

# Approximating Nash Equilibria in General-Sum Games via Meta-Learning

David Sychrovský  
Charles University & EquiLibre  
Technologies  
Prague, Czechia  
sychrovsky@kam.mff.cuni.cz

Christopher Solinas  
University of Alberta  
Edmonton, Canada  
solinas@ualberta.ca

Revan MacQueen  
Alberta Machine Intelligence Institute  
Edmonton, Canada  
revan.macqueen@amii.ca

Kevin Wang  
Brown University  
Providence, United States  
kevin\_a\_wang@brown.edu

James R. Wright  
University of Alberta  
Edmonton, Canada  
james.wright@ualberta.ca

Nathan R. Sturtevant  
University of Alberta  
Edmonton, Canada  
nathanst@ualberta.ca

Michael Bowling  
University of Alberta  
Edmonton, Canada  
mbowling@ualberta.ca

## ABSTRACT

Nash equilibrium is perhaps the best-known solution concept in game theory. Such a solution assigns a strategy to each player which offers no incentive to unilaterally deviate. While a Nash equilibrium is guaranteed to always exist, the problem of finding one in general-sum games is PPAD-complete, generally considered intractable. Regret minimization is an efficient framework for approximating Nash equilibria in two-player zero-sum games. However, in general-sum games, such algorithms are only guaranteed to converge to a coarse-correlated equilibrium (CCE), a solution concept where players can correlate their strategies. In this work, we use meta-learning to minimize the correlations in strategies produced by a regret minimizer. This encourages the regret minimizer to find strategies that are closer to a Nash equilibrium. The meta-learned regret minimizer is still guaranteed to converge to a CCE, but we give a bound on the distance to Nash equilibrium in terms of our meta-loss. We evaluate our approach in general-sum imperfect information games. Our algorithms provide significantly better approximations of Nash equilibria than state-of-the-art regret minimization techniques.

## KEYWORDS

regret minimization; meta-learning; Nash equilibria

### ACM Reference Format:

David Sychrovský, Christopher Solinas, Revan MacQueen, Kevin Wang, James R. Wright, Nathan R. Sturtevant, and Michael Bowling. 2026. Approximating Nash Equilibria in General-Sum Games via Meta-Learning. In

*Proc. of the 25th International Conference on Autonomous Agents and Multiagent Systems (AAMAS 2026), Paphos, Cyprus, May 25 – 29, 2026*, IFAAMAS, 9 pages. <https://doi.org/10.65109/>

## 1 INTRODUCTION

The Nash equilibrium is one of the most influential solution concepts in game theory. A strategy profile is a Nash equilibrium if it has the guarantee that no player can benefit by unilaterally deviating from it. The robustness of this guarantee means that Nash equilibria have applications in many domains ranging from economics [33, 47] to machine learning [22]. Finding an efficient algorithm for computing Nash equilibria has attracted much attention [11, 25, 31, 34, 41]. However, it was shown that, in its full generality, finding a Nash equilibrium is PPAD-complete [13, 39]. Many related decision problems, such as ‘Is a given action in the support of a Nash equilibrium?’, are NP-complete [21].

Despite these negative results, computing Nash equilibria in special classes of games, in particular two-player zero-sum games, is tractable. In this setting, regret minimization has become the dominant approach for finding Nash equilibria [38]. This framework casts each player as an independent online learner who repeatedly interacts with the game, selecting strategies according to dynamics that lead to sublinear growth of their accumulated *regret*. Regret minimizers guarantee convergence to Nash equilibria in two-player zero-sum games, and are the basis for many significant results in imperfect information games [4, 5, 7, 10, 35, 42].

Outside the two-player zero-sum setting, regret minimization algorithms are no longer guaranteed to converge to a Nash equilibrium. Instead, a regret minimizer’s empirical distribution of play converges to a coarse-correlated equilibrium (CCE) [23, 24]. The CCE is a relaxed equilibrium concept, which gives a distribution over the *outcomes* of the game such that it isn’t beneficial for any player to deviate from it. If this distribution is uncorrelated, meaning it can be expressed as a profile of independent strategies, it is also a Nash equilibrium. As such, Nash equilibria form a subset of CCEs, for which the outcome distribution can be marginalized



This work is licensed under a Creative Commons Attribution International 4.0 License.

*Proc. of the 25th International Conference on Autonomous Agents and Multiagent Systems (AAMAS 2026)*, C. Amato, L. Dennis, V. Mascardi, J. Thangarajah (eds.), May 25 – 29, 2026, Paphos, Cyprus. © 2026 International Foundation for Autonomous Agents and Multiagent Systems ([www.ifaamas.org](http://www.ifaamas.org)). <https://doi.org/10.65109/>

into strategies of the individual players. The degree to which a CCE is correlated, or how much a player can infer about the actions of other players given their action, can be formalized by total correlation [48].

A recently proposed *learning not to regret* framework allows one to meta-learn a regret minimizer to optimize a specified objective, while keeping regret minimization guarantees [44]. Their goal was to accelerate the empirical convergence rate on a distribution of black-box tasks. In this work, we meta-learn predictions that optimize an alternative meta-objective: minimizing correlation in the players’ strategies. The resulting algorithm is still guaranteed to converge to a CCE, and is meta-learned to empirically converge to a Nash equilibrium on a distribution of interest. If the support of the distribution doesn’t include all general-sum games, the problem of finding Nash may be tractable even if  $P \neq \text{PPAD}$ . We further show this approach is sound by providing a bound on the distance to a Nash equilibrium in terms of our meta-objective. We evaluate our approach in general-sum imperfect information games. Our algorithms provide significantly better approximations of Nash equilibria than state-of-the-art regret minimization techniques.

## 1.1 Related Work

The Nash equilibrium is one of the oldest solution concepts in game theory. Thanks to its many appealing properties, developing efficient algorithms for approximating Nash equilibria has seen much attention [3, 14–16, 26, 30]. Furthermore, it was shown that, unless  $P = \text{NP}$ , polynomial algorithms for finding all Nash equilibria cannot exist [21]. This negative result suggests that there are games for which finding a Nash equilibrium requires enumerating all possible strategies — an amount exponential in the number of actions.

The Lemke-Howson algorithm [29] is one such algorithm, which provably finds a Nash equilibrium of two-player general-sum games in normal-form. It works by constructing a path on an abstract polyhedron, which is guaranteed to terminate at the Nash equilibrium. Similar to the simplex method [37], the path may be exponentially long in some games. However, such games are empirically rare [12]. Several modifications of the Lemke-Howson algorithm were proposed to improve its empirical performance [12, 19]. However, the algorithm cannot work with games in extensive-form. When converted to normal-form, the size of the game increases exponentially, making these algorithms scale very poorly.

Regret minimization is a powerful framework for online convex optimization [38, 51], with regret matching as one of the most popular algorithms in game applications [24]. Counterfactual regret minimization enables the use of regret matching in sequential decision-making, by decomposing the full regret to individual states [52]. In two-player zero-sum games, regret minimization algorithms are guaranteed to converge to a Nash equilibrium. Many prior works explored modifications of regret matching to speed up its empirical performance in two-player zero-sum games, such as CFR<sup>+</sup> [45], Linear CFR [6], PCFR<sup>+</sup> [17], Discounted CFR [8], and their hyperparameter-scheduled counterparts [50].

