

Multi-agent Markov Entanglement

Anonymous authors

Paper under double-blind review

Abstract

1 Value decomposition has long been a fundamental technique in multi-agent reinforce-
 2 ment learning (RL) and dynamic programming. Specifically, the value function of
 3 a global state (s_1, s_2, \dots, s_N) is often approximated as the sum of local functions:
 4 $V(s_1, s_2, \dots, s_N) \approx \sum_{i=1}^N V_i(s_i)$. This approach has found various applications in
 5 modern RL systems. However, the theoretical justification for why this decomposition
 6 works so effectively remains underexplored. In this paper, we uncover the underly-
 7 ing mathematical structure that enables value decomposition. We demonstrate that a
 8 Markov decision process (MDP) permits value decomposition *if and only if* its transi-
 9 tion matrix is not “entangled”—a concept analogous to quantum entanglement in quan-
 10 tum physics. Drawing inspiration from how physicists measure quantum entanglement,
 11 we introduce how to measure the “Markov entanglement” and show that this measure
 12 can be used to bound the decomposition error in general multi-agent MDPs. Using the
 13 concept of Markov entanglement, we proved that a widely-used class of policies, the
 14 index policy, is weakly-entangled and enjoys a sublinear $\mathcal{O}(\sqrt{N})$ scale of decompo-
 15 sition error for N -agent systems. Finally, we show how Markov entanglement can be
 16 efficiently estimated in practice, providing practitioners with an empirical proxy for the
 17 quality of value decomposition.

18 1 Introduction

19 Learning the value function given certain policy, or *policy evaluation*, is one of the most fundamental
 20 tasks in RL. Significant attention has been paid to single-agent policy evaluation (Sutton & Barto,
 21 2018; Bertsekas & Tsitsiklis, 1996; Tsitsiklis & Van Roy, 1996). However, when it comes to multi-
 22 agent reinforcement learning (MARL), single-agent methodologies typically suffer from *the curse*
 23 *of dimensionality*: the state space of the system scales exponentially with the number of agents. To
 24 tackle this problem, one common technique is value decomposition,

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

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

28 The remaining question is whether this decomposition is effective. This is non-trivial due to the
 29 coupling of agents—individual agent’s action and transition depend on other agents. In the past
 30 several decades, both positive and negative results have been reported. Back to the last century,
 31 Whittle (1988); Weber & Weiss (1990) apply Lagrange relaxations to decompose the global value
 32 and obtain the well-known Whittle index policy. The Lagrange decomposition idea has also been
 33 proved successful in many other important multi-agent tasks such as network revenue management
 34 (Adelman, 2007; Zhang & Adelman, 2009), resource allocation (Kadota et al., 2016; Balseiro et al.,
 35 2023), and online matching (Brown & Zhang, 2022; Shar & Jiang, 2023; Kanoria & Qian, 2024).

36 However, Lagrange decomposition relies on the knowledge of system dynamics and [Adelman &](#)
 37 [Mersereau \(2008\)](#) show its decomposition error can be arbitrarily bad for general multi-agent MDPs.
 38 In more recent days, practitioners apply online (deep) reinforcement learning to train a local value
 39 function for each individual agent. This practice gives birth to state-of-the-art dispatching policies
 40 in ride-hailing platforms and has been well recognized by the operations research community, such
 41 as DiDi Chuxing ([Qin et al. \(2020\)](#), [Daniel H. Wagner Prize, 2020](#)) and Lyft ([Azagirre et al. \(2024\)](#),
 42 [Franz Edelman Laureates, 2024](#)). Intervention policies based on a similar value decomposition idea
 43 also demonstrate substantial empirical advantages and have been deployed by a behavioral health
 44 platform in Kenya ([Baek et al. \(2023\)](#), [Pierskalla Award, 2024](#)). In broader MARL literature, value
 45 decomposition serves as one key component of centralized training and decentralized execution
 46 (CTDE) paradigm, achieving strong empirical performance ([Sunehag et al., 2018](#); [Mahajan et al.,](#)
 47 [2019](#); [Rashid et al., 2020](#)). However, recent research has started reflecting on the invalidity and
 48 potential flaw of value decomposition in practice ([Hong et al., 2022](#); [Dou et al., 2022](#)).

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

51 1.1 This paper

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

57 We start with a trivial example where two agents are independent, i.e. each following independent
 58 MDPs. It’s clear that the global value function can be decomposed as the sum of value functions of
 59 local MDPs. As two agents are independent, it holds $P^\pi(s'_A, s'_B | s_A, s_B) = P^\pi(s'_A | s_A) \cdot P^\pi(s'_B |$
 60 $s_B)$, or in matrix form,

$$P_{AB}^\pi = P_A^\pi \otimes P_B^\pi,$$

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

63 **A Sufficient and Necessary Condition** We introduce a new condition called “Markov Entangle-
 64 ment” 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*

$$P_{AB}^\pi = \sum_{j=1}^K x_j P_A^{(j)} \otimes P_B^{(j)},$$

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

65

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

$$\text{separable } P_{AB}^\pi \iff \text{decomposable } V_{AB}^\pi,$$

69 where V_{AB}^π is decomposable if there exist local value functions V_A, V_B such that $V_{AB}^\pi(s_A, s_B) =$
 70 $V_A(s_A) + V_B(s_B)$ for all (s_A, s_B) . This result sharply unravels the secret structure of system
 71 dynamics governing value decomposition. As a sufficient condition, our finding strictly generalizes
 72 the previous independent subsystem example, extending it to scenarios involving interacting and
 73 coupled agents. As a necessary condition, we prove that exact value decomposition under any

74 reward requires the system dynamics to be separable. Taken together, this result provides a *complete*
 75 *characterization* of when exact value function decomposition is possible in multi-agent MDPs.

76 More interestingly, our Markov entanglement condition turns out to be a mathematical counterpart of
 77 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*

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

then ρ_{AB} is separable; otherwise entangled.

78

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

83 **Decomposition Error in General Multi-agent MDPs** General multi-agent MDPs can exhibit
 84 arbitrary complexity, with agents intricately entangled. This raises a critical question: *can value de-*
 85 *composition serve as a meaningful approximation in such scenarios?* To address this, we introduce
 86 a mathematical quantification to measure the Markov entanglement in general multi-agent MDPs,

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

87 where \mathcal{P}_{SEP} is the set of all separable transition matrices and $d(\cdot, \cdot)$ is some distance measure. In
 88 other words, the degree of Markov entanglement is determined by its distance to the closest separable
 89 transition matrix. This concept also has a counterpart in quantum entanglement measurement.

$$E(\rho_{AB}) := \min_{\rho \in \rho_{\text{SEP}}} d(\rho_{AB}, \rho),$$

90 where ρ_{SEP} is the set of all separable quantum states. In quantum physics, various distance mea-
 91 sures have been designed for density matrices and capture different physical interpretations (Nielsen
 92 & Chuang, 2010). In the Markov world, we analogously design distance measures for transition
 93 matrices and relate them to the value decomposition error,

$$\left\| \text{decomposition error of } \mathbf{V}_{AB}^\pi \right\| = \mathcal{O}\left(E(\mathbf{P}_{AB}^\pi)\right).$$

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

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

106 1.2 Other related work

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

109 Prior theoretical research has extensively investigated the decomposition of optimal value functions
 110 in multi-agent settings. A prominent area involves decomposition via Lagrange relaxation. The per-
 111 agent decomposition error is proven to decay asymptotically to zero (Weber & Weiss, 1990; 1991;
 112 Verloop, 2016) and enjoys a quadratic or exponential rate (Gast et al., 2023; 2024; Brown & Zhang,
 113 2022; Zhang & Frazier, 2021; 2022). Other work generalizes to Weakly-Coupled MDPs (WCMDPs)
 114 (Balseiro et al., 2021; Brown & Zhang, 2025; Gast et al., 2022). Despite these advancements,
 115 characterizing decomposition error for general multi-agent MDPs remains unknown. In contrast,
 116 our Markov entanglement theory analyzes value decomposition for general multi-agent MDPs under
 117 arbitrary policies, including optimal ones.

118 Another line of theoretical work has concentrated on policy optimization via value decomposition.
 119 Despite reported empirical successes, rigorous theoretical analysis remains challenging. Baek et al.
 120 (2023) derived an approximation ratio for a specific index policy on a two-state RMAB. Wang et al.
 121 (2021); Dou et al. (2022) analyzed the convergence of the CTDE paradigm under strong exploration
 122 assumptions, while also highlighting scenarios of divergence. In contrast, our work instead focuses
 123 on policy evaluation rather than optimization. This enables us to derive clear and interpretable
 124 bounds on the decomposition error for general finite-state multi-agent MDPs that only require the
 125 existence of a stationary distribution.

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

129 2 Model

130 We consider a standard two-agent MDP $\mathcal{M}_{AB}(\mathcal{S}, \mathcal{A}, \mathbf{P}, \mathbf{r}_A, \mathbf{r}_B, \gamma)$ with joint state space $\mathcal{S} = \mathcal{S}_A \times$
 131 \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
 132 $|\mathcal{S}_A| = |\mathcal{S}_B| = |\mathcal{S}|$ and $|\mathcal{A}_A| = |\mathcal{A}_B| = |\mathcal{A}|$. For agents at global state $\mathbf{s} = (s_A, s_B)$ with
 133 action $\mathbf{a} = (a_A, a_B)$ taken, the system will transit to $\mathbf{s}' = (s'_A, s'_B)$ according to transition kernel
 134 $\mathbf{s}' \sim \mathbf{P}(\cdot | \mathbf{s}, \mathbf{a})$ and each agent $i \in \{A, B\}$ will receive its local reward $r_i(s_i, a_i)$. The global
 135 reward r_{AB} is defined as the summation of local rewards $r_{AB}(\mathbf{s}, \mathbf{a}) := r_A(s_A, a_A) + r_B(s_B, a_B)$,
 136 or in vector form $\mathbf{r}_{AB} \in \mathbb{R}^{|\mathcal{S}|^2|\mathcal{A}|^2} := \mathbf{r}_A \otimes \mathbf{e} + \mathbf{e} \otimes \mathbf{r}_B$, where \otimes is the tensor product and
 137 $\mathbf{e} = \mathbf{1} \in \mathbb{R}^{|\mathcal{S}|}$ is the vector of all ones.¹ We further assume the local rewards are bounded, i.e. for
 138 agent $i \in \{A, B\}$, $|r_i(s_i, a_i)| \leq r_{\max}^i$ for all (s_i, a_i) .

139 Given any global policy $\pi: \mathcal{S} \rightarrow \Delta(\mathcal{A})$, the global Q-value under policy π is defined as the dis-
 140 counted summation of global rewards $Q_{AB}^\pi(\mathbf{s}, \mathbf{a}) = \mathbb{E}[\sum_{t=0}^{\infty} \gamma^t r_{AB}(\mathbf{s}^t, \mathbf{a}^t) | \pi, (\mathbf{s}^0, \mathbf{a}^0) = (\mathbf{s}, \mathbf{a})]$
 141 where $\gamma \in [0, 1)$ is the discount factor. The value function is then defined as $V_{AB}^\pi(\mathbf{s}) =$
 142 $\mathbb{E}_{\mathbf{a} \sim \pi(\cdot | \mathbf{s})}[Q_{AB}^\pi(\mathbf{s}, \mathbf{a})]$. We denote $\mathbf{P}_{AB}^\pi \in \mathbb{R}^{|\mathcal{S}|^2|\mathcal{A}|^2 \times |\mathcal{S}|^2|\mathcal{A}|^2}$ as the transition matrix induced by
 143 π where $\mathbf{P}_{AB}^\pi(\mathbf{s}', \mathbf{a}' | \mathbf{s}, \mathbf{a}) = \mathbf{P}(\mathbf{s}' | \mathbf{s}, \mathbf{a}) \cdot \pi(\mathbf{a}' | \mathbf{s}')$. Then by the Bellman Equation, we have
 144 $Q_{AB}^\pi = (\mathbf{I} - \gamma \mathbf{P}_{AB}^\pi)^{-1} \mathbf{r}_{AB}$. Our objective is to decompose this global Q-value Q_{AB}^π as the sum-
 145 mation of some local functions Q_A and Q_B , i.e. $Q_{AB}^\pi(\mathbf{s}, \mathbf{a}) = Q_A(s_A, a_A) + Q_B(s_B, a_B)$, or in
 146 vector form,

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

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

151 2.1 Local (Q-)value functions

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

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

154 This meta-algorithm framework is simple and intuitive: each agent independently fits its local Q-
 155 values based on its local observations. Notably, the framework requires no prior knowledge of the
 156 MDP, and learning can be performed in a fully decentralized manner. Furthermore, we use term
 157 *meta* in that we do not pose restrictions on how agents estimate their local Q-values. For tabular
 158 case, one can plug in Temporal Difference (TD) learning (Sutton & Barto, 2018) or its variants. For
 159 large-scale problems, one can apply linear function approximations (Baek et al., 2023; Han et al.,
 160 2022; Bertsekas & Tsitsiklis, 1996) or more sophisticated neural networks (Qin et al., 2020; Sunehag
 161 et al., 2018; Mahajan et al., 2019).

