# **Object Recognition**

Computer Vision Fall 2018 Columbia University

# Project Proposals and Homework 3

• How's it going?

Questions?

# Quick Experiment

Get pen and paper
Draw a coffee cup



# **Two Extremes of Vision**

### Extrapolation problem

Generalization Diagnostic features Interpolation problem Correspondence Finding the differences



# **Tiny Images**



# 80 million tiny images: a large dataset for nonparametric object and scene recognition

Antonio Torralba, Rob Fergus and William T. Freeman. PAMI 2008. http://groups.csail.mit.edu/vision/TinyImages/

### 256x256



## c) Segmentation of 32x32 images



# Given a benchmark, resolution and human scene recognition accuracy increase to a limit



Torralba et al.

## Humans vs. Computers: Car Classification



# Powers of 10

Number of images on my hard drive:

Number of images seen during my first 10 years: (3 images/second \* 60 \* 60 \* 16 \* 365 \* 10 = 630,720,000)

Number of images seen by all humanity: 106,456,367,669 humans<sup>1</sup> \* 60 years \* 3 images/second \* 60 \* 60 \* 16 \* 365 = 1 from http://www.prb.org/Articles/2002/HowManyPeopleHaveEverLivedonEarth.aspx

Number of photons in the universe:

Number of all 32x32 images: 256 32\*32\*3~ 10<sup>7373</sup>

107373

1088

106

108

10<sup>20</sup>

# But not all scenes are so original











## Lots

# Of Images



## Lots

Of Images Target 7,900 790,000

## Lots

Of Images

790,000























79,000,000

# **Application: Automatic Colorization**



Input



**Color Transfer** 



Color Transfer



Matches (gray)



Matches (w/ color)



Avg Color of Match

# **Application: Automatic Colorization**



Input



## **Color Transfer**



## **Color Transfer**



Matches (gray)



Matches (w/ color)



Avg Color of Match

# **Person Recognition**



## .\|

## Exploring the Limits of Weakly Supervised Pretraining Laurens van der Maaten ECCV 2018





Dhruv Mahajan

Ross Girshick Vignesh Ramanathan



anathan Kaiming He



Manohar Paluri



Yixuan Li



Ashwin Bharambe

facebook Artificial Intelligence Research

https://arxiv.org/pdf/1805.00932.pdf

## .\|

#### Research question

# Can we use large amounts of weakly supervised images for pretraining?

#### Highlights

- We pretrain models by predicting relevant hashtags for images
- We pretrain models to predict 17.5K hashtags for 3.5B images
- After finetuning, we beat the state-of-the-art on, e.g., ImageNet

### Hashtag Supervision

- It is easy to get billions of public images and hashtags
- Hashtags are more structured than captions
- Hashtags were often assigned to make images "searchable"



#cheesecake #birthday

facebook Artificial Intelligence Research

### Hashtag Supervision

- But hashtags are not perfect supervision
- . Some hashtags are not visually relevant
- Other hashtags are not in the photo
- And there are many false negatives



#cat #travel #thailand #family #building #fence #...

facebook Artificial Intelligence Research

# 3,500,000,000 images!

### Experiments

- · Select a set of hashtags
- . Download all public Instagram images that has at least one of these hashtags
- · Use WordNet synsets to merge hashtags into canonical form (merge #brownbear and #ursusarctos)
- . The final list has 17,517 hashtags

facebook Artificial Intelligence Research

1	aar	44
2	aardvark	45
3	aardwolf	46
4	aba	47
5	abaca	48
6	abacus	49
7	abalone	50
8	abatis	51
9	abaya	52
10	abbey	53
11	abele	54
12	abelia	55
13	abies	56
14	abila	57
15	abm	58
16	abortus	59
17	abronia	60
18	absinth	61
19	absinthe	62
20	abstraction	63
21	abstractionism	64
22	abutilon	65
23	abutment	66
24	abyss	67
25	abyssinian	68
26	acacia	69
27	acaciadealbata	70
28	academy	71
29	acalypha	72
80	acanthaceae	73
81	acanthurus	74
32	acanthus	75
33	acanthusmollis	76
34	acapulcogold	77
35	acarus	78
86	accelerator	79
37	accelerometer	80
88	access	81
39	accessory	82
10	accident	83
11	accipiter	84
12	accipiternisus	85
13	accipitridae	86

