

Negotiated Binding Agreements*

Malachy James Gavan[†]
University of Liverpool

November, 2023

Current Version: [\[Here\]](#)

Abstract

I study the binding agreements that may result from players negotiating over their behaviour in an underlying strategic environment. I propose a negotiation protocol where, in each round, agents make public proposals of the action they will take in the underlying game. The protocol terminates when these proposals are confirmed. Confirmation results in a binding agreement over the action profile and payoffs are the corresponding ones in the underlying game. I study the outcomes of Negotiated Binding Agreements of the negotiation protocol, which is a refinement of Subgame Perfect Equilibrium that I introduce in this context to obtain both credibility and tractability. A full characterisation of agreement outcomes is provided for two-player games. An outcome can be agreed upon if and only if appropriate individual punishments can be found, where both the outcome and punishments are defined in the underlying game itself. This condition is also sufficient in n -player games, while any outcome of the underlying game that is agreed upon necessarily satisfies an *iterative* individual rationality constraint. Finally, to allow for the possibility that agents make binding agreements over how they will negotiate, I extend the solution concept to allow for cooperative agreements within the negotiation game. Generalisations of the main results hold and refine the set of agreement outcomes.

Keywords: Agreements, Negotiation, Cooperation

JEL Codes: C70, C71, C72

*This paper was previously circulated under the title of “Negotiated Equilibrium”. This paper was chapter 1 of my PhD thesis, supervised by Antonio Penta, and has benefited greatly from his ongoing support. I have also had innumerable useful discussions with Larbi Alaoui, Alexander Frug and Pia Ennuschat, for which I am greatly indebted. I am also particularly grateful to Marco Mariotti and Olivier Bochet. I also thank (in alphabetical order) Nemanja Antic, Josefina Cenzon, Francesco Cerigioni, Vincent Crawford, Faruk Gul, Gianmarco Leon, Gilat Levy, Raquel Lorenzo, Zoel Martín Vilató, Rosemarie Nagel, Maria Ptashkina, Debraj Ray, Danila Smirnov and participants in a number of seminars and conferences. All faults are my own.

[†]*email:* malachy.gavan@liverpool.ac.uk.

1 Introduction

Negotiations and their resulting binding agreements play an important role in the economy.¹ Many such negotiations are over the behaviour that will be taken, rather the division of some abstract surplus. For instance, when prospective employees and employers negotiate, they may do so over pay but also the opportunity for flexible working, parental leave, vacation time, or other work benefits and conditions. In many cases, the non-monetary, or non-transferable, components may be the only aspect of negotiation. This would be true, for instance, if pay is fixed within a range, a common occurrence in public institutions or sectors with a strong union presence.² In this sense, the negotiation can be seen as over an agreement of the behaviour in game representing the provision of benefits by the employer and the acceptance of an offer by the employee, rather than a division of a surplus. Similarly, countries negotiate tariffs, quotas and regulation while committees may negotiate contributions to a public good, but may not have the ability to make direct monetary transfers due to political or legal reasons.³ All such situations are well described as a negotiation over the behaviour that will be taken within an underlying strategic environment.

Negotiations that are not over the “split of the pie” but rather by an agreement over what to play in a game and are clearly prevalent. Despite this, providing theoretical predictions for negotiated binding agreements over such environments has proven difficult. On one hand, some works provide fully specified models of negotiation and agreements, that ensure credibility of behaviour at all stages of the negotiation (Kalai, 1981; Bhaskar, 1989; Chwe, 1994; Mariotti, 1997). However, due to the complexity that these models entail, they do not provide results that can be applied to a broad range of environments. On the other hand, there are models that provide easy-to-use conditions for what can be agreed upon in the underlying environment (Aumann, 1961; Chander and Tulkens, 1997; Currarini and Marini, 2003). However, they abstract from credibility while negotiating, allowing for behaviour that may never be agreed upon to be taken when agents do not negotiate as expected. Presently, the tension between tractability of agreement outcomes and credibility of negotiation behaviour has been difficult to resolve. To bridge this gap, I propose a negotiation protocol for agreements over what to play in the underlying game and a refinement

¹To provide just one example, in 2022 alone, the Office of United States Trade Representatives estimated that trade between U.S. and Mexico was valued at over 800 billion USD, all of which was facilitated under the binding agreement of NAFTA. The latest round of this agreement was negotiated over many years.

²Negotiating over wages is rare in the US school system. An exception is in Wisconsin due to Act 10, which allowed individual negotiation over the employment contracts for teachers, leading to a better understanding of negotiation empirically (Baron 2018; Biasi 2021; Biasi and Sarsons 2021, 2022).

³Limao (2016) points to the fact that direct transfers may not be possible for international trade agreements and an increasing number of trade agreements include agreements surrounding policy, such as IP rights and foreign direct investment. For further discussion and modelling of these non-tariff barriers see Grossman et al. (2021).

1 of subgame perfect equilibrium, ensuring full credibility, while also showing that easy-to-use
2 conditions for agreement outcomes result.

3 The negotiation protocol I consider regarding the behaviour players should take in the
4 underlying game takes the following form. In each period, each agent makes a proposal
5 of the action they will take in the underlying game.⁴ Agents then observe the proposals
6 made by all others and can decide whether to “confirm” their choice or propose a new
7 action. Confirmation is modelled by proposing the same action again. If all agents confirm
8 their choice of action, a binding agreement is made. This binding agreement is to take
9 the confirmed action profile and therefore each agent receives the payoff of said outcome.
10 If any agent does not confirm their choice, the new proposed actions are observed and all
11 agents make the same confirmation or new proposal choice. Agents repeat this process until
12 confirmation, or agreement, is made by all agents. When agents never agree, or there is
13 *perpetual disagreement*, I make a weak assumption on the payoffs that result, consistent
14 with many interpretations, which is discussed in full in section 2.⁵

15 As the aforementioned negotiation protocol defines a dynamic game with complete in-
16 formation, to ensure agents always credibly negotiate in their own best interest, I explore a
17 refinement of Subgame Perfect Equilibrium. I impose two main refinement criteria: firstly, I
18 only consider Subgame Perfect Equilibria of the negotiation game that result in agreement,
19 as the agreement outcomes are the key objects of interest of this paper. Secondly, I impose
20 a *no babbling* condition where agents only make proposals of actions they could agree to.⁶ I
21 refer to this solution concept as *Negotiated Binding Agreements*. The agreement outcomes
22 of the Negotiated Binding Agreements will be referred to as *supported* by a Negotiated
23 Binding Agreement. As the negotiation game has infinitely many histories, with different
24 types of terminal histories, this is a complex object to consider. However, I show that
25 Negotiated Binding Agreements allow for a tractable solution, which I outline next.

26 Firstly, in section 3, I provide a full characterisation for the agreement outcome of
27 Negotiated Binding Agreements in two-player games. I show that any outcome of the
28 underlying game can be supported by a Negotiated Binding Agreement if and only if gives
29 each players a payoff weakly higher than an “individual punishment” profile, also defined
30 within the underlying game. These punishments are used as “threat” agreements, where
31 the punishment of an agent will be agreed to in the case that they do not act as expected

⁴This is similar to the approach of Kalai (1981); Bhaskar (1989); Kimya (2020); Nishihara (2022), etc.

⁵For instance, it is consistent with probabilistic termination, taking this probability of termination to 0 or taking the weighted average of all proposals made. Additionally, I show that the results of the paper are consistent with a number of variations in this *baseline procedure*, including in the timing of proposals, proposing action profiles, and in the payoff of perpetual disagreement, studied in the online appendix.

⁶The no babbling assumption can embed a form of no delay equilibrium used within bargaining games with a large number of players (Chatterjee et al., 1993), imposing that all proposed divisions of surplus can be agreed to.

1 when negotiating. The individual punishment profiles are such that (i) the payoff for any
 2 other players' punishment is weakly better than the payoff of their own punishment and
 3 (ii) when being punished, each player is prescribed to play their best response within the
 4 underlying game to their punishment profile. The logic of the result is simple. A player
 5 makes an agreement if they believe it is better than the worst agreement that could be made
 6 for them. However, the worst agreement must in it self be agreeable. Therefore there must
 7 have no reason to deviate from such an agreement, in this case by having no incentive to
 8 deviate in the underlying game.⁷ This result reflects, although is more restrictive than, the
 9 *player-specific punishments* used in the literature of infinitely repeated games, for example
 10 in Fudenberg and Maskin (1986) and Abreu et al. (1994) - which in it self can be seen as an
 11 agreement to play a strategy based on the threat of future punishment. Additionally, this
 12 provides a link to the Commitment Folk Theorems in the literature on contractable contracts
 13 (Peters and Szentes, 2012). Nonetheless, the characterisation of agreement outcomes can
 14 be substantially more restrictive than both in a number of games. As one of the most
 15 canonical examples in economics, I explore a leading example of a Cournot Duopoly with
 16 linear demand and heterogeneous marginal costs to illustrate the key ideas of the proof, as
 17 well as to demonstrate the key difference to the aforementioned models in a well understood
 18 environment.

19 As for games of n -players, in section 4, I show that the characterisation of agreement
 20 outcomes for two-player games can be used as a sufficient condition for n -player games. I
 21 show that a necessary condition for any action profile supported or proposed in a Negotiated
 22 Binding Agreement must survive *iterated elimination of individually irrational actions* of
 23 the underlying game, which I introduce in this paper. Specifically, an action a is *individually*
 24 *irrational* if, given the most optimistic beliefs the agent can have when evaluating it, the
 25 payoff it induces is still strictly worse than the minimum payoff that an agent can receive
 26 from best responding to some action profile of others. Performing this process iteratively,
 27 deleting *all* individually irrational actions within a round before moving to the next, results
 28 in actions that survive iterated elimination of individually irrational actions. I show that the
 29 minimum payoff that an agent can receive from a Negotiated Binding Agreement outcome
 30 is always weakly higher than the worst best response payoff in the underlying game, taken
 31 over the set of actions that survives iterated elimination of individually irrational actions.
 32 The conditions for deletion and calculating the minimum payoff are simple to implement in
 33 any underlying game that is finite or with smooth utility functions. Further, I show that in
 34 an important class of n -player games the characterisation is tight.

