

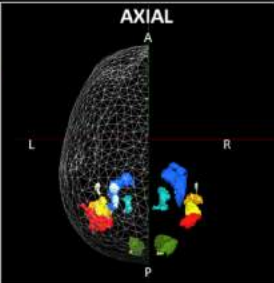
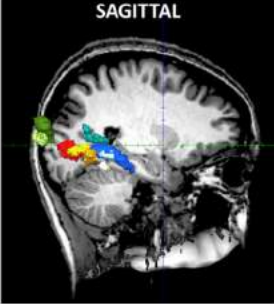
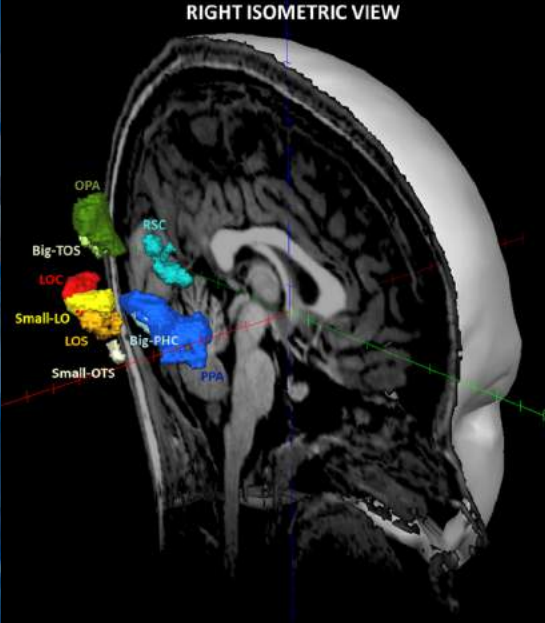
Interpretability and Visualization of Deep Neural Networks

Aude Oliva
MIT



MIT-IBM Watson AI Lab





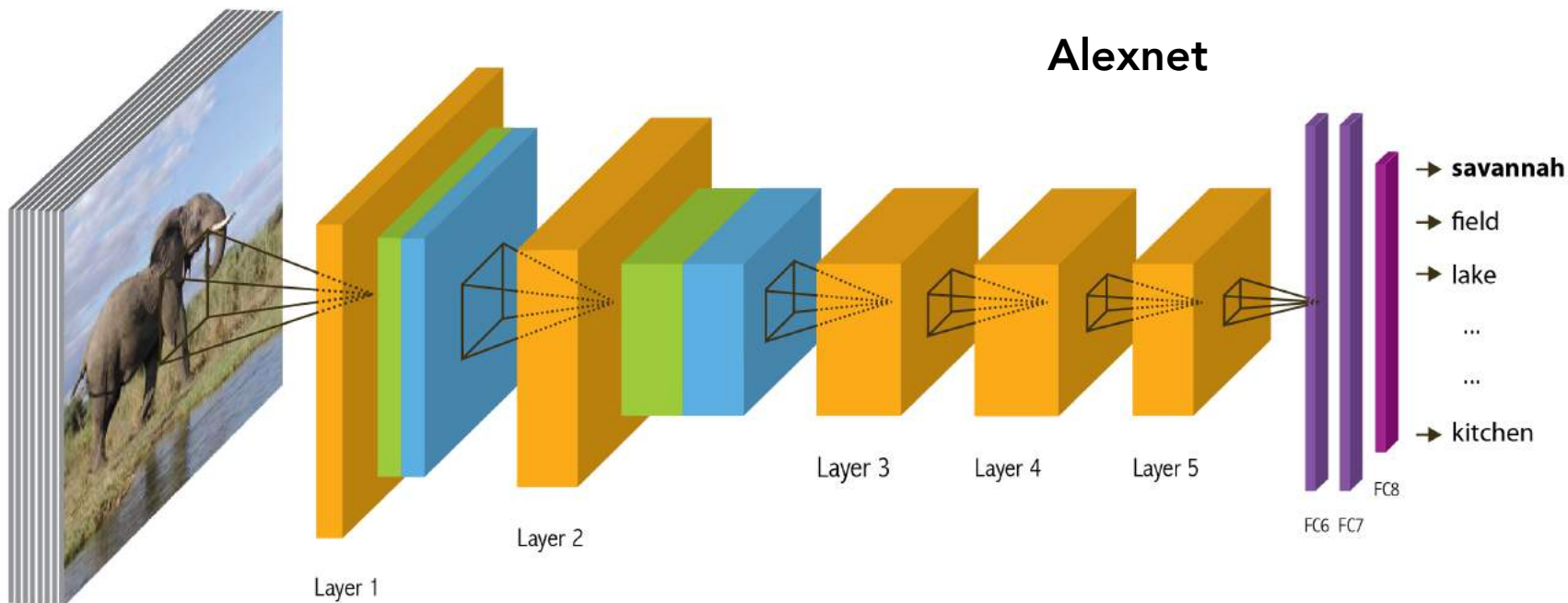
IM  GENET

places 



Convolutional Neural Networks

convolution max-pooling normalization full connected



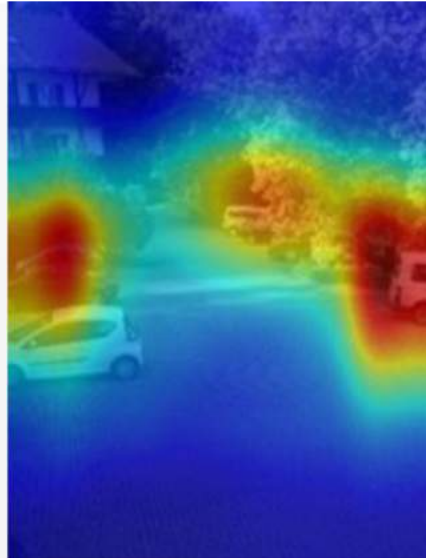
Each layer learns progressively more complex features

places



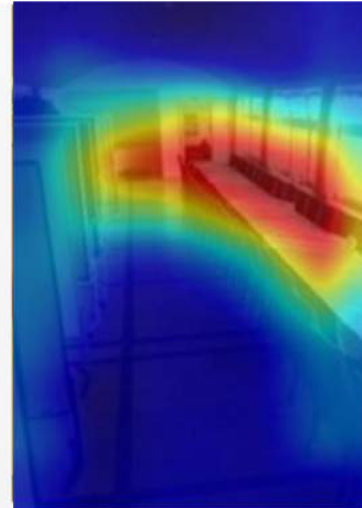
Predictions:

- **Type of environment:** indoor
- **Semantic categories:** restaurant:0.27, coffee_shop:0.23, cafeteria:0.21, food_court:0.12, restaurant_patio:0.09



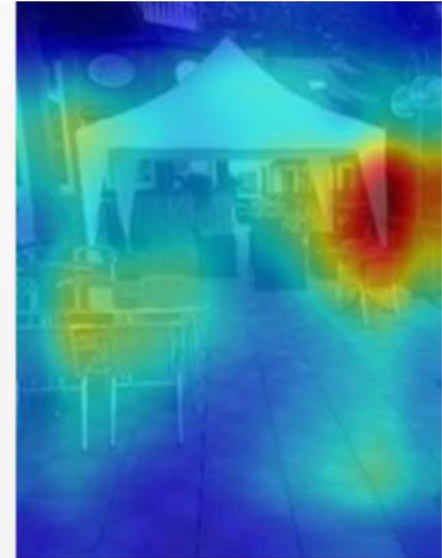
Predictions:

- **Type of environment:** outdoor
- **Semantic categories:** parking_lot:0.46, driveway:0.44,



Predictions:

- **Type of environment:** indoor
- **Semantic categories:** conference_room:0.29, dining_room:0.27, banquet_hall:0.08, classroom:0.06,

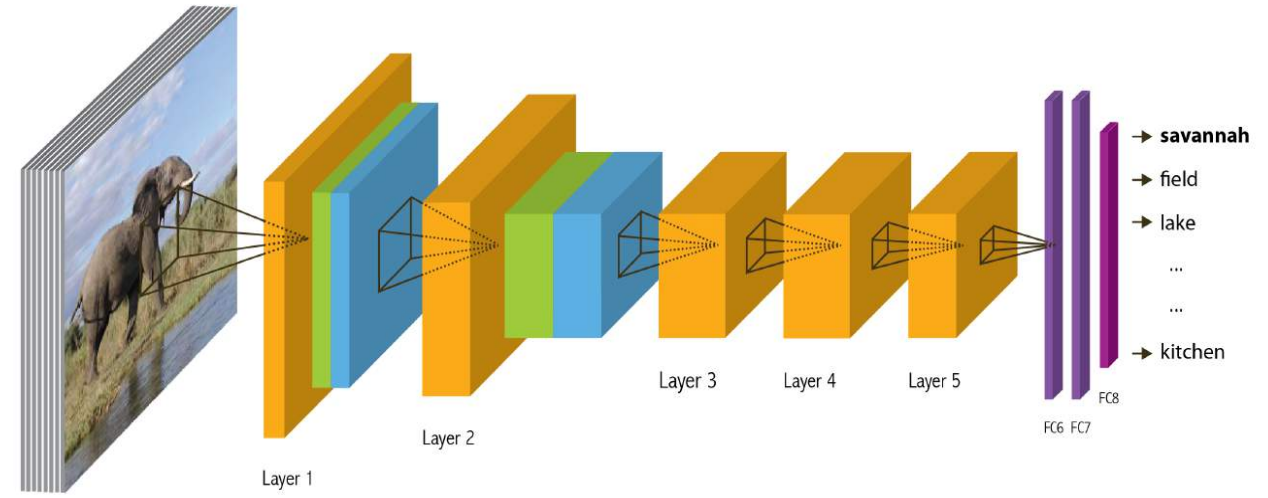
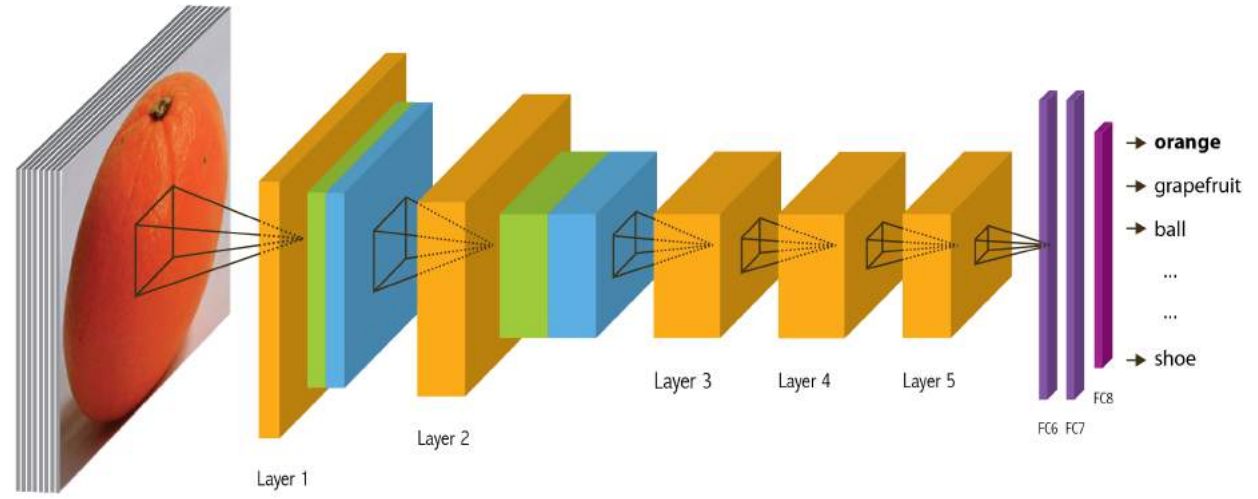


