
HIERARCHICAL DEEP COUNTERFACTUAL REGRET MINIMIZATION

Anonymous authors

Paper under double-blind review

ABSTRACT

Imperfect Information Games (IIGs) are used to model games under uncertainty or lack complete information. Counterfactual Regret Minimization (CFR) is one of the most successful families of algorithms for IIGs. The integration of skill-based strategy learning with CFR could potentially mirror more human-like decision-making and improve learning on complex IIGs. It enables the learning of a hierarchical strategy, wherein low-level components represent skills for solving subgames and the high-level component manages the transition between skills. In this paper, we introduce the first hierarchical version of Deep CFR (HDCFR), an innovative method that boosts learning efficiency in tasks involving extensively large state spaces and deep game trees. Notably, HDCFR enables learning with predefined (human) expertise and extracting skills transferable to similar tasks. We first present the algorithm and establish its theory in a tabular setting, including hierarchical CFR update rules and a variance-reduced Monte Carlo sampling extension for the model-free setting, where backtracking is infeasible. We then extend HDCFR to large-scale tasks via deep learning objectives that match the tabular targets under exact function fitting. Code: <https://anonymous.4open.science/r/HDCFR>.

1 INTRODUCTION

Imperfect Information Games (IIGs) can model various application domains where decision-makers have incomplete information about the state of the environment, such as auctions Noe et al. (2012), diplomacy Bakhtin et al. (2022), cybersecurity Kakkad et al. (2019). Tabular Counterfactual Regret Minimization (CFR) Zinkevich et al. (2007) has been employed in all recent milestones of Poker AI Bowling et al. (2015); Moravčík et al. (2017); Brown & Sandholm (2018), a quintessential benchmark for IIGs, highlighting the robustness of CFR. To improve scalability, researchers have proposed deep learning extensions of CFR Brown et al. (2019); Li et al. (2020); Steinberger et al. (2020), leveraging neural networks as function approximators to operate in extensive state spaces.

Professionals in a field typically possess domain-specific skills, which they compose into strategies for diverse tasks. Integrating skill-based strategy learning with CFR can enable human-like decision-making for IIGs and improve learning on complex tasks with extended decision horizons. To accomplish this, the agent learns a hierarchical strategy: low-level components represent specific skills, and the high-level component coordinates switching among skills. Notably, this is akin to the option framework Sutton et al. (1999) proposed in reinforcement learning (RL), which enables learning or planning at multiple levels of temporal abstraction. Hierarchical strategies are more interpretable, allowing humans to identify specific subgames where AI agents struggle and inject critical skills as solutions. Also, skills acquired in one task can be transferred to similar tasks to facilitate learning in new IIGs.

We propose the first hierarchical extension of Deep CFR (HDCFR), a novel approach that improves learning efficiency in tasks with deep game trees and facilitates knowledge transfer from similar games. We establish the theoretical foundations of our algorithm in the tabular setting and then introduce deep learning extensions for practical applications in large-scale games. Our contributions are as follows: (1) We extend the standard definition of IIGs by incorporating a hierarchical strategy and provide CFR updating rules (i.e., HCFR) for this strategy, along with convergence guarantees. (2) Vanilla CFR relies on a perfect game tree model and requires a complete traversal of the game tree in each training iteration, which restricts its use. Thus, we propose a sample-based model-free

054 extension of HCFR, including unbiased Monte Carlo estimators of counterfactual regrets and a
 055 hierarchical baseline function for variance reduction. Controlling sample variance is vital for tasks
 056 with extended decision horizons, which our algorithm targets. (3) We present HDCFR, where the
 057 hierarchical strategy, regret, and baseline are approximated with neural networks. The training
 058 objectives are demonstrated to be equivalent to those proposed in the tabular setting, i.e., (1) and (2),
 059 when optimality is achieved, thereby preserving the theoretical results while enjoying scalability.

060 We focus on the model-free outcome sampling (OS) setting, where backtracking over the full game
 061 tree is infeasible and the game model is not explicit or traversable. Accordingly, we benchmark against
 062 strong model-free OS baselines in the same setting, including DREAM, OS-MCCFR/OSSDCFR and
 063 NFSP. Our goal is to stabilize deep-horizon learning via temporal abstraction and variance reduction,
 064 without restricting the original action space.

066 2 BACKGROUND

068 **A detailed discussion of related works is provided in Appendix A.**

070 2.1 COUNTERFACTUAL REGRET MINIMIZATION

072 In an IIG Osborne & Rubinstein (1994), players make sequential moves represented by a game
 073 tree. At each non-terminal state, the player in control chooses from a set of available actions; at
 074 each terminal state, each player receives a payoff. With imperfect information, a player may not
 075 know which state they are in (e.g., in poker, a player observes private cards and public boards
 076 but not opponents' hands). Thus, each player acts based on an *information set*—a collection of
 077 indistinguishable states. Formally, the game is $\langle N, H, A, P, \sigma_c, u, \mathcal{I} \rangle$. N is a finite set of players. H
 078 is the set of histories, where each history is a sequence of actions from the start of the game and
 079 corresponds to a state. For $h, h' \in H$, write $h \sqsubseteq h'$ if h is a prefix of h' . The set of actions at $h \in H$
 080 is $A(h)$. If $a \in A(h)$, then $(ha) \in H$ is a successor history. Histories with no successors are terminal,
 081 $H_{TS} \subseteq H$. $P : H \setminus H_{TS} \rightarrow N \cup \{c\}$ maps each non-terminal history to the player in control, where
 082 c is the chance player acting according to $\sigma_c(\cdot|h)$. The utility $u : N \times H_{TS} \rightarrow \mathbb{R}$ assigns a payoff
 083 to each player at every terminal history. For player i , \mathcal{I}_i is a partition of $\{h \in H : P(h) = i\}$; each
 084 $I_i \in \mathcal{I}_i$ is an information set and represents observable information shared by all $h \in I_i$. Due to
 085 indistinguishability, $A(h) = A(I_i)$, $P(h) = P(I_i)$. Our work focuses on the two-player zero-sum
 086 setting, where $N = \{1, 2\}$ and $u_1(h) = -u_2(h)$, $\forall h \in H_{TS}$, as in prior CFR works Brown et al.
 (2019); Davis et al. (2020).

087 Every player $i \in N$ selects actions according to a strategy σ_i that maps $I_i \in \mathcal{I}_i$ to a dis-
 088 tribution over $A(I_i)$, and $\sigma_i(\cdot|h) = \sigma_i(\cdot|I_i)$, $\forall h \in I_i$. CFR aims to find a Nash Equilib-
 089 rium (NE) strategy profile $\sigma^* = \{\sigma_1^*, \dots, \sigma_N^*\}$, where no player has an incentive to deviate:
 090 $u_i(\sigma^*) \geq \max_{\sigma_i} u_i(\{\sigma_i, \sigma_{-i}^*\})$, $\forall i \in N$, where $-i$ denotes the players other than i , and $u_i(\sigma)$ is
 091 the expected payoff of player i :

$$092 \quad u_i(\sigma) = \sum_{h' \in H_{TS}} u_i(h') \pi^\sigma(h'), \quad \pi^\sigma(h') = \prod_{(ha) \sqsubseteq h'} \sigma_{P(h)}(a|I(h)) \quad (1)$$

095 $I(h)$ denotes the information set containing h , and $\pi^\sigma(h)$ is the reach probability of h under σ .
 096 $\pi^\sigma(h)$ can be decomposed as $\prod_{i \in N \cup \{c\}} \pi_i^\sigma(h)$, where $\pi_i^\sigma(h) = \prod_{(ha) \sqsubseteq h', P(h)=i} \sigma_i(a|I(h))$. In
 097 addition, $\pi^\sigma(I) = \sum_{h \in I} \pi^\sigma(h)$ is the reach probability of I .

098 CFR Zinkevich et al. (2007) iteratively accumulates the counterfactual regret $R_i^T(a|I)$ for player i at
 099 each information set $I \in \mathcal{I}_i$:

$$101 \quad R_i^T(a|I) = \frac{1}{T} \sum_{t=1}^T \pi_{-i}^{\sigma^t}(I) (u_i(\sigma^t|_{I \rightarrow a}, I) - u_i(\sigma^t, I)), \quad u_i(\sigma, I) = \sum_{h \in I} \pi_{-i}^\sigma(h) \sum_{h' \in H_{TS}} \pi^\sigma(h, h') u_i(h') / \pi_{-i}^\sigma(I) \quad (2)$$

104 Here, σ^t is the strategy profile at iteration t , $\sigma^t|_{I \rightarrow a}$ is identical to σ^t except always choosing a
 105 at I , and $\pi^\sigma(h, h')$ is the reach probability from h to h' (equals $\frac{\pi^\sigma(h')}{\pi^\sigma(h)}$ if $h \sqsubseteq h'$ and 0 otherwise).
 106 $R_i^T(a|I)$ is the expected regret of not choosing a at I . By regret matching Abernethy et al. (2011),
 107 the next strategy $\sigma_i^{T+1}(\cdot|I)$ is defined as $\sigma_i^{T+1}(a|I) \propto \max(R_i^T(a|I), 0)$. After T iterations, the

average strategy is $\bar{\sigma}_i^T(a|I) \propto \sum_{t=1}^T \pi_i^{\sigma^t}(I) \sigma_i^t(a|I)$. CFR guarantees that the average profile $\bar{\sigma}^T = \{\bar{\sigma}_i^T | i \in N\}$ converges to a Nash Equilibrium as $T \rightarrow \infty$.

2.2 THE OPTION FRAMEWORK

An option Sutton et al. (1999) $z \in \mathcal{Z}$ is specified by three components: an initiation set $Init_z \subseteq \mathcal{S}$, an intra-option policy $\sigma_z(a|s) : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$, and a termination function $\beta_z(s) : \mathcal{S} \rightarrow [0, 1]$. \mathcal{S} , \mathcal{A} , \mathcal{Z} denote the state, action, and option spaces, respectively. Option z is available in state s iff $s \in Init_z$. Once taken, actions follow σ_z until stochastic termination by β_z . A high-level policy $\sigma_{\mathcal{Z}}(z|s) : \mathcal{S} \times \mathcal{Z} \rightarrow [0, 1]$ activates a new option when the previous one terminates. Thus, $\sigma_{\mathcal{Z}}(z|s)$ and $\sigma_z(a|s)$ form a hierarchical policy. Hierarchical policies tend to have superior performance on complex long-horizon tasks that can be decomposed into subtasks.

The one-step option framework Jing et al. (2021a) learns hierarchical policies without explicitly defining $Init_z$ or β_z . First, it assumes every option is available at every state, i.e., $Init_z = \mathcal{S}, \forall z \in \mathcal{Z}$. Second, it redefines the high-level and low-level policies as $\sigma^H(z | s, z')$ and $\sigma^L(a | s, z)$, respectively:

$$\sigma^H(z | s, z') = \beta_{z'}(s) \sigma_{\mathcal{Z}}(z | s) + (1 - \beta_{z'}(s)) \delta_{z=z'}, \quad \sigma^L(a | s, z) = \sigma_z(a | s) \quad (3)$$

where z' is the option in the previous timestep and $\delta_{z=z'}$ is the indicator function. Thus, if the previous option terminates (with probability $\beta_{z'}(s)$), the agent samples a new option via $\sigma_{\mathcal{Z}}(z | s)$; otherwise, it continues with z' . The authors of Li et al. (2021a) implement the high-level policy as an end-to-end neural network with Multi-Head Attention (MHA) Vaswani et al. (2017), enabling temporal extension of options without an explicit β_z . Intuitively, if z' still fits s , $\sigma^H(z | s, z')$ assigns higher attention to z' and tends to continue; otherwise, it samples a more compatible option. In this one-step framework, the option is sampled at each timestep rather than only upon termination, so we train the hierarchical policy σ^H and σ^L directly.

3 METHODOLOGY

We extend CFR to learn a hierarchical strategy with neural networks (NNs) for IIGs with extensive state spaces and deep game trees. Low-level components represent skills (options) over primitive actions, while the high-level strategy coordinates their use; thus, they are learned as separate functions. Given the lack of prior hierarchical CFR, we proceed incrementally. **First**, we define the hierarchical strategy and hierarchical counterfactual regret and give CFR-style updates with a convergence guarantee. **Second**, we propose a Low-Variance Monte Carlo sampling extension to handle vast or unknown game trees—where standard traversal is impractical—without weakening the convergence rate. **Finally**, we develop HDCFR by approximating these hierarchical functions with NNs and training them via objectives that match the tabular updates under exact fitting, preserving the established theory.

3.1 PRELIMINARIES

In the extended game model, at each time step t , each player i makes its decision by selecting a hierarchical action $\tilde{a}_t \triangleq (z_t, a_t)$, which consists of the option and primitive action, based on the observable information, i.e. I_i . With the hierarchical actions, we can redefine the IIG model as $\langle N, H, \tilde{A}, P, \sigma_c, u, \mathcal{I} \rangle$. Here, N , P and u retain the definitions in Section 2.1. H includes all the possible histories, each of which is a sequence of hierarchical actions of all players starting from the first time step. Consequently, \mathcal{I} denotes the collection of information sets induced by the new H , containing all observable history. $\tilde{A}(h) = Z(h) \times A(h)$, where $Z(h)$ and $A(h)$ represent the options and primitive actions available at h , respectively. No action pruning: options only condition the policy and do not change $A(h)$. $\sigma_c((z_c, a)|h) = \sigma_c(a|h)$, where $\sigma_c(a|h)$ is the predefined distribution in the original game model and z_c (a dummy variable) is the only option choice for the chance player.

Each player i possesses a hierarchical strategy $\sigma_i(\tilde{a}_t|I_i)$, which, by the chain rule, equals $\sigma_i^H(z_t|I_i) \cdot \sigma_i^L(a_t|I_i, z_t)$. Note that although I_i includes $z_{1:t-1}$, we follow the conditional independence assumption of the one-step option framework Zhang & Whiteson (2019): $z_t \perp\!\!\!\perp z_{1:(t-2)} \mid z_{t-1}$ and $a_t \perp\!\!\!\perp z_{1:(t-1)} \mid z_t$, thus only $z_{t-1}(z_t)$ is used for $\sigma_i^H(\sigma_i^L)$ to determine $z_t(a_t)$. With the

162 hierarchical strategy, we can redefine the expected payoff and reach probability in Eq (1) by simply
 163 substituting a with \tilde{a} , based on which we have the definition of the average overall regret of player i
 164 at iteration T : (From this point forward, t refers to a learning iteration rather than a time step within
 165 an iteration.)

$$166 R_{full,i}^T = \frac{1}{T} \max_{\sigma_i'} \sum_{t=1}^T (u_i(\{\sigma_i', \sigma_{-i}^t\}) - u_i(\sigma^t)) \quad (4)$$

169 The following theorem (Theorem 2 from Zinkevich et al. (2007)) provides a connection between the
 170 average overall regret and the Nash Equilibrium solution.

171 **Theorem 1.** *In a two-player zero-sum game at iteration T , if both players' average overall regret is*
 172 *less than ϵ , then $\bar{\sigma}^T = \{\bar{\sigma}_1^T, \bar{\sigma}_2^T\}$ is a 2ϵ -Nash Equilibrium.*

174 The average strategy $\bar{\sigma}_i^T$ is defined as Eq (5) ($\forall i \in N, I \in \mathcal{I}_i, \tilde{a} \in \tilde{A}(I)$). An ϵ -Nash Equilibrium
 175 σ approximates an NE, with the property that $u_i(\sigma) + \epsilon \geq \max_{\sigma_i'} u_i(\{\sigma_i', \sigma_{-i}\})$, $\forall i \in N$. Thus,
 176 ϵ measures the distance of σ to the NE in expected payoff. Then, according to Theorem 1, as
 177 $R_{full,i}^T \rightarrow 0$ ($\forall i \in N$), $\bar{\sigma}^T$ converges to an NE. Theorem 1 can be applied directly to our hierarchical
 178 setting, as the only difference in derivations is the replacement of a with \tilde{a} in $R_{full,i}^T$ and $\bar{\sigma}_i^T(\tilde{a}|I)$.

$$180 \bar{\sigma}_i^T(\tilde{a}|I) = \left(\sum_{t=1}^T \pi_i^{\sigma^t}(I) \sigma_i^t(\tilde{a}|I) \right) / \sum_{t=1}^T \pi_i^{\sigma^t}(I) \quad (5)$$

183 3.2 HIERARCHICAL CFR

184 A direct approach is to view $\sigma_i(\tilde{a}|I) = \sigma_i^H(z|I) \cdot \sigma_i^L(a|I, z)$ as a unified strategy on \tilde{A} and run CFR.
 185 However, this hides the separation of levels, hindering skill reuse or human-initialized skills. We
 186 therefore introduce Hierarchical CFR (HCFR) to learn $\sigma_i^H(z|I)$ and $\sigma_i^L(a|I, z)$ separately for all
 187 $I \in \mathcal{I}_i, z \in Z(I), a \in A(I)$, and provide a convergence guarantee.

189 As noted in Section 3.1, achieving an NE requires minimizing the average overall regrets $R_{full,i}^T$.
 190 Instead of directly optimizing $R_{full,i}^T$, we minimize its upper bound (Theorem 2)—the sum of high-
 191 level and low-level **counterfactual regrets** at each information set: $R_i^{T,H}(z|I)$ and $R_i^{T,L}(a|I, z)$.
 192 This lets us decouple learning by independently minimizing $R_i^{T,H}(z|I)$ and $R_i^{T,L}(a|I, z)$ via σ_i^H and
 193 σ_i^L .

194 **Theorem 2.** *With the following definitions of high-level and low-level counterfactual regrets:*

$$196 R_i^{T,H}(z|I) = \frac{1}{T} \sum_{t=1}^T \pi_i^{\sigma^t}(I) (u_i(\sigma^t|_{I \rightarrow z}, I) - u_i(\sigma^t, I))$$

$$198 R_i^{T,L}(a|I, z) = \frac{1}{T} \sum_{t=1}^T \pi_i^{\sigma^t}(I) (u_i(\sigma^t|_{Iz \rightarrow a}, Iz) - u_i(\sigma^t, Iz)) \quad (6)$$

202 we have $R_{full,i}^T \leq \sum_{I \in \mathcal{I}_i} \left[R_{i,+}^{T,H}(I) + \sum_{z \in Z(I)} R_{i,+}^{T,L}(I, z) \right]$.

204 We introduce the new notations here: $\sigma^t|_{Iz \rightarrow a}$ is a hierarchical strategy profile identical to σ^t
 205 except that the intra-option strategy of option z at I is always choosing a ; $u_i(\sigma^t, Iz)$ is the ex-
 206 pected payoff for choosing option z at I ; $R_{i,+}^{T,H}(I) = \max(\max_z R_i^{T,H}(z|I), 0)$, $R_{i,+}^{T,L}(I, z) =$
 207 $\max(\max_a R_i^{T,L}(a|I, z), 0)$. Proof of Theorem 2 is in Appendix B.

209 After obtaining the regrets $R_i^{T,H}$ and $R_i^{T,L}$, we update the high- and low-level strategies for the next
 210 iteration as follows: ($\forall i \in N, I \in \mathcal{I}_i, z \in Z(I), a \in A(I), \mu^H$ and μ^L are normalizing factors.)

$$212 \sigma_i^{T+1,H}(z|I) = \begin{cases} R_{i,+}^{T,H}(z|I)/\mu^H, & \mu^H > 0, \\ 1/|Z(I)|, & \text{o} \setminus w. \end{cases}, \sigma_i^{T+1,L}(a|I, z) = \begin{cases} R_{i,+}^{T,L}(a|I, z)/\mu^L, & \mu^L > 0, \\ 1/|A(I)|, & \text{o} \setminus w. \end{cases} \quad (7)$$

215 Thus, regrets and strategies are iteratively computed (i.e., $\sigma^{1:t} \rightarrow R^t \rightarrow \sigma^{t+1}$, $\sigma^{1:t+1} \rightarrow R^{t+1} \rightarrow$
 σ^{t+2}) with Eq (6) and (7) until convergence (i.e., $R_{full,i}^T \rightarrow 0$). The convergence rate is:

Theorem 3. *If player i selects options and actions according to Eq (7), then $R_{full,i}^T \leq \Delta_{u,i} |\mathcal{I}_i| (\sqrt{|Z_i|} + |Z_i| \sqrt{|A_i|}) / \sqrt{T}$, where $\Delta_{u,i} = \max_{h' \in H_{TS}} u_i(h') - \min_{h' \in H_{TS}} u_i(h')$, $|\mathcal{I}_i|$ is the number of information sets for player i , $|A_i| = \max_{h: P(h)=i} |A(h)|$, $|Z_i| = \max_{h: P(h)=i} |Z(h)|$.*

Thus, as $T \rightarrow \infty$, $R_{full,i}^T \rightarrow 0$. The rate $\mathcal{O}(T^{-0.5})$ matches CFR Zinkevich et al. (2007), so introducing options preserves convergence while enabling skill-based learning. Proof is in Appendix C.

