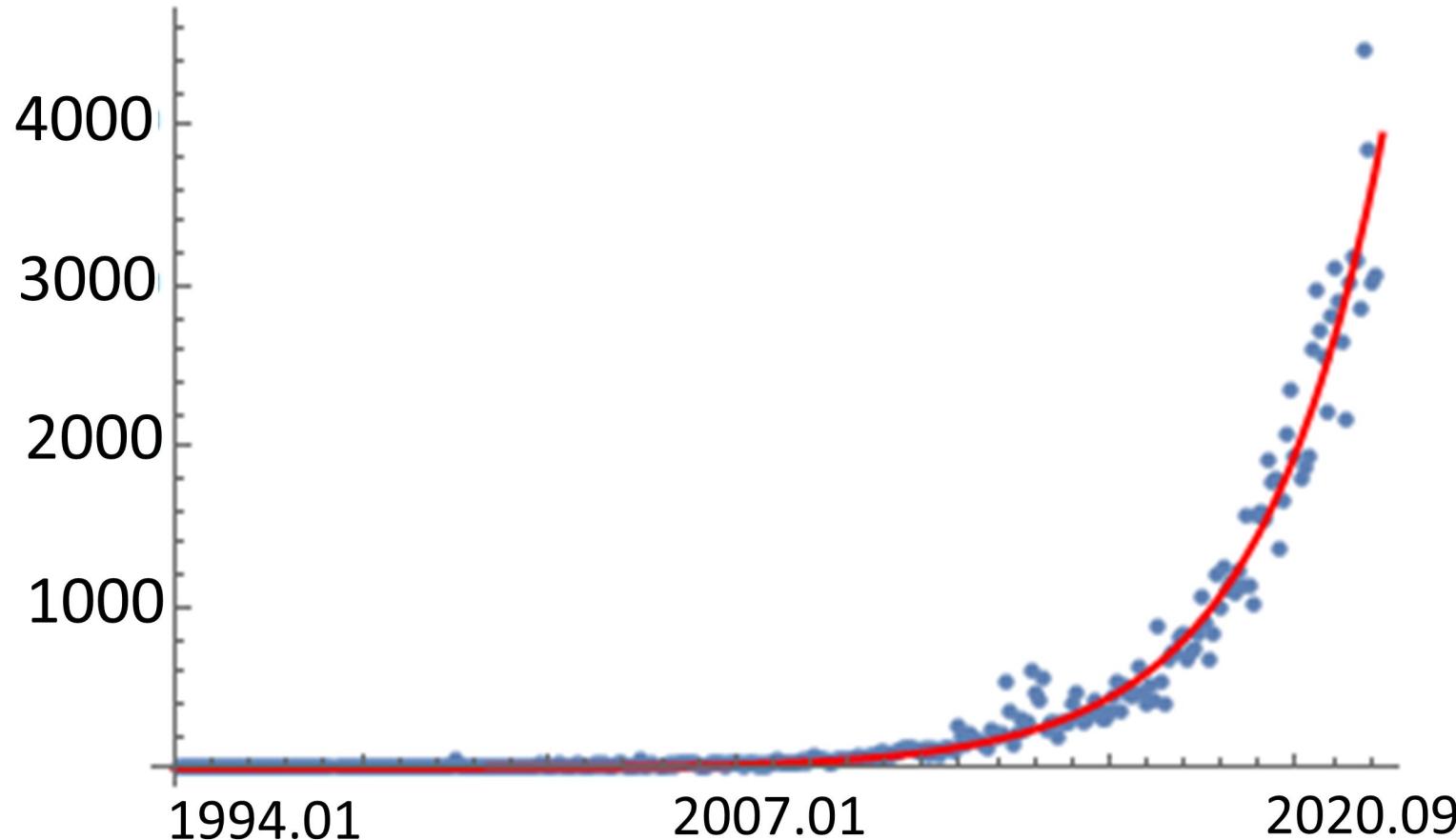


# **Full Stack Deep Learning**

## **Research Directions**

Pieter Abbeel, Sergey Karayev, Josh Tobin

# ML+AI arXiv papers per month



# Outline

---

- Sampling of research directions
- Overall research theme
- How to keep up

# Many Exciting Directions in AI

---

- Unsupervised Learning
- Reinforcement Learning
- Unsupervised RL
- Meta-Reinforcement Learning
- Few-Shot Imitation
- Domain Randomization
- DL for Science and Engineering
- Mitigating Bias
- Multi-modal Learning
- Architecture Search
- Value Alignment
- Scaling Laws
- Human-in-the-Loop
- Explainability

# Many Exciting Directions in AI

- Unsupervised Learning
- Reinforcement Learning
- Unsupervised RL
- Meta-Reinforcement Learning
- Few-Shot Imitation
- Domain Randomization
- DL for Science and Engineering
- Mitigating Bias
- Multi-modal Learning
- Architecture Search
- Value Alignment
- Scaling Laws
- Human-in-the-Loop
- Explainability

# Many Exciting Directions in AI

- ***Unsupervised Learning***
- Reinforcement Learning
- Unsupervised RL
- Meta-Reinforcement Learning
- Few-Shot Imitation
- Domain Randomization
- DL for Science and Engineering
- Mitigating Bias
- Multi-modal Learning
- Architecture Search
- Value Alignment
- Scaling Laws
- Human-in-the-Loop
- Explainability

# Recall: Deep Supervised Learning

---

- Works!
- BUT: requires so much annotated data

# So: Can we learn with less labels?

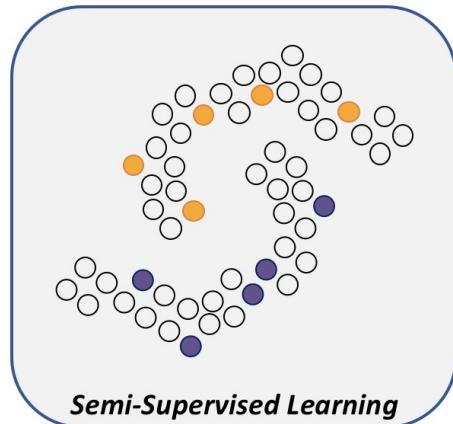
---

Yes!

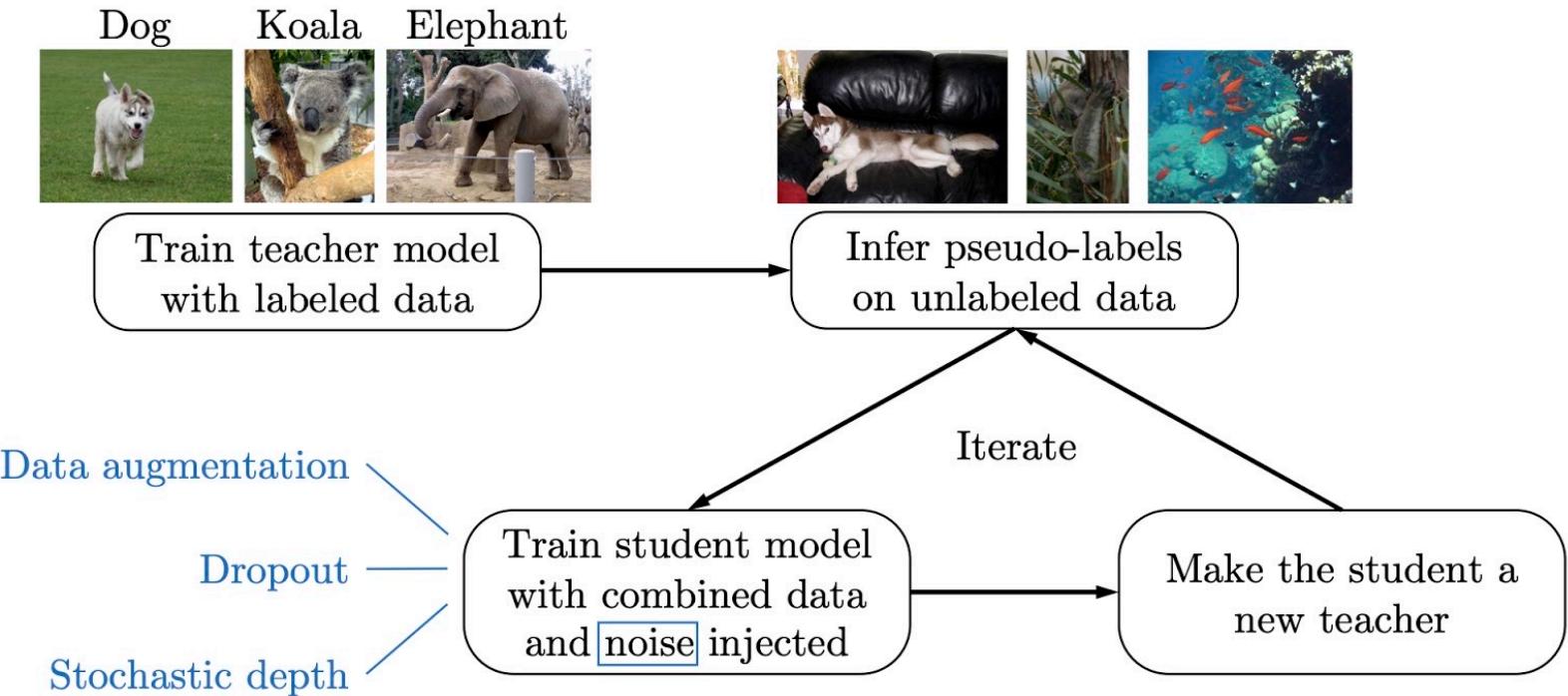
- *Deep semi-supervised learning*
- Deep unsupervised learning

# Semi-supervised Learning

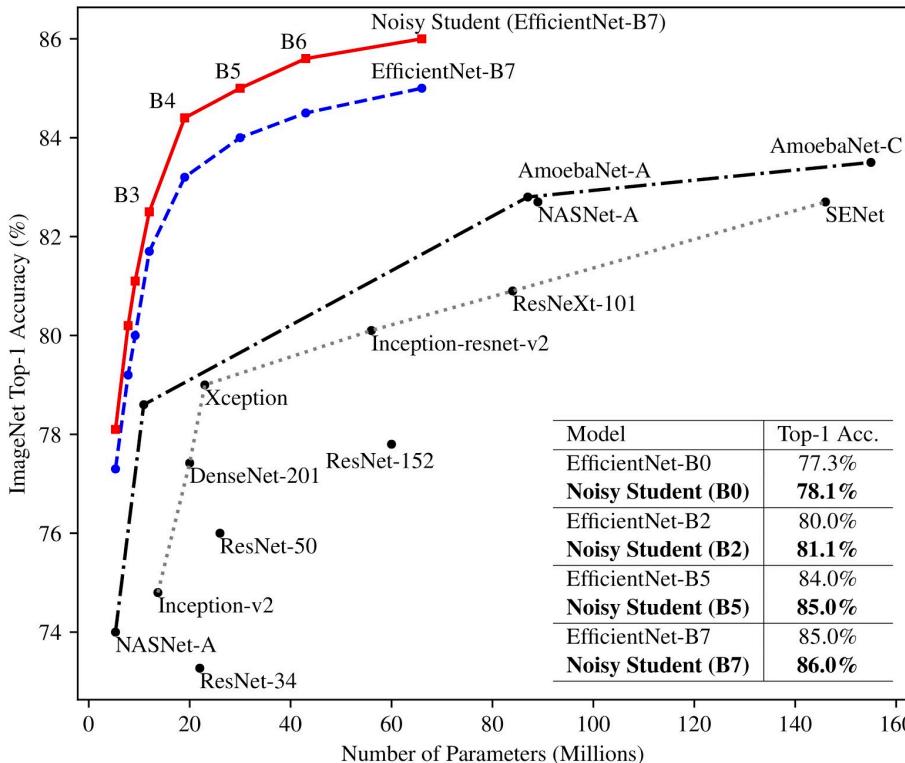
- Assumption:
  - Classification problem
  - Each data-point belongs to one of the classes
- Toy Example: “two moons”



# Noisy-Student



# Noisy-Student -- Results



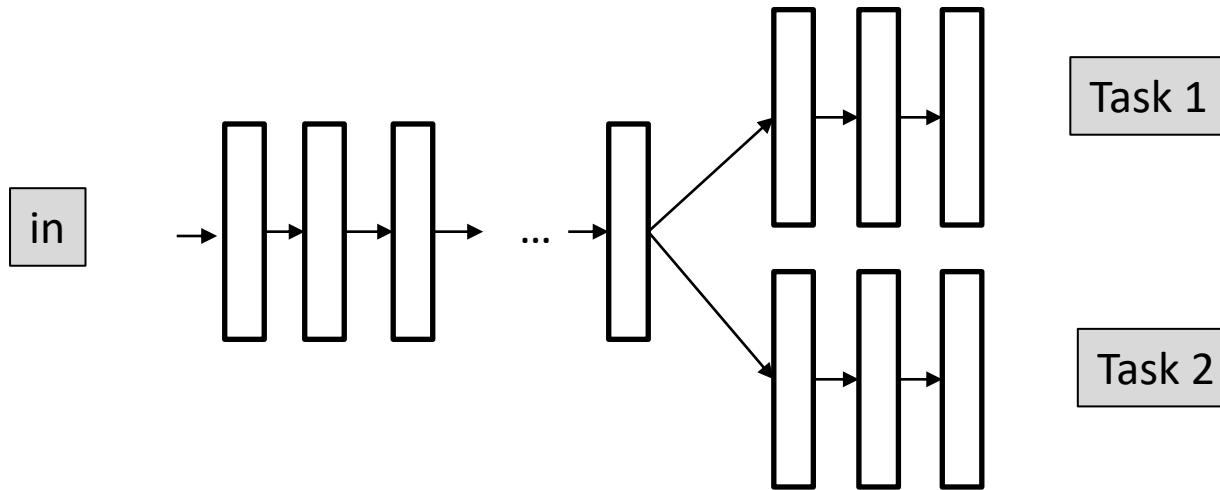
# So: Can we learn with less labels?

---

Yes!

- Deep semi-supervised learning
- *Deep unsupervised learning*

# Transfer with Multi-headed Networks



# Deep Unsupervised Learning

- Key hypothesis:

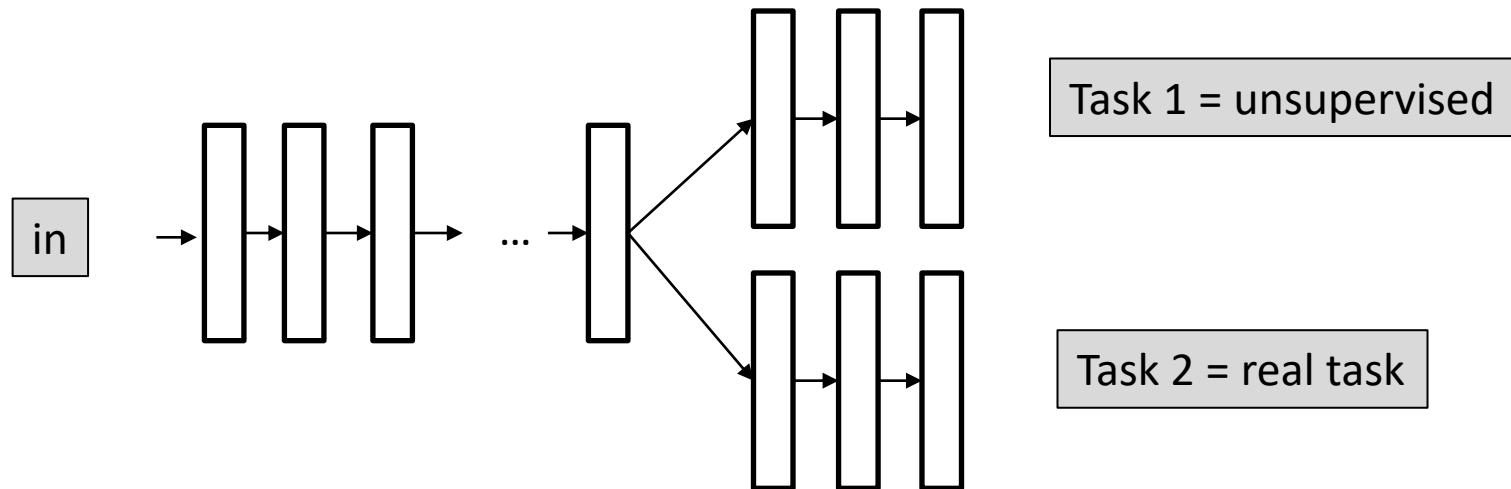
Task 1

- IF neural network smart enough to predict:
  - Next frame in video
  - Next word in sentence
  - Generate realistic images
  - ``Translate'' images
  - ...

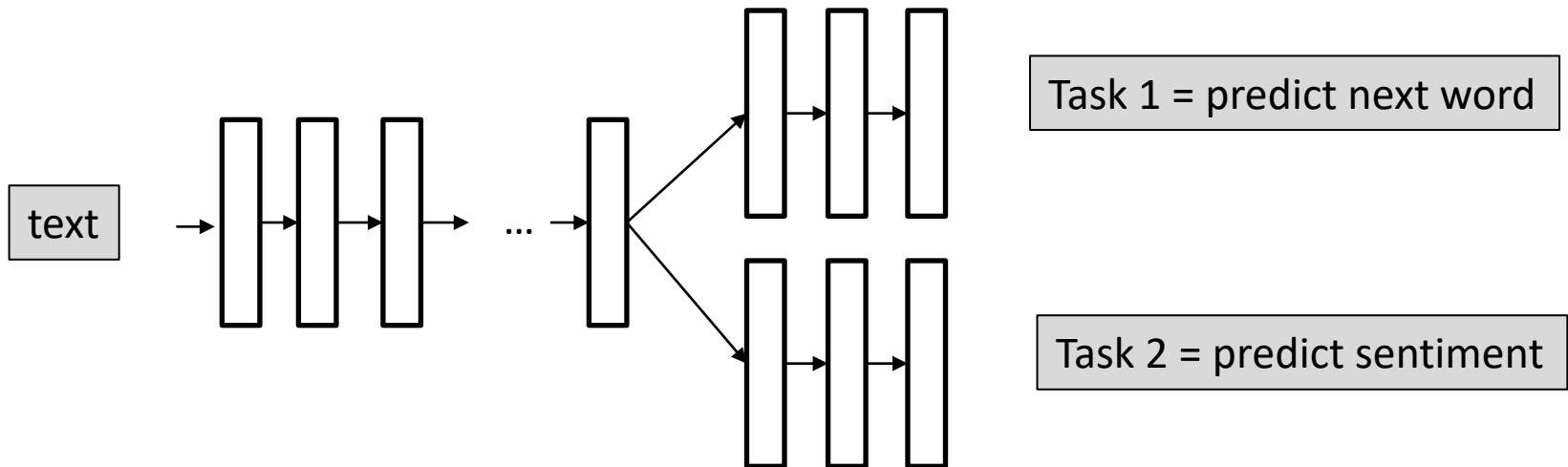
Task 2

- THEN same neural network is ready to do Deep Supervised Learning from very small data-set

# Transfer from Unsupervised Learning



# Example Setting



# Next Word Prediction (OpenAI's GPT-2)

Prompt: A train carriage containing controlled nuclear materials was stolen in Cincinnati today. Its whereabouts are unknown.

GPT-2: *The incident occurred on the downtown train line, which runs from Covington and Ashland stations. In an email to Ohio news outlets, the U.S. Department of Energy said it is working with the Federal Railroad Administration to find the thief.*

*'The theft of this nuclear material will have significant negative consequences on public and environmental health, our workforce and the economy of our nation,' said Tom Hicks, the U.S. Energy Secretary, in a statement. 'Our top priority is to secure the theft and ensure it doesn't happen again.'*

*The stolen material was taken from the University of Cincinnati's Research Triangle Park nuclear research site, according to a news release from Department officials.*

*The Nuclear Regulatory Commission did not immediately release any information. According to the release, the U.S. Department of Energy's Office of Nuclear Material Safety and Security is leading that team's investigation.*

*'The safety of people, the environment and the nation's nuclear stockpile is our highest priority,' Hicks said. 'We will get to the bottom of this and make no excuses.'*

# Next Word Prediction (OpenAI's GPT-2)

---

**Prompt:** Recycling is good for the world, no, you could not be more wrong.

**GPT-2:** Recycling is NOT good for the world. It is bad for the environment, it is bad for our health, and it is bad for our economy. I'm not kidding. Recycling is not good for the environment. It is destructive to the earth and it is a major contributor to global warming. Recycling is not good for our health. It contributes to obesity and diseases like heart disease and cancer. Recycling is bad for our economy. It increases the cost of a product, and in turn, the price of everything that is made with that product. Recycling is not good for our nation. We pay a tremendous price for the privilege of having the world's most advanced and efficient recycling system. Recycling is a huge, colossal waste of time, energy, money, and resources.

# Text Generation I's GPT-2)

SYSTEM PROMPT  
(HUMAN-WRITTEN)

In a shocking find, a remote, previously unknown language was surprising to the scientists.



MODEL COMPLETION  
(MACHINE-WRITTEN,  
10 TRIES)

The scientist named it Ovid's Unicorn. This previously unknown language

Now, after almost two years, this odd phenomenon is finally explained.

Dr. Jorge Pérez, an environmentalist from La Paz, and several colleagues were climbing the Andes Mountains when they found a herd of creatures that looked like humans. Pérez noticed that the creatures had long, thin tails and could shoot streams of water from a natural fountain, surrounded by a ring of snow.

Pérez and the others then realized that the time we reached the top of the mountain with some crystals on top,"

Pérez and his friends were amazed. These creatures could be seen too much to see them – they were covered in horns.

OpenAI built a text generator so good, it's considered too dangerous ...  
TechCrunch - 17 Feb 2019  
OpenAI built a text generator so good, it's considered too dangerous to release ...  
OpenAI said its new natural language model, GPT-2, was trained to ... said, it's ...  
only releasing a smaller version of the language model, citing its ...  
Scientists Developed an AI So Advanced They Say It's Too Dangerous ...  
ScienceAlert - 18 Feb 2019

AI text writing technology too dangerous to release, creators claim  
The Drum - 17 Feb 2019  
This technology could 'absolutely devastate' the internet as we know it  
NEWS.com.au - 17 Feb 2019  
This AI is so good at writing that its creators won't let you use it  
In-Depth - CNN - 18 Feb 2019  
Lord of The Rings, Celebrity Gossip: This AI is So Good at Writing That ...  
In-Depth - News18 - 18 Feb 2019

[View all](#)

When Is Technology Too Dangerous to Release to the Public?  
Slate Magazine - 22 Feb 2019  
If your knowledge of the model, called GPT-2, came solely on headlines ... U.K.  
read, "Elon Musk-Founded OpenAI Builds Artificial Intelligence So ... had trained a  
language model using text from 8 million webpages to predict ...  
AI Weekly: Experts say OpenAI's controversial model is a potential ...  
In-Depth - VentureBeat - 22 Feb 2019

[View all](#)

OpenAI's Text Model so Disruptive it's Deemed Too Dangerous To ...  
Computer Business Review - 15 Feb 2019  
OpenAI's Text Model so Disruptive it's Deemed Too Dangerous To Release ...  
OpenAI has declined to release the full research due to concerns over ... We've ...  
trained an unsupervised language model that can generate ...  
New AI fake text generator may be too dangerous to release, say ...  
Highly Cited - The Guardian - 14 Feb 2019

[View all](#)

bizarre creatures the scientists discovered spoke some fairly regular English. Pérez said, for example, that they have a common language that sounds like a dialect or dialectic."

He said the unicorns may have originated in South America. "Unicorns were believed to be descendants of a species that had lived there before the arrival of humans in South America."

It is currently unclear, some believe that perhaps the first time a human and a unicorn met each other in prehistoric civilization. According to Pérez, such meetings seem to be quite common."

He added that it is likely that the only unicorns that still exist are indeed the descendants of the mythical creature. "But they seem to be able to communicate, which I believe is a sign of a highly developed social organization," said the scientist.

# Unsupervised Sentiment Neuron

---

This is one of Crichton's best books. The characters of Karen Ross, Peter Elliot, Munro, and Amy are beautifully developed and their interactions are exciting, complex, and fast-paced throughout this impressive novel. And about 99.8 percent of that got lost in the film. Seriously, the screenplay AND the directing were horrendous and clearly done by people who could not fathom what was good about the novel. I can't fault the actors because frankly, they never had a chance to make this turkey live up to Crichton's original work. I know good novels, especially those with a science fiction edge, are hard to bring to the screen in a way that lives up to the original. But this may be the absolute worst disparity in quality between novel and screen adaptation ever. The book is really, really good. The movie is just dreadful.

