

FASTER NO-REGRET LEARNING DYNAMICS FOR EXTENSIVE-FORM CORRELATED EQUILIBRIUM

Anonymous authors

Paper under double-blind review

ABSTRACT

A recent emerging trend in the literature on learning in games has been concerned with providing accelerated learning dynamics for correlated and coarse correlated equilibria in normal-form games. Much less is known about the significantly more challenging setting of extensive-form games, which can capture sequential and simultaneous moves, as well as imperfect information. In this paper, we develop faster no-regret learning dynamics for *extensive-form correlated equilibrium (EFCE)* in multiplayer general-sum imperfect-information extensive-form games. When all agents play T repetitions of the game according to the accelerated dynamics, the correlated distribution of play is an $O(T^{-3/4})$ -approximate EFCE. This significantly improves over the best prior rate of $O(T^{-1/2})$. One of our conceptual contributions is to connect predictive (that is, optimistic) regret minimization with the framework of Φ -regret. One of our main technical contributions is to characterize the stability of certain fixed point strategies through a refined perturbation analysis of a structured Markov chain, which may be of independent interest. Finally, experiments on standard benchmarks corroborate our findings.

1 INTRODUCTION

Game-theoretic solution concepts describe how agents should rationally act in games. Over the last two decades there has been tremendous progress in imperfect-information game solving and algorithms based on game-theoretic solution concepts have become the state of the art. Prominent milestones of this were an optimal strategy for Rhode Island hold'em poker (Gilpin & Sandholm, 2007), a near-optimal strategy for limit Texas hold'em (Bowling et al., 2015), and a superhuman strategy for no-limit Texas hold'em (Brown & Sandholm, 2017). In particular, these advances rely on algorithms that approximate *Nash equilibria (NE)* of two-player zero-sum *extensive-form games (EFGs)*. EFGs are a broad class of games that capture sequential and simultaneous interaction, and imperfect information. For two-player zero-sum EFGs, it is by now well-understood how to compute a Nash equilibrium at scale: in theory this can be achieved using accelerated uncoupled no-regret learning dynamics, for example by having each player use an *optimistic* regret minimizer and leveraging suitable *distance-generating functions* (Hoda et al., 2010; Kroer et al., 2020; Farina et al., 2021c) for the EFG decision space. Such a setup converges to an equilibrium at a rate of $O(T^{-1})$. In practice, modern variants of the *counterfactual regret minimization (CFR)* framework typically lead to better performance, although the worst-case convergence rate is $O(T^{-1/2})$ (Zinkevich et al., 2007). CFR is also an uncoupled no-regret learning dynamic.

However, many real-world applications are not two-player zero-sum games, but instead have *general-sum* utilities and often more than two players. In such settings, Nash equilibrium suffers from several drawbacks when used as a prescriptive tool. First, there can be multiple equilibria, and an equilibrium strategy may perform very poorly when played against the “wrong” equilibrium strategies of the other player(s). Thus, the players effectively would need to communicate in order to find an equilibrium, or hope to converge to it via some sort of learning dynamics. Second, finding a Nash equilibrium is computationally hard both in theory (Daskalakis et al., 2006; Etessami & Yannakakis, 2007) and in practice (Berg & Sandholm, 2017). This effectively squashes any hope of developing efficient learning dynamics that converge to general-sum Nash equilibria.

A competing notion of rationality proposed by Aumann (1974) is that of *correlated equilibrium (CE)*, typically modeled via a trusted mediator who privately recommends actions to the players.

Unlike NE, it is known that the latter can be computed in polynomial time and, perhaps even more importantly, it can be attained through *uncoupled* learning dynamics, where the players only need to reason about their own observed utilities. This overcomes the often unreasonable presumption that players have knowledge about the other players’ utilities. At the same time, uncoupled learning algorithms have proven to be a remarkably *scalable* approach for computing equilibria in large-scale games, as described above. The basic CE notion is defined for normal-form games, and there it has long been known that uncoupled no-regret learning dynamics can converge to CE or the *coarse correlated equilibrium* (CCE) variant at a rate of $O(T^{-1/2})$ (Hart & Mas-Colell, 2000; Celli et al., 2019). More recently, it was shown that accelerated uncoupled no-regret learning dynamics can compute CCE and CE at a rate of $O(T^{-3/4})$ (Syrkanis et al., 2015; Chen & Peng, 2020).

In the context of EFGs, the idea of correlation is much more intricate, and there are several notions of correlated equilibrium, based on when the mediator gives recommendations and how the mediator reacts to players who disregard the advice. One of the most compelling notions for EFGs is the *extensive-form correlated equilibrium* (henceforth EFCE) (von Stengel & Forges, 2008) for extensive-form games with *perfect recall*. Because of the sequential nature, the presence of private information in the game, and the gradual revelation of recommendations, the constraints associated with EFCE are significantly more complex than for normal-form games. For these reasons, the question of whether uncoupled learning dynamics can converge to an EFCE was only very recently resolved by Celli et al. (2020). Moreover, in a follow-up work they also established an explicit rate of convergence of $O(T^{-1/2})$ (Farina et al., 2021a). Our paper is concerned with the following fundamental question: *Can one develop faster uncoupled no-regret learning dynamics for EFCE?*

Contributions. Our primary contribution is to answer this question in the positive:

Theorem 1.1. *On any finite perfect-recall general-sum multiplayer extensive-form game, the uncoupled no-regret learning dynamics described in this paper lead to a correlated distribution of play that is an $O(T^{-3/4})$ -approximate EFCE, where the $O(\cdot)$ notation suppresses game-specific parameters polynomial in the size of the game.*

We achieve this result using the framework of *predictive* (also known as *optimistic*) regret minimization (Chiang et al., 2012; Rakhlin & Sridharan, 2013b). One of our conceptual contributions is to connect this line of work with the framework of Φ -regret minimization of Greenwald & Jafari (2003); Gordon et al. (2008), by providing a general template for stable-predictive Φ -regret minimization. The importance of Φ -regret is that it leads to substantially more powerful notions of hindsight rationality, beyond the usual *external* regret (Gordon et al., 2008), including the powerful notion of *swap regret* (Blum & Mansour, 2007). Moreover, one of the primary insights behind the result of Farina et al. (2021a) is to cast convergence to an EFCE as a Φ -regret minimization problem. Given these prior connections, we believe that our stable-predictive Φ template is of independent interest, and could lead to further applications in the future.

Theorem 1.1 extends and strengthens several prior papers in the literature, including the seminal work of Syrgkanis et al. (2015) that provides accelerated dynamics for *coarse* correlated equilibrium in normal-form games, as well as the more recent result of Chen & Peng (2020) which showed $O(T^{-3/4})$ convergence to a correlated equilibrium in normal-form games. For the more challenging class of extensive-form games, accelerated rates were previously known only for finding a *Nash* equilibrium in the special case of *two-player zero-sum* games, where an $O(T^{-3/4})$ rate was achieved via a stable-predictive CFR setup (Farina et al., 2019a) and an $O(T^{-1})$ rate was achieved via optimistic regret minimizers coupled with good distance-generating functions (Farina et al., 2019c).

From a technical standpoint, in order to apply our generic template for accelerated Φ -regret minimization, we establish two separate ingredients. First, we develop a *stable-predictive* external regret minimizer for the set of transformations Φ associated with EFCE. This differs from the construction by Farina et al. (2021a) in that we have to additionally guarantee and preserve the stability—and subsequently the predictivity—throughout the construction. The second component consists of sharply characterizing the stability of fixed points of *trigger deviation functions*. This turns out to be particularly challenging, and direct extensions of prior techniques appear to only give a bound that is *exponential* in the size of the game. In this context, one of our key technical contributions is to provide a refined perturbation analysis for a Markov chain consisting of a rank-one stochastic matrix, employing tools that have not been used before in this line of work, and substantially extending the techniques of Chen & Peng (2020). This leads to a rate of convergence that depends *polynomially*

on the description of the game, which is crucial for the applicability of the accelerated dynamics. Finally, we support our theoretical findings with experiments on several general-sum benchmarks.

Further Related Work. The line of work on accelerated no-regret learning for *Nash* equilibrium was pioneered by Daskalakis et al. (2015), showing that one can bypass the adversarial $\Omega(T^{-1/2})$ barrier for the incurred average regret if *both* players in a zero-sum game employ an uncoupled variant of the excessive gap technique (Nesterov, 2005), leading to a near-optimal rate of $O(\log T/T)$. Subsequently, Rakhlin & Sridharan (2013a) showed that the optimal rate of $O(1/T)$ can be obtained with a remarkably simple variant of Online Mirror Descent which incorporates a *prediction* term in the update step. While these results only hold for zero-sum games, Syrgkanis et al. (2015) showed that $O(T^{-3/4})$ rate can be obtained for multiplayer general-sum normal-form games. In a recent result, Chen & Peng (2020) strengthened the regret bounds of Syrgkanis et al. (2015) from external to swap regret using the celebrated construction of Blum & Mansour (2007). We also acknowledge a recent result of Daskalakis et al. (2021) which establishes a near-optimal rate of convergence of $\tilde{O}(1/T)$ to a coarse correlated equilibrium when all players employ the Optimistic Multiplicative Weights Update (OMWU) algorithm in a normal-form game. Extending their result to extensive-form games presents considerable technical challenges since their analysis crucially hinges on the closed-form softmax-type structure of OMWU on the simplex.

Correlated equilibrium in extensive-form games is much less understood than Nash equilibrium. A feasible EFCE can also be computed efficiently through a variant of the *Ellipsoid algorithm* (Papadimitriou & Roughgarden, 2008; Jiang & Leyton-Brown, 2015), and an alternative sampling-based approach was given by Dudík & Gordon (2009). However, those approaches perform poorly in large-scale problems, and do not allow the players to arrive at EFCE via distributed learning. Celli et al. (2019) devised variants of the CFR algorithm that provably convergence to *normal-form coarse correlated equilibria*, a solution concept much less appealing than EFCE in extensive-form games Gordon et al. (2008). Finally, Morrill et al. (2021a;b) characterize hindsight rationality notions and associate a set of solution concepts with suitable $O(T^{-1/2})$ no-regret learning dynamics.

2 PRELIMINARIES

Extensive-form Games. An extensive-form game is abstracted on a directed and rooted *game tree* \mathcal{T} . The set of nodes of \mathcal{T} is denoted with \mathcal{H} ; non-terminal nodes are referred as *decision nodes*, and are associated with a player who acts by selecting an action from a set of possible actions $\mathcal{A}(h)$, where $h \in \mathcal{H}$ represents the decision node. By convention, the set of players $[n] \cup \{c\}$ includes a *fictitious* agent c who “selects” actions according to fixed probability distributions dictated by the nature of the game (e.g., the roll of a dice); this intends to model external stochastic phenomena occurring during the game. For a player $i \in [n] \cup \{c\}$, we let $\mathcal{H}^{(i)} \subseteq \mathcal{H}$ be the subset of decision nodes wherein a player i makes a decision. The set of *leaves* $\mathcal{Z} \subseteq \mathcal{H}$, or equivalently the *terminal nodes*, correspond to different outcomes; once the game transitions to a terminal node $z \in \mathcal{Z}$, payoffs are assigned to each player based on a set of *normalized* utility functions $\{u^{(i)} : \mathcal{Z} \rightarrow [-1, 1]\}_{i \in [n]}$. It will also be convenient to represent with $p^{(c)}(z)$ the product of probabilities of “chance” moves encountered in the path from the root until the terminal node $z \in \mathcal{Z}$.

Imperfect Information. To model imperfect information, the set of decision nodes $\mathcal{H}^{(i)}$ of player i are partitioned into a collection of sets $\mathcal{J}^{(i)}$, which are called *information sets*. Each information set $j \in \mathcal{J}^{(i)}$ groups nodes which cannot be distinguished by i . Thus, for any nodes $h, h' \in j$ we have $\mathcal{A}(h) = \mathcal{A}(h')$. As usual, we assume that the game satisfies *perfect recall*: players never forget information once acquired. We will also define a partial order \prec on $\mathcal{J}^{(i)}$, so that $j \prec j'$, for $j, j' \in \mathcal{J}^{(i)}$, if there exist nodes $h \in j$ and $h' \in j'$ such that the path from the root to h' passes through h . If $j \prec j'$, we will say that j is an *ancestor* of j' , or equivalently, j is a descendant of j' .

Sequence-form Strategies. For a player $i \in [n]$, an information set $j \in \mathcal{J}^{(i)}$, and an action $a \in \mathcal{A}(j)$, we will denote with $\sigma = (j, a)$ the *sequence* of i ’s actions encountered on the path from the root of the game until (and included) action a . For notational convenience, we will use the special symbol \emptyset to denote the *empty sequence*. Then, i ’s set of sequences is defined as $\Sigma^{(i)} := \{(j, a) : j \in \mathcal{J}^{(i)}, a \in \mathcal{A}(j)\} \cup \{\emptyset\}$; we will also use the notation $\Sigma_*^{(i)} := \Sigma^{(i)} \setminus \{\emptyset\}$. For a given information set $j \in \mathcal{J}^{(i)}$ we will use $\sigma^{(i)}(j) \in \Sigma^{(i)}$ to represent the *parent sequence*; i.e. the last sequence

encountered by player i before reaching any node in the information set j , assuming that it exists. Otherwise, we let $\sigma^{(i)}(j) = \emptyset$, and we say that j is the *root information set* of player i . A *strategy* for a player specifies a probability distribution for every possible information set encountered in the game tree. For perfect-recall EFGs, strategies can be equivalently represented in *sequence-form*:

Definition 2.1 (Sequence-form Polytope). The *sequence-form strategy polytope* for player $i \in [n]$ is defined as the following (convex) polytope:

$$\mathcal{Q}^{(i)} := \left\{ \mathbf{q} \in \mathbb{R}_{\geq 0}^{|\Sigma^{(i)}|} : \mathbf{q}[\emptyset] = 1, \quad \mathbf{q}[\sigma^{(i)}(j)] = \sum_{a \in \mathcal{A}(j)} \mathbf{q}[(j, a)], \quad \forall j \in \mathcal{J}^{(i)} \right\}. \quad (1)$$

Analogously, one can define the sequence-form strategy polytope for the *subtree* of the partially ordered set $(\mathcal{J}^{(i)}, \prec)$ rooted at $j \in \mathcal{J}^{(i)}$, which will be denoted as $\mathcal{Q}_j^{(i)}$. Moreover, the set of *deterministic* sequence-form strategies for player $i \in [n]$ is the set $\Pi^{(i)} = \mathcal{Q}^{(i)} \cap \{0, 1\}^{|\Sigma^{(i)}|}$, and similarly for $\Pi_j^{(i)}$. The *joint* set of deterministic sequence-form strategies of the players will be represented with $\Pi := \times_{i \in [n]} \Pi^{(i)}$. As such, an element $\boldsymbol{\pi} \in \Pi$ is an n -tuple $(\boldsymbol{\pi}^{(1)}, \dots, \boldsymbol{\pi}^{(n)})$ specifying a deterministic sequence-form strategy for every player $i \in [n]$. Finally, the utility of player $i \in [n]$ under a profile $\boldsymbol{\pi} \in \Pi$ can be expressed as

$$u^{(i)}(\boldsymbol{\pi}) := \sum_{z \in \mathcal{Z}} p^{(c)}(z) u^{(i)}(z) \mathbf{1}\{\boldsymbol{\pi}^{(k)}[\sigma^{(k)}(z)] = 1, \forall k \in [n]\}. \quad (2)$$

We summarized in Table 1 the EFG notation that we will be using most often throughout the paper.

An Illustrative Example. To clarify some of the concepts we have introduced, we illustrate a simple two-player EFG in Figure 1. Black nodes belong to player 1, white round nodes to player 2, square nodes are terminal nodes (aka leaves), and the crossed node is a chance node. Player 2 has two information sets, $\mathcal{J}^{(2)} := \{C, D\}$, each containing two nodes. This captures the lack of knowledge regarding the action played by player 1. In contrast, the outcome of the chance move is observed by both players. At the information set C, player 2 has two possible actions, $\mathcal{A}(C) := \{5, 6\}$. Thus, one possible sequence for player 2 is the pair $\sigma = (C, 5) \in \Sigma^{(2)}$.

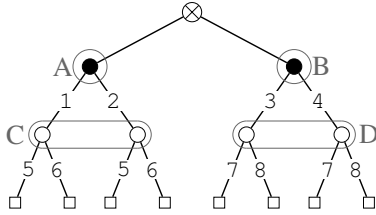


Figure 1: Example of a two-player EFG.

Description	
$\mathcal{J}^{(i)}$	Information sets of player i
$\mathcal{A}(j)$	Actions at information set j
$\Sigma^{(i)}$	Set of sequences of player i
$\mathcal{Q}_j^{(i)}$	Sequence-form strategies rooted at $j \in \mathcal{J}^{(i)}$
$\mathcal{D}^{(i)}$	Maximum depth of any $j \in \mathcal{J}^{(i)}$

Table 1: Summary of the basic notation.

Regret, Φ -Regret and Optimistic Regret Minimization. Consider a convex and compact set $\mathcal{X} \subseteq \mathbb{R}^d$ representing the space of strategies of some agent. In the online decision making framework, a *regret minimizer* \mathcal{R} can be thought of as a black-box device which interacts with the external environment via the following two basic subroutines:

- $\mathcal{R}.$ NEXTSTRATEGY(): The regret minimizer returns the strategy $\mathbf{x}^t \in \mathcal{X}$ at time t ;
- $\mathcal{R}.$ OBSERVEUTILITY(ℓ^t): The regret minimizer receives as feedback a linear utility function $\ell^t : \mathcal{X} \ni \mathbf{x} \mapsto \langle \ell^t, \mathbf{x} \rangle$, and may alter its internal state accordingly.

The decision making is *online* in the sense that the regret minimizer can adapt to previously received information, but no information about future utilities is available. The error of a regret minimizer is typically measured in terms of *external regret*, defined, for a time horizon T , as follows:

$$R^T := \max_{\mathbf{x}^* \in \mathcal{X}} \sum_{t=1}^T \langle \mathbf{x}^*, \ell^t \rangle - \sum_{t=1}^T \langle \mathbf{x}^t, \ell^t \rangle, \quad (3)$$

That is, the performance of the online algorithm is compared with the best *fixed* strategy in *hindsight*.

