IN-CONTEXT LEARNING FOR GAMES

Anonymous authors

Paper under double-blind review

ABSTRACT

Most literature in algorithmic game theory focuses on equilibrium finding, particularly Nash Equilibrium (NE). However, computing NE typically involves repeated computations of best responses (e.g., policy space response oracle (PSRO)), which can be computationally intensive. Moreover, NE strategies may not be ideal in games with more than two players or when facing irrational opponents. Consequently, NE strategies often require further adaptions to effectively address various types of opponents, impeding practical deployments. In contrast, In-Context Learning (ICL), i.e., learning from context examples, plays the core role in the generalizability of large language models (LLMs) to novel tasks without changing parameters. While ICL has been applied to decision-making tasks, e.g., algorithm distillation (AD), existing research primarily focuses on single-agent scenarios, and the ICL for games is largely unexplored. To facilitate the game solving and the practical deployment, the research question investigated in this work is: Can we leverage ICL to learn a model to i) play as **any player** of the game, ii) exploit **any** opponent to maximize the utility, and iii) be used to compute NE, without changing the parameters? In this work, we propose In-Context Exploiter (ICE) to address this question: i) ICE generates the diverse opponents with different capability levels for each player of the game to generate the training datasets, ii) ICE combines the curriculum learning and the ICL for single-agent scenarios (e.g., AD), to train the single model for all players of games, and iii) ICE leverages the pre-trained single model to play as each player of the game against different opponents and integrate with the equilibrium finding framework, e.g., PSRO, to compute NE. Extensive experiments on Kuhn poker, Leduc poker, and Goofspiel demonstrate that ICE can efficiently exploit different opponents as different players of the games and can be seamlessly integrated with PSRO to compute NE without changing the parameters.

034

000

001 002 003

004

005 006 007

008 009

010

011

012

013

014

015

016

017

018

019

021

023

024

025

026

027

028

029

031

032

035

036

1 INTRODUCTION

Multiplayer games provide ideal testbeds of Artificial Intelligence (AI) research (Silver et al., 2018; Brown & Sandholm, 2019), aptly referred as "Drosophila of AI" (Omidshafiei et al., 2020). Most literature in algorithmic game theory mainly focuses on equilibrium finding and Nash Equilibrium 040 (NE) (Nash, 1950) is the canonical solution concept for games, where no player can increase their 041 utility by unilaterally deviating. There are two main issues impeding the practical deployments 042 of NE into real-world scenarios. First, computing NE is typically computationally intensive, e.g., 043 policy space response oracle (PSRO) (Lanctot et al., 2017) requires repeatedly computing best-044 responses, which is resource-demanding even for small-scale games, e.g., Kuhn poker (Kuhn, 1950). Second, NE is not always an ideal solution, particularly in games with more than two players or when facing irrational opponents, therefore, further adaptions of these strategies are needed to handle 046 various opponents, which bring additional complexity and computational burdens to the deployment. 047

In-Context Learning (ICL), i.e., learning from context examples, is the core mechanism of the generalizability of large language models (LLMs) to novel tasks (Dong et al., 2022). Recent work (Laskin et al., 2022; Lee et al., 2024) extends ICL to decision making tasks, where a single policy can efficiently adapt to solve novel tasks from in-context examples, i.e., transitions. However, these works mainly focus on single-agent tasks and the ICL for games is largely unexplored. To facilitate the game solving and the practical deployment, in this work we provide a systematic investigation to ICL for games with the following core research question:

055 056

077

Can we leverage ICL to learn a model to i) play as **any player** of the game, ii) exploit **any opponent** to maximize the utility, and iii) be used to compute NE, without changing the parameters?

057 The three desiderata proposed in this paper aim to significantly improve the efficiency of equilibrium finding and enhance the practicability of the model for real-world deployment. However, there 058 are several challenges to achieving these desiderata. First, we need to sample diverse opponents with different levels of capabilities against different players to generate the training sets, which 060 is extremely important for the trained model to exploit the unknown opponents during inference. 061 Second, the training process needs to be carefully designed, as more capable opponents are more 062 difficult to exploit and the single model would suffer the catastrophic forgetting issues in exploiting 063 various opponents. Third, how to integrate the model trained with ICL into the current equilibrium 064 finding framework is still unclear, including both training and evaluation of the model. Therefore, 065 novel methods are required to tackle ICL for games. 066

To address these issues, we propose the In-Context Exploiter (ICE). The main contributions of 067 **ICE** include: i) **ICE** generates the diverse opponents for each player with the equilibrium finding 068 algorithms, e.g., CFR (Zinkevich et al., 2007), and collects the trajectories of players as the training 069 set, ii) ICE combines the curriculum learning and the ICL for single-agent scenarios, i.e., Algorithm Distillation (AD) (Laskin et al., 2022) and Decision-Pretrained Transformer (DPT) (Lee et al., 2024), 071 to train the model without catastrophic forgetting, and iii) ICE leverages the trained model to play as 072 each player of the game against any opponent and the trained model can also be used to compute NE 073 with equilibrium finding algorithms, e.g., PSRO (Lanctot et al., 2017) during inference. Extensive 074 experiments on Kuhn poker, Leduc poker, and Goofspiel demonstrate that ICE outperforms the NE strategies and RL methods, e.g., PPO, to efficiently exploit different opponents as different players 075 of the games and can be used to compute NE without changing the parameters. 076

2 RELATED WORK

079 In this section, we provide a concise review of the related work of ICL for games. The first line of related work is opponent modeling, which basically uses the prior knowledge and the observations 081 to infer the behaviors of an opponent (Nashed & Zilberstein, 2022). However, these methods usually 082 are based on the explicit model of the opponent and update the belief of the opponents during the 083 playing (Von Der Osten et al., 2017) or require further adaption, i.e., fine-tuning (Wu et al., 2022; 084 Foerster et al., 2017), which brings additional complexities of the methods. Conceptually, ICL can be viewed as an implicit opponent modeling and only adapt the behaviors of the model through 085 changing the in-contexts, which is much simpler than the current opponent modeling methods. The 086 second line of related work is the equilibrium finding, which lies in the core research of game the-087 ory. CFR (Zinkevich et al., 2007) is a no-regret method, which plays the core role in the success 088 on poker (Brown & Sandholm, 2019). Policy Space Response Oracle (PSRO) (Lanctot et al., 2017; 089 Muller et al., 2020) is another popular framework for equilibrium finding, which starts with a re-090 stricted games with limited number of policies for players and iteratively adding new best-responses 091 into consideration. Most equilibrium finding methods require the iterative computation of the best-092 or better-responses, which make these methods extremely computationally extensive. In this work, we intend to apply the ICL to facilitate the equilibrium finding. The third line of related work is in-094 context learning, which is the core mechanism of the remarkable generalizability of large language models (LLMs), e.g., GPT-4 (Achiam et al., 2023). By providing different in-context examples, 095 LLMs can quickly generalize to novel tasks. Recent work (Laskin et al., 2022; Lee et al., 2024) 096 successfully develop novel ICL methods to handle the decision making tasks with generalizability to novel tasks, which motivates us to investigate the ICL for games. 098

