

## Problem Definition and Contribution

**Goal**  
Improving the deep learning model interpretability on the perspective of the quality of saliency predictions, while maintaining the model discriminative power.

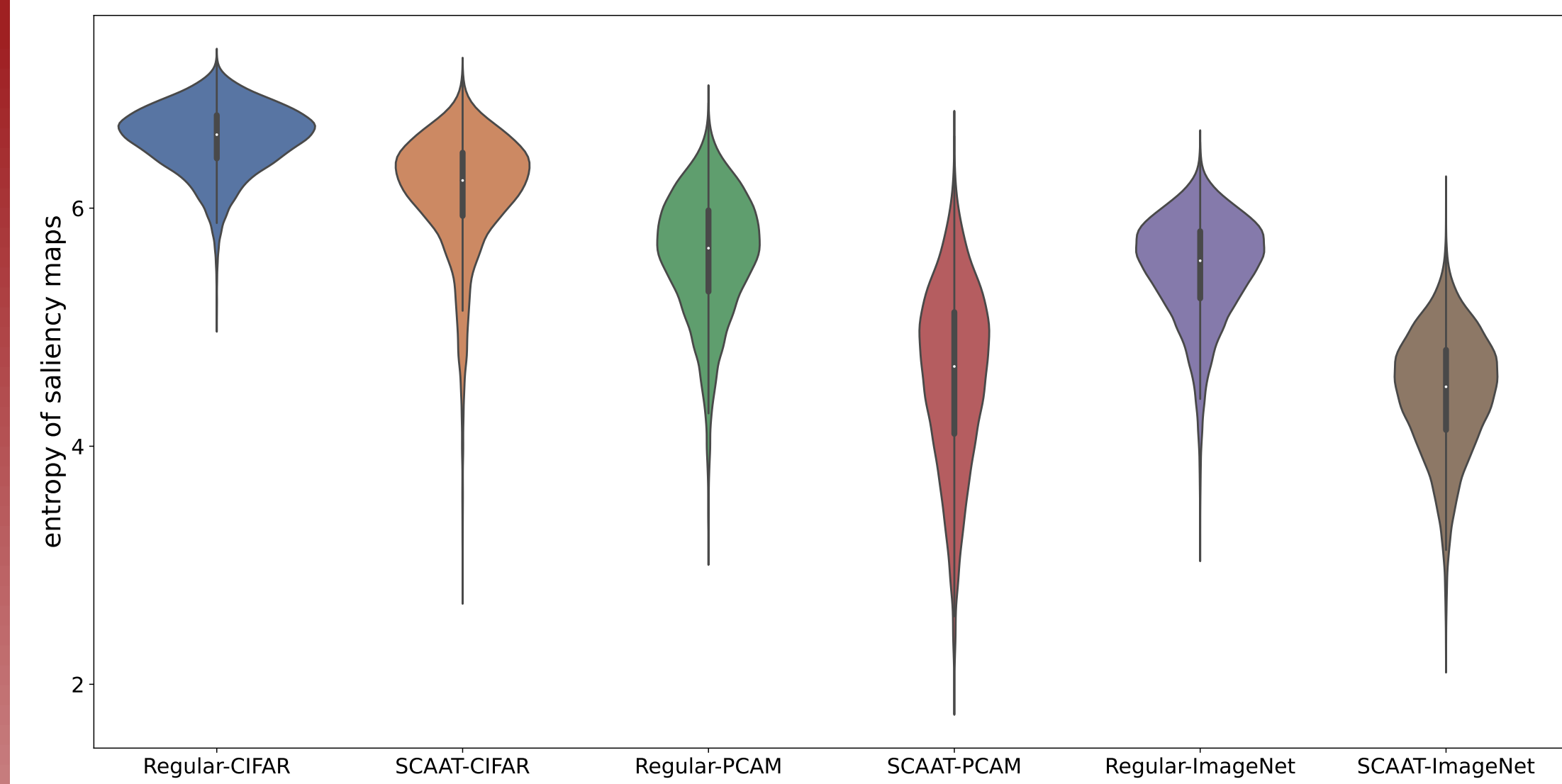
- Key Contributions**
- A novel model-agnostic adaptive adversarial training framework which improves the interpretability of deep neural networks **without changing the network architectures**. The method can be generalized to various models and domains.
  - Adaptive perturbation searching method which can **estimate the ratio of irrelevant features for each sample** then balance the learning performance and the quality of saliency predictions.
  - SCAAT significantly improves the model interpretability in measures of saliency map sparsity and faithfulness on various natural and pathological image datasets, while **barely sacrificing the predictive performance of the models**.

