

521 **A Missing statements and proofs**

522 **A.1 Statements for Section 3.1**

523 **Claim A.1.** Let a two-player Markov game where both players affect the transition. Further, consider  
 524 a correlated policy  $\sigma$  and its corresponding marginalized product policy  $\pi^\sigma = \pi_1^\sigma \times \pi_2^\sigma$ . Then, for  
 525 any  $\pi'_1, \pi'_2$ ,

$$\begin{aligned} V_{k,1}^{\pi'_1, \sigma^{-1}}(s_1) &= V_{k,1}^{\pi'_1, \pi_2^\sigma}(s_1), \\ V_{k,2}^{\sigma^{-2}, \pi'_2}(s_1) &= V_{k,2}^{\pi_1^\sigma, \pi'_2}(s_1). \end{aligned}$$

526 **Proof.** We will effectively show that the problem of best-responding to a correlated policy  $\sigma$  is  
 527 equivalent to best-responding to the marginal policy of  $\sigma$  for the opponent. The proof follows from  
 528 the equivalence of the two MDPs.

529 As a reminder,

$$\begin{aligned} \pi_{1,h}(a|s) &= \sum_{b \in \mathcal{A}_2} \sigma_h(a, b|s) \\ \pi_{2,h}(b|s) &= \sum_{a \in \mathcal{A}_1} \sigma_h(a, b|s) \end{aligned}$$

530 As we have seen in Section 2.1, in the case of unilateral deviation from joint policy  $\sigma$ , an agent  
 531 faces a single agent MDP. More specifically, agent 2, best-responds by optimizing a reward function  
 532  $\bar{r}_{2,h}(s, b)$  under a transition kernel  $\bar{\mathbb{P}}_2$  for which,

$$\bar{r}_{2,h}(s, b) = \mathbb{E}_{b \sim \sigma} [r_{2,h}(s, a, b)] = \mathbb{E}_{b \sim \pi_1^\sigma} [r_{2,h}(s, a, b)] = r_{2,h}(s, \pi_1^\sigma, b).$$

533 Similarly,

$$\bar{r}_{1,h}(s, b) = r_{1,h}(s, a, \pi_2^\sigma).$$

534 Analogously, for each of the transition kernels,

$$\bar{\mathbb{P}}_{2,h}(s'|s, b) = \mathbb{E}_{a \sim \sigma} [\mathbb{P}_{2,h}(s'|s, a, b)] = \mathbb{E}_{a \sim \pi_2^\sigma} [\mathbb{P}_{2,h}(s'|s, a, b)] = \mathbb{P}_{2,h}(s'|s, \pi_1^\sigma, b),$$

535 as for agent 1,

$$\bar{\mathbb{P}}_{1,h}(s'|s, a) = \mathbb{P}_{1,h}(s'|s, a, \pi_2^\sigma).$$

536 Hence, it follows that,  $V_{2,1}^{\sigma^{-2} \times \pi'_2}(s_1) = V_{2,1}^{\pi_1^\sigma \times \pi'_2}(s_1), \forall \pi'_2$  and  $V_{1,1}^{\pi'_1 \times \sigma^{-1}}(s_1) =$   
 537  $V_{1,1}^{\pi'_1 \times \pi_2^\sigma}(s_1), \forall \pi'_1$ .

538 □

539 Before that, given a (possibly correlated) joint policy  $\sigma$  we define a nonlinear program, ( $\text{P}_{\text{BR}}$ ), whose  
 540 optimal solutions are best-response policies of each agent  $k$  to  $\sigma_{-k}$  and the values for each state  $s$   
 541 and timestep  $h$ :

542 **A.2 Proof of Theorem 3.2**

543 **The best-response program.** First, we state the following lemma that will prove useful for several  
 544 of our arguments,

545 **Lemma A.1** (Best-response LP). Let a (possibly correlated) joint policy  $\hat{\sigma}$ . Consider the following  
 546 linear program with variables  $w \in \mathbb{R}^{n \times H \times S}$ ,

$$\begin{aligned}
& \min \quad \sum_{k \in [n]} w_{k,s}(s_1) - \mathbf{e}_{s_1}^\top \sum_{h=1}^H \left( \prod_{\tau=1}^h \mathbb{P}_\tau(\hat{\sigma}_\tau) \right) \mathbf{r}_{k,h}(\hat{\sigma}_h) \\
& \text{s.t. } w_{k,h}(s) \geq r_{k,h}(s, a, \hat{\sigma}_{-k,h}) + \mathbb{P}_h(s, a, \hat{\sigma}_{-k,h}) \mathbf{w}_{k,h+1}, \\
& \quad \forall s \in \mathcal{S}, \forall h \in [H], \forall k \in [n], \forall a \in \mathcal{A}_k; \\
& \quad w_{k,H}(s) = 0, \quad \forall k \in [n], \forall s \in \mathcal{S}.
\end{aligned}
\tag{P_{BR}}$$

548 The optimal solution  $\mathbf{w}^\dagger$  of the program is unique and corresponds to the value function of each  
549 player  $k \in [n]$  when player  $k$  best-responds to  $\hat{\sigma}$ .

550 **Proof.** We observe that the program is separable to  $n$  independent linear programs, each with  
551 variables  $\mathbf{w}_k \in \mathbb{R}^{n \times H}$ ,

$$\begin{aligned}
& \min w_{k,1}(s_1) \\
& \text{s.t. } w_{k,h}(s) \geq r_{k,h}(s, a, \hat{\sigma}_{-k,h}) + \mathbb{P}_h(s, a, \hat{\sigma}_{-k,h}) \mathbf{w}_{k,h+1}, \\
& \quad \forall s \in \mathcal{S}, \forall h \in [H], \forall a \in \mathcal{A}_k; \\
& \quad w_{k,H}(s) = 0, \quad \forall k \in [n], \forall s \in \mathcal{S}.
\end{aligned}$$

552 Each of these linear programs describes the problem of a single agent MDP (Neu and Pike-Burke,  
553 2020, Section 2) —that agent being  $k$ — which, as we have seen in Best-response policies, is  
554 equivalent to the problem of finding a best-response to  $\hat{\sigma}_{-k}$ . It follows that the optimal  $\mathbf{w}_k^\dagger$  for every  
555 program is unique (each program corresponds to a set of Bellman optimality equations).  $\square$

556 **Properties of the NE program.** Second, we need to prove that the minimum value of the objective  
557 function of the program is nonnegative.

