

# OFFLINE EQUILIBRIUM FINDING IN EXTENSIVE-FORM GAMES: DATASETS, METHODS, AND ANALYSIS

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Offline reinforcement learning (Offline RL) brings new methods to tackle real-world decision-making problems by leveraging pre-collected datasets. Despite substantial progress in single-agent scenarios, the application of offline learning to multiplayer games remains largely unexplored. Therefore, we introduce a novel paradigm *offline equilibrium finding* (Offline EF) in extensive-form games (EFGs), which aims at computing equilibrium strategies from offline datasets. The primary challenges of offline EF include i) the absence of a comprehensive dataset of EFGs for evaluation; ii) the inherent difficulties in computing an equilibrium strategy solely from an offline dataset, as equilibrium finding requires referencing all potential action profiles; and iii) the impact of dataset quality and completeness on the effectiveness of the derived strategies. To overcome these challenges, we make four main contributions in this work. First, we construct diverse datasets, encompassing a wide range of games, which form the foundation for the offline EF paradigm and serve as a basis for evaluating the performance of offline EF algorithms. Second, we design a novel framework, BOMB, which integrates the behavior cloning technique within a model-based method. BOMB can seamlessly integrate online equilibrium finding algorithms to the offline setting with minimal modifications. Third, we provide a comprehensive theoretical and empirical analysis of our BOMB framework, offering performance guarantees across various offline datasets. Finally, extensive experiments have been carried out across different games under different offline datasets, and the results not only demonstrate the superiority of our approach compared to traditional offline RL algorithms but also highlight the remarkable efficiency in computing equilibrium strategies offline.

## 1 INTRODUCTION

Extensive-form games (EFGs) provide a versatile framework for modeling the interactions between multiple players under stochastic and imperfect information settings (Nisan et al., 2007). The canonical solution concept is Nash Equilibrium (NE), where no player can increase his own utility by unilaterally deviating. There are various methods designed for solving extensive-form games, including linear programming (Shoham & Leyton-Brown, 2008), double-oracle algorithms (McMahan et al., 2003), counterfactual regret minimization (CFR) (Zinkevich et al., 2007), and policy-space response oracles (PSRO) (Lanctot et al., 2017). These methods have been successfully applied to real-world large-sale EFGs, e.g., pursuit-evasion games (Xue et al., 2021; Li et al., 2023), poker games (Brown & Sandholm, 2018; 2019; Zha et al., 2021) and Stratego (Perolat et al., 2022).

Despite the successes, existing algorithms require continuous interaction with the game environment or an accurate simulator. For example, CFR-based algorithms necessitate traversing the game tree to compute regret values, and PSRO and its variants demand simulations within the game environment to compute the best response oracle and estimate the entries in the meta-game. We call this paradigm to compute NE as “*online equilibrium finding*”. However, in many real-world applications, such as sports games (Liu et al., 2022), network intrusion detection (Khraisat et al., 2019), and automated negotiations (Kiruthika et al., 2020), the immediate interaction with the environment may be expensive and inefficient and the accurate simulator cannot be built. Therefore, offline learning is a preferred option for equilibrium finding in real-world applications.

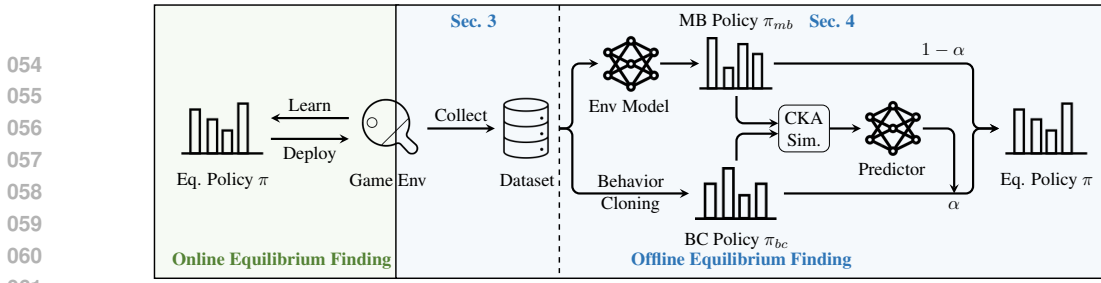


Figure 1: Comparison between online and offline equilibrium finding

Offline reinforcement learning (Offline RL) has successfully tackled numerous real-world problems by leveraging its offline learning paradigm (Levine et al., 2020). These algorithms fall into two categories: model-free and model-based. Model-free approaches, such as Best-Action Imitation Learning (BAIL) (Chen et al., 2020), directly learn optimal policies from datasets. Conversely, model-based approaches, like Model-based Offline Policy Optimization (MOPO) (Yu et al., 2020) first construct a dynamic model from the dataset, then proceed with planning. The success of these algorithms showcases the significant impact of the offline learning paradigm in advancing RL applications. In recent years, there have been several attempts to formalize the offline learning paradigm in the context of games. StarCraft II Unplugged (Mathieu et al., 2021) provides a dataset of human game-plays in a two-player zero-sum symmetric game. Some previous works (Cui & Du, 2022; Zhong et al., 2022) also explore the necessary properties of offline datasets of two-player zero-sum Markov games to successfully infer their NEs. However, these works mainly focus on solving Markov games, leaving a gap in the literature when it comes to solving extensive-form games in the offline setting. Furthermore, to our understanding, there has been no study focusing specifically on multi-player games in an offline setting. More importantly, there is a notable absence of systematic definitions and research efforts aimed at formalizing offline learning within the context of games.

To address this gap, we propose the novel *offline equilibrium finding* (Offline EF) paradigm, which computes the equilibrium strategies using offline datasets. There are several challenges for offline EF. First, the absence of comprehensive benchmarking standards complicates the evaluation and comparison of algorithm performance. Without universally accepted benchmarks, it becomes difficult to objectively measure progress within the field. Second, accurately computing or approximating equilibrium strategies solely from offline datasets is inherently difficult. Specifically, data from just two action profiles are often insufficient for determining proximity to an equilibrium strategy, as equilibrium identification requires all other potential action profiles for reference (Cui & Du, 2022). Third, the quality and completeness of data within offline datasets can significantly impact the effectiveness of derived strategies. Offline datasets fail to cover all possible game states, and this lack of comprehensive coverage can skew the algorithm’s ability to generalize from the available data.

This work presents a comprehensive investigation of Offline EF. Specifically, our contributions are fourfold: i) We curate a collection of diverse offline datasets, including random datasets, expert datasets, learning datasets, and hybrid datasets in different extensive-form games; ii) we propose the BOMB framework, which integrates behavior cloning and model-based methods along with a novel parameter estimation method and the model-based method can incorporate any online EF algorithm, e.g., CFR, into the offline context; iii) we provide a comprehensive theoretical analysis for our BOMB framework, offering performance guarantees under different datasets; and iv) we demonstrate the effectiveness of our BOMB framework in computing equilibrium strategies offline through extensive experiments on various offline datasets.

## 2 PRELIMINARIES

**Imperfect-Information Extensive-Form Games.** We use a tuple to represent an imperfect-information extensive-form game (IIEFG), i.e.,  $\mathcal{G} = (N, H, A, P, \mathcal{I}, u)$  (Shoham & Leyton-Brown, 2008). The set of players is represented by  $N = \{1, \dots, n\}$ , and  $H$  represents the set of histories (i.e., the possible action sequences). Especially, the root node of the game tree is represented by the empty sequence  $\emptyset$ , which is included in  $H$ . Every prefix of a sequence in  $H$  is also included in  $H$ . The set of terminal histories is represented by  $Z$  and belongs to  $H$ , i.e.,  $Z \subseteq H$ .  $A(h) = \{a : (h, a) \in H\}$  is the set of available actions at any non-terminal history  $h \in H \setminus Z$ .  $P$  is the player function, which maps each non-terminal history to a player, i.e.,  $P(h) \mapsto N \cup \{c\}, \forall h \in H \setminus Z$ , where  $c$  is the

“chance player” representing these stochastic events outside of the players’ controls.  $\mathcal{I}$  denotes the set of information set, which forms a partition over the set of histories where player  $i$  takes actions, such that player  $i$  cannot distinguish these histories within the same information set  $I_i$ . Every information set  $I_i \in \mathcal{I}_i$  corresponds to one decision point of player  $i$  which means that  $P(h_1) = P(h_2)$  and  $A(h_1) = A(h_2)$  for any  $h_1, h_2 \in I_i$ . For convenience, we use  $A(I_i)$  and  $P(I_i)$  to represent the action set  $A(h)$  and the player  $P(h)$  for any  $h \in I_i$ . For each player  $i$ , a utility function  $u_i : Z \rightarrow \mathbb{R}$  specifies the payoff of player  $i$  for every terminal history. The behavior strategy of player  $i$ ,  $\sigma_i$ , is a function mapping every information set of player  $i$  to a probability distribution over  $A(I_i)$ , and  $\Sigma_i$  is the set of strategies for player  $i$ . A strategy  $\sigma_i$  is defined as a pure strategy if  $\forall I_i \in \mathcal{I}$  and  $\forall a \in A(I_i), \sigma_i(I_i, a) \in \{0, 1\}$ . It is defined as a mixed strategy if  $\forall I_i \in \mathcal{I}$  and  $\forall a \in A(I_i), \sigma_i(I_i, a) \in [0, 1]$ . Moreover,  $\sigma_i$  is considered a fully mixed strategy if  $\forall I_i \in \mathcal{I}$  and  $\forall a \in A(I_i), \sigma_i(I_i, a) > 0$ . A strategy profile  $\sigma$  is a tuple of strategies, one for each player,  $(\sigma_1, \sigma_2, \dots, \sigma_n)$ , with  $\sigma_{-i}$  referring to all the strategies in  $\sigma$  except  $\sigma_i$ . Let  $\pi^\sigma(h) = \prod_{i \in N \cup \{c\}} \pi_i^\sigma(h)$  be the reaching probability of history  $h$  when all players choose actions according to  $\sigma$ , where  $\pi_i^\sigma(h)$  is the contribution of player  $i$  to this probability. Given a strategy profile  $\sigma$ , the expected value to player  $i$  is the sum of expected payoffs of these resulting terminal nodes,  $u_i(\sigma) = \sum_{z \in Z} \pi^\sigma(z) u_i(z)$ .

**Solution Concepts.** The common solution concept for IIEFGs is Nash equilibrium (NE) (Nash, 1950), where no player can increase their utility by unilaterally deviating. Formally, a strategy profile  $\sigma^*$  forms an NE if it satisfies  $u_i(\sigma^*) = \max_{\sigma'_i \in \Sigma_i} u_i(\sigma'_i, \sigma_{-i}^*), \forall i \in N$ . To measure the distance from the NE, we use the metric  $\text{NASHCONV}(\sigma) = \sum_{i \in N} \text{NASHCONV}_i(\sigma)$ , where  $\text{NASHCONV}_i(\sigma) = \max_{\sigma'_i} u_i(\sigma'_i, \sigma_{-i}) - u_i(\sigma)$ . When  $\text{NASHCONV}(\sigma) = 0$ , it indicates that  $\sigma$  is the NE. Especially, for  $n$ -player general-sum games, apart from NE, (Coarse) Correlated Equilibrium ((C)CE) is also a common solution concept. Similar to the NE, a CE is a joint mixed strategy in which no player has the incentive to deviate (Aumann, 1987). Formally, let  $S_i$  be the strategy space for player  $i$  and  $S$  be the joint strategy space. The strategy profile  $\sigma^*$  forms a CCE if it satisfies for  $\forall i \in N, s_i \in S_i, u_i(\sigma^*) \geq u_i(s_i, \sigma_{-i}^*)$  where  $\sigma_{-i}^*$  is the marginal distribution of  $\sigma^*$  on strategy space  $S_{-i}$ . Analogous to NE, the (C)CE Gap Sum is adopted to measure the distance from the (C)CE (Marris et al., 2021).

**Why Existing Methods Fail?** Offline RL focuses on learning the optimal strategies in single-agent scenarios (Levine et al., 2020), which fails to compute the equilibrium in games with offline datasets. Opponent modeling (OM) (He et al., 2016) are used to predict the opponents’ behavior strategies. However, opponent modeling algorithms aim at computing the best response strategy of one player instead of the equilibrium strategy, and they also need access to the game environment, which is not applicable to Offline EF. The widely used equilibrium finding algorithms, including no-regret methods, e.g., CFR (Zinkevich et al., 2007) and empirical game theoretic analysis (EGTA), e.g., PSRO (Lanctot et al., 2017), require the interactions with the game environments or an accurate simulator (termed as “online EF”) and cannot be applied to Offline EF. A clear comparison of existing methods is presented in Table 1 and App. B provides a detailed discussion of related methods.

**Problem Statement.** To facilitate the widespread application of game theory, we extend the offline learning framework into the extensive-form games and introduce the *offline equilibrium finding* paradigm, which focuses on learning equilibrium strategy from historical game-playing data.

**Definition 2.1** (Offline EF). Let  $\mathcal{D}$  be an offline dataset of an IIEFG  $\mathcal{G}$ , generated by an unknown behavior strategy profile  $\sigma$ . The goal of the *offline equilibrium finding* paradigm is to deduce a strategy profile  $\hat{\sigma}$  from  $\mathcal{D}$  to achieving a minimal gap from the equilibrium strategy  $\sigma^*$ . Formally,  $\hat{\sigma} = \arg \min_{\sigma' \in \Sigma} \text{GAP}(\sigma', \sigma^*)$ , where  $\text{GAP}(\cdot)$  is a metric function that measures the gap between a given strategy and the equilibrium strategy.  $\sigma$  is an  $\epsilon$ -equilibrium if  $\text{GAP}(\sigma, \sigma^*) \leq \epsilon$ .

Building on the definition of the offline EF paradigm, we can instantiate this paradigm by defining a metric for the gap from the equilibrium strategy, such as the NASHCONV for NE (Nash, 1950) and (C)CE Gap Sum for (C)CE (Aumann, 1987). While offline EF shares similarities with offline RL to some extent, it also presents distinct differences and unique challenges. Firstly, unlike offline RL, which aims to compute an optimal strategy (Levine et al., 2020), the offline EF paradigm seeks to

Methods	Work w/o env	Converge to equilibrium
Offline RL	✓	✗
OM	✗	✗
Online EF	✗	✓

Table 1: Issues of Existing Methods.

162 achieve an equilibrium strategy. This objective necessitates an iterative process to calculate the best  
 163 response strategy, introducing distinct complexities. Secondly, the offline EF paradigm involves at  
 164 least two players, making the game dynamics particularly sensitive to distribution shifts and other  
 165 uncertainties – a stark contrast to offline RL. Thirdly, while in offline RL, the data from two actions  
 166 may suffice to determine which action is better, in the offline EF paradigm, simply comparing the  
 167 data of two action tuples is inadequate for identifying which tuple is closer to an equilibrium strategy,  
 168 as equilibrium identification requires other action tuples for references (Cui & Du, 2022).

### 170 3 DATASETS

173 Datasets play a pivotal role in offline learning, however, there are no publicly available datasets  
 174 specifically tailored for the offline EF paradigm. Consequently, we outline our methods to collect  
 175 datasets at different expert levels that will serve as a basis for advancing offline EF research.

176 **Formats.** Before delving into the methods of dataset collection, it is essential  
 177 to outline the data formats of the offline EF dataset for IIEFGs. The offline dataset can be represented by  $\mathcal{D} =$   
 178  $(s_t, a_t, s_{t+1}, u_{t+1}, d_{t+1})$ . Here,  $s_t$  and  $s_{t+1}$  represent the game states at time  
 179 step  $t$  and  $t + 1$  respectively from the game-level perspective. Specifically,  $s_t$   
 180 encompasses all relevant game information at time step  $t$ , which includes  
 181 the information sets for each player and other game information  $GI$  out of the  
 182 control of players, such as the results of chance node  $(I_1^t, I_2^t, \dots, I_n^t, GI)$ , the player who needs to act ( $p^t$ ), the set of available actions for the  
 183 acting player ( $A(I_{p^t}^t)$ ), i.e.,  $s_t = (I_1^t, I_2^t, \dots, I_n^t, GI, p^t, A(I_{p^t}^t))$ . Notable,  $p^t$  may represent a chance  
 184 player  $c$  to include the game’s stochastic events outside of all players’ control. The utility for each  
 185 player at time step  $t + 1$  is represented by  $u_{t+1} = (u_1^{t+1}, u_2^{t+1}, \dots, u_n^{t+1})$ . Finally, the variable  $d_{t+1}$   
 186 indicates whether the game ends at state  $s_{t+1}$ , with a value of 1 if the game ends and 0 otherwise.

187 **Collecting Methods.** Similar to the practices in the offline RL domain, datasets in the offline EF  
 188 area must be diverse to serve as effective benchmarks for developing and evaluating algorithms.  
 189 Many benchmarks in the offline RL area, such as those discussed in (Fujimoto et al., 2019; Gul-  
 190 cehre et al., 2020), collected data from online RL training runs. Additionally, D4RL (Fu et al., 2020)  
 191 incorporates a range of dataset collection methods inspired by real-world applications, including human  
 192 demonstrations, exploratory agents, and hand-coded controllers. Inspired by these benchmarks  
 193 in the offline RL area, we propose several methods for collecting offline datasets at different expert  
 194 levels. The first one, referred to as the *random method*, involves each player adopting a uniform strat-  
 195 egy and participating repeatedly in the game to collect data as the random dataset. This approach is  
 196 motivated by the innate exploratory tendency and mimics a novice’s initial gaming experience. The  
 197 second method, the *learning-based method*, draws inspiration from the player skill improvement  
 198 process. We implement an existing equilibrium finding algorithm, such as CFR (Zinkevich et al.,  
 199 2007) or PSRO (Lanctot et al., 2017), collecting and storing intermediate game interactions to com-  
 200 pile the learning dataset. The final method, the *expert method*, capitalizes on insights gained from  
 201 observing expert players’ strategies. In this approach, each player follows an assigned equilibrium  
 202 strategy and repeatedly engages in the game to generate the expert dataset. Additionally, to enhance  
 203 realism and increase dataset diversity, we propose a hybrid approach that combines the random and  
 204 expert datasets in varying proportions, resulting in a more comprehensive collection of datasets.

