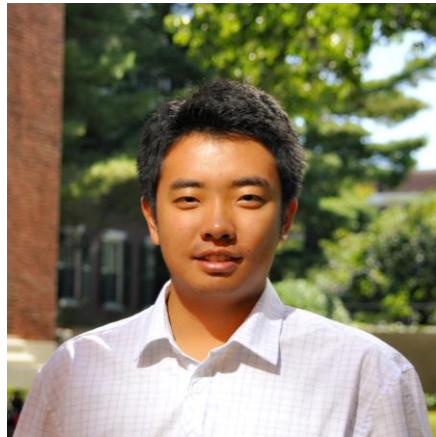




# Bayes-TrEx: A Bayesian Sampling Approach to Transparency by **Example**



Serena Booth\*



Yilun Zhou\*



Ankit Shah



Julie Shah

Consider a Corgi/Bread Classifier

# Consider a Corgi/Bread Classifier



# Consider a Corgi/Bread Classifier

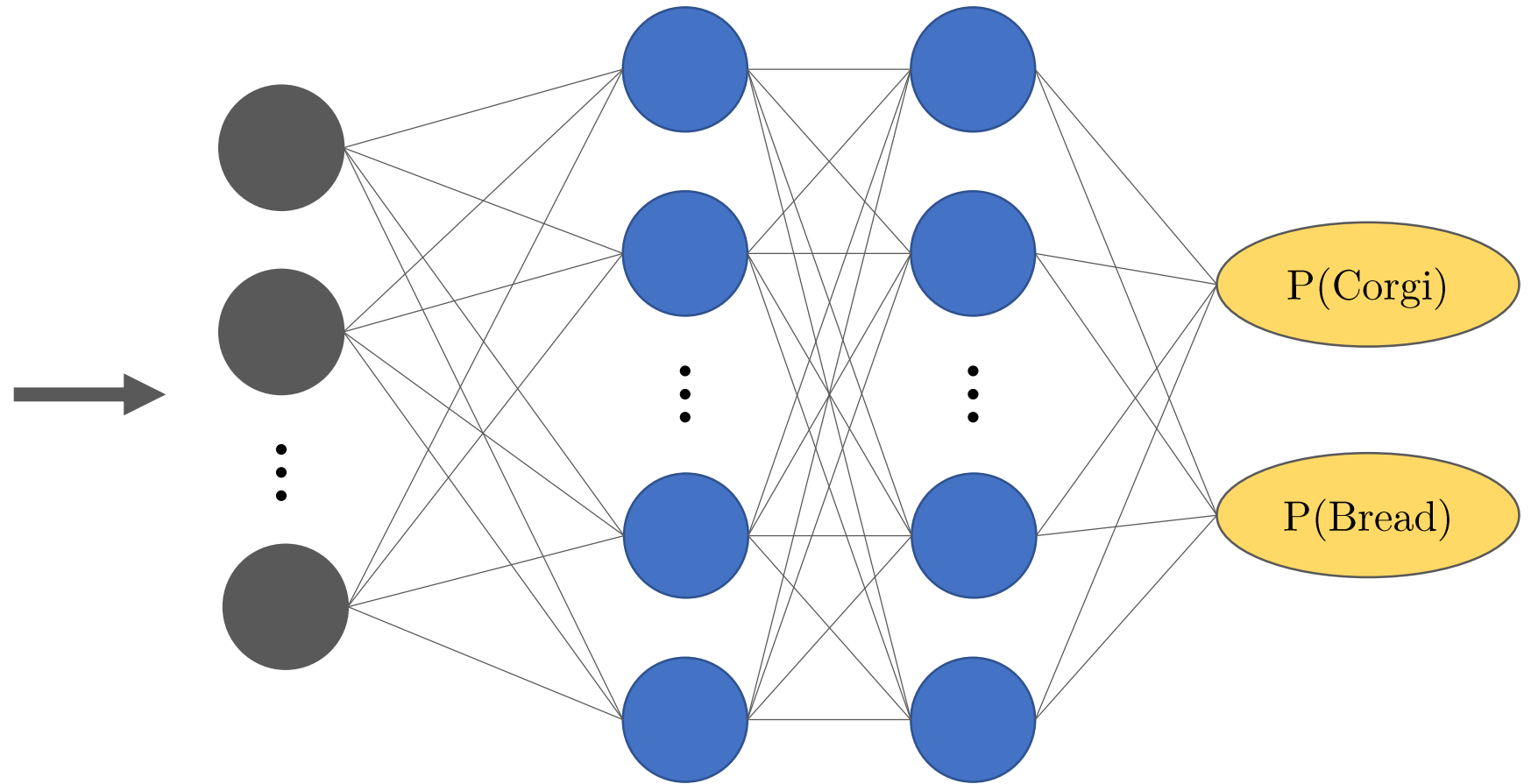


# Consider a Corgi/Bread Classifier



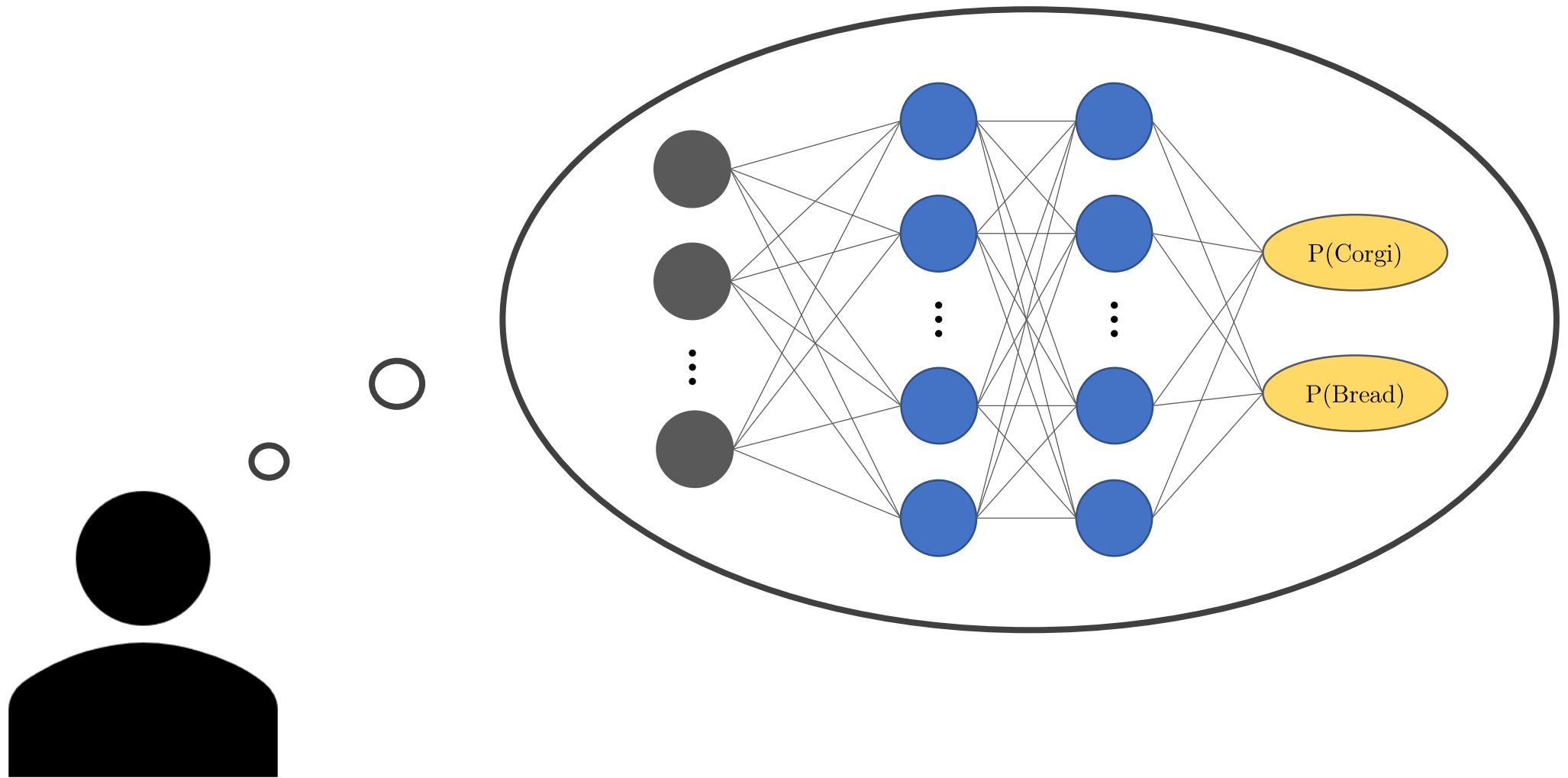


Input Image



Our objective:  
Build a *holistic* understanding of the classifier.

Build a *holistic* understanding of the classifier.



1. Develop an understanding of class boundaries.



# 1. Develop an understanding of class boundaries.

True class

Corgi



51% Corgi,  
49% Bread

Bread



48% Corgi,  
52% Bread

Corgi



50% Corgi,  
50% Bread

Bread



45% Corgi,  
55% Bread

Classifier's prediction

2. Predict classifier behavior on new examples.

## 2. Predict classifier behavior on new examples.

True class

Corgi



Bread



Corgi



Bread



## 2. Predict classifier behavior on new examples.

True class

Corgi



Bread



Corgi



Bread



Corgi

Bread

Bread

Corgi

Classifier's prediction

3. Predict classifier behavior on novel class instances.

### 3. Predict classifier behavior on novel class instances.

True class

Cat



Cake



Croissant



Potatoes



Corgi

Corgi

Bread

Bread

Classifier's prediction

How can we build this understanding?

Search the test set: ~50% Corgi, ~50% Bread

Test set? Many images.



Search the test set: ~50% Corgi, ~50% Bread



Test set? Many images.

Search the test set: ~50% Corgi, ~50% Bread



Test set? Many images.

50% Confident? Few images.

Search the test set: ~50% Corgi, ~50% Bread



Test set? Many images.

50% Confident? Few images.

Problem: sparse data

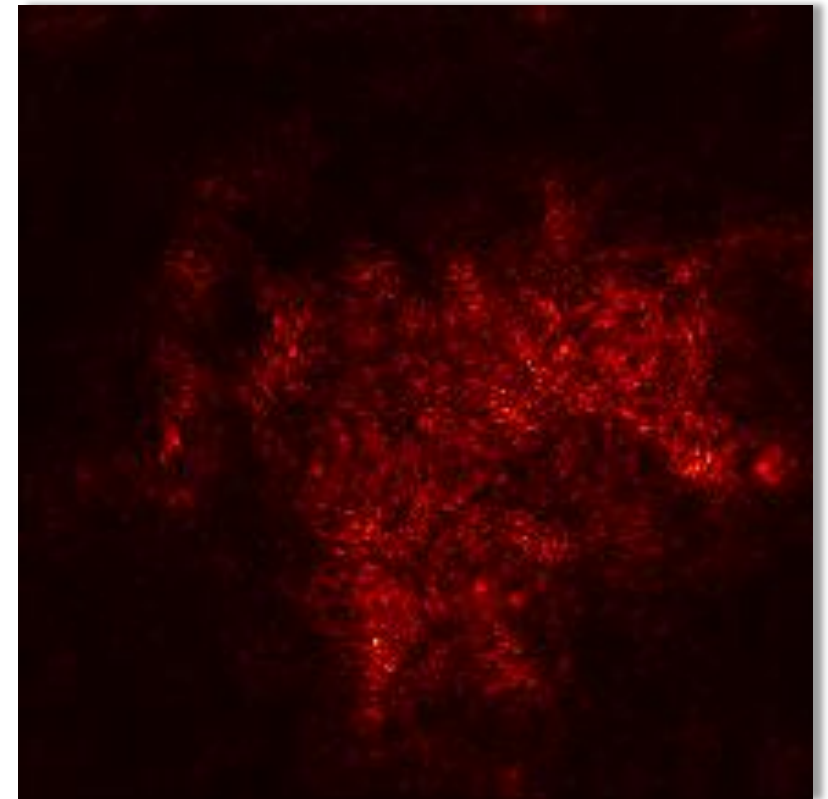
# How can we build a holistic understanding?



Input



Explanation  
pipeline



Saliency Map

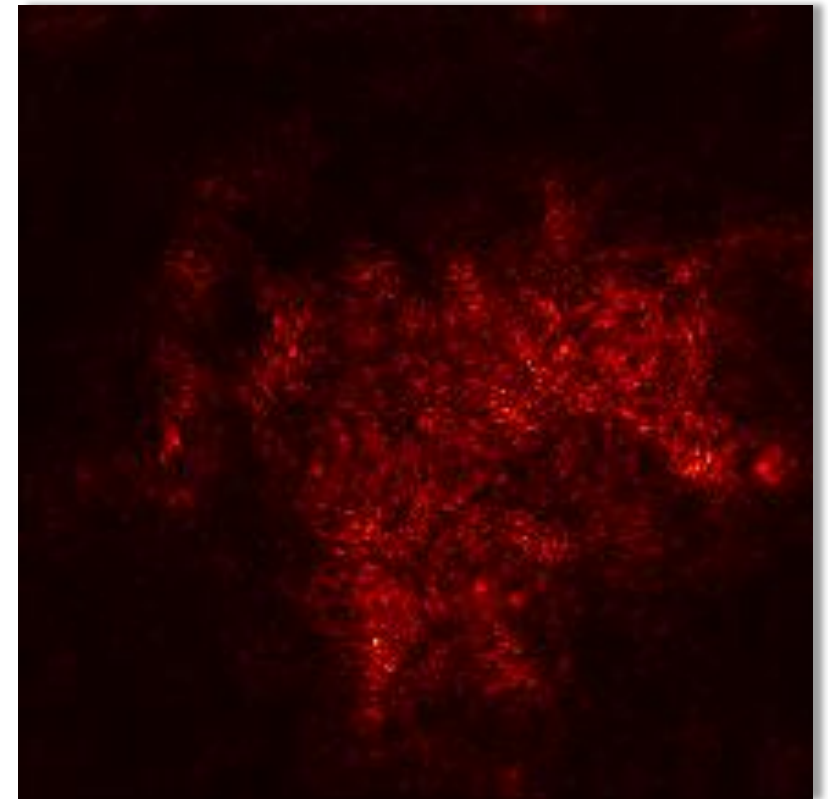
# How can we build a holistic understanding?



Input



Explanation  
pipeline

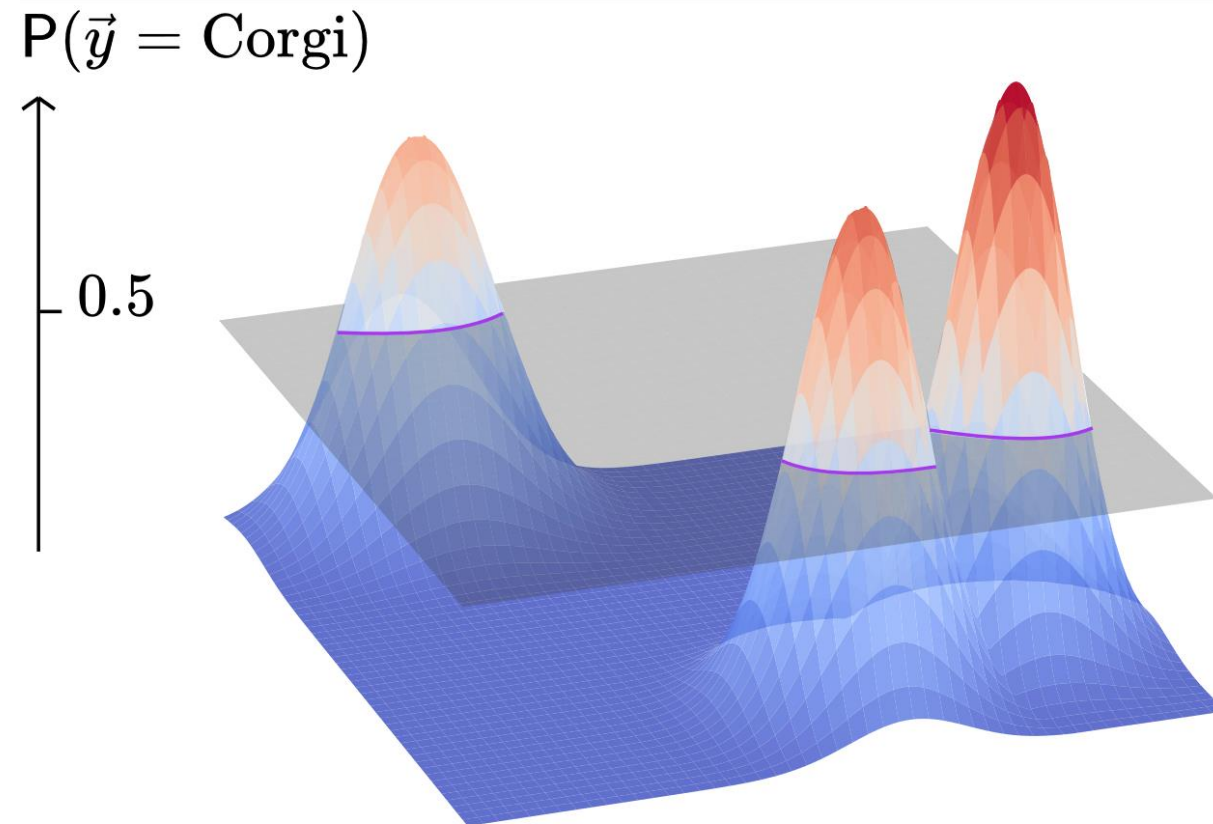


Saliency Map

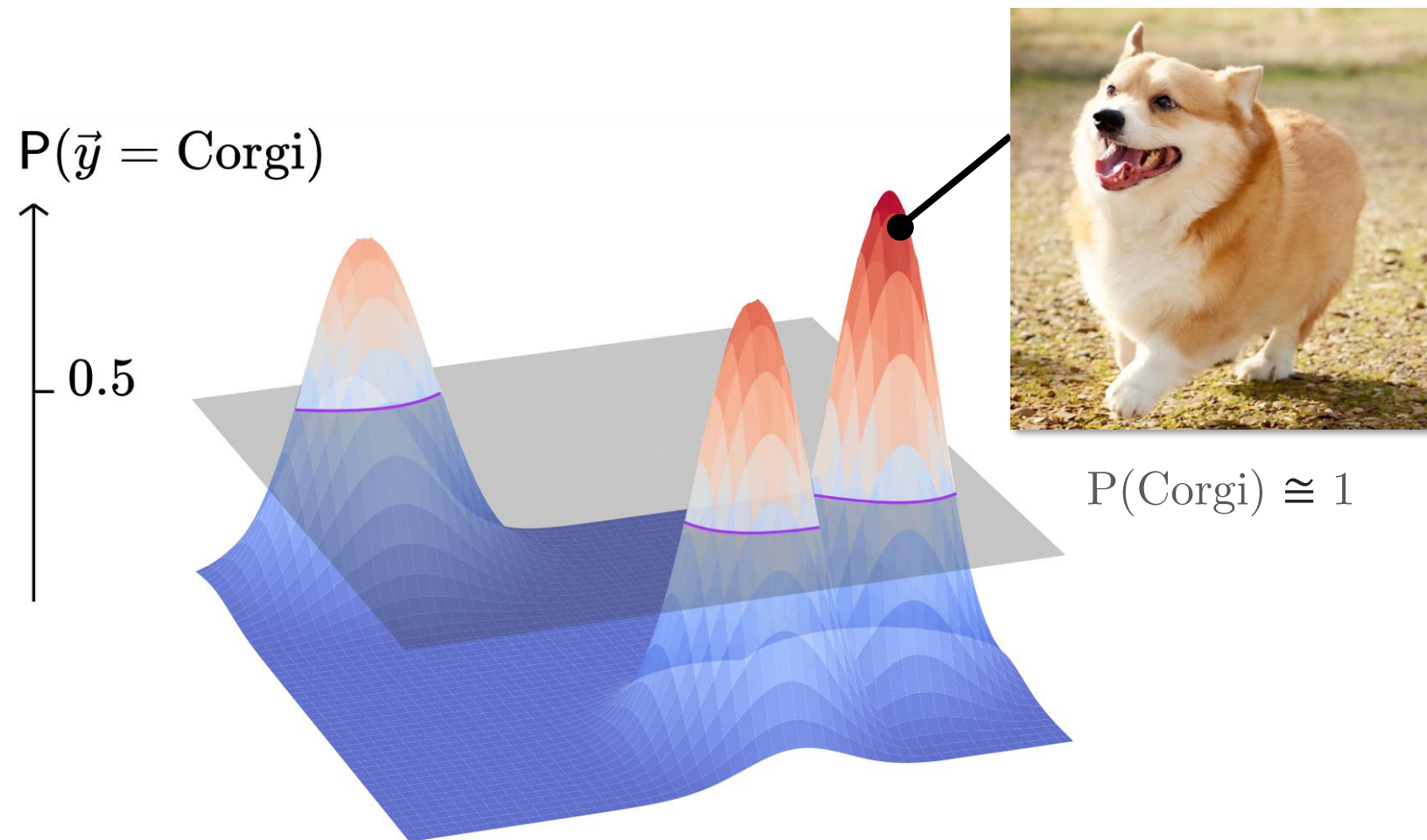
**Problem: where does the test input come from?**