# Benchmarks

**TASK**    **Common Sense Reasoning:** resolution of an ambiguous pronoun

**DATASET**    Winograd Schema Challenge

**EXAMPLES**    *The trophy doesn't fit into the brown suitcase because it is too large.*

**Correct answer:** *it = trophy*

**Model answer:** *it = trophy*

*The trophy doesn't fit into the brown suitcase because it is too small.*

**Correct answer:** *it = suitcase*

**Model answer:** *it = suitcase*

# Benchmarks

TASK      **Question Answering**

DATASET      Natural Questions

EXAMPLES      *Who wrote the book the origin of species?*

**Correct answer:** Charles Darwin

**Model answer:** Charles Darwin

*What is the largest state in the U.S. by land mass?*

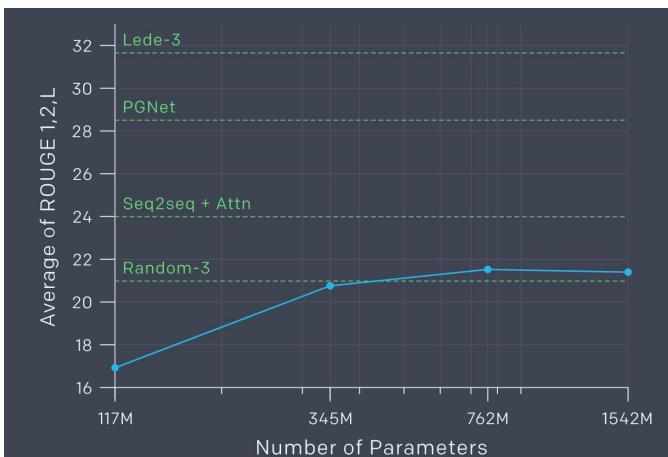
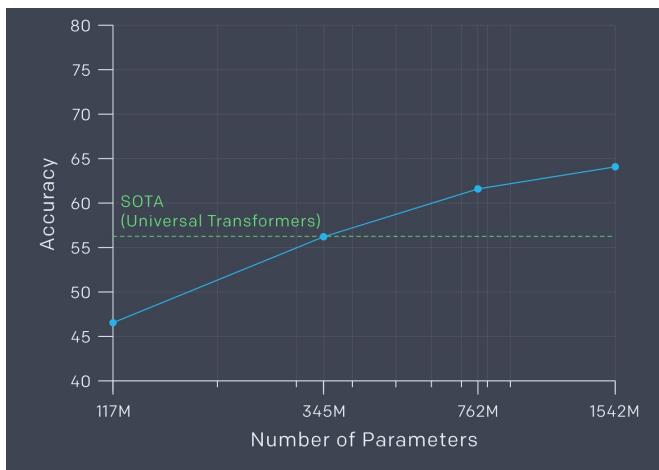
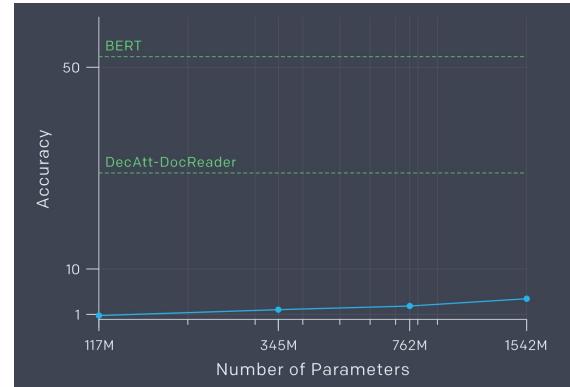
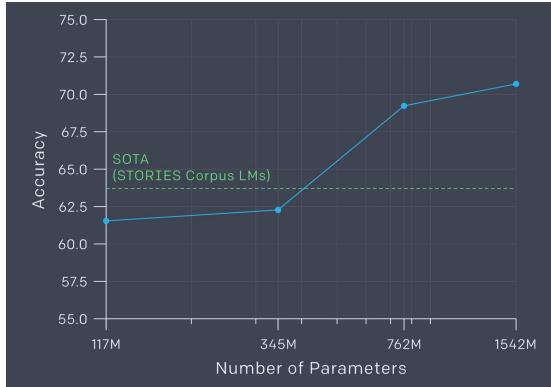
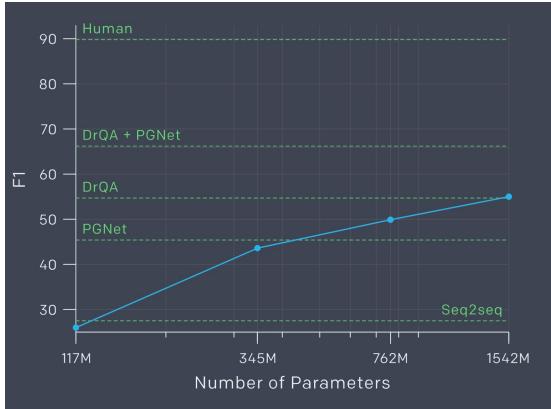
**Correct answer:** Alaska

**Model answer:** California

# Benchmarks

Dataset	Metric	Our Result	Previous Record	Human
Winograd Schema Challenge	accuracy (+)	<b>70.70%</b>	63.7%	92%+
LAMBADA	accuracy (+)	<b>63.24%</b>	59.23%	95%+
LAMBADA	perplexity (-)	<b>8.6</b>	99	~1-2
Children's Book Test Common Nouns (validation accuracy)	accuracy (+)	<b>93.30%</b>	85.7%	96%
Children's Book Test Named Entities (validation accuracy)	accuracy (+)	<b>89.05%</b>	82.3%	92%
Penn Tree Bank	perplexity (-)	<b>35.76</b>	46.54	unknown
WikiText-2	perplexity (-)	<b>18.34</b>	39.14	unknown

# Scaling



# BERT

- Main idea:
    - GPT-2: predict next word / token
    - BERT: predict a word / token that was removed

# BERT

**Sentence A** = The man went to the store.  
**Sentence B** = He bought a gallon of milk.  
**Label** = IsNextSentence

**Sentence A** = The man went to the store.  
**Sentence B** = Penguins are flightless.  
**Label** = NotNextSentence

# BERT

## GLUE Results

System	MNLI-(m/mm)	QQP	QNLI	SST-2	CoLA	STS-B	MRPC	RTE	Average
	392k	363k	108k	67k	8.5k	5.7k	3.5k	2.5k	-
Pre-OpenAI SOTA	80.6/80.1	66.1	82.3	93.2	35.0	81.0	86.0	61.7	74.0
BiLSTM+ELMo+Attn	76.4/76.1	64.8	79.9	90.4	36.0	73.3	84.9	56.8	71.0
OpenAI GPT	82.1/81.4	70.3	88.1	91.3	45.4	80.0	82.3	56.0	75.2
BERT <sub>BASE</sub>	84.6/83.4	71.2	90.1	93.5	52.1	85.8	88.9	66.4	79.6
BERT <sub>LARGE</sub>	<b>86.7/85.9</b>	<b>72.1</b>	<b>91.1</b>	<b>94.9</b>	<b>60.5</b>	<b>86.5</b>	<b>89.3</b>	<b>70.1</b>	<b>81.9</b>

### MultiNLI

Premise: Hills and mountains are especially sanctified in Jainism.

Hypothesis: Jainism hates nature.

Label: Contradiction

### CoLa

Sentence: The wagon rumbled down the road.

Label: Acceptable

Sentence: The car honked down the road.

Label: Unacceptable

Single Grain • Google • What Is the Google BERT Search Algorithm Update?

# What Is the Google BERT Search Algorithm Update?

JOYDEEP BHATTACHARYA

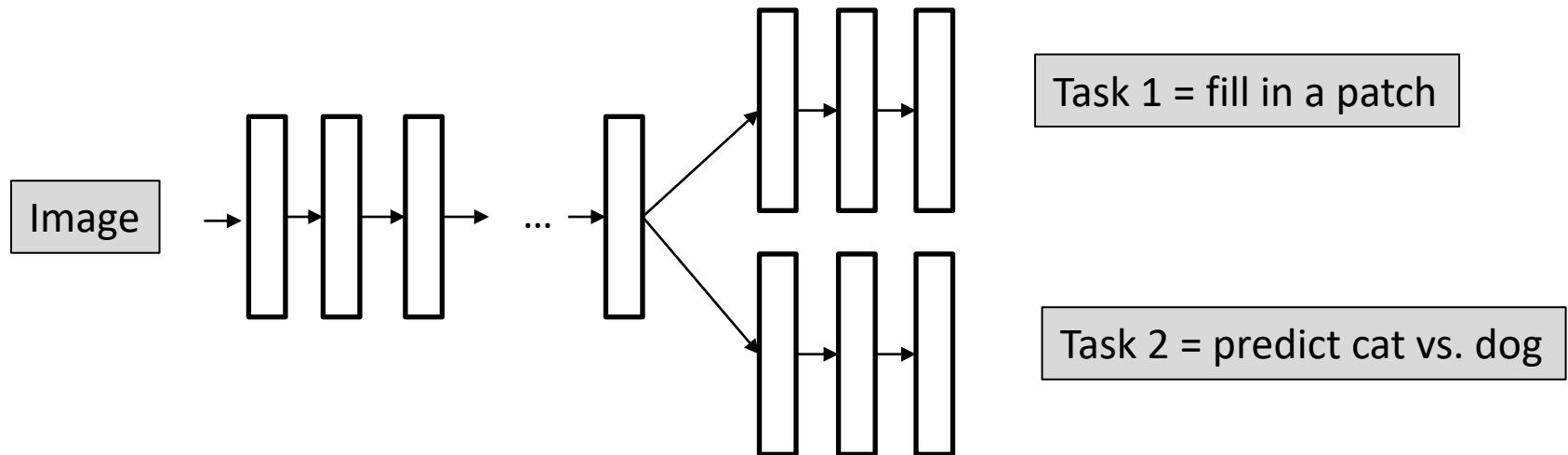
Google BERT stands for **Bidirectional Encoder Representations from Transformers** and is an update to the core search algorithm aimed at improving the language understanding capabilities of Google.

BERT is one of the biggest updates that Google has made since RankBrain in 2015 and has proven successful in comprehending the intent of the searcher behind a search query.

# Downstream Tasks - NLP (BERT Revolution)

Rank	Name	Model	Link	URL	Score	CoLA	SST-2	MRPC	STS-B	QQP	MNLI-m	MNLI-mm	QNLI	RTE	WNLI
1	T5 Team - Google	T5		90.3	71.6	97.5	92.8/90.4	93.1/92.8	75.1/90.6	92.2	91.9	96.9	92.8	94.5	
2	ERNIE Team - Baidu	ERNIE		90.0	72.2	97.5	93.0/90.7	92.9/92.5	75.2/90.8	91.2	90.8	96.0	90.9	94.5	
3	Microsoft D365 AI & MSR AI & GATECHMT-DNN-SMART			89.9	69.5	97.5	93.7/91.6	92.9/92.5	73.9/90.2	91.0	90.8	99.2	89.7	94.5	
+ 4	王玮	ALICE v2 large ensemble (Alibaba DAMO NLP)		89.7	73.2	97.1	93.9/91.9	93.0/92.5	74.8/91.0	90.8	90.6	95.9	87.4	94.5	
+ 5	Microsoft D365 AI & UMD	FreeLB-RoBERTa (ensemble)		88.4	68.0	96.8	93.1/90.8	92.3/92.1	74.8/90.3	91.1	90.7	95.6	88.7	89.0	
6	Junjie Yang	HIRE-RoBERTa		88.3	68.6	97.1	93.0/90.7	92.4/92.0	74.3/90.2	90.7	90.4	95.5	87.9	89.0	
7	Facebook AI	RoBERTa		88.1	67.8	96.7	92.3/89.8	92.2/91.9	74.3/90.2	90.8	90.2	95.4	88.2	89.0	
+ 8	Microsoft D365 AI & MSR AI	MT-DNN-ensemble		87.6	68.4	96.5	92.7/90.3	91.1/90.7	73.7/89.9	87.9	87.4	96.0	86.3	89.0	
9	GLUE Human Baselines	GLUE Human Baselines		87.1	66.4	97.8	86.3/80.8	92.7/92.6	59.5/80.4	92.0	92.8	91.2	93.6	95.9	
10	Stanford Hazy Research	Snorkel MeTaL		83.2	63.8	96.2	91.5/88.5	90.1/89.7	73.1/89.9	87.6	87.2	93.9	80.9	65.1	

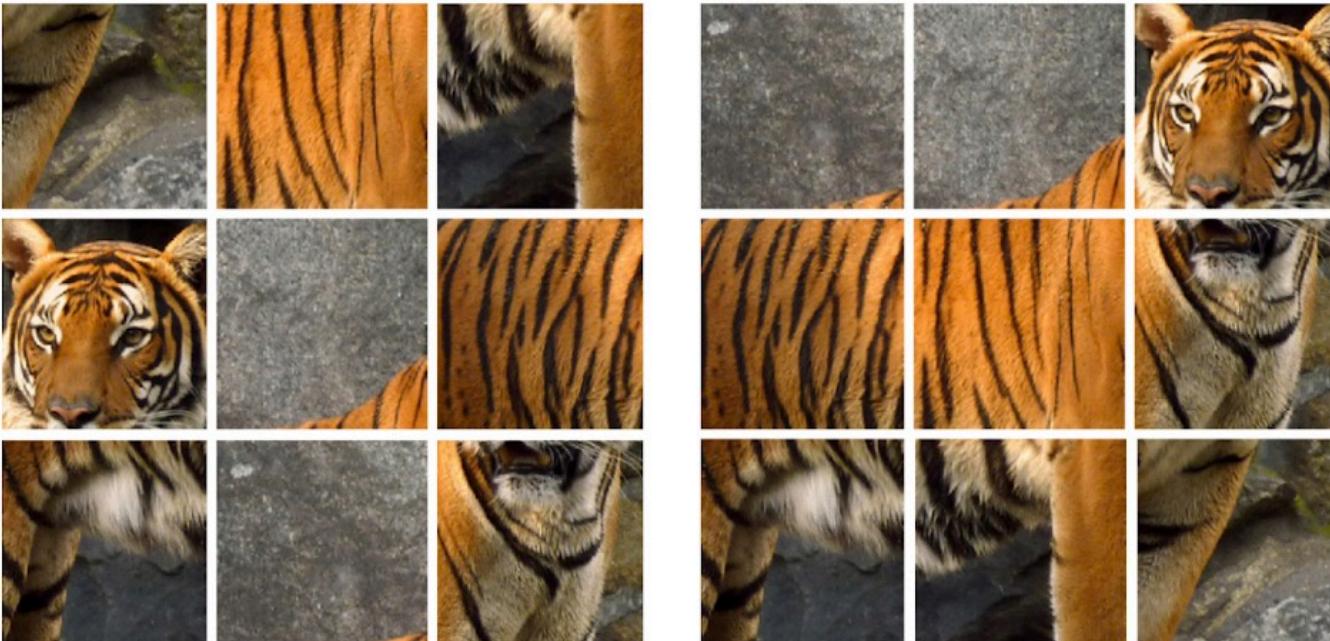
# Unsupervised Learning in Vision



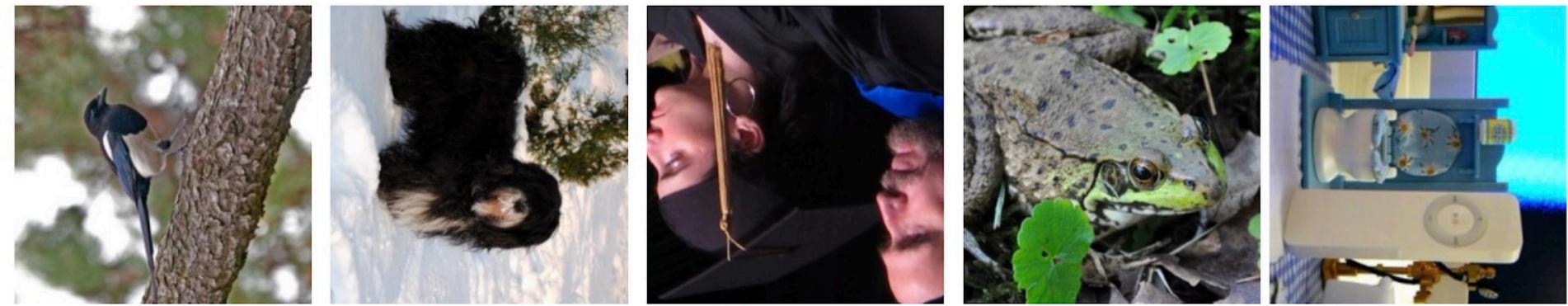
# Predict Missing Patch



# Solving Jigsaw Puzzles



# Rotation Prediction



90° rotation

270° rotation

180° rotation

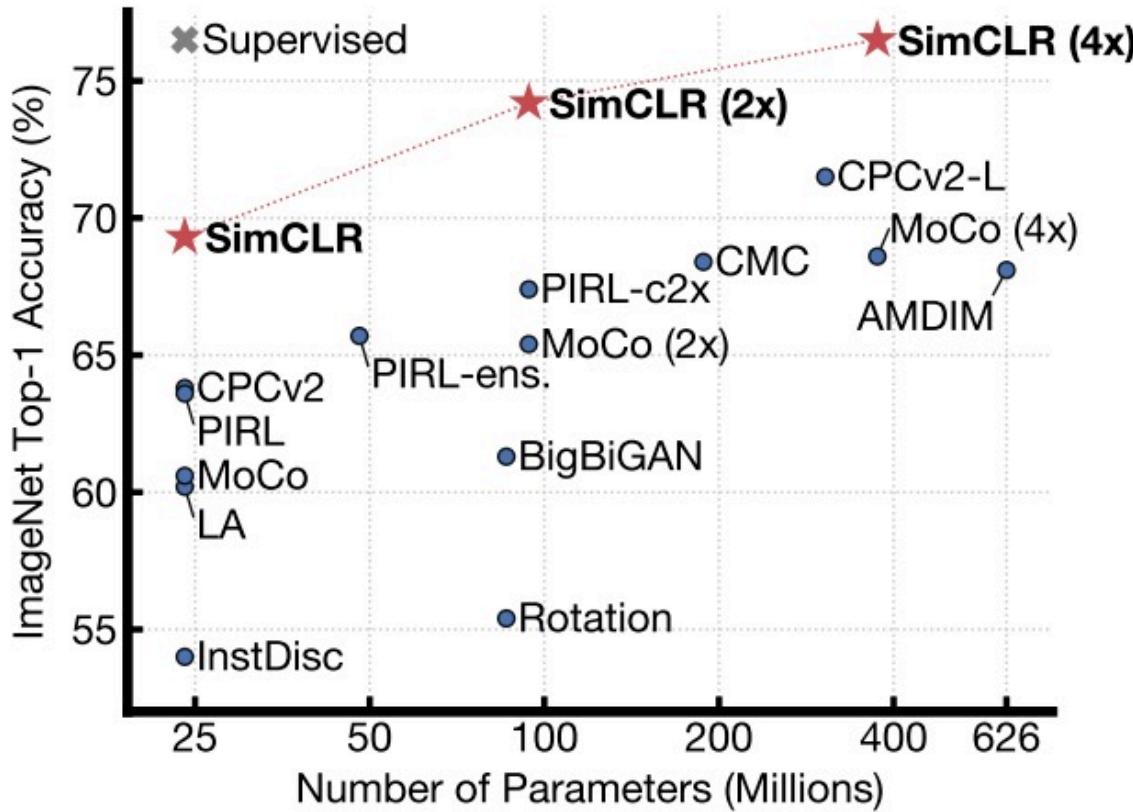
0° rotation

270° rotation

# SimCLR and MoCo



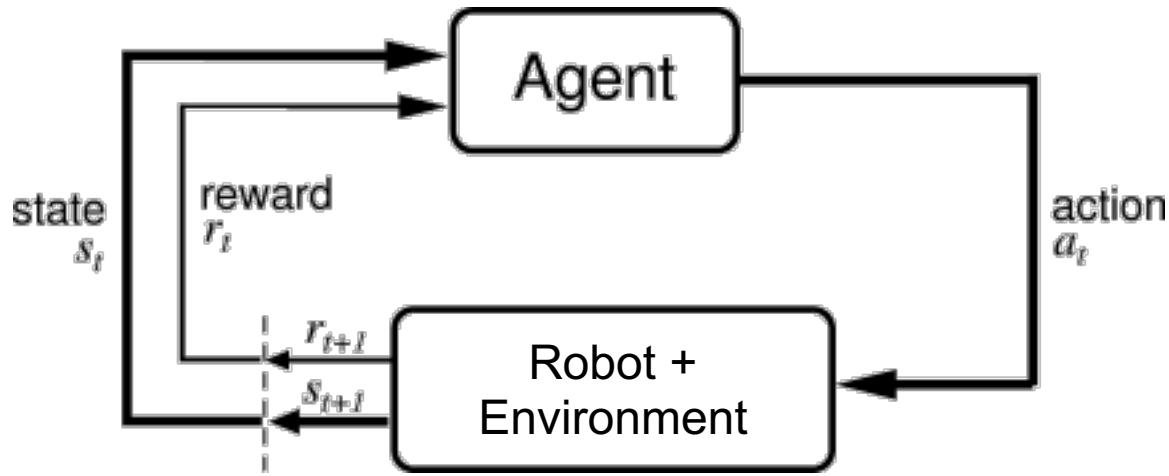
# SimCLR + linear classifier



# Many Exciting Directions in AI

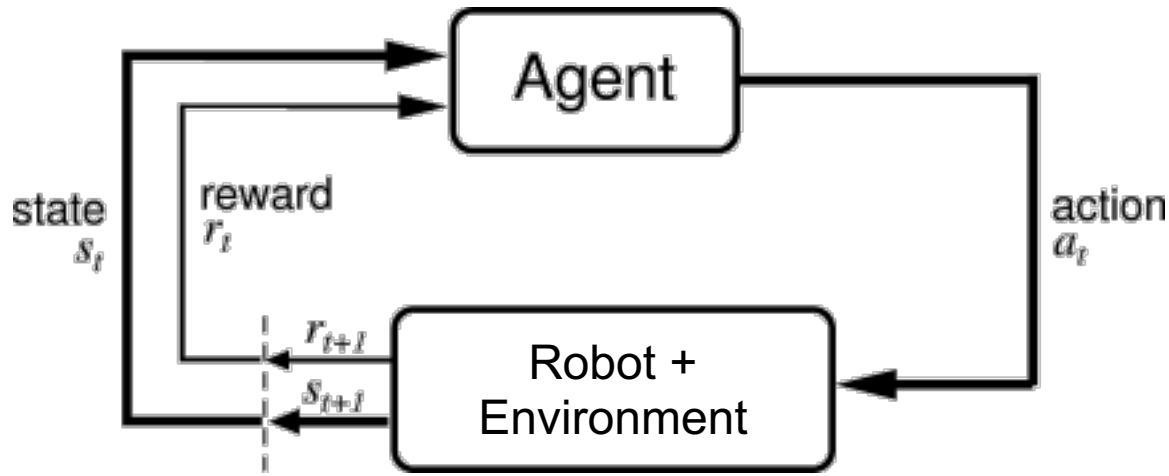
- Unsupervised Learning
- ***Reinforcement Learning***
- Unsupervised RL
- Meta-Reinforcement Learning
- Few-Shot Imitation
- Domain Randomization
- DL for Science and Engineering
- Mitigating Bias
- Multi-modal Learning
- Architecture Search
- Value Alignment
- Scaling Laws
- Human-in-the-Loop
- Explainability

# Reinforcement Learning (RL)



$$\max_{\theta} \mathbb{E} \left[ \sum_{t=0}^H R(s_t) | \pi_{\theta} \right]$$

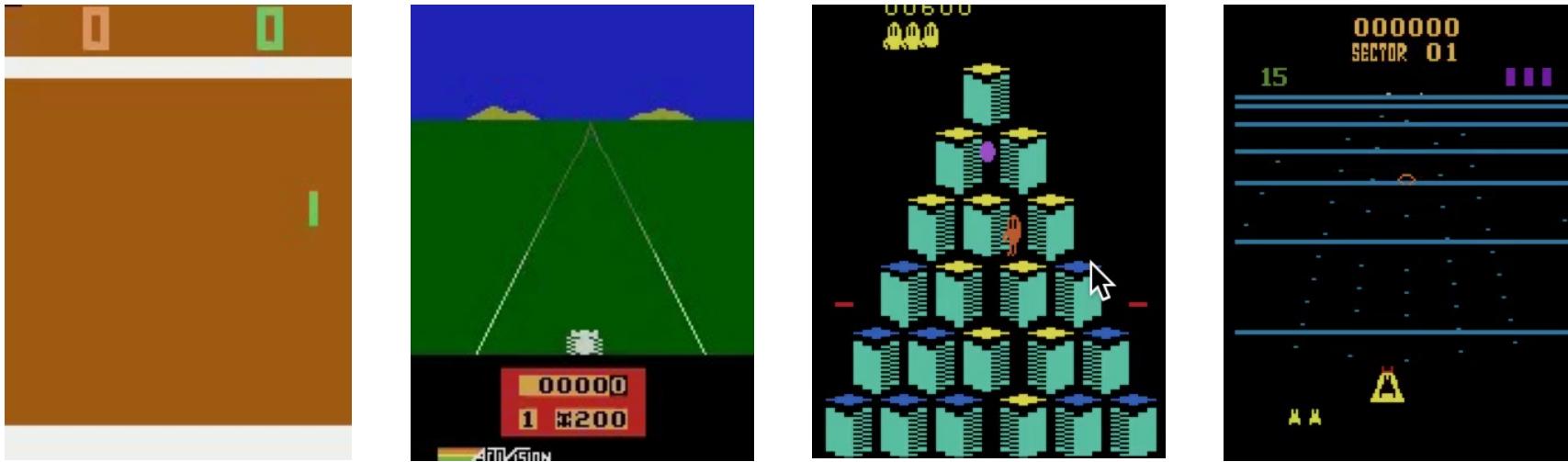
# Reinforcement Learning (RL)



- Compared to supervised learning, additional challenges:
  - Credit assignment
  - Stability
  - Exploration

$$\max_{\theta} \mathbb{E}\left[\sum_{t=0}^H R(s_t) | \pi_{\theta}\right]$$

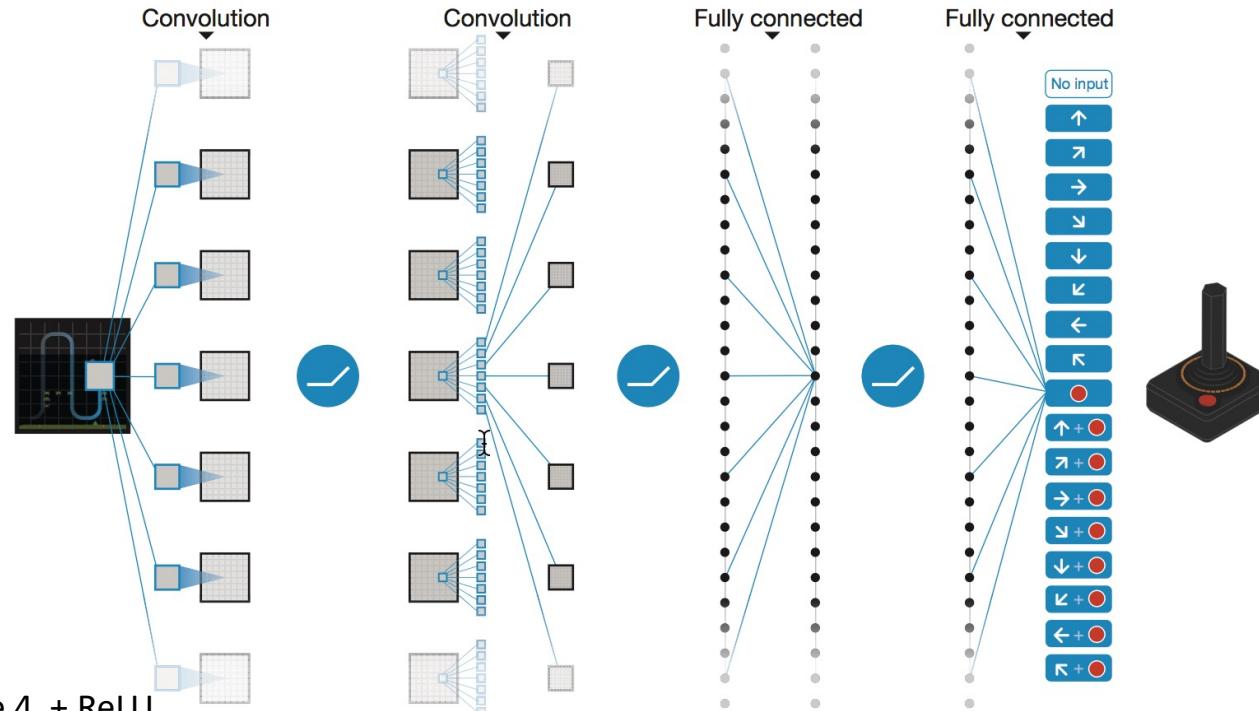
# Deep RL: Atari



DQN Mnih et al, NIPS 2013 / Nature 2015

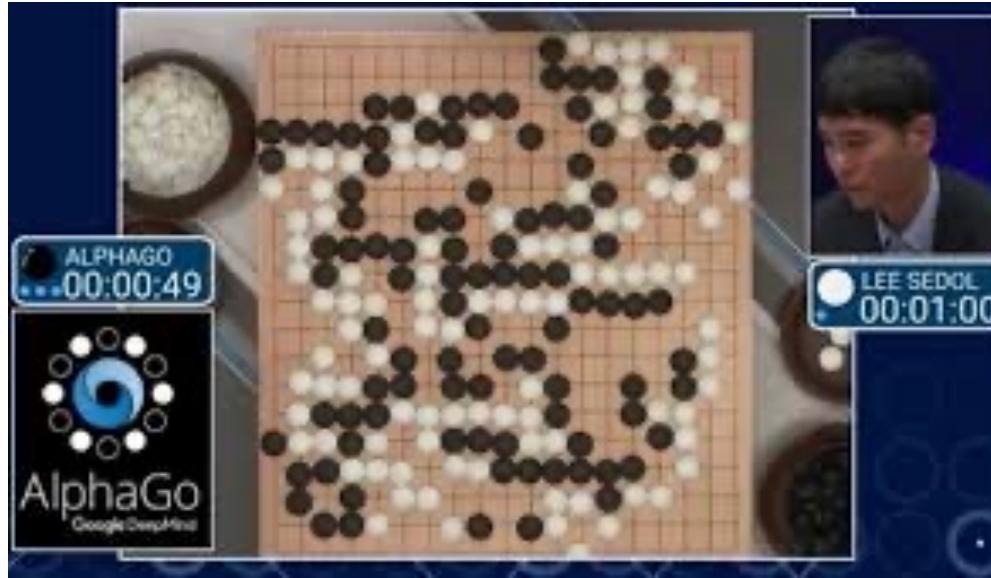
MCTS Guo et al, NIPS 2014; TRPO Schulman, Levine, Moritz, Jordan, Abbeel, ICML 2015; A3C Mnih et al, ICML 2016; Dueling DQN Wang et al ICML 2016; Double DQN van Hasselt et al, AAAI 2016; Prioritized Experience Replay Schaul et al, ICLR 2016; Bootstrapped DQN Osband et al, 2016; Q-Ensembles Chen et al, 2017; Rainbow Hessel et al, 2017; Accelerated Stooke and Abbeel, 2018; ...

