

# CLUE: A Method for Explaining Uncertainty Estimates

Javier Antorán, Umang Bhatt, Tameem Adel, Adrian Weller, José Miguel Hernández-Lobato

{ja666, usb20, tah47, aw665, jmh233}@cam.ac.uk



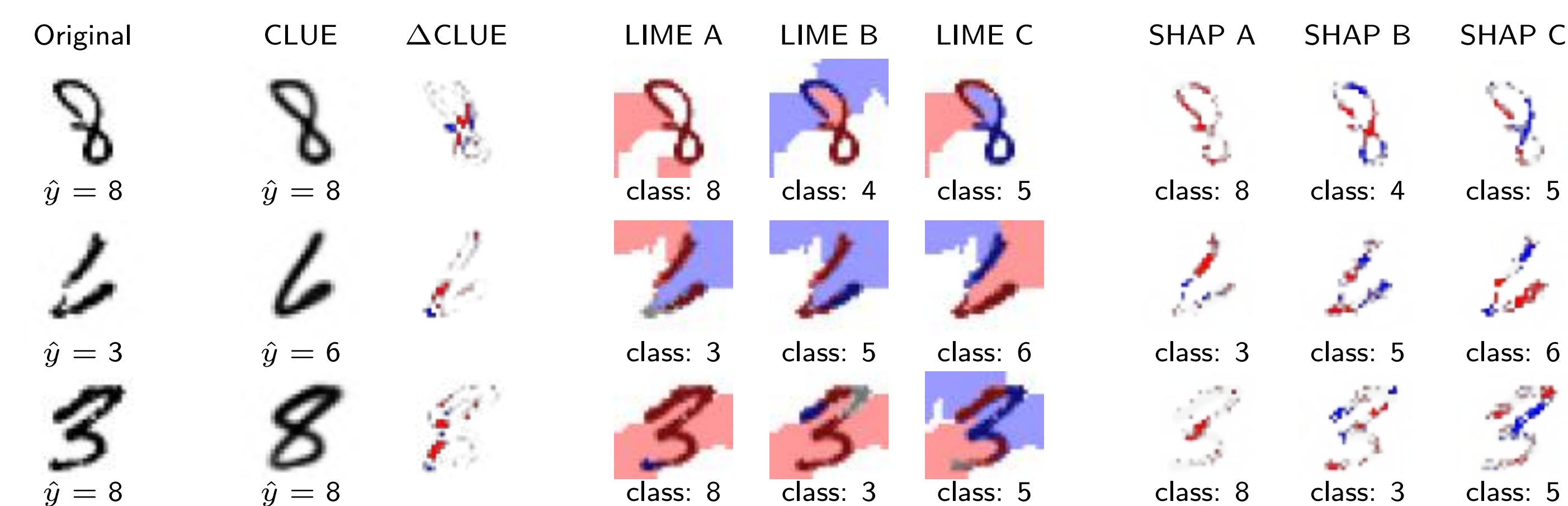
## Intersecting Explainability and Uncertainty in ML

Can we interrogate probabilistic models (e.g. BNNs) about their reasoning when they are uncertain?

