

# Lecture 13: Further Contemporary RL Algorithms

Oliver Wallscheid



- 1 Trust region policy optimization (TRPO)
- 2 Proximal policy optimization (PPO)

# Reinterpreting the stochastic policy gradient (1)

- ▶ In the following we will **focus on stochastic policies**  $\pi(\mathbf{u}|\mathbf{x})$  .
- ▶ First, we rewrite the performance metric (11.7) to obtain

$$J_\pi = \mathbb{E}_\pi \left[ \sum_{k=0}^{\infty} \gamma^k R_k \right]. \quad (13.1)$$

- ▶ Using the advantage  $a_\pi(\mathbf{x}, \mathbf{u}) = q_\pi(\mathbf{x}, \mathbf{u}) - v_\pi(\mathbf{x})$  we can calculate the performance of an updated policy  $\pi \rightarrow \tilde{\pi}$ <sup>1</sup>:

$$J_{\tilde{\pi}} = J_\pi + \int_{\mathcal{X}} p^{\tilde{\pi}}(\mathbf{x}) \int_{\mathcal{U}} \tilde{\pi}(\mathbf{u}|\mathbf{x}) a_\pi(\mathbf{x}, \mathbf{u}). \quad (13.2)$$

- ▶ While for finite MDPs, the policy improvement theorem guaranteed  $J_{\tilde{\pi}} \geq J_\pi$  for each policy update, there might be some states where  $\int_{\mathcal{U}} \tilde{\pi}(\mathbf{u}|\mathbf{x}) a_\pi < 0$  for continuous MDPs using function approximation.

---

<sup>1</sup>proof from: S. Kakade and J. Langford, *Approximately optimal approximate reinforcement learning*, ICML, vol. 2, pp 267-274, 2002

## Reinterpreting the stochastic policy gradient (2)

- ▶ For easier calculation, we introduce a local approximation to (13.2)

$$\mathcal{L}_\pi(\tilde{\pi}) = J_\pi + \int_{\mathcal{X}} p^\pi(\mathbf{x}) \int_{\mathcal{U}} \tilde{\pi}(\mathbf{u}|\mathbf{x}) a_\pi(\mathbf{x}, \mathbf{u}) \quad (13.3)$$

where  $p^\pi(\mathbf{x})$  is used instead of  $p^{\tilde{\pi}}(\mathbf{x})$ , i.e., neglecting the state distribution change due to a policy update.

- ▶ For any parametrized and differentiable policy  $\pi_\theta(\mathbf{u}|\mathbf{x})$ , it can be shown that

$$\begin{aligned} \mathcal{L}(\pi_{\theta_0}) &= J(\pi_{\theta_0}), \\ \nabla_{\theta} \mathcal{L}_{\pi_{\theta_0}}(\pi_{\theta})|_{\theta=\theta_0} &= \nabla_{\theta} J(\pi_{\theta})|_{\theta=\theta_0} \end{aligned} \quad (13.4)$$

for any initial parameter set  $\theta_0$ .

- ▶ For a sufficiently small step size, improving  $\mathcal{L}_{\pi_{\theta_0}}$  will also improve  $J$ .

However, we do not know how much the actual stochastic policy will change while moving through the parameter space. Hence, we do not have a good decision basis to choose the policy gradient step size.

## Adding a trust region constraint (1)

- ▶ From the previous discussion it can be concluded that we want a **metric describing how much a policy is changed in the action space when updating the policy in the parameter space**.
- ▶ Against this background, we make use of the **Kullback-Leibler divergence** (also called relative entropy)

$$D_{\text{KL}}(P \parallel Q) = \int_{-\infty}^{\infty} p(x) \log \left( \frac{p(x)}{q(x)} \right) dx \quad (13.5)$$

defined for continuous distributions  $P$  and  $Q$  with their probability densities  $p$  and  $q$ .

