# ML4Seismic Partners Meeting 2023

# Interventionist Uncertainty in Neural Networks: A Case Study in Prompting

Mohit Prabhushankar, Prithwijit Chowdhury, Mohammad Alotaibi, and Ghassan AlRegib

OLIVES
@GeorgiaTech

GT Georgia Tech

**Prompts allow extracting contextual and relevant information from the model**

Segmentation without Prompting

Segmentation with Prompt



All objects segmented

Manual prompting selects only one segment

OLIVES
@GeorgiaTech

Georgia Tech

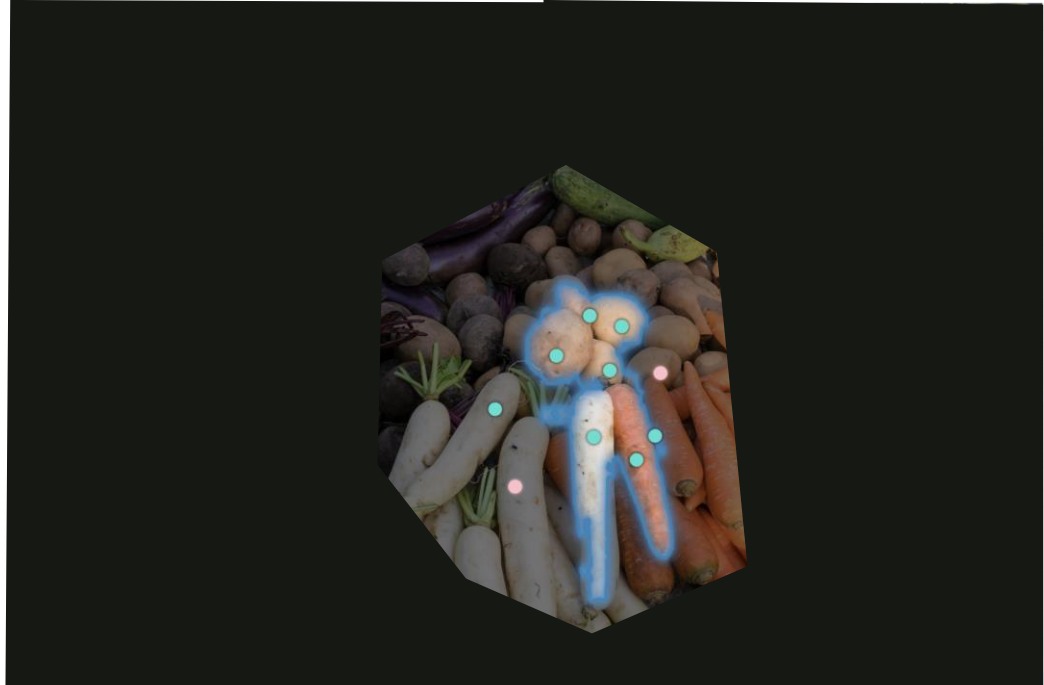**A naïve view of Prompts: Remove irrelevant (as defined by interpreters) <u>data</u> from input**

