Multi-agent Markov Entanglement

Shuze Chen

Graduate School of Business Columbia University New York, NY 10027 shuze.chen@columbia.edu

Tianyi Peng

Graduate School of Business Columbia University New York, NY 10027 tianyi.peng@columbia.edu

Abstract

Value decomposition has long been a fundamental technique in multi-agent reinforcement learning and dynamic programming. Specifically, the value function of a global state (s_1,s_2,\ldots,s_N) is often approximated as the sum of local functions: $V(s_1,s_2,\ldots,s_N) \approx \sum_{i=1}^N V_i(s_i)$. This approach has found various applications in modern reinforcement learning systems. However, the theoretical justification for why this decomposition works so effectively remains underexplored. In this paper, we uncover the underlying mathematical structure that enables value decomposition. We demonstrate that a Markov decision process (MDP) permits value decomposition if and only if its transition matrix is not "entangled"—a concept analogous to quantum entanglement in quantum physics. Drawing inspiration from how physicists measure quantum entanglement, we introduce how to measure the "Markov entanglement" and show that this measure can be used to bound the decomposition error in general multi-agent MDPs. Using the concept of Markov entanglement, we proved that a widely-used class of policies, the index policy, is weakly-entangled and enjoys a sublinear $\mathcal{O}(\sqrt{N})$ scale of decomposition error for N-agent systems. Finally, we show Markov entanglement can be efficiently estimated, guiding practitioners on the feasibility of value decomposition.

1 Introduction

Learning the value function given certain policy, or *policy evaluation*, is one of the most fundamental tasks in RL. Significant attention has been paid to single-agent policy evaluation [39, 8, 40]. However, when it comes to multi-agent reinforcement learning (MARL), single-agent methodologies typically suffer from *the curse of dimensionality*: the state space of the system scales exponentially with the number of agents. To tackle this problem, one common technique is value decomposition,

$$V(s_1, s_2, \dots, s_N) \approx \sum_{i=1}^N V_i(s_i),$$

where V_i is some local function that can be learned independently by each agent. It quickly follows that this decomposition greatly reduces the computation complexity from exponential to linear dependency on the number of agents N.

The remaining question is whether this decomposition is effective. This is non-trivial due to the coupling of agents—individual agent's action and transition depend on other agents. For example, in a ride-hailing platform, if one driver took the order, then other drivers are not allowed fulfill the same order. As a result, value decomposition may lose information and introduce bias without considering the global constraints.

In the past several decades, both positive and negative results have been reported. Back to the last century, [49, 47] apply Lagrange relaxations to decompose the global value and obtain the well-known

Whittle index policy. The Lagrange decomposition idea has also been proved successful in many other important multi-agent tasks such as network revenue management [1, 50], resource allocation [27, 7], and online matching [11, 12, 36, 28]. However, Lagrange decomposition relies on the knowledge of system dynamics and [2] show its decomposition error can be arbitrarily bad for general multi-agent MDPs. In more recent days, practitioners apply online (deep) reinforcement learning to train a local value function for each individual agent. This practice gives birth to state-of-the-art dispatching policies in ride-hailing platforms and has been well recognized by the operations research community, such as DiDi Chuxing [33] (Daniel H. Wagner Prize, 2020) and Lyft [4] (Franz Edelman Laureates, 2024). Intervention policies based on a similar value decomposition idea also demonstrate substantial empirical advantages and have been deployed by a behavioral health platform in Kenya [5] (Pierskalla Award, 2024). In broader MARL literature, value decomposition serves as one key component of centralized training and decentralized execution (CTDE) paradigm, achieving strong empirical performance [38, 29, 35]. However, recent research has started reflecting on the invalidity and potential flaw of value decomposition in practice [25, 16].

Despite all these empirical success and failures, there remains little theoretical understanding of whether and how we can decompose the value function in multi-agent MDPs.

1.1 This paper

In this paper, we will uncover the underlying mathematical structure that enables/disables value decomposition. Our new theoretical framework quantifies the inter-dependence of agents in multiagent MDPs and systematically characterizes the effectiveness of value decomposition. For simplicity, we will demonstrate the main results through two-agent MDPs indexed by agent A and B. We later extend our results to general N-agent MDPs in Appendix H.

We start with a trivial example where two agents are independent, i.e. each following independent MDPs. It's clear that the global value function can be decomposed as the sum of value functions of local MDPs. As two agents are independent, it holds $P^{\pi}(s'_A, s'_B \mid s_A, s_B) = P^{\pi}(s'_A \mid s_A) \cdot P^{\pi}(s'_B \mid s_B)$, or in matrix form,

$$oldsymbol{P}_{AB}^{\pi} = oldsymbol{P}_{A}^{\pi} \otimes oldsymbol{P}_{B}^{\pi}$$
,

where \otimes is the tensor product or Kronecker product of matrices. The important question is whether we can extend beyond this trivial case of independent subsystems.

A Sufficient and Necessary Condition We introduce a new condition called "Markov Entanglement" to describe the intrinsic structure of transition dynamics in multi-agent MDPs.

Definition 1 (Markov Entanglement). *Consider a two-agent MDP with transition* P_{AB}^{π} . *If there exists*

$$\boldsymbol{P}_{AB}^{\pi} = \sum_{j=1}^{K} x_j \boldsymbol{P}_{A}^{(j)} \otimes \boldsymbol{P}_{B}^{(j)},$$

then P_{AB}^{π} is separable; otherwise entangled.

Compared with the preceding example of independent subsystems, Markov entanglement offers an intuitive interpretation: a two-agent MDP is separable if it can be expressed as a *linear combination of independent subsystems*. We then demonstrate,

separable
$$P_{AB}^{\pi} \Longleftrightarrow$$
 decomposable V_{AB}^{π} ,

where V_{AB}^{π} is decomposable if there exist local value functions V_A , V_B such that $V_{AB}^{\pi}(s_A, s_B) = V_A(s_A) + V_B(s_B)$ for all (s_A, s_B) . This result sharply unravels the secret structure of system dynamics governing value decomposition. As a sufficient condition, our finding strictly generalizes the previous independent subsystem example, extending it to scenarios involving interacting and coupled agents. As a necessary condition, we prove that exact value decomposition under any reward kernel requires the system dynamics to be separable. Taken together, this result provides a *complete characterization* of when exact value function decomposition is possible in multi-agent MDPs.

More interestingly, our Markov entanglement condition turns out be a mathematical counterpart of quantum entanglement in quantum physics, whose definition is provided below.

Definition 2 (Quantum Entanglement). Consider a two-party quantum state ρ_{AB} . If there exists

$$ho_{AB} = \sum_{j=1}^K x_j \rho_A^{(j)} \otimes \rho_B^{(j)}, \quad \boldsymbol{x} \geq 0,$$

then ρ_{AB} is separable; otherwise entangled.

The quantum state is represented by a *density matrix*, a positive semidefinite matrix with unit trace, analogous to transition matrix in the Markov world. The concept of quantum entanglement describes the inter-dependence of particles in a quantum system, while Markov entanglement describes that of agents in a Markov system.

Finally, we introduce several novel proof techniques concerning the sufficient and necessary condition, including an "absorbing" technique for separable transition matrices and a novel characterization of the linear space spanned by tensor products of transition matrices. We believe these techniques hold independent interest for the broader RL community.

Decomposition Error in General Multi-agent MDPs Despite the precise characterization of Markov entanglement and exact value decomposition, general multi-agent MDPs can exhibit arbitrary complexity, with agents intricately entangled. This raises a critical question: *can value decomposition serve as a meaningful approximation in such scenarios?* To address this, we introduce a mathematical quantification to measure the Markov entanglement in general multi-agent MDPs,

$$E(\mathbf{P}_{AB}^{\pi}) := \min_{\mathbf{P} \in \mathcal{P}_{SEP}} d(\mathbf{P}_{AB}^{\pi}, \mathbf{P}), \qquad (1)$$

where \mathcal{P}_{SEP} is the set of all separable transition matrices and $d(\cdot, \cdot)$ is some distance measure. In other words, the degree of Markov entanglement is determined by its distance to the closest separable transition matrix. This concept can also find its counterpart in quantum physics, with the measure of quantum entanglement defined as

$$E(\rho_{AB}) \coloneqq \min_{\rho \in \rho_{SEP}} d(\rho_{AB}, \rho),$$

where $\rho_{\rm SEP}$ is the set of all separable quantum states. In quantum physics, various distance measures have been designed for density matrices and capture different physical interpretations [31]. In the Markov world, we analogously design distance measures for transition matrices and relate them to the value decomposition error,

$$\left\| ext{decomposition error of } oldsymbol{V}_{AB}^{\pi}
ight\| = \mathcal{O}\Big(E(oldsymbol{P}_{AB}^{\pi})\Big)$$
 .

where $\|\cdot\|$ depends on the distance we use to measure Markov entanglement. We explore diverse distance measures including the well-known total variation distance and its stationary distribution weighted variant. We also design a novel agent-wise distance incorporating the multi-agent structure, which may be of independent interest to the MARL community. We further demonstrate how different distance measures give birth to the decomposition error in different norms.

Applications of Markov Entanglement Finally, we leverage our Markov entanglement theory to analyze several structured multi-agent MDPs. We prove that a widely-used class of index policies is asymptotically separable, exhibiting a decomposition error that scales as $\mathcal{O}(\sqrt{N})$ with the number of agents N. This result theoretically justifies the practical effectiveness of value decomposition for index-based policies. Our proof builds on innovations that integrate Markov entanglement with mean-field analysis. We also show that Markov entanglement admits an efficient empirical estimation, thus helping practitioners determine when value decomposition is feasible.

1.2 Other related work

In the first section, we have reviewed typical empirical works on value decomposition. Here, we complement that discussion with related literature on theoretical insights.

Prior theoretical research has extensively investigated the decomposition of optimal value functions in multi-agent settings. A prominent area involves Lagrange relaxation, with the Restless Multi-Armed

Bandit (RMAB, [49]) as a foundational model. Lagrange relaxation decouples the constraint of agents, yielding a decomposable value that upper bounds the original value. The per-agent decomposition error is proven to decay asymptotically to zero [47, 48, 41] and enjoys a quadratic or exponential rate [20, 21, 11, 51, 52]. Other work generalizes to Weakly-Coupled MDPs (WCMDPs) [6, 13, 19]. However, [2] showed Lagrange relaxation can have arbitrarily large errors and proposed an alternative decomposition called Approximate Linear Programs (ALP), which is proven to have tighter error [12]. Despite these advancements, characterizing decomposition error for general multi-agent MDPs remains unknown. In contrast, our Markov entanglement theory analyzes value decomposition for general multi-agent MDPs under arbitrary policies, including optimal ones.

Another line of theoretical work has concentrated on policy optimization via value decomposition. Despite reported empirical successes, rigorous theoretical analysis remains challenging. [5] derived an approximation ratio for a specific index policy on a two-state RMAB. [43, 16] analyzed the convergence of the CTDE paradigm under strong exploration assumptions, while also highlighting scenarios of divergence. In contrast, our work instead focuses on policy evaluation rather than optimization. This enables us to derive clear and interpretable bounds on the decomposition error for general finite-state multi-agent MDPs that only require the existence of a stationary distribution.

Notations We abbreviate subscripts $(s) := (s_{1:N}) := (s_1, s_2, \dots, s_N)$. Particularly, for two-agent case, when the context is clear, we abbreviate $(s) := (s_A, s_B) := (s_A, s_B)$. Let $[N] = \{1, 2, \dots, N\}$ and \mathbb{Z}^+ be the set of positive integers.

2 Model

We consider a standard two-agent MDP $\mathcal{M}_{AB}(\mathcal{S},\mathcal{A},P,r_A,r_B,\gamma)$ with joint state space $\mathcal{S}=\mathcal{S}_A\times\mathcal{S}_B$ and joint action space $\mathcal{A}=\mathcal{A}_A\times\mathcal{A}_B$ where A,B represent two agents. For simplicity, let $|\mathcal{S}_A|=|\mathcal{S}_B|=|S|$ and $|\mathcal{A}_A|=|\mathcal{A}_B|=|A|$. For agents at global state $s=(s_A,s_B)$ with action $a=(a_A,a_B)$ taken, the system will transit to $s'=(s'_A,s'_B)$ according to transition kernel $s'\sim P(\cdot\mid s,a)$ and each agent $i\in\{A,B\}$ will receive its local reward $r_i(s_i,a_i)$. The global reward r_{AB} is defined as the summation of local rewards $r_{AB}(s,a):=r_A(s_A,a_A)+r_B(s_B,a_B)$, or in vector form $r_{AB}\in\mathbb{R}^{|S|^2|A|^2}:=r_A\otimes e+e\otimes r_B$, where \otimes is the tensor product and $e=1\in\mathbb{R}^{|S||A|}$ is the vector of all ones. We further assume the local rewards are bounded, i.e. for agent $i\in\{A,B\}$, $|r_i(s_i,a_i)|\leq r_{\max}^i$ for all (s_i,a_i) .

Given any global policy $\pi\colon\mathcal{S}\to\Delta(\mathcal{A})$, the global Q-value under policy π is defined as the discounted summation of global rewards $Q_{AB}^\pi(s,a)=\mathbb{E}\left[\sum_{t=0}^\infty \gamma^t r_{AB}(s^t,a^t)\mid \pi,(s^0,a^0)=(s,a)\right]$ where $\gamma\in[0,1)$ is the discount factor. The value function is then defined as $V_{AB}^\pi(s)=\mathbb{E}_{a\sim\pi(\cdot|s)}\left[Q_{AB}^\pi(s,a)\right]$. We denote $P_{AB}^\pi\in\mathbb{R}^{|S|^2|A|^2\times|S|^2|A|^2}$ as the transition matrix induced by π where $P_{AB}^\pi(s',a'\mid s,a)=P(s'\mid s,a)\cdot\pi(a'\mid s')$. Then by the Bellman Equation, we have $Q_{AB}^\pi=(I-\gamma P_{AB}^\pi)^{-1}r_{AB}$. Our objective is to decompose this global Q-value Q_{AB}^π as the summation of some local functions Q_A and Q_B , i.e. $Q_{AB}^\pi(s,a)=Q_A(s_A,a_A)+Q_B(s_B,a_B)$, or in vector form,

$$Q_{AB}^{\pi} = Q_A \otimes \mathbf{e} + \mathbf{e} \otimes Q_B. \tag{2}$$

Notice we formally introduce our research question using Q-value instead of V-value function as in the introduction. Q-value decomposition is a stronger result that implies V-value function decomposition. It also turns out that Q-value further incorporates action information enabling more general theoretical analysis. More discussions can be found in Appendix B.

2.1 Local (Q-)value functions

Recent literature offers several algorithms for learning local (Q-)values. In this paper, we use a meta-algorithm framework in 1 to summarize their underlying principles.

This meta-algorithm framework is simple and intuitive: each agent independently fits its local Q-values based on its local observations. Notably, the framework requires no prior knowledge of the MDP, and learning can be performed in a fully decentralized manner. Furthermore, we use term *meta* in that we do not pose restrictions on how agents estimate their local Q-values. For tabular case, one

¹In Appendix J.4, we extend our results to multi-agent MDP model where the global cannot be decomposed.

Meta Algorithm 1: Leaning Local Q-value Functions

Require: Global policy π ; horizon length T.

- 1: Execute π for T epochs and obtain $\mathcal{D} = \left\{ (s_{AB}^t, a_{AB}^t, r_{AB}^t, s_{AB}^{t+1}, a_{AB}^{t+1}) \right\}_{t=1}^{T-1}$. 2: Each agent $i \in \{A, B\}$ fits Q_i^{π} using local observations $\mathcal{D}_i = \left\{ (s_i^t, a_i^t, r_i^t, s_i^{t+1}, a_i^{t+1}) \right\}_{t=1}^{T-1}$.

can plug in Temporal Difference (TD) learning [39] or its variants. For large-scale problems, one can apply linear function approximations (e.g. [5, 24, 8]) or more sophisticated neural networks (e.g. [33, 38, 29]).

Despite the flexibility in fitting local value functions, it is helpful to call out a particular approach: TD learning for local Q-values in the tabular case, as it facilitates the analysis and reveals the structure of value decomposition in the next section.

Local TD learning. Although each agent's environment is not Markovian in a local sense (it is, more precisely, partially observed Markovian), one can still define its "marginalized" local transition matrix under the stationary distribution. Mathematically, for agent A, we denote $P_A^{\pi} \in \mathbb{R}^{|S||A| \times |S||A|}$ as its local transition where

$$P_A^{\pi}(s_A', a_A' \mid s_A, a_A) = \sum_{s_B', a_B'} \sum_{s_B, a_B} P_{AB}^{\pi}(s_{AB}', a_{AB}' \mid s_{AB}, a_{AB}) \mu_{AB}^{\pi}(s_B, a_B \mid s_A, a_A).$$
(3)

Here, $\mu_{AB}^{\pi} \in \Delta(\mathcal{S})$ denotes the global stationary distribution under policy π (for convenience, we assume π induces a unichain, i.e. μ_{AB}^{π} is unique and strictly positive).² Given this "marginalized" local transition, the local Q-values obtained by Meta Algorithm 1 using tabular TD learning converge to the solution of the following "marginalized" Bellman equation:

$$Q_A^{\pi} = (\boldsymbol{I} - \gamma \boldsymbol{P}_A^{\pi})^{-1} \boldsymbol{r}_A.$$

By symmetry, we can derive analogous results for agent B, obtaining its transition matrix P_B^{π} and local Q-values Q_B^{π} . Next, we show how Q_A^{π} and Q_B^{π} contribute to the exact value decomposition.

Exact value decomposition

To begin, recall the key condition we identify in the introduction: Markov Entanglement in Definition 1. Our first theorem shows that an MDP with no Markov entanglement is indeed sufficient for the exact value decomposition. More importantly, local TD learning (or Meta Algorithm 1 more generally) is guaranteed to recover such decomposition, i.e. $Q_{AB}^{\pi}=Q_A^{\pi}\otimes e+e\otimes Q_B^{\pi}$.

Theorem 1. Consider a two-agent MDP \mathcal{M}_{AB} and policy $\pi \colon \mathcal{S} \to \Delta(\mathcal{A})$. If two agents are separable, i.e. there exists $K \in \mathbb{Z}^+$, measure $\{x_j\}_{j \in [K]}$, and transition matrices $\{\boldsymbol{P}_A^{(j)}, \boldsymbol{P}_B^{(j)}\}_{j \in [K]}$ such that $\boldsymbol{P}_{AB}^{\pi} = \sum_{j=1}^{K} x_j \boldsymbol{P}_A^{(j)} \otimes \boldsymbol{P}_B^{(j)}$. Then it holds $\boldsymbol{P}_A^{\pi} = \sum_{i=1}^{K} x_j \boldsymbol{P}_A^{(j)}$ and $\boldsymbol{P}_B^{\pi} = \sum_{j=1}^{K} x_j \boldsymbol{P}_B^{(j)}$. Furthermore, the Eq. (2) holds

$$Q_{AB}^{\pi} = Q_A^{\pi} \otimes \boldsymbol{e} + \boldsymbol{e} \otimes Q_B^{\pi}$$
.

This theorem establishes that even when the system is not independent, as long as it can be represented as a linear combination of independent subsystems, the global Q-value admits an exact decomposition.

An illustrative example of coupling and Markov entanglement To elucidate the concept of Markov entanglement, we present an example of two-agent MDP where agents are coupled but not entangled. Consider a two-agent MDP \mathcal{M}_{AB} with $|\mathcal{A}_A| = |\mathcal{A}_B| = 2$, where action 1 means activate and 0 means idle. Each agent $i \in \{A, B\}$ has its own local transition kernel P_i . We examine the following policy: at each time-step, we randomly activate one agent and keep another idle, i.e.

²For $\mu_{AB}^{\pi}(s_B, a_B \mid s_A, a_A)$ to be well-defined, we require $\mu_{AB}^{\pi}(s_A, a_A) > 0$. If $\mu_{AB}^{\pi}(s_A, a_A) = 0$, then action a_A is never taken in state s_A under policy π , and we exclude such pairs by restricting the feasible action set $A(s_A)$. All theoretical results apply to the remaining valid state-action pairs.