accommodation		17474
accompaniment		17475
accordion		17476
accoutrement		17477
accumulator		17478
ace		17479
aceofclubs		17480
aceofdiamonds		17481
aceofhearts		17482
aceofspades		17483
acer		17484
acerjaponicum		17485
acerola		17486
acerpalmatum		17487
acerrubrum		17488
acetaminophen		17489
acetate		17490
acheron		17491
acherontia		17492
acherontiaatropos		17493
achillea		17494
achilleamillefolium		17495
achimenes	•••	17496
acid		17497
acidophilus		17498
acinonyxjubatus		17499
acinus		1/500
ackee		1/501
aconcagua		1/502
aconite		1/503
aconitum		17504
acorn		17586
acornsquash		17500
acousticguitar		17508
acoustics		17500
acrididae		17510
acrobates		17511
acropolis		17512
acropora		17513
acrylic		17514
acrylicpaints		17515
actias		17516
actiasluna		17517

accommodation

zantac zantedeschia zap zapper zarf zea zebra zebrafinch zebrawood zebu zero zeus zhujiang ziggurat zill zimmerframe zinfandel zina zingiber zinnia zipgun zipper zither ziti ziziphus zizz zodiac zoloft zombi zoologicalgarden zoom zooplankton zootsuit zori zoysia zuiderzee zygnema zygocactus zygoptera

yurt

zabaglione zambeziriver zamboni zamia

### Experiments

- Select a set of hashtags
- Download all public Instagram images that has at least one of these hashtags
- Use WordNet synsets to merge hashtags into canonical form (merge #brownbear and #ursusarctos)
- . Final dataset has ~3.5 **billion** images

**facebook** Artificial Intelligence Research



Q Search

 $\oslash \bigcirc \oslash$ 



#brownbear

164,637 posts

Most Recent

















### Experiments

- Select a set of hashtags
- Download all public Instagram images that has at least one of these hashtags
- Use WordNet synsets to merge hashtags into canonical form (merge #brownbear and #ursusarctos)
- . Final dataset has ~3.5 **billion** images





### Experiments

- Train ResNeXt-32xCd convolutional networks
- Use c-of-*K* vector to represent multiple labels
- . Train to minimize multi-class logistic loss
- Distribute training batches across 336 GPUs
- Scale learning rate by batch size (*N*=8,064) after learning rate "warm-up" (Goyal *et al.*, 2017)



**facebook** Artificial Intelligence Research

- Pretrain model on ImageNet or Instagram
- . Finetune on ImageNet



facebook Artificial Intelligence Research







- Pretrain model on ImageNet or Instagram
- . Finetune on ImageNet
- Similar results on larger versions of ImageNet



**facebook** Artificial Intelligence Research

# **Two Extremes of Vision**

### Extrapolation problem

Generalization Diagnostic features Interpolation problem Correspondence Finding the differences



# Exemplar-SVMs



- Learn a separate linear SVM for each instance (exemplar) in the dataset (PASCALVOC)
- Each Exemplar-SVM is trained with a **single** positive instance
- Each Exemplar-SVM is more defined by "what it is not" vs. "what it is similar to"

# Large-scale training

Ex

Exemplar-SVM 2

Each exemplar performs its CPU own hard negative mining

Exemplar-SVM 1

- Solve many convex learning problems
- Parallel training on cluster



 $\mathbf{E}\mathbf{X}_2$ 

Exemplar-SVM N

EXN

CPU

# Exemplar



## Appearance





## Exemplar



## Appearance







## Exemplar



## Appearance


















3D Model













Appearance











#### Appearance



5

i

:.







## What's this?



Photo from Coffee Creek Watershed Preserve

## What's this?



Entry-level categories (Jolicoeur, Gluck, Kosslyn 1984)

- Typical member of a basic-level category are categorized at the expected level
- Atypical members tend to be classified at a subordinate level.



A bird



An ostrich

# **Classical Categorization**

- Group objects by common properties
- What are birds?
  - animals, has wings, has feathers, can fly, chirps





# **Prototype Theory**

Rosch and Lakoff

- According to the prototype view, an object will be classified as an instance of a category if it is sufficiently similar to the prototype.
- Evidence for Prototype:
- **Typicality ratings**: how good are robins as an example of birds
- Production order of exemplars: Name all the kinds of bird you can think of
- Time to verify categorical statements: True or false: a robin is a bird



**Figure 7.3.** Schematic of the prototype model. Although many exemplars are seen, only the prototype is stored. The prototype is updated continually to incorporate more experience with new exemplars.

### **Canonical Perspective**

The "best," most easily identified view of an object. (Palmer, Rosch & Chase, 1981)













Slide by Palmer

# **Dyirbal Indigenous People**



# The perception of function Direct perception (affordances): Gibson



Mediated perception (Categorization)



### **Direct perception**

Some aspects of an object function can be perceived directly

 Functional form: Some forms clearly indicate to a function ("sittable-upon", container, cutting device, ...)



