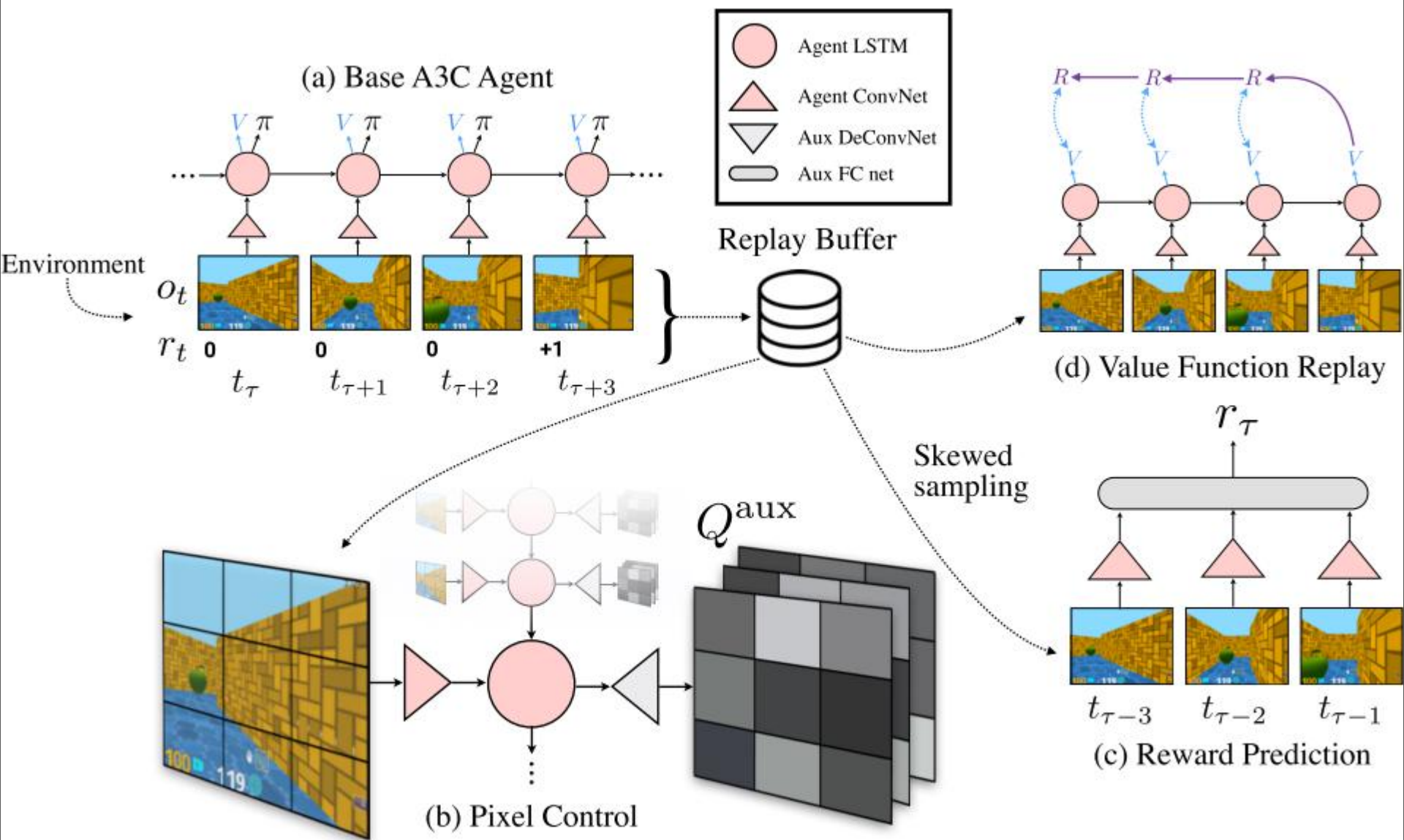


UNREAL

presented by Jason TOKO

算法介绍

- 算法全称是：UNsupervised REinforcement and Auxiliary Learning(UNREAL)
- 核心思想：UNREAL算法在A3C算法的基础上进行改善，通过在训练A3C的同时，训练多个辅助任务（**AUXILIARY TASKS**）来改进算法性能
- 算法背景：Asynchronous Advantage Actor-Critic(A3C) 和 Long Short-Term Memory (LSTM)



辅助任务

- 1、辅助控制任务

- 优化目标:

$$\arg \max_{\theta} \mathbb{E}_{\pi} [R_{1:\infty}] + \lambda_c \sum_{c \in \mathcal{C}} \mathbb{E}_{\pi_c} [R_{1:\infty}^{(c)}],$$

- 其中 $R_{t:t+n}^{(c)} = \sum_{k=1}^n \gamma^k r_t^{(c)}$

- n-step Q learning的损失函数:

$$\mathcal{L}_Q^{(c)} = \mathbb{E} \left[\left(R_{t:t+n} + \gamma^n \max_{a'} Q^{(c)}(s', a', \theta^-) - Q^{(c)}(s, a, \theta) \right)^2 \right]$$

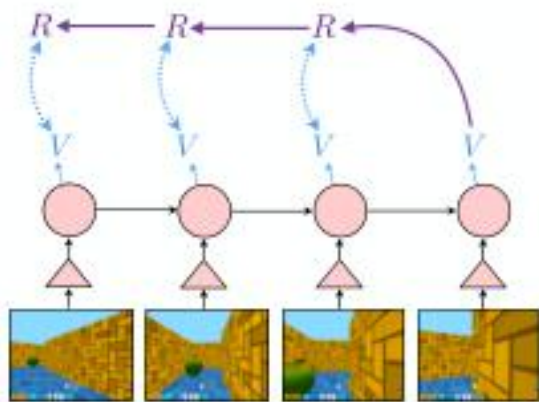
- 两种辅助控制任务

- 像素控制 (pixel control)
 - 特征控制 (feature control)

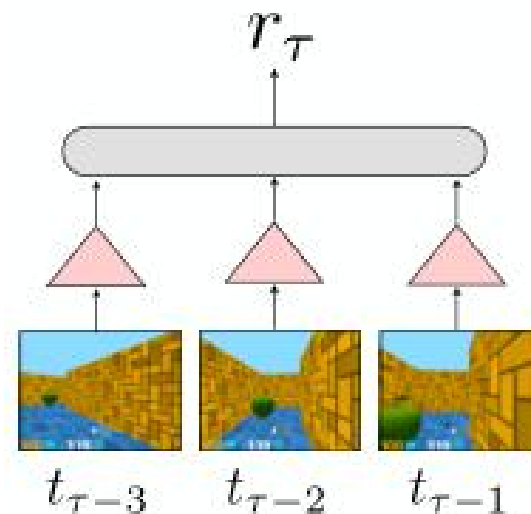
辅助任务

- 2、奖励预测（Reward Prediction）
 - 目的：消除奖励的稀疏性辅助学习，同时不引入偏差
 - 做法：通过序列 $S_\tau = (s_{\tau-k}, s_{\tau-k+1}, \dots, s_{\tau-1})$ 预测奖励 r_τ

- 3、经验回放（Experience Replay）



(d) Value Function Replay



(c) Reward Prediction

- UNREAL 的损失函数:

$$\mathcal{L}_{UNREAL}(\theta) = \mathcal{L}_{A3C} + \lambda_{VR} \mathcal{L}_{VR} + \lambda_{PC} \sum_c \mathcal{L}_Q^{(c)} + \lambda_{RP} \mathcal{L}_{RP}$$

- 算法效果:

演示