Segmentation without Prompting

Segmentation with Prompt



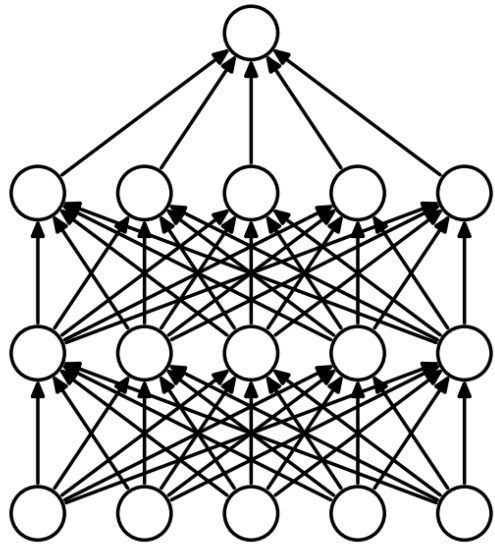All objects segmented

Manual prompting selects only one segment

$$x \rightarrow \{x, P\} \rightarrow S_x$$

OLIVES
@GeorgiaTech

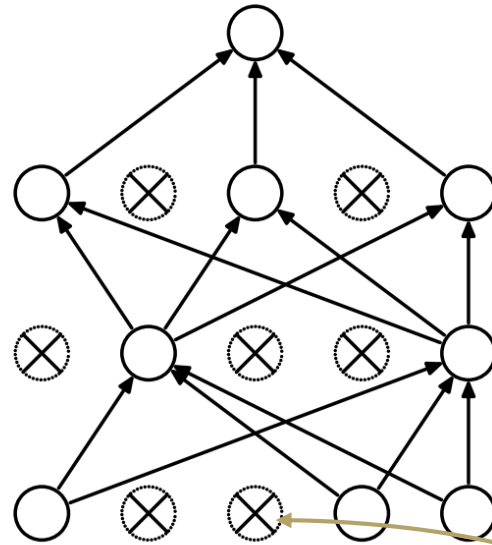Georgia Tech

**A naïve view of Prompts: Remove irrelevant (as defined by interpreters) <u>weights</u> from model**

Segmentation with Prompt



All objects segmented

Manual prompting
selects only one segment

Monte-Carlo Dropout

$$W \rightarrow \{W, P\} \rightarrow W_P$$

Model Uncertainty

[Interventionist Uncertainty] | [Mohit Prabhushankar] | [Nov. 8, 2023]

OLIVES
@GeorgiaTech

Georgia
Tech.

**Objective: To motivate and quantify Prompts as Uncertainty**

## Existing Uncertainty Framework: Flow of Information in <u>one direction</u>



Data Parameters      Seismic Data      Algorithm      Labels

Seismic Images

Sensors

Training/Inference

Data      Images      Interpretations

$$p(X|\xi)$$

$$p(\alpha|X,W)$$

$$\xi \longrightarrow X \longrightarrow \alpha = \text{"Labels"}$$

**Data Uncertainty**      **Model Uncertainty**      **Interpretational Uncertainty**

OLIVES
@GeorgiaTech

Georgia Tech

**Prompts at Inference: Flow of Information is a <u>Loop</u>**

[Interventionist Uncertainty] | [Mohit Prabhushankar] | [Nov. 8, 2023]

**Analyzing Prompting: Via objective mean Intersection Over Union (mIOU)**

Mask 1, Score: 0.321

Mask 1, Score: 0.716

Interpreter 1

Interpreter 2

Interpreter 2 > Interpreter 1

- **Goal:** To delineate region of interest
- **Quantifiable score:** Mean Intersection over Union (mIOU) between prediction and ground truth
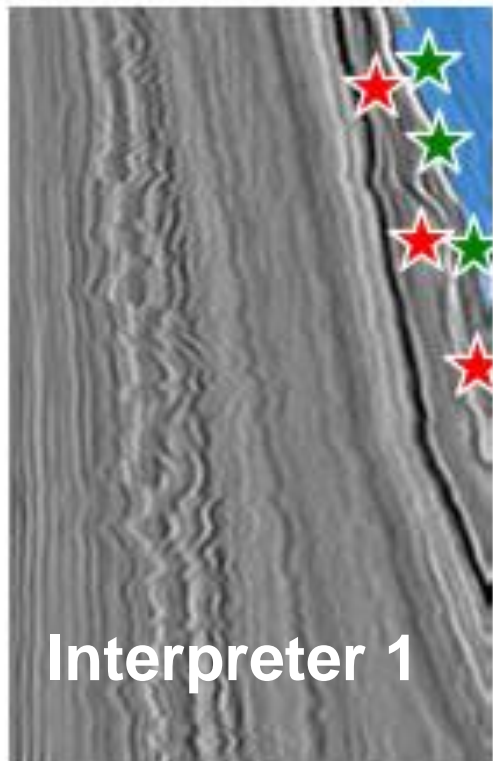
$$E(y|S_{x2}) > E(y|S_{x1})$$

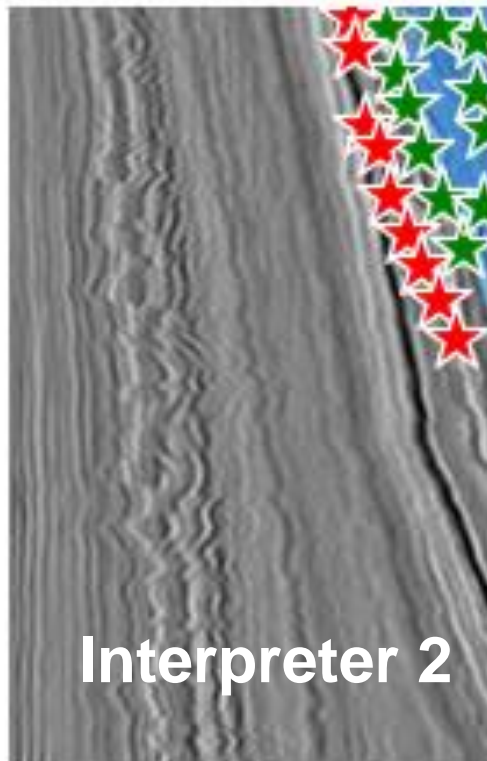[Interventionist Uncertainty] | [Mohit Prabhushankar] | [Nov. 8, 2023]

