# Robotics at OpenAI

May 1, 2017
By Wojciech Zaremba

# Why OpenAI?

- OpenAI's mission is to build safe AGI, and ensure AGI's benefits are as widely and evenly distributed as possible.

Wojciech Zaremba - OpenAI

- OpenAI's mission is to build safe AGI, and ensure AGI's benefits are as widely and evenly distributed as possible.
- Can we build a General Purpose Robot and deploy it in the most beneficial way to humans?
  - We have several ideas
  - But maybe you can help us through collaboration!

Wojciech Zaremba - OpenAI

# Why OpenAI?

- OpenAI's mission is to build safe AGI, and ensure AGI's benefits are as widely and evenly distributed as possible.
- Can we build a General Purpose Robot and deploy it in the most beneficial way to humans?
  - We have several ideas
  - But maybe you can help us through collaboration!
- We're well positioned to do this, due to extraordinary researchers, engineers, and amount of compute

# What is a General Purpose Robot ?

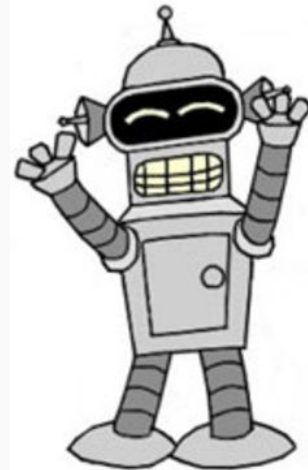**A robot that can solve a variety of tasks without being trained on them**

- Currently, all robots are trained to solve a single task
  - Roomba cannot drive a car or play chess

- Human has general purpose capabilities
  - Human can clean an apartment, drive a car, and play chess

Wojciech Zaremba - OpenAI

**A robot that can solve a variety of tasks without being trained on them**



- Currently, all robots are trained to solve a single task
  - Roomba cannot drive a car or play chess


- Human has general purpose capabilities
  - Human can clean an apartment, drive a car, and play chess

Wojciech Zaremba - OpenAI

# What is a General Purpose Robot ?

We think that the following are critical components of the General
Purpose Robot

- Training on diverse environments

- Obtaining complex behaviours

- Having a way to ask a robot to solve a task of interest
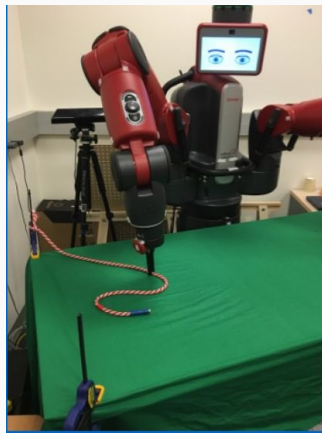
Wojciech Zaremba - OpenAI

# Overview

- Where to get rich, diverse data for robotics?

- How to obtain complex behaviors on robots?
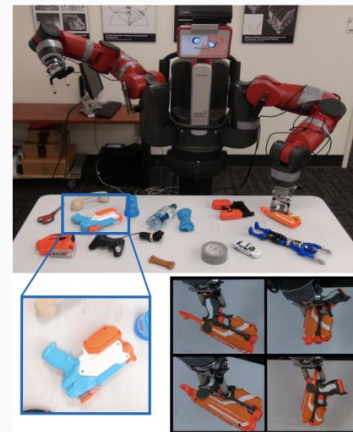
- How to convey the intent of the task to the robot?

Wojciech Zaremba - OpenAI

- ***Where to get rich, diverse data for robotics?***

- How to obtain complex behaviors on robots?

- How to convey the intent of the task to the robot?

# Data from Physical Robots


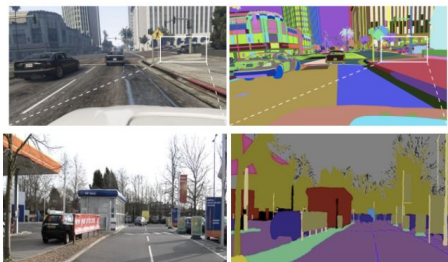Levine et al. 2016


Nair, Chen, Agrawal, et al. 2016


Pinto et al. 2016

- Real data is closest to reality
- Real data is expensive
- Hard to obtain large diversity

Wojciech Zaremba - OpenAI

# Data from Simulation

## Maximally Realistic Simulation



Richter et al. 2016



James et al. 2016

## Fine-tuning



Rusu et al. 2016 (progressive nets)

## Domain Adaptation



Tzeng et al. 2014

Wojciech Zaremba - OpenAI

Do we ever need real data ?

Does our simulation have to be photorealistic ?

Wojciech Zaremba - OpenAI

Do we ever need real data ?

Does our simulation have to be photorealistic ?

*Or, if the model sees enough simulated variation, might the real world may look like the other variation?*

Wojciech Zaremba - OpenAI

- Quadcopter collision avoidance

- ~40-50% of 1000m trajectories are collision-free

**Fereshteh Sadeghi and Sergey Levine. (cad)ˆ2 rl: Real single-image flight without a single real image.**

Wojciech Zaremba - OpenAI

- Quadcopter collision avoidance

Can it be precise enough for manipulation ?

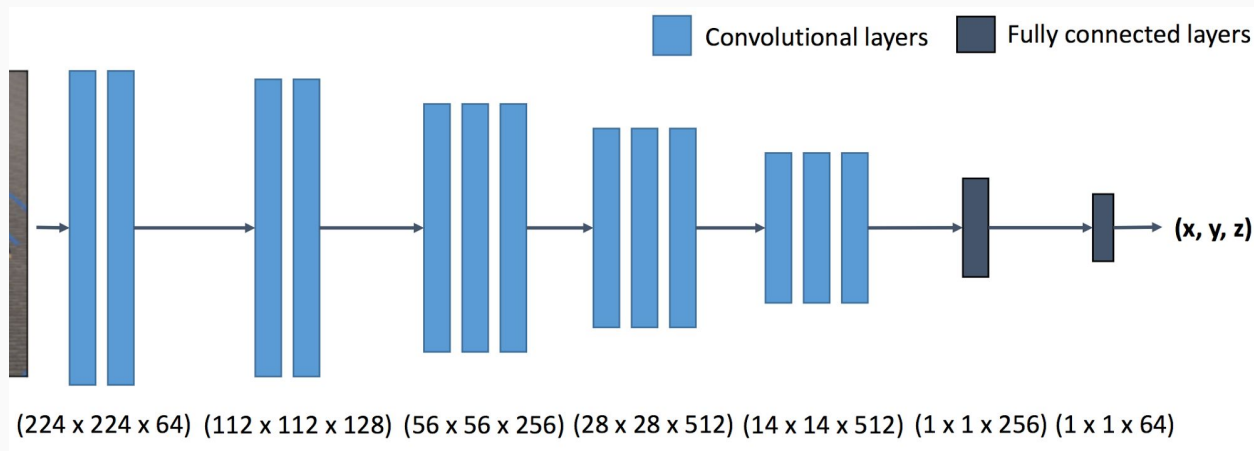Fereshteh Sadeghi and Sergey Levine. (cad)^2 rl: Real single-image flight without a single real image.

