

Net2Vec: Quantifying and Explaining how Concepts are Encoded by Filters in Deep Neural Networks

Ruth Fong and Andrea Vedaldi



UNIVERSITY OF
OXFORD

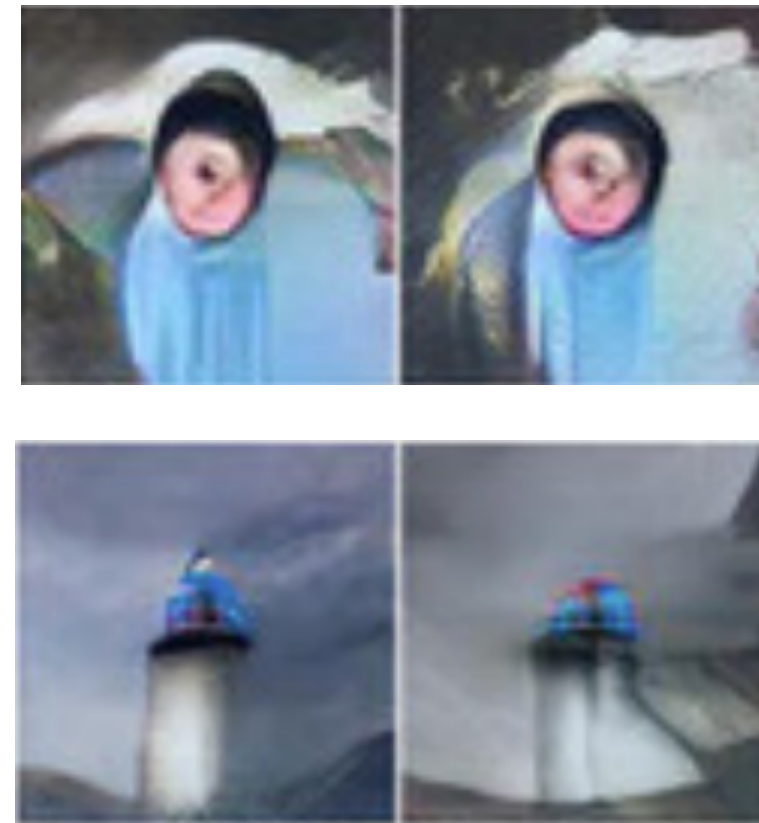


Filter Visualizations

Zeiler & Fergus, ECCV 2014



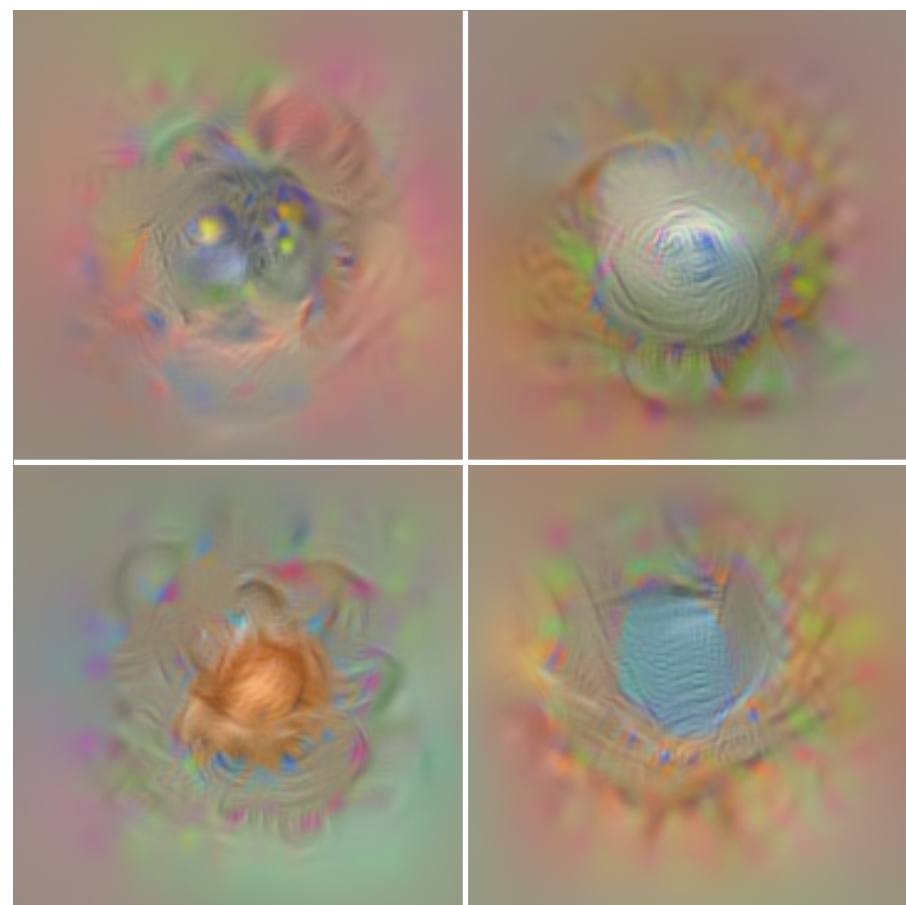
Nyugen et al., NIPS 2016



Zhou et al., ICLR 2015



Mahendran and Vedaldi,
IJCV 2016



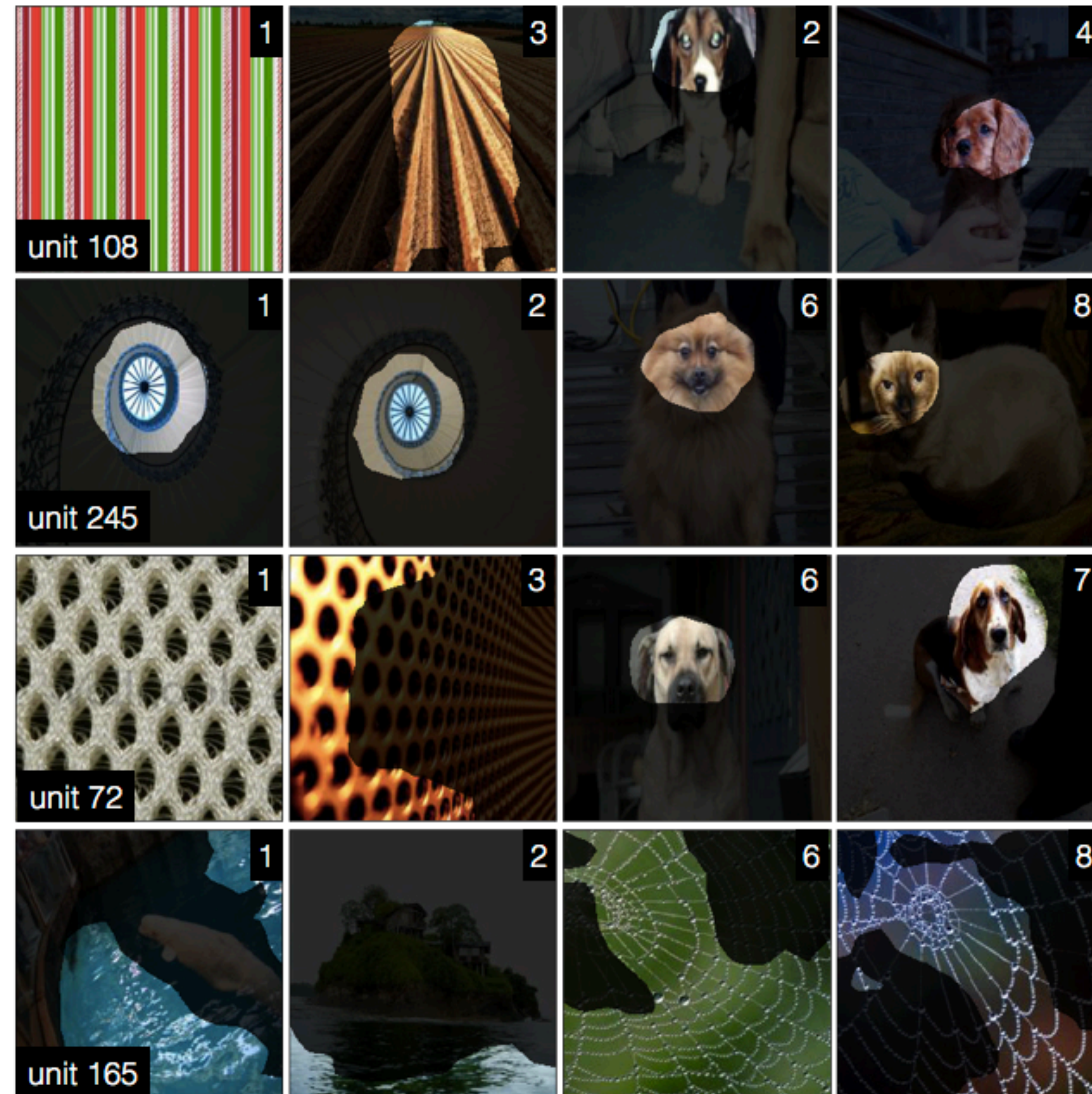
Olah et al., Distill 2017



Bau et al., CVPR 2017

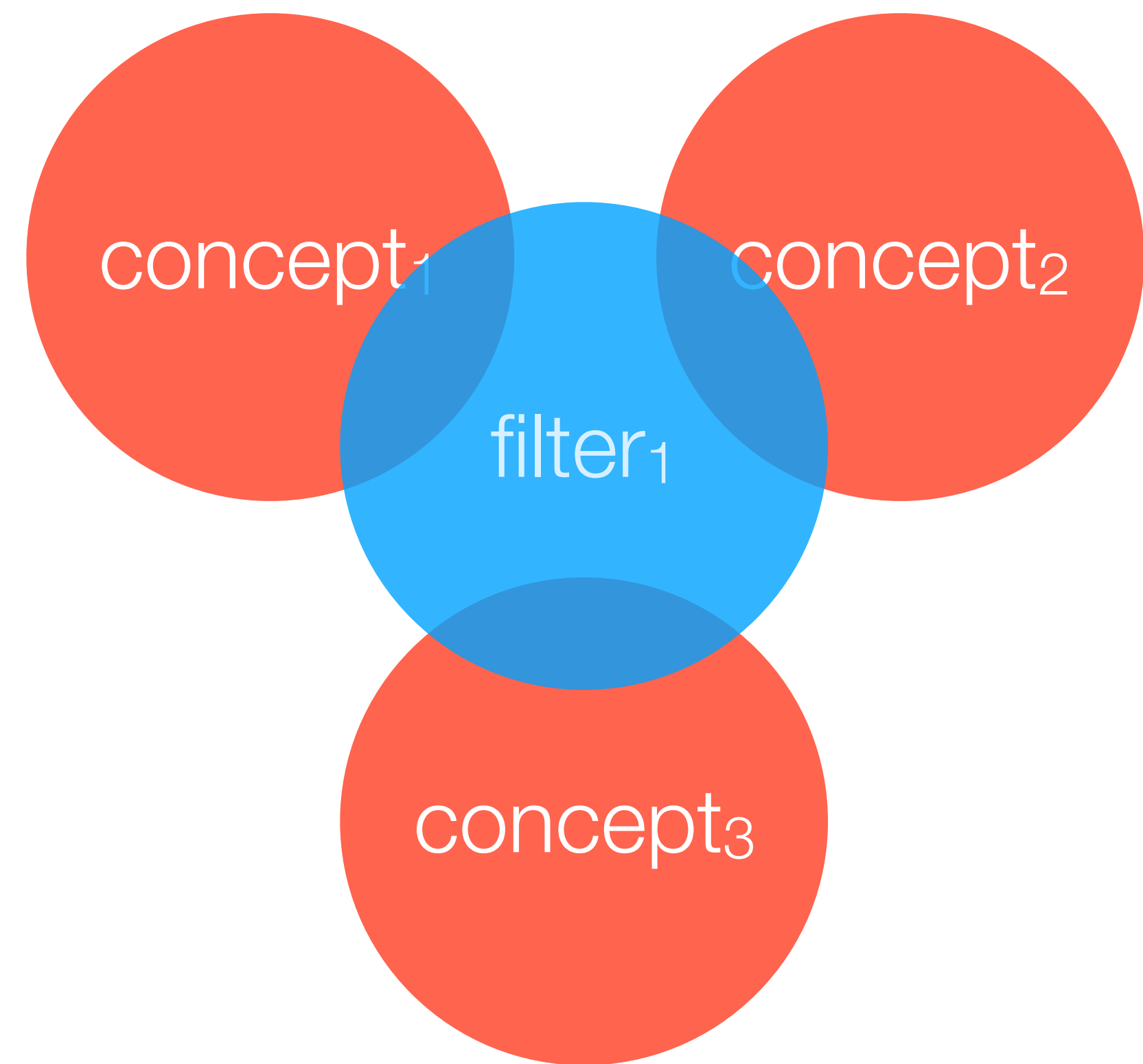
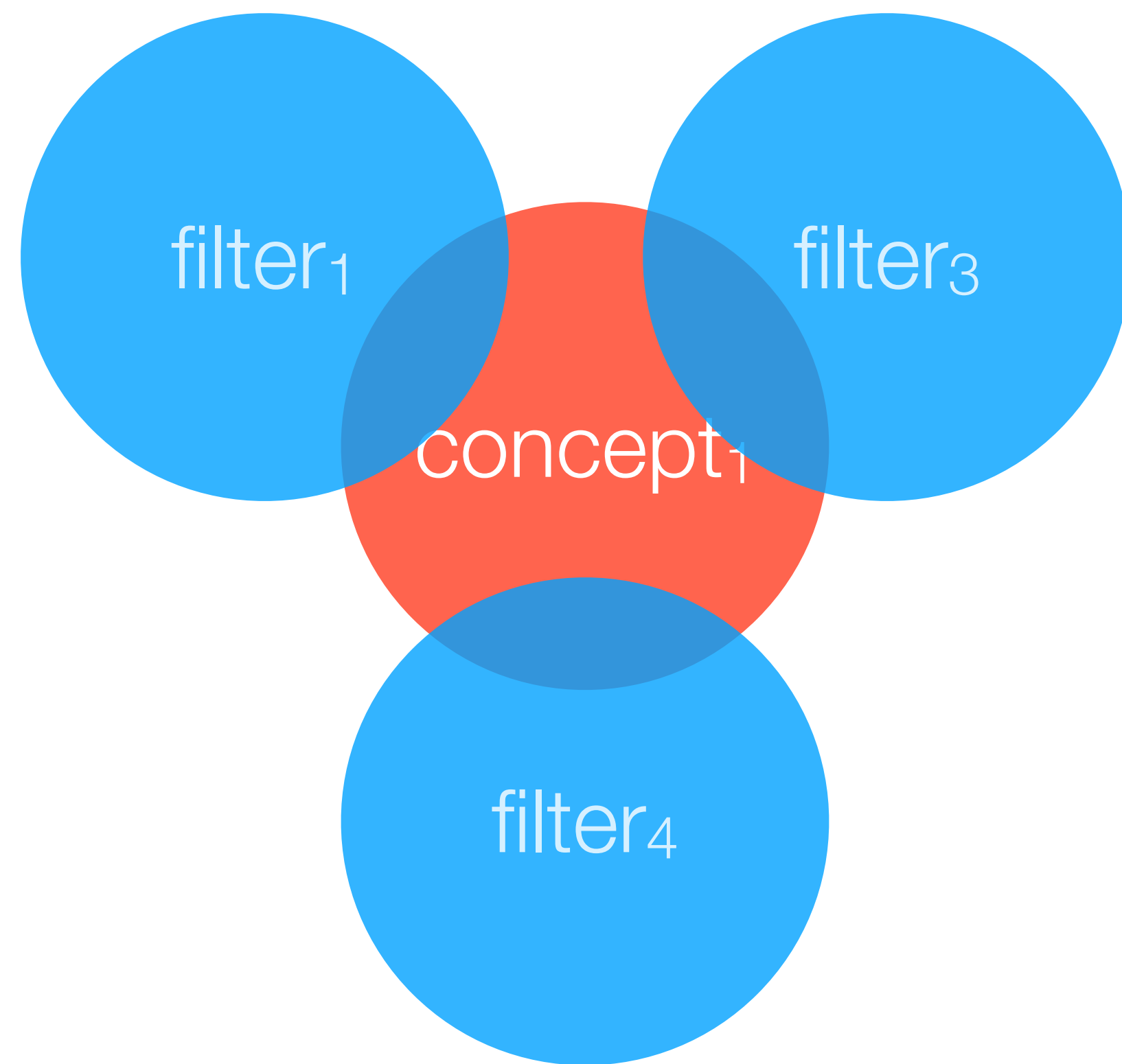


Filter-Concept Overlap



