Coalitions in Repeated Games^{*}

S. Nageeb Ali^{\dagger} Ce Liu^{\ddagger}

February 9, 2025

Abstract

This paper proposes a framework and solution concept for repeated coalitional behavior. We model history-dependent schemes that deter coalitions from blocking using continuation promises and punishments. We evaluate the effectiveness of these schemes across a range of settings. We apply our results to repeated matching and negotiations.

^{*}We thank Dan Barron, Federico Echenique, Jon Eguia, Matt Elliott, Drew Fudenberg, Ben Golub, Yingni Guo, Scott Kominers, Maciej Kotowski, Elliot Lipnowski, George Mailath, Francesco Nava, Ziwei Wang, Alex Wolitzky, Nathan Yoder, and various audiences for helpful comments.

[†]Department of Economics, Pennsylvania State University. Email: nageeb@psu.edu.

[‡]Department of Economics, Michigan State University. Email: celiu@msu.edu.

1 Introduction

The study of repeated games models history-dependent schemes that enable players to cooperate even if each myopically favors defection. This canonical approach focuses on non-cooperative play in which actions are chosen only by individuals. But, in many contexts, analysts have found it more useful to allow groups of players to act jointly. For instance, matching and network theory model "pairwise stable arrangements" from which no pair of players can profitably deviate. Political economy models emphasize the "Condorcet Winner," a policy preferred by a majority of voters to all others. More broadly, the study of cooperative games looks at the "core," an arrangement that no group of players would find it profitable to block. These notions are all modeled for static interactions, without harnessing the power of promises and punishments.

A natural question is how one can marry these two approaches to group behavior. Our motivation is twofold. First, to develop a tractable and portable framework that speaks to how repeated play could shape behavior in matching, voting, and other coalitional games. Second, to evaluate generally when dynamic incentives deter group defection; conversely, our analysis identifies settings in which the ability of groups to defect cripples dynamic incentives.

We illustrate our framework using the *Roommates Problem*. Ann, Bella, and Carol decide who will room together. The hitch is that only two people can share a room, leaving at least one person out. Each person prefers to have a roommate, and each has a favorite; Table 1 depicts their payoffs.

	Ann	Bella	Carol
Ann	1	3	2
Bella	2	1	3
Carol	3	2	1

TABLE 1. Payoffs of Row Player from matching with Column Player (or remaining unmatched).

As is well known, no arrangement is pairwise stable. For instance, were Ann and Bella paired as roommates, Bella and Carol would each gain if they defected and roomed together instead. Our point of entry is to see how punishment and rewards can "solve" this problem.

Suppose that instead of a one-time decision, the trio made choices monthly. Each accrues the flow payoffs described above, and weights them by the per period discount factor δ . As in repeated games, no player can commit to her future behavior on- or



FIGURE 1. A perfect coalitional equilibrium for the Roommates Problem if $\delta \geq 1/2$.

off-path. What long-run stable matches can be supported through continuation play?

Figure 1 depicts a stable scheme. On the path of play, Ann and Bella room together each month, leaving Carol out. While Bella and Carol might like to defect and share a room this month, the scheme assures that they refrain from doing so if δ exceeds 1/2. For Bella anticipates that after the deviation, starting from next month, Ann and Carol would room together and she would then be left out. Her short-term gain from rooming with Carol would not offset her long-term loss, since $(1 - \delta)3 + \delta(1) \leq 2$. Moreover, the punishment is itself credible because the prescription following every history, including those off-path, is self-enforcing.

We model schemes of this form in general games. We consider the repeated play of an abstract stage game that accommodates many settings: (i) a strategic form game in which players and coalitions choose actions, (ii) a characteristic function game in which coalitions can block, (iii) matching games with and without transferable utility, and (iv) voting games in which different coalitions have varying power to push through policies. In this repeated game, we propose history-dependent schemes where no coalition profits from blocking at any history given how it affects continuation play. We call such schemes *perfect coalitional equilibria* (PCE).

PCEs are inherently recursive. In repeated strategic form games, PCE refines subgame perfect equilibria. Outside of that context, PCE offers a tractable way to model how continuation promises shape coalitional behavior; PCE's recursive nature implies that its payoff set can be obtained via self-generation approaches developed by Abreu, Pearce, and Stacchetti (1990). In fact, we show that it can be even simpler in many settings: all PCE-supportable payoffs can be supported by stationary PCE if the stage game exhibits *default-independent power*. This property holds in the characteristic function games studied in cooperative game theory, matching models without externalities, and in models of voting.

Our main results characterize generally when continuation promises and punishments deter coalitional blocking. We offer conditions under which history dependence thwarts coalitional deviations so that the set of PCE-supportable payoffs is large. Conversely, we also identify settings in which coalitional deviations choke the possibility for cooperation, resulting in "anti-folk" theorems. Underlying our results is a simple principle: a coalition can withstand punishments if and only if its members have highly aligned interests. Then, all members of that coalition enjoy a high minmax, considerably above their individual minmax. However, if there is any wedge in members' incentives, player-specific punishments can splinter coalitions. Then, the set of PCE-supportable payoffs virtually coincides with that of subgame perfect equilibria.

Building on this principle, we explore features that align coalitions' interests. We find that one such feature is the use of *strongly symmetric* schemes, where players behave symmetrically after every history. These schemes feature often in the study of repeated games, such as grim-trigger punishments used to support cooperation in a repeated prisoner's dilemma or to sustain collusion among oligopolists. Although these schemes constitute subgame perfect equilibria, we show that they typically are not PCE. The reason is that once players' incentives are aligned, punishments are no longer credible. By contrast, schemes that feature asymmetric play off-path can credibly deter blocking coalitions and support a larger payoff set.

We also study whether the ability to transfer utility across players necessarily aligns incentives within a coalition. One might have expected the answer to be yes: if a coalition achieves a net gain by blocking, it can distribute those gains among its members to ensure that each benefits. However, we show that if all transfers are publicly observed, a PCE can break coalitions apart by conditioning its continuation play on who pays whom. Conversely, if a coalition can make transfers "secretly"—that is, without the transfers being publicly observed—it can entirely align its incentives. We show that such a coalition effectively functions as a unitary agent and secures a high payoff across all PCE. Therefore, secret side-payments limit how much such a coalition can be punished and, consequently, what a PCE can enforce.

We consider a detailed application of our results to repeated labor-market matching, building on Kelso and Crawford (1982)'s seminal work. Here, we study the kinds of matchings, wages, and allocations of surplus that can be supported through repeated play. It turns out that the set of supportable outcomes hinges crucially on the transparency of past wages. Under wage transparency, a vast range of outcomes can be supported, enabling workers or firms to capture much of the surplus. By contrast, if wage terms are observed only by the employer and employee, there is a complete collapse of intertemporal incentives: the supportable payoff set reduces to the core of the stage game. As to who then benefits from wage transparency depends on economic primitives. Workers benefit if they are plentiful or their marginal returns fall quickly. Absent transparency, competitive forces bid down their wages; by contrast, wage transparency enables them to enforce higher wages through collective-bargaining schemes. If workers are scarce or their marginal returns fall slowly, then it is firms who profit from wage transparency because that enables them to collusively suppress wages. Thus, our application highlights a new dimension to the debate on wage transparency, connecting it to the side of the market that it empowers to collude.

We also study repeated negotiations when some players have veto power. In the stage game, those with veto power effectively become dictators: the core entails that the veto players capture the entire surplus. Against that backdrop, we study when history dependence can promote egalitarianism. Using our stationarity result described above, we characterize the set of supportable payoffs for fixed discount factors and show that it typically includes equal splits. However, egalitarian schemes collapse if veto players can make secret side-payments to other coalition members; if every minimal winning coalition can make secret side-payments, veto players return to being de facto dictators.

We briefly discuss related work. Bernheim and Slavov (2009) propose the notion of a Dynamic Condorcet Winner for an infinitely repeated voting game. Our solution concept coincides with theirs when applied to majority-rule voting. While they describe some properties of their concept, they do not characterize its limits. Consequently, many issues central to our study do not feature in their work. Moreover, our analysis goes beyond voting games, applying to general coalitional settings.

We draw on results in repeated games, particularly Fudenberg and Maskin (1986) and Abreu, Dutta, and Smith (1994); the latter's characterization of equivalent utilities turns out to be the appropriate notion of alignment for settings without transfers. One of our results highlights how secret transfers can align coalitional incentives and thereby undermine punishments. That secret side-payments can cripple dynamic incentives features also in Barron and Guo (2021), who model a relational contracting game between a long-run principal and a sequence of short-run agents.

Our findings on wage transparency contribute to a growing literature, surveyed in Cullen (2024). In a bargaining model with incomplete information, Cullen and Pakzad-Hurson (2023) show that wage transparency disadvantages workers. We offer a complementary perspective, studying the implications of wage transparency for the repeated game. Beyond this application, we view incorporating long-run incentives in Kelso and Crawford (1982)'s workhorse model to be of independent interest. In the static setting, this framework has been enriched in various directions (e.g., Hatfield and Milgrom 2005; Hatfield, Kominers, Nichifor, Ostrovsky, and Westkamp 2013) but relatively little is known about how carrots and sticks affect these matching markets.

Several papers model coalitional deviations in repeated games. Aumann (1959) and Rubinstein (1980) study Strong Nash and Strong Perfect Equilibria of infinitely repeated games; these concepts assume that players cannot commit to long-term plans on the equilibrium path but can do so when deviating. DeMarzo (1992) focuses on finite-horizon games, proposing an inductive solution that corresponds to a Strong Nash equilibrium of the reduced normal-form game.¹ A different strand of the literature models questions of renegotiation—see, for instance, Bernheim and Ray (1989), Farrell and Maskin (1989), Miller and Watson (2013), and Safronov and Strulovici (2018)—in which players can collectively rewire their expectations of future play. Our work studies the complementary question of when coalitions refrain from profitably blocking given their rational expectations of future play.

That expectations of future play can shape coalitional behavior underlies the work on farsighted coalition formation, surveyed in Ray (2007). Konishi and Ray (2003) and Gomes and Jehiel (2005) study settings like ours in which payoffs accrue in real time and coalitions evaluate moves based on discounted continuation values. Their analyses assume history independence, precluding the use of punishments and rewards. Vartiainen (2011) models history-dependence in a different setting without real-time payoffs; his focus is on the existence of deterministic absorbing processes. Dynamic considerations also feature in studies of matching in which players account for future play when deciding with whom to match, e.g., Corbae, Temzelides, and Wright (2003), Damiano and Lam (2005), Kadam and Kotowski (2018a,b), Doval (2022), and Kotowski (2024). Rostek and Yoder (2024) propose a stability notion for static matching in which similar considerations emerge from players thinking strategically about others' choices.

¹An alternative way to model coalitional play is through repeated extensive-form games (Mailath, Nocke, and White 2017); Hatfield, Kominers, and Lowery (2020) and Hatfield, Kominers, Lowery, and Barry (2020) use such an approach to model collusion in brokered and syndicated markets.

Since our initial draft, Bardhi, Guo, and Strulovici (2024) use our approach to model stability in labor markets in which firms learn about workers' types, evaluating how early-career discrimination can result in persistent wage gaps. Liu (2023) and Liu, Wang, and Zhang (2024) also build on the approach here to study repeated matching between long-run firms and short-run workers.

This paper proceeds as follows. Section 2 describes the basic framework. Section 3 identifies structural properties of PCE and characterizes its payoff set. Section 4 studies the game augmented with transfers. Section 5 applies our results to matching and distribution problems. Section 6 concludes. All proofs are in appendices.

2 Model

Players $N := \{1, 2, ..., n\}$ interact repeatedly at t = 0, 1, 2, ... A coalition is a nonempty subset of N, and we denote the set of coalitions by $\mathcal{C} := 2^N \setminus \{\emptyset\}$.

The Stage Game. In each period, the players collectively choose an *alternative a* from A, a compact metrizable space. The alternative a generates payoff vector $v(a) := (v_1(a), \ldots, v_n(a)) \in \mathbb{R}^n$ for players, where the mapping $v : A \to \mathbb{R}^n$ is continuous.

Given an alternative, a coalition of players can choose to block or participate in it. Our specification allows for blocking by either a single coalition or multiple disjoint coalitions, but the former is what matters when evaluating stability. If coalition C alone blocks alternative a, then it can choose any alternative in $E_C(a)$. The correspondence $E_C: A \Rightarrow A$ is C's effectivity correspondence and offers a standard approach to model coalitional power (e.g., Rosenthal 1972; Moulin and Peleg 1982; Chwe 1994). We assume that $E_C(\cdot)$ is continuous, compact-valued, and reflexive (i.e., $a \in E_C(a)$). In our analysis, we also assume that larger coalitions can do more: for each alternative a, $E_{C'}(a) \subseteq E_C(a)$ for $C' \subseteq C$. This assumption is for notational convenience; we detail in footnotes, when necessary, how to adapt notation if this assumption fails.

To see what this formulation captures, let us revisit the Roommates Problem described in the introduction. An alternative is a rooming arrangement, and the set of alternatives, $A = \{ab|c, bc|a, ac|b, a|b|c\}$ is that of all arrangements, where ij|k denotes *i* and *j* rooming together leaving *k* out. For the alternative that puts Ann and Bella together, Bella could block as an individual and choose an arrangement in $E_{\text{Bella}}(ab|c) = \{ab|c, a|b|c\}$. This specification is tantamount to her only choice *as an individual* being whether to accept or reject Ann as a roommate. Carol has even less power— $E_{\{\text{Carol}\}}(ab|c) = \{ab|c\}$ —because she cannot room with someone else without that player's consent. But by teaming up and blocking as a pair, Bella and Carol could choose any alternative in $E_{\{\text{Bella}, \text{Carol}\}}(ab|c) = \{bc|a, ab|c, a|b|c\}$ where the first element denotes the pair rooming together.

The Roommates Problem serves as a specific illustration but the abstract form is considerably more general, capturing strategic form games, characteristic function games, voting games, as well as matching. We formalize how to embed the first three games below, deferring the discussion of matching to Section 5.1.

Example 1. Consider a strategic form game in which player i's action set, A_i , is compact: A_i can be either the set of pure actions or the set of mixtures over finite actions. The set of alternatives is the set of action profiles $A := \times_{i=1}^{n} A_i$. The effectivity correspondence is $E_C(a) := \{a' \in A : a'_j = a_j \text{ for all } j \notin C\}$, modeling the possibility for a blocking coalition to choose action profiles in which players outside the coalition do not change their actions. This formulation extends the standard definition for individual deviations that are used to define Nash equilibria.

Example 2. Consider majority voting, as in Bernheim and Slavov (2009). Let W be the set of coalitions that have at least $\lceil \frac{n}{2} \rceil$ players. The effectivity correspondence specifies that for every a, $E_C(a) = A$ if $C \in W$, and $E_C(a) = \{a\}$ otherwise.

Example 3. Consider characteristic function games (N, U) where for each coalition $C \in C$, the mapping $U(C) \subseteq \mathbb{R}^{|C|}$ specifies a set of feasible payoff vectors for coalition C if it forms. An alternative a is now a tuple (π, u) , where π is a partition of N, and $u \in \mathbb{R}^n$ is a payoff vector satisfying $u_C \in U(C)$ for each coalition $C \in \pi$. The effectivity correspondence $E_C(a)$ specifies the set of alternatives to which coalition C may move, and the payoff function is $v((\pi, u)) = u$.

Outcomes, Histories, and Plans. We develop our notation recursively. A plan specifies a default alternative, say a, at the beginning of period t = 0. This default is chosen if no coalition blocks it, in which case we record the stage-game outcome as (a, \emptyset) . If coalitions $\{C_1, \ldots, C_k\}$ block the default, and their moves result in the alternative a', we record the stage-game outcome as $(a', \{C_1, \ldots, C_k\})$. Based on the outcome at t = 0, a plan specifies a default at t = 1, and the game continues recursively.²

 $^{^{2}}$ We assume that coalitional blocking is directly observable so as to hew closely to repeated games with perfect monitoring, which we view to be the natural starting point. In some settings, the identity

Proceeding abstractly, let \mathcal{P} be the set of all partitions over players, and define $\mathcal{B} := \{B \subseteq 2^N : B \subseteq \pi \text{ for some } \pi \in \mathcal{P}\}$, so that each B in \mathcal{B} is a collection of disjoint coalitions. We denote the set of stage-game outcomes by $\mathcal{O} := A \times \mathcal{B}$; an outcome specifies an alternative a and a collection of disjoint blocking coalitions. At the beginning of period t, the history $h := (a^{\tau}, B^{\tau})_{\tau=0}^{t-1}$ records the stage-game outcomes up to the start of period t. We denote the set of all t-period histories by \mathcal{H}^t for $t \ge 1$. The set of all histories is $\mathcal{H} := \bigcup_{t=0}^{\infty} \mathcal{H}^t$, where $\mathcal{H}^0 = \{\emptyset\}$. A plan $\sigma : \mathcal{H} \to A$ specifies a default alternative following each history.

Payoffs. A path $(a^t)_{t=0,1,2,\ldots}$ is an infinite sequence of alternatives; from that path, player *i* accrues a normalized discounted payoff of $(1 - \delta) \sum_{t=0}^{\infty} \delta^t v_i(a^t)$, in which δ in [0, 1) is a common discount factor. After a history *h*, a plan σ results in the path $(\sigma(h), \sigma(h, \sigma(h), \emptyset), \ldots)$ recursively and $U_i(h|\sigma)$ denotes player *i*'s payoff from that path.

Solution Concept. Before describing our solution concept, we restate the "static core" in the language of this model. In the stage game, coalition C profitably blocks alternative a if there exists $a' \in E_C(a)$ such that $v_i(a') > v_i(a)$ for all $i \in C$. An alternative a is a core alternative if it cannot be profitably blocked by any coalition. A payoff vector \tilde{v} is in the core if there exists a core-alternative a such that $\tilde{v} = v(a)$. For example, in the Roommates Problem, ab|c fails to be a core alternative because Bella and Carol together can profitably block it.

We build on this notion in the repeated game: when players contemplate blocking the alternative $\sigma(h)$ specified by plan σ at history h, they care not only about their instantaneous payoffs but also about how their choices today affect future behavior.

Definition 1. Coalition C profitably blocks plan σ at history h if there exists $a' \in E_C(\sigma(h))$ such that for all $i \in C$,

$$(1-\delta)v_i(a') + \delta U_i(h, (a', \{C\}) \mid \sigma) > U_i(h \mid \sigma).$$

Definition 2. A plan σ is a **perfect coalitional equilibrium** (PCE) if it cannot be profitably blocked by any coalition at any history.

of a blocking coalition is implied by the chosen alternative. But, in other settings (e.g., matching), the chosen alternative alone might not be enough to encode who initiated the block. For instance, in the Roommates Problem, if the alternative ab|c is blocked and we only record that a|b|c is chosen instead, it does not distinguish which of Ann and Bella blocked.

A PCE is a "self-enforcing" plan in that, given the continuation play, no coalition finds it profitable to block.³ In the language of self-generation, the alternative specified at each history is *enforced* by continuation promises that themselves are credible given that the requirement is imposed at every history (including those off-path). Thus, the continuation of a PCE at every history must itself be a PCE. This recursive form implies that the set of PCE supportable payoffs is amenable to dynamic programming à la Abreu, Pearce, and Stacchetti (1990). If $\delta = 0$, PCEs implement only core alternatives of the stage game; furthermore, a plan that specifies a core alternative a^* after every history is necessarily a PCE for every $\delta \geq 0$. Below, we describe PCE-supportable payoffs and identify some structural properties of PCE.

3 What Payoffs Are Supported by PCE?

3.1 Coalitional Minmaxes

In the introduction, we mentioned how coalitions can withstand punishments if they share highly aligned interests, which in turn limits the scope of PCE-supportable payoffs. We illustrate this phenomenon using the common-interest game depicted in Table 2(A). Here, we adopt the specification of coalitional moves stipulated in Example 1: each player can adjust her own action and the pair can choose an action profile.



TABLE 2. Payoffs in (A) are perfectly aligned while those in (B) are slightly misaligned ($\epsilon > 0$).

This game has a unique PCE, which prescribes the efficient action profile (U, L)at every history guaranteeing each player a payoff of 1. To see why, let \underline{w} denote the infimum of the normalized discounted payoffs from all PCEs, and consider an arbitrary PCE in which each player accrues $w \in [0, 1]$. Since the continuation of a PCE at any history must itself be a PCE, if the pair blocks the alternative in the first period and chooses (U, L) instead, each player receives at least $(1 - \delta) + \delta \underline{w}$. The pair profits from the deviation unless $w \ge (1 - \delta) + \delta \underline{w}$. Because this inequality must hold for warbitrarily close to \underline{w} , it then follows that $\underline{w} \ge 1$.

 $^{^{3}}$ An alternative definition of profitable blocking might stipulate that each coalition member gains weakly and at least one does so strictly. This modification would not affect our results.

In this common-interest game, there is a gap between PCE and subgame perfect equilibrium (henceforth SPE) payoffs: as (D, R) is a Nash equilibrium of the stage game, all payoffs in [0, 1] can be supported by SPE of the repeated game. It turns out that the complete alignment of preferences is both sufficient and necessary for this gap. We find that coalitions of perfectly "like-minded" players can guarantee themselves a high coalitional payoff across all PCEs, regardless of discount factors. But misaligning preferences ever so slightly disrupts coalitions sufficiently that individual minmaxes again prove relevant. For instance, our analysis shows that for the stage game in Table 2(B), PCE can support payoffs arbitrarily close to 0 when players are patient.

