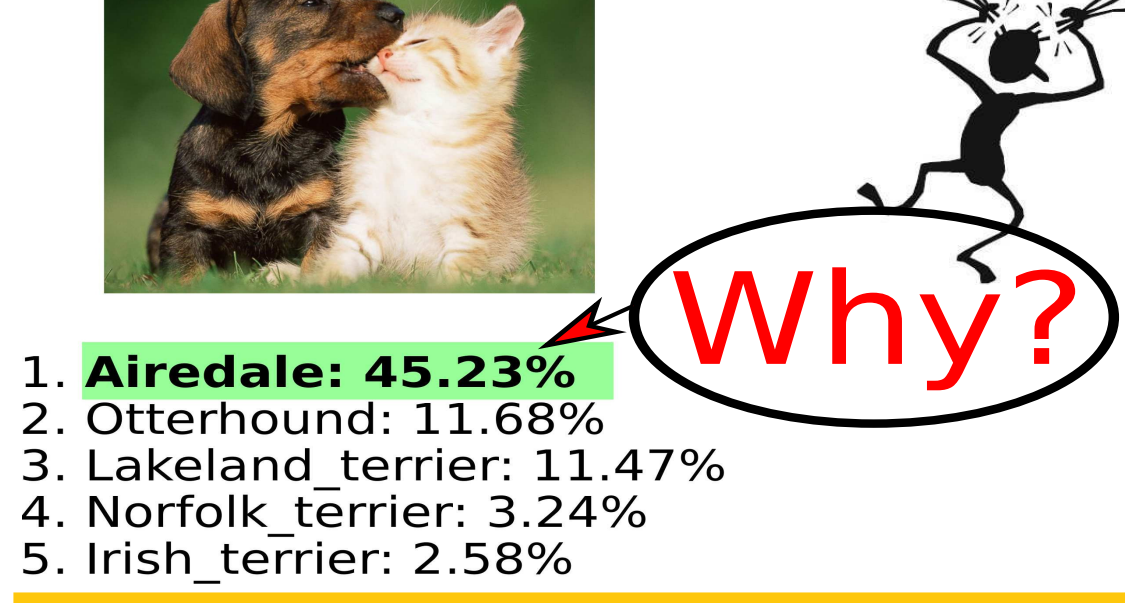


A generalizable saliency-map based interpretation of model outcome

Shailja Thakur and Sebastian Fischmeister

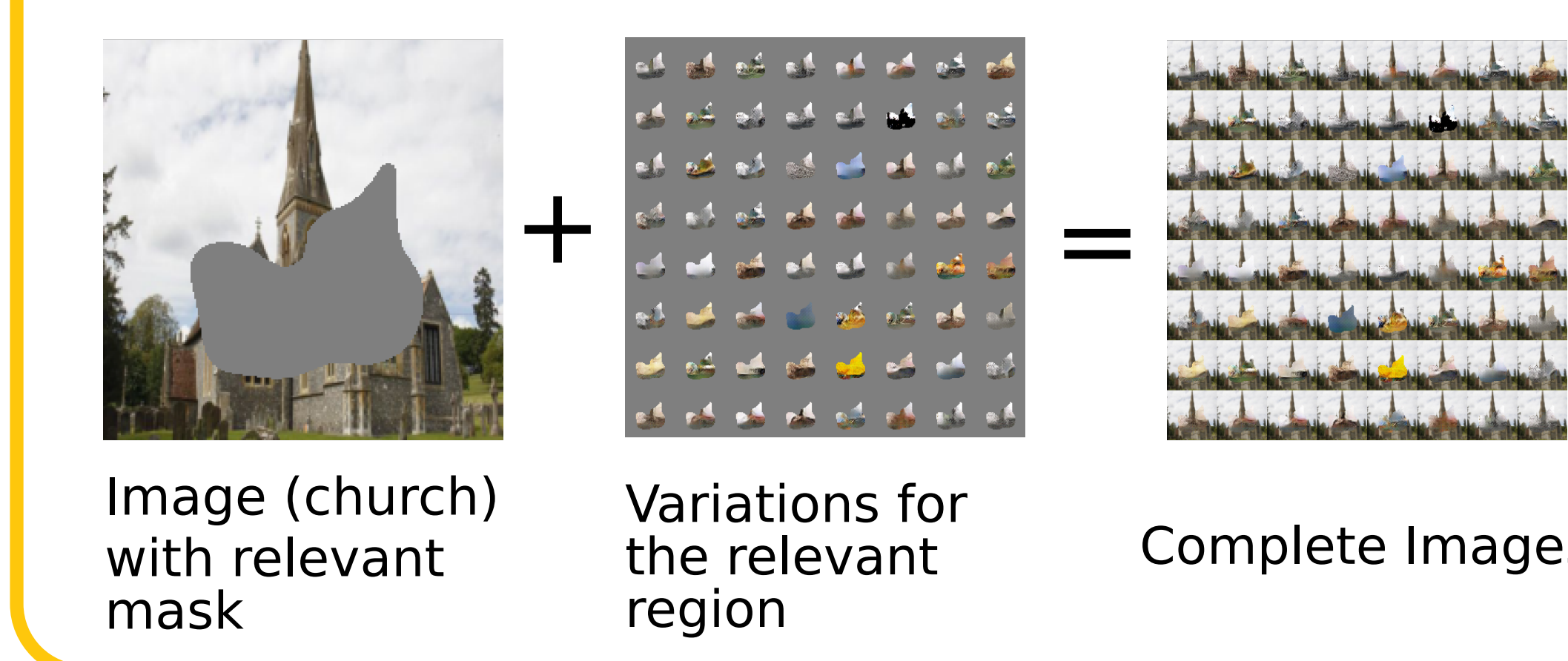
{s7thakur, sfischme}@uwaterloo.ca



Problem

Problem1: Given a classification model, an input vector, and a confidence score for the target, weigh the input samples in the order of their importance for classification of the input to the target class.

Problem2: Given a relevance mask and a generative model, identify the distribution of acceptable variations for the important input samples.

