

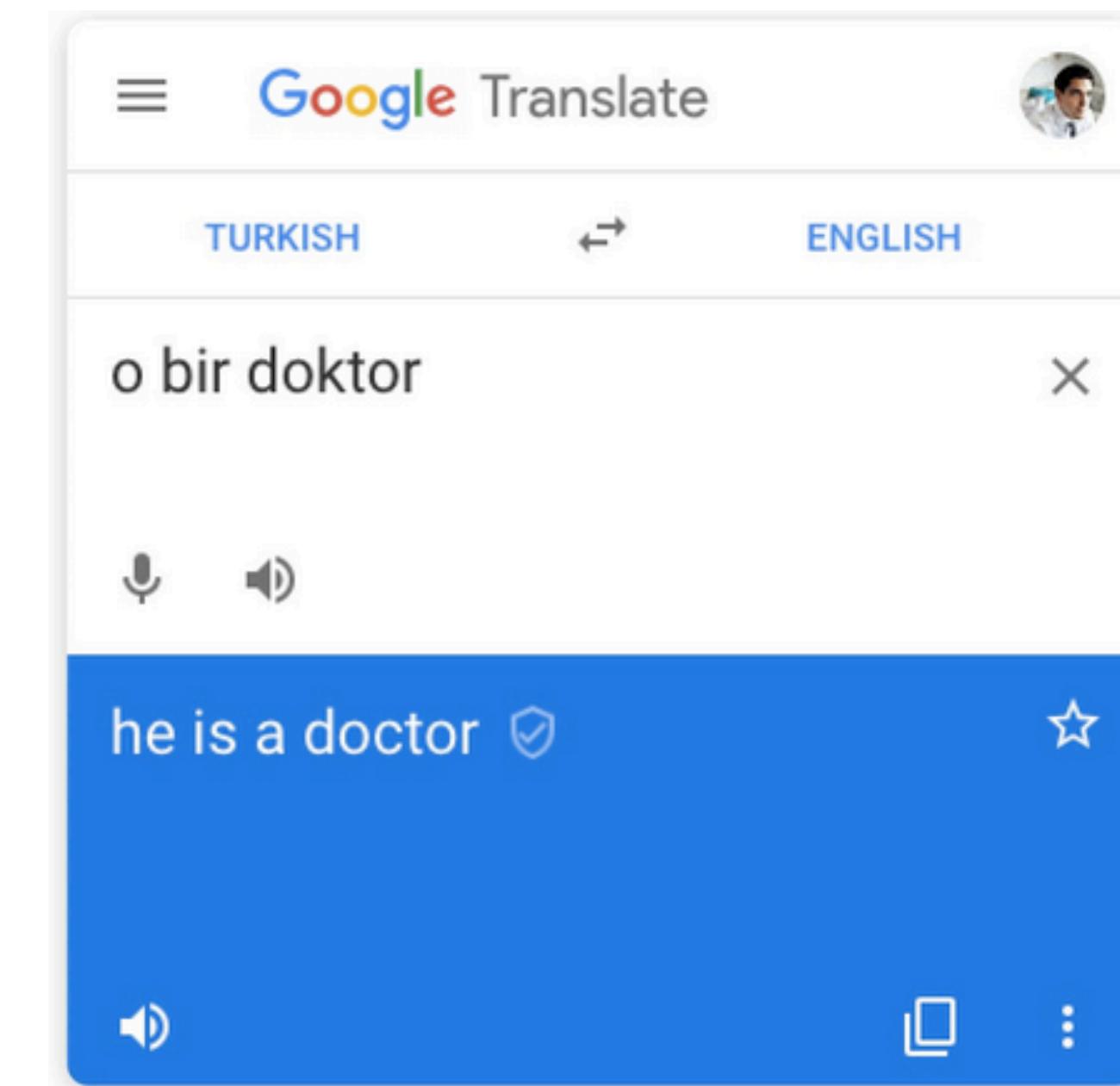
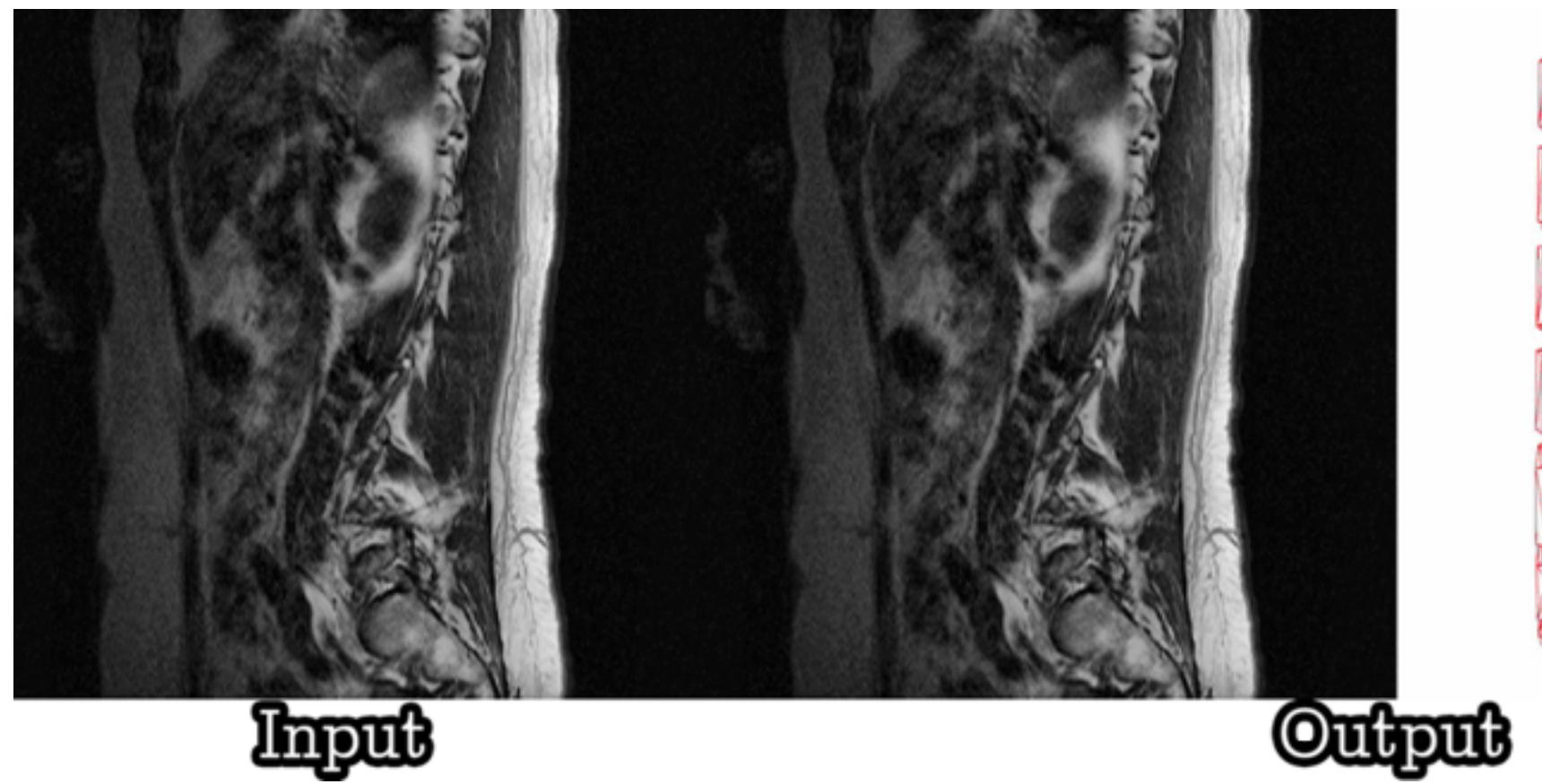
Understanding Deep Neural Networks

Ruth Fong

CVPR 2020 Tutorial on Interpretable Machine Learning for Computer Vision

Saturday, 13 June 2020

Applications of Deep Learning



Interpretability tools are crucial for high-impact, high-risk applications of deep learning.

[Jamaludin et al., 2017; <https://ai.googleblog.com/2018/12/providing-gender-specific-translations.html>] 2

A Brief Primer on Deep Learning

Supervised Learning

(



x

,

"sheepdog"

y

Supervised Learning

(



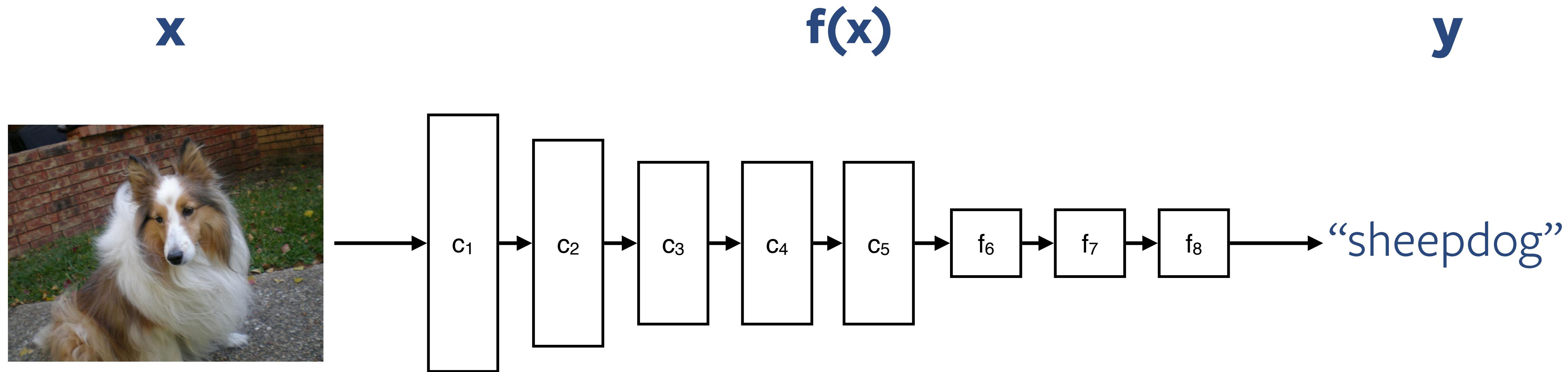
x

,

y

“sheepdog”

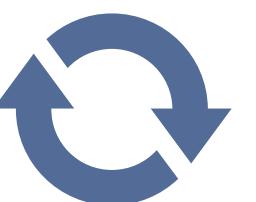
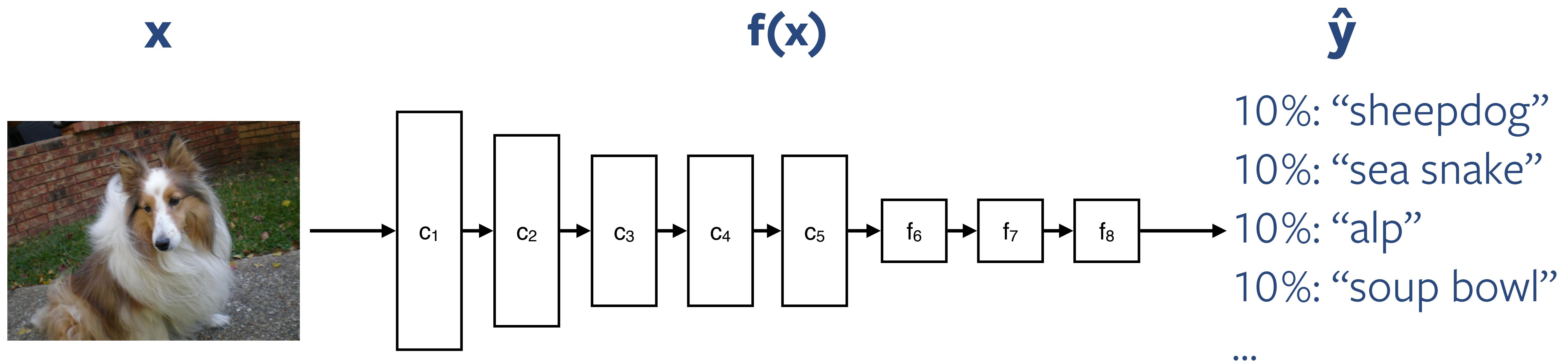
Deep Learning



Network built up of layers, with weights θ connecting one layer to the next

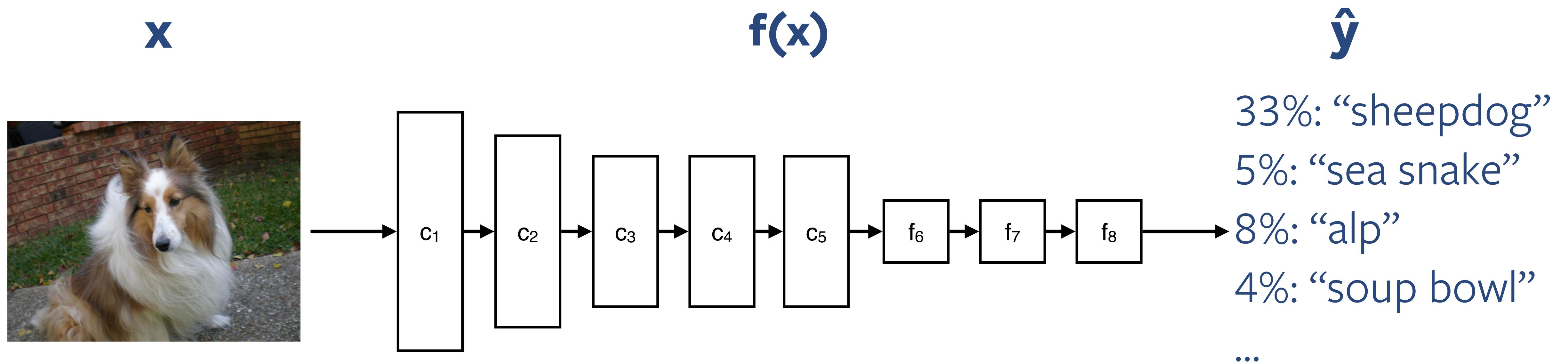
Update rule: $\theta \leftarrow \theta - \eta \frac{dL}{d\theta}$, maximizes probability of correct prediction

Deep Learning



$$\theta \leftarrow \theta - \eta \frac{dL}{d\theta}$$

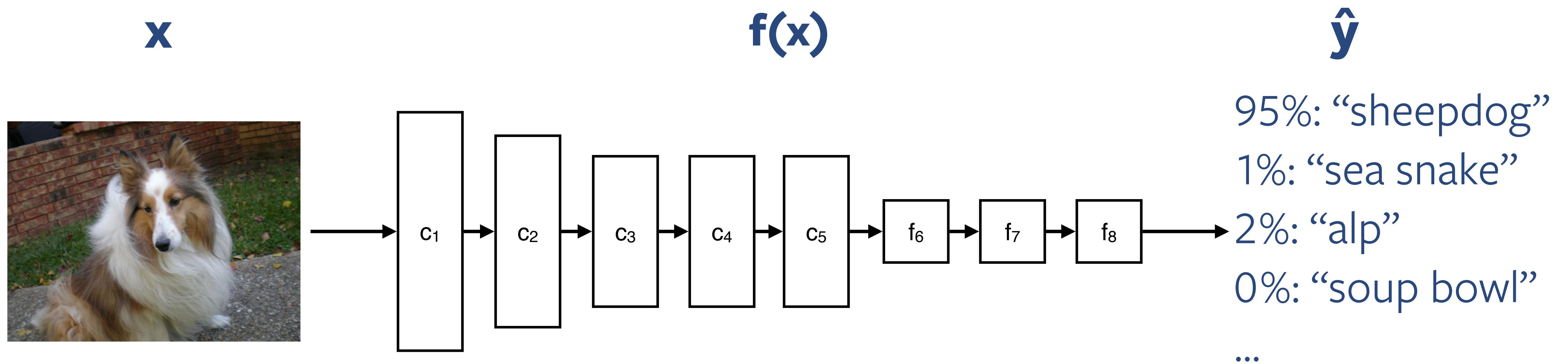
Deep Learning



↻

$$\theta \leftarrow \theta - \eta \frac{dL}{d\theta}$$

Deep Learning



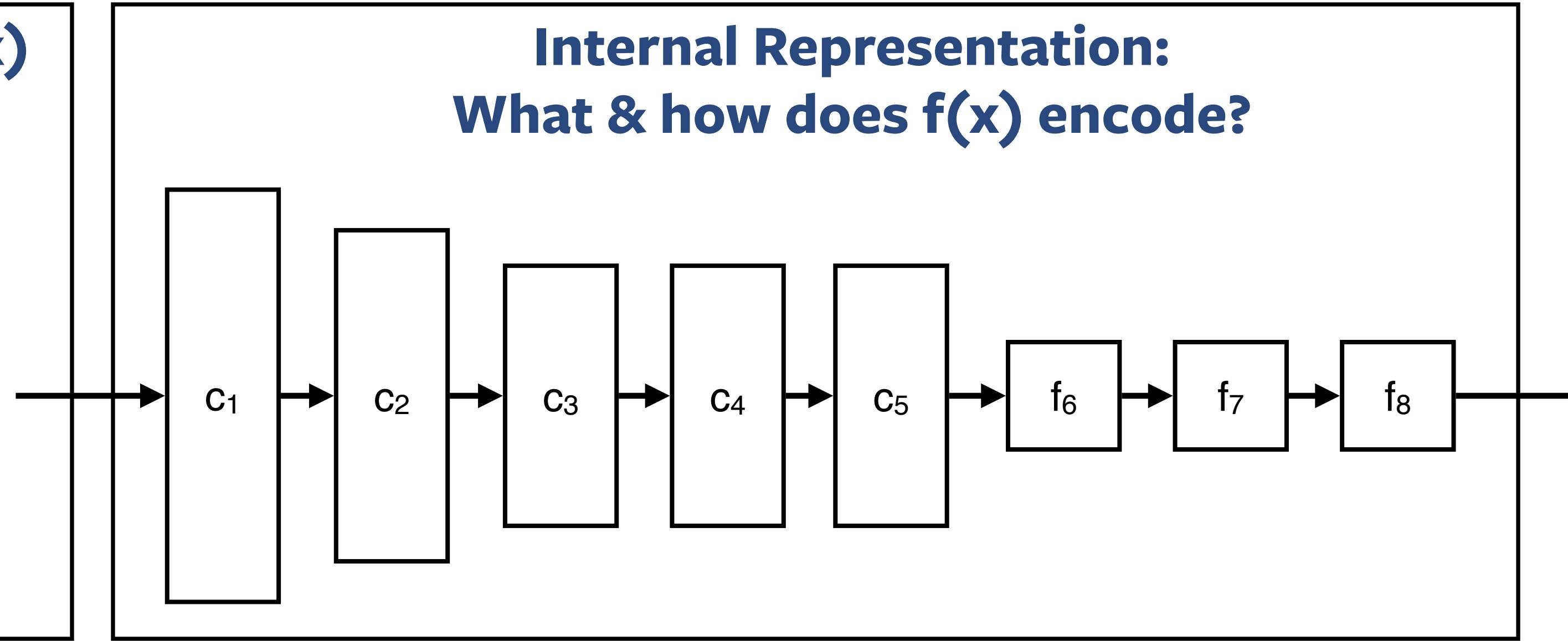
$$\theta \leftarrow \theta - \eta \frac{dL}{d\theta}$$

Research Themes

Inputs: What is $f(x)$ looking at?



**Internal Representation:
What & how does $f(x)$ encode?**



95%: “sheepdog”
1%: “sea snake”
2%: “alp”
0%: “soup bowl”
...

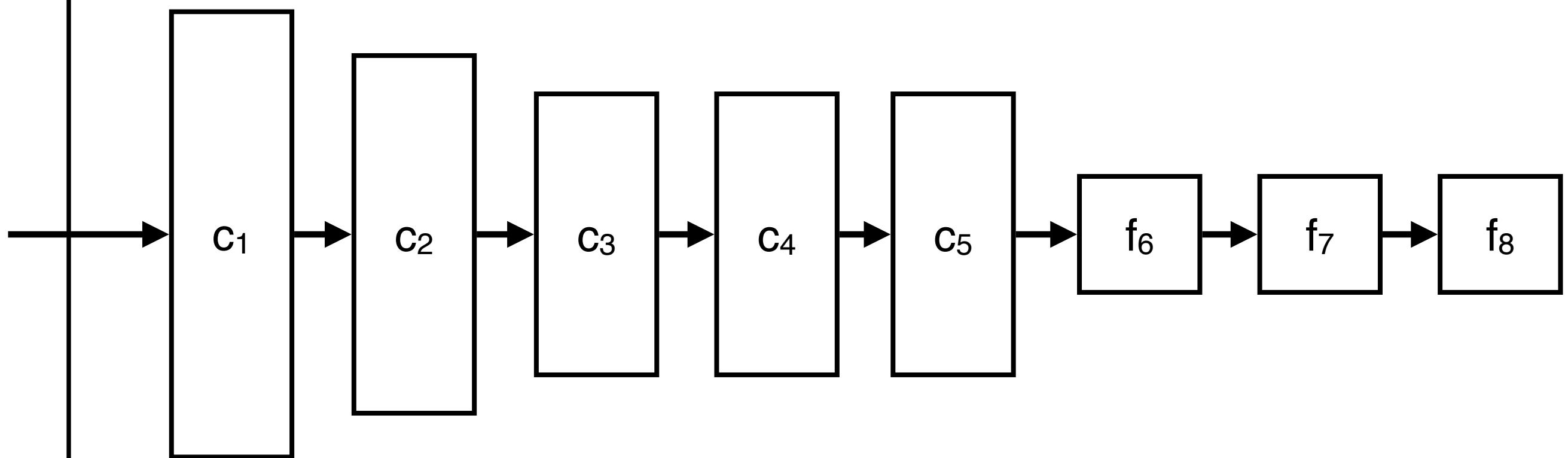
Training Procedure: How can we improve $f(x)$?



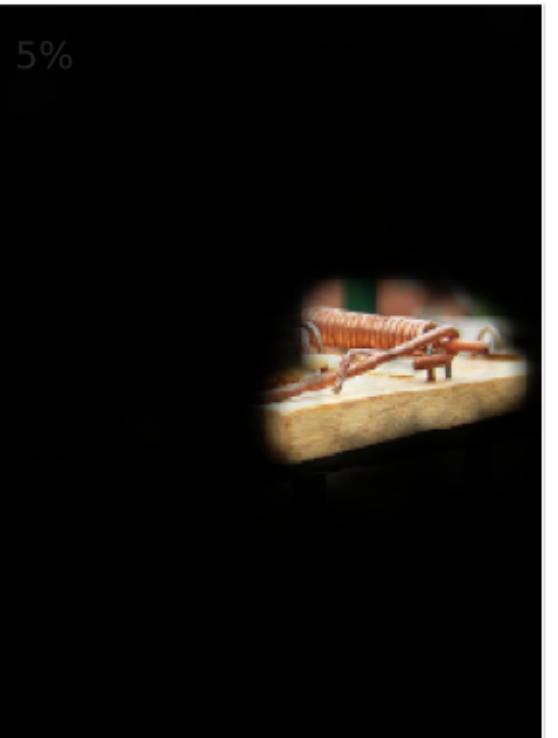
$$\theta \leftarrow \theta - \eta \frac{dL}{d\theta}$$

Research Themes

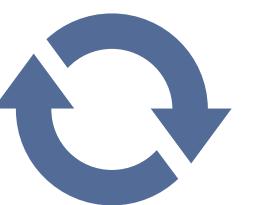
Inputs: What is $f(x)$ looking at?



Fong & Vedaldi, ICCV 2017



Fong et al., ICCV 2019

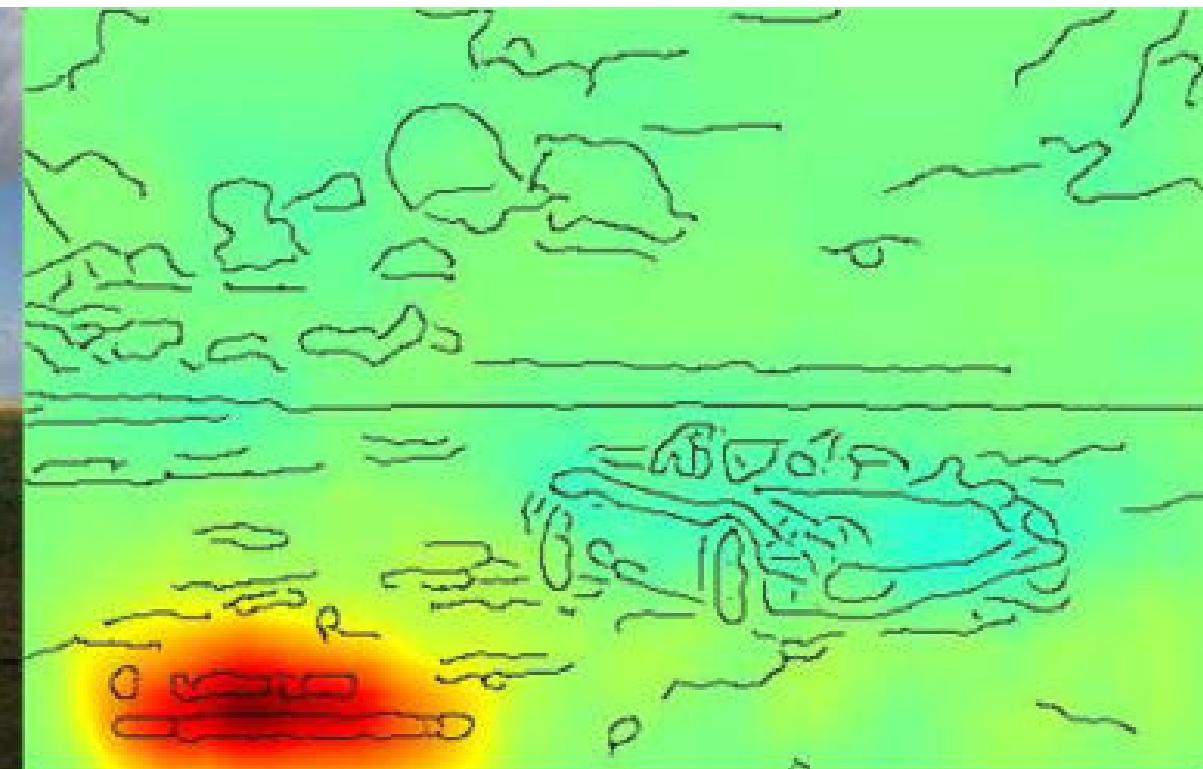
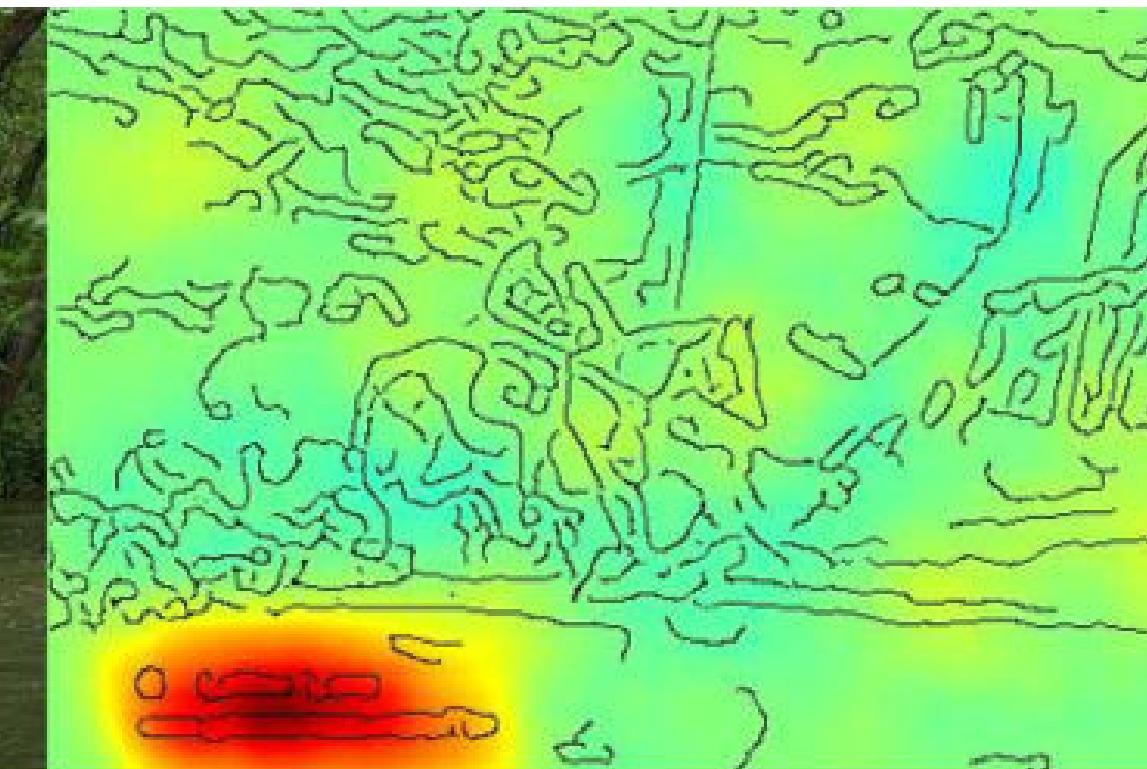


$$\theta \leftarrow \theta - \eta \frac{dL}{d\theta}$$

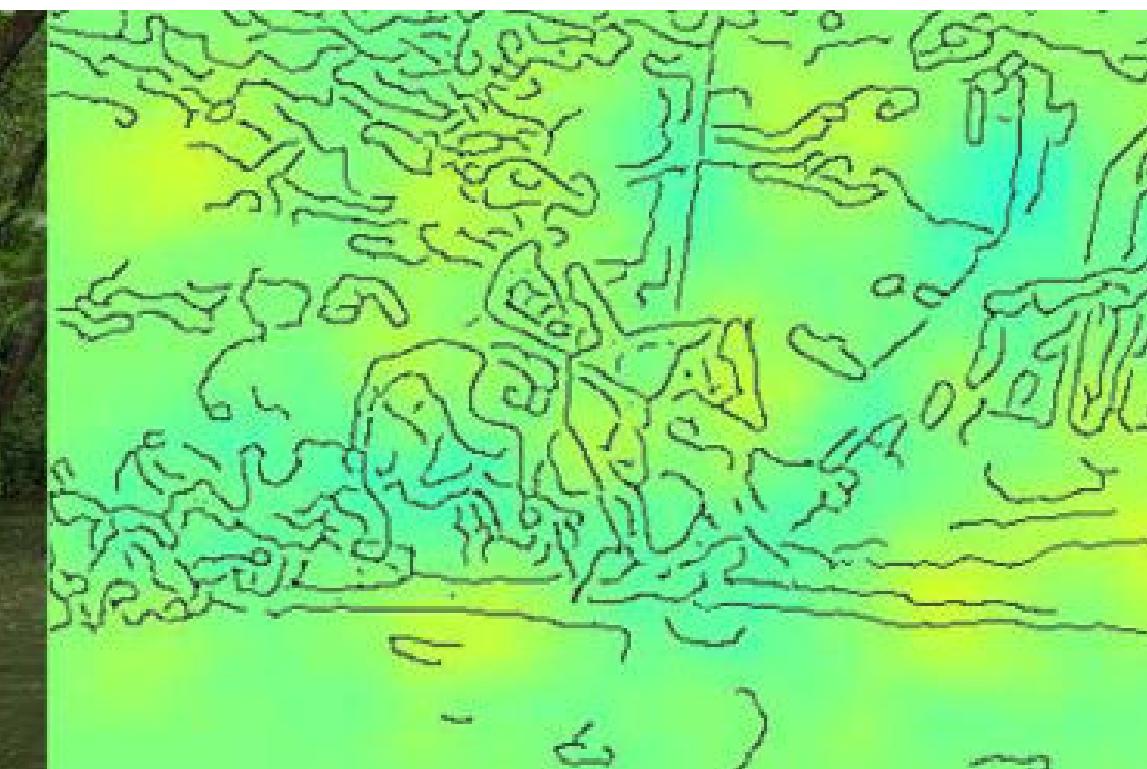


“Math whiz” Clever Hans horse

“horse”

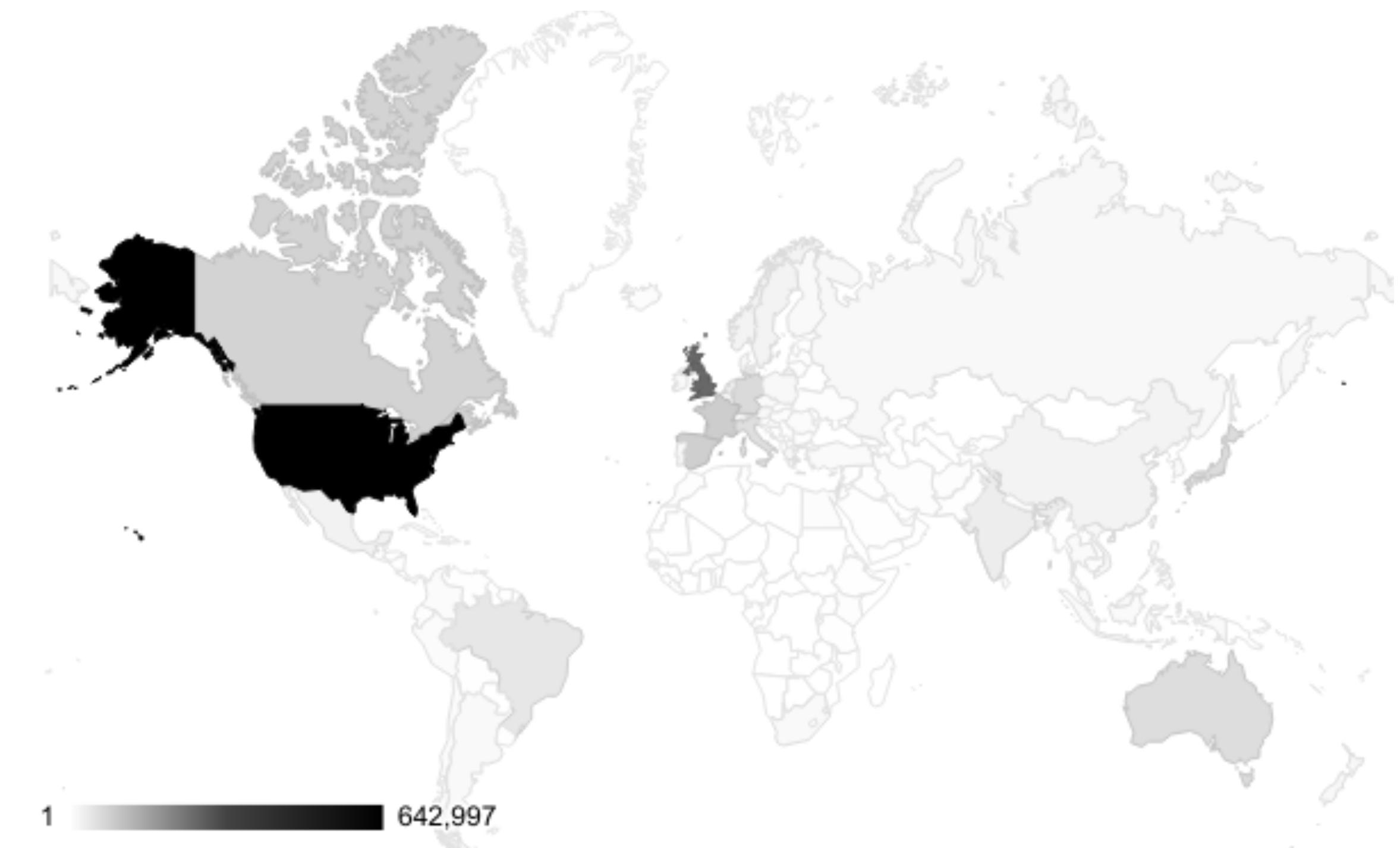
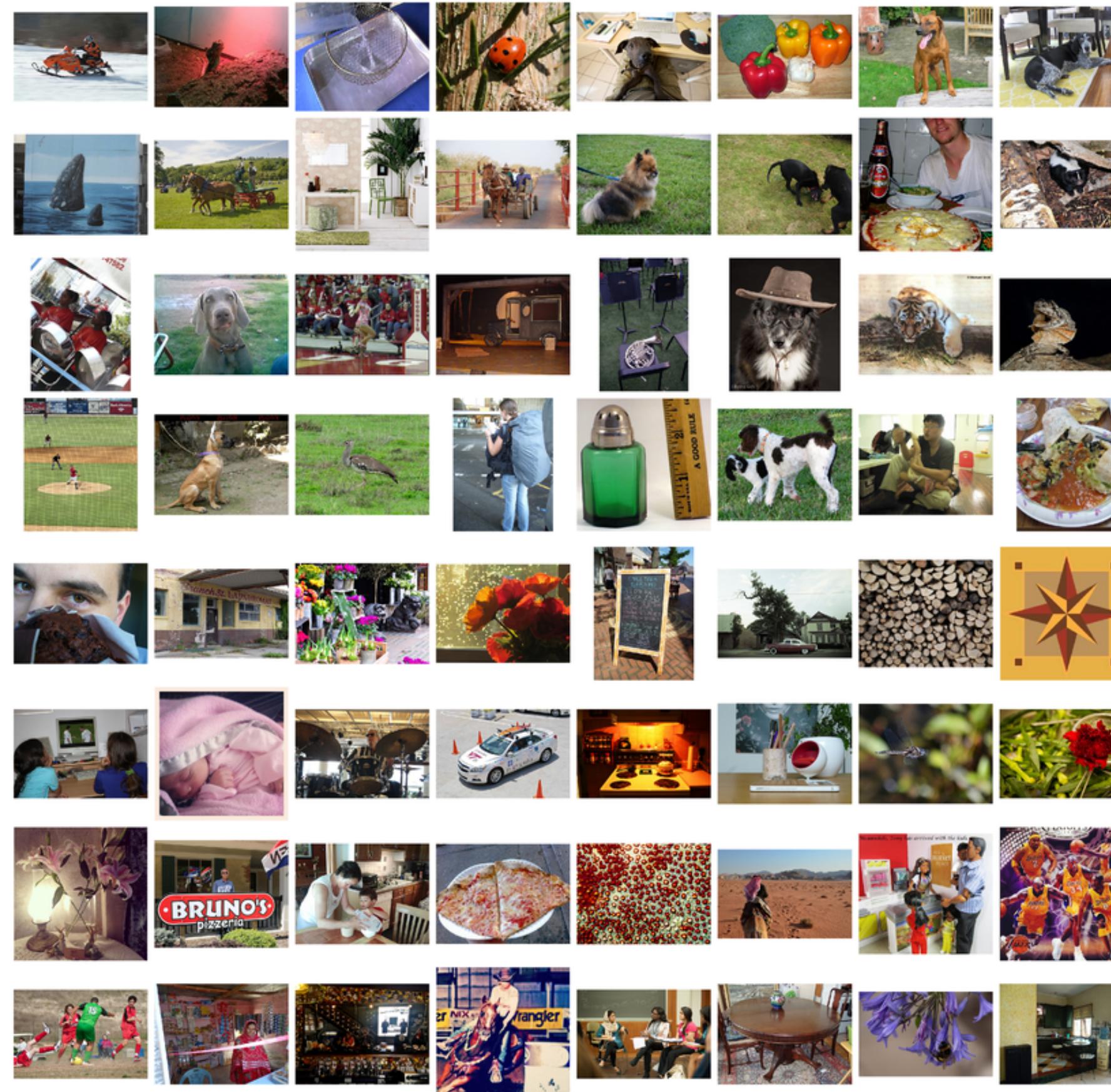