We introduce Bayes-TrEx:  
an approach to transparency by example

# We introduce Bayes-TrEx: an approach to transparency by example

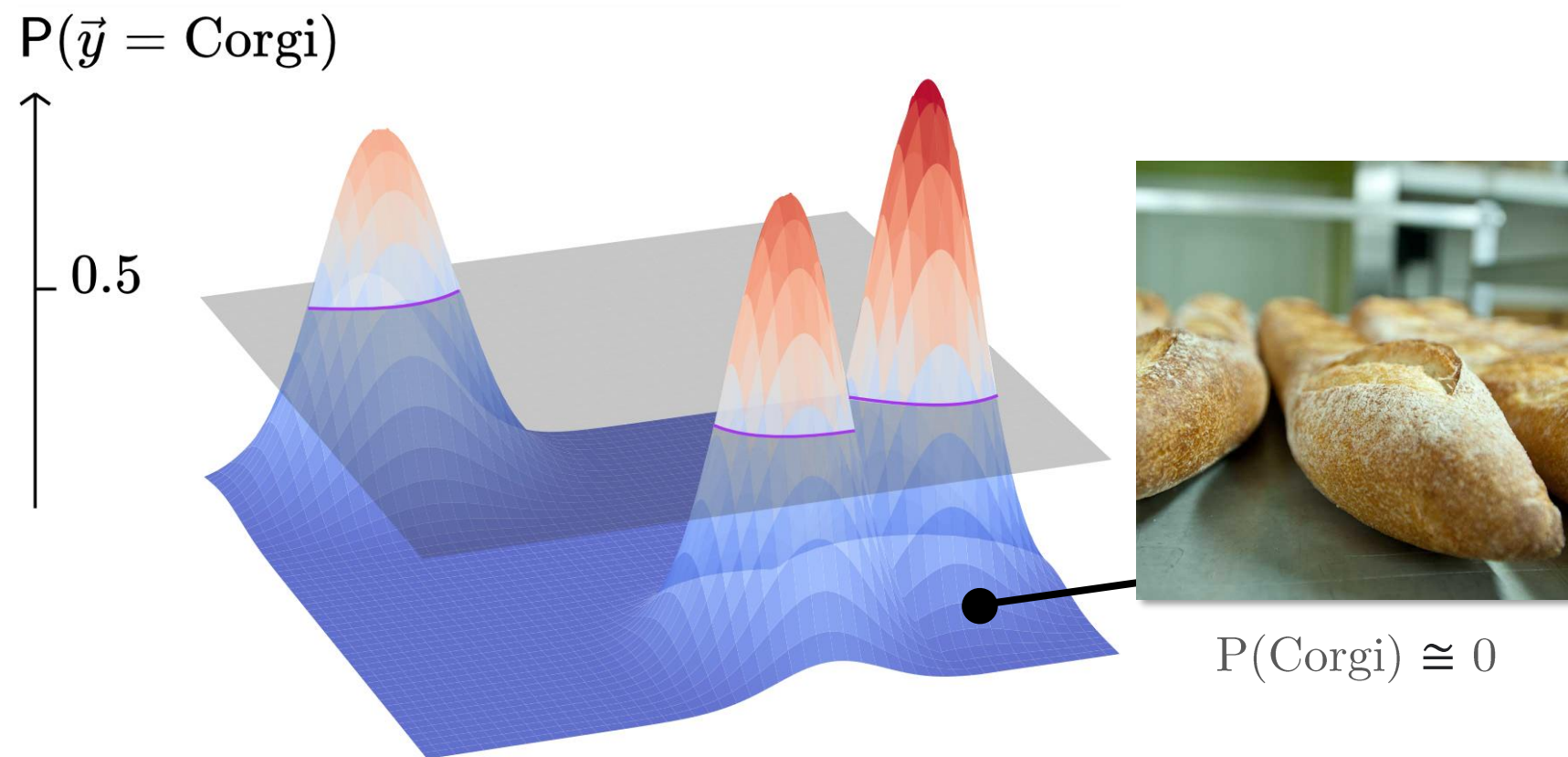


# A Corgi/Bread Decision Surface

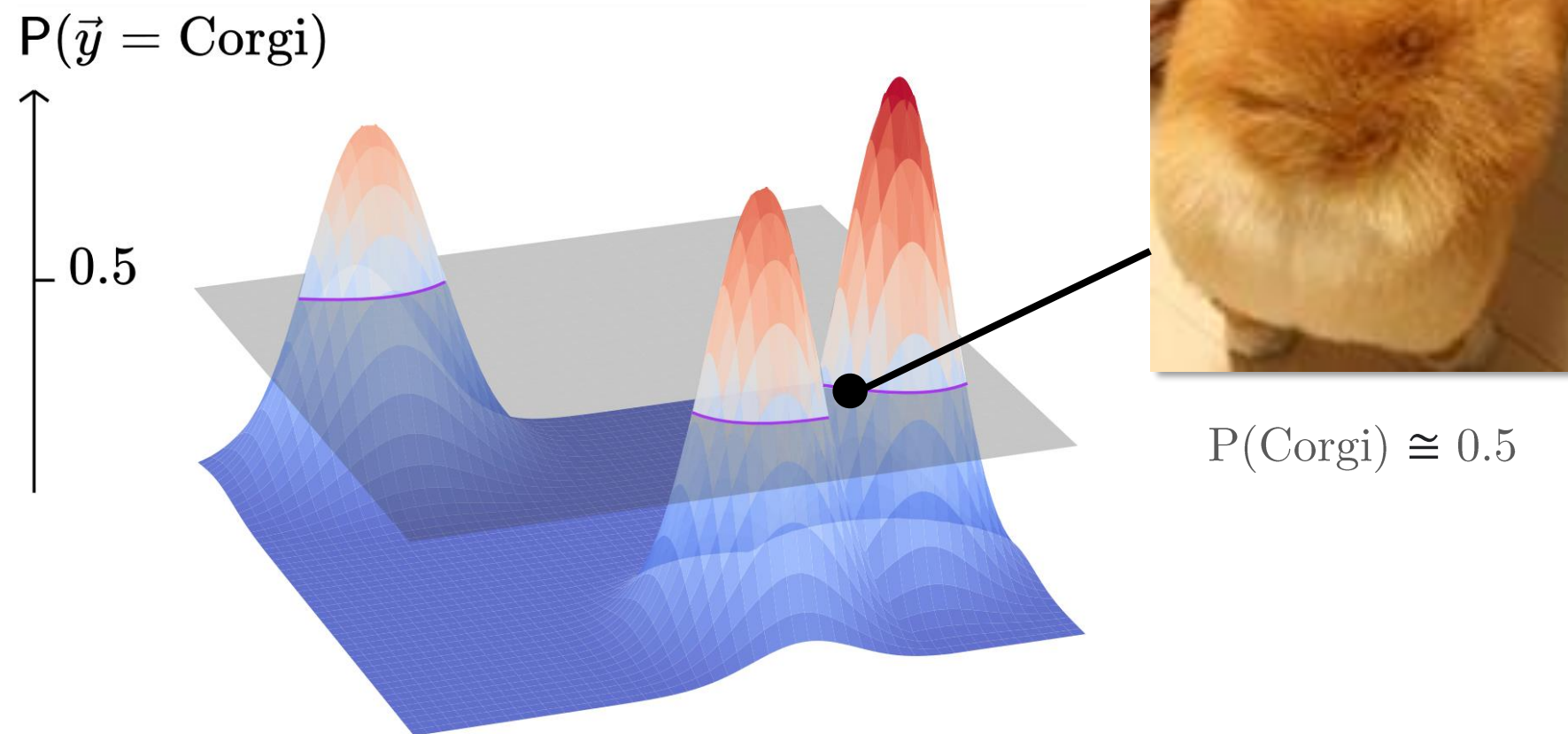




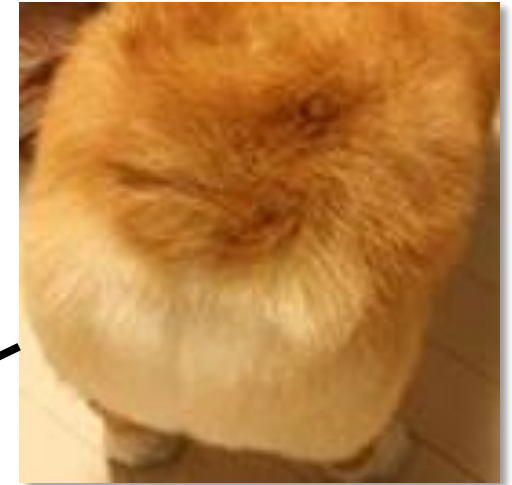
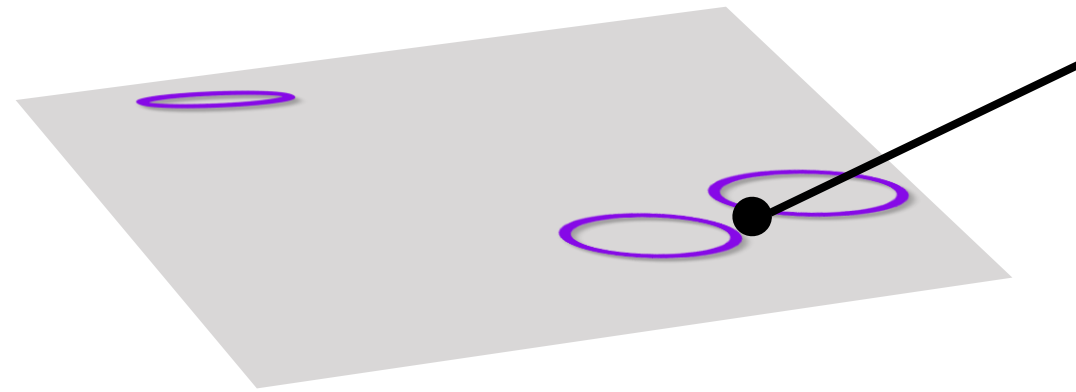
# A Corgi/Bread Decision Surface



# A Corgi/Bread Decision Surface



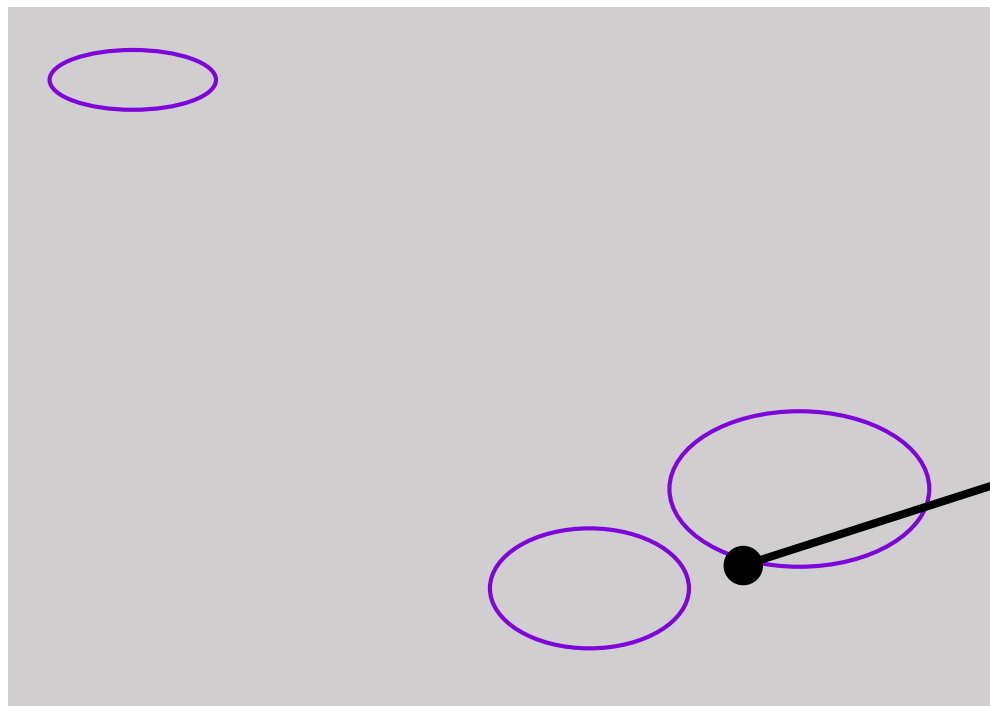
# A Corgi/Bread Decision Surface



$P(\text{Corgi}) \cong 0.5$

$P(\text{Corgi}) = 0.5$  Level Set  
True Posterior

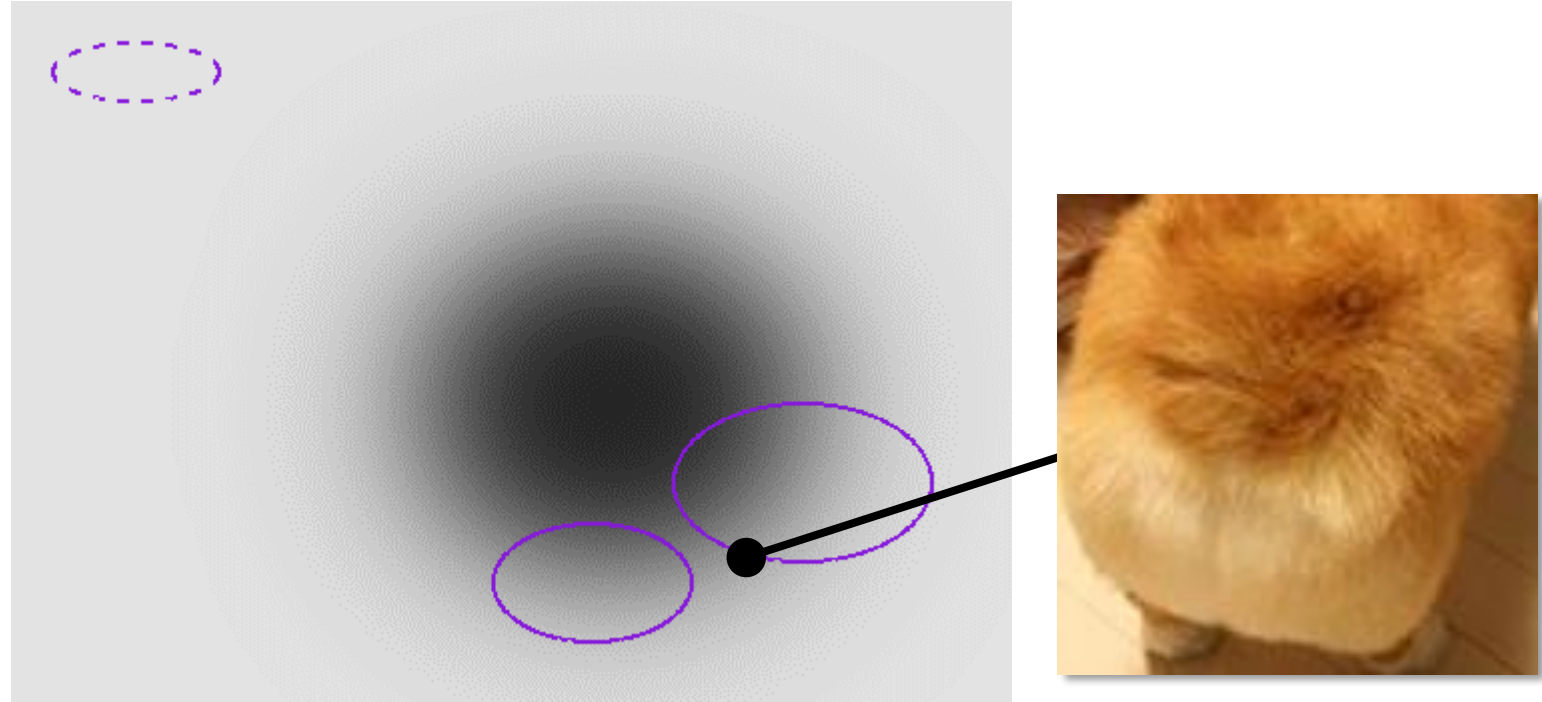
We find prediction-matching examples from **p-level sets**



$P(\text{Corgi}) = 0.5$  Level Set  
True Posterior

$P(\text{Corgi}) \cong 0.5$

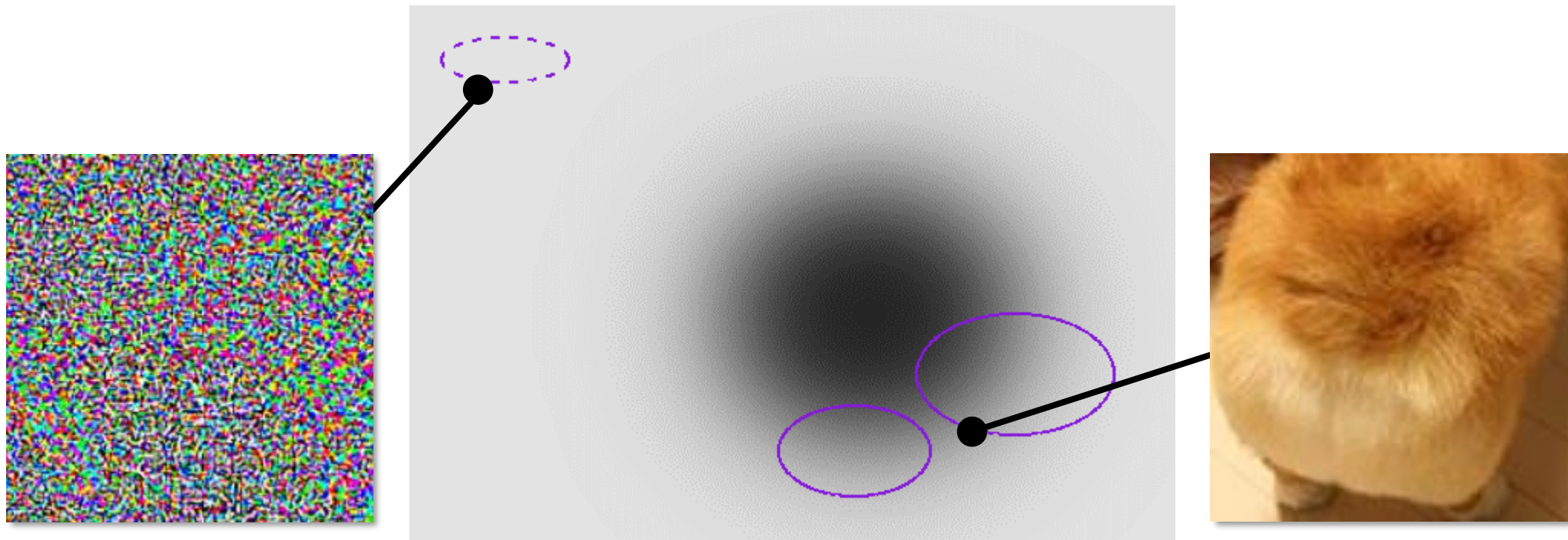
We find prediction-matching examples from **p-level sets**