Φ -Regret. A conceptual generalization of the concept of external regret is the so-called Φ -regret. Specifically, in this framework the performance of the learning algorithm is measured based on a *set of transformations* $\Phi : \mathcal{X} \rightarrow \mathcal{X}$, leading to the notion of cumulative Φ -regret:

$$R^T := \max_{\phi^* \in \Phi} \sum_{t=1}^T \langle \phi^*(\mathbf{x}^t), \ell^t \rangle - \sum_{t=1}^T \langle \mathbf{x}^t, \ell^t \rangle. \quad (4)$$

When the set of transformations Φ coincides with the set of *constant* functions, one recovers the notion of external regret given in Equation (3). However, Φ -regret is substantially more expressive and yields a more appealing notion of hindsight rationality (Gordon et al., 2008), incorporating the notion of *swap regret* (Blum & Mansour, 2007).

We will employ the following definition, which is a slight modification of the RVU property introduced by (Syrkanis et al., 2015, Definition 3).

Definition 2.2 (Stable-predictivity). Let \mathcal{R} be a regret minimizer and let $\|\cdot\|$ be a norm. \mathcal{R} is said to be κ -stable with respect to $\|\cdot\|$ if for all $t \geq 2$, the strategies output by \mathcal{R} satisfy

$$\|\mathbf{x}^t - \mathbf{x}^{t-1}\| \leq \kappa, \quad (5)$$

Moreover, it is said to be (α, β) -predictive with respect to $\|\cdot\|$ if for all $t \geq 1$ its regret R^T satisfies

$$R^T \leq \alpha(T) + \beta \sum_{t=2}^T \|\ell^t - \ell^{t-1}\|_*^2, \quad (6)$$

no matter the sequence of utility vectors ℓ^1, \dots, ℓ^T , where $\|\cdot\|_*$ is the dual norm of $\|\cdot\|$.

Optimistic Follow the Regularized Leader. Let d be a 1-strongly convex function with respect to some norm $\|\cdot\|$, and $\eta > 0$ the *learning rate*. OFTRL's update rule takes the following form:

$$\mathbf{x}^t := \arg \max_{\mathbf{x} \in \mathcal{X}} \left\{ \left\langle \mathbf{x}, 2\ell^{t-1} + \sum_{\tau=1}^{t-2} \ell^\tau \right\rangle - \frac{d(\mathbf{x})}{\eta} \right\}, \quad (\text{OFTRL})$$

where $\mathbf{x}^1 := \arg \min_{\mathbf{x} \in \mathcal{X}} d(\mathbf{x})$. Syrgkanis et al. (2015) established the following property:

Lemma 2.3. (OFTRL) is 2η -stable and $(\Omega_d/\eta, \eta)$ -predictive with respect to any norm $\|\cdot\|$ for which d is 1-strongly convex, where Ω_d is the range of d on \mathcal{X} , that is, $\Omega_d := \max_{\mathbf{x}, \mathbf{x}' \in \mathcal{X}} \{d(\mathbf{x}) - d(\mathbf{x}')\}$.

In this paper, we consider the entropic regularizer with respect to the simplex $d(\mathbf{x}) := \sum_{i=1}^d x_i \log x_i$, which is 1-strongly convex with respect to the ℓ_1 norm. The pair of dual norms in the predictivity bound will therefore be $(\|\cdot\|_1, \|\cdot\|_\infty)$. We call this OFTRL setup *Optimistic Multiplicative Weights Updates* (OMWU).

Extensive-Form Correlated Equilibrium. We will work with the definition of EFCE due to Farina et al. (2019e), which is equivalent to that of von Stengel & Forges (2008). First, let us introduce the concept of a *trigger deviation function*.

Definition 2.4. Consider some player $i \in [n]$, a sequence $\hat{\sigma} = (j, a) \in \Sigma_*^{(i)}$, and joint sequence-form strategies $\pi \in \Pi_j^{(i)}$. A *trigger deviation function* with respect to a *trigger sequence* $\hat{\sigma}$ and *continuation strategy* $\hat{\pi}$ is any linear function $f : \mathbb{R}^{|\Sigma^{(i)}|} \rightarrow \mathbb{R}^{|\Sigma^{(i)}|}$ with the following properties.

- Any strategy $\pi \in \Pi^{(i)}$ which does not prescribe the sequence $\hat{\sigma}$ remains invariant. That is, $f(\pi) = \pi$ for any $\pi \in \Pi^{(i)}$ such that $\pi[\hat{\sigma}] = 0$;
- Otherwise, the prescribed sequence $\hat{\sigma} = (j, a)$ is modified so that the behavior at j , as well as all its descendants is replaced by the behavior specified by the continuation strategy:

$$f(\pi)[\sigma] = \begin{cases} \pi[\sigma] & \text{if } \sigma \not\succeq j; \\ \hat{\pi}[\sigma] & \text{if } \sigma \succeq j, \end{cases} \quad (7)$$

for all $\sigma \in \Sigma^{(i)}$ and $\pi \in \Pi^{(i)}$ such that $\pi[\hat{\sigma}] = 1$.

We will let $\Psi^{(i)} := \{\phi_{\hat{\sigma} \rightarrow \hat{\pi}}^{(i)} : \hat{\sigma} = (j, a) \in \Sigma_*^{(i)}, \hat{\pi} \in \Pi_j^{(i)}\}$ be the set of all possible linear mappings defining trigger deviation functions for player i . We are ready to introduce the concept of EFCE.

Definition 2.5 (EFCE). For $\epsilon \geq 0$, a probability distribution $\mu \in \Delta^{[\Pi]}$ is an ϵ -approximate EFCE if for every player $i \in [n]$ and every trigger deviation function $\phi_{\hat{\sigma} \rightarrow \hat{\pi}}^{(i)} \in \Psi^{(i)}$, it holds that

$$\mathbb{E}_{\pi \sim \mu} \left[u^{(i)} \left(\phi_{\hat{\sigma} \rightarrow \hat{\pi}}^{(i)}(\pi^{(i)}), \pi^{(-i)} \right) - u^{(i)}(\pi) \right] \leq \epsilon, \quad (8)$$

where $\pi = (\pi_1, \dots, \pi_n) \in \Pi$. A probability distribution $\mu \in \Delta^{[\Pi]}$ is an EFCE if it is a 0-EFCE.

Theorem 2.6 (Farina et al. (2021a)). For every player $i \in [n]$, let $\pi^{(i),1}, \dots, \pi^{(i),T} \in \Pi^{(i)}$ be a sequence of deterministic sequence-form strategies whose cumulative $\Psi^{(i)}$ -regret is $R^{(i),T}$ with respect to the sequence of linear utility functions

$$\ell^{(i),t} : \Pi^{(i)} \ni \pi^{(i)} \mapsto u^{(i)} \left(\pi^{(i)}, \pi^{(-i),t} \right). \quad (9)$$

Then, the empirical frequency of play $\mu \in \Delta^{[\Pi]}$ is an ϵ -EFCE, where $\epsilon := \frac{1}{T} \max_{i \in [n]} R^{(i),T}$.

3 ACCELERATING Φ -REGRET MINIMIZATION VIA OPTIMISM

In this section we develop a general template for accelerated Φ -regret minimization for general sets, and then we instantiate the template for dynamics for EFCE. Our approach combines a framework of Gordon et al. (2008) with the framework of stable-predictive (aka. optimistic) regret minimization. As in Gordon et al. (2008), in our template we combine 1) a regret minimizer that outputs a linear transformation $\phi^t \in \Phi$ at every time t , and 2) a fixed-point oracle for each $\phi^t \in \Phi$. However, in our framework, we further require that 2) is stable (in the sense of Definition 2.2). To achieve this, we will focus on regret minimizers that have the following property:

Definition 3.1. Consider a set of functions Φ such that $\phi(\mathcal{X}) \subseteq \mathcal{X}$ for all $\phi \in \Phi$, and a no-regret algorithm \mathcal{R}_Φ for the set of transformations Φ which returns a sequence $\phi^t \in \Phi$. We say that \mathcal{R}_Φ is *fixed point G -stable*, for $G \geq 0$, if the following conditions hold:

- Every ϕ^t admits a fixed point. That is, there exists $\mathbf{x}^t \in \mathcal{X}$ such that $\phi^t(\mathbf{x}^t) = \mathbf{x}^t$.
- For any \mathbf{x}^t such that $\mathbf{x}^t = \phi^t(\mathbf{x}^t)$, there exists \mathbf{x}^{t+1} with $\mathbf{x}^{t+1} = \phi^{t+1}(\mathbf{x}^{t+1})$ such that $\|\mathbf{x}^{t+1} - \mathbf{x}^t\| \leq G$.

We will show how to construct an accelerated Φ -regret minimizer starting from the following:

1. \mathcal{R}_Φ : A κ -stable (α, β) -predictive fixed point G -stable regret minimizer for Φ ;
2. $\text{STABLEFPORACLE}(\phi; \tilde{\mathbf{x}}, G, \epsilon)$: A *stable fixed point oracle* which returns a point $\mathbf{x} \in \mathcal{X}$ such that (i) $\|\phi(\mathbf{x}) - \mathbf{x}\| \leq \epsilon$, and (ii) $\|\mathbf{x} - \tilde{\mathbf{x}}\| \leq G$ (the existence of such a fixed point is guaranteed by the fixed point G -stability assumption for the regret minimizer).

Given these two components, our next theorem builds a stable-predictive Φ -regret minimizer.

Theorem 3.2 (Accelerated Φ -Regret Minimization). Consider a κ -stable (α, β) -predictive regret minimizer \mathcal{R}_Φ for a set of linear transformations Φ , with respect to the ℓ_1 norm $\|\cdot\|_1$. Moreover, assume that \mathcal{R}_Φ is fixed point G -stable with respect to Φ . Then, if we have access to a STABLEFPORACLE , we can construct a G -stable algorithm with Φ -regret R^T bounded as

$$R^T \leq \alpha(T) + 2\beta D_\ell^2 \kappa^2 T + 2\beta \sum_{t=2}^T \|\ell^t - \ell^{t-1}\|_\infty^2 + D_\ell \sum_{t=1}^T \epsilon_t, \quad (10)$$

where ϵ_t is the error of STABLEFPORACLE at time t , and D_ℓ is an upper bound on the ℓ_∞ norm of ℓ^t 's. It is also assumed that $\|\mathbf{x}\|_\infty \leq 1$ for all $\mathbf{x} \in \mathcal{X}$.

The proof is similar to that of Gordon et al. (2008), and is included in Appendix B.

3.1 CONSTRUCTING A STABLE-PREDICTIVE REGRET MINIMIZER FOR $\Psi^{(i)}$

Here we develop a regret minimizer for the set $\text{co } \Psi^{(i)}$, the convex hull of the set of trigger deviation functions. Given that $\text{co } \Psi^{(i)} \supseteq \Psi^{(i)}$, this will immediately imply a regret minimizer for the set

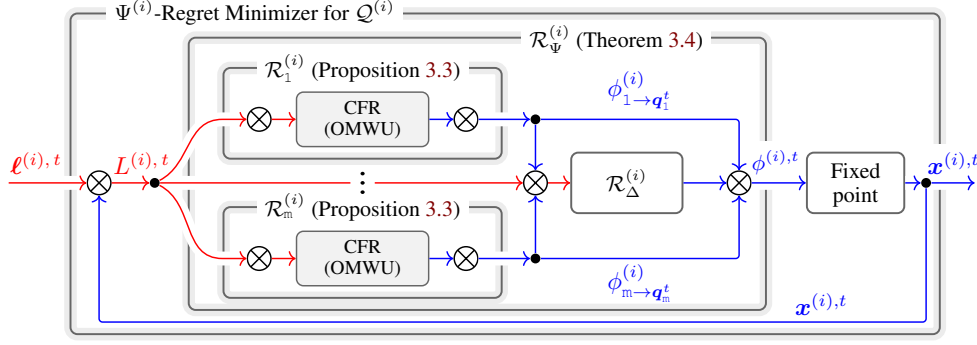


Figure 2: An overview of the overall construction. For notational convenience we have let $\Sigma_*^{(i)} := \{1, 2, \dots, m\}$. The symbol \otimes in the figure denotes a multilinear transformation of the inputs. We also note that blue corresponds to the iterates, while red corresponds to the utilities.

$\Psi^{(i)}$. An overview of the algorithm is given in Figure 2. Farina et al. (2021a) observed that the set $\text{co } \Psi^{(i)}$ can be evaluated in two stages. First, for a fixed sequence $\hat{\sigma} = (j, a) \in \Sigma_*^{(i)}$ we define the set $\Psi_{\hat{\sigma}}^{(i)} := \text{co} \left\{ \phi_{\hat{\sigma} \rightarrow \hat{\pi}} : \hat{\pi} \in \Pi_j^{(i)} \right\}$; then, we take the convex hull of all $\Psi_{\hat{\sigma}}^{(i)}$, that is, $\text{co } \Psi^{(i)} = \text{co} \{ \Psi_{\hat{\sigma}}^{(i)} : \hat{\sigma} \in \Sigma_*^{(i)} \}$. Correspondingly, we first develop a stable-predictive regret minimizer for the set $\Psi_{\hat{\sigma}}^{(i)}$, for any $\hat{\sigma} \in \Sigma_*^{(i)}$, and these individual regret minimizers are then combined using a *regret circuit* to conclude the construction in Theorem 3.4. All the omitted proofs and pseudocode for this section are included in Appendix B.1.

Stable-Predictive Regret Minimizer for the set $\Psi_{\hat{\sigma}}^{(i)}$. Consider a sequence $\hat{\sigma} \in \Sigma_*^{(i)}$. Farina et al. (2021a) observed that the set of transformations $\Psi_{\hat{\sigma}}^{(i)} := \text{co} \left\{ \phi_{\hat{\sigma} \rightarrow \hat{\pi}} : \hat{\pi} \in \Pi_j^{(i)} \right\}$ is the image of $\mathcal{Q}_j^{(i)}$ under the affine mapping $h_{\hat{\sigma}}^{(i)} : q \mapsto \phi_{\hat{\sigma} \rightarrow q}^{(i)}$. Hence, it is well-known that a regret minimizer for $\Psi_{\hat{\sigma}}^{(i)}$ can be constructed starting from a regret minimizer for $\mathcal{Q}_j^{(i)}$. We now show that the same can be said if one restricts to *stable-predictive* regret minimizers. In particular, we have the following.

Proposition 3.3. *Consider a player $i \in [n]$ and any trigger sequence $\hat{\sigma} = (j, a) \in \Sigma_*^{(i)}$. There exists an algorithm which constructs a deterministic regret minimizer $\mathcal{R}_{\hat{\sigma}}^{(i)}$ with access to a K -stable (A_T, B) -predictive deterministic regret minimizer $\mathcal{R}_{\mathcal{Q}}^{(i)}$ for the set $\mathcal{Q}_j^{(i)}$, such that $\mathcal{R}_{\hat{\sigma}}^{(i)}$ is K -stable and (A_T, B) -predictive.*

In Appendix A we describe a stable-predictive variant of CFR for the set $\mathcal{Q}_j^{(i)}$, for each $j \in \mathcal{J}^{(i)}$, following the construction of Farina et al. (2019a).

Stable-Predictive Regret Minimizer for $\text{co } \Psi^{(i)}$. The next step consists of combining the regret minimizers $\Psi_{\hat{\sigma}}^{(i)}$, for all $\hat{\sigma} \in \Sigma_*^{(i)}$, to a composite regret minimizer for the set $\text{co } \Psi^{(i)}$. To this end, we employ *regret circuits* (Farina et al., 2019d), leading to the main result of this section:

Theorem 3.4. *Consider a κ -stable (α, β) -predictive regret minimizer $\mathcal{R}_{\Delta}^{(i)}$ for the the simplex $\Delta^{|\Sigma_*^{(i)}|}$, and K -stable (A, B) -predictive regret minimizers $\mathcal{R}_{\hat{\sigma}}^{(i)}$ for each $\hat{\sigma} \in \Sigma_*^{(i)}$, all with respect to the pair of norms $(\|\cdot\|_1, \|\cdot\|_\infty)$. Then, there exists an algorithm which constructs a regret minimizer $\mathcal{R}_{\Psi}^{(i)}$ for the set $\text{co } \Psi^{(i)}$ such that (i) $\mathcal{R}_{\Psi}^{(i)}$ is $O(K + |\Sigma^{(i)}|\kappa)$ -stable, and (ii) under any sequence of linear utility functions L^1, \dots, L^T the regret incurred can be bounded as*

$$R_{\Psi}^T \leq O(\alpha(T) + A(T) + \beta D_{\mathbf{L}}^2 K^2 T) + O(B + \beta |\Sigma^{(i)}|^2) \sum_{t=2}^T \|\mathbf{L}^t - \mathbf{L}^{t-1}\|_\infty^2, \quad (11)$$

where $\|\mathbf{L}^t\|_\infty \leq D_{\mathbf{L}}$.

3.2 STABILITY OF THE FIXED POINTS

In this subsection we complete the construction of the $\Psi^{(i)}$ -regret minimizer by establishing a *stable* fixed point oracle for any $\phi \in \text{co } \Psi^{(i)}$. All of the proofs of this section are included in Appendix B.2.

Multiplicative Stability. A sequence $\{z^t\}$, with $z^t \in \mathbb{R}_{\geq 0}^d$, is said to be κ -multiplicative-stable if $(1 - \kappa)z_i^{t-1} \leq z_i^t \leq (1 + \kappa)z_i^{t-1}$, for any $i \in [d]$, and for all $t \geq 2$. Importantly, this notion of multiplicative stability is guaranteed by OMWU (see Lemma B.2). Thus, if $\mathfrak{D}^{(i)}$ is the depth of i 's actions and $D_{\mathbf{x}}^{(i)}$ is an upper bound on the ℓ_1 norm in the treeplex, we can show the following:

Lemma 3.5. *When each regret minimizer $\mathcal{R}_{\hat{\sigma}}^{(i)}$ is constructed using predictive CFR instantiated with OMWU with learning rate η (Theorem A.4) such that for all $\hat{\sigma} \in \Sigma_*^{(i)}$, the output sequence is $O(\eta(\mathfrak{D}^{(i)})^2 D_{\mathbf{x}}^{(i)} D_{\ell})$ -multiplicatively-stable. Moreover, if the regret minimizer $\mathcal{R}_{\Delta}^{(i)}$ is realized using OMWU with learning rate η , it will output an $O(\eta|\Sigma^{(i)}|D_{\ell})$ -multiplicatively-stable sequence.*

This characterization will be crucial for establishing the stability of the fixed points. In particular, following the approach of Farina et al. (2021a), let us introduce the following definitions:

Definition 3.6. Consider a player $i \in [n]$ and let $J \subseteq \mathcal{J}^{(i)}$ be a subset of i 's information sets. We say that J is a *trunk* of $\mathcal{J}^{(i)}$ if, for every $j \in J$, all predecessors of j are also in J .

