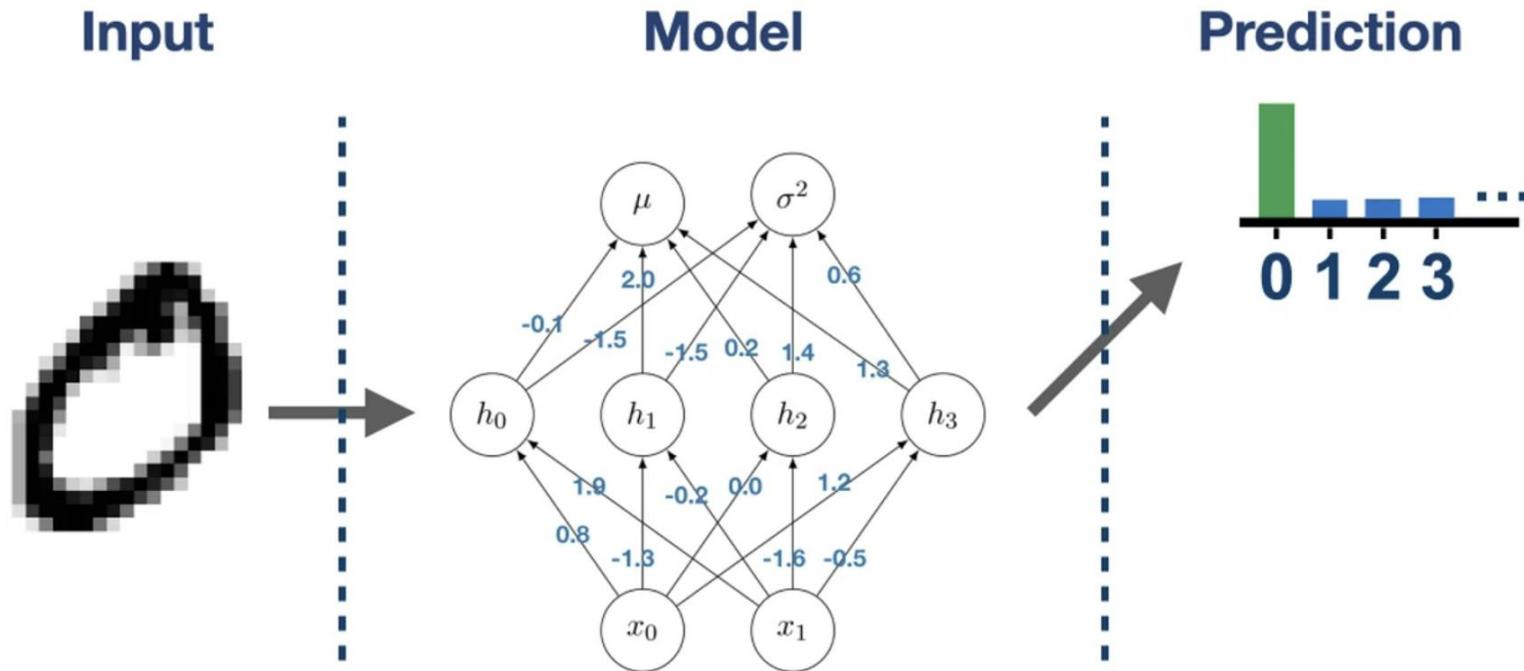# GETTING A CLUE

TANVI SHETTY

SUPERVISOR: JELLE HÜNTELMANN

# Neural Networks
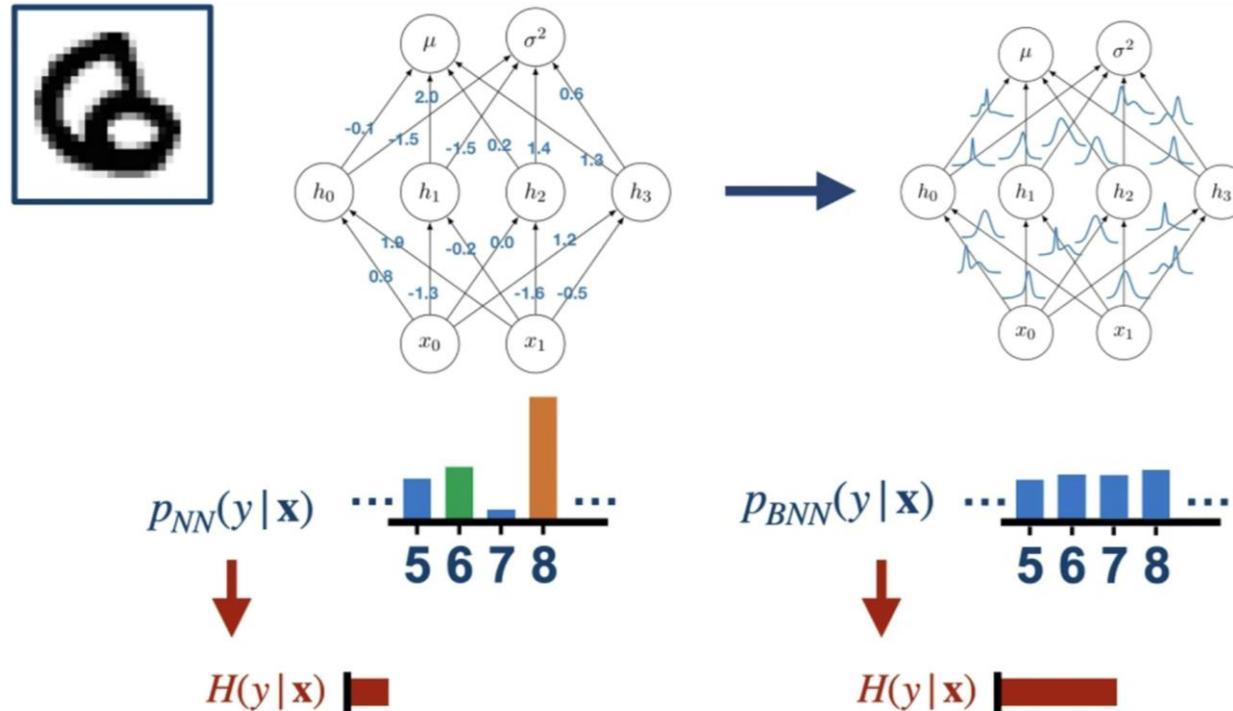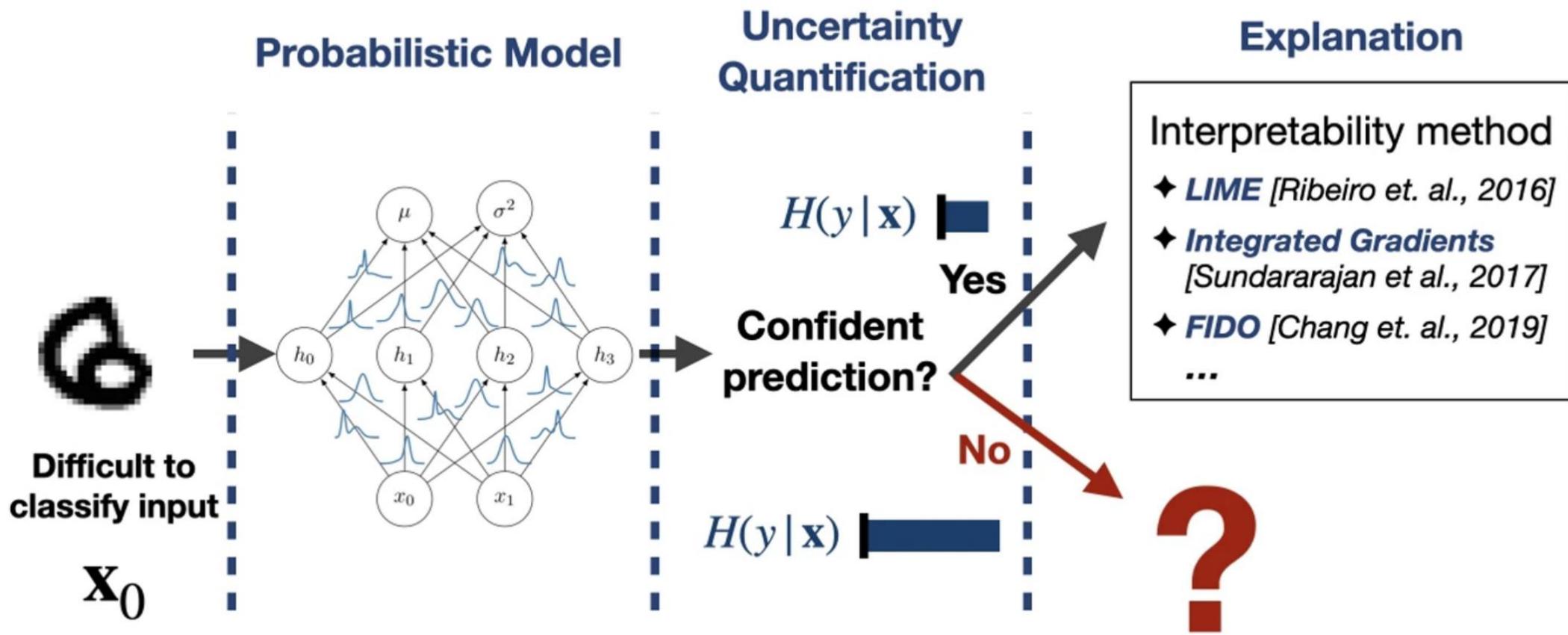
# Neural Networks



**Input**

**Model**

**Prediction**

Drawbacks
- Not easily interpretable
- Overconfident wrong predictions

3

