



Shifting the Frame: The Labors of ImageNet and AI Data

Alex Hanna

Director of Research

Distributed AI Research Institute

The data that transformed AI research—and possibly the world



Stanford professor and Google Cloud chief scientist Fei-Fei Li changed everything.

Image: AP Photo/Jeff Chiu

The ImageNet Dataset

Effort to "map out the entire world of objects"

Over 14 million images

Over 20,000 categories

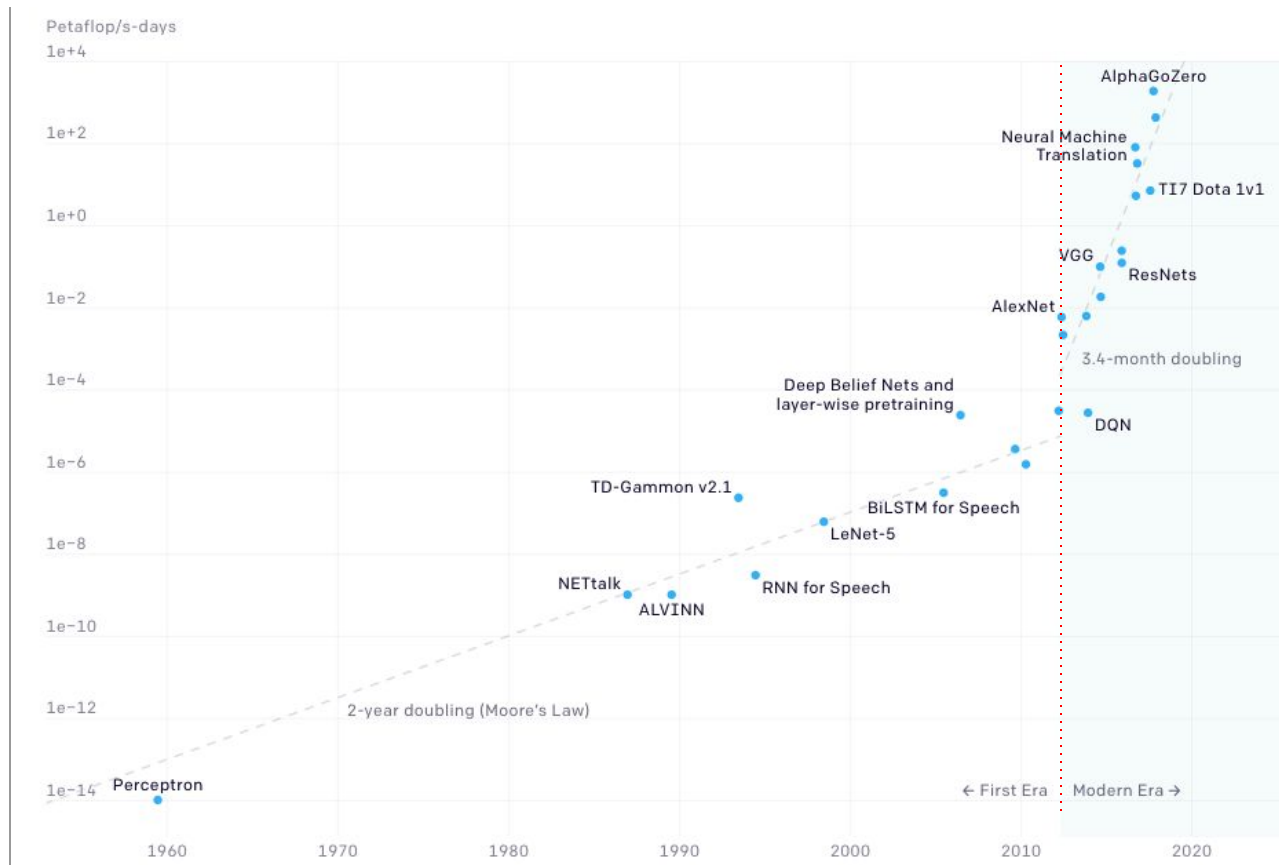
Regarded as a key benchmark



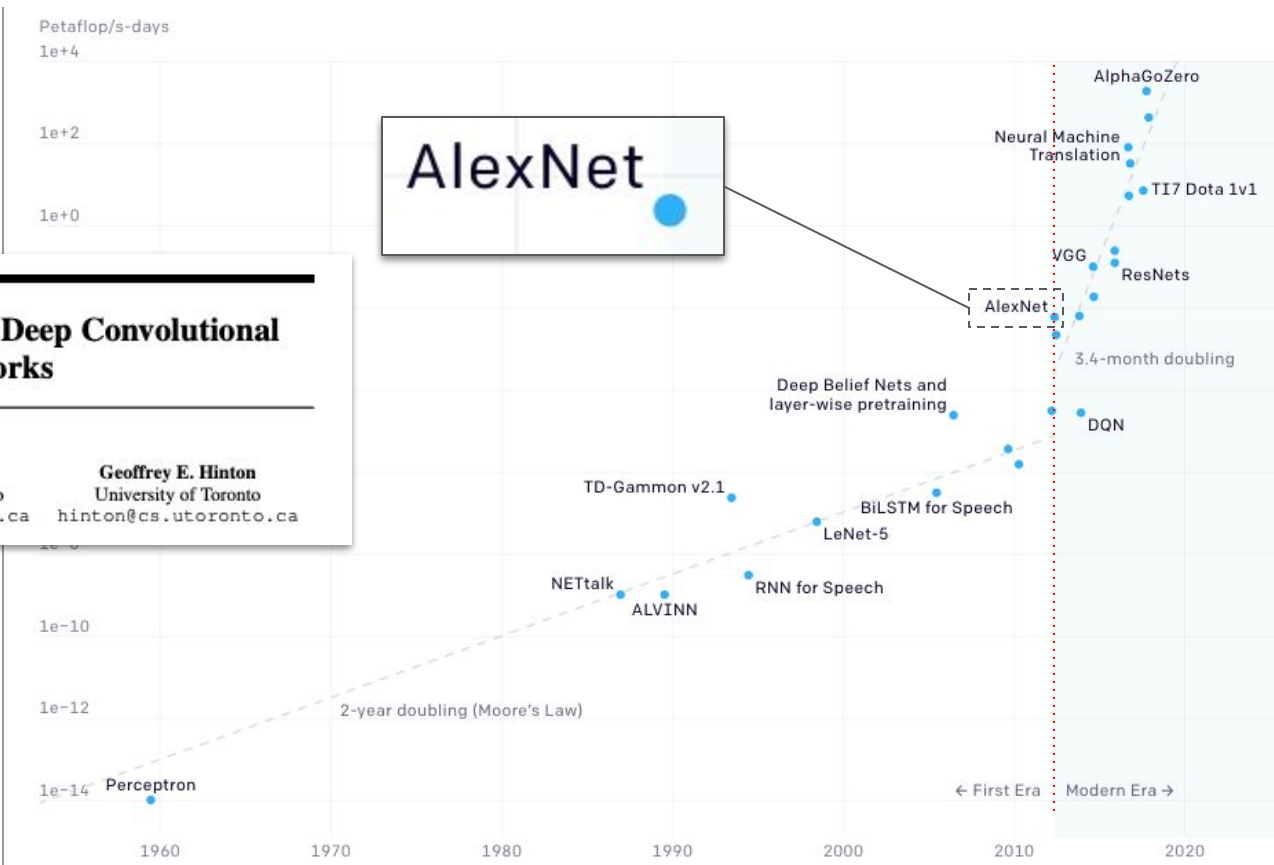
The ImageNet Challenge

Two Distinct Eras of Compute Usage in Training AI Systems

Show Error Bars All Domains ▾



The ImageNet Challenge



ImageNet Classification with Deep Convolutional Neural Networks

Alex Krizhevsky
University of Toronto
kriz@cs.utoronto.ca

Ilya Sutskever
University of Toronto
ilya@cs.utoronto.ca

Geoffrey E. Hinton
University of Toronto
hinton@cs.utoronto.ca

The ImageNet Challenge

Petaflop/s-days

1e+4

1e+2

1e+0

AlexNet

ImageNet Classification with Deep Convolutional Neural Networks

Alex Krizhevsky
University of Toronto
kriz@cs.utoronto.ca

Ilya Sutskever
University of Toronto
ilya@cs.utoronto.ca

Geoffrey E. Hinton
University of Toronto
hinton@cs.utoronto.ca

[PDF] [Imagenet classification with deep convolutional neural networks](#)