# Deep Q-Network (DQN): From Pixels to Joystick Commands



[Source: Mnih et al., Nature 2015 (DeepMind) ]

# Deep RL: Go



**AlphaGo** Silver et al, Nature 2015

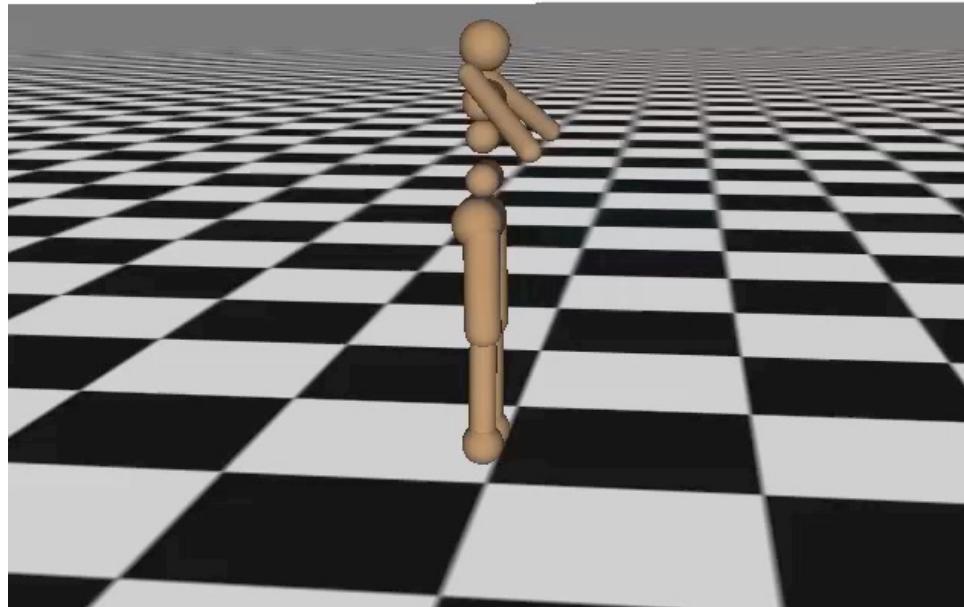
**AlphaGoZero** Silver et al, Nature 2017

**AlphaZero** Silver et al, 2017

Tian et al, 2016; Maddison et al, 2014; Clark et al, 2015

# Deep RL: Robot Locomotion

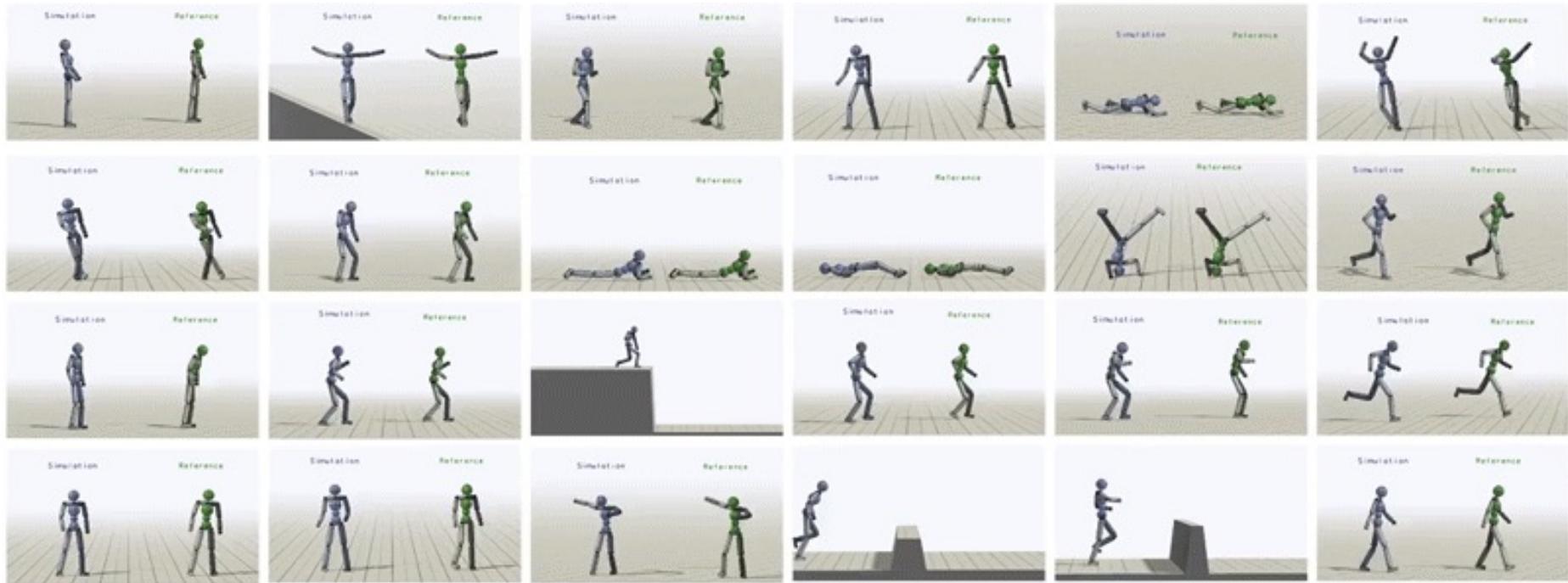
Iteration 0



^ **TRPO** Schulman et al, 2015 + **GAE** Schulman et al, 2016

See also: **DDPG** Lillicrap et al 2015; **SVG** Heess et al, 2015; **Q-Prop** Gu et al, 2016; **Scaling up ES** Salimans et al, 2017; **PPO** Schulman et al, 2017; **Parkour** Heess et al, 2017;

# Deep RL: Robot Locomotion



# Deep RL Success: Dynamic Animation



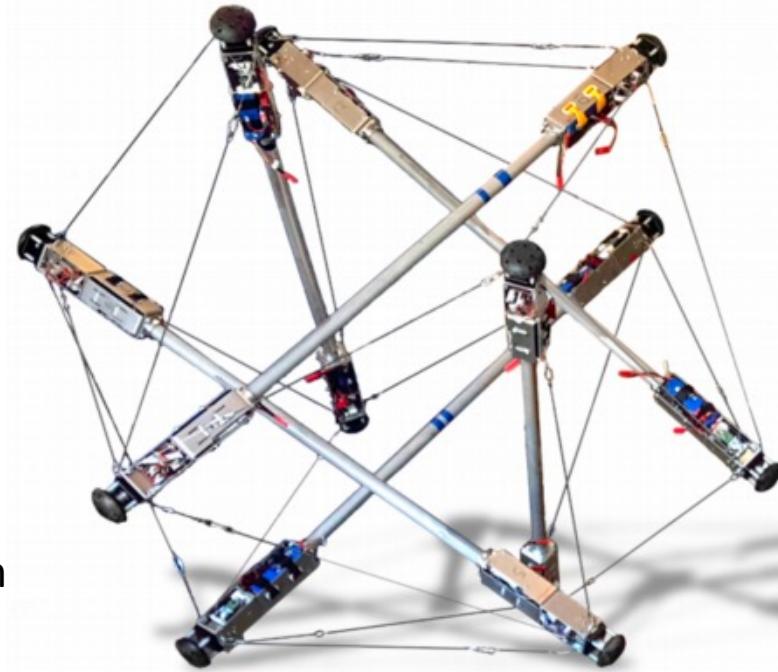
# BRETT: Berkeley Robot for the Elimination of Tedious Tasks

---

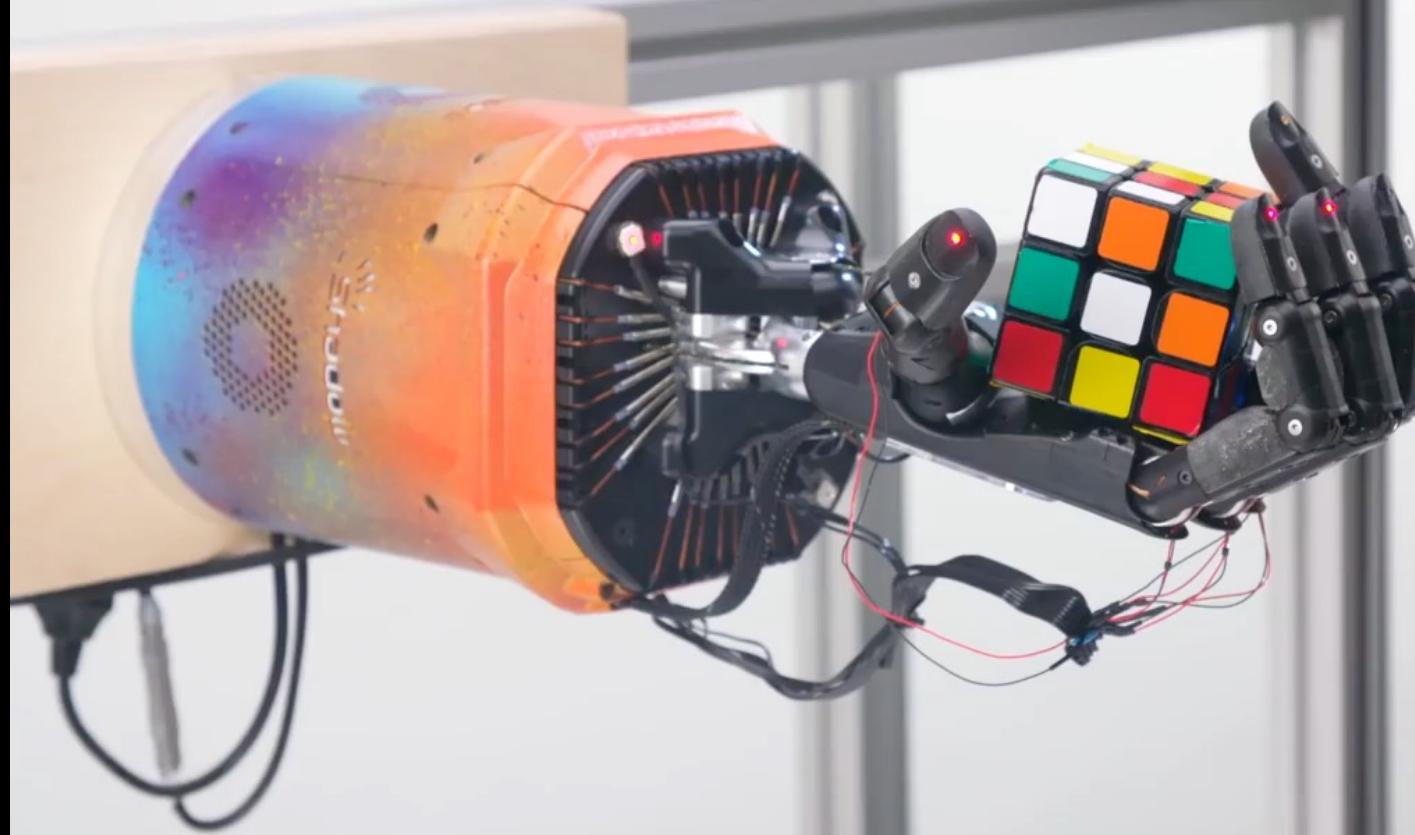


# Tensegrity Robotics: NASA SuperBall

- Rigid rods connected by elastic cables
- Controlled by motors that extend / contract cables
- Properties:
  - Lightweight
  - Low cost
  - Capable of withstanding significant impact
- NASA investigates them for space exploration
- Major challenge: control







# Covariant



Press

January 29, 2020

A Warehouse Robot Learns to Sort Out  
the Tricky Stuff



Press

January 29, 2020

AI Helps Warehouse Robots Pick Up  
New Tricks



Press

May 6, 2020

Logistics AI Startup Covariant Reaps  
\$40 Million in Funding Round

*This is the entry point of AI Robotics into the *real* world*

The New York Times  
29-January-2020

## Autonomous Order Picking



Knapp Pick-It-Easy powered by Covariant Brain

# Many Exciting Directions in AI

- Unsupervised Learning
- Reinforcement Learning
- *Unsupervised RL*
- Meta-Reinforcement Learning
- Few-Shot Imitation
- Domain Randomization
- DL for Science and Engineering
- Mitigating Bias
- Multi-modal Learning
- Architecture Search
- Value Alignment
- Scaling Laws
- Human-in-the-Loop
- Explainability

# Mastery? Yes

**Deep RL (DQN)**

Score: 18.9

vs.

**Human**

Score: 9.3

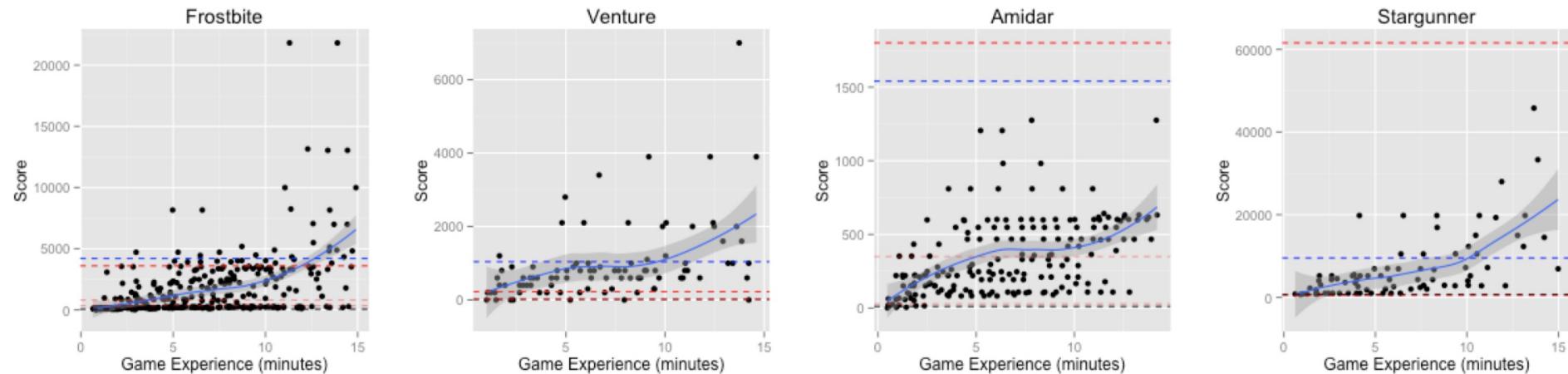


# How fast is the learning itself?

---

# Humans vs. DDQN

Humans after 15 minutes tend to outperform DDQN after 115 hours



Black dots: human play

Blue curve: mean of human play

Blue dashed line: 'expert' human play

Red dashed lines:

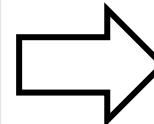
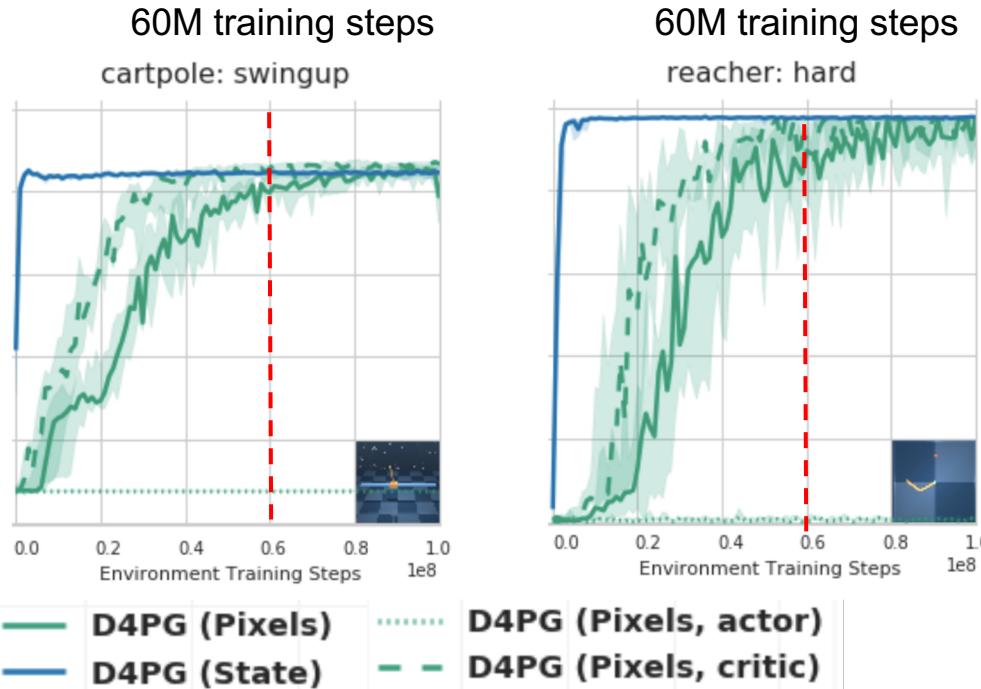
DDQN after 10, 25, 200M frames  
(~ 46, 115, 920 hours)

# How to bridge this gap?

---

# Can visual RL achieve same data-efficiency as RL on state?

- State-based D4PG (blue) vs pixel-based D4PG (green)

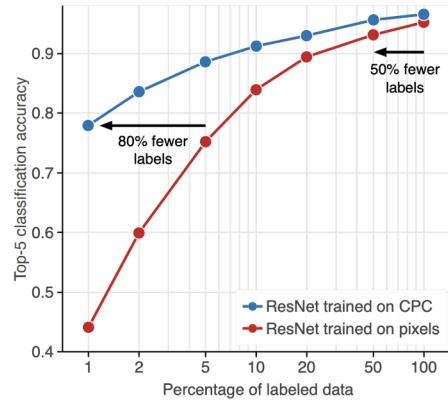


Pixel-based needs > 50M more training steps than state-based to solve same tasks

[Tassa et al., 2018] Tassa, Y., Doron, Y., Muldal, A., Erez, T., Li, Y., Casas, D.D.L., Budden, D., Abdolmaleki, A., Merel, J., Lefrancq, A. and Lillicrap, T. [DeepMind Control Suite](#), arxiv:1801.00690, 2018.

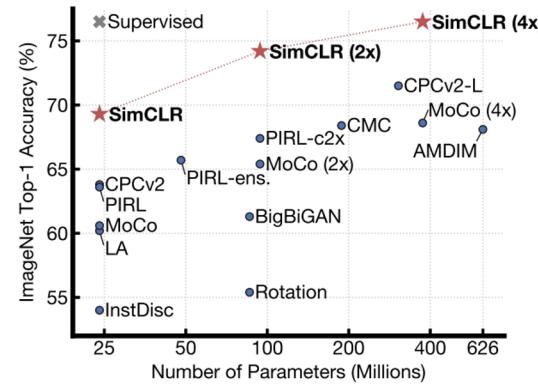
# Contrastive learning: SOTA in computer vision

CPCv2 **top-5** ImageNet accuracy as function of labels



[Henaff, Srinivas et al., 2019]

SimCLR **top-1** ImageNet accuracy as function of # of parameters



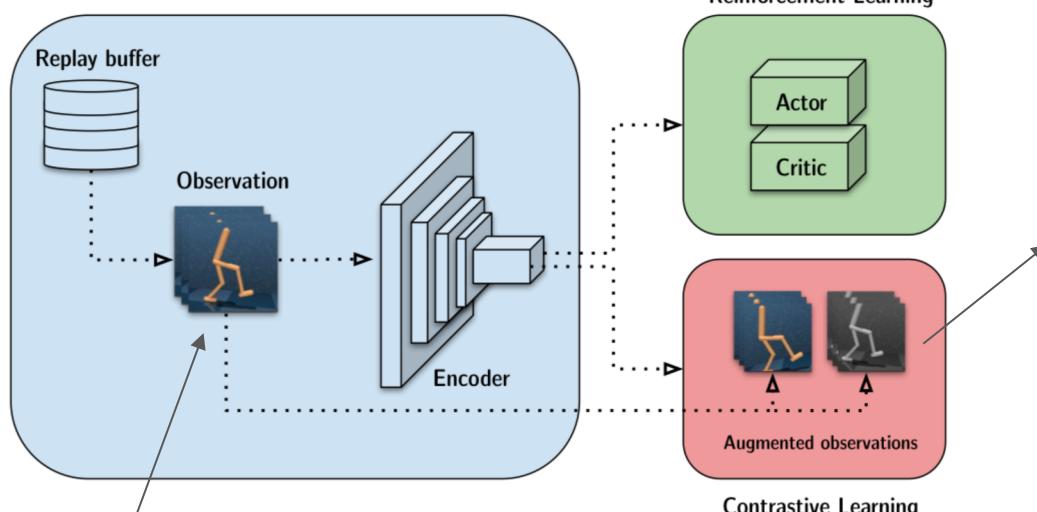
[Chen et al., 2020]

[Henaff et al., 2019] Olivier J. Hénaff, Aravind Srinivas, Jeffrey De Fauw, Ali Razavi, Carl Doersch, S. M. Ali Eslami, Aaron van den Oord [Data-Efficient Image Recognition with Contrastive Coding](#) arxiv:1905.09272, 2019.

[Chen et al., 2020] Chen, T., Kornblith, S., Norouzi, M. and Hinton, G. [A Simple Framework for Contrastive Learning of Visual Representations](#) arxiv:2002.05709, 2020.

