IM & GENET

Where have we been? Where are we going?

LI FEI-FEI







The Beginning: CVPR 2009



Jia Deng, Wei Dong, Richard Scol Dens, of Communer Science	her, Li-Jia Li, Kai Li and Li Fei-Fei Princeton University, USA	1
Itagent, white, reader, Tigl,	il. Buitellibboe.grteenten.edt	200
		- 25
		- 261
Abstract	content-based image scoreh and image sails estimating al go-	1000
an an an an 1976 Sprace and a state of the	stitutes, as well as for previding critical standing and beach-	1.5
The englishment of througe iteas on the Enformer has the pro-	Instanting auto per source argorithma.	- 254
rithms to index, arthreet, organize and internet with pre-	Each meaningful concept in WordPlot, resultin described	6 m
and multimedie data. But enough how such data care	by multiple words or word pleasas, is called a "symmym-	
humenoid and organized remains a critical problem. We	set" or "sympet". There are around \$0,000 nears sympto-	22
induce here a new darphase cathof "IntegeNet", a horge-	to WordNet. In ItrageNet, we also to provide on avail-	1.1.1
the statisticity of statutes had apon the tardition of the	age 505-1000 images to illustrate such synol. Istages of	1.00
the hit 000 ensure of BirrdNet with an exercise of 580-	as described in Sec. 3.2. ImageNet, therefore, will offer	1
Wetcom and full revolution images. This well would in	arm of sollions of algority sound images. In this paper,	
s of settleme of annotated analyse organized by the se-	we report the carriert version of ImageNet, consisting of 12	
atic hiselatchy of WordNee. This paper affere a detailed	"sabbrers": enumeral. Rind, Sub, reprint, aerphabian, soldele	
error of magnetic to to correct state. 12 subserve with P consists and 3.7 million instance in total. We show that	Source that. These solutions contain 5547 counts and 3.7	25.
upoNee to much larger its scale and allversity and much	million images. Fig. 1 shows a snapshet of non-branches of	100
re accursice than the current intege distorers. Construct-	the manerial and valuely subtracy. The database is publicly	
such a targe-scale doubters is a challenging spil, 101	available # http://www.bhapp-net.com	
orthe alter altatu collectives scheme with Assactor Mechani-	The rest of the paper is segurized as follows: My first	
I Tork. Larry, we sharman the anglebars of hengefier	show that therefore is a targe-wate, such the out driverse	100
Characterization and automatic object characters. The holes	continuentication manufactory by continuing the partner fact-	81
t the scale, screency, directity and Nersechical sense-	geNet, mostly the manufail and vehicle subtrant. Our goal	100
of InsegrNet can offer unparalleled opportanistics to re-	in to show that dwargeNet start serve are a use/fal renovance for	100
othere in the computer vision community and beyond.	stout incognition applications such as object recognition,	ske os
	anger characteristication and object investigation. In addition, the	an of
to the distance of the second s	you to be an edy on traditional data collection wethods.	Page 1
Introduction	Sec. 3 describes how fenageNet is unsurrated by leverag-	
The digital are has breaght with it an energies-	ing Amazon Mechanical Terk.	in the
n of data. The latest estimations put a member of every	3. Properties of ImageNet	mie.
a 2 betten planes on Hacke, a similar number of valen	a such such a such set	go at
County france are an even safgre turness for happenet	ImageNet is built upon the hierarchical sincture pre-	esc.
and invadely and algorithms can be proposed by exploit-	contain to the order of 50 million cleanly laboled full second	dise
these images, resulting in homor applications for users	hation assages (500-100) pay surrouts. At the time this paper	
ndex, retrieve, regarize and interact with these data. But	is written, ImageNet consists of 12 subtrees. Most analysis	8 40
city how such data can be eithinid and organized in a	will be based on the transmal and vahicle subtrace.	+ the
an database called "bearsNet", a large scale provings	Single Insurables are served to the most commentance	100
tragges. We believe that a large scale conslogy of imager	and diserve coverage of the image work). The current 17	100
critical resource for developing advanced, keep-scale	referencement of a trul of 3.2 million clearly annuald	2.De
		-fy.
	12	
		in.
of margin. Two only a solitation first, margins are publicly another	jets in mugas should have suitable appearances p	siters.
1111/06/04/17/07/10/07/11/01/01/01/01/01/01/01/01/01/01/01/01/		
		-
and the Wayness' Residences in from and 5 Kings of	Van and included Van Mill denote the antipher Law other Sector	-

J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li and L. Fei-Fei, ImageNet: A Large-Scale Hierarchical Image Database. IEEE Computer Vision and Pattern Recognition (CVPR), 2009.

