

PPO PPOxFamily

第七讲：挖掘黑科技

主办



承办



上海人工智能实验室
Shanghai Artificial Intelligence Laboratory

协办

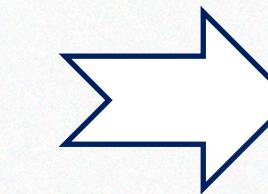
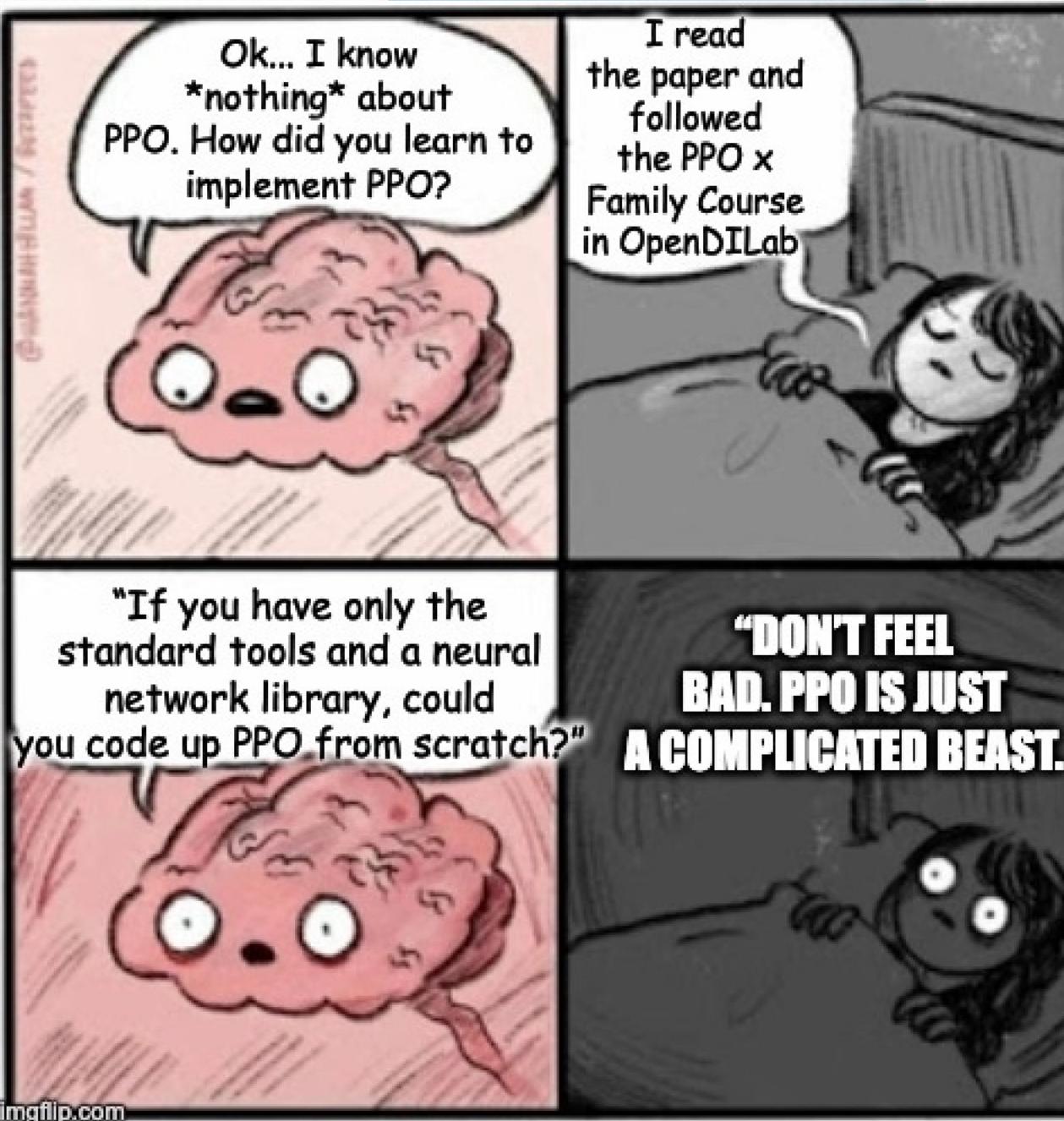


支持

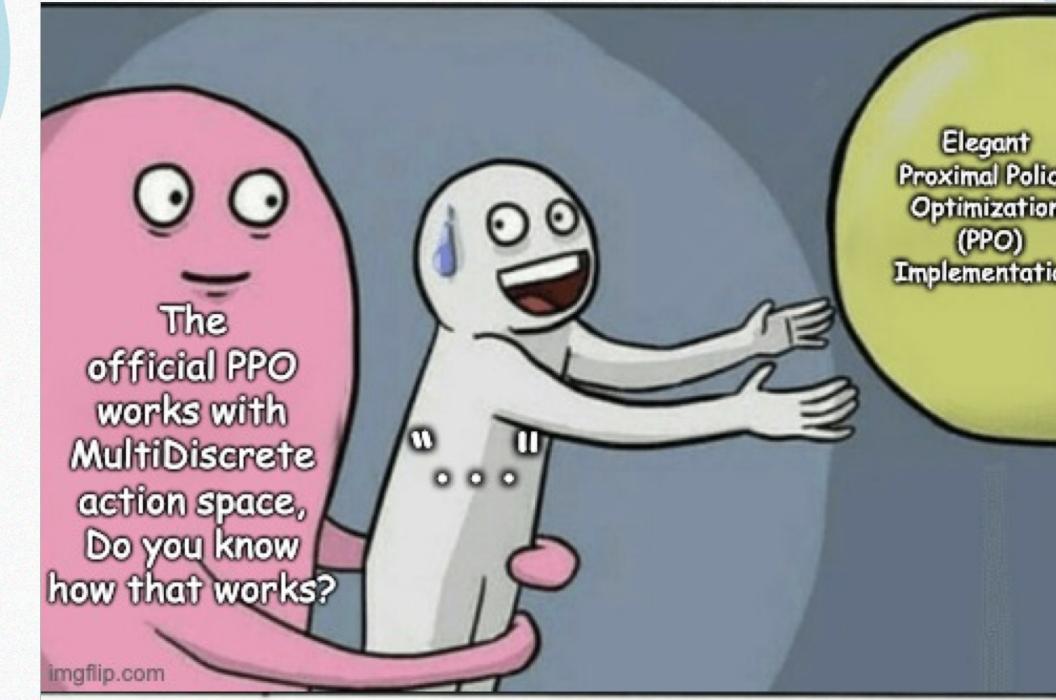
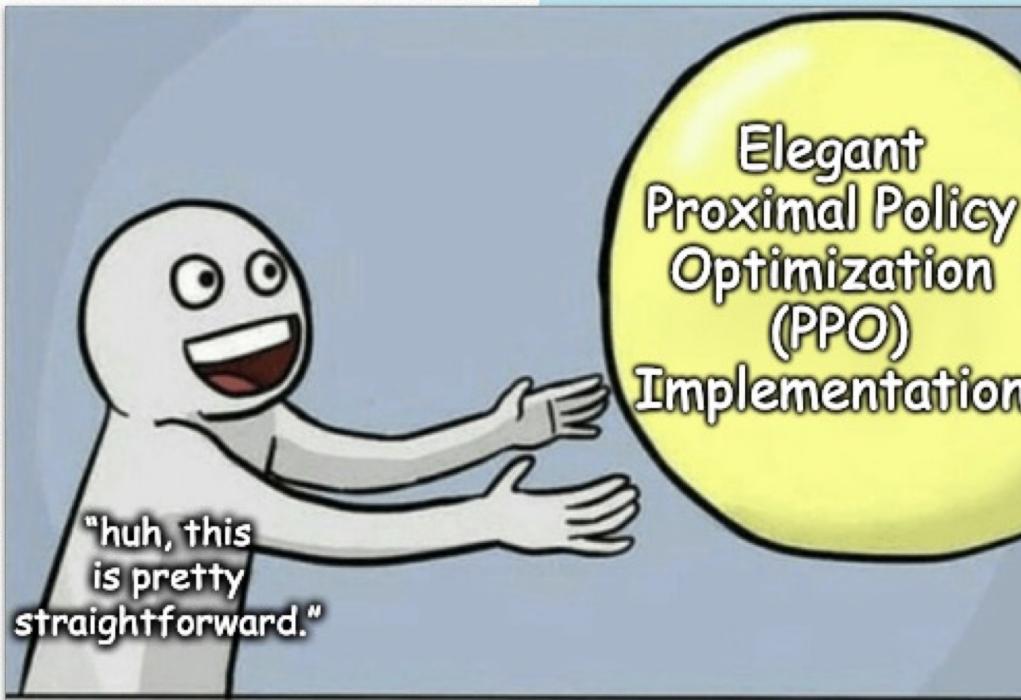




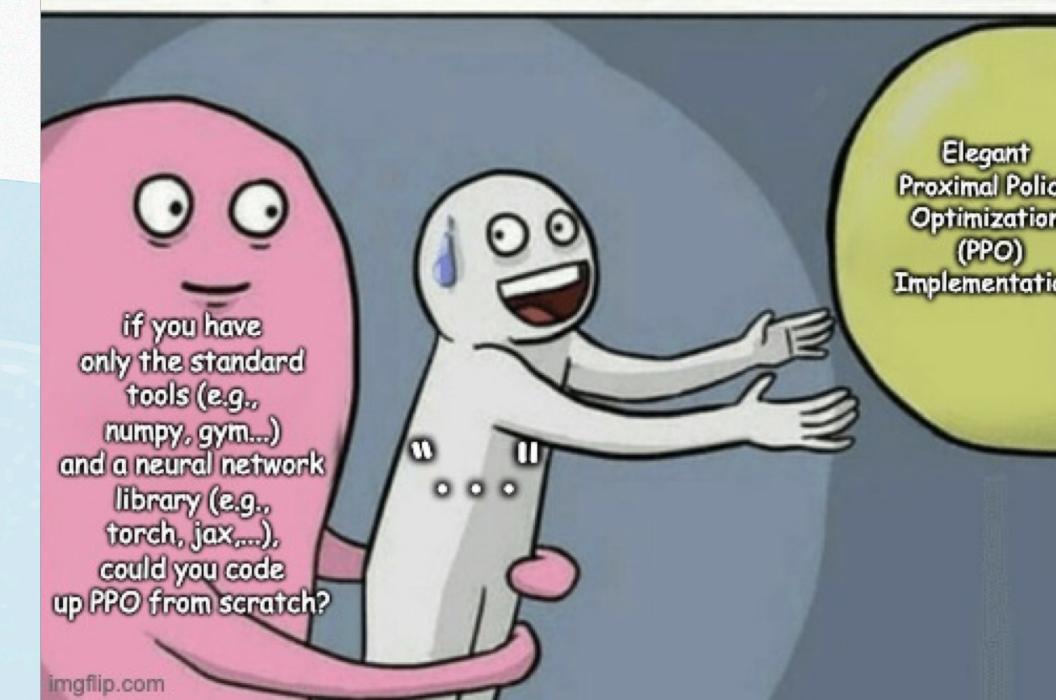
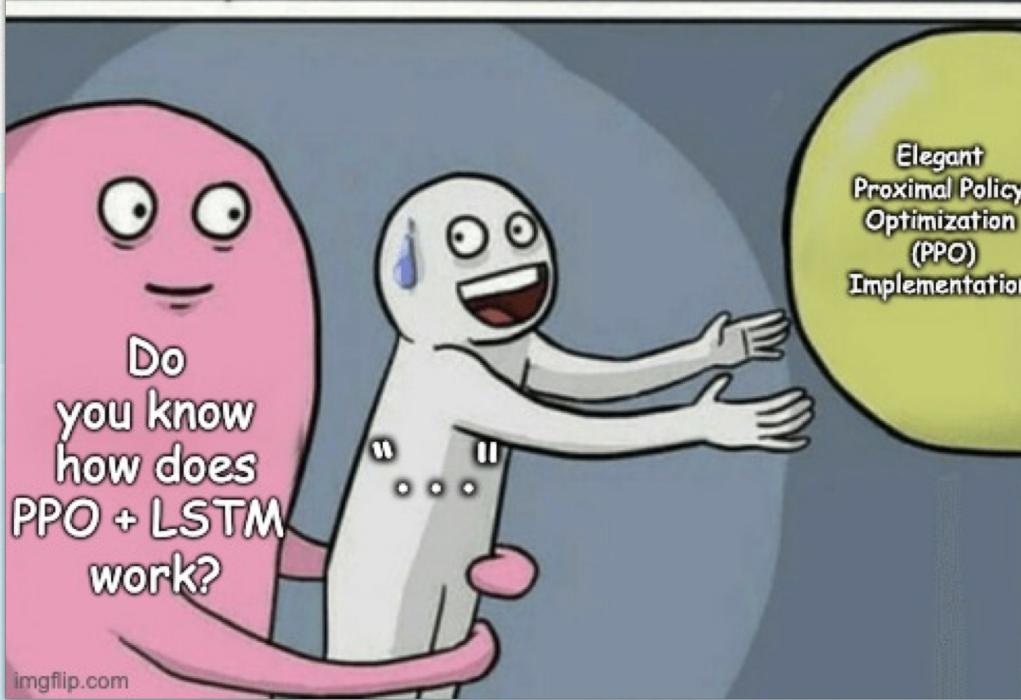
问题：强化学习究竟是不是玄学？