$P(\text{Corgi}) = 0.5$  Level Set  
True Posterior

$P(\text{Corgi}) \cong 0.5$

We find prediction-matching examples from **p-level sets**

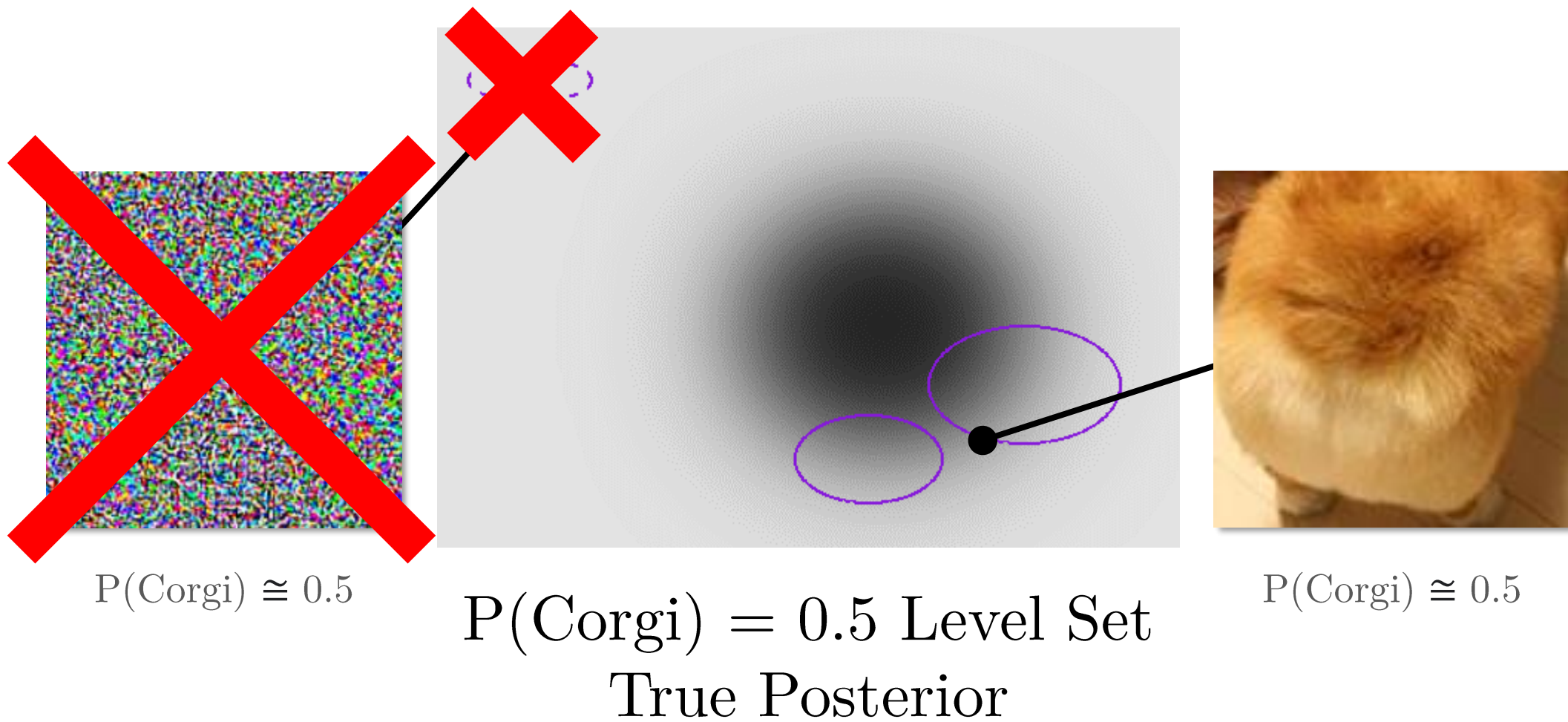


$P(\text{Corgi}) \cong 0.5$

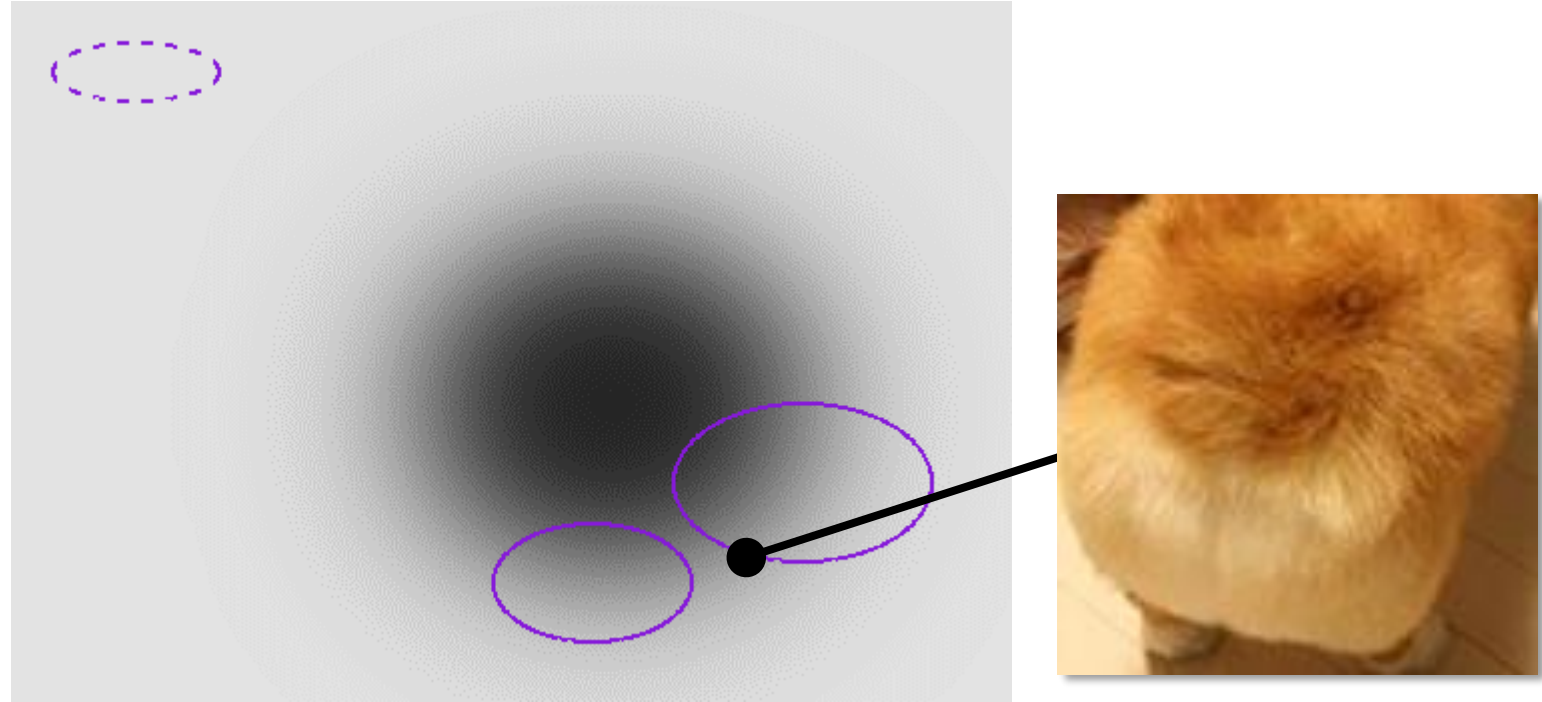
$P(\text{Corgi}) = 0.5$  Level Set  
True Posterior

$P(\text{Corgi}) \cong 0.5$

We find prediction-matching examples from **p-level sets**



We want to find a *natural* example  $\mathbf{x}$  where the classifier  $f(\mathbf{x})$  has confidence  $\mathbf{p}$

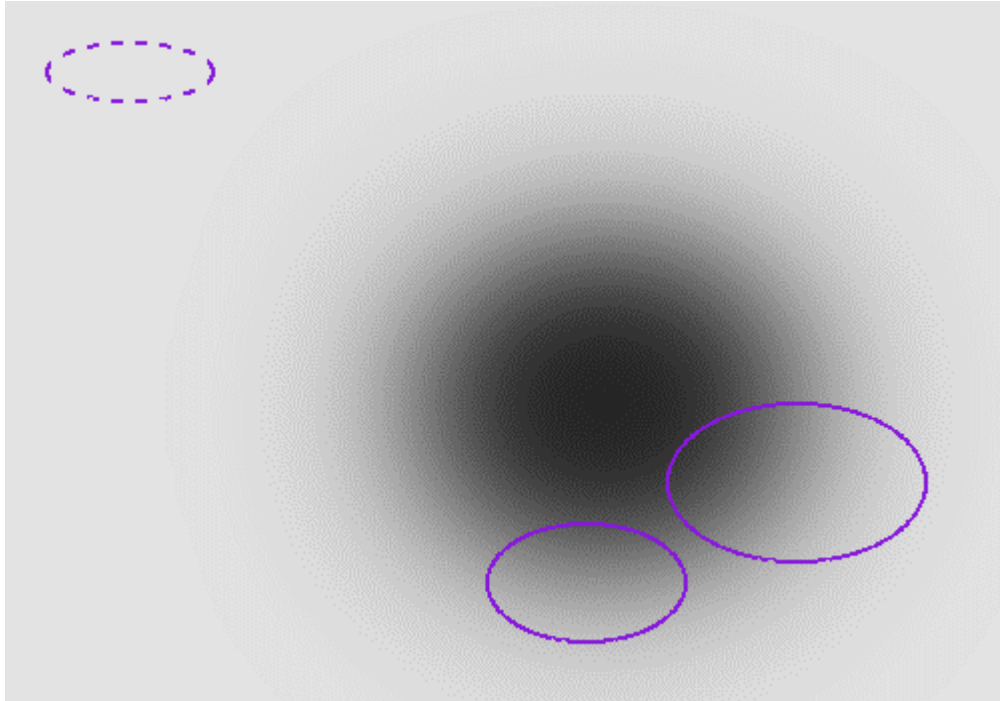


$P(\text{Corgi}) = 0.5$  Level Set  
True Posterior

$P(\text{Corgi}) \cong 0.5$



We want to find a *natural* example  $\mathbf{x}$   
where the classifier  $f(\mathbf{x})$  has confidence  $\mathbf{p}$

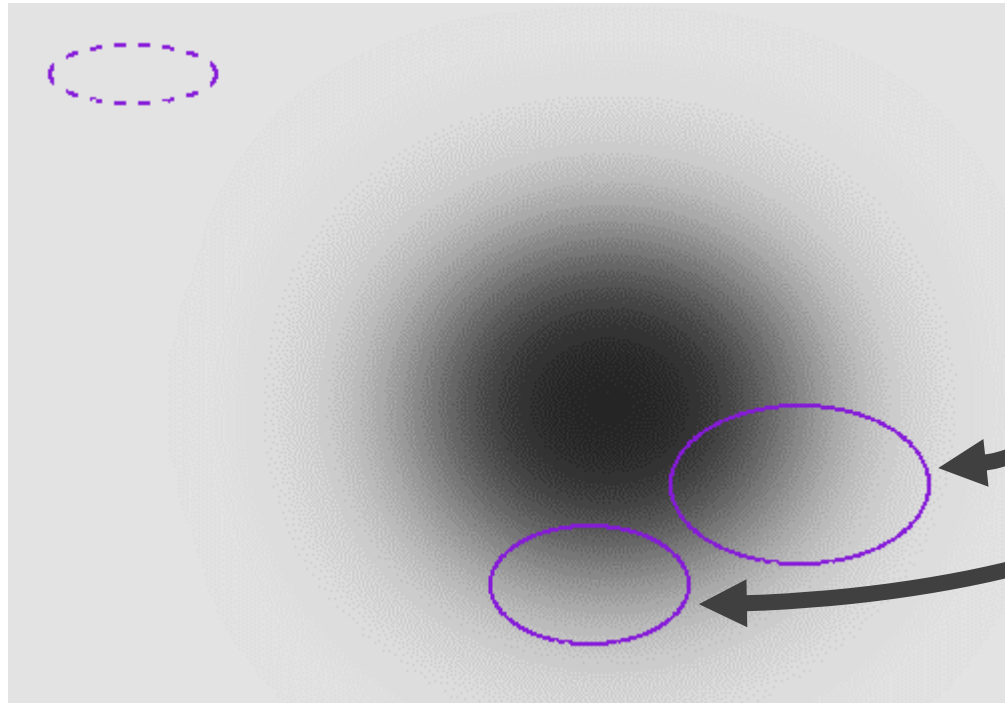


Want to sample from:

$$p(\mathbf{x} | f(\mathbf{x}) = \mathbf{p}) \propto p(\mathbf{x}) p(f(\mathbf{x}) = \mathbf{p} | \mathbf{x})$$

$P(\text{Corgi}) = 0.5$  Level Set  
True Posterior

Want to sample from:  $p(\mathbf{x}|f(\mathbf{x}) = \mathbf{p}) \propto p(\mathbf{x}) p(f(\mathbf{x}) = \mathbf{p}|\mathbf{x})$



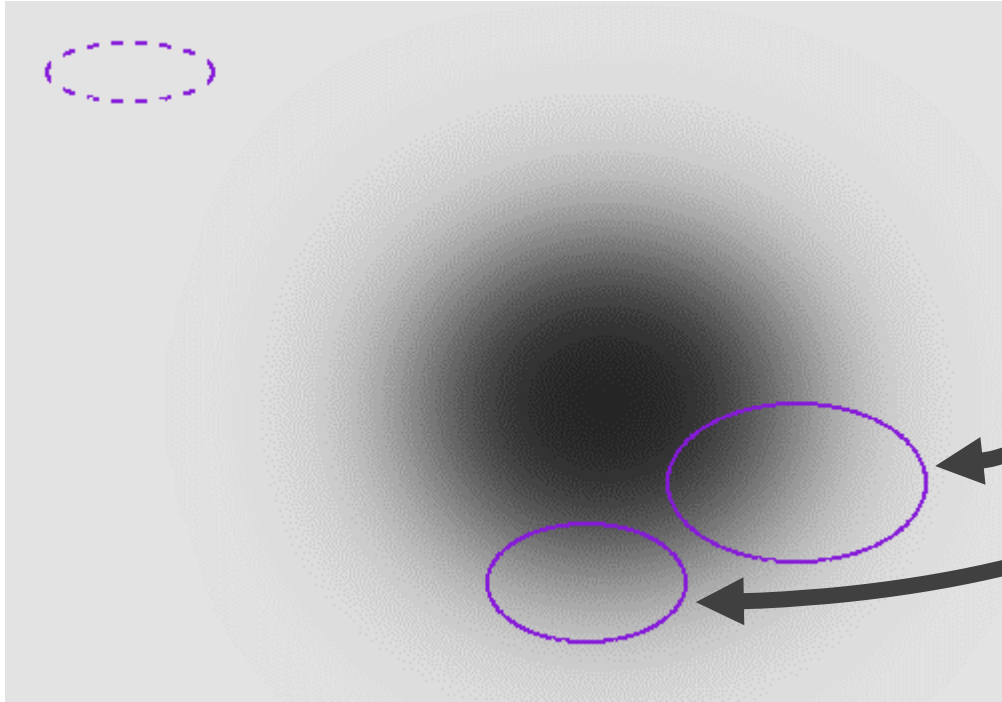
Problem 1:

$$\{\mathbf{x} : f(\mathbf{x}) = \mathbf{p}\}$$

has small or even zero measure

$P(\text{Corgi}) = 0.5$  Level Set  
True Posterior

Want to sample from:  $p(\mathbf{x}|f(\mathbf{x}) = \mathbf{p}) \propto p(\mathbf{x}) p(f(\mathbf{x}) = \mathbf{p}|\mathbf{x})$



$P(\text{Corgi}) = 0.5$  Level Set  
True Posterior

**Problem 1:**

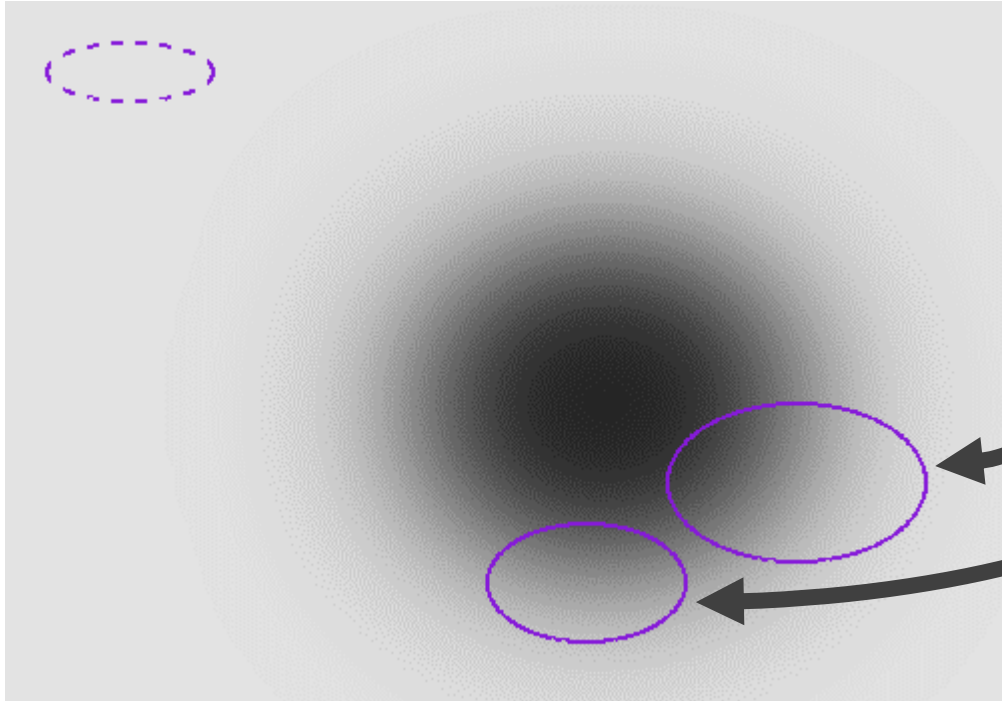
$$\{\mathbf{x} : f(\mathbf{x}) = \mathbf{p}\}$$

has small or even zero measure

**Problem 2:**

$\mathbf{x}$  is too high-dimensional

Want to sample from:  ~~$p(\mathbf{x}|f(\mathbf{x}) = \mathbf{p}) \propto p(\mathbf{x}) p(f(\mathbf{x}) = \mathbf{p}|\mathbf{x})$~~



$P(\text{Corgi}) = 0.5$  Level Set  
True Posterior

Problem 1:

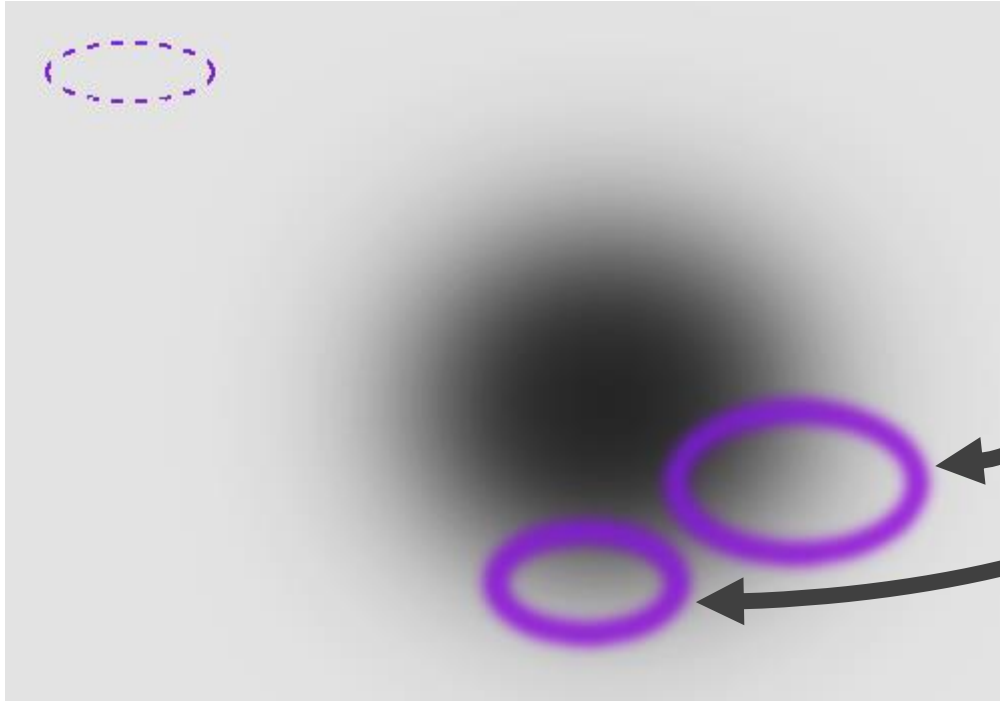
$$\{\mathbf{x} : f(\mathbf{x}) = \mathbf{p}\}$$

has small or even zero measure

Problem 2:

