

10/04/24

Lecture 06 → Recent advances & miscellaneous

① RL as a sequence prediction problem

↳ $s_t, a_t, r_t, s_{t+1}, \dots$

↳ Attraction of using large supervised models

② RL as a search heuristic for hard optimisation problems

↳ solve specific instances from scratch

③ RL in multi-agent settings

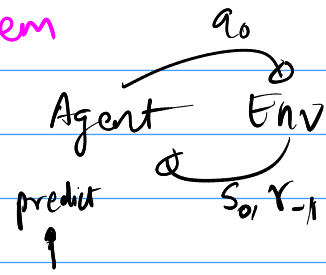
↳ Apart from adversarial ones such as games

④ Improving RL convergence

↳ reward shaping

↳ advanced sampling methods

RL as a sequence prediction problem



Upside-down RL

consider trajectory: $s_0, a_0, r_0, s_1, a_1, r_1, s_2, a_2, \dots$

Traditional RL converts this to an optimisation problem:

$$a_0 = \max_a E[G_0 | s_0]$$

$$[s_0, a_0, r_0, \dots, s_{t-1}, a_{t-1}, r_{t-1}, s_t] [r] [R_T] \rightarrow a_t$$

$r_0 + \gamma r_1 + \dots + \gamma^n [\text{terminal reward}]$

Challenges \rightarrow ① scalar reward, minimal feedback signal
(no error is available)

② Necessity of discount factor for convergence, especially in non-terminating problems

③ Sample efficiency: n step episode gives n samples

④ Difficulty of working with traditional deep architectures

⑤ Generalisation to out-of-distribution areas, rather than learning a narrow behaviour distribution

⑥ credit assignment in case of delayed rewards

Solution: Can we think of the trajectory (above) as a sequence, and use an LLM (essentially causal transformer) to generate actions?

① Conditioning on full trajectory

② Supervised training

③ Credit assignment with self-attention (causal \rightarrow previous tokens only)

④ approx $\frac{n^2}{2}$ samples per episode

⑤ Highly capable pretrained architectures

Diversion: Quick recap of transformers

Assume input encoding $x \rightarrow p$ tokens of n size each $\Rightarrow n \times p$

$$\begin{aligned} \text{Query } Q &= x W_Q \\ \text{Key } K &= x W_K \\ \text{Values } V &= x W_V \end{aligned} \quad \left. \begin{array}{l} \text{matrices} \\ \text{size } n \times m \end{array} \right\}$$