 $\pi(\boldsymbol{a} \mid \boldsymbol{s}) = 1/2$ if $\boldsymbol{a} = (0,1)$ or $\boldsymbol{a} = (1,0)$. Consequently, this policy couples the agents through the constraint $a_A + a_B = 1$ at each timestep. However, we will demonstrate that despite this coupling, there's *no* entanglement. Specifically, we construct the following decomposition

$$\boldsymbol{P}_{AB}^{\pi} = \frac{1}{2} \boldsymbol{P}_{A}^{0} \otimes \boldsymbol{P}_{B}^{1} + \frac{1}{2} \boldsymbol{P}_{A}^{1} \otimes \boldsymbol{P}_{B}^{0}, \qquad (4)$$

where P_i^a refers to the transition matrix of agents $i \in \{A, B\}$ taking action $a \in \{0, 1\}$. Intuitively, the right-hand side of Eq. (4) describes how at each time step, the global system randomly selects between two possible transitions: $P_A^0 \otimes P_B^1$ or $P_A^1 \otimes P_B^0$, each with equal probability (akin to rolling a fair dice). This example thus clearly demonstrates a *coupled* system can still be *separable* and thus admits an exact value decomposition.

Proof of sufficiency Theorem 1 admits a simple proof based on the several basic properties of tensor product. First of all, given $P_{AB}^{\pi} = \sum_{j=1}^{K} x_j P_A^{(j)} \otimes P_B^{(j)}$, we can plug this into the formulation of P_A^{π} in Eq. (3) and quickly verify $P_A^{\pi} = \sum_{i=1}^{K} x_i P_A^{(i)}$. It remains to show Eq. (2). Notice that

$$(\boldsymbol{I} - \gamma \boldsymbol{P}_{AB}^{\pi})^{-1} (\boldsymbol{r}_{A} \otimes \boldsymbol{e}) = \sum_{t=0}^{\infty} \gamma^{t} \left(\sum_{j=1}^{K} x_{j} \boldsymbol{P}_{A}^{(j)} \otimes \boldsymbol{P}_{B}^{(j)} \right)^{t} (\boldsymbol{r}_{A} \otimes \boldsymbol{e})$$

$$\stackrel{(i)}{=} \sum_{t=0}^{\infty} \gamma^{t} \left(\left(\sum_{j=1}^{K} x_{j} \boldsymbol{P}_{A}^{(j)} \right)^{t} \boldsymbol{r}_{A} \right) \otimes \boldsymbol{e} = \left((\boldsymbol{I} - \gamma \boldsymbol{P}_{A}^{\pi})^{-1} \boldsymbol{r}_{A} \right) \otimes \boldsymbol{e} = Q_{A}^{\pi} \otimes \boldsymbol{e}.$$

where we refer to (i) as an "absorbing" technique based on the bilinearity and mixed-product property of tensor product³. Specifically, since Pe = e for any transition matrix P, we have for any t,

$$\left(\sum_{j=1}^{K} x_j \boldsymbol{P}_A^{(j)} \otimes \boldsymbol{P}_B^{(j)}\right)^t (\boldsymbol{r}_A \otimes \boldsymbol{e}) = \left(\sum_{j=1}^{K} x_j \boldsymbol{P}_A^{(j)} \otimes \boldsymbol{P}_B^{(j)}\right)^{t-1} \left(\sum_{j=1}^{K} x_j \left(\boldsymbol{P}_A^{(j)} \boldsymbol{r}_A\right) \otimes \left(\boldsymbol{P}_B^{(j)} \boldsymbol{e}\right)\right)$$

$$= \left(\sum_{j=1}^{K} x_j \boldsymbol{P}_A^{(j)} \otimes \boldsymbol{P}_B^{(j)}\right)^{t-1} \left(\sum_{j=1}^{K} x_j \boldsymbol{P}_A^{(j)} \boldsymbol{r}_A\right) \otimes \boldsymbol{e} = \dots = \left(\left(\sum_{j=1}^{K} x_j \boldsymbol{P}_A^{(j)}\right)^t \boldsymbol{r}_A\right) \otimes \boldsymbol{e}.$$

Similar results can be derived for P_B^{π} such that $(I - \gamma P_{AB}^{\pi})^{-1} (e \otimes r_B) = e \otimes Q_B^{\pi}$. Finally, combining the above results, we have

$$Q_{AB}^{\pi} = \left(\boldsymbol{I} - \gamma \boldsymbol{P}_{AB}^{\pi}\right)^{-1} \boldsymbol{r}_{AB} = \left(\boldsymbol{I} - \gamma \boldsymbol{P}_{AB}^{\pi}\right)^{-1} \left(\boldsymbol{r}_{A} \otimes \boldsymbol{e} + \boldsymbol{e} \otimes \boldsymbol{r}_{B}\right) = Q_{A}^{\pi} \otimes \boldsymbol{e} + \boldsymbol{e} \otimes Q_{B}^{\pi}.$$

3.1 Necessary condition for the exact value decomposition

We then investigate whether Markov entanglement is necessary for the exact Q-value decomposition. The answer is in general no, since one can construct trivial counterexamples such as $r_A = r_B = 0$ or $\gamma = 0$, where the decomposition trivially holds. On the other hand, we focus on a stronger and more general concept of the exact value decomposition that holds under any reward kernel given $\gamma > 0$. Formally, we present the following theorem.

Theorem 2. Consider a two-agent Markov MDP \mathcal{M}_{AB} with discount factor $\gamma > 0$ and $\pi \colon \mathcal{S} \to \Delta(\mathcal{A})$. Suppose there exists local functions $Q_i \colon r_i \to \mathbb{R}^{|S||A|}$ for $i \in \{A, B\}$ such that $Q_{AB}^{\pi} = Q_A(r_A) \otimes e + e \otimes Q_B(r_B)$ holds for any pair of reward r_A, r_B , then A, B must be separable.

Combined with Theorem 1, we conclude Markov entanglement serves as a sufficient and necessary condition for the exact value decomposition. We also emphasize that Theorem 2 considers general local functions Q_i . This generality accommodates all methods for fitting local Q_i , such as deep neural networks, provided that the training relies solely on the local observations of agent i.

There exist other possible ways for value decomposition. For example, [38, 16] consider $Q_{AB}^{\pi}(s, a) = L_A(s_A, a_A, r_{AB}) + L_B(s_B, a_B, r_{AB})$ where L_A, L_B are learned jointly via minimizing the global Bellman error⁴; [35, 29, 37, 42] consider general monotonic operations beyond

³We introduce several basic properties of tensor product in Appendix A.

⁴In Appendix E, we provide an example of entangled MDP that allows for an exact value decomposition where L_A depends on both r_A and r_B .

additive decompositions. These methods introduce possibly richer representations at the cost of more sophisticated implementations and less interpretability, which is beyond the scope of this paper.

Proof sketch of necessity Our proof builds on several novel techniques. Recall \mathcal{P}_{SEP} is the set of all separable transition matrices.

Step 1: Understanding the orthogonal complement. If a transition matrix is entangled, it will have non-zero component in the orthogonal complement of \mathcal{P}_{SEP} , which we construct as

$$\mathcal{P}_{\text{SEP}}^{\perp} = \left\{ \left. \sum_{j=1}^{|S||A|-1} \left(\varepsilon_j \boldsymbol{e}^\top \right) \otimes \boldsymbol{W}_j^1 + \sum_{j=1}^{|S||A|-1} \boldsymbol{W}_j^2 \otimes \left(\varepsilon_j \boldsymbol{e}^\top \right) \, \middle| \, W_{1:j}^1, W_{1:j}^2 \in \mathbb{R}^{|S||A| \times |S||A|} \right\} \,,$$

where $\varepsilon_j = (1,0,\dots,0,-1,0,\dots,0)^{\top}$ with the first element 1 and (j+1)-th element -1. Then, we study an intermediate transition matrix $(1-\gamma)(\boldsymbol{I}-\gamma\boldsymbol{P}_{AB}^{\pi})^{-1}$. We show if it's entangled, we are able to construct $\boldsymbol{r}_A, \boldsymbol{r}_B$ based on its component in $\mathcal{P}_{\text{SEP}}^{\perp}$ such that Q_{AB}^{π} is not decomposable under this pair of rewards. We thus conclude decomposable $Q_{AB}^{\pi} \Longrightarrow \text{separable } (1-\gamma)(\boldsymbol{I}-\gamma\boldsymbol{P}_{AB}^{\pi})^{-1}$.

Step 2: Connecting to "inverse". Finally, we complete the proof via a lemma showing separable $(1-\gamma)(I-\gamma P_{AB}^\pi)^{-1} \Longleftrightarrow$ separable P_{AB}^π . The \Longleftrightarrow side is straightforward since $(I-\gamma P_{AB}^\pi)^{-1}$ is the Neumann series of γP_{AB}^π . For the converse \Longrightarrow , we seek to invert this Neumann series. This is achieved by a careful analysis of the operator norm of $I-(1-\gamma)(I-\gamma P_{AB}^\pi)^{-1}$.

4 Value decomposition error in general two-agent MDPs

In general, the system transition P^{π}_{AB} can be arbitrarily entangled. In these scenarios, we investigate when value decomposition $Q^{\pi}_{A} \otimes e + e \otimes Q^{\pi}_{B}$ is an effective approximation of Q^{π}_{AB} . As mentioned in the introduction, we define the measure of Markov entanglement in Eq. (1) as certain distance between P^{π}_{AB} and its closet separable transition matrix. We will examine several distance measures for transition matrices and relate them to the decomposition error.

4.1 Entry-wise error bound

Total variation distance One widely used metric for transition matrices is Total Variation (TV) distance. Specifically, for two transition matrices $P, P' \in \mathbb{R}^{|S|^2|A|^2 \times |S|^2|A|^2}$, define

$$\|\boldsymbol{P} - \boldsymbol{P}'\|_{\text{TV}} := \max_{(\boldsymbol{s}, \boldsymbol{a}) \in \mathcal{S} \times \mathcal{A}} D_{\text{TV}}(\boldsymbol{P}(\cdot, \cdot \mid \boldsymbol{s}, \boldsymbol{a}), \boldsymbol{P}'(\cdot, \cdot \mid \boldsymbol{s}, \boldsymbol{a})),$$
 (5)

where D_{TV} is the total variation distance between probability measures. While TV distance is straightforward, it does not take into account the inherent multi-agent structure.

Agent-wise distance We thus introduce a more refined distance specially designed for multi-agent MDPs. Formally, the Agent-wise Total Variation (ATV) distance between two transition matrices $P, P' \in \mathbb{R}^{|S|^2|A|^2 \times |S|^2|A|^2}$ w.r.t agent A is defined as

$$\|\boldsymbol{P} - \boldsymbol{P}'\|_{\text{ATV}_{A}} \coloneqq \max_{(\boldsymbol{s}, \boldsymbol{a}) \in \mathcal{S} \times \mathcal{A}} D_{\text{TV}} \left(\sum_{s'_{B}, a'_{B}} \boldsymbol{P}(\cdot, \cdot \mid \boldsymbol{s}, \boldsymbol{a}), \sum_{s'_{B}, a'_{B}} \boldsymbol{P}'(\cdot, \cdot \mid \boldsymbol{s}, \boldsymbol{a}) \right). \tag{6}$$

The ATV distance w.r.t agent B can be defined similarly. Intuitively, compared to TV, ATV focuses on an individual agent and measures the difference between its local transitions. One can also verify ATV is tighter distance, i.e. $\| \boldsymbol{P} - \boldsymbol{P}' \|_{\text{ATV}_A} \leq \| \boldsymbol{P} - \boldsymbol{P}' \|_{\text{TV}}$. We can plug ATV into Eq. (1) and obtain the measure of Markov entanglement w.r.t ATV distance $E_i(\boldsymbol{P}_{AB}^{\pi}) \coloneqq \min_{\boldsymbol{P} \in \mathcal{P}_{\text{SEP}}} \| \boldsymbol{P}_{AB}^{\pi} - \boldsymbol{P} \|_{\text{ATV}_i}$ for $i \in \{A, B\}$. In fact, one can also verify

$$E_{A}(\mathbf{P}_{AB}^{\pi}) = \min_{\mathbf{P}_{A}} \max_{(\mathbf{s}, \mathbf{a}) \in \mathcal{S} \times \mathcal{A}} D_{\text{TV}} \left(\mathbf{P}_{AB}^{\pi}(\cdot, \cdot \mid \mathbf{s}, \mathbf{a}), \mathbf{P}_{A}(\cdot, \cdot \mid s_{A}, a_{A}) \right), \tag{7}$$

The following theorem connects these measures to the value decomposition error.

Theorem 3. Consider a two-agent Markov system \mathcal{M}_{AB} and policy $\pi \colon \mathcal{S} \to \Delta(\mathcal{A})$ with the measure of Markov entanglement $E_A(\mathbf{P}_{AB}^{\pi})$, $E_B(\mathbf{P}_{AB}^{\pi})$ defined in Eq. (7), then the decomposition error is entry-wise bounded by the measure of Markov entanglement,

$$\left\|Q_{AB}^{\pi} - (Q_A^{\pi} \otimes \boldsymbol{e} + \boldsymbol{e} \otimes Q_B^{\pi})\right\|_{\infty} \leq \frac{4\gamma \left(E_A(\boldsymbol{P}_{AB}^{\pi})r_{\max}^A + E_B(\boldsymbol{P}_{AB}^{\pi})r_{\max}^B\right)}{(1-\gamma)^2}.$$

4.2 Error weighted by stationary distribution

Entry-wise error bound is a very strong result for Q-value decomposition. This comes with the entry-wise TV bounds in both TV and ATV distance. An alterative choice is to consider an error weighted by the stationary distribution. Formally, consider

$$\left\|Q_{AB}^{\pi}-(Q_A^{\pi}\otimes \boldsymbol{e}+\boldsymbol{e}\otimes Q_B^{\pi})\right\|_{\mu_{AB}^{\pi}}\coloneqq \sum_{\boldsymbol{s},\boldsymbol{a}}\mu_{AB}^{\pi}(\boldsymbol{s},\boldsymbol{a}) \Big|Q_{AB}^{\pi}(\boldsymbol{s},\boldsymbol{a})-(Q_A^{\pi}(s_A,a_A)+Q_B^{\pi}(s_B,a_B))\Big|\,.$$

We note that this norm is clearly weaker than the entry-wise norm. Nevertheless, a stationary distribution weighted error bound is sufficient in many practical scenarios. Similar ideas are also quite common in policy evaluation literature [14, 40, 9].

Distance weighted by stationary distribution To analyze this μ_{AB}^{π} -weight decomposition error, we analogously propose the μ_{AB}^{π} -weighted distance measure of Markov entanglement. Specifically, we have the following μ_{AB}^{π} -weighted version of Eq. (7).

$$E_A(\mathbf{P}_{AB}^{\pi}) = \min_{\mathbf{P}_A} \sum_{\mathbf{s}, \mathbf{a}} \mu_{AB}^{\pi}(\mathbf{s}, \mathbf{a}) D_{\text{TV}} \left(\mathbf{P}_{AB}^{\pi}(\cdot, \cdot \mid \mathbf{s}, \mathbf{a}), \mathbf{P}_A(\cdot, \cdot \mid \mathbf{s}_A, a_A) \right). \tag{8}$$

Eq. (8) substitutes the μ_{AB}^{π} -weighted average for the maximum operator in Eq. (7). Finally, we have the following variant of Theorem 3.

Theorem 4. Under the same setup as Theorem 3 with μ_{AB}^{π} -weighted measure of Markov entanglement $E_A(\mathbf{P}_{AB}^{\pi})$, $E_B(\mathbf{P}_{AB}^{\pi})$ defined in Eq. (8), the μ_{AB}^{π} -weighted decomposition error is bounded,

$$\left\|Q_{AB}^{\pi} - (Q_A^{\pi} \otimes \boldsymbol{e} + \boldsymbol{e} \otimes Q_B^{\pi})\right\|_{\mu_{AB}^{\pi}} \leq \frac{4\gamma \left(E_A(\boldsymbol{P}_{AB}^{\pi})r_{\max}^A + E_B(\boldsymbol{P}_{AB}^{\pi})r_{\max}^B\right)}{(1-\gamma)^2}.$$

Compared to Theorem 3, Theorem 4 measures a weaker μ_{AB}^{π} -weighted decomposition error, while the condition on P_{AB}^{π} is also relaxed, requiring only a weighted average bound in Eq. (8).

4.3 Multi-agent Markov entanglement

Finally, we extend the results to multi-agent MDPs with the measure of Markov entanglement $E_{1:N}(\mathbf{P}_{1:N}^{\pi})$ for an N-agent MDP. The extension is relatively straightforward. We demonstrate the extension of Theorem 4 below and more details can be found in Appendix H.

Theorem 5. Consider a N-agent MDP $\mathcal{M}_{1:N}$ with the measure of Markov entanglement $E_i(\mathbf{P}_{1:N}^{\pi})$ w.r.t ATV distance, the $\mu_{1:N}^{\pi}$ -weighted decomposition error is bounded by the measure of Markov entanglement,

$$\left\| Q_{1:N}^{\pi}(s, \boldsymbol{a}) - \sum_{i=1}^{N} Q_{i}^{\pi}(s_{i}, a_{i}) \right\|_{\mu_{1:N}^{\pi}} \leq \frac{4\gamma \left(\sum_{i=1}^{N} E_{i}(\boldsymbol{P}_{1:N}^{\pi}) r_{\max}^{i} \right)}{(1 - \gamma)^{2}}.$$

5 Applications of Markov Entanglement

In this section, we apply Markov entanglement and demonstrate a widely-used class of index policies is asymptotically separable. To begin, we introduce the model of Restless Multi-Armed Bandit (RMAB, [49]). In an N-agent RMAB, each agent follows a homogeneous two-action MDP with action 1 meaning activate and 0 idle. A central decision maker will activate $M \leq N$ agents at each timestep and leave other agents idle. In other words, agents transit independently but are coupled under constraint $\sum_{i=1}^{N} a_i = M$. In RMAB, arguably the most classical and widely-used policy is the index policy, which we formally define as

Definition 3 (Index Policy). There exists a priority index ν_s for each local state s. The decision maker will always activate agents in the descending order of the priority until the budget constraint M is met. Ties are resolved fairly via uniform random sampling of agents at the same state.

The index policy traces back to the well-known Gittins Index [46], Whittle Index [49, 47, 20], and fluid-based index policies [41, 21]. [33, 4, 5, 30, 44, 3] apply data-driven method to optimize index policies and report great empirical success in industrial implementations. Understanding the mystery behind such success calls for a theory for general index policies. We then present our main theorem.

Theorem 6. Consider an N-agent restless multi-armed bandit. For any index policy satisfying mild technical conditions, there exists constant C independent of N, such that for any agent $i \in [N]$, its $\mu_{1:N}^{\pi}$ -weighted measure of Markov entanglement is bounded, $E_i(\mathbf{P}_{1:N}^{\pi}) \leq C/\sqrt{N}$.

Theorem 6 requires two standard technical conditions for index policies: non-degenerate and uniform global attractor property, which are used in almost all related theoretical work [47, 41, 20, 21] and are detailed in Appendix I. Theorem 6 justifies index policies are asymptotically separable. Combined with an N-agent version of Theorem 4, we obtain the sublinear decomposition error for index policies

$$\left\| Q_{1:N}^{\pi}(\boldsymbol{s}, \boldsymbol{a}) - \sum_{i=1}^{N} Q_{i}^{\pi}(s_{i}, a_{i}) \right\|_{\mu_{1:N}^{\pi}} \leq \mathcal{O}(\sqrt{N}).$$

This sublinear error result explains why the value decomposition in [33, 4, 5] manages to effectively approximate the global value function in large-scale practical applications.