⁷Note that this logic is distinct from that of the solution to the Rubinstein (1982) bargaining game, where instead agents fear that lack of acceptance will lead to a discounted future agreement or no agreement at all. In this case, there is no cost of delay, and therefore no possibility of discounted future agreement.

To illustrate the logic of this necessary result, I use an underlying game of a simple three-bidder First Price Auction with heterogeneous valuations. In this case, in any profile of bids supported by a Negotiated Binding Agreement the bidder with the highest valuation must receive the good with positive probability. However, in comparison to the Nash equilibria of this underlying game, it is possible that *any* bidder receives the good with positive probability, at many different prices. Nonetheless, these possibilities are restricted by the minimal payoffs that bidders must receive, and therefore it is not the case that any outcome is possible.⁸

Negotiated Binding Agreements as a solution concept only contemplates unilateral deviations, but we may also be interested in the possibility of agents making binding agreements over *how* they will negotiate. To allow for this, I extend the solution concept to allow for cooperative agreements within the very negotiation game. I do so by introducing coalitions of agents to jointly choose a new strategy and will do so if it is profitable for all agents within the coalition. This may include permissible coalitions that overlap.

To capture the possibility of agents acting in such a way within the negotiation procedure, in section 5, I define the concept of \mathcal{C} -Negotiated Binding Agreement. Here, no coalition in a predefined set \mathcal{C} can profitably deviate at any history and a no babbling condition is imposed. I show that the natural extension of the baseline necessary and sufficient conditions for agreement outcomes in n -player games hold. Within \mathcal{C} -Negotiated Binding Agreement, players only make proposals from the set of actions that survives iterated elimination of *coalitionally* irrational actions in the underlying game, defined similarly to individually irrational actions taking coalition-wide preferences into account. Further, for all permissible coalition the outcome of negotiation must satisfy a notion of *coalitional rationality* in the underlying game. These conditions can be viewed as a perturbed version of the cooperative game theoretic notion of the β -core (Aumann, 1961).⁹ I provide sufficient conditions of the outcomes of the underlying game that can be supported using coalition-specific punishments; a further refined version of the β -core. In a simple Cournot model with fixed costs, I show that all the β -core outcomes can be sustained in \mathcal{C} -Negotiated Binding Agreement, while having the backing of a fully specified negotiation procedure, which is not the case for the β -core itself.

Finally, in section 6, I expand on the relations to these and other works within the literature review. I conclude the paper in section 7, pointing to a number of directions for future work.

⁸In the online appendix, I also provide an application of a public goods game with n -players and show that an agent can contribute if and only if a minimal level of aggregate contribution is reached, permitting both full contribution or no contribution at all.

⁹The β -core allows any outcome that is better than the worse-case scenario for any group to be agreed upon, even if these worst-case scenarios make use of non-credible behaviour that could never be agreed upon.

2 Model

Let the underlying game being negotiated over be $G = \langle N, (u_i, A_i)_{i \in N} \rangle$ where $N = \{1, 2, 3, \dots, n\}$ is a finite set of players, A_i is a set of actions for each player with typical element $a_i \in A_i$. $A = \times_{i \in N} A_i$ is the set of action profiles with typical element $a \in A$. u_i is utility function such that $u_i : A \rightarrow \mathbb{R}$ and u_i is bounded for all $i \in N$. Let $A_{-i} = \times_{j \neq i} A_j$.

I now define the *negotiation game* over G . There will be potentially infinitely many periods to reach an agreement and the process will take the following form. In each period, agents make a proposal of the action they will take within the underlying game G . Agents then observe the proposal made by all others. After doing so, they may simultaneously decide whether to “confirm” their choice by proposing the same action again, or alternatively propose a new action. If all agents confirm the proposal, an agreement is made, and that action profile is implemented in a binding way. If not, they continue to the next round and the same process occurs until confirmation is made by all agents, leading to an agreement. If there are infinitely many periods without agreement, I refer to this as *perpetual disagreement*.¹⁰

Formally, let the set of partial histories consists of all $h = (a^1, a^2, \dots, a^k)$ such that $a^t \neq a^{t-1}$ for any $t \leq k$ where $a^t = (a_i^t)_{i \in N}$ denotes the profile of proposals made in period t . I will denote the set of all partial histories by H . Proposals are assumed to be made simultaneously within a period, and therefore no history is such that only some agents have made proposals.¹¹

A history is terminal if, either:

- a) the same action profile is proposed twice in consecutive periods, and no earlier occurrence of consecutive repetition is present. That is, $z = (a^1, \dots, a^{k-1}, a^k)$ is terminal if $a^k = a^{k-1}$ and $a^m \neq a^{m-1}$ for all $m < k$. Let the set of such histories be denoted by Z' and refer to such histories as with *agreement*.
- b) or there is an infinite sequence of proposed action profiles where the same action profile is never proposed consecutively. Let the set of such histories be denoted by Z'' . I will refer to these as histories with *perpetual disagreement*.

Let the set of all terminal histories be given by $Z = Z' \cup Z''$.

Let $U_i : Z \rightarrow \mathbb{R}$ denote the payoff for player $i \in N$ of the negotiation game.

Whenever there is an agreement, it is assumed that the payoff is that of the agreed-upon action profile. Formally, whenever $z = (a^1, \dots, a^k) \in Z'$, that is a history that ends in

¹⁰Formally, this game is similar to that used in the farsighted stable set for games, which is discussed at length in the literature review in section 6.

¹¹See the online appendix for the extension of non-simultaneous proposals.

1 agreement, let $U_i(z) = u_i(a^k)$ for all $i \in N$.

2 Whenever there is perpetual disagreement, the payoff is defined to be between the
 3 \liminf and \limsup of the utility in the underlying game of the proposals made.¹² Formally,
 4 whenever $z = (a^1, a^2, \dots, a^k, \dots) \in Z''$, that is a terminal history with perpetual disagreement,
 5 I assume that $U_i(z) \in [\liminf_{t \rightarrow \infty} u_i(a^t), \limsup_{t \rightarrow \infty} u_i(a^t)]$.

6 This restriction is consistent a the standard model, where the proposal today is imple-
 7 mented with probability $(1 - \delta)$ for each period, while the process continues with probability
 8 δ , if the probability of continuation is taken to 1. Therefore, this can also be interpreted as
 9 a limiting version of the condition used within [Kimya \(2020\)](#), where there is a probability
 10 that the negotiation will end at the currently proposed actions.¹³ This is formalised by the
 11 following lemma, and the proof is provided in the appendix.

12 **Lemma 1.** For $z = (a^1, a^2, \dots, a^t, \dots) \in Z''$

$$\lim_{\delta \rightarrow 1} (1 - \delta) \sum_{t=1}^{\infty} \delta^{t-1} u_i(a^t) \in \left[\liminf_{k \rightarrow \infty} u_i(a^k), \limsup_{k \rightarrow \infty} u_i(a^k) \right]$$

13 By taking the view that the payoff of perpetual disagreement can take on *any* value from
 14 this set it weakens the reliance on the specific method of confirmation for agreement. So
 15 long as this confirmation is simultaneously made by all agents, the results would remain the
 16 same. To see this, notice that *any* payoff in the underlying game that is proposed countably
 17 infinitely many times can be used for the payoff of perpetual disagreement. Equally, if more
 18 than one profile of proposals is made a countably infinite number of times, one can easily
 19 be ignored. With this, it is possible to use a proposal to specifically avoid agreement,
 20 without it being used within the payoff of perpetual disagreement. Therefore, a proposal
 21 could be used to avoid a consecutive repetition, leading to confirmation, without impacting
 22 payoffs. With this, one may consider a single action of the underlying game being used as
 23 an “object” button, while confirmation of the previous choice is seen as an “accept” button,
 24 and unanimity of acceptance is needed for agreement. Therefore, one interpretation of
 25 the payoff of perpetual disagreement is that an there is an ϵ probability of each player
 26 mistakenly pressing accept, and such ϵ is taken to 0. This specification may also embed,
 27 for example, the approach of infinitely repeated games with no discounting: i.e. using the
 28 limit of means criteria when well defined ([Rubinstein, 1994](#); [Aumann and Shapley, 1994](#))
 29 where joint commitment is modelled.

30 The structure of the negotiation game has some similarities to the structure of repeated
 31 games, due to the structure of the partial histories and payoff of perpetual disagreement.

¹²See the online appendix for alternative specifications.

¹³Similar notions also exist in the context of [Rubinstein \(1982\)](#) bargaining, where [Busch and Wen \(1995\)](#) take a game to be played in each rejection phase, which is implemented with probability $1 - \delta$ and continuation occurs to a new proposal happens with probability δ , allowing for an endogenous outside option.

1 There are a few important differences. Firstly, repeated games only have one type of
 2 terminal history, where the underlying game has been repeated the specified number of
 3 times, be that some finite number or infinitely. This negotiation game allows for two
 4 distinct types of terminal histories, those with agreement and those without. Secondly,
 5 repeated games use flow payoffs, receiving a payoff in each period of play to guide strategic
 6 behaviour. This negotiation game only allows for payoffs to be realised upon termination.
 7 Identical disparities between negotiation games and repeated games are common in the
 8 literature (see Kalai 1981; Bhaskar 1989; Kimya 2020; Nishihara 2022, etc.).

