

Multi-Agent RL, PPO, MAPPO

Milan Straka

 January 02, 2023



EUROPEAN UNION
European Structural and Investment Fund
Operational Programme Research,
Development and Education

Charles University in Prague
Faculty of Mathematics and Physics
Institute of Formal and Applied Linguistics



unless otherwise stated

We use the thesis

Cooperative Multi-Agent Reinforcement Learning

<https://dspace.cuni.cz/handle/20.500.11956/127431>

as an introduction text.

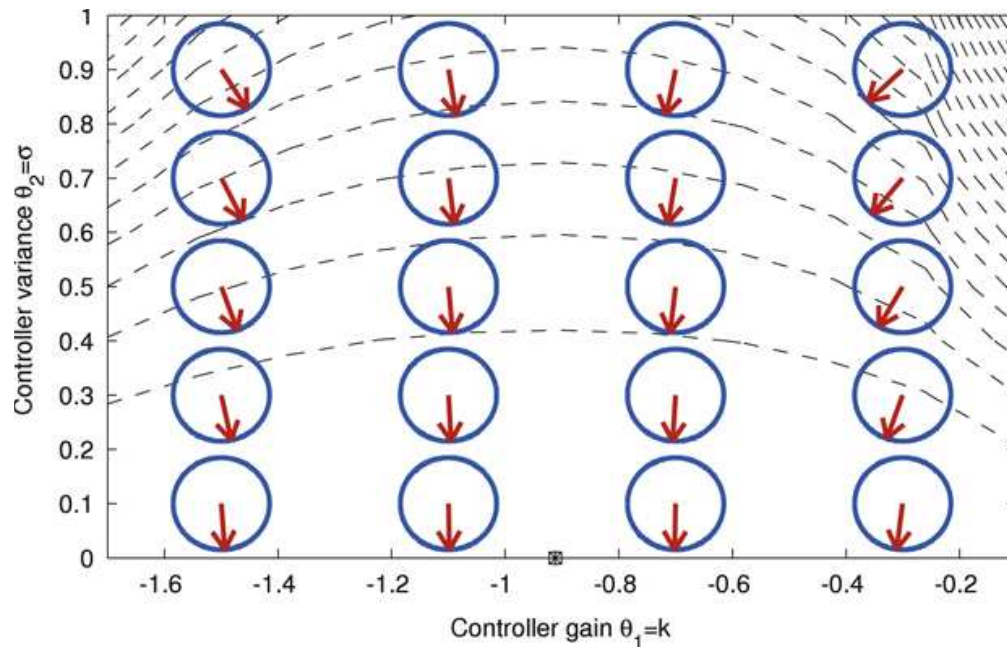
Natural Policy Gradient

The following approach has been introduced by Kakade (2002).

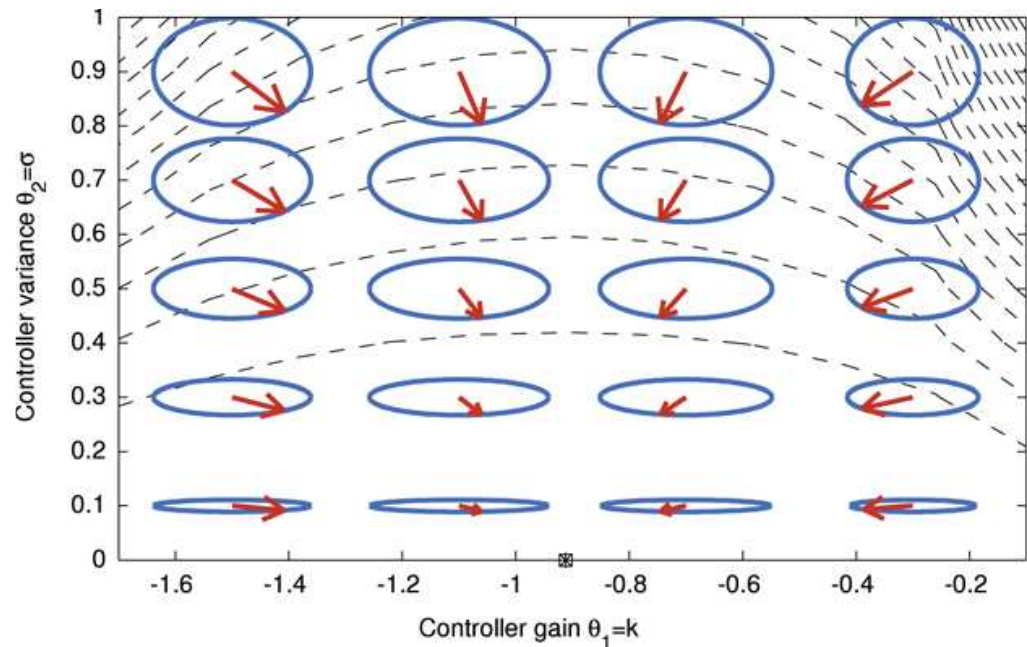
Using policy gradient theorem, we are able to compute ∇v_π . Normally, we update the parameters by using directly this gradient. This choice is justified by the fact that a vector \mathbf{d} which maximizes $v_\pi(\mathbf{s}; \boldsymbol{\theta} + \mathbf{d})$ under the constraint that $\|\mathbf{d}\|^2$ is bounded by a small constant is exactly the gradient ∇v_π .

