# Learning to see

# Antonio Torralba

Computer Science and Artificial Intelligence Laboratory (CSAIL)
Department of Electrical Engineering and Computer Science

# Exciting times for computer vision

# A bit of history…

# The early optimism (1960-1970)

MASSACHUSETTS INSTITUTE OF TECHNOLOGY

PROJECT MAC

Artificial Intelligence Group                    July 7, 1966
Vision Memo. No. 100.

## THE SUMMER VISION PROJECT

Seymour Papert

The summer vision project is an attempt to use our summer workers effectively in the construction of a significant part of a visual system. The particular task was chosen partly because it can be segmented into sub-problems which will allow individuals to work independently and yet participate in the construction of a system complex enough to be a real landmark in the development of "pattern recognition".
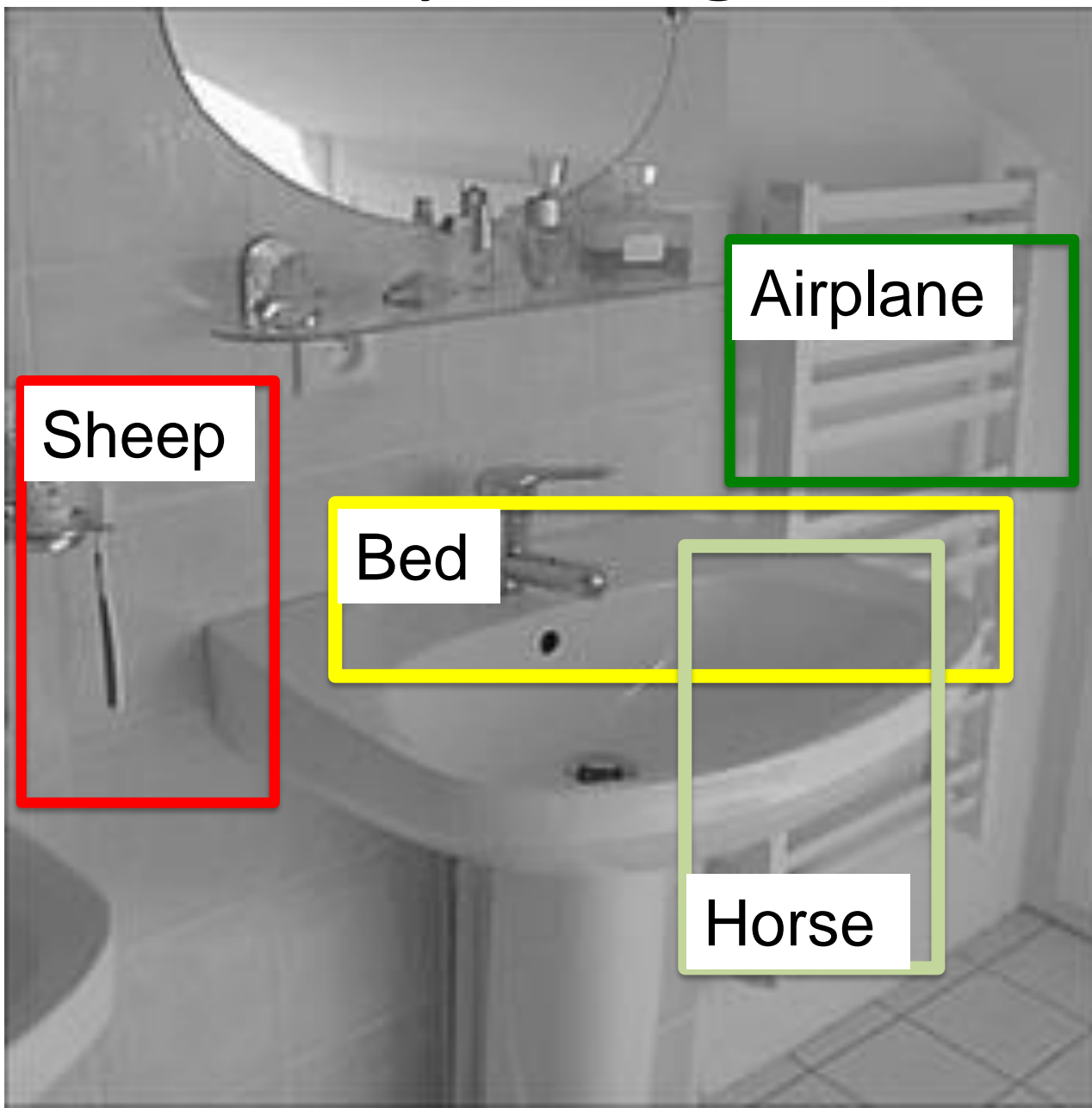
# 50 years ago…

# 50 years ago…

# 50 years ago…



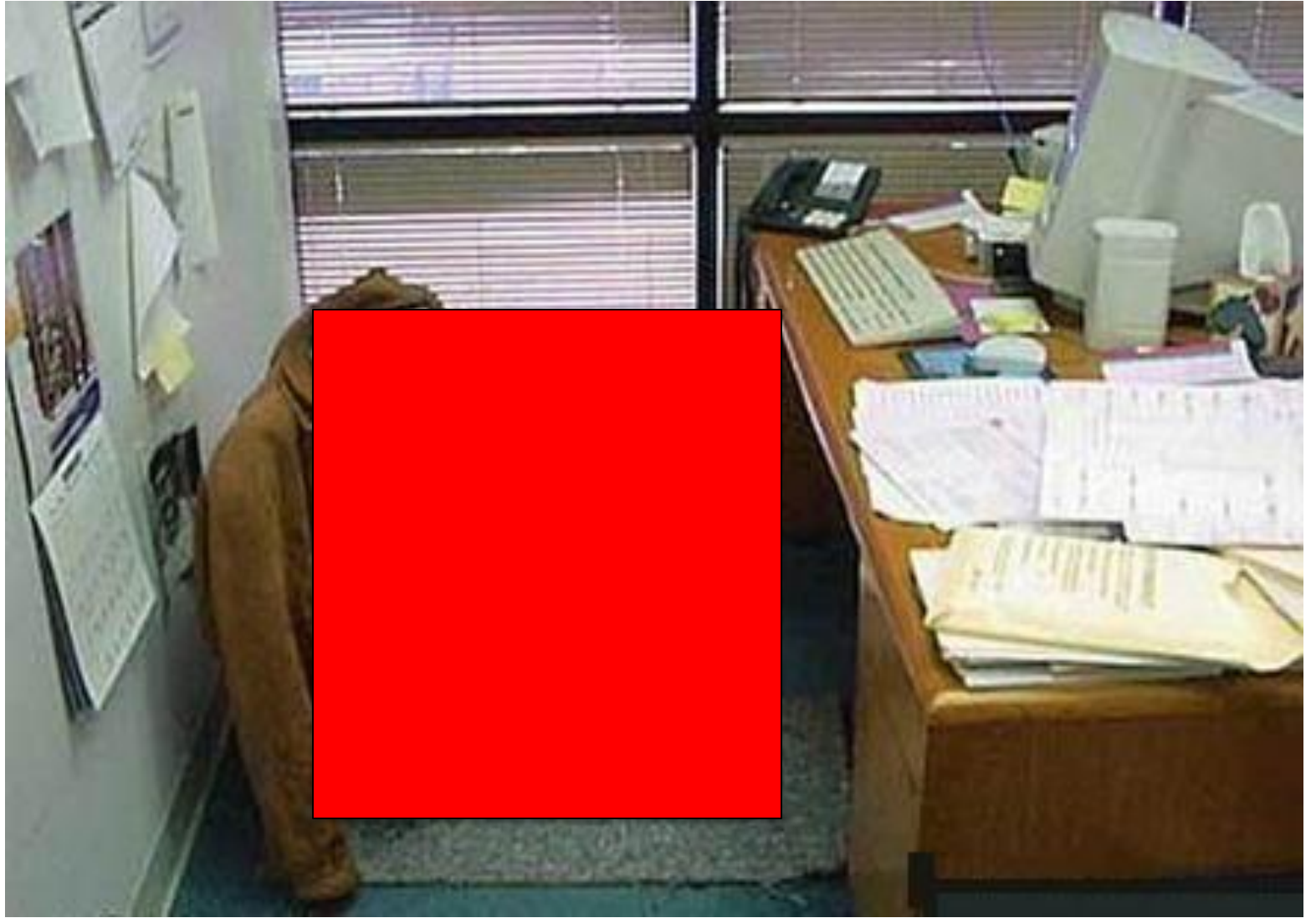Out of memory

# 25 years ago…

# 25 years ago…
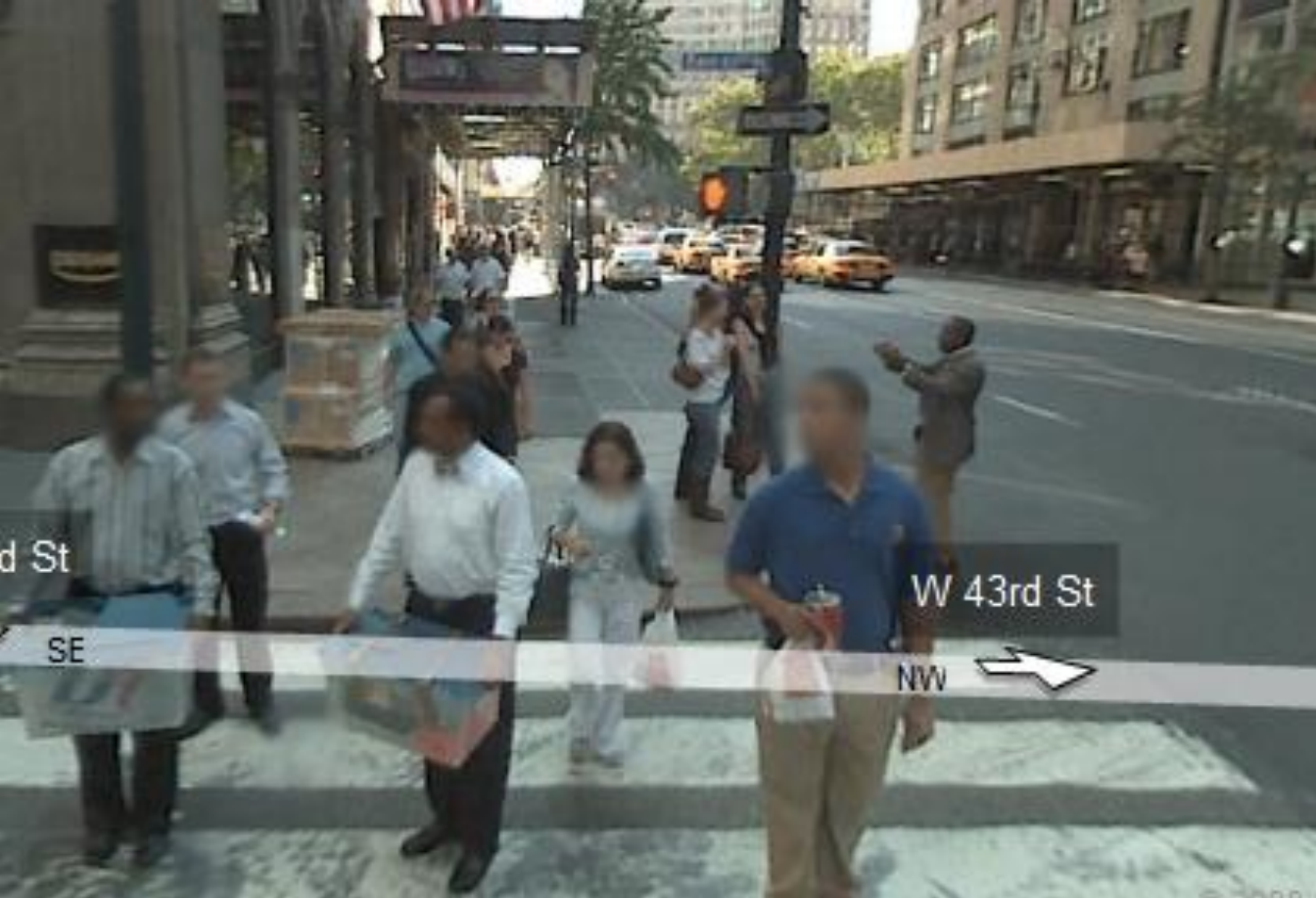
# 25 years ago…
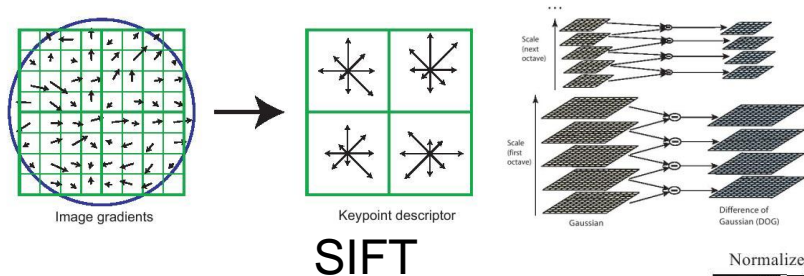
# The vision crisis (1970-2000)

# But 15 years ago…

• The representation and matching of pictorial structures Fischler, Elschlager (1973).
• Face recognition using eigenfaces M. Turk and A. Pentland (1991).
• Human Face Detection in Visual Scenes - Rowley, Baluja, Kanade (1995)
• Graded Learning for Object Detection - Fleuret, Geman (1999)
• Robust Real-time Object Detection - Viola, Jones (2001)
• Feature Reduction and Hierarchy of Classifiers for Fast Object Detection in Video Images - Heisele, Serre, Mukherjee, Poggio (2001)
•….

- The representation and matching of pictorial structures Fischler, Elschlager (1973).
- Face recognition using eigenfaces M. Turk and A. Pentland (1991).
- Human Face Detection in Visual Scenes - Rowley, Baluja, Kanade (1995)
- Graded Learning for Object Detection - Fleuret, Geman (1999)
- Robust Real-time Object Detection - Viola, Jones (2001)
- Feature Reduction and Hierarchy of Classifiers for Fast Object Detection in Video Images - Heisele, Serre, Mukherjee, Poggio (2001)
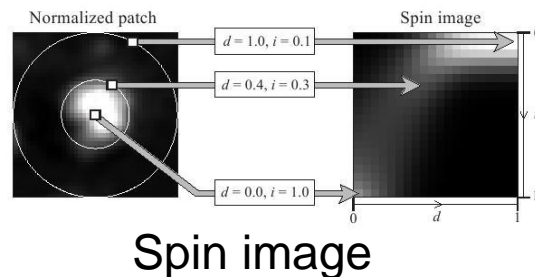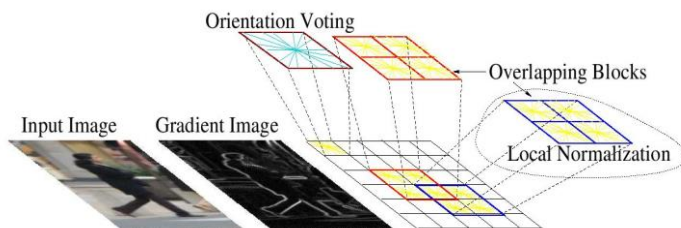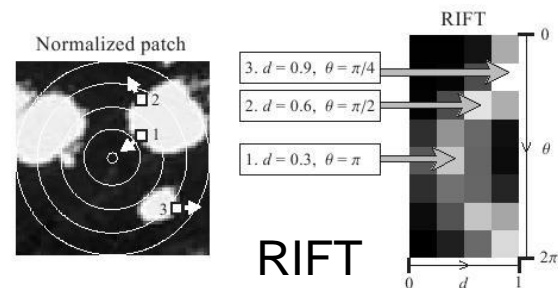- ….

• The representation and matching of pictorial structures Fischler, Elschlager (1973).
• Face recognition using eigenfaces M. Turk and A. Pentland (1991).
• Human Face Detection in Visual Scenes - Rowley, Baluja, Kanade (1995)
• Graded Learning for Object Detection - Fleuret, Geman (1999)
• Robust Real-time Object Detection - Viola, Jones (2001)
• Feature Reduction and Hierarchy of Classifiers for Fast Object Detection in Video Images - Heisele, Serre, Mukherjee, Poggio (2001)
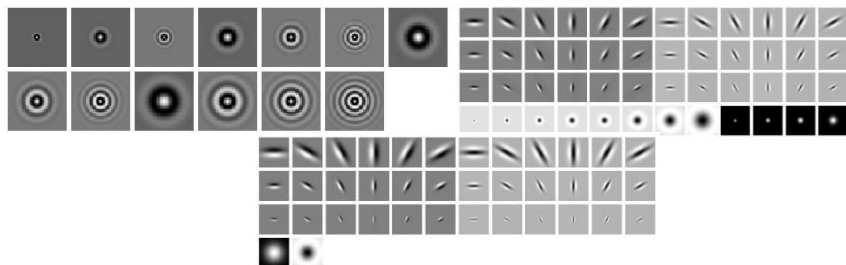• ….

