

PPo^x Family

第三讲：表征多模态观察空间

主办



承办



协办



支持

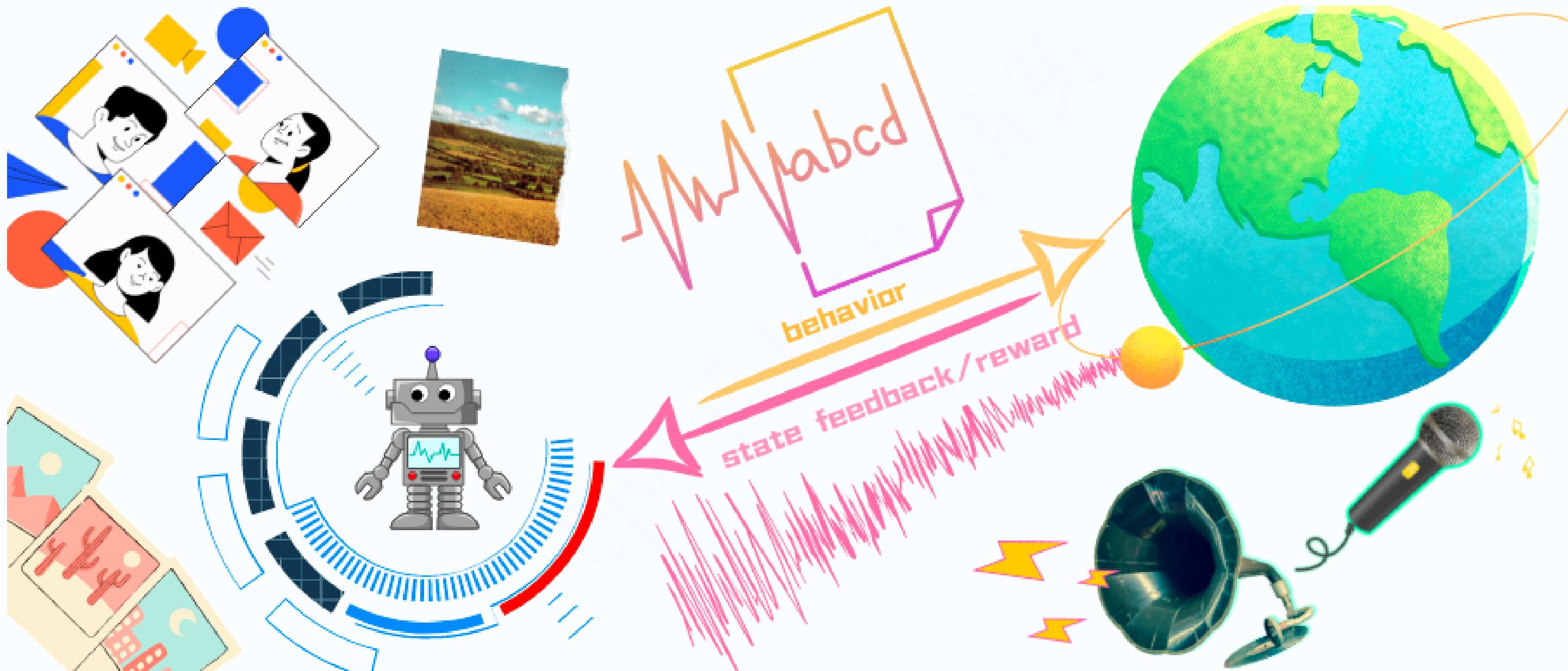


观察空间三部曲

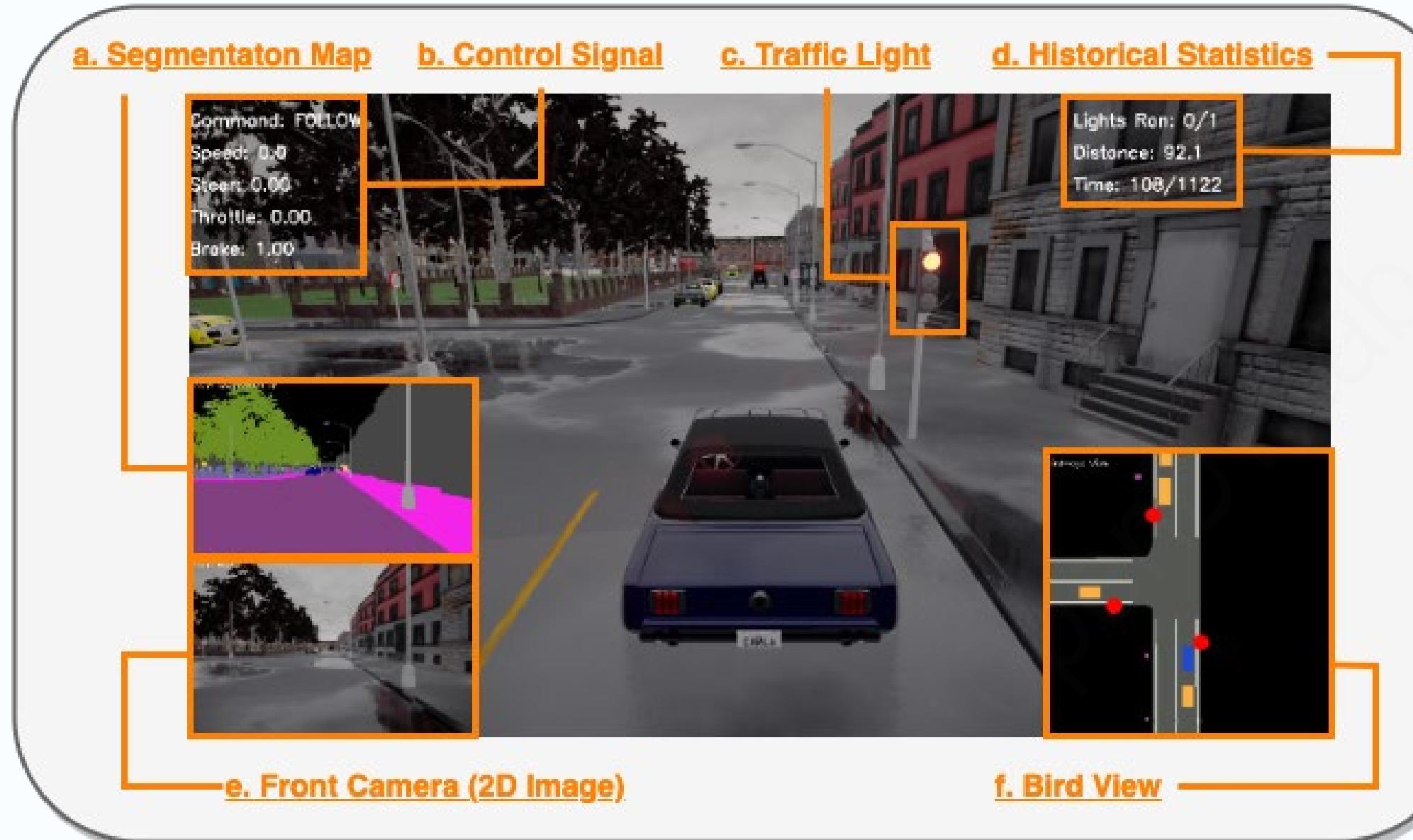


Obs Space Overview

观察空间概述



常见的观察空间类型

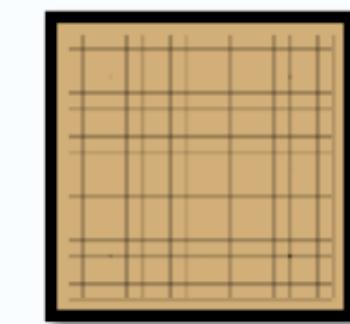
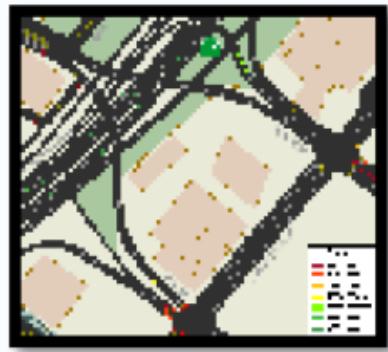


特点：

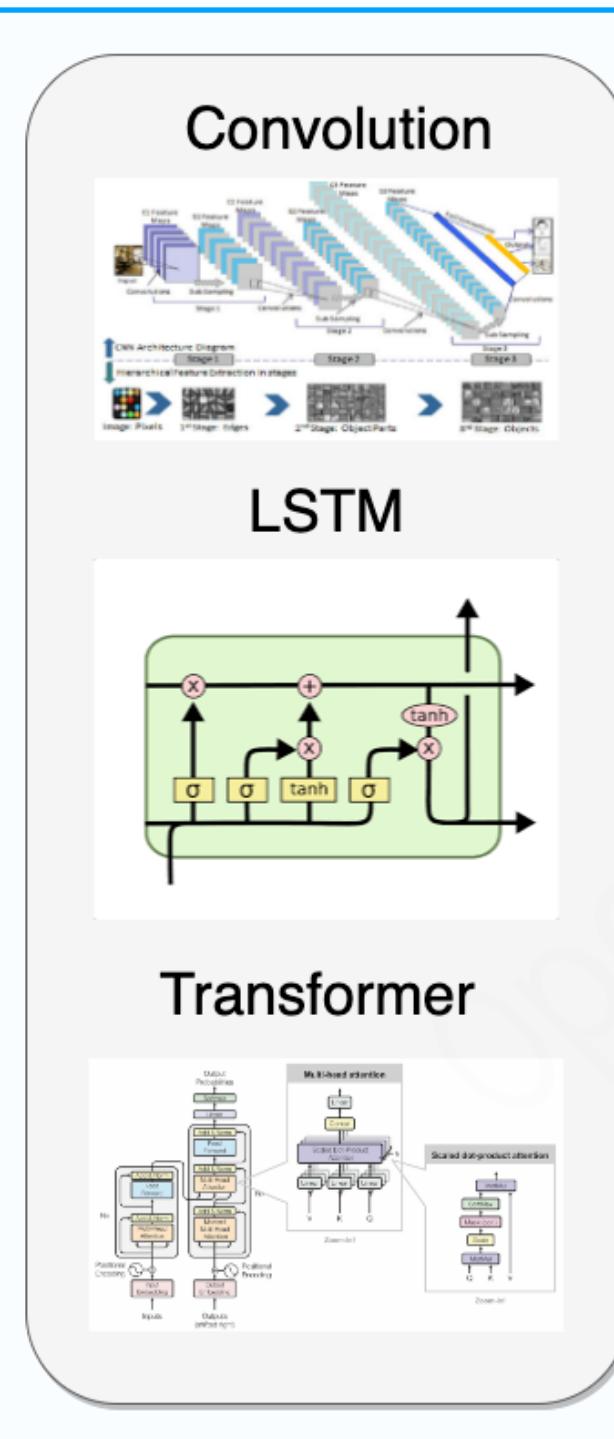
- 特征预处理可以事半功倍
- 高维的观察信息更需要专门设计的神经网络架构
- 复杂的决策问题往往包含多种模态的观察信息

Encode 编码器

多模态观察空间



表征编码器 (Encoder)



决策预测器 (Head)

离散动作预测

连续动作回归

连续动作重参数化

多维离散

层次结构化动作

自回归预测

复杂动作空间

Turn Right



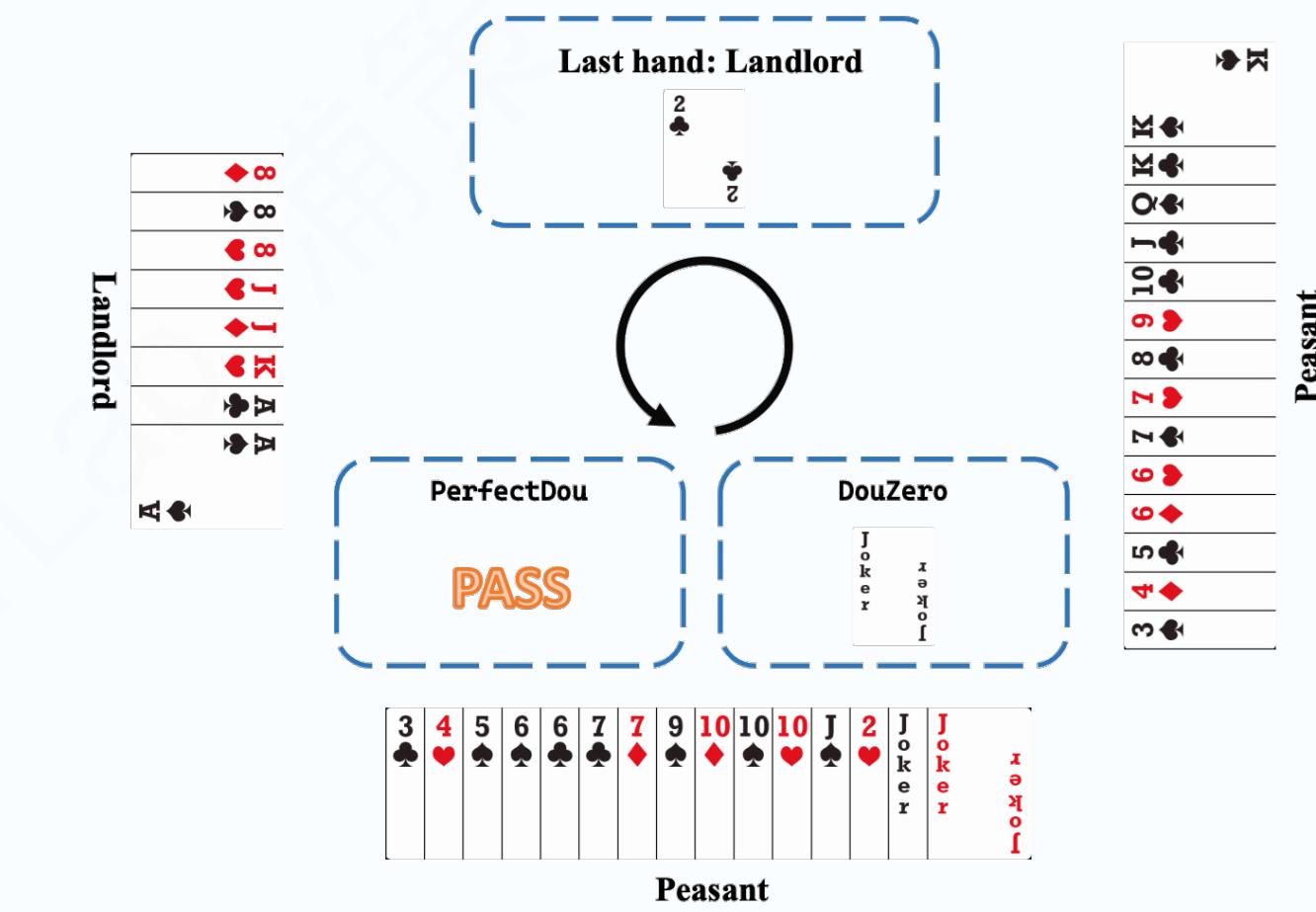
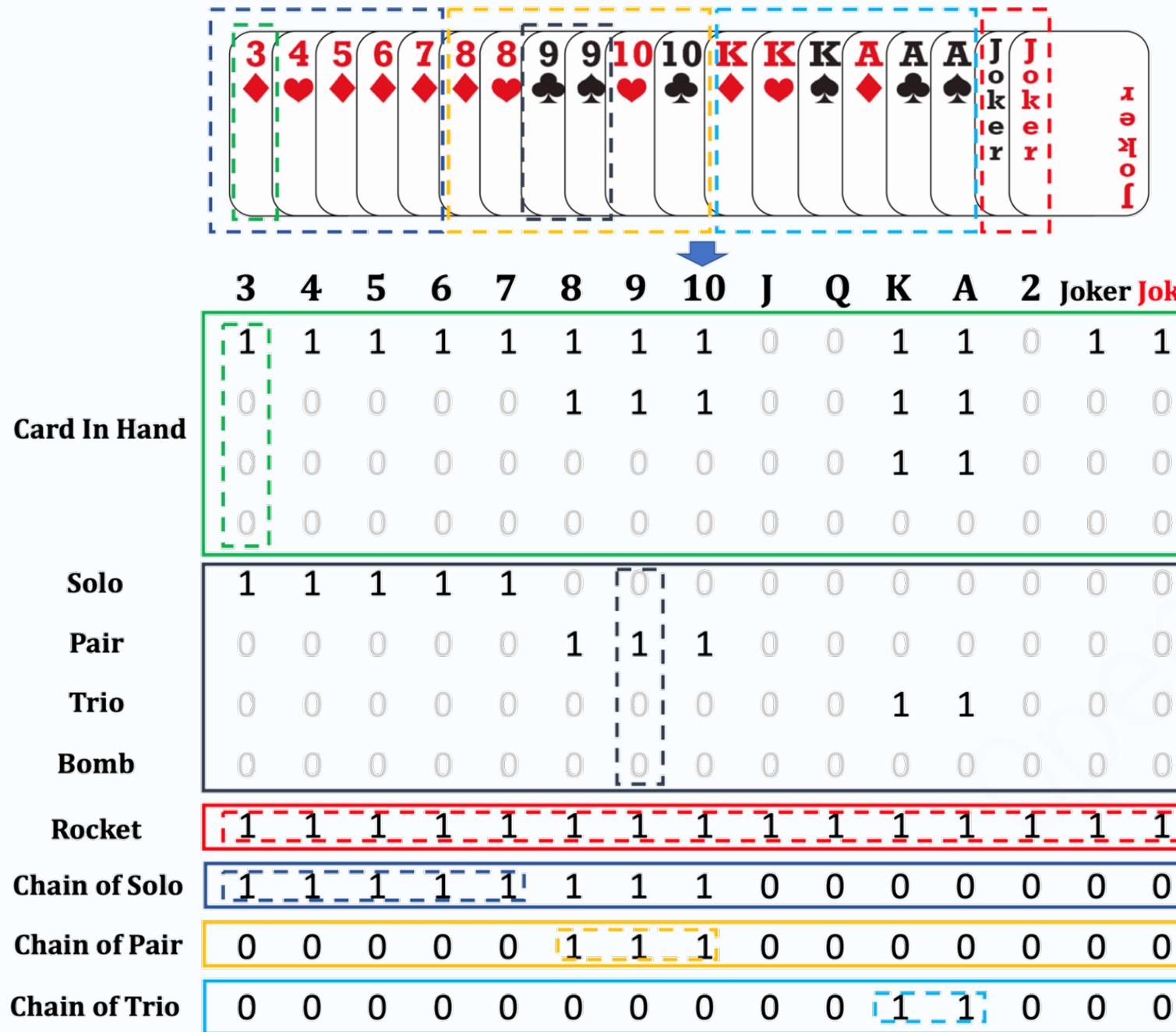
Shoot



