

# Some New Ideas in RL

# 传统算法回顾

Two posts:

- <https://lilianweng.github.io/lil-log/2018/02/19/a-long-peek-into-reinforcement-learning.html>
- <https://lilianweng.github.io/lil-log/2018/04/08/policy-gradient-algorithms.html>

# 传统算法回顾

Two papers

- <https://lil-log.com/2018/02/19/a-long-peek-into-reinforcement-learning/>

- <https://lil-log.com/2018/04/08/policy-gradient-theorem/>

- [What is Reinforcement Learning?](#)
  - [Key Concepts](#)
    - [Model: Transition and Reward](#)
    - [Policy](#)
    - [Value Function](#)
    - [Optimal Value and Policy](#)
  - [Markov Decision Processes](#)
  - [Bellman Equations](#)
    - [Bellman Expectation Equations](#)
    - [Bellman Optimality Equations](#)
- [Common Approaches](#)
  - [Dynamic Programming](#)
    - [Policy Evaluation](#)
    - [Policy Improvement](#)
    - [Policy Iteration](#)
  - [Monte-Carlo Methods](#)
  - [Temporal-Difference Learning](#)
    - [Bootstrapping](#)
    - [Value Estimation](#)
    - [SARSA: On-Policy TD control](#)
    - [Q-Learning: Off-policy TD control](#)
    - [Deep Q-Network](#)
  - [Combining TD and MC Learning](#)
  - [Policy Gradient](#)
    - [Policy Gradient Theorem](#)
    - [REINFORCE](#)
    - [Actor-Critic](#)
    - [A3C](#)
  - [Evolution Strategies](#)
- [Known Problems](#)
  - [Exploration-Exploitation Dilemma](#)
  - [Deadly Triad Issue](#)
- [Case Study: AlphaGo Zero](#)
- [References](#)

<https://lil-log.com/2018/02/19/a-long-peek-into-reinforcement-learning/>

<https://lil-log.com/2018/04/08/policy-gradient-theorem/>

# 传统算法回顾

Two papers

- <https://lil-log.com/2018/01/18/reinforcement-learning/>
- <https://lil-log.com/2018/01/18/policy-gradient-algorithms/>

- What is Reinforcement Learning?
  - Key Concepts
    - Model: Transition and Reward
    - Policy
    - Value Function
    - Optimal Value and Policy
  - Markov Decision Processes
  - Bellman Equations
    - Bellman Expectation Equations
    - Bellman Optimality Equations
  - Common Approaches
    - Dynamic Programming
      - Policy Evaluation
      - Policy Improvement
      - Policy Iteration
    - Monte-Carlo Methods
    - Temporal-Difference Learning
      - Bootstrapping
      - Value Estimation
      - SARSA: On-Policy TD control
      - Q-Learning: Off-policy TD control
      - Deep Q-Network
    - Combining TD and MC Learning
    - Policy Gradient
      - Policy Gradient Theorem
      - REINFORCE
      - Actor-Critic
      - A3C
    - Evolution Strategies
  - Known Problems
    - Exploration-Exploitation Dilemma
    - Deadly Triad Issue
  - Case Study: AlphaGo Zero
  - References

- What is Policy Gradient
  - Notations
  - Policy Gradient
  - Policy Gradient Theorem
  - Proof of Policy Gradient Theorem
- Policy Gradient Algorithms
  - REINFORCE
  - Actor-Critic
  - Off-Policy Policy Gradient
    - A3C
    - A2C
    - DPG
    - DDPG
    - D4PG
    - MADDPG
    - TRPO
    - PPO
    - ACER
    - ACTKR
    - SAC
    - SAC with Automatically Adjusted Temperature
    - TD3
- Quick Summary
- References

# 挑战

- Sparse supervision
- Severe partial observability
- Sample efficiency
- .....