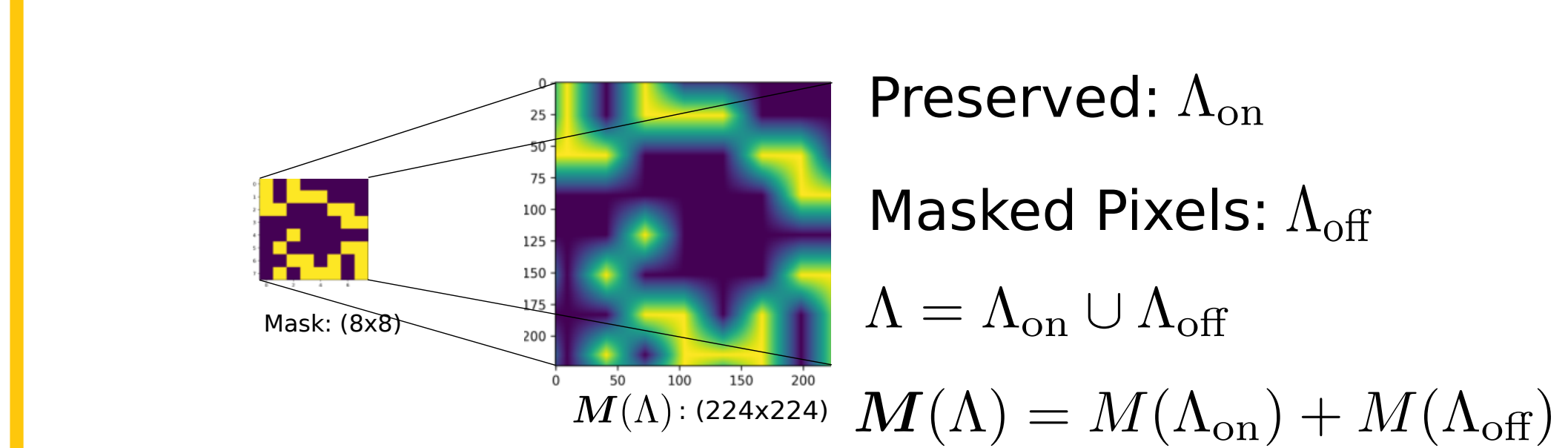


Motivation

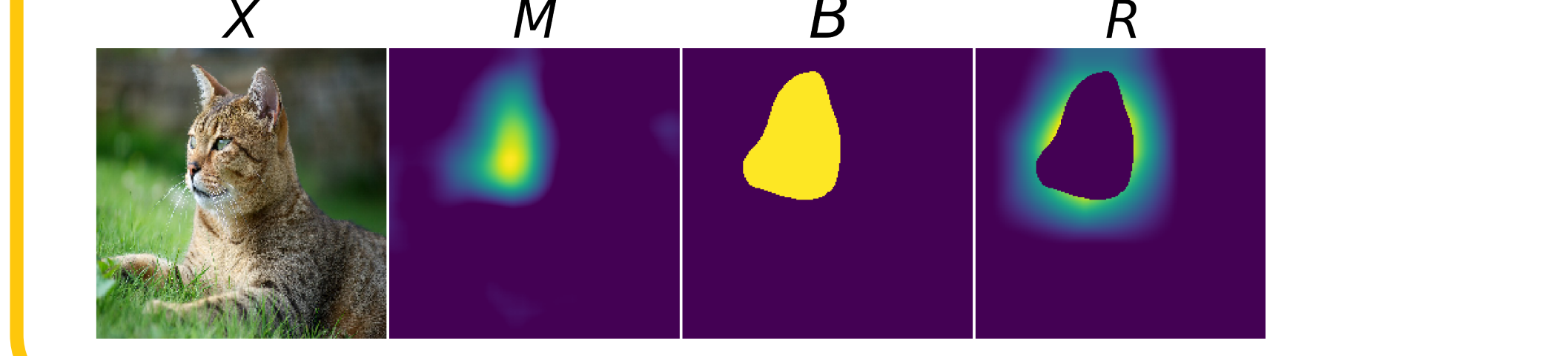
- Interpretability \downarrow as model complexity \uparrow
- Model outcome critical in decision-making
 - Medical applications
 - Self-driving cars
 - Safety-Critical autonomous systems
- Risk to human lives, property, and the environment

Notations & Initializations

Input $X : \Lambda \in \{1, \dots, H\} \times \{1, \dots, W\} \rightarrow \mathbb{R}^3$,