211 **Statistics of Datasets.** We developed a benchmark dataset for offline EF, employing previously out-  
 212 lined collection methods on **eight** commonly used IIEFGs, as depicted in Fig. 2. In total, our offline  
 213 EF dataset comprises approximately **3.8 million** data points, occupying about **11GB** of memory.  
 214 For each game, we have generated three distinct types of datasets: Expert, random, and Learning,  
 215 each reflecting our data collection methods. The proportions of each dataset are visually detailed  
 and comprehensive statistics on the distribution of these datasets are detailed further in App. C.2.

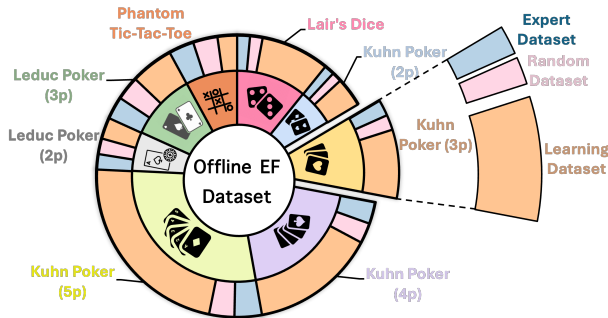


Figure 2: Dataset of Offline EF.



## 4 BOMB: FRAMEWORK AND THEORETICAL ANALYSIS

Inspiring by the success of Offline RL, there are two main directions to develop the algorithmic framework for Offline EF: i) behavior cloning (BC) (Fujimoto & Gu, 2021), which basically imitates the strategies used to collect the offline data with additional exploration, and ii) model-based methods (Yu et al., 2020; Kidambi et al., 2020), which first learns a world model from the offline dataset and then learn the strategy from the world model. However, BC may fail when the collecting policy is a random policy, which can be exploited by the opponent and model-based methods may fail when the collecting strategies are an equilibrium strategy, in which only a small portion of the game state is visited. To mitigate these issues, we propose the BOMB framework which combines Behavior clOning and Model-Based method for offline EF paradigm.

### 4.1 BOMB FRAMEWORK

**BOMB.** Alg. 1 shows the whole framework of BOMB. Given an offline dataset  $\mathcal{D}$ , we first train the policy  $\sigma_\theta$  based on the dataset  $\mathcal{D}$  using a behavior cloning (BC) technique (Line 2). Note  $\mathcal{D} = (s_t, a_t, s_{t+1}, u_{t+1}, d_{t+1})$  and  $s_t = (I_1^t, I_2^t, \dots, I_n^t, GI, p^t, A(I_{p^t}^t))$ . Since the policy network  $\sigma_\theta$  is trained to mimic the behavior strategy, only the information set  $I_{p^t}^t$  and the corresponding action  $a_t$

in  $\mathcal{D}$  are required for training. The cross-entropy loss is taken as the training loss, defined as  $\mathcal{L}_{bc} = -\mathbb{E}_{(I_{p^t}^t, a_t) \sim \mathcal{D}}[a_t \cdot \log(\sigma(I_{p^t}^t; \theta))]$ . On the other hand, inspired by model-based offline RL algorithms, where a dynamic model is trained to simulate the real environment (Kidambi et al., 2020; Yu et al., 2020; Matsushima et al., 2020), we learn an environment model  $E_{\theta_e}$  is trained based on dataset  $\mathcal{D}$  and  $E_{\theta_e}$  is used for learning the MB policy  $\sigma_{mb}$  by any online EF algorithm, e.g., PSRO (Lanctot et al., 2017), (Lines 3-4). Specifically, we use the game state  $s_t$  and corresponding action  $a_t$  as inputs, with the subsequent game state  $s_{t+1}$ , reward  $u_{t+1}$  and the termination variable  $d_{t+1}$  serving as labels. Stochastic gradient descent (SGD) is employed as the optimizer for parameter updates, and the mean squared error loss is used as the training loss, defined as  $\mathcal{L}_{env} = \mathbb{E}_{(s_t, a_t, s_{t+1}, u_{t+1}, d_{t+1}) \sim \mathcal{D}}[\text{MSE}((s_{t+1}, u_{t+1}, d_{t+1}), E(s_t, a_t; \theta_e))]$ . The final policy is obtained to combine the BC and MB policies, i.e.,  $\sigma = \alpha\sigma_\theta + (1 - \alpha)\sigma_{mb}$  where  $\alpha$  denote the weight of the BC policy (Lines 5-6). The estimation method for determining  $\alpha$  is introduced below.

**Estimation of Parameter  $\alpha$ .** Here, we introduce three estimation methods of parameter  $\alpha$ . The simplest method is randomly selecting a value from the interval  $[0, 1]$  as the parameter  $\alpha$ . Although this method can be implemented fully offline, it lacks guarantees for achieving the most effective combined strategy. The second method is the grid search method, in which we define a set of 11 candidate values for  $\alpha$ , i.e.,  $\alpha = \{0, 0.1, \dots, 1\}$ , and these values are used to configure combined policies, which are then tested in a real environment. The value of  $\alpha$  that results in the smallest gap from the equilibrium strategy is selected as optimal. This method can yield the best performance and similar techniques that determine offline parameters or fine-tune offline policies through online interactions are commonly employed in offline RL (Kalashnikov et al., 2018; Lee et al., 2022). To render our approach fully offline while still achieving optimal parameter values, we propose a learning-based method, depicted in Fig. 3. In this method, a predictor is trained to estimate  $\alpha$  based on the difference between the BC and the MB policies. We first use the grid search method to get optimal parameter values as labels. The predictor takes the centered kernel alignment (CKA) (Kornblith et al., 2019) similarity vector between the BC and the MB policies as input and

---

#### Algorithm 1 BOMB Framework

---

- 1: **Input:** an offline dataset  $\mathcal{D}$
  - 2: Train policy  $\sigma_\theta$  based on  $\mathcal{D}$  using BC technique;
  - 3: Train an environment model  $E_{\theta_e}$  based on  $\mathcal{D}$ ;
  - 4: Learn  $\sigma_{mb}$  policy using any EF algorithm on  $E_{\theta_e}$ ;
  - 5: Select  $\alpha$  using parameter estimation method;
  - 6:  $\sigma = \alpha \cdot \sigma_\theta + (1 - \alpha) \cdot \sigma_{mb}$ ;
  - 7: **Output:** Policy  $\sigma$
- 

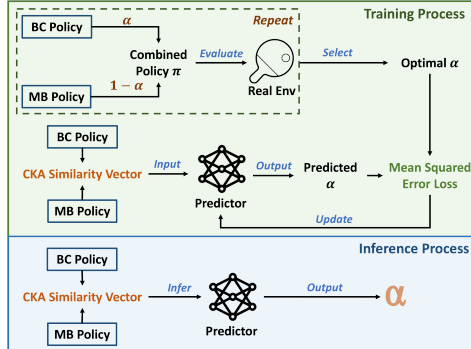


Figure 3: Learning-based estimation method.

Figure 3: Learning-based estimation method. In this method, a predictor is trained to estimate  $\alpha$  based on the difference between the BC and the MB policies. We first use the grid search method to get optimal parameter values as labels. The predictor takes the centered kernel alignment (CKA) (Kornblith et al., 2019) similarity vector between the BC and the MB policies as input and

270 outputs the estimated  $\alpha$ . The predictor can be trained in one game where access to the environment  
 271 is feasible, and reuse the predictor in similar games. Though the predictor can only provide an  
 272 approximate optimal parameter value, it requires no further online interactions once trained.  
 273

274 **Advantages of BOMB.** There are several advantages of BOMB. First, by combining the BC and  
 275 MB, BOMB can work on the datasets collected by any strategies, i.e., either random or equilibrium  
 276 strategies. Second, with the learned world model for games, BOMB can seamlessly integrate the  
 277 online EF algorithms, thus BOMB can generalize to different equilibria. For example, for computing  
 278 NE, we adapt PSRO (Lanctot et al., 2017) and Deep CFR (Brown et al., 2019) methods, referred  
 279 to as MB-PSRO and MB-CFR, respectively. Additionally, we adapt the JPSRO method (Marris  
 280 et al., 2021) (MB-JPSRO) for computing (C)CE. iii) BOMB is game-agnostic, which can learn the  
 281 game rules from the offline datasets and do not rely on the knowledge of the game, which shares the  
 282 similar advantages with MuZero (Schrittwieser et al., 2020)

## 283 4.2 THEORETICAL ANALYSIS

285 In the offline RL area, dataset coverage over the optimal pol-  
 286 icy is sufficient for offline learning (Rashidinejad et al., 2021;  
 287 Xie et al., 2021). However, we found the dataset assumption  
 288 that the dataset generated by the equilibrium strategy is not suffi-  
 289 cient for computing equilibrium strategies in an offline manner.  
 290 It can be confirmed by the counter-example illustrated in Fig. 4.  
 291 In this game, we can easily get NE strategy,  $\sigma^* = (\sigma_1^*, \sigma_2^*) =$   
 292  $(\{I_1 : a_1\}, \{I_2 : b_2\})$ . If we use this equilibrium strategy to  
 293 generate the offline dataset  $\mathcal{D}$ , then  $\mathcal{D}$  would only include the  
 294 data point  $((I_1^{t_1} = I_1, I_2^{t_1} = \emptyset, GI = \emptyset, 1, \{a_1, a_2\}), a_1, (I_1^{t_2} =$   
 295  $I_1 a_1, I_2^{t_2} = \emptyset, GI = \emptyset, -1, \emptyset), (0, 0), 1)$ . Clearly, the dataset  $\mathcal{D}$   
 296 is not sufficient for computing the NE strategy since there is no  
 297 information about Player 2. Another assumption — that the equi-  
 298 librium strategy is covered by the offline dataset — is also insufficient for the offline EF paradigm,  
 299 as we prove in App. D.1. In this section, we outline the necessary and sufficient conditions for  
 300 the coverage of an offline dataset that guarantees the convergence of our methods in IIEFGs with  
 301 perfect recall. We start by introducing two key concepts of dataset coverage: uniform coverage and  
 equilibrium coverage.

302 **Definition 4.1.** An offline dataset  $\mathcal{D}$  is said to be a *uniform coverage* of an IIEFG  $\mathcal{G}$  if and only if the  
 303 offline dataset  $\mathcal{D}$  covers all possible state-action pairs. Formally,  $(s_t, a_t, s_{t+1}, u_{t+1}, d_{t+1}), \forall s_t, a_t \in$   
 304  $A(s_t)$  and  $s_{t+1} \in T(s_t, a_t)$  where  $T$  is the transition function of game  $\mathcal{G}$ .

305 **Definition 4.2.** An offline dataset  $\mathcal{D}$  is said to be an  $\epsilon$ -*equilibrium coverage* over an IIEFG  $\mathcal{G}$  if  
 306 and only if its underlying behavior strategy  $\sigma_{\mathcal{D}}$  satisfies  $\text{GAP}(\sigma_{\mathcal{D}}, \sigma^*) < \epsilon$ , where  $\sigma_{\mathcal{D}}$  is defined  
 307 as  $\sigma_{\mathcal{D}}(s_t, a_t) = \frac{C(s_t, a_t)}{C(s_t)}$  and  $\sigma_{\mathcal{D}}(s_t, a_t) > 0$  for all  $s_t$  and  $a_t \in A(s_t)$ , with  $C(s_t, a_t)$  and  $C(s_t)$   
 308 denoting the counts of data points containing  $(s_t, a_t)$  and  $s_t$  in  $\mathcal{D}$ , respectively.

309 Building on the dataset coverage definitions previously introduced, we now discuss the conditions  
 310 under which our method achieves convergence. To facilitate this analysis, we introduce an assump-  
 311 tion about the error in training neural networks within the algorithm. All subsequent theorems are  
 312 derived under this assumption unless stated otherwise.

313 **Assumption 4.3.** The error in training neural networks within our method is assumed to be smaller  
 314 than an arbitrarily small  $\epsilon$ , provided that the dataset contains a sufficient amount of data.

315 To further support this assumption, we provided a general generalization bound for the training error  
 316 under a dataset with size  $m$  in App. D.2. Then we present our result as follows.

317 **Theorem 4.4.** Let  $\sigma_{MB(\mathcal{D})}$  be the strategy profile learned by our **model-based algorithm** based on  
 318 the offline dataset  $\mathcal{D}$  with sufficient data under Assumption 4.3. Then,  $\sigma_{MB(\mathcal{D})}$  is guaranteed to be  
 319 an  $\epsilon$ -equilibrium strategy of the IIEFG  $\mathcal{G}$  if and only if  $\mathcal{D}$  is a uniform coverage of  $\mathcal{G}$  and  $\sigma_{MB(\mathcal{D})}$  is  
 320 an  $\epsilon$ -equilibrium strategy for the trained environment model within the model-based algorithm.  
 321

322 *Sketch Proof.* According to Assumption 4.3, the error in training the environment game model based  
 323 on  $\mathcal{D}$  can be considered negligible. Consequently, the trained environment game model is identical

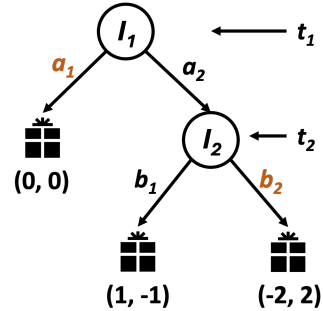


Figure 4: Game example.

to the original game  $\mathcal{G}$ , as the dataset  $\mathcal{D}$  provides full coverage of all state transitions. Therefore, if  $\sigma_{MB(\mathcal{D})}$  is an  $\epsilon$ -equilibrium strategy for the trained environment game model, it is also an  $\epsilon$ -equilibrium strategy for the original game  $\mathcal{G}$ . Any slight violation of these conditions would invalidate the convergence result. A complete proof is provided in App. D.1.  $\square$

**Theorem 4.5.** *Let  $\sigma_{BC(\mathcal{D})}$  be the strategy profile learned by our **behavior cloning algorithm** based on the offline dataset  $\mathcal{D}$  with sufficient data under Assumption 4.3. Then  $\sigma_{BC(\mathcal{D})}$  is guaranteed to be an  $\epsilon$ -equilibrium strategy of IIEFG  $\mathcal{G}$  if and only if the offline dataset  $\mathcal{D}$  is an  $\epsilon$ -equilibrium coverage of the IIEFG  $\mathcal{G}$ .*

*Sketch Proof.* According to Assumption 4.3, the error in training the behavior cloning strategy  $\sigma_{BC(\mathcal{D})}$  from the dataset  $\mathcal{D}$  is negligible. Therefore, by the behavior cloning process,  $\sigma_{BC(\mathcal{D})}$  is identical to the behavior strategy underlying  $\mathcal{D}$ , i.e.,  $\sigma_{BC(\mathcal{D})} = \sigma_{\mathcal{D}}$ . Consequently, if  $\mathcal{D}$  is an  $\epsilon$ -equilibrium coverage of  $\mathcal{G}$ , then  $\sigma_{BC(\mathcal{D})}$  is an  $\epsilon$ -equilibrium strategy for the IIEFG  $\mathcal{G}$ , as  $\text{GAP}(\sigma_{\mathcal{D}}, \sigma^*) < \epsilon$  implies  $\text{GAP}(\sigma_{BC(\mathcal{D})}, \sigma^*) < \epsilon$ . Any slight violation of these conditions would invalidate the convergence result. The full proof is provided in App. D.1.  $\square$

Building on the insights provided by the preceding two theorems, we propose the following theorem concerning the performance of BOMB under a general case where an unknown strategy profile generates the offline dataset. The full proof can be found in App. D.1.

**Theorem 4.6.** *Let  $\sigma_{BOMB(\mathcal{D})}$  represent the strategy profile learned by our **BOMB algorithm** based on the offline dataset  $\mathcal{D}$  with sufficient data under Assumption 4.3,  $\sigma_{\mathcal{D}}$  represent the underlying behavior strategy of  $\mathcal{D}$  and  $\sigma^*$  represent the equilibrium strategy of IIEFG  $\mathcal{G}$ . Then the gap between  $\sigma_{BOMB(\mathcal{D})}$  and  $\sigma^*$  is at most equal to, or smaller than, the gap between  $\sigma_{\mathcal{D}}$  and  $\sigma^*$ , i.e.,  $\text{GAP}(\sigma_{BOMB(\mathcal{D})}, \sigma^*) \leq \text{GAP}(\sigma_{\mathcal{D}}, \sigma^*)$ .*

To better analyze the performance of our algorithm under real-world cases, we first analyze the offline dataset we generated for the offline EF paradigm. Based on dataset collection procedures, we find that the random dataset can be considered as a uniform coverage of the game  $\mathcal{G}$  when the dataset is sufficiently large. This is because the random dataset is collected using a uniform strategy, ensuring that every action is adequately sampled as long as enough data is collected. On the other hand, the expert dataset can be considered as an  $\epsilon$ -equilibrium coverage of game  $\mathcal{G}$ , where  $\epsilon$  decreases as the dataset size increases. Since the expert dataset is generated by an equilibrium strategy, a larger sample size means the underlying behavior strategy of the dataset more closely approximates the equilibrium strategy, resulting in a smaller  $\epsilon$ . Therefore, the above properties of our algorithm hold under these two datasets, as shown in the following experimental results.

## 5 EXPERIMENTAL RESULTS

To assess the performance of our proposed algorithm – BOMB, we conduct the following experiments: i) we compare two offline RL algorithms to our BOMB algorithm; ii) we evaluate the performance of different estimation methods; and iii) we run the BOMB framework on various offline datasets to evaluate its performance in computing different equilibrium strategies.