After T iterations, the average high-level and low-level strategies are: (See Appendix D for proof.)

$$\bar{\sigma}_i^{T,H}(z|I) = \frac{\sum_{t=1}^T \pi_i^{\sigma^t}(I) \sigma_i^{t,H}(z|I)}{\sum_{t=1}^T \pi_i^{\sigma^t}(I)}, \quad \bar{\sigma}_i^{T,L}(a|I, z) = \frac{\sum_{t=1}^T \pi_i^{\sigma^t}(Iz) \sigma_i^{t,L}(a|I, z)}{\sum_{t=1}^T \pi_i^{\sigma^t}(Iz)} \quad (8)$$

where $\pi_i^{\sigma^t}(Iz) = \pi_i^{\sigma^t}(I) \sigma_i^{t,H}(z|I)$. Then:

Proposition 1. *If both players sequentially use their average high-level and low-level strategies following the one-step option model, i.e., $\forall I \in \mathcal{I}_i$, selecting an option z according to $\bar{\sigma}_i^{T,H}(\cdot|I)$ and then selecting the action a according to the corresponding intra-option strategy $\bar{\sigma}_i^{T,L}(\cdot|I, z)$, the resulting strategy profile converges to an NE as $T \rightarrow \infty$.*

3.3 HIERARCHICAL DEEP CFR

In vanilla CFR, regrets are updated for every information set each iteration, which is infeasible at scale. Monte Carlo CFR (MCCFR) Lanctot et al. (2009) instead updates on sampled parts of the tree. OS updates along a single sampled trajectory; ES explores all actions of the traverser and samples a single action for non-traversers, requiring an explicit traversable game model and growing exponentially with horizon. Our setting features deep trees and simulator-only access, so we adopt OS, but it suffers from high variance. **We therefore complete Section 3.2 by adding a low-variance outcome-sampling extension: replace counterfactual regrets with unbiased estimators and introduce baselines to reduce variance (Appendix E).**

Building on Section 3.2 and Appendix E, we present HDCFR, which uses NNs to approximate the regret, strategy, and baseline functions, enabling infinite-scale state spaces. **We define deep objectives that match the tabular targets under exact fitting; in practice, approximation/optimization errors may affect regret. We then present the complete algorithm in pseudo-code.** We represent each skill z as a learnable embedding and implement the subpolicy $\pi_z(a|I)$ as a conditional network on $[I; z]$ that outputs over the full $A(I)$ (with I containing both private and public observations). In particular, we train three types of networks: the counterfactual regret networks $R_{i,\theta}^{t,H}$, $R_{i,\theta}^{t,L}$, average strategy networks $\bar{\sigma}_{i,\phi}^{T,H}$, $\bar{\sigma}_{i,\phi}^{T,L}$, and baseline network b^t . Notably, we do not maintain the baselines b^t for each player. Instead, we leverage the property of two-player zero-sum games where the payoff of the two players offsets each other. That is, $b^t = b_1^t = -b_2^t$.

First, the counterfactual regret networks are trained by minimizing the following two objectives: $\mathcal{L}_{R,i}^{t,H}$ and $\mathcal{L}_{R,i}^{t,L}$.

$$\mathbb{E}_{(I, \hat{r}_i^{t',H}) \sim \tau_R^i} \left[\sum_{z \in Z(I)} (R_{i,\theta}^{t,H}(z|I) - \hat{r}_i^{t',H}(I, z))^2 \right], \quad \mathbb{E}_{(Iz, \hat{r}_i^{t',L}) \sim \tau_R^i} \left[\sum_{a \in A(I)} (R_{i,\theta}^{t,L}(a|I, z) - \hat{r}_i^{t',L}(Iz, a))^2 \right] \quad (9)$$

Here, τ_R^i stores sampled immediate counterfactual regrets $\hat{r}_i^{t'}$ from iterations 1 to t (Appendix E). The core idea of MCCFR is to replace $R_{i,\theta}^{t,H}$ and $R_{i,\theta}^{t,L}$ with their unbiased Monte Carlo estimates. As a justification of this design:

Proposition 2. *Let $R_{i,*}^{t,H}$ and $R_{i,*}^{t,L}$ denote the minimal points of $\mathcal{L}_{R,i}^{t,H}$ and $\mathcal{L}_{R,i}^{t,L}$, respectively. For all $I \in \mathcal{I}_i$, $z \in Z(I)$, $a \in A(I)$, $R_{i,*}^{t,H}(z|I)$ and $R_{i,*}^{t,L}(a|I, z)$ yield unbiased estimations of the true counterfactual regrets scaled by positive constant factors, i.e., $C_1 R_{i,*}^{t,H}(z|I)$ and $C_2 R_{i,*}^{t,L}(a|I, z)$. Specifically, $R_{i,*}^{t,H}(z|I) \rightarrow C_1 R_i^{t,H}(z|I)$ and $R_{i,*}^{t,L}(a|I, z) \rightarrow C_2 R_i^{t,L}(a|I, z)$, as $|\tau_R^i| \rightarrow \infty$.*

Please refer to Appendix I for the proof. Since regrets are only used to compute the next strategy via Eq (7), the positive scales C_1, C_2 cancel in normalization; thus $R_{i,*}^{t,H}$ and $R_{i,*}^{t,L}$ can replace $R_{i,\theta}^{t,H}$ and $R_{i,\theta}^{t,L}$.

Second, the average strategy networks are learned from immediate strategies across iterations 1 to T by minimizing $\mathcal{L}_{\bar{\sigma},i}^H$ and $\mathcal{L}_{\bar{\sigma},i}^L$:

$$\mathbb{E}_{(I, \sigma_i^{t,H}) \sim \tau_{\bar{\sigma}}^i} \left[\sum_{z \in Z(I)} (\bar{\sigma}_{i,\phi}^{T,H}(z|I) - \sigma_i^{t,H}(z|I))^2 \right], \quad \mathbb{E}_{(I, z, \sigma_i^{t,L}) \sim \tau_{\bar{\sigma}}^i} \left[\sum_{a \in A(I)} (\bar{\sigma}_{i,\phi}^{T,L}(a|I, z) - \sigma_i^{t,L}(a|I, z))^2 \right] \quad (10)$$

We adopt a sampling scheme to fulfill the following proposition. Define $q^{t,i}$ as the sample strategy profile at iteration t when i is the traverser; $q_p^{t,i}$ is uniformly random when $p = i$ and equals σ_p^t when $p = 3 - i$. Samples in $\tau_{\bar{\sigma}}^i$ are gathered when the traverser is $3 - i$ (so i uses σ_i^t). Then: (See Appendix J.)

Proposition 3. *Let $\bar{\sigma}_{i,*}^{T,H}$ and $\bar{\sigma}_{i,*}^{T,L}$ be the minimal points of $\mathcal{L}_{\bar{\sigma},i}^H$ and $\mathcal{L}_{\bar{\sigma},i}^L$, and let $\tau_{\bar{\sigma}}^{t,i}$ be the partition of $\tau_{\bar{\sigma}}^i$ at iteration t . If $\tau_{\bar{\sigma}}^{t,i}$ follows the scheme above, then $\bar{\sigma}_{i,*}^{T,H}(z|I) \rightarrow \bar{\sigma}_{i,*}^{T,H}(z|I)$ and $\bar{\sigma}_{i,*}^{T,L}(a|I, z) \rightarrow \bar{\sigma}_{i,*}^{T,L}(a|I, z)$, $\forall I \in \mathcal{I}_i$, $z \in Z(I)$, $a \in A(I)$, as $|\tau_{\bar{\sigma}}^{t,i}| \rightarrow \infty$ ($t \in \{1, \dots, T\}$).*

By Propositions 1 and 3, $\bar{\sigma}_*^{T,H}$ and $\bar{\sigma}_*^{T,L}$ can be returned as an approximate NE.

Last, at the end of each iteration t , we learn the baseline function for iteration $t+1$ to reduce variance by minimizing:

$$\mathcal{L}_b^{t+1} = \mathbb{E}_{h' \sim \tau_b^t} \left[\sum_{hza \sqsubseteq h'} (\hat{b}^{t+1}(h, z, a) - \hat{b}^{t+1}(hza|h'))^2 \right] \quad (11)$$

Here, τ_b^t stores trajectories collected at iteration t when player 1 is the traverser. For each trajectory, we record sampled baselines $\hat{b}^{t+1}(h|h')$, $\forall h \sqsubseteq h'$, recursively defined as: $\hat{b}^{t+1}(h|h') = u_1(h)$ if $h \in H_{TS}$

$$\begin{aligned} \hat{b}^{t+1}(h|h') &= \sum_{z \in Z(h)} \sigma_{P(h)}^{t+1,H}(z|h) \hat{b}^{t+1}(h, z|h'), \quad \hat{b}^{t+1}(hza|h') = \sum_{a \in A(h)} \sigma_{P(h)}^{t+1,L}(a|h, z) \hat{b}^{t+1}(hza|h'); \\ \hat{b}^{t+1}(h, z|h') &= \frac{\delta(hz \sqsubseteq h')}{q^{t,1}(z|h)} [\hat{b}^{t+1}(hza|h') - b^t(h, z)] + b^t(h, z), \\ \hat{b}^{t+1}(hza|h') &= \frac{\delta(hza \sqsubseteq h')}{q^{t,1}(a|h, z)} [\hat{b}^{t+1}(hza|h') - b^t(h, z, a)] + b^t(h, z, a). \end{aligned} \quad (12)$$

The high-level baseline $b^{t+1}(h, z)$ is defined from the low-level baseline $b^{t+1}(h, z, a)$ as $b^{t+1}(h, z) = \sum_{a \in A(h)} \sigma_{P(h)}^{t+1,L}(a|I(h), z) b^{t+1}(h, z, a)$. With these sampled baselines and their relation, we have:

Proposition 4. *Denote $b^{t+1,*}$ as the minimal point of \mathcal{L}_b^{t+1} and consider trajectories in τ_b^t as independent and identically distributed random samples. We have $b^{t+1,*}(h, z, a) \rightarrow v^{t+1,H}(\sigma^{t+1}, hza)$ and $b^{t+1,*}(h, z) = \sum_{a'} \sigma_{P(h)}^{t+1,L}(a'|I(h), z) b^{t+1,*}(h, z, a') \rightarrow v^{t+1,L}(\sigma^{t+1}, hz)$, $\forall h \in H$, $z \in Z(h)$, $a \in A(h)$, as $|\tau_b^t| \rightarrow \infty$.*

Thus, the ideal criteria for baselines (Theorem 5 in Appendix E) are satisfied at the optimum of \mathcal{L}_b^{t+1} (Appendix K). **In particular, the variance of estimating R_i^t is minimized w.r.t. these baselines at the optimum.**

To sum up, we present the pseudo code of HDCFR as Algorithm 1 and 2 in Appendix N. There are T iterations. **(1)** At iteration t , players alternate as traverser and collect K trajectories (Algorithm 1, L6–11) via outcome sampling (Algorithm 2). Along each trajectory, immediate counterfactual regrets for the traverser i (i.e., \hat{r}_i^t) are computed by Eq (26), (27) (Appendix E) and stored in τ_R^i ; non-traverser strategies σ_{3-i}^t are derived from $R_{3-i,\theta}^{t-1}$ using Eq (7) and saved in $\tau_{\bar{\sigma}}^{3-i}$. **(2)** At the end of t , train $R_{i,\theta}^{t-1}$ on τ_R^i via Eq (9) to obtain $R_{i,\theta}^t$ (Algorithm 1, L12–14). Then update the baseline b^t to b^{t+1} by Eq (11) (L22–30). b^{t+1} and $R_{i,\theta}^t$ are used in the next iteration. **(3)** After T iterations, train $\bar{\sigma}_{\phi}^T$ on $\tau_{\bar{\sigma}}$ via Eq (10) (L17–19) and return it as an approximate NE.

4 EVALUATION AND MAIN RESULTS

In Section 4.1, we benchmark HDCFR against leading model-free methods for imperfect-information zero-sum games, including DREAM Steinberger et al. (2020), OSSDCFR (an outcome-sampling

324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377

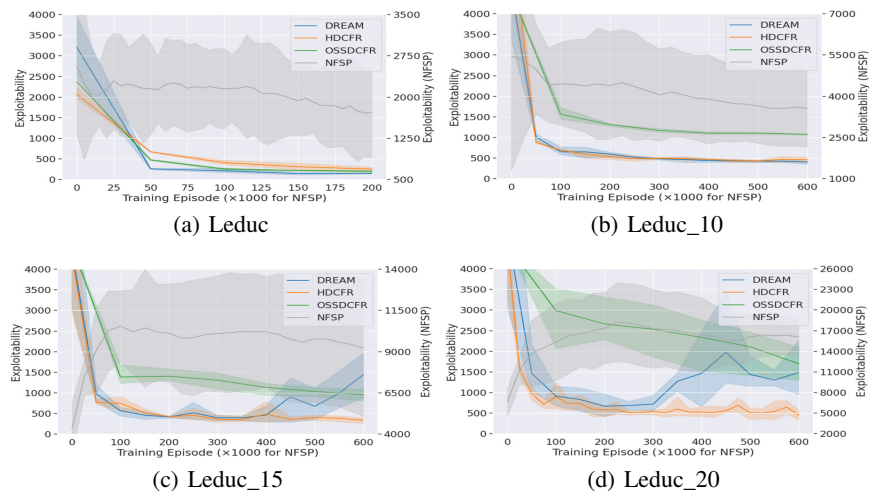


Figure 1: Performance comparison on Leduc poker games. Lower exploitability indicates a closer approximation to the Nash Equilibrium. HDCFR exhibits superior convergent performance as the game’s decision horizon increases; All curves are averaged over 5 random seeds

variant of Deep CFR) Steinberger (2019); Brown et al. (2019), and NFSP Heinrich & Silver (2016). Like HDCFR, these algorithms do not require domain knowledge and can be applied in environments with unknown game models (i.e., the model-free setting). For evaluation benchmarks, as a common practice, we select poker games: Leduc Southey et al. (2005) and heads-up flop hold’em (FHP) Brown et al. (2019). Given its hierarchical design, HDCFR is expected to be superior in tasks involving extended decision-making horizons. **Thus, we elevate complexity of the standard poker benchmarks by raising the number of cards and the cap on the total raises and accordingly increasing the initial stack size for each player, compelling agents to strategize over longer horizons. Details of the SIX benchmarks used are available in Table 3 in Appendix L.** Then, in Section 4.2, we analyze the hierarchical strategy learned by HDCFR. We examine whether the high-level strategy can temporally extend skills and if the low-level strategies (i.e., skills) can be transferred to new tasks to aid learning. Notably, we utilize the baseline and benchmark implementation from Steinberger (2020).

4.1 COMPARISON WITH STATE-OF-THE-ART MODEL-FREE ALGORITHMS ON ZERO-SUM IIGS

For Leduc poker games, we can explicitly compute the best response (BR) function for the learned strategy profile $\sigma = \{\sigma_1, \sigma_2\}$. We then can use the exploitability of σ : $\text{exploitability}(\sigma) = 1/2 \max_{\sigma'} [u_1(\sigma'_1, \sigma_2) + u_2(\sigma_1, \sigma'_2)]$, as the learning performance metric. Commonly used in extensive-form games, exploitability measures the deviation from an NE, so **a lower value is preferred**. For hold’em poker games (like our benchmarks), exploitability is usually quantified in milli big blinds per game (mbb/g). In Figure 1, we show the learning curves of HDCFR and the baselines. Solid lines represent the mean, while shadowed areas indicate the 95% confidence intervals from repeated trials. **(1)** For CFR-based algorithms, the players sample 900 trajectories in each training episode, and visit around 10^7 game states in the learning process. In contrast, the RL-based NFSP algorithm is trained over more episodes ($\times 1000$) and the players visit 10^8 game states in total during training. However, NFSP consistently underperforms across all benchmarks, so it is significantly less sample-efficient compared to the CFR-based algorithms. Note that NFSP uses a separate y-axis.

(2) In the absence of game models, backtracking is not allowed and so the player can sample only one action at each information set, which is known as outcome sampling. Thus, algorithms that require backtracking, like Deep CFR Brown et al. (2019) and DNCFR Li et al. (2020), cannot work directly, unless adapted with the outcome sampling scheme. However, the performance of the resulting algorithm OSSDCFR declines significantly with increasing game complexity, primarily due to the high sample variance associated with outcome sampling. **(3)** With variance reduction techniques, DREAM achieves comparable performance to HDCFR in simpler scenarios. HDCFR, owing to

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

Table 1: HDCFR vs. Baseline Algorithms in Head-to-Head Matchups. The table rows represent the average payoffs (higher is better) of HDCFR against the baselines: Leduc, Leduc_10, Leduc_15, Leduc_20, FHP, and FHP_10.

Benchmark	DREAM	OSSDCFR	NFSP
Leduc	-11.94 ± 53.79	4.11 ± 64.03	596.55 ± 73.46
Leduc_10	-14.22 ± 62.10	500.0 ± 73.22	642.67 ± 109.41
Leduc_15	171.33 ± 70.80	563.75 ± 83.31	1351.5 ± 207.27
Leduc_20	196.89 ± 76.69	587.0 ± 68.83	1725.33 ± 206.01
FHP	184.58 ± 36.75	68.11 ± 36.61	244.61 ± 41.36
FHP_10	282.42 ± 14.20	343.22 ± 15.35	537.39 ± 16.91

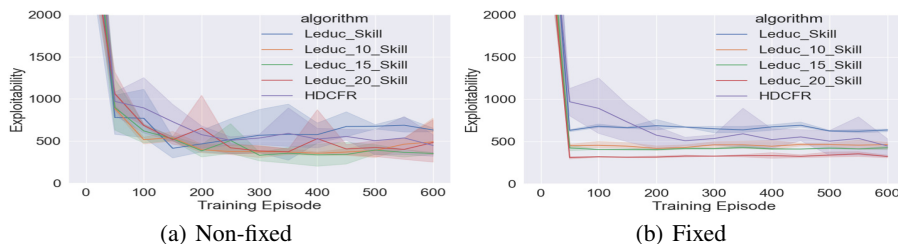


Figure 2: Learning performance on Leduc_20 with transferred skills from other Leduc tasks. The transferred skills can either be fixed or not when learning a hierarchical strategy on the new scenario. As a reference, the learning performance without transferred skills is labeled as HDCFR. By fixing pre-learned skills, the agent focuses on mastering a high-level strategy, thus accelerating learning. However, adjusting these skills alongside the high-level strategy can yield improved results: when using Leduc_15 skills, peaked performance occurs around episode 400.

its hierarchical structure, excels over DREAM in games with extended horizons, where DREAM struggles to converge. HDCFR’s superiority becomes more significant as the game complexity increases. On Leduc_15/20, DREAM’s exploitability occasionally decreases and then increases—this occurs in long-horizon settings where outcome sampling has high variance. Our ablations (Appendix M, Fig. 3a) reproduce a similar pattern when removing the baseline or the multi-head attention module (no temporally extended skills), indicating that variance control and temporal abstraction are crucial for stability in deep horizons.

Further, we conduct head-to-head tournaments between HDCFR and each baseline. We select the top three checkpoints for each algorithm, resulting in total nine pairings for each tournament. Each pair of strategies competes over 1,000 hands. Table 1 shows the average payoff of HDCFR’s strategy profile (calculated as $1/2 [u_1(\sigma_1^{\text{HDCFR}}, \sigma_2^{\text{baseline}}) + u_2(\sigma_1^{\text{baseline}}, \sigma_2^{\text{HDCFR}})]$), along with 95% confidence intervals, measured in mbb/g. **A higher payoff indicates superior decision-making performance.** Observations from Leduc poker games in this table (i.e., rows 1-4) align with aforementioned conclusions (1)-(3). **To further show the superiority of our algorithm, we compare it with baselines on larger-scale FHP games, which boast a game tree exceeding 10^{12} nodes in size.** Due to the immense scale of FHP games, computing best response functions is impractical, so we offer only head-to-head comparison results. Training an instance on FHP games requires 7 days using a device with 8 CPU cores (3rd Gen Intel Xeon) and 128 GB of RAM. We utilized the RAY parallel computing framework Moritz et al. (2018). We can see that HDCFR’s advantage grows as task difficulty goes up.