Predictions:

- **Type of environment:** outdoor
- **Semantic categories:** patio:0.38, restaurant_patio:0.35, restaurant:0.06,

What did the network learn ?

Comparing Object and Scenes CNNs



Data driven approach inspired by Neuroscience: Empirical receptive field



Pipeline for estimating the Receptive Fields:

Goal is to identify which regions of the image lead to the high unit activations.



sliding-window stimuli

5000 occluded versions

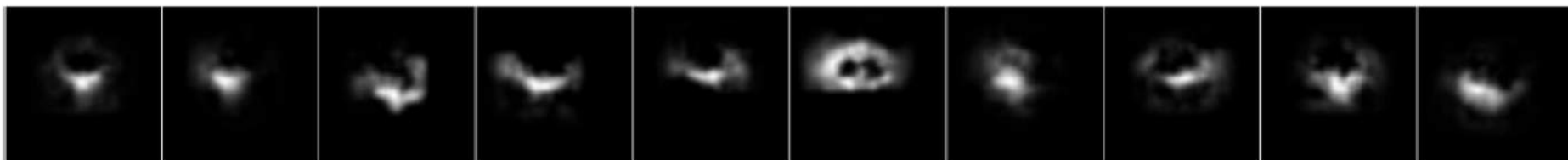
Discrepancy map per unit



Pipeline for estimating the Receptive Fields



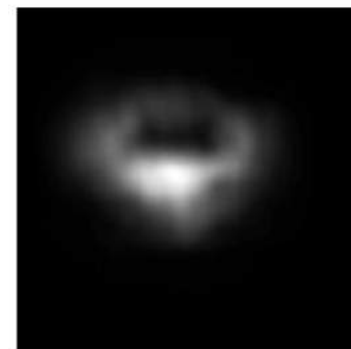
discrepancy maps for top 10 images



calibrated discrepancy maps

To consolidate the information from several images, we center the discrepancy map around the spatial location of the unit that caused the maximum activation for the given image.

Then we average the re-centered discrepancy maps to generate the final receptive field of each given unit.



receptive field

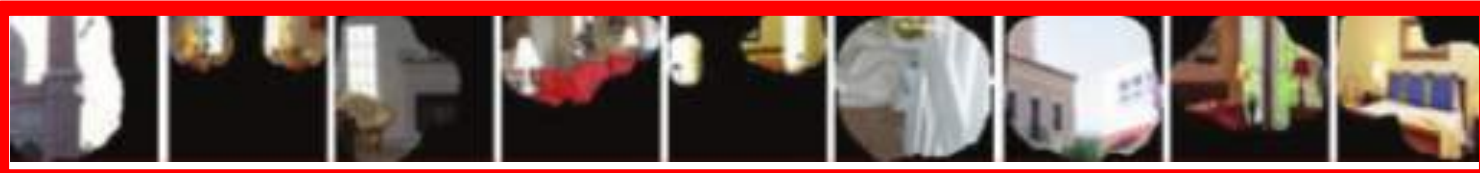
Annotating the Semantics of Units

Pool5, unit 76; Label: ocean; Type: scene; Precision: 93%



Annotating the Semantics of Units

Pool5, unit 13; Label: Lamps; Type: object; Precision: 84%



Annotating the Semantics of Units

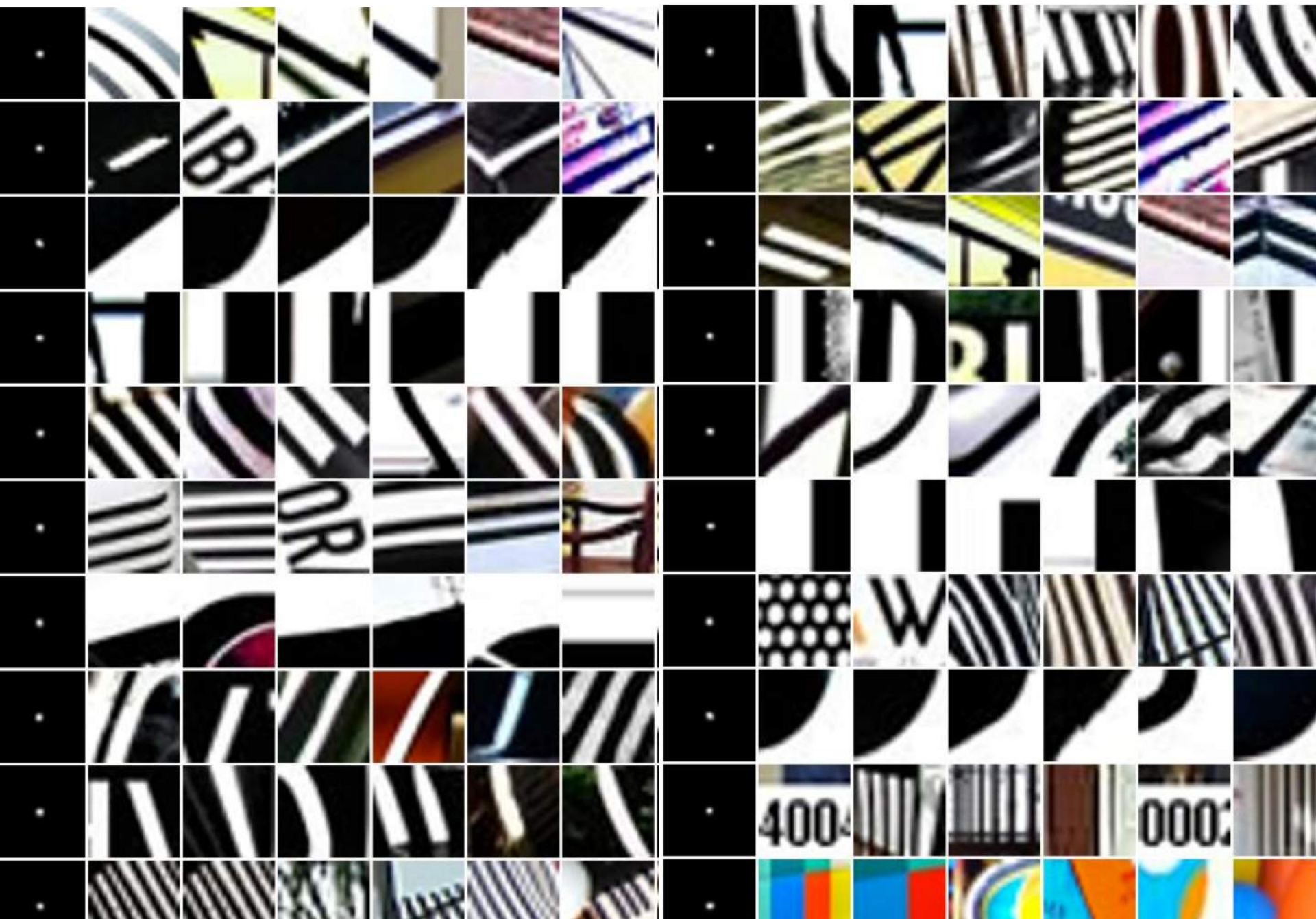
Pool5, unit 77; Label: legs; Type: object part; Precision: 96%