The right general formulation of alignment here uses Abreu, Dutta, and Smith (1994)'s notion of equivalent utilities: players i and j have equivalent utilities if there exist k > 0 and $c \in \mathbb{R}$ such that $v_j(a) = kv_i(a) + c$ for all $a \in A$; otherwise, their utilities are not equivalent. We partition the set of players according to this criterion; let C(i) be the set of players with whom player i shares equivalent utilities.

This set leads to what we find to be player *i*'s *coalitional minmax*, namely the lowest payoff that she can be pushed down to when coalition C(i) collectively best responds.

$$\underline{v}_i^{\circ} := \min_{a \in A} \max_{a' \in E_{C(i)}(a)} v_i(a').$$
 (Player *i*'s coalitional minmax)

This term is well-defined as A is compact, $v(\cdot)$ is continuous, and E_C is continuous and compact-valued.⁴ Using \mathcal{V} to denote the convex hull of stage-game payoffs, we define $\mathcal{V}_{CR} := \{v \in \mathcal{V} : v_i > \underline{v}_i^\circ \text{ for every } i = 1, \ldots, n\}$ as the set of *strictly coalitionally rational payoffs.*⁵ We distinguish this minmax from a player's individual minmax:

$$\underline{v}_i := \min_{a \in A} \max_{a' \in E_{\{i\}}(a)} v_i(a').$$
 (Player *i*'s individual minmax)

Generally, \underline{v}_i° is higher than \underline{v}_i . The two coincide if player *i* does not share equivalent utilities with any other player, i.e., $C(i) = \{i\}$. More strongly, \mathcal{V}_{CR} coincides with the set of strictly individually rational payoffs if no two players have equivalent utilities. The Roommates Problem lies in this class as do non-cooperative games that satisfy the NEU condition or full-dimensionality.

⁴Assuming that E_C is monotone in C simplifies the expression; absent monotonicity, the coalitional minmax would be $\underline{v}_i^{\circ} := \min_{a \in A} \max_{C \subseteq C(i)} \max_{a' \in E_C(a)} v_i(a')$. This expression coincides with that above if $E_C(a) \subseteq E_{C(i)}(a)$ for every player *i*, coalition $C \subseteq C(i)$, and alternative *a*.

⁵Although our setup does not have public randomization, the convex hull is relevant because these payoffs may be reached through intertemporal averaging (Sorin 1986; Fudenberg and Maskin 1991).

3.2 The Power of Scapegoat Schemes

With these preliminaries in place, we state our first result.

Theorem 1. For every $\delta \geq 0$, every PCE gives each player *i* a payoff of at least \underline{v}_i° . Moreover, for every $v \in \mathcal{V}_{CR}$, there is a $\underline{\delta} < 1$ such that for every $\delta \in (\underline{\delta}, 1)$, there exists a PCE with discounted payoff equal to *v*.

The first part of the result identifies \underline{v}_i° as the appropriate minmax when coalitions can block: no PCE can push a player's payoff below her coalitional minmax. The second part shows that every strictly coalitionally rational payoff vector can be supported if players are sufficiently patient.

A natural comparison for Theorem 1 is to the folk theorem for SPE in repeated games with perfect monitoring. For this comparison, suppose that no two players share equivalent payoffs. Then Theorem 1's implication coincides with that of Fudenberg and Maskin (1986) and Abreu, Dutta, and Smith (1994). We highlight two differences. First, the result applies for both repeated "cooperative" and "non-cooperative" games, including settings such as repeated matching. Second, and more crucially, PCE are robust to both coalitional and individual deviations. Against this backdrop, our result clarifies that if players' preferences are misaligned and players are patient, deterring coalitional deviations is not harder than deterring individual deviations.

Why? In the proof, we design punishments that crack coalitions. Because blocking requires all coalition members to agree, we can deter coalitions by singling out and punishing just one member of each coalition—a "scapegoat"—as though she were the sole deviator, while granting amnesty to the rest. This approach assures that if coalition members' utilities are not equivalent, a PCE can push each player's payoff arbitrarily close to her individual minmax.⁶

This divide-and-conquer scheme fails if all members of the coalition share equivalent utilities. A higher minmax then applies. To see why, consider such a coalition C and suppose towards a contradiction that some PCE σ could push the payoffs of players in C below their coalitional minmax. Observe that members of coalition C could guarantee their coalitional minmax if they could somehow commit to a long-run plan in which they collectively best respond to the alternative specified by the plan after

⁶We observe that a coalition can be cracked even if all coalition members share the same *ordinal* rankings over alternatives; a PCE can nevertheless create player-specific punishments and isolate a coalition member as a scapegoat through its choice of how to sequence alternatives.

every history. Because payoffs are equivalent, all gains and losses for C's members move in sync if the coalition were to proceed with this plan. By induction, it then follows that there exists a history at which coalition C could profitably block, which precludes σ from being a PCE.

Thus, Theorem 1 highlights that a coalition withstands the force of repeated games if and only if its members share *completely* aligned preferences. On the one hand, this conclusion might appear to emphasize a "knife-edge" consideration. Yet, it applies to common-interest games, a commonly studied setting, in which we find that only efficient action profiles are chosen in a PCE.⁷ More importantly, as we show in our study of strongly symmetric equilibria (Theorem 3) and secret side-payments (Theorem 5), this theme of alignment emerges even if players do not have equivalent utilities.

3.3 Structural Properties

3.3.1 When Do Stationary PCEs Suffice?

Given its recursive form, PCE payoffs may be obtained using self-generation approaches (Abreu, Pearce, and Stacchetti 1990). Herein, we highlight a property that distinguishes PCE from SPE: in a rich class of games, all PCE payoffs may be achieved using PCE that are stationary. A plan σ is *stationary* if following every history h, the plan σ specifies the same alternative in each period so long as the plan is not blocked, i.e., $\sigma(h, (\sigma(h), \emptyset)) = \sigma(h)$. Blocking induces a transition whereby a stationary plan would then specify a potentially different alternative but would do so again in every subsequent period.

In principle, stationarity could restrict the set of supportable payoffs. However, it turns out not to have any bite in *convex* games with *default-independent power*. As the latter notion is novel, we turn to it first. For each coalition C and alternative a, let $v_C(a) := \{(v_i(a))_{i \in C}\}$ denote the projection of v(a) to C's payoff space.

Definition 3. A stage game exhibits **default-independent power** if for every coalition C and alternatives a and a', $v_C(E_C(a)\setminus\{a\}) = v_C(E_C(a')\setminus\{a'\})$.

Definition 3 asserts that what a coalition can obtain through blocking an alternative does not depend on that alternative. While this notion might appear stringent, several well-studied coalitional games (including our applications) exhibit this property.

⁷For these games, the coalitional minmax therefore departs from—and is generally higher than— Wen (1994)'s effective minmax for SPE, which would be $\min_{a \in A} \max_{j \in C(i)} \max_{a'_j \in A_j} v_i(a_{-j}, a'_j)$.

For instance, consider any characteristic function game studied in the classical cooperative game theory literature (Example 3). In such a game, if a coalition blocks a default partition π , the set of utilities that it achieves does not depend on the partition. As an application here, one might consider a matching model without externalities; therein, the set of utilities that a group of players can obtain from blocking an assignment does not hinge on the assignment.

As another example, consider models of voting in political economy in which a "winning" coalition can block and choose any policy— $E_C(a) = A$ —and every nonwinning coalition is completely powerless (i.e., $E_C(a) = \{a\}$). Such settings include majoritarian rules (where |C| must have at least (n + 1)/2 players to have power) as well as supermajority voting rules, weighted voting rules in which players have unequal power, and voting procedures in which some players have veto power.

Our result identifies how stationary PCEs suffice in default-independent games if the stage game is *convex*, i.e., $\{\tilde{v} \in \mathbb{R}^n : \exists a \in A \text{ such that } \tilde{v} = v(a)\}$ is a convex set.⁸

Theorem 2. If the stage game is convex and exhibits default-independent power, then for every $\delta \geq 0$, the set of PCE-supportable payoffs coincides with that supported by stationary PCEs.

Theorem 2 offers a conclusion that would be unexpected of subgame perfect equilibria of repeated games; optimal penal codes often involve non-stationary play (Abreu 1988). The proof invokes both convexity and default-independent power: the former enables us to replace a non-stationary path of play with a stationary path and the latter assures that the replacement does not affect any coalition's incentives. Our result generalizes Bernheim and Slavov (2009) who obtain this conclusion for Dynamic Condorcet Winners. By clarifying that default-independent power is the key underlying property, Theorem 2 establishes that this conclusion holds much more broadly.

3.3.2 An Anti-Folk Theorem for Strongly Symmetric PCE

Green and Porter (1984) and Fudenberg, Levine, and Maskin (1994) elucidate how with monitoring imperfections, strongly symmetric equilibria are inefficient but asymmetric play can support near-efficient payoffs. The theory here offers a complementary

⁸We view convexity to be suitable for applications in which players can transfer utility or face a pure distribution problem. Alternatively, the set may be convex if players can access public correlation devices and make a choice to block before the realization of those lotteries. We note that this payoff set is compact because A is compact and v is continuous.

rationale for asymmetric play that applies even with perfect monitoring: a strongly symmetric PCE cannot credibly punish players because it aligns their interests.

To make this point, we study a strategic form stage game (Example 1) that is symmetric: $A_i = A_j$ for all players *i* and *j*, and for each permutation μ of $\{1, \ldots, n\}$, $v_i(a_{\mu(1)}, \ldots, a_{\mu(n)}) = v_{\mu(i)}(a_1, \ldots, a_n)$ for every action profile *a* and player *i*. The set of symmetric action profiles is $A^S := \{a \in A : a_i = a_j \text{ for all } i, j \in N\}$ and $V^S :=$ $\{v(a) : a \in A^S\}$ denotes their associated payoffs. Given that V^S is compact and totally ordered, a maximal element exists denoted \hat{v} .

A plan σ is *strongly symmetric* if it specifies a symmetric action profile, $\sigma(h)$ in A^S , after every history h. Theorem 3 characterizes strongly symmetric PCE.

Theorem 3. A strongly symmetric PCE exists if and only if the stage game has a symmetric core alternative \hat{a} such that $v(\hat{a}) = \hat{v}$; moreover, \hat{v} is then the unique payoff supported by a strongly symmetric PCE.

This result reflects a collapse of intertemporal incentives in that a strongly symmetric PCE exists if and only if the highest symmetric payoff lies in the core and could therefore be supported without carrots and sticks altogether. The condition is highly restrictive, ruling out games like the repeated prisoner's dilemma or collusion in oligopolistic markets. For instance, consider the use of grim-trigger strategies to support mutual cooperation in these settings. Although such strategy profiles constitute subgame perfect equilibria, Theorem 3 implies that they do not qualify as PCEs. The challenge is that players would find it profitable to deviate at any history during the punishment phase, rendering the punishments non-credible. Theorem 1 nevertheless asserts that high cooperation payoffs can be supported by PCEs if players are patient. The construction must resort to asymmetric punishments; for example, a PCE in the prisoner's dilemma would punish a player by having her cooperate while the other defects for a specified number of periods before returning to mutual cooperation.

4 Do Transfers Align Incentives?

4.1 Transferable Utility Framework

We turn to the question of whether transfers align incentives. Because we vary the observability of transfers, we model them separately from alternatives; Section 4.2 considers publicly observed transfers and Section 4.3 models secret side-payments.

We represent transfers by $T := [T_{ij}]_{i,j\in N}$ where $T_{ij} \in [0,\infty)$ is the utility that player *i* transfers to player *j*. A player's *experienced payoff* is the sum of the payoff from the chosen alternative and net transfers: $u_i(a,T) := v_i(a) + \sum_{j\in N} T_{ji} - \sum_{j\in N} T_{ij}$. Let \mathcal{T} be the set of all $n \times n$ transfer matrices in which entries along the main diagonal equal 0 (so that a player cannot transfer utility to herself). We denote transfers paid by members of coalition C by $T_C := [T_{ij}]_{i\in C, j\in N}$; \mathcal{T}_C is the set of $|C| \times n$ transfer matrices.

An outcome of the stage game now includes the chosen alternative, the identity of blocking coalitions (if any), and the chosen transfers. The set of stage-game outcomes is $\overline{\mathcal{O}} := A \times \mathcal{B} \times \mathcal{T}$. Histories and paths are defined analogous to the NTU case with the addition of transfers. We denote the set of histories with transfers by $\overline{\mathcal{H}}$. A plan $\sigma: \overline{\mathcal{H}} \to A \times \mathcal{T}$ specifies an alternative and configuration of transfers, based on history. We use $a(h|\sigma)$ and $T(h|\sigma)$ to denote the default alternative and transfers in $\sigma(h)$. We modify the definition of $U_i(h|\sigma)$ to reflect the influence of transfers.

By blocking the default (a, T), a coalition C can choose a different alternative $a' \in E_C(a)$ and change their transfers to any T'_C . A question we have to tackle is: if a coalition blocks, what transfers do players outside the coalition make? Two distinct answers strike us as reasonable. The first hews to a "simultaneous noncooperative" formulation in which the blocking by coalition C surprises players outside the coalition, who therefore make transfers T_{-C} as was specified by the plan. The second models a "cooperative" approach in which if a coalition blocks, its members can transfer utility among themselves but players outside that coalition do not transfer any utility to them. To accommodate both answers, we formulate the transfers of others abstractly.

Assumption 1. For each coalition C, if C blocks a default transfers matrix T, the transfers made by players outside of C is $\chi^C(T)$ where $\chi^C : T \to T_{N \setminus C}$ satisfies:

- 1. For each bounded set $S \subseteq \mathcal{T}$, the image $\chi^{C}(S) \subseteq \mathcal{T}_{N \setminus C}$ is also bounded.
- 2. If T satisfies $T_{ij} = 0$ for all $i \notin C, j \in C$, then $\chi^C_{ij}(T) = 0$ for all $i \notin C, j \in C$.

Assumption 1 encompasses the two specifications described above: the former corresponds to $\chi^C(T) = T_{-C}$ whereas the latter corresponds to $\chi^C_{ij}(T) = T_{ij}$ for all $i, j \in N \setminus C$, and $\chi^C_{ij}(T) = 0$ for all $i \in N \setminus C$ and $j \in C$.

Thus, if coalition C blocks, chooses actions a' and changes transfers to T'_C , the realized outcome is then $(a', \{C\}, T'_C, \chi^C(T))$. We now define the versions of profitable blocking and PCE appropriate for this setting.

Definition 4. Coalition C profitably blocks plan σ at history h if there exists an alternative $a' \in E_C(a(h|\sigma))$ and transfers $T'_C = [T'_{ij}]_{i \in C, j \in N}$ such that for all $i \in C$,

$$(1-\delta)u_i(a', T'_C, \chi^C(T(h|\sigma))) + \delta U_i(h, a', \{C\}, T'_C, \chi^C(T(h|\sigma)) \mid \sigma) > U_i(h|\sigma).$$

Definition 5. A plan σ is a **perfect coalitional equilibrium** if it cannot be profitably blocked by any coalition at any history.

To rule out Ponzi schemes, we make the following technical assumption.

Assumption 2. We consider plans σ such that continuation values are bounded across histories: $\{U(h|\sigma) : h \in \overline{\mathcal{H}}\}$ is a bounded subset of \mathbb{R}^n .

4.2 Publicly Observed Transfers

Transfers allow blocking coalitions to distribute gains among their members. One might intuit that transfers would then align coalition members' incentives. However, we find that public transfers have the opposite effect, undermining coalitions. Even those coalitions whose payoffs would be aligned absent transfers can now be splintered. Our result below establishes that all payoffs that are feasible and strictly *individually* rational can be supported.

To state this result, we re-define the set of feasible payoffs to account for transfers:

$$\mathcal{U} := \operatorname{co}\left(\left\{u \in \mathbb{R}^n : \exists a \in A \text{ such that } \sum_{i \in N} u_i = \sum_{i \in N} v_i(a)\right\}\right).$$

The set of feasible and strictly individually rational payoffs is

$$\mathcal{U}_{IR} := \{ u \in \mathcal{U} : u_i > \underline{v}_i \text{ for every } i = 1, \dots, n \}$$

Theorem 4. For every $\delta \geq 0$, every PCE gives each player *i* a payoff of at least \underline{v}_i . Moreover, for every $u \in \mathcal{U}_{IR}$, there is a $\underline{\delta} < 1$ such that for every $\delta \in (\underline{\delta}, 1)$, there exists a PCE with discounted payoff equal to *u*.

The presence of transfers implies that if a member of a blocking coalition anticipates punishment, she can be bribed by others to still go along with it. But comparing Theorem 4 to Theorem 1 reveals that rather than aligning coalition members' incentives, public transfers actually undermine any existing preference alignment that was present before the transfers. The key idea is that public transfers make the distribution of utilities within any blocking coalition transparent to all. Therefore, a PCE can tailor the selection of a scapegoat in a blocking coalition to these transfers so as to punish the coalition member who benefited the least. Conceptually, players i and jhave misaligned interests (or non-equivalent utilities) when one has to pay transfers to the other. This misalignment allows one to construct player-specific punishments.⁹

4.3 Secret Transfers

In light of the analysis above, we ask: what if some coalitions can make secret sidepayments when they block? Could their incentives then be aligned? We view this question to be of conceptual and practical import given that, in many contexts, transfers within coalitions are not public. For instance, a firm when poaching another firm's employees might offer a contract whose terms are observed by the worker and firm alone. These contracts are often confidential, a point to which we return in Section 5.1 in our discussion of wage transparency. More broadly, groups of players often seek and find ways to transfer money under the table when defecting from a social arrangement. Our analysis here identifies the benefits that coalitions accrue from making secret transfers even if their blocking decision is observable.

We consider a setting in which some but not all coalitions can make secret transfers; $S \subseteq C$ denote the set of coalitions that can. In our leading application, we consider firms that can offer contracts to workers with private wage terms. A secret side-payment is observed within a coalition but not outside it. Aligned with this idea, we define the outcomes that are publicly observed by all parties.

Definition 6. Given the set $S \subseteq C$ and a stage-game outcome $o = (a, B, T) \in \overline{O}$, the public transfers, denoted T^p , exclude those made within any blocking coalition in S:

$$T_{ij}^{p} = \begin{cases} \# & \text{if } \exists C \in \mathcal{S} \cap B \text{ such that } \{i, j\} \subseteq C, \\ T_{ij} & \text{otherwise.} \end{cases}$$

The public component of o, denoted by o_p , is $o_p := (a, B, T^p)$. For any history $h = (o^{\tau})_{\tau=0}^t$, the public component of h is $h_p = (o_p^{\tau})_{\tau=0}^t$.

⁹Safronov and Strulovici (2018) also highlight how transfers can undermine groups in their study of renegotiation; they show that the ability to punish players for proposals and transfers can cause inefficient norms to persist.

The definition stipulates that if coalition C can make secret transfers, then transfers within it are not recorded in the public history whenever it blocks; instead, those transfers are recorded as #, indicating that they are missing. In this setting with imperfect public monitoring, we consider the analog of a *perfect public equilibrium* (Abreu, Pearce, and Stacchetti 1990; Fudenberg, Levine, and Maskin 1994).

Definition 7. A plan σ is **public** if $\sigma(h) = \sigma(h')$ for all $h, h' \in \overline{\mathcal{H}}$ satisfying $h_p = h'_p$. A **public PCE** is a public plan σ that constitutes a PCE of the repeated game.

We argue below that secret transfers empower a coalition to act as if it were a single party. To this end, imagine that some coalition $C \in \mathcal{S}$ were a unitary actor that maximizes the total utility $\sum_{i \in C} v_i(\cdot)$ with an effectivity function $E_C(\cdot)$. We would then define its minmax as

$$\underline{u}_C := \min_{a \in A} \max_{a' \in E_C(a)} \sum_{i \in C} v_i(a')$$
 (Coalition C's minmax)

Treating each coalition C in S in this way would lead to the set of feasible and strictly S-coalitionally rational payoffs

$$\mathcal{U}_{CR}(\mathcal{S}) := \left\{ u \in \mathcal{U} : \begin{array}{c} u_i > \underline{v}_i \text{ for every } i \in N, \\ \sum_{i \in C} u_i > \underline{u}_C \text{ for every } C \in \mathcal{S} \end{array} \right\}.$$

The above set derives the set of feasible and "individually" rational payoffs in a fictitious game in which the set of players is $N \cup S$. Our result below shows that this set characterizes the limits of public PCE.

Theorem 5. For every $\delta \geq 0$, every public PCE gives each coalition $C \in S$ a total payoff of at least \underline{u}_C and every player i a payoff of at least \underline{v}_i . Moreover, for every $u \in \mathcal{U}_{CR}(S)$, there is a $\underline{\delta} < 1$ such that for every $\delta \in (\underline{\delta}, 1)$, there exists a public PCE with a discounted payoff equal to u.

Theorem 5 identifies the significant gains that coalitions accrue from finding a channel to transfer utility secretly; all those in such a coalition can collectively enjoy a higher minmax while those outside a secret coalition can be pushed towards their individually rational payoffs.¹⁰ One might view these high coalitional minmaxes as

¹⁰Note that Theorem 4 corresponds to the special case of Theorem 5 in which S is empty.

conveying an "anti-folk" flavor. Indeed, in our applications in Section 5, we show that secret transfers can reduce the supportable payoff set to the core of the stage game.

To see why Theorem 5 holds, we first explain why each coalition $C \in S$ is assured its minmax. Consider a plan σ and suppose towards a contradiction that coalition Cfailed to achieve \underline{u}_C . Were coalition C a unitary actor, it could guarantee a total payoff \underline{u}_C from executing a long-run plan where it best responds to the default alternative in each period. An argument similar to the one-shot deviation principle then establishes that the total utility of members of coalition C must increase by blocking at some history h. By apportioning that gain across the members of C through secret sidepayments, it can then be assured that each member simultaneously profits from the block at that history without affecting continuation play.

