

# 人类动作模仿和强化

基于模型的方法（model-based methods）和无模型的方法

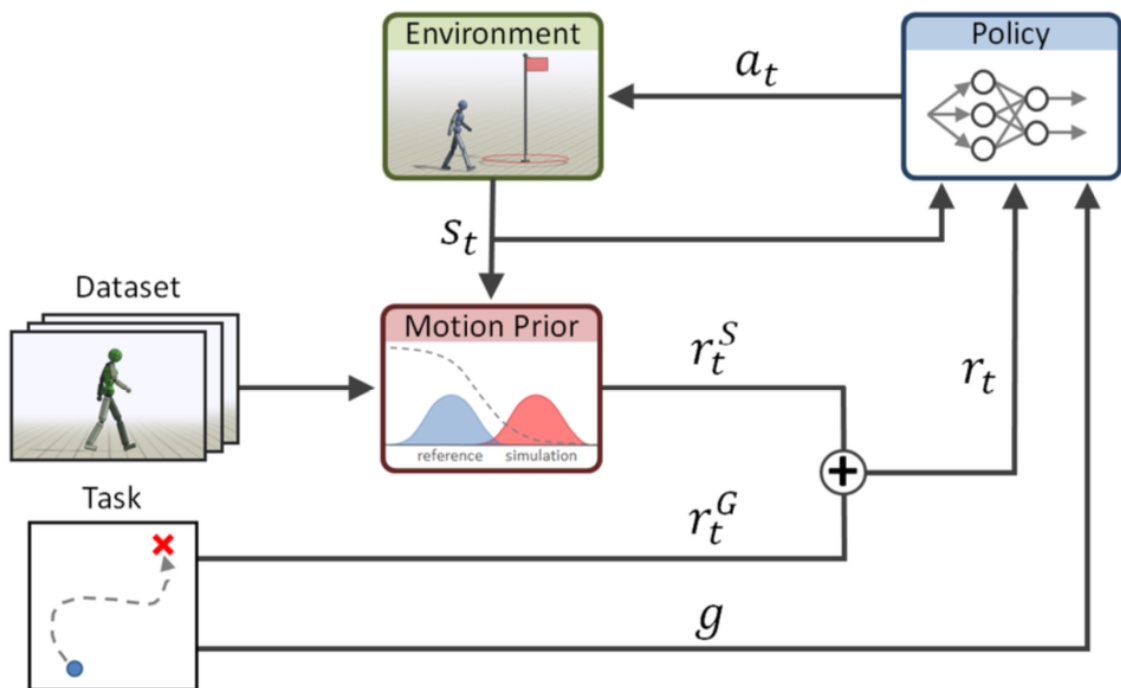
（model-free methods）；其中的无模型方法又可以类为基于奖励函数工程法与生成式对抗强化学习法。

奖励函数方式：

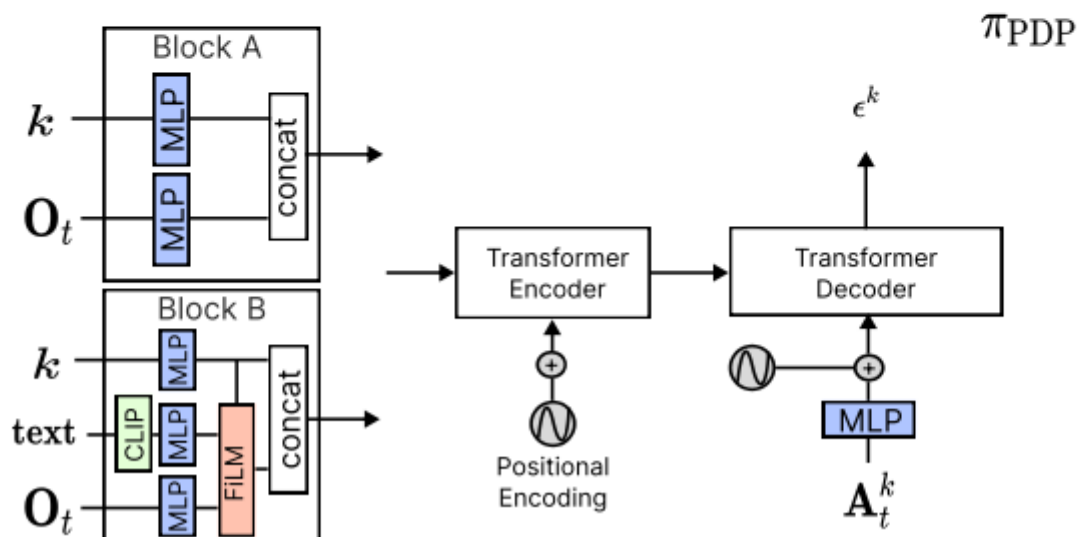
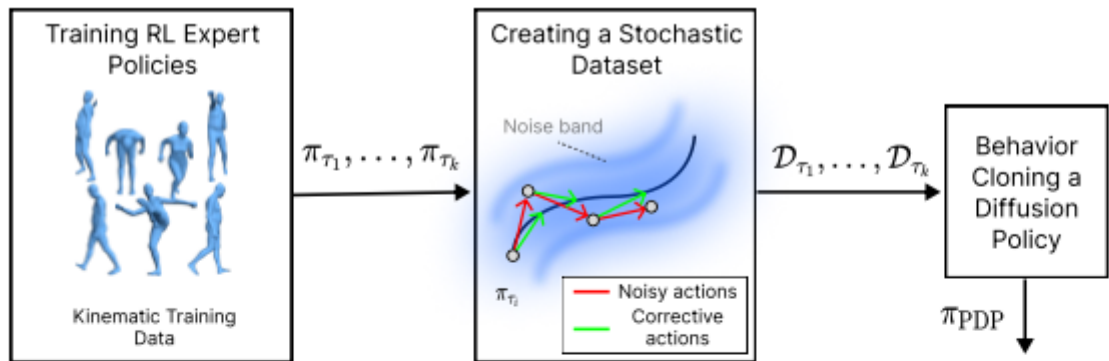
Real- world humanoid locomotion with  
reinforcement learning （已实现）