5.1 Efficient verification of value decomposition

For practitioners, verifying the feasibility of value decomposition is challenging due to the exponential computational complexity of estimating the global Q-value. As a solution, Markov entanglement offers an efficient way to empirically test whether value decomposition can be safely applied. Consider the μ_{AB}^{π} -weighted measure of Markov entanglement in Eq. (8), we have

$$E_{A}(\mathbf{P}_{AB}^{\pi}) \approx \frac{1}{2} \min_{\mathbf{P}_{A}} \frac{1}{T} \sum_{t=1}^{T} \sum_{s'_{A}, a'_{A}} \left| \mathbf{P}_{AB}^{\pi}(s'_{A}, a'_{A} \mid \mathbf{s}^{t}, \mathbf{a}^{t}) - \mathbf{P}_{A}(s'_{A}, a'_{A} \mid s'_{A}, a'_{A}) \right|$$
(9)

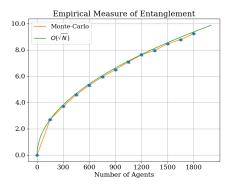
In other words, we can apply a Monte-Carlo estimation for $E_A(P_{AB}^{\pi})$. Notice Eq. (9) is *convex* for P_A , which enables efficient solutions. As a result, Eq. (9) provides an efficient estimation of Markov entanglement via simulation and can be easily extend to N-agent MDPs.

Numerical experiments. Finally, we empirically study the value decomposition for the index policy on a circulant RMAB benchmark [3, 52, 10, 18] that has 4 different states each local agent. As a result, the global state space scales as large as $4^{1800} > 10^{1000}$ for N=1800 agents. The specific transitions and rewards are introduced in Appendix K. For each RMAB instance, we sample a trajectory of length T=5N and use the collected data to i) solve Eq. (9) to estimate the measure of Markov entanglement; ii) train local Q-value decomposition. It quickly follows from the results in Figure 1:

The estimated Markov entanglement decays as $\mathcal{O}(1/\sqrt{N})$ in the left panel, consistent with theoretical predictions. This also implies a low decomposition error scaling of $\mathcal{O}(\sqrt{N})$, as seen in the right panel. Furthermore, the simulated trajectory has a length of T=5N while the global state space has size $|S|^N$, making both entanglement estimation and local Q-value decomposition sample-efficient.

6 Discussions

Comparison with quantum entanglement One notable difference between the definition of Markov and quantum entanglement is that the former does not require coefficients $x \geq 0$. In Appendix C, we show there exist separable two-agent MDPs that can only be represented by linear combinations but not convex combinations of independent subsystems, highlighting a structural difference between Markov and quantum entanglement. Finally, we emphasize that our analogy to quantum entanglement is mostly in the mathematical formulation; there is no clear physical interpretation analogy between Markov and quantum entanglement.



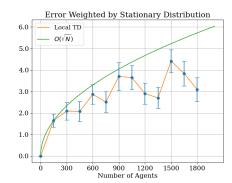


Figure 1: Circulant RMAB under an index policy. *Left:* empirical estimation of Markov entanglement multiplied by the number of agents, $NE_1(P_{1:N}^{\pi})$. *Right:* μ -weighted decomposition error.

Relations to Influenced-based MARL There's another line of MARL research that explicitly models the influence of other agents as intrinsic rewards for exploration [45, 26]. It turns out the mutual information can be viewed as the measure of Markov entanglement under KL-divergence. Specifically, we can rewrite mutual information in [45] as

$$I(S'_{2}, A'_{2}; S_{2}, A_{2}|S_{1}, A_{1}) = \sum_{\boldsymbol{s}, \boldsymbol{a}, s'_{2}, a'_{2}} p^{\pi}(\boldsymbol{s}, \boldsymbol{a}, s'_{2}, a'_{2}) \left(\log \frac{p^{\pi}(s'_{2}, a'_{2}|\boldsymbol{s}, \boldsymbol{a})}{p^{\pi}(s'_{2}, a'_{2}|s_{2}, a_{2})} \right)$$
$$= \sum_{\boldsymbol{s}, \boldsymbol{a}} \mu^{\pi}(\boldsymbol{s}, \boldsymbol{a}) D_{KL} \left(p^{\pi}(\cdot|\boldsymbol{s}, \boldsymbol{a}) ||P_{2}(\cdot|s_{2}, a_{2}) \right).$$

This is highly related to our measure of Markov entanglement under a μ^{π} -weighted agent-wise KL-divergence, which we can define as

$$E_2(\mathbf{P}_{12}) = \min_{\mathbf{P}_2} \sum_{\mathbf{s}, \mathbf{a}} \mu^{\pi}(\mathbf{s}, \mathbf{a}) D_{KL} \left(p^{\pi}(\cdot | \mathbf{s}, \mathbf{a}) || P_2(\cdot | s_2, a_2) \right) .$$

Intuitively, the measure of Markov entanglement can be viewed as how closely one agent can be approximated as an independent subsystem. This characterization aligns naturally with mutual information. Furthermore, since KL-divergence provides an upper bound for total variation distance, it consequently bounds our Markov entanglement measure relative to the ATV distance introduced in our paper. This connection demonstrates that influence-based MARL methods naturally fit within our theoretical framework, corresponding to a specialized distance measure.

7 Conclusion

This paper established the mathematical foundation of value decomposition in MARL. Drawing inspiration from quantum physics, we propose the idea of Markov entanglement and prove that it serves as a sufficient and necessary condition for the exact value decomposition. We further characterize the decomposition error in general multi-agent MDPs through the measure of Markov entanglement. As application examples, we prove widely-used index policies are asymptotically separable and suggest practitioners using Markov entanglement as a proxy for estimating the effectiveness of value decomposition.

Acknowledgments and Disclosure of Funding

We thank all anonymous reviewers for their constructive comments. We are also grateful to Prof. Tongyang Li for valuable and insightful discussions.

References

[1] Daniel Adelman. Dynamic bid prices in revenue management. *Operations Research*, 55(4): 647–661, 2007.

- [2] Daniel Adelman and Adam J. Mersereau. Relaxations of weakly coupled stochastic dynamic programs. *Operations Research*, 56(3):712–727, 2008. doi: 10.1287/opre.1070.0445.
- [3] Konstantin E Avrachenkov and Vivek S Borkar. Whittle index based q-learning for restless bandits with average reward. *Automatica*, 139:110186, 2022.
- [4] Xabi Azagirre, Akshay Balwally, Guillaume Candeli, Nicholas Chamandy, Benjamin Han, Alona King, Hyungjun Lee, Martin Loncaric, Sébastien Martin, Vijay Narasiman, Zhiwei (Tony) Qin, Baptiste Richard, Sara Smoot, Sean Taylor, Garrett van Ryzin, Di Wu, Fei Yu, and Alex Zamoshchin. A better match for drivers and riders: Reinforcement learning at lyft. *INFORMS Journal on Applied Analytics*, 54(1):71–83, 2024.
- [5] Jackie Baek, Justin J Boutilier, Vivek F Farias, Jonas Oddur Jonasson, and Erez Yoeli. Policy optimization for personalized interventions in behavioral health. arXiv preprint arXiv:2303.12206, 2023.
- [6] Santiago R. Balseiro, David B. Brown, and Chen Chen. Dynamic pricing of relocating resources in large networks. *Management Science*, 67(7):4075–4094, 2021.
- [7] Santiago R. Balseiro, Haihao Lu, and Vahab Mirrokni. The best of many worlds: Dual mirror descent for online allocation problems. *Operations Research*, 71(1):101–119, 2023.
- [8] Dimitri Bertsekas and John N Tsitsiklis. Neuro-dynamic programming. Athena Scientific, 1996.
- [9] Jalaj Bhandari, Daniel Russo, and Raghav Singal. A finite time analysis of temporal difference learning with linear function approximation. *Operations Research*, 69(3):950–973, 2021.
- [10] Arpita Biswas, Gaurav Aggarwal, Pradeep Varakantham, and Milind Tambe. Learn to intervene: An adaptive learning policy for restless bandits in application to preventive healthcare. In *Proceedings of the 30th International Joint Conference on Artificial Intelligence (IJCAI 2021)*, pages 4036–4049, 2021.
- [11] David B Brown and Jingwei Zhang. Dynamic programs with shared resources and signals: Dynamic fluid policies and asymptotic optimality. *Operations Research*, 70(5):3015–3033, 2022.
- [12] David B. Brown and Jingwei Zhang. Technical note—on the strength of relaxations of weakly coupled stochastic dynamic programs. *Operations Research*, 71(6):2374–2389, 2023. doi: 10.1287/opre.2022.2287.
- [13] David B. Brown and Jingwei Zhang. Fluid policies, reoptimization, and performance guarantees in dynamic resource allocation. *Operations Research*, 73(2):1029–1045, 2025.
- [14] Qi Cai, Zhuoran Yang, Jason D Lee, and Zhaoran Wang. Neural temporal-difference learning converges to global optima. In Advances in Neural Information Processing Systems, volume 32, 2019.
- [15] Shao-Hung Chan, Zhe Chen, Teng Guo, Han Zhang, Yue Zhang, Daniel Harabor, Sven Koenig, Cathy Wu, and Jingjin Yu. The league of robot runners competition: Goals, designs, and implementation. In *ICAPS 2024 System's Demonstration track*, 2024.
- [16] Zehao Dou, Jakub Grudzien Kuba, and Yaodong Yang. Understanding value decomposition algorithms in deep cooperative multi-agent reinforcement learning. *arXiv preprint arXiv:2202.04868*, 2022.
- [17] Vivek Farias, Hao Li, Tianyi Peng, Xinyuyang Ren, Huawei Zhang, and Andrew Zheng. Correcting for interference in experiments: A case study at douyin. In *Proceedings of the 17th ACM Conference on Recommender Systems*, pages 455–466, 2023.
- [18] Jing Fu, Yoni Nazarathy, Sarat Moka, and Peter G. Taylor. Towards q-learning the whittle index for restless bandits. In 2019 Australian New Zealand Control Conference (ANZCC), 2019.
- [19] Nicolas Gast, Bruno Gaujal, and Chen Yan. Reoptimization nearly solves weakly coupled markov decision processes. *arXiv preprint arXiv:2211.01961*, 2022.

- [20] Nicolas Gast, Bruno Gaujal, and Chen Yan. Exponential asymptotic optimality of whittle index policy. *Queueing Syst. Theory Appl.*, 104(1–2):107–150, may 2023.
- [21] Nicolas Gast, Bruno Gaujal, and Chen Yan. Linear program-based policies for restless bandits: Necessary and sufficient conditions for (exponentially fast) asymptotic optimality. *Mathematics of Operations Research*, 49(4):2468–2491, 2024.
- [22] Carlos Guestrin, Daphne Koller, and Ronald Parr. Multiagent planning with factored mdps. In *Advances in Neural Information Processing Systems*, volume 14. MIT Press, 2001.
- [23] Carlos Guestrin, Daphne Koller, Ronald Parr, and Shobha Venkataraman. Efficient solution algorithms for factored mdps. *Journal of Artificial Intelligence Research*, 19(1):399–468, 2003.
- [24] Benjamin Han, Hyungjun Lee, and Sébastien Martin. Real-time rideshare driver supply values using online reinforcement learning. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, KDD '22, page 2968–2976, 2022.
- [25] Yitian Hong, Yaochu Jin, and Yang Tang. Rethinking individual global max in cooperative multi-agent reinforcement learning. *Advances in neural information processing systems*, 35: 32438–32449, 2022.
- [26] Natasha Jaques, Angeliki Lazaridou, Edward Hughes, Caglar Gulcehre, Pedro Ortega, Dj Strouse, Joel Z. Leibo, and Nando De Freitas. Social influence as intrinsic motivation for multi-agent deep reinforcement learning. In Kamalika Chaudhuri and Ruslan Salakhutdinov, editors, *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 3040–3049. PMLR, 09–15 Jun 2019.
- [27] Igor Kadota, Elif Uysal-Biyikoglu, Rahul Singh, and Eytan Modiano. Minimizing the age of information in broadcast wireless networks. In 2016 54th Annual Allerton Conference on Communication, Control, and Computing (Allerton), pages 844–851. IEEE, 2016.
- [28] Yash Kanoria and Pengyu Qian. Blind dynamic resource allocation in closed networks via mirror backpressure. *Management Science*, 70(8):5445–5462, 2024.
- [29] Anuj Mahajan, Tabish Rashid, Mikayel Samvelyan, and Shimon Whiteson. Maven: Multi-agent variational exploration. *Advances in neural information processing systems*, 32, 2019.
- [30] Khaled Nakhleh, Santosh Ganji, Ping-Chun Hsieh, I-Hong Hou, and Srinivas Shakkottai. Neurwin: Neural whittle index network for restless bandits via deep rl. In *Advances in Neural Information Processing Systems*, volume 34, pages 828–839, 2021.
- [31] Michael A Nielsen and Isaac L Chuang. *Quantum computation and quantum information*. Cambridge university press, 2010.
- [32] Ian Osband and Benjamin Van Roy. Near-optimal reinforcement learning in factored mdps. In *Proceedings of the 28th International Conference on Neural Information Processing Systems*, NIPS'14, page 604–612, 2014.
- [33] Zhiwei (Tony) Qin, Xiaocheng Tang, Yan Jiao, Fan Zhang, Zhe Xu, Hongtu Zhu, and Jieping Ye. Ride-hailing order dispatching at didi via reinforcement learning. *INFORMS Journal on Applied Analytics*, 50(5):272–286, 2020. doi: 10.1287/inte.2020.1047.
- [34] Naveen Janaki Raman, Zheyuan Ryan Shi, and Fei Fang. Global rewards in restless multi-armed bandits. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024.
- [35] Tabish Rashid, Mikayel Samvelyan, Christian Schroeder De Witt, Gregory Farquhar, Jakob Foerster, and Shimon Whiteson. Monotonic value function factorisation for deep multi-agent reinforcement learning. *Journal of Machine Learning Research*, 21(178):1–51, 2020.
- [36] Ibrahim El Shar and Daniel R. Jiang. Weakly coupled deep q-networks. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.

- [37] Kyunghwan Son, Daewoo Kim, Wan Ju Kang, David Earl Hostallero, and Yung Yi. QTRAN: Learning to factorize with transformation for cooperative multi-agent reinforcement learning. In *Proceedings of the 36th International Conference on Machine Learning*, volume 97, pages 5887–5896. PMLR, 2019.
- [38] Peter Sunehag, Guy Lever, Audrunas Gruslys, Wojciech Marian Czarnecki, Vinicius Zambaldi, Max Jaderberg, Marc Lanctot, Nicolas Sonnerat, Joel Z. Leibo, Karl Tuyls, and Thore Graepel. Value-decomposition networks for cooperative multi-agent learning based on team reward. In Proceedings of the 17th International Conference on Autonomous Agents and MultiAgent Systems, AAMAS '18, page 2085–2087, 2018.
- [39] Richard S. Sutton and Andrew G. Barto. Reinforcement Learning: An Introduction. The MIT Press, 2018. ISBN 0262039249.
- [40] John Tsitsiklis and Benjamin Van Roy. Analysis of temporal-diffference learning with function approximation. In Advances in Neural Information Processing Systems, volume 9. MIT Press, 1996.
- [41] Ina Maria Verloop. Asymptotically optimal priority policies for indexable and nonindexable restless bandits. Annals of Applied Probability, 26:1947–1995, 2016.
- [42] Jianhao Wang, Zhizhou Ren, Terry Liu, Yang Yu, and Chongjie Zhang. Qplex: Duplex dueling multi-agent q-learning. *arXiv preprint arXiv:2008.01062*, 2020.
- [43] Jianhao Wang, Zhizhou Ren, Beining Han, Jianing Ye, and Chongjie Zhang. Towards understanding cooperative multi-agent q-learning with value factorization. *Advances in Neural Information Processing Systems*, 34:29142–29155, 2021.
- [44] Kai Wang, Lily Xu, Aparna Taneja, and Milind Tambe. Optimistic whittle index policy: Online learning for restless bandits. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, pages 10131–10139, 2023.
- [45] Tonghan Wang, Jianhao Wang, Yi Wu, and Chongjie Zhang. Influence-based multi-agent exploration. *arXiv preprint arXiv:1910.05512*, 2019.
- [46] Richard Weber. On the gittins index for multiarmed bandits. *The Annals of Applied Probability*, 2(4):1024 1033, 1992.
- [47] Richard R. Weber and Gideon Weiss. On an index policy for restless bandits. *Journal of Applied Probability*, 27(3):637–648, 1990.
- [48] Richard R. Weber and Gideon Weiss. Addendum to 'on an index policy for restless bandits'. *Advances in Applied Probability*, 23(2):429–430, 1991.
- [49] P. Whittle. Restless bandits: Activity allocation in a changing world. *Journal of Applied Probability*, 25:287–298, 1988.
- [50] Dan Zhang and Daniel Adelman. An approximate dynamic programming approach to network revenue management with customer choice. *Transportation Science*, 43(3):381–394, 2009.
- [51] Xiangyu Zhang and Peter I Frazier. Restless bandits with many arms: Beating the central limit theorem. *arXiv preprint arXiv:2107.11911*, 2021.
- [52] Xiangyu Zhang and Peter I Frazier. Near-optimality for infinite-horizon restless bandits with many arms. *arXiv preprint arXiv:2203.15853*, 2022.

Contents

1	Introduction			
	1.1 This paper	2		
	1.2 Other related work	3		
2	Model	4		
	2.1 Local (Q-)value functions	۷		
3	Exact value decomposition	5		
	3.1 Necessary condition for the exact value decomposition	6		
4	Value decomposition error in general two-agent MDPs	7		
	4.1 Entry-wise error bound	7		
	4.2 Error weighted by stationary distribution	8		
	4.3 Multi-agent Markov entanglement	8		
5	Applications of Markov Entanglement	8		
	5.1 Efficient verification of value decomposition	ç		
6	Discussions			
7	Conclusion			
A	Linear algebra with tensor product			
В	B Decompose value functions			
C	C Comparison with quantum entanglement			
D	Proof of Theorem 2			
E	Decomposition via general functions			
F	Proof of Theorem 3			
G	Proof of Theorem 4			
H	Results for multi-agent MDPs			
I	Proof of Theorem 6			
J	Extensions of Markov entanglement	27		
	J.1 (Weakly-)coupled MDPs	27		
	J.2 Coupled MDPs with exogenous information	28		
	I 3 Factored MDPs	30		

	J.4	Fully cooperative Markov games	31
K	Sim	ulation environments	31
	K .1	Monte-Carlo estimation of Markov entanglement	32
	K.2	Learning local Q-values	32
	K .3	Sample Complexity and Computation	33

A Linear algebra with tensor product

We briefly introduce the basic properties of tensor product or Kronecker product. Let $A \in \mathbb{R}^{m_1 \times n_1}$, $B \in \mathbb{R}^{m_2 \times n_2}$, then

$$m{A} \otimes m{B} = \left[egin{array}{cccc} a_{11} m{B} & a_{12} m{B} & \dots & a_{1n_1} m{B} \\ a_{21} m{B} & a_{22} m{B} & \dots & a_{2n_1} m{B} \\ \dots & \dots & \dots & \dots \\ a_{m_11} m{B} & a_{m_12} m{B} & \dots & a_{m_1n_1} m{B} \end{array}
ight] \in \mathbb{R}^{m_1 m_2 \times n_1 n_2} \, .$$

Tensor product satisfies the following basic properties.

- 1. Bilinearity For any matrix A, B, C and constant k, it holds $k(A \otimes B) = (kA) \otimes B = A \otimes (kB), (A+B) \otimes C = A \otimes C + B \otimes C$, and $A \otimes (B+C) = A \otimes B + A \otimes C$.
- 2. Mixed-product Property For any matrix A, B, C, D, if AC and BD form valid matrix product, then $(A \otimes B)(C \otimes D) = (AC) \otimes (BD)$.