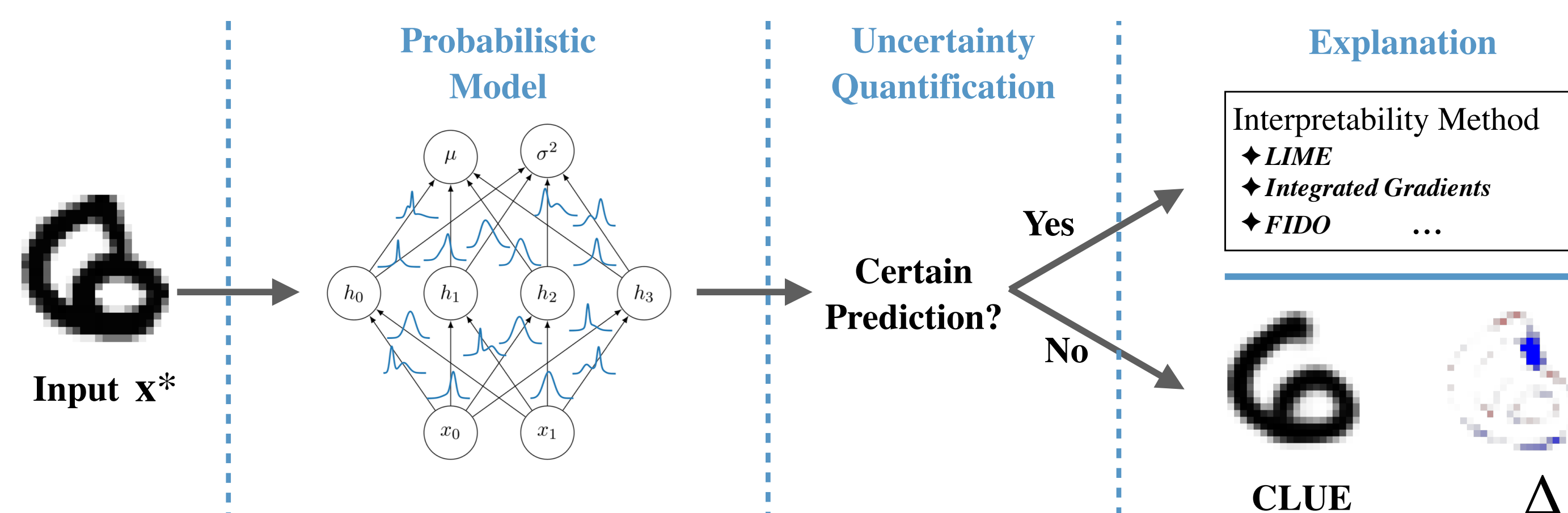
- Help ML practitioners identify input regions where training data is sparse; e.g. identify under-represented groups (by age, gender, etc.)
- Direct end-users' (e.g. medical professionals) attention to anomalous characteristics in uncertain inputs.

## Applying Explainability Techniques to Uncertainty

- **Uncertainty Sensitivity Analysis [1]:** Gradient-based approach. In high dimensions can take points off data manifold.
- **Feature Importance Approaches:** highlight features which provide evidence for a given class. What if there is conflicting evidence for multiple classes or a lack of evidence for any class?



## A Proposed Workflow



## References

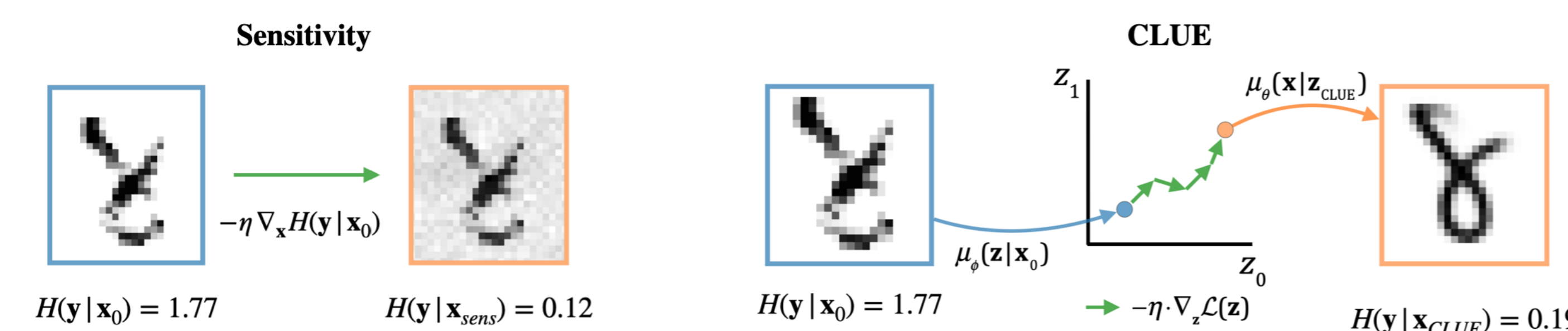
[1] S. Depeweg, J. M. Hernández-Lobato, S. Udluft, and T. A. Runkler. Sensitivity analysis for predictive uncertainty. In *ESANN*, 2017.