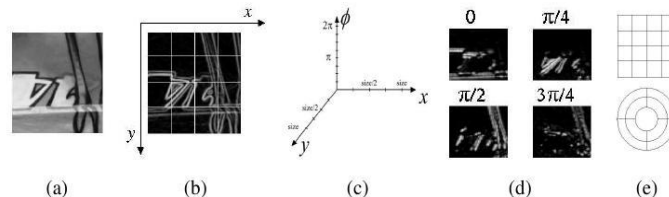# Advances in computer vision


SIFT


GIST


Spin image


HoG


RIFT


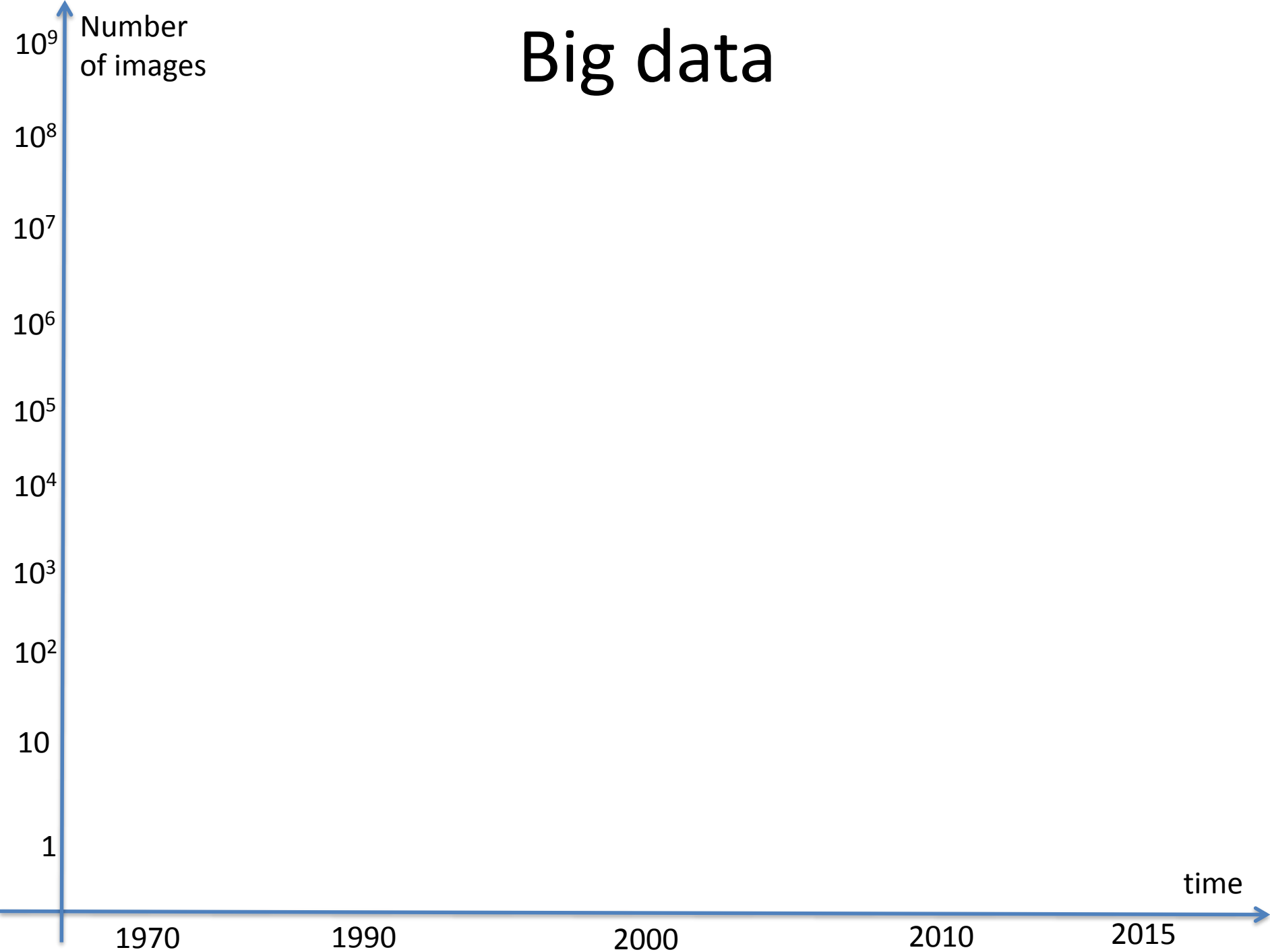Textons
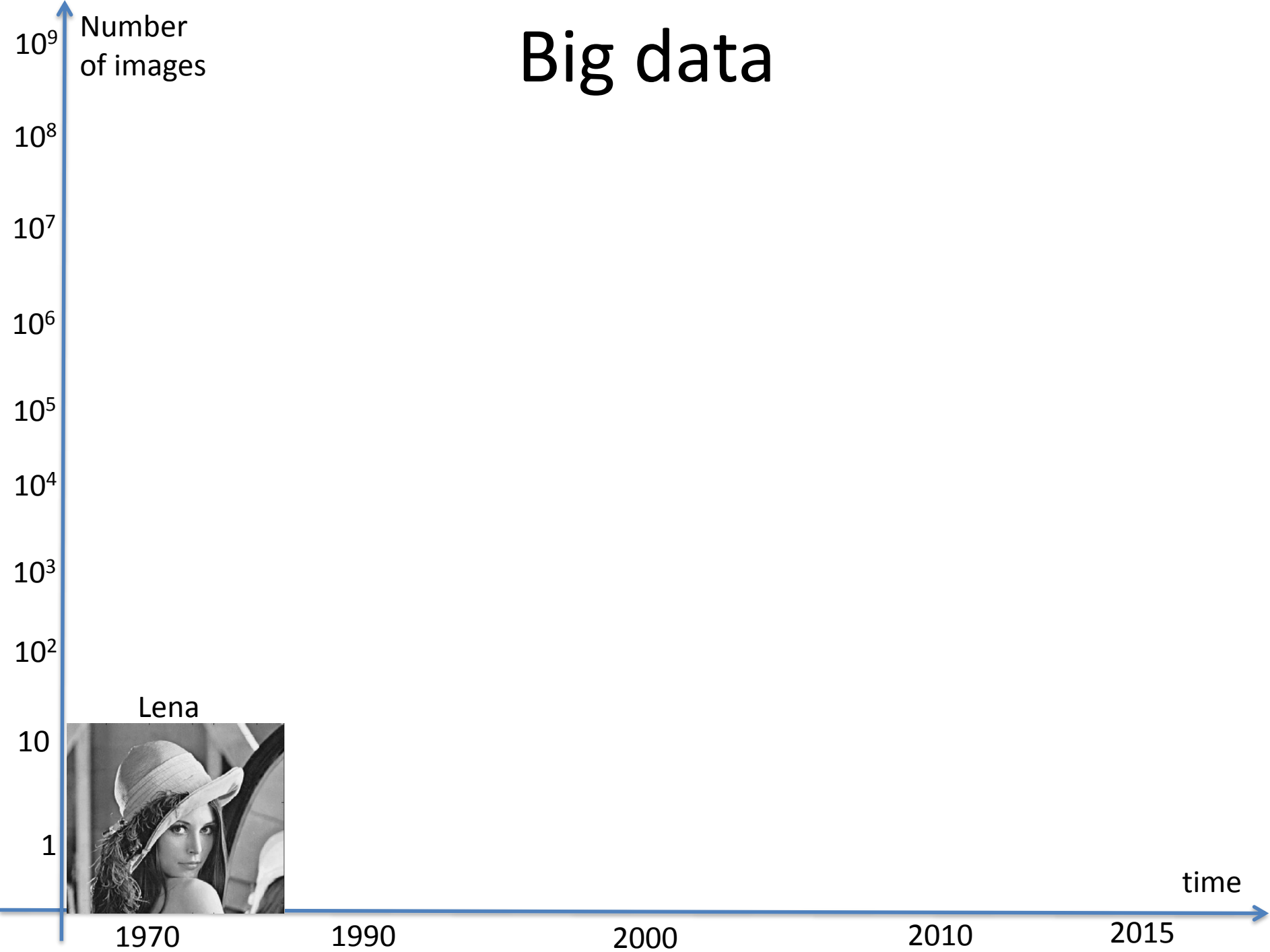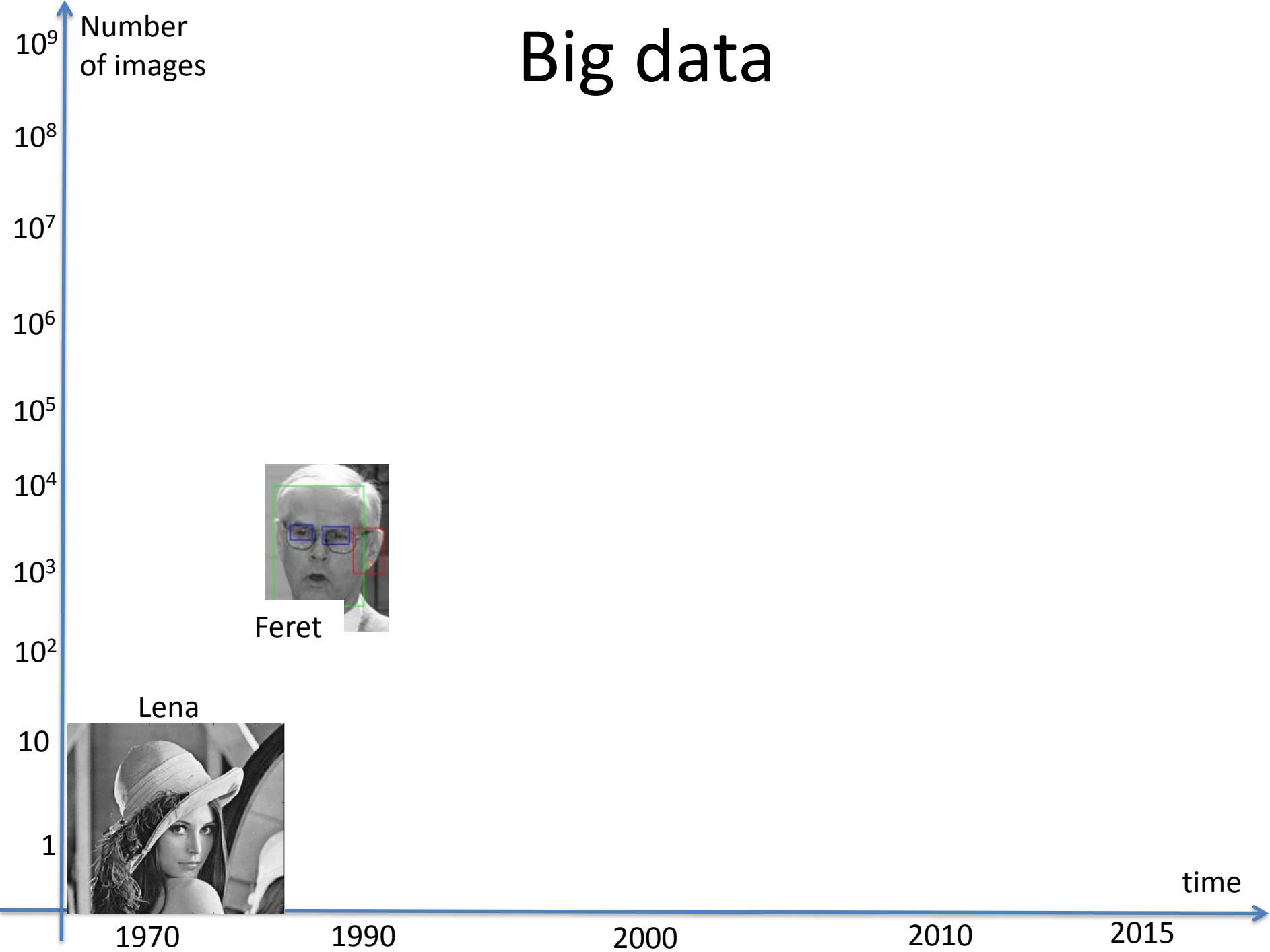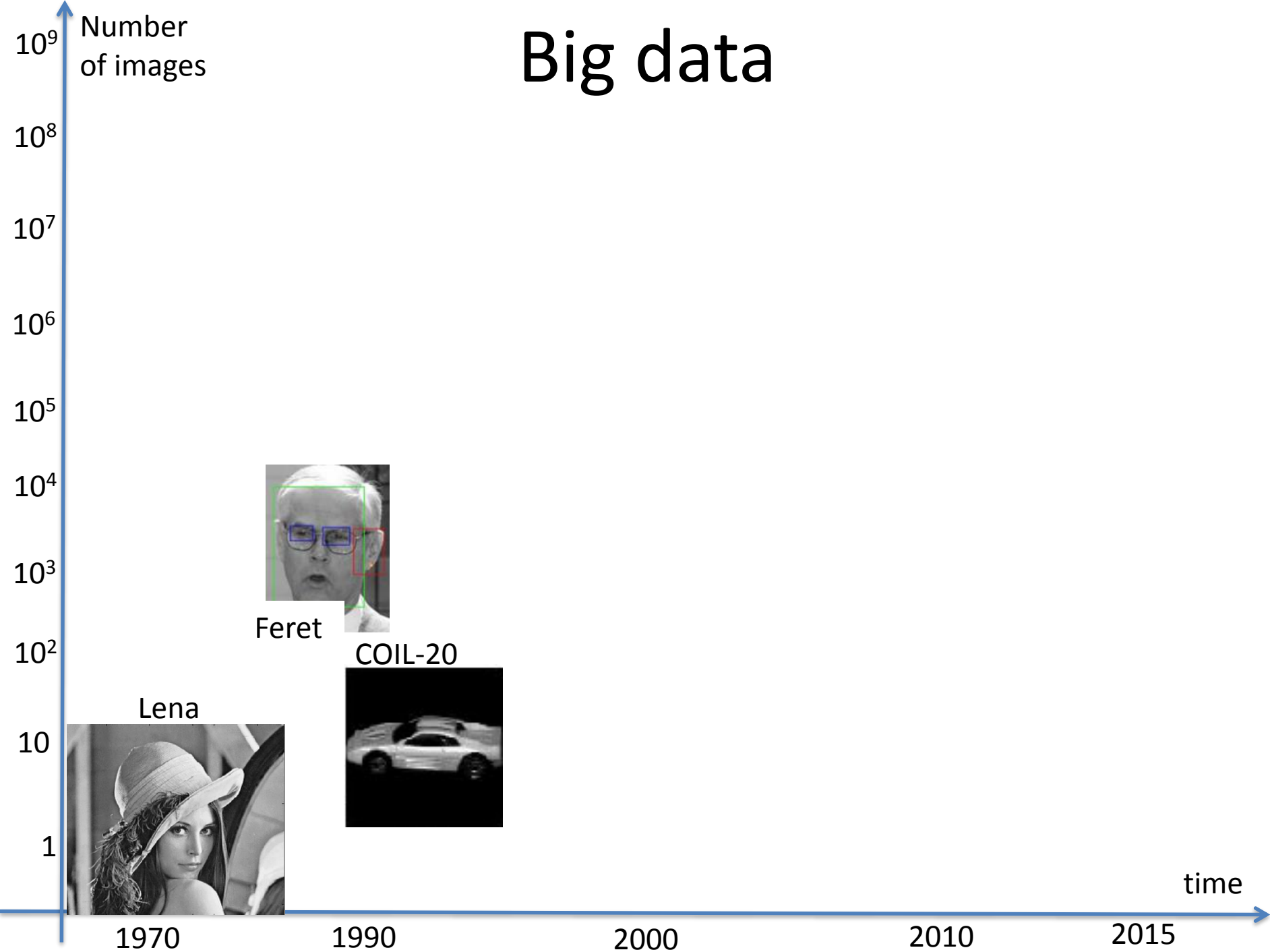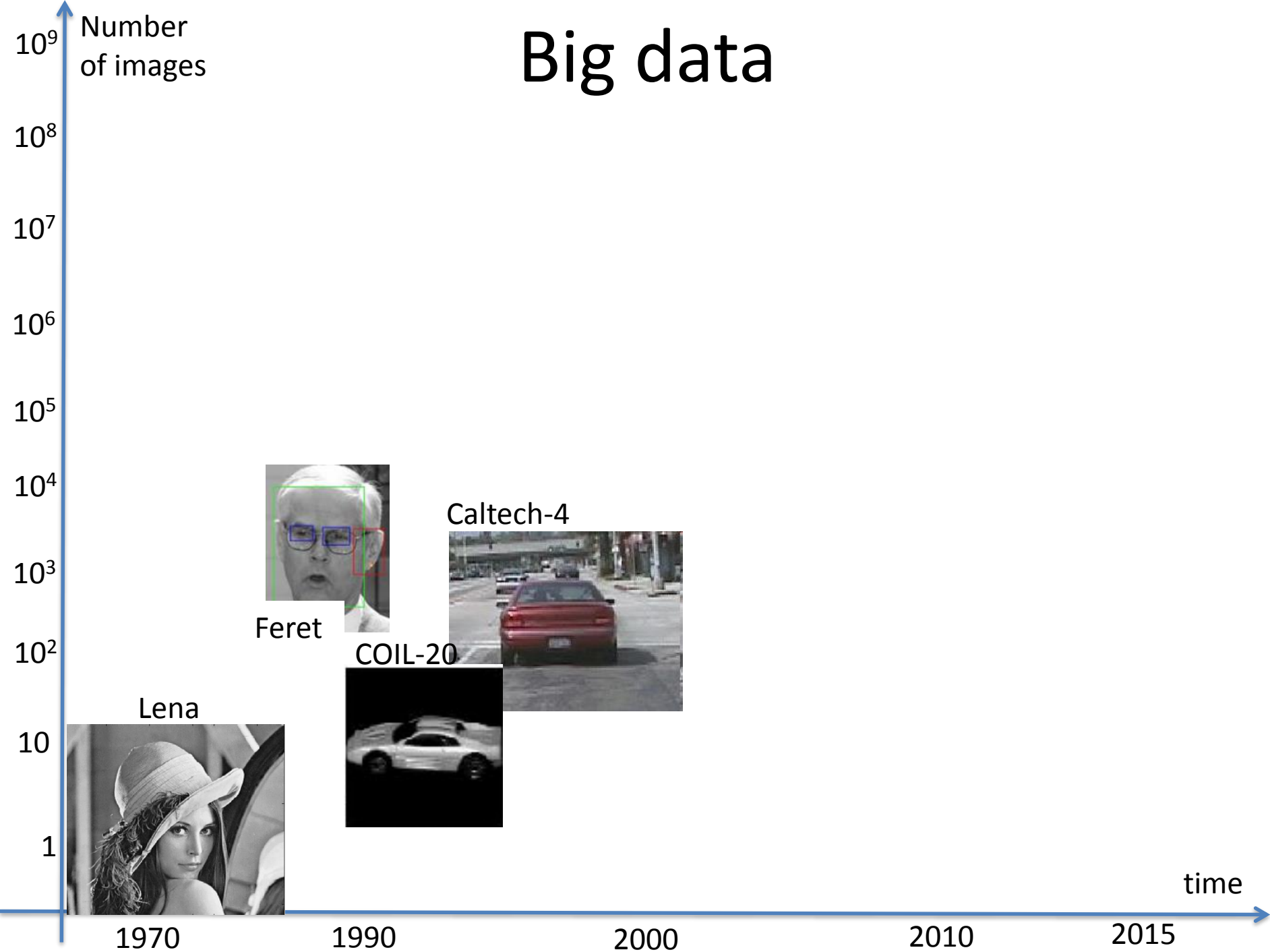

GLOH

# A short story of image databases

Big data

Number of images

$10^9$

$10^8$

$10^7$

$10^6$

$10^5$

$10^4$

$10^3$

$10^2$

10

1

time

1970　　　　　1990　　　　　2000　　　　　2010　　　2015

Big data

Number of images

$10^9$
$10^8$
$10^7$
$10^6$
$10^5$
$10^4$
$10^3$
$10^2$
$10$
$1$

Lena

time

1970    1990    2000    2010    2015

# Big data

**Number of images**

$10^9$

$10^8$

$10^7$

$10^6$

$10^5$

$10^4$

$10^3$

Feret

$10^2$

10

Lena

1

time

1970    1990    2000    2010    2015

Big data

# Big data

Number of images

$10^9$

$10^8$

$10^7$

$10^6$

$10^5$

$10^4$

$10^3$

$10^2$

10

1

Lena

Feret

COIL-20

Caltech-4

time

1970    1990    2000    2010    2015

# Big data

Number of images

$10^9$

$10^8$

$10^7$

$10^6$

$10^5$

$10^4$

$10^3$

$10^2$

10

1

Lena

COIL-20

Feret

Caltech-4

Caltech 101

time

1970      1990      2000      2010      2015

# Big data



Number of images (y-axis): $10^9$, $10^8$, $10^7$, $10^6$, $10^5$, $10^4$, $10^3$, $10^2$, $10$, $1$

time (x-axis): 1970, 1990, 2000, 2010, 2015

Lena

Feret

COIL-20

Caltech-4

Caltech 101

PASCAL

# Big data

10^9
10^8

**80 million images**

10^7

IM**A**GENET

10^6

Caltech 101

10^5

PASCAL

10^4

Caltech-4

10^3

Feret

10^2

COIL-20

10

Lena

1

Number of images

time

1970    1990    2000    2010    2015

# Big data



Number of images

$10^9$

$10^8$

80 million images

$10^7$

places

IMAGENET

$10^6$

Caltech 101

$10^5$

PASCAL

$10^4$

Caltech-4

$10^3$

Feret

$10^2$

COIL-20

10

Lena

1

time

1970        1990        2000        2010        2015

Big data

# The time of big data

In 2010, a new student gets into computer vision…

# In 2010, a new student gets into computer vision…

## Pick one dataset

# In 2010, a new student gets into computer vision…

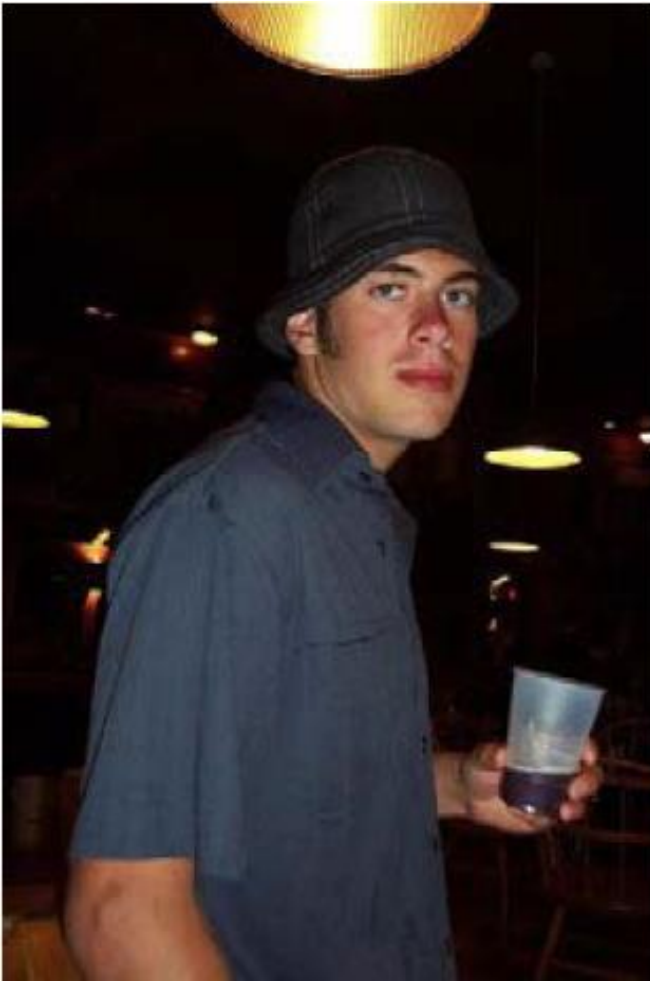## Pick one dataset



## Pick one model
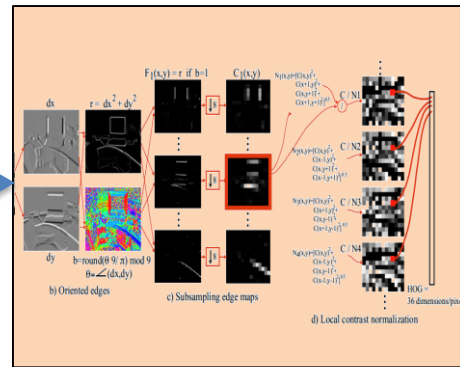
car

# Who's to blame?
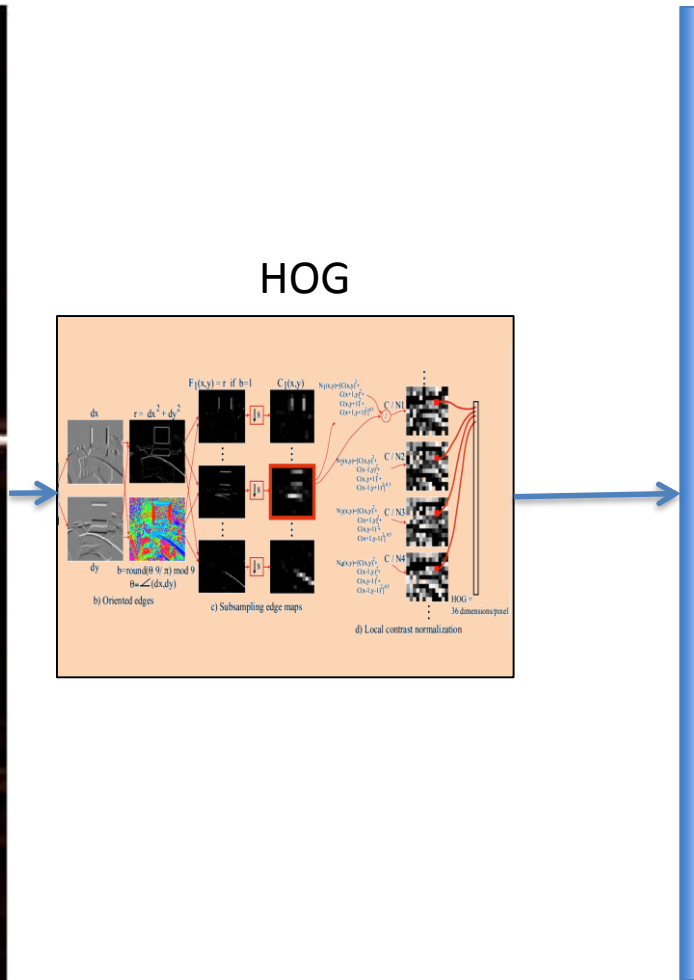


car

- The data
- The features
- The student

# Features for object detection



HOG

# What does a detector sees?



HOG

Can we visualize this output?

Carl Vondrick

Aditya Khosla

# What does a detector sees?



HOG

Vondrick, Khosla, Malisiewicz, Torralba. "Inverting and Visualizing Features for Object Detection.", ICCV 2013

# Person



# Chair



# Car



## Can you tell which ones are not the object?

Vondrick, Khosla, Malisiewicz, Torralba. "Inverting and Visualizing Features for Object Detection."

# Person



# Chair



# Car



Vondrick, Khosla, Malisiewicz, Torralba. "Inverting and Visualizing Features for Object Detection."

# HOG visualization predicts SVM performance



Chair detection test

Legend:
- HOG+Human AP = 0.63
- RGB+Human AP = 0.96
- HOG+DPM AP = 0.51

Vondrick, Khosla, Malisiewicz, Torralba. "Inverting and Visualizing Features for Object Detection."

car

http://mit.edu/vondrick/ihog/

Vondrick, Khosla, Malisiewicz, Torralba. "Inverting and Visualizing Features for Object Detection."
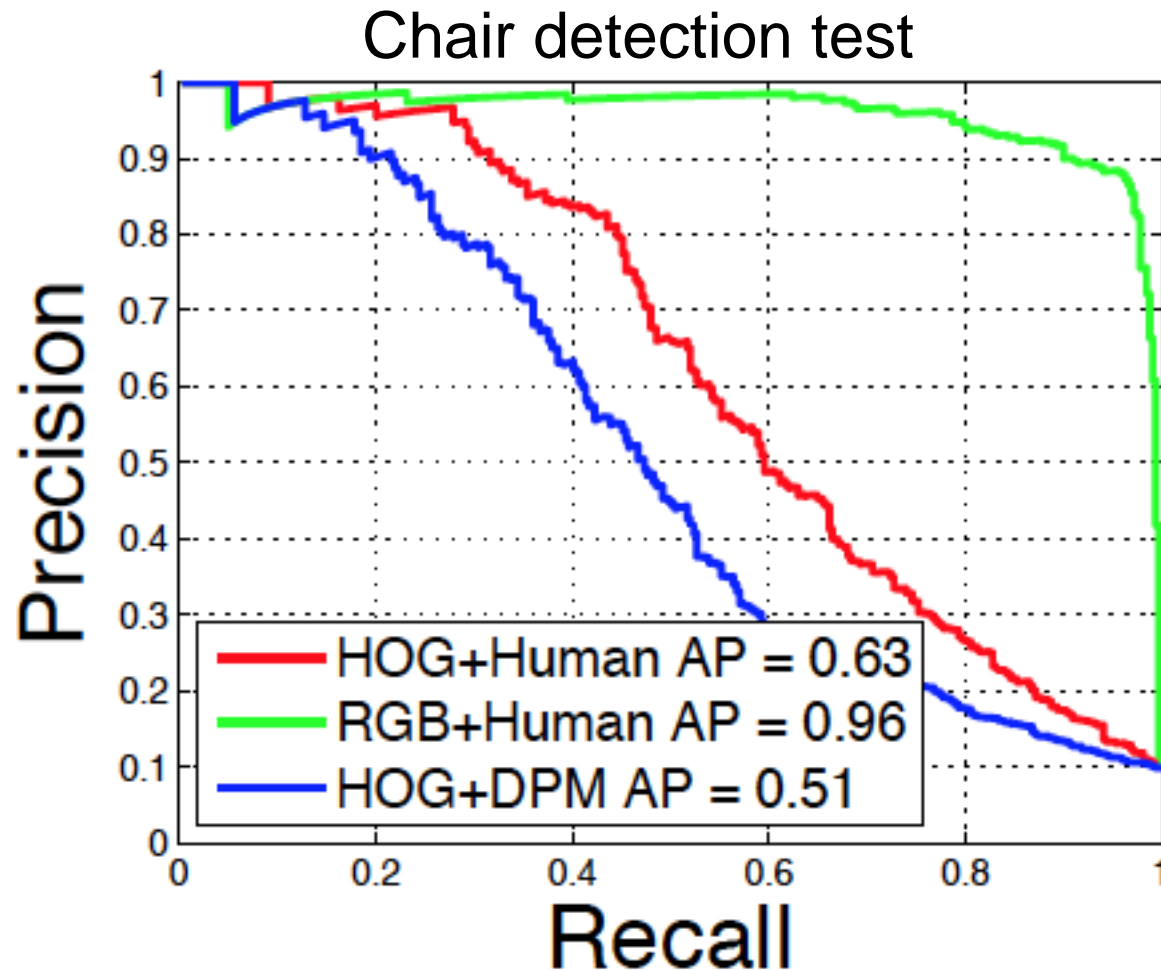
car

The image patch

http://mit.edu/vondrick/ihog/

Vondrick, Khosla, Malisiewicz, Torralba. "Inverting and Visualizing Features for Object Detection."

car

The image patch

What the detector sees

http://mit.edu/vondrick/ihog/
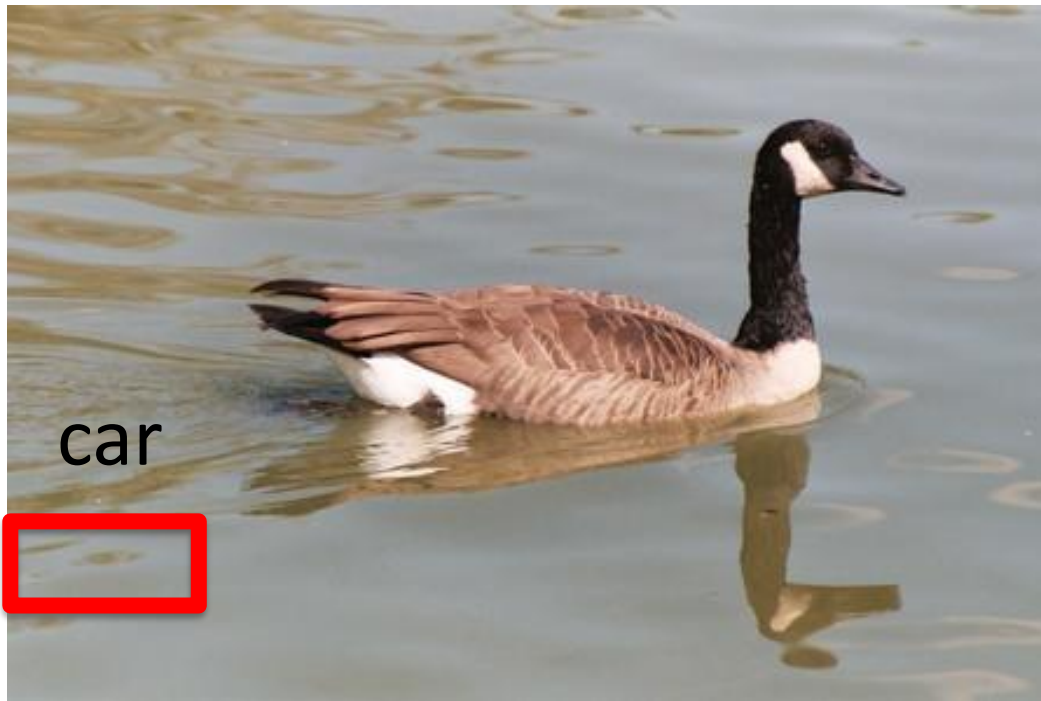
Vondrick, Khosla, Malisiewicz, Torralba. "Inverting and Visualizing Features for Object Detection."

# Deep architectures
## Geoffrey Hinton, Yann LeCun



Input image

Layer 1

Layer 2

Layer 3

Layer 4

Layer 5

FC6 FC7

FC8

orange

grapefruit

ball

...

...

shoe

Classifier output

# Scene recognition demo
## http://places.csail.mit.edu/demo.html



Users report 78% correct results

# http://places.csail.mit.edu/demo.html



places.csail.mit.edu

Take/Choose a photo

**Predictions**:

- **Type of environment:** outdoor
- **Semantic categories:**
  swimming_pool/outdoor:0.74,
  sandbar:0.11,

**Predictions**:

- **Type of environment:** indoor
- **Semantic categories:**
  airport_terminal:0.70,
- **SUN scene attributes:** enclosedarea,
  electricindoorlighting, nohorizon, man-
  made, congregating, cloth, glass,
  socializing, glossy,
  waitinginlinequeuing

the image uploaded follow Creative
Commons licenses.

