

Relational Reasoning via Set Transformers: Provable Efficiency and Applications to MARL

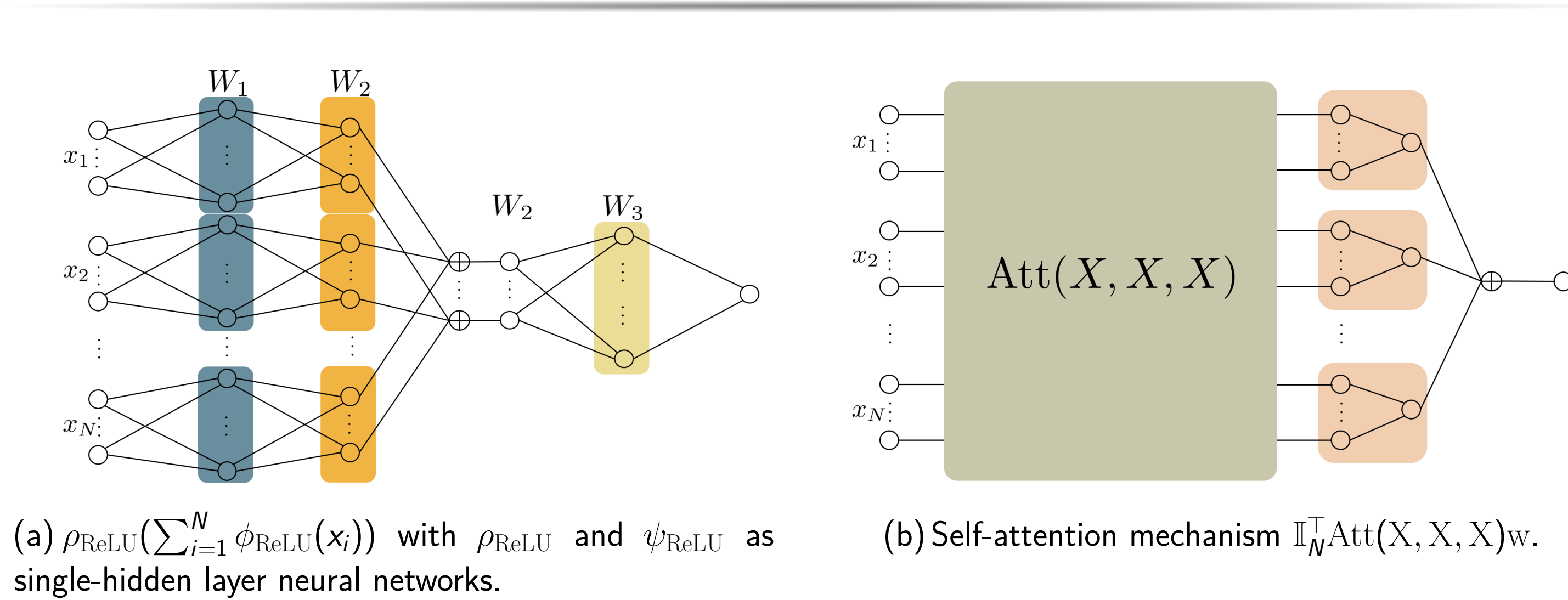
Fengzhuo Zhang[†], Boyi Liu^{*}, Kaixin Wang[†], Vincent Y. F. Tan[†], Zhuoran Yang[‡], Zhaoran Wang^{*}

[†]National University of Singapore, ^{*} Northwestern University, [‡] Yale University

Summary of Results

- MultiLayer Perceptrons (MLP) cannot approximate transformer unless the width is **exponential** in the input dimension of each channel.
- The generalization error bound of transformer function class is **independent** of the number of channels.
- The transformer helps to break “the curse of many agents”.

Superiority of Transformer Over the MLP in Terms of Relational Reasoning



- The function class of the **permutation invariant** MLP can be defined from deepset

$$\mathcal{N}(W) = \left\{ f : \mathbb{R}^{N \times d} \rightarrow \mathbb{R} \mid f(X) = \rho_{\text{ReLU}} \left(\sum_{i=1}^N \phi_{\text{ReLU}}(x_i) \right) \text{ with } \max_{i \in [3]} W_i \leq W \right\},$$

where ρ_{ReLU} and ϕ_{ReLU} as width-constrained ReLU networks with **maximal widths** W_1 and W_3 .