问题：別人家的PPO为什么这么强？



- 集成万千技巧 (Tricks)



- 应对大量环境 (Envs)

The Seven Realms of PPO

调优PPO的七重境界

合理运用随机：
Dynamic Seed

亡羊补牢：
Extra Clip

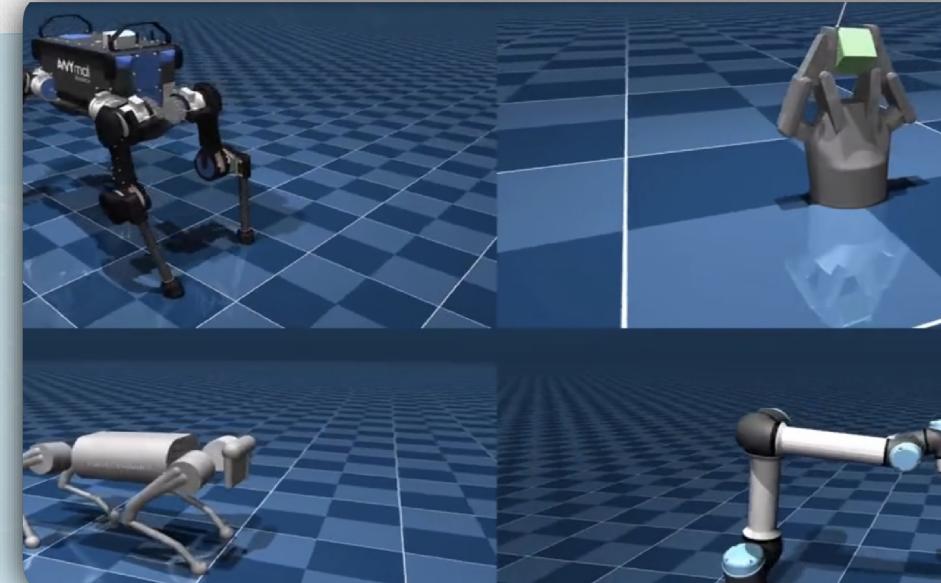
笨鸟先飞：
Network Initialization

免死金牌：
Grad Clip

权衡偏差与方差：
GAE

稳扎稳打：
Recompute Advantage

巧夺天工：
Entropy Balance



双重境界：权衡偏差与方差的 GAE

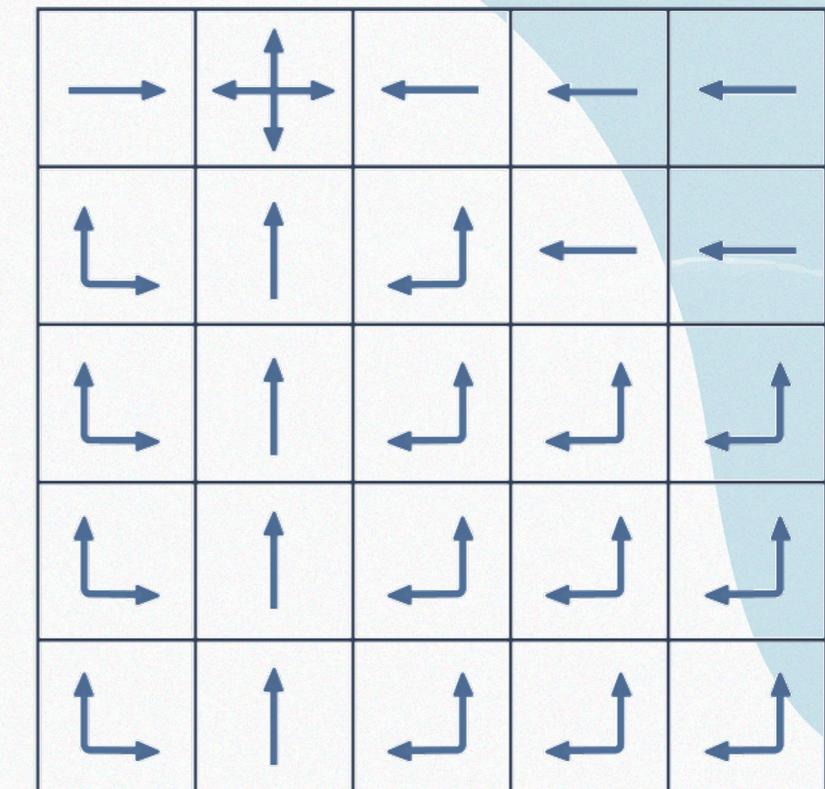
价值函数

$$V^\pi(s) = \mathbb{E}_{\tau \sim \pi}[G(\tau) | s_t = s]$$

$$Q^\pi(s, a) = \mathbb{E}_{\tau \sim \pi}[G(\tau) | s_t = s, a_t = a]$$

示例：格子世界 (Grid World)

22.0	24.4	22.0	19.4	17.5
19.8	22.0	19.8	17.8	16.0
17.8	19.8	17.8	16.0	14.4
16.0	17.8	16.0	14.4	13.0
14.4	16.0	14.4	13.0	11.7



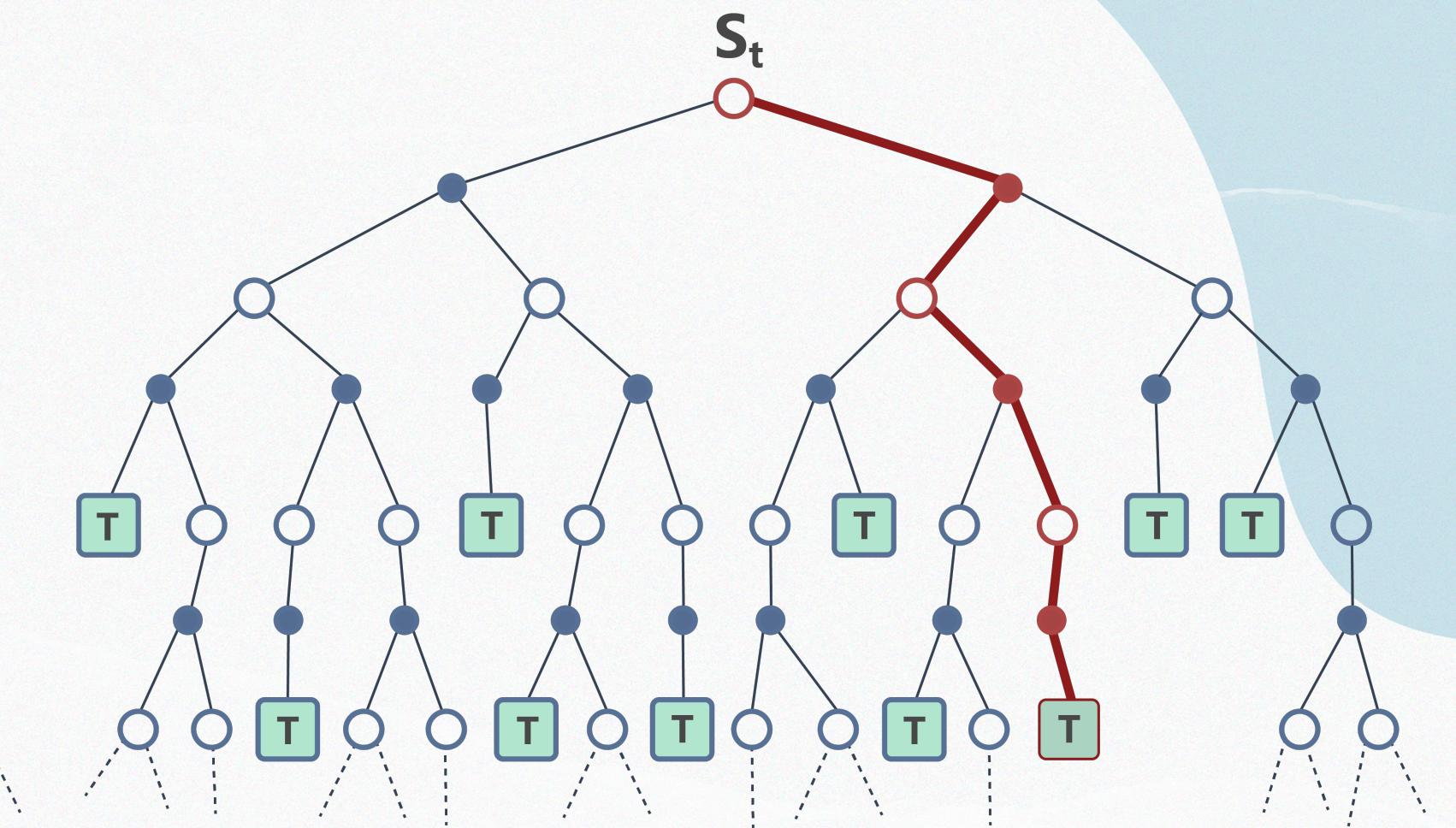
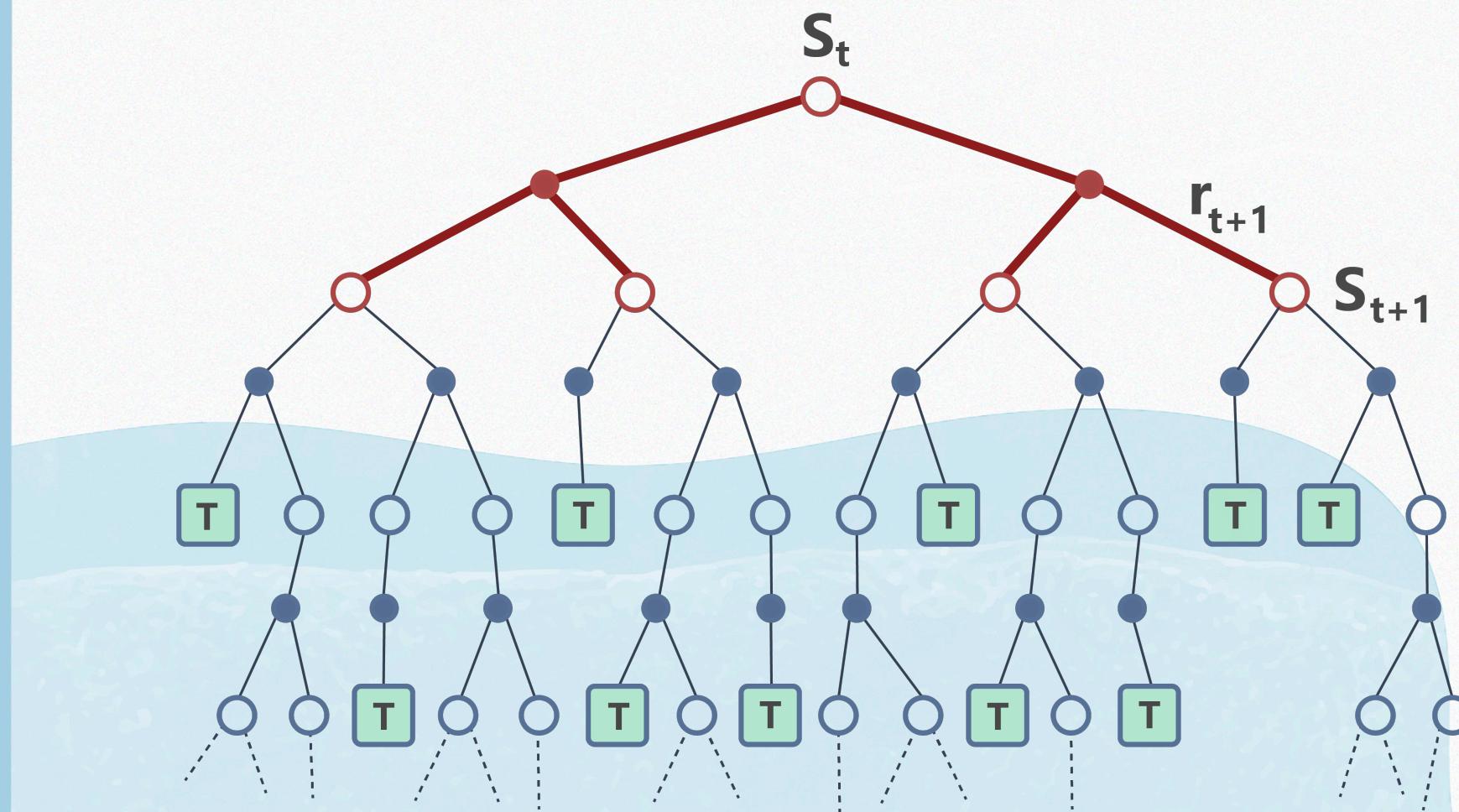
双重境界：权衡偏差与方差的 GAE

动态规划 (DP)