$$\text{attention} = \text{softmax}\left(\frac{QK^T}{\sqrt{m}}\right)$$

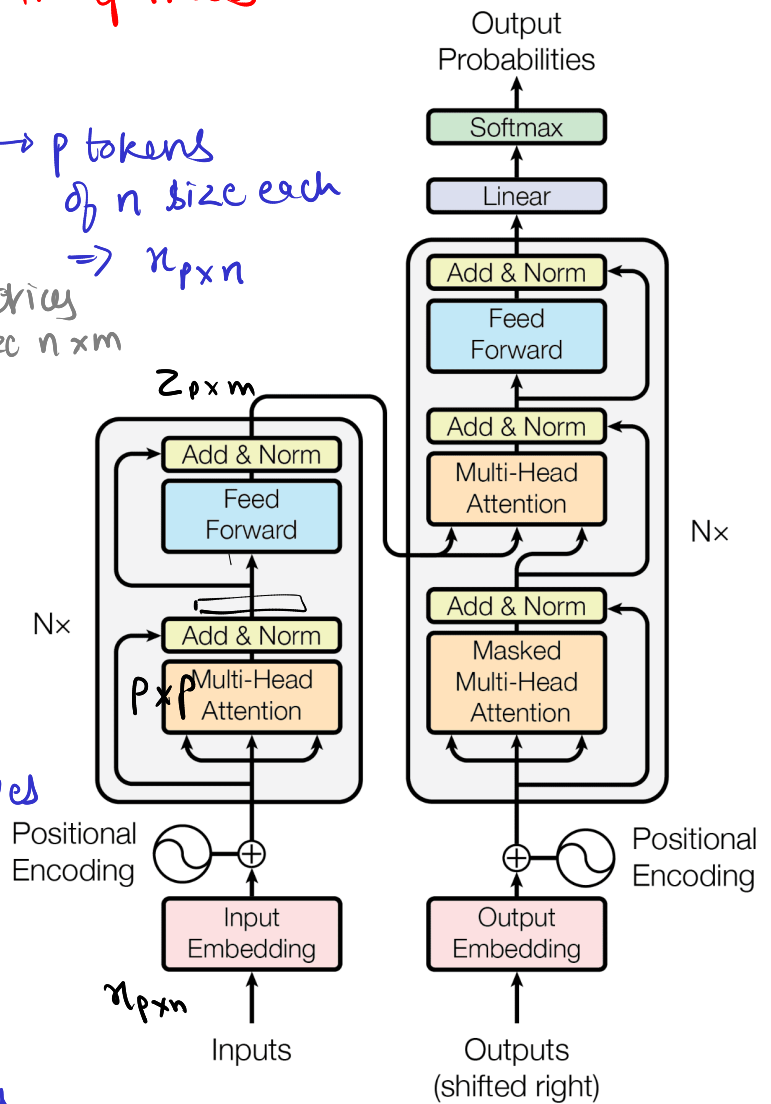
note: attention size $p \times p$

output $z = \text{attention} \times \text{values}$

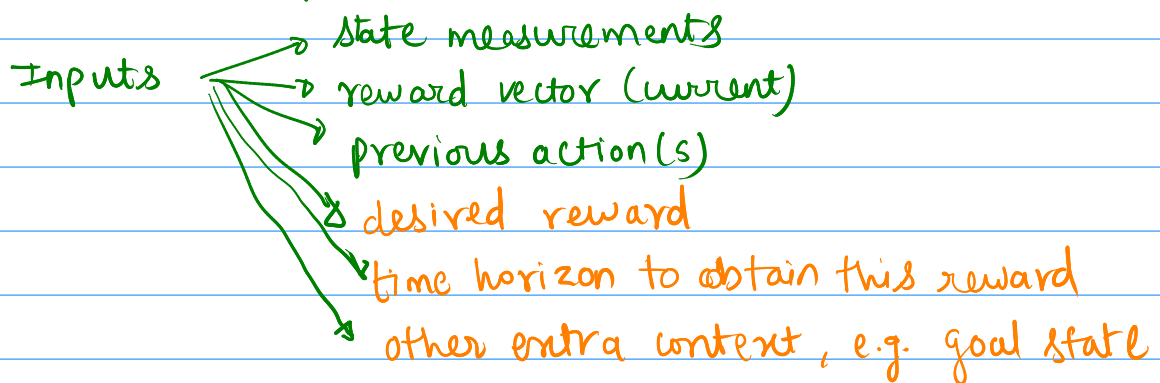
$p \times m$

z is passed to next block

→ Note number of samples stays same, only embedding size may be changed



So how does upside-down RL work?



upside down RL (UDRL) paper uses RNN instead of xformers

Phase 1: Experience collection

inputs \rightarrow action probabilities \rightarrow action \rightarrow new inputs

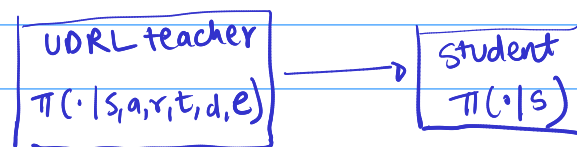
Phase 2: training

2a) Experience re-labeling \rightarrow replace by actual observations instead of desired ones

2b) Trajectory subsampling \rightarrow pick every pair of time instants and compute Σr , Δt

