

# Understanding Deep Neural Networks Via Smooth Masks and Extremal Perturbations

Ruth Fong\*, Mandela Patrick\*, Andrea Vedaldi

ICCV 2019

\* equal contribution

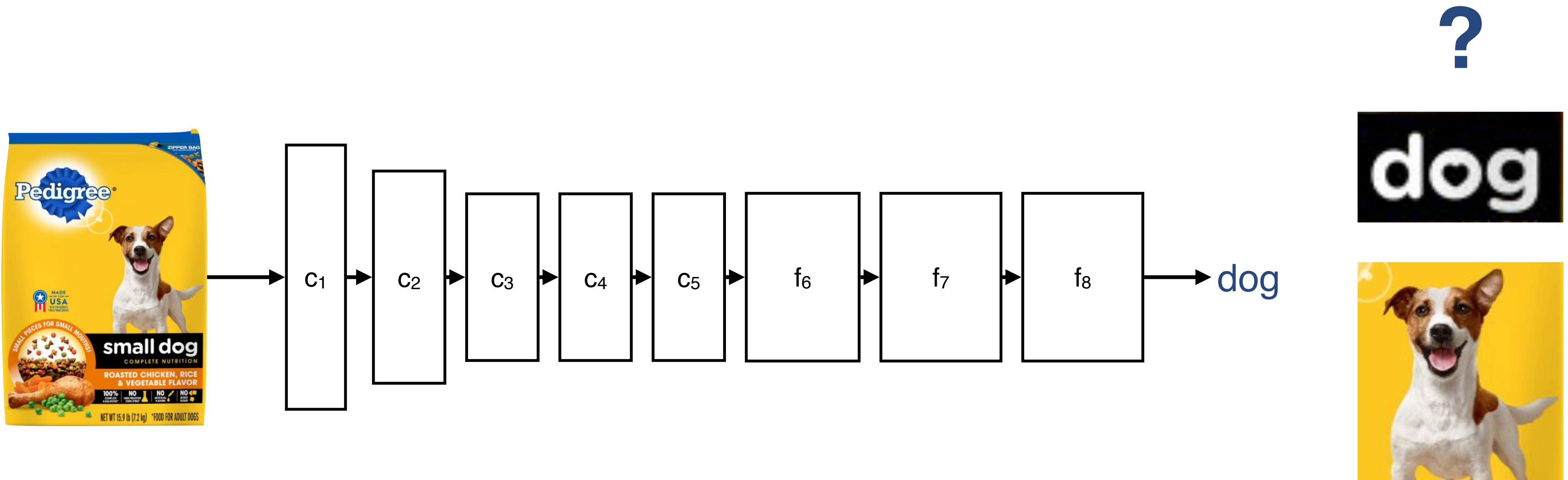


UNIVERSITY OF  
**OXFORD**

**FACEBOOK**

# Attribution

Where is the model looking?