Move

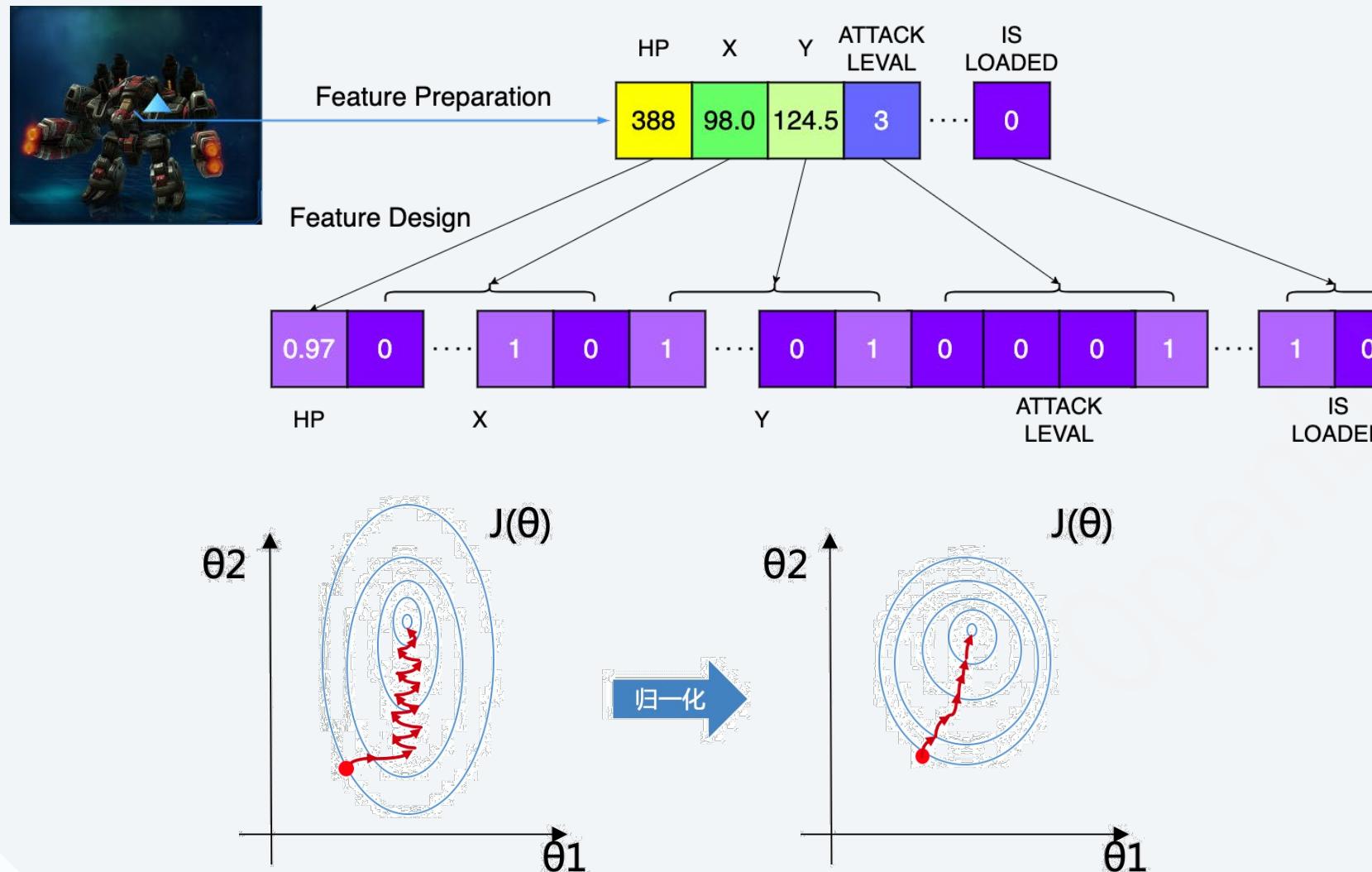


向量观察空间



理论：PPO + 向量观察 ● 预处理

统一量纲



编码设计

One-hot encoding

The diagram shows a list of words and their corresponding one-hot encoded vectors. The words are: Color, can, I, eat, the, pizza. The encoding is shown in a 5x6 grid where each row represents a word and each column represents a dimension d0 to d4. A red arrow labeled 'One-hot encoding' points from the word list to the grid.

	d0	d1	d2	d3	d4
Color	1	0	0	0	0
can	0	0	1	0	0
I	0	1	0	0	0
eat	0	0	0	0	1
the	0	0	0	0	0
pizza	0	0	0	1	0

Word Embedding

The diagram shows a list of words and their corresponding word embeddings. The words are: Color, Man, Woman, King, Queen, Orange, Apple, The embeddings are shown in a 7x5 grid where each row represents a word and each column represents a dimension d0 to d4. A red arrow labeled 'Word Embedding' points from the word list to the grid.

	d0	d1	d2	d3	d4
Color	-1.00	0.01	0.03	0.04	0
Man	1.00	0.02	0.02	0.01	0
Woman	-0.95	0.93	0.70	0.02	0
King	0.97	0.95	0.69	0.01	0
Queen	0.00	-0.01	0.03	0.95	0
Orange	0.01	0.00	-0.02	0.97	0
Apple

理论：PPO + 向量观察 ● 不变性

平移不变性



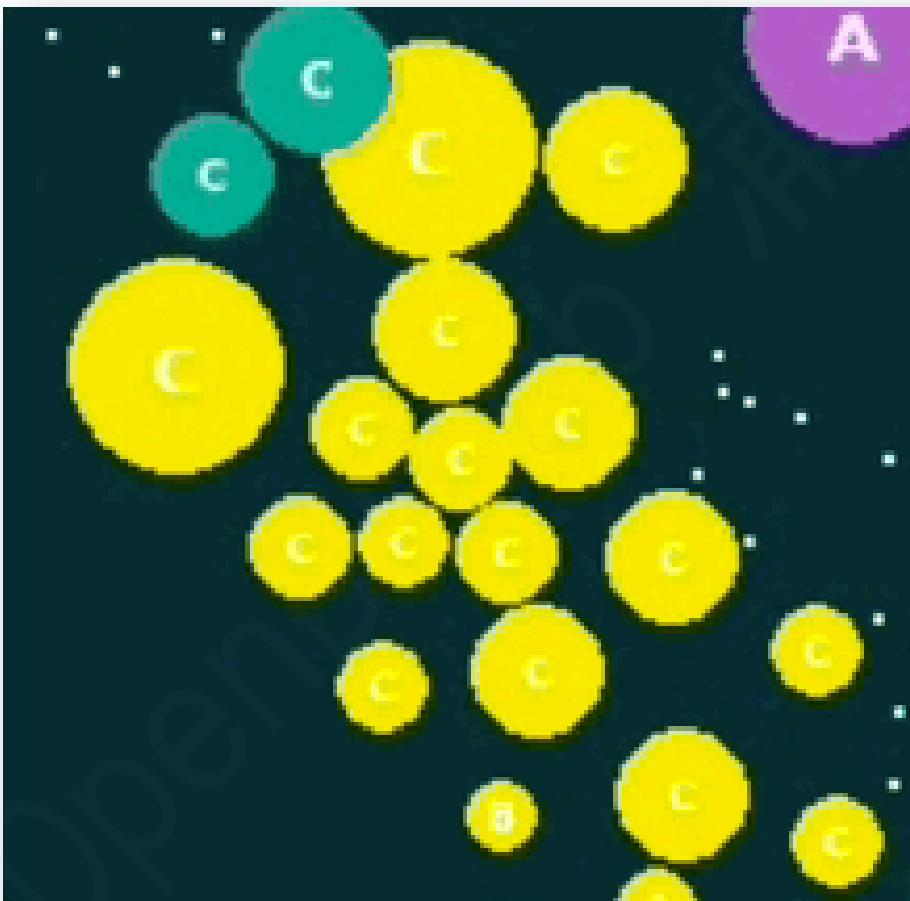
不变性 (invariance)

$$F(X) = F(T(X))$$

等变性 (equivariance) $T(F(X)) = F(T(X))$

F指特征编码或神经网络，T指变换，例如平移，旋转，改变顺序

旋转不变性

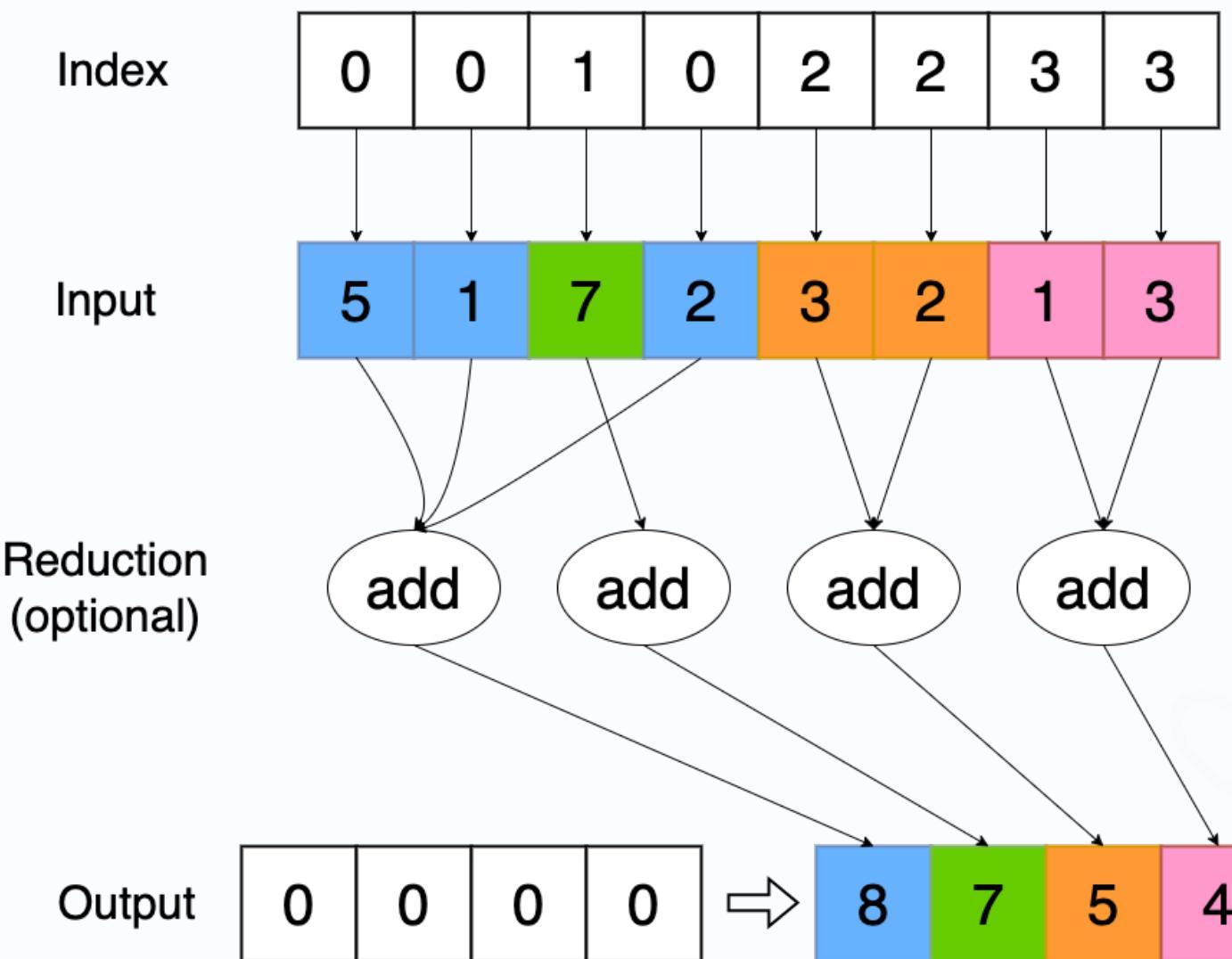


排列不变性

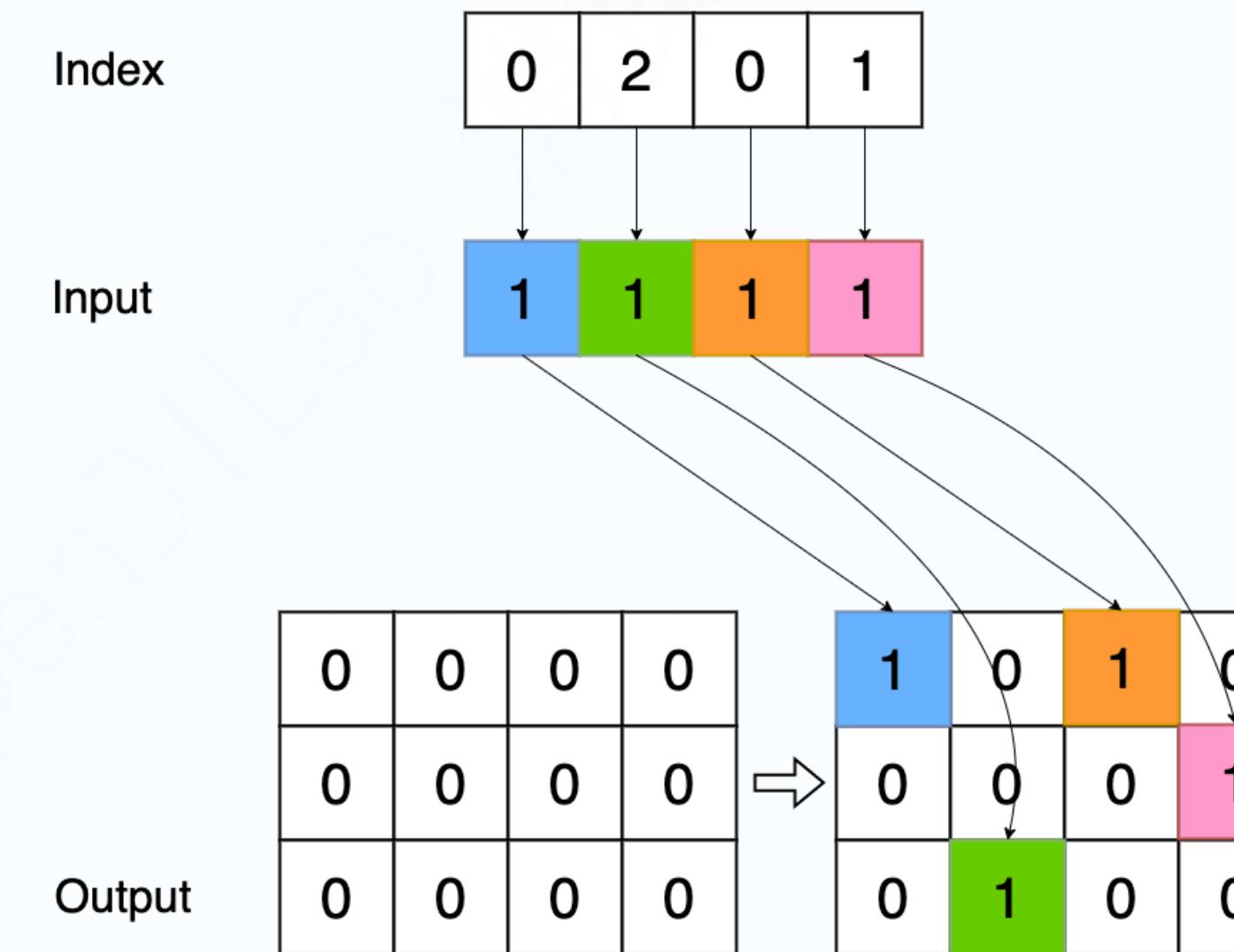


代码：特征工程中的张量操作技巧

$$\text{Scatter: } \text{out}_i = \text{out}_i + \sum_j \text{in}_j$$



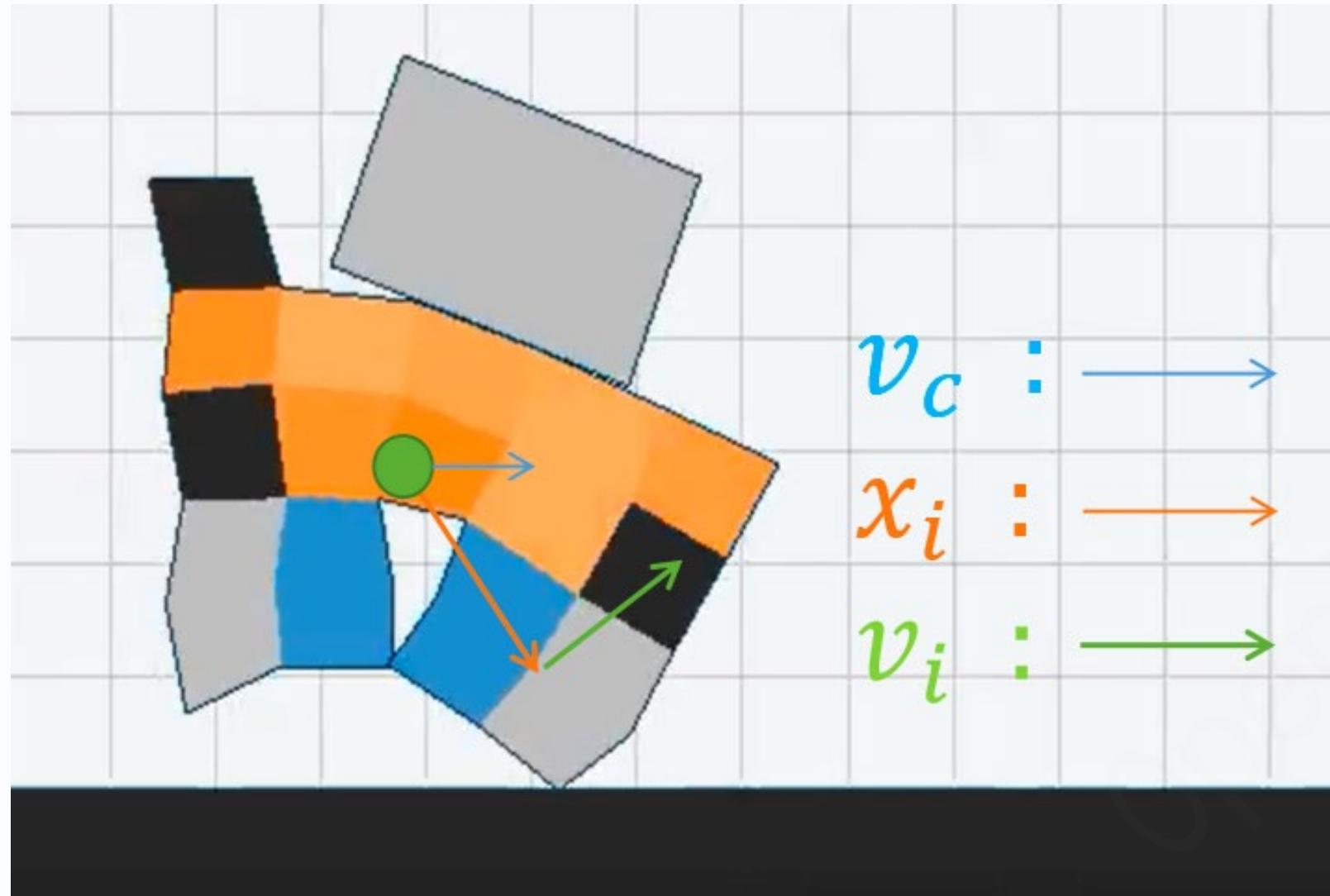
One-Hot Encoding By Scatter



Scatter: https://github.com/rusty1s/pytorch_scatter

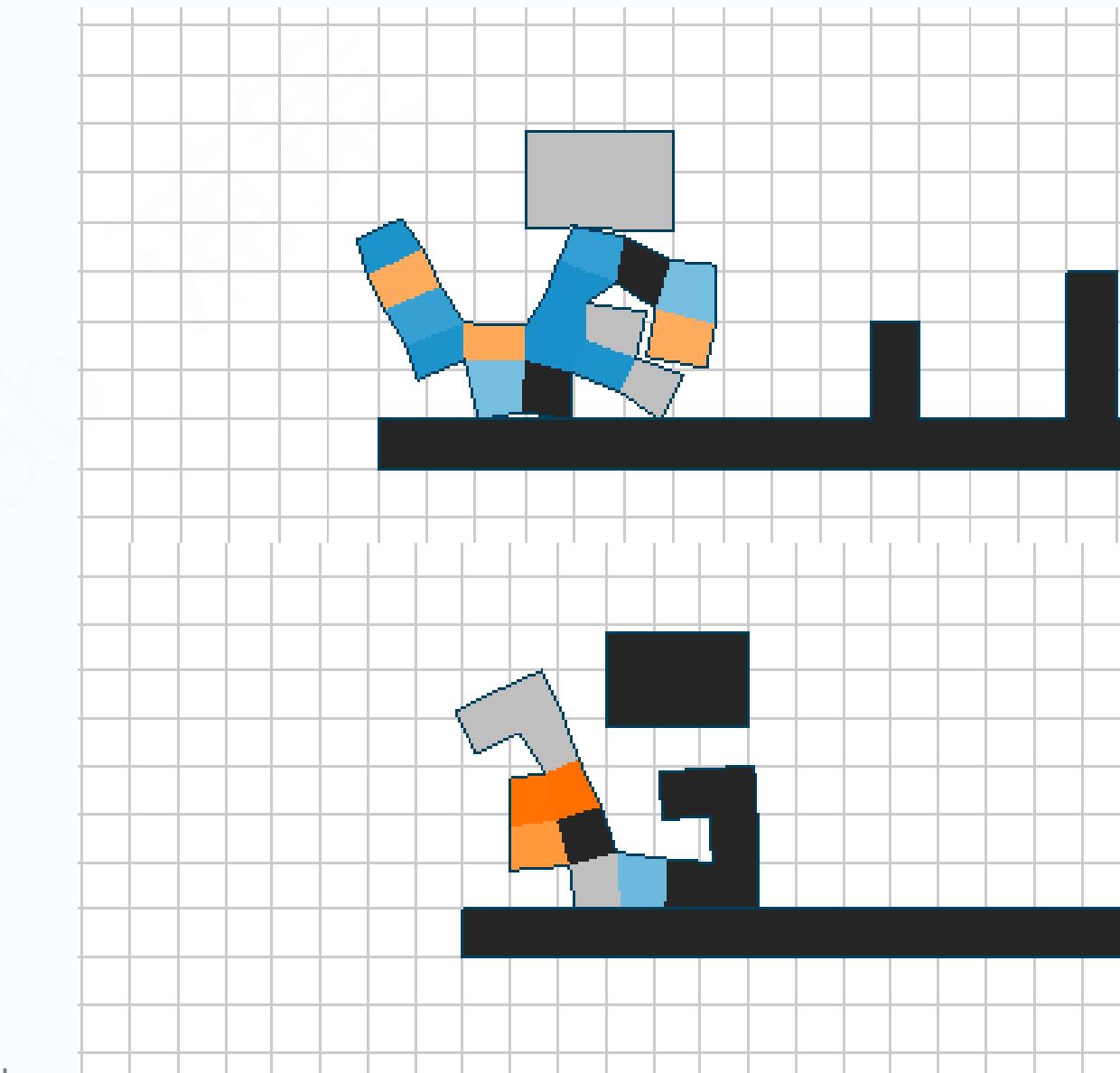
完整示例: <https://opendilab.github.io/PPOxFamily/>

