
**Small perturbations
lead to new outputs**



MNIST LeNet with L2 Regularization

**Smooth Decision
Boundary**

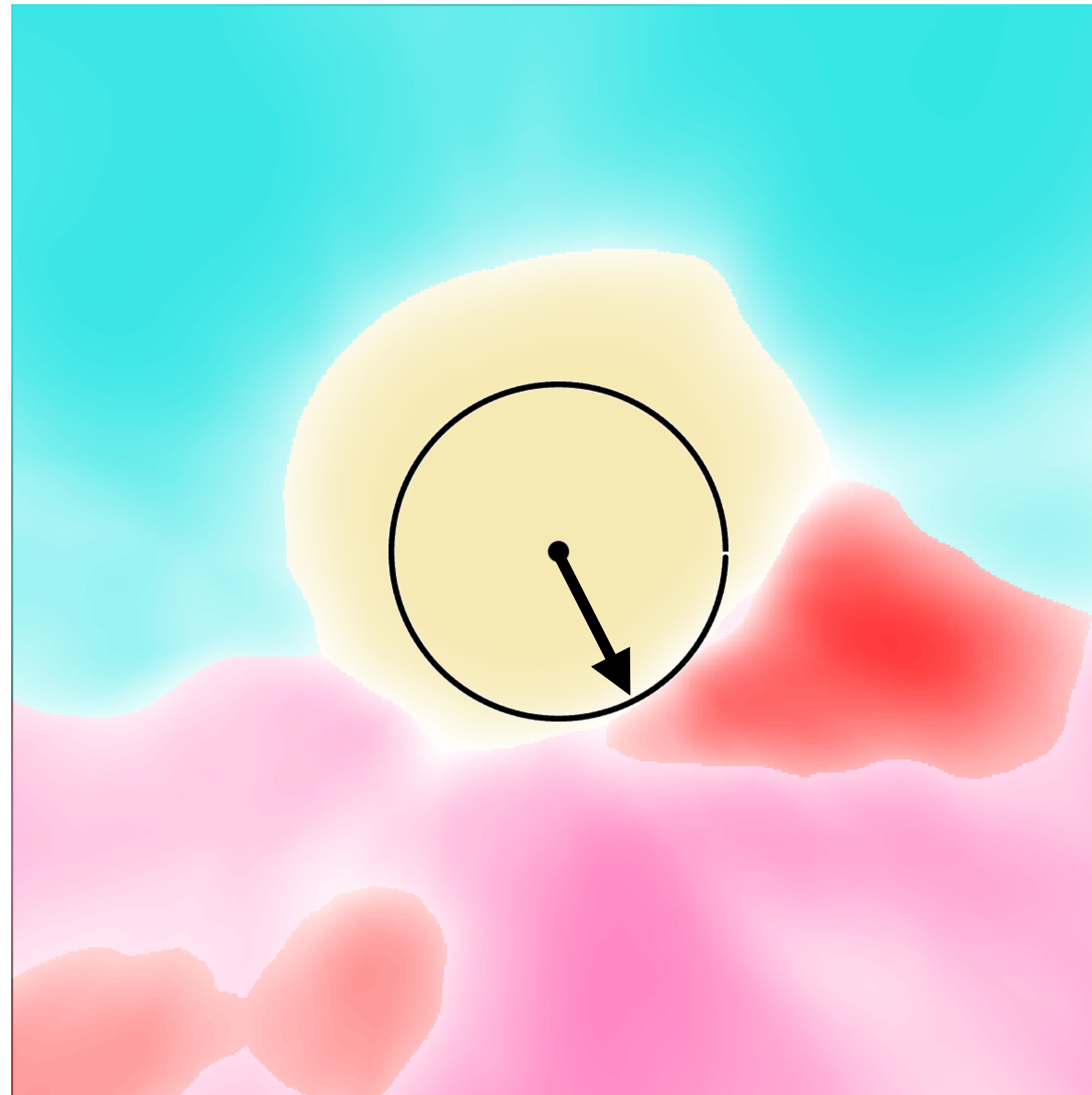
**Small perturbations
lead to new outputs**