# New Ideas

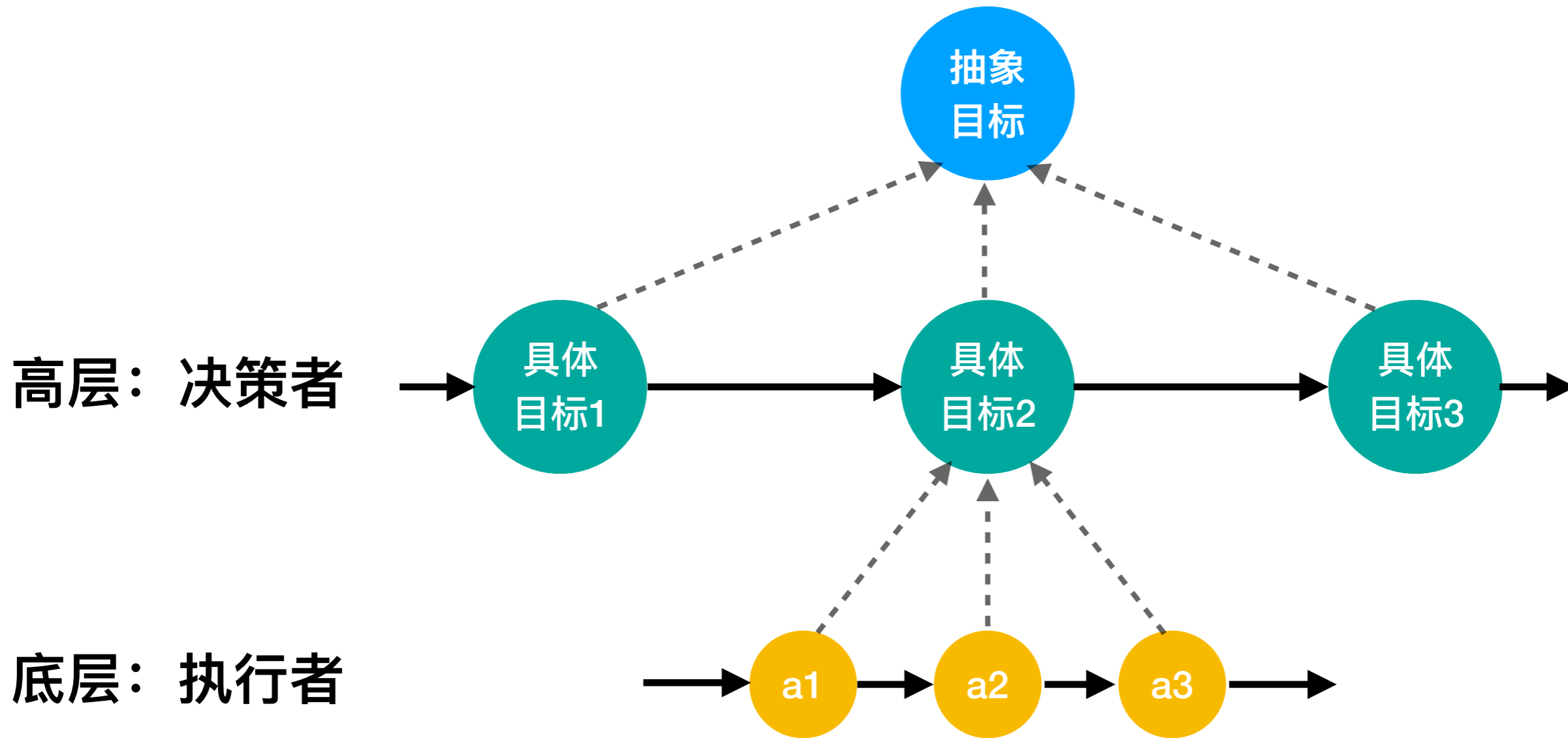
- 分层强化学习
- 记忆和注意力
- 世界模型和想象

# 分层强化学习：HRL

针对复杂任务、长程反馈——多层策略：

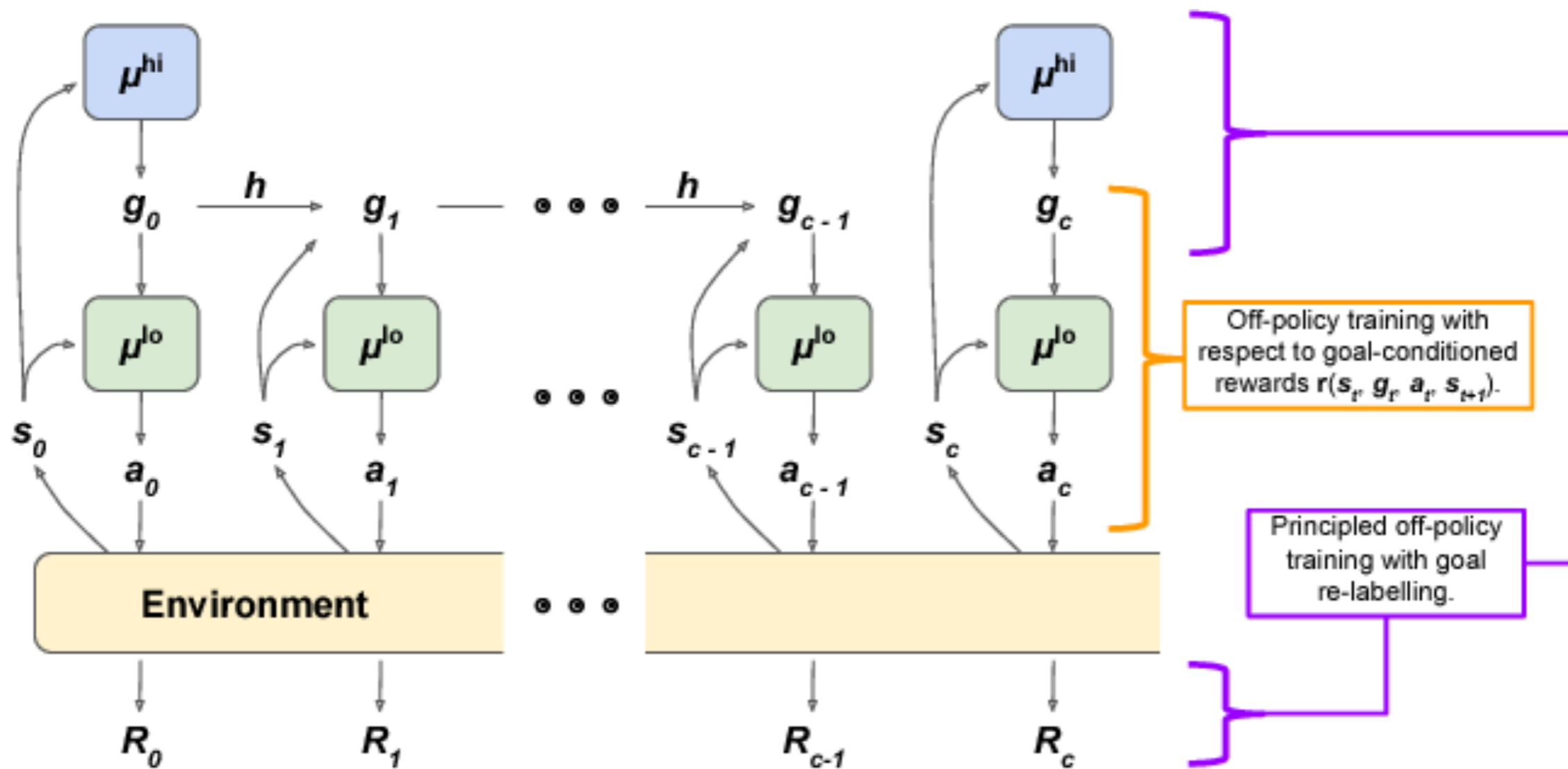