Take/Choose a photo



**Predictions**:

- **Type of environment:** indoor
- **Semantic categories:** cockpit:0.08,
  parking_lot:0.06, playground:0.05,
- **SUN scene attributes:** nohorizon,
  enclosedarea, cloth, man-made,
  electricindoorlighting, working,
  stressful, dry, competing,
  waitinginlinequeuing

**http://places.csail.mit.edu/demo.html**



**Predictions:**
- **Type of environment:** indoor
- **Semantic categories:**
  auditorium:0.61,
  conference_center:0.34,

**Predictions**:

- **Type of environment:** indoor
- **Semantic categories:** bar:0.25, auditorium:0.20, restaurant_kitchen:0.07, coffee_shop:0.05,
- **SUN scene attributes:** enclosedarea, nohorizon, man-made, electricindoorlighting, wood(notpartofatree), working, matte, glass, cloth, conductingbusiness

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

**Take/Choose a photo**



**Predictions**:

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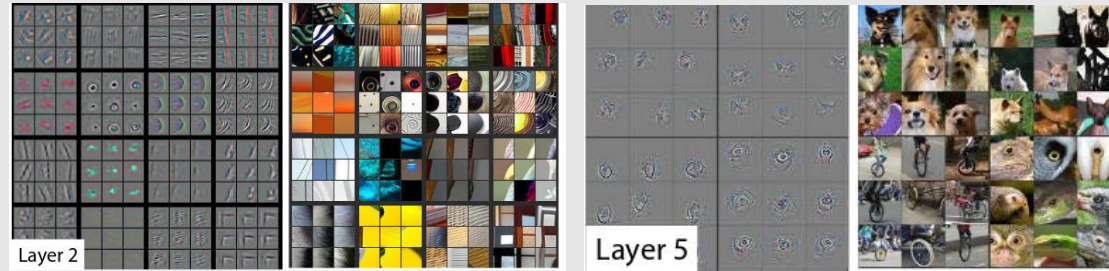  hotel_room:0.50, bedroom:0.47,

**Screenshot 1 (left):**

●●●●○ vodafone ES  3G  10:35 PM    ✈ 48% 🔋

places.csail.mit.edu

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,

**Screenshot 2 (right):**

●●○○○ vodafone ES  3G  10:31 PM    ✈ 49% 🔋⚡

places.csail.mit.edu

**Predictions**:

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

# Why is working so well?



But what is the representation built by the network?

# Visualizing the internal representation

**Deconvolution**



*Zeiler & Fergus, Visualizing and Understanding Convolutional Networks, ECCV 2014.*

**Backpropagation**



*Simonyan et al. Visualizing image classification models and saliency maps. ICLRW, 2014.*

**Strong activations**



*Girshick, et al, Rich feature hierarchies for accurate object detection and semantic segmentation. CVPR 2014.*

# Visualizing and Understanding Convolutional Networks

Matthew D. Zeiler and Rob Fergus

Dept. of Computer Science,
New York University, USA
{zeiler,fergus}@cs.nyu.edu

# Generative Adversarial Nets

**Ian J. Goodfellow,   Jean Pouget-Abadie,[*] Mehdi Mirza,  Bing Xu,  David Warde-Farley,
Sherjil Ozair,[†] Aaron Courville,  Yoshua Bengio[‡]**
Département d'informatique et de recherche opérationnelle
Université de Montréal
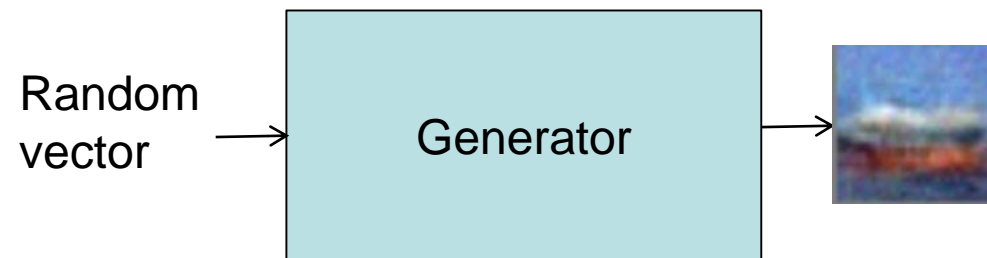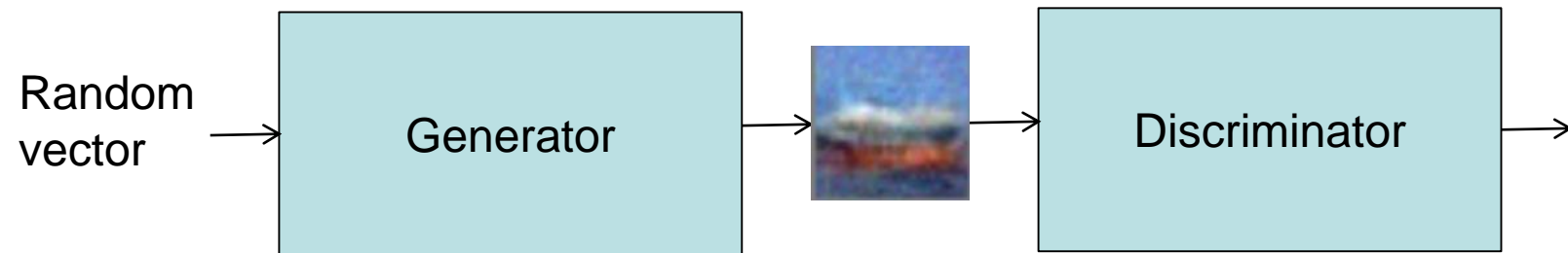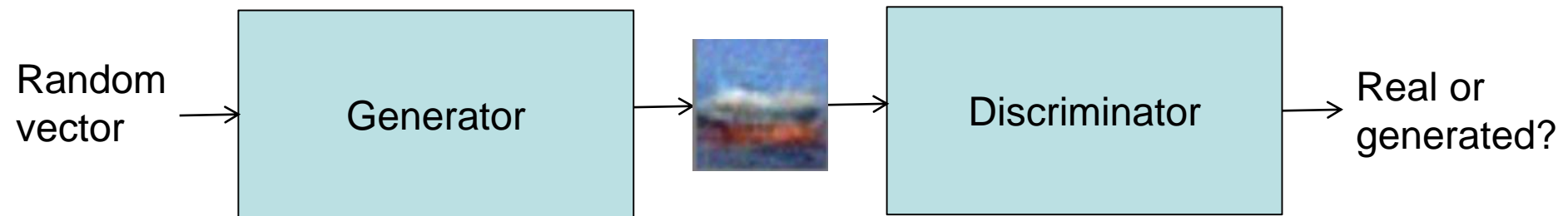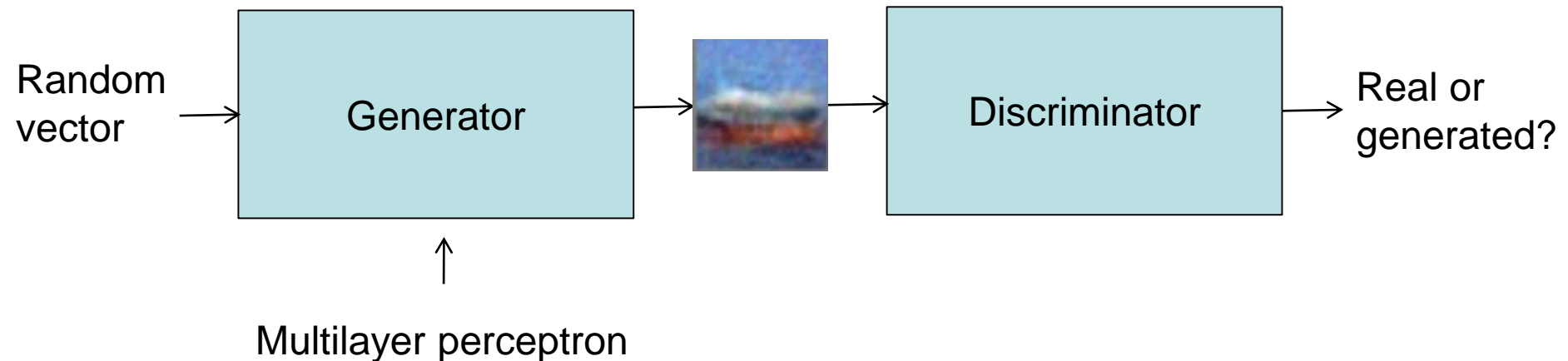Montréal, QC H3C 3J7

# Generative Adversarial Nets

**Ian J. Goodfellow,   Jean Pouget-Abadie,[*]  Mehdi Mirza,  Bing Xu,  David Warde-Farley,
Sherjil Ozair,[†]  Aaron Courville,  Yoshua Bengio[‡]**
Département d'informatique et de recherche opérationnelle
Université de Montréal
Montréal, QC H3C 3J7

Generator

# Generative Adversarial Nets

**Ian J. Goodfellow, Jean Pouget-Abadie,[*] Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair,[†] Aaron Courville, Yoshua Bengio[‡]**
Département d'informatique et de recherche opérationnelle
Université de Montréal
Montréal, QC H3C 3J7

Random vector → **Generator** →

# Generative Adversarial Nets

Ian J. Goodfellow,   Jean Pouget-Abadie,* Mehdi Mirza,  Bing Xu,  David Warde-Farley,
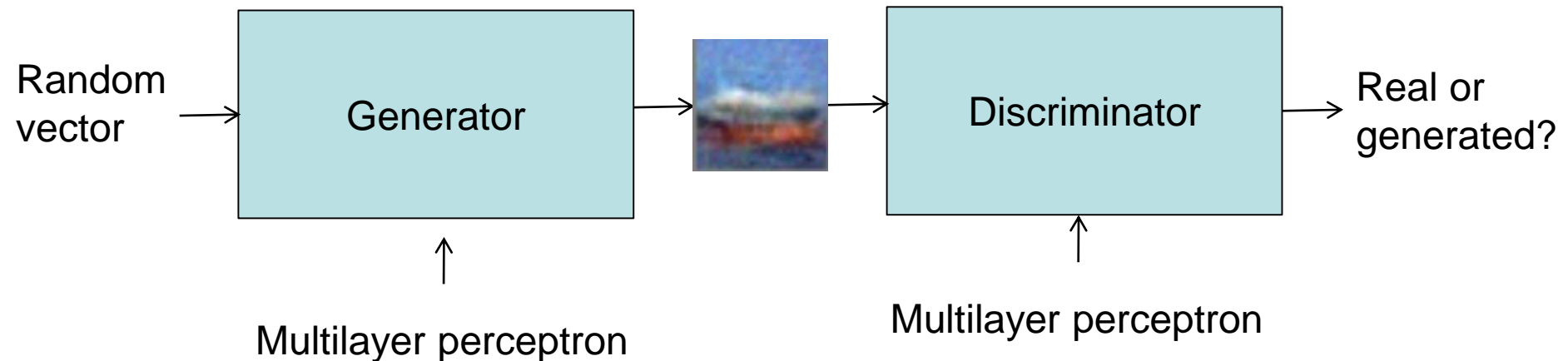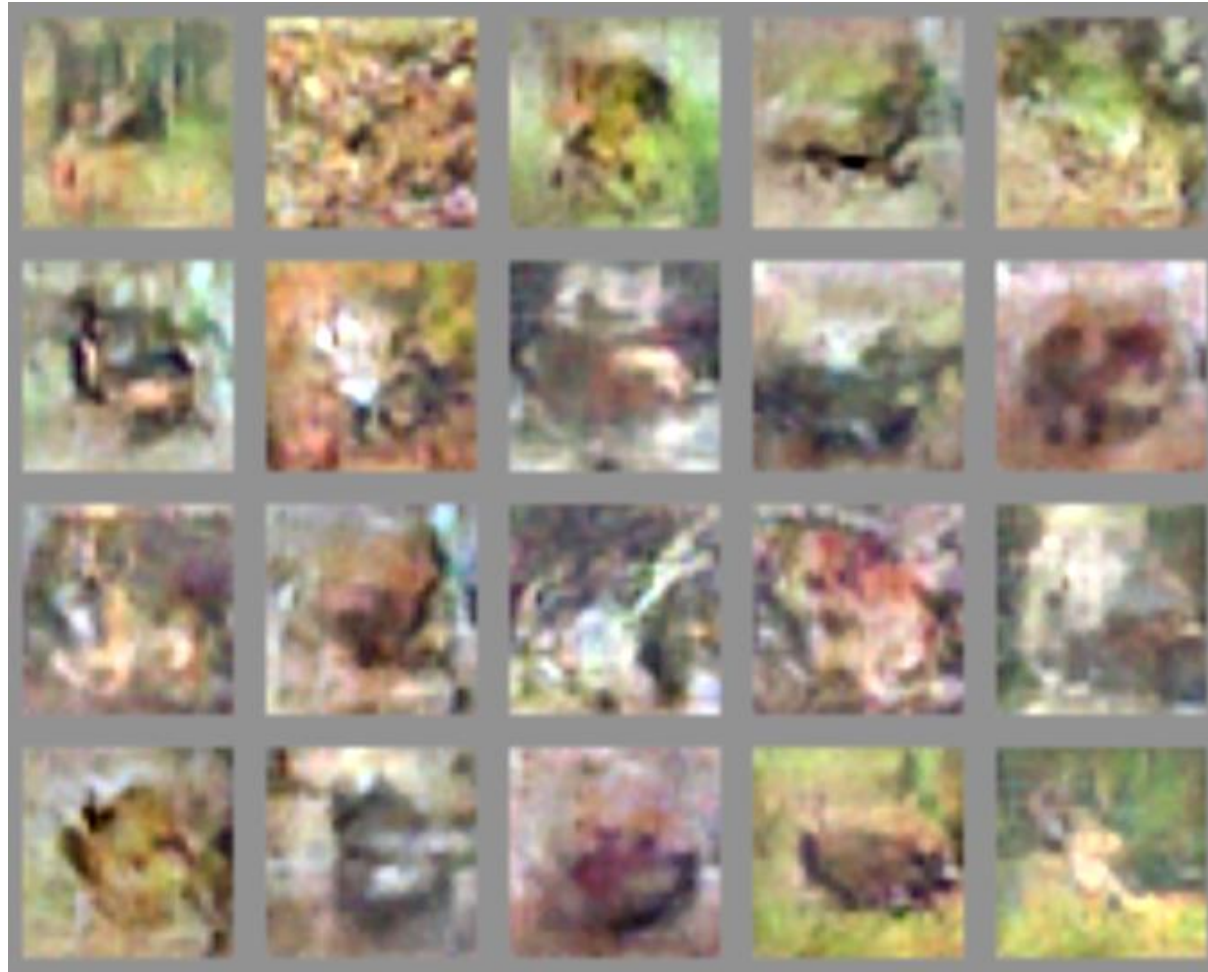Sherjil Ozair,† Aaron Courville,  Yoshua Bengio‡
Département d'informatique et de recherche opérationnelle
Université de Montréal
Montréal, QC H3C 3J7

# Generative Adversarial Nets

**Ian J. Goodfellow,   Jean Pouget-Abadie,[*] Mehdi Mirza,  Bing Xu,  David Warde-Farley,
Sherjil Ozair,[†] Aaron Courville,  Yoshua Bengio[‡]**
Département d'informatique et de recherche opérationnelle
Université de Montréal
Montréal, QC H3C 3J7

# Generative Adversarial Nets

Ian J. Goodfellow, Jean Pouget-Abadie,* Mehdi Mirza, Bing Xu, David Warde-Farley,
Sherjil Ozair,† Aaron Courville, Yoshua Bengio‡
Département d'informatique et de recherche opérationnelle
Université de Montréal
Montréal, QC H3C 3J7

# Generative Adversarial Nets

**Ian J. Goodfellow, Jean Pouget-Abadie,[*] Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair,[†] Aaron Courville, Yoshua Bengio[‡]**
Département d'informatique et de recherche opérationnelle
Université de Montréal
Montréal, QC H3C 3J7

Random vector → Generator → Discriminator → Real or generated?

Multilayer perceptron

# Generative Adversarial Nets

**Ian J. Goodfellow, Jean Pouget-Abadie,\* Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair,† Aaron Courville, Yoshua Bengio‡**
Département d'informatique et de recherche opérationnelle
Université de Montréal
Montréal, QC H3C 3J7

# Generated images



Trained with CIFAR-10

# UNSUPERVISED REPRESENTATION LEARNING WITH DEEP CONVOLUTIONAL GENERATIVE ADVERSARIAL NETWORKS

**Alec Radford & Luke Metz**
indico Research
Boston, MA
{alec,luke}@indico.io

**Soumith Chintala**
Facebook AI Research
New York, NY
soumith@fb.com

Introduced a form of ConvNet more stable under adversarial training than previous attempts.

# Generator

100 z

Project and reshape

1024

4

4

CONV 1

512

8

8

5

5

Stride 2

256

16

16

5

5

Stride 2

128

32

32

5

5

Stride 2

64

64

5

5

Stride 2

3

G(z)

CONV 2

CONV 3

CONV 4

# Generator



Random uniform
vector (100 numbers)

# Generator



Project and reshape

CONV 1

CONV 2

CONV 3

CONV 4

G(z)

100 z

1024

512

256

128

64

4

4

8

8

16

16

32

32

64

64

5

5

5

5

5

5

5

5

3

Stride 2

Stride 2

Stride 2

Stride 2

Random uniform
vector (100 numbers)

# Synthesizing the preferred inputs for neurons in neural networks via deep generator networks

**Anh Nguyen**
anguyen8@uwyo.edu

**Alexey Dosovitskiy**
dosovits@cs.uni-freiburg.de

**Jason Yosinski**
jason@geometricintelligence.com

**Thomas Brox**
brox@cs.uni-freiburg.de

**Jeff Clune**
jeffclune@uwyo.edu

# Two components

## Generator



Project and reshape · CONV 1 · CONV 2 · CONV 3 · CONV 4 · G(z)

## Network to visualize



conv1 · conv2 · conv3 · conv4 · conv5 · fc6 · fc7 · Classification layer · car

# Two components

## Generator



## Network to visualize

# Two components



Generator

Table lamp

Classification layer

conv1  conv2  conv3  conv4  conv5  fc6  fc7

100 z  Project and reshape  CONV 1  CONV 2  CONV 3  CONV 4  G(z)

1024  512  256  128  3

Stride 2  Stride 2  Stride 2  Stride 2

# Two components



Generator

Unit to visualize

Table lamp

# Two components



Generator

Unit to visualize

Table lamp

# Synthesizing Images Preferred by CNN

ImageNet-Alexnet-final units (class units)



Nguyen A, Dosovitskiy A, Yosinski J, Brox T, Clune J. (2016). "Synthesizing the preferred inputs for neurons in neural networks via deep generator networks.". arXiv:1605.09304.

# Object detection vs. Scene recognition

# Object detection vs. Scene recognition

# Object detection vs. Scene recognition

# Object detection vs. Scene recognition



Bird

# Object detection vs. Scene recognition

Bird

Bedroom

# IM🔲GENET

- An ontology of images based on WordNet

- ImageNet currently has
  - 13,000+ categories of visual concepts
  - 10 million human-cleaned images (~700im/categ)
  - 1/3+ is released online @ **www.image-net.org**



IM🔲GENET

$\sim 10^5 +$ nodes
$\sim 10^8 +$ images

animal

shepherd dog, sheep dog

collie

**German shepherd**

Deng, Dong, Socher, Li & Fei-Fei, CVPR 2009

# places

places.csail.mit.edu

Zhou    Lapedriz    Khosla    Xiao    Oliva

# places

places.csail.mit.edu

1. We take all scene words
from a dictionary

**WordNet Dictionary**
A lexical database for the English language

WordNet® is an online lexical
reference system whose design is
inspired by current
psycholinguistic theories of
human lexical memory.

Zhou    Lapedriz    Khosla    Xiao    Oliva

Zhou, Lapedriza, Xiao, Oliva & Torralba (NIPS 2014)

# places

places.csail.mit.edu

1. We take all scene words from a dictionary

**WordNet Dictionary**
A lexical database for the English language

WordNet® is an online lexical reference system whose design is inspired by current psycholinguistic theories of human lexical memory.

2. We download images and clean the categories

Google Image Search

altavista

flickr

Zhou    Lapedriz    Khosla    Xiao    Oliva

# Two large databases, two tasks



IMAGENET

brambling

terrier

**Places Database**

bedroom

mountain

# ImageNet CNN and Places CNN

# ImageNet CNN and Places CNN



**ImageNet CNN for Object Classification**

# ImageNet CNN and Places CNN



**ImageNet CNN for Object Classification**

Layer 1

Layer 2

Layer 3  Layer 4  Layer 5

FC6 FC7  FC8

→ **orange**
→ grapefruit
→ ball
...
...
→ shoe

IM**A**GENET

Same architecture: AlexNet

**Places CNN for Scene Classification**

Layer 1

Layer 2

Layer 3  Layer 4  Layer 5

FC6 FC7  FC8

→ **savannah**
→ field
→ lake
...
...
→ kitchen

**Places**

# Possible internal representations



[Deng et al. CVPR 2009]

**PLACES**

# Learning to Recognize Objects



brambling

terrier

# Learning to Recognize Objects

IM**A**GENET

brambling

terrier

## Possible internal representations:

- Object parts
- Textures
- Attributes

HAIR

EYE    EYE

LEFT
EDGE    RIGHT
EDGE

NOSE

MOUTH

# Learning to Recognize Scenes

# Learning to Recognize Scenes

bedroom

mountain

Possible internal representations:

- Scene parts
- Objects
- Scene attributes
- Object parts
- Textures

# Places and objects

# Places and objects

# Places and objects



Features + SVM

# Places and objects



**Scene datasets**

|  | SUN397 | MIT Indoor67 | Scene15 | SUN Attribute |
|---|---|---|---|---|
| Places-CNN feature | **54.32±0.14** | **68.24** | **90.19±0.34** | **91.29** |
| ImageNet-CNN feature | 42.61±0.16 | 56.79 | 84.23±0.37 | 89.85 |

**Zhou, Lapedriza, Xiao, Torralba & Oliva (NIPS 2014)**

# Places and objects



Features + SVM

**Scene datasets**

|                      | SUN397         | MIT Indoor67 | Scene15        | SUN Attribute |
| -------------------- | -------------- | ------------ | -------------- | ------------- |
| Places-CNN feature   | **54.32±0.14** | **68.24**    | **90.19±0.34** | **91.29**     |
| ImageNet-CNN feature | 42.61±0.16     | 56.79        | 84.23±0.37     | 89.85         |

**Object datasets**

|                      | Caltech101     | Caltech256     | Action40       | Event8         |
| -------------------- | -------------- | -------------- | -------------- | -------------- |
| Places-CNN feature   | 65.18±0.88     | 45.59±0.31     | 42.86±0.25     | 94.12±0.99     |
| ImageNet-CNN feature | **87.22±0.92** | **67.23±0.27** | **54.92±0.33** | **94.42±0.76** |

**Zhou, Lapedriza, Xiao, Torralba & Oliva (NIPS 2014)**

# Preferred images

# Preferred images



ImageNet-CNN

Pool 1

# Preferred images

# Preferred images

# Preferred images

# Preferred images

# Preferred images



ImageNet-CNN

Places-CNN

Pool 1

Pool 2

conv3

# Preferred images

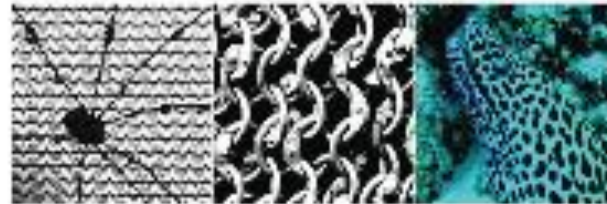# Preferred images



ImageNet-CNN

Places-CNN

Pool 1
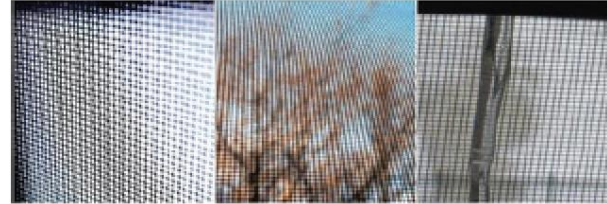
Pool 2
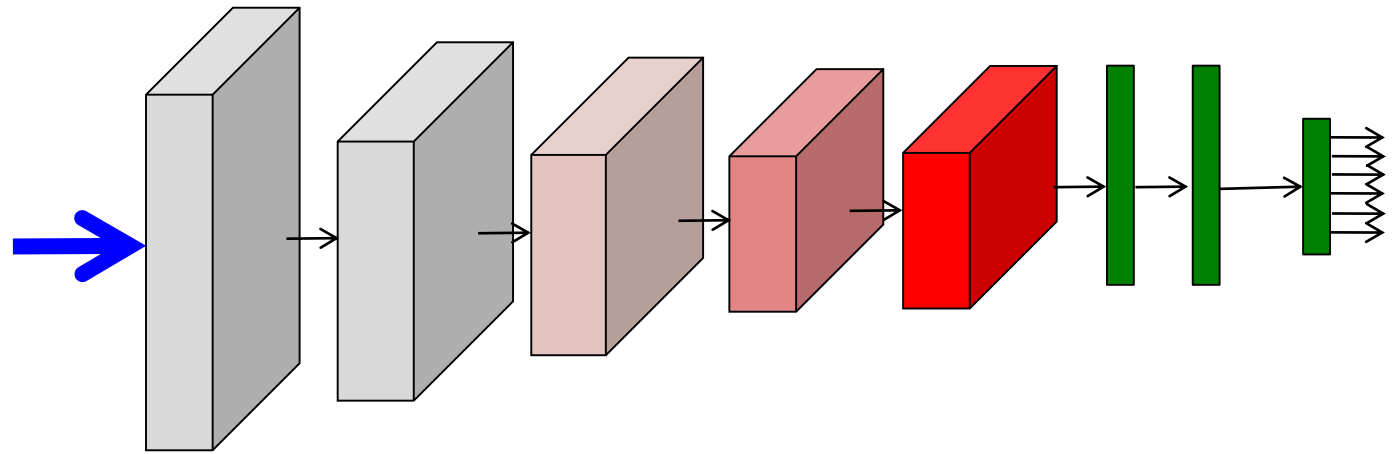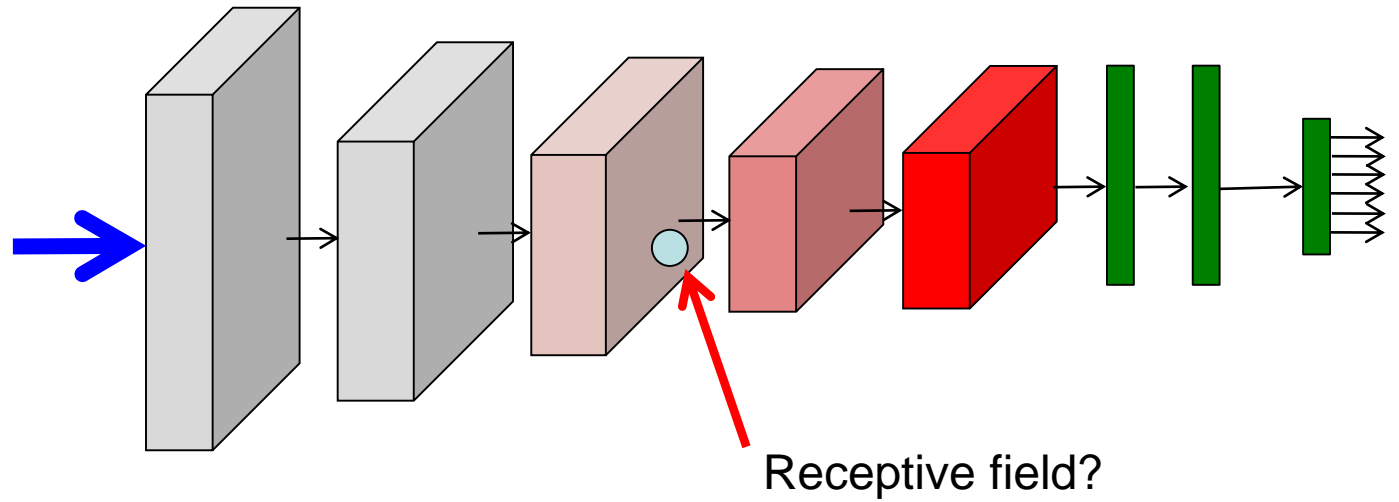
conv3

conv 4
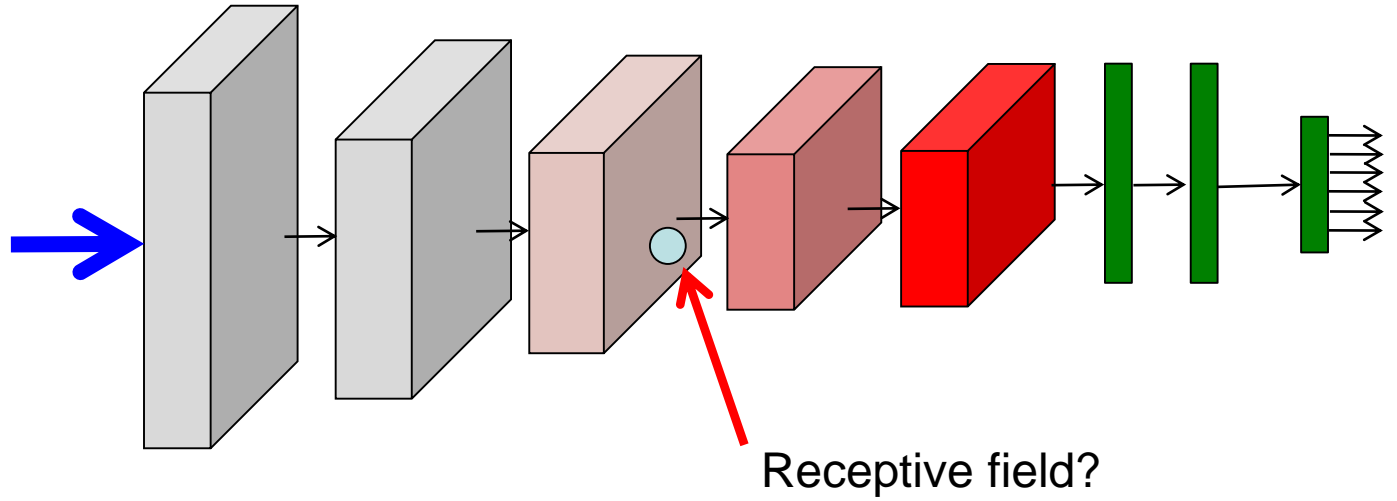
# Preferred images



ImageNet-CNN          Places-CNN

Pool 1

Pool 2

conv3

conv 4

Pool 5

# Preferred images



ImageNet-CNN          Places-CNN

Pool 1

Pool 2

conv3

conv 4

Pool 5

# Estimating the receptive field

# Estimating the receptive field



Receptive field?

