**Diverse, Global and Amortised Counterfactual Explanations for Uncertainty Estimates**
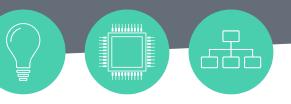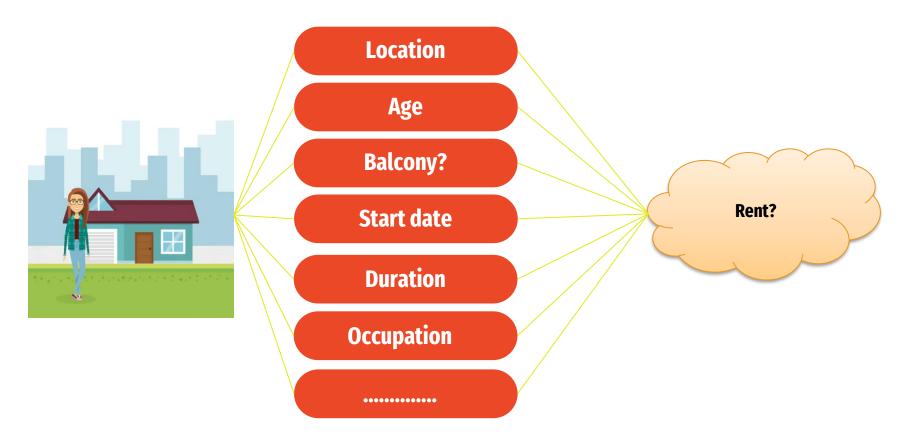
# Overview

# Counterfactual Explanations

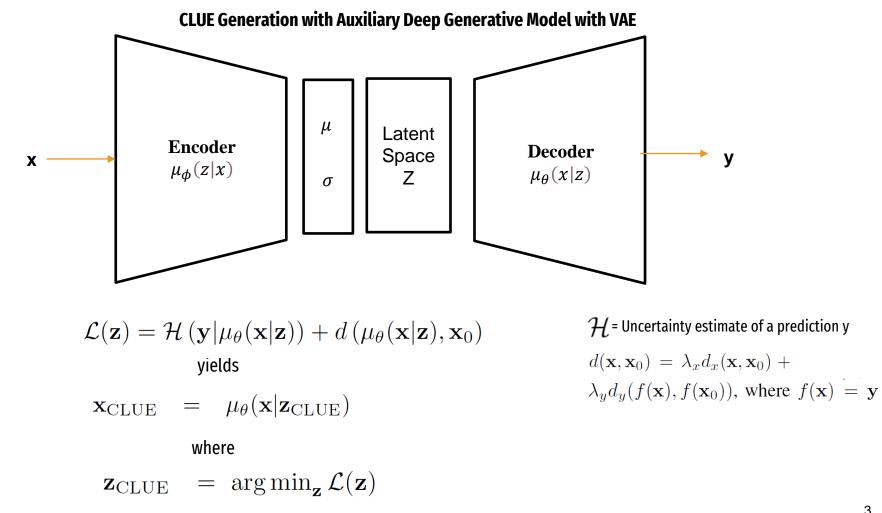# Counterfactual Latent Uncertainty Explanation (CLUE)

*What is the smallest on-manifold change that can be done to an input so that our model becomes more certain*

High Uncertainty

Incorrect prediction

More counterfactuals

Uncertainty explanations are a precedent for model explanation

# CLUE Generation with Auxiliary Deep Generative Model with VAE



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

yields

$$\mathbf{x}_{\mathrm{CLUE}} = \mu_\theta(\mathbf{x}|\mathbf{z}_{\mathrm{CLUE}})$$

where

$$\mathbf{z}_{\mathrm{CLUE}} = \arg\min_{\mathbf{z}} \mathcal{L}(\mathbf{z})$$

$\mathcal{H}$ = Uncertainty estimate of a prediction y

$$d(\mathbf{x}, \mathbf{x}_0) = \lambda_x d_x(\mathbf{x}, \mathbf{x}_0) + \lambda_y d_y(f(\mathbf{x}), f(\mathbf{x}_0)), \text{ where } f(\mathbf{x}) = \mathbf{y}$$

3

# Why this paper?

# δ-CLUE vs CLUE

## δ-CLUE

Multiplicity is achieved by searching randomly in different areas of latent space
- Sampling around an input in latent space
- Gradient descent