099 100

3 PRELIMINARIES

Imperfect-Information Extensive-Form Games. An imperfect-information extensive-form game (EFG) can be represented by a tuple $(N, H, A, P, \mathcal{I}, u)$ (Shoham & Leyton-Brown, 2008). N is the set of players, i.e., $N = \{1, ..., n\}$ and H is the set of histories which is the past action sequence. In particular, when the game starts, the history is an empty sequence \emptyset , representing the root node of the game tree. Additionally, every prefix of any sequence within H is also included in H. There is a set of special histories, called terminal histories, which are sequences that end in the leaf nodes of the game tree. Z is used to represent the set of terminal histories which is a subset of H, i.e., $Z \subset H$. $A(h) = \{a : (h, a) \in H\}$ is the set of available actions at any non-terminal history $h \in H \setminus Z$. P is 108 the player function that maps each non-terminal history to a player, i.e., $P(h) \mapsto N \cup \{c\}$ in which c 109 denotes chance player, representing these stochastic events beyond players' control. The information 110 set, represented by \mathcal{I}_i , forms a partition over the set of histories where i takes action, such that 111 player $i \in N$ cannot distinguish these histories within the same information set I_i . Therefore, each 112 information set $I_i \in \mathcal{I}_i$ corresponds to one decision-making point for 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 can employ $A(I_i)$ 113 and $P(I_i)$ to denote A(h) and P(h) for any history h within I_i . u_i represents the utility function of 114 player *i* that maps every terminal history to real numbers, i.e., $u_i : Z \mapsto \mathbb{R}$. 115

116 **Nash Equilibrium (NE).** The behavior strategy for player i, denoted by σ_i , is a function that maps 117 every information set I_i to a probability distribution over the available action $A(I_i)$. The set of strategies for player *i* is denoted by Σ_i , i.e., $\sigma_i \in \Sigma_i$. Given a strategy profile $\sigma = (\sigma_1, \sigma_2, ..., \sigma_n)$, 118 the expected value to player i is the sum of the expected payoff of these resulting terminal nodes, 119 i.e, $u_i(\sigma) = \sum_{z \in Z} \pi^{\sigma}(z) u_i(z)$. $\pi^{\sigma}(z) = \prod_{i \in N \cup \{c\}} \pi_i^{\sigma}(z)$ is the reaching probability of terminal 120 history z and $\pi_i^{\sigma}(z)$ is the contribution of player i to reach the terminal history z. The common so-121 lution concept for the imperfect-information extensive-form game is Nash equilibrium (NE) (Nash, 122 1950), defined as a strategy profile such that no player can increase their expected utility by unilat-123 erally switching to a different strategy. Formally, a strategy profile σ^* forms an NE if it satisfies 124 $u_i(\sigma^*) = \max_{\sigma'_i \in \Sigma_i} u_i(\sigma'_i, \sigma^*_{-i}), \forall i \in N$, where σ^*_{-i} refers to all the strategies in σ except for σ_i . 125

126 ICL for Decision Making. AD (Laskin et al., 2022) collects the dataset of learning histories gener-127 ated by a source RL algorithm and learns a causal transformer by autoregressively predicting actions given the cross-episode trajectory as context. The model trained by AD can be deployed to efficiently 128 complete novel tasks. DPT (Lee et al., 2024) is another ICL method for decision making, where the 129 context of DPT can be randomly sampled transitions and the prediction of the model is the optimal 130 action of the query state. Both methods demonstrate the impressive ability of ICL to efficiently 131 generalize to novel decision-making tasks without changing the parameters, which motivates us to 132 extend the ICL to games to address the issues illustrated in the motivating example. 133

Motivating Example. Given the rock-paper-scissors (RPS) game, whose payoff table is depicted in Table 1, the only NE is $(\frac{1}{3}, \frac{1}{3}, \frac{1}{3})$ for both players and the expected utility of each player is 0. However, suppose that an opponent always plays rock, the expected utility of playing NE strategy against him is still 0, while the expected utility of always playing paper is 1. Therefore, playing the NE strategy would be a safe option but is not optimal in even two-player zero-sum symmetric games.

Table 1. RUCK-Fapel-Scissors	Table	1:	Rock-Paper-Scissors
------------------------------	-------	----	----------------------------

	R	Р	S
R	(0, 0)	(-1, 1)	(1,-1)
Р	(1, -1)	(0, 0)	(-1, 1)
S	(-1, 1)	(1, -1)	(0, 0)

141 Simplex-NeuPL (Liu et al., 2022) also highlights a similar motivation for deviating from NE strate-142 gies in certain scenarios. Furthermore, computing NE usually requires the iterative computation 143 of best-responses where each iteration in PSRO (Lanctot et al., 2017) is using RL methods, e.g., PPO, to train the policy from scratch for each player, which is extremely time-consuming even for 144 small-scale games, e.g., Kuhn poker. The inefficiency of the equilibrium finding algorithms and the 145 impracticability of the NE strategies for deployment motivates us to investigate the ICL for games 146 to leverage one pre-trained model to i) play as any player, ii) exploit any opponent, and iii) be used 147 to compute to the NE without changing the parameters. 148

149 4 IN-CONTEXT EXPLOITER (ICE)

In this section, we introduce In-Context Exploiter (ICE), specifically crafted to train a model to exploit any unknown opponent and increase its utility through in-context learning. Fig. 1 provides an overview of ICE, which includes three stages: i) collecting interactive histories via any RL algorithm with diverse opponent strategies, ii) training a model within a curriculum learning framework, and iii) employing the trained model to play as any player in opponent exploitation and computing equilibrium in equilibrium finding algorithms. Each of these stages plays a critical role in ensuring the model's adaptability and generalizability.

157

158 4.1 STAGE I: DATA COLLECTION

159

Opponent Strategies Generation. Let \mathcal{D} denotes the set of opponent strategies. The key is that \mathcal{D} should consist of diverse and representative opponent strategies. To this end, we employ two methods: *random generation* and *learning-based generation*, and the corresponding sets are respectively

162 Collecting Training Inference **Opponent Exploiting** 163 ₩ÅŴ $(I_1, a_1, r_1, I'_1) \cdots (I_N, a_N, r_N, I'_N)$ Equilibrium Finding Divers In-Contexts (Meto М 166 Opponents Solve 0 0 Strategy Policy n Game In-Contex BR Policy ulat (Transformer) Dataset 168 Expand Policy (One for All Players) Oracle (Training with AD or DPT) 169 S 170

Figure 1: The three stages of ICE

denoted as \mathcal{D}_r and \mathcal{D}_l . Suppose $\mathcal{D}_r = \{\mathcal{D}_{r1}, \cdots, \mathcal{D}_{rm}\}$ which includes m random opponent strate-173 gies. Specifically, for each opponent's possible information set, we randomly generate a probability 174 distribution over the set of available actions. This randomness ensures that the generated opponent 175 strategies are highly diverse and include a variety of unpredictable opponent strategies, mimicking 176 scenarios where opponents may act irrationally. On the other hand, let $\mathcal{D}_l = \{\mathcal{D}_{l1}, \cdots, \mathcal{D}_{lb}\}$ which 177 consists of b opponent strategies generated via the learning-based method. We take CFR (Zinkevich 178 et al., 2007) as an example. Specifically, we run CFR for n iterations where at each iteration, the 179 opponent's average strategy among all previous strategies is added to the set \mathcal{D}_l . Intuitively, with the 180 progress of the CFR process, the generated average strategies will approach the NE strategy. From the definition of NE, \mathcal{D}_{l1} and \mathcal{D}_{lb} are respectively the easiest and hardest strategies to exploit, which 181 shows that \mathcal{D}_l includes diverse opponent strategies spanning a range of skill levels. 182