实践：PPO+软体机器人



Evogym: <https://github.com/EvolutionGym/evogym>

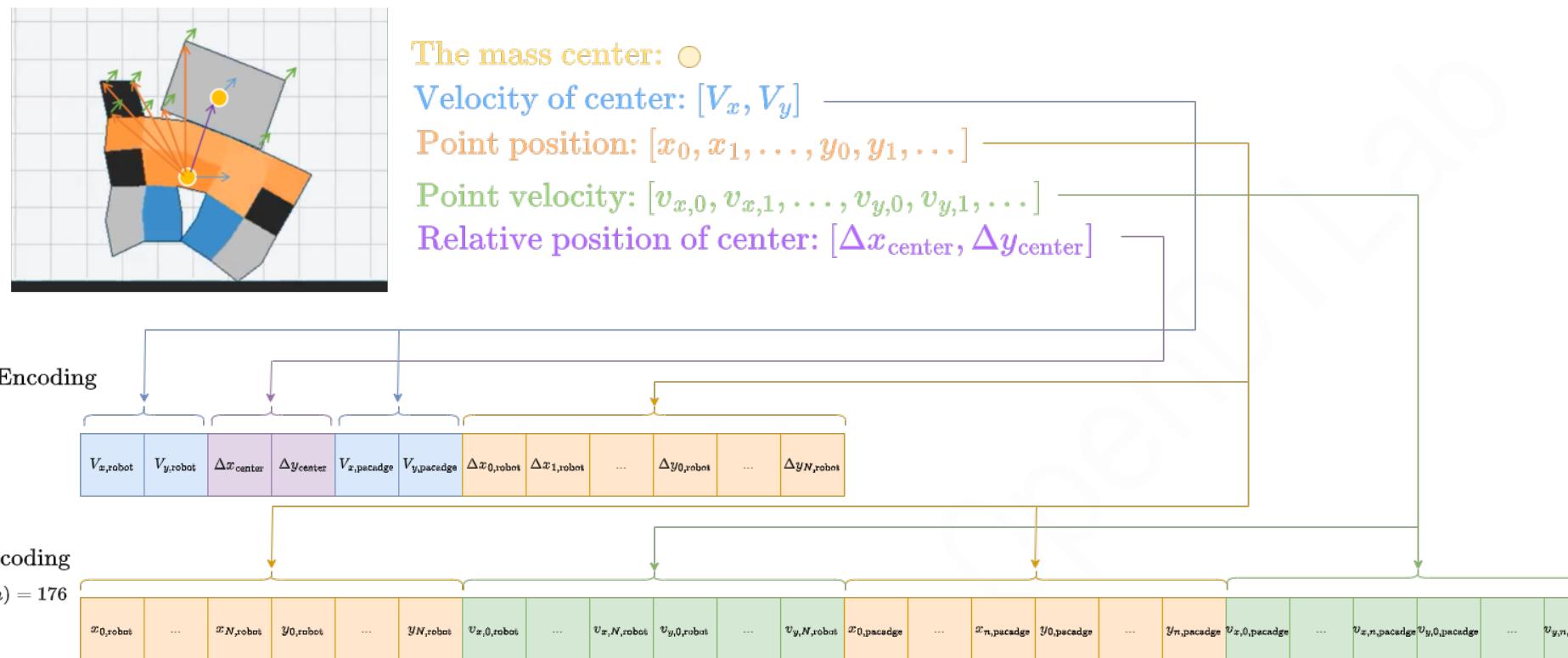
DI-engine + Evogym doc: https://di-engine-docs.readthedocs.io/en/latest/13_envs/evogym.html



实践：PPO+软体机器人

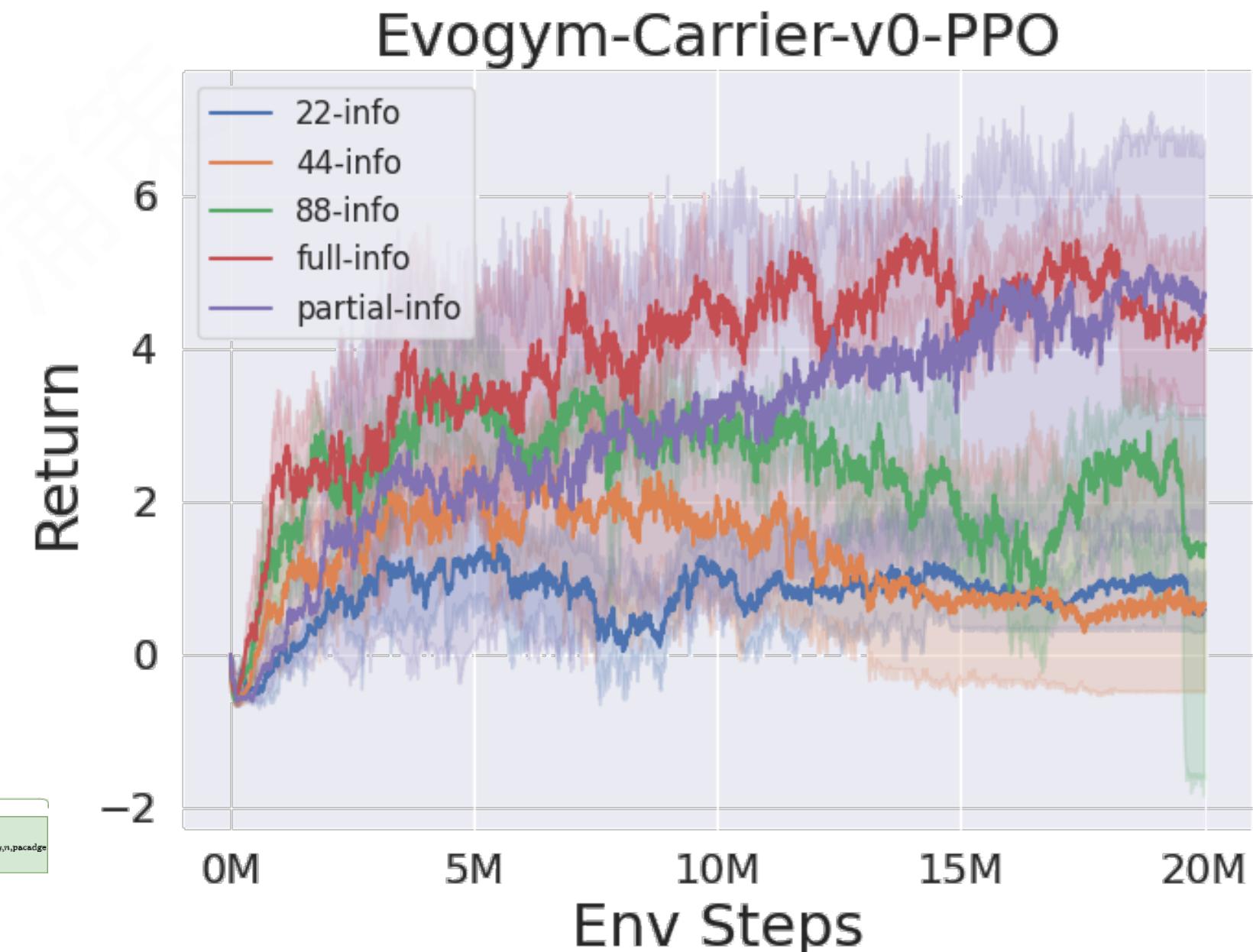
观察空间：

机器人全局和局部信息（例如位置，速度）的编码



实验细节：<https://github.com/opendilab/PPOxFamily/issues/8>

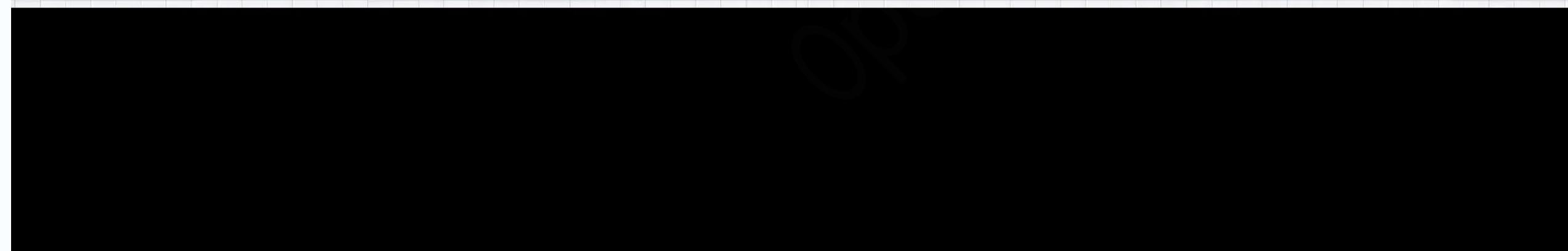
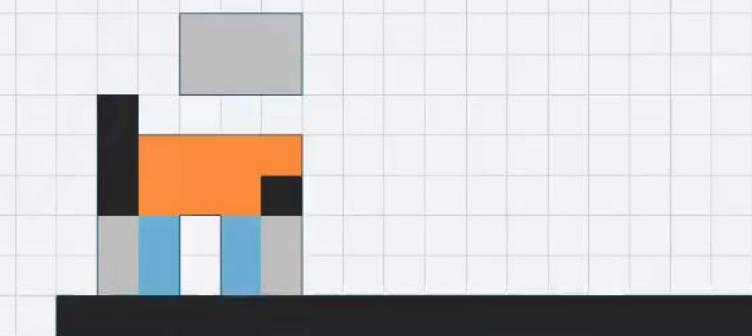
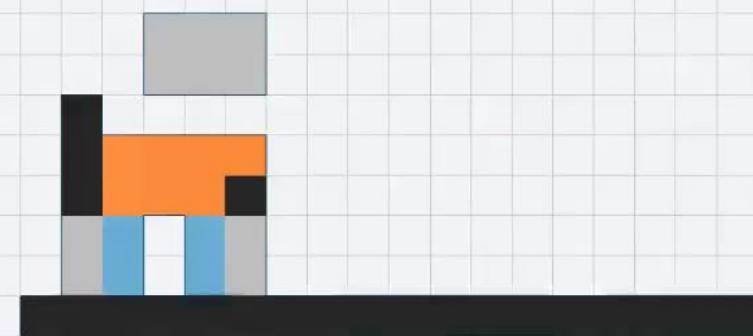
P.S. 88/44/22实验组指只选取软体机器人中部分数量小块的观察信息，数字代表数量



实践：PPO+ 软体机器人

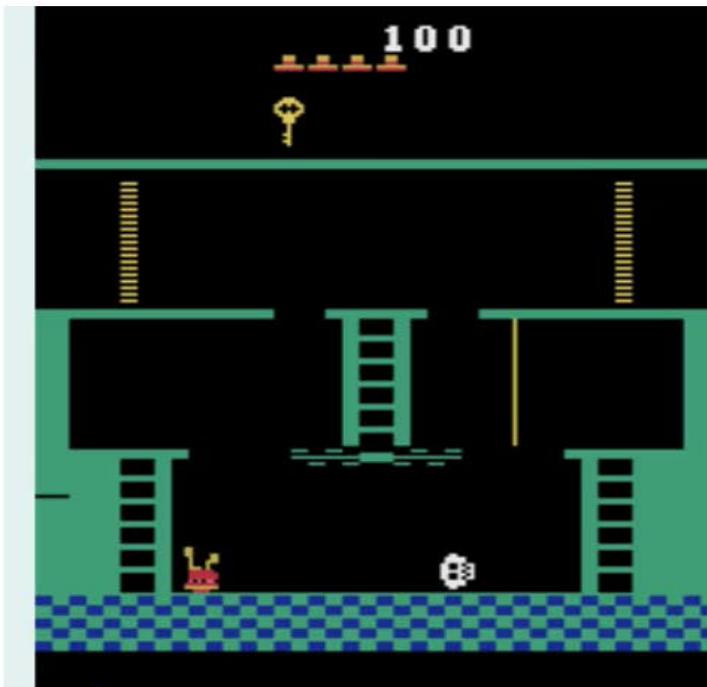
Agent during Early Training

Trained Agent



视频样例：<https://github.com/opendilab/PPOxFamily/issues/8>

图片观察空间



The agent always dies after getting the first key.



Step 83
Reward: 0.001

理论：PPO + 图片观察 ● 预处理

问题：图片数据维度很高但有很多冗余，因此信息密度很低，

Origin



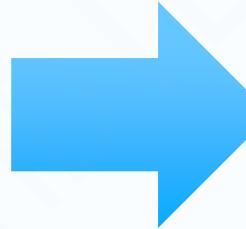
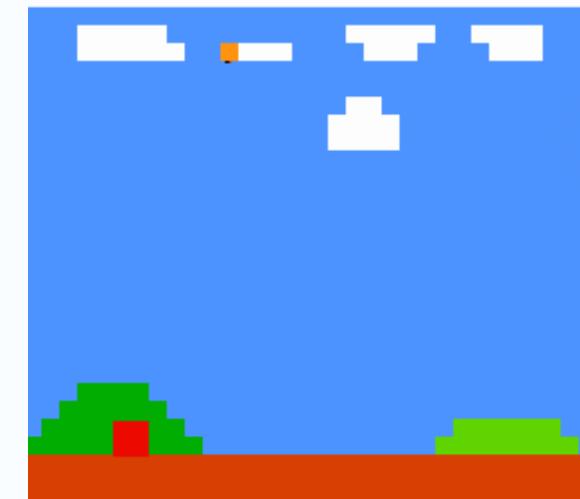
Pixel



Downsample

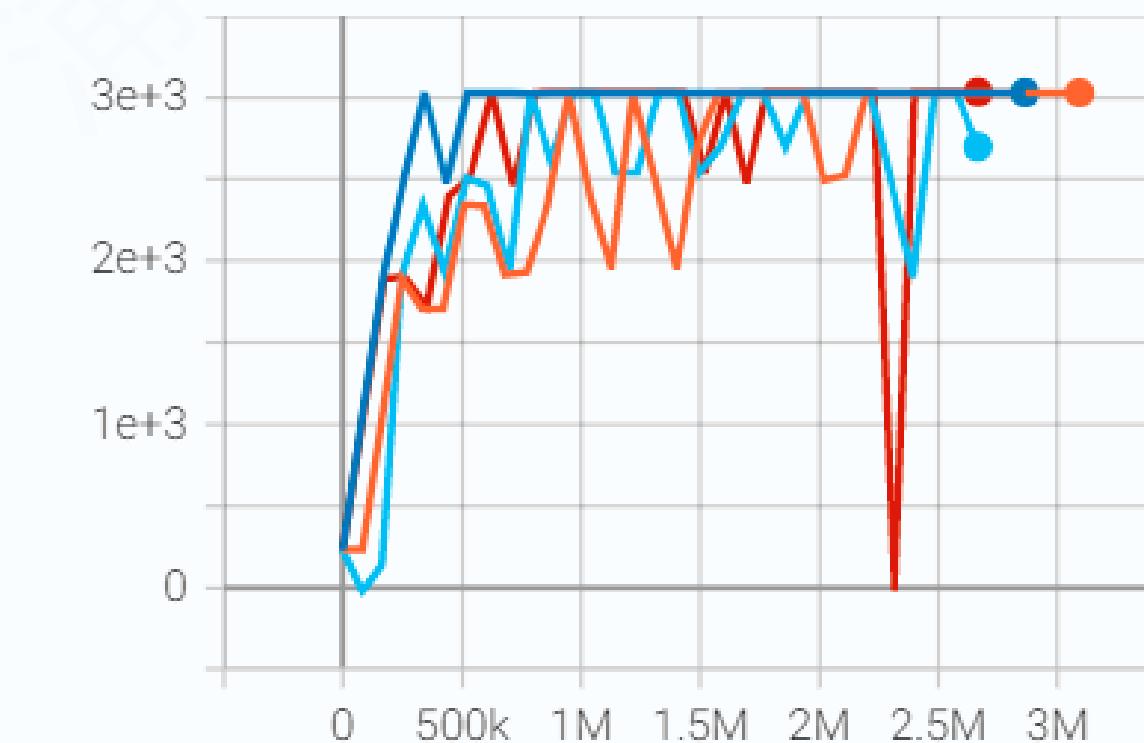


Rectangle



reward_mean

tag: evaluator_step/reward_mean



downsample > pixel >= rectangle >= origin



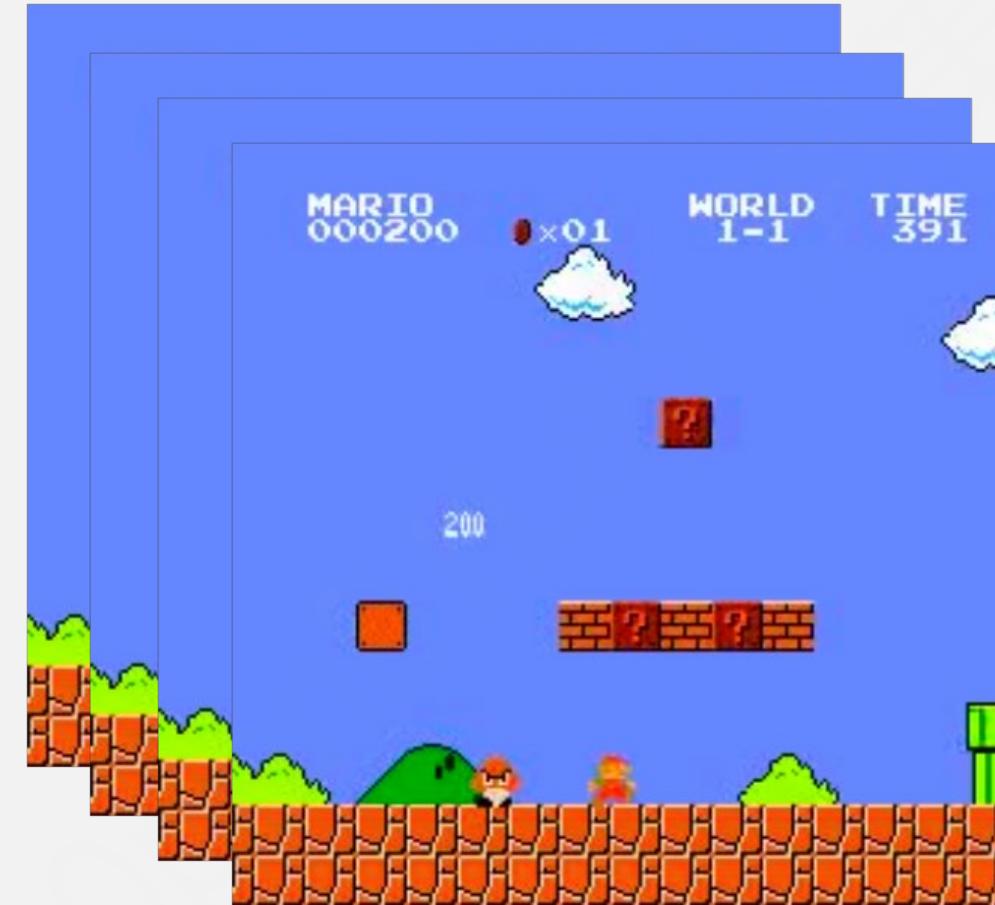
理论：PPO + 图片观察 ● 运动信息

单帧



State Shape: (1, 84, 84)