- ▶ Example: for two multivariate Gaussian distributions of equal dimensions  $d$ , with means  $\boldsymbol{\mu}_0, \boldsymbol{\mu}_1$  and with (non-singular) covariance matrix  $\boldsymbol{\Sigma}_0, \boldsymbol{\Sigma}_1$  we receive

$$D_{\text{KL}}(\mathcal{N}_0 \parallel \mathcal{N}_1) = \frac{1}{2} \left( \text{tr}(\boldsymbol{\Sigma}_1^{-1} \boldsymbol{\Sigma}_0) + (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0)^{\top} \boldsymbol{\Sigma}_1^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0) - d + \ln \left( \frac{\det \boldsymbol{\Sigma}_1}{\det \boldsymbol{\Sigma}_0} \right) \right).$$

## Adding a trust region constraint (2)

- ▶ The **trust region policy optimization (TRPO)** updates the policy parameters while constraining the KL divergence between the new and the old policy distribution:

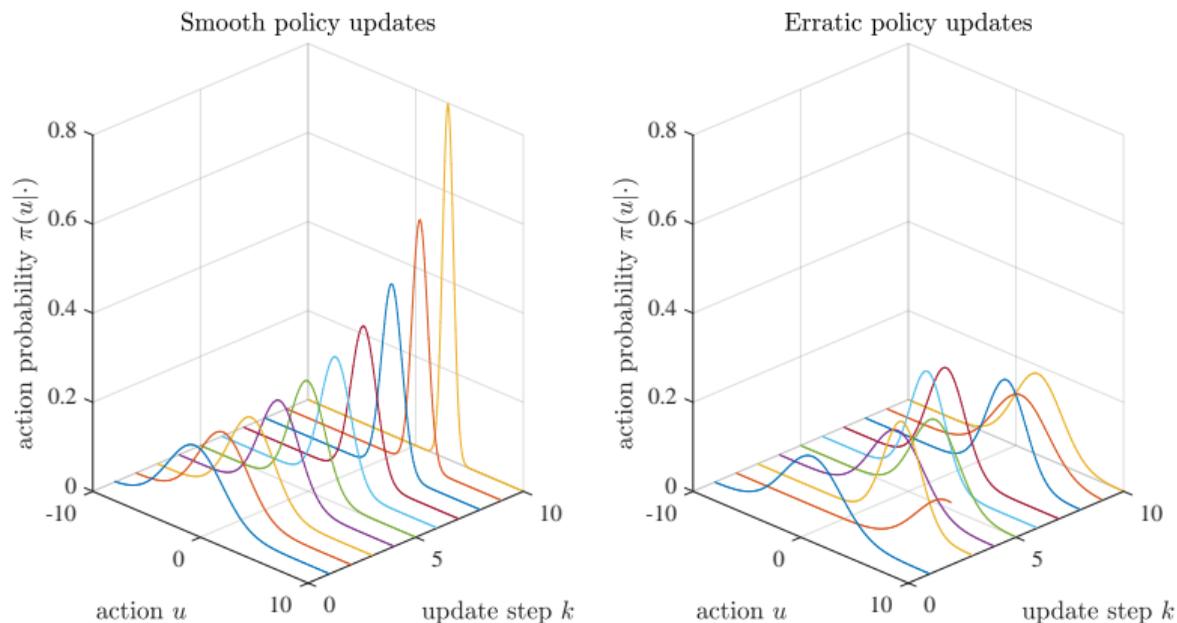
$$\begin{aligned} & \max_{\boldsymbol{\theta}} \mathcal{L}_{\boldsymbol{\theta}_k}(\boldsymbol{\theta}), \\ \text{s.t.} \quad & \bar{D}_{\text{KL}}(\boldsymbol{\theta}_k, \boldsymbol{\theta}) \leq \kappa \end{aligned} \tag{13.6}$$

with

$$\bar{D}_{\text{KL}}(\boldsymbol{\theta}_k, \boldsymbol{\theta}) = \bar{D}_{\text{KL}}(\pi_{\boldsymbol{\theta}_k}, \pi_{\boldsymbol{\theta}}) = \mathbb{E}_{\pi_{\boldsymbol{\theta}_k}} [D_{\text{KL}}(\pi_{\boldsymbol{\theta}_k}(\cdot|\mathbf{X}) \parallel \pi_{\boldsymbol{\theta}}(\cdot|\mathbf{X}))].$$

- ▶ Hence, we want to **limit the average KL divergence w.r.t. the states visited by the old policy**.
- ▶ The constraint  $\kappa$  is a TRPO hyperparameter (typically  $\kappa \ll 1$ ).
- ▶ Although (13.6) does not provide any formal convergence guarantee, we at least have a link between changes in the parameter and policy distribution space. Therefore, **we can use this tool to prevent erratic policy changes**.