# Contrastive + RL

CURL

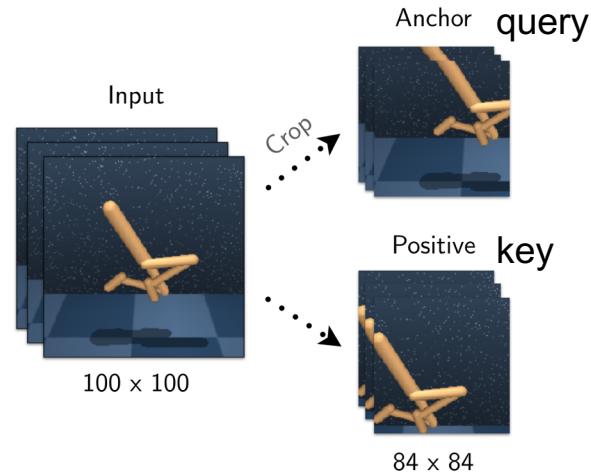


Observations are  
stacked frames

Need to define:

1. query / key pairs
2. similarity measure
3. architecture

# 1. Query / key pairs: random crop



## 2. Bilinear inner product with learned weight matrix

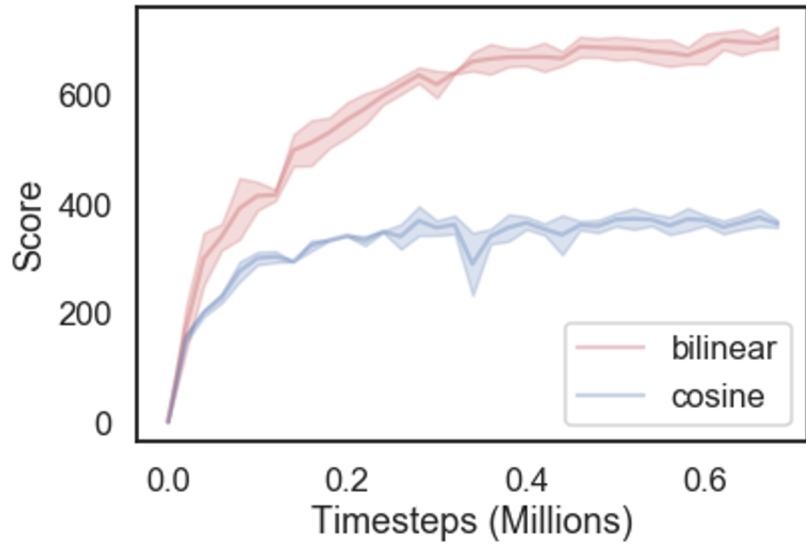
logits

$$\begin{bmatrix} q_0^T W k_0 & q_0^T W k_1 & \dots & q_0^T W k_j \\ q_1^T W k_0 & q_1^T W k_1 & \dots & q_1^T W k_j \\ \vdots & \vdots & \ddots & \vdots \\ q_j^T W k_0 & q_j^T W k_1 & \dots & q_j^T W k_j \end{bmatrix} \begin{bmatrix} 1 & 0 & \dots & 0 \\ 0 & 1 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 1 \end{bmatrix}$$

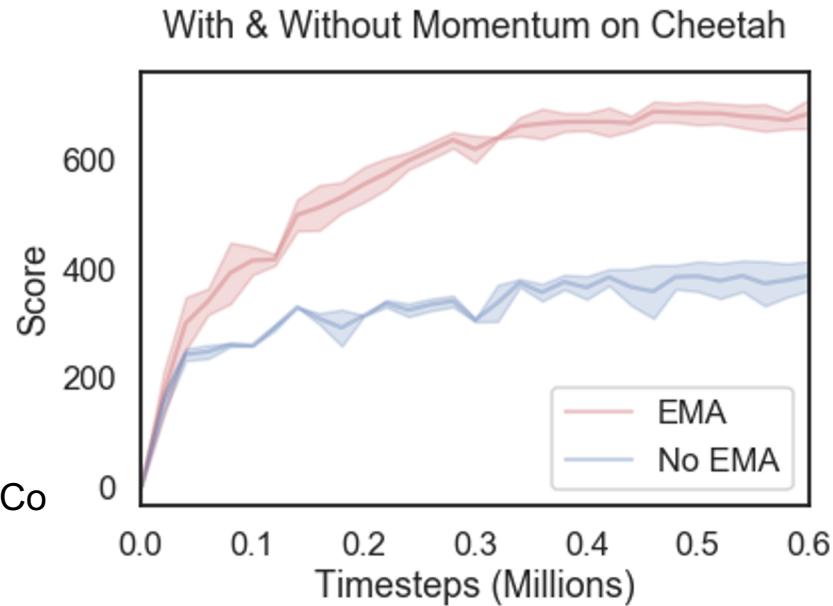
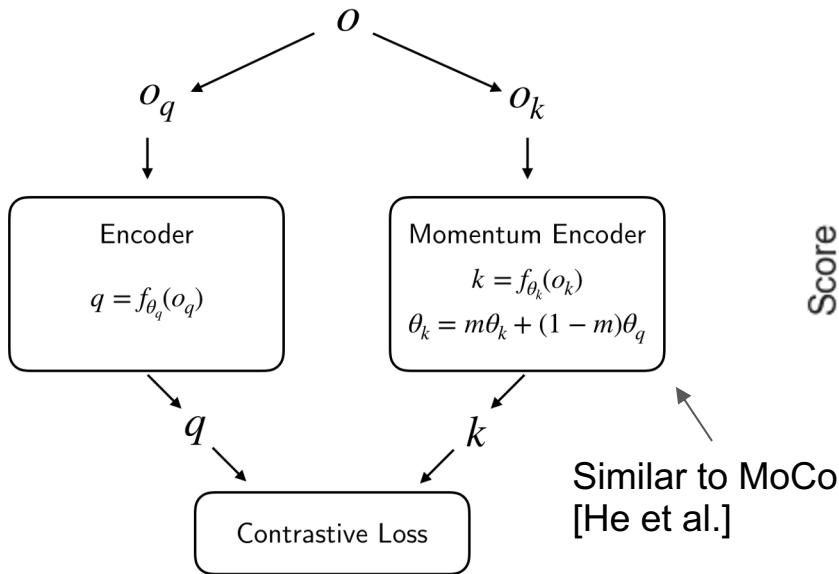
labels

$$\mathcal{L}_q = \frac{\exp(q^T W k_+)}{\exp(q^T W k_+) + \sum_{i=0}^{K-1} \exp(q^T W k_i)}$$

Comparing Similarity Measures on Cheetah

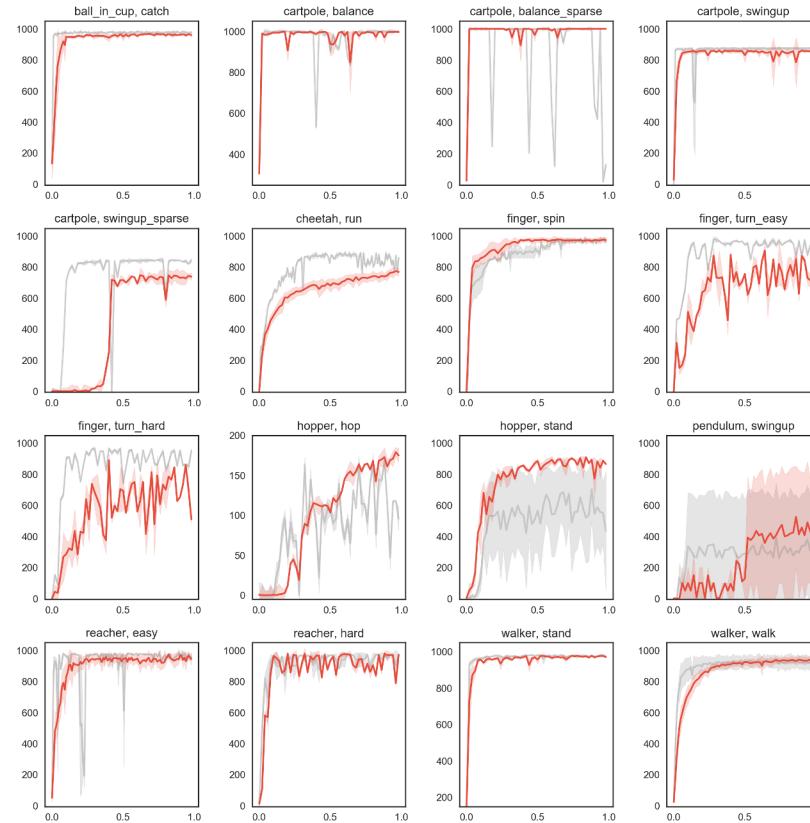


### 3. Keys encoded with momentum

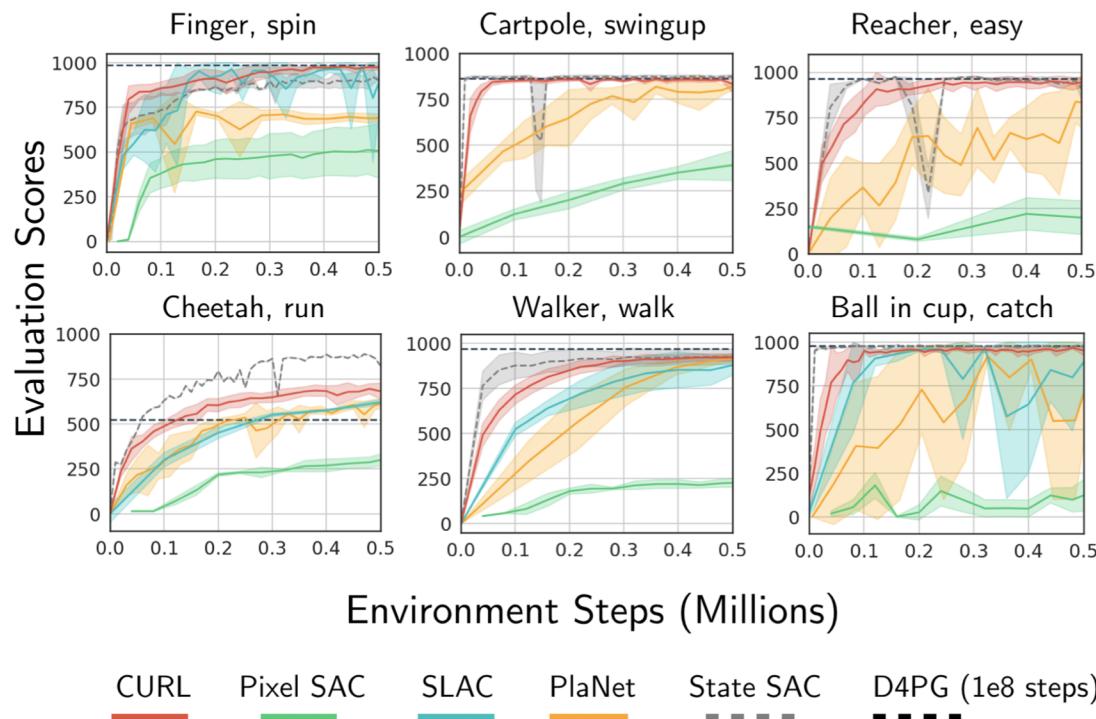


# CURL from pixels matches state-based SAC

**GRAY:** SAC State  
**RED:** CURL



# CURL Comparison: DeepMind Control Suite



# CURL Comparison: Atari

Atari performance benchmarked at 100K frames

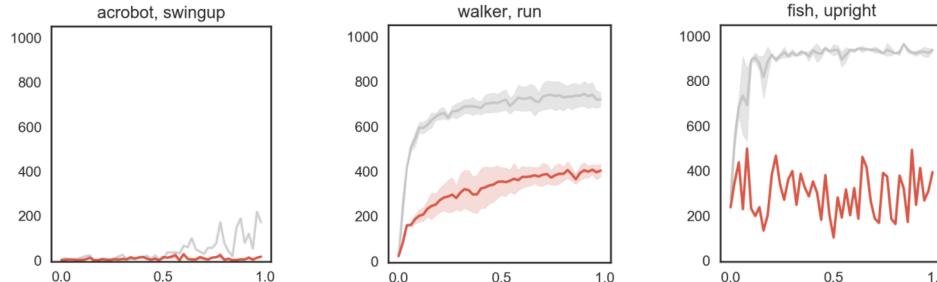
100K STEP SCORES	CURL RAINBOW	SIMPLE	RAINBOW	HUMAN	RANDOM
ALIEN	<b>1148.2</b>	616.9	318.7	6875	184.8
AMIDAR	<b>232</b>	74.3	32.5	1676	11.8
ASSAULT	473	<b>527.2</b>	231	1496	248.8
BATTLEZONE	<b>11208</b>	4031.2	3285.71	37800	2895
FREEWAY	<b>27</b>	16.7	0	29.6	0
FROSTBITE	<b>924</b>	236.9	60.2	4335	74
JAMESBOND	<b>400</b>	100.5	47.4	406.7	29.2
QBERT	<b>1352</b>	1288.8	123.46	13455	166.1
SEAQUEST	408	<b>683.3</b>	131.69	20182	61.1

# Hard Environments

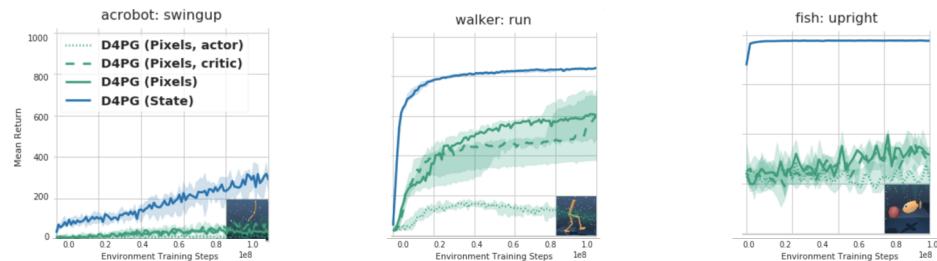
GRAY: SAC State

RED: CURL

Environment training steps 1 = 1M

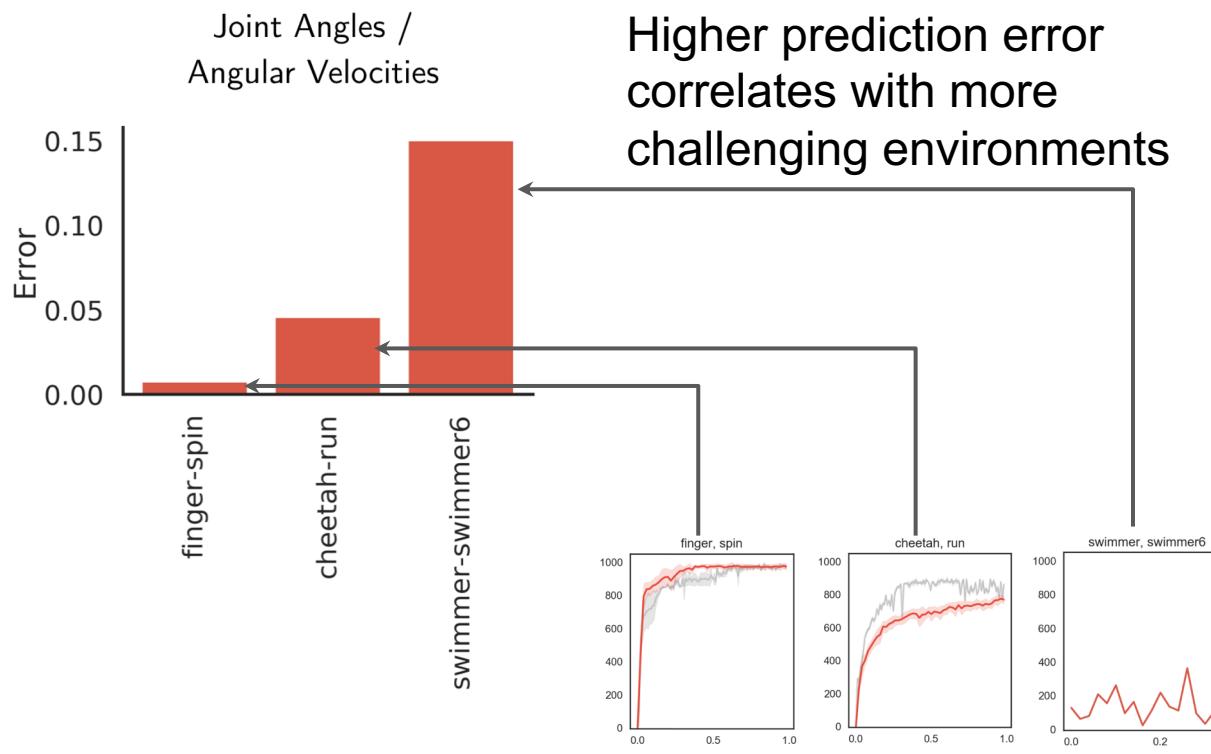


Environment training steps 1 = 100M



**Observation:**  
similar struggle  
asymptotically for  
pixel-based Deep RL

# Predicting state from pixels



# Many Exciting Directions in AI

- Unsupervised Learning
- Reinforcement Learning
- Unsupervised RL
- ***Meta-Reinforcement Learning***
- Few-Shot Imitation
- Domain Randomization
- DL for Science and Engineering
- Mitigating Bias
- Multi-modal Learning
- Architecture Search
- Value Alignment
- Scaling Laws
- Human-in-the-Loop
- Explainability

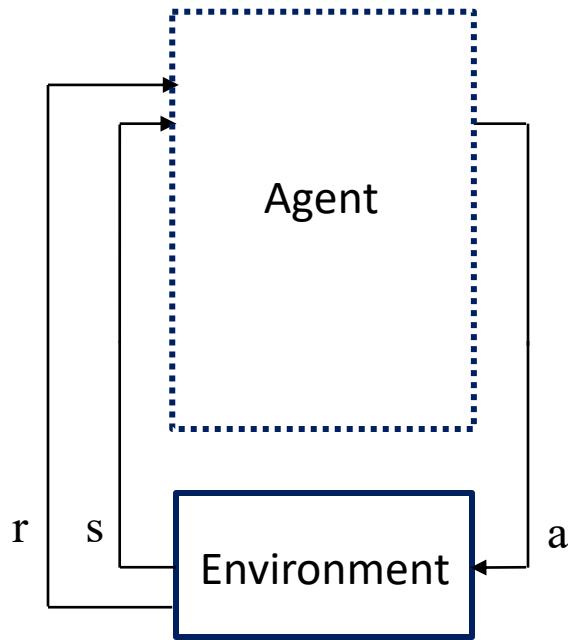
# Starting Observations

---

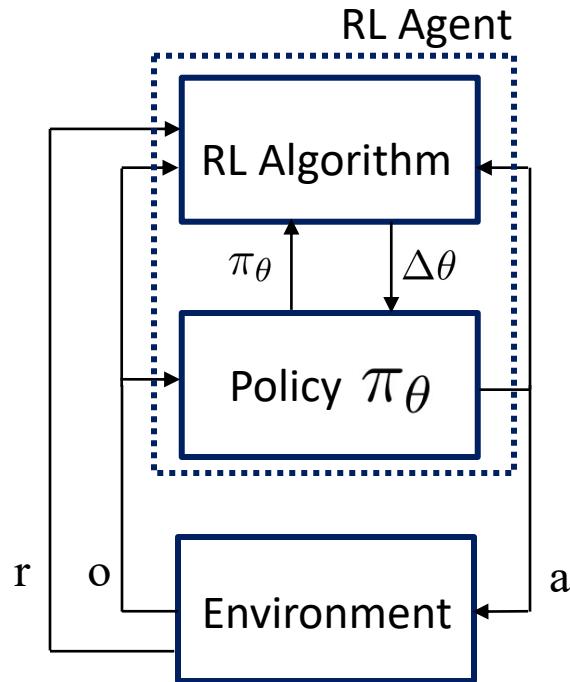
- TRPO, DQN, A3C, DDPG, PPO, Rainbow, ... are fully general RL algorithms
  - i.e., for any environment that can be mathematically defined, these algorithms are equally applicable
- Environments encountered in real world
  - = tiny, tiny subset of all environments that could be defined (e.g. they all satisfy our universe's physics)

**Can we develop “fast” RL algorithms that take advantage of this?**

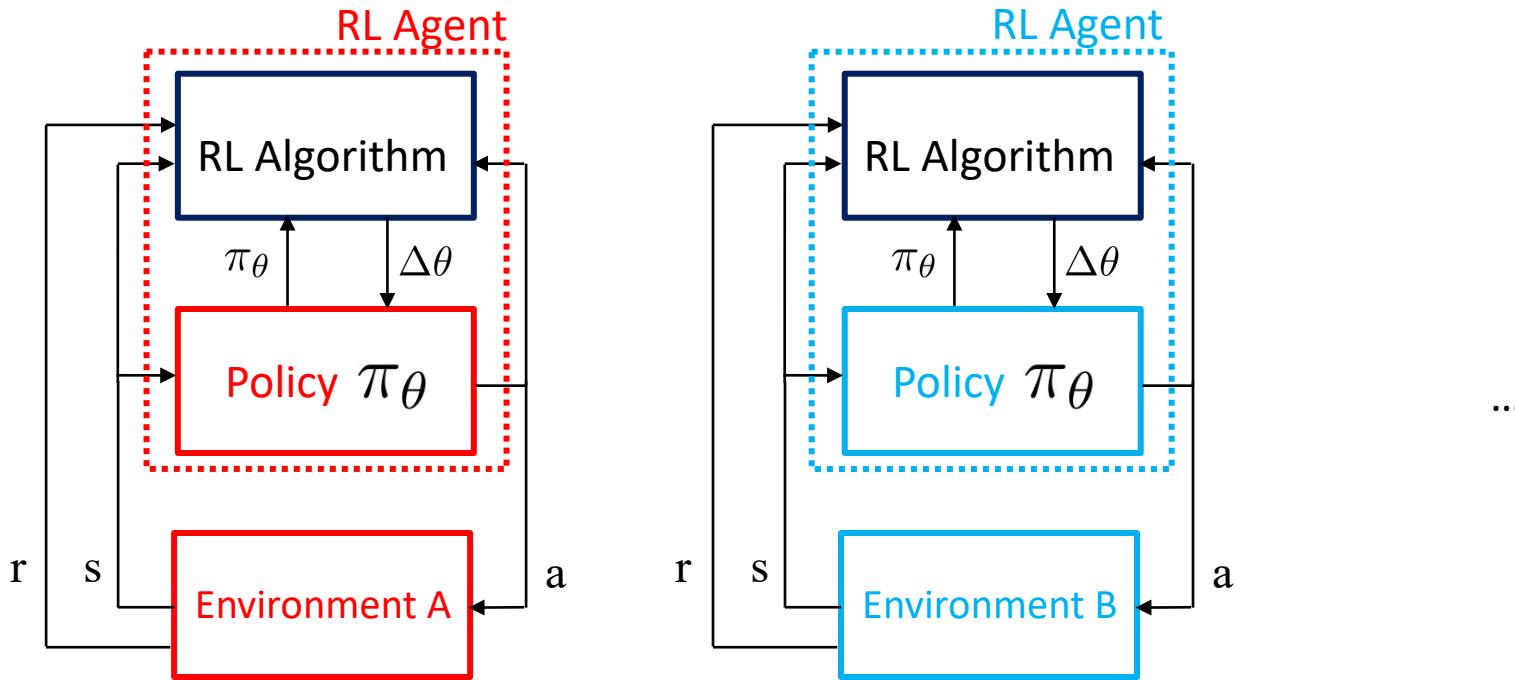
# Reinforcement Learning



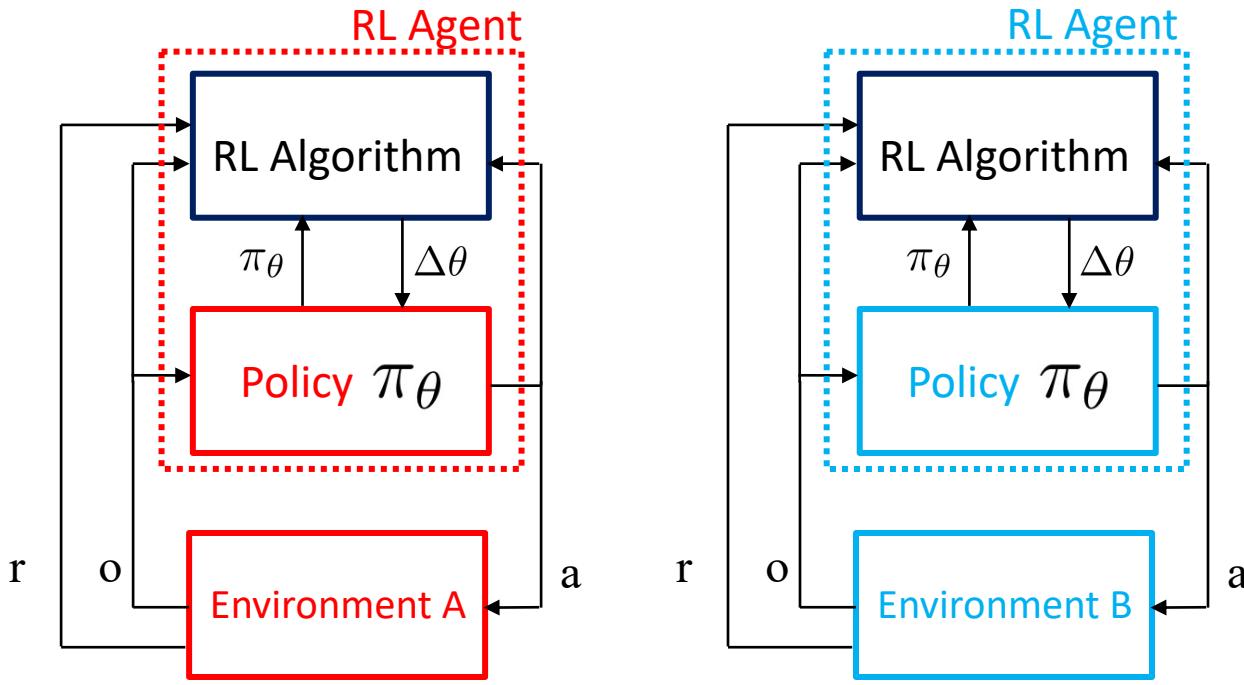
# Reinforcement Learning



# Reinforcement Learning



# Reinforcement Learning



## Traditional RL research:

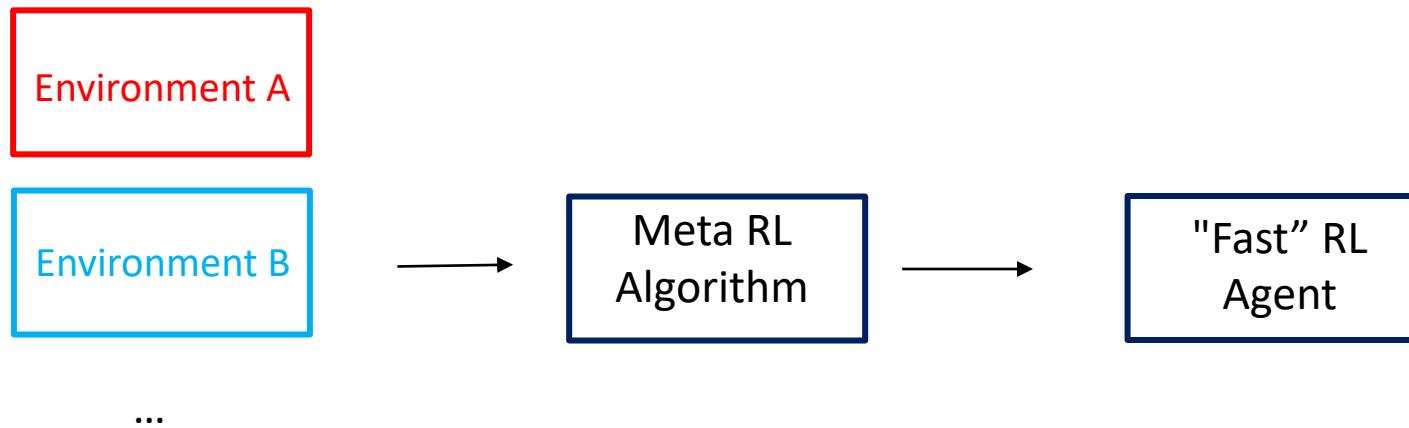
- Human experts develop the RL algorithm
- After many years, still no RL algorithms nearly as good as humans...

## Alternative:

- Could we learn a better RL algorithm?
- Or even learn a better entire agent?

