MIT CSAIL

# Interpreting Deep Visual Representations
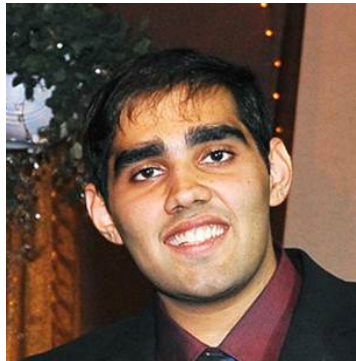
## Bolei Zhou

## MIT

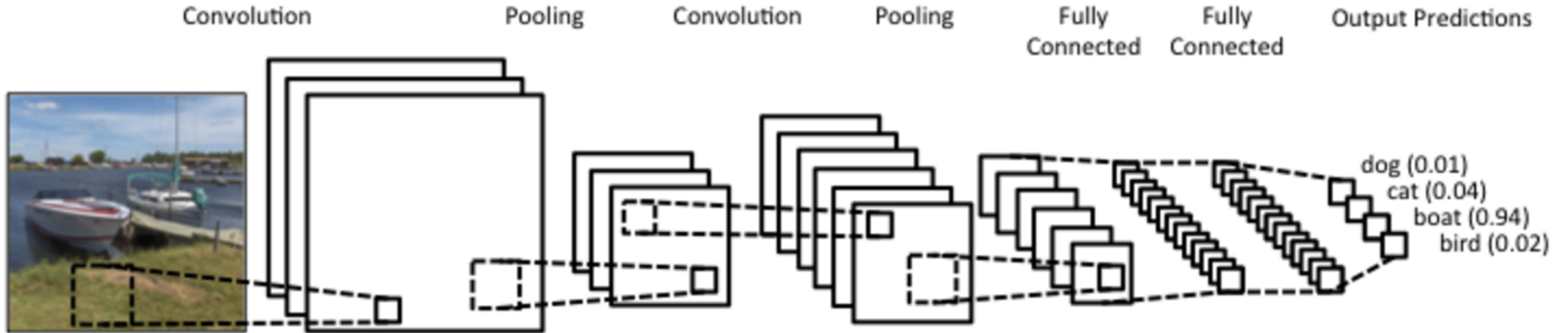David Bau  Aditya Khosla  Aude Oliva  Antonio Torralba
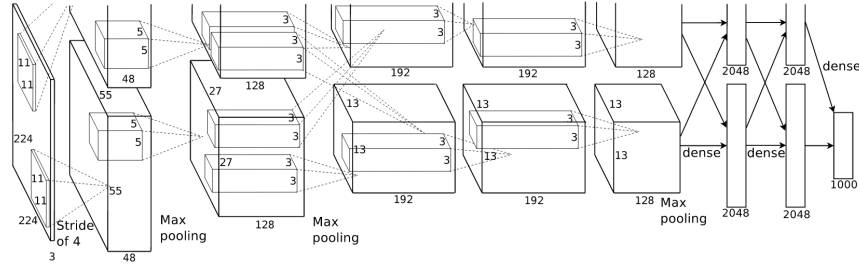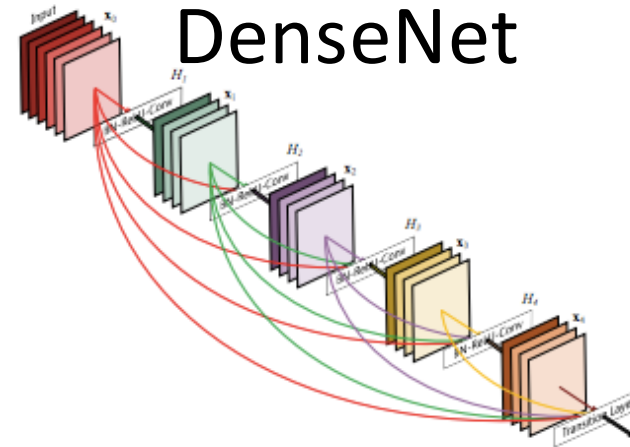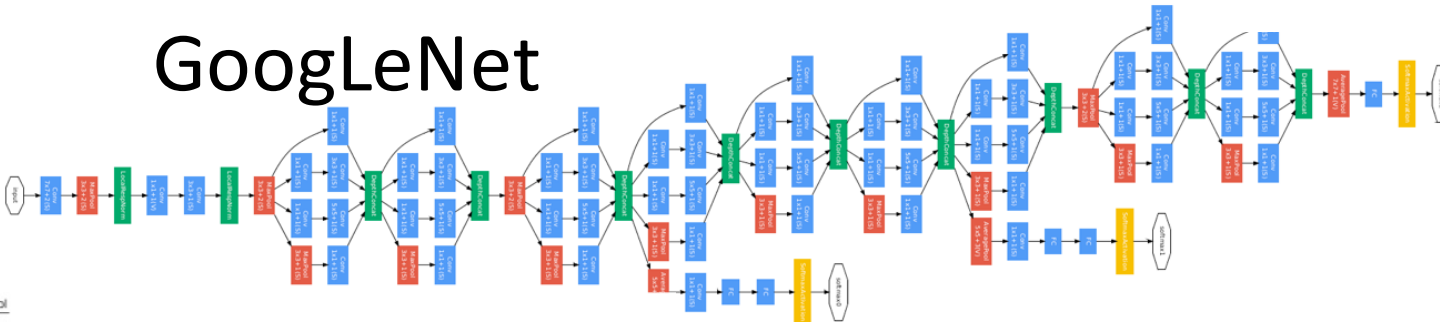
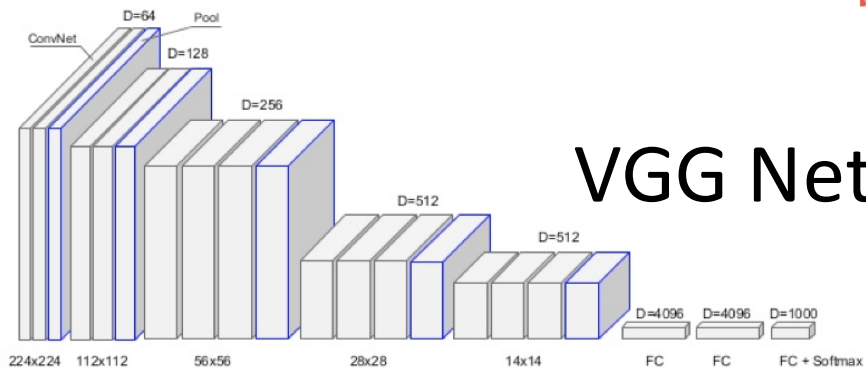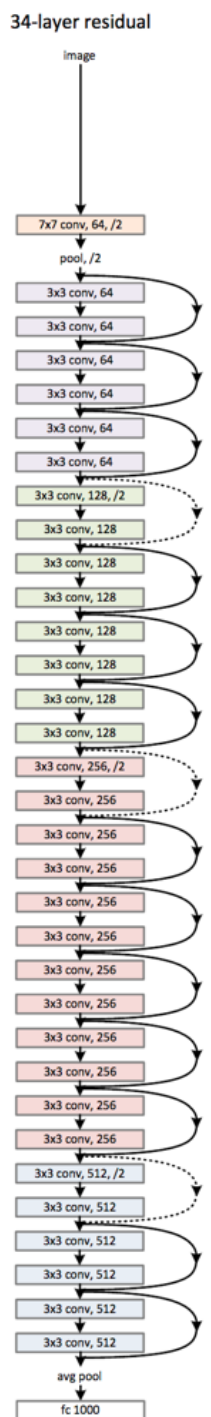# Background

Convolutional Neural Network (ConvNet)

# Many networks

AlexNet

DenseNet

GoogLeNet

VGG Net

ResNet

# Why works so well



Upload your image for scene recognition using **Places-CNN** from MIT.

**Take/Choose a photo**

**Predictions**:

- **type:** indoor
- **semantic categories:**
  hotel_room:0.50, bedroom:0.47,

**Predictions**:

- **type:** indoor
- **semantic categories:**
  hotel_room:0.35, bedroom:0.15,
  living_room:0.09, dorm_room:0.06,
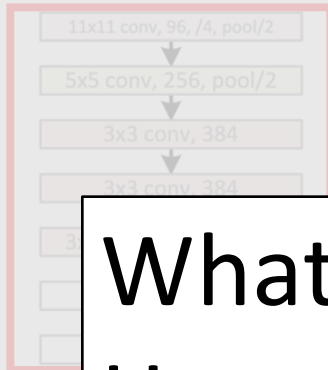  basement:0.05

# When it fails, why is it?



Output: cutting vegetables.
Correct label: gardening



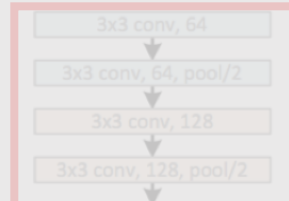Output: washing dishes.
Correct label: brushing

# Deep ConvNet for Visual Recognition
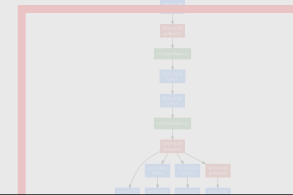
**2012: AlexNet**
5 conv. layers

**2014: VGG**
16 conv. layers

**2015: GoogLeNet**
22 conv. layers

**2016: ResNet**
>100 conv. layers

11x11 conv, 96, /4, pool/2

5x5 conv, 256, pool/2

3x3 conv, 384

3x3 conv, 384

3x3 conv, 64

3x3 conv, 64, pool/2

3x3 conv, 128

3x3 conv, 128, pool/2

## What have been learned inside?
## How to compare the internal representations?

Error: 15.3%

3x3 conv, 512

3x3 conv, 512, pool/2

3x3 conv, 512

3x3 conv, 512

3x3 conv, 512

3x3 conv, 512, pool/2

fc, 4096

fc, 4096

fc, 1000

Error: 8.5%

Error: 7.8%

Error: 4.4%

# Work on Network Visualization

## Deconvolution



Layer 2

Layer 5

Zeiler et al., ECCV 2014.

## Back-propagation



bell pepper　　　lemon　　　husky

Simonyan et al., ICLR 2015



Horizon　　　Trees　　　Leaves

Towers & Pagodas　　Buildings　　Birds & Insects

Inceptionism. Google Blog. June 2015

## Feature inversion



norm1　　conv2　　relu2

relu5　　mpool5　　fc6

Mahendran et al, CVPR 2015
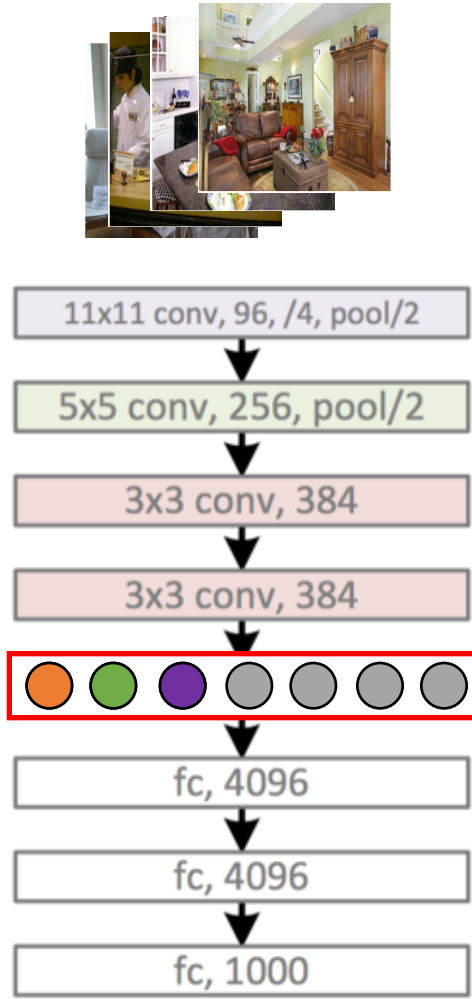
## Top activated images



Girshick et al., CVPR 2014

# Going From Visualization to Interpretation



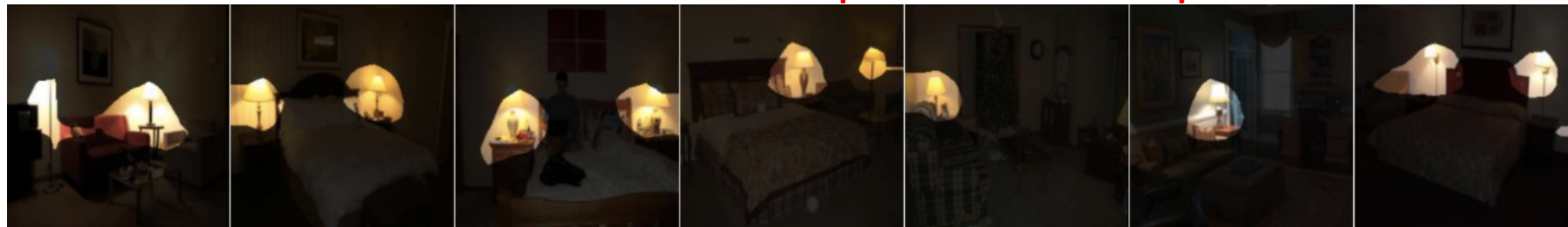Top Activated Images     Interpretation: head     Score: 0.23

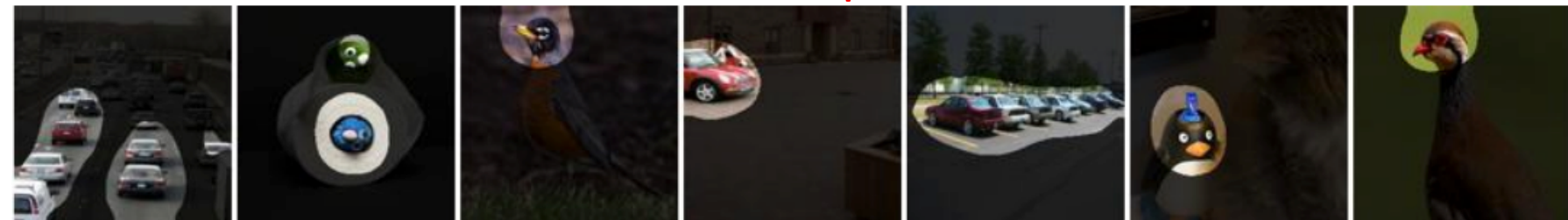Top Activated Images     Interpretation: lamp     Score: 0.15

Top Activated Images     Interpretation: car     Score: 0.02

Unit 1

Unit 2

Unit 3

11x11 conv, 96, /4, pool/2

5x5 conv, 256, pool/2

3x3 conv, 384

3x3 conv, 384

fc, 4096

fc, 4096

fc, 1000

# Solution: Evaluate units for semantic segmentation

Unit 1            Top activated images



Lamp,  Intersection over Union (IoU)= 0.12

$$IoU = \frac{Area\ of\ Overlap}{Area\ of\ Union}$$

Bau*, Zhou*, Khosla, Oliva, Torralba. Network Dissection: quantifying interpretability of deep visual representations. CVPR'17

# Network Dissection
## Framework to interpret the deep visual representations

Bau*, Zhou*, Khosla, Oliva, Torralba. Network Dissection: quantifying interpretability of deep visual representations. CVPR'17