9 At each round of the negotiation game, before agreements have been made, agents
 10 consider all previous proposals, both of themselves and others and decide on a new proposal
 11 to make. With this, strategies map each partial history to a new proposal of what they will
 12 play in an underlying game. Formally, at each partial history, $h \in H$ the strategy of $i \in N$
 13 dictates the proposal i would make in the next round: $s_i : H \rightarrow A_i$. Let S_i be the space of
 14 all such mappings. Let $s : H \rightarrow A$ be the joint strategy, such that $s(h) = (s_i(h))_{i \in N}$.

15 For a partial history $h \in H$ and a joint strategy s let $(s|h)$ denote the continuation
 16 history of h given by s . That is, $(s|h) = z \in Z$ such that $z = (h, a'^1, a'^2, \dots, a'^k, \dots)$
 17 where $a'^1 = s(h)$, $a'^2 = s((h, a'^1))$, $a'^k = s((h, a'^1, a'^2, \dots, a'^{k-1}))$. With some abuse of
 18 notation, let $U_i(s|h) = U_i(z')$ where $z' \in Z'$ is defined as before and $U_i(s|h) = U_i(z'')$, where
 19 $(s|h) = (h, z'') \in Z''$. That is, only take the continuation of the history h for perpetual
 20 disagreement. When $z = (a^1, a^2, \dots, a^k) \in Z'$, i.e. an agreement is made, let $a(z) = a^k$ and
 21 $a_i(z) = a_i^k$.

22 2.1. Solution Concept

23 This negotiation protocol defines a dynamic game with complete information therefore
 24 Subgame Perfect Equilibrium (SPE) is well defined.

25 **Definition** (Subgame Perfect Equilibrium). s^* is Subgame Perfect Equilibrium, if for all
 26 partial histories $h \in H$, for all $i \in N$, $U_i(s^*|h) \geq U_i(s_i, s_{-i}^*|h)$, for all $s_i \in S_i$.

27 Due to the structure of the negotiation protocol, in any SPE agents must receive a payoff
 28 weakly higher than their inf-sup payoff in the underlying game. This is true for any history.
 29 This is formalised by the following lemma.

30 **Lemma 2.** For any Subgame Perfect Equilibrium s^* , for any partial history $h \in H$

$$U_i(s^*|h) \geq \inf_{a_{-i} \in A_{-i}} \sup_{a_i \in A_i} u_i(a_i, a_{-i})$$

31 As I do not require that the utility functions are continuous and actions are from a
 32 compact set, the minimum or maximum need not exist. However, whenever the underlying

1 game being considered has well defined maxima and minima I will refer to them as such.

2 Note that the set of SPE of the negotiation game trivially includes many perpetual
 3 disagreement outcomes. As this work is primarily focused on the agreements that can be
 4 supported by some equilibrium, I focus on the SPE of the negotiation game that reaches an
 5 agreement from the initial history. I will also restrict attention to SPE where proposals only
 6 involve actions that can be part of an agreement. This rules out agents proposing actions
 7 that they would never agree to on or off the path of play. I will refer to such property
 8 as agreements having *no babbling*. Similar concepts have been used within the literature
 9 on bargaining. For instance, no delay equilibrium of Chatterjee et al. (1993), where the
 10 proposals can only be made, at any history, if they could be accepted. I will refer to this
 11 concept as *Negotiated Binding Agreement*. Formally:

12 **Definition 1** (Negotiated Binding Agreement). s^* is a *Negotiated Binding Agreement* if:

13 a) s^* is a *Subgame Perfect Equilibrium* of the negotiation game.

14 b) *No babbling*: $\forall h \in H, \exists h' \in H$ such that $s_i^*(h) = a_i(s^*|h')$.

15 a^* is supported by s^* if $a^* = a(s^*|\emptyset)$.

16 Note that the set of actions that may be supported by a Negotiated Binding Agreement
 17 does not change if the no babbling condition were to be defined as only making proposals
 18 that could be agreed to in *some* Negotiated Binding Agreement, rather than no babbling
 19 requiring proposed actions must be agreed to in the Negotiated Binding Agreement being
 20 considered.

21 3 Negotiated Binding Agreement Outcomes for Two-Player Games

22 In this section, I provide a full characterisation of the Negotiated Binding Agreement out-
 23 comes for two-player games where the underlying action space is compact and the utility
 24 function is continuous. As outlined in the introduction, the logic of the characterisation is
 25 as follows. Each player will be willing to agree to an outcome if it is better than the worst
 26 possible agreement, from their perspective, that could be implemented. This immediately
 27 implies that there is a worst agreement for each individual player, call them \underline{a}^1 and \underline{a}^2 for
 28 the worst agreement for player 1 and player 2 respectively and call this their “punishment”.
 29 Clearly, it must be that $u_i(\underline{a}^i) \leq u_i(a)$ for any agreement outcome a , including \underline{a}^{-i} . Further,
 30 it must be that player i is willing to agree to their worst agreement. It turns out, that it can
 31 be shown that the condition for ensuring the worst agreement for player i is agreeable is that
 32 they best respond to their “punishment” in the underlying game. Therefore, the agreement

outcomes for the Negotiated Binding Agreements can be completely characterised purely with information of the underlying game, in a simple way. To further demonstrate the logic of this result, I now turn to a Cournot Duopoly with Linear Demand and Heterogeneous costs, where I will also discuss the distinction of the Negotiated Binding Agreement outcomes, player specific punishment (Fudenberg and Maskin, 1986; Abreu et al., 1994) and commitment folk Theorems (Peters and Szentes, 2012).

3.1. Leading Example and Preview of Results for two-player games

Consider a simple Cournot Duopoly model as the underlying game, G . Let $q_1, q_2 \in [0, b] = A_i$ be the quantities produced and the inverse demand be given by $p(q_1, q_2) = \max\{b - q_1 - q_2, 0\}$. Let firms have potentially heterogeneous costs, c_1 and c_2 . Without loss of generality let $c_1 \geq c_2 \geq 0$. Assume that firm 1 is a viable competitor: $\frac{b+c_2}{2} \geq c_1$. Profits are given by $\pi_i(q_1, q_2) = q_i(p(q_1, q_2) - c_i)$.

Notice that the best responses in the underlying game are given by:

$$q_i^*(q_{-i}) = \begin{cases} 0 & \text{if } q_{-i} \geq b - c_i \\ \frac{b - c_i - q_{-i}}{2} & \text{if } q_{-i} < b - c_i \end{cases}$$

The Nash equilibrium of this underlying game is given by quantities of $\left(\frac{b+c_2-2c_1}{3}, \frac{b+c_1-2c_2}{3}\right)$, leading to payoffs of $\left(\left(\frac{b+c_2-2c_1}{3}\right)^2, \left(\frac{b+c_1-2c_2}{3}\right)^2\right)$.

Consider supporting (q_1^*, q_2^*) such that $\pi_1(q_1^*, q_2^*) \geq 0$ and $\pi_2(q_1^*, q_2^*) \geq (c_1 - c_2)^2$ in a Negotiated Binding Agreement. Note given the assumption that $\frac{b+c_2}{2} \geq c_1$ it follows that $(\frac{b-c_2}{2})^2 \geq (c_1 - c_2)^2$ and therefore such a profile exists. Consider the following strategies to do so. Take $\underline{q}_2^1 = b - c_1$ and $\underline{q}^2 = (b - 2c_1 + c_2, c_1 - c_2)$.

1. [Firm 1's punishment for deviating] Let $s^*(h') = (0, \underline{q}_2^1)$ whenever $h' = (q^1, q^2, \dots, (q'_1, q_2^*))$, $q'_1 \neq q_1^*$, $h' = (q^1, q^2, \dots, (q''_1, \underline{q}_2^1))$, or $(q^1, q^2, \dots, (q'_1, \underline{q}_2^2))$, $q'_1 \neq \underline{q}_1^2$.
2. [Firm 2's punishment for deviating] Let $s^*(h'') = \underline{q}^2$ whenever $h'' = (q_1, q_2, \dots, (q_1^*, q'_2))$, $q'_2 \neq q_2^*$, $h'' = (q_1, q_2, \dots, (\underline{q}_1^2, q_2''))$, or $h'' = (q^1, q^2, \dots, (0, q_2''))$, $q_2'' \neq \underline{q}_2^1$.
3. [No / multilateral deviations] Otherwise, $s^*(\emptyset) = s^*(h) = (q_1^*, q_2^*)$ for all other h .

The intuition of this Negotiated Binding Agreement is to have each firm propose to flood the market as much as possible whenever the other firm is not acting as expected when negotiating, while maintaining a positive profit, understanding the other agent will propose their best response in the underlying game. Notice that if $c_1 > c_2$ firm 1 cannot entirely flood firm 2 out of the market, while maintaining positive profits, when firm 2 has not negotiated as expected.

To see this is a Negotiated Binding Agreement, notice that no babbling applies as all three rules are absorbing. Therefore all that is left is to check that s^* is a Subgame Perfect Equilibrium of the negotiation game.

First, let us consider firm 1. Firstly, consider the strategy in case 1, where firm 1 faces punishment for deviation. In this case, regardless of what they propose, firm 2 will continue to propose $b - c_1$ for all periods. Given this, firm 1 cannot profitably produce statically, and therefore they cannot improve upon the current strategy. Now let us consider a deviation of firm 1 from the punishment of firm 2. Under the current strategy and rule, no deviation will lead to an agreement of \underline{q}^2 . In which case, firm 1 receives a profit of 0. However, a deviation can only lead to firm 2 proposing $b - c_1$ in every period. With this, firm 1's payoff would be pinned down by $\pi_1(q'_1, b - c_1) \leq 0$. With this, it cannot be profitable to deviate. Finally, consider a deviation from any other history. No deviation will lead to an agreement for q^* , inducing a profit of $\pi_1(q^*) \geq 0$. Again, any deviation can only lead to firm 2 proposing \underline{q}_2^1 for all subsequent periods. Given this, it must be that the payoff of said deviation is again at most 0, and therefore cannot be profitable.