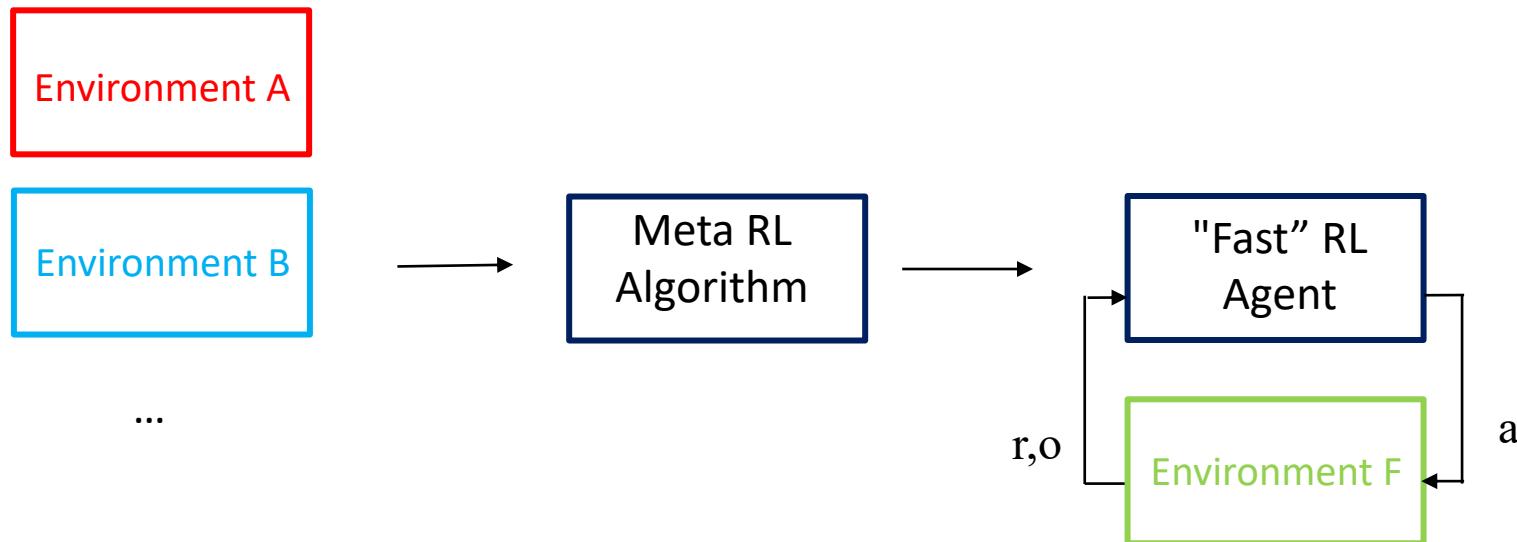
# Meta-Reinforcement Learning

Meta-training environments



# Meta-Reinforcement Learning

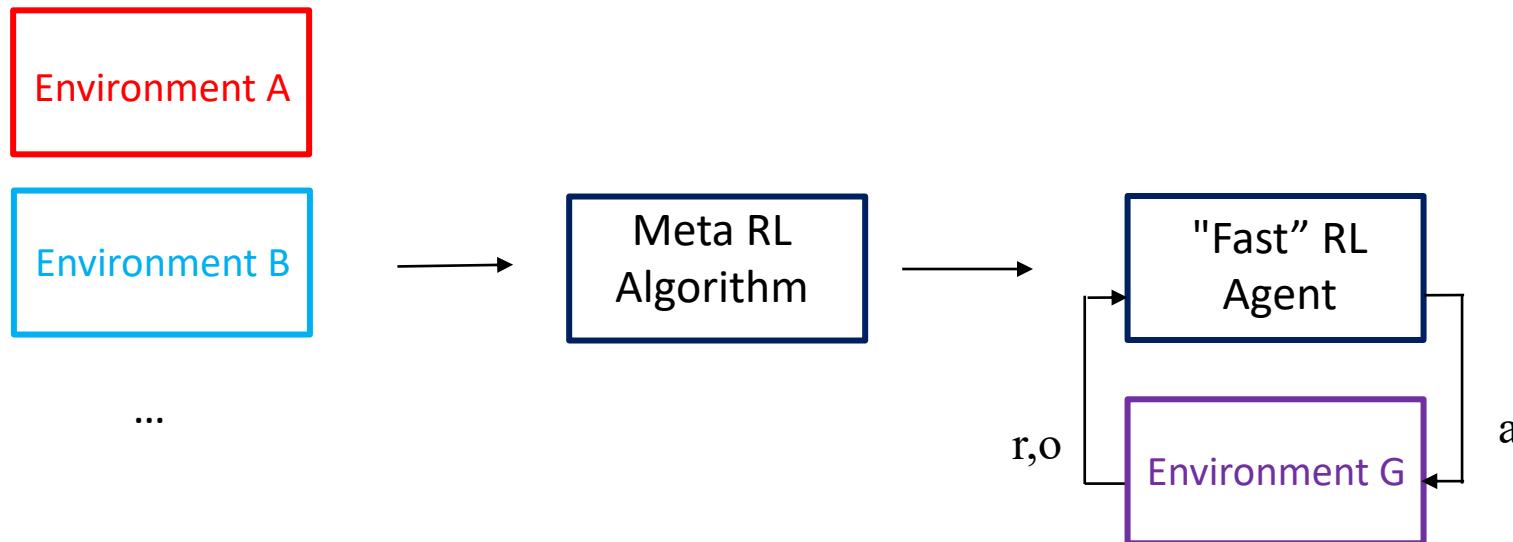
Meta-training environments



Testing environments

# Meta-Reinforcement Learning

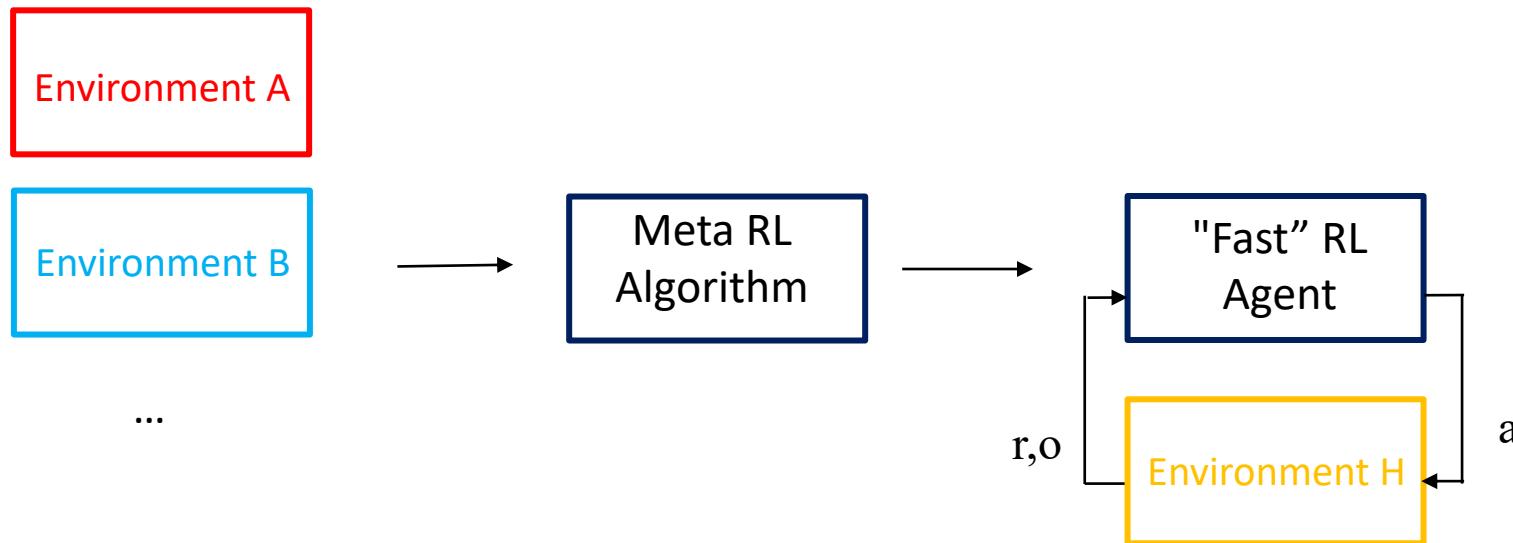
Meta-training environments



Testing environments

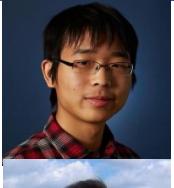
# Meta-Reinforcement Learning

Meta-training environments



Testing environments

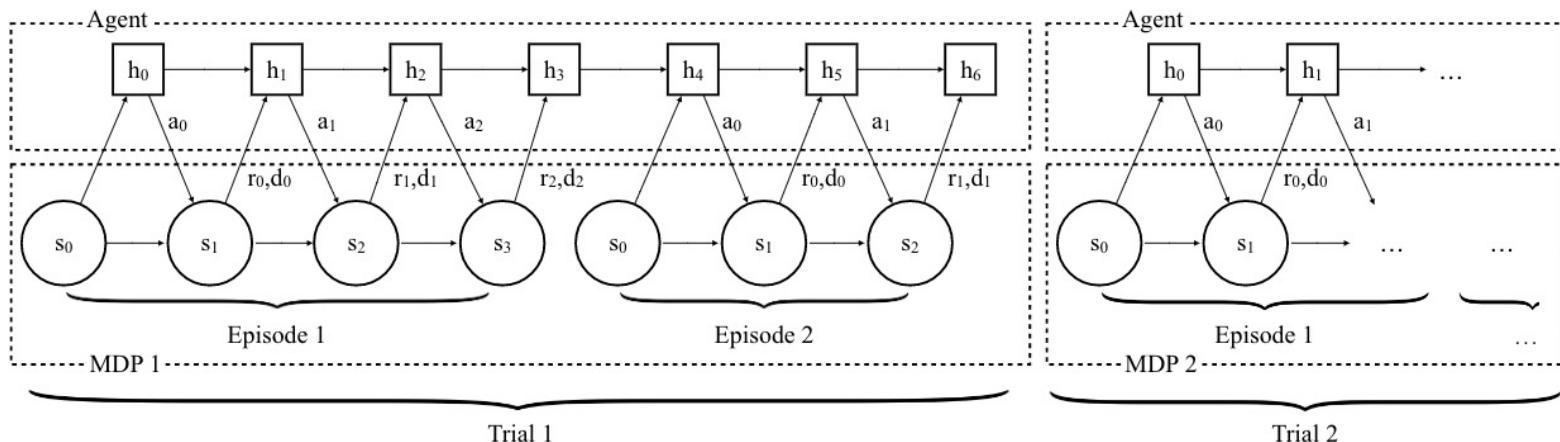
# Formalizing Learning to Reinforcement Learn



$$\max_{\theta} \mathbb{E}_M \mathbb{E}_{\tau_M^{(k)}} \left[ \sum_{k=1}^K R(\tau_M^{(k)}) \mid \text{RLagent}_{\theta} \right]$$

$M$  : sample environment

$\tau_M^{(k)}$  :  $k$ 'th episode in environment  $M$



# Formalizing Learning to Reinforcement Learn



$$\max_{\theta} \mathbb{E}_M \mathbb{E}_{\tau_M^{(k)}} \left[ \sum_{k=1}^K R(\tau_M^{(k)}) \mid \text{RLagent}_{\theta} \right]$$

$M$  : sample MDP

$\tau_M^{(k)}$  :  $k$ 'th trajectory in MDP  $M$

Meta-train:

$$\max_{\theta} \sum_{M \in M_{\text{train}}} \mathbb{E}_{\tau_M^{(k)}} \left[ \sum_{k=1}^K R(\tau_M^{(k)}) \mid \text{RLagent}_{\theta} \right]$$

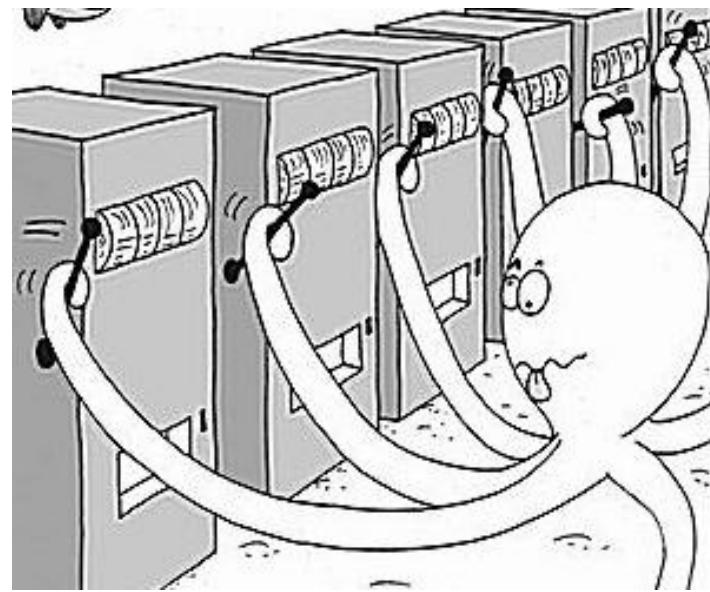
# Representing RLagent $_{\theta}$

$$\max_{\theta} \sum_{M \in M_{\text{train}}} \mathbb{E}_{\tau_M^{(k)}} \left[ \sum_{k=1}^K R(\tau_M^{(k)}) \mid \text{RLagent}_{\theta} \right]$$

- RLagent = RNN = generic computation architecture
  - different weights in the RNN means different RL algorithm and prior
  - different activations in the RNN means different current policy
  - meta-train objective can be optimized with an existing (slow) RL algorithm

# Evaluation: Multi-Armed Bandits

- Multi-Armed Bandits setting
  - Each bandit has its own distribution over pay-outs
  - Each episode = choose 1 bandit
  - Good RL agent should explore bandits sufficiently, yet also exploit the good/best ones
- Provably (asymptotically) optimal RL algorithms have been invented by humans: Gittins index, UCB1, Thompson sampling, ...



# Evaluation: Multi-Armed Bandits

Setup	Random	Gittins	TS	OTS	UCB1	$\epsilon$ -Greedy	Greedy	RL <sup>2</sup>
$n = 10, k = 5$	5.0	<b>6.6</b>	5.7	6.5	<b>6.7</b>	<b>6.6</b>	<b>6.6</b>	<b>6.7</b>
$n = 10, k = 10$	5.0	<b>6.6</b>	5.5	6.2	<b>6.7</b>	<b>6.6</b>	<b>6.6</b>	<b>6.7</b>
$n = 10, k = 50$	5.1	6.5	5.2	5.5	<b>6.6</b>	6.5	6.5	<b>6.8</b>
$n = 100, k = 5$	49.9	<b>78.3</b>	74.7	<b>77.9</b>	<b>78.0</b>	75.4	74.8	<b>78.7</b>
$n = 100, k = 10$	49.9	<b>82.8</b>	76.7	81.4	82.4	77.4	77.1	<b>83.5</b>
$n = 100, k = 50$	49.8	<b>85.2</b>	64.5	67.7	84.3	78.3	78.0	<b>84.9</b>
$n = 500, k = 5$	249.8	<b>405.8</b>	<b>402.0</b>	<b>406.7</b>	<b>405.8</b>	388.2	380.6	<b>401.6</b>
$n = 500, k = 10$	249.0	<b>437.8</b>	429.5	<b>438.9</b>	<b>437.1</b>	408.0	395.0	432.5
$n = 500, k = 50$	249.6	<b>463.7</b>	427.2	437.6	457.6	413.6	402.8	438.9

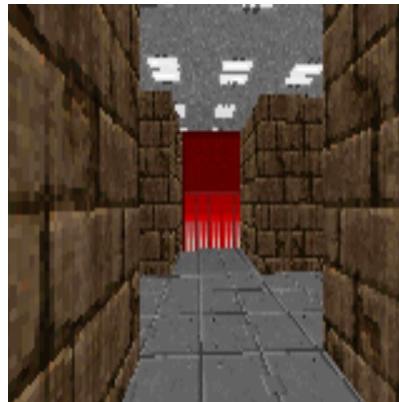
We consider Bayesian evaluation setting. Some of these prior works also have adversarial guarantees, which we don't consider here.

# Evaluation: Visual Navigation

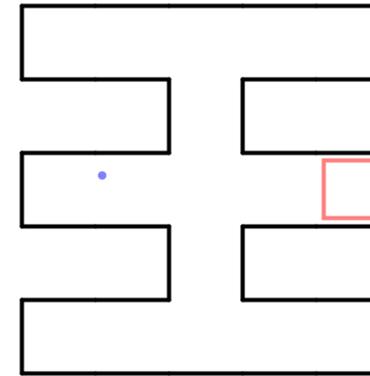
**Agent input:** current image

**Agent action:** straight / 2 degrees left / 2 degrees right

*Map just shown for our purposes, but not available to agent*



Agent's view



Maze

Related work: Mirowski, et al, 2016; Jaderberg et al, 2016; Mnih et al, 2016; Wang et al, 2016

# Evaluation: Visual Navigation

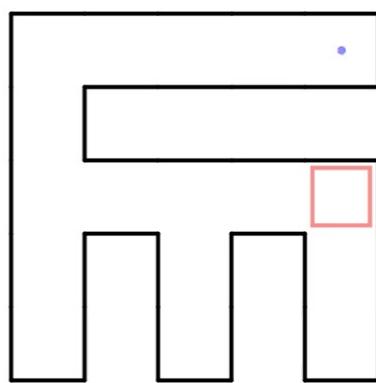
**Agent input:** current image

**Agent action:** straight / 2 degrees left / 2 degrees right

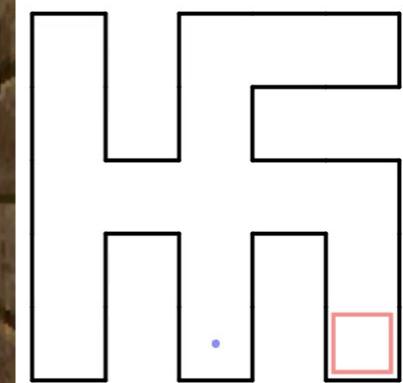
*Map just shown for our purposes, but not available to agent*



Before learning-to-learn



After learning-to-learn



# Meta-Learning Curves

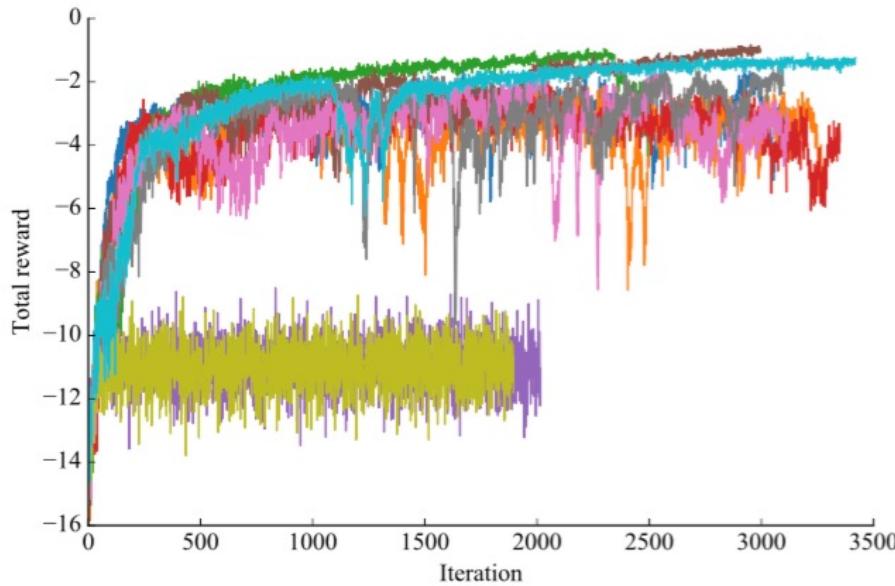


Figure 5:  $\text{RL}^2$  learning curves for visual navigation. Each curve shows a different random initialization of the RNN weights. Performance varies greatly across different initializations.

# Meta Learning for RL

## Task distribution: different environments

- Schmidhuber. Evolutionary principles in self-referential learning. (1987)
- Wiering, Schmidhuber. Solving POMDPs with Levin search and EIRA. (1996)
- Schmidhuber, Zhao, Wiering. Shifting inductive bias with success-story algorithm, adaptive Levin search, and incremental self-improvement. (MLJ 1997)
- Schmidhuber, Zhao, Schraudolph. Reinforcement learning with self-modifying policies (1998)
- Zhao, Schmidhuber. Solving a complex prisoner's dilemma with self-modifying policies. (1998)
- Schmidhuber. A general method for incremental self-improvement and multiagent learning. (1999)
- Singh, Lewis, Barto. Where do rewards come from? (2009)
- Singh, Lewis, Barto. Intrinsically Motivated Reinforcement Learning: An Evolutionary Perspective (2010)
- Niekum, Spector, Barto. Evolution of reward functions for reinforcement learning (2011)
- Duan et al., (2016) RL2: Fast Reinforcement Learning via Slow Reinforcement Learning
- Wang et al., (2016) Learning to Reinforcement Learn
- Finn et al., (2017) Model-Agnostic Meta-Learning (MAML)
- Mishra, Rohinenjad et al., (2017) Simple Neural Attentlve meta-Learner
- Frans et al., (2017) Meta-Learning Shared Hierarchies

# Meta Learning for RL

## Task distribution: different environments

- Schmidhuber. Evolutionary principles in self-referential learning. (1987)
- Wiering, Schmidhuber. Solving POMDPs with Levin search and EIRA. (1996)
- Schmidhuber, Zhao, Wiering. Shifting inductive bias with success-story algorithm, adaptive Levin search, and incremental self-improvement. (MLJ 1997)
- Schmidhuber, Zhao, Schraudolph. Reinforcement learning with self-modifying policies (1998)
- Zhao, Schmidhuber. Solving a complex prisoner's dilemma with self-modifying policies. (1998)
- Schmidhuber. A general method for incremental self-improvement and multiagent learning. (1999)
- Singh, Lewis, Barto. Where do rewards come from? (2009)
- Singh, Lewis, Barto. Intrinsically Motivated Reinforcement Learning: An Evolutionary Perspective (2010)
- Niekum, Spector, Barto. Evolution of reward functions for reinforcement learning (2011)
- Duan et al., (2016) RL2: Fast Reinforcement Learning via Slow Reinforcement Learning
- Wang et al., (2016) Learning to Reinforcement Learn
- ***Finn et al., (2017) Model-Agnostic Meta-Learning (MAML)***
- Mishra, Rohinenjad et al., (2017) Simple Neural Attentlve meta-Learner
- Frans et al., (2017) Meta-Learning Shared Hierarchies

# Many Exciting Directions in AI

- Unsupervised Learning
- Reinforcement Learning
- Unsupervised RL
- Meta-Reinforcement Learning
- ***Few-Shot Imitation***
- Domain Randomization
- DL for Science and Engineering
- Mitigating Bias
- Multi-modal Learning
- Architecture Search
- Value Alignment
- Scaling Laws
- Human-in-the-Loop
- Explainability

# Imitation Learning in Robotics



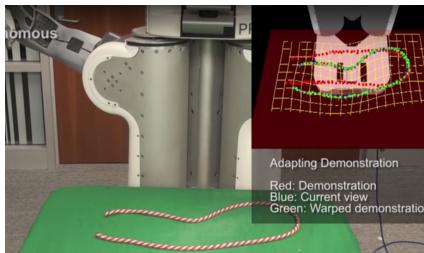
[Abbeel et al. 2008]



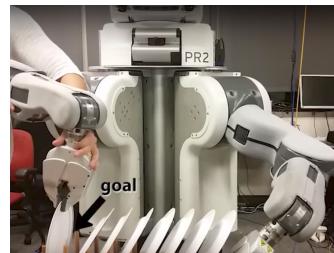
[Kolter et al. 2008]



[Ziebart et al. 2008]

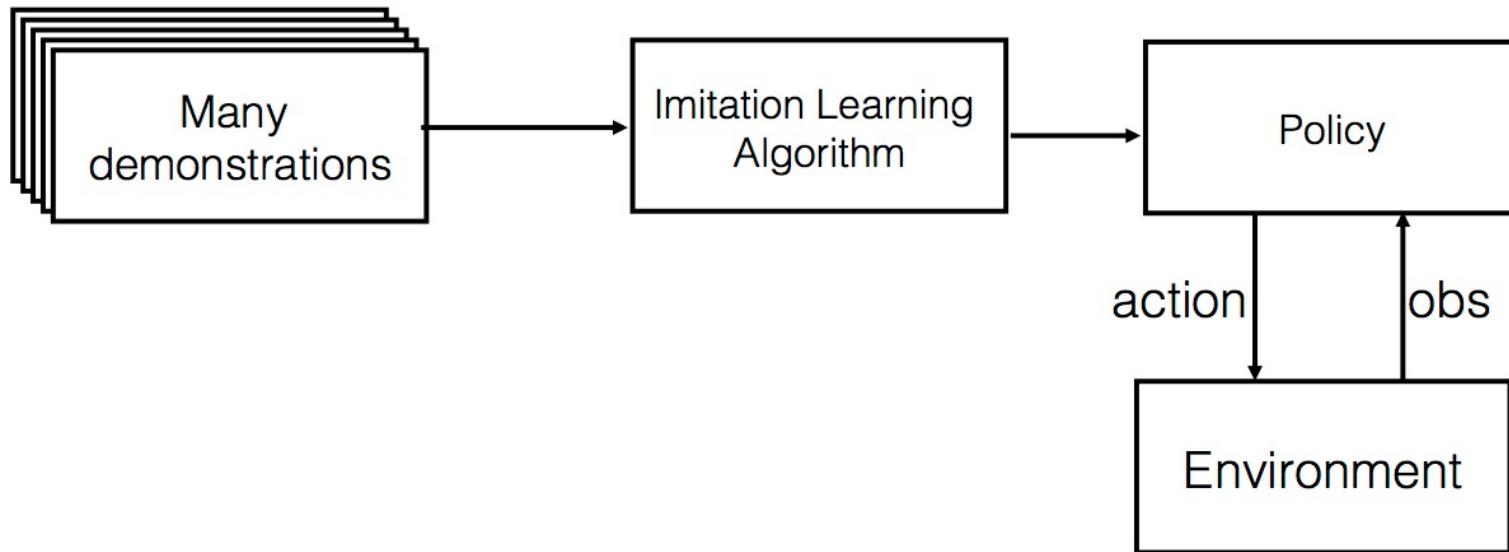


[Schulman et al. 2013]



[Finn et al. 2016]