# Broadly and Densely (**Broden**) Annotated Dataset

**ADE20K**

    Zhou et al, CVPR'17

**Pascal Context**

    Mottaghi et al, CVPR'14

**Pascal Part**

    Chen et al, CVPR'14

**Open-Surfaces**

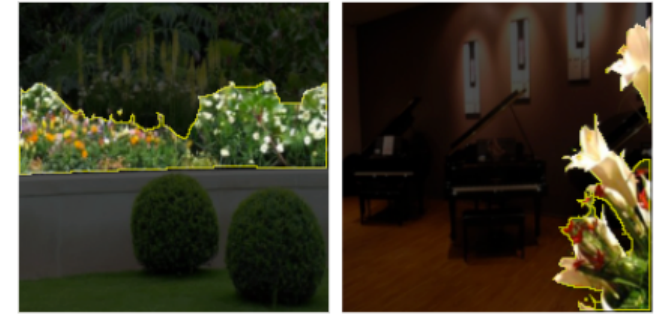    Bell et al, SIGGRAPH'14

**Describable Textures**

    Cimpoi et al, CVPR'14

**Colors**

Total = 63,305 images

    1,197 visual concepts

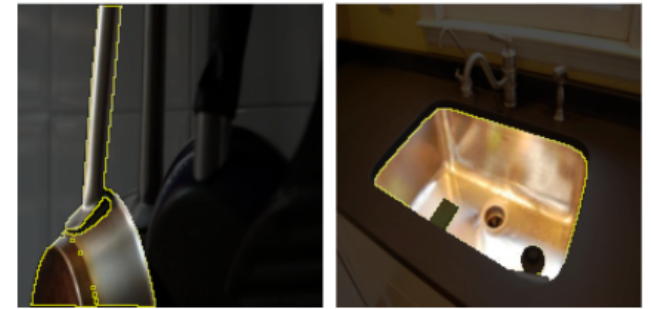

street (scene)

flower (object)
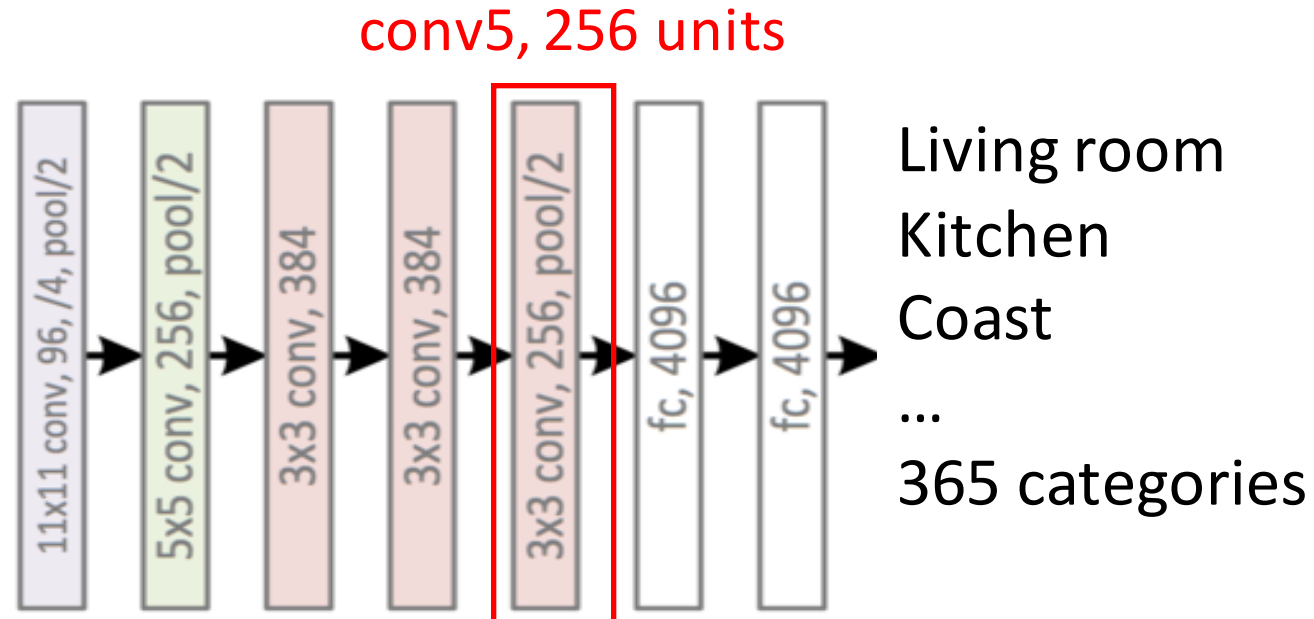
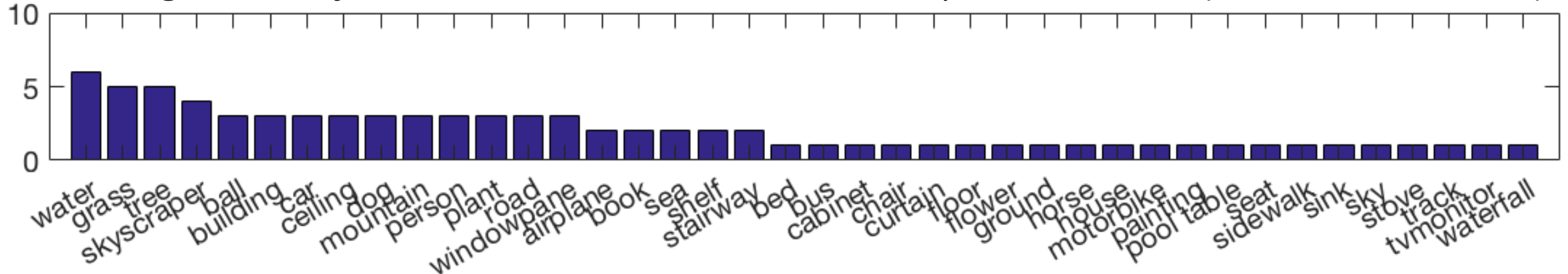headboard (part)

metal (material)

swirly (texture)

pink (color)

# Result of AlexNet trained on places THE SCENE RECOGNITION DATABASE

**conv5, 256 units**



| 11x11 conv, 96, /4, pool/2 | 5x5 conv, 256, pool/2 | 3x3 conv, 384 | 3x3 conv, 384 | 3x3 conv, 256, pool/2 | fc, 4096 | fc, 4096 |

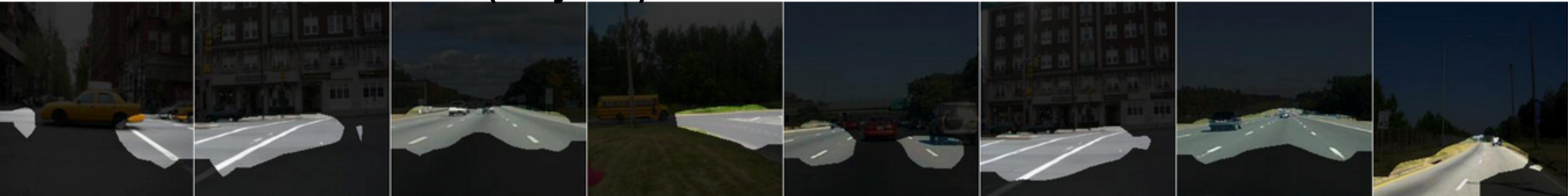Living room
Kitchen
Coast
...
365 categories

**Histogram of object detectors:** Detector:81/256,  Unique Detector:40 (Units with IoU>0.04)

conv5 unit 79  car (object)  IoU=0.13

conv5 unit 107  road (object)  IoU=0.15

**Histogram of object detectors:** Detector:81/256, Unique Detector:40 (Units with IoU>0.04)

conv5 unit 144    mountain (object)    IoU=0.13
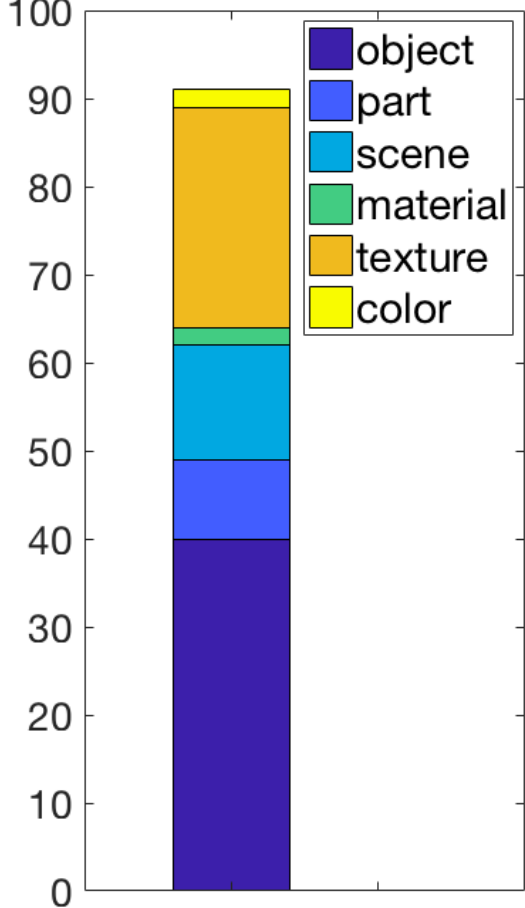
conv5 unit 200    mountain (object)    IoU=0.11

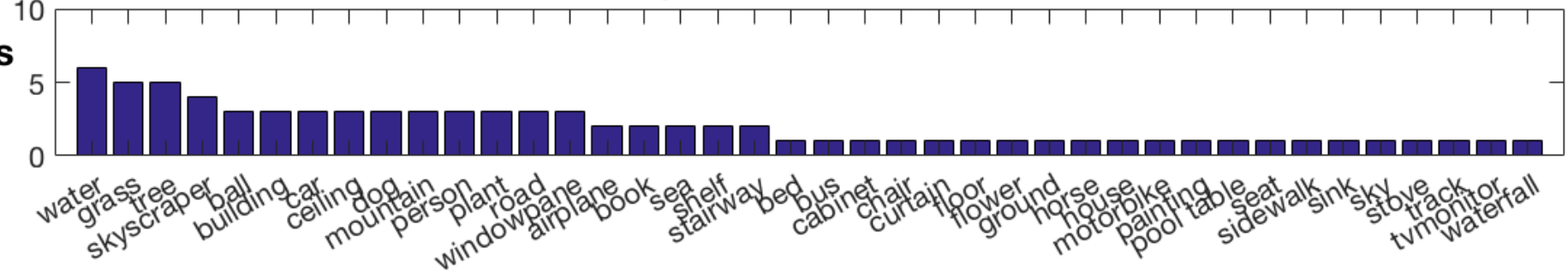**Histogram of object detectors:** Detector:81/256,  Unique Detector:40 (Units with IoU>0.04)
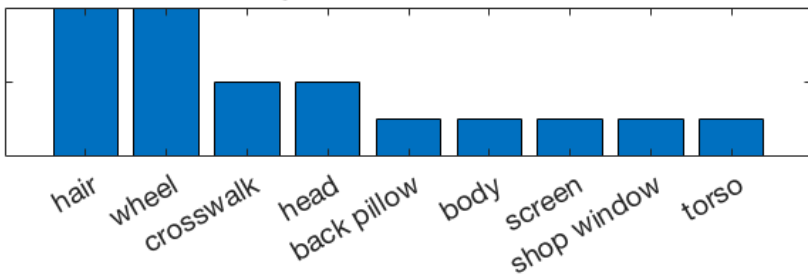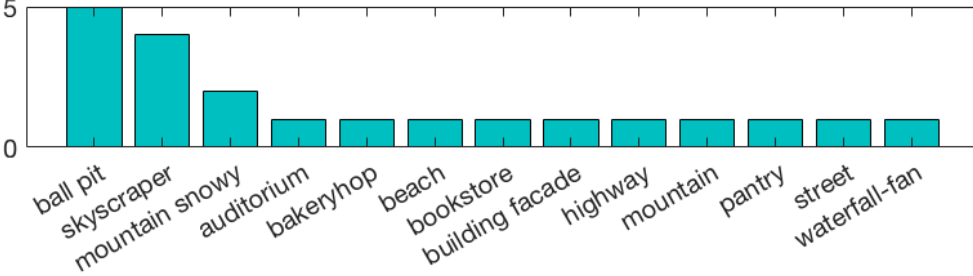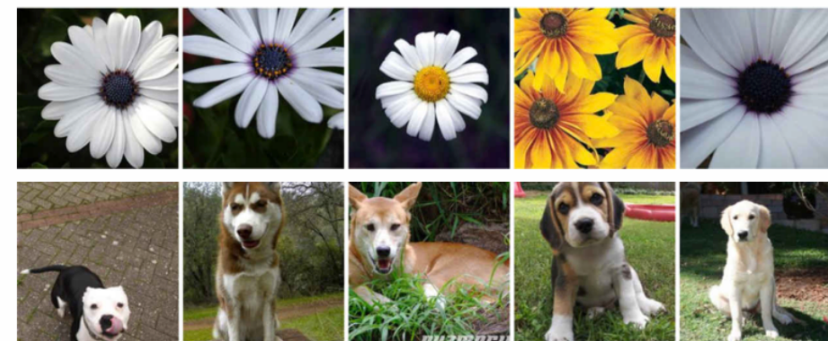
# Dissection Report

