

From Visual Recognition To Visual Understanding Laurens van der Maaten

Yapay Öğrenme Yaz Okulu 2020 Bilkent Üniversitesi

.\|

Mission

"Advancing the state-of-the-art in AI through open research for the benefit of all."

Values

- Openness: Publish and open-source
- Freedom: Researchers have complete control on their agenda
- <u>Collaboration</u>: With internal and external partners
- Excellence: Focus on most impactful projects
- Scale: Operate at large scale

portal







Translation

50%

of users have at least one friend with a different native language



translated posts every day on Facebook



From Visual Recognition...

Convolutional Networks

- Repeatedly filters image
- Filters are learned based on labeled training images
- Final representation is "semantic"













Pretraining Vision Models

 First, train a model on a large "source" dataset (say, ImageNet)

Pretraining Vision Models

 First, train a model on a large "source" dataset (say, ImageNet)

- Finetune on a small "target" dataset
- Measure accuracy on target task

facebook Artificial Intelligence Research

12

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

facebook Artificial Intelligence Research

Dhruv Mahajan

Ross Girshick

Vignesh Ramanathan

Manohar Paluri

Yixuan Li

Ashwin Bharambe

Kaiming He

- 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

• But hashtags are not perfect supervision

#cat #travel #thailand #family

- But hashtags are not perfect supervision
- Some hashtags are not visually relevant

#cat #travel #thailand #family

- But hashtags are not perfect supervision
- Some hashtags are not visually relevant
- Other hashtags are not in the photo

#cat #travel #thailand #family

- 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 #...

- 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
- Is scaling up sufficient to make up for this noise?

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

- 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)

- Select a set of hashtags
- Download all public Instagram images that has at least one of these hashtags

aar

- Use WordNet synsets to merge hashtags into canonical form (merge #brownbear and #ursusarctos)
- The final list has 17,517 hashtags

aar	44	accommodation	17474	yurt
aardvark	45	accompaniment	17475	zabaglione
aardwolf	46	accordion	17476	zambeziriver
aba	47	accoutrement	17477	zamboni
abaca	48	accumulator	17478	zamia
abacus	49	ace	17479	zantac
abalone	50	aceofclubs	17480	zantedeschia
abatis	51	aceofdiamonds	17481	zap
abaya	52	aceofhearts	17482	zapper
abbey	53	aceofspades	17483	zarf
abele	54	acer	17484	zea
abelia	55	acerjaponicum	17485	zebra
abies	56	acerola	17486	zebrafinch
abila	57	acerpalmatum	17487	zebrawood
abm	58	acerrubrum	17488	zebu
abortus	59	acetaminophen	17489	zero
abronia	60	acetate	17490	zeus
absinth	61	acheron	17491	zhujiang
absinthe	62	acherontia	17492	ziggurat
abstraction	63	acherontiaatropos	17493	zill
abstractionism	64	achillea	17494	zimmerframe
abutilon	65	achilleamillefolium	17495	zinfandel
abutment	66	achimenes	17496	zing
abyss	67	acid	17497	zingiber
abyssinian	68	acidophilus	17498	zinnia
acacia	69	acinonyxjubatus	17499	zipgun
acaciadealbata	70	acinus	17500	zipper
academy	71	ackee	17501	zither
acalvoha	72	aconcagua	1/502	2111
acanthaceae	73	aconite	1/503	zizipnus
acanthurus	74	aconitum	17504	2122 Todios
acanthus	75	acorn	1/505	zoloft
acanthusmollis	76	acornsquash	17500	zorbi
acapulcogold	77	acousticquitar	17509	zoologicalgarden
acarus	78	acoustics	17500	Zoom
accelerator	79	acrididae	17510	zoonlankton
accelerometer	80	acrobates	17510	zootsuit
access	81	acropolis	17512	zori
accessory	82	acropora	17513	zovsia
accident	83	acrylic	17514	zuiderzee
accipiter	84	acrylicpaints	17515	zvonema
accipiternisus	85	actias	17516	zygocactus
accipitridae	86	actiasluna	17517	zygoptera

- 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)

I cannot show you the images, but you can look...

https://www.instagram.com/explore/tags/brownbear/ https://www.instagram.com/explore/tags/crane/ https://www.instagram.com/explore/tags/...