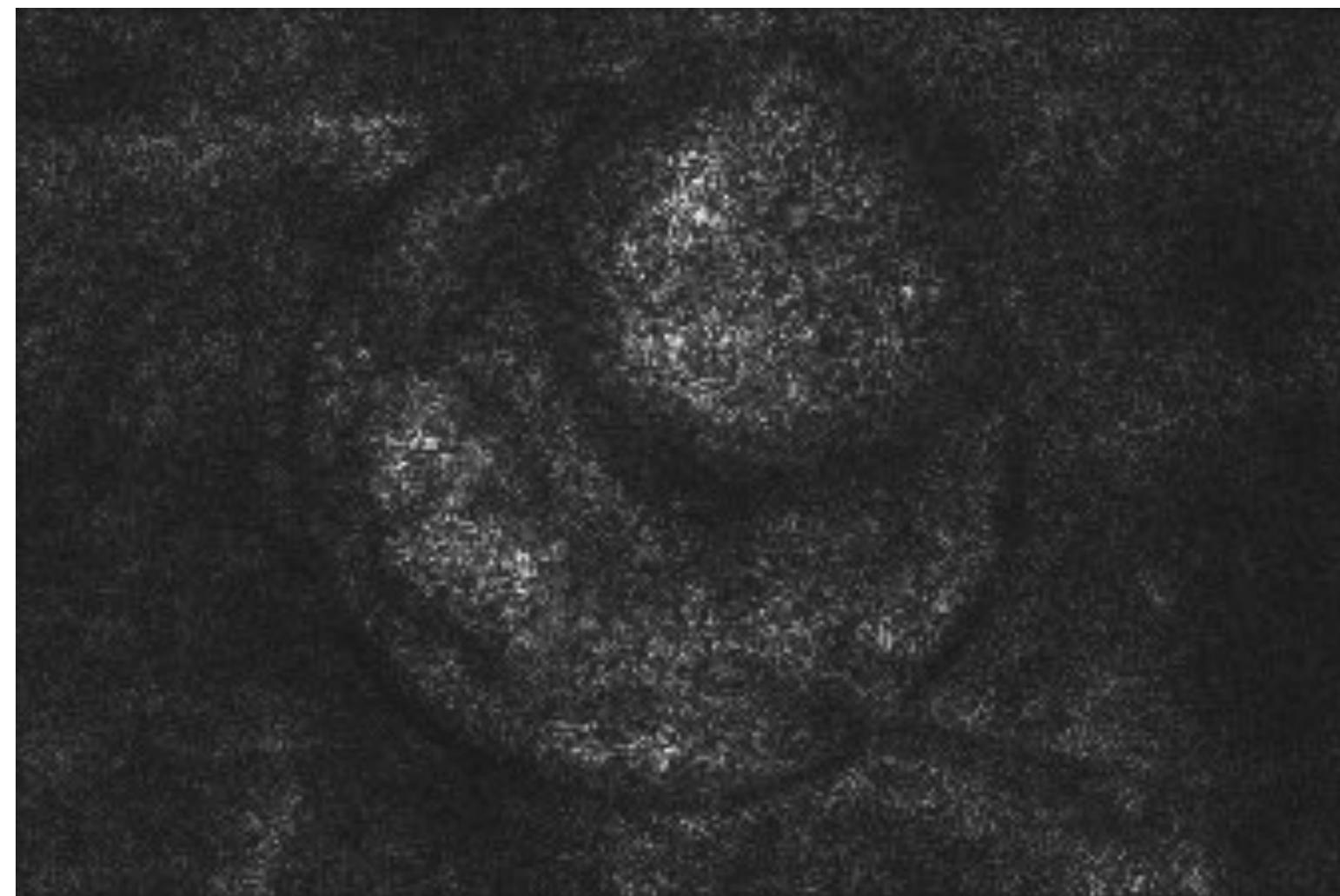
# Backpropogation

Combine network activations and gradients

Input



Gradient



GRAD-CAM



Fast, but difficult to characterize

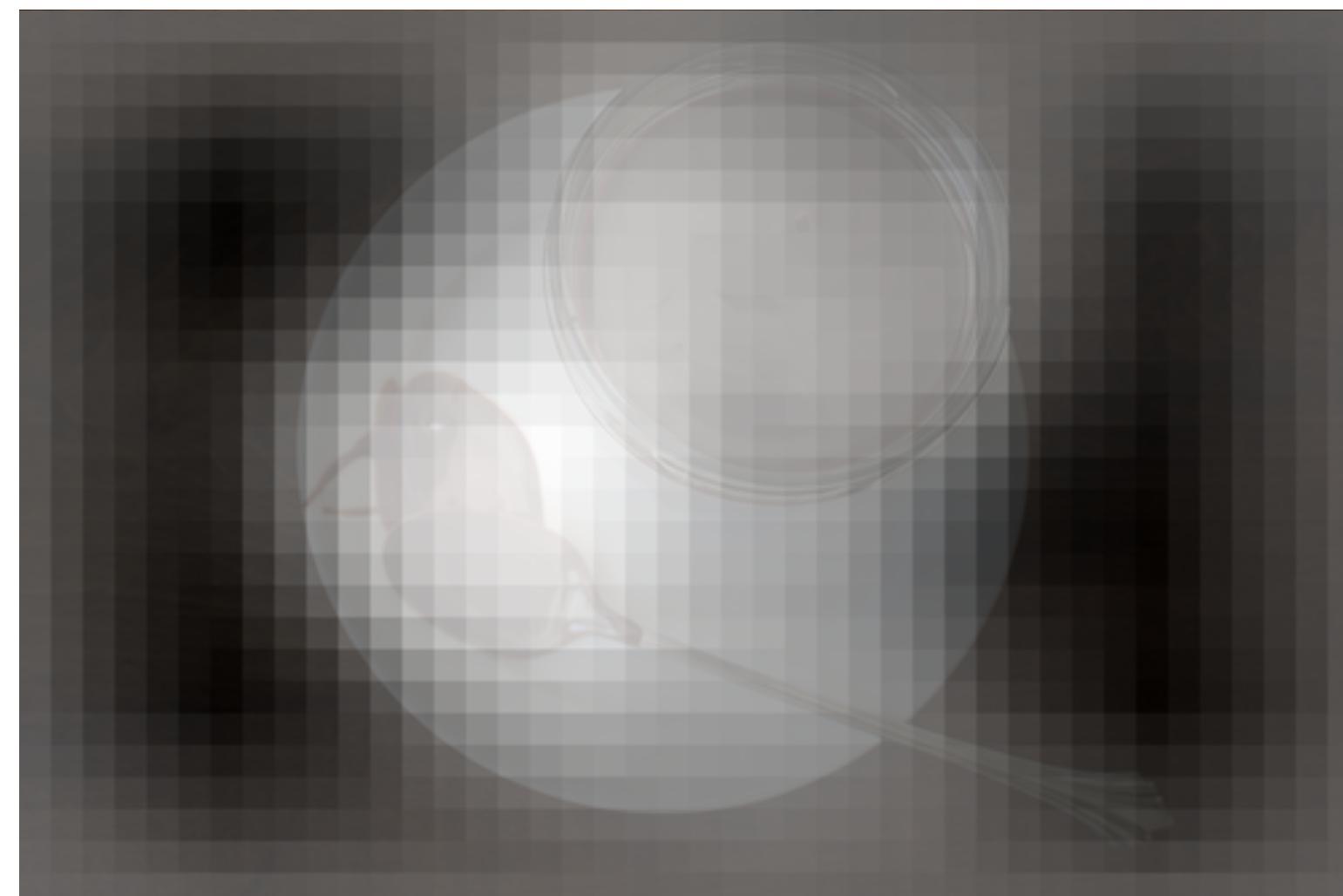
# Perturbation

Change the input and observe the effect on the output

Input



Occlusion

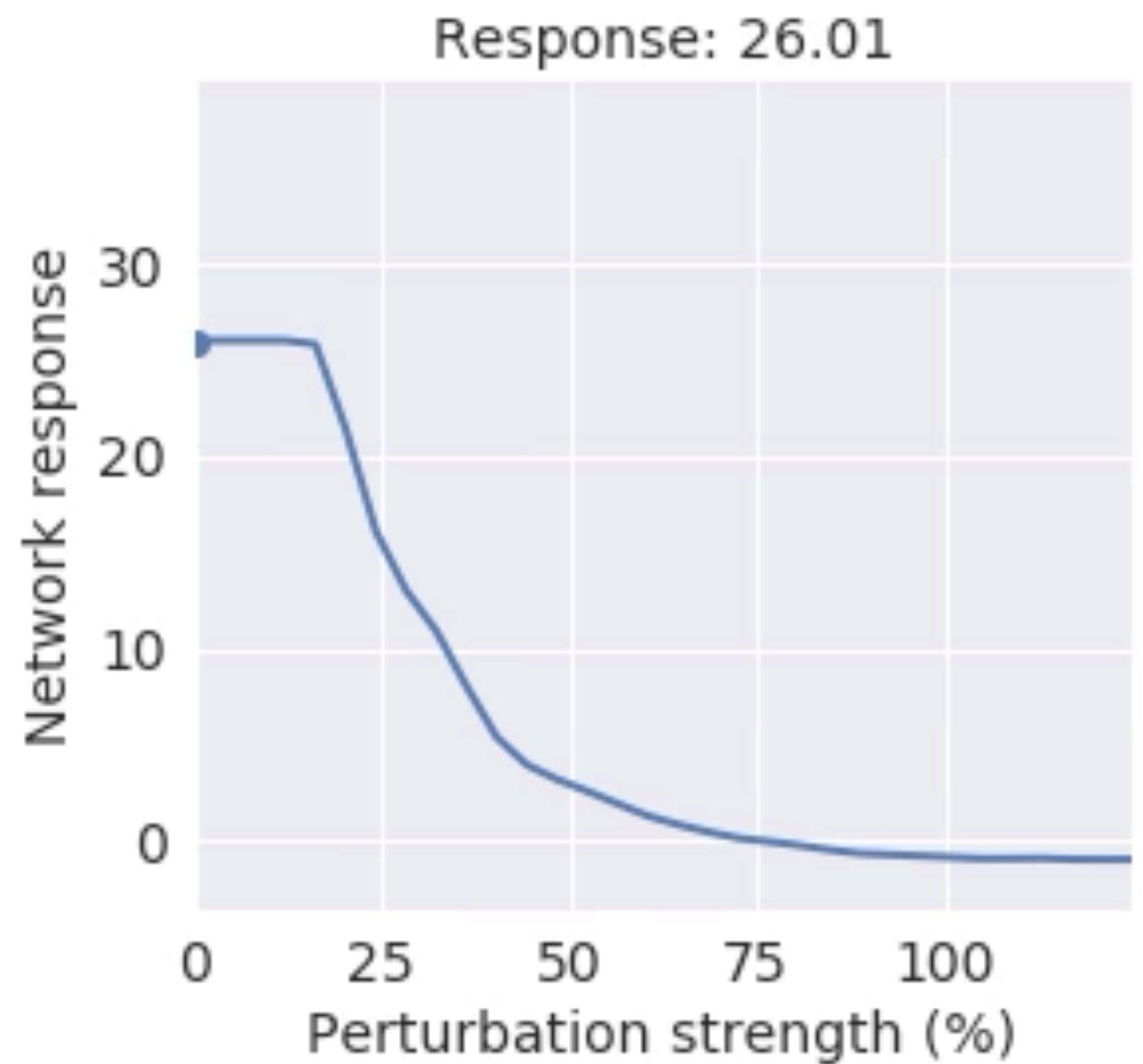
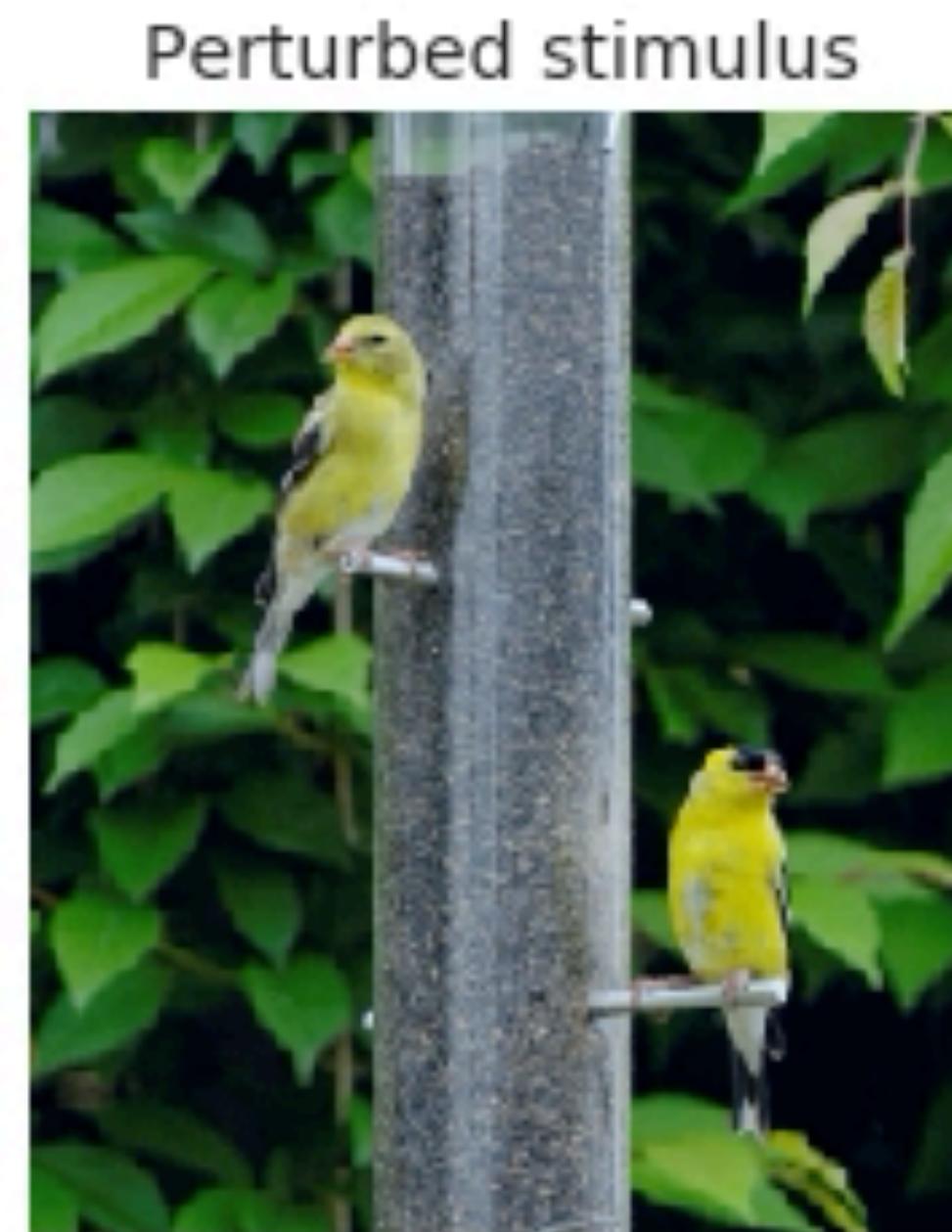