- Quadcopter collision avoidance

Can it be precise enough for manipulation ?
How realistic do textures need to be ?

Fereshteh Sadeghi and Sergey Levine. (cad)^2 rl: Real single-image flight without a single real image.

- Quadcopter collision avoidance

Can it be precise enough for manipulation ?
How realistic do textures need to be ?
Do we need pretraining on real data ?

Fereshteh Sadeghi and Sergey Levine. (cad)^2 rl: Real single-image flight without a single real image.

Wojciech Zaremba - OpenAI

Convolutional layers | Fully connected layers

(224 x 224 x 64) (112 x 112 x 128) (56 x 56 x 256) (28 x 28 x 512) (14 x 14 x 512) (1 x 1 x 256) (1 x 1 x 64)

$(x, y, z)$

Randomized 100k scenes
- lighting
- textures (checkerboards and solid)
- camera position

**By Josh Tobin, Rachel Fong, Alex Ray, Jonas Schneider, Wojciech Zaremba, Pieter Abbeel**

Wojciech Zaremba - OpenAI

4.5X SPEED

Wojciech Zaremba - OpenAI
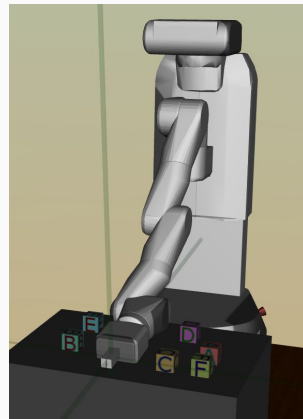
# How does it work?   More Data = Better

# More Textures = Better



Wojciech Zaremba - OpenAI

- Use multiple cameras, depth sensors and higher resolution images

- More randomization

- Apply to large number of tasks and complex generated works



Wojciech Zaremba - OpenAI

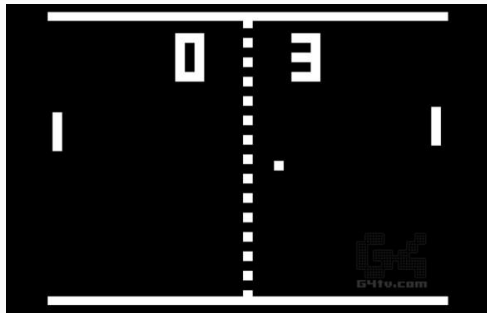- Where to get rich, diverse data for robotics ?
  Approach: Domain randomization



- ***How to obtain complex behaviors on robots ?***

- How to convey the intent of the task to the robot?
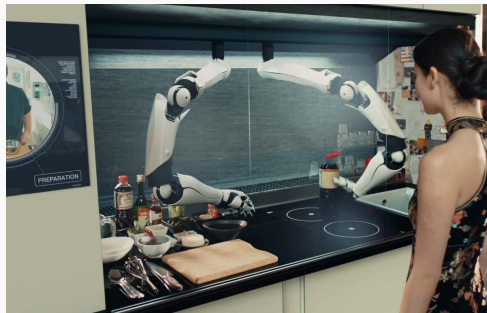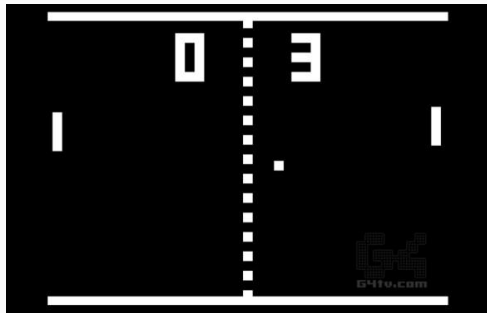
Wojciech Zaremba - OpenAI

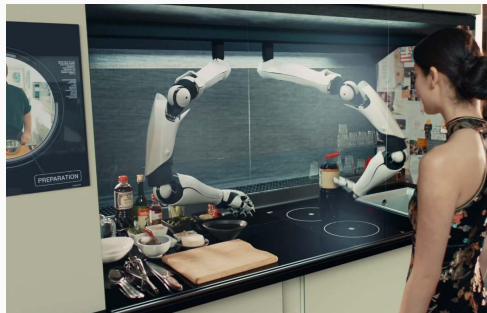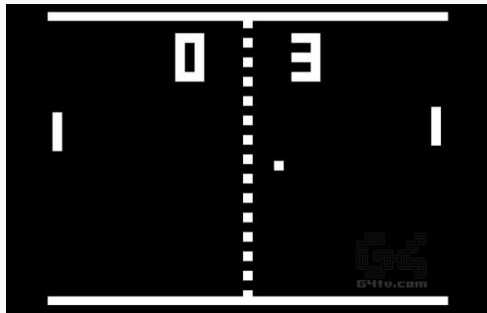- Initial DeepMind Atari results - 1 week of training