“not horse”



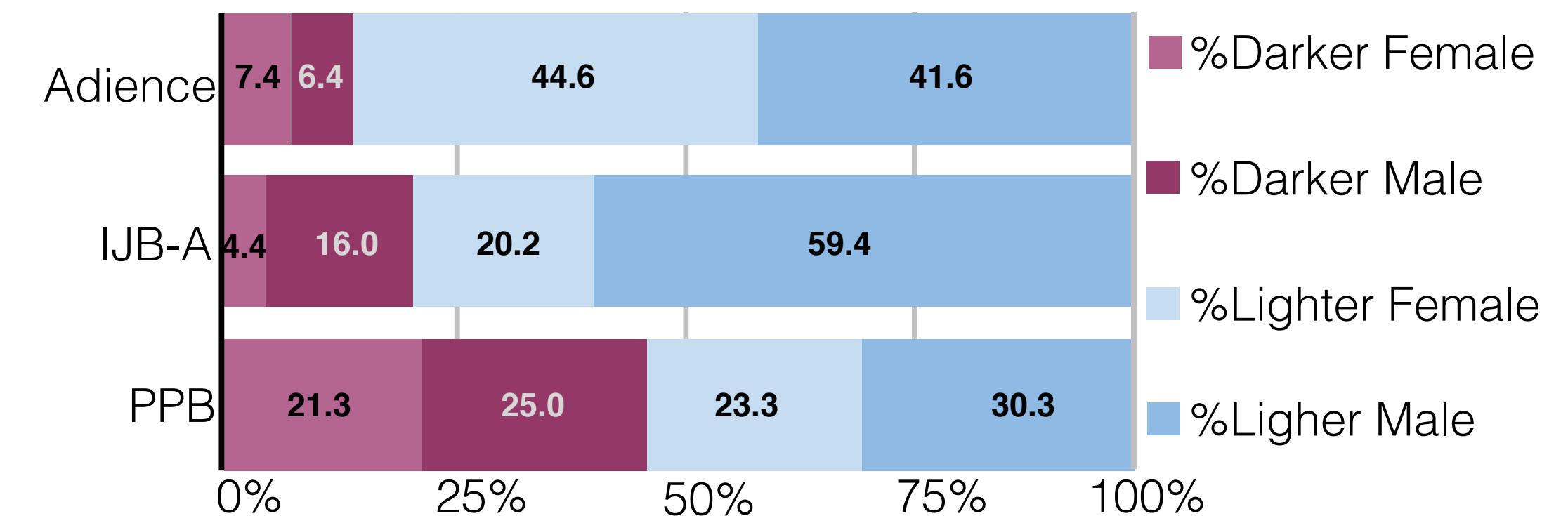
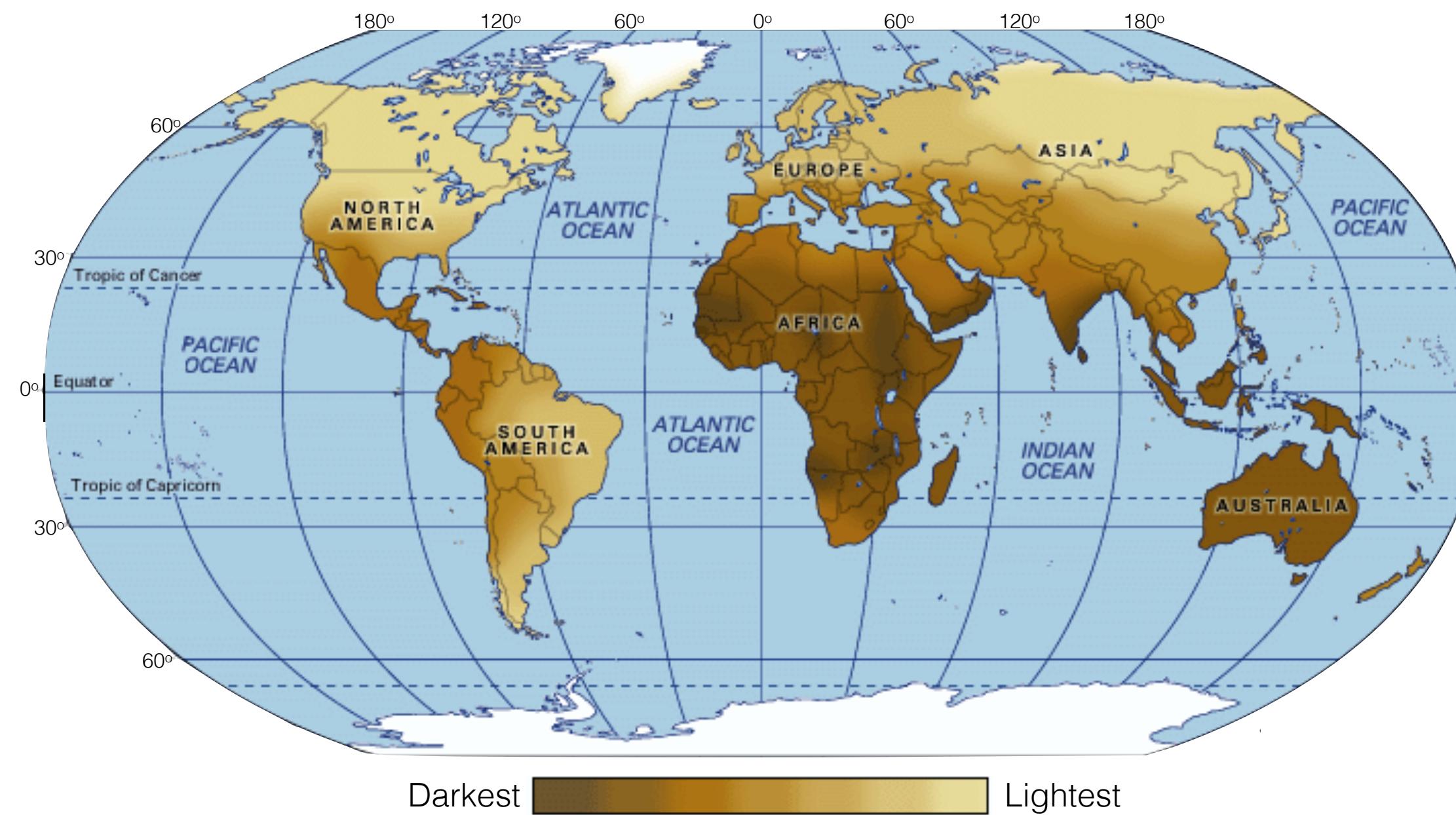
PASCAL object detection dataset

[Everingham et al., IJCV 2010; Lapuschkin et al., Nat. Commun. 2019]



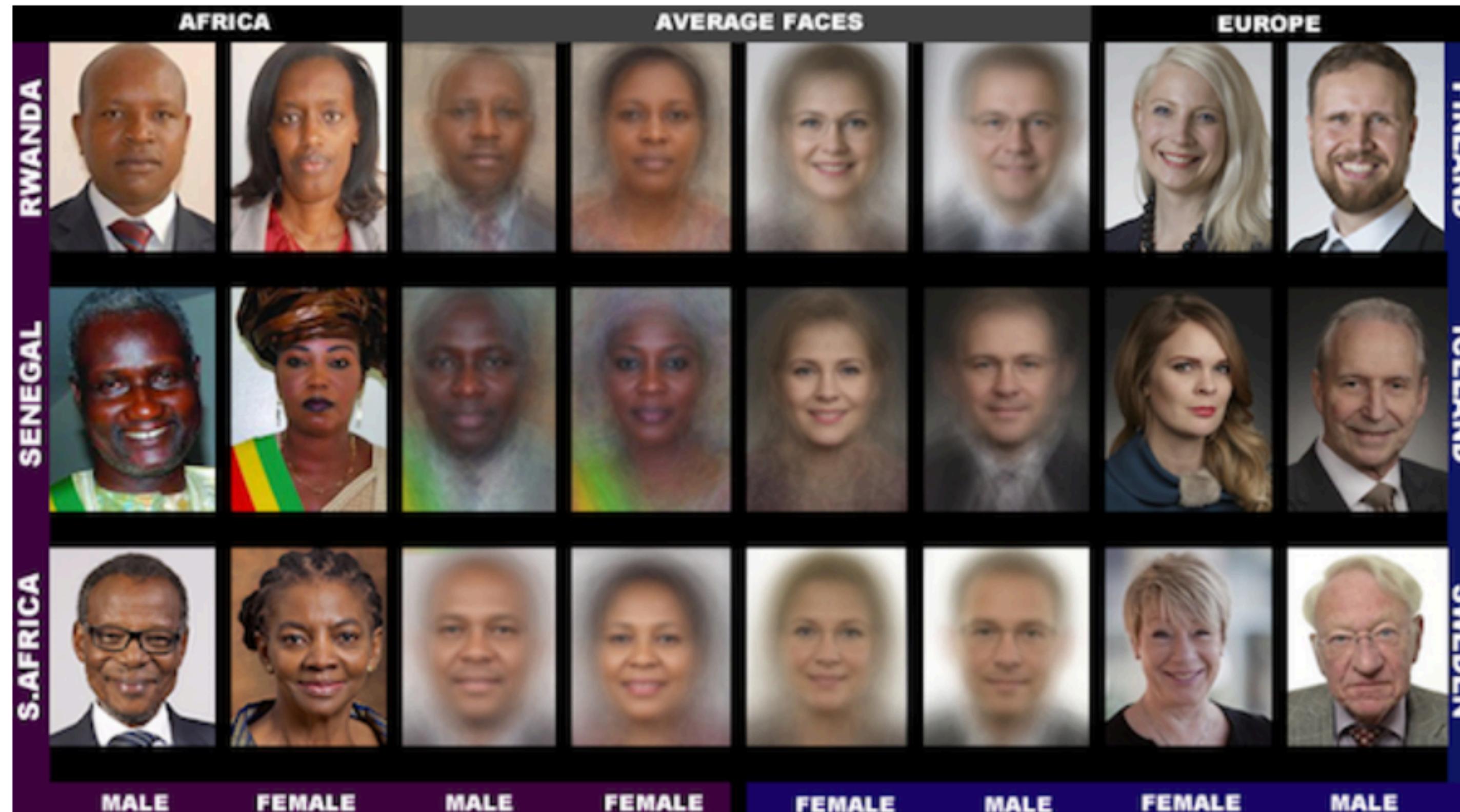
ImageNet object recognition dataset

[Russakovsky et al., IJCV 2015; Shankar et al., NeurIPS Workshop 2017]



Face datasets

[Buolamwini & Gebru, JMLR 2018; globe image from Encyclopedia Britannica]



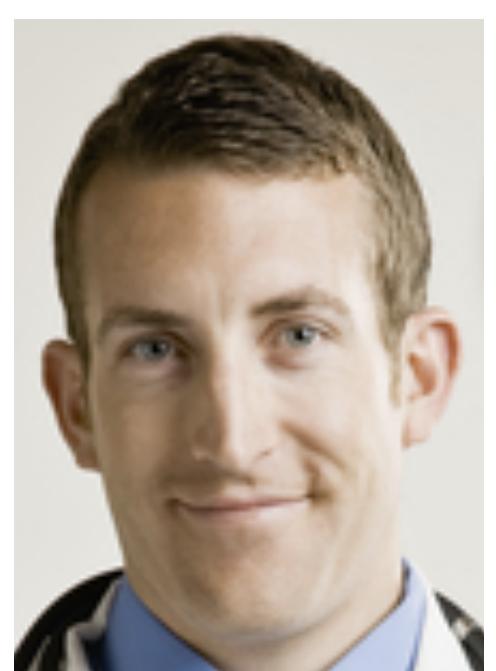
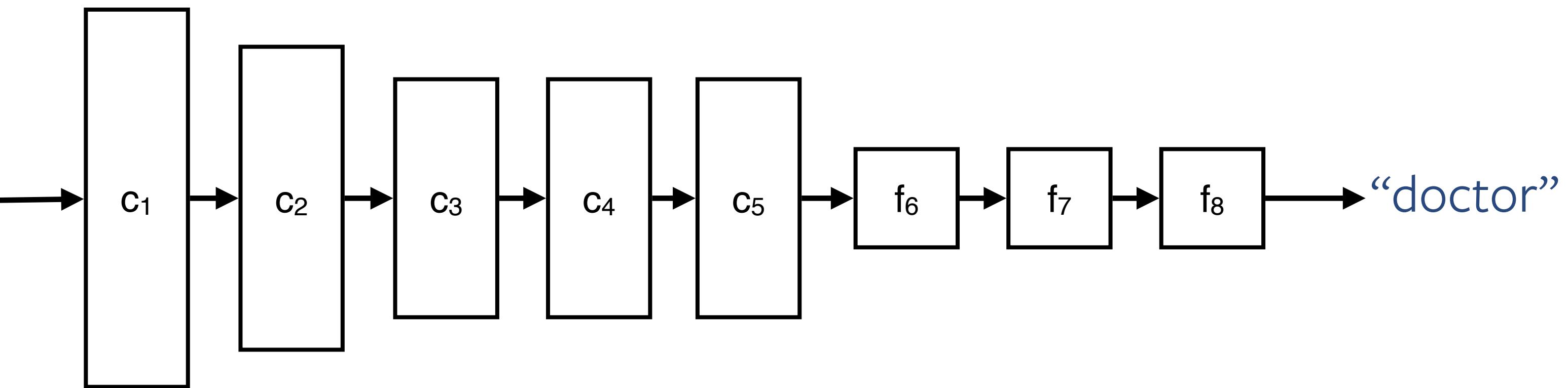
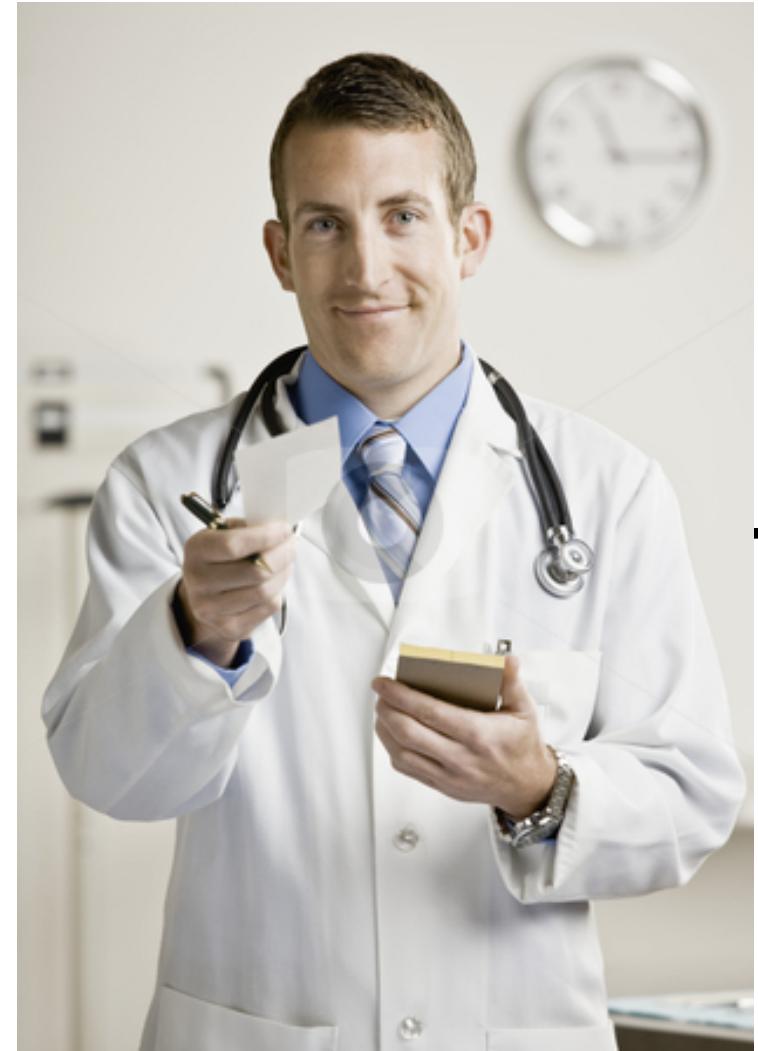
Face datasets

[Buolamwini & Gebru, JMLR 2018; globe image from Encyclopedia Britannica]

Attribution

Identify input features responsible for model decision

?



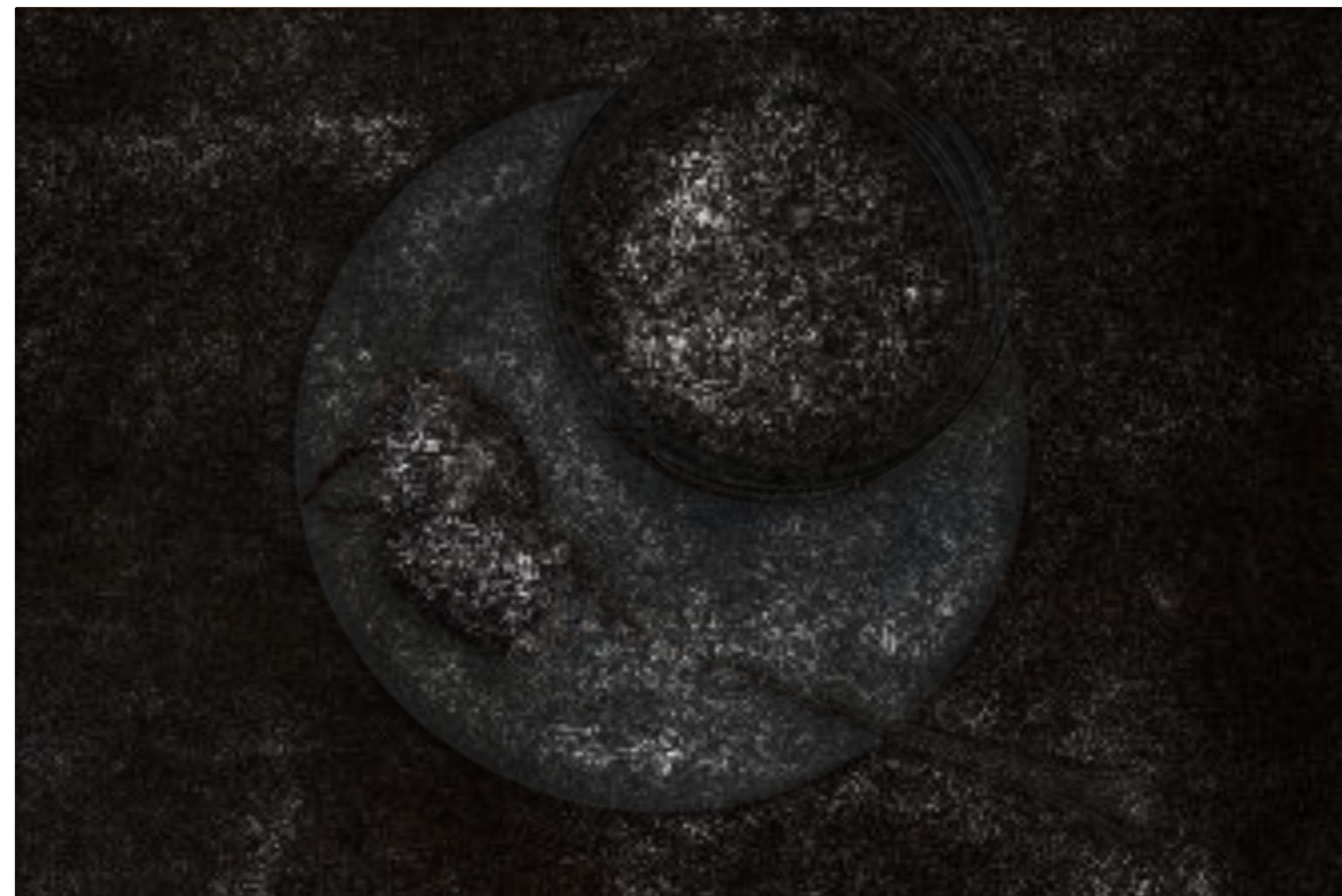
Prior Work: Propagation-based methods

Combine network activations and gradients

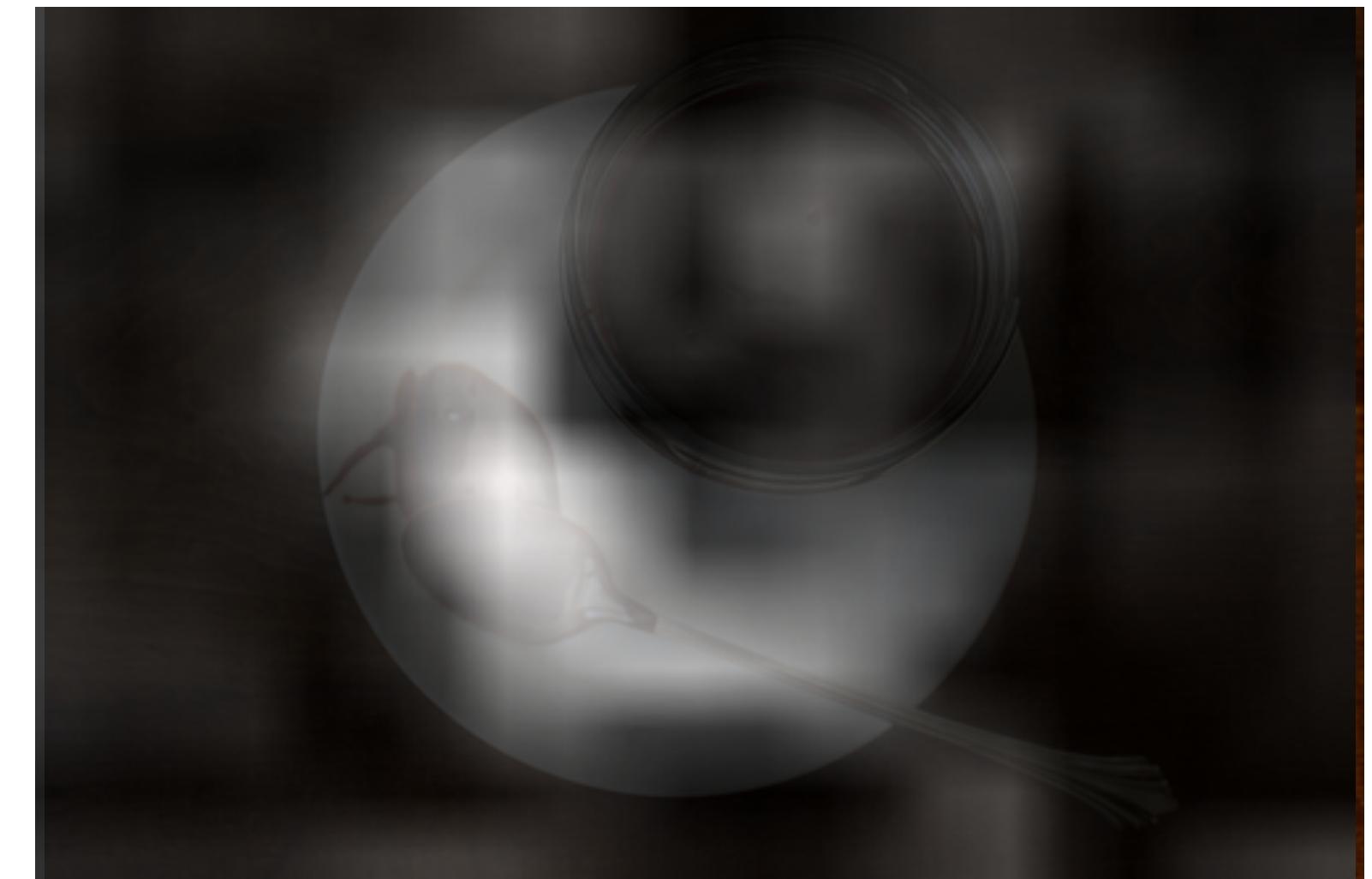
Input



Gradient



Grad-CAM



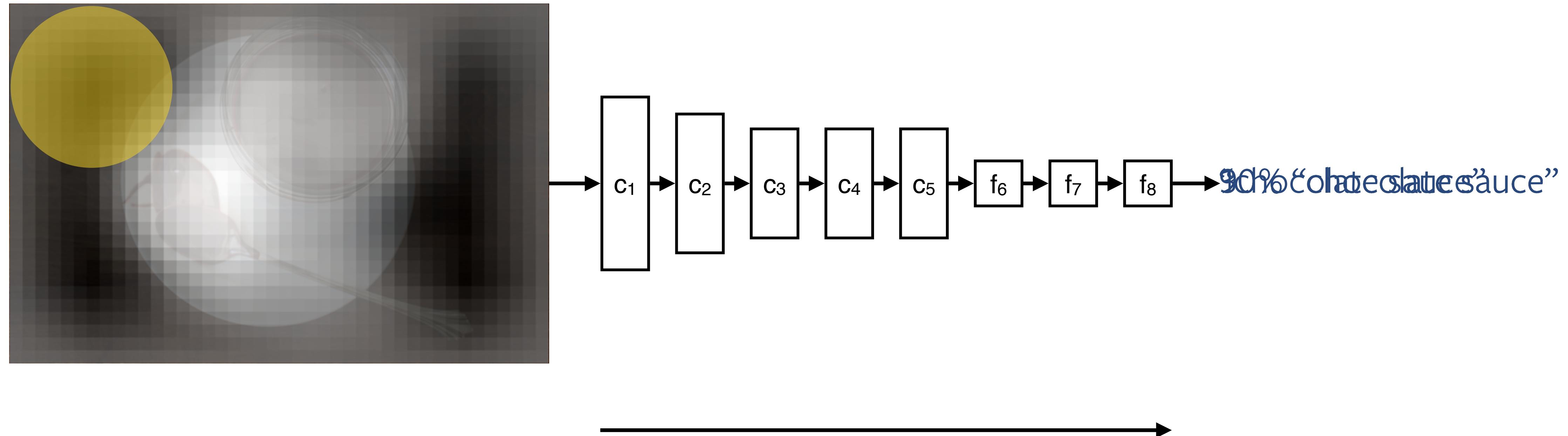
Fast, but difficult to interpret

[Simonyan et al., ICLR Workshop 2014; Selvaraju et al., ICCV 2017]
[Mahendran and Vedaldi, ECCV 2016; Adebayo et al., NeurIPS 2018]

Prior Work: Perturbation Approaches

Change the input and observe the effect on the output

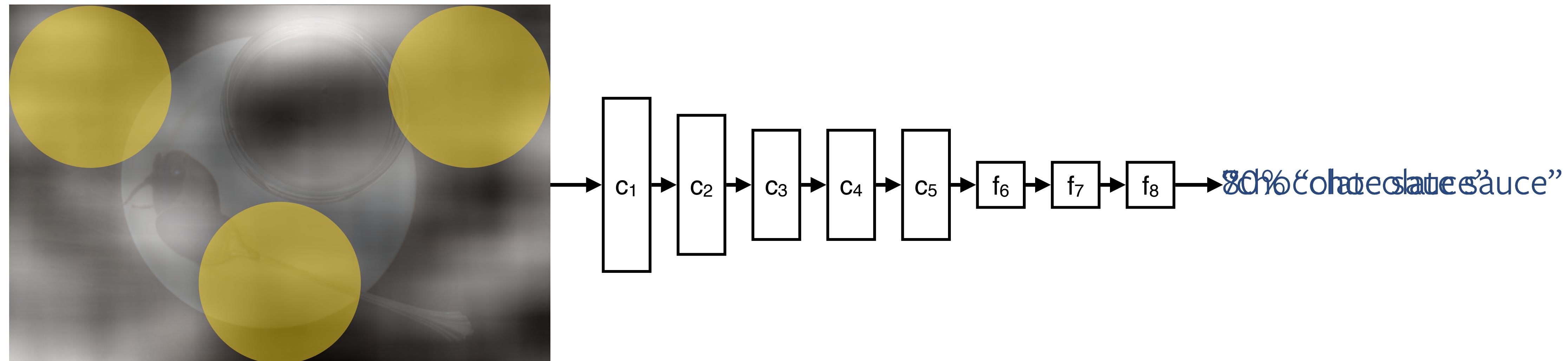
Occlusion



Prior Work: Perturbation Approaches

Change the input and observe the effect on the output

RISE



Clear meaning, but can only test a small range of occlusions →

Desired Approach

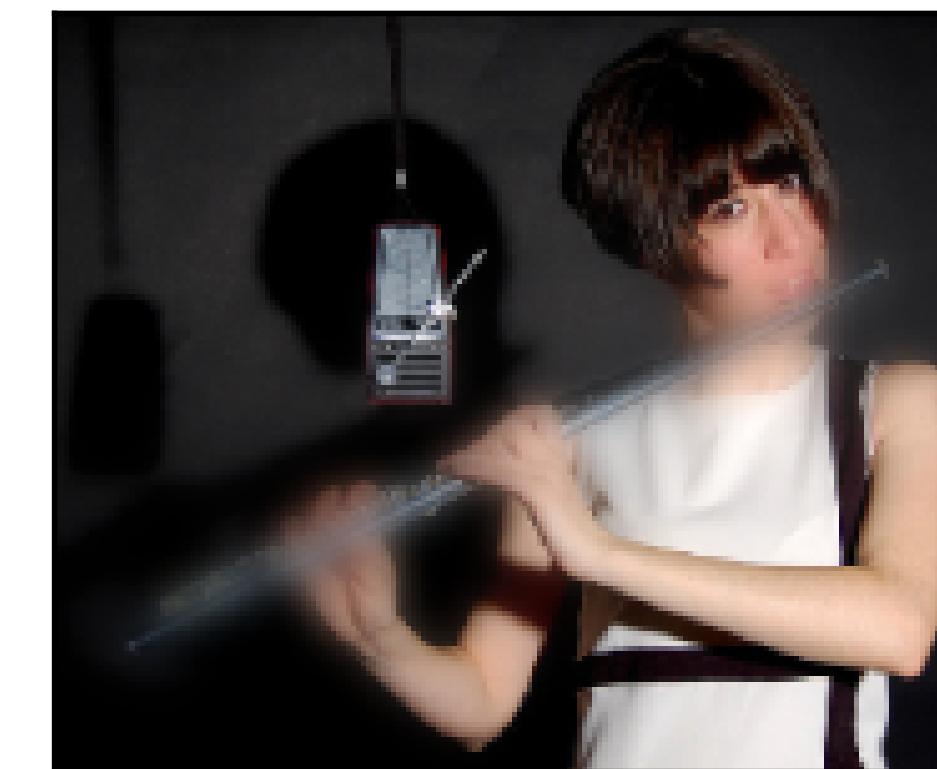


Our Approach: Meaningful Perturbations

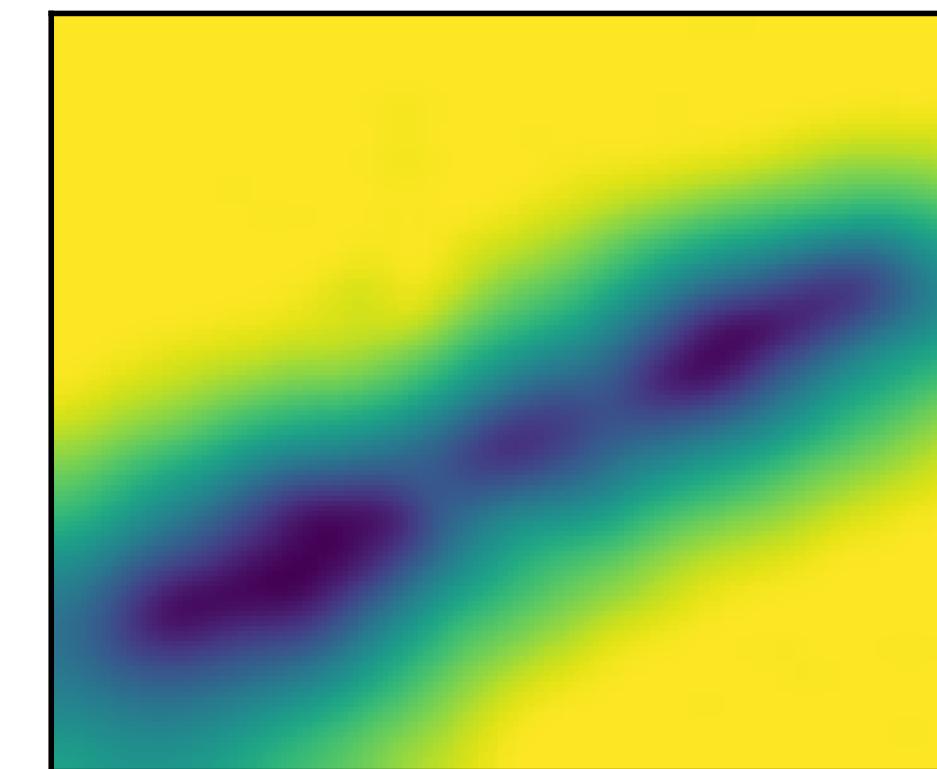
flute: 0.9973



flute: 0.0007



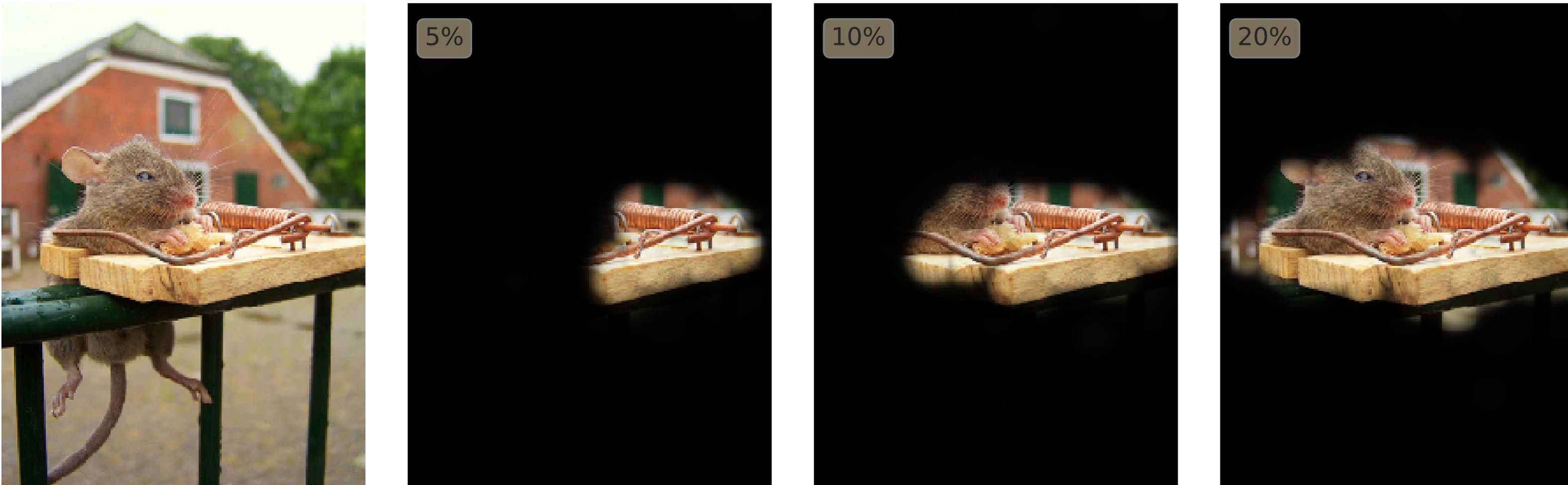
Learned Mask



Learn a **minimal** mask \mathbf{m} to perturb input \mathbf{x} that
maximally affects the network's output

Our method considers a wide range of occlusion sizes and shapes.