Despite the lack of theoretical guarantees in general-sum games, regret minimization algorithms empirically converge close to Nash equilibria on many standard benchmarks [10, 20, 40]. Recently,

some theoretical advancements have been made to understand this empirical performance. If the game has a special ‘pair-wise zero-sum’ structure, then the regret minimizers are guaranteed to find a Nash equilibrium [11]. Moreover, if a game is ‘close’ to such ‘pair-wise zero-sum’ games, the regret minimizers converge ‘close’ to a Nash equilibrium [32].

A recently introduced extension of regret matching, predictive regret matching [18], forms a continuous class of algorithms with regret minimization guarantees. Subsequently, [44] introduced the ‘learning not to regret’ framework—a way to meta-learn the predictions while keeping regret minimization guarantees.

## 1.2 Main Contribution

In this work, we extend the *learning not to regret* framework to encourage convergence to Nash equilibria in general-sum games. Our approach penalizes correlations in the average empirical strategy profile found by the regret minimizer. While our meta-learned algorithms do not guarantee convergence to a Nash equilibrium, we find that our algorithms empirically converge to CCEs with low correlations in the players’ strategies, and provide significantly better approximations of Nash equilibria than prior regret minimization algorithms.

We demonstrate the feasibility of our approach by conducting experiments in multiplayer general-sum games. We start with a distribution of normal-form games, where prior regret minimization algorithms overwhelmingly converge to a strictly correlated CCE. Next, we shift our attention to Leduc poker, a standard extensive-form imperfect information benchmark. We show that, after a small modification of the rules (to make the game general-sum), prior regret minimizers no longer reliably converge to a Nash equilibrium. When trained on this distribution, our meta-learning framework produces a regret minimizer that reach significantly closer to a Nash equilibrium. Finally, we demonstrate that our framework can even be used to obtain better approximations of a Nash equilibrium on a single general-sum game rather than just a family of games. We choose the three-player Leduc poker, obtaining, to our best knowledge, the closest approximation of a Nash equilibrium of this game.

## 2 PRELIMINARIES

We briefly introduce the formalism of incomplete information games we will use. Next, we describe regret minimization, a general online convex optimization framework. Finally, we discuss how regret minimization can be used to find equilibria of these games.

### 2.1 Games

We work within a formalism based on factored-observation stochastic games [27] with terminal utilities.

**Definition 1.** A game is a tuple  $\langle \mathcal{N}, \mathcal{W}, w^0, \mathcal{A}, \mathcal{T}, u, O \rangle$ , where

- $\mathcal{N} = \{1, \dots, n\}$  is a **player set**. We use symbol  $i$  for a player and  $-i$  for its opponents.
- $\mathcal{W}$  is a set of **world states** and  $w^0 \in \mathcal{W}$  is a unique initial world state.
- $\mathcal{A} = \mathcal{A}_1 \times \dots \times \mathcal{A}_n$  is a space of **joint actions**. A world state with no legal actions is **terminal**. We denote the set of terminal world states as  $\mathcal{Z}$ .

- After taking a (legal) joint action  $a$  at  $w$ , the **transition function**  $\mathcal{T}$  determines the next world state  $w'$ , drawn from the probability distribution  $\mathcal{T}(w, a) \in \Delta(\mathcal{W})$ .
- $u_i(z)$  is the **utility** player  $i$  receives when a terminal state  $z \in \mathcal{Z}$  is reached.
- $O = (O_1, \dots, O_n)$  is the **observation function** specifying both the private and public observation that players receive upon the state transition.

The space  $\mathcal{S}_i$  of all action-observation sequences can be viewed as the infostate tree of player  $i$ . A **strategy profile** is a tuple  $\sigma = (\sigma_1, \dots, \sigma_n)$ , where each player's **strategy**  $\sigma_i : s_i \in \mathcal{S}_i \mapsto \sigma_i(s_i) \in \Delta^{|\mathcal{A}_i(s_i)|}$  specifies the probability distribution from which player  $i$  draws their next action conditional on having information  $s_i$ . We denote the space of all strategy profiles as  $\Sigma$ . A **pure strategy**  $\rho_i$  is a deterministic strategy: i.e.  $\sigma_i(s_i, a_i) = 1$  for some  $a_i \in \mathcal{A}_i(s_i)$ . A selection of pure strategies for all players  $\rho = (\rho_1, \dots, \rho_n)$  is a **pure strategy profile** and the set of all pure strategy profiles is  $\mathbf{P}$ .

Let  $\Delta(X)$  denote the set of distributions over a domain  $X$ . A **joint strategy profile**  $\delta \in \Delta(\mathbf{P})$  is a distribution over pure strategy profiles. As such, every strategy profile is also a joint strategy profile. However, the opposite is not true in general: only *some* joint strategy profiles are “marginalizable” into an equivalent strategy profile, while those with correlations between players’ strategies are not.

The expected **utility** under a joint strategy profile  $\delta$  is  $u_i(\delta) = \mathbb{E}_{z \sim \delta} u_i(z)$ , where the expectation is over the terminal states  $z \in \mathcal{Z}$  and their reach probability under  $\delta$ . The **best-response** to the joint strategy of the other players is  $br(\delta_{-i}) \in \arg \max_{\sigma_i} u_i(\sigma_i, \delta_{-i})$ , where  $\delta_{-i}(\rho_{-i}) = \sum_{\rho_i \in \mathcal{A}_i} \delta(\rho_i, \rho_{-i})$ .

We may measure the distance of a strategy profile  $\sigma$  from a Nash equilibrium by its **NashGap**: the maximum gain any player can obtain by unilaterally deviating from  $\sigma$

$$\text{NashGap}(\sigma) = \max_{i \in N} [u_i(br(\sigma_{-i}), \sigma_i) - u_i(\sigma)].$$

A strategy profile is a Nash equilibrium if its NashGap is zero.<sup>1</sup>

The coarse correlated equilibrium (CCE) [36, 38] is a generalization of Nash equilibrium to joint strategy profiles that allows for correlation between players’ strategies. A CCE is a joint strategy profile such that any unilateral deviation by any player doesn’t increase that player’s utility, while other players continue to play according to the joint strategy. We define the **CCE Gap** as

$$\text{CCE Gap}(\delta) = \max_{i \in N} [u_i(br(\delta_{-i}), \delta_i) - u_i(\delta)].$$

A joint strategy profile  $\delta$  is a CCE if and only if its CCE Gap is non-negative. If a joint strategy profile has zero CCE Gap, and can be written in terms of its marginal strategies for each player  $\delta = (\sigma_1, \dots, \sigma_n)$ , then its marginals  $\sigma_i$  are a Nash equilibrium. In general, CCEs do not admit this player-wise decomposition of the joint strategy profile—see Section 4.1 for an example.

## 2.2 Regret Minimization

An **online algorithm**  $m$  for the regret minimization task repeatedly interacts with an **environment** through available actions  $\mathcal{A}_i$ . The goal of a regret minimization algorithm is to maximize its hindsight performance (i.e., to minimize regret). For reasons discussed in the

<sup>1</sup>This is because then the individual strategy profiles are mutual best-responses.

