# B3C: A Minimalist Approach to Offline Multi-Agent Reinforcement Learning

**Woojun Kim** [1]  **Katia Sycara** [1]

## Abstract

Overestimation arising from selecting unseen actions during policy evaluation is a major challenge in offline reinforcement learning (RL). A minimalist approach in the single-agent setting—adding behavior cloning (BC) regularization to existing online RL algorithms—has been shown to be effective; however, this approach is understudied in multi-agent settings. In particular, overestimation becomes worse in multi-agent settings due to the presence of multiple actions, resulting in the BC regularization-based approach easily suffering from either over-regularization or critic divergence. To address this, we propose a simple yet effective method, Behavior Cloning regularization with Critic Clipping (B3C), which clips the target critic value in policy evaluation based on the maximum return in the dataset and pushes the limit of the weight on the RL objective over BC regularization, thereby improving performance. Additionally, we leverage existing value factorization techniques, particularly non-linear factorization, which is understudied in offline settings. Integrated with non-linear value factorization, B3C outperforms state-of-the-art algorithms on various offline multi-agent benchmarks.

## 1. Introduction

Reinforcement learning (RL) trains agents to maximize cumulative rewards through interaction with an environment in an online manner; however, such interaction can be time-consuming due to the sample inefficiency of RL and costly in safety-critical environments due to inherent risks. Furthermore, only existing datasets are available for use in practice. These challenges have motivated a shift toward offline RL, which trains agents using pre-collected datasets without requiring interaction.

A major challenge of offline RL is that bootstrapping from unseen actions in policy evaluation induces extrapolation error, leading to accumulated errors (Fujimoto & Gu, 2021), due to the inability to generate data. This issue consequently causes overestimation of value functions and adversely affects policy improvement. To address this, a variety of techniques, such as conservative value estimation (Kumar et al., 2020) and modeling behavior policy (Fujimoto et al., 2019), have been introduced. However, these approaches increase algorithmic complexity, hinder reproducibility, and amplify hyperparameter sensitivity. This motivates minimalist approaches, which introduce minimal changes to existing RL algorithms (Fujimoto & Gu, 2021; Tarasov et al., 2024). One representative example is TD3+BC, which integrates behavior cloning (BC) regularization into an existing online RL algorithm called TD3 (Fujimoto & Gu, 2021). This simple approach does not require any additional components and introduces only one more hyperparameter for regularization, but it surprisingly achieves SOTA performance. This line of research—minimalist approaches—has been extensively investigated in single-agent settings (Fujimoto & Gu, 2021; Tarasov et al., 2024), but it has barely been explored in multi-agent settings, which entail more algorithmic complexity due to the presence of multiple agents.

In this paper, we aim to develop a minimalist approach that introduces minimal changes to existing multi-agent RL algorithms. One might ask: is adding BC regularization sufficient in multi-agent settings? We argue that BC regularization alone is not enough in multi-agent settings. This is because the overestimation becomes more severe in multi-agent settings. It arises from the use of a centralized critic, which conditions on multiple actions, increasing the likelihood of encountering unseen actions. Consequently, critic learning often becomes unstable, necessitating a greater reliance on BC regularization over the RL objective—resulting in over-regularization. This over-regularization inherently limits performance, as it becomes heavily dependent on the quality of the dataset. In addition, recent research on offline multi-agent RL overlooks existing value factorization techniques, such as non-linear decomposition—they either rely on linear factorization or disregard factorization entirely.

To address the aforementioned limitations, we propose a simple yet effective technique for offline multi-agent RL: **B**ehavior **C**loning regularization with **C**ritic **C**lipping (**B3C**).

---

[1]Robotics Institute, Carnegie Mellon University, Pittsburgh, United States. Correspondence to: Woojun Kim <woojunk@andrew.cmu.edu>.

The proposed method alleviates overestimation by clipping the critic value during policy evaluation, using the maximum return from the given dataset, in addition to BC regularization. Critic clipping enables the use of a higher weight for the RL objective relative to BC regularization, resulting in improved performance. Furthermore, we integrate B3C with an existing online multi-agent RL algorithm called FACMAC (Peng et al., 2021), which utilizes factored critics and value factorization techniques, resulting in **FACMAC+B3C**. Notably, we provide a design choice for value factorization in the offline setting, empirically demonstrating that non-monotonic factorization performs better than monotonic and linear factorization, which contrasts with findings in the online setting (Peng et al., 2021).

Despite its simplicity, the effectiveness of FACMAC+B3C is empirically demonstrated across a variety of offline multi-agent RL benchmarks, including multi-agent Mujoco and particle environments, with various types of dataset. We provide a thorough analysis demonstrating how the proposed method outperforms the baselines in terms of performance.

## 2. Background and Related Works

### 2.1. Multi-Agent RL

We consider fully cooperative multi-agent tasks, which can be modeled as an N-agent decentralized partially observable MDP (Dec-POMDP) (Oliehoek, 2012). At each timestep $t$, each agent selects its own action, $a_t^i$, based on its local information such as the partial observation $o_t^i$. This forms a joint action $\boldsymbol{a}_t$ yielding the next state $s_{t+1}$, the joint observations $\boldsymbol{o}_{t+1} = (o_{t+1}^1, \cdots, o_{t+1}^N)$, and a shared reward $r_t$. The goal is to find the optimal joint policy that maximizes the team return, $\mathbb{E}[\sum_{l=0}^{\infty} \gamma^l r_l]$. For this, recent multi-agent RL algorithms adopt a framework of centralized training with decentralized execution (CTDE) (Wang et al., 2020; Zhang et al., 2021; Jeon et al., 2022; Kim et al., 2023), where individual policies execute actions based on local information, e.g., Agent $i$'s action-observation history $\tau^i$, but are trained with global information such as the state.

**FACMAC** (Peng et al., 2021) is a representative multi-agent RL algorithm that trains a centralized but factored critic, which is decomposed into a nonlinear combination of individual critics, and uses it to train decentralized deterministic policies. The key idea is factorizing the centralized critic into individual critics via a non-linear function, where the centralized critic is decomposed as $Q_{jt}(s, \boldsymbol{\tau}, \mathbf{a}) = f_{mixer}(s, Q^1(\tau^1, a^1), \cdots, Q^N(\tau^N, a^N))$, where $f_{mixer}$ is a mixer network that encodes state information and learns the weights of individual critics, and $Q^i(\tau^i, a^i)$ is the individual critic of Agent $i$. FACMAC considers two value factorization methods: (a) monotonic factorization (*mono*), which constrains the parameters of

$f_{mixer}$ to enforce $\partial Q_{jt}/\partial Q^i \geq 0$, as in QMIX (Rashid et al., 2020); and (b) non-monotonic factorization (*non-mono*), which removes the constraint from (a) to increase representational capacity. In addition to these two methods (c) linear factorization (Sunehag et al., 2017) (*vdn*) was also considered. It was shown that *non-mono* and *mono* approaches are complementary: *non-mono* performs better on tasks requiring greater representational capacity, while *mono* outperforms on others. However, *vdn* underperforms both and performs drastically worse, particularly in multi-agent Mujoco environment.