• The final image set has ~3.5B images

- Train ResNeXt-32xCd convolutional networks
- Use *c*-of-*K* vector to represent multiple labels
- Train to minimize multi-class logistic loss

- Train ResNeXt-32xCd convolutional networks
- Use *c*-of-*K* vector to represent multiple labels
- Train to minimize multi-class logistic loss

most experiments use ResNeXt-101 32x16d

- Pretrain model on ImageNet or Instagram
- Finetune on ImageNet

facebook Artificial Intelligence Research

Artificial Intelligence Research

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

facebook Artificial Intelligence Research

Fix Data; Vary Model

- Increasing model capacity has a larger positive effect
- Even lower error rates may be possible?

Fix Data; Vary Model

- Increasing model capacity has a larger positive effect
- Even lower error rates may be possible?

Learning Curves

- Accuracy on target task improves (almost) loglinearly with data size
- Matching hashtags to target task may help

facebook Artificial Intelligence Research

.\|

Research question

Does image recognition work for everyone?

Highlights

- We test cloud services for image recognition on household items
- Images used for test come from across the world
- Services work better in some countries than others

Terrance DeVries Changhan Wang Ishan Misra

facebook Artificial Intelligence Research

Image Classification

• Is image classification solved? Yes?

Soap

Country of Origin: UK Prediction: Toiletry

Spices

Toothpaste

Country of Origin: USA Prediction: Toothpaste

6:23

Image Classification

 Is image classification solved? Yes? No! Soap

Country of Origin: UK Prediction: Toiletry Spices

Country of Origin: USA Prediction: Spice Toothpaste

Country of Origin: USA Prediction: Toothpaste

Country of Origin: Philippines Prediction: Beer

facebook Artificial Intelligence Research

Country of Origin: Nepal Prediction: Food

Country of Origin: Burundi Prediction: Wood

35

Dollar Street

- Photo collection gathered by Gapminder to show how people across the world live
- Photos are annotated with object class, country, and family income
- We used ~20,000 photos from 117 classes

https://www.gapminder.org/dollar-street
Results

- Average top-5 accuracy of all image recognition systems:
- Amazon Rekognition
- Google Cloud Vision
- Clarifai
- Microsoft Azure
- IBM Watson
- Accuracy varies per country
- Results are consistent across all services analyzed



Red: 60%; Yellow: 75%; Green: 90%

Results

- Average top-5 accuracy of all image recognition systems:
- Amazon Rekognition
- Google Cloud Vision
- Clarifai

facebook

Artificial Intelligence Research

- Microsoft Azure
- IBM Watson
- Accuracy varies per country
- Results are consistent across all services analyzed





Results

 Analysis of only India shows this is not only due to income correlating with location



Main problems

 Dataset collection relies on services that are primarily popular in the West

Main problems

- Dataset collection relies on services that are primarily popular in the West
- Most dataset collections start with English queries

Top results for "Wedding"



Top results for "शादी" (Wedding in Hindi)



Main problems

- Dataset collection relies on services that are primarily popular in the West
- Most dataset collections start with English queries

Top results for "Wedding"



Top results for "शादी" (Wedding in Hindi)



Top results for "Spices"



Top results for "मसाले" (Spices in Hindi)









facebook Artificial Intelligence Research

... to Visual Understanding

.\|

Research question

How do we measure image understanding?

Highlights

• We find problems in common benchmarks for visual understanding





Allan Jabri



Armand Joulin

Visual Question Answering

• How do we measure understanding?

Visual Question Answering

- How do we measure understanding?
- <u>Proposal</u>: Given an image, answer questions about that image



What color is the jacket? -Red and blue. -Yellow. -Black. -Orange.



How many cars are parked? -Four. -Three. -Five. -Six.



What event is this?

- -A wedding. -Graduation. -A funeral.
- -A picnic.



When is this scene taking place? -Day time.

46

- -Night time.
- -Evening.
- -Morning.

Models for VQA

facebook

Artificial Intelligence Research

 Feed image through convolutional network







Models for VQA

- Feed image through convolutional network
- Feed question through recurrent network





Models for VQA

- Feed image through
 convolutional network
- Feed question through recurrent network
- Apply multi-class logistic regressor to predict answer











Visual Question Answering

- Train on a collection of 70K multiple-choice questions
- Measure accuracy on 40K held-out questions

	Percentage of Correct Answers
Encode question with LSTM Encode image with conv. network	52.1%



Are we building horses?

- Take multiple-choice visual question
- Throw away both image and question
- Encode answer using word2vec features
- Train binary classifier to predict whether or not answer is correct
- Test time: Predict highest-scoring answer

Visual Question Answering

• Our simple baseline outperforms the state-of-the-art

	Percentage of Correct Answers		
Encode question with LSTM Encode image with conv. network	52.1%		
Simple model	52.9%		

Visual Question Answering

 Does looking at the image and question help? Sure.