- 高层：分解高层目标为抽象的低层目标
- 底层：针对低层目标输出环境动作

# HRL

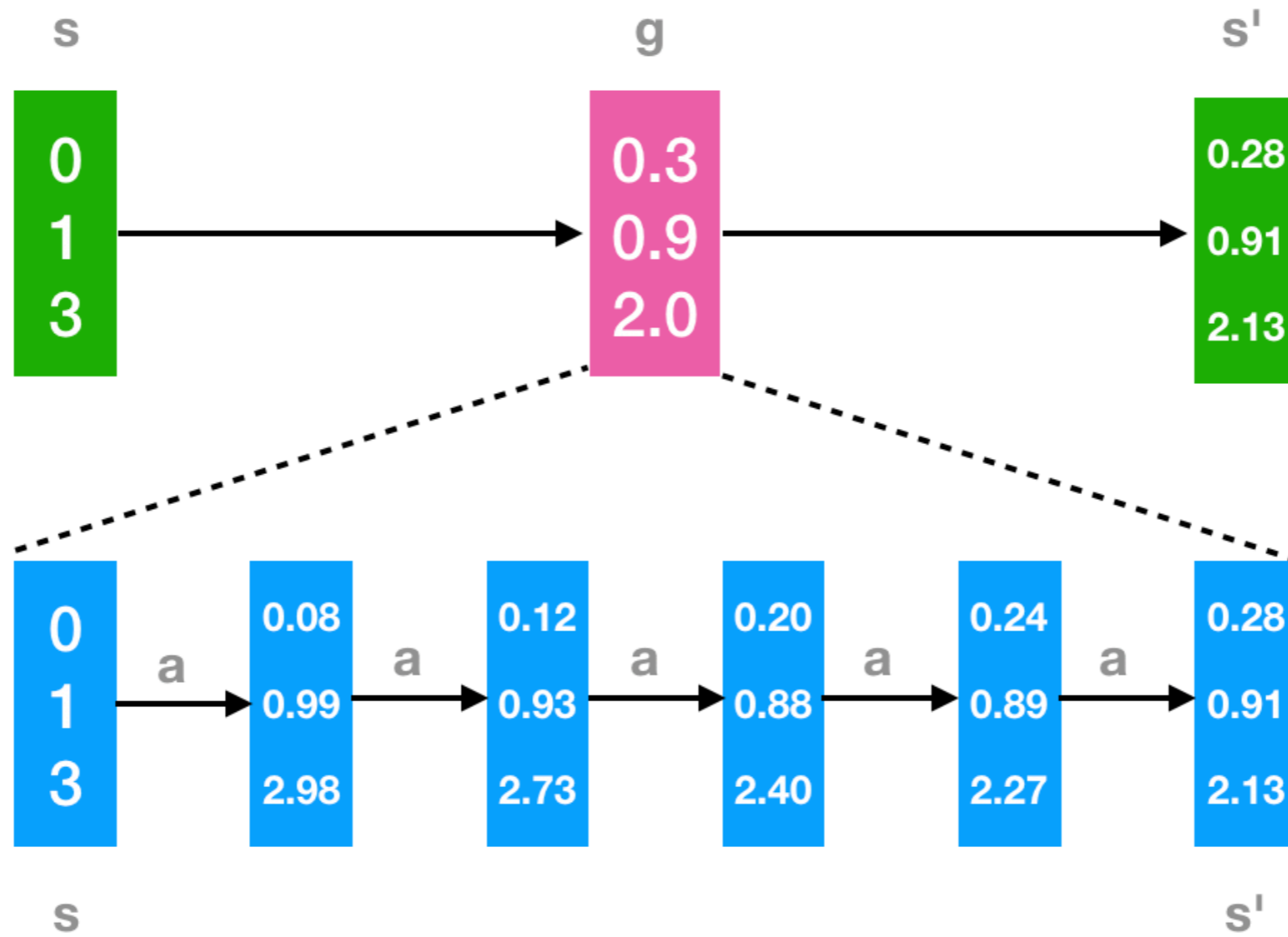




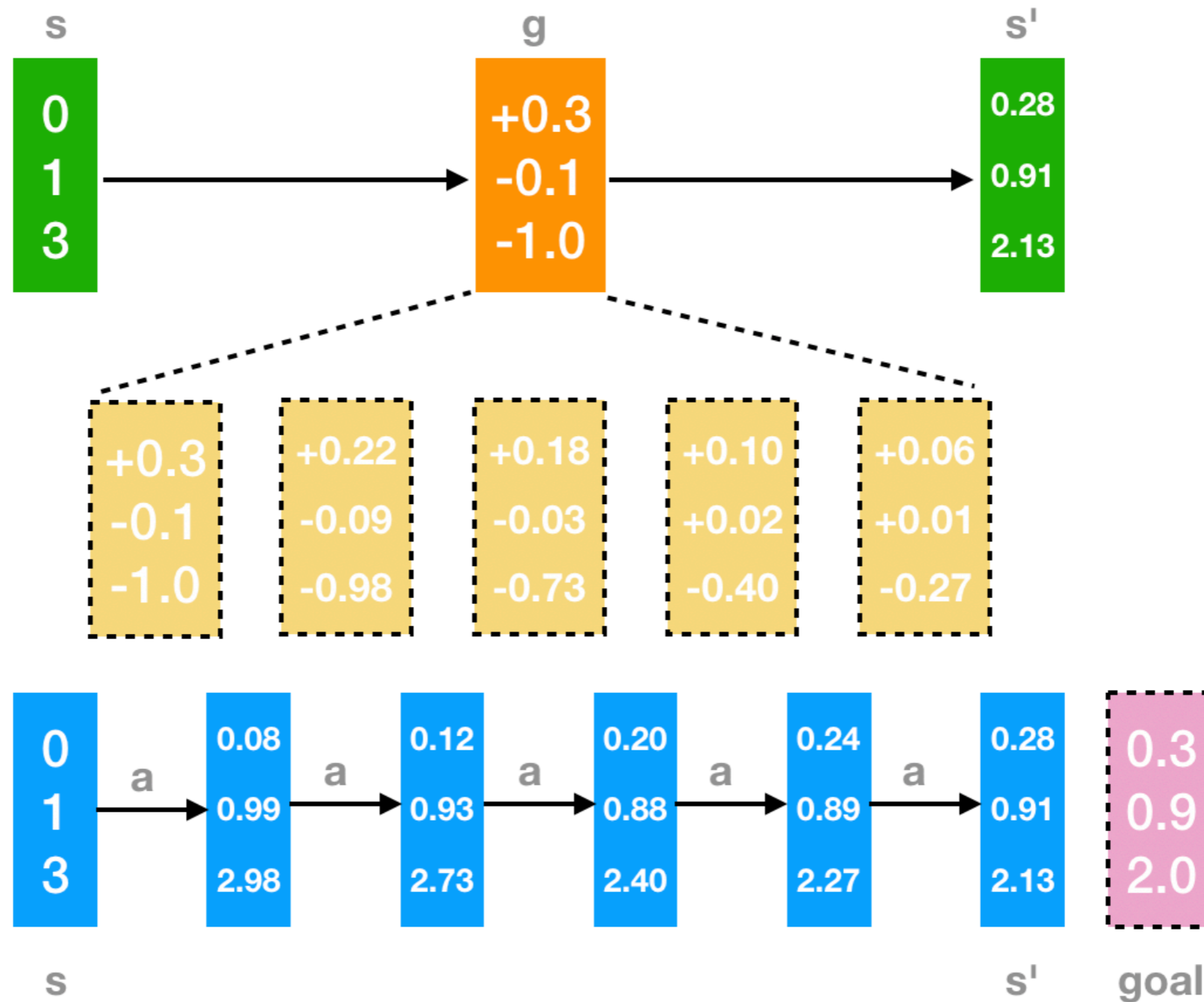
# 实现：HIRO



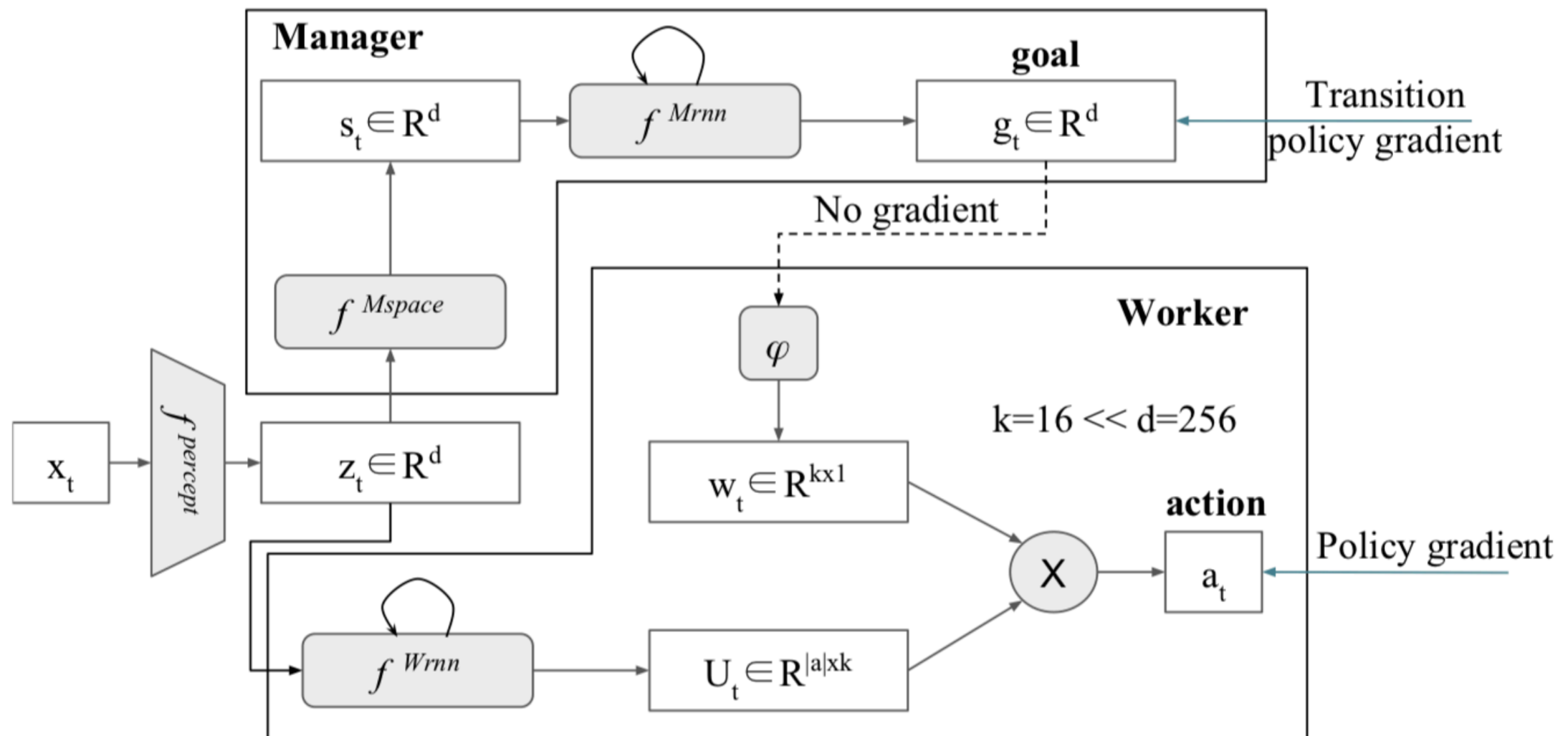