The centralized but factored critic, parameterized by $\boldsymbol{\theta_Q}$, is trained to minimize the temporal-difference (TD) error:

$$\mathcal{L}_{\text{FACMAC}}(\boldsymbol{\theta_Q}) = \mathbb{E}_{\mathcal{D}} \left[ \left( y^{jt} - Q_{jt}(s, \boldsymbol{\tau}, \boldsymbol{a}; \boldsymbol{\theta_Q}) \right)^2 \right],$$
$$\text{where } y^{jt} = r + \gamma Q_{jt}(s', \boldsymbol{\tau'}, \boldsymbol{\pi}(\boldsymbol{\tau'}); \boldsymbol{\theta_Q^-}), \quad (1)$$

$\mathcal{D}$ is the replay buffer, $\boldsymbol{\pi}$ is the deterministic joint policy, and $\boldsymbol{\theta_Q'}$ are the parameters of the target critic networks. Based on this centralized critic, the decentralized policies, parameterized by $\boldsymbol{\theta_\pi}$, are trained using the deterministic policy gradient (Silver et al., 2014) where the loss function is written as $\mathcal{L}_{\text{FACMAC}}(\boldsymbol{\theta_\pi}) =$

$$\mathbb{E}_{\mathcal{D}} \left[ -Q_{jt}(s, \boldsymbol{\tau}, \pi^1(\tau^1; \theta_\pi^1), \cdots, \pi^N(\tau^N; \theta_\pi^N)) \right], \quad (2)$$

where $\theta_\pi^i$ is Agent $i$'s actor parameter.

A main benefit of FACMAC over prior works such as MAD-DPG (Lowe et al., 2017) is leveraging non-linear factorization. However, in offline settings, value factorization techniques remain underexplored—recent research still relies on linear factorization (Yang et al., 2021; Wang et al., 2024) or omits it by using a non-factored critic. In this paper, to investigate value factorization for offline settings, we adopt FACMAC to build upon the recent achievements.

### 2.2. Offline RL

Offline RL aims to learn a policy that maximizes the expected return from a fixed dataset $\mathcal{D}$, consisting of trajectories generated by arbitrary behavior policies, without additional environment interactions (Kumar et al., 2020; Kostrikov et al., 2022; Kim et al., 2024a). This inability to generate additional data introduces a major challenge in offline RL—*extrapolation error* in policy evaluation, i.e., the target value in the Bellman equation becomes inaccurate if actions from the learning policy are not included in the dataset. This extrapolation error often leads to overestimation and degrades performance. To address this challenge, several approaches have been proposed, such as behavior-constrained policy optimization, which regularizes the policy to stay close to the behavior policy (Wu et al., 2019; Fujimoto & Gu, 2021), and conservative value estimation (Kumar et al., 2020; Kostrikov et al., 2021), which

lower-bounds the value function by penalizing the values of unseen actions. These approaches have demonstrated effectiveness; however, they often require significant modifications or additional components to existing RL algorithms, such as generative models. These additions introduce more hyperparameters, which can affect performance and impact stability and reproducibility. This necessitates an offline RL algorithm with minimal changes from the current online RL algorithms. As an example, (Fujimoto & Gu, 2021) proposed TD3 with Behavior Cloning (TD3+BC), which adds a BC loss function as a regularizer on top of the TD3 loss (Fujimoto et al., 2018) to update the policy. This encourages the action generated by the learning policy to be the same as the action in the dataset. The policy of TD3+BC is updated to minimize the following loss function.

$$\mathcal{L}_{\text{TD3+BC}}(\pi) = \mathbb{E}_{\mathcal{D}}\Big[ - \underbrace{\boldsymbol{w}\, Q(s, \pi(s))}_{TD3} + \underbrace{(\pi(s) - a)^2}_{BC}\Big],$$
$$\text{where} \quad \boldsymbol{w} = \frac{\alpha}{\frac{1}{N}\sum |Q(s, \pi(s))|} \quad (3)$$

where $\alpha$ is a hyperparameter that controls the balance between RL and BC. TD3+BC has been widely used due to its simplicity and performance comparable to SOTA algorithms (Kim et al., 2024b).

### 2.3. Offline Multi-Agent RL

A natural approach to offline multi-agent RL is to extend the previously mentioned single-agent offline RL to the multi-agent setting. However, naively extending these methods has been shown to be insufficient (Pan et al., 2022; Shao et al., 2024). For example, naive conservative value estimation for the joint policy can lead to excessively large penalties, significantly degrading performance as the number of agents increases (Shao et al., 2024). Furthermore, the problem of extrapolation error worsens as the size of the action space grows with the number of agents (Yang et al., 2021), making it essential to balance conservatism and regularization with extrapolation error. To mitigate this problem, CFCQL (Shao et al., 2024) introduces counterfactual regularization to each agent individually to alleviate excessive conservatism. OMAR (Pan et al., 2022) introduces zeroth-order optimization to address policies getting stuck in poor local optima under conservative value function training. Note that both OMAR and CFCQL are based on conservative value estimation, which penalizes value functions for unseen action spaces, thereby forcing them to underestimate. OMIGA (Wang et al., 2024) introduces a framework that converts global regularization into implicit local-level regularization using linear value decomposition, while employing in-sampling techniques to avoid querying unseen actions. MADIFF (Zhu et al., 2023) presents a diffusion-based generative approach for multi-agent systems, enabling decentralized policy development through
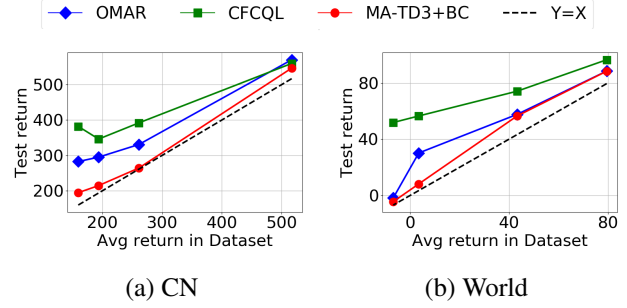


*Figure 1.* Test return with respect to dataset quality (average return of the dataset) in the Cooperative Navigation (CN) and World environments, reported in Shao et al. (2024). Performance close to $y = x$ indicates a strong dependency on data quality.

effective teammate modeling.

Although the aforementioned offline multi-agent RL algorithms perform reasonably, they still require significant modifications from existing multi-agent RL algorithms, as described in single-agent settings in Sec. 2.2, with even more adjustments needed due to the presence of multiple agents. Additionally, they have not leveraged the currently existing value factorization techniques, which are central to the recent success of multi-agent RL. Therefore, inspired by (Fujimoto & Gu, 2021), we aim to develop a minimalist approach that minimally modifies an existing multi-agent RL algorithm, focusing on simplicity and practicality while incorporating existing multi-agent RL techniques.

## 3. Methodology

To achieve the aforementioned goal, we focus on a BC regularization-based approach. In multi-agent settings, the presence of multiple agents amplifies the issue of overestimation, as mentioned in Sec.2.3. Consequently, the BC regularization-based approach is susceptible to either over-regularization or unstable critic learning, which often stems from insufficient regularization. We propose a simple yet effective regularization method called **B**ehavior **C**loning with **C**ritic **C**lipping (**B3C**), which prevents over-regularization while ensuring stability, allowing a focus on the RL objective and thereby significantly improving performance. Additionally, we explore the application of value factorization techniques with B3C, which remains understudied.