# Smooth policy updates via TRPO



**Fig. 13.1:** Simplified representation of the policy evolution for a scalar action given some fixed state. Left: TRPO-style updates finding the optimal action with increasing probability. Right: Unmonitored policy distributions not converging towards an optimal policy ('policy chattering').

# Sample-based objective and constraint estimation (1)

- ▶ To actually solve (13.6) we will make use of samplings from **Monte Carlo rollouts**.
- ▶ Expanding the objective yields

$$\max_{\theta} \mathcal{L}_{\theta_k}(\theta) = \max_{\theta} J_{\pi_k} + \int_{\mathcal{X}} p^{\pi_k}(\mathbf{x}) \int_{\mathcal{U}} \pi_{\theta}(\mathbf{u}|\mathbf{x}) a_{\pi_k}(\mathbf{x}, \mathbf{u}). \quad (13.7)$$

- ▶ The first term  $J_{\pi_k}$  can be dropped, since it is irrelevant for the optimization result (constant).
- ▶ Using samples we can approximate  $\int_{\mathcal{X}} p^{\pi_k}(\mathbf{x}) \approx \frac{1}{1-\gamma} \mathbb{E}_{\pi_{\theta_k}}[\mathbf{X}]$ .
- ▶ Moreover,  $\int_{\mathcal{U}} \pi_{\theta}(\mathbf{u}|\mathbf{x}) a_{\pi_k}(\mathbf{x}, \mathbf{u}) \approx \mathbb{E}_{\pi_{\theta_k}} \left[ \frac{\pi_{\theta}(\mathbf{U}|\mathbf{X})}{\pi_{\theta_k}(\mathbf{U}|\mathbf{X})} a_{\pi_k}(\mathbf{X}, \mathbf{U}) \right]$  is also approximated applying importance sampling based on data from the old policy.
- ▶ Hence, the sampled objective is

$$\max_{\theta} \mathbb{E}_{\pi_{\theta_k}} \left[ \frac{\pi_{\theta}(\mathbf{U}|\mathbf{X})}{\pi_{\theta_k}(\mathbf{U}|\mathbf{X})} a_{\pi_k}(\mathbf{X}, \mathbf{U}) \right]. \quad (13.8)$$

## Sample-based objective and constraint estimation (2)

- ▶ Applying the previous sample-based estimation we obtain

$$\begin{aligned} \boldsymbol{\theta}_{k+1} &= \arg \max_{\boldsymbol{\theta}} \mathbb{E}_{\pi_{\boldsymbol{\theta}_k}} \left[ \frac{\pi_{\boldsymbol{\theta}}(\mathbf{U}|\mathbf{X})}{\pi_{\boldsymbol{\theta}_k}(\mathbf{U}|\mathbf{X})} a_{\pi_k}(\mathbf{X}, \mathbf{U}) \right], \\ \text{s.t.} \quad &\mathbb{E}_{\pi_{\boldsymbol{\theta}_k}} [D_{\text{KL}}(\pi_{\boldsymbol{\theta}_k}(\cdot|\mathbf{X}) \parallel \pi_{\boldsymbol{\theta}}(\cdot|\mathbf{X}))] \leq \kappa. \end{aligned} \tag{13.9}$$

- ▶ Hence, we have a **three-step procedure** for each TRPO update:

- 1 Use Monte Carlo simulations based on the old policy to obtain data.
- 2 Use the data to construct (13.9).
- 3 Solve the constrained optimization problem to update the policy parameter vector.

Solving (13.9) is generally a nonlinear optimization problem. The original TRPO implementation uses a local objective and constraint approximation together with conjugate gradient and line search algorithms. However, many other constrained-nonlinear solvers are also applicable.

# Generalized advantage estimation

- ▶ Having data  $\langle \mathbf{x}, \mathbf{u}, r, \mathbf{x}' \rangle$  in  $\mathcal{D}$  from a Monte Carlo rollout available, an important problem is to estimate  $a_{\pi_k}(\mathbf{x}, \mathbf{u})$  in (13.9).
- ▶ A particular suggestion in the TRPO context is to use a **generalized advantage estimator (GAE)**<sup>1</sup> defined as

$$\hat{a}_k^{(\gamma, \lambda)} = \sum_{i=0}^{\infty} (\gamma \lambda)^i \delta_{k+i}. \quad (13.10)$$

- ▶ Here,  $\delta_k = r_k + \gamma v(\mathbf{x}_{k+1}) - v(\mathbf{x}_k)$  is a single advantage sample.
- ▶ Hence, the GAE is the exponentially-weighted average of the discounted advantage samples with an additional weighting  $\lambda$ .
- ▶ Similar formulation compared to TD( $\lambda$ ) but the estimator's target is the advantage.
- ▶ The choice of  $(\gamma \lambda)$  trade-offs the bias and variance of the estimator.