4.2 ANALYSIS ON THE LEARNED HIERARCHICAL STRATEGY

One key benefit of hierarchical learning is to use prelearned skills as building blocks for strategy learning, providing a manner for integrating expert knowledge. Even in the absence of domain knowledge, where rule-based skills can’t be provided as expert guidance, we can leverage skills learned from similar games. Skills, as policy segments, often exhibit greater transferability than complete strategies. In Figure 2, we demonstrate the transfer of skills from various Leduc games to Leduc_20 and present the learning outcomes. For comparison, we also include the performance

Table 2: Comparison of the switch frequencies when using skills from different source tasks.

Source Task	Leduc	Leduc_10	Leduc_15	Leduc_20
Switch Frequency	0.1363 $\pm 4.02 \times 10^{-4}$	0.1256 $\pm 3.43 \times 10^{-4}$	0.1088 $\pm 2.16 \times 10^{-4}$	0.1016 $\pm 3.64 \times 10^{-4}$

without transferred skills, labeled as HDCFR. We define skill-switch frequency as the ratio between the number of time steps at which the selected option changes and the total number of time steps; its inverse approximates the average option duration. These prelearned skills can either remain static (Figure 2(b)) or be trained alongside the high-level strategy (Figure 2(a)). When kept static, the agent can focus on mastering its high-level strategy to select among a set of effective skills, resulting in quicker convergence. Notably, the superior convergent performance observed in Figure 2(b) positively correlates with the similarity between the source task of the skills and Leduc_20. On the other hand, if the skills are adjusted with the high-level strategy, the improvement on the convergence speed may not be obvious, but skills can be more customized for the current task and better performance may be achieved. For instance, with Leduc_15 skills, peak performance is achieved around episode 400; for Leduc skills, training with dynamic skills (Figure 2(a)) results in better performance compared to static skills (Figure 2(b)). However, for Leduc_20 skills, fixed skills works better. This could be because they originate from the same task, eliminating the need for further adaptation.

Next, we provide an analysis of the learned high-level strategies. As depicted in Figure 2(b), when utilizing fixed skills from different source tasks, corresponding high-level strategies can be learned. To evaluate whether these high-level strategies can extend skills temporally (with the attention mechanism) – instead of frequently toggling between skills – we calculate the frequency of skill switches in the game tree of Leduc_20 (containing 113954 nodes in total), considering all possible hands of cards and five repeated experiments. Table 2 reports the means and 95% confidence intervals of the skill switch frequency. As the decision horizon of the skills’ source task increases, the switch frequency decreases, reflecting longer durations of single-skill utilization. Notably, for Leduc_20 skills, skill switches occur only about 10% of the time. This indicates the agent’s preference for decision-making at an extended-skill level, with an average skill duration of approximately 10 steps, rather than at the level of primitive actions, aligning with our expectations. **An ablation study is provided in Appendix M to highlight the importance of each component within our algorithm.**

5 CONCLUSION

We introduce the first hierarchical extension of CFR based on the option framework. We first establish the theoretical foundations in a tabular setting and then extend them using neural networks as function approximators, resulting in a theoretically grounded deep learning algorithm – HDCFR. Evaluations in complex two-player zero-sum games show that HDCFR outperforms leading algorithms in this field, with the advantage increasing as the decision horizon grows, underscoring its potential for tasks involving deep game trees. Moreover, we show that the learned high-level strategy can *temporally* extend skills to exploit hierarchical subtask structure in long-horizon tasks, and that the learned skills can be transferred to related tasks to facilitate learning. Our algorithm provides a novel framework to learn with pre-learned skills in zero-sum IIGs.

Ethics Statement This work complies with the ICLR Code of Ethics. While our methods are general, they may be applied in contexts with societal implications, including risks related to bias, fairness, and privacy. We encourage responsible use and declare no conflicts of interest.

Reproducibility Statement We provide detailed descriptions of our methodology, datasets, model configurations, and evaluation metrics in both the main text and the Appendix. In addition, the complete source code is included in the supplementary materials to facilitate reproducibility.

REFERENCES

Jacob Abernethy, Peter L Bartlett, and Elad Hazan. Blackwell approachability and no-regret learning are equivalent. In *Proceedings of the 24th Annual Conference on Learning Theory*, pp. 27–46. JMLR Workshop and Conference Proceedings, 2011.

486 Anton Bakhtin, Noam Brown, Emily Dinan, Gabriele Farina, Colin Flaherty, Daniel Fried, Andrew
487 Goff, Jonathan Gray, Hengyuan Hu, Athul Paul Jacob, Mojtaba Komeili, Karthik Konath, Minae
488 Kwon, Adam Lerer, Mike Lewis, Alexander H. Miller, Sasha Mitts, Adithya Renduchintala,
489 Stephen Roller, Dirk Rowe, Weiyan Shi, Joe Spisak, Alexander Wei, David Wu, Hugh Zhang, and
490 Markus Zijlstra. Human-level play in the game of diplomacy by combining language models with
491 strategic reasoning. *Science*, pp. 1067–1074, 2022.

492 Michael Bowling, Neil Burch, Michael Johanson, and Oskari Tammelin. Heads-up limit hold'em
493 poker is solved. *Science*, 347(6218):145–149, 2015.

494

495 Noam Brown and Tuomas Sandholm. Superhuman ai for heads-up no-limit poker: Libratus beats top
496 professionals. *Science*, 359(6374):418–424, 2018.

497

498 Noam Brown, Adam Lerer, Sam Gross, and Tuomas Sandholm. Deep counterfactual regret mini-
499 mization. In *Proceedings of the 36th International Conference on Machine Learning*, volume 97
500 of *Proceedings of Machine Learning Research*, pp. 793–802. PMLR, 2019.

501 Jiayu Chen, Vaneet Aggarwal, and Tian Lan. ODPP: A unified algorithm framework for unsupervised
502 option discovery based on determinantal point process. *CoRR*, abs/2212.00211, 2022a.

503

504 Jiayu Chen, Jingdi Chen, Tian Lan, and Vaneet Aggarwal. Learning multi-agent options for tabular
505 reinforcement learning using factor graphs. *IEEE Transactions on Artificial Intelligence*, pp. 1–13,
506 2022b. doi: 10.1109/tai.2022.3195818.

507 Jiayu Chen, Jingdi Chen, Tian Lan, and Vaneet Aggarwal. Multi-agent covering option discovery
508 based on kronecker product of factor graphs. *IEEE Transactions on Artificial Intelligence*, 2022c.

509

510 Jiayu Chen, Tian Lan, and Vaneet Aggarwal. Option-aware adversarial inverse reinforcement learning
511 for robotic control. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*,
512 pp. 5902–5908. IEEE, 2023a.

513 Jiayu Chen, Dipesh Tamboli, Tian Lan, and Vaneet Aggarwal. Multi-task hierarchical adversarial
514 inverse reinforcement learning. In *Proceedings of the 40th International Conference on Machine
515 Learning*, volume 202 of *Proceedings of Machine Learning Research*, pp. 4895–4920. PMLR,
516 2023b.

517

518 Trevor Davis, Martin Schmid, and Michael Bowling. Low-variance and zero-variance baselines for
519 extensive-form games. In *Proceedings of the 37th International Conference on Machine Learning*,
520 volume 119 of *Proceedings of Machine Learning Research*, pp. 2392–2401. PMLR, 2020.

521 Benjamin Eysenbach, Abhishek Gupta, Julian Ibarz, and Sergey Levine. Diversity is all you need:
522 Learning skills without a reward function. In *Proceedings of the 7th International Conference on
523 Learning Representations*. OpenReview.net, 2019.

524

525 Sam Ganzfried and Tuomas Sandholm. Potential-aware imperfect-recall abstraction with earth
526 mover’s distance in imperfect-information games. In *Proceedings of 28th the AAAI Conference on
527 Artificial Intelligence*, volume 28, 2014.

528 Richard G. Gibson, Marc Lanctot, Neil Burch, Duane Szafron, and Michael Bowling. Generalized
529 sampling and variance in counterfactual regret minimization. In *Proceedings of the 26th AAAI
530 Conference on Artificial Intelligence*. AAAI Press, 2012.

531

532 Johannes Heinrich and David Silver. Deep reinforcement learning from self-play in imperfect-
533 information games. *CoRR*, abs/1603.01121, 2016.

534

535 Mingxuan Jing, Wenbing Huang, Fuchun Sun, Xiaojian Ma, Tao Kong, Chuang Gan, and Lei Li.
536 Adversarial option-aware hierarchical imitation learning. In *Proceedings of the 38th International
537 Conference on Machine Learning*, pp. 5097–5106. PMLR, 2021a.

538

539 Mingxuan Jing, Wenbing Huang, Fuchun Sun, Xiaojian Ma, Tao Kong, Chuang Gan, and Lei Li.
Adversarial option-aware hierarchical imitation learning. In *Proceedings of the 38th International
Conference on Machine Learning*, pp. 5097–5106. PMLR, 2021b.

540 Yuu Jinnai, Jee Won Park, Marlos C. Machado, and George Dimitri Konidaris. Exploration in
541 reinforcement learning with deep covering options. In *Proceedings of the 8th International*
542 *Conference on Learning Representations*. OpenReview.net, 2020.

543

544 Vishruti Kakkad, Hitarth Shah, Reema Patel, and Nishant Doshi. A comparative study of applications
545 of game theory in cyber security and cloud computing. *Procedia Computer Science*, 155:680–685,
546 2019.

547 Marc Lanctot, Kevin Waugh, Martin Zinkevich, and Michael H. Bowling. Monte carlo sampling for
548 regret minimization in extensive games. In *Proceedings of the 23rd Annual Conference on Neural*
549 *Information Processing Systems*, pp. 1078–1086. Curran Associates, Inc., 2009.

550

551 Marc Lanctot, Vinícius Flores Zambaldi, Audrunas Gruslys, Angeliki Lazaridou, Karl Tuyls, Julien
552 Pérolat, David Silver, and Thore Graepel. A unified game-theoretic approach to multiagent
553 reinforcement learning. In *Advances in Neural Information Processing Systems 30*, pp. 4190–4203,
554 2017.

555 Chang Li, Dongjin Song, and Dacheng Tao. The skill-action architecture: Learning abstract action
556 embeddings for reinforcement learning. In *Submissions of the 9th International Conference on*
557 *Learning Representations*, 2021a.

558

559 Huale Li, Xuan Wang, Zengyue Guo, Jiajia Zhang, and Shuhan Qi. D2cfr: Minimize counterfactual
560 regret with deep dueling neural network. *arXiv preprint arXiv:2105.12328*, 2021b.

561 Hui Li, Kailiang Hu, Shaohua Zhang, Yuan Qi, and Le Song. Double neural counterfactual regret
562 minimization. In *Proceedings of the 8th International Conference on Learning Representations*.
563 OpenReview.net, 2020.

564

565 Matej Moravčík, Martin Schmid, Neil Burch, Viliam Lisý, Dustin Morrill, Nolan Bard, Trevor
566 Davis, Kevin Waugh, Michael Johanson, and Michael Bowling. Deepstack: Expert-level artificial
567 intelligence in heads-up no-limit poker. *Science*, 356(6337):508–513, 2017.

568 Philipp Moritz, Robert Nishihara, Stephanie Wang, Alexey Tumanov, Richard Liaw, Eric Liang,
569 Melih Elibol, Zongheng Yang, William Paul, Michael I. Jordan, and Ion Stoica. Ray: A distributed
570 framework for emerging AI applications. In *13th USENIX Symposium on Operating Systems*
571 *Design and Implementation*, pp. 561–577. USENIX Association, 2018.

572

573 Thomas H. Noe, Michael Rebelló, and Jun Wang. Learning to bid: The design of auctions under
574 uncertainty and adaptation. *Games and Economic Behavior*, pp. 620–636, 2012.

575

576 Martin J Osborne and Ariel Rubinstein. *A course in game theory*. MIT press, 1994.

577

578 Martin Schmid, Neil Burch, Marc Lanctot, Matej Moravcik, Rudolf Kadlec, and Michael Bowling.
579 Variance reduction in monte carlo counterfactual regret minimization (VR-MCCFR) for extensive
580 form games using baselines. In *Proceedings of 33rd the AAAI Conference on Artificial Intelligence*,
581 pp. 2157–2164. AAAI Press, 2019.

582

583 Finnegan Southey, Michael H. Bowling, Bryce Larson, Carmelo Piccione, Neil Burch, Darse Billings,
584 and D. Chris Rayner. Bayes’ bluff: Opponent modelling in poker. In *Proceedings of the 21st*
Conference in Uncertainty in Artificial Intelligence, pp. 550–558, 2005.

585

586 Sriram Srinivasan, Marc Lanctot, Vinícius Flores Zambaldi, Julien Pérolat, Karl Tuyls, Rémi Munos,
587 and Michael Bowling. Actor-critic policy optimization in partially observable multiagent environ-
588 ments. In *Proceedings of the 32nd Annual Conference on Neural Information Processing Systems*,
pp. 3426–3439, 2018.

589

590 Eric Steinberger. Single deep counterfactual regret minimization. *CoRR*, abs/1901.07621, 2019.

591

592 Eric Steinberger. Pokerrl. <https://github.com/EricSteinberger/DREAM>, 2020.

593

Eric Steinberger, Adam Lerer, and Noam Brown. DREAM: deep regret minimization with advantage
baselines and model-free learning. *CoRR*, abs/2006.10410, 2020.

594 Richard S. Sutton, Doina Precup, and Satinder Singh. Between mdps and semi-mdps: A framework
595 for temporal abstraction in reinforcement learning. *Artificial Intelligence*, 112(1):181–211, 1999.
596 ISSN 0004-3702. doi: [https://doi.org/10.1016/S0004-3702\(99\)00052-1](https://doi.org/10.1016/S0004-3702(99)00052-1).
597

598 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez,
599 Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Proceedings of the 31st Annual
600 Conference on Neural Information Processing Systems*, pp. 5998–6008, 2017.

601 Kevin Waugh, Dustin Morrill, James Andrew Bagnell, and Michael H. Bowling. Solving games with
602 functional regret estimation. In *Proceedings of the 29th AAAI Conference on Artificial Intelligence*,
603 pp. 2138–2145. AAAI Press, 2015.

604 Fuxiang Zhang, Chengxing Jia, Yi-Chen Li, Lei Yuan, Yang Yu, and Zongzhang Zhang. Discovering
605 generalizable multi-agent coordination skills from multi-task offline data. In *Proceedings of the
606 11th International Conference on Learning Representations*, 2022.
607

608 Shangdong Zhang and Shimon Whiteson. DAC: the double actor-critic architecture for learning
609 options. In *Proceedings of the 33rd Annual Conference on Neural Information Processing Systems*,
610 pp. 2010–2020, 2019.

611 Martin Zinkevich, Michael Johanson, Michael H. Bowling, and Carmelo Piccione. Regret minimiza-
612 tion in games with incomplete information. In *Proceedings of the 21st Annual Conference on
613 Neural Information Processing Systems*, pp. 1729–1736. Curran Associates, Inc., 2007.
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647

A RELATED WORKS

Counterfactual Regret Minimization (CFR) Zinkevich et al. (2007) is an algorithm for learning Nash Equilibria in extensive-form games through iterative self-play. As part of this process, it must traverse the entire game tree during each learning iteration, which is prohibitive for large-scale games. This motivates the development of Monte Carlo CFR (MCCFR) Lanctot et al. (2009), which samples trajectories traversing part of the tree to allow for significantly faster iterations. However, the variance of Monte Carlo sampling can become a significant issue, particularly for long sample trajectories. The authors of Schmid et al. (2019); Davis et al. (2020) then propose to introduce baseline functions for variance reduction. Notably, all methods mentioned above are tabular-based. For games with large state space, domain-specific abstraction schemes Ganzfried & Sandholm (2014); Moravčík et al. (2017) are required to shrink them to a manageable size by clustering states into buckets, which necessitates expert knowledge and is not applicable to all games.

To obviate the need of abstractions, several CFR variants utilizing function approximators have been proposed. Pioneering this was Regression CFR Waugh et al. (2015), which adopts regression trees to model cumulative regrets but relies on hand-crafted features and full traversals of the game tree. Subsequently, several works Brown et al. (2019); Li et al. (2020); Steinberger (2019); Li et al. (2021b) propose to model the cumulative counterfactual regrets and average strategies in MCCFR as neural networks to enhance the scalability. **However, all these methods rely on knowledge of the game model to realize backtracking (i.e., sampling multiple actions at an information set) for regret estimation.** As a model-free approach, Neural Fictitious Self-Play (NFSP) Heinrich & Silver (2016) is the first deep reinforcement learning algorithm capable of learning a Nash Equilibrium in two-player imperfect information games through self-play. Since its introduction, various policy gradient and actor-critic methods have demonstrated similar convergence properties when appropriately tuned Lanctot et al. (2017); Srinivasan et al. (2018). However, fictitious play empirically converges slower than CFR-based approaches in many settings. DREAM Steinberger et al. (2020) extends Deep CFR with variance-reduction techniques from Davis et al. (2020) and represents the state-of-the-art in model-free algorithms of this area. **Compared with DREAM, our algorithm provides solid theoretical justification, enables hierarchical learning with (prelearned) skills, and empirically shows enhanced performance on longer-horizon games.**

As another important module of HDCFR, the option framework Sutton et al. (1999) enables learning and planning at multiple temporal levels and has been widely adopted in reinforcement learning (RL). Multiple research areas centered on this framework have been developed. Unsupervised Option Discovery focuses on identifying skills that are diverse and effective for (various) downstream tasks without relying on task-specific reward signals. Algorithms have been developed for both **single-agent settings** Eysenbach et al. (2019); Jinnai et al. (2020); Chen et al. (2022a) and **collaborative multi-agent scenarios** Chen et al. (2022c;b); Zhang et al. (2022). Hierarchical Reinforcement Learning Zhang & Whiteson (2019); Li et al. (2021a) and Hierarchical Imitation Learning Jing et al. (2021b); Chen et al. (2023a;b), on the other hand, aim at directly learning a hierarchical policy that incorporates skills, either through interactions with the environment or from expert demonstrations.

Our proposed algorithm – HDCFR, is a pioneering effort to amalgamate options with CFR and demonstrates the superiority of hierarchical learning in **competitive multi-agent scenarios** (more specifically, zero-sum games).

