

Deep Learning (1470)

Randall Balestriero

Class 21: TRPO and PPO

Recap!

On-Policy and Off-Policy Learning

RL algorithms collect experiences and learn from these experiences

On-Policy Algorithms have to collect experiences with the policy they are learning

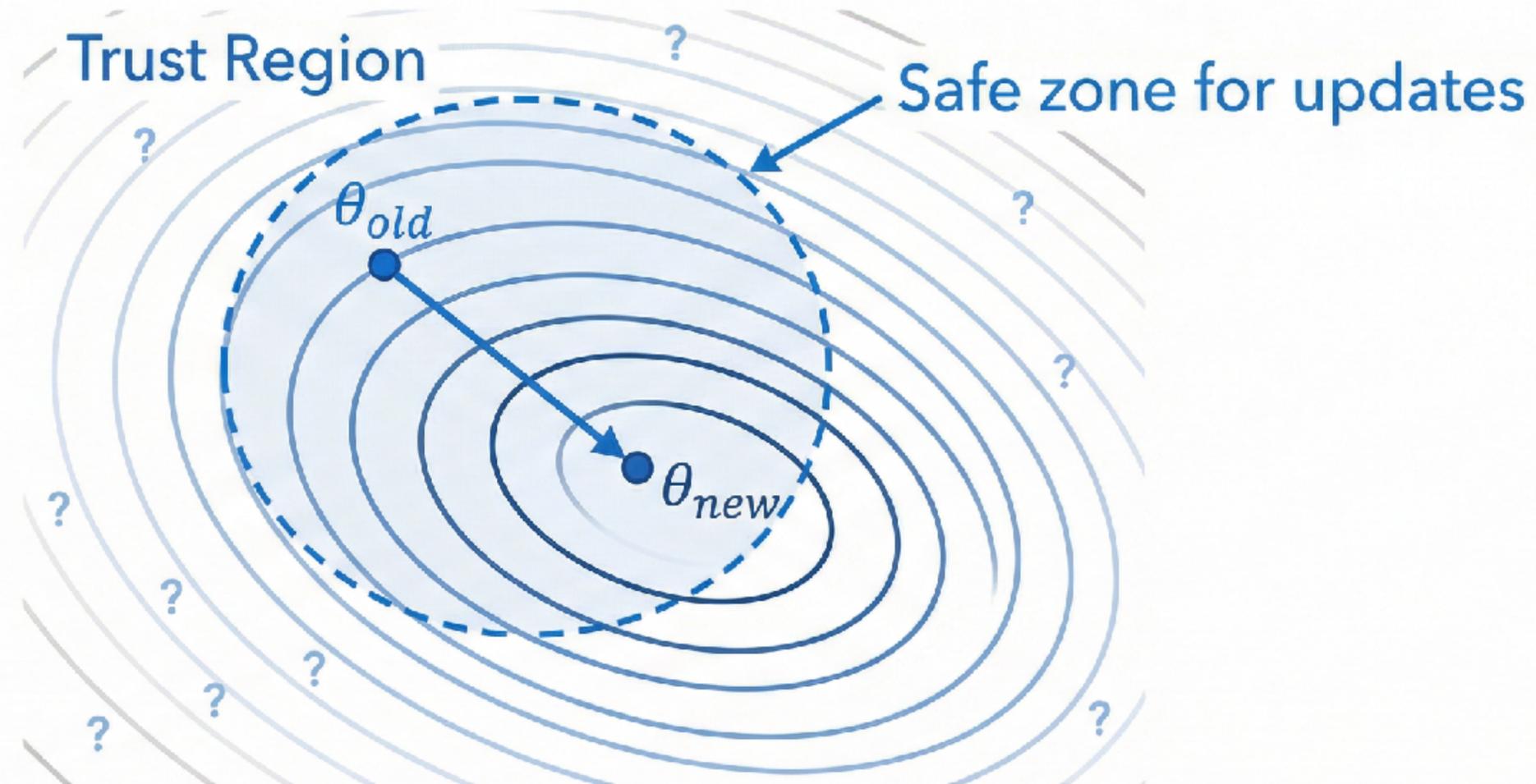
Off-Policy Algorithms can use **any** policy to collect experiences

Review: On + Off-Policy Learning

	On-Policy	Off Policy
Summary	Learns policy/value function based on policy used during training	Learns policy independent of policy used to collect experiences during training
Algorithms	SARSA, Policy Gradient, Actor Critic, PPO	Q-Learning, Off-policy Actor-Critic, Deep Deterministic Policy Gradient (DDPG)

Any Idea?