2c) supervised training of policy output

final point on running in real world without reward conditioning



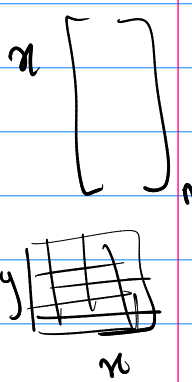
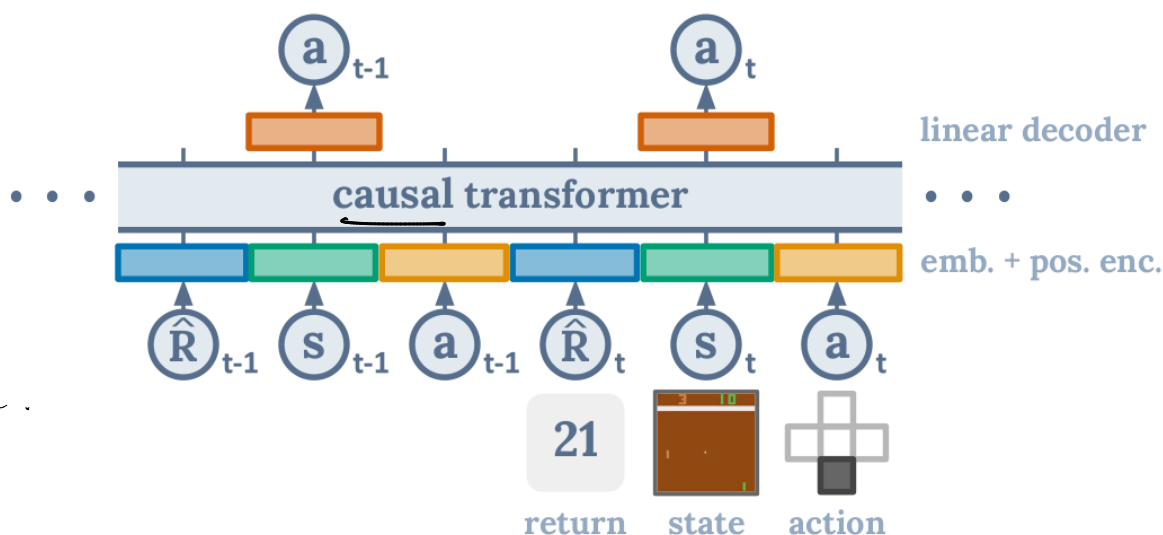
use the student network (trained in supervised fashion)

Use in learning by demonstration

- ① let robot take random actions
- ② get humans to imitate, record video
- ③ Train using UDRL with ② as "desired"
- ④ Take videos of actually useful demonstrations
- ⑤ Fine-tune using UDRL

Decision transformer & Trajectory transformer

Think of r, s, a, r, s, a, \dots as a sequence, and use a transformer to produce the next token.



Decision transformer

Set returns target $\hat{R}_t = \sum_{t'=t}^T r_{t'}$

Trajectory $\tau = (\hat{R}_1, s_1, a_1, \dots)$

K time steps = $3K$ tokens

s, a, r have separate embeddings for matching dimensions

In-context learning
Algo distillation

Output = action embedding, to be decoded to action

Trajectory transformer

Use step rewards

Trajectory $\tau = (s_t^1, s_t^2, \dots, a_t^1, a_t^2, \dots, r_{t,1}, r_{t,2}, \dots)$

Each token is one (scalar) dimension of state/action or reward, discretised
 \Rightarrow No embedding needed apart from normalisation

output \rightarrow similar to NLP models, generate T length trajectories by sampling and choose action corresponding to most likely trajectory (beam search)

Discovering matrix multiplication using RL

conceptual steps:

① Simplest matrix multiplication has 3 'for' loops, hence around k^3 scalar multiplication operations

② Multiplication operations more expensive than addition operations

③ # multiplications can be reduced by clever & selective multiplications of elements, followed by simple addition/subtraction

↳ generic for any $(n \times m) \times (m \times p)$ multiplication
↳ every set (n, m, p) has unique formula

we want to discover this

④ The multiplication of matrices, $C_{n \times p} = A_{n \times m} \times B_{m \times p}$ can be considered to be a tensor of dimension $n \times m \times p$

⑤ Just like a matrix of rank R has R independent eigenvectors, this tensor $\mathcal{T}_{n,m,p}$ also has R independent terms. For a square operation,

$$\mathcal{T}_{n,n,p} := \mathcal{T}_n = \sum_{r=1}^R U^{(r)} \otimes V^{(r)} \otimes W^{(r)}$$

outer product

each a vector of length n^2

⑥ Tensor $R \leq n^2$, since trivially you could have individual $a_i b_j$ terms (total n^2) combined to give elements of C

⑦ Goal of RL is to select $[U^{(t)}, V^{(t)}, W^{(t)}]$ in each step, to go from \mathcal{T}_n to 0 in as few steps as possible.

