

Attention-Based Recurrence for Multi-Agent Reinforcement Learning under Stochastic Partial Observability

Thomy Phan^{1,2}, Fabian Ritz², Philipp Altmann², Maximilian Zorn², Jonas Nüßlein²,
Michael Kölle², Thomas Gabor², Claudia Linnhoff-Popien²

¹University of Southern California,

²LMU Munich



Multi-Agent Reinforcement Learning in Dec-POMDPs

Most state-of-the-art MARL approaches assume deterministic observability and focus on

$$Q_{MPD}^*(s, a) = \mathcal{R}(s, a) + \gamma \max_{a'} Q_{MPD}^*(s', a')$$

However, the true optimal value function Q^* in Dec-POMDPs is actually defined by

$$Q^*(\tau, a) = b(s|\tau) \left(\mathcal{R}(s, a) + \gamma \sum_{s'} \sum_{z'} \mathcal{P}(s'|s, a) \Omega(z'|\mathbf{a}, s') Q^*(\tau', \pi^*(\tau')) \right)$$

Popular benchmarks do not exhibit stochastic partial observability unlike general Dec-POMDPs



SMAC



Multi-Agent Recurrence

$$\pi^* = \operatorname{argmax}_{\pi} \sum_{t=0}^{T-1} \sum_{\tau_t \in (\mathcal{Z}^N \times \mathcal{A})^t} \underbrace{C^{\pi}(\tau_t) \mathbf{P}^{\pi}(\tau_t | b_0) Q^*(\cdot)}_{\text{multi-agent recurrence}}$$

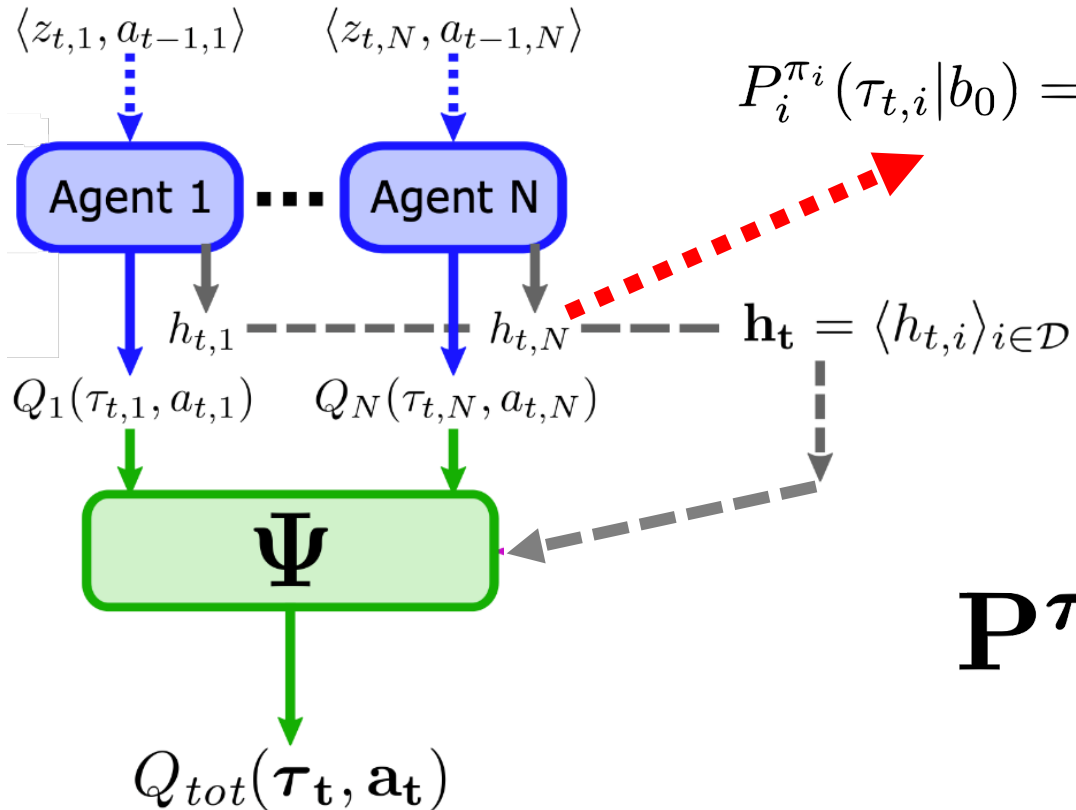
← multi-agent recurrence

$$\begin{aligned} \mathbf{P}^{\pi}(\tau_t | b_0) &= \mathbf{P}(\mathbf{z}_0 | b_0) \prod_{c=1}^t \mathbf{P}(\mathbf{z}_c | \tau_{c-1}, \pi) \\ &= \mathbf{P}(\mathbf{z}_0 | b_0) \prod_{c=1}^t \sum_{s_c \in \mathcal{S}} \sum_{s_{c-1} \in \mathcal{S}} \mathcal{T}(\cdot) \Omega(\cdot) \end{aligned}$$

AERIAL: Attention-based Embeddings of Recurrence In multi-Agent Learning

Value factorization approach considering ...