---

### Algorithm 1: Neural Predictive Regret Matching [44]

---

```

1  $R^0 \leftarrow \mathbf{0} \in \mathbb{R}^{|\mathcal{A}|}$ ,  $\mathbf{x}^0 \leftarrow \mathbf{0} \in \mathbb{R}^{|\mathcal{A}|}$ 
2  $\mathbf{e}_s \leftarrow$  embedding of state  $s$ 
3 function NEXTSTRATEGY()
4    $\xi^t \leftarrow [R^{t-1} + \mathbf{p}^t]^+$ 
5   if  $\|\xi^t\|_1 > 0$ 
6     return  $\sigma^t \leftarrow \xi^t / \|\xi^t\|_1$ 
7   return  $\sigma^t \leftarrow$  arbitrary point in  $\Delta^{|\mathcal{A}|}$ 
8 function OBSERVE REWARD( $\mathbf{x}^t, \mathbf{e}_s$ )
9    $R^t \leftarrow R^{t-1} + r(\sigma^t, \mathbf{x}^t)$ 
10   $\mathbf{p}^{t+1} \leftarrow \alpha(r(\sigma^t, \mathbf{x}^t) + \pi(r(\sigma^t, \mathbf{x}^t), R^t, \mathbf{e}_s, |\theta|))$ 

```

---

following section, we will describe the formalism from the point of view of player  $i$  acting at an infostate  $s \in \mathcal{S}_i$ .

Formally, at each step  $t \leq T$ , the algorithm submits a **strategy**  $\sigma_i^t(s) \in \Delta^{|\mathcal{A}_i(s)|}$ . Subsequently, it observes the expected **reward**  $\mathbf{x}_i^t \in \mathbb{R}^{|\mathcal{A}_i(s)|}$  at the state  $s$  for each of the actions from the environment, which depends on the strategy in the rest of the game. The difference in reward obtained under  $\sigma_i^t(s)$  and any fixed action strategy is called the **instantaneous regret**  $r_i(\sigma^t, s) = \mathbf{x}_i^t(\sigma^t) - \langle \sigma_i^t(s), \mathbf{x}_i^t(\sigma^t) \rangle \mathbf{1}$ . The **cumulative regret** throughout time  $t$  is  $R_i^t(s) = \sum_{\tau=1}^t r_i(\sigma^\tau, s)$ .

The goal of a regret minimization algorithm is to ensure that the regret grows sublinearly for any sequence of rewards. One way to do that is for  $m$  to select  $\sigma_i^{t+1}(s)$  proportionally to the positive parts of  $R_i^t(s)$ , known as regret matching [2].

## 2.3 Connection Between Games and Regret Minimization

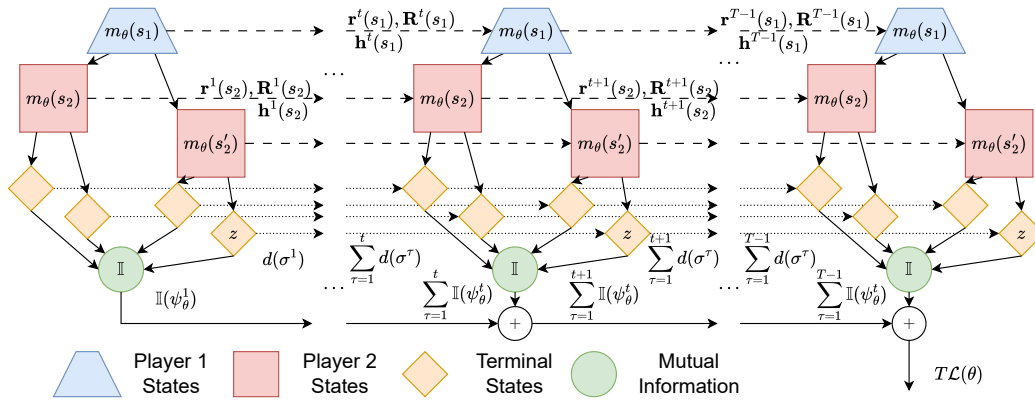
In normal-form games, or when  $\mathcal{S}_i$  is a singleton, if the **external regret**  $R_i^{\text{ext}, T} = \max_{a \in \mathcal{A}_i} R_i^T(a)$  grows as  $\mathcal{O}(\sqrt{T})$  for all players, then the empirical average joint strategy profile  $\bar{\delta}^T \stackrel{\text{def}}{=} \frac{1}{T} \sum_{t=1}^T \sigma_1^t \times \dots \times \sigma_n^t$  converges to a CCE as  $\mathcal{O}(1/\sqrt{T})$  [38].

In extensive-form games, in order to obtain the external regret, we would need to convert the game to normal-form. However, the size of the normal-form representation is exponential in the size extensive-form representation. Thankfully, one can upper-bound the normal-form regret by individual (i.e. per-infostate) **counterfactual regrets** [52]

$$\sum_{i \in N} R_i^{\text{ext}, T} \leq \sum_{i \in N} \sum_{s \in \mathcal{S}_i} \max \{ \|\mathbf{R}_i^T(s)\|_\infty, 0 \}.$$

The counterfactual regret is defined with respect to the **counterfactual reward**. At an infostate  $s \in \mathcal{S}_i$ , the counterfactual rewards measure the expected utility the player would obtain in the game when playing to reach  $s$ . In other words, it is the expected utility of  $i$  at  $s$ , multiplied by the opponent’s and chance’s contribution to the probability of reaching  $s$ . We can treat each infostate as a separate environment, and minimize their counterfactual regrets independently. This approach converges to a CCE [52].

In two-player zero-sum games, the empirical average strategy  $\bar{\sigma}$  is guaranteed to converge to a Nash equilibrium [52]. In fact, any



**Figure 1: Computational graph of NPCFR<sup>(+)</sup> for a simple extensive form game.** The algorithm  $m_\theta$  produces a strategy in each infostate using the regret  $r^t, R^t$ , and its hidden state  $h^t$ , see Algorithm 1. Each terminal state  $z \in \mathcal{Z}$  accumulates its empirical average reach probability  $\frac{1}{t} \sum_{\tau=1}^t d(\sigma^\tau)(z)$ . Marginalizability  $\mathbb{I}$  is computed between this accumulated average reach and the reach probability under the empirical average strategy profile in the game tree. The meta-loss is the average mutual information experienced over  $T$  steps, according to (1). Its gradient is propagated through all edges.

CCE of a two-player zero-sum game is guaranteed to be marginalizable [38]. Intuitively, any correlations will be beneficial for one of the players, which makes it irrational for the opponent to follow it.

### 3 META-LEARNING FRAMEWORK

We aim to find a regret minimization algorithm  $m_\theta$  with some parameterization  $\theta$  which tends to converge close to a Nash equilibrium on a distribution of games  $G$ . In this section, we describe the regret minimization algorithm and formalize our meta-learning objective.

#### 3.1 Neural Predictive Counterfactual Regret Minimization (NPCFR)

We work in the learning not to regret framework [44], which is built on the predictive regret matching (PRM) [18]. PRM is an extension of regret matching [24] which additionally uses a predictor about future reward. PRM provably enjoys  $\mathcal{O}(\sqrt{T})$  bound on the external regret for arbitrary bounded predictions [18].