Modified Network Dissection visualization of AlexNet conv5 filters [Bau, et al., 2017]

Filter-Concept Overlap



Method

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

Bau et al, CVPR 2017. “Network Dissection.”

Agrawal et al., ECCV 2014. “Analyzing the performance of neural networks.”

Alain & Bengio, ICLR Workshop 2017. “Understanding intermediate layers using linear probes.”

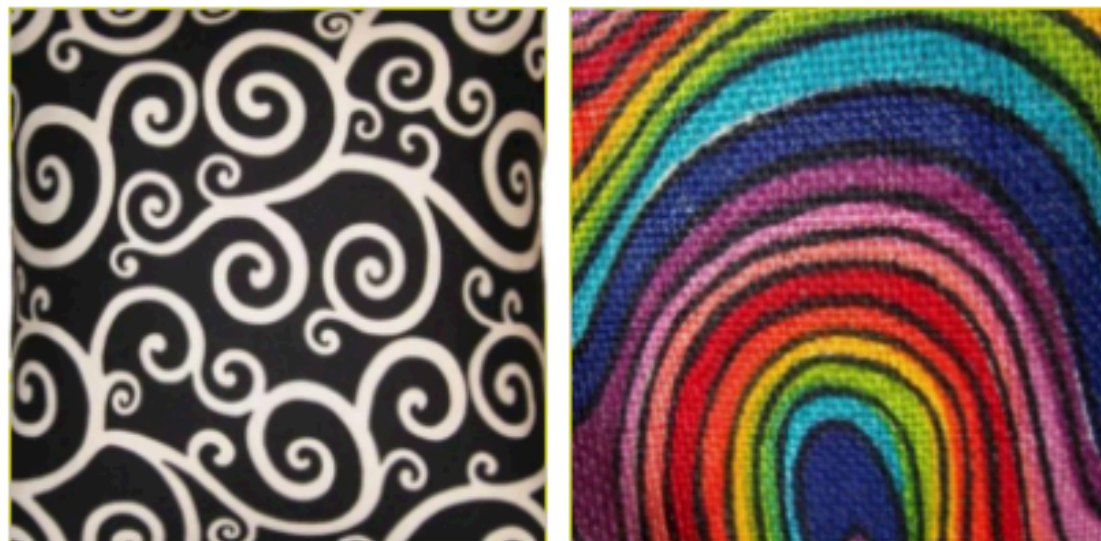
BRODEN Dataset

Image-level Annotations

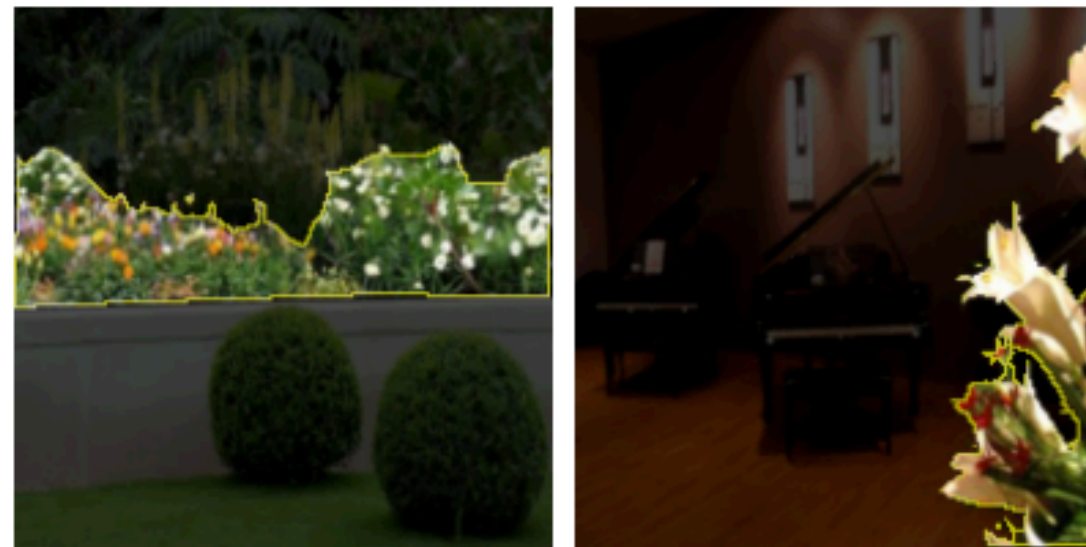
street (scene)



swirly (texture)



flower (object)



pink (color)

