

Yuta Tsuboi <tsuboi@preferred.jp>
Preferred Networks, Inc.

Data

Interactively Picking Real-World Objects with Unconstrained Spoken Language Instructions [Hatori+, 2017] <https://arxiv.org/abs/1710.06280>





Data is king

Can you beat the baseline method
using 10x more data?

#citations: ImageNet vs. Dropout

Note: Figures are retrieved from Google scholar

TITLE	CITED BY	YEAR
-------	----------	------

Imagenet: A large-scale hierarchical image database	5315	2009
--	------	------

J Deng, W Dong, R Socher, LJ Li, K Li, L Fei-Fei
Computer Vision and Pattern Recognition, 2009. CVPR 2009. IEEE Conference on ...

Imagenet large scale visual recognition challenge	3973	2015
--	------	------

O Russakovsky, J Deng, H Su, J Krause, S Satheesh, S Ma, Z Huang, ...
International Journal of Computer Vision 115 (3), 211-252

TITLE	CITED BY	YEAR
-------	----------	------

Dropout: a simple way to prevent neural networks from overfitting.	4272	2014
---	------	------

N Srivastava, GE Hinton, A Krizhevsky, I Sutskever, R Salakhutdinov
Journal of machine learning research 15 (1), 1929-1958

Improving neural networks by preventing co-adaptation of feature detectors	2295	2012
---	------	------

GE Hinton, N Srivastava, A Krizhevsky, I Sutskever, RR Salakhutdinov
arXiv preprint arXiv:1207.0580

ImageNet
9200+

Dropout
6500+

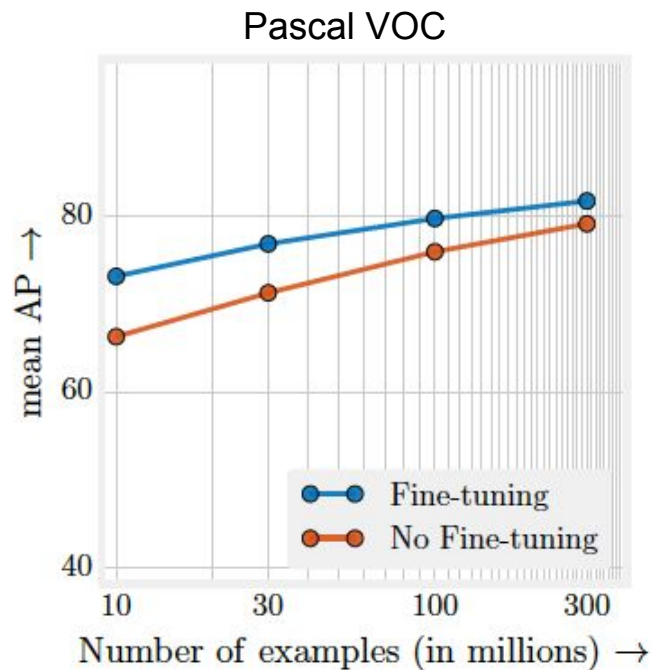
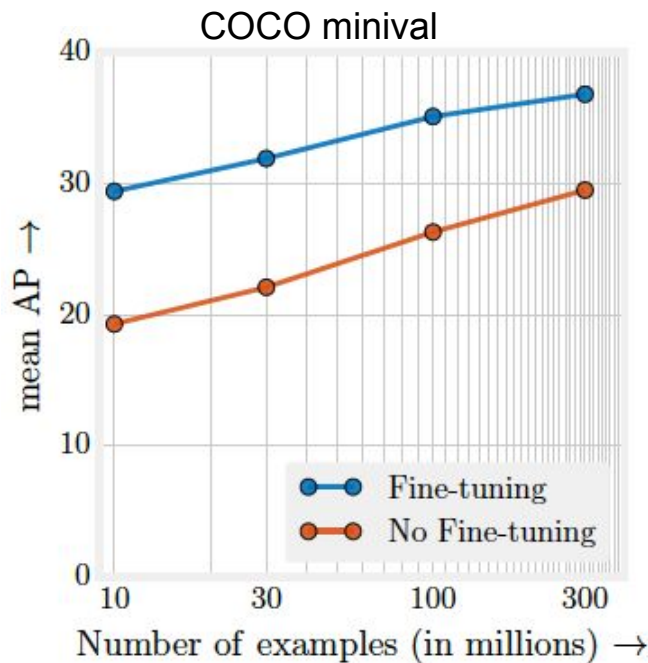
Datasets vs. Algorithms