B PROOF OF THEOREM 2

Define $D(I)$ to be the information sets of player i reachable from I (including I), and $\sigma|_{D(I) \rightarrow \sigma'_i}$ to be a strategy profile equal to σ except that player i adopts σ'_i in the information sets contained in $D(I)$. Then, the average overall regret starting from I ($I \in \mathcal{I}_i$) can be defined as:

$$R_{full,i}^T(I) = \frac{1}{T} \max_{\sigma'_i} \sum_{t=1}^T \pi_{-i}^{\sigma^t}(I) (u_i(\sigma^t|_{D(I) \rightarrow \sigma'_i}, I) - u_i(\sigma^t, I)) \quad (13)$$

Further, we define $S_i(I, \tilde{a})$ to be the set of all possible next information sets of player i given that action $\tilde{a} \in \tilde{A}(I)$ was just selected at I and define $S_i(I) = \bigcup_{\tilde{a} \in \tilde{A}(I)} S_i(I, \tilde{a})$, $S_i(Iz) = \bigcup_{a \in A(I)} S_i(I, za)$. Then, we have the following lemma:

Lemma 1. $R_{full,i}^{T,+}(I) \leq R_{i,+}^{T,H}(I) + \sum_{z \in Z(I)} R_{i,+}^{T,L}(I, z) + \sum_{I' \in S_i(I)} R_{full,i}^{T,+}(I')$

Proof.

$$\begin{aligned}
R_{full,i}^T(I) &= \frac{1}{T} \max_{\sigma'_i} \sum_{t=1}^T \pi_{-i}^{\sigma^t}(I) (u_i(\sigma^t|_{D(I) \rightarrow \sigma'_i}, I) - u_i(\sigma^t, I)) \\
&= \frac{1}{T} \max_{z \in Z(I)} \max_{\sigma'_i} \sum_{t=1}^T \pi_{-i}^{\sigma^t}(I) \left[(u_i(\sigma^t|_{I \rightarrow z}, I) - u_i(\sigma^t, I)) + (u_i(\sigma^t|_{D(Iz) \rightarrow \sigma'_i}, Iz) - u_i(\sigma^t, Iz)) \right] \\
&\leq \frac{1}{T} \max_{z \in Z(I)} \sum_{t=1}^T \pi_{-i}^{\sigma^t}(I) (u_i(\sigma^t|_{I \rightarrow z}, I) - u_i(\sigma^t, I)) + \frac{1}{T} \max_{z \in Z(I)} \max_{\sigma'_i} \sum_{t=1}^T \pi_{-i}^{\sigma^t}(I) (u_i(\sigma^t|_{D(Iz) \rightarrow \sigma'_i}, Iz) - u_i(\sigma^t, Iz)) \\
&= R_i^{T,H}(I) + \frac{1}{T} \max_{z \in Z(I)} \max_{\sigma'_i} \sum_{t=1}^T \pi_{-i}^{\sigma^t}(Iz) (u_i(\sigma^t|_{D(Iz) \rightarrow \sigma'_i}, Iz) - u_i(\sigma^t, Iz)) \\
&\leq R_i^{T,H}(I) + \sum_{z \in Z(I)} \left[\frac{1}{T} \max_{\sigma'_i} \sum_{t=1}^T \pi_{-i}^{\sigma^t}(Iz) (u_i(\sigma^t|_{D(Iz) \rightarrow \sigma'_i}, Iz) - u_i(\sigma^t, Iz)) \right]^+ \\
&= R_i^{T,H}(I) + \sum_{z \in Z(I)} R_{full,i}^{T,+}(Iz)
\end{aligned} \tag{14}$$

$$\begin{aligned}
R_{full,i}^T(Iz) &= \frac{1}{T} \max_{\sigma'_i} \sum_{t=1}^T \pi_{-i}^{\sigma^t}(Iz) (u_i(\sigma^t|_{D(Iz) \rightarrow \sigma'_i}, Iz) - u_i(\sigma^t, Iz)) \\
&= \frac{1}{T} \max_{a \in A(I)} \max_{\sigma'_i} \sum_{t=1}^T \pi_{-i}^{\sigma^t}(Iz) [(u_i(\sigma^t|_{Iz \rightarrow a}, Iz) - u_i(\sigma^t, Iz)) + \\
&\quad \sum_{I' \in S_i(I, za)} P_{\sigma^t_{-i}}(I'|I, za) (u_i(\sigma^t|_{D(I') \rightarrow \sigma'_i}, I') - u_i(\sigma^t, I'))] \\
&= R_i^{T,L}(I, z) + \max_{a \in A(I)} \sum_{I' \in S_i(I, za)} \frac{1}{T} \max_{\sigma'_i} \sum_{t=1}^T \pi_{-i}^{\sigma^t}(I') (u_i(\sigma^t|_{D(I') \rightarrow \sigma'_i}, I') - u_i(\sigma^t, I')) \\
&= R_i^{T,L}(I, z) + \max_{a \in A(I)} \sum_{I' \in S_i(I, za)} R_{full,i}^T(I') \\
&\leq R_i^{T,L}(I, z) + \sum_{I' \in S_i(Iz)} R_{full,i}^{T,+}(I')
\end{aligned} \tag{15}$$

In Equation (14) and (15), we employ the one-step look-ahead expansion (Equation (10) in Zinkevich et al. (2007)) for the second line. At iteration t , when player i selects a hierarchical action $\tilde{a} = (za)$, it will transit to the subsequent information set $I' \in S_i(I, za)$ with a probability of $P_{\sigma^t_{-i}}(I'|I, za)$, since only player $-i$ will act between I and I' according to σ^t_{-i} . According to the definition of the reach probability, $\pi_{-i}(Iz) = \pi_{-i}(I)$ (since z and a are executed by player i) and $\pi_{-i}(I)P_{\sigma^t_{-i}}(I'|I, za) = \pi_{-i}(I')$. Combining Equation (14) and (15), we can get:

$$\begin{aligned}
R_{full,i}^T(I) &\leq R_i^{T,H}(I) + \sum_{z \in Z(I)} R_{full,i}^{T,+}(Iz) \\
&\leq R_i^{T,H}(I) + \sum_{z \in Z(I)} \left[R_{i,+}^{T,L}(I, z) + \sum_{I' \in S_i(Iz)} R_{full,i}^{T,+}(I') \right] \\
&= R_i^{T,H}(I) + \sum_{z \in Z(I)} R_{i,+}^{T,L}(I, z) + \sum_{z \in Z(I)} \sum_{I' \in S_i(Iz)} R_{full,i}^{T,+}(I') \\
&= R_i^{T,H}(I) + \sum_{z \in Z(I)} R_{i,+}^{T,L}(I, z) + \sum_{I' \in S_i(I)} R_{full,i}^{T,+}(I')
\end{aligned} \tag{16}$$

In previous derivations, we have repeatedly employed the inequality $\max(a + b, 0) \leq \max(a, 0) + \max(b, 0)$, which holds for all $a, b \in \mathbb{R}$, as in the last inequality of Equation (14) and (15). By applying this inequality once more to Equation (16), we can obtain Lemma 1. \square

Lemma 2. $R_{full,i}^{T,+}(I) \leq \sum_{I' \in D(I)} \left[R_{i,+}^{T,H}(I') + \sum_{z \in Z(I')} R_{i,+}^{T,L}(I', z) \right]$

Proof. We prove this lemma by induction on the height of the information set I on the game tree. When the height is 1, i.e., $S_i(I) = \emptyset$, $D(I) = \{I\}$, then Lemma 1 implies Lemma 2. Now, for the general case:

$$\begin{aligned}
R_{full,i}^{T,+}(I) &\leq R_{i,+}^{T,H}(I) + \sum_{z \in Z(I)} R_{i,+}^{T,L}(I, z) + \sum_{I' \in S_i(I)} R_{full,i}^{T,+}(I') \\
&\leq R_{i,+}^{T,H}(I) + \sum_{z \in Z(I)} R_{i,+}^{T,L}(I, z) + \sum_{I' \in S_i(I)} \sum_{I'' \in D(I')} \left[R_{i,+}^{T,H}(I'') + \sum_{z \in Z(I'')} R_{i,+}^{T,L}(I'', z) \right] \quad (17) \\
&= \sum_{I' \in D(I)} \left[R_{i,+}^{T,H}(I') + \sum_{z \in Z(I')} R_{i,+}^{T,L}(I', z) \right]
\end{aligned}$$

In the second line, we employ the induction hypothesis. In the third line, we use the following facts: $D(I) = \{I\} \cup \bigcup_{I' \in S_i(I)} D(I')$, $\{I\} \cap \bigcup_{I' \in S_i(I)} D(I') = \emptyset$, and $D(I') \cap D(I'') = \emptyset$ for all distinct $I', I'' \in S_i(I)$. The third fact here is derived from the perfect recall property of the game: all players can recall their previous (hierarchical) actions and the corresponding information sets. Then, $D(I') \cap D(I'') = \emptyset$ because elements from the two sets possess distinct prefixes (i.e., I' and I''). \square

Last, for the average overall regret, we have $R_{full,i}^T = R_{full,i}^T(\emptyset)$, where \emptyset corresponds to the start of the game tree and $D(\emptyset) = \mathcal{I}_i$. Applying Lemma 2, we can get the theorem: $R_{full,i}^T \leq R_{full,i}^{T,+}(\emptyset) \leq \sum_{I \in \mathcal{I}_i} \left[R_{i,+}^{T,H}(I) + \sum_{z \in Z(I)} R_{i,+}^{T,L}(I, z) \right]$.

C PROOF OF THEOREM 3

Regret matching can be defined in a domain where a fixed set of actions A and a payoff function $u^t : A \rightarrow \mathbb{R}$ exist. At each iteration t , a distribution over the actions, σ^t , is chosen based on the cumulative regret $R^t : A \rightarrow \mathbb{R}$. Specifically, the cumulative regret at iteration T for not playing action a is defined as:

$$R^T(a) = \frac{1}{T} \sum_{t=1}^T \left[u^t(a) - \sum_{a' \in A} \sigma^t(a') u^t(a') \right] \quad (18)$$

where $\sigma^t(a)$ is obtained by:

$$\sigma^t(a) = \begin{cases} R^{t-1,+}(a)/\mu, & \mu > 0, \\ 1/|A|, & \mu \leq 0. \end{cases} \quad \mu = \sum_{a' \in A} R^{t-1,+}(a') \quad (19)$$

Then, we have the following lemma (Theorem 8 in Zinkevich et al. (2007)):

Lemma 3. $\max_{a \in A} R^T(a) \leq \frac{\Delta_u \sqrt{|A|}}{\sqrt{T}}$, where $\Delta_u = \max_{t \in \{1, \dots, T\}} \max_{a, a' \in A} (u^t(a) - u^t(a'))$.

To apply this lemma, we must transform the definitions of $R_i^{T,H}$ and $R_i^{T,L}$ in Equation (6) to a form resembling Equation (18). With Equation (6) and (2), we can get:

$$\begin{aligned}
R_i^{T,H}(z|I) &= \frac{1}{T} \sum_{t=1}^T \left[\sum_{h \in I} \pi_{-i}^{\sigma^t}(h) \sum_{h' \in H_{TS}} u_i(h') \pi^{\sigma^t}(hz, h') - \right. \\
&\quad \left. \sum_{h \in I} \pi_{-i}^{\sigma^t}(h) \sum_{h' \in H_{TS}} u_i(h') \sum_{z' \in Z(h)} \sigma_t^H(z'|h) \pi^{\sigma^t}(hz', h') \right] \\
&= \frac{1}{T} \sum_{t=1}^T \left[\sum_{h \in I} \pi_{-i}^{\sigma^t}(h) \sum_{h' \in H_{TS}} u_i(h') \pi^{\sigma^t}(hz, h') - \right. \\
&\quad \left. \sum_{z' \in Z(I)} \sigma_t^H(z'|I) \sum_{h \in I} \pi_{-i}^{\sigma^t}(h) \sum_{h' \in H_{TS}} u_i(h') \pi^{\sigma^t}(hz', h') \right] \\
&= \frac{1}{T} \sum_{t=1}^T \left[v_t^H(z) - \sum_{z' \in Z(I)} \sigma_t^H(z'|I) v_t^H(z') \right]
\end{aligned} \tag{20}$$

Applying the same process on $R_i^{T,L}(a|I, z)$, we can get:

$$\begin{aligned}
R_i^{T,L}(a|I, z) &= \frac{1}{T} \sum_{t=1}^T \left[v_t^L(a) - \sum_{a' \in A(I)} \sigma_t^L(a'|I, z) v_t^L(a') \right] \\
v_t^L(a) &= \sum_{h \in I} \pi_{-i}^{\sigma^t}(h) \sum_{h' \in H_{TS}} u_i(h') \pi^{\sigma^t}(hza, h')
\end{aligned} \tag{21}$$

Then, we can apply Lemma 3 and obtain:

$$\begin{aligned}
\max_{z \in Z(I)} R_i^{T,H}(z|I) &= R_i^{T,H}(I) \leq \frac{\Delta_{v^H} \sqrt{|Z(I)|}}{\sqrt{T}} \leq \frac{\Delta_{u,i} \sqrt{|Z(I)|}}{\sqrt{T}} \\
\max_{a \in A(I)} R_i^{T,L}(a|I, z) &= R_i^{T,L}(I, z) \leq \frac{\Delta_{v^L} \sqrt{|A(I)|}}{\sqrt{T}} \leq \frac{\Delta_{u,i} \sqrt{|A(I)|}}{\sqrt{T}}
\end{aligned} \tag{22}$$

Here, $\Delta_{u,i} = \max_{h' \in H_{TS}} u_i(h') - \min_{h' \in H_{TS}} u_i(h')$ is the range of the payoff function for i , which covers Δ_{v^H} and Δ_{v^L} . We can directly apply Lemma 3, because the regret matching is adopted at each information set independently as defined in Equation (7). By integrating Equation 22 and Theorem 2, we then get:

$$\begin{aligned}
R_{full,i}^T &\leq \sum_{I \in \mathcal{I}_i} \left[\frac{\Delta_{u,i} \sqrt{|Z(I)|}}{\sqrt{T}} + \sum_{z \in Z(I)} \frac{\Delta_{u,i} \sqrt{|A(I)|}}{\sqrt{T}} \right] \\
&\leq \frac{\Delta_{u,i} |\mathcal{I}_i|}{\sqrt{T}} (\sqrt{|Z_i|} + |Z_i| \sqrt{|A_i|})
\end{aligned} \tag{23}$$

where $|\mathcal{I}_i|$ is the number of information sets for player i , $|A_i| = \max_{h:P(h)=i} |A(h)|$, $|Z_i| = \max_{h:P(h)=i} |Z(h)|$.

D PROOF OF PROPOSITION 1

According to Theorem 1 and 3, as $T \rightarrow \infty$, $R_{full,i}^T \rightarrow 0$, and thus the average strategy $\bar{\sigma}_i^T(\tilde{a}|I)$ converges to a Nash Equilibrium. We claim that $\bar{\sigma}_i^T(\tilde{a}|I) = \bar{\sigma}_i^{T,H}(z|I) \cdot \bar{\sigma}_i^{T,L}(a|I, z)$.

864 *Proof.*

$$\begin{aligned}
865 \quad \bar{\sigma}_i^{T,H}(z|I) \cdot \bar{\sigma}_i^{T,L}(a|I, z) &= \frac{\sum_{t=1}^T \pi_i^{\sigma^t}(I) \sigma_i^{t,H}(z|I)}{\sum_{t=1}^T \pi_i^{\sigma^t}(I)} \frac{\sum_{t=1}^T \pi_i^{\sigma^t}(Iz) \sigma_i^{t,L}(a|I, z)}{\sum_{t=1}^T \pi_i^{\sigma^t}(Iz)} \\
866 &= \frac{\sum_{t=1}^T \pi_i^{\sigma^t}(I) \sigma_i^{t,H}(z|I)}{\sum_{t=1}^T \pi_i^{\sigma^t}(I)} \frac{\sum_{t=1}^T \pi_i^{\sigma^t}(Iz) \sigma_i^{t,L}(a|I, z)}{\sum_{t=1}^T \pi_i^{\sigma^t}(I) \sigma_i^{t,H}(z|I)} \\
867 &= \frac{\sum_{t=1}^T \pi_i^{\sigma^t}(Iz) \sigma_i^{t,L}(a|I, z)}{\sum_{t=1}^T \pi_i^{\sigma^t}(I)} = \frac{\sum_{t=1}^T \pi_i^{\sigma^t}(I) \sigma_i^{t,H}(z|I) \sigma_i^{t,L}(a|I, z)}{\sum_{t=1}^T \pi_i^{\sigma^t}(I)} \\
868 &= \frac{\sum_{t=1}^T \pi_i^{\sigma^t}(I) \sigma_i^t((z, a)|I)}{\sum_{t=1}^T \pi_i^{\sigma^t}(I)} = \bar{\sigma}_i^T(\tilde{a}|I)
\end{aligned} \tag{24}$$

876 \square

877 Given this equivalence, we can infer that if both players adhere to the one-step option model for each I
878 – selecting an option z based on $\bar{\sigma}_i^{T,H}(\cdot|I)$ and subsequently choosing the action a in accordance with
879 the corresponding intra-option strategy $\bar{\sigma}_i^{T,L}(\cdot|I, z)$, an approximate NE solution can be acquired.

882 E LOW-VARIANCE MONTE CARLO SAMPLING EXTENSION

883 MCCFR’s main insight is substituting the counterfactual regrets R_i^T in CFR with their unbiased
884 estimations, while maintaining the other learning rules (as in Section 3.2). This approach allows
885 for updating functions only at information sets within the sample trajectories, bypassing the need to
886 traverse the full game tree. With MCCFR, the average overall regret $R_{full,i}^T \rightarrow 0$ as $T \rightarrow \infty$ at the
887 same convergence rate as vanilla CFR, with high probability, as stated in Theorem 5 of Lanctot et al.
888 (2009). Therefore, to apply the Monte Carlo sampling extension, we propose unbiased estimations of
889 $R_i^{T,H}(z|I)$ and $R_i^{T,L}(a|I, z)$, $\forall i \in N, I \in \mathcal{I}_i, z \in Z(I), a \in A(I)$.

891 First, we define $R_i^{T,H}(z|I)$ and $R_i^{T,L}(a|I, z)$ with the immediate counterfactual regrets r_i^t and values
892 v_i^t : ($v_i^{t,H}(\sigma^t, h) = u_i(h)$, $\forall h \in H_{TS}$)

$$\begin{aligned}
894 \quad R_i^{T,H}(z|I) &= \frac{1}{T} \sum_{t=1}^T r_i^{t,H}(I, z), \quad r_i^{t,H}(I, z) = \sum_{h \in I} \pi_{-i}^{\sigma^t}(h) \left[v_i^{t,L}(\sigma^t, hz) - v_i^{t,H}(\sigma^t, h) \right] \\
895 & \\
896 & \\
897 \quad R_i^{T,L}(a|I, z) &= \frac{1}{T} \sum_{t=1}^T r_i^{t,L}(Iz, a), \quad r_i^{t,L}(Iz, a) = \sum_{h \in I} \pi_{-i}^{\sigma^t}(h) \left[v_i^{t,H}(\sigma^t, hza) - v_i^{t,L}(\sigma^t, hz) \right] \\
898 & \\
899 & \\
900 \quad v_i^{t,H}(\sigma^t, h) &= \sum_{z \in Z(h)} \sigma_{P(h)}^{t,H}(z|h) v_i^{t,L}(\sigma^t, hz), \quad v_i^{t,L}(\sigma^t, hz) = \sum_{a \in A(h)} \sigma_{P(h)}^{t,L}(a|h, z) v_i^{t,H}(\sigma^t, hza) \\
901 & \\
902 & \\
903 & \tag{25}
\end{aligned}$$

903 The equivalence between Equation (25) and (6) is proved in Appendix F.

904 Next, we propose to collect trajectories $h' \in H_{TS}$ with the sample strategy q^t at each iteration t , and
905 compute the corresponding **sampld** immediate counterfactual regrets \hat{r}_i^t and values \hat{v}_i^t as follows:

$$\begin{aligned}
906 \quad \hat{r}_i^{t,H}(I, z|h') &= \sum_{h \in I} \frac{\pi_{-i}^{\sigma^t}(h)}{\pi^{q^t}(h)} \left[\hat{v}_i^{t,H}(\sigma^t, h, z|h') - \hat{v}_i^{t,H}(\sigma^t, h|h') \right] \\
907 & \\
908 & \\
909 & \tag{26} \\
910 \quad \hat{r}_i^{t,L}(Iz, a|h') &= \sum_{h \in I} \frac{\pi_{-i}^{\sigma^t}(h)}{\pi^{q^t}(hz)} \left[\hat{v}_i^{t,L}(\sigma^t, hz, a|h') - \hat{v}_i^{t,L}(\sigma^t, hz|h') \right] \\
911 & \\
912 &
\end{aligned}$$

912 Here, inspired by Davis et al. (2020), $\hat{v}_i^{t,H}(\sigma^t, h, z|h')$ and $\hat{v}_i^{t,L}(\sigma^t, hz, a|h')$ are incorporated with
913 the baseline function b_i^t for variance reduction: ($\hat{v}_i^{t,H}(\sigma^t, h'|h') = u_i(h')$)

$$\begin{aligned}
914 \quad \hat{v}_i^{t,H}(\sigma^t, h, z|h') &= \frac{\delta(hz \sqsubseteq h')}{q^t(z|h')} \left[\hat{v}_i^{t,L}(\sigma^t, hz|h') - b_i^t(h, z) \right] + b_i^t(h, z) \\
915 & \\
916 & \\
917 \quad \hat{v}_i^{t,L}(\sigma^t, hz, a|h') &= \frac{\delta(hza \sqsubseteq h')}{q^t(a|h, z)} \left[\hat{v}_i^{t,H}(\sigma^t, hza|h') - b_i^t(h, z, a) \right] + b_i^t(h, z, a)
\end{aligned} \tag{27}$$

where $\delta(\cdot)$ is the indicator function. Accordingly, $\hat{v}_i^{t,H}(\sigma^t, h|h') = \sum_{z \in Z(h)} \sigma_{P(h)}^{t,H}(z|h) \hat{v}_i^{t,H}(\sigma^t, h, z|h')$ and $\hat{v}_i^{t,L}(\sigma^t, hz|h') = \sum_{a \in A(h)} \sigma_{P(h)}^{t,L}(a|h, z) \hat{v}_i^{t,L}(\sigma^t, hz, a|h')$. For superscripts on \hat{r} and \hat{v} , we use H when the agent is in state h or hza , corresponding to high-level option choices, and L when the agent is in state hz , corresponding to low-level action decisions.