B Decompose value functions

Compared to the decomposition of Q-value, the value function further requires the reward to be *state-dependent*. To illustrate, notice by Bellman equation,

$$V_{AB}^{\pi} = (\boldsymbol{I} - \gamma \boldsymbol{P}_{AB}^{\pi})^{-1} \boldsymbol{r}_{AB}^{\pi} ,$$

where we abuse notation and denote $P_{AB}^{\pi}(s'\mid s) = \sum_{a} \pi(a\mid s) P(s'\mid s, a)$ and reward $r_{AB}^{\pi}(s) = \sum_{a} \pi(a\mid s) r_{AB}(s, a)$. A key subtlety arises because r_{AB}^{π} may not be decomposable—even when r_{AB} is decomposable—unless the reward r_{AB} is state-dependent. Consequently, we cannot directly apply the "absorbing" equation as in the proof of Theorem 1.

On the other hand, Q-value decomposition bypasses the state-dependence assumption and provides a stronger condition that directly implies value function decomposition. As a result, while learning local value functions may seem more intuitive, we recommend learning local Q-values instead and using them to approximate the global value function.

C Comparison with quantum entanglement

It turns out that our Markov entanglement condition serves as a mathematical counterpart of quantum entanglement in quantum physics. We provide the formal definition of the latter for comparison.

Definition 4 (Two-party Quantum Entanglement). Consider a two-party quantum system composed of two subsystems A and B. The joint state ρ_{AB} is **separable** if there exists $K \in \mathbb{Z}^+$, a probability measure $\{x_j\}_{j \in [K]}$, and density matrices $\{\rho_A^{(j)}, \rho_B^{(j)}\}_{j \in [K]}$ such that

$$\rho_{AB} = \sum_{j=1}^{K} x_j \rho_A^{(j)} \otimes \rho_B^{(j)}.$$

If there exists no such decomposition, ρ_{AB} is entangled.

The density matrices are square matrices satisfying certain properties such as positive semidefiniteness and trace normalization, which can be viewed as the counterparts of transition matrices in the Markov world. Despite the similarities in mathematical form, quantum entanglement imposes an additional constraint requiring $\{x_j\}_{j\in[K]}$ to be a probability measure, i.e. $x\geq 0$. In contrast, our Markov entanglement defined in Definition 1 permits general linear coefficients $\{x_j\}_{j\in[K]}$ as long as $\sum_{j=1}^k x_j = 1$. This distinction raises the important question of whether negative coefficients are indeed necessary in characterizing Markov entanglement.

To start with, we introduce the set of all separable transition matrices

$$\mathcal{P}_{\text{SEP}} = \left\{ oldsymbol{P} \geq 0 \; \middle| \; oldsymbol{P} = \sum_{j=1}^K x_j oldsymbol{P}_A^{(j)} \otimes oldsymbol{P}_B^{(j)} \; , \; \sum_{j=1}^K x_j = 1
ight\} \; ,$$

where $K \in \mathbb{Z}^+$ and $\left\{ \boldsymbol{P}_A^{(j)}, \boldsymbol{P}_B^{(j)} \right\}_{j \in [K]}$ are transition matrices. $\boldsymbol{P} \geq 0$ calls for every element of \mathcal{P}_{SEP} to be a valid transition matrix. It's clear that a transition matrix $\boldsymbol{P}_{AB}^{\pi}$ is separable if and only if $\boldsymbol{P}_{AB}^{\pi} \in \mathcal{P}_{\text{SEP}}$. On the other hand, a direct analogy of quantum entanglement gives us the following set that further requires non-negative coefficients,

$$\mathcal{P}_{\text{SEP}}^+ = \left\{ oldsymbol{P} \geq 0 \;\middle|\; oldsymbol{P} = \sum_{j=1}^K x_j oldsymbol{P}_A^{(j)} \otimes oldsymbol{P}_B^{(j)} \;,\; \sum_{j=1}^K x_j = 1 \;,\; oldsymbol{x} \geq 0
ight\} \;.$$

Interestingly, it turns out $\mathcal{P}_{SEP}^+ \nsubseteq \mathcal{P}_{SEP}$. In other words, there exist separable two-agent MDPs that can only be represented by linear combinations but not convex combinations of independent subsystems. Specifically, consider the following basis

$$m{E}_{00} = \left(egin{array}{cc} 1 & 0 \\ 1 & 0 \end{array}
ight), \quad m{E}_{01} = \left(egin{array}{cc} 1 & 0 \\ 0 & 1 \end{array}
ight), \quad m{E}_{10} = \left(egin{array}{cc} 0 & 1 \\ 1 & 0 \end{array}
ight), \quad m{E}_{11} = \left(egin{array}{cc} 0 & 1 \\ 0 & 1 \end{array}
ight)$$

And the corresponding transition matrix we provide is

$$m{P} = \left(egin{array}{cccc} 0.5 & 0 & 0 & 0.5 \ 0.5 & 0 & 0 & 0.5 \ 0.5 & 0 & 0 & 0.5 \ 0 & 0.5 & 0.5 & 0 \end{array}
ight) = rac{1}{2}m{E}_{00} \otimes m{E}_{00} + rac{1}{2}m{E}_{10} \otimes m{E}_{11} + rac{1}{2}m{E}_{11} \otimes m{E}_{10} - rac{1}{2}m{E}_{10} \otimes m{E}_{10}$$

One can also verify P can not be represented by the convex combination of tensor products of these basis. This result justifies the necessity of negative coefficients in x and highlights a structural difference between Markov entanglement and quantum entanglement

D Proof of Theorem 2

We provide the full proof of Theorem 2 in this section.

Step 1: Characterize the Orthogonal Complement. To start with, we consider the smallest subspace containing all transition matrices $\Omega_P := \operatorname{span}(P)$ where P are the set of all transition matrices in $\mathbb{R}^{m \times m}$. We then study the dimension of Ω_P .

Lemma 1. The dimension of Ω_P is $\dim(\Omega_P) = m^2 - m + 1$.

Proof. Let $Z_{ij} \in \mathbb{R}^{m \times m}$ such that

$$\boldsymbol{Z}_{ij}(a,b) = \begin{cases} 1 & (a=i \wedge b=j) \vee (a=b) \\ 0 & o.w. \end{cases}.$$

One basis for all transition matrices is given by $\{Z_{ij}\}_{i,j\in[m]}$ whose cardinarlity is m^2-m+1 . \square

Let $\Omega_{P^{\otimes 2}} \coloneqq \operatorname{span}(P_1 \otimes P_2)$ be the minimal subspace containing all separable transition matrices. It quickly follows that

$$\dim(\Omega_{P^{\otimes 2}}) = (\dim(\Omega_P))^2.$$

We then construct the orthogonal complement of $\Omega_{P^{\otimes 2}}$ under Frobenius inner product. Let $\{\varepsilon_j\}_{j\in[m-1]}$ be a set of vector in \mathbb{R}^m such that $\varepsilon_j=(1,0,\ldots,0,-1,0,\ldots,0)^{\top}$ with the first element 1 and j+1-th element -1. Notice that

$$\operatorname{Tr}\left(\boldsymbol{e}\varepsilon_{j}^{\top}\boldsymbol{P}\right) = \operatorname{Tr}\left(\varepsilon_{j}^{\top}\boldsymbol{P}\boldsymbol{e}\right) = 0,$$

for all ε_i . Consider the following subspace

$$\Omega' = \left\{ \sum_{j=1}^{m-1} \left(arepsilon_j oldsymbol{e}^{ op}
ight) \otimes oldsymbol{W}_j^1 + \sum_{j=1}^{m-1} oldsymbol{W}_j^2 \otimes \left(arepsilon_j oldsymbol{e}^{ op}
ight) \mid W_{1:j}^1, W_{1:j}^2 \in \mathbb{R}^{m imes m}
ight\} \,.$$

We then show Ω' is exactly the orthogonal complement of $\Omega_{P^{\otimes 2}}$. First, notice that

$$\dim(\Omega') = 2(m-1)m^2 - (m-1)^2$$
.

and thus $\dim(\Omega') + \dim(\Omega_{P^{\otimes 2}}) = m^4$. Moreover, one can verify for any $X \in \Omega_{P^{\otimes 2}}$ and $Y \in \Omega'$, $\operatorname{Tr}(X^\top Y) = 0$. As a result, it holds

$$\Omega' = \Omega_{P^{\otimes 2}}^{\perp}$$
.

Step 2: Connection to "Inverse" The decomposition of Q-value ultimately concerns with the properties of $(I - \gamma P_{AB}^{\pi})^{-1}$. The following lemma bridges this gap.

Lemma 2. Given any transition matrix P and $\gamma > 0$, P is separable if and only if $(1-\gamma)(I-\gamma P)^{-1}$ is separable.

Proof. (\Rightarrow) One can verify that $(\boldsymbol{I} - \gamma \boldsymbol{P})\boldsymbol{e} = (1 - \gamma)\boldsymbol{e}$, which implies $(1 - \gamma)(\boldsymbol{I} - \gamma \boldsymbol{P})^{-1}$ is a transition matrix. Moreover, $(1 - \gamma)(\boldsymbol{I} - \gamma \boldsymbol{P})^{-1} = (1 - \gamma)\sum_{i=0}^{\infty}(\gamma \boldsymbol{P})^{i}$ falls in $\Omega_{P^{\otimes 2}}$ as $\boldsymbol{P} \in \Omega_{P^{\otimes 2}}$.

 (\Leftarrow) This side is more involved. Denote $U \coloneqq (1-\gamma)(I-\gamma P)^{-1}$. Then if the spectral radius $\rho(I-U) < 1$, then

$$U^{-1} = (I - (I - U))^{-1} = \sum_{i=0}^{\infty} (I - U)^i \in \Omega_{P^{\otimes 2}}.$$

This implies $U^{-1} = \frac{1}{1-\gamma}(I - \gamma P) \in \Omega_{P^{\otimes 2}}$ and thus $P \in \Omega_{P^{\otimes 2}}$, finishing the proof. It then suffices to show $\rho(I - U) < 1$. Notice that

$$\lambda_i(\mathbf{I} - \mathbf{U}) = 1 - \lambda_i(\mathbf{U}) = 1 - \frac{1 - \gamma}{\lambda(\mathbf{I} - \gamma \mathbf{P})} = 1 - \frac{1 - \gamma}{1 - \gamma \lambda_i(\mathbf{P})}$$

Let $\lambda_i(\mathbf{P}) = a + bi$ and taking modulus for both side

$$\begin{aligned} |\lambda_i(\boldsymbol{I} - \boldsymbol{U})| &= \left| \frac{\gamma - \gamma \lambda_i(\boldsymbol{P})}{1 - \gamma \lambda_i(\boldsymbol{P})} \right| \\ &= \frac{|\gamma - \gamma \lambda_i(\boldsymbol{P})|}{|1 - \gamma \lambda_i(\boldsymbol{P})|} \\ &= \sqrt{\frac{\gamma^2 (1 - a)^2 + \gamma^2 b^2}{(1 - \gamma a)^2 + \gamma^2 b^2}} \\ &= \sqrt{1 + \frac{(1 - \gamma)(2a\gamma - \gamma - 1)}{(1 - \gamma a)^2 + \gamma^2 b^2}} \\ &\leq \sqrt{1 - \frac{(1 - \gamma)^2}{(1 - \gamma a)^2 + \gamma^2 b^2}} < 1 \, . \end{aligned}$$

We conclude the proof given $\rho(I - U) = \max_i |\lambda_i(I - U)| < 1$.

Step 3: Put it together By Lemma 2, if P_{AB}^{π} is entangled, then $(1-\gamma)(I-\gamma P_{AB}^{\pi})^{-1}$ is also entangled. Then there exists $\mathbf{Y} \in \Omega' \neq \mathbf{0}$ such that $\mathrm{Tr}(\mathbf{Y}^{\top}(I-\gamma P_{AB}^{\pi})^{-1}) \neq 0$. We apply singular value decomposition to all $W_{1:j}^1, W_{1:j}^2$ and conclude there exists some j and $\mathbf{u}, \mathbf{v} \in \mathbb{R}^m$ such that either $\mathrm{Tr}((e\varepsilon_j^{\top}) \otimes (v\mathbf{u}^{\top}) (I-\gamma P_{AB}^{\pi})^{-1}) \neq 0$ or $\mathrm{Tr}((v\mathbf{u}^{\top}) \otimes (e\varepsilon_j^{\top}) (I-\gamma P_{AB}^{\pi})^{-1}) \neq 0$. We assume the former without loss of generality, it holds

$$(\varepsilon_i^{\top} \otimes \boldsymbol{u}^{\top})(\boldsymbol{I} - \gamma \boldsymbol{P}_{AB}^{\pi})^{-1}(\boldsymbol{e} \otimes \boldsymbol{v}) \neq 0.$$

Now set $r_A = 0$ and $r_B = v$. Since Q_{AB}^{π} is decomposable, there exists some local function Q_A, Q_B such that

$$(I - \gamma P_{AB}^{\pi})^{-1}(e \otimes v) = Q_A(\mathbf{0}) \otimes e + e \otimes Q_B(v).$$

Left multiply by $(\varepsilon_i^\top \otimes \boldsymbol{u}^\top)$, we have

$$(\varepsilon_i^{\top} \otimes \boldsymbol{u}^{\top})(\boldsymbol{I} - \gamma \boldsymbol{P}_{AB}^{\pi})^{-1}(\boldsymbol{e} \otimes \boldsymbol{v}) = (\varepsilon_i^{\top} \otimes \boldsymbol{u}^{\top})(Q_A(\boldsymbol{0}) \otimes \boldsymbol{e}) \neq 0,$$

Then set $r_A = 0$ and $r_B = -v$, we can similarly derive

$$-(\varepsilon_j^\top \otimes \boldsymbol{u}^\top)(\boldsymbol{I} - \gamma \boldsymbol{P}_{AB}^\pi)^{-1}(\boldsymbol{e} \otimes \boldsymbol{v}) = (\varepsilon_j^\top \otimes \boldsymbol{u}^\top)(Q_A(\boldsymbol{0}) \otimes \boldsymbol{e}) \neq 0,$$

This gives use $(\varepsilon_j^\top \otimes \boldsymbol{u}^\top)(Q_A(\boldsymbol{0}) \otimes \boldsymbol{e}) = 0$, which is a contradiction.

E Decomposition via general functions

Entangled ${m P}$ precludes the local decomposition with local value functions, but may admit decompositions with more general functions. Consider ${m P}=\frac{1}{4}\left(ee^{\top}\right)\otimes\left(ee^{\top}\right)+\delta\left(\epsilon e^{\top}\right)\otimes\left(e\epsilon^{\top}\right)$, where $e=[1,1],\epsilon=[1-1]$. Clearly such ${m P}$ is entangled. We also have ${m P}^k=\frac{1}{4}\left(ee^{\top}\right)\otimes\left(ee^{\top}\right)$, for $k\geq 2$. Then $(I-\gamma P)^{-1}={m I}+\frac{\gamma+\gamma^2}{4}\left(ee^{\top}\right)\otimes\left(ee^{\top}\right)+\delta\gamma\left(\epsilon e^{\top}\right)\otimes\left(e\epsilon^{\top}\right)$. Then for any ${m r}_A,{m r}_B$, we have $({m I}-\gamma {m P})^{-1}\left({m r}_A\otimes e+e\otimes {m r}_B\right)={m r}_A\otimes e+h_A\left(\gamma+\gamma^2\right)/2e\otimes e+{m r}_B\otimes e+h_B\left(\gamma+\gamma^2\right)/2e\otimes e+2\delta\gamma\left(\epsilon^{\top}{m r}_B\right)$ $\epsilon\otimes e$ where $h_A=e^{\top}{m r}_A,h_B=e^{\top}{m r}_B$.

F Proof of Theorem 3

Additional Notations For (semi-)norm $\|\cdot\|_{\alpha}$ and norm $\|\cdot\|_{\beta}$, we define the α, β -norm for matrix A as

$$\|\boldsymbol{A}\|_{lpha,eta} = \sup_{\|\boldsymbol{x}\|_{eta}=1} \|\boldsymbol{A}\boldsymbol{x}\|_{lpha}.$$

We further abbreviate $\|A\|_{\alpha} := \|A\|_{\alpha,\alpha}$. Moreover, we define the operator |x| taking the absolute value of each element of vector or matrix x.

To prove the theorem, we introduce the key technique of analyzing perturbation bounds of the transition matrix, which is also used in [17].

Lemma 3 (Lemma 1 in [17]). Let $P, P' \in \mathbb{R}^{n \times n}$ such that $(I - P)^{-1}$ and $(I - P')^{-1}$ exist. Then it holds

$$(I - P')^{-1} = (I - P)^{-1} + (I - P')^{-1}(P' - P)(I - P)^{-1}$$
.

We are then ready to prove the main theorem.

Proof of Theorem 3. Let P_A , P_B be the optimal solution to Eq. (7) w.r.t agent A, B. For any subset of state-action pairs of agent A, $\mathcal{F} \subseteq \mathcal{S}_A \times \mathcal{A}_A$, we have

$$\begin{vmatrix} \sum_{s'_A, a'_A \in \mathcal{F}} (\mathbf{P}_A^{\pi} - \mathbf{P}_A)_{(s'_A, a'_A \mid s_A, a_A)} \end{vmatrix}$$

$$= \begin{vmatrix} \sum_{s'_A, a'_A \in \mathcal{F}} \sum_{s'_B, a'_B} \sum_{s_B, a_B} (\mathbf{P}_{AB}^{\pi} - \mathbf{P}_A \otimes \mathbf{P}_B)_{(s', \mathbf{a}' \mid s, \mathbf{a})} \mu_{AB}^{\pi}(s_B, a_B \mid s_A, a_A) \end{vmatrix}$$

$$\leq \sum_{s_B, a_B} \begin{vmatrix} \sum_{s'_A, a'_A \in \mathcal{F}} \sum_{s'_B, a'_B} (\mathbf{P}_{AB}^{\pi} - \mathbf{P}_A \otimes \mathbf{P}_B)_{(s', \mathbf{a}' \mid s, \mathbf{a})} \end{pmatrix} \mu_{AB}^{\pi}(s_B, a_B \mid s_A, a_A)$$

$$\leq \sum_{s_B, a_B} E_A(\mathbf{P}_{AB}^{\pi}) \mu_{AB}^{\pi}(s_B, a_B \mid s_A, a_A) = E_A(\mathbf{P}_{AB}^{\pi})$$

where the last inequality follows from the definition of agent-wise total variation distance. Since the result holds for any \mathcal{F} and $(s_A, a_A) \in \mathcal{S}_A \times \mathcal{A}_A$, we have

$$\|\boldsymbol{P}_{A}^{\pi}-\boldsymbol{P}_{A}\|_{\mathrm{TV}}\leq E_{A}(\boldsymbol{P}_{AB}^{\pi}),$$

and similar results hold for P_B^{π} .