叠帧



State Shape: (4, 84, 84)

光流



State Shape: (3, 84, 84)

问题：单帧图片无法表示运动信息（例如方向，速度）

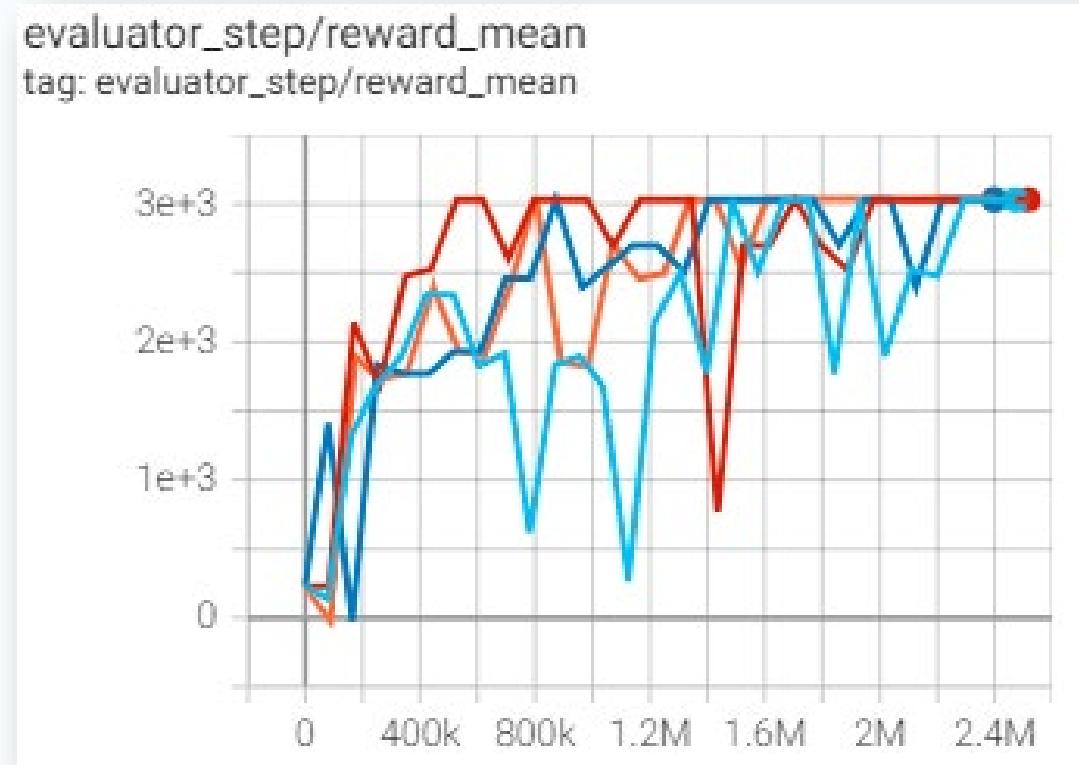
解决思路：堆叠连续4帧环境帧为一个动作帧 or 提取光流



理论：PPO + 图片观察 ● 运动信息

- origin
- downsample
- rectangle
- pixel

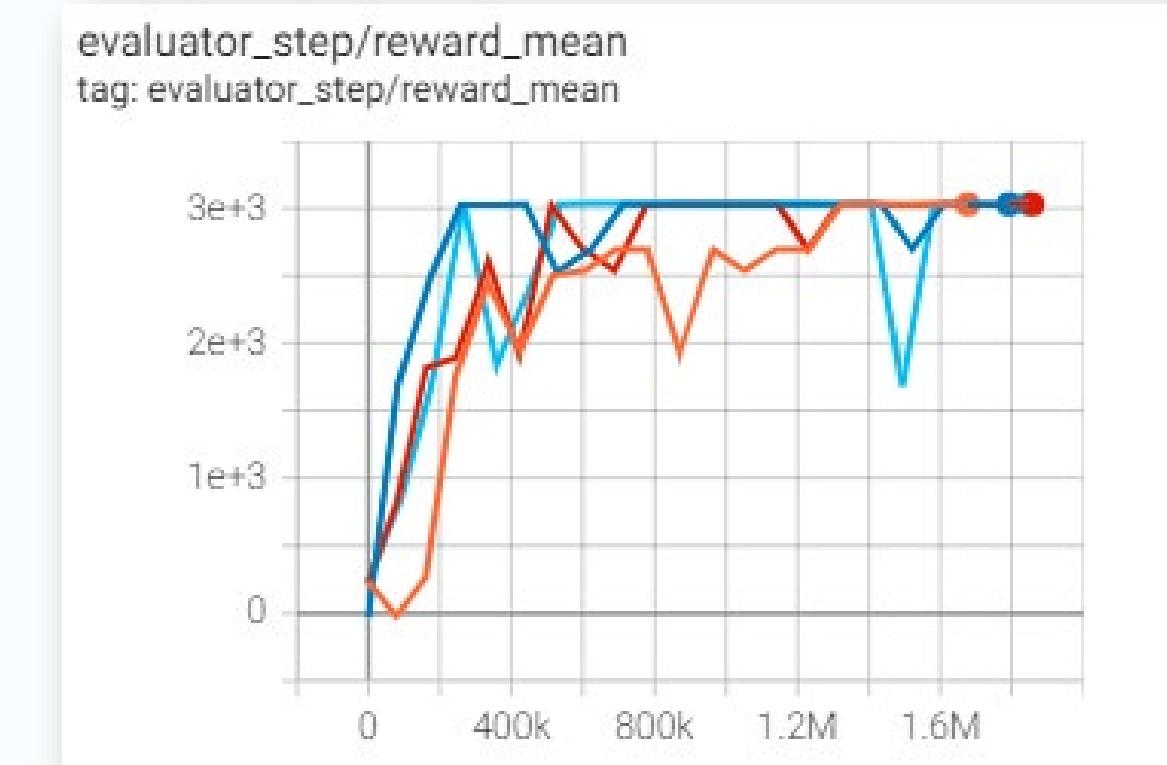
单帧



叠帧



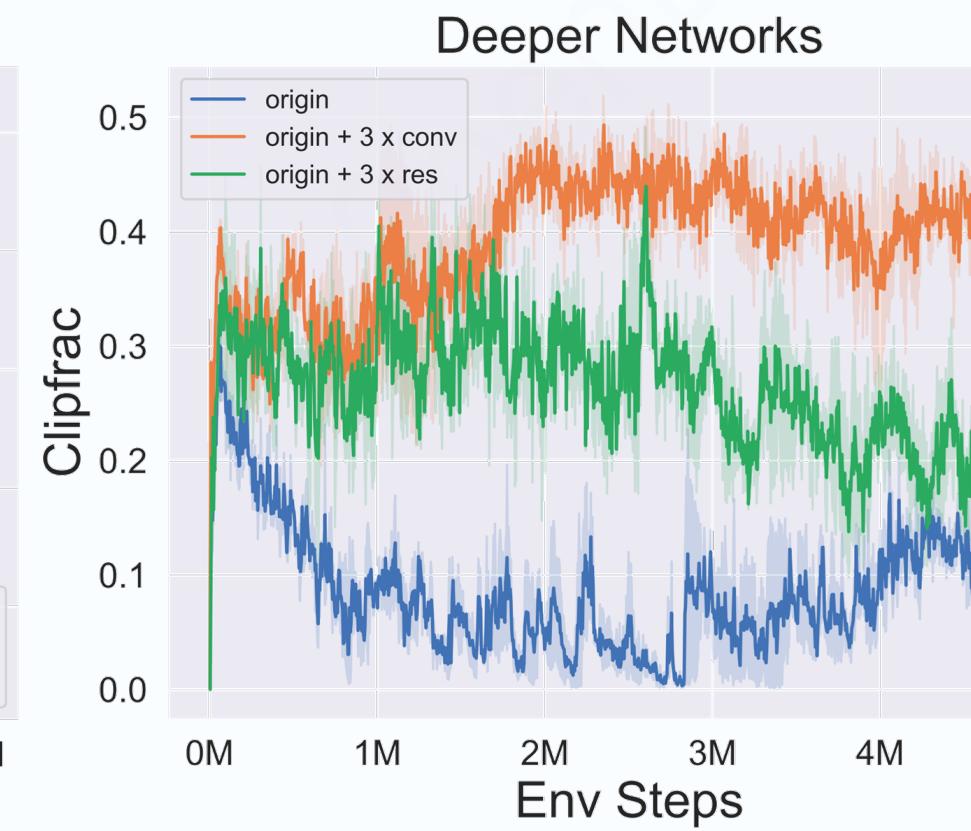
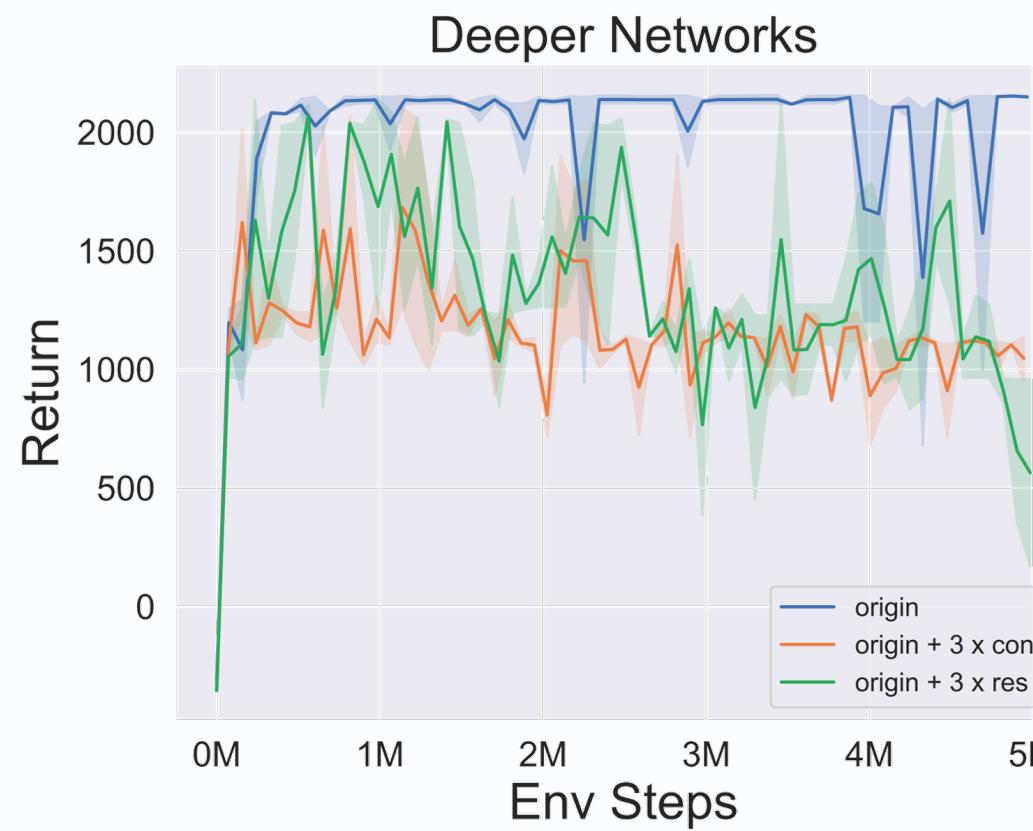
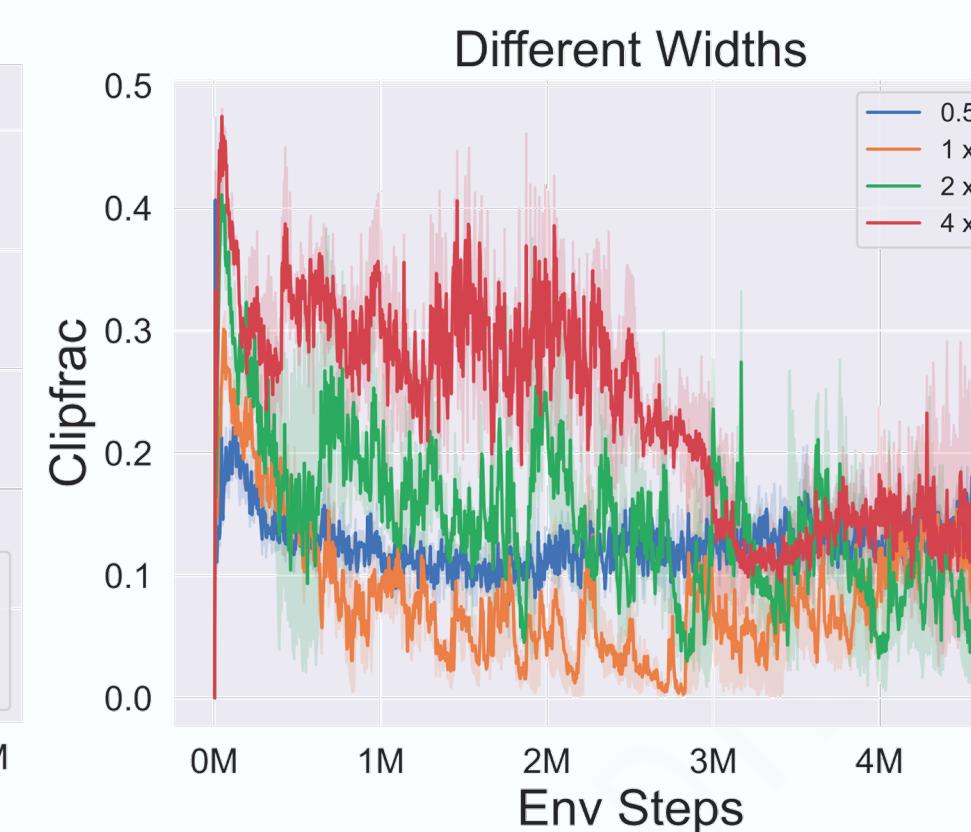
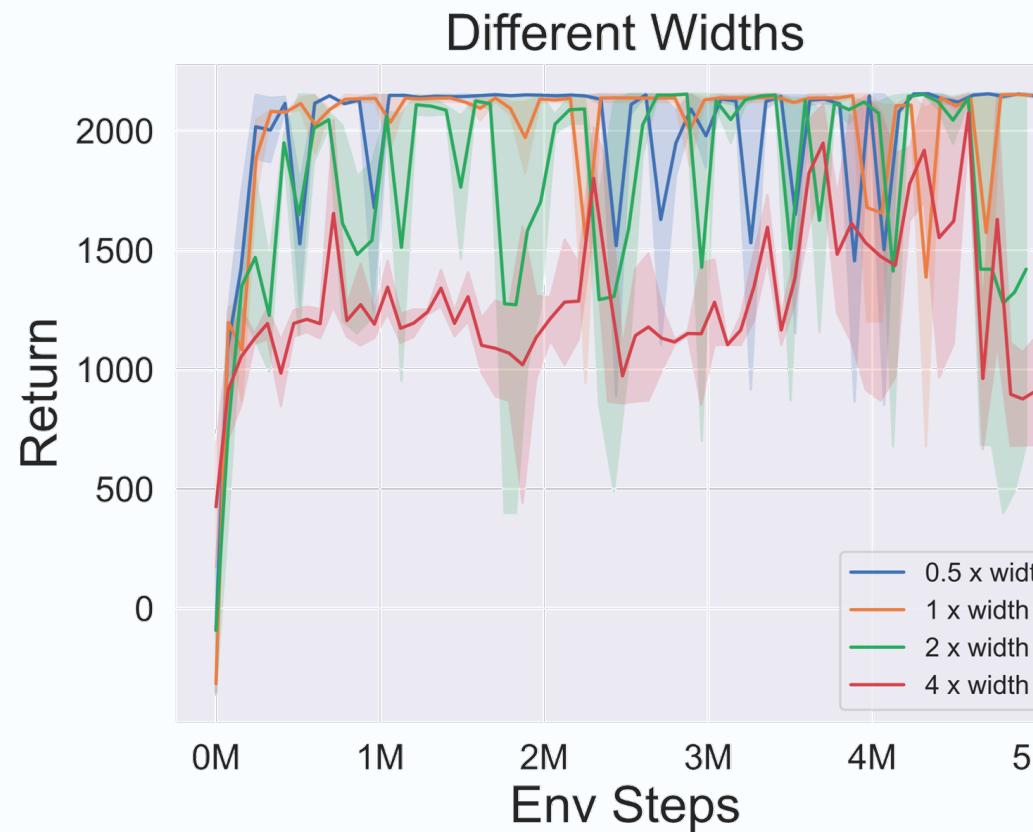
光流



实验细节：<https://github.com/opendilab/PPOxFamily/issues/8>

理论：PPO + 图片观察

● 网络架构



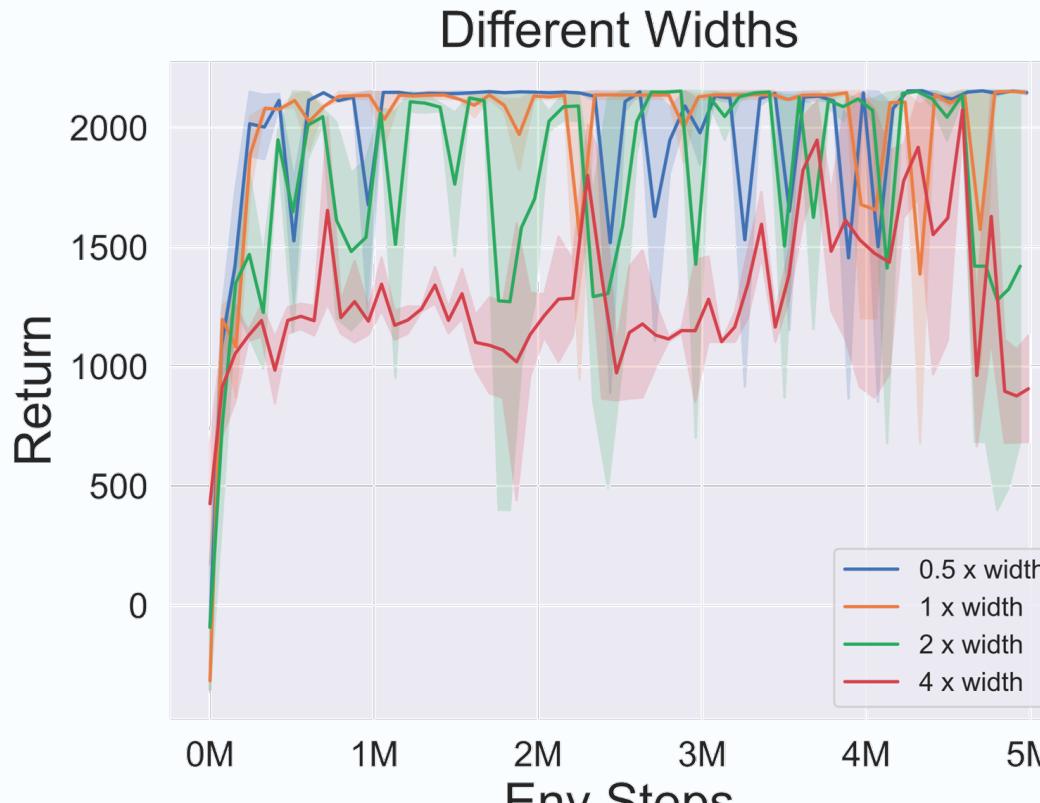
Key Facts:

- 单纯的加深 or 加宽网络并不总是有效的

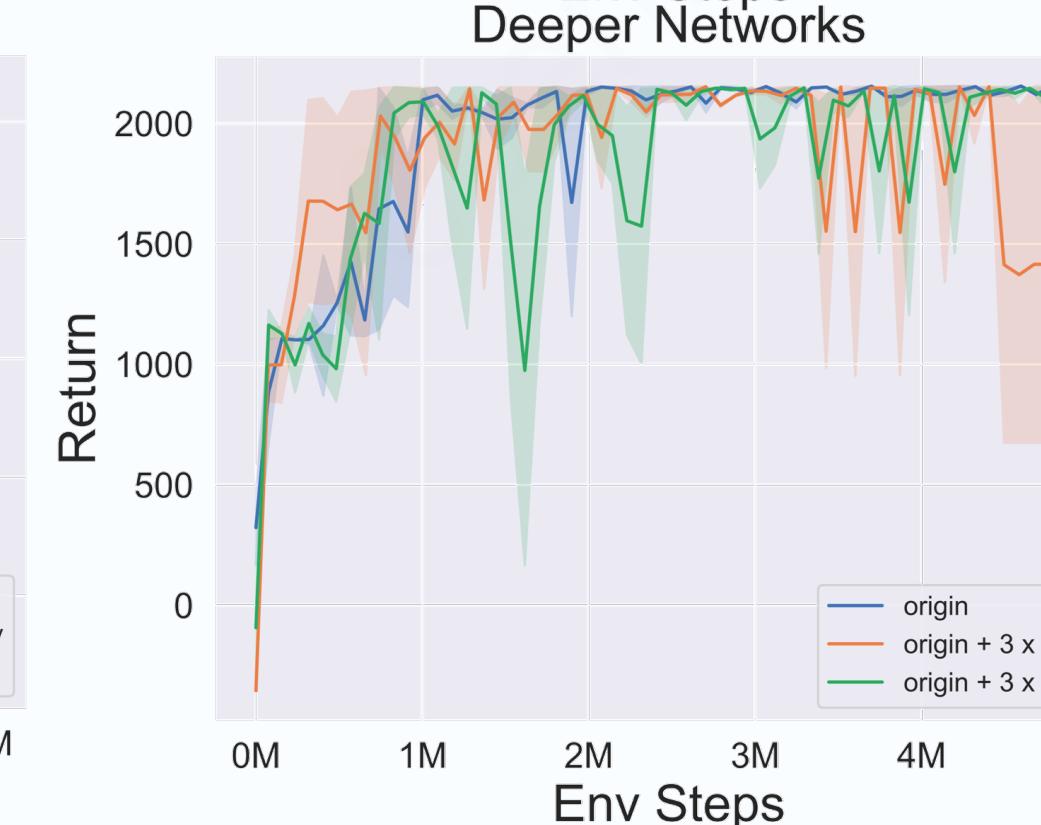
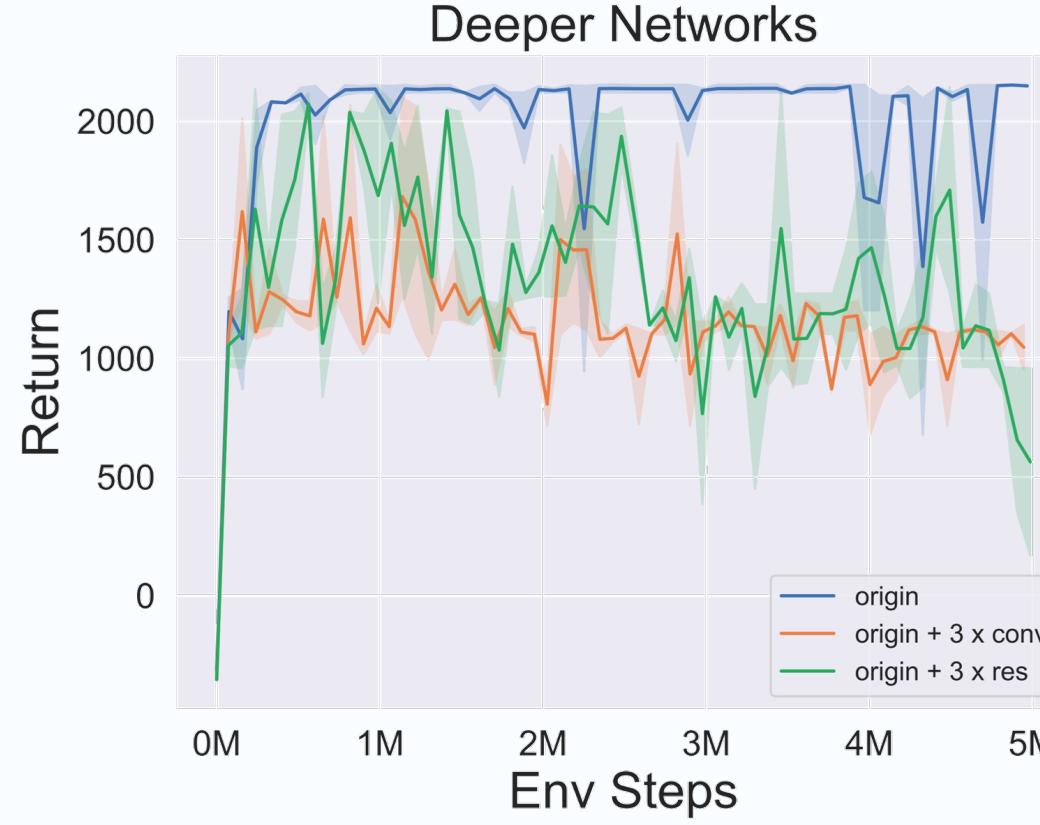
理论：PPO + 图片观察

网络架构

without LN

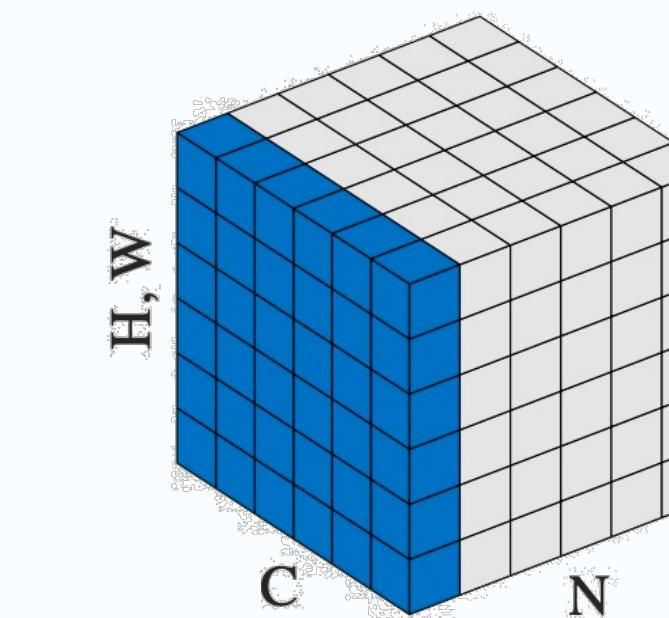


with LN



Key Facts:

- 网络一旦变大，就需要类似 LN 这样的模块来稳定训练



代码：巧用 Env Wrapper

Stars 1.1k bilibili video course Follow @opendilab 44
View code on GitHub

PyTorch 代码实现：用于 PPO 训练的超级马里奥兄弟仿真环境设置
Related Link

概述
创建超级马里奥环境，并使用一系列 Env Wrapper 来实现各种预处理变换

```
1 from ding.envs import DingEnvWrapper
2 from ding.envs.env_wrappers import MaxAndSkipWrapper, WarpFrameWrapper, ScaledFloatFrameWrapper,
3     FrameStackWrapper, \
4     EvalEpisodeReturnEnv
5 import gym_super_mario_bros
6 from nes_py.wrappers import JoypadSpace
7
8 def wrapped_mario_env():

9     return DingEnvWrapper(
10         JoypadSpace(gym_super_mario_bros.make("SuperMarioBros-1-1-v0"), [
11             ["right"], ["right", "A"]]),
12         cfg={
13             'env_wrapper': [
14                 lambda env: MaxAndSkipWrapper(env, skip=4),
15                 lambda env: WarpFrameWrapper(env, size=84),
16             ]
17         }
18     )
19
20     eval_episode_return = EvalEpisodeReturnEnv(wrapped_mario_env())
21
22     return eval_episode_return
```

转换为 DI-engine 的环境接口格式

将动作空间简化为2维的离散空间（向右走，向右上跳）

添加各类 Env Wrapper 的相关配置

添加环境帧采样 Wrapper，每隔四帧采样一帧作为智能体的观测信息

添加图片格式和尺寸统一化 Wrapper

完整示例：<https://opendilab.github.io/PPOxFamily/>

实践：PPO + 超级马里奥

Level 1-1



Level 1-4

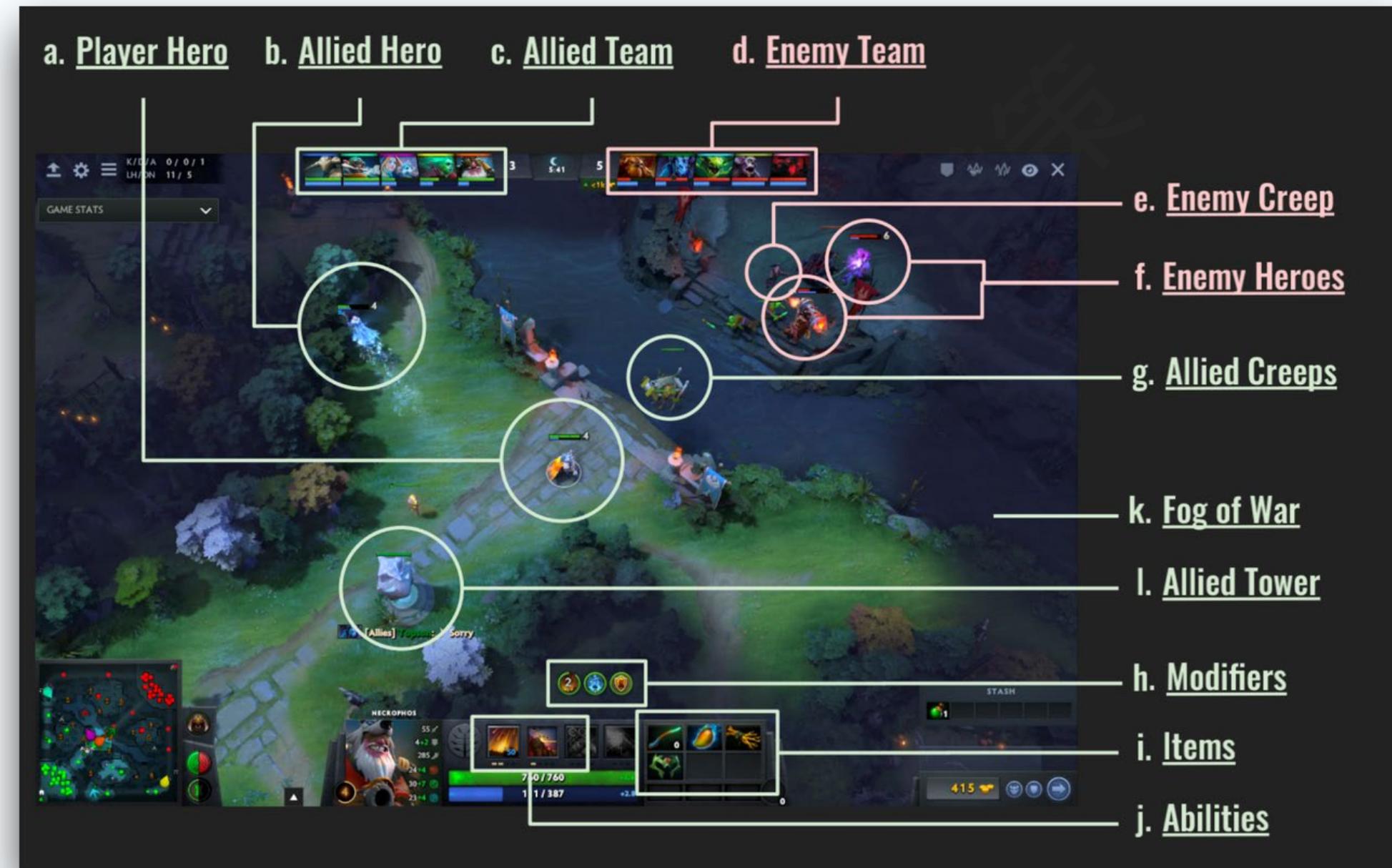


Super Mario Bros: <https://github.com/Kautenja/gym-super-mario-bros>

DI-engine + Mario Doc: https://di-engine-docs.readthedocs.io/zh_CN/latest/13_envs/gym_super_mario_bros_zh.html

视频样例: <https://github.com/opendilab/PPOxFamily/issues/8>

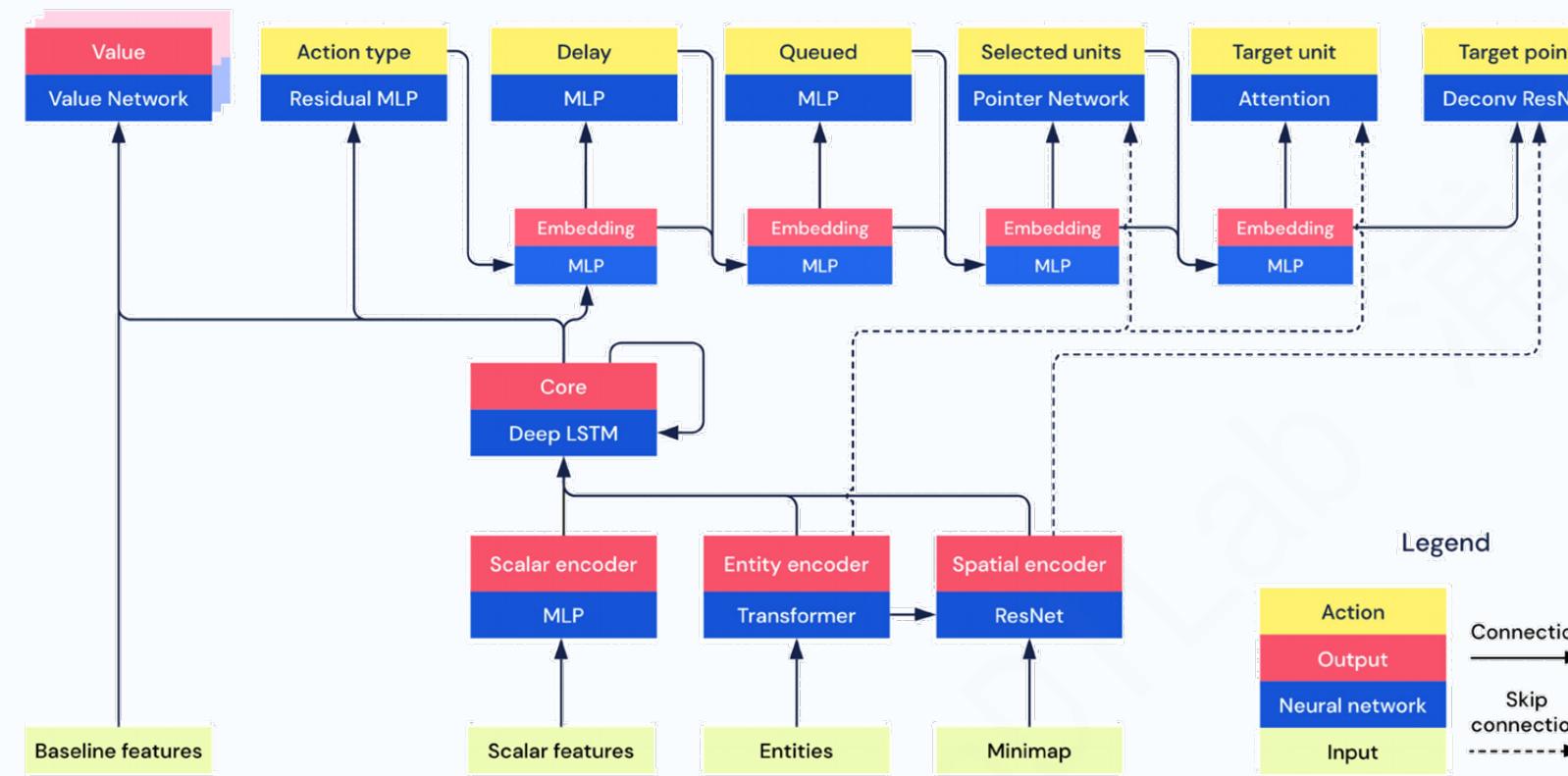
复杂结构化观察空间



OpenAI Five: <https://openai.com/five/>



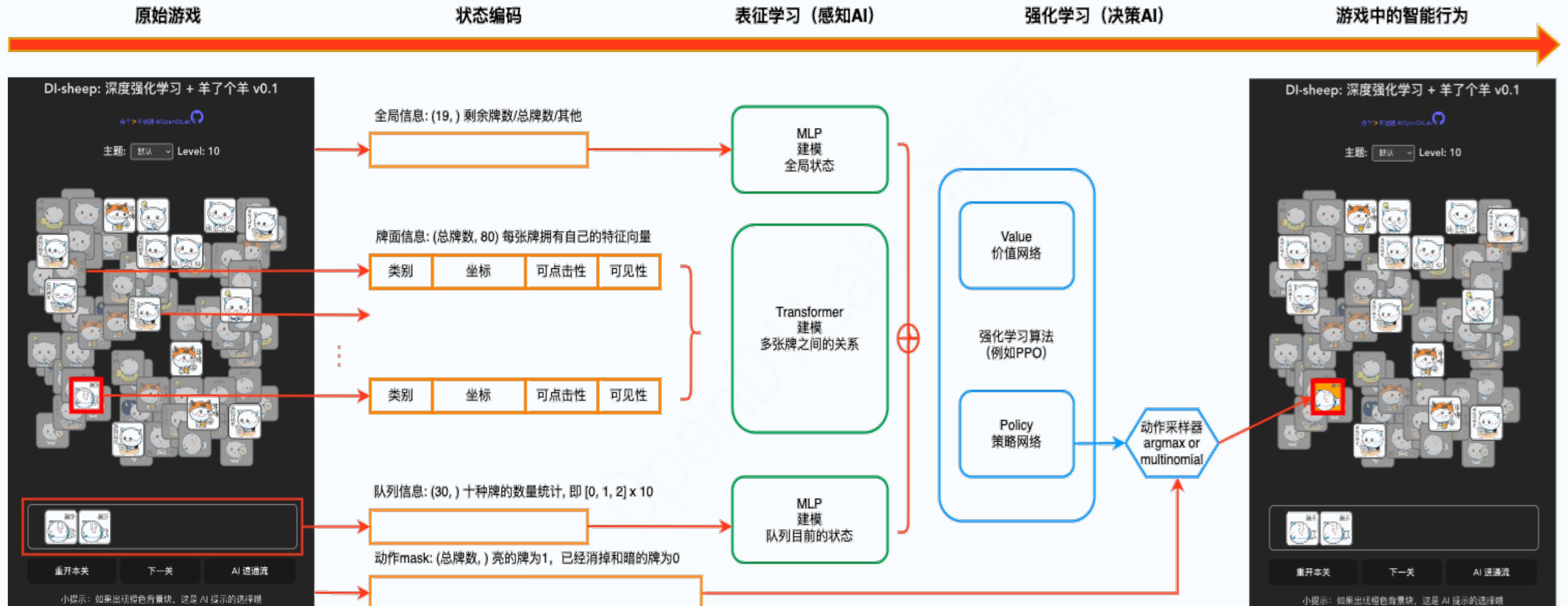
理论：PPO + 结构化观察 ● 概述



解决方案：

- 不同模态的数据需要结合对应领域经典的建模方法和网络结构设计
- 重要的信息利用 skip connection 进行补充
- 要根据决策任务的特点灵活组合各种网络结构设计方法

