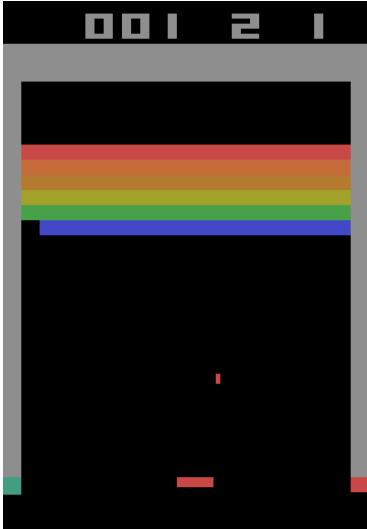


# Benchmarking Deep Reinforcement Learning for Continuous Control

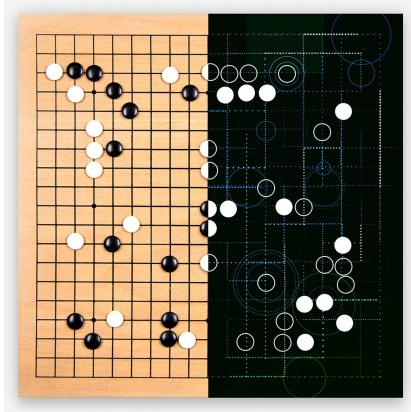
Yan (Rocky) Duan<sup>†</sup>, Xi Chen<sup>†</sup>, Rein Houthooft<sup>†‡</sup>, John Schulman<sup>†§</sup>, Pieter Abbeel<sup>†</sup>

<sup>†</sup>Berkeley Artificial Intelligence Research (BAIR) laboratory, <sup>‡</sup>Ghent University, <sup>§</sup>OpenAI

# Deep Reinforcement Learning



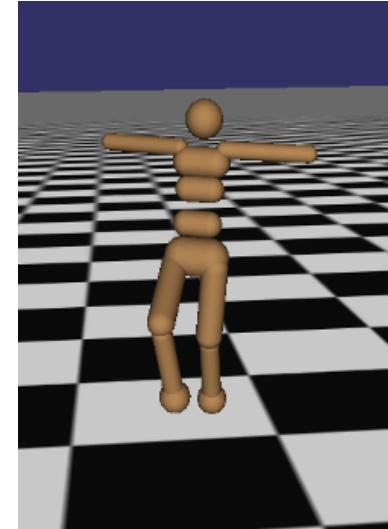
[Mnih et al. 2013; Guo et al. 2014; Mnih et al. 2015; Oh et al. 2015; Nair et al. 2015; Van Hasselt et al. 2015; Schulman et al. 2015; Parisotto et al. 2015; Wang et al. 2015; Kulkarni et al. 2016]



[Silver et al. 2016; Tian et al. 2016; Maddison et al. 2014; Clark et al. 2015]



[Levine et al. 2015; Finn et al. 2016; Watter et al. 2015; Tzeng et al. 2015]



[Schulman et al. 2015; Lillicrap et al. 2015; Heess et al. 2015; Gu et al. 2016]

# Deep Reinforcement Learning

---

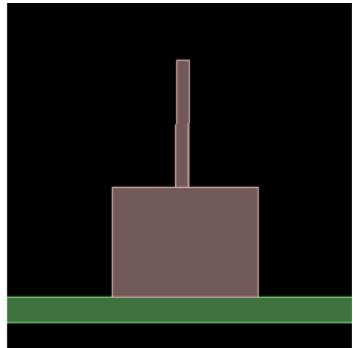
- High-dimensional state space + discrete actions:
  - Arcade Learning Environment (ALE) [Bellemare et al., 2012], Minecraft [Tessler et al., 2016; Abel et al., 2016; Oh et al., 2016]...
- High-dimensional continuous actions
  - ?