Output A confidence score in the target class $c \in \mathbb{R}^C$, **Classifier** $f : X \rightarrow \mathbb{R}^C$



M: Estimates saliency-map, **B:** Binary bounding-box for **M** using a threshold, **R:** Inverted and convolved (kernel $(s \times s)$) version of **B**



Saliency-Map Algorithm

Key: The classifier's output is sensitive to the changes in the input.

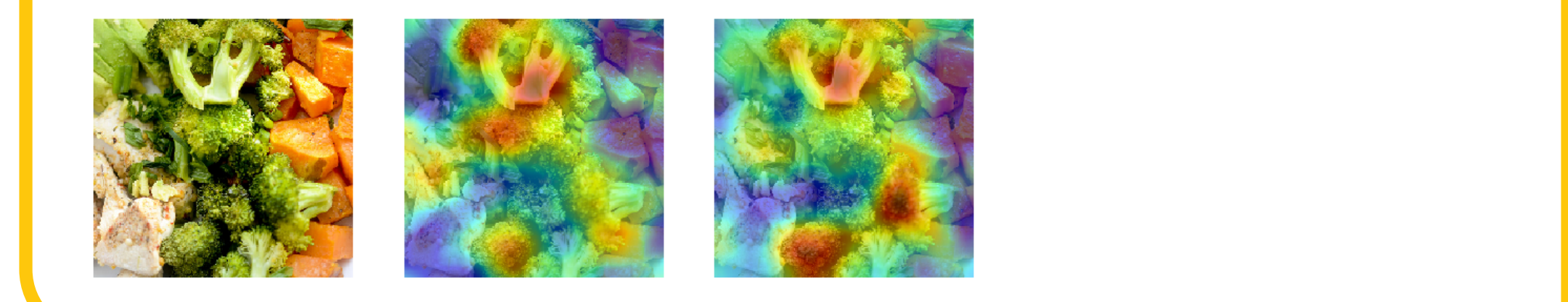
Mask Estimation Algorithm 1

Input: Input image $I \in \mathcal{I}$, Target class $c \in \mathcal{C}$
Output: saliency map M