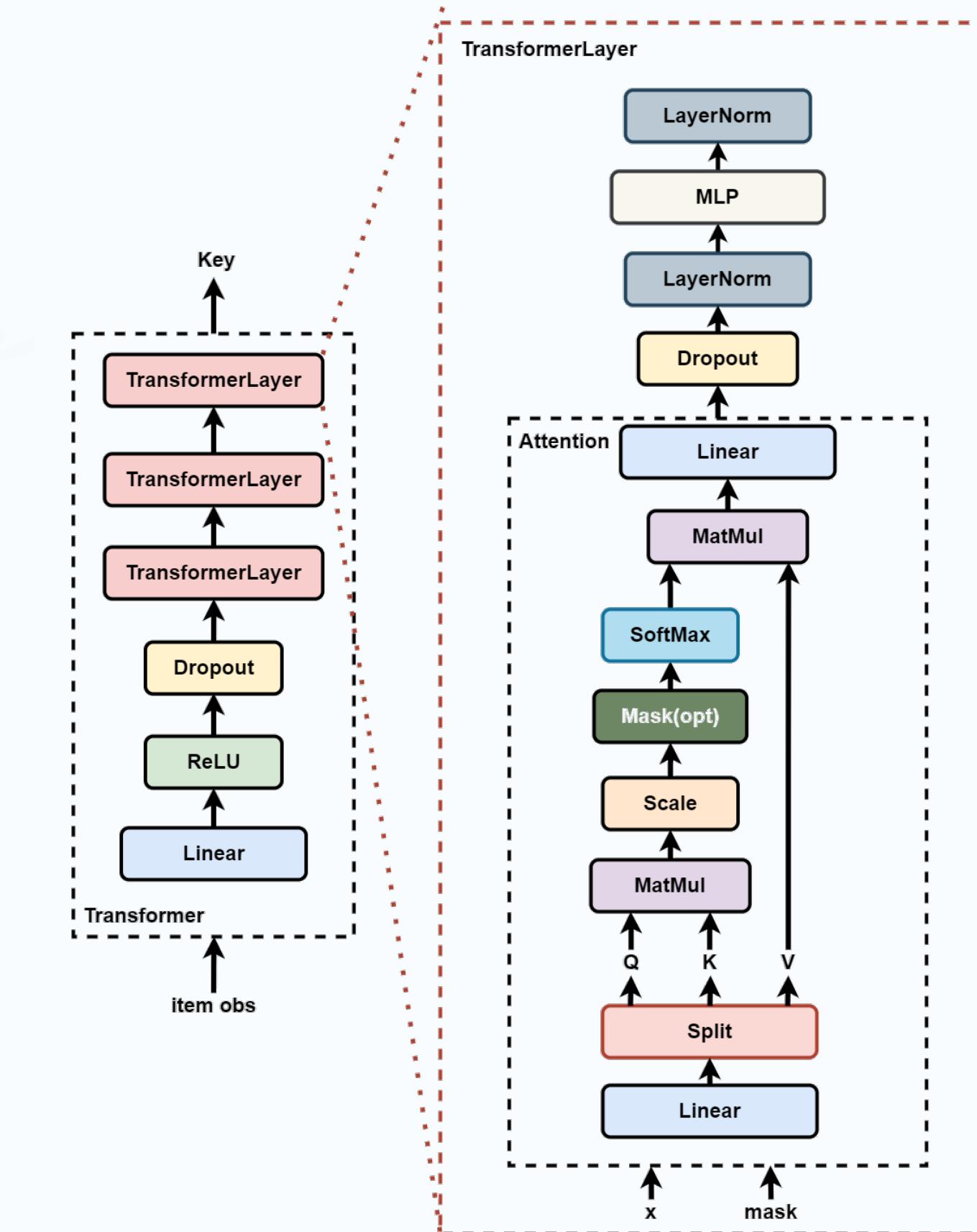
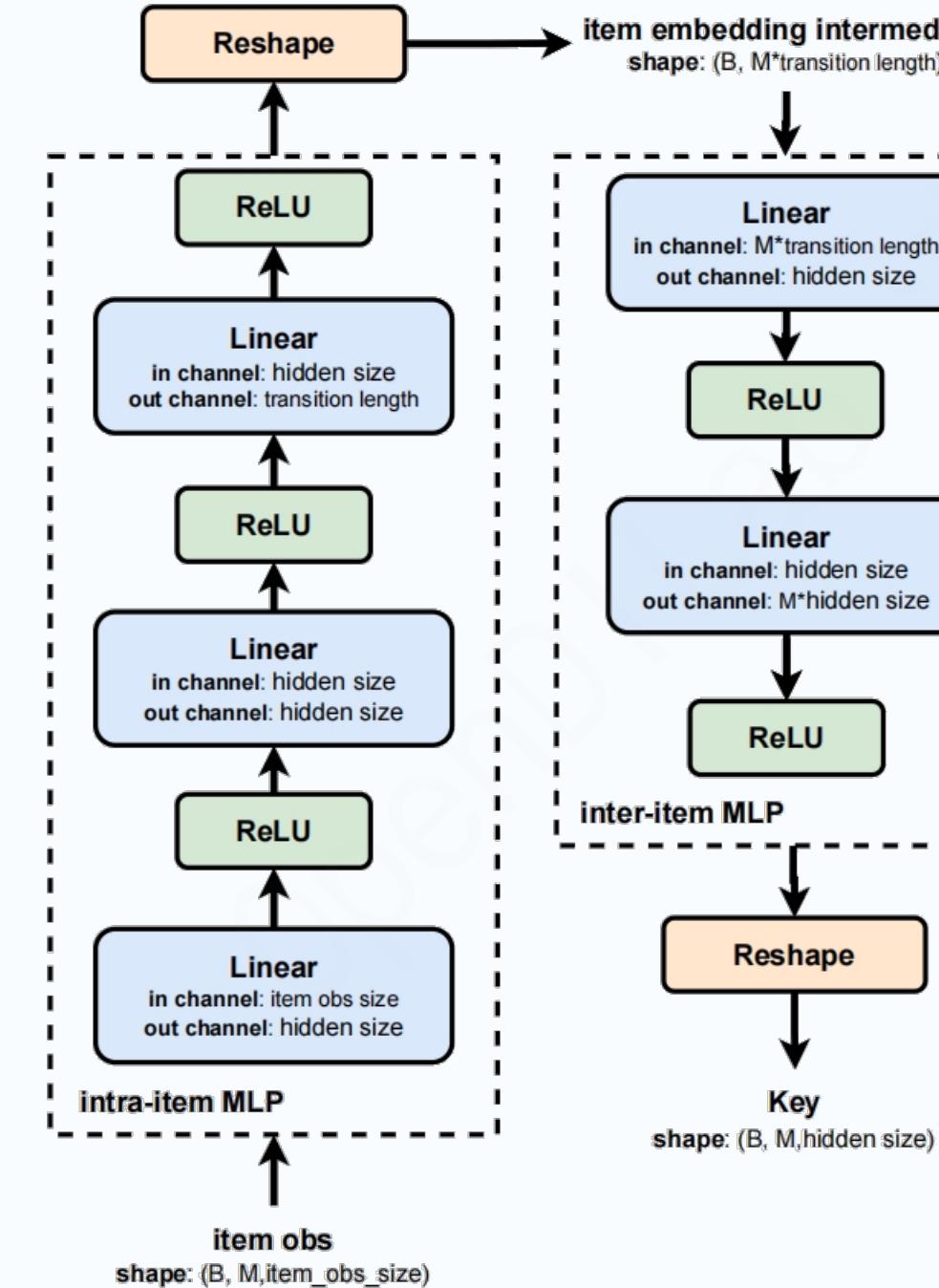
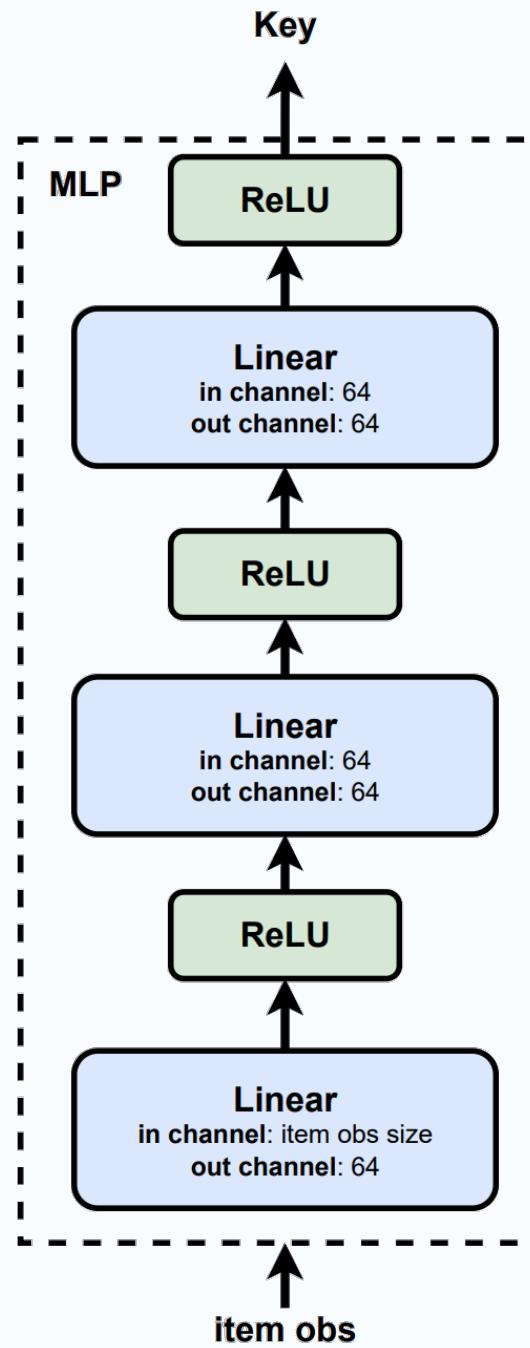
Struct 实践： PPO + 羊了个羊



DI-sheep: <https://github.com/opendilab/DI-sheep>

实践：PPO + 羊了个羊

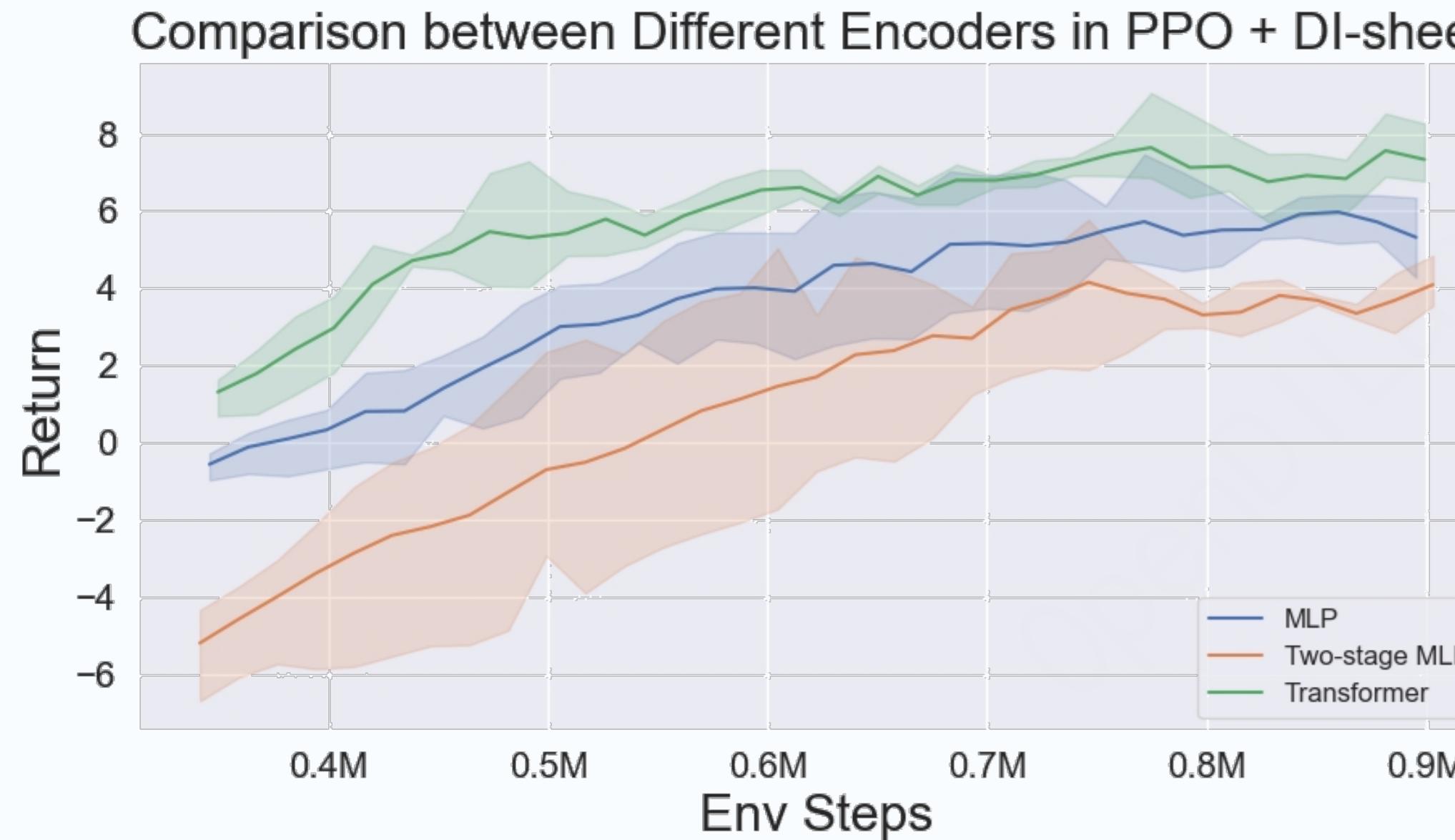
不同的牌面信息编码器

**MLP****Two-Stages MLP****Transformer**



实践：PPO + 羊了个羊

(实验对比) 不同的牌面信息编码器



实验细节：<https://github.com/opendilab/PPOxFamily/issues/8>

Key Facts:

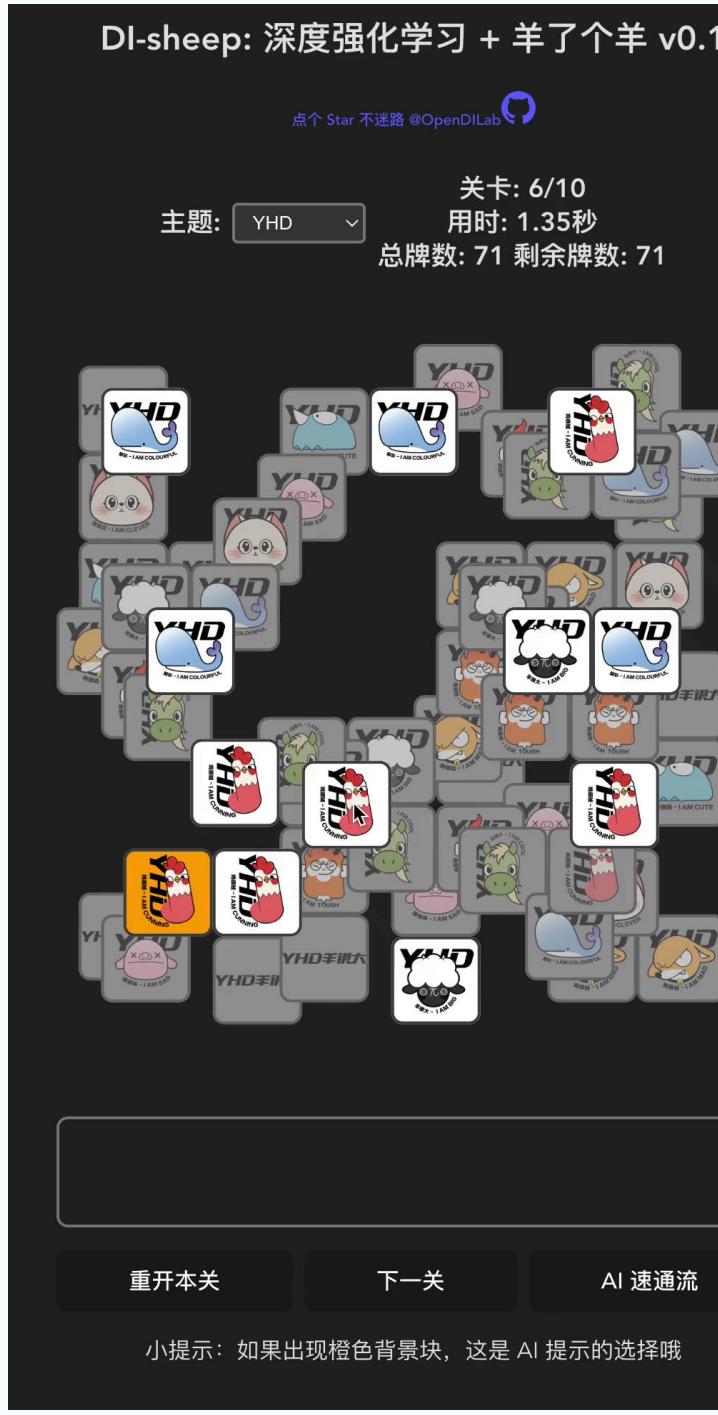
- 环境角度：多张牌之间的关系建模很重要
- 计算角度：参数量大小不一定重要，合理利用参数进行计算很重要
- 算法角度：网络结构的设计代表问题的上限，网络的训练方式决定问题的下限