We use OpenSpiel<sup>1</sup> (Lanctot et al., 2019) as our experimental platform, as it offers a well-established collection of environments and algorithms for game research, thereby facilitating future replicability. We select several poker games, Liar’s Dice and Phantom Tic-Tac-Toe, which are all widely used in previous works (Lisý et al., 2015; Brown et al., 2019). Experiments are conducted on a workstation with a ten-core 3.3GHz Intel i9-9820X CPU and NVIDIA RTX 2080Ti GPU. All results are averaged over three seeds and error bars are also reported. To demonstrate the performance of our algorithm, we present our results by answering the following research questions (RQs).

**RQ1:** *Can the BOMB framework outperform offline RL methods?*

To support this claim that offline RL algorithms are insufficient for the offline EF paradigm, we choose one model-free algorithm–Best-Action Imitation Learning (BAIL) (Chen et al., 2020) and one model-based algorithm–Model-based Offline Policy Optimization (MOPO) (Yu et al., 2020) as

<sup>1</sup>[https://github.com/deepmind/open\\_spiel](https://github.com/deepmind/open_spiel)

378  
379  
380  
381  
382  
383  
384  
385  
386  
387  
388  
389  
390  
391  
392  
393  
394  
395  
396  
397  
398  
399  
400  
401  
402  
403  
404  
405  
406  
407  
408  
409  
410  
411  
412  
413  
414  
415  
416  
417  
418  
419  
420  
421  
422  
423  
424  
425  
426  
427  
428  
429  
430  
431

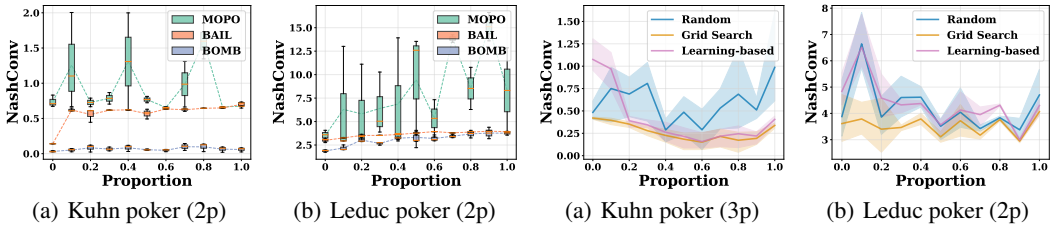


Figure 5: Comparison with offline RL. Figure 6: Results of different estimation methods.

the representative of offline RL algorithms. Fig. 5 shows the comparison results in two-player Kuhn poker and Leduc poker games under hybrid datasets. The x-axis represents the proportion of data from the random datasets in the hybrid dataset. When the ratio is zero, the hybrid dataset is equivalent to the expert dataset; conversely, when the ratio is one, the hybrid dataset is reduced to the random dataset. We found that BOMB outperforms both offline RL algorithms in all cases. It means that neither of these offline RL algorithms can produce a strategy profile close enough to the equilibrium strategy, which might be attributed to the players’ policies being optimized independently.

**RQ2:** How do different parameter estimation methods perform?

As introduced previously, we propose three combination methods: random method, grid search method, and learning-based method. To evaluate the performance of different combination methods, we conduct experiments on poker games. For the learning-based method, we train the parameter predictor on the two-player Kuhn poker game. And then, we test the parameter predictor in other poker games. Fig. 6 shows the performance results of three combination methods on three-player Kuhn poker and two-player Leduc poker games. We can find that the grid search method achieves the best performance and the learning-based method performs similarly to the grid search method on the three-player Kuhn poker while it performs slightly worse on the two-player Leduc poker game. It implies that the performance of the parameter predictor mainly depends on the difference between the test game and the game used to train the predictor. The most interesting result is that the random method performs well in many cases, which means that even a simple combination works well. In the rest of the experiments, we use the grid search method as the parameter estimation method.

**RQ3:** Can the BOMB framework compute NE?

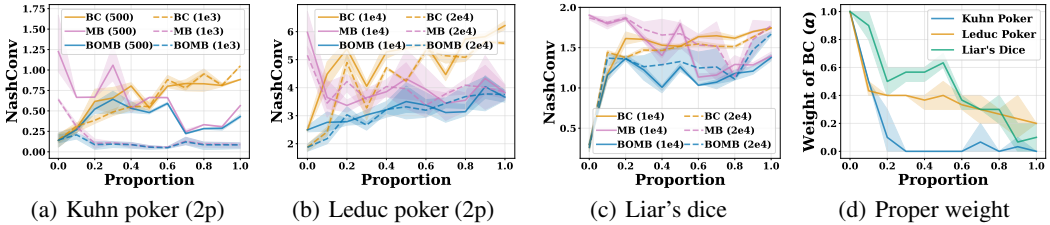


Figure 7: Experimental results on computing NE in two-player games.

To answer this question, we conduct extensive experiments covering two-player cases, multi-player cases, and real-world scenarios simulated using learning datasets. This comprehensive approach allows for an adequate evaluation of our method’s performance in computing the NE strategy.

**Two-Player Cases.** We first move to evaluate the performance of our algorithm, BOMB, in computing the NE strategy. In addition to performing the BOMB framework, we also assess the individual performance of the behavior cloning technique and the model-based algorithm. This assessment not only helps in understanding the strengths and weaknesses of each component but also provides a comprehensive insight into the efficacy of the BOMB framework in computing the equilibrium offline. Figs. 7(a)-7(c) show results on some two-player games under different sizes of offline datasets. Here, we use MB-CFR or MB-PSRO to compute the NE strategy. The MB framework’s perfor-



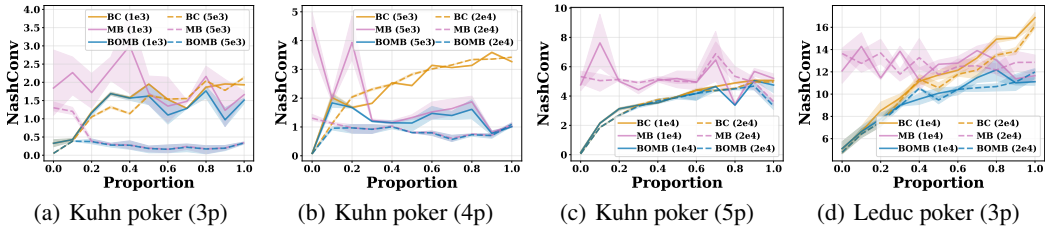


Figure 8: Experimental results on computing NE in multi-player games.

mance is independent of the algorithm used to compute the NE strategy (shown in App. F). As the proportion of the random dataset increases, we observe that the performance of BC decreases while the performance of MB slightly increases. Additionally, we notice that as the size of offline data increases, the improvement of the BC’s performance is not significant and the MB’s performance improves. It means that the performance of BC mainly depends on the quality of datasets, i.e., the quality of the behavior policy generating the dataset, and the performance of MB relies on the similarity between the environment model and the actual environment. These figures show that our algorithm, BOMB, outperforms both BC and MB methods in all cases, demonstrating its effectiveness in computing NE strategy for two-player imperfect-information extensive-form games.

To further analyze the performance of the BOMB method in two-player games, we plot the parameter  $\alpha$  for these combined policies, as illustrated in Fig. 7(d). The results show that as the proportion of the random dataset in the hybrid dataset increases, the weight of the BC policy decreases. It confirms that the BC policy performs better under the expert dataset while the MB policy performs better under the random dataset from another side.

**Multi-Player Cases.** We also conduct experiments on multi-player games, specifically evaluating the performance of our method in computing the NE strategy across several multi-player games, as shown in Fig. 8. The results demonstrate that our BOMB framework consistently performs as well as or better than both BC and MB algorithms, similar to the findings in two-player scenarios. Furthermore, as the proportion of the random dataset increases, the performance of BC decreases, while MB shows instability with a slight downward trend. It is important to note that we adopt MB-CFR as the model-based algorithm, and since CFR-based algorithms do not guarantee convergence to the NE strategy in multi-player games, the performance of MB may be affected. Additionally, the performance of the model-based method also relies on the accuracy of the trained environment game model. Consequently, the underperformance of the MB algorithm may be due to either an inadequately trained environment game model or the limitations of the CFR-based algorithm in multi-player settings. Therefore, developing an effective equilibrium-finding algorithm and training an accurate environment game model are both key challenges for offline EF in multi-player imperfect-information extensive-form games.

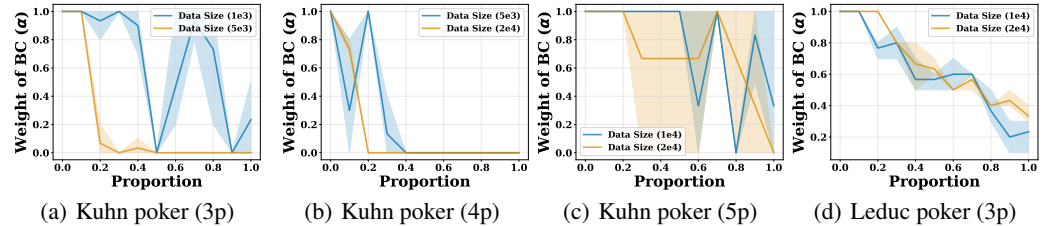


Figure 9: Proper weight of BC policy in multi-player games.

The appropriate weights of the BC policy ( $\alpha$ ) within the BOMB framework across different hybrid datasets are presented in Fig. 9. In three-player and four-player Kuhn poker games, we observe that the weight of the BC policy quickly drops to zero as the proportion of the random dataset in the hybrid dataset increases, indicating that the MB method generally outperforms the BC method, except when the random dataset proportion is low. In contrast, in five-player Kuhn poker and three-

486  
487  
488  
489  
490  
491  
492  
493  
494  
495  
496  
497  
498  
499  
500  
501  
502  
503  
504  
505  
506  
507  
508  
509  
510  
511  
512  
513  
514  
515  
516  
517  
518  
519  
520  
521  
522  
523  
524  
525  
526  
527  
528  
529  
530  
531  
532  
533  
534  
535  
536  
537  
538  
539

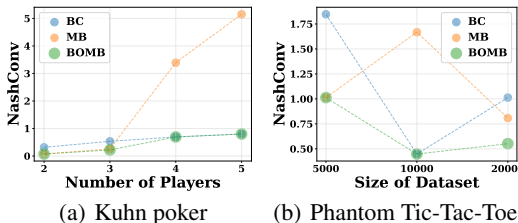


Figure 10: Results on learning dataset.

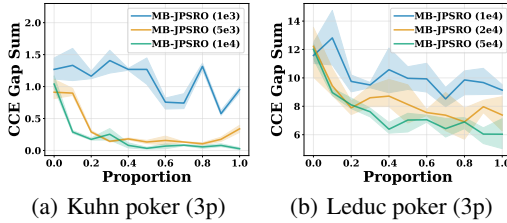


Figure 11: Results on computing CCE.

player Leduc poker games, the weight of the BC policy remains high in most cases, except when the proportion of the random dataset is high. This may be due to the poor performance of the MB method in these games, highlighting the challenge of learning an approximate equilibrium strategy using the MB method in complex, multi-player games, where both developing an effective equilibrium-finding algorithm and training an accurate environment game model are particularly difficult.

**Simulating Real-World Cases.** We also conduct experiments on the learning dataset, which closely approximates real-world conditions. Fig. 10(a) shows the results of Kuhn poker games with different numbers of players and Fig. 10(b) shows the results of Phantom Tic-Tac-Toe under different numbers of offline data. It indicates that given an offline dataset generated by an unknown strategy, our algorithm can also perform better than BC and MB in approximating the NE strategy.

**RQ4: Can the BOMB framework compute CCE?**

We proceed to evaluate the performance of the model-based method in computing the CCE strategy. We do not perform the BC technique and the BOMB framework to compute the CCE strategy since the offline dataset is collected using an independent strategy for each player, rather than a joint strategy. Fig. 11 shows the results of performing the MB-JPSRO algorithm on three-player Kuhn poker and Leduc poker games. We can observe that as the size of the offline data increases, the performance of MB-JPSRO improves. This further supports the notion that the performance of the model-based method primarily depends on the quality of the trained environment model and also highlights its significance in computing equilibrium strategy offline.

6 CONCLUSION

We investigated the paradigm of offline equilibrium finding (Offline EF) in extensive-form games, which focuses on finding equilibrium strategies from offline datasets. To be specific, we first created the offline EF datasets using several established data-collecting methods, which solves the challenge of the absence of a comprehensive dataset for evaluation. Then, we proposed a novel algorithm, BOMB, which combines the behavior cloning technique with a model-based approach that can adapt regular online equilibrium finding algorithms to the offline setting by introducing an environment model. To better understand the algorithm, we provide a comprehensive theoretical and empirical analysis, providing performance guarantees of our algorithms across different offline datasets. Finally, extensive experimental results further validated the superiority of the BOMB framework over existing offline RL algorithms, affirming its efficacy for computing equilibrium strategies in an offline manner. We hope our efforts can open up new avenues in equilibrium finding and accelerate research in large-scale game theory.

**Limitations and Future Work.** There are several limitation of this work. First, the games considered are relatively small-scale, and the large-scale games including Texas Hold'em poker (Brown & Sandholm, 2018) and football games (Liu et al., 2022) will be included in the future work. Second, this work primarily focus on NE and CCE, more solution concepts will be considered such as quantal response equilibrium (QRE) (McKelvey & Palfrey, 1995) and  $\alpha$ -rank (Omidshafiei et al., 2019). We will investigate the generalizability of both the datasets and the BOMB framework to novel solution concepts in the future work. Third, the relationships between the datasets and the offline EF algorithms can be further investigated, where instead of collecting the datasets by researchers, we can apply the offline EF to human-play datasets toward real deployment (Wang et al., 2024).

## REFERENCES

- 540  
541  
542 Abbas Abdolmaleki, Jost Tobias Springenberg, Yuval Tassa, Remi Munos, Nicolas Heess, and Martin Riedmiller. Maximum a posteriori policy optimisation. *arXiv preprint arXiv:1806.06920*, 2018.  
544
- 545 Maruan Al-Shedivat, Trapit Bansal, Yuri Burda, Ilya Sutskever, Igor Mordatch, and Pieter Abbeel. Continuous adaptation via meta-learning in nonstationary and competitive environments. *arXiv preprint arXiv:1710.03641*, 2017.  
546  
547  
548
- 549 Robert J Aumann. Correlated equilibrium as an expression of Bayesian rationality. *Econometrica: Journal of the Econometric Society*, pp. 1–18, 1987.  
550
- 551 Noam Brown and Tuomas Sandholm. Superhuman AI for heads-up no-limit poker: Libratus beats top professionals. *Science*, 359(6374):418–424, 2018.  
552  
553
- 554 Noam Brown and Tuomas Sandholm. Superhuman ai for multiplayer poker. *Science*, 365(6456):885–890, 2019.  
555
- 556 Noam Brown, Adam Lerer, Sam Gross, and Tuomas Sandholm. Deep counterfactual regret minimization. In *ICML*, pp. 793–802, 2019.  
557  
558
- 559 Xinyue Chen, Zijian Zhou, Zheng Wang, Che Wang, Yanqiu Wu, and Keith Ross. Bail: Best-action imitation learning for batch deep reinforcement learning. In *NeurIPS*, pp. 18353–18363, 2020.  
560  
561
- 562 Qiwen Cui and Simon S. Du. When is offline two-player zero-sum markov game solvable? *arXiv preprint arXiv:2201.03522*, 2022.  
563
- 564 Jakob N Foerster, Richard Y Chen, Maruan Al-Shedivat, Shimon Whiteson, Pieter Abbeel, and Igor Mordatch. Learning with opponent-learning awareness. *arXiv preprint arXiv:1709.04326*, 2017.  
565  
566
- 567 Justin Fu, Aviral Kumar, Ofir Nachum, George Tucker, and Sergey Levine. D4rl: Datasets for deep data-driven reinforcement learning. *arXiv preprint arXiv:2004.07219*, 2020.  
568
- 569 Scott Fujimoto and Shixiang Gu. A minimalist approach to offline reinforcement learning. In *NeurIPS*, pp. 20132–20145, 2021.  
570  
571
- 572 Scott Fujimoto, Edoardo Conti, Mohammad Ghavamzadeh, and Joelle Pineau. Benchmarking batch deep reinforcement learning algorithms. *arXiv preprint arXiv:1910.01708*, 2019.  
573  
574
- 575 Richard Gibson, Marc Lanctot, Neil Burch, Duane Szafron, and Michael Bowling. Generalized sampling and variance in counterfactual regret minimization. In *AAAI*, pp. 1355–1361, 2012.  
576
- 577 Aditya Grover, Maruan Al-Shedivat, Jayesh Gupta, Yuri Burda, and Harrison Edwards. Learning policy representations in multiagent systems. In *ICML*, pp. 1802–1811, 2018.  
578  
579
- 580 Caglar Gulcehre, Ziyu Wang, Alexander Novikov, Thomas Paine, Sergio Gómez, Konrad Zolna, Rishabh Agarwal, Josh S Merel, Daniel J Mankowitz, Cosmin Paduraru, et al. Rl unplugged: A suite of benchmarks for offline reinforcement learning. *Advances in Neural Information Processing Systems*, 33:7248–7259, 2020.  
581  
582  
583
- 584 Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In *ICML*, pp. 1861–1870, 2018.  
585  
586
- 587 He He, Jordan Boyd-Graber, Kevin Kwok, and Hal Daumé III. Opponent modeling in deep reinforcement learning. In *ICML*, pp. 1804–1813, 2016.  
588  
589
- 590 Zhang-Wei Hong, Shih-Yang Su, Tzu-Yun Shann, Yi-Hsiang Chang, and Chun-Yi Lee. A deep policy inference q-network for multi-agent systems. *arXiv preprint arXiv:1712.07893*, 2017.  
591  
592
- 593 Patrick R Jordan, L Julian Schvartzman, and Michael P Wellman. Strategy exploration in empirical games. In *AAMAS*, pp. 1131–1138, 2010.