- The self-attention function class is

$$\mathcal{F} = \{ f : \mathbb{R}^{N \times d} \rightarrow \mathbb{R} \mid f(X) = \mathbb{I}_N^{\top} \text{Att}(X, X, X)w \text{ for some } w \in [0, 1]^d \}.$$

- Let $W^*(\xi, d, \mathcal{F})$ be the smallest width of the neural network such that

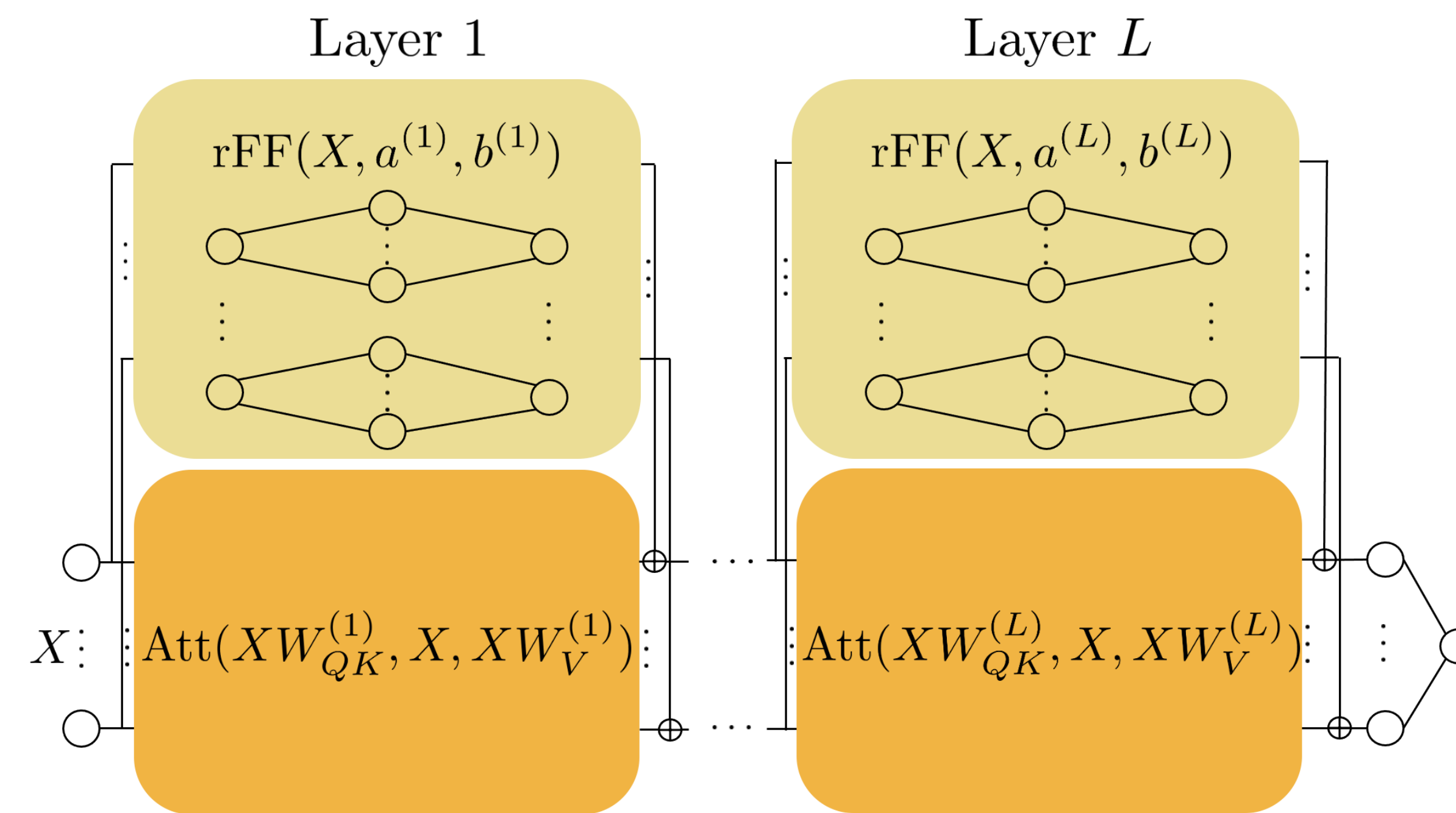
$$\forall f \in \mathcal{F}, \exists g \in \mathcal{N}(W) \text{ s.t. } \sup_{X \in [0, 1]^{N \times d}} |f(X) - g(X)| \leq \xi.$$

With sufficient number of channels N , it holds that $W^*(\xi, d, \mathcal{F}) = \Omega(\exp(cd)\xi^{-1/4})$ for some $c > 0$.

Intuition: The deep set only adopts a single-hidden layer net to reason the relationship among inputs. It is too **coarse** compared with the self-attention.