实践：PPO + 羊了个羊

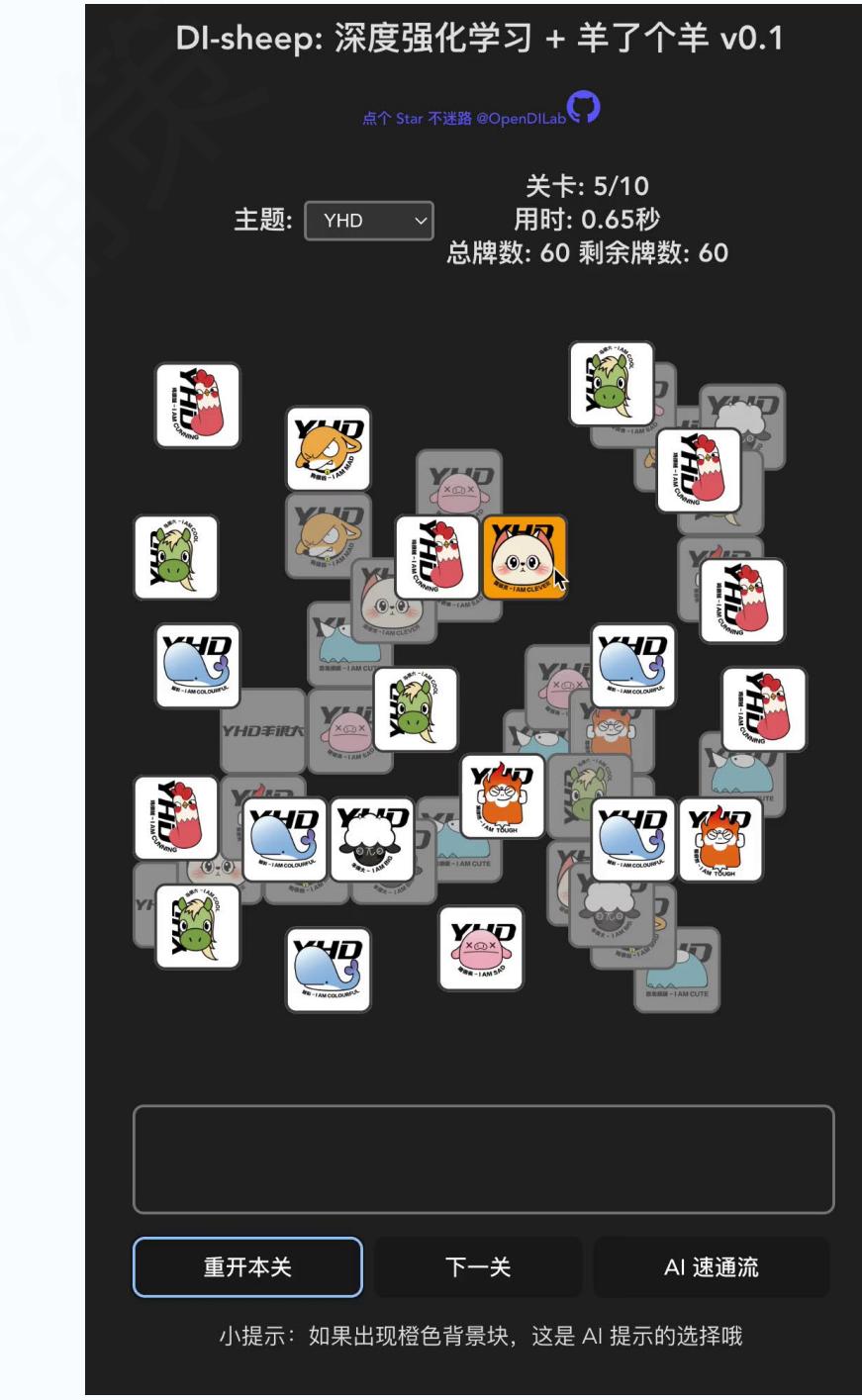
乖巧型



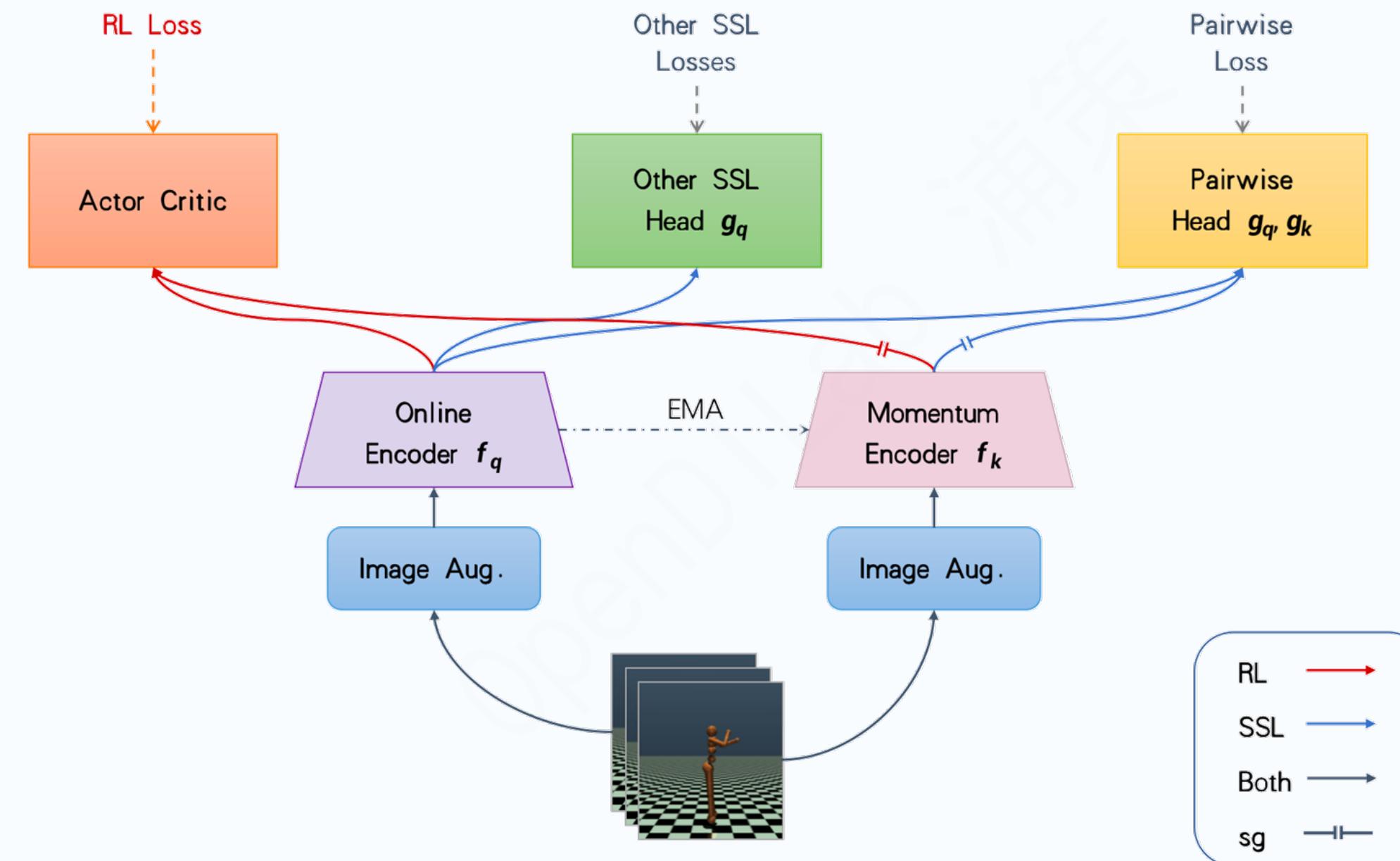
叛逆型



合作型



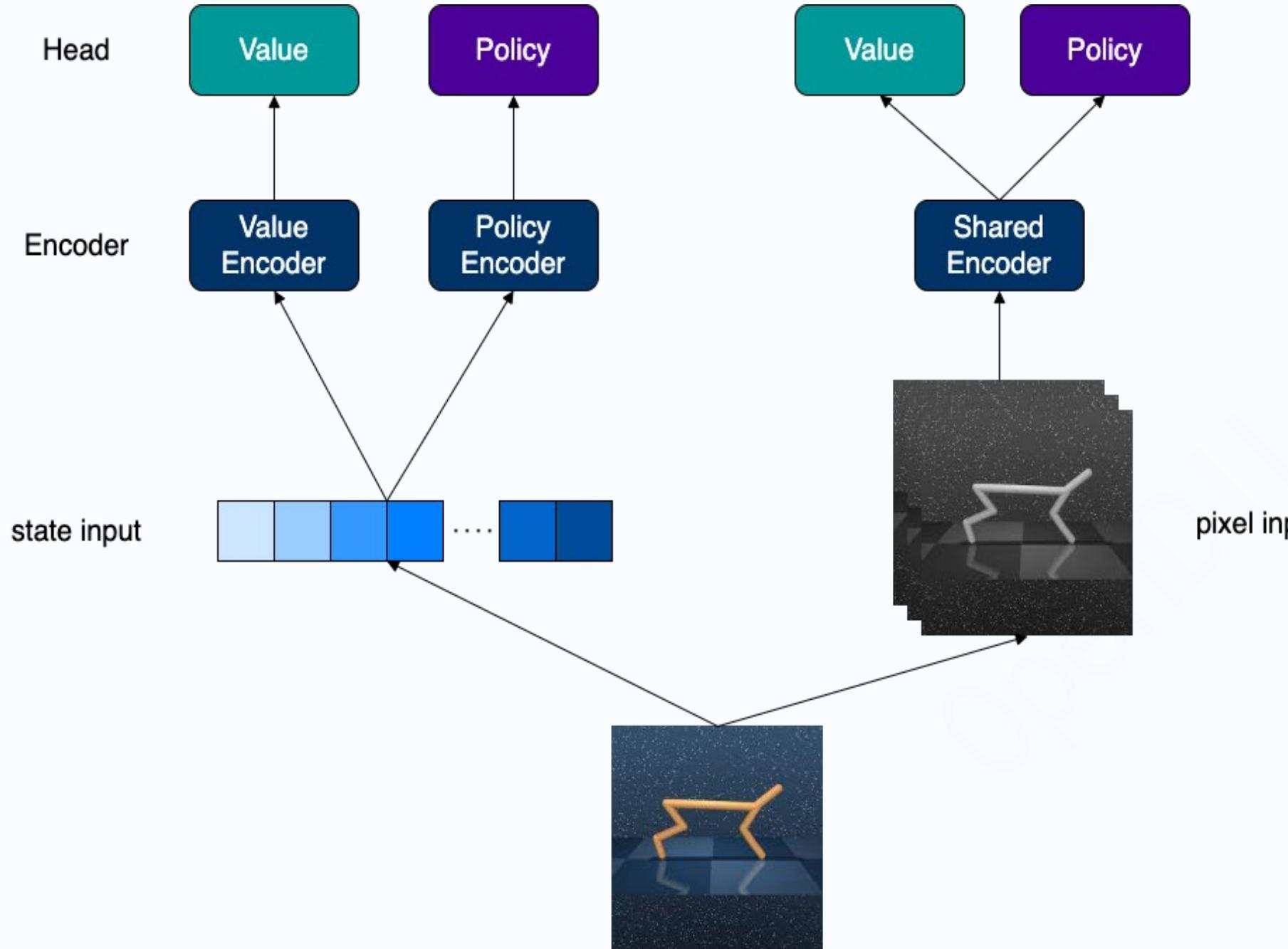
通用观察空间训练方法



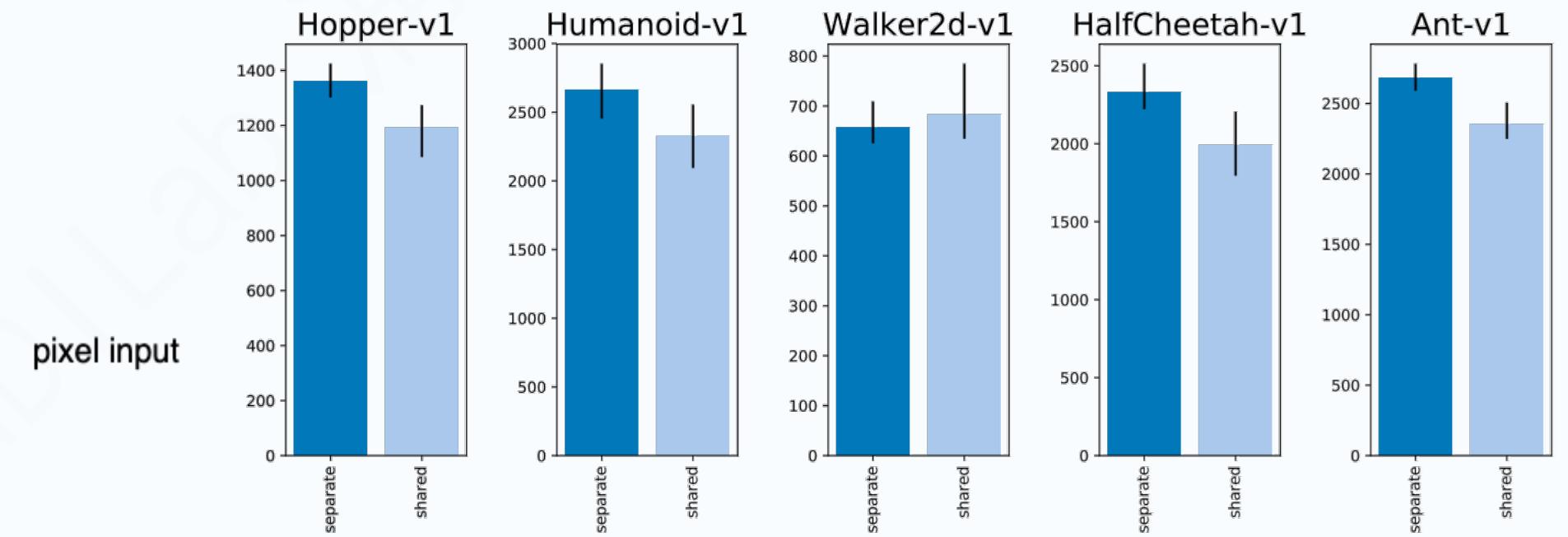
Does Self-supervised Learning Really Improve Reinforcement Learning from Pixels? <https://openreview.net/forum?id=fVslVNBfjd8>

观察空间表征学习补充材料: https://github.com/opendilab/PPOxFamily/blob/main/chapter3_obs/chapter3_supp_representation.pdf