[A Krizhevsky](#), [I Sutskever](#), [GE Hinton](#) - *Advances in neural information ...*, 2012 - [kr.nvidia.com](#)

We trained a large, **deep** convolutional neural network to **classify** the 1.2 million high-resolution images in the **ImageNet** LSVRC-2010 contest into the 1000 different classes. On the test data, we achieved top-1 and top-5 error rates of 37.5% and 17.0% which is ...

☆ 📄 Cited by **77971** Related articles All 121 versions 📄

1960

1970

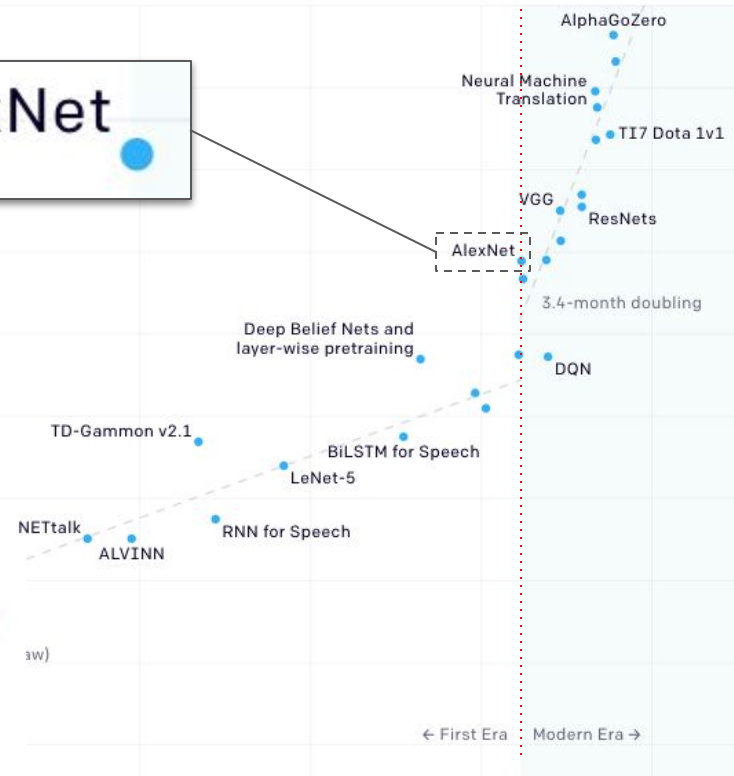
1980

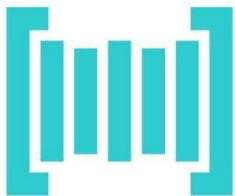
1990

2000

2010

2020





Papers with Code



CIFAR-10

The CIFAR-10 dataset (Canadian Institute for Advanced Research, 10 classes) is a subset of the Tiny Images dataset and consists of 60000 32x32 color images. The images are labell...

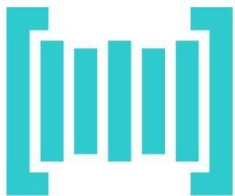
11,018 PAPERS • 69 BENCHMARKS



ImageNet

The ImageNet dataset contains 14,197,122 annotated images according to the WordNet hierarchy. Since 2010 the dataset is used in the ImageNet Large Scale Visual Recognition...

10,614 PAPERS • 102 BENCHMARKS



Papers with Code



CIFAR-10

The CIFAR-10 dataset (Canadian Institute for Advanced Research) is a subset of the Tiny Images dataset and consists of 100 classes of images.