# Reinforcement Learning

- Initial DeepMind Atari results - 1 week of training



- We would like to train faster *and* significantly more complicated tasks



Wojciech Zaremba - OpenAI

- Initial DeepMind Atari results - 1 week of training



- We would like to train faster *and* significantly more complicated tasks



## Solution: Parallelization?

# Distributed Reinforcement Learning

- GOogle ReInforcement Learning Architecture (Gorila)
  [Nair et al, 2015]
  - Parallel acting
  - Distributed replay memory
  - Parallel learning
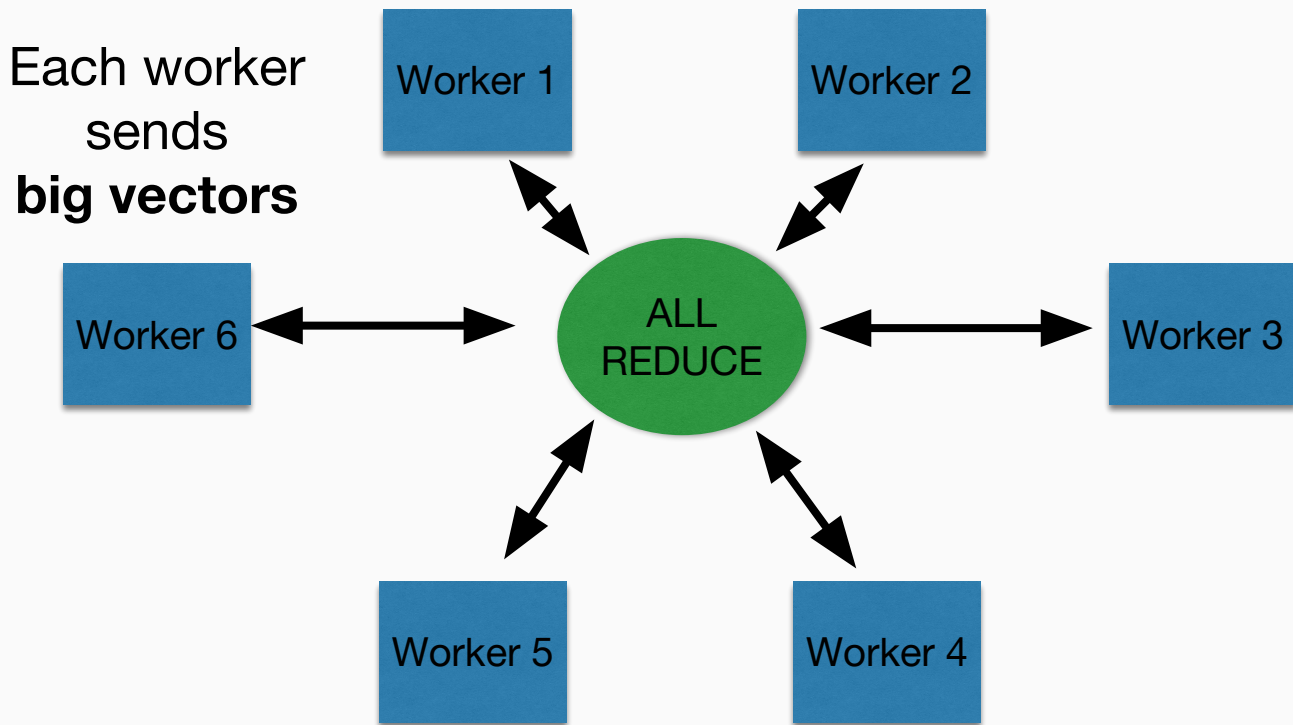  - Distributed neural network
  - Quite complex
- Asynchronous Advantage Actor Critic (A3C)
  [Mnih et al, 2016]

- Network communication is a bottleneck



Each worker sends **big vectors**

Worker 1, Worker 2, Worker 3, Worker 4, Worker 5, Worker 6

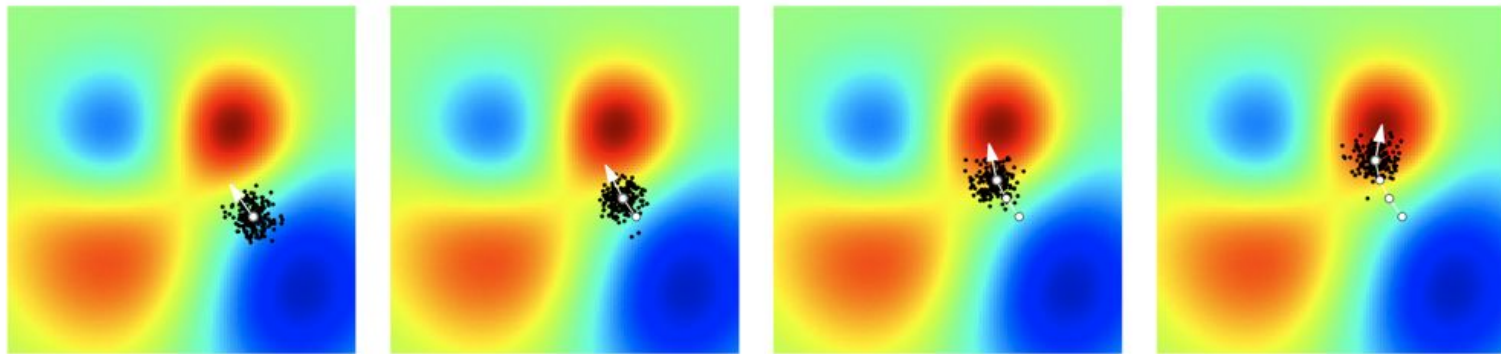ALL REDUCE

Wojciech Zaremba - OpenAI

# Do we have to communicate all parameters ?

Wojciech Zaremba - OpenAI

- Simplest algorithm imaginable:
  - Add perturbation to the parameters
  - If the result improves, keep the change
  - Repeat



**By Tim Salimans, Jonathan Ho, Peter Chen, Ilya Sutskever**

## **Classical RL**                    ## **Evolution**

- Sample *action* perturbations

## **Classical RL**

- Sample *action* perturbations

## **Evolution**

- Sample *parameter* perturbations

Wojciech Zaremba - OpenAI

## **Classical RL**

- Sample *action* perturbations
- Communicate *gradients/latest parameters*
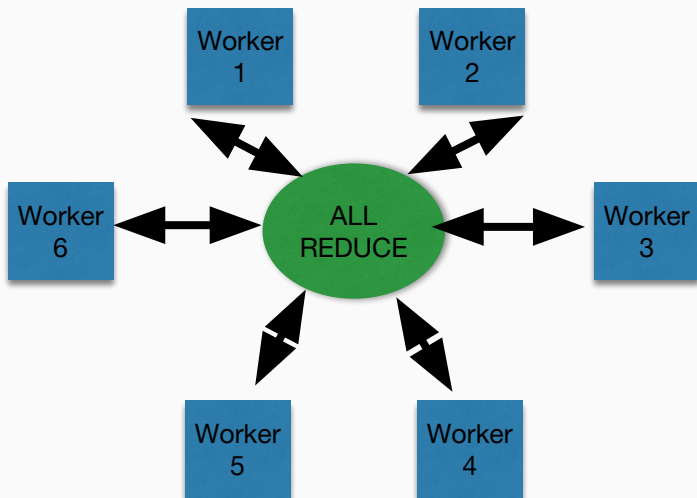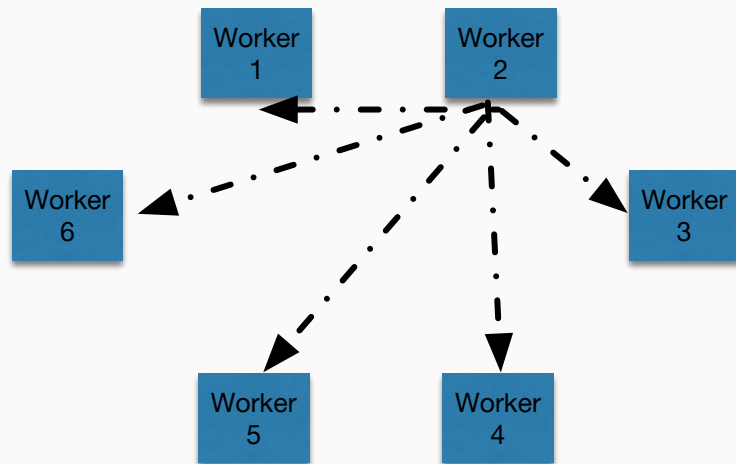
## **Evolution**