The Key Insight: Trust Regions

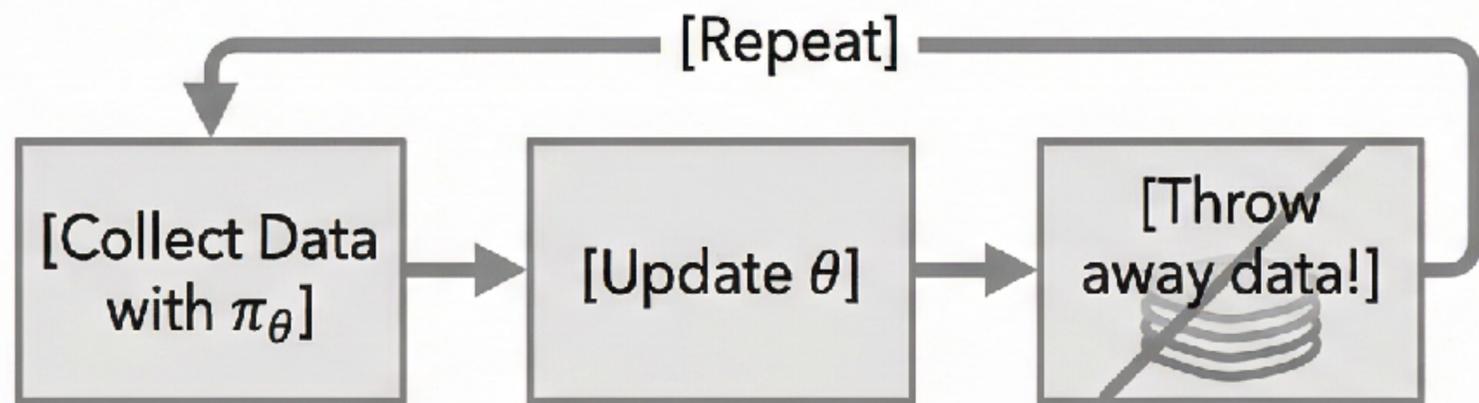


Core Idea: Only trust our local approximation within a small region

- Local linear/quadratic approximations are only accurate nearby
- Constrain updates to stay within the "trusted" neighborhood
- Guarantees monotonic improvement (with high probability)