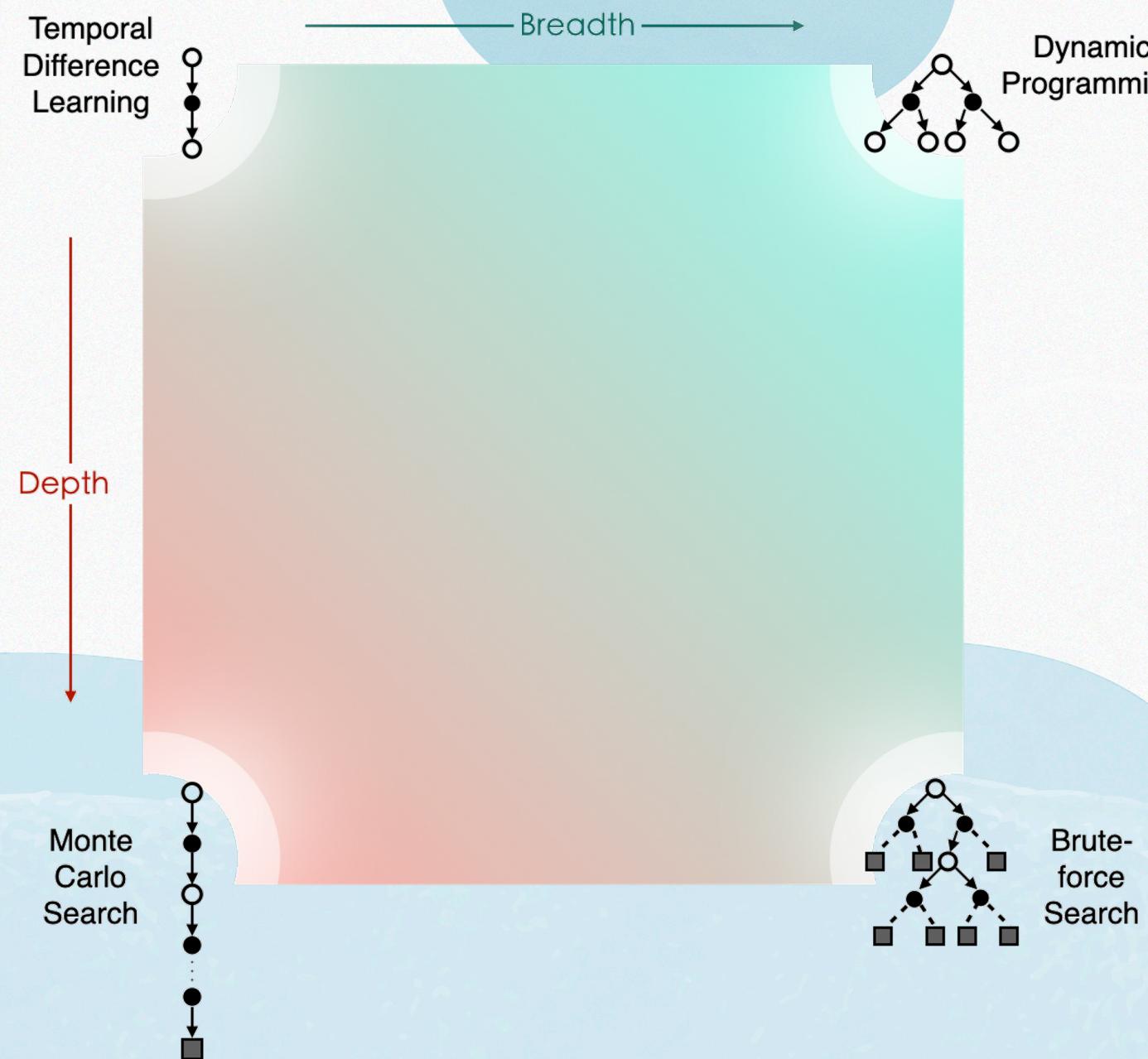
$$V(s_t) \leftarrow \mathbb{E}_{\tau \sim \pi}[r_t + \gamma V(s_{t+1})]$$

蒙特卡洛 (MC)

$$V(s_t) \leftarrow V(s_t) + \alpha(G_t(\tau) - V(s_t))$$



双重境界：权衡偏差与方差的 GAE



$$V(S_t) \leftarrow V(S_t) + \alpha [R_{t+1} + \gamma V(S_{t+1}) - V(S_t)]$$

Legend:

- New value of state t
- Former estimation of value of state t
- Learning Rate
- Reward
- Discounted value of next state

Method Name	Bootstrapping	Sampling
DP	✓	✗
MC	✗	✓
TD	✓	✓

双重境界：权衡偏差与方差的 GAE

偏差 (Bias)

1-step TD
and TD(0)



2-step TD



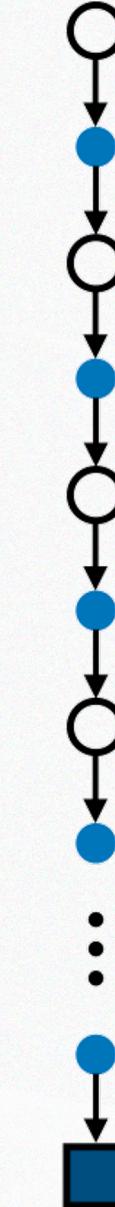
3-step TD



n-step TD



∞ -step TD
and Monte Carlo



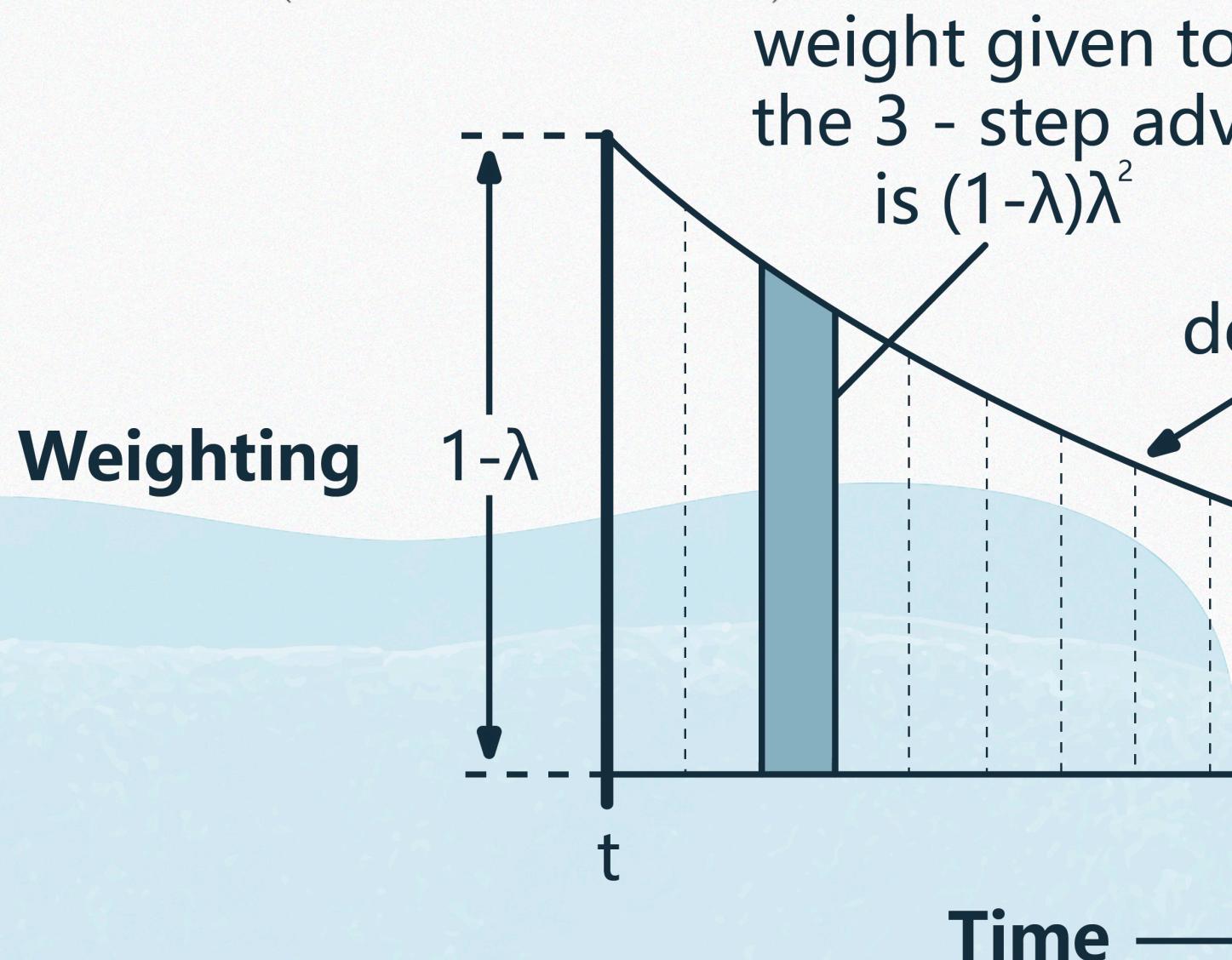
方差
(Variance)

- 如何更好地平衡方差和偏差？

双重境界：权衡偏差与方差的 GAE

$$\hat{A}_t^{(k)} := \sum_{l=0}^{k-1} \gamma^l \delta_{t+l}^V = -V(s_t) + r_t + \gamma r_{t+1} + \cdots + \gamma^{k-1} r_{t+k-1} + \gamma^k V(s_{t+k})$$

$$\hat{A}_t^{\text{GAE}(\gamma, \lambda)} := (1 - \lambda) \left(\hat{A}_t^{(1)} + \lambda \hat{A}_t^{(2)} + \lambda^2 \hat{A}_t^{(3)} + \dots \right)$$



$$\text{GAE}(\gamma, 0) : \quad \hat{A}_t := \delta_t = r_t + \gamma V(s_{t+1}) - V(s_t)$$

$$\text{GAE}(\gamma, 1) : \quad \hat{A}_t := \sum_{l=0}^{\infty} \gamma^l \delta_{t+l} = \sum_{l=0}^{\infty} \gamma^l r_{t+l} - V(s_t)$$

双重境界：权衡偏差与方差的 GAE

Algorithm 1 Proximal Policy Optimization (PPO)

- 1: Input: initialize policy (actor) parameters θ and value (critic) parameters ϕ . Training epochs per collect E , trajectory length T , batch size B .
- 2: **for** $k = 0, 1, 2, \dots$ **do**
- 3: Collect set of trajectories $\mathcal{D}_k = \{\tau_i\}$ by running policy π_{θ_k} interaction with environment.
- 4: Compute trajectory target return estimates \hat{R}_t . (e.g. n-step TD or other methods)
- 5: Compute advantage estimates \hat{A}^{θ_k} with value V_{ϕ_k} . (e.g. GAE or other methods)
- 6: **for** $e = 0, 1, \dots, E - 1$ **do**
- 7: **for** minibatch: $b \in \mathcal{D}_k$ **do**
- 8: Update the policy by maximizing the PPO-Clip objective with SGD:

$$L_\theta = \frac{1}{B \cdot T} \sum_{\tau \in b} \sum_{t=0}^{T-1} \min(r(\theta) \hat{A}^{\theta_k}(s_t, a_t), \text{clip}(r(\theta), 1 - \epsilon, 1 + \epsilon) \hat{A}^{\theta_k}(s_t, a_t))$$

$$r(\theta) = \frac{\pi_\theta(a_t | s_t)}{\pi_{\theta_k}(a_t | s_t)}$$

- 9: Fit the value by regression on mean-squared error with SGD:

$$L_\phi = \frac{1}{B \cdot T} \sum_{\tau \in b} \sum_{t=0}^{T-1} (V_\phi(s_t) - \hat{R}_t)^2$$

- 10: **end for**
- 11: **end for**
- 12: **end for**

- 收集数据需要保留时间顺序
- 在收集完毕，训练之前计算
- 多个 trajectory 拼接后计算

```

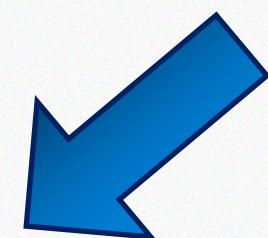
next_value *= (1 - done)
delta = reward + gamma * next_value - value
factor = gamma * lambda_ * (1 - traj_flag)
adv = torch.zeros_like(value)
gae_item = torch.zeros_like(value[0])

for t in reversed(range(reward.shape[0])):
    gae_item = delta[t] + factor[t] * gae_item
    adv[t] = gae_item

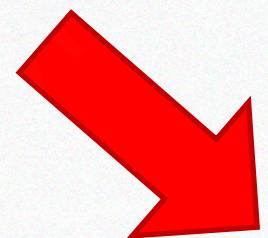
```

二重境界：稳扎稳打之 Recompute

$$L_{\theta} = \frac{1}{B \cdot T} \sum_{\tau \in \mathbf{b}} \sum_{t=0}^{T-1} \min(r(\theta) \hat{A}^{\theta_k}(s_t, a_t), \text{clip}(r(\theta), 1 - \epsilon, 1 + \epsilon) \hat{A}^{\theta_k}(s_t, a_t))$$



- 新老策略尽可能地接近，IS 和 adv 可能存在的偏差就更小



- PPO 常用实现中包含两重循环
 - 第一重：训练多个epoch
 - 第二重，一个 epoch 中划分 N 个 minibatch
- 收集一次数据，训练 epoch \times N 个迭代
- 只要更新一个迭代，新老策略就已经产生了不一致，这与 PPO 的设计初衷相冲突