**Vs**

## CLUE also does this, but:

- Finds minima in a limited region of space
- Might strays far away from Counterfactuals

$$\mathbf{x}_{\delta-\mathrm{CLUE}} = \mu_\theta \left( \mathbf{x} | \mathbf{z}_{\delta-\mathrm{CLUE}} \right) \text{ where } \mathbf{z}_{\delta-\mathrm{CLUE}} = \arg\min_{\mathbf{z}: \ \rho(\mathbf{z}, \mathbf{z}_0) \leq \delta} \mathcal{L}(\mathbf{z})$$

$$\mathbf{z_0} = \mu_\phi(\mathbf{z} | \mathbf{x_0})$$
$$\rho(\mathbf{z}, \mathbf{z}_0) = \|\mathbf{z} - \mathbf{z}_0\|_2$$

**Algorithm 3:** $\delta$-CLUE

**Inputs:** $\delta$, $k$, $\mathcal{S}$, $r$, $\mathbf{x}_0$, $d$, $\rho$, $\mathcal{H}$, $\mu_\theta$, $\mu_\phi$

1   Initialise $\varnothing$ of CLUEs: $X_{\text{CLUE}} = \{\}$;
2   Set $\delta$-ball centre of $\mathbf{z}_0 = \mu_\phi(\mathbf{z}|\mathbf{x}_0)$;
3   **for** $1 \leq i \leq k$ **do**
4     Set initial value of $\mathbf{z}_i = \mathcal{S}(\mathbf{z}_0, r, i, k)$;
5     **while** *loss $\mathcal{L}$ has not converged* **do**
6       Decode: $\mathbf{x} = \mu_\theta(\mathbf{x}|\mathbf{z}_i)$;
7       Use predictor to obtain $\mathcal{H}(\mathbf{y}|\mathbf{x})$;
8       $\mathcal{L} = \mathcal{H}(\mathbf{y}|\mathbf{x}) + d(\mathbf{x}, \mathbf{x}_0)$;
9       Update $\mathbf{z}_i$ with $\nabla_{\mathbf{z}}\mathcal{L}$;
10      **if** $\rho(\mathbf{z}_i, \mathbf{z}_0) > \delta$ **then**
11        Project $\mathbf{z}_i$ onto the surface of the $\delta$-ball as $\mathbf{z}_i = \delta \times \frac{\mathbf{z}_i - \mathbf{z}_0}{\rho(\mathbf{z}_i, \mathbf{z}_0)}$;
12      **end if**
13     **end while**
14     Decode explanation: $\mathbf{x}_{\delta-\text{CLUE}} = \mu_\theta(\mathbf{x}|\mathbf{z}_i)$;
15     **if** $\mathcal{H}(\mathbf{y}|\mathbf{x}_{\delta_{\text{CLUE}}}) < \mathcal{H}_{\text{threshold}}$ **then**
16      $X_{\text{CLUE}} \leftarrow X_{\text{CLUE}} \cup \mathbf{x}_{\delta_{\text{CLUE}}}$;
17     **end if**
18   **end for**

**Outputs:** $X_{\text{CLUE}}$, a set of $n \leq k$ CLUEs



ORIGINAL CLUE     $\delta$-CLUE

High $\mathcal{L}(\mathbf{z})$    Low $\mathcal{L}(\mathbf{z})$

# Different trials on δ-CLUE Algorithm

**01** — Range of δ values from 0.5 to 3.5
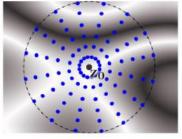
**02** — Two latent space loss functions:
- Uncertainty : $\mathcal{L}_{\mathcal{H}} = \mathcal{H}$
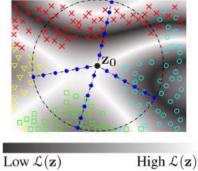- Distance : $\mathcal{L}_{\mathcal{H}+d} = \mathcal{H} + d$

**03** — Two initialisation schemes like:
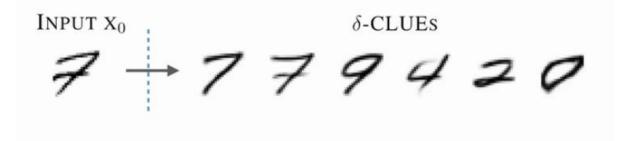- Radially Uniform
- Nearest Neighbour



$\mathcal{S}_1$: RADIALLY UNIFORM

$\mathcal{S}_2$: NEAREST NEIGHBOUR PATH

Low $\mathcal{L}(\mathbf{z})$    High $\mathcal{L}(\mathbf{z})$

Source: Ley, Dan, Umang Bhatt, and Adrian Weller. "Diverse, Global and Amortised Counterfactual Explanations for Uncertainty Estimates." *Proceedings of the AAAI Conference on Artificial Intelligence*. Vol. 36. No. 7. 2022.