Initialisation: $M_0 \in \mathbb{R}^{h \times w}$

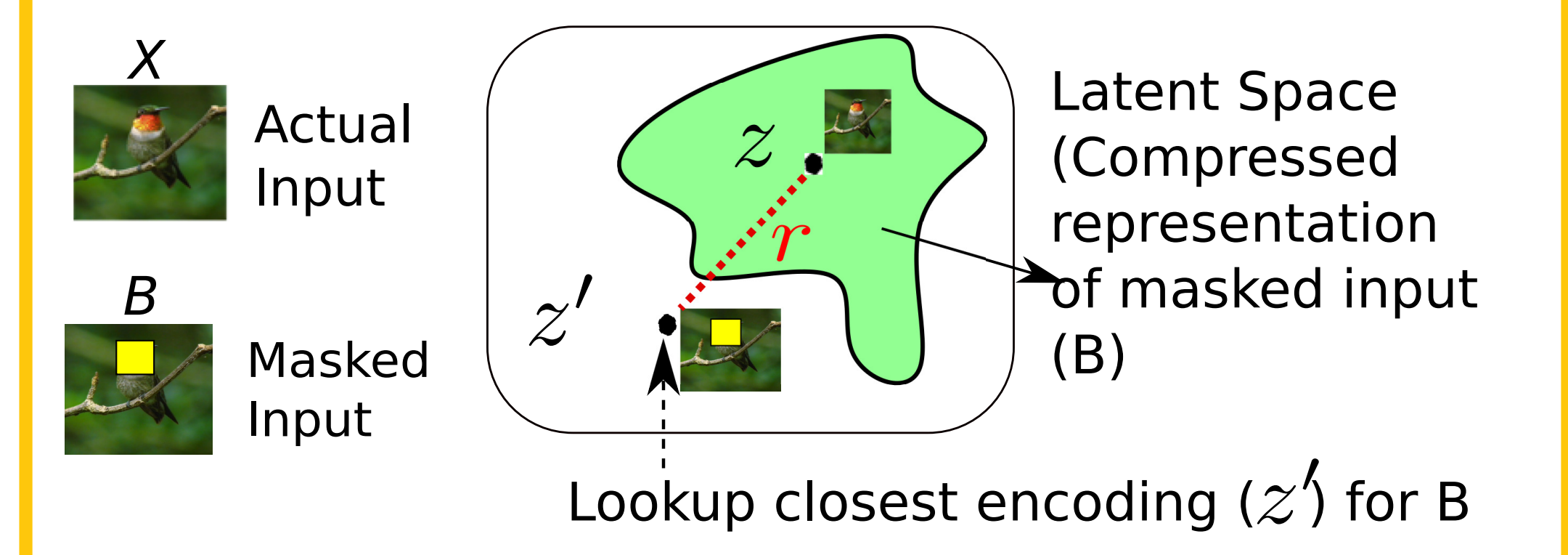
- for $i = 1$ to N do
- $M \leftarrow$ upsample M_{i-1}
- $p_i(c|(I \odot M)) = f(I \odot M)$ **Compute model's confidence (p_i) in target class (c) using masked input ($I \odot M$)**
- $\Delta p \leftarrow |p_i - p_{i-1}|$
- $\Lambda_1 \leftarrow$ randomly select $n_1 p_i$ pixels of Λ_{on}
- $\Lambda_2 \leftarrow$ randomly select $n_2(1 - p_i)$ pixels of Λ_{off} **Revise the set of preserved (Λ_{on}) and masked pixel (Λ_{off})**
- $\Lambda_{on} \leftarrow \Lambda_1 \cup \Lambda_2$
- $\Lambda_{off} \leftarrow \Lambda \setminus \Lambda_{on}$
- if $(\lambda \in \Lambda_{on})$ then
- $M_i(\lambda) \leftarrow M_{i-1}(\lambda) \Delta p$
- end if
- if $(\lambda \in \Lambda_{off})$ then
- $M_i(\lambda) \leftarrow 0$
- end if
- $V = \sum_{x,y} \|M_{i-1}^{xy+1} - M_{i-1}^{xy}\|^2 + \sum_{x,y} \|M_{i-1}^{x+1y} - M_{i-1}^{xy}\|^2$
- $M_i \leftarrow M_i + \eta V M_{i-1}$ **Update Mask**
- end for
- return M **Estimated Saliency-Map**