双重境界：稳扎稳打之 Recompute

Algorithm 1 Proximal Policy Optimization (PPO)

```

1: Input: initialize policy (actor) parameters  $\theta$  and value (critic) parameters  $\phi$ . Training epochs per
   collect  $E$ , trajectory length  $T$ , batch size  $B$ .
2: for  $k = 0, 1, 2, \dots$  do
3:   Collect set of trajectories  $\mathcal{D}_k = \{\tau_i\}$  by running policy  $\pi_{\theta_k}$  interaction with environment.
4:   Compute trajectory target return estimates  $\hat{R}_t$ . (e.g. n-step TD or other methods)
5:   for  $e = 0, 1, \dots, E - 1$  do
6:     Recompute advantage estimates  $\hat{A}^{\theta_k}$  with the latest value  $V_\phi$  (e.g. GAE or other methods)
7:     for minibatch:  $b \in \mathcal{D}_k$  do
8:       Update the policy by maximizing the PPO-Clip objective with SGD:

$$L_\theta = \frac{1}{B \cdot T} \sum_{\tau \in b} \sum_{t=0}^{T-1} \min(r(\theta) \hat{A}^{\theta_k}(s_t, a_t), \text{clip}(r(\theta), 1 - \epsilon, 1 + \epsilon) \hat{A}^{\theta_k}(s_t, a_t))$$


$$r(\theta) = \frac{\pi_\theta(a_t | s_t)}{\pi_{\theta_k}(a_t | s_t)}$$

9:     Fit the value by regression on mean-squared error with SGD:

$$L_\phi = \frac{1}{B \cdot T} \sum_{\tau \in b} \sum_{t=0}^{T-1} (V_\phi(s_t) - \hat{R}_t)^2$$

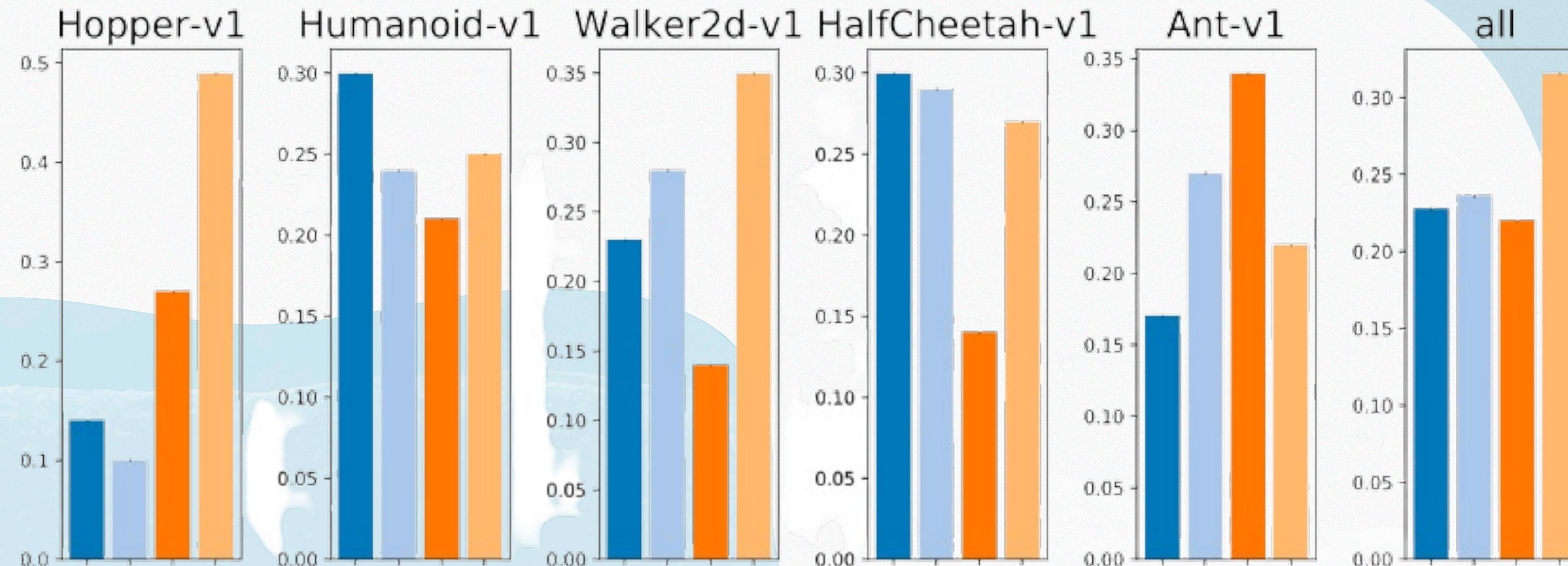
10:    end for
11:  end for
12: end for

```

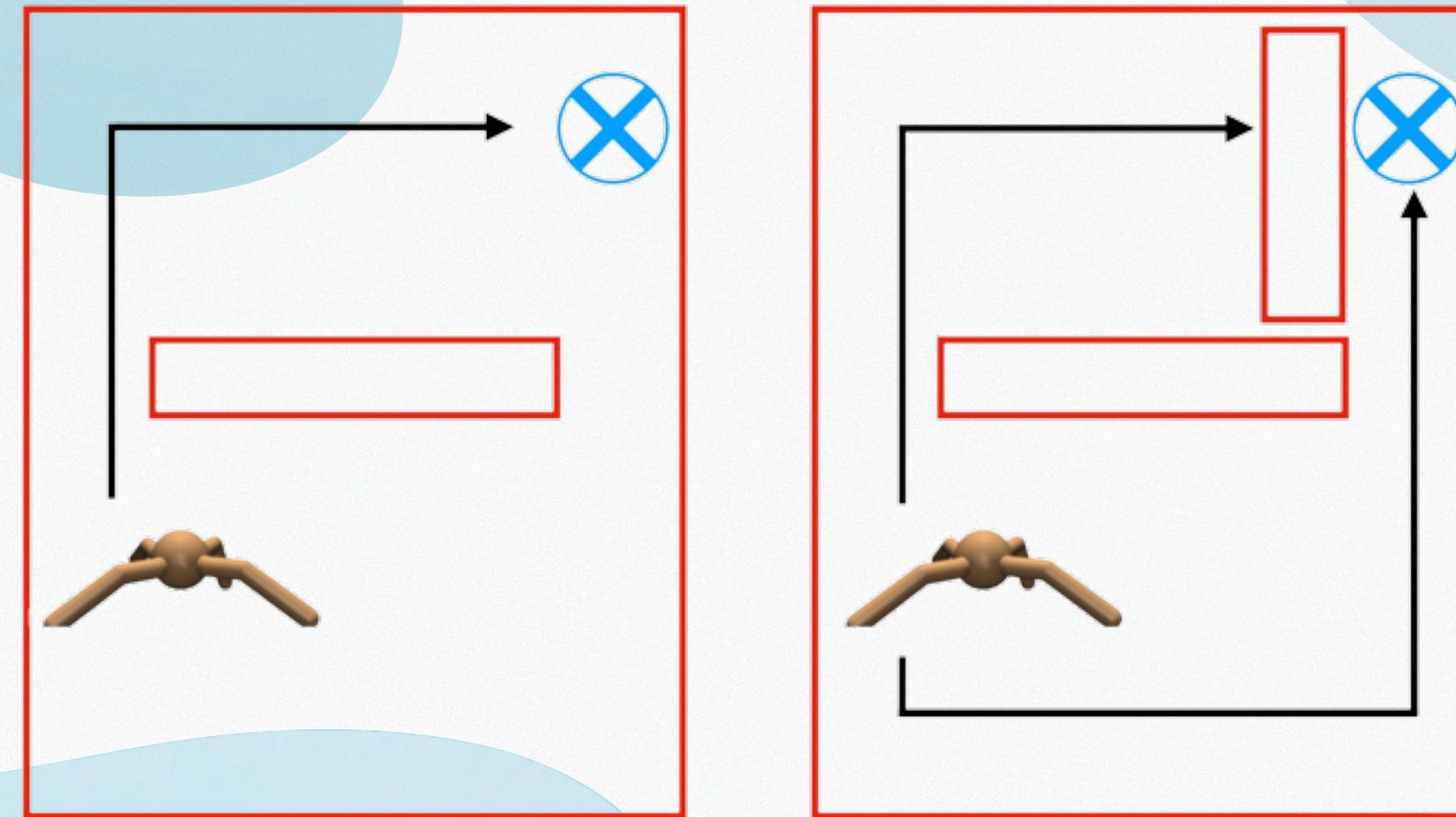
- Recompute 步骤无需梯度计算，节省时间和显存
- 在第一重循环内部整体计算。不需要额外维护时间顺序关系，每次计算完成后先 shuffle 数据再切分 minibatch 即可
- 在 epoch 较大时尤其有用

二重境界：稳扎稳打之 Recompute

从左到右依次是 : Fixed、Shuffle Trajectories、Shuffle Transitions、Shuffle Transitions + Recompute (整体效果最好)