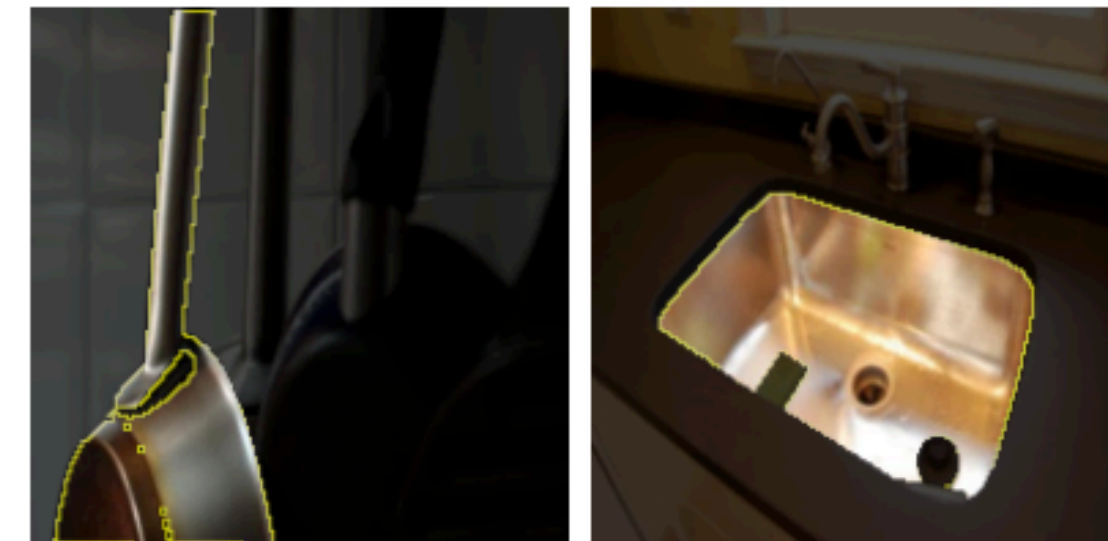


Pixel-level Annotations

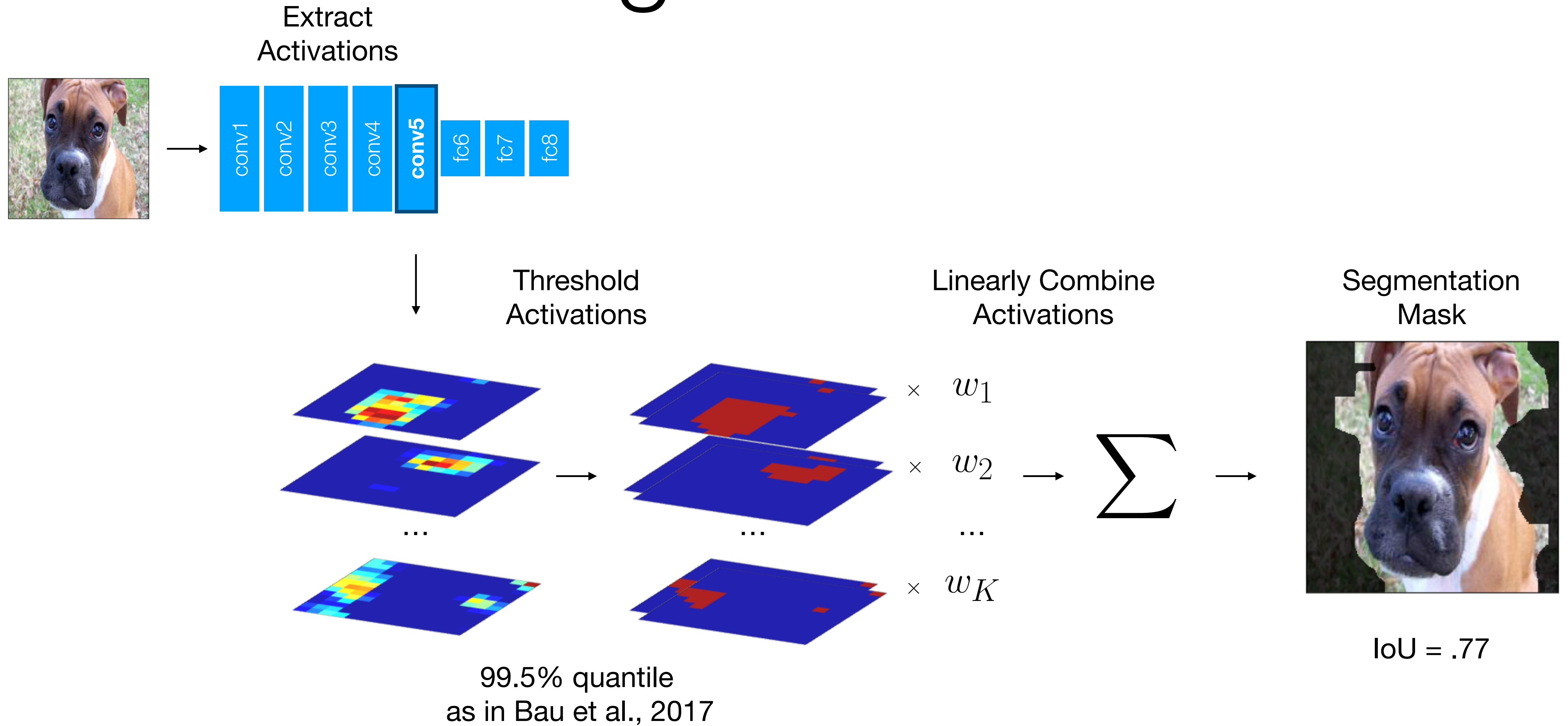
headboard (part)



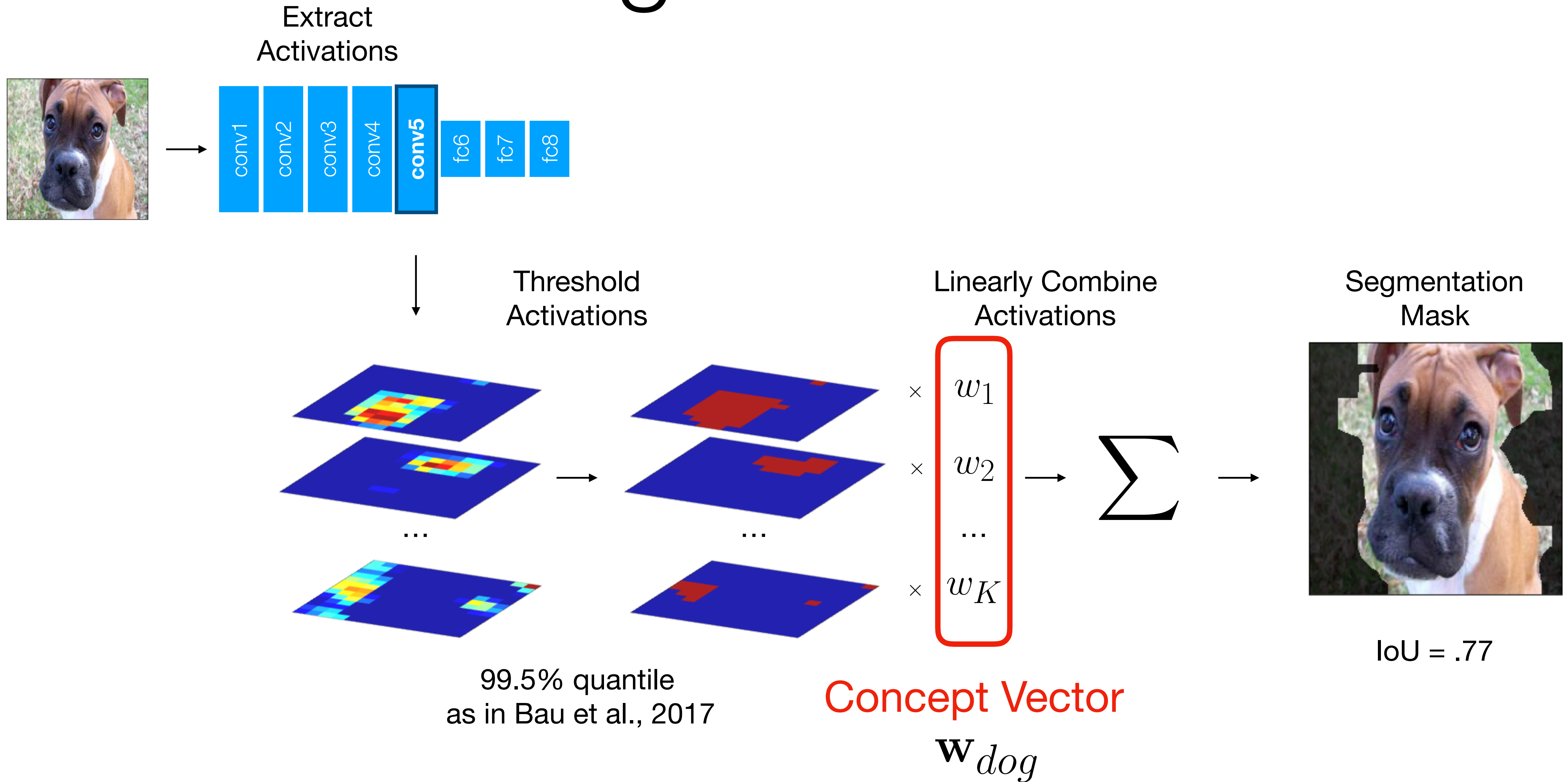
metal (material)