RISE

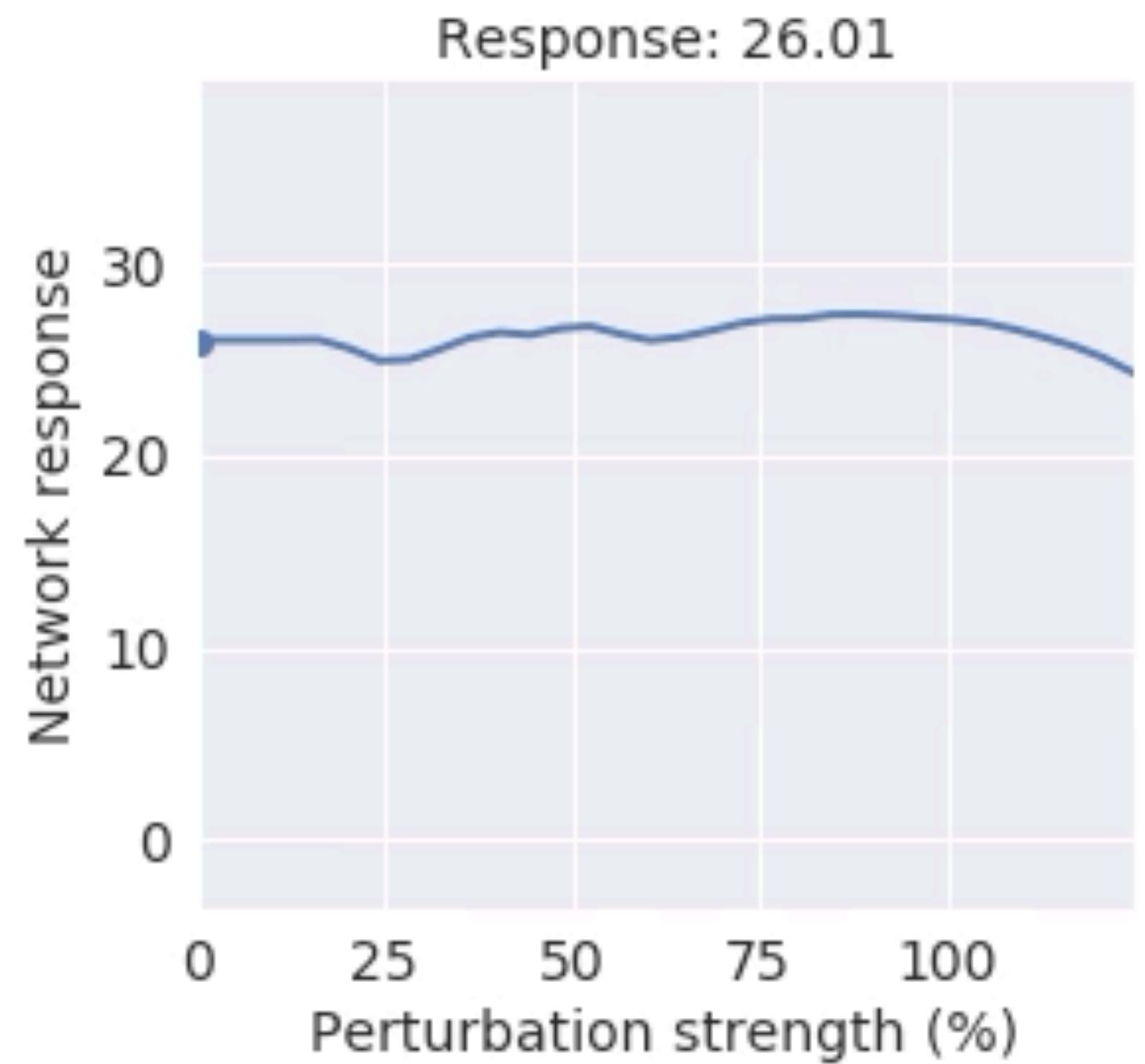


Clear meaning, but can only test a small number of occlusion patterns

## Blur everywhere => response suppressed



# Preserve 10% => response preserved



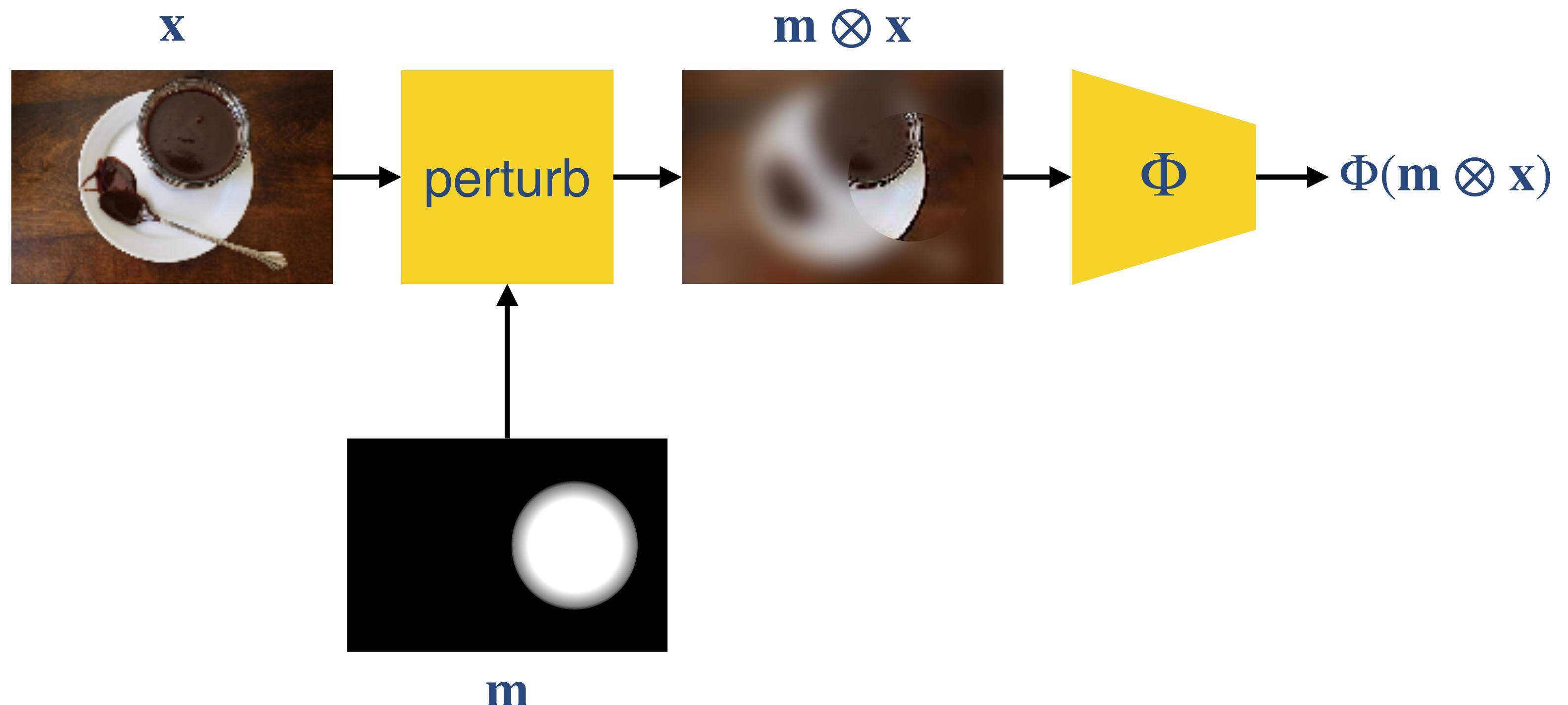
# Extremal Perturbations

A mask is optimized to

maximally excite the network:

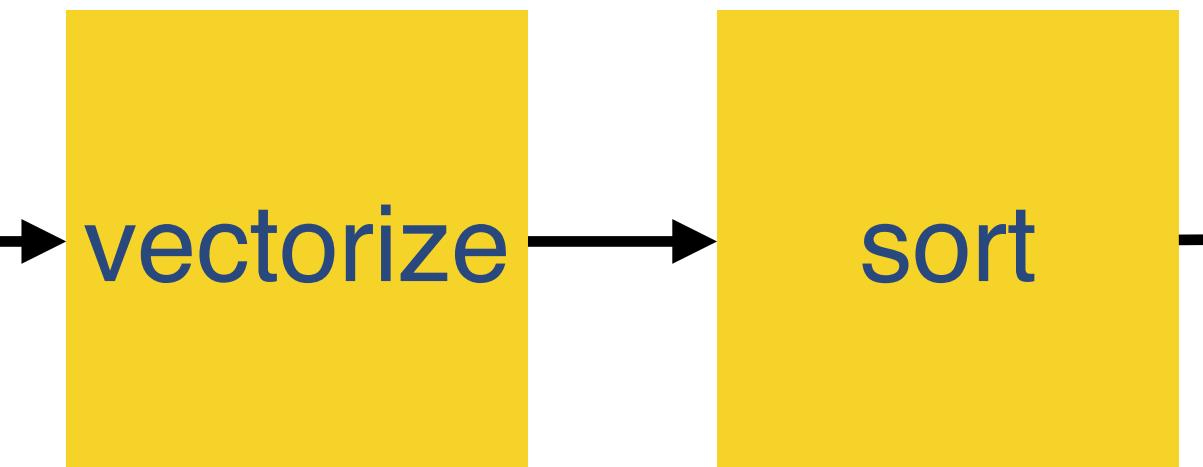
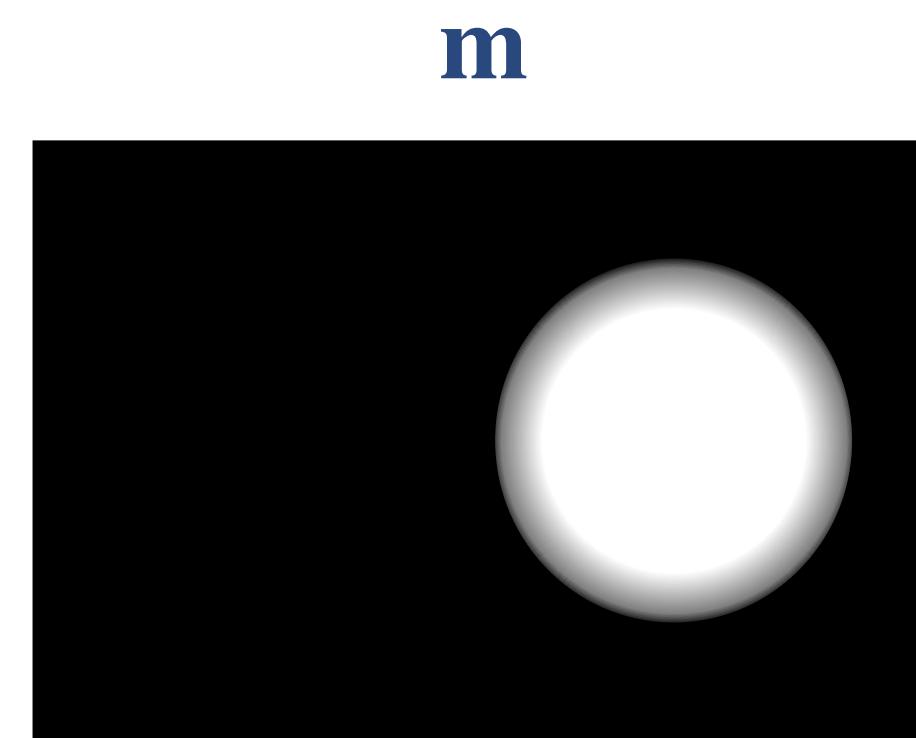
$$\underset{\mathbf{m}}{\operatorname{argmax}} \Phi(\mathbf{m} \otimes \mathbf{x})$$

subject to  $\text{area}(\mathbf{m}) = a$



## Area Constraint

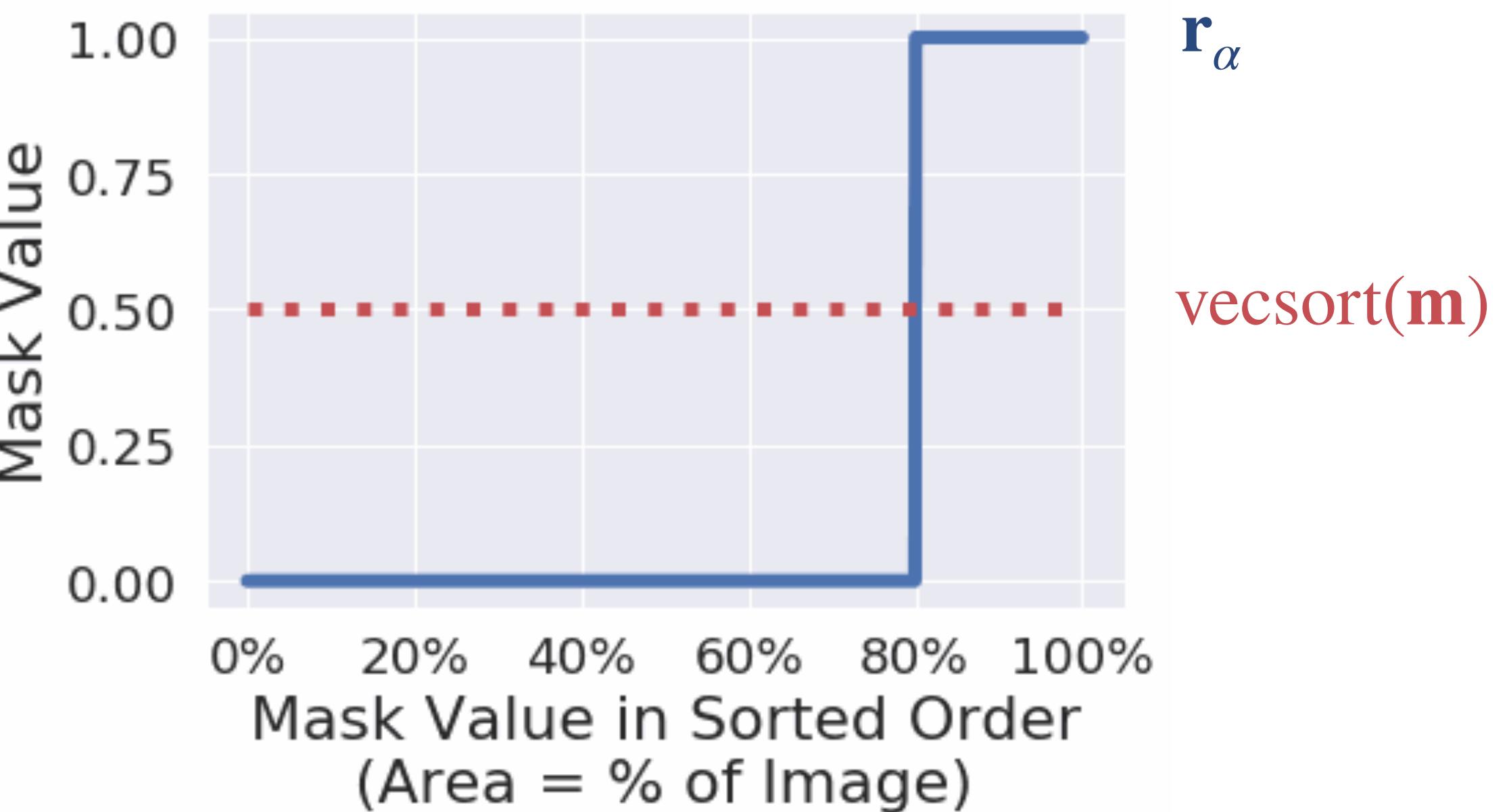
Optimizing w.r.t. to an area constraint is challenging



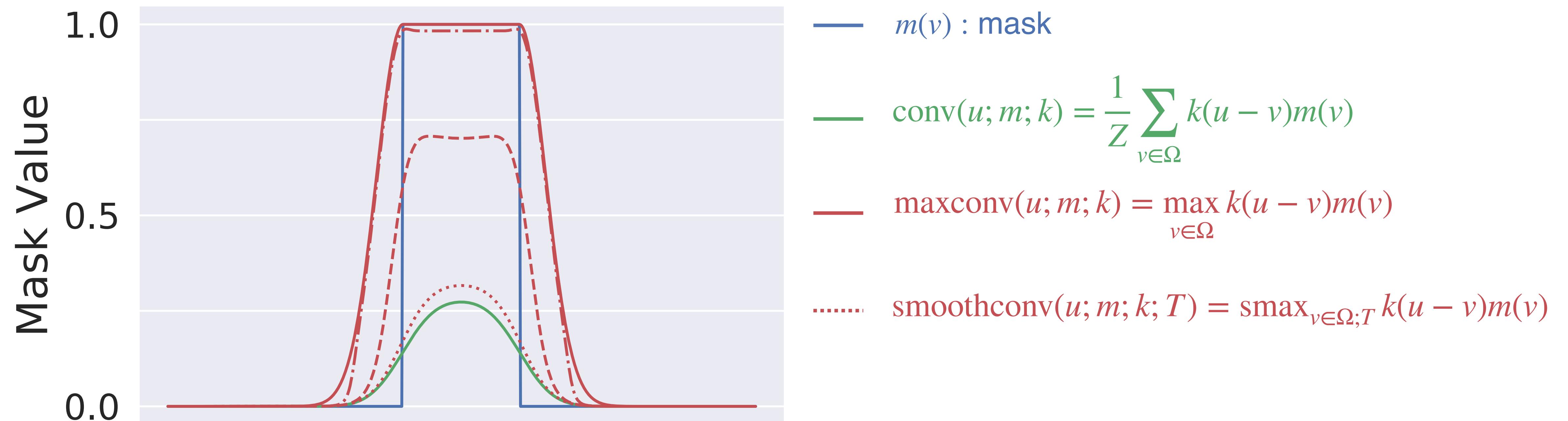
subject to  $\text{area}(\mathbf{m}) = a$