## Counterfactual Latent Uncertainty Explanations

We ensure our counterfactuals represent plausible inputs by introducing an auxiliary latent variable DGM. **CLUE finds points in latent space which generate inputs similar to an original observation  $x_0$  but are assigned low uncertainty  $\mathcal{H}$ :**

$$\mathcal{L}(z) = \mathcal{H}(y|\mu_\theta(x|z)) + d(\mu_\theta(x|z), x_0), \quad (1)$$

$$x_{\text{CLUE}} = \mu_\theta(x|z_{\text{CLUE}}) \text{ where } z_{\text{CLUE}} = \arg \min_z \mathcal{L}(z). \quad (2)$$



### Algorithm 1: CLUE

**Inputs:** datapoint  $x_0$ , distance  $d(\cdot)$ , Uncertainty estimator  $\mathcal{H}$ , DGM decoder  $\mu_\theta(\cdot)$ , DGM encoder  $\mu_\phi(\cdot)$

- 1 Set initial value  $z = \mu_\phi(x|x_0)$ ;
  - 2 **while** loss  $\mathcal{L}$  is not converged **do**
  - 3   Decode:  $x = \mu_\theta(x|z)$ ;
  - 4   Use predictor to obtain  $\mathcal{H}(y|x)$ ;
  - 5    $\mathcal{L} = \mathcal{H}(y|x) + d(x, x_0)$ ;
  - 6   Update  $z$  with  $\nabla_z \mathcal{L}$ ;
  - 7 **end**
- Output:**  $x_{\text{CLUE}} = \mu_\theta(x|z)$