理论：Actor和Critic是否共享编码器

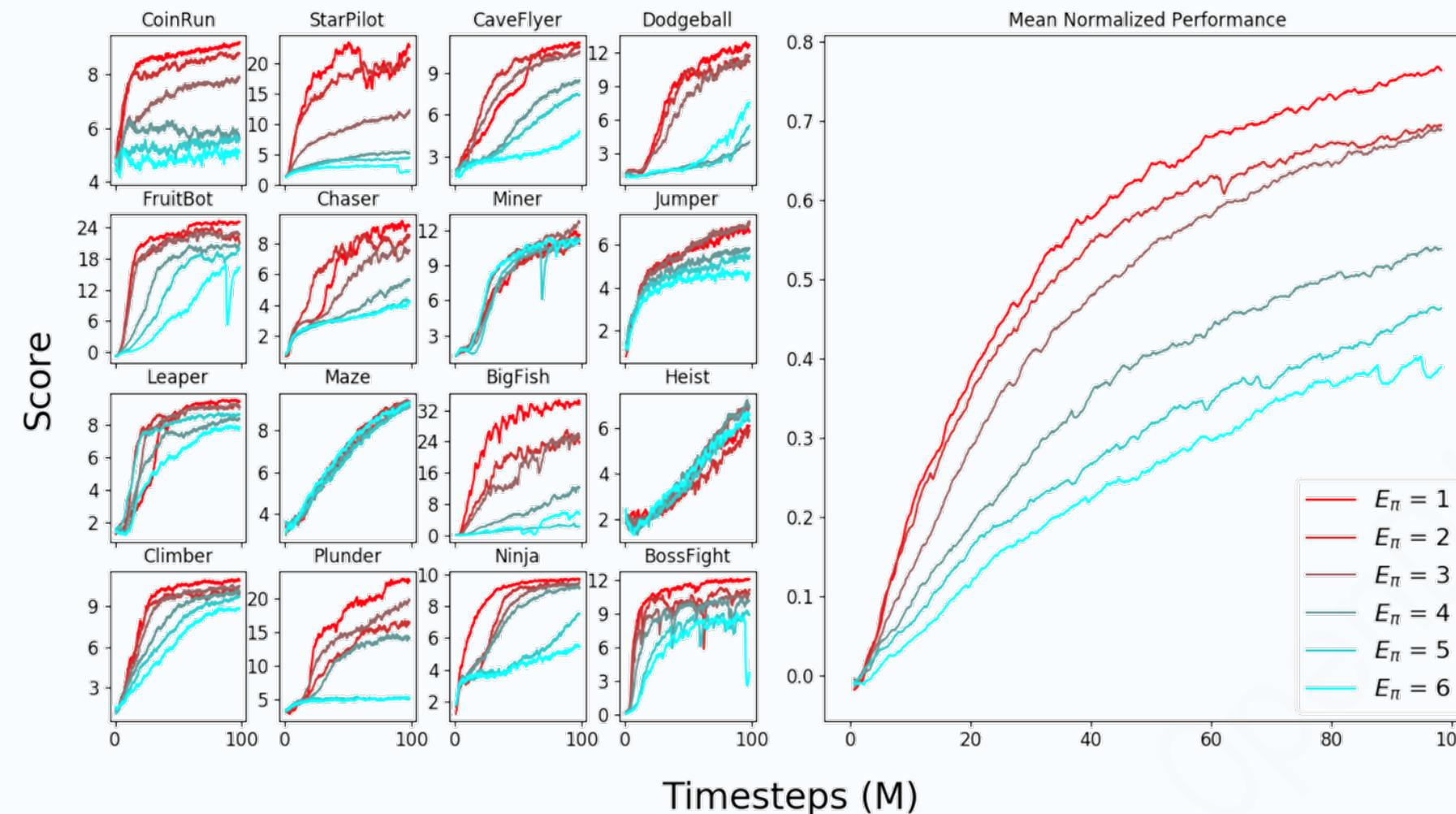


经验原则：高维共享，低维不共享

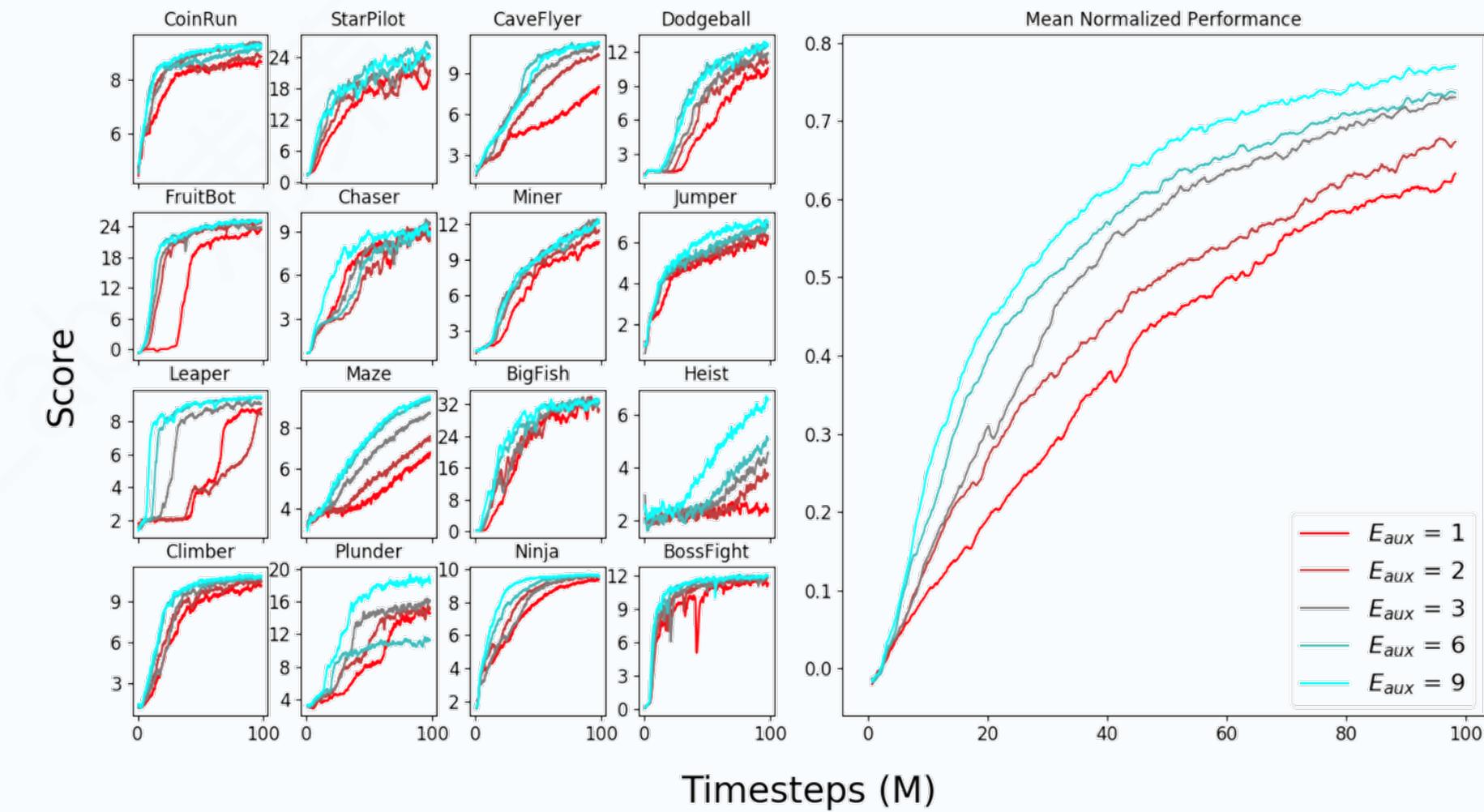


理论：Actor和Critic训练的不同特性

Policy Reuse



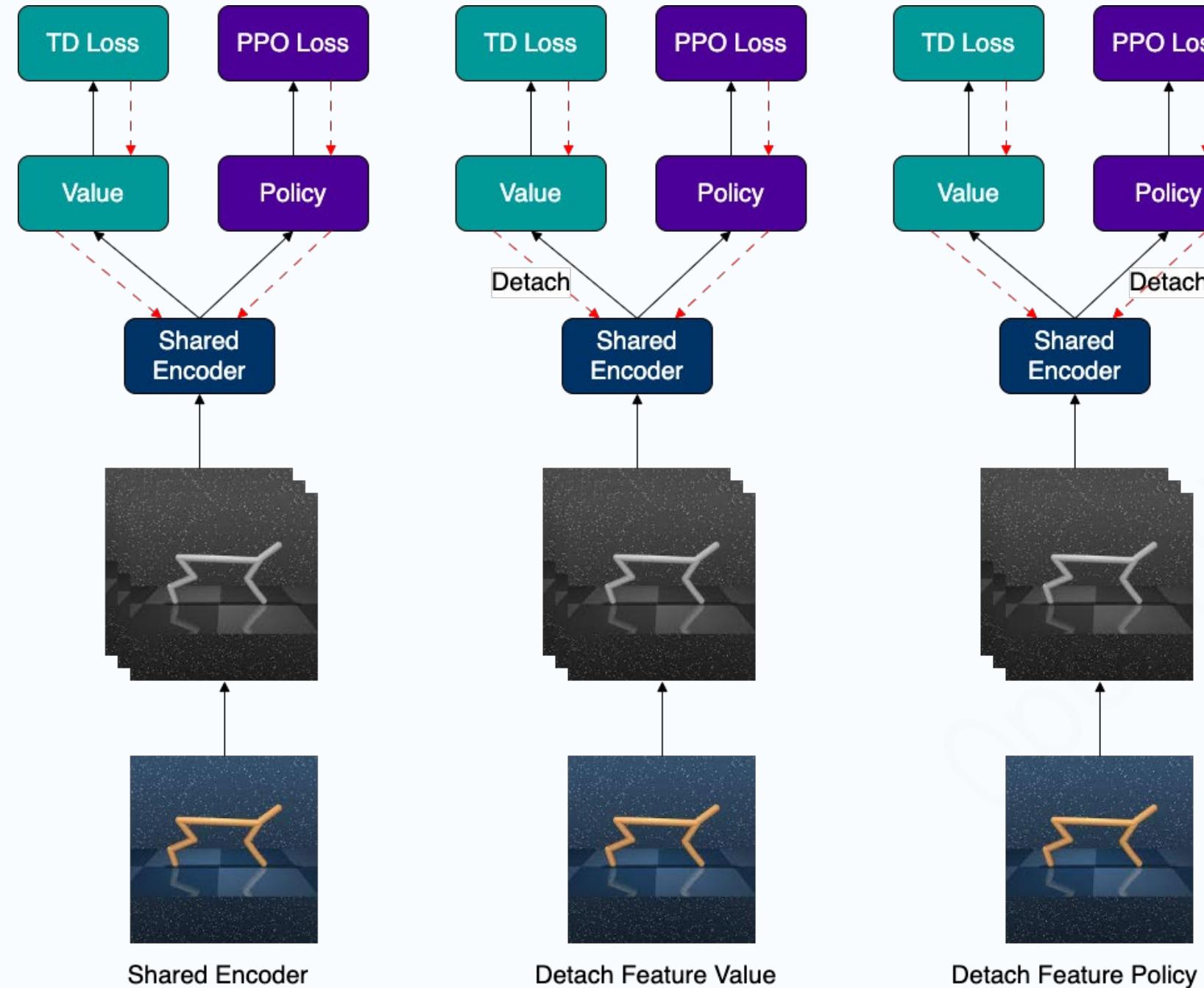
Value Reuse



策略函数 (Policy) : 对数据新旧程度和复用次数很敏感，尽可能每个数据只用一次

价值函数 (Value) : 需要更高的数据多样性，可以通过复用数据来提升性能

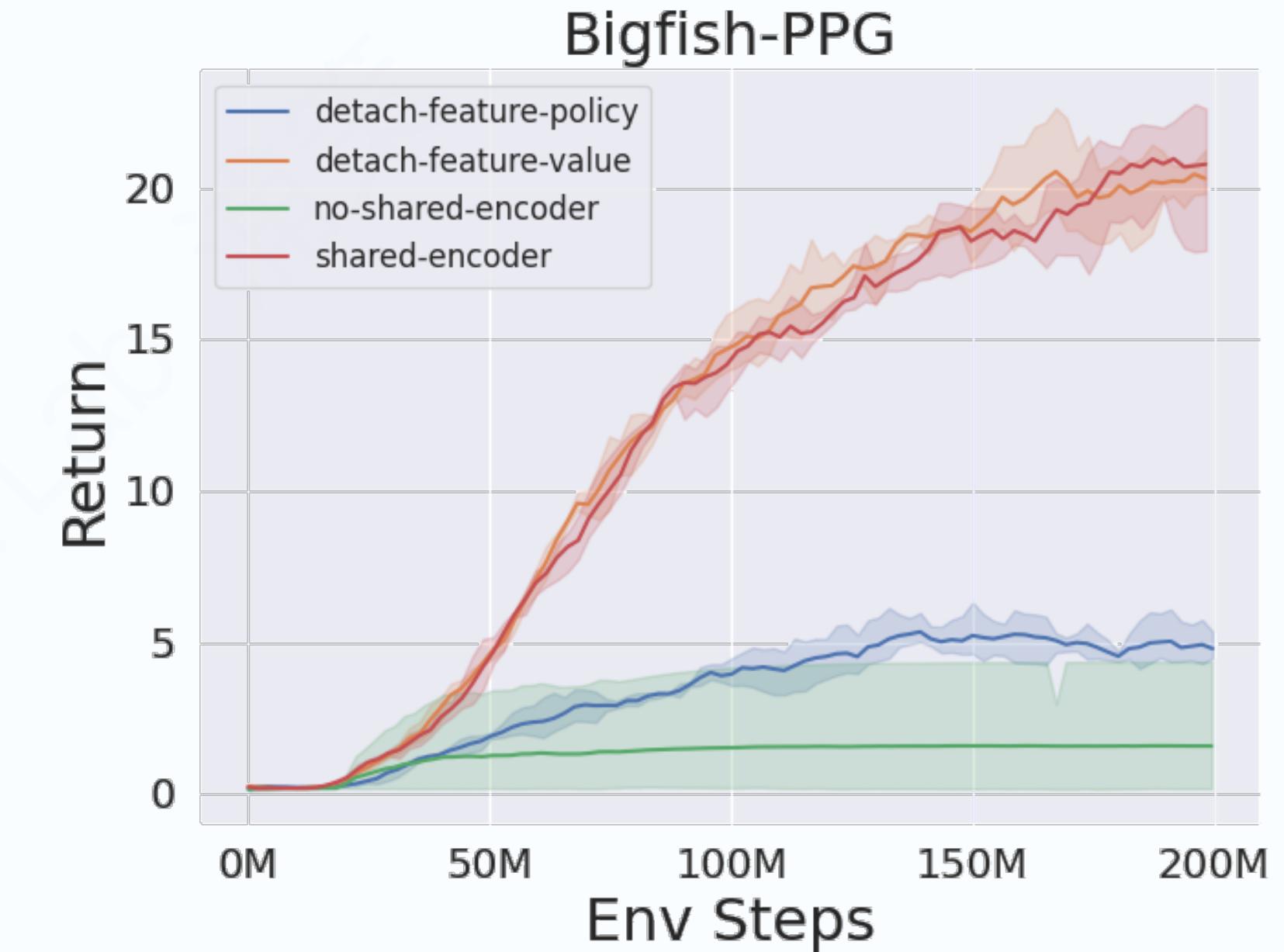
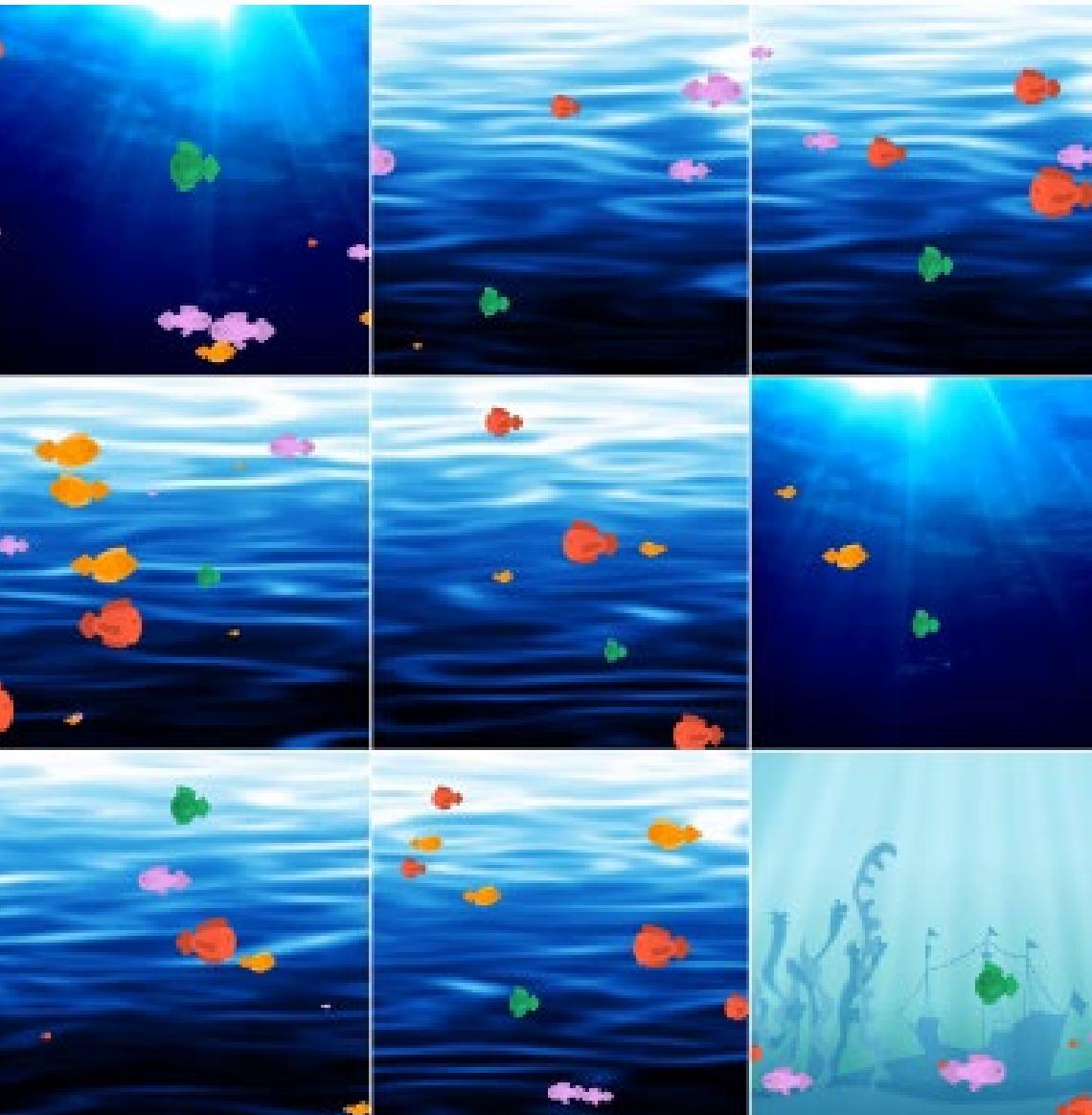
理论：Actor和Critic分离编码器训练



对于高维观察（比如图片）空间：

- Policy 和 Value 共享编码器是非常重要的
- 但 Policy 和 Value 对于编码器具体有什么样的影响关系，二者是缺一不可，还是各有不同的作用？

实践：PPO+大鱼吃小鱼(Procgen)



Procgen: <https://openai.com/five/>

DI-engine + Procgen Doc: https://di-engine-docs.readthedocs.io/zh_CN/latest/13_envs/procgen_zh.html

实验细节: <https://github.com/opendilab/PPOxFamily/issues/8>

代码：如何理解计算图中的梯度流动

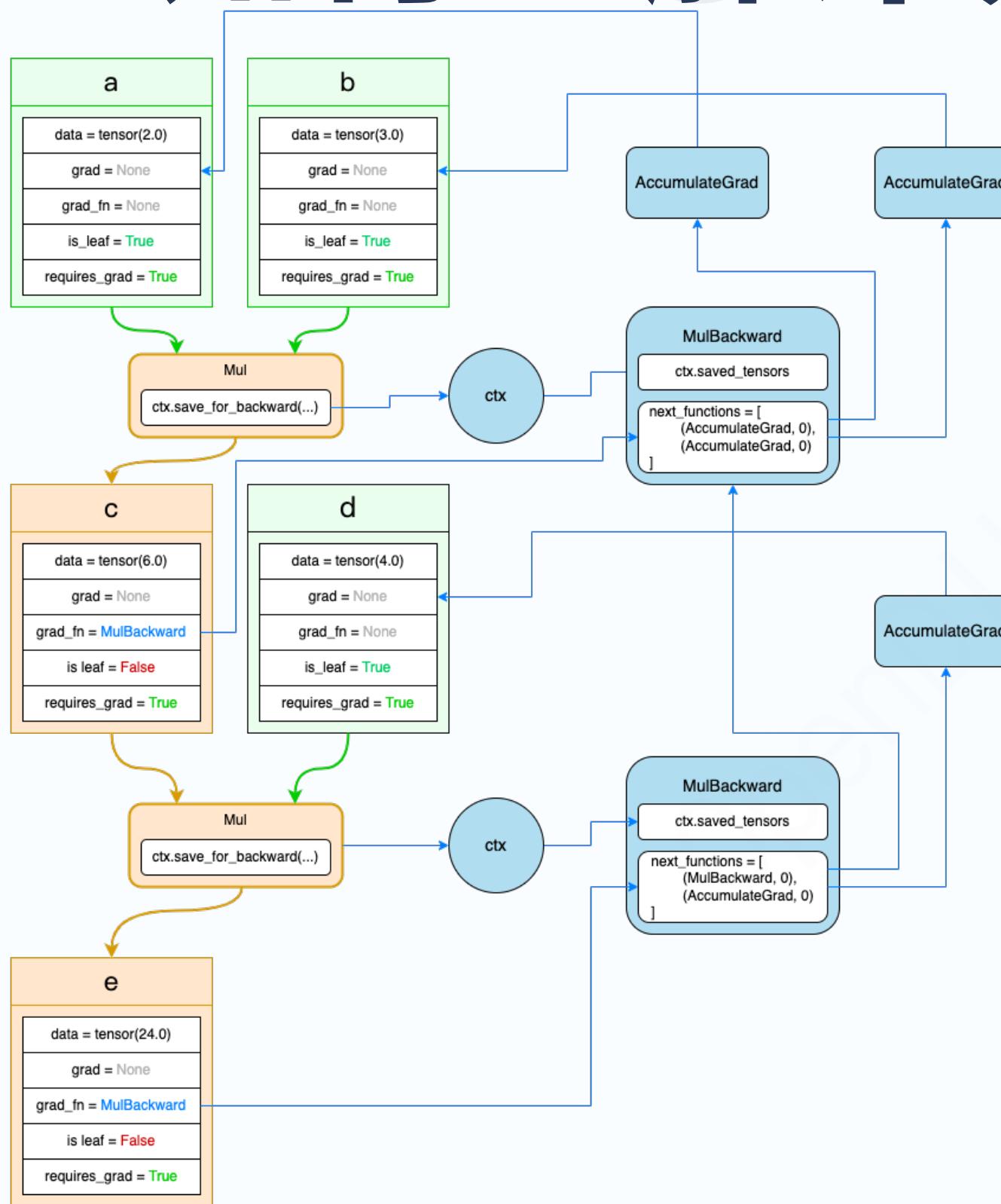
```
a = torch.tensor(2.0, requires_grad=True)
```

```
b = torch.tensor(3.0, requires_grad=True)
```

$c = a * b$

```
d = torch.tensor(4.0, requires_grad=True)
```

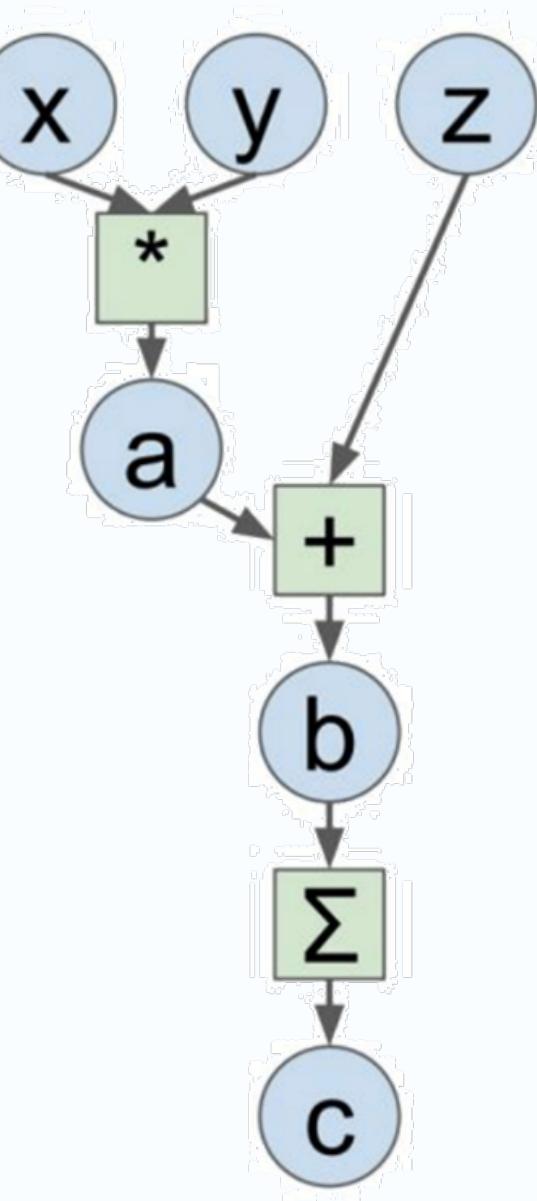
$e = c * d$



```

1 import numpy as np
2 np.random.seed(0)
3
4
5 # generate data
6 B, D = 3, 4
7
8 x = np.random.randn(B, D)
9 y = np.random.randn(B, D)
10 z = np.random.randn(B, D)
11
12 # forward
13 a = x * y
14 b = a + z
15 c = np.sum(b)
16
17 # backward
18 grad_c = 1.0
19 grad_b = grad_c * np.ones((B, D))
20 grad_a = grad_b.copy()
21 grad_z = grad_b.copy()
22 grad_x = grad_a * y
23 grad_y = grad_a * x

```



总结：PPO + 观察空间

小节	技术要点	代表决策任务
向量观察空间	<ul style="list-style-type: none"> 统一量纲+特征编码设计 不变性与等变性概述 	机器人控制 交通信号控制
图片观察空间	<ul style="list-style-type: none"> 视觉信息预处理+运动信息建模 网络架构设计 (LN) 	视频游戏 (超级马里奥) 棋类游戏 (围棋)
复杂结构化观察空间	<ul style="list-style-type: none"> 应用对应模态数据经典的网络架构 灵活利用 skip connection 根据数据模态特性定制化设计 	自动驾驶 三消类游戏 (羊了个羊) 即时战略游戏 (星际争霸2)
通用训练方法	<ul style="list-style-type: none"> Actor和Critic是否共享编码器 表征学习 + 强化学习联合优化 	/

下节预告

(四) 探索之法：解密稀疏奖励空间

- 稀疏奖励空间的万能钥匙
- 多尺度奖励空间的锦囊妙计