

Introduction

Comparing the underlying decision making process of Transformers and CNNs(ResNet50 and VGG19) by causally perturbing the image regions and observing the change in the output confidence.

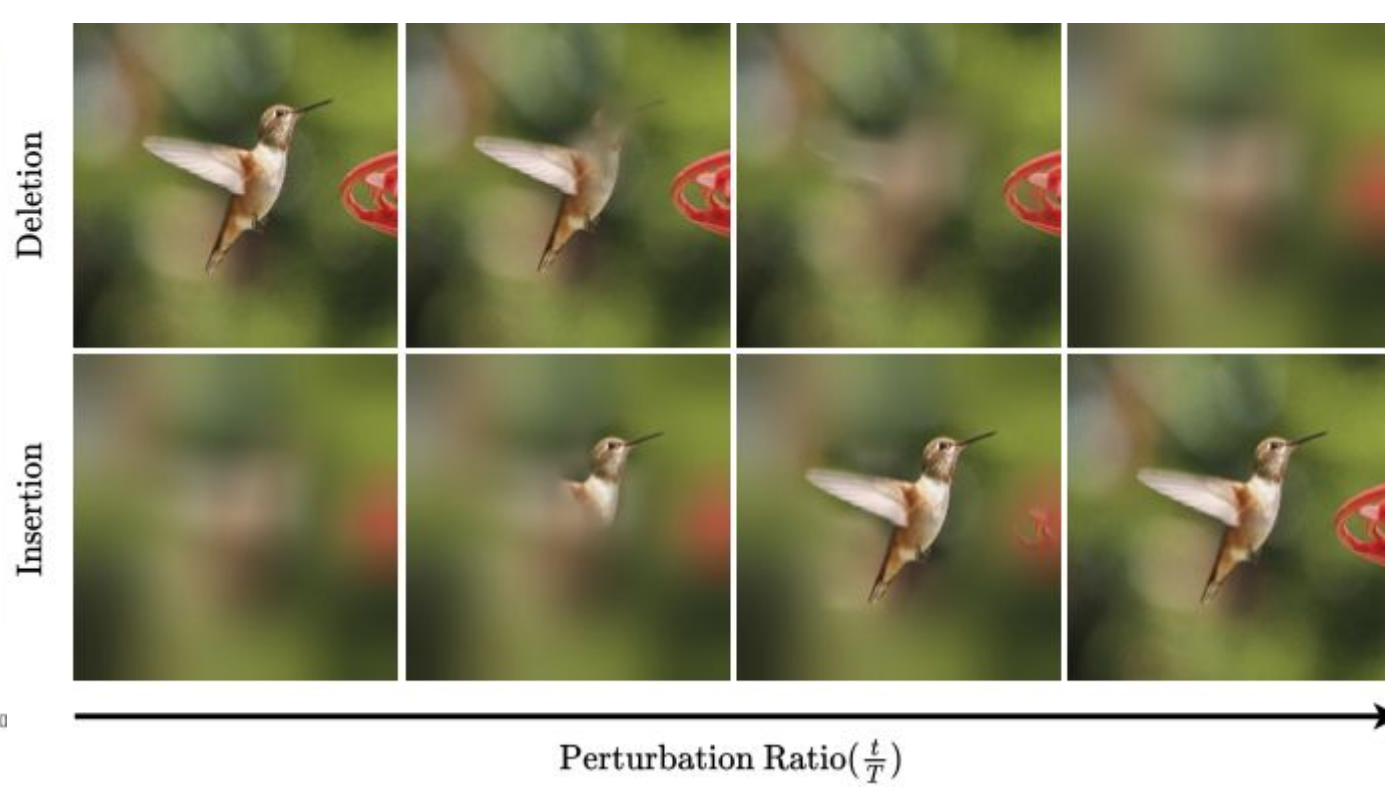
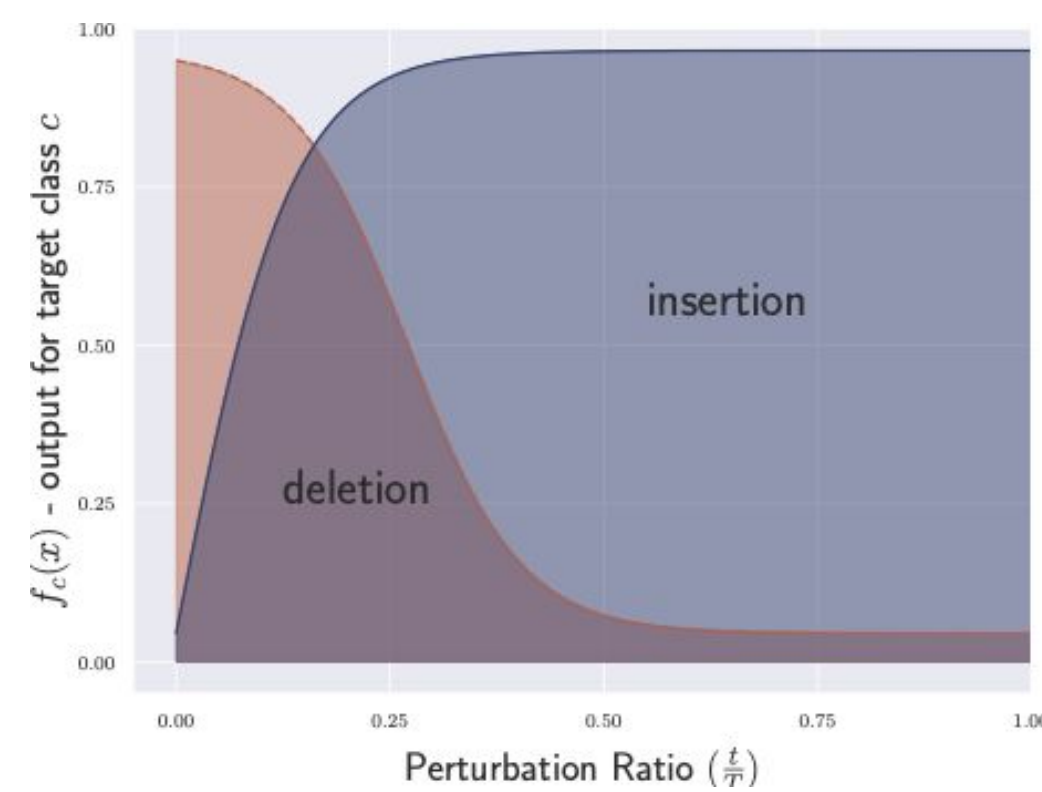
Method