$\mathbf{x}$  is too high-dimensional

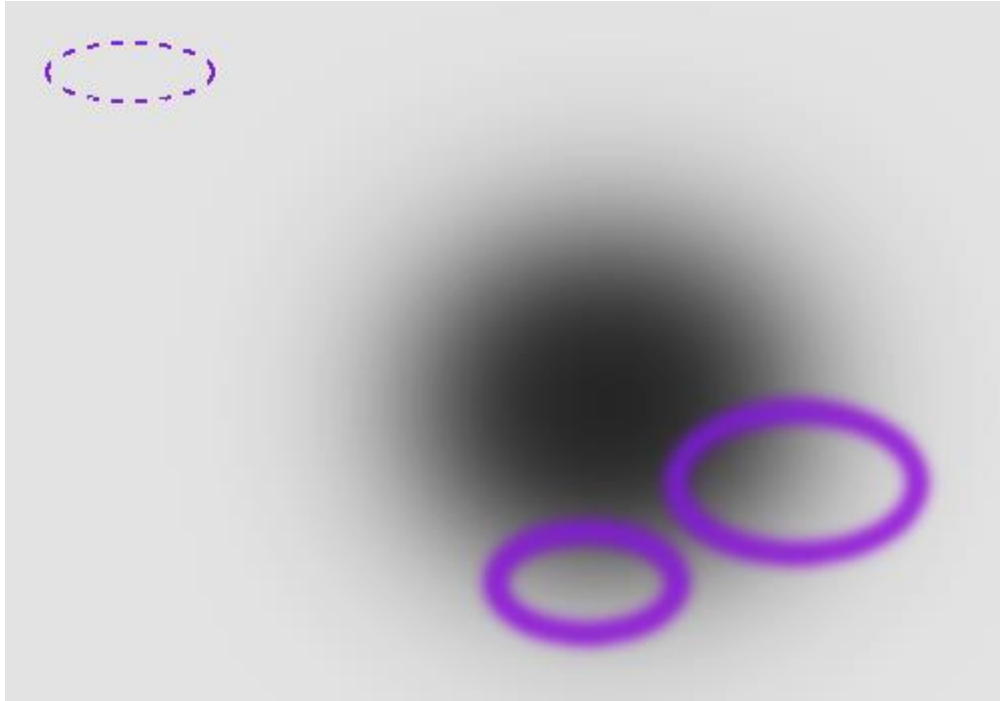
# Make inference tractable



Relax the formulation by  
“widening” the level set

$P(\text{Corgi}) = 0.5$  Level Set  
Relaxed Posterior

# Relaxed Formulation

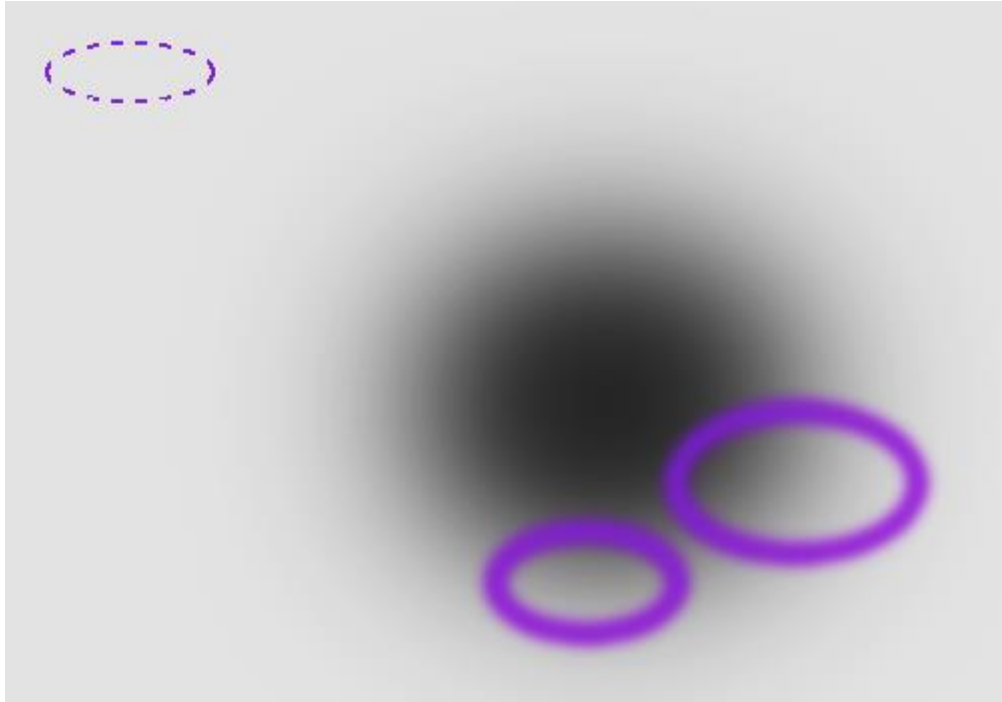


Introduce a random vector:

$$\mathbf{u}|\mathbf{x} \sim \mathcal{N}(f(\mathbf{x}), \sigma^2)$$

$P(\text{Corgi}) = 0.5$  Level Set  
Relaxed Posterior

# Relaxed Formulation



Introduce a random vector:

$$\mathbf{u}|\mathbf{x} \sim \mathcal{N}(f(\mathbf{x}), \sigma^2)$$

And sample from the new posterior:

$$p(\mathbf{x}|\mathbf{u} = \mathbf{u}^*) \propto p(\mathbf{x})p(\mathbf{u} = \mathbf{u}^*|\mathbf{x})$$

$$\mathbf{u}^* = \mathbf{p}$$

$P(\text{Corgi}) = 0.5$  Level Set  
Relaxed Posterior

But, how can we sample an image  $\mathbf{x}$ ?



But, how can we sample an image  $\mathbf{x}$ ?

We sample  $\mathbf{z}$  from a latent space  $Z$ , instead.

$$\mathbf{x} = g(\mathbf{z})$$

But, how can we sample an image  $\mathbf{x}$ ?

We sample  $\mathbf{z}$  from a latent space  $Z$ , instead.

$$\mathbf{x} = g(\mathbf{z})$$

$$\mathbf{u}|\mathbf{z} \sim N(f(\mathbf{x}), \sigma^2)$$

$$p(\mathbf{z}|\mathbf{u} = \mathbf{u}^*) \propto p(\mathbf{z})p(\mathbf{u} = \mathbf{u}^*|\mathbf{z})$$

To use Bayes-TrEx, we need 3 requirements:

To use Bayes-TrEx, we need 3 requirements:

1. A classifier which outputs class probabilities

To use Bayes-TrEx, we need 3 requirements:

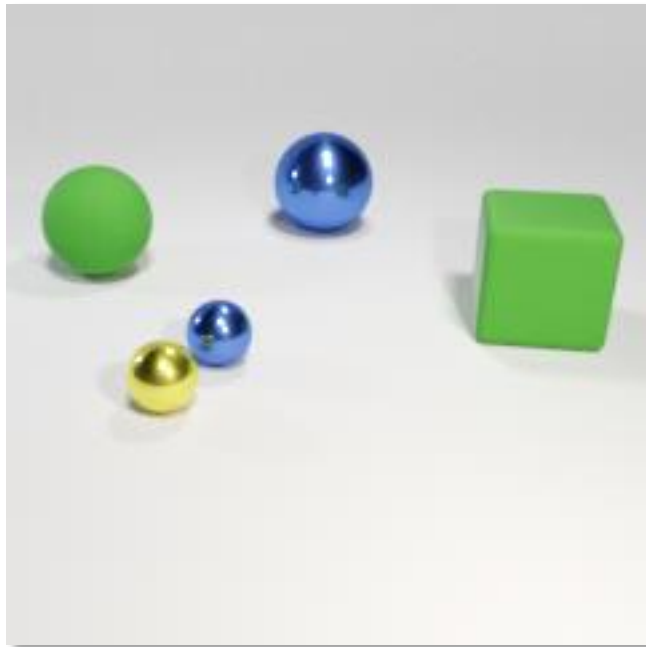
1. A classifier which outputs class probabilities
2. A user-specified confidence target

To use Bayes-TrEx, we need 3 requirements:

1. A classifier which outputs class probabilities
2. A user-specified confidence target
3. A data distribution we can sample from

# Experiments

CLEVR



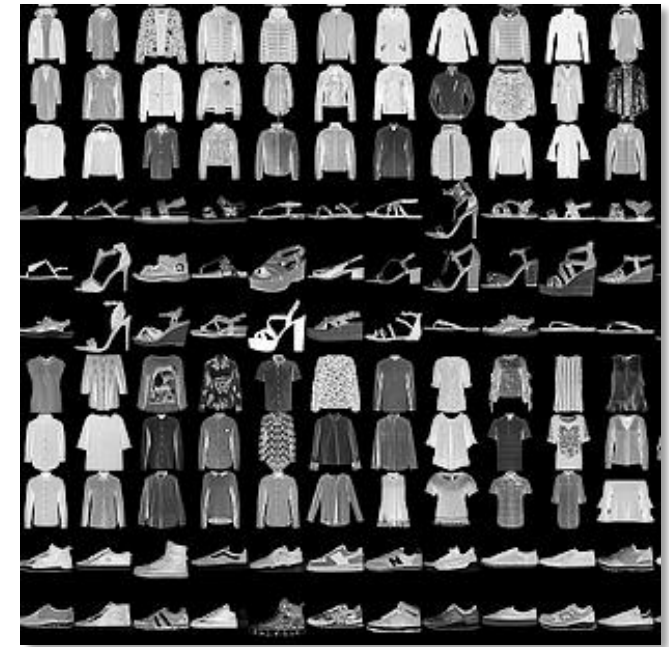
Scene graph

MNIST



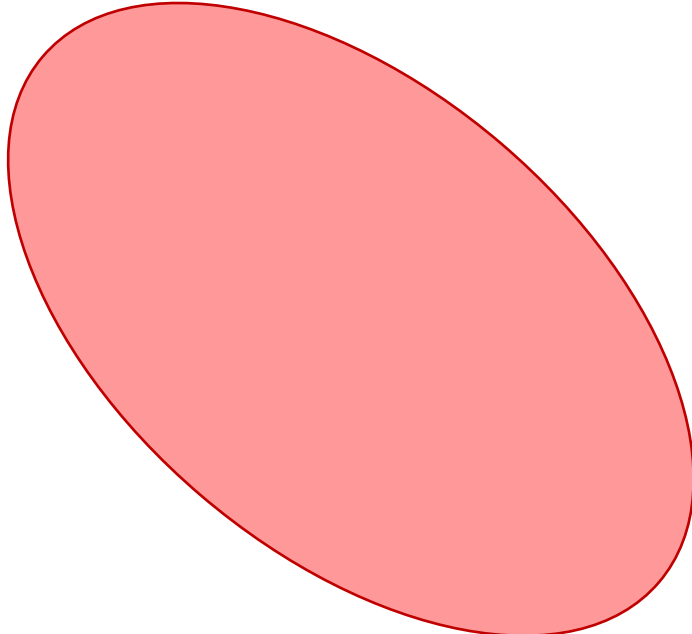
VAEs, GANs

Fashion-MNIST



VAEs, GANs

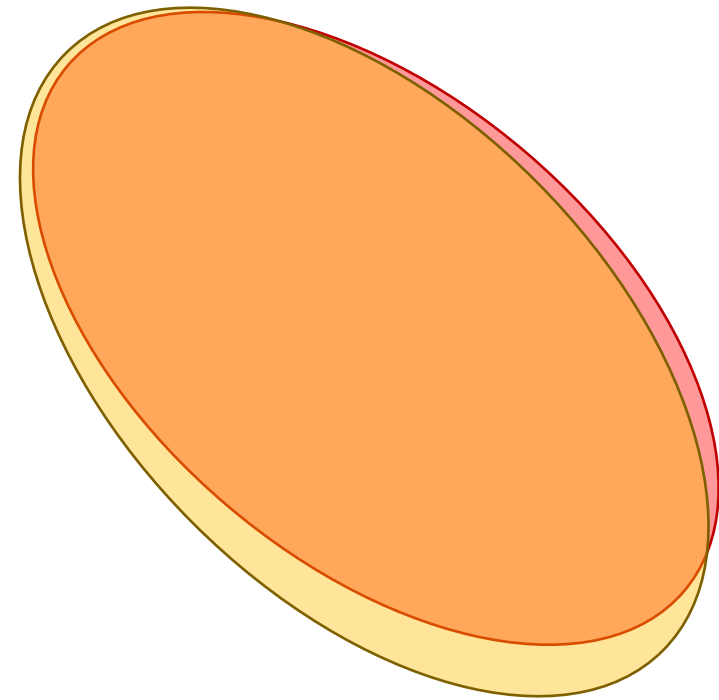
Smoke Test: High Confidence Examples ( $p_i=1, p_{\neg i}=0$ )



$\mathbb{P}_C$ , classifier training distribution, red



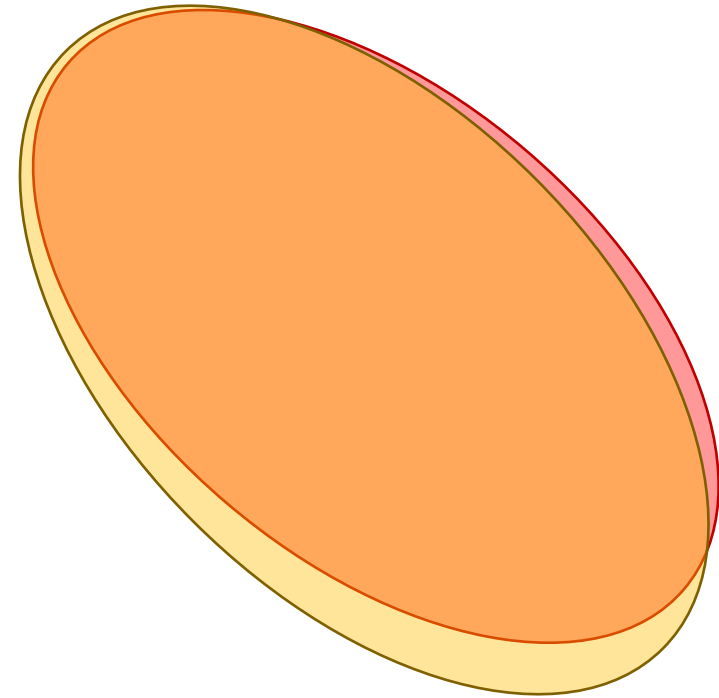
# Smoke Test: High Confidence Examples ( $p_i=1, p_{\neg i}=0$ )



$\mathbb{P}_C$ , classifier training distribution, red

$\mathbb{P}_D$ , Bayes-TrEx's data distribution, yellow

# Smoke Test: High Confidence Examples ( $p_i=1, p_{\neg i}=0$ )

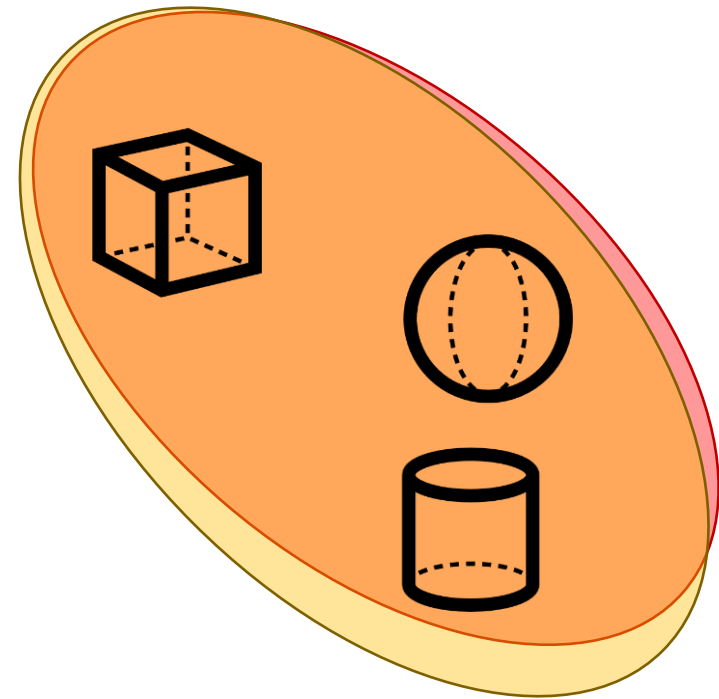


$\mathbb{P}_C$ , classifier training distribution, red

$\mathbb{P}_D$ , Bayes-TrEx's data distribution, yellow

$$\mathbb{P}_D \approx \mathbb{P}_C$$

# Smoke Test: High Confidence Examples ( $p_i=1, p_{\neg i}=0$ )

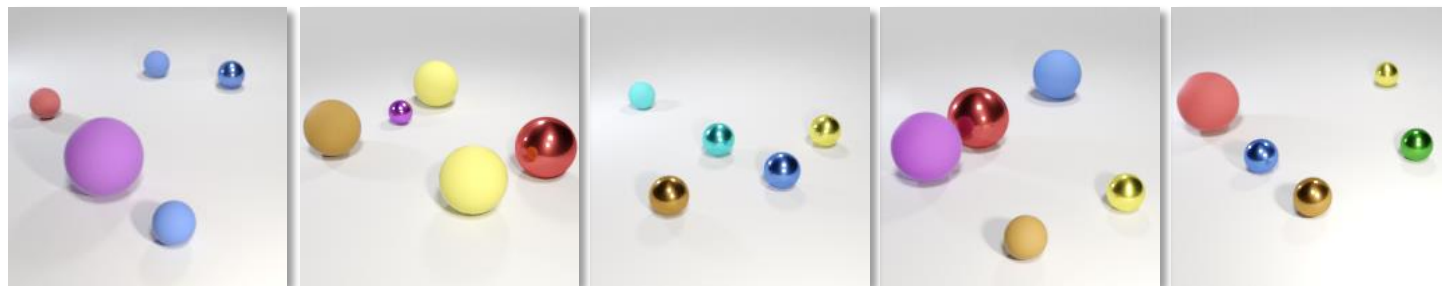


$\mathbb{P}_C$ , classifier training distribution, red

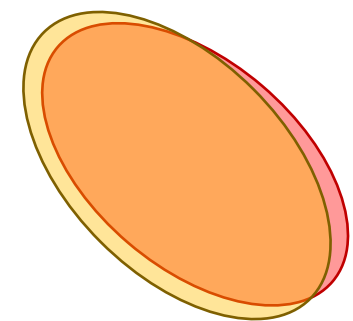
$\mathbb{P}_D$ , Bayes-TrEx's data distribution, yellow

$$\mathbb{P}_D \approx \mathbb{P}_C$$

# Smoke Test: High Confidence Examples ( $p_i=1, p_{\neg i}=0$ )

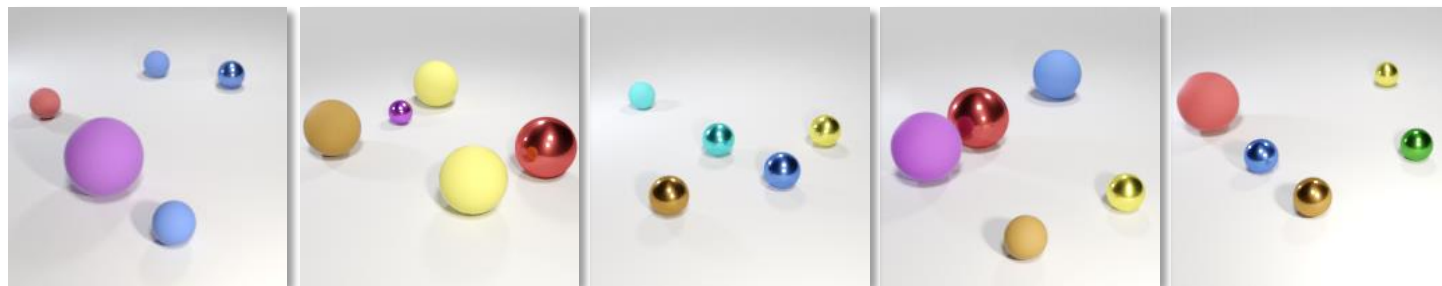


CLEVR:  
5 Spheres



$$\mathbb{P}_D \approx \mathbb{P}_C$$

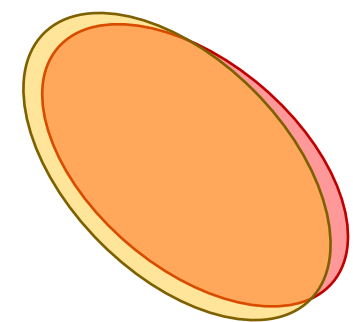
# Smoke Test: High Confidence Examples ( $p_i=1, p_{\neg i}=0$ )