Our Approach: Extremal Perturbations



Learn a **fixed-sized** mask **m** to perturb input **x** that
maximally **preserves** the network's output

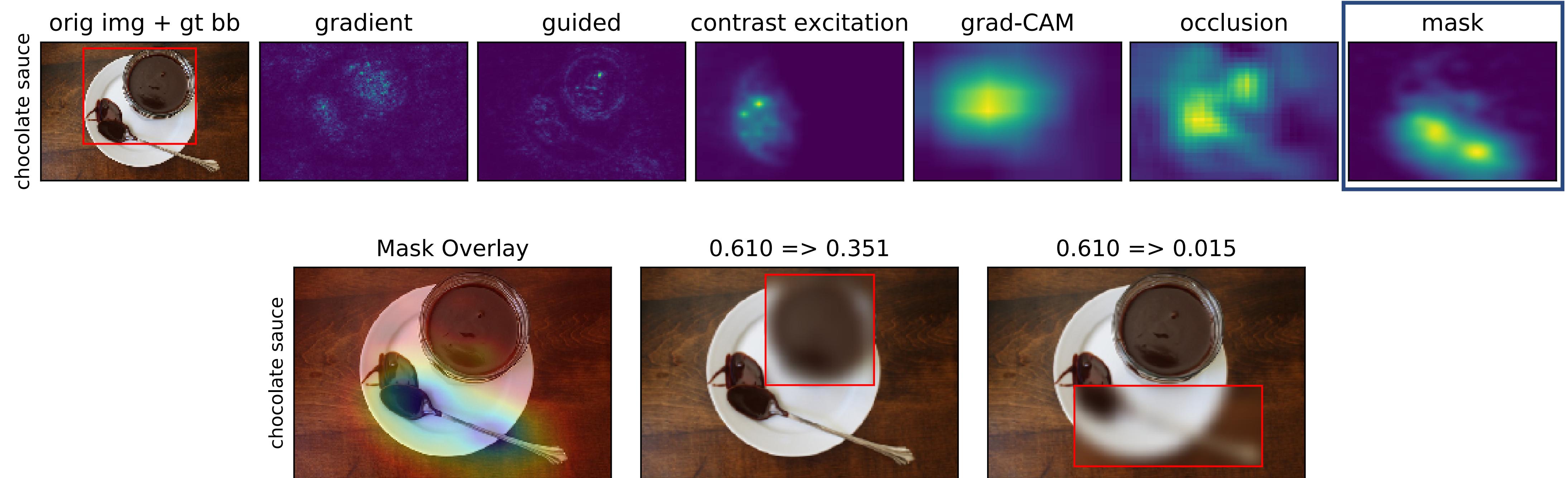
[Fong et al., ICCV 2019]

23

Concurrent work: [Kapishnikov et al., ICCV 2019]

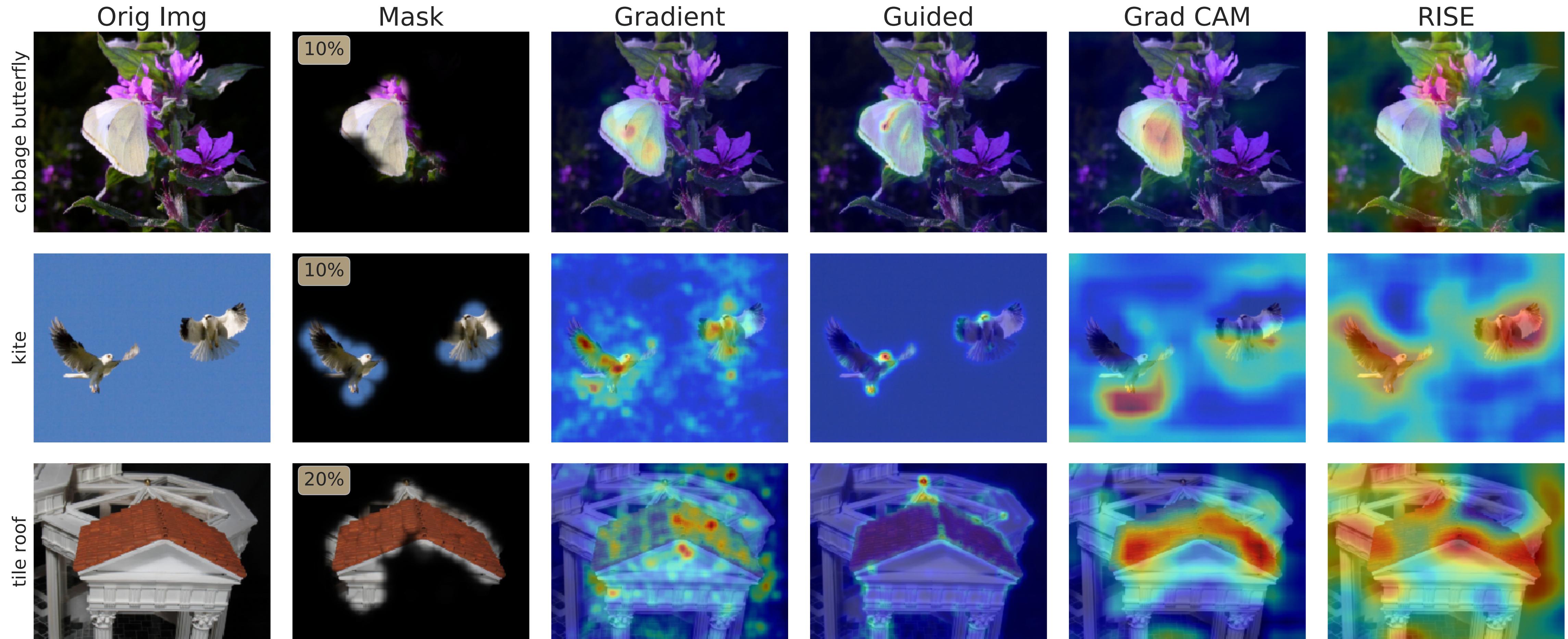
Results

Interpretability

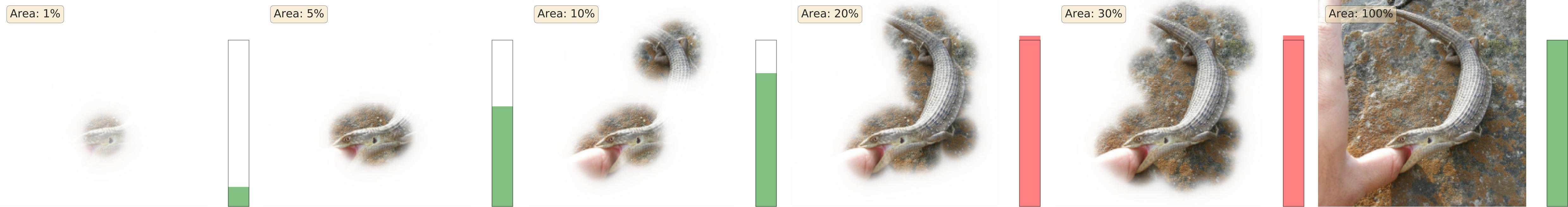
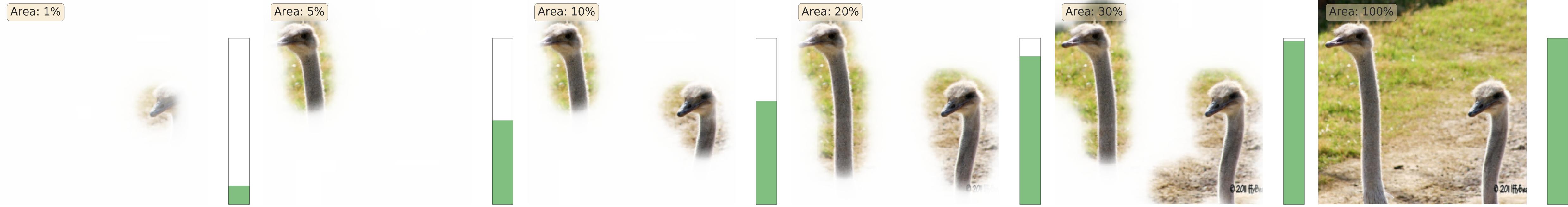


An explanation should be **falsifiable**.

Comparison



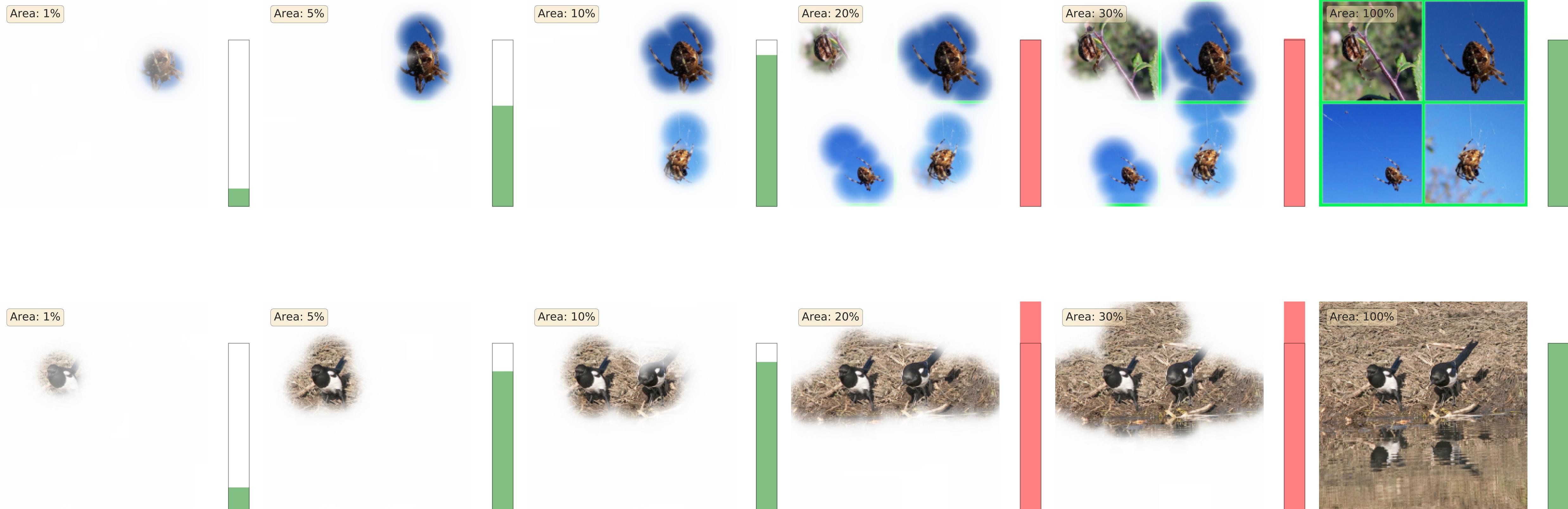
Foreground evidence is usually sufficient



Large objects are recognized by their details



Multiple objects contribute cumulatively



Suppressing the background may overdrive the network

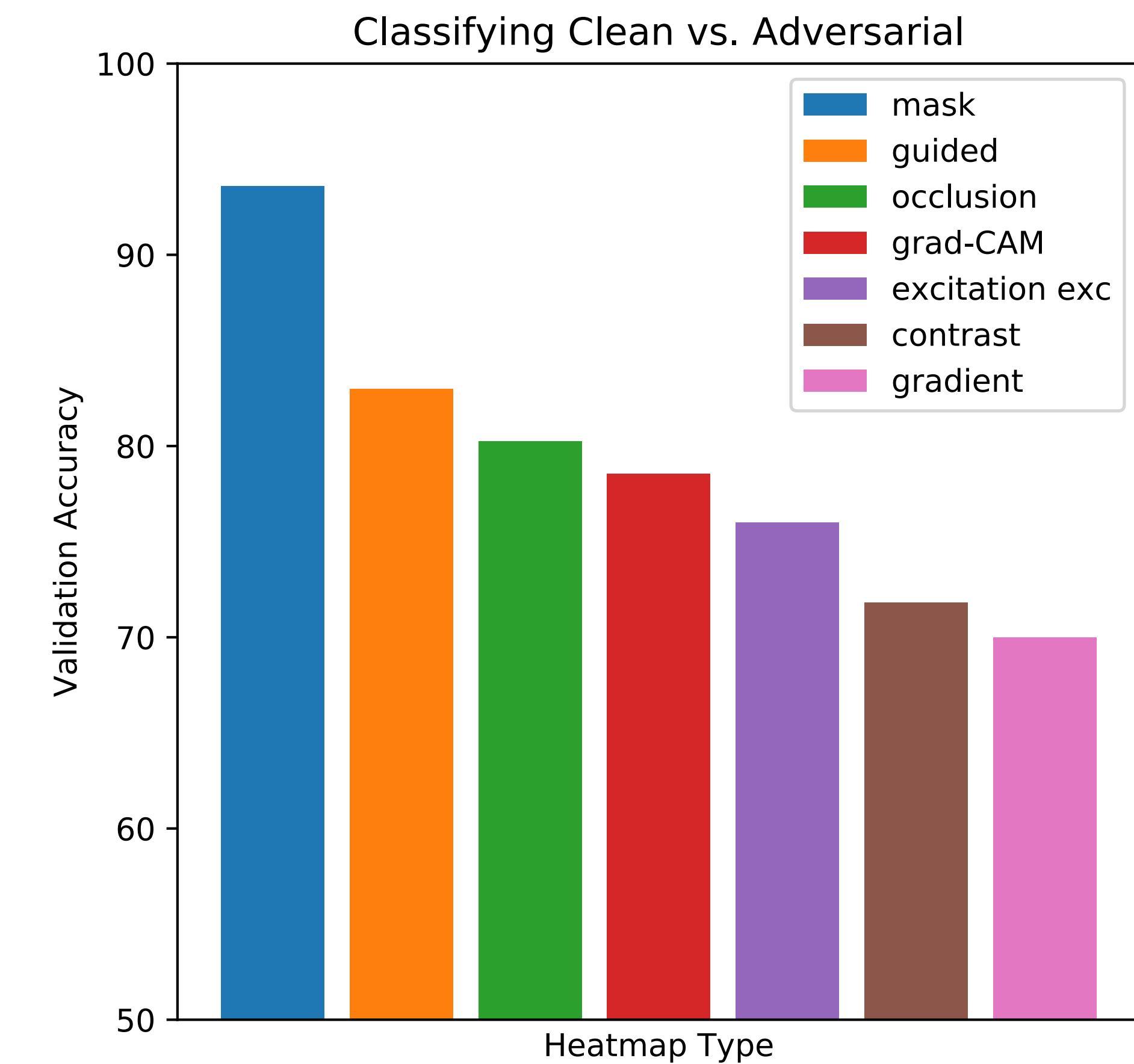
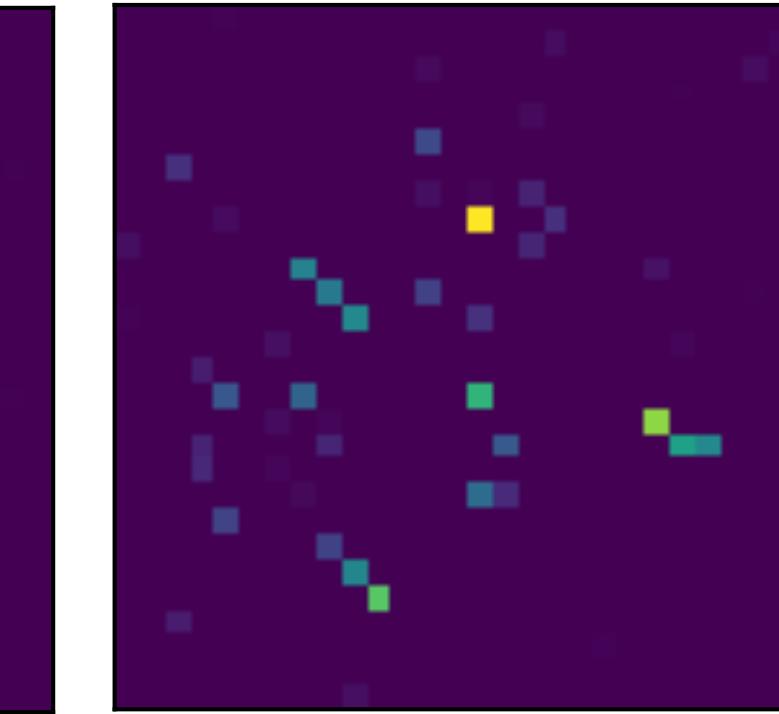


Adversarial Defense

Mask on
Clean Image



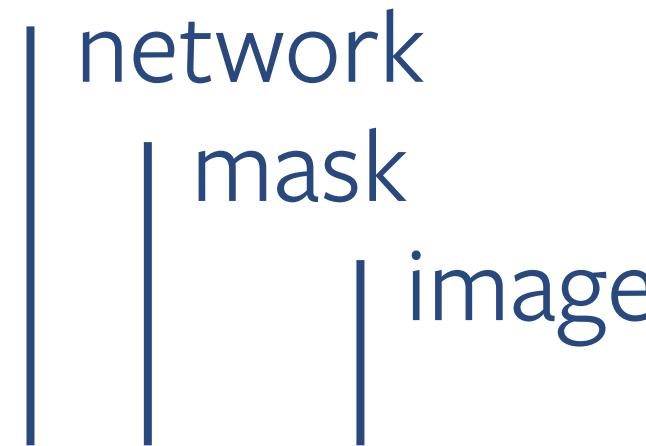
Mask on
Adv. Image



Our method allows us to defend **any model** against adversarial attacks.

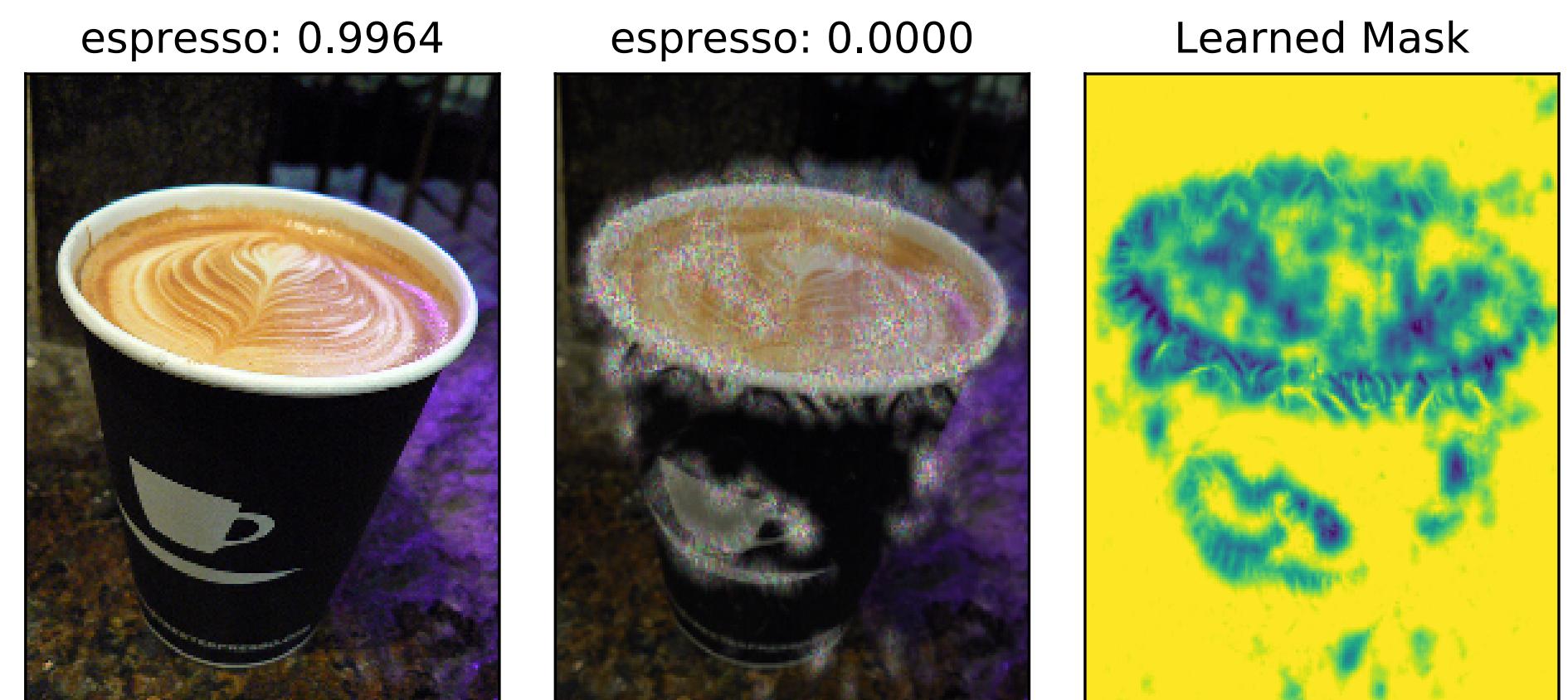
Details

Regularization to mitigate artifacts



$$v1: \mathbf{m}^*(\lambda) = \operatorname{argmin}_{\mathbf{m}} \Phi(\mathbf{m} \otimes \mathbf{x}) + \lambda \operatorname{area}(\mathbf{m})$$

$$v2: \mathbf{m}^*(\lambda_1, \lambda_2) = \operatorname{argmin}_{\mathbf{m}} \mathbb{E}_{\text{jitter}}[\Phi(M_{\text{upsample}}(\mathbf{m}) \otimes \mathbf{x})] + \lambda_1 \operatorname{area}(\mathbf{m}) + \lambda_2 \operatorname{smooth}(\mathbf{m})$$



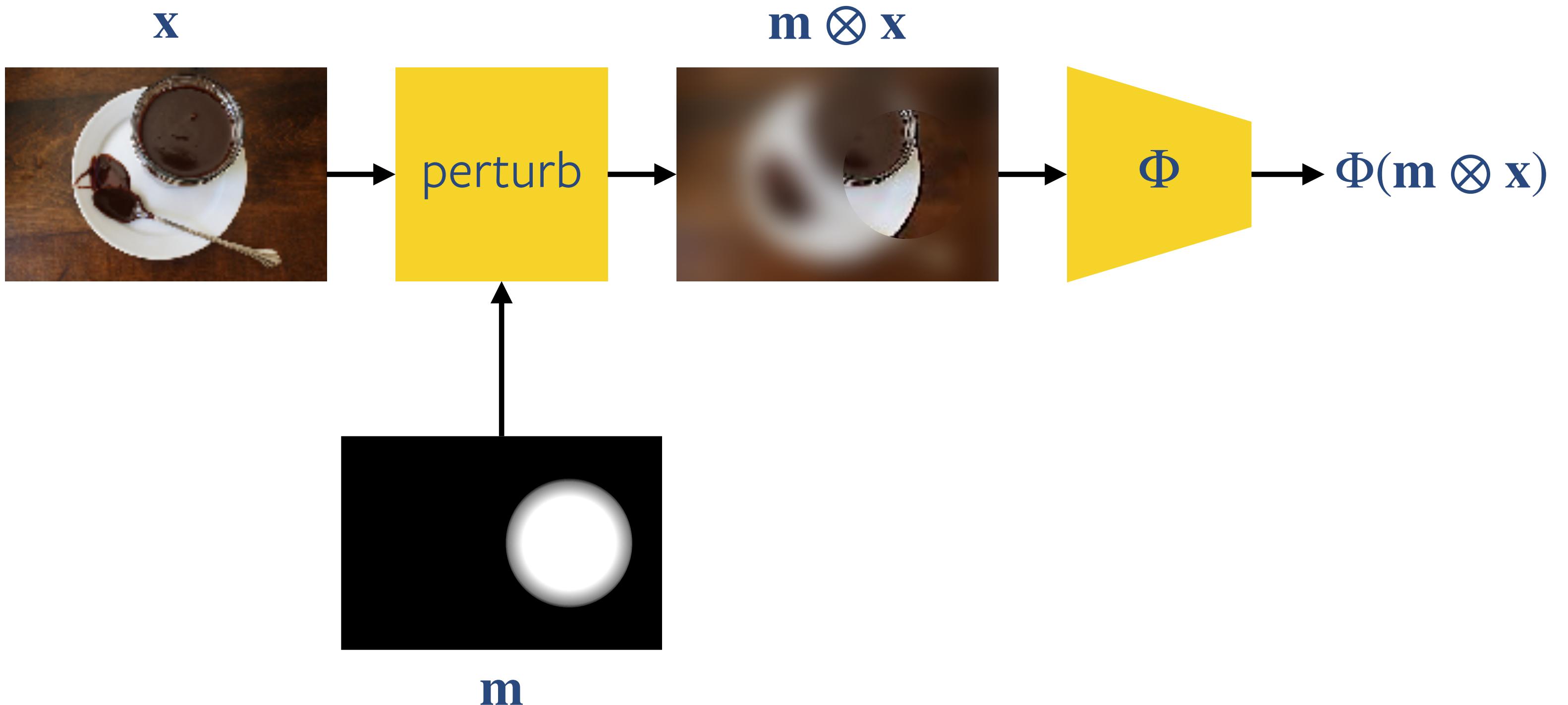
Tradeoff between **attribution objective** and **regularization**

Extremal Perturbations

A mask is optimized to maximally excite the network:

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

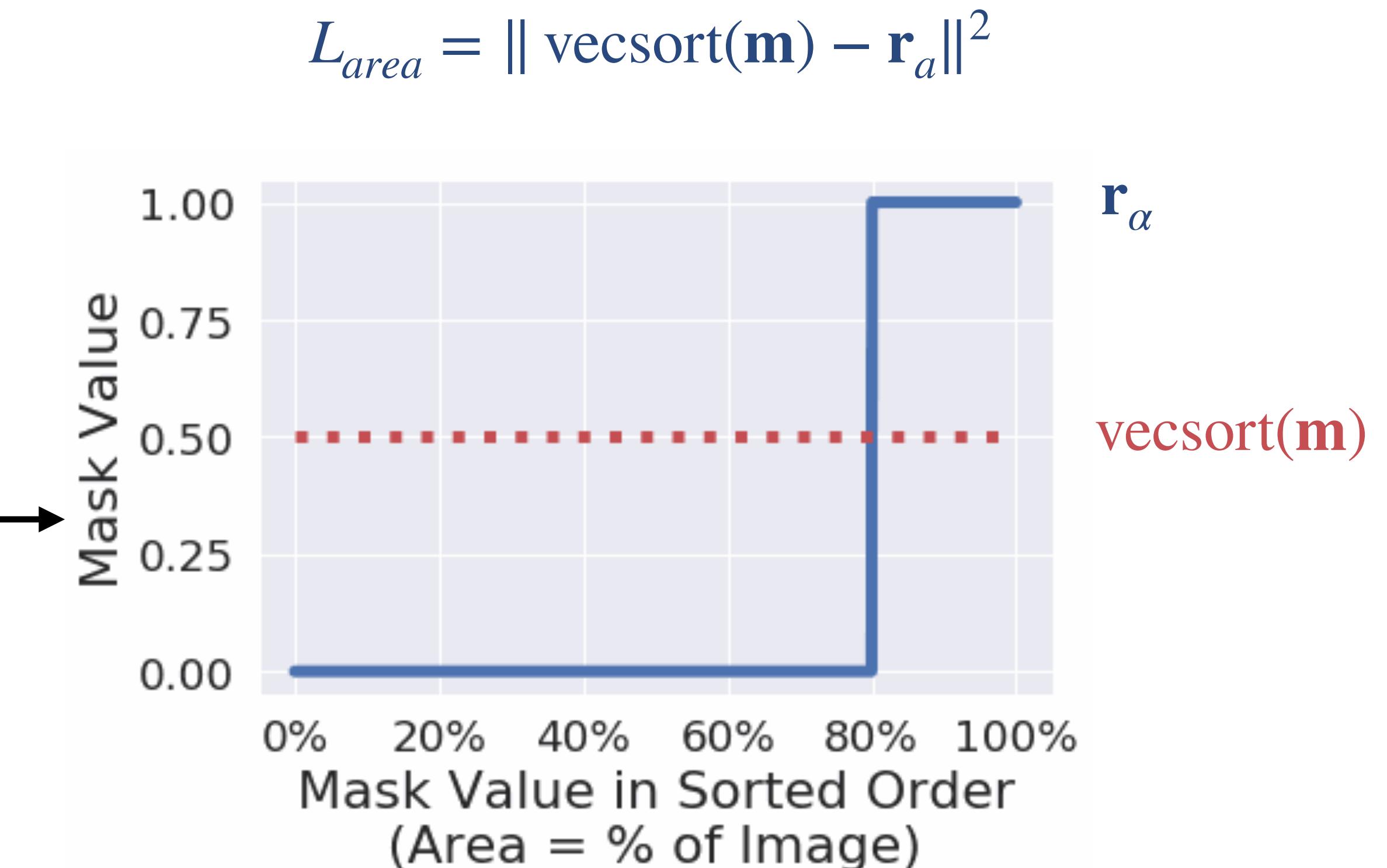
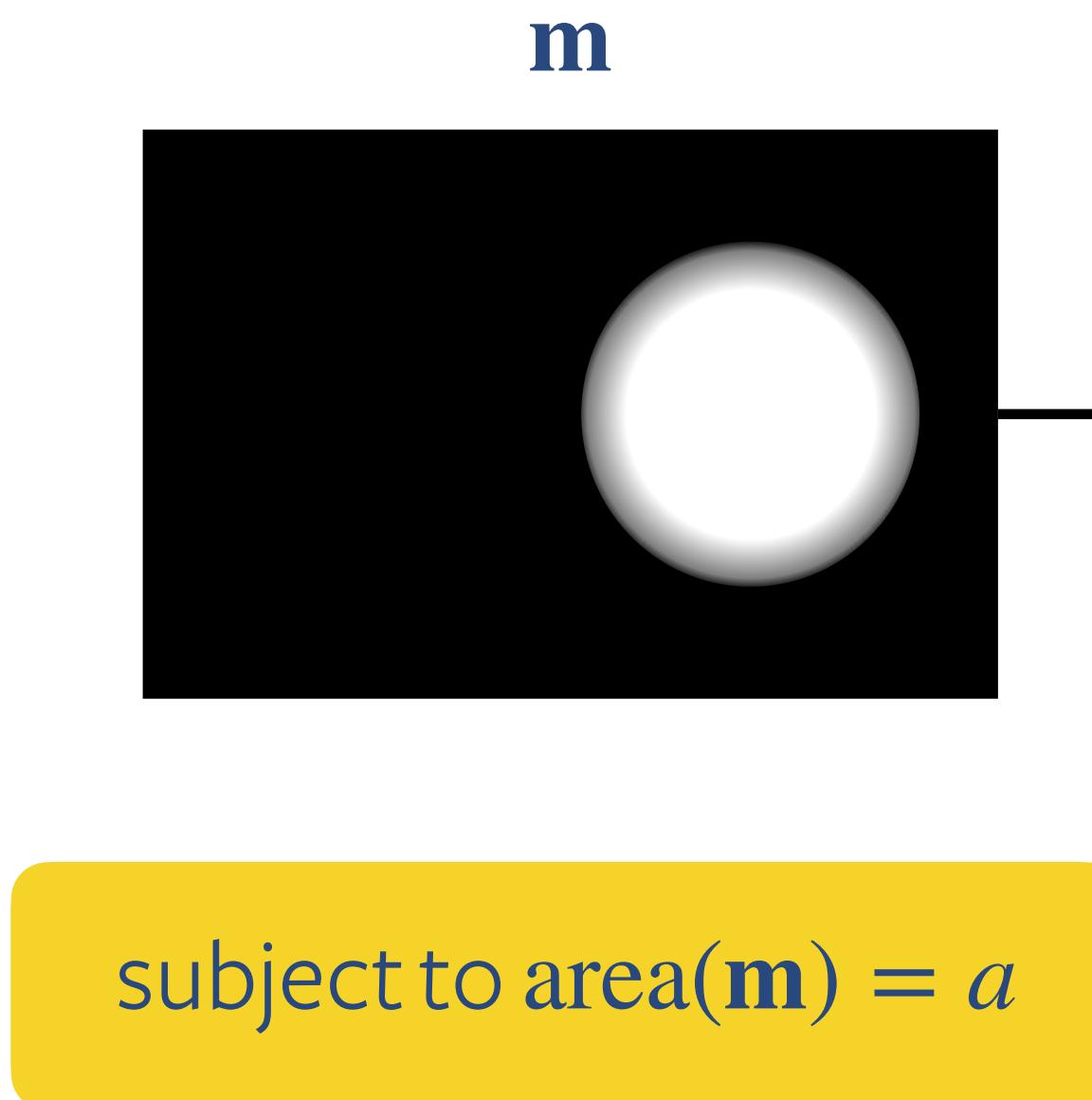
subject to $\operatorname{area}(\mathbf{m}) = a$