Here we re-formulate it as matching a **rank statistics**

$$L_{area} = \| \text{vecsrt}(\mathbf{m}) - \mathbf{r}_a \|^2$$

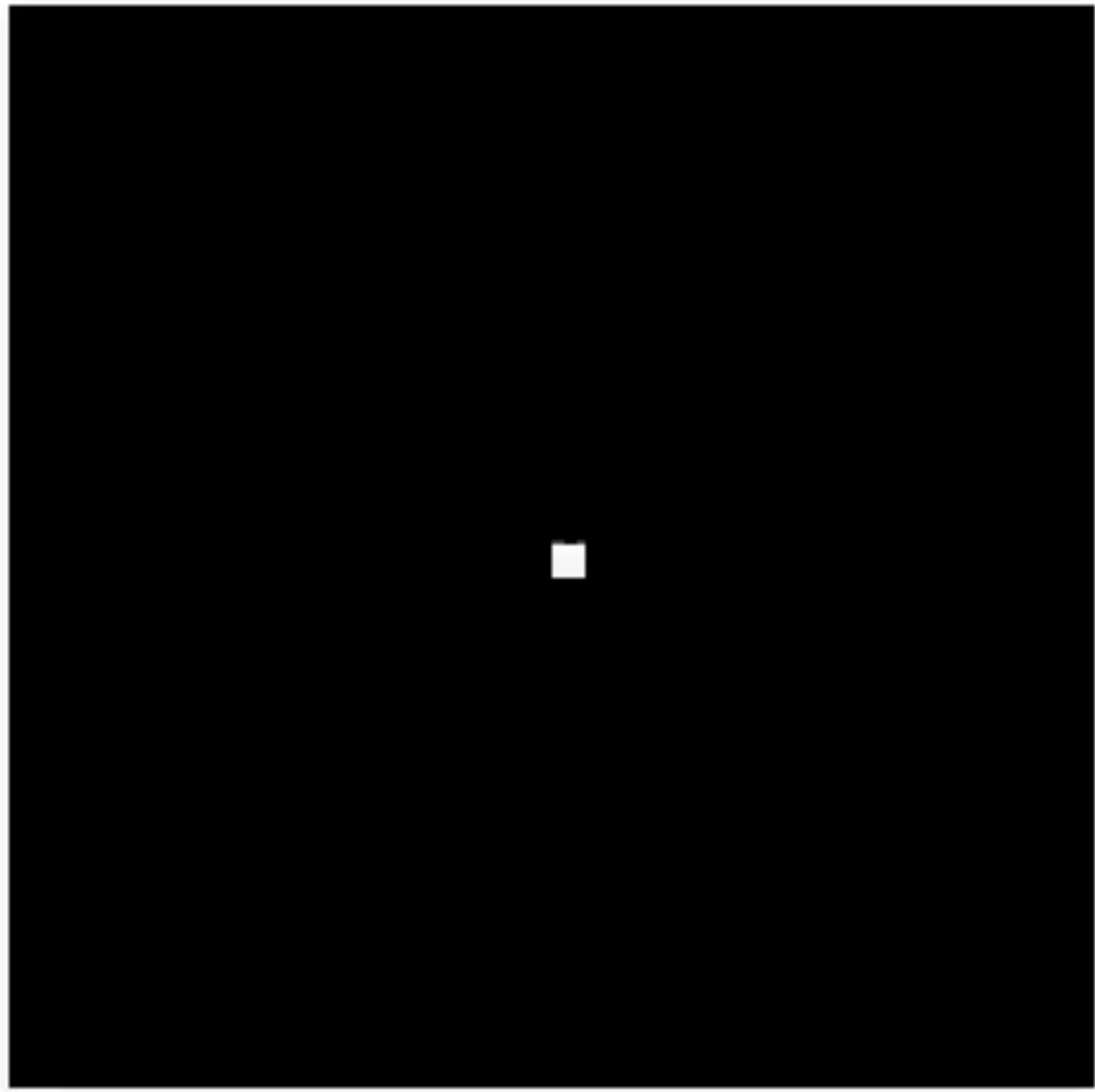


# Smooth Masks

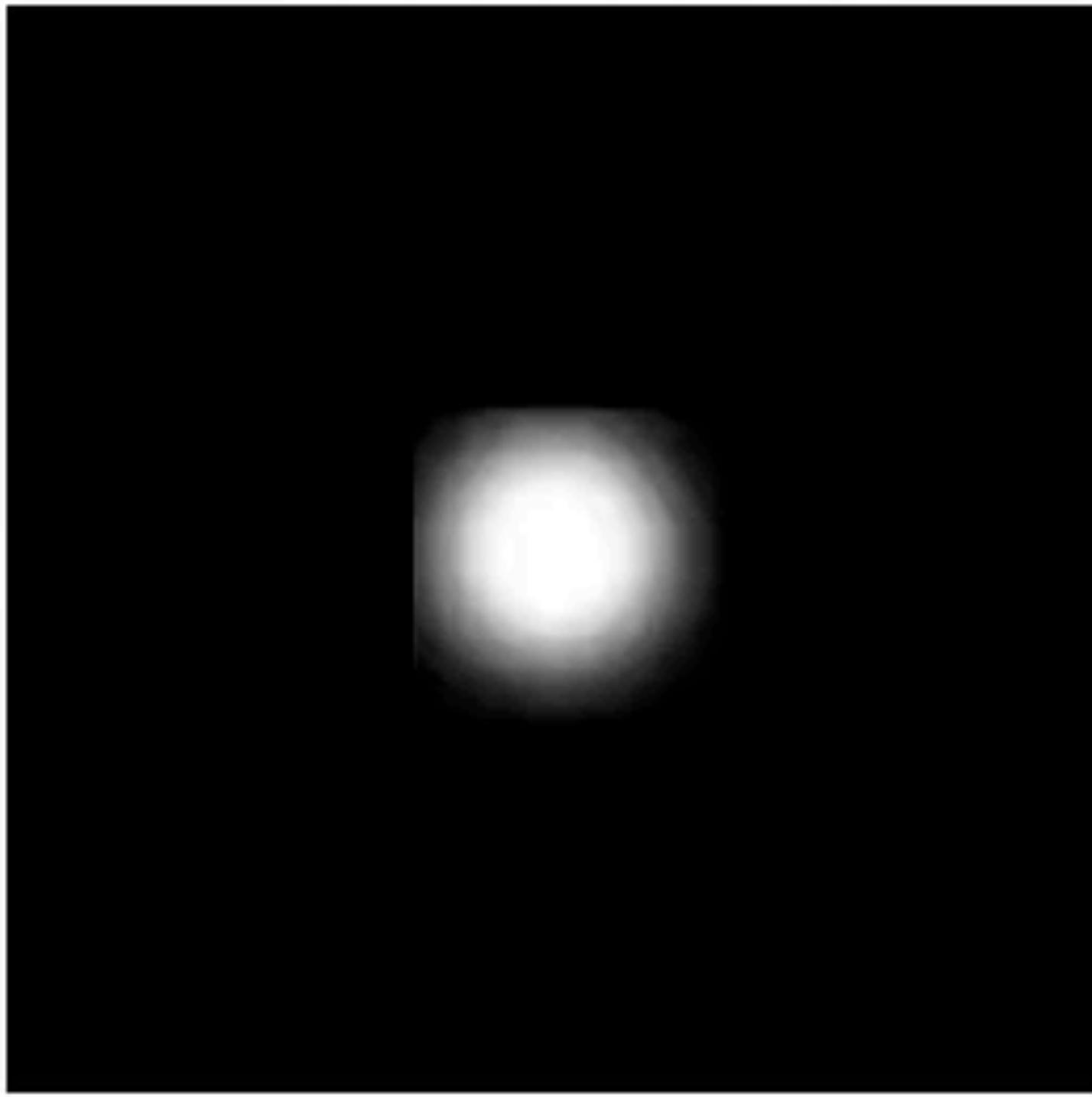


# Smooth Masks

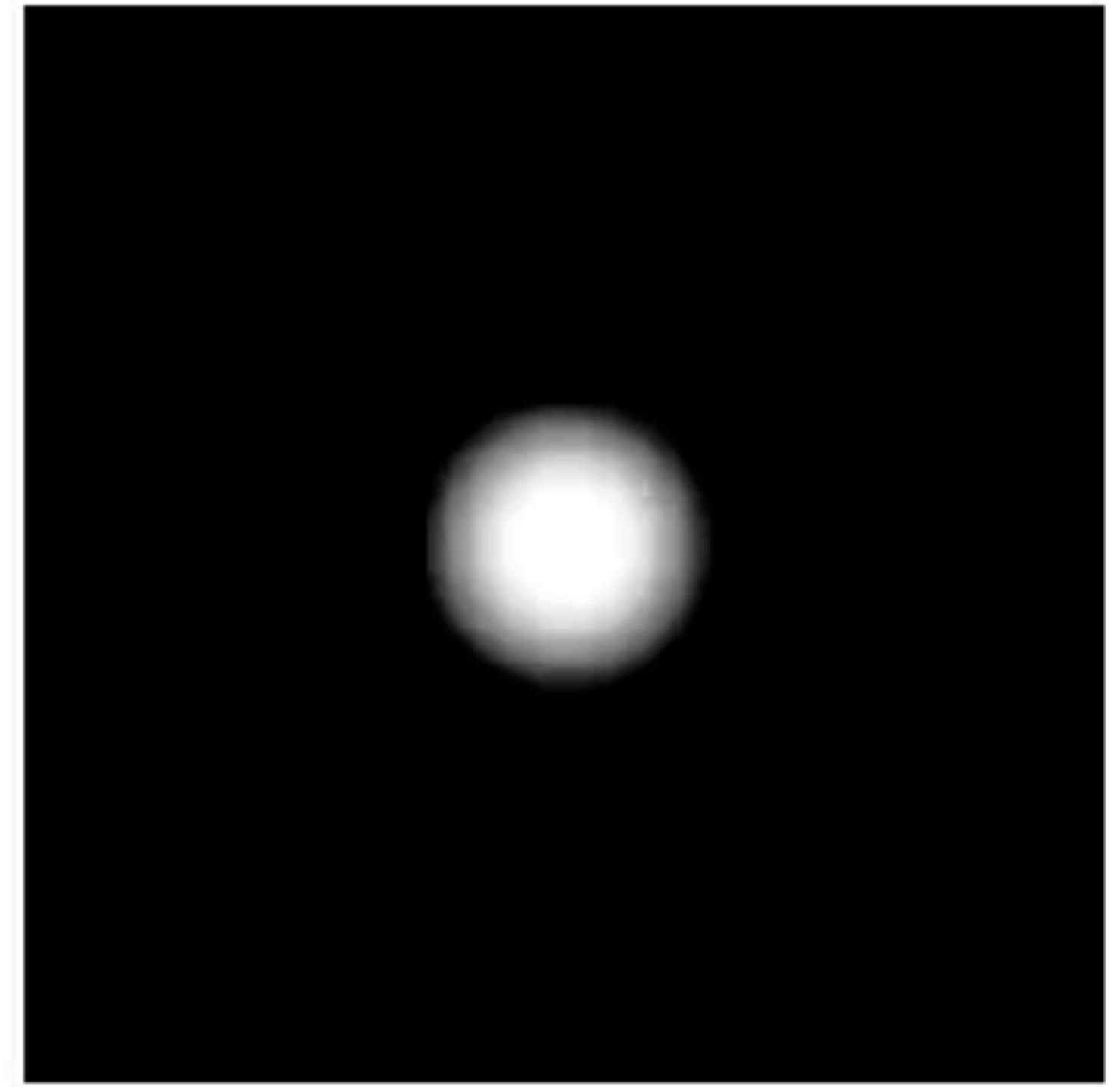
Mask parameters



Gaussian smoothing



Max-conv smoothing



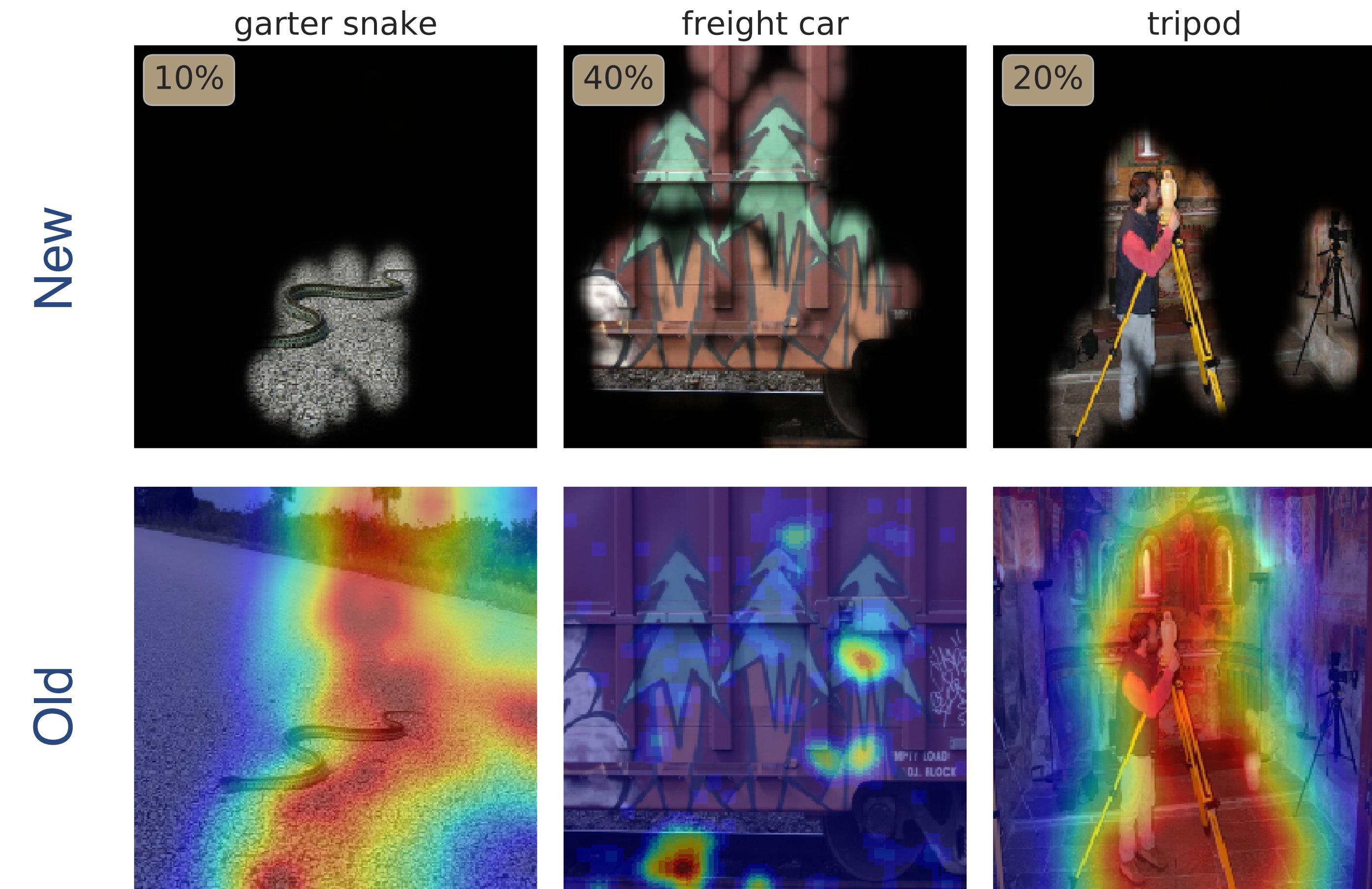
# Algorithm

1. Pick an area  $a$