# How to avoid overconfident predictions?

# How to capture uncertainty?

**Probabilistic Model**

**Uncertainty Quantification**

**Explanation**

Difficult to classify input

$\mathbf{x}_0$

$\mu$  $\sigma^2$

$h_0$  $h_1$  $h_2$  $h_3$

$x_0$  $x_1$

$H(y \mid \mathbf{x})$

**Confident prediction?**

**Yes**

**No**

$H(y \mid \mathbf{x})$

Interpretability method

✦ *LIME* [Ribeiro et. al., 2016]

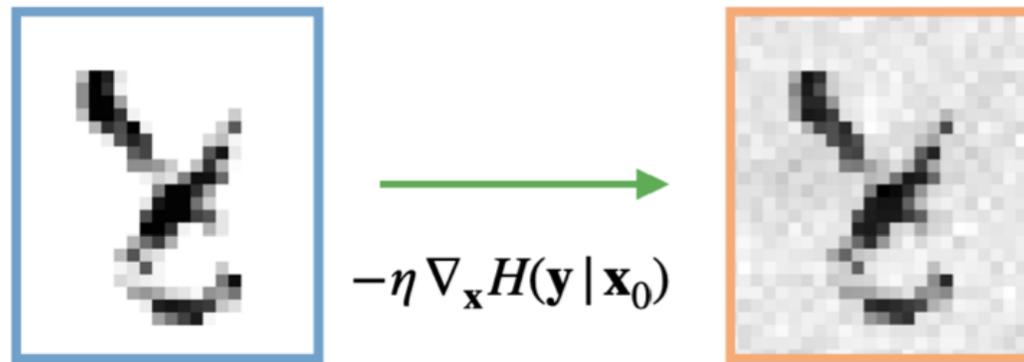✦ *Integrated Gradients* [Sundararajan et al., 2017]

✦ *FIDO* [Chang et. al., 2019]

...

?

# Why were the predictions uncertain?

# Uncertainty sensitivity analysis

- Does not scale well with high dimensional data.