CLEVR:  
5 Spheres

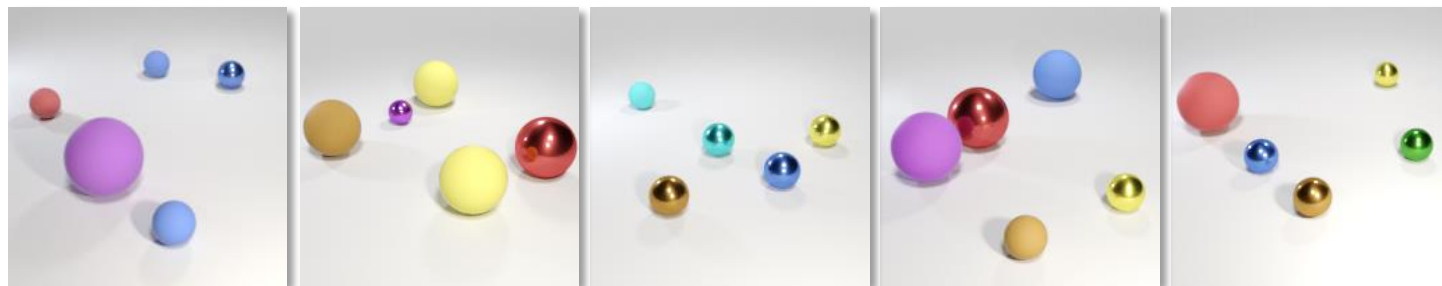


MNIST



$$\mathbb{P}_D \approx \mathbb{P}_C$$

# Smoke Test: High Confidence Examples ( $p_i=1, p_{\neg i}=0$ )



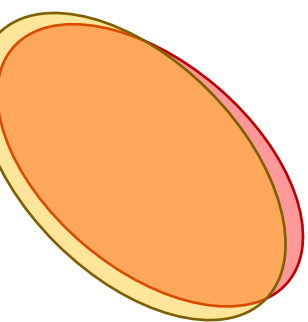
CLEVR:  
5 Spheres



MNIST

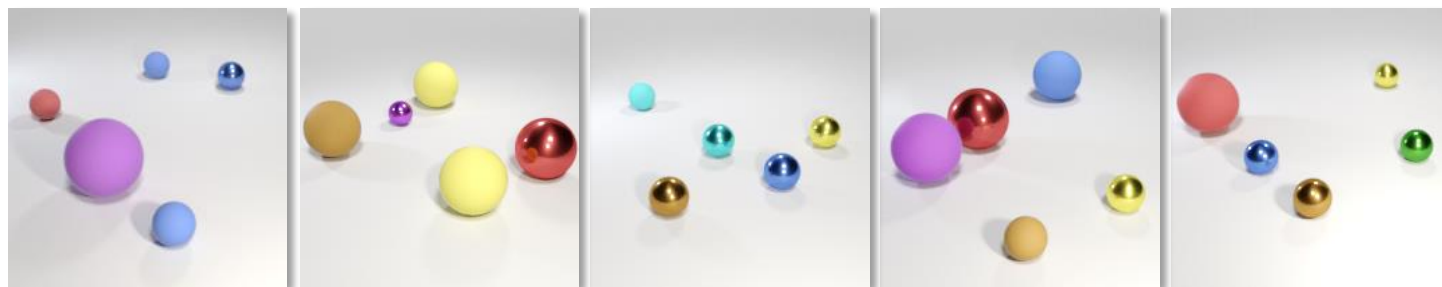


Fashion-  
MNIST



$$\mathbb{P}_D \approx \mathbb{P}_C$$

# Smoke Test: High Confidence Examples ( $p_i=1, p_{\neg i}=0$ )



$\approx 0.943$

CLEVR:  
5 Spheres



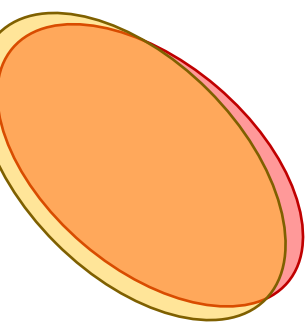
$\approx 0.999$

MNIST



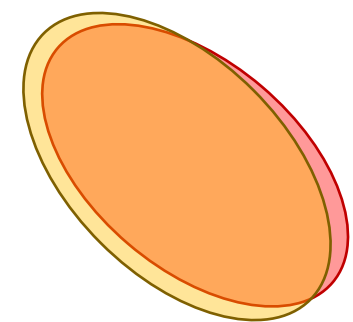
$\approx 0.998$

Fashion-  
MNIST

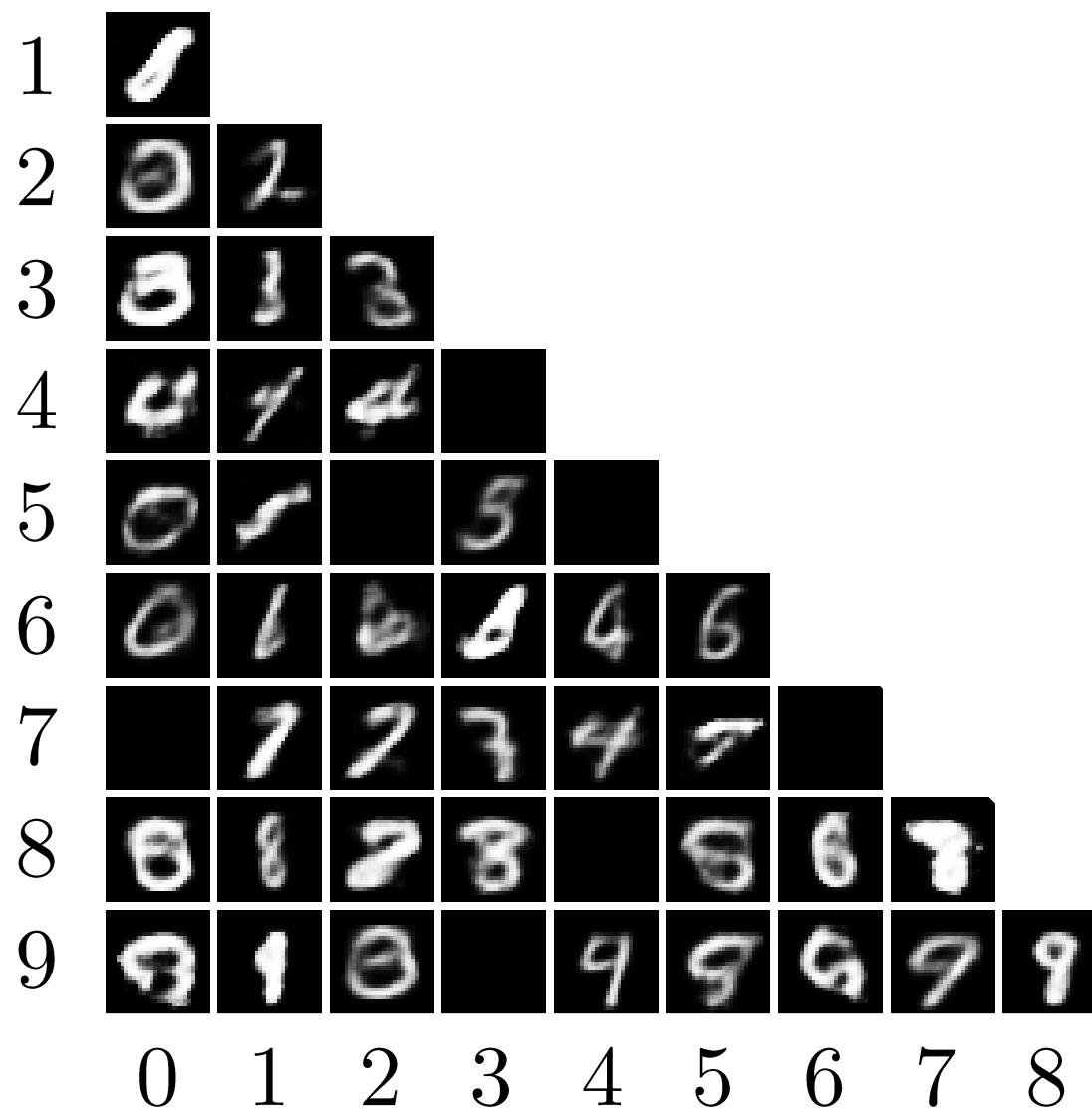


$$\mathbb{P}_D \approx \mathbb{P}_C$$

# Class Boundaries ( $p_i=0.5, p_j=0.5, p_{\neg i, \neg j}=0$ )

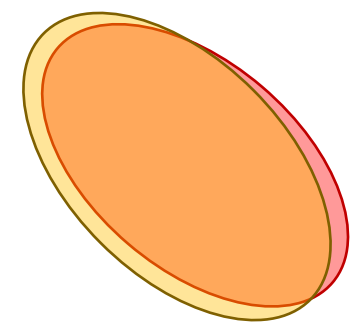


$$\mathbb{P}_D \approx \mathbb{P}_C$$

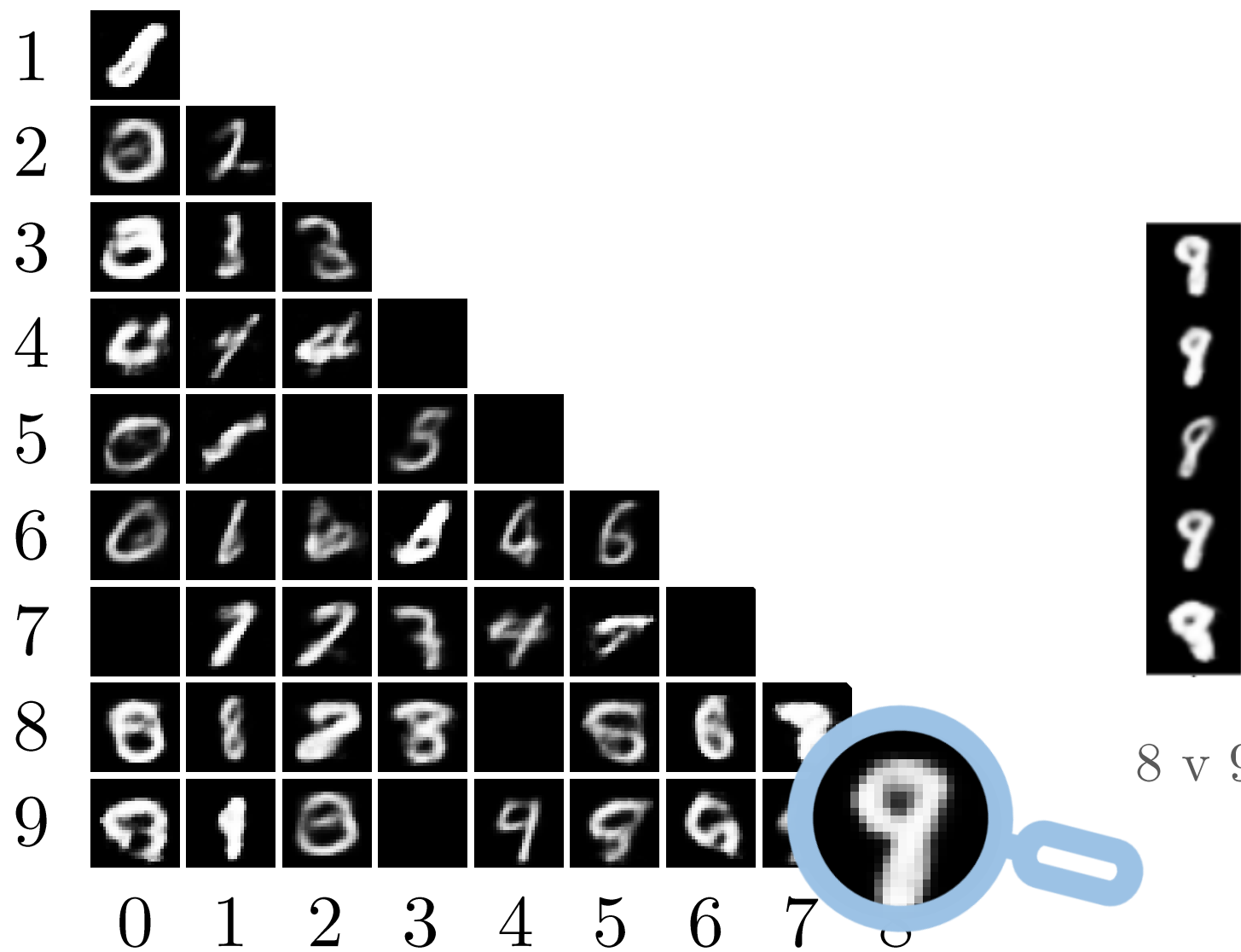




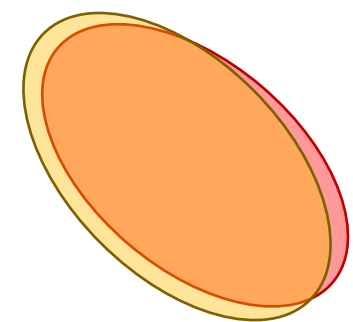
# Class Boundaries ( $p_i=0.5, p_j=0.5, p_{\neg i, \neg j}=0$ )



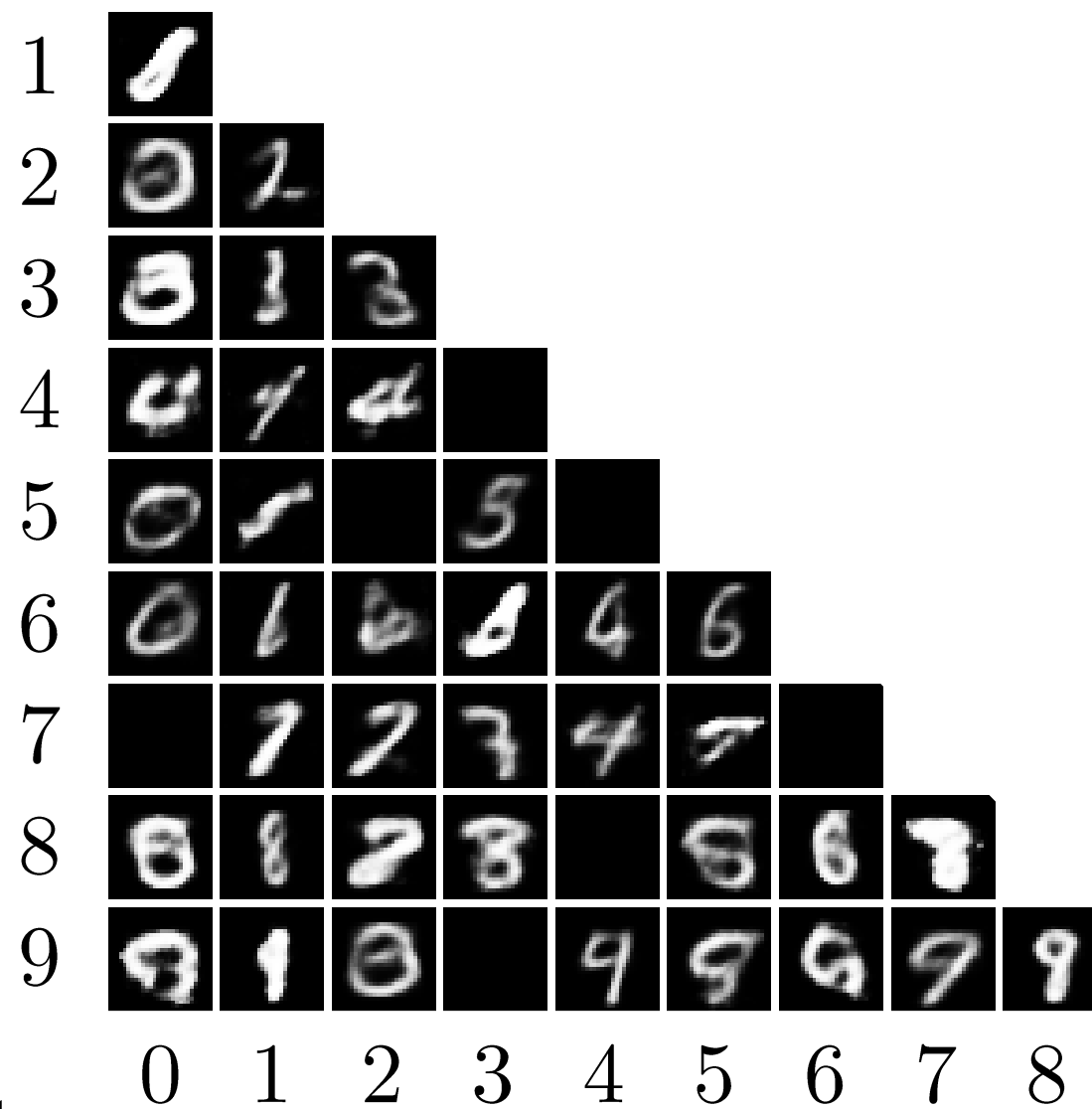
$$\mathbb{P}_D \approx \mathbb{P}_C$$



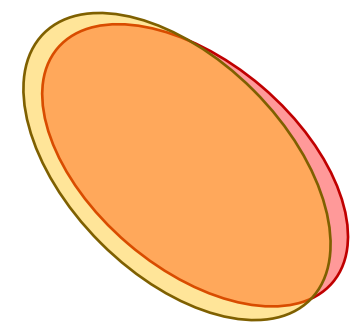
# Class Boundaries ( $p_i=0.5, p_j=0.5, p_{\neg i, \neg j}=0$ )