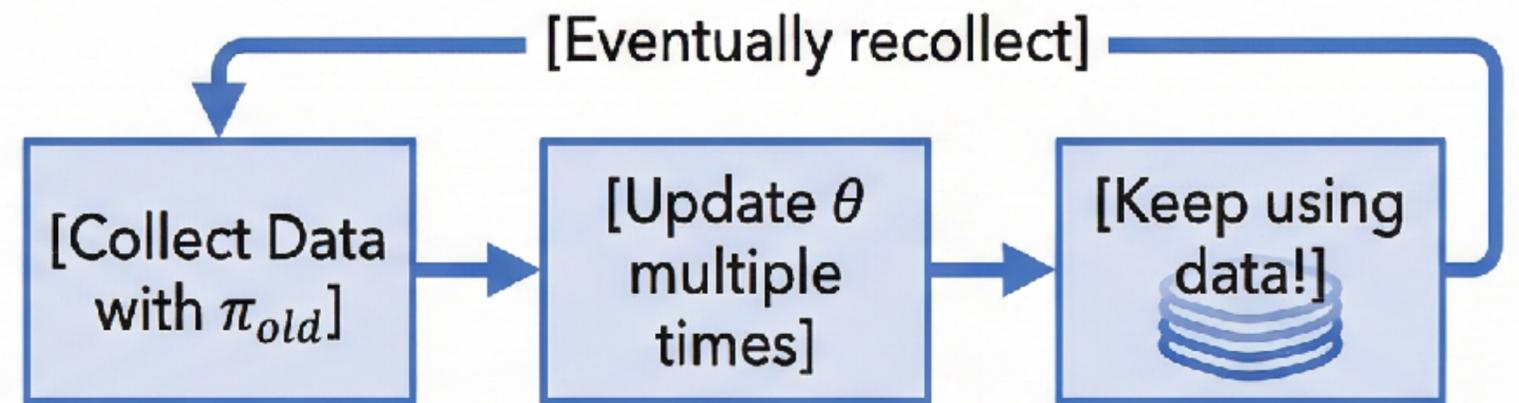
The Off-Policy Trick: Reusing Old Data

Pure On-Policy (e.g., REINFORCE, A2C)



- Must collect new data after EVERY update
- Old trajectories are "stale" - from wrong distribution
- Sample inefficient
- Simple but wasteful

Off-Policy via Importance Sampling (PPO/TRPO)

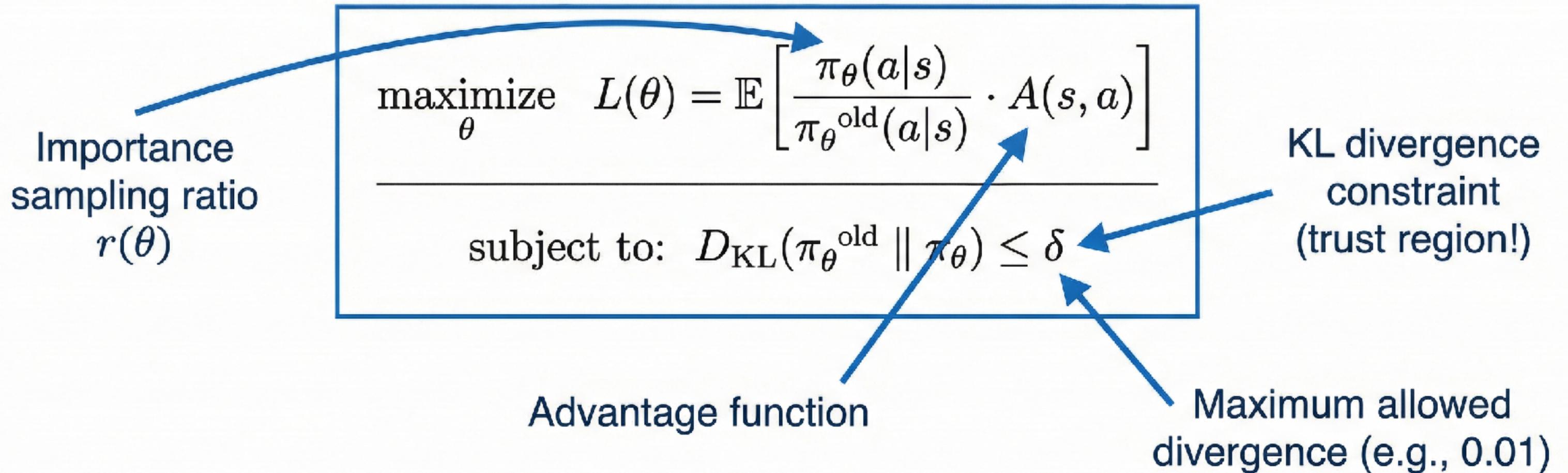


- Reuse data for K epochs of updates
- Correct distribution mismatch via importance ratio
- $r(\theta) = \pi_\theta(a|s) / \pi_{old}(a|s)$
- Much more sample efficient!