The Impact of IM GENET

IM GENET on Google Scholar

4,386 Citations

Imagenet: A large-scale hierarchical image database J Deng, W Dong, R Socher, LJ Li, K Li... - Computer Vision and ..., 2009 - ieeexplore.ieee.org Abstract: The explosion of image data on the Internet has the potential to foster more sophisticated and robust models and algorithms to index, retrieve, organize and interact with images and multimedia data. But exactly how such data can be harnessed and organized Cited by 4386 Related articles All 30 versions Cite Save

2,847 Citations Imagenet large scale visual recognition challenge O Russakovsky, J Deng, H Su, J Krause... - International Journal of ..., 2015 - Springer Abstract The ImageNet Large Scale Visual Recognition Challenge is a benchmark in object category classification and detection on hundreds of object categories and millions of images. The challenge has been run annually from 2010 to present, attracting participation Cited by 2847 Related articles All 17 versions Cite Save

...and many more.

From IMAGENET Challenge Contestants to Startups



clarifai







DNNresearch VUNO

A Revolution in Deep Learning

FORTUNE

Why Deen Laarning is Suddanly

The New York Times

By Roger Parloff

Changin;

The Great Artificial Intelligence Awakeni QUARTZ

By Gideon Lewis-K

The data that transformed AI research—and possibly the world

By Dave Gershgorn

"The IM GENET of x"



SpaceNet DigitalGlobe, CosmiQ Works, NVIDIA



MusicNet J. Thickstun et al, 2017



ShapeNet A.Chang et al, 2015



EventNet G. Ye et al, 2015



Medical ImageNet Stanford Radiology, 2017



ActivityNet F. Heilbron et al, 2015

An Explosion of Datasets





"Datasets—not algorithms—might be the key limiting factor to development of human-level artificial intelligence."

ALEXANDER WISSNER-GROSS Edge.org, 2016

The Untold History of IM GENET

Hardly the First Image Dataset



Segmentation (2001) D. Martin, C. Fowlkes, D. Tal, J. Malik.



KTH human action (2004) I. Leptev & B. Caputo



CAVIAR Tracking (2005) R. Fisher, J. Santos-Victor J. Crowley



CMU/VASC Faces (1998) H. Rowley, S. Baluja, T. Kanade



Sign Language (2008) P. Buehler, M. Everingham, A. Zisserman



Middlebury Stereo (2002) D. Scharstein R. Szeliski



MSRC (2006) Shotton et al. 2006



PASCAL (2007) Everingham et al, 2009



FERET Faces (1998) P. Phillips, H. Wechsler, J. Huang, P. Raus



UIUC Cars (2004) S. Agarwal, A. Awan, D. Roth



CalTech 101/256 (2005) Fei-Fei et al, 2004 GriffIn et al. 2007



Lotus Hill (2007)Yao et al, 2007



COIL Objects (1996) S. Nene, S. Nayar, H. Murase



3D Textures (2005) S. Lazebnik, C. Schmid, J. Ponce

14
10
1

LabelMe (2005) Russell et al, 2005





MNIST digits (1998-10) Y LeCun & C. Cortes



CuRRET Textures (1999) K. Dana B. Van Ginneken S. Nayar J. Koenderink

205 The Str came		1-	Dog
The second second		and the second second	Leash
	=	Bar March	German
and a	Ξ.	10 T 11	Shepard
The survey of th	1.2	Can 1 (Ca) .	Standing
		- Lima	Canine

ESP (2006) Ahn et al. 2006



TinyImage (2008) Torralba et al. 2008

A Profound Machine Learning Problem Within Visual Learning



Machine Learning 101: Complexity, Generalization, Overfitting



One-Shot Learning



Fei-Fei et al, 2003, 2004



Fei-Fei et al, 2003, 2004

How Children Learn to See



A new way of thinking...