Definition 3.7. Consider a player $i \in [n]$, a trunk $J \subseteq \mathcal{J}^{(i)}$, and $\phi \in \text{co } \Psi^{(i)}$. A vector $\mathbf{x} \in \mathbb{R}_{\geq 0}^{|\Sigma^{(i)}|}$ is a J -partial fixed point of ϕ if the following conditions hold:

- $\mathbf{x}[\emptyset] = 1$ and $\mathbf{x}[\sigma^{(i)}(j)] = \sum_{a \in \mathcal{A}(j)} \mathbf{x}[(j, a)]$, for all $j \in J$;
- $\phi(\mathbf{x})[\emptyset] = \mathbf{x}[\emptyset] = 1$, and $\phi(\mathbf{x})[(j, a)] = \mathbf{x}[(j, a)]$, for all $j \in J$, and $a \in \mathcal{A}(j)$.

An important property is that a J -partial fixed point can be efficiently “promoted” to a $J \cup \{j^*\}$ -partial fixed point by computing the stationary distribution of a certain Markov chain. However, a significant concern is whether this fixed point operation can potentially cause a substantial degradation in terms of stability. One of our key results is that the associated Markov chain has a particular structure, which enables us to substantially improve the stability bound and thereby obtain a polynomial degradation in stability. More precisely, this boils down to the following technical lemma.

Lemma 3.8. *Let \mathbf{M} and \mathbf{M}' be transition matrices of m -state Markov chains such that $\mathbf{M} = \mathbf{v}\mathbf{1}^\top + \mathbf{C}$ and $\mathbf{M}' = \mathbf{v}'\mathbf{1}^\top + \mathbf{C}'$, where $\mathbf{C}, \mathbf{C}', \mathbf{v}, \mathbf{v}'$ have strictly positive entries. Moreover, let π and π' be the (unique) stationary distributions of \mathbf{M} and \mathbf{M}' respectively. Then, if (i) the entries of the matrices \mathbf{C} and \mathbf{C}' are κ -multiplicatively-close, (ii) the entries of the vectors \mathbf{v} and \mathbf{v}' are γ -multiplicatively-close, and (iii) the sum of the entries of \mathbf{v} and \mathbf{v}' are κ -multiplicatively-close, then π and π' are $(\gamma + O(\kappa m))$ -multiplicatively-close, for a sufficiently small $\kappa = O(1/m)$.*

Using a slightly more general result (Corollary B.10), we manage to obtain the following:

Proposition 3.9. *Consider a player $i \in [n]$, and let $\phi = \sum_{\hat{\sigma} \in \Sigma_*^{(i)}} \lambda[\hat{\sigma}] \phi_{\hat{\sigma} \rightarrow \mathbf{q}_{\hat{\sigma}}}^{(i)}$ be a transformation in $\text{co } \Psi^{(i)}$ such that the sequence of λ^t 's and $\mathbf{q}_{\hat{\sigma}}^t$'s is κ -multiplicatively-stable, for all $\hat{\sigma} \in \Sigma_*^{(i)}$. If \mathbf{x}^t is a γ -multiplicatively-stable J -partial fixed point sequence, there is an algorithm which computes a $(J \cup \{j^*\})$ -partial fixed point $(\mathbf{x}^t)'$ of ϕ such that the sequence of $(\mathbf{x}^t)'$'s is $(\gamma + O(\kappa|\mathcal{A}(j^*)|))$ -multiplicatively-stable, for any sufficiently small $\kappa = O(1/|\mathcal{A}(j^*)|)$.*

Thus, using our technical lemma, we manage to bypass the substantial overhead of the term $\gamma|\mathcal{A}(j^*)|$, which would follow using techniques similar to Chen & Peng (2020). This turns out to be crucial for obtaining a polynomial dependence on the size of the game. Finally, we can inductively employ this proposition to show the overall stability of the fixed points:

Theorem 3.10. *Consider a player $i \in [n]$, and let $\phi = \sum_{\hat{\sigma} \in \Sigma_*^{(i)}} \lambda[\hat{\sigma}] \phi_{\hat{\sigma} \rightarrow \mathbf{q}_{\hat{\sigma}}}^{(i)}$ be a transformation in $\text{co } \Psi^{(i)}$ such that the sequence of λ^t 's and $\mathbf{q}_{\hat{\sigma}}^t$'s is κ -multiplicatively-stable, for all $\hat{\sigma} \in \Sigma_*^{(i)}$. Then, there exists an algorithm which computes a fixed point $\mathbf{q}^t \in \mathcal{Q}^{(i)}$ of ϕ such that the sequence of \mathbf{q}^t 's is $O(\kappa|\mathcal{A}^{(i)}|\mathfrak{D}^{(i)})$ -multiplicatively-stable, where $|\mathcal{A}^{(i)}| := \max_{j \in \mathcal{J}^{(i)}} |\mathcal{A}(j)|$, and for a sufficiently small $\kappa = O(1/(|\mathcal{A}^{(i)}|\mathfrak{D}^{(i)}))$.*

Finally, if we use the stability values derived in Lemma 3.5, we arrive at the following conclusion:

Corollary 3.11. For $\kappa = O((D_{\mathbf{x}}^{(i)}(\mathfrak{D}^{(i)})^2 + |\Sigma^{(i)}|)|\mathcal{A}^{(i)}|\mathfrak{D}^{(i)}D_{\ell})$, the sequence of fixed points will be $(\eta\kappa)$ -multiplicatively-stable, for any sufficiently small $\eta = O(1/\kappa)$.

Putting Everything Together. Finally, having established these ingredients, we can use the template of Theorem 3.2 to obtain Theorem 1.1, as we formally show in Appendix B.3.

4 EXPERIMENTS

In this section we experimentally investigate the performance of our stable-predictive algorithm compared to two other popular approaches based on a CFR-style decomposition of regrets into local regret-minimization problems: the existing algorithm by Farina et al. (2021a) instantiated with (i) *regret matching*⁺ (RM⁺) (Tammelin, 2014) for each simplex (in place of regret matching), and (ii) using the vanilla MWU algorithm for each simplex. In accordance to the theoretical predictions, the stepsize for OMWU is set as $\eta_t = \tau \cdot t^{-1/4}$ (cf. Corollary B.13), and for MWU it is set as $\eta_t = \tau \cdot t^{-1/2}$, where the parameter τ is chosen by picking the best-performing value among $\{0.01, 0.1, 1, 10, 100\}$. In particular, we evaluate their performance based on the following popular benchmark games: (i) a three-player variant of *Kuhn poker* (Kuhn, 1950); (ii) a two-player bargaining game known as *Sheriff* (Farina et al., 2019e)—a benchmark game introduced specifically for the study of correlated equilibria; and (iii) a three-player version of *Liar’s dice* (Lisý et al., 2015). A detailed description of each of the three game instances is available in Appendix D.

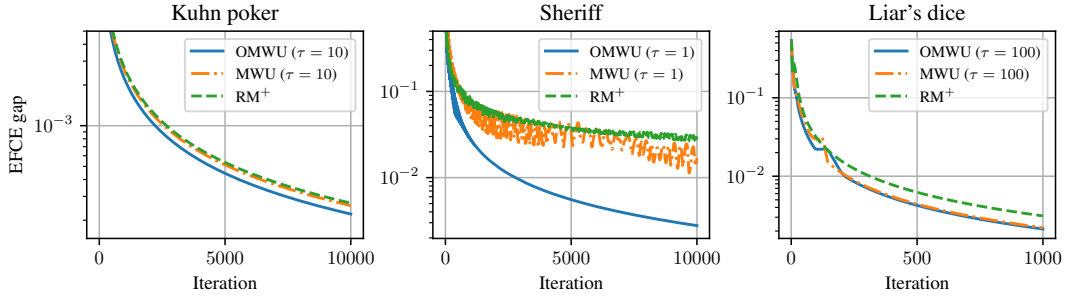


Figure 3: The performance of MWU, OMWU, and RM⁺ on three general-sum EFGs.

Figure 3 shows the performance of each of the three learning dynamics for computing EFCE. On the x -axis we plot the number of iterations performed by each algorithm, and on the y -axis we plot the EFCE gap, defined as the maximum advantage that any player can gain by defecting optimally from the mediator’s recommendations. **It should be noted that one iteration costs the same for every algorithm, up to constant factors.** We see that on every game, OMWU performs better than or on par with RM⁺ and MWU. On Sheriff, OMWU performs significantly better than both RM⁺ and MWU, by about an order of magnitude. One caveat to these results is that we did not use two tricks that help CFR⁺ in two-player zero-sum EFG solving: alternation and linear averaging. These tricks are known to retain convergence guarantees in that context (Tammelin et al., 2015; Farina et al., 2019b; Burch et al., 2019), but it is unclear if they still guarantee convergence in the EFCE setting.

5 CONCLUSIONS

We described uncoupled no-regret learning dynamics so that if all agents play T repetitions of the game according to the dynamics, the correlated distribution of play is an $O(T^{-3/4})$ -approximate EFCE. This substantially improves over the prior best rate of $O(T^{-1/2})$. One of our conceptual contributions is to connect the line of work on optimistic regret minimization with the framework of Φ -regret. One of our main technical contributions is to characterize the stability of the fixed points associated with trigger deviation functions through a refined perturbation analysis of a certain structured Markov chain, which may be of independent interest. Finally, experiments conducted on standard benchmarks corroborated our theoretical findings.

REFERENCES

- Robert Aumann. Subjectivity and correlation in randomized strategies. *Journal of Mathematical Economics*, 1:67–96, 1974.
- Kimmo Berg and Tuomas Sandholm. Exclusion method for finding nash equilibrium in multiplayer games. In *AAAI Conference on Artificial Intelligence (AAAI)*, 2017.
- Avrim Blum and Yishay Mansour. From external to internal regret. *J. Mach. Learn. Res.*, 8:1307–1324, 2007.
- Michael Bowling, Neil Burch, Michael Johanson, and Oskari Tammelin. Heads-up limit hold’em poker is solved. *Science*, 347(6218), January 2015.
- Noam Brown and Tuomas Sandholm. Superhuman AI for heads-up no-limit poker: Libratus beats top professionals. *Science*, pp. eaao1733, Dec. 2017.
- Neil Burch, Matej Moravcik, and Martin Schmid. Revisiting CFR+ and alternating updates. *Journal of Artificial Intelligence Research*, 64:429–443, 2019.
- Andrea Celli, Alberto Marchesi, Tommaso Bianchi, and Nicola Gatti. Learning to correlate in multi-player general-sum sequential games. In *Proceedings of the Annual Conference on Neural Information Processing Systems (NeurIPS)*, volume 32, 2019.
- Andrea Celli, Alberto Marchesi, Gabriele Farina, and Nicola Gatti. No-regret learning dynamics for extensive-form correlated equilibrium. In Hugo Larochelle, Marc’Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin (eds.), *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020*, 2020.
- Xi Chen and Binghui Peng. Hedging in games: Faster convergence of external and swap regrets. In *Proceedings of the Annual Conference on Neural Information Processing Systems (NeurIPS)*, 2020.
- Chao-Kai Chiang, Tianbao Yang, Chia-Jung Lee, Mehrdad Mahdavi, Chi-Jen Lu, Rong Jin, and Shenghuo Zhu. Online optimization with gradual variations. In *Conference on Learning Theory*, pp. 6–1, 2012.
- Constantinos Daskalakis, Paul Goldberg, and Christos Papadimitriou. The complexity of computing a Nash equilibrium. In *Proceedings of the Annual Symposium on Theory of Computing (STOC)*, 2006.
- Constantinos Daskalakis, Alan Deckelbaum, and Anthony Kim. Near-optimal no-regret algorithms for zero-sum games. *Games and Economic Behavior*, 92:327–348, 2015.
- Constantinos Daskalakis, Maxwell Fishelson, and Noah Golowich. Near-optimal no-regret learning in general games. *CoRR*, abs/2108.06924, 2021.
- Miroslav Dudík and Geoffrey J. Gordon. A sampling-based approach to computing equilibria in succinct extensive-form games. In Jeff A. Bilmes and Andrew Y. Ng (eds.), *UAI 2009, Proceedings of the Twenty-Fifth Conference on Uncertainty in Artificial Intelligence, Montreal, QC, Canada, June 18-21, 2009*, pp. 151–160. AUAI Press, 2009.
- Kousha Etessami and Mihalis Yannakakis. On the complexity of Nash equilibria and other fixed points (extended abstract). In *Proceedings of the Annual Symposium on Foundations of Computer Science (FOCS)*, pp. 113–123, 2007.
- Gabriele Farina, Christian Kroer, Noam Brown, and Tuomas Sandholm. Stable-predictive optimistic counterfactual regret minimization. In *International Conference on Machine Learning (ICML)*, 2019a.
- Gabriele Farina, Christian Kroer, and Tuomas Sandholm. Online convex optimization for sequential decision processes and extensive-form games. In *AAAI Conference on Artificial Intelligence*, 2019b.

- Gabriele Farina, Christian Kroer, and Tuomas Sandholm. Optimistic regret minimization for extensive-form games via dilated distance-generating functions. In *Advances in Neural Information Processing Systems, NeurIPS 2019*, pp. 5222–5232, 2019c.
- Gabriele Farina, Christian Kroer, and Tuomas Sandholm. Regret circuits: Composability of regret minimizers. In *International Conference on Machine Learning*, pp. 1863–1872, 2019d.
- Gabriele Farina, Chun Kai Ling, Fei Fang, and Tuomas Sandholm. Correlation in extensive-form games: Saddle-point formulation and benchmarks. In *Conference on Neural Information Processing Systems (NeurIPS)*, 2019e.
- Gabriele Farina, Andrea Celli, Alberto Marchesi, and Nicola Gatti. Simple uncoupled no-regret learning dynamics for extensive-form correlated equilibrium, 2021a.
- Gabriele Farina, Andrea Celli, and Tuomas Sandholm. Efficient decentralized learning dynamics for extensive-form coarse correlated equilibrium: No expensive computation of stationary distributions required. ArXiv preprint, 2021b.
- Gabriele Farina, Christian Kroer, and Tuomas Sandholm. Better regularization for sequential decision spaces: Fast convergence rates for Nash, correlated, and team equilibria. In *ACM Conference on Economics and Computation*, 2021c.
- Andrew Gilpin and Tuomas Sandholm. Lossless abstraction of imperfect information games. *Journal of the ACM*, 54(5), 2007.
- Geoffrey J Gordon, Amy Greenwald, and Casey Marks. No-regret learning in convex games. In *Proceedings of the 25th international conference on Machine learning*, pp. 360–367. ACM, 2008.
- Amy Greenwald and Amir Jafari. A general class of no-regret learning algorithms and game-theoretic equilibria. In *Conference on Learning Theory (COLT)*, Washington, D.C., 2003.
- Sergiu Hart and Andreu Mas-Colell. A simple adaptive procedure leading to correlated equilibrium. *Econometrica*, 68:1127–1150, 2000.
- Samid Hoda, Andrew Gilpin, Javier Peña, and Tuomas Sandholm. Smoothing techniques for computing Nash equilibria of sequential games. *Mathematics of Operations Research*, 35(2), 2010.
- Albert Xin Jiang and Kevin Leyton-Brown. Polynomial-time computation of exact correlated equilibrium in compact games. *Games Econ. Behav.*, 91:347–359, 2015.
- Christian Kroer, Kevin Waugh, Fatma Kılınç-Karzan, and Tuomas Sandholm. Faster algorithms for extensive-form game solving via improved smoothing functions. *Mathematical Programming*, 2020.
- Alex Kruckman, Amy Greenwald, and John R. Wicks. An elementary proof of the Markov chain tree theorem. Technical Report 10-04, Brown University, 2010.
- H. W. Kuhn. A simplified two-person poker. In H. W. Kuhn and A. W. Tucker (eds.), *Contributions to the Theory of Games*, volume 1 of *Annals of Mathematics Studies*, 24, pp. 97–103. Princeton University Press, Princeton, New Jersey, 1950.
- Viliam Lisý, Marc Lanctot, and Michael Bowling. Online Monte Carlo counterfactual regret minimization for search in imperfect information games. In *Autonomous Agents and Multi-Agent Systems*, pp. 27–36, 2015.
- Dustin Morrill, Ryan D’Orazio, Marc Lanctot, James R. Wright, Michael Bowling, and Amy R. Greenwald. Efficient deviation types and learning for hindsight rationality in extensive-form games. In Marina Meila and Tong Zhang (eds.), *Proceedings of the 38th International Conference on Machine Learning, ICML 2021*, volume 139 of *Proceedings of Machine Learning Research*, pp. 7818–7828. PMLR, 2021a.

- Dustin Morrill, Ryan D’Orazio, Reza Sarfati, Marc Lanctot, James R. Wright, Amy R. Greenwald, and Michael Bowling. Hindsight and sequential rationality of correlated play. In *Thirty-Fifth AAAI Conference on Artificial Intelligence, AAAI 2021*, pp. 5584–5594. AAAI Press, 2021b.
- Yurii Nesterov. Excessive gap technique in nonsmooth convex minimization. *SIAM Journal of Optimization*, 16(1), 2005.
- Christos H. Papadimitriou and Tim Roughgarden. Computing correlated equilibria in multi-player games. *J. ACM*, 55(3):14:1–14:29, 2008.
- Alexander Rakhlin and Karthik Sridharan. Online learning with predictable sequences. In *Conference on Learning Theory*, pp. 993–1019, 2013a.
- Alexander Rakhlin and Karthik Sridharan. Optimization, learning, and games with predictable sequences. In *Advances in Neural Information Processing Systems*, pp. 3066–3074, 2013b.
- Vasilis Syrgkanis, Alekh Agarwal, Haipeng Luo, and Robert E Schapire. Fast convergence of regularized learning in games. In *Advances in Neural Information Processing Systems*, pp. 2989–2997, 2015.
- Oskari Tammelin. Solving large imperfect information games using CFR+. arXiv preprint, 2014.
- Oskari Tammelin, Neil Burch, Michael Johanson, and Michael Bowling. Solving heads-up limit Texas hold’em. In *Proceedings of the 24th International Joint Conference on Artificial Intelligence (IJCAI)*, 2015.
- Bernhard von Stengel and Françoise Forges. Extensive-form correlated equilibrium: Definition and computational complexity. *Mathematics of Operations Research*, 33(4):1002–1022, 2008.
- Martin Zinkevich, Michael Bowling, Michael Johanson, and Carmelo Piccione. Regret minimization in games with incomplete information. In *Proceedings of the Annual Conference on Neural Information Processing Systems (NIPS)*, 2007.