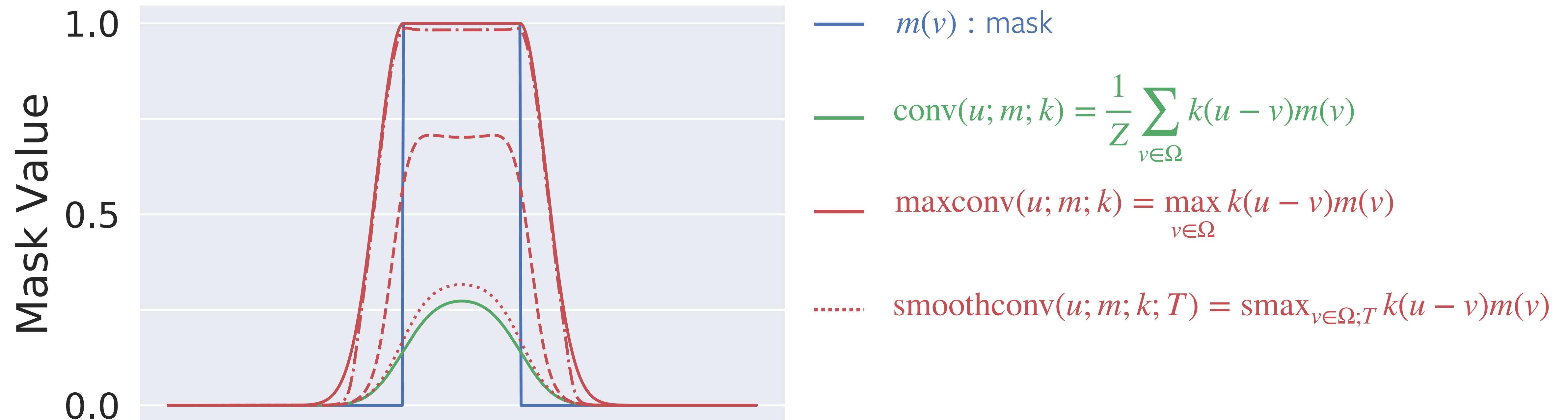
Area Constraint

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

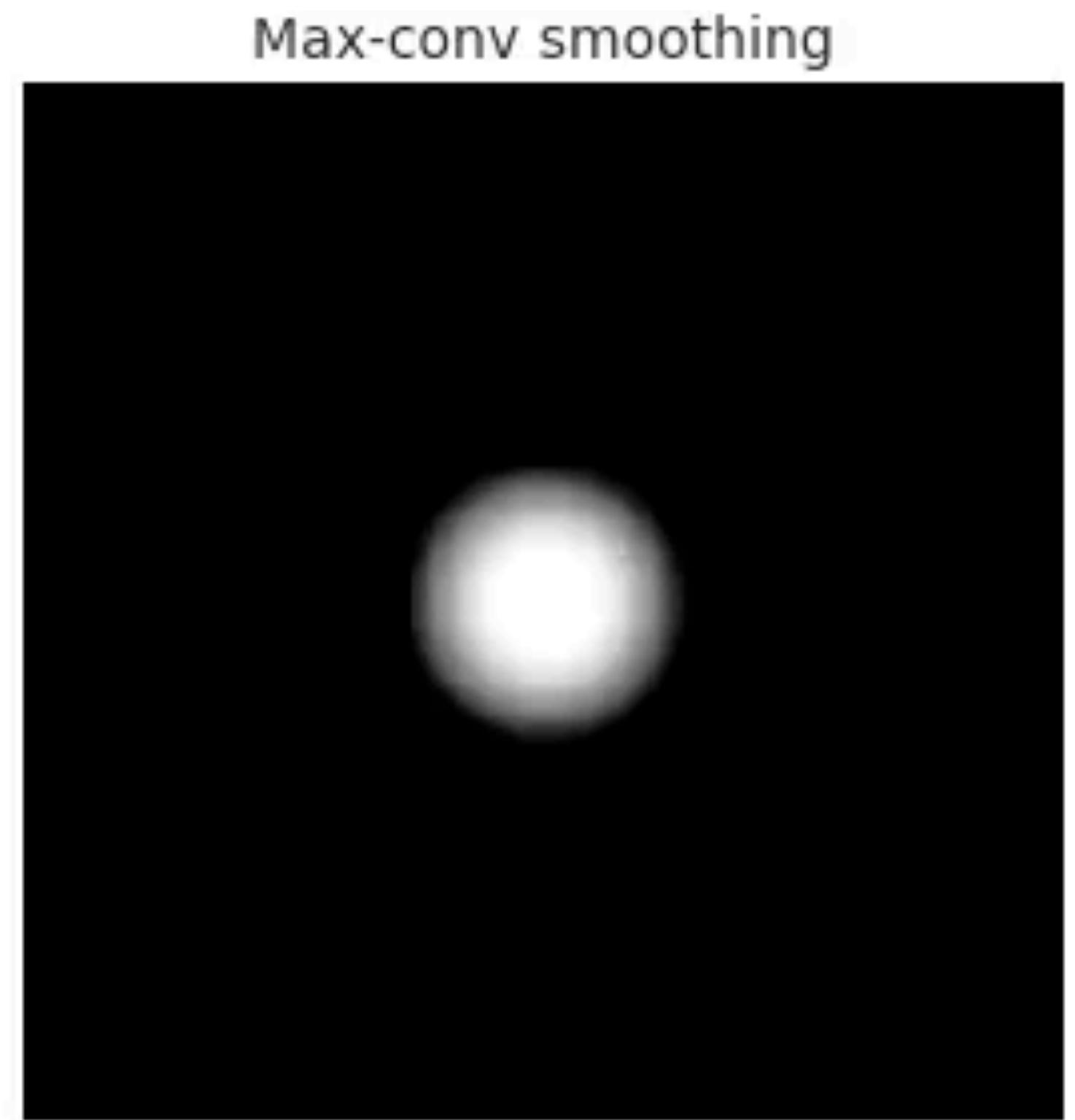
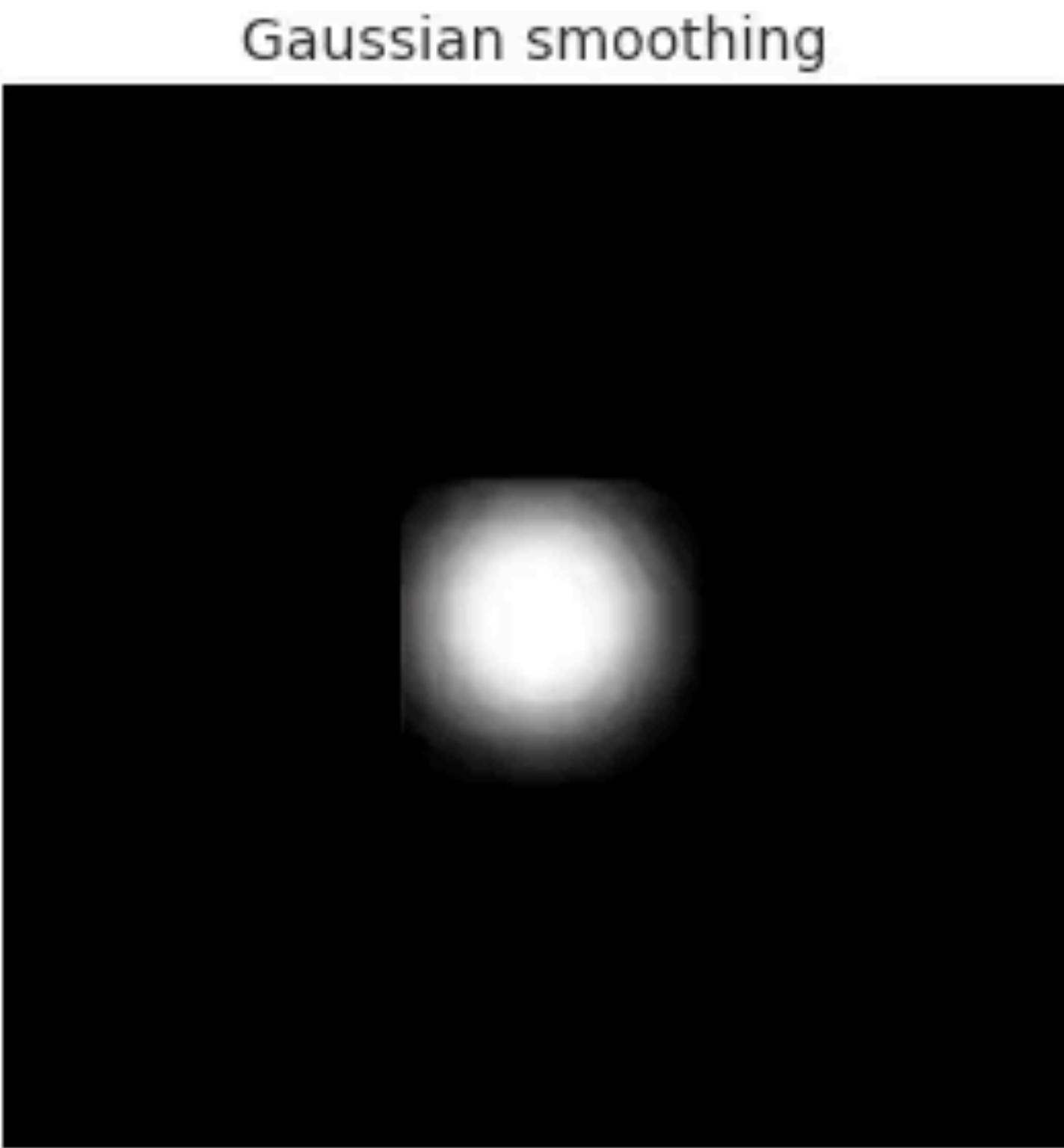
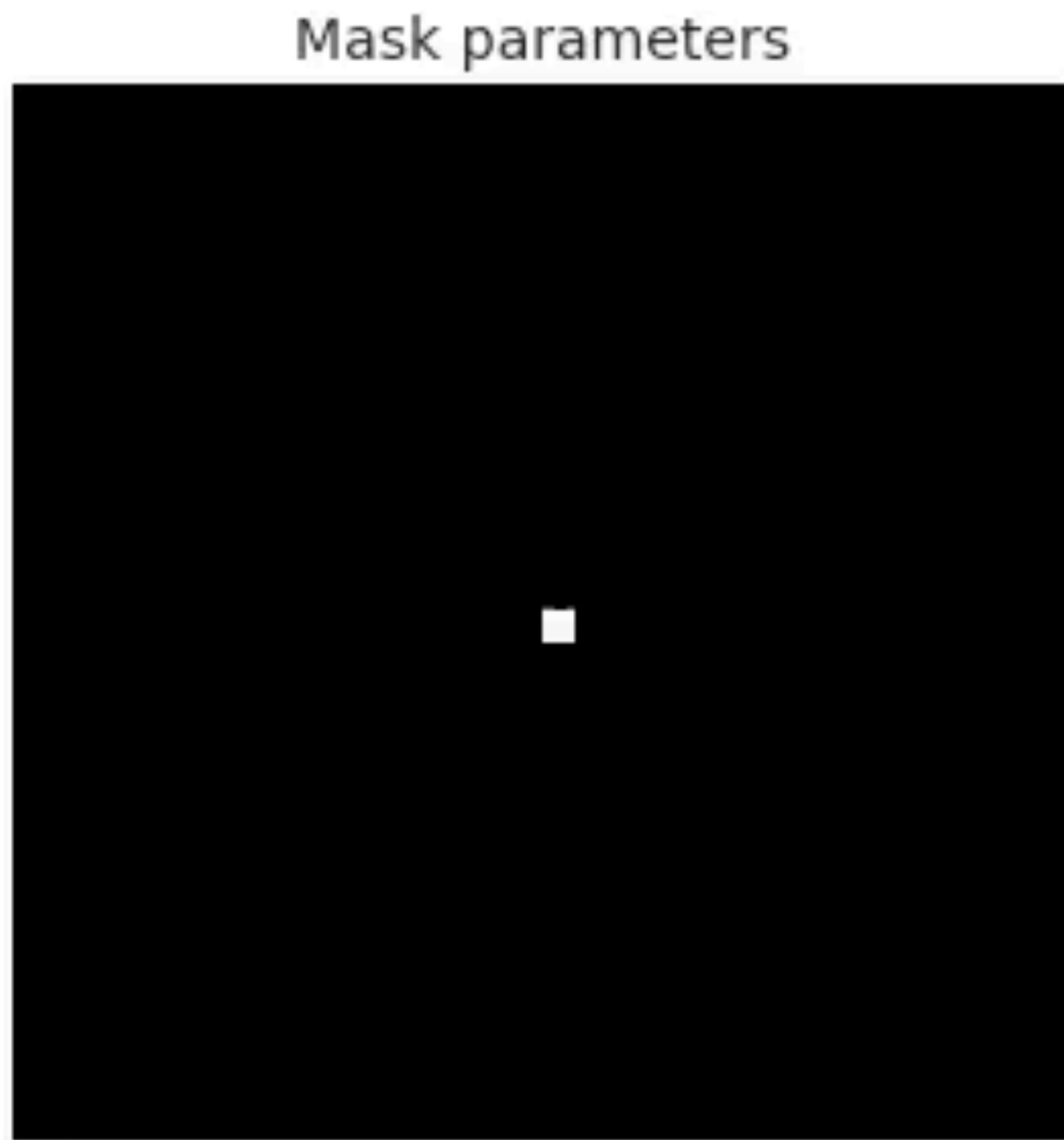
Here we re-formulate it as matching **rank statistics**



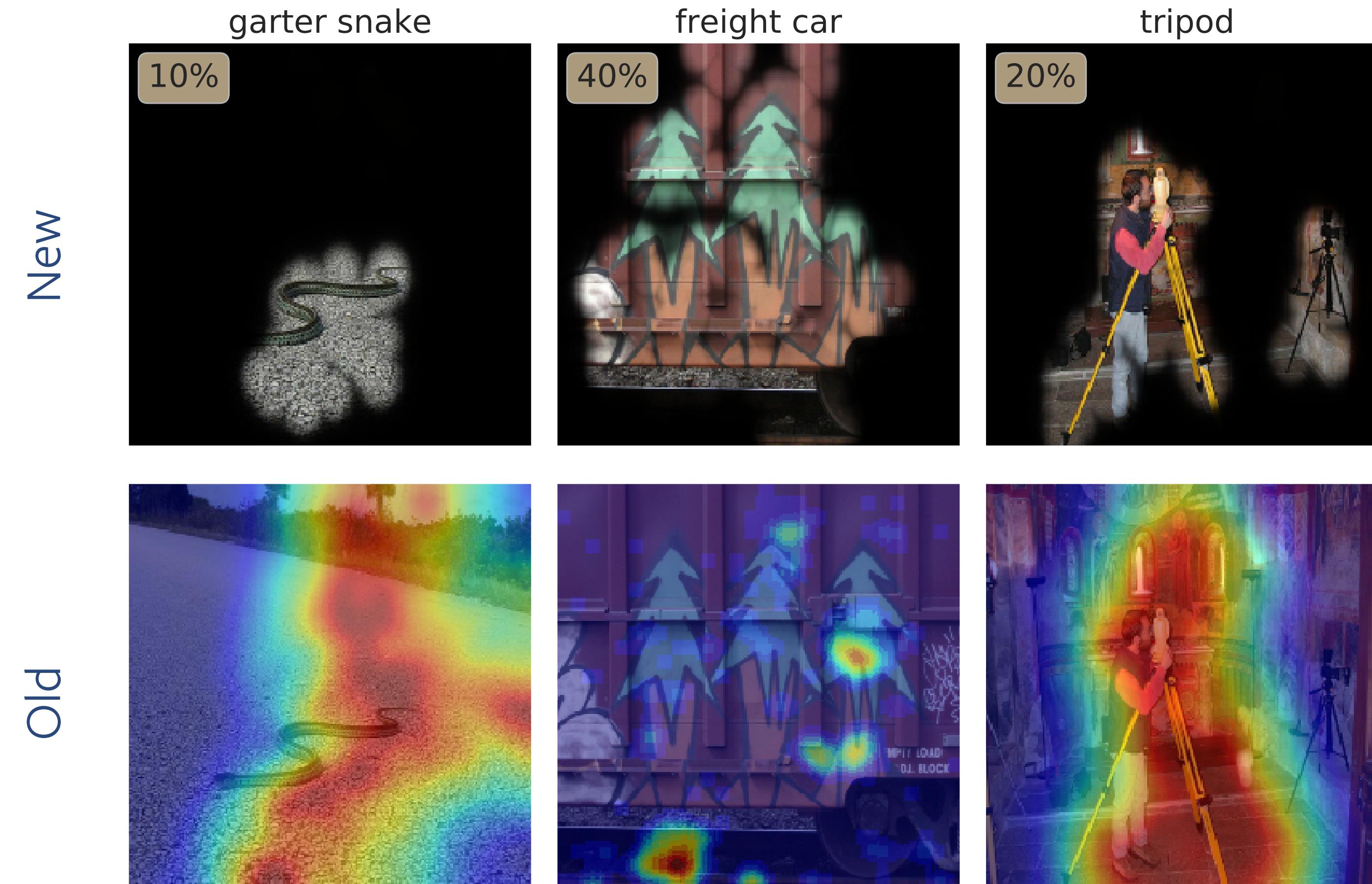
Smooth Masks



Smooth Masks



Comparison with Prior Work

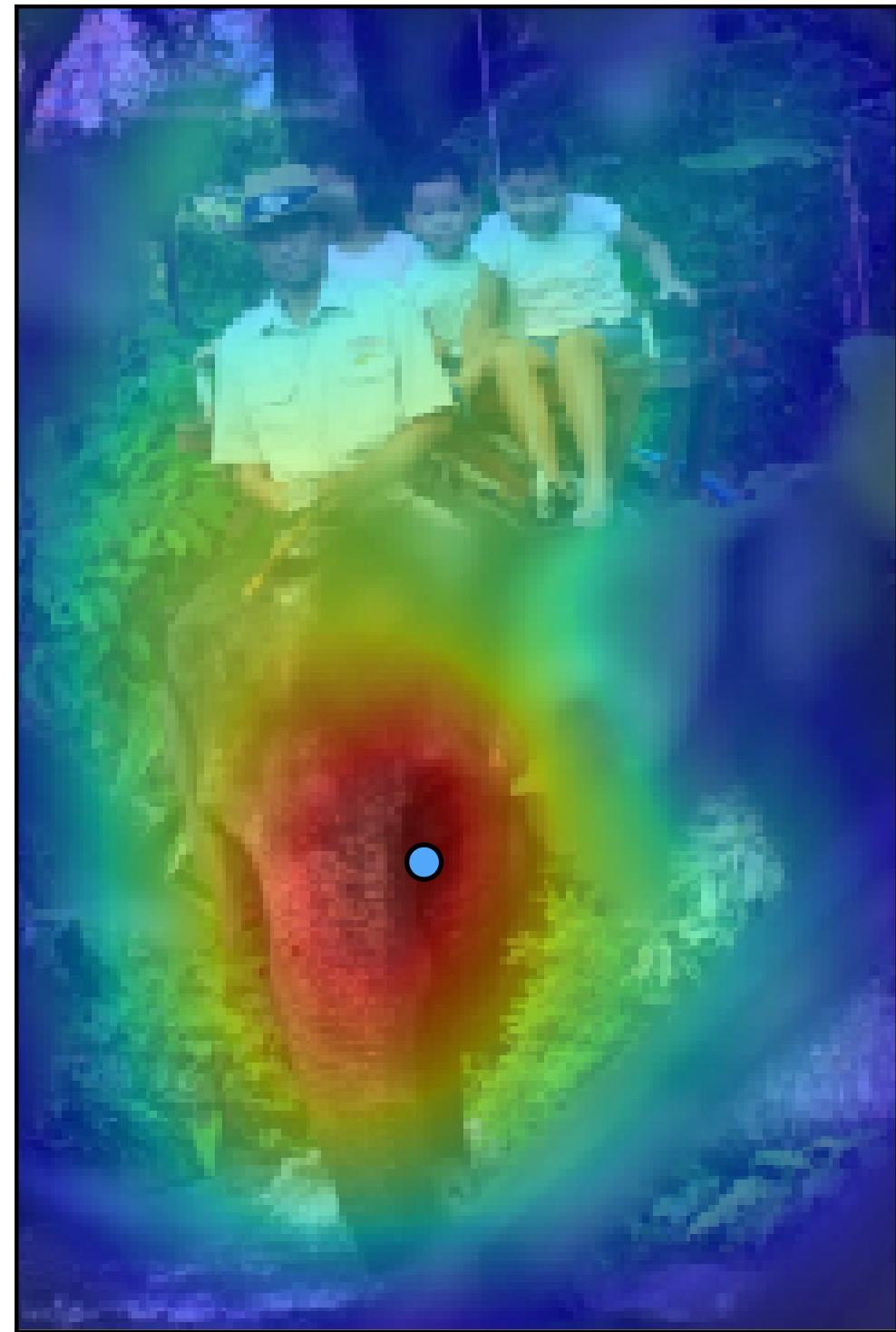


[Fong & Vedaldi, 2017; Fong et al., ICCV 2019] 38

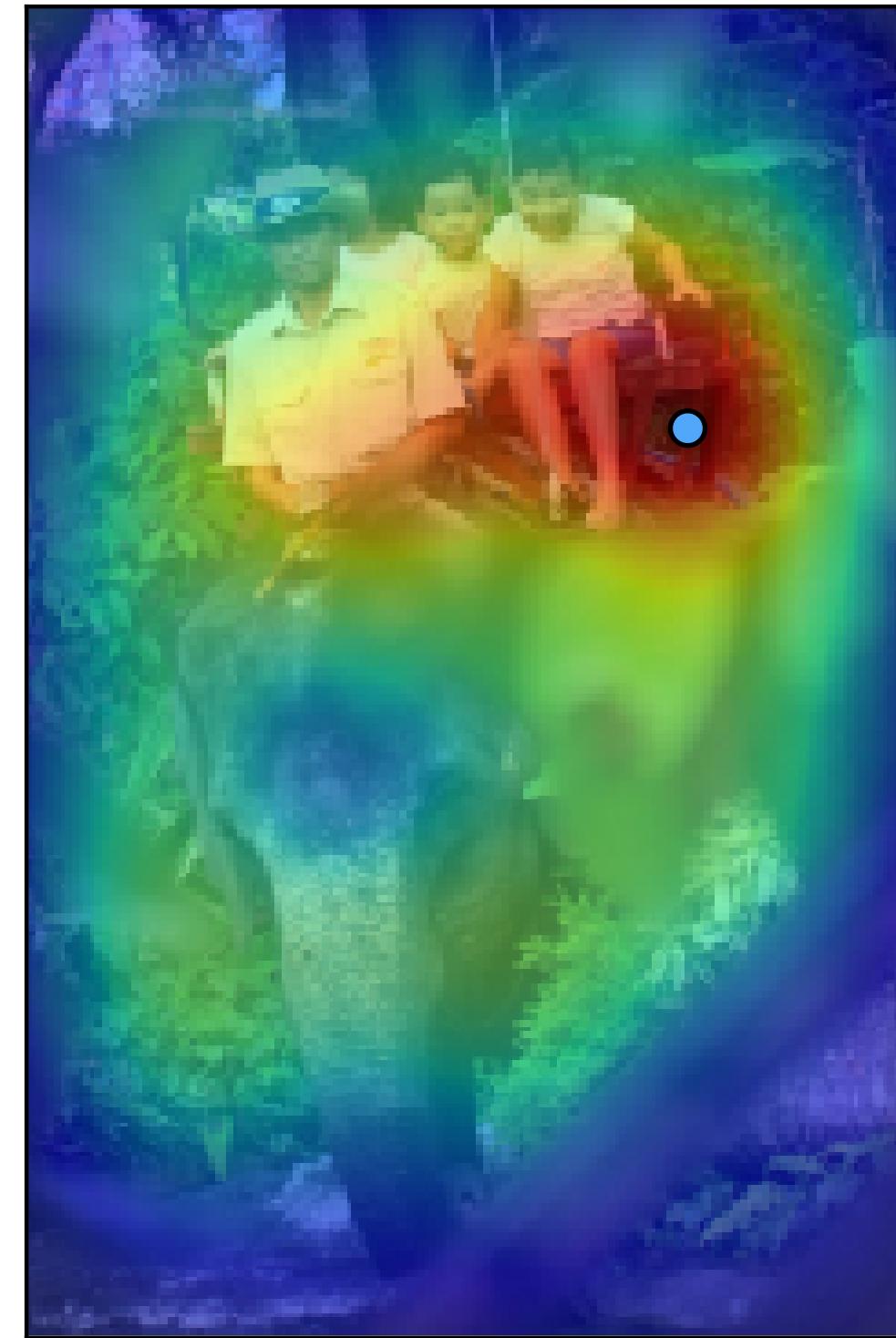
Evaluating and using attribution heatmaps

Measure Performance on Weak Localization

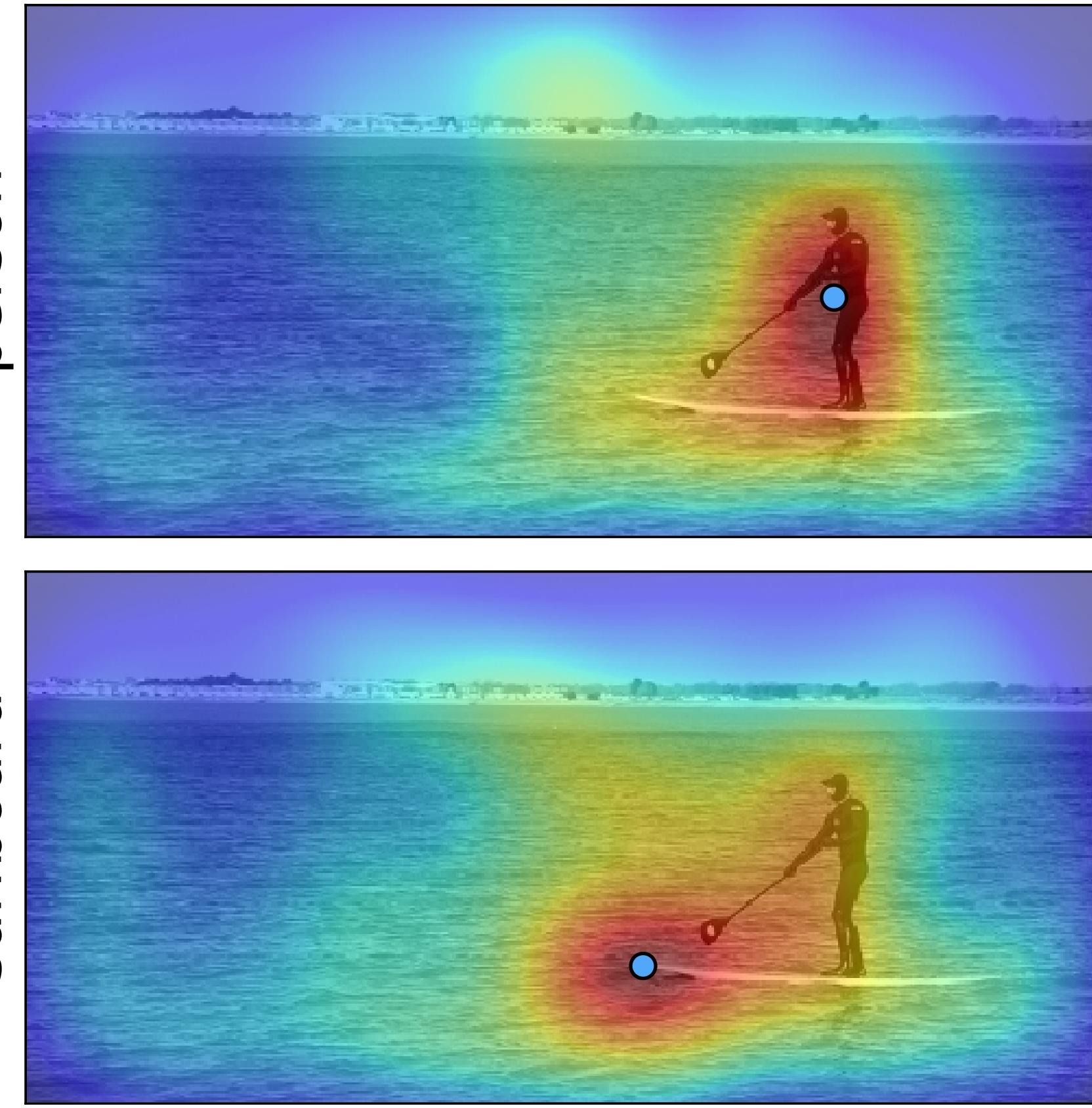
elephant



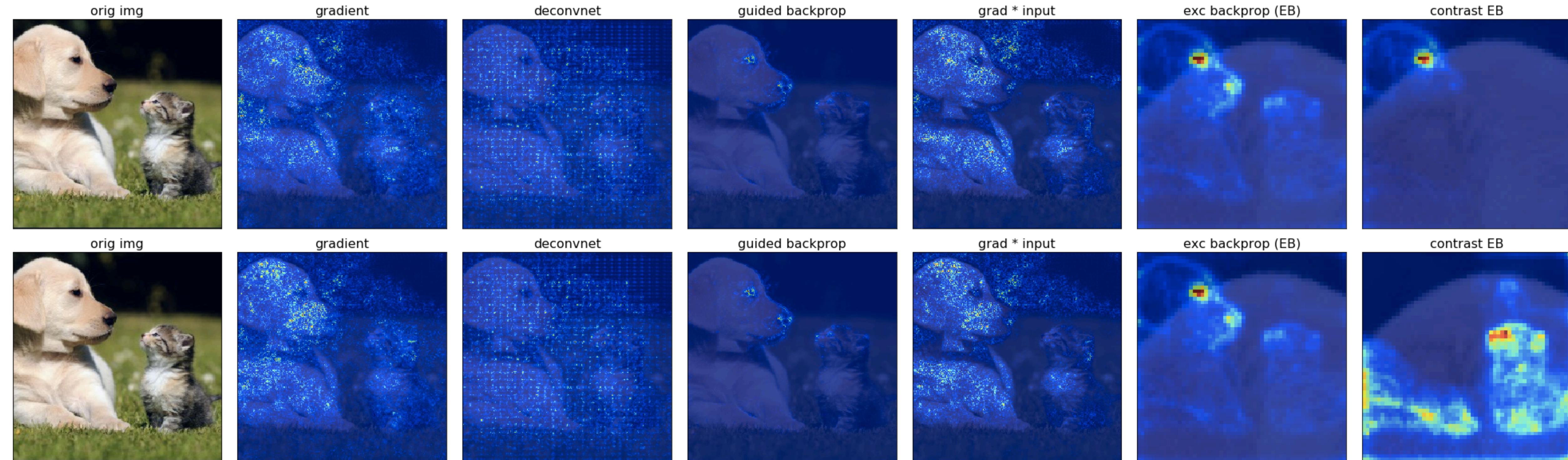
bench



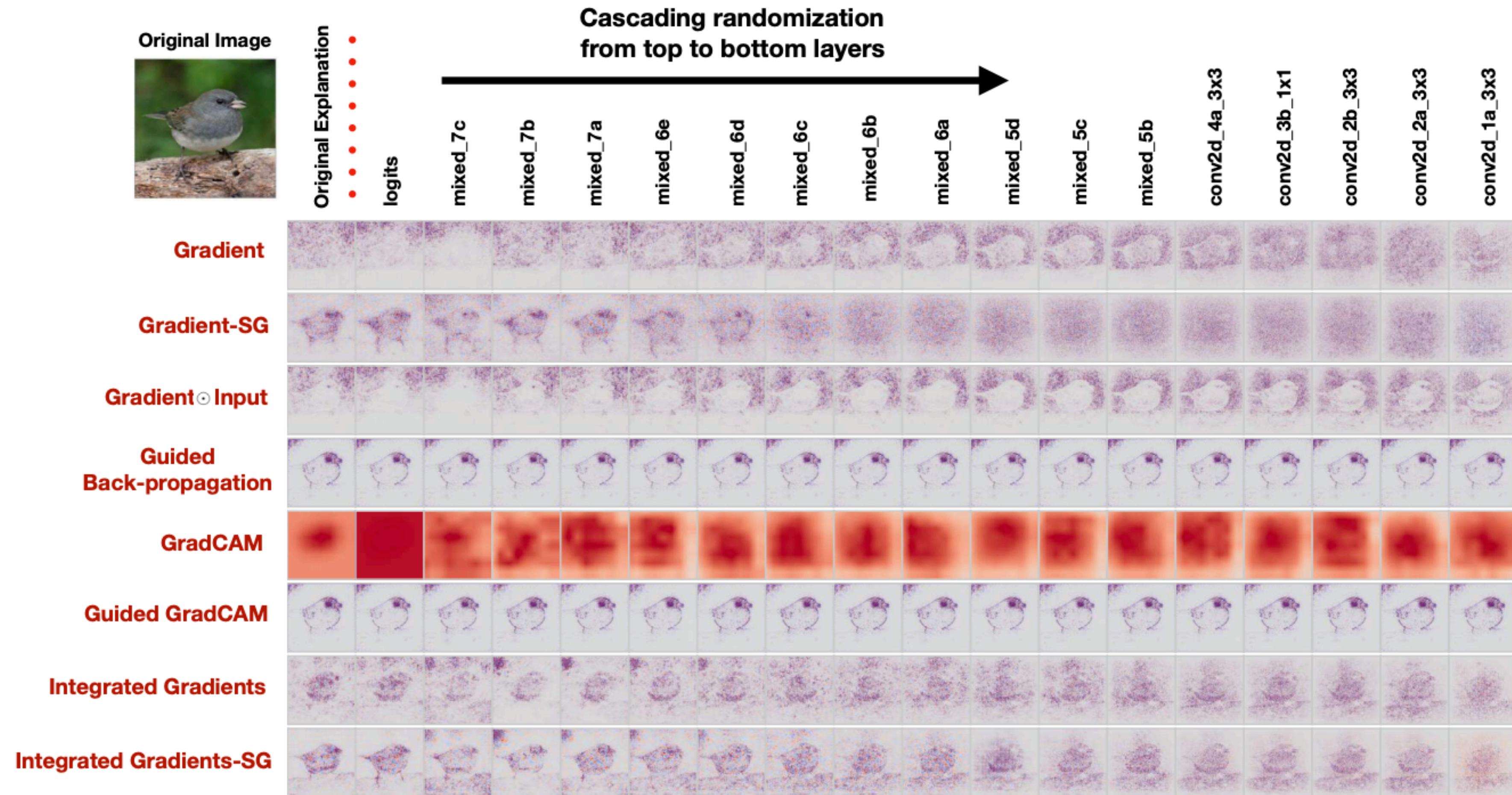
person
surfboard



Selectivity to Output Class



Sensitive to Model Parameters



Research Development: Critically design and evaluate attribution methods

General Usage: Assume a model has failures and use attribution methods to understand them