## Dataset

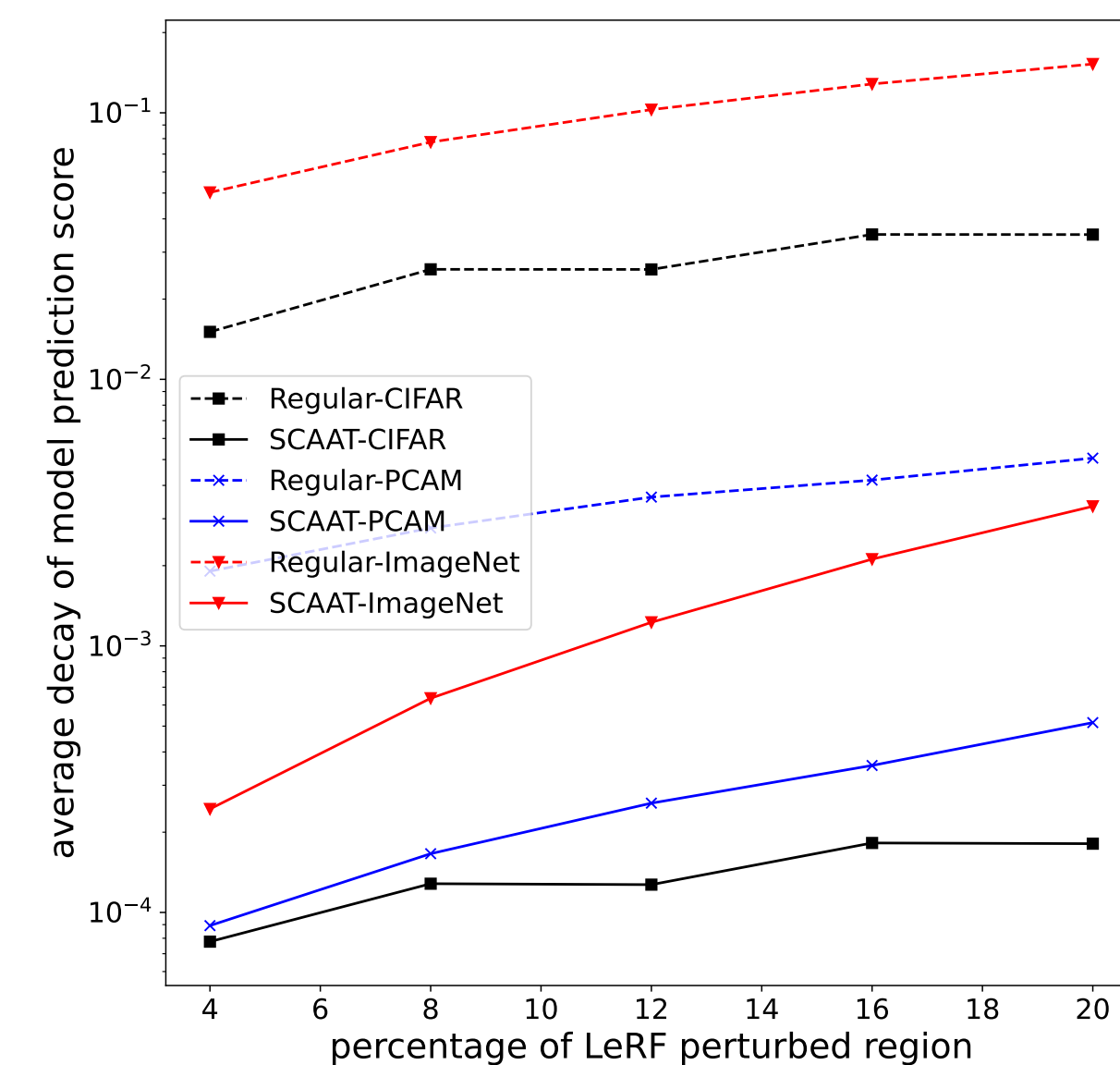
We conducted detailed experiments on CIFAR-10 and ImageNet-1k dataset for natural images, and PCAM dataset for pathological domain.