# Proposed Benchmark

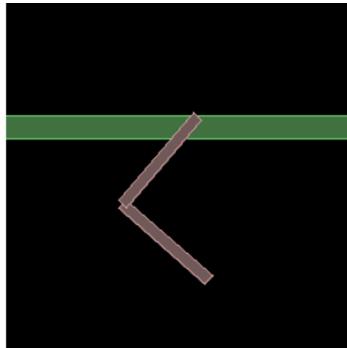
---

- 5 classic tasks
- 7 locomotion tasks
- 15 partially observable tasks
- 4 hierarchical tasks
- Also: reference implementations of 8 algorithms
- ***Everything open source***

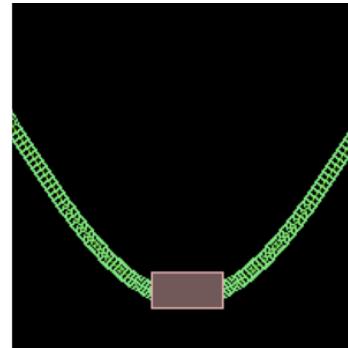
# Classic Tasks



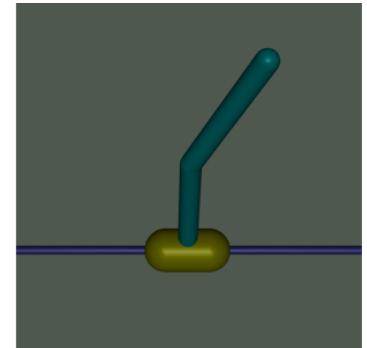
Cart Pole Balancing  
+  
Inverted Pendulum