# Estimating the receptive field



Receptive field?

# Estimating the receptive field



Receptive field?

# Estimating the receptive field



Receptive field?

# Estimating the receptive field

**Theoretical size**

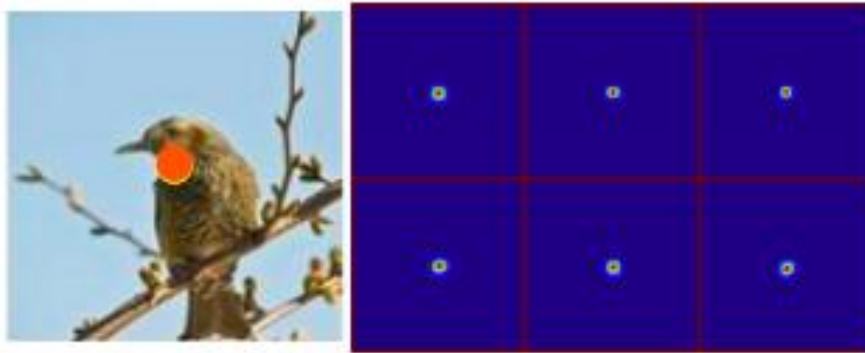**Actual size**

# Estimating the receptive field
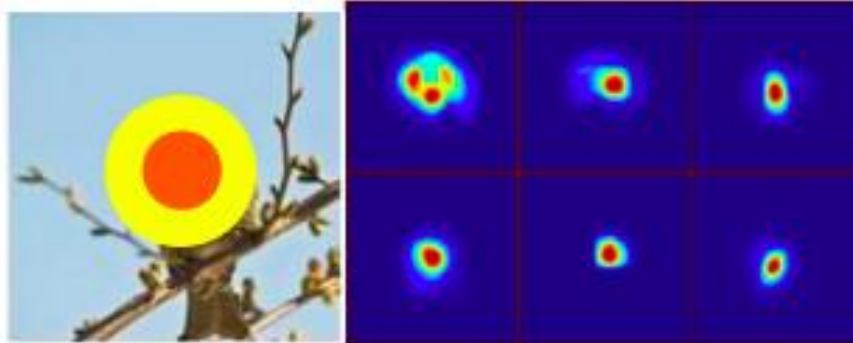


Layer 1

**Theoretical size**

**Actual size**

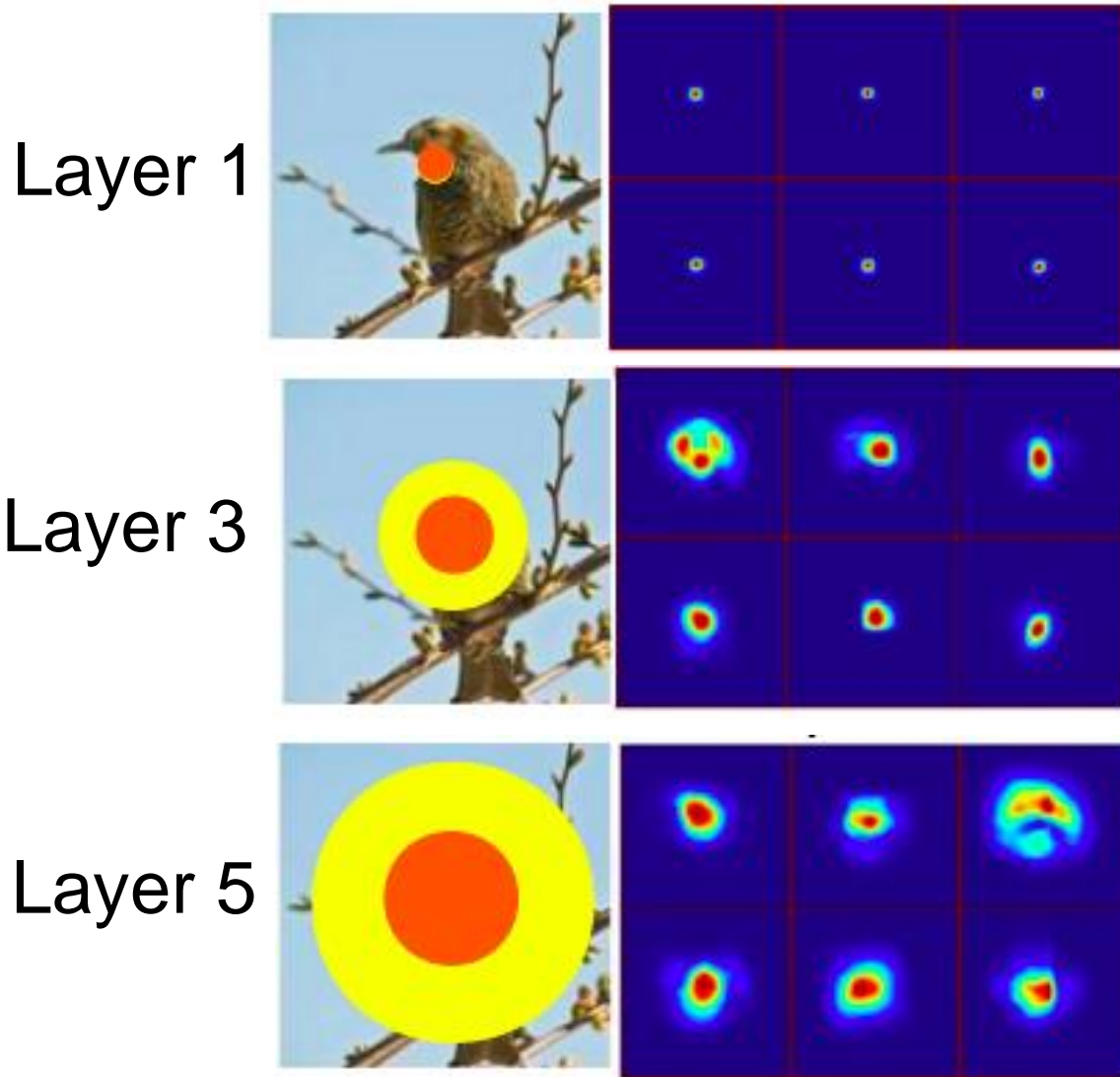# Estimating the receptive field

Layer 1

Layer 3

**Theoretical size**

**Actual size**

# Estimating the receptive field

Layer 1

Layer 3

Layer 5

**Theoretical size**

**Actual size**

# Generating segmentations

| | | |
|---|---|---|
| 0.1 | 0.7 | 0.5 |
| 0.3 | 0.6 | 0.4 |
| 0 | 0.1 | 0 |
| 0 | 0 | 0 |

feature map

# Generating segmentations

| | | |
|---|---|---|
| 0.1 | 0.7 | 0.5 |
| 0.3 | 0.6 | 0.4 |
| 0 | 0.1 | 0 |
| 0 | 0 | 0 |

feature map

0.1*        +    

# Generating segmentations

| 0.1 | 0.7 | 0.5 |
|-----|-----|-----|
| 0.3 | 0.6 | 0.4 |
| 0   | 0.1 | 0   |
| 0   | 0   | 0   |

feature map

0.7*  +

# Generating segmentations

| | | |
|---|---|---|
| 0.1 | 0.7 | 0.5 |
| 0.3 | 0.6 | 0.4 |
| 0 | 0.1 | 0 |
| 0 | 0 | 0 |

feature map

0.5*

# Generating segmentations

| | | |
|---|---|---|
| 0.1 | 0.7 | 0.5 |
| 0.3 | 0.6 | 0.4 |
| 0 | 0.1 | 0 |
| 0 | 0 | 0 |

feature map

# Crowdsourcing units



**Task 1**

Word/Short description:

lighthouse

**Task 2**

Mark (by clicking on them) the images which don't correspond to the short description you just wrote

# Crowdsourcing units



## Task 3

Which category does your short description mostly belong to?

○ Scene (kitchen, corridor, street, beach, ...)

○ Region or surface (road, grass, wall, floor, sky, ...)

◉ Object (bed, car, building, tree, ...)

○ Object part (leg, head, wheel, roof, ...)

○ Texture or material (striped, rugged, wooden, plastic, ...)

○ Simple elements or colors (vertical line, curved line, color blue, ....)

# Annotating the Semantics of Units

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

# Annotating the Semantics of Units

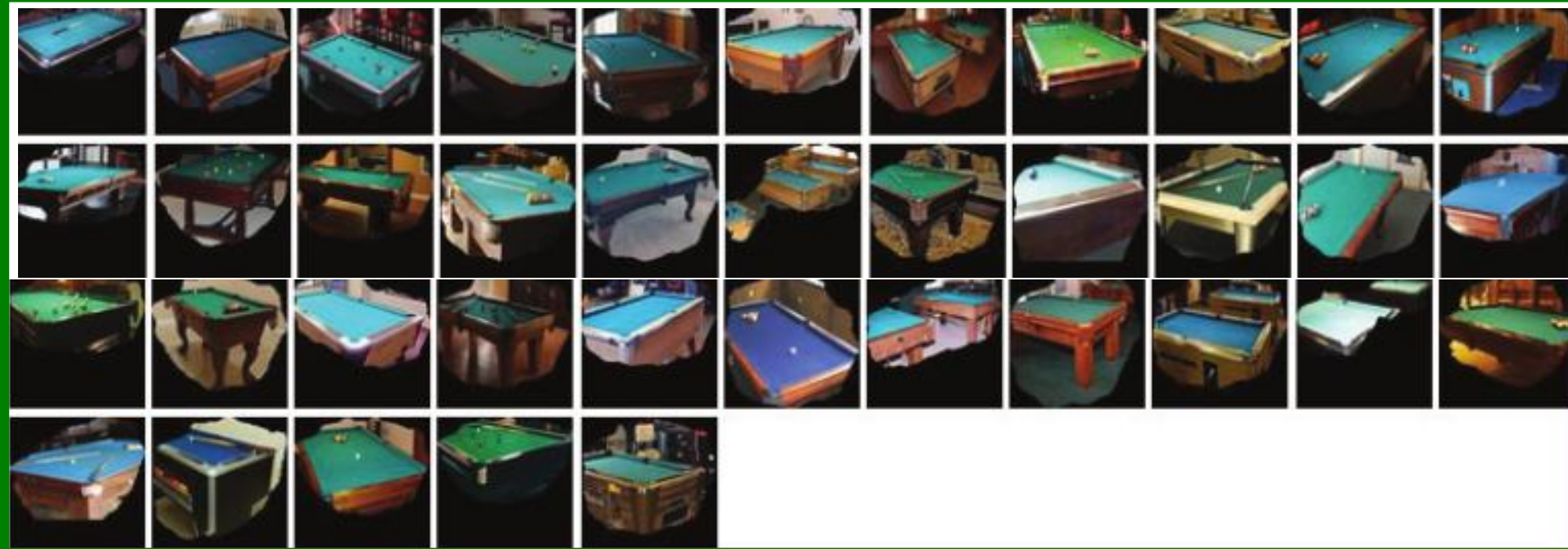Pool5, unit 13; Label: Lamps; Type: object; Precision: 84%

# Annotating the Semantics of Units

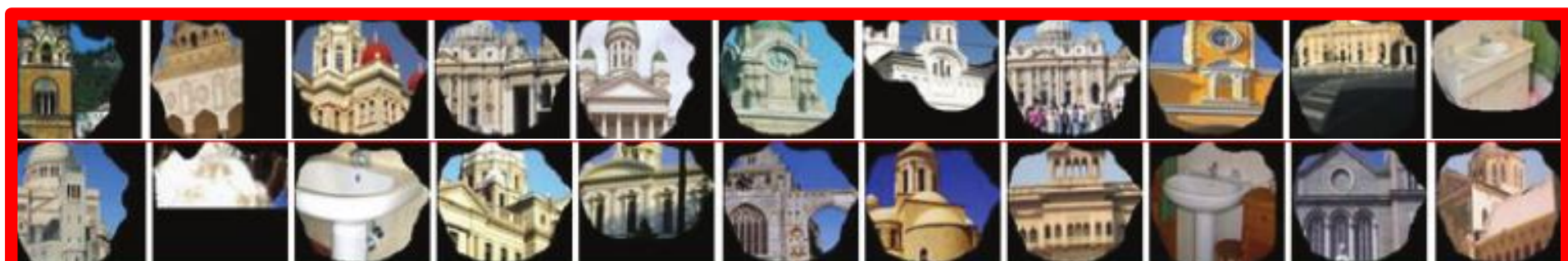Pool5, unit 77; Label: legs; Type: object part; Precision: 96%

# Annotating the Semantics of Units

Pool5, unit 112; Label: pool table; Type: object; Precision: 70%
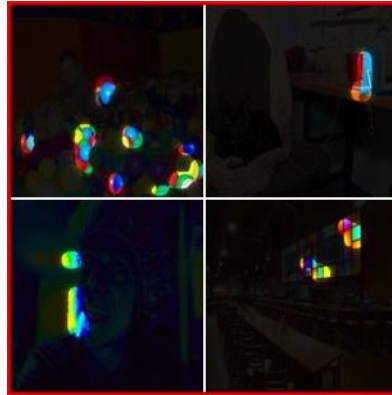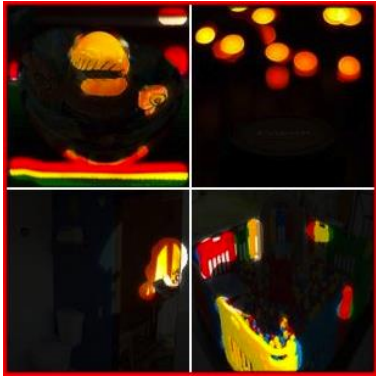
# Annotating the Semantics of Units

Pool5, unit 22; Label: dinner table; Type: scene; Precision: 60%

# Distribution of semantic types at each layer
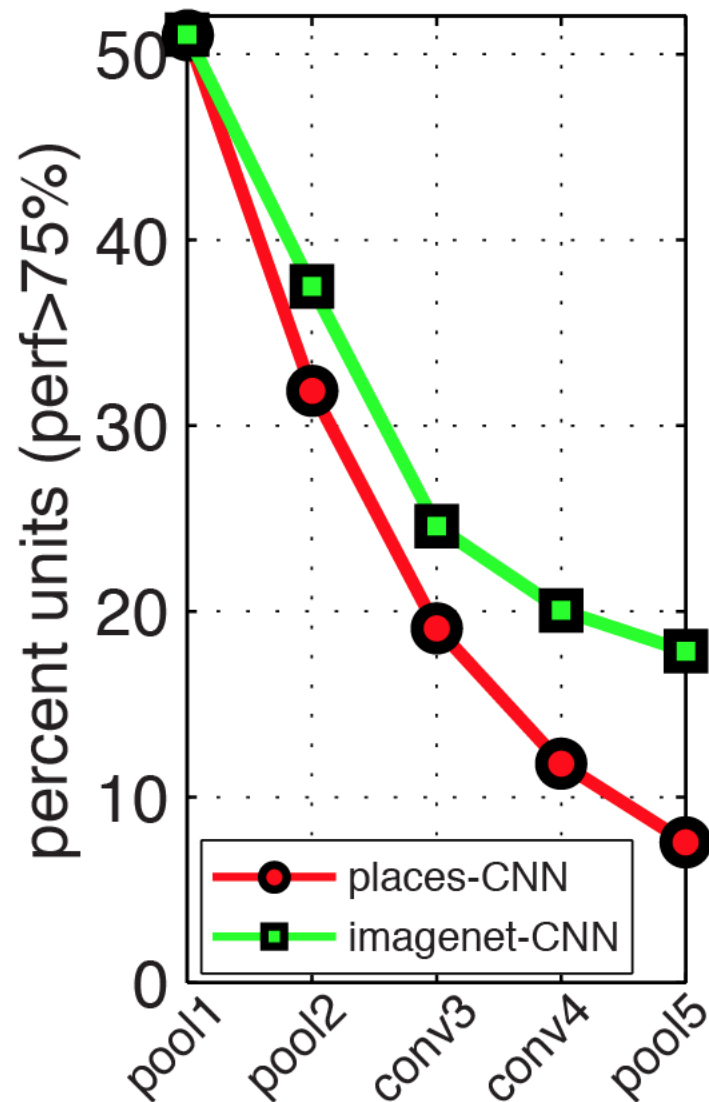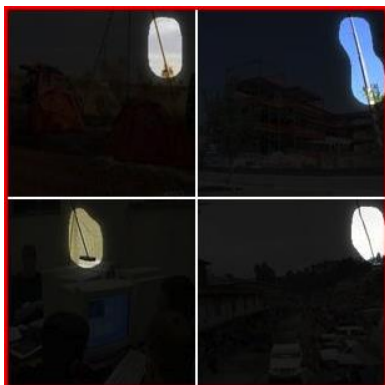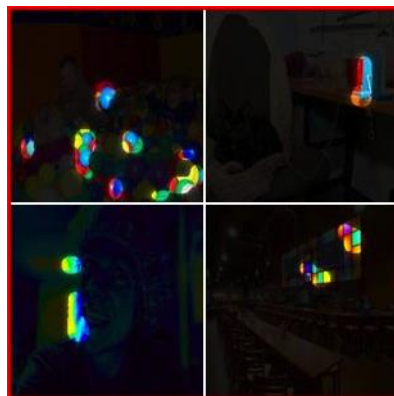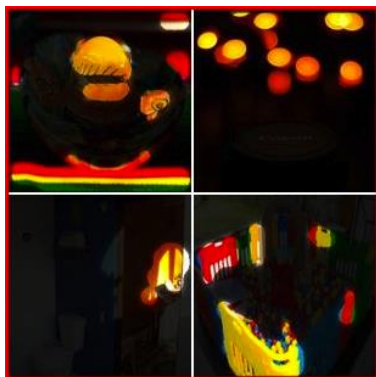# 1 - Simple elements and colors

Ex: vertical line, curved line, color blue, ….

# Distribution of semantic types at each layer
# 1 - Simple elements and colors

Ex: vertical line, curved line, color blue, ….

# Distribution of semantic types at each layer

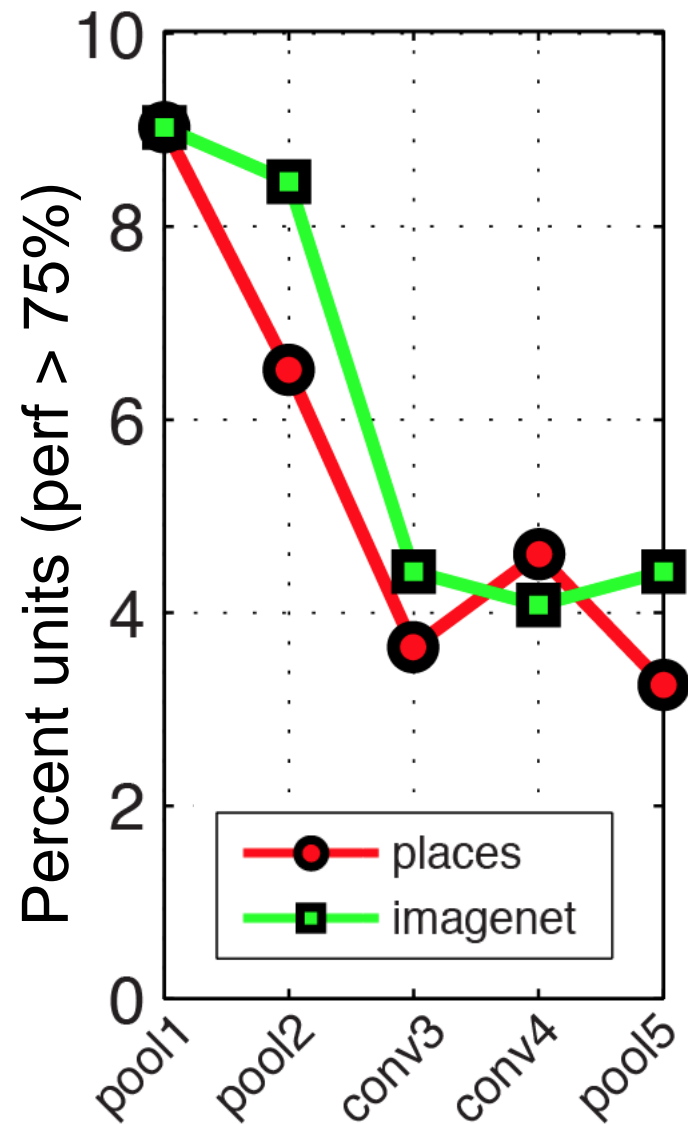## 2 - Texture or materials

Ex: stripes, wooden, plastic, ...



Percent units (perf > 75%)

places
imagenet

# Distribution of semantic types at each layer

# 2 - Texture or materials

Ex: stripes, wooden, plastic, ...

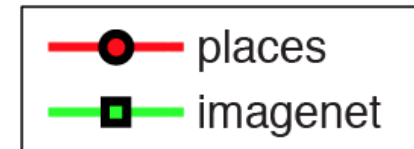# Distribution of semantic types at each layer

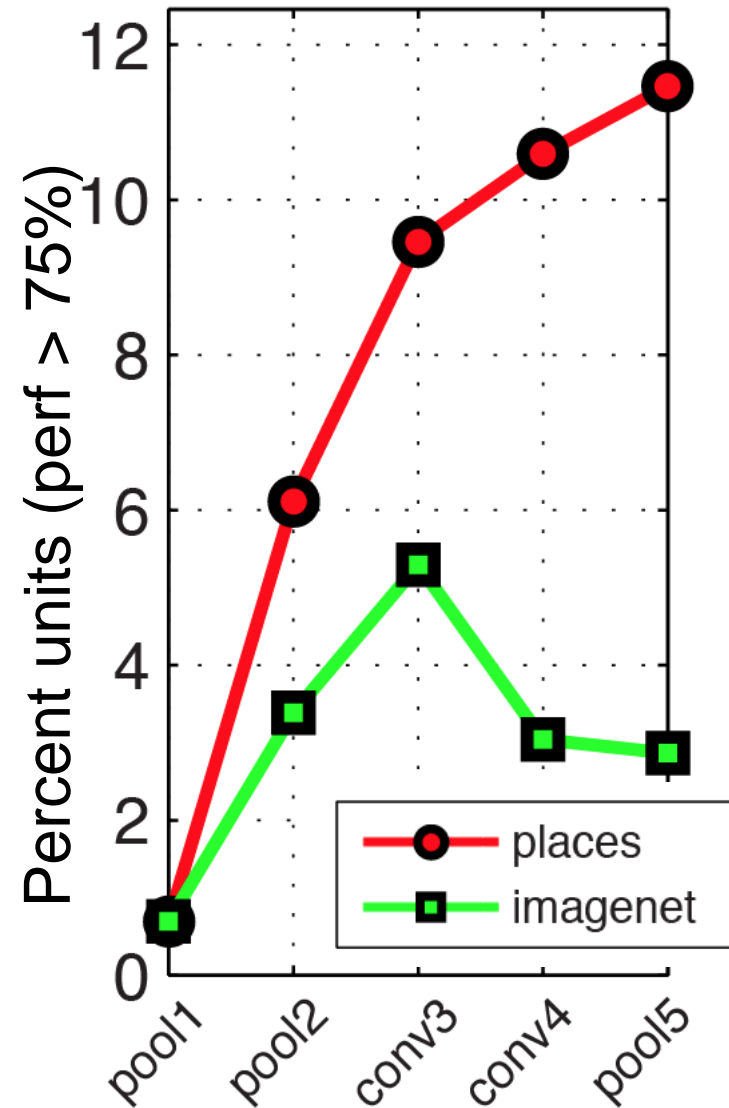# 3 - Regions and surfaces

Ex: Road, grass, wall, floor, sky, ….