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

165 **Local TD learning.** Although each agent’s environment is not Markovian in a local sense (it is, more
 166 precisely, partially observed Markovian), one can still define its “marginalized” local transition ma-
 167 trix under the stationary distribution. Mathematically, for agent A , we denote $\mathbf{P}_A^\pi \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}| \times |\mathcal{S}||\mathcal{A}|}$
 168 as its local transition where

$$P_A^\pi(s'_A, a'_A | s_A, a_A) = \sum_{s'_B, a'_B} \sum_{s_B, a_B} P_{AB}^\pi(s'_{AB}, a'_{AB} | s_{AB}, a_{AB}) \mu_{AB}^\pi(s_B, a_B | s_A, a_A). \quad (3)$$

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

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

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

175 3 Exact value decomposition

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

180 **Theorem 1.** Consider a two-agent MDP \mathcal{M}_{AB} and policy $\pi: \mathcal{S} \rightarrow \Delta(\mathcal{A})$. If two agents are
 181 separable, then the Eq. (2) holds

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

182 This theorem shows that even when the system is not independent, as long as it can be represented as
 183 a *linear combination of independent subsystems*, the global Q-value admits an exact decomposition.
 184 In Appendix D, we provide an MDP instance where agents are coupled but not entangled.

²For $\mu_{AB}^\pi(s_B, a_B | s_A, a_A)$ 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 $\mathcal{A}(s_A)$. All theoretical results apply to the remaining valid state-action pairs.

185 3.1 Necessary condition for the exact value decomposition

186 We then investigate whether Markov entanglement is necessary for the exact Q-value decomposition.
 187 The answer is in general no, since one can construct trivial counterexamples such as $\mathbf{r}_A = \mathbf{r}_B = \mathbf{0}$
 188 or $\gamma = 0$, where the decomposition trivially holds. On the other hand, we focus on a stronger and
 189 more general concept of the exact value decomposition that holds under any reward kernel given
 190 $\gamma > 0$. Formally, we present the following theorem.

191 **Theorem 2.** Consider a two-agent Markov MDP \mathcal{M}_{AB} with discount factor $\gamma > 0$ and $\pi: \mathcal{S} \rightarrow$
 192 $\Delta(\mathcal{A})$. Suppose there exists local functions $Q_i: \mathbf{r}_i \rightarrow \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}$ for $i \in \{A, B\}$ such that $Q_{AB}^\pi =$
 193 $Q_A(\mathbf{r}_A) \otimes \mathbf{e} + \mathbf{e} \otimes Q_B(\mathbf{r}_B)$ holds for any pair of reward $\mathbf{r}_A, \mathbf{r}_B$, then A, B must be separable.

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

198 4 Value decomposition error in general two-agent MDPs

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

204 4.1 Entry-wise error bound

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

$$\|\mathbf{P} - \mathbf{P}'\|_{\text{TV}} := \max_{(\mathbf{s}, \mathbf{a}) \in \mathcal{S} \times \mathcal{A}} D_{\text{TV}}(\mathbf{P}(\cdot, \cdot | \mathbf{s}, \mathbf{a}), \mathbf{P}'(\cdot, \cdot | \mathbf{s}, \mathbf{a})), \quad (4)$$

207 where D_{TV} is the total variation distance between probability measures.

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

$$\|\mathbf{P} - \mathbf{P}'\|_{\text{ATV}_A} := \max_{(\mathbf{s}, \mathbf{a}) \in \mathcal{S} \times \mathcal{A}} D_{\text{TV}} \left(\sum_{s'_B, a'_B} \mathbf{P}(\cdot, \cdot | \mathbf{s}, \mathbf{a}), \sum_{s'_B, a'_B} \mathbf{P}'(\cdot, \cdot | \mathbf{s}, \mathbf{a}) \right). \quad (5)$$

211 The ATV distance w.r.t agent B can be defined similarly. Intuitively, compared to TV, ATV fo-
 212 cuses on an individual agent and measures the difference between its local transitions. One can
 213 also verify ATV is tighter distance, i.e. $\|\mathbf{P} - \mathbf{P}'\|_{\text{ATV}_A} \leq \|\mathbf{P} - \mathbf{P}'\|_{\text{TV}}$. We can plug ATV
 214 into Eq. (1) and obtain the measure of Markov entanglement w.r.t ATV distance $E_i(\mathbf{P}_{AB}^\pi) :=$
 215 $\min_{\mathbf{P} \in \mathcal{P}_{\text{SEP}}} \|\mathbf{P}_{AB}^\pi - \mathbf{P}\|_{\text{ATV}_i}$ for $i \in \{A, B\}$. In fact, one can also verify

$$E_A(\mathbf{P}_{AB}^\pi) = \min_{\mathbf{P}_A} \max_{(\mathbf{s}, \mathbf{a}) \in \mathcal{S} \times \mathcal{A}} D_{\text{TV}} \left(\mathbf{P}_{AB}^\pi(\cdot, \cdot | \mathbf{s}, \mathbf{a}), \mathbf{P}_A(\cdot, \cdot | s_A, a_A) \right), \quad (6)$$

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

217 **Theorem 3.** Consider a two-agent Markov system \mathcal{M}_{AB} and policy $\pi: \mathcal{S} \rightarrow \Delta(\mathcal{A})$ with the mea-
 218 sure of Markov entanglement $E_A(\mathbf{P}_{AB}^\pi), E_B(\mathbf{P}_{AB}^\pi)$ defined in Eq. (6), then the decomposition error
 219 is entry-wise bounded by the measure of Markov entanglement,

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

220 4.2 Error weighted by stationary distribution

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

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

224 A stationary distribution weighted error bound is common in policy evaluation literature (Cai et al.,
 225 2019; Tsitsiklis & Van Roy, 1996; Bhandari et al., 2021).

226 **Distance weighted by stationary distribution** To analyze this μ_{AB}^π -weight decomposition error,
 227 we analogously propose the μ_{AB}^π -weighted distance measure of Markov entanglement. Specifically,
 228 we have the following μ_{AB}^π -weighted version of Eq. (6).

$$E_A(\mathbf{P}_{AB}^\pi) = \min_{\mathbf{P}_A} \sum_{\mathbf{s}, \mathbf{a}} \mu_{AB}^\pi(\mathbf{s}, \mathbf{a}) D_{\text{TV}} \left(\mathbf{P}_{AB}^\pi(\cdot, \cdot | \mathbf{s}, \mathbf{a}), \mathbf{P}_A(\cdot, \cdot | s_A, a_A) \right). \quad (7)$$

229 Eq. (7) substitutes the μ_{AB}^π -weighted average for the maximum operator in Eq. (6). Finally, we have
 230 the following variant of Theorem 3.

231 **Theorem 4.** *Under the same setup as Theorem 3 with μ_{AB}^π -weighted measure of Markov entangle-*
 232 *ment $E_A(\mathbf{P}_{AB}^\pi), E_B(\mathbf{P}_{AB}^\pi)$ defined in Eq. (7), the μ_{AB}^π -weighted decomposition error is bounded,*

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

233 Finally it's straightforward to extend our results to multi-agent MDPs, detailed in Appendix J.

234 5 Applications of Markov Entanglement

235 In this section, we apply Markov entanglement and demonstrate a widely-used class of index policies
 236 is asymptotically separable. To begin, we introduce the model of Restless Multi-Armed Bandit
 237 (RMAB, Whittle (1988)). In an N -agent RMAB, each agent follows a homogeneous two-action
 238 MDP with action 1 meaning activate and 0 idle. A central decision maker will activate $M \leq N$
 239 agents at each timestep and leave other agents idle. In other words, agents transit independently but
 240 are coupled under constraint $\sum_{i=1}^N a_i = M$. In RMAB, arguably the most classical and widely-used
 241 policy is the index policy, which we formally define as

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

245 The index policy traces back to the well-known Gittins Index (Weber, 1992), Whittle Index (Whittle,
 246 1988; Weber & Weiss, 1990; Gast et al., 2023), and fluid-based index policies (Verloop, 2016; Gast
 247 et al., 2024). Qin et al. (2020); Azagirre et al. (2024); Baek et al. (2023); Nakhleh et al. (2021); Wang
 248 et al. (2023); Avrachenkov & Borkar (2022) apply data-driven method to optimize index policies
 249 and report great empirical success in industrial implementations. Understanding the mystery behind
 250 such success calls for a theory for general index policies. We then present our main theorem.

251 **Theorem 5.** *Consider an N -agent restless multi-armed bandit. For any index policy satisfying mild*
 252 *technical conditions, there exists constant C independent of N , such that for any agent $i \in [N]$, its*
 253 *$\mu_{1:N}^\pi$ -weighted measure of Markov entanglement is bounded, $E_i(\mathbf{P}_{1:N}^\pi) \leq C/\sqrt{N}$.*

254 Theorem 5 requires two standard technical conditions for index policies: non-degenerate and uni-
 255 form global attractor property, detailed in Appendix K. Theorem 5 justifies index policies are asymp-
 256 totically separable. Combined with an N -agent version of Theorem 4, we obtain the sublinear
 257 decomposition error for index policies

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

258 This sublinear error result explains why the value decomposition in [Qin et al. \(2020\)](#); [Azagirre et al.](#)
 259 [\(2024\)](#); [Baek et al. \(2023\)](#) manages to effectively approximate the global value function in large-
 260 scale practical applications.

261 5.1 Efficient verification of value decomposition

262 For practitioners, verifying the feasibility of value decomposition is challenging due to the exponen-
 263 tial computational complexity of estimating the global Q-value. As a solution, Markov entanglement
 264 offers an efficient way to empirically test whether value decomposition can be safely applied. Con-
 265 sider the μ_{AB}^π -weighted measure of Markov entanglement in Eq. (7), we have

$$E_A(\mathbf{P}_{AB}^\pi) \approx \frac{1}{2} \min_{\mathbf{P}_A} \frac{1}{T} \sum_{t=1}^T \sum_{s'_A, a'_A} |\mathbf{P}_{AB}^\pi(s'_A, a'_A | \mathbf{s}^t, \mathbf{a}^t) - \mathbf{P}_A(s'_A, a'_A | \mathbf{s}^t, \mathbf{a}^t)| \quad (8)$$

266 In other words, we can apply a Monte-Carlo estimation for $E_A(\mathbf{P}_{AB}^\pi)$. Notice Eq. (8) is *convex* for
 267 \mathbf{P}_A , which enables efficient solutions. As a result, Eq. (8) provides an efficient estimation of Markov
 268 entanglement via simulation and can be easily extend to N -agent MDPs.

269 **Numerical experiments.** Finally, we empirically study the value decomposition for the index policy
 270 on a circulant RMAB benchmark ([Avrachenkov & Borkar, 2022](#); [Zhang & Frazier, 2022](#); [Biswas](#)
 271 [et al., 2021](#); [Fu et al., 2019](#)) that has 4 different states each local agent. As a result, the global state
 272 space scales as large as $4^{1800} > 10^{1000}$ for $N = 1800$ agents. The specific transitions and rewards
 273 are introduced in Appendix M. For each RMAB instance, we sample a trajectory of length $T = 5N$
 274 and use the collected data to i) solve Eq. (8) to estimate the measure of Markov entanglement; ii)
 275 train local Q-value decomposition. It quickly follows from the results in Figure 1:

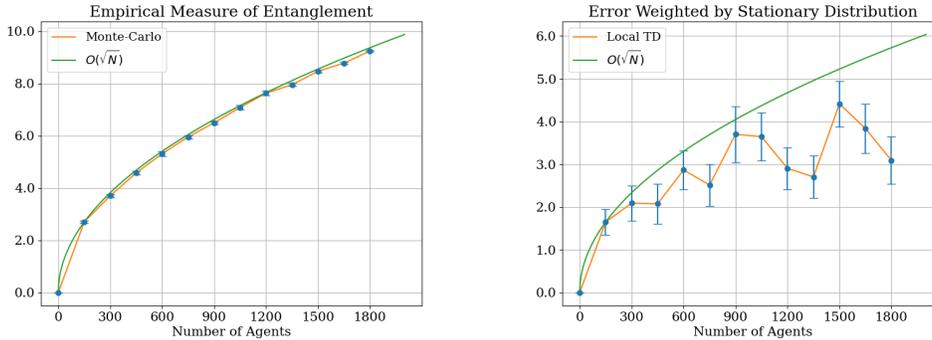


Figure 1: Circulant RMAB under an index policy. *Left*: empirical estimation of Markov entangle-
 ment multiplied by the number of agents, $NE_1(\mathbf{P}_{1:N}^\pi)$. *Right*: μ -weighted decomposition error.