11,018 PAPERS • 69 BENCHMARKS

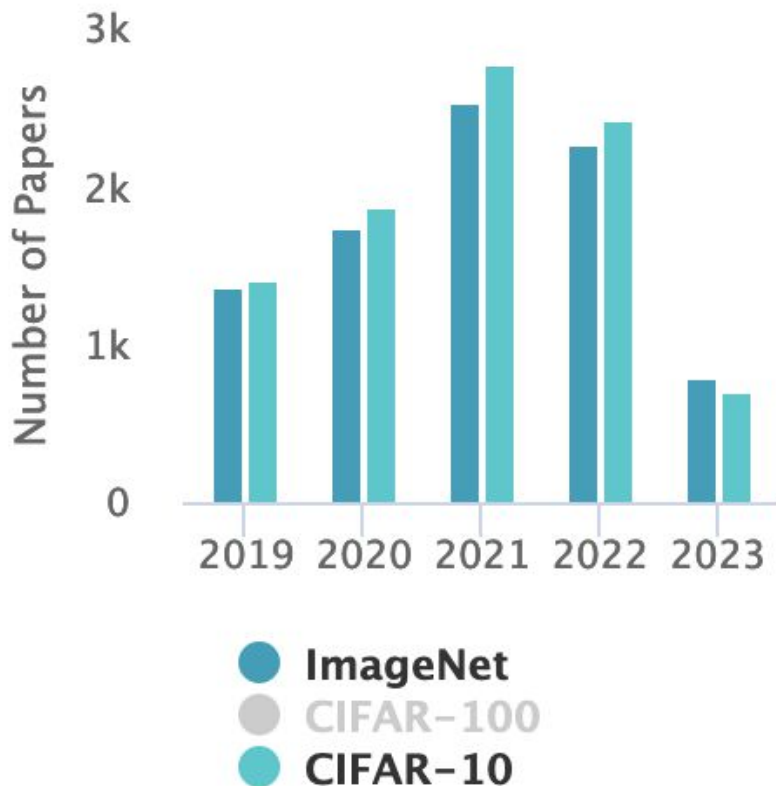


ImageNet

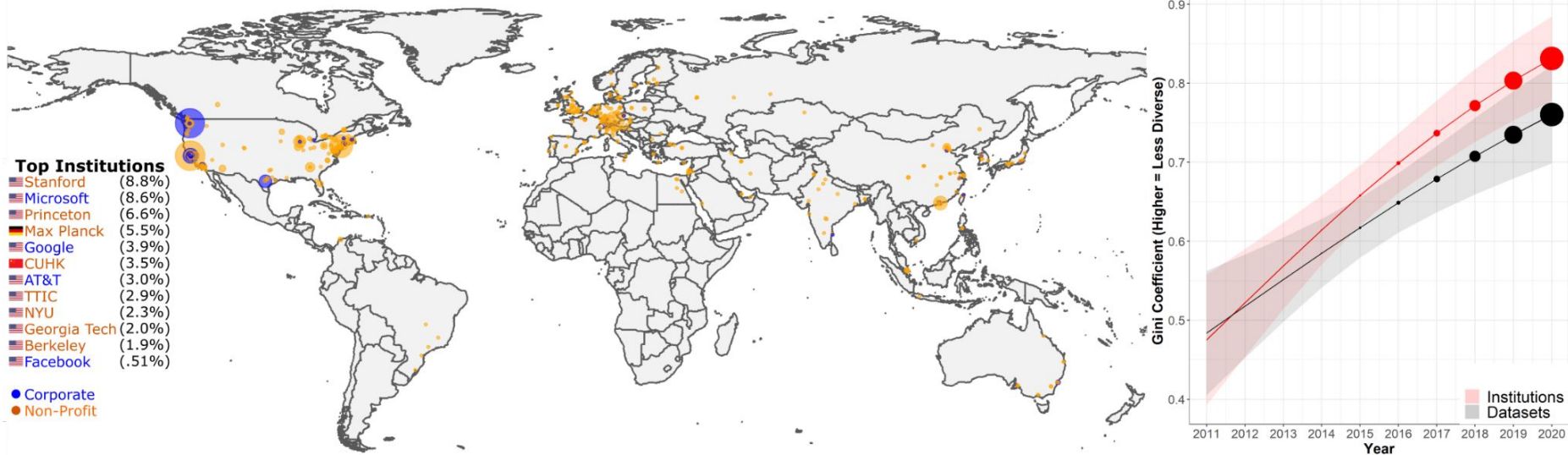
The ImageNet dataset contains 14,197 classes of images. Since 2010 the dataset is used for a wide range of computer vision tasks.

10,614 PAPERS • 102 BENCHMARKS

Usage




Concentration of Benchmark Creation



Koch et al. 2021. "Reduced, Reused and Recycled: The Life of a Dataset in Machine Learning Research." *NeurIPS (Data and Benchmark Track)*.

“The Unreasonable Effectiveness of Data”

“The IMAGENET of x ”



SpaceNet
DigitalGlobe, CosmiQ Works, NVIDIA

MusicNet
J. Thickstun et al, 2017

Medical ImageNet
Stanford Radiology, 2017

ShapeNet
A.Chang et al, 2015

EventNet
G. Ye et al, 2015

ActivityNet
F. Heilbron et al, 2015


Denton, Hanna, Amironesei, et al. 2021. “On the Genealogy of Machine Learning Datasets: A Critical History of ImageNet.” *Big Data & Society*.

“The Unreasonable
Effectiveness of Data”

The Steep Cost of Capture

o Meredith Whittaker, New York University

“The IMAGENET of x ”



SpaceNet
DigitalGlobe, CosmiQ Works, NVIDIA

MusicNet
J. Thickstun et al, 2017

Medical ImageNet
Stanford Radiology, 2017

ShapeNet
A.Chang et al, 2015

EventNet
G. Ye et al, 2015

ActivityNet
F. Heilbron et al, 2015

Denton, Hanna, Amironesei, et al. 2021. “On the Genealogy of Machine Learning Datasets: A Critical History of ImageNet.” *Big Data & Society*.

Agenda

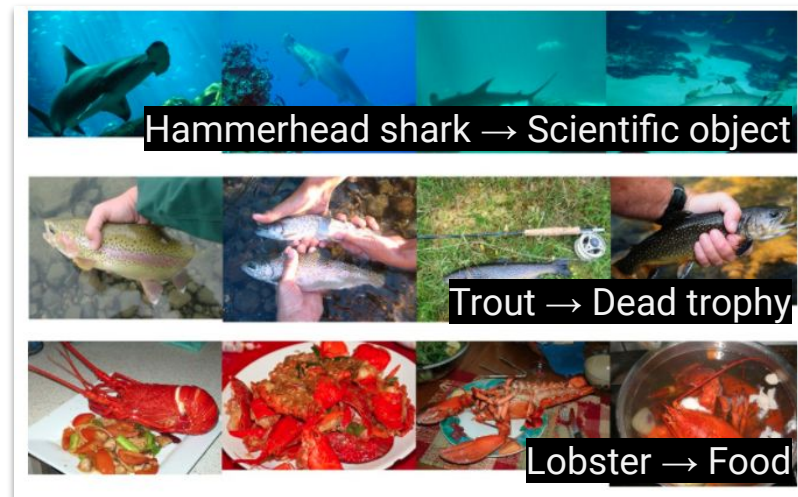
Ontologies Upon Ontologies

“Heroes Emerge Only in Times of Great Need!”