558 **Lemma A.2** (Feasibility of  $(P'_{NE})$  and global optimum). The nonlinear program  $(P'_{NE})$  is feasible,  
559 has a nonnegative objective value, and its global minimum is equal to 0.

560 **Proof.** Analogously to the finite-horizon case, for the feasibility of the nonlinear program, we invoke  
561 the theorem of the existence of a Nash equilibrium. We let a NE product policy,  $\boldsymbol{\pi}^*$ , and a vector  
562  $\mathbf{w}^* \in \mathbb{R}^{n \times S}$  such that  $w_k^*(s) = V_k^{\dagger, \boldsymbol{\pi}^*} (s)$ ,  $\forall k \in [n] \times \mathcal{S}$ .

563 By Lemma A.1, we know that  $(\boldsymbol{\pi}^*, \mathbf{w}^*)$  satisfies all the constraints of  $(P_{NE})$ . Additionally, because  
564  $\boldsymbol{\pi}^*$  is a NE,  $V_{k,h}^{\boldsymbol{\pi}^*}(s_1) = V_{k,h}^{\dagger, \boldsymbol{\pi}^*} (s_1)$  for all  $k \in [n]$ . Observing that,

$$w_{k,1}^*(s_1) - \mathbf{e}_{s_1}^\top \sum_{h=1}^H \left( \prod_{\tau=1}^h \mathbb{P}_\tau(\boldsymbol{\pi}_\tau^*) \right) \mathbf{r}_{k,h}(\boldsymbol{\pi}_h^*) = V_{k,h}^{\dagger, \boldsymbol{\pi}^*} (s_1) - V_{k,h}^{\boldsymbol{\pi}^*} (s_1) = 0,$$

565 concludes the argument that a NE attains an objective value equal to 0.

566 Continuing, we observe that due to (1) the objective function can be equivalently rewritten as,

$$\begin{aligned}
& \sum_{k \in [n]} \left( w_{k,1}(s_1) - \mathbf{e}_{s_1}^\top \sum_{h=1}^H \left( \prod_{\tau=1}^h \mathbb{P}_\tau(\boldsymbol{\pi}_\tau) \right) \mathbf{r}_{k,h}(\boldsymbol{\pi}_h) \right) \\
& = \sum_{k \in [n]} w_{k,1}(s_1) - \mathbf{e}_{s_1}^\top \sum_{h=1}^H \left( \prod_{\tau=1}^h \mathbb{P}_\tau(\boldsymbol{\pi}_\tau) \right) \sum_{k \in [n]} \mathbf{r}_{k,h}(\boldsymbol{\pi}_h) \\
& = \sum_{k \in [n]} w_{k,1}(s_1).
\end{aligned}$$

567 Next, we focus on the inequality constraint

$$w_{k,h}(s) \geq r_{k,h}(s, a, \boldsymbol{\pi}_{-k,h}) + \mathbb{P}_h(s, a, \boldsymbol{\pi}_{-k,h}) \mathbf{w}_{k,h+1}$$

568 which holds for all  $s \in \mathcal{S}$ , all players  $k \in [n]$ , all  $a \in \mathcal{A}_k$ , and all timesteps  $h \in [H - 1]$ .

569 By summing over  $a \in \mathcal{A}_k$  while multiplying each term with a corresponding coefficient  $\pi_{k,h}(a|s)$ ,  
 570 the display written in an equivalent element-wise vector inequality reads:

$$\mathbf{w}_{k,h} \geq \mathbf{r}_{k,h}(\boldsymbol{\pi}_h) + \mathbb{P}_h(\boldsymbol{\pi}_h)\mathbf{w}_{k,h+1}.$$

571 Finally, after consecutively substituting  $\mathbf{w}_{k,h+1}$  with the element-wise lesser term  $\mathbf{r}_{k,h+1}(\boldsymbol{\pi}_{h+1}) +$   
 572  $\mathbb{P}_{h+1}(\boldsymbol{\pi}_{h+1})\mathbf{w}_{k,h+2}$ , we end up with the inequality:

$$\mathbf{w}_{k,1} \geq \sum_{h=1}^H \left( \prod_{\tau=1}^h \mathbb{P}_\tau(\boldsymbol{\pi}_\tau) \right) \mathbf{r}_{k,h}(\boldsymbol{\pi}_h). \quad (5)$$

573 Summing over  $k$ , it holds for the  $s_1$ -th entry of the inequality,

$$\sum_{k \in [n]} w_{k,1} \geq \sum_{k \in [n]} \sum_{h=1}^H \left( \prod_{\tau=1}^h \mathbb{P}_\tau(\boldsymbol{\pi}_\tau) \right) \mathbf{r}_{k,h}(\boldsymbol{\pi}_h) = 0.$$

574 Where the equality holds due to the zero-sum property, (1).  $\square$

575 **An approximate NE is an approximate global minimum.** We show that an  $\epsilon$ -approximate NE,  
 576  $\boldsymbol{\pi}^*$ , achieves an  $n\epsilon$ -approximate global minimum of the program. Utilizing Lemma A.1, setting  
 577  $w_{k,1}^*(s_1) = V_{k,1}^{\dagger, \boldsymbol{\pi}^*} (s_1)$ , and the definition of an  $\epsilon$ -approximate NE we see that,

$$\begin{aligned} \sum_{k \in [n]} \left( w_{k,1}^*(s_1) - \mathbf{e}_{s_1}^\top \sum_{h=1}^H \left( \prod_{\tau=1}^h \mathbb{P}_\tau(\boldsymbol{\pi}_\tau^*) \right) \mathbf{r}_{k,h}(\boldsymbol{\pi}_h^*) \right) &= \sum_{k \in [n]} \left( w_{k,1}^*(s_1) - V_{k,1}^{\boldsymbol{\pi}^*}(s_1) \right) \\ &\leq \sum_{k \in [n]} \epsilon = n\epsilon. \end{aligned}$$

578 Indeed, this means that  $\boldsymbol{\pi}^*, \mathbf{w}^*$  is an  $n\epsilon$ -approximate global minimizer of (P<sub>NE</sub>).