# HIRO



# HIRO



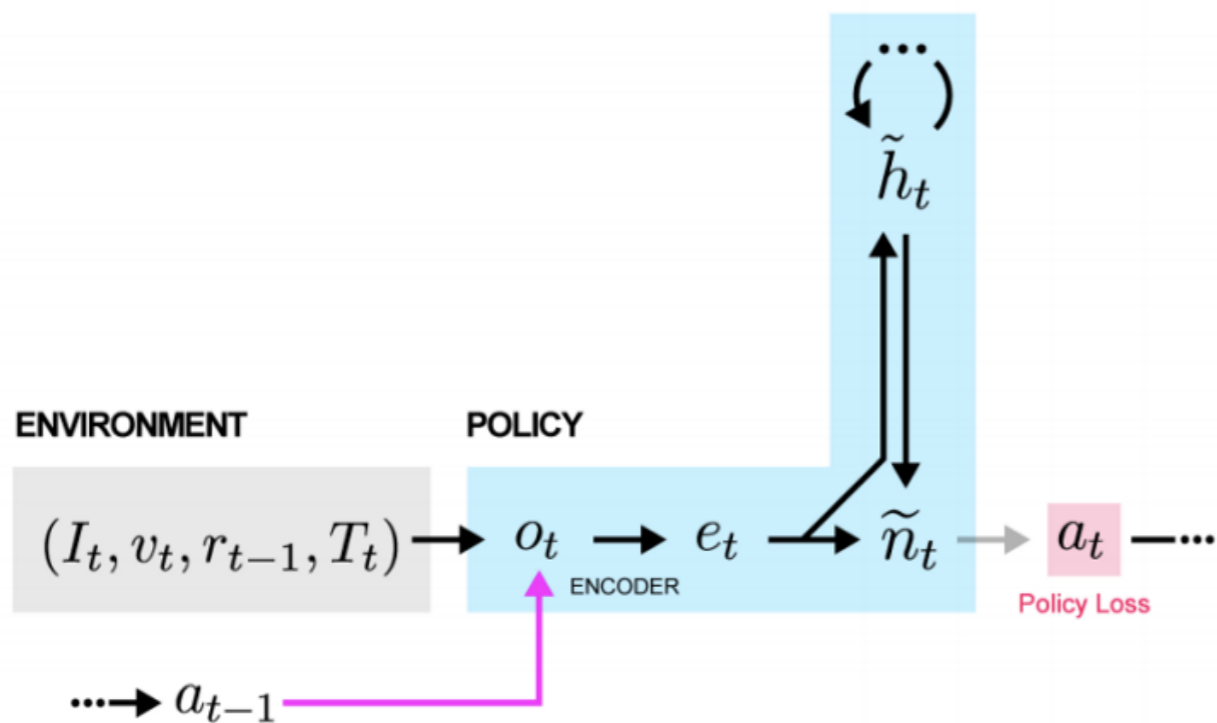
# FeUdal Networks



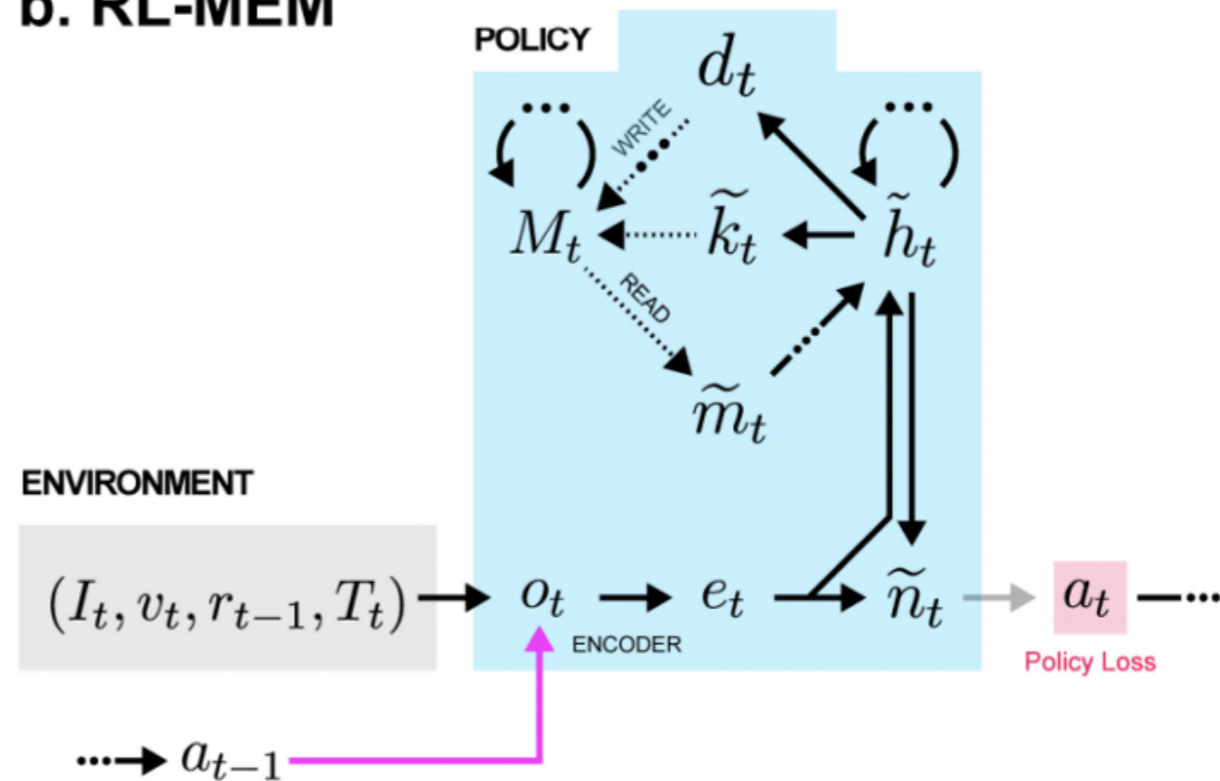
# 记忆和注意力

- 部分可见问题
- 将观察融入记忆