183 **Interactive History Collection.** With a slight abuse of notation, we use \mathcal{D}_{ri} , $1 \leq i \leq m$, (resp., 184 $\mathcal{D}_{li}, 1 \leq i \leq b$) to denote the in-context dataset corresponding to the opponent's randomly generated 185 strategy \mathcal{D}_{ri} (resp., learning-based generated strategy \mathcal{D}_{li}) as it is clear from the context. Following 186 AD (Laskin et al., 2022), we can utilize any RL algorithm to play against the opponent strategy \mathcal{D}_{ri} (resp., \mathcal{D}_{li}) and record the corresponding learning trajectories, each of which consists of the se-187 quence of information sets, actions taken, and the resulting utilities: $(I_0, a_0, r_0, \dots, I_T, a_T, r_T)$. In 188 our experiments, we employ proximal policy optimization (PPO) (Schulman et al., 2017) to collect 189 the trajectories. For DPT (Lee et al., 2024), the in-context training dataset consists of three compo-190 nents: in-context data (i.e., interactive history data), query states, and optimal actions. Unlike AD, 191 DPT requires an optimal strategy to give the optimal action for the corresponding query state. For 192 this purpose, we can leverage the strategy learned via PPO as mentioned previously as the optimal 193 strategy when collecting in-context datasets for DPT.

194 195 196

164

171

172

4.2 STAGE II: TRAINING

197 Recall that we aim to train a Transformer model such that it can play as each player of a game against 198 any opponent and can be used to compute NE with equilibrium finding algorithms during inference. 199 The most straightforward idea is to directly train the model on the collected datasets $\mathcal{D} = \mathcal{D}_r \cup \mathcal{D}_l$, 200 which could be unstable and inefficient since the model must adapt to a wide range of behaviors and tactics represented by these datasets. To overcome this issue, we devise a novel curriculum training 201 framework to stabilize the training process and improve the training efficiency. 202

203 Curriculum Generation. As mentioned in the previous section, for the learning-based generated 204 opponent strategies, \mathcal{D}_{l1} is the easiest task while \mathcal{D}_{lb} is the most difficult task. For the randomly gen-205 erated opponent strategies \mathcal{D}_r , while it is possible to sort them based on the gap between these strate-206 gies and the NE strategy, computing the gaps would be time-consuming due to the large number of 207 opponent strategies. Nevertheless, we observe that most of the randomly generated opponent strategies tend to be relatively simple to exploit, i.e., simple tasks. Therefore, we intersperse them into the 208 learning-based generated opponent strategies to generate an effective curriculum \mathcal{D}_{o} , which is shown 209 in Algorithm 1. Specifically, for every g learning-based generated opponent strategies, we add a ran-210 dom opponent strategy into the curriculum, and finally we have $\mathcal{D}_o = [\mathcal{D}_{l1}, \cdots, \mathcal{D}_{lq}, \mathcal{D}_{r1}, \cdots]$. For 211 notation convince, in the following, we use $\mathcal{D}_o = [\mathcal{D}_1, \cdots, \mathcal{D}_B]$ to denote the generated curriculum 212 where B = b + m denote the total number of datasets (opponent strategies). 213

Loss Functions. Before delving into the detailed training process, we first present the loss func-214 tions for training the Transformer model M_{θ} parameterized with θ . For the AD algorithm, for 215 each dataset $\mathcal{D}_i \in \mathcal{D}_o$, it includes K rounds and each round k has T_k steps, i.e., $(s_0^{(k)})$

226 227

229

231

233

236

237

239

256

257

 $(I_0^{(k)}, a_0^{(k)}, r_0^{(k)}), \dots, s_{T_k}^{(k)} = (I_{T_k}^{(k)}, a_{T_k}^{(k)}, r_{T_k}^{(k)})).$ These data are concatenated into one *long* history, i.e., $\mathcal{D}_i = (s_0^{(1)}, \dots, s_{T_1}^{(1)}, \dots, s_0^{(K)}, \dots, s_{T_k}^{(K)}).$ For convenience, we renumber this dataset, 216 217 218 $\mathcal{D}_i = (s_0^{(i)}, \cdots, s_T^{(i)})$, where $T = \sum_{k=0}^{K} T_k$. Then the loss function for training the model M_{θ} with context length L is given as 219 220

$$\mathcal{L}^{\rm AD}(\theta) = -\mathbb{E}_{h_{t,L}\sim\mathcal{D}_i} \sum_{l=0}^{L} \log M_{\theta}(a_{t+l}^{(j)}|I_{t+l}^{(j)}, h_{t,l}^{(j)}) \tag{1}$$

222 where $h_{t,l} = (s_t, ..., s_{t+l})$ is the context data. Similarly, in the DPT algorithm, for each dataset $\mathcal{D}_i \in \mathcal{D}_o$, we have $\mathcal{D}_i = \{D_1 = (s_0^{(i)}, \cdots, s_T^{(i)}), D_2 = \{(I_j^q, a_j^*)\}_{j=0}^K\}$, where D_1 is the same as 223 224 dataset with AD and (I_i^q, a_i^*) is the query information set and its corresponding optimal action. 225

$$\mathcal{L}^{\text{DPT}}(\theta) = -\mathbb{E}_{(h_{t,L}, I^{q}, a^{*}) \sim \mathcal{D}_{i}} \sum_{l=0}^{L} \log M_{\theta}(a^{*} | I^{q}, h_{t,l}^{(j)}).$$
(2)

Curriculum Training Process. Given the generated curriculum \mathcal{D}_{o} , we sequentially train the Trans-228 former M_{θ} on the datasets using the loss functions defined before, which is depicted in Algorithm 2. In this training process, a critical problem is catastrophic forgetting which is one of the common 230 issues when learning on multiple tasks represented by the different datasets. To mitigate this issue, we propose to periodically retrain the model on the previously visited datasets. Intuitively, this re-232 view mechanism could ensure that the model retains its proficiency in the earlier learned tasks while simultaneously acquiring new capabilities for the new tasks. Let \mathcal{D} denote the set of datasets that 234 have been visited. When $|\mathcal{D}| < N$, we apply the review mechanism to determine the dataset in 235 the current iteration t. With a slight abuse notation, we use \mathcal{D}_t to denote the first dataset that has not been visited in \mathcal{D}_{o} . Then, at the current iteration t, we randomly sample a dataset from \mathcal{D} with probability $1 - \sigma$ while training on the current dataset \mathcal{D}_t with probability σ . After all the datasets have been visited, i.e., $|\overline{\mathcal{D}}| = N$, we randomly sample a dataset from \mathcal{D}_{o} at each training iteration. 238

