

## 469 A Proofs

### 470 A.1 Proof of Theorem 4.1

471 *Proof.* For any two global state value functions  $V_{tot}^1$  and  $V_{tot}^2$ , let  $\pi_{tot}^1$  be the optimal global policy  
472 under  $\mathcal{T}_f^* V_{tot}^1$ , and we can get:

$$\begin{aligned}
& (\mathcal{T}_f^* V_{tot}^1)(\mathbf{o}) - (\mathcal{T}_f^* V_{tot}^2)(\mathbf{o}) \\
&= \sum_{\mathbf{a}} \pi_{tot}^1(\mathbf{a}|\mathbf{o}) \left[ r + \gamma \mathbb{E}_{\mathbf{o}'} [V_{tot}^1(\mathbf{o}')] - \alpha \log \left( \frac{\pi_{tot}^1(\mathbf{a}|\mathbf{o})}{\mu_{tot}(\mathbf{a}|\mathbf{o})} \right) \right] - \max_{\mathbf{a}} \sum_{\mathbf{a}} \pi_{tot}(\mathbf{a}|\mathbf{o}) \left[ r + \gamma \mathbb{E}_{\mathbf{o}'} [V_{tot}^2(\mathbf{o}')] - \alpha \log \left( \frac{\pi_{tot}(\mathbf{a}|\mathbf{o})}{\mu_{tot}(\mathbf{a}|\mathbf{o})} \right) \right] \\
&\leq \sum_{\mathbf{a}} \pi_{tot}^1(\mathbf{a}|\mathbf{o}) \left[ r + \gamma \mathbb{E}_{\mathbf{o}'} [V_{tot}^1(\mathbf{o}')] - \alpha \log \left( \frac{\pi_{tot}^1(\mathbf{a}|\mathbf{o})}{\mu_{tot}(\mathbf{a}|\mathbf{o})} \right) \right] - \sum_{\mathbf{a}} \pi_{tot}^1(\mathbf{a}|\mathbf{o}) \left[ r + \gamma \mathbb{E}_{\mathbf{o}'} [V_{tot}^2(\mathbf{o}')] - \alpha \log \left( \frac{\pi_{tot}^1(\mathbf{a}|\mathbf{o})}{\mu_{tot}(\mathbf{a}|\mathbf{o})} \right) \right] \\
&= \gamma \sum \pi_{tot}^1(\mathbf{a}|\mathbf{o}) \mathbb{E}_{\mathbf{o}'} [V_{tot}^1(\mathbf{o}') - V_{tot}^2(\mathbf{o}')] \\
&\leq \gamma \|V_{tot}^1 - V_{tot}^2\|_{\infty}
\end{aligned}$$

473 Therefore, it follows that:

$$\|\mathcal{T}_f^* V_{tot}^1 - \mathcal{T}_f^* V_{tot}^2\|_{\infty} \leq \gamma \|V_{tot}^1 - V_{tot}^2\|_{\infty}$$

474

□

### 475 A.2 Proof of Proposition 4.2

476 *Proof.* For a behavior-regularized Dec-POMDP with  $f(\pi_{tot}, \mu_{tot}) = \log(\pi_{tot}/\mu_{tot})$ , the learning  
477 objective can be written as  $\max_{\pi_{tot}} \mathbb{E} [\sum_{t=0}^{\infty} \gamma^t (r(\mathbf{o}_t, \mathbf{a}_t) - \alpha \log(\pi_{tot}(\mathbf{a}_t|\mathbf{o}_t)/\mu_{tot}(\mathbf{a}_t|\mathbf{o}_t)))]$ . Its  
478 Lagrangian function can obtain when the optimal global policy is written as follows:

$$\begin{aligned}
L(\pi_{tot}, \beta, u) &= \sum_{\mathbf{o}} d_{\pi_{tot}}(\mathbf{o}) \sum_{\mathbf{a}} \pi_{tot}(\mathbf{a}|\mathbf{o}) \left( Q_{tot}(\mathbf{o}, \mathbf{a}) - \alpha \log \left( \frac{\pi_{tot}(\mathbf{a}|\mathbf{o})}{\mu_{tot}(\mathbf{a}|\mathbf{o})} \right) \right) \\
&\quad - \sum_{\mathbf{o}} d_{\pi_{tot}}(\mathbf{o}) \left[ u(\mathbf{o}) \left( \sum_{\mathbf{a}} \pi_{tot}(\mathbf{a}|\mathbf{o}) - 1 \right) + \sum_{\mathbf{a}} \beta(\mathbf{a}|\mathbf{o}) \pi_{tot}(\mathbf{a}|\mathbf{o}) \right],
\end{aligned}$$