Figure is retrieved from <http://www.spacemachine.net/views/2016/3/datasets-over-algorithms>

Year	Breakthroughs in AI	Datasets (First Available)	Algorithms (First Proposed)
1994	Human-level spontaneous speech recognition	Spoken Wall Street Journal articles and other texts (1991)	Hidden Markov Model (1984)
1997	IBM Deep Blue defeated Garry Kasparov	700,000 Grandmaster chess games, aka "The Extended Book" (1991)	Negascout planning algorithm (1983)
2005	Google's Arabic- and Chinese-to-English translation	1.8 trillion tokens from Google Web and News pages (collected in 2005)	Statistical machine translation algorithm (1988)
2011	IBM Watson became the world Jeopardy! champion	8.6 million documents from Wikipedia, Wiktionary, Wikiquote, and Project Gutenberg (updated in 2010)	Mixture-of-Experts algorithm (1991)
2014	Google's GoogLeNet object classification at near-human performance	ImageNet corpus of 1.5 million labeled images and 1,000 object categories (2010)	Convolution neural network algorithm (1989)
2015	Google's Deepmind achieved human parity in playing 29 Atari games by learning general control from video	Arcade Learning Environment dataset of over 50 Atari games (2013)	Q-learning algorithm (1992)
Average No. of Years to Breakthrough:		3 years	18 years

Performance increases linearly with orders of magnitude of training data!

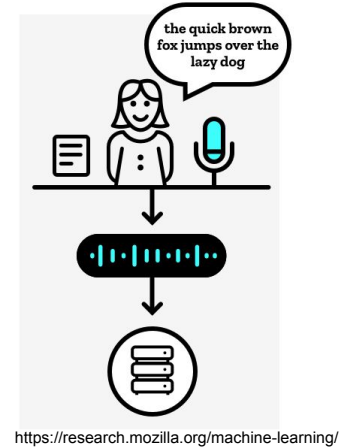
Revisiting Unreasonable Effectiveness of Data in Deep Learning Era [Sun+, 2017]



Figures are retrieved from the original paper

Supervised Data

- Manual annotations (the first wave)



- Logs of human/business activities (the second wave)

amazon.com

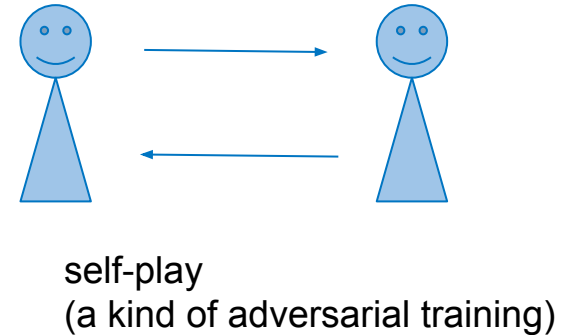
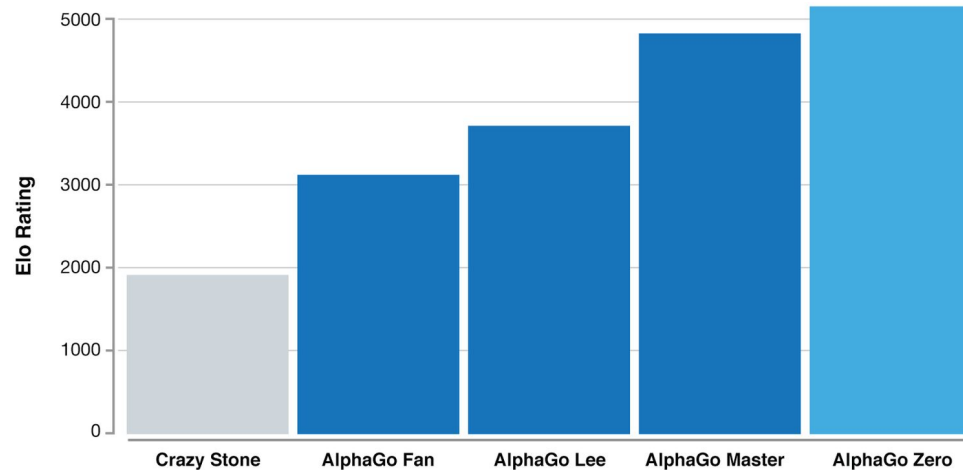