Now instead consider firm 2. Firstly, consider the first case above, and consider whether there could be a profitable deviation from punishing firm 1. If firm 2 does not deviate, this will lead to a profit of $\pi_2(0, \underline{q}_2^1) = (b - c_1)(c_1 - c_2)$. However, a deviation will lead to firm 1 proposing \underline{q}_1^2 in all subsequent periods. With this, a deviation will lead to a payoff at most the static best response to \underline{q}_1^2 , $\pi_2(\underline{q}_1^2) = (c_1 - c_2)^2$. However, this can not be profitable due to the viable competitor assumption. In a similar vein as firm 1, it cannot be that it is profitable for 2 to deviate from their punishment, due to statically best responding to their own punishment, nor case 3, as this would lead to an agreement for q^* , which provides them with a higher payoff than $(c_1 - c_2)^2$.

Now we will study why the Negotiated Binding Agreement outcome must be necessarily better than that of the outlined punishment \underline{q}^i . Firstly, notice that due to both firms being restricted to only making proposals that they can agree to, in any Negotiated Binding Agreement $s^{*,'}$, at any history $s^{*,'}(h) \in Q^*$, where Q^* is some set of agreement outcome quantities. Now notice that in any Negotiated Binding Agreement $s^{*,'}$ it is not possible that some firm i receives a lower payoff than the one prescribed by $\max_{q_i \in [0, b]} \min_{q_{-i} \in Q_{-i}^*} \pi_i(q_i, q_{-i})$, as they could elect to best respond statically in each period.¹⁴ With this, the agreement payoff is still bounded below by their worst best response payoff in Q_{-i}^* . This is true for both firm 1 and 2. Further, notice that if a profile is included in Q^* , then any profile that provides both players with a higher payoff must also be included in Q^* . All that is left to show is that \underline{q}^i prescribes the worst best response payoff for each player in Q^* . As $\pi_1(\underline{q}^1) = 0$, it is clear

¹⁴I leave the argument that Q_{-i}^* is compact for the formal proof of theorem 1.

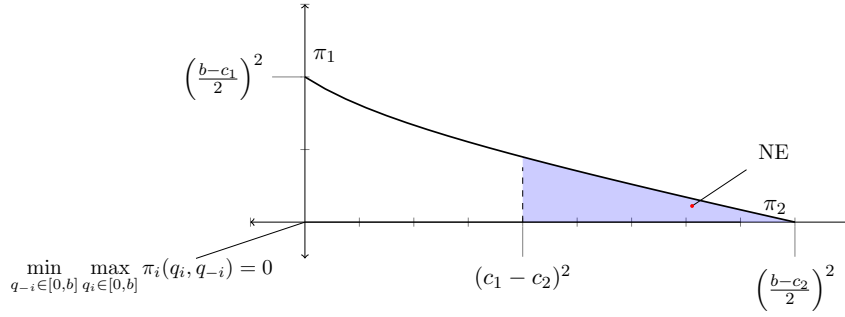


Figure 1: The payoff space of the Cournot Duopoly in example 3.1., where costs are such that $0 \leq c_2 < c_1 < \frac{c_2+b}{2}$. The black curve represents the payoff frontier. The shaded blue area represents that payoffs that can be sustained in Negotiated Binding Agreements.

that this is the lowest best response payoff. Further, as any agreement outcome must guarantee firm 1 a payoff of at least 0, even with firm 2's quantity taken into account, it must be that in any agreement outcome $q^{*,'} \in Q^*$ $\pi_1(q^{*,'}) \geq 0$. With this, taking into account firm 2's best response, to account for $\max_{q_2 \in [0, b]} \min_{q_1 \in Q_1^*} \pi_2(q_1, q_2)$, firm 1 can produce at most $b - 2c_1 + c_2$. This leads to a minimum payoff for firm 2 of $(c_1 - c_2)^2$. Showing that above strategy fully characterises the q^* s that may be supported under Negotiated Binding Agreements.

The payoff space is represented by figure 1.

Note that the construction provides us with some natural comparative statics and comparison to player specific punishments of Fudenberg and Maskin (1986); Abreu et al. (1994) and the commitment folk theorems of Peters and Szentes (2012). If $c_1 = c_2$, then both players may agree to an outcome that provides any payoff above their individually rational payoff of 0. In this case, the resulting set of agreement outcomes is identical to that of the commitment folk theorem and the payoff space of individual punishments. However, if it is not the case, and $c_1 > c_2$, then the agreement outcome of the commitment folk theorems and the payoff space of individual punishments remains unchanged. However, under Negotiated Binding Agreements the space is restricted to reflect the additional bargaining power firm 2 has due to firm 1 not proposing or agreeing to outcomes known to be bad for themselves. In the most extreme case, when firm 1 is no longer a viable competitor, when $c_1 = \frac{b+c_2}{2}$, we conclude that firm 2's profit is $\pi_2(q^2) = \left(\frac{b-c_2}{2}\right)^2$, their monopoly profit.¹⁵

¹⁵Note that if $c_1 > \frac{b+c_2}{2}$ then the only outcome that can be supported by a Negotiated Binding Agreement is $q_1^* = 0$, $q_2^* = \frac{b-c_2}{2}$, while under commitment folk theorems and individual punishments all individually rational payoffs would still be supported.

3.2. Results

In this subsection, I provide show that in two-player games the result generalises beyond the Cournot application. Informally, the result states that in two-player games, a target profile of the underlying game a^* is supported by a Negotiated Binding Agreement **if and only if** there is a punishment for each player such that: a) Each player best responds to their punishment in G b) They prefer the other's punishments to their own and c) They prefer the target to their punishment.¹⁶ In essence, this relies on player-specific punishment strategies, that have been used for sufficiency for SPE in infinitely repeated games (Fudenberg and Maskin, 1986; Abreu et al., 1994). The requirements in my paper are more stringent, as the profile used to punish i must use i 's best response to the punishment in the baseline game. This is because there is no future payoff to compensate for abiding to the punishment, as the agreement to play such an action is binding and the negotiation game terminates.

Theorem 1 (Full Characterisation for Two-Player Games). *For any game G such that $N = \{1, 2\}$, A_i is a compact subset of a metric space for $i = 1, 2$ and u_i is continuous, then a^* is supported by a Negotiated Binding Agreement if and only if $\exists \{\underline{a}^1, \underline{a}^2\} \subseteq A$ such that:*

1. $\underline{a}_i^i \in \arg \max_{a_i \in A_i} u_i(a_i, \underline{a}_{-i}^i)$.
2. $u_i(\underline{a}^j) \geq u_i(\underline{a}^i)$ for all $i \neq j$.
3. $u_i(a^*) \geq u_i(\underline{a}^i)$.

It is worth noting that any pure Nash equilibrium of the game G is supported by a Negotiated Binding Agreement. Further, any action profile that Pareto dominates a pure Nash equilibrium in the underlying game can be sustained by this reasoning. However, in underlying games where no pure Nash equilibrium exists there may exist a Negotiated Binding Agreement due to the above characterisation.

Example 1. Consider the underlying game, G , being the following two-player game. For clarity, I have underlined the corresponding best responses in the baseline game.

1\2	L	C	R
T	7,7	<u>4,4</u>	0, <u>12</u>
M	4, <u>4</u>	0,0	<u>2</u> ,3
D	<u>12</u> ,0	3, <u>2</u>	1,1

Notice that there is no pure Nash equilibrium in this underlying game. However, there exists a Negotiated Binding Agreement. Specifically, take $a^* = (T, L)$, while tak-

¹⁶These conditions can be viewed as similar to *player contingent threats* of Greenberg (1990), as the negotiation game can be viewed as a special case of a social situation.

ing $\underline{a}^1 = (M, R)$ and $\underline{a}^2 = (D, C)$, which satisfies the assumptions. Therefore there exists a Negotiated Binding Agreement that supports (T, L) , while there is no pure Nash equilibrium in the underlying game. ▼

4 Negotiated Binding Agreement Outcomes for n-Player Games

In this section I will explore the necessary and sufficient conditions for n -player games. First, I show that the idea of the characterisation for two-player games is still sufficient for n -player games. However, it is no longer necessary, as the no babbling condition does not impose strong coordination on the action *profile* proposed by players after a deviation. Therefore it is not possible to look for a punishment that is best responded to with strong conditions on coordination. However, when strong conditions on coordination of agreement outcome *profiles*, i.e. not only does a player have to propose an action they would agree to, but the profile of actions proposed by any set of agents is such that they would jointly agree, then this condition returns to being necessary. Outside of imposing this condition, I show that an iterative individual rationality constraint on the underlying game, which I call *iterated elimination of individually irrational actions* is necessary for agreement outcomes to be supported by a Negotiated Binding Agreement.

4.1. Sufficiency

The logic of the sufficient condition for agreement outcomes of two-player games can be generalised to n -player games. Specifically, we require that each agent has a specific punishment action profile in the underlying game, which I will denote \underline{a}^i . For this action profile, i will best respond to \underline{a}_{-i}^i in the baseline game G . The action of the underlying game that is sustained in Negotiated Binding Agreement a^* must, for each player i , give a weakly higher payoff than \underline{a}^i . Further, I will require that the punishment of other agents gives a weakly higher payoff than the punishment for i in the underlying game. If such a collection of action profiles exist, then a^* can be supported by a Negotiated Binding Agreement. I state this formally in the following theorem.