Segmentation



Segmentation



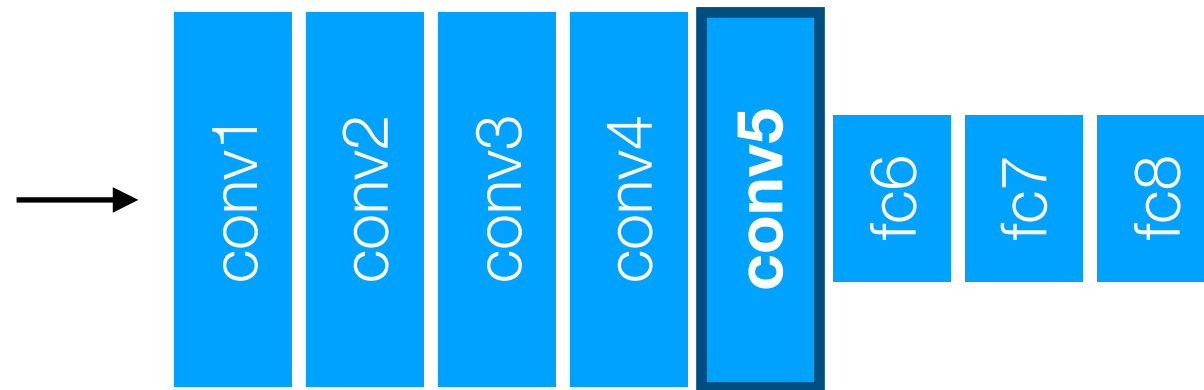
Segmentation

Subset:
Only use top F filters,
chosen by magnitude
($F = 4$)

Subset selection follows
Agrawal et al., 2014

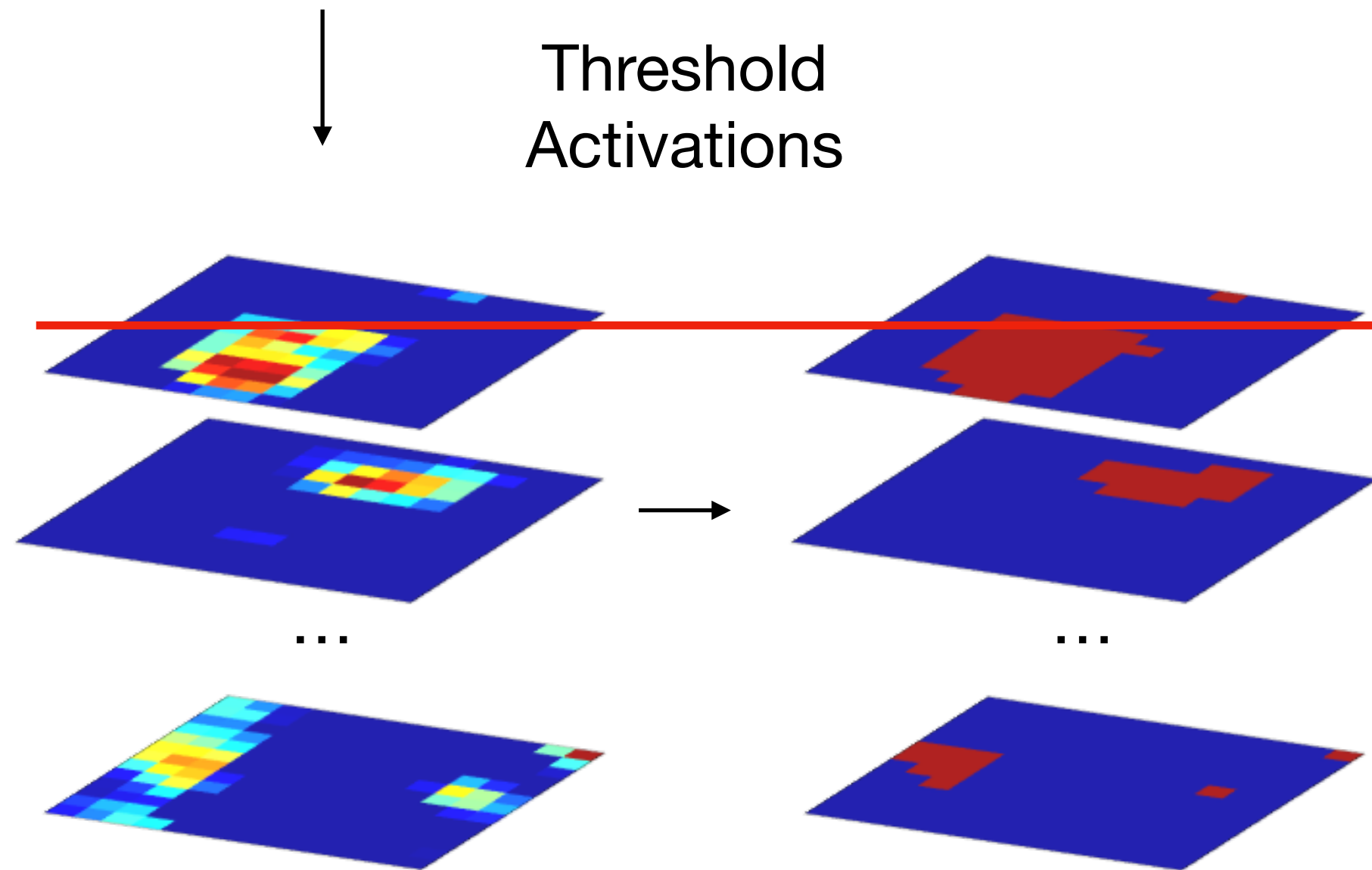


Extract
Activations



$$\mathbf{w}_{dog,F} = \begin{bmatrix} w_1 \\ \dots \\ w_F \end{bmatrix}$$

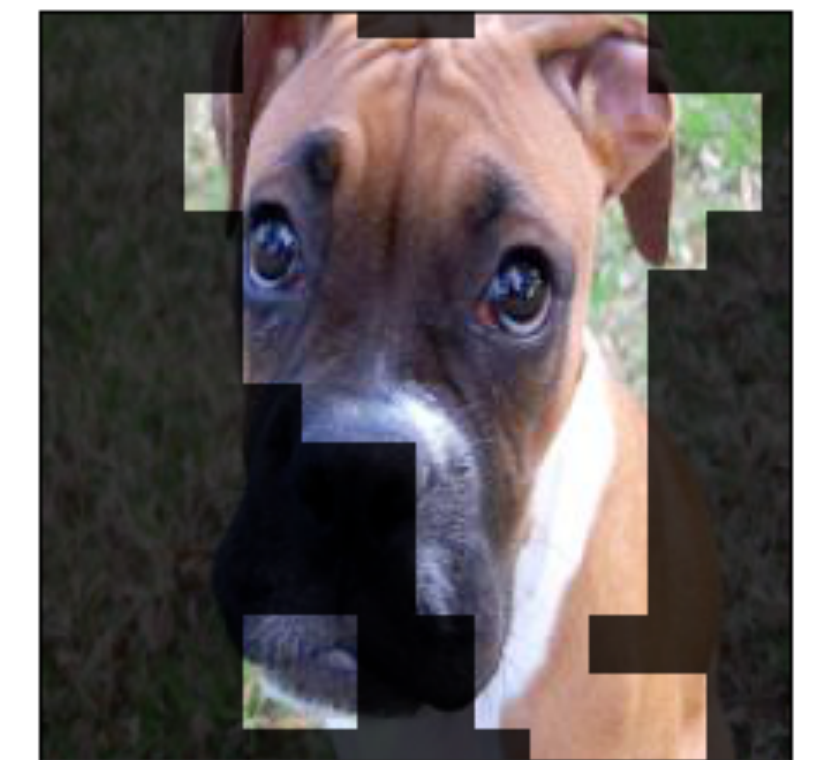
Threshold
Activations



Linearly Combine
Activations

$$\begin{matrix} \times w_1 \\ \dots \\ \times w_F \end{matrix} \rightarrow \Sigma \rightarrow$$