Acrobot

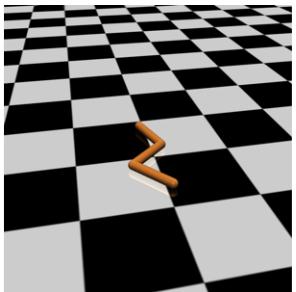


Mountain Car

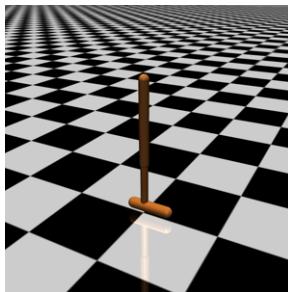


Double Inverted Pendulum

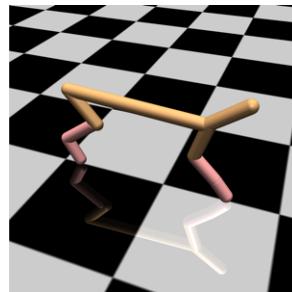
# Locomotion Tasks



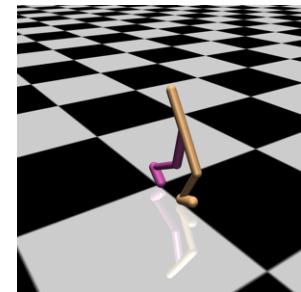
Swimmer



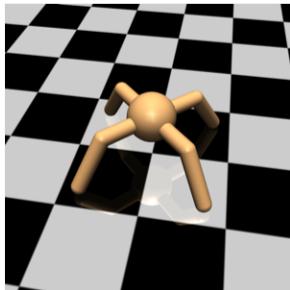
Hopper



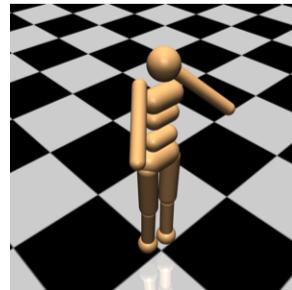
Half Cheetah



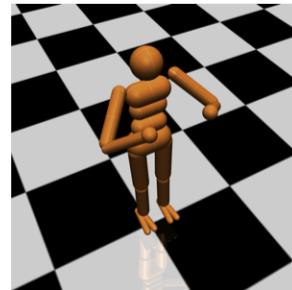
Walker



Ant



Simplified Humanoid



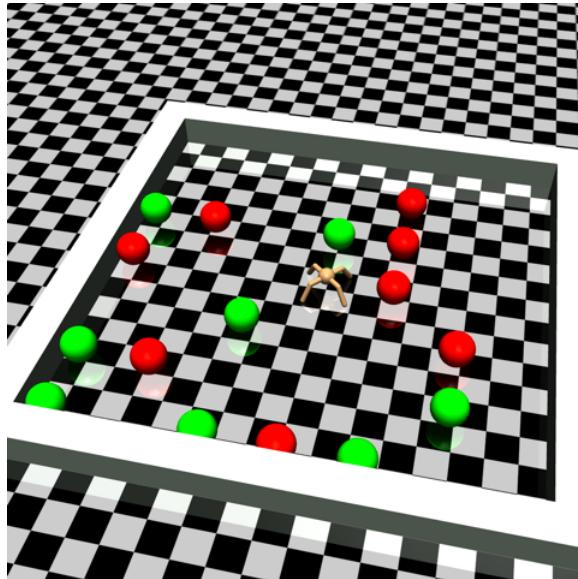
Full Humanoid

# Partially Observable Tasks

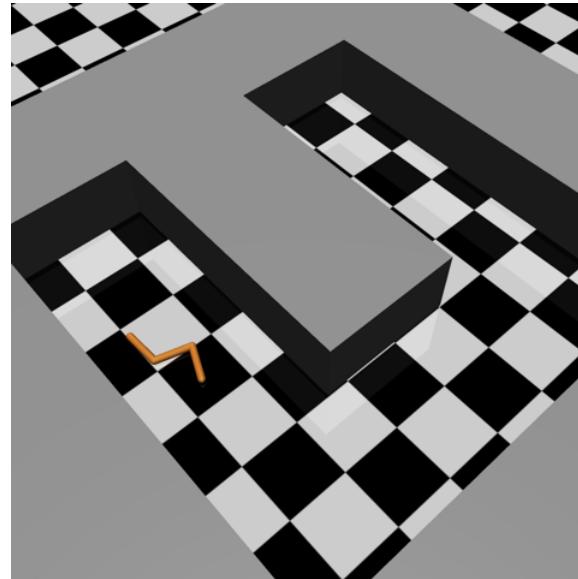
---

- Limited Sensors (LS)
- Noisy Observations and Delayed Actions (NO)
- System Identification (SI)

# Hierarchical Tasks



Ant + Gather  
(Also: Swimmer + Gather)



Swimmer + Maze  
(Also: Ant + Maze)

# Algorithms

---

- REINFORCE [Williams, 1992] + ADAM [Kingma & Ba, 2014]
- TNPG [Kakade, 2002; Peters et al., 2003; Bagnell & Schneider, 2003; Schulman et al., 2015]
- RWR [Peters & Schaal, 2007; Kober & Peters, 2009]
- REPS [Peters et al., 2010]
- TRPO [Schulman et al. 2015]
- CEM [Rubinstein, 1999; Szita & Lorincz, 2006]
- CMA-ES [Hansen & Ostermeier, 2001]
- DDPG [Silver et al., 2014; Lillicrap et al., 2015]

# Batch, Gradient-Based Algorithms

- REINFORCE [Williams, 1992] + ADAM [Kingma & Ba, 2014]

$$\hat{g} \leftarrow \text{PolicyGradient}$$

$$\Delta\theta \leftarrow \text{AdamUpdate}(\theta, \hat{g})$$

- NPG [Kakade, 2002; Peters et al., 2003; Bagnell & Schneider, 2003]

$$\Delta\theta \propto \hat{F}(\theta)^{-1} \hat{g}$$

- TNPG [Schulman et al., 2015]      Compute descent direction via truncated conjugate gradient.
- TRPO [Schulman et al., 2015]      Solve constrained optimization problem

$$\Delta\theta = \arg \max_{\Delta\theta} \eta(\theta + \Delta\theta)$$

$$\text{s.t. } D_{KL}(\pi_\theta \| \pi_{\theta + \Delta\theta}) \leq \delta$$

# Batch, Gradient-Based Algorithms

---

- RWR [Peters & Schaal, 2007; Kober & Peters, 2009]

Formulated as Expectation-Maximization problem.

- REPS [Peters et al., 2010]

Limits the loss of information in the state-action distribution.

Formulated as dual optimization problem.

# Batch, Gradient-Free Algorithms

---

- CEM [Rubinstein, 1999; Szita & Lorincz, 2006]

Black-box algorithm which maintains and updates a factored Gaussian distribution over the parameters.

- CMA-ES [Hansen & Ostermeier, 2001]

Similar to CEM, but maintains a full covariance matrix.

# Online Algorithms

---

- DDPG [Silver et al., 2014; Lillicrap et al., 2015]

$$\Delta\theta \propto \frac{\partial}{\partial a} \hat{Q}^\pi(s, a) \Big|_{a=\pi_\theta(s)} \frac{\partial}{\partial \theta} \pi_\theta(s)$$

# Policy Representation

---

- All Markov environments
  - 3-layer MLP, 100-50-25 units with relu nonlinearity
- Partially observable tasks
  - Single-layer LSTM with 32 hidden units
- DDPG: same architecture as in [Lillicrap et al. 2015]