240	Algorithm 1 Curriculum Generation	Algorithm 2 Curriculum Training Framework
241	1: Input: Learning-based generation $\mathcal{D}_l =$	1: Input: Datasets $\mathcal{D}_{o} = [\mathcal{D}_{1},, \mathcal{D}_{N}]$
242	$[\mathcal{D}_{l1}, \cdots, \mathcal{D}_{ln}]$, random generation $\mathcal{D}_r =$	2: Initialize Transformer M_{θ} , σ and $\vec{\mathcal{D}} \leftarrow \emptyset$;
243	$[\mathcal{D}_{r1},\cdots,\mathcal{D}_{rm}], g, \mathcal{D}_o \leftarrow \emptyset$	3: for iteration $t = 1$ to T do
244	2: for $i = 1$ to <i>N</i> do	4: for train episode $p = 1$ to M do
245	3: if $i \mod g = 0$ then	5: if $ \overline{\mathcal{D}} < N$ then
246	4: $\mathcal{D}_o \leftarrow \mathcal{D}_o \cup \mathcal{D}_r[0], \mathcal{D}_r \leftarrow \mathcal{D}_r \setminus \mathcal{D}_r[0];$	6: Sample \mathcal{D}_t from $\overline{\mathcal{D}}$ with prob. $1 - \sigma$
247	5: else	or use \mathcal{D}_t in \mathcal{D}_o with prob. σ ;
248	6: $\mathcal{D}_o \leftarrow \mathcal{D}_o \cup \mathcal{D}_l[0], \mathcal{D}_l \leftarrow \mathcal{D}_l \setminus \mathcal{D}_l[0];$	7: else
249	7: end if	8: Sample a dataset \mathcal{D}_t from $\overline{\mathcal{D}}$;
250	8: if $ \mathcal{D}_r = 0$ or $ \mathcal{D}_l = 0$ then	9: end if
250	9: $\mathcal{D}_o \leftarrow \mathcal{D}_o \cup \mathcal{D}_l \text{ or } \mathcal{D}_o \leftarrow \mathcal{D}_o \cup \mathcal{D}_r;$	10: Train M_{θ} on \mathcal{D}_t using Eq. (1) or (2);
201	10: Early Break;	11: end_for
252	11: end if	12: if $ \mathcal{D} < N$ then $\mathcal{D} \leftarrow \mathcal{D} \cup \{\mathcal{D}_t\}$; end if;
253	12: end for	13: end for
254	13: Output: The curriculum \mathcal{D}_o	14: Output: The Transformer M_{θ}
255		

4.3 STAGE III: INFERENCE

258 After the training process, we freeze the Transformer model and then utilize it for exploiting oppo-259 nents as any player and computing the NE in equilibrium finding algorithms.

260 **Opponent Exploitation.** Given M_{θ} , a direct application is to use the model to play as any player 261 $p \in N$ and exploit any unknown opponent in an online manner over K episodes of interaction. Let 262 C_p denote the context of player p, initialized to the empty set $C_p \leftarrow \emptyset$. For simplicity, let E^j denote 263 the complete trajectory of j-th episode. Then, in episode k, at each time step t, the context consists 264 of the trajectories of the previous episodes and the interaction histories of the current episode $C_p =$ 265 $(E^1, \dots, E^{k-1}, I_0, a_0, r_0, \dots, r_{t-1})$. Next, an action is sampled according to $a_t \sim M_{\theta}(\cdot | I_t, C_p)$ 266 for the current information set I_t . Finally, the interaction tuple (I_t, a_t, r_t) is added to the context 267 C_p for predicting the next action in the next information set. After the K episodes, we can use the final context C_p to denote the "learned" strategy against the given opponent via in-context learning. 268 For brevity, we use $C_p \leftarrow OPPEXP(M_{\theta}, C_{-p}, K)$ to denote the process of opponent exploitation against the opponent with strategy C_{-p} and return C_p the player p's strategy. 269

270 Equilibrium Finding. We consider two common equilibrium-finding algorithms: self-play and 271 PSRO. In the self-play framework, the model M_{θ} is used to play as each player in the game. Self-272 play takes turns to update the strategy of each player $p \in N$ through the above opponent exploitation 273 process $C_p \leftarrow OPPEXP(M_{\theta}, C_{-p}, K)$ given the strategies of other players C_{-p} . However, the selfplay framework does not guarantee convergence. Instead, a more popular framework is PSRO. 274 We present the new PSRO framework in Algorithm 3, in which there are two primary differences 275 compared to the conventional PSRO: i) As the model M_{θ} is frozen, the policy of each player $p \in N$ 276 is represented by the context C_p , and thus, the (restricted) policy space of the player p is a set 277 of contexts Π_p ; ii) Learning the best response policy of each player at each PSRO iteration is an 278 opponent exploitation process given above. 279

280 281

282

283

284

285

286

287

288

289 290 291

292 293

295

296 297

298

Algorithm 3 PSRO with Opponent Exploitation

1: **Input:** Transformer model M_{θ}

Initialize players' policy spaces Π_p = {C_p ← Ø}, ∀p ∈ N;
 repeat
 Update the meta-game payoff matrix U^Π based on policy spaces Π via simulation;
 Compute players' meta-strategies δ_p for U^Π using any meta-solver, ∀p ∈ N;
 Expand policy space: Π_p ← Π_p ∪ {OPPEXP(M_θ, C_{-p}, K)} where C_{-p} ~ δ_{-p}, ∀p ∈ N;
 until convergence
 Output: The context sets and the meta-strategies of all players, Π_p and δ_p, ∀p ∈ N

5 EXPERIMENTS

To assess the effectiveness of **ICE** algorithm, we conduct comprehensive experiments on several popular extensive-form games. We start by outlining our experimental setting, followed by a detailed analysis of the results, structured around answering several key research questions.

5.1 EXPERIMENTAL SETTING

299 **Experimental Setup.** We selected a variety of poker games as test subjects, including both two-300 player and three-player versions of Kuhn Poker, Leduc Poker, and Goofspiel (with five cards). These 301 games serve as diverse platforms to evaluate our algorithm's performance. i) To rigorously evaluate 302 the performance of **ICE** algorithm in exploiting unknown opponents, we constructed three distinct 303 types of testbeds by randomly sampling different opponents: *in-distribution*, *out-of-distribution*, and 304 *NE opponent*. For the in-distribution testbed, we selected approximately 30 opponent tasks from the task dataset utilized during training. For the out-of-distribution testbed, we randomly sampled 305 20 opponent strategies to create a diverse set of test tasks. Finally, for the NE opponent testbed, 306 we specifically configured the opponent's strategy to align with the NE strategy, thereby forming 307 test tasks that directly reflect NE opponents. ii) To assess the performance of ICE algorithm in 308 computing the NE strategy, we employ NASHCONV to measure the gap between the learned strategy 309 and the NE strategy. Here, we focus specifically on two-player games as testbeds, as finding NEs in 310 multi-player games is notably challenging. 311

Baselines. We conduct online testing against different opponents to evaluate the performance of ICE 312 algorithm in exploiting unknown opponents. In this setting, we first select two widely used strategies 313 as baselines: the Best Response (BR) strategy and the Nash Equilibrium (NE) strategy. Notably, BR 314 is theoretically optimal in an online setting, as it is tailored to exploit a known opponent's strategy, 315 making it a theoretical upper-bound benchmark. The NE strategy is included as a baseline due to its 316 robustness against any opponent. Additionally, since this is an online setting, we select the proximal 317 policy optimization (PPO) algorithm (Schulman et al., 2017) as another baseline for comparison. 318 Furthermore, we include a multi-task pre-training with the fine-tuning framework as a baseline, as 319 exploiting different opponents can be framed as a multi-task learning problem. Since BR and NE 320 strategies are fixed once the opponent is given, we directly simulate these strategies against the 321 opponent's strategy to evaluate their performance. For the PPO algorithm, multi-task pre-training with fine-tuning, and ICE algorithm, we conduct evaluations under a limited number of online 322 interactions with the opponent. This limitation is deliberate, as our goal is to assess the capability of 323 these algorithms to quickly and effectively exploit different unknown opponents.