INPUT $X_0$

INPUT $X_0$      $\delta$-CLUEs

Source: Ley, Dan, Umang Bhatt, and Adrian Weller. "{\delta}-CLUE: Diverse Sets of Explanations for Uncertainty Estimates." *arXiv preprint arXiv:2104.06323* (2021).

# Uncertainty vs Distance Trade-off



The hyperparameters $(\lambda_x, \lambda_y)$ controls this trade-off

Source: Ley, Dan, Umang Bhatt, and Adrian Weller. "Diverse, Global and Amortised Counterfactual Explanations for Uncertainty Estimates." *Proceedings of the AAAI Conference on Artificial Intelligence*. Vol. 36. No. 7. 2022.

# Diversity Metrics (D)

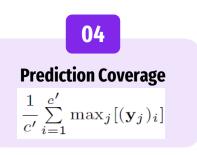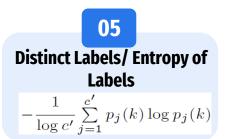**01**

### Determinantal Point Process

$$\det(\mathbf{K}) \text{ where } \mathbf{K}_{i,j} = \frac{1}{1 + d(\mathbf{x}_i, \mathbf{x}_j)}$$

**02**

### Average Pairwise Distance

$$\frac{1}{\binom{k}{2}} \sum_{i=1}^{k-1} \sum_{j=i+1}^{k} d(\mathbf{x}_i, \mathbf{x}_j)$$

**03**

### Coverage

$$\frac{1}{d'} \sum_{i=1}^{d'} \left( \max_j (\mathbf{x}_j - \mathbf{x}_0)_i + \max_j (\mathbf{x}_0 - \mathbf{x}_j)_i \right)$$