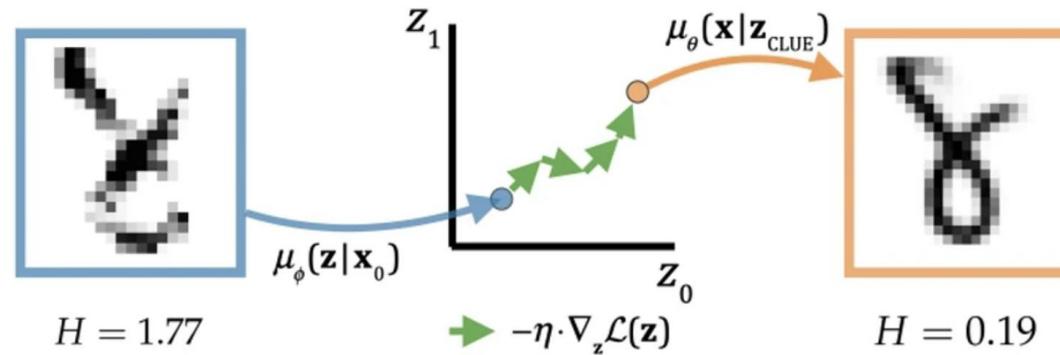$$-\eta \nabla_{\mathbf{x}} H(\mathbf{y} \,|\, \mathbf{x}_0)$$
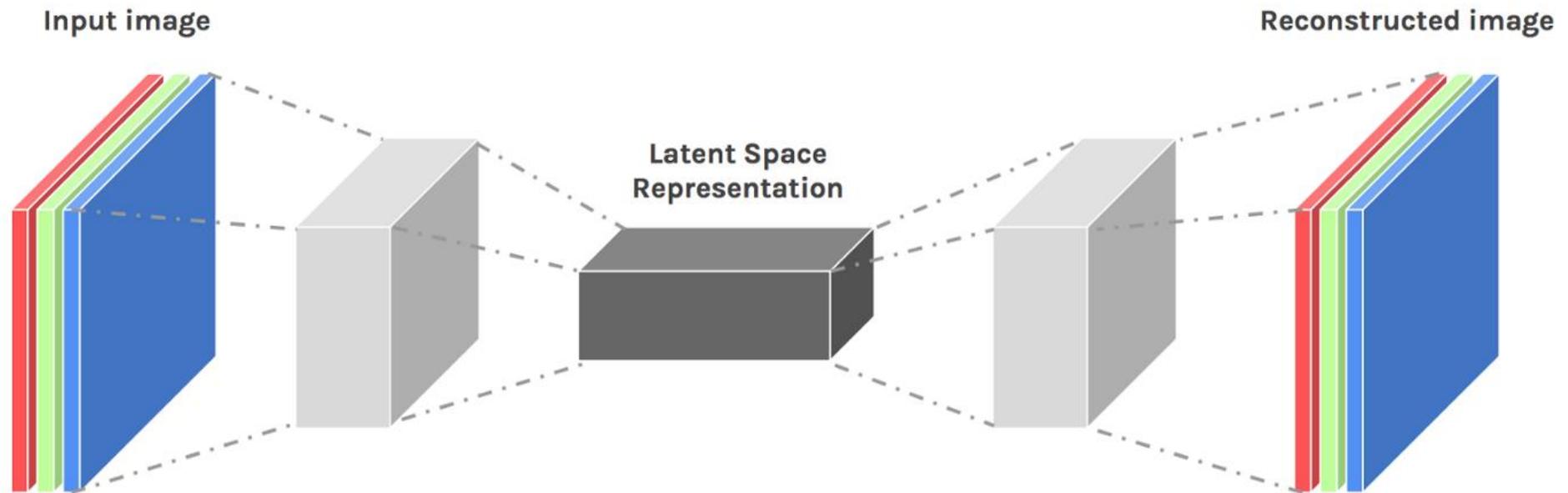
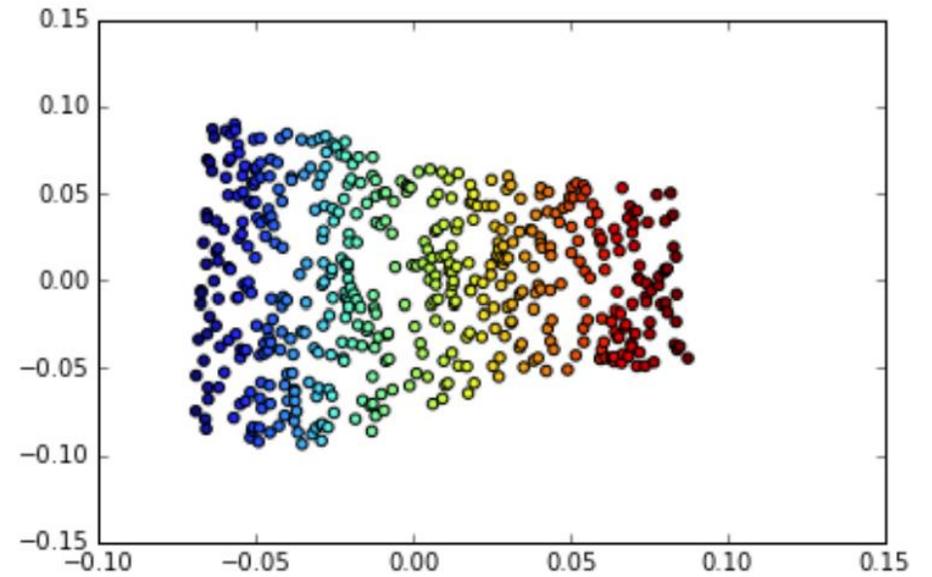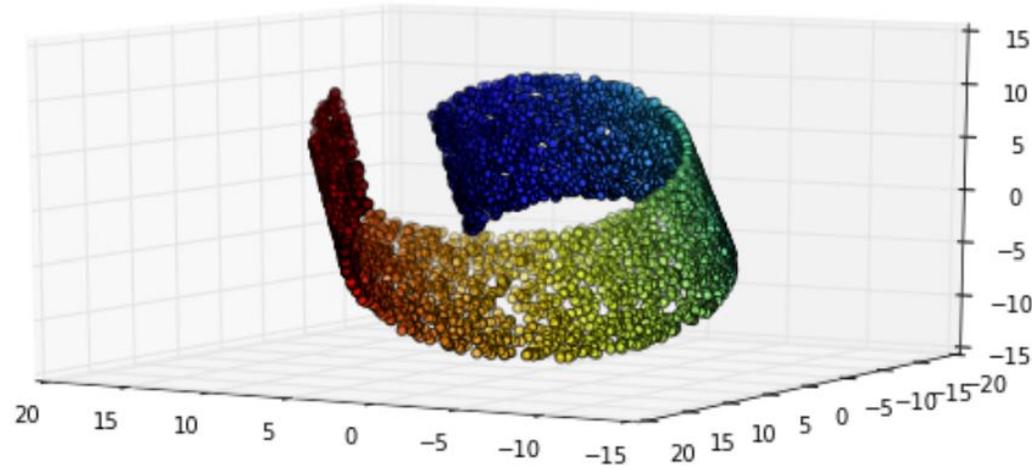# How do we overcome this?