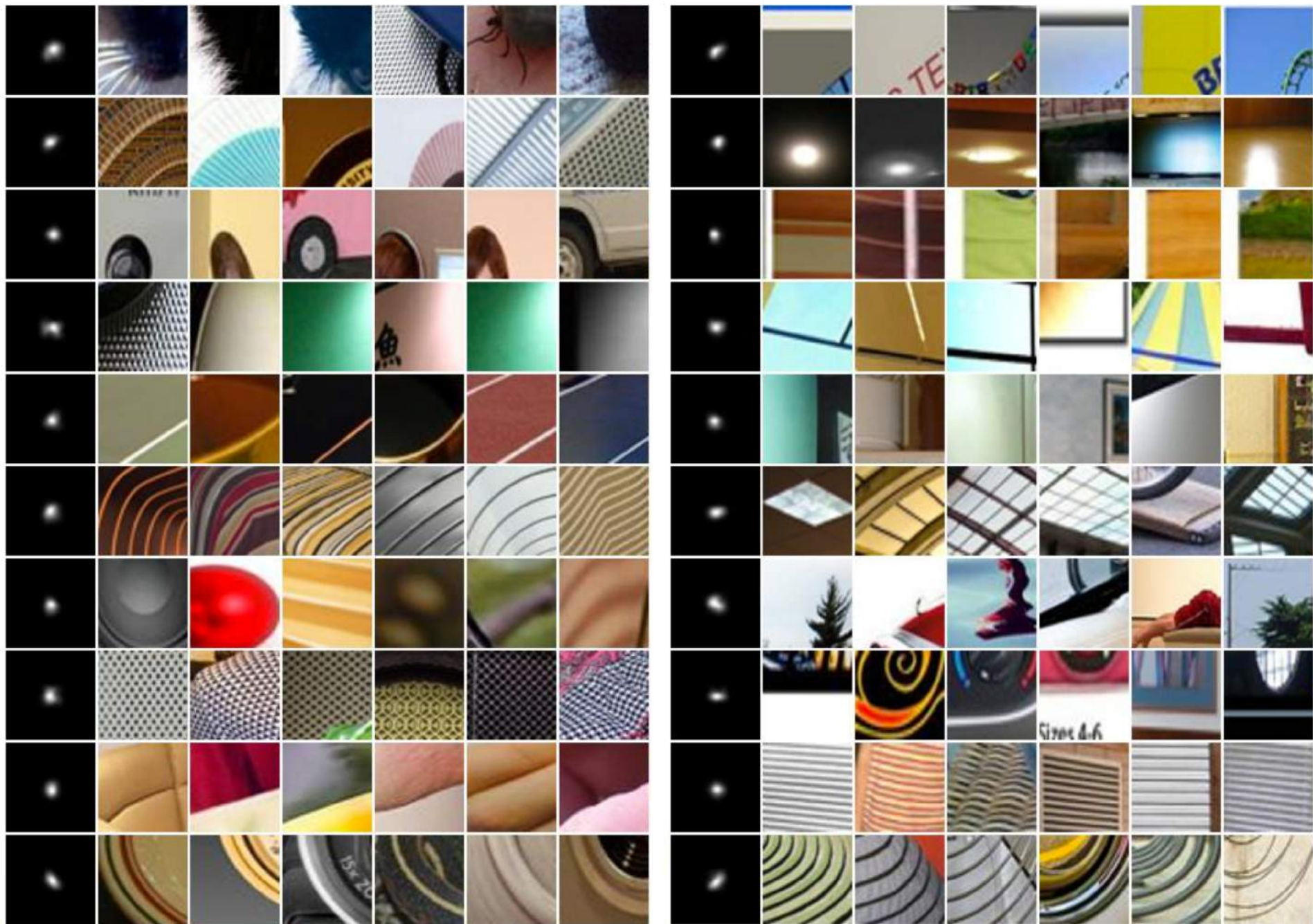


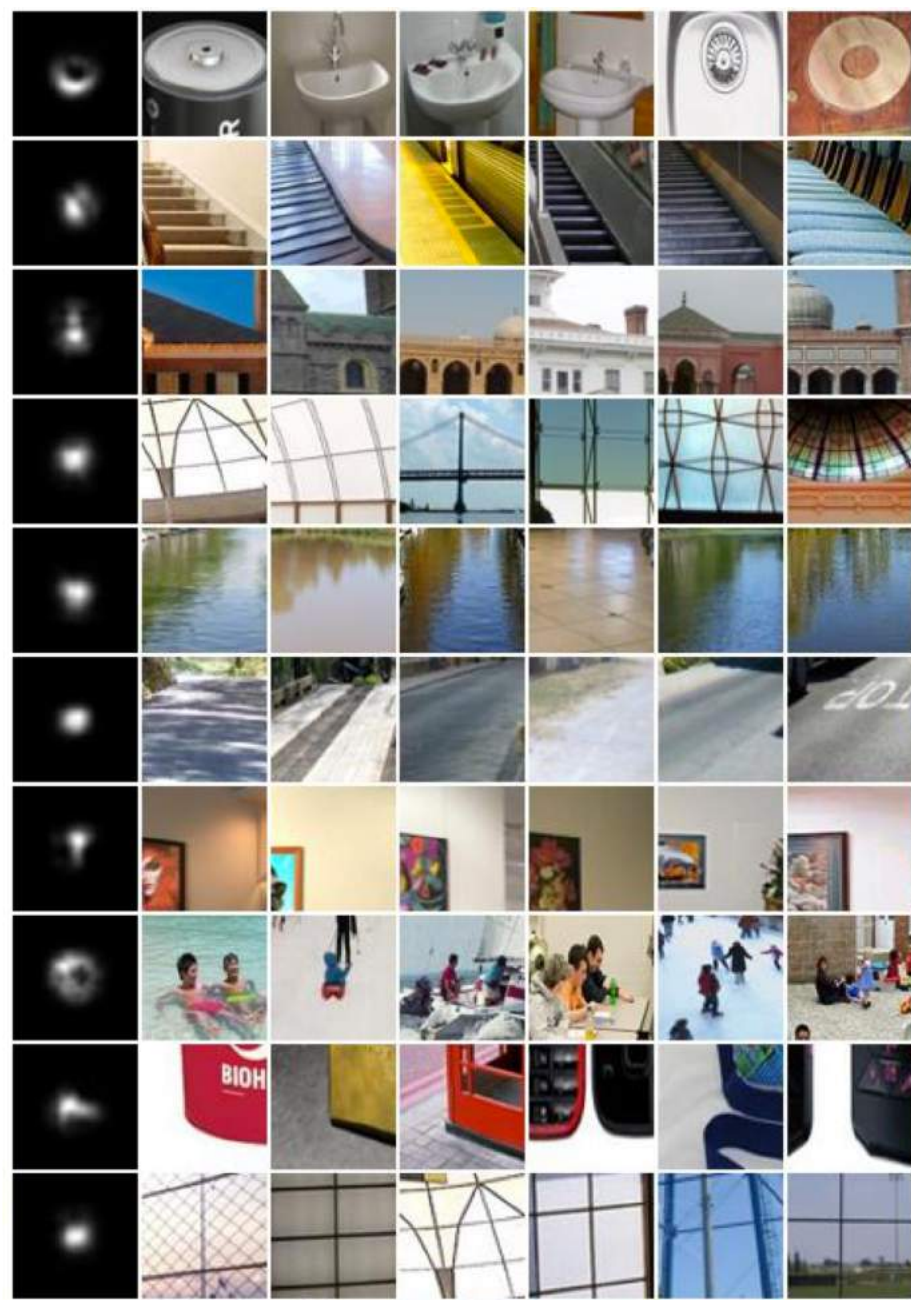
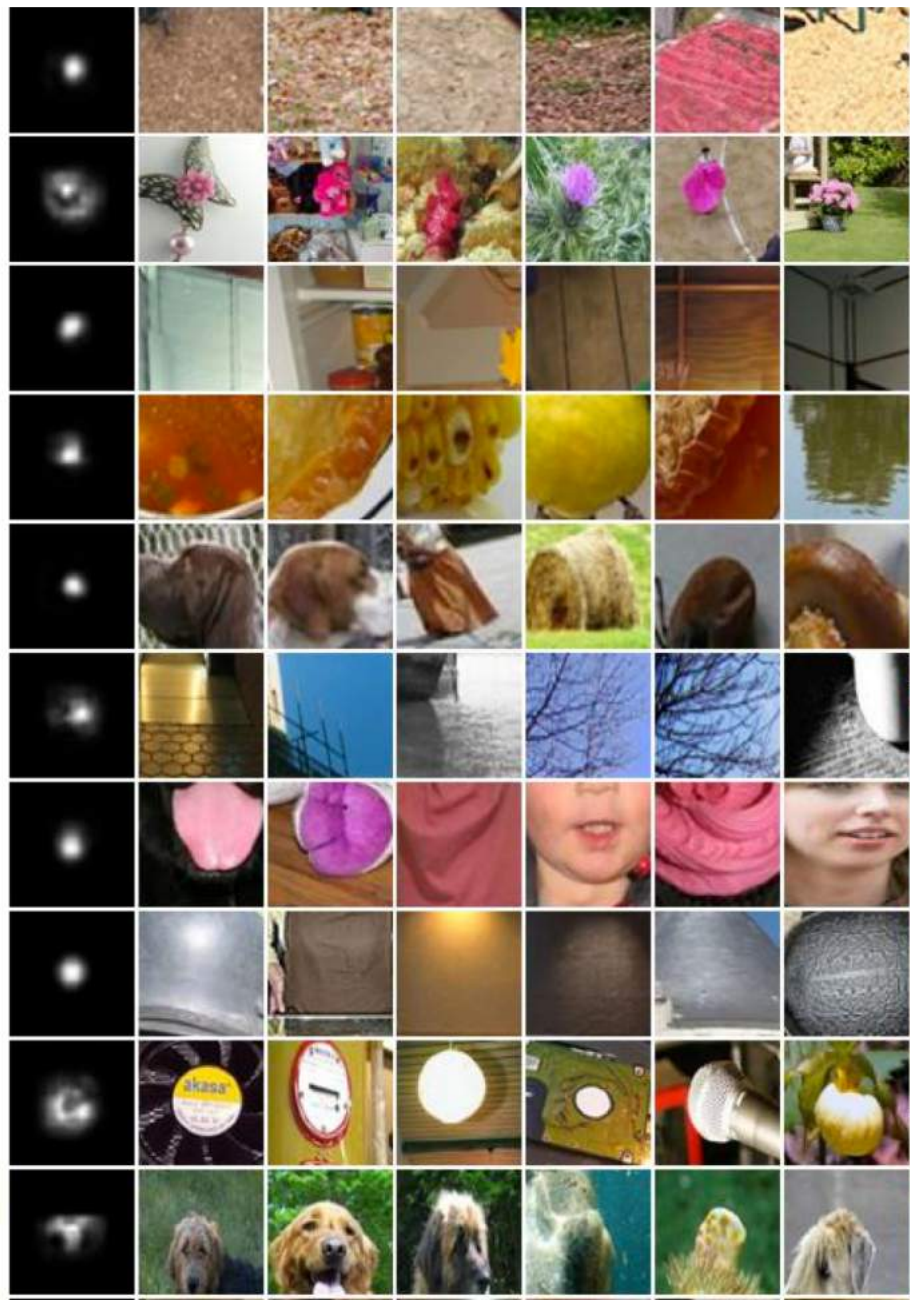
IMAGENET