579 **An approximate global minimum is an approximate NE.** For the opposite direction, we let a  
 580 feasible  $\epsilon$ -approximate global minimizer of the program (P<sub>NE</sub>),  $(\boldsymbol{\pi}^*, \mathbf{w}^*)$ . Because a global minimum  
 581 of the program is equal to 0, an  $\epsilon$ -approximate global optimum must be at most  $\epsilon > 0$ . We observe  
 582 that for every  $k \in [n]$ ,

$$w_{k,1}^*(s_1) \geq \mathbf{e}_{s_1}^\top \sum_{h=1}^H \left( \prod_{\tau=1}^h \mathbb{P}_\tau(\boldsymbol{\pi}_\tau^*) \right) \mathbf{r}_{k,h}(\boldsymbol{\pi}_h^*), \quad (6)$$

583 which follows from induction on the inequality constraint over all  $h$  similar to (5).

584 Consequently, the assumption that

$$\epsilon \geq \sum_{k \in [n]} \left( w_{k,1}^*(s_1) - \mathbf{e}_{s_1}^\top \sum_{h=1}^H \left( \prod_{\tau=1}^h \mathbb{P}_\tau(\boldsymbol{\pi}_\tau^*) \right) \mathbf{r}_{k,h}(\boldsymbol{\pi}_h^*) \right),$$

585 and Equation (6), yields the fact that

$$\begin{aligned} \epsilon &\geq w_{k,1}^*(s_1) - \mathbf{e}_{s_1}^\top \sum_{h=1}^H \left( \prod_{\tau=1}^h \mathbb{P}_\tau(\boldsymbol{\pi}_\tau^*) \right) \mathbf{r}_{k,h}(\boldsymbol{\pi}_h^*) \\ &\geq V_{k,1}^{\dagger, \boldsymbol{\pi}^*} (s_1) - V_{k,1}^{\boldsymbol{\pi}^*} (s_1), \end{aligned}$$

586 where the second inequality holds from the fact that  $\mathbf{w}^*$  is feasible for (P<sub>BR</sub>). The latter concludes  
 587 the proof, as the display coincides with the definition of an  $\epsilon$ -approximate NE.

588 **A.3 Proof of Claim 3.1**

589 **Proof.** The value function of  $s_1$  for  $h = 1$  of players 1 and 2 read:

$$\begin{aligned} V_{1,1}^\sigma(s_1) &= e_{s_1}^\top (\mathbf{r}_1(\boldsymbol{\sigma}) + \mathbb{P}(\boldsymbol{\sigma})\mathbf{r}_1(\boldsymbol{\sigma})) \\ &= -\frac{9\sigma(a_1, b_1|s_1)}{20} + \frac{\sigma(a_1, b_2|s_1)}{20} + \frac{(1 - \sigma(a_1, b_1|s_1))(\sigma(a_1, b_1|s_2) + \sigma(a_1, b_2|s_2))}{20}, \end{aligned}$$

590 and,

$$\begin{aligned} V_{2,1}^\sigma(s_1) &= e_{s_1}^\top (\mathbf{r}_2(\boldsymbol{\sigma}) + \mathbb{P}(\boldsymbol{\sigma})\mathbf{r}_2(\boldsymbol{\sigma})) \\ &= -\frac{9\sigma(a_1, b_1|s_1)}{20} + \frac{\sigma(a_2, b_2|s_1)}{20} + \frac{(1 - \sigma(a_1, b_1|s_1))(\sigma(a_1, b_1|s_2) + \sigma(a_2, b_1|s_2))}{20}. \end{aligned}$$

591 We are indifferent to the corresponding value function of player 3 as they only have one available  
592 action per state and hence, cannot affect their rewards. For the joint policy  $\boldsymbol{\sigma}$ , the corresponding  
593 value functions of both players 1 and 2 are  $V_{1,1}^\sigma(s_1) = V_{2,1}^\sigma(s_1) = \frac{1}{20}$ .

594 **Deviations.** We will now prove that no deviation of player 1 manages to accumulate a reward  
595 greater than  $\frac{1}{20}$ . The same follows for player 2 due to symmetry.

596 When a player deviates unilaterally from a joint policy, they experience a single agent Markov  
597 decision process (MDP). It is well-known that MDPs always have a deterministic optimal policy.  
598 As such, it suffices to check whether  $V_{1,1}^{\pi_1, \sigma^{-1}}(s_1)$  is greater than  $\frac{1}{20}$  for any of the four possible  
599 deterministic policies:

$$\begin{aligned} 600 \quad & \bullet \pi_1(s_1) = \pi_1(s_2) = (1 \ 0), & 602 \quad & \bullet \pi_1(s_1) = (1 \ 0), \pi_1(s_2) = (0 \ 1), \\ 601 \quad & \bullet \pi_1(s_1) = \pi_1(s_2) = (0 \ 1), & 603 \quad & \bullet \pi_1(s_1) = (0 \ 1), \pi_1(s_2) = (1 \ 0). \end{aligned}$$

604 Finally, the value function of any deviation  $\pi'_1$  writes,

$$V_{1,1}^{\pi'_1 \times \sigma^{-1}}(s_1) = -\frac{\pi'_1(a_1|s_1)}{5} - \frac{\pi'_1(a_1|s_2)(\pi'_1(a_1|s_1) - 2)}{40}.$$