Being a Data Subject in Data-ful Times

Datasets as Infrastructure

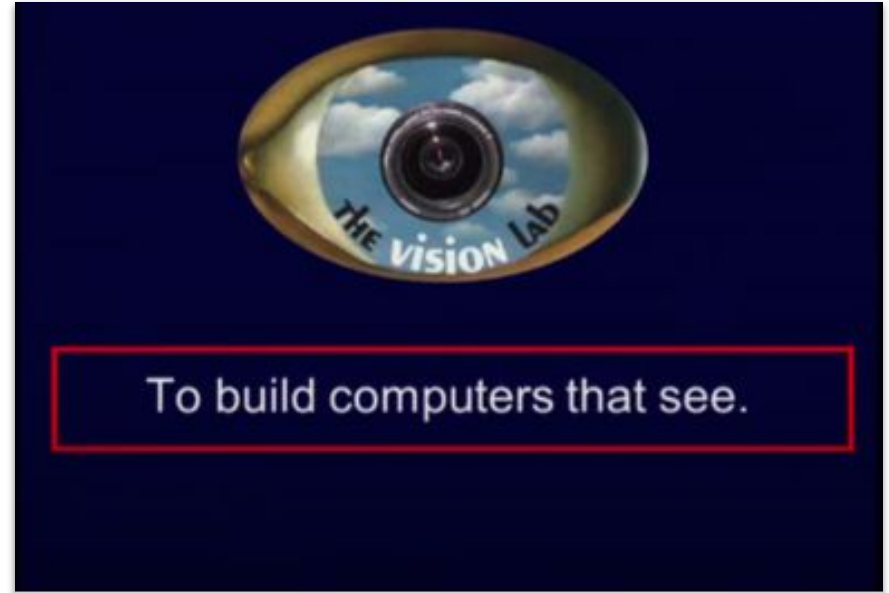
- 1) **Datasets determine what a model learns**
- 2) Datasets benchmark algorithms
- 3) Datasets serve as model organisms
- 4) Datasets provide methodological grounds for model development in industry contexts



Malevé. 2019. [An Introduction to Image Datasets.](#)

Denton and Hanna et al. 2020. "[Bringing the People Back In: Contesting Benchmark Machine Learning Datasets.](#)" PAML Workshop, *ICML*.

Computational Construction of Meaning and Understanding




Denton, Hanna, Amironesei, et al. 2021. "On the Genealogy of Machine Learning Datasets: A Critical History of ImageNet." *Big Data & Society*.

What is WordNet?

Any opinions, findings, and conclusions or recommendations expressed in this material are those of the creators of WordNet and do not necessarily reflect the views of any funding agency or Princeton University.

When writing a paper or producing a software application, tool, or interface based on WordNet, it is necessary to properly [cite the source](#). Citation figures are critical to WordNet funding.

About WordNet

WordNet® is a large lexical database of English. Nouns, verbs, adjectives and adverbs are grouped into sets of cognitive synonyms (synsets), each expressing a distinct concept. Synsets are interlinked by means of conceptual-semantic and lexical relations. The resulting network of meaningfully related words and concepts can be navigated with the [browser](#) . WordNet is also freely and publicly available for [download](#). WordNet's structure makes it a useful tool for computational linguistics and natural language processing.

WordNet Search - 3.1

- [WordNet home page](#) - [Glossary](#) - [Help](#)

Word to search for:

Display Options:

Key: "S:" = Show Synset (semantic) relations, "W:" = Show Word (lexical) relations

Where do you get the definitions for WordNet? (short answer)

Our lexicographers write them.

Where do you get the definitions for WordNet? (long answer)

From the foreword to [WordNet: An Electronic Lexical Database](#), pp. xviii-xix:

People sometimes ask, "Where did you get your words?" We began in 1985 with the words in Kučera and Francis's Standard Corpus of Present-Day Edited English (familarly known as the Brown Corpus), principally because they provided frequencies for the different parts of speech. We were well launched into that list when Henry Kučera warned us that, although he and Francis owned the Brown Corpus, the syntactic tagging data had been sold to Houghton Mifflin. We therefore dropped our plan to use their frequency counts (in 1988 Richard Beckwith developed a polysemy index that we use instead). We also incorporated all the adjectives pairs that Charles Osgood had used to develop the semantic differential. And since synonyms were critically important to us, we looked words up in various thesauruses: for example, Laurence Urdang's little "Basic Book of Synonyms and Antonyms" (1978), Urdang's revision of Rodale's "The Synonym Finder" (1978), and Robert Chapman's 4th edition of "Roget's International Thesaurus" (1977) -- in such works, one word quickly leads on to others. Late in 1986 we received a list of words compiled by Fred Chang at the Naval Personnel Research and Development Center, which we compared with our own list; we were dismayed to find only 15% overlap.

So Chang's list became input. And in 1993 we obtained the list of 39,143 words that Ralph Grishman and his colleagues at New York University included in their common lexicon, COMLEX; this time we were dismayed that WordNet contained only 74% of the COMLEX words. But that list, too, became input. In short, a variety of sources have contributed; we were not well disciplined in building our vocabulary. The fact is that the English lexicon is very large, and we were lucky that our sponsors were patient with us as we slowly crawled up the mountain.

Where do you get the definitions for WordNet? (short answer)

Our lexicographers write them.

Where do you get the definitions for WordNet? (long answer)

From the foreword to [WordNet: An Electronic Lexical Database](#), pp. xviii-xix:

People sometimes ask, "Where did you get your words?" We began in 1985 with the words in Kučera and Francis's Standard Corpus of Present-Day Edited English (famously known as **the Brown Corpus**), principally because they provided frequencies for the different parts of speech. We were well launched into that list when Henry Kučera warned us that, although he and Francis owned the Brown Corpus, the syntactic tagging data had been sold to Houghton Mifflin. We therefore dropped our plan to use their frequency counts (in 1988 Richard Beckwith developed a polysemy index that we use instead). We also incorporated all the adjectives pairs that **Charles Osgood had** used to develop the semantic differential. And since synonyms were critically important to us, we looked words up in **various thesauruses**: for example, Laurence Urdang's little "Basic Book of Synonyms and Antonyms" (1978), Urdang's revision of Rodale's "The Synonym Finder" (1978), and Robert Chapman's 4th edition of "Roget's International Thesaurus" (1977) -- in such works, one word quickly leads on to others. Late in 1986 we received a list of words compiled **by Fred Chang at the** Naval Personnel Research and Development Center, which we compared with our own list; we were dismayed to find only 15% overlap.