Regarding estimators proposed in Eqs (26) and (27), we have the following theorems:

Theorem 4. For all $i \in N$, $I \in \mathcal{I}_i$, $z \in Z(I)$, $a \in A(I)$, we have:

$$\mathbb{E}_{h' \sim \pi^{qt}(\cdot)} \left[\hat{r}_i^{t,H}(I, z|h') \right] = r_i^{t,H}(I, z), \quad \mathbb{E}_{h' \sim \pi^{qt}(\cdot)} \left[\hat{r}_i^{t,L}(Iz, a|h') \right] = r_i^{t,L}(Iz, a) \quad (28)$$

Therefore, we can acquire unbiased estimations of R_i^T by substituting r_i^t with \hat{r}_i^t in Equation (25). This theorem is proved in Appendix G. Notably, Theorem 4 doesn't prescribe any specific form for the baseline function b_i^t . Yet, the baseline design can affect the sample variance of these unbiased estimators. As posited in Gibson et al. (2012), given a fixed $\epsilon > 0$, estimators with reduced variance necessitate fewer iterations to converge to an ϵ -Nash equilibrium. Hence, we propose the following ideal criteria for the baseline function to minimize the sample variance:

Theorem 5. If $b_i^t(h, z, a) = v_i^{t,H}(\sigma^t, hza)$ and $b_i^t(h, z) = v_i^{t,L}(\sigma^t, hz)$, for all $h \in H \setminus H_{TS}$, $z \in Z(h)$, $a \in A(h)$, we have:

$$\text{Var}_{h'} \left[\hat{v}_i^{t,H}(\sigma^t, h, z|h') | h' \supseteq h \right] = \text{Var}_{h'} \left[\hat{v}_i^{t,L}(\sigma^t, hz, a|h') | h' \supseteq hz \right] = 0 \quad (29)$$

Consequently, $\text{Var}_{h' \sim \pi^{qt}(\cdot)} \left[\hat{r}_i^{t,H}(I, z|h') \right]$ and $\text{Var}_{h' \sim \pi^{qt}(\cdot)} \left[\hat{r}_i^{t,L}(Iz, a|h') \right]$ are minimized with respect to b_i^t for all $I \in \mathcal{I}_i$, $z \in Z(I)$, $a \in A(I)$.

The proof can be found in Appendix H. **The ideal criteria for the baseline function proposed in Theorem 5 is incorporated into our objective design in Section 3.3.**

To sum up, by utilizing the immediate counterfactual regret estimators defined in Equations (26) and (27) and selecting baseline functions that meet the ideal criteria, we can enhance the adaptability and learning efficiency of our method (i.e., HCFR proposed in Section 3.2) through a low-variance outcome Monte Carlo sampling extension.

F PROOF OF EQUIVALENCE BETWEEN EQUATIONS (25) AND (6)

Through induction on the height of h on the game tree, one can easily prove that:

$$v_i^{t,H}(\sigma^t, h) = \sum_{h' \in H_{TS}} \pi^{\sigma^t}(h, h') u_i(h'), \quad v_i^{t,L}(\sigma^t, hz) = \sum_{h' \in H_{TS}} \pi^{\sigma^t}(hz, h') u_i(h') \quad (30)$$

Thus, we have:

$$\begin{aligned} r_i^{t,H}(I, z) &= \sum_{h \in I} \pi_{-i}^{\sigma^t}(h) \sum_{h' \in H_{TS}} \pi^{\sigma^t}(hz, h') u_i(h') - \sum_{h \in I} \pi_{-i}^{\sigma^t}(h) \sum_{h' \in H_{TS}} \pi^{\sigma^t}(h, h') u_i(h') \\ &= \pi_{-i}^{\sigma^t}(I) (u_i(\sigma^t |_{I \rightarrow z}, I) - u_i(\sigma^t, I)) \\ r_i^{t,L}(Iz, a) &= \sum_{h \in I} \pi_{-i}^{\sigma^t}(h) \sum_{h' \in H_{TS}} \pi^{\sigma^t}(hza, h') u_i(h') - \sum_{h \in I} \pi_{-i}^{\sigma^t}(h) \sum_{h' \in H_{TS}} \pi^{\sigma^t}(hz, h') u_i(h') \\ &= \pi_{-i}^{\sigma^t}(I) (u_i(\sigma^t |_{Iz \rightarrow a}, Iz) - u_i(\sigma^t, Iz)) \end{aligned} \quad (31)$$

The equation above connects the definitions of $R_i^{T,H}$ and $R_i^{T,L}$ in Equation (25) and (6).

G PROOF OF THEOREM 4

Lemma 4. For all $h \in H \setminus H_{TS}$, $z \in Z(h)$, $a \in A(h)$:

$$\begin{aligned} E_{h'} \left[\hat{v}_i^{t,H}(\sigma^t, h, z|h') | h' \supseteq h \right] &= E_{h'} \left[\hat{v}_i^{t,L}(\sigma^t, hz|h') | h' \supseteq hz \right] \\ E_{h'} \left[\hat{v}_i^{t,L}(\sigma^t, hz, a|h') | h' \supseteq hz \right] &= E_{h'} \left[\hat{v}_i^{t,H}(\sigma^t, hza|h') | h' \supseteq hza \right] \end{aligned} \quad (32)$$

972 *Proof.*

$$\begin{aligned}
973 & E_{h'} \left[\hat{v}_i^{t,H}(\sigma^t, h, z|h') | h' \supseteq h \right] \\
974 & = E_{h'} \left[\frac{\delta(hz \sqsubseteq h')}{q^t(z|h)} \left[\hat{v}_i^{t,L}(\sigma^t, hz|h') - b_i^t(h, z) \right] + b_i^t(h, z) | h' \supseteq h \right] \\
975 & = P(hz \sqsubseteq h' | h' \supseteq h) E_{h'} \left[\frac{1}{q^t(z|h)} \left[\hat{v}_i^{t,L}(\sigma^t, hz|h') - b_i^t(h, z) \right] + b_i^t(h, z) | h' \supseteq hz \right] + \\
976 & \quad P(hz \not\sqsubseteq h' | h' \supseteq h) b_i^t(h, z) \\
977 & = q^t(z|h) \left[\frac{1}{q^t(z|h)} \left[E_{h'}(\hat{v}_i^{t,L}(\sigma^t, hz|h') | h' \supseteq hz) - b_i^t(h, z) \right] + b_i^t(h, z) \right] + (1 - q^t(z|h)) b_i^t(h, z) \\
978 & = E_{h'} \left[\hat{v}_i^{t,L}(\sigma^t, hz|h') | h' \supseteq hz \right] \\
979 & \hspace{20em} (33)
\end{aligned}$$

980 Using the definition of $\hat{v}_i^{t,L}(\sigma^t, hz, a|h')$ in Equation (27) and following the same process as above,
981 we can get the second part of the lemma. \square

982 Now, we present the proof of the first part of Theorem 4.

$$\begin{aligned}
983 & \mathbb{E}_{h' \sim \pi^{q^t}(\cdot)} \left[\hat{r}_i^{t,H}(I, z|h') \right] = \sum_{h \in I} \frac{\pi^{\sigma^t}(h)}{\pi^{q^t}(h)} \left[\mathbb{E}_{h'} \left[\hat{v}_i^{t,H}(\sigma^t, h, z|h') \right] - \mathbb{E}_{h'} \left[\hat{v}_i^{t,H}(\sigma^t, h|h') \right] \right] \\
984 & = \sum_{h \in I} \pi^{\sigma^t}(h) \left[\mathbb{E}_{h'} \left[\hat{v}_i^{t,H}(\sigma^t, h, z|h') | h' \supseteq h \right] - \mathbb{E}_{h'} \left[\hat{v}_i^{t,H}(\sigma^t, h|h') | h' \supseteq h \right] \right] \\
985 & \hspace{20em} (34)
\end{aligned}$$

986 For the second equality, we use the following fact:

$$\begin{aligned}
987 & \mathbb{E}_{h'} \left[\hat{v}_i^{t,H}(\sigma^t, h|h') \right] = P(h' \supseteq h) \mathbb{E}_{h'} \left[\hat{v}_i^{t,H}(\sigma^t, h|h') | h' \supseteq h \right] + P(h' \not\supseteq h) \mathbb{E}_{h'} \left[\hat{v}_i^{t,H}(\sigma^t, h|h') | h' \not\supseteq h \right] \\
988 & = \pi^{q^t}(h) \mathbb{E}_{h'} \left[\hat{v}_i^{t,H}(\sigma^t, h|h') | h' \supseteq h \right] \\
989 & \hspace{20em} (35)
\end{aligned}$$

990 Based on Equation (27), $\mathbb{E}_{h'} \left[\hat{v}_i^{t,H}(\sigma^t, h|h') | h' \not\supseteq h \right] = \mathbb{E}_{h'} \left[\hat{v}_i^{t,H}(\sigma^t, h, z|h') | h' \not\supseteq h \right] = 0$. Similar
991 with Equation (35), we can get $\mathbb{E}_{h'} \left[\hat{v}_i^{t,H}(\sigma^t, h, z|h') \right] = \pi^{q^t}(h) \mathbb{E}_{h'} \left[\hat{v}_i^{t,H}(\sigma^t, h, z|h') | h' \supseteq h \right]$,
992 which completes the proof of Equation (34).

993 Equation (25) and (34) show that, to prove $\mathbb{E}_{h' \sim \pi^{q^t}(\cdot)} \left[\hat{r}_i^{t,H}(I, z|h') \right] = r_i^{t,H}(I, z)$, we only need to
994 show the following lemma:

995 **Lemma 5.** For all $h \in H$, $z \in Z(h)$:

$$996 E_{h'} \left[\hat{v}_i^{t,H}(\sigma^t, h, z|h') | h' \supseteq h \right] = v_i^{t,L}(\sigma^t, hz), \quad E_{h'} \left[\hat{v}_i^{t,H}(\sigma^t, h|h') | h' \supseteq h \right] = v_i^{t,H}(\sigma^t, h) \quad (36)$$

997 *Proof.* We prove this lemma by induction on the height of h on the game tree. For the base case, if
998 $(hza) \in H_{TS}$, we have:

$$\begin{aligned}
999 & E_{h'} \left[\hat{v}_i^{t,H}(\sigma^t, h, z|h') | h' \supseteq h \right] = E_{h'} \left[\hat{v}_i^{t,L}(\sigma^t, hz|h') | h' \supseteq hz \right] \\
1000 & = \sum_{a \in A(h)} \sigma_{P(h)}^{t,L}(a|h, z) E_{h'} \left[\hat{v}_i^{t,L}(\sigma^t, hz, a|h') | h' \supseteq hz \right] \\
1001 & = \sum_{a \in A(h)} \sigma_{P(h)}^{t,L}(a|h, z) E_{h'} \left[\hat{v}_i^{t,H}(\sigma^t, hza|h') | h' \supseteq hza \right] \\
1002 & = \sum_{a \in A(h)} \sigma_{P(h)}^{t,L}(a|h, z) u_i(hza) = v_i^{t,L}(\sigma^t, hz) \\
1003 & \hspace{20em} (37)
\end{aligned}$$

Here, the first and third equality are due to Lemma 4, and the others are based on the corresponding definitions. Still, for this base case, we have:

$$\begin{aligned} E_{h'} \left[\hat{v}_i^{t,H}(\sigma^t, h|h') | h' \supseteq h \right] &= \sum_{z \in Z(h)} \sigma_{P(h)}^{t,H}(z|h) E_{h'} \left[\hat{v}_i^{t,H}(\sigma^t, h, z|h') | h' \supseteq h \right] \\ &= \sum_{z \in Z(h)} \sigma_{P(h)}^{t,H}(z|h) v_i^{t,L}(\sigma^t, hz) = v_i^{t,H}(\sigma^t, h) \end{aligned} \quad (38)$$

where the second equality comes for Equation (37). Then, we can move on to the general case, with the hypothesis that Lemma 5 holds for the nodes lower than h on the game tree:

$$\begin{aligned} E_{h'} \left[\hat{v}_i^{t,H}(\sigma^t, h, z|h') | h' \supseteq h \right] &= E_{h'} \left[\hat{v}_i^{t,L}(\sigma^t, hz|h') | h' \supseteq hz \right] \\ &= \sum_{a \in A(h)} \sigma_{P(h)}^{t,L}(a|h, z) E_{h'} \left[\hat{v}_i^{t,L}(\sigma^t, hz, a|h') | h' \supseteq hz \right] \\ &= \sum_{a \in A(h)} \sigma_{P(h)}^{t,L}(a|h, z) E_{h'} \left[\hat{v}_i^{t,H}(\sigma^t, hza|h') | h' \supseteq hza \right] \\ &= \sum_{a \in A(h)} \sigma_{P(h)}^{t,L}(a|h, z) v_i^{t,H}(\sigma^t, hza) = v_i^{t,L}(\sigma^t, hz) \end{aligned} \quad (39)$$

where the induction hypothesis is adopted for the fourth equality. Equation (39) and (38) imply that $E_{h'} \left[\hat{v}_i^{t,H}(\sigma^t, h|h') | h' \supseteq h \right] = v_i^{t,H}(\sigma^t, h)$ holds for the general case. \square

So far, we have proved the first part of Theorem 4, i.e., $\mathbb{E}_{h' \sim \pi^{qt}(\cdot)} \left[\hat{r}_i^{t,H}(I, z|h') \right] = r_i^{t,H}(I, z)$. The second part, $\mathbb{E}_{h' \sim \pi^{qt}(\cdot)} \left[\hat{r}_i^{t,L}(Iz, a|h') \right] = r_i^{t,L}(Iz, a)$, can be proved with the same process as above based on Lemma 4, so we skip the complete proof and only present the following lemma within it.

Lemma 6. For all $h \in H \setminus H_{TS}$, $z \in Z(h)$, $a \in A(h)$:

$$E_{h'} \left[\hat{v}_i^{t,L}(\sigma^t, hz, a|h') | h' \supseteq hz \right] = v_i^{t,H}(\sigma^t, hza), \quad E_{h'} \left[\hat{v}_i^{t,L}(\sigma^t, hz|h') | h' \supseteq hz \right] = v_i^{t,L}(\sigma^t, hz) \quad (40)$$

H PROOF OF THEOREM 5

Part I:

First, we can apply the law of total variance to $\text{Var}_{h' \sim \pi^{qt}(\cdot)} \left[\hat{r}_i^{t,H}(I, z|h') \right]$, conditioning on $\delta(h' \supseteq I)$ (i.e., if h' is reachable from I), and get:

$$\text{Var}_{h' \sim \pi^{qt}(\cdot)} \left[\hat{r}_i^{t,H}(I, z|h') \right] = \mathbb{E} \left[\text{Var}_{h'} \left[\hat{r}_i^{t,H}(I, z|h') | \delta(h' \supseteq I) \right] \right] + \text{Var} \left[\mathbb{E}_{h'} \left[\hat{r}_i^{t,H}(I, z|h') | \delta(h' \supseteq I) \right] \right] \quad (41)$$

The first term can be expanded as follows, where the second equality is due to $\hat{r}_i^{t,H}(I, z|h') = 0$ when $h' \not\supseteq I$.

$$\begin{aligned} &\mathbb{E} \left[\text{Var}_{h'} \left[\hat{r}_i^{t,H}(I, z|h') | \delta(h' \supseteq I) \right] \right] \\ &= P(h' \supseteq I) \text{Var}_{h'} \left[\hat{r}_i^{t,H}(I, z|h') | h' \supseteq I \right] + P(h' \not\supseteq I) \text{Var}_{h'} \left[\hat{r}_i^{t,H}(I, z|h') | h' \not\supseteq I \right] \\ &= P(h' \supseteq I) \text{Var}_{h'} \left[\hat{r}_i^{t,H}(I, z|h') | h' \supseteq I \right] \end{aligned} \quad (42)$$

The second term can be converted as follows, based on the fact that $\mathbb{E}_{h'}(\hat{r}_i^{t,H}(I, z|h') | \delta(h' \supseteq I)) = \frac{r_i^{t,H}(I, z)}{P(h' \supseteq I)}$ (i.e., $\mathbb{E}_{h'}(\hat{r}_i^{t,H}(I, z|h') | h' \supseteq I)$) with probability $P(h' \supseteq I)$, and $\mathbb{E}_{h'}(\hat{r}_i^{t,H}(I, z|h') | \delta(h' \not\supseteq I)) = 0$.

$\delta(h' \supseteq I) = 0$ (i.e., $\mathbb{E}_{h'}(\hat{r}_i^{t,H}(I, z|h')|h' \not\supseteq I)$) with probability $1 - P(h' \supseteq I)$.

$$\begin{aligned} & \text{Var} \left[\mathbb{E}_{h'} \left[\hat{r}_i^{t,H}(I, z|h') | \delta(h' \supseteq I) \right] \right] \\ &= \mathbb{E} \left[\left[\mathbb{E}_{h'}(\hat{r}_i^{t,H}(I, z|h') | \delta(h' \supseteq I)) \right]^2 \right] - \left[\mathbb{E} \left[\mathbb{E}_{h'}(\hat{r}_i^{t,H}(I, z|h') | \delta(h' \supseteq I)) \right] \right]^2 \\ &= \frac{1 - P(h' \supseteq I)}{P(h' \supseteq I)} (r_i^{t,H}(I, z))^2 \end{aligned} \quad (43)$$

Note that $\frac{1 - P(h' \supseteq I)}{P(h' \supseteq I)} (r_i^{t,H}(I, z))^2$ and $P(h' \supseteq I)$ is not affected by b_i^t , so we focus on $\text{Var}_{h'} \left[\hat{r}_i^{t,H}(I, z|h') | h' \supseteq I \right]$ in Equation (42). Applying the law of total variance:

$$\begin{aligned} & \text{Var}_{h'} \left[\hat{r}_i^{t,H}(I, z|h') | h' \supseteq I \right] \\ &= \mathbb{E}_{h \in I} \left[\text{Var}_{h'} \left[\hat{r}_i^{t,H}(I, z|h') | h' \supseteq h \right] \right] + \text{Var}_{h \in I} \left[\mathbb{E}_{h'} \left[\hat{r}_i^{t,H}(I, z|h') | h' \supseteq h \right] \right] \\ &\geq \text{Var}_{h \in I} \left[\mathbb{E}_{h'} \left[\hat{r}_i^{t,H}(I, z|h') | h' \supseteq h \right] \right] \end{aligned} \quad (44)$$

Fix $h \in I$, $\mathbb{E}_{h'} \left[\hat{r}_i^{t,H}(I, z|h') | h' \supseteq h \right] = \frac{\pi_{-i}^{\sigma^t}(h)}{\pi^{q^t}(h)} \left[v_i^{t,L}(\sigma^t, hz) - v_i^{t,H}(\sigma^t, h) \right]$, based on the definition of $\hat{r}_i^{t,H}(I, z|h')$ and Lemma 5. Thus, the second term in Equation (44) is irrelevant to b_i^t . According to Equation (41)-(44), we conclude that the minimum of $\text{Var}_{h' \sim \pi^{q^t}(\cdot)} \left[\hat{r}_i^{t,H}(I, z|h') \right]$ with respect to b_i^t can be achieved when $\mathbb{E}_{h \in I} \left[\text{Var}_{h'} \left[\hat{r}_i^{t,H}(I, z|h') | h' \supseteq h \right] \right] = 0$. Following the same process, we can show that the minimum of $\text{Var}_{h' \sim \pi^{q^t}(\cdot)} \left[\hat{r}_i^{t,L}(Iz, a|h') \right]$ with respect to b_i^t can be achieved when $\mathbb{E}_{h \in I} \left[\text{Var}_{h'} \left[\hat{r}_i^{t,L}(Iz, a|h') | h' \supseteq hz \right] \right] = 0$.