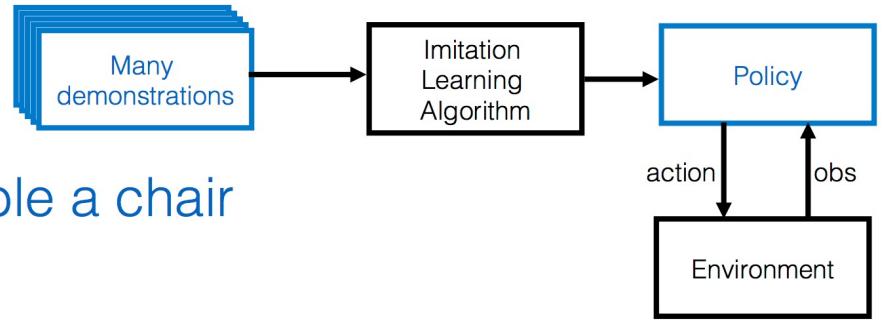
# Imitation Learning



# Imitation Learning

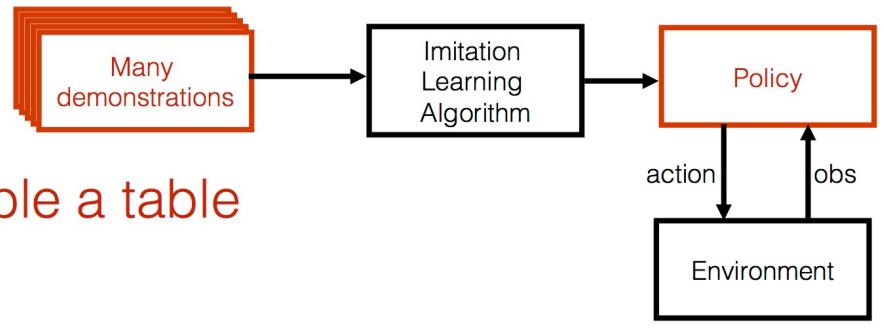
Task A

e.g. assemble a chair

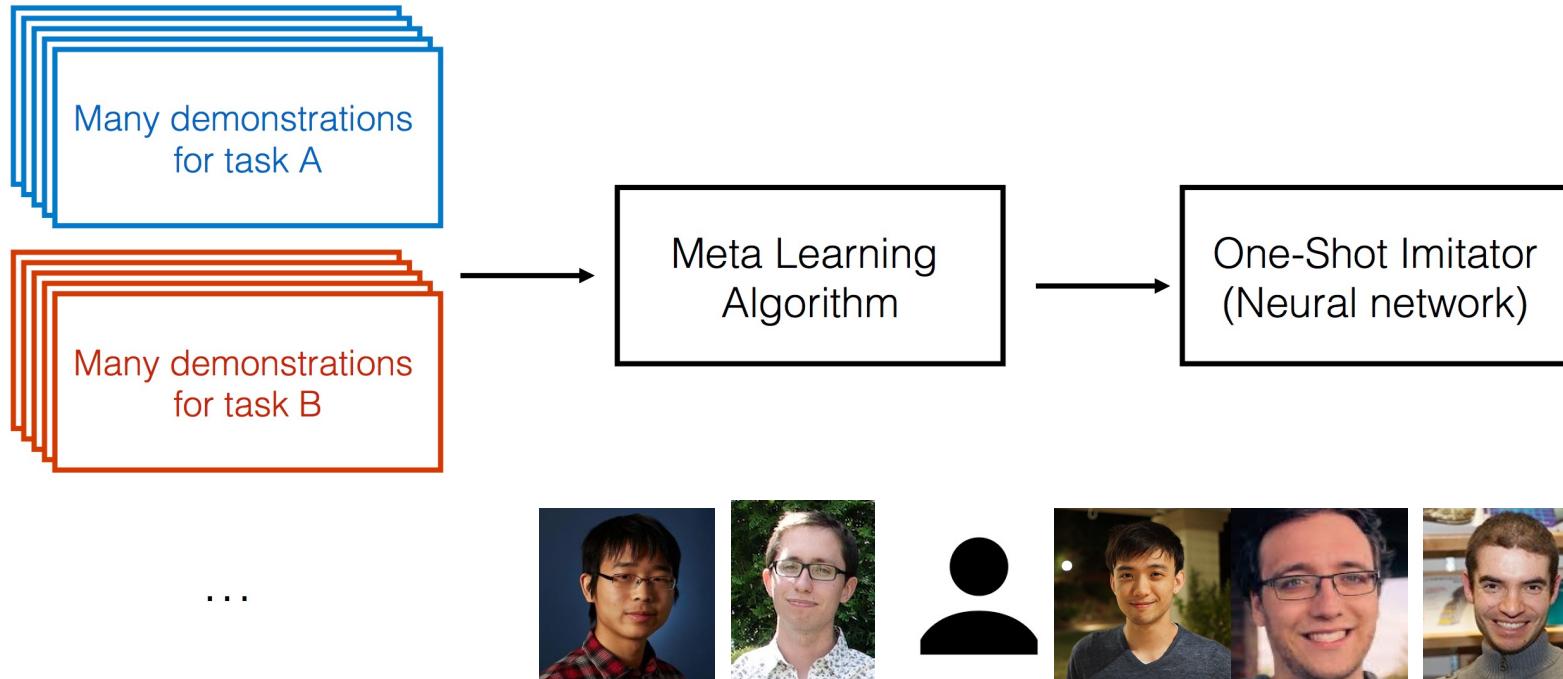


Task B

e.g. assemble a table

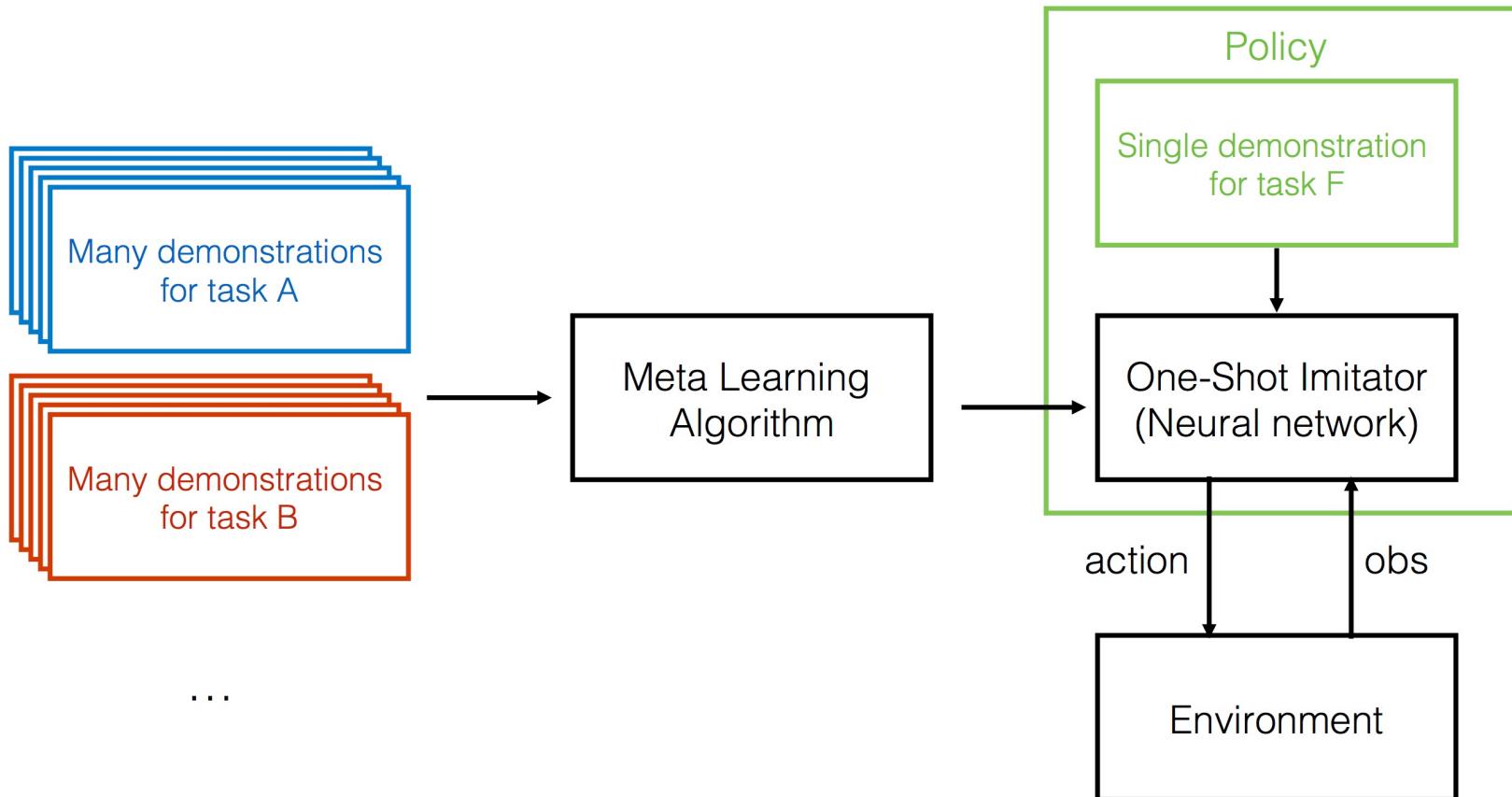


# One-Shot Imitation Learning

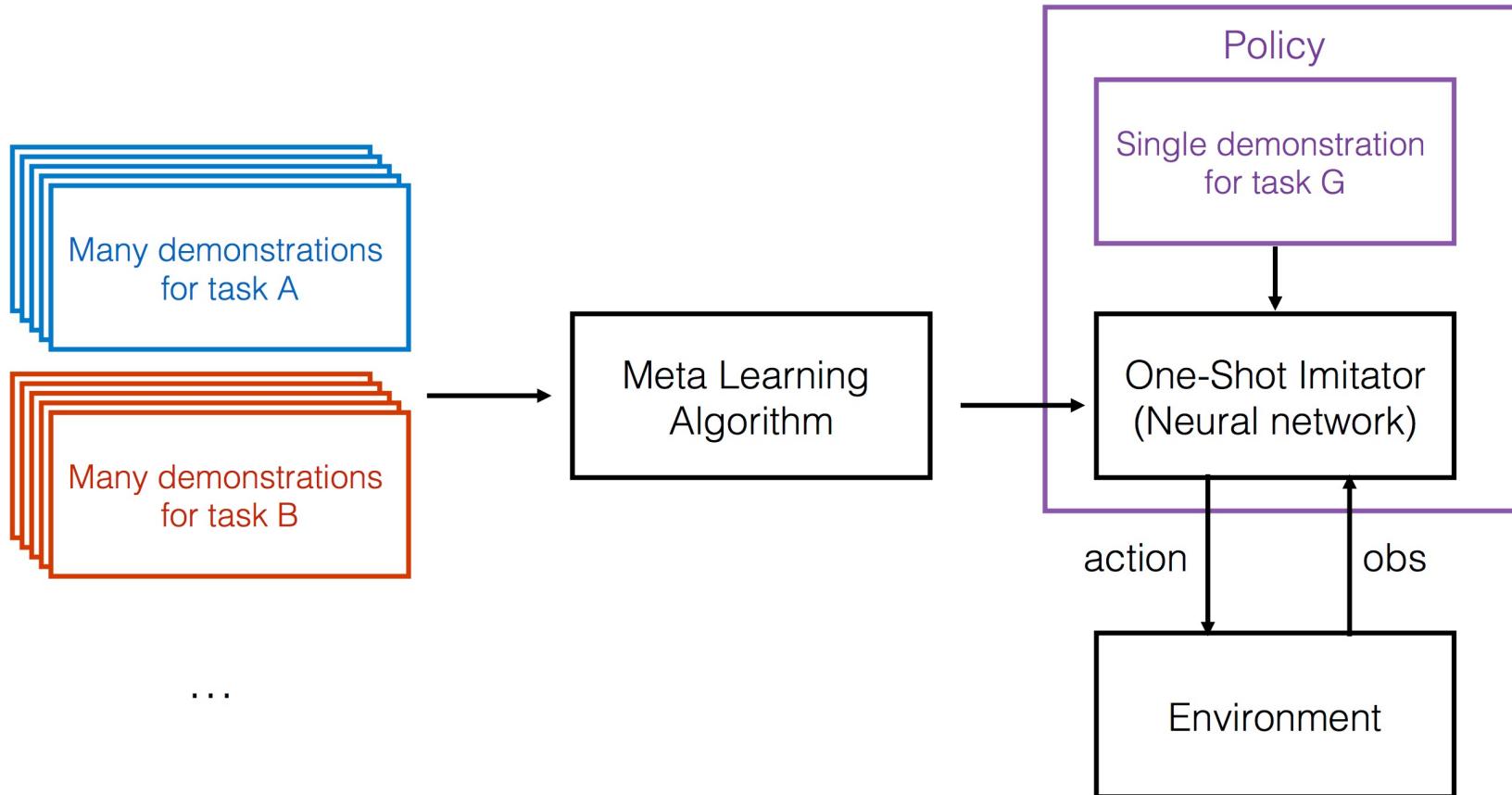


[Duan, Andrychowicz, Stadie, Ho, Schneider, Sutskever, Abbeel, Zaremba, 2017]

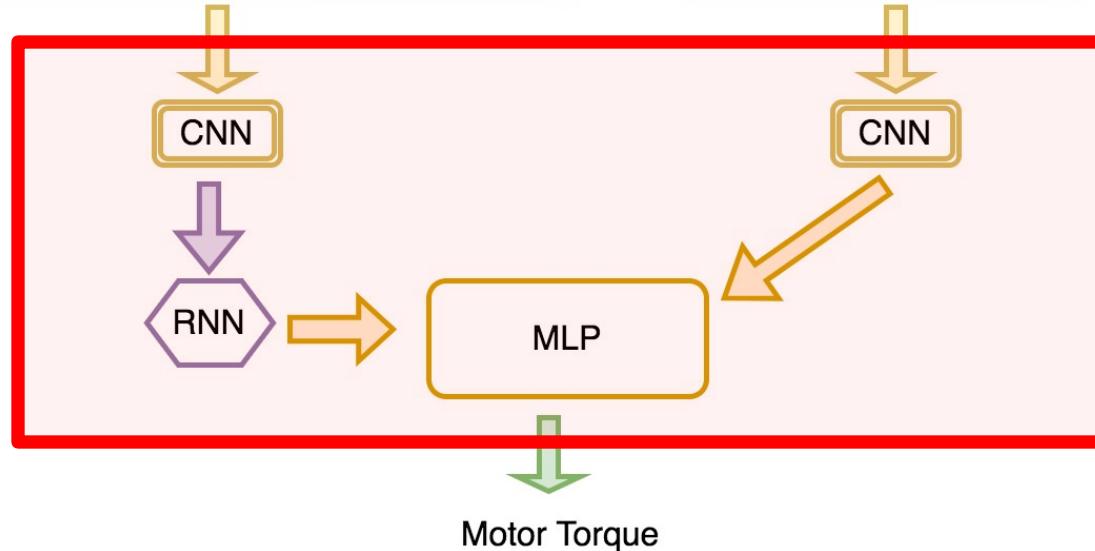
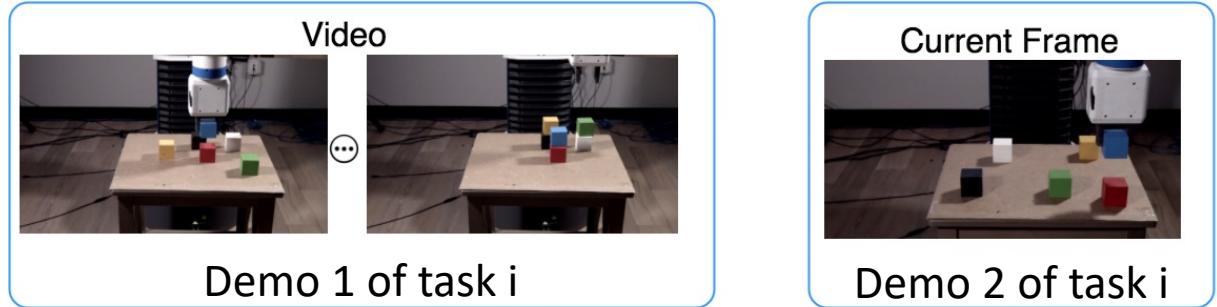
# One-Shot Imitation Learning



# One-Shot Imitation Learning



# Learning a One-Shot Imitator

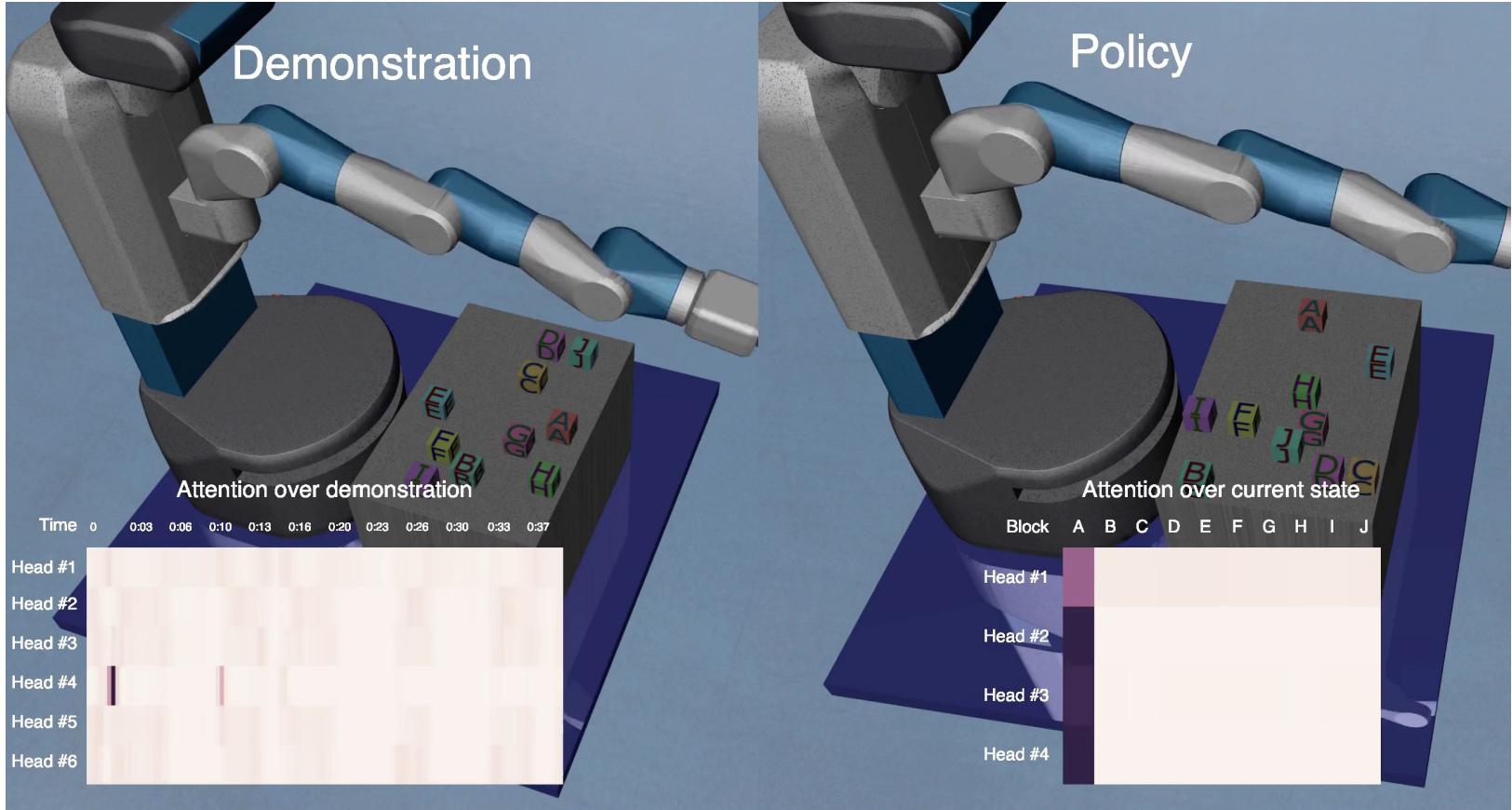


# Proof-of-concept: Block Stacking

- Each task is specified by a desired final layout
  - Example: abcd
    - “Place c on top of d, place b on top of c, place a on top of b.”



# Evaluation



# Many Exciting Directions in AI

- Unsupervised Learning
- Reinforcement Learning
- Unsupervised RL
- Meta-Reinforcement Learning
- Few-Shot Imitation
- ***Domain Randomization***
- DL for Science and Engineering
- Mitigating Bias
- Multi-modal Learning
- Architecture Search
- Value Alignment
- Scaling Laws
- Human-in-the-Loop
- Explainability

# Motivation for Simulation

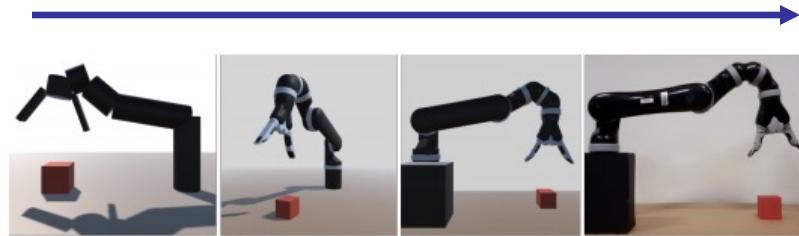
---

Compared to the real world, simulated data collection is...

- Less expensive
- Faster / more scalable
- Less dangerous
- Easier to label

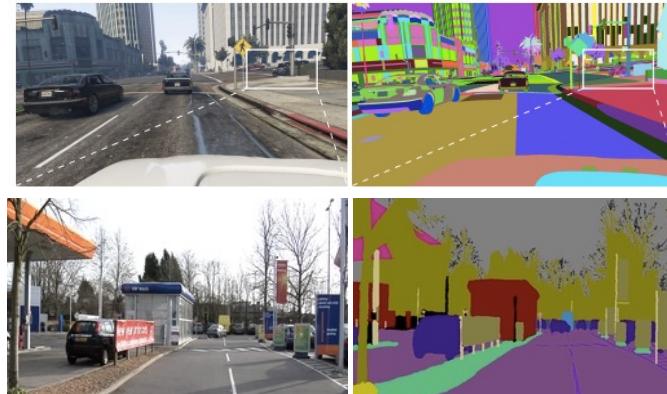
How can we learn useful real-world skills in the simulator?

# Approach 1 – Use Realistic Simulated Data



Simulation

Real world



GTA V

Real world

Carefully match the simulation to the world [1,2,3,4]

[1] Stephen James, Edward Johns. *3d simulation for robot arm control with deep q-learning* (2016)

[2] Johns, Leutenegger, Davison. *Deep learning a grasp function for grasping under gripper pose uncertainty* (2016)

[3] Mahler et al, Dex-Net 3.0 (2017)

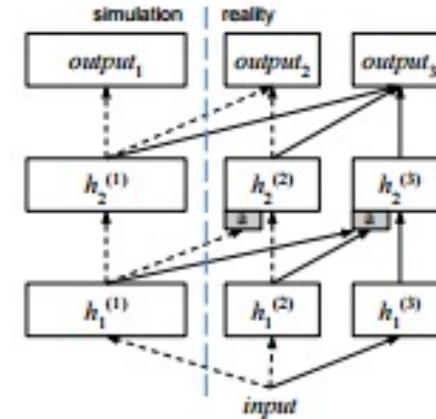
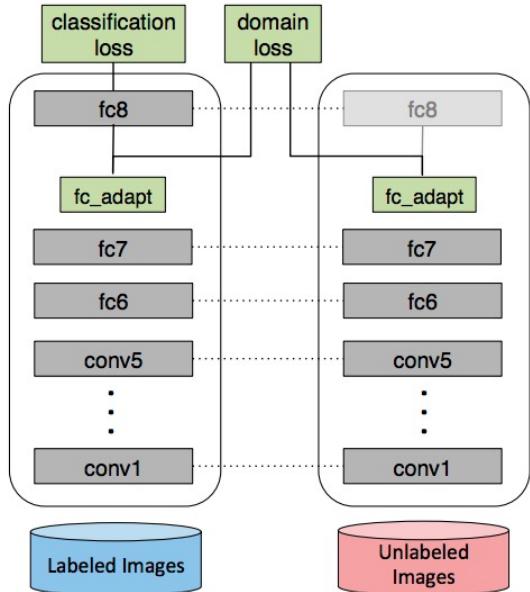
[4] Koenemann et al. *Whole-body model-predictive control applied to the HRP-2 humanoid.* (2015)

Augment simulated data with real data [5,6]

[5] Stephan R Richter, Vibhav Vineet, Stefan Roth, and Vladlen Koltun. *Playing for data: Ground truth from computer games* (2016)

[6] Bousmalis et al. *Using simulation and domain adaptation to improve efficiency of robotic grasping* (2017)

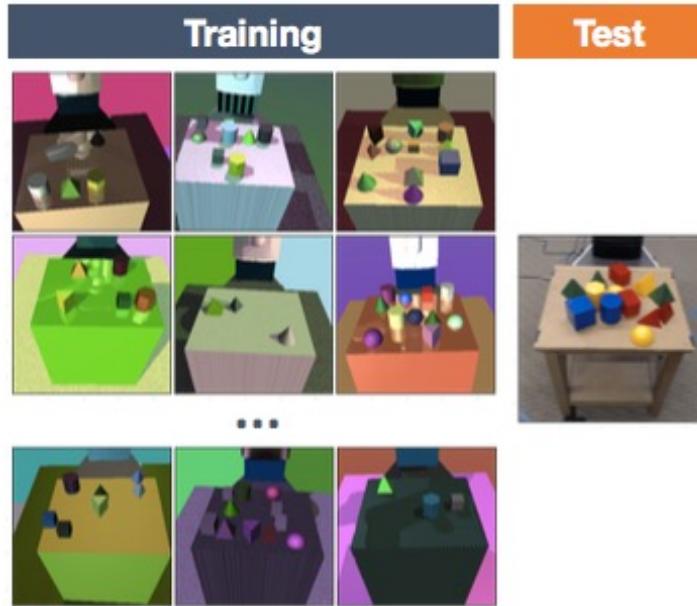
# Approach 2 – Domain Confusion / Adaptation



Eric Tzeng, Judy Hoffman, Ning Zhang, Kate Saenko, Trevor Darrell. Deep Domain Confusion: Maximizing for Domain Invariance. *arXiv preprint arXiv:1412.3474*, 2014.

Andrei A Rusu, Matej Vecerik, Thomas Rothořl, Nicolas Heess, Razvan Pascanu, and Raia Hadsell. Sim-to-real robot learning from pixels with progressive nets. *arXiv preprint arXiv:1610.04286*, 2016.