So Chang's list became input. And in 1993 we obtained the list of 39,143 words that **Ralph Grishman and his colleagues at New York University** included in their common lexicon, COMLEX; this time we were dismayed that WordNet contained only 74% of the COMLEX words. But that list, too, became input. In short, a variety of sources have contributed; we were not well disciplined in building our vocabulary. The fact is that the English lexicon is very large, and we were lucky that our sponsors were patient with us as we slowly crawled up the mountain.

BROWN CORPUS MAUNAL

BROWN CORPUS MAUNAL

LIST OF SAMPLES



A01	Atlanta Constitution	Political Reportage
A02	Dallas Morning News	Political Reportage
	Chicago Tribune	Political Reportage
A03	Chicago Tribune	Political Reportage
A04	Christian Science Monitor	Political Reportage
A05	Providence Journal	Political Reportage
A06	Newark Evening News	Political Reportage
A07	New York Times	Political Reportage
A08	Times-Picayune, New Orleans	Political Reportage
A09	Philadelphia Inquirer	Political Reportage
	Chicago Tribune	Political Reportage
A10	Oregonian, Portland	Political Reportage
A11	Sun, Baltimore	Sports Reportage
A12	Dallas Morning News	Sports Reportage
A13	Rocky Mountain News	Sports Reportage
	Dallas Morning News	Sports Reportage.
A14	New York Times	Sports Reportage.
A15	St. Louis Post-Dispatch	Sports Reportage
A16	Chicago Tribune	Society Reportage
A17	Rocky Mountain News	Society Reportage
	Dallas Morning News	Society Reportage
A18	Philadelphia Inquirer	Society Reportage
	Times-Picayune, New Orleans	Society Reportage

BROWN CORPUS MAUNAL

LIST OF SAMPLES



[A01](#) Atlanta Constitution
[A02](#) Dallas Morning News
Chicago Tribune
[A03](#) Chicago Tribune
[A04](#) Christian Science Monitor
[A05](#) Providence Journal
[A06](#) Newark Evening News
[A07](#) New York Times
[A08](#) Times-Picayune, New Orleans
[A09](#) Philadelphia Inquirer
Chicago Tribune
[A10](#) Oregonian, Portland
[A11](#) Sun, Baltimore
[A12](#) Dallas Morning News
[A13](#) Rocky Mountain News
Dallas Morning News
[A14](#) New York Times
[A15](#) St. Louis Post-Dispatch
[A16](#) Chicago Tribune
[A17](#) Rocky Mountain News
Dallas Morning News
[A18](#) Philadelphia Inquirer
Times-Picayune, New Orleans

[D01](#) William Pollard
[D02](#) Schubert Ogden
[D03](#) Edward E. Kelly
[D04](#) Jaroslav Pelikan
[D05](#) Perry Miller
[D06](#) A Howard Kelly
[D06B](#) Shirley Schuyler
[D06C](#) Nathanael Olson
[D07](#) Peter Eldersveld
[D08](#) Schuyler Cammann
[D09](#) Eugene E. Golay
[D10](#) Huston Smith
[D11](#) Paul Ramsey

Sports Reportage
Sports Reportage
Sports Reportage.
Sports Reportage.
Sports Reportage
Society Reportage
Society Reportage
Society Reportage
Society Reportage
Society Reportage

Physicist and Christian
Christ Without Myth
Christian Unity in England
The Shape of Death
Theodore Parker: Apostasy With in Liberalism
Out of Doubt into Faith
Not as the World Giveth
Are You in Orbit?
Faith Amid Fear
The Magic Square of Three
Organizing the Local Church
Interfaith Communication: The Contemporary Scene
War & the Christian Conscience

BROWN CORPUS MAUNAL

LIST OF SAMPLES



[A01](#) Atlanta Constitution
[A02](#) Dallas Morning News
Chicago Tribune
[A03](#) Chicago Tribune
[A04](#) Christian Science Monitor
[A05](#) Providence Journal
[A06](#) Newark Evening News
[A07](#) New York Times
[A08](#) Times-Picayune, New Orleans
[A09](#) Philadelphia Inquirer
Chicago Tribune
[A10](#) Oregonian, Portland
[A11](#) Sun, Baltimore
[A12](#) Dallas Morning News
[A13](#) Rocky Mountain News
Dallas Morning News
[A14](#) New York Times
[A15](#) St. Louis Post-Dispatch
[A16](#) Chicago Tribune
[A17](#) Rocky Mountain News
Dallas Morning News
[A18](#) Philadelphia Inquirer
Times-Picayune, New Orleans

[D01](#) William Pollard
[D02](#) Schubert Ogden
[D03](#) Edward E. Kelly
[D04](#) Jaroslav Pelikan
[D05](#) Perry Miller
[D06](#) A Howard Kelly
[D06B](#) Shirley Schuyler
[D06C](#) Nathanael Olson
[D07](#) Peter E
[D08](#) Schuyler
[D09](#) Eugene
[D10](#) Huston
[D11](#) Paul R

[F01](#) Rosemary Blackmon
[F02](#) Glenn Infield
[F03](#) Nathan Rapport
[F04](#) Ruth F. Rosevear
[F05](#) Richard S. Allen
[F06A](#) Alice Ho Austin
[F06B](#) Harold P. Winchester
[F07A](#) Marvin Sentnor and Stephen Hult
[F07B](#) Ho Walter Yoder
[F08](#) Philip Reaves
[F09A](#) David Martinson
[F09B](#) Isel D. Rugget
[F10](#) Jack Kaplan
[F11](#) Lillian Pompian
[F12](#) Marian Neater
[F13](#) Orlin J. Scoville
[F14](#) Harold Rosenberg
[F15](#) John A. O'Brien

Physicist and Christian
Christ Without Myth
Christian Unity in England
The Shape of Death
Theodore Parker: Apostasy With in Liberalism
Out of Doubt into Faith
Not as the World Giveth
Are You in Orbit?