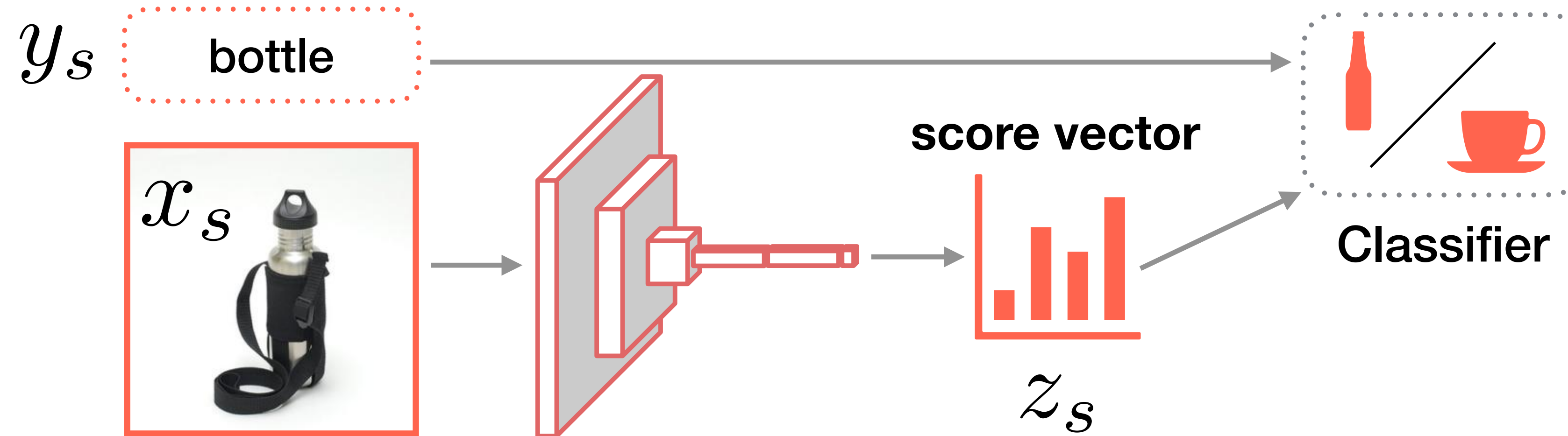
MNIST LeNet with L2 Regularization

Smooth Decision Boundary

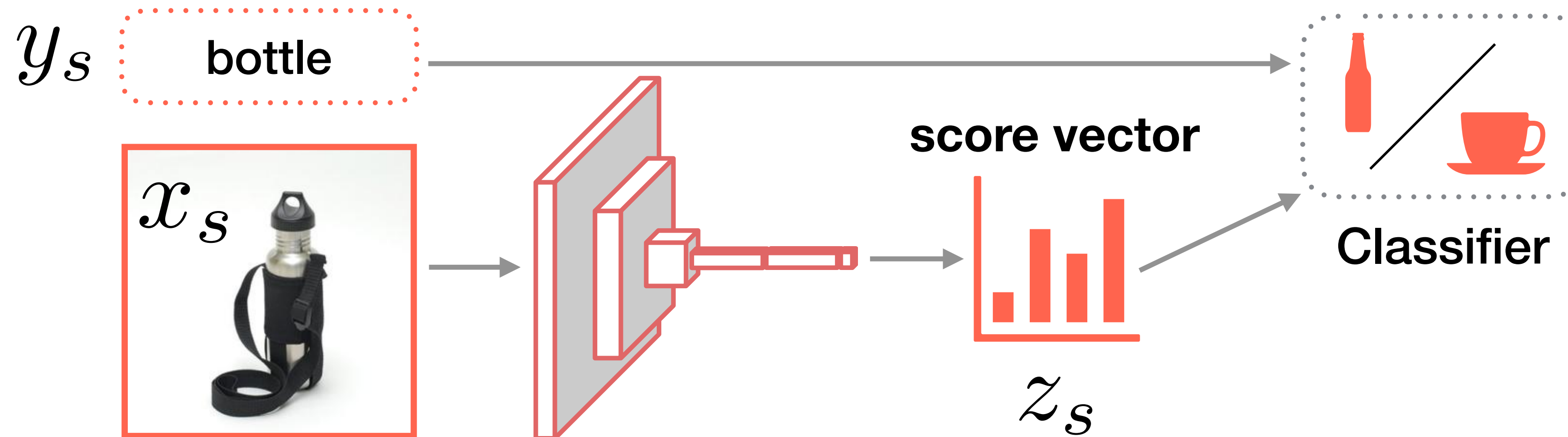
Small perturbations lead to new outputs



Jacobian Regularization



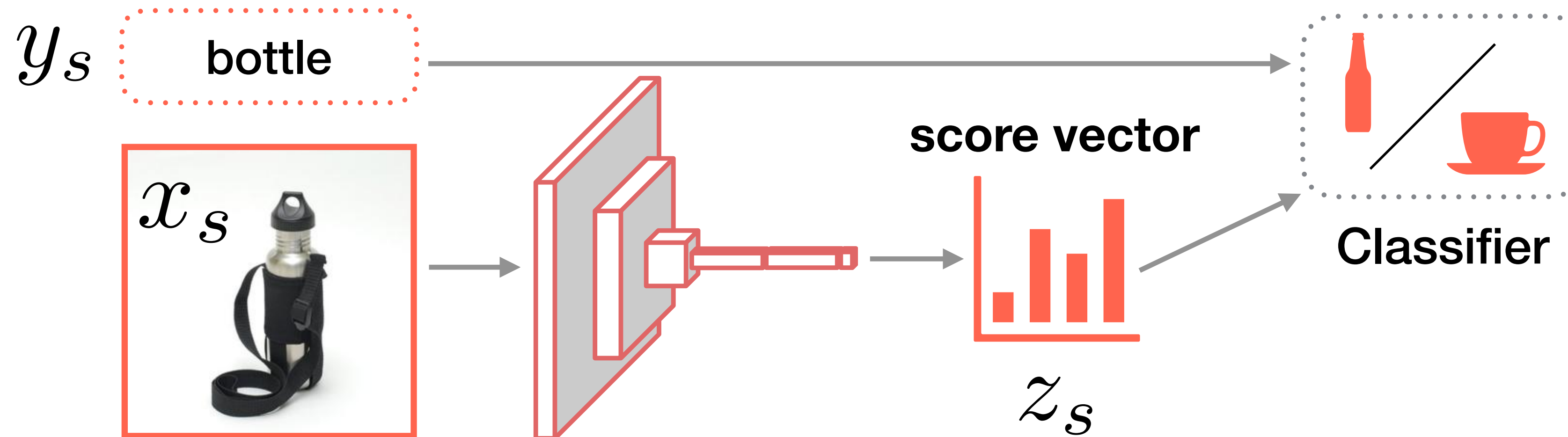
Jacobian Regularization



**Input-output
Jacobian matrix**

$$J_{c,i} = \frac{\partial z_c}{\partial x_i}$$