Neural predictive regret matching is an extension of PRM which uses a predictor  $\pi$ , parameterized by a neural network  $\theta$  [44]; see Algorithm 1. At each step  $t$  and each infostate  $s \in \mathcal{S}_i, i \in \mathcal{N}$ , the predictor  $\pi(\cdot|\theta)$  makes a prediction about the next observed regret  $r^{t+1}$ . This prediction is then used when selecting the strategy, as if that regret was in fact observed. The strategy is then selected as if this predicted regret was observed. Network parameters  $\theta$  are shared across all infostates  $s \in \mathcal{S}_i, i \in \mathcal{N}$ , and  $\alpha \in \mathbb{R}$  is a hyperparameter, see Appendix B in the full version for more details. The  $e_s$  denotes some embedding of the infostate  $s$ ; see Section 4.

Since we make the predictions bounded, the predictor can be meta-learned to minimize a desired objective while maintaining the regret minimization guarantees [44], which makes the algorithm converge to a CCE. We use a novel meta-objective, which is introduced in the following section, to encourage the algorithm to

converge to a Nash equilibrium. Applying the algorithm to counterfactual regrets at each infostate allows us to use it on extensive-form games. This setup is referred to as neural predictive counterfactual regret minimization (NPCFR).

#### 3.2 Meta-Loss Function

Any instance of NPCFR is a regret minimizer and is therefore guaranteed to converge to a CCE. Since any Nash equilibrium is a CCE for which player strategies are uncorrelated, we propose a meta-loss objective that penalizes correlation in the CCE found by NPCFR. Informally, these correlations measure the mutual dependence of players' strategies. Or in other words, how much a player can infer about the actions of other players given their action.

One could express this measure of correlation as the *mutual information* of the CCE.<sup>2</sup> However, for extensive-form games, this leads to an exponential blow-up in the size of the game, since there are exponentially more pure strategies than infostates. Instead, we exploit the structure of extensive-form games to define an equivalent meta-loss that does not suffer from this blow-up.

Formally, let  $\psi^T = (\sigma^t)_{t=1}^T$  be a sequence of strategy profiles selected by a regret minimizer. Let  $d(\sigma)$  be the distribution of reach probabilities of terminals  $z \in \mathcal{Z}$  under  $\sigma$ , where  $d(\sigma)(z)$  is the reach probability of  $z$ .  $d(\sigma)$  can be decomposed into a product of player's (and chance's) contribution of reaching  $z$ :  $d(\sigma)(z) = d_c(z) \prod_{i \in \mathcal{N}} d_i(\sigma)(z)$  where  $d_c(z)$  is chance's contribution to reaching  $z$  and  $d_i(\sigma)(z)$  is the product of  $\sigma_i(s, a)$  for infostates  $s \in \mathcal{S}_i$  on the path to  $z$ .

The average distribution over terminals across  $\psi^T$  is  $d(\psi^T) \stackrel{\text{def}}{=} \frac{1}{T} \sum_{t=1}^T d(\sigma^t)$ . We define the *marginal across terminals*  $\mu(\psi^T)$  for  $\psi^T$  as a distribution across terminals under the empirical average

<sup>2</sup>We describe this measure in more detail in Appendix A in the full version.

NashGap	CFR <sup>(+)</sup>		PCFR <sup>(+)</sup>		DCFR	LCFR	SPCFR <sup>(+)</sup>		Hedge <sup>(+)</sup>		NPCFR <sup>(+)</sup>	
10 <sup>-2</sup>	0.78	0.09	<b>1</b>	0.09	0.09	0.42	<b>1</b>	0.09	<b>1</b>	0.36	<b>1</b>	<b>1</b>
10 <sup>-3</sup>	0.09	0.02	0.91	0.02	0.02	0.02	<b>1</b>	0.02	<b>1</b>	0.06	<b>1</b>	<b>1</b>
10 <sup>-5</sup>	0	0	0.02	0	0	0	0.11	0	0.25	0	0.14	<b>1</b>

**Table 1: The fraction of games from biased\_shapley each algorithm can solve to a given NashGap within 2<sup>14</sup> = 16,384 steps. For the algorithms marked <sup>(+)</sup>, the left column shows the standard version, while the right shows the ‘plus’. See also Table 3 in Appendix D in the full version.**

strategy in the game. Formally,

$$\mu(\psi^T)(z) \stackrel{\text{def}}{=} d_c(z) \prod_{i \in \mathcal{N}} \frac{1}{T} \sum_{t=1}^T d_i(\sigma^t)(z).$$

In words, this is the distribution on terminals induced by each player’s empirical average strategy in the game tree. The sequence  $\psi^T$  is uncorrelated if  $d(\psi^T)$  and  $\mu(\psi^T)$  have no mutual dependence. This is formally captured by taking the KL divergence across terminals between  $d(\psi^T)$  and  $\mu(\psi^T)$ . We denote this KL as  $\mathbb{I}(\psi)$ , since it is equal to mutual information for the two-player case and total correlation for the  $n$ -player case [48].

**Definition 2.** We say that  $\psi^T$  is  $\epsilon$ -extensive-form marginalizable ( $\epsilon$ -EFM) if

$$\mathbb{I}(\psi^T) \stackrel{\text{def}}{=} D_{\text{KL}}\left(d(\psi^T) \parallel \mu(\psi^T)\right) \leq \epsilon. \tag{1}$$

When a sequence of strategies of a regret minimizer is close to extensive-form marginalizable, it provably converges close to a Nash equilibrium. Formally, let  $\bar{\sigma}^T$  be the average strategy profile of  $\psi^T$ .

**THEOREM 1.** If  $\psi^T$  was produced by an external regret minimizer with regret bounded by  $O(\sqrt{T})$  after  $T$  iterations and  $\psi^T$  is  $\epsilon$ -EFM, then

$$\text{NashGap}(\bar{\sigma}^T) \leq O(1/\sqrt{T}) + 2M\sqrt{2\epsilon}, \tag{2}$$

where  $M = \max_{i \in \mathcal{N}} \max_{z \in \mathcal{Z}} |u_i(z)|$ .

**PROOF SKETCH.** In normal-form, the statement can be obtained by using Pinsker’s inequality. For an extensive-form game consider its normal-form representation. This representation is exponentially larger, and many normal-form strategies have the same extensive-form equivalent. However, some of these strategies introduce additional correlations not present in the extensive-form. Considering the normal-form representation with the minimum mutual information extends the theorem to extensive-form. See Appendix A in the full version for the complete proof.  $\square$

Importantly, since the extensive-form marginalizability can be computed fully in the extensive-form, we can avoid the exponential blow-up caused by converting the game to normal-form. This is a major disadvantage of the classical algorithms, which are limited to the normal-form.