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

280 281 6 Conclusion

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

289 **References**

- 290 Daniel Adelman. Dynamic bid prices in revenue management. *Operations Research*, 55(4):647–
291 661, 2007.
- 292 Daniel Adelman and Adam J. Mersereau. Relaxations of weakly coupled stochastic dynamic pro-
293 grams. *Operations Research*, 56(3):712–727, 2008. DOI: 10.1287/opre.1070.0445.
- 294 Konstantin E Avrachenkov and Vivek S Borkar. Whittle index based q-learning for restless bandits
295 with average reward. *Automatica*, 139:110186, 2022.
- 296 Xabi Azaguirre, Akshay Balwally, Guillaume Candeli, Nicholas Chamandy, Benjamin Han, Alona
297 King, Hyungjun Lee, Martin Loncaric, Sébastien Martin, Vijay Narasiman, Zhiwei (Tony) Qin,
298 Baptiste Richard, Sara Smoot, Sean Taylor, Garrett van Ryzin, Di Wu, Fei Yu, and Alex
299 Zamoshchin. A better match for drivers and riders: Reinforcement learning at lyft. *INFORMS*
300 *Journal on Applied Analytics*, 54(1):71–83, 2024.
- 301 Jackie Baek, Justin J Boutilier, Vivek F Farias, Jonas Oddur Jonasson, and Erez Yoeli. Policy op-
302 timization for personalized interventions in behavioral health. *arXiv preprint arXiv:2303.12206*,
303 2023.
- 304 Santiago R. Balseiro, David B. Brown, and Chen Chen. Dynamic pricing of relocating resources in
305 large networks. *Management Science*, 67(7):4075–4094, 2021.
- 306 Santiago R. Balseiro, Haihao Lu, and Vahab Mirrokni. The best of many worlds: Dual mirror
307 descent for online allocation problems. *Operations Research*, 71(1):101–119, 2023.
- 308 Dimitri Bertsekas and John N Tsitsiklis. *Neuro-dynamic programming*. Athena Scientific, 1996.
- 309 Jalaj Bhandari, Daniel Russo, and Raghav Singal. A finite time analysis of temporal difference
310 learning with linear function approximation. *Operations Research*, 69(3):950–973, 2021.
- 311 Arpita Biswas, Gaurav Aggarwal, Pradeep Varakantham, and Milind Tambe. Learn to intervene: An
312 adaptive learning policy for restless bandits in application to preventive healthcare. In *Proceedings*
313 *of the 30th International Joint Conference on Artificial Intelligence (IJCAI 2021)*, pp. 4036–4049,
314 2021.
- 315 David B Brown and Jingwei Zhang. Dynamic programs with shared resources and signals: Dynamic
316 fluid policies and asymptotic optimality. *Operations Research*, 70(5):3015–3033, 2022.
- 317 David B. Brown and Jingwei Zhang. Technical note—on the strength of relaxations of weakly
318 coupled stochastic dynamic programs. *Operations Research*, 71(6):2374–2389, 2023. DOI: 10.
319 1287/opre.2022.2287.
- 320 David B. Brown and Jingwei Zhang. Fluid policies, reoptimization, and performance guarantees in
321 dynamic resource allocation. *Operations Research*, 73(2):1029–1045, 2025.
- 322 Qi Cai, Zhuoran Yang, Jason D Lee, and Zhaoran Wang. Neural temporal-difference learning con-
323 verges to global optima. In *Advances in Neural Information Processing Systems*, volume 32,
324 2019.
- 325 Shao-Hung Chan, Zhe Chen, Teng Guo, Han Zhang, Yue Zhang, Daniel Harabor, Sven Koenig,
326 Cathy Wu, and Jingjin Yu. The league of robot runners competition: Goals, designs, and imple-
327 mentation. In *ICAPS 2024 System’s Demonstration track*, 2024.
- 328 Zehao Dou, Jakub Grudzien Kuba, and Yaodong Yang. Understanding value decomposition algo-
329 rithms in deep cooperative multi-agent reinforcement learning. *arXiv preprint arXiv:2202.04868*,
330 2022.

- 331 Vivek Farias, Hao Li, Tianyi Peng, Xinyuyang Ren, Huawei Zhang, and Andrew Zheng. Correct-
332 ing for interference in experiments: A case study at douyin. In *Proceedings of the 17th ACM*
333 *Conference on Recommender Systems*, pp. 455–466, 2023.
- 334 Jing Fu, Yoni Nazarathy, Sarat Moka, and Peter G. Taylor. Towards q-learning the whittle index for
335 restless bandits. In *2019 Australian New Zealand Control Conference (ANZCC)*, 2019.
- 336 Nicolas Gast, Bruno Gaujal, and Chen Yan. Reoptimization nearly solves weakly coupled markov
337 decision processes. *arXiv preprint arXiv:2211.01961*, 2022.
- 338 Nicolas Gast, Bruno Gaujal, and Chen Yan. Exponential asymptotic optimality of whittle index
339 policy. *Queueing Syst. Theory Appl.*, 104(1–2):107–150, may 2023.
- 340 Nicolas Gast, Bruno Gaujal, and Chen Yan. Linear program-based policies for restless bandits:
341 Necessary and sufficient conditions for (exponentially fast) asymptotic optimality. *Mathematics*
342 *of Operations Research*, 49(4):2468–2491, 2024.
- 343 Carlos Guestrin, Daphne Koller, and Ronald Parr. Multiagent planning with factored mdps. In
344 *Advances in Neural Information Processing Systems*, volume 14. MIT Press, 2001.
- 345 Carlos Guestrin, Daphne Koller, Ronald Parr, and Shobha Venkataraman. Efficient solution algo-
346 rithms for factored mdps. *Journal of Artificial Intelligence Research*, 19(1):399–468, 2003.
- 347 Benjamin Han, Hyungjun Lee, and Sébastien Martin. Real-time rideshare driver supply values
348 using online reinforcement learning. In *Proceedings of the 28th ACM SIGKDD Conference on*
349 *Knowledge Discovery and Data Mining*, KDD ’22, pp. 2968–2976, 2022.
- 350 Yitian Hong, Yaochu Jin, and Yang Tang. Rethinking individual global max in cooperative multi-
351 agent reinforcement learning. *Advances in neural information processing systems*, 35:32438–
352 32449, 2022.
- 353 Igor Kadota, Elif Uysal-Biyikoglu, Rahul Singh, and Eytan Modiano. Minimizing the age of infor-
354 mation in broadcast wireless networks. In *2016 54th Annual Allerton Conference on Communi-*
355 *cation, Control, and Computing (Allerton)*, pp. 844–851. IEEE, 2016.
- 356 Yash Kanoria and Pengyu Qian. Blind dynamic resource allocation in closed networks via mirror
357 backpressure. *Management Science*, 70(8):5445–5462, 2024.
- 358 Anuj Mahajan, Tabish Rashid, Mikayel Samvelyan, and Shimon Whiteson. Maven: Multi-agent
359 variational exploration. *Advances in neural information processing systems*, 32, 2019.
- 360 Khaled Nakhleh, Santosh Ganji, Ping-Chun Hsieh, I-Hong Hou, and Srinivas Shakkottai. Neurwin:
361 Neural whittle index network for restless bandits via deep rl. In *Advances in Neural Information*
362 *Processing Systems*, volume 34, pp. 828–839, 2021.
- 363 Michael A Nielsen and Isaac L Chuang. *Quantum computation and quantum information*. Cam-
364 bridge university press, 2010.
- 365 Ian Osband and Benjamin Van Roy. Near-optimal reinforcement learning in factored mdps. In
366 *Proceedings of the 28th International Conference on Neural Information Processing Systems*,
367 NIPS’14, pp. 604–612, 2014.
- 368 Zhiwei (Tony) Qin, Xiaocheng Tang, Yan Jiao, Fan Zhang, Zhe Xu, Hongtu Zhu, and Jieping Ye.
369 Ride-hailing order dispatching at didi via reinforcement learning. *INFORMS Journal on Applied*
370 *Analytics*, 50(5):272–286, 2020. DOI: 10.1287/inte.2020.1047.
- 371 Naveen Janaki Raman, Zheyuan Ryan Shi, and Fei Fang. Global rewards in restless multi-armed
372 bandits. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*,
373 2024.

- 374 Tabish Rashid, Mikayel Samvelyan, Christian Schroeder De Witt, Gregory Farquhar, Jakob Foerster,
375 and Shimon Whiteson. Monotonic value function factorisation for deep multi-agent reinforcement
376 learning. *Journal of Machine Learning Research*, 21(178):1–51, 2020.
- 377 Ibrahim El Shar and Daniel R. Jiang. Weakly coupled deep q-networks. In *Thirty-seventh Confer-*
378 *ence on Neural Information Processing Systems*, 2023.
- 379 Kyunghwan Son, Daewoo Kim, Wan Ju Kang, David Earl Hostallero, and Yung Yi. QTRAN:
380 Learning to factorize with transformation for cooperative multi-agent reinforcement learning. In
381 *Proceedings of the 36th International Conference on Machine Learning*, volume 97, pp. 5887–
382 5896. PMLR, 2019.
- 383 Peter Sunehag, Guy Lever, Audrunas Gruslys, Wojciech Marian Czarnecki, Vinicius Zambaldi, Max
384 Jaderberg, Marc Lanctot, Nicolas Sonnerat, Joel Z. Leibo, Karl Tuyls, and Thore Graepel. Value-
385 decomposition networks for cooperative multi-agent learning based on team reward. In *Pro-*
386 *ceedings of the 17th International Conference on Autonomous Agents and MultiAgent Systems*,
387 AAMAS '18, pp. 2085–2087, 2018.
- 388 Richard S. Sutton and Andrew G. Barto. *Reinforcement Learning: An Introduction*. The MIT Press,
389 2018. ISBN 0262039249.
- 390 John Tsitsiklis and Benjamin Van Roy. Analysis of temporal-difference learning with function
391 approximation. In *Advances in Neural Information Processing Systems*, volume 9. MIT Press,
392 1996.
- 393 Ina Maria Verloop. Asymptotically optimal priority policies for indexable and nonindexable restless
394 bandits. *Annals of Applied Probability*, 26:1947–1995, 2016.
- 395 Jianhao Wang, Zhizhou Ren, Terry Liu, Yang Yu, and Chongjie Zhang. Qplex: Duplex dueling
396 multi-agent q-learning. *arXiv preprint arXiv:2008.01062*, 2020.
- 397 Jianhao Wang, Zhizhou Ren, Beining Han, Jianing Ye, and Chongjie Zhang. Towards understanding
398 cooperative multi-agent q-learning with value factorization. *Advances in Neural Information*
399 *Processing Systems*, 34:29142–29155, 2021.
- 400 Kai Wang, Lily Xu, Aparna Taneja, and Milind Tambe. Optimistic whittle index policy: Online
401 learning for restless bandits. In *Proceedings of the AAAI Conference on Artificial Intelligence*,
402 volume 37, pp. 10131–10139, 2023.
- 403 Richard Weber. On the gittins index for multiarmed bandits. *The Annals of Applied Probability*, 2
404 (4):1024 – 1033, 1992.
- 405 Richard R. Weber and Gideon Weiss. On an index policy for restless bandits. *Journal of Applied*
406 *Probability*, 27(3):637–648, 1990.
- 407 Richard R. Weber and Gideon Weiss. Addendum to ‘on an index policy for restless bandits’. *Ad-*
408 *vances in Applied Probability*, 23(2):429–430, 1991.
- 409 P. Whittle. Restless bandits: Activity allocation in a changing world. *Journal of Applied Probability*,
410 25:287–298, 1988.
- 411 Dan Zhang and Daniel Adelman. An approximate dynamic programming approach to network
412 revenue management with customer choice. *Transportation Science*, 43(3):381–394, 2009.
- 413 Xiangyu Zhang and Peter I Frazier. Restless bandits with many arms: Beating the central limit
414 theorem. *arXiv preprint arXiv:2107.11911*, 2021.
- 415 Xiangyu Zhang and Peter I Frazier. Near-optimality for infinite-horizon restless bandits with many
416 arms. *arXiv preprint arXiv:2203.15853*, 2022.

417 A Linear algebra with tensor product

418 We briefly introduce the basic properties of tensor product or Kronecker product. Let $A \in$
419 $\mathbb{R}^{m_1 \times n_1}$, $B \in \mathbb{R}^{m_2 \times n_2}$, then

$$A \otimes B = \begin{bmatrix} a_{11}B & a_{12}B & \dots & a_{1n_1}B \\ a_{21}B & a_{22}B & \dots & a_{2n_1}B \\ \dots & \dots & \dots & \dots \\ a_{m_11}B & a_{m_12}B & \dots & a_{m_1n_1}B \end{bmatrix} \in \mathbb{R}^{m_1 m_2 \times n_1 n_2}.$$

420 Tensor product satisfies the following basic properties,

- 421 • **1. Bilinearity** For any matrix A, B, C and constant k , it holds $k(A \otimes B) = (kA) \otimes B =$
422 $A \otimes (kB)$, $(A + B) \otimes C = A \otimes C + B \otimes C$, and $A \otimes (B + C) = A \otimes B + A \otimes C$.
- 423 • **2. Mixed-product Property** For any matrix A, B, C, D , if AC and BD form valid matrix
424 product, then $(A \otimes B)(C \otimes D) = (AC) \otimes (BD)$.

425 B Decompose value functions

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

$$V_{AB}^\pi = (\mathbf{I} - \gamma P_{AB}^\pi)^{-1} r_{AB}^\pi,$$

428 where we abuse notation and denote $P_{AB}^\pi(s' | s) = \sum_{\mathbf{a}} \pi(\mathbf{a} | s) P(s' | s, \mathbf{a})$ and reward $r_{AB}^\pi(s) =$
429 $\sum_{\mathbf{a}} \pi(\mathbf{a} | s) r_{AB}(s, \mathbf{a})$. A key subtlety arises because r_{AB}^π may not be decomposable—even when
430 r_{AB} is decomposable—unless the reward r_{AB} is state-dependent. Consequently, we cannot directly
431 apply the "absorbing" equation as in the proof of Theorem 1.