# Are the emerging concepts real?

**Szegedy et al. Intriguing properties of neural networks. arXiv.2014**

- "No distinction between individual high level units and random linear combinations of high level unit"
- "It suggests that it is the space, rather than the individual units, that contains the semantic information in network"
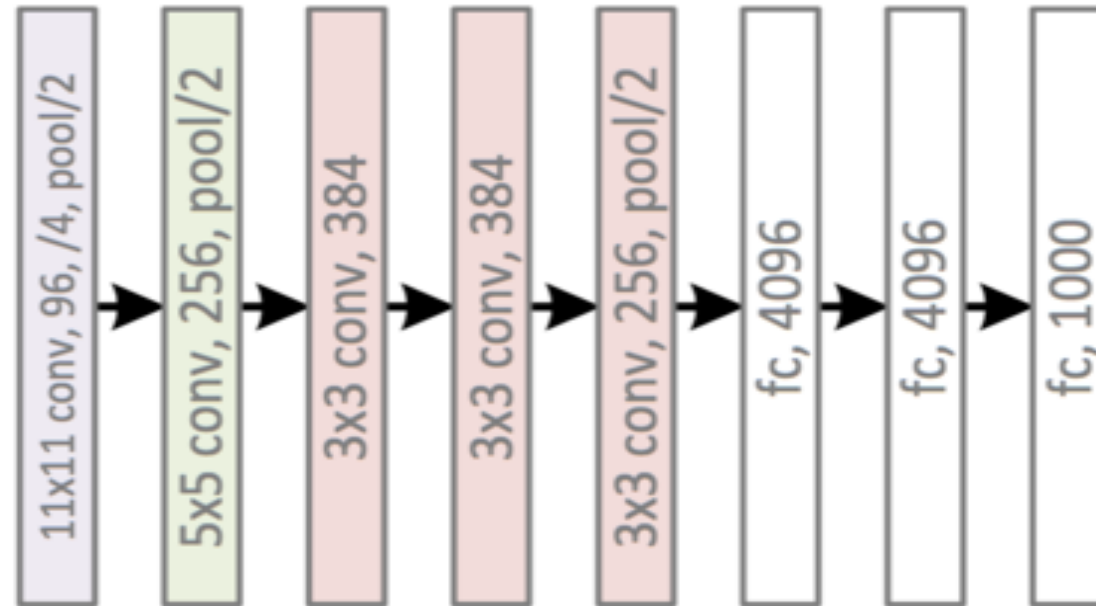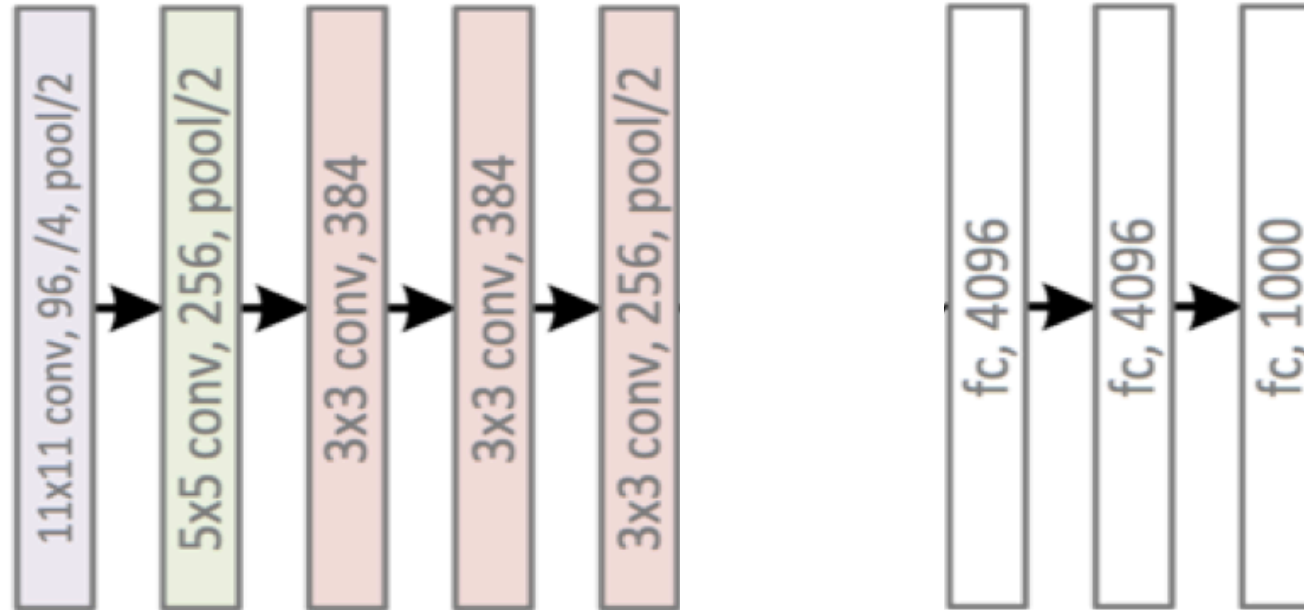


Single Neuron
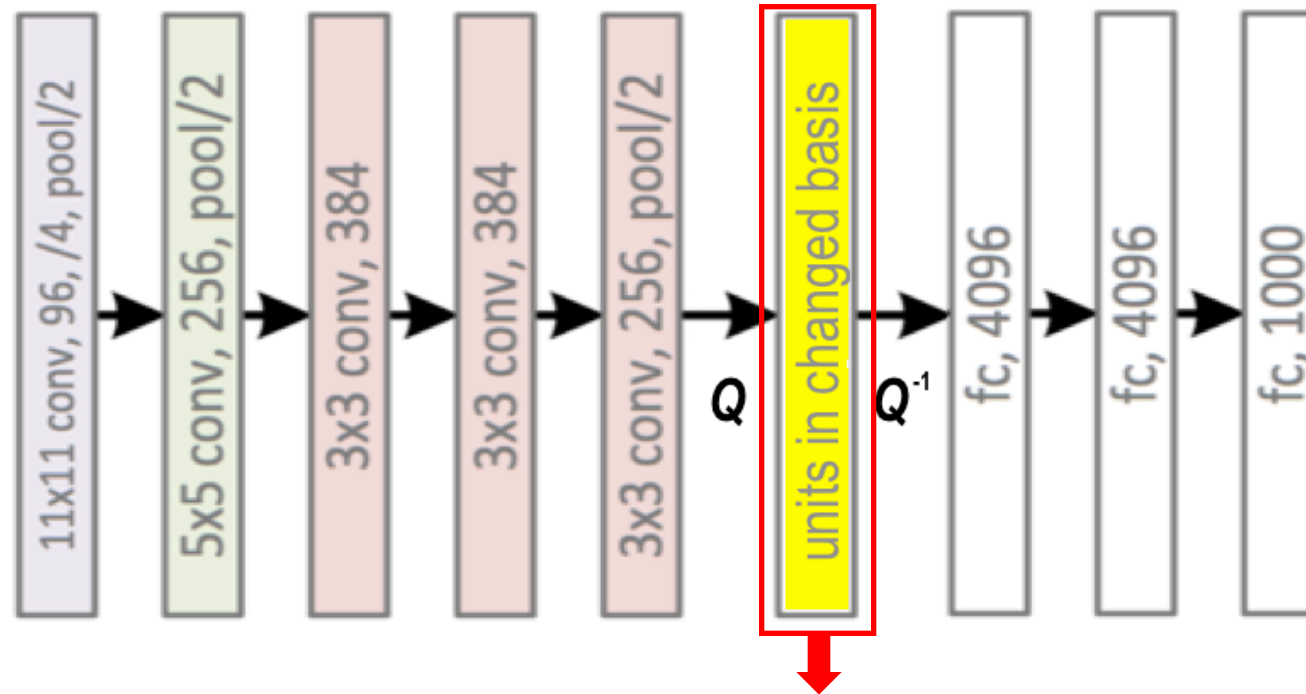
Random Projection

# Are the emerging concepts real?

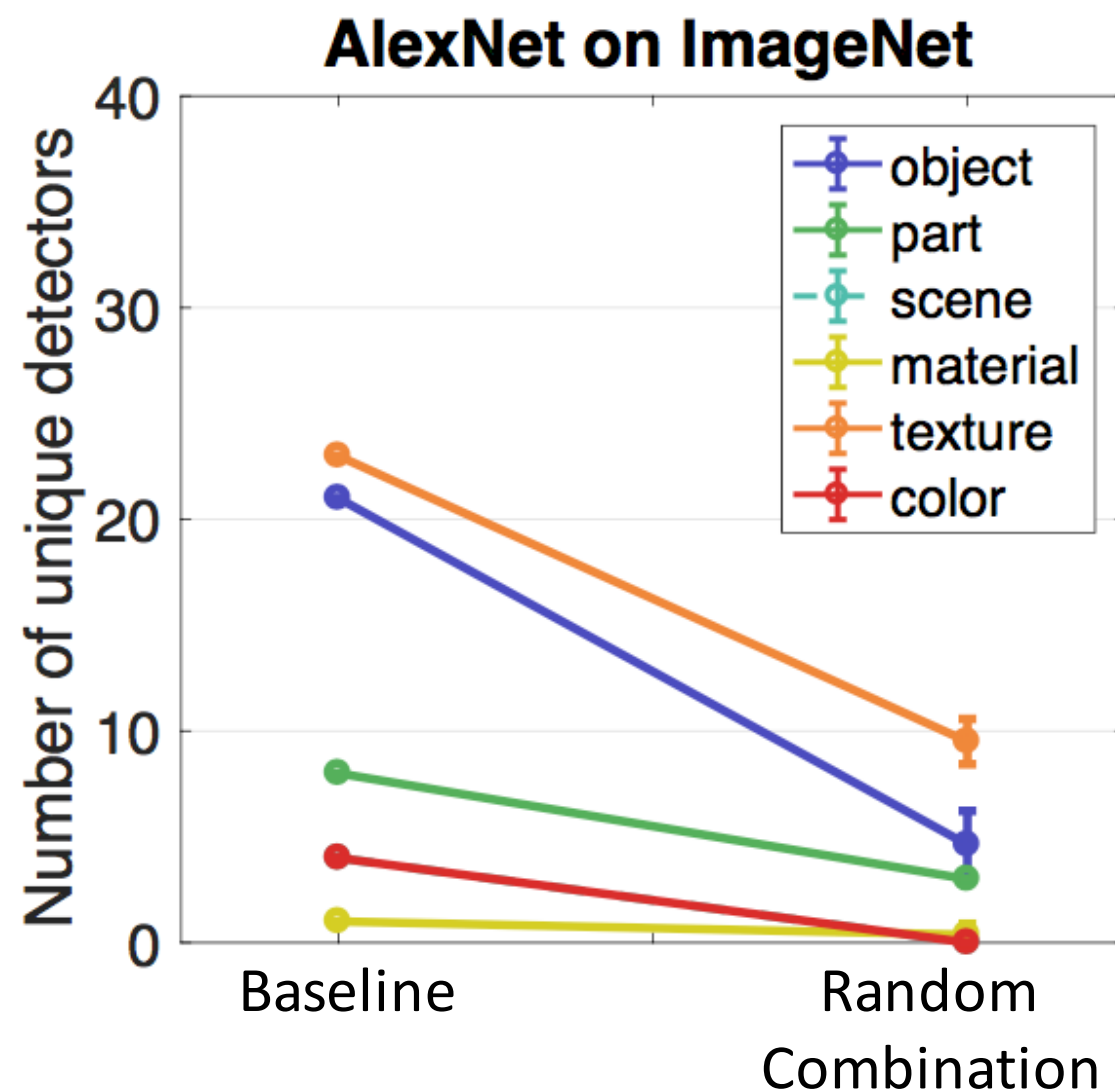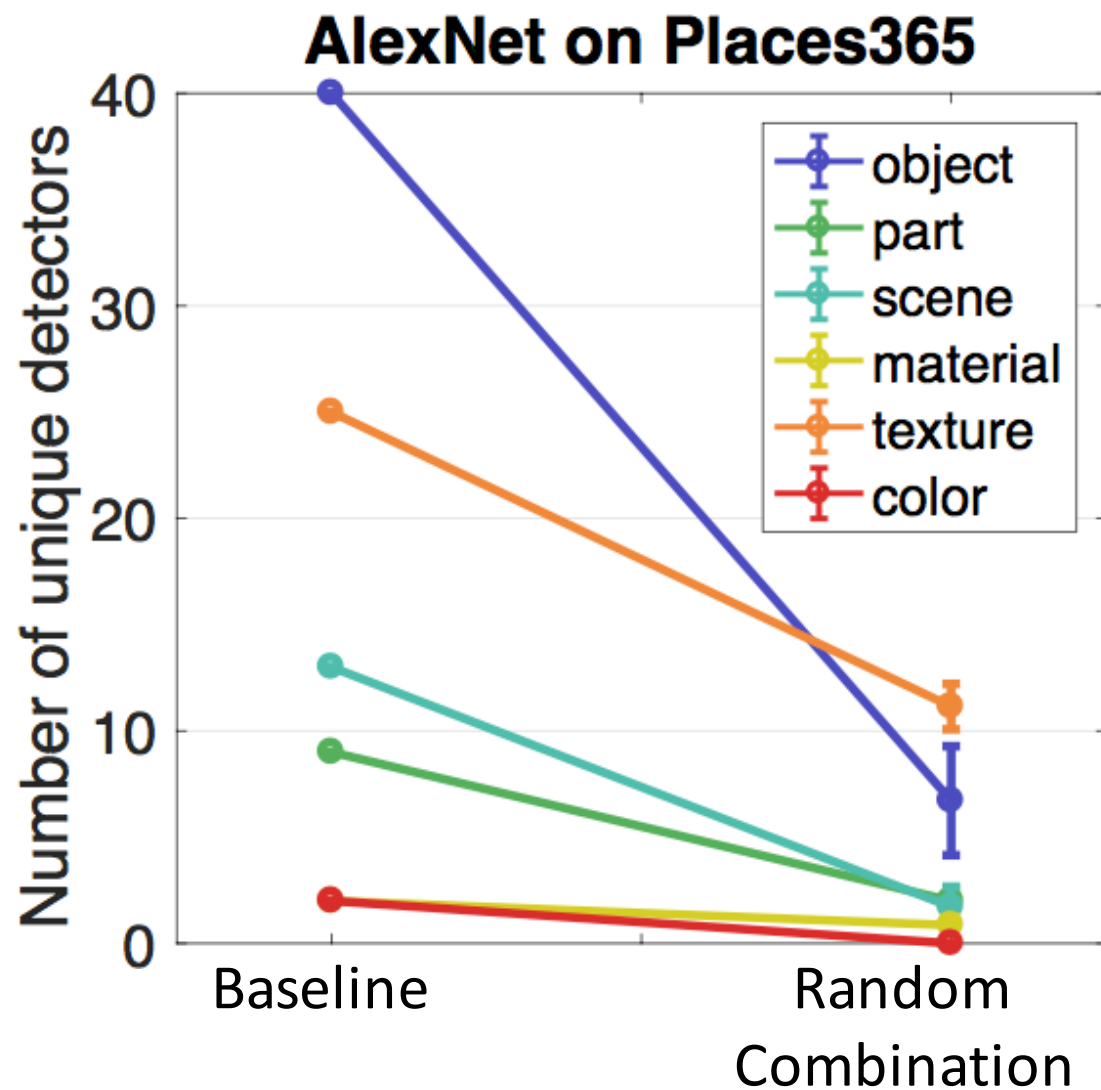# Are the emerging concepts real?

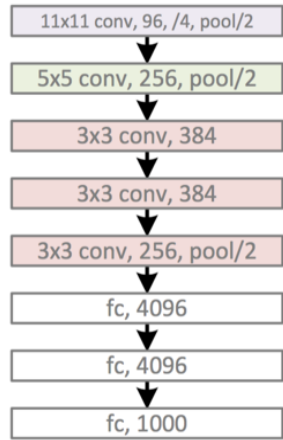# Are the emerging concepts real?