## Experiments and Results

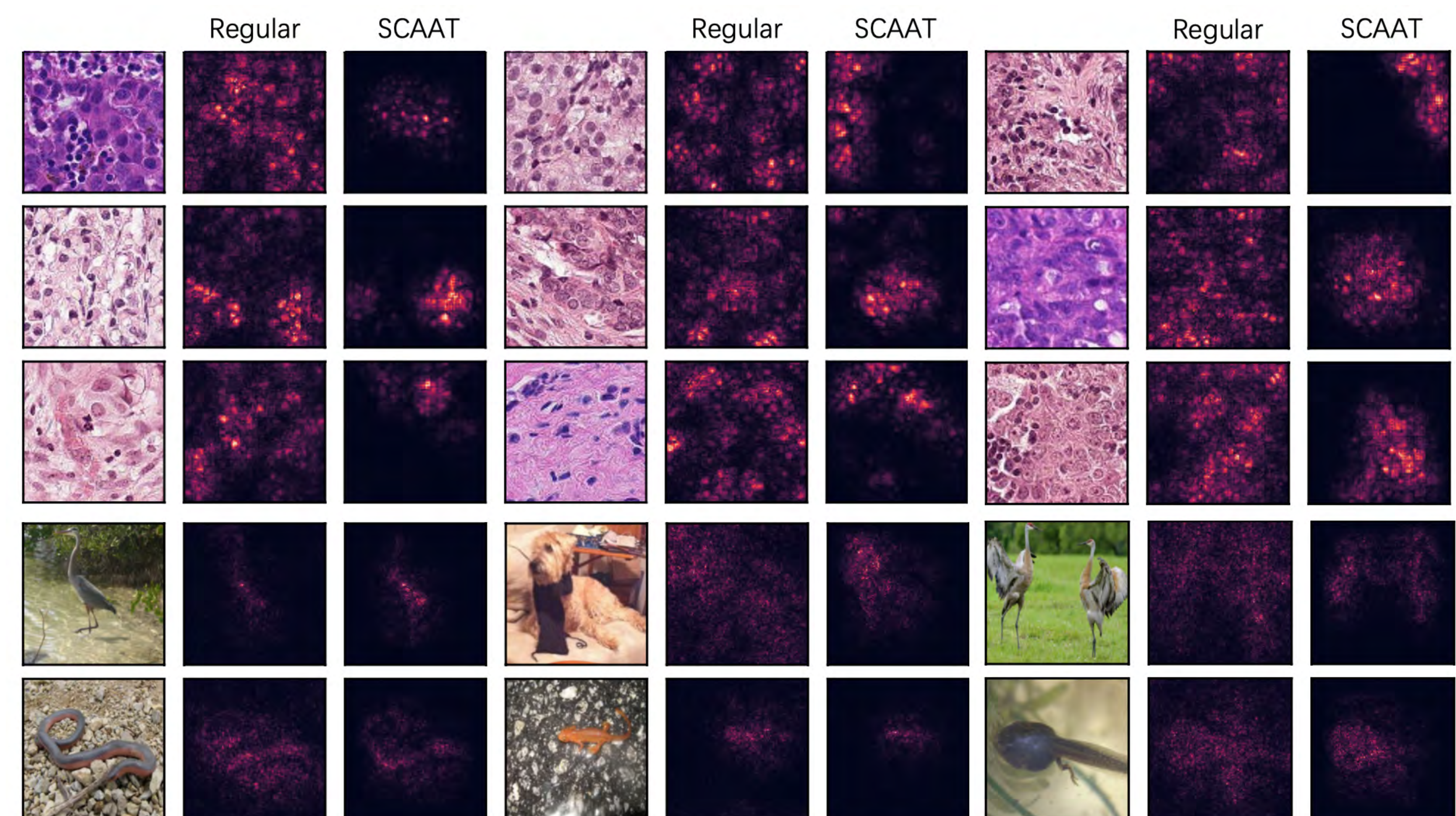
Saliency Map Entropy Distributions



Low-saliency Region Robustness Curve



Saliency Map Visualization on PCAM and ImageNet

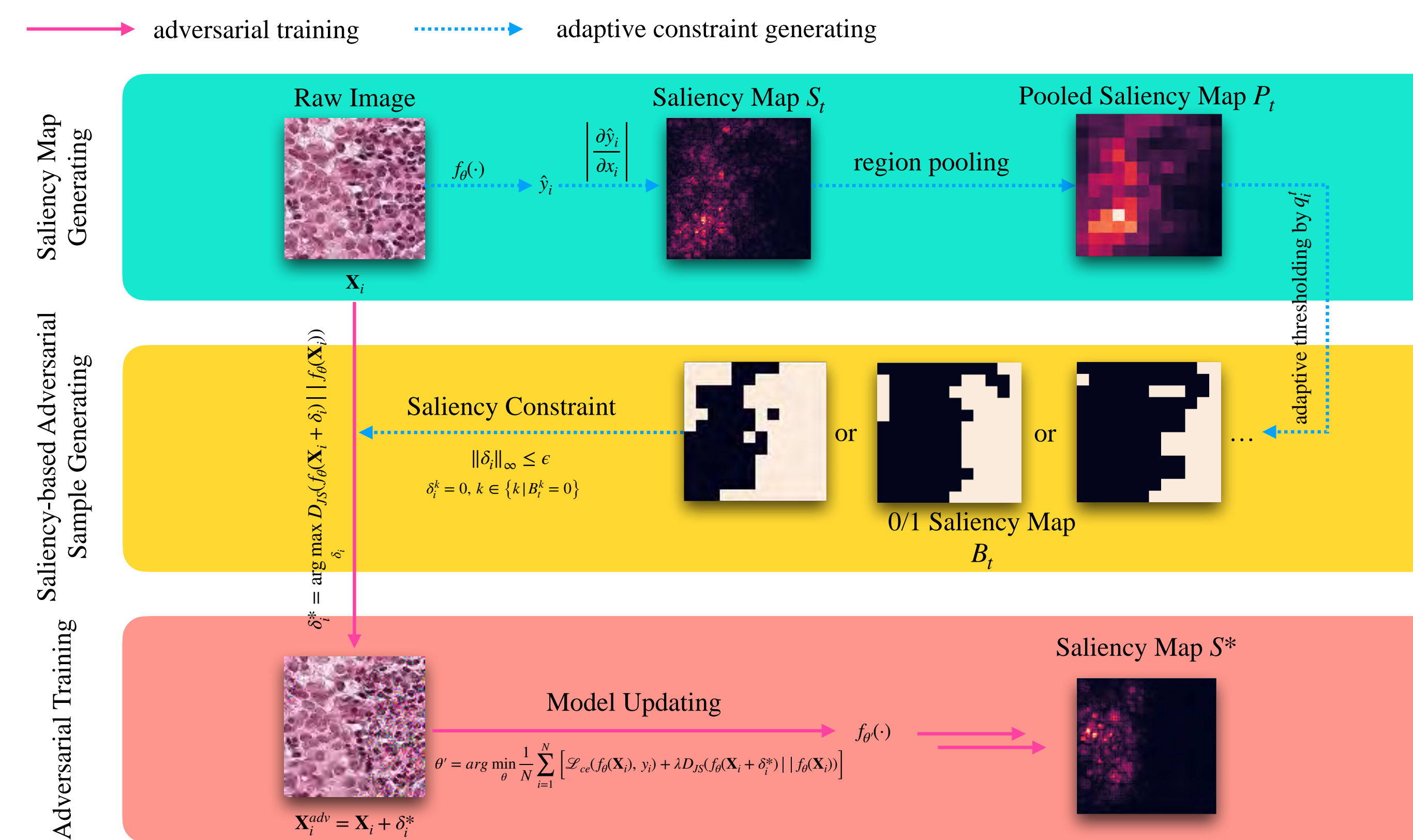


Saliency Map Quality Comparison Across Datasets and Saliency Methods

	Evaluation Metric	Vallina Grad		Smooth Grad		Integrated Grad	
		Regular	SCAAT	Regular	SCAAT	Regular	SCAAT
PCAM	Saliency Entropy ↓	5.61	<b>4.56</b>	5.60	<b>4.54</b>	4.93	<b>4.43</b>
	Saliency Size (Kbyte) ↓	2.48	<b>1.61</b>	2.45	<b>1.61</b>	2.23	<b>1.52</b>
	AOPC <sub>leRF</sub> ↓ (10 <sup>-3</sup> )	3.20	<b>0.23</b>	2.89	<b>0.23</b>	8.94	<b>0.21</b>
	AOPC <sub>rel</sub> ↑	78.1	<b>1030</b>	90.0	<b>982</b>	24.6	<b>938</b>
ImageNet-1k	Saliency Entropy ↓	5.49	<b>4.45</b>	5.12	<b>4.23</b>	4.98	<b>4.15</b>