605 We can now check that for all deterministic policies  $V_{1,1}^{\pi'_1 \times \sigma^{-1}}(s_1) \leq \frac{1}{20}$ . By symmetry, it follows  
606 that  $V_{2,1}^{\pi'_2 \times \sigma^{-2}}(s_1) \leq \frac{1}{20}$  and as such  $\boldsymbol{\sigma}$  is indeed a CCE.  $\square$

607 **A.4 Proof of Claim 3.2**

608 **Proof.** In general, the value functions of each player 1 and 2 are:

$$V_{1,1}^{\pi_1 \times \pi_2}(s_1) = -\frac{\pi_1(a_1|s_1)\pi_2(b_1|s_1)}{2} + \frac{\pi_1(a_1|s_1)}{20} - \frac{\pi_1(a_1|s_2)(\pi_1(a_1|s_1)\pi_2(b_1|s_1) - 1)}{20},$$

609 and

$$V_{2,1}^{\pi_1 \times \pi_2}(s_1) = -\frac{\pi_1(a_1|s_1)\pi_2(b_1|s_1)}{2} + \frac{\pi_1(b_1|s_1)}{20} - \frac{\pi_1(b_1|s_2)(\pi_1(a_1|s_1)\pi_2(b_1|s_1) - 1)}{20}.$$

610 Plugging in  $\pi_1^\sigma, \pi_2^\sigma$  yields  $V_{1,1}^{\pi_1^\sigma \times \pi_2^\sigma}(s_1) = V_{2,1}^{\pi_1^\sigma \times \pi_2^\sigma}(s_1) = -\frac{13}{160}$ . But, if player 1 deviates to say  
611  $\pi'_1(s_1) = \pi'_1(s_2) = (0 \ 1)$ , they get a value equal to 0 which is clearly greater than  $-\frac{13}{160}$ . Hence,  
612  $\pi_1^\sigma \times \pi_2^\sigma$  is not a NE.  $\square$

613 **A.5 Proof of Theorem 3.4**

614 **Proof.** The proof follows from the game of Example 1, and Claims 3.1 and 3.2.  $\square$

## 615 B Proofs for infinite-horizon Zero-Sum Polymatrix Markov Games

616 In this section we will explicitly state definitions, theorems and proofs relating to the infinite-horizon  
617 discounted zero-sum polymatrix Markov games.

### 618 B.1 Definitions of equilibria for the infinite-horizon

619 Let us restate the definition specifically for infinite-horizon Markov games. They are defined as a  
620 tuple  $\Gamma(H, \mathcal{S}, \{\mathcal{A}_k\}_{k \in [n]}, \mathbb{P}, \{r_k\}_{k \in [n]}, \gamma, \rho)$ .

- 621 •  $H = \infty$  denotes the *time horizon*
- 622 •  $\mathcal{S}$ , with cardinality  $S := |\mathcal{S}|$ , stands for the state space,
- 623 •  $\{\mathcal{A}_k\}_{k \in [n]}$  is the collection of every player's action space, while  $\mathcal{A} := \mathcal{A}_1 \times \dots \times \mathcal{A}_n$   
624 denotes the *joint action space*; further, an element of that set—a joint action—is generally  
625 noted as  $\mathbf{a} = (a_1, \dots, a_n) \in \mathcal{A}$ ,
- 626 •  $\mathbb{P} : \mathcal{S} \times \mathcal{A} \rightarrow \Delta(\mathcal{S})$  is the transition probability function,
- 627 •  $r_k : \mathcal{S}, \mathcal{A} \rightarrow [-1, 1]$  yields the reward of player  $k$  at a given state and joint action,
- 628 • a discount factor  $0 < \gamma < 1$ ,
- 629 • an initial state distribution  $\rho \in \Delta(\mathcal{S})$ .

630 **Policies and value functions.** In infinite-horizon Markov games policies can still be distinguished  
631 in two main ways, *Markovian/non-Markovian* and *stationary/nonstationary*. Moreover, a joint policy  
632 can be a *correlated* policy or a *product* policy.

633 *Markovian* policies attribute a probability over the simplex of actions solely depending on the running  
634 state  $s$  of the game. On the other hand, *non-Markovian* policies attribute a probability over the  
635 simplex of actions that depends on any subset of the history of the game. *I.e.*, they can depend on any  
636 sub-sequence of actions and states up until the running timestep of the horizon.

637 *Stationary* policies are those that will attribute the same probability distribution over the simplex  
638 of actions for every timestep of the horizon. *Nonstationary* policies, on the contrary can change  
639 depending on the timestep of the horizon.

640 A joint Markovian stationary policy  $\sigma$  is said to be *correlated* when for every state  $s \in \mathcal{S}$ , attributes  
641 a probability distribution over the simplex of joint actions  $\mathcal{A}$  for all players, *i.e.*,  $\sigma(s) \in \Delta(\mathcal{A})$ .  
642 A Markovian stationary policy  $\pi$  is said to be a *product* policy when for every  $s \in \mathcal{S}$ ,  $\pi(s) \in$   
643  $\prod_{k=1}^n \Delta(\mathcal{A}_k)$ . It is rather easy to define *correlated/product* policies for the case of non-Markovian  
644 and nonstationary policies.

645 Given a Markovian stationary policy  $\pi$ , the value function for an infinite-horizon discounted game is  
646 defined as,

$$V_k^\pi(s_1) = \mathbb{E}_\pi \left[ \sum_{h=1}^H \gamma^{h-1} r_{k,h}(s_h, \mathbf{a}_h) \mid s_1 \right] = e_{s_1}^\top \sum_{h=1}^H \left( \gamma^{h-1} \prod_{\tau=1}^h \mathbb{P}_\tau(\pi_\tau) \right) \mathbf{r}_{k,h}(\pi_h).$$

647 It is possible to express the value function of each player  $k$  in the following way,

$$V_k^\pi(s_1) = e_{s_1}^\top (\mathbf{I} - \gamma \mathbb{P}(\pi))^{-1} \mathbf{r}(\pi).$$

648 Where  $\mathbf{I}$  is the identity matrix of appropriate dimensions. Also, when the initial state is drawn from  
649 the initial state distribution, we denote, the value function reads  $V_k^\pi(\rho) = \rho^\top (\mathbf{I} - \gamma \mathbb{P}(\pi))^{-1} \mathbf{r}(\pi)$ .