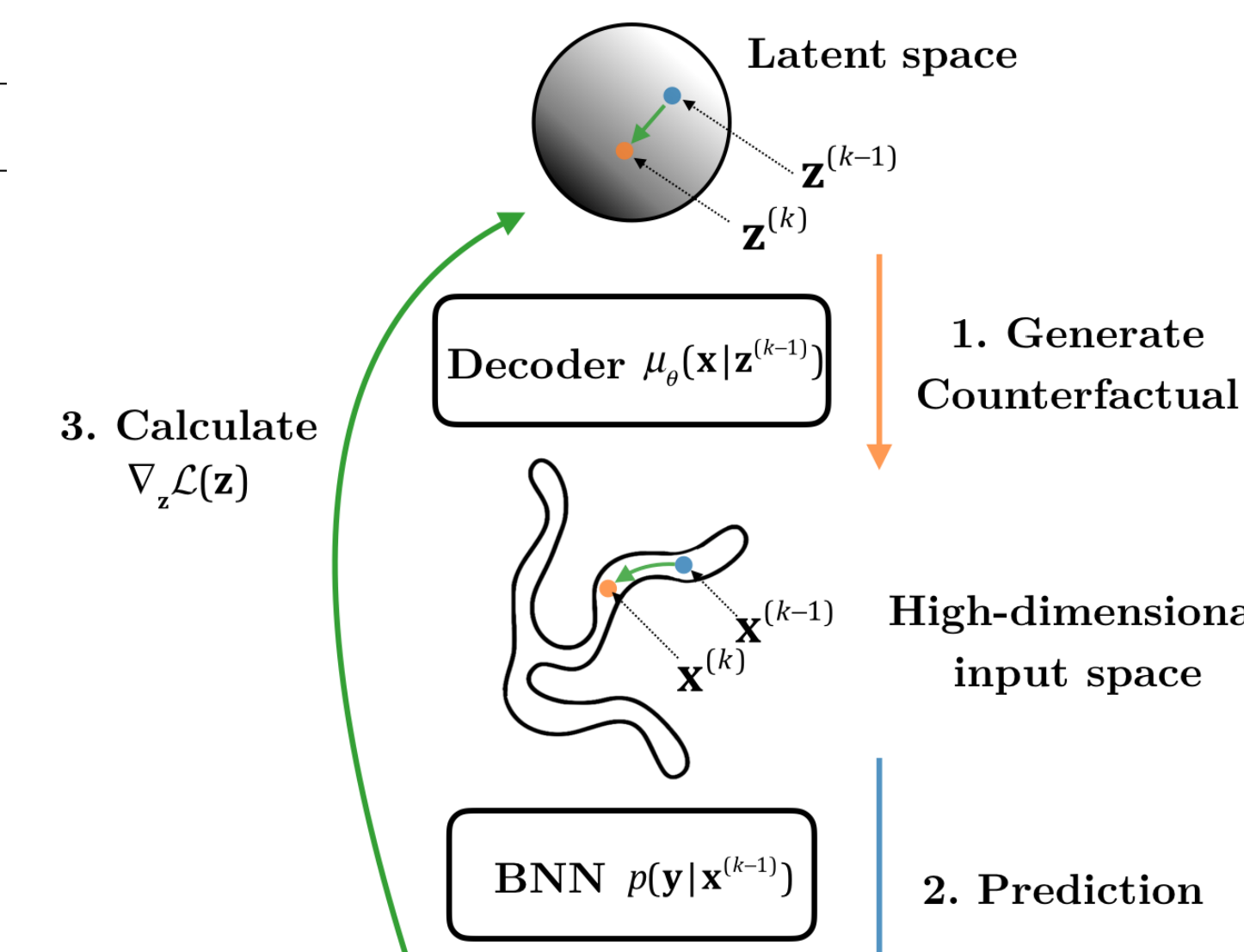
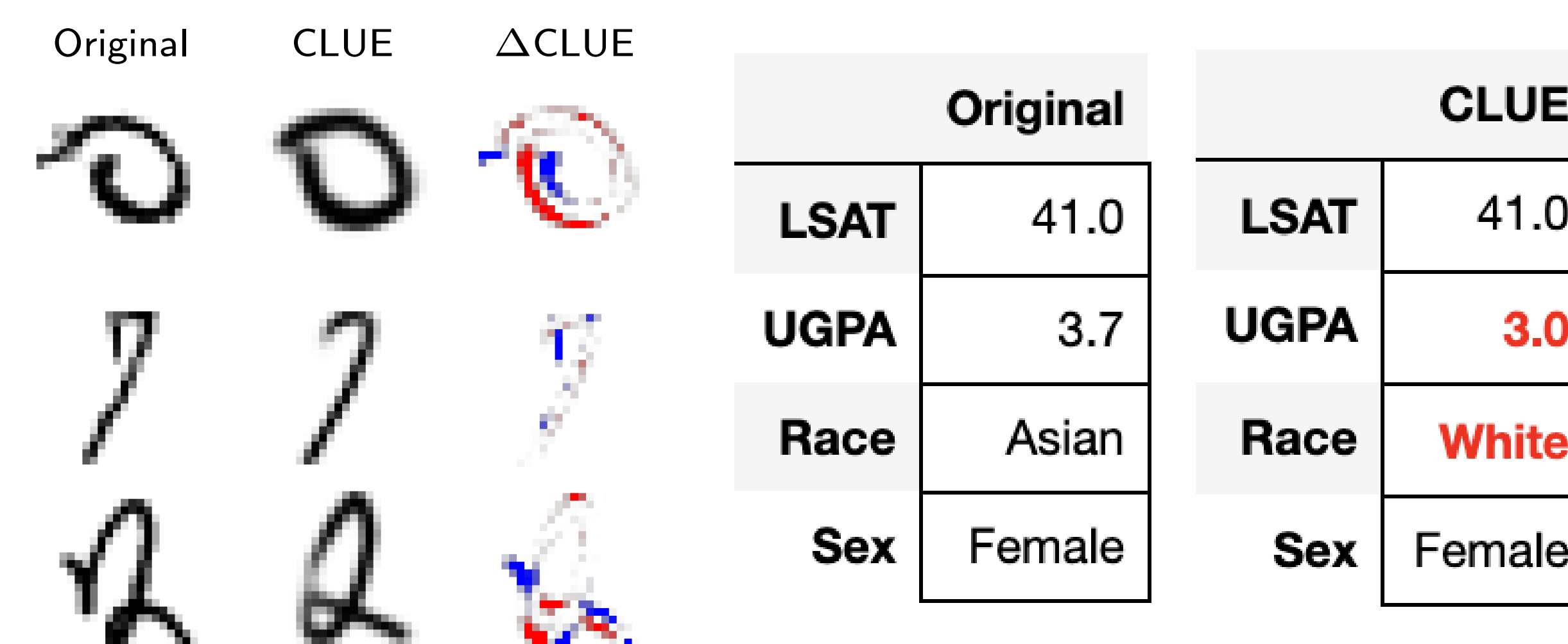


Figure: Latent codes are decoded into inputs for which a BNN generates uncertainty estimates; Our objective's gradients are back-propagated to latent space.