Lemma 7. *If $\text{Var}_{h'} \left[\hat{v}_i^{t,H}(\sigma^t, h, z|h') | h' \supseteq h \right] = \text{Var}_{h'} \left[\hat{v}_i^{t,L}(\sigma^t, hz, a|h') | h' \supseteq hz \right] = 0$, for all $h \in H \setminus H_{TS}$, $z \in Z(h)$, $a \in A(h)$, then $\mathbb{E}_{h \in I} \left[\text{Var}_{h'} \left[\hat{r}_i^{t,H}(I, z|h') | h' \supseteq h \right] \right] = \mathbb{E}_{h \in I} \left[\text{Var}_{h'} \left[\hat{r}_i^{t,L}(Iz, a|h') | h' \supseteq hz \right] \right] = 0, \forall I \in \mathcal{I}_i, z \in Z(I), a \in A(I)$.*

Proof. Pick any $h \in H \setminus H_{TS}$, $z \in Z(h)$. Based on Lemma 5, $E_{h'} \left[\hat{v}_i^{t,H}(\sigma^t, h, z|h') | h' \supseteq h \right] = v_i^{t,L}(\sigma^t, hz)$. If $\text{Var}_{h'} \left[\hat{v}_i^{t,H}(\sigma^t, h, z|h') | h' \supseteq h \right] = 0$, then $\hat{v}_i^{t,H}(\sigma^t, h, z|h') = v_i^{t,L}(\sigma^t, hz), \forall h' \supseteq h$. It follows that $\hat{v}_i^{t,H}(\sigma^t, h|h') = v_i^{t,H}(\sigma^t, h), \forall h' \supseteq h$, based on the definitions of $v_i^{t,H}(\sigma^t, h)$ and $\hat{v}_i^{t,H}(\sigma^t, h|h')$. Now, for any $I \in \mathcal{I}_i, h \in I, h' \supseteq h$:

$$\begin{aligned} \hat{r}_i^{t,H}(I, z|h') &= \sum_{h'' \in I} \frac{\pi_{-i}^{\sigma^t}(h'')}{\pi^{q^t}(h'')} \left[\hat{v}_i^{t,H}(\sigma^t, h'', z|h') - \hat{v}_i^{t,H}(\sigma^t, h''|h') \right] \\ &= \frac{\pi_{-i}^{\sigma^t}(h)}{\pi^{q^t}(h)} \left[\hat{v}_i^{t,H}(\sigma^t, h, z|h') - \hat{v}_i^{t,H}(\sigma^t, h|h') \right] \\ &= \frac{\pi_{-i}^{\sigma^t}(h)}{\pi^{q^t}(h)} \left[v_i^{t,L}(\sigma^t, hz) - v_i^{t,H}(\sigma^t, h) \right] \end{aligned} \quad (45)$$

Thus, $\text{Var}_{h'} \left[\hat{r}_i^{t,H}(I, z|h') | h' \supseteq h \right] = 0, \forall I \in \mathcal{I}_i, h \in I$. Then, it follows that for any $I, \mathbb{E}_{h \in I} \left[\text{Var}_{h'} \left[\hat{r}_i^{t,H}(I, z|h') | h' \supseteq h \right] \right] = 0$. With the same process as above, we can show the second part of Lemma 7. \square

Given the discussions above, to complete the proof of Theorem 5, we need to further show that, $\forall i \in N$, if $b_i^t(h, z, a) = v_i^{t,H}(\sigma^t, hza)$ and $b_i^t(h, z) = v_i^{t,L}(\sigma^t, hz)$, for all $h \in H \setminus H_{TS}, z \in Z(h), a \in$

1134 $A(h)$, we have $\text{Var}_{h'} \left[\hat{v}_i^{t,H}(\sigma^t, h, z|h')|h' \supseteq h \right] = \text{Var}_{h'} \left[\hat{v}_i^{t,L}(\sigma^t, hz, a|h')|h' \supseteq hz \right] = 0$, for all
 1135 $h \in H \setminus H_{TS}$, $z \in Z(h)$, $a \in A(h)$.
 1136

1137 **Part II:**

1138 **Lemma 8.** For any $i \in N$, $h \in H \setminus H_{TS}$, $z \in Z(h)$, $a \in A(h)$ and any baseline function b_i^t :

$$\begin{aligned}
 1140 \quad \text{Var}_{h'} \left[\hat{v}_i^{t,H}(\sigma^t, h|h')|h' \supseteq h \right] &= \sum_{z \in Z(h)} \frac{(\sigma_{P(h)}^{t,H}(z|h))^2}{q^t(z|h)} \text{Var}_{h'} \left[\hat{v}_i^{t,L}(\sigma^t, hz|h')|h' \supseteq hz \right] \\
 1141 &+ \text{Var}_{z \sim q^t(\cdot|h)} \left[\frac{\sigma_{P(h)}^{t,H}(z|h)}{q^t(z|h)} (v_i^{t,L}(\sigma^t, hz) - b_i^t(h, z)) \right] \\
 1142 & \\
 1143 \quad \text{Var}_{h'} \left[\hat{v}_i^{t,L}(\sigma^t, hz|h')|h' \supseteq hz \right] &= \sum_{a \in A(h)} \frac{(\sigma_{P(h)}^{t,L}(a|h, z))^2}{q^t(a|h, z)} \text{Var}_{h'} \left[\hat{v}_i^{t,H}(\sigma^t, hza|h')|h' \supseteq hza \right] \\
 1144 &+ \text{Var}_a \left[\frac{\sigma_{P(h)}^{t,L}(a|h, z)}{q^t(a|h, z)} (v_i^{t,H}(\sigma^t, hza) - b_i^t(h, z, a)) \right] \\
 1145 & \\
 1146 & \\
 1147 & \\
 1148 & \\
 1149 & \\
 1150 & \\
 1151 & \\
 1152 & \tag{46}
 \end{aligned}$$

1153 *Proof.* By conditioning on the option choice at h , we apply the law of total variance to
 1154 $\text{Var}_{h'} \left[\hat{v}_i^{t,H}(\sigma^t, h|h')|h' \supseteq h \right]$:

$$\begin{aligned}
 1155 \quad \text{Var}_{h'} \left[\hat{v}_i^{t,H}(\sigma^t, h|h')|h' \supseteq h \right] &= \mathbb{E}_{z \sim q^t(\cdot|h)} \left[\text{Var}_{h'} \left[\hat{v}_i^{t,H}(\sigma^t, h|h')|h' \supseteq hz \right] \right] + \\
 1156 & \\
 1157 & \text{Var}_{z \sim q^t(\cdot|h)} \left[\mathbb{E}_{h'} \left[\hat{v}_i^{t,H}(\sigma^t, h|h')|h' \supseteq hz \right] \right] \\
 1158 & \\
 1159 & \tag{47}
 \end{aligned}$$

1160 According to the definition of $\hat{v}_i^{t,H}(\sigma^t, h|h')$ and the fact that $h' \supseteq hz$, we have:

$$\begin{aligned}
 1161 \quad \text{Var}_{h'} \left[\hat{v}_i^{t,H}(\sigma^t, h|h')|h' \supseteq hz \right] & \\
 1162 &= \text{Var}_{h'} \left[\frac{\sigma_{P(h)}^{t,H}(z|h)}{q^t(z|h)} \left[\hat{v}_i^{t,L}(\sigma^t, hz|h') - b_i^t(h, z) \right] + \sum_{z' \in Z(h)} \sigma_{P(h)}^{t,H}(z'|h) b_i^t(h, z') |h' \supseteq hz \right] \\
 1163 & \\
 1164 &= \left[\frac{\sigma_{P(h)}^{t,H}(z|h)}{q^t(z|h)} \right]^2 \text{Var}_{h'} \left[\hat{v}_i^{t,L}(\sigma^t, hz|h')|h' \supseteq hz \right] \\
 1165 & \\
 1166 & \\
 1167 & \\
 1168 & \\
 1169 & \\
 1170 & \\
 1171 \quad \mathbb{E}_{z \sim q^t(\cdot|h)} \left[\text{Var}_{h'} \left[\hat{v}_i^{t,H}(\sigma^t, h|h')|h' \supseteq hz \right] \right] &= \sum_{z \in Z(h)} \frac{(\sigma_{P(h)}^{t,H}(z|h))^2}{q^t(z|h)} \text{Var}_{h'} \left[\hat{v}_i^{t,L}(\sigma^t, hz|h')|h' \supseteq hz \right] \\
 1172 & \\
 1173 & \tag{48}
 \end{aligned}$$

1174 Then, we analyze the second term in Equation (47):

$$\begin{aligned}
 1175 \quad \mathbb{E}_{h'} \left[\hat{v}_i^{t,H}(\sigma^t, h|h')|h' \supseteq hz \right] & \\
 1176 &= \frac{\sigma_{P(h)}^{t,H}(z|h)}{q^t(z|h)} \left[\mathbb{E}_{h'} \left[\hat{v}_i^{t,L}(\sigma^t, hz|h')|h' \supseteq hz \right] - b_i^t(h, z) \right] + \sum_{z' \in Z(h)} \sigma_{P(h)}^{t,H}(z'|h) b_i^t(h, z') \\
 1177 & \\
 1178 &= \frac{\sigma_{P(h)}^{t,H}(z|h)}{q^t(z|h)} \left[v_i^{t,L}(\sigma^t, hz) - b_i^t(h, z) \right] + \sum_{z' \in Z(h)} \sigma_{P(h)}^{t,H}(z'|h) b_i^t(h, z') \\
 1179 & \\
 1180 & \\
 1181 & \\
 1182 & \\
 1183 & \\
 1184 \quad \text{Var}_{z \sim q^t(\cdot|h)} \left[\mathbb{E}_{h'} \left[\hat{v}_i^{t,H}(\sigma^t, h|h')|h' \supseteq hz \right] \right] &= \text{Var}_{z \sim q^t(\cdot|h)} \left[\frac{\sigma_{P(h)}^{t,H}(z|h)}{q^t(z|h)} \left[v_i^{t,L}(\sigma^t, hz) - b_i^t(h, z) \right] \right] \\
 1185 & \\
 1186 & \tag{49}
 \end{aligned}$$

1187 Based on Equation (47)-(49), we can get the first part of Lemma 8. The second part can be obtained similarly. \square

Lemma 8 illustrates the outcome of a single-step lookahead from state h . Employing this in an inductive manner, we can derive the complete expansion of $\text{Var}_{h'} \left[\hat{v}_i^{t,H}(\sigma^t, h|h') | h' \supseteq h \right]$ on the game tree as the following lemma:

Lemma 9. For any $i \in N$, $h \in H$ and any baseline function b_i^t :

$$\begin{aligned} \text{Var}_{h'} \left[\hat{v}_i^{t,H}(\sigma^t, h|h') | h' \supseteq h \right] &= \sum_{h'' \supseteq h} \frac{(\pi^{\sigma^t}(h, h''))^2}{\pi^{q^t}(h, h'')} f(h'') \\ f(h'') &= \text{Var}_z \left[\frac{\sigma_{P(h'')}^{t,H}(z|h'')}{q^t(z|h'')} (v_i^{t,L}(\sigma^t, h''z) - b_i^t(h'', z)) \right] + \\ &\quad \sum_{z \in Z(h'')} \frac{(\sigma_{P(h'')}^{t,H}(z|h''))^2}{q^t(z|h'')} \text{Var}_a \left[\frac{\sigma_{P(h'')}^{t,L}(a|h'', z)}{q^t(a|h'', z)} (v_i^{t,H}(\sigma^t, h''za) - b_i^t(h'', z, a)) \right] \end{aligned} \quad (50)$$

Proof. We proof this lemma through an induction on the height of h on the game tree. For the base case, $h \in H_{TS}$, then $Z(h) = A(h) = \emptyset$, so $f(h'') = 0$, $\forall h'' \supseteq h$. In addition, we have $\text{Var}_{h'} \left[\hat{v}_i^{t,H}(\sigma^t, h|h') | h' \supseteq h \right] = \text{Var}_{h'} \left[\hat{v}_i^{t,H}(\sigma^t, h|h') | h' = h \right] = 0$. Thus, the lemma holds for the base case.

For the general case, $h \in H \setminus H_{TS}$, we apply Lemma 8 and get:

$$\begin{aligned} \text{Var}_{h'} \left[\hat{v}_i^{t,H}(\sigma^t, h|h') | h' \supseteq h \right] &= \text{Var}_{z \sim q^t(\cdot|h)} \left[\frac{\sigma_{P(h)}^{t,H}(z|h)}{q^t(z|h)} (v_i^{t,L}(\sigma^t, hz) - b_i^t(h, z)) \right] \\ &+ \sum_{z \in Z(h)} \frac{(\sigma_{P(h)}^{t,H}(z|h))^2}{q^t(z|h)} \text{Var}_{a \sim q^t(\cdot|h,z)} \left[\frac{\sigma_{P(h)}^{t,L}(a|h, z)}{q^t(a|h, z)} (v_i^{t,H}(\sigma^t, hza) - b_i^t(h, z, a)) \right] \\ &+ \sum_{(z,a) \in \tilde{A}(h)} \frac{(\sigma_{P(h)}^t((z,a)|h))^2}{q^t((z,a)|h)} \text{Var}_{h'} \left[\hat{v}_i^{t,H}(\sigma^t, hza|h') | h' \supseteq hza \right] \\ &= f(h) + \sum_{(z,a) \in \tilde{A}(h)} \frac{(\sigma_{P(h)}^t((z,a)|h))^2}{q^t((z,a)|h)} \text{Var}_{h'} \left[\hat{v}_i^{t,H}(\sigma^t, hza|h') | h' \supseteq hza \right] \end{aligned} \quad (51)$$

where the first equality is the result of the sequential use of the two formulas in Lemma 8 and the second equality is based on the definition of $f(h)$. Next, we apply the induction hypothesis on hza , i.e., a node lower than h on the game tree, and get:

$$\text{Var}_{h'} \left[\hat{v}_i^{t,H}(\sigma^t, hza|h') | h' \supseteq hza \right] = \sum_{h'' \supseteq hza} \frac{(\pi^{\sigma^t}(hza, h''))^2}{\pi^{q^t}(hza, h'')} f(h'') \quad (52)$$

By integrating Equation (51) and (52), we can get:

$$\begin{aligned} \text{Var}_{h'} \left[\hat{v}_i^{t,H}(\sigma^t, h|h') | h' \supseteq h \right] &= f(h) + \sum_{(z,a)} \frac{(\sigma_{P(h)}^t((z,a)|h))^2}{q^t((z,a)|h)} \sum_{h'' \supseteq hza} \frac{(\pi^{\sigma^t}(hza, h''))^2}{\pi^{q^t}(hza, h'')} f(h'') \\ &= f(h) + \sum_{h'' \supseteq h} \frac{(\pi^{\sigma^t}(h, h''))^2}{\pi^{q^t}(h, h'')} f(h'') \\ &= \sum_{h'' \supseteq h} \frac{(\pi^{\sigma^t}(h, h''))^2}{\pi^{q^t}(h, h'')} f(h'') \end{aligned} \quad (53)$$

For the second equality, we use the definitions of $\pi^{\sigma^t}(h, h'')$ and $\pi^{q^t}(h, h'')$, and the fact that they equal 1 when $h'' = h$. \square

Before moving to the final proof, we introduce another lemma as follows.

Lemma 10. For any $i \in N$, $h \in H \setminus H_{TS}$, $z \in Z(h)$ and any baseline function b_i^t :

$$\begin{aligned} \text{Var}_{h'} \left[\hat{v}_i^{t,L}(\sigma^t, hz|h') | h' \supseteq hz \right] &\leq \sum_{\substack{a \in A(h), \\ h'' \supseteq hza, \\ z'' \in Z(h'')}} \frac{(\pi^{\sigma^t}(hz, h''z''))^2}{\pi^{q^t}(hz, h''z'')} \left[v_i^{t,L}(\sigma^t, h''z'') - b_i^t(h'', z'') \right]^2 \\ &+ \sum_{\substack{h''z'' \supseteq hz, \\ a'' \in A(h'')}} \frac{(\pi^{\sigma^t}(hz, h''z''a''))^2}{\pi^{q^t}(hz, h''z''a'')} \left[v_i^{t,H}(\sigma^t, h''z''a'') - b_i^t(h'', z'', a'') \right]^2 \end{aligned} \quad (54)$$

Proof. Applying the fact $\text{Var}(X) = \mathbb{E}(X^2) - (\mathbb{E}(X))^2 \leq \mathbb{E}(X^2)$ to both variance terms of Equation (50) and after rearranging the terms, we arrive at the following expression:

$$\begin{aligned} \text{Var}_{h'} \left[\hat{v}_i^{t,H}(\sigma^t, h|h') | h' \supseteq h \right] &\leq \sum_{\substack{h'' \supseteq h, \\ z'' \in Z(h'')}} \frac{(\pi^{\sigma^t}(h, h''z''))^2}{\pi^{q^t}(h, h''z'')} \left[v_i^{t,L}(\sigma^t, h''z'') - b_i^t(h'', z'') \right]^2 \\ &+ \sum_{\substack{h'' \supseteq h, \\ (z'', a'') \in \bar{A}(h'')}} \frac{(\pi^{\sigma^t}(h, h''z''a''))^2}{\pi^{q^t}(h, h''z''a'')} \left[v_i^{t,H}(\sigma^t, h''z''a'') - b_i^t(h'', z'', a'') \right]^2 \end{aligned} \quad (55)$$

Note that the equation above holds for any $h \in H$. Then, to get an upper bound of $\text{Var}_{h'} \left[\hat{v}_i^{t,L}(\sigma^t, hz|h') | h' \supseteq hz \right]$, we go back to Lemma 8 and apply Equation (55) and $\text{Var}(X) \leq \mathbb{E}(X^2)$ to its first and second term, respectively. After rearranging, we can get:

$$\begin{aligned} \text{Var}_{h'} \left[\hat{v}_i^{t,L}(\sigma^t, hz|h') | h' \supseteq hz \right] &\leq \sum_{\substack{a \in A(h), \\ h'' \supseteq hza, \\ z'' \in Z(h'')}} \frac{(\pi^{\sigma^t}(hz, h''z''))^2}{\pi^{q^t}(hz, h''z'')} \left[v_i^{t,L}(\sigma^t, h''z'') - b_i^t(h'', z'') \right]^2 \\ &+ \sum_{\substack{a \in A(h), \\ h'' \supseteq hza, \\ (z'', a'') \in \bar{A}(h'')}} \frac{(\pi^{\sigma^t}(hz, h''z''a''))^2}{\pi^{q^t}(hz, h''z''a'')} \left[v_i^{t,H}(\sigma^t, h''z''a'') - b_i^t(h'', z'', a'') \right]^2 \\ &+ \sum_{a \in A(h)} \frac{(\sigma_{P(h)}^{t,L}(a|h, z))^2}{q^t(a|h, z)} \left[v_i^{t,H}(\sigma^t, hza) - b_i^t(h, z, a) \right]^2 \end{aligned} \quad (56)$$

We note that the second term of Equation (54) can be obtained by combining the last two terms of Equation (56). The second and third term of Equation (56) correspond to the sum over $h''z'' \supseteq hz$, $a'' \in A(h'')$ and $h''z'' = hz$, $a'' \in A(h'')$, respectively. \square

Based on the discussions above, we give out the upper bound of $\text{Var}_{h'} \left[\hat{v}_i^{t,H}(\sigma^t, h, z|h') | h' \supseteq h \right]$ as the following lemma:

Lemma 11. For any $i \in N$, $h \in H \setminus H_{TS}$, $z \in Z(h)$ and any baseline function b_i^t :

$$\begin{aligned} \text{Var}_{h'} \left[\hat{v}_i^{t,H}(\sigma^t, h, z|h') | h' \supseteq h \right] &\leq \frac{1}{q^t(z|h)} \sum_{h''z'' \supseteq hz} \frac{(\pi^{\sigma^t}(hz, h''z''))^2}{\pi^{q^t}(hz, h''z'')} \left[v_i^{t,L}(\sigma^t, h''z'') - b_i^t(h'', z'') \right]^2 \\ &+ \frac{1}{q^t(z|h)} \sum_{\substack{h''z'' \supseteq hz, \\ a'' \in A(h'')}} \frac{(\pi^{\sigma^t}(hz, h''z''a''))^2}{\pi^{q^t}(hz, h''z''a'')} \left[v_i^{t,H}(\sigma^t, h''z''a'') - b_i^t(h'', z'', a'') \right]^2 \end{aligned} \quad (57)$$

1296

1297

Proof.