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

436 C Proof of Sufficiency

437 Theorem 1 admits a simple proof based on the several basic properties of tensor product. First of
438 all, given $P_{AB}^\pi = \sum_{j=1}^K x_j P_A^{(j)} \otimes P_B^{(j)}$, we have

$$P_{AB}^\pi(s'_A, s'_B, a'_A, a'_B | s_A, s_B, a_A, a_B) = \sum_{j=1}^K x_j P_A^{(j)}(s'_A, a'_A | s_A, a_A) P_B^{(j)}(s'_B, s'_B | s_B, a_B).$$

439 Recall P_A^π in Eq. (3), it's evident that

$$\begin{aligned} P_A^\pi(s'_A, a'_A | s_A, a_A) &= \sum_{s'_B, a'_B} \sum_{s_B, a_B} \sum_{j=1}^K x_j P_A^{(j)}(s'_A, a'_A | s_A, a_A) P_B^{(j)}(s'_B, s'_B | s_B, a_B) \mu_{AB}^\pi(s_B, a_B | s_A, a_A) \\ &= \sum_{j=1}^K x_j P_A^{(j)}(s'_A, a'_A | s_A, a_A) \sum_{s_B, a_B} \mu_{AB}^\pi(s_B, a_B | s_A, a_A) \sum_{s'_B, a'_B} P_B^{(j)}(s'_B, s'_B | s_B, a_B) \\ &= \sum_{j=1}^K x_j P_A^{(j)}(s'_A, a'_A | s_A, a_A), \end{aligned}$$

440 where the second last equation holds by rearranging the summation. This leads to $\mathbf{P}_A^\pi =$
 441 $\sum_{i=1}^K x_i \mathbf{P}_A^{(i)}$. It remains to show Eq. (2), and notice that

$$\begin{aligned} (\mathbf{I} - \gamma \mathbf{P}_{AB}^\pi)^{-1} (\mathbf{r}_A \otimes \mathbf{e}) &= \sum_{t=0}^{\infty} \gamma^t \left(\sum_{j=1}^K x_j \mathbf{P}_A^{(j)} \otimes \mathbf{P}_B^{(j)} \right)^t (\mathbf{r}_A \otimes \mathbf{e}) \\ &\stackrel{(i)}{=} \sum_{t=0}^{\infty} \gamma^t \left(\left(\sum_{j=1}^K x_j \mathbf{P}_A^{(j)} \right)^t \mathbf{r}_A \right) \otimes \mathbf{e} \\ &= \left((\mathbf{I} - \gamma \mathbf{P}_A^\pi)^{-1} \mathbf{r}_A \right) \otimes \mathbf{e} = \mathbf{Q}_A^\pi \otimes \mathbf{e}, \end{aligned}$$

442 where we refer to (i) as an ‘‘absorbing’’ technique based on the bilinearity and mixed-product prop-
 443 erty of tensor product³. Specifically, since $\mathbf{P}\mathbf{e} = \mathbf{e}$ for any transition matrix \mathbf{P} , we have for any
 444 t ,

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

445 Similar results can be derived for \mathbf{P}_B^π such that $(\mathbf{I} - \gamma \mathbf{P}_{AB}^\pi)^{-1} (\mathbf{e} \otimes \mathbf{r}_B) = \mathbf{e} \otimes \mathbf{Q}_B^\pi$. Finally,
 446 combining the above results, we have

$$\mathbf{Q}_{AB}^\pi = (\mathbf{I} - \gamma \mathbf{P}_{AB}^\pi)^{-1} \mathbf{r}_{AB} = (\mathbf{I} - \gamma \mathbf{P}_{AB}^\pi)^{-1} (\mathbf{r}_A \otimes \mathbf{e} + \mathbf{e} \otimes \mathbf{r}_B) = \mathbf{Q}_A^\pi \otimes \mathbf{e} + \mathbf{e} \otimes \mathbf{Q}_B^\pi.$$

447 D An illustrative example of coupling and Markov entanglement

448 To elucidate the concept of Markov entanglement, we present an example of two-agent MDP where
 449 agents are coupled but not entangled. Consider a two-agent MDP \mathcal{M}_{AB} with $|\mathcal{A}_A| = |\mathcal{A}_B| = 2$,
 450 where action 1 means activate and 0 means idle. Each agent $i \in \{A, B\}$ has its own local transition
 451 kernel \mathbf{P}_i . We examine the following policy: at each time-step, we randomly activate one agent
 452 and keep another idle, i.e. $\pi(\mathbf{a} \mid \mathbf{s}) = 1/2$ if $\mathbf{a} = (0, 1)$ or $\mathbf{a} = (1, 0)$. Consequently, this
 453 policy couples the agents through the constraint $a_A + a_B = 1$ at each timestep. However, we
 454 will demonstrate that despite this coupling, there’s *no* entanglement. Specifically, we construct the
 455 following decomposition

$$\mathbf{P}_{AB}^\pi = \frac{1}{2} \mathbf{P}_A^0 \otimes \mathbf{P}_B^1 + \frac{1}{2} \mathbf{P}_A^1 \otimes \mathbf{P}_B^0, \quad (9)$$

456 where \mathbf{P}_i^a refers to the transition matrix of agents $i \in \{A, B\}$ taking action $a \in \{0, 1\}$. Intuitively,
 457 the right-hand side of Eq. (9) describes how at each time step, the global system randomly selects
 458 between two possible transitions: $\mathbf{P}_A^0 \otimes \mathbf{P}_B^1$ or $\mathbf{P}_A^1 \otimes \mathbf{P}_B^0$. This example thus clearly demonstrates
 459 a *coupled* system can still be *separable* and thus admits an exact value decomposition.

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

460 E Comparison with quantum entanglement

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

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

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

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

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

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

$$\mathcal{P}_{\text{SEP}} = \left\{ \mathbf{P} \geq 0 \mid \mathbf{P} = \sum_{j=1}^K x_j \mathbf{P}_A^{(j)} \otimes \mathbf{P}_B^{(j)}, \sum_{j=1}^K x_j = 1 \right\},$$

475 where $K \in \mathbb{Z}^+$ and $\{\mathbf{P}_A^{(j)}, \mathbf{P}_B^{(j)}\}_{j \in [K]}$ are transition matrices. $\mathbf{P} \geq 0$ calls for every element of
476 \mathcal{P}_{SEP} to be a valid transition matrix. It's clear that a transition matrix \mathbf{P}_{AB}^π is separable if and only if
477 $\mathbf{P}_{AB}^\pi \in \mathcal{P}_{\text{SEP}}$. On the other hand, a direct analogy of quantum entanglement gives us the following
478 set that further requires non-negative coefficients,

$$\mathcal{P}_{\text{SEP}}^+ = \left\{ \mathbf{P} \geq 0 \mid \mathbf{P} = \sum_{j=1}^K x_j \mathbf{P}_A^{(j)} \otimes \mathbf{P}_B^{(j)}, \sum_{j=1}^K x_j = 1, \mathbf{x} \geq 0 \right\}.$$

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

$$\mathbf{E}_{00} = \begin{pmatrix} 1 & 0 \\ 1 & 0 \end{pmatrix}, \quad \mathbf{E}_{01} = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}, \quad \mathbf{E}_{10} = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}, \quad \mathbf{E}_{11} = \begin{pmatrix} 0 & 1 \\ 0 & 1 \end{pmatrix}$$

482 And the corresponding transition matrix we provide is

$$\mathbf{P} = \begin{pmatrix} 0.5 & 0 & 0 & 0.5 \\ 0.5 & 0 & 0 & 0.5 \\ 0.5 & 0 & 0 & 0.5 \\ 0 & 0.5 & 0.5 & 0 \end{pmatrix} = \frac{1}{2} \mathbf{E}_{00} \otimes \mathbf{E}_{00} + \frac{1}{2} \mathbf{E}_{10} \otimes \mathbf{E}_{11} + \frac{1}{2} \mathbf{E}_{11} \otimes \mathbf{E}_{10} - \frac{1}{2} \mathbf{E}_{10} \otimes \mathbf{E}_{10}$$

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

486 **F Proof of Theorem 2**

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

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

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

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

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

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

 494 Let $\Omega_{P \otimes 2} := \text{span}(P_1 \otimes P_2)$ be the minimal subspace containing all separable transition matrices.
 495 It quickly follows that

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

 496 We then construct the orthogonal complement of $\Omega_{P \otimes 2}$ under Frobenius inner product. Let
 497 $\{\varepsilon_j\}_{j \in [m-1]}$ be a set of vector in \mathbb{R}^m such that $\varepsilon_j = (1, 0, \dots, 0, -1, 0, \dots, 0)^\top$ with the first
 498 element 1 and $j + 1$ -th element -1 . Notice that

$$\text{Tr}(e \varepsilon_j^\top P) = \text{Tr}(\varepsilon_j^\top P e) = 0 ,$$

 499 for all ε_j . Consider the following subspace

$$\Omega' = \left\{ \sum_{j=1}^{m-1} (\varepsilon_j e^\top) \otimes W_j^1 + \sum_{j=1}^{m-1} W_j^2 \otimes (\varepsilon_j e^\top) \mid W_{1:j}^1, W_{1:j}^2 \in \mathbb{R}^{m \times m} \right\} .$$

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

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

 501 and thus $\dim(\Omega') + \dim(\Omega_{P \otimes 2}) = m^4$. Moreover, one can verify for any $X \in \Omega_{P \otimes 2}$ and $Y \in \Omega'$,
 502 $\text{Tr}(X^\top Y) = 0$. As a result, it holds

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

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

 505 **Lemma 2.** *Given any transition matrix P and $\gamma > 0$, P is separable if and only if $(1 - \gamma)(I -$
 506 $\gamma P)^{-1}$ is separable.*

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

 509 (\Leftarrow) This side is more involved. Denote $U := (1 - \gamma)(I - \gamma P)^{-1}$. Then if the spectral radius
 510 $\rho(I - U) < 1$, then

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

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

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

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

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

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

515 **Step 3: Put it together** By Lemma 2, if \mathbf{P}_{AB}^π is entangled, then $(1 - \gamma)(\mathbf{I} - \gamma\mathbf{P}_{AB}^\pi)^{-1}$ is also
 516 entangled. Then there exists $\mathbf{Y} \in \Omega' \neq \mathbf{0}$ such that $\text{Tr}(\mathbf{Y}^\top (\mathbf{I} - \gamma\mathbf{P}_{AB}^\pi)^{-1}) \neq 0$. We apply singular
 517 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
 518 either $\text{Tr}((\mathbf{e}\varepsilon_j^\top) \otimes (\mathbf{v}\mathbf{u}^\top) (\mathbf{I} - \gamma\mathbf{P}_{AB}^\pi)^{-1}) \neq 0$ or $\text{Tr}((\mathbf{v}\mathbf{u}^\top) \otimes (\mathbf{e}\varepsilon_j^\top) (\mathbf{I} - \gamma\mathbf{P}_{AB}^\pi)^{-1}) \neq 0$. We
 519 assume the former without loss of generality, it holds

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

520 Now set $\mathbf{r}_A = \mathbf{0}$ and $\mathbf{r}_B = \mathbf{v}$. Since Q_{AB}^π is decomposable, there exists some local function
 521 Q_A, Q_B such that

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

522 Left multiply by $(\varepsilon_j^\top \otimes \mathbf{u}^\top)$, we have

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

523 Then set $\mathbf{r}_A = \mathbf{0}$ and $\mathbf{r}_B = -\mathbf{v}$, we can similarly derive

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

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

525 G Decomposition via general functions

526 Entangled \mathbf{P} precludes the local decomposition with local value functions, but may admit decom-
 527 positions with more general functions.

528 There exist other possible ways for value decomposition. For example, [Sunhag et al. \(2018\)](#);
 529 [Dou et al. \(2022\)](#) consider $Q_{AB}^\pi(\mathbf{s}, \mathbf{a}) = L_A(s_A, a_A, \mathbf{r}_{AB}) + L_B(s_B, a_B, \mathbf{r}_{AB})$ where L_A, L_B
 530 are learned jointly via minimizing the global Bellman error⁴; [Rashid et al. \(2020\)](#); [Mahajan et al.](#)

⁴In Appendix G, we provide an example of entangled MDP that allows for an exact value decomposition where L_A depends on both \mathbf{r}_A and \mathbf{r}_B .

531 (2019); Son et al. (2019); Wang et al. (2020) consider general monotonic operations beyond addi-
 532 tive decompositions. These methods introduce possibly richer representations at the cost of more
 533 sophisticated implementations and less interpretability, which is beyond the scope of this paper.

534 Consider $\mathbf{P} = \frac{1}{4}(ee^\top) \otimes (ee^\top) + \delta(\epsilon e^\top) \otimes (\epsilon e^\top)$, where $e = [1, 1], \epsilon = [1 - 1]$. Clearly such \mathbf{P} is
 535 entangled. We also have $\mathbf{P}^k = \frac{1}{4}(ee^\top) \otimes (ee^\top)$ for $k \geq 2$. Then $(\mathbf{I} - \gamma\mathbf{P})^{-1} = \mathbf{I} + \frac{\gamma + \gamma^2}{4}(ee^\top) \otimes$
 536 $(ee^\top) + \delta\gamma(\epsilon e^\top) \otimes (\epsilon e^\top)$. Then for any $\mathbf{r}_A, \mathbf{r}_B$, we have $(\mathbf{I} - \gamma\mathbf{P})^{-1}(\mathbf{r}_A \otimes e + e \otimes \mathbf{r}_B) =$
 537 $\mathbf{r}_A \otimes e + h_A(\gamma + \gamma^2)/2e \otimes e + \mathbf{r}_B \otimes e + h_B(\gamma + \gamma^2)/2e \otimes e + 2\delta\gamma(\epsilon^\top \mathbf{r}_B) \epsilon \otimes e$ where $h_A =$
 538 $e^\top \mathbf{r}_A, h_B = e^\top \mathbf{r}_B$.