Jacobian Regularization



**Input-output
Jacobian matrix**

$$J_{c,i} = \frac{\partial z_c}{\partial x_i}$$

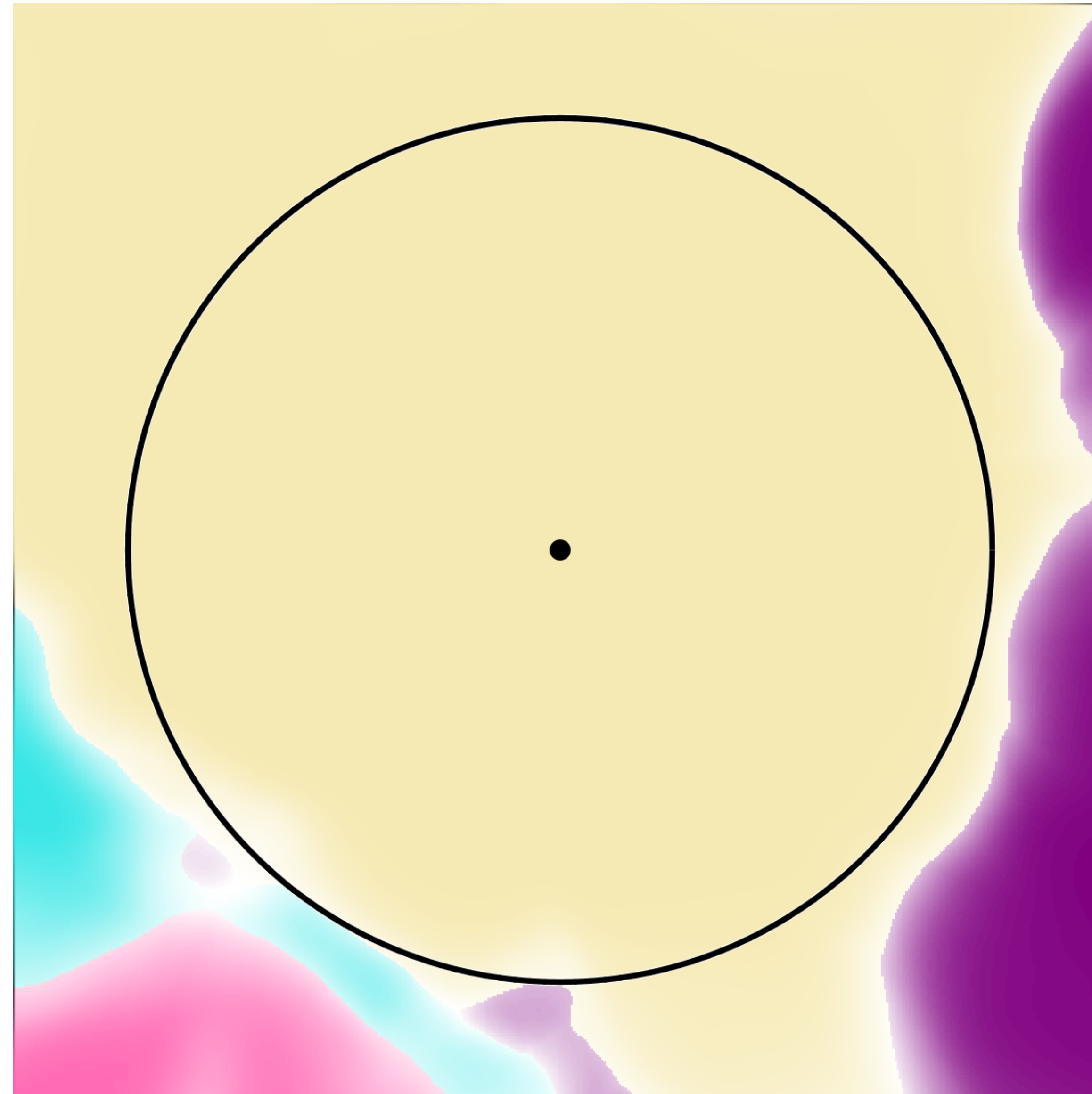
**Minimize
Frobenius Norm**

$$\|J\|_F^2$$

MNIST LeNet with Jacobian Regularization

**Mostly Smooth
Decision Boundary**

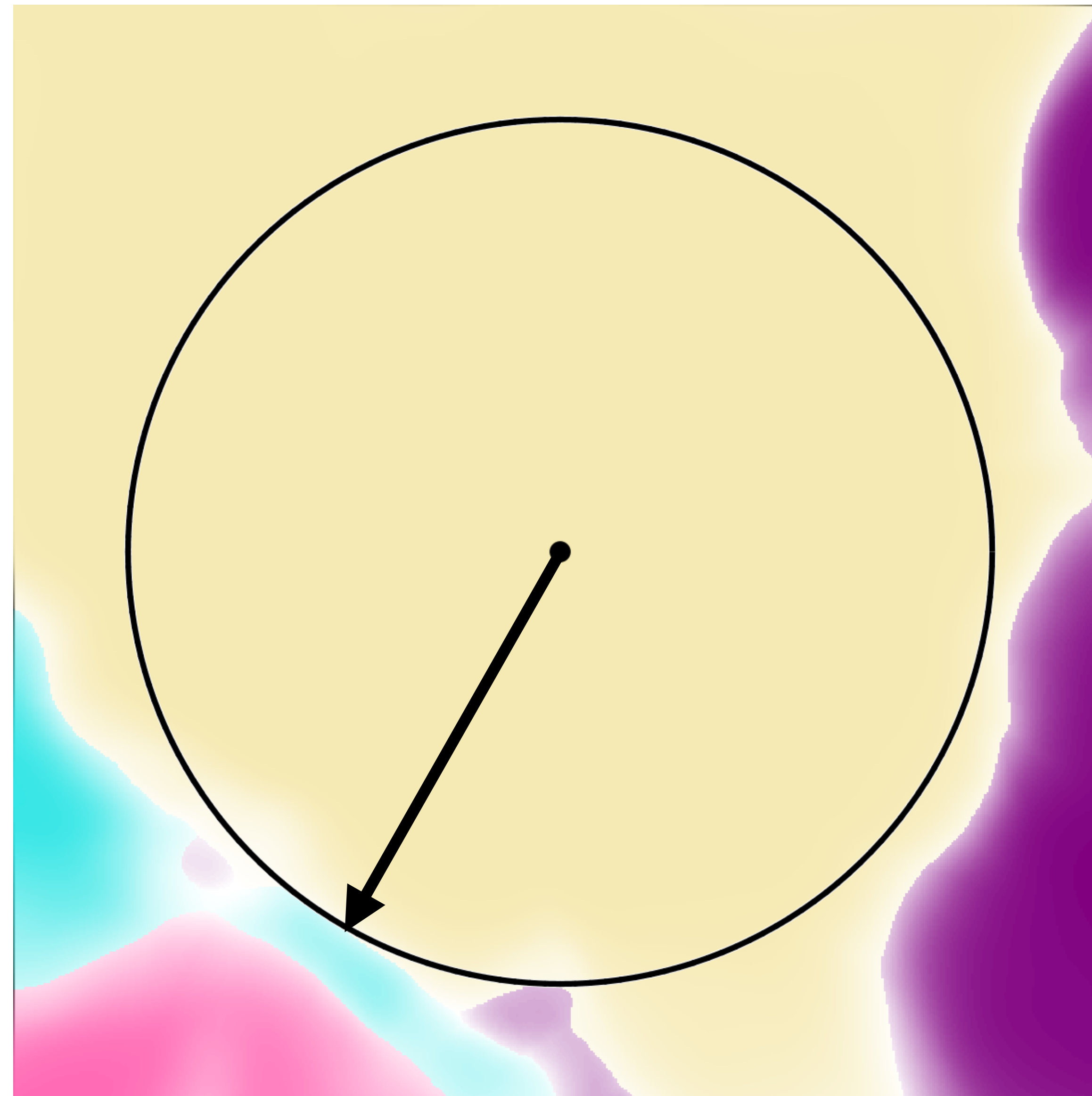
**Larger perturbations
needed to lead to
new outputs**