Cross-Testing

- Utilize one specific deep model to generate an attribution map by iGOS++, and use another deep model to assess the insertion/deletion metrics
- Insertion metric(Deletion is input image I and the baseline image \tilde{I} are swapped):

$$insertion = \frac{1}{T} \left(\sum_{t=0}^{T-1} \frac{1}{2} (f_c(\phi^{(t)}(I, \tilde{I}, M)) + f_c(\phi^{(t+1)}(\tilde{I}, I, M))) \right)_{p_{data}}$$

where



- Fair comparison: normalization $\overline{score} = \frac{(score - b)}{(t - b)}$ where t/b are the top-1/fully-blurred confidences

Minimal Sufficient Explanations(MSEs)

- Using beam search method of SAG to find diverse set of MSEs:a minimal conjunction /region that achieves a certain high classification confidence P_h

$$f_c(N_i) \geq P_h f_c(I), \max_{n_j \in N_i} f_c(n_j) < P_h f_c(I)$$

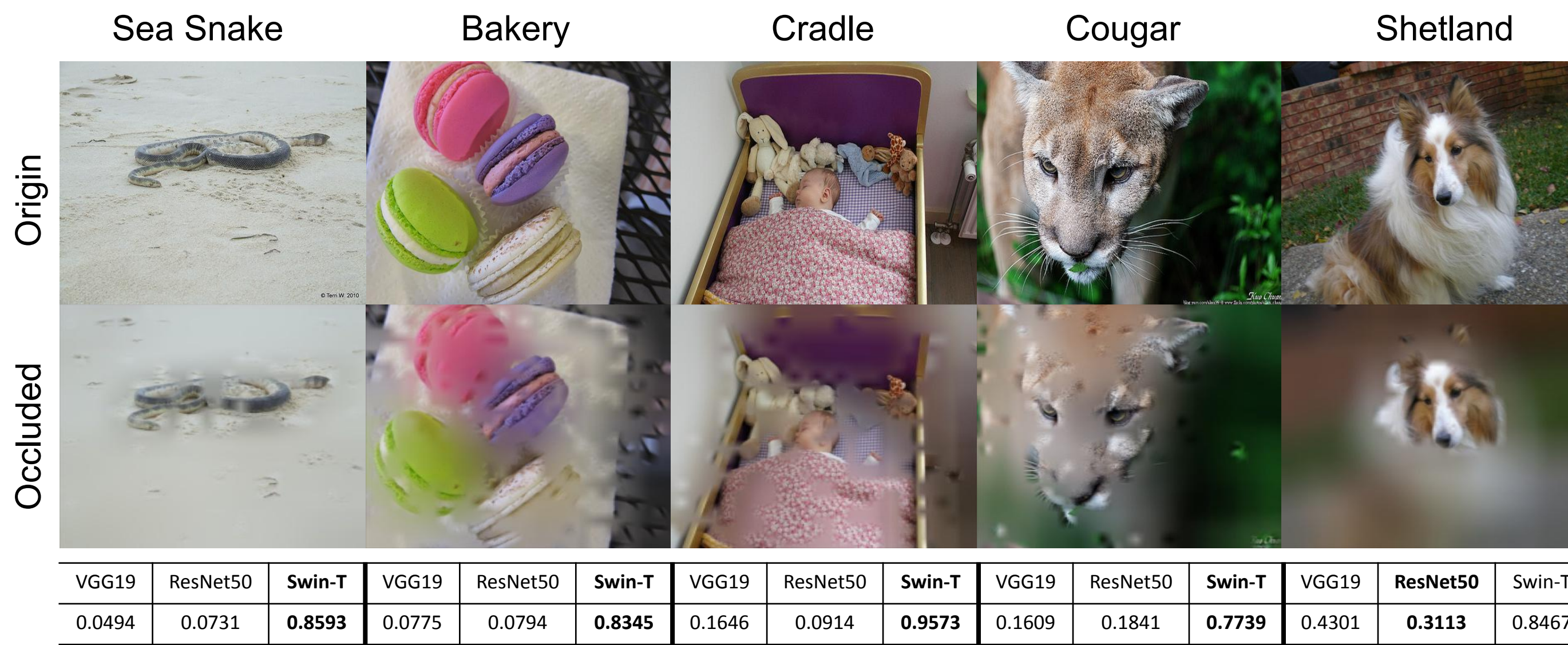
Sub-Explanation Counting

- Construct a tree for each MSE by deleting one patch at a time from a parent node to generate child nodes.
- Stop expansion when the nodes are with a confidence less than 50% compared to the classification confidence on the original image.