### 3.1. Revisiting Behavior Cloning: Over-regularization

We begin by revisiting BC regularization in offline multi-agent RL. Building on the success of BC regularization in the single-agent setting, recent studies on offline multi-agent RL (Pan et al., 2022; Shao et al., 2024) have introduced an extension of TD3+BC, called MA-TD3+BC, as a baseline

by adapting Eq. 3 to the multi-agent setting. Although MA-TD3+BC achieves reasonable performance, it is less effective than the proposed methods, particularly when applied to low-quality datasets. Fig. 1 illustrates the test returns of OMAR (Pan et al., 2022), CFCQL (Shao et al., 2024), and MA-TD3+BC with respect to the average return of the dataset, reported in (Shao et al., 2024). MA-TD3+BC achieves performance close to dataset-level quality, but with a small margin. In contrast, OMAR and CFCQL show significant improvements over dataset quality, particularly for low-quality datasets. This observation, that MA-TD3+BC heavily depends on data quality, indicates that it is over-regularized by behavior cloning.

One might ask: what if we place more weight on the RL objective in Eq. 3, which can be achieved by increasing $\alpha$ in Eq. 3, to reduce dependence on data quality? We observed that increasing $\alpha$ sometimes leads to better performance compared to prior work; however, it often fails to converge. The critic is updated unstably and even frequently diverges due to extrapolation error, as we will discuss in Sec. 4.3.1. Therefore, to achieve stable performance, $\alpha$ must be kept low, which results in over-regularization and limits performance based on dataset quality.

### 3.2. Critic Clipping: Alleviate Overestimation

In order to simultaneously avoid unstable critic training and over-regularization, we propose a simple technique called Critic Clipping (CC) for training the centralized critic, which integrates with BC regularization for training policies. CC clips the critic target value when calculating the target value during policy evaluation, which can lead to overestimation based on the maximum return in the dataset. Unlike conservative value estimation, which requires additional sampling or modeling of the behavior policy and introduces underestimation bias, CC is easy to implement and proves effective by limiting overestimation to below a certain threshold. Here, the reason for using the maximum return rather than the average return as the clipping standard is as follows: The learned RL policy is expected to outperform the behavior policy, and therefore, the expected return, which the critic approximates, is expected to be larger than the average return in the dataset. For simplicity, we use the maximum return, representing the best episode in the dataset. To allow further flexibility, we can use a scaled value of the maximum return for the clipping value. Note that we observe using the maximum return is sufficient empirically, as we will discuss in Sec.4.3.2. The corresponding loss function of the critic is given by $\mathcal{L}_{CC}(\boldsymbol{\phi}, \psi) =$

$$\mathbb{E}_{(s,\boldsymbol{\tau},\boldsymbol{a},r,s',\boldsymbol{\tau}')\sim\mathcal{D}}\left[\left(y^{jt} - Q_{jt}(s,\boldsymbol{\tau},\boldsymbol{a};\boldsymbol{\phi})\right)^2\right], \quad \text{where}$$

$$y^{jt} = r + \gamma \underbrace{\text{Max}\Big[Q_{jt}(s',\boldsymbol{\tau}',\boldsymbol{\pi}(\boldsymbol{\tau}';\boldsymbol{\theta}^-);\boldsymbol{\phi}^-), \ R^*\Big]}_{\text{Critic Clipping}}$$

$$\text{and} \ \ R^* = M \times \max_d \sum_{t=1}^{T} r_{d,t}. \tag{4}$$

Here, $M$ is a scalar value, and $\max_d \sum_{t=1}^{T} r_{d,t}$ represents the maximum return in the dataset where $r_{d,t}$ is the reward at time-step $t$ of $d$-th episode in the dataset. For simplicity, we use $M = 1$ in most cases during the experiments. We observed that this straightforward use of $M$ is sufficient in most cases, reducing the burden of hyperparameter tuning. We provide an ablation study on this value in Sec. 4.3.2.

### 3.3. B3C with Value Factorization

We integrate CC with BC regularization, introducing the proposed method, **B3C**. The main advantage of CC is its flexibility, which reduces sensitivity to BC regularization and allows for maximizing the weight of the RL objective. This ultimately leads to improved performance. Additionally, we adopt FACMAC, which uses a centralized but factored critic with non-linear value factorization. Accordingly, the policy loss function of **FACMAC+B3C** is given by $\mathcal{L}_{\text{FACMAC+B3C}}(\boldsymbol{\theta_\pi}) =$

$$\mathbb{E}_{\mathcal{D}}\Big[ - \boldsymbol{w}\underbrace{Q_{jt}(s,\boldsymbol{\tau},\pi^1(\tau^1),\cdots,\pi^N(\tau^N))}_{FACMAC} \ + \tag{5}$$

$$\beta \sum_{i=1}^{N}\underbrace{(\pi^i(\tau^i) - a^i)^2}_{BC}\Big], \text{where } \boldsymbol{w} = \frac{\alpha}{\frac{1}{N}\sum|Q_{jt}(s,\boldsymbol{\tau},\boldsymbol{a})|},$$

$\boldsymbol{\pi} = (\pi^1,\cdots,\pi^N)$ and $\boldsymbol{a} = (a^1,\cdots,a^N)$ represent the joint policy and the joint action, respectively. Here, $\boldsymbol{\theta_\pi}$ denotes the parameters of the joint policy, and $\mathcal{D}$ represents the given offline dataset. We normalize the Q-value to ensure robustness to scale, following Fujimoto & Gu (2021). Note that there are two important hyperparameters, $\alpha$ and $\beta$, which balance the trade-off between RL and BC. We will refer to $\alpha$ and $\beta$ as the RL coefficient and the BC coefficient, respectively. A key difference from Fujimoto & Gu (2021) is the introduction of the BC coefficient, $\beta$, alongside the RL coefficient, $\alpha$. For hyperparameter tuning, we recommend practitioners start by fixing $\beta = 1$ and tuning $\alpha$, following Fujimoto & Gu (2021). If over-regularization persists, $\beta$ can then be adjusted, which has been shown to be effective in multi-agent particle environments.

**Empirical Observation regarding value factorization:** We consider three value factorization methods discussed in Sec.2.1: *non-mono*, *mono*, and *vdn*. We empirically observe that *non-mono* outperforms *mono* in most cases, which

differs from the findings in the online setting discussed in Sec. 2.1. We argue that limiting representational capacity, which is inherent to monotonic factorization, often leads to reduced performance in offline settings. We provide an ablation study for this in Sec. 4.3.3.

As summarized above, we update the joint policy to maximize the critic with BC regularization by minimizing Eq. 5, and we update the centralized critic to minimize the TD error with clipping, as shown in Eq. 4. We refer to the proposed algorithm—FACMAC integrated with **B**ehavior **C**loning regularization and **C**ritic **C**lipping—as FACMAC+B3C. Additionally, we also apply B3C to MA-TD3, referred to as MA-TD3+B3C.

# 4. Experiments

## 4.1. Experimental Setup

**Environments and Offline Dataset.** We consider various multi-agent environments categorized by their types, the number of agents, and levels of partial observability. Additionally, we use offline datasets generated by different MARL algorithms, covering various levels.