How Much Do You Tell When You Talk?
America's Secret Poison Gas Tragedy
I've Been Here Before
North Country School Cares for the Whole Child
When Fogg Flew the Mail
Let's Discuss Retirement
What It Means to be Creative
How to Have a Successful Honeymoon
Attitudes Toward Nudity
Who Rules the Marriage Bed?
Fantastic Life & Death of the Golden Prostitute.
When It Comes to Carpets
Therapy by Witchcraft
Tooth-Straightening Today
New Methods of Parapsychology.
Part-time Farming
The Trial and Eichmann
Let's Take Birth Control Out of Politics



src



brown.py ×



brown.py > ...

```
1 import nltk
2 from nltk.corpus import brown
3
4 ## Who Rules the Marriage Bed?
5 [' '.join(x) for x in brown.sents(fileids=['cf08'])]
6
```

OUTPUT DEBUG CONSOLE PROBLEMS TERMINAL

python3.8 + ▾ □ ☒ ... ^ ×

In [44]: [' '.join(x) for x in brown.sents(fileids=['cf08'])]

Out[44]:

```
'In tradition and in poetry , the marriage bed is a place of unity and harmony .',
'The partners each bring to it unselfish love , and each takes away an equal share of pleasure and joy .',
'At its most ecstatic moments , husband and wife are elevated far above worldly cares .',
'Everything else is closed away .',
'This is the ideal .',
'But marriage experts say that such mutual contribution and mutual joy are seldom achieved .',
'Instead one partner or the other dominates the sexual relationship .',
'In the past , it has been the husband who has been dominant and the wife passive .',
'But today there are signs that these roles are being reversed .',
'In a growing number of American homes , marriage counselors report , the wife is taking a commanding role
in sexual relationships .',
'It is she who decides the time , the place , the surroundings , and the frequency of the sexual act .',
'It is she who says aye or nay to the intimate questions of sexual technique and mechanics — not the husband'
```


Ball-buster, ball-breaker

A demanding woman who destroys men's confidence

49 pictures 21.64% Popularity Wordnet ID#

Treemap Visualization Images of the Synset Downloads

- mother figure (0)
- yellow woman (0)
- white woman (0)
- jezebel (0)
- Black woman (0)
- enchantress, temptress, sylph (0)
- nymphet (0)
- B-girl, bar girl (0)
- matriarch, materfamilias
- Wac (0)
- divorcee, grass widow (0)
- vestal (0)
- debutante, deb (0)
- Cinderella (0)
- gold digger (0)
- amazon, virago (0)
- ball-buster, ball-breaker (1)**
- cat (0)
- nymph, houri (0)
- mesliza (0)
- maenad (0)
- maenad (0)
- bridesmaid, maid of honor
- nulipara (0)
- girlfriend (0)
- shiksa, shikse (0)
- dame, madam, ma'am, it
- girl wonder (0)
- foster-sister, foster sister (0)
- female offspring (2)
- woman (0)

Images of children synsets are not included. All images shown are thumbnails. Images may be subject to copyright.

Prev 1 2 Next

excavating.ai
(Crawford and Paglen, 2020)



Selections from the "Person" classes, ImageNet

Agenda

Ontologies Upon Ontologies

“Heroes Emerge Only in Times of Great Need!”

Being a Data Subject in Data-ful Times

3rd Attempt: A Godsend Emerges

**ImageNet PhD
Students**



**Crowdsourced
Labor**

amazon **mechanical turk™**
Artificial Artificial Intelligence

49k Workers from 167 Countries
2007-2010

Fei-Fei L 2017. Imagenet: Where have we
gone? Where are we going?

So are we exploiting chained prisoners?

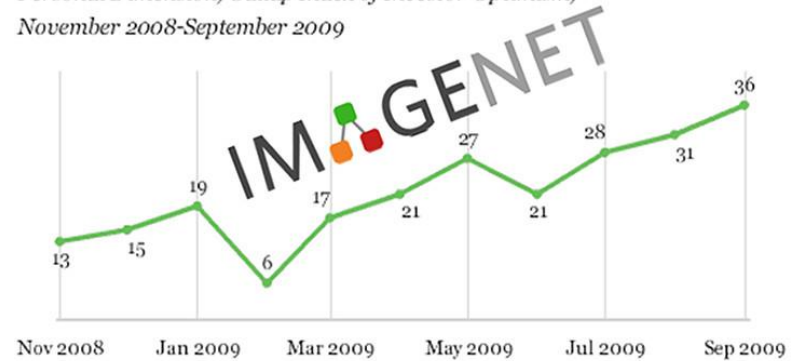


amazonmechanical turk
beta Artificial Artificial Intelligence

(a)

U.S. economy 2008 - 2009

Personal Dimension, Gallup Index of Investor Optimism,
November 2008-September 2009



IMAGENET hired more than 25,000 AMT workers in this period of time!!

(b)

Denton, Hanna, Amironesei, et al. 2021. "On the Genealogy of Machine Learning Datasets: A Critical History of ImageNet." *Big Data & Society*.

MEDIA WALL

Mimi Onuoha The Future Is Here!

The Future Is Here! is a two-part, four-screen installation that explores the idea of the future and the role of technology in our lives. The work is a response to the question: "What does the future hold for us?"

The first part of the work is a four-screen installation that explores the idea of the future and the role of technology in our lives. The work is a response to the question: "What does the future hold for us?"

The second part of the work is a four-screen installation that explores the idea of the future and the role of technology in our lives. The work is a response to the question: "What does the future hold for us?"

Copyright © 2019 Mimi Onuoha. All rights reserved. This work is a response to the question: "What does the future hold for us?"



Mimi Onuoha. 2019. "The Future Is Here!"

Turkopticon

[REQUESTER LIST](#)

[REVIEWS](#)

[ABOUT](#)

[RULES](#)

[FAQ](#)

[BLOG](#)

[REVIEW](#)

[DONATE](#)

Tweets by [@turkopticon](#)

Our mission is to organize mutual aid, resources, and advocacy to improve conditions for all people using Amazon's Mechanical Turk (AMT) platform while striving to make this work a good job for all.




Turkopticon was founded by Lilly Irani and Six Silberman as a review website to provide AMT workers with a space to share information about bad requesters and tasks. This forum also serves as a safety net. Workers are still coming to the website ten years later to check requesters' reviews before they accept tasks. **BUT, WAIT! THERE IS MORE!**

In 2019 it was decided that Turkopticon could be more than a review site. Platform workers came together to form a team with Lilly and Six. Graduate students and tech workers stepped forward to help with software needs and more. **We are so excited to be able to say that we are now a worker-led non-profit organization!**

ImageNet AMT

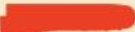


A3C8DPCYBZH618

[Averages »](#)[Requester on AMT »](#)[Review Requester »](#)[Contact the requester »](#)

FAIR: 4 / 5 
FAST: 3 / 5 
COMM: NO DATA
PAY: 1 / 5 

*Play a game of identifying birds and get paid! and *Look at car images and identify cars.

- Both of them awful underpaid. I think this requester is mocking us. I did just two hits but no way to continue wasting so much time for him.