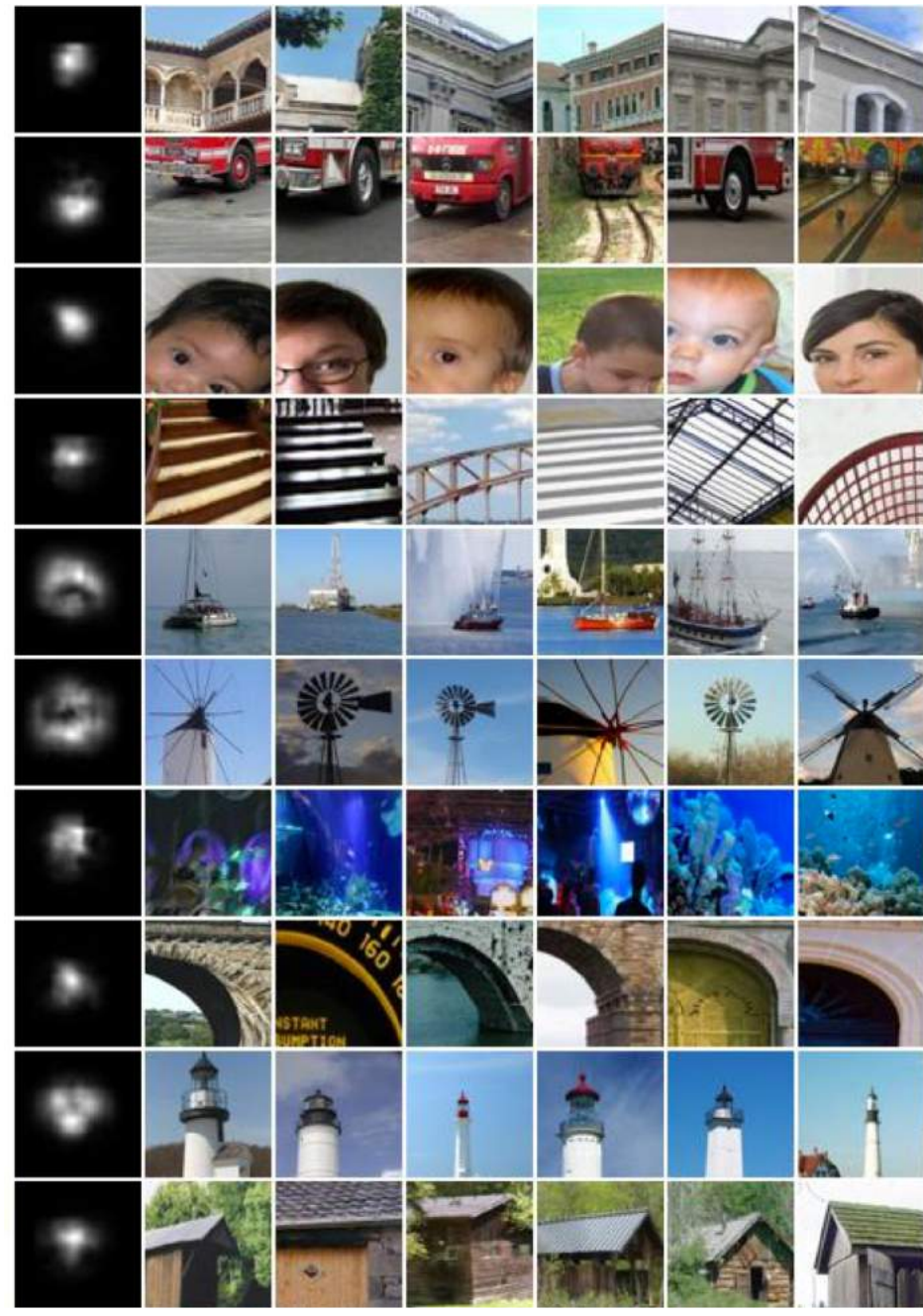
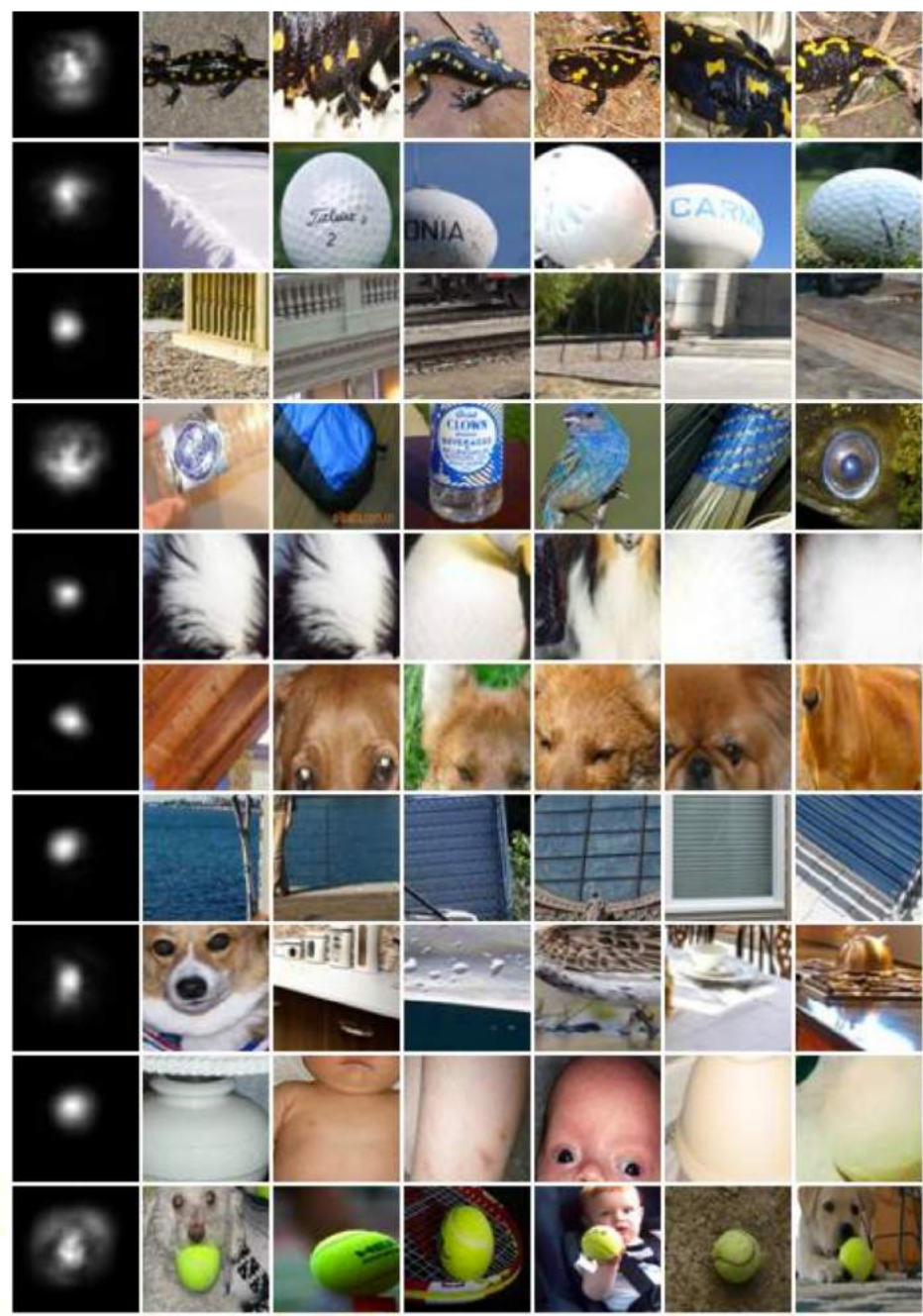
Layer 1