- 结合**观察和记忆**采取决策
- **相关记忆**：注意力机制

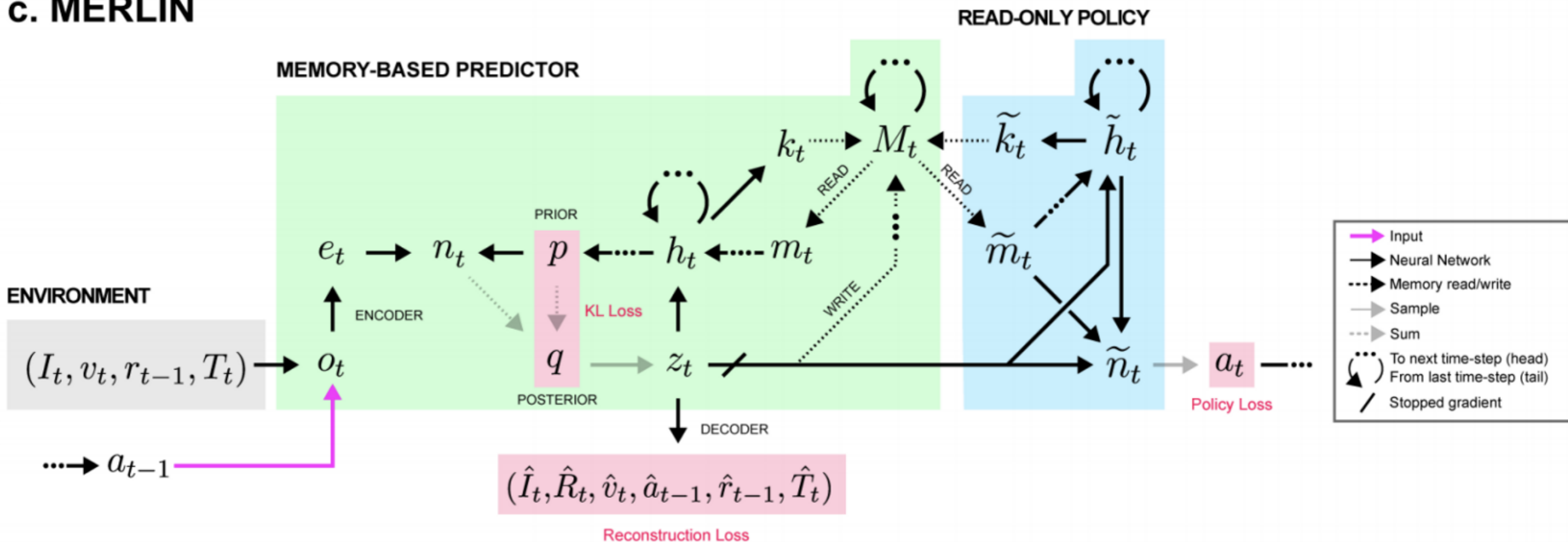
### a. RL-LSTM



### b. RL-MEM



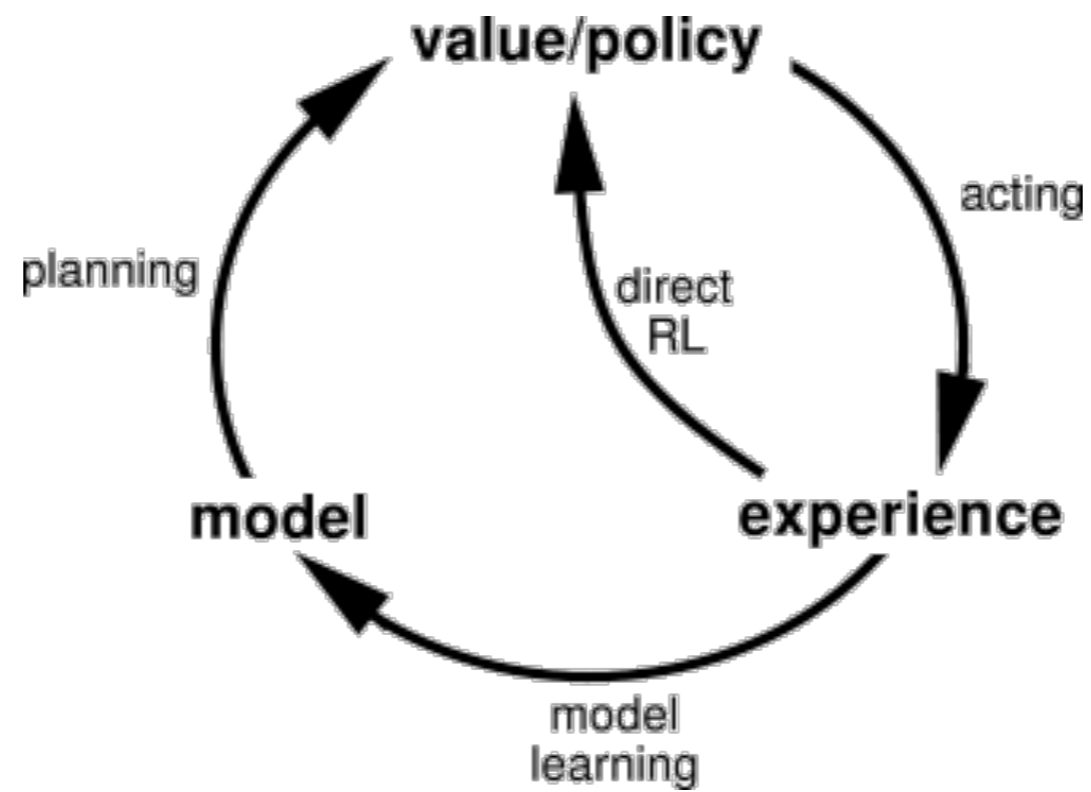