We turn to why every payoff in $\mathcal{U}_{CR}(\mathcal{S})$ can be attained for patient players. Consider the fictitious game in which the set of players is $N \cup \mathcal{S}$. In this game, we directly construct "player-specific" punishments for each player; one can see that such punishments exist because the payoffs across the players in this fictitious game satisfy the NEU condition. Using these punishments, the payoff of each player in $N \cup \mathcal{S}$ can be pushed arbitrarily close to its minmax.

We contrast this result with Theorem 4. Therein, we could crack coalitions by fine-tuning the selection of the scapegoat to the details of who pays whom. Such an approach fails here because the continuation play cannot condition the punishment on these fine details for the coalitions in S. Not only do these scapegoat schemes unravel but so do any other that pushes the payoff of one of these coalitions below its minmax.

5 Applications

5.1 Labor Market Matching and Wage Transparency

Many practices in labor markets, such as collective wage bargaining and firms' collusive wage-setting, are fundamentally driven by long-run incentives. We incorporate these considerations into the canonical model of Kelso and Crawford (1982) [henceforth KC82]. In the process, we obtain a new perspective on when and how wage transparency benefits workers.

In this stage game, the set of players is $N := \mathcal{F} \cup \mathcal{W}$, where \mathcal{F} is the set of firms and \mathcal{W} is the set of workers. We use f to denote a generic firm, w to denote a generic worker; i and j denote generic players who could be workers or firms. Each firm can hire multiple workers. An assignment ϕ is a mapping on $\mathcal{F} \cup \mathcal{W}$ such that (i) every worker w is assigned to a firm or herself, $\phi(w) \in \mathcal{F} \cup \{w\}$; (ii) every firm f is assigned to a (potentially empty) set of workers, $\phi(f) \subseteq \mathcal{W}$; and (iii) $w \in \phi(f)$ if and only if $\phi(w) = f$. The set of alternatives A comprises all assignments between firms and workers. A matching is an assignment of workers to firms and a specification of transfers made between players. Following KC82, we allow non-zero transfers to occur only between employers and their employees. Therefore, the set of matchings is

$$\mathcal{M} := \{ (\phi, T) \in A \times \mathcal{T} : T_{ij} \neq 0 \text{ only if } i = \phi(j) \text{ or } i \in \phi(j) \}.$$

Each firm f has a revenue function $v_f : 2^{\mathcal{W}} \to \mathbb{R}$, with $v_f(\emptyset)$ normalized to 0; similarly, worker w has a premuneration utility function $v_w : \mathcal{F} \cup \{w\} \to \mathbb{R}$, where the payoff of being unemployed, $v_w(\{w\})$, is normalized to 0. Abusing notation, we use v_i to also denote the utility player i receives from an assignment, so $v_i(\phi) = v_i(\phi(i))$. Given a matching (ϕ, T) , player i's experienced payoff is $u_i(\phi, T) := v_i(\phi) + \sum_{j \neq i} T_{ji} - \sum_{j \neq i} T_{ij}$.

Following KC82, we focus on three kinds of blocking: a worker can reject her match, a firm can fire all its workers, or a firm and a set of workers choose to match even if that departs from the original assignment. While none of our results change if other coalitions of players could also block, we make this assumption to match KC82 and because it embodies the realistic setting in which all contracting is between a firm and a set of workers. Formally, let $\mathcal{E} := \{\{f\} \cup W : f \in \mathcal{F}, W \subseteq \mathcal{W}\}$ denote all *essential* coalitions, i.e., those comprising a single firm and a set of workers. If C is a singleton or essential coalition, then for all $\phi \in A$, $E_C(\phi) = \{\phi, \phi'\}$,¹¹ where

- 1. A worker w can reject her match: if $C = \{w\}$, then $\phi'(w) = w$, and $\phi'(w') = \phi(w')$ for all $w' \in \mathcal{W} \setminus \{w\}$,
- 2. A firm f can fires its workers: If $C = \{f\}$, then $\phi'(f) = \emptyset$, and $\phi'(f') = \phi(f')$ for all $f' \in \mathcal{F} \setminus \{f\}$, and
- 3. A firm f and set of workers W can choose to match: If $C = \{f\} \cup W$, then $\phi'(f) = W$, and $\phi'(f') = \phi(f') \setminus W$ for all $f' \in \mathcal{F} \setminus \{f\}$.

This formulation specifies that if an assignment ϕ is blocked by coalition C, the resulting assignment coincides with ϕ apart from the departure made by coalition C. In other words, all untouched workers remain matched with their assigned partners.

¹¹For all other coalitions, $E_C(\phi) = \{\phi\}$.

Given that transfers happen only between matched players, those outside of C make no transfers to those in C if C blocks. This specification adheres to the "budgetbalance" case described in Section 4.1; the mapping $\{\chi^C\}_{C\in\mathcal{C}}$ denotes the transfers made across players outside of coalition C.

We now state the definitions of profitable blocking and core used in KC82.

Definition 8. A matching (ϕ, T) is **profitably blocked by coalition** C if there exists an alternative assignment $\phi' \in E_C(\phi)$ and transfers $T'_C = [T'_{ij}]_{i \in C, j \in N}$ such that all in C are better off from the matching $(\phi', T'_C, \chi^C(T))$:

$$u_i(\phi', T'_C, \chi^C(T)) > u_i(\phi, T)$$
 for all $i \in C$.

A matching (ϕ, T) is a **core allocation** if it cannot be profitably blocked by any coalition. The **stage-game core**, denoted by \mathcal{K} , are the payoffs of core allocations.

KC82 show that if firms' revenue functions satisfy gross substitutes, the core is nonempty; we assume the same condition and define it formally in this footnote.¹²

Having described the stage game, we now consider the implications of repetition, using the framework and analyses of Section 4. The concept of PCE defined in Definition 5 naturally extends to this setting, where a plan specifies a stage-game matching at every history. The set of feasible payoffs in this repeated game is $\mathcal{U}^{\mathcal{M}} := \operatorname{co}\left(\left\{u \in \mathbb{R}^n : \exists (\phi, T) \in \mathcal{M} \text{ such that } u = u(\phi, T)\right\}\right)$. Player *i*'s individual minmax payoff is $\underline{v}_i = 0$, which is achieved through a matching that ostracizes her. Thus, the set of feasible and individually rational payoffs is $\mathcal{U}_{IR}^{\mathcal{M}} := \left\{u \in \mathcal{U}^{\mathcal{M}} : u_i > 0 \text{ for all } i \in N\right\}$.

Public vs. Private Wages. The set of matchings that can be supported in the repeated game hinges on whether past wage terms are publicly or privately observed. In the former, we find that many outcomes may be supported; in the latter, we see a collapse of intertemporal incentives leading to only payoffs in the core being tenable.

Suppose all transfers are public. Then, an argument identical to Theorem 4's yields that all feasible and individually rational payoffs can be supported for patient players.

 $[\]frac{1^{2} \text{For a vector of wages from firm } f, T_{f} = (T_{fw})_{w \in \mathcal{W}}, \text{ define } Ch_{f}(T_{f}) := \arg \max_{W \subseteq \mathcal{W}} (v_{f}(W) - \sum_{w \in W} T_{fw}). \text{ For every set of workers } W \text{ and pair of wage vectors } T_{f} \text{ and } T'_{f} \text{ such that } T'_{fw} \geq T_{fw} \text{ for all } w \in \mathcal{W}, \text{ define } E(W, T_{f}, T'_{f}) := \{w \in W : T'_{fw} = T_{fw}\}. \text{ Firm } f\text{'s revenue function satisfies } gross substitutes if } \widehat{W} \in Ch_{f}(T_{f}) \text{ implies that there exists } \widehat{W}' \in Ch_{f}(T'_{f}) \text{ such that } E(\widehat{W}, T_{f}, T'_{f}) \subseteq \widehat{W}'.$

Proposition 1. For every $\delta \geq 0$, every PCE gives each player *i* a payoff of at least 0. Moreover, for every $u \in \mathcal{U}_{IR}^{\mathcal{M}}$, there is a $\underline{\delta} < 1$ such that for every $\delta \in (\underline{\delta}, 1)$, there exists a PCE with discounted payoff equal to *u*.

Although we assume gross substitutes to maintain comparability with KC82, Proposition 1 itself does not require gross substitutes or the absence of externalities. Thus, even if the stage-game core is empty, there exist stable schemes when players are sufficiently patient.¹³

Now suppose each firm can hire and offer private wage terms to a group of workers. Formally, the set of secret coalitions, S, includes all essential coalitions, \mathcal{E} . In this setting, we obtain a conclusion sharper than Theorem 5: all payoff vectors outside the core are untenable regardless of the players' patience.

Proposition 2. Suppose the set of essential coalitions \mathcal{E} can make secret transfers. For every $\delta \geq 0$, a public PCE supports a discounted payoff vector if and only if that payoff vector is in \mathcal{K} .

Proposition 2 is an anti-folk theorem that asserts that empowering essential coalitions to make secret transfers cripples a PCE's ability to go beyond the stage-game core. The "if" direction is immediate as the infinite repetition of a core allocation constitutes a public PCE. For the "only if" direction, observe that by Theorem 5, every essential coalition is assured its minmax payoff in a public PCE. In other words, every firm f and group of workers W must achieve a total utility of at least what they would get from matching together, namely, $v_f(W) + \sum_{i \in W} v_w(f)$, which is this coalition's value. In our proof, we show that all payoff vectors that assure that each coalition obtains at least its value lie in the stage-game core.¹⁴

Who Benefits from Wage Transparency? To answer this question, we specialize to a setting in which workers are homogeneous (which KC82 also consider). Suppose that all workers have the same payoff function, $v_w(f) = \lambda(f)$ for each firm f and

 $^{^{13}}$ In this sense, our results offer a microfoundation for this payoff set that complements that of Rostek and Yoder (2024); they develop notions of strategic consistent choices and beliefs for a static game that can render some outcomes stable even if the core is empty.

¹⁴We comment here on a subtle detail of our analysis. In the setup of Section 4.3, the public history does not record wage offers *only* when a coalition blocks. One may be interested in the setting in which wages are never recorded in the public history, both on- and off-path. Proposition 2 would remain true in that setting: the "if" direction holds as a PCE can support core allocations without observing any past wage offers and the "only if" direction holds because making all wage offers private imposes a further restriction on the public PCE.

worker w. Additionally, each firm's revenue depends only on the number of workers it hires: $v_f(W) = \tilde{v}_f(|W|)$. Let $\rho(f, l) := \lambda(f) + \tilde{v}_f(l) - \tilde{v}_f(l-1)$ be the surplus generated from assigning the l^{th} worker to firm f. We continue to assume that firm revenues satisfy gross substitutes, which KC82 show translates into a condition on diminishing marginal returns: $\rho(f, l)$ is then weakly decreasing in l for each f.

In this setting, the assignment ϕ^* that maximizes total social surplus is found by greedily assigning workers to firm slots in order of their contribution to total surplus; henceforth, we refer to ϕ^* as the efficient assignment. Formally, let $L := |\mathcal{W}|$ be the total labor supply and $\eta(\ell)$ be the ℓ^{th} highest value of $\{\rho(f,l) : f \in \mathcal{F}, l \geq 1\}$ for ℓ in $\{1, \ldots, L\}$, which represents the marginal value of assigning the ℓ^{th} worker optimally. To simplify our exposition, we assume that the set $\{\rho(f,l) : f \in \mathcal{F}, l \geq 1\}$ has no ties and excludes 0; assignment ϕ^* then fills "slots" $\{(f,l) : \rho(f,l) \geq \max\{0,\eta(L)\}\}$ leaving all others vacant. Finally, we assume that it would be inefficient for a single firm to hire all workers so that each firm faces some competition.

The set of utilities compatible with assignment ϕ^* in which each player obtains more than her minmax is $\mathcal{U}^+ := \{ \widetilde{u} = u(\phi^*, T) : (\phi^*, T) \in \mathcal{M}, \widetilde{u}_i > 0 \text{ for all } i \in \mathcal{F} \cup \mathcal{W} \}$. Among these, we consider a specific surplus division in which all workers obtain an identical payoff. Let $\eta(L+1)$ be the $(L+1)^{\text{th}}$ highest value of $\rho(f, l)$ assuming that there were an additional worker in the economy. Then we define:

$$\mathcal{U}^* := \{ \widetilde{u} \in \mathcal{U}^+ : \max\{0, \eta(L+1)\} \le \widetilde{u}_w = \widetilde{u}_{w'} \le \max\{0, \eta(L)\} \text{ for all } w, w' \in \mathcal{W} \}.$$

In these surplus divisions, each worker obtains a net utility of approximately the "marginal product" of the last employee in the economy while firms are residual claimants. We show that $\mathcal{K} = \mathcal{U}^*$, which yields the following conclusion.

Proposition 3. If wages are public, for each $u \in U^+$, there exists $\underline{\delta} < 1$ such that for every $\delta \in (\underline{\delta}, 1)$, there exists a PCE with discounted payoff equal to u. By contrast, if wages are private, for every $\delta \geq 0$, the set of payoffs supported by public PCEs is U^* .

Proposition 3 asserts that any surplus division from the efficient assignment ϕ^* in which individual rationality conditions hold can be supported if wages are public. Firms could collude to extract nearly all surplus from workers; alternatively, workers can collectively bargain to retain almost the entire surplus. By contrast, if wages are private, workers accrue the value of the marginal product of the least productive employee, and firms capture the remaining surplus.



FIGURE 2. (A) and (B) show the distribution of surplus under private wages when the marginal productivity falls slowly or quickly. In the latter case, workers have more to gain from wage transparency.

Given Proposition 3, workers favor wage transparency if they are plentiful—i.e., $\eta(L) < 0$ —or their marginal product falls quickly. Without transparency, workers compete intensely for slots and thereby drive their earnings to near 0. By contrast, wage transparency enables them to use collective bargaining to obtain higher wages for them all. In such a scheme, were a firm to try to poach workers in a way that is mutually profitable, a PCE would deter workers from accepting those offers by reverting to the stage-game core from the next period onwards. Thus, workers recognize that the future promise of high wages—and the continued success of their collective bargaining efforts—requires them to reject offers that are tempting today.

By contrast, if workers are scarce or the marginal product of workers falls slowly i.e., $\eta(1) \approx \eta(L)$ —it is firms who favor wage transparency. All PCEs under private wages result in high wages, as firms compete heavily for workers. Wage transparency allows firms to collusively suppress wages, with all of them setting low wages and agreeing not to poach each other's workers. Such an agreement is viable given the continuation play in which poaching today triggers a "salary war" tomorrow.

We depict this prediction in Figure 2: (A) shows the distribution of worker and firm surplus under private wages when the marginal product of labor falls slowly and (B) shows the same when the marginal product falls quickly. As the figure shows, workers are worse off absolutely and relatively in the latter case. Were wages transparent, workers or firms could obtain better terms. In (A), workers have little to gain but much to lose from wage transparency as firms could then suppress wages; by contrast, in (B), it is workers who can use wage transparency to secure a larger share of the pie. Formalizing this comparative statics prediction, consider two markets M_1 and M_2 that are identical in all respects but one: they differ in the productivity of labor as captured in firms' revenues. In market *i*, firm *f*'s revenue function is $\tilde{v}_{f,i}$, and the marginal value of assigning worker ℓ optimally is then $\eta_i(\ell)$. We assume that labor is valuable in each market, in that $\eta_i(L+1) > 0$ for each *i*, and that the two markets accrue the same gain from hiring the first worker, $\eta_1(1) = \eta_2(1)$.

Definition 9. Market M_2 exhibits more steeply decreasing returns to labor than market M_1 if for every ℓ in $\{1, \ldots, L\}$,

$$\eta_2(\ell) - \eta_2(\ell+1) \ge \eta_1(\ell) - \eta_1(\ell+1).$$

If the inequality is strict for some ℓ , then M_2 exhibits **strictly** more steeply decreasing returns to labor.

In each market, given gross substitutes, the marginal product of labor falls with each incremental worker; Definition 9 asserts that this fall is always more pronounced in M_2 . Modulo integer issues, this definition translates into a standard condition on the second derivative of the total product being more negative in M_2 .¹⁵

We turn to the implications for how surplus is divided between workers and firms. Let $\Pi_i := \sum_{\ell=1}^L \eta_i(\ell)$ denote the total surplus from the efficient assignment in market M_i . Let $\Pi_i^{\mathcal{W}} := [L\eta_i(L+1), L\eta_i(L)]$ denote the set of potential workers' total surplus under private wages; recall from Proposition 3 that each worker is paid the same, which is around the marginal product of the least productive worker. Firms capture the gap between total surplus and that taken by workers; $\Pi_i^{\mathcal{F}}$ denotes the set of potential firms' total surplus. We compare these surplus divisions between the two markets; when comparing sets, we use the strong set order denoted \succeq_{SSO} .

Proposition 4. Suppose M_2 exhibits more steeply decreasing returns to labor than M_1 . Then the following hold about the distribution of surplus under private wages:

- (a) The total surplus in market M_1 is higher: $\Pi_1 \ge \Pi_2$.
- (b) Worker surplus in market M_1 must be higher: $\Pi_1^{\mathcal{W}} \succeq_{SSO} \Pi_2^{\mathcal{W}}$.
- (c) Firm surplus in market M_1 must be lower: $\Pi_1^{\mathcal{F}} \preccurlyeq_{SSO} \Pi_2^{\mathcal{F}}$.

¹⁵The "total product" in each market with ℓ units of labor would be $\widehat{\Pi}_i(\ell) := \sum_{l=1}^{\ell} \eta_i(l)$. As the first worker in markets M_1 and M_2 generates the same gain, Definition 9 implies that $\widehat{\Pi}_2(\cdot)$ must be a concave transformation of $\widehat{\Pi}_1(\cdot)$, which is tantamount to the standard Arrow-Pratt comparison.

Furthermore, if M_2 exhibits strictly more steeply decreasing returns to labor, then all the orders above are strict.

Proposition 4 identifies an interesting property: while a more steeply decreasing returns reduces both total and worker surplus under private wages, it has a more pronounced effect on the latter. Hence, as seen in (c), the residual surplus captured by firms is actually higher in M_2 . If wage transparency enables workers to capture firms' profits, then workers have more to gain (and less to potentially lose) in M_2 than M_1 .

5.2 Distributive Politics

Herein, we study a repeated distribution problem, in which the players repeatedly choose how to divide a dollar. Such division problems feature prominently in the political economy literature (e.g. Baron and Ferejohn 1989) and relate to the *simple* games (Von Neumann and Morgenstern 1945) studied in cooperative game theory. The set of alternatives A are divisions of the dollar, $\{a \in \mathbb{R}^N_+ : \sum_{i \in N} a_i = 1\}$, where player i's payoff from alternative a is a_i . Divisions are chosen by a "winning" coalition: \mathcal{W} is a set of coalitions such that for every coalition C in \mathcal{W} , $E_C(a) = A$ for every division a, and for every coalition C not in \mathcal{W} , $E_C(a) = \{a\}$. As standard, \mathcal{W} is monotone and proper.¹⁶ A simple-majority rule protocol corresponds to \mathcal{W} comprising all coalitions that have at least (n + 1)/2 players. This formulation also allows for veto power: if a player belongs to every winning coalition $(\bigcap_{C \in \mathcal{W}} C)$, then effectively no block can happen without her approval. We denote the set of veto players by $D := \bigcap_{C \in \mathcal{W}} C$.

Bernheim and Slavov (2009) approach this setting with simple majority rule in mind, emphasizing how Dynamic Condorcet Winners exist although the stage game lacks a Condorcet Winner. We focus instead on settings with at least one veto and one non-veto player, and in which veto players are not dictators ($D \notin W$). Absent history dependence, these settings are prone to highly unequal splits: the veto players steal the entire dollar, emerging as *de facto* dictators of the game. Formally, the set of core alternatives of the stage game is $\mathcal{K} := \{a \in \mathbb{R}^N_+ : \sum_{i \in D} a_i = 1\}$. The logic is that any division that gives a positive share to a non-veto player would be profitably blocked by a winning coalition who would extract that share and divide it among themselves.

Against this backdrop, we evaluate how history dependence can counter this tendency towards unequal splits. Consider a three-player example in which player 1 alone

¹⁶In other words, if C is in \mathcal{W} , then \mathcal{W} contains every superset of C but not its complement.



FIGURE 3. (A) depicts the set of supportable outcome. The red region depicts payoffs supported by core-reversion, and the blue region illustrates those from other PCE. (B) shows the set of supportable payoffs once coalition $\{1,2\}$ can make secret transfers; player 3 then obtains 0.

has veto power; however, she needs the support of at least one other player to block. In the core of this stage game, player 1 captures the entire dollar. Nevertheless, relatively simple schemes in the repeated game can promote equal splits. Consider a core reversion plan that prescribes $(\frac{1}{3}, \frac{1}{3}, \frac{1}{3})$ every period if that has been the division up to now and switches to the stage-game core otherwise. On the equilibrium path, even if player 1 offers the entire dollar to either player 2 or 3, neither finds it profitable to block with her if $(1 - \delta)(1) + \delta(0) \leq \frac{1}{3}$. Going further, core-reversion can support any division in the triangle formed by the vertices $\{(2\delta - 1, 1 - \delta, 1 - \delta), (0, \delta, 1 - \delta), (0, 1 - \delta, \delta)\}$, which converges to the unit simplex as $\delta \to 1$.

One could go beyond core-reversion to characterize all PCE payoffs. Because the game is convex and exhibits default-independent power, Theorem 2 implies that all PCE payoffs can be supported by stationary PCE. Using this result, we find that if players are sufficiently patient, then every payoff in which each non-veto player obtains up to δ can be supported in a PCE. We depict these outcomes in Figure 3(A).