---

<sup>1</sup>cf. J. Schulman et al., *High Dimensional Continuous Control Using Generalized Advantage Estimation*, <https://arxiv.org/abs/1506.02438>, 2015

# TRPO summary

The TRPO's key facts are:

- ▶ The TRPO constrains policy distribution changes when updating the policy parameters (for stochastic policies and on-policy learning).
- ▶ The objective is to enable a monotonically improving learning process.
- ▶ Using trust regions, erratic policy updates should be prevented.

The TRPO's main hurdles are:

- ▶ Constructing the objective function and constraint requires Monte Carlo rollouts (time consuming, data inefficient).
- ▶ When the sampled optimization problem is set up, a nonlinear and constrained optimization step is required (no simple policy gradient, computational costly).

We will not provide any specific TRPO implementation suggestion at this point, since this is rather cumbersome. Instead we will move forward to a similar algorithm which is pursuing the same goal (prevent erratic policy changes) with a much simpler implementation.

# Table of contents

- 1 Trust region policy optimization (TRPO)
- 2 Proximal policy optimization (PPO)

# Background and motivation

- ▶ The upcoming **proximal policy optimization (PPO)** algorithm tries to mimic the constrained TRPO problem

$$\begin{aligned} \boldsymbol{\theta}_{k+1} &= \arg \max_{\boldsymbol{\theta}} \mathbb{E}_{\pi_{\boldsymbol{\theta}_k}} \left[ \frac{\pi_{\boldsymbol{\theta}}(\mathbf{U}|\mathbf{X})}{\pi_{\boldsymbol{\theta}_k}(\mathbf{U}|\mathbf{X})} a_{\pi_k}(\mathbf{X}, \mathbf{U}) \right], \\ \text{s.t.} \quad &\mathbb{E}_{\pi_{\boldsymbol{\theta}_k}} [D_{\text{KL}}(\pi_{\boldsymbol{\theta}_k}(\cdot|\mathbf{X}) \parallel \pi_{\boldsymbol{\theta}}(\cdot|\mathbf{X}))] \leq \kappa. \end{aligned}$$

based on related unconstrained problems.

- ▶ Hence, the objective will be reformulated to incorporate mechanisms preventing excessively large variations of the policy distribution during a parameter update (leading to an updated policy with sufficient proximity to the old one).
- ▶ Moreover, **PPO incorporates two variants** which we will discuss:
  - 1 Clipping the surrogate objective,
  - 2 Adaptive tuning of a KL-associated penalty coefficient.

# Clipped surrogate objective

- ▶ The first approach is based on the following **clipped objective**:

$$\mathbb{E}_{\pi_{\theta_k}} \left[ \min \left\{ \frac{\pi_{\theta}(\mathbf{U}|\mathbf{X})}{\pi_{\theta_k}(\mathbf{U}|\mathbf{X})} a_{\pi_k}(\mathbf{X}, \mathbf{U}), \text{clip} \left( \frac{\pi_{\theta}(\mathbf{U}|\mathbf{X})}{\pi_{\theta_k}(\mathbf{U}|\mathbf{X})}, 1 - \epsilon, 1 + \epsilon \right) a_{\pi_k}(\mathbf{X}, \mathbf{U}) \right\} \right]. \quad (13.11)$$

- ▶ Above,  $\epsilon < 1$  is a PPO hyperparameter serving as a regularizer.
- ▶ The first element of  $\min\{\cdot\}$  is the previous TRPO objective.
- ▶ The second element of  $\min\{\cdot\}$  modifies the surrogate objective by clipping the importance sampling ratio  $\pi_{\theta}/\pi_{\theta_k}$ .
- ▶ The latter should remove the incentive for moving the importance sampling ratio outside of the interval  $[1 - \epsilon, 1 + \epsilon]$ .
- ▶ The modified objective is therefore a lower bound of the unclipped TRPO objective.