For a given horizon  $T$ , we define the meta-loss of NPCFR to be the average mutual information of the average terminal reach of

the strategies selected up to  $T$  on games  $g \sim G$

$$\mathcal{L}(\theta) = \mathbb{E}_{g \in G} \left[ \frac{1}{T} \sum_{t=1}^T \mathbb{I}(\psi_\theta^t) \right]. \tag{3}$$

Note minimizing this loss is different from directly minimizing the extensive form marginalizability after  $T$  steps. We do this to encourage the iterates to be marginalizable as well. This is analogous to minimizing  $\sum_{t=1}^T f(x^t)$  rather than  $f(x^T)$  as in [1], where the authors meta-learned a function optimizer. The computational graph of NPCFR is shown in Figure 1. The gradient of (3) originates in the cumulative mutual information and propagates through the game tree, the regrets  $r^t, R^t$  and the hidden states  $h^t$ . The gradients accumulate in the predictor  $\pi(\cdot|\theta)$ , which is used by the algorithms  $m_\theta$  at every information state  $s \in \mathcal{S}_i$  and every step  $t$ , see Algorithm 1.

## 4 EXPERIMENTS

We conduct our experiments in general-sum games where regret minimizers are not guaranteed to converge to a Nash equilibrium. Starting in the normal-form setting, we present a distribution of games for which standard regret minimization algorithms converge to a strictly correlated CCE. We then apply our meta-learning framework to the extensive-form settings, showing we can obtain much better approximate Nash equilibria than prior algorithms. Finally, we illustrate that the meta-learned algorithms may lose their empirical performance when used out-of-distribution.

We minimize (3) for  $T = 32$  iterations over 256 epochs using the Adam optimizer. The neural network uses two LSTM layers followed by a fully-connected layer. We performed a small grid search over relevant hyperparameters, see Appendix B in the full version. The meta-learning can be completed in about ten minutes for the normal-form experiments, and ten hours for the extensive-form games on a single CPU. See Table 4 in the full version for the memory requirements of all algorithms used.

We compare the meta-learned algorithms to a selection of current and former state-of-the-art regret minimization algorithms. Each algorithm is used to minimize counterfactual regret at each infostate of the game tree [52]. Specifically, we use regret matching (CFR) [24], predictive regret matching (PCFR) [18], smooth predictive regret matching (SPCFR) [17], discounted and linear regret minimization (DCFR, LCFR) [9], and Hedge [28]. Whenever applicable, we also investigate the ‘plus’ version of each algorithm [46].

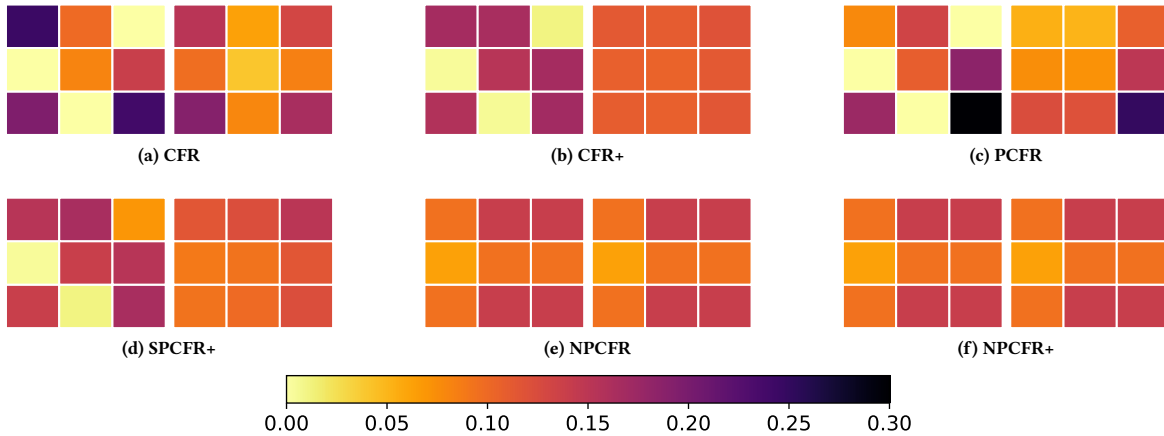


Figure 2: The empirical average joint strategy profiles found by regret minimizers  $\bar{\delta}^T$  (left) and its marginalized version (right) found on a random sample drawn from  $\text{biased\_shapley}(0, 1/2)$  after  $T = 2^{14}$  steps; see Eq. (5). Darker colors indicate higher probability under  $\bar{\delta}^T$ , and minimal differences between left and right figures imply the joint strategy is marginalizable. The remaining algorithms are shown in Figure 5 in Appendix C.1 in the full version.

### 4.1 Normal-Form Games

The Shapley game

$$u_1(\sigma) = \sigma_1^\top \cdot \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix} \cdot \sigma_2, \quad u_2(\sigma) = \sigma_1^\top \cdot \begin{pmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \end{pmatrix} \cdot \sigma_2, \tag{4}$$

was used as a simple example where the best-response dynamics doesn't stabilize [43]. Indeed, it cycles on the elements which are non-zero for one player. The empirical average joint-strategy converges to a CCE

$$\delta^* = \frac{1}{6} \begin{pmatrix} 1 & 1 & 0 \\ 0 & 1 & 1 \\ 1 & 0 & 1 \end{pmatrix}. \tag{5}$$

Clearly,  $\delta^*$  is not a Nash equilibrium, as it cannot be written as  $\sigma_1 \sigma_2^\top$ . However, thanks to the symmetry of the game, the marginals of  $\delta^*$ , or the uniform strategy, turn out to be a Nash equilibrium.

In order to break the symmetry, we perturb the utility of one of the outcomes of the game. Specifically, we give payoff  $\eta \in \mathbb{R}$  to both players when the first player selects the first, and the second their last action, see Appendix C.1 in the full version. To preserve that  $\delta^*$  is a CCE, the perturbation  $\eta$  needs to be bounded. We show in Appendix C.1 in the full version that for  $\eta \leq 1/2$ ,  $\delta^*$  is a CCE. Furthermore, there is a unique Nash equilibrium, which is non-uniform for  $\eta \neq 0$ . We denote the distribution over biased Shapley games for  $\eta \sim \mathcal{U}(a, b)$  as  $\text{biased\_shapley}(a, b)$ .

To quantify the performance of the regret minimization algorithms, we study the chance that they find a solution with a given NashGap. We present our results in Table 1. All the prior regret minimization algorithms fail to reliably find the Nash equilibrium. The 'plus' non-meta-learned algorithms exhibit particularly poor

performance in this regime, typically converging to a strictly correlated CCE. However, they don't all converge to  $\delta^*$  either, see Figure 2 for an illustration of the joint strategy profiles each algorithm converges to. In contrast, NPCFR<sup>(+)</sup> exhibit fast convergence and remarkable generalization. We show the convergence comparison of the regret minimization algorithms on  $\text{biased\_shapley}(0, 1/2)$  in Figure 4 in Appendix D.1 in the full version. Despite being trained only for  $T = 32$  steps, our meta-learned algorithms are able to minimize NashGap past  $10^4$  steps.

### 4.2 Extensive-Form Games

In the sequential setting, we use the Leduc poker [49], see also Appendix C.2 in the full version.

**4.2.1 Two-Player Leduc Poker.** Since Leduc poker is a zero-sum game, regret minimizers are guaranteed to converge to a Nash equilibrium in the two-player version. Under standard rules, players split the pot in the case of a tie, receiving a payoff equal to their total amount bet. We break the zero-sum property by modifying tie payoffs such that players only receive a  $\beta$ -fraction of their bets. This change disincentives betting to increase the size of the pot, but only if the players have the same card ranks, potentially leading to correlations in players' strategies.