## Explaining Uncertainty with Counterfactuals

CLUEs can be generated for classification and regression tasks on both tabular and image data. We highlight changes  $\Delta_{\text{CLUE}} = (x_{\text{CLUE}} - x_0)^2$ .



(a) MNIST

(b) LSAT

Figure: Example image and tabular CLUEs together with their corresponding  $\Delta_{\text{CLUE}}$ .

## Diversity in Plausible Counterfactuals

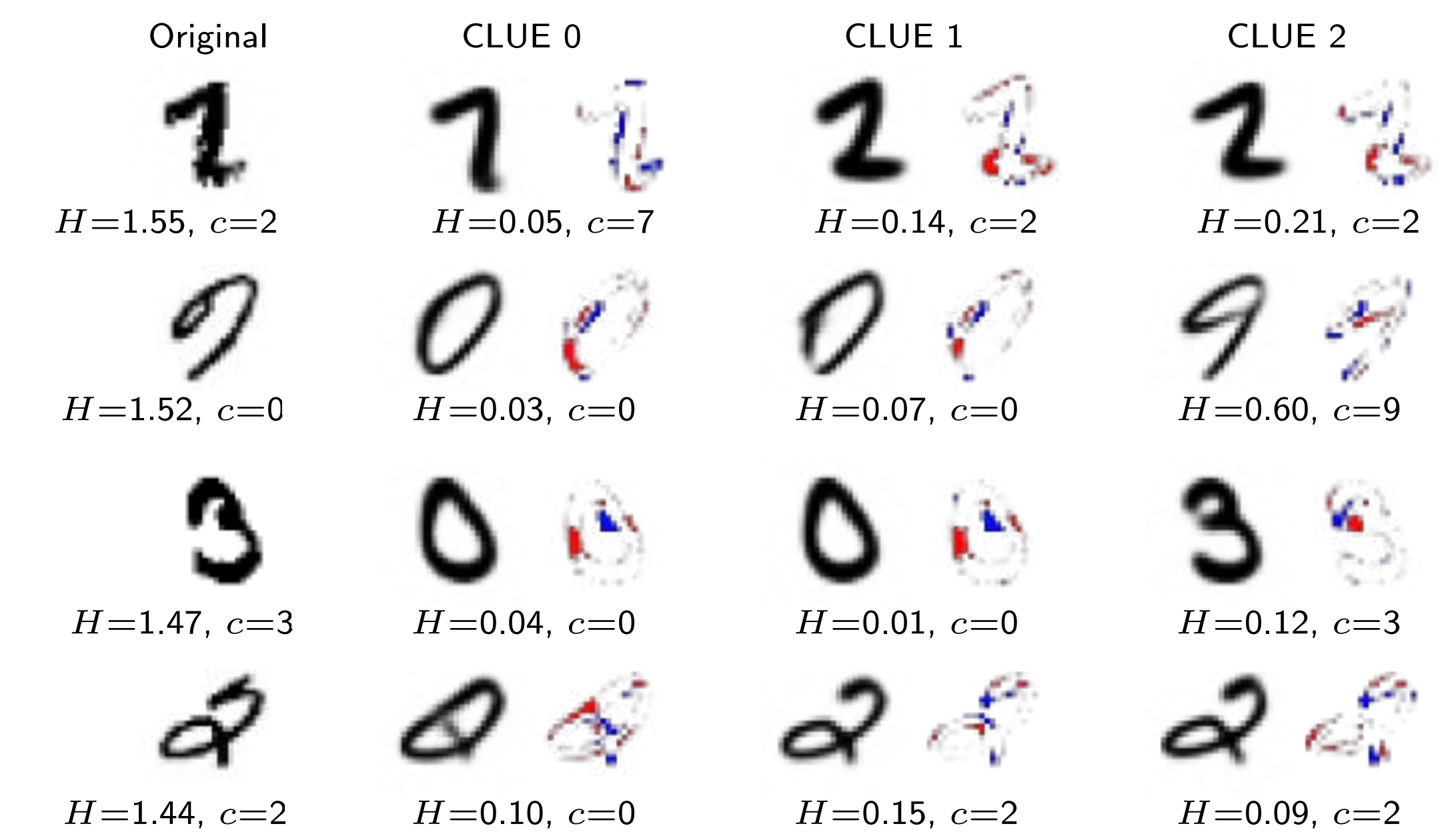


Figure: CLUEs for MNIST digits found from different random initialisations of  $z_0$ .