These schemes collapse if the veto player can make and receive secret side-payments. Suppose players 1 and 2 can transfer utility under the table. Theorem 5 implies that player 3 then obtains 0 in all PCE payoffs, as illustrated in Figure 3(B). Even worse, player 1 takes the entire dollar in every period if she can make secret side-payments with each player.

These intuitions generalize to n-player games in which there are at least one veto and one non-veto player, and veto players are not dictators. We call a coalition C a minimal winning coalition if C is a winning coalition and every proper subset is not. **Proposition 5.** The following hold:

- (a) Absent secret transfers, there exists $\underline{\delta} \ge 0$ such that if $\underline{\delta} \ge \underline{\delta}$, the set of supportable payoffs are those that give at least (1δ) to each winning coalition.
- (b) A winning coalition C obtains the entire dollar in every period in every PCE, regardless of δ, if it can make secret transfers.
- (c) The veto players obtain the entire dollar in every period in every PCE, regardless of δ , if every minimal winning coalition can make secret transfers.

Proposition 5(a) highlights how egalitarian schemes can be supported by history dependence. We use Theorem 2 to obtain this fixed discount factor characterization; it turns out that $\underline{\delta} = 0$ if there are at least two veto players so the characterization then is complete. Proposition 5(b) and (c) elucidate how secret side-payments destabilize egalitarian schemes: the veto players regain *de facto* dictatorial power if every minimal winning coalition can make secret transfers.

6 Conclusion

This paper develops a portable framework for coalitional repeated games, which enables us to evaluate the role of dynamic incentives in coalitional behavior across a broad range of settings. Our analysis uncovers the importance of alignment: history dependence keeps coalitions in line if coalition members' interests are even slightly misaligned. Simple scapegoat schemes then deter coalitional deviations. However, if players in a coalition have completely aligned interests, they can secure a higher minmax payoff by effectively acting as a unitary agent.

This perspective delivers additional insights. Strongly symmetric schemes do not deter coalitional deviations, pushing towards the use of asymmetric punishments. Being able to transfer utility alone does not align interests; to the contrary, publicly observed transfers create a wedge between coalition partners and thereby undermine coalitions. However, the ability to make transfers under the table forges strong ties: a coalition that can do so is assured a high net payoff across PCE.

In our applications, these secret side-payments cripple intertemporal incentives, reducing the set of supportable outcomes to the stage-game core. We use these results to identify conditions under which workers favor wage transparency in repeated labor-market matching. We also study repeated negotiations to evaluate how history dependence can counter the tendency of veto players to become de facto dictators.

Our setup models a purely repeated game in which choices today have no direct bearing on future payoffs. A natural direction for future research would study settings like dynamic public good provision, natural resource depletion, or experimentation in which actions in one period directly impact options (or beliefs) in the next.¹⁷

References

- Abreu, Dilip. 1988. "On the Theory of Infinitely Repeated Games with Discounting." *Econometrica* 56 (2):383–396.
- Abreu, Dilip, Prajit K. Dutta, and Lones Smith. 1994. "The Folk Theorem for Repeated Games: A NEU Condition." *Econometrica* 62 (4):939–948.
- Abreu, Dilip, David Pearce, and Ennio Stacchetti. 1990. "Toward A Theory of Discounted Repeated Games with Imperfect Monitoring." *Econometrica* :1041–1063.
- Aumann, Robert J. 1959. "Acceptable Points in General Cooperative n-Person Games." In Contributions to the Theory of Games IV, vol. 4, edited by H. W. Kuhn and R. D. Luce. Princeton, NJ: Princeton University Press, 287.
- Bardhi, Arjada, Yingni Guo, and Bruno Strulovici. 2024. "Early-Career Discrimination: Spiraling or Self-Correcting?" Working Paper.
- Baron, David P. and John A. Ferejohn. 1989. "Bargaining in Legislatures." American Political Science Review 83 (4):1181–1206.
- Barron, Daniel and Yingni Guo. 2021. "The use and misuse of coordinated punishments." The Quarterly Journal of Economics 136 (1):471–504.
- Bernheim, B. Douglas and Debraj Ray. 1989. "Collective Dynamic Consistency in Repeated Games." Games and Economic Behavior 1 (4):295–326.
- Bernheim, B. Douglas and Sita N. Slavov. 2009. "A Solution Concept for Majority Rule in Dynamic Settings." *Review of Economic Studies* 76 (1):33–62.
- Chwe, Michael. 1994. "Farsighted Coalitional Stability." *Journal of Economic Theory* 63 (2):299–325.

¹⁷In this vein, Bardhi, Guo, and Strulovici (2024) apply our solution concept to a dynamic matching labor market in which firms learn about worker types.

- Corbae, Dean, Ted Temzelides, and Randall Wright. 2003. "Directed Matching and Monetary Exchange." *Econometrica* 71 (3):731–756.
- Cullen, Zoë. 2024. "Is Pay Transparency Good?" Journal of Economic Perspectives 38 (1):153–80.
- Cullen, Zoë B and Bobak Pakzad-Hurson. 2023. "Equilibrium effects of pay transparency." *Econometrica* 91 (3):765–802.
- Damiano, Ettore and Ricky Lam. 2005. "Stability in Dynamic Matching Markets." Games and Economic Behavior 52 (1):34–53.
- DeMarzo, Peter M. 1992. "Coalitions, Leadership, and Social Norms: The Power of Suggestion in Games." Games and Economic Behavior 4 (1):72–100.
- Doval, Laura. 2022. "Dynamically Stable Matching." *Theoretical Economics* 17 (2):687–724.
- Farrell, Joseph and Eric Maskin. 1989. "Renegotiation in Repeated Games." Games and Economic Behavior 1 (4):327–360.
- Fudenberg, Drew, David Levine, and Eric Maskin. 1994. "The Folk Theorem with Imperfect Public Monitoring." *Econometrica* 62 (5):997–1039.
- Fudenberg, Drew and Eric Maskin. 1986. "The Folk Theorem in Repeated Games with Discounting or with Incomplete Information." *Econometrica* :533–554.
- ———. 1991. "On the Dispensability of Public Randomization in Discounted Repeated Games." *Journal of Economic Theory* 53 (2):428—438.
- Gomes, Armando and Philippe Jehiel. 2005. "Dynamic Processes of Social and Economic Interactions: On the Persistence of Inefficiencies." Journal of Political Economy 113 (3):626–667.
- Green, Edward J and Robert H Porter. 1984. "Noncooperative collusion under imperfect price information." *Econometrica* :87–100.
- Hatfield, John William, Scott Duke Kominers, and Richard Lowery. 2020. "Collusion in brokered markets." Working Paper.
- Hatfield, John William, Scott Duke Kominers, Richard Lowery, and Jordan M. Barry. 2020. "Collusion in Markets with Syndication." *Journal of Political Economy* 128 (10):3779–3819.
- Hatfield, John William, Scott Duke Kominers, Alexandru Nichifor, Michael Ostrovsky,

and Alexander Westkamp. 2013. "Stability and competitive equilibrium in trading networks." *Journal of Political Economy* 121 (5):966–1005.

- Hatfield, John William and Paul R Milgrom. 2005. "Matching with Contracts." American Economic Review 95 (4):913–935.
- Kadam, Sangram V. and Maciej H. Kotowski. 2018a. "Multiperiod Matching." International Economic Review 59 (4):1927–1947.
- ———. 2018b. "Time Horizons, Lattice Structures, and Welfare in Multi-Period Matching Markets." *Games and Economic Behavior* 112:1–20.
- Kelso, Alexander S and Vincent P Crawford. 1982. "Job Matching, Coalition Formation, and Gross Substitutes." *Econometrica* :1483–1504.
- Konishi, Hideo and Debraj Ray. 2003. "Coalition Formation as A Dynamic Process." Journal of Economic Theory 110 (1):1–41.
- Kotowski, Maciej H. 2024. "A perfectly robust approach to multiperiod matching problems." *Journal of Economic Theory* 222:105919.
- Liu, Ce. 2023. "Stability in Repeated Matching Markets." *Theoretical Economics* 18 (4):1711–1757.
- Liu, Ce, Ziwei Wang, and Hanzhe Zhang. 2024. "Self-Enforced Job Matching." Working Paper.
- Mailath, George J., Volker Nocke, and Lucy White. 2017. "When and How The Punishment Must Fit The Crime." *International Economic Review* 58 (2):315–330.
- Miller, David A and Joel Watson. 2013. "A theory of disagreement in repeated games with bargaining." *Econometrica* 81 (6):2303–2350.
- Moulin, Herve and Bezalel Peleg. 1982. "Cores of effectivity functions and implementation theory." *Journal of Mathematical Economics* 10 (1):115–145.
- Ray, Debraj. 2007. A Game-Theoretic Perspective on Coalition Formation. New York, NY: Oxford University Press.
- Rosenthal, Robert W. 1972. "Cooperative Games in Effectiveness Form." Journal of Economic Theory 5 (1):88–101.
- Rostek, Marzena J and Nathan Yoder. 2024. "Matching with Strategic Consistency." Working Paper.

- Rubinstein, Ariel. 1980. "Strong Perfect Equilibrium in Supergames." International Journal of Game Theory 9 (1):1–12.
- Safronov, Mikhail and Bruno Strulovici. 2018. "Contestable norms." Working Paper.
- Sorin, Sylvain. 1986. "On Repeated Games with Complete Information." *Mathematics* of Operations Research 11 (1):147–160.
- Vartiainen, Hannu. 2011. "Dynamic Coalitional Equilibrium." Journal of Economic Theory 146 (2):672–698.
- Von Neumann, John and Oskar Morgenstern. 1945. Theory of Games and Economic Behavior. Princeton University Press Princeton, NJ.
- Wen, Quan. 1994. "The "Folk Theorem" for Repeated Games with Complete Information." *Econometrica* :949–954.

A Appendix

The main appendix contains proofs for Theorems 1, 2, 3, and 5. All other proofs are in the Supplementary Appendix. Throughout our analysis, we use sequences of play to convexify payoffs, following standard arguments from Sorin (1986) and Fudenberg and Maskin (1991). Below, we reproduce the statement that we invoke in our arguments.

Lemma 1. (Lemma 2 of Fudenberg and Maskin 1991) Let X be a convex polytope in \mathbb{R}^n with vertices x^1, \ldots, x^K . For all $\epsilon > 0$, there exists a $\underline{\delta} < 1$ such that for all $\underline{\delta} < \delta < 1$, and any $x \in X$, there exits a sequence $\{x_{\tau}\}_{\tau=0}^{\infty}$ drawn from $\{x^1, \ldots, x^K\}$, such that $(1 - \delta) \sum_{\tau=0}^{\infty} \delta^{\tau} x_{\tau} = x$ and at any t, $||x - (1 - \delta) \sum_{\tau=t}^{\infty} \delta^{\tau-t} x_{\tau}|| < \epsilon$.

A.1 Proof of Theorem 1 on p. 11

A Preliminary Result. A blocking plan by coalition C from a plan σ is a function $\alpha : \mathcal{H} \to A$ such that $\alpha(h) \in E_C(\sigma(h))$ for every history $h \in \mathcal{H}$. After each history h, the blocking plan α generates a path $(\alpha(h), \alpha(h, \alpha(h), \{C\}), \ldots)$ that is distinct from the one generated by σ . We use $U_i(h|\alpha)$ to denote player *i*'s normalized discounted payoff from that path. The blocking plan α is profitable if there exists a history h such that $U_i(h|\alpha) > U_i(h|\sigma)$ for all $i \in C$. Below, we say that a coalition is in the alignment partition if it corresponds to a coalition C(i) for some player *i*.

Lemma 2. If σ is a PCE, then no coalition in the alignment partition has a profitable blocking plan.

Proof. Consider a plan σ from which coalition $C(i^*)$ has a profitable blocking plan α for some $i^* \in N$. In particular, there exists a history $h \in \mathcal{H}$ such that $U_i(h|\alpha) > U_i(h|\sigma)$ for every $i \in C(i^*)$. We show that coalition $C(i^*)$ must then have a profitable block from the plan σ at some history, so σ is not a PCE.

Since the set of alternatives A is compact and $v : A \to \mathbb{R}^n$ is continuous, the plan σ has bounded continuation values for all players. Given discounting, the standard one-shot deviation principle applies. Therefore, there exists a history $\hat{h} \in \mathcal{H}$ such that

$$(1-\delta)u_{i^*}(\alpha(\widehat{h})) + \delta U_{i^*}(\widehat{h}, \alpha(\widehat{h}), \{C(i^*)\} | \sigma) > U_{i^*}(\widehat{h} | \sigma).$$

Since players in $C(i^*)$ have equivalent payoffs, for each $j \in C(i^*)$ there exists $\lambda_{ji^*} > 0$ and $\mu_{ji^*} \in \mathbb{R}$ such that $u_j(a) = \lambda_{ji^*} u_{i^*}(a) + \mu_{ji^*}$ and all alternatives $a \in A$; in addition, for every $j \in C(i^*)$, the discounted payoffs satisfy $U_j(h|\sigma) = \lambda_{ji^*}U_{i^*}(h|\sigma) + \mu_{ji^*}$ at every history $h \in \mathcal{H}$. Substituting into the inequality above, we have

$$(1-\delta)u_j(\alpha(\widehat{h})) + \delta U_j(\widehat{h}, \alpha(\widehat{h}), \{C(i^*)\} | \sigma) > U_j(\widehat{h} | \sigma) \text{ for all } j \in C(i^*).$$

Therefore, coalition $C(i^*)$ has a profitable block at history \hat{h} .

Proof of Theorem 1.

<u>Part 1</u>: For every $\delta \geq 0$, every PCE gives each player *i* a payoff of at least \underline{v}_i° . We establish this claim by proving its contrapositive: let σ be a plan, and suppose there exists a player *i*^{*} that satisfies $U_{i^*}(\emptyset|\sigma) < \underline{v}_{i^*}^{\circ}$. We show that σ cannot be a PCE. Given Lemma 2, it suffices to show that coalition $C(i^*)$ has a profitable blocking plan.

Consider the following blocking plan α for coalition $C(i^*)$: at every history h, coalition $C(i^*)$ chooses its myopic best response to the default alternative, $\alpha(h) \in$ $\arg \max_{a' \in E_{C(i^*)}(\sigma(h))} v_{i^*}(a')$. By the definition of $\underline{v}_{i^*}^{\circ}$, $v_{i^*}(\alpha(h)) \geq \underline{v}_{i^*}^{\circ}$ for every history h, so player i^* 's continuation value from period 0 must be higher: $U_{i^*}(\emptyset|\alpha) > U_{i^*}(\emptyset|\sigma)$. Given that all players $j \in C(i^*)$ have equivalent utilities, $U_j(\emptyset|\alpha) > U_j(\emptyset|\sigma)$ for all $j \in C(i^*)$, so α is a profitable blocking plan for coalition $C(i^*)$.

<u>Part 2</u>: For every $v \in \mathcal{V}_{CR}$, there is a $\underline{\delta} < 1$ such that for every $\delta \in (\underline{\delta}, 1)$, there exists a PCE with discounted payoff equal to v.

Fix $v^* \in \mathcal{V}_{CR}$. First, observe that for any pair of players i, j such that $j \notin C(i)$, their payoffs satisfy the non-equivalent utilities (NEU) condition. By Lemma 1 and Lemma 2 of Abreu, Dutta, and Smith (1994), we can find *coalition-specific punishments* for v^* : there exist payoff vectors $\{v^{C(i)}\}_{i=1}^n \subseteq \mathcal{V}_{CR}$ such that $v_i^{C(i)} < v_i^*$ for all $i \in N$, and $v_i^{C(i)} > v_j^{C(i)}$ for all $j \notin C(i)$.

Second, let us define coalitional minmaxing alternatives: for each coalition C(i), let $\underline{a}_{C(i)}^{\circ} \in \arg\min_{a \in A} \max_{a' \in E_{C(i)}(a)} v_j(a')$ for some $j \in C(i)$ —note that the specific choice of $j \in C(i)$ in the definition does not matter given the equivalent payoffs within C(i)—as the alternative that will be used to minmax coalition C(i). Since A is compact, v is continuous, and $E_{C(i)}(\cdot)$ is continuous and compact-valued, by Berge's maximum theorem, $\underline{a}_{C(i)}^{\circ}$ is well-defined for each $i \in N$. By construction, $v_i(a') \leq \underline{v}_i^{\circ}$ for all $i \in N$ and $a' \in E_{C(i)}(\underline{a}_{C(i)}^{\circ})$ and in particular, $v_i(\underline{a}_{C(i)}^{\circ}) \leq \underline{v}_i^{\circ}$.

Given these payoffs and punishments, let $\kappa \in (0, 1)$ be such that for every $\tilde{\kappa} \in [\kappa, 1]$, the following is true for every *i*:

(

$$1 - \widetilde{\kappa})v_i(\underline{a}_{C(i)}^\circ) + \widetilde{\kappa}v_i^{C(i)} > \underline{v}_i^\circ \tag{1}$$

For every
$$i \in N$$
 and $j \notin C(i)$: $(1 - \tilde{\kappa})v_j(\underline{a}_{C(i)}^\circ) + \tilde{\kappa}v_j^{C(i)} > (1 - \tilde{\kappa})\underline{v}_j^\circ + \tilde{\kappa}v_j^{C(j)}$ (2)

Inequality (1) implies that every player $j \in C(i)$ is willing to bear the cost of $v_j(\underline{a}_{C(i)}^\circ)$ with the promise of transitioning into their coalition-specific punishment rather than staying at their coalitional minmax, where the promise is discounted at $\tilde{\kappa}$. Similarly, inequality (2) implies that player j is willing to bear the cost of minmaxing any coalition with whom j does not share equivalent utilities, given the promise of transitioning into coalition C(i)'s specific punishment rather than her own, when the post-minmaxing phase payoffs are discounted at $\tilde{\kappa}$. Each inequality holds at $\tilde{\kappa} = 1$ for all i and $j \notin C(i)$. Since the set of players is finite, there exists a value of $\kappa \in (0, 1)$ such that the inequality holds for all $\tilde{\kappa} \in [\kappa, 1], i \in N$ and $j \notin C(i)$.

Let $L(\delta) := \left\lceil \frac{\log \kappa}{\log \delta} \right\rceil$ where $\lceil \cdot \rceil$ is the ceiling function. Observe that $\delta^{L(\delta)} \in [\delta^{\frac{\log \kappa}{\log \delta}+1}, \delta^{\frac{\log \kappa}{\log \delta}}] = [\delta\kappa, \kappa]$. Therefore, $\lim_{\delta \to 1} \delta^{L(\delta)} = \kappa$.

Since $\{v^{C(i)}\}_{i=1}^n \cup \{v^*\} \subseteq \operatorname{co}\{v(a) : a \in A\} \subseteq \mathbb{R}^n$, by Carathéodory's theorem, there exist $\{\widehat{a}^1, \ldots, \widehat{a}^K\} \subseteq A$ for some integer K, such that $\{v^{C(i)}\}_{i=1}^n \cup \{v^*\} \subseteq \operatorname{co}\{v(\widehat{a}^k) : k = 1, \ldots, K\}$. Define $\mathcal{I} := \{C(i)\}_{i=1}^n$, and $\widehat{\mathcal{I}} := \{C(i)\}_{i=1}^n \cup \{*\}$. Lemma 1 then guarantees that for any $\epsilon > 0$, there exists $\underline{\delta} \in (0, 1)$ such that for all $\delta \in (\underline{\delta}, 1)$, there exist sequences $\{\{a^{S,\tau}\}_{\tau=0}^\infty : S \in \widehat{\mathcal{I}}\}$ such that for each $S \in \widehat{\mathcal{I}}$ and t,

$$\begin{split} (1-\delta)\sum_{\tau=0}^{\infty}\delta^{\tau}v(a^{S,\tau}) &= v^{S} \text{ and } \left| \left| v^{S} - (1-\delta)\sum_{\tau=t}^{\infty}\delta^{\tau}v(a^{S,\tau}) \right| \right| < \epsilon. \text{ We fix an} \\ \epsilon &< (1-\kappa)\min\Big\{ \min_{S\in\widehat{\mathcal{I}}, i\in N\setminus S} (v_{i}^{S} - v_{i}^{C(i)}), \ \min_{i\in N} v_{i}^{C(i)} - \underline{v}_{i}^{\circ} \Big\}, \end{split}$$

and given that ϵ , consider δ exceeding the appropriate $\underline{\delta}$.