(1) Three multi-agent particle environments (provided by Pan et al. (2022) and generated using MA-TD3 (Ackermann et al., 2019)): (a) Cooperative Navigation (CN), where three agents cover three landmarks without colliding, requiring coordination and teamwork; (b) Predator-Prey (PP), where three slower predators cooperate to capture a faster, pretrained prey; and (c) World, where four predators attempt to catch two prey that aim to eat food while strategically using forest hiding spots for protection. For each task, evaluation is performed using four datasets of varying quality: expert, medium, medium-replay, and random.

(2) Three fully observable multi-agent MuJoCo environments (provided in Wang et al. (2024) and generated using HAPPO (Kuba et al., 2021)): We consider three tasks including 3-Agent Hopper (HC), 2-Agent Ant (Ant), and 6-Agent HalfCheetah (HC). For each task, evaluation is performed using four datasets of varying quality: expert, medium, medium-replay, and medium-expert.

(3) Two partially observable multi-agent MuJoCo environments: 6-Agent HalfCheetah and 5-Agent Swimmer, each with differing degrees of partial observability and performance levels. We generated the offline datasets using ADER (Kim & Sung, 2023) with varying levels of performance: expert, medium-1, medium-2, and their combinations. In both environments, each agent can observe its $K$ nearest neighbors. Here, $K$ determines the degree of partial observability. We consider $K = 0, 1$ for HC, which is denoted as HC-k$K$, and $K = 0$ for Swimmer (SW).

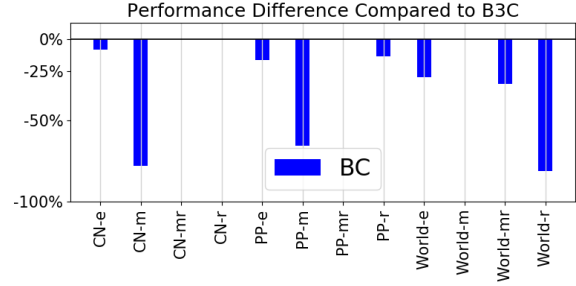**Baselines.** We compare the proposed method against vari-



*Figure 2.* Performance difference between the worst-performing seed of MA-TD3+BC and that of MA-TD3+B3C, with negative values indicating that MA-TD3+BC performs worse than MA-TD3+B3C. A value of $-5\%$ represents that BC performs 5% worse than B3C.

ous baselines for each environment. For the multi-agent particle environments, we use MA-TD3+BC (Pan et al., 2022), OMAR (Pan et al., 2022), MADIFF (Zhu et al., 2023), and CFCQL (Shao et al., 2024) as baselines. For both multi-agent Mujoco environments, we use the multi-agent versions of BCQ (Fujimoto et al., 2019) and CQL (Kumar et al., 2020), as well as IQL (Yang et al., 2021), OMAR (Pan et al., 2022), CFCQL (Shao et al., 2024), and OMIGA (Wang et al., 2024).

**Implementation and Hyperparameters.** To ensure a fair comparison, we implement the proposed method using the released CFCQL (Shao et al., 2024) code for the multi-agent particle environments and the released OMIGA (Wang et al., 2024) code for the multi-agent Mujoco environments. The RL and BC coefficients in Eq. 5 are critical hyperparameters in our approach. As mentioned in Sec. 3.3, we initially set $\beta = 1$ and tune $\alpha$ for each task, as done in Fujimoto & Gu (2021). We observe that $\beta = 1$ is sufficient for multi-agent Mujoco, but tuning $\beta$ improves performance in the multi-agent particle environments. For another hyperparameter, the scalar clipping value $M$ in Eq. 4, we use $M = 1$ for simplicity, except for the Ant task with the medium-replay dataset, where a lower $M$ significantly improves performance. An ablation study on this value is provided in Sec. 4.3.2. The detailed hyperparameters used are listed in Appendix A.

## 4.2. Numerical Results

**Multi-agent particle environments.** The results are shown in Table 1. As discussed in Sec. 3.1, MA-TD3+BC from prior works suffers from over-regularization, which is evident from its performance on medium-replay and random datasets, where it significantly underperforms compared to the baselines. Thus, we first further tune the coefficients of RL and BC, which is referred to as MA-TD3+BC (ours) in Table 1. This simple tuned version substantially im-

| Task-Dataset | Avg R | Max R | OMAR | CFCQL | MADIFF | MA-TD3+BC | MA-TD3+BC (ours) | MA-TD3+B3C (ours) |
|---|---|---|---|---|---|---|---|---|
| CN-e | 100.0 | 163.1 | **114.9±2.6** | 112.0±4.0 | 95.0±5.3 | 108.3±3.3 | 100.6±8.3 | 102.3±6.7 |
| CN-mr | 9.5 | 123.2 | 37.9±12.3 | **52.2±9.6** | 30.3±2.5 | 15.4±5.6 | **53.8±4.8** | **53.8±4.8** |
| CN-m | 27.8 | 145.0 | 47.9±18.9 | **65.0±10.2** | 64.9±7.7 | 29.3±4.8 | 55.9±25.8 | **63.5±9.7** |
| CN-r | 0.0 | 103.7 | 6.9±3.1 | 62.2±8.1 | 6.9±3.1 | 9.8±4.9 | **72.6±7.0** | **73.3±6.6** |
| PP-e | 100.0 | 275.7 | 116.2±19.8 | 118.2±13.1 | 120.9±14.6 | 115.2±12.5 | 108.1±22.9 | **124.6±21.2** |
| PP-mr | 9.6 | 128.0 | 47.1±15.3 | 71.1±6.0 | 62.3±9.2 | 28.7±20.9 | **75.2±6.8** | **76.5±8.6** |
| PP-m | 48.9 | 196.2 | 66.7±23.2 | 68.5±21.8 | 77.2±10.4 | 65.1±29.5 | 91.3±35.7 | **104.3±11.9** |
| PP-r | 0.0 | 61.4 | 11.1±2.8 | 78.5±15.6 | 3.2±4.0 | 5.7±3.5 | **105.0±16.7** | **105.5±10.0** |
| World-e | 100.0 | 270.7 | 110.4±25.7 | 119.7±26.4 | **122.6±14.4** | 110.3±21.3 | 120.7±17.9 | **124.3±5.1** |
| World-mr | 12.0 | 138.7 | 42.9±19.5 | **73.4±23.2** | 57.1±10.7 | 17.4±8.1 | **75.2±23.1** | **75.8±10.9** |
| World-m | 58.1 | 206.8 | 74.6±11.5 | 93.8±31.8 | **123.5±4.5** | 73.4±9.3 | 119.3±18.2 | 119.7±16.7 |
| World-r | 0.0 | 63.1 | 5.9±5.2 | 68.0±20.8 | 2.0±3.0 | 2.8±5.5 | 98.3±49.7 | **101.0±16.8** |