- Sample *parameter* perturbations

- Neural networks have millions of parameters
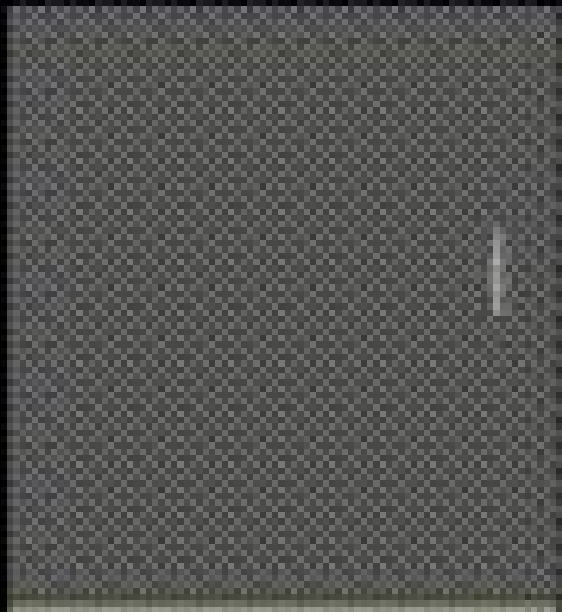- Folk Wisdom: There's no chance for this kind of random hillclimbing to succeed

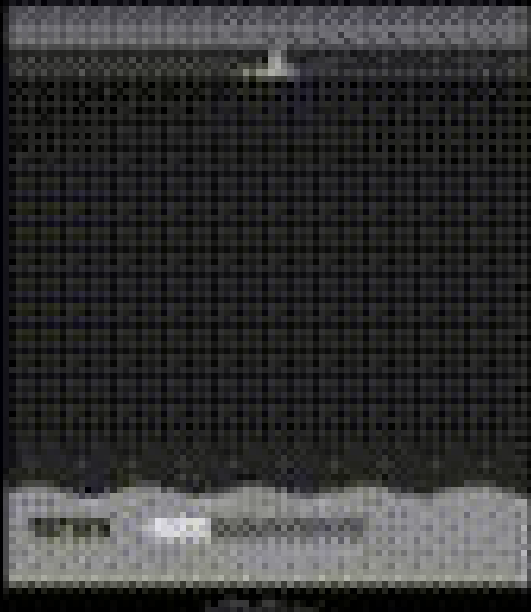Evolution Strategies is competitive with today's RL algorithms on standard benchmarks

# Evolution Atari Results

Pong    Seaquest    Beamrider
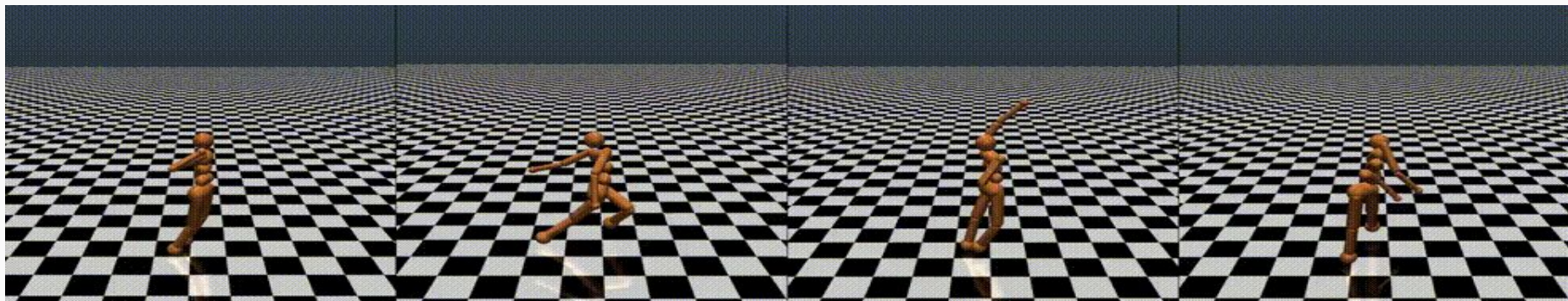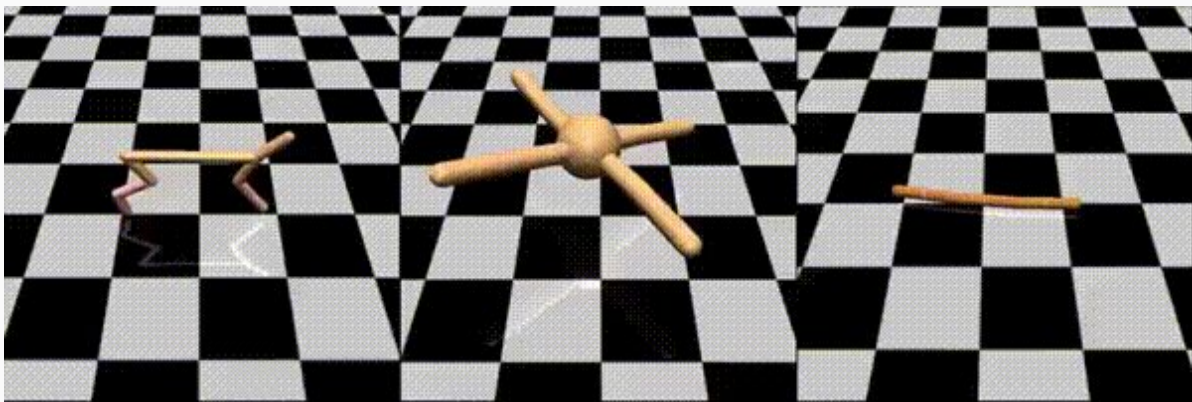
Wojciech Zaremba - OpenAI

# Evolution Atari Results

- Prior state-of-the-art on Atari in distributed RL: A3C [Mnih et al '16]
  - Training time 1 day