1298

1299

1300

1301

1302

1303

1304

1305

1306

$$\begin{aligned}
& \text{Var}_{h'} \left[\hat{v}_i^{t,H}(\sigma^t, h, z|h') | h' \supseteq h \right] \\
&= \text{Var}_{h'} \left[\frac{\delta(hz \sqsubseteq h')}{q^t(z|h)} \left[\hat{v}_i^{t,L}(\sigma^t, hz|h') - b_i^t(h, z) \right] | h' \supseteq h \right] \\
&= \mathbb{E} \left[\text{Var}_{h'} \left[\frac{\delta(hz \sqsubseteq h')}{q^t(z|h)} \left[\hat{v}_i^{t,L}(\sigma^t, hz|h') - b_i^t(h, z) \right] | h' \supseteq h, \delta(hz \sqsubseteq h') \right] \right] + \\
& \quad \text{Var} \left[\mathbb{E}_{h'} \left[\frac{\delta(hz \sqsubseteq h')}{q^t(z|h)} \left[\hat{v}_i^{t,L}(\sigma^t, hz|h') - b_i^t(h, z) \right] | h' \supseteq h, \delta(hz \sqsubseteq h') \right] \right]
\end{aligned} \tag{58}$$

1307

1308

1309

Here, we apply the definition of $\hat{v}_i^{t,H}(\sigma^t, h, z|h')$ to get the first equality, and the law of total variance conditioned on $\delta(hz \sqsubseteq h')$ (given $h \sqsubseteq h'$) to get the second equality. Next, we analyze the two terms in the third and fourth line of Equation (58) separately.

1310

1311

1312

1313

1314

1315

1316

1317

1318

$$\begin{aligned}
& \mathbb{E} \left[\text{Var}_{h'} \left[\frac{\delta(hz \sqsubseteq h')}{q^t(z|h)} \left[\hat{v}_i^{t,L}(\sigma^t, hz|h') - b_i^t(h, z) \right] | h' \supseteq h, \delta(hz \sqsubseteq h') \right] \right] \\
&= P(hz \sqsubseteq h' | h \sqsubseteq h') \text{Var}_{h'} \left[\frac{1}{q^t(z|h)} \left[\hat{v}_i^{t,L}(\sigma^t, hz|h') - b_i^t(h, z) \right] | h' \supseteq hz \right] \\
&= q^t(z|h) \text{Var}_{h'} \left[\frac{1}{q^t(z|h)} \left[\hat{v}_i^{t,L}(\sigma^t, hz|h') - b_i^t(h, z) \right] | h' \supseteq hz \right] \\
&= \frac{1}{q^t(z|h)} \text{Var}_{h'} \left[\hat{v}_i^{t,L}(\sigma^t, hz|h') | h' \supseteq hz \right]
\end{aligned} \tag{59}$$

1319

1320

1321

Note that $\delta(hz \sqsubseteq h')$ can be 0 or 1 (with probability $P(hz \sqsubseteq h' | h \sqsubseteq h')$), and the variance equals 0 when $\delta(hz \sqsubseteq h') = 0$, so we get the first equality in Equation (59). Similarly, we can get:

1322

1323

1324

1325

1326

1327

1328

1329

1330

1331

1332

1333

$$\begin{aligned}
& \text{Var} \left[\mathbb{E}_{h'} \left[\frac{\delta(hz \sqsubseteq h')}{q^t(z|h)} \left[\hat{v}_i^{t,L}(\sigma^t, hz|h') - b_i^t(h, z) \right] | h' \supseteq h, \delta(hz \sqsubseteq h') \right] \right] \\
&\leq \mathbb{E} \left[\left[\mathbb{E}_{h'} \left[\frac{\delta(hz \sqsubseteq h')}{q^t(z|h)} \left[\hat{v}_i^{t,L}(\sigma^t, hz|h') - b_i^t(h, z) \right] | h' \supseteq h, \delta(hz \sqsubseteq h') \right] \right]^2 \right] \\
&= q^t(z|h) \left[\mathbb{E}_{h'} \left[\frac{1}{q^t(z|h)} \left[\hat{v}_i^{t,L}(\sigma^t, hz|h') - b_i^t(h, z) \right] | h' \supseteq hz \right] \right]^2 \\
&= \frac{1}{q^t(z|h)} \left[\mathbb{E}_{h'} \left[\hat{v}_i^{t,L}(\sigma^t, hz|h') | h' \supseteq hz \right] - b_i^t(h, z) \right]^2 \\
&= \frac{1}{q^t(z|h)} \left[v_i^{t,L}(\sigma^t, hz) - b_i^t(h, z) \right]^2
\end{aligned} \tag{60}$$

1334

Integrating Equation (58)-(60) and utilizing the upper bound proposed in Lemma 10, we can get:

1335

1336

1337

1338

1339

1340

1341

1342

1343

1344

1345

1346

1347

1348

1349

$$\begin{aligned}
& \text{Var}_{h'} \left[\hat{v}_i^{t,H}(\sigma^t, h, z|h') | h' \supseteq h \right] \leq \frac{1}{q^t(z|h)} \left[v_i^{t,L}(\sigma^t, hz) - b_i^t(h, z) \right]^2 + \\
& \frac{1}{q^t(z|h)} \sum_{\substack{\alpha \in A(h), \\ h'' \supseteq hza, \\ z'' \in Z(h'')}} \frac{(\pi^{\sigma^t}(hz, h''z''))^2}{\pi^{q^t}(hz, h''z'')} \left[v_i^{t,L}(\sigma^t, h''z'') - b_i^t(h'', z'') \right]^2 + \\
& \frac{1}{q^t(z|h)} \sum_{\substack{h''z'' \supseteq hz, \\ a'' \in A(h'')}} \frac{(\pi^{\sigma^t}(hz, h''z''a''))^2}{\pi^{q^t}(hz, h''z''a'')} \left[v_i^{t,H}(\sigma^t, h''z''a'') - b_i^t(h'', z'', a'') \right]^2
\end{aligned} \tag{61}$$

Note that the sum of the first two terms of Equation (61) equals the first term of Equation (57). The first and second term of Equation (61) correspond to the sum over $h''z'' = hz$ and $h''z'' \supseteq hz$, respectively. \square

Similarly, we can derive the upper bound for $\text{Var}_{h'} \left[\hat{v}_i^{t,L}(\sigma^t, hz, a|h') | h' \supseteq hz \right]$ shown as follows.

Lemma 12. For any $i \in N$, $h \in H \setminus H_{TS}$, $z \in Z(h)$, $a \in A(h)$ and any baseline function b_i^t :

$$\begin{aligned} & \text{Var}_{h'} \left[\hat{v}_i^{t,L}(\sigma^t, hz, a|h') | h' \supseteq hz \right] \\ & \leq \frac{1}{q^t(a|h, z)} \sum_{h''z''a'' \supseteq hza} \frac{(\pi^{\sigma^t}(hza, h''z''a''))^2}{\pi^{q^t}(hza, h''z''a'')} \left[v_i^{t,H}(\sigma^t, h''z''a'') - b_i^t(h'', z'', a'') \right]^2 + \\ & \quad \frac{1}{q^t(a|h, z)} \sum_{\substack{h'' \supseteq hza, \\ z'' \in Z(h'')}} \frac{(\pi^{\sigma^t}(hza, h''z''))^2}{\pi^{q^t}(hza, h''z'')} \left[v_i^{t,L}(\sigma^t, h''z'') - b_i^t(h'', z'') \right]^2 \end{aligned} \quad (62)$$

Proof. By applying the law of total variance conditioned on $\delta(hza \sqsubseteq h')$ (given $hz \sqsubseteq h'$) and following the same process as Equation (58)-(60), we can get:

$$\begin{aligned} \text{Var}_{h'} \left[\hat{v}_i^{t,L}(\sigma^t, hz, a|h') | h' \supseteq hz \right] & \leq \frac{1}{q^t(a|h, z)} \left[v_i^{t,H}(\sigma^t, hza) - b_i^t(h, z, a) \right]^2 + \\ & \quad \frac{1}{q^t(a|h, z)} \text{Var}_{h'} \left[\hat{v}_i^{t,H}(\sigma^t, hza|h') | h' \supseteq hza \right] \end{aligned} \quad (63)$$

Then, we can apply the upper bound shown as Equation (55) and get:

$$\begin{aligned} \text{Var}_{h'} \left[\hat{v}_i^{t,L}(\sigma^t, hz, a|h') | h' \supseteq hz \right] & \leq \frac{1}{q^t(a|h, z)} \left[v_i^{t,H}(\sigma^t, hza) - b_i^t(h, z, a) \right]^2 + \\ & \quad \frac{1}{q^t(a|h, z)} \sum_{\substack{h'' \supseteq hza, \\ (z'', a'') \in \bar{A}(h'')}} \frac{(\pi^{\sigma^t}(hza, h''z''a''))^2}{\pi^{q^t}(hza, h''z''a'')} \left[v_i^{t,H}(\sigma^t, h''z''a'') - b_i^t(h'', z'', a'') \right]^2 + \\ & \quad \frac{1}{q^t(a|h, z)} \sum_{\substack{h'' \supseteq hza, \\ z'' \in Z(h'')}} \frac{(\pi^{\sigma^t}(hza, h''z''))^2}{\pi^{q^t}(hza, h''z'')} \left[v_i^{t,L}(\sigma^t, h''z'') - b_i^t(h'', z'') \right]^2 \end{aligned} \quad (64)$$

Again, we can combine the first two terms of Equation (64) and get the first term of the right hand side of Equation (62), since the first term of Equation (64) corresponds to the case that $h'' = h$, $h''z''a'' \supseteq hza$ and the second term is equivalent to the sum over $h'' \supseteq h$, $h''z''a'' \supseteq hza$. \square

Finally, with Lemma 11 - 12 and the fact that variance cannot be negative, we can claim: $\forall i \in N$, if $b_i^t(h, z, a) = v_i^{t,H}(\sigma^t, hza)$ and $b_i^t(h, z) = v_i^{t,L}(\sigma^t, hz)$, for all $h \in H \setminus H_{TS}$, $z \in Z(h)$, $a \in A(h)$, we have $\text{Var}_{h'} \left[\hat{v}_i^{t,H}(\sigma^t, h, z|h') | h' \supseteq h \right] = \text{Var}_{h'} \left[\hat{v}_i^{t,L}(\sigma^t, hz, a|h') | h' \supseteq hz \right] = 0$, for all $h \in H \setminus H_{TS}$, $z \in Z(h)$, $a \in A(h)$. This completes the proof of Theorem 5, according to the last paragraph of Part I in this section.

I PROOF OF PROPOSITION 2

We start from the definition:

$$\begin{aligned} \mathcal{L}_{R,i}^{t,H} &= \mathcal{L}(R_{i,\theta}^{t,H}) = \mathbb{E}_{(I, \hat{r}_i^{t',H}) \sim \tau_R^i} \left[\sum_{z \in Z(I)} (R_{i,\theta}^{t,H}(z|I) - \hat{r}_i^{t',H}(I, z))^2 \right] \\ &= \frac{1}{\text{norm}} \sum_{t'=1}^t \sum_{I \in \mathcal{I}_i} \sum_{k=1}^K x_{t'}^k(I) \left[\sum_{z \in Z(I)} (R_{i,\theta}^{t,H}(z|I) - \hat{r}_i^{t',H}(I, z))^2 \right] \end{aligned} \quad (65)$$

Here, $x_{t'}^k(I)$ denotes whether I is visited in the k -th sampled trajectory at iteration t' , and $\text{norm} = \sum_{t'=1}^t \sum_{I \in \mathcal{I}_i} \sum_{k=1}^K x_{t'}^k(I)$ serves as the normalizing factor.

Let $R_{i,*}^{t,H}$ denote a minimal point of $\mathcal{L}(R_{i,\theta}^{t,H})$. Utilizing the first-order necessary condition for optimality, we obtain: $\nabla \mathcal{L}(R_{i,*}^{t,H}) = 0$. Thus, for the (I, z) entry of $R_{i,*}^{t,H}$, we deduce:

$$\frac{\partial \mathcal{L}(R_{i,*}^{t,H})}{\partial R_{i,\theta}^{t,H}(z|I)} = \frac{2}{norm} \sum_{t'=1}^t \sum_{k=1}^K x_{t'}^k(I) (R_{i,*}^{t,H}(z|I) - \hat{r}_i^{t',H}(I, z)) = 0 \quad (66)$$

$$R_{i,*}^{t,H}(z|I) = \frac{1}{norm'} \sum_{t'=1}^t \sum_{k=1}^K x_{t'}^k(I) \hat{r}_i^{t',H}(I, z) = \frac{1}{norm'} \sum_{t'=1}^t \sum_{k=1}^K \hat{r}_i^{t',H}(I, z|h'_{t',k})$$

where $norm' = \sum_{t'=1}^t \sum_{k=1}^K x_{t'}^k(I)$ denotes the normalizing factor, which is a positive constant for a certain memory $\tau_{R,i}^i$, and $h'_{t',k}$ is the termination state of the k -th sampled trajectory at iteration t' . In the second line of Equation (66), the second equality is valid based on the definition of sampled counterfactual regret (Equation (26) and (27)), which assigns non-zero values exclusively to information sets along the sampled trajectory. Now, we consider the expectation of $R_{i,*}^{t,H}(z|I)$ on the set of sampled trajectories $\{h'_{t',k}\}$:

$$\begin{aligned} \mathbb{E}_{\{h'_{t',k}\}} [R_{i,*}^{t,H}(z|I)] &= \frac{1}{norm'} \sum_{t'=1}^t \sum_{k=1}^K \mathbb{E}_{h'_{t',k}} [\hat{r}_i^{t',H}(I, z|h'_{t',k})] \\ &= \frac{1}{norm'} \sum_{t'=1}^t \sum_{k=1}^K r_i^{t',H}(I, z) = C_1 R_i^{t,H}(z|I) \end{aligned} \quad (67)$$

where $C_1 = \frac{T}{K \times norm'}$ and the second equality holds due to Theorem 4. The second part of Proposition 2, i.e., $\mathbb{E}_{\{h'_{t',k}\}} [R_{i,*}^{t,L}(a|I, z)] = C_2 R_i^{t,L}(a|I, z)$, can be demonstrated similarly.

J PROOF OF PROPOSITION 3

According to the definition of $\mathcal{L}_{\bar{\sigma},i}^H$ in Equation (10) and following the same process as Equation (65) - (66), we can obtain:

$$\bar{\sigma}_{i,*}^{T,H}(z|I) = \frac{1}{norm'} \sum_{t=1}^T \sum_{k=1}^K x_t^k(I) \sigma_i^{t,H}(z|I) \quad (68)$$

According to the law of large numbers, as $|\tau_{\bar{\sigma},i}^{t,i}| \rightarrow \infty$ ($t \in \{1, \dots, T\}$), we have:

$$\bar{\sigma}_{i,*}^{T,H}(z|I) \rightarrow \frac{\sum_{t=1}^T \pi^{q^{t,3-i}}(I) \sigma_i^{t,H}(z|I)}{\sum_{t=1}^T \pi^{q^{t,3-i}}(I)} \quad (69)$$

The equation above is based on the fact that, at a certain iteration t , the samples for updating the strategy of player i are collected when $3-i$ is the traverser, so the probability to visit a certain information set I is $\pi^{q^{t,3-i}}(I)$. Ideally, to apply the law of large numbers, we should randomly select a single information set for each randomly-sampled trajectory and add its strategy distribution to the memory $\tau_{\bar{\sigma},i}^{t,i}$. This guarantees that occurrences of information sets within each iteration t are independent and identically distributed, as the sampling strategy within an iteration remains consistent. However, in practice (Algorithm 2), we gather strategy distributions of all information sets for the non-traverser along each sampled trajectory to enhance sample efficiency, which has been empirically proven to be effective.

To connect the convergence result in Equation (69) and the definition of $\bar{\sigma}_i^{T,H}$ in Equation (8), we need to show that $\forall I \in \mathcal{I}_i$, $t \in \{1, \dots, T-1\}$, $\frac{\pi^{q^{t,3-i}}(I)}{\pi^{q^{t+1,3-i}}(I)} = \frac{\pi_i^{\sigma^t}(I)}{\pi_i^{\sigma^{t+1}}(I)}$. According to the sampling scheme, $q_p^{t,3-i}$ is a uniformly random strategy when $p = 3-i$, and it is equal to σ_p^t when $p = i$. Therefore, we have:

$$\frac{\pi^{q^{t,3-i}}(I)}{\pi^{q^{t+1,3-i}}(I)} = \frac{\sum_{h \in I} \pi_{3-i}^{Unif}(h) \pi_i^{\sigma^t}(h)}{\sum_{h \in I} \pi_{3-i}^{Unif}(h) \pi_i^{\sigma^{t+1}}(h)} = \frac{\sum_{h \in I} \pi_i^{\sigma^t}(h)}{\sum_{h \in I} \pi_i^{\sigma^{t+1}}(h)} = \frac{\pi_i^{\sigma^t}(I)}{\pi_i^{\sigma^{t+1}}(I)} \quad (70)$$

1458 It is satisfied in our/usual game settings that $\pi_{3-i}^{Unif}(h)$ remains consistent for all $h \in I$. This
1459 is attributable to the fact that histories within a single information set have the same height, and
1460 player 3 – i consistently employs a uniformly random strategy. Similarly, we can deduce that
1461 $\bar{\sigma}_{i,*}^{T,L}(a|I, z) \rightarrow \bar{\sigma}_i^{T,L}(a|I, z)$ using the aforementioned procedure.
1462

1463 K PROOF OF PROPOSITION 4

1464
1465 First, we present a lemma concerning the sampled baseline values $\hat{b}^{t+1}(h|h')$, as defined in Equation
1466 (12). This definition closely resembles that of the sampled counterfactual values in Equation (27),
1467 with two key distinctions: (1) \hat{b}^{t+1} is replaced with b^t , as b^{t+1} is not yet available; and (2) q^{t+1} is
1468 substituted with q^t , enabling the reuse of trajectories sampled with q^t for updating b^{t+1} , thereby
1469 enhancing efficiency.

1470 **Lemma 13.** *For all $h \in H$, $z \in Z(h)$, $a \in A(h)$, we have:*

$$1471 \mathbb{E}_{h'} \left[\hat{b}^{t+1}(h|h') | h' \supseteq h \right] = v^{t+1,H}(\sigma^{t+1}, h) \quad (71)$$

1472
1473
1474 *Proof.* Given the similarity between \hat{b}^{t+1} and $\hat{v}^{t+1,H}$, we can follow the proof of Lemma 4 and 5 to
1475 justify the lemma here.

$$1476 \begin{aligned} & E_{h'} \left[\hat{b}^{t+1}(h, z|h') | h' \supseteq h \right] \\ 1477 &= E_{h'} \left[\frac{\delta(hz \sqsubseteq h')}{q^{t,1}(z|h)} \left[\hat{b}^{t+1}(hz|h') - b^t(h, z) \right] + b^t(h, z) | h' \supseteq h \right] \\ 1478 &= P(hz \sqsubseteq h' | h' \supseteq h) E_{h'} \left[\frac{1}{q^{t,1}(z|h)} \left[\hat{b}^{t+1}(hz|h') - b^t(h, z) \right] + b^t(h, z) | h' \supseteq hz \right] + \\ 1479 & \quad P(hz \not\sqsubseteq h' | h' \supseteq h) b^t(h, z) \\ 1480 &= q^{t,1}(z|h) \left[\frac{1}{q^{t,1}(z|h)} \left[E_{h'}(\hat{b}^{t+1}(hz|h') | h' \supseteq hz) - b^t(h, z) \right] + b^t(h, z) \right] + \\ 1481 & \quad (1 - q^{t,1}(z|h)) b^t(h, z) \\ 1482 &= E_{h'} \left[\hat{b}^{t+1}(hz|h') | h' \supseteq hz \right] \end{aligned} \quad (72)$$

1483
1484 According to Algorithm 2, the trajectories for updating \hat{b}^{t+1} are sampled at iteration t when
1485 player 1 is the traverser, so $P(hz \sqsubseteq h' | h' \supseteq h) = q^{t,1}(z|h)$. Similarly, we can obtain
1486 $E_{h'} \left[\hat{b}^{t+1}(hz, a|h') | h' \supseteq hz \right] = E_{h'} \left[\hat{b}^{t+1}(hza|h') | h' \supseteq hza \right]$.