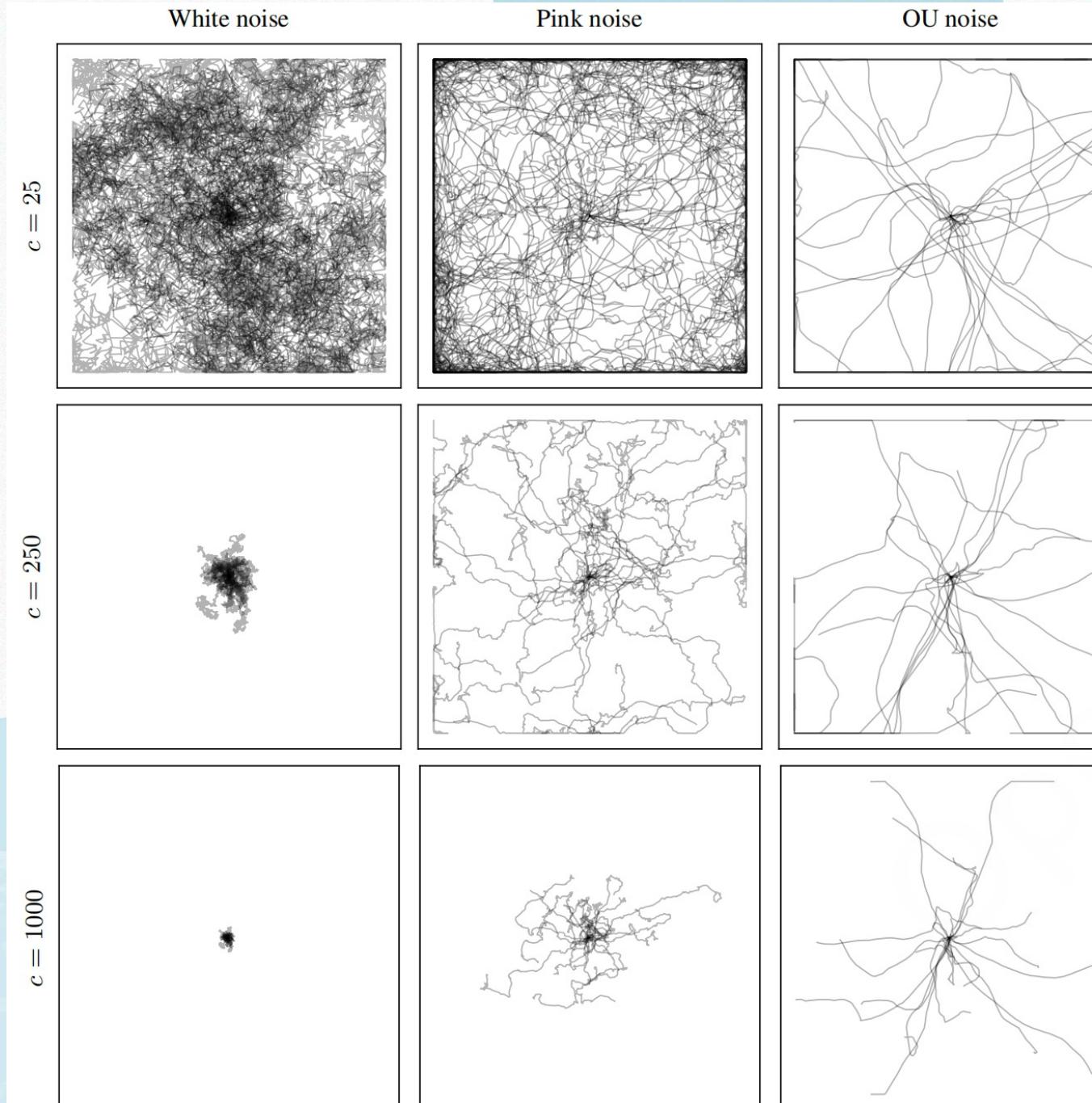


三重境界：巧夺天工之Entropy

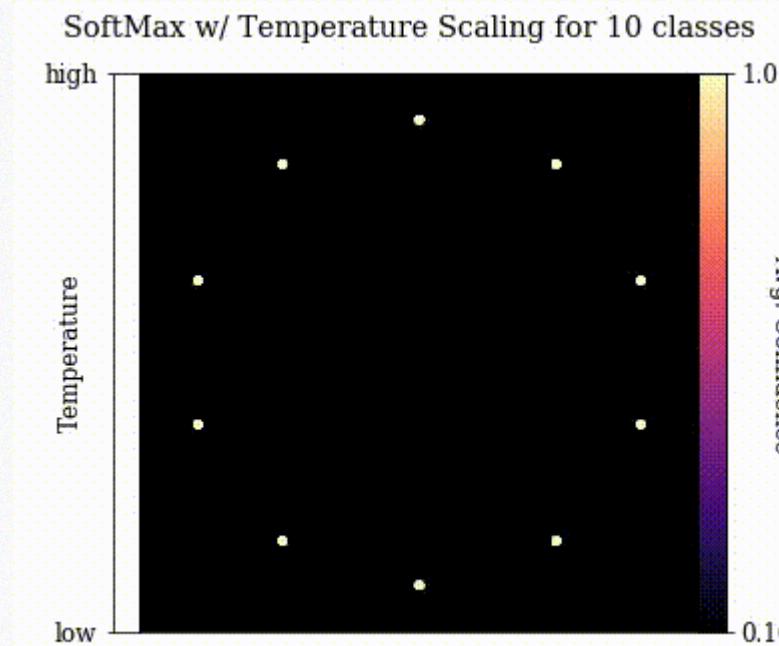


策略早熟：
策略以某种“投机取巧”的方式
找到了局部最优解

三重境界：巧夺天工之Entropy



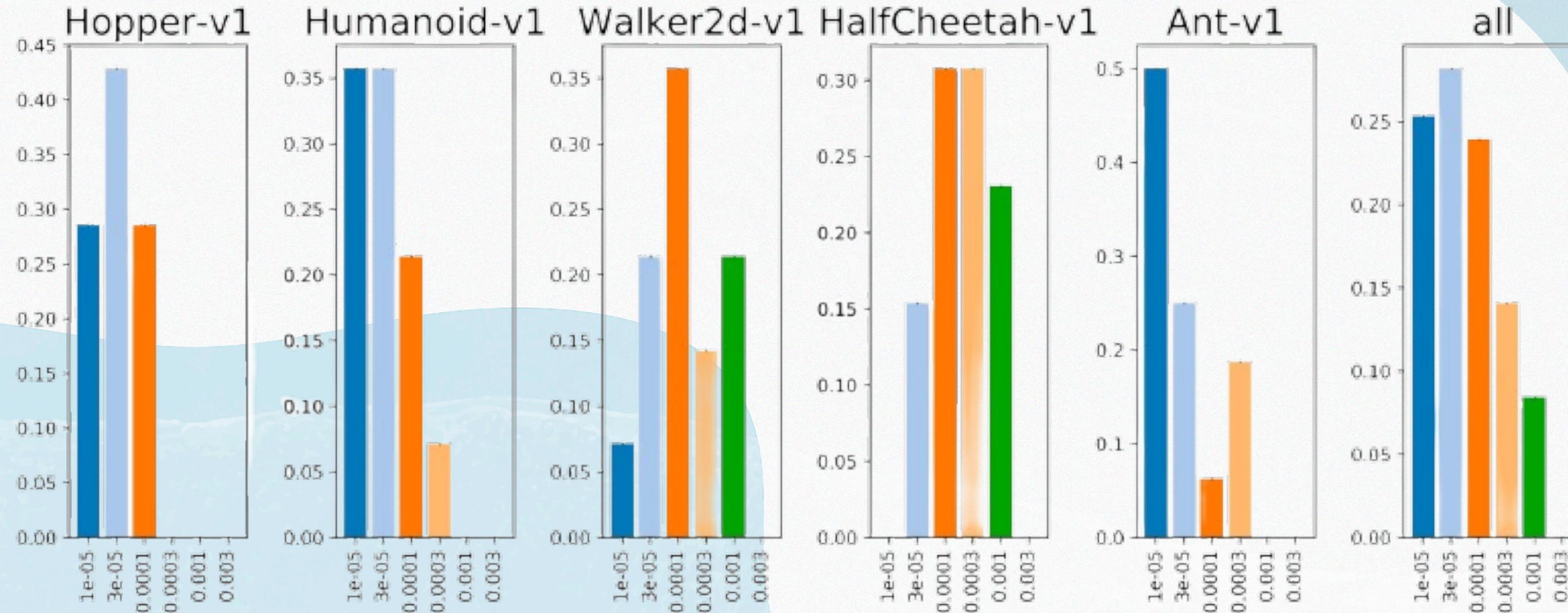
- 不同的环境类型，不同的环境设置，所需的策略探索能力不同
- PPO中控制策略探索能力的两个关键点：
 - 选择动作时的 **temperature**
 - 计算损失函数时的 **Entropy Bonus**



三重境界：巧夺天工之Entropy

合理的 Entropy Bonus 才能有效防止策略“早熟”并高效稳定收敛

更进一步，可以在训练周期内设置相应的 bonus 衰减（指数衰减 or 类比 SAC）

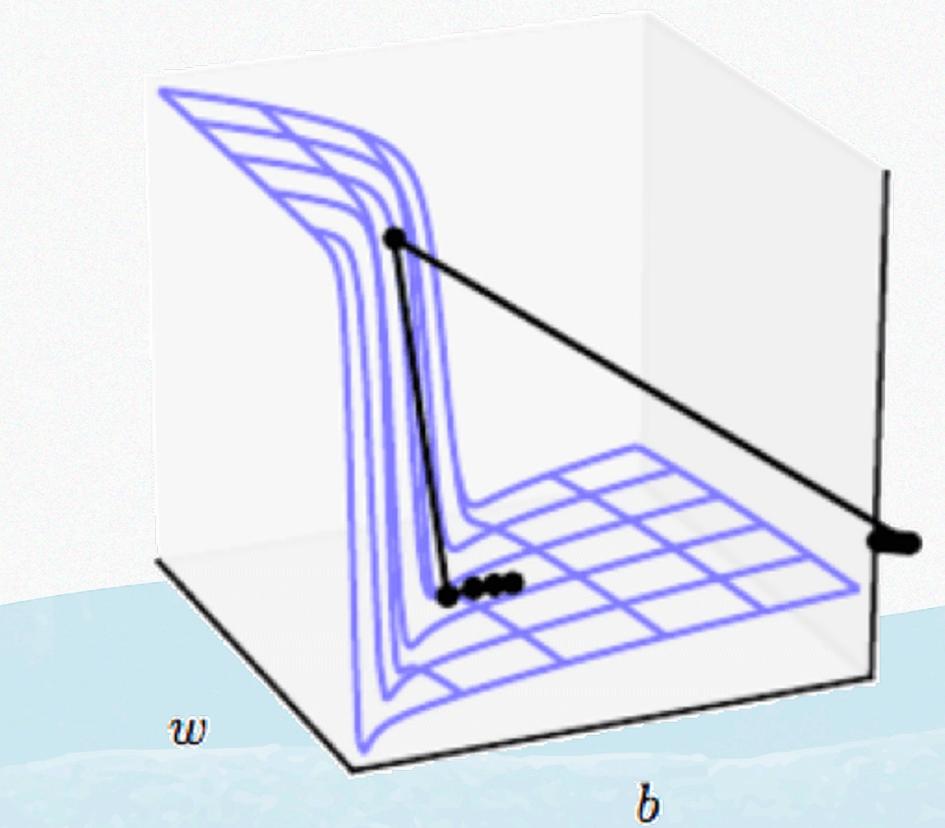


四重境界：免死金牌之 Grad Clip

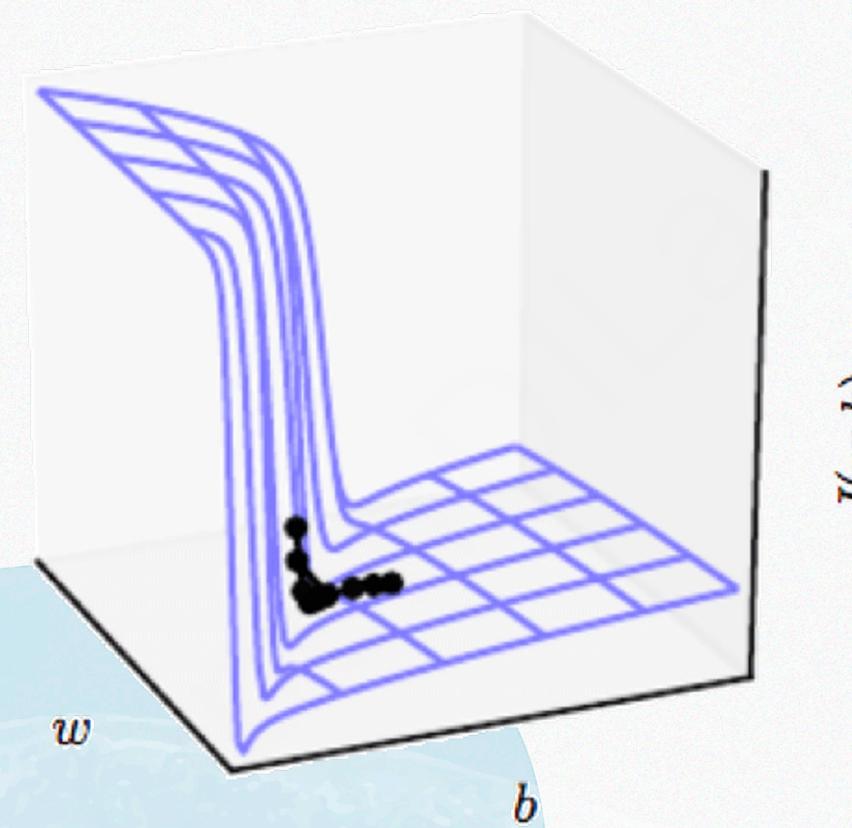


四重境界：免死金牌之 Grad Clip

Without clipping



With clipping



Standard gradient

$$g(\theta) = \frac{1}{N} \sum_{n=1}^N \nabla l_\theta(x_n, y_n)$$

Clipped gradient

$$\bar{g}_\tau(\theta) = \text{clip}_\tau(g(\theta))$$

Clip norm to τ

四重境界：免死金牌之Grad Clip

梯度的最大范数 (max norm) :

$$\text{total_norm}^\infty = \max_{\theta_i \in \Theta} |\text{grad}(\theta_i)|$$

梯度的 p-范数 (p-norm) :

$$\begin{aligned}\text{total_norm} &= \left(\sum_{\theta \in \Theta} \left(\sum_{\theta_i} \text{grad}(\theta_i)^p \right)^{\frac{1}{p}} \right)^p \\ &= \left(\sum_{\theta \in \Theta} \left(\sum_{\theta_i} \text{grad}(\theta_i)^p \right) \right)^{\frac{1}{p}}\end{aligned}$$

裁减系数 (1e-6用于避免分母为零):

$$\text{clip_coef} = \frac{\text{max_norm}}{\text{total_norm}}$$

将系数的最大值固定为1

如果 $\text{total_norm} < \text{max_norm}$, `torch.clamp` 操作将 `clip_coef` 的最大值固定为1, 所以 `clip_coef_clamped = 1`, 这样 `total_norm` 将保持不变。

如果 $\text{total_norm} > \text{max_norm}$, 将对原来的梯度进行裁减, 使得裁减后的梯度对应的 `total_norm` 的大小为 `max_norm`:

$$\begin{aligned}\text{total_norm}' &= \left(\sum_{\theta \in \Theta} \left(\sum_{\theta_i} \left(\text{grad}(\theta_i) \cdot \frac{\text{max_norm}}{\text{total_norm}} \right)^p \right)^{\frac{1}{p}} \right)^p \\ &= \frac{\left(\sum_{\theta \in \Theta} \left(\sum_{\theta_i} \text{grad}(\theta_i)^p \right)^{\frac{1}{p}} \right)^p}{\text{total_norm}} \cdot \text{max_norm} \\ &= \text{max_norm}\end{aligned}$$

```
16     if norm_type == inf:
17         norms = [g.detach().abs().max().to(device) for g in grads]
18         total_norm = norms[0] if len(norms) == 1 else torch.max(torch.stack(norms))
```

```
19     else:
20         total_norm = torch.norm(torch.stack([torch.norm(g.detach(), norm_type).to(device) for g in grads]), norm_type)
```

```
21     clip_coef = max_norm / (total_norm + 1e-6)
```

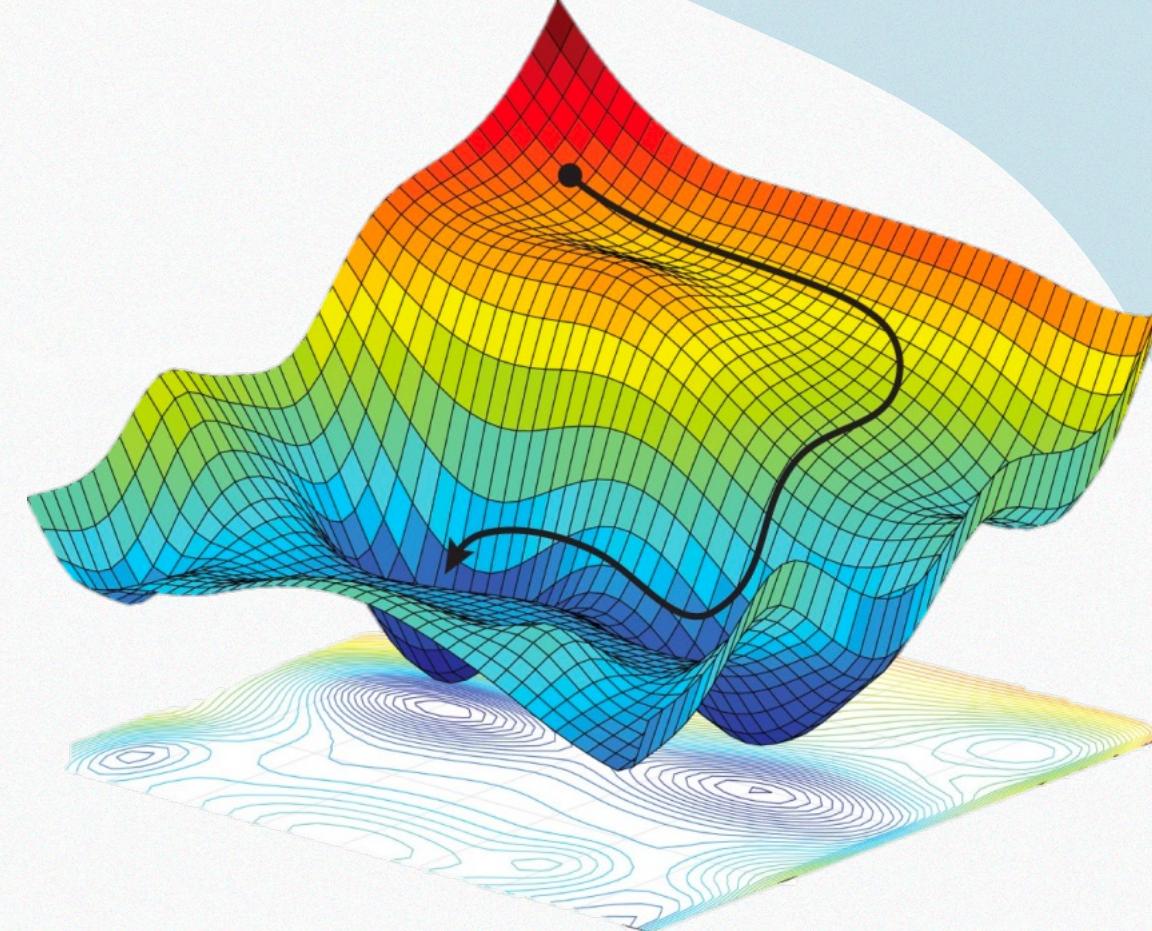
```
22     clip_coef_clamped = torch.clamp(clip_coef, max=1.0)
```

```
23     for g in grads:
24         g.detach_.mul_(clip_coef_clamped.to(g.device))
25     return total_norm
26
27
```