### **Direct perception**

- Some aspects of an object function can be perceived directly
- Observer relativity: Function is observer dependent



#### Limitations of Direct Perception

Objects of similar structure might have very different functions



**Figure 9.1.2** Objects with similar structure but uncreated and tions. Mailboxes afford letter mailing, whereas trash cans do not, even though they have many similar physical features, such as size, location, and presence of an opening large enough to insert letters and medium-sized packages.



Not all functions seem to be available from direct visual information only.

The functions are the same at some level of description: we can put things inside in both and somebody will come later to empty them. However, we are not expected to put inside the same kinds of things...

#### Segmentation: Where *really* are the people?















## Solution 1: Image pyramids



Learning Hierarchical Features for Scene Labeling. Clement Farabet, Camille Couprie, Laurent Najman, Yann LeCun. In *TPAMI*, 2013.

Slide credit: Bharath Hariharan

## Solution 2: Skip connections



Slide credit: Bharath Hariharan

## Skip connections



Fully convolutional networks for semantic segmentation. Evan Shelhamer, Jon Long, Trevor Darrell. In CVPR 2015 Slide credit: Bharath Hariharan

## Skip connections

• Problem: early layers not semantic



Visualizations from : M. Zeiler and R. Fergus. Visualizing and Understanding Convolutional Networks. In ECCV 2014.

Slide credit: Bharath Hariharan

## Solution 3: Dilation

- Need subsampling to allow convolutional layers to capture large regions with small filters
  - Can we do this without subsampling?



## Solution 3: Dilation

- Need subsampling to allow convolutional layers to capture large regions with small filters
  - Can we do this without subsampling?



## Solution 3: Dilation

- Need subsampling to allow convolutional layers to capture large regions with small filters
  - Can we do this without subsampling?



### Solution 4: Conditional Random Fields

- Idea: take convolutional network prediction and sharpen using classic techniques
- Conditional Random Field

$$\mathbf{y}^* = \arg\min_{\mathbf{y}} \sum_{(i,j)} E_{data}(y(i,j)) + \sum_{(i,j),(k,l) \in \mathcal{N}} E_{smooth}(y(i,j), y(k,l))$$

 $E_{smooth}(y(i,j), y(k,l)) = \mathbb{I}(y(i,j) \neq y(k,l))w(i,j,k,l)$ 

## Fully Connected CRFs

- Typically, only adjacent pixels connected
  - Fewer connections => Easier to optimize
- Dense connectivity: every pixel connected to everything else
- Intractable to optimize except if pairwise potential takes specific form

$$E_{smooth}(y(i,j),y(k,l)) = \mathbb{I}(y(i,j) \neq y(k,l))w(i,j,k,l)$$
$$w(i,j,k,l) = \sum_{m} w_{m}e^{-\|\mathbf{f}_{m}(i,j) - \mathbf{f}_{m}(k,l)\|^{2}}$$

Efficient Inference in Fully Connected CRFs with Gaussian Edge Potentials. Philipp Krahenbuhl, Vladlen Koltun. In NIPS, 2011. Slide credit: Bharath Hariharan



http://www.visualqa.org/challenge.html



What color are her eyes? What is the mustache made of?



Is this person expecting company? What is just under the tree?



How many slices of pizza are there? Is this a vegetarian pizza?



Does it appear to be rainy? Does this person have 20/20 vision?

Fig. 1: Examples of free-form, open-ended questions collected for images via Amazon Mechanical Turk. Note that commonsense knowledge is needed along with a visual understanding of the scene to answer many questions.

#### Questions and answers collected with Amazon Mechanical Turk



Is something under the sink broken?	yes yes yes	no no no
What number do you see?	33 33 33	5 6 7



	yes
	no
Is this man crying?	no
	no

no yes yes



Can you park here?	no no no	no no yes
What color is the hydrant?	white and orange white and orange white and orange	red red yellow



Has the pizza been baked?	yes yes yes	yes yes yes
What kind of cheese is topped on this pizza?	feta feta ricotta	mozzarella mozzarella mozzarella



What kind of store is this?	bakery bakery pastry	art supplies grocery grocery
Is the display case as full as it could be?	no no no	no yes yes



How many pickles are on the plate?	1 1 1	1 1 1
What is the shape of the plate?	circle round round	circle round round

Fig. 2: Examples of questions (black), (a subset of the) answers given when looking at the image (green), and answers given when not looking at the image (blue) for numerous representative examples of the dataset. See the appendix for more examples.
### Architecture



### Words

• Need ways to compare words

Next to the 'sofa' is a desk, and a 'person' is sitting behind it. 'armchair' 'man' 'bench' 'woman' 'chair' 'child' 'teenager' 'deck chair' 'ottoman' 'girl' 'seat' 'boy' 'stool' 'baby' 'daughter' 'swivel chair' 'son' 'loveseat' • • • . . .