Next we have

$$(I - \gamma P_{AB}^{\pi})^{-1} (r_A \otimes e) - ((I - \gamma P_A^{\pi})^{-1} r_A) \otimes e$$

$$= (I - \gamma P_{AB}^{\pi})^{-1} (r_A \otimes e) - (I - \gamma P_A \otimes P_B)^{-1} (r_A \otimes e)$$

$$+ (I - \gamma P_A \otimes P_B)^{-1} (r_A \otimes e) - ((I - \gamma P_A^{\pi})^{-1} r_A) \otimes e$$

$$\stackrel{(i)}{=} \underbrace{(I - \gamma P_{AB}^{\pi})^{-1} (r_A \otimes e) - (I - \gamma P_A \otimes P_B)^{-1} (r_A \otimes e)}_{(I)}$$

$$+ \underbrace{((I - \gamma P_A)^{-1} r_A) \otimes e - ((I - \gamma P_A^{\pi})^{-1} r_A) \otimes e}_{(II)}$$

where (i) also follows the same "absorbing" technique in the proof of Theorem 1. For (I), apply Lemma 3, it holds

$$\begin{aligned} & \left\| (\boldsymbol{I} - \gamma \boldsymbol{P}_{AB}^{\pi})^{-1} (\boldsymbol{r}_{A} \otimes \boldsymbol{e}) - (\boldsymbol{I} - \gamma \boldsymbol{P}_{A} \otimes \boldsymbol{P}_{B})^{-1} (\boldsymbol{r}_{A} \otimes \boldsymbol{e}) \right\|_{\infty} \\ &= \left\| (\boldsymbol{I} - \gamma \boldsymbol{P}_{AB}^{\pi})^{-1} (\gamma \boldsymbol{P}_{AB}^{\pi} - \gamma \boldsymbol{P}_{A} \otimes \boldsymbol{P}_{B}) (\boldsymbol{I} - \gamma \boldsymbol{P}_{A} \otimes \boldsymbol{P}_{B})^{-1} (\boldsymbol{r}_{A} \otimes \boldsymbol{e}) \right\|_{\infty} \\ &\leq \left\| (\boldsymbol{I} - \gamma \boldsymbol{P}_{AB}^{\pi})^{-1} \right\|_{\infty} \left\| (\gamma \boldsymbol{P}_{AB}^{\pi} - \gamma \boldsymbol{P}_{A} \otimes \boldsymbol{P}_{B}) \left((\boldsymbol{I} - \gamma \boldsymbol{P}_{A})^{-1} \boldsymbol{r}_{A} \right) \otimes \boldsymbol{e} \right\|_{\infty} \\ &\leq \left\| (\boldsymbol{I} - \gamma \boldsymbol{P}_{AB}^{\pi})^{-1} \right\|_{\infty} 2\gamma E_{A} (\boldsymbol{P}_{AB}^{\pi}) \left\| (\boldsymbol{I} - \gamma \boldsymbol{P}_{A})^{-1} \boldsymbol{r}_{A} \right\|_{\infty} \\ &\leq \frac{2\gamma E_{A} (\boldsymbol{P}_{AB}^{\pi}) r_{\max}^{A}}{1 - \gamma} \left\| (\boldsymbol{I} - \gamma \boldsymbol{P}_{AB}^{\pi})^{-1} \right\|_{\infty} \leq \frac{2\gamma E_{A} (\boldsymbol{P}_{AB}^{\pi}) r_{\max}^{A}}{(1 - \gamma)^{2}}, \end{aligned}$$

where (i) follows by the definition of agent-wise total variation distance when $\|\mathbf{r}_A\|_{\infty} \neq 0$, and also trivially hold when $\|\mathbf{r}_A\|_{\infty} = 0$. Similarly, for (II) we have

$$\begin{aligned} & \left\| \left(\left(\boldsymbol{I} - \gamma \boldsymbol{P}_{A} \right)^{-1} \boldsymbol{r}_{A} \right) \otimes \boldsymbol{e} - \left(\left(\boldsymbol{I} - \gamma \boldsymbol{P}_{A}^{\pi} \right)^{-1} \boldsymbol{r}_{A} \right) \otimes \boldsymbol{e} \right\|_{\infty} \\ &= \left\| \left(\left(\boldsymbol{I} - \gamma \boldsymbol{P}_{A} \right)^{-1} - \left(\boldsymbol{I} - \gamma \boldsymbol{P}_{A}^{\pi} \right)^{-1} \right) \boldsymbol{r}_{A} \right\|_{\infty} \\ &= \left\| \left(\boldsymbol{I} - \gamma \boldsymbol{P}_{A}^{\pi} \right)^{-1} \left(\gamma \boldsymbol{P}_{A}^{\pi} - \gamma \boldsymbol{P}_{A} \right) \left(\boldsymbol{I} - \gamma \boldsymbol{P}_{A} \right)^{-1} \boldsymbol{r}_{A} \right\|_{\infty} \\ &\leq \frac{2\gamma E_{A} (\boldsymbol{P}_{AB}^{\pi}) r_{\max}^{A}}{(1 - \gamma)^{2}} \, . \end{aligned}$$

Then we have

$$\left\| \left(\boldsymbol{I} - \gamma \boldsymbol{P}_{AB}^{\pi} \right)^{-1} \left(\boldsymbol{r}_A \otimes \boldsymbol{e} \right) - \left(\left(\boldsymbol{I} - \gamma \boldsymbol{P}_A^{\pi} \right)^{-1} \boldsymbol{r}_A \right) \otimes \boldsymbol{e} \right\|_{\infty} \leq \frac{4 \gamma E_A (\boldsymbol{P}_{AB}^{\pi}) r_{\max}^A}{(1 - \gamma)^2} \,.$$

We can derive similar results for agent B, i.e.,

$$\left\| \left(\boldsymbol{I} - \gamma \boldsymbol{P}_{AB}^{\pi} \right)^{-1} \left(\boldsymbol{e} \otimes \boldsymbol{r}_{B} \right) - \boldsymbol{e} \otimes \left(\left(\boldsymbol{I} - \gamma \boldsymbol{P}_{B}^{\pi} \right)^{-1} \boldsymbol{r}_{B} \right) \right\|_{\infty} \leq \frac{4 \gamma E_{B} (\boldsymbol{P}_{AB}^{\pi}) r_{\max}^{B}}{(1 - \gamma)^{2}}.$$

Put it all together we have

$$\left\|Q_{AB}^{\pi} - (Q_A^{\pi} \otimes \boldsymbol{e} + \boldsymbol{e} \otimes Q_B^{\pi})\right\|_{\infty} \leq \frac{4\gamma (E_A(\boldsymbol{P}_{AB}^{\pi})r_{\max}^A + E_B(\boldsymbol{P}_{AB}^{\pi})r_{\max}^B)}{(1-\gamma)^2}.$$

G Proof of Theorem 4

We first introduce the μ -weighted ATV distance Formally, we introduce the following norm.

Definition 5 (μ -norm). Given a transition matrix $P \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|} \times |\mathcal{S}||\mathcal{A}|$ with occupancy measure $\mu \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}$, for any vector $\mathbf{x} \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}$ the μ -norm is defined as

$$\|\boldsymbol{x}\|_{\mu} := \sum_{(s,a)\in\mathcal{S}\times\mathcal{A}} \mu(s,a) |x(s,a)| = \mu^{\top} |\boldsymbol{x}|.$$
 (10)

One can verify that μ -norm satisfies triangle inequality and is a valid norm when $\mu(s,a) > 0$ for all (s,a). Otherwise μ -norm is a *semi-norm* in general. We then introduce the distance

Definition 6 (μ -weighted Agent-wise Total Variation Distance). Given probability distribution $\mu \in \mathbb{R}^{|S|^2|A|^2}$, the μ -weighted total variation distance between two transition matrices $P, P' \in \mathbb{R}^{|S|^2|A|^2 \times |S|^2|A|^2}$ w.r.t agent A is defined as

$$\|\boldsymbol{P} - \boldsymbol{P}'\|_{\mu - \text{ATV}_{A}} = \frac{1}{2} \sup_{\|\boldsymbol{x}\|_{\infty} = 1} \|(\boldsymbol{P} - \boldsymbol{P}')(\boldsymbol{x} \otimes \boldsymbol{e})\|_{\mu}.$$

The μ -weighted ATV distance w.r.t agent B can be defined similarly. We claim that the μ -weighted ATV is also a counterpart of ATV distance in Definition 6. This follows from the constrained optimization formulation of ATV

$$\|\boldsymbol{P} - \boldsymbol{P}'\|_{\text{ATV}_{A}} = \frac{1}{2} \sup_{\|\boldsymbol{x}\|_{\infty} = 1} \|(\boldsymbol{P} - \boldsymbol{P}')(\boldsymbol{x} \otimes \boldsymbol{e})\|_{\infty}.$$
(11)

Thus μ -ATV substitutes μ -norm for the original ℓ_{∞} -norm. We plug μ -weighted ATV into Eq. (1) and obtain the corresponding measure of Markov entanglement $E(\boldsymbol{P}_{AB}^{\pi})$ and $E_A(\boldsymbol{P}_{AB}^{\pi})$. Similar to ATV in Eq. (7), this μ -weighted version of $E_A(\boldsymbol{P}_{AB}^{\pi})$ admits the following formulation

$$E_A(\mathbf{P}_{AB}^{\pi}) \leq \min_{\mathbf{P}_A} \sum_{\mathbf{s}, \mathbf{a}} \rho_{AB}^{\pi}(\mathbf{s}, \mathbf{a}) D_{\text{TV}} \Big(\mathbf{P}_{AB}^{\pi}(\cdot, \cdot \mid \mathbf{s}, \mathbf{a}), \mathbf{P}_A(\cdot, \cdot \mid s_A, a_A) \Big).$$
(12)

⁵Since $\mu \in \mathbb{R}^{|S||\mathcal{A}|}$ is the stationary distribution of $\mathbf{P} \in \mathbb{R}^{|S||\mathcal{A}| \times |S||\mathcal{A}|}$, we use "stationary distribution" and "occupancy measure" exchangeably when the context is clear.

This recovers Eq. (8) that substitutes the μ -weighted average for the maximum operator in Eq. (7). Thus intuitively, $E(\mathbf{P}_{AB}^{\pi})$ w.r.t μ -weighted ATV distance measures how closely agent A can be approximated as an independent subsystem under the stationary distribution.

We provide the proof for two agents here, one can easily generalize the proof to multi-agent scenarios. Compared to the proof of Theorem 3, this proof follows similar framework and differs in several details.

The first one is the following lemma for the "localized" stationary distribution

Lemma 4. P_A^{π} has stationary distribution μ_A^{π} with

$$\forall (s_A, a_A), \, \mu_A^{\pi}(s_A, a_A) = \sum_{s_B, a_B} \mu_{AB}^{\pi}(s_A, s_B, a_A, a_B).$$

In other words, the local stationary distribution of each agent is exactly the marginal distribution of global μ_{AB}^{π} .

Proof of Lemma 4. We proof by verify the definition of stationary distribution. For any (s'_A, a'_A) , it holds

$$\begin{split} &\sum_{s_{A},a_{A}} \left(\sum_{s_{B},a_{B}} \mu_{AB}^{\pi}(s_{A},s_{B},a_{A},a_{B}) \right) P^{\pi}(s_{A}',a_{A}' \mid s_{A},a_{A}) \\ &= \sum_{s_{A},a_{A}} \sum_{s_{B},a_{B}} \mu_{AB}^{\pi}(s_{A},s_{B},a_{A},a_{B}) \sum_{s_{B}',a_{B}'} \sum_{s_{B}'',a_{B}''} P^{\pi}\left(s_{A}',s_{B}',a_{A}',a_{B}' \mid s_{A},s_{B}'',a_{A}''\right) \mu_{AB}^{\pi}(s_{B}'',a_{B}'' \mid s_{A},a_{A}) \\ &= \sum_{s_{A},a_{A}} \sum_{s_{B},a_{B}} \mu_{AB}^{\pi}(s_{B},a_{B} \mid s_{A},a_{A}) \sum_{s_{B}',a_{B}'} \sum_{s_{B}'',a_{B}''} P^{\pi}\left(s_{A}',s_{B}',a_{A}',a_{B}' \mid s_{A},s_{B}'',a_{A},a_{B}''\right) \mu_{AB}^{\pi}(s_{A},s_{B}'',a_{A},a_{B}'') \\ &= \sum_{s_{A},a_{A}} \sum_{s_{B}',a_{B}'} \sum_{s_{B}'',a_{B}''} P^{\pi}\left(s_{A}',s_{B}',a_{A}',a_{B}' \mid s_{A},s_{B}'',a_{A},a_{B}''\right) \mu_{AB}^{\pi}(s_{A},s_{B}'',a_{A},a_{B}'') \\ &= \sum_{s_{A},a_{A}} \mu_{AB}^{\pi}(s_{A}',s_{B}',a_{A}',a_{B}') \\ &= \sum_{s_{B}',a_{B}'} \mu_{AB}^{\pi}(s_{A}',s_{B}',a_{A}',a_{A}') \\ &= \sum_{s_{B}',a_{B}'} \mu_{AB}^{\pi}(s_{A}',s_{A}',a_{A}',a_{A}') \\ &= \sum_{s_{B}',a_{B}'} \mu_{AB}^{\pi}(s_{A}',s_{A}',a_{A}',a_{A}') \\ &= \sum_{s_{$$

where the last equation follows from the definition of μ_{AB}^{π} . Hence we conclude that $\sum_{s_B,a_B} \mu_{AB}^{\pi}(s_A,s_B,a_A,a_B)$ is a stationary distribution of P_A^{π} .

We are then ready to prove Theorem 4. We first note that similar to ATV distance in Eq. (7), the optimal solution to $E_A(\boldsymbol{P}_{AB}^{\pi})$ w.r.t μ_{AB}^{π} -weighted ATV distance also only depends on \boldsymbol{P}_A . Thus, let $\boldsymbol{P}_A, \boldsymbol{P}_B$ be the optimal solutions to $E_A(\boldsymbol{P}_{AB}^{\pi}), E_B(\boldsymbol{P}_{AB}^{\pi})$ respectively.

Let $x \in \mathbb{R}^{|\mathcal{S}_A||\mathcal{A}_A|}$ with $||x||_{\infty} = 1$. Following the same technique in the proof of Theorem 4, we have

$$\mu_{A}^{\pi^{\top}} | (\mathbf{P}_{A}^{\pi} - \mathbf{P}_{A}) \mathbf{x} | \\
= \sum_{s_{A}, a_{A}} \mu_{A}^{\pi} (s_{A}, a_{A}) \left| \sum_{s'_{A}, a'_{A}} (\mathbf{P}_{A}^{\pi} - \mathbf{P}_{A})_{(s'_{A}, a'_{A} | s_{A}, a_{A})} \mathbf{x} (s'_{A}, a'_{A}) \right| \\
= \sum_{s_{A}, a_{A}} \mu_{A}^{\pi} (s_{A}, a_{A}) \left| \sum_{s'_{A}, a'_{A}} \mathbf{x} (s'_{A}, a'_{A}) \sum_{s'_{B}, a'_{B}} \sum_{s_{B}, a_{B}} (\mathbf{P}_{AB}^{\pi} - \mathbf{P}_{A} \otimes \mathbf{P}_{B})_{(\mathbf{s}', \mathbf{a}' | \mathbf{s}, \mathbf{a})} \mu_{AB}^{\pi} (s_{B}, a_{B} | s_{A}, a_{A}) \right| \\
\leq \sum_{s, \mathbf{a}} \left| \sum_{s'_{A}, a'_{A}} \mathbf{x} (s'_{A}, a'_{A}) \sum_{s'_{B}, a'_{B}} (\mathbf{P}_{AB}^{\pi} - \mathbf{P}_{A} \otimes \mathbf{P}_{B})_{(\mathbf{s}', \mathbf{a}' | \mathbf{s}, \mathbf{a})} \right| \mu_{AB}^{\pi} (s, \mathbf{a}) \leq 2E_{A} (\mathbf{P}_{AB}^{\pi})$$

where the second last inequality follows from Lemma 4. We then conclude

$$\|\boldsymbol{P}_{A}^{\pi}-\boldsymbol{P}_{A}\|_{\mu,\infty}\leq 2E_{A}(\boldsymbol{P}_{AB}^{\pi}),$$

and similar results hold for P_B^{π} . We then apply the decomposition

$$(\boldsymbol{I} - \gamma \boldsymbol{P}_{AB}^{\pi})^{-1} (\boldsymbol{r}_{A} \otimes \boldsymbol{e}) - ((\boldsymbol{I} - \gamma \boldsymbol{P}_{A}^{\pi})^{-1} \boldsymbol{r}_{A}) \otimes \boldsymbol{e}$$

$$= \underbrace{(\boldsymbol{I} - \gamma \boldsymbol{P}_{AB}^{\pi})^{-1} (\boldsymbol{r}_{A} \otimes \boldsymbol{e}) - (\boldsymbol{I} - \gamma \boldsymbol{P}_{A} \otimes \boldsymbol{P}_{B})^{-1} (\boldsymbol{r}_{A} \otimes \boldsymbol{e})}_{(I)} + \underbrace{((\boldsymbol{I} - \gamma \boldsymbol{P}_{A})^{-1} \boldsymbol{r}_{A}) \otimes \boldsymbol{e} - ((\boldsymbol{I} - \gamma \boldsymbol{P}_{A}^{\pi})^{-1} \boldsymbol{r}_{A}) \otimes \boldsymbol{e}}_{(II)}$$

For (I), we have

$$\begin{split} & \left\| (\boldsymbol{I} - \gamma \boldsymbol{P}_{AB}^{\pi})^{-1} (\boldsymbol{r}_{A} \otimes \boldsymbol{e}) - (\boldsymbol{I} - \gamma \boldsymbol{P}_{A} \otimes \boldsymbol{P}_{B})^{-1} (\boldsymbol{r}_{A} \otimes \boldsymbol{e}) \right\|_{\mu_{AB}^{\pi}} \\ &= \left\| (\boldsymbol{I} - \gamma \boldsymbol{P}_{AB}^{\pi})^{-1} (\gamma \boldsymbol{P}_{AB}^{\pi} - \gamma \boldsymbol{P}_{A} \otimes \boldsymbol{P}_{B}) (\boldsymbol{I} - \gamma \boldsymbol{P}_{A} \otimes \boldsymbol{P}_{B})^{-1} (\boldsymbol{r}_{A} \otimes \boldsymbol{e}) \right\|_{\mu_{AB}^{\pi}} \\ &\leq \frac{1}{1 - \gamma} \left\| \left((\gamma \boldsymbol{P}_{AB}^{\pi} - \gamma \boldsymbol{P}_{A} \otimes \boldsymbol{P}_{B}) (\boldsymbol{I} - \gamma \boldsymbol{P}_{A})^{-1} \boldsymbol{r}_{A} \right) \otimes \boldsymbol{e} \right\|_{\mu_{AB}^{\pi}} \\ &\leq \frac{2\gamma E(\pi)}{1 - \gamma} \left\| (\boldsymbol{I} - \gamma \boldsymbol{P}_{A})^{-1} \boldsymbol{r}_{A} \right\|_{\infty} \leq \frac{2\gamma E(\pi) r_{\text{max}}}{(1 - \gamma)^{2}}, \end{split}$$

where (i) follows from the fact that for any x

$$\|\boldsymbol{P}\boldsymbol{x}\|_{\mu} = \mu^{\top}|\boldsymbol{P}\boldsymbol{x}| \leq \mu^{\top}\boldsymbol{P}|\boldsymbol{x}| = \mu^{\top}|\boldsymbol{x}| = \|\boldsymbol{x}\|_{\mu}.$$

For (II) one can use Lemma 4 to verify

$$\begin{split} & \left\| \left(\left(\boldsymbol{I} - \gamma \boldsymbol{P}_{\!A} \right)^{-1} \boldsymbol{r}_{\!A} \right) \otimes \boldsymbol{e} - \left(\left(\boldsymbol{I} - \gamma \boldsymbol{P}_{\!A}^{\pi} \right)^{-1} \boldsymbol{r}_{\!A} \right) \otimes \boldsymbol{e} \right\|_{\mu_{AB}^{\pi}} \\ & = \left\| \left(\boldsymbol{I} - \gamma \boldsymbol{P}_{\!A} \right)^{-1} \boldsymbol{r}_{\!A} - \left(\boldsymbol{I} - \gamma \boldsymbol{P}_{\!A}^{\pi} \right)^{-1} \boldsymbol{r}_{\!A} \right\|_{\mu_{A}^{\pi}} \end{split}$$

And similar results to (I) holds. We then conclude the proof of Theorem 4.

H Results for multi-agent MDPs

In quantum physics, the concept of quantum entanglement of two-party system can be well extended to multi-party system. In this section, we demonstrate a similar extension of two-agent Markov entanglement to multi-agent settings. We begin with the model of multi-agent MDPs.