Reactor 五重境界：笨鸟先飞之Init

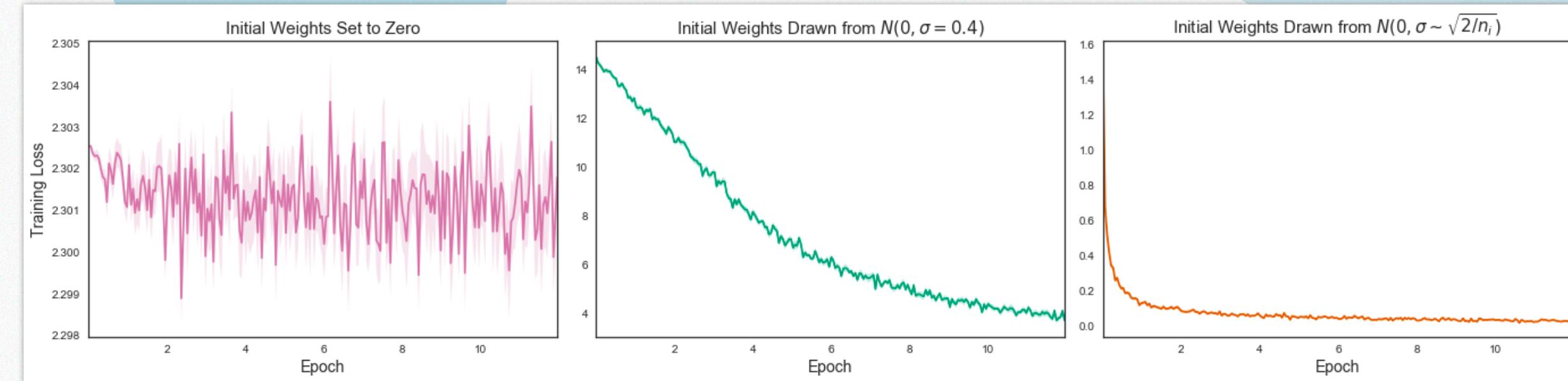
神经网络初始化设计的五大因素

- 输入层数量
- 输出层数量
- 非线性函数（激活函数）的类型
- 网络结构的类型（Linear/Conv/MHA）
- 网络的设计使用目的

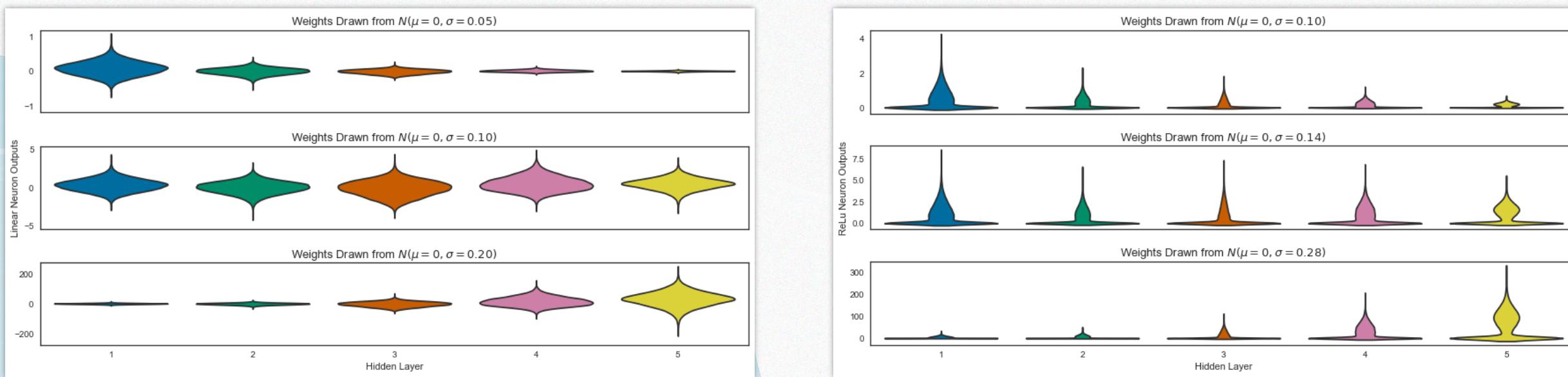


Rea⁵五重境界：笨鸟先飞之Init

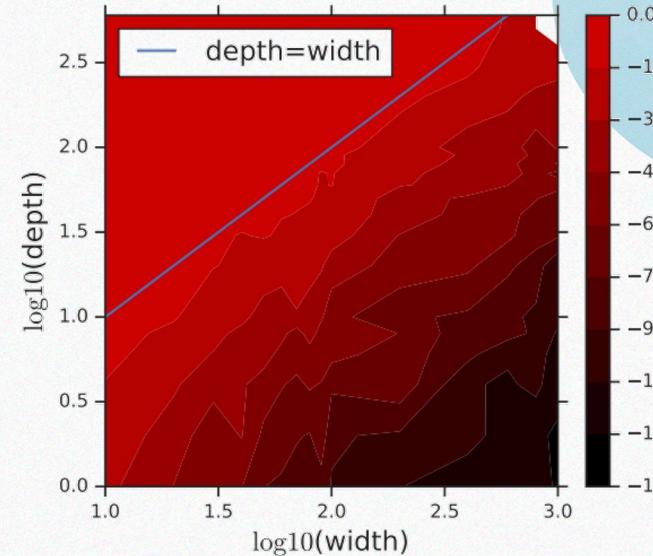
- 为什么不使用全零或固定方差的初始化



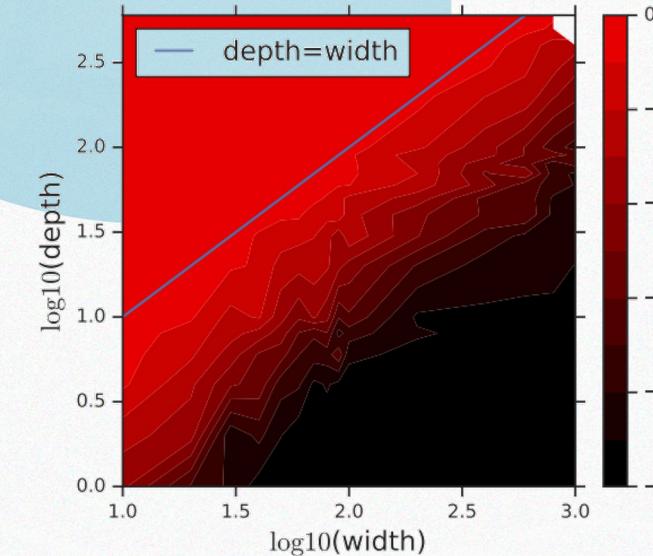
- 为什么要初始化考虑非线性函数的类型



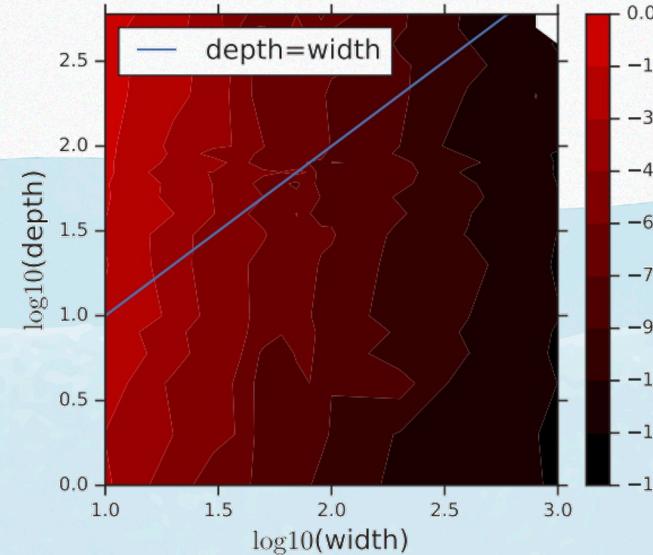
Reinforcement Learning 五重境界：笨鸟先飞之Init



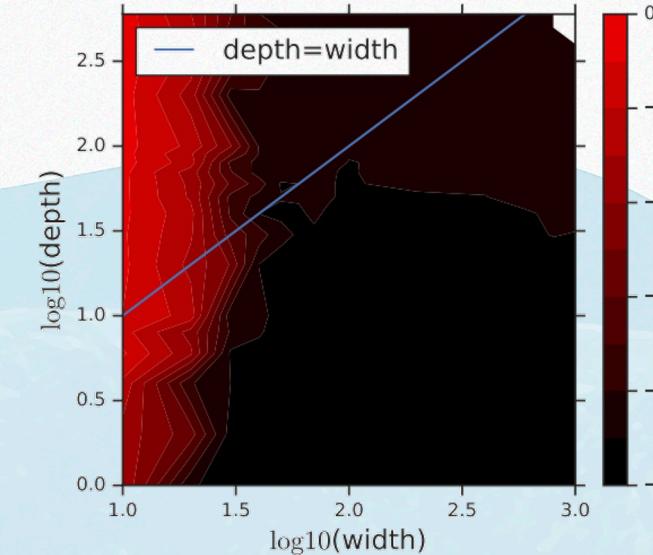
(a) Gaussian, steps=1258



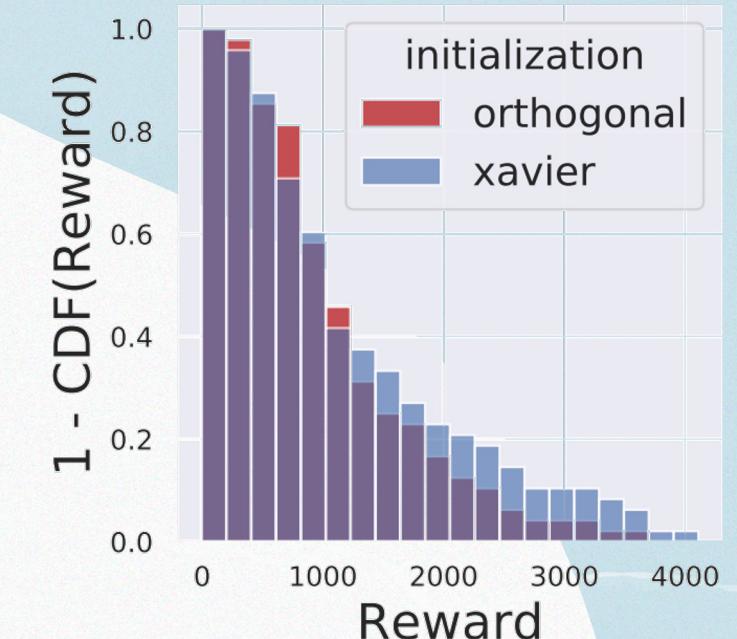
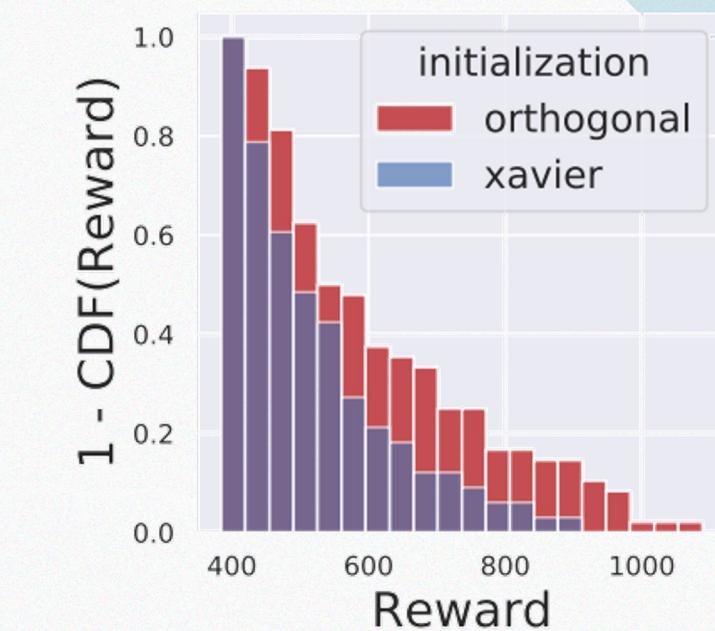
(b) Gaussian, steps=10000