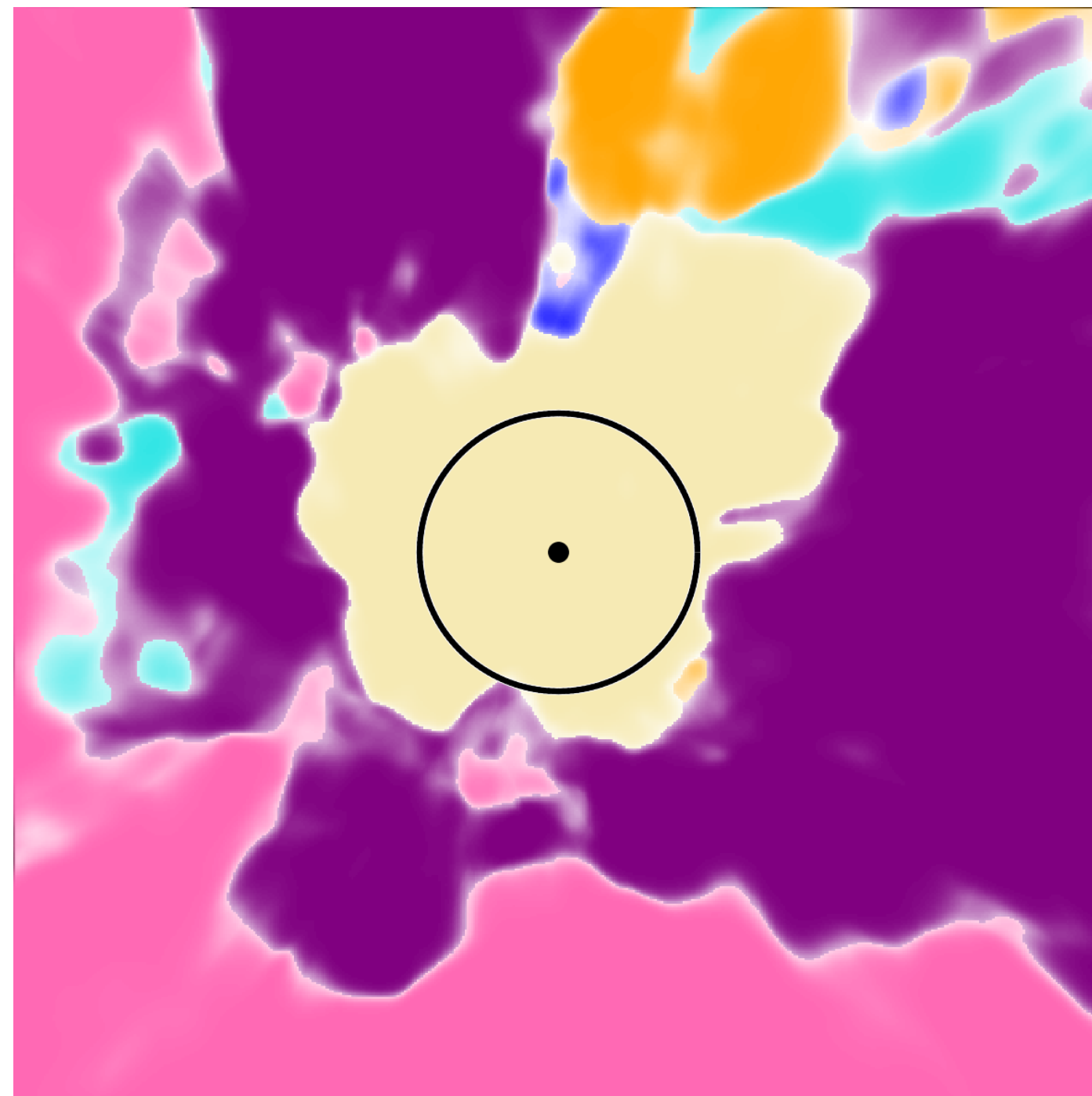
MNIST LeNet with Jacobian Regularization