To shift the focus of Machine Learning for visual recognition

from modeling...

...to data. Lots of data.

Internet Data Growth 1990-2010



Source: Cisco

What is WordNet?



Original paper by [George Miller, et al 1990] cited over 5,000 times Organizes over 150,000 words into 117,000 categories called *synsets*. Establishes ontological and lexical relationships in NLP and related tasks.

Christiane Fellbaum

Senior Research Scholar Computer Science Department, Princeton President, Global WordNet Consortium

Individually Illustrated WordNet Nodes



jacket: a short coat



German shepherd: breed of large shepherd dogs used in police work and as a guide for the blind.



microwave: kitchen appliance that cooks food by passing an electromagnetic wave through it.



mountain: a land mass that projects well above its surroundings; higher than a hill.



A massive ontology of images to transform computer vision



IM GENET Comrades



Prof. Kai Li Princeton



Jia Deng 1st Ph.D. student Princeton



Step 1: Ontological structure based on WordNet





Dog

German Shepherd **Step 2:** Populate categories with thousands of images from the Internet





Dog

Shepherd

hand

Three Attempts at Launching IMAGENET

1st Attempt: The Psychophysics Experiment



1st Attempt: The Psychophysics Experiment

- # of synsets: **40,000** (subject to: imageability analysis)
- # of candidate images to label per synset: **10,000**
- *#* of people needed to verify: **2-5**
- Speed of human labeling: 2 images/sec (one fixation: ~200msec)
- Massive parallelism (N ~ 10^2-3)

40,000 × 10,000 × 3 / 2 = 6000,000,000 sec

≈ 19 years N

2nd Attempt: Human-in-the-Loop Solutions

Towards scalable dataset construction: An active learning approach

Brendan Collins, Jia Deng, Kai {bmcollin, dengjia, li, feifei

Department of Computer Science, Princeton

Abstract. As computer vision research co and greater variation within object categor more exhaustive datasets are necessary. He ing such datasets is laborious and monoto in which many images have been automa category (typically by automatic internet s relevant images from noise. We present a d which employs active, online learning to with minimal user input. The principle adv vious endeavors is its scalability. We demon superior to the state-of-the-art, with scala work.

1 Introduction

Though it is difficult to foresee the future of co that its trajectory will include examining a gr (such as objects or scenes), that the complexit

OPTIMOL: automatic Online Picture collecTion via Incremental MOdel Learning

Li-Jia Li1, Gang Wang1 and Li Fei-Fei2

¹ Dept. of Electrical and Computer Engineering, University of Illinois Urbana-Champaign, USA ² Dept. of Computer Science, Princeton University, USA jail3@uiuc.edu, gwang6@uiuc.edu, feifeili@cs.princeton.edu

Abstract

A well-built dataset is a necessary starting point for advanced computer vision research. It plays a crucial role in evaluation and provides a continuous challenge to stateof-the-art algorithms. Dataset collection is, however, a tedious and time-consuming task. This paper presents a novel automatic dataset collecting and model learning approach



2nd Attempt: Human-in-the-Loop Solutions



Machine-generated datasets can only match the best algorithms of the time.



Human-generated datasets transcend algorithmic limitations, leading to better machine perception.



The Result: IM GENET Goes Live in 2009





What We Did Right

While Others Targeted Detail...



LabelMe

Per-Object Regions and Labels Russell et al, 2005



Lotus Hill

Hand-Traced Parse Trees Yao et al, 2007
...We Targeted Scale

SUN, 131K [Xiao et al. '10]

LabelMe, 37K

[Russell et al. '07]

PASCAL VOC, 30K

[Everingham et al. '06-'12]

Caltech101, 9K

[Fei-Fei, Fergus, Perona, '03]



Additional IM GENET Goals



Carnivore

- Canine
 - Dog
 - Working Dog
 - Husky



High Resolution

To better replicate human visual acuity

High-Quality Annotation

To create a benchmarking dataset and advance the state of machine perception, not merely reflect it

Free of Charge

To ensure immediate application and a sense of community

An Emphasis on Community and Achievement

IM GENET

Large Scale Visual Recognition Challenge (ILSVRC 2010-2017)