A SEQUENTIAL DECISION MAKING AND STABLE-PREDICTIVE CFR

In this section we introduce the basic setting of *sequential decision making*, and we subsequently provide a stable-predictive variant of CFR following the construction of Farina et al. (2019a). In particular, the main result is captured in Theorem A.4.

A sequential decision process can be represented using a tree consisting of two types of nodes: *decision nodes* and *observation nodes*. The set of all decision nodes will be denoted with \mathcal{J} , while the set of observation nodes with \mathcal{K} . At every decision node $j \in \mathcal{J}$ the agent has to select a strategy in the form of a probability distribution over all possible actions $\mathcal{A}(j)$. On the other hand, the agent receives some type of feedback at each observation point $k \in \mathcal{K}$. More precisely, it is assumed that the agent may receive a signal from a set $\mathcal{S}(k)$.

Now at each decision point $j \in \mathcal{J}$ of the sequential decision process, the strategy $\mathbf{x}_j \in \Delta^{|\mathcal{A}(j)|}$ obtains a (linear) utility of the form $\langle \ell_j, \mathbf{x}_j \rangle$, for some utility vector ℓ_j . The expected utility throughout the entire decision process can be expressed as $\sum_{j \in \mathcal{J}} \pi_j \langle \ell_j, \mathbf{x}_j \rangle$, where π_j is the probability that the agent reaches the decision point j . Before we proceed with the representation of the strategies, it is important to point out that in all extensive-form games of *perfect recall* the agents face a sequential decision process.

Decomposition of Sequence-Form Representation. Our construction will rely on the following recursive decomposition of the sequence-form space \mathcal{X}^Δ :

- Consider an observation node $k \in \mathcal{K}$, and let \mathcal{C}_k be the children decision points of k . Then, \mathcal{X}_k^Δ can be decomposed as the following Cartesian product:

$$\mathcal{X}_k^\Delta := \bigtimes_{j \in \mathcal{C}_k} \mathcal{X}_j^\Delta; \quad (12)$$

- Consider a decision point $j \in \mathcal{J}$, and let $\mathcal{C}_j = \{k_1, \dots, k_{n_j}\}$ be the children observation points of j , with $n_j = |\mathcal{A}(j)|$. Then, \mathcal{X}_j^Δ can be decomposed as follows:

$$\mathcal{X}_j^\Delta := \left\{ \begin{pmatrix} \lambda_1 \\ \vdots \\ \lambda_{n_j} \\ \lambda_1 \mathbf{x}_1 \\ \vdots \\ \lambda_{n_j} \mathbf{x}_{n_j} \end{pmatrix} : (\lambda_1, \dots, \lambda_{n_j}) \in \Delta^{n_j}, \mathbf{x}_1 \in \mathcal{X}_{k_1}^\Delta, \dots, \mathbf{x}_{n_j} \in \mathcal{X}_{k_{n_j}}^\Delta \right\}. \quad (13)$$

In view of this decomposition, our regret minimizer for the sequence-form strategy space \mathcal{X}^Δ will be established based on *localized* regret minimizers. Specifically, the basic ingredients for the overall construction are given in Proposition A.1 and Proposition A.2. We should note that the stable-predictive property will be tacitly considered with respect to the pair of norms $(\|\cdot\|_1, \|\cdot\|_\infty)$.

Proposition A.1. *Consider an observation node $k \in \mathcal{K}$, and assume access to a κ_j -multiplicatively-stable (α_j, β_j) -predictive regret minimizer \mathcal{R}_j^Δ over the sequence-form strategy space \mathcal{X}_j^Δ , for each $j \in \mathcal{C}_k$. Then, we can construct a κ -multiplicatively-stable (A, B) -predictive regret minimizer \mathcal{R}_k^Δ for the sequence-form strategy space \mathcal{X}_k^Δ , where $A = \sum_{j \in \mathcal{C}_k} \alpha_j$ and $B = \sum_{j \in \mathcal{C}_k} \beta_j$.*

Proof. Given the decomposition of Equation (12), the composite regret minimizer can be constructed based on the *regret circuit* for the Cartesian product given by Farina et al. (2019d). In particular, it is direct to verify that $R_k^{\Delta, T} = \sum_{j \in \mathcal{C}_k} R_j^{\Delta, T}$, where $R_k^{\Delta, T}$ is the regret accumulated by the composite regret minimizer, and $R_j^{\Delta, T}$ the cumulative regret of each individual regret minimizer. In particular, by assumption we know that

$$R_j^{\Delta, T} \leq \alpha_j + \beta_j \sum_{t=2}^T \|\ell_j^{\Delta, t} + \ell_j^{\Delta, t-1}\|_\infty^2.$$

As a result, we can conclude that

$$R_k^{\Delta,T} \leq \left(\sum_{j \in \mathcal{C}_k} \alpha_j \right) + \left(\sum_{j \in \mathcal{C}_k} \beta_j \right) \sum_{t=2}^T \|\ell_k^{\Delta,t} - \ell_k^{\Delta,t-1}\|_\infty^2,$$

where we used that $\|\ell_j^{\Delta,t} - \ell_j^{\Delta,t-1}\|_\infty \leq \|\ell_k^{\Delta,t} - \ell_k^{\Delta,t-1}\|_\infty$. Finally, the κ -multiplicative-stability of \mathcal{R}_k^Δ follows directly from the κ -multiplicative-stability of each \mathcal{R}_j^Δ . \square

Proposition A.2. Consider a decision node $j \in \mathcal{J}$, and assume access to a K -multiplicatively-stable (α_k, β_k) -predictive regret minimizer \mathcal{R}_k^Δ over the sequence-form strategy space \mathcal{X}_k^Δ , for each $k \in \mathcal{C}_j$. Moreover, assume access to a κ -multiplicatively-stable (α, β) -predictive regret minimizer \mathcal{R}_Δ over the simplex $\Delta^{|\mathcal{A}(j)|}$. Then, for a sufficiently small constant κ , we can construct an $(2\kappa + K)$ -multiplicatively-stable (A, B) -predictive regret minimizer \mathcal{R}_j^Δ for the sequence-form strategy space \mathcal{X}_j^Δ , where

$$\begin{aligned} A(T) &= \alpha + \max_{k \in \mathcal{C}_j} \{\alpha_k\} + 2\beta D_\ell^2 D_x^2 K^2 T; \\ B &= \max_{k \in \mathcal{C}_j} \{\beta_k\} + 2\beta D_x^2, \end{aligned} \quad (14)$$

with D_ℓ being an upper bound on the ℓ_∞ norm of the sequence of ℓ^t 's, and D_x an upper bound on the ℓ_1 norm of all $\mathbf{x} \in \mathcal{X}^\Delta$.

Proof. For this construction we will use the regret circuit for the convex hull, presented in Proposition B.1. In particular, first note that by assumption the regret $R_k^{\Delta,T}$ accumulated by each regret minimizer \mathcal{R}_k^Δ can be bounded as

$$R_k^{\Delta,T} \leq \alpha_k + \beta_k \sum_{t=2}^T \|\ell_k^{\Delta,t} - \ell_k^{\Delta,t-1}\|_\infty^2.$$

Moreover, by construction each regret minimizer \mathcal{R}_k^Δ receives the same utility as \mathcal{R}_j^Δ ; this, along with the guarantee of Proposition B.1, imply that

$$R_j^{\Delta,T} \leq \alpha + \max_{k \in \mathcal{C}_j} \{\alpha_k\} + \max_{k \in \mathcal{C}_j} \{\beta_k\} \sum_{t=2}^T \|\ell_j^{\Delta,t} - \ell_j^{\Delta,t-1}\|_\infty^2 + \beta \sum_{t=2}^T \|\mathbf{\Lambda}^t - \mathbf{\Lambda}^{t-1}\|_\infty^2, \quad (15)$$

where $\mathbf{\Lambda}^t$ represents the utility function received as input by \mathcal{R}_Δ . Next, similarly to the analysis of Theorem 3.4 we can deduce that for some $k \in \mathcal{C}_j$,

$$\begin{aligned} \|\mathbf{\Lambda}^t - \mathbf{\Lambda}^{t-1}\|_\infty^2 &\leq 2\|\ell^t - \ell^{t-1}\|_\infty^2 \|\mathbf{x}_k^t\|_1^2 + 2\|\ell^{t-1}\|_\infty^2 \|\mathbf{x}_k^t - \mathbf{x}_k^{t-1}\|_1^2 \\ &\leq 2D_x^2 \|\ell^t - \ell^{t-1}\|_\infty^2 + 2D_\ell^2 D_x^2 K^2, \end{aligned}$$

where we used that $\|\mathbf{x}_k^t\|_1 \leq D_x$, and the fact that $\|\mathbf{x}_k^t - \mathbf{x}_k^{t-1}\|_1 \leq K \|\mathbf{x}_k^{t-1}\|_1 \leq K D_x$ by K -multiplicative-stability of the sequence of \mathbf{x}_k^t 's. As a result, if we plug-in this bound to Equation (15) we can conclude that

$$R_j^{\Delta,T} \leq \left(\alpha + \max_{k \in \mathcal{C}_j} \{\alpha_k\} + 2\beta D_\ell^2 D_x^2 K^2 T \right) + \left(\max_{k \in \mathcal{C}_j} \{\beta_k\} + 2\beta D_x^2 \right) \sum_{t=2}^T \|\ell_j^{\Delta,t} - \ell_j^{\Delta,t-1}\|_\infty^2.$$

Finally, the $(2\kappa + K)$ -multiplicative-stability of \mathcal{R}_j^Δ can be directly verified from the decomposition given in Equation (13). \square

Remark A.3. Given the decomposition provided in Equation (13), the regret circuit for the convex hull should operate on the appropriate “lifted” subspace for each \mathcal{X}_k^Δ , which does not essentially alter the argument since the augmented entries remain invariant; see (Farina et al., 2019d, Figure 7).

Finally, we can inductively combine Proposition A.1 and Proposition A.2 in order to establish a stable-predictive variant of CFR:

Theorem A.4 (Optimistic CFR). *If every local regret minimizer \mathcal{R}_j^Δ is updated using OMWU with learning rate η , for each $j \in \mathcal{J}$, we can construct an $O(\eta \mathfrak{D} D_{\mathbf{x}}^2 D_{\ell})$ -stable (A_T, B) -predictive regret minimizer \mathcal{R}^Δ for the space of sequence-form strategies \mathcal{X}^Δ , such that*

$$\begin{aligned} A_T &= O\left(\frac{\log |\mathcal{A}|}{\eta} |\mathcal{J}| + \eta^3 D_{\ell}^4 D_{\mathbf{x}}^6 \mathfrak{D}^2 T |\mathcal{J}|\right); \\ B &= O(\eta \mathfrak{D} D_{\mathbf{x}}^2 |\mathcal{J}|), \end{aligned} \quad (16)$$

where $|\mathcal{A}| := \max_{j \in \mathcal{J}} |\mathcal{A}(j)|$, D_{ℓ} is an upper bound on the ℓ_{∞} norm of the utility functions, $D_{\mathbf{x}}$ is an upper bound on the ℓ_1 norm of all $\mathbf{x} \in \mathcal{X}^\Delta$, and \mathfrak{D} is the depth of the decision process. Moreover, \mathcal{R}^Δ is $O(\eta \mathfrak{D} D_{\mathbf{x}} D_{\ell})$ -multiplicatively-stable.

Proof. First of all, it is easy to see that all the (localized) *counterfactual* losses have ℓ_{∞} norm bounded by $O(D_{\mathbf{x}} D_{\ell})$, where recall that $D_{\mathbf{x}}$ is an upper bound on the ℓ_1 norm of all $\mathbf{x} \in \mathcal{X}^\Delta$. As a result, we know from Lemma B.2 that the output of each local regret minimizer \mathcal{R}_j^Δ under OMWU with learning rate η is $O(\eta D_{\mathbf{x}} D_{\ell})$ -multiplicatively-stable. Along with Proposition A.2, we can inductively infer that the output of \mathcal{R} is $O(\eta \mathfrak{D} D_{\mathbf{x}} D_{\ell})$ -multiplicatively-stable, where \mathfrak{D} is the depth of the decision process. Moreover, we can conclude from this property that \mathcal{R} is $O(\eta \mathfrak{D} D_{\mathbf{x}}^2 D_{\ell})$ -stable (in the additive sense of Definition 2.2) given that the ℓ_1 norm of $\mathbf{x} \in \mathcal{X}^\Delta$ is bounded by $D_{\mathbf{x}}$.

For the predictivity bound, first note that the range of the entropic regularizer on the m -dimensional simplex is $\log m$. Thus, by Lemma 2.3 we know that each local regret minimizer at the information set $j \in \mathcal{J}$ instantiated with OMWU with learning rate η will be $(\log(|\mathcal{A}_j|/\eta, \eta)$ -predictive. Moreover, we previously argued that the stability parameter K appearing in the regret bound of Proposition A.2 will be bounded by $O(\eta \mathfrak{D} D_{\mathbf{x}}^2 D_{\ell})$. As a result, our predictivity bound follows directly from Proposition A.2. \square

Naturally, the same bounds apply for constructing a regret minimizer for the subspace \mathcal{X}_j^Δ , for any $j \in \mathcal{J}$, as required in Proposition 3.3.

B PROOFS

In this section we include all of the omitted proofs. First, let us introduce some additional useful notation. In particular, it will be convenient to instantiate a trigger deviation function in the form of a linear mapping $\phi_{\hat{\sigma} \rightarrow \hat{\pi}}^{(i)} : \mathbb{R}^{|\Sigma^{(i)}|} \ni \mathbf{x} \mapsto \mathbf{M}_{\hat{\sigma} \rightarrow \hat{\pi}}^{(i)} \mathbf{x}$, where $\mathbf{M}_{\hat{\sigma} \rightarrow \hat{\pi}}^{(i)}$ is such that for any $\sigma_r, \sigma_c \in \Sigma^{(i)}$,

$$\mathbf{M}_{\hat{\sigma} \rightarrow \hat{\pi}}^{(i)}[\sigma_r, \sigma_c] = \begin{cases} 1 & \text{if } \sigma_c \not\preceq \hat{\sigma} \quad \& \quad \sigma_r = \sigma_c; \\ \hat{\pi}[\sigma_r] & \text{if } \sigma_c = \hat{\sigma} \quad \& \quad \sigma_r \succeq j; \\ 0 & \text{otherwise,} \end{cases} \quad (17)$$

where $\hat{\sigma} = (j, a) \in \Sigma_*^{(i)}$. It is not hard to show that the linear mapping described in Equation (17) is indeed a trigger deviation function in the sense of Definition 2.4. We will also sometimes use the notation $\mathbf{x} \otimes \mathbf{y} = \mathbf{x} \mathbf{y}^\top$ to denote the *outer product* of the vectors \mathbf{x} and \mathbf{y} . Moreover, we will write $(\mathbf{M})^b$ to represent the *vectorization* of the matrix \mathbf{M} .

Proof of Theorem 3.2. Fix some iteration $t \geq 2$. The first step is to determine the next strategy of \mathcal{R}_Φ : $\phi^t = \mathcal{R}_\Phi.\text{NEXTSTRATEGY}()$. Then, our regret minimizer \mathcal{R} will simply output the strategy \mathbf{x}^t such that $\mathbf{x}^t = \text{STABLEFPORACLE}(\phi^t; \mathbf{x}^{t-1}, G, \epsilon_t)$.¹ By assumption (recall Definition 3.1) we know that this is indeed well-defined, and \mathbf{x}^t will be such that (i) $\|\phi^t(\mathbf{x}^t) - \mathbf{x}^t\| \leq \epsilon_t$, and (ii) $\|\mathbf{x}^t - \mathbf{x}^{t-1}\| \leq G$. This immediately implies that \mathcal{R} will be G -stable.