OLIVES
@GeorgiaTech

Georgia Tech

## Analyzing Interventionist Uncertainty: Via objective output metrics



Mask 1, Score: 0.853 — Interpreter 1
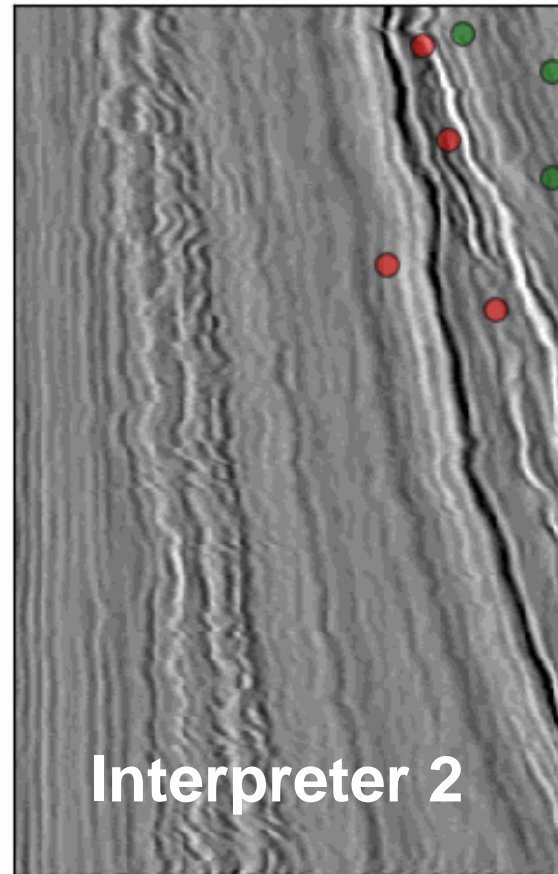
Mask 1, Score: 0.841 — Interpreter 2

Interpreter 2 ❓ Interpreter 1

- **Goal:** To delineate region of interest
- **Quantifiable score:** Mean Intersection over Union (mIOU) between prediction and ground truth
  1. **mIOU** between the two interpreters **is the same**
  2. However, the prompts (**numbers and locations**) are different
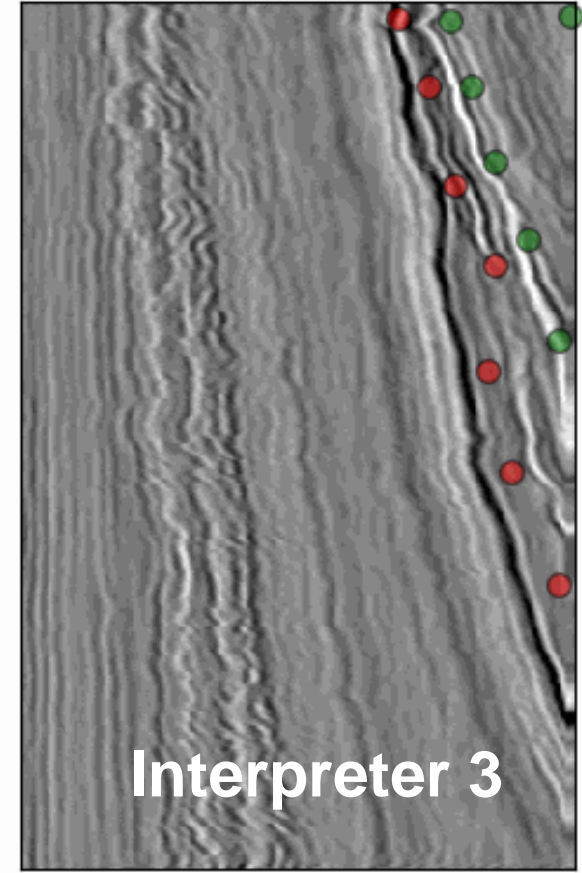
OLIVES @GeorgiaTech