1487
1488 Next, we can employ these two equations to perform induction based on the height of h within the
1489 game tree. If $h \in H_{TS}$, $E_{h'} \left[\hat{b}^{t+1}(h|h') | h' \supseteq h \right] = \hat{b}^{t+1}(h|h') = u_1(h) = v^{t+1,H}(\sigma^{t+1}, h)$ based
1490 on the definition. If $hza \in H_{TS}$, we have:

$$1491 \begin{aligned} & E_{h'} \left[\hat{b}^{t+1}(h, z|h') | h' \supseteq h \right] = E_{h'} \left[\hat{b}^{t+1}(hz|h') | h' \supseteq hz \right] \\ 1492 &= \sum_{a \in A(h)} \sigma_{P(h)}^{t+1,L}(a|h, z) E_{h'} \left[\hat{b}^{t+1}(hz, a|h') | h' \supseteq hz \right] \\ 1493 &= \sum_{a \in A(h)} \sigma_{P(h)}^{t+1,L}(a|h, z) E_{h'} \left[\hat{b}^{t+1}(hza|h') | h' \supseteq hza \right] \\ 1494 &= \sum_{a \in A(h)} \sigma_{P(h)}^{t+1,L}(a|h, z) v^{t+1,H}(\sigma^{t+1}, hza) = v^{t+1,L}(\sigma^{t+1}, hz) \end{aligned} \quad (73)$$

1495
1496 Here, we employ the induction hypothesis in the fourth equivalence, and incorporate pertinent
1497 definitions for the remaining equivalences. It follows that:

$$1498 \begin{aligned} & E_{h'} \left[\hat{b}^{t+1}(h|h') | h' \supseteq h \right] = \sum_{z \in Z(h)} \sigma_{P(h)}^{t+1,H}(z|h) E_{h'} \left[\hat{b}^{t+1}(h, z|h') | h' \supseteq h \right] \\ 1499 &= \sum_{z \in Z(h)} \sigma_{P(h)}^{t+1,H}(z|h) v^{t+1,L}(\sigma^{t+1}, hz) = v^{t+1,H}(\sigma^{t+1}, h) \end{aligned} \quad (74)$$

By repeating the two equations above, we can show that $E_{h'} [\hat{b}^{t+1}(h|h')|h' \supseteq h] = v^{t+1,H}(\sigma^{t+1}, h)$ holds for a general $h \notin H_{TS}$. \square

Next, we complete the proof of Proposition 4.

$$\begin{aligned} \mathcal{L}_b^{t+1} = \mathcal{L}(b^{t+1}) &= \mathbb{E}_{h' \sim \tau_b^t} \left[\sum_{hza \sqsubseteq h'} (b^{t+1}(h, z, a) - \hat{b}^{t+1}(hza|h'))^2 \right] \\ &= \frac{\sum_{h'} N(h') \sum_{hza \sqsubseteq h'} (b^{t+1}(h, z, a) - \hat{b}^{t+1}(hza|h'))^2}{\sum_{h'} N(h')} \end{aligned} \quad (75)$$

Here, $N(h')$ denotes the number of occurrences of h' in the memory τ_b^t . Let $b^{t+1,*}$ denote a minimal point of \mathcal{L}_b^{t+1} . Utilizing the first-order necessary condition for optimality, we obtain: $\nabla \mathcal{L}(b^{t+1,*}) = 0$. Thus, for the (h, z, a) entry of $b^{t+1,*}$, we deduce:

$$\begin{aligned} \frac{\partial \mathcal{L}(b^{t+1,*})}{\partial b^{t+1}(h, z, a)} &= \frac{2 \sum_{h' \supseteq hza} N(h') (b^{t+1,*}(h, z, a) - \hat{b}^{t+1}(hza|h'))}{\sum_{h'} N(h')} = 0 \\ b^{t+1,*}(h, z, a) &= \frac{\sum_{h' \supseteq hza} N(h') \hat{b}^{t+1}(hza|h')}{\sum_{h' \supseteq hza} N(h')} \end{aligned} \quad (76)$$

The trajectories in τ_b^t can be considered as a sequence of independent and identically distributed random variables, since they are independently sampled with the same sample strategy $q^{t,1}$. Then, according to the law of large numbers, as $|\tau_b^t| \rightarrow \infty$, we conclude:

$$\begin{aligned} b^{t+1,*}(h, z, a) &\rightarrow \frac{\sum_{h' \supseteq hza} \pi^{q^{t,1}}(h') \hat{b}^{t+1}(hza|h')}{\sum_{h' \supseteq hza} \pi^{q^{t,1}}(h')} \\ &= \mathbb{E}_{h'} [\hat{b}^{t+1}(hza|h') | h' \supseteq hza] = v^{t+1,H}(\sigma^{t+1}, hza) \end{aligned} \quad (77)$$

where the last equality comes from Lemma 13. It follows:

$$\begin{aligned} b^{t+1,*}(h, z) &= \sum_a \sigma_{P(h)}^{t+1,L}(a|I(h), z) b^{t+1,*}(h, z, a) \\ &\rightarrow \sum_a \sigma_{P(h)}^{t+1,L}(a|I(h), z) v^{t+1,H}(\sigma^{t+1}, hza) = v^{t+1,L}(\sigma^{t+1}, hz) \end{aligned} \quad (78)$$

L BENCHMARK INFORMATION

Table 3: Details of the Selected Benchmarks: For each benchmark, we provide its decision horizon, the number of nodes in the game tree, and the initial stack size for each player.

Benchmark	Leduc	Leduc_10	Leduc_15	Leduc_20	FHP	FHP_10
Stack Size	13	60	80	100	2000	4000
Horizon	4	20	30	40	8	20
# of Nodes	464	31814	67556	113954	2.58×10^{12}	3.17×10^{13}

Details of the benchmarks are summarized in Table 3.

M ABLATION STUDY RESULTS

HDCFR integrates the one-step option framework (Section 2.2) and variance-reduced Monte Carlo CFR (Section 3.2 and Appendix E). This section offers an ablation analysis highlighting each crucial

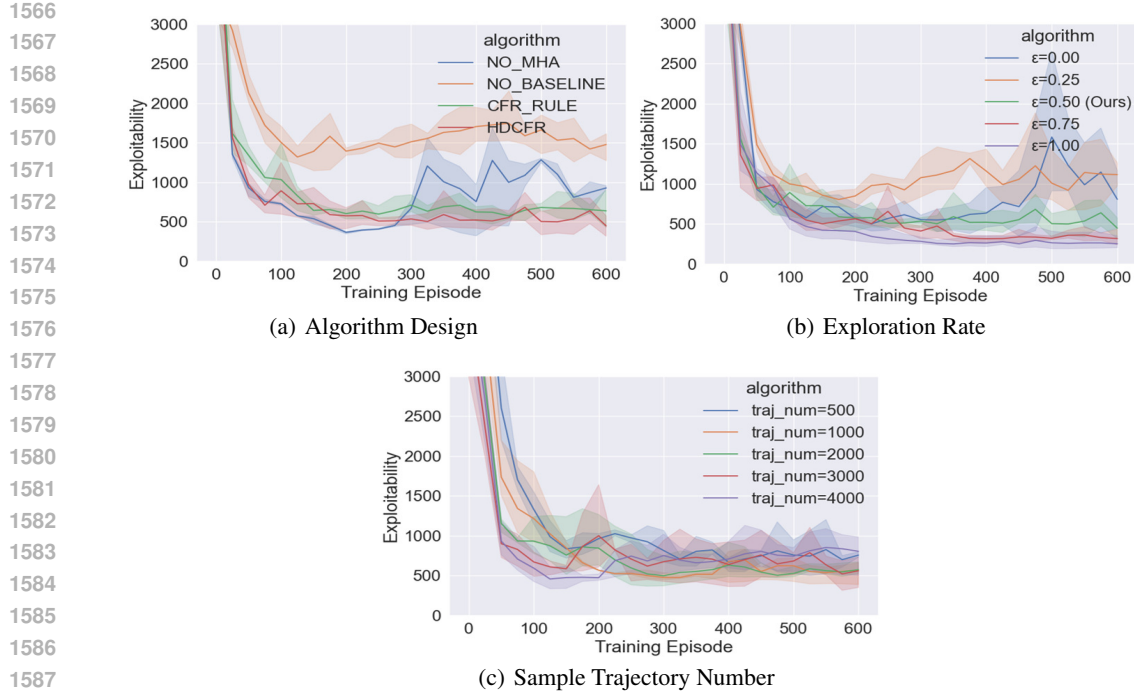


Figure 3: Learning curves of different ablations on Leduc_20. (a) Without the MHA component in the high-level strategy (NO_MHA) or the baseline function for variance reduction (NO_BASELINE), the learning performance degrades significantly. Following the CFR rule (Equation (7)) results in slightly slower convergence. (b) Increased randomness in the traverser’s sample strategy enhances learning. (c) More sampled trajectories in each training episode boost initial convergence speed without affecting the convergent performance.

element of our algorithm: the option framework, variance reduction, Monte Carlo sampling, and CFR.

The option framework: The key component of the one-step option framework is the Multi-Head Attention (MHA) mechanism which enables the agent to temporarily extend skills and so form a hierarchical policy in the learning process. Without this component in the high-level strategy (i.e., NO_MHA in Figure 3(a)), the agent struggles to converge at the final stage, akin to the behavior observed for DREAM in Figure 1(d).

Variance reduction: In HDCFR, we incorporate a baseline function to reduce variance. This function proves pivotal for extended-horizon tasks where sampling variance can escalate. Excluding the baseline function from the learning process, as marked by NO_BASELINE in Figure 3(a), results in a substantial performance decline.

Monte Carlo sampling: During Monte Carlo sampling, as outlined in Section 3.3, the traverser should use a uniformly random sampling strategy. Yet, for fair comparisons with DREAM, we employ a weighted average of a uniformly random strategy (with the weight ϵ) and the traverser’s current strategy (σ_i^t). The controlling weight ϵ for HDCFR in the other experiments is set as 0.5, aligning with the configuration for the baselines. Figure 3(b) indicates that as ϵ increases, approximately there is a correlating rise in learning performance. Notably, our original design, i.e., setting $\epsilon = 1$, delivers the best result, amplifying the performance depicted in Figure 1(d). Another key aspect of Monte Carlo sampling is the number of sampled trajectories per training episode (i.e., the hyperparameter K in Algorithm 1). According to Figure 3(c), increasing this count accelerates convergence during the initial training phase. However, it does not necessarily improve the final convergent performance and instead proportionally increases the overall training time.

CFR: As indicated by Brown et al. (2019) and Steinberger et al. (2020), slightly modifying the CFR updating rule (Equation (7)), that is, to greedily select the action with the largest regret rather than

1620 use a random one when the sum $\mu^H, \mu^L \leq 0$, can speed up the convergence. We adopt the same trick
 1621 and find that it can improve the convergence speed slightly, as compared to the original setting (i.e.,
 1622 CFR_RULE in Figure 3(a)).
 1623

1624 N PSEUDO CODE OF HDCFR

1625 **Algorithm 1** Hierarchical Deep Counterfactual Regret Minimization (HDCFR)

1626
 1627
 1628 1: **Initialize** the counterfactual regret networks $R_{i,\theta}^{0,H}, R_{i,\theta}^{0,L}, \forall i \in \{1, 2\}$ (collectively denoted as
 1629 R_{θ}^0), and the baseline network b^1 , so that they return 0 for all inputs
 1630 2: **Initialize** the average strategy networks $\bar{\sigma}_{i,\phi}^{T,H}, \bar{\sigma}_{i,\phi}^{T,L}, \forall i \in \{1, 2\}$ with random parameters
 1631 3: **Initialize** the replay buffer for the counterfactual regrets and average strategies, i.e., $\tau_R^i, \tau_{\bar{\sigma}}^i, \forall i \in$
 1632 $\{1, 2\}$ as empty sets
 1633 4: **for** $t = \{1, \dots, T\}$ **do**
 1634 5: **Initialize** the the replay buffer for the baseline function at iteration t : $\tau_b^t = \emptyset$
 1635 6: **for** $i = \{1, 2\}$ **do**
 1636 7: **Define** the sample strategy profile at t with i being the traverser, i.e., $q^{t,i}$
 1637 8: **for** traversal $k = \{1, \dots, K\}$ **do**
 1638 9: $HighRollout(\emptyset, R_{\theta}^{t-1}, \tau_R^i, \tau_{\bar{\sigma}}^{3-i}, \tau_b^t, q^{t,i}, b^t)$
 1639 10: **end for**
 1640 11: **end for**
 1641 12: **for** $i = \{1, 2\}$ **do**
 1642 13: **Train** $R_{i,\theta}^{t,H}, R_{i,\theta}^{t,L}$ from scratch by minimizing Equation (9)
 1643 14: **end for**
 1644 15: $b^{t+1} = BaselineTraining(b^t, \tau_b^t, R_{\theta}^t, q^{t,1})$
 1645 16: **end for**
 1646 17: **for** $i = \{1, 2\}$ **do**
 1647 18: **Obtain** $\bar{\sigma}_{i,\phi}^{T,H}, \bar{\sigma}_{i,\phi}^{T,L}$ by minimizing Equation (10)
 1648 19: **end for**
 1649 20: **Return** $\{(\bar{\sigma}_{1,\phi}^{T,H}, \bar{\sigma}_{1,\phi}^{T,L}), (\bar{\sigma}_{2,\phi}^{T,H}, \bar{\sigma}_{2,\phi}^{T,L})\}$, i.e., the approximate Nash Equilibrium hierarchical
 1650 strategy profile
 1651 21:
 1652 22: **function** $BaselineTraining(b^t, \tau_b^t, R_{\theta}^t, q^{t,1})$
 1653 23: **for** h' in τ_b^t **do**
 1654 24: **for** $hza \sqsubseteq h'$ **do** (tracing back from h' to its initial state)
 1655 25: **Compute** $\hat{b}^{t+1}(hza|h')$ using b^t, R_{θ}^t , and $q^{t,1}$, following Equation (12), where
 1656 R_{θ}^t indicates σ^{t+1} according to Equation (7)
 1657 26: **end for**
 1658 27: **end for**
 1659 28: **Train** b^{t+1} by minimizing Equation (11)
 1660 29: **Return** b^{t+1}
 1661 30: **end function**

1662
 1663
 1664
 1665
 1666
 1667
 1668
 1669
 1670
 1671
 1672
 1673

1674
1675
1676
1677
1678
1679
1680
1681
1682
1683
1684
1685
1686
1687
1688
1689
1690
1691
1692
1693
1694
1695
1696
1697
1698
1699
1700
1701
1702
1703
1704
1705
1706
1707
1708
1709
1710
1711
1712
1713
1714
1715
1716
1717
1718
1719
1720
1721
1722
1723
1724
1725
1726
1727

Algorithm 2 Hierarchical Deep Counterfactual Regret Minimization (HDCFR) Continued

```

1: function HighRollout( $h, R_\theta^{t-1}, \tau_R^i, \tau_{\bar{\sigma}}^{3-i}, \tau_b^t, q^{t,i}, b^t$ )
2:   if  $h \in H_{TS}$  then
3:     Assign  $h' = h$ 
4:     if  $i == 1$  then
5:       Add  $h'$  to  $\tau_b^t$ 
6:     end if
7:     Return  $u_1(h')$ 
8:   end if
9:    $I = I(h), p = P(h)$ 
10:  Sample an option  $z \sim q^{t,i}(\cdot|h)$ 
11:   $\hat{v}^{t,L}(\sigma^t, hz|h') = \text{LowRollout}(h, z, R_\theta^{t-1}, \tau_R^i, \tau_{\bar{\sigma}}^{3-i}, \tau_b^t, q^{t,i}, b^t)$ 
12:   $\hat{v}^{t,H}(\sigma^t, h, z'|h') = b^t(h, z'), \forall z' \neq z$ 
13:   $\hat{v}^{t,H}(\sigma^t, h, z|h') = \frac{1}{q^{t,i}(z|h)} [\hat{v}^{t,L}(\sigma^t, hz|h') - b^t(h, z)] + b^t(h, z)$ 
14:   $\hat{v}^{t,H}(\sigma^t, h|h') = \sum_{z \in Z(h)} \sigma_p^{t,H}(z|h) \hat{v}^{t,H}(\sigma^t, h, z|h')$ 
15:  if  $p == i$  then
16:     $\hat{r}_i^{t,H}(I, \cdot|h') = (-1)^{i+1} \frac{\pi_{3-i}^{\sigma^t}(h)}{\pi^{q^{t,i}}(h)} [\hat{v}^{t,H}(\sigma^t, h, \cdot|h') - \hat{v}^{t,H}(\sigma^t, h|h')]$ 
17:    Add  $(I, t, \hat{r}_i^{t,H}(I, \cdot|h'))$  to  $\tau_R^i$ 
18:  else if  $p == 3 - i$  then
19:    Compute  $\sigma_{3-i}^{t,H}(\cdot|I)$  based on  $R_{3-i}^{t-1,H}(\cdot|I)$  following Equation (7)
20:    Add  $(I, t, \sigma_{3-i}^{t,H}(\cdot|I))$  to  $\tau_{\bar{\sigma}}^{3-i}$ 
21:  end if
22:  Return  $\hat{v}^{t,H}(\sigma^t, h|h')$ 
23: end function
24:
25: function LowRollout( $h, z, R_\theta^{t-1}, \tau_R^i, \tau_{\bar{\sigma}}^{3-i}, \tau_b^t, q^{t,i}, b^t$ )
26:   $I = I(h), p = P(h)$ 
27:  Sample an action  $a \sim q^{t,i}(\cdot|h, z)$ 
28:   $\hat{v}^{t,H}(\sigma^t, hza|h') = \text{HighRollout}(hza, R_\theta^{t-1}, \tau_R^i, \tau_{\bar{\sigma}}^{3-i}, \tau_b^t, q^{t,i}, b^t)$ 
29:   $\hat{v}^{t,L}(\sigma^t, hz, a'|h') = b^t(h, z, a'), \forall a' \neq a$ 
30:   $\hat{v}^{t,L}(\sigma^t, hz, a|h') = \frac{1}{q^{t,i}(a|h,z)} [\hat{v}^{t,H}(\sigma^t, hza|h') - b^t(h, z, a)] + b^t(h, z, a)$ 
31:   $\hat{v}^{t,L}(\sigma^t, hz|h') = \sum_{a \in A(h)} \sigma_p^{t,L}(a|h, z) \hat{v}^{t,L}(\sigma^t, hz, a|h')$ 
32:  if  $p == i$  then
33:     $\hat{r}_i^{t,L}(Iz, \cdot|h') = (-1)^{i+1} \frac{\pi_{3-i}^{\sigma^t}(h)}{\pi^{q^{t,i}}(hz)} [\hat{v}^{t,L}(\sigma^t, hz, \cdot|h') - \hat{v}^{t,L}(\sigma^t, hz|h')]$ 
34:    Add  $(Iz, t, \hat{r}_i^{t,L}(Iz, \cdot|h'))$  to  $\tau_R^i$ 
35:  else if  $p == 3 - i$  then
36:    Compute  $\sigma_{3-i}^{t,L}(\cdot|I, z)$  based on  $R_{3-i}^{t-1,L}(\cdot|I, z)$  following Equation (7)
37:    Add  $(Iz, t, \sigma_{3-i}^{t,L}(\cdot|I, z))$  to  $\tau_{\bar{\sigma}}^{3-i}$ 
38:  end if
39:  Return  $\hat{v}^{t,L}(\sigma^t, hz|h')$ 
40: end function

```

1728
1729
1730
1731
1732
1733
1734
1735
1736
1737
1738
1739
1740
1741
1742
1743
1744
1745
1746
1747
1748
1749
1750
1751
1752
1753
1754
1755
1756
1757
1758
1759
1760
1761
1762
1763
1764
1765
1766
1767
1768
1769
1770
1771
1772
1773
1774
1775
1776
1777
1778
1779
1780
1781

O THE USE OF LARGE LANGUAGE MODELS

Large Language Models (LLMs) were *not* used for the research idea, related-work selection, algorithm or experiment design, proofs or derivations, figure/table generation, or analysis/interpretation of results. We used LLMs only for lightly improving grammar and phrasing on a small number of sentences. No technical content was generated by LLMs; All suggestions were reviewed and, when appropriate, rewritten by the authors. No confidential or proprietary data were provided to LLMs. The authors take full responsibility for the contents of this paper, and LLMs are not authors.