## Clipped surrogate objective: positive advantage

- ▶ Consider a single sample  $(\mathbf{x}, \mathbf{u})$  with a **positive advantage**  $a_{\pi_k}(\mathbf{x}, \mathbf{u})$ :

$$\max_{\theta} \min \left\{ \frac{\pi_{\theta}(\mathbf{u}|\mathbf{x})}{\pi_{\theta_k}(\mathbf{u}|\mathbf{x})} a_{\pi_k}(\mathbf{x}, \mathbf{u}), \text{clip} \left( \frac{\pi_{\theta}(\mathbf{u}|\mathbf{x})}{\pi_{\theta_k}(\mathbf{u}|\mathbf{x})}, 1 - \epsilon, 1 + \epsilon \right) a_{\pi_k}(\mathbf{x}, \mathbf{u}) \right\}.$$

- ▶ Because the advantage is positive, the objective will increase if the action becomes more likely, i.e., if  $\pi_{\theta}(\mathbf{u}|\mathbf{x})$  increases.
- ▶ If  $\pi_{\theta}(\mathbf{u}|\mathbf{x}) > (1 + \epsilon)\pi_{\theta_k}(\mathbf{u}|\mathbf{x})$  the clipping becomes active.
- ▶ Hence, the objective reduces to

$$\max_{\theta} \min \left\{ \frac{\pi_{\theta}(\mathbf{u}|\mathbf{x})}{\pi_{\theta_k}(\mathbf{u}|\mathbf{x})}, 1 + \epsilon \right\} a_{\pi_k}(\mathbf{x}, \mathbf{u}).$$

- ▶ Due to the  $\min\{\cdot\}$  operator, the entire objective is therefore limited to  $(1 + \epsilon)a_{\pi_k}(\mathbf{x}, \mathbf{u})$ .
- ▶ Interpretation: the new policy does not benefit from going very away from the old policy distribution.

## Clipped surrogate objective: negative advantage

- ▶ Consider a single sample  $(\mathbf{x}, \mathbf{u})$  with a **negative advantage**  $a_{\pi_k}(\mathbf{x}, \mathbf{u})$ :

$$\max_{\boldsymbol{\theta}} \min \left\{ \frac{\pi_{\boldsymbol{\theta}}(\mathbf{u}|\mathbf{x})}{\pi_{\boldsymbol{\theta}_k}(\mathbf{u}|\mathbf{x})} a_{\pi_k}(\mathbf{x}, \mathbf{u}), \text{clip} \left( \frac{\pi_{\boldsymbol{\theta}}(\mathbf{u}|\mathbf{x})}{\pi_{\boldsymbol{\theta}_k}(\mathbf{u}|\mathbf{x})}, 1 - \epsilon, 1 + \epsilon \right) a_{\pi_k}(\mathbf{x}, \mathbf{u}) \right\}.$$

- ▶ Because the advantage is negative, the objective will increase if the action becomes less likely, i.e., if  $\pi_{\boldsymbol{\theta}}(\mathbf{u}|\mathbf{x})$  decreases.
- ▶ If  $\pi_{\boldsymbol{\theta}}(\mathbf{u}|\mathbf{x}) < (1 - \epsilon)\pi_{\boldsymbol{\theta}_k}(\mathbf{u}|\mathbf{x})$  the clipping becomes active.
- ▶ Hence, the objective reduces to

$$\max_{\boldsymbol{\theta}} \max \left\{ \frac{\pi_{\boldsymbol{\theta}}(\mathbf{u}|\mathbf{x})}{\pi_{\boldsymbol{\theta}_k}(\mathbf{u}|\mathbf{x})}, 1 - \epsilon \right\} a_{\pi_k}(\mathbf{x}, \mathbf{u}).$$

- ▶ Due to the  $\max\{\cdot\}$  operator, the entire objective is limited to  $(1 - \epsilon)a_{\pi_k}(\mathbf{x}, \mathbf{u})$ .