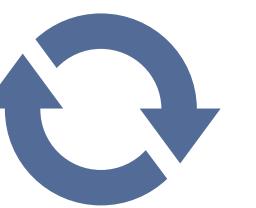
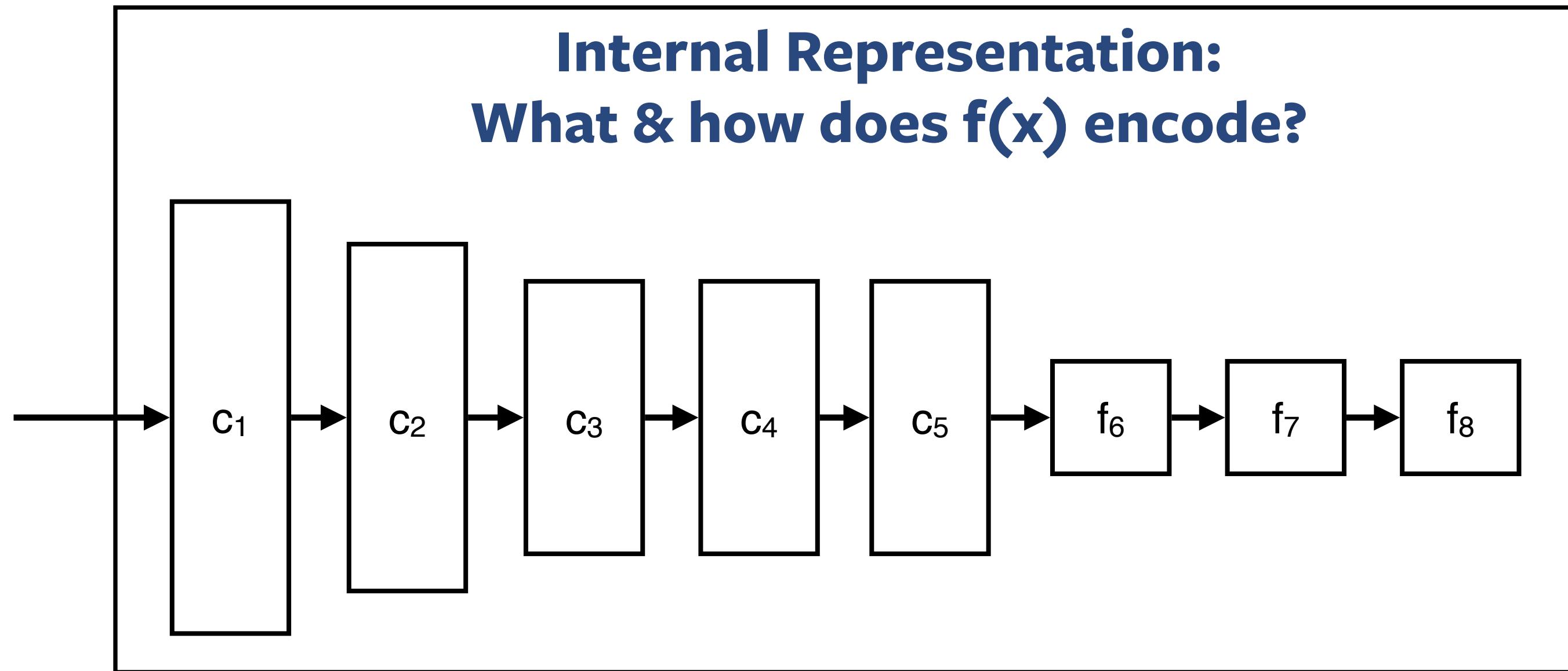
TorchRay: PyTorch interpretability library

github.com/facebookresearch/torchray

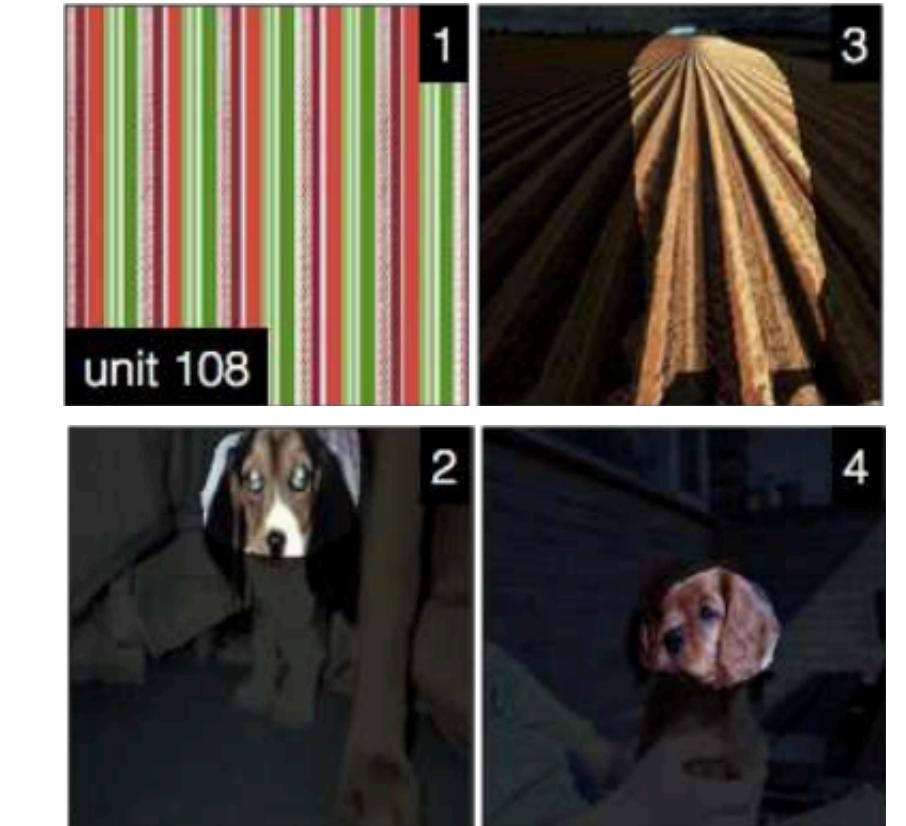


Research Themes

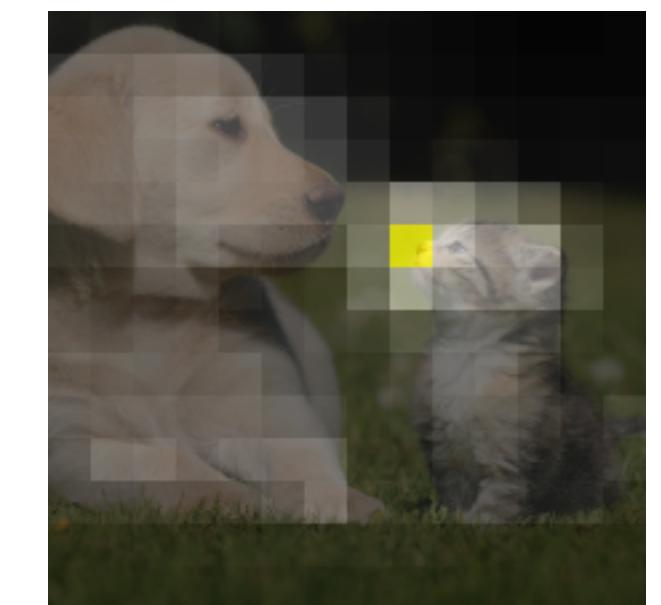
Fong et al., ICCV 2019



$$\theta \leftarrow \theta - \eta \frac{dL}{d\theta}$$



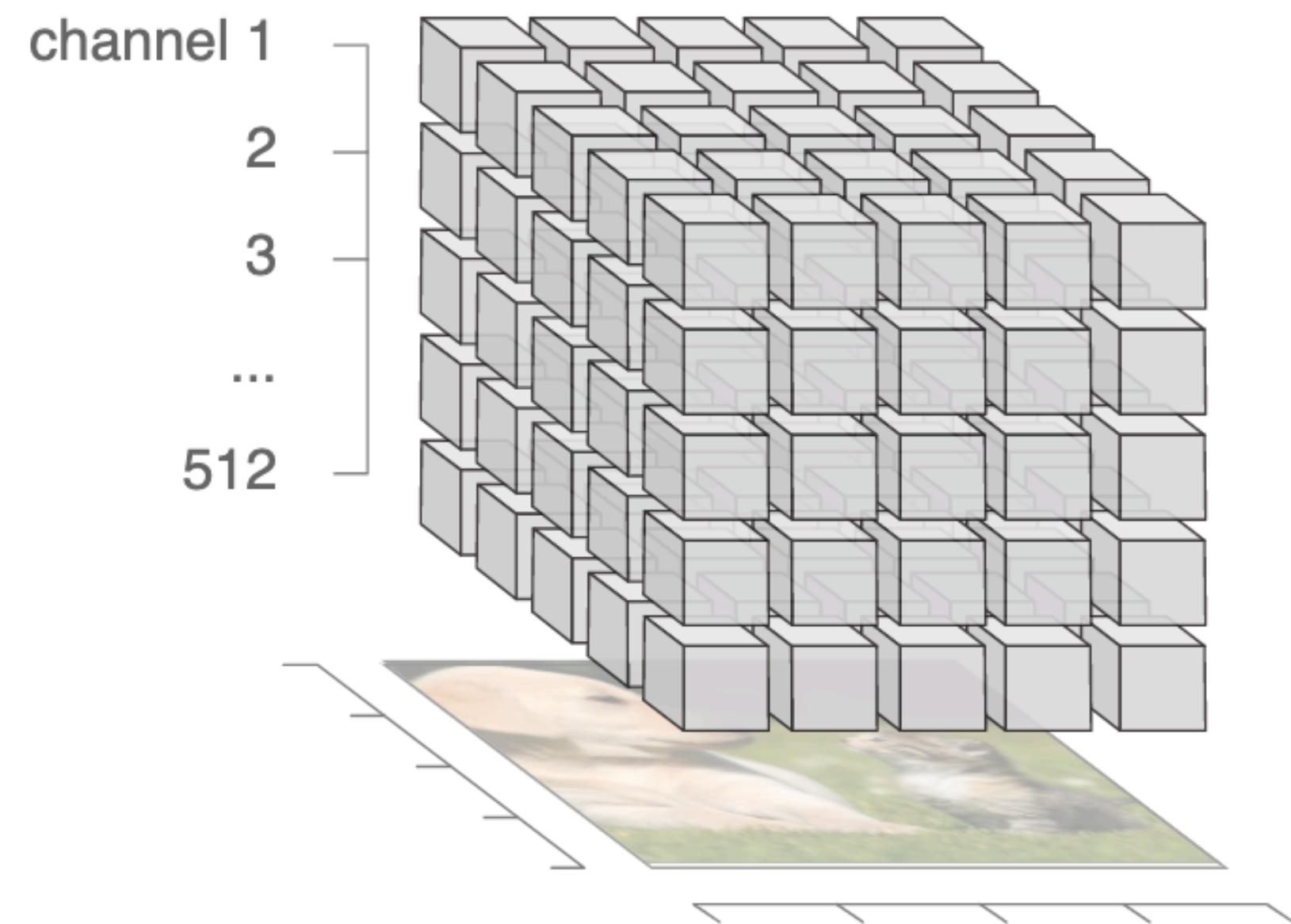
Fong & Vedaldi, CVPR 2018



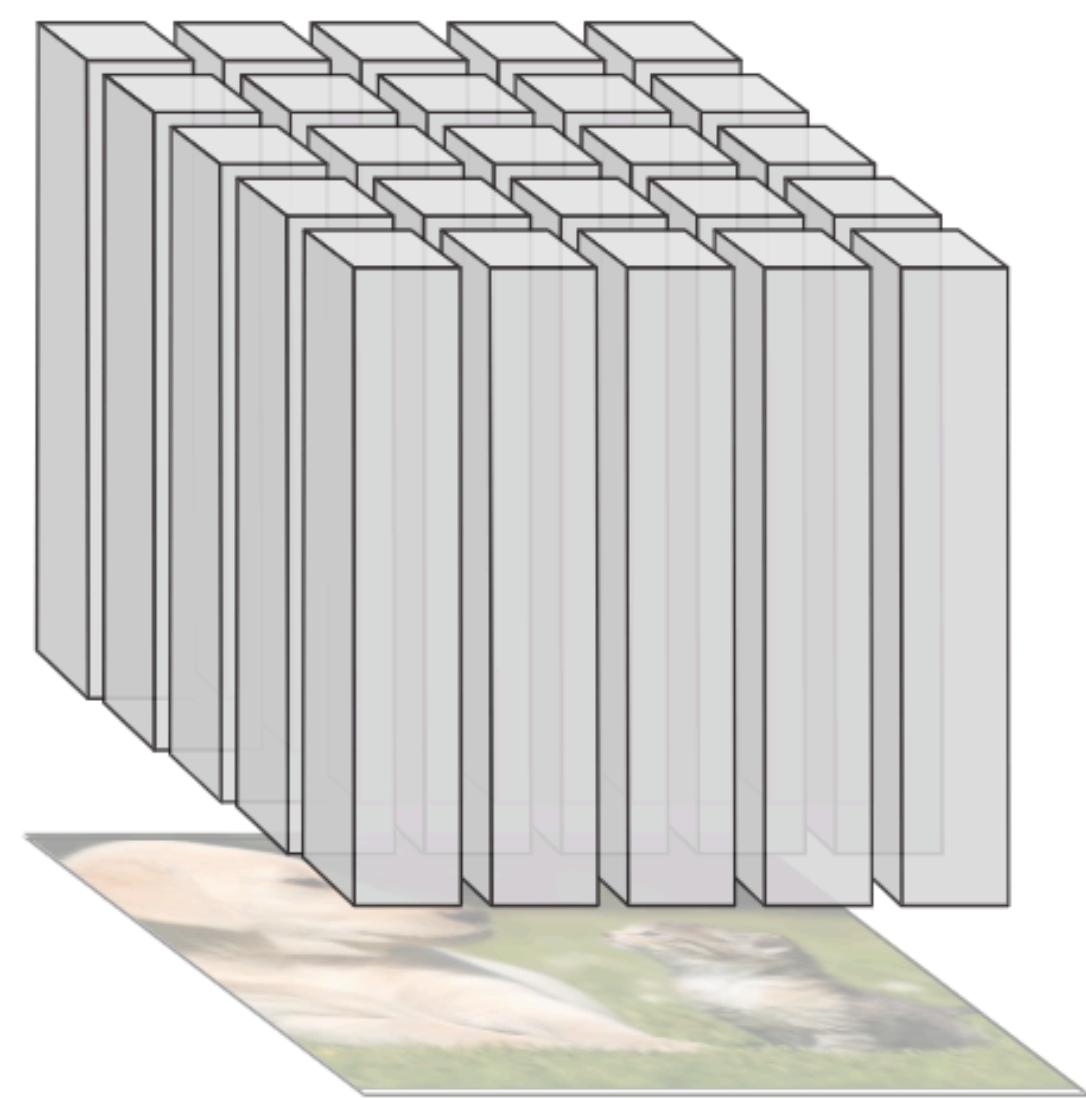
Fong et al., 2020 (in prep.)

Intermediate Activations

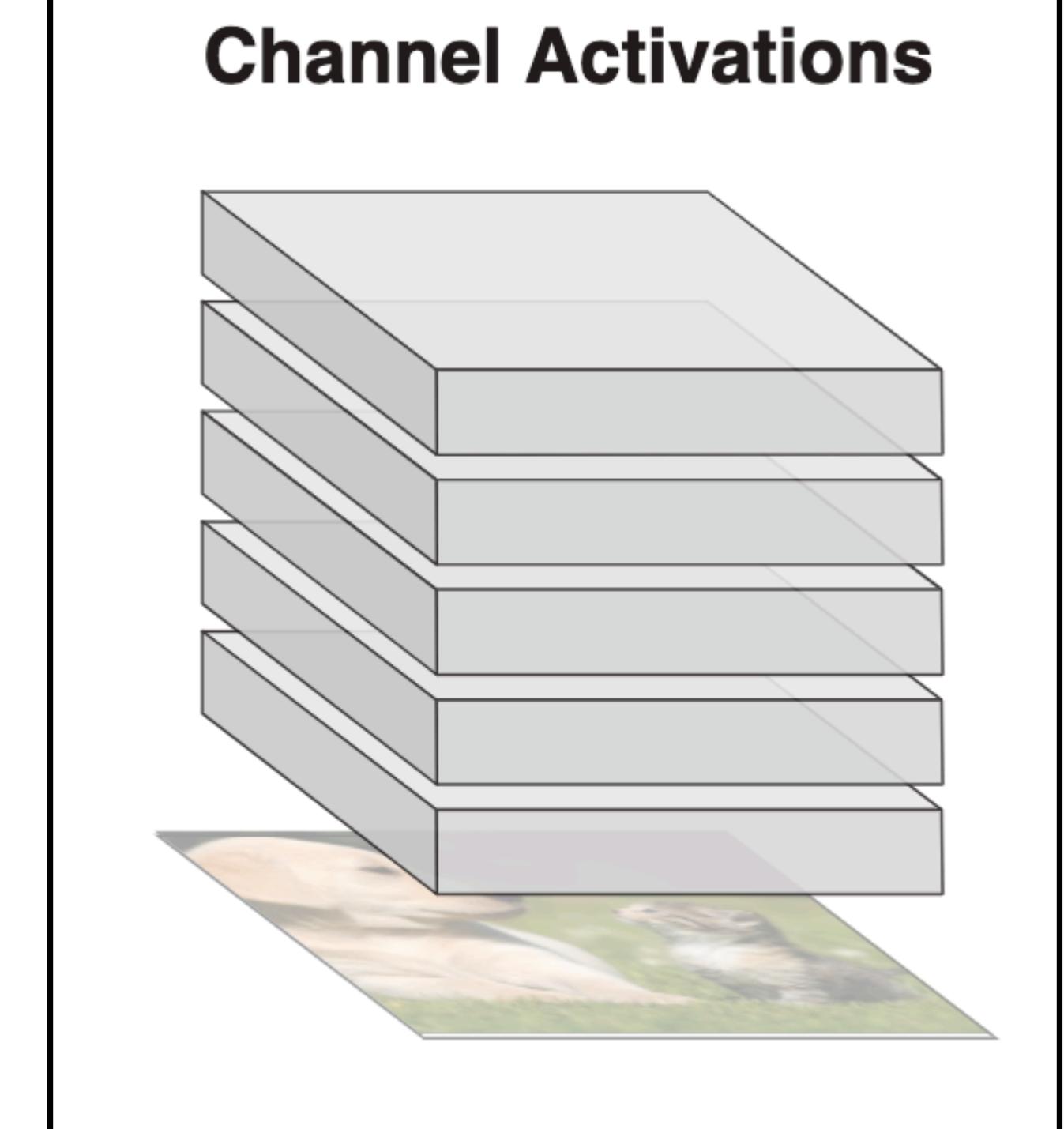
Individual Neurons



Spatial Activations



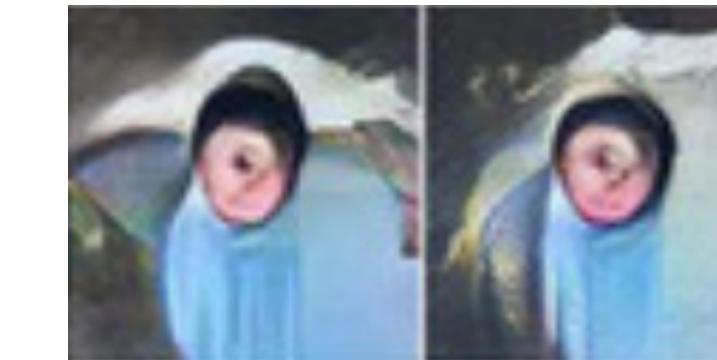
Channel Activations



Zeiler & Fergus, ECCV 2014



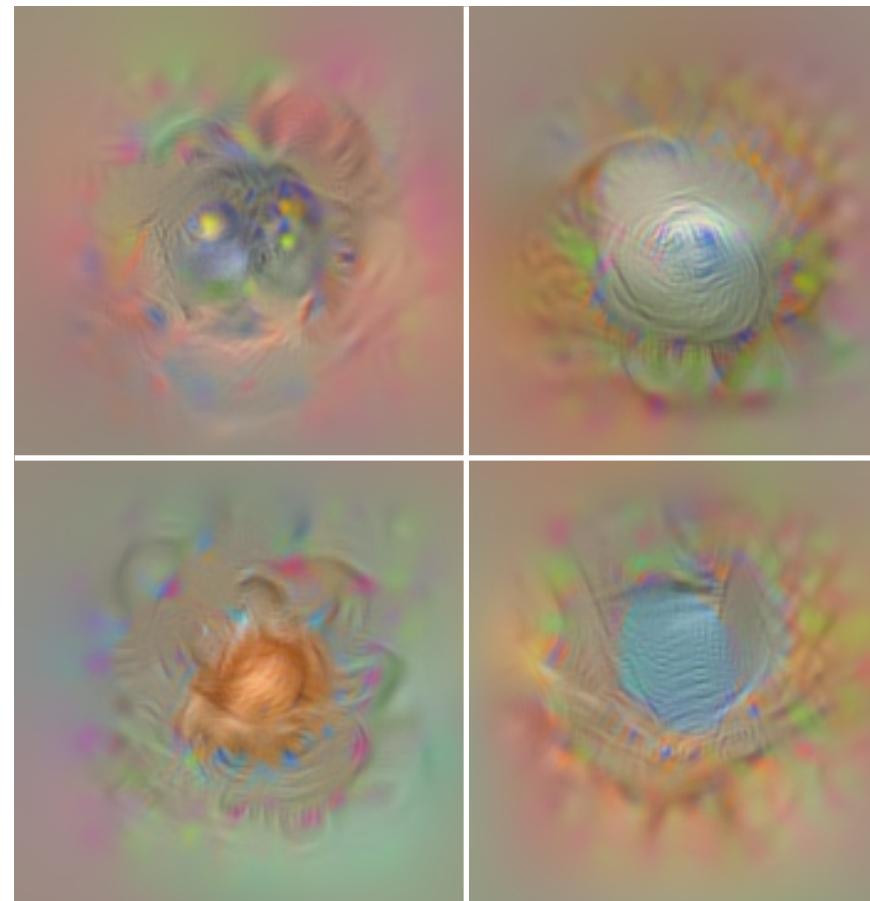
Nyugen et al., NIPS 2016



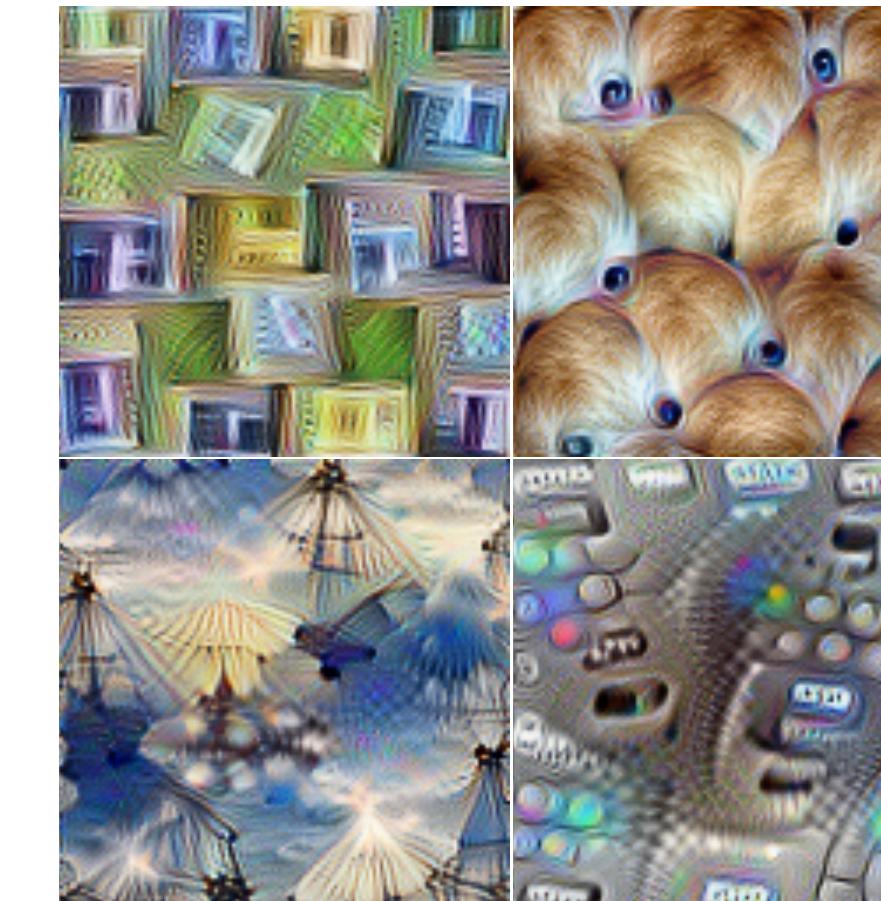
Zhou et al., ICLR 2015



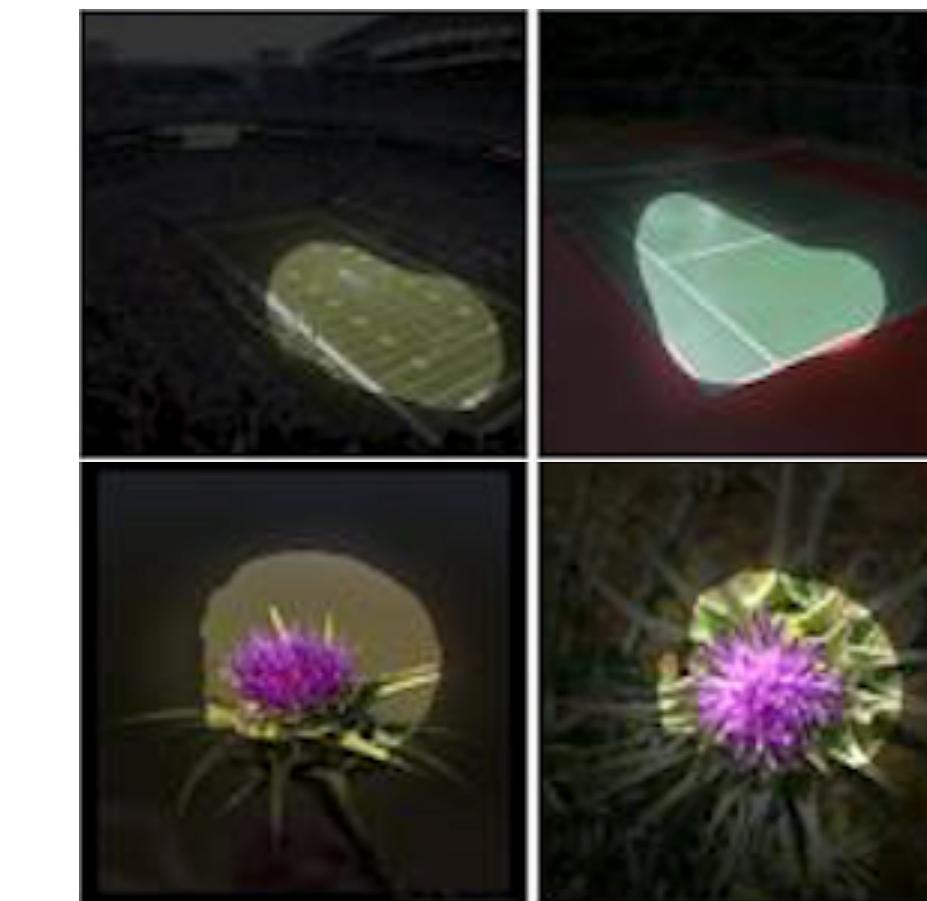
Mahendran & Vedaldi, IJCV 2016



Olah et al., Distill 2017



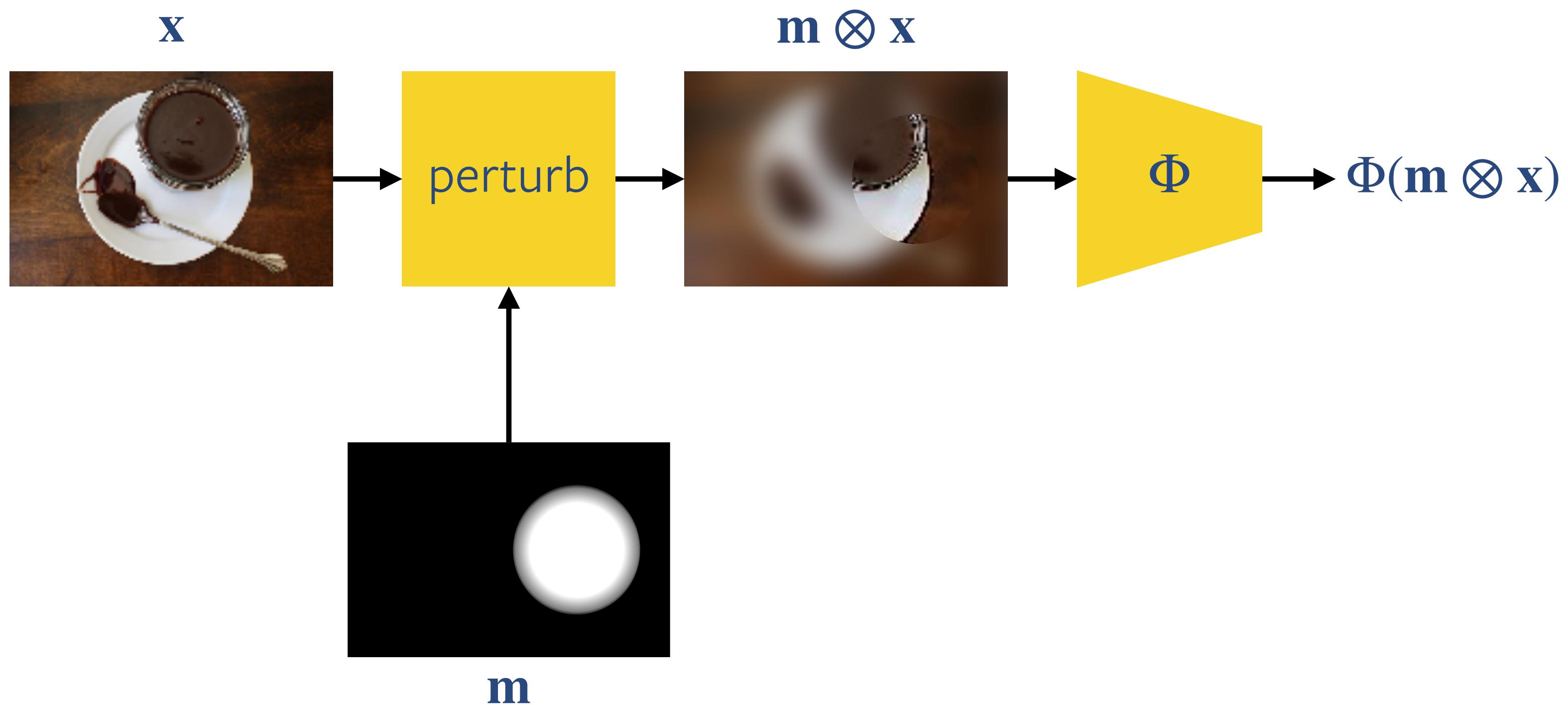
Bau et al., CVPR 2017



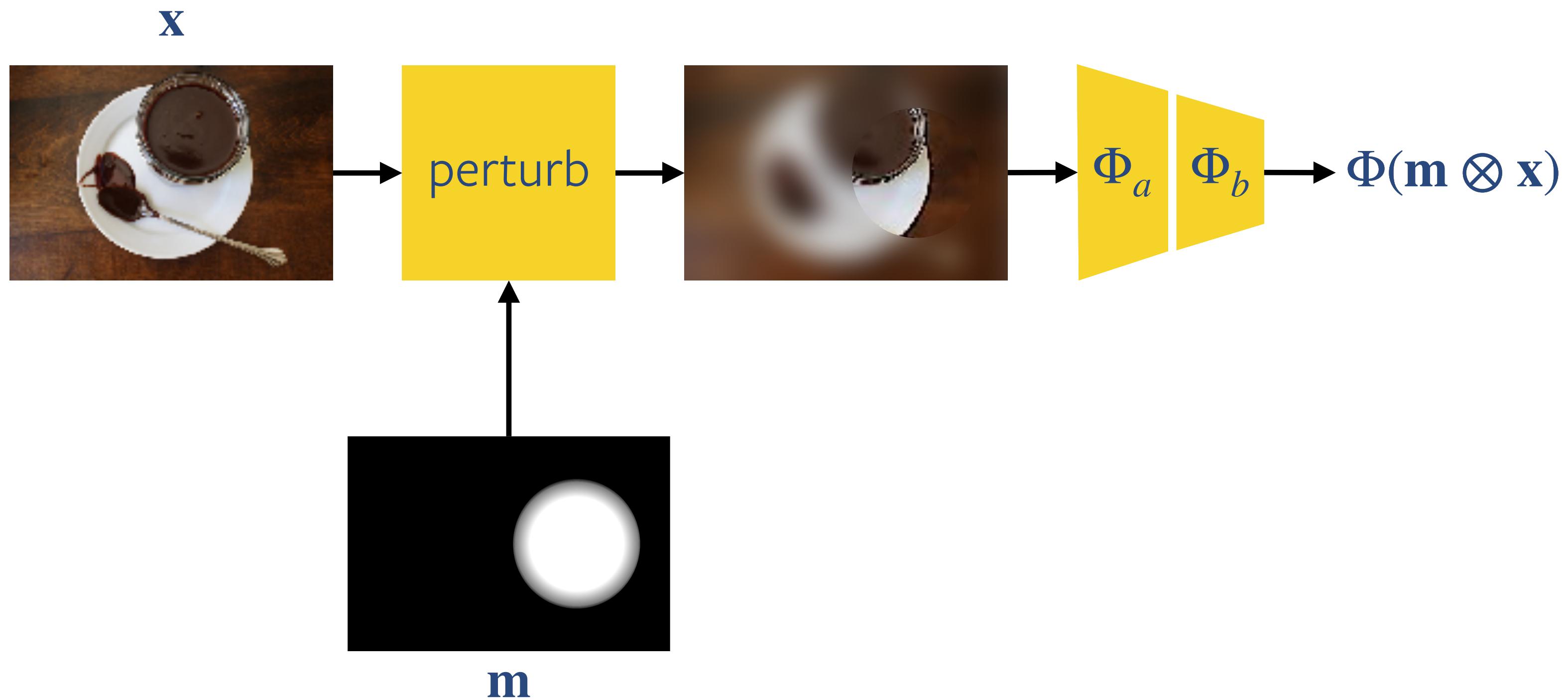
Most prior work focuses on visualizing **single channels**.