Percent units (perf > 75%)

places
imagenet

# 3 - Regions and surfaces



Ex: Road, grass, wall, floor, sky, ….

# 4 - Object parts

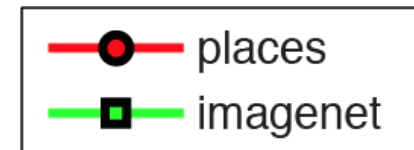Ex: leg, head, wheel, roof, ….



Percent units (perf > 75%)

places
imagenet
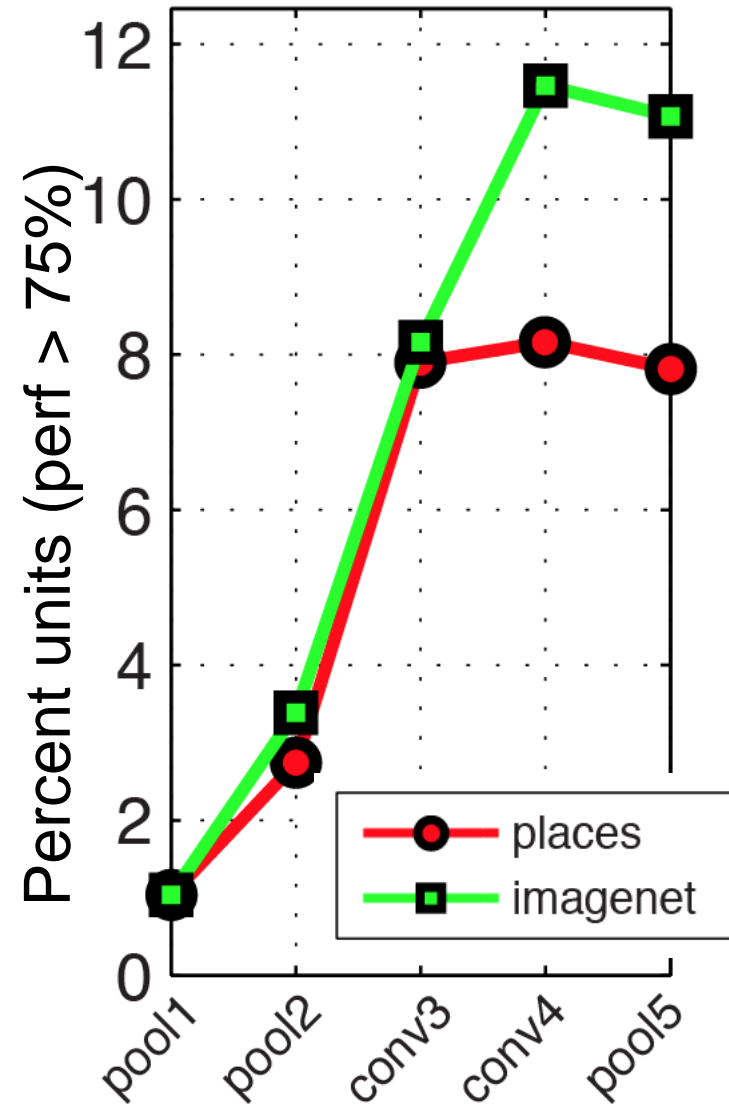
# Distribution of semantic types at each layer

# 4 - Object parts

Ex: leg, head, wheel, roof, ….

# Distribution of semantic types at each layer
# 5 - Objects

Ex: bed, car, building, tree, ….
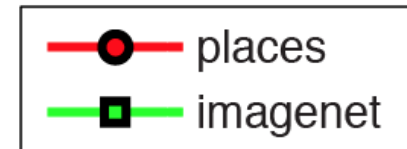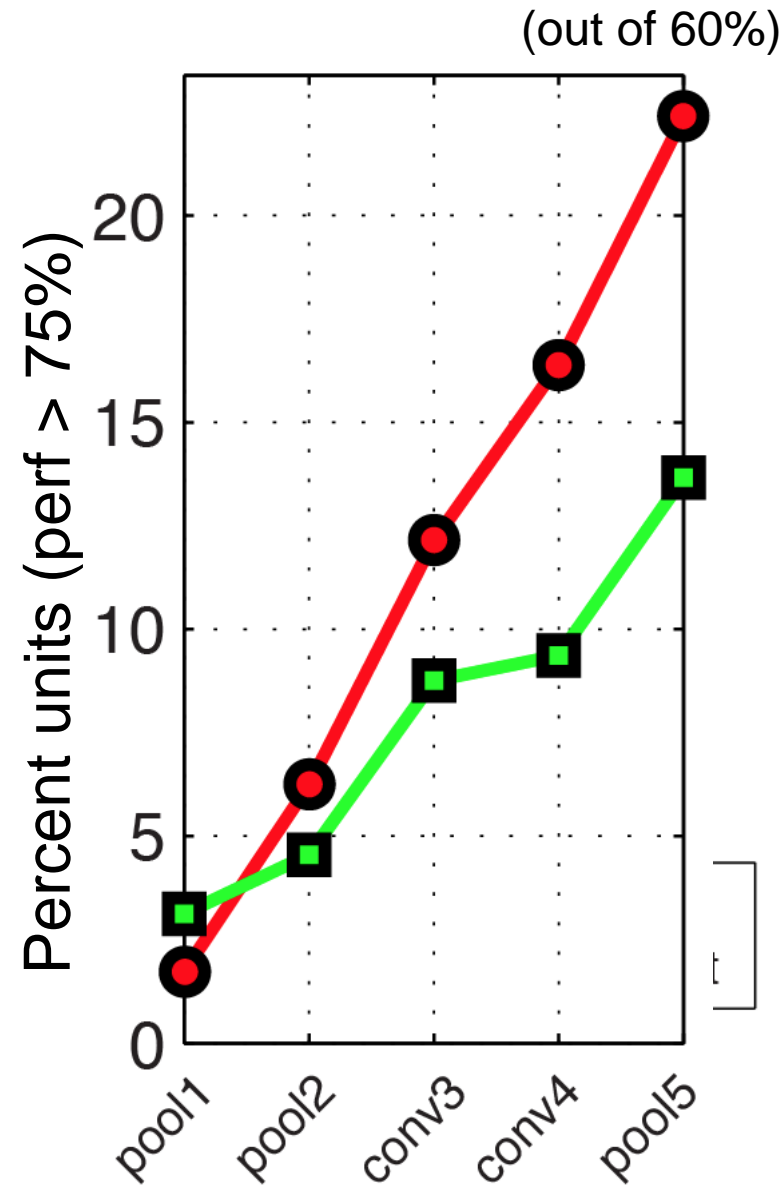


Percent units (perf > 75%)

places
imagenet

# Distribution of semantic types at each layer

## 5 - Objects

Ex: bed, car, building, tree, ....

(out of 60%)

# Distribution of semantic types at each layer
# 6 - Scenes

Ex: kitchen, corridor, street, beach, ….
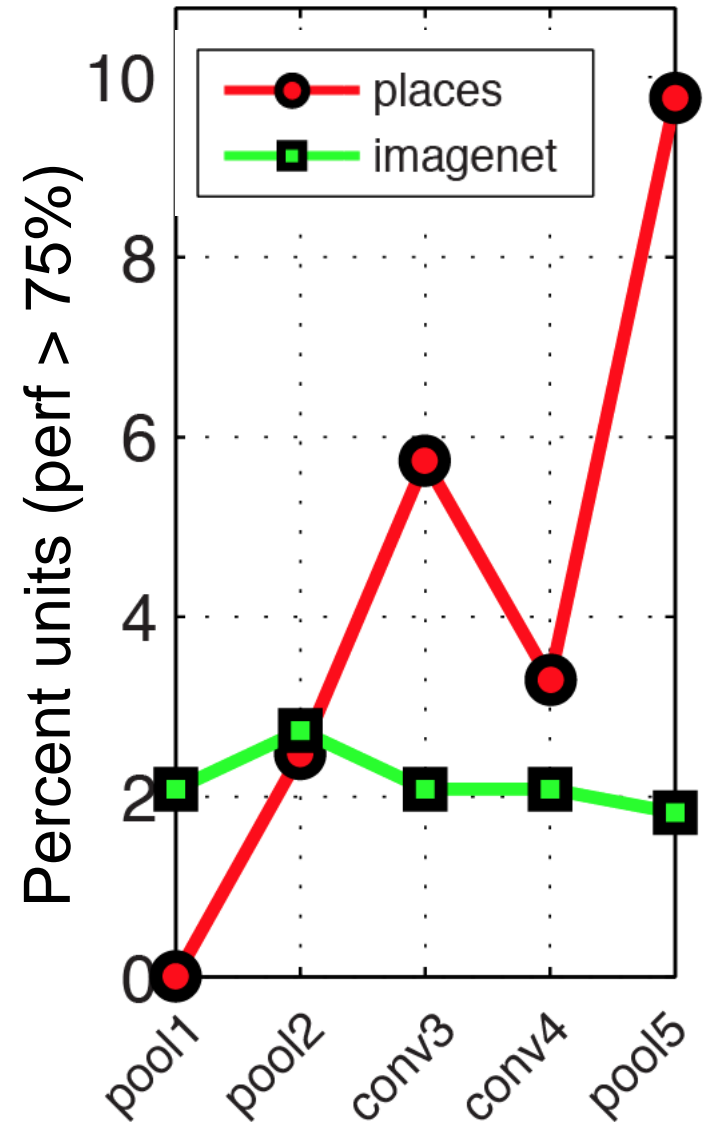


Percent units (perf > 75%)

places
imagenet

# Distribution of semantic types at each layer

# **6 - Scenes**

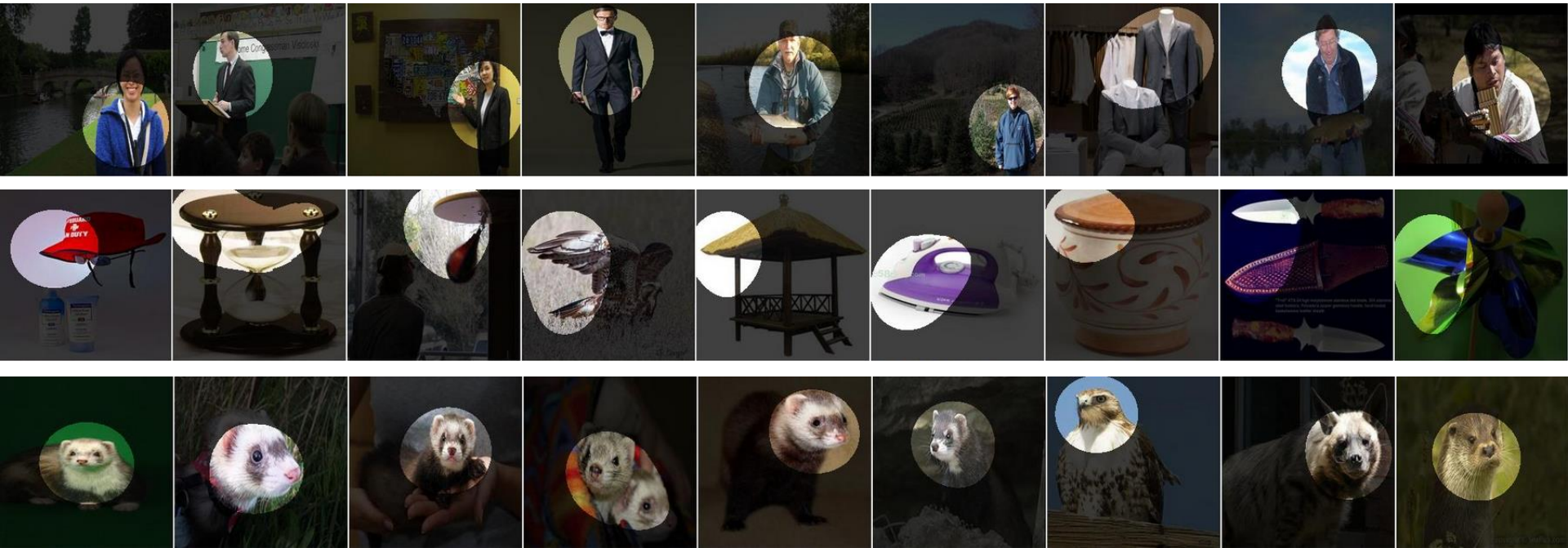Ex: kitchen, corridor, street, beach, ….

# What objects are found?
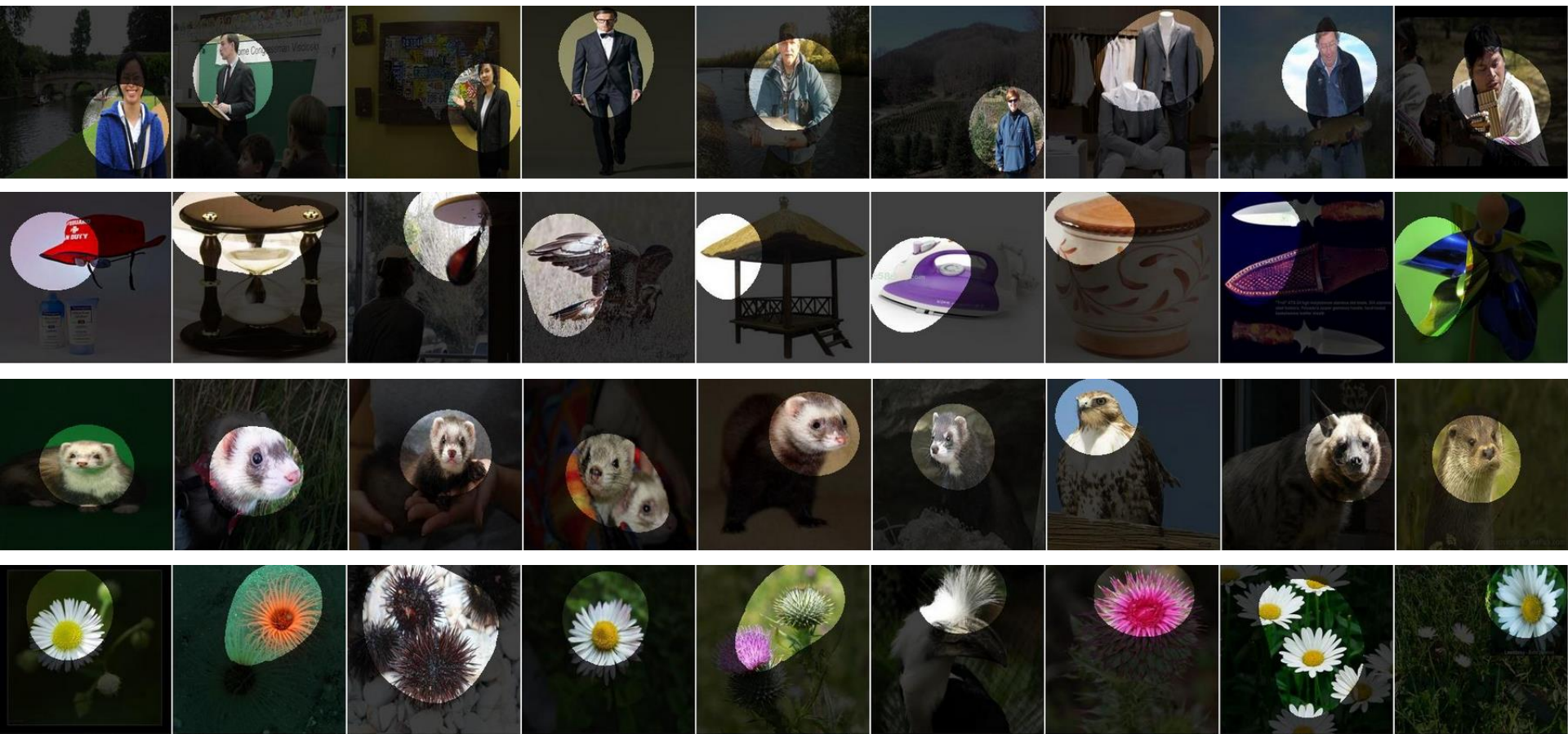
# ImageNet-CNN Units
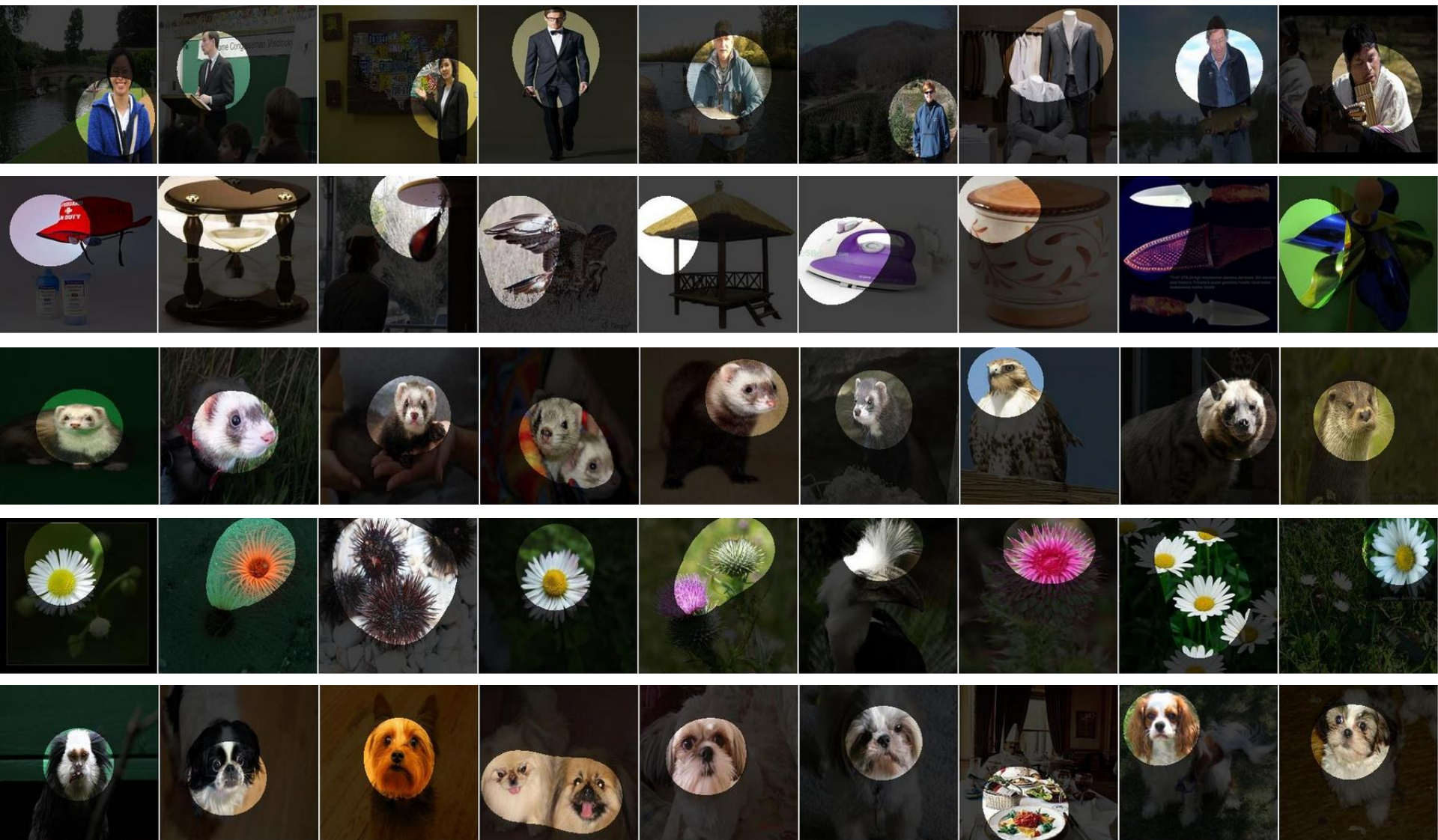
# ImageNet-CNN Units