- 594 Dmitry Kalashnikov, Alex Irpan, Peter Pastor, Julian Ibarz, Alexander Herzog, Eric Jang, Deirdre  
595 Quillen, Ethan Holly, Mrinal Kalakrishnan, Vincent Vanhoucke, et al. Scalable deep reinforce-  
596 ment learning for vision-based robotic manipulation. In *Conference on robot learning*, pp. 651–  
597 673. PMLR, 2018.
- 598 Ansam Khraisat, Iqbal Gondal, Peter Vamplew, and Joarder Kamruzzaman. Survey of intrusion  
599 detection systems: Techniques, datasets and challenges. *Cybersecurity*, 2(1):1–22, 2019.
- 600 Rahul Kidambi, Aravind Rajeswaran, Praneeth Netrapalli, and Thorsten Joachims. Morel: Model-  
601 based offline reinforcement learning. In *NeurIPS*, pp. 21810–21823, 2020.
- 602 Dong Ki Kim, Miao Liu, Matthew D Riemer, Chuangchuang Sun, Marwa Abdulhai, Golnaz Habibi,  
603 Sebastian Lopez-Cot, Gerald Tesauro, and Jonathan How. A policy gradient algorithm for learning  
604 to learn in multiagent reinforcement learning. In *ICML*, pp. 5541–5550, 2021.
- 605 B Ravi Kiran, Ibrahim Sobh, Victor Talpaert, Patrick Mannion, Ahmad A Al Sallab, Senthil Yoga-  
606 mani, and Patrick Pérez. Deep reinforcement learning for autonomous driving: A survey. *IEEE*  
607 *Transactions on Intelligent Transportation Systems*, pp. 4909–4926, 2022.
- 608 Usha Kiruthika, Thamarai Selvi Somasundaram, and S Kanaga Suba Raja. Lifecycle model of a ne-  
609 gotiation agent: A survey of automated negotiation techniques. *Group Decision and Negotiation*,  
610 29:1239–1262, 2020.
- 611 Simon Kornblith, Mohammad Norouzi, Honglak Lee, and Geoffrey Hinton. Similarity of neural  
612 network representations revisited. In *ICML*, pp. 3519–3529, 2019.
- 613 Aviral Kumar, Aurick Zhou, George Tucker, and Sergey Levine. Conservative q-learning for offline  
614 reinforcement learning. In *NeurIPS*, pp. 1179–1191, 2020.
- 615 Marc Lanctot, Kevin Waugh, Martin Zinkevich, and Michael Bowling. Monte Carlo sampling for  
616 regret minimization in extensive games. In *NeurIPS*, pp. 1078–1086, 2009.
- 617 Marc Lanctot, Vinicius Zambaldi, Audrūnas Gruslys, Angeliki Lazaridou, Karl Tuyls, Julien  
618 Pérolat, David Silver, and Thore Graepel. A unified game-theoretic approach to multiagent rein-  
619 forcement learning. In *NeurIPS*, pp. 4193–4206, 2017.
- 620 Marc Lanctot, Edward Lockhart, Jean-Baptiste Lespiau, Vinicius Zambaldi, Satyaki Upadhyay,  
621 Julien Pérolat, Sriram Srinivasan, Finbarr Timbers, Karl Tuyls, Shayegan Omidshafiei, et al.  
622 OpenSpiel: A framework for reinforcement learning in games. *arXiv preprint arXiv:1908.09453*,  
623 2019.
- 624 Seunghyun Lee, Younggyo Seo, Kimin Lee, Pieter Abbeel, and Jinwoo Shin. Offline-to-online  
625 reinforcement learning via balanced replay and pessimistic q-ensemble. In *Conference on Robot*  
626 *Learning*, pp. 1702–1712. PMLR, 2022.
- 627 Sergey Levine, Aviral Kumar, George Tucker, and Justin Fu. Offline reinforcement learning: Tuto-  
628 rial, review, and perspectives on open problems. *arXiv preprint arXiv:2005.01643*, 2020.
- 629 Hui Li, Kailiang Hu, Shaohua Zhang, Yuan Qi, and Le Song. Double neural counterfactual regret  
630 minimization. In *ICLR*, 2019.
- 631 Shuxin Li, Xinrun Wang, Youzhi Zhang, Wanqi Xue, Jakub Černý, and Bo An. Solving large-scale  
632 pursuit-evasion games using pre-trained strategies. In *Proceedings of the AAAI Conference on*  
633 *Artificial Intelligence*, volume 37, pp. 11586–11594, 2023.
- 634 Viliam Lisý, Marc Lanctot, and Michael Bowling. Online Monte Carlo counterfactual regret mini-  
635 mization for search in imperfect information games. In *AAMAS*, pp. 27–36, 2015.
- 636 Siqi Liu, Kay Choong See, Kee Yuan Ngiam, Leo Anthony Celi, Xingzhi Sun, Mengling Feng,  
637 et al. Reinforcement learning for clinical decision support in critical care: Comprehensive review.  
638 *Journal of Medical Internet Research*, 22(7):e18477, 2020.
- 639 Siqi Liu, Guy Lever, Zhe Wang, Josh Merel, SM Ali Eslami, Daniel Hennes, Wojciech M Czarnecki,  
640 Yuval Tassa, Shayegan Omidshafiei, Abbas Abdolmaleki, et al. From motor control to team play  
641 in simulated humanoid football. *Science Robotics*, 7(69):eabo0235, 2022.



- 648 Luke Marris, Paul Muller, Marc Lanctot, Karl Tuyls, and Thore Graepel. Multi-agent training  
649 beyond zero-sum with correlated equilibrium meta-solvers. *arXiv preprint arXiv:2106.09435*,  
650 2021.
- 651 Michael Mathieu, Sherjil Ozair, Srivatsan Srinivasan, Caglar Gulcehre, Shangdong Zhang, Ray  
652 Jiang, Tom Le Paine, Konrad Zolna, Richard Powell, Julian Schrittwieser, et al. StarCraft II  
653 unplugged: Large scale offline reinforcement learning. In *Deep RL Workshop NeurIPS 2021*,  
654 2021.
- 655 Tatsuya Matsushima, Hiroki Furuta, Yutaka Matsuo, Ofir Nachum, and Shixiang Gu. Deployment-  
656 efficient reinforcement learning via model-based offline optimization. *arXiv preprint*  
657 *arXiv:2006.03647*, 2020.
- 659 Stephen McAleer, John Lanier, Roy Fox, and Pierre Baldi. Pipeline psro: A scalable approach for  
660 finding approximate nash equilibria in large games. In *NeurIPS*, pp. 20238–20248, 2020.
- 661 Stephen McAleer, Kevin Wang, Marc Lanctot, John Lanier, Pierre Baldi, and Roy Fox. Anytime  
662 optimal psro for two-player zero-sum games. *arXiv preprint arXiv:2201.07700*, 2022.
- 664 Stephen Marcus McAleer, John Banister Lanier, Kevin Wang, Pierre Baldi, and Roy Fox. Xdo: A  
665 double oracle algorithm for extensive-form games. In *NeurIPS*, pp. 23128–23139, 2021.
- 666 Richard D McKelvey and Thomas R Palfrey. Quantal response equilibria for normal form games.  
667 *Games and economic behavior*, 10(1):6–38, 1995.
- 669 H Brendan McMahan, Geoffrey J Gordon, and Avrim Blum. Planning in the presence of cost  
670 functions controlled by an adversary. In *ICML*, pp. 536–543, 2003.
- 671 John F Nash. Equilibrium points in n-person games. *PNAS*, 36(1):48–49, 1950.
- 672 Noam Nisan, Tim Roughgarden, Eva Tardos, and Vijay V Vazirani. *Algorithmic Game Theory*.  
673 Cambridge University Press, 2007.
- 675 Shayegan Omidshafiei, Christos Papadimitriou, Georgios Piliouras, Karl Tuyls, Mark Rowland,  
676 Jean-Baptiste Lespiau, Wojciech M Czarnecki, Marc Lanctot, Julien Perolat, and Remi Munos.  
677  $\alpha$ -rank: Multi-agent evaluation by evolution. *Scientific Reports*, 9(1):1–29, 2019.
- 678 Julien Perolat, Bart De Vylder, Daniel Hennes, Eugene Tarassov, Florian Strub, Vincent de Boer,  
679 Paul Muller, Jerome T Connor, Neil Burch, Thomas Anthony, et al. Mastering the game of  
680 stratego with model-free multiagent reinforcement learning. *Science*, 378(6623):990–996, 2022.
- 682 Rafael Figueiredo Prudencio, Marcos ROA Maximo, and Esther Luna Colombini. A survey  
683 on offline reinforcement learning: Taxonomy, review, and open problems. *arXiv preprint*  
684 *arXiv:2203.01387*, 2022.
- 685 Paria Rashidinejad, Banghua Zhu, Cong Ma, Jiantao Jiao, and Stuart Russell. Bridging offline rein-  
686 forcement learning and imitation learning: A tale of pessimism. *Advances in Neural Information*  
687 *Processing Systems*, 34:11702–11716, 2021.
- 688 Martin Schmid, Neil Burch, Marc Lanctot, Matej Moravcik, Rudolf Kadlec, and Michael Bowling.  
689 Variance reduction in Monte Carlo counterfactual regret minimization (VR-MCCFR) for  
690 extensive form games using baselines. In *AAAI*, pp. 2157–2164, 2019.
- 692 Julian Schrittwieser, Ioannis Antonoglou, Thomas Hubert, Karen Simonyan, Laurent Sifre, Simon  
693 Schmitt, Arthur Guez, Edward Lockhart, Demis Hassabis, Thore Graepel, et al. Mastering Atari,  
694 Go, chess and shogi by planning with a learned model. *Nature*, 588(7839):604–609, 2020.
- 695 L Julian Schvartzman and Michael P Wellman. Exploring large strategy spaces in empirical game  
696 modeling. *Agent Mediated Electronic Commerce (AMEC 2009)*, pp. 139, 2009a.
- 697 L Julian Schvartzman and Michael P Wellman. Stronger cda strategies through empirical game-  
698 theoretic analysis and reinforcement learning. In *AAMAS*, pp. 249–256, 2009b.
- 699 Shai Shalev-Shwartz and Shai Ben-David. *Understanding machine learning: From theory to algo-*  
700 *ritms*. Cambridge university press, 2014.

- 702 Yoav Shoham and Kevin Leyton-Brown. *Multiagent Systems: Algorithmic, Game-Theoretic, and*  
703 *Logical Foundations*. Cambridge University Press, 2008.
- 704
- 705 Noah Y Siegel, Jost Tobias Springenberg, Felix Berkenkamp, Abbas Abdolmaleki, Michael Ne-  
706 unert, Thomas Lampe, Roland Hafner, Nicolas Heess, and Martin Riedmiller. Keep doing  
707 what worked: Behavioral modelling priors for offline reinforcement learning. *arXiv preprint*  
708 *arXiv:2002.08396*, 2020.
- 709 Bharat Singh, Rajesh Kumar, and Vinay Pratap Singh. Reinforcement learning in robotic applica-  
710 tions: A comprehensive survey. *Artificial Intelligence Review*, pp. 1–46, 2021.
- 711
- 712 Adish Singla, Anna N Rafferty, Goran Radanovic, and Neil T Heffernan. Reinforcement learning  
713 for education: Opportunities and challenges. *arXiv preprint arXiv:2107.08828*, 2021.
- 714 Eric Steinberger. Single deep counterfactual regret minimization. *arXiv preprint arXiv:1901.07621*,  
715 2019.
- 716
- 717 Robert Tarjan. Depth-first search and linear graph algorithms. *SIAM journal on computing*, 1(2):  
718 146–160, 1972.
- 719 Zhe Wang, Petar Veličković, Daniel Hennes, Nenad Tomašev, Laurel Prince, Michael Kaisers,  
720 Yoram Bachrach, Romuald Elie, Li Kevin Wenliang, Federico Piccinini, et al. TacticAI: an AI  
721 assistant for football tactics. *Nature communications*, 15(1):1906, 2024.
- 722
- 723 Kevin Waugh, David Schnizlein, Michael H Bowling, and Duane Szafron. Abstraction pathologies  
724 in extensive games. In *AAMAS*, pp. 781–788, 2009.
- 725 Michael P Wellman. Methods for empirical game-theoretic analysis. In *AAAI*, pp. 1552–1556, 2006.
- 726
- 727 Tengyang Xie, Ching-An Cheng, Nan Jiang, Paul Mineiro, and Alekh Agarwal. Bellman-consistent  
728 pessimism for offline reinforcement learning. *Advances in neural information processing systems*,  
729 34:6683–6694, 2021.
- 730 Wanqi Xue, Youzhi Zhang, Shuxin Li, Xinrun Wang, Bo An, and Chai Kiat Yeo. Solving large-scale  
731 extensive-form network security games via neural fictitious self-play. In *IJCAI*, pp. 3713–3720,  
732 2021.
- 733
- 734 Tianhe Yu, Garrett Thomas, Lantao Yu, Stefano Ermon, James Y Zou, Sergey Levine, Chelsea Finn,  
735 and Tengyu Ma. MOPO: Model-based offline policy optimization. In *NeurIPS*, pp. 14129–14142,  
736 2020.
- 737 Tianhe Yu, Aviral Kumar, Rafael Rafailov, Aravind Rajeswaran, Sergey Levine, and Chelsea Finn.  
738 Combo: Conservative offline model-based policy optimization. *Advances in Neural Information*  
739 *Processing Systems*, 34, 2021a.
- 740 Xiaopeng Yu, Jiechuan Jiang, Haobin Jiang, and Zongqing Lu. Model-based opponent modeling.  
741 *arXiv preprint arXiv:2108.01843*, 2021b.
- 742
- 743 Daochen Zha, Jingru Xie, Wenye Ma, Sheng Zhang, Xiangru Lian, Xia Hu, and Ji Liu. Douzero:  
744 Mastering doudizhu with self-play deep reinforcement learning. In *ICML*, pp. 12333–12344,  
745 2021.
- 746
- 747 Yan Zheng, Zhaopeng Meng, Jianye Hao, Zongzhang Zhang, Tianpei Yang, and Changjie Fan. A  
748 deep bayesian policy reuse approach against non-stationary agents. In *NeurIPS*, pp. 962–972,  
749 2018.
- 750 Han Zhong, Wei Xiong, Jiyuan Tan, Liwei Wang, Tong Zhang, Zhaoran Wang, and Zhuoran Yang.  
751 Pessimistic minimax value iteration: Provably efficient equilibrium learning from offline datasets.  
752 In *ICML*, pp. 27117–27142, 2022.
- 753 Martin Zinkevich, Michael Johanson, Michael Bowling, and Carmelo Piccione. Regret minimization  
754 in games with incomplete information. In *NeurIPS*, pp. 1729–1736, 2007.
- 755

## A FREQUENTLY ASKED QUESTIONS

### Q1: What are the potential impacts of this work?

This work fills the lack of offline learning in the game theory field. Benefiting from the offline learning framework, we anticipate that our offline equilibrium finding setting could pave a path to solving real-world problems using these game theory-based methods and inspire new research directions in equilibrium finding. Furthermore, the equilibrium strategy is more robust compared with just the optimal strategy in some security-related scenarios. Consequently, offline EF plays a crucial role in obtaining more robust strategies for tackling these competitive real-world problems.

### Q2: Why offline EF is important and is more difficult than offline cooperative MARL?

Utilizing offline EF algorithms specifically designed for adversarial environments is crucial in strictly competitive games, such as security games. This setting fundamentally differs from offline multi-agent reinforcement learning, which generally focuses on cooperation between agents rather than strict competition. For instance, consider the class of pursuit-evasion games, where the pursuer (defender) chases the evader (attacker). In this scenario, we cannot make any assumptions about the attacker’s strategy beforehand, as the attacker is strategic and capable of learning. Employing a vanilla offline RL algorithm to learn the defender’s optimal strategy based solely on historical data might lead to a significant utility loss, as the defender’s optimal strategy could be exploitable. In other words, the attacker may switch to the best response against the computed strategy of the defender instead of adhering to their past behavior estimated from the data. Therefore, achieving Nash Equilibrium (NE) may be a more suitable solution, as NE strategies are non-exploitable.

To be more specific, traditional offline RL focuses on learning the optimal strategy, i.e., obtaining the highest utility, for an agent acting in a dynamic environment modeled as a single MDP, which does not depend on the actions of other agents. In contrast, in two-player games, the dynamics for one player depend not only on the environment but also on the strategy of the opponent. In other words, the MDP in which a player acts in games is determined by both the game and the fixed strategy of the opponent, and hence a change in the opponent’s strategy instigates a corresponding change in the MDP. This makes computing the best strategy for the defender against a strategic opponent using offline RL significantly more difficult. The framework of offline EF we introduced provides methods for computing a player’s NE strategy, which is their optimal strategy against a strategic opponent (i.e., the worst case for the player).

### Q3: What are the differences between Offline EF and EGTA?

1) As described in (Wellman, 2006), EGTA takes the game simulator as input and performs strategic reasoning through interleaved simulation and game-theoretic analysis. Therefore, **the game simulator is required in EGTA**. However, only the offline dataset is available in the offline EF paradigm and the game simulator is not required. 2) The estimated game model (empirical game) in EGTA is built based on the simulation’s results, which are obtained by performing **known strategies** on the simulator. In contrast, in the offline EF paradigm, the offline dataset is generated with an **unknown strategy**. Although we use different behavior strategies to generate several offline datasets, we do not utilize these strategies in the offline EF paradigm.