Segmentation
Mask



IoU = .63

Segmentation

Single Filter

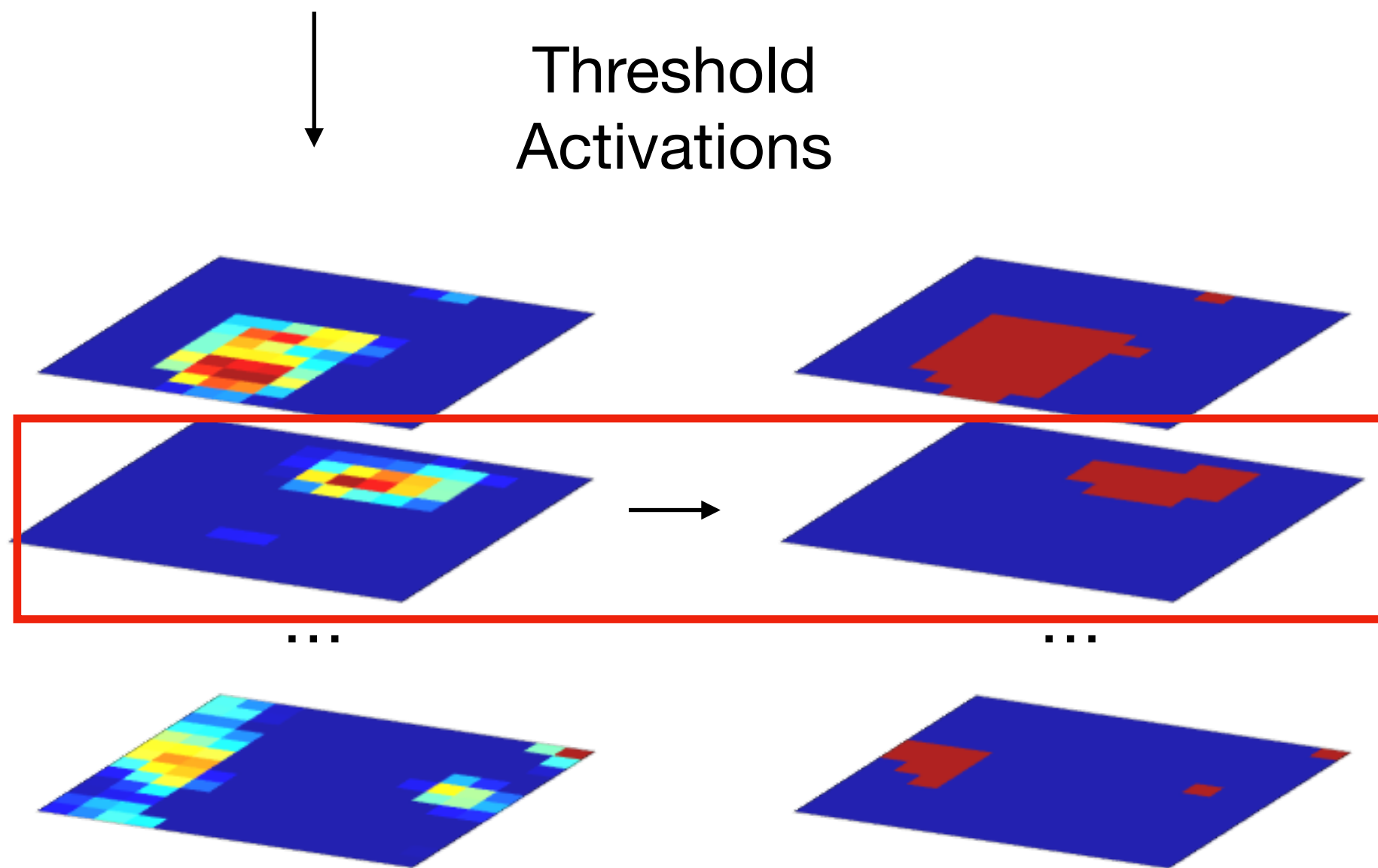


Extract
Activations



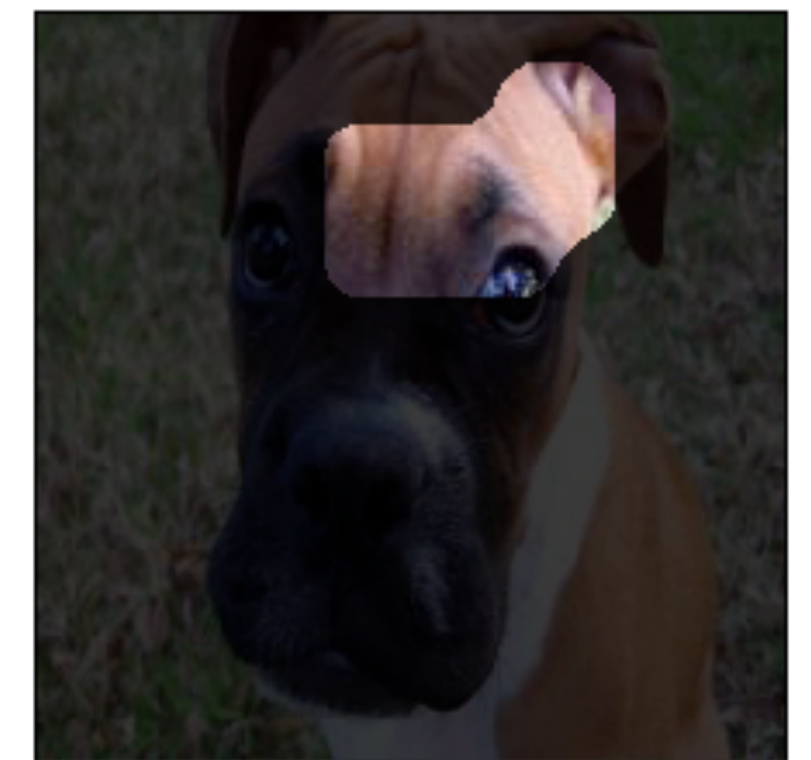
Threshold
Activations

Filter 169



Choose
Best Filter

Segmentation
Mask



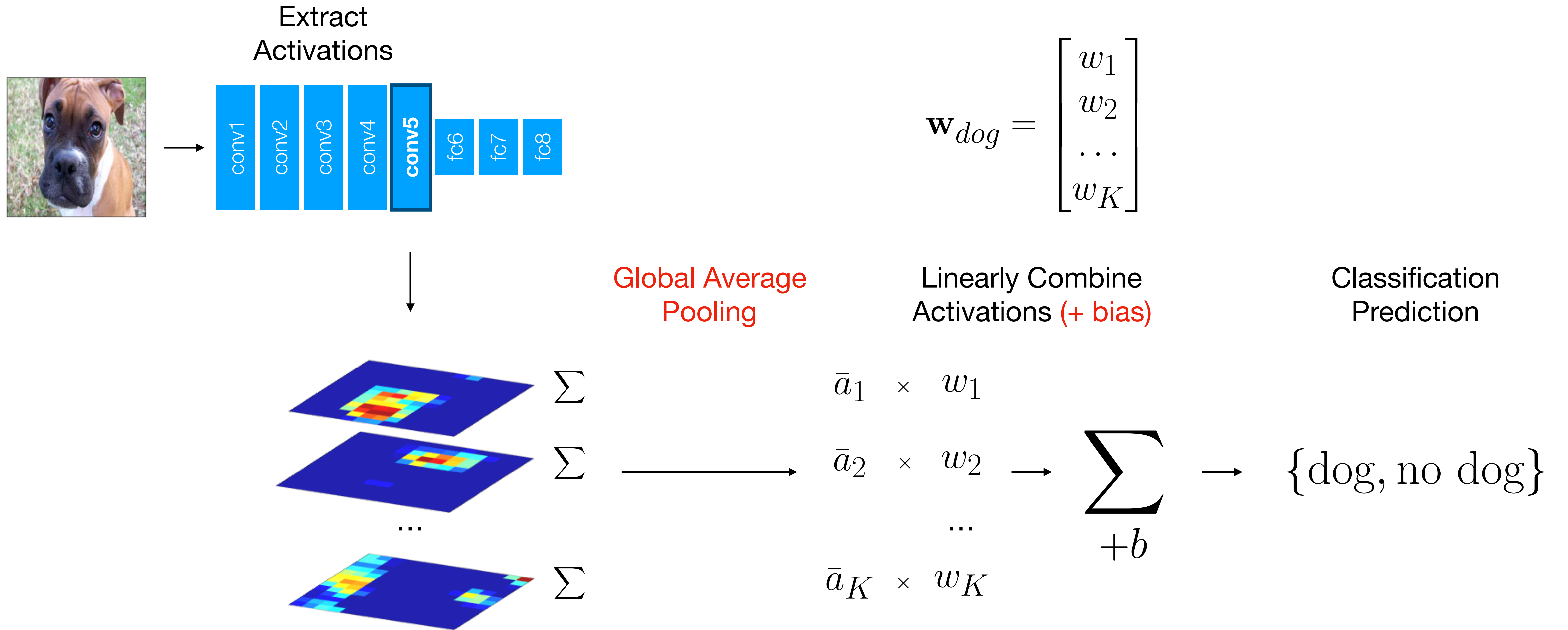
IoU = .18

Near equivalent to Bau et al., 2017

$$\text{IoU}_{\text{set}}(c; M, s) = \frac{\sum_{\mathbf{x} \in X_{s,c}} |M(\mathbf{x}) \cap L_c(\mathbf{x})|}{\sum_{\mathbf{x} \in X_{s,c}} |M(\mathbf{x}) \cup L_c(\mathbf{x})|}$$

$$\text{IoU}_{\text{ind}}(\mathbf{x}, c; M) = \frac{|M(\mathbf{x}) \cap L_c(\mathbf{x})|}{|M(\mathbf{x}) \cup L_c(\mathbf{x})|}$$