539 H Proof of Theorem 3

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

$$\|\mathbf{A}\|_{\alpha, \beta} = \sup_{\|\mathbf{x}\|_\beta=1} \|\mathbf{A}\mathbf{x}\|_\alpha.$$

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

544 To prove the theorem, we introduce the key technique of analyzing perturbation bounds of the tran-
 545 sition matrix, which is also used in Farias et al. (2023).

546 **Lemma 3** (Lemma 1 in Farias et al. (2023)). *Let $\mathbf{P}, \mathbf{P}' \in \mathbb{R}^{n \times n}$ such that $(\mathbf{I} - \mathbf{P})^{-1}$ and $(\mathbf{I} - \mathbf{P}')^{-1}$
 547 exist. Then it holds*

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

548 We are then ready to prove the main theorem.

549 *Proof of Theorem 3.* Let $\mathbf{P}_A, \mathbf{P}_B$ be the optimal solution to Eq. (6) w.r.t agent A, B . For any subset
 550 of state-action pairs of agent A , $\mathcal{F} \subseteq \mathcal{S}_A \times \mathcal{A}_A$, we have

$$\begin{aligned} & \left| \sum_{s'_A, a'_A \in \mathcal{F}} (\mathbf{P}_A^\pi - \mathbf{P}_A)_{(s'_A, a'_A | s_A, a_A)} \right| \\ &= \left| \sum_{s'_A, a'_A \in \mathcal{F}} \sum_{s'_B, a'_B} \sum_{s_B, a_B} (\mathbf{P}_{AB}^\pi - \mathbf{P}_A \otimes \mathbf{P}_B)_{(s', a' | s, a)} \mu_{AB}^\pi(s_B, a_B | s_A, a_A) \right| \\ &\leq \sum_{s_B, a_B} \left| \sum_{s'_A, a'_A \in \mathcal{F}} \sum_{s'_B, a'_B} (\mathbf{P}_{AB}^\pi - \mathbf{P}_A \otimes \mathbf{P}_B)_{(s', a' | s, a)} \right| \mu_{AB}^\pi(s_B, a_B | s_A, a_A) \\ &\leq \sum_{s_B, a_B} E_A(\mathbf{P}_{AB}^\pi) \mu_{AB}^\pi(s_B, a_B | s_A, a_A) = E_A(\mathbf{P}_{AB}^\pi) \end{aligned}$$

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

$$\|\mathbf{P}_A^\pi - \mathbf{P}_A\|_{\text{TV}} \leq E_A(\mathbf{P}_{AB}^\pi),$$

553 and similar results hold for \mathbf{P}_B^π .

554 Next we have

$$\begin{aligned}
& (\mathbf{I} - \gamma \mathbf{P}_{AB}^\pi)^{-1} (\mathbf{r}_A \otimes \mathbf{e}) - \left((\mathbf{I} - \gamma \mathbf{P}_A^\pi)^{-1} \mathbf{r}_A \right) \otimes \mathbf{e} \\
&= (\mathbf{I} - \gamma \mathbf{P}_{AB}^\pi)^{-1} (\mathbf{r}_A \otimes \mathbf{e}) - (\mathbf{I} - \gamma \mathbf{P}_A \otimes \mathbf{P}_B)^{-1} (\mathbf{r}_A \otimes \mathbf{e}) \\
&\quad + (\mathbf{I} - \gamma \mathbf{P}_A \otimes \mathbf{P}_B)^{-1} (\mathbf{r}_A \otimes \mathbf{e}) - \left((\mathbf{I} - \gamma \mathbf{P}_A^\pi)^{-1} \mathbf{r}_A \right) \otimes \mathbf{e} \\
&\stackrel{(i)}{=} \underbrace{(\mathbf{I} - \gamma \mathbf{P}_{AB}^\pi)^{-1} (\mathbf{r}_A \otimes \mathbf{e}) - (\mathbf{I} - \gamma \mathbf{P}_A \otimes \mathbf{P}_B)^{-1} (\mathbf{r}_A \otimes \mathbf{e})}_{(I)} \\
&\quad + \underbrace{\left((\mathbf{I} - \gamma \mathbf{P}_A)^{-1} \mathbf{r}_A \right) \otimes \mathbf{e} - \left((\mathbf{I} - \gamma \mathbf{P}_A^\pi)^{-1} \mathbf{r}_A \right) \otimes \mathbf{e}}_{(II)}
\end{aligned}$$

555 where (i) also follows the same ‘‘absorbing’’ technique in the proof of Theorem 1.

556 For (I), apply Lemma 3, it holds

$$\begin{aligned}
& \left\| (\mathbf{I} - \gamma \mathbf{P}_{AB}^\pi)^{-1} (\mathbf{r}_A \otimes \mathbf{e}) - (\mathbf{I} - \gamma \mathbf{P}_A \otimes \mathbf{P}_B)^{-1} (\mathbf{r}_A \otimes \mathbf{e}) \right\|_\infty \\
&= \left\| (\mathbf{I} - \gamma \mathbf{P}_{AB}^\pi)^{-1} (\gamma \mathbf{P}_{AB}^\pi - \gamma \mathbf{P}_A \otimes \mathbf{P}_B) (\mathbf{I} - \gamma \mathbf{P}_A \otimes \mathbf{P}_B)^{-1} (\mathbf{r}_A \otimes \mathbf{e}) \right\|_\infty \\
&\leq \left\| (\mathbf{I} - \gamma \mathbf{P}_{AB}^\pi)^{-1} \right\|_\infty \left\| (\gamma \mathbf{P}_{AB}^\pi - \gamma \mathbf{P}_A \otimes \mathbf{P}_B) \left((\mathbf{I} - \gamma \mathbf{P}_A)^{-1} \mathbf{r}_A \right) \otimes \mathbf{e} \right\|_\infty \\
&\stackrel{(i)}{\leq} \left\| (\mathbf{I} - \gamma \mathbf{P}_{AB}^\pi)^{-1} \right\|_\infty 2\gamma E_A(\mathbf{P}_{AB}^\pi) \left\| (\mathbf{I} - \gamma \mathbf{P}_A)^{-1} \mathbf{r}_A \right\|_\infty \\
&\leq \frac{2\gamma E_A(\mathbf{P}_{AB}^\pi) r_{\max}^A}{1 - \gamma} \left\| (\mathbf{I} - \gamma \mathbf{P}_{AB}^\pi)^{-1} \right\|_\infty \leq \frac{2\gamma E_A(\mathbf{P}_{AB}^\pi) r_{\max}^A}{(1 - \gamma)^2},
\end{aligned}$$

557 where (i) follows by the definition of agent-wise total variation distance when $\|\mathbf{r}_A\|_\infty \neq 0$, and also

558 trivially hold when $\|\mathbf{r}_A\|_\infty = 0$. Similarly, for (II) we have

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

559 Then we have

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

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

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

561 Put it all together we have

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

562

□

563 **I Proof of Theorem 4**

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

 565 **Definition 5** (μ -norm). *Given a transition matrix $\mathbf{P} \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}| \times |\mathcal{S}||\mathcal{A}|}$ with occupancy measure⁵*
 566 *$\mu \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}$, for any vector $\mathbf{x} \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}$ the μ -norm is defined as*

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

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

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

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

 572 The μ -weighted ATV distance w.r.t agent B can be defined similarly. We claim that the μ -weighted
 573 ATV is also a counterpart of ATV distance in Definition 5. This follows from the constrained
 574 optimization formulation of ATV

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

 575 Thus μ -ATV substitutes μ -norm for the original ℓ_∞ -norm. We plug μ -weighted ATV into Eq. (1)
 576 and obtain the corresponding measure of Markov entanglement $E(\mathbf{P}_{AB}^\pi)$ and $E_A(\mathbf{P}_{AB}^\pi)$. Similar to
 577 ATV in Eq. (6), this μ -weighted version of $E_A(\mathbf{P}_{AB}^\pi)$ admits the following formulation

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

 578 This recovers Eq. (7) that substitutes the μ -weighted average for the maximum operator in Eq. (6).
 579 Thus intuitively, $E(\mathbf{P}_{AB}^\pi)$ w.r.t μ -weighted ATV distance measures *how closely agent A can be*
 580 *approximated as an independent subsystem under the stationary distribution.*

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

584 The first one is the following lemma for the “localized” stationary distribution

 585 **Lemma 4.** \mathbf{P}_A^π has stationary distribution μ_A^π with

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

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

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

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

$$\begin{aligned}
 & \sum_{s_A, a_A} \left(\sum_{s_B, a_B} \mu_{AB}^\pi(s_A, s_B, a_A, a_B) \right) P^\pi(s'_A, a'_A \mid s_A, a_A) \\
 = & \sum_{s_A, a_A} \sum_{s_B, a_B} \mu_{AB}^\pi(s_A, s_B, a_A, a_B) \sum_{s'_B, a'_B} \sum_{s''_B, a''_B} P^\pi(s'_A, s'_B, a'_A, a'_B \mid s_A, s''_B, a_A, a''_B) \mu_{AB}^\pi(s''_B, a''_B \mid s_A, a_A) \\
 = & \sum_{s_A, a_A} \sum_{s_B, a_B} \mu_{AB}^\pi(s_B, a_B \mid s_A, a_A) \sum_{s'_B, a'_B} \sum_{s''_B, a''_B} P^\pi(s'_A, s'_B, a'_A, a'_B \mid s_A, s''_B, a_A, a''_B) \mu_{AB}^\pi(s_A, s''_B, a_A, a''_B) \\
 = & \sum_{s_A, a_A} \sum_{s'_B, a'_B} \sum_{s''_B, a''_B} P^\pi(s'_A, s'_B, a'_A, a'_B \mid s_A, s''_B, a_A, a''_B) \mu_{AB}^\pi(s_A, s''_B, a_A, a''_B) \\
 = & \sum_{s'_B, a'_B} \mu_{AB}^\pi(s'_A, s'_B, a'_A, a'_B).
 \end{aligned}$$

590 where the last equation follows from the definition of μ_{AB}^π . Hence we conclude that
 591 $\sum_{s_B, a_B} \mu_{AB}^\pi(s_A, s_B, a_A, a_B)$ is a stationary distribution of \mathbf{P}_A^π . \square

592 We are then ready to prove Theorem 4. We first note that similar to ATV distance in Eq. (6), the
 593 optimal solution to $E_A(\mathbf{P}_{AB}^\pi)$ w.r.t μ_{AB}^π -weighted ATV distance also only depends on \mathbf{P}_A . Thus,
 594 let $\mathbf{P}_A, \mathbf{P}_B$ be the optimal solutions to $E_A(\mathbf{P}_{AB}^\pi), E_B(\mathbf{P}_{AB}^\pi)$ respectively.

595 Let $\mathbf{x} \in \mathbb{R}^{|S_A| |A_A|}$ with $\|\mathbf{x}\|_\infty = 1$. Following the same technique in the proof of Theorem 4, we
 596 have

$$\begin{aligned}
 & \mu_A^{\pi^\top} |(\mathbf{P}_A^\pi - \mathbf{P}_A) \mathbf{x}| \\
 = & \sum_{s_A, a_A} \mu_A^\pi(s_A, a_A) \left| \sum_{s'_A, a'_A} (\mathbf{P}_A^\pi - \mathbf{P}_A)_{(s'_A, a'_A | s_A, a_A)} \mathbf{x}(s'_A, a'_A) \right| \\
 = & \sum_{s_A, a_A} \mu_A^\pi(s_A, a_A) \left| \sum_{s'_A, a'_A} \mathbf{x}(s'_A, a'_A) \sum_{s'_B, a'_B} \sum_{s_B, a_B} (\mathbf{P}_{AB}^\pi - \mathbf{P}_A \otimes \mathbf{P}_B)_{(s', a' | s, a)} \mu_{AB}^\pi(s_B, a_B \mid s_A, a_A) \right| \\
 \leq & \sum_{s, a} \left| \sum_{s'_A, a'_A} \mathbf{x}(s'_A, a'_A) \sum_{s'_B, a'_B} (\mathbf{P}_{AB}^\pi - \mathbf{P}_A \otimes \mathbf{P}_B)_{(s', a' | s, a)} \right| \mu_{AB}^\pi(s, a) \leq 2E_A(\mathbf{P}_{AB}^\pi)
 \end{aligned}$$

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

$$\|\mathbf{P}_A^\pi - \mathbf{P}_A\|_{\mu, \infty} \leq 2E_A(\mathbf{P}_{AB}^\pi),$$

598 and similar results hold for \mathbf{P}_B^π . We then apply the decomposition

$$\begin{aligned}
 & (\mathbf{I} - \gamma \mathbf{P}_{AB}^\pi)^{-1} (\mathbf{r}_A \otimes \mathbf{e}) - \left((\mathbf{I} - \gamma \mathbf{P}_A^\pi)^{-1} \mathbf{r}_A \right) \otimes \mathbf{e} \\
 = & \underbrace{(\mathbf{I} - \gamma \mathbf{P}_{AB}^\pi)^{-1} (\mathbf{r}_A \otimes \mathbf{e}) - (\mathbf{I} - \gamma \mathbf{P}_A \otimes \mathbf{P}_B)^{-1} (\mathbf{r}_A \otimes \mathbf{e})}_{(I)} \\
 & + \underbrace{\left((\mathbf{I} - \gamma \mathbf{P}_A)^{-1} \mathbf{r}_A \right) \otimes \mathbf{e} - \left((\mathbf{I} - \gamma \mathbf{P}_A^\pi)^{-1} \mathbf{r}_A \right) \otimes \mathbf{e}}_{(II)}
 \end{aligned}$$