- Evolution Strategies
  - 1 hour with 720 cores matches A3C
  - 3x-10x more data
  - No backward pass
    - no need to store activations in memory
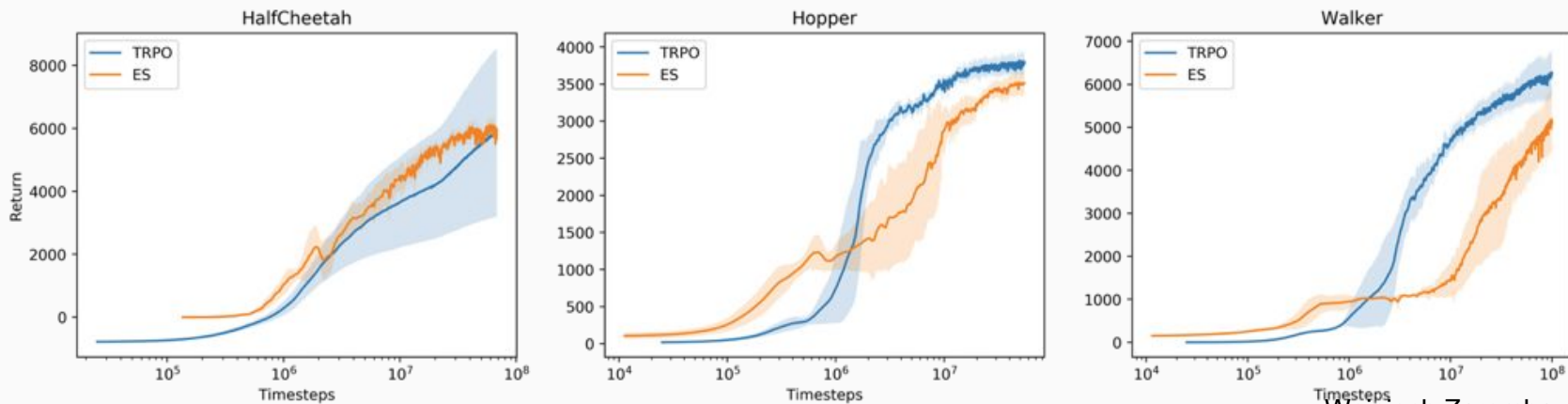    - reduces compute per episode by 2/3

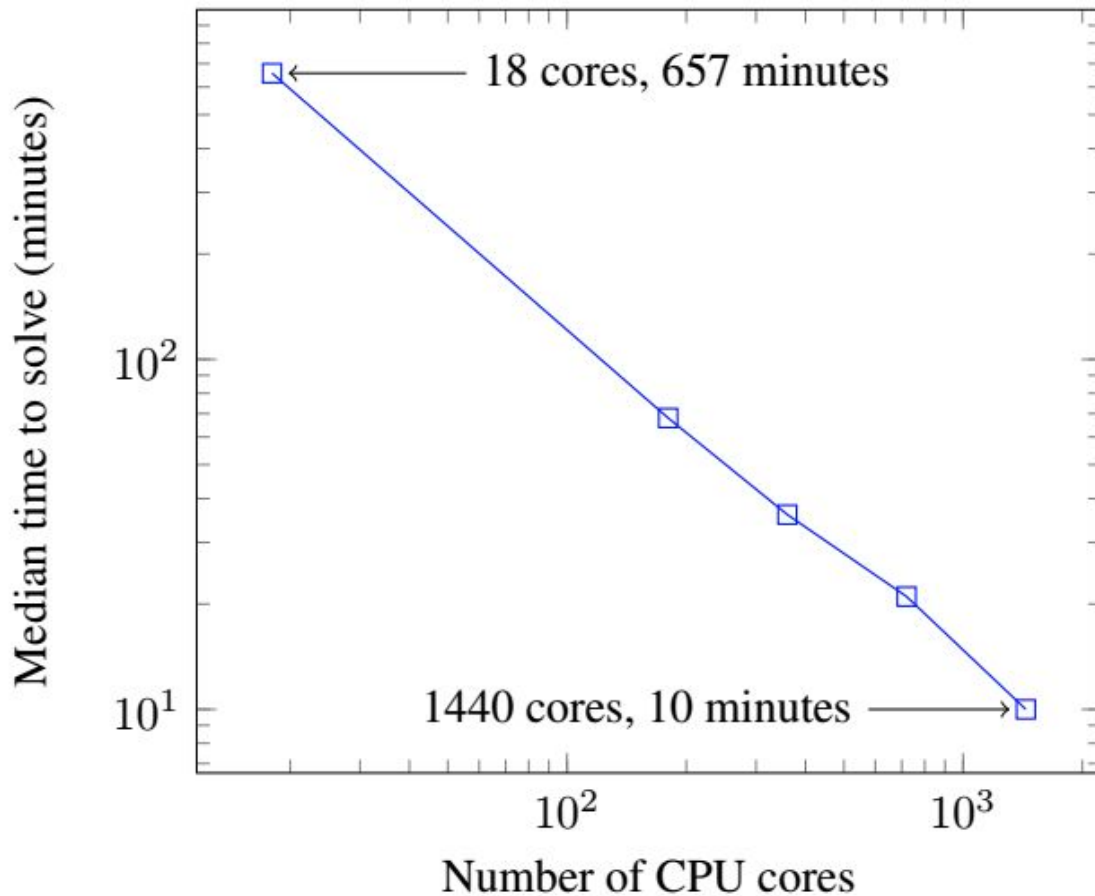Wojciech Zaremba - OpenAI

Wojciech Zaremba - OpenAI

- Evolution needs more data, but it achieves nearly the same result
- If we use 1440 cores, we need 10 minutes to solve the humanoid task, which takes 1 day with TRPO [Schulman et al., 2015] on a single machine



Wojciech Zaremba - OpenAI

Wojciech Zaremba - OpenAI
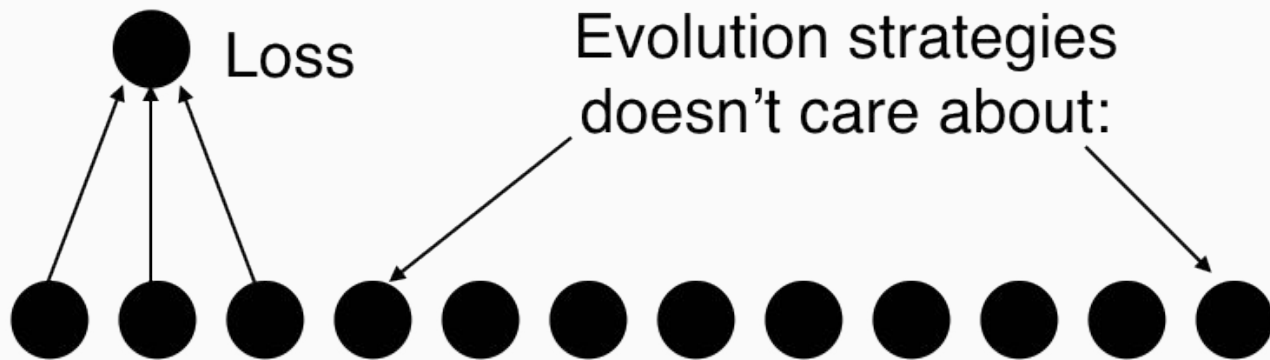
Fact: the speed of Evolution Strategies depends on the intrinsic dimensionality of the problem, not on the actual dimensionality

Loss

Evolution strategies doesn't care about:

- Evolution strategies *automatically discards* the irrelevant dimensions — even when they live on a complicated subspace!