(c) Orthogonal, steps=1258



(d) Orthogonal, steps=10000

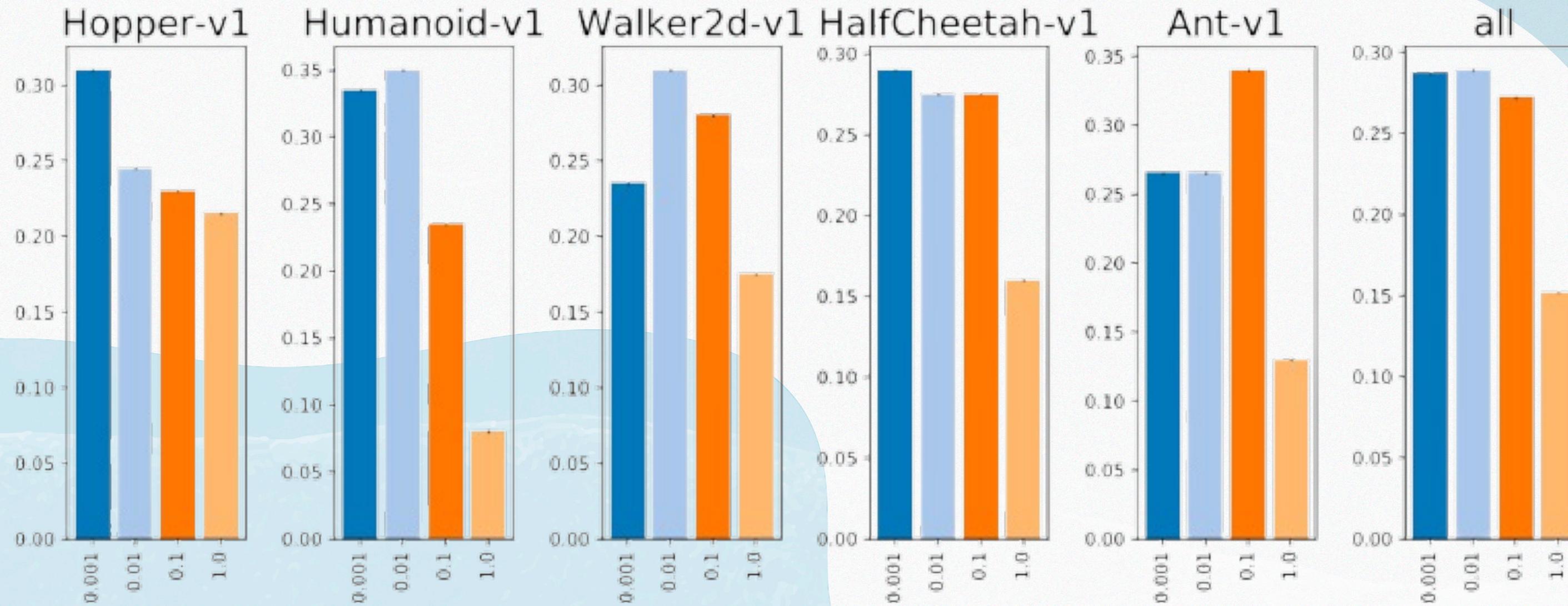


```

# Compute the qr factorization
q, r = torch.linalg.qr(flattened)
# Make Q uniform according to
# https://arxiv.org/pdf/math-ph/0609050.pdf
d = torch.diag(r, 0)
ph = d.sign()
q *= ph
# Initialize tensor
with torch.no_grad():
    tensor.view_as(q).copy_(q)
    tensor.mul_(gain)
  
```

Reinforcement Learning 五重境界：笨鸟先飞之Init

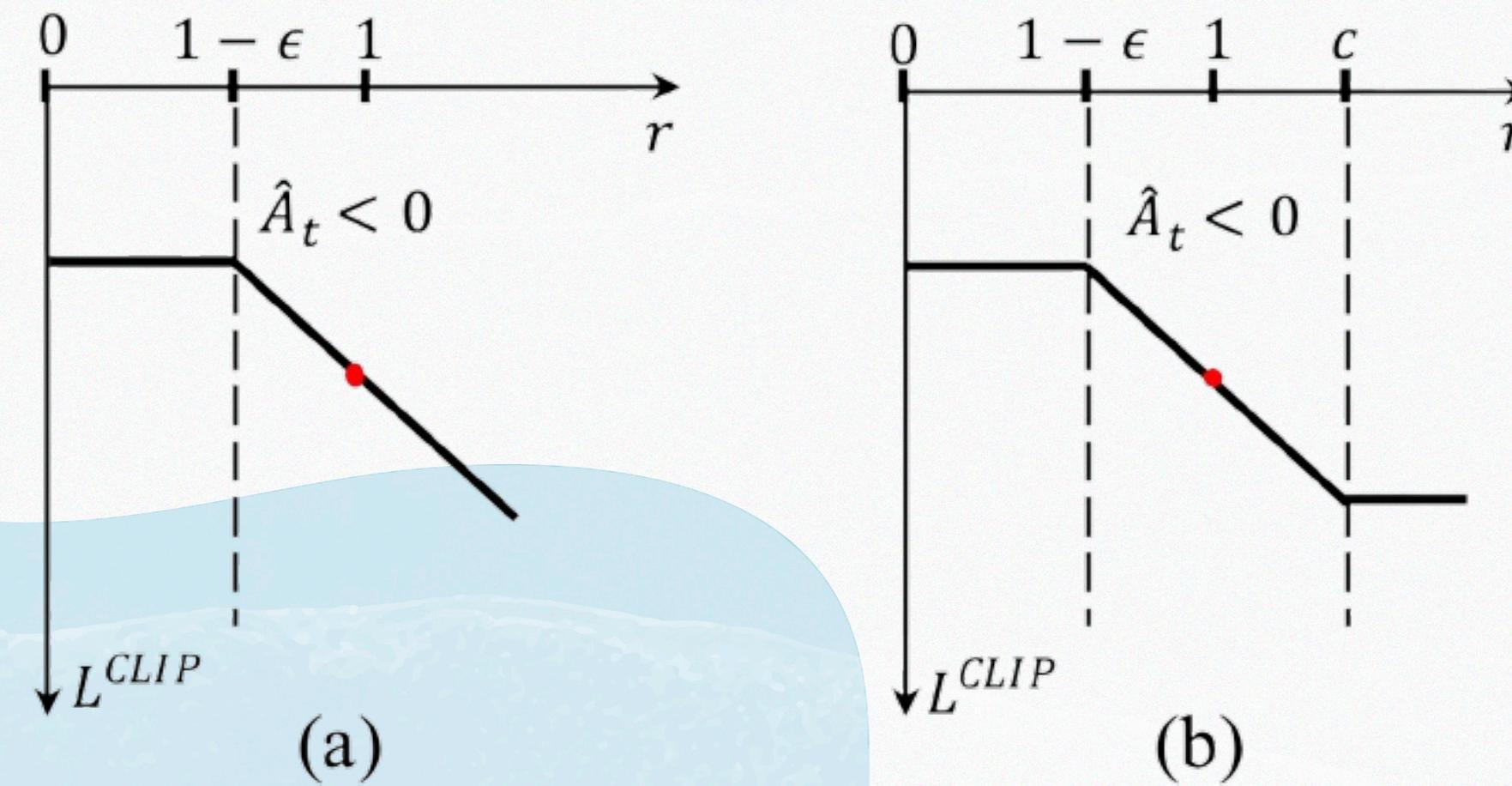
策略输出层的特殊初始化，保证训练初期的均匀探索



六重境界：亡羊补牢之Extra Clip

问题：当 $A < 0$ 时，如果新旧策略差异非常大（例如分布式训练情形），IS会异常偏大

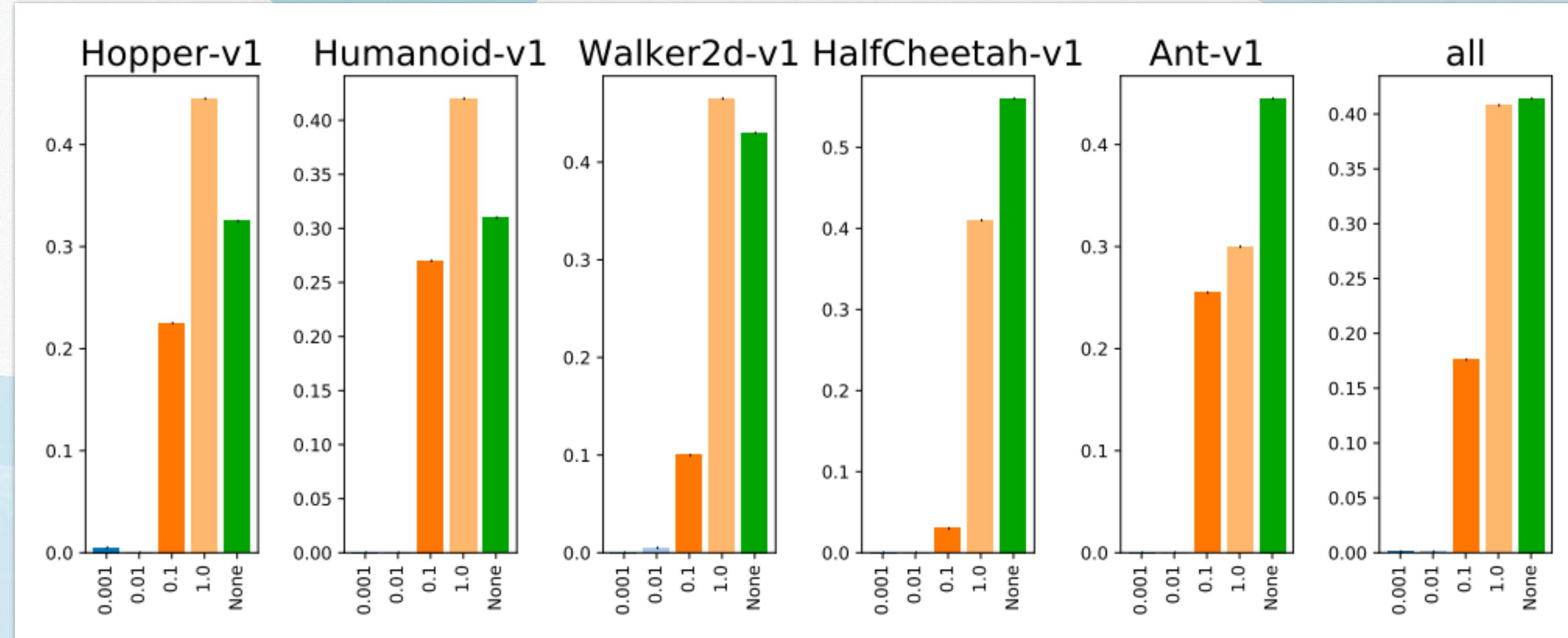
解法：对 $A < 0$ 的情况额外增加一个 clip 操作，即 dual clip 形式



$$L_\theta = \frac{1}{B \cdot T} \sum_{\tau \in b} \sum_{t=0}^{T-1} \max \left(\min(r(\theta) \hat{A}^{\theta_k}, \text{clip}(r(\theta), 1 - \epsilon, 1 + \epsilon) \hat{A}^{\theta_k}), c \hat{A}^{\theta_k} \right)$$