Georgia Tech

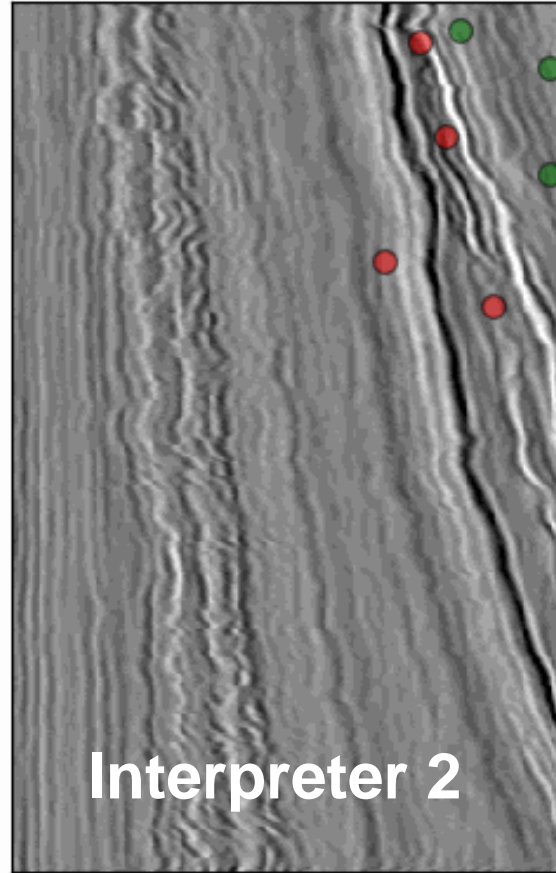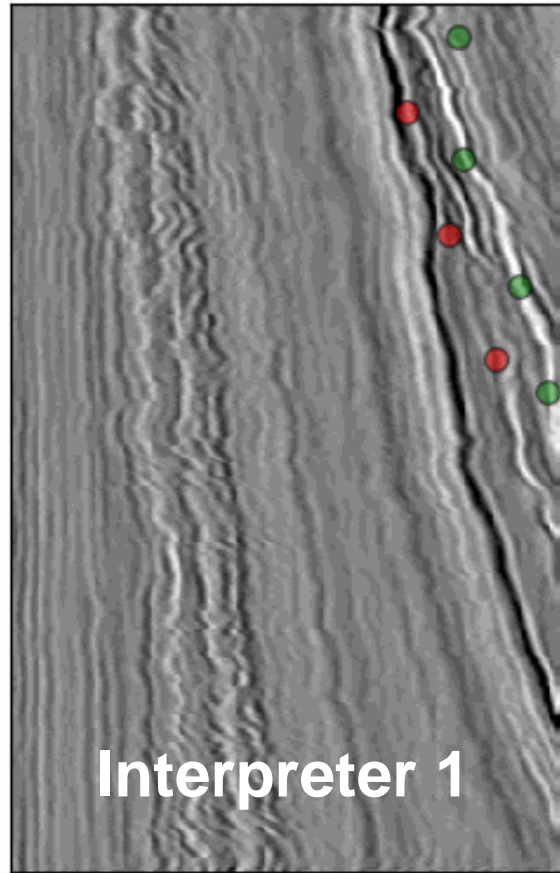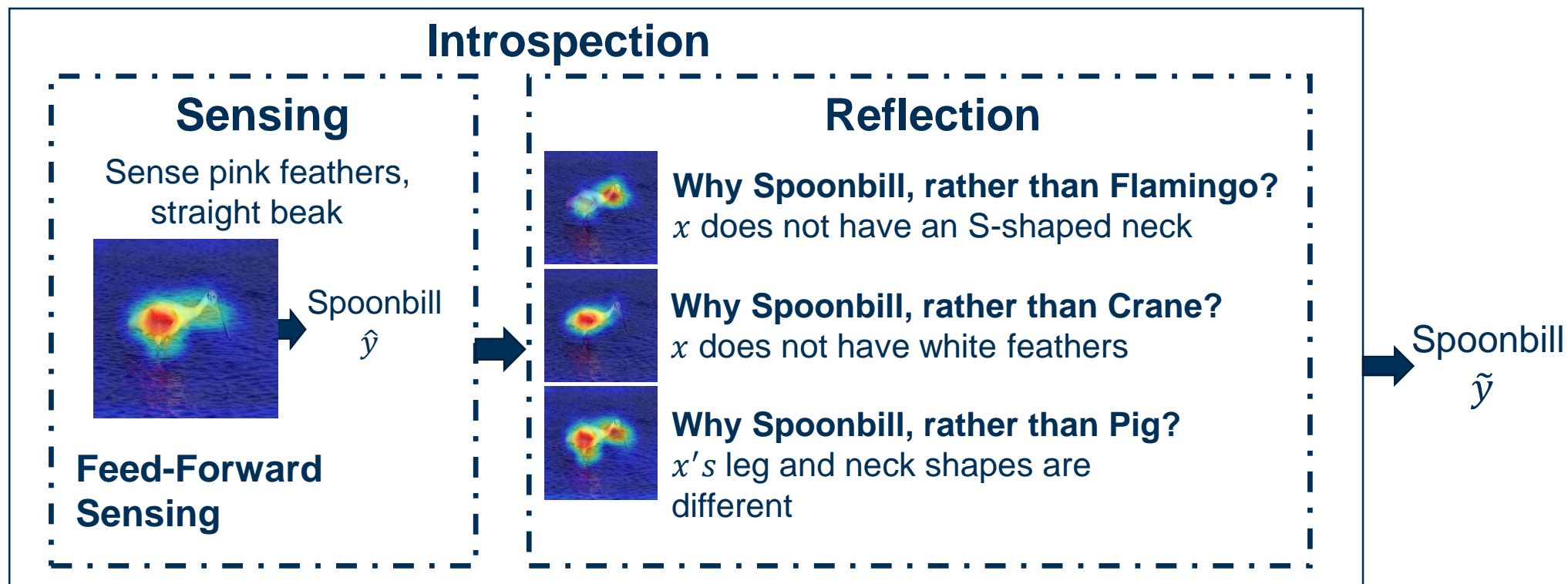**Prompting