# Main Idea

# Main Idea



- Encode our input to a latent space.

- Perform some optimization that aims to minimize uncertainty.

- Decode into some resulting input for which our model is more certain.

# Latent space



Input image

Reconstructed image

Latent Space Representation

# Manifold

# Main Idea



- Encode our input to a latent space.

- Perform some optimization that aims to minimize uncertainty.

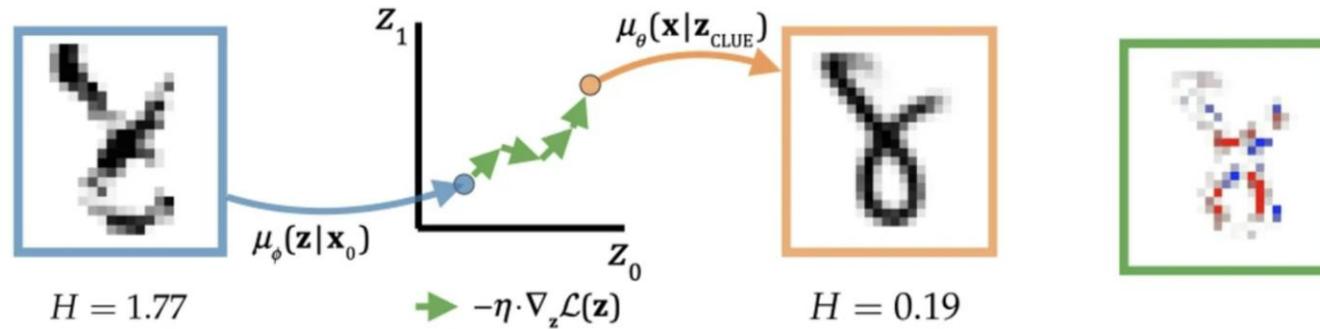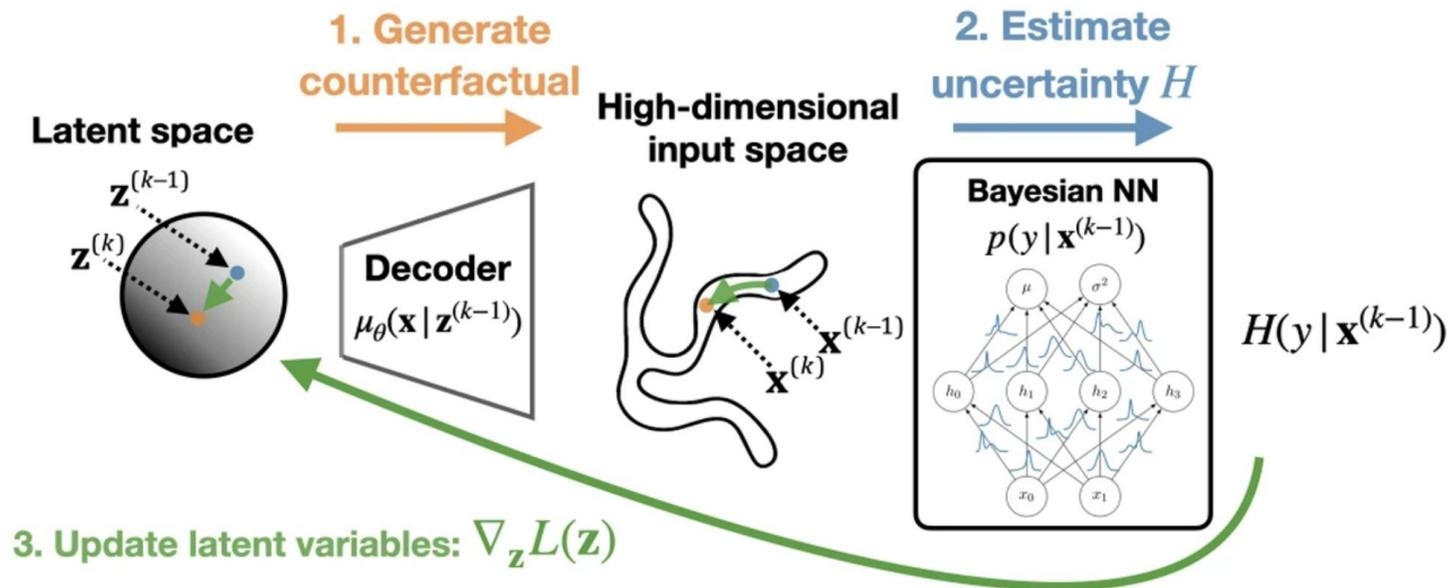- Decode into some resulting input for which our model is more certain.