# Results

Task	Random	REINFORCE	TNPG	RWR	REPS	TRPO	CEM	CMA-ES	DDPG
Cart-Pole Balancing	$77.1 \pm 0.0$	$4693.7 \pm 14.0$	<b><math>3986.4 \pm 748.9</math></b>	<b><math>4861.5 \pm 12.3</math></b>	$565.6 \pm 137.6$	<b><math>4869.8 \pm 37.6</math></b>	$4815.4 \pm 4.8$	$2440.4 \pm 568.3$	$4634.4 \pm 87.8$
Inverted Pendulum*	$-153.4 \pm 0.2$	$13.4 \pm 18.0$	<b><math>209.7 \pm 55.5</math></b>	$84.7 \pm 13.8$	$-113.3 \pm 4.6$	<b><math>247.2 \pm 76.1</math></b>	$38.2 \pm 25.7$	$-40.1 \pm 5.7$	$40.0 \pm 244.6$
Mountain Car	$-415.4 \pm 0.0$	$-67.1 \pm 1.0$	<b><math>-66.5 \pm 4.5</math></b>	$-79.4 \pm 1.1$	$-275.6 \pm 166.3$	<b><math>-61.7 \pm 0.9</math></b>	$-66.0 \pm 2.4$	$-85.0 \pm 7.7$	$-288.4 \pm 170.3$
Acrobot	$-1904.5 \pm 1.0$	$-508.1 \pm 91.0$	$-395.8 \pm 121.2$	$-352.7 \pm 35.9$	$-1001.5 \pm 10.8$	$-326.0 \pm 24.4$	$-436.8 \pm 14.7$	$-785.6 \pm 13.1$	<b><math>-223.6 \pm 5.8</math></b>
Double Inverted Pendulum*	$149.7 \pm 0.1$	$4116.5 \pm 65.2$	<b><math>4455.4 \pm 37.6</math></b>	$3614.8 \pm 368.1$	$446.7 \pm 114.8$	<b><math>4412.4 \pm 50.4</math></b>	$2566.2 \pm 178.9$	$1576.1 \pm 51.3$	$2863.4 \pm 154.0$
Swimmer*	$-1.7 \pm 0.1$	$92.3 \pm 0.1$	<b><math>96.0 \pm 0.2</math></b>	$60.7 \pm 5.5$	$3.8 \pm 3.3$	<b><math>96.0 \pm 0.2</math></b>	$68.8 \pm 2.4$	$64.9 \pm 1.4$	$85.8 \pm 1.8$
Hopper	$8.4 \pm 0.0$	$714.0 \pm 29.3$	<b><math>1155.1 \pm 57.9</math></b>	$553.2 \pm 71.0$	$86.7 \pm 17.6$	<b><math>1183.3 \pm 150.0</math></b>	$63.1 \pm 7.8$	$20.3 \pm 14.3$	$267.1 \pm 43.5$
2D Walker	$-1.7 \pm 0.0$	$506.5 \pm 78.8$	<b><math>1382.6 \pm 108.2</math></b>	$136.0 \pm 15.9$	$-37.0 \pm 38.1$	<b><math>1353.8 \pm 85.0</math></b>	$84.5 \pm 19.2$	$77.1 \pm 24.3$	$318.4 \pm 181.6$
Half-Cheetah	$-90.8 \pm 0.3$	$1183.1 \pm 69.2$	<b><math>1729.5 \pm 184.6</math></b>	$376.1 \pm 28.2$	$34.5 \pm 38.0$	<b><math>1914.0 \pm 120.1</math></b>	$330.4 \pm 274.8$	$441.3 \pm 107.6$	<b><math>2148.6 \pm 702.7</math></b>
Ant*	$13.4 \pm 0.7$	$548.3 \pm 55.5$	<b><math>706.0 \pm 127.7</math></b>	$37.6 \pm 3.1$	$39.0 \pm 9.8$	<b><math>730.2 \pm 61.3</math></b>	$49.2 \pm 5.9$	$17.8 \pm 15.5$	$326.2 \pm 20.8$
Simple Humanoid	$41.5 \pm 0.2$	$128.1 \pm 34.0$	<b><math>255.0 \pm 24.5</math></b>	$93.3 \pm 17.4$	$28.3 \pm 4.7$	<b><math>269.7 \pm 40.3</math></b>	$60.6 \pm 12.9$	$28.7 \pm 3.9$	$99.4 \pm 28.1$
Full Humanoid	$13.2 \pm 0.1$	$262.2 \pm 10.5$	<b><math>288.4 \pm 25.2</math></b>	$46.7 \pm 5.6$	$41.7 \pm 6.1$	<b><math>287.0 \pm 23.4</math></b>	$36.9 \pm 2.9$	N/A ± N/A	$119.0 \pm 31.2$
Cart-Pole Balancing (LS)*	$77.1 \pm 0.0$	$420.9 \pm 265.5$	<b><math>945.1 \pm 27.8</math></b>	$68.9 \pm 1.5$	$898.1 \pm 22.1$	<b><math>960.2 \pm 46.0</math></b>	$227.0 \pm 223.0$	$68.0 \pm 1.6$	