六重境界：亡羊补牢之Extra Clip

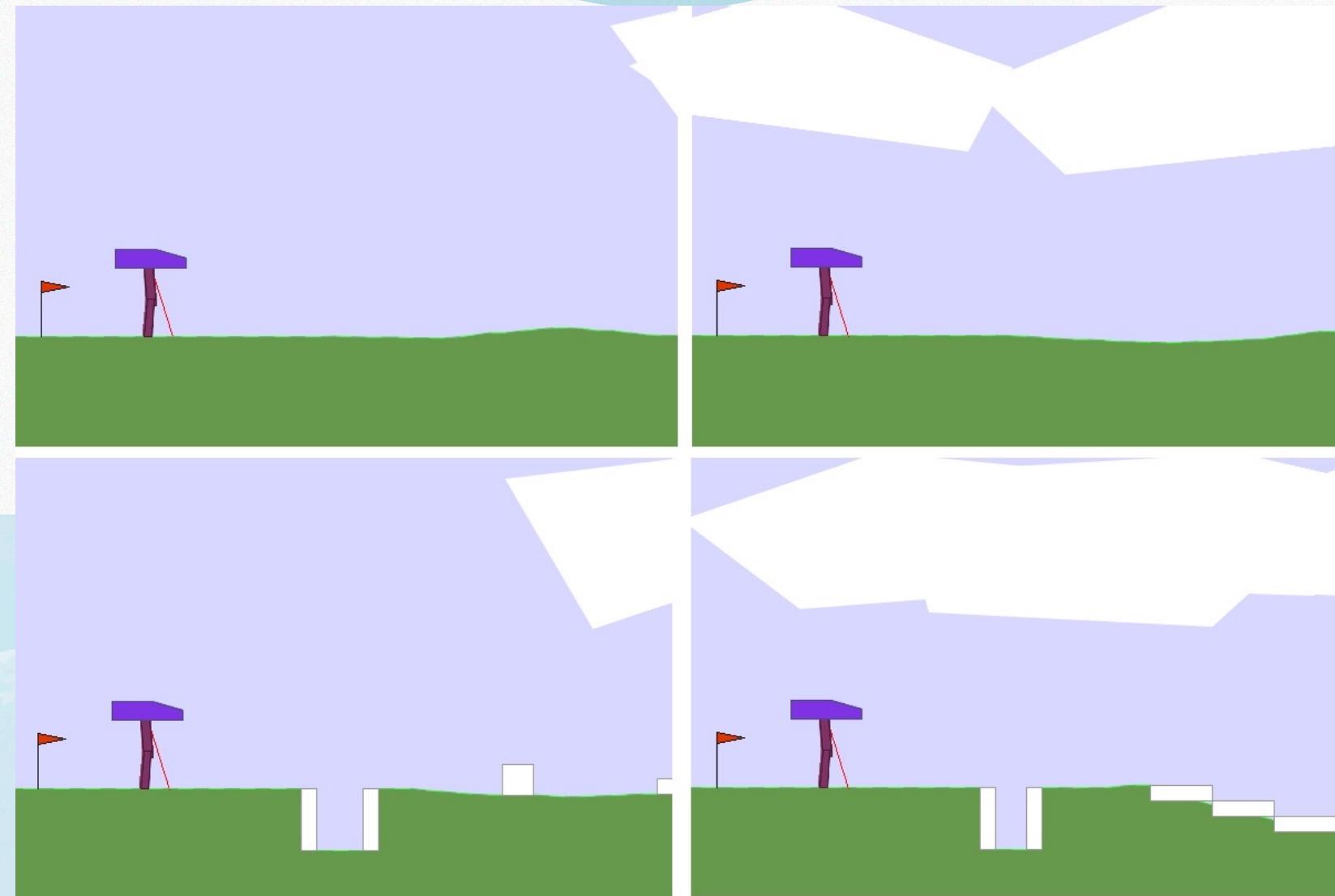
做法：对 Value Network 的优化也添加 clip 操作



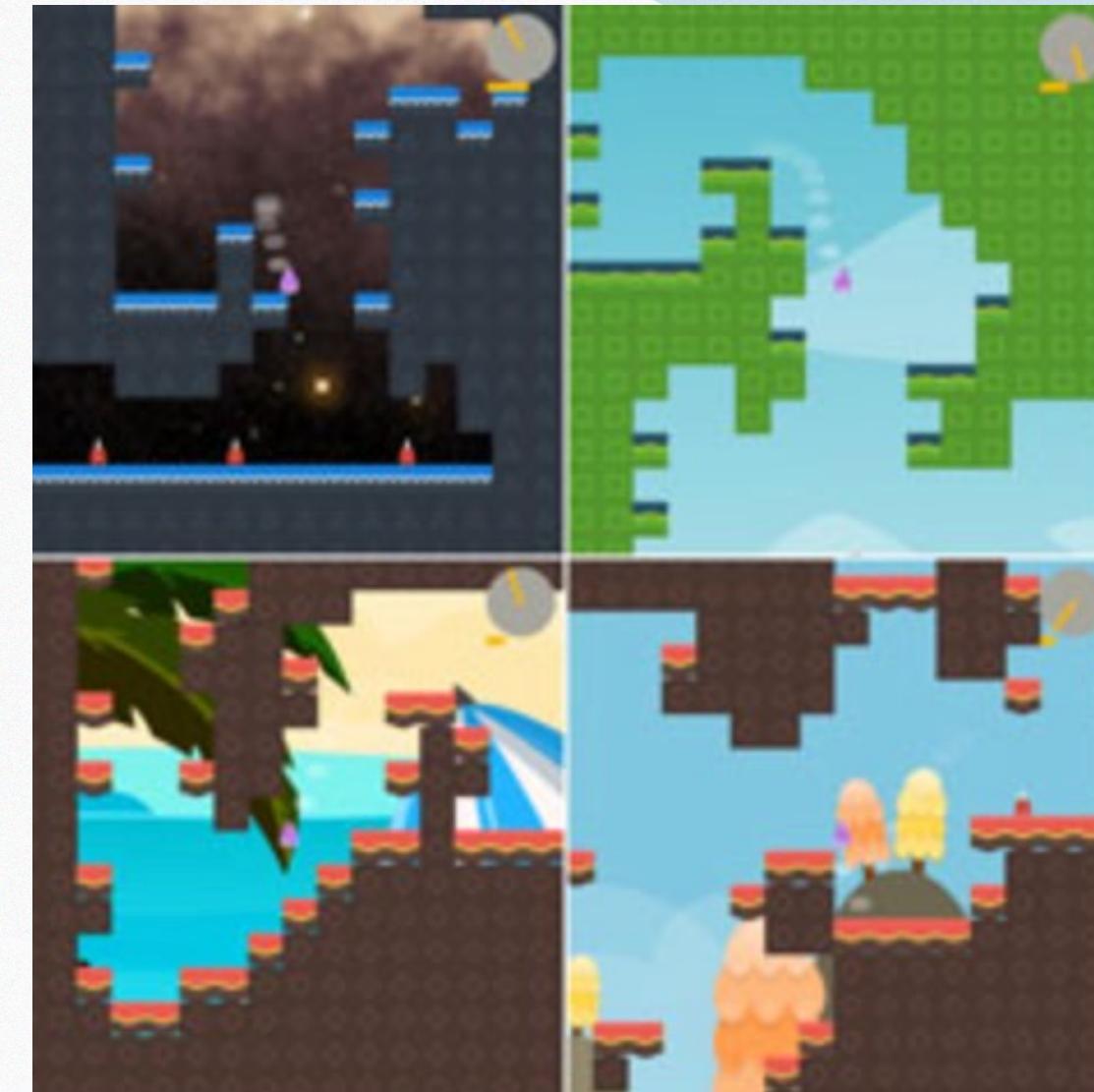
RealEnv 七重境界：合理利用随机

环境的随机种子：影响环境的状态空间（布局，内容）和状态转移（物理仿真参数）

Bipedalwalker



Progen: procedurally-generated



Realm 七重境界：合理利用随机

工具库导致的随机

收集数据时随机选择动作

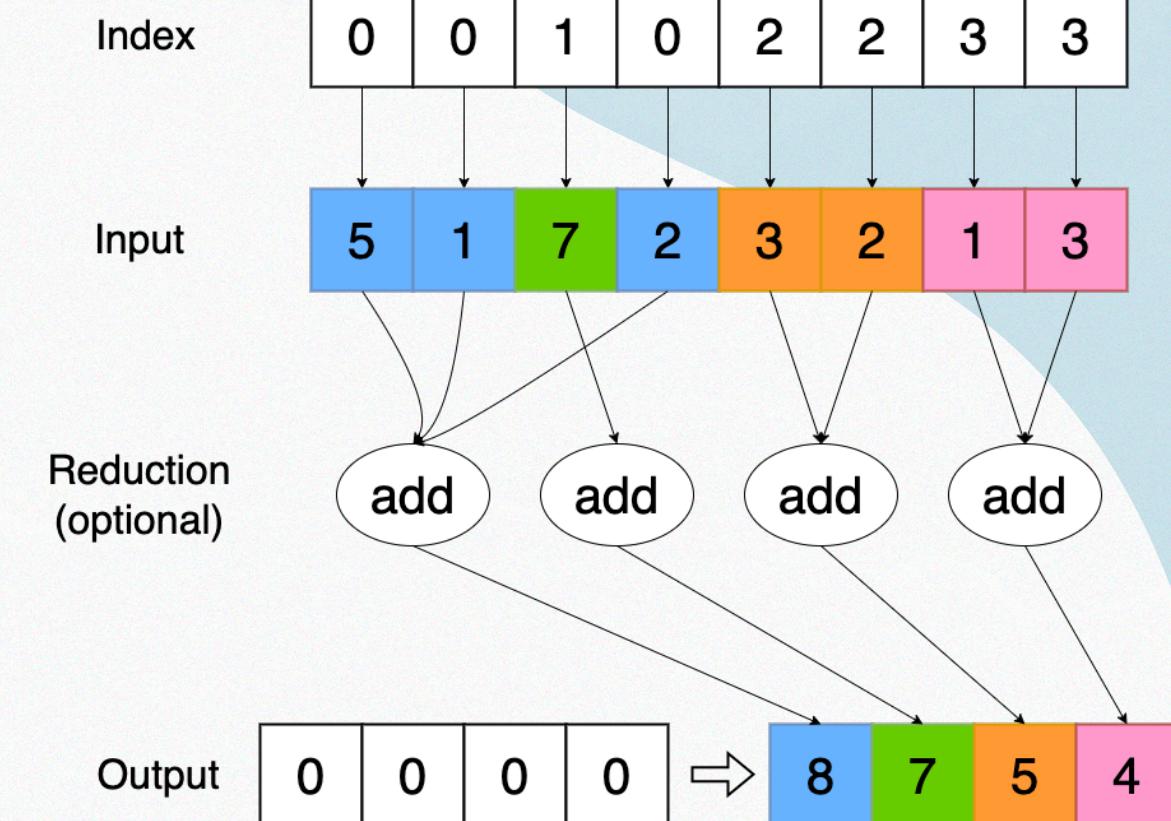
训练前随机从 buffer 中采样数据

神经网络参数的初始化

数据处理/变换带来的随机

不同计算设备引发的随机

异步/并行计算带来的随机



`random.seed(seed)`

`np.random.seed(seed)`

`torch.manual_seed(seed)`

`torch.cuda.manual_seed(seed)`

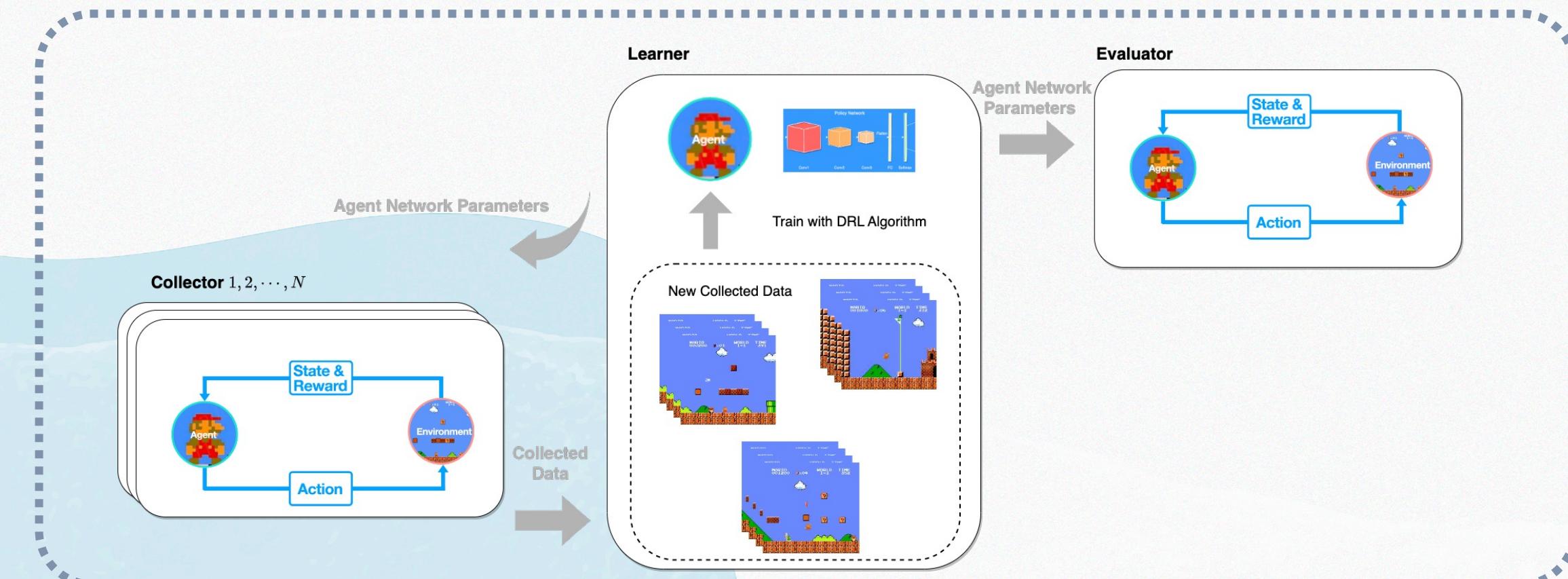
RealIn 七重境界：合理利用随机

收集数据 (collect) :

- 目的是收集训练数据。欢迎各类合理的随机性，增加数据的多样性，从而提升策略通用性

评估测试 (eval) :

- 目的是评估智能体的性能。单次测试要足够覆盖环境的随机性，但多次测试之间要可比



RealIn 七重境界：合理利用随机

		分散程度 (Dispersion)	风险度量 (Risk)
训练	多次训练时刻 (在一次训练之内)	滑动窗口平均后，算法表现的 IQR	短期指标：算法表现的一阶差分的 CVaR 长期指标：算法表现最大回撤的 CVaR
	多次训练	低通滤波后的算法表现的 IQR	算法表现的 CVaR
测试	固定策略的多次试验	对于某个固定策略的多次试验中算法表现的 IQR	对于某个固定策略的多次试验中算法表现的 CVaR

IQR:

inter-quartile range

即 75 和 25 分位点之差

CVaR:

conditional value at risk

即 最坏情况下的期望回撤

下节预告

(八) 突破智能体终极界限

- PPO + 分布式训练
- PPO + 大语言模型 (LLM)
- 神秘终极挑战任务