We define  $\text{biased\_2p\_leduc}$  as a distribution over such games, where  $\beta \sim \mathcal{U}(0, 1/2)$ . To quantify the performance of regret minimization algorithms, we plot the expected NashGap for each algorithm on  $\text{biased\_2p\_leduc}$  in Figure 6 in the full version. While the performance averaged over the domain is similar for all algorithms, the meta-learned algorithms obtain much better approximations of Nash equilibria in each run. To show this, we investigate the chance that they find a solution with at most a given NashGap. Table 2 shows the chance for thresholds  $10^{-2}$ ,  $10^{-3}$ , and  $10^{-5}$ . With some exceptions, non-meta-learned algorithms generally fail to find a solution with a NashGap of  $10^{-2}$ . The 'plus' variants perform better

NashGap	CFR <sup>(+)</sup>		PCFR <sup>(+)</sup>		DCFR	LCFR	SPCFR <sup>(+)</sup>		Hedge <sup>(+)</sup>		NPCFR <sup>(+)</sup>	
10 <sup>-2</sup>	0	<b>1</b>	0.03	<b>1</b>	0.13	0	0.54	<b>1</b>	0	0.29	0.84	<b>1</b>
10 <sup>-3</sup>	0	0	0	0.87	0	0	0	0.72	0	0	0.73	<b>0.98</b>
10 <sup>-5</sup>	0	0	0	0.16	0	0	0	0.11	0	0	0.73	<b>0.96</b>

**Table 2: The fraction of games from `biased_2p_leduc` each algorithm can solve to a given NashGap within  $2^{18} = 262,144$  steps. For the algorithms marked <sup>(+)</sup>, the left column shows the standard version, while the right shows the ‘plus’. See also Table 3 in Appendix D.1 in the full version.**

empirically but still struggle to obtain solutions close to a Nash equilibrium as reliably as NPCFR<sup>(+)</sup>. NPCFR<sup>(+)</sup> performs the best overall.

**4.2.2 Three-Player Leduc Poker.** Generally, meta-learning is applied over a distribution of problem instances. However, in our setting, it is appealing even to apply it to a single instance of a game. This is because regret minimization algorithms are not guaranteed to converge to a Nash equilibrium in general-sum games. However, our meta-learning framework allows us to obtain better approximations of Nash equilibrium.

We demonstrate this approach on the three-player version of the Leduc poker; see Appendix C.2 for its description. We refer to the game as `three_player_leduc`. There have been conflicting reports in the literature as to the ability of regret minimization algorithms to converge to a Nash equilibrium in this game [32, 40]. We found the performance of non-meta-learned algorithms varied significantly, with those using alternating updates giving approx. 4–6-times better results. The best approximation of a Nash equilibrium we found among non-meta-learned algorithms using alternating updates<sup>3</sup> was `NashGap` = 0.004, produced by CFR<sup>(+)</sup>. Without alternating updates, we found `NashGap` = 0.027, produced by CFR. Our meta-learned algorithms have been able to find a strategy with `NashGap` = 0.012 for NPCFR, and `NashGap` = 0.001 for NPCFR<sup>(+)</sup>; see Table 6 and Figure 7 in Appendix D.3 in the full version for details. To the best of our knowledge, this is the closest approximation of Nash equilibrium of `three_player_leduc`.

To the best of our knowledge, the only theoretically sound way to find a Nash equilibrium in this game is to use support-enumeration-based algorithms such as the Lemke-Howson [29]. First, we would need to transform it into a two-player general-sum game. This can be done by having one of the players always best-respond, and treating them as a part of chance.<sup>4</sup> However, all of these algorithms work with the game in normal-form. For `three_player_leduc`, the number of pure strategies per player is  $\approx 10^{472}$ , making these approaches unusable in practice.

### 4.3 Out-of-Distribution Convergence

To illustrate that the meta-learned algorithms are tailored to a specific domain, we evaluate them out-of-distribution. Specifically, we run NPCFR<sup>(+)</sup>, which were trained on `biased_shapley(0, 1/2)`, on `biased_shapley(-1, 0)`. When evaluated out-of-distribution,

<sup>3</sup>Among the algorithms we consider, this includes the ‘plus’ algorithms and DCFR. DCFR is similar to CFR<sup>(+)</sup>, and was shown to outperform CFR<sup>(+)</sup> on two-player poker [8].

<sup>4</sup>This is the ‘inverse’ of the dummy player argument, which is normally used to show that  $n$ -player zero-sum games are as hard to solve as  $n - 1$ -player general-sum games.

the meta-learned algorithms lose the ability to converge to a Nash equilibrium. See Figure 8 in Appendix D.4 in the full version for more details.

## 5 CONCLUSION

We present a novel framework for approximating Nash equilibria in general-sum games. We apply regret minimization, which is a family of efficient algorithms, guaranteed to converge to a coarse-correlated equilibrium (CCE). This weaker solution concept allows players to correlate their strategies. We use meta-learning to search a class of predictive regret minimization algorithms, minimizing the correlations in the CCE found by the algorithm. The resulting algorithm is still guaranteed to converge to a CCE, and is meta-learned to empirically find close approximations of Nash equilibria. Experiments in general-sum games, including large imperfect-information games, reveal our algorithms can considerably outperform other regret minimization algorithms.

*Future Work.* Our meta-learning framework might be useful for finding CCEs with desired properties. For example, one can search for welfare maximizing equilibria by setting the meta-loss to the negative total utility of all players. We also see other domains, such as auctions, as a promising field where our approach can be used. One limitation of our approach is that it can be quite memory demanding, especially for larger horizons. Training on abstractions of the games is promising.

## ACKNOWLEDGMENTS

This work was supported by the Horizon Europe Programme under Grant Agreement No. 101183743 (AGATE) and by the Czech Science Foundation (GAČR) under grant No. 25-18031S. Computational resources were supplied by the project “e-Infrastruktura CZ” (e-INFRA LM2018140) provided within the program Projects of Large Research, Development and Innovations Infrastructures.

## REFERENCES

- [1] Marcin Andrychowicz, Misha Denil, Sergio Gómez Colmenarejo, Matthew W. Hoffman, David Pfau, Tom Schaul, Brendan Shillingford, and Nando de Freitas. 2016. Learning to learn by gradient descent by gradient descent. In *Proceedings of the 30th International Conference on Neural Information Processing Systems (Barcelona, Spain) (NeurIPS’16)*. Curran Associates Inc., Red Hook, NY, USA, 3988–3996.
- [2] David Blackwell et al. 1956. An analog of the minimax theorem for vector payoffs. *Pacific J. Math.* 6, 1 (1956), 1–8. <https://doi.org/10.2140/PJM.1956.6.1>
- [3] Hartwig Bosse, Jaroslaw Byrka, and Evangelos Markakis. 2010. New algorithms for approximate Nash equilibria in bimatrix games. *Theoretical Computer Science* 411, 1 (2010), 164–173. <https://doi.org/10.1016/j.tcs.2009.09.023>