**Mostly Smooth
Decision Boundary**

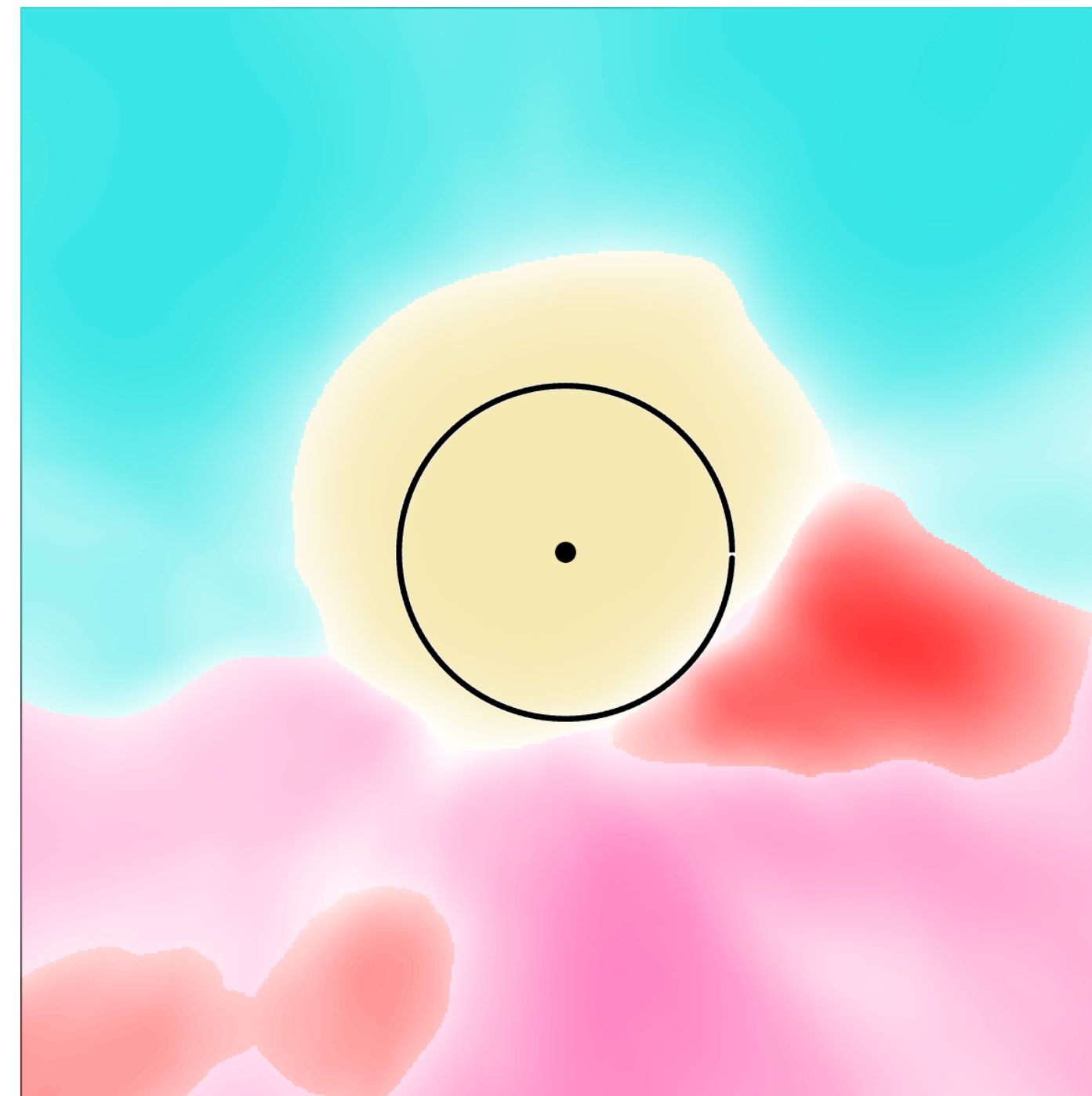
**Larger perturbations
needed to lead to
new outputs**