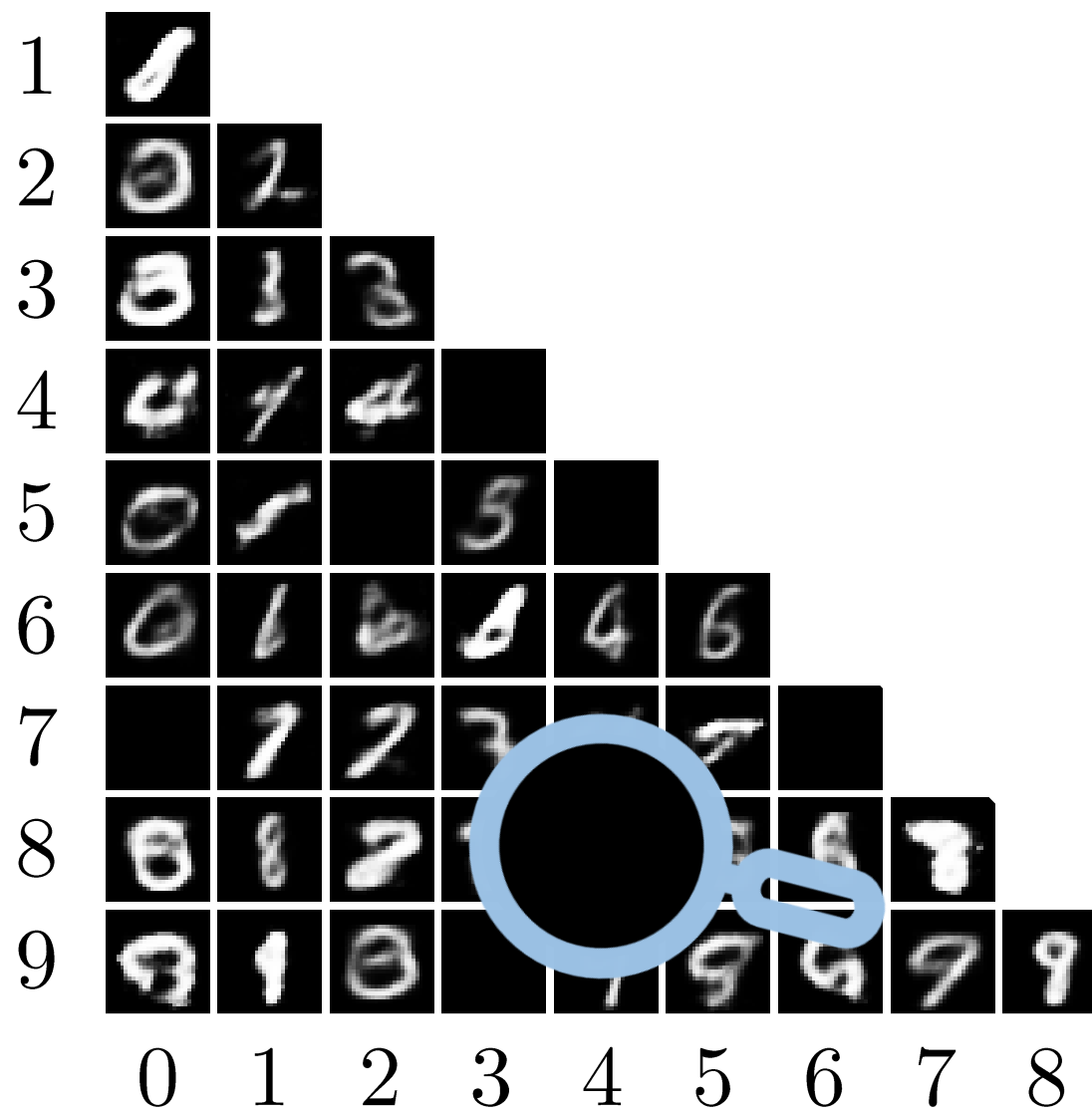
$$\mathbb{P}_D \approx \mathbb{P}_C$$



# Class Boundaries ( $p_i=0.5, p_j=0.5, p_{\neg i, \neg j}=0$ )

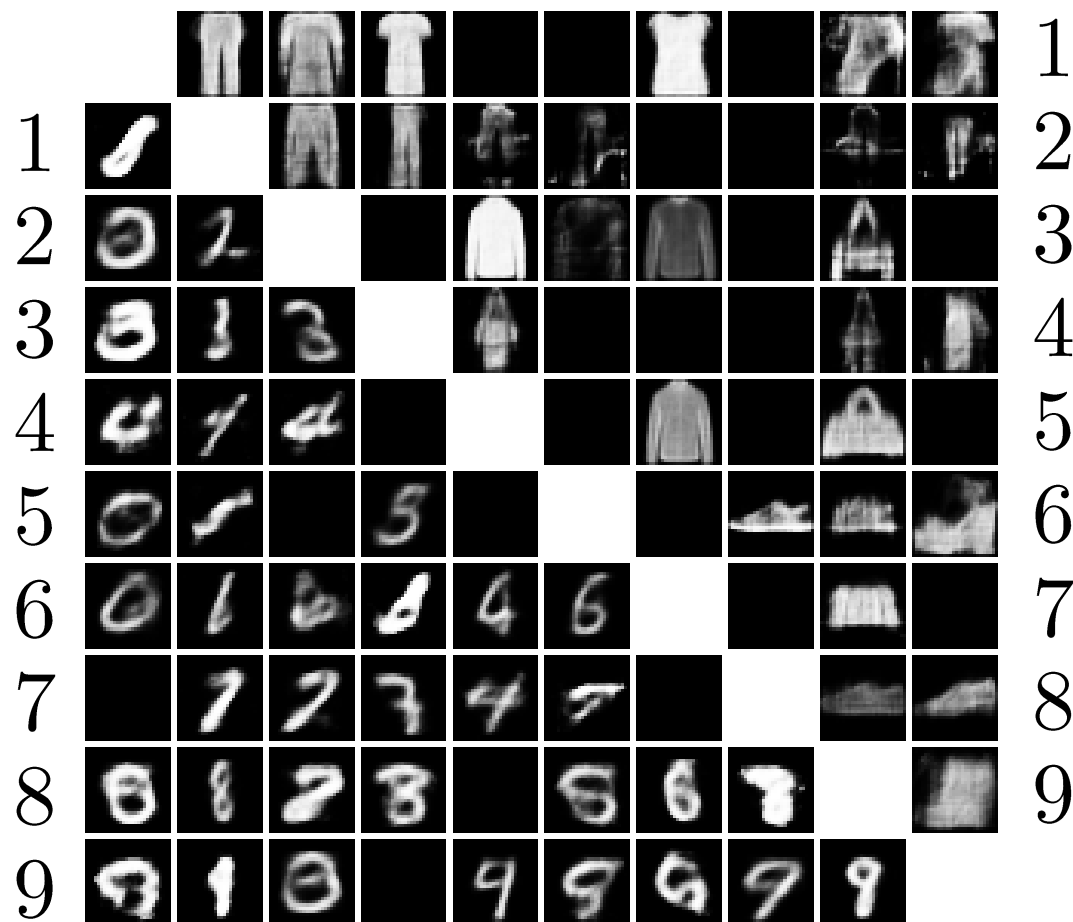


$$\mathbb{P}_D \approx \mathbb{P}_C$$



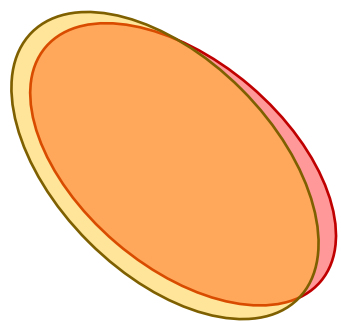
# Class Boundaries ( $p_i=0.5, p_j=0.5, p_{\neg i, \neg j}=0$ )

0 1 2 3 4 5 6 7 8



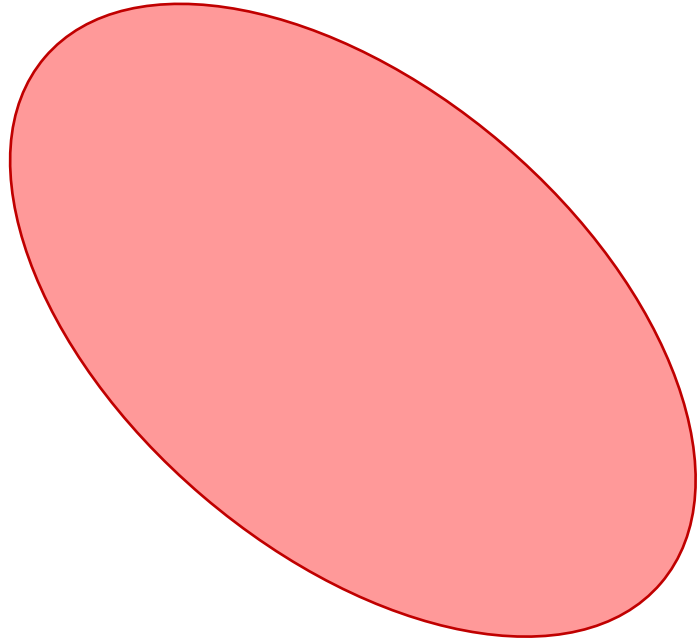
1  
2  
3  
4  
5  
6  
7  
8  
9

0 1 2 3 4 5 6 7 8



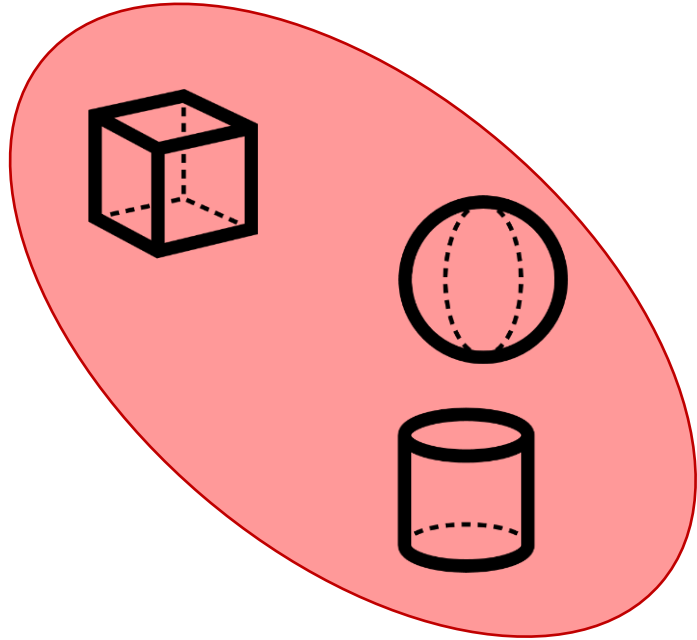
$$\mathbb{P}_D \approx \mathbb{P}_C$$

# High Confidence Failures



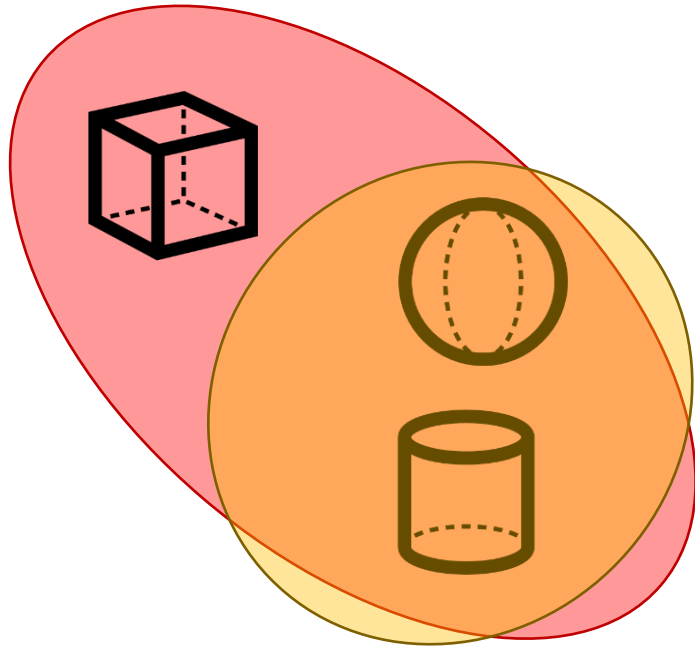
$\mathbb{P}_C$ , classifier training distribution, red

# High Confidence Failures



$\mathbb{P}_C$ , classifier training distribution, red

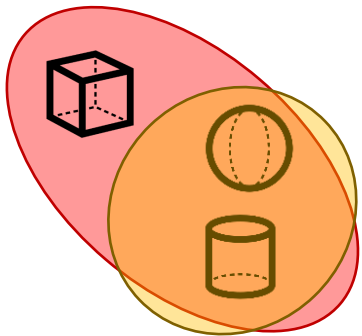
# High Confidence Failures



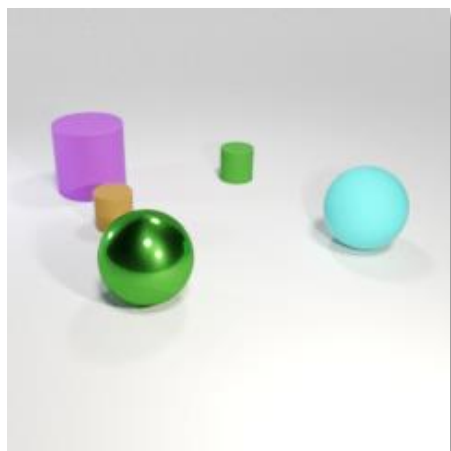
$\mathbb{P}_C$ , classifier training distribution, red

$\mathbb{P}_D$ , Bayes-TrEx's data distribution, yellow

$$\mathbb{P}_D \subsetneq \mathbb{P}_C$$



# High Confidence Failures



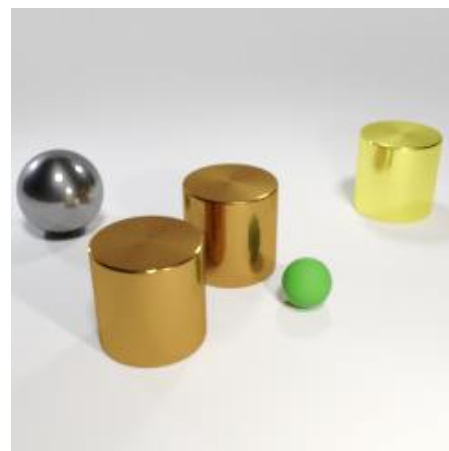
97.2%



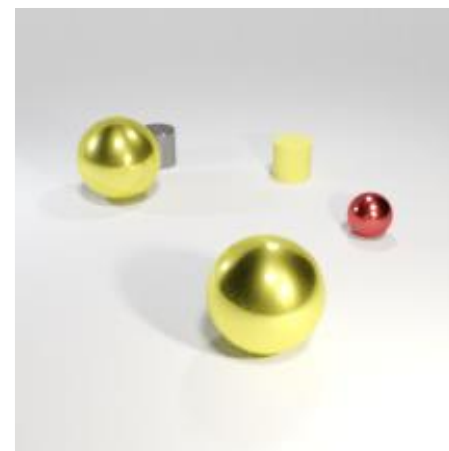
96.0%



94.5%



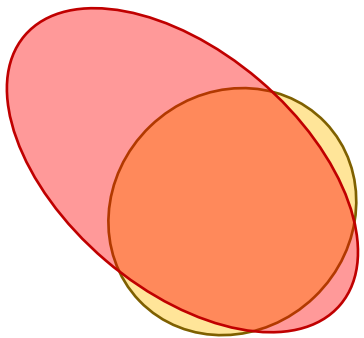
67.3%



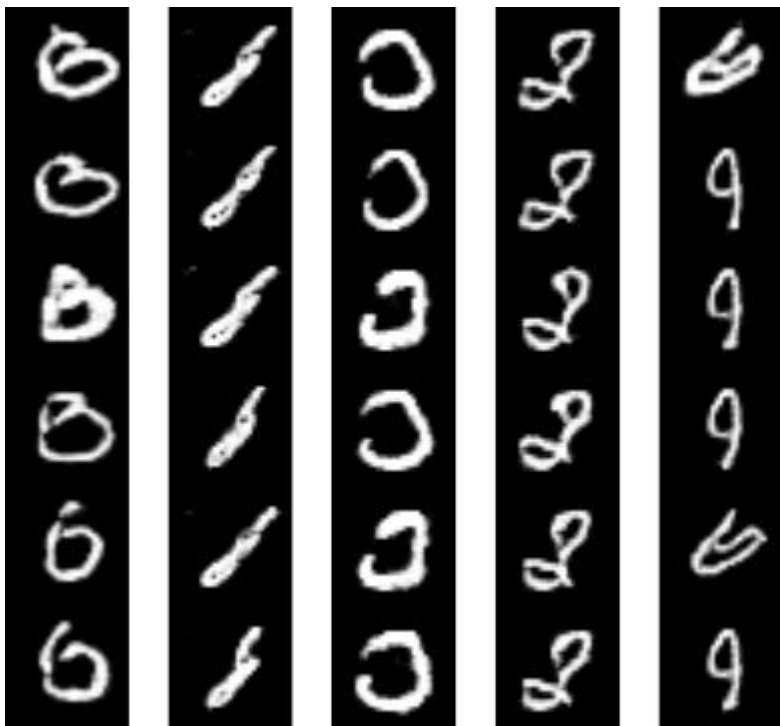
93.5%

Contains 1 Cube

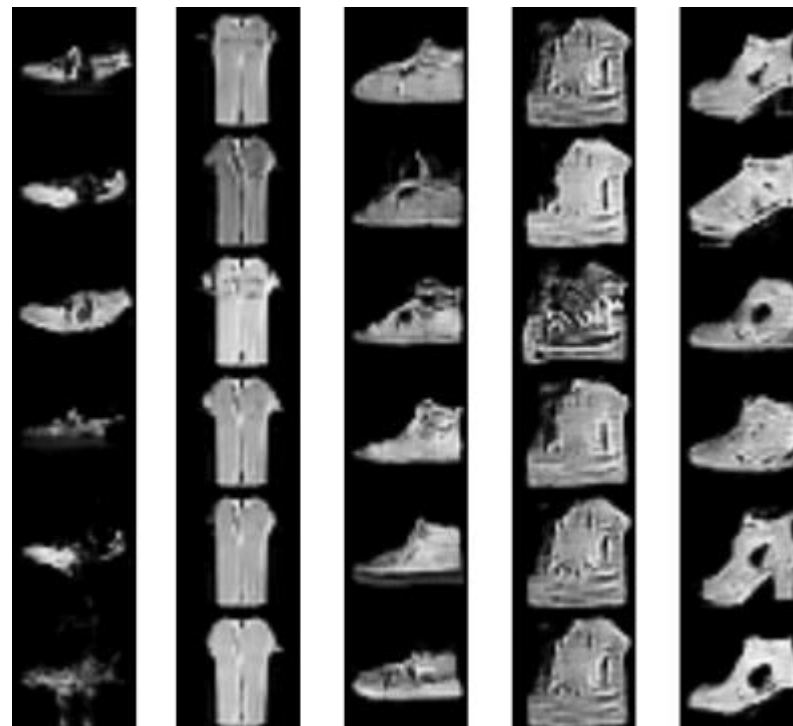




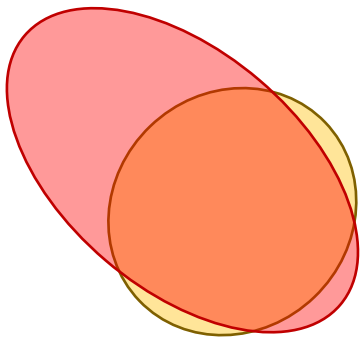
# High Confidence Failures



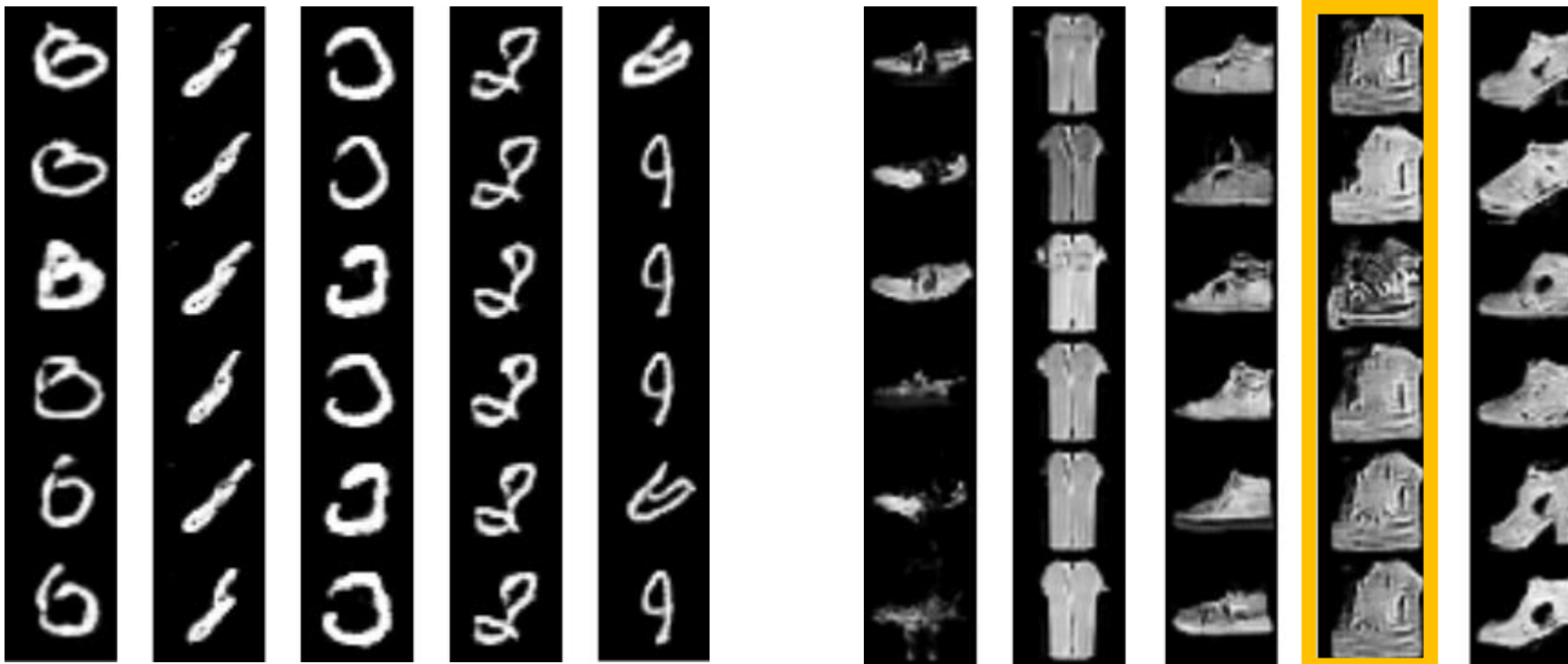
0:  $0.981 \pm 0.027$   
1:  $0.953 \pm 0.028$   
2:  $0.968 \pm 0.028$   
3:  $0.969 \pm 0.027$   
4:  $0.955 \pm 0.030$



Sandal:  $0.986 \pm 0.030$   
Shirt:  $0.938 \pm 0.032$   
Sneaker:  $0.969 \pm 0.028$   
Bag:  $0.967 \pm 0.026$   
Ankle boot:  $0.971 \pm 0.027$



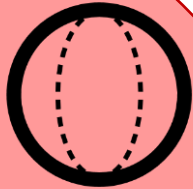
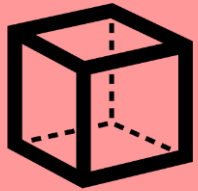
# High Confidence Failures