Afterwards, we receive feedback from the environment in the form of a utility vector ℓ^t . This vector is used to construct the utility function $L^t : \phi \mapsto \langle \ell^t, \phi(\mathbf{x}^t) \rangle$, with \mathbf{L}^t being the corresponding canonical vector of this linear mapping. Then, this function is given as feedback to \mathcal{R}_Φ ; that is, we

¹For $t = 1$ it suffices to return any \mathbf{x}^t such that $\mathbf{x}^1 = \phi^1(\mathbf{x}^1)$.

invoke the subroutine \mathcal{R}_Φ . OBSERVEUTILITY(L^t). As a result, the (external) regret of \mathcal{R}_Φ can be expressed as

$$R_\Phi^T = \max_{\phi^* \in \Phi} \sum_{t=1}^T \langle \ell^t, \phi^*(\mathbf{x}^t) \rangle - \sum_{t=1}^T \langle \ell^t, \phi^t(\mathbf{x}^t) \rangle. \quad (18)$$

In particular, if R^T is the Φ -regret of \mathcal{R} , we have that

$$\begin{aligned} R^T - R_\Phi^T &= \sum_{t=1}^T \langle \ell^t, \phi^t(\mathbf{x}^t) \rangle - \sum_{t=1}^T \langle \ell^t, \mathbf{x}^t \rangle = \sum_{t=1}^T \langle \ell^t, \phi^t(\mathbf{x}^t) - \mathbf{x}^t \rangle \\ &\leq \sum_{t=1}^T \|\ell^t\|_* \|\phi^t(\mathbf{x}^t) - \mathbf{x}^t\| \leq D_\ell \sum_{t=1}^T \epsilon_t, \end{aligned} \quad (19)$$

where we used the Cauchy-Schwarz inequality, as well as the fact that $\|\phi^t(\mathbf{x}^t) - \mathbf{x}^t\| \leq \epsilon_t$. Next, we will bound the factor $\|\mathbf{L}^t - \mathbf{L}^{t-1}\|_\infty$ arising from the stable-predictive bound on R_Φ^T in terms of $\|\ell^t - \ell^{t-1}\|_\infty$. In particular, it follows that

$$\begin{aligned} \|\mathbf{L}^t - \mathbf{L}^{t-1}\|_\infty^2 &= \|(\ell^t \otimes \mathbf{x}^t)^\flat - (\ell^{t-1} \otimes \mathbf{x}^{t-1})^\flat\|_\infty^2 \\ &= \|(\ell^t \otimes \mathbf{x}^t)^\flat - (\ell^{t-1} \otimes \mathbf{x}^t)^\flat + (\ell^{t-1} \otimes \mathbf{x}^t)^\flat - (\ell^{t-1} \otimes \mathbf{x}^{t-1})^\flat\|_\infty^2 \\ &= \|((\ell^t - \ell^{t-1}) \otimes \mathbf{x}^t)^\flat + (\ell^{t-1} \otimes (\mathbf{x}^t - \mathbf{x}^{t-1}))^\flat\|_\infty^2 \\ &\leq 2\|((\ell^t - \ell^{t-1}) \otimes \mathbf{x}^t)^\flat\|_\infty^2 + 2\|(\ell^{t-1} \otimes (\mathbf{x}^t - \mathbf{x}^{t-1}))^\flat\|_\infty^2 \end{aligned} \quad (20)$$

$$= 2\|\ell^t - \ell^{t-1}\|_\infty^2 \|\mathbf{x}^t\|_\infty^2 + 2\|\ell^{t-1}\|_\infty^2 \|\mathbf{x}^t - \mathbf{x}^{t-1}\|_\infty^2 \quad (21)$$

$$\leq 2\|\ell^t - \ell^{t-1}\|_\infty^2 + 2D_\ell^2 \kappa^2, \quad (22)$$

where we used the triangle inequality together with Young's inequality in (20), the property that $\|(\mathbf{w} \otimes \mathbf{z})^\flat\|_\infty = \|\mathbf{w}\|_\infty \|\mathbf{z}\|_\infty$ in (21), and the stability property $\|\mathbf{x}^t - \mathbf{x}^{t-1}\|_\infty \leq \|\mathbf{x}^t - \mathbf{x}^{t-1}\|_1 \leq \kappa$ in (22). As a result, if we plug in the derived bound in Equation (19) we can conclude that

$$\begin{aligned} R^T &\leq \alpha_T + D_\ell \sum_{t=1}^T \epsilon_t + \beta \sum_{t=1}^T (2\|\ell^t - \ell^{t-1}\|_\infty^2 + 2D_\ell^2 \kappa^2) \\ &= \alpha_T + 2\beta D_\ell^2 \kappa^2 T + D_\ell \sum_{t=1}^T \epsilon_t + 2\beta \sum_{t=1}^T \|\ell^t - \ell^{t-1}\|_\infty^2, \end{aligned}$$

concluding the proof. \square

B.1 PROOFS FOR SECTION 3.1

In this section we present all the omitted proofs from Section 3.1. Specifically, we commence with the proof of Proposition 3.3, established via Algorithm 2. We note that a similar construction appears in (Farina et al., 2021a).

Proof of Proposition 3.3. As suggested in Algorithm 1, let us consider the linear function $g_\sigma^{(i),t} : \mathbb{R}^{|\Sigma_j^{(i)}|} \ni \mathbf{x} \mapsto L^t(h_\sigma^{(i)}(\mathbf{x})) - L(h_\sigma^{(i)}(\mathbf{0}))$, with $\mathbf{g}_\sigma^{(i)} \in \mathbb{R}^{|\Sigma_j^{(i)}|}$ being the corresponding canonical vector. In Algorithm 1 the observed utility function L^t at time t is used to construct $g_\sigma^{(i),t}$, and the latter linear function is given as input to $\mathcal{R}_Q^{(i)}$. As a result, it follows that

$$\sup_{\phi^* \in \Psi_\sigma^{(i)}} \sum_{t=1}^T L^t(\phi^*) - \sum_{t=1}^T L^t(\phi_{\sigma \rightarrow \mathbf{q}_\sigma^*}^{(i)}) = \sup_{\mathbf{q}_\sigma^* \in \mathcal{Q}_j^{(i)}} \sum_{t=1}^T g_\sigma^{(i),t}(\mathbf{q}_\sigma^*) - \sum_{t=1}^T g_\sigma^{(i),t}(\mathbf{q}_\sigma^t). \quad (23)$$

That is, the cumulative regret incurred by $\mathcal{R}_\sigma^{(i)}$ under the sequence of utility functions L^1, \dots, L^T is equal to the regret incurred by $\mathcal{R}_Q^{(i)}$ under the sequence of utility functions $g_\sigma^{(i),t}$. As a result, if

we invoke the (A_T, B) -predictive assumption for the regret minimizer $\mathcal{R}_Q^{(i)}$, we can infer that the cumulative regret R^T of $\mathcal{R}_\sigma^{(i)}$ can be bounded as

$$R^T \leq A_T + B \sum_{t=2}^T \|\mathbf{g}_\sigma^{(i),t} - \mathbf{g}_\sigma^{(i),t-1}\|_\infty^2 \leq A_T + B \sum_{t=2}^T \|\mathbf{L}^t - \mathbf{L}^{t-1}\|_\infty^2, \quad (24)$$

where the bound follows from the fact that $\mathbf{g}_\sigma^{(i),t} = (\mathbf{L}^t[\sigma_r, \hat{\sigma}])_{\sigma_r \succeq j}$. It is also direct to verify that the stability is preserved since $\|(\mathbf{M}_{\hat{\sigma} \rightarrow \mathbf{q}_\sigma^t}^{(i)})^b - (\mathbf{M}_{\hat{\sigma} \rightarrow \mathbf{q}_\sigma^{t-1}}^{(i)})^b\|_1 = \|\mathbf{q}_\sigma^t - \mathbf{q}_\sigma^{t-1}\|_1 \leq K$, by K -stability of $\mathcal{R}_Q^{(i)}$. Finally, the claim regarding the complexity of Algorithm 1 follows directly given that we can store the vector $\mathbf{g}_\sigma^{(i),t}$ in $O(|\Sigma_j^{(i)}|)$ time. \square

Algorithm 1: Stable-Predictive Regret Minimizer $\mathcal{R}_\sigma^{(i)}$ for the set $\Psi_\sigma^{(i)}$

Input:

- Player $i \in [n]$
- A trigger sequence $\hat{\sigma} = (j, a) \in \Sigma_*^{(i)}$
- A (deterministic) K -stable (A_T, B) -predictive regret minimizer $\mathcal{R}_Q^{(i)}$ for $\mathcal{Q}_j^{(i)}$

1 **function** NEXTSTRATEGY():

2 $\mathbf{q}_\sigma^t \leftarrow \mathcal{R}_Q^{(i)} \cdot \text{NEXTSTRATEGY}()$

3 **return** $\phi_{\hat{\sigma} \leftarrow \mathbf{q}_\sigma^t}^{(i)}$

4 **function** OBSERVEUTILITY(L^t):

5 Construct the linear function $g_\sigma^{(i),t} : \mathbb{R}^{|\Sigma_j^{(i)}|} \ni \mathbf{x} \mapsto L^t(h_\sigma^{(i)}(\mathbf{x})) - L^t(h_\sigma^{(i)}(\mathbf{0}))$

6 $\mathcal{R}_Q^{(i)} \cdot \text{OBSERVEUTILITY}(g_\sigma^{(i),t})$

Next, we conclude the construction by combining the individual regret minimizers via a regret circuit for the convex hull; to be more precise, we will employ the following guarantee.

Proposition B.1 (Farina et al. (2019d)). *Consider a collection of sets $\mathcal{X}_1, \dots, \mathcal{X}_m$, and let \mathcal{R}_i be a regret minimizer for the set \mathcal{X}_i , for each $i \in [m]$. Moreover, let \mathcal{R}_Δ be a regret minimizer for the m -simplex Δ^m . A regret minimizer \mathcal{R}_{co} for the set $\text{co}\{\mathcal{X}_1, \dots, \mathcal{X}_m\}$ can be constructed as follows:*

- $\mathcal{R}_{\text{co}} \cdot \text{NEXTSTRATEGY}$ obtains the next strategy \mathbf{x}_i^t of each regret minimizer \mathcal{R}_i , as well as the next strategy $\boldsymbol{\lambda}^t = (\lambda_1^t, \dots, \lambda_m^t) \in \Delta^m$ of \mathcal{R}_Δ , and returns the strategy $\lambda_1^t \mathbf{x}_1^t + \dots + \lambda_m^t \mathbf{x}_m^t$.
- $\mathcal{R}_{\text{co}} \cdot \text{OBSERVEUTILITY}(L^t)$ forwards the function L^t to each of the regret minimizers $\mathcal{R}_1, \dots, \mathcal{R}_m$, while it also forwards the utility $(\lambda_1, \dots, \lambda_m) \mapsto \lambda_1 L^t(\mathbf{x}_1^t) + \dots + \lambda_m L^t(\mathbf{x}_m^t)$.

Then, if R_1^T, \dots, R_m^T is the regret accumulated by the regret minimizers $\mathcal{R}_1, \dots, \mathcal{R}_m$, and R_Δ^T is the regret of \mathcal{R}_Δ , then the cumulative regret of the composite regret minimizers \mathcal{R}_{co} can be bounded as

$$R_{\text{co}}^T \leq R_\Delta^T + \max\{R_1^T, \dots, R_m^T\}. \quad (25)$$

In our setting, this proposition can be cast in the form of Algorithm 2, and the stable-predictive properties of the induced regret minimizer are analyzed in the following proof.

Proof of Theorem 3.4. First of all, Proposition B.1 implies that the accumulated regret of the regret circuit for the convex hull can be bounded as

$$R_\Psi^T \leq \alpha + A_T + B \sum_{t=2}^T \|\mathbf{L}^t - \mathbf{L}^{t-1}\|_\infty^2 + \beta \sum_{t=2}^T \|\boldsymbol{\Lambda}^t - \boldsymbol{\Lambda}^{t-1}\|_\infty^2, \quad (26)$$

where we used the fact that each regret minimizer $\mathcal{R}_\sigma^{(i)}$ obtains as input the same utility function as $\mathcal{R}_\Psi^{(i)}$, while we also used the notation $\boldsymbol{\lambda}^t \in \mathbb{R}^{|\Sigma_*^{(i)}|}$ to represent the utility function given to $\mathcal{R}_\Delta^{(i)}$, as

predicted by Proposition B.1. Next, let us focus on bounding the norm $\|\mathbf{\Lambda}^t - \mathbf{\Lambda}^{t-1}\|_\infty^2$; in particular, it follows that for some index $s \in \{1, \dots, |\Sigma_*^{(i)}|\}$,

$$\|\mathbf{\Lambda}^t - \mathbf{\Lambda}^{t-1}\|_\infty^2 = (\langle \mathbf{L}^t, \mathbf{x}_s^t \rangle - \langle \mathbf{L}^{t-1}, \mathbf{x}_s^{t-1} \rangle)^2 \quad (27)$$

$$\begin{aligned} &= (\langle \mathbf{L}^t, \mathbf{x}_s^t \rangle - \langle \mathbf{L}^{t-1}, \mathbf{x}_s^t \rangle + \langle \mathbf{L}^{t-1}, \mathbf{x}_s^t \rangle - \langle \mathbf{L}^{t-1}, \mathbf{x}_s^{t-1} \rangle)^2 \\ &= (\langle \mathbf{L}^t - \mathbf{L}^{t-1}, \mathbf{x}_s^t \rangle + \langle \mathbf{L}^{t-1}, \mathbf{x}_s^t - \mathbf{x}_s^{t-1} \rangle)^2 \\ &\leq 2(\langle \mathbf{L}^t - \mathbf{L}^{t-1}, \mathbf{x}_s^t \rangle)^2 + 2(\langle \mathbf{L}^{t-1}, \mathbf{x}_s^t - \mathbf{x}_s^{t-1} \rangle)^2 \end{aligned} \quad (28)$$

$$\leq 2\|\mathbf{L}^t - \mathbf{L}^{t-1}\|_\infty^2 \|\mathbf{x}_s^t\|_1^2 + 2\|\mathbf{L}^{t-1}\|_\infty^2 \|\mathbf{x}_s^t - \mathbf{x}_s^{t-1}\|_1^2 \quad (29)$$

$$\leq 2D_x^2 \|\mathbf{L}^t - \mathbf{L}^{t-1}\|_\infty^2 + 2D_L^2 K^2, \quad (30)$$

where (27) follows from the definition of the ℓ_∞ norm, (28) by Young's inequality, and (29) by Cauchy-Schwarz. Note that we also used the notation D_x to represent an upper bound for the ℓ_1 norm of $\|\mathbf{x}_s^t\|$. As a result, if we plug the bound of Equation (30) to (26) we can conclude that

$$R_\Psi^T \leq \alpha + A_T + (B + 2\beta D_x^2) \sum_{t=2}^T \|\mathbf{L}^t - \mathbf{L}^{t-1}\|_\infty^2 + 2\beta D_L^2 K^2 T. \quad (31)$$

Moreover, regarding the stability of the composite regret minimizer we observe that

$$\begin{aligned} \left\| \sum_{k=1}^{|\Sigma_*^{(i)}|} \lambda_k^t \mathbf{x}_k^t - \sum_{k=1}^{|\Sigma_*^{(i)}|} \lambda_k^{t-1} \mathbf{x}_k^{t-1} \right\|_1 &\leq \sum_{k=1}^{|\Sigma_*^{(i)}|} \|\lambda_k^t \mathbf{x}_k^t - \lambda_k^{t-1} \mathbf{x}_k^{t-1}\|_1 \\ &\leq \sum_{k=1}^{|\Sigma_*^{(i)}|} \|\lambda_k^t \mathbf{x}_k^t - \lambda_k^t \mathbf{x}_k^{t-1} + \lambda_k^t \mathbf{x}_k^{t-1} - \lambda_k^{t-1} \mathbf{x}_k^{t-1}\|_1 \quad (32) \\ &\leq \sum_{k=1}^{|\Sigma_*^{(i)}|} (\|\lambda_k^t \mathbf{x}_k^t - \lambda_k^t \mathbf{x}_k^{t-1}\|_1 + \|\lambda_k^t \mathbf{x}_k^{t-1} - \lambda_k^{t-1} \mathbf{x}_k^{t-1}\|_1) \\ &\leq \sum_{k=1}^{|\Sigma_*^{(i)}|} \lambda_k^t \|\mathbf{x}_k^t - \mathbf{x}_k^{t-1}\|_1 + \sum_{k=1}^{|\Sigma_*^{(i)}|} \|\mathbf{x}_k^{t-1}\|_1 |\lambda_k^t - \lambda_k^{t-1}| \\ &\leq K + D_x \|\boldsymbol{\lambda}^t - \boldsymbol{\lambda}^{t-1}\|_1 \quad (33) \\ &\leq K + D_x \kappa, \quad (34) \end{aligned}$$

where we used the triangle inequality in (32), the ℓ_1 stability of \mathbf{x}_k^t 's in (33), and the ℓ_1 stability of $\boldsymbol{\lambda}^t$'s in (34). Given that $D_x = O(|\Sigma^{(i)}|)$, this verifies our claim about the stability and the predictivity of $\mathcal{R}_\Psi^{(i)}$. Finally, the complexity analysis for the NEXTSTRATEGY function follows directly since the NEXTSTRATEGY operation of each individual regret minimizer runs in $O(|\Sigma^{(i)}|)$, while the analysis of the OBSERVEUTILITY routine follows similarly to (Farina et al., 2021a, Theorem 4.6), and it is therefore omitted. \square

B.2 PROOFS FOR SECTION 3.2

We commence this section with the proof that OMWU guarantees multiplicative-stability.

Lemma B.2 (Multiplicative Stability of OMWU on the Simplex). *Consider the OMWU algorithm \mathcal{R}_Δ on the m -dimensional simplex with learning rate $\eta > 0$. If all the utility functions $\ell^t \in \mathbb{R}^m$ are such that $\|\ell^t\|_\infty \leq L$, and $\eta < 1/(12L)$, then the sequence of \mathbf{x}^t 's produced by OMWU is $O(\eta L)$ -multiplicatively-stable.*

Proof. It is well-known that the update rule of OMWU on the simplex has the following form:

$$\mathbf{x}_i^t = \frac{e^{2\eta \ell_i^{t-1} - \eta \ell_i^{t-2}}}{\sum_{k=1}^m e^{2\eta \ell_k^{t-1} - \eta \ell_k^{t-2}}} \mathbf{x}_i^{t-1}. \quad (35)$$

Algorithm 2: Stable-Predictive Regret Minimizer $\mathcal{R}_{\Psi}^{(i)}$ for the set $\text{co } \Psi^{(i)}$

Input:

- Player $i \in [n]$
- A K -stable (A_T, B) -predictive regret minimizer $\mathcal{R}_{\hat{\sigma}}^{(i)}$ for $\Psi_{\hat{\sigma}}^{(i)}$, for each $\hat{\sigma} \in \Sigma_*^{(i)}$
- A deterministic κ -stable (α, β) -predictive regret minimizer $\mathcal{R}_{\Delta}^{(i)}$ for $\Delta^{|\Sigma_*^{(i)}|}$

1 **Function** NEXTSTRATEGY():

2 $\lambda^t \leftarrow \mathcal{R}_{\Delta}^{(i)}. \text{NEXTSTRATEGY}()$

3 **for** $\hat{\sigma} \in \Sigma_*^{(i)}$ **do**

4 $\phi_{\hat{\sigma} \rightarrow q_{\hat{\sigma}}^t} \leftarrow \mathcal{R}_{\hat{\sigma}}^{(i)}. \text{NEXTSTRATEGY}()$

5 **return** $\sum_{\hat{\sigma} \in \Sigma_*^{(i)}} \lambda^t[\hat{\sigma}] \phi_{\hat{\sigma} \rightarrow q_{\hat{\sigma}}^t}$ represented implicitly as $\{\lambda^t[\hat{\sigma}], q_{\hat{\sigma}}^t\}_{\hat{\sigma} \in \Sigma_*^{(i)}}$

6 **Function** OBSERVEUTILITY(L^t):

7 **for** $\hat{\sigma} \in \Sigma_*^{(i)}$ **do**

8 $\mathcal{R}_{\hat{\sigma}}^{(i)}. \text{OBSERVEUTILITY}(L^t)$

9 Construct the linear function $\ell_{\lambda}^t : \lambda \mapsto \sum_{\hat{\sigma} \in \Sigma_*^{(i)}} \lambda[\hat{\sigma}] L^t(\phi_{\hat{\sigma} \rightarrow q_{\hat{\sigma}}^t})$