A. Attributing channels in intermediate activations

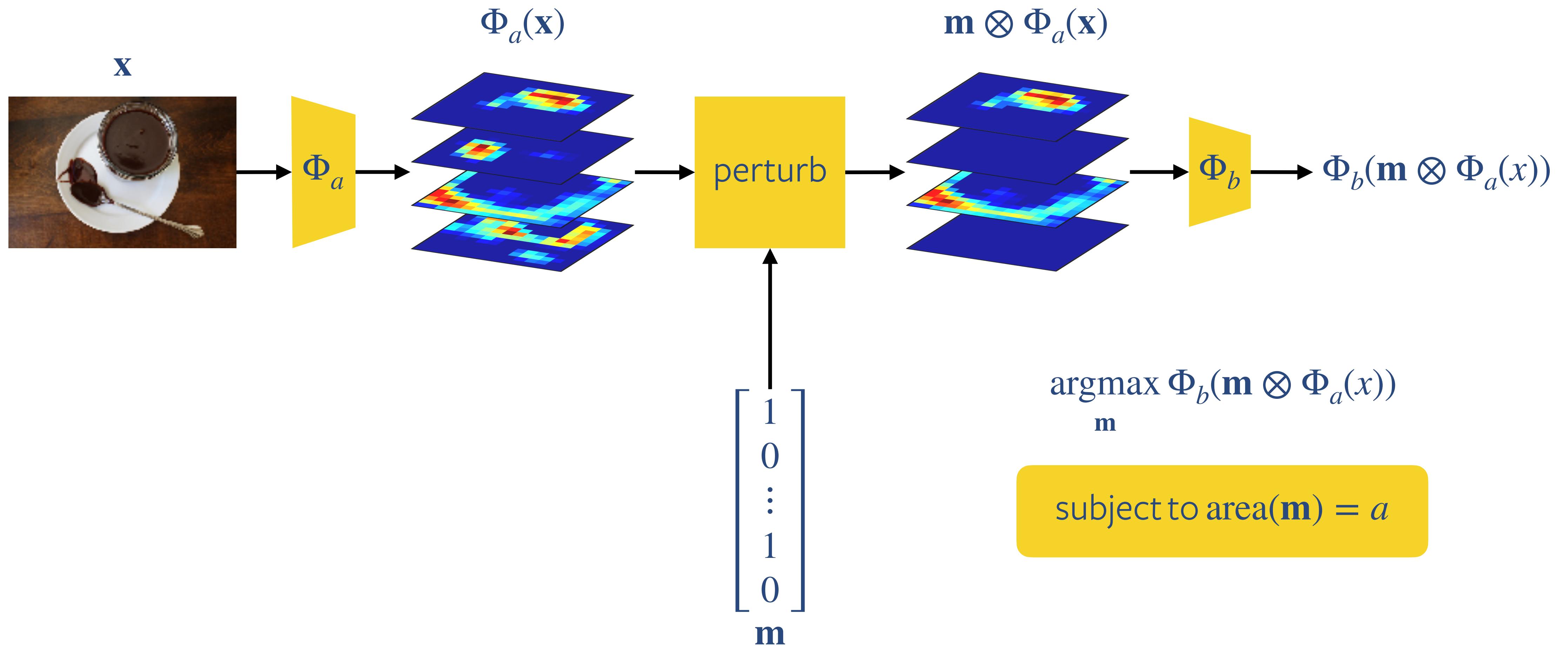
Spatial Attribution



Channel Attribution



Channel Attribution



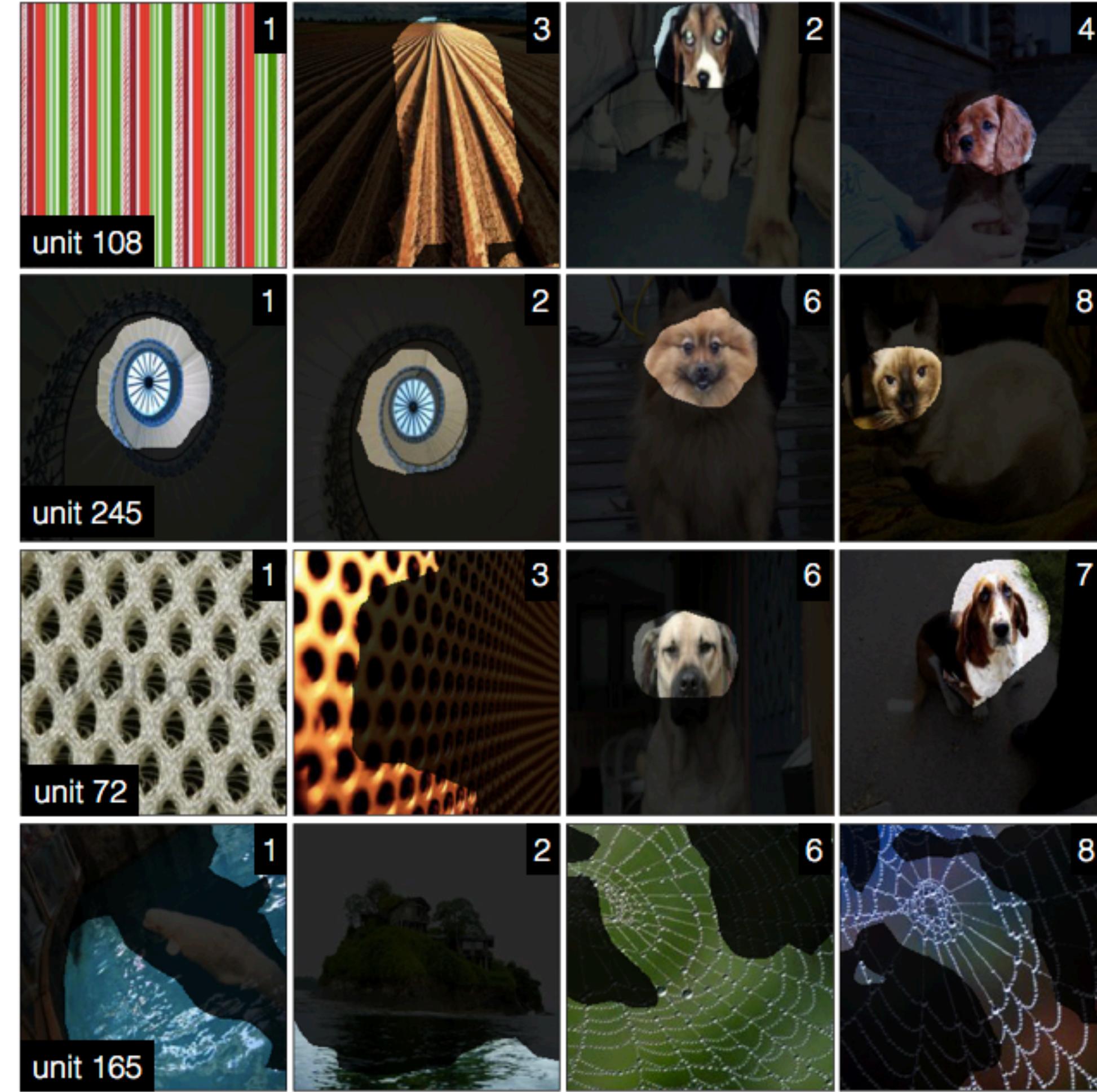
Activation “Diffing”



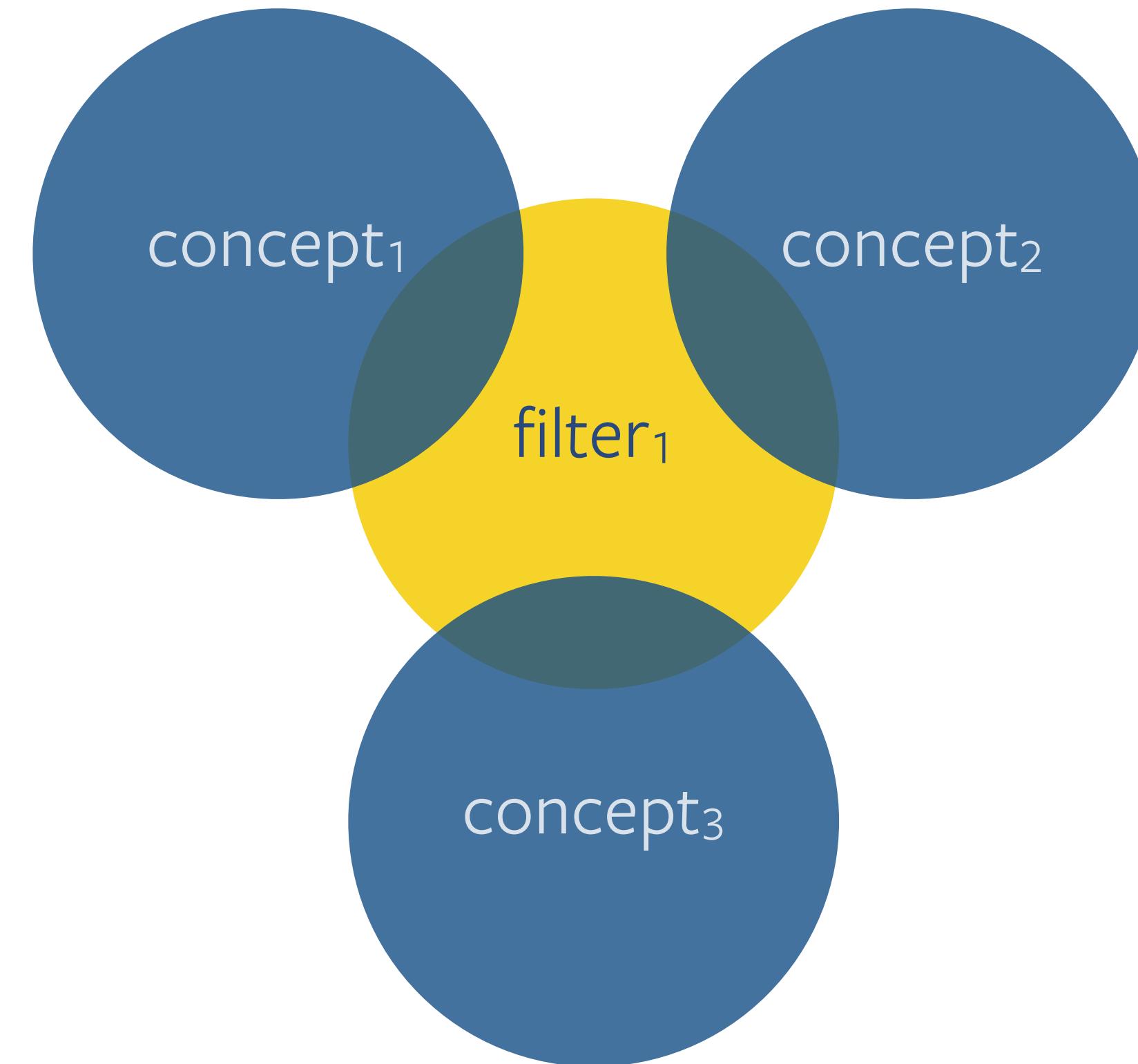
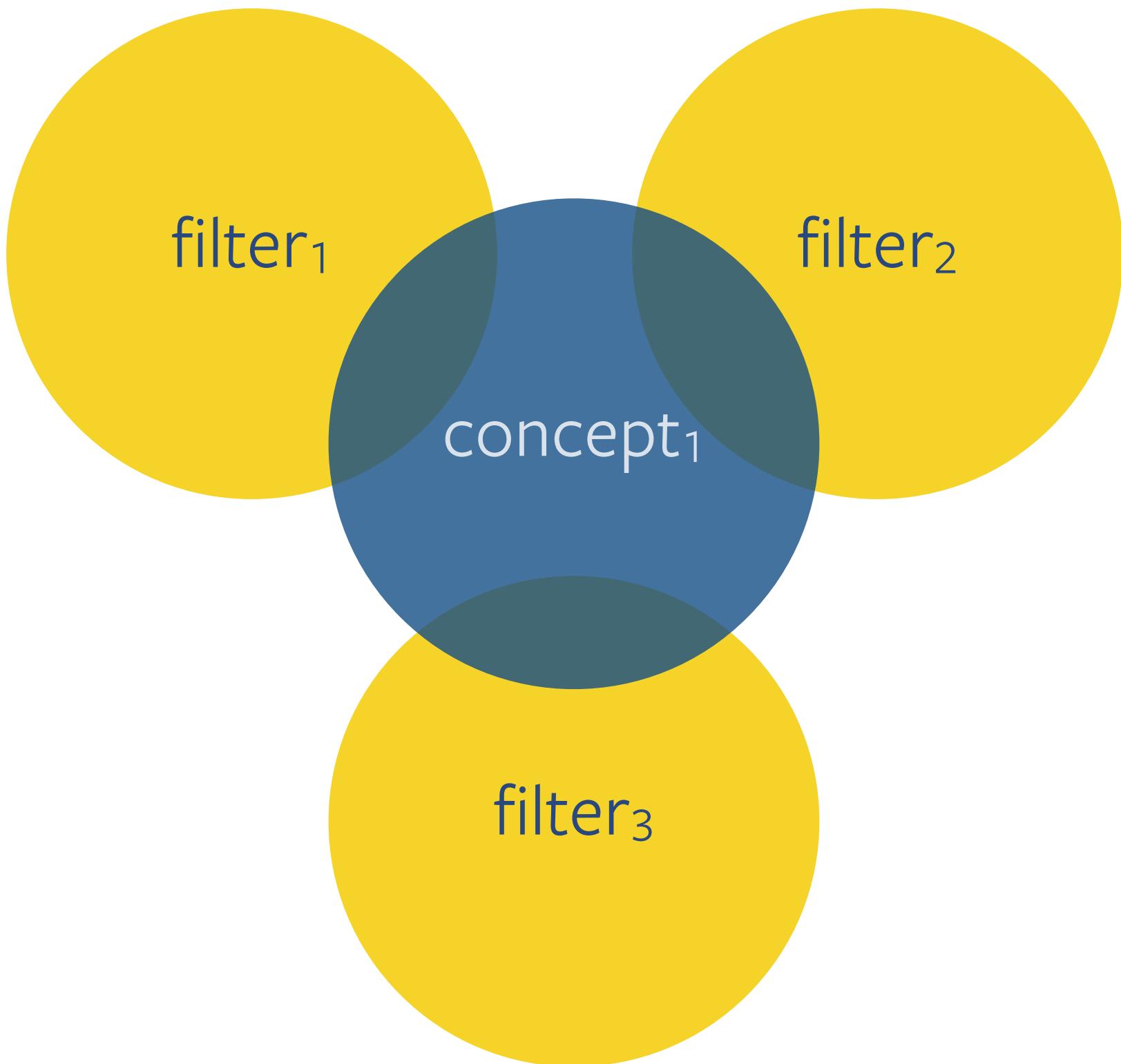
B. Understanding how semantic concepts are encoded

Filter-Concept Overlap

Top activated patches for specific units in AlexNet conv5 filters points to a **packing phenomenon**.



Filter-Concept Overlap



Net2Vec

Learn **concept vectors** that describe how a concept is encoded **across** channels.

Probe a **network** with a **concept dataset** and learn to perform a **task** using **activations** at a given layer.

Results

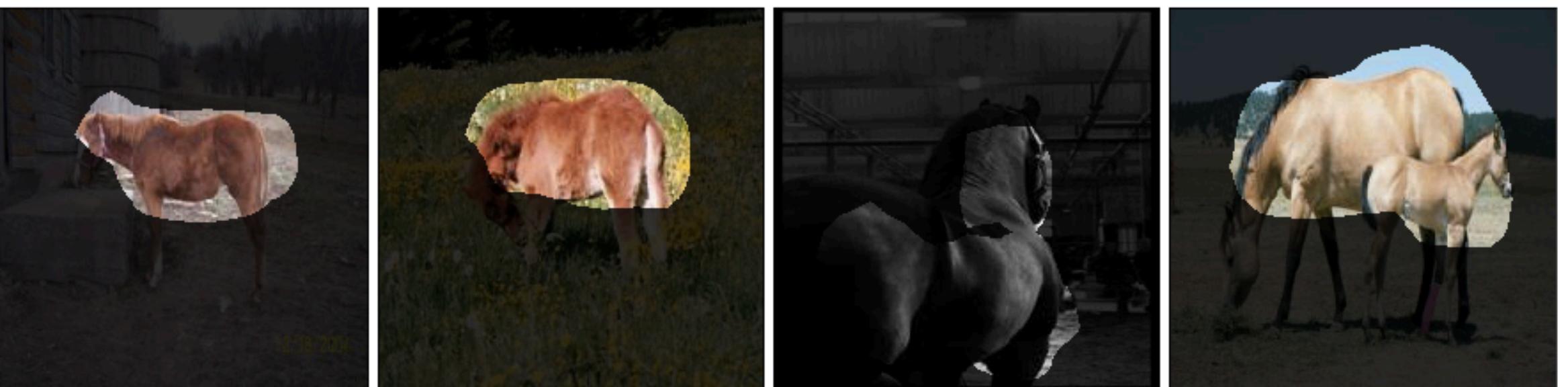
Concepts per Filter

AlexNet conv5 unit 66 is highly selective for various farm animals

Sheep
(IoU_{set} = .21)



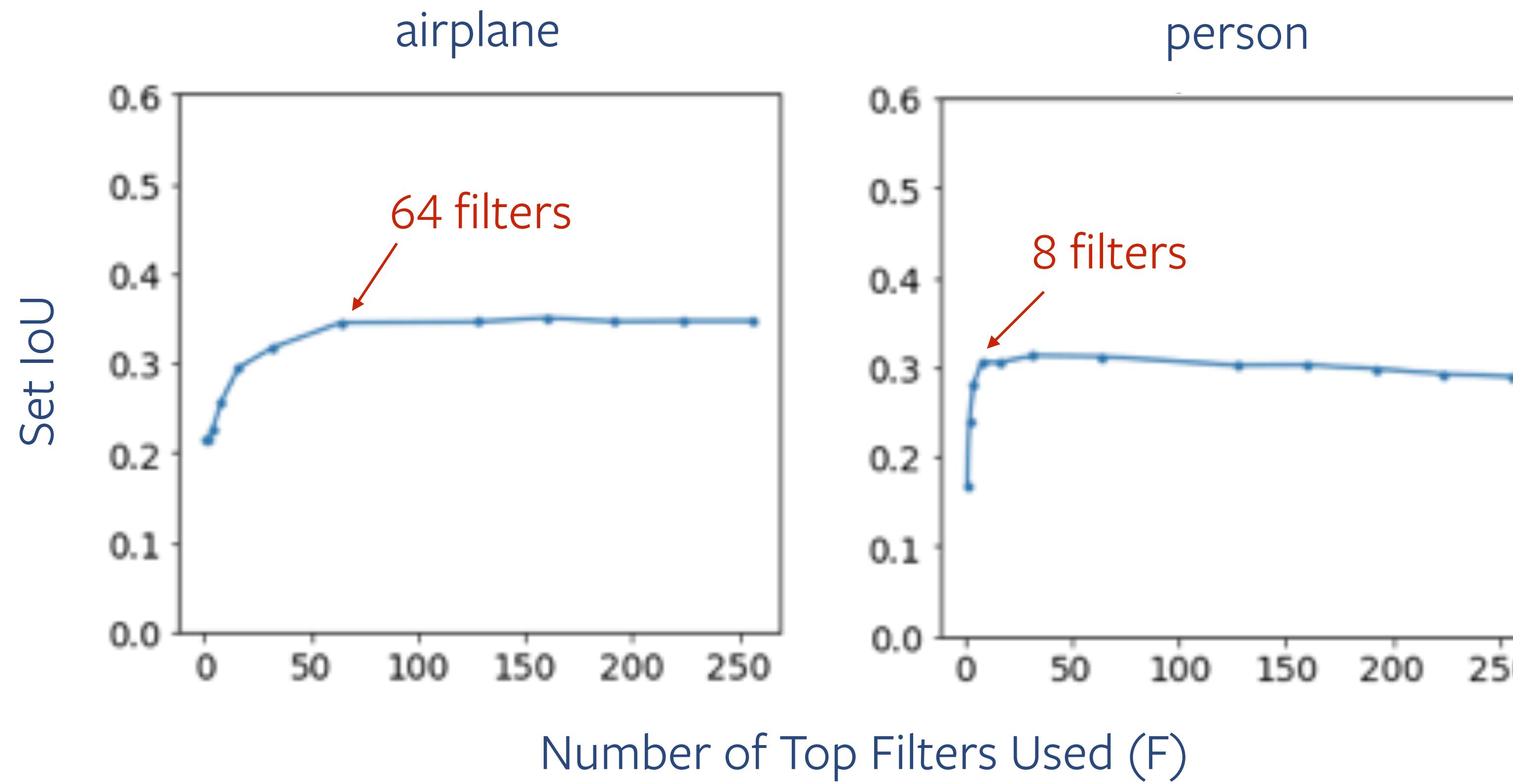
Horse
(IoU_{set} = .21)



Cow
(IoU_{set} = .20)



Filters per Concept



Different concepts require different number of filters for encoding.

Self-supervised learning

(



,

“sheepdog”

)

x

y

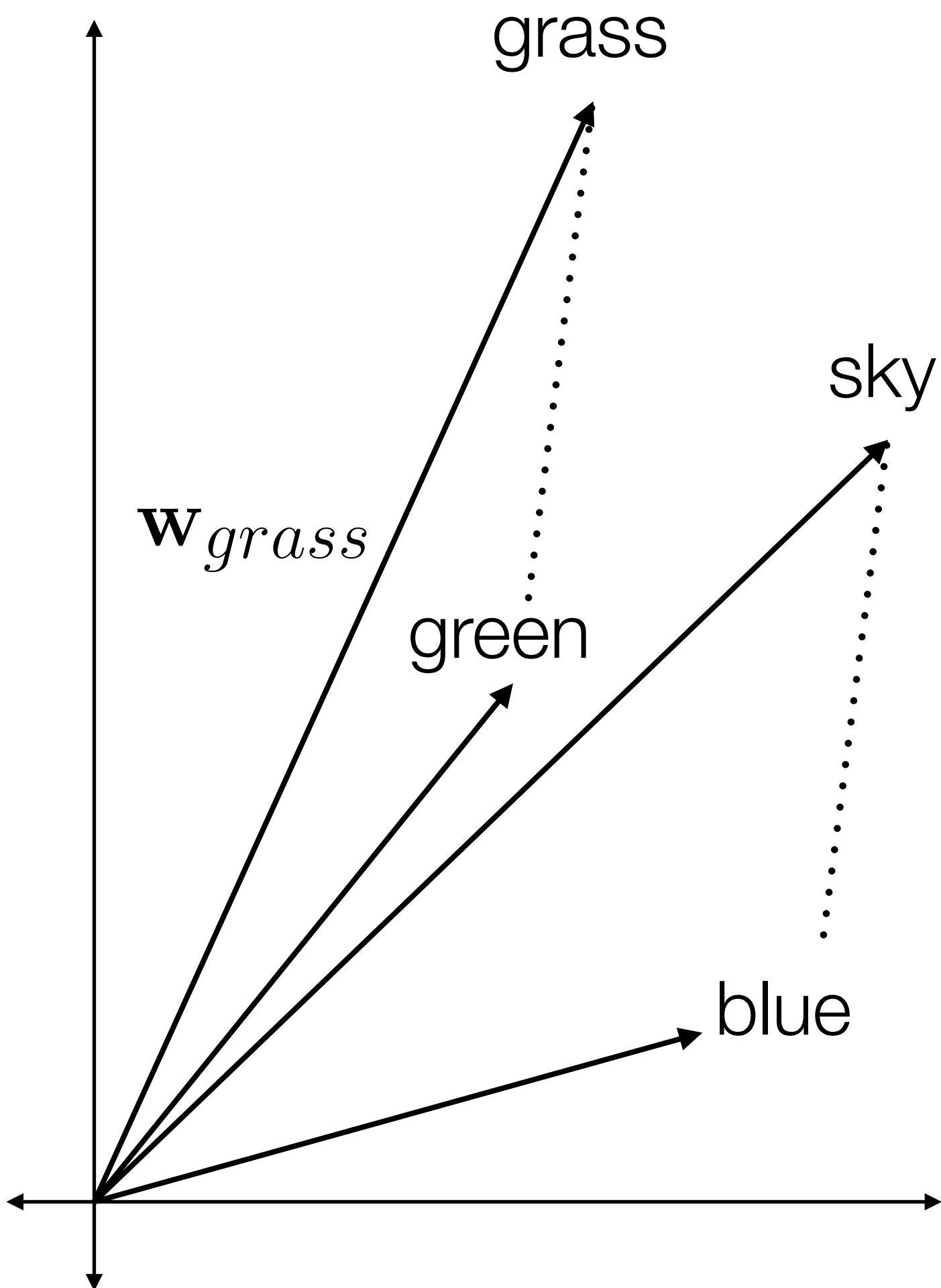
Filters: Supervised vs. Self-Supervised

Performance Improvement (Single Filter → All Filters):

- Self-supervised networks: 5-6x
- Fully-supervised networks: 2-4x

Self-supervised networks encode connects more distributively.

Comparing Concept Embeddings

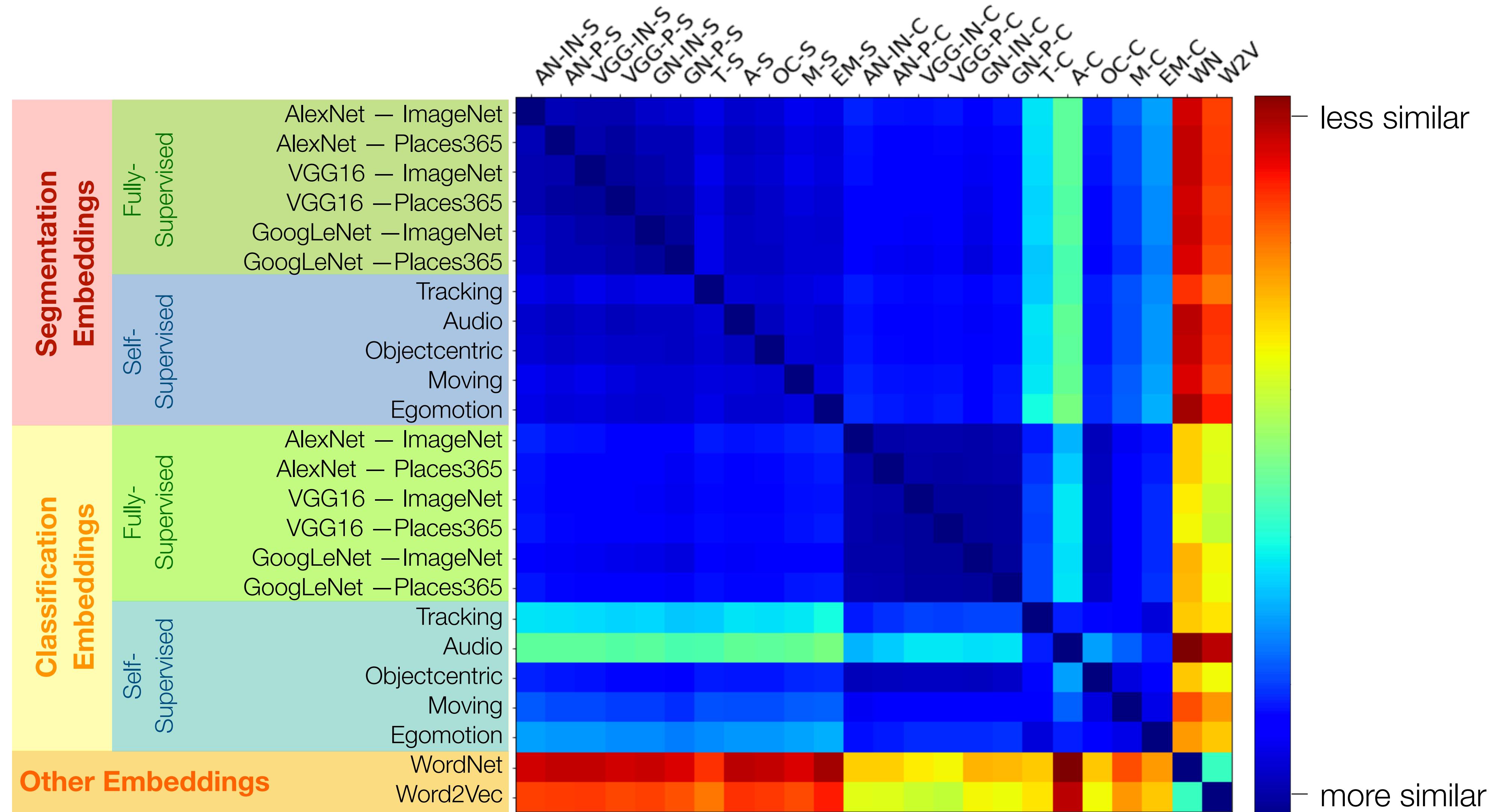


$$\text{grass} + \text{blue} - \text{green} = \text{sky}$$

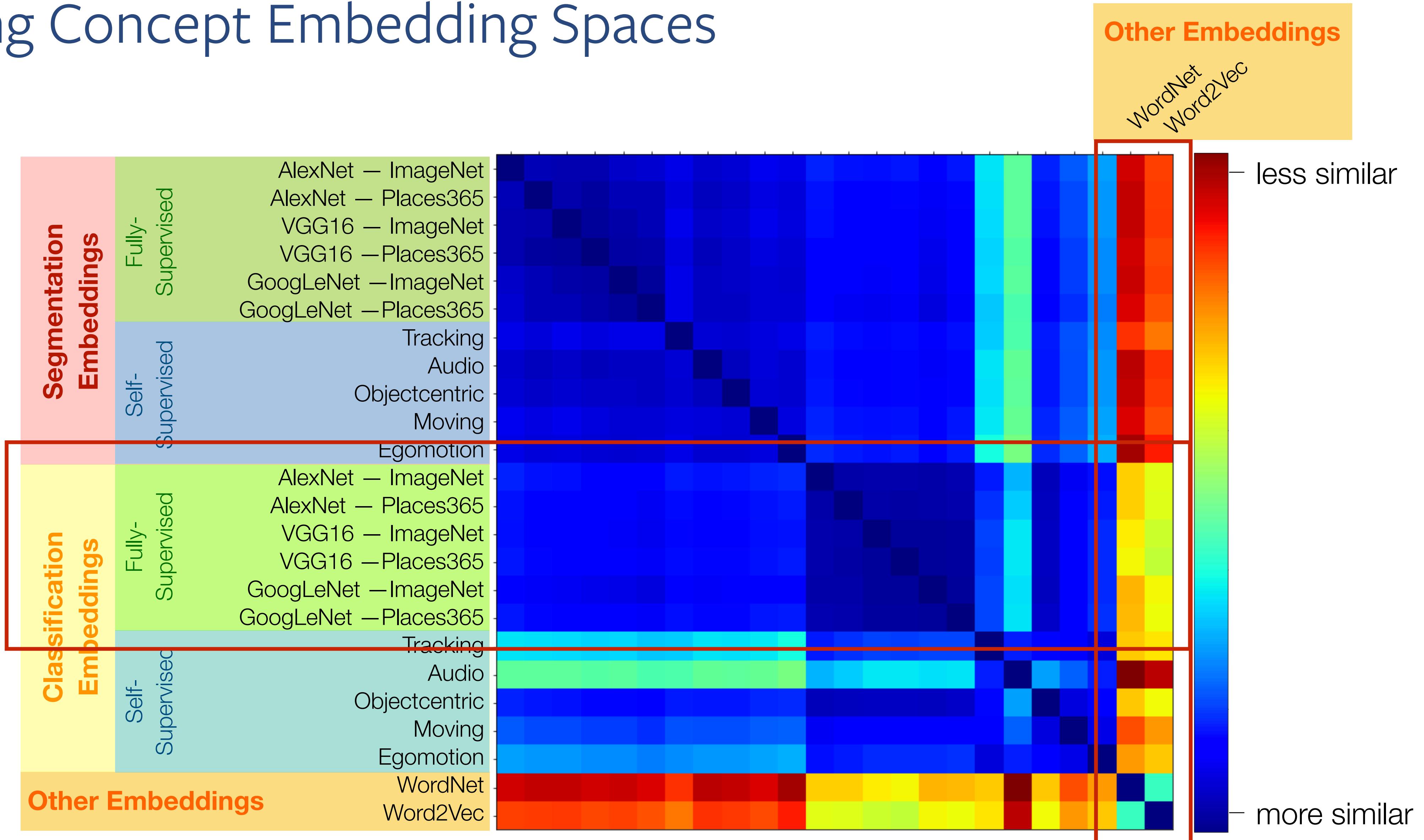
$$\text{tree} - \text{wood} = \text{plant}$$

$$\text{person} - \text{torso} = \text{foot}$$

Comparing Concept Embedding Spaces



Comparing Concept Embedding Spaces



Details

BRODEN dataset

Image-level Annotations

street (scene)

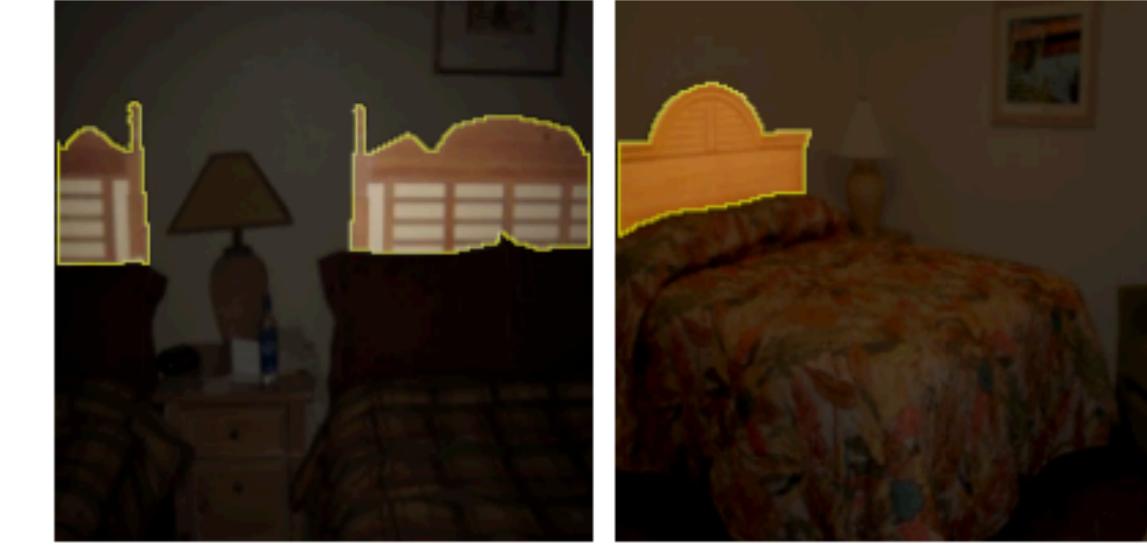


Pixel-level Annotations

flower (object)



headboard (part)



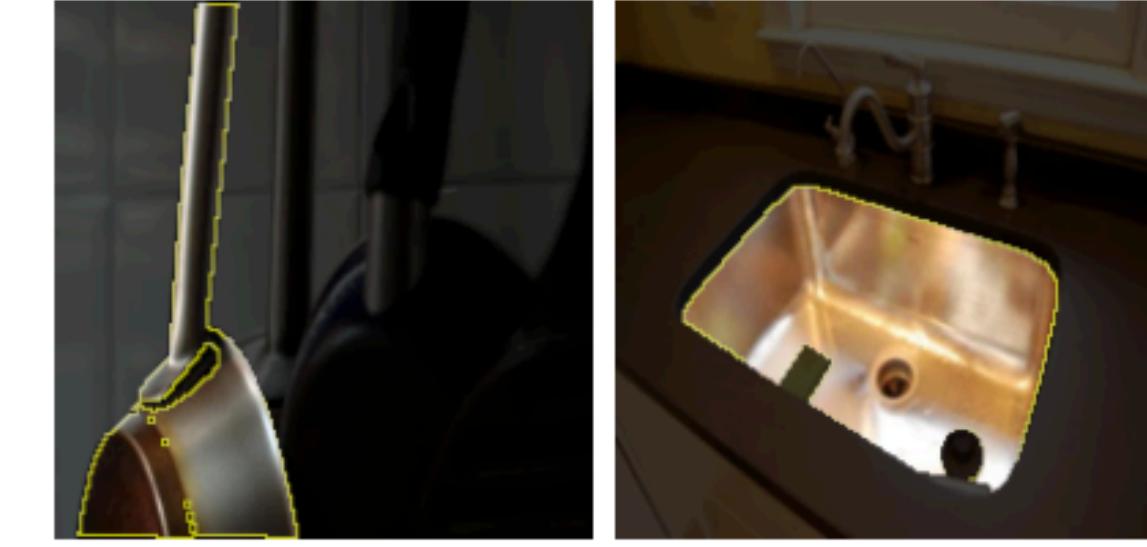
swirly (texture)



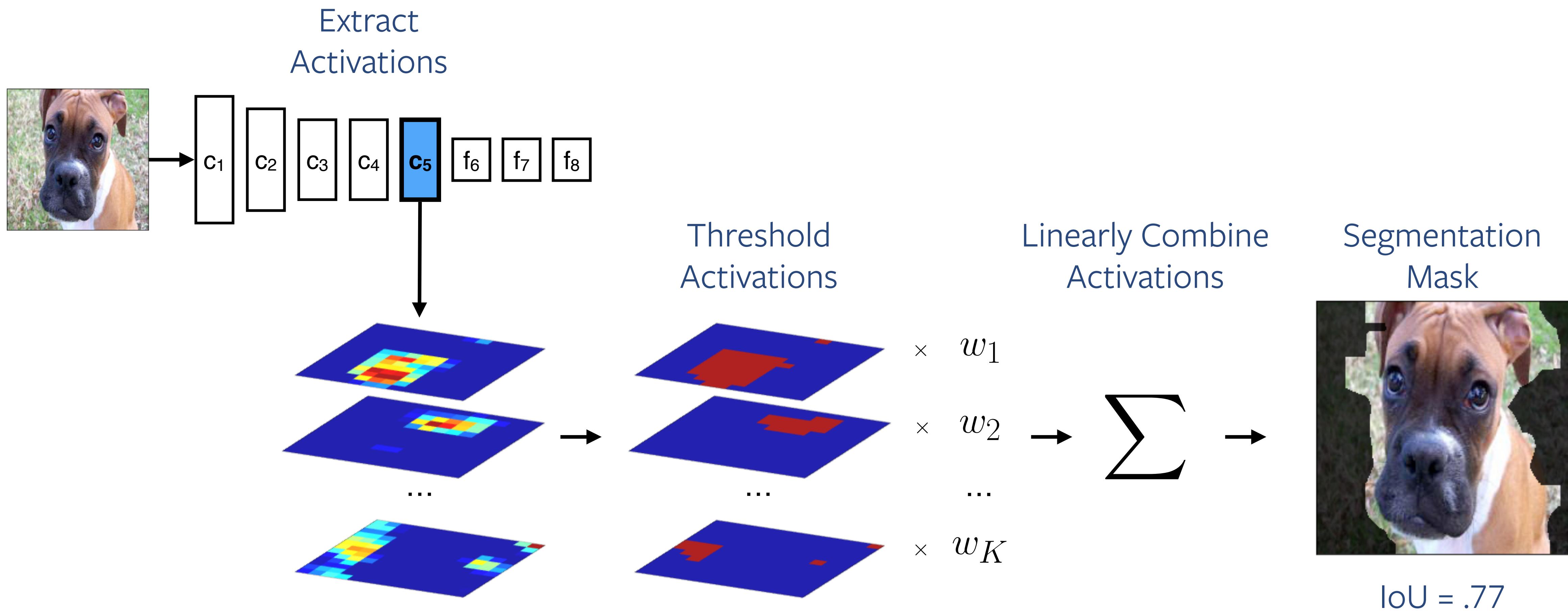
pink (color)