Robustness in terms of Insertion and Deletion

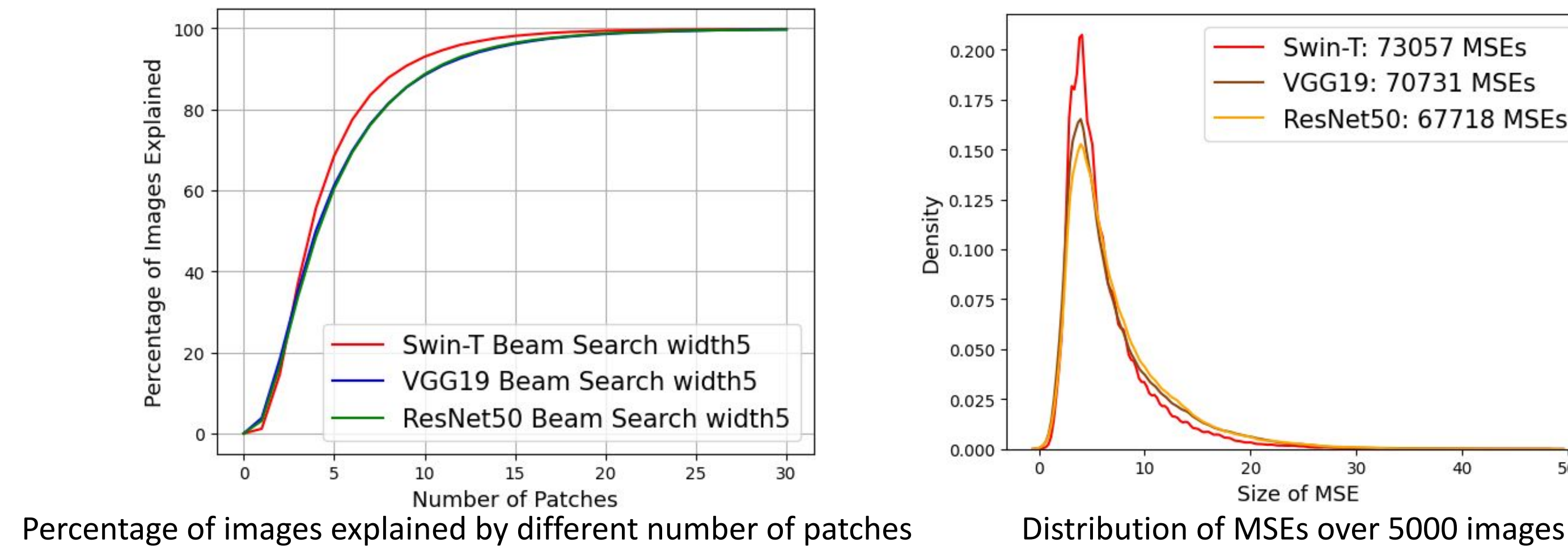
Model	Deletion ↓	Insertion ↑	Model Pair	Deletion ↓	Insertion ↑
Swin-T	0.1291	0.9433	Swin-T → ResNet50	0.1771 (+0.0480)	0.7670 (-0.1763)
ResNet50	0.1109	0.9437	ResNet50 → Swin-T	0.2322 (+0.1068)	0.8179 (-0.1076)
VGG19	0.1254	0.9255	Swin-T → VGG19	0.1731 (+0.0440)	0.7151 (-0.2282)
			VGG19 → Swin-T	0.2449 (+0.1340)	0.8236 (-0.1019)