479 where  $d_{\pi_{tot}}$  is the stationary joint observation distribution of the global policy  $\pi_{tot}$ .  $u$  and  $\beta$  are  
480 Lagrangian multipliers for the equality and inequality constraints.

481 According to the Karush-Kuhn-Tucker (KKT) conditions where the derivative of the Lagrangian  
482 objective function with respect to the global policy is zero at the optimal solution, it follows that:

$$\begin{aligned}
Q_{tot}(\mathbf{o}, \mathbf{a}) - \alpha \left( \log \left( \frac{\pi_{tot}(\mathbf{a}|\mathbf{o})}{\mu_{tot}(\mathbf{a}|\mathbf{o})} \right) + 1 \right) - u(\mathbf{o}) + \beta(\mathbf{a}|\mathbf{o}) &= 0 \quad (16) \\
\sum_{\mathbf{a}} \pi_{tot}(\mathbf{a}|\mathbf{o}) &= 1 \\
\beta(\mathbf{a}|\mathbf{o}) \pi_{tot}(\mathbf{a}|\mathbf{o}) &= 0 \\
0 \leq \pi_{tot}(\mathbf{a}|\mathbf{o}) \leq 1 \text{ and } 0 \leq \beta(\mathbf{a}|\mathbf{o}) &
\end{aligned}$$

483 From Eq. (16), we can further solve the optimal global policy as:

$$\pi_{tot}(\mathbf{a}|\mathbf{o}) = \mu_{tot}(\mathbf{a}|\mathbf{o}) \cdot \exp \left( \frac{Q_{tot}(\mathbf{o}, \mathbf{a}) - u(\mathbf{o}) + \beta(\mathbf{a}|\mathbf{o})}{\alpha} - 1 \right)$$

484 The above formula can be further simplified.  $\beta$  is the Lagrangian multiplier, and meets comple-  
485 mentary slackness  $\beta(\mathbf{a}|\mathbf{o}) \pi_{tot}(\mathbf{a}|\mathbf{o}) = 0$ . Considering the joint observation  $\mathbf{o}$  is fixed,  
486  $\exp \left( \frac{Q_{tot}(\mathbf{o}, \mathbf{a}) - u(\mathbf{o}) + \beta(\mathbf{a}|\mathbf{o})}{\alpha} - 1 \right)$  is always larger than 0. Therefore, for any positive probability  
487 action, its corresponding Lagrangian multiplier  $\beta(\mathbf{a}|\mathbf{o})$  is 0. Therefore,  $\pi_{tot}(\mathbf{a}|\mathbf{o})$  can be reformulated  
488 as:

$$\pi_{tot}(\mathbf{a}|\mathbf{o}) = \mu_{tot}(\mathbf{a}|\mathbf{o}) \cdot \exp \left( \frac{Q_{tot}(\mathbf{o}, \mathbf{a}) - u(\mathbf{o})}{\alpha} - 1 \right) \quad (17)$$

489 Bringing Eq. (17) into  $\sum_{\mathbf{a}} \pi_{tot}(\mathbf{a}|\mathbf{o}) = 1$ , we have:

$$\mathbb{E}_{\mathbf{a} \sim \mu_{tot}} \left[ \exp \left( \frac{Q_{tot}(\mathbf{o}, \mathbf{a}) - u(\mathbf{o})}{\alpha} - 1 \right) \right] = 1 \quad (18)$$

490 The left side of Eq. (18) can be seen as a continuous and monotonic function of  $u$ , so it has only one  
 491 solution denoted as  $u^*$ , and we denote the corresponding policy  $\pi_{tot}$  as  $\pi_{tot}^*$ .