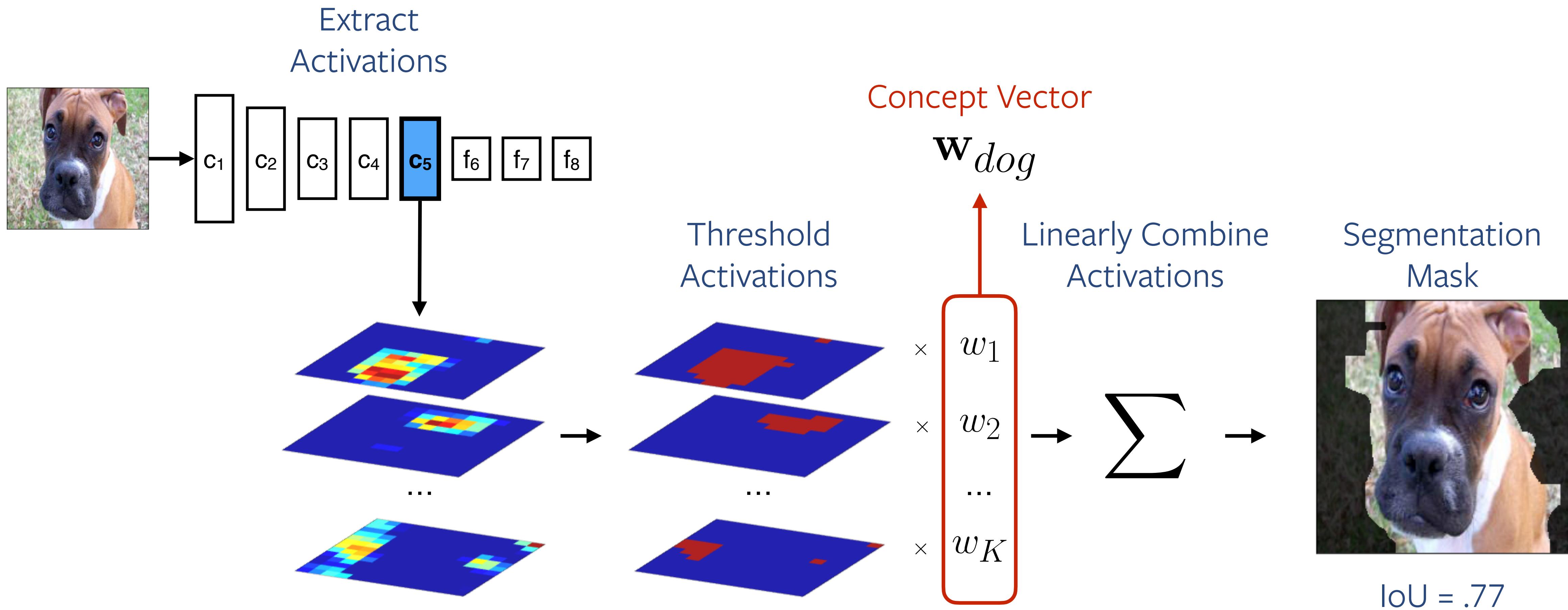
metal (material)



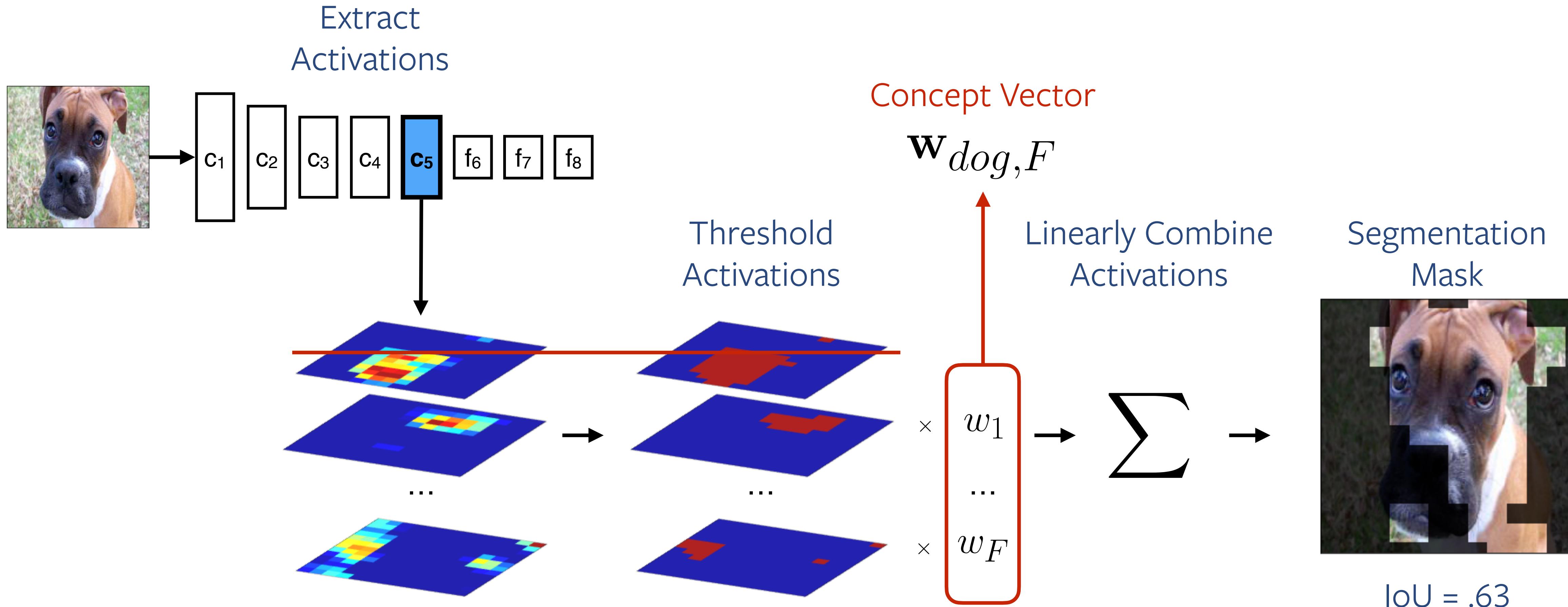
Segmentation



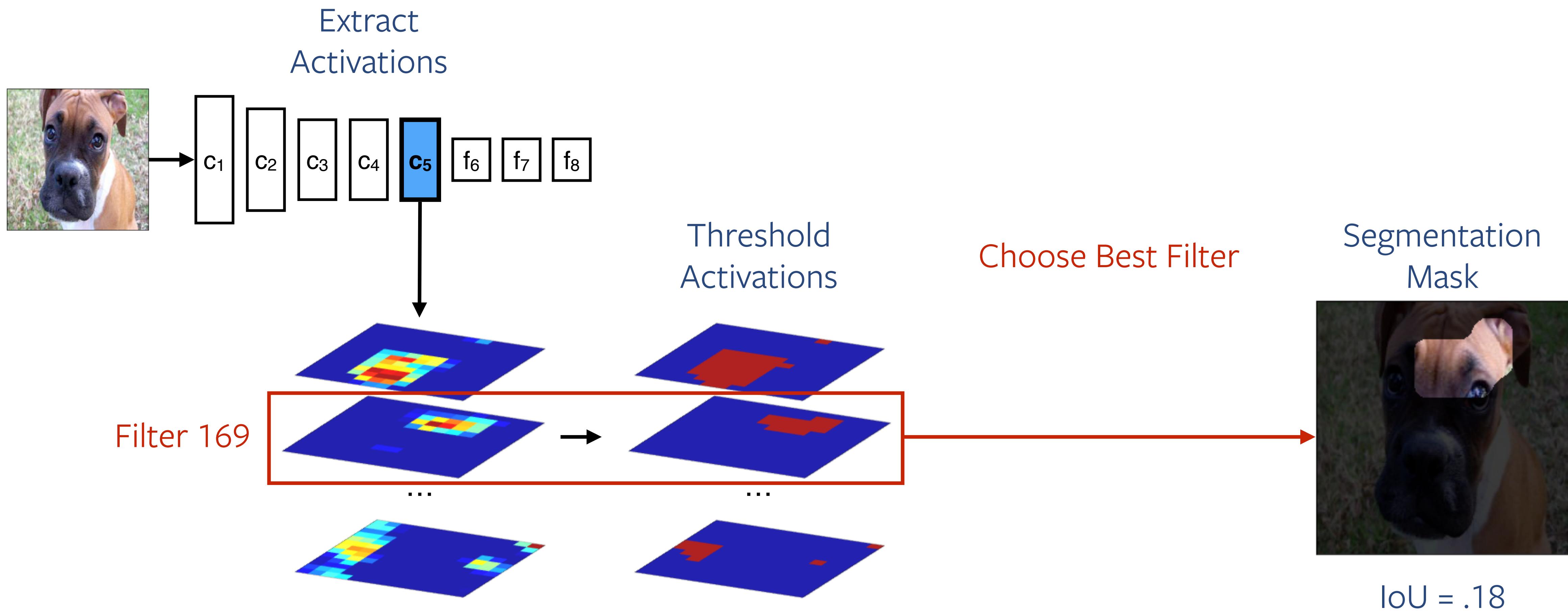
Segmentation



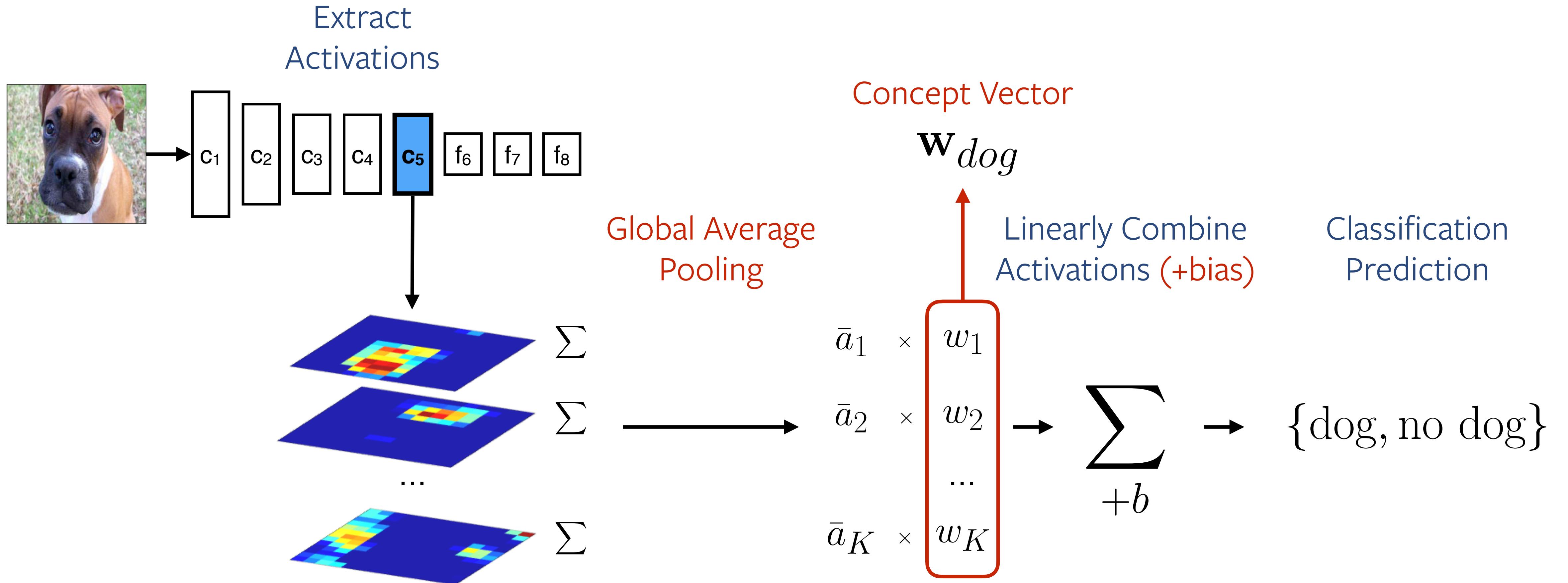
Segmentation



Segmentation

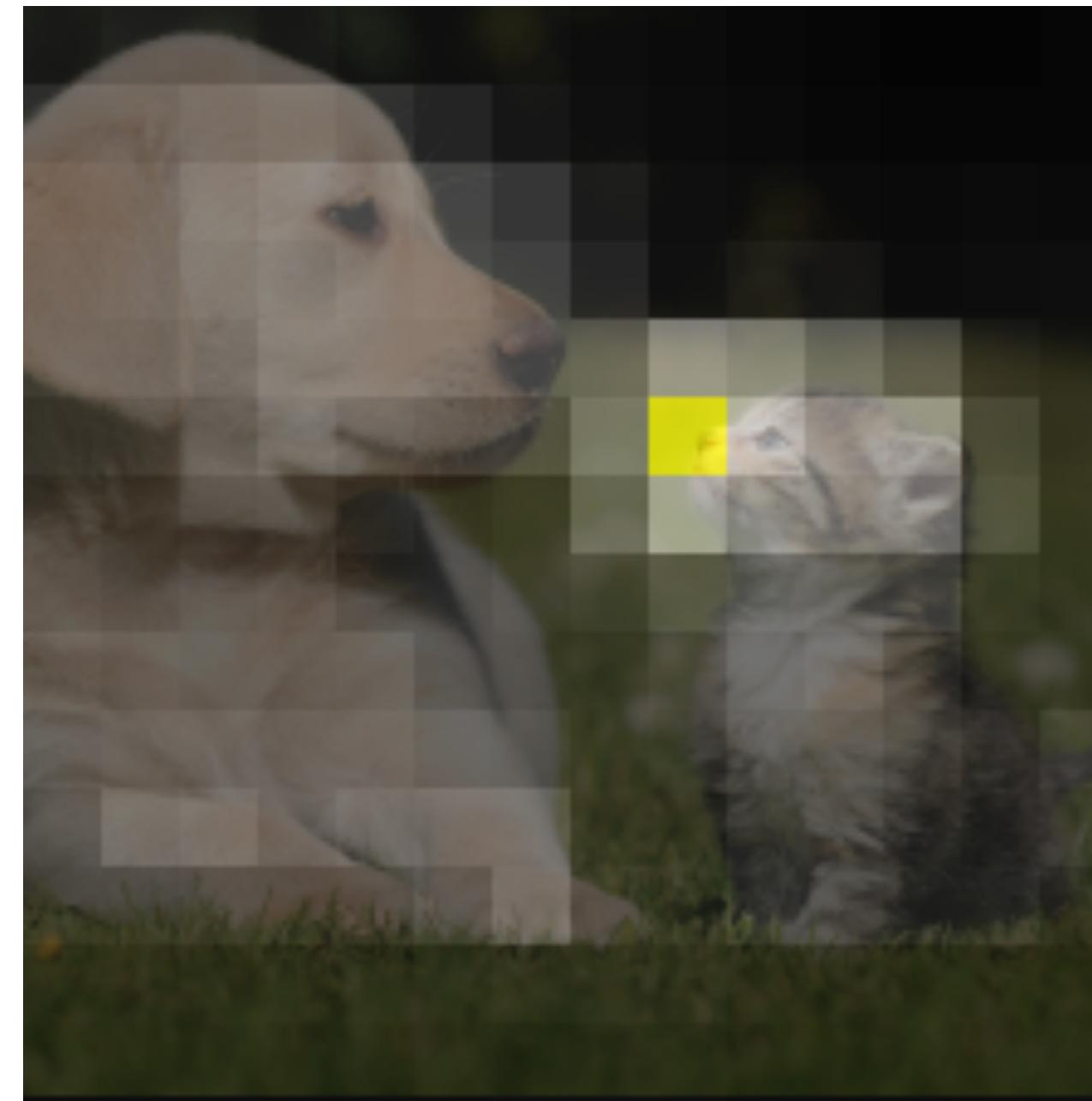


Classification



C. Exploring activations via interactive visualizations

Preview: Interactive Similarity Overlays



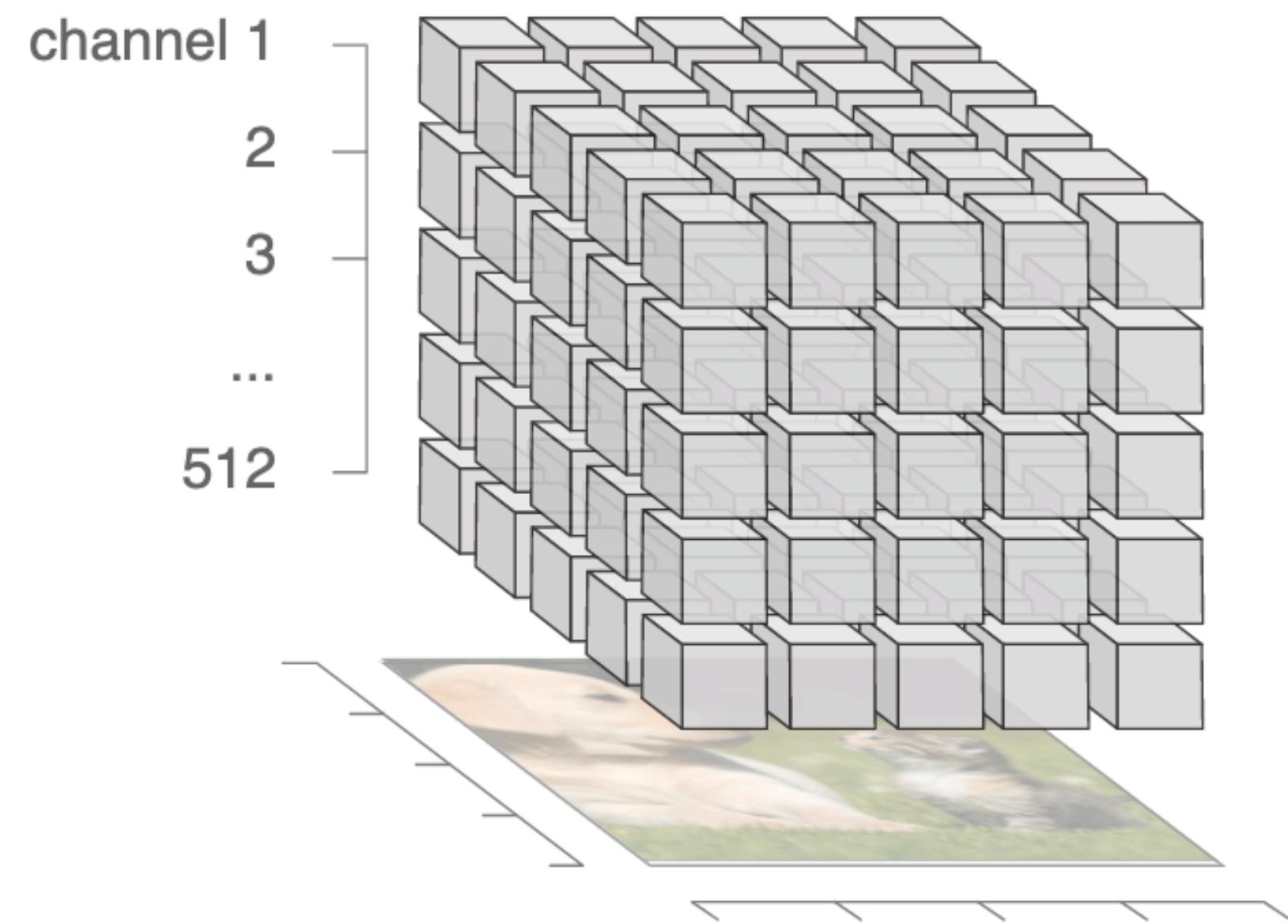
Interactive visualizations empower practitioners to easily understand model behavior.

[Fong et al., 2020 (in prep.)] ⁷³

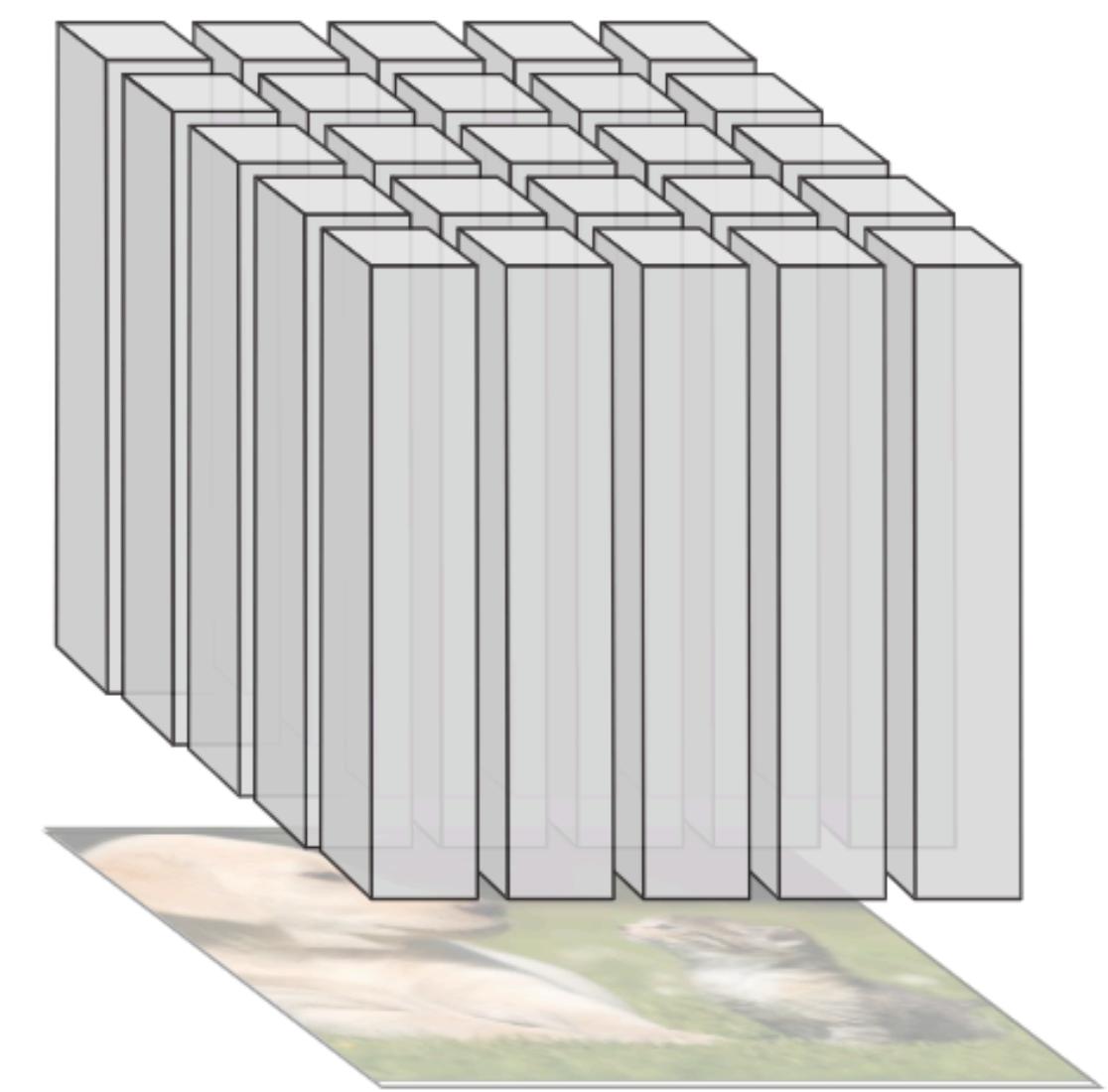
Live Preview

Intermediate Activations

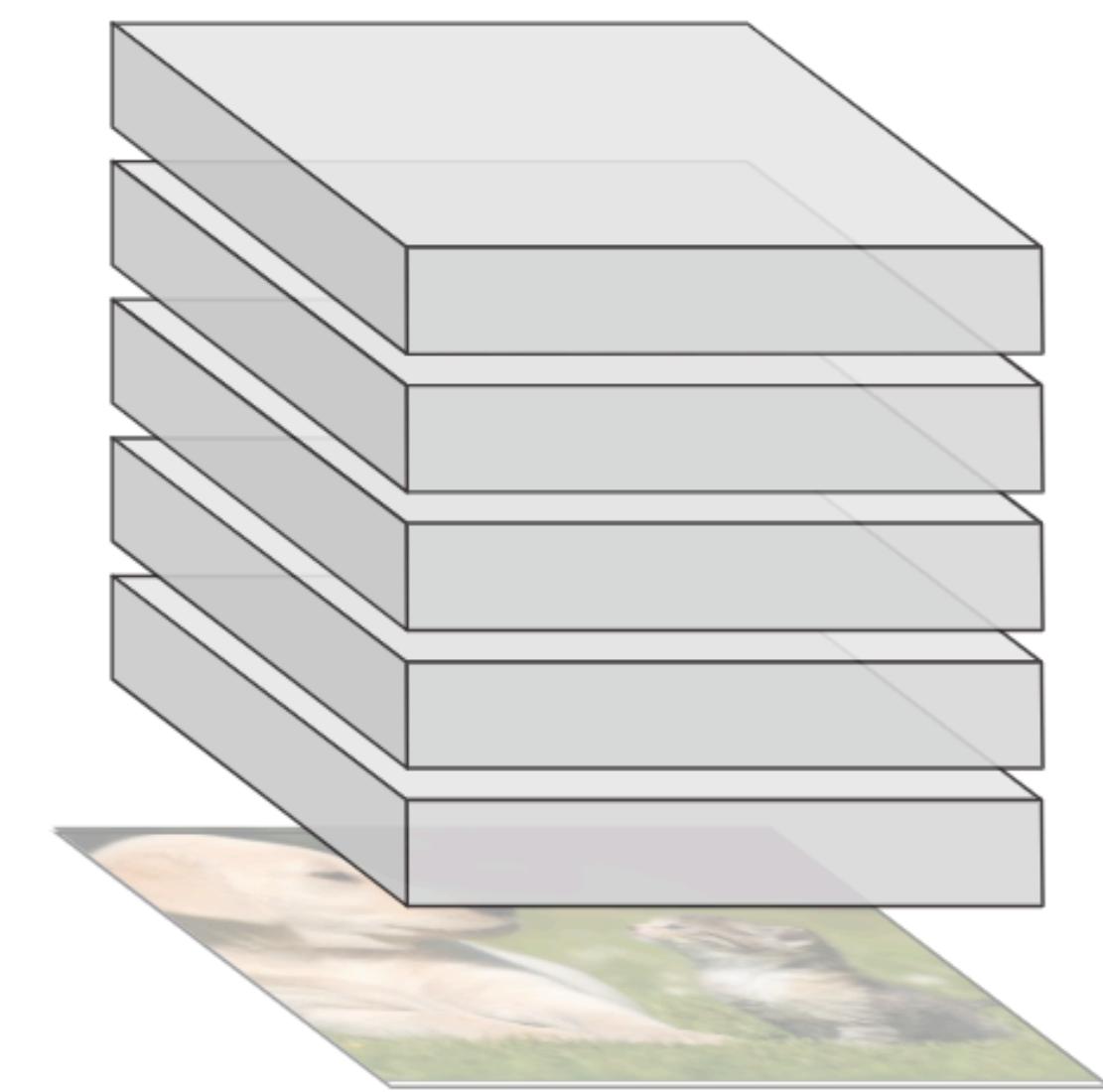
Individual Neurons



Spatial Activations



Channel Activations

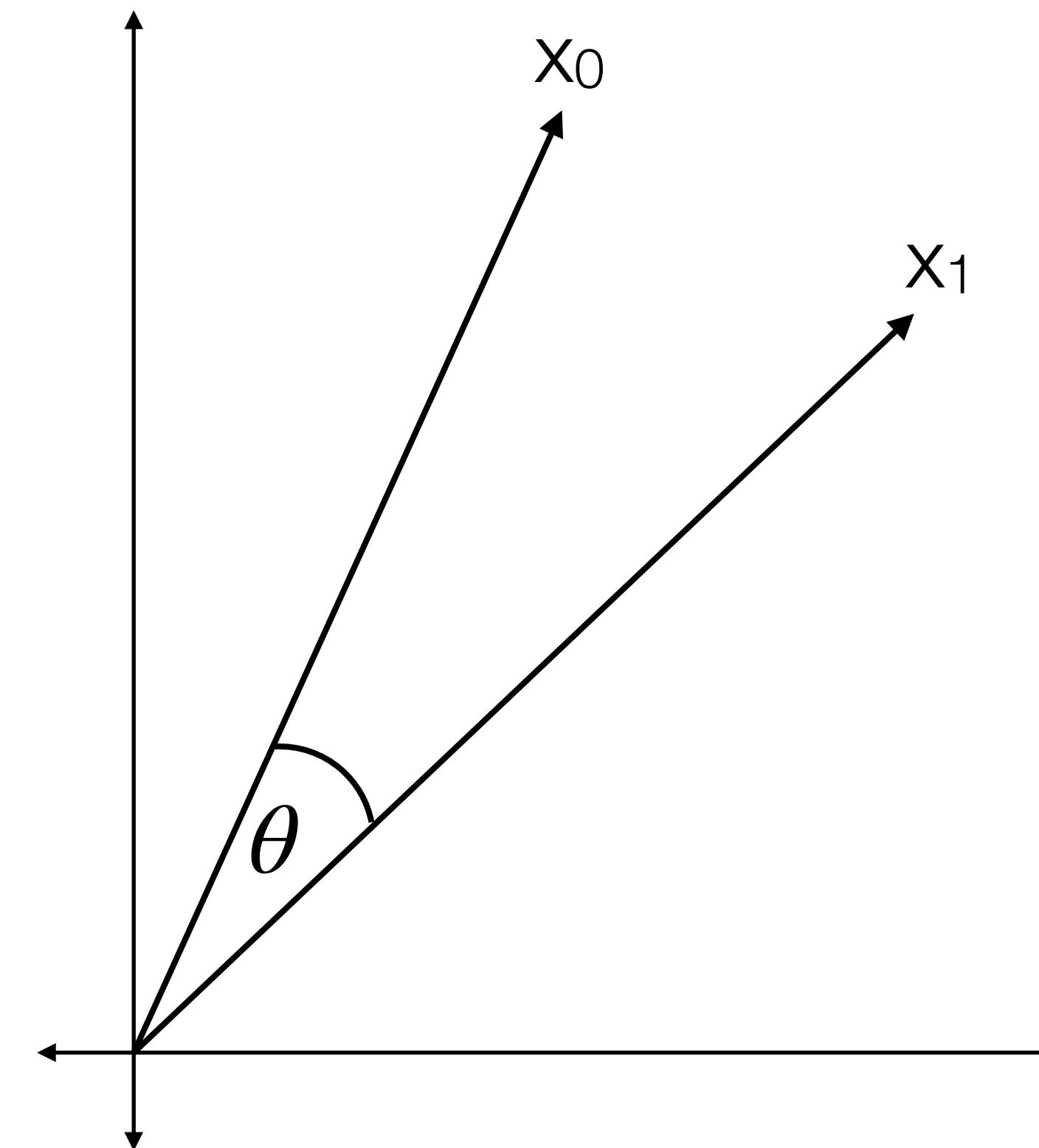
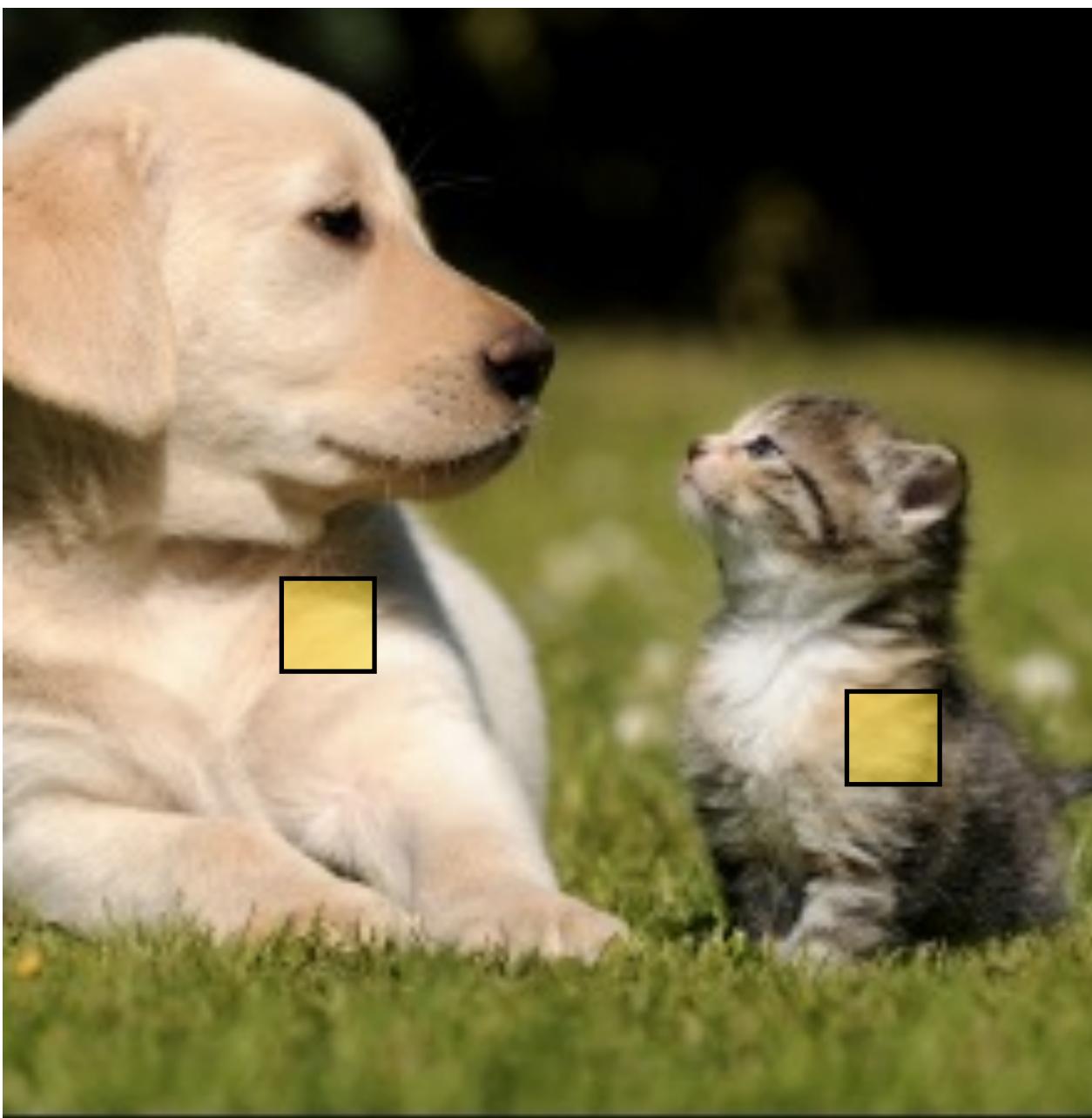