We now describe the plan that supports v^* . Consider the automaton $(W, w(*, 0), f, \gamma)$:

- $W := \{w(d,\tau) | d \in \widehat{\mathcal{I}}, \tau \ge 0\} \cup \{\underline{w}(S,\tau) | S \in \mathcal{I}, 0 \le \tau < L(\delta)\}$ is the set of possible states and w(*,0) is the initial state;
- $f: W \to \mathcal{O}$ is the output function, where $f(w(d, \tau)) = (a^{d,\tau}, \emptyset)$ and $f(\underline{w}(S, \tau)) = (\underline{a}_{S}^{\circ}, \emptyset)$.
- $\gamma : W \times \mathcal{O} \to W$ is the transition function. For any collection of blocking coalitions $B \in \mathcal{B}$, let $\widehat{C}(B) = \bigcup_{C \in B} C$ denote their union. For states of the form $w(d, \tau)$, the transition is

$$\gamma(w(d,\tau),(a,B)) = \begin{cases} \underline{w}(C(j^*),0) & \text{if } B \neq \emptyset \text{, where } j^* = \min \widehat{C}(B) \\ w(d,\tau+1) & \text{otherwise} \end{cases}$$

For states of the form $\{\underline{w}(S,\tau)|S \in \mathcal{I}, 0 \leq \tau < L(\delta) - 1\},\$

$$\gamma\big(\underline{w}(S,\tau),(a,B)\big) = \begin{cases} \underline{w}(C(j^*),0) & \text{if } \widehat{C}(B) \nsubseteq S \text{ , where } j^* = \min(\widehat{C}(B) \backslash S) \\ \underline{w}(S,0) & \text{if } \widehat{C}(B) \subseteq S \text{ and } \widehat{C}(B) \neq \emptyset \\ \underline{w}(S,\tau+1) & \text{otherwise} \end{cases}$$

For states of the form $\{\underline{w}(S, L(\delta) - 1) | S \in \mathcal{I}\}$, the transition is

$$\gamma(\underline{w}(S, L(\delta) - 1), (a, B)) = \begin{cases} \underline{w}(C(j^*), 0) & \text{if } \widehat{C}(B) \nsubseteq S, \\ & \text{where } j^* = \min(\widehat{C}(B) \backslash S) \\ \underline{w}(S, 0) & \text{if } \widehat{C}(B) \subseteq S \text{ and } \widehat{C}(B) \neq \emptyset \\ w(S, 0) & \text{otherwise} \end{cases}$$

The plan represented by the above automaton yields payoff profile v^* . By construction,

 $||v^d - V(w(d,\tau))|| < \epsilon$ for all (d,τ) ; in addition, for $\tau = 0, 1, \dots, L(\delta)$,

$$V(\underline{w}(S,\tau)) = (1 - \delta^{L(\delta) - \tau})v(\underline{a}_S^\circ) + \delta^{L(\delta) - \tau}V(w(S,0)).$$

Below, we show that this plan is a PCE by showing that no coalition can profitably block in any state of this automaton.

States of the form $w(d, \tau)$: Set $b > \max_{a \in A, i \in N} v_i(a)$. Consider coalition C blocking and implementing the alternative a. Let $j^* = \min C$. For all τ , without the blocking j^* obtains a payoff greater than $v_{j^*}^d - \epsilon$. By participating in the blocking, j^* obtains a payoff less than $(1 - \delta)b + \delta V_{j^*}(\underline{w}(C(j^*), 0)) = (1 - \delta)b + \delta[(1 - \delta^{L(\delta)})v_{j^*}(\underline{a}_{C(j^*)}^{\circ}) + \delta^{L(\delta)}v_{j^*}^{\circ})]$. The blocking is not profitable if the preceding term is no more than $v_{j^*}^d - \epsilon$. We prove that this is true in two separate cases.

First consider the case where $d \in \widehat{\mathcal{I}} \setminus \{C(j^*)\}$. Observe that

$$\lim_{\delta \to 1} (1 - \delta)b + \delta \left[(1 - \delta^{L(\delta)}) v_{j^*}(\underline{a}^{\circ}_{C(j^*)}) + \delta^{L(\delta)} v_{j^*}^{C(j^*)} \right]$$
$$= \lim_{\delta \to 1} \left[(1 - \delta^{L(\delta)}) v_{j^*}(\underline{a}^{\circ}_{C(j^*)}) + \delta^{L(\delta)} v_{j^*}^{C(j^*)} \right] < v_{j^*}^{C(j^*)},$$

where the inequality follows from $v_{j^*}(\underline{a}_{C(j^*)}^{\circ}) \leq \underline{v}_j^{\circ} < v_{j^*}^{C(j^*)}$. Because ϵ by construction is strictly less than $v_{j^*}^d - v_{j^*}^{C(j^*)}$, it follows that payoff from blocking is less than $v_{j^*}^d - \epsilon$ when δ is sufficiently large.

Now suppose that $d = C(j^*)$. The blocking payoff being less than $v_{j^*}^{C(j^*)} - \epsilon$ can be re-written as $(1 - \delta)(b - v_{j^*}^{C(j^*)}) + \epsilon \leq \delta(1 - \delta^{L(\delta)})(v_{j^*}^{C(j^*)} - v_{j^*}(\underline{a}_{C(j^*)}^{\circ})))$. As $\delta \to 1$, the LHS converges to ϵ . Because $\lim_{\delta \to 1} \delta^{L(\delta)} = \kappa$, the RHS converges to $(1 - \kappa)(v_{j^*}^{C(j^*)} - v_{j^*}(\underline{a}_{C(j^*)}^{\circ})))$. By definition of ϵ , the above inequality holds, and therefore, no coalition can profitably block if δ is sufficiently high.

States of the form $\underline{w}(S,\tau)$: We first consider the case where $C \subseteq S$ and $C \neq \emptyset$. Choose an arbitrary $i \in C$ and we will show that the blocking is not profitable for i. By the definition $\underline{a}_{S}^{\circ}$, coalition C cannot generate a payoff of more than $\underline{v}_{i}^{\circ}$ for player i, so i finds the blocking to be unprofitable if

$$(1 - \delta^{L(\delta) - \tau})v_i(\underline{a}_S^\circ) + \delta^{L(\delta) - \tau}v_i^S \ge (1 - \delta)\underline{v}_i^\circ + \delta(1 - \delta^{L(\delta)})v_i(\underline{a}_S^\circ) + \delta^{L(\delta) + 1}v_i^S.$$
(3)

Because $v_i^S > \underline{v}_i^\circ \ge v_i(\underline{a}_S^\circ)$, it suffices to show that the inequality above holds at $\tau = 0$. Re-arranging terms yields that $(1-\delta)(1-\delta^{L(\delta)})v_i(\underline{a}_S^\circ) + (1-\delta)\delta^{L(\delta)}v_i^S \ge (1-\delta)\underline{v}_i^\circ$, and then dividing by $(1 - \delta)$ yields $(1 - \delta^{L(\delta)})v_i(\underline{a}_S^\circ) + \delta^{L(\delta)}v_i^S \geq \underline{v}_i^\circ$. Let us verify that this inequality holds for sufficiently high δ . Taking $\delta \to 1$ yields (1), which is true. Hence (3) holds for sufficiently high δ .

Next we consider the case where $C \nsubseteq S$. By construction, $j^* \notin S$. Player j^* finds blocking to be unprofitable if

$$(1 - \delta^{L(\delta) - \tau})v_{j^*}(\underline{a}_S^{\circ}) + \delta^{L(\delta) - \tau}v_{j^*}^S \ge (1 - \delta)b + \delta(1 - \delta^{L(\delta)})v_{j^*}(\underline{a}_{C(j^*)}^{\circ}) + \delta^{L(\delta) + 1}v_{j^*}^{C(j^*)}.$$
 (4)

We prove that this inequality is satisfied if δ is sufficiently high. Examining the LHS, observe that for all τ such that $0 \leq \tau \leq L(\delta) - 1$,

$$\lim_{\delta \to 1} \left[(1 - \delta^{L(\delta) - \tau}) v_{j^*}(\underline{a}_S^\circ) + \delta^{L(\delta) - \tau} v_{j^*}^S \right] = \lim_{\delta \to 1} \left[\left(1 - \frac{\kappa}{\delta^\tau} \right) v_{j^*}(\underline{a}_S^\circ) + \frac{\kappa}{\delta^\tau} v_{j^*}^S \right]$$
$$= (1 - \widetilde{\kappa}) v_{j^*}(\underline{a}_S^\circ) + \widetilde{\kappa} v_{j^*}^S$$

for some $\tilde{\kappa} \in [\kappa, 1]$. Examining the RHS of (4), observe that

$$\lim_{\delta \to 1} \left[(1 - \delta)b + \delta(1 - \delta^{L(\delta)})v_{j^*}(\underline{a}_{C(j^*)}^{\circ}) + \delta^{L(\delta)+1}v_{j^*}^{C(j^*)} \right] \\
= \lim_{\delta \to 1} \left[(1 - \delta^{L(\delta)})v_{j^*}(\underline{a}_{C(j^*)}^{\circ}) + \delta^{L(\delta)}v_{j^*}^{C(j^*)} \right] \\
= (1 - \kappa)v_{j^*}(\underline{a}_{C(j^*)}^{\circ}) + \kappa v_{j^*}^{C(j^*)} \leq (1 - \kappa)\underline{v}_{j^*}^{\circ} + \kappa v_{j^*}^{C(j^*)} \leq (1 - \widetilde{\kappa})\underline{v}_{j^*}^{\circ} + \widetilde{\kappa} v_{j^*}^{C(j^*)},$$

where the first equality follows from taking limits, the second from $\lim_{\delta \to 1} \delta^{L(\delta)} = \kappa$, the first weak inequality follows from $v_{j^*}(\underline{a}_{C(j^*)}^\circ) \leq \underline{v}_{j^*}^\circ$, the second weak inequality follows from $\widetilde{\kappa} \geq \kappa$ and $\underline{v}_{j^*}^\circ < v_{j^*}^{C(j^*)}$. Since $\widetilde{\kappa} \in [\kappa, 1]$ and $j^* \notin S$, (2) delivers that $(1 - \widetilde{\kappa})v_{j^*}(\underline{a}_S^\circ) + \widetilde{\kappa}v_{j^*}^S$ is strictly higher than $(1 - \widetilde{\kappa})\underline{v}_{j^*}^\circ + \widetilde{\kappa}v_{j^*}^{C(j^*)}$. This term guarantees that (4) holds for sufficiently high δ .

A.2 Proof of Theorem 2 on p. 13

Since the stage game exhibits default-independent power, the set $v_C(E_C(a)\setminus\{a\})$ is independent of the default alternative a for each coalition C; we can therefore define $D(C) := v_C(E_C(a)\setminus\{a\})$ for some $a \in A$. To study stationary PCEs, we define the analogue of the self-generation map (Abreu, Pearce, and Stacchetti 1990). In a stationary PCE, at any history the prescribed current-period payoff and continuation value are the same. Accordingly, we define the stationary self-generation map as follows. For any set $Y \subseteq \mathcal{V} \subseteq \mathbb{R}^n$, define

$$\Phi_{\delta}(Y) := \left\{ y \in Y : \forall C \in \mathcal{C} \text{ and } z_C \in D(C), \exists y' \in Y \text{ and } i \in C \text{ s.t. } y_i \ge (1-\delta)z_i + \delta y'_i \right\}.$$

The self-generation map identifies the set of discounted payoffs that are supportable when continuation payoffs must lie in the set Y, so for any blocking coalition C contemplating payoff $z_C \in D(C)$ for its members, there is an alternative continuation payoff y' that deters C from doing so.

Lemma 3. If $Y \subseteq \mathcal{V}$ and $Y \subseteq \Phi_{\delta}(Y)$, then every $y \in Y$ can be supported by a stationary PCE.

Proof. Consider any $y \in Y \subseteq \Phi_{\delta}(Y)$. We will construct a stationary PCE σ : $\bigcup_{\tau=0}^{\infty} \mathcal{H}^{\tau} \to A$ such that $U(\emptyset|\sigma) = y$. Since $y \in \mathcal{V}$, there exists $\tilde{a} \in A$ such that $v(\tilde{a}) = y$. Define $\sigma(\emptyset) = \tilde{a}$. We will extend σ 's domain to $\bigcup_{\tau=0}^{\infty} \mathcal{H}^{\tau}$ while making sure that for all $\tau \geq 0, h^{\tau} \in H^{\tau}, \sigma$ satisfies (i) stationarity: $\sigma(h^{\tau}, \sigma(h^{\tau}), \emptyset) = \sigma(h^{\tau})$ and $v(\sigma(h^{\tau})) \in Y$, and (ii) no profitable block: for all $C \in \mathcal{C}$ and $a' \in E_C(\sigma(h^{\tau}))$, there exists $i \in C$ such that $v_i(\sigma(h^{\tau})) \geq (1 - \delta)v_i(a') + \delta v_i(\sigma(h^{\tau}, a', \{C\}))$.

Since $y \in \Phi_{\delta}(Y)$, we know that for all $C \in \mathcal{C}$, and $a \in E_C(\tilde{a})$, there exists $y'[a, C] \in Y$ such that $y_i \geq (1 - \delta)v_i(a) + \delta y'_i[a, C]$ for some $i \in C$. Furthermore since $y'[a, C] \in Y \subseteq \mathcal{V}$, this implies the existence of $a'[a, C] \in A$ such that y'[a, C] = v(a'[a, C]). We extend σ to the domain $\{\emptyset\} \cup \mathcal{H}^1$ as follows: $\sigma(\tilde{a}, B) = \tilde{a}$ if B is either an empty set or comprises more than one coalition; and $\sigma(a, \{C\}) = a'[a, C]$ for all $C \in \mathcal{C}$ and $a \in E_C(\tilde{a})$. Clearly, this satisfies properties (i) and (ii) for $\tau = 0$.

Now we complete the definition of the plan σ through induction on t. Fix t > 1, and assume we've defined the function $\sigma : \bigcup_{\tau=0}^{t-1} \mathcal{H}^{\tau} \to A$ satisfying properties (i) and (ii) for $\tau = 0, \ldots, t - 1$ and all $h^{\tau} \in \mathcal{H}^{\tau}$. Consider any $h^{t-1} \in \mathcal{H}^{t-1}$. Since $v(\sigma(h^{t-1})) \in Y \subseteq \Phi_{\delta}(Y)$, we know that for all $C \in \mathcal{C}$, and $a \in E_C(\sigma(h^{t-1}))$, there exists $y^{h^{t-1}}[a, C] \in Y$ such that $v_i(\sigma(h^{t-1})) \ge (1-\delta)v_i(a) + \delta y_i^{h^{t-1}}[a, C]$ for some $i \in C$. In addition, since $y^{h^{t-1}}[a, C] \in Y \subseteq \mathcal{V}$, this implies the existence of $a^{h^{t-1}}[a, C]$ such that $y^{h^{t-1}}[a, C] = v(a^{h^{t-1}}[a, C])$. Extend σ to the domain \mathcal{H}^t by defining $\sigma(h^{t-1}, a, B)$ as follows: $\sigma(h^{t-1}, \sigma(h^{t-1}), B) = \sigma(h^{t-1})$ if B is either an empty set or comprises more than one coalition; and $\sigma(h^{t-1}, a, \{C\}) = a^{h^{t-1}}[a, C]$ for all $h^{t-1} \in \mathcal{H}^{t-1}$, $C \in \mathcal{C}$ and $a \in E_C(\sigma(h^{t-1}))$. Note that, by construction, the function satisfies properties (i) and (ii) for $\tau = 0, \ldots, t$ and $h^{\tau} \in \mathcal{H}^{\tau}$. This completes the induction step.

By property (i), σ is stationary, which implies that at any history h^t it delivers the

discounted payoff $U(h^t|\sigma) = v(\sigma(h^t))$. In particular, $U(\emptyset|\sigma) = y$, as required. Property (ii) then implies that σ is a PCE.

Proof of Theorem 2. We show that any payoff that can be supported by a PCE can be supported by a stationary PCE (given that the converse holds by definition). Take a PCE σ , and let $\mathcal{U}(\sigma) := \{U(h|\sigma) : h \in \mathcal{H}\}$ denote the set of continuation values associated with σ . Since \mathcal{V} is convex, it follows that $\mathcal{U}(\sigma) \subseteq \mathcal{V}$. Given Lemma 3, it suffices to show that $\mathcal{U}(\sigma) \subseteq \Phi_{\delta}(\mathcal{U}(\sigma))$.

Consider any $y \in \mathcal{U}(\sigma)$ and let h^t be some t-history such that $U(h^t|\sigma) = y$. Since σ is a PCE, we know that for any $C \in \mathcal{C}$ and any $a' \in E_C(\sigma(h^t))$, there exists $y' \in \mathcal{U}(\sigma)$ and $i \in C$ such that $(1 - \delta)v_i(a) + \delta y'_i \leq y_i$. Since the stage game exhibits default-independent power, this is equivalent to the statement that for all $C \in \mathcal{C}$ and $z_C \in D(C)$, there exists $y' \in \mathcal{U}(\sigma)$ and $i \in C$ such that $y_i \geq (1 - \delta)z_i + \delta y'_i$, which implies $\mathcal{U}(\sigma) \subseteq \Phi_{\delta}(\mathcal{U}(\sigma))$. The claim then follows from Lemma 3.

A.3 Proof of Theorem 3 on p. 14

A Preliminary Result. Let \mathcal{U}^S denote the set of discounted payoff profiles from strongly symmetric PCEs. The following lemma is useful for proving Theorem 3.

Lemma 4. If \mathcal{U}^S is nonempty, then \mathcal{U}^S is the singleton set $\{\hat{v}\}$.

Proof. Suppose \mathcal{U}^S is nonempty. Since players accrue identical payoffs from symmetric action profiles, we have $u_i = u_j$ for all players i, j and $u \in \mathcal{U}^S$. Let $\hat{x} := \max_{a \in A^S} v_1(a)$ be the highest feasible symmetric payoff, so $\hat{v} = (\hat{x}, \ldots, \hat{x})$. Let $\underline{x} := \inf\{x : (x, \ldots, x) \in \mathcal{U}^S\}$ denote the lowest symmetric PCE payoff.

Consider a sequence $\{(x^k, \ldots, x^k)\}_{k=1}^{\infty} \subseteq \mathcal{U}^S$ that converges to $(\underline{x}, \ldots, \underline{x})$ and let σ^k be the PCE that supports payoff profile (x^k, \ldots, x^k) . As a PCE, σ^k cannot be profitably blocked by the grand coalition N choosing $\widehat{a} \in \arg \max_{a \in A^S} v_1(a)$, which would generate the stage-game payoff profile $(\widehat{x}, \ldots, \widehat{x})$. So for each k we have $x^k \ge (1 - \delta)\widehat{x} + \delta \underline{x}$. Since $x^k \to \underline{x}$, it follows that $\underline{x} \ge \widehat{x}$. However, by definitions $\underline{x} \le \widehat{x}$, so $\mathcal{U}^S = \{\widehat{v}\}$. \Box

Proof of Theorem 3. For the "only if" direction: Suppose there exists a strongly symmetric PCE σ . By Lemma 4, σ 's continuation values satisfy $U(h|\sigma) = \hat{v}$ for all $h \in \mathcal{H}$. Since \hat{v} is the maximal feasible payoff from symmetric action profiles, and σ must prescribe symmetric action profiles after every history, it follows that the default action profile $\sigma(h)$ must satisfy $u(\sigma(h)) = \hat{v}$ for all $h \in \mathcal{H}$.

Take an arbitrary $\hat{h} \in \mathcal{H}$ and let $\hat{a} := \sigma(\hat{h})$ be the default action profile. As argued above \hat{a} is symmetric and $u(\hat{a}) = \hat{v}$, so it remains to show that \hat{a} is a core alternative. This must be true since otherwise, there exists coalition C and $a' \in E_C(\hat{a})$ such that $v_i(a') > v_i(\hat{a})$ for all $i \in C$, which would imply

$$U_{i}(\widehat{h}|\sigma) = (1-\delta)v_{i}(\widehat{a}) + \delta U_{i}(\widehat{h},\widehat{a},\emptyset) = (1-\delta)v_{i}(\widehat{a}) + \delta \widehat{v}_{i}$$

$$< (1-\delta)v_{i}(a') + \delta \widehat{v}_{i} = (1-\delta)v_{i}(a') + \delta U_{i}(\widehat{h},a',\{C\}|\sigma)$$

for all $i \in C$, which contradicts σ being a PCE.

For the remainder of the theorem, suppose $\hat{v} = v(\hat{a})$ for some symmetric core alternative \hat{a} . A plan that specifies \hat{a} as default at all histories is a strongly symmetric PCE that supports the discounted payoff \hat{v} . Uniqueness follows from Lemma 4. \Box

A.4 Proof of Theorem 5 on p. 18

Preliminary Results. To prove our claim, we first introduce the transferable-utility analogue of the concept of a blocking plan, as defined in Appendix A.1.

A (transferable-utility) blocking plan by coalition C from a plan σ is a pair (α, β) , where $\alpha : \overline{\mathcal{H}} \to A$ and $\beta : \overline{\mathcal{H}} \to \mathcal{T}$ satisfy $\alpha(h) \in E_C(a(h|\sigma))$ and $\beta_{-C}(h) = \chi^C(T(h|\sigma))$ for every history $h \in \overline{\mathcal{H}}$. After each history, the blocking plan (α, β) generates a path

$$(\alpha(h),\beta(h), \alpha(h,\alpha(h),\{C\},\beta(h)),\beta(h,\alpha(h),\{C\},\beta(h)),\dots)$$

that is distinct from the one generated by σ . We will use $U_i(h|\alpha,\beta)$ to denote player *i*'s discounted payoff from that path. The blocking plan (α,β) is profitable if there exists a history *h* such that $U_i(h|\alpha,\beta) > U_i(h|\sigma)$ for all $i \in C$.

Lemma 5. If a plan σ is a public PCE, then no coalition $C \in S \cup N$ has a profitable blocking plan.

Proof. Consider a public plan σ from which coalition $C \in S \cup N$ has a profitable blocking plan (α, β) . In particular, there exists a history $h \in \mathcal{H}$ such that $U_i(h|\alpha, \beta) > U_i(h|\sigma)$ for every $i \in C$. We will show that this implies that coalition C has a profitable block from the plan σ , so σ is not a PCE.

By Assumption 2, the plan σ has bounded continuation values. Moreover, as proven in Lemma 7 in the Supplementary Appendix, it is without loss to assume that the blocking plan (α, β) also has bounded continuation values. Treating coalition C a fictitious player whose payoff is the sum of those of its members, we can see that C faces a decision tree with bounded values. Applying the standard one-shot deviation principle to the fictitious player C yields the existence of $\hat{h} \in \overline{\mathcal{H}}$ such that such that

$$(1-\delta)\sum_{i\in C}u_i\Big(\alpha(\widehat{h}),\ \beta(\widehat{h})\Big)+\delta\sum_{i\in C}U_i\Big(\widehat{h},\ \alpha(\widehat{h}),\ \{C\},\ \beta(\widehat{h})\Big|\sigma\Big)>\sum_{i\in C}U_i(\widehat{h}|\sigma).$$

To show that C can profitably block σ at \hat{h} amounts to showing that this total payoff can be divided so that every individual member can be made better off.