## Random combination of units



Do concepts associate with individual units or the whole feature space?
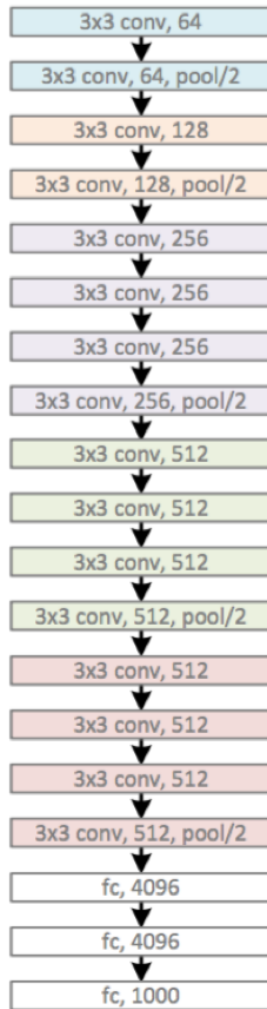
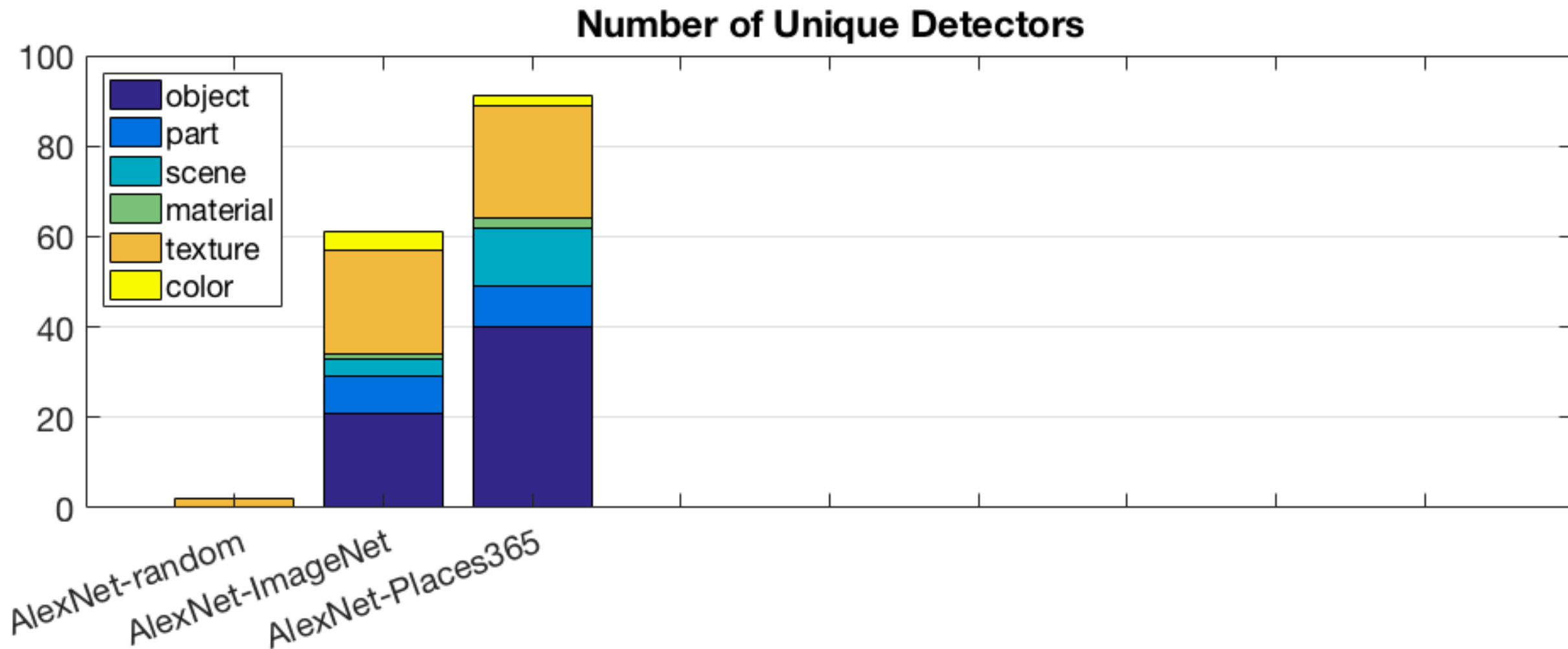# Are the emerging concepts real?

# Architectures

## Datasets



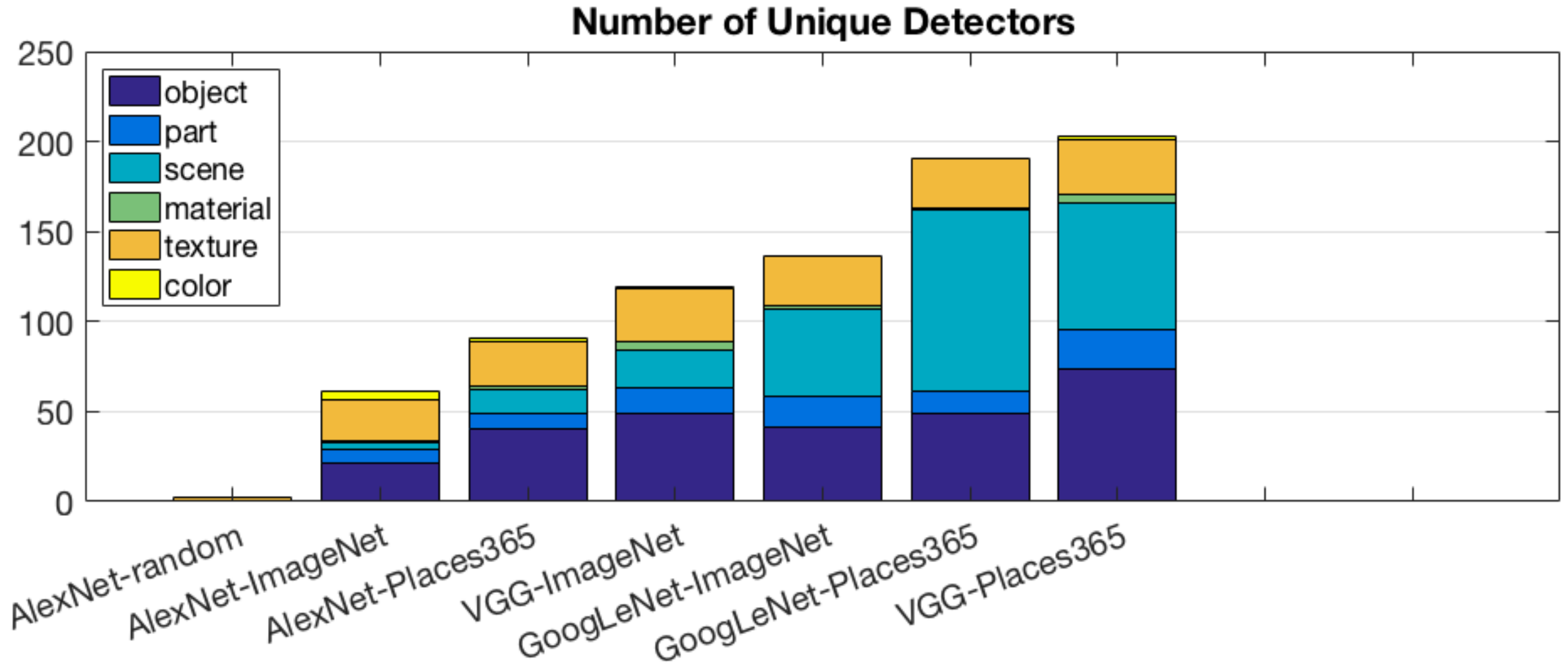| 11x11 conv, 96, /4, pool/2 |
| 5x5 conv, 256, pool/2 |
| 3x3 conv, 384 |
| 3x3 conv, 384 |
| 3x3 conv, 256, pool/2 |
| fc, 4096 |
| fc, 4096 |
| fc, 1000 |

**AlexNet**

| 3x3 conv, 64 |
| 3x3 conv, 64, pool/2 |
| 3x3 conv, 128 |
| 3x3 conv, 128, pool/2 |
| 3x3 conv, 256 |
| 3x3 conv, 256 |
| 3x3 conv, 256 |
| 3x3 conv, 256, pool/2 |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| 3x3 conv, 512, pool/2 |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| 3x3 conv, 512, pool/2 |
| fc, 4096 |
| fc, 4096 |
| fc, 1000 |

**VGG**          **GoogLeNet**          **ResNet**

IMAGENET

places
THE SCENE RECOGNITION DATABASE

# Interpretable Units in Different Architectures

# Interpretable Units in Different Architectures



**Number of Unique Detectors**

# Interpretable Units in Different Architectures



**Number of Unique Detectors**

Legend:
- object
- part
- scene
- material
- texture
- color

X-axis categories: AlexNet-random, AlexNet-ImageNet, AlexNet-Places365, VGG-ImageNet, GoogLeNet-ImageNet, GoogLeNet-Places365, VGG-Places365, ResNet152-ImageNet, ResNet152-Places365

House

Airplane

AlexNet — conv5 unit 36 — IoU=0.053

conv5 unit 13 — IoU=0.101

VGG — conv5_3 unit 243 — IoU=0.070

conv5_3 unit 151 — IoU=0.150

GoogLeNet — inception_4e unit 789 — IoU=0.137

inception_4e unit 92 — IoU=0.164

ResNet — res5c unit 1410 — IoU=0.142

res5c unit 1243 — IoU=0.172

Train

Plant

AlexNet — conv4 unit 180 — IoU=0.047

AlexNet — conv5 unit 55 — IoU=0.087

VGG — conv5_3 unit 463 — IoU=0.126

VGG — conv5_3 unit 85 — IoU=0.086

GoogLeNet — inception_5b unit 626 — IoU=0.145

GoogLeNet — inception_4e unit 714 — IoU=0.105

ResNet — res5c unit 924 — IoU=0.293

ResNet — res5c unit 264 — IoU=0.126

**AlexNet** Detector: 81   Unique Detector: 40

Object units built from Places

**ResNet** Detector: 774   Unique Detector: 84

**AlexNet** Detector: 49   Unique Detector: 21

Object units built from ImageNet

**ResNet** Detector: 858   Unique Detector: 75

# Interpretable Units over Layers



AlexNet on Places365

# Interpretable Units over Layers



**AlexNet on Places365**

**AlexNet on ImageNet**

# Interpretable Units over Layers

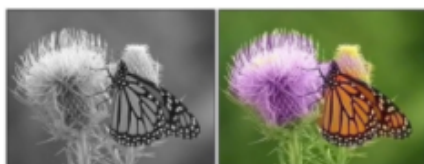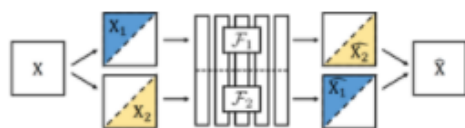# CNNs Trained from Self-supervised Learning
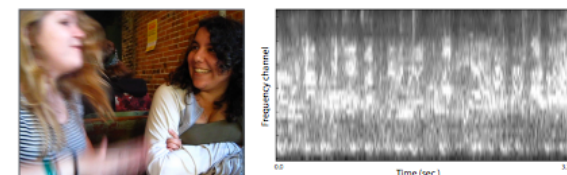
Training CNN without image labels.


Context prediction, ICCV'15


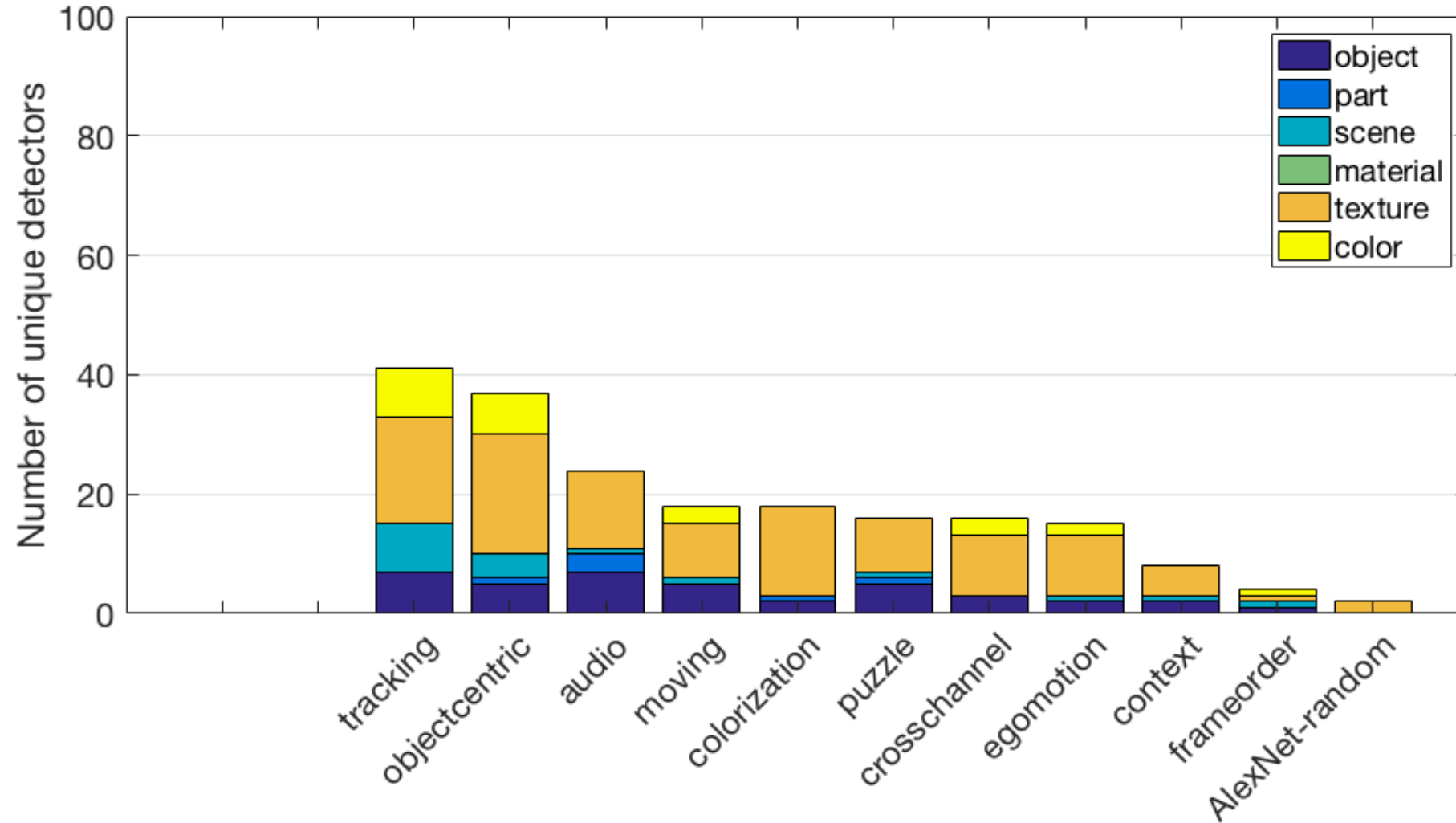Solving puzzle, ECCV'16


Colorization, ECCV'16 and CVPR'17


Audio prediction, ECCV'16
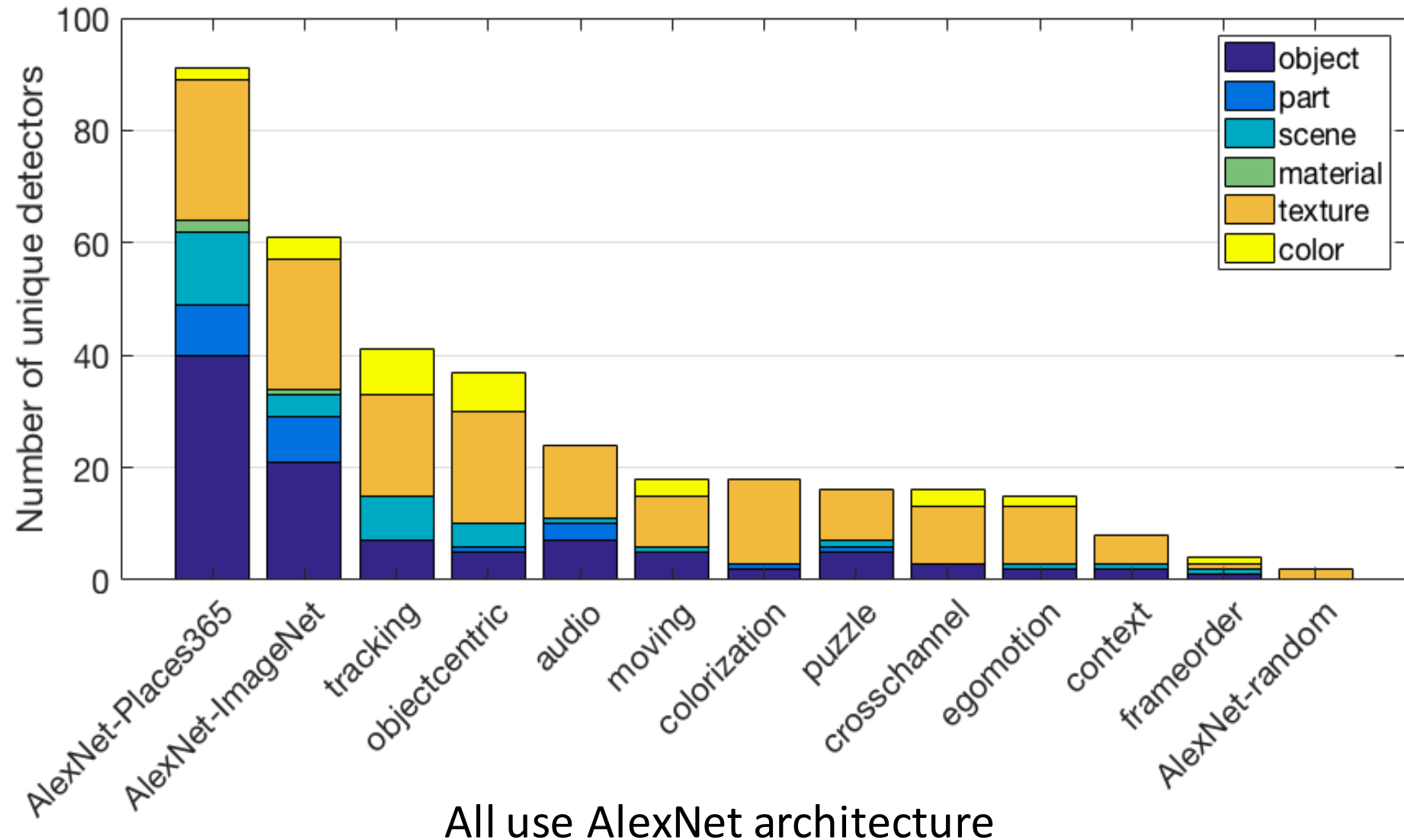
# Comparison of Supervisions