places 

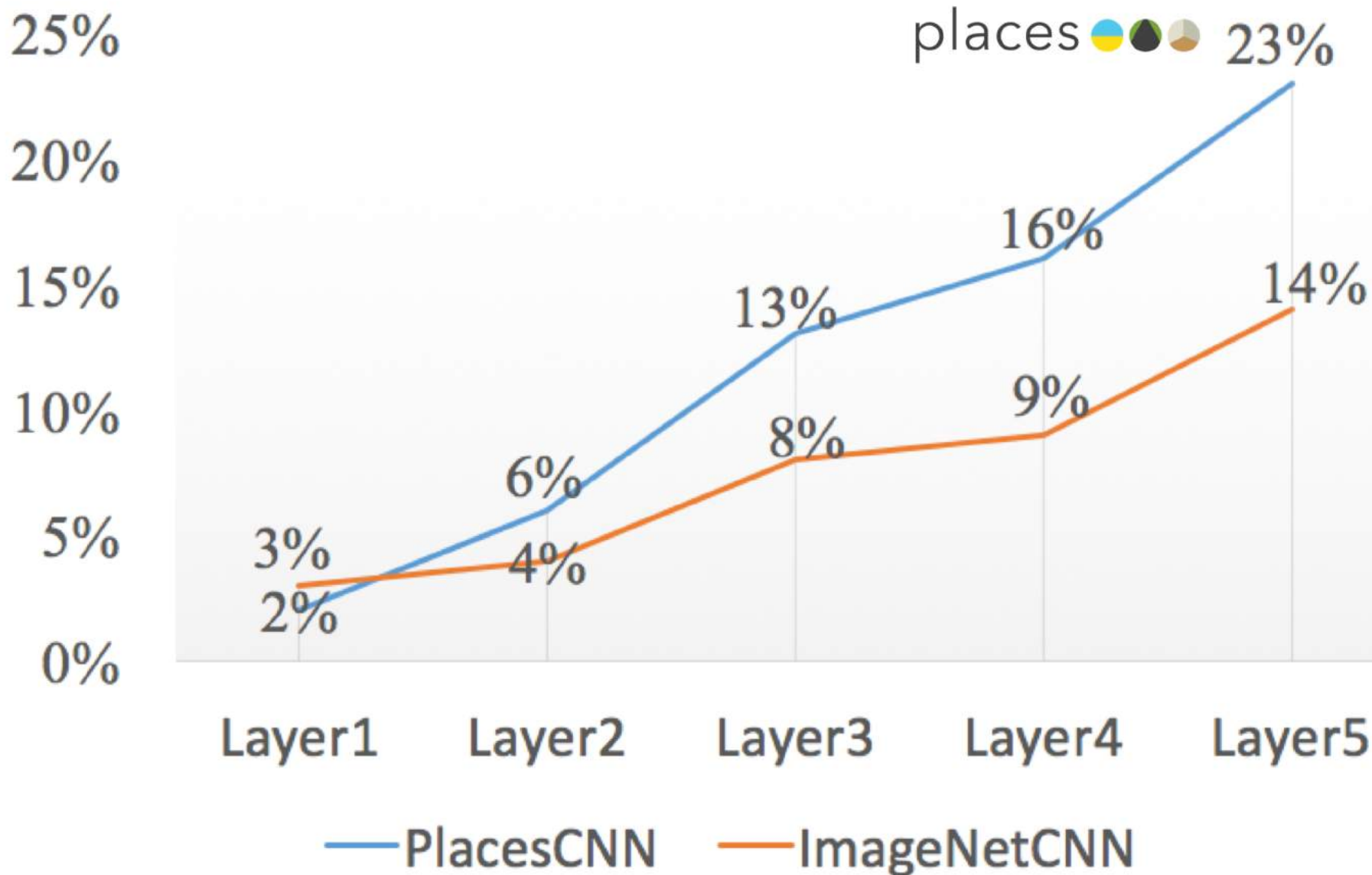








% Units as Detectors for Objects

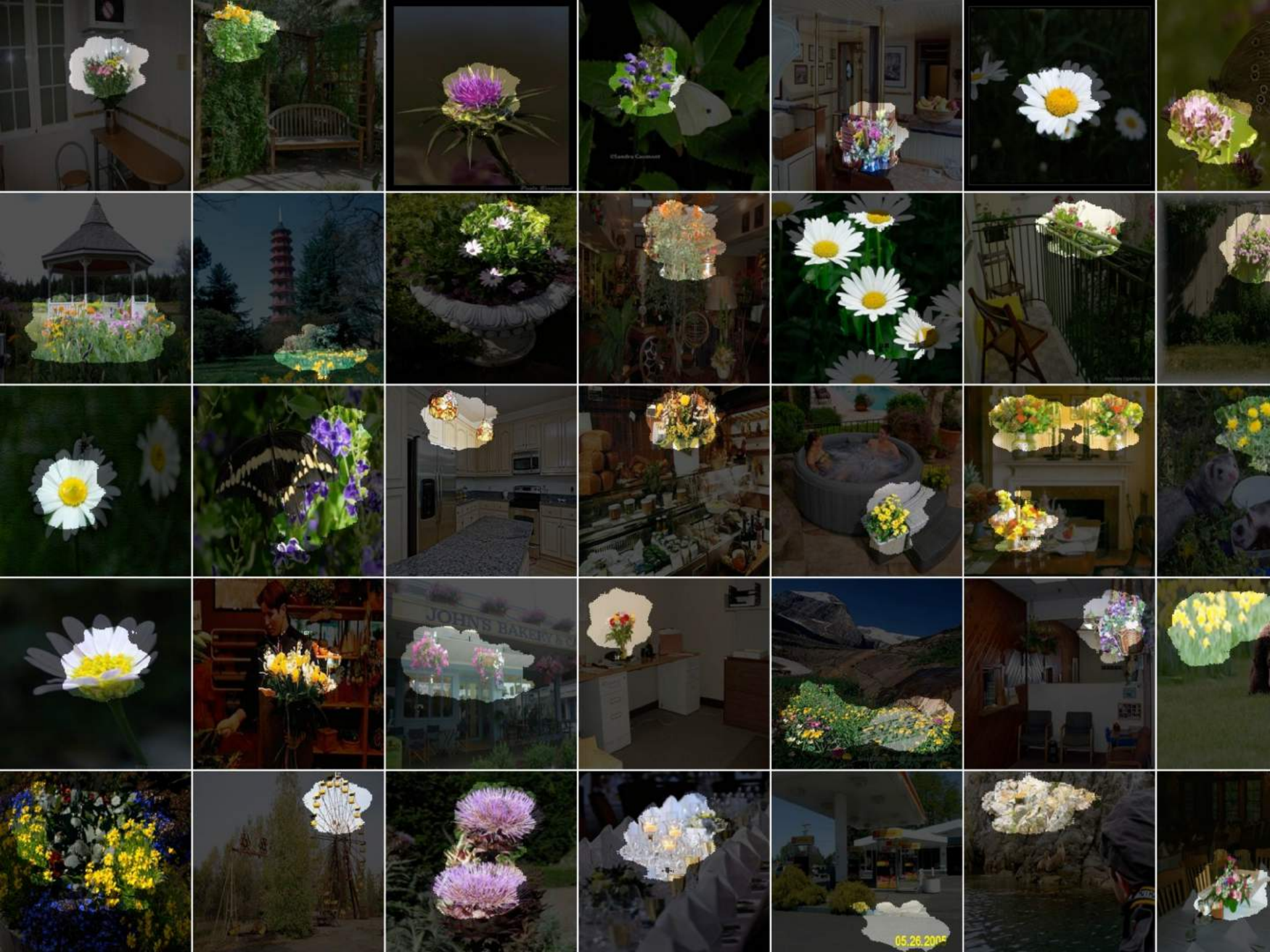




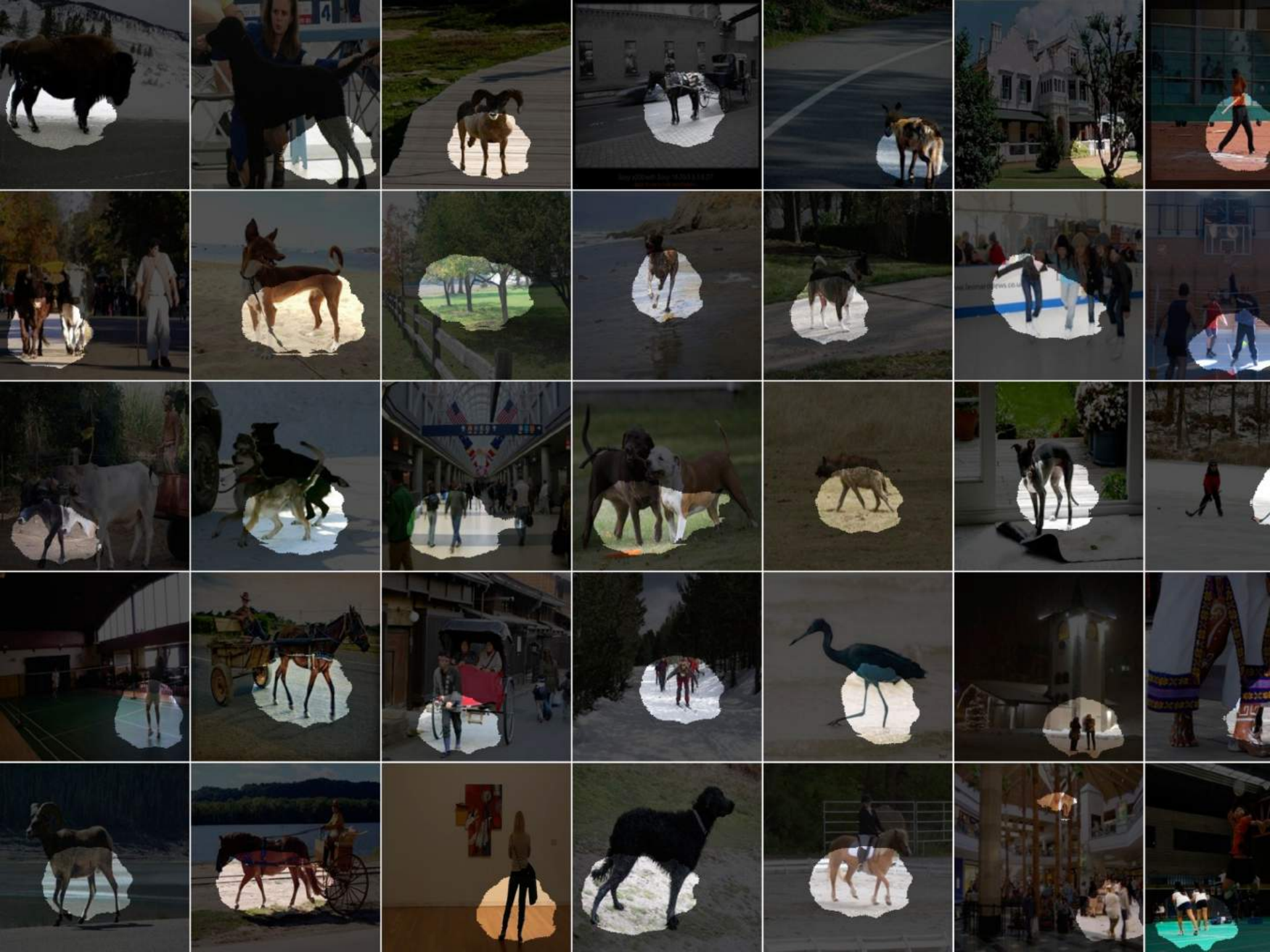


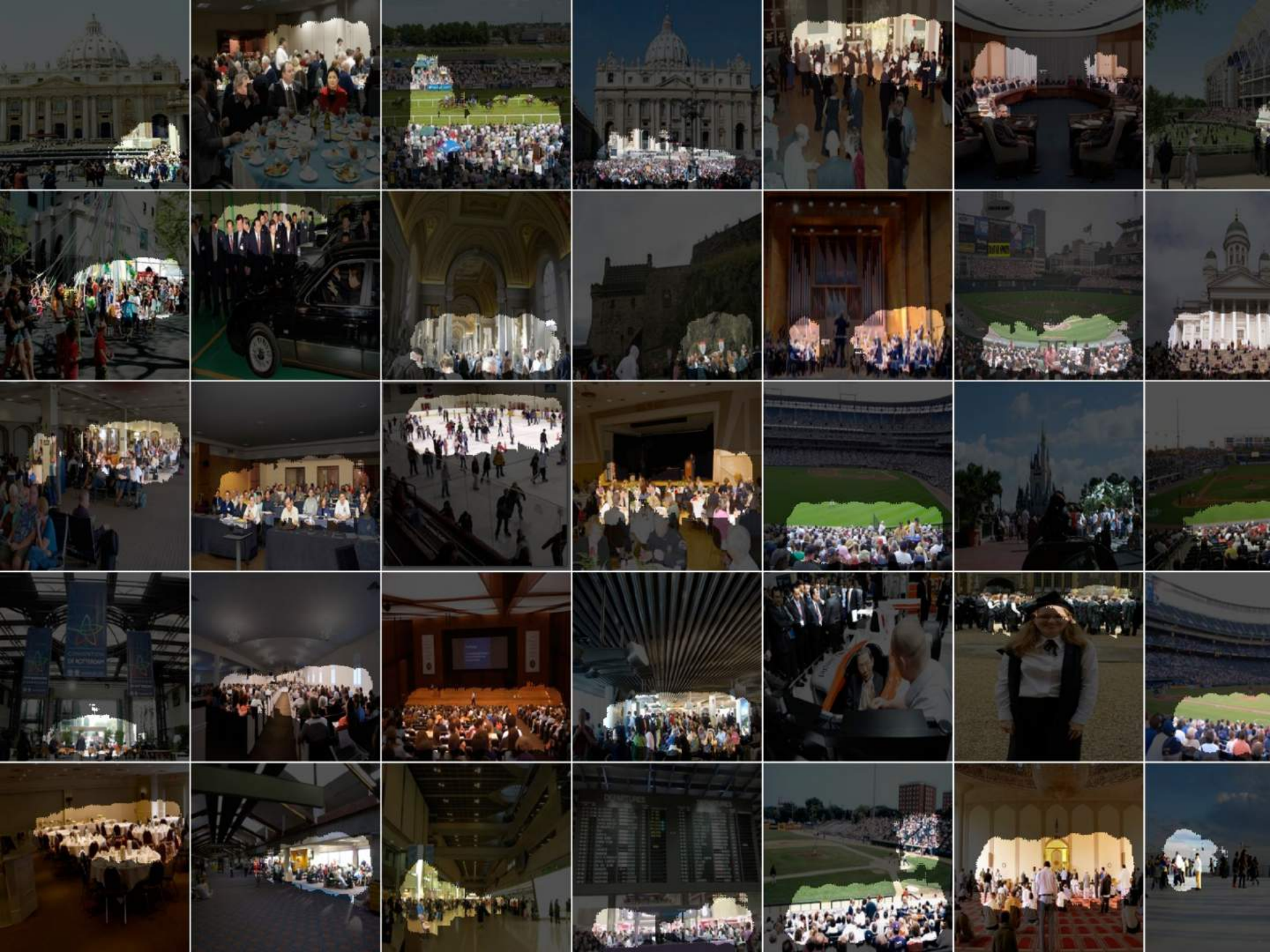


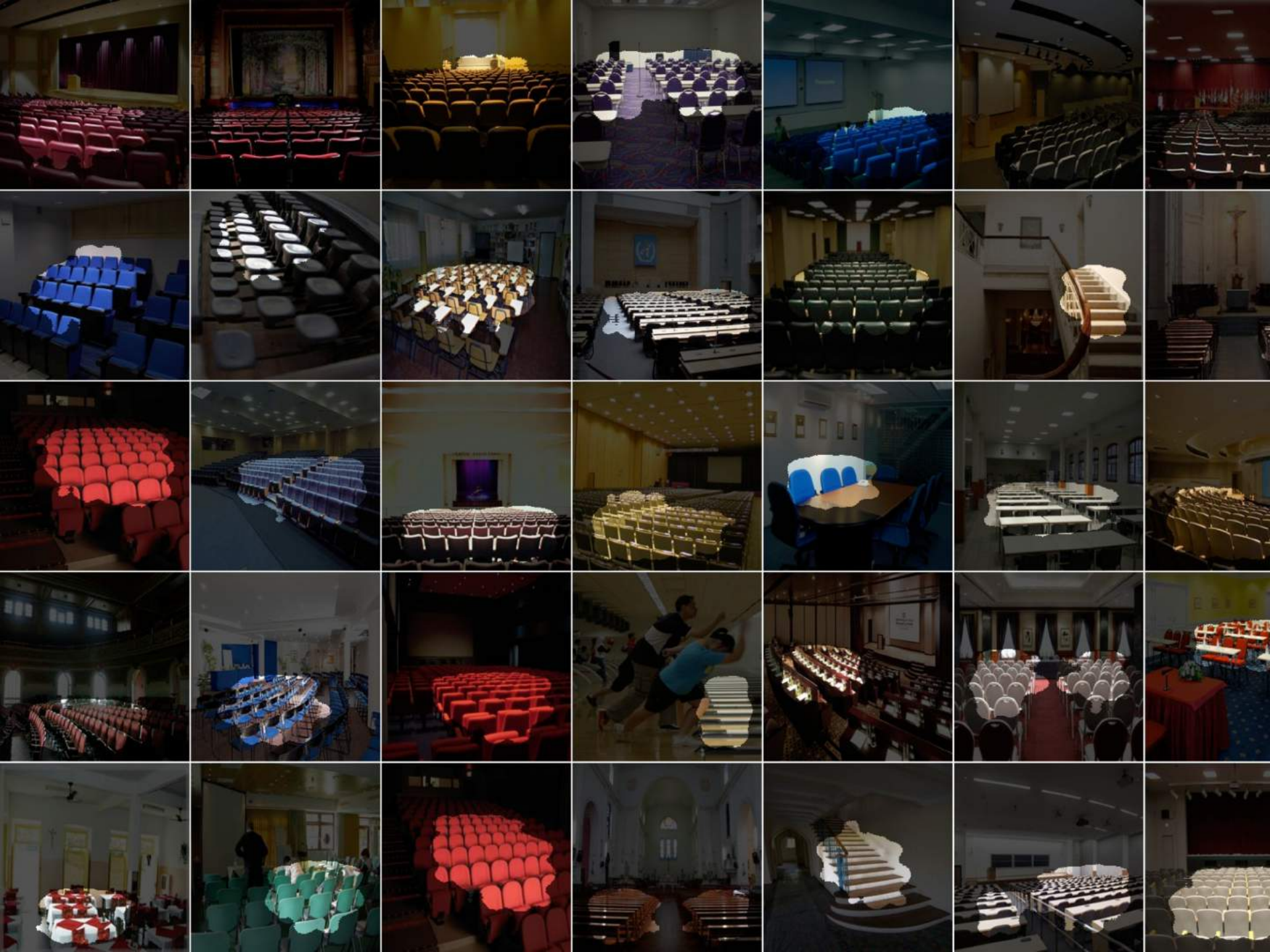












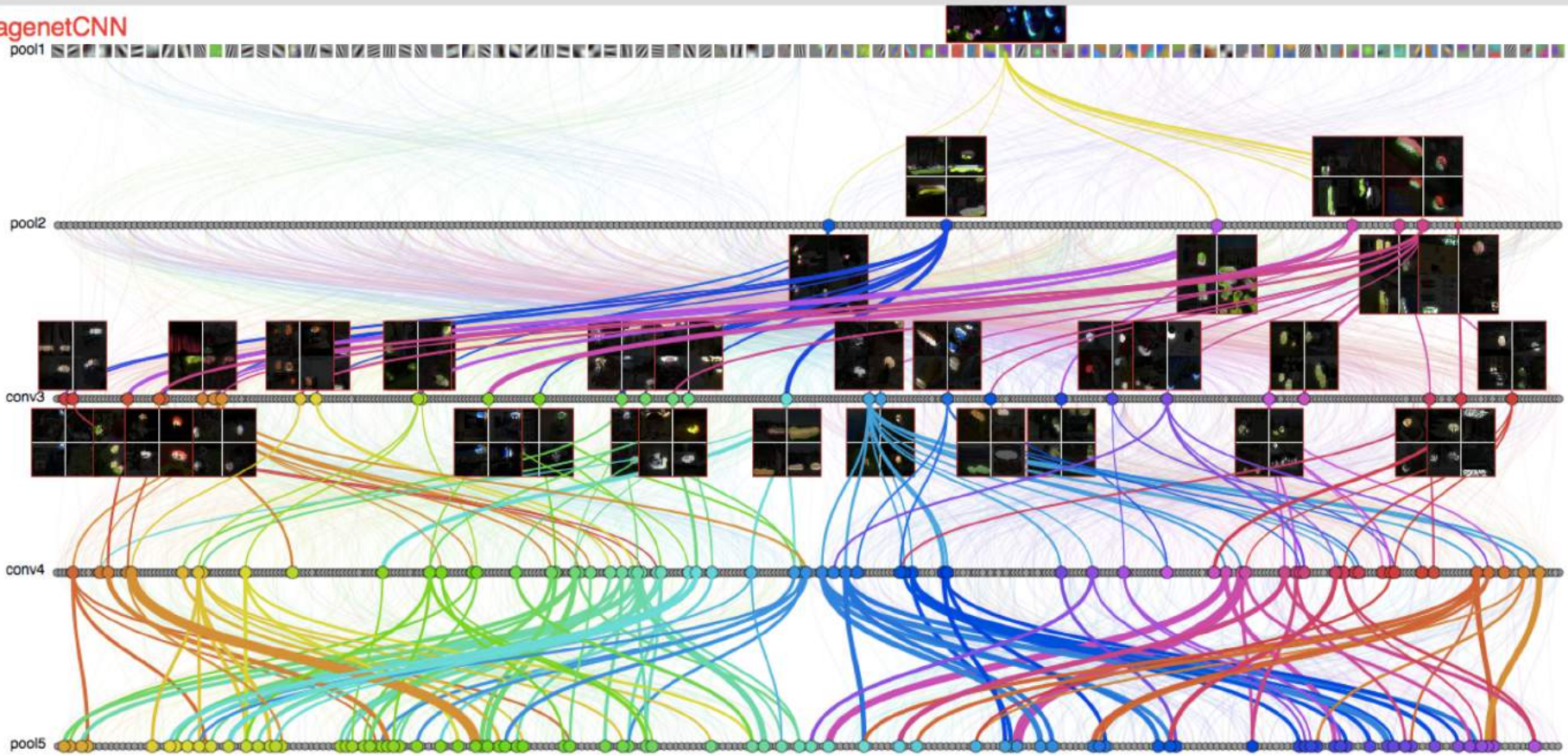