650 **Best-response policies.** Given an arbitrary joint policy  $\sigma$  (which can be either a correlated or  
651 product policy), a best-response policy of a player  $k$  is defined to be  $\pi_k^\dagger \in \Delta(\mathcal{A}_k)^S$  such that  
652  $\pi_k^\dagger \in \arg \max_{\pi_k'} V_k^{\pi_k' \times \sigma^{-k}}(s)$ . Also, we will denote  $V_k^{\dagger, \sigma^{-k}}(s) = \max_{\pi_k'} V_k^{\pi_k', \sigma^{-k}}(s)$ . It is rather  
653 straightforward to see that the problem of computing a best-response to a given policy is equivalent  
654 to solving a single-agent MDP problem.

655 **Notions of equilibria.** Now that best-response policies have been defined, it is straightforward to  
 656 define the different notions of equilibria. First, we define the notion of a coarse-correlated equilibrium.

657 **Definition B.1** (CCE—infinite-horizon). *A joint (potentially correlated) policy  $\sigma \in \Delta(\mathcal{A})^S$  is an*  
 658  *$\epsilon$ -approximate coarse-correlated equilibrium if it holds that for an  $\epsilon$ ,*

$$V_k^\dagger, \sigma^{-k}(\rho) - V_k^\sigma(\rho) \leq \epsilon, \forall k \in [n].$$

659 Second, we define the notion of a Nash equilibrium. The main difference of the definition of the  
 660 coarse-correlated equilibrium, is the fact that a NE Markovian stationary policy is a *product policy*.

661 **Definition B.2** (NE—infinite-horizon). *A joint (potentially correlated) policy  $\pi \in \prod_{k \in [n]} \Delta(\mathcal{A}_k)^S$*   
 662 *is an  $\epsilon$ -approximate coarse-correlated equilibrium if it holds that for an  $\epsilon$ ,*

$$V_k^\dagger, \pi^{-k}(\rho) - V_k^\pi(\rho) \leq \epsilon, \forall k \in [n].$$

663 As it is folklore by now, infinite-horizon discounted Markov games have a stationary Markovian Nash  
 664 equilibrium.

## 665 C Main results for infinite-horizon games

666 The workhorse of our arguments in the following results is still the following nonlinear program with  
 667 variables  $\pi, \mathbf{w}$ ,

$$\begin{aligned} & \min \sum_{k \in [n]} \rho^\top (\mathbf{w}_k - (\mathbf{I} - \gamma \mathbb{P}(\pi))^{-1} \mathbf{r}_k(\pi)) \\ & \text{s.t. } w_k(s) \geq r_k(s, a, \pi_{-k}) + \gamma \mathbb{P}(s, a, \pi_{-k}) \mathbf{w}_k, \\ & \quad \forall s \in \mathcal{S}, \forall k \in [n], \forall a \in \mathcal{A}_k; \\ & \quad \pi_k(s) \in \Delta(\mathcal{A}_k), \\ & \quad \forall s \in \mathcal{S}, \forall k \in [n], \forall a \in \mathcal{A}_k. \end{aligned}$$

668 (P'\_{NE})

669 As we will prove, approximate NE's correspond to approximate global minima of (P'\_{NE}) and vice-  
 670 versa. Before that, we need some intermediate lemmas. The first lemma we prove is about the  
 671 best-response program.

672 **The best-response program.** Even for the infinite-horizon, we can define a linear program for the  
 673 best-responses of all players. That program is the following, with variables  $\mathbf{w}$ ,

$$\begin{aligned} & \min \sum_{k \in [n]} \rho^\top (\mathbf{w}_k - (\mathbf{I} - \gamma \mathbb{P}(\hat{\sigma}))^{-1} \mathbf{r}_k(\hat{\sigma})) \\ & \text{s.t. } w_k(s) \geq r_k(s, a, \hat{\sigma}_{-k}) + \mathbb{P}(s, a, \hat{\sigma}_{-k}) \mathbf{w}_k, \\ & \quad \forall s \in \mathcal{S}, \forall k \in [n], \forall a \in \mathcal{A}_k. \end{aligned}$$

674 (P'\_{BR})

675 **Lemma C.1** (Best-response LP—infinite-horizon). Let a (possibly correlated) joint policy  $\hat{\sigma}$ . Con-  
 676 sider the linear program (P'\_{BR}). The optimal solution  $\mathbf{w}^\dagger$  of the program is unique and corresponds  
 677 to the value function of each player  $k \in [n]$  when player  $k$  best-responds to  $\hat{\sigma}$ .

678 **Proof.** We observe that the program is separable to  $n$  independent linear programs, each with  
 679 variables  $\mathbf{w}_k \in \mathbb{R}^n$ ,

$$\begin{aligned} & \min \rho^\top \mathbf{w}_k \\ & \text{s.t. } w_k(s) \geq r_k(s, a, \hat{\sigma}_{-k}) + \gamma \mathbb{P}(s, a, \hat{\sigma}_{-k}) \mathbf{w}_k, \\ & \quad \forall s \in \mathcal{S}, \forall a \in \mathcal{A}_k. \end{aligned}$$

680 Each of these linear programs describes the problem of a single agent MDP—that agent being  $k$ .  
 681 It follows that the optimal  $\mathbf{w}_k^\dagger$  for every program is unique (each program corresponds to a set of  
 682 Bellman optimality equations).  $\square$

683 **Properties of the NE program.** Second, we need to prove that the minimum value of the objective  
 684 function of the program is nonnegative.

685 **Lemma C.2** (Feasibility of  $(P'_{NE})$  and global optimum). The nonlinear program  $(P'_{NE})$  is feasible,  
 686 has a nonnegative objective value, and its global minimum is equal to 0.