All use AlexNet architecture

# Comparison of Supervisions
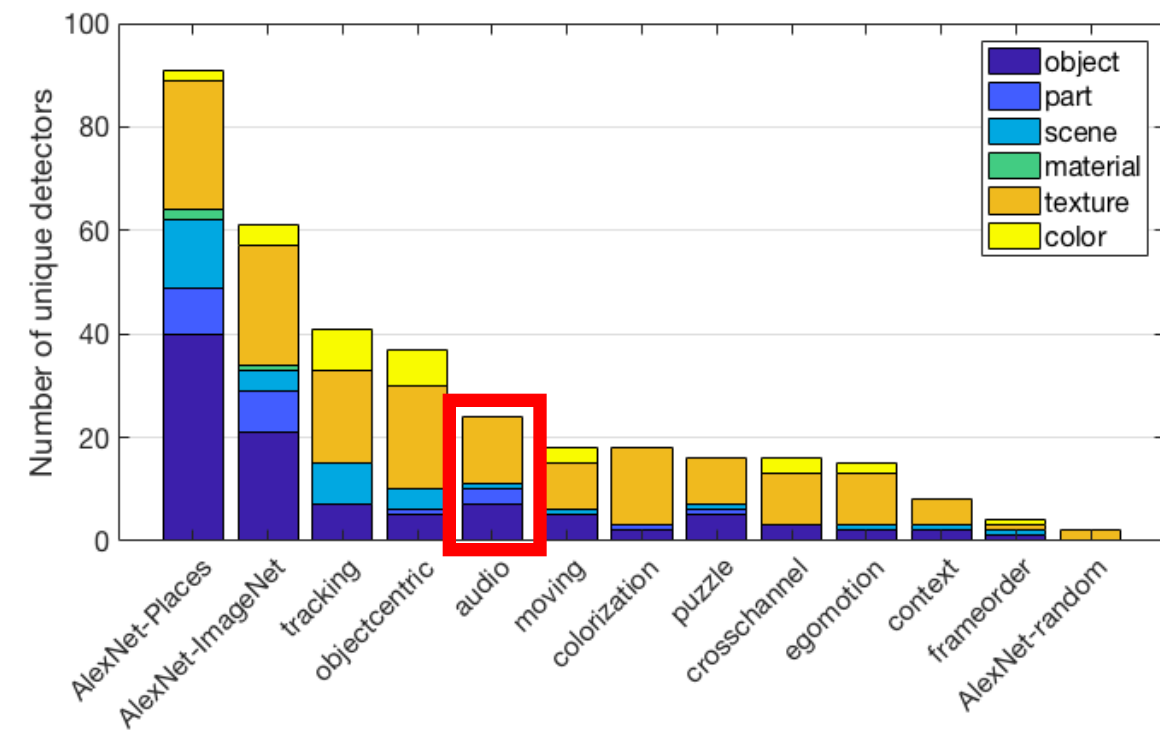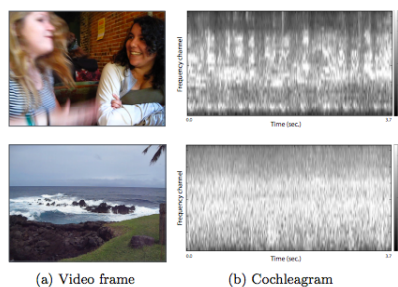


All use AlexNet architecture

# Interpretable Units in Self-supervised Networks

Predict audio from video frames.  ECCV'16 Owens et al.



(a) Video frame　　(b) Cochleagram
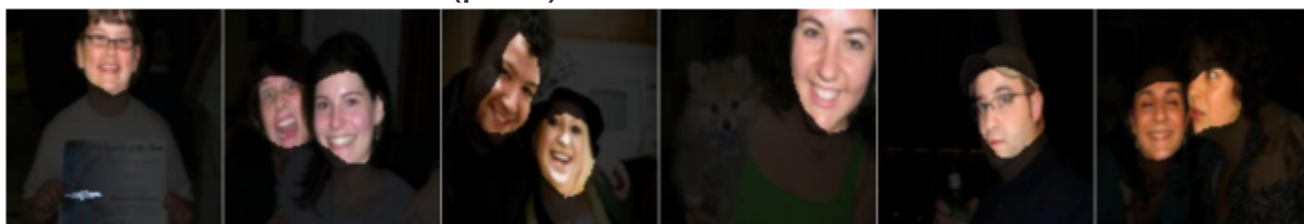
conv5 unit 205: car (object)        IoU=0.063

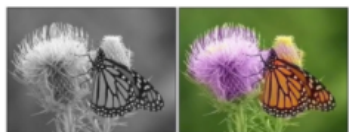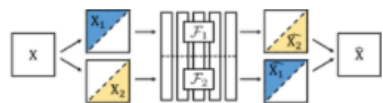conv5 unit 124: creek (scene)        IoU=0.031

conv5 unit 51: head (part)        IoU=0.061

# Interpretable Units in Self-supervised Networks

Colorize grey images
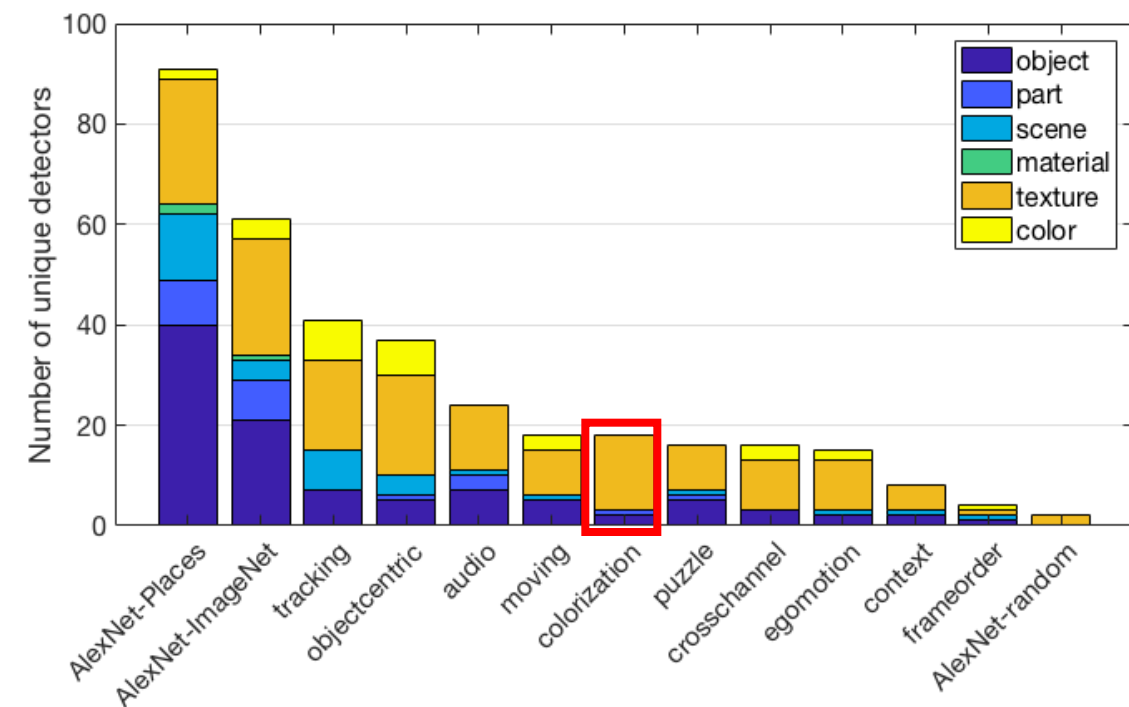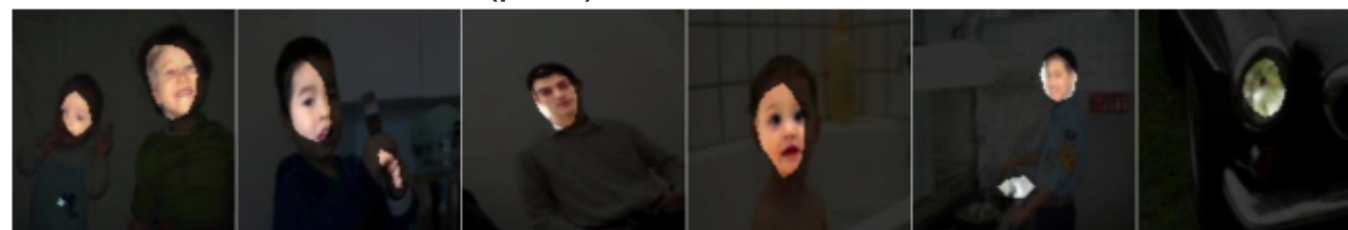ECCV'16. Zhang et al.

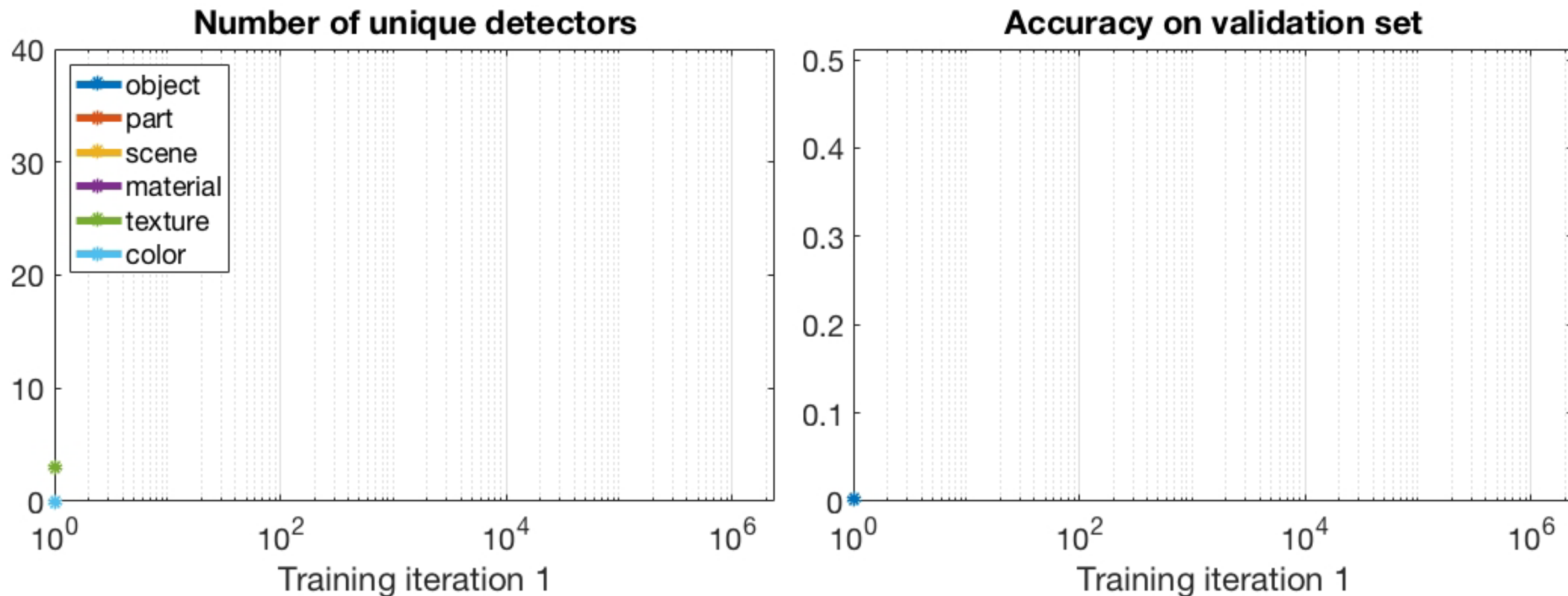conv5 unit 15: banded (texture)    IoU=0.13