The importance ratio $r(\theta)$ corrects for the fact that data came from π_{old} , not π_θ

TRPO: Trust Region Policy Optimization

Schulman et al., 2015



The KL constraint ensures the new policy doesn't deviate too far from the old one

Importance Sampling: The Math

We want to estimate:

$$\mathbb{E}_{a \sim \pi_{\theta}}[A(s, a)]$$

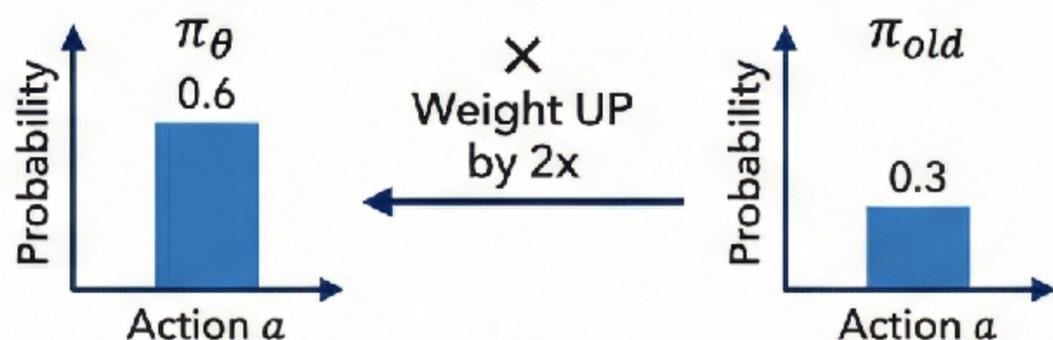
⚠ But our data came from π_{old} !

✅ Solution: Importance Sampling

$$\begin{aligned}\mathbb{E}_{a \sim \pi_{\theta}}[A(s, a)] &= \mathbb{E}_{a \sim \pi_{old}} \left[\frac{\pi_{\theta}(a|s)}{\pi_{old}(a|s)} \cdot A(s, a) \right] \\ &= \mathbb{E}_{a \sim \pi_{old}} [r(\theta) \cdot A(s, a)]\end{aligned}$$

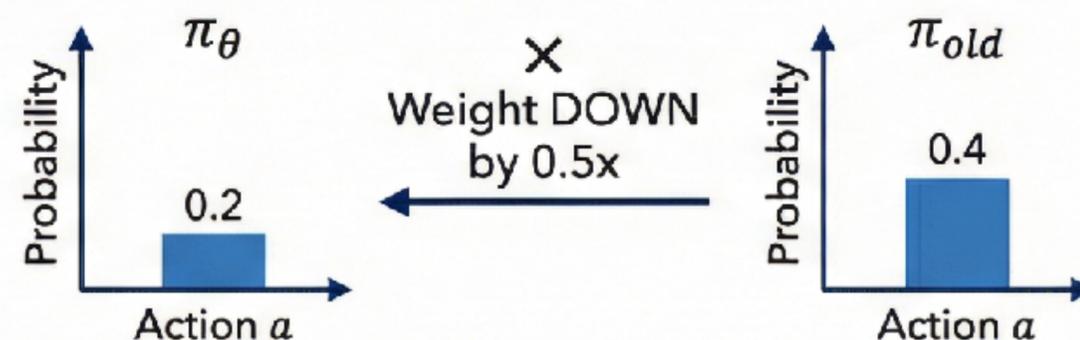
Intuition

If $\pi_{\theta}(a) = 0.6$ and $\pi_{old}(a) = 0.3$, then $r(\theta) = 2$



Action a was sampled less often than π_{θ} would
→ weight it UP by 2x

If $\pi_{\theta}(a) = 0.2$ and $\pi_{old}(a) = 0.4$, then $r(\theta) = 0.5$



Action a was sampled more often than π_{θ} would
→ weight it DOWN by 0.5x

⚠ Danger: If π_{θ} drifts too far from π_{old} , importance weights explode!
That's why we need clipping (PPO) or KL constraint (TRPO)

Enter PPO: A Simpler Approach

Schulman et al., 2017

TRPO Recap

- Constrained optimization
- Second-order methods (Fisher matrix)
- Conjugate gradient solver
- Line search
- Complex implementation

Simplification



PPO Solution

- First-order methods only!
- Simple clipping mechanism
- No KL constraint needed
- Standard SGD/Adam
- Easy to implement