687 **Proof.** For the feasibility of the nonlinear program, we invoke the theorem of the existence of  
 688 a Nash equilibrium. *i.e.*, let a NE product policy,  $\pi^*$ , and a vector  $w^* \in \mathbb{R}^{n \times H \times S}$  such that  
 689  $w_{k,s}^*(s) = V_k^{\dagger, \pi^*} (s)$ ,  $\forall k \in [n] \times \mathcal{S}$ .

690 By Lemma C.1, we know that  $(\pi^*, w^*)$  satisfies all the constraints of  $(P'_{NE})$ . Additionally, because  
 691  $\pi^*$  is a NE,  $V_k^{\pi^*}(\rho) = V_k^{\dagger, \pi^*}(\rho)$  for all  $k \in [n]$ . Observing that,

$$\rho^\top (w_k^* - (\mathbf{I} - \gamma \mathbb{P}(\pi^*))^{-1} r_k(\pi^*)) = V_k^{\dagger, \pi^*}(\rho) - V_k^{\pi^*}(\rho) = 0,$$

692 concludes the argument that a NE attains an objective value equal to 0.

693 Continuing, we observe that due to (1) the objective function can be equivalently rewritten as,

$$\begin{aligned} & \sum_{k \in [n]} (\rho^\top w_k - \rho^\top (\mathbf{I} - \gamma \mathbb{P}(\pi))^{-1} r_k(\pi)) \\ &= \sum_{k \in [n]} \rho^\top w_k - \rho^\top (\mathbf{I} - \gamma \mathbb{P}(\pi))^{-1} \sum_{k \in [n]} r_k(\pi) \\ &= \sum_{k \in [n]} \rho^\top w_k. \end{aligned}$$

694 Next, we focus on the inequality constraint

$$w_k(s) \geq r_k(s, a, \pi_{-k}) + \gamma \mathbb{P}(s, a, \pi_{-k}) w_k$$

695 which holds for all  $s \in \mathcal{S}$ , all players  $k \in [n]$ , and all  $a \in \mathcal{A}_k$ .

696 By summing over  $a \in \mathcal{A}_k$  while multiplying each term with a corresponding coefficient  $\pi_k(a|s)$ , the  
 697 display written in an equivalent element-wise vector inequality reads:

$$w_k \geq r_{k,h}(\pi) + \gamma \mathbb{P}(\pi) w_k.$$

698 Finally, after consecutively substituting  $w_k$  with the element-wise lesser term  $r_k(\pi) + \gamma \mathbb{P}(\pi) w_k$ ,  
 699 we end up with the inequality:

$$w_k \geq (\mathbf{I} - \gamma \mathbb{P}(\pi))^{-1} r_k(\pi). \quad (9)$$

700 We note that  $\mathbf{I} + \gamma \mathbb{P}(\pi) + \gamma^2 \mathbb{P}^2(\pi) + \dots = (\mathbf{I} - \gamma \mathbb{P}(\pi))^{-1}$ .

701 Summing over  $k$ , it holds for the  $s_1$ -th entry of the inequality,

$$\sum_{k \in [n]} w_k \geq \sum_{k \in [n]} (\mathbf{I} - \gamma \mathbb{P}(\pi))^{-1} r_k(\pi) = (\mathbf{I} - \gamma \mathbb{P}(\pi))^{-1} \sum_{k \in [n]} r_k(\pi) = 0.$$

702 Where the equality holds due to the zero-sum property, (1).  $\square$

703 **Theorem C.1** (NE and global optima of  $(P'_{NE})$ —infinite-horizon). *If  $(\pi^*, w^*)$  yields an  $\epsilon$ -*  
 704 *approximate global minimum of  $(P'_{NE})$ , then  $\pi^*$  is an  $n\epsilon$ -approximate NE of the infinite-horizon*  
 705 *zero-sum polymatrix switching controller MG,  $\Gamma$ . Conversely, if  $\pi^*$  is an  $\epsilon$ -approximate NE of the*  
 706 *MG  $\Gamma$  with corresponding value function vector  $w^*$  such that  $w_k^*(s) = V_k^{\pi^*}(s) \forall (k, s) \in [n] \times \mathcal{S}$ ,*  
 707 *then  $(\pi^*, w^*)$  attains an  $\epsilon$ -approximate global minimum of  $(P'_{NE})$ .*

708 **Proof.**

709 **An approximate NE is an approximate global minimum.** We show that an  $\epsilon$ -approximate NE,  
 710  $\pi^*$ , achieves an  $n\epsilon$ -approximate global minimum of the program. Utilizing Lemma C.1 by setting  
 711  $w_k^* = V_k^{\dagger, \pi^*}(\rho)$ , feasibility, and the definition of an  $\epsilon$ -approximate NE we see that,

$$\begin{aligned} \sum_{k \in [n]} \left( \rho^\top \mathbf{w}_k^* - \rho^\top (\mathbf{I} - \gamma \mathbb{P}(\boldsymbol{\pi}^*))^{-1} \mathbf{r}_k(\boldsymbol{\pi}^*) \right) &= \sum_{k \in [n]} \left( \rho^\top \mathbf{w}_k^* - V_k^{\boldsymbol{\pi}^*}(\rho) \right) \\ &\leq \sum_{k \in [n]} \epsilon = n\epsilon. \end{aligned}$$

712 Indeed, this means that  $\boldsymbol{\pi}^*, \mathbf{w}^*$  is an  $n\epsilon$ -approximate global minimizer of  $(\mathbf{P}'_{\text{NE}})$ .

713 **An approximate global minimum is an approximate NE.** For this direction, we let a feasible  
714  $\epsilon$ -approximate global minimizer of the program  $(\mathbf{P}'_{\text{NE}})$ ,  $(\boldsymbol{\pi}^*, \mathbf{w}^*)$ . Because a global minimum of the  
715 program is equal to 0, an  $\epsilon$ -approximate global optimum must be at most  $\epsilon > 0$ . We observe that for  
716 every  $k \in [n]$ ,

$$\rho^\top \mathbf{w}_k^* \geq \rho^\top (\mathbf{I} - \gamma \mathbb{P}(\boldsymbol{\pi}^*))^{-1} \mathbf{r}_k(\boldsymbol{\pi}^*), \quad (10)$$

717 which follows from induction on the inequality constraint (9).