conv5 unit 159: tree (object)    IoU=0.039
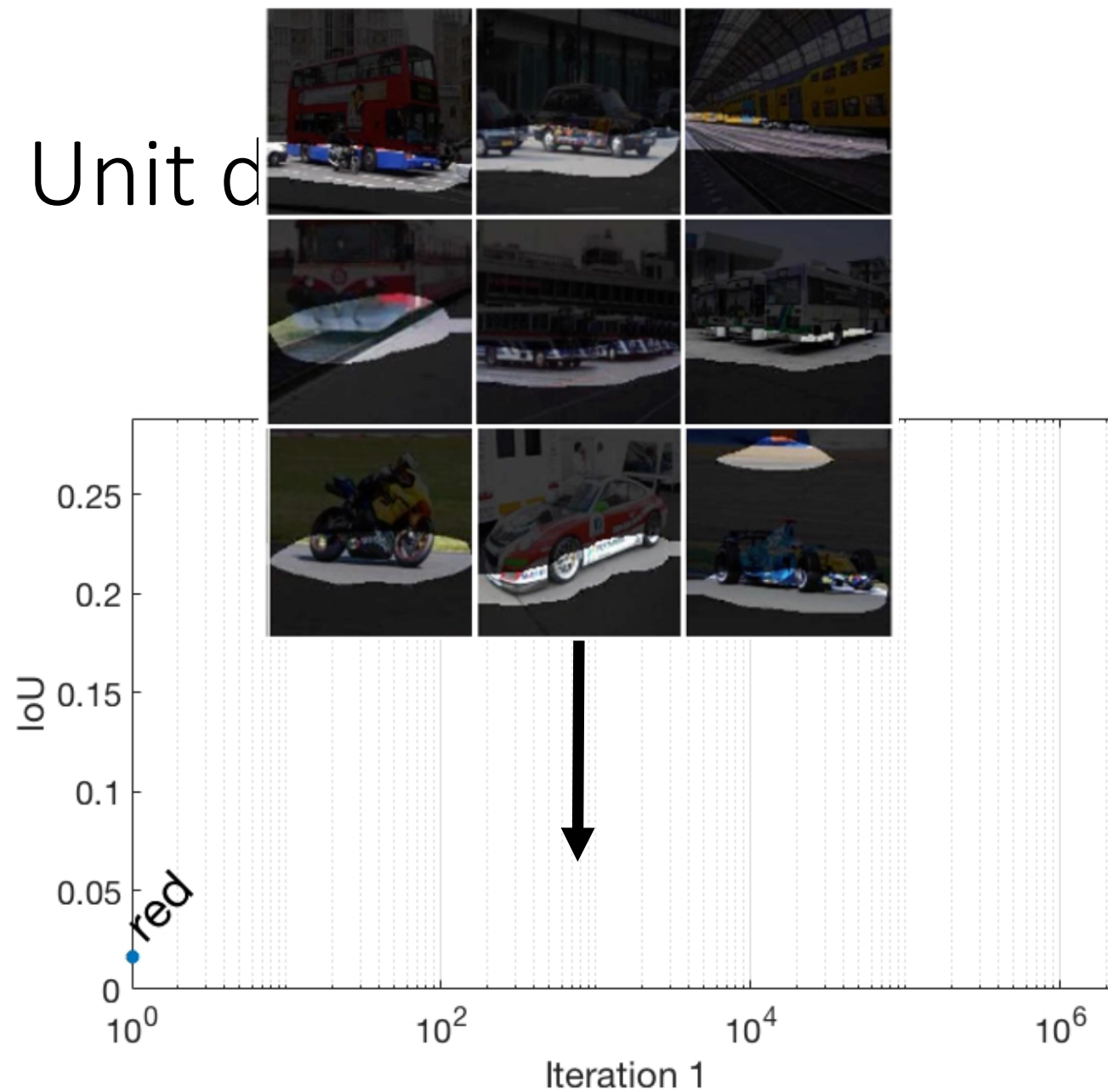
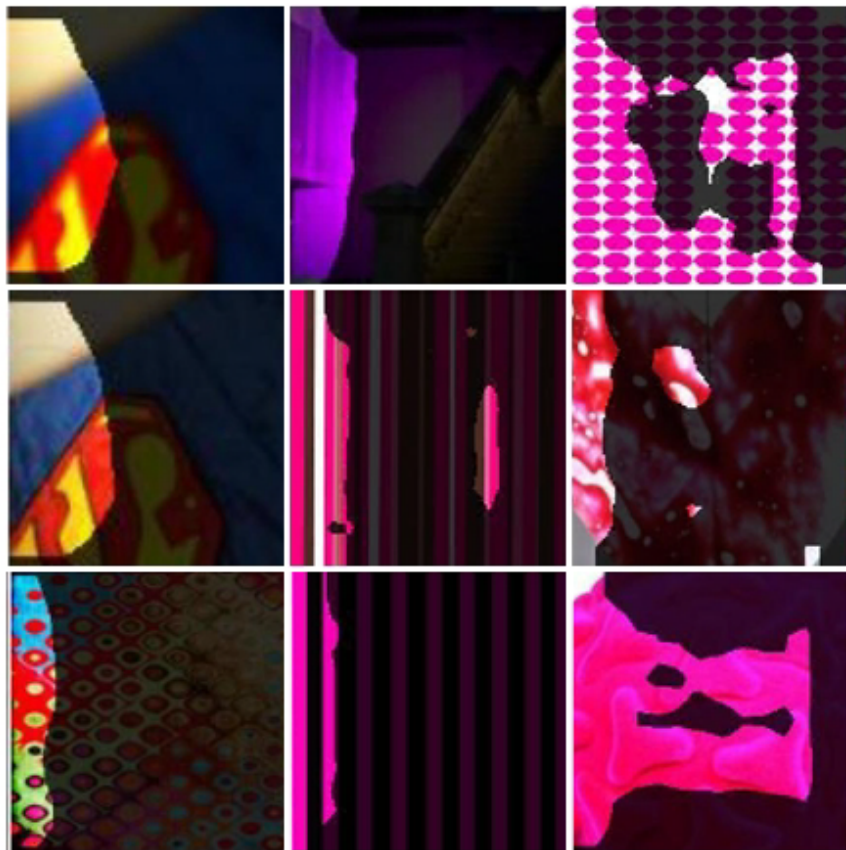conv5 unit 210: head (part)    IoU=0.038

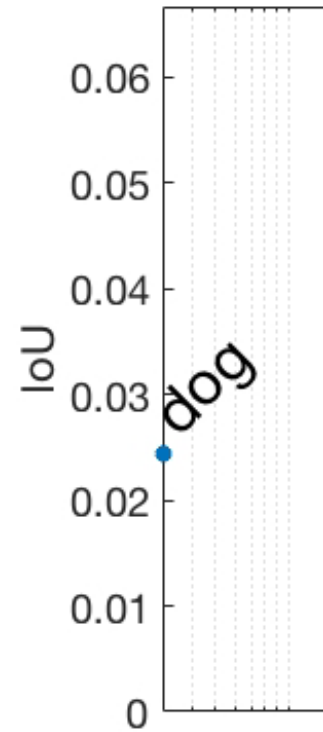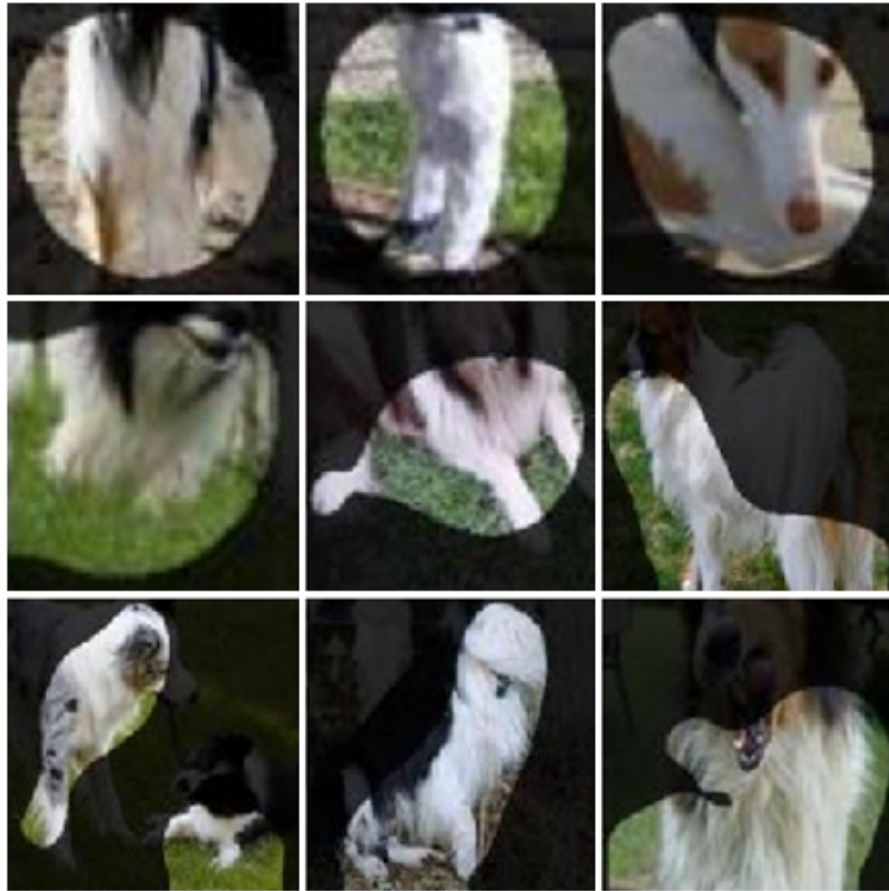# Emergence of Interpretable Units during Training

# Individual Unit d

## Unit 23 at conv5 layer



red

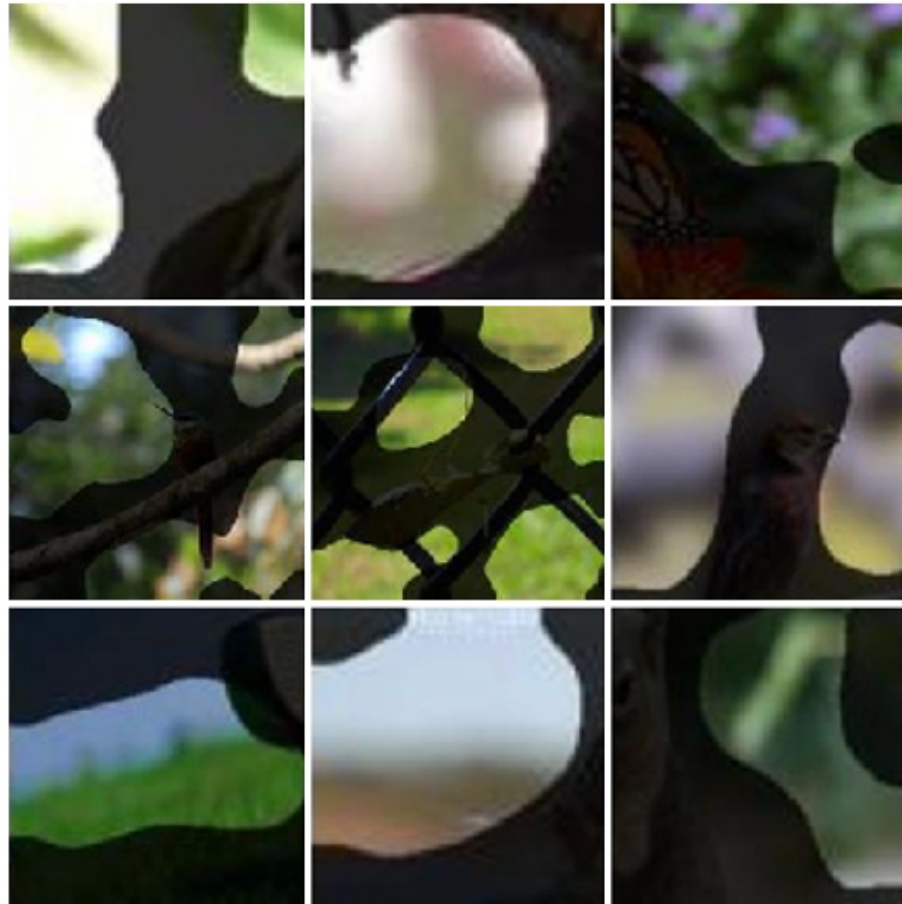# Fine-tuning from ImageNet to Places
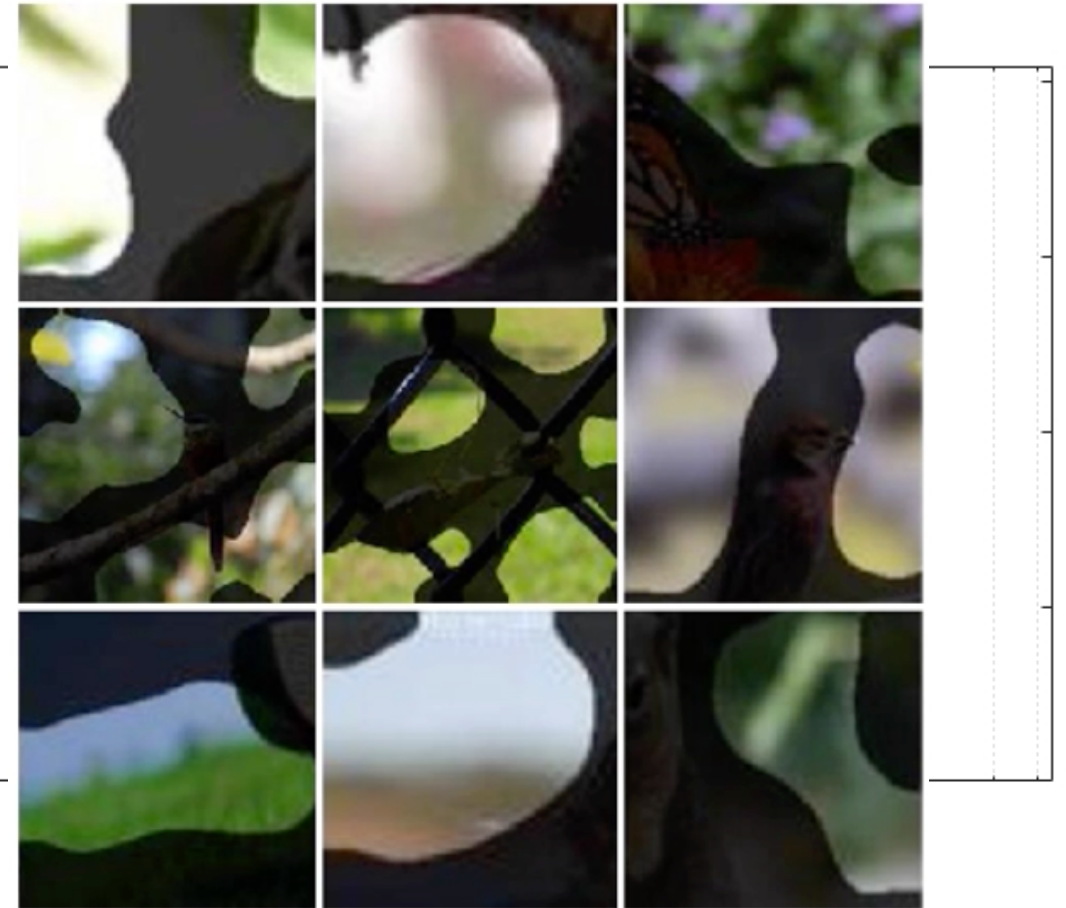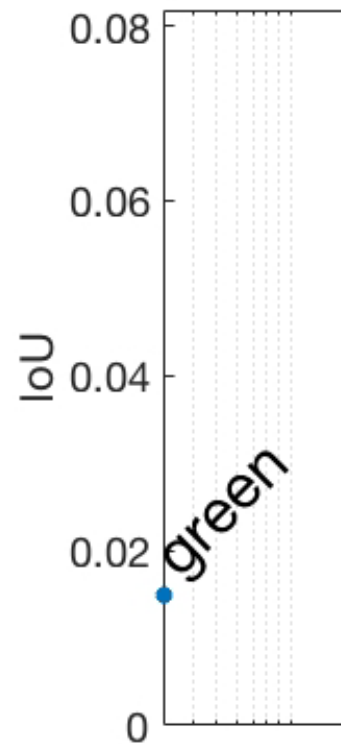
Unit 8 at conv5 layer