across Sections: Notice the change in prompts in consecutive sections**
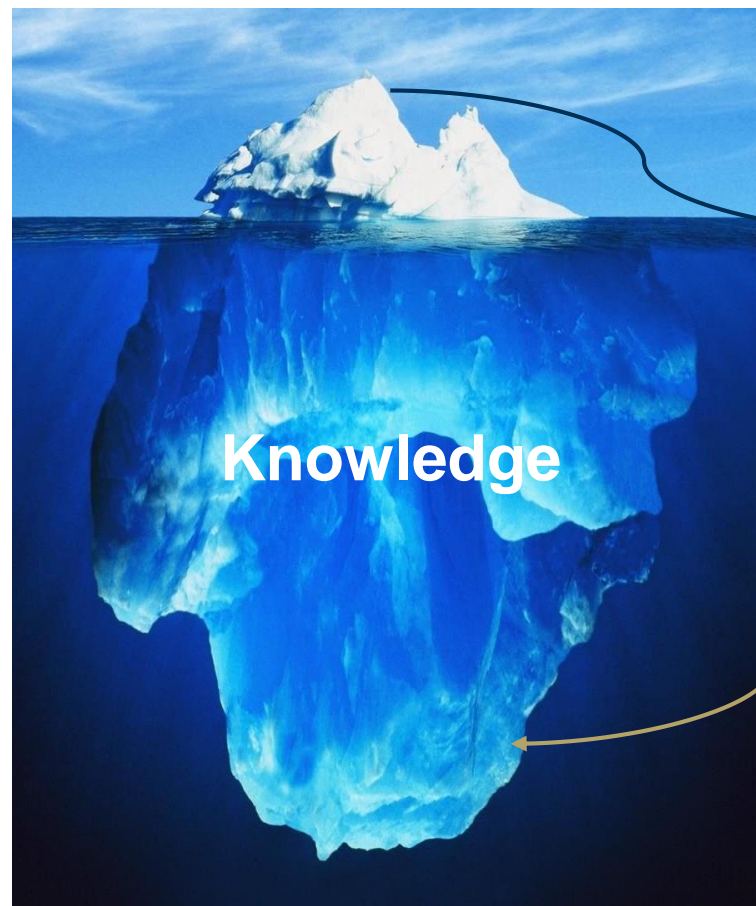