	Saliency Size (Kbyte) ↓	13.2	<b>7.12</b>	12.9	<b>6.94</b>	12.8	<b>6.80</b>
	AOPC <sub>leRF</sub> ↓ (10 <sup>-3</sup> )	85.2	<b>0.98</b>	72.5	<b>0.93</b>	43.2	<b>1.21</b>
	AOPC <sub>rel</sub> ↑	3.84	<b>321</b>	4.66	<b>346</b>	4.21	<b>305</b>

## Method

Overall Pipeline of SCAAT



Irrelevant Feature Ratio Searching

**Require:**  $i$ : Index of sample,  $iter$ : Iteration index;  $N^{warm-up}$ : Warm-up iterations;  
**Require:**  $q_{max}, q_{min}$ : Boundary values for  $q$ ;  $\gamma$ : Discretization for  $q$  searching.  
**if**  $iter \leq N^{warm-up}$  **then**  
    Set  $q_i' = q_i$   
**else if**  $f_{\theta}(X_i^{adv})$  predicts as  $y_i$  **then**  
    Set  $q_i' = q_i + \gamma$   
**else**  
    Set  $q_i' = q_i - \gamma$   
**end if**  
Set  $q_i'' = \min(\max(q_i', q_{min}), q_{max})$   
**return**  $q_i''$