ILSVRC Contributors



Alex Berg UNC Chapel Hill



Jia Deng Univ. of Michigan



Zhiheng Huang Stanford



Aditya Khosla Stanford





Fei-Fei Li Stanford



Wei Liu **UNC Chapel Hill**



Sean Ma Stanford



Eunbyung Park UNC Chapel Hill



Olga Russakovsky Sanjeev Satheesh Stanford Stanford





Hao Su Stanford

Jonathan Krause Stanford

Our Inspiration: PASCAL VOC



2005-2012

Our Inspiration: PASCAL VOC

Mark Everingham 1973-2012



Mark Everingham Prize @ ECCV 2016

IM GENET

Alex Berg, Jia Deng, Fei-Fei Li, Wei Liu, Olga Russakovsky

Participation and Performance



Entries

Participation and Performance



Participation and Performance



What we did to make IMAGENET better

Lack of Details



Lack of Details...ILSVRC Detection Challenge





Evaluation of ILSVRC Detection

Need to annotate the presence of all classes (to penalize false detections)

Table	Chair	Horse	Dog	Cat	Bird
+	+	-	-	-	-
+	-	-	-	+	-
+	+	-	-	-	-

images: 400K
classes: 200
annotations = 80M!

Evaluation of ILSVRC Detection

Hierarchical annotation



J. Deng, O. Russakovsky, J. Krause, M. Bernstein, A. Berg, & L. Fei-Fei. CHI, 2014

What does classifying 10K+ classes tell us?



Fine-Grained Recognition





Fine-Grained Recognition



[Gebru, Krause, Deng, Fei-Fei, CHI 2017]





2567 classes 700k images

Expected Outcomes



ImageNet becomes a benchmark



Breakthroughs in object recognition



Machine learning advances and changes dramatically

Unexpected Outcomes

Neural Nets are Cool Again!



13,259 Citations Imagenet classification with deep convolutional neural networks A Krizhevsky, I Sutskever, GE Hinton - Advances in neural ..., 2012 - papers.nips.cc Abstract We trained a large, deep convolutional neural network to classify the 1.3 million high-resolution images in the LSVRC-2010 ImageNet training set into the 1000 different classes. On the test data, we achieved top-1 and top-5 error rates of 39.7\% and 18.9\% Cited by 13259 Related articles All 95 versions Cite Save

Krizhevsky, Sutskever & Hinton, NIPS 2012

...And Cooler and Cooler ③

"AlexNet"



0	

"GoogLeNet"

Vet	9.
	Vet







[Krizhevsky et al. NIPS 2012]

[Szegedy et al. CVPR 2015]

[Simonyan & Zisserman, ICLR 2015]

[He et al. CVPR 2016]



Ontological Structure Structure Not Used as Much





Deng, Krause, Berg & Fei-Fei, CVPR 2012



Deng, Krause, Berg & Fei-Fei, CVPR 2012



Optimizing with a Knowledge Ontology Results in Big Gains in Information at Arbitrary Accuracy



Deng, Krause, Berg & Fei-Fei, CVPR 2012

Relatively Few Works Have Used Ontology



Kuettel, Guillaumin, Ferrari. Segmentation Propagation in ImageNet. ECCV 2012

About 93 results (0.07 sec)

Hedging your bets: Optimizing accuracy-specificity trade-offs in large scale visual recognition

Search within citing articles

Youtube2text: Recognizing and describing arbitrary activities using semantic hierarchies and zero-shot recognition

<u>Scuedarrama</u>, N Krishnamoorthy... - Proceedings of the ..., 2013 - cv-foundation org Abstract Despite a recent push towards large-scale object recognition, activity recognition remains limited to narrow domains and small vocabularies of actions. In this paper, we tackle the chollenge of recognizing and describing activities' in-the-wild'. We present a Cited by 129. Related articles All 13 versions Cite. Save

Reasoning about object affordances in a knowledge base representation Y Zhu, A Fathi, L Fell-Fel, European conference on computer vision, 2014 - Springer Adstract Reasoning about objects and their affordances is a fundamental problem for visual intelligence. Most of the previous work casts this problem as a classification task where separate classifiers are trained to label objects, recognize attributes, or assign affordances. Citad by 78 Related articles Al7 versions. Cita Save

[PDF] TREETALK: Composition and Compression of Trees for Image Descriptions.