Wojciech Zaremba - OpenAI

- Evolution Strategies was proposed in 1977 by Rechenberg & Eigen

- Entire journals devoted to Evolution, e.g. Evolutionary Computation Journal

- Showed that evolution is competitive with today's existing RL algorithms on standard RL benchmarks

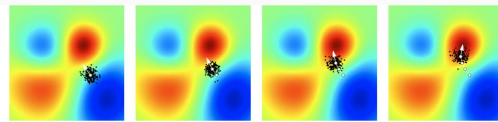- Showed that evolution parallelizes extremely well

- Where to get rich, diverse data for robotics ?
  Approach: Domain randomization



- How to obtain complex behaviors on robots?
  Approach: Evolution



- **How to convey the intent of the task to the robot ?**

Wojciech Zaremba - OpenAI

- Language seems to be one option
  - Limits robot to tasks involving words that it knows
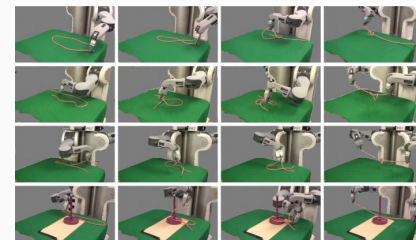
- Language seems to be one option
  - Limits robot to tasks involving words that it knows


- Alternative is to show the task

Wojciech Zaremba - OpenAI

- Abbeel, Coates, Ng: "Autonomous Helicopter Aerobatics through Apprenticeship Learning"



- Schulman et al. "Learning from Demonstrations Through the Use of Non-Rigid Registration"



- van den Berg et al. "Superhuman Performance of Surgical Tasks by Robots using Iterative Learning from Human-Guided Demonstrations"



Wojciech Zaremba - OpenAI

Prior work, proposes various imitation algorithms to learn from multiple demonstrations
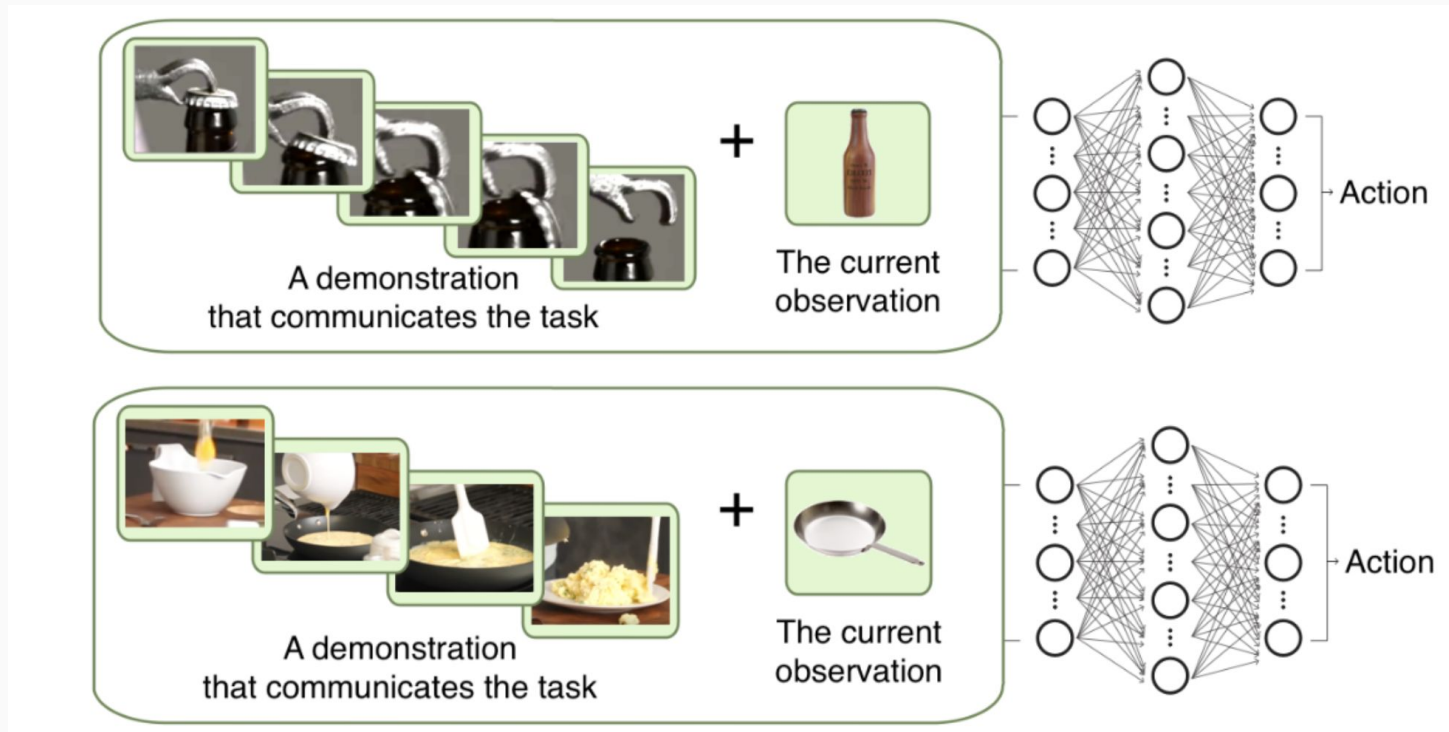
Wojciech Zaremba - OpenAI

Prior work, proposes various imitation algorithms to learn from multiple demonstrations