2. Use SGD to solve the optimization problem for a large  $\lambda$ :

$$\underset{\mathbf{m}}{\operatorname{argmax}} \Phi(\operatorname{smooth}(\mathbf{m}) \otimes \mathbf{x}) - \lambda \|\operatorname{vecsorth}(\operatorname{smooth}(\mathbf{m})) - \mathbf{r}_a\|^2$$

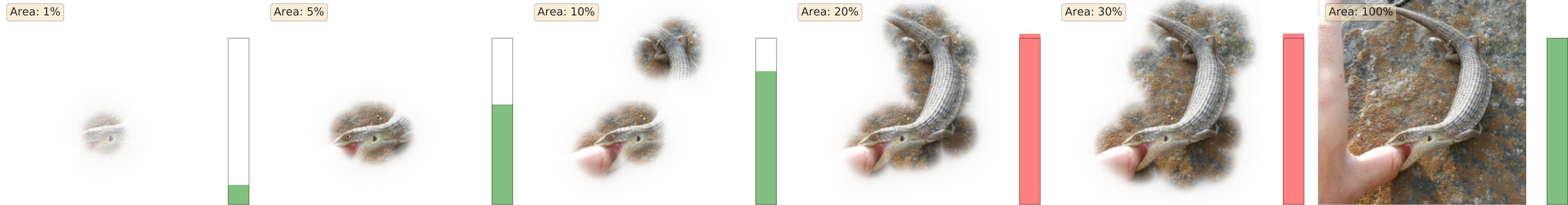
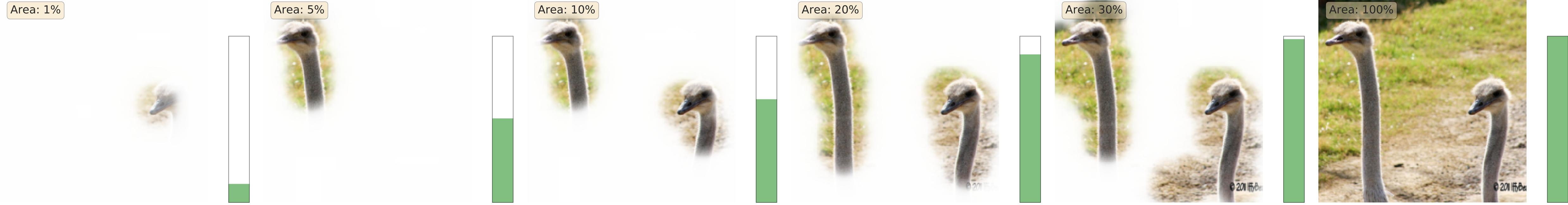
3. If needed, sweep  $a$  and repeat

