

# **DANCE: Enhancing saliency maps using decoys**

Yang Lu<sup>1,\*</sup> Wenbo Guo<sup>2,\*</sup> Xinyu Xing<sup>2</sup>, William Staffold Noble<sup>1,3</sup>

\* Equal contribution 1. Department of Genome Sciences, University of Washington, Seattle, WA 2. College of Information Sciences and Technology, The Pennsylvania State University, State College, PA 3. Paul G. Allen School of Computer Science and Engineering, University of Washington, Seattle, WA



## Overview

DANCE (Decoy enhANCEd saliency) is a saliency method which aims to tackle the following limitations with theoretical guarantee:

Gradient saturation

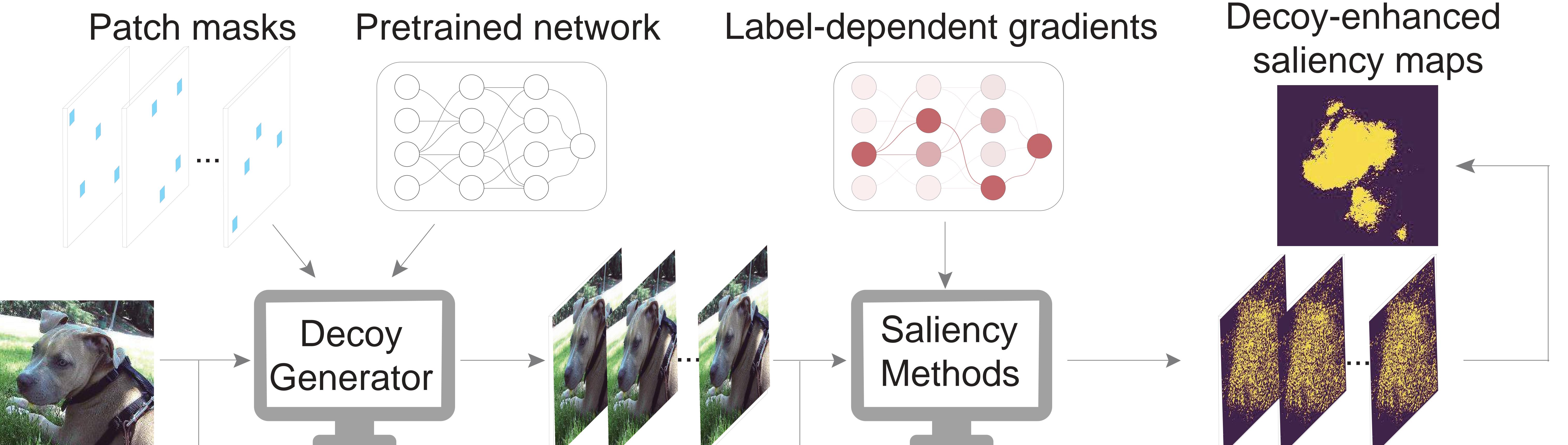
Strong joint evidences together with other neighbor pixels. Weak marginal evidence for individual pixel.

Importance isolation

The gradient is calculated by fixing other features. Smoothness in input doesn't hold in saliency maps.

Perturbation sensitivity

Imperceivable noises can drastically change the saliency maps.