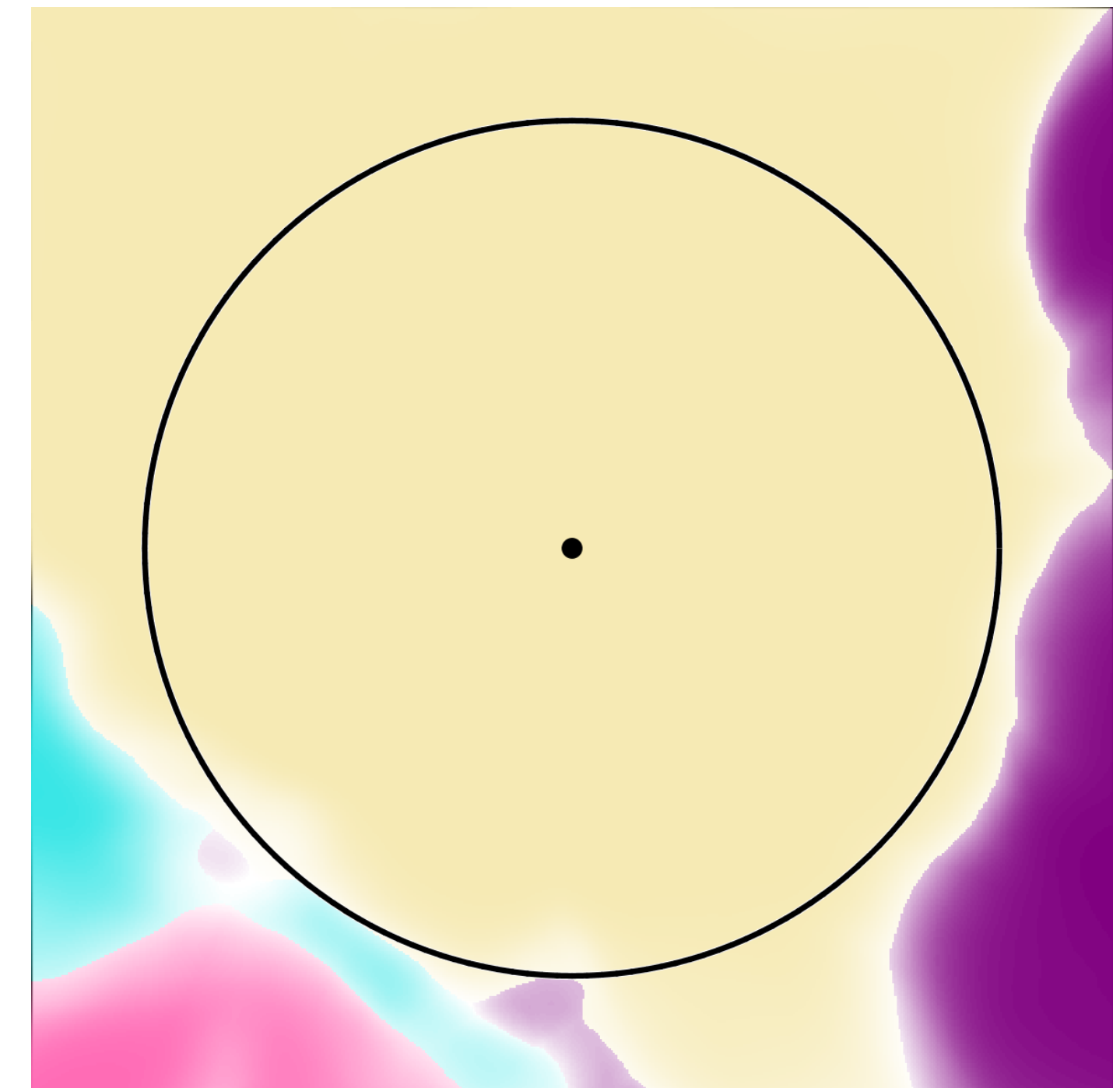
Decision Boundary Comparison



**No
Regularization**



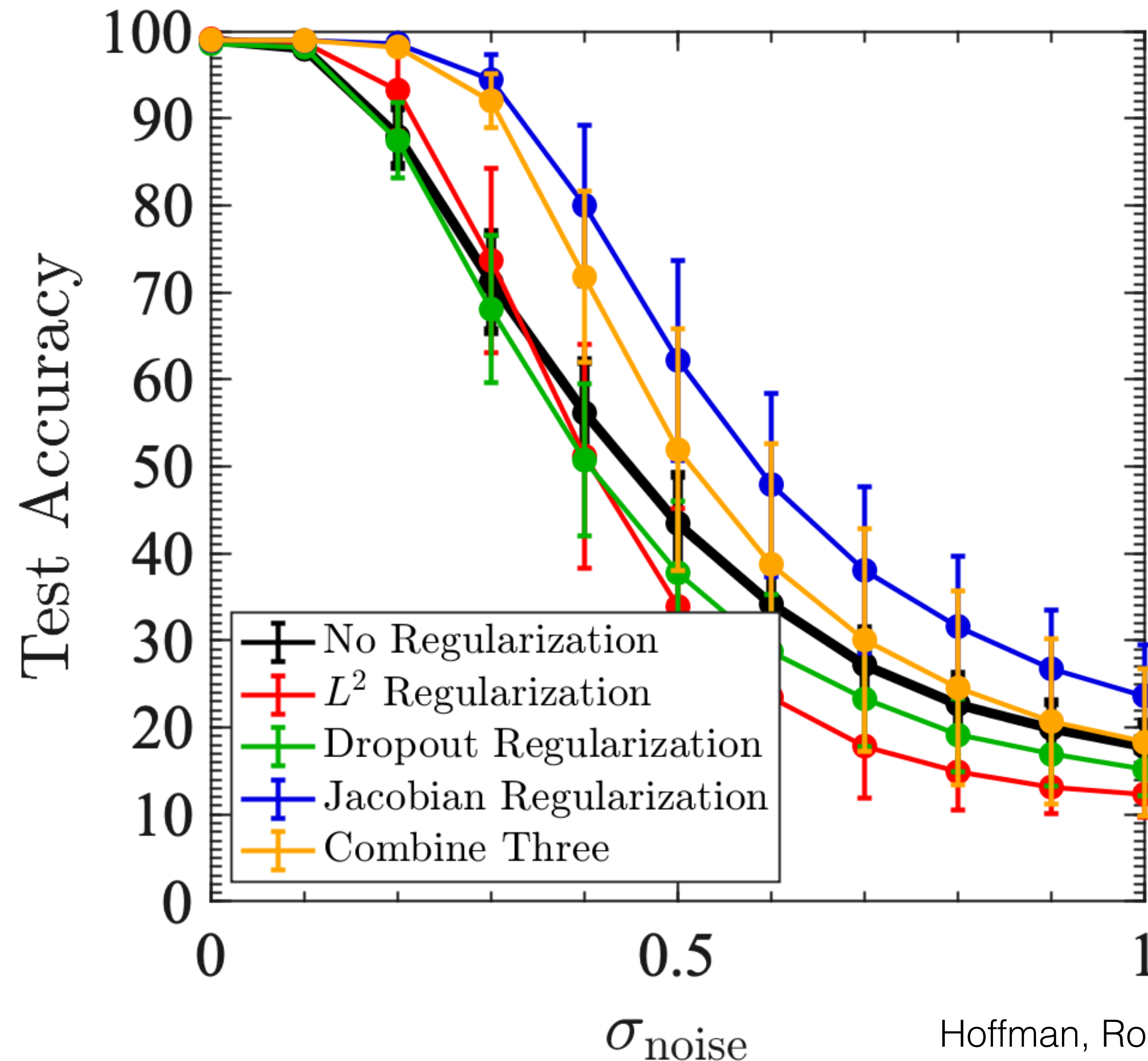