(a) Ant2x4-return

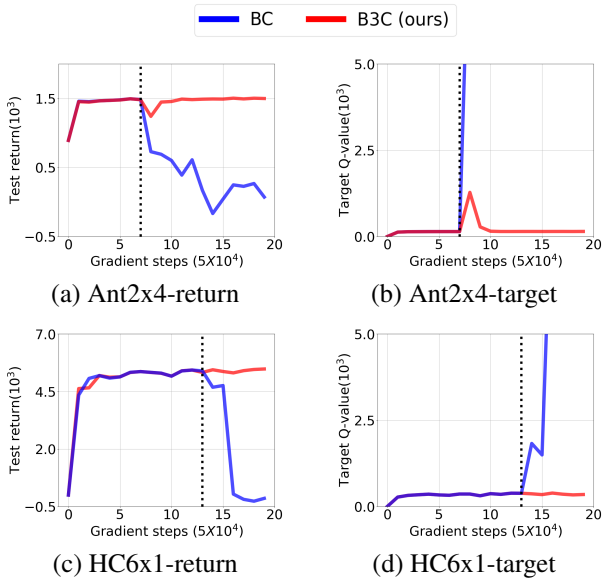(b) Ant2x4-target

(c) HC6x1-return

(d) HC6x1-target

*Figure 3.* Test return and target value during policy evaluation in the training of BC (blue) and B3C (red). The black dotted line indicates the moment when the target value starts to diverge.

proves performance and even outperforms the baselines in some cases. However, this simple hyperparameter tuning is insufficient: *BC regularization alone is often unstable*. We claim that the proposed CC can address this instability. For this, Fig. 2 illustrates the performance difference between the worst-performing seed among 5 seeds of MA-TD3+BC and that of MA-TD3+B3C, where $-5\%$ indicates that MA-TD3+BC performs $5\%$ worse than MA-TD3+B3C. As shown in Fig. 2, using BC alone results in a significant performance drop in the worst-performing seed compared to **B3C** due to its instability, demonstrating that **B3C**, which incorporates critic clipping, stabilizes learning.

**Multi-agent Mujoco.** In this environment, we evaluate FACMAC+B3C, which incorporates non-monotonic

value factorization with B3C. For fully observable cases, we use the performance results of BCQ-MA, CQL-MA, IQL, OMAR, and OMIGA as reported in (Wang et al., 2024), while we generate results for CFCQL and FACMAC+B3C. For partially observable cases, we generate all results using their public code. As shown in Table 2, the proposed method, FACMAC+B3C, outperforms the baselines in most cases. Notably, across all datasets for the HalfCheetah environment and several datasets for the Hopper and Ant environments, FACMAC+B3C even outperforms the maximum return of the datasets.

### 4.3. Analysis

To gain deeper insights into the proposed method, B3C, and value factorization techniques for offline learning, we present several ablation studies in the following sections.

#### 4.3.1. CRITIC CLIPPING: B3C VERSUS BC

**Critic clipping alleviates overestimation** by directly constraining the target value in policy evaluation. As mentioned earlier, BC regularization alone often makes learning unstable, and even the critic diverges. We have already shown instability problem in Sec. 4.2 through the worst performance of MA-TD3+BC in the multi-agent particle environments. Here, we provide further analysis of how critic clipping addresses this issue for better understanding in multi-agent Mujoco environments. Fig. 3 presents the test return and the target value, $y^{jt}$ from Eq. 4, during the training of FAC-MAC+BC and FACMAC+B3C. The results indicate that FACMAC+BC, which relies solely on BC regularization, results in diverging value estimates. When this divergence occurs, marked by the dotted line in Fig. 3, the performance correspondingly degrades. However, B3C effectively mitigates this divergence—it occasionally tends to overestimate but either stabilizes quickly or avoids divergence entirely, ensuring more consistent performance. Note that this phenomenon does not always occur but is more likely when
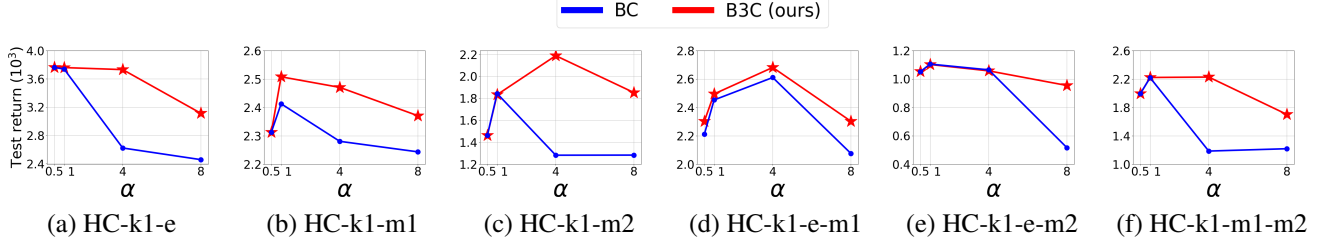
Figure 4. Performance with respect to the RL coefficient $\alpha$ in the Halfcheetah environment.

*Table 2.* Results of FACMAC+B3C and the considered baselines across six tasks and several datasets in multi-agent Mujoco. Maximum and average returns are provided for each dataset. Dataset qualities are abbreviated as follows: 'expert' as 'e,' 'medium1' as 'm1,' and 'medium2' as 'm2.' Combinations of datasets are denoted as 'e-m1' for 'e' and 'm1,' and 'm1-m2' for 'm1' and 'm2.' Boldface numbers indicate the highest score or a comparable one among the algorithms.