0:  $0.981 \pm 0.027$   
1:  $0.953 \pm 0.028$   
2:  $0.968 \pm 0.028$   
3:  $0.969 \pm 0.027$   
4:  $0.955 \pm 0.030$

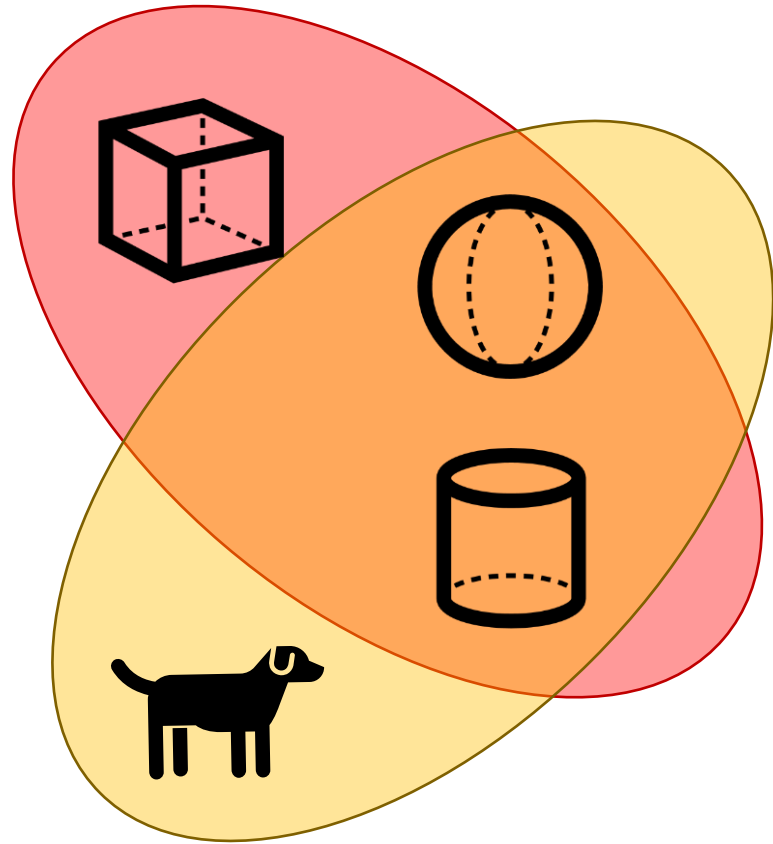
Sandal:  $0.986 \pm 0.030$   
Shirt:  $0.938 \pm 0.032$   
Sneaker:  $0.969 \pm 0.028$   
Bag:  $0.967 \pm 0.026$   
Ankle boot:  $0.971 \pm 0.027$

# Novel Class Extrapolation



$\mathbb{P}_C$ , classifier training distribution, red

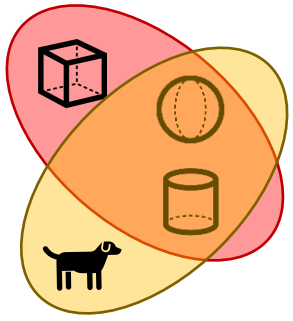
# Novel Class Extrapolation



$\mathbb{P}_C$ , classifier training distribution, red

$\mathbb{P}_D$ , Bayes-TrEx's data distribution, yellow

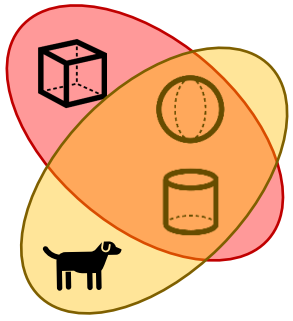
$$\mathbb{P}_D \cap \mathbb{P}_C \neq \mathbb{P}_D \neq \mathbb{P}_C$$



# Novel Class Extrapolation



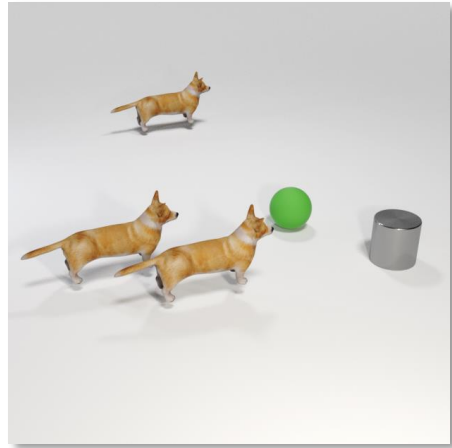
Contains 5 Cubes: 89.3%



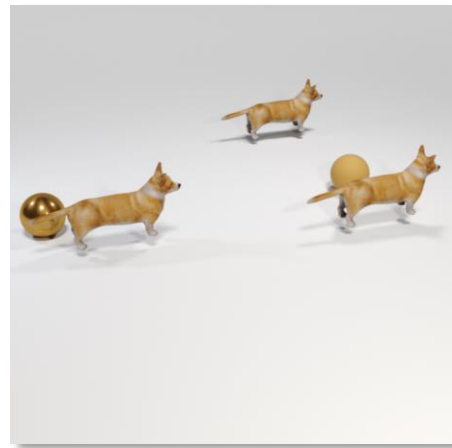
# Novel Class Extrapolation



89.3%



81.2%



83.5%



90.4%



90.5%

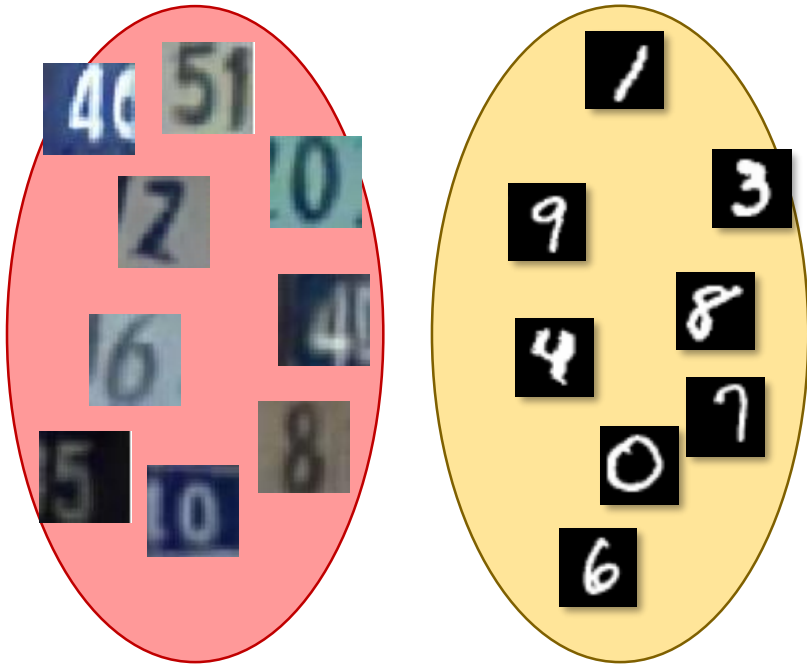
Contains 5 Cubes

# Domain Adaptation



$\mathbb{P}_C$ , classifier training distribution, red

# Domain Adaptation



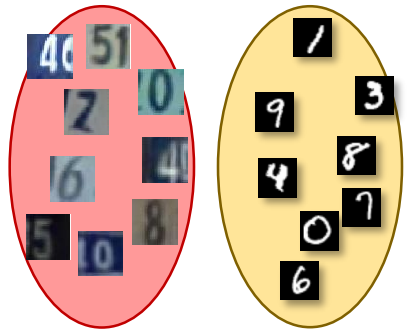
$$\mathbb{P}_D \cap \mathbb{P}_C = \emptyset$$

$\mathbb{P}_C$ , classifier training distribution, red

$\mathbb{P}_D$ , Bayes-TrEx's data distribution, yellow



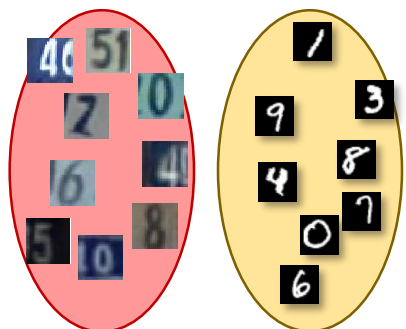
# Domain Adaptation: ADDA



Overall:

- Baseline: 61% accuracy
- ADDA: 71% accuracy

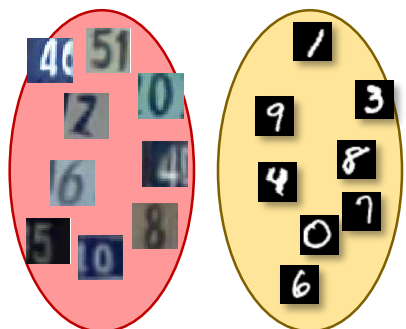
# Domain Adaptation: ADDA



Overall:

- Baseline: 61% accuracy
- **ADDA: 71% accuracy**

# Domain Adaptation: ADDA



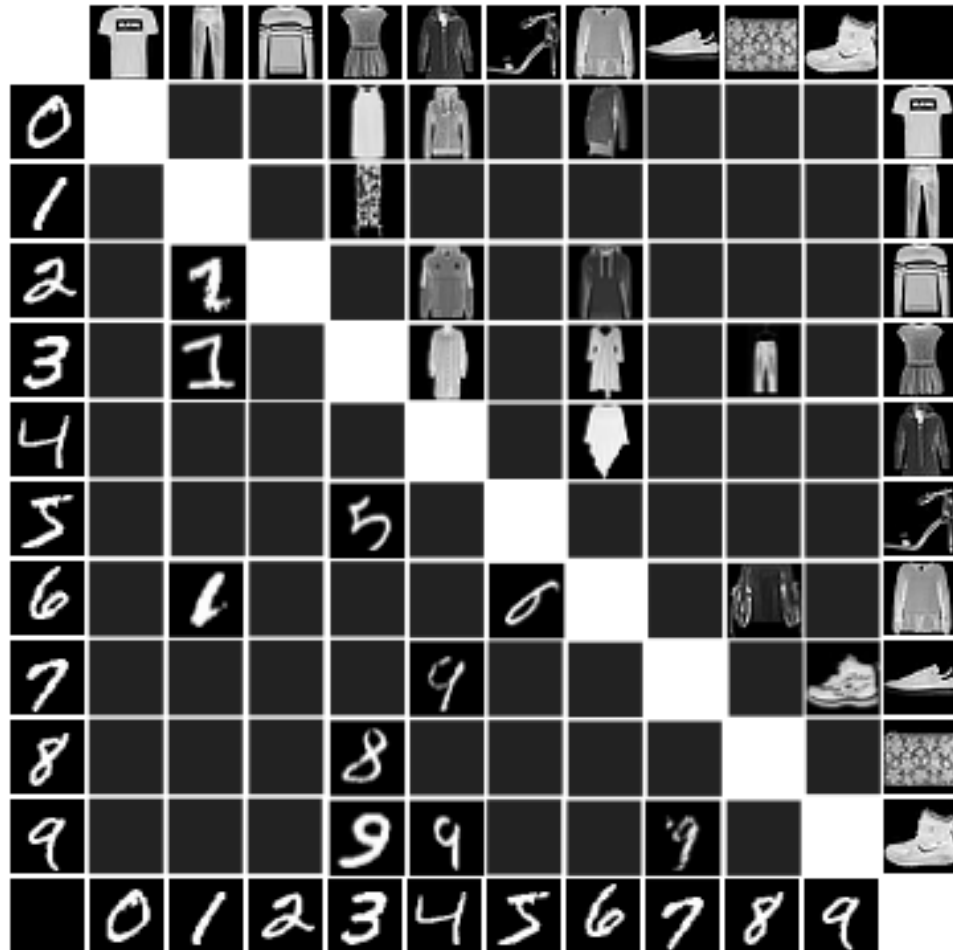
Overall:

- Baseline: 61% accuracy
- **ADDA: 71% accuracy**

High Confidence Examples:

- **Baseline: 80% accuracy**
- **ADDA: 72% accuracy**

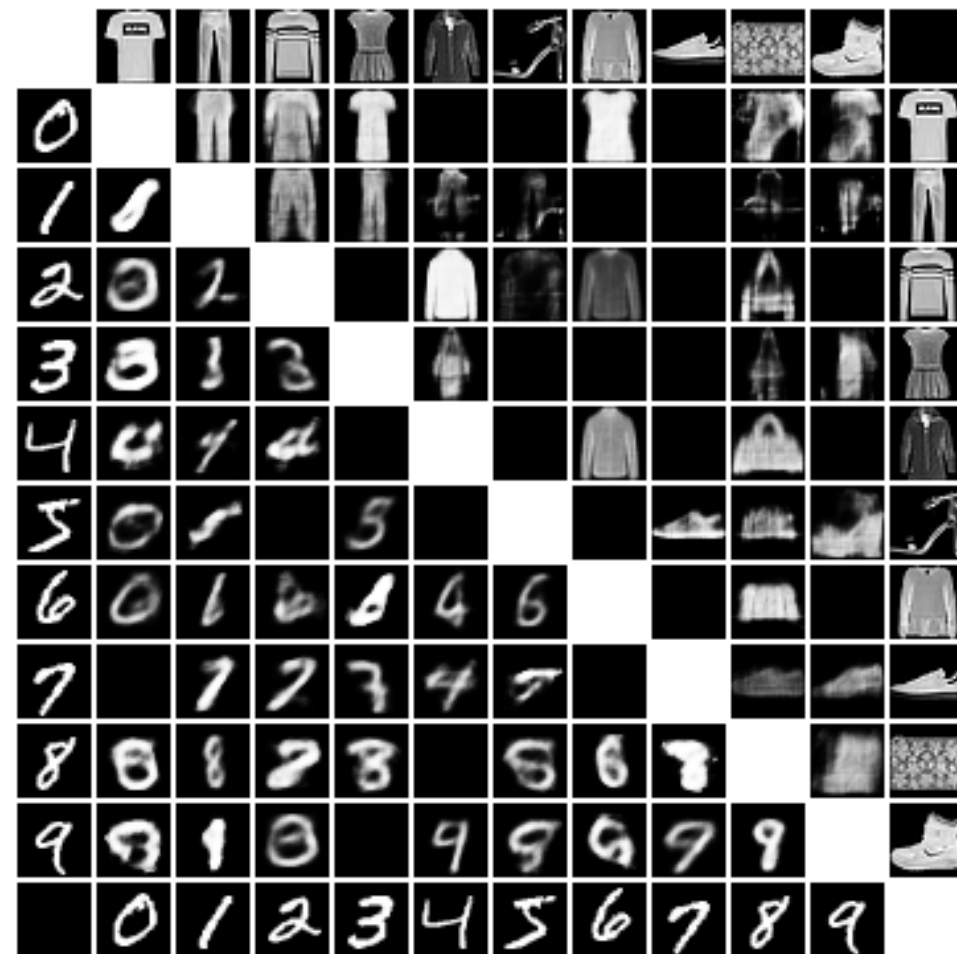
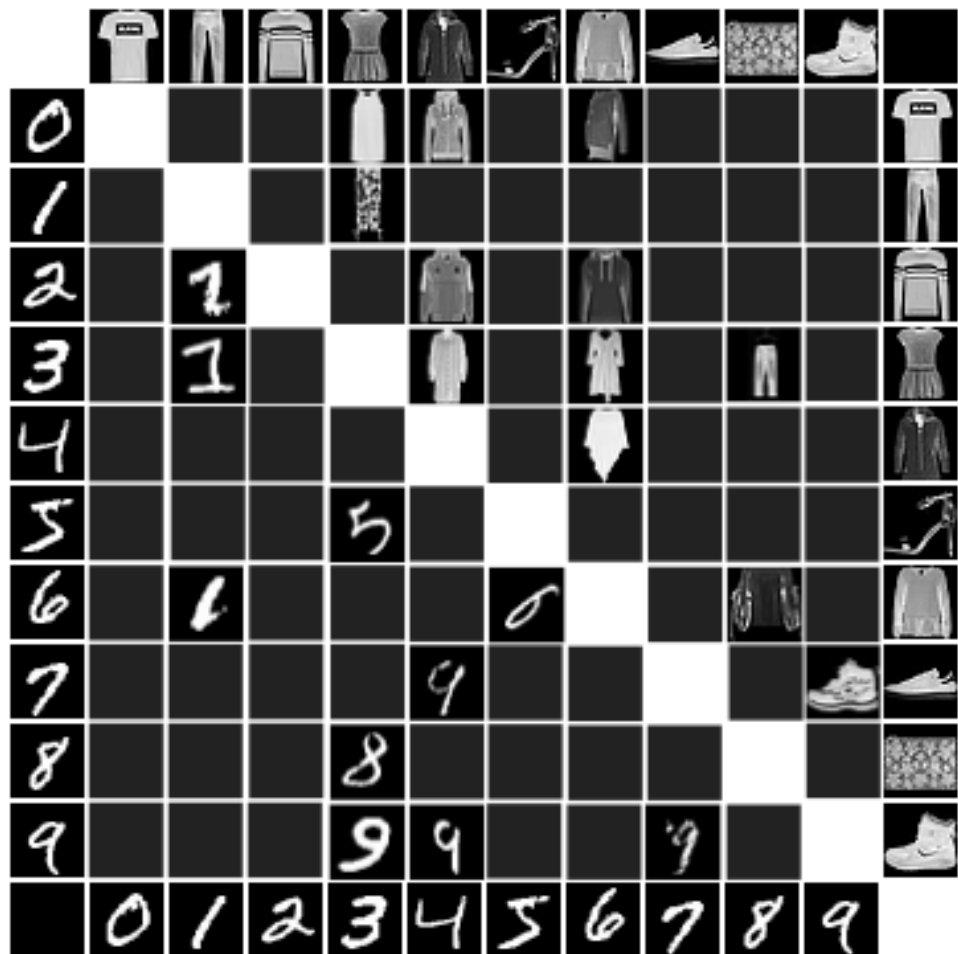
# How does Bayes-TrEx compare to the test set?



Test set:

22 Ambiguous Pairings

# How does Bayes-TrEx compare to the test set?



How does Bayes-TrEx compare to the test set?

Test set exposes more *mislabelings* than true classification failures.

# How does Bayes-TrEx compare to the test set?

Class: 2

Mislabeled

Misclassified

Test set exposes more *mislabelings* than true classification failures.

# How does Bayes-TrEx compare to the test set?

Class: 2

Mislabeled



Misclassified

Test set exposes more *mislabelings* than true classification failures.



# How does Bayes-TrEx compare to the test set?

Class: 2

Mislabeled



Misclassified



Test set exposes more *mislabelings* than true classification failures.

# How does Bayes-TrEx compare to the test set?

Class: 2

Mislabeled

Misclassified



Test set exposes more *mislabelings* than true classification failures.

# How does Bayes-TrEx compare to the test set?

Class: 2

Mislabeled

Misclassified



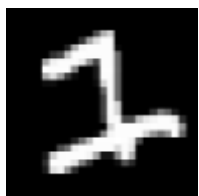
Test set exposes more *mislabelings* than true classification failures.

# How does Bayes-TrEx compare to the test set?

Class: 2

Mislabeled

Misclassified



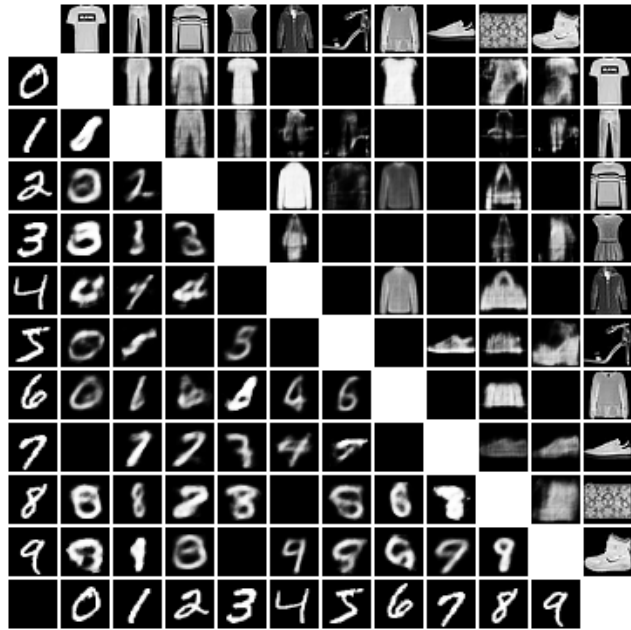
Test set exposes more *mislabelings* than true classification failures.