Explanation using saliency based approach lacks consistency in the detected relevant regions across runs.



Alternate Variations

Key: Latent space (z) is invariant to small perturbations (rotation, inversion, scaling, and shear).
Generative Adversarial Network



Use backpropagation, learn the encoding (z') by minimizing the loss $L(z') = L_{Reconstruction} + L_{Contextual} = r$.

$L_{Reconstruction}$ is the MSE(x, x') and $L_{Contextual}$ determines whether x' is realistic.

Reconstructed variant of the input is,

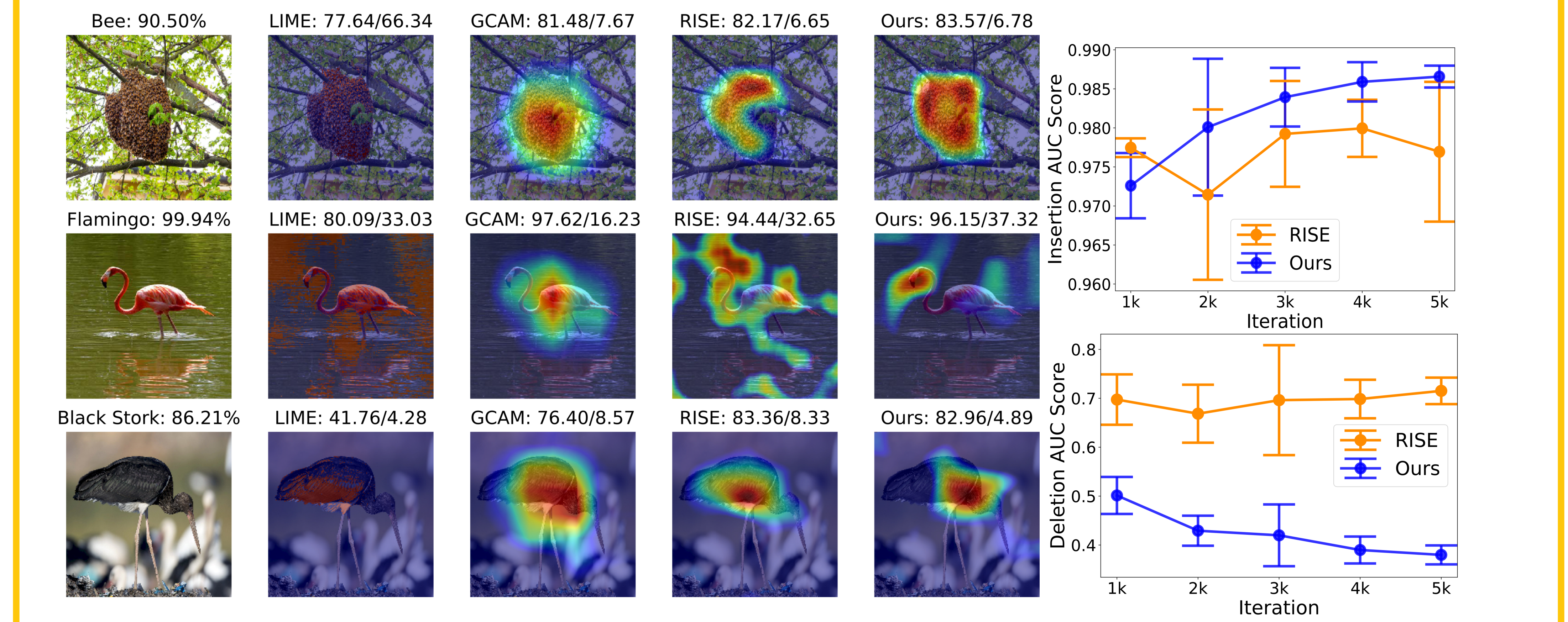
$$X' = (X \odot (1 - B)) + (B \odot G(z')) \quad (1)$$