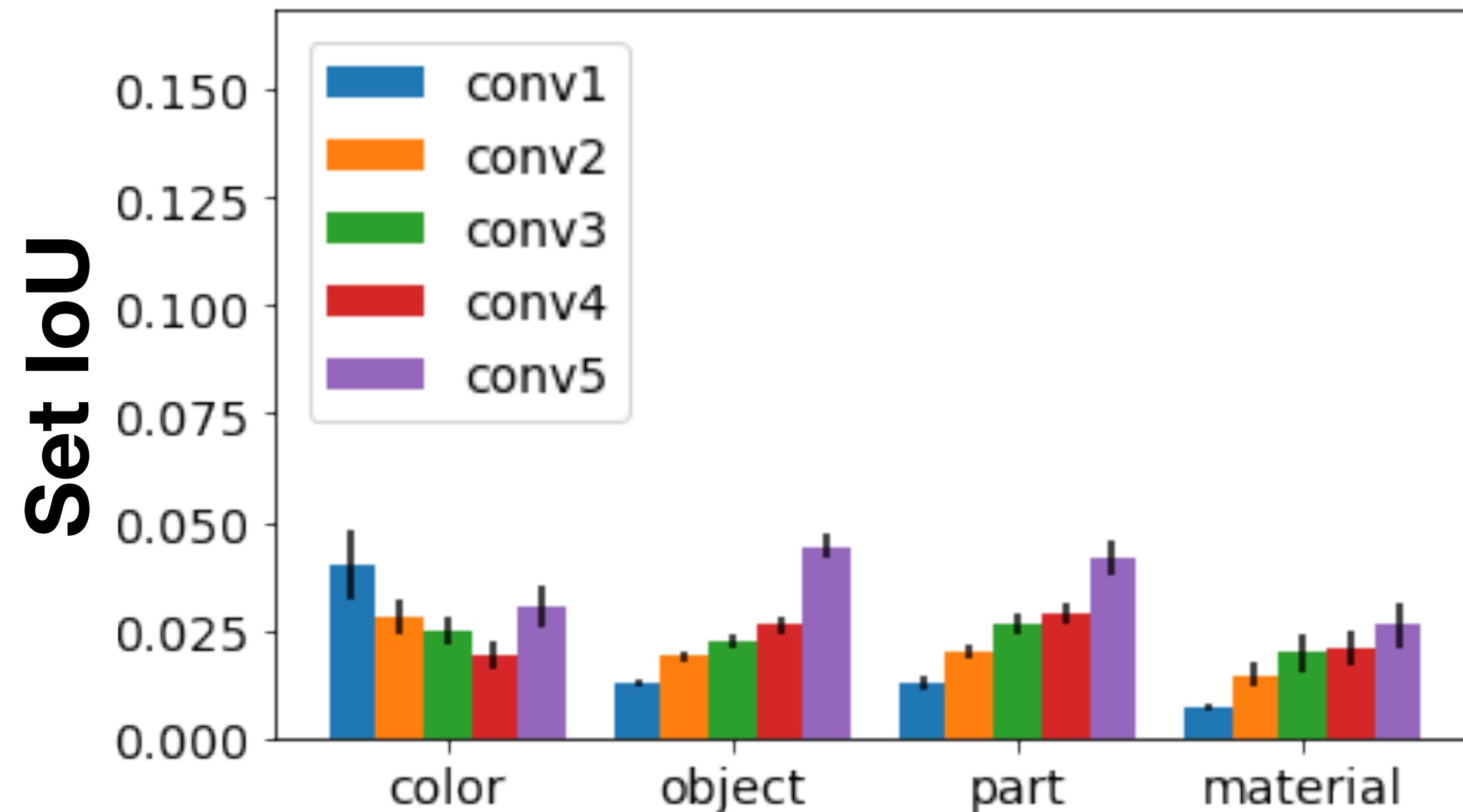
Classification



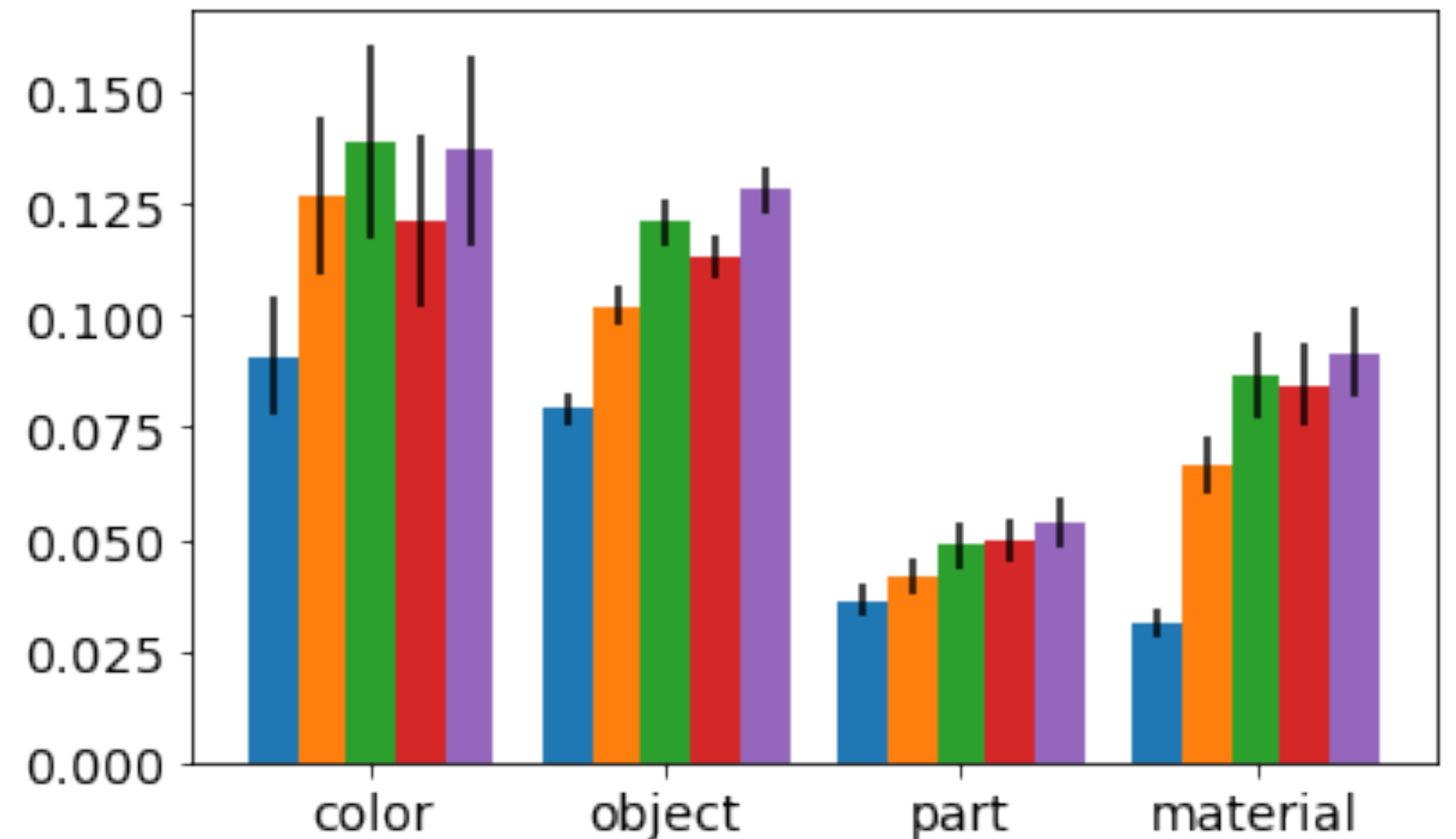
A Few Results

Single vs. All Filters

Single Filter

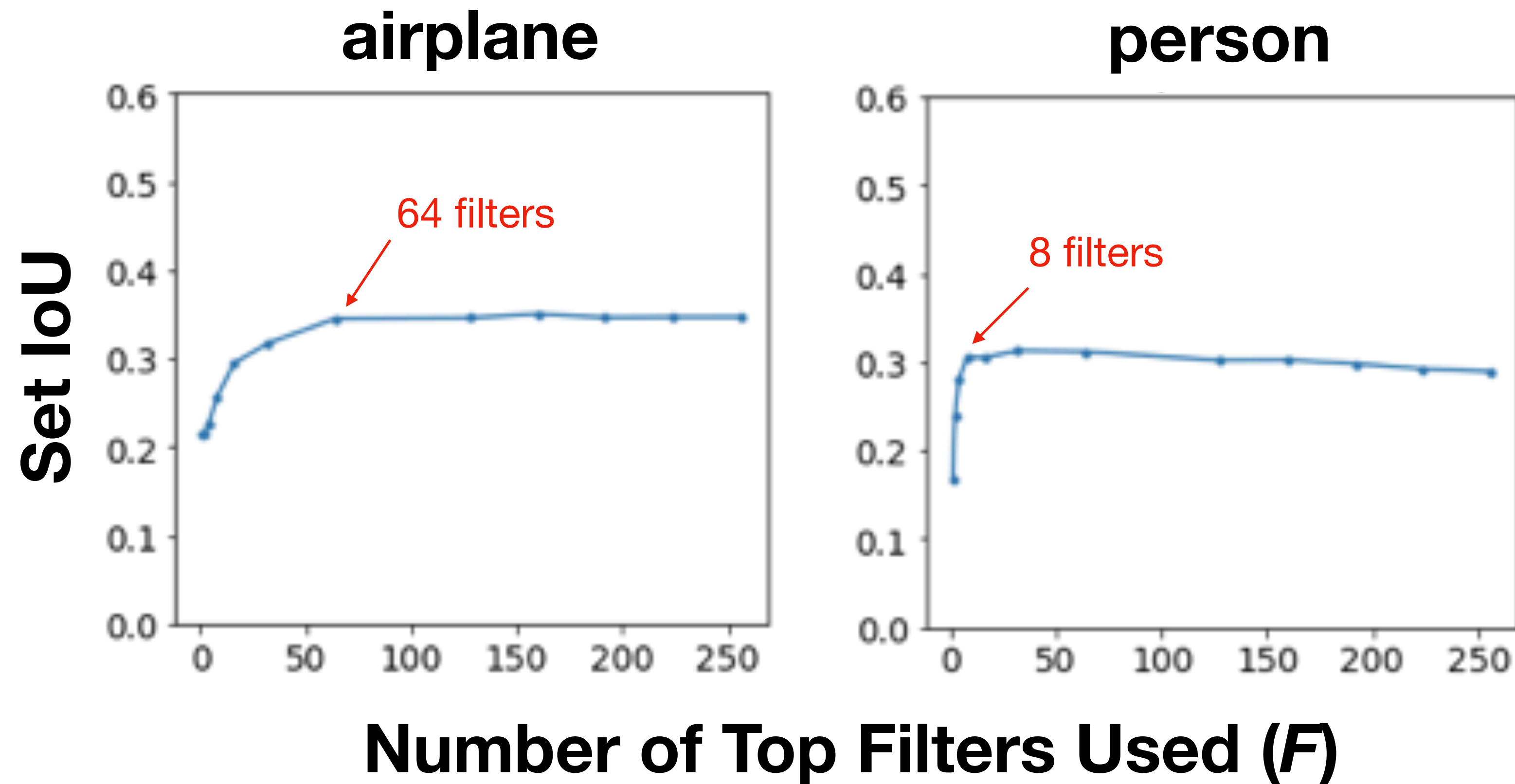


All Filters



Concepts are encoded better when using multiple filters.

Filters Per Concept



Different concepts require different number of filters for encoding.

Filters: Supervised vs. Self-Supervised

Performance Improvement (Single Filter → All Filters):

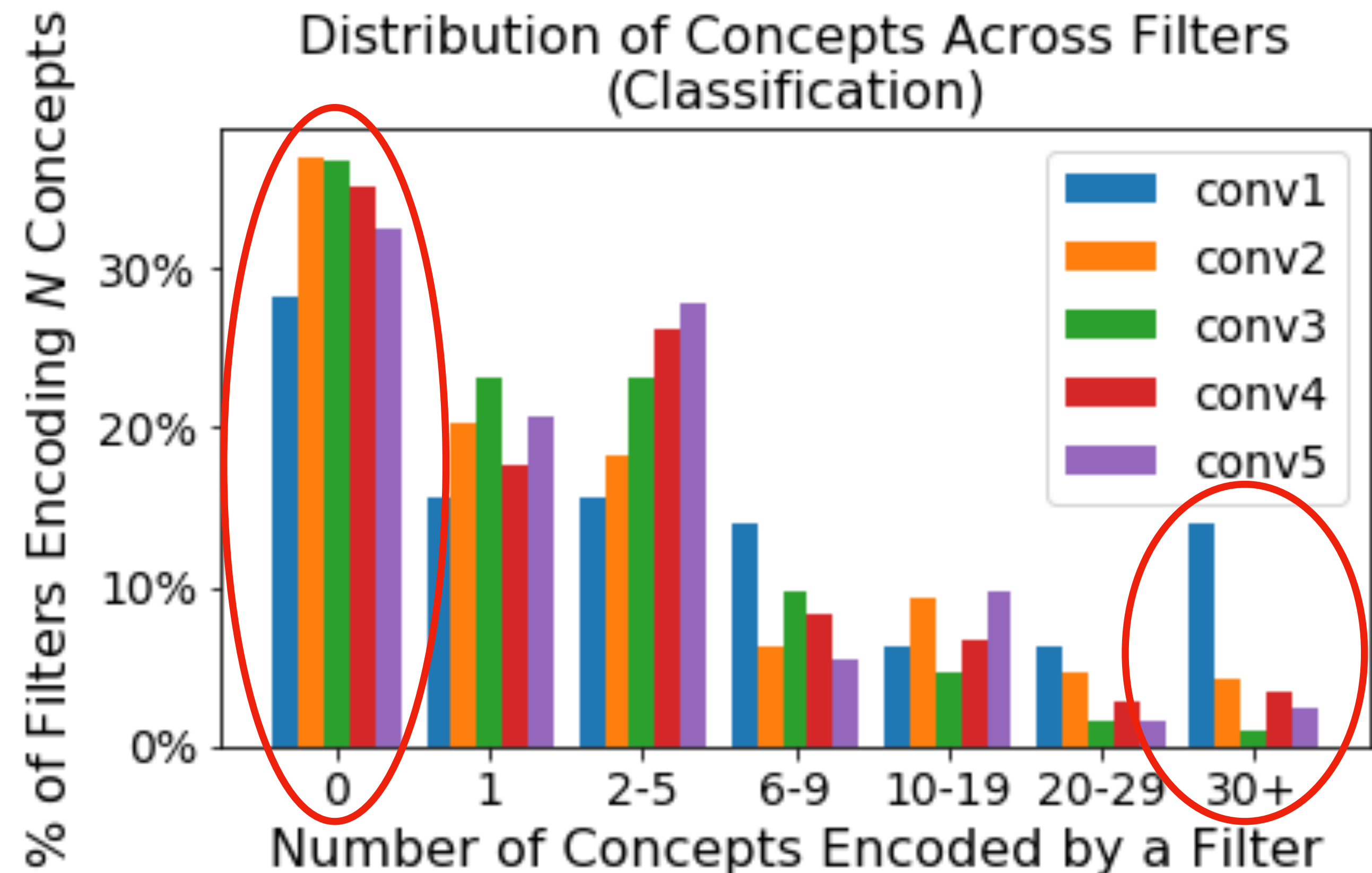
- Self-supervised networks: 5-6x
- Fully-supervised networks: 2-4x