Multi-Agent RL

Multi-Agent Actor-Critic for Mixed Cooperative-Competitive Environments

Ryan Lowe*
McGill University
OpenAI

Yi Wu*
UC Berkeley

Aviv Tamar
UC Berkeley

Jean Harb
McGill University
OpenAI

Pieter Abbeel
UC Berkeley
OpenAI

Igor Mordatch
OpenAI

Multi-agent RL similar to stochastic games

For N agents, MDP $\rightarrow (S, \{a_n\}, \{r_n\}, P, \gamma)$

$P: S \times A^N \rightarrow S$ aggregate state transition

$\{a_n\}$ individual actions

$\{r_n\}$ local rewards

Markov Games In this work, we consider a multi-agent extension of Markov decision processes (MDPs) called partially observable Markov games [19]. A Markov game for N agents is defined by a set of states \mathcal{S} describing the possible configurations of all agents, a set of actions $\mathcal{A}_1, \dots, \mathcal{A}_N$ and a set of observations $\mathcal{O}_1, \dots, \mathcal{O}_N$ for each agent. To choose actions, each agent i uses a stochastic policy $\pi_{\theta_i} : \mathcal{O}_i \times \mathcal{A}_i \mapsto [0, 1]$, which produces the next state according to the state transition function $\mathcal{T} : \mathcal{S} \times \mathcal{A}_1 \times \dots \times \mathcal{A}_N \mapsto \mathcal{S}$.^[2] Each agent i obtains rewards as a function of the state and agent's action $r_i : \mathcal{S} \times \mathcal{A}_i \mapsto \mathbb{R}$, and receives a private observation correlated with the state $o_i : \mathcal{S} \mapsto \mathcal{O}_i$. The initial states are determined by a distribution $\rho : \mathcal{S} \mapsto [0, 1]$. Each agent i aims to maximize its own total expected return $R_i = \sum_{t=0}^T \gamma^t r_i^t$ where γ is a discount factor and T is the time horizon.

Individual returns

local observation

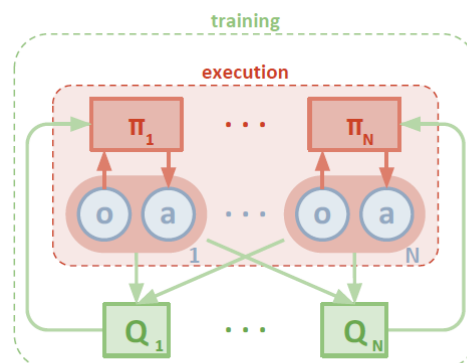


Figure 1: Overview of our multi-agent decentralized actor, centralized critic approach.

Gradient for agent i :

$$\nabla_{\theta_i} J(\theta_i) = \mathbb{E}_{s \sim p^\mu, a_i \sim \pi_i} [\nabla_{\theta_i} \log \pi_i(a_i | o_i) Q_i^\pi(\mathbf{x}, a_1, \dots, a_N)].$$

policy parameters for agent i , resulting in policy $\pi_i(\cdot | o_i)$

actions depend only on o_i

true state

Q values have access to all agent actions and true state

However note that each agent i builds its own Q_i , allowing us to handle non-cooperative / adversarial settings

In multi-agent DDPG (MA-DDPG), replace $\pi_i(\cdot | o_i)$ by $\mu_i(o_i)$

$$\mathcal{L}(\theta_i) = \mathbb{E}_{\mathbf{x}, a, r, \mathbf{x}'} [(Q_i^\mu(\mathbf{x}, a_1, \dots, a_N) - y)^2], \quad y = r_i + \gamma Q_i^{\mu'}(\mathbf{x}', a'_1, \dots, a'_N) \big|_{a'_j = \mu'_j(o_j)}, \quad (6)$$

where $\mu' = \{\mu_{\theta'_1}, \dots, \mu_{\theta'_N}\}$ is the set of target policies with delayed parameters θ'_i . As shown in

This is a problem, since \mathcal{L} depends on knowing all agent policies

To remove the assumption of knowing other agents' policies, as required in Eq. [6], each agent i can additionally maintain an approximation $\hat{\mu}_{\phi_i^j}$ (where ϕ are the parameters of the approximation; henceforth $\hat{\mu}_i^j$) to the true policy of agent j , μ_j . This approximate policy is learned by maximizing the log probability of agent j 's actions, with an entropy regularizer:

$$\mathcal{L}(\phi_i^j) = -\mathbb{E}_{o_j, a_j} [\log \hat{\mu}_i^j(a_j | o_j) + \lambda H(\hat{\mu}_i^j)], \quad (7)$$

where H is the entropy of the policy distribution. With the approximate policies, y in Eq. [6] can be replaced by an approximate value \hat{y} calculated as follows:

$$\hat{y} = r_i + \gamma Q_i^{\mu'}(\mathbf{x}', \hat{\mu}_i^{j1}(o_1), \dots, \mu_i'(o_i), \dots, \hat{\mu}_i^{jN}(o_N)), \quad (8)$$

Policy invariance under reward transformations

Theory and application to reward shaping

Reward shaping

In general, changing reward \Rightarrow changing optimal policy
 \downarrow
 not desirable