Preview: Interactive Similarity Overlays

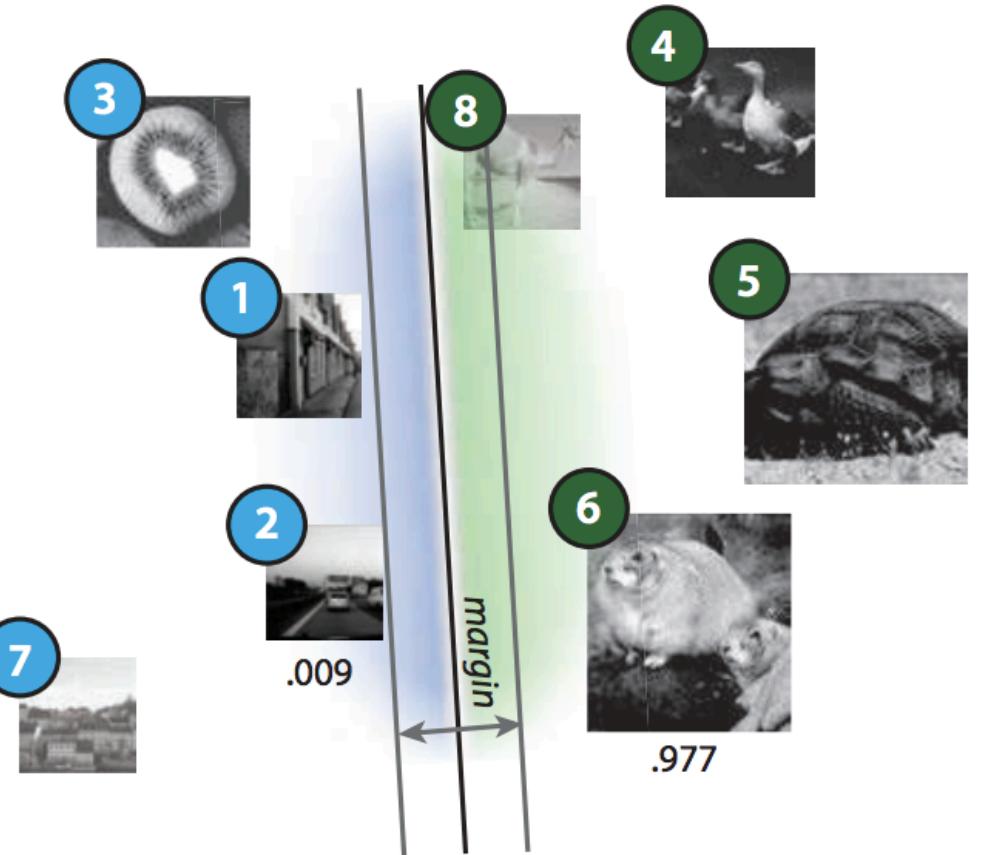
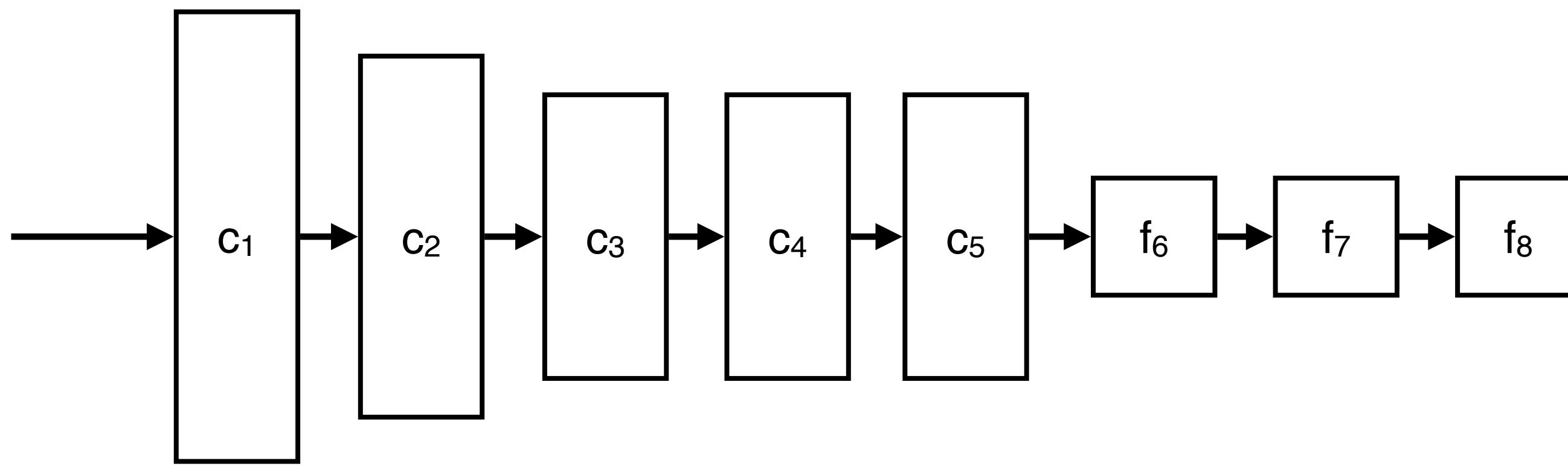


$a_{6,5} = [17.7, 0, 103.4, 6.81, 0, 0, 0, 0, 32.0, 0, 0, 0, \dots]$

Preview: Interactive Similarity Overlays



Research Themes



Fong et al., Sci. Reports 2018

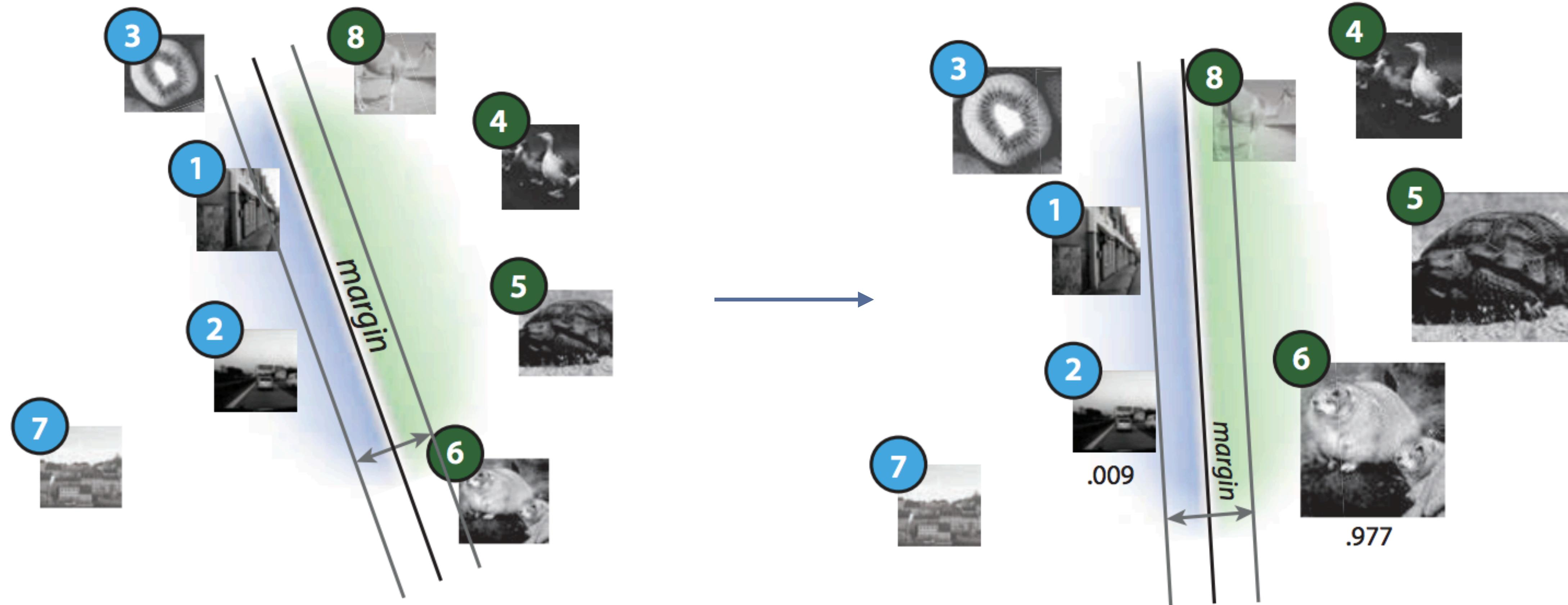
Training Procedure: How can we improve $f(x)$?

$$\theta \leftarrow \theta - \eta \frac{dL}{d\theta}$$



Fong & Vedaldi, ICCVW 2019

Human-Guided Machine Learning

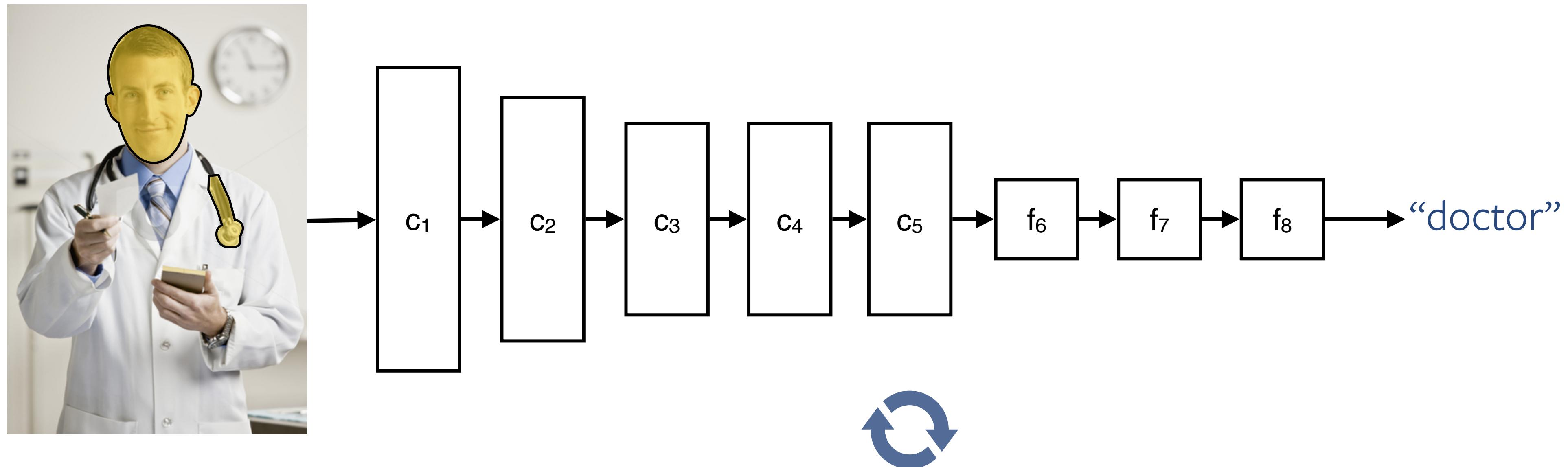


Align machine decisions with human decisions from brain activity.

[Fong et al., Sci. Reports 2018] ⁷⁹

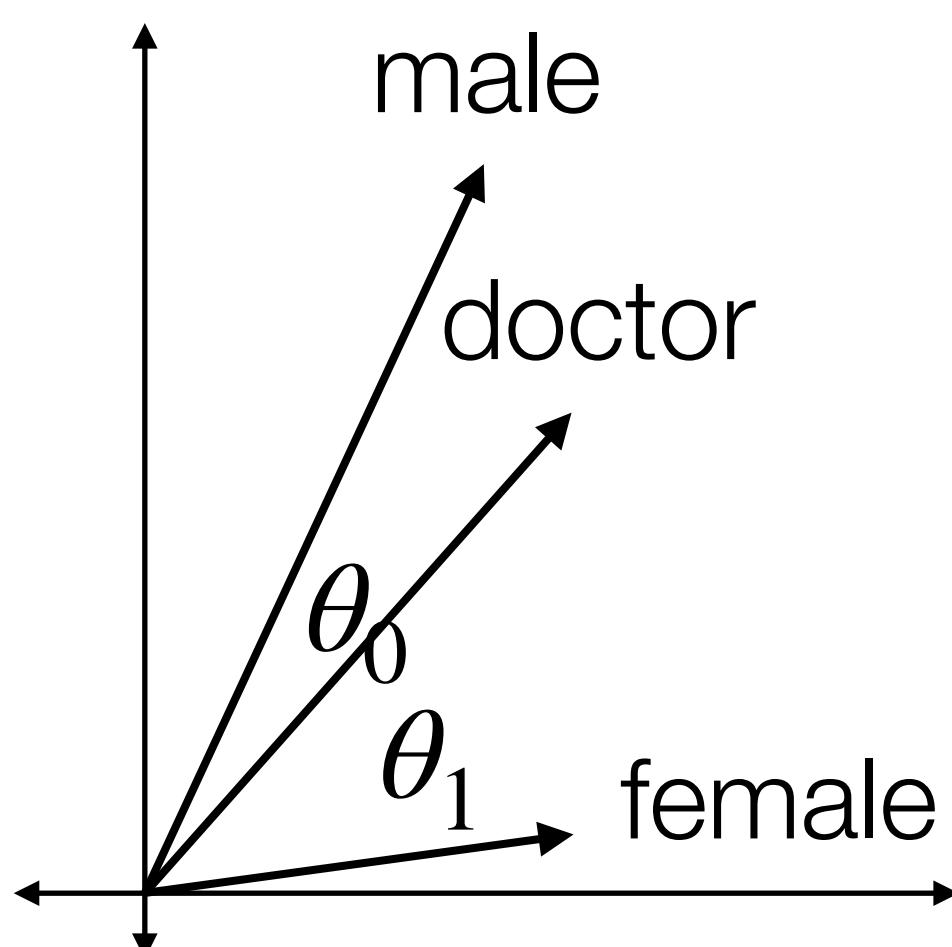
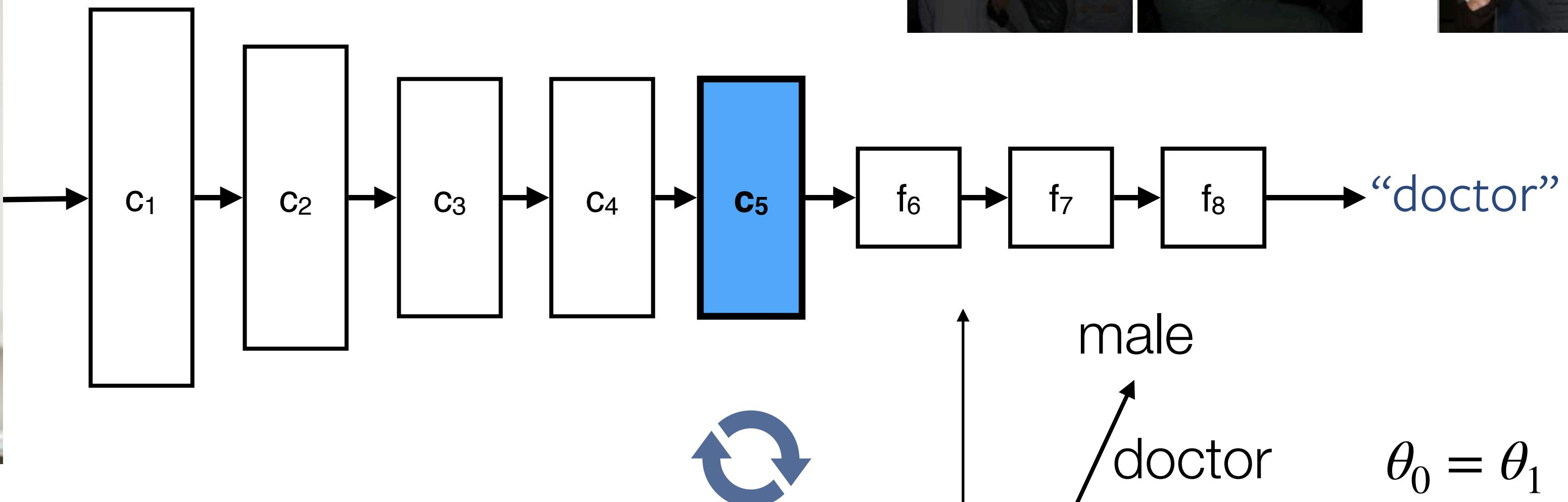
Future Work: Model Debugging

Identify and correct systematic mistakes



Future Work: Model Debugging

Identify and correct systematic mistakes



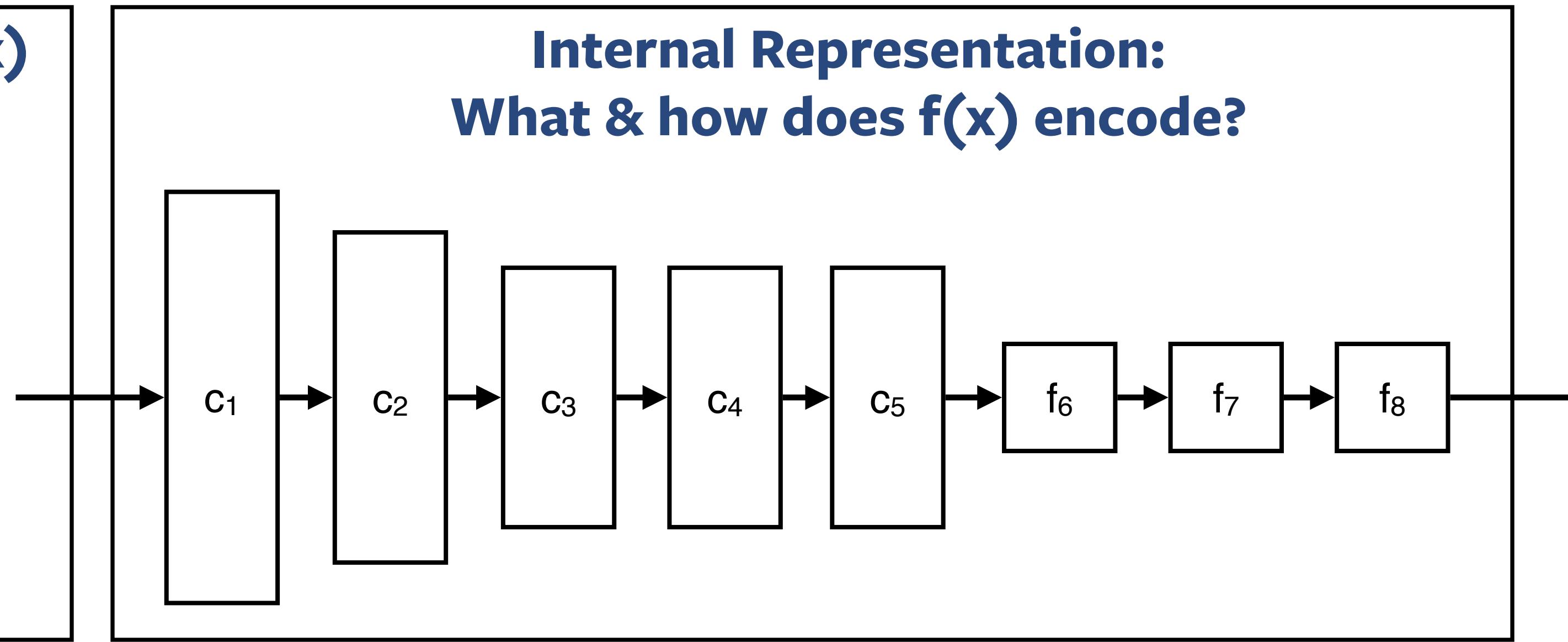
[Bolukbasi et al., NeurIPS 2016;
David et al., (in prep)]

Research Themes

Inputs: What is $f(x)$ looking at?



**Internal Representation:
What & how does $f(x)$ encode?**



95%: “sheepdog”
1%: “sea snake”
2%: “alp”
0%: “soup bowl”
...

Training Procedure: How can we improve $f(x)$?



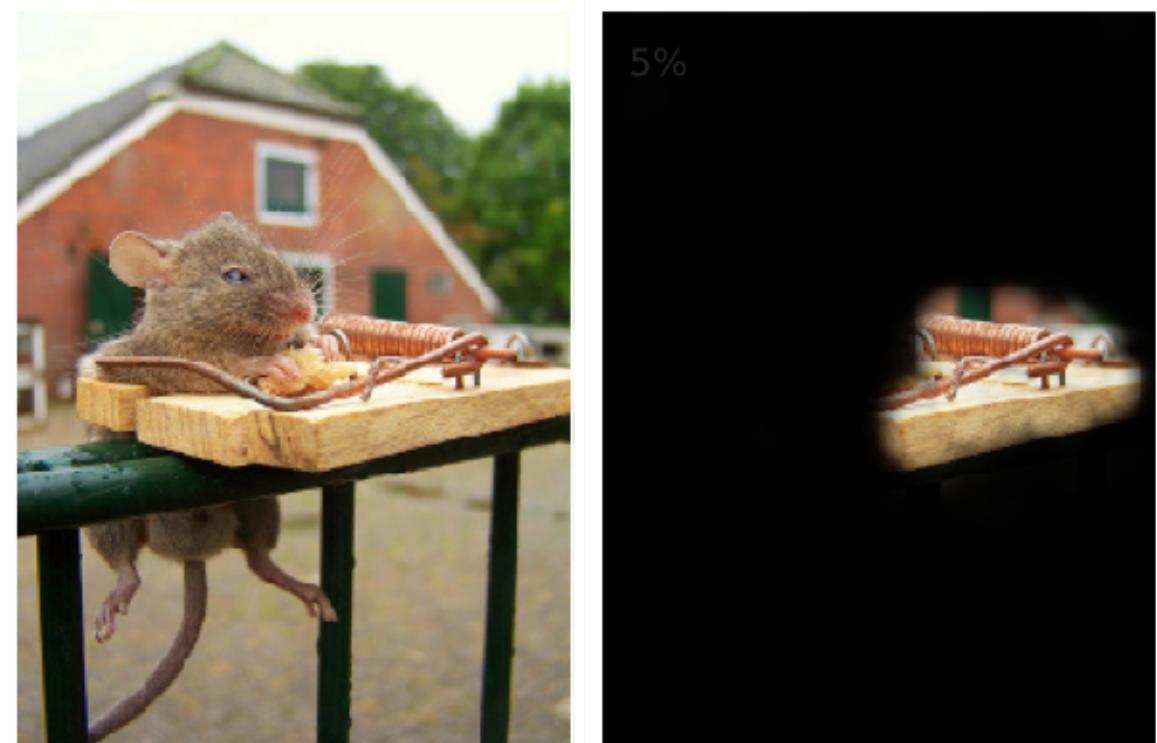
$$\theta \leftarrow \theta - \eta \frac{dL}{d\theta}$$

Research Themes

What is $f(x)$ looking at?

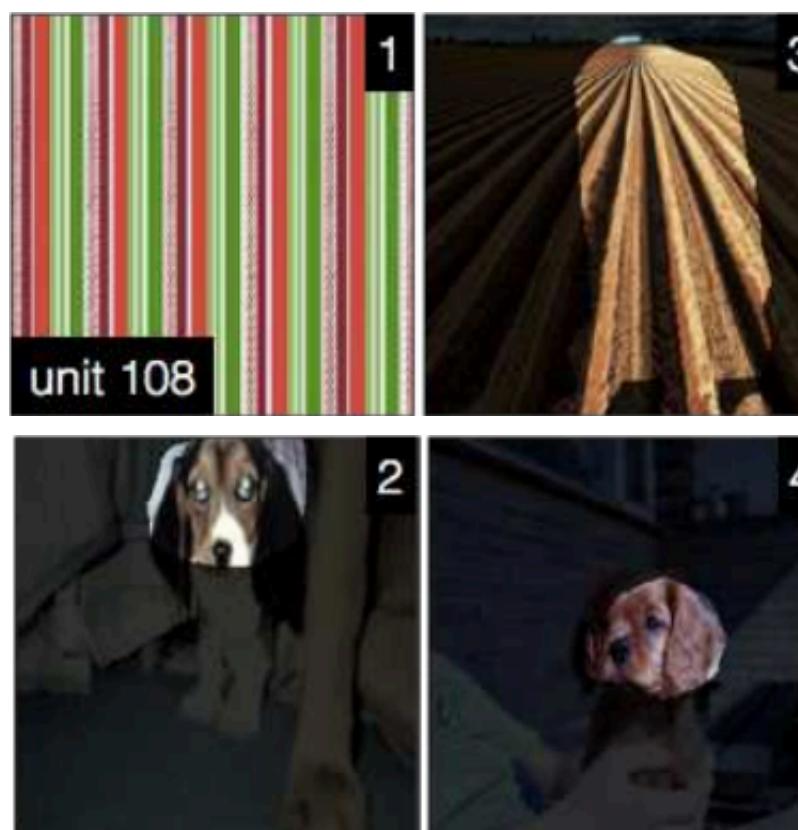


Fong & Vedaldi, ICCV 2017

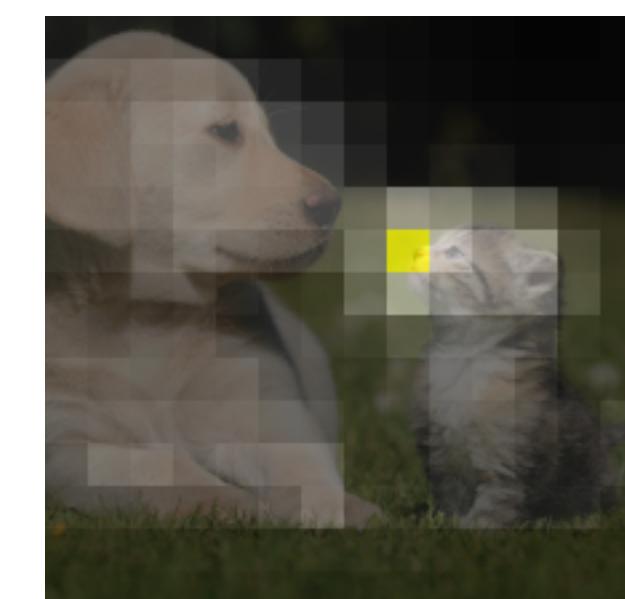


Fong et al., ICCV 2019

What & how does $f(x)$ encode?

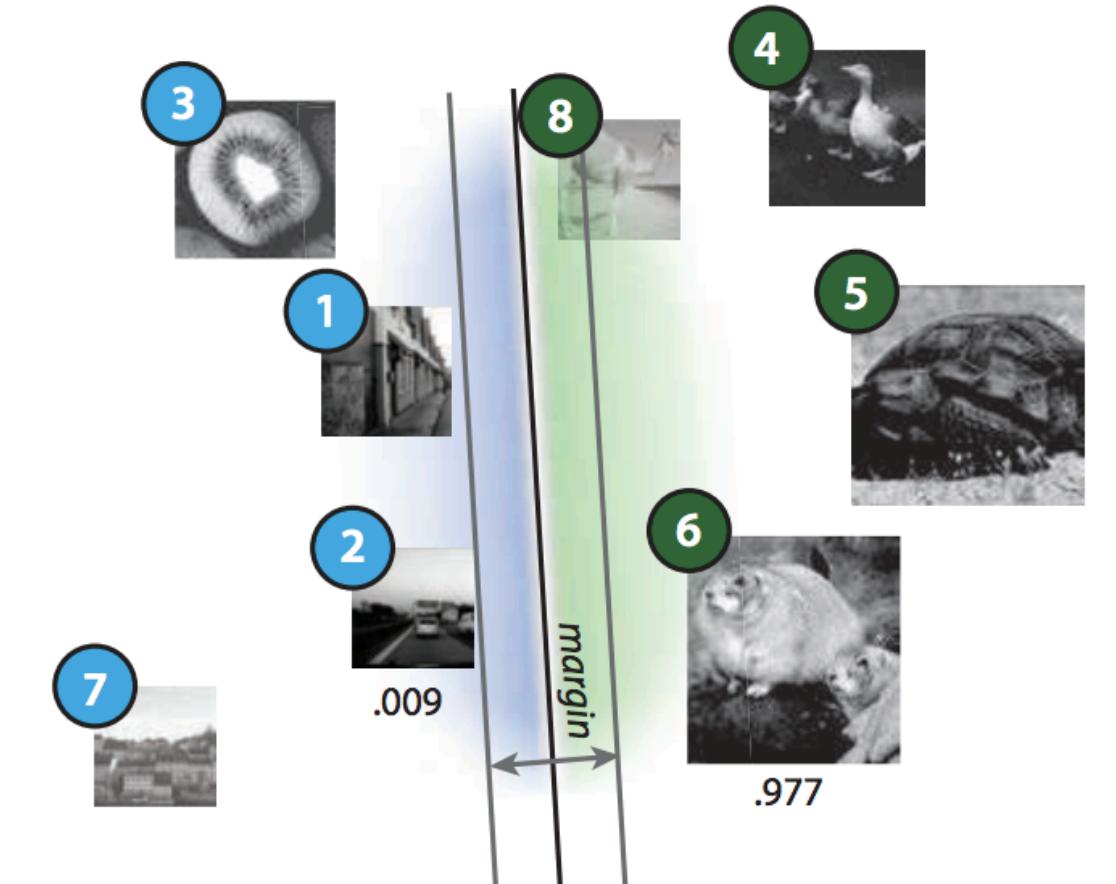


Fong & Vedaldi, CVPR 2018

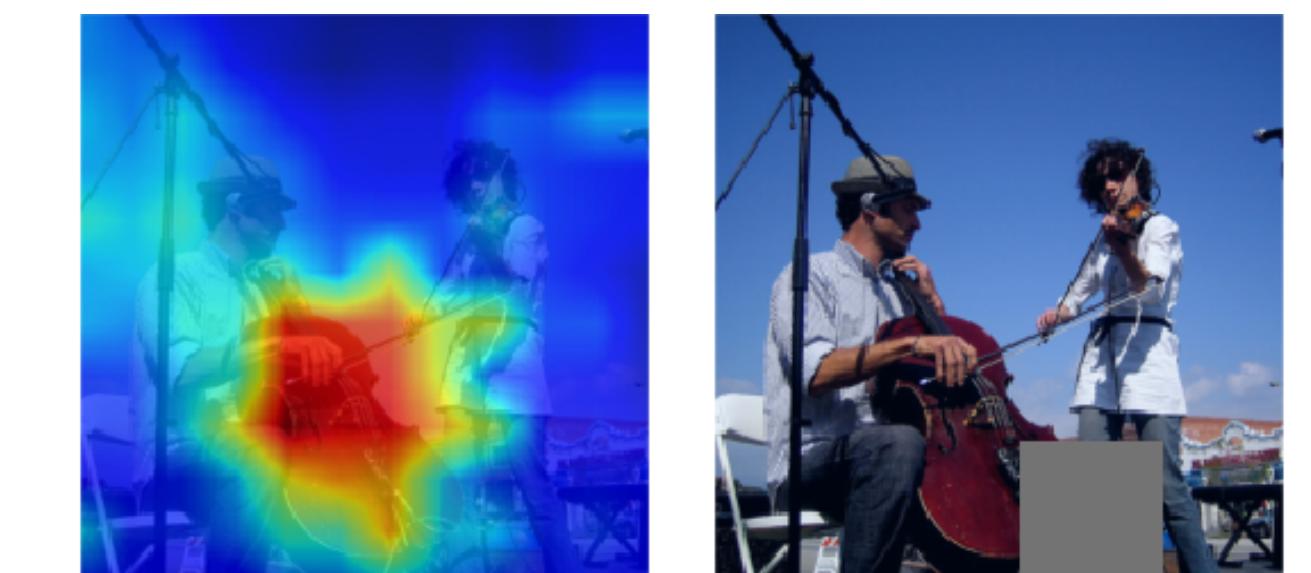


Fong et al., 2020 (in prep.)

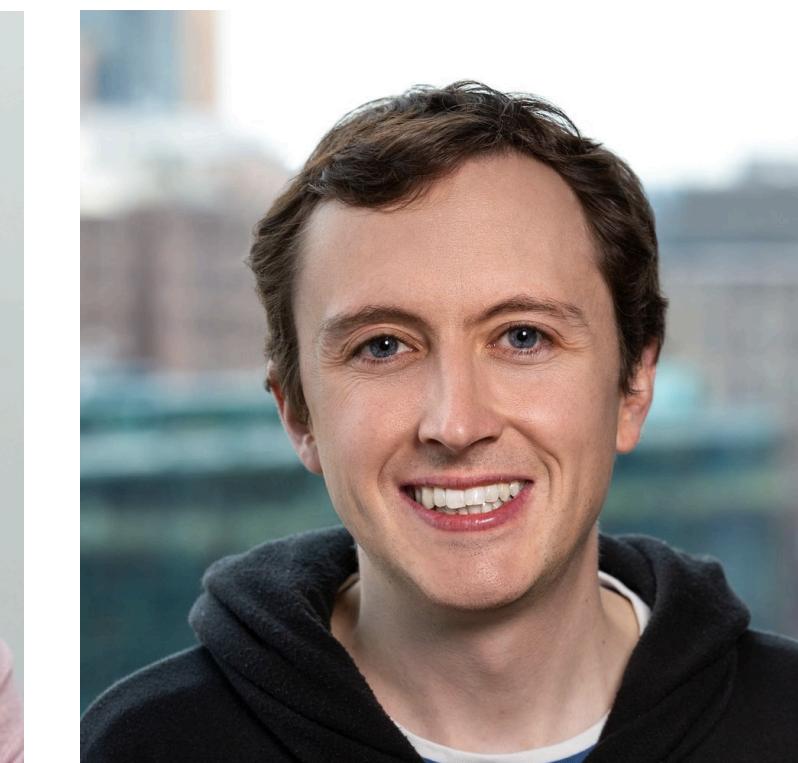
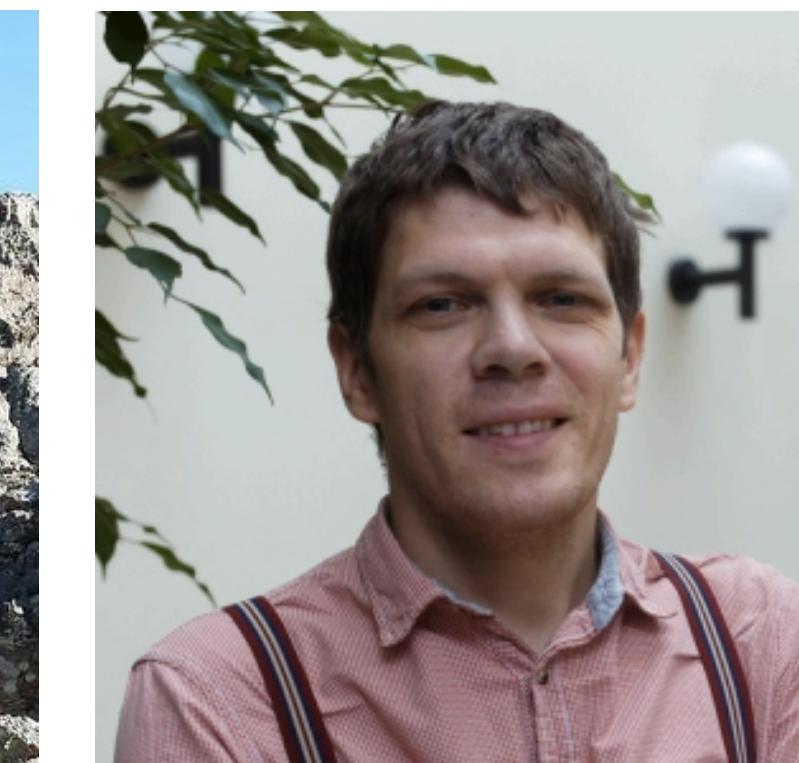
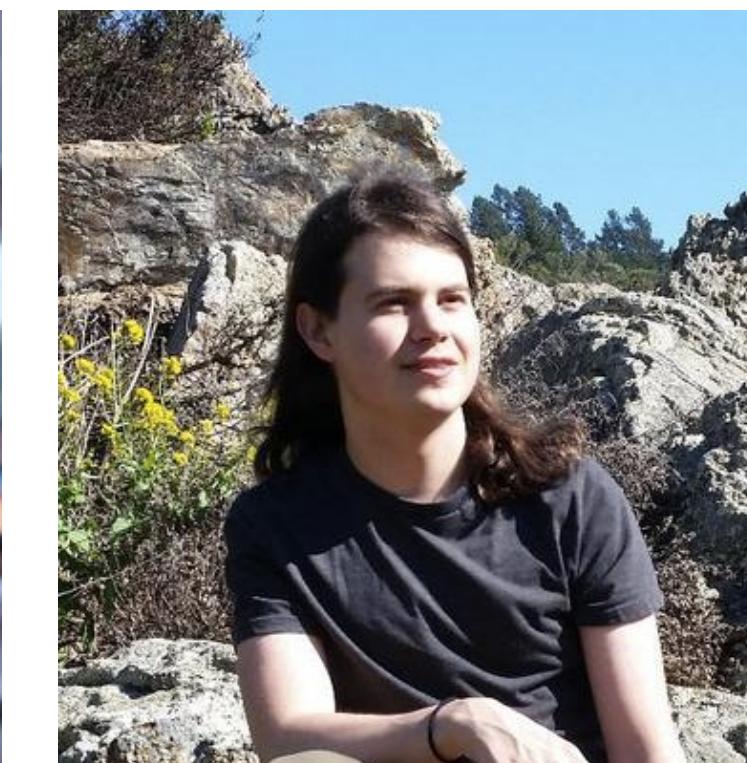
How can we improve $f(x)$?



Fong et al., Sci. Reports 2018



Fong & Vedaldi, ICCVW 2019



Andrea Vedaldi

Mandela Patrick

Chris Olah

Alexander
Mordvintsev

Walter Scheirer

David Cox



Thank you