- memory representations of agents



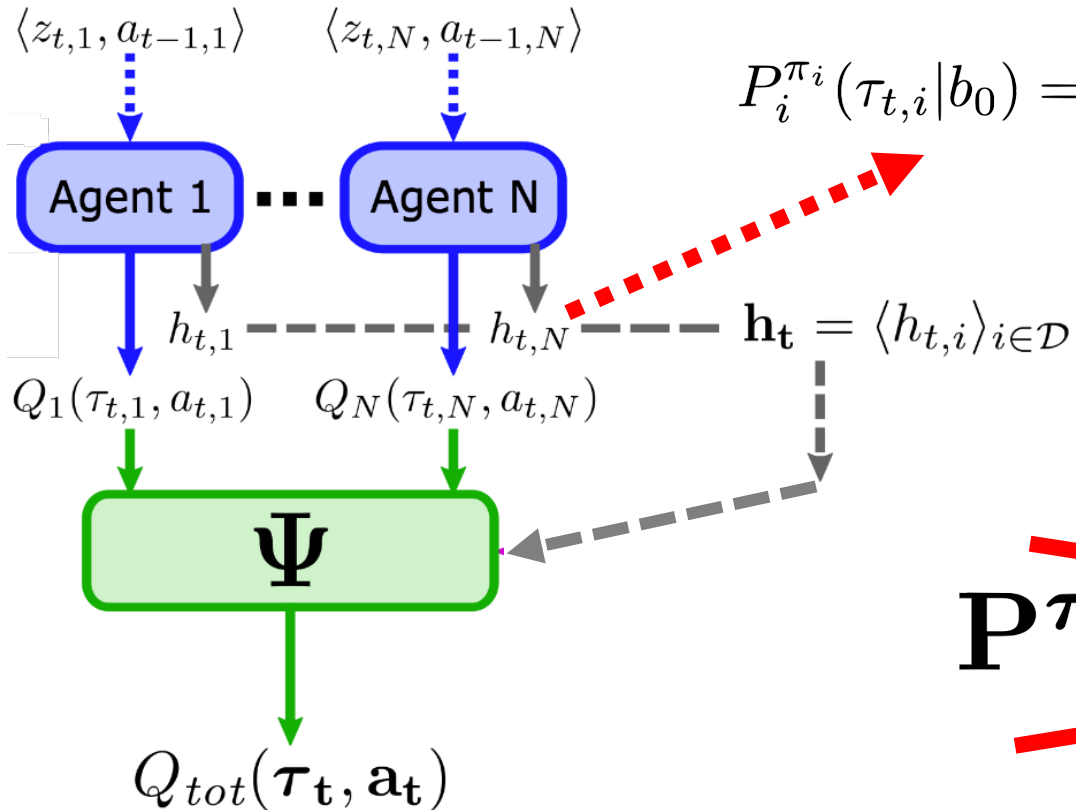
$$P_i^{\pi_i}(\tau_{t,i} | b_0) = P_i(z_{0,i} | b_0) \prod_{c=1}^t P_i(z_{c,i} | \tau_{c-1,i}, \pi_i)$$

$$\mathbf{P}^{\pi}(\tau_t | b_0) = \prod_{i=1}^N P_i^{\pi_i}(\tau_{t,i} | b_0)$$

AERIAL: Attention-based Embeddings of Recurrence In multi-Agent Learning

Value factorization approach considering ...

- memory representations of agents



$$P_i^{\pi_i}(\tau_{t,i} | b_0) = P_i(z_{0,i} | b_0) \prod_{c=1}^t P_i(z_{c,i} | \tau_{c-1,i}, \pi_i)$$