Self-supervised networks encode BRODEN concepts more distributively.

Concepts per Filters

Found a wide range in filter capacity to encode concepts:

- Many filters aren't selective for any concepts
- A few filters are selective for many concepts



Concepts per Filters

Sheep
(IoU_{set} = .21)



AlexNet conv5 unit 66
is highly selective for
various farm animals

Horse
(IoU_{set} = .21)



Cow
(IoU_{set} = .20)



Visualizing Non-Maximal Examples

IoU_{ind}

All Filters

$\text{IoU}_{\text{set}} = .35$

[ours]

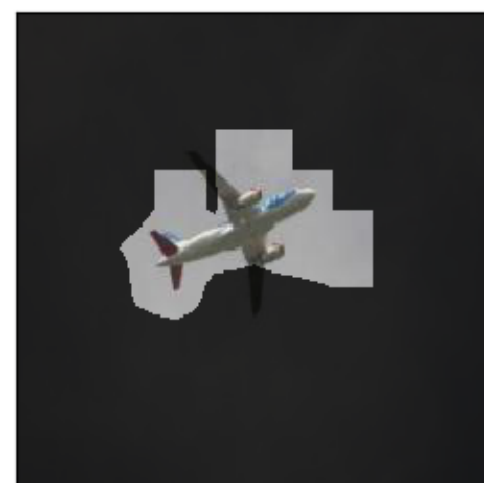
.01



.02



.24



.28



.33



.37



.39



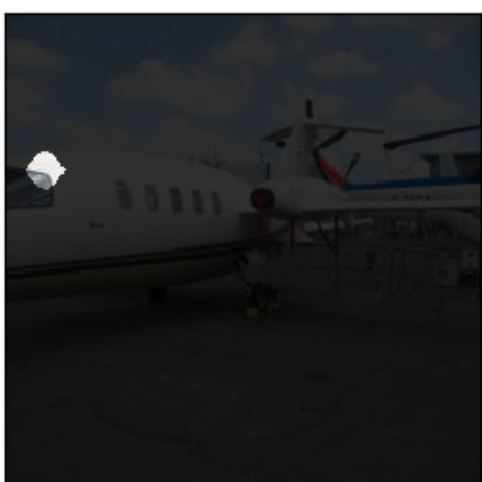
.43



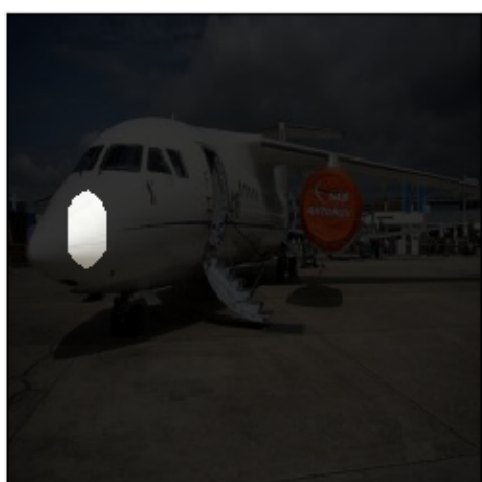
.49



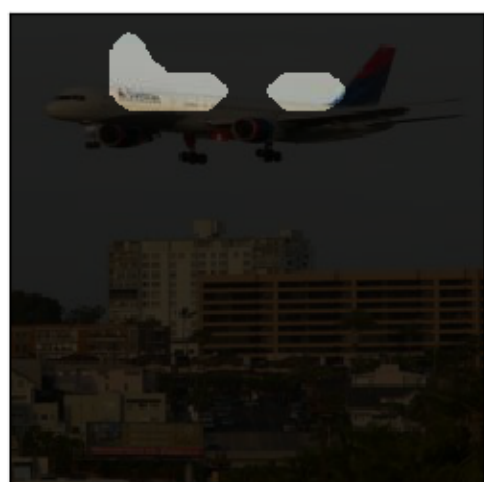