# Let's get a CLUE

# Optimizing the objective

$$L(\mathbf{z}) = H(y \mid \mu_\theta(\mathbf{x} \mid \mathbf{z})) + d(\mu_\theta(\mathbf{x} \mid \mathbf{z}), \mathbf{x}_0)$$

$$\mathbf{x}_{\text{CLUE}} = \mu_\theta(\mathbf{x}|\mathbf{z}_{\text{CLUE}}) \quad \text{where} \quad \mathbf{z}_{\text{CLUE}} = \arg\min_{\mathbf{z}} \mathcal{L}(\mathbf{z})$$

# Multiplicity of CLUEs



$H=1.52,\ c=0$  $H=0.05,\ c=0$  $H=0.14,\ c=0$  $H=0.03,\ c=0$  $H=0.07,\ c=0$  $H=0.60,\ c=9$
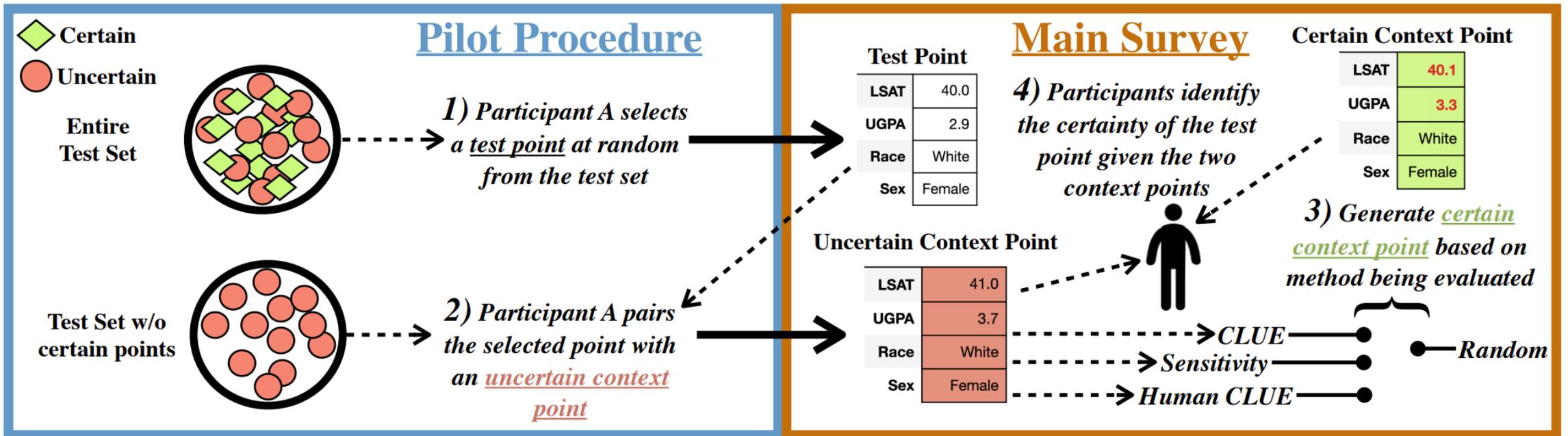
# User Study