生成对抗方式：

AMP:



PDP:



AMP 和 PDP 确实主要关注角色自身的动作拟真性和基础物理属性（如刚度和阻尼），但对复杂环境交互（如动态物体操作、非刚性接触、地形适应等）的支持有限。以下从环境交互的复杂性、现有方法的局限性、改进方向三个维度展开分析：

## 1. 环境交互的复杂性挑战

复杂环境交互通常涉及以下问题：

- **动态物体操作：**推/拉物体、抓取、避障等需要实时感知和反馈控制。
- **非刚性接触：**与软体（如布料、泥地）交互时，接触力学模型复杂。
- **地形适应：**不同材质（冰面、沙地）、坡度、障碍物需调整步态。
- **多目标耦合：**同时完成多个任务（如搬运物体时保持平衡）。
- **长期规划：**在动态环境中预测未来状态并提前规划动作。

这些场景要求策略不仅生成拟人动作，还需 动态感知环境变化、建模物理交互、实时调整策略。

## 2. AMP 和 PDP 的局限性

### (1) AMP 的短板

- 静态环境假设：**AMP 的对抗训练依赖固定数据集，无法泛化到未见过的动态环境。
- 奖励设计局限：**任务奖励  $(r^G_t)$  通常针对单一目标（如到达目标点），难以建模多步交互（如“绕过障碍→推门→拾取物体”）。
- 判别器误导风险：**若环境交互动作未在专家数据中出现，判别器可能抑制策略的创新性动作。

### (2) PDP 的短板

- 专家策略覆盖性：**若 RL 专家未训练复杂交互任务（如开门、攀爬），扩散模型无法生成相关动作。
- 物理简化：**依赖刚体动力学假设，难以处理软体交互或流体动力学。
- 实时性不足：**扩散模型的迭代采样难以满足动态环境中的毫秒级响应需求。

## 3. 改进方向与前沿方案

### (1) 环境感知与状态表征

- 多模态传感器融合：**引入视觉、触觉、力觉等多模态输入，构建环境状态编码（如物体位置、材质属性）。
- 图神经网络（GNN）：**建模环境中物体间的拓扑关系（如“手-门把手-门”的交互链）。

### (2) 分层强化学习（HRL）

- 高层规划+底层控制：**
  - 高层策略：**生成交互目标（如“绕过桌子→接近门把手”）。

- **底层策略**：执行具体动作（如 AMP/PDP 生成的步态）。
- **案例**：Meta 的 **HSSD** 框架通过层级结构实现复杂物体操作。

### (3) 物理增强的模仿学习

- **混合仿真**：在训练中引入高保真物理引擎（如 MuJoCo 的软体扩展、NVIDIA Flex 的流体模拟），生成包含复杂交互的专家数据。
- **交互奖励设计**：在 AMP 的奖励函数中增加交互项（如接触力匹配度、物体位移奖励）。

### (4) 基于世界模型的预测

- **环境动力学建模**：学习环境的状态转移模型  $p(s_{t+1}|s_t,a_t)$ ，预测交互后果（如“推箱子后是否倾倒”）。
- **案例**：DeepMind 的 **Dreamer** 系列通过 latent world model 实现长时程规划。

### (5) 元学习与自适应策略

- **快速适应**：在策略中嵌入元学习机制，使角色能在线适应新环境（如从平地到沙地的步态调整）。
- **案例**：OpenAI 的 **PPO+Meta-Learning** 在仿真机器人中实现地形自适应。

## 4. 结合 AMP/PDP 的潜在方案

方法	改进思路
AMP + 分层规划	高层策略生成交互目标，底层 AMP 策略生成动作，任务奖励 $(r^G_t)$ 包含交互目标完成度。
PDP + 物理增强数据	在专家训练阶段引入复杂交互任务（如攀爬、搬运），扩散模型蒸馏多模态交互策略。
AMP/PDP + 世界模型	使用世界模型预测环境变化，调整动作生成（如提前避开滚落的石块）。

---

## 5. 总结与展望

AMP 和 PDP 在拟人化和基础物理控制上取得了进展，但复杂环境交互仍需突破：

- **核心问题：**现有方法缺乏对 **环境状态显式建模** 和 **多步物理因果推理** 的能力。
- **未来方向：**
  - 融合 **感知-规划-控制** 的全流程交互框架。
  - 利用 **神经物理引擎** 实现高保真交互仿真。
  - 开发 **通用交互奖励函数**，减少人工设计依赖。

最终目标是为角色赋予类似人类的 **环境理解与物理直觉**，使其在开放世界中自主完成复杂任务。