**Instead, we *learn an imitation algorithm* that imitates based on a single demonstration**
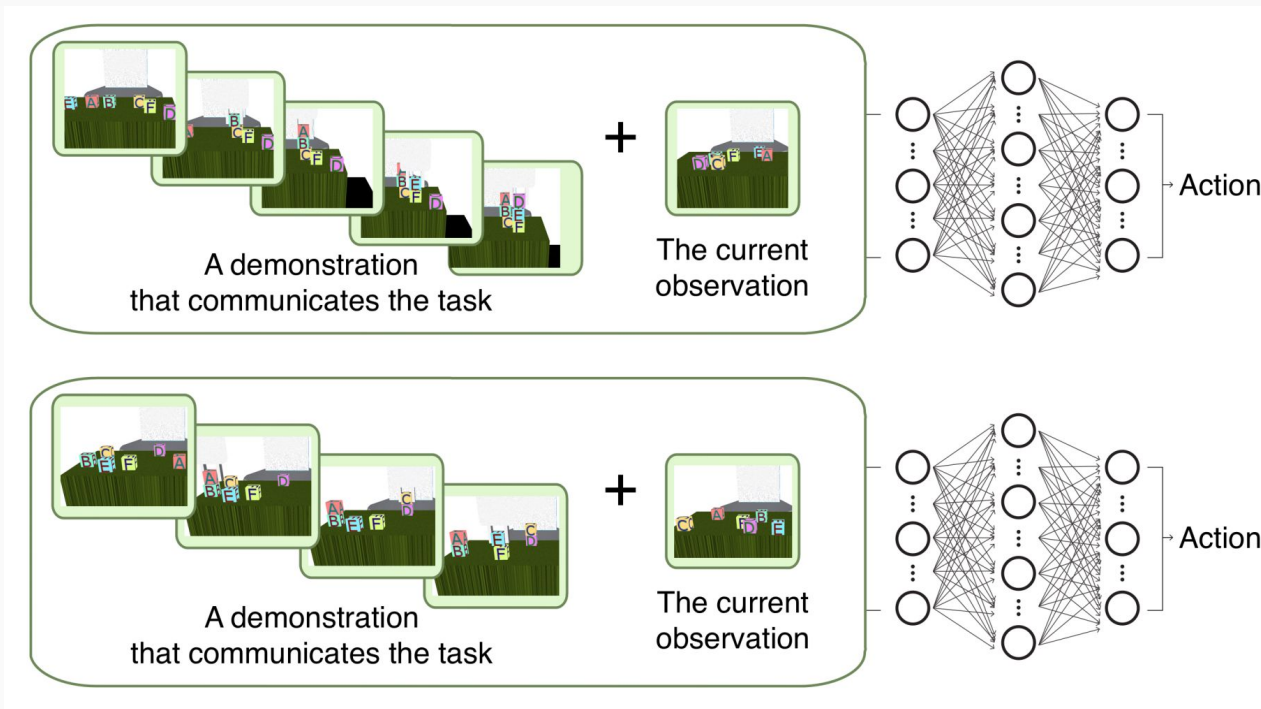
- Sample a task
- Sample an input demonstration from the task
- Sample a target demonstration from the task (in different initial condition)
- Train network given input demonstration to predict the target demonstration

# Learning the Imitation Algorithm - Our Setup



**By Rocky Duan, Marcin Andrychowicz, Bradly Stadie, Jonathan Ho, Jonas Schneider, Ilya Sutskever, Pieter Abbeel, Wojciech Zaremba**

- Each task is specified by the desired final layout

- Example: *abcd*
  - "Place *c* on top of *d*,
    place b on top of *c*,
    place a on top of *b*"

- Each task is specified by the desired final layout

- Example: *abc def gh*
  - "Place *b* on top of *c*; a on top of *b*;"
  - "Place *e* on top of *f*; d on top of *e*;"
  - "Place *g* on top of *h*."

Size of dataset

- Number of blocks vary from 2 to 10
- 183 distinct tasks, not counting equivalent permutations
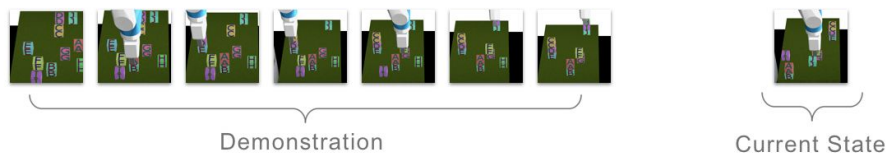- 140 tasks for training, and 43 tasks for testing

- Works with demonstrations of different size

- Works with very long demonstrations
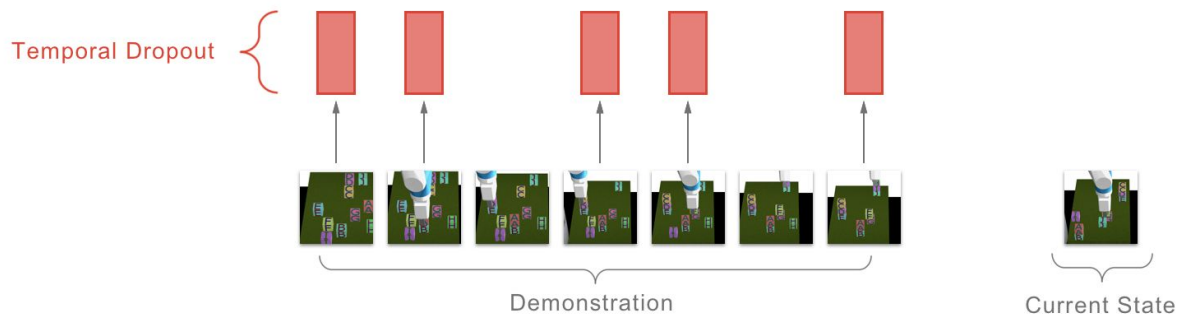
- Works with variable number of blocks