Interpreter 2

**Prompting across Interpreters: Notice the change in prompts between interpreters**



Interpreter 1

Interpreter 2

Interpreter 3

OLIVES
@GeorgiaTech

Georgia Tech

## Introspective Learning using Contrastive Questions: Uncertainty resides in contrastive questions



**Introspection**

**Sensing**

Sense pink feathers, straight beak

Spoonbill $\hat{y}$

**Feed-Forward Sensing**

**Reflection**

**Why Spoonbill, rather than Flamingo?**
$x$ does not have an S-shaped neck

**Why Spoonbill, rather than Crane?**
$x$ does not have white feathers

**Why Spoonbill, rather than Pig?**
$x's$ leg and neck shapes are different

Spoonbill $\tilde{y}$

OLIVES @GeorgiaTech

Georgia Tech

**Trained Neural Nets have hidden knowledge. Goal is to prompt it at Inference.**



Optimization

Inference

Why P, rather than Q?

Data Collection



Knowledge

Traditional *Why P?*

*What if? Why P, rather than Q?*

[Interventionist Uncertainty] | [Mohit Prabhushankar] | [Nov. 8, 2023]

OLIVES
@GeorgiaTech

Georgia Tech

**Our Analysis: Interventionist Uncertainty via Contrastive Questions**



Mask 1, Score: 0.853    Mask 1, Score: 0.841

Why Prompt $P_1$, rather than $P_2$?

OLIVES
@GeorgiaTech

Georgia Tech

**Variance Decomposition of Uncertainty under intervention**

$$V[y|S_x] = V[E(y|S_x)] + E(V[(y|S_x)])$$

$y$ = Prediction
$V[y]$ = Variance of prediction (Predictive Uncertainty)
$S_x$ = Intervened Data
$E(Y|S_x)$ = Expectation of class under intervention
$V(Y|S_x)$ = Variance of class under all residuals

[Interventionist Uncertainty] | [Mohit Prabhushankar] | [Nov. 8, 2023]

## Variance Decomposition of Uncertainty under intervention

$$V[y|S_x] = V[E(y|S_x)] + E(V[(y|S_x)])$$



Mask 1, Score: 0.853     Mask 1, Score: 0.841

zero

$y$ = Prediction
$V[y]$ = Variance of prediction (Predictive Uncertainty)
$S_x$ = Intervened Data
$E(Y|S_x)$ = Expectation of class under intervention
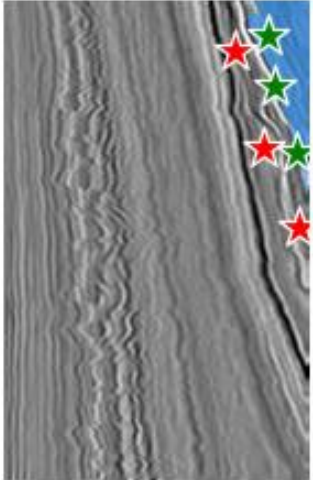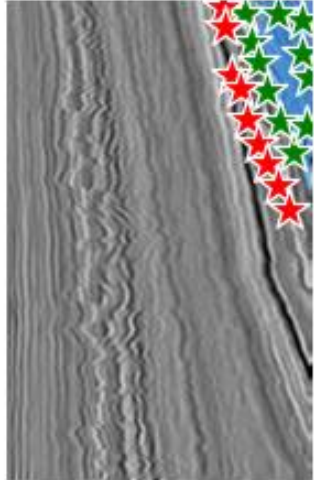$V(Y|S_x)$ = Variance of class under all residuals

$$E(y|S_{x2}) = E(y|S_{x1})$$

[Interventionist Uncertainty] | [Mohit Prabhushankar] | [Nov. 8, 2023]

OLIVES
@GeorgiaTech

Georgia Tech

**Variance Decomposition of Uncertainty under intervention**

$$V[y|S_x] = V[E(y|S_x)] + E(V[(y|S_x)])$$

Given an intervention $S_x$,
find alternative
interventions $S_x{'}$ that result
in non-zero $V[(y|S_x)]$

$y$ = Prediction
$V[y]$ = Variance of prediction (Predictive Uncertainty)
$S_x$ = Intervened Data
$E(Y|S_x)$ = Expectation of class under intervention
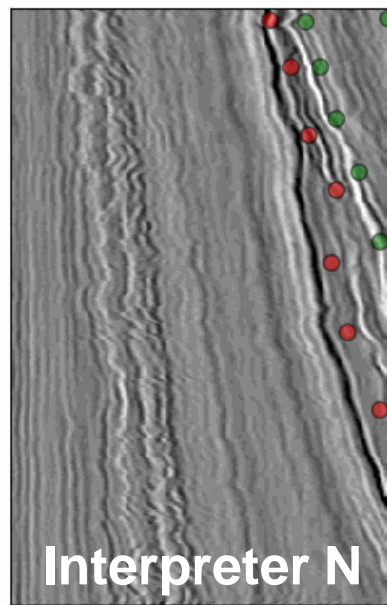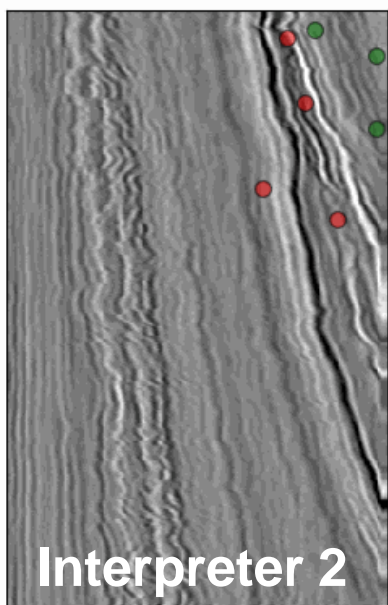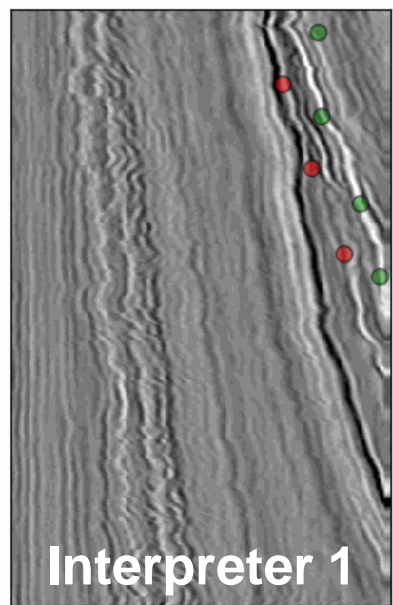$V(Y|S_x)$ = Variance of class under all residuals