Inverted Pendulum (LS)	$-122.1 \pm 0.1$	$-13.4 \pm 3.2$	<b><math>0.7 \pm 6.1</math></b>	$-107.4 \pm 0.2$	$-87.2 \pm 8.0$	<b><math>4.5 \pm 4.1</math></b>	$-81.2 \pm 33.2$	$-62.4 \pm 3.4$	
Mountain Car (LS)	$-83.0 \pm 0.0$	$-81.2 \pm 0.6$	<b><math>-65.7 \pm 9.0</math></b>	$-81.7 \pm 0.1$	$-82.6 \pm 0.4$	<b><math>-64.2 \pm 9.5</math></b>	<b><math>-68.9 \pm 1.3</math></b>	<b><math>-73.2 \pm 0.6</math></b>	
Acrobot (LS)*	$-393.2 \pm 0.0$	$-128.9 \pm 11.6$	<b><math>-84.6 \pm 2.9</math></b>	$-235.9 \pm 5.3$	$-379.5 \pm 1.4$	<b><math>-83.3 \pm 9.9</math></b>	$-149.5 \pm 15.3$	$-159.9 \pm 7.5$	
Cart-Pole Balancing (NO)*	$101.4 \pm 0.1$	$616.0 \pm 210.8$	<b><math>916.3 \pm 23.0</math></b>	$93.8 \pm 1.2$	$99.6 \pm 7.2$	$606.2 \pm 122.2$	$181.4 \pm 32.1$	$104.4 \pm 16.0$	
Inverted Pendulum (NO)	$-122.2 \pm 0.1$	$6.5 \pm 1.1$	<b><math>11.5 \pm 0.5</math></b>	$-110.0 \pm 1.4$	$-119.3 \pm 4.2$	<b><math>10.4 \pm 2.2</math></b>	$-55.6 \pm 16.7$	$-80.3 \pm 2.8$	
Mountain Car (NO)	$-83.0 \pm 0.0$	$-74.7 \pm 7.8$	<b><math>-64.5 \pm 8.6</math></b>	$-81.7 \pm 0.1$	$-82.9 \pm 0.1$	<b><math>-60.2 \pm 2.0</math></b>	$-67.4 \pm 1.4$	$-73.5 \pm 0.5$	
Acrobot (NO)*	$-393.5 \pm 0.0$	<b><math>-186.7 \pm 31.3</math></b>	<b><math>-164.5 \pm 13.4</math></b>	$-233.1 \pm 0.4$	$-258.5 \pm 14.0$	<b><math>-149.6 \pm 8.6</math></b>	$-213.4 \pm 6.3$	$-236.6 \pm 6.2$	
Cart-Pole Balancing (SI)*	$76.3 \pm 0.1$	$431.7 \pm 274.1$	<b><math>980.5 \pm 7.3</math></b>	$69.0 \pm 2.8$	$702.4 \pm 196.4$	<b><math>980.3 \pm 5.1</math></b>	$746.6 \pm 93.2$	$71.6 \pm 2.9$	
Inverted Pendulum (SI)	$-121.8 \pm 0.2$	$-5.3 \pm 5.6$	<b><math>14.8 \pm 1.7</math></b>	$-108.7 \pm 4.7$	$-92.8 \pm 23.9$	<b><math>14.1 \pm 0.9</math></b>	$-51.8 \pm 10.6$	$-63.1 \pm 4.8$	
Mountain Car (SI)	$-82.7 \pm 0.0$	$-63.9 \pm 0.2$	<b><math>-61.8 \pm 0.4</math></b>	$-81.4 \pm 0.1$	$-80.7 \pm 2.3$	<b><math>-61.6 \pm 0.4</math></b>	$-63.9 \pm 1.0$	$-66.9 \pm 0.6$	
Acrobot (SI)*	$-387.8 \pm 1.0$	<b><math>-169.1 \pm 32.3</math></b>	<b><math>-156.6 \pm 38.9</math></b>	$-233.2 \pm 2.6$	$-216.1 \pm 7.7$	<b><math>-170.9 \pm 40.3</math></b>	$-250.2 \pm 13.7$	$-245.0 \pm 5.5$	
Swimmer + Gathering	$0.0 \pm 0.0$	$0.0 \pm 0.0$	$0.0 \pm 0.0$	$0.0 \pm 0.0$	$0.0 \pm 0.0$	$0.0 \pm 0.0$	$0.0 \pm 0.0$	$0.0 \pm 0.0$	$0.0 \pm 0.0$
Ant + Gathering	$-5.8 \pm 5.0$	$-0.1 \pm 0.1$	$-0.4 \pm 0.1$	$-5.5 \pm 0.5$	$-6.7 \pm 0.7$	$-0.4 \pm 0.0$	$-4.7 \pm 0.7$	N/A ± N/A	$-0.3 \pm 0.3$
Swimmer + Maze	$0.0 \pm 0.0$	$0.0 \pm 0.0$	$0.0 \pm 0.0$	$0.0 \pm 0.0$	$0.0 \pm 0.0$	$0.0 \pm 0.0$	$0.0 \pm 0.0$	$0.0 \pm 0.0$	$0.0 \pm 0.0$
Ant + Maze	$0.0 \pm 0.0$	$0.0 \pm 0.0$	$0.0 \pm 0.0$	$0.0 \pm 0.0$	$0.0 \pm 0.0$	$0.0 \pm 0.0$	$0.0 \pm 0.0$	N/A ± N/A	$0.0 \pm 0.0$