718 Consequently, the assumption that

$$\epsilon \geq \rho^\top \mathbf{w}_k^* - \rho^\top (\mathbf{I} - \gamma \mathbb{P}(\boldsymbol{\pi}^*))^{-1} \mathbf{r}_k(\boldsymbol{\pi}^*)$$

719 and Equation (10), yields the fact that

$$\begin{aligned} \epsilon &\geq \rho^\top \mathbf{w}_k^* - \rho^\top (\mathbf{I} - \gamma \mathbb{P}(\boldsymbol{\pi}^*))^{-1} \mathbf{r}_k(\boldsymbol{\pi}^*) \\ &\geq V_k^{\dagger, \boldsymbol{\pi}^*}(\rho) - V_k^{\boldsymbol{\pi}^*}(\rho), \end{aligned}$$

720 where the second inequality holds from the fact that  $\mathbf{w}^*$  is also feasible for  $(\mathbf{P}'_{\text{BR}})$ . The latter  
721 concludes the proof, as the display coincides with the definition of an  $\epsilon$ -approximate NE.  $\square$

722 **Theorem C.2** (CCE collapse to NE in polymatrix MG—infinite-horizon). *Let a zero-sum polymatrix*  
723 *switching-control Markov game, i.e., a Markov game for which Assumptions 1 and 2 hold. Further,*  
724 *let an  $\epsilon$ -approximate CCE of that game  $\boldsymbol{\sigma}$ . Then, the marginal product policy  $\boldsymbol{\pi}^\sigma$ , with  $\boldsymbol{\pi}_k^\sigma(a|s) =$   
725  $\sum_{\mathbf{a}_{-k} \in \mathcal{A}_{-k}} \boldsymbol{\sigma}(a, \mathbf{a}_{-k}), \forall k \in [n]$  is an  $n\epsilon$ -approximate NE.*

726 **Proof.** Let an  $\epsilon$ -approximate CCE policy,  $\boldsymbol{\sigma}$ , of game  $\Gamma$ . Moreover, let the best-response value-vectors  
727 of each agent  $k$  to joint policy  $\boldsymbol{\sigma}_{-k}, \mathbf{w}_k^\dagger$ .

728 Now, we observe that due to Assumption 1,

$$\begin{aligned} \mathbf{w}_k^\dagger(s) &\geq r_k(s, a, \boldsymbol{\sigma}_{-k}) + \mathbb{P}_h(s, a, \boldsymbol{\sigma}_{-k}) \mathbf{w}_k^\dagger \\ &= \sum_{j \in \text{adj}(k)} r_{(k,j),h}(s, a, \boldsymbol{\pi}_j^\sigma) + \mathbb{P}(s, a, \boldsymbol{\sigma}_{-k}) \mathbf{w}_k^\dagger. \end{aligned}$$

729 Further, due to Assumption 2,

$$\mathbb{P}(s, a, \boldsymbol{\sigma}_{-k}) \mathbf{w}_k^\dagger = \mathbb{P}(s, a, \boldsymbol{\pi}_{\text{argctrl}(s)}^\sigma) \mathbf{w}_k^\dagger,$$

730 or,

$$\mathbb{P}(s, a, \boldsymbol{\sigma}_{-k}) \mathbf{w}_k^\dagger = \mathbb{P}(s, a, \boldsymbol{\pi}^\sigma) \mathbf{w}_k^\dagger.$$

731 Putting these pieces together, we reach the conclusion that  $(\boldsymbol{\pi}^\sigma, \mathbf{w}^\dagger)$  is feasible for the nonlinear  
732 program  $(\mathbf{P}'_{\text{NE}})$ .

733 What is left is to prove that it is also an  $\epsilon$ -approximate global minimum. Indeed, if  $\sum_k \rho^\top \mathbf{w}_k^\dagger \leq \epsilon$   
734 (by assumption of an  $\epsilon$ -approximate CCE), then the objective function of  $(\mathbf{P}'_{\text{NE}})$  will attain an  
735  $\epsilon$ -approximate global minimum. In turn, due to Theorem C.1 the latter implies that  $\boldsymbol{\pi}^\sigma$  is an  
736  $n\epsilon$ -approximate NE.  $\square$

### 737 C.1 No equilibrium collapse with more than one controllers per-state

738 **Example 2.** *We consider the following 3-player Markov game that takes place for a time horizon*  
739  *$H = 3$ . There exist three states,  $s_1, s_2$ , and  $s_3$  and the game starts at state  $s_1$ . Player 3 has a*  
740 *single action in every state, while players 1 and 2 have two available actions  $\{a_1, a_2\}$  and  $\{b_1, b_2\}$*   
741 *respectively in every state. The initial state distribution  $\boldsymbol{\rho}$  is the uniform probability distribution over*  
742  *$\mathcal{S}$ .*

743 **Reward functions.** If player 1 (respectively, player 2) takes action  $a_1$  (resp.,  $b_1$ ), in either of the  
 744 states  $s_1$  or  $s_2$ , they get a reward equal to  $\frac{1}{20}$ . In state  $s_3$ , both players get a reward equal to  $-\frac{1}{2}$   
 745 regardless of the action they select. Player 3 always gets a reward that is equal to the negative sum  
 746 of the reward of the other two players. This way, the zero-sum polymatrix property of the game is  
 747 ensured (Assumption 1).

748 **Transition probabilities.** If players 1 and 2 select the joint action  $(a_1, b_1)$  in state  $s_1$ , the game  
 749 will transition to state  $s_2$ . In any other case, it will transition to state  $s_3$ . The converse happens if  
 750 in state  $s_2$  they take joint action  $(a_1, b_1)$ ; the game will transition to state  $s_3$ . For any other joint  
 751 action, it will transition to state  $s_1$ . From state  $s_3$ , the game transition to state  $s_1$  or  $s_2$  uniformly  
 752 at random.