Before fine-tuning

# Fine-tuning from ImageNet to Places

Unit 52 at conv5 layer

Before fine-tuning

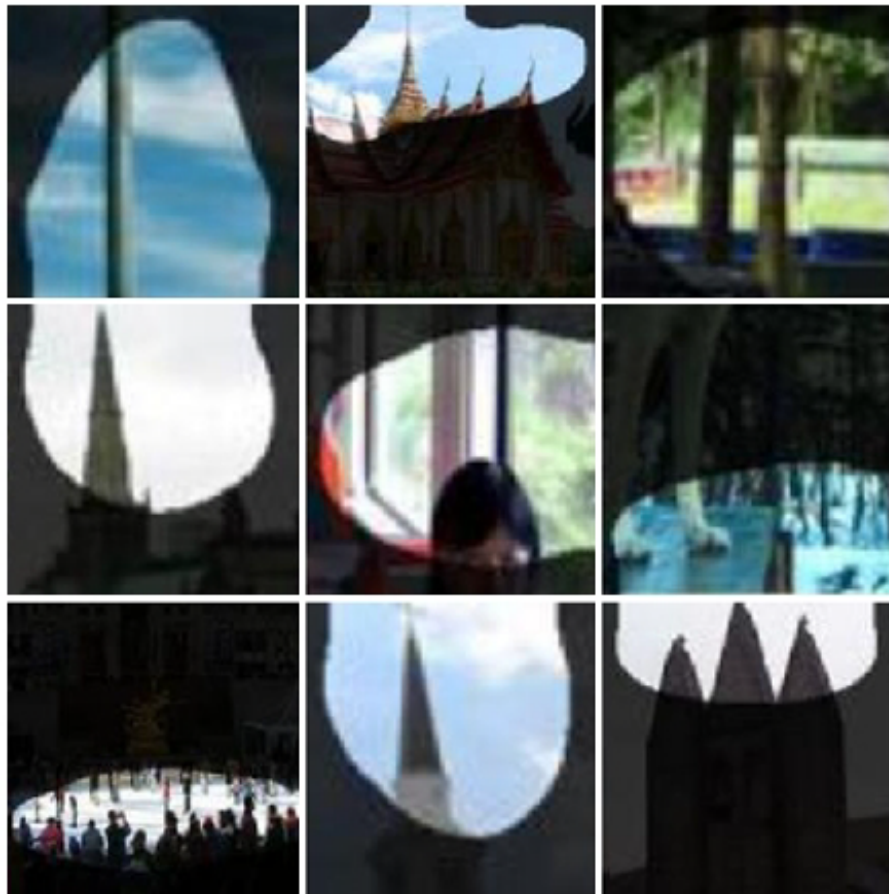# Fine-tuning from Places to ImageNet

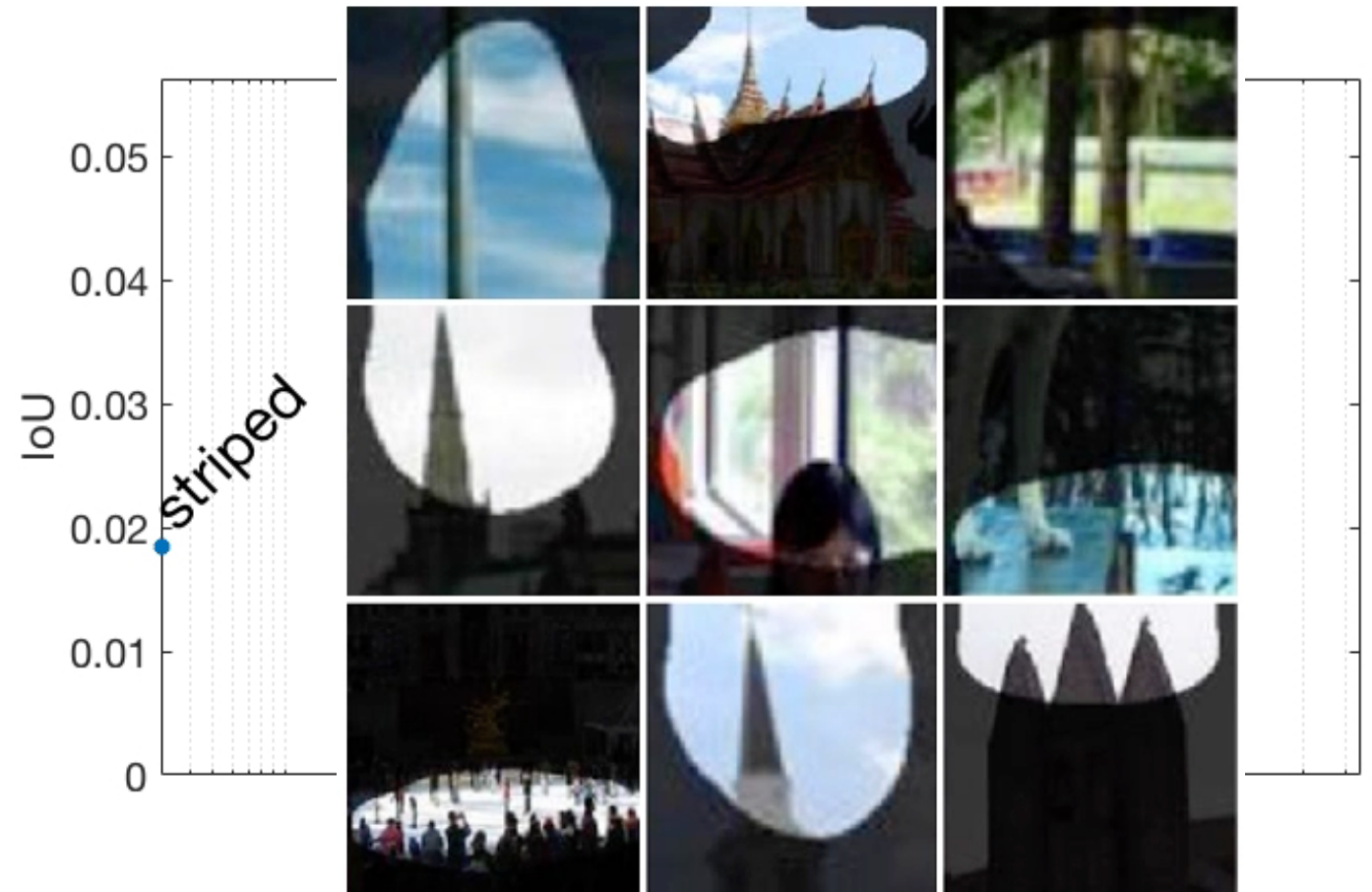Unit 35 at conv5 layer

Before fine-tuning

# Fine-tuning from Places to ImageNet
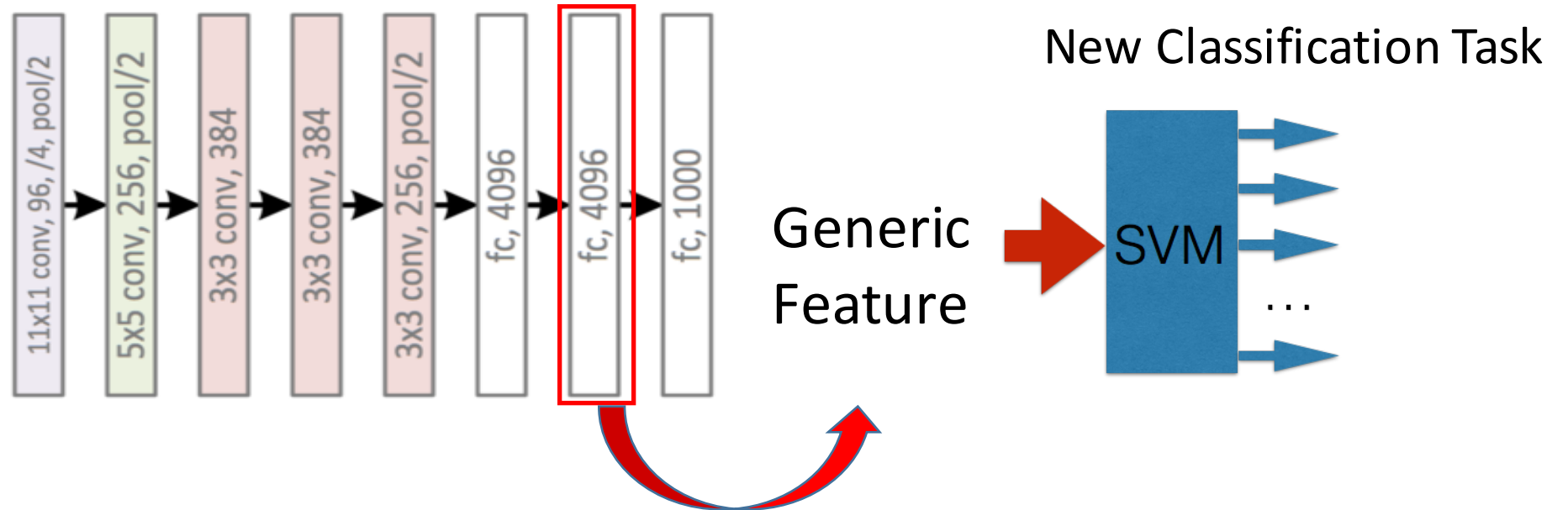
Unit 103 at conv5 layer
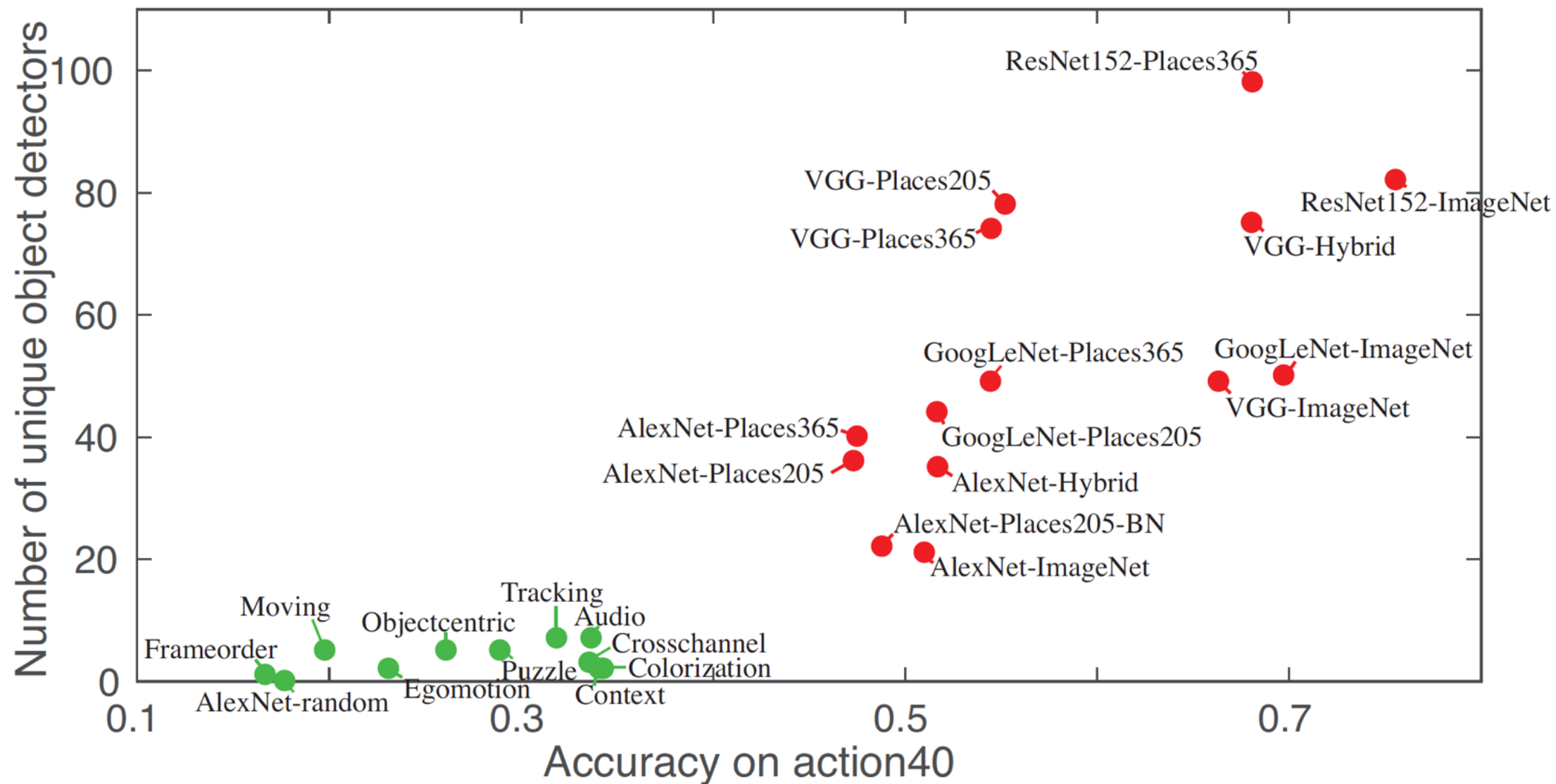
Before fine-tuning

# Explainable Deep Features

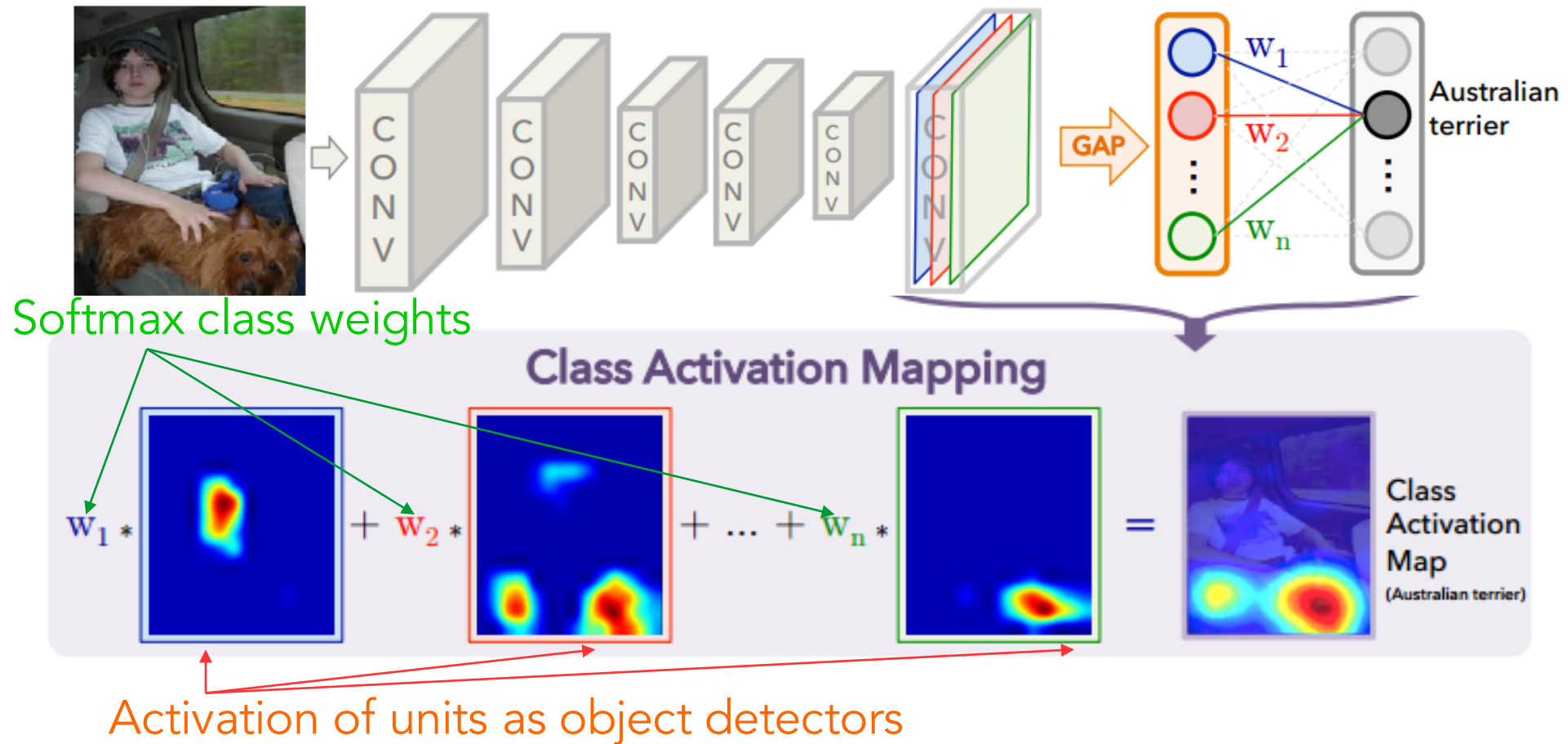## Activations from CNN as generic visual feature