Normally, the length $\|\mathbf{d}\|^2$ is computed using Euclidean metric. But in general, any metric could be used. Representing a metric using a positive-definite matrix \mathbf{G} (identity matrix for Euclidean metric), we can compute the distance as $\|\mathbf{d}\|_{\mathbf{G}}^2 = \sum_{ij} G_{ij} d_i d_j = \mathbf{d}^T \mathbf{G} \mathbf{d}$. The steepest ascent direction is then given by $\mathbf{G}^{-1} \nabla v_\pi$.

Note that when \mathbf{G} is the Hessian $\mathbf{H} v_\pi$, the above process is exactly Newton's method.



(a) Vanilla policy gradient.



(b) Natural policy gradient.

Figure 3 of "Reinforcement learning of motor skills with policy gradients" by Jan Peters et al.

Natural Policy Gradient

A suitable choice for the metric is *Fisher information matrix*, which is defined as a *covariance matrix* of the **score function** $\nabla_{\theta} \log \pi(a|s; \theta)$. The expectation of the score function is 0:

$$\mathbb{E}_{\pi(a|s; \theta)} \nabla_{\theta} \log \pi(a|s; \theta) = \int \pi(a|s; \theta) \frac{\nabla_{\theta} \pi(a|s; \theta)}{\pi(a|s; \theta)} da = \nabla_{\theta} \int \pi(a|s; \theta) da = \nabla_{\theta} 1 = 0.$$

The Fisher information matrix is therefore

$$F_s(\theta) \stackrel{\text{def}}{=} \mathbb{E}_{\pi(a|s; \theta)} \left[(\nabla_{\theta} \log \pi(a|s; \theta)) (\nabla_{\theta} \log \pi(a|s; \theta))^T \right].$$

It can be shown that the Fisher information metric is the only Riemannian metric (up to rescaling) invariant to change of parameters under sufficient statistic.

The Fisher information matrix is also a Hessian of the $D_{\text{KL}}(\pi(a|s; \theta) \parallel \pi(a|s; \theta'))$:

$$F_s(\theta) = \frac{\partial^2}{\partial \theta'_i \partial \theta'_j} D_{\text{KL}}(\pi(a|s; \theta) \parallel \pi(a|s; \theta')) \Big|_{\theta' = \theta}.$$

Using the metric

$$F(\boldsymbol{\theta}) = \mathbb{E}_{s \sim \mu_{\boldsymbol{\theta}}} F_s(\boldsymbol{\theta})$$

we want to update the parameters using $\mathbf{d}_F \stackrel{\text{def}}{=} F(\boldsymbol{\theta})^{-1} \nabla v_{\pi}$.

An interesting property of using the \mathbf{d}_F to update the parameters is that

- updating $\boldsymbol{\theta}$ using ∇v_{π} will choose an arbitrary *better* action in state s ;
- updating $\boldsymbol{\theta}$ using $F(\boldsymbol{\theta})^{-1} \nabla v_{\pi}$ chooses the *best* action (maximizing expected return), similarly to tabular greedy policy improvement.

However, computing \mathbf{d}_F in a straightforward way is too costly.

Duan et al. (2016) in paper *Benchmarking Deep Reinforcement Learning for Continuous Control* propose a modification to the NPG to efficiently compute \mathbf{d}_F .