- [4] Michael Bowling, Neil Burch, Michael Johanson, and Oskari Tammelin. 2015. Heads-up limit hold'em poker is solved. *Science* 347, 6218 (2015), 145–149. <https://doi.org/10.1126/science.1259433>
- [5] Noam Brown, Anton Bakhtin, Adam Lerer, and Qucheng Gong. 2020. Combining deep reinforcement learning and search for imperfect-information games. *Advances in Neural Information Processing Systems* 33 (2020), 17057–17069.
- [6] Noam Brown, Adam Lerer, Sam Gross, and Tuomas Sandholm. 2019. Deep counterfactual regret minimization. In *International conference on machine learning*. PMLR, AAAI press, 1101 Pennsylvania Ave, NW Washington, DC USA, 793–802.
- [7] Noam Brown and Tuomas Sandholm. 2018. Superhuman AI for heads-up no-limit poker: Libratus beats top professionals. *Science* 359, 6374 (2018), 418–424. <https://doi.org/10.1126/science.aao1733>
- [8] Noam Brown and Tuomas Sandholm. 2019. Solving imperfect-information games via discounted regret minimization. In *Proceedings of the Thirty-Third AAAI Conference on Artificial Intelligence* (Honolulu, Hawaii, USA) (AAAI'19). AAAI Press, 1101 Pennsylvania Ave, NW Washington, DC USA, Article 225, 8 pages. <https://doi.org/10.1609/aaai.v33i01.33011829>
- [9] Noam Brown and Tuomas Sandholm. 2019. Solving Imperfect-Information Games via Discounted Regret Minimization. *Proceedings of the AAAI Conference on Artificial Intelligence* 33, 01 (Jul. 2019), 1829–1836. <https://doi.org/10.1609/aaai.v33i01.33011829>
- [10] Noam Brown and Tuomas Sandholm. 2019. Superhuman AI for multiplayer poker. *Science* 365, 6456 (2019), 885–890. <https://doi.org/10.1126/science.aay2400>
- [11] Yang Cai and Constantinos Daskalakis. 2011. On Minmax Theorems for Multiplayer Games. In Proceedings of the Twenty-Second Annual ACM-SIAM Symposium on Discrete Algorithms, SODA. *Proceedings of the Annual ACM-SIAM Symposium on Discrete Algorithms*, 217–234. <https://doi.org/10.1137/1.9781611973082.20>
- [12] Bruno Codenotti, Stefano De Rossi, and Marino Pagan. 2008. An experimental analysis of Lemke-Howson algorithm. arXiv:0811.3247 [cs.DS]
- [13] Constantinos Daskalakis, Paul W. Goldberg, and Christos H. Papadimitriou. 2009. The Complexity of Computing a Nash Equilibrium. *SIAM J. Comput.* 39, 1 (2009), 195–259. <https://doi.org/10.1137/070699652> arXiv:<https://doi.org/10.1137/070699652>
- [14] Constantinos Daskalakis, Aranyak Mehta, and Christos Papadimitriou. 2007. Progress in approximate nash equilibria. In *Proceedings of the 8th ACM Conference on Electronic Commerce* (San Diego, California, USA) (EC '07). Association for Computing Machinery, New York, NY, USA, 355–358. <https://doi.org/10.1145/1250910.1250962>
- [15] Constantinos Daskalakis, Aranyak Mehta, and Christos Papadimitriou. 2009. A note on approximate Nash equilibria. *Theoretical Computer Science* 410, 17 (2009), 1581–1588. <https://doi.org/10.1016/j.tcs.2008.12.031> Internet and Network Economics.
- [16] Argyrios Deligkas, Michail Fasoulakis, and Evangelos Markakis. 2023. A Polynomial-Time Algorithm for 1/3-Approximate Nash Equilibria in Bimatrix Games. *ACM Trans. Algorithms* 19, 4, Article 31 (aug 2023), 17 pages. <https://doi.org/10.1145/3606697>
- [17] Gabriele Farina, Julien Grand-Clément, Christian Kroer, Chung-Wei Lee, and Haipeng Luo. 2023. Regret Matching+: (In)Stability and Fast Convergence in Games. In *Advances in Neural Information Processing Systems*, A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine (Eds.), Vol. 36. Curran Associates, Inc., 57 Morehouse Lane Red Hook, NY 12571 USA, 61546–61572. [https://proceedings.neurips.cc/paper\\_files/paper/2023/file/c209cd57e13f3344a4cad4ce84d0ee1b-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2023/file/c209cd57e13f3344a4cad4ce84d0ee1b-Paper-Conference.pdf)
- [18] Gabriele Farina, Christian Kroer, and Tuomas Sandholm. 2021. Faster Game Solving via Predictive Blackwell Approachability: Connecting Regret Matching and Mirror Descent. *Proceedings of the AAAI Conference on Artificial Intelligence* 35, 6 (May 2021), 5363–5371. <https://doi.org/10.1609/aaai.v35i6.16676>
- [19] Nicola Gatti, Giorgio Patrini, Marco Rocco, and Tuomas Sandholm. 2012. Combining local search techniques and path following for bimatrix games. <http://arxiv.org/abs/1210.4858>
- [20] Richard Gibson. 2014. *Regret minimization in games and the development of champion multiplayer computer poker-playing agents*. Ph.D. Dissertation. University of Alberta.
- [21] Itzhak Gilboa and Eitan Zemel. 1989. Nash and correlated equilibria: Some complexity considerations. *Games and Economic Behavior* 1, 1 (1989), 80–93. [https://doi.org/10.1016/0899-8256\(89\)90006-7](https://doi.org/10.1016/0899-8256(89)90006-7)
- [22] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. 2014. Generative Adversarial Nets. In *Advances in Neural Information Processing Systems*, Z. Ghahramani, M. Welling, C. Cortes, N. Lawrence, and K.Q. Weinberger (Eds.), Vol. 27. Curran Associates, Inc., 57 Morehouse Lane Red Hook, NY 12571 USA. [https://proceedings.neurips.cc/paper\\_files/paper/2014/file/5ca3e9b122f61f8f06494c97b1afce3-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2014/file/5ca3e9b122f61f8f06494c97b1afce3-Paper.pdf)
- [23] James Hannan. 1957. Approximation to Bayes risk in repeated play. *Contributions to the Theory of Games* 3 (1957), 97–139.
- [24] Sergiu Hart and Andreu Mas-Colell. 2000. A simple adaptive procedure leading to correlated equilibrium. *Econometrica* 68, 5 (2000), 1127–1150. <https://doi.org/10.1111/1468-0262.00153>
- [25] Michael Kearns, Michael L. Littman, and Satinder Singh. 2001. Graphical models for game theory. In *Proceedings of the Seventeenth Conference on Uncertainty in Artificial Intelligence* (Seattle, Washington) (UAI'01). Morgan Kaufmann Publishers Inc., San Francisco, CA, USA, 253–260.
- [26] Spyros C. Kontogiannis, Panagiota N. Panagopoulou, and Paul G. Spirakis. 2009. Polynomial algorithms for approximating Nash equilibria of bimatrix games. *Theoretical Computer Science* 410, 17 (2009), 1599–1606. <https://doi.org/10.1016/j.tcs.2008.12.033> Internet and Network Economics.
- [27] Vojtěch Kovařík, Martin Schmid, Neil Burch, Michael Bowling, and Viliam Lisý. 2022. Rethinking formal models of partially observable multiagent decision making. *Artificial Intelligence* 303 (2022), 103645.
- [28] Tor Lattimore and Csaba Szepesvári. 2020. *Bandit algorithms*. Cambridge University Press, 1 Liberty Plaza, New York, NY 10006, USA.
- [29] C. E. Lemke and J. T. Howson, Jr. 1964. Equilibrium Points of Bimatrix Games. *J. Soc. Indust. Appl. Math.* 12, 2 (1964), 413–423. <https://doi.org/10.1137/0112033> arXiv:<https://doi.org/10.1137/0112033>
- [30] Hanyu Li, Wenhan Huang, Zhijian Duan, David Henry Mguni, Kun Shao, Jun Wang, and Xiaotie Deng. 2024. A survey on algorithms for Nash equilibria in finite normal-form games. *Computer Science Review* 51 (2024), 100613. <https://doi.org/10.1016/j.cosrev.2023.100613>
- [31] Michael L. Littman and Peter Stone. 2005. A polynomial-time Nash equilibrium algorithm for repeated games. *Decision Support Systems* 39, 1 (2005), 55–66. <https://doi.org/10.1016/j.dss.2004.08.007> The Fourth ACM Conference on Electronic Commerce.
- [32] Revan MacQueen and James Wright. 2024. Guarantees for Self-Play in Multiplayer Games via Polymatrix Decomposability. *Advances in Neural Information Processing Systems* 36 (2024).
- [33] Paul R. Milgrom and Robert J. Weber. 1982. A Theory of Auctions and Competitive Bidding. *Econometrica* 50, 5 (1982), 1089–1122. <http://www.jstor.org/stable/1911865>
- [34] Dov Monderer and Lloyd S Shapley. 1996. Potential games. *Games and economic behavior* 14, 1 (1996), 124–143.
- [35] Matej Moravčík, Martin Schmid, Neil Burch, Viliam Lisý, Dustin Morrill, Nolan Bard, Trevor Davis, Kevin Waugh, Michael Johanson, and Michael Bowling. 2017. Deepstack: Expert-level artificial intelligence in heads-up no-limit poker. *Science* 356, 6337 (2017), 508–513. <https://doi.org/10.1126/science.aam6960>
- [36] H. Moulin and J.P. Vial. 1975. Strategically zero-sum games: The class of games whose completely mixed equilibria cannot be improved upon. *International Journal of Game Theory* 7 (9 1975), 201–221. <https://doi.org/10.1007/BF01769190>
- [37] K.G. Murty. 1984. *Linear Programming*. Wiley, Hoboken, 111 River St, United States. <https://books.google.ca/books?id=ibQJvAEACAAJ>
- [38] Noam Nisan, Tim Roughgarden, Eva Tardos, and Vijay V Vazirani. 2007. *Algorithmic Game Theory*. Cambridge University Press, New York, NY, USA.
- [39] Christos H. Papadimitriou. 1994. On the complexity of the parity argument and other inefficient proofs of existence. *J. Comput. System Sci.* 48, 3 (1994), 498–532. [https://doi.org/10.1016/S0022-0000\(05\)80063-7](https://doi.org/10.1016/S0022-0000(05)80063-7)
- [40] Nick Abou Risk and Duane Szafron. 2010. Using counterfactual regret minimization to create competitive multiplayer poker agents. In *Proceedings of the 9th International Conference on Autonomous Agents and Multiagent Systems: Volume 1 - Volume 1* (Toronto, Canada) (AAMAS '10). International Foundation for Autonomous Agents and Multiagent Systems, Richland, SC, 159–166.
- [41] Robert W. Rosenthal. 1973. A class of games possessing pure-strategy Nash equilibria. *Int. J. Game Theory* 2, 1 (dec 1973), 65–67. <https://doi.org/10.1007/BF01737559>
- [42] Martin Schmid, Matej Moravčík, Neil Burch, Rudolf Kadlec, Josh Davidson, Kevin Waugh, Nolan Bard, Finbarr Timbers, Marc Lanctot, G. Zacharias Holland, Elnaz Davoodi, Alden Christianson, and Michael Bowling. 2023. Student of Games: A unified learning algorithm for both perfect and imperfect information games. *Science Advances* 9, 46 (2023), eadg3256. <https://doi.org/10.1126/sciadv.adg3256> arXiv:<https://www.science.org/doi/pdf/10.1126/sciadv.adg3256>
- [43] L. S. Shapley. 1964. *1. Some Topics in Two-Person Games*. Princeton University Press, Princeton, 1–28. <https://doi.org/10.1515/9781400882014-002>
- [44] David Sychrovský, Michal Šustr, Elnaz Davoodi, Michael Bowling, Marc Lanctot, and Martin Schmid. 2024. Learning Not to Regret. *Proceedings of the AAAI Conference on Artificial Intelligence* 38, 14 (Mar. 2024), 15202–15210. <https://doi.org/10.1609/aaai.v38i14.29443>
- [45] Oskari Tammelin. 2014. Solving large imperfect information games using CFR+. <https://arxiv.org/abs/1407.5042>
- [46] Oskari Tammelin, Neil Burch, Michael Johanson, and Michael Bowling. 2015. Solving heads-up limit Texas Hold'em. In *Proceedings of the 24th International Conference on Artificial Intelligence* (Buenos Aires, Argentina) (IJCAI'15). AAAI Press, 1101 Pennsylvania Ave, NW, Suite 300 Washington, DC 20004, 645–652.
- [47] William Vickrey. 1961. Counterspeculation, Auctions, and Competitive Sealed Tenders. *The Journal of Finance* 16, 1 (1961), 8–37. <http://www.jstor.org/stable/2977633>
- [48] Satoshi Watanabe. 1960. Information Theoretical Analysis of Multivariate Correlation. *IBM Journal of Research and Development* 4, 1 (1960), 66–82. <https://doi.org/10.1147/rd.41.0066>