599 For (I), we have

$$\begin{aligned}
 & \left\| (\mathbf{I} - \gamma \mathbf{P}_{AB}^\pi)^{-1} (\mathbf{r}_A \otimes \mathbf{e}) - (\mathbf{I} - \gamma \mathbf{P}_A \otimes \mathbf{P}_B)^{-1} (\mathbf{r}_A \otimes \mathbf{e}) \right\|_{\mu_{AB}^\pi} \\
 &= \left\| (\mathbf{I} - \gamma \mathbf{P}_{AB}^\pi)^{-1} (\gamma \mathbf{P}_{AB}^\pi - \gamma \mathbf{P}_A \otimes \mathbf{P}_B) (\mathbf{I} - \gamma \mathbf{P}_A \otimes \mathbf{P}_B)^{-1} (\mathbf{r}_A \otimes \mathbf{e}) \right\|_{\mu_{AB}^\pi} \\
 &\stackrel{(i)}{\leq} \frac{1}{1-\gamma} \left\| \left((\gamma \mathbf{P}_{AB}^\pi - \gamma \mathbf{P}_A \otimes \mathbf{P}_B) (\mathbf{I} - \gamma \mathbf{P}_A)^{-1} \mathbf{r}_A \right) \otimes \mathbf{e} \right\|_{\mu_{AB}^\pi} \\
 &\leq \frac{2\gamma E(\pi)}{1-\gamma} \left\| (\mathbf{I} - \gamma \mathbf{P}_A)^{-1} \mathbf{r}_A \right\|_\infty \leq \frac{2\gamma E(\pi) r_{\max}}{(1-\gamma)^2},
 \end{aligned}$$

600 where (i) follows from the fact that for any \mathbf{x}

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

601 For (II) one can use Lemma 4 to verify

$$\begin{aligned}
 & \left\| \left((\mathbf{I} - \gamma \mathbf{P}_A)^{-1} \mathbf{r}_A \right) \otimes \mathbf{e} - \left((\mathbf{I} - \gamma \mathbf{P}_A^\pi)^{-1} \mathbf{r}_A \right) \otimes \mathbf{e} \right\|_{\mu_{AB}^\pi} \\
 &= \left\| (\mathbf{I} - \gamma \mathbf{P}_A)^{-1} \mathbf{r}_A - (\mathbf{I} - \gamma \mathbf{P}_A^\pi)^{-1} \mathbf{r}_A \right\|_{\mu_A^\pi}
 \end{aligned}$$

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

603 J Results for multi-agent MDPs

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

607 Consider an N -agent MDP $\mathcal{M}_{1:N}(\mathcal{S}, \mathcal{A}, \mathbf{P}, \mathbf{r}_{1:N}, \gamma)$ with joint state space $\mathcal{S} = \times_{i=1}^N \mathcal{S}_i$ and joint
 608 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 .
 609 For agents at global state $\mathbf{s} = (s_1, s_2, \dots, s_N)$ with action $\mathbf{a} = (a_1, a_2, \dots, a_N)$ taken, the system
 610 will transit to $\mathbf{s}' = (s'_1, s'_2, \dots, s'_N)$ according to transition kernel $\mathbf{s}' \sim \mathbf{P}(\cdot | \mathbf{s}, \mathbf{a})$ and each agent
 611 $i \in [N]$ will receive its local reward $r_i(s_i, a_i)$. The global reward $r_{1:N}$ is defined as the summation
 612 of local rewards $r_{1:N}(\mathbf{s}, \mathbf{a}) := \sum_{i=1}^N r_i(s_i, a_i)$, or in vector form,

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

613 We further assume the local rewards are bounded, i.e. for agent $i \in [N]$, $|r_i(s_i, a_i)| \leq$
 614 r_{\max}^i for all (s_i, a_i) . Given any global policy $\pi: \mathcal{S} \rightarrow \Delta(\mathcal{A})$, we denote $\mathbf{P}_{1:N}^\pi \in$
 615 $\mathbb{R}^{|\mathcal{S}|^N |\mathcal{A}|^N \times |\mathcal{S}|^N |\mathcal{A}|^N}$ as the transition matrix induced by π where $P_{1:N}^\pi(s'_{1:N}, a'_{1:N} | s_{1:N}, a_{1:N}) :=$
 616 $\mathbf{P}(s'_{1:n} | s_{1:n}, a_{1:n}) \pi(a'_{1:n} | s'_{1:n})$. Then the global Q-value is defined by Bellman Equation
 617 $Q_{1:N}^\pi = (\mathbf{I} - \gamma \mathbf{P}_{1:N}^\pi)^{-1} \mathbf{r}_{1:N}$. The local Q-values follow the similar framework to Meta Algorithm 1
 618 where each agent $i \in [N]$ fits Q_i^π using its local observations. We then sum up local Q-values to
 619 approximate the global Q-value, i.e.

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

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

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

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

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

626 Analogously, we define the Multi-agent Markov Entanglement,

627 **Definition 8** (Multi-agent Markov Entanglement). Consider a N -agent Markov system $\mathcal{M}_{1:N}$ and
 628 policy $\pi: \mathcal{S} \rightarrow \Delta(\mathcal{A})$, the agents are **separable** under policy π if there exists $K \in \mathbb{Z}^+$, measure
 629 $\{x_j\}_{j \in [K]}$ satisfying $\sum_{j=1}^K x_j = 1$, and transition matrices $\{\mathbf{P}_{1:N}^{(j)}\}_{j \in [K]}$ such that

$$\mathbf{P}_{1:N}^\pi = \sum_{j=1}^K x_j \mathbf{P}_1^{(j)} \otimes \mathbf{P}_2^{(j)} \otimes \cdots \otimes \mathbf{P}_N^{(j)}.$$

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

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

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

$$P_i^\pi(s'_i, a'_i | s_i, a_i) = \sum_{s'_{-i}, a'_{-i}} \sum_{s_{-i}, a_{-i}} P_{1:N}^\pi(s'_{1:N}, a'_{1:N} | s_{1:N}, a_{1:N}) \rho_{1:N}^\pi(s_{-i}, a_{-i} | s_i, a_i).$$

635 **Definition 9** (Measure of Multi-agent Markov Entanglement). Consider a N -agent Markov system
 636 $\mathcal{M}_{1:N}$ 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
 637 $\pi: \mathcal{S} \rightarrow \Delta(\mathcal{A})$, the measure of Markov entanglement of N agents is

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

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

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

641 **Theorem 6.** 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
 642 $\mathcal{A} = \times_{i=1}^N \mathcal{A}_i$. Given any policy $\pi: \mathcal{S} \rightarrow \Delta(\mathcal{A})$ with the measure of Markov entanglement $E_i(\pi)$
 643 w.r.t ATV distance, it holds for any agent i ,

$$\|\mathbf{P}_i^\pi - \mathbf{P}_i\|_\infty \leq 2_i E(\pi).$$

644 where \mathbf{P}_i is the optimal solution of Eq. (13). Furthermore, the decomposition error is entry-wise
 645 bounded by the measure of Markov entanglement,

$$\left\| Q_{1:N}^\pi(\mathbf{s}, \mathbf{a}) - \sum_{i=1}^N Q_i^\pi(s_i, a_i) \right\|_\infty \leq \frac{4\gamma \left(\sum_{i=1}^N E_i(\mathbf{P}_{1:N}^\pi) \tau_{\max}^i \right)}{(1-\gamma)^2}.$$

646 The proof mainly follows the following lemma, which generalizes the key technique used in Theo-
 647 rem 1.

648 **Lemma 5.** For any agent i , it holds

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

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

651 **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
652 $\mathcal{A} = \times_{i=1}^N \mathcal{A}_i$. Given any policy $\pi: \mathcal{S} \rightarrow \Delta(\mathcal{A})$ with the measure of Markov entanglement $E_i(\mathbf{P}_{1:N}^\pi)$
653 w.r.t the $\mu_{1:N}^\pi$ -weighted agent-wise total variation distance, it holds for any agent i ,

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

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

$$\left\| Q_{1:N}^\pi(\mathbf{s}, \mathbf{a}) - \sum_{i=1}^N Q_i^\pi(s_i, a_i) \right\|_{\mu_{1:N}^\pi} \leq \frac{4\gamma \left(\sum_{i=1}^N E_i(\mathbf{P}_{1:N}^\pi) r_{\max}^i \right)}{(1-\gamma)^2}.$$

657 K Proof of Theorem 5

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

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

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

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

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

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

669 We thus focus on the deviation from \mathbf{m} to \mathbf{m}^* . We extend the concentration analysis from [Gast
670 et al. \(2023; 2024\)](#) to derive a new stability bound for the RHS. Specifically, we finishing the proof
671 via demonstrating the deviation decays at the rate $\mathcal{O}(1/\sqrt{N})$.

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

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

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

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

676 **Assumption A** (Uniform Global Attractor Property (UGAP)). There exists a uniform global attrac-
677 tor \mathbf{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 $\mathbf{m} \in \Delta^{|\mathcal{S}|}$,
678 one has $\|\Phi_t(\mathbf{m}) - \mathbf{m}^*\|_\infty < \varepsilon$.

679 The UGAP assumption ensures the uniqueness of \mathbf{m}^* and guarantees fast convergence from any
680 initial \mathbf{m} to \mathbf{m}^* .

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

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

684 Non-degenerate and UGAP are two standard technical assumptions for the index policy, which re-
685 strict chaotic behavior in asymptotic regime and will be further introduced in subsequent sections.
686 We note here these two assumptions are also used in almost all theoretical work on index policies
687 [Weber & Weiss \(1990\)](#); [Verloop \(2016\)](#); [Gast et al. \(2023; 2024\)](#).

688 *Proof of Theorem 5.* In the subsequent proof, we let $\nu_1 > \nu_2 > \nu_3 > \dots > \nu_{|\mathcal{S}|}$. This does not lose
689 generality in that we can always exchange state index. The proof consists of several steps

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

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

691 Notice that the definition of ϕ^π does not depend on N . We adapt from Lemma B.1 in [Gast et al.](#)
692 [\(2023\)](#) defined specially for Whittle Index,

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

695 When the context is clear, we abbreviate ϕ^π as ϕ . For any $\mathbf{m} \in \Delta^{|\mathcal{S}|}$, define $s(\mathbf{m}) \in [|\mathcal{S}|]$ be
696 the state such that $\sum_{i=1}^{s(\mathbf{m})-1} m_i \leq \alpha < \sum_{i=1}^{s(\mathbf{m})} m_i$. Lemma 7 characterizes for any $\mathbf{m} \in \mathcal{Z}_i :=$
697 $\{\mathbf{m} \in \Delta^{|\mathcal{S}|} \mid s(\mathbf{m}) = i\}$, there exists $\mathbf{K}_{s(\mathbf{m})}, \mathbf{b}_{s(\mathbf{m})}$ such that

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

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

701 **Step 2: Connecting policy entanglement with the deviation of stationary distribution** Com-
702 bine Proposition 8 with the RMAB model, we have

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

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

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

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

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

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

$$|\pi(a_i = a \mid s_i = s, \mathbf{m}) - \pi^*(a \mid s)| = |\pi(a_j = a \mid s_j = s, \mathbf{m}) - \pi^*(a \mid s)|.$$

708 It then suffices to prove $\sum_{s_i=s, s_{1:N}=\mathbf{m}} \mu(s_{1:N}) = \sum_{s_j=s, s_{1:N}=\mathbf{m}} \mu(s_{1:N})$. If
709 $\sum_{s_i=s, s_{1:N}=\mathbf{m}} \mu(s_{1:N}) \leq \sum_{s_j=s, s_{1:N}=\mathbf{m}} \mu(s_{1:N})$, we can exchange the agent index of i and j . This

710 will result in the same stationary distribution and $\sum_{s_i=s, s_{1:N}=\mathbf{m}} \mu(s_{1:N}) \geq \sum_{s_j=s, s_{1:N}=\mathbf{m}} \mu(s_{1:N})$
 711 and thus the equation. We then rewrite the bound in Proposition 8,

$$\begin{aligned} E(\pi) &\leq \frac{1}{2} \sup_i \sum_{s_{1:N}} \mu^\pi(s_{1:N}) \sum_{a_i} |\pi(a_i | s_{1:N}) - \pi^*(a_i | s_i)| \\ &= \sup_i \sum_{s_{1:N}} \mu^\pi(s_{1:N}) |\pi(a_i = 1 | s_{1:N}) - \pi^*(a_i = 1 | s_i)| \\ &= \frac{1}{N} \sum_{s_{1:N}} \mu^\pi(s_{1:N}) \sum_{i=1}^N |\pi(a_i = 1 | s_{1:N}) - \pi^*(a_i = 1 | s_i)| \\ &= \sum_{\mathbf{m}} \mu^\pi(\mathbf{m}) \sum_{s \in \mathcal{S}} \mathbf{m}_s |\pi(a = 1 | s, \mathbf{m}) - \pi^*(a = 1 | s)| \end{aligned}$$

712 For any configuration \mathbf{m} and state s , we have

$$\begin{aligned} &\mathbf{m}_s |\pi(a = 1 | s, \mathbf{m}) - \pi^*(a = 1 | s)| \\ &= \mathbf{m}_s \left| \frac{\pi^*(a = 1 | s) \mathbf{m}_s^* N + k_s}{\mathbf{m}_s^* N + \ell_s} - \pi^*(a = 1 | s) \right| \\ &= \frac{\mathbf{m}_s^* N + \ell_s}{N} \left| \frac{k_s - \ell_s \pi^*(a = 1 | s)}{\mathbf{m}_s^* N + \ell_s} \right| \\ &\leq |\mathcal{S}| \|\mathbf{m} - \mathbf{m}^*\|_\infty, \end{aligned}$$

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

716 □

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

719 **Lemma 9** (One-step Concentration, Lemma 1 in [Gast et al. \(2024\)](#)). *Let $\epsilon[1] = \mathbf{m}[1] - \phi(\mathbf{m}[0])$, it*
 720 *holds*

$$\mathbb{E}[\|\epsilon[1]\|_1 | \mathbf{m}[0]] \leq \sqrt{\frac{|\mathcal{S}|}{N}}.$$

721 **Lemma 10** (Multi-step Concentration, Lemma C.4 in [Gast et al. \(2023\)](#)). *There exists a positive*
 722 *constant K such that for all $t \in \mathbb{N}$ and $\delta > 0$,*

$$\Pr[\|\mathbf{m}[t] - \Phi_t(\mathbf{m})\|_\infty \geq (1 + K + K^2 + \dots + K^t)\delta | \mathbf{m}[0] = \mathbf{m}] \leq t|\mathcal{S}|e^{-2N\delta^2}$$

723 **Lemma 11** (Local Stability, Lemma C.5 in [Gast et al. \(2023\)](#)). *Under non-degenerate and UGAP:*

724 (i) $\mathbf{K}_{s(\mathbf{m}^*)}$ is a stable matrix, i.e. its spectral radius is strictly less than 1.