Consider an N-agent MDP $\mathcal{M}_{1:N}(\mathcal{S},\mathcal{A},\boldsymbol{P},\boldsymbol{r}_{1:N},\gamma)$ with joint state space $\mathcal{S}=\times_{i=1}^N\mathcal{S}_i$ and joint action space $\mathcal{A}=\times_{i=1}^N\mathcal{A}_i$. For simplicity, we assume $|\mathcal{S}_i|=|\mathcal{S}|$ and $|\mathcal{A}_i|=|\mathcal{A}|$ for each agent i. For agents at global state $\boldsymbol{s}=(s_1,s_2,\ldots,s_N)$ with action $\boldsymbol{a}=(a_1,a_2,\ldots,a_N)$ taken, the system will transit to $\boldsymbol{s}'=(s_1',s_2',\ldots,s_N')$ according to transition kernel $\boldsymbol{s}'\sim\boldsymbol{P}(\cdot\mid\boldsymbol{s},\boldsymbol{a})$ and each agent $i\in[N]$ will receive its local reward $r_i(s_i,a_i)$. The global reward $r_{1:N}$ is defined as the summation of local rewards $r_{1:N}(\boldsymbol{s},\boldsymbol{a})\coloneqq\sum_{i=1}^Nr_i(s_i,a_i)$, or in vector form,

$$oldsymbol{r}_{1:N} \in \mathbb{R}^{|\mathcal{S}|^N|\mathcal{A}|^N} := \sum_{i=1}^N (oldsymbol{e}\otimes)^{i-1} oldsymbol{r}_i (\otimes oldsymbol{e})^{N-i} \,.$$

We further assume the local rewards are bounded, i.e. for agent $i \in [N]$, $|r_i(s_i,a_i)| \leq r_{\max}^i$ for all (s_i,a_i) . Given any global policy $\pi \colon \mathcal{S} \to \Delta(\mathcal{A})$, we denote $P_{1:N}^{\pi} \in \mathbb{R}^{|\mathcal{S}|^N|\mathcal{A}|^N \times |\mathcal{S}|^N|\mathcal{A}|^N}$ as the transition matrix induced by π where $P_{1:N}^{\pi}(s'_{1:N},a'_{1:N}|s_{1:N},a_{1:N}) \coloneqq P(s'_{1:n}|s_{1:N},a_{1:N})\pi(a'_{1:N}|s'_{1:N})$. Then the global Q-value is defined by Bellman Equation $Q_{1:N}^{\pi} = (I - \gamma P_{1:N}^{\pi})^{-1} r_{1:N}$. The local Q-values follow the similar framework to Meta Algorithm 1 where each agent $i \in [N]$ fits Q_i^{π} using its local observations. We then sum up local Q-values to approximate the global Q-value, i.e.

$$Q_{1:N}^{\pi}(s, a) \approx \sum_{i=1}^{N} Q_i^{\pi}(s_i, a_i).$$

To illustrate the extension, we first provide the definition of multi-party quantum entanglement here for reference.

Definition 7 (Multi-party Quantum Entanglement). Consider a multi-party quantum system composed of N subsystems, indexed by [N]. The joint state $\rho_{1:N}$ is **separable** if there exists $K \in \mathbb{Z}^+$, probability distribution $\{x_i\}_{i \in [K]}$, and density matrices $\{\rho_{1:N}^{(j)}\}_{i \in [K]}$ such that

$$\rho_{1:N} = \sum_{j=1}^K x_j \rho_1^{(j)} \otimes \rho_2^{(j)} \otimes \cdots \otimes \rho_N^{(j)}.$$

If there exists no such decomposition, $\rho_{1:N}$ is called **entangled**.

Analogically, we define the Multi-agent Markov Entanglement,

Definition 8 (Multi-agent Markov Entanglement). Consider a N-agent Markov system $\mathcal{M}_{1:N}$ and policy $\pi \colon \mathcal{S} \to \Delta(\mathcal{A})$, the agents are **separable** under policy π if there exists $K \in \mathbb{Z}^+$, measure $\{x_j\}_{j \in [K]}$ satisfying $\sum_{j=1}^K x_j = 1$, and transition matrices $\{P_{1:N}^{(j)}\}_{j \in [K]}$ such that

$$\boldsymbol{P}_{1:N}^{\pi} = \sum_{j=1}^{K} x_{j} \boldsymbol{P}_{1}^{(j)} \otimes \boldsymbol{P}_{2}^{(j)} \otimes \cdots \otimes \boldsymbol{P}_{N}^{(j)}.$$

If there exists no such decomposition, the agents are **entangled** under policy π .

For clarity, we use superscript s^i to denote the *i*-th element in state space and subscript s_i to represent the state at *i*-th arm. Furthermore, we denote $S^{-i} := S \setminus s^i$ and $s := s_{1:N} := \{s_1, s_2, \dots, s_N\}$ is the profile of N-arms.

Given any global policy π , for any agent $i \in [N]$,

$$P_i^{\pi}(s_i', a_i' \mid s_i, a_i) = \sum_{s_{-i}', a_{-i}'} \sum_{s_{-i}, a_{-i}} P_{1:N}^{\pi}\left(s_{1:N}', a_{1:N}' \mid s_{1:N}, a_{1:N}\right) \rho_{1:N}^{\pi}(s_{-i}, a_{-i} \mid s_i, a_i) \,.$$

Definition 9 (Measure of Multi-agent Markov Entanglement). Consider a N-agent Markov system $\mathcal{M}_{1:N}$ with joint state space $\mathcal{S} = \times_{i=1}^N \mathcal{S}_i$ and action space $\mathcal{A} = \times_{i=1}^N \mathcal{A}_i$. Given any policy $\pi \colon \mathcal{S} \to \Delta(\mathcal{A})$, the measure of Markov entanglement of N agents is

$$E(\mathbf{P}_{1:N}^{\pi}) = \min_{\mathbf{P} \in \mathcal{P}_{SFP}} d(\mathbf{P}_{1:N}^{\pi}, \mathbf{P}), \qquad (13)$$

where $d(\cdot, \cdot)$ is some distance measure.

The following theorem generalizes the results of value-decomposition for two-agent Markov systems in Theorem 3 to multi-agent Markov systems.

Theorem 7. Consider a N-agent MDP $\mathcal{M}_{1:N}$ with joint state space $\mathcal{S} = \times_{i=1}^{N} \mathcal{S}_i$ and action space $\mathcal{A} = \times_{i=1}^{N} \mathcal{A}_i$. Given any policy $\pi : \mathcal{S} \to \Delta(\mathcal{A})$ with the measure of Markov entanglement $E_i(\mathbf{P}_{1:N}^{\pi})$ w.r.t ATV distance, it holds for any agent i,

$$\|\boldsymbol{P}_i^{\pi} - \boldsymbol{P}_i\|_{\infty} \leq 2_i E(\boldsymbol{P}_{1:N}^{\pi}).$$

where P_i is the optimal solution of Eq. (13). Furthermore, the decomposition error is entry-wise bounded by the measure of Markov entanglement,

$$\left\| Q_{1:N}^{\pi}(\boldsymbol{s}, \boldsymbol{a}) - \sum_{i=1}^{N} Q_{i}^{\pi}(s_{i}, a_{i}) \right\|_{\infty} \leq \frac{4\gamma \left(\sum_{i=1}^{N} E_{i}(\boldsymbol{P}_{1:N}^{\pi}) r_{\max}^{i} \right)}{(1 - \gamma)^{2}}.$$

The proof mainly follows the following lemma, which generalizes the key technique used in Theorem 1.

Lemma 5. For any agent i, it holds

$$\left(\sum_{j=1}^{K} x_j \mathbf{P}_1^{(j)} \otimes \mathbf{P}_2^{(j)} \otimes \cdots \otimes \mathbf{P}_N^{(j)}\right) \cdot \left((\mathbf{e} \otimes)^{i-1} \mathbf{r}_i (\otimes \mathbf{e})^{N-i}\right) = (\mathbf{e} \otimes)^{i-1} \left(\sum_{j=1}^{K} x_j \mathbf{P}_i^{(j)} \mathbf{r}_i\right) (\otimes \mathbf{e})^{N-i}.$$
(14)

The lemma follows from the property of tensor product. We can also extend Theorem 4 to multi-agent MDPs.

Theorem 8. Consider a N-agent MDP $\mathcal{M}_{1:N}$ with joint state space $\mathcal{S} = \times_{i=1}^{N} \mathcal{S}_i$ and action space $\mathcal{A} = \times_{i=1}^{N} \mathcal{A}_i$. Given any policy $\pi \colon \mathcal{S} \to \Delta(\mathcal{A})$ with the measure of Markov entanglement $E_i(\mathbf{P}_{1:N}^{\pi})$ w.r.t the $\mu_{1:N}^{\pi}$ -weighted agent-wise total variation distance, it holds for any agent i,

$$\|\boldsymbol{P}_i^{\pi} - \boldsymbol{P}_i\|_{\mu_i^{\pi},\infty} \leq 2E_i(\boldsymbol{P}_{1:N}^{\pi}).$$

where P_i is the optimal solution of Eq. (13) and μ_i^{π} is the stationary distribution of the projected transition P_i^{π} . Furthermore, the $\mu_{1:N}^{\pi}$ -weighted decomposition error is bounded by the measure of Markov entanglement,

$$\left\| Q_{1:N}^{\pi}(s, \boldsymbol{a}) - \sum_{i=1}^{N} Q_{i}^{\pi}(s_{i}, a_{i}) \right\|_{\mu_{1:N}^{\pi}} \leq \frac{4\gamma \left(\sum_{i=1}^{N} E_{i}(\boldsymbol{P}_{1:N}^{\pi}) r_{\max}^{i} \right)}{(1 - \gamma)^{2}}.$$

I Proof of Theorem 6

We first provide an overview of the proof and introduce the technical assumptions.

To begin, we consider the system configuration $\boldsymbol{m} \in \Delta^{|\mathcal{S}|}$ where $\boldsymbol{m}_s = \frac{1}{N} \sharp \{\text{Agents in state s}\}$ is the proportion of agents in state s. When $N \to \infty$, the transition between configurations will become deterministic under index policy and \boldsymbol{m} will approach its mean-field limit \boldsymbol{m}^* . Furthermore, in the mean-field, each agent's local transition will only depend its local state. As a result, the system will de-couple and become separable as $N \to \infty$.

To formalize this intuition, we introduce the following lemma that connects Markov entanglement measure with the mean-field analysis

Lemma 6. The measure of Markov entanglement w.r.t $\mu_{1:N}^{\pi}$ -weighted ATV distance is bounded by the deviation of mean-field configuration,

$$E_i(\pi) \leq |\mathcal{S}|^2 \cdot \mathbb{E}\left[\|\boldsymbol{m} - \boldsymbol{m}^*\|_{\infty}\right],$$

where the expectation is taking over the stationary distribution $m{m} \sim \mu_{1:N}^{\pi}$.

We thus focus on the deviation from m to m^* . We extend the concentration analysis from [20, 21] to derive a new stability bound for the RHS. Specifically, we finishing the proof via demonstrating the deviation decays at the rate $\mathcal{O}(1/\sqrt{N})$.

One caveat here is that we have to restrict chaotic behaviors in the mean-field limit. We thus introduce two technical assumptions.

We first define the transition of configuration under index policy π as $\phi^{\pi}: \Delta^{|\mathcal{S}|} \to \Delta^{|\mathcal{S}|}$ such that

$$\phi^{\pi}(\boldsymbol{m}) = \mathbb{E}\left[\boldsymbol{m}[t+1] \mid \boldsymbol{m}[t] = \boldsymbol{m}, \pi\right].$$

For t > 0, we denote $\Phi_t := (\phi^{\pi})^t$ apply the transition mapping for t rounds.

Assumption A (Uniform Global Attractor Property (UGAP)). There exists a uniform global attractor m^* of $\phi^{\pi}(\cdot)$, i.e. for all $\varepsilon > 0$, there exists $T(\varepsilon)$ such that for all $t \geq T(\varepsilon)$ and all $m \in \Delta^{|\mathcal{S}|}$, one has $\|\Phi_t(m) - m^*\|_{\infty} < \varepsilon$.

The UGAP assumption ensures the uniqueness of m^* and guarantees fast convergence from any initial m to m^* .

Assumption B (Non-degenerate RMAB). There exists state $s \in \mathcal{S}$ such that $0 < \pi^*(s, 0) < 1$, where π^* is the policy under m^* .

The non-degenerate assumption further restricts cyclic behavior in the mean-field limit.

Non-degenerate and UGAP are two standard technical assumptions for the index policy, which restrict chaotic behavior in asymptotic regime and will be further introduced in subsequent sections. We note here these two assumptions are also used in almost all theoretical work on index policies [47, 41, 20, 21].

Proof of Theorem 6. In the subsequent proof, we let $\nu_1 > \nu_2 > \nu_3 > \cdots > \nu_{|S|}$. This does not lose generality in that we can always exchange state index. The proof consists of several steps

Step 1: Find m^* Recall the transition mapping for configurations $\phi^{\pi} : \Delta^{|\mathcal{S}|} \to \Delta^{|\mathcal{S}|}$,

$$\phi^{\pi}(\boldsymbol{m}) = \mathbb{E}\left[\boldsymbol{m}[t+1] \mid \boldsymbol{m}[t] = \boldsymbol{m}, \pi\right].$$

Notice that the definition of ϕ^{π} does not depend on N. We adapt from Lemma B.1 in [20] defined specially for Whittle Index,

Lemma 7 (Piecewise Affine). Given any index policy π , ϕ^{π} is a piecewise affine continuous function with |S| affine pieces.

When the context is clear, we abbreviate ϕ^{π} as ϕ . For any $m \in \Delta^{|\mathcal{S}|}$, define $s(m) \in [|\mathcal{S}|]$ be the state such that $\sum_{i=1}^{s(m)-1} m_i \leq \alpha < \sum_{i=1}^{s(m)} m_i$. Lemma 7 characterizes for any $m \in \mathcal{Z}_i := \{m \in \Delta^{|\mathcal{S}|} \mid s(m) = i\}$, there exists $K_{s(m)}, b_{s(m)}$ such that

$$\phi(\boldsymbol{m}) = \boldsymbol{K}_{s(\boldsymbol{m})} \boldsymbol{m} + \boldsymbol{b}_{s(\boldsymbol{m})}.$$

By Brouwer fixed point theorem, there exists a fixed point m^* such that $\phi(m^*) = m^*$. The UGAP condition guarantees the uniqueness of m^* . Our choice of π^* is the corresponding policy under m^* .

Step 2: Connecting policy entanglement with the deviation of stationary distribution Combine Proposition 9 with the RMAB model, we have

Lemma 8. The measure of Markov entanglement w.r.t $\mu_{1:N}^{\pi}$ -weighted ATV distance is bounded by the deviation of mean-field configuration,

$$E_i(\pi) \leq |\mathcal{S}|^2 \cdot \mathbb{E}\left[\|\boldsymbol{m} - \boldsymbol{m}^*\|_{\infty}\right],$$

where the expectation is taking over the stationary distribution $m{m} \sim \mu_{1:N}^{\pi}$.

Proof. Given the homogeneity of agents, we first demonstrate for any two agent i, j, it holds

$$\sum_{s_{1:N}} \mu^{\pi}(s_{1:N}) \left| \pi(a_i = a \mid s_{1:N}) - \pi^*(a_i = a \mid s_i) \right| = \sum_{s_{1:N}} \mu^{\pi}(s_{1:N}) \left| \pi(a_j = a \mid s_{1:N}) - \pi^*(a_j = a \mid s_i) \right| .$$

To see this, we first notice by the definition of index policy

$$|\pi(a_i=a\mid s_i=s, \boldsymbol{m})-\pi^*(a\mid s)|=|\pi(a_j=a\mid s_j=s, \boldsymbol{m})-\pi^*(a\mid s)|\;.$$
 It then suffices to prove
$$\sum_{s_i=s,s_{1:N}=\boldsymbol{m}}\mu(s_{1:N})=\sum_{s_j=s,s_{1:N}=\boldsymbol{m}}\mu(s_{1:N}).$$
 If
$$\sum_{s_i=s,s_{1:N}=\boldsymbol{m}}\mu(s_{1:N})\leq \sum_{s_j=s,s_{1:N}=\boldsymbol{m}}\mu(s_{1:N}), \text{ we can exchange the agent index of } i \text{ and } j. \text{ This will result in the same stationary distribution and } \sum_{s_i=s,s_{1:N}=\boldsymbol{m}}\mu(s_{1:N})\geq \sum_{s_j=s,s_{1:N}=\boldsymbol{m}}\mu(s_{1:N})$$
 and thus the equation. We then rewrite the bound in Proposition 9,

$$E(\pi) \leq \frac{1}{2} \sup_{i} \sum_{s_{1:N}} \mu^{\pi}(s_{1:N}) \sum_{a_{i}} |\pi(a_{i} | s_{1:N}) - \pi^{*}(a_{i} | s_{i})|$$

$$= \sup_{i} \sum_{s_{1:N}} \mu^{\pi}(s_{1:N}) |\pi(a_{i} = 1 | s_{1:N}) - \pi^{*}(a_{i} = 1 | s_{i})|$$

$$= \frac{1}{N} \sum_{s_{1:N}} \mu^{\pi}(s_{1:N}) \sum_{i=1}^{N} |\pi(a_{i} = 1 | s_{1:N}) - \pi^{*}(a_{i} = 1 | s_{i})|$$

$$= \sum_{m} \mu^{\pi}(m) \sum_{s \in \mathcal{S}} m_{s} |\pi(a = 1 | s, m) - \pi^{*}(a = 1 | s)|$$

For any configuration m and state s, we have

$$\begin{aligned} & m_{s} \left| \pi(a = 1 \mid s, m) - \pi^{*}(a = 1 \mid s) \right| \\ &= m_{s} \left| \frac{\pi^{*}(a = 1 \mid s) m_{s}^{*} N + k_{s}}{m_{s}^{*} N + \ell_{s}} - \pi^{*}(a = 1 \mid s) \right| \\ &= \frac{m_{s}^{*} N + \ell_{s}}{N} \left| \frac{k_{s} - \ell_{s} \pi^{*}(a = 1 \mid s)}{m_{s}^{*} N + \ell_{s}} \right| \\ &\leq |\mathcal{S}| ||m - m^{*}||_{\infty}. \end{aligned}$$

where $|k_s| \leq (|\mathcal{S}|-1) \|\boldsymbol{m}-\boldsymbol{m}^*\|_{\infty} N$ representing the additional fraction of state s to be activated due to the deviation from m^* and $|\ell_s| \leq \|\boldsymbol{m}-\boldsymbol{m}^*\|_{\infty} N$ representing the deviation of \boldsymbol{m}_s from \boldsymbol{m}_s^* . The results then hold by taking summation over s and expectation over \boldsymbol{m} .

Step 3: Concentrations and local stability To bound $\mathbb{E}[\|m-m^*\|_{\infty}]$, we start with several technical lemmas from previous RMAB literature. We use the same notation $\Phi_t = \phi(\Phi_{t-1})$.

Lemma 9 (One-step Concentration, Lemma 1 in [21]). Let $\epsilon[1] = m[1] - \phi(m[0])$, it holds

$$\mathbb{E}\left[\|\epsilon[1]\|_1 \mid m{m}[0]
ight] \leq \sqrt{rac{|\mathcal{S}|}{N}}$$
 .

Lemma 10 (Multi-step Concentration, Lemma C.4 in [20]). There exists a positive constant K such that for all $t \in \mathbb{N}$ and $\delta > 0$,

$$\Pr\left[\|\boldsymbol{m}[t] - \Phi_t(\boldsymbol{m})\|_{\infty} \ge (1 + K + K^2 + \dots + K^t)\delta \mid \boldsymbol{m}[0] = \boldsymbol{m}\right] \le t|\mathcal{S}|e^{-2N\delta^2}$$

Lemma 11 (Local Stability, Lemma C.5 in [20]). *Under non-degenerate and UGAP:*

- (i) $K_{s(m^*)}$ is a stable matrix, i.e. its spectral radius is strictly less than 1.
- (ii) For any ϵ , there exists $T(\epsilon) > 0$ such that for all $m \in \Delta^{|\mathcal{S}|}$, $\|\Phi_{T(\epsilon)}(m) m^*\|_{\infty} < \epsilon$.