**L2
Regularization**



**Jacobian
Regularization**

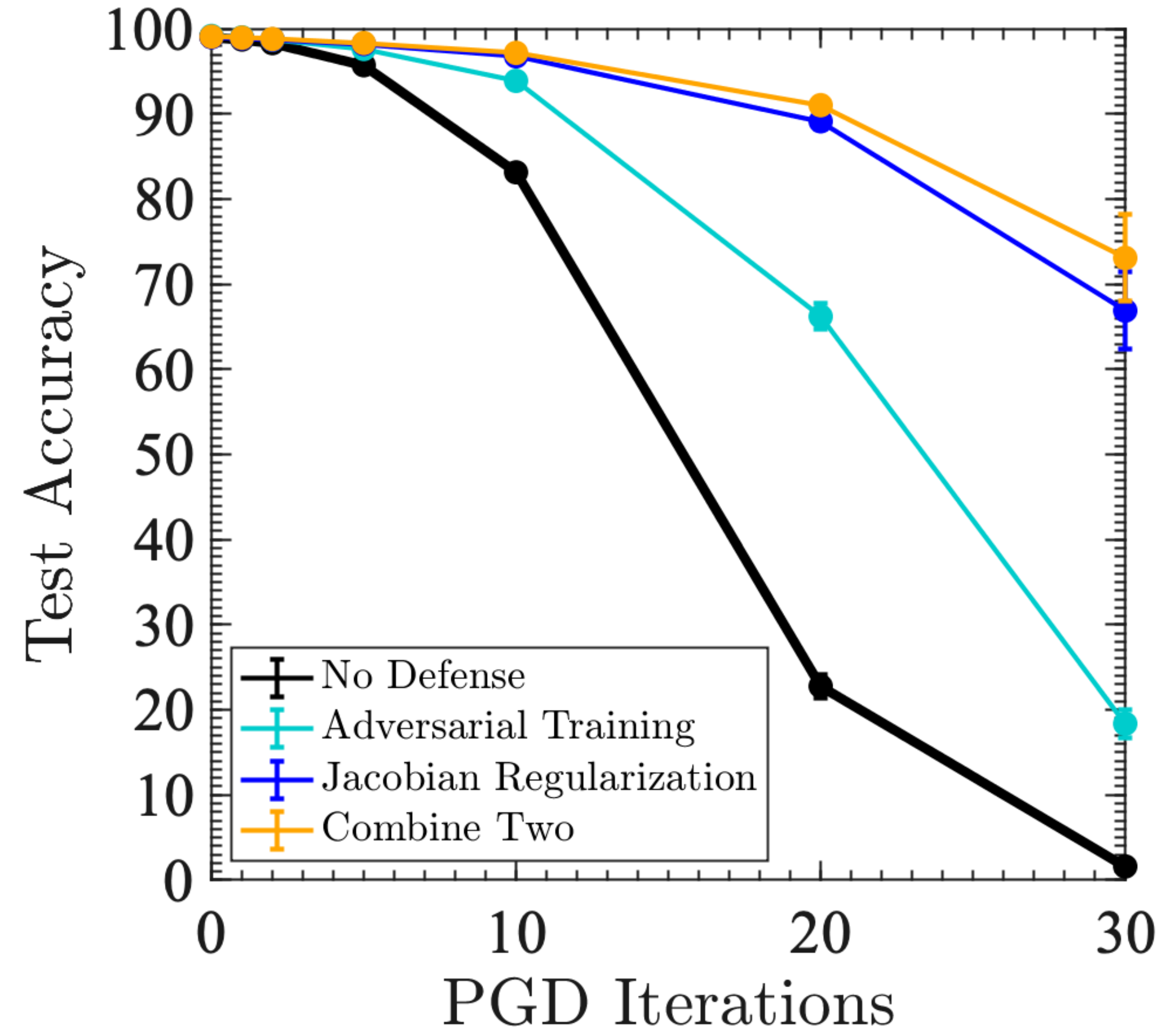
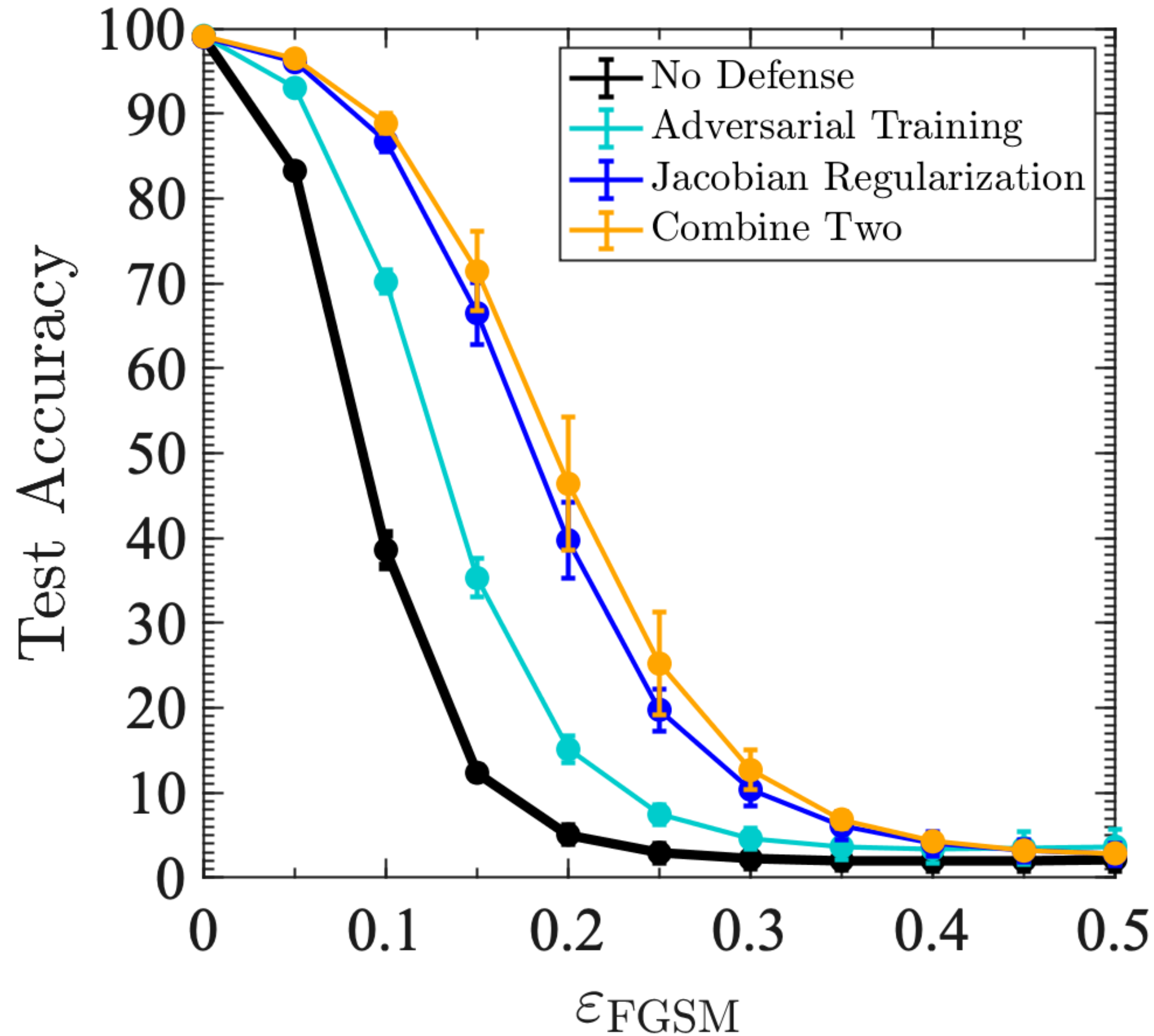