Original image

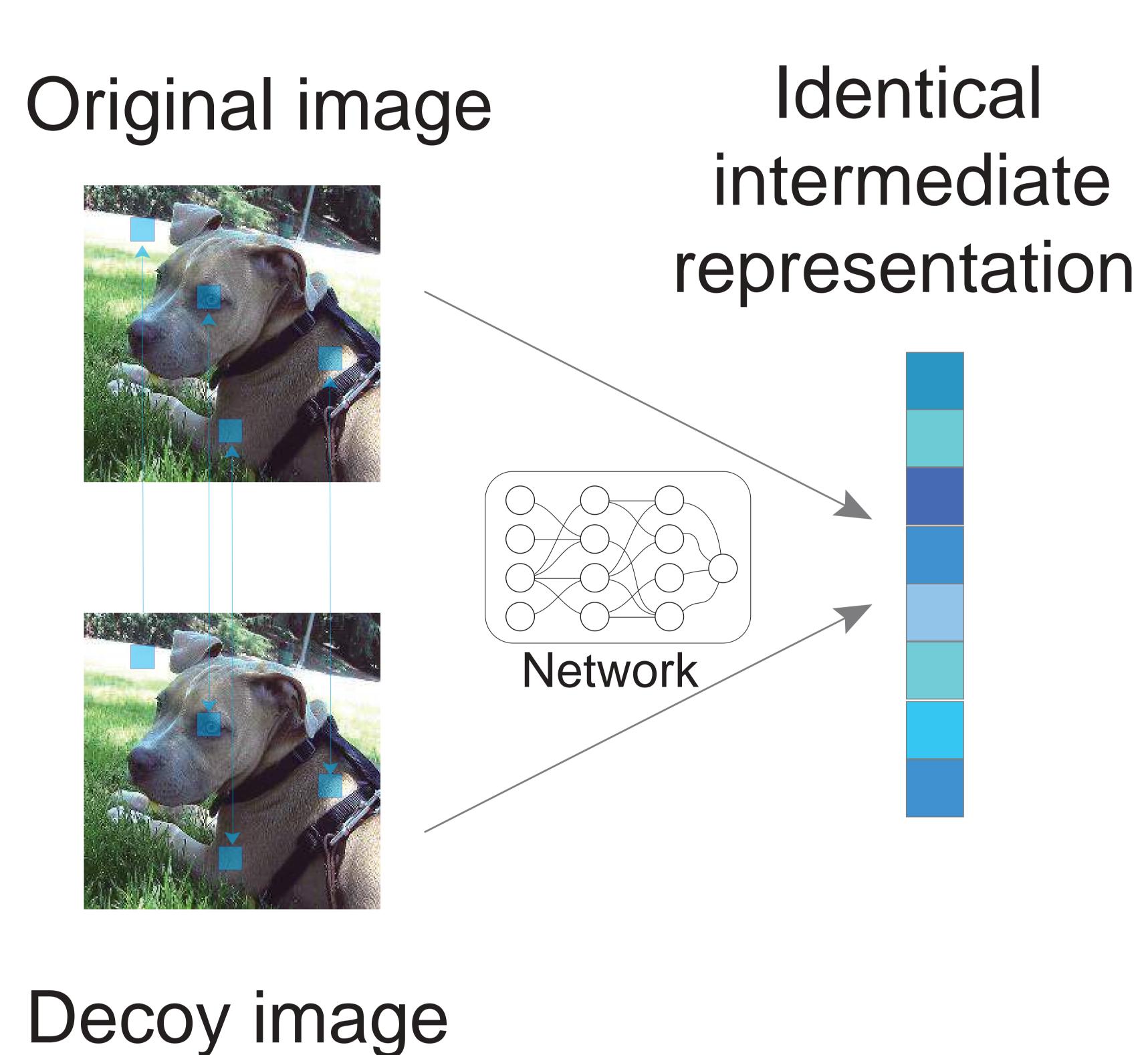
## Decoy images

Saliency maps of decoys

### **Decoy definition**

Decoy is a perturbed variant of the input sample which is designed to preserve the the neural network's intermediate representation.