324 5.2 EXPERIMENTAL RESULTS

To better illustrate our findings, we show them by answering the following research questions (RQs). **RQ1:** *Can ICE act as any player in the game?*



Figure 2: In-distribution results when acting as any player

We claim that **ICE** algorithm is capable of training a model to perform as any player in the game. To substantiate this claim, we conduct evaluations by positioning the model trained by the **ICE** algorithm in various player roles within the game. Fig. 2 presents the results for both two-player and three-player Kuhn poker, assessed using the in-distribution testbed. Our experimental results reveal that **ICE** outperforms both the NE strategy and the PPO algorithm when acting as any player of the game, whether in two-player or three-player games. Notably, **ICE** exhibits the capacity to self-improve and closely approximate the BR strategy, leveraging its in-context learning ability. These results show the effectiveness of our method in adjusting to various strategic roles. Consequently, we only show the results from the perspective of one player, as a representation of our algorithm's ability to adapt to and perform in any given role.

RQ2: Can **ICE** adaptively exploit any opponent?



Figure 3: Results (left-Kuhn, middle-Leduc, right-Goofspiel) Figure

piel) Figure 4: NE opponent results

To answer this question, we perform experiments on the three distinct testbeds we previously intro-duced, which simulate different opponents including NE opponents. This diverse range of testing environments is crucial to comprehensively evaluate the adaptability and effectiveness of ICE al-gorithm in confronting any type of opponent. Fig. 3 and Fig. 4 display the results of playing three two-player games against different opponents. We have also carried out experiments for three-player games, with those results included in the Appendix due to page constraints. From Fig. 3, it is evident that ICE algorithm effectively demonstrates its in-context learning capability. Within a limited num-ber of interactions, **ICE** surpasses both the NE strategy and the PPO algorithm. In simpler cases, the PPO algorithm may reach performance levels similar to that of **ICE** algorithm. A key distinction, however, is that unlike PPO and other RL algorithms which require retraining from scratch for each new opponent, ICE algorithm achieves this without any parameter updates. It is worth noting that the out-of-distribution results are worse than the in-distribution results, highlighting a key area for further improvement in enhancing the model's generalization capabilities. In Fig. 4, we observe

that ICE algorithm is capable of achieving results comparable to those of the NE and BR strategies,
 which means that ICE can achieve high rewards against the NE opponent.



RQ3: *Can ICE* be used to compute NE without changing the parameter?

To answer this question, we conduct experiments on two-player Kuhn, Leduc, and Goofspiel poker games. We use the DPT algorithm to train the Transformer model, as its inference process does not impose sequence requirements on the context, making it convenient for evaluating the gap between the resulting strategy and the NE strategy, i.e., NASHCOV. Fig. 5 presents results using both the self-play and PSRO frameworks. From these results, we observe that, compared to the performance in Kuhn Poker and Leduc Poker, only the results in the Goofspiel game achieve low exploitability, regardless of whether the self-play or PSRO framework is used. This may be due to the existence of a pure strategy equilibrium in Goofspiel with 3 cards. Additionally, we observe that using the self-play framework to compute the NE strategy leads to unstable performance in Kuhn and Leduc Poker games. When using the PSRO framework, the NASHCONV decreases progressively in the Kuhn Poker game, indicating improved convergence, whereas it remains unstable in the Leduc Poker game. It may be the high dynamic in the Leduc Poker game. These findings represent a promising first step toward leveraging in-context learning for equilibrium computation without parameter updates. Nevertheless, the approach still struggles in games with complex dynamics, highlighting a clear direction for future work to improve robustness and adaptability in such scenarios.

RQ4: How does **ICE** perform compared with multi-task pre-training with fine-tuning framework?





Figure 7: Kuhn poker (2P) (play as P2)

Recent work has shown that multi-task pre-training with fine-tuning on new tasks performs equally or better than meta-learning pre-training with meta adaptation in RL tasks (Mandi et al., 2022). It indicates that pre-training with fine-tuning can quickly adapt to new tasks. In this paper, we compare this framework with ICE algorithm. Firstly, we pre-train a model using tasks generated from opponents' strategies, the same as those used in the ICE algorithm. Then, we evaluate its performance by fine-tuning based on the interactions with the opponent. The results for a two-player Leduc poker game are depicted in Fig. 6. Our findings reveal that ICE outperforms pre-training with fine-tuning approach in both in-distribution and out-of-distribution testbeds. Notably, pre-training with fine-tuning performs even underperforms compared to the PPO algorithm in the out-of-distribution testbed. It might be attributed to the extensive potential opponent strategies, where pre-training cannot encompass all opponent types, leading to slower adaptation to new tasks. Additionally, the conflict in training direction for different player roles in zero-sum games could further hinder the effectiveness of pre-training with fine-tuning.

RQ5: *Can our curriculum learning (CL) framework enhance the performance?*

ICE algorithm incorporates a curriculum learning (CL) framework for training the Transformer model. To explore the significance of CL, we conducted a comparative analysis by training the

432 Transformer model under different setups, including without the CL framework, which is trained 433 on a randomly ordered sequence of tasks, training first on random opponents, followed by learning-434 based opponents ("Random First"), and training first on learning-based opponents, followed by 435 random opponents ("Learned First"). The results are presented in Fig. 7. Interestingly, ICE, Ran-436 dom First, and Learned First achieve similar overall performance, which may be attributed to the inherent curriculum-like nature of learning-based opponents. However, we observe that Learned 437 First performs better in in-distribution cases, likely due to the stronger initial focus on structured 438 strategies and Random First performs better in out-of-distribution cases, benefiting from early ex-439 posure to diverse random strategies. Our ICE method demonstrates more stable performance across 440 both in-distribution and out-of-distribution settings. This differential in performance underscores 441 the significant contribution of the CL framework in boosting the effectiveness of ICE algorithm. 442

443 444

445

446 447

448

449

450

451 452

453

RQ6: Does the context length influence the performance of ICE?



Figure 8: Results of different context lengths

454 To investigate how the pre-defined context length affects performance, we conducted experiments 455 with various context lengths in our game scenarios. The results for two-player Kuhn poker and 456 three-player Goofspiel are shown in Fig. 8. For Kuhn poker, we observe that context length has 457 minimal impact on performance. This could be attributed to the simplicity of the game, where 458 even a short context is sufficient for effective in-context learning. Additionally, we note that, in the early stages, a larger context length may initially underperform compared to a shorter one, possibly 459 because more interactions are required to fully leverage the extended context. Conversely, for the 460 Goofspiel game, the results indicate that a larger context length improves performance. It implies 461 that in more complex games, a large context, which includes more interaction history information, 462 can significantly enhance the decision-making process. The extended context length provides a 463 broader historical perspective, which is particularly beneficial in complex strategic environments 464 where previous interactions greatly influence future decisions.