- ▶ The second PPO variant makes use of the following **KL-penalized objective**

$$\mathbb{E}_{\pi_{\theta_k}} \left[ \frac{\pi_{\theta}(\mathbf{U}|\mathbf{X})}{\pi_{\theta_k}(\mathbf{U}|\mathbf{X})} a_{\pi_k}(\mathbf{X}, \mathbf{U}) - \beta D_{\text{KL}}(\pi_{\theta_k}(\cdot|\mathbf{X}) \parallel \pi_{\theta}(\cdot|\mathbf{X})) \right]. \quad (13.12)$$

- ▶ Transfers the KL-based constraint into a penalty for large policy distribution changes.
- ▶ The parameter  $\beta$  weights the penalty against the policy improvement.
- ▶ The original PPO implementation suggests an adaptive tuning of  $\beta$  w.r.t. the sampled average KL divergence  $\bar{D}_{\text{KL}}(\theta_k, \theta)$  estimated from previous experience

$$\begin{aligned} \bar{D}_{\text{KL}}(\theta_k, \theta) < \bar{D}_{\text{KL}}^* &: \beta \leftarrow \beta/2, \\ \bar{D}_{\text{KL}}(\theta_k, \theta) > \bar{D}_{\text{KL}}^* &: \beta \leftarrow \beta \cdot 2. \end{aligned} \quad (13.13)$$

with some target value of the KL divergence  $\bar{D}_{\text{KL}}^*$  (additional hyperparameter).

# Algo. implementation: PPO

**input:** diff. stochastic policy fct.  $\pi(\mathbf{u}|\mathbf{x}, \boldsymbol{\theta})$  and value fct.  $\hat{v}(\mathbf{x}, \mathbf{w})$   
**parameter:** step sizes  $\{\alpha_w, \alpha_\theta\} \in \{\mathbb{R} | 0 < \alpha\}$   
**init:** weights  $\mathbf{w} \in \mathbb{R}^\zeta$  and  $\boldsymbol{\theta} \in \mathbb{R}^d$  arbitrarily, memory  $\mathcal{D}$   
**for**  $j = 1, 2, \dots$ , (*sub-*)*episodes* **do**  
    initialize  $\mathbf{x}_0$  (if new episode);  
    collect a set of tuples  $\langle \mathbf{x}_k, \mathbf{u}_k, r_{k+1}, \mathbf{x}_{k+1} \rangle$  by a rollout using  $\pi(\mathbf{u}|\mathbf{x}, \boldsymbol{\theta}_j)$ ;  
    store them in  $\mathcal{D}$ ;  
    estimate the advantage  $\hat{a}_{\pi_j}(\mathbf{x}, \mathbf{u})$  based on  $\hat{v}(\mathbf{x}, \mathbf{w}_j)$  and  $\mathcal{D}$  (e.g., GAE);  
     $\boldsymbol{\theta}_{j+1} \leftarrow$  policy gradient update based on the PPO variant (13.11) or (13.12);  
     $\mathbf{w}_{j+1} \leftarrow$  minimizing the mean-squared TD errors using  $\mathcal{D}$  (critic);  
    delete entries in  $\mathcal{D}$  (due to on-policy learning);

**Algo. 13.1:** Proximal policy optimization (output: parameter vectors  $\boldsymbol{\theta}^*$  for  $\pi^*(\mathbf{u}|\mathbf{x}, \boldsymbol{\theta}^*)$  and  $\mathbf{w}^*$  for  $\hat{v}^*(\mathbf{x}, \mathbf{w}^*)$ )

## Some PPO remarks

- ▶ Clipping the surrogate objective (13.11) was reported to achieve higher performances than the KL penalty (13.12).<sup>1</sup>
- ▶ Like TRPO, PPO is an on-policy algorithm. Hence, the memory  $\mathcal{D}$  is not a rolling replay buffer (cf. off-policy algorithms like DQN, DDPG or TD3) but a **rollout buffer** using one fixed policy.
- ▶ These rollouts are likely to result in an increased sample demand either using a simulator or a real experiment.

Although PPO is derived from a TRPO background pursuing monotonically increasing policy performance, its realization is based on multiple heuristics and approximations. Hence, there is no guarantee on achieving this goal and the specific performance of the PPO algorithm must be evaluated empirically given a certain application.