# Comparison to prior work on “meaningful perturbations”



# Results

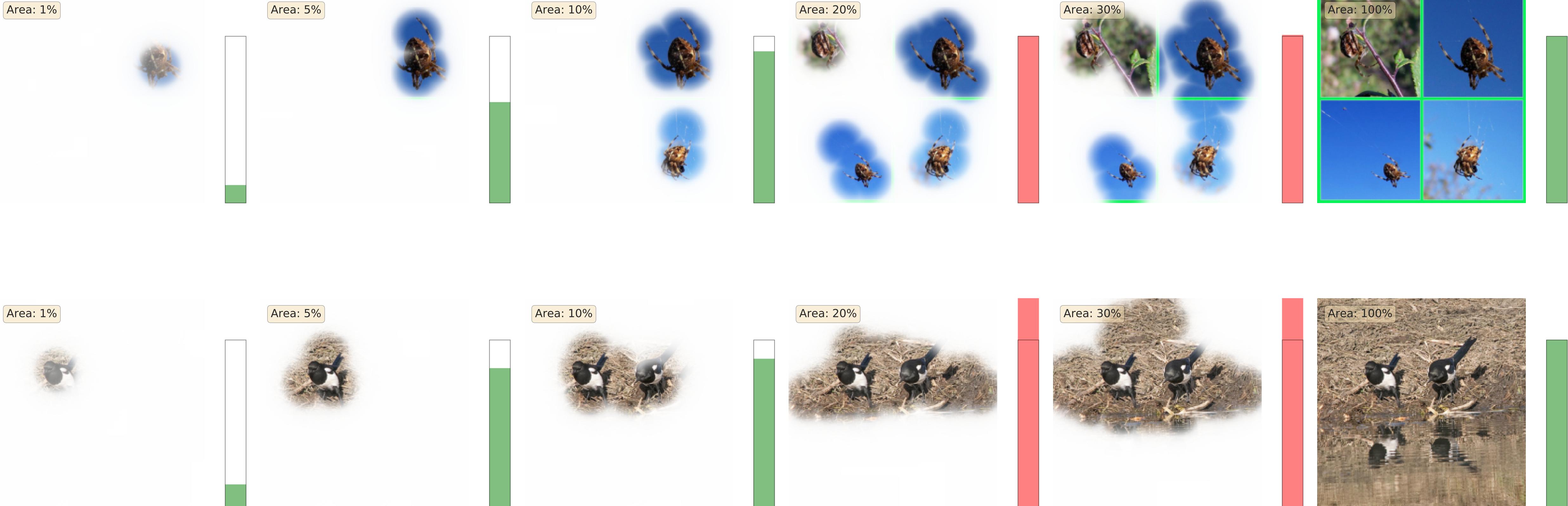
# Foreground evidence is usually sufficient