# ImageNet-CNN Units

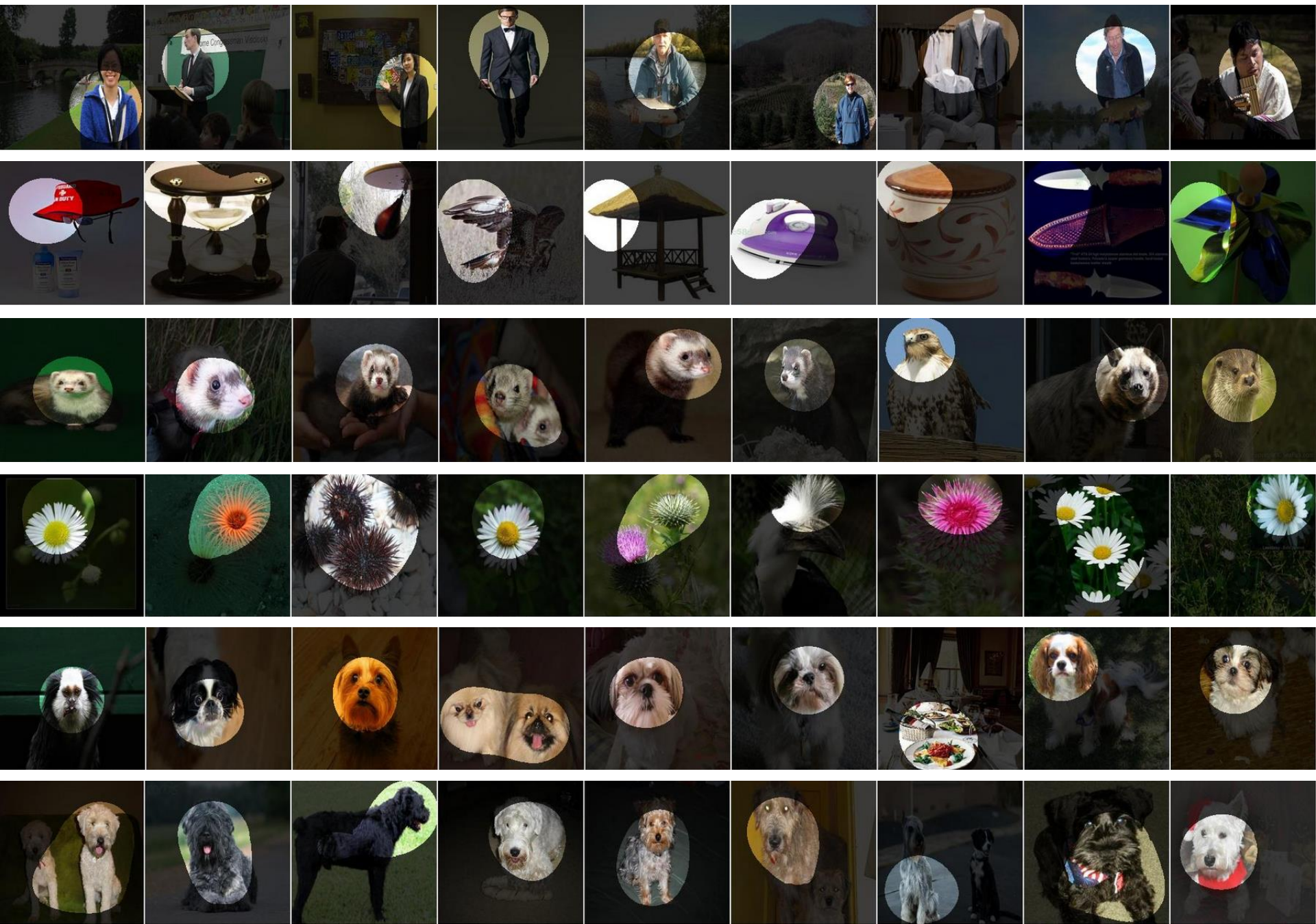# ImageNet-CNN Units

# ImageNet-CNN Units

# ImageNet-CNN Units

# ImageNet-CNN Units

# Places-CNN Units

# Places-CNN Units
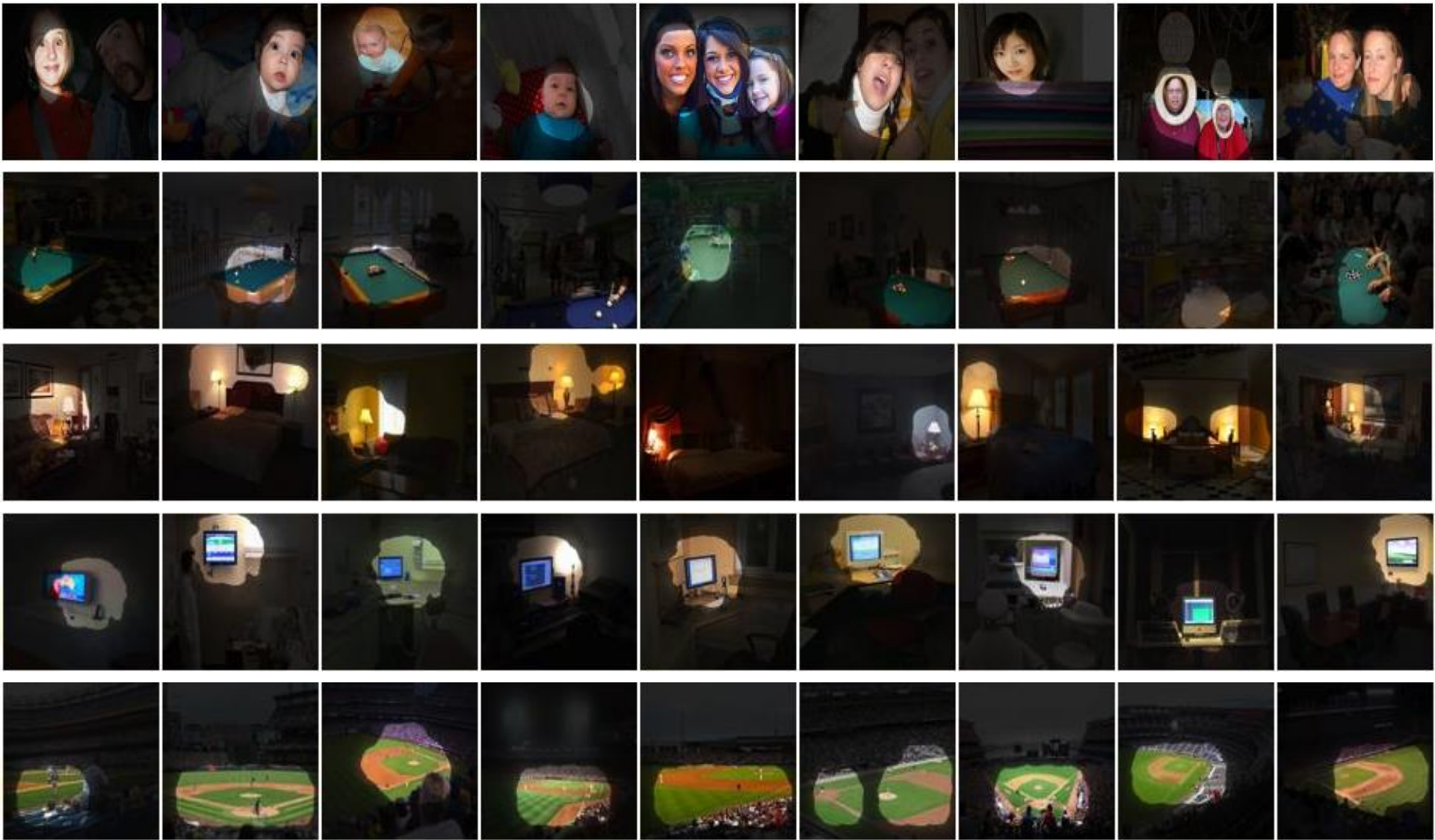
# Places-CNN Units

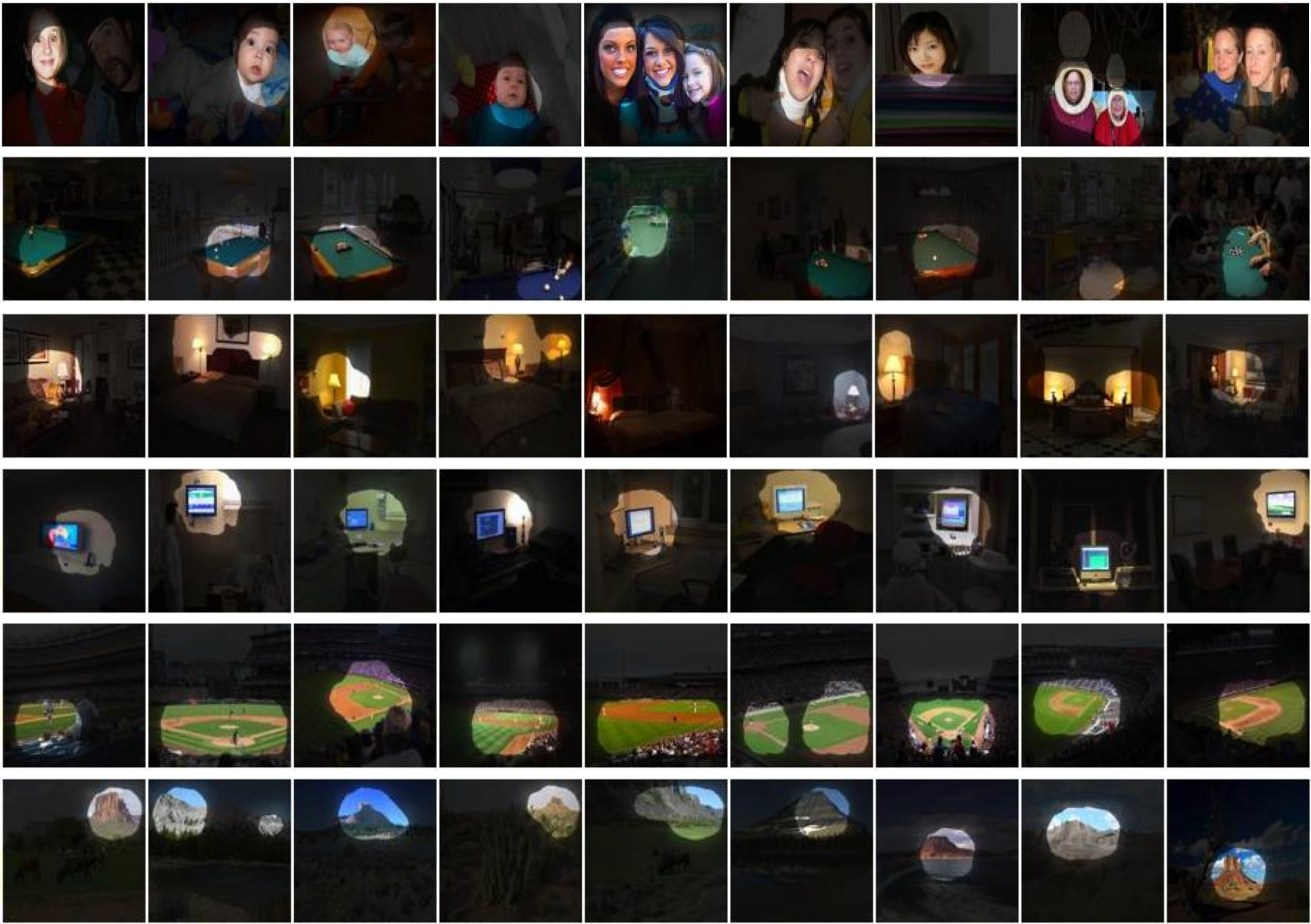# Places-CNN Units

# Places-CNN Units

# Places-CNN Units

# Places-CNN Units

# Histogram of Emerged Objects in Pool5

ImageNet-CNN (59/256)



Includes: Objects, nameable parts, and regions

# Histogram of Emerged Objects in Pool5



Includes: Objects, nameable parts, and regions

# Histogram of Emerged Objects in Pool5

ImageNet-CNN (59/256)



Includes: Objects, nameable parts, and regions

# Histogram of Emerged Objects in Pool5

ImageNet-CNN (59/256)



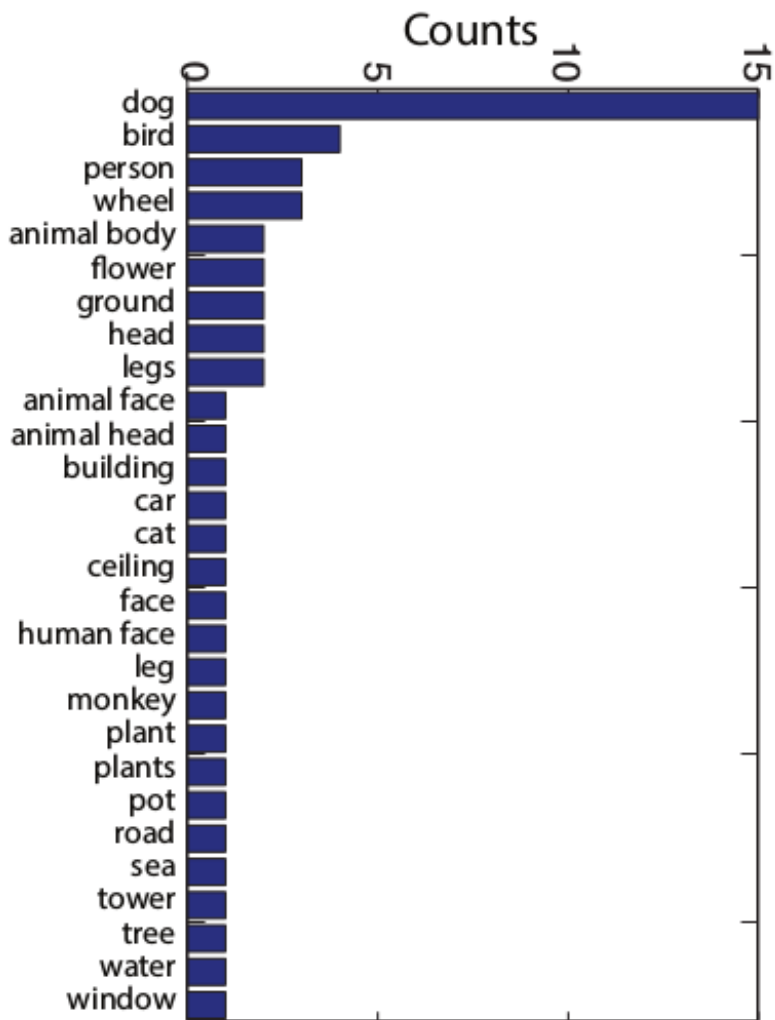Includes: Objects, nameable parts, and regions

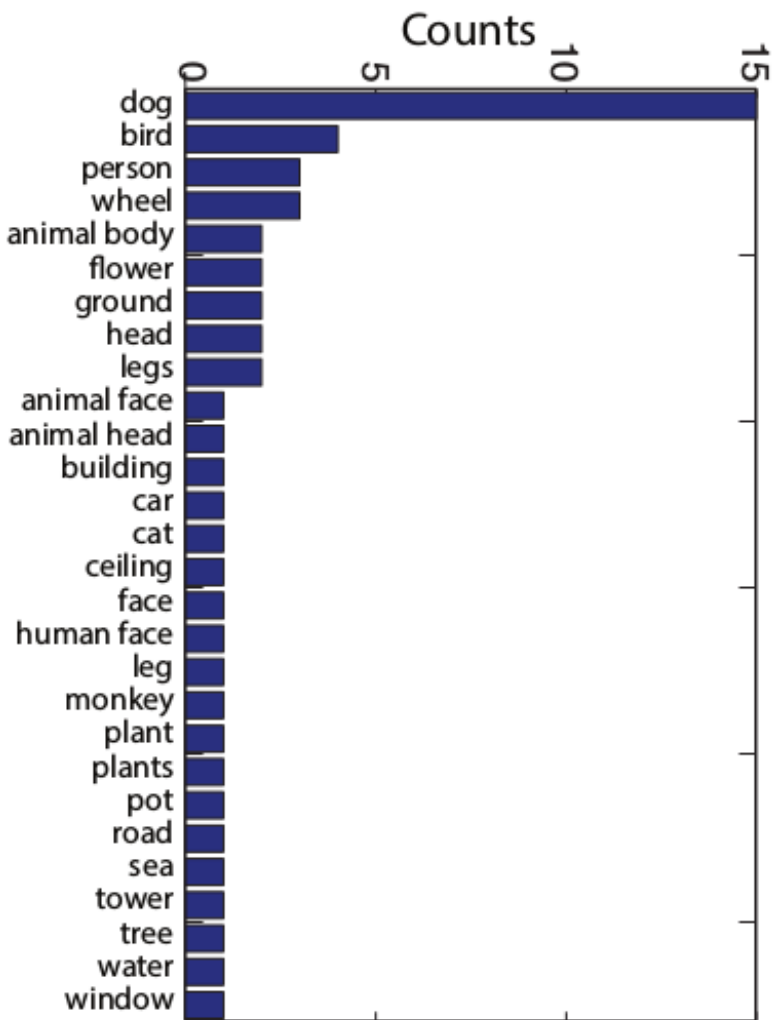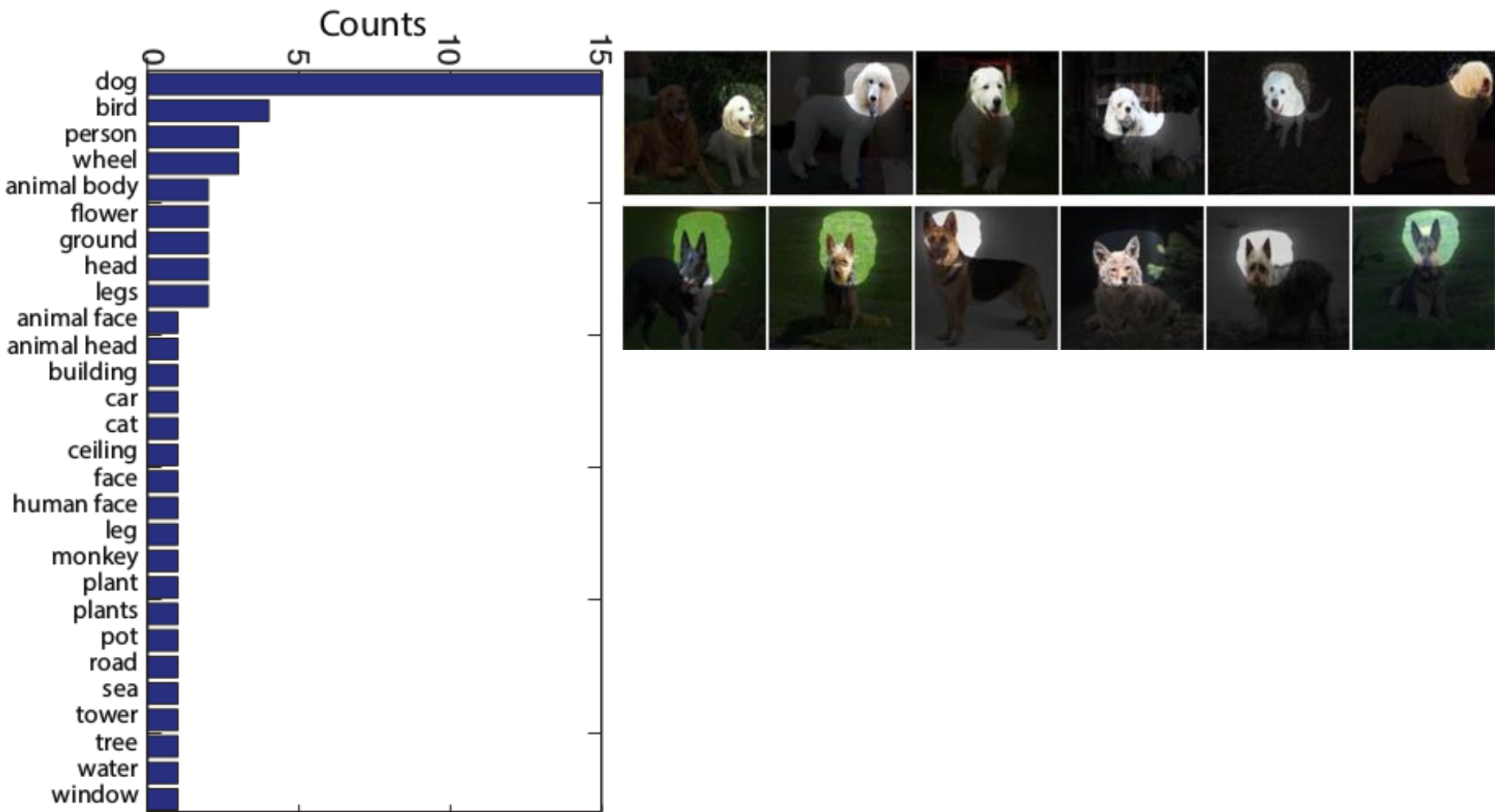# Histogram of Emerged Objects in Pool5

ImageNet-CNN (59/256)



Includes: Objects, nameable parts, and regions

# Histogram of Emerged Objects in Pool5

ImageNet-CNN (59/256)



Includes: Objects, nameable parts, and regions
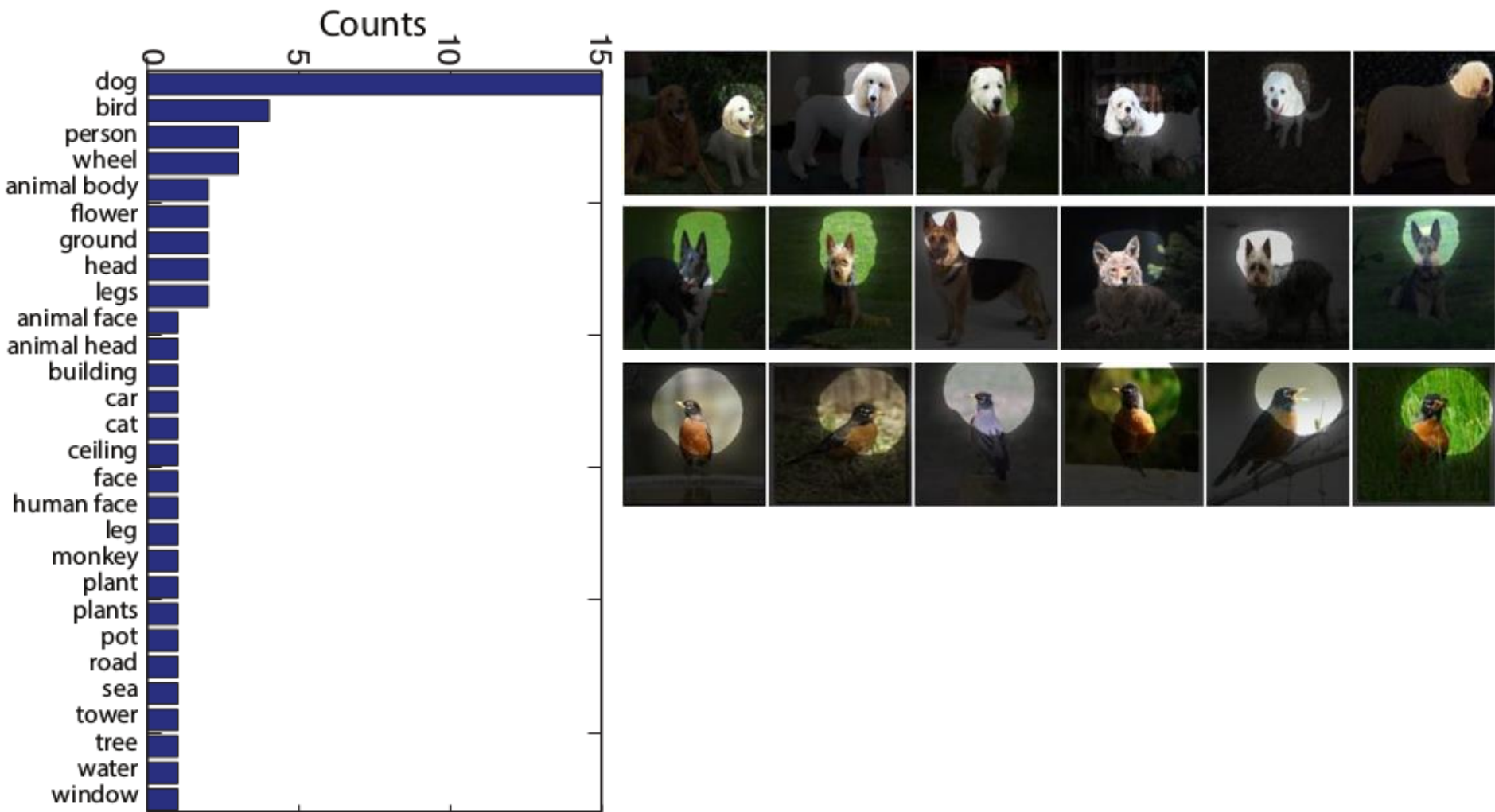
# Histogram of Emerged Objects in Pool5
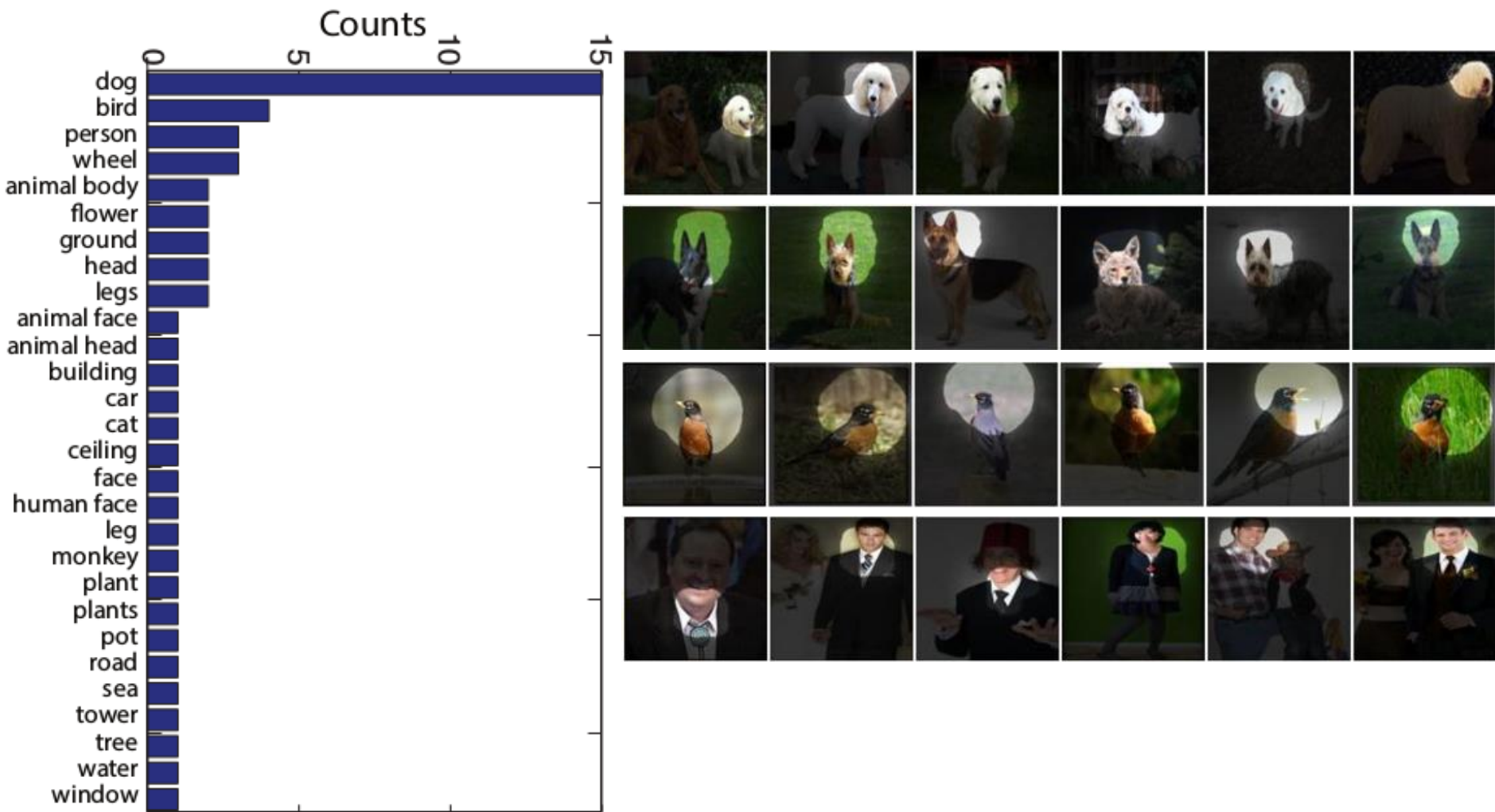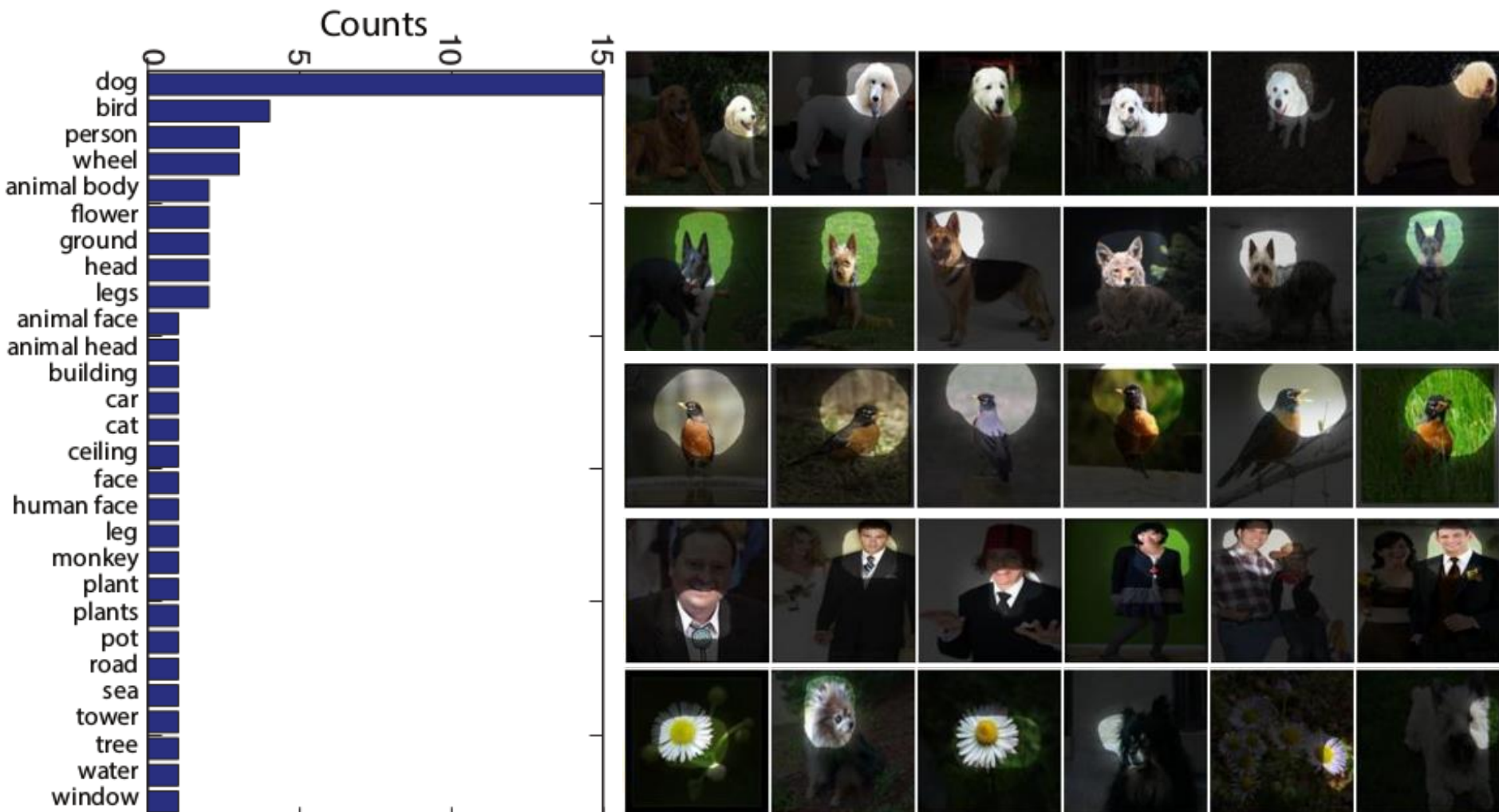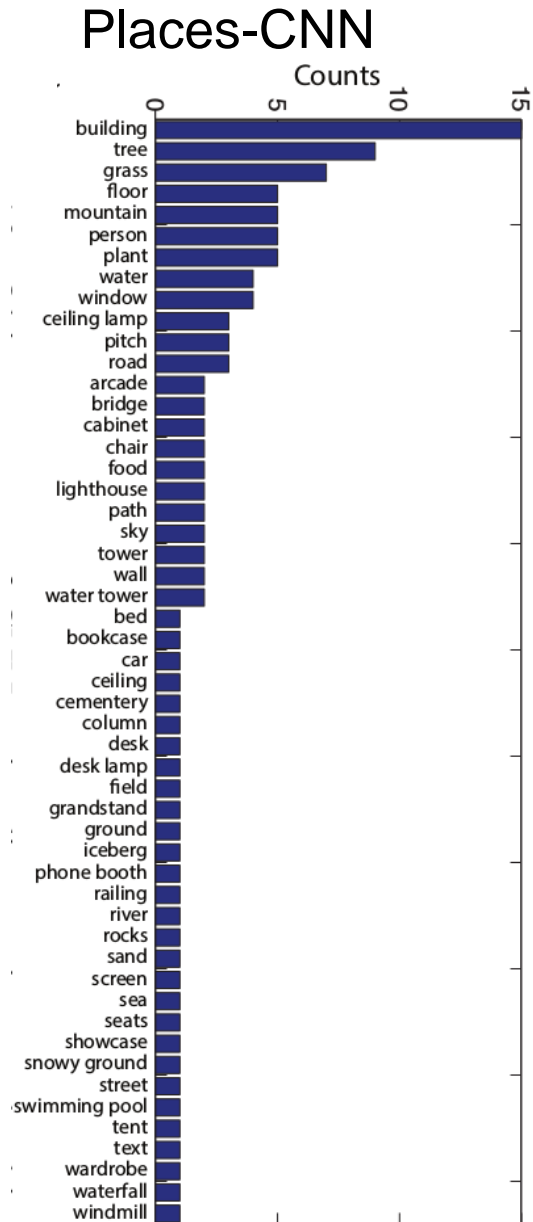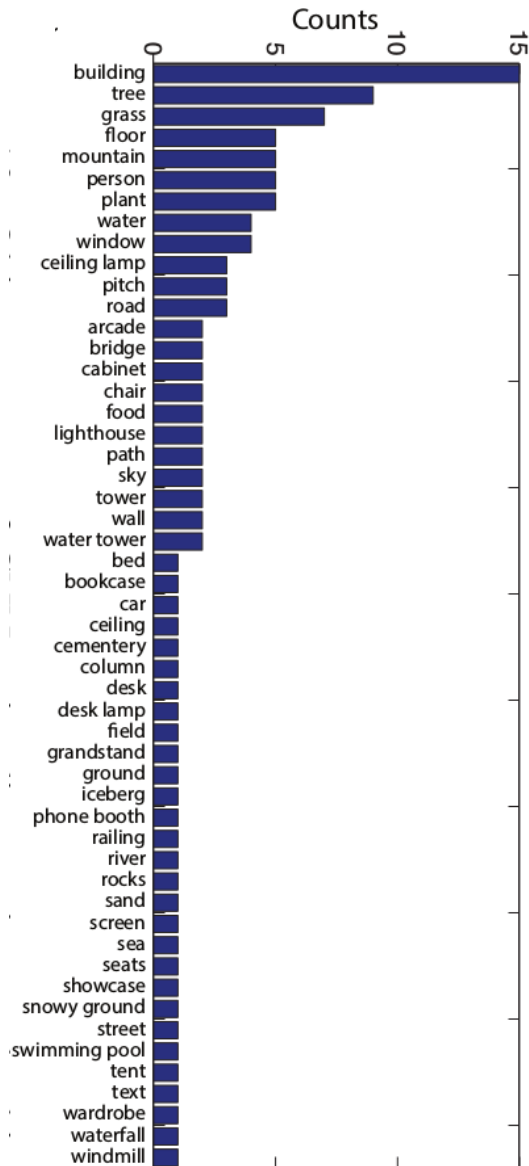
# Histogram of Emerged Objects in Pool5