Theorem 2. *Take any underlying game, G , such that $\exists \{a^*, \underline{a}^1, \dots, \underline{a}^n\} \subseteq A$ such that:*

1. $\underline{a}_i^i \in \operatorname{argmax}_{a_i \in A_i} u_i(a_i, \underline{a}_{-i}^i)$
2. $u_i(a^*) \geq u_i(\underline{a}^i)$
3. $u_i(\underline{a}^j) \geq u_i(\underline{a}^i)$ for all $i, j \in N$

Then a^ can be supported in a Negotiated Binding Agreement.*

4.1.1. A Refinement of Negotiated Binding Agreements where outcomes are Fully Characterised

Further justification for the general sufficient conditions can be found. For a refinement of Negotiated Binding Agreements, where the focus is upon SPE that end in immediate agreement following from each history, the sufficient conditions for agreement outcomes are also necessary for underlying games where the action space is a compact subset of a metric space and utility is continuous. This No Delay condition applies for all possible histories, and therefore applies both on and off the path. I refer to this solution as No Delay Negotiated Agreements and is similar to the no delay equilibrium proposed by Chatterjee et al. (1993). Therefore, for the class of No Delay Negotiated Binding Agreements, I fully characterise the set of outcomes that can be supported. Here I formally define No Delay Negotiated Binding Agreement and state the formal result.

Definition 2 (No Delay Negotiated Binding Agreement). s^* is a No Delay Negotiated Binding Agreement supporting $a^* = a(s^*|\emptyset)$ if:

- a) s^* is a Subgame Perfect Equilibrium of the negotiation game.
- b) No Delay: For all partial histories $h \in H$, $s^*(h) = s^*(h, s^*(h)) = a^*(s^*|h)$.

Proposition 1. For any underlying game G such that A_i is a compact subset of a metric space and u_i is continuous for all $i \in N$, a^* is supported by a No Delay Negotiated Binding Agreement, s^* , if and only if $\exists \{\underline{a}^1, \dots, \underline{a}^n\} \subseteq A$ such that:

- 1. $\underline{a}_i^i \in \operatorname{argmax}_{a_i \in A_i} u_i(a_i, \underline{a}_{-i}^i)$
- 2. $u_i(a^*) \geq u_i(\underline{a}^i)$
- 3. $u_i(\underline{a}^j) \geq u_i(\underline{a}^i)$ for all $i, j \in N$

Finally, note that within the literature on agreements it is common to use the notion of Perfect Equilibrium of Selten (1988), for instance in Kalai (1981) and Bhaskar (1989). This is a Subgame Perfect Equilibrium that does not make use of weakly dominated strategies at any history. Notice that this does not have a significant change in the results, and to ensure the sufficient conditions for agreement outcomes remain true for this refinement, as well as the no babbling and agreement for all histories condition, the only check is to ensure that the action \underline{a}_i^i is not weakly dominated in the underlying game G .

Before moving to the necessary conditions for the Negotiated Binding Agreement outcomes, I turn to the following example to preview the logic.

4.2. Preview of Necessary Conditions n -player games

Here the underlying game G is a 3 player single unit First Price Auction with heterogeneous valuations. Specifically, there are three bidders, $N = \{1, 2, 3\}$. Each bidder has a value for the good, v_i . It is assumed that $v_1 = 6$, while $v_2 = 5$ and $v_3 = 2$. Each bidder may bid an integer from 0 to 7, $b_i \in \{0, 1, \dots, 7\}$.¹⁷ The highest bidder wins the good with uniform probability and pay their bid. Bidders who do not win the good receive a utility of 0. Therefore utility is given by their probability of winning, multiplied by their value minus their bid. Formally,

$$u_i(b) = \begin{cases} \frac{v_i - b_i}{|\arg\min_{j \in \{1, 2, 3\}} b_j|} & \text{if } i \in \arg\min_{j \in \{1, 2, 3\}} b_j \\ 0 & \text{if } i \notin \arg\min_{j \in \{1, 2, 3\}} b_j \end{cases}$$

The Nash equilibria of this underlying game are such that either a) $b^* = (5, 5, b_3^*)$ with $b_3^* \leq 4$ leading to payoffs of $(1/2, 0, 0)$, b) $b^* = (5, b_2^*, b_3^*)$ with $\max\{b_2^*, b_3^*\} = 4$ with payoffs of $(1, 0, 0)$ or c) $b^* = (4, 4, b_3^*)$ with $b_3^* \leq 3$, leading to payoffs of $(1, 1/2, 0)$. Notice those within a) make use of weakly dominated actions. The lowest payoffs in any Nash equilibria are $(1/2, 0, 0)$, with $b^* = (5, 5, b_3^*)$.

Firstly, can it be that any bidder agrees to the maximal bid, $b_i = 7$, in a Negotiated Binding Agreement? If this were the case, bidder i would receive a strictly negative utility, as they would certainly win the auction with positive probability and at a price above their valuation. However, they could avoid such an outcome by deciding to propose their own valuation in every round of negotiation, $s_i(h) = v_i$ for all $h \in H$. If they did so, regardless of whether the negotiation game ended in agreement or perpetual disagreement, they would receive a payoff of 0. This is because the payoff can only be pinned down by losing the auction, or by winning at their valuation, leading to a payoff of 0. More concretely, bidding $b_i = 7$ is *individually irrational* in the underlying game, which will be formalised in the next section, as they can guarantee themselves a higher payoff. With this, it cannot be that agreeing to bid $b_i = 7$ is supported by a Negotiated Binding Agreement, as such a strategy cannot be a Subgame Perfect Equilibrium of the negotiation game. Further, by no babbling, it cannot be that proposing to bid 7 occurs in *any* Negotiated Binding Agreement.

Now consider whether it is the case that bidders 2 or 3 could agree to bid 6 in a Negotiated Binding Agreement. By the previous argument, we conclude that agreeing to bid 6 will result in winning the good with positive probability, as we know no bidder will ever bid 7. With this, as the valuations of bidders 2 and 3 are below 6, it must be they receive a strictly negative payoff from such an agreement. However, we can again consider a

¹⁷The maximal bid being 7 is not important for the analysis, we only need to ensure payoffs are bounded by including some maximum bid.

deviation of these firms in the negotiation game to always propose their valuation, ensuring a payoff of 0. More concretely, bidding 6 is *individually irrational* for bidders 2 and 3 in the underlying game, again formalised in the next section, as they can guarantee themselves a higher payoff, given that 7 cannot be bid, and therefore cannot be agreed to. With this, we conclude that such an agreement cannot be a Negotiated Binding Agreement, as it would not be a Subgame Perfect Equilibrium of the negotiation game. Further, by the no babbling condition, we conclude that in no Negotiated Binding Agreement can bidding 6 *ever* be proposed by bidders 2 and 3.

We can continue this induction, concluding that bidder 1 would also never bid 6 once bidders 2 and 3 will not. Bidder 3 would never bid 5, due to this bid being *iteratively individually irrational* in the underlying game.

By the same argument as ruling out such bids, we conclude that any Negotiated Binding Agreement must provide bidders 2 and 3 with a payoff of at least 0. Also notice in any Negotiated Binding Agreement it must be that bidder 1 receives a payoff of at least $1/2$. To see this, notice that the worst possible stream of proposals for bidder 1 is that bidders 2 and 3 bid their highest possible bid in every round of the negotiation game, 5 and 4 respectively. Given this, bidder 1 can simply respond by bidding 5 in every round of the negotiation game, guaranteeing a payoff of $1/2$.¹⁸

With this, I move on to provide general necessary conditions, which this example has already pointed to.

4.3. Necessary Conditions

Within this section, I characterise a number of necessary conditions for a Negotiated Binding Agreement outcomes and strategies for n -player games. First, I will show that any action proposed at some history of Negotiated Binding Agreement must survive a procedure of *iterated deletion of individually irrational actions* in the underlying game, which to my knowledge is a novel definition. This procedure works inductively as follows. If an individual's action, regardless of the action profile of other agents chosen, always provides a payoff that is not individually rational, in the sense of inf-sup utility, then it is individually irrational. In the iterated elimination we can therefore remove said actions from consideration. Now, upon deleting such actions, we proceed inductively. If an individual's action, regardless of the action profile of other agents chosen *within* the set that has survived iterated deletion of individually irrational actions, always provides a payoff that is not individually rational, in the sense of inf-sup utility, where the inf is taken *over the set of actions that survives iterated individual rationality*, then it does not survive iterated deletion

¹⁸It is worth noting that the sufficient conditions would imply the same, but for different reasoning.

1 of individually irrational actions. The formal definition of individual irrational actions and
 2 iterated deletion of individually irrational actions are formally defined below.

3 **Definition 3** (Individually Irrational actions given $C_{-i} \subseteq A_{-i}$). For a game G , $a_i \in A_i$ is
 4 individually irrational given $C_{-i} \subseteq A_{-i}$ if:

$$\inf_{a'_{-i} \in C_{-i}} \sup_{a'_i \in A_i} u_i(a'_i, a'_{-i}) > \sup_{a_{-i} \in C_{-i}} u_i(a_i, a_{-i})$$

5 Denote the set of actions that are individually irrational given C_{-i} by $D_i(C_{-i})$.

6 This notion is similar to the notion of absolute dominance by [Salcedo \(2017\)](#), simulta-
 7 neously developed in [Halpern and Pass \(2018\)](#), who instead compare the best case of one
 8 action and the worst case of another, whereas I compare based on the best case of an action
 9 compared to the inf-sup.¹⁹ Therefore the set that survives elimination of individually irra-
 10 tional actions is smaller. Note that, if in a normal form game there is a single action that
 11 is not absolutely dominated given A_{-i} , then this action is an obviously dominant strategy
 12 as defined by [Li \(2017\)](#). Therefore if a single action is not individually irrational it is also
 13 obviously dominant.

14 **Definition 4** (Iterated Deletion of Individually Irrational Actions). For a game G , let
 15 $\tilde{A}_i^0 = A_i$ for all $i \in N$. Let $\tilde{A}_{-i}^0 = A_{-i}$. Then for all $m > 0$ let $\tilde{A}_i^m = \tilde{A}_i^{m-1} \setminus D_i(\tilde{A}_{-i}^{m-1})$
 16 where $\tilde{A}_{-i}^{m-1} = \times_{j \neq i} \tilde{A}_j^{m-1}$.