# Deep features as generic visual descriptor

# Explaining the Output



**Softmax class weights**

**Class Activation Mapping**

$$w_1 * \ \ \ \ + \ w_2 * \ \ \ \ + ... + \ w_n * \ \ \ \ = $$

Class Activation Map (Australian terrier)

**Activation of units as object detectors**

Zhou et al. **Learning Deep Features for Discriminative Localization.** CVPR *2016*
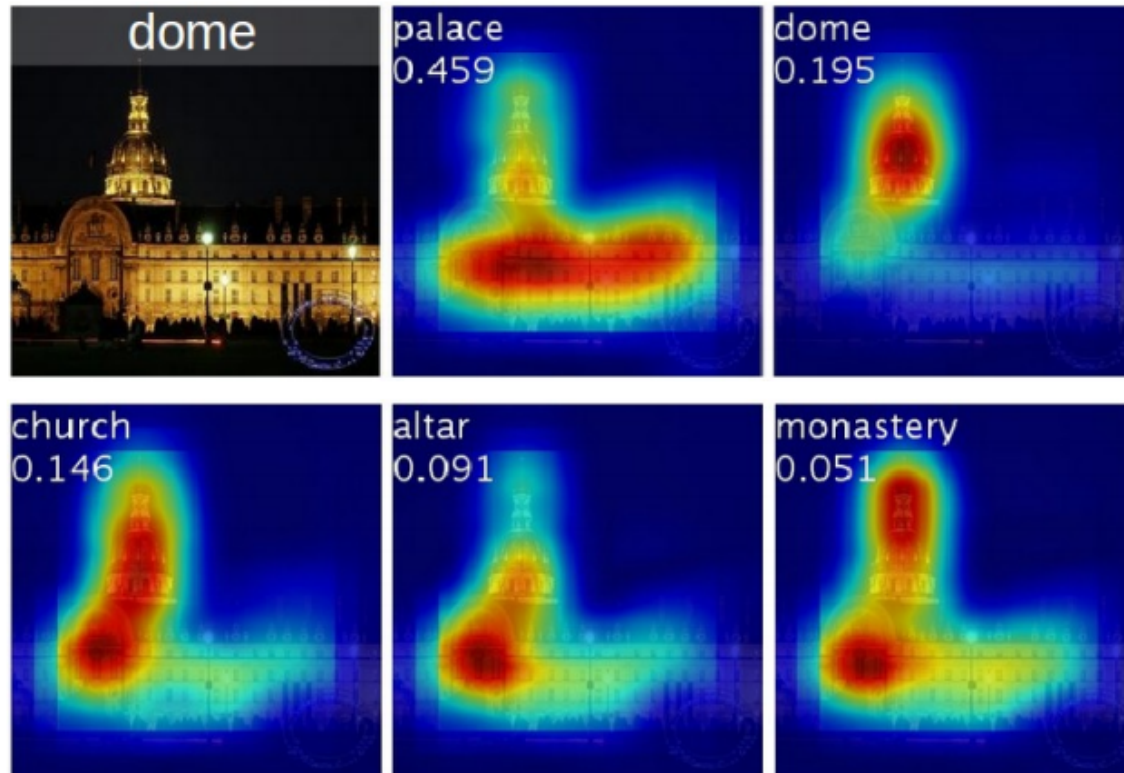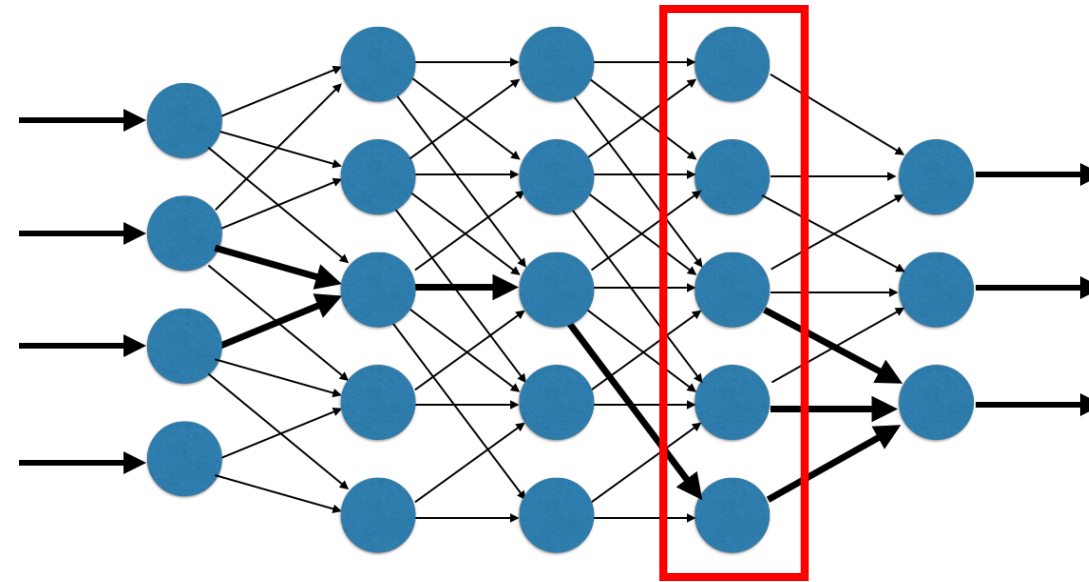
# Explaining the Output

- Class Activation Maps (CAM) for the top5 predictions:
  palace, dome, church, altar, monastery



Zhou et al. **Learning Deep Features for Discriminative Localization.** CVPR *2016*

# Explaining the Output by Unit Interpretations



Walking the dog

Top activated units

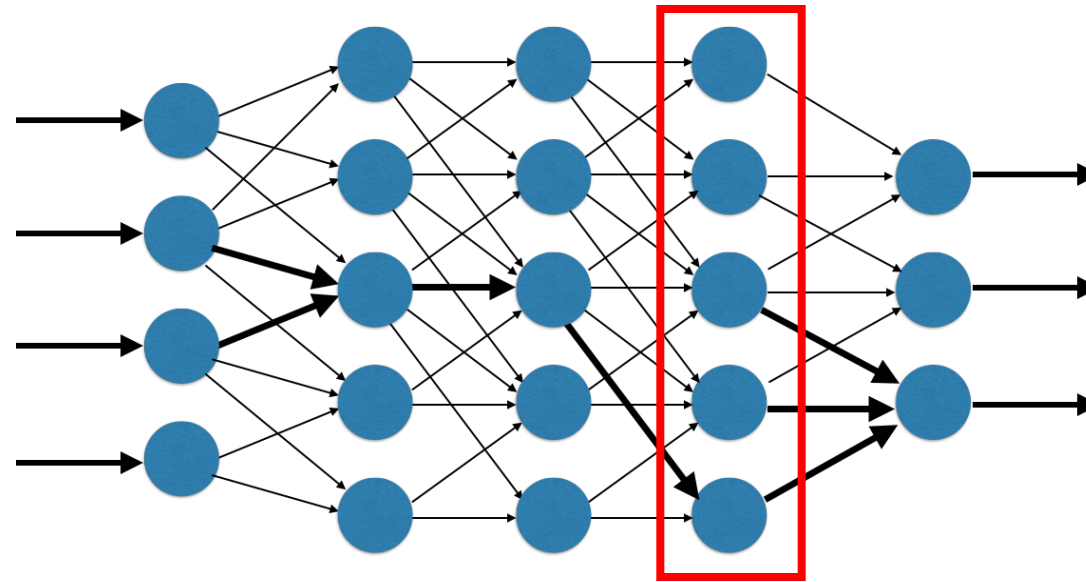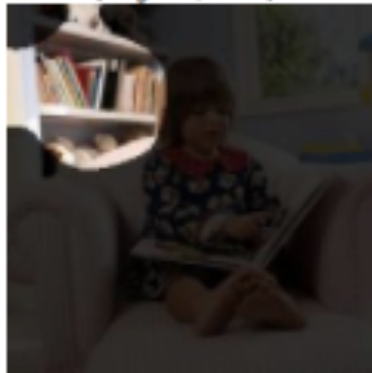| unit 20 | unit 1349 | unit 757 | unit 25 | unit 1647 |
| dog (object,0.04) | leg (part,0.07) | person (object,0.10) | dog (object,0.09) | dog (object,0.02) |

# Explaining the Output by Unit Interpretations



Reading

Top activated units

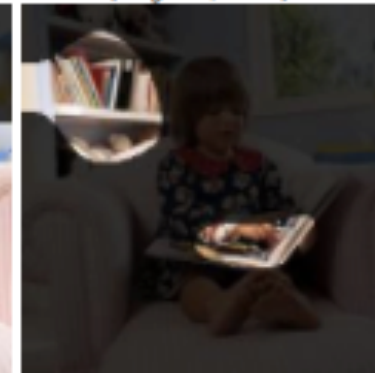unit 362
book (object,0.15)

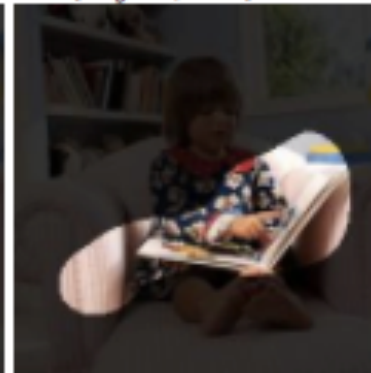unit 1226
person (object,0.13)
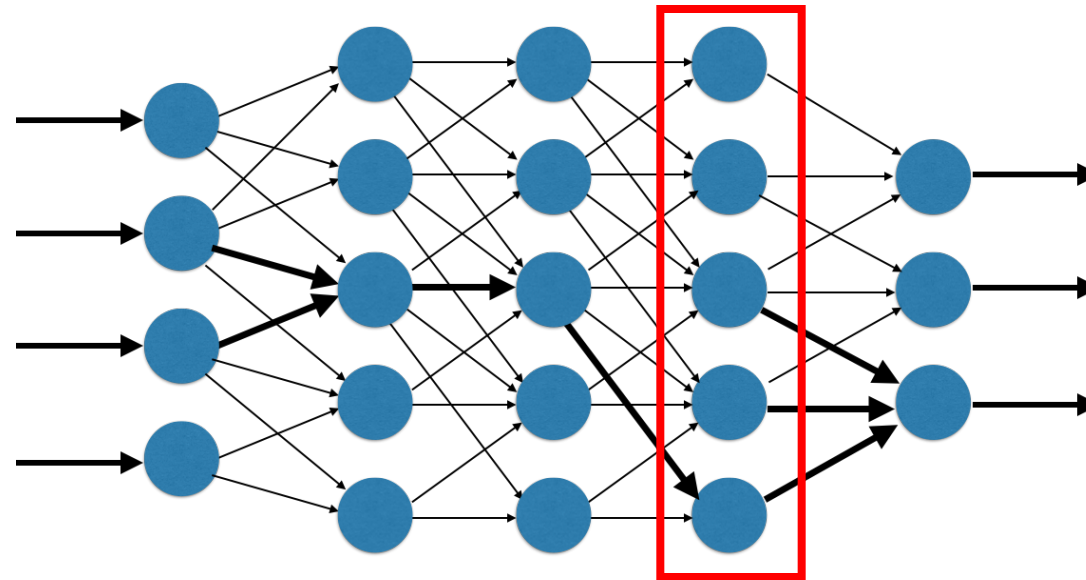
unit 246
back pillow (part,0.04)

unit 365
book (object,0.07)

unit 927
car (object,0.18)

# Explaining the Output by Unit Interpretations



Correct label:
Gardening

Cutting vegetables

Top activated units

unit 1927
arm (part,0.06)

unit 575
table (object,0.03)

unit 1230
pottedplant (object,0.10)

unit 1618
sandbox (scene,0.03)

unit 961
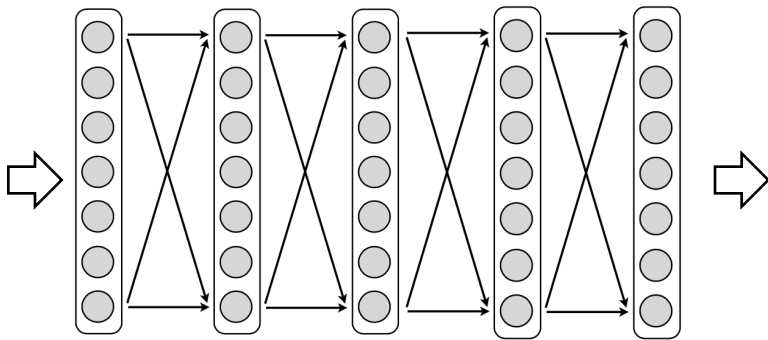sea (object,0.06)

# Conclusion


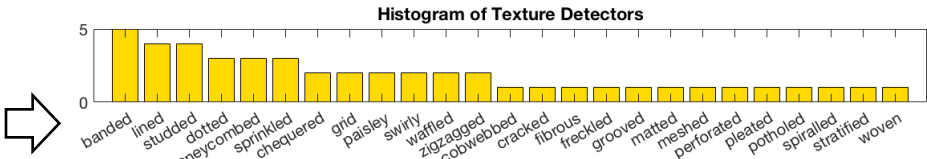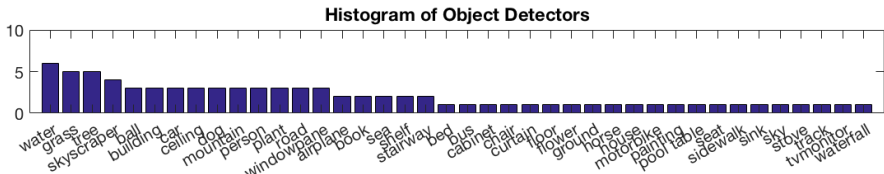
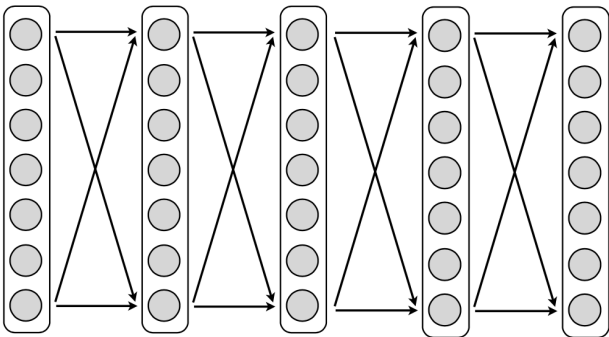Living room
Kitchen
Coast
Theater
...

## Interpretability Report

### Network Dissection



**Histogram of Object Detectors**

**Histogram of Texture Detectors**

unit 79 car, IoU=0.13          unit 107 road, IoU=0.15