10 $\mathcal{R}_{\Delta}^{(i)}. \text{OBSERVEUTILITY}(\ell_{\lambda}^t)$

As a result, it follows that

$$\mathbf{x}_i^t \leq \frac{e^{3\eta L}}{\sum_{k=1}^m e^{-3\eta L} \mathbf{x}_k^{t-1}} \mathbf{x}_i^{t-1} = e^{6\eta L} \mathbf{x}_i^{t-1} \leq (1 + 12\eta L) \mathbf{x}_i^{t-1},$$

where we used that $\ell_i^t \in [-L, L]$, the fact that $\mathbf{x}^{t-1} \in \Delta^m$, and that $e^x \leq 1 + 2x$, for all $x \in [0, 1/2]$. Similarly, we have that

$$\mathbf{x}_i^t \geq \frac{e^{-3\eta L}}{\sum_{k=1}^m e^{3\eta L} \mathbf{x}_k^{t-1}} \mathbf{x}_i^{t-1} = e^{-6\eta L} \mathbf{x}_i^{t-1} \geq (1 - 6\eta L) \mathbf{x}_i^{t-1}.$$

□

Proof of Lemma 3.5. Let us start with the regret minimizer $\mathcal{R}_{\Delta}^{(i)}$. By Lemma B.2 it suffices to bound the ℓ_{∞} norm of the utility vectors. In particular, if $\mathbf{\Lambda}$ is the associated utility vector predicted from Proposition B.1, it follows that $\|\mathbf{\Lambda}\|_{\infty} = |\langle \mathbf{L}, \mathbf{x}_k \rangle| \leq \|\mathbf{L}\|_{\infty} \|\mathbf{x}_k\|_1 = O(D_{\mathbf{L}} |\Sigma^{(i)}|)$, for some $k \in \{1, \dots, |\Sigma_*^{(i)}|\}$, where we used the (generalized) Cauchy–Schwarz inequality, and the fact that $\|\mathbf{x}_k\|_1 = O(|\Sigma^{(i)}|)$. Moreover, for $\mathbf{L}^t = (\mathbf{x}^t \otimes \ell^t)^b$ and $\|\mathbf{x}^t\|_{\infty} \leq 1$ it follows that $D_{\mathbf{L}} \leq D_{\ell}$, and the claimed bound follows directly from Lemma B.2. Finally, for each regret minimizer $\mathcal{R}_{\hat{\sigma}}$ the multiplicative-stability bound can only be a factor of $\mathfrak{D}^{(i)}$ from the derived bound in Theorem A.4, concluding the proof. □

Next, we focus on the proof of Theorem 3.10. In particular, let us first introduce some important tools that we will require. Our approach is based on the techniques of Kruckman et al. (2010), who provided an alternative proof of the classic Markov chain tree theorem using linear-algebraic techniques. We commence by stating some elementary properties of the determinant.

Fact B.3. *The following properties hold:*

- The determinant is a multilinear function with respect to the rows and columns of the matrix. That is,

$$\det(\mathbf{u}_1, \dots, \alpha \mathbf{u}_k + \beta \mathbf{u}'_k, \dots, \mathbf{u}_m) = \alpha \det(\mathbf{u}_1, \dots, \mathbf{u}_k, \dots, \mathbf{u}_m) + \beta \det(\mathbf{u}_1, \dots, \mathbf{u}'_k, \dots, \mathbf{u}_m),$$

for any $\mathbf{u}_1, \dots, \mathbf{u}_m \in \mathbb{R}^m$, $\mathbf{u}'_k \in \mathbb{R}^m$, and $\alpha, \beta \in \mathbb{R}$;

- If any two rows or columns of a square matrix \mathbf{A} are equal, then $\det(\mathbf{A}) = 0$;
- The determinant remains invariant under permutations.

Given a matrix \mathbf{A} , the minor $\text{mn}^{(i,j)}(\mathbf{A})$ is the matrix formed from \mathbf{A} after deleting the i -th row and the j -th column. Then, the *cofactor* is defined as $\text{co}^{(i,j)}(\mathbf{A}) = (-1)^{i+j} \det(\text{mn}^{(i,j)}(\mathbf{A}))$, while the *adjugate* (or adjoint) matrix $\text{adj}(\mathbf{A})^\top$ is the matrix with entries the corresponding cofactors of \mathbf{A} . With this notation at hand, we are ready to state the following characterization due to (Kruckman et al., 2010, Theorem 3.4):

Theorem B.4 (Kruckman et al. (2010)). *Consider an ergodic m -state Markov chain with transition matrix \mathbf{M} . If $\mathbf{x} \in \mathbb{R}^m$ is such that $\mathbf{x}_i := \text{adj}(\mathcal{L})[(i, i)]$, where $\mathcal{L} := \mathbf{M} - \mathbf{I}_m$ is the Laplacian of the system, \mathbf{x} is an eigenvector of \mathcal{M} with corresponding eigenvalue of 1. That is, $\mathbf{M}\mathbf{x} = \mathbf{x}$.*

An key ingredient of our proof for Theorem 3.10 relies on this theorem in order to characterize the stationary distribution of a certain (ergodic) Markov chain. Incidentally, an alternative characterization can be provided using the classic Markov chain tree theorem. In particular, a central component of this theorem is the notion of a *directed tree*:

Definition B.5 (Directed Tree). A graph $G = (V, E)$ is said to be a *directed tree* rooted at $u \in V$ if (i) it does not contain any cycles, and (ii) u has no outgoing edges, while every other node has exactly one outgoing edge.

We will represent with \mathcal{D}_i the set of all graphs which have property (ii) with respect to a node $i \in [m]$. Moreover, we will use \mathcal{T}_i to represent the subset of \mathcal{D}_i which also has property (i) of Definition B.5. For a matrix $\mathbf{D} \in \mathcal{D}_i$, we define a matrix $\text{mp}(\mathbf{D})$ so that $\text{mp}(\mathbf{D})_{(j,k)} = 1$ if $(k, j) \in E(\mathbf{D})$, and 0 otherwise. The following lemma will be of particular use for our purposes:

Lemma B.6 (Kruckman et al. (2010)). *Consider some $m \times m$ matrix $\mathbf{D} \in \mathcal{D}_i$, and let R_i be the determinant of the Laplacian matrix $\mathcal{L} := \text{mp}(\mathbf{D}) - \mathbf{I}$ after replacing the i -th column with the i -th standard unit vector \mathbf{e}_i . Then, $R_i = (-1)^{m-1}$ if $\mathbf{D} \in \mathcal{T}_i$, i.e. \mathbf{D} contains no (directed) cycles. Otherwise, $R_i = 0$.*

Before we proceed with the main technical proof, we also state another useful elementary fact.

Fact B.7. *The adjugate matrix at (i, i) is equal to the determinant of \mathbf{A} after we replace the i -th column with the vector \mathbf{e}_i .*

Lemma B.8. *Let \mathbf{M} be the transition matrix of an m -state Markov chain such that $\mathbf{M} := \mathbf{v}\mathbf{1}^\top + \mathbf{C}$, where \mathbf{C} is a matrix with strictly positive entries and columns summing to $1 - \lambda$, and \mathbf{v} is a vector with strictly positive entries summing to λ . Then, if $\boldsymbol{\pi}$ is the stationary distribution of \mathbf{M} , there exists, for each $i \in [m]$, a (non-empty) finite set F_i and $F = \bigcup_i F_i$, and corresponding parameters $b_j \in \{0, 1\}$, $0 \leq p_j \leq m - 2$, $|S_j| = m - p_j - b_j - 1$, for each $j \in F_i$, such that*

$$\pi_i = \frac{\sum_{j \in F_i} \lambda^{p_j+1} (\mathbf{v}[q_j])^{b_j} \prod_{(s,w) \in S_j} \mathbf{C}[(s, w)]}{\sum_{j \in F} C_j \lambda^{p_j+b_j} \prod_{(s,w) \in S_j} \mathbf{C}[(s, w)]}, \quad (36)$$

where $C_j = C_j(m)$ is a positive constant.

Before we proceed with the proof, let us stress that the main takeaway of this lemma is that the stationary distribution has an *affine* dependence on the vector \mathbf{v} , which will be of particular importance in order to obtain polynomial bounds in terms of stability.

Proof of Lemma B.8. Let us consider the Laplacian matrix $\mathcal{L} = \mathbf{M} - \mathbf{I}_m$, and the quantities $\Sigma_i := \text{adj}(\mathcal{L})[(i, i)]$. Our first goal is to characterize the behavior of Σ_i 's. In particular, we can focus without loss of generality on the term Σ_1 . We know from Fact B.7 that Σ_1 can be expressed as

$$\Sigma_1 = \det(\mathbf{e}_1, \mathbf{v} + \mathbf{c}_2 - \mathbf{e}_2, \dots, \mathbf{v} + \mathbf{c}_m - \mathbf{e}_m), \quad (37)$$

where \mathbf{c}_j represents the j -th column of \mathbf{C} . Now if $\mathbf{e}_{j,k} := \mathbf{e}_j - \mathbf{e}_k$, given that \mathbf{M} is column-stochastic we can observe that

$$\mathbf{e}_j - \mathbf{v} - \mathbf{c}_j = \sum_{k=1}^m (\mathbf{e}_j - \mathbf{e}_k) \mathbf{v}[k] + \sum_{k=1}^m (\mathbf{e}_j - \mathbf{e}_k) \mathbf{c}_j[k] = \sum_{k=1}^m \mathbf{e}_{j,k} \mathbf{v}[k] + \sum_{k=1}^m \mathbf{e}_{j,k} \mathbf{c}_j[k].$$

Next, if we plug-in this expansion to Equation (37) it follows that

$$\Sigma_1 = \det \left(\mathbf{e}_1, \sum_{k=1}^m \mathbf{e}_{k,2} \mathbf{v}[k] + \sum_{k=1}^m \mathbf{e}_{k,2} \mathbf{c}_2[k], \dots, \sum_{k=1}^m \mathbf{e}_{k,m} \mathbf{v}[k] + \sum_{k=1}^m \mathbf{e}_{k,m} \mathbf{c}_m[k] \right). \quad (38)$$

By multilinearity of the determinant (Fact B.3), Σ_1 can be expressed as the sum of terms, with a single term of the form

$$\det \left(e_1, \sum_{k=1}^m e_{k,2} c_2[k], \dots, \sum_{k=1}^m e_{k,m} c_m[k] \right), \quad (39)$$

independent on v , while any other term can be expressed in the form

$$\det \left(e_1, z_2, \dots, \sum_{k=1}^m e_{k,j} v[k], \dots, z_m \right), \quad (40)$$

for some index j , where z_ℓ is either $\sum_{k=1}^m e_{k,\ell} v[k]$ or $\sum_{k=1}^m e_{k,\ell} c_\ell[k]$. Now let us focus on each term of Equation (40). Specifically, we will show that it can be equivalently expressed so that the vector v appears only in a single column. Indeed, consider any other column in the matrix involved in the determinant of (40), expressed in the form $\sum_{k=1}^m e_{k,\ell} v[k]$, for some index $\ell \neq j$, if such column exists. Then, if we subtract the second column from that column it would take the form

$$\sum_{k=1}^m e_{k,\ell} v[k] - \sum_{k=1}^m e_{k,j} v[k] = \sum_{k=1}^m (e_j - e_\ell) v[k] = \lambda e_{j,\ell},$$

where recall that λ is the sum of the entries of vector v , while this operation does not modify the value of the underlying determinant. Thus, by multilinearity the determinant of (40) is equal to

$$\lambda^p \det \left(e_1, z'_2, \dots, \sum_{k=1}^m e_{k,j} v[k], \dots, z'_m \right), \quad (41)$$

where z'_ℓ is either $\sum_{k=1}^m e_{k,\ell} c_\ell[k]$, or $e_{j,\ell}$, and $0 \leq p \leq m-2$. Next, if we use again the multilinearity property, the term in (41) can be expressed as a sum of terms each of which has the form

$$\left(\lambda^p v[q] \prod_{(s,w) \in S} C[(s,w)] \right) \det(e_1, e_{\cdot,2}, \dots, e_{\cdot,m}),$$

where $|S| = m - p - 2$. In particular, the induced determinant $\det(e_1, e_{\cdot,2}, \dots, e_{\cdot,m})$ matches after a suitable permutation the form of Lemma B.6, associated with some matrix $\mathbf{D} \in \mathcal{D}_i$. Thus, it can either be 0 or $(-1)^{m-1}$, depending on whether the corresponding graph has a (directed) cycle. Similar reasoning applies for the determinant in (39), which can be expressed as a sum of terms

$$(-1)^{m-1} \prod_{(s,w) \in S} C[(s,w)],$$

where $|S| = m - 1$. Overall, we have shown that each Σ_i (due to symmetry) can be expressed in the form

$$(-1)^{m-1} \sum_{j \in F_i} \lambda^{p_j} (v[q_j])^{b_j} \prod_{(s,w) \in S_j} C[(s,w)], \quad (42)$$

where for all j it holds that $b_j \in \{0, 1\}$, and $|S_j| = m - p_j - b_j - 1$. Next, we will focus on characterizing $\Sigma := \sum_{i=1}^m \Sigma_i$. In particular, the stationary distribution π of \mathbf{M} is such that

$$(\mathbf{C} + v \mathbf{1}^\top) \pi = \pi \iff \mathbf{C} \pi + v = \pi \iff (\mathbf{I}_m - \mathbf{C}) \pi = v, \quad (43)$$

where we used that $\mathbf{1}^\top \pi = 1$ since $\pi \in \Delta^m$. Moreover, we claim that the matrix $\mathbf{I}_m - \mathbf{C}$ is invertible. Indeed, the sum of the columns of \mathbf{C} is $1 - \lambda$, and subsequently it follows that the maximum eigenvalue of \mathbf{C} is $(1 - \lambda)$. In turn, this implies that all the eigenvalues of $\mathbf{I}_m - \mathbf{C}$ are at least $\lambda > 0$. As a result, we can use Cramer's rule to obtain an explicit formula for the solution of the linear system with respect to the first coordinate of π :

$$\pi_1 = \frac{\det(v, e_2 - c_2, \dots, e_m - c_m)}{\det(e_1 - c_1, e_2 - c_2, \dots, e_m - c_m)}. \quad (44)$$

Moreover, it follows that

$$\begin{aligned}
\pi_1 &= \det(\mathbf{v}, \mathbf{e}_2 - \mathbf{c}_2, \dots, \mathbf{e}_m - \mathbf{c}_m) \\
&= \det(\mathbf{v}, \mathbf{e}_2 - \mathbf{c}_2 - \mathbf{v}, \dots, \mathbf{e}_m - \mathbf{c}_m - \mathbf{v}) \\
&= \det(\mathbf{v} + (\lambda \mathbf{e}_1 - \mathbf{v}), \mathbf{e}_2 - \mathbf{c}_2 - \mathbf{v}, \dots, \mathbf{e}_m - \mathbf{c}_m - \mathbf{v}) \\
&= \lambda \det(\mathbf{e}_1, \mathbf{e}_2 - \mathbf{c}_2 - \mathbf{v}, \dots, \mathbf{e}_m - \mathbf{c}_m - \mathbf{v}),
\end{aligned} \tag{45}$$

where in (45) we used the fact that $\det(\lambda \mathbf{e}_1 - \mathbf{v}, \dots, \mathbf{e}_m - \mathbf{c}_m - \mathbf{v}) = 0$. Thus, if we use the definition of Σ_1 , Fact B.7, and Equation (44), we arrive at the following conclusion:

$$\pi_1 = \lambda \frac{\Sigma_1}{\det(\mathbf{I}_m - \mathbf{C})}.$$

But we can also infer from Theorem B.4 that $\pi_1 = \Sigma_1/\Sigma$, implying the following identity:

$$\det(\mathbf{I}_m - \mathbf{C}) = \lambda \sum_{i=1}^m \Sigma_i. \tag{46}$$

In fact, we have shown this formula for *any* vector $\lambda \mathbf{p}$, where \mathbf{p} is a probability distribution, and $\lambda > 0$. Thus, it must also hold for $\mathbf{v} := \frac{\lambda}{m} \mathbf{1}$. That is,

$$\det(\mathbf{I}_m - \mathbf{C}) = \lambda(-1)^{m-1} \sum_{j \in F} C_j \lambda^{p_j+b_j} \prod_{(s,w) \in S_j} \mathbf{C}[(s,w)], \tag{47}$$

where $|S_j| \leq m-1-p_j$, $C_j = C_j(m)$ is a positive parameter independent on the entries of \mathbf{v} and \mathbf{C} , and $F = \bigcup_i F_i$. Finally, given that the vector $\boldsymbol{\pi} \in \Delta^m$ with $\pi_i = \Sigma_i/\Sigma$ is the (unique) stationary distribution of \mathbf{M} , the lemma follows directly from (42), (46), and (47). \square

Corollary B.9. *Let \mathbf{M} be the transition matrix of an m -state Markov chain such that $\mathbf{M} = \mathbf{v}\mathbf{1}^\top + \mathbf{C}$, where \mathbf{C} is a matrix with strictly positive entries and columns summing to $1-\lambda$, and \mathbf{v} is a vector with strictly positive entries summing to λ . Moreover, let $\mathbf{v} = \mathbf{r}/l$, for some $l > 0$. Then, if $\boldsymbol{\pi}$ is the stationary distribution of \mathbf{M} , there exists a set F_i , for each $i \in [m]$, a (non-empty) finite set F_i and $F = \bigcup_i F_i$, and corresponding parameters $b_j \in \{0, 1\}$, $0 \leq p_j \leq m-2$, $|S_j| = m-p_j-b_j-1$, for each $j \in F_i$, such that the i -th coordinate of the vector $\mathbf{w} := l\boldsymbol{\pi}$ can be expressed as*

$$\mathbf{w}_i = \frac{\sum_{j \in F_i} \lambda^{p_j+1} (\mathbf{r}[q_j])^{b_j} l^{1-b_j} \prod_{(s,w) \in S_j} \mathbf{C}[(s,w)]}{\sum_{j \in F} C_j \lambda^{p_j+b_j} \prod_{(s,w) \in S_j} \mathbf{C}[(s,w)]}, \tag{48}$$

where $C_j = C_j(m)$ is a positive constant.