**04**

### Prediction Coverage

$$\frac{1}{c'} \sum_{i=1}^{c'} \max_j [(\mathbf{y}_j)_i]$$

**05**

### Distinct Labels / Entropy of Labels
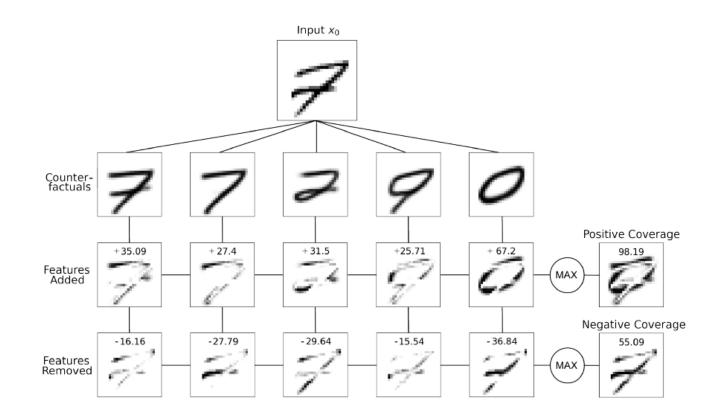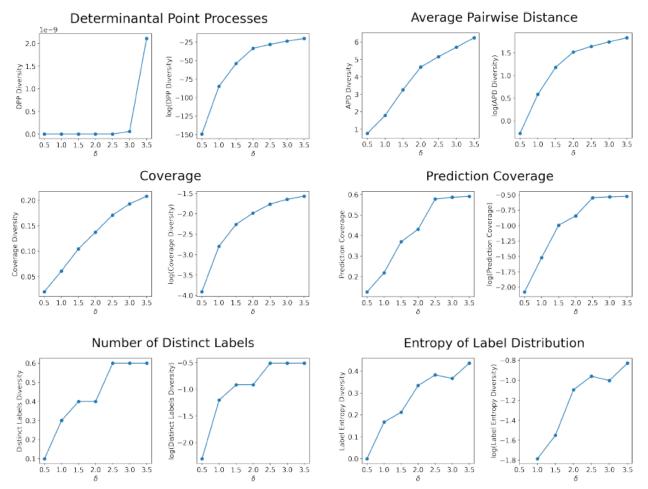
$$-\frac{1}{\log c'} \sum_{j=1}^{c'} p_j(k) \log p_j(k)$$
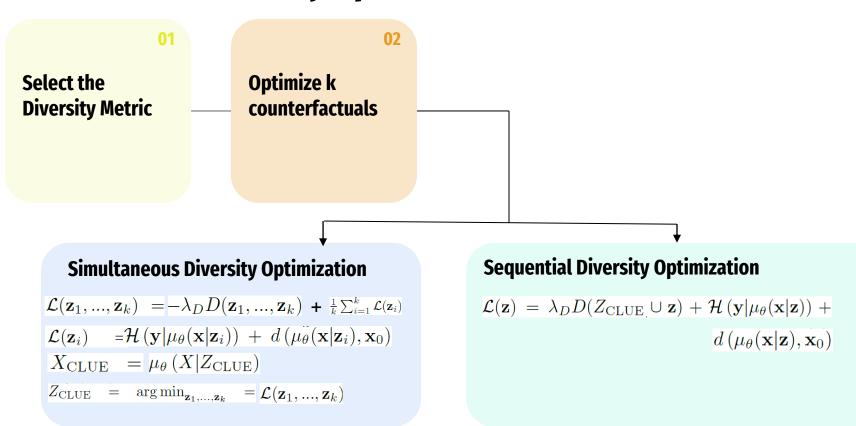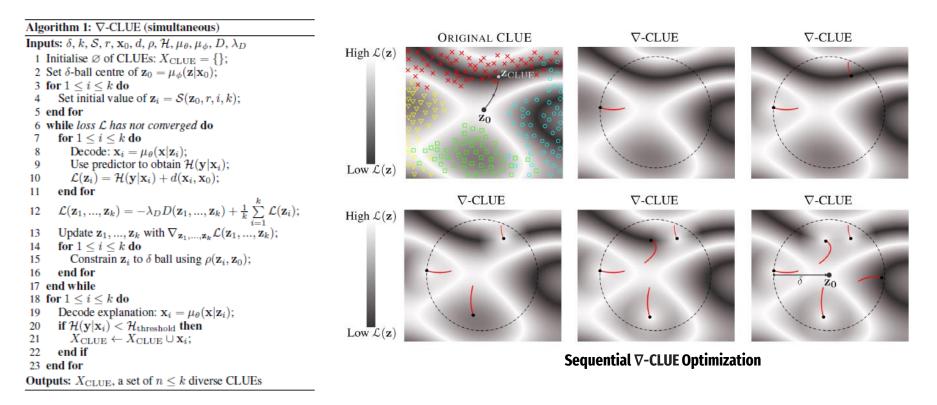
# Coverage as a Metric



Source: Ley, Dan, Umang Bhatt, and Adrian Weller. "Diverse, Global and Amortised Counterfactual Explanations for Uncertainty Estimates." *Proceedings of the AAAI Conference on Artificial Intelligence*. Vol. 36. No. 7. 2022.