- [49] Kevin Waugh, Nolan Bard, and Michael Bowling. 2009. Strategy Grafting in Extensive Games. In *Advances in Neural Information Processing Systems*, Y. Bengio, D. Schuurmans, J. Lafferty, C. Williams, and A. Culotta (Eds.), Vol. 22. Curran Associates, Inc., 57 Morehouse Lane Red Hook, NY 12571 USA. [https://proceedings.neurips.cc/paper\\_files/paper/2009/file/e0ec453e28e061cc58ac43f91dc2f3f0-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2009/file/e0ec453e28e061cc58ac43f91dc2f3f0-Paper.pdf)
- [50] Naifeng Zhang, Stephen McAleer, and Tuomas Sandholm. 2024. Faster Game Solving via Hyperparameter Schedules. arXiv:2404.09097 [cs.GT]
- [51] Martin Zinkevich. 2003. Online convex programming and generalized infinitesimal gradient ascent. In *Proceedings of the Twentieth International Conference on International Conference on Machine Learning (Washington, DC, USA) (ICML'03)*. AAAI Press, 1101 Pennsylvania Ave, NW Washington, DC USA, 928–935.
- [52] Martin Zinkevich, Michael Johanson, Michael Bowling, and Carmelo Piccione. 2007. Regret minimization in games with incomplete information. In *Advances in neural information processing systems*. Advances in Neural Information Processing Systems 20, 340 Pine Street, Sixth Floor. San Francisco, CA, 1729–1736.