Visualizing Units & Connections

drawNet

imagenetCNN



<http://people.csail.mit.edu/torralba/research/drawCNN/drawNet.html>

Network Dissection: Quantifying Interpretability of Deep Visual Representations

David Bau*, Bolei Zhou*, Aditya Khosla, Aude Oliva, Antonio Torralba
Massachusetts Institute of Technology

<http://netdissect.csail.mit.edu/>



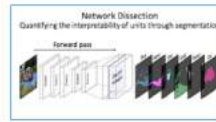
Code and Data



CVPR 2017 paper



CVPR 2017 poster



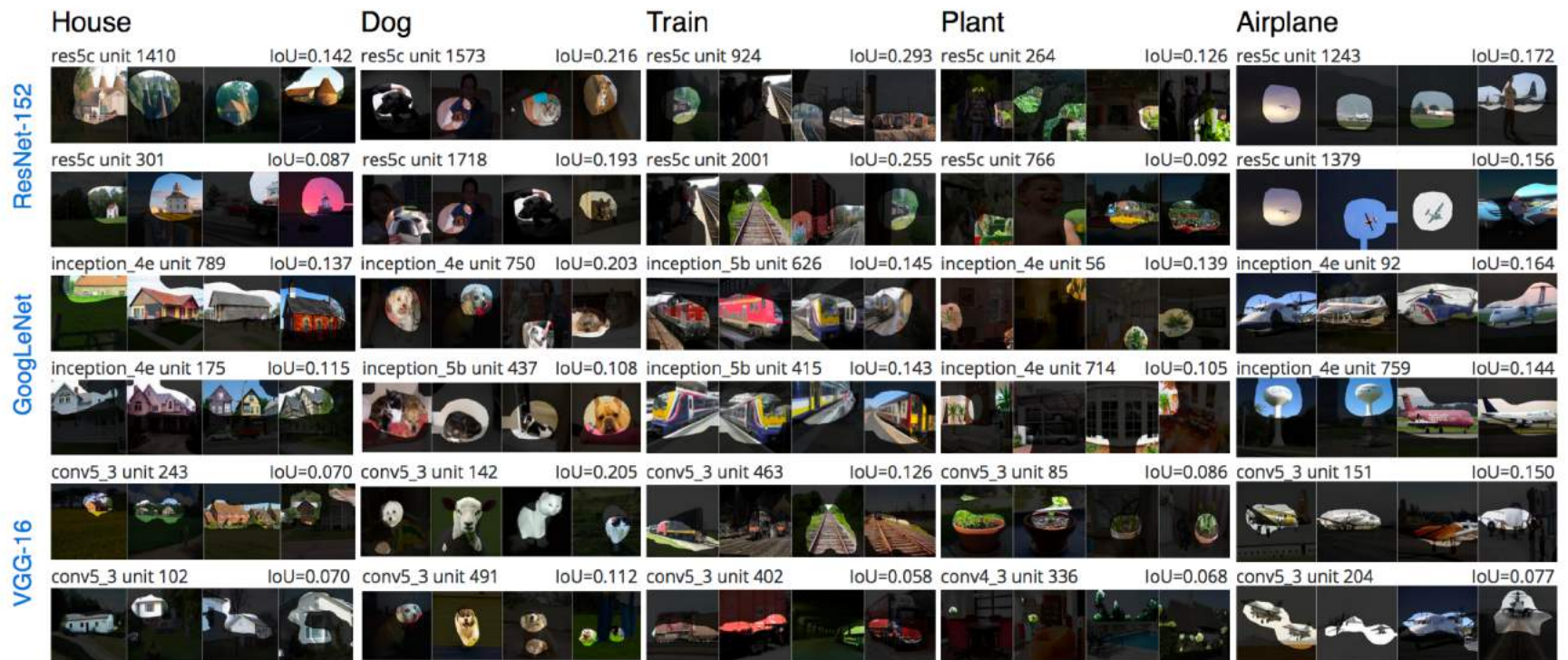
CVPR 2017 slides



CVPR 2017 oral

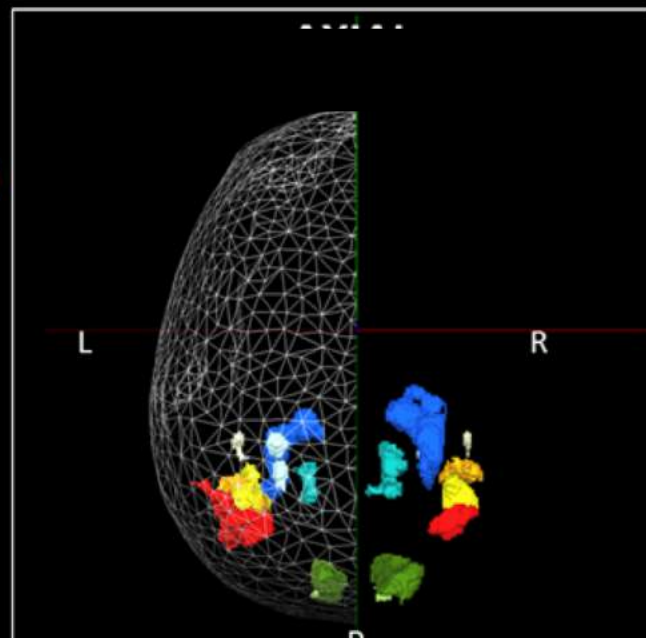
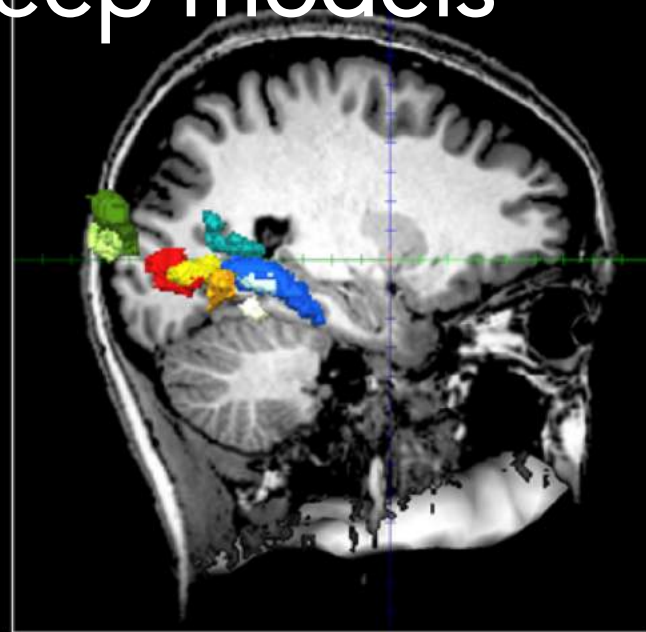
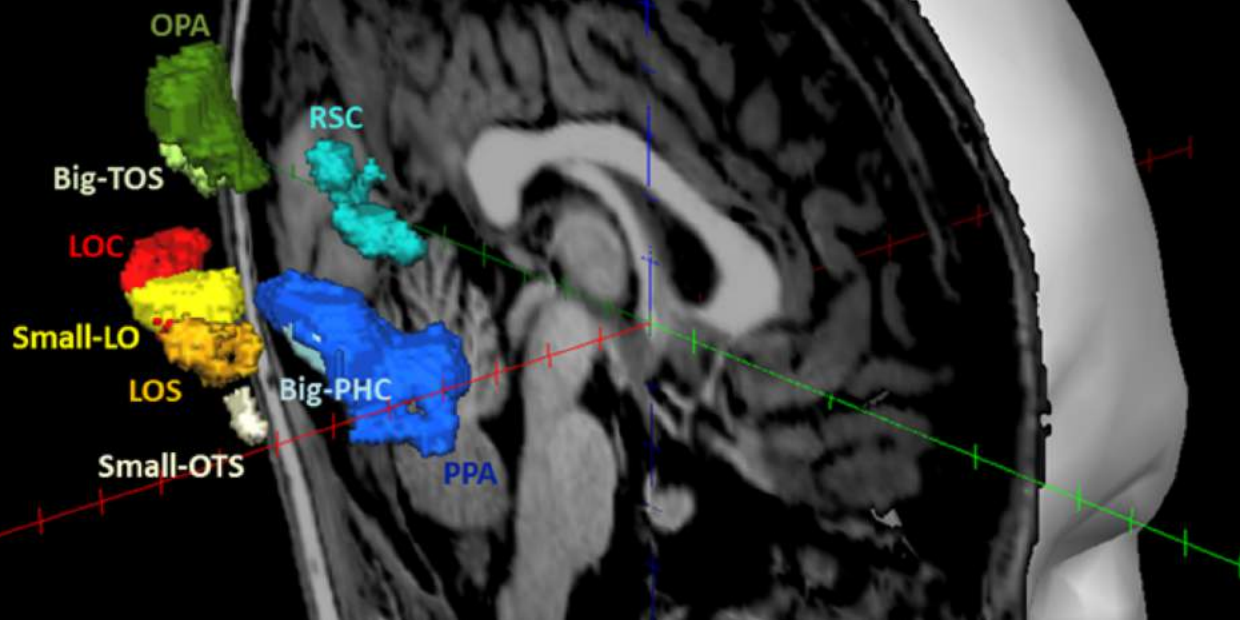