# Large objects are recognised by their details

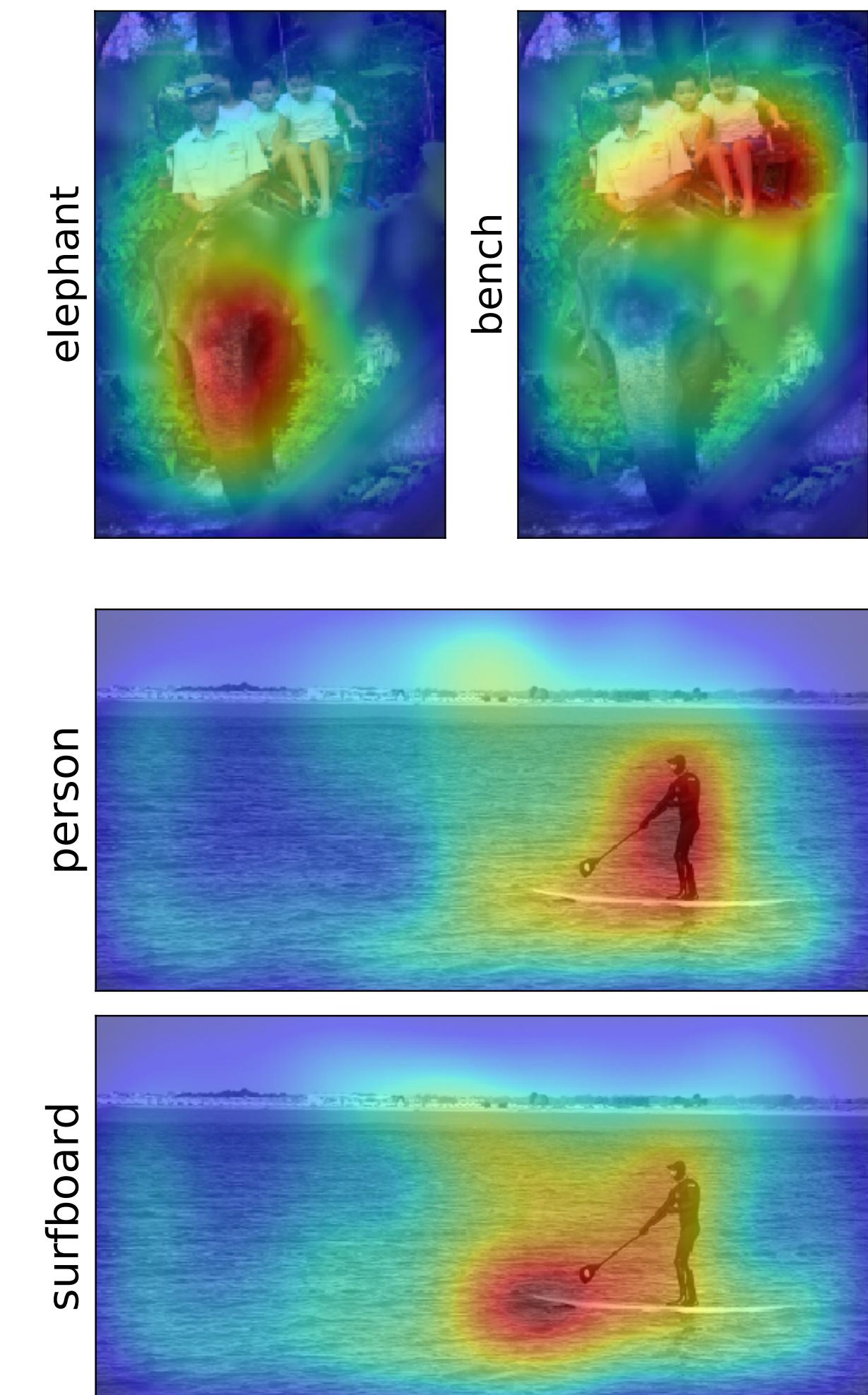


# Multiple objects contribute cumulatively



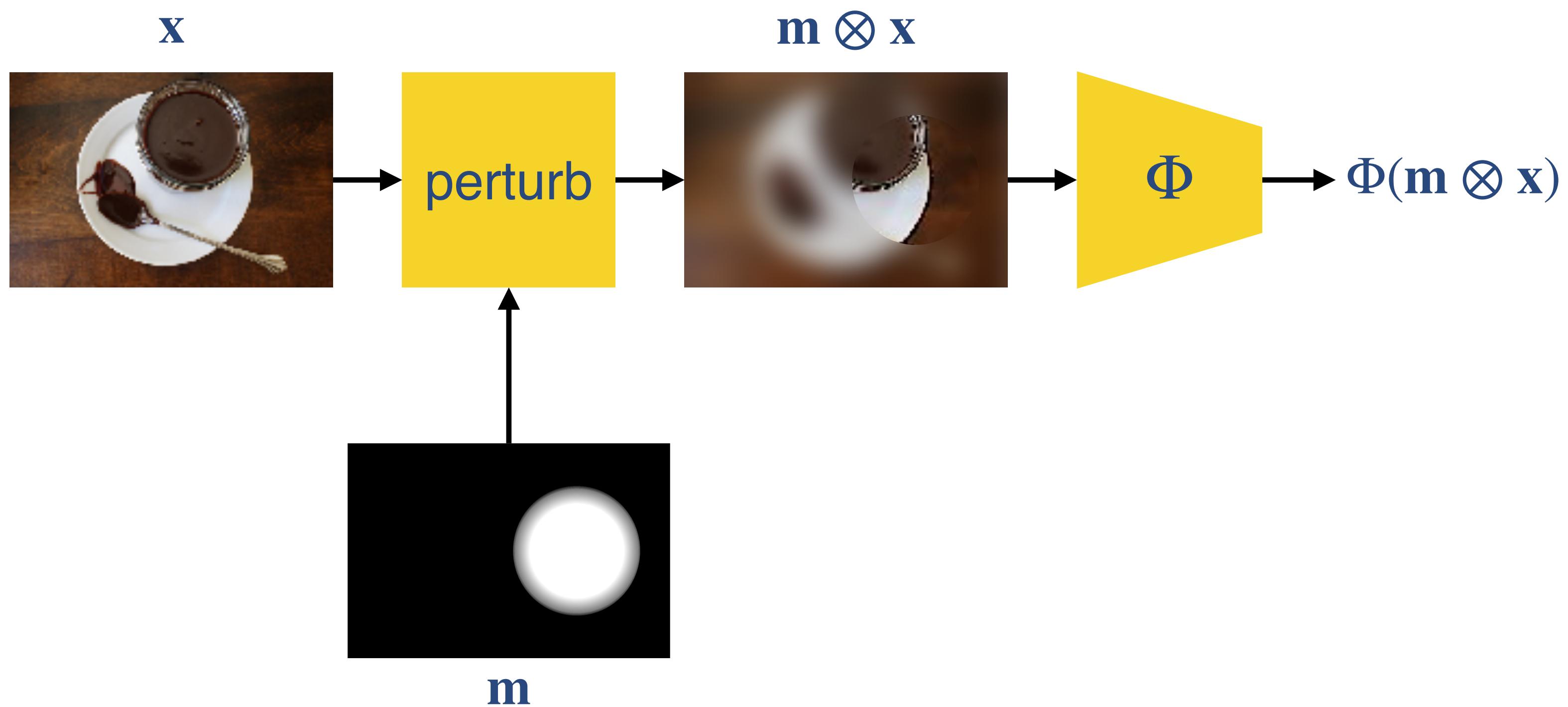
## Pointing game: weak localization

Method	<i>VOC07 Test (All/Diff)</i>		<i>COCO14 Val (All/Diff)</i>	
	<i>VGG16</i>	<i>ResNet50</i>	<i>VGG16</i>	<i>ResNet50</i>
Grad	76.3/56.9	72.3/56.8	37.7/31.4	35.0/29.4
DConv	67.5/44.1	68.6/44.7	30.7/23.0	30.0/21.9
Guid.	75.9/53.0	77.2/59.5	39.1/31.4	42.1/35.3
MWP	77.1/56.6	84.4/70.8	39.8/32.8	49.6/43.9
cMWP	79.9/66.5	<b>90.6/82.2</b>	49.7/44.3	<b>58.5/53.6</b>
RISE*	<u>87.3/—</u>	88.9/—	50.7/—	55.6/—
GCAM	86.6/74.0	<u>90.4/82.3</u>	<b>54.2/49.0</b>	<u>57.3/52.3</u>
Ours	<b>88.7/75.5</b>	86.3/73.4	<u>53.4/47.7</u>	55.7/46.9

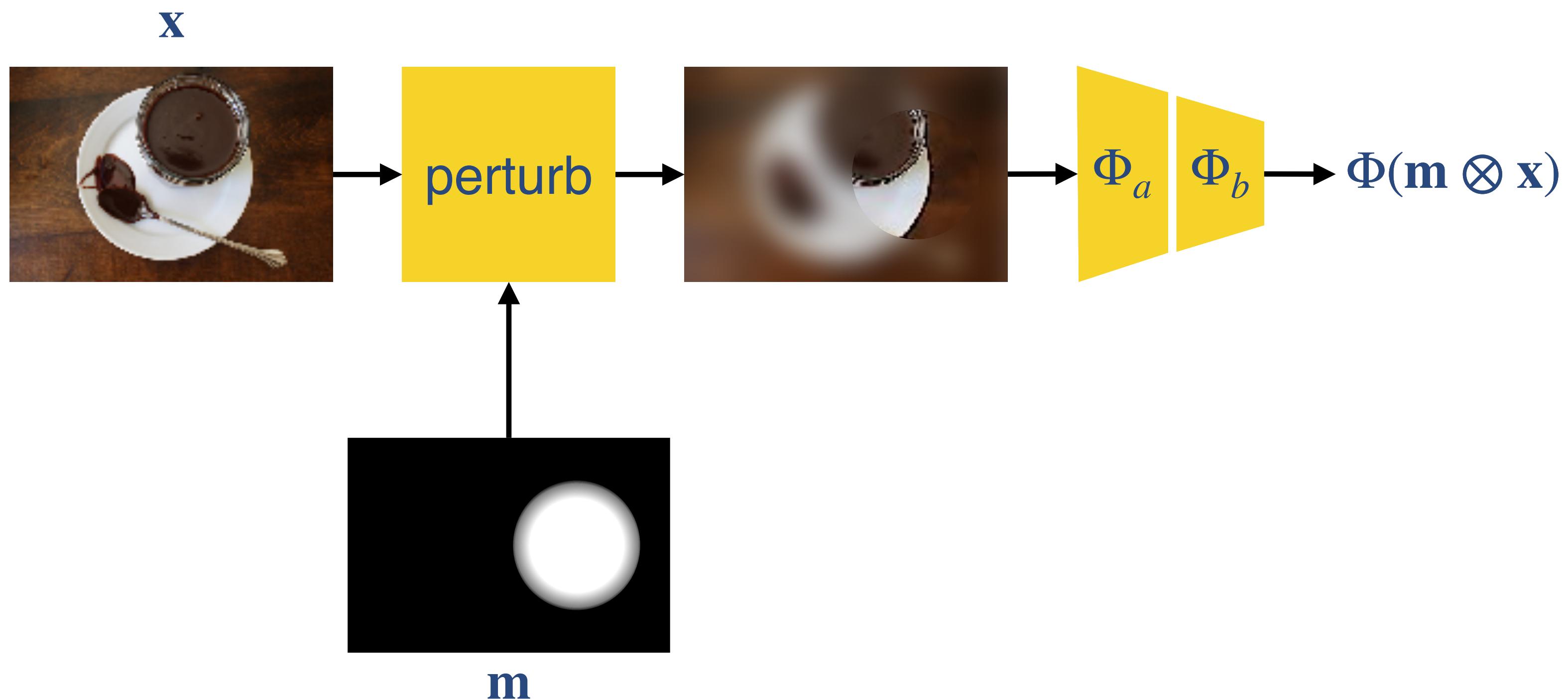


# Attributing channels at intermediate layers

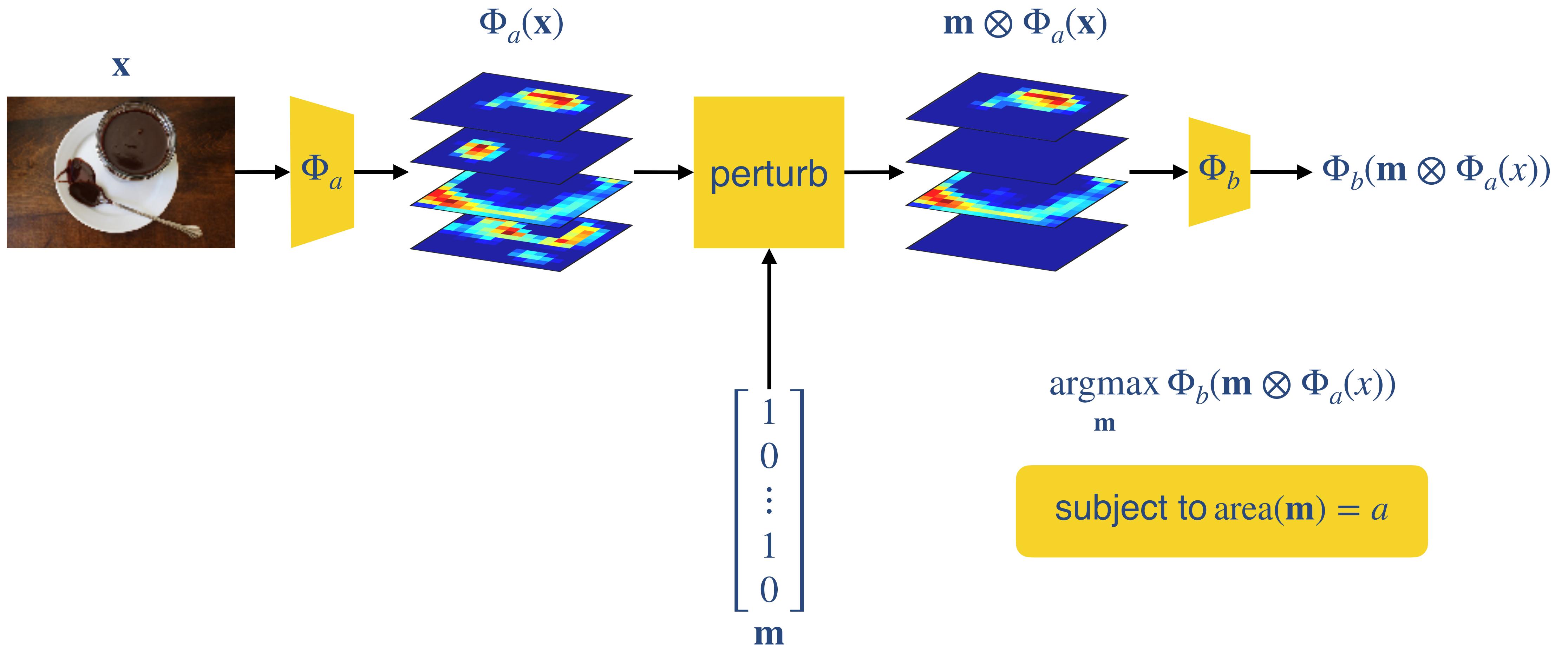
# Spatial Attribution



# Channel Attribution



# Channel Attribution



# Activation “diffing”