.02



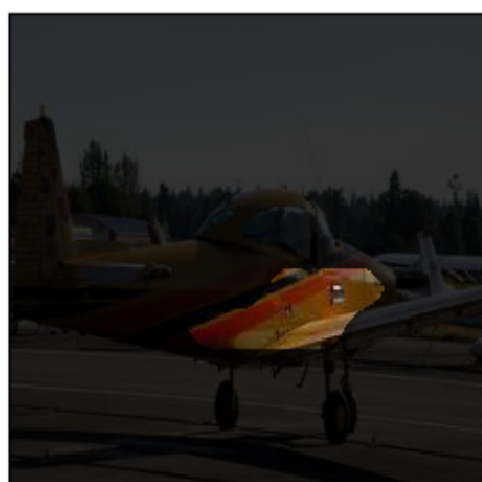
.05



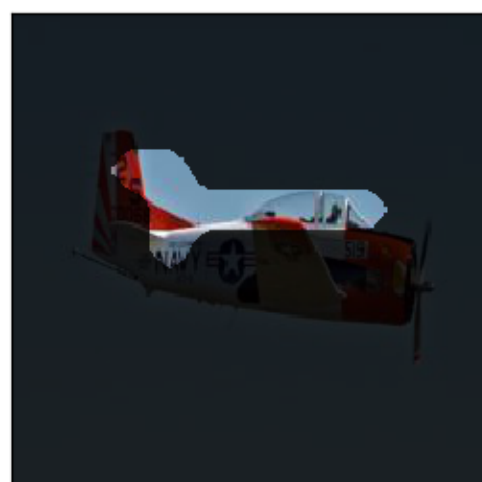
.10



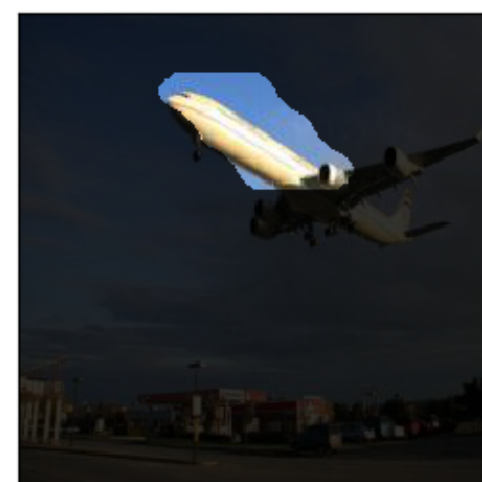
.16



.19



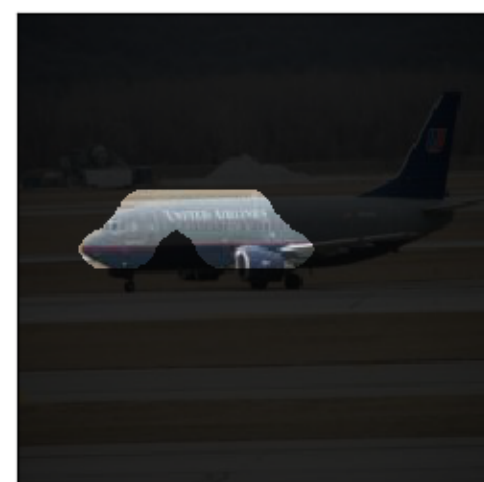
.27



.31



.37



.46



10%

20%

30%

40%

50%

60%

70%

80%

90%

Spanning Deciles



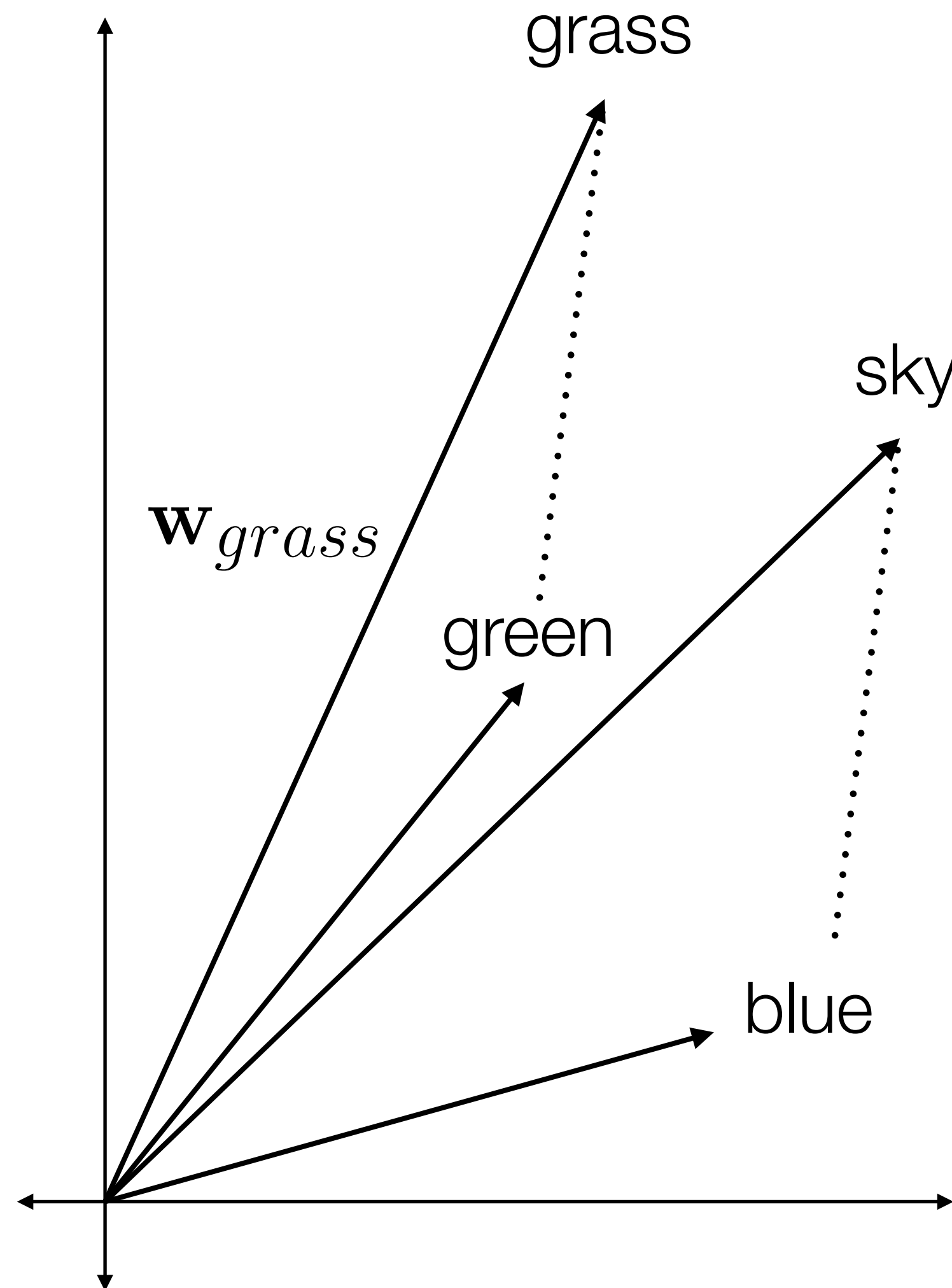
Best Filter

$\text{IoU}_{\text{set}} = .14$

[Bau et al.,

2017]

Comparing Concept Embeddings

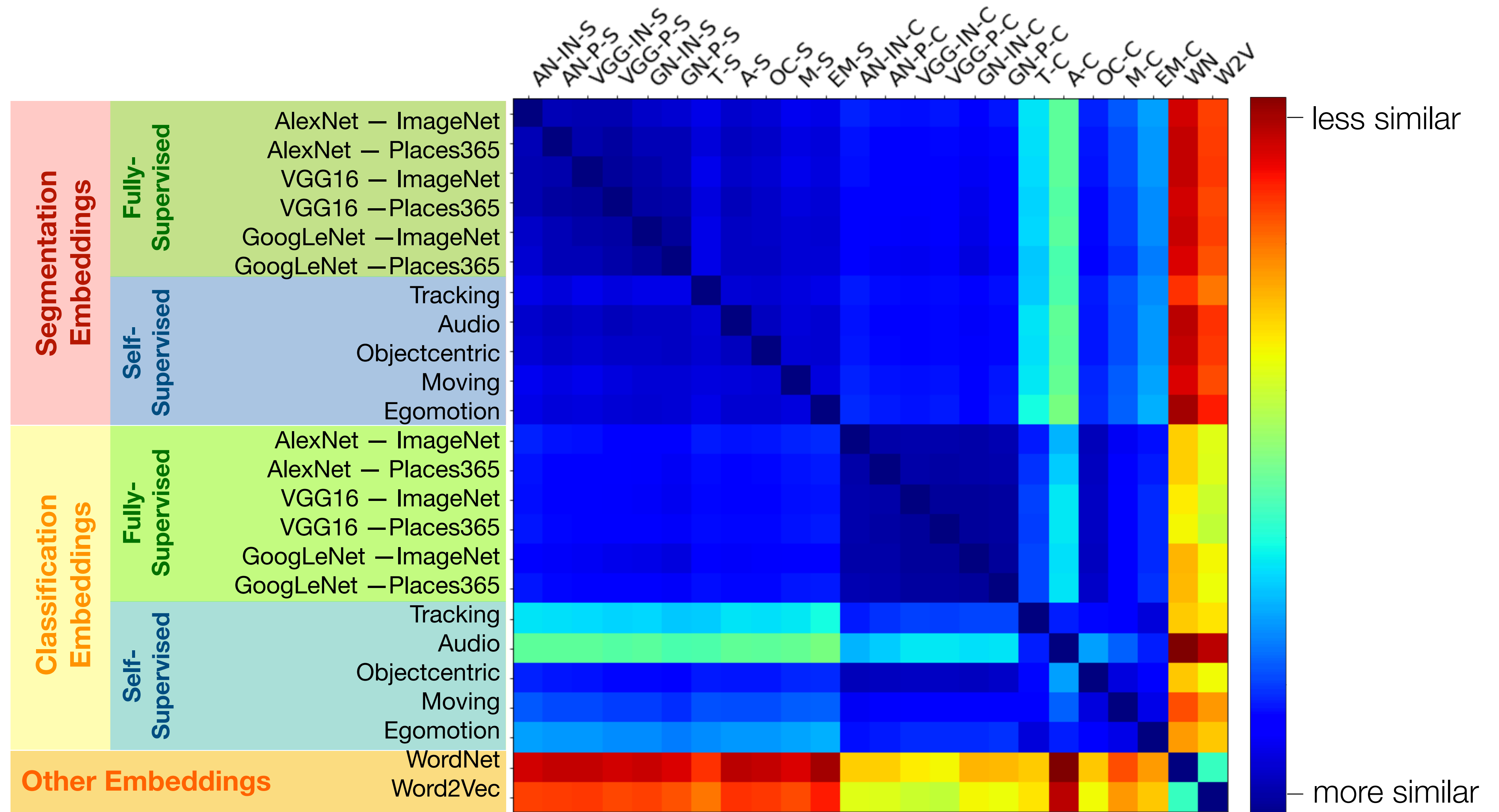


grass + blue — green = sky

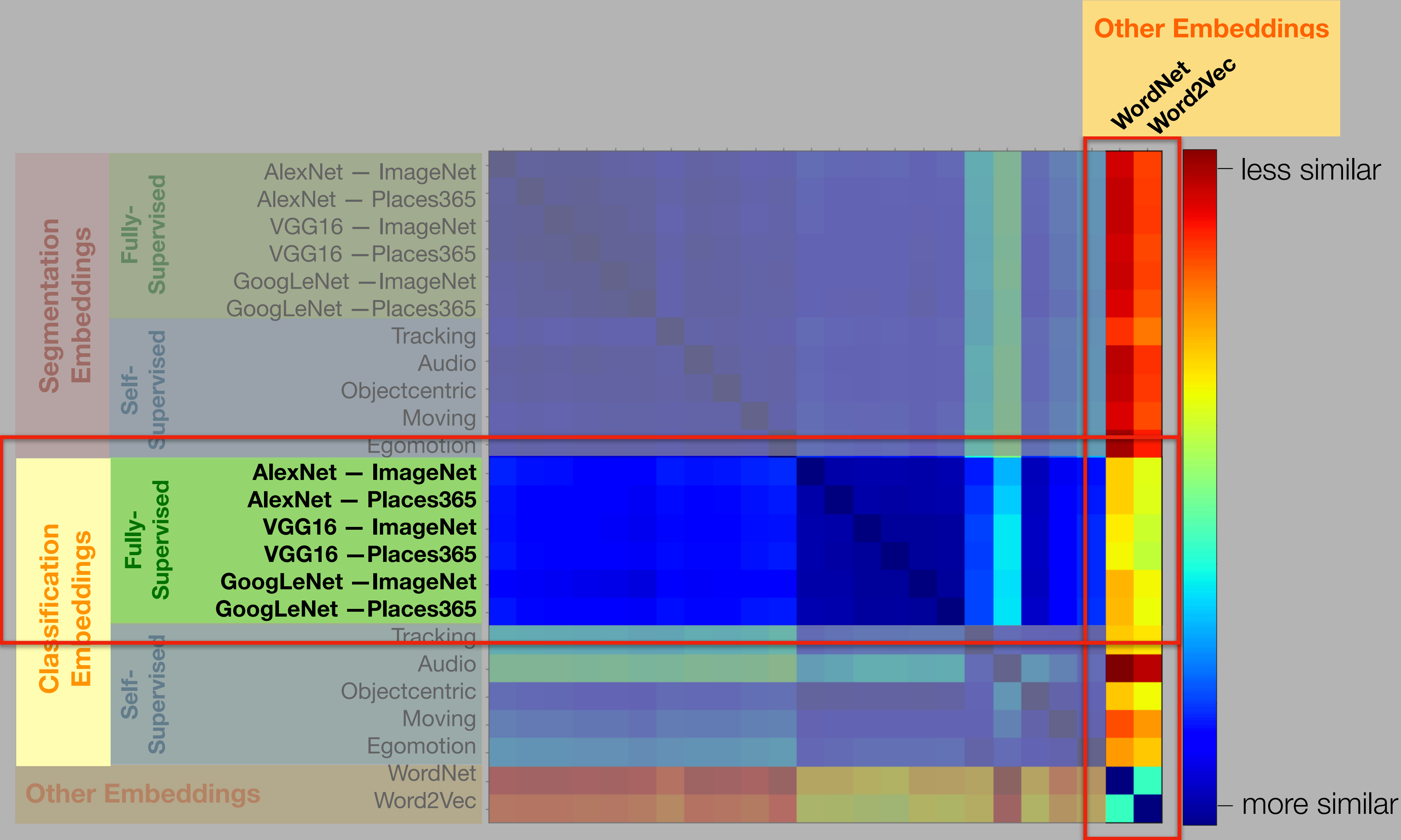
tree — wood = plant

person — torso = foot

Comparing Concept Embedding **Spaces**



Comparing Concept Embedding Spaces



Chat more at poster E9!

Code: <https://github.com/ruthcfong/net2vec>