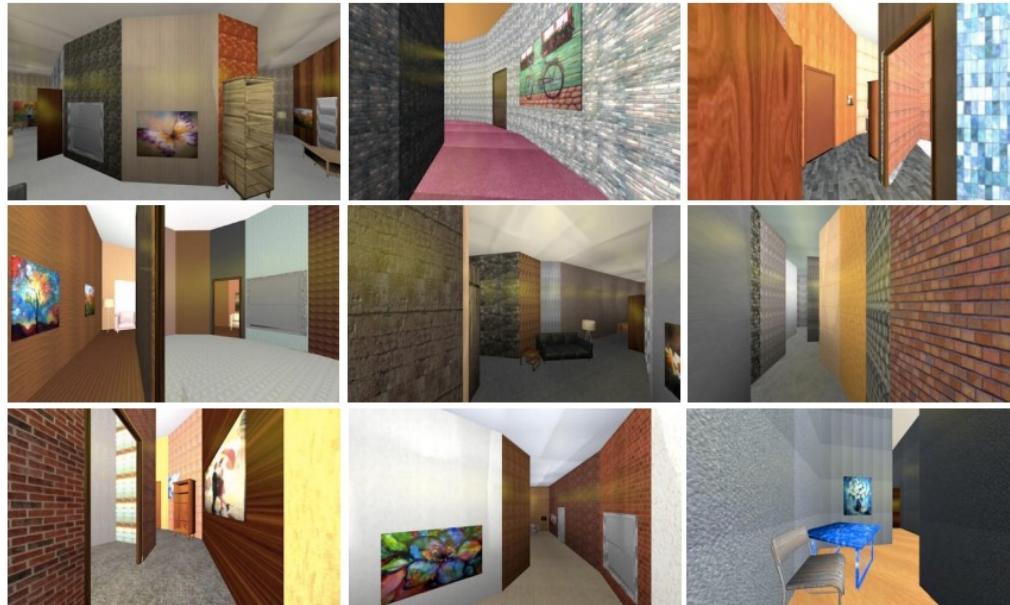
# Approach 3 – Domain Randomization



If the model sees enough simulated variation, the real world may look like just the next simulator

# Domain Randomization

**(cad)<sup>2</sup> rl: Real Single-Image Flight Without a Single Real Image.**



- Quadcopter collision avoidance
- ~500 semi-realistic textures, 12 floor plans
- ~40-50% of 1000m trajectories are collision-free

[3] Fereshteh Sadeghi and Sergey Levine. (cad)<sup>2</sup> rl: Real single-image flight without a single real image. *arXiv preprint arXiv:1611.04201*, 2016.

# Domain Randomization for Pose Estimation

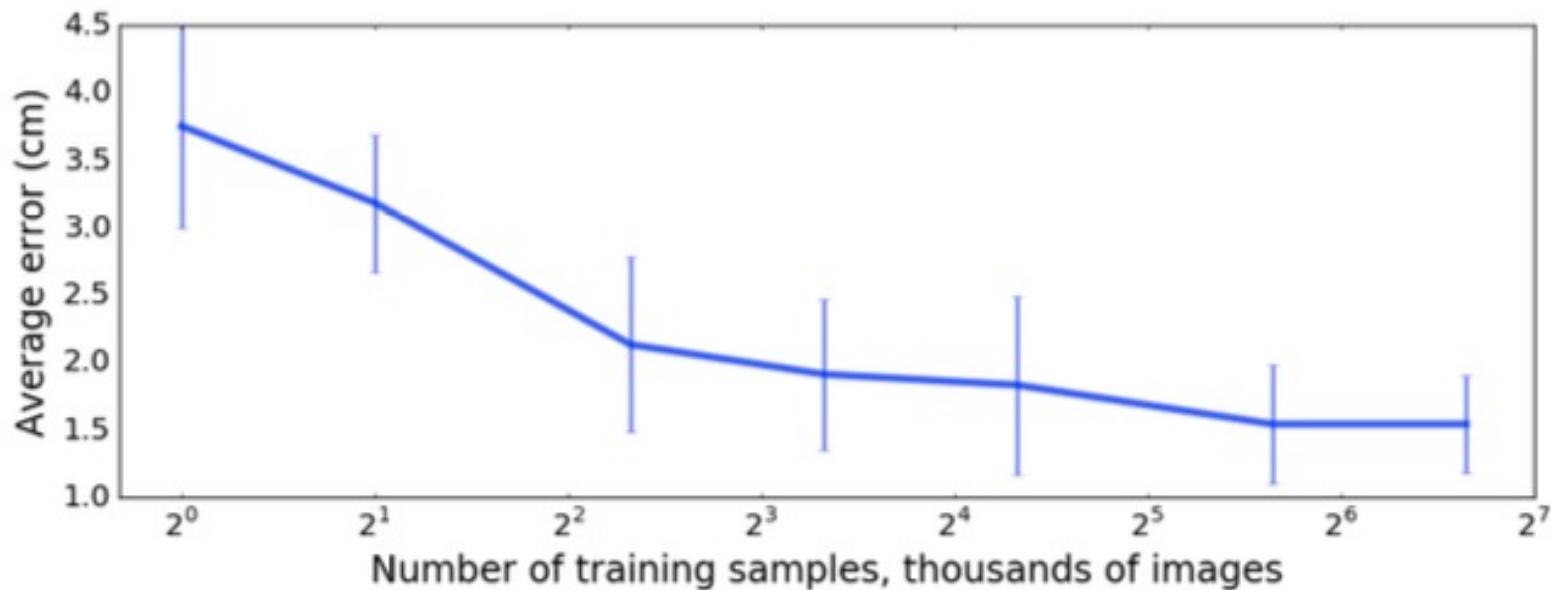


- Precise object pose localization
- 100K images with simple randomly generated textures

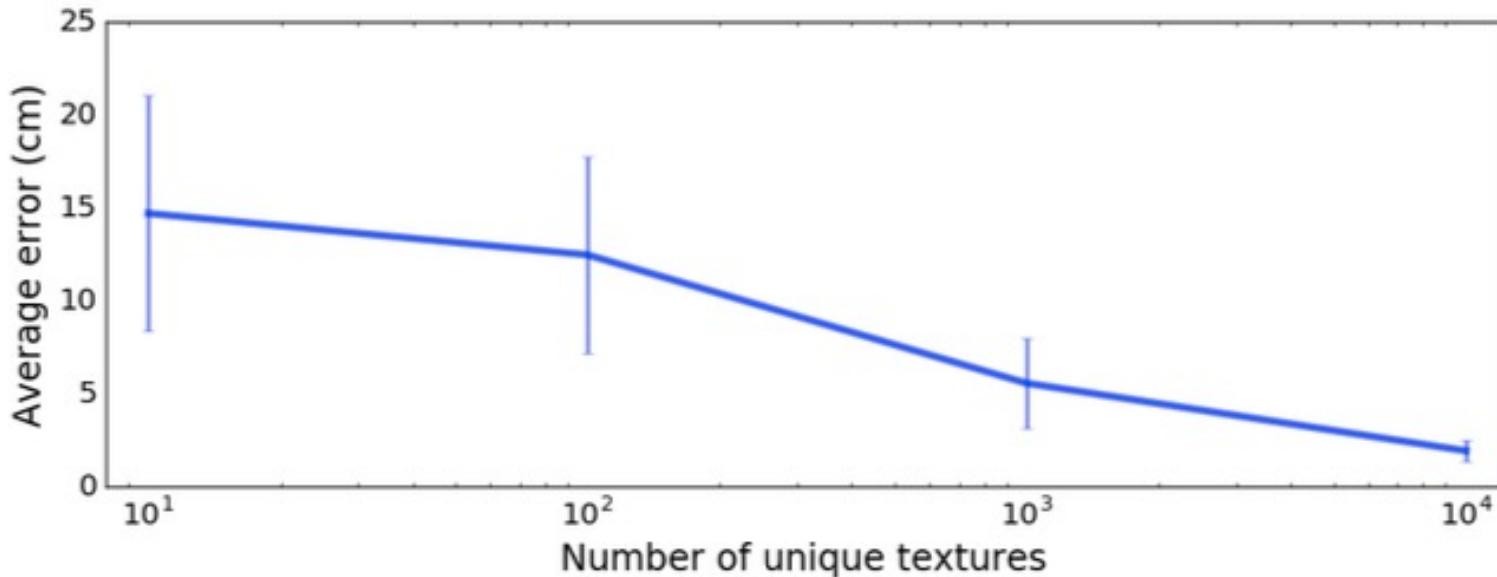


[Tobin, Fong, Ray, Schneider, Zaremba, Abbeel, 2017]

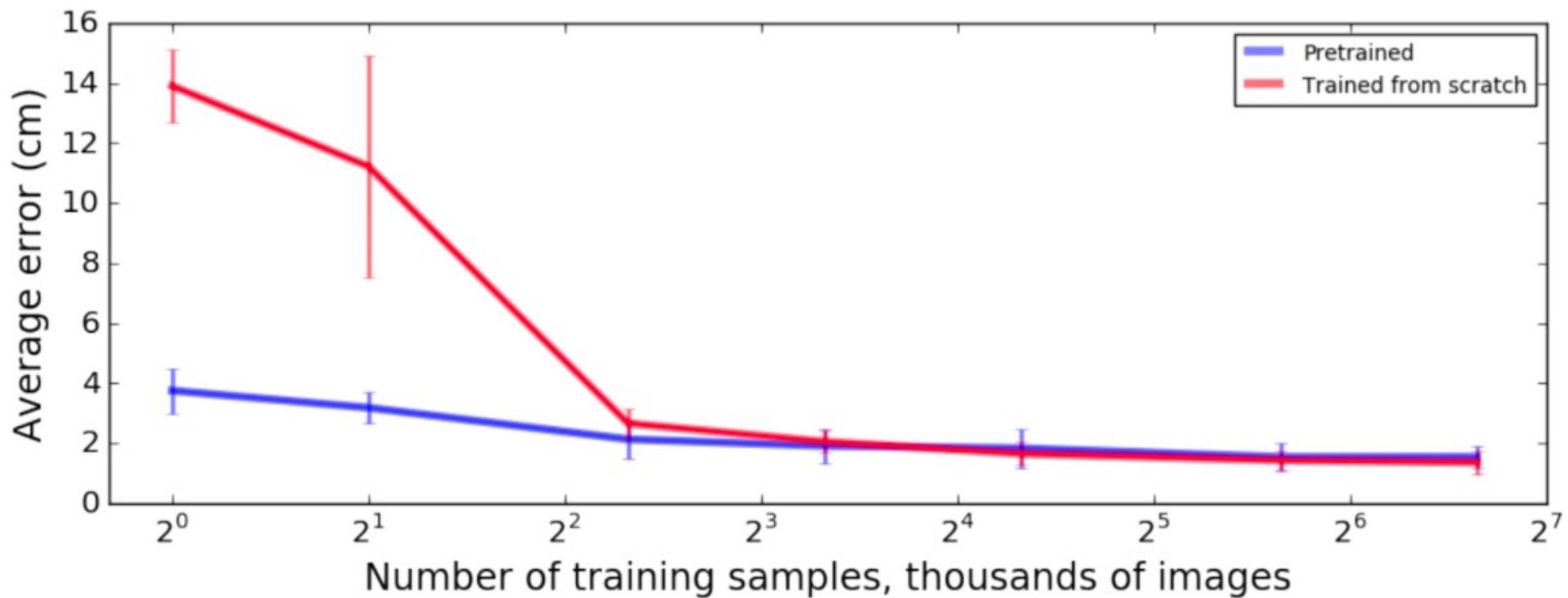
# How does it work? More Data = Better



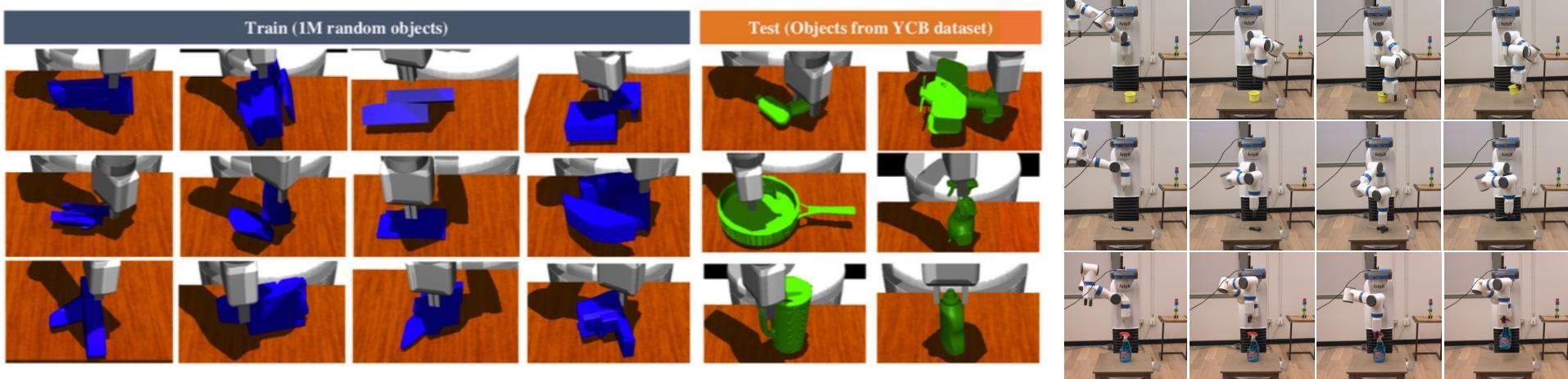
# More Textures = Better



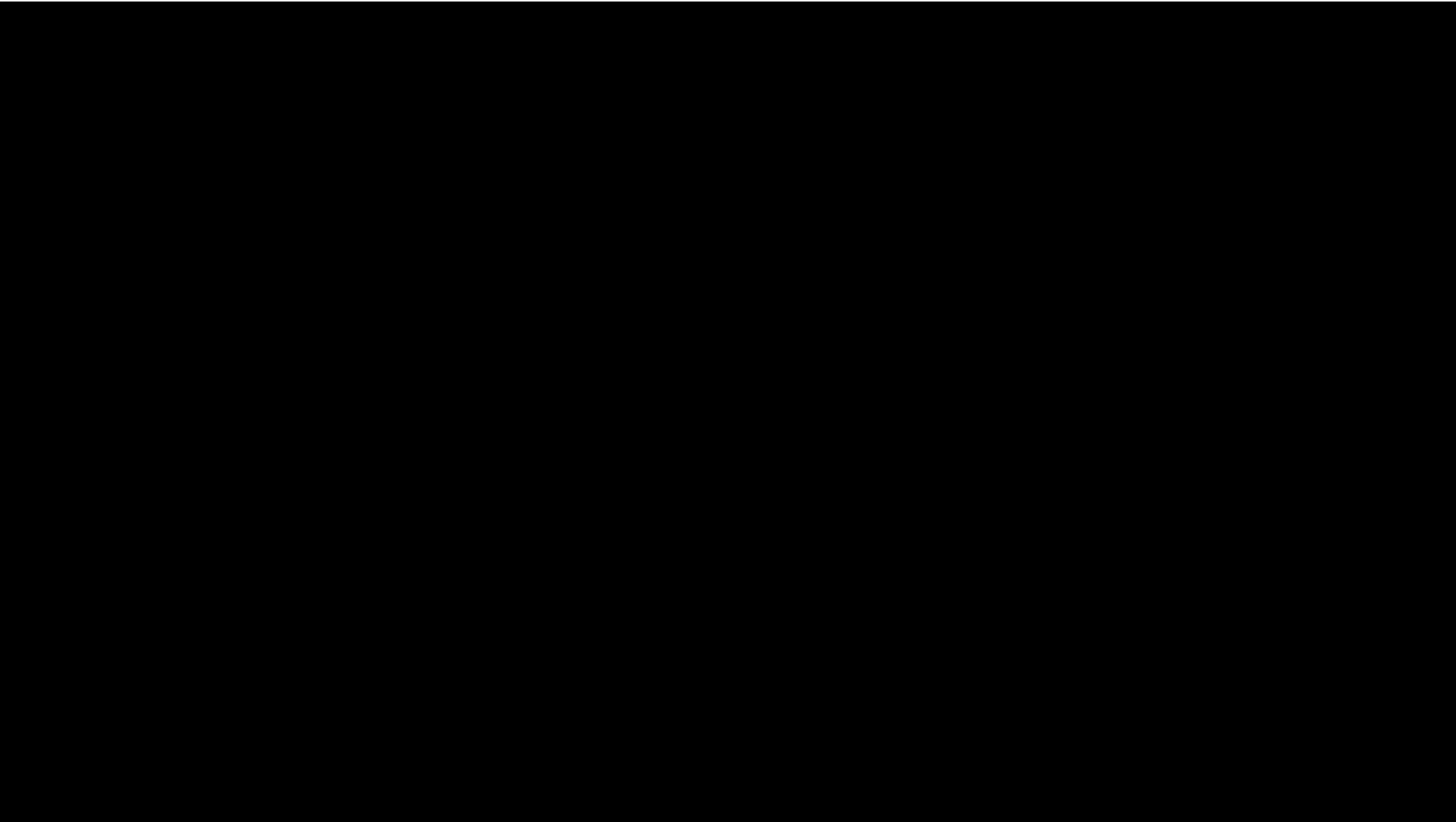
# Pre-Training is not Necessary

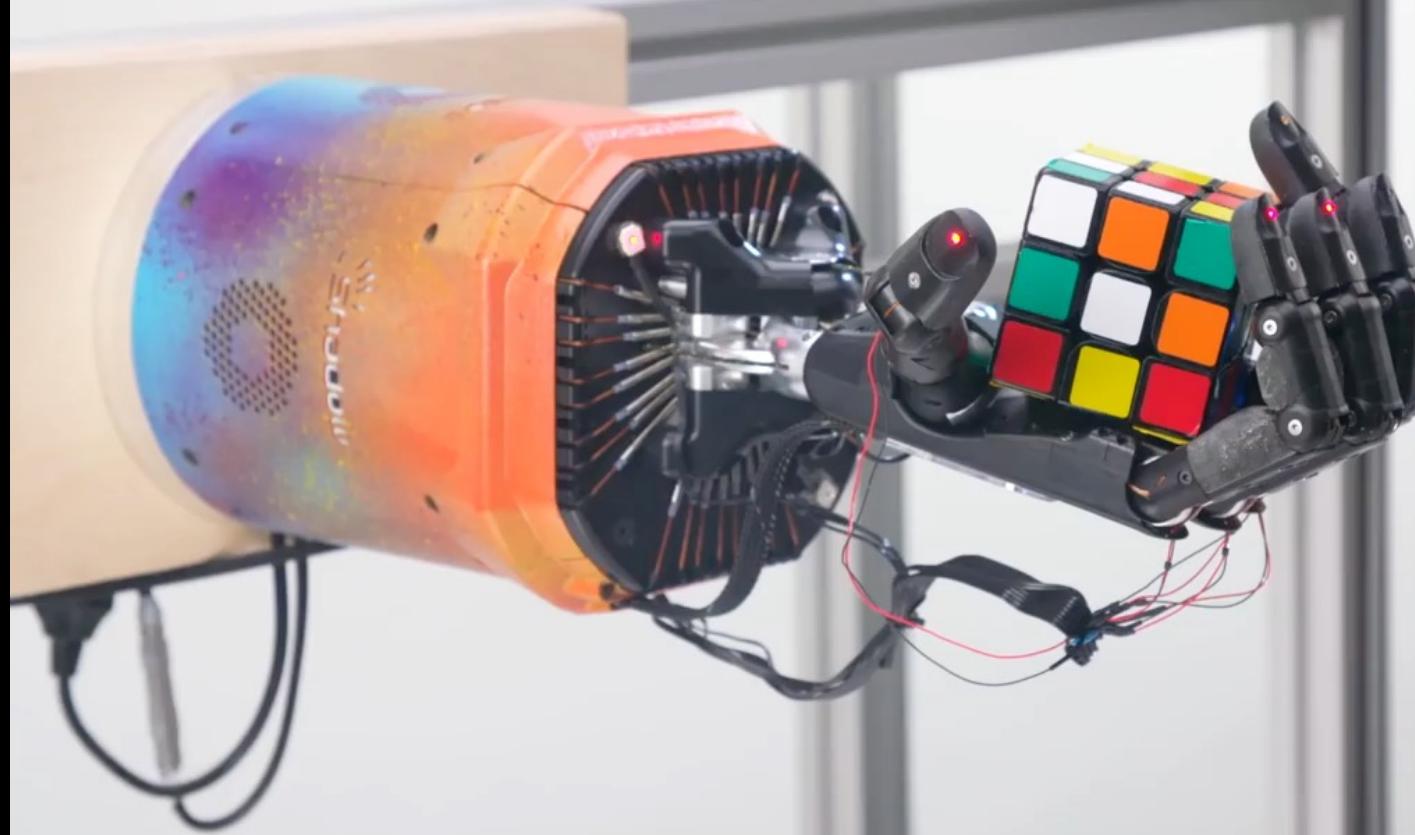


# Domain Randomization for Grasping



**Hypothesis:** Training on a diverse array of procedurally generated objects can produce comparable performance to training on realistic object meshes.





# Many Exciting Directions in AI

- Unsupervised Learning
- Reinforcement Learning
- Unsupervised RL
- Meta-Reinforcement Learning
- Few-Shot Imitation
- Domain Randomization
- ***DL for Science and Engineering***
- Mitigating Bias
- Multi-modal Learning
- Architecture Search
- Value Alignment
- Scaling Laws
- Human-in-the-Loop
- Explainability

NEWS · 30 NOVEMBER 2020

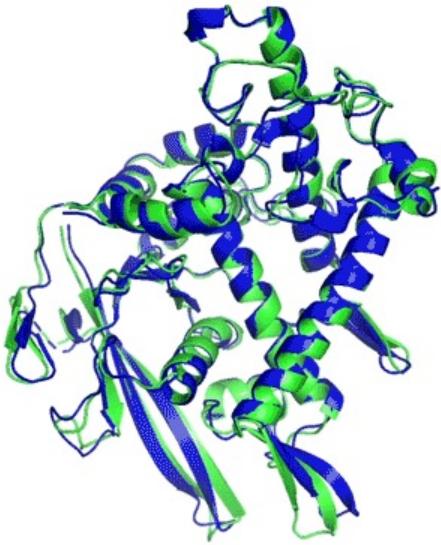
# 'It will change everything': DeepMind's AI makes gigantic leap in solving protein structures

Google's deep-learning program for determining the 3D shapes of proteins stands to transform biology, say scientists.

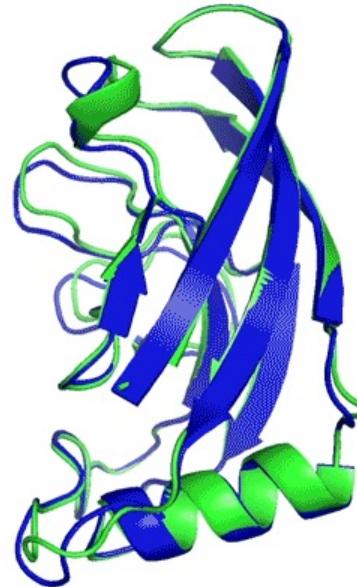
Ewen Callaway



A protein's function is determined by its 3D shape. Credit: DeepMind



**T1037 / 6vr4**  
90.7 GDT  
(RNA polymerase domain)

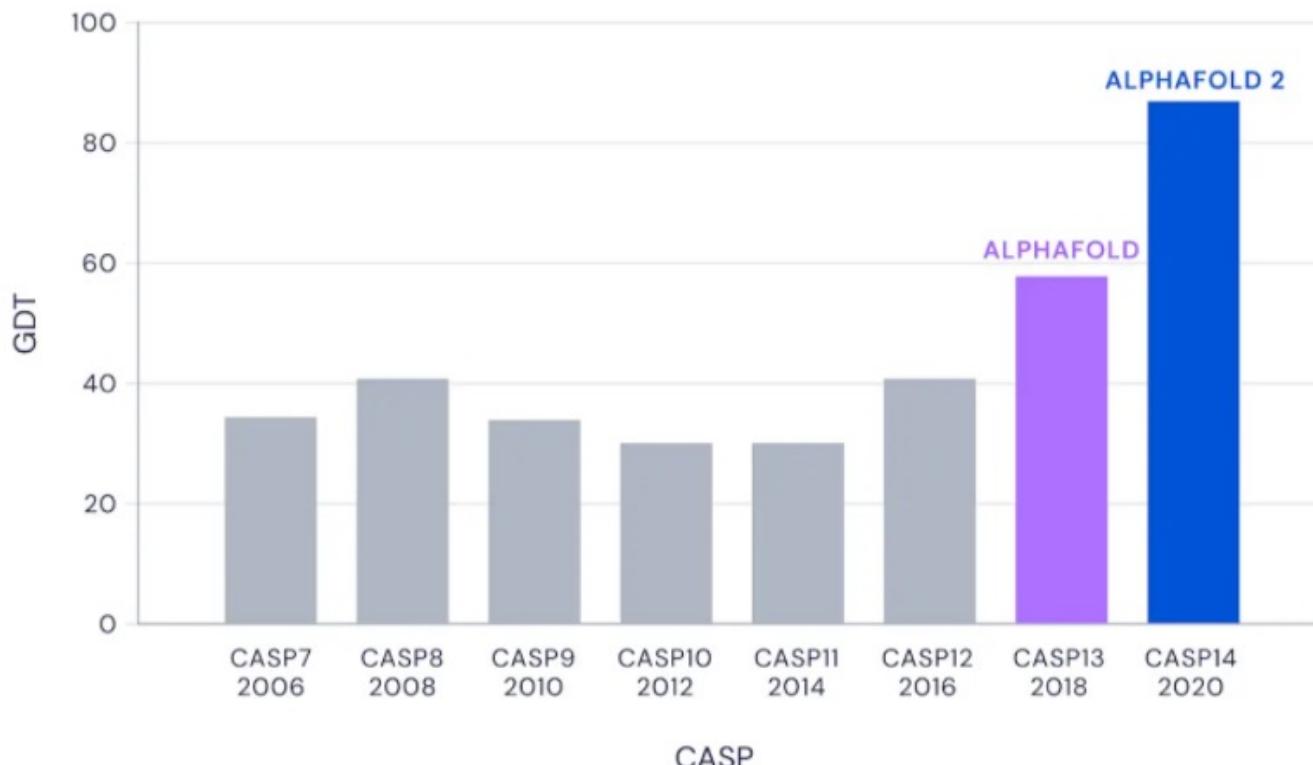


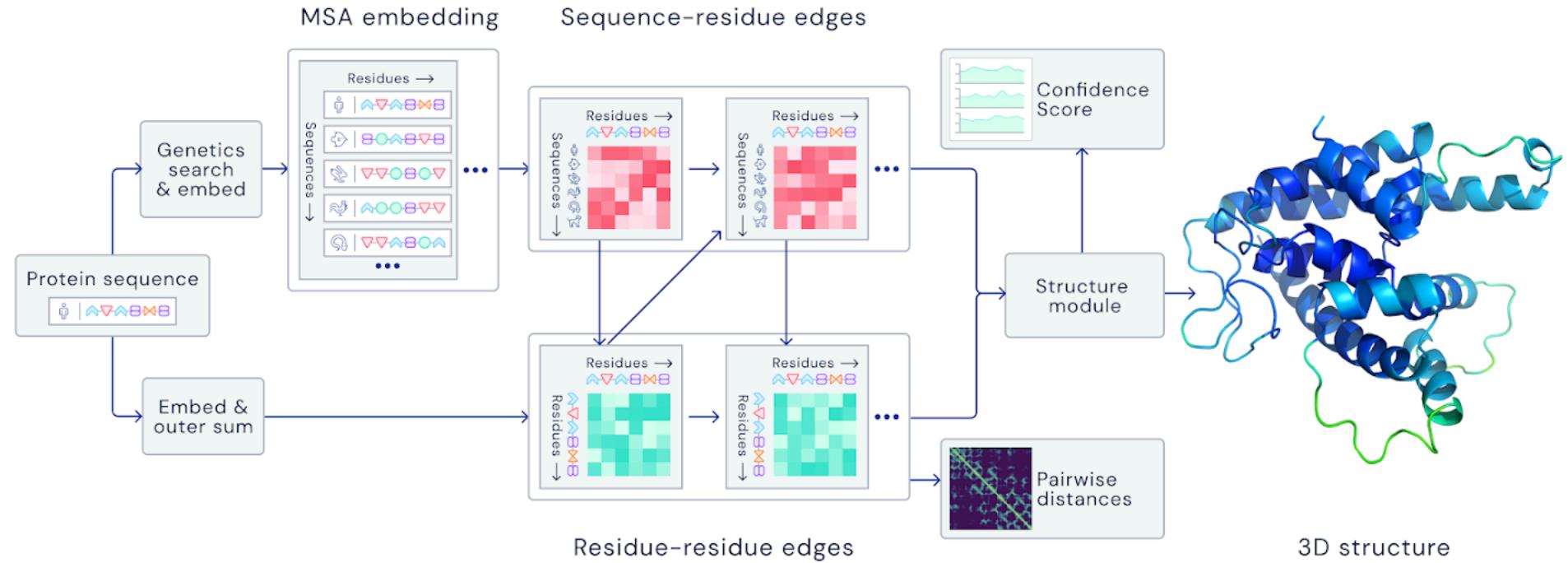
**T1049 / 6y4f**  
93.3 GDT  
(adhesin tip)