Following Schulman et al. (2015), they suggest to use *conjugate gradient algorithm*, which can solve a system of linear equations $\mathbf{A}\mathbf{x} = \mathbf{b}$ in an iterative manner, by using \mathbf{A} only to compute products $\mathbf{A}\mathbf{v}$ for a suitable \mathbf{v} .

Therefore, \mathbf{d}_F is found as a solution of

$$F(\boldsymbol{\theta})\mathbf{d}_F = \nabla v_\pi$$

and using only 10 iterations of the algorithm seem to suffice according to the experiments.

Furthermore, Duan et al. suggest to use a specific learning rate suggested by Peters et al (2008) of

$$\frac{\alpha}{\sqrt{(\nabla v_\pi)^T F(\boldsymbol{\theta})^{-1} \nabla v_\pi}}.$$

Schulman et al. in 2015 wrote an influential paper introducing TRPO as an improved variant of NPG.

Considering two policies $\pi, \tilde{\pi}$, we can write

$$v_{\tilde{\pi}} = v_{\pi} + \mathbb{E}_{s \sim \mu(\tilde{\pi})} \mathbb{E}_{a \sim \tilde{\pi}(a|s)} a_{\pi}(a|s),$$

where $a_{\pi}(a|s)$ is the advantage function $q_{\pi}(a|s) - v_{\pi}(s)$ and $\mu(\tilde{\pi})$ is the on-policy distribution of the policy $\tilde{\pi}$.

Analogously to policy improvement, we see that if $a_{\pi}(a|s) \geq 0$, policy $\tilde{\pi}$ performance increases (or stays the same if the advantages are zero everywhere).

However, sampling states $s \sim \mu(\tilde{\pi})$ is costly. Therefore, we instead consider

$$L_{\pi}(\tilde{\pi}) = v_{\pi} + \mathbb{E}_{s \sim \mu(\pi)} \mathbb{E}_{a \sim \tilde{\pi}(a|s)} a_{\pi}(a|s).$$

$$L_{\pi}(\tilde{\pi}) = v_{\pi} + \mathbb{E}_{s \sim \mu(\pi)} \mathbb{E}_{a \sim \tilde{\pi}(a|s)} a_{\pi}(a|s)$$

It can be shown that for parametrized $\pi(a|s; \theta)$ the $L_{\pi}(\tilde{\pi})$ matches $v_{\tilde{\pi}}$ to the first order.

Schulman et al. additionally proves that if we denote $\alpha = D_{\text{KL}}^{\max}(\pi_{\text{old}} \parallel \pi_{\text{new}}) = \max_s D_{\text{KL}}(\pi_{\text{old}}(\cdot|s) \parallel \pi_{\text{new}}(\cdot|s))$, then

$$v_{\pi_{\text{new}}} \geq L_{\pi_{\text{old}}}(\pi_{\text{new}}) - \frac{4\varepsilon\gamma}{(1-\gamma)^2} \alpha \quad \text{where} \quad \varepsilon = \max_{s,a} |a_{\pi}(s,a)|.$$

Therefore, TRPO maximizes $L_{\pi_{\theta_0}}(\pi_{\theta})$ subject to $D_{\text{KL}}^{\theta_0}(\pi_{\theta_0} \parallel \pi_{\theta}) < \delta$, where

- $D_{\text{KL}}^{\theta_0}(\pi_{\theta_0} \parallel \pi_{\theta}) = \mathbb{E}_{s \sim \mu(\pi_{\theta_0})} [D_{\text{KL}}(\pi_{\text{old}}(\cdot|s) \parallel \pi_{\text{new}}(\cdot|s))]$ is used instead of D_{KL}^{\max} for performance reasons;
- δ is a constant found empirically, as the one implied by the above equation is too small;
- importance sampling is used to account for sampling actions from π .

maximize $L_{\pi_{\theta_0}}(\pi_{\theta}) = \mathbb{E}_{s \sim \mu(\pi_{\theta_0}), a \sim \pi_{\theta_0}(a|s)} \left[\frac{\pi_{\theta}(a|s)}{\pi_{\theta_0}(a|s)} a_{\pi_{\theta_0}}(a|s) \right]$ subject to $D_{\text{KL}}^{\theta_0}(\pi_{\theta_0} \parallel \pi_{\theta}) < \delta$

The parameters are updated using $\mathbf{d}_F = F(\boldsymbol{\theta})^{-1} \nabla L_{\pi_{\theta_0}}(\pi_{\theta})$, utilizing the conjugate gradient algorithm as described earlier for TNPG (note that the algorithm was designed originally for TRPO and only later employed for TNPG).