### Q4: What are the novelties of the proposed Offline EF algorithm – BOMB?

To our knowledge, we are **the first ones to propose an empirical algorithm** for computing the equilibrium strategy from the offline dataset, i.e., the offline EF paradigm. Unlike traditional offline RL algorithms, which belong to either model-based or model-free categories, our algorithm combines the advantages of both model-based and model-free approaches to efficiently compute equilibrium strategies in an offline manner. Our BOMB framework integrates the behavior cloning technique with a model-based method, equipping novel parameter estimation methods. We introduce an environment model to design the model-based method that can generalize regular online equilibrium finding algorithms to the offline setting. Furthermore, we proposed several different methods to determine the combination parameter value. In different scenarios, according to whether the online interaction is available, there are corresponding algorithms to determine the parameter value. Finally, experimental results show that BC and MB cannot perform consistently well and BOMB outperforms them in all cases. It indicated that our BOMB framework takes advantage of both algorithms and performs well in computing equilibrium strategies in an offline manner.

## B MORE RELATED WORK

**Equilibrium Finding Algorithms.** As described in the main paper, the contemporary state-of-the-art algorithms for solving IIEFGs may be roughly divided into two groups: no-regret methods derived from CFR (Zinkevich et al., 2007), and incremental strategy-space generation methods of the PSRO framework (Lanctot et al., 2017). Next, we will introduce these two classes of algorithms.

For the first group, CFR is a family of iterative methods for approximately solving imperfect-information extensive-form games. Let  $\sigma_i^t$  be the strategy used by player  $i$  in iteration  $t$ . We use  $u_i(\sigma, h)$  to define the expected utility of player  $i$  given that the history  $h$  is reached and all players act according to strategy  $\sigma$  from that point on. Accordingly,  $u_i(\sigma, h \cdot a)$  is used to define the expected utility of player  $i$  given that the history  $h$  is reached and all players play according to strategy  $\sigma$  except player  $i$  selects action  $a$  in history  $h$ . Formally,  $u_i(\sigma, h) = \sum_{z \in Z} \pi^\sigma(h, z) u_i(z)$  and  $u_i(\sigma, h \cdot a) = \sum_{z \in Z} \pi^\sigma(h \cdot a, z) u_i(z)$ . The *counterfactual value* of the information set  $I$ ,  $v_i^\sigma(I)$ , is the expected value of information set  $I$  given that player  $i$  attempts to reach it. This value is the weighted average of the expected utility of each history in the information set. The weight is proportional to the contribution of all players except player  $i$  to reach each history. Thus,  $v_i^\sigma(I) = \sum_{h \in I} \pi_{-i}^\sigma(h) u_i(\sigma, h)$ . For any action  $a \in A(I)$ , the counterfactual value of action  $a$  is  $v_i^\sigma(I, a) = \sum_{h \in I} \pi_{-i}^\sigma(h) u_i(\sigma, h \cdot a)$ . The *instantaneous counterfactual regret* for an action  $a$  in information set  $I$  during iteration  $t$  is  $r^t(I, a) = v_{P(I)}^\sigma(I, a) - v_{P(I)}^\sigma(I)$ . Therefore, the counterfactual regret for an action  $a$  in information set  $I$  on iteration  $T$  is  $R^T(I, a) = \sum_{t=1}^T r^t(I, a)$ . In vanilla CFR, players use *Regret Matching* to pick a distribution over available actions in an information set proportional to the cumulative regret of those actions. Formally, in iteration  $T + 1$ , player  $i$  selects action  $a \in A(I)$  according to probabilities

$$\sigma^{T+1}(I, a) = \begin{cases} \frac{R_+^T(I, a)}{\sum_{b \in A(I)} R_+^T(I, b)} & \text{if } \sum_{b \in A(I)} R_+^T(I, b) > 0, \\ \frac{1}{|A(I)|} & \text{otherwise,} \end{cases}$$

where  $R_+^T(I, a) = \max\{R^T(I, a), 0\}$  is the position portion of the regret value since we often are most concerned about the cumulative regret when it is positive. If a player acts according to regret matching in the information set  $I$  on every iteration, then in iteration  $T$ ,  $R^T(I) \leq \Delta_i \sqrt{|A_i|} \sqrt{T}$  where  $\Delta_i = \max_z u_i(z) - \min_z u_i(z)$  is the range of utilities of player  $i$ . Moreover,  $R_i^T \leq \sum_{I \in \mathcal{I}_i} R^T(I) \leq |\mathcal{I}_i| \Delta_i \sqrt{|A_i|} \sqrt{T}$ . Therefore,  $\lim_{T \rightarrow \infty} \frac{R_i^T}{T} = 0$ . In two-player zero-sum games, if both players' average regret  $\frac{R_i^T}{T} \leq \epsilon$ , their average strategies  $(\bar{\sigma}_1^T, \bar{\sigma}_2^T)$  over all iterations form a  $2\epsilon$ -equilibrium (Wagh et al., 2009). Some CFR-based variants are proposed to solve large-scale imperfect-information extensive-form games. Some sampling-based CFR variants (Lanctot et al., 2009; Gibson et al., 2012; Schmid et al., 2019) are proposed to effectively solve large-scale games by traversing a subset of the game tree instead of the whole game tree. With the development of deep learning techniques, neural network function approximation can be applied to the CFR algorithm. Deep CFR (Brown et al., 2019), Single Deep CFR (Steinberger, 2019), and Double Neural CFR (Li et al., 2019) are algorithms using deep neural networks to replace the tabular representation.

For the second group, PSRO is a general framework that scales Double Oracle (DO) (McMahan et al., 2003) to large extensive-form games via using reinforcement learning to compute the best response strategy approximately. To make PSRO more effective in solving large-scale games, Pipeline PSRO (P2SRO) (McAleer et al., 2020) is proposed by parallelizing PSRO with convergence guarantees. Extensive-Form Double Oracle (XDO) (McAleer et al., 2021) is a version of PSRO where the restricted game allows mixing population strategies not only at the root of the game but every information set. It can guarantee to converge to an approximate NE in a number of iterations that are linear in the number of information sets, while PSRO may require a number of iterations exponential in the number of information sets. Neural XDO (NXDO) as a neural version of XDO learns approximate best response strategies through any deep reinforcement learning algorithm. Recently, Anytime Double Oracle (ADO) (McAleer et al., 2022), a tabular double oracle algorithm for two-player zero-sum games is proposed to converge to an NE while decreasing exploitability from one iteration to the next. Anytime PSRO (APPRO) as a version of ADO calculates best responses via reinforcement learning algorithms. Except for NEs, we also consider (Coarse) Correlated equilibrium ((C)CE). Joint Policy Space Response Oracles (JPSRO) (Marris et al., 2021) is proposed for training



864 agents in n-player, general-sum extensive-form games, which provably converges to (C)CEs. The  
865 excellent performance of these equilibrium-finding algorithms depends on the interactions with the  
866 actual game environment or a precise simulator. Therefore, these algorithms cannot directly be ap-  
867 plied to the offline EF paradigm. In our paper, we propose a model-based method that can adapt  
868 existing equilibrium finding algorithms to the offline context.

869 **Opponent Modeling.** Opponent modeling algorithm is necessary for multi-agent settings where  
870 secondary agents with competing goals also adapt their strategies, yet it remains challenging because  
871 policies interact with each other and change (He et al., 2016). One simple idea of opponent modeling  
872 is to build a model each time a new opponent or group of opponents is encountered (Zheng et al.,  
873 2018). However, it is infeasible to learn a model every time. A better approach is to represent  
874 an opponent’s policy with an embedding vector. Grover et al. (2018) use a neural network as an  
875 encoder, taking the trajectory of one agent as input. Imitation learning and contrastive learning  
876 are also used to train the encoder. Then, the learned encoder can be combined with reinforcement  
877 learning algorithms by feeding the generated representation into the policy and/or value network.  
878 DRON (He et al., 2016) and DPIQN (Hong et al., 2017) are two algorithms based on DQN, which  
879 use a secondary network that takes observations as input and predicts opponents’ actions. However,  
880 if the opponents can also learn, these methods become unstable. Therefore, it is necessary to take the  
881 learning process of opponents into account. Foerster et al. (2017) propose a method named Learning  
882 with Opponent-Learning Awareness (LOLA), in which each agent shapes the anticipated learning  
883 of the other agents in the environment. Further, the opponents may still be learning continuously  
884 during execution. Therefore, Al-Shedivat et al. (2017) propose a method based on a meta-policy  
885 gradient named Meta-MPG. It uses trajectories from current opponents to perform multiple meta-  
886 gradient steps and constructs a policy that favors updating the opponents. Meta-MAPG (Kim et al.,  
887 2021) extends Meta-MPG by including an additional term that accounts for the impact of the agent’s  
888 current policy on the future policies of opponents, similar to LOLA. Yu et al. (2021b) propose model-  
889 based opponent modeling (MBOM), which employs the environment model to adapt to various  
890 opponents. In our offline EF paradigm, our goal is to compute the equilibrium strategy based on the  
891 offline dataset. Applying opponent modeling is not enough for the offline EF paradigm since it only  
892 aims at computing the best response strategy instead of the equilibrium strategy.

892 **Empirical Game Theoretic Analysis.** Empirical game theoretic analysis (EGTA) is an empirical  
893 methodology that bridges the gap between game theory and simulation for practical strategic reason-  
894 ing (Wellman, 2006). In EGTA, game models are iteratively extended through a process of generat-  
895 ing new strategies based on learning from experience with prior strategies. The strategy exploration  
896 problem (Jordan et al., 2010) that how to efficiently assemble an efficient portfolio of policies for  
897 EGTA is the most challenging problem. Schwartzman & Wellman (2009b) deploy tabular RL as a  
898 best-response oracle in EGTA for strategy generation. They also build the general problem of strat-  
899 egy exploration in EGTA and investigate whether better options exist beyond best-responding to an  
900 equilibrium (Schwartzman & Wellman, 2009a). Investigation of strategy exploration was advanced  
901 significantly by the introduction of the Policy Space Response Oracle (PSRO) framework (Lanctot  
902 et al., 2017) which is a flexible framework for iterative EGTA, where at each iteration, new strategies  
903 are generated through reinforcement learning. Note that when employing NE as the meta-strategy  
904 solver, PSRO reduces to the double oracle (DO) algorithm (McMahan et al., 2003). In EGTA, a  
905 space of strategies is examined through simulation, which means that it needs a simulator, and the  
906 policies are known in advance. However, in the offline EF paradigm, only an offline dataset is  
907 provided. Therefore, techniques in EGTA cannot be directly applied to the offline EF paradigm.

907 **Offline Reinforcement Learning.** Offline reinforcement learning (offline RL) is a *data-driven*  
908 paradigm that learns exclusively from static datasets of previously collected interactions, mak-  
909 ing it feasible to extract policies from large and diverse training datasets (Levine et al., 2020).  
910 This paradigm can be extremely valuable in settings where online interaction is impractical, either  
911 because data collection is expensive or dangerous (e.g., in robotics (Singh et al., 2021), educa-  
912 tion (Singla et al., 2021), healthcare (Liu et al., 2020), and autonomous driving (Kiran et al., 2022)).  
913 Therefore, efficient offline RL algorithms have a much broader range of applications than online  
914 RL and are particularly appealing for real-world applications (Prudencio et al., 2022). Due to its  
915 attractive characteristics, there have been a lot of recent studies. Here, we can divide the research of  
916 offline RL into two categories: model-based algorithm and model-free algorithm.

917 Model-free offline RL algorithms learn a good policy directly from the offline dataset. To do this,  
there are two types of algorithms: actor-critic and imitation learning methods. Those actor-critic

918 algorithms focus on implementing policy regularization and value regularization based on existing  
919 reinforcement learning algorithms. [Haarnoja et al. \(2018\)](#) propose soft actor-critic (SAC) by adding  
920 an entropy regularization term to the policy gradient objective. This work mainly focuses on policy  
921 regularization. For the research of value regularization, an offline RL method named Constrained  
922 Q-Learning (CQL) ([Kumar et al., 2020](#)) learns a lower bound of the true Q-function by adding value  
923 regularization terms to its objective. Another line of model-free offline RL research is imitation  
924 learning which mimics the behavior policy based on the offline dataset. [Chen et al. \(2020\)](#) propose  
925 a method named Best-Action Imitation Learning (BAIL), which fits a value function, then uses it to  
926 select the best actions. Meanwhile, [Siegel et al. \(2020\)](#) propose a method that learns an Advantage-  
927 weighted Behavior Model (ABM) and uses it as a prior in performing Maximum a-posteriori Policy  
928 Optimization (MPO) ([Abdolmaleki et al., 2018](#)). It consists of multiple iterations of policy evalua-  
929 tion and prior learning until they finally perform a policy improvement step using their learned prior  
930 to extracting the best possible policy.

931 Model-based algorithms rely on the offline dataset to learn a dynamics model or a trajectory dis-  
932 tribution used for planning. The trajectory distribution induced by models is used to determine the  
933 best set of actions to take at each given time step. [Kidambi et al. \(2020\)](#) propose a method named  
934 Model-based Offline Reinforcement Learning (MOREL), which measures their model’s epistemic  
935 uncertainty through an ensemble of dynamics models. Meanwhile, [Yu et al. \(2020\)](#) propose an-  
936 other method named Model-based Offline Policy Optimization (MOPO), which uses the maximum  
937 prediction uncertainty from an ensemble of models. Concurrently, [Matsushima et al. \(2020\)](#) pro-  
938 pose the BehaviorREGularized Model-ENsemble (BREMEN) method, which learns an ensemble of  
939 models of the behavior MDP, as opposed to a pessimistic MDP. In addition, it implicitly constrains  
940 the policy to be close to the behavior policy through trust-region policy updates. More recently,  
941 [Yu et al. \(2021a\)](#) proposed a method named Conservative Offline Model-Based policy Optimization  
942 (COMBO), a model-based version of CQL. The main advantage of COMBO concerning MOREL  
943 and MOPO is that it removes the need for uncertainty quantification in model-based offline RL ap-  
944 proaches, which is challenging and often unreliable. However, these above offline RL algorithms  
945 cannot be directly applied to the offline EF paradigm, which we have described in Section 2 and  
946 experimental results empirically verify this claim.

947  
948  
949  
950  
951  
952  
953  
954  
955  
956  
957  
958  
959  
960  
961  
962  
963  
964  
965  
966  
967  
968  
969  
970  
971

C DATASETS

C.1 DATASET FORMAT

To better introduce the format of our offline EF dataset, we provide an example to show the composition of the offline dataset. According to the data format introduced in the main paper, the data point would be  $(s_t, a_t, s_{t+1}, u_{t+1}, d_{t+1})$  and  $s_t = (I_1^t, I_2^t, \dots, I_n^t, GI, p^t, A(I_{p^t}^t))$ . Especially, if the  $s_{t+1}$  is the terminal state, i.e.,  $d_{t+1} = 1$ , then we define the  $p^{t+1} = -1$  to identify that there is no player need to decide this state. Fig. 12 shows one two-player imperfect-information extensive-form game  $\mathcal{G}$ .  $I_1$  and  $I_2$  are information set for Player 1 and Player 2, respectively. If an offline dataset  $\mathcal{D}$  covers all state-action pairs, then  $\mathcal{D}$  would include the following data points:

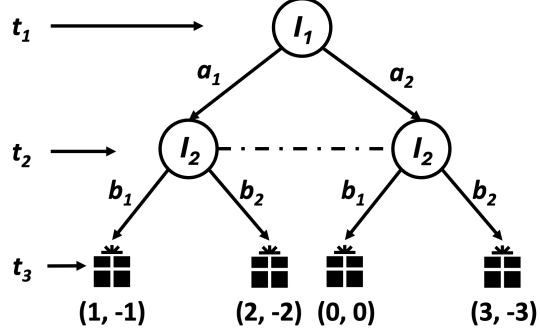


Figure 12: An example game.

$((I_1^{t_1} = I_1, I_2^{t_1} = \emptyset, GI = \emptyset, 1, \{a_1, a_2\}), a_1, (I_1^{t_2} = I_1 a_1, I_2^{t_2} = I_2, GI = \emptyset, 2, \{b_1, b_2\}), (0, 0), 0),$   
 $((I_1^{t_1} = I_1, I_2^{t_1} = \emptyset, GI = \emptyset, 1, \{a_1, a_2\}), a_2, (I_1^{t_2} = I_1 a_2, I_2^{t_2} = I_2, GI = \emptyset, 2, \{b_1, b_2\}), (0, 0), 0),$   
 $((I_1^{t_2} = I_1 a_1, I_2^{t_2} = I_2, GI = \emptyset, 2, \{b_1, b_2\}), b_1, (I_1^{t_3} = I_1 a_1, I_2^{t_3} = I_2 b_1, GI = \emptyset, -1, \emptyset), (1, -1), 1),$   
 $((I_1^{t_2} = I_1 a_1, I_2^{t_2} = I_2, GI = \emptyset, 2, \{b_1, b_2\}), b_2, (I_1^{t_3} = I_1 a_1, I_2^{t_3} = I_2 b_2, GI = \emptyset, -1, \emptyset), (2, -2), 1),$   
 $((I_1^{t_2} = I_1 a_2, I_2^{t_2} = I_2, GI = \emptyset, 2, \{b_1, b_2\}), b_1, (I_1^{t_3} = I_1 a_2, I_2^{t_3} = I_2 b_1, GI = \emptyset, -1, \emptyset), (0, 0), 1),$   
 $((I_1^{t_2} = I_1 a_2, I_2^{t_2} = I_2, GI = \emptyset, 2, \{b_1, b_2\}), b_2, (I_1^{t_3} = I_1 a_2, I_2^{t_3} = I_2 b_2, GI = \emptyset, -1, \emptyset), (3, -3), 1).$

We can find that in states  $(I_1^{t_2} = I_1 a_1, I_2^{t_2} = I_2, GI = \emptyset, 2, \{b_1, b_2\})$  and  $(I_1^{t_2} = I_1 a_2, I_2^{t_2} = I_2, GI = \emptyset, 2, \{b_1, b_2\})$ , the information set for Player 2 is the same, as shown in the game tree. Since our dataset is collected from the perspective of the game, we can still distinguish them through the game information of other players and the game information  $GI$ . Note that there is no chance node in this game, the game information  $GI$  is an empty set here. If there is a chance node in the game, the results of the chance node would be recorded into game information  $GI$  within the game state and we can distinguish these game states through game information  $GI$ .