---

<sup>1</sup>cf. original PPO paper results by J. Schulman et al., *Proximal Policy Optimization Algorithms*, <https://arxiv.org/abs/1707.06347>, 2017

# Exemplary performance comparison

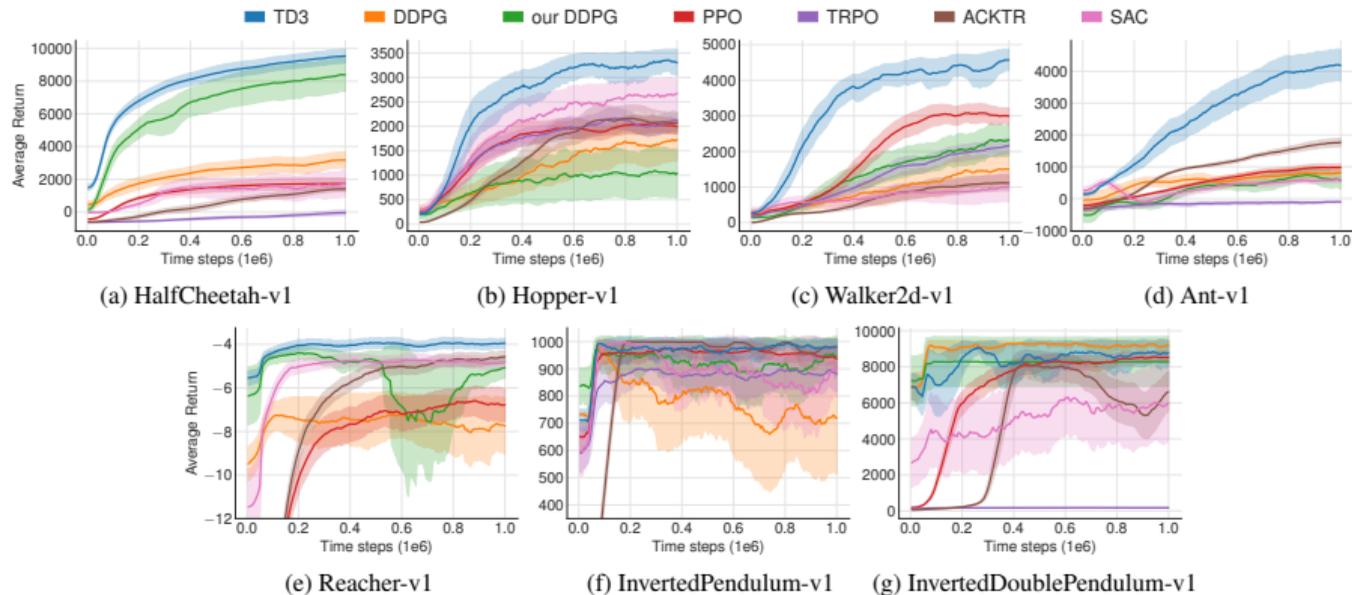


Fig. 13.2: Learning curves for OpenAI Gym continuous control tasks. The shaded region represents half a standard deviation of the average evaluation over ten trials (source: S. Fujimoto et al., *Addressing Function Approximation Error in Actor-Critic Methods*, 2018).

# Algorithmic outlook: other contemporary model-free algorithms (1)

The selection of algorithms appears endless:

- ▶ DQN variants such as
  - ▶ (Prioritized) dueling DQN
  - ▶ Noisy DQN
  - ▶ Distributional DQN
- ▶ Rainbow (combining multiple DQN extensions)
- ▶ Soft actor-critic (SAC)
- ▶ Actor critic using Kronecker-factored trust region (ACKTR)
- ▶ Asynchronous advantage actor-critic (A3C)
- ▶ ....

Remarks:

- ▶ You have already learned the basic building blocks in order to make yourself familiar with any value-/policy-based model-free RL approach.
- ▶ Use this knowledge!
- ▶ Focus on primary scientific literature for self-studying and not on unreliable sources!

## Algorithmic outlook: other contemporary model-free algorithms (2)

Algorithm collections with tutorial-style documentation:

- ▶ Intel Reinforcement Learning Coach
- ▶ OpenAI Spinning Up

Algorithm collections with decent application-oriented documentation:

- ▶ RLlib (Ray)
- ▶ Stable Baselines3
- ▶ Acme
- ▶ Garage
- ▶ Google Dopamine
- ▶ Tensorforce
- ▶ TF-Agents
- ▶ ...

## Summary: what you've learned today

- ▶ Trust region policy optimization (TRPO) pursues monotonically increasing policy performance by limiting policy distribution changes.
- ▶ This results in a nonlinear constrained optimization problem adding computational complexity (no simple policy gradients).
- ▶ Proximal policy optimization (PPO) converts the TRPO idea into an unconstrained optimization problem by a modified objective. Likewise, the PPO's objective is to prevent erratic policy distribution changes.