492 Integrating Eq. (17) into the expression of optimal global state value, we can get:

$$\begin{aligned} V_{tot}^*(\mathbf{o}) &= \mathcal{T}_f^* V_{tot}^*(\mathbf{o}) \\ &= \sum_{\mathbf{a}} \pi_{tot}^*(\mathbf{a}|\mathbf{o}) \left( Q_{tot}^*(\mathbf{o}, \mathbf{a}) - \alpha \log \left( \frac{\pi_{tot}^*(\mathbf{a}|\mathbf{o})}{\mu_{tot}(\mathbf{a}|\mathbf{o})} \right) \right) \\ &= \sum_{\mathbf{a}} \pi_{tot}^*(\mathbf{a}|\mathbf{o}) (u^*(\mathbf{o}) + \alpha) \\ &= u^*(\mathbf{o}) + \alpha \end{aligned}$$

493 To summarize, we obtain the optimality condition of the behavior regularized MDP with Reverse KL  
 494 divergence as follows:

$$\begin{aligned} Q_{tot}^*(\mathbf{o}, \mathbf{a}) &= r(\mathbf{o}, \mathbf{a}) + \gamma \mathbb{E}_{\mathbf{o}'|\mathbf{o}, \mathbf{a}} [V_{tot}^*(\mathbf{o}')] \\ V_{tot}^*(\mathbf{o}) &= u^*(\mathbf{o}) + \alpha \\ \pi_{tot}^*(\mathbf{a}|\mathbf{o}) &= \mu_{tot}(\mathbf{a}|\mathbf{o}) \cdot \exp \left( \frac{Q_{tot}^*(\mathbf{o}, \mathbf{a}) - u^*(\mathbf{o})}{\alpha} - 1 \right) \end{aligned}$$

495 where  $u(\mathbf{o})$  is a normalization term and has a optimal value  $u^*$  that makes the corresponding optimal  
 496 policy  $\pi_{tot}^*$  satisfy  $\sum_{\mathbf{a} \in \mathcal{A}^n} \pi_{tot}^*(\mathbf{a}|\mathbf{o}) = 1$ .

497 □

## 498 B Experiment Settings

### 499 B.1 Multi-Agent MuJoCo

500 Multi-agent Mujoco [7] is a benchmark framework developed for assessing and comparing the  
 501 effectiveness of algorithms in continuous multi-agent robotic control. Within this framework, a  
 502 robotic system is partitioned into independent agents, each tasked with controlling a specific set of  
 503 joints. The agents collaborate harmoniously to accomplish shared objectives, such as acquiring the  
 504 ability to walk through an environment, with the ultimate goal of maximizing the cumulative reward.  
 505 Multi-agent MuJoCo environment consists of multiple different robot configurations, and it is often  
 506 used for the study of novel MARL algorithms for decentralized coordination in isolation.

507 To generate the dataset transitions, we captured the interactions between the environment and trained  
 508 online MARL algorithms. Specifically, we use HAPPO [16] algorithm to collect data. The expert  
 509 dataset is generated by employing the converged HAPPO algorithm. This involves training the  
 510 algorithm until it reaches a state of convergence, where the agents have learned optimal policies.  
 511 The medium dataset is generated by first training a policy online using HAPPO, early-stopping  
 512 the training, and collecting samples from this partially-trained policy. The medium-replay dataset  
 513 consists of recording all samples in the replay buffer observed during training until the policy reaches  
 514 the medium level of performance. The medium-expert dataset by mixing equal amounts of expert  
 515 demonstrations and suboptimal data. For all datasets, the hyperparameter `env_args.agent_obsk`  
 516 (determines up to which connection distance agents will be able to form observations) is set to 1. The  
 517 reward distribution of our datasets is listed in Table 2.