465 466 467

468

6 CONCLUSION

In this paper, we investigate an important research question: Can we leverage ICL to learn a model 469 to i) play as any player of the game, ii) exploit any opponent to maximize the utility, and iii) be used 470 to compute NE, without changing the parameters? To this end, we propose In-Context Exploiter 471 (ICE), which aims to train a single model that can satisfy all the desiderata presented in the previous 472 question. ICE delivers three main contributions: i) it generates diverse opponents via the equilib-473 rium finding algorithms, e.g., CFR and PSRO, and collects the trajectories of players using PPO as 474 the training dataset; ii) it combines the curriculum learning and ICL for single-agent scenarios (AD 475 and DPT) to train the model and employs a revisiting mechanism to preventing catastrophic forget-476 ting when training the model on multiple datasets; iii) it leverages the trained model to play as each 477 player in opponent exploitation and integrates the trained model into equilibrium finding algorithms, e.g., PSRO, to compute NE of the game. Extensive experimental results demonstrate that ICE can 478 efficiently exploit different opponents and can be used to compute NE without updating parameters. 479

Limitations and Future Work. There are several limitations to this work. First, we only focus
 on the games which are relatively small-scale. The ICL for large-scale games would require large
 models with longer in-context lengths and training time. Second, for the equilibrium finding with the
 trained model, we only focus on NE. Other solution concepts such as (coarse) correlated equilibrium
 (CCE) will be considered for multiplayer games. Third, the theoretical analysis of the ICL for games
 is largely unexplored, which will be investigated in future work. More detailed discussions on the
 limitations and future directions can be found in Appendix A.4.

486 REFERENCES

493

523

524

525 526

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. GPT-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- 491 Stefano V Albrecht and Peter Stone. Autonomous agents modelling other agents: A comprehensive
 492 survey and open problems. *Artificial Intelligence*, 258:66–95, 2018.
- 494 Noam Brown and Tuomas Sandholm. Superhuman AI for multiplayer poker. *Science*, 365(6456):
 495 885–890, 2019.
- 496
 497
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
 498
- Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation
 of deep networks. In *International conference on machine learning*, pp. 1126–1135, 2017.
- 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.
- He He, Jordan Boyd-Graber, Kevin Kwok, and Hal Daumé III. Opponent modeling in deep rein forcement learning. In *International Conference on Machine Learning*, pp. 1804–1813, 2016.
- Harold W Kuhn. A simplified two-person poker. *Contributions to the Theory of Games*, 1(97-103): 2, 1950.
- Marc Lanctot, Vinicius Zambaldi, Audrūnas Gruslys, Angeliki Lazaridou, Karl Tuyls, Julien
 Pérolat, David Silver, and Thore Graepel. A unified game-theoretic approach to multiagent rein forcement learning. In *Proceedings of the 31st International Conference on Neural Information Processing Systems*, pp. 4193–4206, 2017.
- 513 Michael Laskin, Luyu Wang, Junhyuk Oh, Emilio Parisotto, Stephen Spencer, Richie Steigerwald,
 514 DJ Strouse, Steven Hansen, Angelos Filos, Ethan Brooks, et al. In-context reinforcement learning
 515 with algorithm distillation. *arXiv preprint arXiv:2210.14215*, 2022.
- Jonathan Lee, Annie Xie, Aldo Pacchiano, Yash Chandak, Chelsea Finn, Ofir Nachum, and Emma Brunskill. Supervised pretraining can learn in-context reinforcement learning. *Advances in Neural Information Processing Systems*, 36, 2024.
- Siqi Liu, Marc Lanctot, Luke Marris, and Nicolas Heess. Simplex neural population learning: Any mixture bayes-optimality in symmetric zero-sum games. In *International Conference on Machine Learning*, pp. 13793–13806. PMLR, 2022.
 - Zhao Mandi, Pieter Abbeel, and Stephen James. On the effectiveness of fine-tuning versus metareinforcement learning. *arXiv preprint arXiv:2206.03271*, 2022.
 - Richard D McKelvey and Thomas R Palfrey. Quantal response equilibria for normal form games. *Games and economic behavior*, 10(1):6–38, 1995.
- Paul Muller, Shayegan Omidshafiei, Mark Rowland, Karl Tuyls, Julien Perolat, Siqi Liu, Daniel Hennes, Luke Marris, Marc Lanctot, Edward Hughes, Zhe Wang, Guy Lever, Nicolas Heess, Thore Graepel, and Remi Munos. A generalized training approach for multiagent learning. In *ICLR*, 2020.
- JF Nash. Equilibrium points in n-person games. *Proceedings of the National Academy of Sciences* of the United States of America, 36(1):48–49, 1950.
- Samer Nashed and Shlomo Zilberstein. A survey of opponent modeling in adversarial domains.
 Journal of Artificial Intelligence Research, 73:277–327, 2022.
- Shayegan Omidshafiei, Karl Tuyls, Wojciech M Czarnecki, Francisco C Santos, Mark Rowland,
 Jerome Connor, Daniel Hennes, Paul Muller, Julien Pérolat, Bart De Vylder, et al. Navigating the
 landscape of multiplayer games. *Nature Communications*, 11(1):5603, 2020.

John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. arXiv preprint arXiv:1707.06347, 2017. Shai Shalev-Shwartz et al. Online learning and online convex optimization. Foundations and *Trends*® *in Machine Learning*, 4(2):107–194, 2012. Yoav Shoham and Kevin Leyton-Brown. Multiagent Systems: Algorithmic, Game-Theoretic, and Logical Foundations. Cambridge University Press, 2008. David Silver, Thomas Hubert, Julian Schrittwieser, Ioannis Antonoglou, Matthew Lai, Arthur Guez, Marc Lanctot, Laurent Sifre, Dharshan Kumaran, Thore Graepel, et al. A general reinforcement learning algorithm that masters chess, shogi, and Go through self-play. Science, 362(6419):1140-1144, 2018. Friedrich Burkhard Von Der Osten, Michael Kirley, and Tim Miller. The minds of many: Opponent modeling in a stochastic game. In IJCAI, pp. 3845-3851, 2017. Zhe Wang, Petar Veličković, Daniel Hennes, Nenad Tomašev, Laurel Prince, Michael Kaisers, Yoram Bachrach, Romuald Elie, Li Kevin Wenliang, Federico Piccinini, et al. TacticAI: an AI assistant for football tactics. Nature communications, 15(1):1906, 2024. Zhe Wu, Kai Li, Hang Xu, Yifan Zang, Bo An, and Junliang Xing. L2e: Learning to exploit your opponent. In 2022 International Joint Conference on Neural Networks (IJCNN), pp. 1-8, 2022. Martin Zinkevich, Michael Johanson, Michael Bowling, and Carmelo Piccione. Regret minimization in games with incomplete information. In Proceedings of the 20th International Conference on *Neural Information Processing Systems*, pp. 1729–1736, 2007.