<u>P. Kuznetsova</u>, <u>V. Ordonez</u>, <u>TL. Berg, Y. Chol</u> - TACL, 2014 - pdfs semanticacholar org Abstract We present a new tree based approach to composing expressive image descriptions that makes use of naturally occurring web images with captions. We investigate two related tasks: image caption generalization and generation, where the former is an Citede by 65 Related articles A112 versions Cite. Save More

From large scale image categorization to entry-level categories

<u>V Ordonez, J Deng, V Choi, AC Berg.</u>. - Proceedings of the IEEE ..., 2013 - or foundation ang Abstract Entry level categories the labels people will use to name an object were originally defined and studied by psychologists in the 1980s. In this paper we study entrylevel categories at a large scale and learn the first models for predicting entry-level categories for Clicked by 53. Related anticles A148 versions. Citle: Save

[PDF] Integrating Language and Vision to Generate Natural Language Descriptions of Videos in the Wild.

<u>J Thomson</u>, <u>S Venugapalan</u>, ..., 2014. al2-a2-ad5 a3 amazonawa com Abstract This paper integrates techniques in natural language processing and computer vision to improve recognition and description of entities and addivities in real-world videos. We propose a strategy for generaling textual descriptions of videos by using a factor graph Cated by 59. Related articles M 12 versions. Citle: Save Mere

[PDF] Hierarchical Semantic Labeling for Task-Relevant RGB-D Perception.

C Wu, I Lercy, A Saxena - Robotics: Science and systems, 2014 - pdfs semanticscholar.org Abstract—Semantic labeling of ROB-D scenes is very important in enabling robots to perform mobile manipulation tasks, but different tasks may require entirely different sets of labels. For example, when navigating to an object, we may need only a single label Cited by 47. Related articles All 3 versions. Cite Save: More

ECCV 2012 **Best paper Award**

Most works still use 1M images

to do pre-training

Nemages

15M Images Total

"First, we find that the performance on vision tasks still increases linearly with orders of magnitude of training data size."

C. Sun et al, 2017

2017 Jul 0 -> Ù cn U. arXiv:1707.02968v1

Revisiting Unreasonable Effectiveness of Data in Deep Learning Era

Chen Sun¹, Abhinav Shrivastava^{1,2}, Saurabh Singh¹, and Abhinav Gupta^{1,2}

¹Google Research ²Carnegie Mellon University

Abstract

The success of deep learning in vision can be attributed to: (a) models with high capacity; (b) increased computational power; and (c) availability of large-scale labeled data. Since 2012, there have been significant advances in representation capabilities of the models and computational capabilities of GPUs. But the size of the biggest dataset has surprisingly remained constant. What will happen if we increase the dataset size by 10× or 100×? This paper takes a step towards clearing the clouds of mystery surrounding the relationship between 'enormous data' and deep learning. By exploiting the JFT-300M dataset which has more than 375M noisy labels for 300M images, we investigate how the performance of current vision tasks would change if this data was used for representation learning. Our paper delivers some surprising (and some expected) findings. First, we find that the performance on vision tasks still increases linearly with orders of magnitude of training data size. Second, we show that representation learning (or pretraining) still holds a lot of promise. One can improve performance on any vision tasks by just training a better base model. Finally, as expected, we present new state-of-theart results for different vision tasks including image classification, object detection, semantic segmentation and human pose estimation. Our sincere hope is that this inspires vision community to not undervalue the data and develop collective efforts in building larger datasets.



Figure 1. The Curious Case of Vision Datasets: While GPU computation power and model sizes have continued to increase over the last five years, size of the largest training dataset has surprisingly remained constant. Why is that? What would have happened if we have used our resources to increase dataset size as well? This paper provides a sneak-peek into what could be if the dataset sizes are increased dramatically.

ously, while both GPUs and model capacity have continued to grow, datasets to train these models have remained stormat. First a 101 large RecNet with significantly more

How Humans Compare

STANFOR/ ENGINEERIN

Andrej Karpathy. http://karpathy.github.io/2014/09/02/what-i-learned-from-competing-against-a-convnet-on-imagenet/

How Humans Compare

Human

5.1% Top-5 error rate

Susceptible to:

- Fine-grained recognition
- Class unawareness
- Insufficient training data

GoogLeNet



Susceptible to:

- Small, thin objects
- Image filters
- Abstract representations
- Miscellaneous sources

Andrej Karpathy. http://karpathy.github.io/2014/09/02/what-i-learned-from-competing-against-a-convnet-on-imagenet/

What Lies Ahead

Moving from object recognition...