# Results

	TNPG	TRPO
Inverted Pendulum	<b>209.7±55.5</b>	<b>247.2±76.1</b>
Mountain Car	<b>-66.5±4.5</b>	<b>-61.7±0.9</b>
Swimmer	<b>96±0.2</b>	<b>96±0.2</b>
Half-Cheetah	<b>1729.5±184.6</b>	<b>1914±120.1</b>
Cart-Pole Balancing (LS)	<b>945.1±27.8</b>	<b>960.2±46</b>
Mountain Car (NO)	<b>-64.5±8.6</b>	<b>-60.2±2</b>
Inverted Pendulum (SI)	<b>14.8±1.7</b>	<b>14.1±0.9</b>

# Results

	<b>REINFORCE</b>	<b>TRPO</b>
Cart-Pole Balancing	4693.7±14	<b>4869.8±37.6</b>
Inverted Pendulum	<b>13.4±18</b>	<b>247.2±76.1</b>
Mountain Car	-67.1±1	<b>-61.7±0.9</b>
Double Inverted Pendulum	4116.5±65.2	<b>4412.4±50.4</b>
Swimmer	92.3±0.1	<b>96±0.2</b>
Hopper	<b>714±29.3</b>	<b>1183.3±150</b>
Half-Cheetah	<b>1183.1±69.2</b>	<b>1914±120.1</b>