alternatives

[Interventionist Uncertainty] | [Mohit Prabhushankar] | [Nov. 8, 2023]

OLIVES
@GeorgiaTech

Georgia Tech

**Take variance across outputs derived from N prompts**

$$V[y|S_x] = V[E(y|S_x)] + E(V[(y|S_x)])$$



Interpreter 1    Interpreter 2    Interpreter N

$Y_1|S_{x1}$
$Y_2|S_{x2}$
$Y_3|S_{x3}$
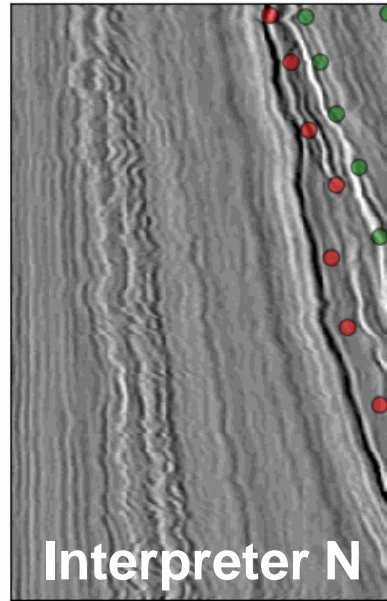$Y_4|S_{x4}$
$Y_5|S_{x5}$
.
.
$Y_N|S_{xN}$

Variance

Variance

[Interventionist Uncertainty] | [Mohit Prabhushankar] | [Nov. 8, 2023]

**Uncertainty resides near the prompt boundaries**

$$V[y|S_x] = V[E(y|S_x)] + E(V[(y|S_x)])$$



Interpreter 1   Interpreter 2   Interpreter N

Variance

**Goal: Among given 7 prompts, choose the best prompts for each section**

Best Prompts = **<u>Highest mIOU</u>** against ground truth

Not available at prompting

Best Guess Prompts = Highest mIOU against **<u>uncertainty</u>**

Accuracy(Best prompts, Best Guess Prompts) = 34.66%
(random = 14%)