# Diversity Optimization : $\nabla$-CLUE

**Algorithm 1: ∇-CLUE (simultaneous)**

**Inputs:** $\delta, k, \mathcal{S}, r, \mathbf{x}_0, d, \rho, \mathcal{H}, \mu_\theta, \mu_\phi, D, \lambda_D$

1 Initialise $\varnothing$ of CLUEs: $X_{\text{CLUE}} = \{\}$;
2 Set $\delta$-ball centre of $\mathbf{z}_0 = \mu_\phi(\mathbf{z}|\mathbf{x}_0)$;
3 **for** $1 \leq i \leq k$ **do**
4     Set initial value of $\mathbf{z}_i = \mathcal{S}(\mathbf{z}_0, r, i, k)$;
5 **end for**
6 **while** *loss $\mathcal{L}$ has not converged* **do**
7     **for** $1 \leq i \leq k$ **do**
8         Decode: $\mathbf{x}_i = \mu_\theta(\mathbf{x}|\mathbf{z}_i)$;
9         Use predictor to obtain $\mathcal{H}(\mathbf{y}|\mathbf{x}_i)$;
10         $\mathcal{L}(\mathbf{z}_i) = \mathcal{H}(\mathbf{y}|\mathbf{x}_i) + d(\mathbf{x}_i, \mathbf{x}_0)$;
11     **end for**
12     $\mathcal{L}(\mathbf{z}_1, ..., \mathbf{z}_k) = -\lambda_D D(\mathbf{z}_1, ..., \mathbf{z}_k) + \frac{1}{k}\sum_{i=1}^{k}\mathcal{L}(\mathbf{z}_i)$;
13     Update $\mathbf{z}_1, ..., \mathbf{z}_k$ with $\nabla_{\mathbf{z}_1,...,\mathbf{z}_k}\mathcal{L}(\mathbf{z}_1, ..., \mathbf{z}_k)$;
14     **for** $1 \leq i \leq k$ **do**
15         Constrain $\mathbf{z}_i$ to $\delta$ ball using $\rho(\mathbf{z}_i, \mathbf{z}_0)$;
16     **end for**
17 **end while**
18 **for** $1 \leq i \leq k$ **do**
19     Decode explanation: $\mathbf{x}_i = \mu_\theta(\mathbf{x}|\mathbf{z}_i)$;
20     **if** $\mathcal{H}(\mathbf{y}|\mathbf{x}_i) < \mathcal{H}_{\text{threshold}}$ **then**
21         $X_{\text{CLUE}} \leftarrow X_{\text{CLUE}} \cup \mathbf{x}_i$;
22     **end if**
23 **end for**

**Outputs:** $X_{\text{CLUE}}$, a set of $n \leq k$ diverse CLUEs
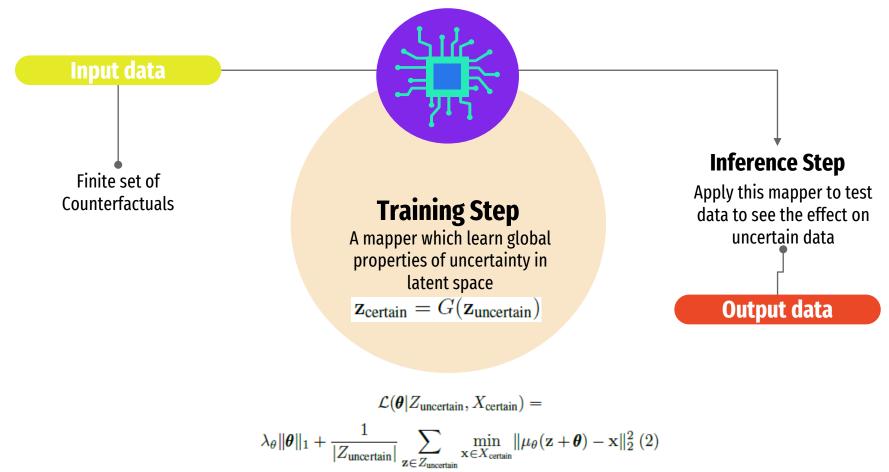


Sequential ∇-CLUE Optimization

Source: Ley, Dan, Umang Bhatt, and Adrian Weller. "Diverse, Global and Amortised Counterfactual Explanations for Uncertainty Estimates." *Proceedings of the AAAI Conference on Artificial Intelligence*. Vol. 36. No. 7. 2022.

# GLobal AMortised CLUE (GLAM-CLUE)



**Input data**

Finite set of Counterfactuals

**Training Step**

A mapper which learn global properties of uncertainty in latent space

$$\mathbf{z}_{\text{certain}} = G(\mathbf{z}_{\text{uncertain}})$$

**Inference Step**

Apply this mapper to test data to see the effect on uncertain data

**Output data**

$$\mathcal{L}(\boldsymbol{\theta}|Z_{\text{uncertain}}, X_{\text{certain}}) =$$

$$\lambda_\theta \|\boldsymbol{\theta}\|_1 + \frac{1}{|Z_{\text{uncertain}}|} \sum_{\mathbf{z} \in Z_{\text{uncertain}}} \min_{\mathbf{x} \in X_{\text{certain}}} \|\mu_\theta(\mathbf{z} + \boldsymbol{\theta}) - \mathbf{x}\|_2^2 \ (2)$$