The first result implies there exists some matrix norm $\|\cdot\|_{\beta}$ such that $\|\boldsymbol{K}_{s(\boldsymbol{m}^*)}\|_{\beta} < 1$. By the equivalence of norms, there exists constant $C^1_{\beta}, C^2_{\beta} > 0$ such that for all $\boldsymbol{x} \in \mathbb{R}^{|\mathcal{S}|}$

$$C_{\beta}^1 \| \boldsymbol{x} \|_{\beta} \leq \| \boldsymbol{x} \|_{\infty} \leq C_{\beta}^2 \| \boldsymbol{x} \|_{\beta}$$
.

Combine the second result of Lemma 11 and non-degenerate condition, we can construct a neighborhood $\mathcal N$ of $\boldsymbol m^*$ such that $\mathcal N=\mathcal B(\boldsymbol m^*,\epsilon)\cap\Delta^{|\mathcal S|}\in\mathcal Z_{s(\boldsymbol m^*)}$ where $\epsilon>0$ and $\mathcal B(\boldsymbol m^*,\epsilon)=\{\boldsymbol m\mid \|\boldsymbol m-\boldsymbol m^*\|_\infty<\epsilon\}$ is an open ball. We next show that $\boldsymbol m[0]$ under stationary distribution will concentrate in $\mathcal N$ with high probability. Let $\tilde T=T(\epsilon/2)$ such that for all $\boldsymbol m\in\Delta^{|\mathcal S|}$, $\|\Phi_{\tilde T}(\boldsymbol m)-\boldsymbol m^*\|_\infty<\epsilon/2$. It holds

$$\begin{split} \Pr\left[\boldsymbol{m}[0] \neq \mathcal{N}\right] &= \Pr\left[\|\boldsymbol{m}[0] - \boldsymbol{m}^*\|_{\infty} \geq \epsilon\right] \\ &\stackrel{(i)}{=} \Pr\left[\left\|\boldsymbol{m}[\tilde{T}] - \boldsymbol{m}^*\right\|_{\infty} \geq \epsilon \mid \boldsymbol{m}[0] = \boldsymbol{m}\right] \\ &\leq \Pr\left[\left\|\boldsymbol{m}[\tilde{T}] - \Phi_{\tilde{T}}(\boldsymbol{m})\right\|_{\infty} \geq \frac{\epsilon}{2} \mid \boldsymbol{m}[0] = \boldsymbol{m}\right] + \Pr\left[\left\|\Phi_{\tilde{T}}(\boldsymbol{m}) - \boldsymbol{m}^*\right\|_{\infty} \geq \frac{\epsilon}{2}\right] \\ &= \Pr\left[\left\|\boldsymbol{m}[\tilde{T}] - \Phi_{\tilde{T}}(\boldsymbol{m})\right\|_{\infty} \geq \frac{\epsilon}{2} \mid \boldsymbol{m}[0] = \boldsymbol{m}\right] \leq \tilde{T}|\mathcal{S}|e^{-2uN} \end{split}$$

where (i) follows from the stationarity $m[\tilde{T}]$ and m[0] are i.i.d and the constant $u = \left(\frac{\epsilon}{2(1+K+K^2+\cdots+K^{\tilde{T}})}\right)^2$ does not depend on N.

Step 4: Put it together Finally, we are ready to bound $\mathbb{E}[\|\boldsymbol{m} - \boldsymbol{m}^*\|_{\infty}]$. Notice for all $\boldsymbol{m}[0] \in \mathcal{N}$, we have

$$m[1] - m^* = \phi(m[0]) + \epsilon[1] - m^*$$
$$= K_{s(m^*)}(m[0] - m^*) + \epsilon[1].$$

Taking $\|\cdot\|_{\beta}$ on both side,

$$\|\boldsymbol{m}[1] - \boldsymbol{m}^*\|_{\beta} \le \|\boldsymbol{K}_{s(\boldsymbol{m}^*)} (\boldsymbol{m}[0] - \boldsymbol{m}^*)\|_{\beta} + \|\epsilon[1]\|_{\beta}$$

 $\le \|\boldsymbol{K}_{s(\boldsymbol{m}^*)}\|_{\beta} \|\boldsymbol{m}[0] - \boldsymbol{m}^*\|_{\beta} + \|\epsilon[1]\|_{\beta}.$

Taking expectation on both side,

$$\mathbb{E}\left[\left\|\boldsymbol{m}[1]-\boldsymbol{m}^*\right\|_{\beta}\right]$$

$$=\mathbb{E}\left[\left\|\phi(\boldsymbol{m}[0])-\boldsymbol{m}^*\right\|_{\beta}\cdot\mathbf{1}\left\{\boldsymbol{m}[0]\in\mathcal{N}\right\}\right]+\mathbb{E}\left[\left\|\phi(\boldsymbol{m}[0])-\boldsymbol{m}^*\right\|_{\beta}\cdot\mathbf{1}\left\{\boldsymbol{m}[0]\notin\mathcal{N}\right\}\right]+\mathbb{E}\left[\left\|\epsilon[1]\right\|_{\beta}\right]$$

$$\leq\left\|\boldsymbol{K}_{s(\boldsymbol{m}^*)}\right\|_{\beta}\mathbb{E}\left[\left\|\boldsymbol{m}[0]-\boldsymbol{m}^*\right\|_{\beta}\cdot\mathbf{1}\left\{\boldsymbol{m}[0]\in\mathcal{N}\right\}\right]+\Pr\left[\boldsymbol{m}[0]\notin\mathcal{N}\right]\sup_{\boldsymbol{m}[0]}\left\|\phi(\boldsymbol{m}[0])-\boldsymbol{m}^*\right\|_{\beta}+\mathbb{E}\left[\left\|\epsilon[1]\right\|_{\beta}\right]$$

$$\leq\left\|\boldsymbol{K}_{s(\boldsymbol{m}^*)}\right\|_{\beta}\mathbb{E}\left[\left\|\boldsymbol{m}[0]-\boldsymbol{m}^*\right\|_{\beta}\right]+\Pr\left[\boldsymbol{m}[0]\notin\mathcal{N}\right]\sup_{\boldsymbol{m}[0]}\left\|\phi(\boldsymbol{m}[0])-\boldsymbol{m}^*\right\|_{\beta}+\mathbb{E}\left[\left\|\epsilon[1]\right\|_{\beta}\right].$$

By stationarity, one have $\mathbb{E}\left[\|\boldsymbol{m}[1]-\boldsymbol{m}^*\|_{\beta}\right]=\mathbb{E}\left[\|\boldsymbol{m}[0]-\boldsymbol{m}^*\|_{\beta}\right]$. This refines the above inequality,

$$\mathbb{E}\left[\left\|\boldsymbol{m}[0] - \boldsymbol{m}^*\right\|_{\infty}\right] \leq \frac{C_{\beta}^2}{1 - \left\|\boldsymbol{K}_{s(\boldsymbol{m}^*)}\right\|_{\beta}} \left(\sup_{\boldsymbol{m}[0]} \Pr\left[\boldsymbol{m}[0] \notin \mathcal{N}\right] \left\|\phi(\boldsymbol{m}[0]) - \boldsymbol{m}^*\right\|_{\beta} + \mathbb{E}\left[\left\|\epsilon[1]\right\|_{\beta}\right]\right) \\
\leq \frac{C_{\beta}^2}{C_{\beta}^1 (1 - \left\|\boldsymbol{K}_{s(\boldsymbol{m}^*)}\right\|_{\beta})} \left(\Pr\left[\boldsymbol{m}[0] \notin \mathcal{N}\right] + \mathbb{E}\left[\left\|\epsilon[1]\right\|_{\infty}\right]\right) \\
\leq \frac{C_{\beta}^2}{C_{\beta}^1 (1 - \left\|\boldsymbol{K}_{s(\boldsymbol{m}^*)}\right\|_{\beta})} \left(\tilde{T}|\mathcal{S}|e^{-2uN} + \frac{\sqrt{|\mathcal{S}|}}{\sqrt{N}}\right).$$

We combine Lemma 8 and conclude the proof of Theorem 6.

J Extensions of Markov entanglement

J.1 (Weakly-)coupled MDPs

Weakly-coupled MDPs (WCMDP) are a rich class of multi-agent model that capture many real-world applications such as supply chain management, queuing network and resource allocations [2, 12, 36]. Compared to general multi-agent MDP, WCMDP further ensures each agent follow its local transition while the agents' actions are coupled with each other. Formally,

Definition 10 (Weakly-coupled MDPs). An N-agent MDP $\mathcal{M}_{1:N}(\mathcal{S}, \mathcal{A}, \boldsymbol{P}, \boldsymbol{r}_{1:N}, \gamma)$ is a weakly-coupled MDP if

- Each agent has local transition kernel P_i such that $\forall s, a, s', P(s' \mid s, a) = \prod_{i=1}^N P_i(s_i' \mid s_i, a_i)$.
- At global state s, agents' joint actions a are subject to m coupling constraints $\sum_{i=1}^{N} d(s_i, a_i) \leq b \in \mathbb{R}^m$.

We then demonstrate that this weakly-coupled structure can further refine the analysis of Markov entanglement measure.

Proposition 9. Consider a N-agent weakly-coupled MDP $\mathcal{M}_{1:N}(\mathcal{S}, \mathcal{A}, \boldsymbol{P}, \boldsymbol{r}_{1:N}, \gamma)$. Given any policy $\pi \colon \mathcal{S} \to \Delta(\mathcal{A})$ with measure of Markov entanglement $E_i(\boldsymbol{P}_{1:N}^{\pi})$ w.r.t the $\mu_{1:N}^{\pi}$ -weighted agent-wise total variation distance, it holds for $i \in [N]$,

$$E_i(\mathbf{P}_{1:N}^{\pi}) \leq \min_{\pi'} \frac{1}{2} \sum_{\mathbf{s}} \mu_{1:N}^{\pi}(\mathbf{s}) \sum_{a_i} |\pi(a_i \mid \mathbf{s}) - \pi'(a_i \mid s_i)| ,$$

where $\pi': \mathcal{S}_i \to \mathcal{A}_i$ is any local policy for agent i.

Proof of Proposition 9. We demonstrate the proof for two-agent WCMDP and the generalization to multi-agent WCMDP is straightforward. Consider $P_A^{\pi'}$ be the transition of agent A under local policy

 π' . We focus on agent A

$$\begin{split} &E_{A}(\boldsymbol{P}_{AB}^{\pi}) \\ &\leq \frac{1}{2} \sum_{\boldsymbol{s}, \boldsymbol{a}} \mu_{AB}^{\pi}(\boldsymbol{s}, \boldsymbol{a}) \sum_{s_{A}', a_{A}'} \left| P_{AB}^{\pi}(s_{A}', a_{A}' \mid \boldsymbol{s}, \boldsymbol{a}) - P_{A}^{\pi'}(s_{A}', a_{A}' \mid s_{A}, a_{A}) \right| \\ &= \frac{1}{2} \sum_{\boldsymbol{s}, \boldsymbol{a}} \mu_{AB}^{\pi}(\boldsymbol{s}, \boldsymbol{a}) \sum_{s_{A}', a_{A}'} \left| \sum_{s_{B}'} P_{AB}^{\pi}(\boldsymbol{s}', a_{A} \mid \boldsymbol{s}, \boldsymbol{a}) - P_{A}^{\pi'}(s_{A}' \mid s_{A}, a_{A}) \pi'(a_{A}' \mid s_{A}') \right| \\ &\stackrel{(i)}{=} \frac{1}{2} \sum_{\boldsymbol{s}, \boldsymbol{a}} \mu_{AB}^{\pi}(\boldsymbol{s}, \boldsymbol{a}) \sum_{s_{A}', a_{A}'} \left| \sum_{s_{B}'} P_{AB}^{\pi}(\boldsymbol{s}', a_{A} \mid \boldsymbol{s}, \boldsymbol{a}) - \sum_{s_{B}'} P(\boldsymbol{s}' \mid \boldsymbol{s}, \boldsymbol{a}) \pi'(a_{A}' \mid s_{A}') \right| \\ &= \frac{1}{2} \sum_{\boldsymbol{s}, \boldsymbol{a}} \mu_{AB}^{\pi}(\boldsymbol{s}, \boldsymbol{a}) \sum_{s_{A}', a_{A}'} \left| \sum_{s_{B}'} P(\boldsymbol{s}' \mid \boldsymbol{s}, \boldsymbol{a}) \left(\pi(a_{A}' \mid \boldsymbol{s}') - \pi'(a_{A}' \mid s_{A}') \right) \right| \\ &\leq \frac{1}{2} \sum_{\boldsymbol{s}, \boldsymbol{a}} \mu_{AB}^{\pi}(\boldsymbol{s}, \boldsymbol{a}) \sum_{s_{A}'} P(\boldsymbol{s}' \mid \boldsymbol{s}, \boldsymbol{a}) \sum_{a_{A}'} \left| \pi(a_{A}' \mid \boldsymbol{s}') - \pi'(a_{A}' \mid s_{A}') \right| \\ &\stackrel{(ii)}{=} \frac{1}{2} \sum_{\boldsymbol{s}'} \mu_{AB}^{\pi}(\boldsymbol{s}') \sum_{a_{A}'} \left| \pi(a_{A}' \mid \boldsymbol{s}') - \pi'(a_{A}' \mid s_{A}') \right| \\ &\stackrel{(ii)}{=} \frac{1}{2} \sum_{\boldsymbol{s}'} \mu_{AB}^{\pi}(\boldsymbol{s}') \sum_{a_{A}'} \left| \pi(a_{A}' \mid \boldsymbol{s}') - \pi'(a_{A}' \mid s_{A}') \right| . \end{split}$$

where (i) follows from the transition structure of weakly coupled MDP $P(s'\mid s, a) = P(s'_A\mid s_A, a_A)\cdot P(s'_B\mid s_B, a_B);$ and (ii) comes from the fact that $P^\pi(s'\mid s) = \sum_{\boldsymbol{a}} \pi(\boldsymbol{a}\mid s)P(s'\mid s, a)$ and $\sum_{\boldsymbol{s}} \mu^\pi(s)P^\pi(s'\mid s) = \mu^\pi(s').$

Proposition 9 establishes an upper bound for Markov entanglement in WCMDP. Intuitively, this bound characterizes *how agent i can be viewed as making independent decisions*. It takes advantage of the weakly-coupled structure and shaves off the transition in Markov entanglement measure.

J.2 Coupled MDPs with exogenous information

In many practical scenarios, the agents' transitions and actions are coupled by a shared exogenous signal. For example, in ride-hailing platforms, the specific dispatch is related to the exogenous order at the current moment [33, 24, 4]; in warehouse routing, the scheduling of robots is also related to the exogenous task revealed so far [15].

We will then enrich our framework by incorporating these exogenous information. At each timestep t, there will an exogenous information z_t revealed to the decision maker. z_t is assumed to evolve following a Markov chain independent of the action and transition of agents. We assume $z_t \in \mathcal{Z}$ and \mathcal{Z} is finite.

Given the current state s and exogenous information z, the policy is given by $\pi: \mathcal{S} \times \mathcal{Z} \to \Delta(\tilde{\mathcal{A}})$, where $\tilde{\mathcal{A}}$ refers to the set of feasible actions. We then have the global transition depending on exogenous information z,

$$P^{\pi}_{ABz}(\boldsymbol{s}',\boldsymbol{a}',z'\mid \boldsymbol{s},\boldsymbol{a},z) = P(\boldsymbol{s}'\mid \boldsymbol{s},\boldsymbol{a},z) \cdot \pi(\boldsymbol{a}'\mid \boldsymbol{s}',z') \cdot P(z'\mid z) \,.$$

and global Q-value $Q_{ABz}^{\pi} \in \mathbb{R}^{|\mathcal{S}|^N |\mathcal{A}|^N |\mathcal{Z}|}$,

$$Q_{AB}^{\pi}(s, \boldsymbol{a}, z) = \mathbb{E}\left[\sum_{t=0}^{\infty} \sum_{i=1}^{N} r(s_{i,t}, a_{i,t}, z_t) \mid s_0 = s, \boldsymbol{a}_0 = \boldsymbol{a}, z_0 = z\right].$$

We assume the system is unichain and the stationary distribution is μ_{ABz}^{π} . Then we can derive the local transition under new algorithm by

$$P_{Az}(s_A^\prime, a_A^\prime, z^\prime \mid s_A, a_A, z) = \sum_{s_B, a_B} \mu_{ABz}^\pi(s_B, a_B \mid s_A, a_A, z) \sum_{s_B^\prime, a_B^\prime} P_{ABz}^\pi(s^\prime, \boldsymbol{a}^\prime, z^\prime \mid \boldsymbol{s}, \boldsymbol{a}, z) ,$$

Given the local transition, we have the local value $m{Q}_{Az}^{\pi}=(m{I}-\gammam{P}_{Az})^{-1}(m{r}_{Az})$ via Bellman Equation.

Combined with exogenous information, we consider the following value decomposition

$$Q_{AB}^{\pi}(s, \boldsymbol{a}, z) = Q_{A}^{\pi}(s_A, a_A, z) + Q_{B}^{\pi}(s_B, a_B, z)$$
.

We start by introducing agent-wise Markov entanglement defined for each agent

$$P_{ABz}^{\pi} = \sum_{j=1}^{K} x_j P_{Az}^{(j)} \otimes P_B^{(j)}.$$
 (15)

Proposition 10. If the system is agent-wise separable for all agents, then

$$oldsymbol{Q}_{ABz}^{\pi} = oldsymbol{Q}_{Az}^{\pi} \otimes oldsymbol{e}_{|\mathcal{S}||\mathcal{A}|} + oldsymbol{e}_{|\mathcal{S}||\mathcal{A}|} \otimes oldsymbol{Q}_{Bz}^{\pi}$$
 .

Proof. The proof is basically the same as Theorem 1. One can first quickly show that $P_{Az} = \sum_{j=1}^{K} x_j P_{Az}^{(j)}$. And then it holds

$$\left(\sum_{j=1}^{K} x_{j} \mathbf{P}_{Az}^{(j)} \otimes \mathbf{P}_{B}^{(j)}\right)^{t} \left(\mathbf{r}_{A} \otimes \mathbf{e}_{|z|} \otimes \mathbf{e}_{|\mathcal{S}||\mathcal{A}|}\right)$$

$$= \left(\sum_{j=1}^{K} x_{j} \mathbf{P}_{Az}^{(j)} \otimes \mathbf{P}_{B}^{(j)}\right)^{t-1} \left(\sum_{j=1}^{K} x_{j} \left(\mathbf{P}_{Az}^{(j)}(\mathbf{r}_{A} \otimes \mathbf{e}_{|z|})\right) \otimes \left(\mathbf{P}_{B}^{(j)} \mathbf{e}\right)\right)$$

$$= \left(\sum_{j=1}^{K} x_{j} \mathbf{P}_{Az}^{(j)} \otimes \mathbf{P}_{B}^{(j)}\right)^{t-1} \left(\sum_{j=1}^{K} x_{j} \mathbf{P}_{Az}^{(j)}(\mathbf{r}_{A} \otimes \mathbf{e}_{|z|})\right) \otimes \mathbf{e}$$

$$= \dots = \left(\left(\sum_{j=1}^{K} x_{j} \mathbf{P}_{Az}^{(j)}\right)^{t} \left(\mathbf{r}_{A} \otimes \mathbf{e}_{|z|}\right)\right) \otimes \mathbf{e}.$$

We then provide the measure of Markov entanglement with exogenous information w.r.t agent-wise total variation distance.

$$E_{A}(\mathbf{P}_{AB}^{\pi}, \mathcal{Z}) := \min \frac{1}{2} \left\| \mathbf{P}_{ABz}^{\pi} - \sum_{j=1}^{K} x_{j} \mathbf{P}_{Az}^{(j)} \otimes \mathbf{P}_{B}^{(j)} \right\|_{ATV_{1}}$$

$$= \min_{\mathbf{P}_{Az}} \max_{\mathbf{s}, \mathbf{a}, z} \frac{1}{2} \sum_{s'_{A}, a'_{A}, z'} \left| \mathbf{P}_{ABz}^{\pi}(s'_{A}, a'_{A}, z' \mid \mathbf{s}, \mathbf{a}, z) - P_{Az}(s'_{A}, a'_{A}, z' \mid s_{A}, a_{A}, z) \right|.$$
(16)

Similar to Theorem 3, we can connect this measure of Markov entanglement with the value decomposition error.