To guarantee improvement and respect the D_{KL} constraint, a line search is in fact performed.

We start by the learning rate of $\sqrt{\delta / (\mathbf{d}_F^T F(\boldsymbol{\theta})^{-1} \mathbf{d}_F)}$ and shrink it exponentially until the constraint is satisfied and the objective improves.

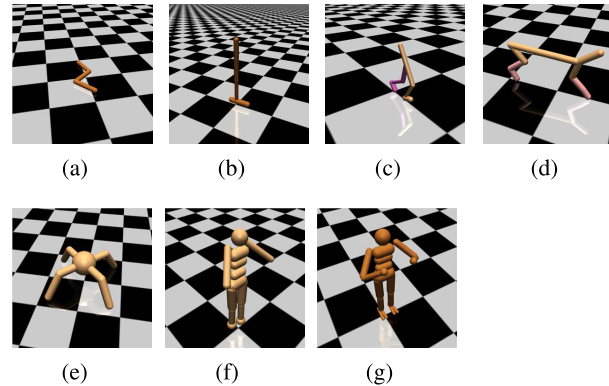


Figure 1. Illustration of locomotion tasks: (a) Swimmer; (b) Hopper; (c) Walker; (d) Half-Cheetah; (e) Ant; (f) Simple Humanoid; and (g) Full Humanoid.

Figure 1 of "Benchmarking Deep Reinforcement Learning for Continuous Control" by Duan et al.

Task	Random	REINFORCE	TNPG	RWR	REPS	TRPO	CEM	CMA-ES	DDPG
Cart-Pole Balancing	77.1 ± 0.0	4693.7 ± 14.0	3986.4 ± 748.9	4861.5 ± 12.3	565.6 ± 137.6	4869.8 ± 37.6	4815.4 ± 4.8	2440.4 ± 568.3	4634.4 ± 87.8
Inverted Pendulum*	-153.4 ± 0.2	13.4 ± 18.0	209.7 ± 55.5	84.7 ± 13.8	-113.3 ± 4.6	247.2 ± 76.1	38.2 ± 25.7	-40.1 ± 5.7	40.0 ± 244.6
Mountain Car	-415.4 ± 0.0	-67.1 ± 1.0	-66.5 ± 4.5	-79.4 ± 1.1	-275.6 ± 166.3	-61.7 ± 0.9	-66.0 ± 2.4	-85.0 ± 7.7	-288.4 ± 170.3
Acrobot	-1904.5 ± 1.0	-508.1 ± 91.0	-395.8 ± 121.2	-352.7 ± 35.9	-1001.5 ± 10.8	-326.0 ± 24.4	-436.8 ± 14.7	-785.6 ± 13.1	-223.6 ± 5.8
Double Inverted Pendulum*	149.7 ± 0.1	4116.5 ± 65.2	4455.4 ± 37.6	3614.8 ± 368.1	446.7 ± 114.8	4412.4 ± 50.4	2566.2 ± 178.9	1576.1 ± 51.3	2863.4 ± 154.0
Swimmer*	-1.7 ± 0.1	92.3 ± 0.1	96.0 ± 0.2	60.7 ± 5.5	3.8 ± 3.3	96.0 ± 0.2	68.8 ± 2.4	64.9 ± 1.4	85.8 ± 1.8
Hopper	8.4 ± 0.0	714.0 ± 29.3	1155.1 ± 57.9	553.2 ± 71.0	86.7 ± 17.6	1183.3 ± 150.0	63.1 ± 7.8	20.3 ± 14.3	267.1 ± 43.5
2D Walker	-1.7 ± 0.0	506.5 ± 78.8	1382.6 ± 108.2	136.0 ± 15.9	-37.0 ± 38.1	1353.8 ± 85.0	84.5 ± 19.2	77.1 ± 24.3	318.4 ± 181.6
Half-Cheetah	-90.8 ± 0.3	1183.1 ± 69.2	1729.5 ± 184.6	376.1 ± 28.2	34.5 ± 38.0	1914.0 ± 120.1	330.4 ± 274.8	441.3 ± 107.6	2148.6 ± 702.7
Ant*	13.4 ± 0.7	548.3 ± 55.5	706.0 ± 127.7	37.6 ± 3.1	39.0 ± 9.8	730.2 ± 61.3	49.2 ± 5.9	17.8 ± 15.5	326.2 ± 20.8
Simple Humanoid	41.5 ± 0.2	128.1 ± 34.0	255.0 ± 24.5	93.3 ± 17.4	28.3 ± 4.7	269.7 ± 40.3	60.6 ± 12.9	28.7 ± 3.9	99.4 ± 28.1
Full Humanoid	13.2 ± 0.1	262.2 ± 10.5	288.4 ± 25.2	46.7 ± 5.6	41.7 ± 6.1	287.0 ± 23.4	36.9 ± 2.9	N/A ± N/A	119.0 ± 31.2