“PPO attains the data efficiency and reliable performance of TRPO while using only first-order optimization”

The Complete PPO Objective

$$L^{\text{PPO}}(\theta) = \mathbb{E} \left[L^{\text{CLIP}}(\theta) - c_1 \cdot L^{\text{VF}}(\theta) + c_2 \cdot \mathcal{S}[\pi_\theta] \right]$$

$L^{\text{CLIP}}(\theta)$ - Policy Loss

The clipped surrogate objective

Encourages good actions, discourages bad

$L^{\text{VF}}(\theta)$ - Value Function Loss

$$L^{\text{VF}} = (V_\theta(s) - V_{\text{target}})^2$$

Trains the critic network

c_1 typically = 0.5

$\mathcal{S}[\pi_\theta]$ - Entropy Bonus

$$\mathcal{S} = - \sum \pi(a|s) \log \pi(a|s)$$

Encourages exploration

c_2 typically = 0.01

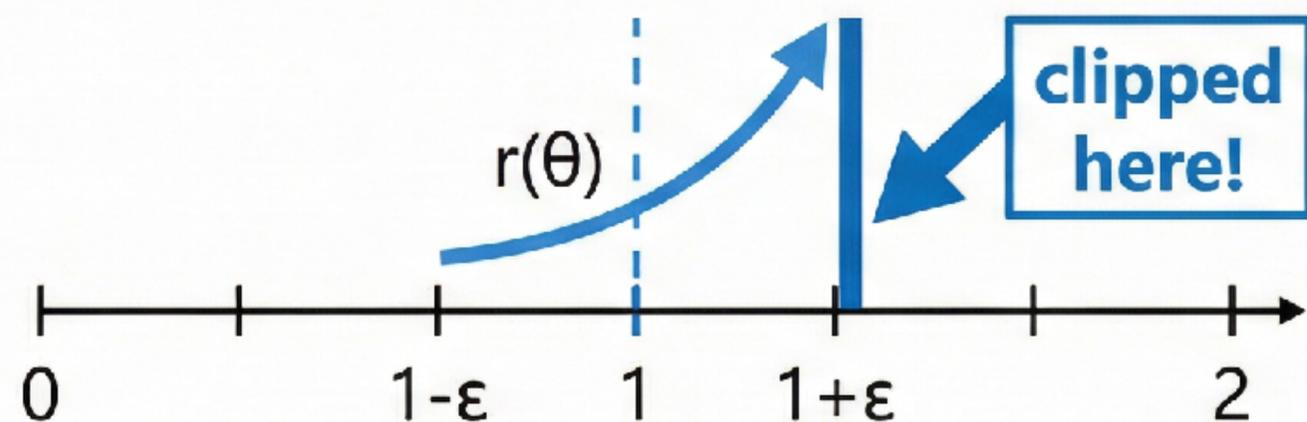
In actor-critic setup, policy and value networks often share parameters

PPO: The Clipped Objective

$$\mathcal{L}^{\text{CLIP}}(\theta) = \mathbb{E} \left[\min(r(\theta) \cdot A, \text{clip}(r(\theta), 1 - \varepsilon, 1 + \varepsilon) \cdot A) \right]$$