Notable exception: Potential based rewards

$$r_{\text{new}} = ar_{\text{old}} + b$$

where $a > 0$

timal policy. It is shown that, besides the positive linear transformation familiar from utility theory, one can add a reward for transitions between states that is expressible as the difference in value of an arbitrary potential function applied to those states. Further-

than running our reinforcement learning algorithm on $M = (S, A, T, \gamma, R)$, we will run it on some transformed MDP $M' = (S, A, T, \gamma, R')$, where $R' = R + F$ is the reward function in the transformed MDP, and $F : S \times A \times S \mapsto \mathbb{R}$ is a bounded real-valued function called the **shaping reward function**. (Simi-

For any fixed MDP and assuming additive, memory-less shaping reward functions, this $R' = R + F$ is the most general possible form of shaping rewards.³ Moreover, they cover a fairly large range of possible shaping rewards one might come up with. For example, to encourage moving towards a goal, a shaping-reward function that one might choose is $F(s, a, s') = r$ whenever s' is closer (in whatever appropriate sense) to the goal than s , and $F(s, a, s') = 0$ otherwise, where r is some positive reward. Or, to encourage taking action a_1 in some set of states S_0 , one might set $F(s, a, s') = r$ whenever $a = a_1, s \in S_0$, and $F(s, a, s') = 0$ otherwise.

Theorem 1 Let any S, A, γ , and any shaping reward function $F : S \times A \times S \mapsto \mathbb{R}$ be given. We say F is a **potential-based** shaping function if there exists a real-valued function $\Phi : S \mapsto \mathbb{R}$ such that for all $s \in S - \{s_0\}, a \in A, s' \in S$,

$$F(s, a, s') = \gamma \Phi(s') - \Phi(s), \quad (2)$$

(where $S - \{s_0\} = S$ if $\gamma < 1$). Then, that F is a potential-based shaping function is a necessary and sufficient condition for it to guarantee consistency with the optimal policy (when learning from $M' = (S, A, T, \gamma, R + F)$ rather than from $M = (S, A, T, \gamma, R)$), in the following sense:

- (Sufficiency) If F is a potential-based shaping function, then every optimal policy in M' will also be an optimal policy in M (and vice versa).
- (Necessity) If F is not a potential-based shaping function (e.g. no such Φ exists satisfying Equation (2)), then there exist (proper) transition functions T and a reward function $R : S \times A \mapsto \mathbb{R}$, such that no optimal policy in M' is optimal in M .

How? $Q^*(s, a) = R + \gamma \max_{a'} Q^*(s', a')$

$$\Rightarrow \underbrace{Q^*(s, a) - \Phi(s)}_{\hat{Q}^*(s, a)} = R + \underbrace{\gamma \Phi(s') - \Phi(s)}_F + \gamma \underbrace{[\max_{a'} Q(s', a') - \Phi(s')]}_{\hat{Q}^*(s', a')}$$

whatever policy
 maximises \hat{Q}^*
 must also
 maximise \hat{Q}^*

STAS: Spatial-Temporal Return Decomposition for Solving Sparse Rewards Problems in Multi-agent Reinforcement Learning

tionality to model complicated relations of the delayed global reward in the temporal dimension and suffer from inefficiencies. To tackle this, we introduce Spatial-Temporal Attention with Shapley (STAS), a novel method that learns credit assignment in both temporal and spatial dimensions. It first decomposes the global return back to each time step, then utilizes the **Shapley Value** to redistribute the individual payoff from the decomposed global reward. To mitigate the com-

$R_{\tau} \rightarrow r_0, r_1, \dots, r_T$
time decomposition
 \downarrow
 $r_{i,0}, r_{i,1}, \dots$
Agent level decomposition

the marginal contribution of player i in C can be defined as:

$$v_i(C) = v(C \cup \{i\}) - v(C). \quad (1)$$

The **Shapley Value** of player i can be computed from the marginal contribution of player i in all subsets of N :

$$\Phi_v(i) = \sum_{C \subseteq N \setminus \{i\}} \frac{|C|!(|N| - |C| - 1)!}{|N|!} v_i(C). \quad (2)$$