60/84 MNIST

42/93 Fashion-MNIST

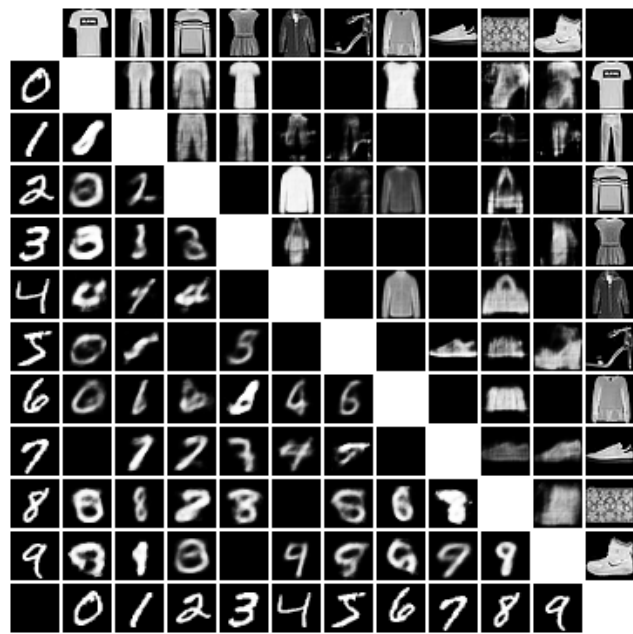
# Bayes-TrEx

## Class Boundaries

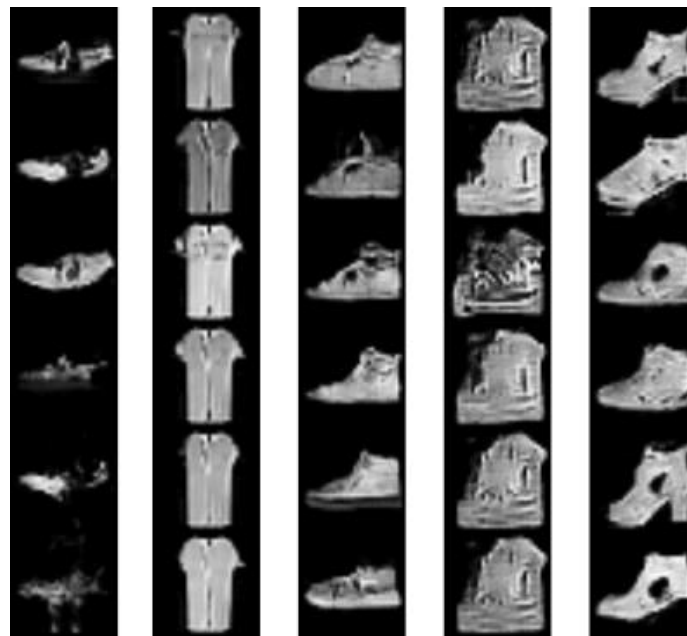


# Bayes-TrEx

Class  
Boundaries

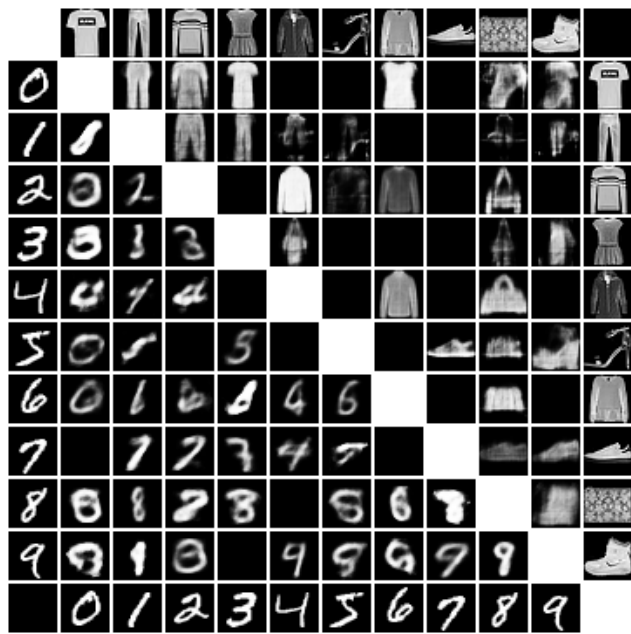


High Confidence  
Failures

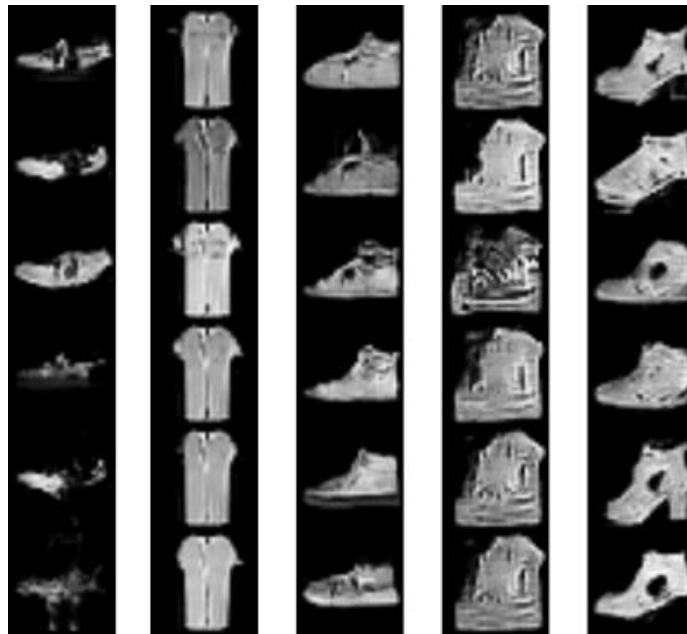


# Bayes-TrEx

Class  
Boundaries



High Confidence  
Failures

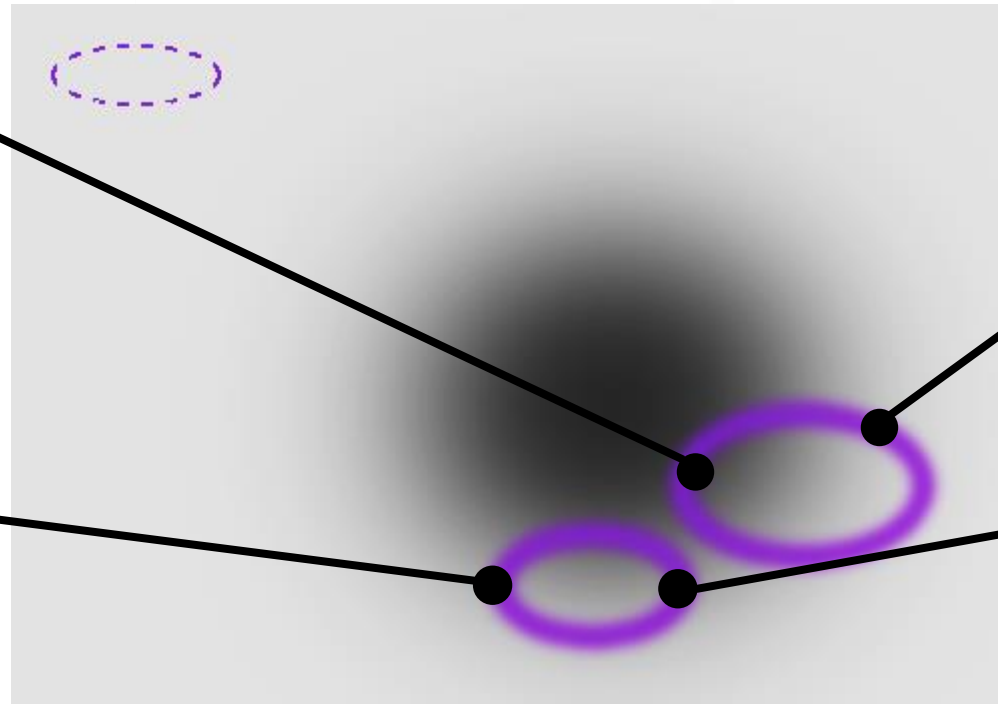
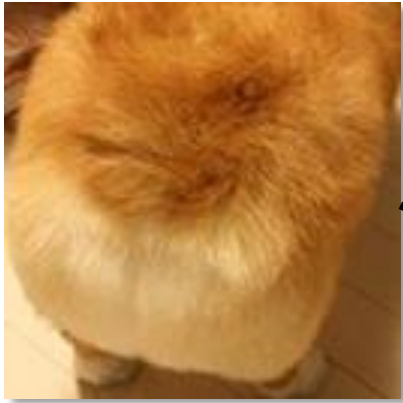


Novel  
Classes



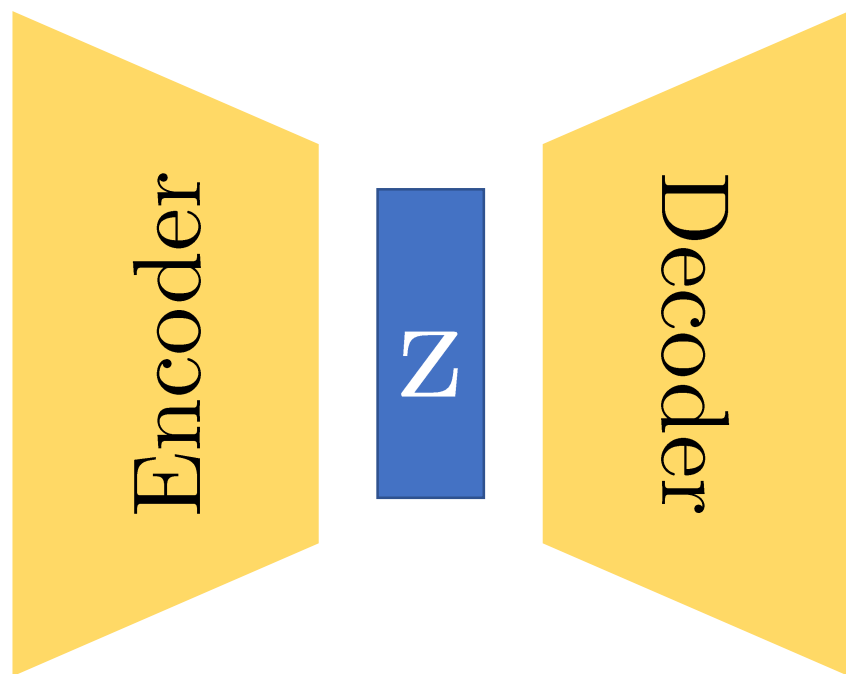
# Limitation 1: Trends & Coverage Measures

$P(\vec{y} = \text{Corgi}) = 0.5$  Level Set  
Relaxed Formulation





# Limitation 2: Latent Space Dimensionality



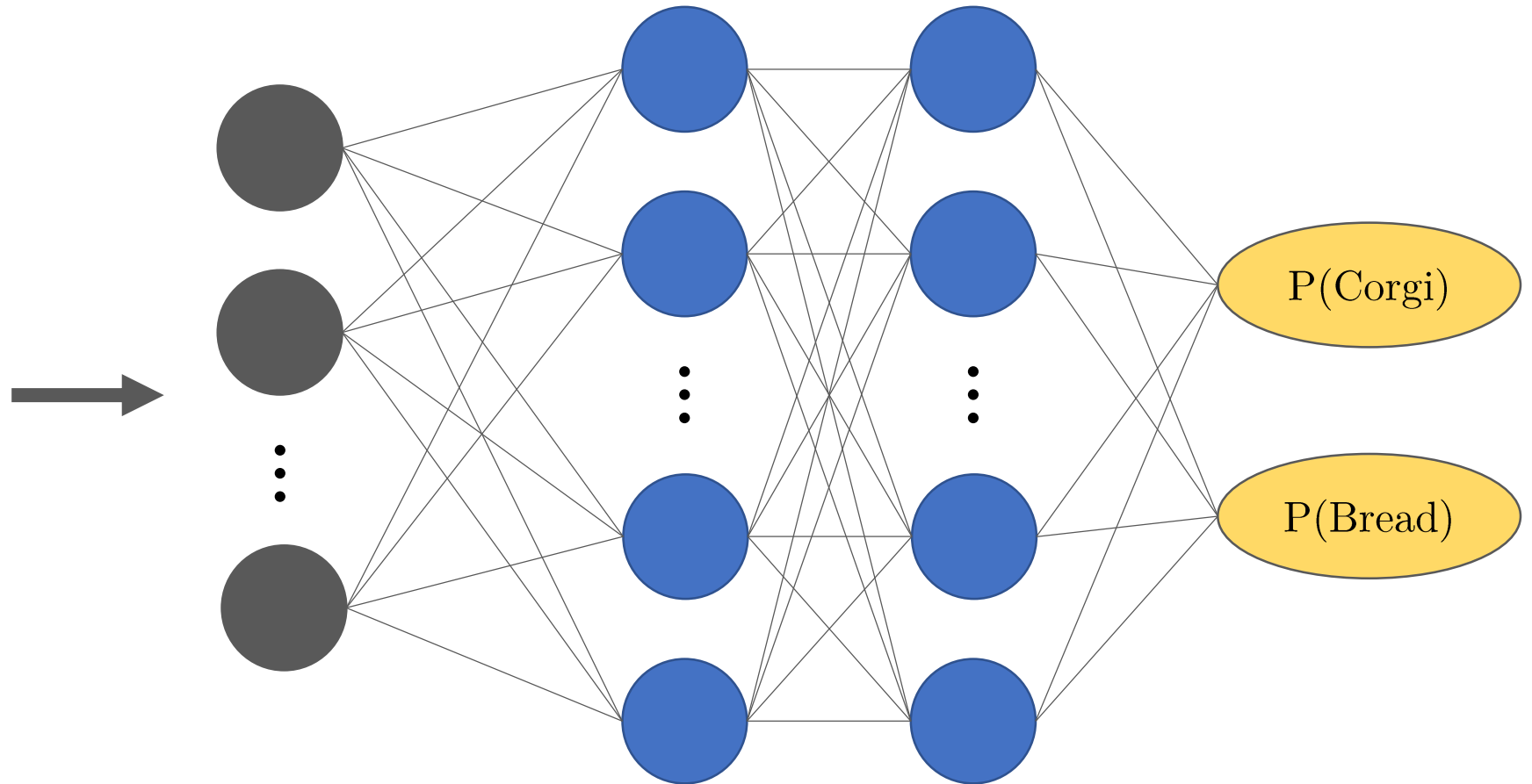
Latent: 10-dim

```
Object 1: {  
  shape: {cube, sphere, cylinder}  
  color: {cyan, ..., green},  
  material: {rubber, metal},  
  x-coord: [0,1],  
  y-coord: [0,1],  
  r-coord: [0, 2pi)  
},  
...  
Object 5: {  
  shape: {cube, sphere, cylinder}  
  color: {cyan, ..., green},  
  material: {rubber, metal},  
  x-coord: [0,1],  
  y-coord: [0,1],  
  r-coord: [0, 2pi)  
}
```

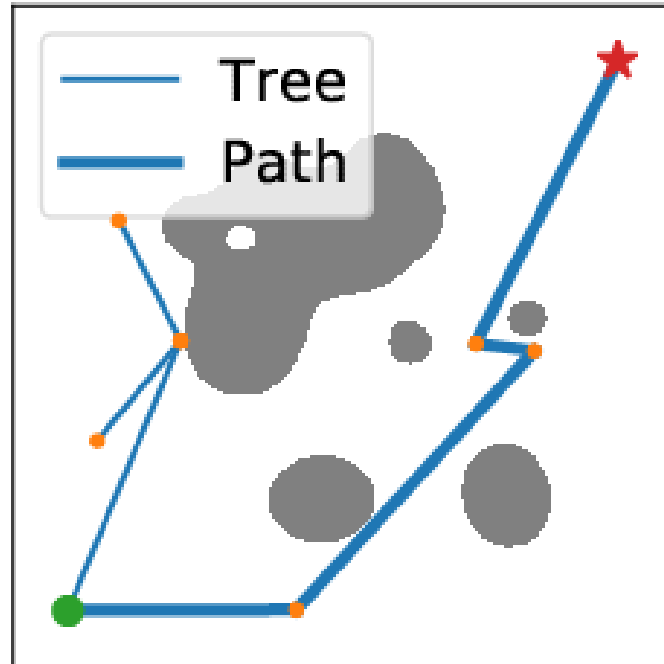
# Limitation 3: Transparency for Classification Tasks



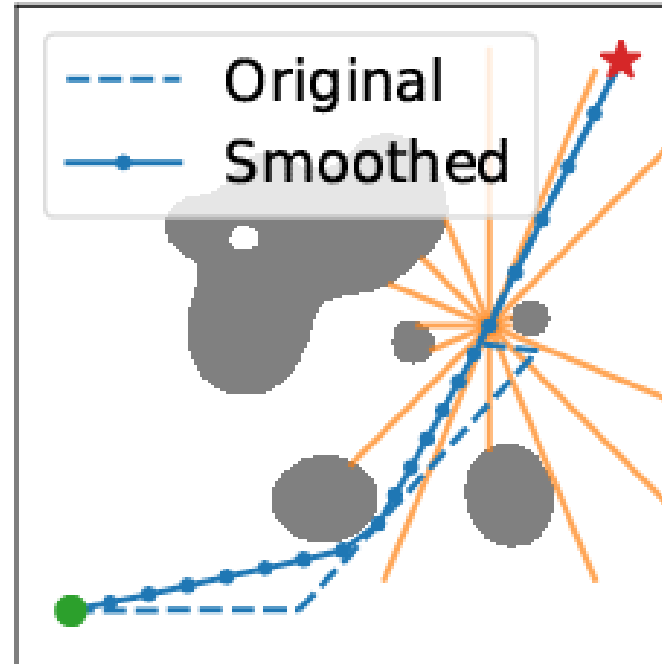
Input Image



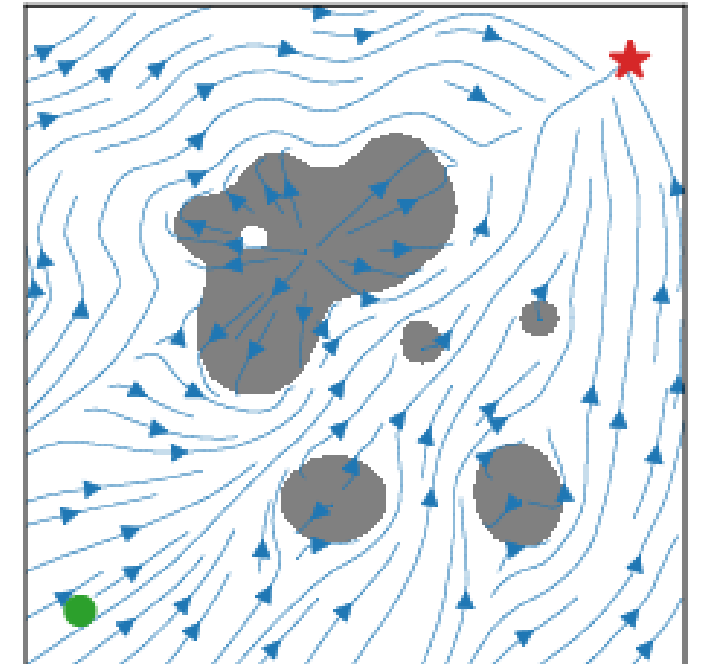
# ~~Limitation 3: Transparency for Classification Tasks~~



RRT



Smoothing and Lidar



DS Modulation

In Review & On ArXiv ([2012.13615](https://arxiv.org/abs/2012.13615)):

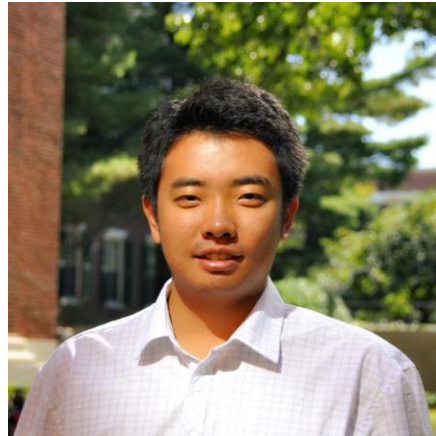
Zhou et al., *RoCUS: Robot Controller Understanding via Sampling*

# Bayes-TrEx

Paper, code, & many more experiments:  
[github.com/serenabooth/bayes-trex](https://github.com/serenabooth/bayes-trex)



Serena Booth\*  
@SerenaLBooth



Yilun Zhou\*  
yilun@mit.edu



Ankit Shah  
@ankitjs



Julie Shah  
@julie\_a\_shah

Follow up (robotics): [arxiv.org/abs/2012.13615](https://arxiv.org/abs/2012.13615)