Let T^* be the transfers matrix such that for all $(j,k) \notin C \times C$, $T_{jk}^* = \beta_{jk}(\widehat{h})$; but for $(j,k) \in C \times C$, T_{jk}^* satisfies for every $i \in C$,

$$(1-\delta)u_i\Big(\alpha(\widehat{h}), T^*\Big) + \delta U_i\Big(\widehat{h}, \alpha(\widehat{h}), \{C\}, \beta(\widehat{h})\Big|\sigma\Big) > U_i(\widehat{h}|\sigma).$$
(5)

Consider the two histories $h^1 := (\widehat{h}, \ \alpha(\widehat{h}), \ \{C\}, \ \beta(\widehat{h}))$ and $h^2 := (\widehat{h}, \ \alpha(\widehat{h}|\sigma'), \ \{C\}, \ T^*)$.

By the construction of T^* and the fact that $C \in S \cup N$, h^1 and h^2 share the same public component $h_p^1 = h_p^2$. Since the plan σ is public, it follows that for all $i \in N$, $U_i(\hat{h}, \alpha(\hat{h}), \{C\}, \beta(\hat{h})|\sigma) = U_i(\hat{h}, \alpha(\hat{h}), \{C\}, T^*|\sigma)$. Inequality (5) can therefore be re-written as $(1 - \delta)u_i(\alpha(\hat{h}), T^*) + \delta U_i(\hat{h}, \alpha(\hat{h}), \{C\}, T^*|\sigma) > U_i(\hat{h}|\sigma)$ for every $i \in C$, which implies that σ is not a PCE.

The next result shows that for any payoff profile in $\mathcal{U}_{CR}(\mathcal{S})$, we can construct " $(\mathcal{S} \cup N)$ -specific punishments" for all coalitions in $\mathcal{S} \cup N$.

Lemma 6. For any $u^* \in \mathcal{U}_{CR}(\mathcal{S})$, there exist $(\mathcal{S} \cup N)$ -specific punishments $\{u^C : C \in \mathcal{S} \cup N\} \subseteq \mathcal{U}_{CR}(\mathcal{S})$ such that $\sum_{i \in C} u_i^C < \sum_{i \in C} u_i^*$ for all $C \in \mathcal{S} \cup N$, and $\sum_{i \in C} u_i^C < \sum_{i \in C} u_i^C$ for all $C, C' \in \mathcal{S} \cup N, C' \neq C$.

Proof. For any coalition $C \in \mathcal{S} \cup N$, consider the vector u^C defined by

$$u_i^C = \begin{cases} u_i^* - \frac{\epsilon}{|C|} & i \in C \\ u_i^* + \frac{\epsilon}{|N \setminus C|} & i \notin C \end{cases}$$

Compared to the payoff vector u^* , in u^C every player in C is taxed equally, with a total summing up to ϵ ; by contrast, players outside of C are paid equally, with a total also summing up to ϵ . The ϵ may be set sufficiently small to ensure all u^C 's are in $\mathcal{U}_{CR}(\mathcal{S})$.

We show that these vectors satisfy the required inequalities. By construction, $\sum_{i\in C} u_i^C = \sum_{i\in C} u_i^* - \epsilon < \sum_{i\in C} u_i^*$ for all $C \in \mathcal{S} \cup N$, which verifies the first set of inequalities in the lemma.

Now consider two coalitions $C, C' \in S \cup N$ with $C \neq C'$. Coalition C can be partitioned as $C = (C \setminus C') \cup (C \cap C')$. Compared to $u^*, u^{C'}$ gives everyone outside C'an extra $\frac{\epsilon}{|N \setminus C'|}$, while lowering the payoff of everyone inside C' by $\frac{\epsilon}{|C'|}$, so

$$\sum_{i \in C} u_i^{C'} = \sum_{i \in C \setminus C'} u_i^{C'} + \sum_{i \in C \cap C'} u_i^{C'} = \left[\sum_{i \in C \setminus C'} u_i^* + \frac{|C \setminus C'|}{|N \setminus C'|}\epsilon\right] + \left[\sum_{i \in C \cap C'} u_i^* - \frac{|C \cap C'|}{|C'|}\epsilon\right].$$

Comining terms above yields $\sum_{i \in C} u_i^{C'} = \sum_{i \in C} u_i^* - \left[\frac{|C \cap C'|}{|C'|} - \frac{|C \setminus C'|}{|N \setminus C'|}\right] \epsilon$. Note that since $C \neq C'$, either $C \setminus C' \neq \emptyset$ or $C \cap C' \neq C'$ (or both) must be true; in other words, either $\frac{|C \setminus C'|}{|N \setminus C'|} > 0$ or $\frac{|C \cap C'|}{|C'|} < 1$. In either case, $\sum_{i \in C} u_i^{C'} > \sum_{i \in C} u_i^* - \epsilon = \sum_{i \in C} u_i^C$, which gives us the second set of inequalities in the lemma.

Proof of Theorem 5.

<u>Part 1</u>: For all $\delta \geq 0$, public PCEs give each $C \in S \cup N$ a total payoff of at least \underline{u}_C . We prove a stronger statement: every public PCE σ guarantees that for every coalition $C \in S \cup N$ and every history $h \in \overline{\mathcal{H}}$, $\sum_{i \in C} U_i(h|\sigma) \geq \underline{u}_C$. Towards a contradiction, suppose a public plan σ such that there exists a coalition $C \in S \cup N$ and history \hat{h} such that $\sum_{i \in C} U_i(\hat{h}|\sigma) < \underline{u}_C$. We prove that σ must not be a PCE.

To this end, we construct a profitable blocking plan from σ for coalition C. At every history $h \in \overline{\mathcal{H}}$, let $a(h|\sigma)$ denote the default and $\alpha(h) \in \arg \max_{a \in E_C(a(h|\sigma))} \sum_{i \in C} v_i(a)$ be an alternative in coalition C's "best response." By the definition of \underline{u}_C , it follows that $\sum_{i \in C} v_i(\alpha(h)) \geq \underline{u}_C > \sum_{i \in C} U_i(\hat{h}|\sigma)$, so we can find transfers among players in Csuch that when combined with $\alpha(h)$, these transfers give each $i \in C$ higher payoff than $U_i(\hat{h}|\sigma)$. Formally, at every history $h \in \overline{\mathcal{H}}$, there exist transfers $\widetilde{T}_C(h) := [\widetilde{T}_{ij}(h)]_{i \in C, j \in N}$ such that $\widetilde{T}_{ij}(h) = 0$ for all $j \in N \setminus C$, and $v_i(\alpha(h)) + \sum_{j \in C} \widetilde{T}_{ji}(h) - \sum_{j \in C} \widetilde{T}_{ij}(h) >$ $U_i(\hat{h}|\sigma)$ for all $i \in C$. As a result, for each player $i \in C$, the experienced payoff from the stage-game outcome satisfies

$$u_i\Big(\alpha(h), \ \widetilde{T}_C(h), \chi^C\big(T(h|\sigma)\big)\Big) \ge v_i(\alpha(h)) + \sum_{j \in C} \widetilde{T}_{ji}(h) - \sum_{j \in C} \widetilde{T}_{ij}(h) > U_i(\widehat{h}|\sigma),$$

where the weak inequality follows because $\chi_{ji}^C(T(h|\sigma)) \ge 0$ for all $j \in N$, and $\widetilde{T}_{ij}(h) = 0$ for all $j \in N \setminus C$. Observe that the inequality above can hold at every history,

including \hat{h} and those that follow. These steps prove that the blocking plan (α, β) by coalition C, where $\beta(h) := [\tilde{T}_C(h), \chi^C(T(h|\sigma))]$ for every history $h \in \overline{\mathcal{H}}$, is profitable: $U_i(\hat{h}|\alpha,\beta) > U_i(\hat{h}|\sigma)$ for every $i \in C$. Lemma 5 then implies that σ is not a PCE.

<u>Part 2</u>: For every $u \in \mathcal{U}_{CR}(\mathcal{S})$, there $\underline{\delta} < 1$ such that for all $\delta \in (\underline{\delta}, 1)$, there exists a public PCE supporting u.

For every $C \in \mathcal{S} \cup N$, let $\underline{a}_C \in \arg\min_{a \in A} \max_{a' \in E_C(a)} \sum_{i \in C} v_i(a')$ be an alternative that can be used to minmax C. Note that by construction, $\sum_{i \in C} v_i(\underline{a}_C) \leq \underline{u}_C$.

Fix any payoff vector $u^* \in \mathcal{U}_{CR}(\mathcal{S})$, and let $\{u^C : C \in \mathcal{S} \cup N\}$ be the $(\mathcal{S} \cup N)$ -specific punishments from Lemma 6. Given these punishments, let $\kappa \in (0, 1)$ be such that for every $\tilde{\kappa} \in [\kappa, 1]$, the following is true for all $C \in \mathcal{S} \cup N$ and $C' \in \mathcal{S} \cup N \setminus \{C\}$:

$$(1 - \widetilde{\kappa})\sum_{i \in C} v_i(\underline{a}_C) + \widetilde{\kappa}\sum_{i \in C} u_i^C > \underline{u}_C$$
(6)

$$(1 - \widetilde{\kappa})\sum_{i \in C'} v_i(\underline{a}_C) + \widetilde{\kappa}\sum_{i \in C'} u_i^C > (1 - \widetilde{\kappa})\sum_{i \in C'} v_i(\underline{a}_{C'}) + \widetilde{\kappa}\sum_{i \in C'} u_i^{C'}.$$
(7)

By an argument identical to that in Theorem 1, there exists $\kappa \in (0, 1)$ such that the inequalities above hold for all $\tilde{\kappa} \in [\kappa, 1]$, $C \in \mathcal{S} \cup N$ and $C' \in \mathcal{S} \cup N \setminus \{C\}$. Let $L(\delta) := \left\lceil \frac{\log \kappa}{\log \delta} \right\rceil$. As before, we use the property that $\lim_{\delta \to 1} \delta^{L(\delta)} = \kappa$.

For each alternative $a \in A$ let $\mathcal{U}(a) := \{u \in \mathbb{R}^n : \sum_i u_i = \sum_i v_i(a)\}$ denote the set of payoff profiles that can be generated by playing alternative a and redistributing through transfers. Let $\overline{a} \in \arg \max_{a \in A} \sum_{i \in N} v_i(a)$ and $\underline{a} \in \arg \min_{a \in A} \sum_{i \in N} v_i(a)$ be alternatives that maximize and minimize total payoffs, respectively. Since $\mathcal{U}_{CR}(\mathcal{S}) \subseteq$ \mathcal{U}_{IR} , by Lemma 8 in the Supp. Appendix, there exist payoff vectors $\{\widetilde{u}^1, \ldots, \widetilde{u}^M\} \subseteq$ $\mathcal{U}(\overline{a}) \cup \mathcal{U}(\underline{a})$ such that $\mathcal{U}_{CR}(\mathcal{S}) \subseteq \operatorname{co}(\widetilde{u}^1, \ldots, \widetilde{u}^M)$, where each $\widetilde{u}^m = u(\widetilde{a}^m, \widetilde{T}^m)$ for some $\widetilde{a}^m \in \{\overline{a}, \underline{a}\}$ and \widetilde{T}^m . Let $\widetilde{\mathcal{T}} := \{\widetilde{T}^m\}_{m=1}^M$ be the set comprising these transfer matrices.

By Lemma 1, for any $\epsilon > 0$, there exists $\underline{\delta} \in (0, 1)$ such that for all $\delta \in (\underline{\delta}, 1)$, there exist sequences $\{(a^{d,\tau}, T^{d,\tau})_{\tau=0}^{\infty} : d \in \mathcal{S} \cup N \cup \{*\}\} \subseteq \{\overline{a}, \underline{a}\} \times \widetilde{\mathcal{T}}$ such that for each d and t, $(1 - \delta) \sum_{\tau=0}^{\infty} \delta^{\tau} u(a^{d,\tau}, T^{d,\tau}) = u^d$ and $||u^d - (1 - \delta) \sum_{\tau=t}^{\infty} \delta^{\tau} u(a^{d,\tau}, T^{d,\tau})|| < \epsilon$. We fix an ϵ such that

$$\epsilon < (1-\kappa) \min\Big\{\min_{d \in \mathcal{S} \cup N, d' \in \mathcal{S} \cup N \cup \{*\}, d' \neq d} \Big(\sum_{i \in d} u_i^{d'} - \sum_{i \in d} u_i^d\Big), \min_{d \in \mathcal{S} \cup N} \sum_{i \in d} u_d^d - \underline{v}_d\Big\},$$

and given that ϵ , consider δ exceeding the appropriate $\underline{\delta}$.

We now describe the public plan that we use to support u^* . Let **0** denote the

transfer matrix where all players make no transfers. Consider the plan represented by the automaton $(W, w(0, 0), f, \gamma)$, where

- $W := \{w(d,\tau) | d \in S \cup N \cup \{*\}, \tau \ge 0\} \cup \{\underline{w}(S,\tau) | S \in S \cup N, 0 \le \tau < L(\delta)\}$ is the set of possible states and w(*,0) is the initial state;
- $f : W \to \overline{\mathcal{O}}$ is the output function, where $f(w(d,\tau)) = (a^{d,\tau}, \emptyset, T^{d,\tau})$ and $f(\underline{w}(S,\tau)) = (\underline{a}_S, \emptyset, \mathbf{0});$
- $\gamma : W \times \overline{\mathcal{O}} \to W$ is the transition function. For any collection of blocking coalitions $B \in \mathcal{B}$, define $\widehat{C}(B) = (B \cap \mathcal{S}) \cup (\cup_{C \in B \setminus \mathcal{S}} C)$. Note that $(B \cap \mathcal{S})$ is the set of secret coalitions in B, while $\cup_{C \in B \setminus \mathcal{S}} C$ are the members of the non-secret blocking coalitions, so $\widehat{C}(B)$ is the collection of "players" in B if coalitions in \mathcal{S} are treated as fictitious players. For each $S \in \mathcal{S}$ let $u_S(a,T) = \sum_{i \in S} u_i(a,T)$ denote the total utility accruing to S.

For states of the form $\{\underline{w}(S,\tau)|0 \leq \tau < L(\delta) - 1, S \in N \cup S\}$, the transition is

$$\gamma(\underline{w}(S,\tau),(a,B,T)) = \begin{cases} \underline{w}(S^*,0) & \text{where } S^* \in \arg\min_{\widehat{C}(B) \setminus \{S\}} u_{S'}(a,T), \\ & \text{if } B \neq \emptyset \text{ but either } \{S \notin \widehat{C}(B)\} \\ & \text{or } \{u_S(a,T) > \underline{u}_S\} \text{ is true.} \end{cases}$$
$$\underbrace{\underline{w}(S,0) & \text{if } B \neq \emptyset \text{ and both } \{S \in \widehat{C}(B)\} \\ & \text{and } \{u_S(a,T) \leq \underline{u}_S\} \text{ are true.} \\ & \underline{w}(S,\tau+1) & \text{if } B = \emptyset. \end{cases}$$

For states of the form $\{\underline{w}(S, L(\delta) - 1) | S \in N \cup S\}$, the transition is

$$\gamma(\underline{w}(S, L(\delta) - 1), (a, B, T)) = \begin{cases} \underline{w}(S^*, 0) & \text{where } S^* \in \arg\min_{\widehat{C}(B) \setminus \{S\}} u_{S'}(a, T), \\ & \text{if } B \neq \emptyset \text{ but either } \{S \notin \widehat{C}(B)\} \\ & \text{or } \{u_S(a, T) > \underline{u}_S\} \text{ is true.} \end{cases}$$
$$\frac{\underline{w}(S, 0) & \text{if } B \neq \emptyset \text{ and both } \{S \in \widehat{C}(B)\} \\ & \text{and } \{u_S(a, T) \leq \underline{u}_S\} \text{ are true.} \\ w(S, 0) & \text{if } B = \emptyset. \end{cases}$$

For states of the form $w(d, \tau)$, the transition is

$$\gamma \big(w(d,\tau), (a,B,T) \big) = \begin{cases} \underline{w}(S^*,0) & \text{if } B \neq \emptyset, \\ & \text{where } S^* \in \arg\min_{S' \in \widehat{C}(B)} u_{S'}(a,T). \\ & w(d,\tau+1) & \text{if } B = \emptyset. \end{cases}$$

The plan represented by this automaton yields payoff profile u^* . The plan is also public since the transition relies only on B and $\{u_S(a,T): S \in \widehat{C}(B)\}$, both of which are public information. By construction, $||u^d - V(w(d,\tau))|| < \epsilon$ and $V(\underline{w}(S,\tau)) =$ $(1 - \delta^{L(\delta)-\tau})v(\underline{a}_S) + \delta^{L(\delta)-\tau}V(w(S,0))$ for all τ in $\{0, \ldots, L(\delta-1)\}$ and $S \in S \cup N$. As the arguments from here on are standard, we verify in the Supplementary Appendix that no coalition can profitably block in any state of this automaton. \Box

B Supplementary Appendix

The Supplementary Appendix completes the proof of Theorem 5 and contains the proofs of all propositions.

B.1 Preliminary Results

Below, we prove a few preliminary results used in the proof of Theorem 5. The first result shows that when checking profitable blocking plans, we can WLOG focus on those with bounded total continuation values.

Lemma 7. Let σ be a PCE. Suppose coalition C has blocking plan (α, β) such that $\sum_{i \in C} U_i(\overline{h}|\alpha, \beta) > \sum_{i \in C} U_i(\overline{h}|\sigma)$ for some $\overline{h} \in \overline{\mathcal{H}}$, then C has blocking plan (α', β') such that $\sum_{i \in C} U_i(\overline{h}|\alpha', \beta') > \sum_{i \in C} U_i(\overline{h}|\sigma)$, and $\{\sum_{i \in C} U_i(h|\alpha', \beta') : h \in \overline{\mathcal{H}}\}$ is bounded.

Proof. We break this argument into two parts.

<u>Part 1:</u> We show that the set $\{\sum_{i\in C} U_i(h|\alpha,\beta) : h \in \overline{\mathcal{H}}\}$ is bounded from above. To this end, it suffices to show that the set of stage-game payoffs from the blocking plan, $\{\sum_{i\in C} u_i(\alpha(h),\beta(h)) : h\in\overline{\mathcal{H}}\}$ is bounded from above.

Consider an arbitrary coalition $C \in \mathcal{C}$ and an arbitrary history $h \in \overline{\mathcal{H}}$. Let $\tilde{a} = a(h|\sigma)$ denote the default alternative specified by σ and $\tilde{T} = T(h|\sigma)$ denote the default transfers. By the definition of a blocking plan, $\alpha(h) \in E_C(\tilde{a})$ and $\beta(h) = (T'_C, \chi^C(\tilde{T}))$ for some $T'_C \in \mathcal{T}_C$. Since the transfers T'_C may involve nonzero transfers to players outside of C, we have

$$\sum_{i \in C} u_i(\alpha(h), \beta(h)) = \sum_{i \in C} v_i(\widetilde{a}) + \sum_{i \in C, j \notin C} \chi_{ji}^C(\widetilde{T}) - \sum_{i \in C, j \notin C} T'_{ij}$$
$$\leq \sum_{i \in C} v_i(\widetilde{a}) + \sum_{i \in C, j \notin C} \chi_{ji}^C(\widetilde{T})$$
(8)

Now suppose the coalition C blocks at history h and chooses alternative $\alpha(h) \in E_C(\tilde{a})$; however, instead of $\beta(h) = (T'_C, \chi^C(\tilde{T})), C$ chooses transfers $(T''_C, \chi^C(\tilde{T}))$, where the transfers T''_C are such that members of C make zero payment to players outside of C while splitting the total payoff within C evenly. If C carries out this block, each member $i \in C$ obtains a discounted utility of at least

$$(1-\delta)\frac{1}{|C|} \Big[\sum_{i\in C} v_i(\alpha(h)) + \sum_{i\in C, j\notin C} \chi_{ji}^C(\widetilde{T})\Big] + \delta \inf_{h\in\overline{\mathcal{H}}, i\in N} U_i(h|\sigma),$$

whereas adhering to σ at h yields each member at most $\sup_{h\in\overline{\mathcal{H}},i\in N} U_i(h|\sigma)$. Since σ is a PCE, $(\alpha(h), T_C'', \chi^C(\widetilde{T}))$ cannot be a profitable block for C, so it must be true that

$$(1-\delta)\frac{1}{|C|}\Big[\sum_{i\in C}v_i(\widetilde{a}) + \sum_{i\in C, j\notin C}\chi_{ji}^C(\widetilde{T})\Big] + \delta \inf_{h\in\overline{\mathcal{H}}, i\in N}U_i(h|\sigma) \le \sup_{h\in\overline{\mathcal{H}}, i\in N}U_i(h|\sigma).$$

Combining the inequality above with (8) yields

$$(1-\delta)\frac{1}{|C|} \Big[\sum_{i\in C} u_i(\alpha(h),\beta(h))\Big] + \delta \inf_{h\in\overline{\mathcal{H}},i\in N} U_i(h|\sigma) \le \sup_{h\in\overline{\mathcal{H}},i\in N} U_i(h|\sigma).$$

Rearranging terms, we have

$$\begin{split} \sum_{i \in C} u_i(\alpha(h), \beta(h)) &\leq \frac{|C|}{1 - \delta} \left[\sup_{h \in \overline{\mathcal{H}}, i \in N} U_i(h|\sigma) - \delta \inf_{h \in \overline{\mathcal{H}}, i \in N} U_i(h|\sigma) \right] \\ &\leq \frac{|C|}{1 - \delta} \left| \sup_{h \in \overline{\mathcal{H}}, i \in N} U_i(h|\sigma) \right| + \frac{|C|\delta}{1 - \delta} \left| \inf_{h \in \overline{\mathcal{H}}, i \in N} U_i(h|\sigma) \right|. \end{split}$$

Since $\{U(h|\sigma) : h \in \overline{\mathcal{H}}\}$ is bounded by Assumption 2, there exists L > 0 such that

$$\left|\sup_{h\in\overline{\mathcal{H}},i\in N}U_i(h|\sigma)\right|\leq L \quad \text{and} \quad \left|\inf_{h\in\overline{\mathcal{H}},i\in N}U_i(h|\sigma)\right|\leq L.$$

Therefore,

$$\sum_{i \in C} u_i(\alpha(h), \beta(h)) \le \frac{1+\delta}{1-\delta} |C| L.$$

Note that the inequality above holds for all $h \in \overline{\mathcal{H}}$ while the right hand side does not depend on h, so our claim follows.