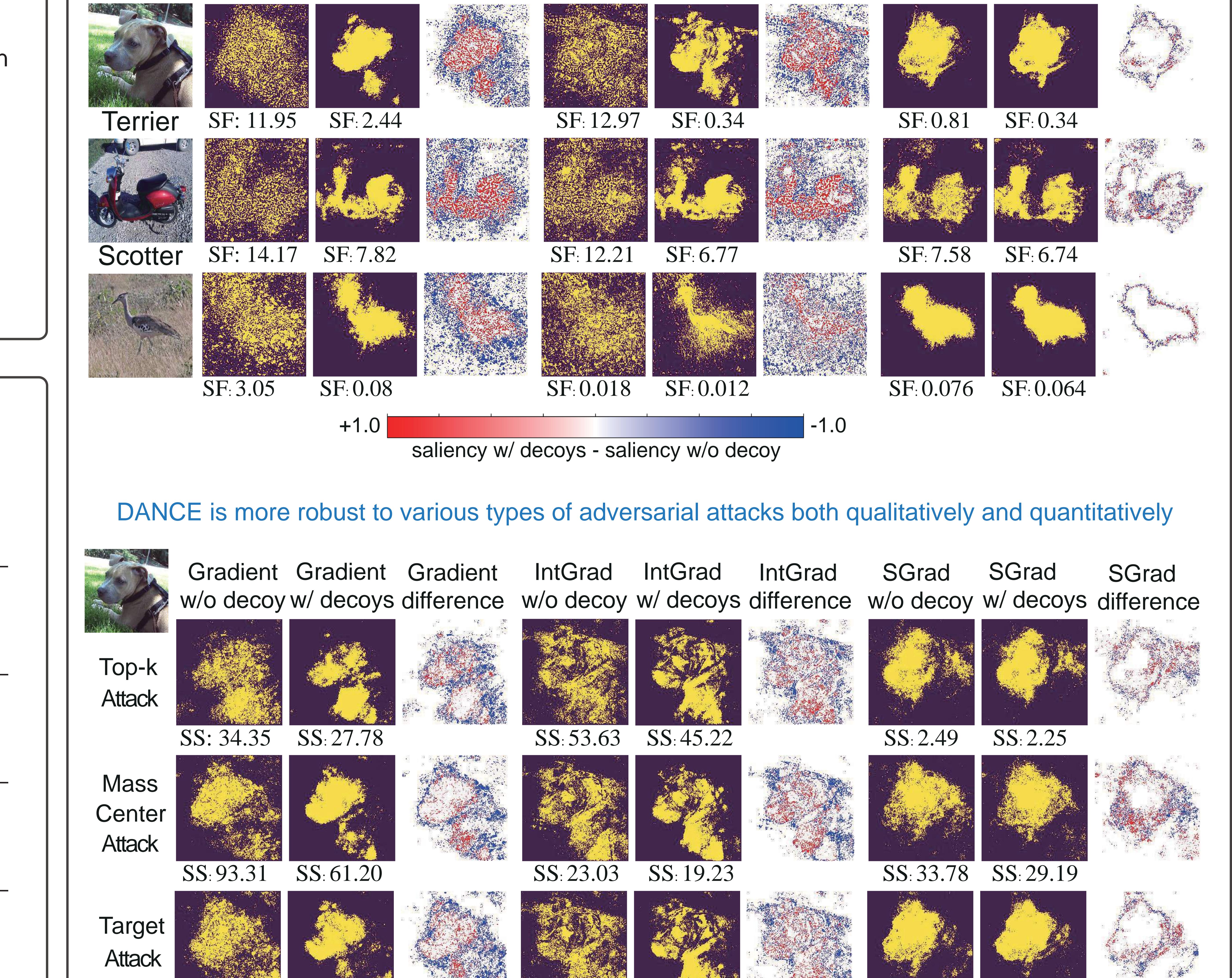
- Decoy aims to model the input variation from sensor noise or adversarial attacks.
- Both input and decoys are indistinguishable to the model prediction.
- Decoy cannot be out-of-distribution.
- Decoy can be efficiently constructed by solving an optimization problem.