Quantitative results of deletion/insertion metrics using themselves and their counterpart heat maps



Qualitative Cross-Testing Results

Minimal Sufficient Explanations



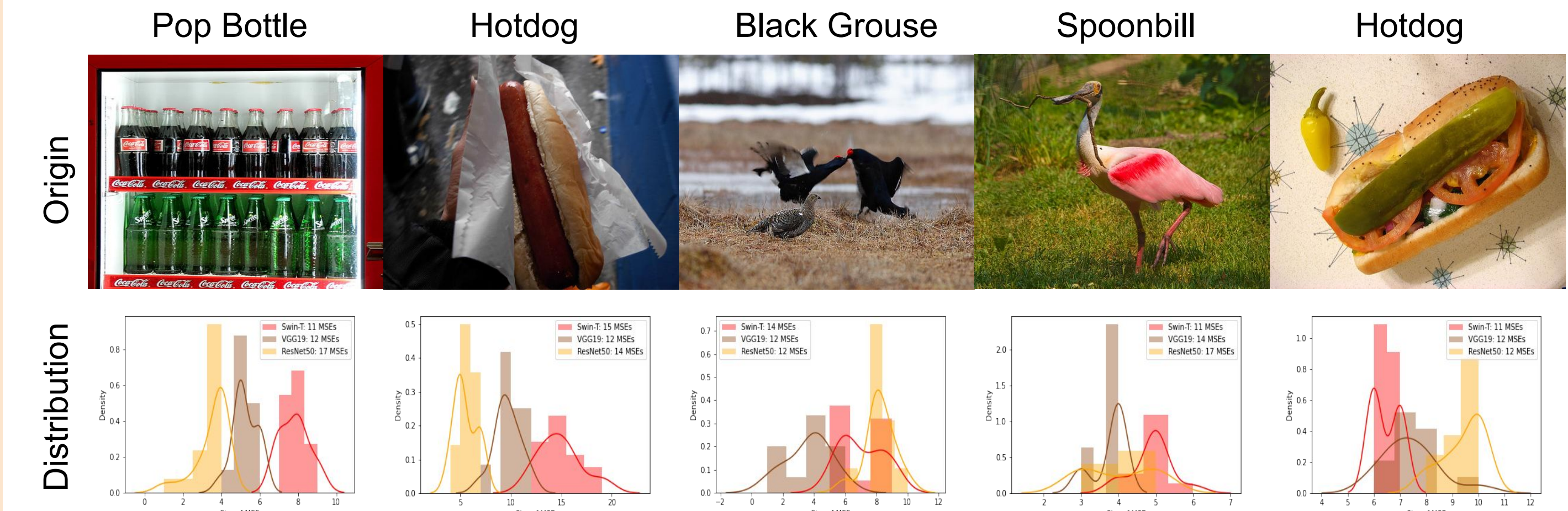
Model	MSEs			Overlap=0			Overlap=1		
	Mean	Std	Median	Mean	Std	Median	Mean	Std	Median
Swin-T	14.66	5.12	13.00	1.26	0.67	1.00	2.12	2.04	1.00
ResNet50	13.52	4.42	12.00	1.25	0.66	1.00	1.98	1.96	1.00
VGG19	14.10	4.97	12.00	1.27	0.69	1.00	2.08	2.12	1.00