- Experimental result
- Computational prediction

# CASP 2020 Competition

Median Free-Modelling Accuracy



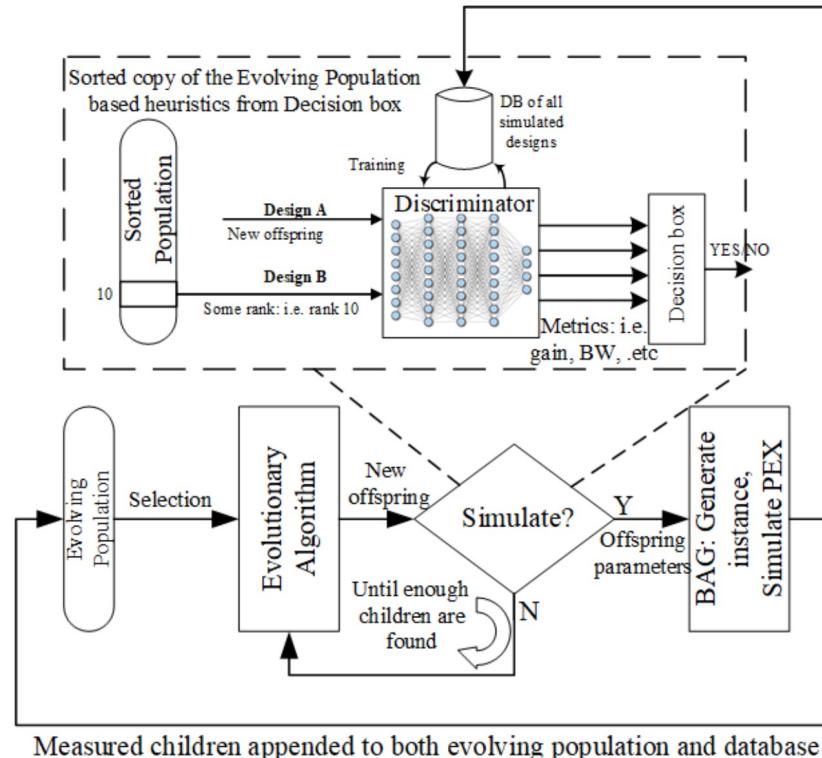


# Molecule Design

---

- DL to predict properties
  - optimize molecule against desired properties
- Synthesis
  - DL to propose synthesis steps

# A General Approach to Speed Up Design



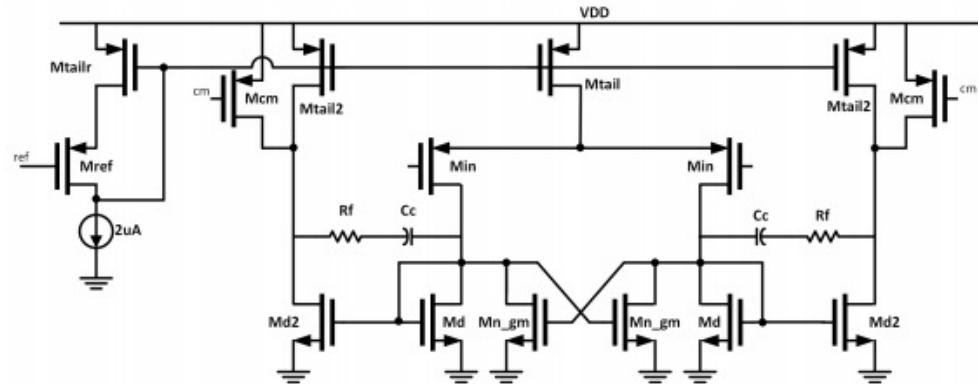
**BagNet: Berkeley Analog Generator with Layout Optimizer Boosted with Deep Neural Networks**

K. Hakhamaneshi, N. Werblun, P. Abbeel, V. Stojanovic.

IEEE/ACM International Conference on Computer-Aided Design (ICAD), Westminster, Colorado, November 2019.

<https://arxiv.org/abs/1907.10515>

# A General Approach to Speed Up Design



$R \rightarrow [\text{value} \times 60]$   
 $Cc \rightarrow [\text{value} \times 75]$   
 $Min \rightarrow [n\_finger \times 20]$   
 $Mn\_gm \rightarrow [n\_finger \times 20]$   
 $Md \rightarrow [n\_finger \times 20]$

$Md2 \rightarrow [n\_finger \times 20]$   
 $Mtail \rightarrow [n\_finger \times 20]$   
 $Mtail2 \rightarrow [n\_finger \times 20]$   
 $Mtailr \rightarrow [n\_finger \times 20]$   
 $Mcm \rightarrow [n\_finger \times 20]$   
 $Mref \rightarrow [n\_finger \times 20]$

Table 3: Performance of expert design methodology and our approach

	Requirement	Expert	Ours
$f_{unity}$	>100 MHz	382 MHz	159 MHz
$pm$	> 60°	64°	75°
gain	>100 (for ours)	42	105

Figure 5: Two stage op-amp with negative  $g_m$  load

# A General Approach to Speed Up Design

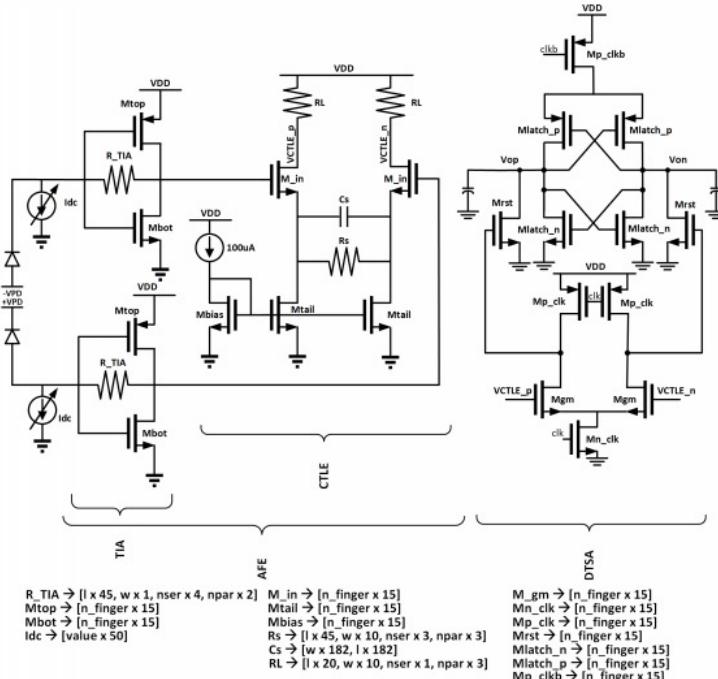


Figure 6: Optical receiver schematic

BagNet: Berkeley Analog Generator with Layout Optimizer Boosted with Deep Neural Networks

K. Hakhamaneshi, N. Werblun, P. Abbeel, V. Stojanovic.

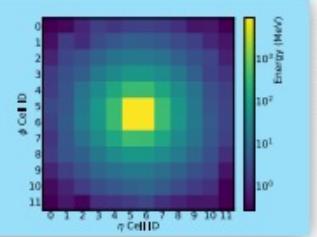
IEEE/ACM International Conference on Computer-Aided Design (ICAD), Westminster, Colorado, November 2019.

<https://arxiv.org/abs/1907.10515>

# GANs in High-Energy Physics

## Accelerating simulations

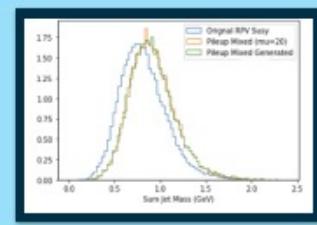
*replace or augment physics simulator*



M. Paganini, L. de Oliveira, BPP

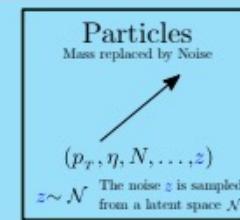
## Saving disk space

*replace libraries with on-the-fly generation*



W. Bhimani, W. Blair, S. Farnell, BPP

## Unbinned, high-dimensional interpolation



J. Lin, W. Bhimani, BPP

# Symbolic Math: Integrals and ODEs

Equation	Solution
$y' = \frac{16x^3 - 42x^2 + 2x}{(-16x^8 + 112x^7 - 204x^6 + 28x^5 - x^4 + 1)^{1/2}}$	$y = \sin^{-1}(4x^4 - 14x^3 + x^2)$
$3xy \cos(x) - \sqrt{9x^2 \sin(x)^2 + 1}y' + 3y \sin(x) = 0$	$y = c \exp(\sinh^{-1}(3x \sin(x)))$
$4x^4yy'' - 8x^4y'^2 - 8x^3yy' - 3x^3y'' - 8x^2y^2 - 6x^2y' - 3x^2y'' - 9xy' - 3y = 0$	$y = \frac{c_1 + 3x + 3 \log(x)}{x(c_2 + 4x)}$

Table 4: Examples of problems that our model is able to solve, on which Mathematica and Matlab were not able to find a solution. For each equation, our model finds a valid solution with greedy decoding.

# Symbolic Math: Integrals and ODEs

	Integration (BWD)	ODE (order 1)	ODE (order 2)
Mathematica (30s)	84.0	77.2	61.6
Matlab	65.2	-	-
Maple	67.4	-	-
Beam size 1	98.4	81.2	40.8
Beam size 10	99.6	94.0	73.2
Beam size 50	99.6	97.0	81.0

# Learn more about Deep Learning for Science/Engineering?

- <https://deepmind.com/blog/article/alphafold-a-solution-to-a-50-year-old-grand-challenge-in-biology>
- **AlphaFold: Improved protein structure prediction using potentials from deep learning**  
Deepmind (Senior et al)  
<https://www.nature.com/articles/s41586-019-1923-7>
- **BagNet: Berkeley Analog Generator with Layout Optimizer Boosted with Deep Neural Networks**  
K. Hakhamaneshi, N. Werblun, P. Abbeel, V. Stojanovic.  
IEEE/ACM International Conference on Computer-Aided Design (ICAD), Westminster, Colorado, November 2019.  
<https://arxiv.org/abs/1907.10515>
- **Evaluating Protein Transfer Learning with TAPE**  
R. Rao, N. Bhattacharya, N. Thomas, Y. Duan, X. Chen, J. Canny, P. Abbeel, Y. Song  
<https://www.biorxiv.org/content/10.1101/676825v1>
- **Opening the black box: the anatomy of a deep learning atomistic potential**  
Justin Smith  
<https://drive.google.com/file/d/1f1iiXKzxNbNz5lL5x2ob9xGYLHNHhxU/view>
- **Exploring Machine Learning Applications to Enable Next-Generation Chemistry**  
Jennifer Wei (Google)  
[https://docs.google.com/presentation/d/1zGoxOMWmid25hgtSgVkJYu1-GP9PT3EQ2\\_tzOlQZcF4/edit#slide=id.p1](https://docs.google.com/presentation/d/1zGoxOMWmid25hgtSgVkJYu1-GP9PT3EQ2_tzOlQZcF4/edit#slide=id.p1)
- **GANs for HEP**  
Ben Nachman  
<https://drive.google.com/file/d/1op6Q6OuVZvJ4VbtLkemi3oJWtSBx5FCc/view>
- **Deep Learning for Symbolic Mathematics**  
G. Lample and F. Charton  
<https://openreview.net/pdf?id=S1eZYeHFDS>
- **A Survey of Deep Learning for Scientific Discovery**  
Maithra Raghu, Eric Schmidt  
<https://arxiv.org/abs/2003.11755>

# Outline

---

- Sampling of research directions
- *Overarching research trend*
- How to keep up

# Compute Increasing Rapidly

## Companies Developing Deep Learning Chips

	Company	HQ	Story
Public	Ambarella	United States	Developing computer vision chips for autonomous cars
	AMD	United States	GPU based deep learning
	Facebook	United States	Forming a team to build SoCs and perhaps inference chips
	Google	United States	Custom designed TPU deployed in Google Cloud
	Intel	United States	Developing a Neural Network chip based on Nervana acquisition
	Nvidia	United States	Current market leader using GPU based deep learning
	Tesla	United States	Developing a custom AI chip for autonomous driving
Private	Bitmain	China	Top maker of Bitcoin mining chips
	Cambricon	China	China's state-backed startup with \$1B+ valuation.
	Cerebras Systems	United States	Ex-AMD team backed by Benchmark Capital
	DeePhi	China	China based startup with a focus on video analysis
	GraphCore	United Kingdom	Building a 16nm deep learning chip for training and inference
	Groq	United States	Ex-Google TPU team backed by Social Capital
	Horizon Robotics	China	Ex-Baidu team. Embedded / computer vision focus
	KnuEdge	United States	Headed by former NASA CTO
	Mythic	United States	In-memory inference for IoT backed by DFJ
	Tenstorrent	Canada	Toronto based chip startup
	Thinci	United States	Computer vision / auto focus
	Wave Computing	United States	Makes data flow acceleration servers

Source: ARK Invest  
[20-April-2018]

# Let's Calibrate Compute Scale

Architecture	Num neurons	Num synapses

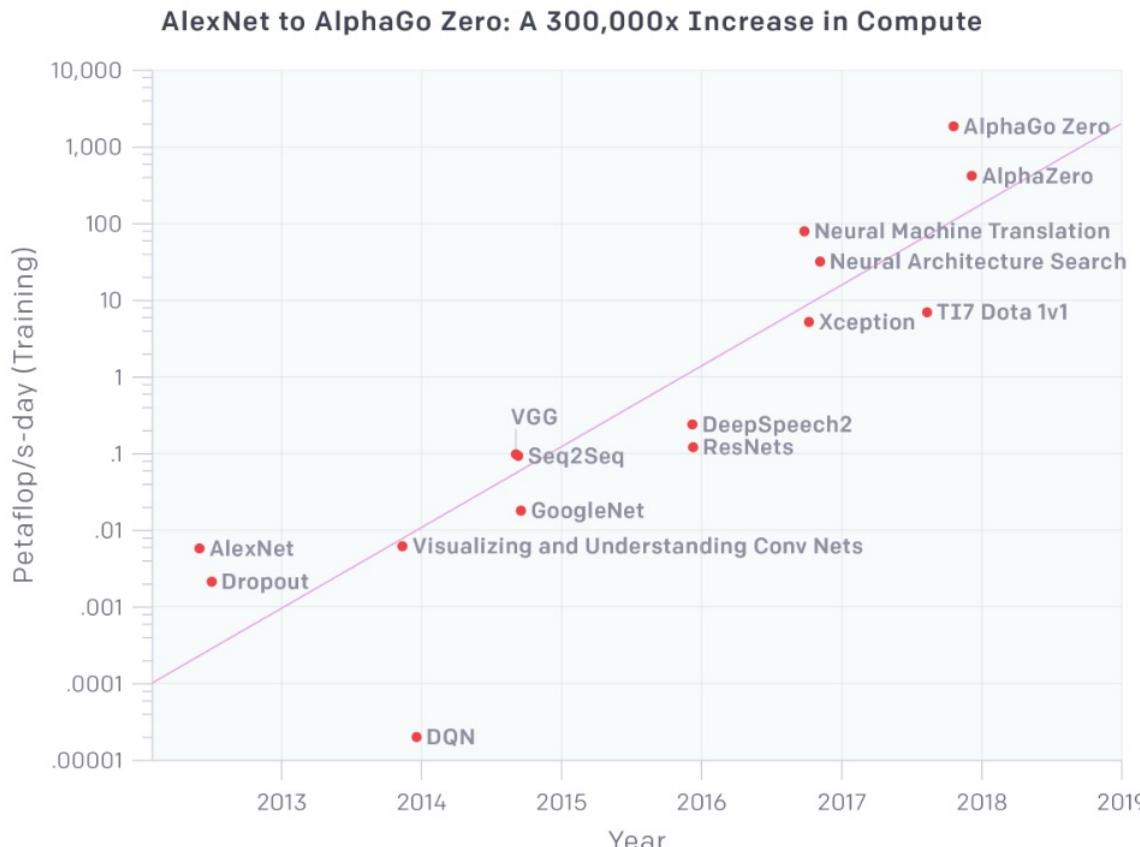
If each synapse is 1 FLOP (i.e., can fire / not fire once per second),

Then human brain requires  $10^{15}$  flops = 1 petaflop.

NVIDIA's DGX-2 = 2 petaflops in one server rack! (\$400,000 price tag)

2 TPU-v3 slices = ~1 petaflop in Google Cloud at \$16/hour (\$5/hour if pre-emptible)

# Compute per Major Result

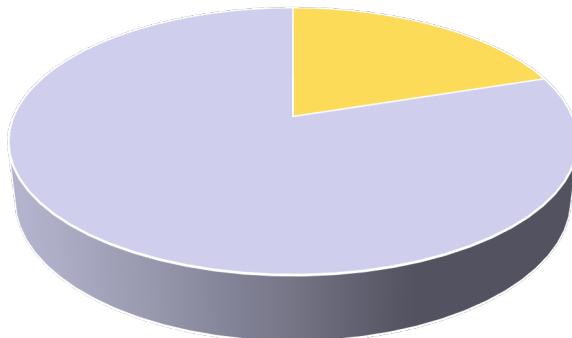


# Implications for research agenda?

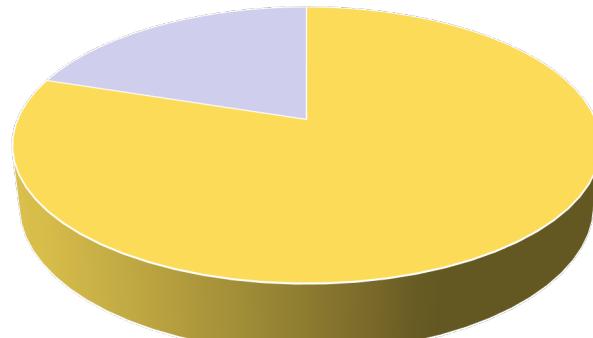
---

# Key Enablers for AI Progress?

DATA  
COMPUTE  
HUMAN INGENUITY



Problem Territory 1



Problem Territory 2

E.g., Deep Learning to Learn

# Outline

- Sampling of research directions
- Overarching research theme
- ***How to keep up***



# How to Keep Up

---

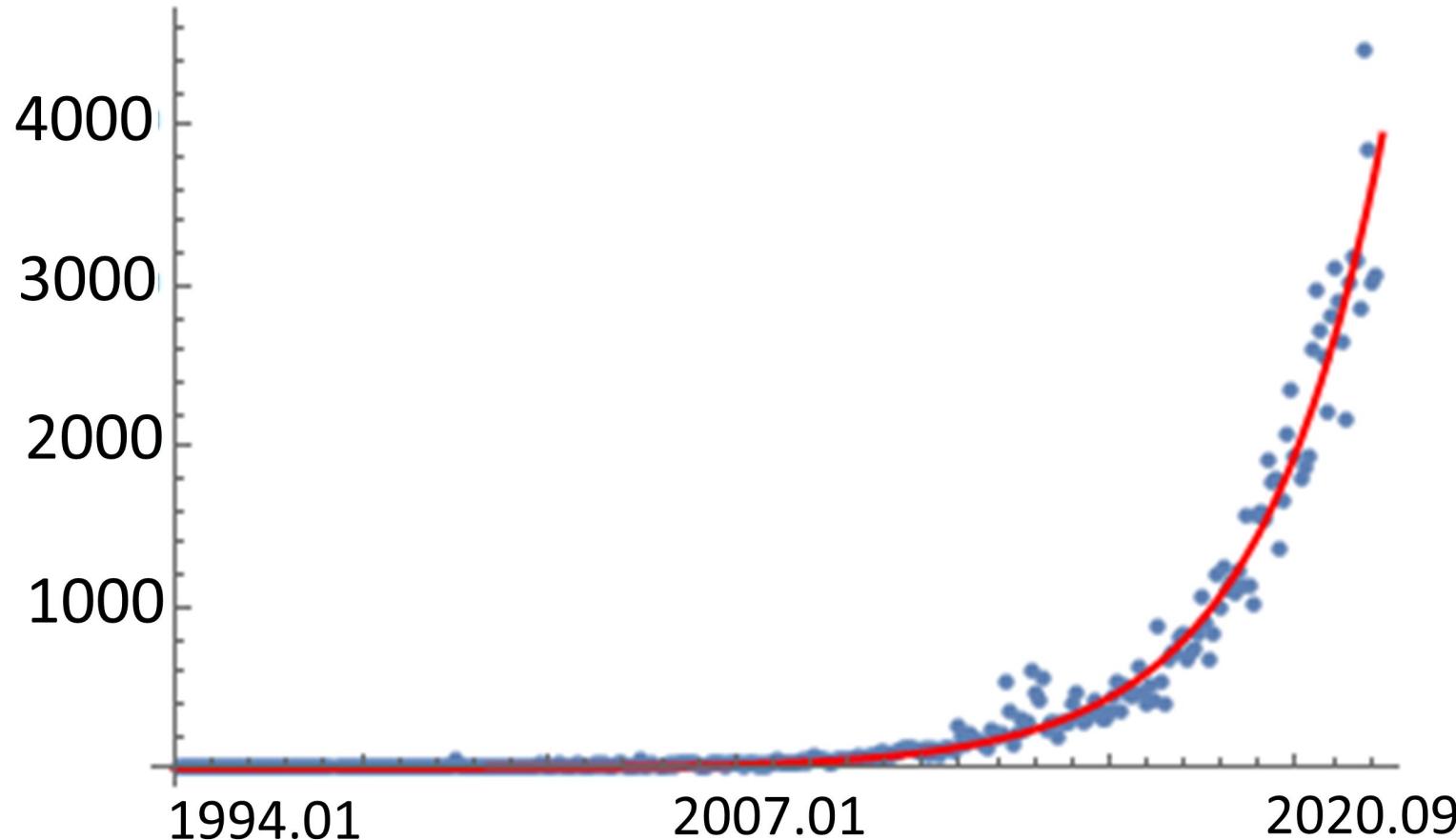
- (mostly) keeping up with (mostly) not reading papers
- when you decide to read papers

# How to Keep Up

---

- ***(mostly) keeping up with (mostly) not reading papers***
- when you decide to read papers

# ML+AI arXiv papers per month



# Great Resources

---

- Tutorials at conferences (e.g., ICML, NeurIPS)
- Graduate courses and seminars
- Yannic Kilcher youtube channel – explains papers
- Two Minute Papers youtube channel
- The Batch – newsletter by Andrew Ng
- Import AI – newsletter by Jack Clark

# How to Keep Up

---

- (mostly) keeping up with (mostly) not reading papers
- *when you decide to read papers*

# Which papers to read? (1/2)

---

- Tutorials at conferences (e.g., ICML, NeurIPS)
- Graduate courses and seminars
- Yannic Kilcher youtube channel – explains papers
- Two Minute Papers youtube channel
- The Batch – newsletter by Andrew Ng
- Import AI – newsletter by Jack Clark

# Which papers to read? (2/2)

## ■ Arxiv sanity (by Karpathy)



## ■ Twitter

- Jack Clark, Karpathy, Ian Goodfellow, hardmaru, smerity, pabbeel, ...

## ■ AI/DL Facebook Group: Very active group where members post anything from news articles to blog posts to general ML questions

## ■ ML Subreddit

# How to Read a Paper?

---

- Read the title, abstract, section headers, and figures
- Try and find slides or a video on the paper (these do not have to be by the authors).
- Read the introduction (Jennifer Widom)
  - What is the problem?
  - Why is it interesting and important?
  - Why is it hard? (e.g., why do naive approaches fail?)
  - Why hasn't it been solved before? (Or, what's wrong with previous proposed solutions? How does this paper differ?)
  - What are the key components of this approach and results? Any limitations?
- Skim the related work.
  - Is there any related work you are familiar with?
  - If so, how does this paper relate to those works?
- Skim the technical section.
  - Where are the novelties?
  - What are the assumptions?
- Read the experiments.
  - What questions are the experiments answering?
  - What questions are the experiments not answering?
  - What baselines do they compare against?
  - How strong are these baselines?
  - Is the experiment methodology sound?
  - Do the results support their claims?
- Read the conclusion/discussion.
- Read the technical section.
  - Read in "good faith" (i.e., assume the authors are correct)
  - Skip over confusing parts that don't seem fundamental
  - If any important part is confusing, see if the material is in class slides or prior papers
- Read the paper as you see fit.

# Reading Group

---

- ***Read papers together with friends:***
  - Mode 1: everyone reads, then discusses
  - Mode 2: 1 (or 2) people read, give small tutorial to others

# PS: Why do a PhD? (or not)

---

- Become one of the world's experts in a topic you really care about
- Technically deep and demanding
- Crudely: develop new tools/techniques rather than use existing tools/techniques

Thank you