Generalization Error is Indep. of the Number of Channels

- We consider the transformer with N channels and L depth.
- The norms of W_{QK} and W_V in attention and weights in fully-connected layer are **bounded** by a tuple B in our transformer function class $\mathcal{F}_{\text{tf}}(B)$.
- Implement **layer-normalization** at the output of each layer.



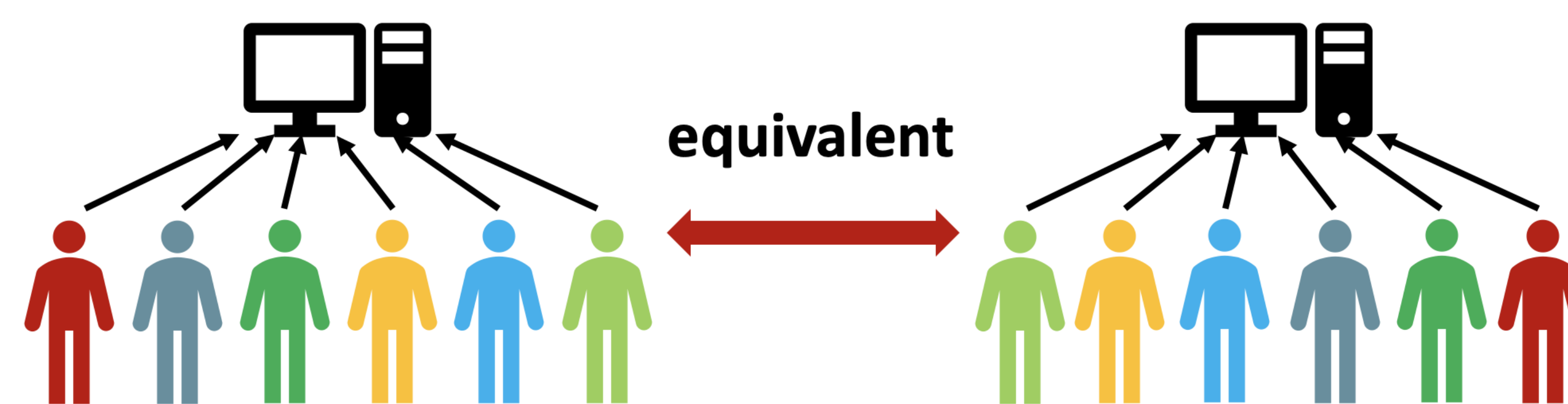
- Aim to predict the value of the response variable $y \in \mathbb{R}$ from the observation matrix $X \in \mathbb{R}^{N \times d}$, where $(X, y) \sim \nu$, and $|y| \leq V$.
- Estimate $f : \mathbb{R}^{N \times d} \rightarrow \mathbb{R}$ from i.i.d. samples $\mathcal{D}_{\text{reg}} = \{(X_i, y_i)\}_{i=1}^n$.
- The **risk** of using $f \in \mathcal{F}_{\text{tf}}(B)$ as a regressor on sample (X, y) is defined as $(f(X) - y)^2$.

For all $f \in \mathcal{F}_{\text{tf}}(B)$, with probability at least $1 - \delta$, we have

$$\begin{aligned} & \left| \mathbb{E}_{\nu}[(f(X) - y)^2] - \frac{1}{n} \sum_{i=1}^n (f(X_i) - y_i)^2 \right| \\ & \leq \frac{1}{2} \mathbb{E}_{\nu}[(f(X) - y)^2] + O\left(\frac{V^2}{n} \left[mL^2 d^2 \log \frac{m d L \bar{B} n}{V} + \log \frac{1}{\delta} \right]\right). \end{aligned}$$

- **Independent** of the number of channels N .
- **Polynomial** in the depth of the network L .

Homogeneous MARL



- We consider the Homogeneous MARL where the reward function and transition kernel are **permutation invariant**, i.e.,
- $$P^*(\bar{S}' | \bar{S}, \bar{A}) = P^*(\psi(\bar{S}') | \psi(\bar{S}), \psi(\bar{A})) \quad \text{and} \quad r(\bar{S}, \bar{A}) = r(\psi(\bar{S}), \psi(\bar{A})).$$
- There exists an **optimal policy** that is **permutation invariant**.
 - For any permutation invariant policy π , the corresponding **value function** V^{π} and **action-value function** Q^{π} are permutation invariant.
 - We only consider the class of **permutation invariant policies** Π , where $\pi(\bar{A} | \bar{S}) = \pi(\psi(\bar{A}) | \psi(\bar{S}))$ for all permutations ψ .