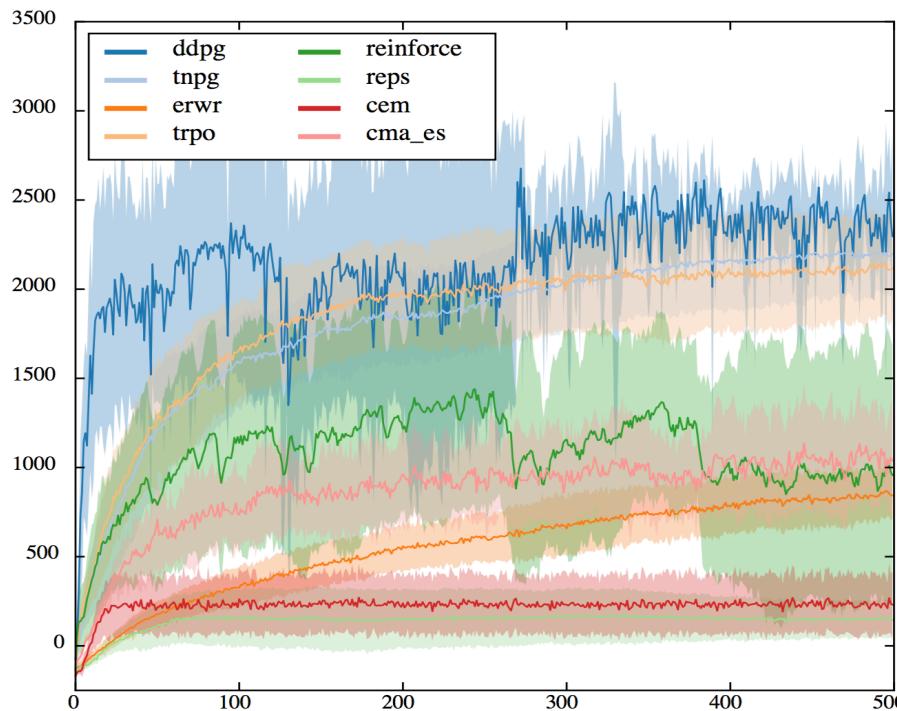
# Results

	<b>CEM</b>	<b>TRPO</b>
Cart-Pole Balancing	$4815.4 \pm 4.8$	<b><math>4869.8 \pm 37.6</math></b>
Mountain Car	$-66 \pm 2.4$	<b><math>-61.7 \pm 0.9</math></b>
Swimmer	$68.8 \pm 2.4$	<b><math>96 \pm 0.2</math></b>
Half-Cheetah	<b><math>330.4 \pm 274.8</math></b>	<b><math>1914 \pm 120.1</math></b>

# Results

	DDPG	TRPO
Cart-Pole Balancing	$4634.4 \pm 87.8$	<b><math>4869.8 \pm 37.6</math></b>
Inverted Pendulum	<b><math>40 \pm 244.6</math></b>	<b><math>247.2 \pm 76.1</math></b>
Acrobot	<b><math>-223.6 \pm 5.8</math></b>	<b><math>-326 \pm 24.4</math></b>
Hopper	<b><math>267.1 \pm 43.5</math></b>	<b><math>1183.3 \pm 150</math></b>
Half-Cheetah	<b><math>2148.6 \pm 702.7</math></b>	$1914 \pm 120.1$