Extended paper



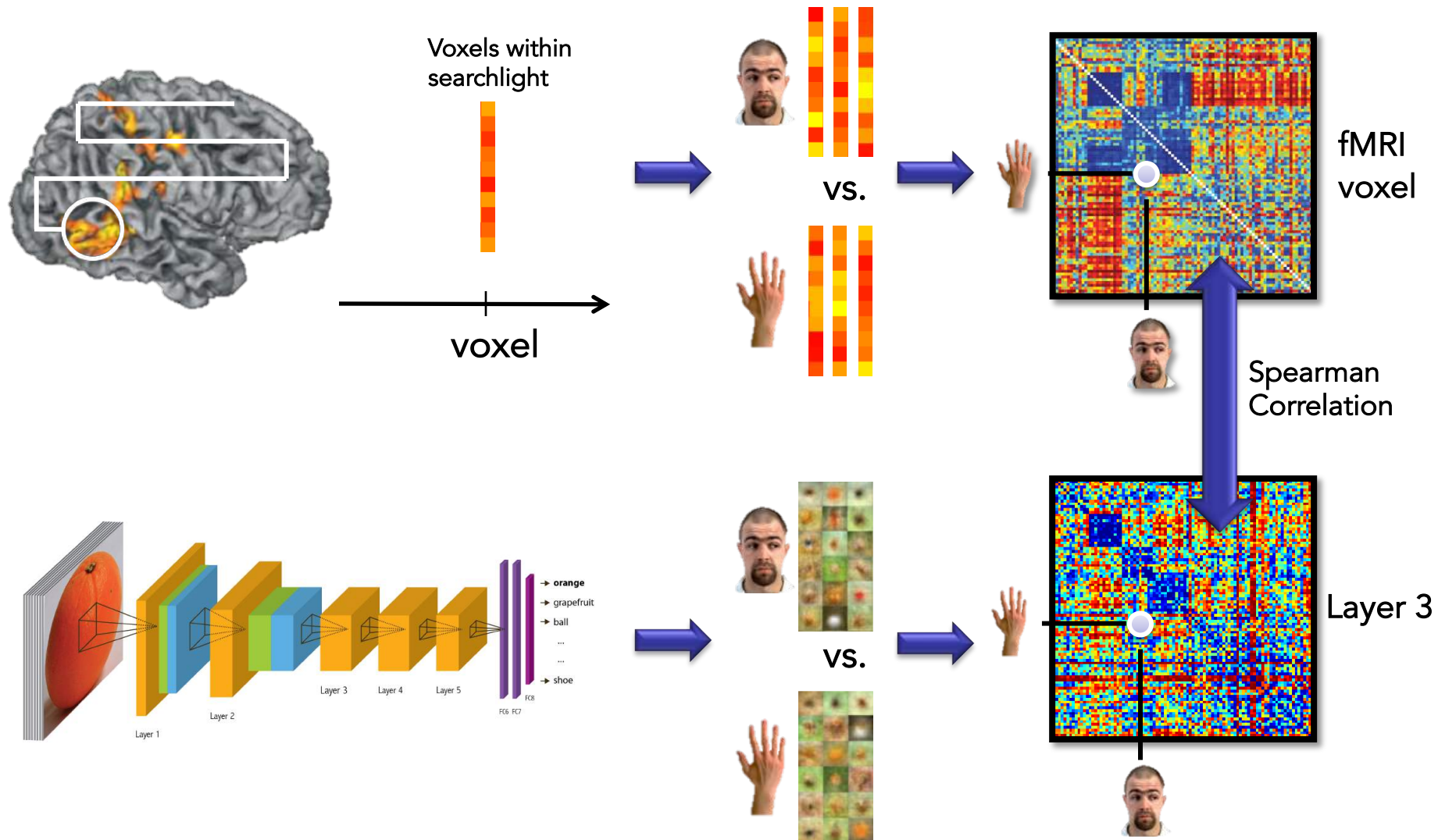
Selected units are shown from three state-of-the-art network architectures when trained to classify images of places (places-365). Many individual units respond to specific high-level concepts (object segmentations) that are not directly represented in the training set (scene classifications).

Correspondence between deep models and human brain ?



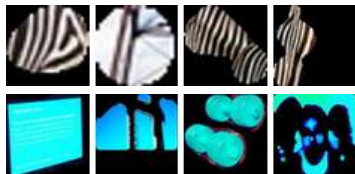
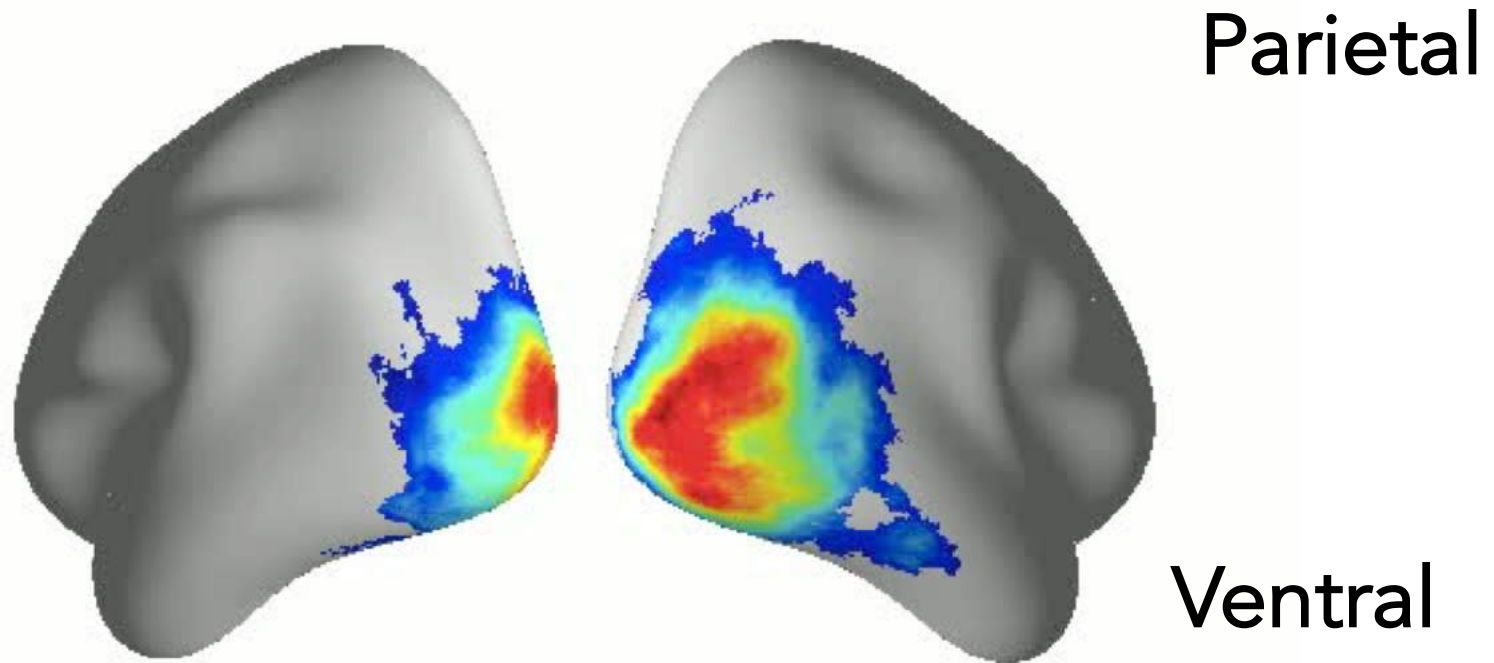
Algorithmic-specific fMRI searchlight analysis

A spatially unbiased view of the relations in similarity structure between models and fMRI



Spatiotemporal maps of correlations between human brain and model layers

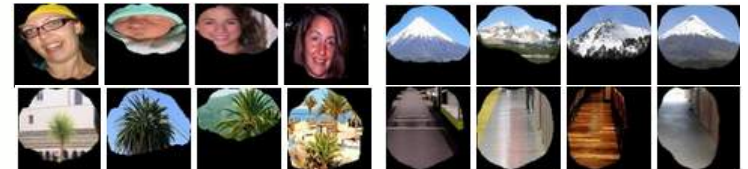
Layer 1



Layers 1-2



Layers 2-4



Layers 5-8

Comparing Natural and Artificial Deep Neural Networks

- **New fields of expertise:** Cognitive / Clinical / Social / Perceptual Computational **Experimentalist**
- **Studying the implementation** that works best for performing specific tasks
- **Characterizing the network behavior** when it is adapting, compromised or enhanced
- **Exploring the alternatives** that have not been taken by biological systems



Bolei
Zhou



David
Bau



Aditya
Khosla



Radoslaw
Cichy



Antonio
Torralba