725 (ii) For any ϵ , there exists $T(\epsilon) > 0$ such that for all $\mathbf{m} \in \Delta^{|\mathcal{S}|}$, $\|\Phi_{T(\epsilon)}(\mathbf{m}) - \mathbf{m}^*\|_\infty < \epsilon$.

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

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

728 Combine the second result of Lemma 11 and non-degenerate condition, we can construct a neigh-
 729 borhood \mathcal{N} of \mathbf{m}^* such that $\mathcal{N} = \mathcal{B}(\mathbf{m}^*, \epsilon) \cap \Delta^{|\mathcal{S}|} \in \mathcal{Z}_{s(\mathbf{m}^*)}$ where $\epsilon > 0$ and $\mathcal{B}(\mathbf{m}^*, \epsilon) =$
 730 $\{\mathbf{m} | \|\mathbf{m} - \mathbf{m}^*\|_\infty < \epsilon\}$ is an open ball. We next show that $\mathbf{m}[0]$ under stationary distribution

731 will concentrate in \mathcal{N} with high probability. Let $\tilde{T} = T(\epsilon/2)$ such that for all $\mathbf{m} \in \Delta^{|\mathcal{S}|}$,
 732 $\|\Phi_{\tilde{T}}(\mathbf{m}) - \mathbf{m}^*\|_\infty < \epsilon/2$. It holds

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

733 where (i) follows from the stationarity $\mathbf{m}[\tilde{T}]$ and $\mathbf{m}[0]$ are *i.i.d* and the constant $u =$
 734 $\left(\frac{\epsilon}{2(1+K+K^2+\dots+K^{\tilde{T}})}\right)^2$ does not depend on N .

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

$$\begin{aligned} \mathbf{m}[1] - \mathbf{m}^* &= \phi(\mathbf{m}[0]) + \epsilon[1] - \mathbf{m}^* \\ &= \mathbf{K}_{s(\mathbf{m}^*)}(\mathbf{m}[0] - \mathbf{m}^*) + \epsilon[1]. \end{aligned}$$

737 Taking $\|\cdot\|_\beta$ on both side,

$$\begin{aligned} \|\mathbf{m}[1] - \mathbf{m}^*\|_\beta &\leq \|\mathbf{K}_{s(\mathbf{m}^*)}(\mathbf{m}[0] - \mathbf{m}^*)\|_\beta + \|\epsilon[1]\|_\beta \\ &\leq \|\mathbf{K}_{s(\mathbf{m}^*)}\|_\beta \|\mathbf{m}[0] - \mathbf{m}^*\|_\beta + \|\epsilon[1]\|_\beta. \end{aligned}$$

738 Taking expectation on both side,

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

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

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

741 We combine Lemma 8 and conclude the proof of Theorem 5.

742 L Extensions of Markov entanglement

743 L.1 (Weakly-)coupled MDPs

744 Weakly-coupled MDPs (WCMDP) are a rich class of multi-agent model that capture many real-
 745 world applications such as supply chain management, queuing network and resource allocations

746 Adelman & Mersereau (2008); Brown & Zhang (2023); Shar & Jiang (2023). Compared to general
 747 multi-agent MDP, WCMDP further ensures each agent follow its local transition while the agents’
 748 actions are coupled with each other. Formally,

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

- 751 • Each agent has local transition kernel \mathbf{P}_i such that $\forall \mathbf{s}, \mathbf{a}, \mathbf{s}', P(\mathbf{s}' | \mathbf{s}, \mathbf{a}) = \prod_{i=1}^N P_i(s'_i | s_i, a_i)$.
- 752 • At global state \mathbf{s} , agents’ joint actions \mathbf{a} are subject to m coupling constraints $\sum_{i=1}^N \mathbf{d}(s_i, a_i) \leq$
 753 $\mathbf{b} \in \mathbb{R}^m$.

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

756 **Proposition 8.** *Consider a N -agent weakly-coupled MDP $\mathcal{M}_{1:N}(\mathcal{S}, \mathcal{A}, \mathbf{P}, \mathbf{r}_{1:N}, \gamma)$. Given any*
 757 *policy $\pi : \mathcal{S} \rightarrow \Delta(\mathcal{A})$ with measure of Markov entanglement $E_i(\mathbf{P}_{1:N}^\pi)$ w.r.t the $\mu_{1:N}^\pi$ -weighted*
 758 *agent-wise total variation distance, it holds for $i \in [N]$,*

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

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

760 *Proof of Proposition 8.* We demonstrate the proof for two-agent WCMDP and the generalization
 761 to multi-agent WCMDP is straightforward. Consider \mathbf{P}_A^π be the transition of agent A under local
 762 policy π' . We focus on agent A

$$\begin{aligned} & E_A(\mathbf{P}_{AB}^\pi) \\ & \leq \frac{1}{2} \sum_{\mathbf{s}, \mathbf{a}} \mu_{AB}^\pi(\mathbf{s}, \mathbf{a}) \sum_{s'_A, a'_A} \left| P_{AB}^\pi(s'_A, a'_A | \mathbf{s}, \mathbf{a}) - P_A^{\pi'}(s'_A, a'_A | s_A, a_A) \right| \\ & = \frac{1}{2} \sum_{\mathbf{s}, \mathbf{a}} \mu_{AB}^\pi(\mathbf{s}, \mathbf{a}) \sum_{s'_A, a'_A} \left| \sum_{s'_B} P_{AB}^\pi(s', a_A | \mathbf{s}, \mathbf{a}) - P_A^{\pi'}(s'_A | s_A, a_A) \pi'(a'_A | s'_A) \right| \\ & \stackrel{(i)}{=} \frac{1}{2} \sum_{\mathbf{s}, \mathbf{a}} \mu_{AB}^\pi(\mathbf{s}, \mathbf{a}) \sum_{s'_A, a'_A} \left| \sum_{s'_B} P_{AB}^\pi(s', a_A | \mathbf{s}, \mathbf{a}) - \sum_{s'_B} P(s' | \mathbf{s}, \mathbf{a}) \pi'(a'_A | s'_A) \right| \\ & = \frac{1}{2} \sum_{\mathbf{s}, \mathbf{a}} \mu_{AB}^\pi(\mathbf{s}, \mathbf{a}) \sum_{s'_A, a'_A} \left| \sum_{s'_B} P(s' | \mathbf{s}, \mathbf{a}) (\pi(a'_A | s') - \pi'(a'_A | s'_A)) \right| \\ & \leq \frac{1}{2} \sum_{\mathbf{s}, \mathbf{a}} \mu_{AB}^\pi(\mathbf{s}, \mathbf{a}) \sum_{\mathbf{s}'} P(\mathbf{s}' | \mathbf{s}, \mathbf{a}) \sum_{a'_A} |\pi(a'_A | \mathbf{s}') - \pi'(a'_A | s'_A)| \\ & \stackrel{(ii)}{=} \frac{1}{2} \sum_{\mathbf{s}'} \mu_{AB}^\pi(\mathbf{s}') \sum_{a'_A} |\pi(a'_A | \mathbf{s}') - \pi'(a'_A | s'_A)|. \end{aligned}$$

763 where (i) follows from the transition structure of weakly coupled MDP $P(\mathbf{s}' | \mathbf{s}, \mathbf{a}) = P(s'_A |$
 764 $s_A, a_A) \cdot P(s'_B | s_B, a_B)$; and (ii) comes from the fact that $P^\pi(\mathbf{s}' | \mathbf{s}) = \sum_{\mathbf{a}} \pi(\mathbf{a} | \mathbf{s}) P(\mathbf{s}' | \mathbf{s}, \mathbf{a})$
 765 and $\sum_{\mathbf{s}} \mu^\pi(\mathbf{s}) P^\pi(\mathbf{s}' | \mathbf{s}) = \mu^\pi(\mathbf{s}')$. \square

766 Proposition 8 establishes an upper bound for Markov entanglement in WCMDP. Intuitively, this
 767 bound characterizes how agent i can be viewed as making independent decisions. It takes advantage
 768 of the weakly-coupled structure and shaves off the transition in Markov entanglement measure.

769 **L.2 Coupled MDPs with exogenous information**

770 In many practical scenarios, the agents' transitions and actions are coupled by a shared exogenous
 771 signal. For example, in ride-hailing platforms, the specific dispatch is related to the exogenous order
 772 at the current moment [Qin et al. \(2020\)](#); [Han et al. \(2022\)](#); [Azagirre et al. \(2024\)](#); in warehouse
 773 routing, the scheduling of robots is also related to the exogenous task revealed so far [Chan et al.](#)
 774 [\(2024\)](#).

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

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

$$P_{ABz}^\pi(\mathbf{s}', \mathbf{a}', z' | \mathbf{s}, \mathbf{a}, z) = P(\mathbf{s}' | \mathbf{s}, \mathbf{a}, z) \cdot \pi(\mathbf{a}' | \mathbf{s}', z') \cdot P(z' | z).$$

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

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

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

$$P_{Az}(s'_A, a'_A, z' | s_A, a_A, z) = \sum_{s_B, a_B} \mu_{ABz}^\pi(s_B, a_B | s_A, a_A, z) \sum_{s'_B, a'_B} P_{ABz}^\pi(s', \mathbf{a}', z' | \mathbf{s}, \mathbf{a}, z),$$

785 Given the local transition, we have the local value $\mathbf{Q}_{Az}^\pi = (\mathbf{I} - \gamma \mathbf{P}_{Az})^{-1}(\mathbf{r}_{Az})$ via Bellman Equa-
 786 tion.

787 Combined with exogenous information, we consider the following value decomposition

$$Q_{AB}^\pi(\mathbf{s}, \mathbf{a}, z) = Q_A^\pi(s_A, a_A, z) + Q_B^\pi(s_B, a_B, z).$$

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

$$\mathbf{P}_{ABz}^\pi = \sum_{j=1}^K x_j \mathbf{P}_{Az}^{(j)} \otimes \mathbf{P}_B^{(j)}. \quad (15)$$

789 **Proposition 9.** *If the system is agent-wise separable for all agents, then*

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

790 *Proof.* The proof is basically the same as Theorem 1. One can first quickly show that $P_{Az} =$
 791 $\sum_{j=1}^K x_j P_{Az}^{(j)}$. And then it holds

$$\begin{aligned}
 & \left(\sum_{j=1}^K x_j P_{Az}^{(j)} \otimes P_B^{(j)} \right)^t (\mathbf{r}_A \otimes \mathbf{e}_{|z|} \otimes \mathbf{e}_{|S||A|}) \\
 &= \left(\sum_{j=1}^K x_j P_{Az}^{(j)} \otimes P_B^{(j)} \right)^{t-1} \left(\sum_{j=1}^K x_j \left(P_{Az}^{(j)} (\mathbf{r}_A \otimes \mathbf{e}_{|z|}) \right) \otimes \left(P_B^{(j)} \mathbf{e} \right) \right) \\
 &= \left(\sum_{j=1}^K x_j P_{Az}^{(j)} \otimes P_B^{(j)} \right)^{t-1} \left(\sum_{j=1}^K x_j P_{Az}^{(j)} (\mathbf{r}_A \otimes \mathbf{e}_{|z|}) \right) \otimes \mathbf{e} \\
 &= \dots = \left(\left(\sum_{j=1}^K x_j P_{Az}^{(j)} \right)^t (\mathbf{r}_A \otimes \mathbf{e}_{|z|}) \right) \otimes \mathbf{e}.
 \end{aligned}$$

792

□

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

$$\begin{aligned}
 E_A(\mathbf{P}_{AB}^\pi, \mathcal{Z}) &:= \min \frac{1}{2} \left\| \mathbf{P}_{ABz}^\pi - \sum_{j=1}^K x_j P_{Az}^{(j)} \otimes P_B^{(j)} \right\|_{\text{ATV}_1} \\
 &= \min_{P_{Az}} \max_{\mathbf{s}, \mathbf{a}, z} \frac{1}{2} \sum_{s'_A, a'_A, z'} |P_{ABz}^\pi(s'_A, a'_A, z' | \mathbf{s}, \mathbf{a}, z) - P_{Az}(s'_A, a'_A, z' | s_A, a_A, z)|.
 \end{aligned} \tag{16}$$

795 Similar to Theorem 3, we can connect this measure of Markov entanglement with the value decom-
 796 position error.

797 **Theorem 10.** Consider a N -agent Markov system $\mathcal{M}_{1:N}$. Given any policy $\pi: \mathcal{S} \rightarrow \Delta(\mathcal{A})$ with the
 798 measure of Markov entanglement $E_i(\mathbf{P}_{1:N}^\pi, \mathcal{Z})$ w.r.t the agent-wise total variation distance, it holds
 799 for any agent i ,

$$\left\| \mathbf{P}_{iz}^\pi - \sum_{j=1}^K x_j P_{iz}^{(j)} \right\|_\infty \leq 2E_i(\mathbf{P}_{1:N}^\pi, \mathcal{Z}).$$

800 Furthermore, the decomposition error is entry-wise bounded by the measure of Markov entangle-
 801 ment,

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

802 In practice, exogenous information is often discussed in the context of (weakly-)coupled MDPs,
 803 where each agent independent evolves by $P_i(s_{i+1} | s_i, a_i, z)$. Interestingly, we can derive a similar
 804 result to Proposition 8 that shaves off the transition in entanglement analysis.