<u>Part 2:</u> It is without loss to assume that $\{\sum_{i\in C} U_i(h|\alpha,\beta) : h\in\overline{\mathcal{H}}\}\$ is bounded from below. If not, we can construct another blocking plan (α',β') such that $\sum_{i\in C} U_i(\overline{h}|\alpha',\beta') > \sum_{i\in C} U_i(\overline{h}|\sigma)\$ while ensuring $\{\sum_{i\in C} U_i(h|\alpha',\beta') : h\in\overline{\mathcal{H}}\}\$ is bounded from below: if $\sum_{i\in C} U_i(\widehat{h}|\alpha,\beta)\$ falls below $\min_{a\in A} \sum_{i\in C} v_i(a)\$ for some history $\widehat{h}\in\overline{\mathcal{H}}$, we will ask C to refuse all outgoing transfers at all histories following \widehat{h} .

Formally, for a history $\hat{h} \in \overline{\mathcal{H}}$, let $F(\hat{h}) := \{h\hat{h} : h \in \overline{\mathcal{H}}\}$ denote the set of histories that can follow from \hat{h} . Let $\underline{H}_C := \{h \in \overline{\mathcal{H}} : \sum_{i \in C} U_i(h|\alpha, \beta) < \min_{a \in A} \sum_{i \in C} v_i(a)\}$. Let $\mathbf{0}_C$ denote the vector of zero-valued transfers made from players in C. Set $\alpha' = \alpha$, and define

$$\beta'(h) = \begin{cases} \left(\mathbf{0}_C, \chi^C(T(h|\sigma))\right) & \forall h \in F(\widehat{h}) \text{ for some } \widehat{h} \in \underline{H}_C, \\ \beta(h) & \text{otherwise.} \end{cases}$$

By construction, the blocking plan (α', β') has continuation values bounded below by $\min_{a \in A} \sum_{i \in C} v_i(a)$. In addition, compared to (α, β) , the blocking plan (α', β') gives coalition C weakly higher total continuation value following any history, so $\sum_{i \in C} U_i(\overline{h} | \alpha', \beta') > \sum_{i \in C} U_i(\overline{h} | \sigma)$.

Next we argue that there exists a finite set of payoff vectors whose convex hull contains the set \mathcal{U}_{IR} .

Lemma 8. For each alternative $a \in A$ let $\mathcal{U}(a) := \{u \in \mathbb{R}^n : \sum_i u_i = \sum_i v_i(a)\}$ denote the set of payoff profiles that can be generated by playing alternative a and redistributing through transfers. Let $\overline{a} \in \arg \max_{a \in A} \sum_{i \in N} v_i(a)$ and $\underline{a} \in \arg \min_{a \in A} \sum_{i \in N} v_i(a)$ be two alternatives that maximize and minimize players' total generated payoffs, respectively. There exist payoff vectors $\{\widetilde{u}^1, \ldots, \widetilde{u}^M\} \subseteq \mathcal{U}(\overline{a}) \cup \mathcal{U}(\underline{a})$, such that $\mathcal{U}_{IR} \subseteq co(\widetilde{u}^1, \ldots, \widetilde{u}^M)$.

Proof. By definition,

$$\mathcal{U}_{IR} \subseteq \overline{\mathcal{U}}_{IR} := \left\{ u \in \mathbb{R}^n : \sum_{i \in N} v_i(\underline{a}) \le \sum_{i \in N} u_i \le \sum_{i \in N} v_i(\overline{a}) \text{ and } u_i \ge \underline{v}_i \forall i \in N \right\}.$$

Since $\overline{\mathcal{U}}_{IR}$ is a bounded polyhedron, it is also a polytope. Let x^1, \ldots, x^K be its vertices. Any point inside \mathcal{U}_{IR} can then be expressed as convex combinations of these vertices. Since $x^k \in \operatorname{co}(\mathcal{U}(\overline{a}) \cup \mathcal{U}(\underline{a}))$ for all $1 \leq k \leq K$, for each k, there exist $\{\widetilde{u}^{k,1}, \ldots, \widetilde{u}^{k,m_k}\} \subseteq \mathcal{U}(\overline{a}) \cup \mathcal{U}(\underline{a})$ such that $x^k \subseteq \operatorname{co}(\widetilde{u}^{k,1}, \ldots, \widetilde{u}^{k,m_k})$. As a result $\mathcal{U}_{IR} \subseteq \operatorname{co}(\bigcup_{1 \leq k \leq K} \{\widetilde{u}^{k,1}, \ldots, \widetilde{u}^{k,m_k}\})$.

B.2 Completing the Proof of Theorem 5

Recall that $\widetilde{\mathcal{T}} = {\{\widetilde{T}^m\}_{m=1}^M}$ and ${\{(a^{d,\tau}, T^{d,\tau})_{\tau=0}^\infty : d \in \mathcal{S} \cup N \cup \{*\}\}} \subseteq {\{\overline{a}, \underline{a}\} \times \widetilde{\mathcal{T}}}$, so all default transfers from the plan are selected from a finite set. By Assumption 1, when a coalition $C \in \mathcal{C}$ blocks a default transfers matrix $T \in \widetilde{\mathcal{T}}$, there exists $\widetilde{b} > 0$ such that

$$\sum_{i \notin C, j \in C} \chi_{ij}^C(T) \le \widetilde{b} \text{ for all } T \in \widetilde{\mathcal{T}} \text{ and } C \in \mathcal{C}.$$
(9)

In addition, since A is compact and v(.) is continuous, there exists \hat{b} such that

$$\max_{C \in \mathcal{C}} \max_{a \in A} \sum_{i \in C} v_i(a) \le \widehat{b}.$$
(10)

We verify the incentives in the automaton states below.

States of the form $w(d, \tau)$: Suppose coalition C blocks and the outcome $(\hat{a}, \{C\}, \hat{T})$ is realized. Note that if $C \in S \cup N$, then $\widehat{C}(\{C\}) = \{C\}$ is a singleton set containing Cas a unitary player. However, if $C \notin S \cup N$, then $\widehat{C}(\{C\}) = C$, which is a non-singleton set consisting of players in C. The plan punishes $S^* \in \arg\min_{S' \in \widehat{C}(\{C\})} u_{S'}(a, T)$, so S^* is either C as a unitary player or some player $i \in C$. In either case, the (total) stage-game payoff for S^* satisfies

$$u_{S^*}(\widehat{a},\widehat{T}) \leq \frac{1}{|\widehat{C}(\{C\})|} \sum_{S' \in \widehat{C}(\{C\})} u_{S'}(\widehat{a},\widehat{T}) \leq \max_{a \in A} \sum_{j \in C} v_j(a) + \sum_{j \in C} \sum_{k \notin C} \chi^C_{kj}(T^{d,\tau}) \leq \widehat{b} + \widetilde{b},$$

where the first inequality follows since the minimum among a set of numbers is less than their average; the second inequality follows since C's total payoff comes from the generated payoffs plus the net transfers paid by players outside of C; lastly, the third inequality follows from (9) and (10) and the fact that all $T^{d,\tau}$ are drawn from $\{\tilde{T}^m\}_{m=1}^M$.

Thus, we can find a uniform bound $\overline{b}_1 := \widehat{b} + \widetilde{b}$ such that the total stage-game payoff of S', satisfies $u_{S'}(\widehat{a}, \widehat{T}) \leq \overline{b}$ for all C, δ and (d, τ) . Following the same steps as those in the analogous part of Theorem 1, we can show that (when viewed as a unitary player if S^* is not a single player) S^* obtains lower total payoff after coalition C blocks. If S^* is a player in C, then this is not profitable block for C; if $S^* = C$, then there exists player $i \in S^* = C$ who is not better off, so again this is not a profitable block for C.

States of the form $\underline{w}(S,\tau)$ where $S \in N \cup S$: Suppose coalition C blocks and the outcome is $(\hat{a}, \{C\}, \hat{T})$. Just like above, depending on whether $C \in S \cup N$, $\hat{C}(\{C\})$ is either $\{C\}$ containing C as a unitary player or the set C containing all its members. There are 2 cases to consider.

Case I: $S \in \widehat{C}(\{C\})$, and $u_S(\widehat{a}, \widehat{T}) \leq \underline{u}_S$. In this case, the plan punishes the current scapegoat S, where S is either C or a member of C. Using (6) for sufficiently high δ and following steps identical to the analogous argument in Theorem 1, we can show

that

$$(1 - \delta^{L(\delta) - \tau})v_S(\underline{a}_S) + \delta^{L(\delta) - \tau}u_S^S \ge (1 - \delta)\underline{u}_S + \delta(1 - \delta^{L(\delta)})v_S(\underline{a}_S) + \delta^{L(\delta) + 1}u_S^S,$$

where $v_S(.) = \sum_{i \in S} v_i(.)$ and $u_S^S = \sum_{i \in S} u_i^S$ denote the sum of the members' payoffs in case S is a coalition. If S is a member of C, then the inequality above shows that the blocking is not profitable for C; if S is C itself, it implies that blocking does not improve C's total value, so there exist $i \in C$ who is not better off, so again the blocking is not profitable for C.

 $\underbrace{ \text{Case II: either } S \notin \widehat{C}(\{C\}) \text{ or } u_S(\widehat{a}, \widehat{T}) > \underline{u}_S}_{\text{arg min}_{S' \in \widehat{C}(\{C\}) \setminus \{S\}} u_{S'}(\widehat{a}, \widehat{T}) \text{ as scapegoat.} }$ In this case the plan punishes $S^* \in \underbrace{ \operatorname{Res}_{S' \in \widehat{C}(\{C\}) \setminus \{S\}} u_{S'}(\widehat{a}, \widehat{T}) \text{ as scapegoat.} }$

First observe that if C blocks in state $\underline{w}(S, \tau)$ and the stage-game payoff satisfies $u_S(\hat{a}, \hat{T}) > \underline{u}_S$, then no matter if $S \in N$ or $S \in \mathcal{S}$, it must be that $C \neq S$ and therefore $\widehat{C}(\{C\}) \neq \{S\}$; otherwise the definition of \underline{u}_S would ensure $u_S(\hat{a}, \hat{T}) \leq \underline{u}_S$. As a result, under either of the conditions defining the current case (i.e. $S \notin \widehat{C}(\{C\})$ or $u_S(\hat{a}, \hat{T}) > \underline{u}_S$), $\widehat{C}(\{C\}) \neq \{S\}$ must be true, so the scapegoat $S^* \in \arg\min_{S' \in \widehat{C}(\{C\}) \setminus \{S\}} u_{S'}(\hat{a}, \hat{T})$ is well defined.

Next we show that the (total) stage-game payoff of S^* is bounded. If $S \notin \widehat{C}(\{C\})$, then

$$S^* \in \underset{S' \in \widehat{C}(\{C\}) \setminus \{S\}}{\operatorname{arg\,min}} u_{S'}(\widehat{a}, \widehat{T}) = \underset{S' \in \widehat{C}(\{C\})}{\operatorname{arg\,min}} u_{S'}(\widehat{a}, \widehat{T})$$

and

$$u_{S^*}(\widehat{a},\widehat{T}) \le \frac{1}{|\widehat{C}(C)|} \sum_{S'\in\widehat{C}(C)} u_{S'}(\widehat{a},\widehat{T}) \le \frac{1}{|\widehat{C}(C)|} \max_{a\in A} \sum_{j\in C} v_j(a), \tag{11}$$

where the last inequality above follows from Assumption 1 and the fact that the default transfers are 0 in states $\underline{w}(S,\tau)$. Alternatively, if $S \in \widehat{C}(\{C\})$ and $u_S(\widehat{a},\widehat{T}) > \underline{u}_S$, then

$$u_{S^*}(\widehat{a},\widehat{T}) \leq \frac{1}{|\widehat{C}(\{C\})| - 1} \sum_{\substack{S' \in \widehat{C}(\{C\}) \setminus \{S\}}} u_{S'}(\widehat{a},\widehat{T})$$

$$= \frac{1}{|\widehat{C}(\{C\})| - 1} \Big[\sum_{\substack{S' \in \widehat{C}(\{C\})}} u_{S'}(\widehat{a},\widehat{T}) - u_{S}(\widehat{a},\widehat{T}) \Big]$$

$$\leq \frac{1}{|\widehat{C}(\{C\})| - 1} \Big[\sum_{\substack{S' \in \widehat{C}(\{C\})}} u_{S'}(\widehat{a},\widehat{T}) - \underline{u}_{S} \Big],$$

where the last inequality follows because we are considering the case $u_S(\hat{a}, \hat{T}) > \underline{u}_S$. Since the plan specifies zero default transfers, Assumption 1 ensures

$$u_{S^*}(\widehat{a},\widehat{T}) \le \frac{1}{|\widehat{C}(\{C\})| - 1} \Big[\max_{a \in A} \sum_{i \in C} v_i(a) - \underline{u}_S\Big].$$
(12)

Comparing the RHS of (11) and (12) to the bounds obtained in (9) and (10), it follows that we can find \overline{b}_2 such that $u_{S^*}(\widehat{a},\widehat{T}) < \overline{b}_2$.

Finally, to show that the blocking by C is not profitable, note that the (total) payoff of S^* is not improved by blocking if

$$(1 - \delta^{L(\delta) - \tau})v_{S^*}(\underline{a}_S) + \delta^{L(\delta) - \tau}u_{S^*}^S \ge (1 - \delta)\overline{b}_2 + \delta(1 - \delta^{L(\delta)})v_{S^*}(\underline{a}_{S^*}) + \delta^{L(\delta) + 1}u_{S^*}^{S^*}.$$

This inequality follows for sufficiently high δ from the same steps as that of the analogous part of Theorem 1. Based on the same arguments as in previous cases, the blocking is not profitable for C.

B.3 Proof of Proposition 1 on p. 21

Note that for each player $i \in \mathcal{F} \cup \mathcal{W}$, the individual minmax is $\underline{v}_i = 0$. The result then follows from applying Theorem 4.

B.4 Proof of Proposition 2 on p. 22

Preliminary Results.

Lemma 9. Let $\hat{x} = \max_{\phi \in A} \sum_{i \in N} v_i(\phi)$. If $u \in \mathbb{R}^n$ satisfies $\sum_{i \in N} u_i \leq \hat{x}$ and $\sum_{i \in C} u_i \geq \underline{u}_C$ for all $C \in \mathcal{E}$, then $\sum_{i \in N} u_i = \hat{x}$.

Proof. Take any $\tilde{u} \in \mathbb{R}^n$ satisfying $\sum_{i \in N} \tilde{u}_i \leq \hat{x}$ and $\sum_{i \in C} \tilde{u}_i \geq \underline{u}_C$ for all $C \in \mathcal{E}$. Towards a contradiction, suppose that \tilde{u} is not utilitarian efficient; that is, suppose $\sum_{i \in N} \tilde{u}_i < \hat{x}$. Then there exists an assignment $\phi' \in A$ such that $\sum_{i \in N} v_i(\phi') > \sum_{i \in N} \tilde{u}_i$. Let π' denote the partition of players into essential coalitions induced by the matching ϕ' , so $\pi' \subseteq \mathcal{E}$. It follows that there exists $C' \in \pi' \subseteq \mathcal{E}$ such that $\underline{u}_{C'} = \sum_{i \in C'} v_i(\phi') > \sum_{i \in C'} \tilde{u}_i$, which is a contradiction to the assumption that $\sum_{i \in C} \tilde{u} \geq \underline{u}_C$ for all $C \in \mathcal{E}$. So \tilde{u} must be utilitarian efficient. **Lemma 10.** Let $\hat{x} = \max_{\phi \in A} \sum_{i \in N} v_i(\phi)$. The set \mathcal{K} is characterized by

$$\mathcal{K} = \{ u \in \mathbb{R}^n : \sum_{i \in N} u_i = \widehat{x}, \sum_{i \in C} u_i \ge \underline{u}_C \text{ for all } C \in \mathcal{E} \}.$$
(13)

Proof. Take any $\tilde{u} \in \mathcal{K}$. Suppose, for the sake of contradiction, that there exists some $C \in \mathcal{E}$ such that $\sum_{i \in C} \tilde{u}_i < \underline{u}_C$, then \tilde{u} would be blocked by C, which contradicts the assumption that $\tilde{u} \in \mathcal{K}$. So $\sum_{i \in C} \tilde{u}_i \geq \underline{u}_C$ must hold for all $C \in \mathcal{E}$. Lemma 9 then implies that \tilde{u} is utilitarian efficient, so \tilde{u} satisfies the conditions in (13).

For the converse, take any \tilde{u} that satisfies the conditions in (13). We will show that $\tilde{u} \in \mathcal{K}$, i.e., there exists a core allocation (ϕ, T) such that $\tilde{u} = u(\phi, T)$. Since \mathcal{K} is nonempty, there exists a core alternative $(\tilde{\phi}, \tilde{T})$, which by the arguments above must satisfy $\sum_{i \in N} v_i(\tilde{\phi}) = \hat{x}$. Since $\sum_{i \in N} \tilde{u}_i = \hat{x}$, there exists $\tilde{T}' \in \mathcal{T}$ such that $\tilde{u} = u(\tilde{\phi}, \tilde{T}')$. Note however that \tilde{T}' may involve nonzero transfers between players who are not in an employment relationship, so $(\tilde{\phi}, \tilde{T}')$ may not constitute a matching. Nevertheless, let $\tilde{\pi}$ denote the partition of players induced by $\tilde{\phi}$. For every $C \in \tilde{\pi}$, it must hold that

$$\sum_{i \in C, j \notin C} \widetilde{T}'_{ij} - \sum_{i \in C, j \notin C} \widetilde{T}'_{ji} = 0,$$

for otherwise we would have $\sum_{i \in C'} \widetilde{u}_i < \sum_{i \in C'} v_i(\widetilde{\phi})$ for some $C' \in \widetilde{\pi}$, contradicting the fact that \widetilde{u} satisfies (13). Therefore, we can construct $\widetilde{T}'' \in \mathcal{T}$ such that

$$\widetilde{u} = u(\widetilde{\phi}, \widetilde{T}''), \quad \text{and} \quad \widetilde{T}''_{ij} \neq 0 \text{ only if } i = \widetilde{\phi}(j) \text{ or } i \in \widetilde{\phi}(j),$$

so $(\tilde{\phi}, \tilde{T}'')$ is a matching that induces payoff profile \tilde{u} . Since $\sum_{i \in C} \tilde{u}_i \geq \underline{u}_C$ for all $C \in \mathcal{E}$, $(\tilde{\phi}, \tilde{T}'')$ cannot be blocked by any coalition, so $(\tilde{\phi}, \tilde{T}'')$ is a core allocation, and therefore $\tilde{u} \in \mathcal{K}$.

Lemma 11. Let

$$\mathcal{U}^{\mathcal{M}} := co\Big(\Big\{u \in \mathbb{R}^n : \exists (\phi, T) \in \mathcal{M} \text{ such that } u = u(\phi, T)\Big\}\Big)$$

denote the convex hull of all feasible matching payoffs. Then

$$\left\{ u \in \mathcal{U}^{\mathcal{M}} : \sum_{i \in C} u_i \ge \underline{u}_C \text{ for all } C \in \mathcal{E} \right\} = \mathcal{K}.$$

Proof. The fact that $\mathcal{K} \subseteq \{u \in \mathcal{U}^{\mathcal{M}} : \sum_{i \in C} u_i \geq \underline{u}_C \text{ for all } C \in \mathcal{E}\}$ follows from the definition of \mathcal{K} .

To show $\{u \in \mathcal{U}^{\mathcal{M}} : \sum_{i \in C} u_i \geq \underline{u}_C \text{ for all } C \in \mathcal{E}\} \subseteq \mathcal{K}$, take any $\widetilde{u} \in \mathcal{U}^{\mathcal{M}}$, since \widetilde{u} is a convex combination of feasible payoff vectors, it must be that

$$\sum_{i \in N} \widetilde{u}_i \le \widehat{x} := \max_{\phi \in A} \sum_{i \in N} v_i(\phi).$$

Lemma 9 then implies that $\sum_{i \in N} \widetilde{u}_i = \widehat{x}$, so $\widetilde{u} \in \mathcal{K}$ by Lemma 10.

Proof of Proposition 2. We first prove that every payoff vector in \mathcal{K} can be supported by a public PCE. For any $\tilde{u} \in \mathcal{K}$ there exists core allocation (ϕ, T) such that $\tilde{u} = u(\phi, T)$. Consider the plan $\tilde{\sigma}$ defined by $\tilde{\sigma}(h) = (\phi, T)$ for all $h \in \overline{\mathcal{H}}$. The plan $\tilde{\sigma}$ is obviously public and produces discounted payoff profile u. Given that (ϕ, T) is a core allocation, $\tilde{\sigma}$ is also a PCE.