| Task-Dataset | Avg R | Max R | BCQ-MA | CQL-MA | IQL | OMAR | CFCQL | OMIGA | FACMAC+B3C (ours) |
|---|---|---|---|---|---|---|---|---|---|
| Fully observable multi-agent MuJoCo provided in OMIGA (Wang et al., 2024) and generated by HAPPO (Kuba et al., 2021) |
| Hop-e | 2452.0 | 3762.7 | 77.9 | 159.1 | 755.7 | 2.4 | 718.5 | 859.6 | **3619.7±1.6** |
| Hop-m | 723.6 | 2776.5 | 45.6 | 401.2 | 501.8 | 21.3 | 674.2 | 1189.3 | **3242.7±129.1** |
| Hop-mr | 746.4 | 2801.2 | 26.5 | 31.4 | 195.4 | 3.3 | **1380.2** | 774.2 | 736.8±469.4 |
| Hop-me | 1190.6 | 3762.7 | 54.3 | 64.8 | 355.4 | 1.4 | 383.0 | 709.0 | **3328.0±369.8** |
| Ant-e | 2055.1 | 2124.2 | 1317.7 | 1042.4 | 2050.0 | 313.5 | 1756.1 | 2055.5 | **2162.8±46.0** |
| Ant-m | 1418.7 | 1473.9 | 1059.6 | 533.9 | 1412.4 | -1710.0 | 1159.6 | 1418.4 | **1516.5±14.8** |
| Ant-mr | 1029.5 | 1517.1 | 950.8 | 234.6 | 1016.7 | -2014.2 | 1052.9 | 1105.1 | **1259.8±302.4** |
| Ant-me | 1736.9 | 2124.2 | 1020.9 | 800.2 | 1590.2 | -2992.8 | 613.2 | 1720.3 | **2077.6±194.2** |
| HC-e | 2785.1 | 3866.1 | 2992.7 | 1189.5 | 2955.9 | -206.7 | 4999.2 | 3383.6 | **5403.5±169.7** |
| HC-m | 1415.7 | 2113.5 | 2590.5 | 1011.4 | 2549.3 | -265.7 | 4345.0 | 3608.1 | **4756.6±56.3** |
| HC-mr | 655.8 | 2132.6 | -333.6 | 1998.7 | 1922.4 | -235.4 | 3655.3 | 2504.7 | **4602.6±150.2** |
| HC-me | 2105.4 | 3866.1 | 3543.7 | 1194.2 | 2834.0 | -253.8 | 5030.9 | 2948.5 | **5413.7±99.4** |
| Partial observable multi-agent MuJoCo generated by ADER (Kim & Sung, 2023) |
| HC-k0-e | 1394.3 | 1520.8 | - | - | - | 197.2 | 750.7 | 1390.1 | **1396.8±4.5** |
| HC-k0-m1 | 1103.3 | 1332.9 | - | - | - | 189.6 | 443.7 | 1106.3 | **1141.6±18.9** |
| HC-k0-m2 | 840.8 | 1157.4 | - | - | - | 839.1 | 766.5 | 847.6 | **1195.0±51.9** |
| HC-k0-e-m1 | 1252.3 | 1520.8 | - | - | - | 136.4 | 542.6 | 1199.5 | **1307.1±27.3** |
| HC-k0-e-m2 | 1121.7 | 1520.8 | - | - | - | 299.2 | 398.4 | 1097.9 | **1230.0±40.5** |
| HC-k0-m1-m2 | 976.2 | 1332.9 | - | - | - | 709.1 | 1186.1 | 1027.3 | **1210.1±22.8** |
| HC-k1-e | 3766.0 | 3863.3 | - | - | - | 3232.6 | 3390.2 | **3760.8** | 3760.5±24.2 |
| HC-k1-m1 | 1976.2 | 2350.0 | - | - | - | 2312.4 | 2356.0 | 2017.3 | **2508.1±55.7** |
| HC-k1-m2 | 1223.9 | 1758.4 | - | - | - | 1108.9 | 1440.6 | 1196.5 | **2187.8±66.7** |
| HC-k1-e-m1 | 2873.2 | 3863.3 | - | - | - | 2296.4 | 1949.6 | 2112.3 | **2682.0±175.4** |
| HC-k1-e-m2 | 2486.0 | 3863.3 | - | - | - | 524.8 | 864.7 | **1088.2** | 1102.8±349.8 |
| HC-k1-m1-m2 | 1591.6 | 2237.2 | - | - | - | 1410.6 | 1862.2 | 1633.5 | **2222.6±84.1** |
| Sw-e | 430.9 | 438.9 | - | - | - | 395.3 | 403.3 | **430.7** | 430.3±2.6 |
| Sw-m1 | 290.4 | 306.3 | - | - | - | 268.4 | 277.1 | **288.3** | 289.5±1.1 |
| Sw-m2 | 142.7 | 158.5 | - | - | - | 97.9 | 130.0 | 142.6 | **155.3±1.7** |
| Sw-e-m1 | 360.6 | 438.9 | - | - | - | 216.5 | **292.1** | 261.4 | 297.1±53.8 |
| Sw-e-m2 | 286.8 | 438.9 | - | - | - | 113.8 | 138.6 | 148.2 | **211.9±8.4** |
| Sw-m1-m2 | 216.3 | 306.3 | - | - | - | 226.4 | **222.7** | 205.8 | 189.9±7.7 |

the RL objective is weighted more heavily than BC regularization; however, critic clipping can prevent it and enable stable learning.

**Critic clipping improves performance and robustness** by enabling a higher weight on the RL objective, overcoming over-regularization while ensuring stability. As discussed in Sec. 3.1, prior works employing BC regularization-based approaches suffered from over-regularization, which con-strained performance to the quality of the data. This issue is inevitable without the proposed critic clipping, as increasing the weight on the RL objective induces instability, making divergence more likely, as observed in the worst-performing comparison between BC and B3C shown in Fig. 2. However, CC allows for reduced BC regularization, i.e., a higher weight on the RL objective, resulting in improved perfor-mance. We investigate this by comparing FACMAC+BC and FACMAC+B3C with respect to the RL objective weight
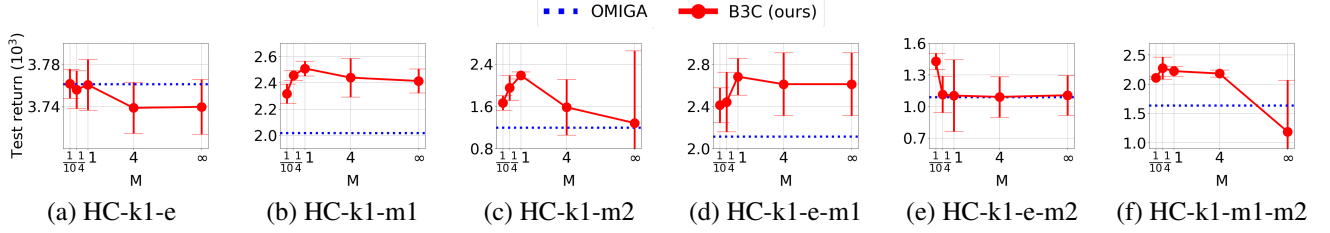
(a) HC-k1-e   (b) HC-k1-m1   (c) HC-k1-m2   (d) HC-k1-e-m1   (e) HC-k1-e-m2   (f) HC-k1-m1-m2

*Figure 5.* Ablation study on the clipping value, $M$: Performance of FACMAC+B3C with respect to $M$. $M = \infty$ (no clip) corresponds to FACMAC+BC. The performance of OMIGA is also included as a blue dotted line. Error bars represent one standard deviation.
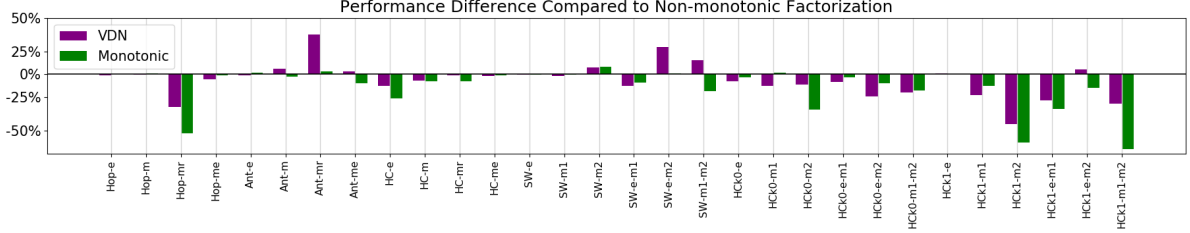


*Figure 6.* Ablation study on value factorization: Performance differences of VDN and monotonic factorization compared to non-monotonic factorization.

$\alpha$, while fixing $\beta = 1$. The corresponding results are shown in Fig. 4. In most cases, FACMAC+BC performs worse when $\alpha > 1$, whereas FACMAC+B3C performs better or equal when $\alpha = 4$ in most cases. Additionally, although FACMAC+B3C experiences a performance drop when $\alpha = 8$, it remains more stable than FACMAC+BC, demonstrating the robustness of B3C across different values of $\alpha$. Additionally, we provide the averaged performance difference between BC and B3C in multi-agent Mujoco in the Appendix B.