Table 1 of "Benchmarking Deep Reinforcement Learning for Continuous Control" by Duan et al.

Proximal Policy Optimization

A simplification of TRPO which can be implemented using a few lines of code.

Let $r_t(\boldsymbol{\theta}) \stackrel{\text{def}}{=} \frac{\pi(A_t|S_t;\boldsymbol{\theta})}{\pi(A_t|S_t;\boldsymbol{\theta}_{\text{old}})}$. PPO maximizes the objective (i.e., you should minimize its negation)

$$L^{\text{CLIP}}(\boldsymbol{\theta}) \stackrel{\text{def}}{=} \mathbb{E}_t \left[\min \left(r_t(\boldsymbol{\theta}) \hat{A}_t, \text{clip}(r_t(\boldsymbol{\theta}), 1 - \epsilon, 1 + \epsilon) \hat{A}_t \right) \right].$$

Such a $L^{\text{CLIP}}(\boldsymbol{\theta})$ is a lower (pessimistic) bound.

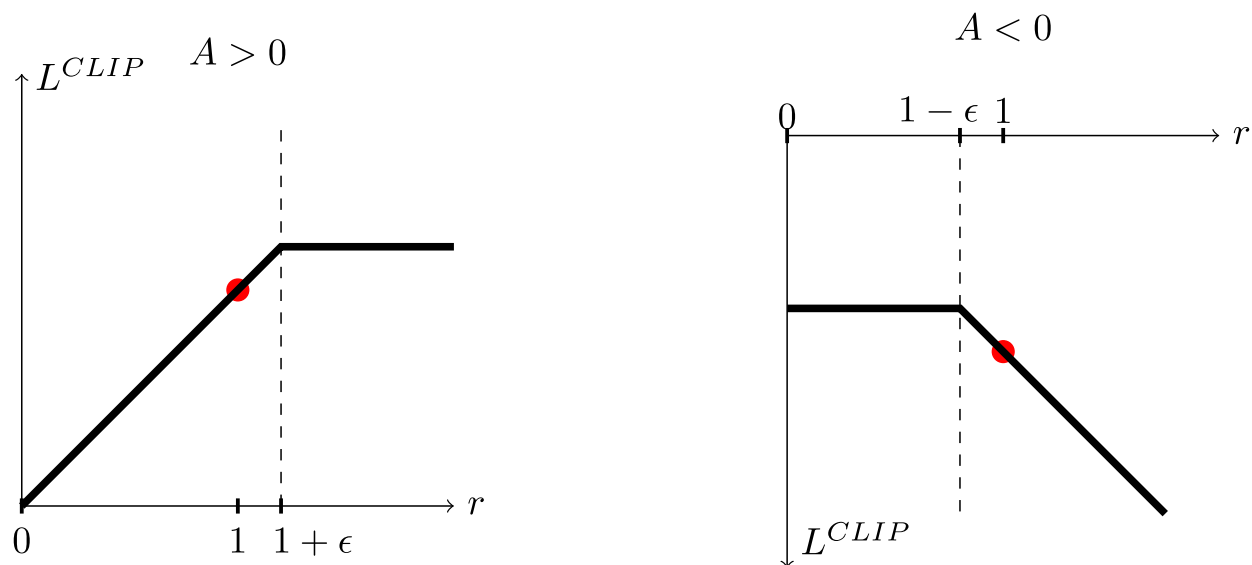


Figure 1 of "Proximal Policy Optimization Algorithms" by Schulman et al.