Oct 17 2012 |     |

ImageNet AMT

A3C8DPCYBZH618

[Averages »](#)[Requester on AMT »](#)[Review Requester »](#)[Contact the requester »](#)

FAIR: 1 / 5 
FAST: 1 / 5 
COMM: 1 / 5 
PAY: 1 / 5 

One of the all time worst requesters on AMT has changed their requester ID to run from all the negative reviews. DO NOT BE FOOLED THIS IS THE SAME ImageNet AMT. See the other profile for real reviews. <http://turkopticon.differenceengines.com/ADLB46N5J8NCP>

Oct 16 2012 |     |

 Aha, now I understand! That's it, "Old habits die hard" !

Oct 18 2012 |     |

ImageNet AMT

A3C8DPCYBZH618

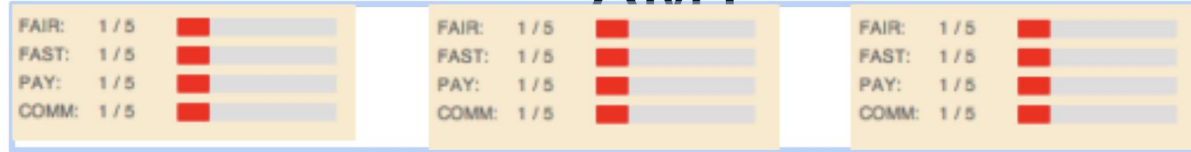
[Averages »](#)[Requester on AMT »](#)[Review Requester »](#)[Contact the requester »](#)

FAIR: NO DATA
FAST: NO DATA
COMM: NO DATA
PAY: 1 / 5 

Wow... they expect you to do 6 pages of tedious work identifying specific years & models of cars for \$0.05! I started working on one & returned it after I had already wasted 10 mins on their qualification test that is included with your first HIT.

Sep 28 2012 |     |

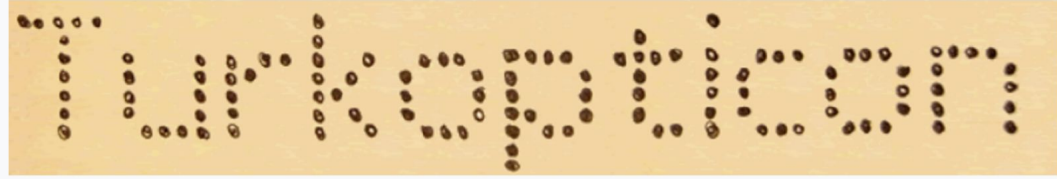
We got terrible reviews on AMT



“apparently the virtually indistinguishable cars are super distinguishable to them”

“I paid very close attention to the pictures and let the requester know that I feel these rejects were unfair. Will update when I get communication back.”

Amazon: Stop the Mass Rejections!



Amazon Mechanical Turk Workers have been organizing for years to prevent mass-rejections on the platform. These mass rejections are unfair, hurt workers, and must end. Turkopticon has organized a mass rejection protection platform that **must** be adopted to protect workers!

Mass rejections

admin December 21, 2020 1 Comment

First, allow me to introduce myself. I'm Sherry, lead organizer for Turkopticon. I have been a Turker for going on 6 years. You may ask why did I finally decide to step up and work to make Turk better? Two words: **MASS REJECTION**. My personal story involves having a great history of finding good new requesters and having plenty of wiggle room in my approval rating so I decided to work for a new requester in the same hopes of finding good work. That is when my Mturk experience drastically changed and I knew I had to stand up.

Agenda

Ontologies Upon Ontologies

“Heroes Emerge Only in Times of Great Need!”

Being a Data Subject in Data-ful Times

IBM Research Releases 'Diversity in Faces' Dataset to Advance Study of Fairness in Facial Recognition Systems



IBM Research Releases 'Diversity in Faces' Dataset to Advance Study of Fairness in Facial Recognition Systems



Facial recognition's 'dirty little secret': Millions of online photos scraped without consent

People's faces are being used without their permission, in order to power technology that could eventually be used to surveil them, legal experts say.



IBM Research Releases 'Diversity in Faces' Dataset to Advance Study of Fairness in Facial Recognition Systems



Facial recognition's 'dirty little secret': Millions of online photos scraped without consent

People's faces are being used without their permission, in order to power technology that could eventually be used to surveil them, legal experts say.



How Photos of Your Kids Are Powering Surveillance Technology

Millions of Flickr images were sucked into a database called MegaFace. Now some of those faces may have the ability to sue.

By Kashmir Hill and Aaron Krolik

IN THE UNITED STATES DISTRICT COURT
FOR THE NORTHERN DISTRICT OF ILLINOIS
EASTERN DIVISION

STEVEN VANCE, for himself and others)	
similarly situated,)	
)	
Plaintiff,)	Case No. _____
)	
v.)	
)	CLASS ACTION COMPLAINT
INTERNATIONAL BUSINESS MACHINES)	
CORPORATION,)	JURY TRIAL DEMANDED
)	
Defendant.)	INJUNCTIVE RELIEF DEMANDED
)	
)	
)	
)	
)	

CLASS ACTION COMPLAINT

Plaintiff STEVEN VANCE, on behalf of himself and all other similarly situated individuals (“Plaintiff”), by and through his attorneys, brings this Class Action Complaint against Defendant INTERNATIONAL BUSINESS MACHINES CORPORATION (“IBM”) and alleges the following:

IN THE UNITED STATES DISTRICT COURT
FOR THE NORTHERN DISTRICT OF ILLINOIS
EASTERN DIVISION

STEVEN VANCE, for himself and others)

for third parties to connect the individual whose biometrics were collected to other photos in which they appeared and other individuals appearing in those photos with them, **subjecting them to increased surveillance, stalking, identity theft, and other invasions of privacy and fraud.**

CLASS ACTION COMPLAINT

Plaintiff STEVEN VANCE, on behalf of himself and all other similarly situated individuals (“Plaintiff”), by and through his attorneys, brings this Class Action Complaint against Defendant INTERNATIONAL BUSINESS MACHINES CORPORATION (“IBM”) and alleges the following:

Excavating “Excavating AI”:
The Elephant in the Gallery
(Lyons, 2020)

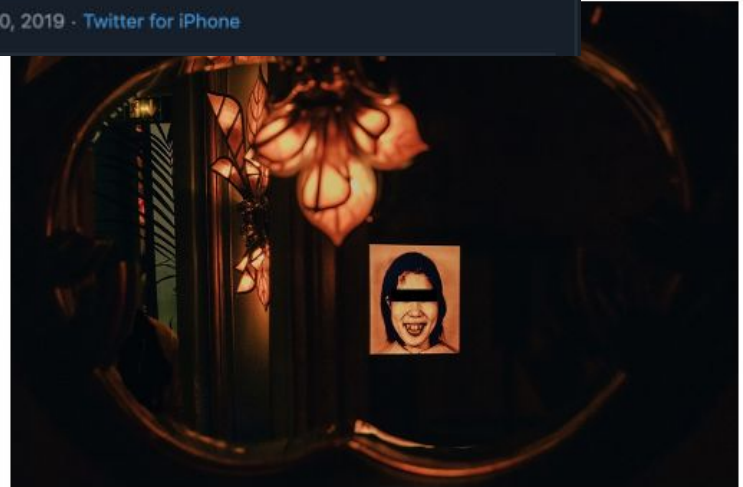


Figure 2: “Making Faces” exhibition at Maxim’s, Prada Mode Paris.