17 The set of actions that survive iterated deletion of individually irrational actions, or
 18 those that are iteratively individually rational, for i is given by $IIR_i = \bigcap_{m \geq 0} \tilde{A}_i^m$. Let
 19 $IIR = \times_{i \in N} IIR_i$.

20 Given these definitions, we can present the first necessary condition of Negotiated Bind-
 21 ing Agreement in n -player games, which states that any proposal, at any history - on and
 22 off the path of play, must survive iterated elimination of individually irrational actions in
 23 the underlying game. This exact process was used in order to find the possible proposals in
 24 the first price auction with heterogeneous values.

25 **Theorem 3.** If s^* is a Negotiated Binding Agreement, then for all $h \in H$, $s_i^*(h) \in IIR_i$.

26 To better understand the set of actions that survives iterated elimination of individually
 27 irrational actions, note the following. In a large class of games, non-emptiness of the set of
 28 actions that are iteratively individually rational is implied by the fact that the set of actions
 29 that survive iterated elimination of never best responses to pure actions, a refinement of
 30 rationalizable strategies as defined by [Bernheim \(1984\)](#); [Pearce \(1984\)](#), also survive iterated

¹⁹The notion of absolute dominance was more recently used by [Doval and Ely \(2020\)](#), who extend this concept to incomplete information.

1 elimination of individually irrational actions. This is formalised in the following definition
2 and lemma.

3 **Definition 5.** Let $a_i \in A_i$ be a never best response to a pure action in $C_{-i} \subseteq A_{-i}$ if, for
4 all $a_{-i} \in C_{-i}$ there is some $a'_i \in A_i$ for which $u_i(a'_i, a_{-i}) > u_i(a_i, a_{-i})$. Denote the set of
5 actions that are never best responses to pure actions in C_{-i} by $NBR_i(C_{-i})$.

6 Let $B_i^0 = A_i$. Let $B_i^k = B_i^{k-1} \setminus NBR_i(A_{-i}^{k-1})$. Let $B^k = \times_{i \in N} B_i^k$ and $B_{-i}^k = \times_{j \neq i} B_j^k$. Let
7 the set of actions that survive iterated elimination of never best responses to pure actions
8 be given by $IENBR = \bigcap_{k \geq 1} B^k$.

9 **Lemma 3.** The set of actions that survive iterated elimination of never best responses to
10 pure actions also survives iterated elimination of iterated deletion of individually irrational
11 actions: $IENBR \subseteq IIR$.

12 Note that the set of actions that survives iterated elimination of never best responses
13 is necessarily non-empty in finite games. Typically even more profiles may survive iterated
14 elimination of individually irrational actions than never best responses to pure actions.
15 To see this, consider the following underlying game.²⁰

16 **Example 2.** Let the underlying game, G , be the following prisoners' dilemma.

1\2	C	D
C	3,3	0,4
D	4,0	1,1

17 D is strictly dominant for both players, hence (D, D) is the only profile that survives
18 survive iterated elimination of never best responses to pure actions. Yet, in IIR , all action
19 profiles survive. This is as the maximum payoff for playing C given by 3. The individually
20 rational payoff is given by 1. Therefore C is not individually irrational. ▼

21 Any action profile satisfying the conditions of the sufficient conditions will be held in
22 IIR , and therefore all pure Nash equilibria must be included.

23 The next result provides further necessary conditions, shows the relation to Negoti-
24 ated Binding Agreement payoffs with individual rationality considerations in the underlying
25 game, when taken over the set of actions that survive iterated elimination of individually
26 irrational actions.

27 **Theorem 4.** if s^* is a Negotiated Binding Agreement then:

$$U_i(s^*|h) \geq \inf_{a'_{-i} \in IIR_{-i}} \sup_{a'_i \in A_i} u_i(a'_i, a'_{-i})$$

28 for all $h \in H$ and $i \in N$.

²⁰For simplicity of exposition, I focus on introducing these definitions and concepts with two-player games.

I illustrate the use of this result with the same underlying prisoner's dilemma game as in example 2.

Example 2. revisited Again consider the underlying game, G , to be that of example 2. In this case, no actions are individually irrational for any player, as previously argued. However, notice that the min-max payoff for each player is 1. The min-max is given by 1, as the worst outcome is the other player selecting D . Therefore we conclude that no Negotiated Binding Agreement can support the action profile (D, C) or (C, D) . However, the necessary conditions do not rule out the possibility of (C, C) . ▼

Note that for any underlying game the inf-sup restricted to the set of actions that survives iterated elimination of individually irrational actions is always weakly higher than the inf-sup without this restriction.

Remark 1. For any underlying game, G , such that \underline{u}_i is well defined the following inequality holds:

$$\inf_{a'_{-i} \in IIR_{-i}} \sup_{a'_i \in A_i} u_i(a'_i, a'_{-i}) \geq \inf_{a_{-i} \in A_{-i}} \sup_{a_i \in A_i} u_i(a_i, a_{-i})$$

Notice this inequality holds strictly within the leading example: the min-max payoff for bidder 1 is 0, via other firms setting bids of 7, however the min-max payoff when we restrict ourselves to IIR is $1/2$.

The results of this section bear resemblance to the analysis of infinitely repeated games, where individual rationality constraints must be satisfied. However, this iterated version can be substantially more restrictive. For instance, in the First Price auction it would only rule out bidders having a net negative valuation, and would not provide a lower bound on the surplus of bidder 1.

Before moving forward, I point to the following corollary, which provides a class of game for which the Negotiated Binding Agreements are fully characterised.

Corollary 1. If a^{NE} is a pure Nash equilibrium of the underlying game G such that:

$$u_i(a^{NE}) = \min_{a_{-i} \in IIR_{-i}} \max_{a_i \in A_i} u_i(a_i, a_{-i})$$

i.e. the IIR min-max profiles are mutual, then a^* can be supported by a Negotiated Binding Agreement if and only if $u_i(a^*) \geq u_i(a^{NE})$.

This is a direct implication of theorems 2 and 4. This provides a class of games for which the Negotiated Binding Agreements are fully characterised by action profiles that Pareto Dominate a Nash equilibrium in the underlying game. Specifically, if a Nash equilibrium provides agents with their individually rational payoffs over the set of actions that survives iterated deletion of individually irrational actions in the underlying game, then an action

1 profile can be supported by a Negotiated Binding Agreement if and only if said action profile
 2 Pareto Dominates this Nash equilibrium of the underlying game. This is the case in the
 3 three bidder first price auction used as a leading example for this section.

4 5 Coalitional Deviations

5 In principle, a negotiation may be susceptible to a collection of agents making binding
 6 agreements over how they will act *within* the negotiation process itself. This is particularly
 7 important given that the negotiation protocol does not generically lead to a unique and
 8 efficient outcome.²¹ To address the concern of susceptibility to groups agreeing to deviate,
 9 I now extend the analysis to allow for this possibility. To do so, I include a collection
 10 of permissible coalitions, where a coalition may jointly deviate. The richest of all such
 11 possibilities is the power set of N , which allows *any* possible subset of players to jointly
 12 deviate.

13 In this analysis, I will look for the most robust form of equilibrium, that prevents
 14 any permissible coalition from deviating, where coalitions are permitted to agree to any
 15 deviation. This can be seen as stronger than necessary, as we may wish for the deviations
 16 to face the same criticism of stability, where these deviations must be the result of some
 17 agreement.²² However, if it were possible to make a binding agreement to not make new
 18 binding agreements, agents may take this option upon deviating. Therefore, in the context
 19 of robust binding agreements, if we do not wish to make assumptions surrounding the game
 20 that is induced to negotiate over when a deviation occurs then this approach ensures no
 21 misspecification. That is, do we allow for agents within a coalition to have veto power? Do
 22 we allow agents to make agreements over what can be within the agreement in the sense
 23 that they pre-commit to rule out some options? This can potentially allow for different
 24 conclusions in the outcome of the negotiation game. Nonetheless, if all deviations of a
 25 coalition are permitted, this includes the outcomes of processes, and therefore if we have
 26 an equilibrium that allows for all possible deviations we certainly have an equilibrium when
 27 all such deviations are not allowed. In this sense, the aim of this analysis differs from the
 28 previous sections in that I will provide easy-to-check conditions for a robust Negotiated

²¹Many negotiation and bargaining protocols do not lead to efficiency. The 2-player Rubinstein (1979) model does not when cost of time is constant, rather than hyperbolic. In the hyperbolic discounting case, when the outside option of this model is taken to be endogenous as in Busch and Wen (1995), a folk theorem is obtained. When there are more than two-players, additional restrictions on the equilibrium notion are needed to regain efficiency (Chatterjee et al., 1993). The work of Harstad (2022) shows that a pledge-and-review bargaining game for contributions to a public good may also lead to inefficient outcomes. Additionally, inefficiencies are common in the contracting literature, for instance in contractable contracts (Tennenholtz, 2004; Kalai et al., 2010; Peters and Szentes, 2012) and strategic contract setting (Jackson and Wilkie, 2005; Yamada, 2003; Ellingsen and Paltseva, 2016).

²²This would renegotiation proofness as in Farrell and Maskin (1989); Bernheim and Ray (1989).

- 1 Binding Agreement, where no specified coalition could deviate, rather than searching for
- 2 all possible Negotiated Binding Agreements that could occur when coalitions can deviate
- 3 in a specific and specified way.

4 5.1. Definitions

5 I first introduce the notation of a coalition and coalition configuration. A coalition configura-
 6 tion defines the set of coalitions that may make a binding agreement within the negotiation.
 7 I let a coalition configuration be denoted by \mathcal{C} , and only restrict \mathcal{C} to be a cover of N . That
 8 is, for all $i \in N$, there is some coalition $C \in \mathcal{C}$ such that $i \in C$. For a coalition configuration
 9 \mathcal{C} , if $C \in \mathcal{C}$ I will refer to C as permissible.