Theorem 11. Consider a N-agent Markov system $\mathcal{M}_{1:N}$. Given any policy $\pi \colon \mathcal{S} \to \Delta(\mathcal{A})$ with the measure of Markov entanglement $E_i(\mathbf{P}_{1:N}^{\pi}, \mathcal{Z})$ w.r.t the agent-wise total variation distance, it holds for any agent i,

$$\left\| \boldsymbol{P}_{iz}^{\pi} - \sum_{j=1}^{K} x_{j} \boldsymbol{P}_{iz}^{(j)} \right\|_{2} \leq 2E_{i}(\boldsymbol{P}_{1:N}^{\pi}, \mathcal{Z}).$$

Furthermore, the decomposition error is entry-wise bounded by the measure of Markov entanglement,

$$\left\| Q_{1:N}^{\pi}(\boldsymbol{s}, \boldsymbol{a}, z) - \sum_{i=1}^{N} Q_{iz}^{\pi}(s_i, a_i, z) \right\|_{20} \leq \frac{4\gamma \left(\sum_{i=1}^{N} E_i(\boldsymbol{P}_{1:N}^{\pi}, \boldsymbol{\mathcal{Z}}) r_{\max}^i \right)}{(1 - \gamma)^2}.$$

In practice, exogenous information is often discussed in the context of (weakly-)coupled MDPs, where each agent independent evolves by $P_i(s_{i+1} \mid s_i, a_i, z)$. Interestingly, we can derive a similar result to Proposition 9 that shaves off the transition in entanglement analysis.

Proposition 12. Consider a N-agent Weakly Coupled Markov system $\mathcal{M}_{1:N}$. Given any policy $\pi \colon \mathcal{S} \to \Delta(\mathcal{A})$ and its measure of Markov entanglement $E_i(\mathbf{P}_{1:N}^{\pi}, \mathcal{Z})$ w.r.t the $\mu_{1:N}^{\pi}$ -weighted agent-wise total variation distance, it holds

$$E_i(\mathbf{P}_{1:N}^{\pi}, \mathcal{Z}) \leq \frac{1}{2} \sum_{s_{1:N}, z} \mu^{\pi}(s_{1:N}, z) \sum_{a_i} |\pi(a_i \mid s_{1:N}, z) - \pi'(a_i \mid s_i, z)|,$$

for any policies π' .

Proof. We provide the proof for two-agent MDP, which can be easily generalized to N-agent case.

$$\begin{split} & E_{A}(\boldsymbol{P}^{\pi}_{AB},\mathcal{Z}) \\ & \leq \frac{1}{2} \sum_{\boldsymbol{s},\boldsymbol{a},z} \mu(\boldsymbol{s},\boldsymbol{a},z) \sum_{s'_{A},a'_{A},z'} |P^{\pi}_{ABz}(s'_{A},a'_{A},z'\mid \boldsymbol{s},\boldsymbol{a},z) - P_{Az}(s'_{A},a'_{A},z'\mid s_{A},a_{A},z)| \\ & = \frac{1}{2} \sum_{\boldsymbol{s},\boldsymbol{a},z} \mu(\boldsymbol{s},\boldsymbol{a},z) \sum_{s'_{A},a'_{A},z'} \left| \sum_{s'_{B}} P^{\pi}_{ABz}(\boldsymbol{s}',a_{A},z'\mid \boldsymbol{s},\boldsymbol{a},z) - P_{Az}(s'_{A},z'\mid s_{A},a_{A},z) \pi'(a'_{A}\mid s'_{A},z') \right| \\ & = \frac{1}{2} \sum_{\boldsymbol{s},\boldsymbol{a},z} \mu(\boldsymbol{s},\boldsymbol{a},z) \sum_{s'_{A},a'_{A},z'} \left| \sum_{s'_{B}} P^{\pi}_{ABz}(\boldsymbol{s}',a_{A},z'\mid \boldsymbol{s},\boldsymbol{a},z) - \sum_{s'_{B}} P(\boldsymbol{s}',z'\mid \boldsymbol{s},\boldsymbol{a},z) \pi'(a'_{A}\mid s'_{A},z') \right| \\ & = \frac{1}{2} \sum_{\boldsymbol{s},\boldsymbol{a},z} \mu(\boldsymbol{s},\boldsymbol{a},z) \sum_{s'_{A},a'_{A},z'} \left| \sum_{s'_{B}} P(\boldsymbol{s}',z'\mid \boldsymbol{s},\boldsymbol{a},z) \left(\pi(a'_{A}\mid \boldsymbol{s}',z') - \pi'(a'_{A}\mid s'_{A},z') \right) \right| \\ & \leq \frac{1}{2} \sum_{\boldsymbol{s},\boldsymbol{a},z} \mu(\boldsymbol{s},\boldsymbol{a},z) \sum_{\boldsymbol{s}',z'} P(\boldsymbol{s}',z'\mid \boldsymbol{s},\boldsymbol{a},z) \sum_{a'_{A}} |\pi(a'_{A}\mid \boldsymbol{s}',z') - \pi'(a'_{A}\mid s'_{A},z')| \\ & = \frac{1}{2} \sum_{\boldsymbol{s}',z'} \mu(\boldsymbol{s}',z') \sum_{a'_{A}} |\pi(a'_{A}\mid \boldsymbol{s}',z') - \pi'(a'_{A}\mid s'_{A},z')| \ . \end{split}$$

J.3 Factored MDPs

Another common class of multi-agent MDPs is Factored MDPs (FMDPs, [22, 23, 32]), which explicitly model the structured dependencies in state transitions. For instance, in a server cluster, the state transition of each server depends only on its neighboring servers. Formally, we define

Definition 11 (Factored MDPs). An N-agent MDP $\mathcal{M}_{1:N}(\mathcal{S}, \mathcal{A}, \boldsymbol{P}, \boldsymbol{r}_{1:N}, \gamma)$ is a factored MDP if each agent i has neighbor set $Z_i \in [N]$ such that its transition is affected by all its neighbors, i.e. $P(s_i' \mid \boldsymbol{s}, \boldsymbol{a}) = P(s_i' \mid s_{Z_i}, a_{Z_i})$.

The neighbor set $|Z_i|$ is often assumed to be much smaller compared to the number of agents N. This helps to encode exponentially large system very compactly. We show this idea can also be captured in Markov entanglement. Consider the measure of Markov entanglement w.r.t ATV distance in Eq. (7),

$$E_{A}(\boldsymbol{P}_{AB}^{\pi}) = \min_{\boldsymbol{P}_{A}} \max_{(\boldsymbol{s},\boldsymbol{a}) \in \mathcal{S} \times \mathcal{A}} D_{\text{TV}} \Big(\boldsymbol{P}_{AB}^{\pi}(\cdot, \cdot \mid \boldsymbol{s}, \boldsymbol{a}), \boldsymbol{P}_{A}(\cdot, \cdot \mid s_{A}, a_{A}) \Big)$$

$$= \min_{\boldsymbol{P}_{A}} \max_{(\boldsymbol{s},\boldsymbol{a}) \in \mathcal{S} \times \mathcal{A}} D_{\text{TV}} \Big(\boldsymbol{P}_{AB}^{\pi}(\cdot, \cdot \mid s_{Z_{A}}, a_{Z_{A}}), \boldsymbol{P}_{A}(\cdot, \cdot \mid s_{A}, a_{A}) \Big).$$

Thus we conclude the agent-wise Markov entanglement will only depend on its neighbor set.

Meta Algorithm 2: Q-value Decomposition with Shared Reward

Require: Global policy π ; horizon length T.

- 1: Execute π for T epochs and obtain $\mathcal{D} = \left\{ (s_{AB}^t, a_{AB}^t, r_{AB}^t, s_{AB}^{t+1}, a_{AB}^{t+1}) \right\}_{t=1}^{T-1}$. 2: Each agent $i \in \{A, B\}$ fits Q_i^π using local observations $\mathcal{D}_i = \left\{ (s_i^t, a_i^t, r_i, s_i^{t+1}, a_i^{t+1}) \right\}_{t=1}^{T-1}$ where the local reward (r_A, r_B) is learned via solving

$$\min_{\boldsymbol{r}_A, \boldsymbol{r}_B} \sum_{t=1}^T \left(r_{AB}^t(\boldsymbol{s}, \boldsymbol{a}) - \left(r_A(s_A^t, a_A^t) + r_B(s_B^t, a_B^t) \right) \right)^2.$$

Fully cooperative Markov games

In fully cooperative settings, only a global reward will be reviewed to all agents. Unlike the modeling in section 2, this global reward may not necessarily be decomposed as the summation of local rewards. In this case, we propose meta algorithm 2 as an extension of meta algorithm 1.

This algorithm follows similar framework of meta algorithm 1 and differs at we now learn the closet local reward decomposition from data. When the reward is completely decomposable, meta algorithm 2 recovers meta algorithm 1. Thus intuitively, the more accurate we can decompose the global reward, the less decomposition error we have. Formally, we define the measure of reward entanglement

$$e(\mathbf{r}_{AB}) \coloneqq \min_{\mathbf{r}_{A}, \mathbf{r}_{B}} \|\mathbf{r}_{AB} - (\mathbf{r}_{A} \otimes \mathbf{e} + \mathbf{e} \otimes \mathbf{r}_{B})\|_{\mu_{AB}^{\pi}}.$$
 (17)

This measure characterizes how accurate we can decompose the global reward under stationary distribution. We then obtain an extension of Theorem 4

Proposition 13. Consider a fully cooperative two-agent Markov system \mathcal{M}_{AB} . Given any policy $\pi \colon \mathcal{S} \to \Delta(\mathcal{A})$ with the measure of Markov entanglement $E_A(P_{AB}^{\pi}), E_B(P_{AB}^{\pi})$ w.r.t the μ_{AB}^{π} weighted agent-wise total variation distance and the measure of reward entanglement $e(r_{AB})$, it

$$\left\|Q_{AB}^{\pi} - \left(Q_A^{\pi} \otimes \boldsymbol{e} + \boldsymbol{e} \otimes Q_B^{\pi}\right)\right\|_{\mu_{AB}^{\pi}} \leq \frac{e(\boldsymbol{r}_{AB})}{1 - \gamma} + \frac{4\gamma \left(E_A(\boldsymbol{P}_{AB}^{\pi})r_{\max}^A + E_B(\boldsymbol{P}_{AB}^{\pi})r_{\max}^B\right)}{(1 - \gamma)^2},$$

where r_{max}^A , r_{max}^B is the bound of optimal solution of Eq. (17).

Although Proposition 1 offers a theoretical guarantee for general two-agent fully cooperative Markov games, its utility is greatest in systems with low reward and transition entanglement. Fully cooperative settings remain inherently challenging-for instance, even the asymptotically optimal Whittle Index may achieve only a $\frac{1}{N}$ -approximation ratio for RMABs with global rewards [34]. In practice, most research [38, 35] relies on sophisticated deep neural networks to learn decompositions in such settings. We thus defer a more refined analysis of fully cooperative scenarios to future work.

Simulation environments

In this section, we empirically study the value decomposition for index policies. Our simulations build on a circulant RMAB benchmark, which is widely used in the literature [3, 52, 10, 18].

Circulant RMAB A circulant RMAB has four states indexed by $\{0, 1, 2, 3\}$. Transition kernels $P_a = p(s, 0, s')_{s, s' \in S}$ for action a = 0 and a = 1 are given by

$$\boldsymbol{P}_0 = \begin{pmatrix} 1/2 & 0 & 0 & 1/2 \\ 1/2 & 1/2 & 0 & 0 \\ 0 & 1/2 & 1/2 & 0 \\ 0 & 0 & 1/2 & 1/2 \end{pmatrix}, \ \boldsymbol{P}_1 = \begin{pmatrix} 1/2 & 1/2 & 0 & 0 \\ 0 & 1/2 & 1/2 & 0 \\ 0 & 0 & 1/2 & 1/2 \\ 1/2 & 0 & 0 & 1/2 \end{pmatrix}.$$

The reward solely depends on the state and is unaffected by the action:

$$r(0, a) = -1, r(1, a) = 0, r(2, a) = 0, r(3, a) = 1; \forall a \in \{0, 1\}.$$

We set the discount factor to $\gamma=0.5$ and require N/5 arms to be pulled per period. Initially, there are N/6 arms in state 0, N/3 arms in state 1 and N/2 arms in state 2, the same as [52]. We then test an index policy with priority: state 2 > state 1 > state 0 > state 3 >

K.1 Monte-Carlo estimation of Markov entanglement

For each RMAB instance, we simulate a trajectory of length T=6N and collect data for the later 5N epochs. Notice RMAB is a special instance of WCMDP, we thus apply the result in Proposition 9

$$E_{i}(\mathbf{P}_{1:N}^{\pi}) \leq \frac{1}{2} \min_{\pi'} \sum_{\mathbf{s}} \mu_{1:N}^{\pi}(\mathbf{s}) \sum_{a_{i}} |\pi(a_{i} \mid \mathbf{s}) - \pi'(a_{i} \mid s_{i})|$$

$$\approx \frac{1}{2} \min_{\pi'} \frac{1}{T} \sum_{t=1}^{T} \sum_{a_{i}} |\pi(a_{i} \mid \mathbf{s}) - \pi'(a_{i} \mid s_{i})|$$
(18)

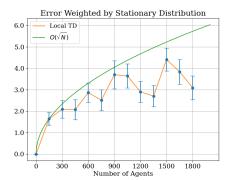
Notice Eq. (18) is *convex* for π' and π' only takes support of size |S||A|=8, we thus apply efficient convex optimization solvers. We replicate this experiment for 10 independent runs to obtain the mean estimation and standard error in the left panel of Figure 1.

K.2 Learning local Q-values

For each RMAB instance, we simulate a trajectory of length T=6N, reserving the later T=5N epochs as the training phase for each agent to fit local Q-value functions. During testing, we estimate the μ -weighted decomposition error using 50 simulations sampled from the stationary distribution.

The ground-truth $Q_{1:N}^{\pi}$ is approximated via Monte Carlo learning [39], with each estimate derived from 30-step simulations averaged over 3N independent runs. Due to the high computational cost of Monte Carlo methods—especially for very large RMABs—we limit the training phase to 10 independent runs and use the mean local Q-value as an approximation. Error bars represent the standard error for both Monte Carlo estimates and μ -weighted decomposition errors.

In addition to μ -weighted error, we also introduce a concept of relative error, defined as $\left\|Q_{1:N}^{\pi}(s, \boldsymbol{a}) - \sum_{i=1}^{N} Q_{i}^{\pi}(s_{i}, a_{i})\right\|_{\mu_{1:N}^{\pi}} / \left\|Q_{1:N}^{\pi}\right\|_{\mu_{1:N}^{\pi}}$. This relative error reflects the approximate ratio of our value decomposition. We present our simulation results below.



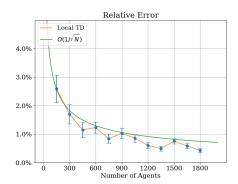


Figure 2: Value Decomposition error in circulant RMAB under an index policy. Left: μ -weighted decomposition error. Right: Relative error, $\|$ decomposition error $\|_{\mu} / \|Q_{1:N}^{\pi}\|_{\mu}$

It immediately follows that the μ -weighted error grows at a sublinear rate $\mathcal{O}(\sqrt{N})$ and the relative error decays at rate $\mathcal{O}(1/\sqrt{N})$. This justifies our theoretical guarantees in Theorem 6. Furthermore, we notice the relative error is no larger than 3% over all data points. As a result, the meta algorithm 1 is able to provide a very close approximation especially for large-scale MDPs even with small amount of training data T=5N while the global state space has size $|S|^N$.

K.3 Sample Complexity and Computation

While each RMAB instance has an exponentially large state space $|S|^N$, we show that our empirical estimation of Markov entanglement—along with the decomposition error—converges quickly. Specifically, we illustrate these errors for an RMAB instance with with 900 agents in Figure 3. As exhibits in Figure 3, both errors decay and converges within T=3N samples. Furthermore,

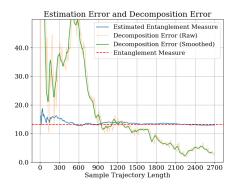


Figure 3: Different errors in RMAB with 900 agents: empirical estimation of Markov entanglement (blue); $\mu_{1:N}^{\pi}$ -weighted decomposition error (green); the true measure of Markov estimated with T=10N samples (red dashed line).

the empirical estimation of Markov entanglement converges in T < N samples, demonstrating its efficiency. Finally, we use standard convex optimization solvers to compute Markov entanglement, which can be run efficiently on a single CPU.

NeurIPS Paper Checklist

1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: [Yes]

Justification: Yes. Our main contributions are also detailed in section 1.1. Also see Appendix J, H for more theoretical results.

Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the
 contributions made in the paper and important assumptions and limitations. A No or
 NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals
 are not attained by the paper.

2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [Yes]

Justification: Yes, we discuss other possible value decompositions in section 3.1 and Appendix E, J.4.

Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

3. Theory assumptions and proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [Yes]

Justification: Theorem 1, 2, 3 and 4 hold for general multi-agent as long as a stationary distribution exists (see section 2). Theorem 6 relies on standard assumptions for index polices, detailed in Appendix I. All proofs are included in main text or Appendix.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and crossreferenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental result reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: Our empirical results build on a publicly-accessible RMAB benchmark in [3, 52, 10, 18], detailed in Appendix K. We upload the codes and instructions to recover the results.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
- (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
- (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
- (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
- (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes]

Justification: Our empirical results build on a publicly-accessible RMAB benchmark in [3, 52, 10, 18], detailed in Appendix K. We upload the codes and instructions to recover the results.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.
- While we encourage the release of code and data, we understand that this might not be
 possible, so "No" is an acceptable answer. Papers cannot be rejected simply for not
 including code, unless this is central to the contribution (e.g., for a new open-source
 benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. Experimental setting/details

Question: Does the paper specify all the training and test details (e.g., data splits, hyperparameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: See Appendix K.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

7. Experiment statistical significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [Yes]

Justification: We detail the calculation of error bars in Appendix K.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.

- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error
 of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

8. Experiments compute resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: See Appendix K.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

9. Code of ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics https://neurips.cc/public/EthicsGuidelines?

Answer: [Yes]

Justification: We followed the NeurIPS Code of Ethics.

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

10. Broader impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [No]

Justification: This work focuses on establishing a new mathematical foundation for MARL. This work is not related to any private or personal data, and there's no explicit negative social impacts.

Guidelines:

• The answer NA means that there is no societal impact of the work performed.

- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: We do not foresee any high risk for misuse of this work.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [NA]

Justification: We do not use existing assets and our empirical simulations are based on synthetic models, whose proposers have been appropriately cited.

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.

- If assets are released, the license, copyright information, and terms of use in the
 package should be provided. For popular datasets, paperswithcode.com/datasets
 has curated licenses for some datasets. Their licensing guide can help determine the
 license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

13. New assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [NA]

Justification: The paper does not release new assets.

Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

14. Crowdsourcing and research with human subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA]

Justification: The paper does not involve crowdsourcing nor research with human subjects.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

15. Institutional review board (IRB) approvals or equivalent for research with human subjects

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: The paper does not involve crowdsourcing nor research with human subjects. Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.

- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

16. Declaration of LLM usage

Question: Does the paper describe the usage of LLMs if it is an important, original, or non-standard component of the core methods in this research? Note that if the LLM is used only for writing, editing, or formatting purposes and does not impact the core methodology, scientific rigorousness, or originality of the research, declaration is not required.

Answer: [NA]

Justification: We do not involve LLMs as any important, original, or non-standard components.

Guidelines:

- The answer NA means that the core method development in this research does not involve LLMs as any important, original, or non-standard components.
- Please refer to our LLM policy (https://neurips.cc/Conferences/2025/LLM) for what should or should not be described.