Proximal Policy Optimization

The advantages \hat{A}_t are additionally estimated using the so-called *generalized advantage estimation*, which is just an analogue of the truncated n-step lambda-return:

$$\hat{A}_t = \sum_{i=0}^{n-1} \gamma^i \lambda^i (R_{t+1+i} + \gamma V(S_{t+i+1}) - V(S_{t+i})).$$

Algorithm 1 PPO, Actor-Critic Style

```

for iteration=1, 2, ... do
  for actor=1, 2, ..., N do
    Run policy  $\pi_{\theta_{\text{old}}}$  in environment for  $T$  timesteps
    Compute advantage estimates  $\hat{A}_1, \dots, \hat{A}_T$ 
  end for
  Optimize surrogate  $L$  wrt  $\theta$ , with  $K$  epochs and minibatch size  $M \leq NT$ 
   $\theta_{\text{old}} \leftarrow \theta$ 
end for

```

Algorithm 1 of "Proximal Policy Optimization Algorithms" by Schulman et al.

Proximal Policy Optimization

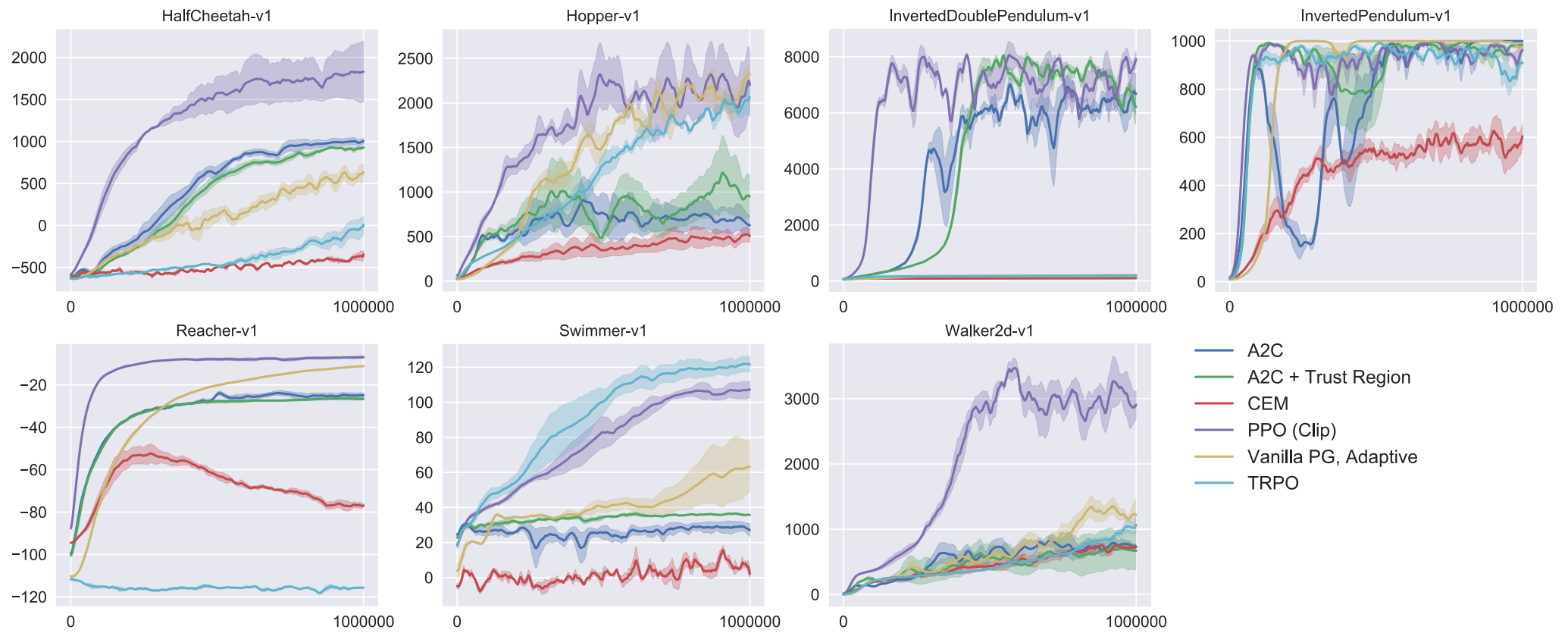


Figure 3: Comparison of several algorithms on several MuJoCo environments, training for one million timesteps.

Figure 3 of "Proximal Policy Optimization Algorithms" by Schulman et al.

As another example, consider <https://openai.com/blog/emergent-tool-use/>.