...to human-level understanding.



Inverse Graphics



Image credit: https://www.youtube.com/watch?v=ip-KIzQmcBo (Oliver Villar)
IM GENET





ImageNet: Deng et al. 2009; COCO: Lin et al. 2014





"A lady in pink dress is skiing."

COCO: Lin et al. 2014



"A lady in pink dress is skiing."

"A man standing." "A clear blue sky at a ski resort." "A snowy hill is in front of pine trees." "There are several pine trees." "A group of people getting ready to ski."

Q: What is the man in the center doing? *A: Standing on a ski*.Q: What is the color of the sky? *A: Blue* Q: Where are the pine trees? A: Behind the hill.







head





[Johnson et al., CVPR 2015]

Over Senator Visual Genome Dataset

A dataset, a knowledge base, an ongoing effort to connect structural image concepts to language.

Specs

- 108,249 images (COCO images)
- 4.2M image descriptions
- 1.8M Visual QA (7W)
- 1.4M objects, 75.7K obj. classes
- 1.5M relationships, 40.5K rel. classes
- 1.7M attributes, 40.5K attr. classes
- Vision and language correspondences
- Everything mapped to WordNet Synset

Goals

- Beyond nouns
 - Objects, verbs, attibutes
- Beyond object classification
 - Relationships and contexts
- Sentences and QAs
- From Perception to Cognition

Over Senator Visual Genome Dataset

A dataset, a knowledge base, an ongoing effort to connect structural image concepts to language.

DenseCap & Paragraph Generation Karpathy et al. CVPR'16 Krause et al. CVPR'17

Relationship Prediction Krishna et al. ECCV'16

Image Retrieval w/ Scene Graphs Johnson et al. CVPR'15 Xu et al. CVPR'17 **Visual Q&A** Zhu et al. CVPR'16









Q: What is the person sitting on the right of the elephant wearing? *A: A blue shirt.*

Krishna et al. IJCV 2016

Over Senate Over Contract Senate Contract Senate Contract Senate Contract Senate Sena

A dataset, a knowledge base, an ongoing effort to connect structural image concepts to language.

Workshop on Visual Understanding by Learning from Web Data 2017

26 July 2017 | *Honolulu, Hawaii in conjunction with CVPR 2017*

http://www.vision.ee.ethz.ch/webvision/workshop.html







Q: What is the person sitting on the right of the elephant wearing? *A: A blue shirt.*

Krishna et al. IJCV 2016

The Future of Vision and Intelligence



Agency: The integration of perception, understanding and action

Eight Years of Competitions

IM GENET

2010-2017

10× reduction of image classification error



What Happens Now?

IM GENET + Kaggle

We're passing the baton to **Kaggle**: a community of more than 1M data scientists.

Why? democratizing data is vital to democratizing AI.

image-net.org remains live at Stanford.

What Happens Now?

IM GENET + Kagge

ImageNet **Object Localization** Challenge ImageNet **Object Detection** Challenge ImageNet **Object Detection from Video** Challenge

IM GENET Contributors/Friends/Advisors

Alex Berg Michael Bernstein Edward Chang Brendan Collins Jia Deng Minh Do Wei Dong Alexei Efros Mark Everingham **Christiane Fellbaum** Adam Finkelstein **Thomas Funkhouser** Timnit Gebru

Derek Hoiem Zhiheng Huang Andrej Karpathy Aditya Khosla Jonathan Krause Fei-Fei Li Kai Li Li-Jia Li Wei Liu Sean Ma Xiaojuan Ma Jitendra Malik Dan Osherson

Eunbyung Park Chuck Rosenberg Olga Russakovksy Sanjeev Satheesh Richard Socher Hao Su Zhe Wang Andrew Zisserman

49k Amazon Mechanical Turk Workers











"This is not the end. It is not even the beginning of the end. But it is, perhaps, the end of the beginning."

WINSTON CHURCHILL