Robustness to Random Perturbations

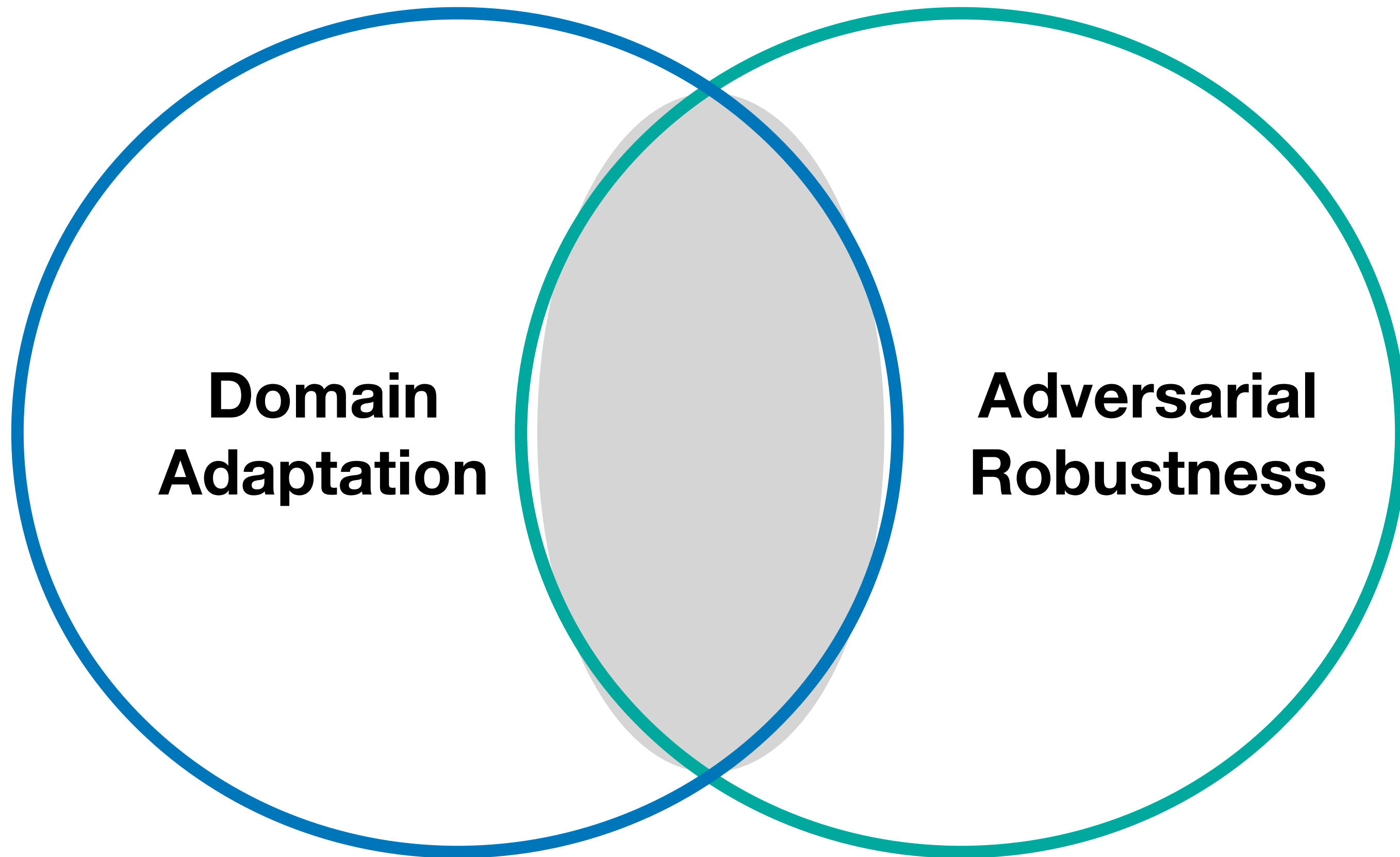


MNIST
LeNet Model

Robustness to Adversarial Perturbations



Next Steps



Jacobian regularizer as
unsupervised adaptive loss?

Adaptation to an adversarial
domain?

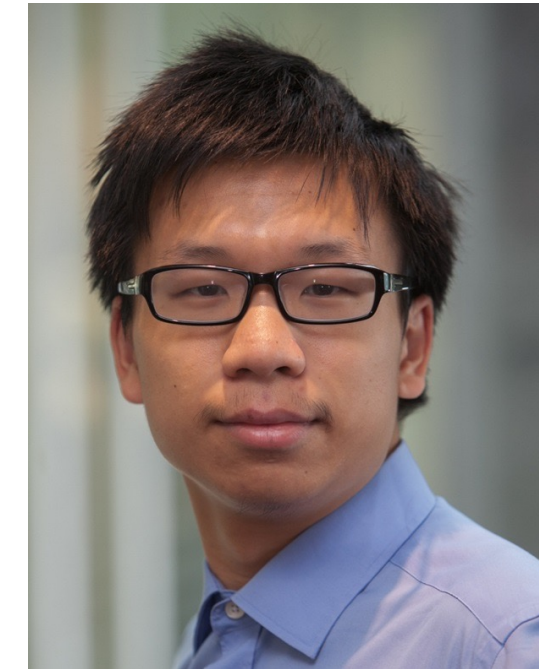
Thank you



Taesung Park
UC Berkeley



Eric Tzeng
UC Berkeley



Jun-Yan Zhu
MIT



Dan Roberts
Diffeo



Phil Isola
MIT



Kate Saenko
Boston University



Trevor Darrell
UC Berkeley



Alyosha Efros
UC Berkeley



Sho Yaida
FAIR



Judy Hoffman
judyhoffman.io