### I parked the car in a nearby street. It is a red car with two doors, ...

# I parked the vehicle in a nearby street...

T. Mikolov, K. Chen, G. Corrado, J. Dean. Efficient Estimation of Word Representations in Vector Space. arXiv:1301.3781, 2013

### word2vec

### I parked the car in a nearby street. It is a red car with two doors, ...



T. Mikolov, K. Chen, G. Corrado, J. Dean. Efficient Estimation of Word Representations in Vector Space. arXiv:1301.3781, 2013

### word2vec



T. Mikolov, K. Chen, G. Corrado, J. Dean. Efficient Estimation of Word Representations in Vector Space. arXiv:1301.3781, 2013

# Algebraic operations with the vector representation of words

X = Vector("Paris") - vector("France") + vector("Italy")

Closest nearest neighbor to X is vector("Rome")

### Architecture



### Architecture



There are 1000 possible answers in this system. Questions are unlimited.



What red objects in front are almost covered by snow?	meters parking meters parking meters	car cars shoes
Is it winter?	yes yes yes	no yes yes
	200	
Is this photo taken	10	

in Antarctica?	no no	yes yes
Overcast or sunny?	overcast overcast overcast	overcast overcast sunny



_	Does the car have a license plate?	yes yes yes	yes yes yes
	Could the truck have a camper?	yes yes yes	yes yes yes



Is the picture hanging<br/>straight?no<br/>yes<br/>yes<br/>yesno<br/>yes<br/>yesHow many cabinets are<br/>on the piece of<br/>furniture?436



Is the woman on the back of the bicycle pedaling?	no no yes	no no yes
Why is the woman holding an umbrella?	sunny to block sun uncertain	it's raining it's raining to stay dry

What type of trees are here?	palm asr palm oak palm pine		
Is the skateboard airborne?	yes yes yes	no yes yes	

Fig. 27: Random examples of questions (black), (a subset of the) answers given when looking at the image (green), and answers given when not looking at the image (blue) for numerous representative examples of the real image dataset.



#### what is on the ground?

#### Submit



#### what is on the ground?

#### Submit

#### Predicted top-5 answers with confidence:

#### sand

0.528%

90.748% SNOW 2.858% beach 1 418% surfboards 0.677% water



#### what color is the umbrella?

#### Submit



#### what color is the umbrella?

#### Submit

#### Predicted top-5 answers with confidence:

#### yellow

	95.090%	
white		
<mark>1.8</mark> 11%		
black		
0.663%		
blue		
0.541%		
gray		

0.362%



#### are we alone in the universe?

#### Submit



#### are we alone in the universe?

#### Submit

#### Predicted top-5 answers with confidence:

#### no

	78.234%	
ves		
J = =		
21.763%		
noonlo		
heohie		
0.001%		
hirde		
DIIUS		
0.000%		
a t		
OUL		
0.000%		



what is the meaning of life?

#### Submit



what is the meaning of life?

#### Submit

#### Predicted top-5 answers with confidence:

#### beach

15.262%			
sand			
8 537%			
0.007 /0			
cooquill			
seayuii			
4.708 <mark>%</mark>			
tower			
0 202%			
2.393 %			
rooko			
TUCKS			
1.746%			



#### what is the yellow thing?

#### Submit

Predicted top-5 answers with confidence:

#### frisbee

1.252%

79.844%

surfboard		
7.319%		
banana		
<mark>2.8</mark> 44%		
lemon		
<mark>2.4</mark> 38%		
surfboards		



how many trains are in the picture?

#### Submit

Predicted top-5 answers with confidence:



### iBOWIMG



### iBOWIMG

Table 2: Performance comparison on test-standard.

	Open-Ended			Multiple-Choice				
	Overall	yes/no	number	others	Overall	yes/no	number	others
LSTMIMG [2]	54.06	-	-	-	_	-	-	-
NMN+LSTM [1]	55.10	-	-	-	-	-	-	-
ACK [16]	55.98	79.05	36.10	40.61	-	-	-	-
DPPnet [11]	57.36	80.28	36.92	42.24	62.69	80.35	38.79	52.79
iBOWIMG	55.89	76.76	34.98	42.62	61.97	76.86	37.30	54.60

# **Two Extremes of Vision**

#### Extrapolation problem

Generalization Diagnostic features Interpolation problem Correspondence Finding the differences