Pessimistic Model-Free Offline Reinforcement Learning

- Learn from the **i.i.d.** dataset $\mathcal{D} = \{(\bar{S}_i, \bar{A}_i, r_i, \bar{S}'_i)\}_{i=1}^n$.
- **Bellman error** of $f \in \mathcal{F}_{\text{tf}}(B)$ for π on \mathcal{D} is denoted as $\mathcal{E}(f, \pi; \mathcal{D})$.

Algorithm:

$$\hat{\pi} = \argmax_{\pi \in \Pi} \min_{f \in \mathcal{F}(\pi, \varepsilon)} f(\bar{S}_0, \pi), \text{ where } \mathcal{F}(\pi, \varepsilon) = \{f \in \mathcal{F}_{\text{tf}}(B) \mid \mathcal{E}(f, \pi; \mathcal{D}) \leq \varepsilon\}.$$

The coefficient $C_{\mathcal{F}_{\text{tf}}}$ measures the **coverage** of the dataset

$$C_{\mathcal{F}_{\text{tf}}} = \max_{f \in \mathcal{F}_{\text{tf}}} \mathbb{E}_{d_{p^*}}[(f(\bar{S}, \bar{A}) - \mathcal{T}^{\pi^*} f(\bar{S}, \bar{A}))^2] / \mathbb{E}_{\nu}[(f(\bar{S}, \bar{A}) - \mathcal{T}^{\pi^*} f(\bar{S}, \bar{A}))^2].$$

Assumptions:

- **Realizability:** $\inf_{f \in \mathcal{F}_{\text{tf}}} \sup_{\mu \in d_{\Pi}} \mathbb{E}_{\mu}[(f(\bar{S}, \bar{A}) - \mathcal{T}^{\pi} f(\bar{S}, \bar{A}))^2] \leq \varepsilon_{\mathcal{F}} < \infty$
- **Completeness:** $\sup_{f \in \mathcal{F}_{\text{tf}}} \inf_{\tilde{f} \in \mathcal{F}_{\text{tf}}} \mathbb{E}_{\nu}[(\tilde{f}(\bar{S}, \bar{A}) - \mathcal{T}^{\pi} f(\bar{S}, \bar{A}))^2] \leq \varepsilon_{\mathcal{F}, \mathcal{F}} < \infty$
- The coefficient $C_{\mathcal{F}_{\text{tf}}}$ is finite for the sampling distribution ν .

With probability at least $1 - \delta$, the suboptimality gap of the policy derived in the model-free algorithm is bounded as

$$V_{p^*}^*(\bar{S}_0) - V_{\hat{p}^*}^*(\bar{S}_0) \leq O(\text{Independent of the number of agents } N)$$

- Broken “the curse of many agents”.

Pessimistic Model-Based Offline Reinforcement Learning

- **System dynamics:** $\bar{S}' = F^*(\bar{S}, \bar{A}) + \bar{\varepsilon}$, where F^* is a nonlinear function, and $\bar{\varepsilon}_i \sim \mathcal{N}(0, \sigma^2 I_{d \times d})$ for $i \in [N]$ are independent random vectors.
- Learn the system dynamics with transformer function class.

Algorithm:

$$\hat{F}_{\text{MLE}} = \argmin_{F \in \mathcal{M}_{\text{tf}}} \frac{1}{n} \sum_{i=1}^n \|\bar{S}'_i - F(\bar{S}_i, \bar{A}_i)\|_F^2 \quad \text{and} \quad \hat{\pi} = \argmax_{\pi \in \Pi} \min_{F \in \mathcal{M}_{\text{MLE}}(\zeta)} V_{P_F}^{\pi}(\bar{S}_0).$$

$\mathcal{M}_{\text{MLE}}(\zeta) = \{F \in \mathcal{M}_{\text{tf}}(B') \mid \frac{1}{n} \sum_{i=1}^n \text{TV}(P_F(\cdot | \bar{S}_i, \bar{A}_i), \hat{P}_{\text{MLE}}(\cdot | \bar{S}_i, \bar{A}_i))^2 \leq \zeta\}$ is the confidence region.

Assumption

- **Realizability:** $F^* \in \mathcal{M}_{\text{tf}}(B')$.
- The coefficient $C_{\mathcal{M}_{\text{tf}}}$ is finite for the sampling distribution ν .

With probability at least $1 - \delta$, the suboptimality gap of the policy learned in the model-based algorithm is upper bounded as

$$V_{p^*}^*(\bar{S}_0) - V_{\hat{p}^*}^*(\bar{S}_0) \leq O(\text{Logarithmic in the number of agents } N)$$

Full Paper is Available at:



<https://arxiv.org/abs/2209.09845>