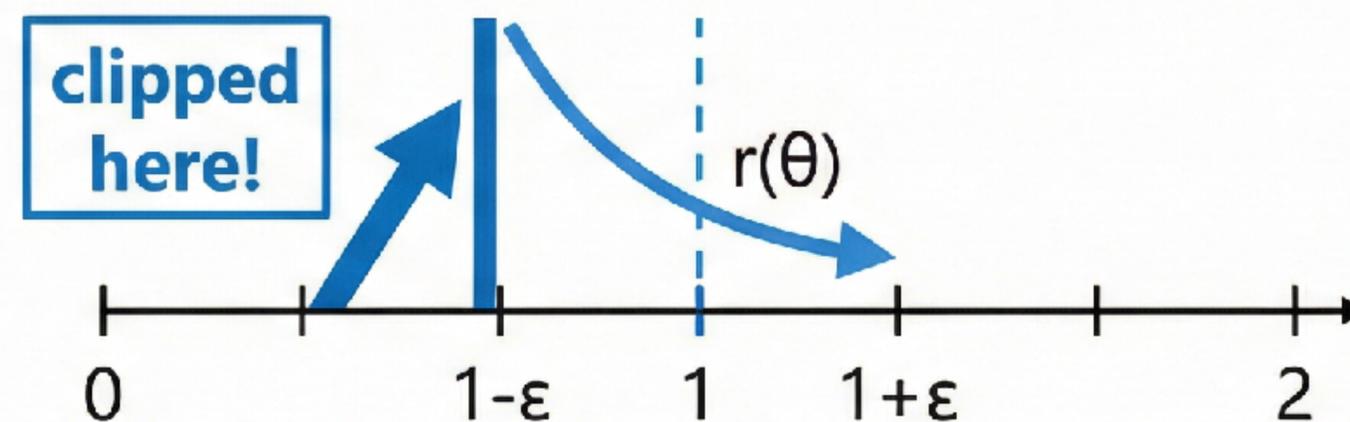
$$\text{where } r(\theta) = \frac{\pi_{\theta}(a|s)}{\pi_{\theta_{\text{old}}}(a|s)}$$

When $A > 0$ (good action)



Can't get more credit than $(1 + \varepsilon) \cdot A$ ✓

When $A < 0$ (bad action)