### 518 B.2 The StarCraft Multi-Agent Challenge

519 The StarCraft Multi-Agent Challenge (SMAC) benchmark is chosen as our testing environment. Due  
 520 to its high control complexity, SMAC is a popular multi-agent cooperative control environment for  
 521 evaluating advanced MARL methods. It consists of a collection of StarCraft II microscenarios in

Table 2: The multi-agent MuJoCo dataset

Scenario	Quality	Reward Distribution
2-Agent Ant	expert	2.06±0.35
	medium	1.42±0.37
	medium-expert	1.74±0.48
	medium-replay	1.03±0.21
3-Agent Hopper	expert	3.64±0.79
	medium	3.16±1.00
	medium-expert	3.41±0.93
	medium-replay	2.37±0.69
6-Agent HalfCheetah	expert	2.79±1.76
	medium	1.43±1.36
	medium-expert	2.11±1.71
	medium-replay	0.66±1.09

522 which two groups of units engage in combat. Agents based on the MARL algorithm control the first  
523 group’s units, while a built-in heuristic game AI bot with different difficulties controls the second  
524 group’s units. Scenarios vary in terms of the initial location, number and type of units, and elevated  
525 or impassable terrain. The available actions for each agent include no operation, move[direction],  
526 attack [enemy id], and stop. The reward that each agent receives is the same. The hit-point damage  
527 dealt and received determines the agents’ share of the reward. SMAC consists of several StarCraft II  
528 multi-agent micromanagement maps. We consider 4 representative battle maps, including 2 hard map  
529 (5m\_vs\_6m, 2c\_vs\_64zg), and 2 super hard maps (6h\_vs\_8z, corridor), as our experiment tasks. The  
530 task type and other details of the maps are listed in the Table 3.

Table 3: SMAC maps for experiments.

Map Name	Ally Units	Enemy Units	Type
5m_vs_6m	5 Marines	6 Marines	homogeneous & asymmetric
2c_vs_64zg	2 Colossi	64 Zerglings	micro-trick: positioning
6h_vs_8z	6 Hydralisks	8 Zealots	micro-trick: focus fire
corridor	6 Zealots	24 Zerglings	micro-trick: wall off

531 The offline SMAC dataset used in this study is provided by [22], which is the largest open offline  
532 dataset on SMAC. Different from single-agent offline datasets, it considers the property of Dec-  
533 POMDP, which owns local observations and available actions for each agent. The dataset is collected  
534 from the trained MAPPO agent, and includes three quality levels: good, medium, and poor. For each  
535 original large dataset, we randomly sample 1000 episodes as our dataset.

## 536 C Implementation Details

### 537 C.1 Details of OMIGA

538 The local Q-value, state value networks and policy networks of OMIGA are represented by 3-layers  
539 ReLU activated MLPs with 256 units for each hidden layer. For the weight network, we use 2-layer  
540 ReLU activated MLPs with 64 units for each hidden layer. All the networks are optimized by Adam  
541 optimizer.

### 542 C.2 Details of baselines

543 We compare OMIGA against four recent offline MARL algorithms: ICQ [40], OMAR [25], BCQ-MA  
544 and CQL-MA. For the ICQ and OMAR, we implement them based on the algorithm description  
545 in their papers. BCQ-MA is the multi-agent version of BCQ, and CQL-MA is the multi-agent  
546 version of CQL. BCQ-MA and CQL-MA use linear weighted value decomposition structure as  
547  $Q_{tot} = \sum_{i=1}^n w_i(o)Q_i(o_i, a_i) + b(o)$ ,  $w^i \geq 0$  for the multi-agent setting. The policy constrain of  
548 BCQ-MA and the value regularization of CQL-MA are both imposed on the local Q-value.

549 In this paper, all experiments are implemented with pytorch and executed on NVIDIA V100 GPUs.

### 550 C.3 Hyperparameters

551 For multi-agent MuJoCo, the hyperparameters of OMIGA and baselines are listed in Table 4.  
 552 An important hyperparameter of OMIGA is the regularization hyperparameter  $\alpha$ . The higher  $\alpha$   
 553 encourages OMIGA staying near the behavioral distribution, and lower  $\alpha$  makes OMIGA more  
 554 optimistic. On most tasks, we use  $\alpha = 10$  to ensure good regularization effect. On the medium  
 555 quality dataset of HalfCheetah task, we choose  $\alpha = 1$ .