Proof. The proof follows directly from the formula derived in Lemma B.8. \square

This expression for the stationary distribution was derived specifically to characterize the multiplicative stability of the fixed points associated with EFCE. In particular, this will be shown directly from the following immediate corollary:

Corollary B.10. *Let \mathbf{M} be the transition matrix of an m -state Markov chain such that $\mathbf{M} = \mathbf{v}\mathbf{1}^\top + \mathbf{C}$, where \mathbf{C} is a matrix with strictly positive entries, and \mathbf{v} is a vector with strictly positive entries such that $\mathbf{v} = \mathbf{r}/l$, for some $l > 0$. If $\boldsymbol{\pi}$ is the stationary distribution of \mathbf{M} , let $\mathbf{w} := l\boldsymbol{\pi}$. Similarly, let \mathbf{M}' be the transition matrix of another m -state Markov chain such that $\mathbf{M}' = \mathbf{v}'\mathbf{1}^\top + \mathbf{C}'$, where \mathbf{C}' is a matrix with strictly positive entries, and \mathbf{v}' is a vector with strictly positive entries such that $\mathbf{v}' = \mathbf{r}'/l'$, for some $l' > 0$. If $\boldsymbol{\pi}'$ is the stationary distribution of \mathbf{M}' , let $\mathbf{w}' := l'\boldsymbol{\pi}'$. Finally, let λ and λ' be the sum of the entries of \mathbf{v} and \mathbf{v}' respectively. Then, if (i) the matrices \mathbf{C} and \mathbf{C}' are κ -multiplicatively-close, (ii) the numbers λ and λ' are κ -multiplicatively-close, (iii) the vectors \mathbf{r} and \mathbf{r}' are γ -multiplicatively-close, and (iv) the numbers l and l' are also γ -multiplicatively-close, then the vectors \mathbf{w} and \mathbf{w}' are $(\gamma + O(\kappa m))$ -multiplicatively-close, assuming that $\kappa = O(1/m)$.*

Proof. Consider some coordinate $i \in [m]$, and let

$$V_j := \lambda^{p_j+1} (\mathbf{r}[q_j])^{b_j} l^{1-b_j} \prod_{(s,w) \in S_j} \mathbf{C}[(s,w)],$$

for some $j \in F_i$. Also let V'_j be the corresponding quantity with respect to \mathbf{M}' . Then, by assumption we have that

$$V'_j \leq (1 + \kappa)^{p_j+1} (1 + \gamma) (1 + \kappa)^{|S_j|} V_j = (1 + \gamma) (1 + \kappa)^m V_j,$$

where we used the fact that $|S_j| + p_j + 1 \leq m$ by Corollary B.9. Moreover, for a sufficiently small $\kappa = O(1/m)$, we can infer that $V'_j \leq (1 + \gamma) (1 + O(\kappa m)) V'_j = (1 + (\gamma + O(\kappa m))) V_j$. Thus, this implies that $\sum_{j \in F_i} V'_j \leq (1 + (\gamma + O(\kappa m))) \sum_{j \in F_i} V_j$. Moreover, we can show that the denominator of Equation (48) induces an extra additive factor of $O(\kappa m)$ in the multiplicative stability, implying that $\mathbf{w}'_i \leq (1 + (\gamma + O(\kappa m))) \mathbf{w}_i$, for any $i \in [m]$. Similarly, it follows that $\mathbf{w}'_i \geq (1 - (\gamma + O(\kappa m))) \mathbf{w}_i$. \square

Next, we will use this statement to prove Proposition 3.9. We note that it is tacitly assumed that the vectors $\boldsymbol{\lambda}^t$, \mathbf{q}_σ^t , and $\mathbf{x}_{(j \in J)}$, involved in Proposition 3.9, lie on the interior; this will be indeed the case when OMWU is used to construct the individual regret minimizers.

Proof of Proposition 3.9. Let us focus on the stability analysis of Algorithm 3, as the rest of the claim follows from (Farina et al., 2021a, Proposition 4.14). In particular, for consistency with the terminology of Corollary B.10, let us define

$$\mathbf{C}[(a_r, a_c)] := \boldsymbol{\lambda}[(j^*, a_c)] \mathbf{q}_{(j^*, a_c)}[(j^*, a_r)] + \left(1 - \sum_{\tilde{\sigma} \preceq (j^*, a_c)} \boldsymbol{\lambda}[\tilde{\sigma}] \right) \mathbf{1}\{a_r = a_c\};$$

and $l := \mathbf{x}[\sigma_p]$. We will show that the requirements of Corollary B.10 are satisfied:

- (i) The entries of \mathbf{C} are $O(\kappa)$ -multiplicatively-stable. In particular, this follows from the fact that $1 - \sum_{\tilde{\sigma} \preceq (j^*, a_c)} \boldsymbol{\lambda}[\tilde{\sigma}] = \sum_{\tilde{\sigma} \in \tilde{\Sigma}} \boldsymbol{\lambda}[\tilde{\sigma}]$, and the latter term is clearly κ -multiplicatively-stable;
- (ii) The sum of the entries of $\mathbf{v}^t := \mathbf{r}^t / l^t$ is κ -multiplicatively-stable. To see this, note that the sum of each column of \mathbf{C} can be expressed as $\sum_{\tilde{\sigma} \in \tilde{\Sigma}} \boldsymbol{\lambda}[\tilde{\sigma}]$, and as a result, since the matrix $\mathbf{C} + \frac{1}{l} \mathbf{r} \mathbf{1}^\top$ is stochastic, we can infer that the sum of the entries of \mathbf{v} can also be expressed as $\sum_{\tilde{\sigma} \in \tilde{\Sigma}} \boldsymbol{\lambda}[\tilde{\sigma}]$, since $\boldsymbol{\lambda}$ is a vector on the simplex. But the latter term is clearly κ -multiplicatively-stable, as desired;
- (iii) The sequence of vectors \mathbf{r}^t is $\gamma + O(\kappa)$ -multiplicatively-stable. This assertion can be directly verified from the definition of \mathbf{r} in Algorithm 3;
- (iv) The sequence of numbers of l^t is γ -multiplicatively-stable. Indeed, this follows directly from the assumption that the sequence of \mathbf{x}^t 's is γ -multiplicatively-stable.

As a result, we can apply Corollary B.10 to conclude the proof. \square

Proof of Theorem 3.10. Our argument proceeds inductively. For a root information set $j \in \mathcal{J}^{(i)}$, Proposition 3.9 implies $O(\kappa|\mathcal{A}|)$ -multiplicative-stability for any induced partial fixed point; this follows given that the \emptyset -partial fixed point is trivially 0-multiplicatively-stable. Next, the theorem follows inductively given that by Proposition 3.9 each sequence can only incur an additive factor of $O(\kappa|\mathcal{A}|)$ in the multiplicative stability bound with respect to the preceding sequences. \square

More precisely, if $F^{(i)} := \max_{j_1 \prec j_2 \prec \dots \prec j_d} \sum_{i=1}^d |\mathcal{A}(j_i)|$, with $j_1, \dots, j_d \in \mathcal{J}^{(i)}$, we can show that the sequence of fixed points is $O(\kappa F^{(i)})$ -multiplicatively-stable. Observe that $F^{(i)}$ can be trivially upper bounded by $|\mathcal{A}^{(i)}| \mathfrak{D}^{(i)}$, as well as the number of sequences $|\Sigma^{(i)}|$.

B.3 PUTTING EVERYTHING TOGETHER

First of all, if we combine Theorem A.4 with Theorem 3.4, instantiating the regret minimizer $\mathcal{R}_\Delta^{(i)}$ of Theorem 3.4 using OMWU with learning rate η , we can construct a regret minimizer for the set $\text{co } \Psi^{(i)}$ with regret R_Ψ^T bounded as

$$R_\Psi^T \leq \frac{\mathcal{P}^{(i)}}{\eta} + \eta^3 D_\ell^4 \mathcal{V}^{(i)} T + \eta \mathcal{B}^{(i)} \sum_{t=2}^T \|\mathbf{L}^t - \mathbf{L}^{t-1}\|_\infty^2, \quad (49)$$

Algorithm 3: $\text{EXTEND}(\phi, J, j^*, \mathbf{x})$; (Farina et al., 2021a)

Input:

- $\phi = \sum_{\hat{\sigma} \in \Sigma_*^{(i)}} \lambda[\hat{\sigma}] \phi_{\hat{\sigma} \rightarrow \mathbf{q}_{\hat{\sigma}}}^{(i)} \in \text{co } \Psi^{(i)}$
- $J \subseteq \mathcal{J}^{(i)}$ trunk for player i
- $j^* \in \mathcal{J}^{(i)}$ information set not in J with an immediate predecessor in J
- $\mathbf{x} \in \mathbb{R}_{\geq 0}^{|\Sigma^{(i)}|}$ J -partial fixed point of ϕ

Output: $\mathbf{x}' \in \mathbb{R}_{\geq 0}^{|\Sigma^{(i)}|}$ $(J \cup \{j^*\})$ -partial fixed point of ϕ

- 1 $\sigma_p \leftarrow \sigma^{(i)}(j^*)$
- 2 Let $\mathbf{r} \in \mathbb{R}_{\geq 0}^{|\mathcal{A}(j^*)|}$ be defined as

$$\mathbf{r}[a] := \sum_{j' \preceq \sigma_p} \sum_{a' \in \mathcal{A}(j')} \lambda[(j', a')] \mathbf{q}_{(j', a')}[(j^*, a)] \mathbf{x}[(j', a')]$$
- 3 Let $\mathbf{W} \in \mathbb{R}^{|\mathcal{A}(j^*)| \times |\mathcal{A}(j^*)|}$ be the matrix with entries $\mathbf{W}[a_r, a_c]$ defined, for $a_r, a_c \in \mathcal{A}(j^*)$, as

$$\mathbf{r}[a_r] + \left(\lambda[(j^*, a_c)] \mathbf{q}_{(j^*, a_c)}[(j^*, a_r)] + \left(1 - \sum_{\hat{\sigma} \preceq (j^*, a_c)} \lambda[\hat{\sigma}] \right) \mathbf{1}\{a_r = a_c\} \right) \mathbf{x}[\sigma_p]$$
- 4 **if** $\mathbf{x}[\sigma_p] = 0$ **then**
 - 5 | $\mathbf{w} \leftarrow \mathbf{0} \in \mathbb{R}_{\geq 0}^{|\mathcal{A}(j^*)|}$
- 6 **else**
 - 7 | $\mathbf{b} \in \Delta^{|\mathcal{A}(j^*)|} \leftarrow$ stationary distribution of $\frac{1}{\mathbf{x}[\sigma_p]} \mathbf{W}$
 - 8 | $\mathbf{w} \rightarrow \mathbf{x}[\sigma_p] \mathbf{b}$
- 9 $\mathbf{x}' \leftarrow \mathbf{x}$
- 10 **for** $a \in \mathcal{A}(j^*)$ **do**
 - 11 | $\mathbf{x}'[(j^*, a)] \leftarrow \mathbf{w}[(j^*, a)]$

Algorithm 4: $\text{FIXEDPOINT}(\phi)$; (Farina et al., 2021a)

Input: $\phi = \sum_{\hat{\sigma} \in \Sigma_*^{(i)}} \lambda[\hat{\sigma}] \phi_{\hat{\sigma} \rightarrow \mathbf{q}_{\hat{\sigma}}}^{(i)} \in \text{co } \Psi^{(i)}$

Output: $\mathbf{q}^* \in \mathcal{Q}^{(i)}$ such that $\mathbf{q} = \phi(\mathbf{q}^*)$

- 1 $\mathbf{q} \leftarrow \mathbf{0} \in \mathbb{R}^{|\Sigma^{(i)}|}, \mathbf{q}[\emptyset] \leftarrow \emptyset$
- 2 $J \leftarrow \emptyset$
- 3 **for** $j \in \mathcal{J}^{(i)}$ **in top-down order do**
 - 4 | $\mathbf{q}^* \leftarrow \text{EXTEND}(\phi, J, j, \mathbf{q}^*)$
 - 5 | $J = J \cup \{j\}$
- 6 **return** \mathbf{q}

where $\mathcal{P}^{(i)}$, $\mathcal{V}^{(i)}$, and $\mathcal{B}^{(i)}$ are game-specific parameters *polynomial* on the size of the game, and independent of T and η , and D_ℓ is an upper bound on the ℓ_∞ norm of the utility sequences. In addition, this regret minimizer will be $\eta \mathcal{K}^{(i)}$ -stable, for a game-specific parameter $\mathcal{K}^{(i)}$ polynomial on the size of the game, and independent on η and T . Thus, this stable-predictive regret minimizer for the set of transformations $\Psi^{(i)}$ realizes the first requirement of Theorem 3.2.

Moreover, let $\eta \mathcal{G}^{(i)}$ be the parameter associated with the stability of the fixed points, as predicted by Corollary 3.11, for some game-specific parameter $\mathcal{G}^{(i)}$ polynomial on the size of the game, and independent on η and T . Now observe that this component realizes the STABLEFPORACLE; the second requirement of Theorem 3.2. As a result, we can apply Theorem 3.2 to conclude the following:

Corollary B.11. *Consider any player $i \in [n]$. There exists a learning algorithm such that under any sequence of utility vectors ℓ^1, \dots, ℓ^T , with ℓ_∞ norm bounded by D_ℓ , the accumulated $\Psi^{(i)}$ -regret R^T can be bounded as*

$$R^T \leq \frac{\mathcal{P}^{(i)}}{\eta} + \eta^3 D_\ell^4 \mathcal{V}^{(i)} T + \eta \mathcal{B}^{(i)} \sum_{t=2}^T \|\ell^t - \ell^{t-1}\|_\infty^2. \quad (50)$$

In addition, the regret minimizer is $\eta\mathcal{G}^{(i)}$ -stable, for a sufficiently small learning rate η . All of the parameters $\mathcal{P}^{(i)}$, $\mathcal{V}^{(i)}$, $\mathcal{B}^{(i)}$, and $\mathcal{G}^{(i)}$ are game-specific parameters, polynomial on the size of the game, and independent on η and T .

We note that the parameters appearing in this corollary are in general different from the ones in Equation (49), and we use the same symbols with an abuse of notation. So far we have focused on bounding the regret of a player without any assumptions on the observed utility functions. A crucial observation is that if all players employ a regularized (or smooth) learning algorithm, then the observed utility functions change slowly over time. To formalize this observation, we start with the following auxiliary claim:

Claim B.12. For any player $i \in [n]$ the observed utilities satisfy

$$\|\ell^{(i),t} - \ell^{(i),t-1}\|_\infty \leq |\mathcal{Z}| \sum_{k \neq i}^n \|\mathbf{q}^{(k),t} - \mathbf{q}^{(k),t-1}\|_1. \quad (51)$$

Proof. For a profile of mixed sequence-form strategies $(\mathbf{q}^{(1)}, \dots, \mathbf{q}^{(n)})$, the utility of player i can be expressed as

$$u^{(i)}(\mathbf{q}^{(1)}, \dots, \mathbf{q}^{(n)}) = \sum_{z \in \mathcal{Z}} p^{(c)}(z) u^{(i)}(z) \prod_{k=1}^n \mathbf{q}^{(k)}(\sigma^{(k)}(z)).$$

As a result, given that $|u^{(i)}(z)| \leq 1$ for all $z \in \mathcal{Z}$, it follows that

$$\begin{aligned} \|\ell^{(i),t} - \ell^{(i),t-1}\|_\infty &\leq \sum_{z \in \mathcal{Z}} \left| \prod_{k \neq i}^n \mathbf{q}^{(k),t}(\sigma^{(k)}(z)) - \prod_{k \neq i}^n \mathbf{q}^{(k),t-1}(\sigma^{(k)}(z)) \right| \\ &\leq \sum_{z \in \mathcal{Z}} \sum_{k \neq i}^n \left| \mathbf{q}^{(k),t}(\sigma^{(k)}(z)) - \mathbf{q}^{(k),t-1}(\sigma^{(k)}(z)) \right|, \end{aligned} \quad (52)$$

where in the last bound we used the inequality

$$|(a_1 a_2 \dots a_m) - (b_1 b_2 \dots b_m)| \leq \sum_{i=1}^m |a_i - b_i| (a_1 \dots a_{i-1}) (b_{i+1} \dots b_m) \leq \sum_{i=1}^m |a_i - b_i|,$$

for any $a_1, \dots, a_m, b_1, \dots, b_m \in [0, 1]$. Finally, from (52) we can conclude that

$$\|\ell^{(i),t} - \ell^{(i),t-1}\|_\infty \leq \sum_{k \neq i}^n \sum_{z \in \mathcal{Z}} \left| \mathbf{q}^{(k),t}(\sigma^{(k)}(z)) - \mathbf{q}^{(k),t-1}(\sigma^{(k)}(z)) \right| \leq |\mathcal{Z}| \sum_{k \neq i}^n \|\mathbf{q}^{(k),t} - \mathbf{q}^{(k),t-1}\|_1,$$

as desired. \square

As a result, if *all* players employ a stable algorithm to update their strategies, the observed utilities will also be stable. This leads to the following conclusion:

Corollary B.13. Assume that all agents play according to the dynamics associated with Corollary B.11. Then, for learning rate $\eta = O(T^{-1/4})$, the $\Psi^{(i)}$ -regret of every player i is bounded by $O(T^{1/4})$, where the $O(\cdot)$ notation suppresses game-specific parameters polynomial in the size of the game.

Proof. Let us use the notation $\mathcal{P} := \max_i \mathcal{P}^{(i)}$, $\mathcal{V} := \max_i \mathcal{V}^{(i)}$, $\mathcal{B} := \max_i \mathcal{B}^{(i)}$, and $\mathcal{G} := \max_i \mathcal{G}^{(i)}$. For any player $i \in [n]$, Claim B.12 implies that $\|\ell^{(i),t} - \ell^{(i),t-1}\|_\infty \leq |\mathcal{Z}|(n-1)\eta\mathcal{G}$, since every player updates her strategy through an $\eta\mathcal{G}$ -stable learning algorithm. Thus, by the guarantee of Corollary B.11 we can conclude that for any player i ,

$$R^{(i),T} \leq \frac{\mathcal{P}}{\eta} + \eta^3 \mathcal{V} T + \eta^3 (n-1)^2 \mathcal{B} |\mathcal{Z}|^2 \mathcal{G}^2 T,$$

where we used the fact that $D_\ell \leq 1$, which follows from the normalization assumption on the utilities. Finally, taking $\eta = O(T^{-1/4})$ yields that $R^{(i),T} = O(T^{1/4})$, as desired. \square

Finally, we conclude with the proof of Theorem 1.1.