## Performance on ImageNet dataset

## DANCE achieves more coherent saliency maps both qualitatively and quantitatively

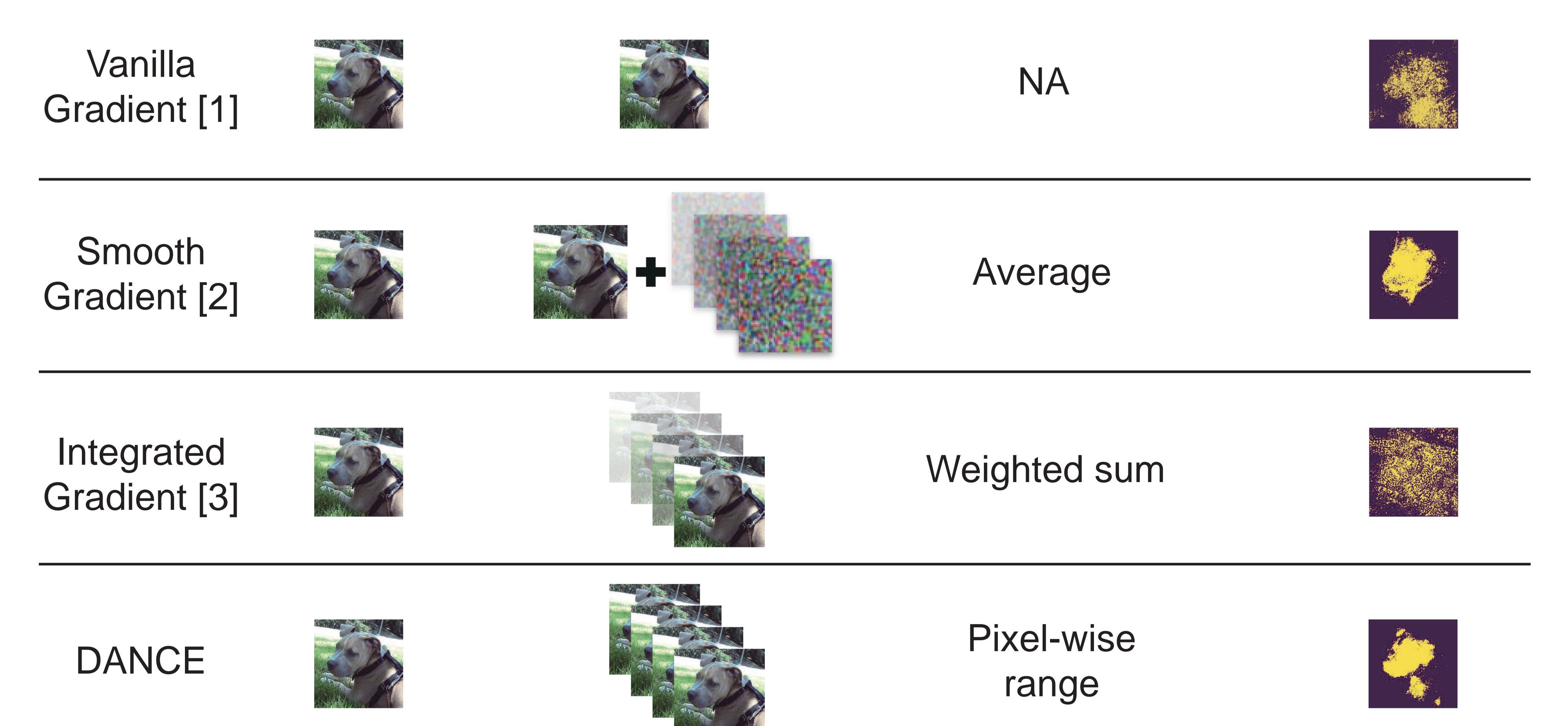
Gradient Gradient Gradient IntGrad IntGrad IntGrad SGrad SGrad SGrad w/o decoy w/ decoys difference w/o decoy w/ decoys difference w/o decoy w/ decoys difference



#### Decoy aggregation

DANCE, as well as other saliency methods, can be summarized into a "variations-and-aggregation" paradigm.

```
Input —> Input variants —> Saliency aggregation —> Output
Method
```



	$\overline{SS: 24.10}$ $\overline{SS: 16.21}$	SS: 30.50 SS: 24.81	SS: 27.15 $SS: 22.70$	
Evaluation metric	A control study on DANCE			
	Two components in DANCE, both decoys and range-based aggregation, are essential.			