Table 4: Hyperparameters of OMIGA and baselines for multi-agent MuJoCo

Hyperparameter	Value
<b>Shared parameters</b>	
Q-value network learning rate	5e-4
Policy network learning rate	5e-4
Optimizer	Adam
Target update rate	0.005
Batch size	128
Discount factor	0.99
Hidden dimension	256
Weight network hidden dimension	64
<b>OMIGA</b>	
State value network learning rate	5e-4
Regularization parameter $\alpha$	1 or 10
<b>Others</b>	
Lagrangian coefficient (ICQ)	10
Tradeoff factor $\alpha$ (OMAR, CQL-MA)	1

556 For SMAC, the hyperparameters of OMIGA and baselines are listed in Table 5. On most tasks, we  
 557 use  $\alpha = 10$ . On the poor dataset of 6h\_vs\_8z map, the quality of the dataset is relatively poor. It  
 558 does not make much sense to make the policy close to the behavioral policy, so we choose  $\alpha = 2$  to  
 559 make the algorithm more radical.

560 On the most SMAC maps, the learning rate of all networks is set to 5e-4. The exception is the map  
 561 2c\_vs\_64zg, on this map, the learning rate of all networks is set to 1e-4.

Table 5: Hyperparameters of OMIGA and baselines for SMAC

Hyperparameter	Value
<b>Shared parameters</b>	
Q-value network learning rate	5e-4 or 1e-4
Policy network learning rate	5e-4 or 1e-4
Optimizer	Adam
Target update rate	0.005
Batch size	128
Discount factor	0.99
Hidden dimension	256
Weight network hidden dimension	64
<b>OMIGA</b>	
State value network learning rate	5e-4 or 1e-4
Regularization parameter $\alpha$	2 or 10
<b>Others</b>	
Lagrangian coefficient (ICQ)	10
Threshold (BCQ-MA)	0.3
Tradeoff factor $\alpha$ (OMAR, CQL-MA)	1

Table 6: Average scores and standard deviations over 5 random seeds on the mixed offline SMAC datasets

Map	Dataset	BCQ-MA	CQL-MA	ICQ	OMAR	OMIGA(ours)
6h_vs_8z	good-poor	11.41±0.44	9.56±0.25	11.00±0.36	9.17±0.19	<b>11.88±0.27</b>
6h_vs_8z	good-medium	11.79±0.29	10.08±0.26	11.18±0.25	10.02±0.16	<b>12.05±0.47</b>
6h_vs_8z	medium-poor	11.18±0.41	10.73±0.38	11.25±0.35	10.42±0.19	<b>11.85±0.35</b>
corridor	good-poor	12.37±1.36	4.88±0.35	11.78±1.53	5.54±0.75	<b>13.01±0.89</b>
corridor	good-medium	13.32±0.71	5.77±1.30	12.98±0.62	6.63±0.74	<b>14.02±1.04</b>
corridor	medium-poor	8.11±0.35	6.18±0.59	8.27±0.48	6.25 ±0.48	<b>9.70±1.40</b>

## 562 D Additional Results

### 563 D.1 Results on mixed datasets

564 We want to investigate whether OMIGA has superior performance when the datasets are mixed.  
 565 Unlike BCQ-MA and OMAR, OMIGA doesn't need to learn a behavior policy. We choose two  
 566 original datasets on the SMAC super hard maps 6h\_vs\_8z and corridor, and make mixed datasets by  
 567 combining these SMAC datasets of different quality, including good-poor, good-medium, medium-  
 568 poor datasets. Each mixed dataset is blended by 50% of each of the two original datasets. On these  
 569 mixed suboptimal datasets, the behavior policy is heterogeneous. Therefore, it is more difficult for  
 570 algorithms such as BCQ-MA and OMAR to learn an accurate behavior policy, making implicit value  
 571 learning with the regularization framework of OMIGA more appealing.

572 Table 6 shows the results that OMIGA consistently outperforms other offline MARL baselines  
 573 under all different mixed dataset experiments. Compared with the results on the original datasets,  
 574 the performance of OMIGA has become more leading, indicating the benefits of implicit value  
 575 regularization of OMIGA.