Proof of Theorem 1.1. The guarantee of Corollary B.13 along with Theorem 2.6 suffice in order to establish Theorem 1.1. In particular, for a player i let $\mu^{(i),t}$ be any probability distribution on the set $\Pi^{(i)}$ such that $\mathbb{E}_{\pi \sim \mu^{(i),t}}[\pi] = q^{(i),t}$, where $q^{(i),t}$ is the output of the regret minimizer, on the *mixed* sequence-form strategy polytope $\mathcal{Q}^{(i)}$, realized with the dynamics of Corollary B.13. Moreover, let $\mu^t := \mu^{(1),t} \otimes \dots \otimes \mu^{(n),t}$ be the joint probability distribution, and $\bar{\mu} := \frac{1}{T} \sum_{t=1}^T \mu^t$. Then, it follows by the linearity of expectation, as well as the linearity of the expression in Definition 2.5 (recall that the set of transformations $\Psi^{(i)}$ is linear) that $\bar{\mu}$ is an ϵ -EFCE, where $\epsilon := \frac{1}{T} \max_i R^{(i),T}$, and $R^{(i),T}$ is the cumulative $\Psi^{(i)}$ -regret with respect to the set $\mathcal{Q}^{(i)}$. Finally, the proof follows given that by virtue of Corollary B.13, $R^{(i),T} = O(T^{1/4})$ for all $i \in [n]$. \square

We also remark that the complexity of each iteration of the described dynamics is analogous to that in (Farina et al., 2021a).

C ACCELERATION FOR EXTENSIVE-FORM COARSE CORRELATED EQUILIBRIUM

In this section we develop accelerated uncoupled no-regret dynamics converging to an extensive-form *coarse* correlated equilibrium (EFCE). In particular, we will provide a slightly better bound compared to EFCE. At the heart of this improvement is an algorithm due to Farina et al. (2021b) which circumvents the need of computing the stationary distribution of a Markov chain; this will allow us to show more favorable stability properties. First, we introduce some basic concepts, commencing with that of *coarse* trigger deviation functions:

Definition C.1 (Coarse Trigger Deviation Functions). Consider some player $i \in [n]$, some information set $j \in \mathcal{J}^{(i)}$, and a continuation strategy $\hat{\pi} \in \Pi_j^{(i)}$. A *coarse trigger deviation function* with respect to the information set j and the continuation strategy $\hat{\pi}$ is any linear function $f : \mathbb{R}^{|\Sigma^{(i)}|} \rightarrow \mathbb{R}^{|\Sigma^{(i)}|}$ with the following properties:

- $f(\pi) = \pi$ for any $\pi \in \Pi^{(i)}$ such that $\pi[\sigma^{(i)}(j)] = 0$;
- For any $\sigma \in \Sigma^{(i)}$ and $\pi \in \Pi^{(i)}$ such that $\pi[\sigma^{(i)}(j)] = 1$,

$$f(\pi)[\sigma] = \begin{cases} \pi[\sigma] & \text{if } \sigma \not\succeq j; \\ \hat{\pi}[\sigma] & \text{if } \sigma \succeq j. \end{cases} \quad (53)$$

It will be convenient to instantiate a coarse trigger deviation function in the form of a linear mapping $\phi_{j \rightarrow \hat{\pi}}^{(i)} : \mathbb{R}^{|\Sigma^{(i)}|} \ni \mathbf{x} \mapsto \mathbf{M}_{j \rightarrow \hat{\pi}}^{(i)} \mathbf{x}$, where $\mathbf{M}_{j \rightarrow \hat{\pi}}^{(i)}$ is such that for any $\sigma_r, \sigma_c \in \Sigma^{(i)}$,

$$\mathbf{M}_{j \rightarrow \hat{\pi}}^{(i)}[\sigma_r, \sigma_c] = \begin{cases} 1 & \text{if } \sigma_c \not\succeq j \text{ \& } \sigma_r = \sigma_c; \\ \hat{\pi}[\sigma_r] & \text{if } \sigma_c = \sigma^{(i)}(j) \text{ \& } \sigma_r \succeq j; \\ 0 & \text{otherwise.} \end{cases} \quad (54)$$

We will also let $\tilde{\Psi}^{(i)} := \{\phi_{j \rightarrow \hat{\pi}}^{(i)} : j \in \mathcal{J}^{(i)}, \hat{\pi} \in \Pi_j^{(i)}\}$ be the set of all possible linear mappings defining trigger deviations functions for player i .

Definition C.2 (Coarse Trigger Regret Minimizer). A coarse trigger regret minimizer for a player $i \in [n]$ is any $\tilde{\Psi}^{(i)}$ -regret minimizer for the set of deterministic sequence-form strategies $\Pi^{(i)}$.

Similarly to EFCE, we can define extensive-form coarse correlated equilibria based on *coarse* trigger deviation functions:

Definition C.3 (EFCE). For $\epsilon \geq 0$, a probability distribution $\mu \in \Delta^{|\Pi|}$ is an ϵ -approximate EFCE if for every player $i \in [n]$ and every coarse trigger deviation function $\phi_{j \rightarrow \hat{\pi}}^{(i)} \in \tilde{\Psi}^{(i)}$, it holds that

$$\mathbb{E}_{\pi \sim \mu} \left[u^{(i)} \left(\phi_{j \rightarrow \hat{\pi}}^{(i)}(\pi^{(i)}), \pi^{(-i)} \right) - u^{(i)}(\pi) \right] \leq \epsilon, \quad (55)$$

where $\pi = (\pi_1, \dots, \pi_n) \in \Pi$. A probability distribution $\mu \in \Delta^{|\Pi|}$ is an EFCCE if it is a 0-EFCCE.

In particular, this equilibrium concept requires that the expected utility of any trigger agent $(j, \hat{\pi})$ is never larger (by more than an amount of ϵ) than the expected utility of the player when the mediator's recommendations are followed.

Theorem C.4 (Farina et al. (2021b)). *For every player $i \in [n]$, let $\pi^{(i),1}, \dots, \pi^{(i),T} \in \Pi^{(i)}$ be a sequence of deterministic sequence-form strategies whose cumulative $\tilde{\Psi}^{(i)}$ -regret is $R^{(i),T}$ with respect to the sequence of linear utility functions*

$$\ell^{(i),t} : \Pi^{(i)} \ni \pi^{(i)} \mapsto u^{(i)}(\pi^{(i)}, \pi^{(-i),t}). \quad (56)$$

Then, the empirical frequency of play $\mu \in \Delta^{|\Pi|}$ is an ϵ -EFCCE, where $\epsilon := \frac{1}{T} \max_{i \in [n]} R^{(i),T}$.

In this context, we can employ the general template we presented in Section 3, and in particular Theorem 3.2, in order to obtain accelerated no-regret dynamics for EFCCE. In fact, observe that a stable-predictive (external) regret minimizer for the set $\tilde{\Psi}^{(i)}$ can be directly obtained from Theorem 3.4, given that $\tilde{\Psi}^{(i)} \subseteq \Psi^{(i)}$. Thus, for the rest of the section we will focus on the second necessary ingredient, namely characterizing the stability of the fixed points.

C.1 STABILITY OF THE FIXED POINTS

Unlike all known methods for computing fixed points related to EFCE, it was recently shown by Farina et al. (2021b) (see also Morrill et al. (2021a)) that one can bypass the computation of a stationary distribution of a Markov chain using a succinct closed-form solution. In particular, this algorithm is presented in Algorithm 5.

Claim C.5 (Theorem 3, Farina et al. (2021b)). *Algorithm 5 returns a fixed point of $\phi \in \text{co } \tilde{\Psi}^{(i)}$, and it runs in time $O(|\Sigma^{(i)}| \mathfrak{D}^{(i)})$.*

Algorithm 5: FIXEDPOINT(ϕ) for $\phi \in \text{co } \tilde{\Psi}^{(i)}$

Input: $\phi = \sum_{j \in \mathcal{J}^{(i)}} \lambda[j] \phi_{j \rightarrow q_j} \in \text{co } \tilde{\Psi}^{(i)}$
Output: $q^* \in \mathcal{Q}^{(i)}$ such that $\phi(q^*) = q^*$

```

1  $q^* \leftarrow \mathbf{0} \in \mathbb{R}_{\geq 0}^{|\Sigma^{(i)}|}$ ,  $q^*[\emptyset] \leftarrow 1$ 
2 for  $\sigma = (j, a) \in \Sigma_*^{(i)}$  in top-down ( $\prec$ ) order do
3    $d_\sigma \leftarrow \sum_{j' \preceq j} \lambda[j']$ 
4   if  $d_\sigma = 0$  then
5      $q^*[\sigma] \leftarrow \frac{q^*[\sigma^{(i)}(j)]}{|\mathcal{A}(j)|}$ 
6   else
7      $q^*[\sigma] \leftarrow \frac{1}{d_\sigma} \sum_{j' \preceq j} \lambda[j'] q_{j'}[\sigma] q^*[\sigma^{(i)}(j')]$ 
8   return  $q^*$ 
```

Proposition C.6. *Under the conditions of Lemma 3.5, Algorithm 5 yields a sequence of $O(\kappa \mathfrak{D}^{(i)})$ -multiplicatively-stable strategies, where $\kappa = O(\eta D_\ell(|\Sigma^{(i)}| + D_x^{(i)}(\mathfrak{D}^{(i)})^2))$ for a sufficiently small $\kappa = O(1/\mathfrak{D}^{(i)})$.*

Proof. Observe that the sequence of λ^t 's and q_σ^t 's is κ -multiplicatively-stable (Lemma 3.5). Moreover, given that the sequence of λ^t 's is updated using OMWU, it will always be the case that $d_\sigma > 0$, for any $\sigma \in \Sigma^{(i)}$; in other words, Algorithm 5 will never visit the first “if” branch. Now fix any $t \geq 2$. We will show by induction that $q^{*,t}[\sigma]$ is such that $q^{*,t}[\sigma] \leq (1 + 2\kappa)^{3\mathfrak{D}[\sigma]-2} q^{*,t-1}[\sigma]$ and $q^{*,t}[\sigma] \geq (1 - \kappa)^{3\mathfrak{D}^{(i)}[\sigma]-2} q^{*,t-1}[\sigma]$, where $\mathfrak{D}^{(i)}[\sigma] \geq 1$ is the depth of sequence $\sigma \in \Sigma_*^{(i)}$ (considering only actions of player i). For the base case, let $\sigma = (j, a)$ such that $j \in \mathcal{J}^{(i)}$ corresponds to a root information set for player i . Then, it follows that $d_\sigma = \lambda[j]$, and subsequently, $q^{*,t} = q_j^t[\sigma]$.

Thus, $\mathbf{q}^{*,t}$ is indeed κ -multiplicatively stable. Next, consider some sequence $\sigma = (j, a)$ at depth $\mathfrak{D}^{(i)}[\sigma] \geq 2$ such that all ancestor sequences—i.e. all $\sigma^{(i)}(j')$ for $j' \preceq j$ —satisfy the inductive hypothesis. Then, it follows that

$$\mathbf{q}^{*,t}[\sigma] = \frac{\sum_{j' \preceq j} \lambda^t[j'] \mathbf{q}_{j'}^t[\sigma] \mathbf{q}^{*,t}[\sigma^{(i)}(j')]}{\sum_{j' \preceq j} \lambda^t[j']} \quad (57)$$

$$\leq \frac{(1 + \kappa)^2}{1 - \kappa} \frac{\sum_{j' \preceq j} \lambda^{t-1}[j'] \mathbf{q}_{j'}^{t-1}[\sigma] \mathbf{q}^{*,t}[\sigma^{(i)}(j')]}{\sum_{j' \preceq j} \lambda^{t-1}[j']} \quad (58)$$

$$\leq (1 + 2\kappa)^3 (1 + 2\kappa)^{3\mathfrak{D}^{(i)}[\sigma] - 5} \mathbf{q}^{*,t-1}[\sigma] \quad (59)$$

$$\leq (1 + 2\kappa)^{3\mathfrak{D}^{(i)}[\sigma] - 2} \mathbf{q}^{*,t-1}[\sigma],$$

where (57) is derived from Algorithm 5, (58) uses the κ -multiplicative stability of λ^t 's and $\mathbf{q}_{j'}^t$'s, and (59) leverages the inductive hypothesis, as well as the fact that $1/(1 - \kappa) \leq (1 + 2\kappa)$ for $\kappa \leq 1/2$. Similar reasoning yields:

$$\begin{aligned} \mathbf{q}^{*,t}[\sigma] &= \frac{\sum_{j' \preceq j} \lambda^t[j'] \mathbf{q}_{j'}^t[\sigma] \mathbf{q}^{*,t}[\sigma^{(i)}(j')]}{\sum_{j' \preceq j} \lambda^t[j']} \\ &\geq \frac{(1 - \kappa)^2}{1 + \kappa} \frac{\sum_{j' \preceq j} \lambda^{t-1}[j'] \mathbf{q}_{j'}^{t-1}[\sigma] \mathbf{q}^{*,t}[\sigma^{(i)}(j')]}{\sum_{j' \preceq j} \lambda^{t-1}[j']} \\ &\geq (1 - \kappa)^3 (1 - \kappa)^{3\mathfrak{D}^{(i)}[\sigma] - 5} \mathbf{q}^{*,t-1}[\sigma] \\ &\geq (1 - \kappa)^{3\mathfrak{D}^{(i)}[\sigma] - 2} \mathbf{q}^{*,t-1}[\sigma]. \end{aligned}$$

As a result, if $\mathfrak{D}^{(i)}$ is the maximum depth of i 's sequences, we can conclude that $\mathbf{q}^{*,t}[\sigma] \leq (1 + 2\kappa)^{3\mathfrak{D}^{(i)} - 2} \mathbf{q}^{*,t-1}[\sigma] \leq e^{6\mathfrak{D}^{(i)}\kappa - 4\kappa} \mathbf{q}^{*,t-1}[\sigma] \leq (1 + O(\mathfrak{D}^{(i)}\kappa)) \mathbf{q}^{*,t-1}[\sigma]$, where we used the inequality $e^x \leq 1 + 2x$ for $x \in [0, 1/2]$. Moreover, we have that $\mathbf{q}^{*,t}[\sigma] \geq (1 - \kappa)^{3\mathfrak{D}^{(i)} - 2} \mathbf{q}^{*,t-1}[\sigma] \geq (1 - O(\mathfrak{D}^{(i)}\kappa)) \mathbf{q}^{*,t-1}[\sigma]$, by Bernoulli's inequality. \square

Observe that the derived bound on stability is slightly better compared to that for EFCE (Corollary 3.11). Consequently, having established the stability of the fixed points, we can apply Theorem 3.2 to derive a $\tilde{\Psi}^{(i)}$ stable-predictive regret minimizer, for each player $i \in [n]$. Then, using the same steps as in Appendix B.3 we arrive at the following conclusion:

Corollary C.7 (Accelerated Convergence to EFCCE). *There exist uncoupled no-regret learning dynamics so that after T iterations the correlated distribution of play converges to an extensive-form coarse correlated equilibrium at a rate of $O(T^{-3/4})$.*

D DESCRIPTION OF GAME INSTANCES USED IN THE EXPERIMENTS

In this section we provide a description of the game instances used in our experiments. The parameters associated with each game are summarized in Table 2.

Kuhn poker. First, we experimented on a *three-player* variant of the popular benchmark game known as *Kuhn poker* (Kuhn, 1950). In our version, a deck of three cards—a Jack, a Queen, and a King—is employed. Players initially commit a single chip to the pot, and privately receive a single card. The first player can either *check* or *bet* (i.e. place an extra chip). Then, the second player can in turn check or bet if the first player checked, or *fold/call* in response to the first player's bet. If no betting occurred in the previous rounds, the third player can either check or bet. In the contrary case, the player can either fold or call. Following a bet of the second player (or respectively the third player), the first player (or respectively the first and the second players) has to decide whether to fold or to call. At the *showdown*, the player with the *highest* card, who has not folded in a previous round, gets to win all the chips committed in the pot.

Sheriff. Our second benchmark is a bargaining game inspired by the board game *Sheriff of Nottingham*, and introduced by Farina et al. (2019e). In particular, we used the *baseline* version of the

game described. In particular, this game consists of two players: the *Smuggler* and the *Sheriff*. The smuggler must originally come up with a number $n \in \{0, 1, 2, 3\}$ which corresponds to the number of illegal items to be loaded in the cargo. It is assumed that each illegal item has a fixed value of 1. Subsequently, 2 rounds of bargaining between the two players follow. At each round, the Smuggler decides on a bribe ranging from 0 to 3, and the Sheriff must decide whether or not the cargo will be inspected, given the bribe amount. The Sheriff’s decision is binding only in the last round of bargaining. In particular, if during the last round the Sheriff accepts the bribe, the game stops with the Smuggler obtaining a utility of n minus the bribe amount b that was proposed in the last bargaining round, while the Sheriff receives a utility equal to b . On the other hand, if the Sheriff does not accept the bribe in last bargaining round and decides to inspect the cargo, there are two possible alternatives. If the cargo has no illegal items (i.e. $n = 0$), the Smuggler receives the fixed amount of 3, while the utility of the Sheriff is set to be -3 . In the contrary case, the utility of the smuggler is assumed to be $-2n$, while the utility of the Sheriff is $2n$.

Liar’s dice. The final benchmark we experimented on is the game of *Liar’s dice*, introduced by Lisý et al. (2015). In the three-player variant, the beginning of the game sees each of the three players privately roll an unbiased 3-face die. Then, the players have to sequentially make claims about their private information. In particular, the first player may announce any face value up to 3, as well as the minimum number of dice that the player claims are showing that value among the dice of *all* players. Then, each player can either make a higher bid, or challenge the previous claim by declaring the previous agent a “liar”. More precisely, it is assumed that a bid is higher than the previous one if either the face value is higher, or if the claimed number of dice is greater. If the current claim is challenged, all the dice must be revealed. If the claim was valid, the last bidder wins and receives a reward of $+1$, while the challenger suffers a negative payoff of -1 . Otherwise, the utilities obtained are reversed. Any other player will receive 0 utility.

Game	Players	Decision points	Sequences	Leaves
Kuhn poker	3	36	75	78
Sheriff	2	73	222	256
Liar’s dice	3	1536	3069	13797

Table 2: The parameters associated with each game.