Results

IOU (Intersection Over Union) Score: Amount of overlap between saliency region and the annotated box

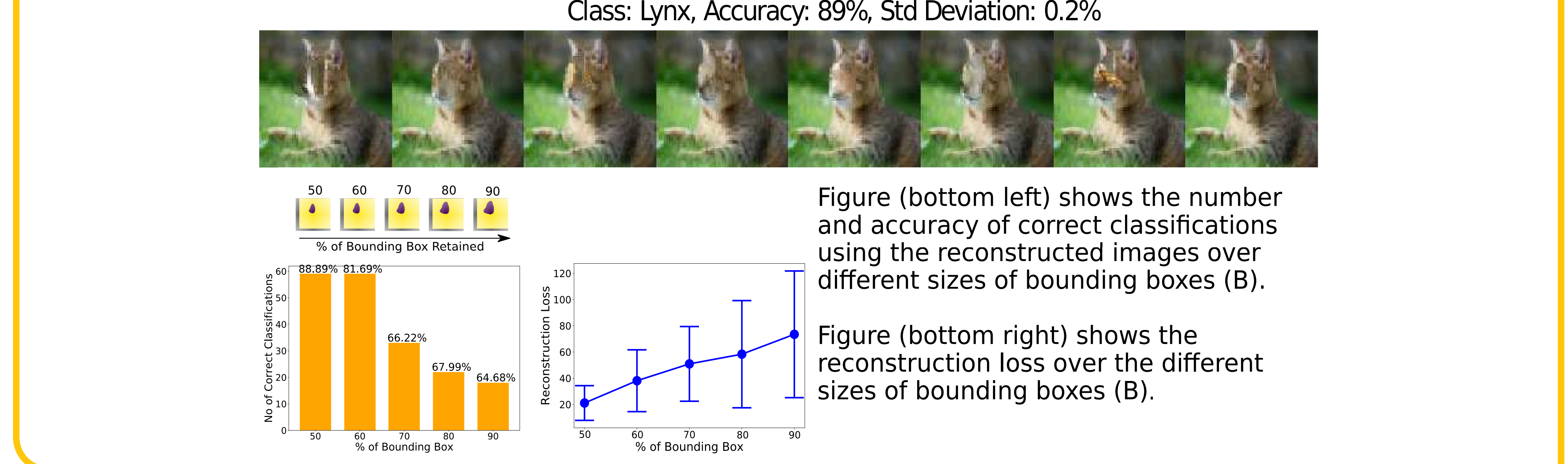


Insertion/Deletion Metric: This metric captures the sensitivity of the model to the insertion of the pixels from the relevant region of the input using an average AUC (Area Under the Curve) score.



First figure, Saliency map generated for the target-specific image classification using our approach, RISE [1], GCAM [2], and LIME [3] and the AUC scores (%) of insertion/deletion metrics [1]. Second figure shows the convergence of the AUC score of insertion/deletion for the saliency map of an input image using our approach and RISE [1] over the iterations.

Acceptable variations (X')



Open Questions

- Can the set of acceptable variations be expanded to identify adversarial examples?
- Can the approach be applied for input of type time-series?

References

- V. Petsiuk, A. Das, K. Saenko. RISE: Randomized Input Sampling for Explanation of Black-box Models 2018
- R. Selvaraju, M. Cogswell, A. Das, R. Vedantam, D. Parikh, D. Batra Grad-CAM: Visual Explanations From Deep Networks via Gradient-Based Localization 2017
- M. Ribeiro, S. Singh, C. Guestrin "Why Should I Trust You?": Explaining the Predictions of Any Classifier
- D. Pathak, P. Krähenbühl, J. Donahue, T. Darrell, A. Efros Context Encoders: Feature Learning by Inpainting