We now prove that for every $\delta \geq 0$, every public PCE implements a discounted payoff profile in \mathcal{K} . By Theorem 5, for every $\delta \geq 0$, every discounted payoff profile \widetilde{u} produced by a public PCE must satisfy $\sum_{i \in C} \widetilde{u}_i \geq \underline{u}_C$ for all $C \in \mathcal{E}$, so $\widetilde{u} \in \{u \in \mathcal{U}^{\mathcal{M}} : \sum_{i \in C} u_i \geq \underline{u}_C \text{ for all } C \in \mathcal{E}\}$. By Lemma 11, \widetilde{u} then must be an element of \mathcal{K} . \Box

B.5 Proof of Proposition 3 on p. 23

Preliminary Results.

Lemma 12. All static stable matchings fill slots $\{(f, l) : \rho(f, l) \ge \max\{0, \eta(L)\}\}$ while leaving other slots vacant; more over, all workers receive the same payoff r where $\max\{0, \eta(L+1)\} \le r \le \max\{0, \eta(L)\}.$

Proof. We break down the proof into two parts.

<u>Part 1</u>: All static stable matchings fill slots $\{(f,l) : \rho(f,l) \ge \max\{0,\eta(L)\}\}$ while leaving other slots vacant.

Let *m* be any static stable matching. We first show that *m* must be utilitarian efficient. Suppose, for the sake of contradiction, that *m* is not utilitarian efficient. Then there exists a reassignment of players that increases players' total payoff, which implies the existence of $f \in \mathcal{F}$ and $W \subseteq \mathcal{W}$ such that $v_f(W) + \sum_{w \in W} v_w(f) > v_f(m) + \sum_{w \in W} v_w(m)$. But this implies that *m* is profitably blocked by (f, W), contradicting the stability of *m*.

Next, we show that since m is utilitarian efficient, it fills all slots in $\{(f,l) : \rho(f,l) \ge \max\{0,\eta(L)\}\}$. Suppose, for the sake of contradiction, that there exists a slot $(\tilde{f},\tilde{l}) \in \{(f,l) : \rho(f,l) \ge \max\{0,\eta(L)\}\}$ that is not filled. Let $\tilde{l}^* = \min\{l : (\tilde{f},l)$ is unfilled} denote the first unfilled position at firm \tilde{f} . Since firms have diminishing marginal products, we have $\rho(\tilde{f},\tilde{l}^*) \ge \rho(\tilde{f},\tilde{l})$, so (\tilde{f},\tilde{l}^*) is an open slot in $\{(f,l) : \rho(f,l) \ge \max\{0,\eta(L)\}\}$ that is immediately accessible by workers. Since there are L workers in total, if not all slots in $\{(f,l) : \rho(f,l) \ge \max\{0,\eta(L)\}\}$ are filled, there exists $w' \in W$ who is either unemployed or filling a slot outside of $\{(f,l) : \rho(f,l) \ge \max\{0,\eta(L)\}\}$. In the first scenario, matching w' to the unfilled slot (\tilde{f},\tilde{l}^*) would strictly increase the total surplus. In the second scenario, let (\hat{f},\hat{l}) be the slot filled by w', and let $\hat{l}^* = \max\{l : (\hat{f}, l)$ is filled} denote the last occupied slot at firm \hat{f} , and \hat{w}^* denote the worker filling (\hat{f},\hat{l}^*) . It follows from decreasing marginal product that (\hat{f},\hat{l}^*) is also outside of $\{(f,l) : \rho(f,l) \ge \max\{0,\eta(L)\}\}$, so matching \hat{w}^* to the unfilled slot $\rho(\tilde{f},\tilde{l}^*)$ instead would strictly increase the total surplus, again contradicting the utilitarian efficiency of m. Thus, all stable matchings must fill the slots in $\{(f,l) : \rho(f,l) \ge \max\{0,\eta(L)\}\}$.

To show that all slots outside of $\{(f,l): \rho(f,l) \ge \max\{0,\eta(L)\}\}$ are vacant, there are two cases to consider. If $\eta(L) > 0$, we know from the arguments above that the L slots in $\{(f,l): \rho(f,l) \ge \eta(L)\}$ are filled, so all other slots must be vacant. If $\eta(L) < 0$, then the set $\{(f,l): \rho(f,l) \ge \max\{0,\eta(L)\}\)$ becomes $\{(f,l): \rho(f,l) \ge 0\}$, and let us suppose, for the sake of a contradiction, that some slot $(\tilde{f},\tilde{l})\)$ with $\rho(\tilde{f},\tilde{l}) < 0$ is filled. Let $\tilde{l}^* = \max\{l: (\tilde{f},l)\)$ is filled} denote the last filled slot at firm \tilde{f} , and let \tilde{w} denote the worker matched to this position. Due to decreasing marginal returns, we have $\rho(\tilde{f},\tilde{l}^*) < 0$ as well, so simply unmatching \tilde{w} from $(\tilde{f},\tilde{l}^*)\)$ will increase the total surplus. This contradicts the efficiency of m, which implies that m would not not stable. It follows that all stable matchings must leave slots outside $\{(f,l): \rho(f,l) \ge \max\{0,\eta(L)\}\}\)$

<u>Part 2</u>: All workers are paid the same wage r, where $\max\{0, \eta(L+1)\} \leq r \leq \max\{0, \eta(L)\}$.

First we establish that all workers have the same wage. Take any static stable matching m. From Part 1, all positions in $\{(f, l) : \rho(f, l) \ge \max\{0, \eta(L)\}\}$ are filled. We prove that all workers have the same wage under two separate cases.

First, suppose $\eta(L) < 0$. It follows that $|\{(f,l) : \rho(f,l) \ge \max\{0,\eta(L)\}\}| < L$, so in the stable matching *m* there exists a worker \widetilde{w} who is unmatched. This means \widetilde{w} receives 0 payoff in the stable matching *m*. It then follows that any other employed worker must also receive 0 payoff, since otherwise there is a profitable block where their employer replaces them with \tilde{w} .

Second, suppose $\eta(L) > 0$. Since by assumption $\rho(f, L) < \max\{0, \eta(L)\}$, there exists f_1 and f_2 such that both f_1 and f_2 employ workers in m. Since workers are identical, each worker working for f_1 must receive the same payoff as any worker at f_2 in m. This implies that workers at f_1 and f_2 all have the same payoff. The same argument applies to workers employed by any other firm, so all workers receive the same payoff.

Let r denote the payoff that workers receive, we next show that $\max\{0, \eta(L+1)\} \leq r \leq \max\{0, \eta(L)\}$. It is obvious that $r \geq 0$ by workers' individual rationality. To complete the arguments, it suffices to demonstrate the validity of three statements: A. $r \geq \eta(L+1)$ if $\eta(L+1) > 0$; B. $r \leq \eta(L)$ if $\eta(L) > 0$; and C. r = 0 if $\eta(L) \leq 0$.

Statement A: if $\eta(L+1) > 0$, then decreasing marginal return implies $\eta(L) > 0$, so from Part 1 we know all L workers are assigned to $\{(f,l) : \rho(f,l) \ge \max\{0,\eta(L)\}\} =$ $\{(f,l) : \rho(f,l) \ge \eta(L)\}$. Let (\tilde{f},\tilde{l}) denote the slot with value $\rho(\tilde{f},\tilde{l}) = \eta(L+1)$. By decreasing marginal return, any slot $\{(\tilde{f},l) : l < \tilde{l}\}$ at \tilde{f} is in $\{(f,l) : \rho(f,l) \ge \eta(L)\}$ and already filled. It follows that $r \ge \eta(L+1)$, since otherwise \tilde{f} can profitably block m by poaching a worker from other firms, which generates additional surplus $\eta(L+1)$, while offering wage r' satisfying $\eta(L+1) > r' > r$.

Statement B: if $\eta(L) > 0$, again from Part 1 we know that all L workers are assigned to $\{(f,l): \rho(f,l) \ge \eta(L)\}$. Let (\tilde{f},\tilde{l}) be the slot such that $\rho(\tilde{f},\tilde{l}) = \eta(L)$. By decreasing marginal return we know that (\tilde{f},\tilde{l}) must be the last filled slot at firm \tilde{f} . It follows that workers' payoff is no more than $\eta(L)$ since otherwise \tilde{f} can profitably block by firing the worker matched to the slot (\tilde{f},\tilde{l}) .

Statement C: if $\eta(L) < 0$, then there are at most (L-1) slots with a positive surplus, which by Part 1 implies that in any stable matching there exists a worker \tilde{w} who is unmatched. In this case, workers' payoff must be 0 since otherwise the matching is profitably blocked by a firm replacing one of its employees with worker \tilde{w} .

Combining statements A, B, and C lets us conclude that $\max\{0, \eta(L+1)\} \le r \le \max\{0, \eta(L)\}.$

Proof of Proposition 3. The first half of Proposition 3 follows from Proposition 1, while the second half of Proposition 3 follows from combining Proposition 2 and Lemma 12. \Box

B.6 Proof of Proposition 4 on p. 25

Since by assumption both markets M_1 and M_2 satisfy $\eta_i(L+1) > 0$, Lemma 12 implies that in each market M_i , where i = 1 or 2, all static stable matchings fill the slots in $\{(f,l) : \rho_i(f,l) \ge \max\{0,\eta_i(L)\}\} = \{(f,l) : \rho_i(f,l) \ge \eta_i(L)\}$. Moreover, the workers' payoff r in market M_i satisfies $\eta_i(L+1) \le r \le \eta_i(L)$. Recall that the total surplus is $\Pi_i := \sum_{\ell=1}^L \eta_i(\ell)$, while the set of potential workers' total surplus is $\Pi_i^{\mathcal{W}} = [L\eta_i(L+1), L\eta_i(L)]$, and the set of potential firms' total surplus is

$$\Pi_{i}^{\mathcal{F}} = \Pi_{i} - \Pi_{i}^{\mathcal{W}} = \left[\sum_{\ell=1}^{L} \eta_{i}(\ell) - L\eta_{i}(L), \sum_{\ell=1}^{L} \eta_{i}(\ell) - L\eta_{i}(L+1)\right].$$

To simplify notation let us denote $\underline{b}_i^{\mathcal{W}} := L\eta_i(L+1)$ and $\overline{b}_i^{\mathcal{W}} := L\eta_i(L)$, so $\Pi_i^{\mathcal{W}} = [\underline{b}_i^{\mathcal{W}}, \overline{b}_i^{\mathcal{W}}]$. Similarly, let $\underline{b}_i^{\mathcal{F}} := \sum_{\ell=1}^L \eta_i(\ell) - L\eta_i(L)$ and $\overline{b}_i^{\mathcal{F}} = \sum_{\ell=1}^L \eta_i(\ell) - L\eta_i(L+1)$, so $\Pi_i^{\mathcal{F}} = [\underline{b}_i^{\mathcal{F}}, \overline{b}_i^{\mathcal{F}}]$.

Let $s := \eta_1(1) = \eta_2(1)$. For each $2 \le \ell \le L$, define $\Delta_\ell^i := \eta_i(\ell - 1) - \eta_i(\ell)$, so $\eta_i(\ell) = s - \sum_{k=2}^{\ell} \Delta_k^i$ for all $\ell \ge 2$. It follows that

$$\sum_{\ell=1}^{L} \eta_i(\ell) = sL - \sum_{\ell=2}^{L} (L+1-\ell)\Delta_{\ell}^i,$$

$$L\eta_i(L) = sL - L\sum_{\ell=2}^{L} \Delta_{\ell}^i, \text{ and } L\eta_i(L+1) = sL - L\sum_{\ell=2}^{L+1} \Delta_{\ell}^i.$$

This allows us to express the bounds for firms' and workers' aggregate surplus in terms of s and Δ_{ℓ}^{i} 's, yielding

$$\underline{b}_{i}^{\mathcal{W}} = sL - L\sum_{\ell=2}^{L+1} \Delta_{\ell}^{i}, \text{ and } \overline{b}_{i}^{\mathcal{W}} = sL - L\sum_{\ell=2}^{L} \Delta_{\ell}^{i},$$
(14)

$$\underline{b}_{i}^{\mathcal{F}} = \sum_{\ell=2}^{L} (\ell-1)\Delta_{\ell}^{i}, \text{ and } \overline{b}_{i}^{\mathcal{F}} = \sum_{\ell=2}^{L+1} (\ell-1)\Delta_{\ell}^{i}.$$

$$(15)$$

Market M_2 exhibits more more steeply decreasing returns than M_1 is equivalent to $\Delta_{\ell}^2 \ge \Delta_{\ell}^1$ for all $2 \le \ell \le L$, which implies $\eta_2(\ell) \le \eta_1(\ell)$ for all $1 \le \ell \le L$, so $\Pi_2 \le \Pi_1$.

In (14), all the Δ_{ℓ}^{i} 's enter the bounds for worker surplus with negative coefficients, so $\underline{b}_{2}^{\mathcal{W}} \leq \underline{b}_{1}^{\mathcal{W}}$ and $\overline{b}_{2}^{\mathcal{W}} \leq \overline{b}_{1}^{\mathcal{W}}$, where the inequalities are strict if M_{2} has strictly more steeply decreasing returns than M_{1} . By contrast, in (15) the Δ_{ℓ}^{i} terms enter the bounds with

positive coefficients, so $\underline{b}_2^{\mathcal{F}} \geq \underline{b}_1^{\mathcal{F}}$ and $\overline{b}_2^{\mathcal{F}} \geq \overline{b}_1^{\mathcal{F}}$, where, again, the inequalities are strict if M_2 has strictly more steeply decreasing returns than M_1 . Together, the directions of change for these bounds imply $\Pi_2^{\mathcal{W}} \preccurlyeq_S \Pi_1^{\mathcal{W}}$ and $\Pi_2^{\mathcal{F}} \succcurlyeq_S \Pi_1^{\mathcal{F}}$, with strict set orders if M_2 has strictly more steeply decreasing returns than M_1 .

B.7 Proof of Proposition 5 on p. 28

Preliminary Results. We will use an alternative $a \in A$ to also represent its generated payoff profile v(a). We establish two preliminary results. Lemma 13 establish the existence of "punishment PCEs" $\{\sigma^i\}_{i=1}^n$ that guarantee $U_i(\emptyset|\sigma^i) = 0$ for each player *i*. Lemma 14 proves that any PCE can be enforced by punishments where every member of a deviating coalition simultaneously obtains 0.

Lemma 13. Under perfect monitoring, for every player $i \in N$, there is a PCE σ^i such that $U_i(\emptyset|\sigma^i) = 0$ when $\delta > \frac{n-2}{n-1}$.

Proof. We consider two case, |D| = 1 and $|D| \ge 2$. The case where |D| = 1 requires the discount factor to be sufficiently high. The case where there are two or more veto players $(|D| \ge 2)$ applies for every discount factor.

Case 1: |D| = 1. Suppose without loss of generality that D consists of player 1. Let $\widehat{a} := (1, 0, \dots, 0)$ denote the unique core alternative, and $\overline{a} := (0, \frac{1}{n-1}, \dots, \frac{1}{n-1})$ denote the alternative that equally divides the total payoff among all non-veto players.

For $i \neq 1$, let σ^i be the plan that specifies the core alternative \hat{a} as default after every history, so each σ^i is a PCE that satisfies $U_i(\emptyset | \sigma^i) = 0$

For i = 1, let σ^1 be the plan that specifies \overline{a} on path, and \widehat{a} at any history where blocking has occurred in the past. Note that $U_1(\emptyset|\sigma^1) = 0$. We will verify that σ^1 is a PCE. No coalition can profitably block once continuation play reverts back to the core alternative. On the path of play, consider a winning coalition $C \in \mathcal{W}$ blocking and choosing alternative a'. Since the game is non-dictatorial, if C is a winning coalition, player 1 cannot be its only member. Let $j \neq 1$ be a player in C. Since $a'_i \leq 1$, we have

$$(1-\delta)a'_j + \delta 0 \le 1 - \delta \le \frac{1}{n-1}$$

so player j prefers following the plan σ^1 over blocking and reverting to the core alternative. As a result, no coalition C can profitably block the plan σ^1 at any history, so σ^1 is a PCE. Case 2: $|D| \ge 2$. Without loss of generality, suppose $\{1, 2\} \subseteq D$. Let $a^1 := (1, 0, ..., 0)$ and $a^2 := (0, 1, 0, ..., 0)$ be two alternatives that allocate all payoff to player 1 and 2, respectively. It follows that both a^1 and a^2 are core alternatives.

Let σ^1 be the plan that specifies a^2 at all histories; for all $i \neq 1$, let σ^i be the plan that specifies a^1 at all histories. Each σ^i is a PCE, and $U_i(\emptyset | \sigma^i) = 0$ for every $i \in N$.

Lemma 14. Suppose \mathcal{U} is the set of PCE-supportable payoff profiles. For each player $i \in N$, let $\underline{u}_i := \min_{u \in \mathcal{U}} u_i$ be player *i*'s smallest possible payoff from PCEs. There is a stationary PCE with payoff profile a if and only if for every coalition C and alternative $a' \in E_C(a)$, there is a player $i \in C$ such that

$$(1-\delta)a_i' + \delta \underline{u}_i \le a_i \tag{16}$$

Proof. To see the "only if" direction, suppose (16) fails for some coalition C and $a' \in E_C(a)$. In other words, suppose there exists a coalition C and alternative a' such that

$$(1-\delta)a'_i + \delta \underline{u}_i > a_i$$
 for all $i \in C$.

Towards a contradiction, suppose also that there exists a stationary PCE σ that supports payoff a. Since σ is a PCE, it follows that $U_i(h|\sigma) \geq \underline{u}_i$ for every $i \in C$ and all $h \in \mathcal{H}$. As a result, for every $i \in C$,

$$(1-\delta)a'_i + \delta U_i(a', \{C\}|\sigma) \ge (1-\delta)a'_i + \delta \underline{u}_i > a_i$$

Moreover, since σ is stationary, it always plays a on path. The inequality above then implies that (a', C) is a profitable block for coalition C on path, contradicting σ being a stationary PCE.

For the "if" direction, (16) implies that for every coalition C and alternative $a' \in E_C(a)$, there exits a player i[a', C] and a PCE $\sigma[a', C]$ such that

$$(1-\delta)a'_{i[a',C]} + \delta U_{i[a',C]} (\emptyset \,|\, \sigma[a',C]) \le a_{i[a',C]}.$$
(17)

Since the stage game exhibits default-independent power, by Theorem 2, we can without loss assume that each $\sigma[a', C]$ is a stationary PCE.

Consider a plan σ^* that specifies a on path, but switches to $\sigma[a', C)$ if coalition

C blocks to implement a'. Inequality (17) implies that on path, no coalition can find profitably block. In addition, the fact that each $\sigma[a', C]$ is a PCE ensures that after any off-path history, no coalition can profitably block. Finally, σ is also stationary since it is stationary on path, and each $\sigma[a', C]$ is also stationary. Therefore, σ is a stationary PCE that supports payoff a.

Proof of Proposition 5.

Statement (a). Set $\underline{\delta} = \frac{n-2}{n-1}$. By Lemma 13, there exist PCEs $\{\sigma^i : i \in N\}$ satisfying $U_i(\emptyset|\sigma^i) = 0$ for all $i \in N$. It is straightforward to see that no players shared aligned payoffs in the stage game; in addition, no single player can form a winning coalition since the game is non-dictatorial. It follows that each player i's individual minmax is $\underline{v}_i = 0$. Moreover, this minmax payoff is achieved by the PCE σ^i .

By Lemma 14, in order for a payoff profile u to be supported by a stationary PCE, it is necessary and sufficient that for every winning coalition $C \in \mathcal{W}$, there exist no alternative $a' \in E_C(a)$ such that

$$(1-\delta)a'_i + \delta \cdot 0 = (1-\delta)a'_i > u_i \text{ for all } i \in C.$$
(18)

Note that the condition above is equivalent to $\sum_{i \in C} u_i \ge 1 - \delta$ for every $C \in \mathcal{W}$, since if $\sum_{i \in C} u_i < (1 - \delta) \cdot 1$ for some coalition $C \in \mathcal{W}$, there would be a certain $a' \in E_C(u)$ representing a division of total payoff 1 among players in C, such that (18) holds for every $i \in C$. It follows that a payoff profile u is supportable by a stationary PCE if and only if $\sum_{i \in C} u_i \ge 1 - \delta$ for every $C \in \mathcal{W}$. Finally, Theorem 2 implies that this same set is also the set of PCE-supportable payoff profiles.

Statement (b). Because a winning coalition can obtain the entire dollar by blocking, its minmax value is 1. This statement then follows immediately from Theorem 5.

Statement (c). Let $\widehat{\mathcal{W}}$ denote the set of minimal winning coalitions. By definition, $\widehat{\mathcal{W}} \subseteq \mathcal{W}$ so $\cap_{C \in \mathcal{W}} C \subseteq \cap_{C \in \widehat{\mathcal{W}}} C$. Furthermore, $\cap_{C \in \widehat{\mathcal{W}}} C \subseteq \cap_{C \in \mathcal{W}} C$, since otherwise there exists $i \in \cap_{C \in \widehat{\mathcal{W}}} C$ and $\widetilde{C} \in \mathcal{W}$ such that $i \notin \widetilde{C}$, but this would lead to a contradiction since \widetilde{C} must contains a winning coalition \widehat{C} , and $i \in \widehat{C}$. So $\cap_{C \in \widehat{\mathcal{W}}} C = \cap_{C \in \mathcal{W}} C = D$. By Theorem 5, every $C \in \widehat{\mathcal{W}}$ obtains total payoff 1. This implies the total payoff for players in D is 1.