10 Further to this, for a non-empty coalition $C \in \mathcal{C}$, let $a_C = (a_i)_{i \in C}$, $A_C = \times_{i \in C} A_i$,
 11 $s_C = (s_i)_{i \in C}$ and $S_C = \times_{i \in C} S_i$. Let $a_{-C} = (a_i)_{i \notin C}$, $A_{-C} = \times_{i \notin C} A_i$, $s_{-C} = (s_i)_{i \notin C}$
 12 and $S_{-C} = \times_{i \notin C} S_i$. For a set $B \subset A$, which may or may not have a product structure,
 13 let $B_C = \{a_C \in A_C \mid \exists a'_{-C} \in A_{-C} \text{ s.t. } (a_C, a'_{-C}) \in B\}$ and $B_{-C} = \{a_{-C} \in A_{-C} \mid \exists a_C \in$
 14 $A_C \text{ s.t. } (a_C, a_{-C}) \in B\}$.

15 With this, I go on to define the natural extension of Subgame Perfect Equilibrium when
 16 coalitions are permitted to jointly deviate. This will be referred to as \mathcal{C} -Subgame Perfect
 17 Equilibrium and will require that strategies are such that, at no history of the negotiation
 18 game, is there a way for *any* permissible coalition of players, $C \in \mathcal{C}$, to jointly deviate and
 19 improve the utility of all players within that coalition.²³

20 **Definition** (\mathcal{C} -Subgame Perfect Equilibrium). *s^* is a \mathcal{C} -Subgame Perfect Equilibrium if,*
 21 *for all partial histories $h \in H$, there does not exist a non-empty coalition $C \in \mathcal{C}$ and a joint*
 22 *strategy $s_C \in \times_{i \in C} S_i$, such that $u_i(s_C, s_{-C}^* | h) > U_i(s^* | h)$ for all $i \in C$.*

23 This concept generalises a number of solution concepts, which I outline here:

- 24 1. Firstly, whenever $\mathcal{C} = \{\{i\}_{i \in N}\}$, \mathcal{C} -Subgame Perfect Equilibrium and Subgame Perfect
 25 Equilibrium of [Selten \(1965\)](#) coincide. Further to this, whenever $\{\{i\}_{i \in N}\} \subset \mathcal{C}$, \mathcal{C} -
 26 Subgame Perfect Equilibrium is a refinement of Subgame Perfect Equilibrium.
- 27 2. Whenever $\mathcal{C} = 2^N \setminus \{\emptyset\}$, \mathcal{C} -Subgame Perfect Equilibrium coincides with the concept of
 28 strong perfect equilibrium of [Rubinstein \(1980\)](#). Whenever $\mathcal{C} = 2^N \setminus \{\emptyset\}$ I will refer to
 29 this concept as strong in its place. Note that any strong Subgame Perfect Equilibrium
 30 would also be a \mathcal{C} -Subgame Perfect Equilibrium for any \mathcal{C} .

²³In essence, this is assuming that, at any history, any permissible coalition may write a private binding agreement that dictates the behaviour they will take going forward. If the agreements were public, the concept would be closer to a coalitional version of [Tennenholtz \(2004\)](#)'s program equilibrium.

3. Finally, when \mathcal{C} is a partition of N , \mathcal{C} -Subgame Perfect Equilibrium can be seen as the extension of coalitional equilibrium of Ray and Vohra (1997) to extensive form games.

Although specified for any coalition configuration, I will take $\{i\}_{i \in N} \subseteq \mathcal{C}$ as implicit within the discussion, although it is not necessary for the formal results. I will also pay particular attention to the grand coalition being permitted; $N \in \mathcal{C}$.

We can now extend the notion of Negotiated Binding Agreements to \mathcal{C} -Negotiated Binding Agreements, where similarly we will require that we have a \mathcal{C} -Subgame Perfect Equilibrium of the negotiation game and require a no babbling condition. Note that the use of \mathcal{C} -Subgame Perfect Equilibria of the negotiation game when $N \in \mathcal{C}$, gives further justification for no babbling agreements, and indeed no delay agreements. To see this, suppose that there was some $\epsilon > 0$ cost for delay for all agents. If this were the case, then there would be no \mathcal{C} -Subgame Perfect Equilibrium of the negotiation game that ends in more than two periods as a joint deviation could reduce the cost of delay.

Definition 6 (\mathcal{C} -Negotiated Binding Agreement). *s^* is a \mathcal{C} -Negotiated Binding Agreement if:*

1. *s^* is a \mathcal{C} -Subgame Perfect Equilibrium of the negotiation game*

2. *\mathcal{C} -no babbling: $\forall h \in H, \exists h' \in H$ such that $s_C^*(h) = a_C(s^*|h')$.*

a^ is supporting by s^* if $a^* = a(s^*|\emptyset)$.*

When $\mathcal{C} = 2^N \setminus \{\emptyset\}$ I refer to this as a strong Negotiated Binding Agreement.

Whenever $\{i\}_{i \in N} \subset \mathcal{C}$, \mathcal{C} -Negotiated Binding Agreement are a subset of Negotiated Binding Agreement and therefore necessary conditions still hold. However, we can strengthen these conditions, and provide conditions that hold for a general coalition configuration \mathcal{C} . I show that natural extensions of the necessary and sufficient conditions used for Negotiated Binding Agreement hold for \mathcal{C} -Negotiated Binding Agreement.

5.2. \mathcal{C} -Negotiated Binding Agreement Outcomes

5.2.1. Necessary Conditions

First, I will show that in any \mathcal{C} -Negotiated Binding Agreement any action proposed in the negotiation game must survive a procedure of *iterated deletion of coalitionally irrational actions* on the underlying game. This procedure works inductively as follows. Consider some joint action of those within a coalition $C \in \mathcal{C}$ in the underlying game, a_C . If, for a

1 coalition $C \in \mathcal{C}$ there is some function, that maps the joint action of those outside of the
 2 coalition to a joint action of the coalition, which, even in the worst case said function can
 3 provide a higher payoff than the joint action a_C , then a_C is a coalitionally irrational joint
 4 action. This generalises the notion of individual rationality.²⁴ Notice this is exactly the
 5 notion of [Aumann \(1961\)](#)'s β -core. We may proceed inductively. Remove all coalitionally
 6 irrational actions for all coalitions $C \in \mathcal{C}$ in the underlying game. Consider some joint action
 7 of those within a coalition $C \in \mathcal{C}$, a_C , which survives iterated elimination of coalitionally
 8 irrational actions up to some iteration k and so on. This provides a recursive version of
 9 [Aumann \(1961\)](#)'s β -core, where the "punishments" themselves must be justified. This,
 10 therefore, provides one answer to the question posed by [Scarf \(1971\)](#), providing a notion of
 11 the core for normal form games that is fully justified.²⁵

12 **Definition 7.** For any underlying game G , for a coalition C , a joint action $a_C \in A_C$ is
 13 coalitionally irrational with respect to $B_{-C} \subseteq A_{-C}$ if, for some $a'_C : B_{-C} \rightarrow A_C$:

$$\inf_{a_{-C} \in B_{-C}} u_i(a'_C(a_{-C}), a_{-C}) > \sup_{a_{-C} \in B_{-C}} u_i(a_C, a_{-C}) \quad \forall i \in C$$

14 Denote the set of joint actions that are coalitionally irrational with respect to B_{-C} by
 15 $D_C(B_{-C})$.

16 **Definition 8** (Iterated Elimination of Coalitionally Irrationality actions with respect to
 17 \mathcal{C}). For any game G , let $\tilde{A}^0(\mathcal{C}) = A$. For $m > 0$ let:

$$\tilde{A}^m(\mathcal{C}) = \tilde{A}^{m-1}(\mathcal{C}) \setminus \left[\bigcup_{C \in \mathcal{C}} [D_C(\tilde{A}^{m-1}(\mathcal{C})_{-C})] \times A_{-C} \right]$$

18 Let the set of action profiles that survive iterated elimination of coalitionally irrational
 19 actions, or those that are iteratively coalitionally rational, with respect to \mathcal{C} be denoted by
 20 $ICIR(\mathcal{C})$ where $ICIR(\mathcal{C}) = \bigcap_{m>0} \tilde{A}^m(\mathcal{C})$.

21 Note, unlike iterated elimination of individually irrational actions, iterated elimination
 22 of coalitionally irrational actions may be empty, even in finite games. To see this, consider
 23 the following example.

24 **Example 3.** Consider the following 2 player game to be the underlying game G . Let
 25 $\mathcal{C} = \{\{1, 2\}, \{1\}, \{2\}\}$.

²⁴For underling games with compact action spaces and continuous utility they are identical when $\mathcal{C} = \{\{i\}_{i \in N}\}$.

²⁵[Chakrabarti \(1988\)](#) offers a different solution to this question by taking the punishments to be such that they cannot be coalitionally dominated for any action, i.e. extending the typical notion of a dominated action with the caveat that coalitions are the unit of decision making, and shows that this can be connected to an refinement of the strong equilibria of an infinitely repeated game using the limit of means criteria.

1\2	L	C	R
T	20,0	20,0	20,0
M	0,7.5	0,7.5	30,5
D	10,10	0,0	0,0

1 Notice that only (M, R) and (D, L) survive iterated elimination of coalitionally irrational
2 actions for the coalition $C = \{1, 2\}$. However, D cannot survive elimination of individually
3 irrational actions for player 1, as the maximum payoff of D is 10 while the min-max utility
4 for player 1 is 20. Therefore we conclude that within the first round of iterated elimination
5 of coalitionally irrational actions only (M, R) survives. However, this implies that R is
6 individually irrational with respect to M for player 2, as the profile (M, R) gives a payoff
7 of 5 while the min-max utility, when restricting attention to player 1 playing R is 7.5.
8 Therefore $ICIR(\mathcal{C}) = \emptyset$. ▼

9 However, it may be non-empty, even when a rich set of coalitions are permitted. Before
10 doing so, notice the following. If $\mathcal{C}' \subset \mathcal{C}$, then $ICIR(\mathcal{C}) \subseteq ICIR(\mathcal{C}')$. Given this, if some
11 action profile survives $ICIR(2^N \setminus \{\emptyset\})$ then it survives any other \mathcal{C} .