On Lacework: watching an entire machine-learning dataset
(Pipkin, 2020)

Agenda

Ontologies Upon Ontologies

“Heroes Emerge Only in Times of Great Need!”

Being a Data Subject in Data-ful Times



The cultural labor of AI data maintenance



Thank You

Alex Hanna



alex@dair-institute.org



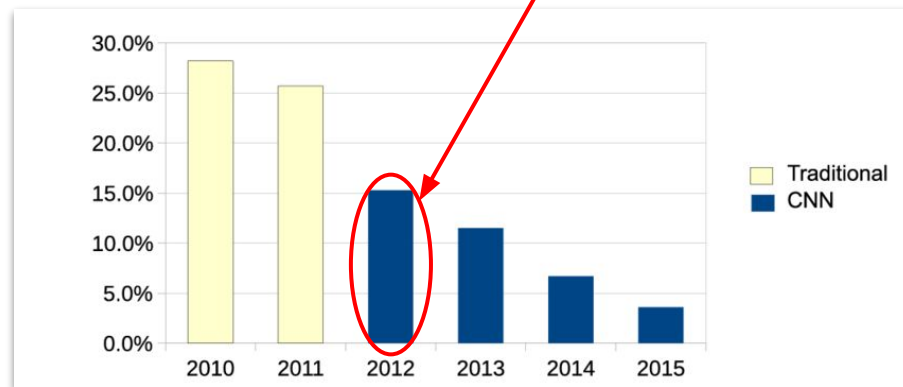
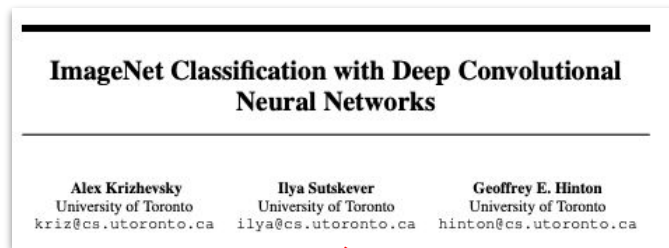
[@alexhanna](https://twitter.com/alexhanna)



[@alex@dair-community.social](https://dair-community.social/@alex)

Datasets as Infrastructure

- 1) Datasets determine what a model learns
- 2) **Datasets benchmark algorithms**
- 3) Datasets serve as model organisms
- 4) Datasets provide methodological grounds for model development in industry contexts



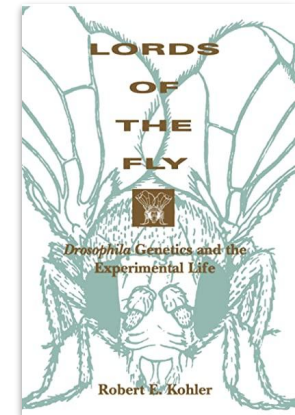
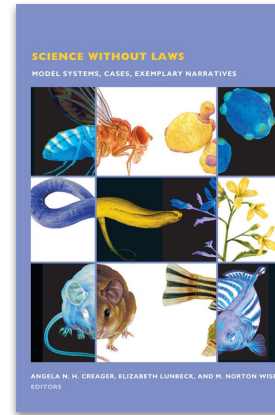
Denton and Hanna et al. 2020. "[Bringing the People Back In: Contesting Benchmark Machine Learning Datasets.](#)" PAML Workshop, *ICML*.

Datasets as Infrastructure

- 1) Datasets determine what a model learns
- 2) Datasets benchmark algorithms
- 3) **Datasets work as model organisms**
- 4) Datasets provide methodological grounds for model development in industry contexts

```
# Load ImageNet data
import tensorflow_datasets as tfds
ds = tfds.load('imagenet2012', split = 'train')

# Load neural network with ResNet50 architecture trained on ImageNet data
from tensorflow.keras.applications.resnet50 import ResNet50
model = ResNet50(weights='imagenet')
```

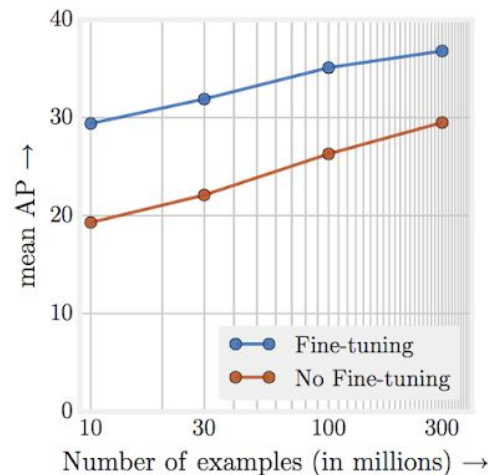


Denton and Hanna et al. 2020. "[Bringing the People Back In: Contesting Benchmark Machine Learning Datasets.](#)" PAML Workshop, *ICML*.

Datasets as Infrastructure

- 1) Datasets determine what a model learns
- 2) Datasets benchmark algorithms
- 3) Datasets work as model organisms
- 4) **Datasets provide methodological grounds for model development in industry contexts**

- Performance increases logarithmically based on volume of training data. We find there is a logarithmic relationship between performance on vision tasks and the amount of training data used for representation learning.



Object detection performance when pre-trained on different subsets of JFT-300M from scratch. x-axis is the dataset size in log-scale, y-axis is the detection performance in mAP@[.5,.95] on COCO-minival subset.

Sun et al. 2017. [Revisiting the Unreasonable Effectiveness of Data.](#)

Denton and Hanna et al. 2020. "[Bringing the People Back In: Contesting Benchmark Machine Learning Datasets.](#)" PAML Workshop, ICML.