# Histogram of Emerged Objects in Pool5

Places-CNN

# Histogram of Emerged Objects in Pool5



Places-CNN

# Histogram of Emerged Objects in Pool5

# Histogram of Emerged Objects in Pool5

Places-CNN

# Object detectors emerge inside the CNN

**Buildings**

56) building

120) arcade

8) bridge

123) building

119) building

9) lighthouse

**Scenes**

145) cementery

127) street

218) pitch

**Indoor objects**

182) food

46) painting

106) screen

53) staircase

107) wardrobe

**People**

3) person

49) person

138) person

100) person

**Furniture**

18) billard table

155) bookcase

116) bed

38) cabinet

85) chair

**Lighting**

55) ceiling lamp

174) ceiling lamp

223) ceiling lamp

13) desk lamp

**Outdoor objects**

87) car

61) road

96) swimming pool

28) water tower

6) windmill

**Nature**

195) grass

89) iceberg

140) mountain

159) sand

unitID 106

unitID 107

unitID 108

unitID 109

Nguyen et al, 2016

conv1

conv2

conv3

conv4

conv5

fc6

fc7

Classification layer

bedroom

Mostly task independent

Representation (objects)

Final classification

# Strategies for training for new tasks



conv1　conv2　conv3　conv4　conv5　fc6　fc7　Classification layer

Freeze all these parameters
(trained with ImageNet or Places)

Just train
final classifier

# Strategies for training for new tasks



conv1

conv2

conv3

conv4

conv5

fc6

fc7

Classification layer

Freeze all these parameters
(trained with ImageNet or Places)

Train upper layers to get
a better representation

# But what if you keep the task but change the input modality?



From Devi's webpage: "Abstract images provide several advantages. They allow for the direct study of how to infer high-level semantic information, since they remove the reliance on noisy low-level object, attribute and relation detectors, or the tedious hand-labeling of images."

Bringing Semantics Into Focus Using Visual Abstraction (CVPR), 2013. Zitnick and Parikh.
Learning the Visual Interpretation of Sentences (ICCV), 2013. Zitnick, Parikh, and Vanderwende
Adopting Abstract Images for Semantic Scene Understanding (PAMI), 2015. Zitnick, Vedantam and Parikh

Learning Aligned Cross-Modal Representations from Weakly Aligned Data. Ll. Castrejón*, Y. Aytar*, C. Vondrick, H. Pirsiavash and A. Torralba. CVPR 2016

# From crowdsourcing

## Line drawings

# Line drawings

# Line drawings

# Line drawings

# Line drawings

# Line drawings

# Aquarium

# Library

# Localized words

# Localized words

# or descriptions

There is a bed with a striped bedspread. Beside this is a nightstand with a drawer. There is also a tall dresser and a chair with a blue cushion. On the dresser is a jewelry box and a clock.

# Descriptions

# (Auditorium)

I'm looking forward to seeing this speaker and hearing his story today. I want to get in before all the seats are filled, because he is quite popular with the students and faculty. I don't want to sit way in the back where the sound may not carry as well to.

# Descriptions

## (Classroom)

This room is where students attend and are taught by a teacher on a variety of subjects. Each student seats in a desk which allows him to place books, and write on notebooks or sheets of paper. The teacher presides this room, and usually writes on a blackboard which occupies most of the front wall.

We collected a dataset formed by examples of 205 scene types in five different modalities:

**Line drawings**: 6,644 training – 2,050 validation examples



**Descriptions**: 4,307 training – 2,050 validation examples

There is a bed with a striped bedspread. Beside this is a nightstand with a drawer. There is also a tall dresser and a chair with a blue cushion. On the dresser is a jewelry box and a clock.

I am inside a room surrounded by my favorite things. This room is filled with pillows and a comfortable bed. There are stuffed animals everywhere. I have posters on the walls. My jewelry box is on the dresser.

There are brightly colored wooden tables with little chairs. There is a rug in one corner with ABC blocks on it. There is a bookcase with picture books, a larger teacher's desk and a chalkboard.

**Clipart**: 11,372 training – 1,954 validation examples



**Spatial Text**: 456,300 training – 2,050 validation examples



**Natural images (Places dataset)**: ~ 2M training – 20,500 validation examples

conv1

conv2

conv3

conv4

conv5

fc6

fc7

Classification layer

bedroom

conv1

conv2

conv3

conv4

conv5

fc6

fc7

Classification layer

bedroom

Mostly task independent

conv1

conv2

conv3

conv4

conv5

fc6    fc7

Classification layer

bedroom

Mostly task independent

Task dependent

conv1

conv2

conv3

conv4

conv5

fc6    fc7

Classification
layer

bedroom

Mostly task independent

Task dependent

Modality dependent

conv1

conv2

conv3

conv4

conv5

fc6    fc7

Classification layer

bedroom

Mostly task independent

Modality dependent

Task dependent

Modality independent

conv5    fc6    fc7    Classification layer

bedroom

conv5　fc6　fc7　Classification layer

bedroom

Freeze parameters trained with natural images

conv1

conv2

conv3

conv4

conv5

fc6

fc7

Classification layer

bedroom

Train with new modality

Learning Aligned Cross-Modal Representations from Weakly Aligned Data. Ll. Castrejón*, Y. Aytar*, C. Vondrick, H. Pirsiavash and A. Torralba. CVPR 2016

conv1

conv2

conv3

conv4

conv5

fc6　fc7

Classification layer

bedroom

Freeze parameters trained with natural images

Train with new modality

Learning Aligned Cross-Modal Representations from Weakly Aligned Data. Ll. Castrejón*, Y. Aytar*, C. Vondrick, H. Pirsiavash and A. Torralba. CVPR 2016

conv5

fc6  fc7

Classification
layer

bedroom

conv5

fc6 fc7

Classification layer

bedroom

Freeze parameters trained with natural images

Learning Aligned Cross-Modal Representations from Weakly Aligned Data. Ll.
Castrejón*, Y. Aytar*, C. Vondrick, H. Pirsiavash and A. Torralba. CVPR 2016

conv5

fc6  fc7

Classification layer

bedroom

Freeze parameters trained with natural images

Learning Aligned Cross-Modal Representations from Weakly Aligned Data. Ll. Castrejón*, Y. Aytar*, C. Vondrick, H. Pirsiavash and A. Torralba. CVPR 2016

conv5

fc6

fc7

Classification
layer

bedroom

Freeze parameters
trained with natural
images

Learning Aligned Cross-Modal Representations from Weakly Aligned Data. Ll.
Castrejón*, Y. Aytar*, C. Vondrick, H. Pirsiavash and A. Torralba. CVPR 2016

The room is predominately filled with a large bed and some dressers. There is also a desk with a compute chair and a laptop. On the far wall is a door to the closet.

conv5    fc6    fc7    Classification layer

bedroom

Freeze parameters trained with natural images

Learning Aligned Cross-Modal Representations from Weakly Aligned Data. Ll. Castrejón*, Y. Aytar*, C. Vondrick, H. Pirsiavash and A. Torralba. CVPR 2016

Unit 115
(Bed)

Unit 115
(Bed)

Unit 115
(Bed)

# Unit 115 (Bed)

# Unit 115 (Bed)



ice, terrain, plane, cold, i, nightstand, inside, beds, two, movement

# Units in pool5 become multimodal



Real

Clip art

Unit 31
(Fountain)

Sketches

Spatial text

Descriptions

we, water, fishes, you,
drink, formed, greek,
would, ball, have

# Generating across modalities



A. Dosovitskiy and T. Brox. Inverting convolutional networks with convolutional networks. arXiv, 2015

# Cross-modal learning

## Description (eg, Wikipedia article)

### Snares penguin

From Wikipedia, the free encyclopedia

The **Snares penguin** (*Eudyptes robustus*), also known as the **Snares crested penguin** and the **Snares Islands penguin**, is a penguin from New Zealand. The species breeds on The Snares, a group of islands off the southern coast of the South Island. This is a medium-small, yellow-crested penguin, at a size of 50–70 cm (19.5–27.5 in) and a weight of 2.5–4 kg (5.5–8.8 lb). It has dark blue-black upperparts and white underparts. It has a bright yellow eyebrow-stripe which extends over the eye to form a drooping, bushy crest. It has bare pink skin at the base of its large red-brown bill.

- Lots of descriptions/entries in Wikipedia available

## Images

# Zero-shot Learning

Description (eg, Wikipedia article)

## Cardinal (bird)

From Wikipedia, the free encyclopedia

*This article is about the bird family. For other uses, see Cardinal.*

**Cardinals**, in the family **Cardinalidae**, are passerine birds found in North and South America. They are also known as cardinal-grosbeaks and cardinal-buntings. The South American cardinals in the genus *Paroaria* are placed in another family, the Thraupidae (previously placed in Emberizidae).

Can we predict an image classifier from a description alone?

# Zero-shot Learning

Description (eg, Wikipedia article)

## Cardinal (bird)

From Wikipedia, the free encyclopedia

*This article is about the bird family. For other uses, see Cardinal.*

**Cardinals**, in the family **Cardinalidae**, are passerine birds found in North and South America. They are also known as cardinal-grosbeaks and cardinal-buntings. The South American cardinals in the genus *Paroaria* are placed in another family, the Thraupidae (previously placed in Emberizidae).

Can we predict an image classifier from a description alone?

Assume:

- In training we have access to wiki articles and labeled images
- For test classes we only have wiki articles
- We want to classify a new image (it can belong to any class)

# Zero-shot Learning

- Goal: learn to predict an image classifier from a description
- Linear binary 1-vs-all classifier:

$$y_c = w_c^T \, x$$

- x   …   image feature vector
- w_c   ...   classifier weight vector for class c

# Zero-shot Learning

- Goal: learn to predict an image classifier from a description
- Linear binary 1-vs-all classifier:

$$y_c = w_c^T x$$
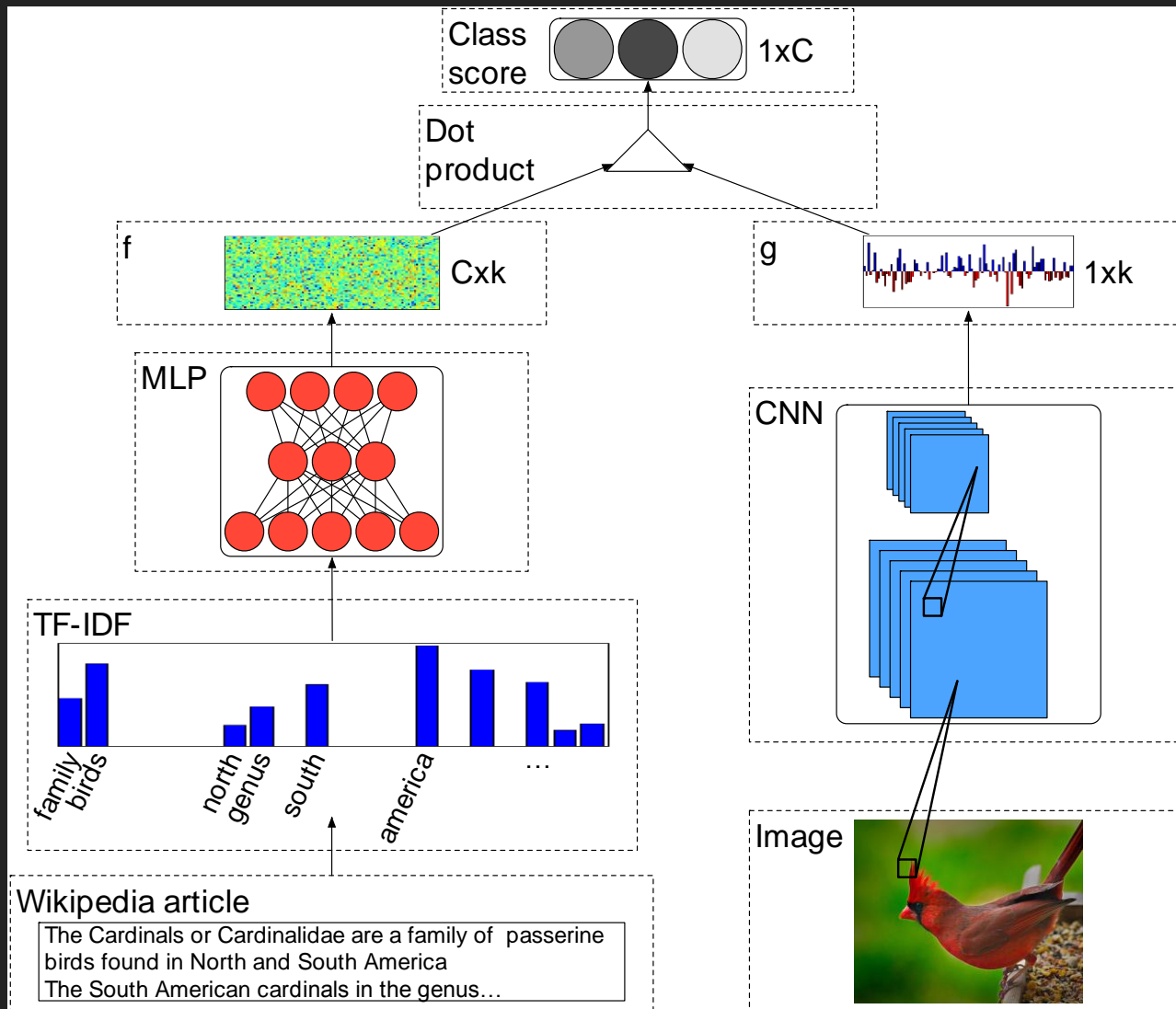
- x … image feature vector
- w_c ... classifier weight vector for class c

- We are also given t_c, a vector representing a textual description about class c
- We want:

$$w_c = f_t(t_c)$$

- f_c … a mapping $\mathbb{R}^p \to \mathbb{R}^d$ that transforms text features to the visual image feature space

# Zero-shot Learning

- f_t can be a neural network



g used to compress x to a k<<d dim

# Red faced Cormorant

The Red-faced Cormorant, Red-faced Shag or Violet Shag, Phalacrocorax urile, is a species of cormorant that is found in the far north of the Pacific Ocean and Bering Sea, from the eastern tip of Hokkaidō in Japan, via the Kuril Islands, the southern tip of the Kamchatka Peninsula and the Aleutian Islands to the Alaska Peninsula and Gulf of Alaska. The Red-faced Cormorant is closely related to the Pelagic Cormorant P. pelagicus, which has a similar range, and like the Pelagic Cormorant is placed by some authors (e.g. Johnsgaard) in a genus Leucocarbo. Where it nests alongside the Pelagic Cormorant, the Red-faced Cormorant generally breeds the more successfully of the two species, and it is currently increasing in numbers, at least in the easterly parts of its range. It is however listed as being of conservation concern{Verify source|date=September 2009}, partly because relatively little is so far known about it.

The adult bird has glossy plumage that is a deep greenish blue in colour, becoming purplish or bronze on the back and sides. In breeding condition it has a double crest,

......

# Red faced Cormorant

The Red-faced Cormorant, Red-faced Shag or Violet Shag, Phalacrocorax urile, is a species of cormorant that is found in the far north of the Pacific Ocean and Bering Sea, from the eastern tip of Hokkaidō in Japan, via the Kuril Islands, the southern tip of the Kamchatka Peninsula and the Aleutian Islands to the Alaska Peninsula and Gulf of Alaska. The Red-faced Cormorant is closely related to the Pelagic Cormorant P. pelagicus, which has a similar range, and like the Pelagic Cormorant is placed by some authors (e.g. Johnsgaard) in a genus Leucocarbo. Where it nests alongside the Pelagic Cormorant, the Red-faced Cormorant generally breeds the more successfully of the two species, and it is currently increasing in numbers, at least in the easterly parts of its range. It is however listed as being of conservation concern{Verify source|date=September 2009}, partly because relatively little is so far known about it.
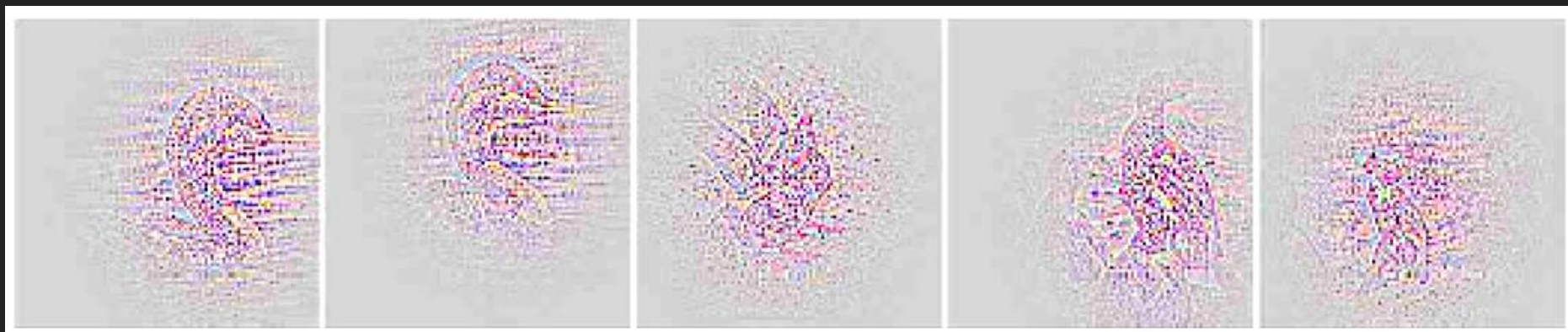
The adult bird has glossy plumage that is a deep greenish blue in colour, becoming purplish or bronze on the back and sides. In breeding condition it has a double crest,
......

# Red faced Cormorant

The Red-faced Cormorant, Red-faced Shag or Violet Shag, Phalacrocorax urile, is a species of cormorant that is found in the far north of the Pacific Ocean and Bering Sea, from the eastern tip of Hokkaidō in Japan, via the Kuril Islands, the southern tip of the Kamchatka Peninsula and the Aleutian Islands to the Alaska Peninsula and Gulf of Alaska. The Red-faced Cormorant is closely related to the Pelagic Cormorant P. pelagicus, which has a similar range, and like the Pelagic Cormorant is placed by some authors (e.g. Johnsgaard) in a genus Leucocarbo. Where it nests alongside the Pelagic Cormorant, the Red-faced Cormorant generally breeds the more successfully of the two species, and it is currently increasing in numbers, at least in the easterly parts of its range. It is however listed as being of conservation concern{Verify source|date=September 2009}, partly because relatively little is so far known about it.

The adult bird has glossy plumage that is a deep greenish blue in colour, becoming purplish or bronze on the back and sides. In breeding condition it has a double crest,
......

→



visualization by Zeiler & Fergus, ECCV'14.

# Learning to see

It is all about the data…

## Strong supervision



## Pixel wise labeling

# Learning to see

It is all about the data…

## Weak supervision



Bird



Bedroom

## Short captions

# Cross modal: text and images



Man holding a metal bowl at the table.

from Microsoft CoCo

Q: Is everyone of these four holding a wine glass? A: No
Q: How many men are there? A: 3
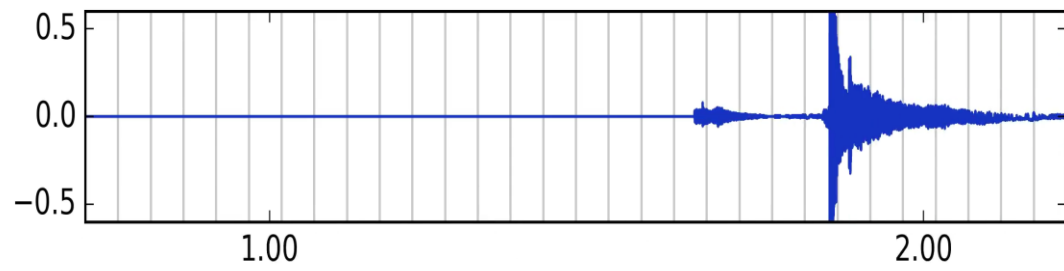Q: Does the window have blinds? A: yes

From http://visualqa.org/index.html

Hard

Soft

Crinkly

# Visually Indicated Sounds



Andrew Owens

Phillip Isola

Josh McDermott

Antonio Torralba

Ted Adelson

Bill Freeman

# Collecting a dataset of physical interactions

# Collecting a dataset of physical interactions



## The Greatest Hits dataset

- 977 videos, 35 sec. long

- 46,577 segmented hits and scratches

- Material, action, reaction labels

The *Greatest Hits* Dataset
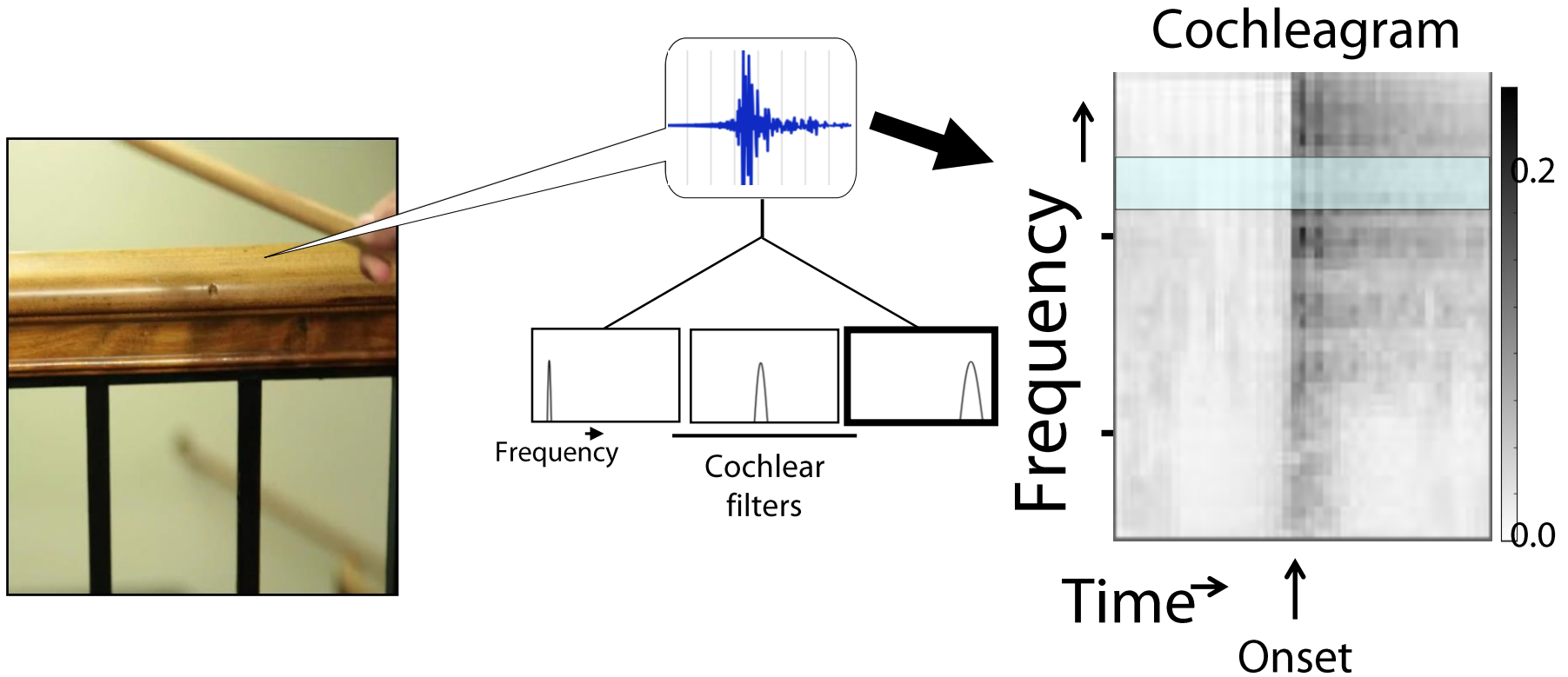
Can we predict material properties from sound?