Can't be penalized more than $(1 - \varepsilon) \cdot A$ ✗

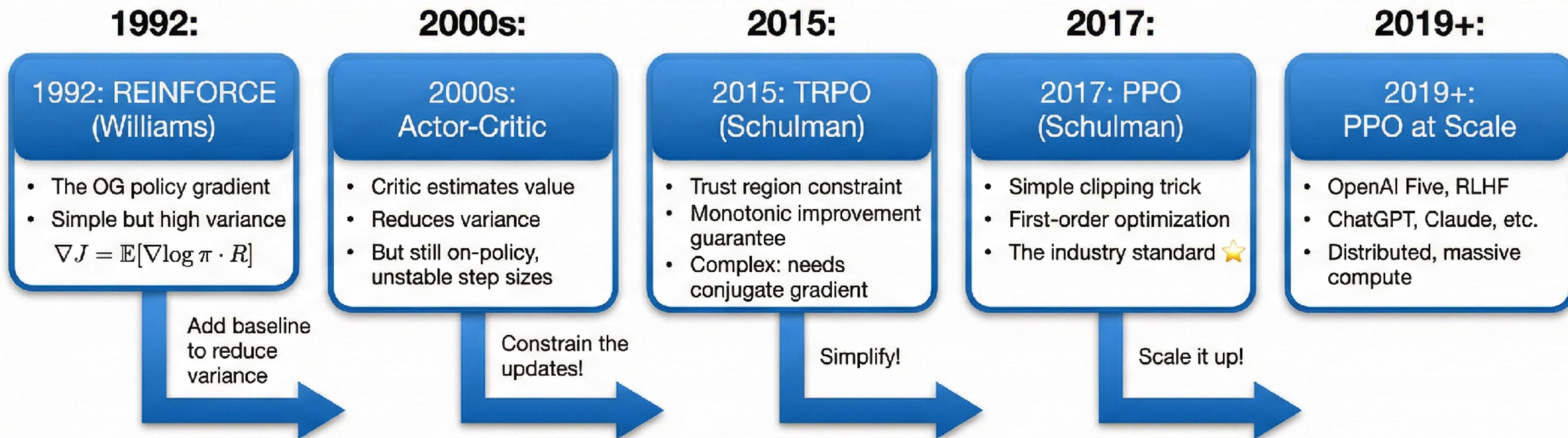
Typical value: $\varepsilon = 0.2$

TRPO vs PPO: Head-to-Head

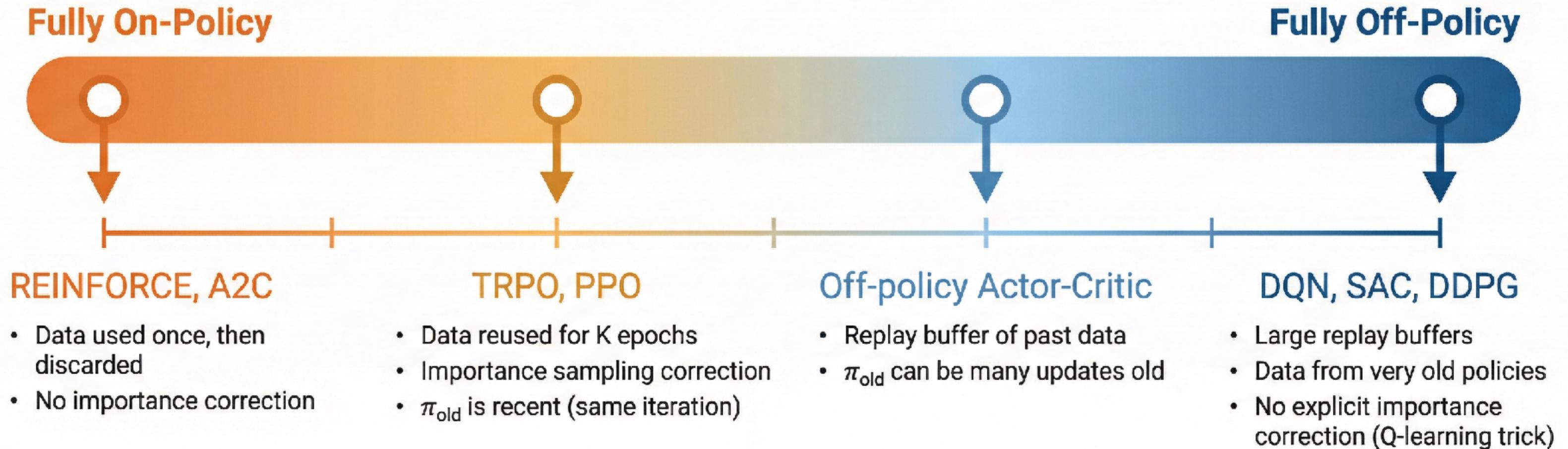
	TRPO	PPO
Constraint	Hard KL constraint $D_{KL} \leq \delta$	Soft clipping $r(\theta) \in [1-\epsilon, 1+\epsilon]$
Optimization	Second-order (Natural gradient)	First-order (Adam/SGD)
Implementation	Complex ~500 lines	Simple ~50 lines
Computation	Expensive (Fisher matrix, CG)	Cheap (Standard backprop)
Sample Efficiency	Similar	Similar
Performance	Strong	Strong (often slightly better)
Popularity	Research	Industry standard (OpenAI, etc.)

PPO: Same benefits as TRPO, much simpler to implement and tune

The Evolution of Policy Optimization



The On-Policy to Off-Policy Spectrum



Key trade-off: More off-policy → more sample efficient, but harder to stabilize



Algorithm Selection Guide: Which One to Use?

Algorithm	Action Space	Sample Eff.	Stability	Best For
DQN	Discrete ONLY	High ★★ ★ (replay)	Medium	<ul style="list-style-type: none">Atari, discrete gamesSimple to implement
A2C	Both	Low ★ (on-policy)	High ★★ ★	<ul style="list-style-type: none">Fast prototypingParallel environments
PPO	Both	Medium ★★	High ★★ ★	<ul style="list-style-type: none">General purpose ✓DEFAULT CHOICE
TRPO	Both	Medium ★★	High ★★ ★	<ul style="list-style-type: none">Research/theoryWhen guarantees needed
DDPG/TD3 SAC	Continuous ONLY	High ★★ ★ (replay)	Low Med-High	<ul style="list-style-type: none">Robotics, controlContinuous actions

🎯 Default Choice: Start with PPO

- Works for both discrete and continuous
- Stable and reliable
- Easy to tune
- Industry standard

See you on Friday!