C.2 VISUALIZATION

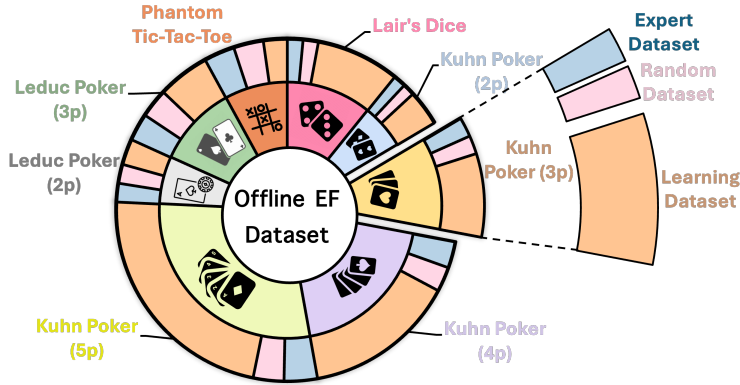


Figure 13: Visualization of the offline EF dataset

Fig. 13 shows a full view of our offline EF dataset. In our offline EF dataset, we collected data for eight games, including two-player Kuhn poker, three-player Kuhn poker, four-player Kuhn poker,

1026 five-player Kuhn poker, two-player Leduc poker, three-player Leduc poker, Phantom Tic-Tac-Toe,  
 1027 and Lair’s Dice games. For each game, we generated three datasets, a random dataset, an expert dat-  
 1028 set, and a learning dataset, following our data collection methods. To validate the diversity of these  
 1029 collected offline datasets and gain insights into them, we also introduce a visualization method for  
 1030 comparing them. Firstly, we generate the game tree for the corresponding game. Subsequently, we  
 1031 traverse the game tree using depth-first search (DFS) (Tarjan, 1972) and assign an index to each leaf  
 1032 node based on the DFS results. Then, we count the frequency of each leaf node within the dataset.  
 1033 The reason why we do this is that each leaf node represents a unique sampled trajectory originating  
 1034 from the root node of the game tree. As a result, the frequency of leaf nodes can effectively capture  
 1035 the distribution of the dataset. Finally, these frequency data can be plotted to visualize. Fig. 14  
 1036 visualizes some datasets of some games. From these figures, we can find that in the random dataset,  
 1037 the frequency of leaf nodes is nearly uniform, whereas, in the expert dataset, the frequency of leaf  
 1038 nodes is uneven. The distribution of the learning dataset and the hybrid dataset falls between that of  
 1039 the expert dataset and the random dataset. These observations confirm that the distribution of these  
 1040 datasets differs, thus validating the diversity of our proposed offline datasets.

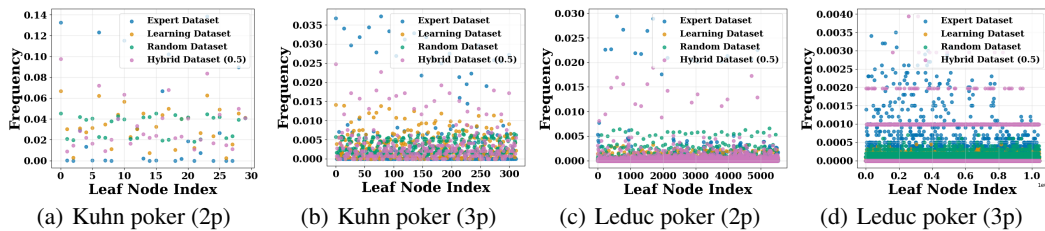


Figure 14: Frequency of leaf node in different offline datasets.



D THEORETICAL ANALYSIS

In this section, we provide a comprehensive theoretical analysis of the offline EF paradigm and our BOMB framework to facilitate the understanding of the offline EF paradigm and BOMB framework. We first provide the minimal dataset assumption that is sufficient to compute the equilibrium strategy in the offline setting. Then we provide a general generalization bound for training neural network models. Finally, we give the performance guarantee for our algorithm. In the following sections, we assume that all extensive-form games discussed here are perfect recall and timetable.

D.1 MINIMAL DATASET ASSUMPTION FOR OFFLINE EF

As demonstrated in offline RL papers (Rashidinejad et al., 2021; Xie et al., 2021), a dataset coverage condition over the optimal policy is sufficient for offline learning. Therefore, it is straightforward to extend this dataset coverage assumption to the offline EF paradigm. In the main paper, we have proved that the dataset generated by the equilibrium strategy is not sufficient for computing the equilibrium strategy in an offline manner by providing a counter-example. Furthermore, we also provide another dataset assumption related to the equilibrium strategy, shown in the following assumption.

**Assumption D.1.** (Single Strategy Coverage) The offline dataset  $\mathcal{D}$  is said to be *single strategy coverage* if the equilibrium strategy profile  $\sigma^*$  is covered by the offline dataset  $\mathcal{D}$ , i.e., for each player  $i$ , each information set  $I_i$ , and action  $a_i$  with  $\sigma_i^*(I_i, a_i) > 0$ , there is a corresponding state-action pair  $(s_t, a_i)$  in  $\mathcal{D}$ .

Subsequently, a question arises: *is the single strategy coverage assumption also sufficient for computing equilibrium strategy in the offline setting?* We employ the following theorem to answer this question and elucidate the rationale behind this.

**Theorem D.2.** *Single strategy coverage assumption over offline dataset  $\mathcal{D}$  is not sufficient for computing computing an  $\epsilon$ -equilibrium for an arbitrarily small  $\epsilon$  in the offline setting.*

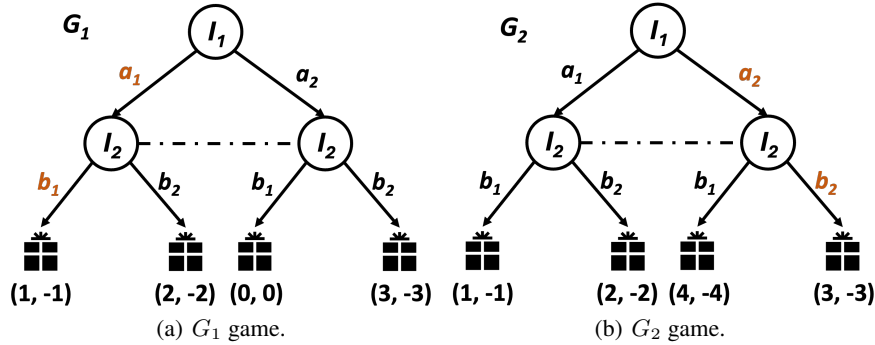


Figure 15: Counter-example for proving Theorem D.2.

*Proof.* We prove this theorem by providing a counter-example. To this end, we consider two two-player IIEFGs  $G_1$  and  $G_2$ , represented in Fig. 15. We can easily find that the NE of the game  $G_1$  is strategy profile  $\sigma^1 = (\sigma_1^1, \sigma_2^1) = (\{I_1 : a_1\}, \{I_2 : b_1\})$ , i.e., Player 1 plays  $a_1$  at information set  $I_1$  and Player 2 plays  $b_1$  at information set  $I_2$ . The NE of the game  $G_2$  is strategy profile  $\sigma^2 = (\sigma_1^2, \sigma_2^2) = (\{I_1 : a_2\}, \{I_2 : b_2\})$ . Now we consider an offline dataset  $\mathcal{D}$  which is generated using a strategy profile  $\sigma_{\mathcal{D}}$ . The  $\sigma_{\mathcal{D}}$  is set to be the uniform distribution over the strategy profiles  $\sigma^1$  and  $\sigma^2$ , which means that dataset  $\mathcal{D}$  covers both  $\sigma^1$  and  $\sigma^2$ . Therefore, the offline dataset  $\mathcal{D}$  satisfies the single strategy coverage assumption for these two games  $G_1$  and  $G_2$ . However, no algorithm can distinguish these two extensive-form games only based on dataset  $\mathcal{D}$  since these two games are both consistent on dataset  $\mathcal{D}$ . In conclusion, the single strategy converges assumption is not sufficient for computing an  $\epsilon$ -equilibrium for an arbitrarily small  $\epsilon$  in the offline setting.  $\square$

From the above proof, we know that the single strategy coverage assumption is sufficient for computing the optimal strategy in the offline RL setting while it is not sufficient for computing an NE

strategy in the offline setting. The intuition behind this is that in an offline RL setting, we can easily use the data of two actions to decide which action is better, whereas, in the offline EF paradigm, we cannot use data from only two action pairs to know which action pair is closer to NE, because identifying NE requires other action pairs as inferences. Based on this analysis, Cui & Du (2022) provide a minimal coverage assumption which is sufficient for computing an NE strategy in the two-player zero-sum Markov games, which is defined as follows,

**Assumption D.3.** (Deterministic Unilateral Coverage) For all deterministic strategy  $\sigma_i$  for player  $i$ ,  $(\sigma_i, \sigma_{-i}^*)$  are covered by the dataset, where  $(\sigma_1^*, \dots, \sigma_n^*)$  is one NE strategy.

**Assumption D.4.** (Unilateral Coverage) For all (possible stochastic) strategy  $\sigma_i$  for all player  $i$ ,  $(\sigma_i, \sigma_{-i}^*)$  are covered by the dataset, where  $(\sigma_1^*, \dots, \sigma_n^*)$  is one NE strategy.

Note that deterministic unilateral coverage assumption is equivalent to unilateral coverage assumption. The intuition behind this is that any mixed strategy can be represented by a combination of deterministic strategies. Therefore, if all deterministic strategies are covered by the dataset, then all mixed strategies are also covered. Based on this finding, in the following proof, we only consider all deterministic strategies. Previously, Cui & Du (2022) established that unilateral coverage assumption is the minimal sufficient condition for computing an NE strategy in the two-player zero-sum Markov games. However, this unilateral coverage assumption over the offline dataset is **not sufficient** for our model-based method to compute the equilibrium strategy in the offline setting. We formally proved this limitation through the following theorem.

**Theorem D.5.** *The unilateral coverage assumption over offline dataset  $\mathcal{D}$  is not sufficient for our model-based method to converge to an  $\epsilon$ -equilibrium for an arbitrarily small  $\epsilon$  in the offline setting.*

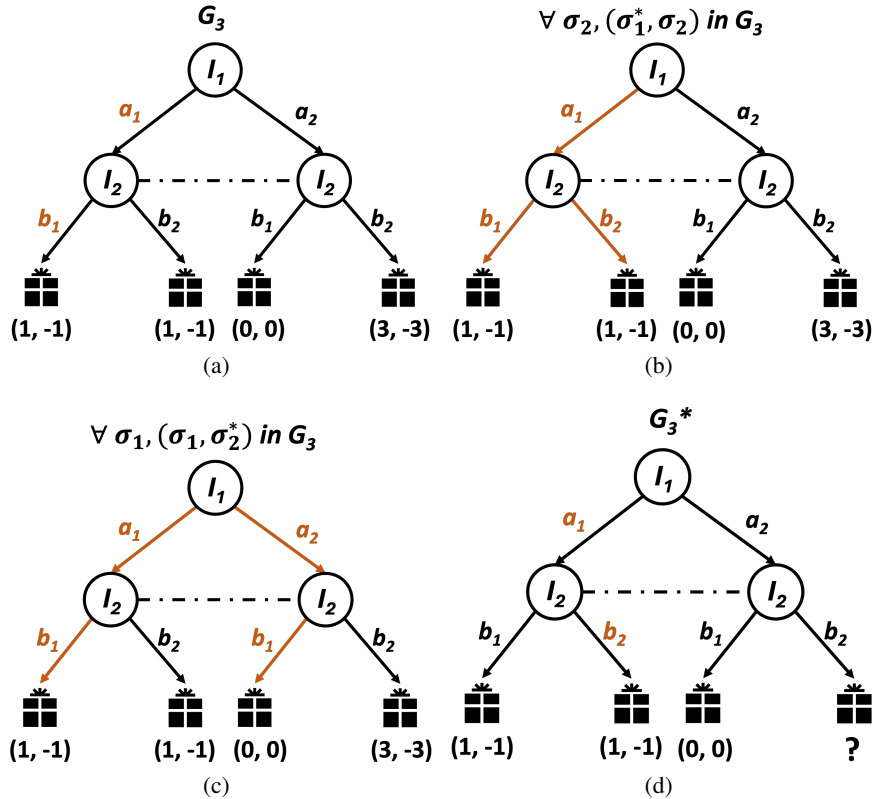


Figure 16: Counter-example for proving Theorem D.5.

*Proof.* We prove this theorem by providing a counter-example. First, we consider an IIEFG  $M_3$ , represented in Fig. 16(a). We can easily find that the NE strategy of game  $G_3$  is strategy profile  $\sigma^* = (\sigma_1, \sigma_2) = (\{I_1 : a_1\}, \{I_2 : b_1\})$ . To build a dataset  $\mathcal{D}$  satisfying the unilateral coverage assumption, the dataset needs to cover  $(\sigma_1^*, \sigma_2)$  for all deterministic strategy  $\sigma_2$  and  $(\sigma_1, \sigma_2^*)$  for

all deterministic strategy  $\sigma_1$ . We show the state-action pairs covered by these strategy profiles in Figs. 16(b)-16(c). It means that if the dataset  $\mathcal{D}$  satisfies the unilateral coverage assumption, then the dataset  $\mathcal{D}$  would cover these state-action pairs marked by these orange lines. When applying our model-based method on the dataset  $\mathcal{D}$ , the first step is to train an environment model based on the dataset  $\mathcal{D}$ . Assume that the environment model can be trained well (i.e., Assumption 4.3 holds) which means that the environment model can precisely represent all game information in the dataset. Therefore, the game represented by the trained environment model would be  $G_3^*$  in Fig. 16(d). Note that there is some missing data in the game. Although our trained environment model can give approximate results for these missing data, it may result in a different equilibrium strategy. For example, if the missing value in  $G_3^*$  is  $(0, 0)$  or  $(-1, 1)$ , then the strategy profile  $\sigma = (\sigma_1, \sigma_2) = (\{I_1 : a_1\}, \{I_2 : b_2\})$  would be the NE strategy of game  $G_3^*$ . However, the strategy profile  $\sigma$  is not the NE strategy for the original game  $G_3$ . Therefore, the unilateral coverage assumption over the dataset is not sufficient for our model-based method to converge to an  $\epsilon$ -equilibrium for an arbitrarily small  $\epsilon$ .  $\square$

To guarantee the convergence of our model-based method, we provide a minimal dataset coverage assumption for our model-based method to converge to the equilibrium strategy of the original game under the offline setting.

**Definition D.6** (Definition 4.1). An offline dataset  $\mathcal{D}$  is said to be a *uniform coverage* of an IIEFG  $\mathcal{G}$  if and only if the offline dataset  $\mathcal{D}$  covers all possible state-action pairs. Formally,  $(s_t, a_t, s_{t+1}, u_{t+1}, d_{t+1}), \forall s_t, a_t \in A(s_t)$  and  $s_{t+1} \in T(s_t, a_t)$  where  $T$  is the transition function of game  $\mathcal{G}$ .

**Theorem D.7** (Theorem 4.4). Let  $\sigma_{MB(\mathcal{D})}$  be the strategy profile learned by our *model-based algorithm* based on the offline dataset  $\mathcal{D}$  with sufficient data under Assumption 4.3. Then,  $\sigma_{MB(\mathcal{D})}$  is guaranteed to be an  $\epsilon$ -equilibrium strategy of the IIEFG  $\mathcal{G}$  if and only if  $\mathcal{D}$  is a uniform coverage of  $\mathcal{G}$  and  $\sigma_{MB(\mathcal{D})}$  is an  $\epsilon$ -equilibrium strategy for the trained environment model within the model-based algorithm.

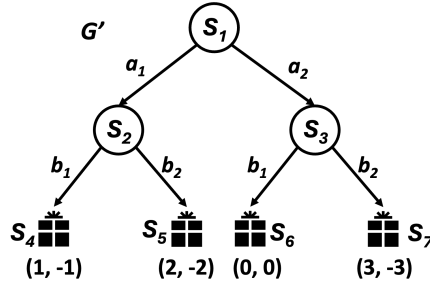


Figure 17:  $G'$  Game.

*Proof.* From the example in the proof of Theorem D.5, we find that a slight violation of the uniform coverage assumption, i.e., only one state-action pair is missing, will impede the computation of the equilibrium strategy using our model-based method. In other words, any state-action pair that is not covered by the dataset may cause failure in computing the equilibrium strategy of the original game using our model-based method.