OLIVES
@GeorgiaTech

Georgia Tech

**Signal-to-Noise Ratio (SNR) of interventional uncertainty follows IoU**

Mean of uncertainty map
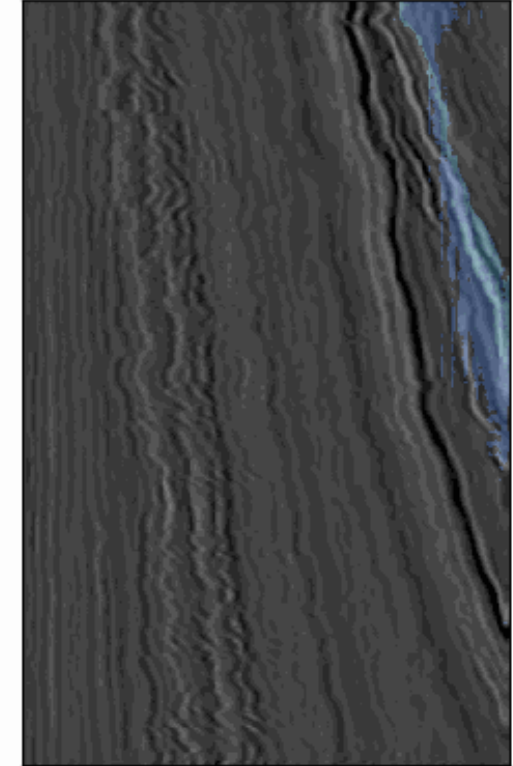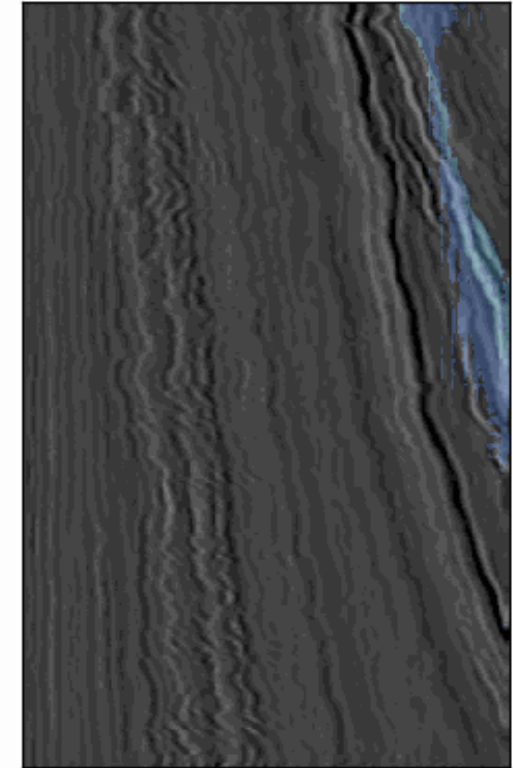
$$SNR = \frac{\mu(V(y|S_x))}{\sigma(V(y|S_x))}$$

Standard deviation of uncertainty map

Cosine Similarity(SNR, IoU(Best prompt, GT)) = 0.83

Even without knowing ground truth (GT), we can estimate how well the best prompt will perform

[Interventionist Uncertainty] | [Mohit Prabhushankar] | [Nov. 8, 2023]

OLIVES
@GeorgiaTech

Georgia Tech.

- **Conclusion 1: Higher the IoU, higher the uncertainty in explanation (or less trustworthy is the explanation)**

- **Conclusion 2: Higher the SNR of uncertainty, more is the dispersal (or less trustworthy is the explanation)**

$$x \qquad S_{x_1} \qquad S_{x_2}$$



$S_{x1}, S_{x2}$ are explanations from GradCAM and RISE

[Interventionist Uncertainty] | [Mohit Prabhushankar] | [Nov. 8, 2023]

OLIVES
@GeorgiaTech

Georgia Tech

VOICE: Variance of Induced Contrastive Explanations to quantify Uncertainty in Neural Network Interpretability

**Framework utilizes gradients at Inference**

OLIVES
@GeorgiaTech

Georgia Tech

# Conclusion
## Key Takeaways

- Interventional Uncertainty provides a framework for analyzing prompts

- Contrastive analysis provides best guess prompts **among available prompts**
  - Answers `*Why Prompt P, rather than Prompt Q?'*

- Contrastive analysis **does not** provide objectively **best prompts**
  - Does not say the exact location of where to prompt

- Future work to include model parameters within the analysis framework

**For more OLIVES content, please visit:**

GitHub

Publications

SCAN ME

SCAN ME

[Interventionist Uncertainty] | [Mohit Prabhushankar] | [Nov. 8, 2023]

OLIVES
@GeorgiaTech

Georgia Tech

# Tutorials and Short Courses
## Completed and Upcoming Tutorials

## Latest tutorial delivered at ICIP 2023



https://alregib.ece.gatech.edu/ieee-icip-2023-tutorial/

## Upcoming tutorials/short courses

- **Dec 5 – 7 (Virtual):** 10 hr Short Course on Explainability – Invited by IEEE Signal Processing Society

- **Dec 15-18 (IEEE Big Data 2023):** 3 hr Tutorial on Robustness of Neural Networks

- **Jan 7-8 (WACV 2024):** 3 hr Tutorial on Explainability, Uncertainty, and Intervenability

- **Jan 22 (EI 2024):** 3 hr Tutorial on Explainability, Uncertainty, and Intervenability

- **Feb 21-22 (AAAI 2024):** 3.5 hr Tutorial on Explainability, Uncertainty, and Intervenability

[Interventionist Uncertainty] | [Mohit Prabhushankar] | [Nov. 8, 2023]

OLIVES
@GeorgiaTech

Georgia Tech