# Can we predict material properties from sound?

# Can we predict material properties from sound?



Cochleagram

Frequency

Time → Onset

- 40 bandpass filters (+ high/low pass)
- 3 samples per frame (90 Hz)
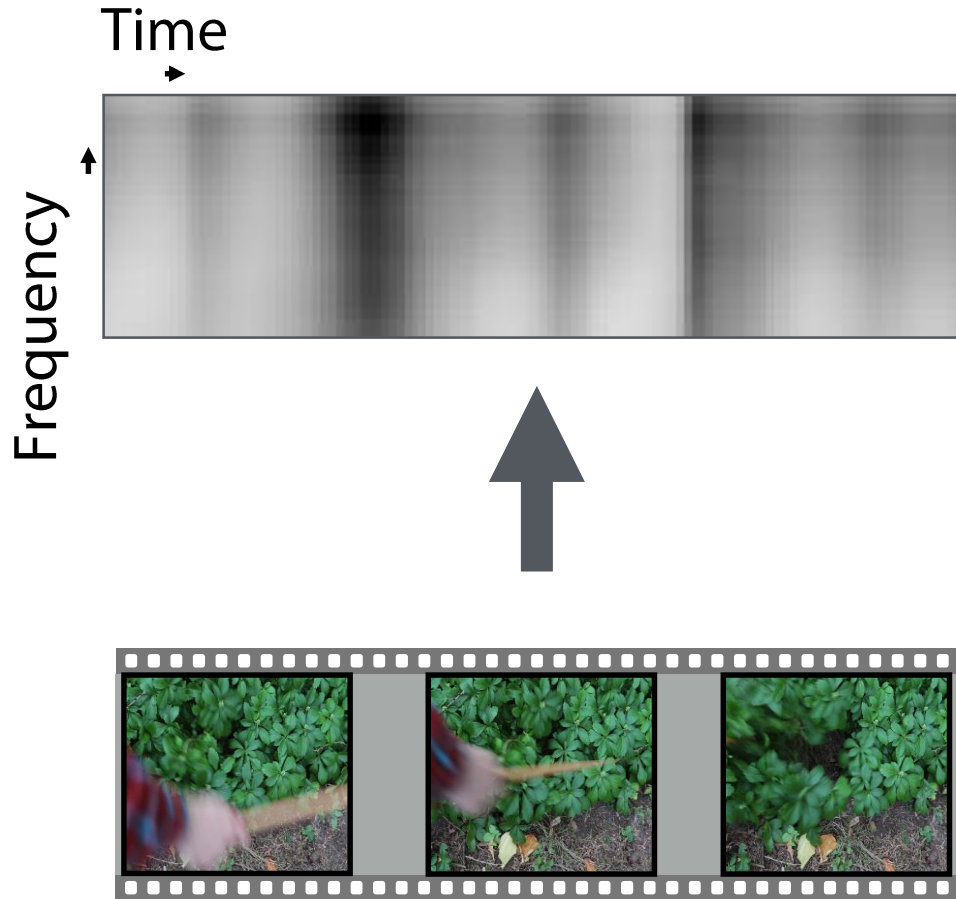
# Can we predict material properties from sound?



Cochleagram

Frequency

Time →

Onset

Frequency

Cochlear filters

- 40 bandpass filters (+ high/low pass)
- 3 samples per frame (90 Hz)

# Can we predict material properties from sound?

## Mean sound features per category

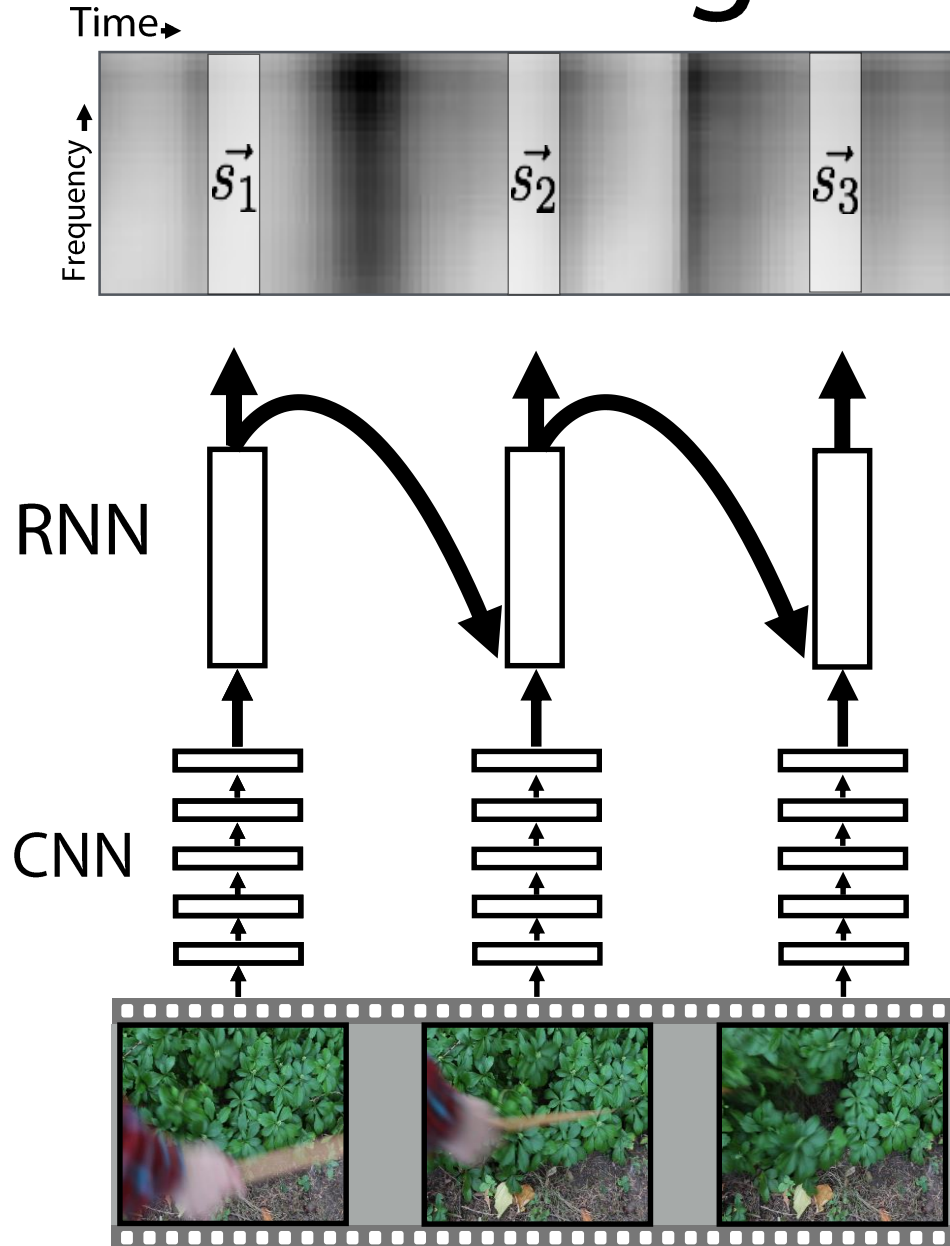# Can we predict material properties from sound?
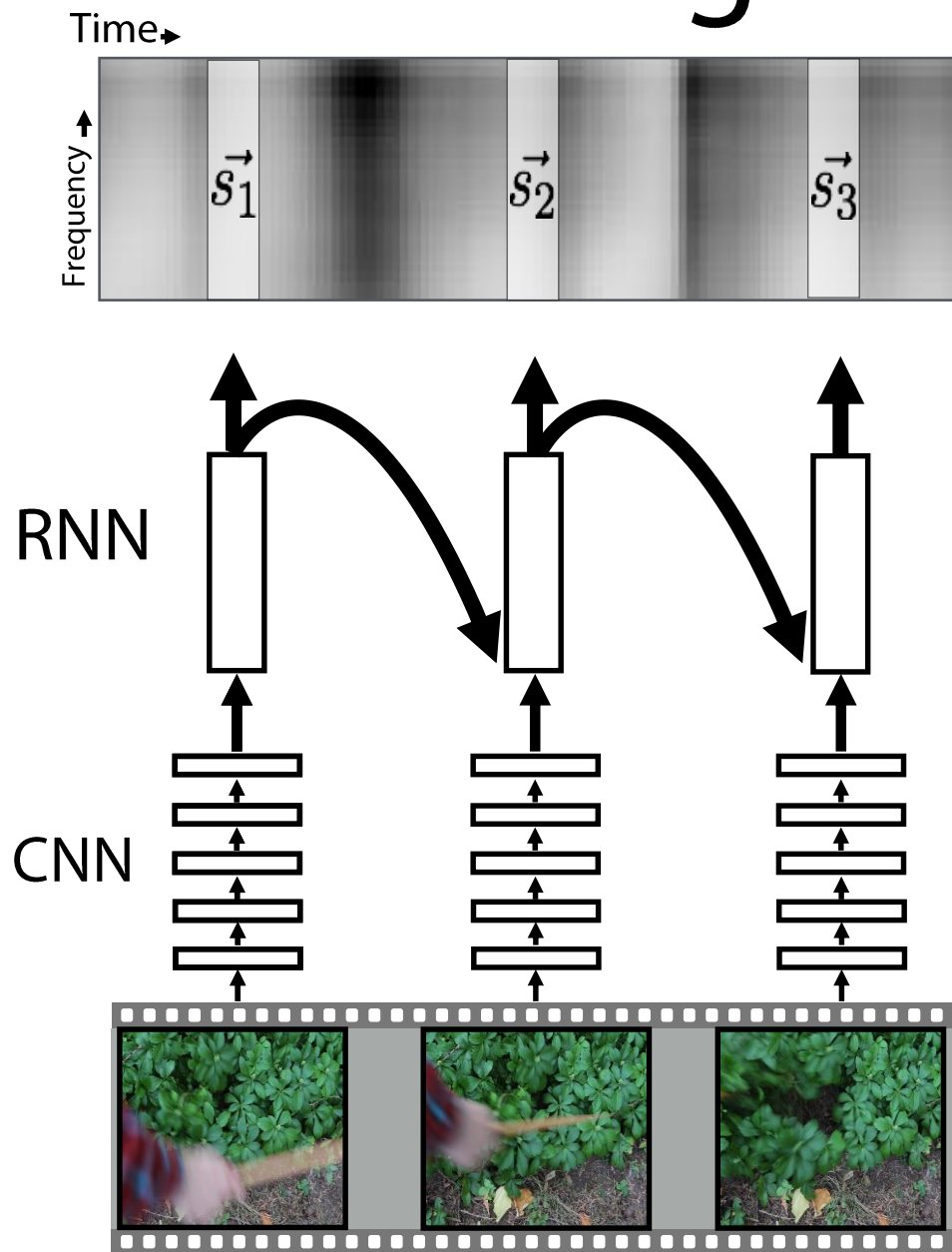
## Mean sound features per category

# Predicting audio features

# Predicting audio features

# Predicting audio features



Time →

Frequency ↑

$\vec{s_1}$  $\vec{s_2}$  $\vec{s_3}$

RNN

CNN

Regression loss

Ground truth
↓

$$\sum_{t=1}^{T} \rho(\|\vec{s}_t - \tilde{\vec{s}}_t\|)$$

where $\rho(r) = \log(\epsilon + r^2)$

- 3D CNN in time domain

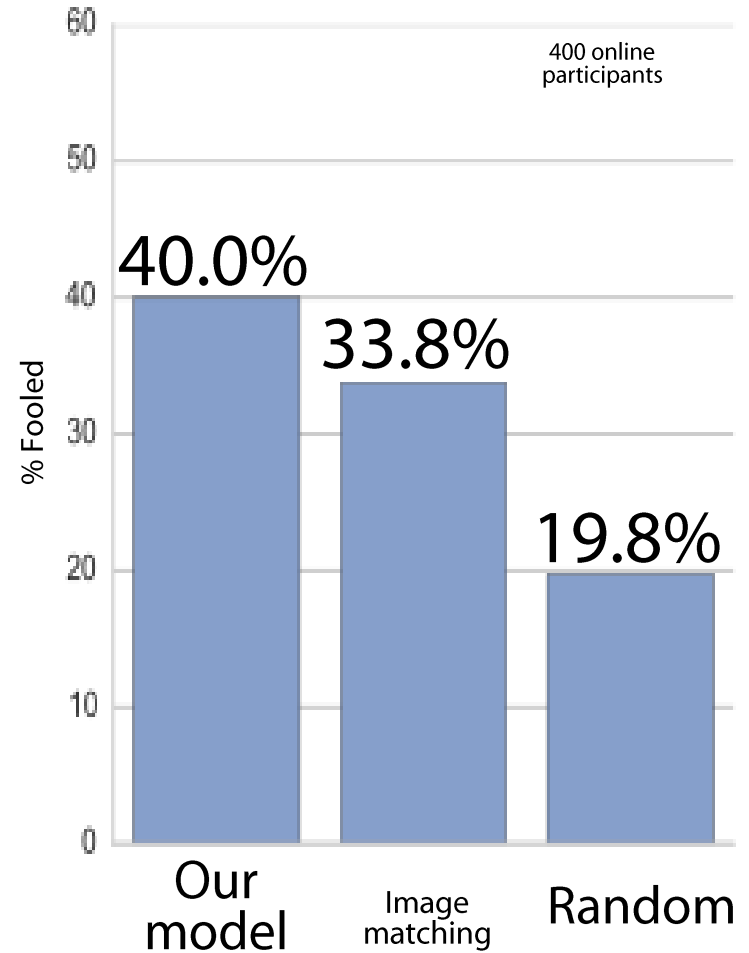- Pretrain from ImageNet

- Long short-term memory

# Real-or-fake study

# Real-or-fake study

# Real-or-fake study



Predicted



Real

# Real-or-fake study


Predicted


Real

Frequency that human participants were fooled.



400 online participants

% Fooled

40.0% — Our model

33.8% — Image matching

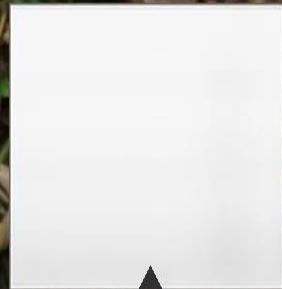19.8% — Random

# Adding soundtracks to silent videos

Predicted
cochleagram

Predicted sound

Predicted sound

Predicted cochleagram

Input video

Transferred
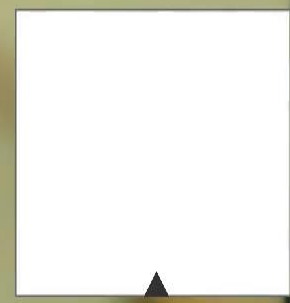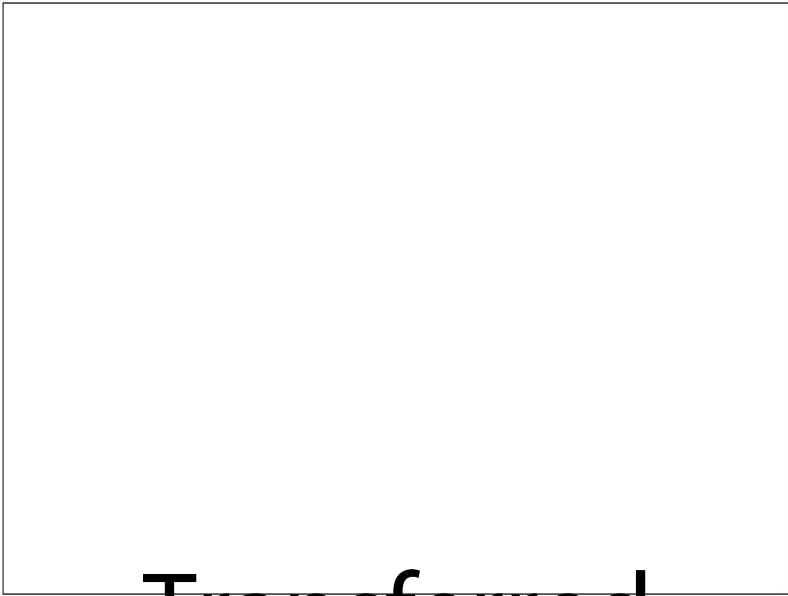audio clips

Predicted sound

Input video

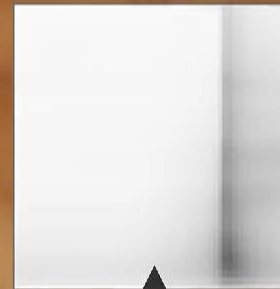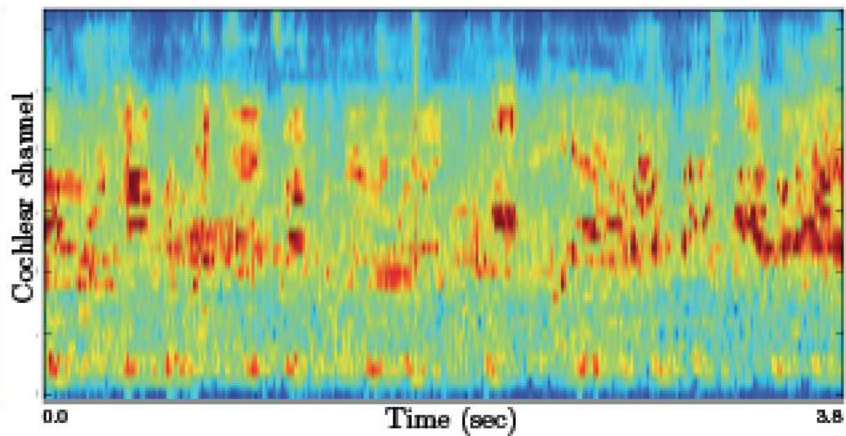Transferred audio clips

Predicted
cochleagra
m

Predicted sound

Input video

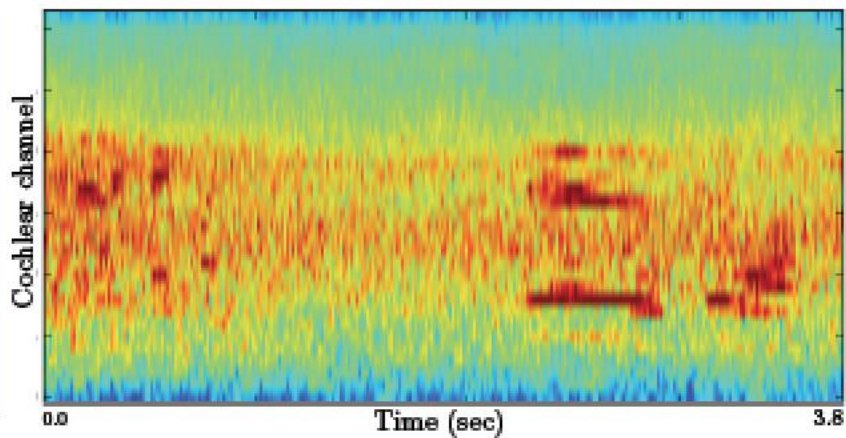Transferred audio clips

Predicted
cochleagra
m

Predicted sound

Predicted cochleagram
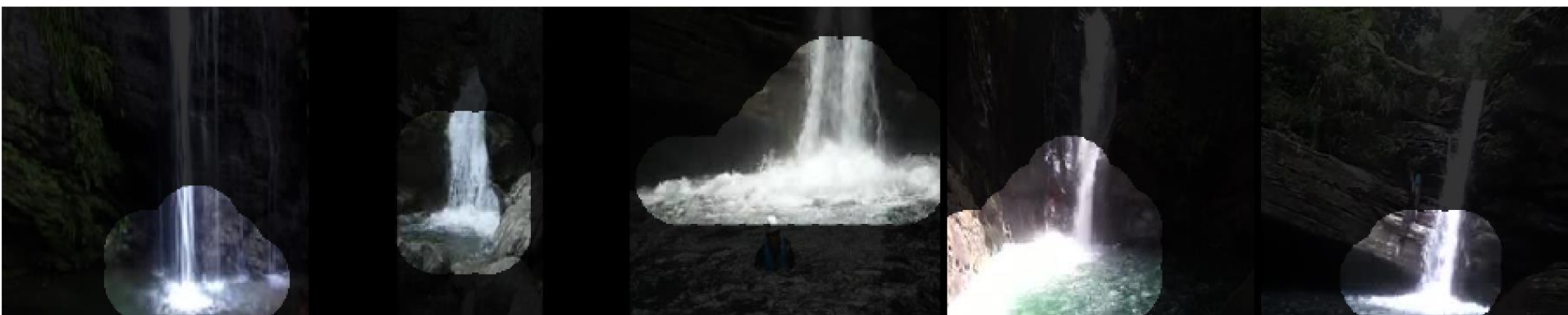
Predicted sound

# Ambient sound



(a) Video frame

(b) Cochleagram

# 99 waterfall

# 99 waterfall



# 194 crowd

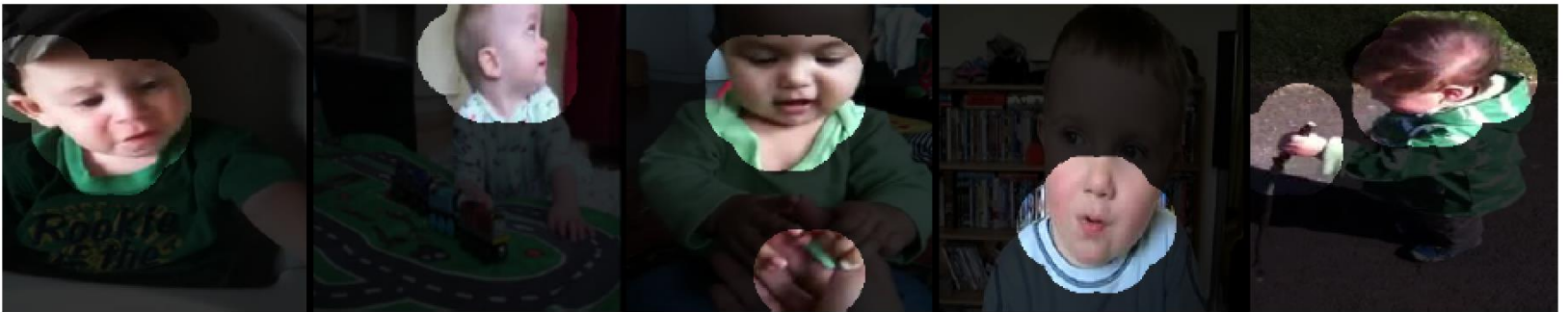# 111 baby

# 111 baby



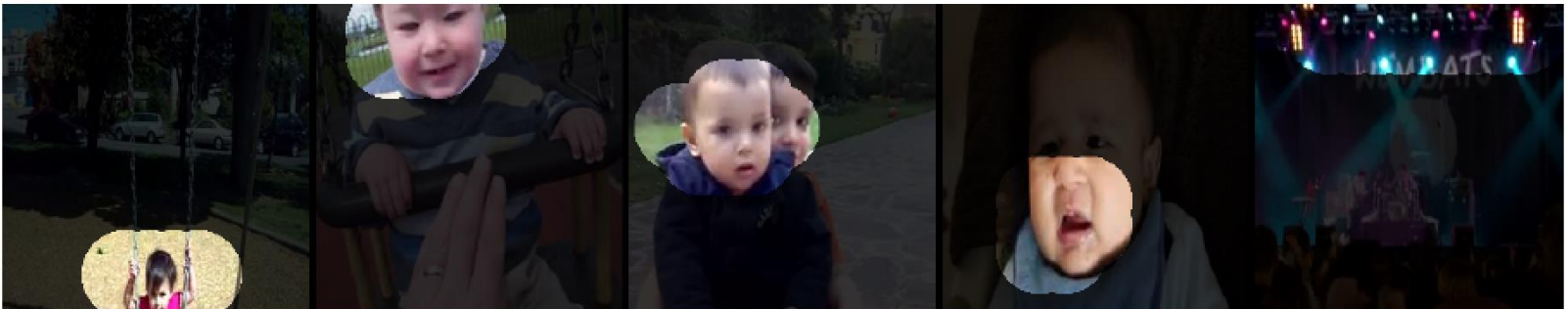# 171 baby

# 111 baby



# 171 baby



# 153 baby

# Neuron visualizations of the network trained by **sound**

## 14 field

## 31 sky
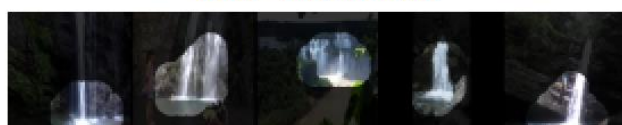
## 67 grass

## 84 snowy ground

## 141 ceiling lamp

## 183 car

## 99 waterfall

## 103 waterfall

## 186 sea

## 111 baby

## 153 baby

## 171 baby

## 15 person

## 20 person
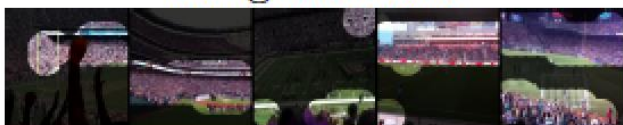
## 37 person

## 194 crowd

## 202 crowd

## 239 crowd

## 150 grandstand

## 163 grandstand

## 218 grandstand

# Learning about the world by hitting things with a drumstick and listening

- Sound is a ubiquitous training signal
- Predicted sounds convey material properties
- Objects make characteristic noises