15

**Algorithm 2:** GLAM-CLUE (Training Step)

**Inputs:** Inputs $X_{\text{uncertain}}$, $X_{\text{certain}}$, groups $Y_{\text{uncertain}}$, $Y_{\text{certain}}$, DGM encoder $\mu_\phi$, loss $\mathcal{L}$, trainable parameters $\boldsymbol{\theta}$

1  **for all** groups $(i \to j)$ in $(Y_{\text{uncertain}}, Y_{\text{certain}})$ **do**
2      Select $X_i$ from $X_{\text{uncertain}}, Y_{\text{uncertain}}$;
3      Select $X_j$ from $X_{\text{certain}}, Y_{\text{certain}}$;
4      Encode: $Z_i = \mu_\phi(Z|X_i)$;
5      **while** *loss $\mathcal{L}$ has not converged* **do**
6          Update $\boldsymbol{\theta}_{i \to j}$ with $\nabla_{\boldsymbol{\theta}_{i \to j}}\mathcal{L}(\boldsymbol{\theta}_{i \to j}|Z_i, X_j)$;
7      **end while**
8  **end for**

**Outputs:** A collection of mapping parameters $\boldsymbol{\theta}_{i \to j}$ for given mappers $G_{i \to j}$ that take uncertain inputs from group $i$ and produce nearby certain outputs in group $j$

$$\mathbf{z}_j = G_{i \to j}(\mathbf{z}_i) = \mathbf{z}_i + \boldsymbol{\theta}_{i \to j}$$

Source: Ley, Dan, Umang Bhatt, and Adrian Weller. "Diverse, Global and Amortised Counterfactual Explanations for Uncertainty Estimates." *Proceedings of the AAAI Conference on Artificial Intelligence*. Vol. 36. No. 7. 2022.

# Performance Test

**01 Difference Between Means (DBM)**

Uncertain data to certain data in input or latent space

**02 Nearest Neighbours (NN)**

Used in high certainty training data in input or latent space



Latent DBM Mapping On Unseen Test Data

- certain test data
- uncertain test data
- latent DBM mapping

Latent NN Mapping On Unseen Test Data

- certain test data
- uncertain test data
- latent NN mapping

Source: Ley, Dan, Umang Bhatt, and Adrian Weller. "Diverse, Global and Amortised Counterfactual Explanations for Uncertainty Estimates." *Proceedings of the AAAI Conference on Artificial Intelligence*. Vol. 36. No. 7. 2022.

# Performance Comparison

| Input DBM | Latent DBM | Input NN |
|-----------|------------|----------|
| 0.0306 | 0.0262 | 0.0236 |

| Latent NN | GLAM-CLUE | CLUE |
|-----------|-----------|------|
| 0.0245 | 0.0238 | 4.68 |

GLAM-CLUE outperforms these baselines
......*almost* **200** *times faster*

Source: Ley, Dan, Umang Bhatt, and Adrian Weller. "Diverse, Global and Amortised Counterfactual Explanations for Uncertainty Estimates." *Proceedings of the AAAI Conference on Artificial Intelligence*. Vol. 36. No. 7. 2022.

# Future Work

**01** **Data set dimensions**
Using higher dimensional data set

**02** **Introduce different metric**
Use FID scores to replace simple distance metric in evaluation and optimisation

**03** **Use different DGMs**
Use DGM alternative like GANs instead of VAEs

# Conclusion

**01**      **Making CLUE more useful in practice**

**02**      **Proposed δ-CLUE and ∇-CLUE to tackle the multiplicity and diversity issues**

**03**      **Introduced GLAM-CLUE which tackles the computational inefficiency caused on**

             **large data sets with δ-CLUE and ∇-CLUE**

# References

- https://towardsdatascience.com/understanding-variational-autoencoders-vaes-f70510919f73

- Antorán, Javier, et al. "Getting a clue: A method for explaining uncertainty estimates." *arXiv preprint arXiv:2006.06848* (2020).

- Ley, Dan, Umang Bhatt, and Adrian Weller. "{\delta}-CLUE: Diverse Sets of Explanations for Uncertainty Estimates." *arXiv preprint arXiv:2104.06323* (2021).

- https://slideslive.com/38955757/deltaclue-diverse-sets-of-explanations-for-uncertainty-estimates?ref=recommended

# THANK YOU

**Sruthi Aikkara**
Matriculation number: **229386**

Supervisor: **Jelle Hüntelmann**