Demonstration

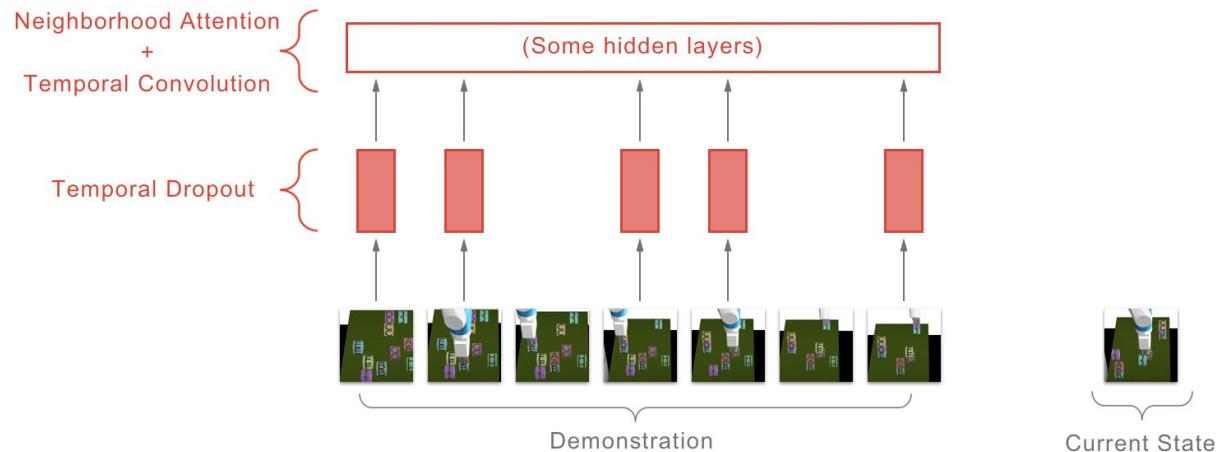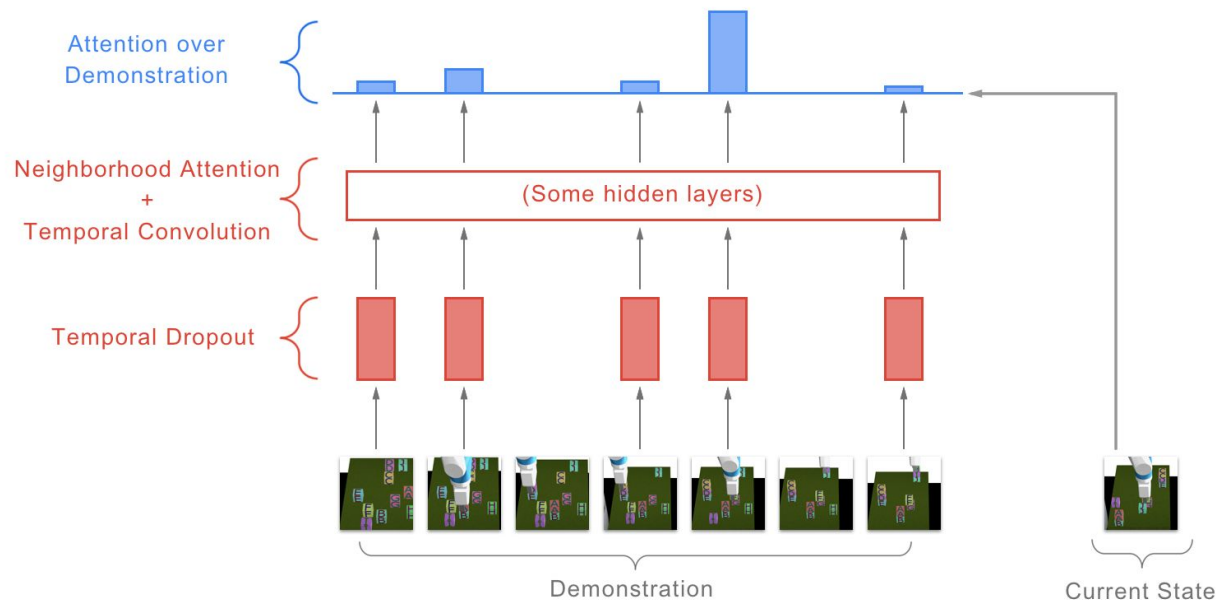Current State

Temporal Dropout

Demonstration

Current State

Neighborhood Attention
+
Temporal Convolution

(Some hidden layers)

Temporal Dropout

Demonstration

Current State

Attention over Demonstration

Neighborhood Attention
+
Temporal Convolution

(Some hidden layers)

Temporal Dropout

Demonstration

Current State

Context Embedding

Attention over Demonstration

Neighborhood Attention + Temporal Convolution

(Some hidden layers)

Temporal Dropout

Demonstration

Current State

Context Embedding

Attention over Current State

Block# A B C D E F G H I J

Attention over Demonstration

Neighborhood Attention + Temporal Convolution

(Some hidden layers)

Temporal Dropout

Demonstration

Current State

# Architecture

Wojciech Zaremba - OpenAI
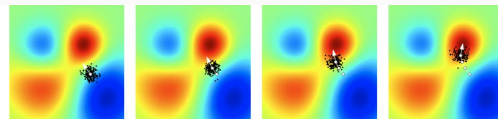
# One-Shot Imitation — Numerical Results

- Where to get rich, diverse data for robotics ?
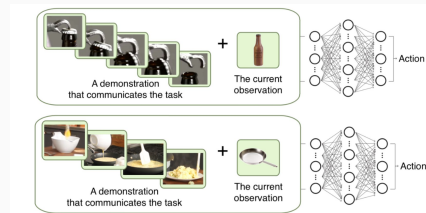  Our approach: Domain randomization

- How to obtain complex behaviors on robots?
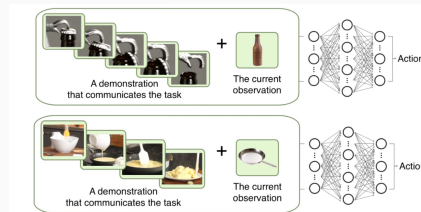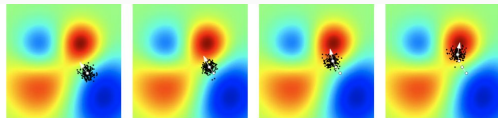  An approach: Evolution

- How to convey the intent of the task to the robot ? An approach: One-shot imitation

Wojciech Zaremba - OpenAI

# Summary



diverse data + scalable training + one-shot imitation
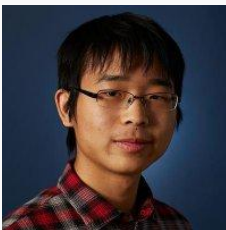
## Still many limitations:

- Methods have not been evaluated on complex tasks such as cooking or cleaning
- Methods might break when simulated data oversimplifies real world
- Parallel gripper is relatively simple to control even without neural networks
- So far, all experiments are just a proof of concept

Wojciech Zaremba - OpenAI
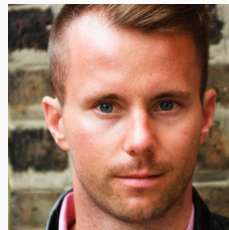
# The robotics team


Marcin Andrychowicz


Rocky Duan


Bradly Stadie


Filip Wolski


Alex Ray


Jonas Schneider


Rachel Fong


Peter Welinder


Ankur Handa


Josh Tobin


Lukas Biewald


Pieter Abbeel


Erika Reinhardt


Bob McGrew


Vikash Kumar

# Thank you

Wojciech Zaremba - OpenAI

# Ablation for one-shot