### c. MERLIN



# 世界模型和想象

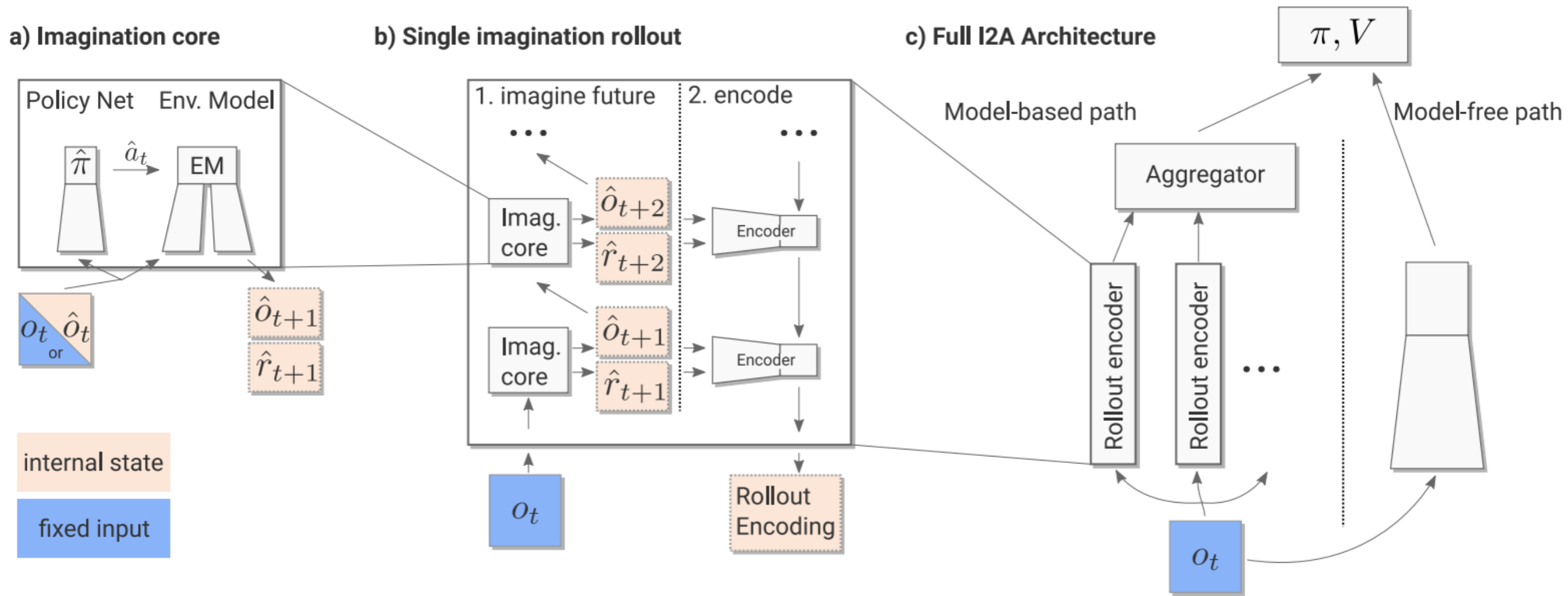
- 世界模型：对环境理解（解释？）
- 想象：世界模型上的规划，improve sample efficiency
- 辅助决策

# Framework





# I2A



# Reference

- <https://towardsdatascience.com/advanced-reinforcement-learning-6d769f529eb3>
- <http://karpathy.github.io/2016/05/31/rl/>
- <https://arxiv.org/abs/1703.01161>
- <https://arxiv.org/pdf/1805.08296.pdf>
- <https://arxiv.org/pdf/1803.10760.pdf>
- <https://arxiv.org/abs/1707.06203>
- <https://deepmind.com/blog/agents-imagine-and-plan/>