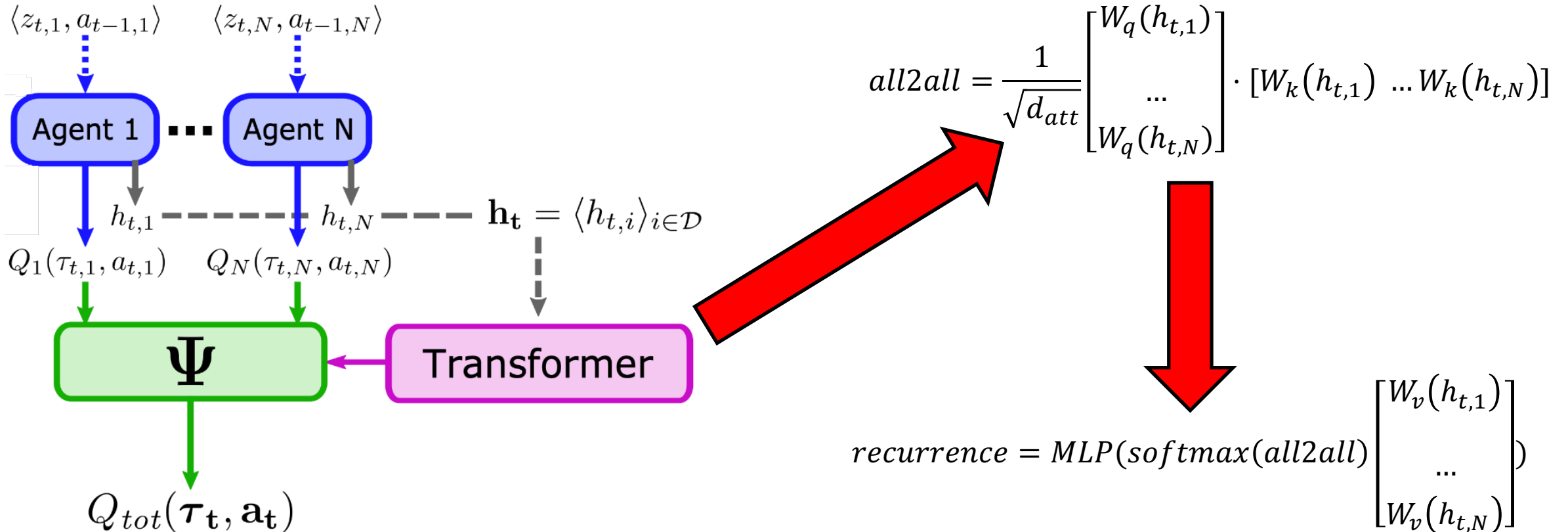
~~$$\mathbf{P}^{\pi}(\tau_t | b_0) = \prod_{i=1}^N P_i^{\pi_i}(\tau_{t,i} | b_0)$$~~

Individual recurrences not independent in general ☹️

AERIAL: Attention-based Embeddings of Recurrence In multi-Agent Learning

Value factorization approach considering ...

- memory representations of agents
- statistical dependence of these representations



MessySMAC

Modified variant of SMAC with ...

- observation stochasticity
- initialization stochasticity

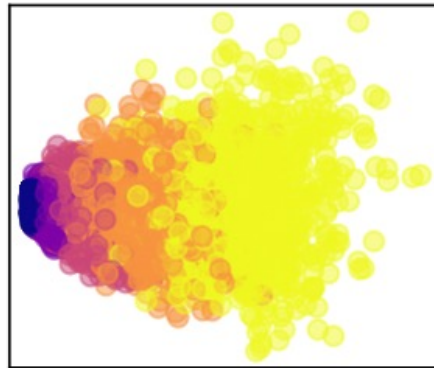
Code



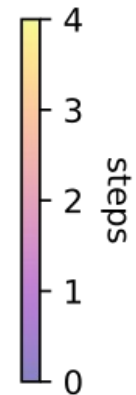
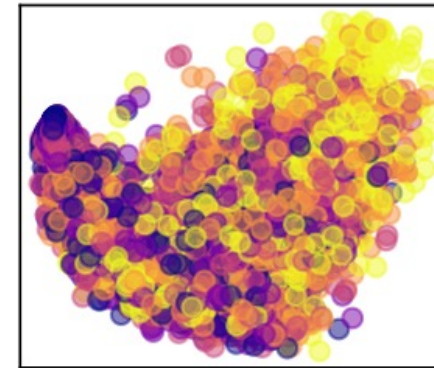
10m_vs_11m



K=0



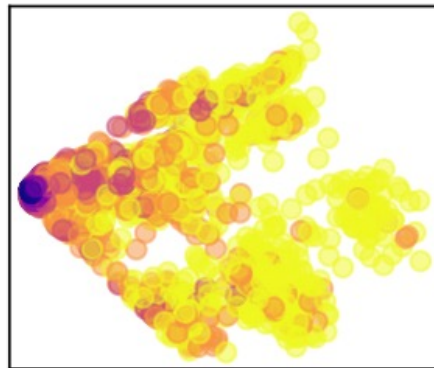
K=10



3s_vs_5z



K=0



K=10