Original		CLUE		CLUE		CLUE	
Age	Greater than 45	Age	Greater than 45	Age	25-45	Age	25-45
Race	African-American	Race	African-American	Race	African-American	Race	African-American
Sex	Female	Sex	Female	Sex	Male	Sex	Male
Current Charge	Felony	Current Charge	Misdemeanour	Current Charge	Misdemeanour	Current Charge	Felony
Reoffended Before	No	Reoffended Before	No	Reoffended Before	No	Reoffended Before	No
Prior Convictions	1	Prior Convictions	0	Prior Convictions	0	Prior Convictions	1
Days Served	0	Days Served	0	Days Served	0	Days Served	0

Figure: CLUEs for COMPAS points found from different random initialisations of  $z_0$ .

## User Study

We run a forward simulation task. Users are told to predict if a BNN will be certain or uncertain about a new point after having been shown context points together with corresponding counterfactual explanations. We employ the COMPAS and LSAT tabular datasets. 40 ML graduate students take our survey, answering 18 questions, 9 per dataset.

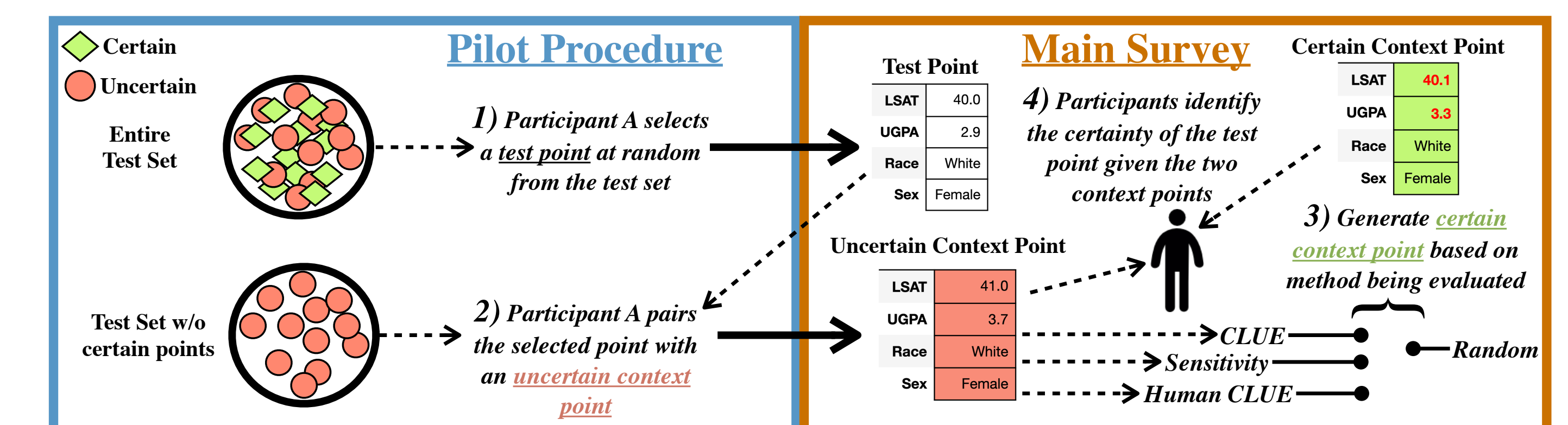


Figure: Experimental workflow for our tabular user study.

Table: Average participant accuracy by variant. Human selected refers to pilot users selecting counterfactuals points manually. We include random certain points as a baseline. We compare CLUE vs baselines using unpaired Wilcoxon signed-rank tests.

	CLUE	Human Choice	Uncertainty	Sensitivity	Random	Explanation
Accuracy (%)	82.22	62.22	52.78	61.67		
p-value	-	2.34e-5	2.60e-9	1.47e-5		