594 A DISCUSSION

596

597

611 612

613 614

615 616

621

622

623 624

625 626

A.1 ICE vs. Opponent Modeling

598 ICE can be viewed as a method with implicit opponent modeling (He et al., 2016; Albrecht & Stone, 2018), where the modeling of the opponent is implicitly encoded into the parameters of the model. There are several advantages of ICE over opponent modeling: i) ICE does not need an 600 explicit model for the opponents, where the explicit model in the opponent modeling may restrict 601 the generalizability of the methods, ii) ICE can exploit different opponents without changing the 602 parameters, where the opponent modeling may need to fit the parameters of the opponent model 603 during game play and then make the decision in response to the opponent. To summarize, ICE is 604 simpler, more efficient, and more generalizable. ICE also has the disadvantage, i.e., the ability to 605 model the opponents is largely determined by the length of the in-context. With longer in-context, 606 ICE can model more opponents, while the model will also be larger and the training cost will be 607 increased. We will discuss the methods to reduce the length of the in-context in the next section. On 608 the other hand, we can also introduce an explicit model for opponents into ICE, where the parameters 609 of the opponent model can be fitted through in-context learning. The explicit opponent model can help us to understand the internal mechanism of ICE. 610

A.2 ICE vs. Equilibrium Finding





The ICE approach and traditional equilibrium finding share the goal of developing strategies for 627 extensive-form games but diverge significantly in their methodology and capabilities. Equilibrium 628 finding methods focus on computing an NE strategy profile π that performs robustly against any 629 potential opponent. This is illustrated in Fig. 9 (left), where the NE strategy remains stable over time, 630 providing consistent performance across a wide range of opponents. In contrast, the ICE framework 631 is designed to exploit any opponent effectively by adapting its strategy through in-context learning. 632 As shown in Fig. 9 (right), ICE leverages the trained model M to exploit different opponents only 633 through in-context learning, leading to performance that improves dynamically over time as more 634 interactions are observed. This adaptability allows ICE to achieve better performance against other 635 irrational opponents compared to the fixed NE strategy.

636 637

638

A.3 ICE vs. Online Learning, Multitask Learning, and Meta Learning

ICE, as well as other in-context learning methods (Laskin et al., 2022; Lee et al., 2024), is similar 639 to online learning methods, e.g., no-regret learning (Shalev-Shwartz et al., 2012). However, ICE 640 does not change the parameters of the model during the game play with the opponents, which differs 641 from online learning. We believe that online learning, especially no-regret learning, can be used 642 to analyze the behaviors of ICE and in-context learning methods, which will be explored in future 643 works. We also consider online learning methods as our baselines. We note that PPO is an online 644 learning and on-policy method and PPO is scalable and widely used. Therefore, we include PPO as 645 the baseline in our experiments. 646

647 The training of ICE is also similar to multitask learning (Mandi et al., 2022) and meta-learning (Finn et al., 2017), where multi-task learning learns a policy for different tasks, and meta-learning enables

the fast adaption of the learned policy on specific tasks. ICE also learns a policy for different tasks,
 where the model parameters are not changed, but the behaviors are changed during game play. In the
 experiment section, we choose the PPO method initialized with a pre-trained policy to benchmark
 multi-task and meta-learning methods, as shown in Figure 6.

653 A.4 LIMITATIONS AND FUTURE WORKS

In this section, we provide a mode detailed discussion about the limitations of this work.

Scalability of ICE. We only focus on relatively small-scale games in this work, including Kuhn
poker and Leduc poker. However, the large-scale games will have more relevance for real-world
deployment, such as Texas Hold'em poker (Brown & Sandholm, 2019) and sport games, e.g., football game (Wang et al., 2024). We will consider to scale up the current ICE methods to tackle
large-scale games, which may require larger models with longer in-context lengths and longer training time. We would also consider to train the foundation model of ICL for games, i.e., one model
trained to generalize to different games to further improve the generalizability of ICL.

Other Solution Concepts. We only focus on exploiting the opponents and computing NE in this
 work. However, there are many other solution concepts such as (coarse) correlated equilibrium
 (CCE) and quantal response equilibrium (QRE) (McKelvey & Palfrey, 1995), which are also important concepts in game theory. ICE has the potential to compute other solution concepts with one
 single pre-trained model, which will be tacked in future work.

Theoretical Analysis. The ICL for games is a new research area and the theoretical analysis is required for further investigations, including the optimality of the converged solutions and the convergence of the equilibrium finding with the ICL models. In the future, we will conduct a systematic theoretical analysis of the ICL for games, similar to the analysis for DPT (Lee et al., 2024).

702 B IMPLEMENTATION DETAILS

704

705

In this section, we provide the experimental details of ICE algorithm from its three main stages.

Opponent Generation. In this paper, we employ two methods, as introduced in the main paper, to 706 generate a diverse range of opponent strategies. To implement the random generation method, we 707 traverse through all the information sets of an opponent and assign a randomly generated strategy 708 to each information set. This approach allows us to generate various opponents exhibiting random 709 behaviors. To implement the learning-based generation method, we utilize a well-known algorithm, 710 Counterfactual Regret Minimization (CFR) (Zinkevich et al., 2007), as the equilibrium-finding al-711 gorithm. By applying CFR to solve the game, we record the average strategy for each player at 712 each iteration. This process generates a series of opponent strategies that evolve from random to 713 increasingly robust over time. These two methods collectively ensure that our dataset includes a 714 wide spectrum of opponent strategies, ranging from entirely unpredictable to highly strategic. Such a comprehensive dataset is instrumental in training our model to adapt and respond effectively to 715 various levels of opponent sophistication and strategy. 716

717 Interactive History Collection. It's important to recognize that when an opponent's strategy is
718 known, the task of exploiting that opponent to maximize utility effectively becomes a reinforcement
19 learning (RL) problem. Consequently, each distinct opponent strategy corresponds to a unique RL
720 task. For the AD algorithm, to collect interactive history data from our diverse opponent strategies
721 for training purposes, we adopt the Proximal Policy Optimization (PPO) algorithm (Schulman et al.,
722 2017) to address each of these RL tasks. During this process, we systematically record the learning
reaction of the PPO algorithm, specifically capturing the contents of the reply buffer used by PPO.

For the DPT algorithm, the training dataset consists of three components: in-context data (i.e., interactive history data), query states, and optimal actions. To generate the in-context data, we use a random strategy to play against various opponents and record the resulting interactive history. The optimal strategies used against different opponents are derived from the PPO policies that were trained earlier. Subsequently, we randomly sample query states and utilize the optimal strategy to obtain the corresponding optimal actions.

730 Curriclum Learning. The curriculum learning framework is crucial for effectively training the 731 Transformer model, with the curriculum's design being its core component. While the main paper 732 provides a comprehensive explanation of the curriculum generation process and the overall learning framework structure, this section will not revisit those specifics. However, it is important to high-733 light that the thoughtful design of the curriculum is key to the success of the model's training. By 734 gradually increasing task complexity and progressing through structured stages, the model is able to 735 incrementally build its understanding and capabilities without being overwhelmed. This approach 736 aligns well with the principles of in-context learning, allowing the Transformer to adapt and respond 737 efficiently to a broad spectrum of strategic scenarios. 738

Inference. In the main paper, we introduce how our trained model is used to exploit opponents and compute equilibrium strategy. Here, we provide a detailed description of the inference process. We first focus on how to use our model to exploit opponents, as outlined in Algorithm 4. For a given opponent, the model M_{θ} is inferred based on the current context to make decisions, while simultaneously updating the context with the previous interactive history data.