805 **Proposition 11.** Consider a N -agent Weakly Coupled Markov system $\mathcal{M}_{1:N}$. Given any policy
 806 $\pi: \mathcal{S} \rightarrow \Delta(\mathcal{A})$ 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-
 807 wise total variation distance, it holds

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

808 for any policies π' .

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

$$\begin{aligned}
& E_A(\mathbf{P}_{AB}^\pi, \mathcal{Z}) \\
& \leq \frac{1}{2} \sum_{\mathbf{s}, \mathbf{a}, z} \mu(\mathbf{s}, \mathbf{a}, z) \sum_{s'_A, a'_A, z'} |P_{ABz}^\pi(s'_A, a'_A, z' | \mathbf{s}, \mathbf{a}, z) - P_{Az}(s'_A, a'_A, z' | s_A, a_A, z)| \\
& = \frac{1}{2} \sum_{\mathbf{s}, \mathbf{a}, z} \mu(\mathbf{s}, \mathbf{a}, z) \sum_{s'_A, a'_A, z'} \left| \sum_{s'_B} P_{ABz}^\pi(s', a_A, z' | \mathbf{s}, \mathbf{a}, z) - P_{Az}(s'_A, z' | s_A, a_A, z) \pi'(a'_A | s'_A, z') \right| \\
& = \frac{1}{2} \sum_{\mathbf{s}, \mathbf{a}, z} \mu(\mathbf{s}, \mathbf{a}, z) \sum_{s'_A, a'_A, z'} \left| \sum_{s'_B} P_{ABz}^\pi(s', a_A, z' | \mathbf{s}, \mathbf{a}, z) - \sum_{s'_B} P(s', z' | \mathbf{s}, \mathbf{a}, z) \pi'(a'_A | s'_A, z') \right| \\
& = \frac{1}{2} \sum_{\mathbf{s}, \mathbf{a}, z} \mu(\mathbf{s}, \mathbf{a}, z) \sum_{s'_A, a'_A, z'} \left| \sum_{s'_B} P(s', z' | \mathbf{s}, \mathbf{a}, z) (\pi(a'_A | s', z') - \pi'(a'_A | s'_A, z')) \right| \\
& \leq \frac{1}{2} \sum_{\mathbf{s}, \mathbf{a}, z} \mu(\mathbf{s}, \mathbf{a}, z) \sum_{s', z'} P(s', z' | \mathbf{s}, \mathbf{a}, z) \sum_{a'_A} |\pi(a'_A | s', z') - \pi'(a'_A | s'_A, z')| \\
& = \frac{1}{2} \sum_{s', z'} \mu(s', z') \sum_{a'_A} |\pi(a'_A | s', z') - \pi'(a'_A | s'_A, z')|.
\end{aligned}$$

810 □

811 L.3 Factored MDPs

812 Another common class of multi-agent MDPs is Factored MDPs (FMDPs, [Guestrin et al. \(2001;](#)
813 [2003\)](#); [Osband & Roy \(2014\)](#)), which explicitly model the structured dependencies in state transi-
814 tions. For instance, in a server cluster, the state transition of each server depends only on its
815 neighboring servers. Formally, we define

816 **Definition 11** (Factored MDPs). *An N -agent MDP $\mathcal{M}_{1:N}(\mathcal{S}, \mathcal{A}, \mathbf{P}, \mathbf{r}_{1:N}, \gamma)$ is a factored MDP if*
817 *each agent i has neighbor set $Z_i \in [N]$ such that its transition is affected by all its neighbors, i.e.*
818 $P(s'_i | \mathbf{s}, \mathbf{a}) = P(s'_i | s_{Z_i}, a_{Z_i})$.

819 The neighbor set $|Z_i|$ is often assumed to be much smaller compared to the number of agents N .
820 This helps to encode exponentially large system very compactly. We show this idea can also be cap-
821 tured in Markov entanglement. Consider the measure of Markov entanglement w.r.t ATV distance
822 in Eq. (6),

$$\begin{aligned}
E_A(\mathbf{P}_{AB}^\pi) &= \min_{\mathbf{P}_A} \max_{(\mathbf{s}, \mathbf{a}) \in \mathcal{S} \times \mathcal{A}} D_{\text{TV}} \left(\mathbf{P}_{AB}^\pi(\cdot, \cdot | \mathbf{s}, \mathbf{a}), \mathbf{P}_A(\cdot, \cdot | s_A, a_A) \right) \\
&= \min_{\mathbf{P}_A} \max_{(\mathbf{s}, \mathbf{a}) \in \mathcal{S} \times \mathcal{A}} D_{\text{TV}} \left(\mathbf{P}_{AB}^\pi(\cdot, \cdot | s_{Z_A}, a_{Z_A}), \mathbf{P}_A(\cdot, \cdot | s_A, a_A) \right).
\end{aligned}$$

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

824 L.4 Fully cooperative Markov games

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

828 This algorithm follows similar framework of meta algorithm 1 and differs at we now learn the
829 closet local reward decomposition from data. When the reward is completely decomposable, meta
830 algorithm 2 recovers meta algorithm 1. Thus intuitively, the more accurate we can decompose the

Meta Algorithm 2: Q-value Decomposition with Shared Reward

Require: Global policy π ; horizon length T .

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

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

831 global reward, the less decomposition error we have. Formally, we define the measure of reward
 832 entanglement

$$e(\mathbf{r}_{AB}) := \min_{\mathbf{r}_A, \mathbf{r}_B} \|\mathbf{r}_{AB} - (\mathbf{r}_A \otimes \mathbf{e} + \mathbf{e} \otimes \mathbf{r}_B)\|_{\mu_{AB}^\pi}. \quad (17)$$

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

835 **Proposition 12.** Consider a fully cooperative two-agent Markov system \mathcal{M}_{AB} . Given any policy
 836 $\pi: \mathcal{S} \rightarrow \Delta(\mathcal{A})$ with the measure of Markov entanglement $E_A(\mathbf{P}_{AB}^\pi), E_B(\mathbf{P}_{AB}^\pi)$ w.r.t the μ_{AB}^π -
 837 weighted agent-wise total variation distance and the measure of reward entanglement $e(\mathbf{r}_{AB})$, it
 838 holds

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

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

840 Although Proposition 1 offers a theoretical guarantee for general two-agent fully cooperative
 841 Markov games, its utility is greatest in systems with low reward and transition entanglement. Fully
 842 cooperative settings remain inherently challenging—for instance, even the asymptotically optimal
 843 Whittle Index may achieve only a $\frac{1}{N}$ -approximation ratio for RMABs with global rewards [Raman
 844 et al. \(2024\)](#). In practice, most research [Sunehag et al. \(2018\)](#); [Rashid et al. \(2020\)](#) relies on sophis-
 845 ticated deep neural networks to learn decompositions in such settings. We thus defer a more refined
 846 analysis of fully cooperative scenarios to future work.

847 M Simulation environments

848 In this section, we empirically study the value decomposition for index policies. Our simulations
 849 build on a circulant RMAB benchmark, which is widely used in the literature [Avrachenkov & Borkar
 850 \(2022\)](#); [Zhang & Frazier \(2022\)](#); [Biswas et al. \(2021\)](#); [Fu et al. \(2019\)](#).

851 **Circulant RMAB** A circulant RMAB has four states indexed by $\{0, 1, 2, 3\}$. Transition kernels
 852 $P_a = p(s, s')_{s, s' \in \mathcal{S}}$ for action $a = 0$ and $a = 1$ are given by

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

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

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

854 We set the discount factor to $\gamma = 0.5$ and require $N/5$ arms to be pulled per period. Initially, there
 855 are $N/6$ arms in state 0, $N/3$ arms in state 1 and $N/2$ arms in state 2, the same as [Zhang & Frazier
 856 \(2022\)](#). We then test an index policy with priority: state 2 > state 1 > state 0 > state 3.

857 **M.1 Monte-Carlo estimation of Markov entanglement**

858 For each RMAB instance, we simulate a trajectory of length $T = 6N$ and collect data for the later
859 $5N$ epochs. Notice RMAB is a special instance of WCMDP, we thus apply the result in Proposition 8

$$\begin{aligned} E_i(\mathbf{P}_{1:N}^\pi) &\leq \frac{1}{2} \min_{\pi'} \sum_{\mathbf{s}} \mu_{1:N}^\pi(\mathbf{s}) \sum_{a_i} |\pi(a_i | \mathbf{s}) - \pi'(a_i | s_i)| \\ &\approx \frac{1}{2} \min_{\pi'} \frac{1}{T} \sum_{t=1}^T \sum_{a_i} |\pi(a_i | \mathbf{s}) - \pi'(a_i | s_i)| \end{aligned} \quad (18)$$

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

863 **M.2 Learning local Q-values**

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

867 The ground-truth $Q_{1:N}^\pi$ is approximated via Monte Carlo learning Sutton & Barto (2018), with
868 each estimate derived from 30-step simulations averaged over $3N$ independent runs. Due to the
869 high computational cost of Monte Carlo methods—especially for very large RMABs—we limit the
870 training phase to 10 independent runs and use the mean local Q-value as an approximation. Error
871 bars represent the standard error for both Monte Carlo estimates and μ -weighted decomposition
872 errors.

873 In addition to μ -weighted error, we also introduce a concept of relative error, defined as
874 $\left\| Q_{1:N}^\pi(\mathbf{s}, \mathbf{a}) - \sum_{i=1}^N Q_i^\pi(s_i, a_i) \right\|_{\mu_{1:N}^\pi} / \|Q_{1:N}^\pi\|_{\mu_{1:N}^\pi}$. This relative error reflects the approximate
875 ratio of our value decomposition. We present our simulation results below.

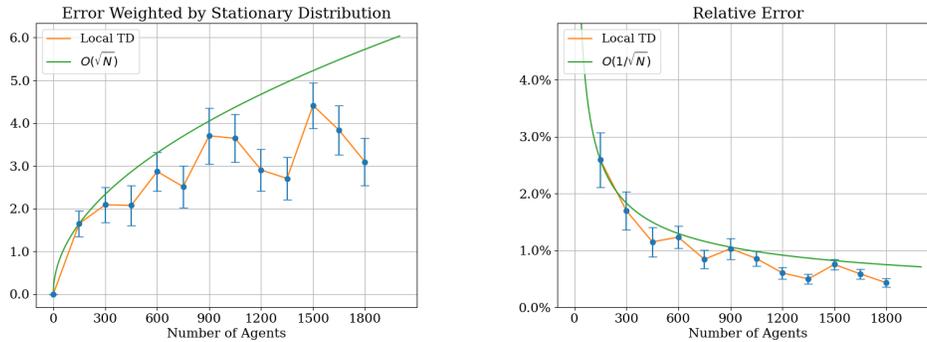


Figure 2: Value Decomposition error in circulant RMAB under an index policy. *Left*: μ -weighted decomposition error. *Right*: Relative error, $\|\text{decomposition error}\|_{\mu} / \|Q_{1:N}^\pi\|_{\mu}$

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

881 **M.3 Sample Complexity and Computation**

882 While each RMAB instance has an exponentially large state space $|S|^N$, we show that our empir-
 883 ical 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

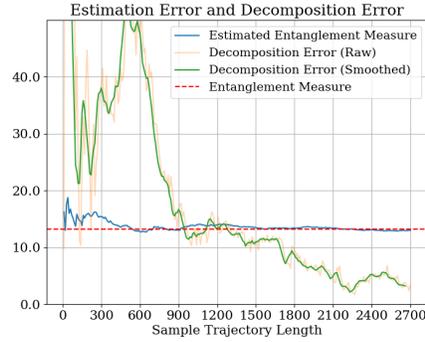


Figure 3: Different errors in RMAB with 900 agents: empirical estimation of Markov entanglement (blue); $\mu_{1:N}^\pi$ -weighted decomposition error (green); the true measure of Markov estimated with $T = 10N$ samples (red dashed line).

884 exhibits in Figure 3, both errors decay and converges within $T = 3N$ samples. Furthermore, the
 885 empirical estimation of Markov entanglement converges in $T < N$ samples, demonstrating its ef-
 886 ficiency. Finally, we use standard convex optimization solvers to compute Markov entanglement,
 887 which can be run efficiently on a single CPU.
 888