753 At this point, it is important to notice that two players control the transition probability from one state  
 754 to another. In other words, Assumption 2 does not hold.

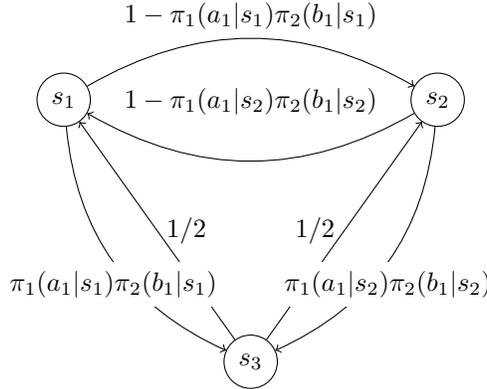


Figure 2: A graph of the state space with transition probabilities parametrized with respect to the policy of each player.

755 Next, we consider the joint policy  $\sigma$ ,

$$\sigma(s_1) = \sigma(s_2) = \begin{matrix} & b_1 & b_2 \\ a_1 & \begin{pmatrix} 0 & 1/2 \end{pmatrix} \\ a_2 & \begin{pmatrix} 1/2 & 0 \end{pmatrix} \end{matrix}.$$

756 **Claim C.1.** The joint policy  $\sigma$  that assigns probability  $\frac{1}{2}$  to the joint actions  $(a_1, b_2)$  and  $(a_2, b_1)$  in  
 757 both states  $s_1, s_2$  is a CCE and  $V_1^\sigma(\rho) = V_2^\sigma(\rho) = -\frac{1}{10}$ .

**Proof.**

$$\begin{aligned} V_1^\sigma(\rho) &= \rho^\top (\mathbf{I} - \gamma \mathbb{P}(\sigma))^{-1} \mathbf{r}_1(\sigma) \\ &= \begin{pmatrix} \frac{1}{3} & \frac{1}{3} & \frac{1}{3} \end{pmatrix} \begin{pmatrix} \frac{9}{5} & \frac{6}{5} & 0 \\ \frac{6}{5} & \frac{9}{5} & 0 \\ 1 & 1 & 1 \end{pmatrix} \begin{pmatrix} \frac{1}{40} \\ \frac{1}{40} \\ -\frac{1}{2} \end{pmatrix} \\ &= -\frac{1}{10}. \end{aligned}$$

758 We check every deviation,

759 •  $\pi_1(s_1) = \pi_1(s_2) = (1 \ 0), V^{\pi_1 \times \sigma^{-1}}(\rho) = -\frac{2}{5},$

760 •  $\pi_1(s_1) = \pi_1(s_2) = (0 \ 1), V^{\pi_1 \times \sigma^{-1}}(\rho) = -\frac{1}{6},$

761 •  $\pi_1(s_1) = (1 \ 0), \pi_1(s_2) = (0 \ 1), V^{\pi_1 \times \sigma^{-1}}(\rho) = -\frac{5}{16},$

762 •  $\pi_1(s_1) = (0 \ 1), \pi_1(s_2) = (1 \ 0), V^{\pi_1 \times \sigma^{-1}}(\rho) = -\frac{5}{16}.$

763 For every such deviation the value of player 1 is smaller than  $-\frac{1}{10}$ . For player 2, the same follows by  
 764 symmetry. Hence,  $\sigma$  is indeed a CCE. □

765

766 *Yet, the marginalized product policy of  $\sigma$  which we note as  $\pi_1^\sigma \times \pi_2^\sigma$  does not constitute a NE. The*  
 767 *components of this policy are,*

$$\left\{ \begin{array}{l} \pi_1^\sigma(s_1) = \pi_1^\sigma(s_2) = \begin{pmatrix} a_1 & a_2 \\ 1/2 & 1/2 \end{pmatrix}, \\ \pi_2^\sigma(s_1) = \pi_2^\sigma(s_2) = \begin{pmatrix} b_1 & b_2 \\ 1/2 & 1/2 \end{pmatrix}. \end{array} \right.$$

768 *I.e., the product policy  $\pi_1^\sigma \times \pi_2^\sigma$  selects any of the two actions of each player in states  $s_1, s_2$*   
 769 *independently and uniformly at random. With the following claim, it can be concluded that in*  
 770 *general when more than one player control the transition the set of equilibria do not collapse.*

771 **Claim C.2.** The product policy  $\pi_1^\sigma \times \pi_2^\sigma$  is not a NE.

772 **Proof.** For  $\pi^\sigma = \pi_1^\sigma \times \pi_2^\sigma$  we get,

$$\begin{aligned} V_1^{\pi^\sigma} &= \rho^\top (\mathbf{I} - \gamma \mathbb{P}(\pi^\sigma))^{-1} \mathbf{r}_1(\pi^\sigma) \\ &= \begin{pmatrix} 1/3 & 1/3 & 1/3 \end{pmatrix} \begin{pmatrix} 34 & 20 & 3 \\ 21 & 34 & 3 \\ 21 & 21 & 7 \\ 6 & 6 & 7 \end{pmatrix} \begin{pmatrix} 1 \\ 40 \\ 40 \\ -1/2 \end{pmatrix} \\ &= -\frac{3}{10}. \end{aligned}$$

773 But, for the deviation  $\pi_1(a_1|s_1) = \pi_1(a_1|s_2) = 0$ , the value function of player 1, is equal to  $-\frac{1}{6}$ .  
 774 Hence,  $\pi^\sigma$  is not a NE. □

775 *In conclusion, Assumption 1 does not suffice to ensure equilibrium collapse.*

776 **Theorem C.3** (No collapse—infinite-horizon). *There exists a zero-sum polymatrix Markov game*  
 777 *(Assumption 2 is not satisfied) that has a CCE which does not collapse to a NE.*

778 **Proof.** The proof follows from the game of Example 2, and Claims C.1 and C.2. □