744 Next, we introduce two frameworks for computing equilibrium strategies, detailed. The first frame-745 work is self-play, where the model essentially plays against itself by taking on the roles of all players in the game. As described in Algorithm 5, the model plays each player in turn to exploit the oppos-746 ing strategies. It is important to note that the context C is sufficient for recording strategies, as the 747 strategy for any state s can be queried through $M_{\theta}(\cdot|s, C)$. Therefore, only the most recent context 748 needs to be stored. The second framework is PSRO, a widely-used algorithm for solving imperfect-749 information extensive-form games. In this framework, we substitute the best response oracle with 750 the opponent-exploiting inference process. The context is recorded to represent the learned best 751 response strategy effectively. 752

Parameter Setting. Here, we list the parameters used in the ICE algorithm for all games in Tab.2. In this table, the previous rate σ is used to control the blend of new and prior tasks to prevent catastrophic forgetting and the number of trains per task refers to the number of training for each selected task (i.e., *M* in Algorithm 2).

	gorithm 4 Opponent Exploiting	(OPPEX	$\operatorname{KP}(M_{\theta}, C)$	(-p, K))			
1:	Input: Transformer model M_{θ}	, player	id p				
2:	Initialize the context $C = \emptyset$;						
3:	for $t = 1$ to K do						
4:	s is initialized as game's init	tial state	and s_{pre} ,	$a_{pre} = None$, None;		
5:	while TRUE do						
6: 7.	If s is the P_i 's turn then if s is not None then	•					
/: g.	ADD(e a r(e))	d(e)	to C .				
0. Q.	end if	(s), (a(s))	ιο Ο,				
10:	$a = M_{\theta}(\cdot C, s)$:						
11:	$s_{pre} = s, a_{pre} = a;$						
12:	else						
13:	# get the opponent action	on					
14:	$a = M_{\theta}(\cdot C_{-p}, s);$						
15:	end if			6			
16:	# get the next game state $T(z)$	based on	transitio	n function			
17:	s' = T(s, a);						
18:	If s' is the end state then $ADD(a = a = r(a))$	d(a) to	C.				
19. 20·	ADD $(S_{pre}, a_{pre}, r(S))$, Break:	u(s)) to	С,				
20.21	end if						
22:	end while						
23:	end for						
24:	Output: Context C						
Alg 1: 2: 3: 4: 5: 6:	gorithm 5 Self-Play Framework Input: Transformer model M_{θ} for $i = 1$ to T do for p in $\{1,, n\}$ do $C_p = OPPEXP(M_{\theta}, C_{-p},$ end for end for	, Iteratio K);	on numbe	r K			
Alg 1: 2: 3: 4: 5: 6:	gorithm 5 Self-Play Framework Input: Transformer model M_{θ} for $i = 1$ to T do for p in $\{1,, n\}$ do $C_p = OPPEXP(M_{\theta}, C_{-p},$ end for end for T	(k); able 2: F	on numbe	for ICE(AD)			
Alg 1: 2: 3: 4: 5: 6:	gorithm 5 Self-Play Framework Input: Transformer model M_{θ} for $i = 1$ to T do for p in $\{1,, n\}$ do $C_p = OPPEXP(M_{\theta}, C_{-p},$ end for end for T Games	, Iteratio K); able 2: F	on numbe	r K for ICE(AD)	Kuhn	Leduc	Goofsp
Alg 1: 2: 3: 4: 5: 6:	gorithm 5 Self-Play Framework Input: Transformer model M_{θ} for $i = 1$ to T do for p in $\{1,, n\}$ do $C_p = OPPEXP(M_{\theta}, C_{-p},$ end for end for T Games Number of Player	able 2: F	on numbe	r <i>K</i> for ICE(AD) Goofspiel 2	Kuhn 3	Leduc 3	Goofspi 3
Alg 1: 2: 3: 4: 5: 6:	gorithm 5 Self-Play Framework Input: Transformer model M_{θ} for $i = 1$ to T do for p in $\{1,, n\}$ do $C_p = OPPEXP(M_{\theta}, C_{-p},$ end for end for $C_p = OPPEXP(M_{\theta}, C_{-p}, C_$	able 2: F K); K); Able 2: F Kuhn 2 0.1	Parameter Leduc 2 0.3	r K for ICE(AD) Goofspiel 2 0.3	Kuhn 3 0.1	Leduc 3 0.3	Goofspi 3 0.3
Alg 1: 2: 3: 4: 5: 6:	gorithm 5 Self-Play Framework Input: Transformer model M_6 for $i = 1$ to T do for p in $\{1,, n\}$ do $C_p = OPPEXP(M_{\theta}, C_{-p},$ end for end for T Games Number of Player Previous Rate σ Number of Train per Task M	able 2: F K); K); Able 2: F Kuhn 2 0.1 10	Parameter Parameter Leduc 2 0.3 10	for ICE(AD) Goofspiel 2 0.3 10	Kuhn 3 0.1 10	Leduc 3 0.3 30	Goofspi 3 0.3 30

C ADDITIONAL EXPERIMENTAL RESULTS.

In this section, we present further experimental results to substantiate the effectiveness of **ICE** algorithm in opponent exploiting. While the main paper provided the performance of **ICE** in three two-player game scenarios evaluated across three distinct testbeds, here we extend our analysis to include results from experiments conducted on three different three-player games.



Figure 10: In-distribution results of three-player games

Firstly, we present the results from the in-distribution testbed in Fig.10. In these three-player games, it is evident that the model trained using the ICE algorithm successfully functions as any player in the game, demonstrating in-context learning ability with increasing iterations. The Best Response (BR) strategy, while theoretically the optimal approach since it is tailored against a known opponent's strategy, isn't practical in real-world scenarios where an opponent's strategy isn't known in advance. In our results, the BR strategy's performance is included merely as a theoretical bench-mark. Notably, while the ICE-trained model doesn't achieve the theoretical optimal values of the BR strategy, it consistently surpasses both the NE strategy and the PPO algorithm. This observation is significant as it indicates that the ICE-trained model can exploit the opponents more effectively than the NE strategy, which is generally considered the most conservative approach. The ability of ICE to outperform in these multi-player game scenarios demonstrates its potential as a powerful tool for strategic decision-making in complex, real-world situations.

Next, Fig.11 shows the results from the out-of-distribution testbed, where we observe trends similar to those in the in-distribution testbed. The key distinction here is that the opponents in the out-of-distribution testbed are randomly generated, which often results in simpler strategic scenarios. In contrast, the in-distribution testbed encompasses a mix of randomly generated and learning-



Lastly, we discuss the results against NE opponents, as shown in Fig.12. Our findings reveal that 906 the ICE algorithm achieves better or comparable performance to the NE strategy only in the three-907 player Kuhn poker game. However, in other game scenarios, while ICE does not outperform the NE 908 strategy, it still maintains a higher level of performance than the PPO algorithm. The less optimal 909 performance of **ICE** in these cases can be attributed to the highly dynamic game environment and 910 stability of the NE opponents. In three-player games, the player faces two opponents simultaneously, 911 and if both adopt the conservative NE strategy, exploiting them concurrently becomes significantly 912 challenging. This observation highlights an area for future development. Improving the ICE algo-913 rithm to more effectively handle situations where multiple opponents employ highly conservative 914 strategies, such as the NE, will be a focus of our future research.

- 915
- 916



Figure 12: Results of three-player games against NE opponent