# Results



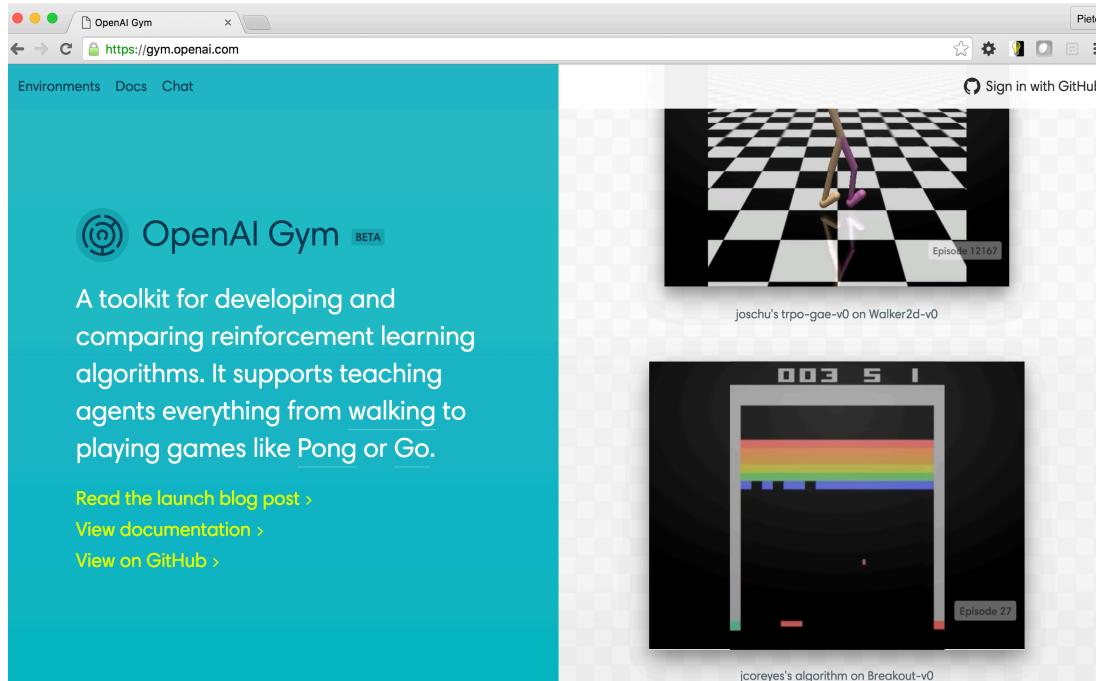
# Results

---

	TNPG	TRPO	DDPG	Others
Best Performance	40%	52%	8%	0%

# Results

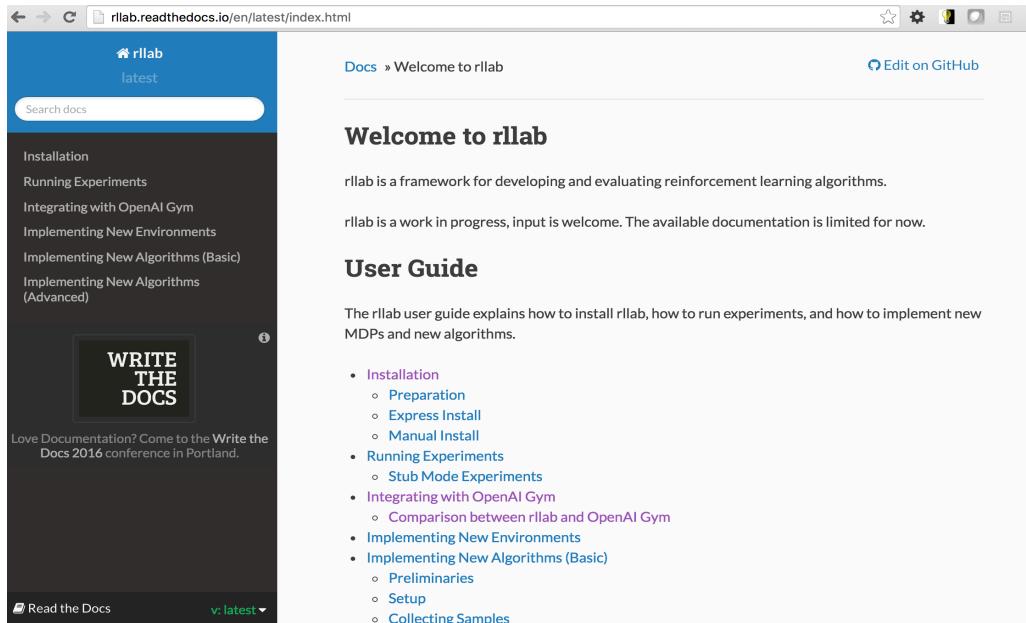
# OpenAI Gym



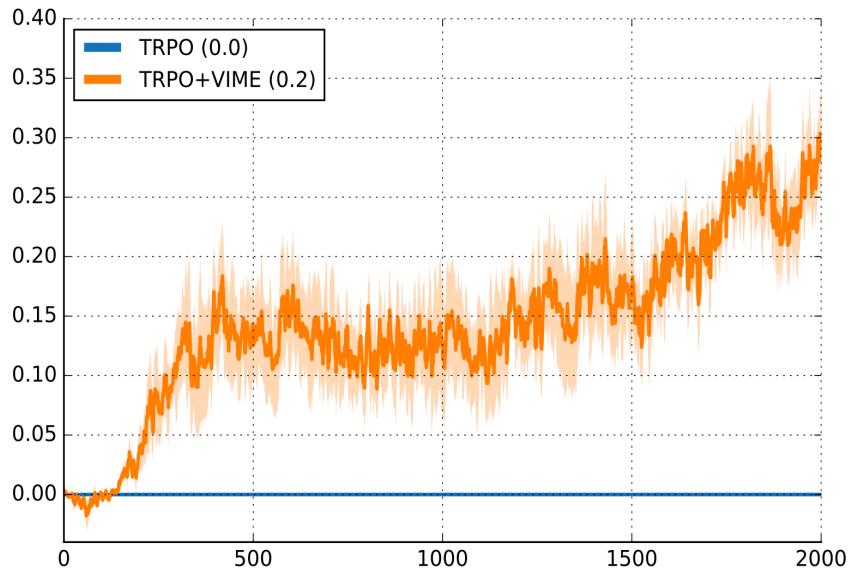
[Brockman et al. 2016]

# rllab

<https://github.com/rllab/rllab>



# VIME



[Houthooft et al. 2016]

# Thank you

Poster session: Monday 3pm-7pm #11