1. Show the subject some uncertain context point.

2. Generate a counterfactual explanation for the context point.

3. Show the subject the generated certain context point using either CLUE, Uncertainty Sensitivity, Human choice, Random.

4. The user to has classify a new unseen point as certain or uncertain.
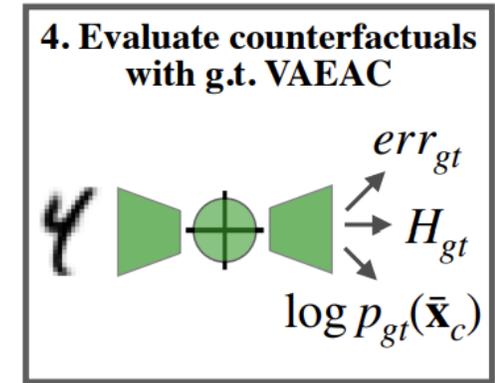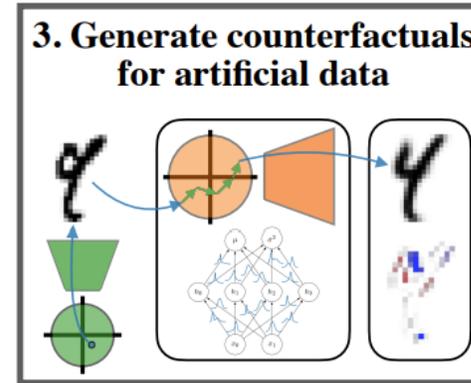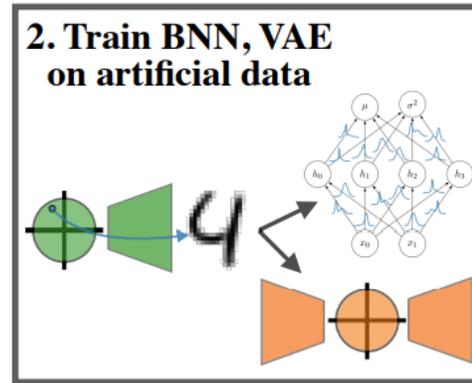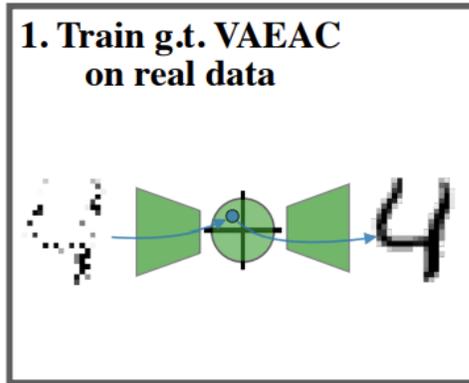
# User Study

# User Study: Results

| | Combined | LSAT | COMPAS |
|---|---|---|---|
| CLUE | **82.22** | **83.33** | **81.11** |
| *Human CLUE* | 62.22 | 61.11 | 63.33 |
| Random | 61.67 | 62.22 | 61.11 |
| Local Sensitivity | 52.78 | 56.67 | 48.89 |

# Evaluate Counterfactual Explanations

"What is the smallest change that could be made to an input,

while keeping it in distribution,

so that our model becomes certain in its decision for said input?"

# Thank you!