### 4.3.2. CRITIC CLIPPING: CLIPPING VALUE

CC clips the Q-function in the target value by the scaled maximum return in the dataset. This introduces a hyperparameter—the clipping value, $M$ in Eq. 4. As discussed in Sec. 3.2, $M = 1$ is sufficient since the maximum return in the dataset is already high enough to handle overestimation. We conduct an ablation study regarding the clipping value. Fig. 5 shows the performance of FACMAC+B3C by varying $M$ across $1/10$, $1/4$, $1$, and $4$. Additionally, we include the performance of FACMAC+BC, which does not use critic clipping, as $M = \infty$, and OMIGA (Wang et al., 2024) as a blue dotted line. It is seen that $M = 1$ achieves the best performance in all cases except for the HC-k1-e-m2 environment. While extremely small or large values of $M$ perform worse than $M = 1$, FACMAC+B3C with any of the tested $M$ values still outperforms both FACMAC+BC and the SOTA algorithm, OMIGA. In the HC-k1-e-m2 environment, where FACMAC+B3C performs suboptimally compared to the dataset quality, reducing $M$ improves performance. Over-

all, $M = 1$ is a reasonable choice, providing the best performance while minimizing the burden of hyperparameter tuning.

### 4.3.3. VALUE FACTORIZATION

We here provide an ablation study of the empirical observation described in Sec. 3.3 that non-monotonic factorization is better than monotonic factorization in offline settings. Fig. 6 shows the performance differences of VDN and monotonic factorization compared to non-monotonic factorization. It is seen that `non-mono` generally performs better than both `mono` and `vdn`, though linear factorization shows better performance in certain cases. However, `mono` performs poorly in offline settings, which differs from the online settings (Peng et al., 2021).

## 5. Conclusion

In this paper, we have proposed a simple regularization technique named behavior cloning with critic clipping (B3C) for offline multi-agent RL that addresses overestimation and over-regularization while ensuring stability. In addition, we have investigated existing value factorization techniques in offline settings, which are understudied, providing the observation that non-monotonic factorization performs better than monotonic factorization. Numerical results show that B3C with non-monotonic factorization yields outstanding performance across the considered environments with diverse dataset levels. We have also provided several analyses and ablation studies to illustrate how the proposed method operates and affects learning.

8

## Impact Statement

Potential societal impacts of this research may arise in areas where multi-agent systems play a critical role, such as autonomous vehicles, distributed resource management, and multi-robot collaboration. Our approach contributes to safer and more reliable decision-making processes in these systems, thereby promoting ethical AI development.

## References

Ackermann, J., Gabler, V., Osa, T., and Sugiyama, M. Reducing overestimation bias in multi-agent domains using double centralized critics. *arXiv preprint arXiv:1910.01465*, 2019.

Fujimoto, S. and Gu, S. S. A minimalist approach to offline reinforcement learning. *Advances in neural information processing systems*, 34:20132–20145, 2021.

Fujimoto, S., Hoof, H., and Meger, D. Addressing function approximation error in actor-critic methods. In *International conference on machine learning*, pp. 1587–1596. PMLR, 2018.

Fujimoto, S., Meger, D., and Precup, D. Off-policy deep reinforcement learning without exploration. In *International conference on machine learning*, pp. 2052–2062. PMLR, 2019.

Jeon, J., Kim, W., Jung, W., and Sung, Y. Maser: Multi-agent reinforcement learning with subgoals generated from experience replay buffer. In *International Conference on Machine Learning*, pp. 10041–10052. PMLR, 2022.

Kim, J., Lee, S., Kim, W., and Sung, Y. Adaptive $q$-aid for conditional supervised learning in offline reinforcement learning. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024a.

Kim, J., Lee, S., Kim, W., and Sung, Y. Decision convformer: Local filtering in metaformer is sufficient for decision making. In *The Twelfth International Conference on Learning Representations*, 2024b.

Kim, W. and Sung, Y. An adaptive entropy-regularization framework for multi-agent reinforcement learning. In *International Conference on Machine Learning*, pp. 16829–16852. PMLR, 2023.

Kim, W., Jung, W., Cho, M., and Sung, Y. A variational approach to mutual information-based coordination for multi-agent reinforcement learning. In *Proceedings of the 2023 International Conference on Autonomous Agents and Multiagent Systems*, pp. 40–48, 2023.

Kostrikov, I., Fergus, R., Tompson, J., and Nachum, O. Offline reinforcement learning with fisher divergence critic regularization. In *International Conference on Machine Learning*, pp. 5774–5783. PMLR, 2021.

Kostrikov, I., Nair, A., and Levine, S. Offline reinforcement learning with implicit q-learning. In *International Conference on Learning Representations*, 2022.

Kuba, J. G., Chen, R., Wen, M., Wen, Y., Sun, F., Wang, J., and Yang, Y. Trust region policy optimisation in multi-agent reinforcement learning. *arXiv preprint arXiv:2109.11251*, 2021.

Kumar, A., Zhou, A., Tucker, G., and Levine, S. Conservative q-learning for offline reinforcement learning. *Advances in Neural Information Processing Systems*, 33:1179–1191, 2020.

Lowe, R., Wu, Y. I., Tamar, A., Harb, J., Pieter Abbeel, O., and Mordatch, I. Multi-agent actor-critic for mixed cooperative-competitive environments. *Advances in neural information processing systems*, 30, 2017.

Oliehoek, F. A. Decentralized pomdps. In *Reinforcement learning: state-of-the-art*, pp. 471–503. Springer, 2012.

Pan, L., Huang, L., Ma, T., and Xu, H. Plan better amid conservatism: Offline multi-agent reinforcement learning with actor rectification. In *International conference on machine learning*, pp. 17221–17237. PMLR, 2022.

Peng, B., Rashid, T., Schroeder de Witt, C., Kamienny, P.-A., Torr, P., Böhmer, W., and Whiteson, S. Facmac: Factored multi-agent centralised policy gradients. *Advances in Neural Information Processing Systems*, 34:12208–12221, 2021.

Rashid, T., Samvelyan, M., De Witt, C. S., Farquhar, G., Foerster, J., and Whiteson, S. Monotonic value function factorisation for deep multi-agent reinforcement learning. *Journal of Machine Learning Research*, 21(178):1–51, 2020.

Shao, J., Qu, Y., Chen, C., Zhang, H., and Ji, X. Counterfactual conservative q learning for offline multi-agent reinforcement learning. *Advances in Neural Information Processing Systems*, 36, 2024.

Silver, D., Lever, G., Heess, N., Degris, T., Wierstra, D., and Riedmiller, M. Deterministic policy gradient algorithms. In *International conference on machine learning*, pp. 387–395. Pmlr, 2014.

Sunehag, P., Lever, G., Gruslys, A., Czarnecki, W. M., Zambaldi, V., Jaderberg, M., Lanctot, M., Sonnerat, N., Leibo, J. Z., Tuyls, K., et al. Value-decomposition networks for cooperative multi-agent learning. *arXiv preprint arXiv:1706.05296*, 2017.

Tarasov, D., Kurenkov, V., Nikulin, A., and Kolesnikov, S. Revisiting the minimalist approach to offline reinforcement learning. *Advances in Neural Information Processing Systems*, 36, 2024.

Wang, X., Xu, H., Zheng, Y., and Zhan, X. Offline multi-agent reinforcement learning with implicit global-to-local value regularization. *Advances in Neural Information Processing Systems*, 36, 2024.

Wang, Y., Han, B., Wang, T., Dong, H., and Zhang, C. Dop: Off-policy multi-agent decomposed policy gradients. In *International conference on learning representations*, 2020.

Wu, Y., Tucker, G., and Nachum, O. Behavior regularized offline reinforcement learning. *arXiv preprint arXiv:1911.11361*, 2019.