Next, we need to prove that the dataset satisfying the uniform coverage assumption can guarantee the convergence to the equilibrium strategy of the original game using our model-based method. In our model-based method, we need to train an environment model based on the offline dataset. Therefore, to prove the convergence guarantee under the uniform coverage dataset assumption, we need to verify whether the game reconstructed from the dataset satisfying the uniform coverage assumption is the same as the original game. Here, we reuse the example in the App. C.1. In that example, the offline dataset  $\mathcal{D}$  of the IIEFG  $\mathcal{G}$  covers all state-action pairs. Therefore, the offline dataset  $\mathcal{D}$  satisfies the uniform coverage dataset assumption. From the offline dataset  $\mathcal{D}$ , we can

1242 easily rebuild the game  $G'$ , as shown in Fig. 17. In the game  $G'$ ,

$$1243 \quad S_1 = (I_1^{t_1} = I_1, I_2^{t_1} = \emptyset, GI = \emptyset, 1, \{a_1, a_2\}), S_2 = (I_1^{t_2} = I_1 a_1, I_2^{t_2} = I_2, GI = \emptyset, 2, \{b_1, b_2\}),$$

$$1244 \quad S_3 = (I_1^{t_2} = I_1 a_2, I_2^{t_2} = I_2, GI = \emptyset, 2, \{b_1, b_2\}), S_4 = (I_1^{t_3} = I_1 a_1, I_2^{t_3} = I_2 b_1, GI = \emptyset, -1, \emptyset),$$

$$1245 \quad S_5 = (I_1^{t_3} = I_1 a_1, I_2^{t_3} = I_2 b_2, GI = \emptyset, -1, \emptyset), S_6 = (I_1^{t_3} = I_1 a_2, I_2^{t_3} = I_2 b_1, GI = \emptyset, -1, \emptyset),$$

$$1246 \quad S_7 = (I_1^{t_3} = I_1 a_2, I_2^{t_3} = I_2 b_2, GI = \emptyset, -1, \emptyset).$$

1249 Especially, for game states  $S_2$  and  $S_3$ , the player acting is both Player 2 and the information set for  
 1250 Player 2 is the same. Therefore, these two game states correspond to different game nodes under  
 1251 the same information set. Although Player 2 cannot distinguish these two game states, from the  
 1252 perspective of the game, we can still distinguish them by the information set of Player 1. Particularly,  
 1253 if there is a chance node in the game, the result of the chance node would be recorded in  $GI$  within  
 1254 the game state  $S$ . Therefore, we can still distinguish these game states by game information  $GI$ .  
 1255 Since the dataset satisfying the uniform coverage assumption covers all state-action pairs, the links  
 1256 between game states can be built following these data points in the dataset. According to Assumption  
 1257 4.3, the error in training the environment game model based on  $\mathcal{D}$  can be considered negligible.  
 1258 Consequently, the trained environment game model is identical to the original game  $\mathcal{G}$ , as the dataset  
 1259  $\mathcal{D}$  provides full coverage of all state transitions. Therefore, we can find that the reconstructed game  
 1260 tree has the same game states and the same transition function as the original game, thereby the same  
 1261 equilibrium strategy. Therefore, our reconstructed game model can provide the same information  
 1262 as the underlying game of the offline dataset. Then applying our model-based equilibrium finding  
 1263 algorithm to the reconstructed game model definitely can converge to the equilibrium strategy of  
 1264 the underlying game in the offline setting. Formally, if  $\sigma_{MB(\mathcal{D})}$  is an  $\epsilon$ -equilibrium strategy for the  
 1265 trained environment game model, it is also an  $\epsilon$ -equilibrium strategy for the original game  $\mathcal{G}$ .  $\square$

1266 So far, we have proved that the uniform dataset coverage assumption is sufficient for our model-  
 1267 based method to converge to the equilibrium strategy under the offline setting. For our behavior  
 1268 cloning method, these dataset coverage assumptions may not be sufficient to converge to the equi-  
 1269 librium strategy since its performance mainly depends on the underlying behavior strategy of the  
 1270 dataset. In the following theorem, we provide a minimal dataset coverage assumption for our behav-  
 1271 ior cloning method to converge to the equilibrium strategy in the offline setting.

1272 **Definition D.8** (Definition 4.2). An offline dataset  $\mathcal{D}$  is said to be an  $\epsilon$ -equilibrium coverage over an  
 1273 IIEFG  $\mathcal{G}$  if and only if its underlying behavior strategy  $\sigma_{\mathcal{D}}$  satisfies  $\text{GAP}(\sigma_{\mathcal{D}}, \sigma^*) < \epsilon$ , where  $\sigma_{\mathcal{D}}$   
 1274 is defined as  $\sigma_{\mathcal{D}}(s_t, a_t) = \frac{C(s_t, a_t)}{C(s_t)}$  and  $\sigma_{\mathcal{D}}(s_t, a_t) > 0$  for all  $s_t$  and  $a_t \in A(s_t)$ , with  $C(s_t, a_t)$  and  
 1275  $C(s_t)$  denoting the counts of data points containing  $(s_t, a_t)$  and  $s_t$  in  $\mathcal{D}$ , respectively.

1276 This definition ensures that the unique correspondence relationship between the equilibrium-covered  
 1277 dataset and the equilibrium strategy. Specifically, the dataset is generated by the equilibrium strategy  
 1278 and the strategy represented by the dataset would be the same as the equilibrium strategy.

1280 **Theorem D.9** (Theorem 4.5). Let  $\sigma_{BC(\mathcal{D})}$  be the strategy profile learned by our **behavior cloning**  
 1281 **algorithm** based on the offline dataset  $\mathcal{D}$  with sufficient data under Assumption 4.3. Then  $\sigma_{BC(\mathcal{D})}$   
 1282 is guaranteed to be an  $\epsilon$ -equilibrium strategy of IIEFG  $\mathcal{G}$  if and only if the offline dataset  $\mathcal{D}$  is an  
 1283  $\epsilon$ -equilibrium coverage of the IIEFG  $\mathcal{G}$ .

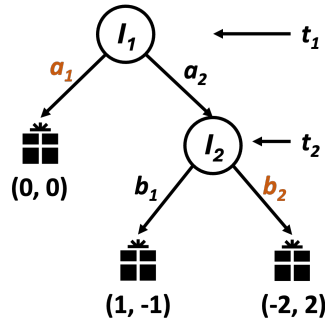


Figure 18: Game example.

1296 *Proof.* According to Assumption 4.3, the error in training the behavior cloning strategy  $\sigma_{BC(\mathcal{D})}$   
 1297 from the dataset  $\mathcal{D}$  is negligible. Therefore, by the behavior cloning process,  $\sigma_{BC(\mathcal{D})}$  is identical  
 1298 to the behavior strategy underlying  $\mathcal{D}$ , i.e.,  $\sigma_{BC(\mathcal{D})} = \sigma_{\mathcal{D}}$ . Consequently, if  $\mathcal{D}$  is an  $\epsilon$ -  
 1299 equilibrium coverage of  $\mathcal{G}$ , then  $\sigma_{BC(\mathcal{D})}$  is an  $\epsilon$ -equilibrium strategy for the IIEFG  $\mathcal{G}$ , and vice  
 1300 visa, as  $\text{GAP}(\sigma_{\mathcal{D}}, \sigma^*) < \epsilon$  if and only if  $\text{GAP}(\sigma_{BC(\mathcal{D})}, \sigma^*) < \epsilon$ . Next, we prove that any slight  
 1301 violation of these conditions would invalidate the convergence result.

1302 Here, we reuse the example in Section 4.2, as shown in Fig. 18. Note that the NE strategy of the game  
 1303 is a pure strategy, i.e.,  $\sigma^* = (\sigma_1^*, \sigma_2^*) = (\{I_1 : a_1\}, \{I_2 : b_2\})$ . If we use this equilibrium strategy to  
 1304 generate the offline dataset  $\mathcal{D}$ , then  $\mathcal{D}$  would only include the data point  $((I_1^{t_1} = I_1, I_2^{t_1} = \emptyset, GI =$   
 1305  $\emptyset, 1, \{a_1, a_2\}), a_1, (I_1^{t_2} = I_1 a_1, I_2^{t_2} = \emptyset, GI = \emptyset, -1, \emptyset), (0, 0), 1)$ . We cannot get the equilibrium  
 1306 strategy only from  $\mathcal{D}$ . In this example game, the offline dataset  $\mathcal{D}$  is generated by a pure equilibrium  
 1307 strategy instead of a fully mixed equilibrium strategy, and the behavior cloning method cannot get  
 1308 the equilibrium strategy from the offline dataset  $\mathcal{D}$  since there is no information about Player 2.  
 1309 Another example is the dataset  $D'$  covering the equilibrium strategy  $\sigma^*$ , i.e., the  $D'$  includes the  
 1310 data points

$$1311 ((I_1^{t_1} = I_1, I_2^{t_1} = \emptyset, GI = \emptyset, 1, \{a_1, a_2\}), a_1, (I_1^{t_2} = I_1 a_1, I_2^{t_2} = \emptyset, GI = \emptyset, -1, \emptyset), (0, 0), 1),$$

$$1312 ((I_1^{t_1} = I_1, I_2^{t_1} = \emptyset, GI = \emptyset, 1, \{a_1, a_2\}), a_2, (I_1^{t_2} = I_1 a_2, I_2^{t_2} = I_2, GI = \emptyset, 2, \{b_1, b_2\}), (0, 0), 0),$$

$$1313 ((I_1^{t_2} = I_1 a_2, I_2^{t_2} = I_2, GI = \emptyset, 2, \{b_1, b_2\}), b_2, (I_1^{t_1} = I_1 a_2, I_2^{t_1} = I_2 b_2, GI = \emptyset, -1, \emptyset), (-2, 2), 1).$$

1316 To cover the equilibrium strategy of Player 2, the data point  $((I_1^{t_1} = I_1, I_2^{t_1} = \emptyset, GI =$   
 1317  $\emptyset, 1, \{a_1, a_2\}), a_2, (I_1^{t_2} = I_1 a_2, I_2^{t_2} = I_2, GI = \emptyset, 2, \{b_1, b_2\}), (0, 0), 0)$  should also be visited.  
 1318 Although  $D'$  covers the equilibrium strategy,  $D'$  does not satisfy the  $\epsilon$ -equilibrium coverage assumption  
 1319 since the  $D'$  does not contain the data point  $((I_1^{t_1} = I_1, I_2^{t_1} = \emptyset, GI = \emptyset, 1, \{a_1, a_2\}), a_2, (I_1^{t_2} =$   
 1320  $I_1 a_2, I_2^{t_2} = I_2, GI = \emptyset, 2, \{b_1, b_2\}), (0, 0), 0)$ . Therefore, a slight violation of the equilibrium  
 1321 coverage assumption would cause failure in computing the  $\epsilon$ -equilibrium strategy of the original  
 1322 game using our behavior cloning method. In conclusion, the equilibrium coverage assumption is the  
 1323 minimal dataset coverage assumption that guarantees the convergence to the equilibrium strategy  
 1324 of the original game using our behavior cloning method. Formally,  $\sigma_{BC(\mathcal{D})}$  is guaranteed to be an  
 1325  $\epsilon$ -equilibrium strategy of IIEFG  $\mathcal{G}$  if and only if the offline dataset  $\mathcal{D}$  is an  $\epsilon$ -equilibrium coverage  
 1326 of the IIEFG  $\mathcal{G}$ .  $\square$

1329 **Theorem D.10** (Theorem 4.6). *Let  $\sigma_{BOMB(\mathcal{D})}$  represent the strategy profile learned by our **BOMB**  
 1330 **algorithm** based on the offline dataset  $\mathcal{D}$  with sufficient data under Assumption 4.3,  $\sigma_{\mathcal{D}}$  represent  
 1331 the underlying behavior strategy of  $\mathcal{D}$  and  $\sigma^*$  represent the equilibrium strategy of IIEFG  $\mathcal{G}$ . Then  
 1332 the gap between  $\sigma_{BOMB(\mathcal{D})}$  and  $\sigma^*$  is at most equal to, or smaller than, the gap between  $\sigma_{\mathcal{D}}$  and  
 1333  $\sigma^*$ , i.e.,  $\text{GAP}(\sigma_{BOMB(\mathcal{D})}, \sigma^*) \leq \text{GAP}(\sigma_{\mathcal{D}}, \sigma^*)$ .*

1335 *Proof.* According to Assumption 4.3, the error in training the behavior cloning strategy  $\sigma_{BC(\mathcal{D})}$   
 1336 from the dataset  $\mathcal{D}$  is negligible. Therefore, by the behavior cloning process,  $\sigma_{BC(\mathcal{D})}$  is identical  
 1337 to the behavior strategy underlying  $\mathcal{D}$ , i.e.,  $\sigma_{BC(\mathcal{D})} = \sigma_{\mathcal{D}}$ . Then  $\text{GAP}(\sigma_{BOMB(\mathcal{D})}, \sigma^*) =$   
 1338  $\text{GAP}(\sigma_{\mathcal{D}}, \sigma^*)$  if  $\alpha = 1$  in our **BOMB algorithm**.

1339 If the dataset satisfies the uniform coverage, by Theorem 4.4,  $\text{GAP}(\sigma_{BOMB(\mathcal{D})}, \sigma^*) \leq$   
 1340  $\text{GAP}(\sigma_{\mathcal{D}}, \sigma^*)$  if  $\alpha = 0$  in our **BOMB algorithm**.

1342 Therefore, in general case,  $\text{GAP}(\sigma_{BOMB(\mathcal{D})}, \sigma^*) \leq \text{GAP}(\sigma_{\mathcal{D}}, \sigma^*)$ .  $\square$

## 1344 D.2 GENERALIZATION BOUND FOR TRAINING MODEL

1346 As described in the main paper, to conduct the BOMB framework, we need to train one behavior  
 1347 cloning policy and an environment model which are both neural network models. Furthermore,  
 1348 these two models are trained in a supervised learning manner with different loss functions based  
 1349 on the offline EF dataset. Here, we provide a general generalization bound for training such neural  
 network models facilitating the following analysis of the BOMB framework.



As we know, the supervised learning framework includes a data-generation distribution  $\sigma$ , a hypothesis class  $\mathcal{H}$  of the neural network approximator, a training dataset  $\mathcal{D}$ , and evaluation metrics to evaluate the performance of any approximator. Here, we can use the loss function  $l$  to evaluate the performance of any approximation. The learning framework aims to minimize the true risk function  $L_\sigma(h)$  which is the expected loss function of  $h \in \mathcal{H}$  under the distribution  $\sigma$ ,

$$L_\sigma(h) = \mathbb{E}_{d \sim \sigma}[l(h(d), d)].$$

Accordingly, the empirical risk function  $L_{\mathcal{D}}(h)$  on the training dataset  $\mathcal{D}$  can be defined as:

$$L_{\mathcal{D}} = \frac{1}{|\mathcal{D}|} \sum_{d \sim \mathcal{D}} [l(h(d), d)].$$

To get a generalization bound, we use an auxiliary lemma from (Shalev-Shwartz & Ben-David, 2014). Therefore, we can measure the capacity of the composition function class  $l \circ \mathcal{H}$  using the empirical Rademacher complexity on the training set  $\mathcal{D}$  with size  $m$ , which is defined as:

$$\mathcal{R}_{\mathcal{D}}(l \circ \mathcal{H}) = \frac{1}{m} \mathbb{E}_{\mathbf{x} \sim \{+1, -1\}^m} \left[ \sup_{h \in \mathcal{H}} \sum_{i=1}^m x_i \cdot l(h(d_i), d_i) \right]$$

where  $\mathbf{x}$  is distributed i.i.d. according to uniform distribution in  $\{+1, -1\}$ . Before providing the generalization bound, we first provide the distance between two different approximators and one common theorem to facilitate the proof of the generalization bound.

**Definition D.11.** (*r-cover*) We say function class  $\mathcal{H}_r$  *r-cover*  $\mathcal{H}$  under  $\ell_{\infty,1}$ -distance if  $\forall h, h_r \in \mathcal{H}$ , there exists  $h_r$  in  $\mathcal{H}_r$  such that  $\|h - h_r\|_{\infty,1} = \max_{x \in \mathcal{D}} \|h(x) - h_r(x)\|_1 \leq r$ .

**Definition D.12.** (*r-covering number*) The *r-covering number* of  $\mathcal{H}$ ,  $\mathcal{N}_{\infty,1}(\mathcal{H}, r)$ , is the cardinality of the smallest function class  $\mathcal{H}_r$  that *r-covers*  $\mathcal{H}$  under  $\ell_{\infty,1}$ -distance.

**Theorem D.13.** (Shalev-Shwartz & Ben-David, 2014) Let  $\mathcal{D}$  be a training set of size  $m$  drawn i.i.d. from distribution  $\sigma$ . Then with probability of at least  $1 - \delta$  over draw of  $\mathcal{D}$  from  $\sigma$ , for all  $h \in \mathcal{H}$ ,

$$L_\sigma(h) - L_{\mathcal{D}}(h) \leq 2\mathcal{R}_{\mathcal{D}}(l \circ \mathcal{H}) + 4\sqrt{\frac{2 \ln(4/\delta)}{m}}.$$

We provide the bound to measure the generalizability of the trained approximator in a training dataset with size  $m$ .