The third wave

- Self-play / self-supervised data

AlphaGo Zero [Silver+, 2017] → Self-play data

- Learning from scratch: learns to play simply by playing games against itself
 - AlphaGo was initially trained on thousands of human games



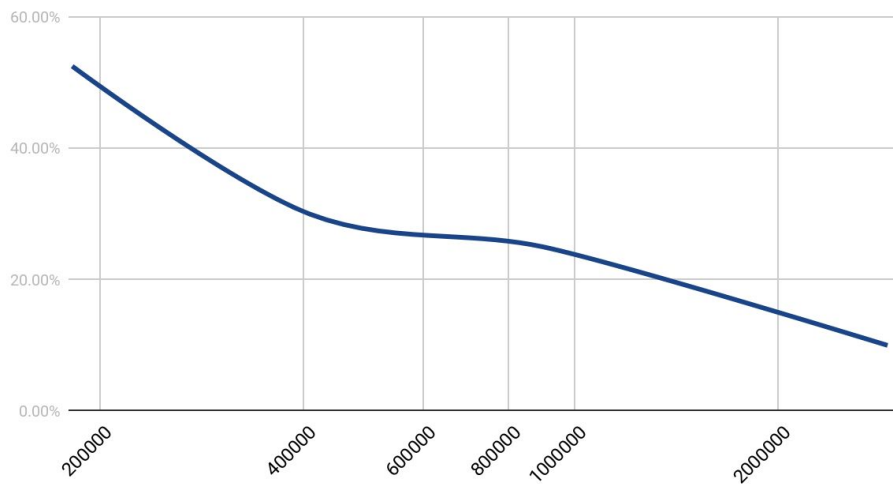
Self-supervised data for visual robot control

Two months using 6 - 14 robotic manipulators

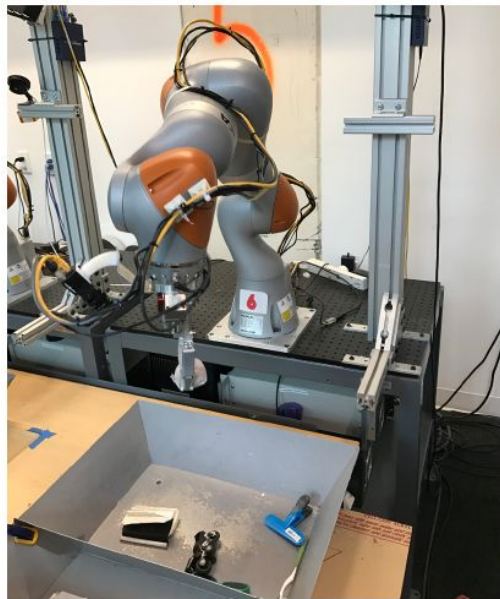
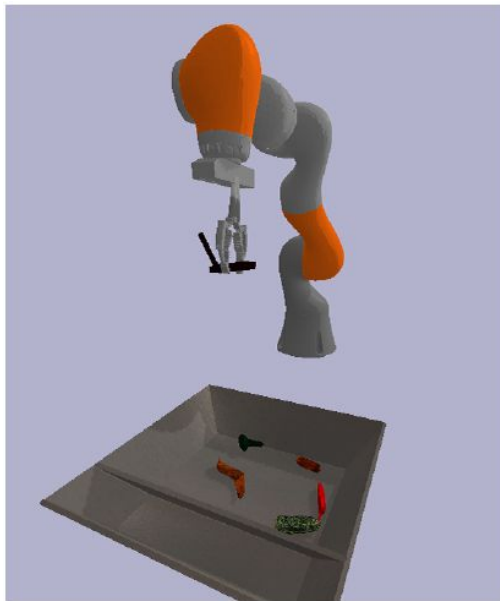
Learning Hand-Eye Coordination for Robotic Grasping with Deep Learning and Large-Scale Data Collection [Levin+, 2016]