Yang, Y., Ma, X., Li, C., Zheng, Z., Zhang, Q., Huang, G., Yang, J., and Zhao, Q. Believe what you see: Implicit constraint approach for offline multi-agent reinforcement learning. *Advances in Neural Information Processing Systems*, 34:10299–10312, 2021.

Zhang, T., Li, Y., Wang, C., Xie, G., and Lu, Z. Fop: Factorizing optimal joint policy of maximum-entropy multi-agent reinforcement learning. In *International conference on machine learning*, pp. 12491–12500. PMLR, 2021.

Zhu, Z., Liu, M., Mao, L., Kang, B., Xu, M., Yu, Y., Ermon, S., and Zhang, W. Madiff: Offline multi-agent learning with diffusion models. *arXiv preprint arXiv:2305.17330*, 2023.

# A. Hyperparameters

We here provide the hyperparameters used for the proposed method. We have three hyperparameters: the RL coefficient $\alpha$, the BC coefficient $\beta$, and the clipping value $M$. As discussed in the main paper, we set $M = 1$ except for one task (Ant, medium-replay dataset). For multi-agent particle environments, the values for $\alpha$ and $\beta$ are provided in Table 3. We use $M = 1$ for all cases. Next, for multi-agent Mujoco environments, we fix $\beta = 1$. The values for $\alpha$ and $M$ are provided in Table. 4.

*Table 3.* Hyperparameters in multi-agent particle environments

| CN | $\alpha$ | $\beta$ | Tag | $\alpha$ | $\beta$ | World | $\alpha$ | $\beta$ |
|---|---|---|---|---|---|---|---|---|
| expert-s0 | 8 | 1 | expert-s0 | 8 | 1 | expert-s0 | 4 | 1 |
| expert-s1 | 8 | 1 | expert-s1 | 1 | 1 | expert-s1 | 2 | 1 |
| expert-s2 | 1 | 1 | expert-s2 | 1 | 1 | expert-s2 | 4 | 1 |
| expert-s3 | 4 | 1 | expert-s3 | 1 | 1 | expert-s3 | 4 | 1 |
| expert-s4 | 1 | 1 | expert-s4 | 4 | 1 | expert-s4 | 2 | 1 |
| medium-s0 | 8 | 0.001 | medium-s0 | 8 | 1 | medium-s0 | 4 | 1 |
| medium-s1 | 16 | 0.001 | medium-s1 | 8 | 1 | medium-s1 | 4 | 1 |
| medium-s2 | 8 | 0.001 | medium-s2 | 8 | 0.001 | medium-s2 | 4 | 1 |
| medium-s3 | 8 | 1 | medium-s3 | 8 | 1 | medium-s3 | 4 | 1 |
| medium-s4 | 4 | 0.001 | medium-s4 | 8 | 1 | medium-s4 | 4 | 1 |
| medium-replay-s0 | 8 | 0.001 | medium-replay-s0 | 1 | 0.001 | medium-replay-s0 | 8 | 0.01 |
| medium-replay-s1 | 8 | 0.001 | medium-replay-s1 | 1 | 1 | medium-replay-s1 | 8 | 0.01 |
| medium-replay-s2 | 4 | 0.001 | medium-replay-s2 | 4 | 0.001 | medium-replay-s2 | 8 | 0.01 |
| medium-replay-s3 | 32 | 0.001 | medium-replay-s3 | 4 | 0.001 | medium-replay-s3 | 8 | 0.01 |
| medium-replay-s4 | 8 | 0.001 | medium-replay-s4 | 1 | 1 | medium-replay-s4 | 16 | 0.1 |
| random-s0 | 8 | 0.001 | random-s0 | 16 | 0.001 | random-s0 | 16 | 0.001 |
| random-s1 | 8 | 0.001 | random-s1 | 16 | 0.001 | random-s1 | 16 | 0.01 |
| random-s2 | 16 | 0.001 | random-s2 | 16 | 0.001 | random-s2 | 16 | 0.001 |
| random-s3 | 8 | 0.001 | random-s3 | 16 | 0.001 | random-s3 | 16 | 0.01 |
| random-s4 | 16 | 0.001 | random-s4 | 16 | 0.001 | random-s4 | 16 | 0.001 |

*Table 4.* Hyperparameters for Various Environments

| Environment | Dataset | $\alpha$ | $M$ |
|---|---|---|---|
| Hopper | expert | 16 | 1 |
| | medium | 8 | 1 |
| | medium-replay | 0.25 | 0.1 |
| | medium-expert | 1 | 1 |
| Ant | expert | 1 | 1 |
| | medium | 1 | 1 |
| | medium-replay | 1 | 1 |
| | medium-expert | 0.5 | 1 |
| HalfCheetah | expert | 16 | 1 |
| | medium | 16 | 1 |
| | medium-replay | 16 | 1 |
| | medium-expert | 16 | 1 |
| HC_obsk0 | expert | 1 | 1 |
| | medium1 | 1 | 1 |
| | medium2 | 4 | 1 |
| | expert-medium1 | 1 | 1 |
| | expert-medium2 | 4 | 1 |
| | medium1-medium2 | 4 | 1 |
| HC_obsk1 | expert | 1 | 1 |
| | medium1 | 1 | 1 |
| | medium2 | 4 | 1 |
| | expert-medium1 | 4 | 1 |
| | expert-medium2 | 1 | 1 |
| | medium1-medium2 | 4 | 1 |
| Swimmer | expert | 1 | 1 |
| | medium1 | 1 | 1 |
| | medium2 | 1 | 1 |
| | expert-medium1 | 4 | 1 |
| | expert-medium2 | 1 | 1 |
| | medium1-medium2 | 0.5 | 1 |

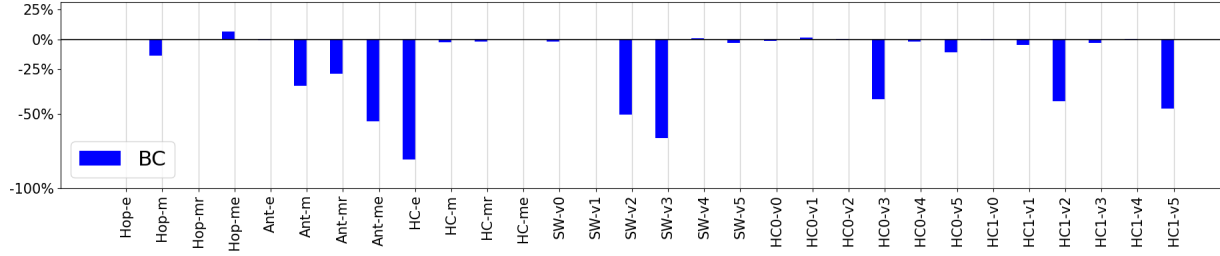## B. Comparison between B3C and BC in the Multi-agent Mujoco environment.



*Figure 7.* Performance difference FACMAC+BC compared to FACMAC+B3C, with negative values indicating that FACMAC+BC performs worse than FACMAC+B3C. A value of $-5\%$ represents that BC performs $5\%$ worse than B3C.

Fig. 7 shows the performance difference between FACMAC+BC and FACMAC+B3C. It is observed that B3C outperforms BC in terms of the average test return in the multi-agent Mujoco environments.