	Percentage of Correct Answers
Encode question with LSTM Encode image with conv. network	52.1%
Simple model	52.9%
Simple model++	67.1%

CLEVR

• Visual reasoning benchmark that cannot be solved by a "horse"

facebook Artificial Intelligence Research



Q: Are there an equal number of large things and metal spheres? Q: What size is the cylinder that is left of the brown metal thing that is left of the big sphere? Q: There is a sphere with the same size as the metal cube; is it made of the same material as the small red sphere? Q: How many objects are either small cylinders or metal things?

Research question

How do we develop better benchmarks for image captioning?

Highlights

- Image captioning systems are difficult to evaluate
- We propose binary image selection as an alternative evaluation







Ishan Misra

facebook Artificial Intelligence Research

 \mathbf{X}

- Evaluating relevance of captions to images is very difficult
- Does a captioning system really possess visual understanding?

- Evaluating relevance of captions to images is very difficult
- Does a captioning system really possess visual understanding?



- Evaluating relevance of captions to images is very difficult
- Does a captioning system really possess visual understanding?



A bunch of luggage sitting on top of a floor.

- Evaluating relevance of captions to images is very difficult
- Does a captioning system really possess visual understanding?



A bunch of luggage sitting on top of a floor. CIDEr-D: 55.96

- Evaluating relevance of captions to images is very difficult
- Does a captioning system really possess visual understanding?



A pile of garbage sitting next to a trash can.

- Evaluating relevance of captions to images is very difficult
- Does a captioning system really possess visual understanding?



A pile of garbage sitting next to a trash can.

CIDEr-D: 0.14

• Captioning scores correlate poorly with human evaluations of correctness:

facebook Artificial Intelligence Research

* Annotations gathered using COCO guidelines for human evaluation.

Image Retrieval

- Is text-based image retrieval a good alternative?
- No: does not provide true negatives.

Text query: A person wearing a banana headdress and necklace.

Correct image

Text query: There is a green clock in the street.

facebook Artificial Intelligence Research

Retrieved image

Correct image

BISON: Binary Image SelectiON

- Given a text query, pick one of two images
- BISON is designed such that both images are similar, but one is a negative for the query

Text query: Plates filled with carrots and beets on a white table.

Negative image

Positive image

Text query: Yellow shirted tennis player looking for incoming ball.

Positive image

Negative image

Analyzing Systems

 Performance of captioning and retrieval systems on COCO-val

	BLEU-4	CIDEr-D	Recall@1	BISON
Show & Tell (Vinyals <i>et al.</i> , 2015)	32.35	97.20	-	78.59
Show & Tell + Attention (Xu et al., 2015)	33.49	101.55	-	82.04
UpDown (Anderson <i>et al.</i> , 2018)	34.53	105.40	-	84.04
Convnet + BoW	-	-	45.19	81.48
Convnet + BiGRU (Faghri <i>et al.</i> , 2018)	-	-	49.34	85.46
SCAN i2t (Lee <i>et al.</i> , 2018)	-	-	52.35	86.40
SCAN t2i (Lee <i>et al.</i> , 2018)	-	-	54.10	87.50
Human	21.70*	85.40*	-	100.0

COCO-validation

* Human scores computed on COCO test set.

COCO-BISON

facebook Artificial Intelligence Research

Try it yourself at http://github.com/facebookresearch/binary-image-selection 66

.\|

Research question

How do we move beyond image understanding to physical reasoning?

Highlights

• PHYRE is a new benchmark for physical reasoning

Anton Bakhtin Laura Gustafson

Gustafson Ross Girshick

Justin Johnson

• New benchmark for test ability of AI agent to perform physical reasoning:

• Overview of the stages in the PHYRE benchmark:

Make the green ball touch the purple jar by adding a red ball

• Solution strategies the agent needs to use are very diverse:

Make the green ball touch the blue/purple object by adding red objects

facebook Artificial Intelligence Research

http://www.phyre.ai

Tasks cannot be solved well by random search

PHYRE

• Agents are still far away from doing optimal ranking (which is non-optimal itself):



.\|

Conclusion

Visual recognition works but we still have a way to go towards visual understanding.

Details

- Large-scale training of modern convolutional networks works great
- Care is needed to prevent networks from having undesired biases
- Visual understanding is still difficult and work-in-progress

Thank you!