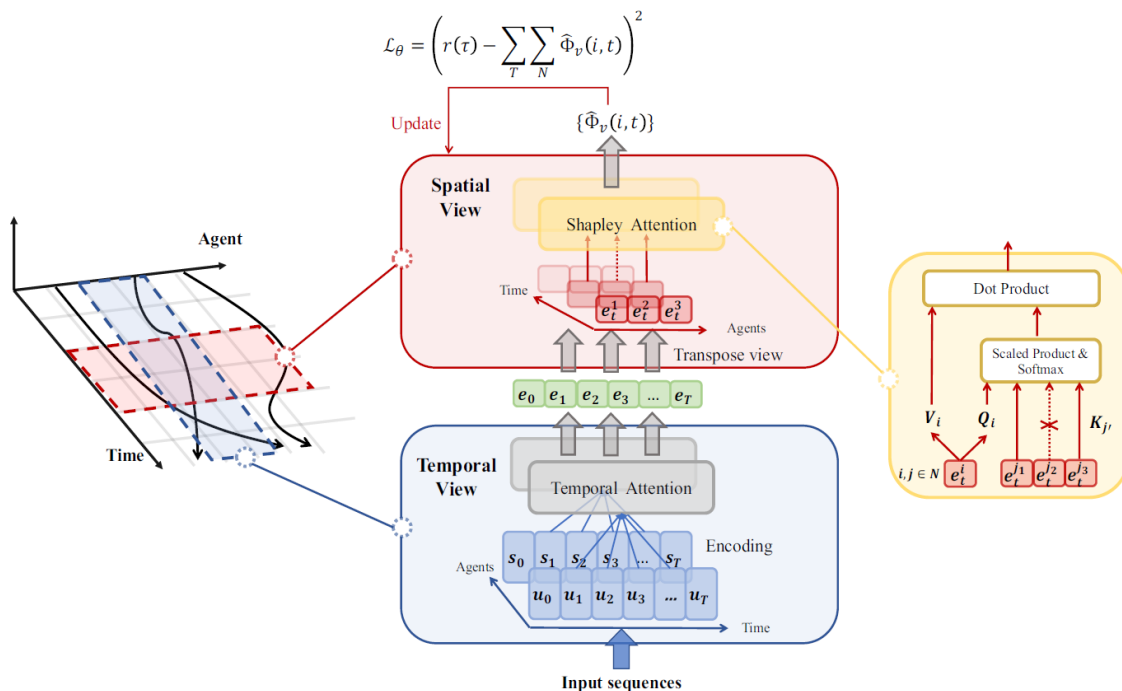
True return

$$r_{\text{ep}}(\tau) \approx \hat{r}_{\text{ep}}(\tau) = \sum_{t=0}^{T-1} \hat{r}(s_t, u_t). \quad r_{\text{ep}}(\tau) = \sum_{t=0}^{T-1} r(s_t, u_t) = \sum_{t=0}^{T-1} \sum_{i=1}^N \Phi_v(i, t).$$

shaped return

Assumed to be sum decomposable

further decompose to agents



Using Contrastive Samples for Identifying and Leveraging Possible Causal Relationships in Reinforcement Learning

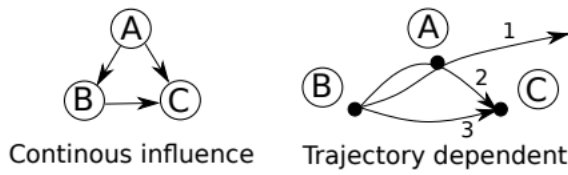


Fig. 1. Difference between a one-step inference problem and a sequential decision-making problem. In the latter case, the influence of A on the long term rewards is more subtle.

Illustration with 5x5 grid: experiments on larger sizes

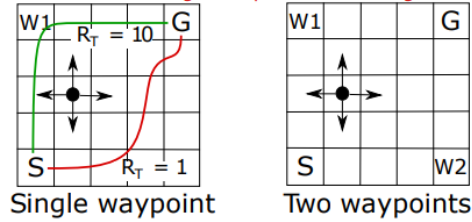


Fig. 2. Schematic of the environments used in this work. Waypoints W1 and W2 influence the terminal rewards if they are visited on the way to G. Size of the environment is illustrative; experiments use 8x8 and 12x12 grids.

Create two buffers

regular
memory

contrastive
memory

$$\text{If } \Delta(s_t, s_{t+1}) \geq \delta$$

$$\text{then } (s_t, a_t, r_t, s_{t+1}) \rightarrow \text{CER}$$

Additionally,

$$(s_t, a_t^* \neq a_t, r_t^*, s_{t+1}^*) \rightarrow \text{CER}$$

for training,
sample fixed
proportion from
regular buffer
and CER