Number of diverse MSEs obtained and by allowing for different degrees of overlap

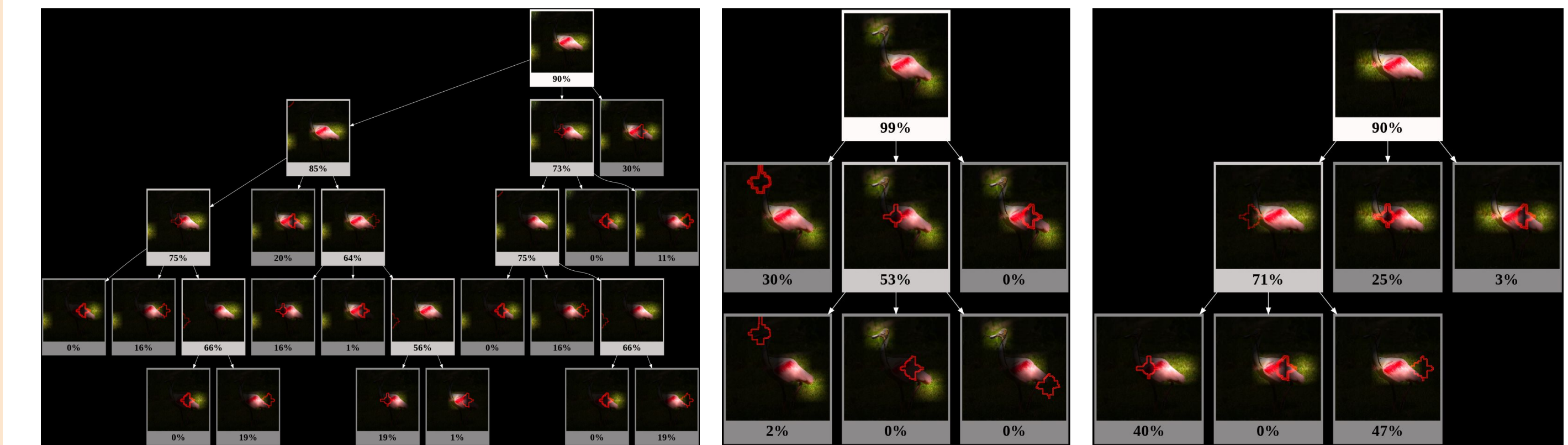
Sub-Explanation Counting

Model	Confidence				
	≥ 90 %	≥ 80 %	≥ 70 %	≥ 60 %	≥ 50 %
Swin-T	3.16	21.01	85.19	255.47	601.72
ResNet50	3.21	11.17	40.70	145.19	373.76
VGG19	3.19	13.86	57.59	158.42	439.52

Quantitative results of Sub-Explanation Counting



A few example distributions of MSE sizes for different algorithms on random images



(a) Swin-T (b) VGG19 (c) ResNet50
Sub-explanations of different models on an image of the **Spoonbill** class

Conclusion

- Swin-T is more robust to occlusion than CNNs using informative regions generated by iGOS++ and beam search method of SAG
- Swin-T can handle CNNs regions but CNNs can not handle Swin-T regions

This research is supported by National Science Foundation (1751402)