Failure rates of grasp for varying dataset sizes



Self-supervised data generated by simulator



Physical experiments are extremely time-consuming and expensive
→ Using simulator to generate synthetic experience (2K simulated robots)

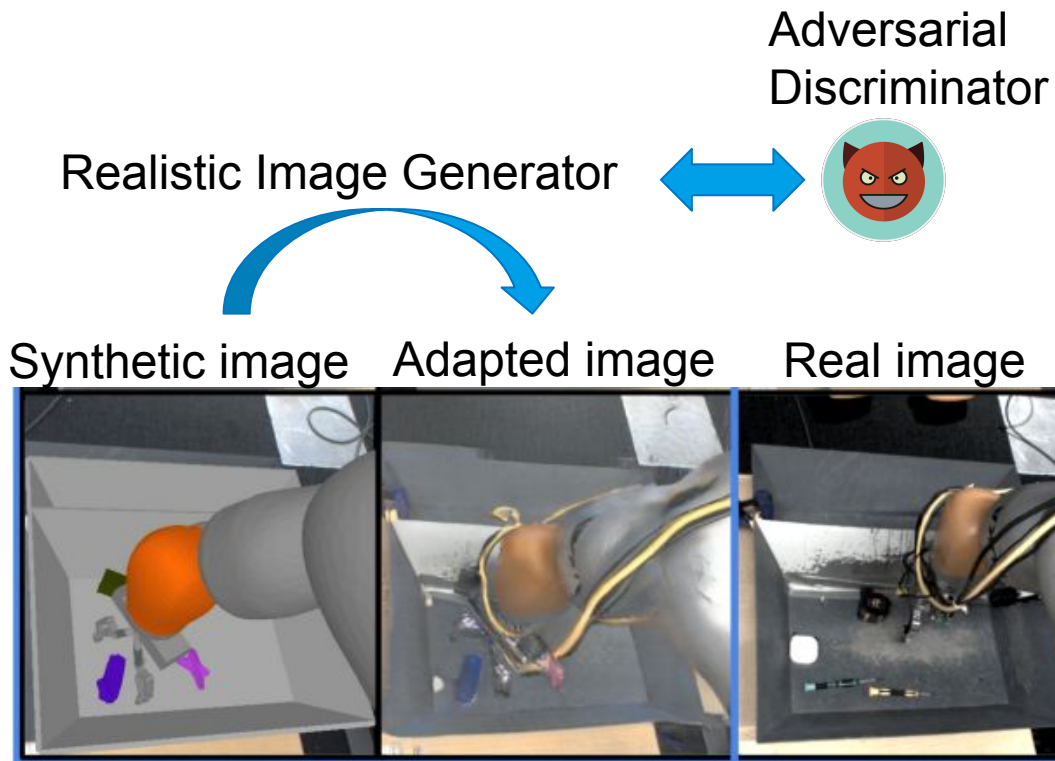
Bridging the reality gap by domain adaptation techniques

[Bousmalis+, 2017]

Procedurally generated
random objects



Simulator



Figures are retrieved from the original paper

Take-home messages

- Data might be the key limiting factor to development of AI
- Self-play or self-supervised data could be the third wave
- Simulators could exceed the limit of physical environment

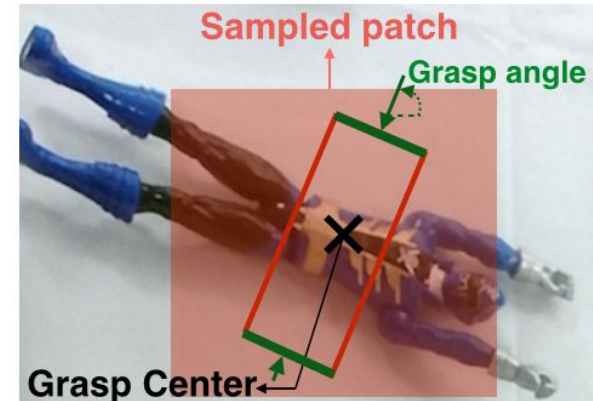


Figure is retrieved from [Pinto and Gupta, 2015]