**Theorem D.14** (Generalization bound). Assume that the loss function  $l$  is  $T$ -Lipschitz continuous, then for hypothesis class  $\mathcal{H}$  of approximator and distribution  $\sigma$ , with probability at least  $1 - \delta$  over draw of the training set  $\mathcal{D}$  with size  $m$  from  $\sigma$ , for all  $h \in \mathcal{H}$ , we have

$$L_\sigma(h) - L_{\mathcal{D}}(h) \leq 2 \cdot \inf_{r>0} \left[ \frac{\sqrt{2 \log \mathcal{N}_{\infty,1}(\mathcal{H}, r)}}{m} + Tr \right] + 4\sqrt{\frac{2 \ln(4/\delta)}{m}}.$$

*Proof.* According to Theorem D.13, we have

$$L_\sigma(h) - L_{\mathcal{D}}(h) \leq 2\mathcal{R}_{\mathcal{D}}(l \circ \mathcal{H}) + 4\sqrt{\frac{2 \ln(4/\delta)}{m}}.$$

According to the assumption, the loss function  $l(x, y)$  is  $T$ -Lipschitz continuous under  $\ell_k$ -distance, i.e.,  $|l(x, y) - l(x', y)| \leq T\|x - x'\|_k$ , where  $\|\cdot\|_k$  is the  $k$ -norm. Let  $\mathcal{H}_r$  be the function class that *r-cover*  $\mathcal{H}$  for some  $r > 0$  and  $|\mathcal{H}_r| = \mathcal{N}_{\infty,1}(\mathcal{H}, r)$  be the *r-covering number* of  $\mathcal{H}_r$ . For all  $h \in \mathcal{H}$ ,  $h_r \in \mathcal{H}_r$  is denoted to be the function approximator that *r-covers*  $h$ . Based on above equation, we have

$$|l(h(x), y) - l(h_r(x), y)| \leq T\|h(x) - h_r(x)\|_k \leq Tr.$$

1404 Then we have

$$1405 \mathcal{R}_{\mathcal{D}}(l \circ \mathcal{H}) = \frac{1}{m} \mathbb{E}_{\mathbf{x} \sim \{+1, -1\}^m} \left[ \sup_{h \in \mathcal{H}} \sum_{i=1}^m x_i \cdot l(h(d_i), d_i) \right] \quad (1)$$

$$1408 = \frac{1}{m} \mathbb{E}_{\mathbf{x} \sim \{+1, -1\}^m} \left[ \sup_{h \in \mathcal{H}} \sum_{i=1}^m x_i \cdot (l(h_r(d_i), d_i) + l(h(d_i), d_i) - l(h_r(d_i), d_i)) \right] \quad (2)$$

$$1411 \leq \frac{1}{m} \mathbb{E}_{\mathbf{x} \sim \{+1, -1\}^m} \left[ \sup_{h_r \in \mathcal{H}_r} \sum_{i=1}^m x_i \cdot l(h_r(d_i), d_i) \right] + \frac{1}{m} \mathbb{E}_{\mathbf{x} \sim \{+1, -1\}^m} \left[ \sup_{h \in \mathcal{H}} \sum_{i=1}^m |x_i \cdot Tr| \right] \quad (3)$$

$$1415 \leq \sup_{h_r \in \mathcal{H}_r} \sqrt{\sum_{i=1}^m (\ell(h_r, d_i))^2} \cdot \frac{\sqrt{2 \log \mathcal{N}_{\infty, 1}(\mathcal{H}, r)}}{m} + \frac{Tr}{m} \mathbb{E}_{\mathbf{x}} \|\mathbf{x}\|_1 \quad (4)$$

$$1418 \leq \frac{\sqrt{2 \log \mathcal{N}_{\infty, 1}(\mathcal{H}, r)}}{m} + Tr \quad (5)$$

1421 The reduction from Eq. 3 to Eq. 4 is based on Massart’s lemma (Shalev-Shwartz & Ben-David, 2014). Finally,

$$1423 L_{\sigma}(h) - L_{\mathbf{D}}(h) \leq 2\mathcal{R}_{\mathcal{D}}(l \circ \mathcal{H}) + 4\sqrt{\frac{2 \ln(4/\delta)}{m}} \leq 2 \cdot \inf_{r>0} \left[ \frac{\sqrt{2 \log \mathcal{N}_{\infty, 1}(\mathcal{H}, r)}}{m} + Tr \right] + 4\sqrt{\frac{2 \ln(4/\delta)}{m}}$$

1426  $\square$

1427

1428 Therefore, given a training dataset with size  $m$ , we can have a generalization bound for the error depending on the characteristic of the loss function. In this paper, we follow the supervised learning framework to train the behavior cloning policy and environment model. Therefore, we can provide the following assumptions for the trained policy and environment models based on the above theorem.

1433 **Assumption D.15.** Suppose the error for training the behavior cloning policy is less than an extremely small  $\epsilon$  on the dataset with enough data (the size of data can be computed according to the above theorem). In that case, we consider that the trained behavior cloning policy is the same as the underlying behavior strategy of the dataset.

1437 **Assumption D.16.** Suppose the error for training the environment model is less than an extremely small  $\epsilon$  on the dataset with enough data. In that case, we consider that the trained environment model can provide the full information for the underlying game of the dataset.

1440

1441

1442

1443

1444

1445

1446

1447

1448

1449

1450

1451

1452

1453

1454

1455

1456

1457

## E IMPLEMENTATION DETAILS

Here, we provide the details for the model-based method by introducing our instantiate algorithms: MB-PSRO and MB-CFR, which are adaptations from two widely-used online equilibrium finding algorithms PSRO and Deep CFR.

### E.1 MB-PSRO

---

#### Algorithm 2 MB-PSRO

---

```

1: Input: Trained environment model  $E_{\theta_e}$ 
2: Initial policy sets  $\Pi$  for all players;
3: Compute expected rewards  $U^\Pi$  for each strategy  $\pi \in \Pi$  based on the environment model  $E_{\theta_e}$ ;
4: Initialize mate-strategies  $\sigma_i = \text{UNIFORM}(\Pi_i), \forall i$ ;
5: repeat
6:   for each player  $i \in [1, \dots, n]$  do
7:     for best response episodes  $t \in [1, \dots, T]$  do
8:       Sample  $\pi_{-i} \sim \sigma_{-i}$ ;
9:       Train best response policy  $\pi'_i$  over  $\rho \sim (\pi'_i, \pi_{-i})$ , which samples on the environment
       model  $E_{\theta_e}$ ;
10:    end for
11:    add the best response policy  $\pi'_i$  to policy set  $\Pi_i$ ;
12:  end for
13:  Compute missing entries in  $U^\Pi$  based on the environment model  $E_{\theta_e}$ ;
14:  Compute the meta-strategy  $\sigma$  using any meta-solver;
15: until Meet the convergence condition
16: Output: Policy set  $\Pi$  and meta-strategy  $\sigma$ 

```

---

We present the whole framework in Alg. 2. In the beginning, we need the well-trained environment model  $E_{\theta_e}$  as input to replace the function of the actual environment. Firstly, we initialize policy sets  $\Pi$  for all players using random strategies. Then, we estimate the expected utilities for each strategy profile based on the environment model  $E_{\theta_e}$  to form the meta-game matrix. In vanilla PSRO, this process needs to interact with the actual game environment. However, in the offline setting, the actual game environment is not available. Therefore, we use the well-trained environment model  $E_{\theta_e}$  to replace the actual game environment to provide the information needed in the algorithm. After building the meta-game matrix, the meta-strategy is initialized by a uniform strategy. Next, we compute the best response policy for every player and add these trained best response policies to their policy sets. When training the best response policy oracle using DQN or other RL algorithms, we sample the training data based on the environment model  $E_{\theta_e}$ . After adding these trained best response policies, we compute missing entries in the meta-game matrix still based on the trained environment model  $E_{\theta_e}$ . Then, the meta-strategy  $\sigma$  of the meta-game matrix can be computed using any meta-solver, such as Nash solver or  $\alpha$ -rank algorithm. For games with more than two players, the  $\alpha$ -rank algorithm is taken as the meta-solver. Finally, we repeat the above processes until meeting the convergence condition and output the policy set and meta-strategy as the approximate equilibrium strategy.

To compute the CCE strategy, we also instantiate one algorithm: MB-JPSRO, an adaptation from the JPSRO algorithm. The process of JPSRO is similar to PSRO except for the best response computation and meta solver. Therefore, MB-JPSRO is also similar to MB-PSRO. For this reason, we do not cover MB-JPSRO in detail here.

### E.2 MB-CFR

Alg. 3 shows the process of MB-CFR, which is adapted from the Deep CFR algorithm. It also needs the well-trained environment model  $E_{\theta_e}$  as input for the MB-CFR algorithm. We first initialize regret and strategy networks for each player and then initialize regret and strategy memories for each player (Lines 2-4). Then we need to update the regret network for every player. To do this, we perform a traverse function to collect corresponding training data. The traverse function can be any sampling-based CFR algorithm. Here, we use the external sampling algorithm as the traverse

**Algorithm 3** MB-CFR

---

```

1512 Input: Trained environment model  $E_{\theta_e}$ 
1513
1514 1: Initialize regret network  $R(I, a|\theta_{r,p})$  for all players;
1515 2: Initialize strategy network  $S(I|\theta_{\pi,p})$  for all players;
1516 3: Initialize regret memory  $M_{r,p}$  and strategy memory  $M_{\pi,p}$  for every player  $p$ ;
1517 4: for iteration  $t = 1$  to  $T$  do
1518 5:   for player  $p \in [1, \dots, n]$  do
1519 6:     for traverse episodes  $k \in [1, \dots, K]$  do
1520 7:       TRVERSE( $\phi, p, \theta_{r,p}, \theta_{\pi,-p}, M_{r,p}, M_{\pi,-p}, E_{\theta_e}$ );
1521 8:       # Use sample algorithm to traverse the game tree and record regret and strategy training
1522 9:       data
1523 10:     end for
1524 11:   Train  $\theta_{r,p}$  from scratch based on regret memory  $M_{r,p}$  for every player  $p$ ;
1525 12: end for
1526 13: Train  $\theta_{\pi,p}$  based on strategy memory  $M_{\pi,p}$  for every player  $p$ ;
1527 14: Output:  $\theta_{\pi,p}$  for every player  $p$ 

```

---

**Algorithm 4** TRVERSE( $s, p, \theta_{r,p}, \theta_{\pi,-p}, M_{r,p}, M_{\pi,-p}, E_{\theta_e}$ )-External Sampling Algorithm

---

```

1529 1: if  $s$  is terminal state then
1530 2:   Get the utility  $u_p(s)$  from the environment model  $E_{\theta_e}$ ;
1531 3:   Output:  $u_p(s)$ 
1532 4: else if  $s$  is a chance state then
1533 5:   Sample an action  $a$  from the available actions, which is obtained from model  $E_{\theta_e}$ ;
1534 6:    $s' = E_{\theta_e}(s, a)$ ;
1535 7:   Output: TRVERSE( $s', p, \theta_{r,p}, \theta_{\pi,-p}, M_{r,p}, M_{\pi,-p}, E_{\theta_e}$ )
1536 8: else if  $P(s) = p$  then
1537 9:    $I \leftarrow s[p]$ ; # Get the corresponding information set from the game state
1538 10:   $\sigma(I) \leftarrow$  strategy of  $I$  computed using regret values  $R(I, a|\theta_{r,p})$  based on regret matching;
1539 11:  for  $a \in A(s)$  do
1540 12:     $s' = E_{\theta_e}(s, a)$ ;
1541 13:     $u(a) \leftarrow$  TRVERSE( $s', p, \theta_{r,p}, \theta_{\pi,-p}, M_{r,p}, M_{\pi,-p}, E_{\theta_e}$ );
1542 14:  end for
1543 15:   $u_\sigma \leftarrow \sum_{a \in A(s)} \sigma(I, a)u(a)$ ;
1544 16:  for  $a \in A(s)$  do
1545 17:     $r(I, a) \leftarrow u(a) - u_\sigma$ ;
1546 18:  end for
1547 19:  Insert the infoset and its action regret values ( $I, r(I)$ ) into regret memory  $M_{r,p}$ ;
1548 20:  Output:  $u_\sigma$ 
1549 21: else
1550 22:   $I \leftarrow s[p]$ ;
1551 23:   $\sigma(s) \leftarrow$  strategy of  $I$  computed using regret value  $R(I, a|\theta_{r,-p})$  based on regret matching;
1552 24:  Insert the infoset and its strategy ( $I, \sigma(s)$ ) into strategy memory  $M_{\pi,-p}$ ;
1553 25:  Sample an action  $a$  from distribution  $\sigma(s)$ ;
1554 26:   $s' = E_{\theta_e}(s, a)$ ;
1555 27:  Output: TRVERSE( $s', p, \theta_{r,p}, \theta_{\pi,-p}, M_{r,p}, M_{\pi,-p}, E_{\theta_e}$ );
1556 28: end if

```

---

method to collect training data, and the process of external sampling is shown in Alg. 4. In this traverse function, we collect the regret training data of the traveler, and the strategy training data of other players are also gathered. After performing the traverse function several times, the regret network can be updated based on the regret memory. The above processes are repeated for  $T$  times. Then the average strategy network for every player is trained based on its corresponding strategy memory. Finally, the trained average strategy networks are output as the approximate equilibrium strategy.

## F ADDITIONAL EXPERIMENTAL RESULTS

In this section, we provide more experimental results and an ablation study. Finally, we provide the main parameters we used in our experiments.

### F.1 EXPERIMENTAL RESULTS

Here, we first verify that the performance of the model-based approach is independent of the algorithm used for computing equilibrium strategy. To this end, we perform both MB-CFR and MB-PSRO algorithms in the two-player Kuhn poker game under different sizes of offline datasets. Fig. 19 shows the results. We can find that under the same size of an offline dataset, MB-PSRO and MB-CFR achieve nearly identical results. When the size of the offline dataset increases, the performance of both algorithms becomes better. It may be caused by the environment model being well-trained with more data. These observations indicate that the performance of the model-based algorithm is independent of the algorithm used to compute the equilibrium strategy and mainly relies on the similarity between the trained environment model and the actual environment.

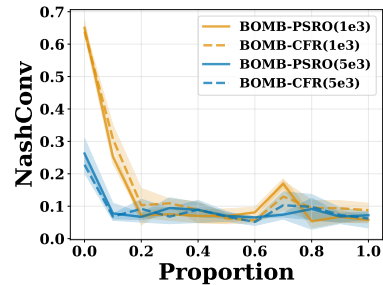
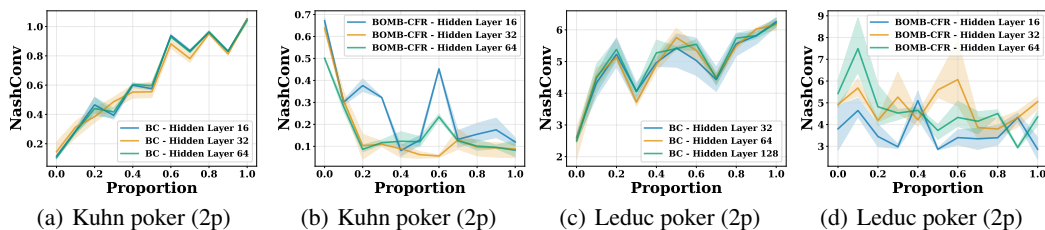


Figure 19: Results of different MB methods

### F.2 ABLATION STUDY



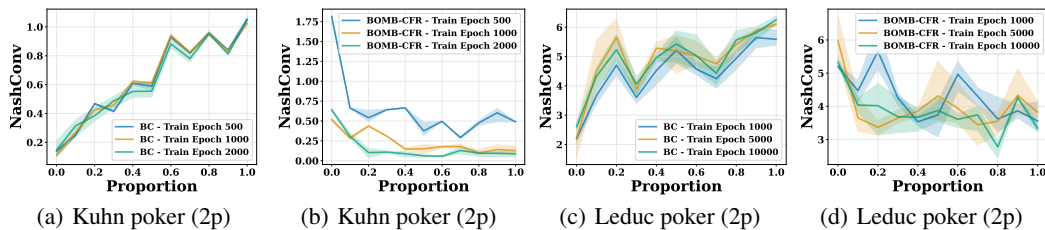
(a) Kuhn poker (2p)

(b) Kuhn poker (2p)

(c) Leduc poker (2p)

(d) Leduc poker (2p)

Figure 20: Ablation results for different hidden layer sizes.



(a) Kuhn poker (2p)

(b) Kuhn poker (2p)

(c) Leduc poker (2p)

(d) Leduc poker (2p)

Figure 21: Ablation results for different train epochs.

To investigate the influence of hyperparameters, we conduct several ablation experiments on two-player Kuhn poker and Leduc poker games. We consider different model structures with various numbers of hidden layers. Specifically, for the 2-Player Kuhn poker game, we use different environment models with 16, 32, and 64 hidden layers. For the 2-Player Leduc poker game, which is a more complicated game, the numbers of hidden layers for different models are 32, 64, and 128. In addition, we train the environment models for different epochs to evaluate the robustness of our approach. Figs. 20-21 show these ablation results. We find that the number of hidden layers and the number of training epochs have little effect on the performance of the BC algorithm. These results further verify that the performance of the BC algorithm primarily depends on the quality of



1620 the dataset. As we know, the performance of the model-based method mainly depends on the trained  
1621 environment model. Since the number of the hidden layer and the number of training epochs influ-  
1622 ence the training phase of the environment model, the number of the hidden layer and the number  
1623 of train epochs have a slight impact on the performance of the model-based method. As long as the  
1624 size of the hidden layer and the number of training epochs can guarantee that the environment model  
1625 is trained accurately, the performance of the model-based method will not be affected.

1626  
1627  
1628  
1629  
1630  
1631  
1632  
1633  
1634  
1635  
1636  
1637  
1638  
1639  
1640  
1641  
1642  
1643  
1644  
1645  
1646  
1647  
1648  
1649  
1650  
1651  
1652  
1653  
1654  
1655  
1656  
1657  
1658  
1659  
1660  
1661  
1662  
1663  
1664  
1665  
1666  
1667  
1668  
1669  
1670  
1671  
1672  
1673

## F.3 PARAMETER SETTING

We list the parameters used to train the behavior cloning policy and environment model for all games used in our experiments in Tab. 2.

Methods	Behavior Cloning Algorithm				Environment Model Training			
Games	Kuhn Poker (2p)		Kuhn Poker (3p)		Kuhn Poker (2p)		Kuhn Poker (3p)	
Data size	500	1000	1000	5000	500	1000	1000	5000
Hidden layer	32	32	32	32	32	32	32	32
Batch size	32	32	32	32	32	32	32	32
Train epoch	1000	2000	5000	5000	1000	2000	2000	5000
Games	Kuhn Poker (4p)		Kuhn Poker (5p)		Kuhn Poker (4p)		Kuhn Poker (5p)	
Data size	5000	20000	10000	20000	5000	20000	10000	20000
Hidden layer	64	64	64	64	64	64	64	64
Batch size	64	128	128	128	64	128	128	128
Train epoch	5000	5000	5000	5000	5000	5000	5000	5000
Games	Leduc Poker (2p)		Leduc Poker (3p)		Leduc Poker (2p)		Leduc Poker (3p)	
Data size	10000	20000	10000	20000	10000	20000	10000	20000
Hidden layer	128	128	128	128	64	128	128	128
Batch size	128	128	128	128	64	128	128	128
Train epoch	10000	10000	10000	5000	10000	10000	10000	10000
Games	Liar's Dice		Phantom TTT		Liar's Dice		Phantom TTT	
Data size	10000	20000	10000	20000	10000	20000	10000	20000
Hidden layer	64	64	128	128	64	64	128	128
Batch size	128	128	128	128	64	128	128	128
Train epoch	5000	5000	5000	5000	5000	5000	5000	5000

Table 2: Parameters for Behavior Cloning algorithm