12 **Example 4.** Consider the following 2 player game as the underlying game, G . Let
13 $\mathcal{C} = \{\{1, 2\}, \{1\}, \{2\}\}$.

1\2	L	C	R
T	2,7	2,8	0,6
M	1,4	0,8	2,3
D	1,9	0,8	20,7.5

14 Notice that (D, R) , and (D, L) and (T, C) are the set of Pareto efficient outcomes,
15 therefore, as $\{1, 2\} \in \mathcal{C}$, it must be all other action profiles are rules out in $\tilde{A}^1(\mathcal{C})$. Further,
16 R is individually irrational for 2 as it provides a payoff of at most 7.5, while the min-max
17 payoff is 8. We conclude that $\tilde{A}^1(\mathcal{C}) = \{(D, L), (T, C)\}$. Now notice that D is individually
18 irrational for 1 with respect to \tilde{A}_{-1}^1 , where $\tilde{A}_{-1}^1 = \{L, C\}$, as the highest payoff that D can
19 provide is 1 while the min-max payoff over this set is 2. We conclude that $\tilde{A}^2(\mathcal{C}) = \{(T, C)\}$.
20 Finally, note that neither T or C are individually irrational given $B_{-1} = \{C\}$ and $B_{-2} = \{T\}$
21 respectively. Therefore $ICIR(\mathcal{C}) = \{(T, C)\}$. ▼

22 One condition that ensures non-emptiness of $ICIR(\mathcal{C})$, regardless of the coalition con-
23 figuration, is the existence of a strong Nash equilibrium.²⁶

24 **Lemma 4.** For any Strong Nash equilibrium a^{SNE} of G , $a^{SNE} \in ICIR(\mathcal{C})$ regardless of \mathcal{C} .

²⁶Recall a strong Nash equilibrium is an action profile a^{SNE} such that for all $C \in 2^N \setminus \{\emptyset\}$ $\nexists a_C \in A_C$ such that $u_i(a_C, a_{-C}^{SNE}) > u_i(a^{SNE})$ for all $i \in C$.

A similar necessary condition to theorem 3 holds, linking $ICIR(\mathcal{C})$ of the underlying game to the proposals made in \mathcal{C} -Negotiated Binding Agreement of the negotiation game.

Theorem 5. *For any \mathcal{C} -Negotiated Binding Agreement, s^* , and any $h \in H$, $s^*(h) \in ICIR(\mathcal{C})$.*

Notice once again that this holds for all histories. Further to this, by the definition of $ICIR(\mathcal{C})$, whenever $N \in \mathcal{C}$, it follows that no proposal is coalitionally irrational for the coalition N . This implies that only proposals that are weakly Pareto optimal in the underlying game may be used.

The following corollary links the observation surrounding the potential emptiness of $ICIR(\mathcal{C})$ of the underlying game to the emptiness of \mathcal{C} -Negotiated Binding Agreement.

Corollary 2. *If $ICIR(\mathcal{C}) = \emptyset$ then no \mathcal{C} -Negotiated Binding Agreement can exist.*

This is an immediate implication of theorem 5. Note that this is possible, i.e. in example 3, and may imply that there is no Negotiated Binding Agreement that is robust to the concerns of coalitions for a specific coalition structure \mathcal{C} .

A result analogous to theorem 4 also holds. This result will state that at any history h , a \mathcal{C} -Negotiated Binding Agreement must give a payoff that is coalitionally rational for any coalition C in the underlying game, with respect to $[ICIR(\mathcal{C})]_{-C}$. A payoff is not coalitionally rational, with respect to $[ICIR(\mathcal{C})]_{-C}$, if, for any punishment a coalition can find some joint action $a_C \in A_C$ such that the utility is higher for all agents. To understand the implications of this result more fully, I define a notion of the β -core Aumann (1961), which I refer to as the β -core with respect to $ICIR(\mathcal{C})$.

Definition 9. *$a^* \in A$ is in the β -core with respect to $ICIR(\mathcal{C})$ if, there is no $C \in \mathcal{C}$ and $a_C : [ICIR(\mathcal{C})]_{-C} \rightarrow A_C$ such that $\inf_{a_{-C} \in [ICIR(\mathcal{C})]_{-C}} u_i(a'_C(a_{-C}), a_{-C}) > u_i(a^*)$ for all $i \in C$.*

For an action profile to be in the β -core the payoff of this profile must be higher than the coalitional rational with respect to A_{-i} , in the sense that a coalition understands that they can only be punished for a deviation with a specific profile of actions. However, the actions used to prevent deviations are not necessarily justifiable. The β -core with respect to $ICIR(\mathcal{C})$ partially resolves this problem, as upon deviating the actions of others are restricted to a set of actions that is consistent with respect to itself and is defined in a similar way to the β -core restriction itself.

With this, I formalise the result connecting \mathcal{C} -Negotiated Binding Agreement to the β -core with respect to $ICIR(\mathcal{C})$.

Theorem 6. For any \mathcal{C} -Negotiated Binding Agreement s^* must be such that, for any history h , and for any coalition $C \in \mathcal{C}$, there is no $a'_C : [ICIR(\mathcal{C})]_{-C} \rightarrow A_C$ such that:

$$\inf_{a_{-C} \in [ICIR(\mathcal{C})]_{-C}} u_i(a'_C(a_{-C}), a_{-C}) > U_i(s^*|h)$$

for all $i \in C$.

In other words, $a(s^*|h)$ must be in the β -core with respect to $ICIR(\mathcal{C})$ for all histories.

Note that it may be that an outcome is both Pareto efficient and individually rational in the underlying game, yet it is not possible to sustain such an outcome via a \mathcal{C} -Negotiated Binding Agreement for $\{N, \{i\}_{i \in N}\} \subseteq \mathcal{C}$.

Example 5. Let the following two-player game be the underlying game G . Consider the richest set of coalitions $\mathcal{C} = \{\{1\}, \{2\}, \{1, 2\}\} = 2^N \setminus \{\emptyset\}$.

1\2	LL	L	R	RR
TT	6,6	0,4	1,12	0,0
T	4,0	0,0	7,2	<u>1,1</u>
D	12,1	2,7	4,4	0, <u>8</u>
DD	0,0	1, <u>1</u>	<u>8</u> ,0	0,0

I have labelled the weakly Pareto efficient outcomes of G in bold blue font, and therefore must be the only actions in \tilde{A}^1 are $\{(TT, LL), (TT, R), (T, R), (D, LL), (D, L)\}$. No further deletion can take place therefore:

$$ICIR(2^N \setminus \{\emptyset\}) = \{(TT, LL), (TT, R), (T, R), (D, LL), (D, L)\}$$

(TT, R) necessarily cannot be sustained in a strong Negotiated Binding Agreement, as it provides a payoff of 1, while the min-max payoff, given that player 2 must choose from $[ICIR(2^N \setminus \{\emptyset\})]_2 = \{LL, L, R\}$, is given by 2. Therefore we conclude that despite the fact that (TT, R) is Pareto efficient in G , and provides a higher payoff than the min-max over all possible profiles it cannot be sustained in a strong Negotiated Binding Agreement. ▼

With these results, I now turn to providing sufficient conditions for \mathcal{C} -Negotiated Binding Agreement.

5.2.2. Sufficient Conditions

To provide sufficient conditions for the outcomes of a \mathcal{C} -Negotiated Binding Agreement, as with theorem 2, I will rely on conditions of the underlying game G . To provide these conditions, I again rely on a structure that does not focus on the deviation that a coalition takes, but only on the deviating coalition. In this case, a coalition must prefer the punishment

1 of others to their own and a coalition must not be able to improve all members' utility
 2 by changing their action profile in G , holding the punishment used against them constant.
 3 Note, due to the rich deletion that can take place, the inclusion of such profiles in $ICIR(\mathcal{C})$
 4 is now required and not implied.

5 **Theorem 7.** *Take any underlying game such that there is some $a^* = \underline{a}^N \in ICIR(\mathcal{C})$ and*
 6 *for all $C \in \mathcal{C} \setminus N \exists \underline{a}^C \in ICIR(\mathcal{C})$ such that:*

- 7 1. $\nexists a'_C \in A_C$ such that $u_i(a'_C, \underline{a}_{-C}^C) > u_i(\underline{a}^C)$ for all $i \in C$
- 8 2. for all $C \in \mathcal{C}$ there is some $i \in C$ such that $u_i(a^*) \geq u_i(\underline{a}^C)$
- 9 3. For all $C, C' \in \mathcal{C}$ there is some $i \in C$ such that $u_i(\underline{a}^{C'}) \geq u_i(\underline{a}^C)$

10 *Then a^* can be supported in a \mathcal{C} -Negotiated Binding Agreement.*

11 Combining this result with the result of lemma 4, which states that if a strong Nash
 12 equilibrium of G exists it is within $ICIR(\mathcal{C})$, implies that any strong Nash equilibrium
 13 of G can be supported in a \mathcal{C} -Negotiated Binding Agreement. However, these conditions
 14 can apply in underlying games with no strong Nash equilibrium, and therefore are a more
 15 general set of conditions.²⁷ To see this, consider the following example.

16 **Example 5. revisited** Consider again the following two-player game as the underlying
 17 game, G , given in example 5. All possible coalitions are permitted, $\mathcal{C} = 2^N \setminus \{\emptyset\}$.

18 Here there is no strong Nash equilibrium of G . In fact, as there is no pure Nash
 19 equilibrium in G , there is no pure coalition proof Nash equilibrium. However, the conditions
 20 of theorem 7 apply. Given the previous analysis we may take $\underline{a}^N = a^* = (TT, LL)$, $\underline{a}^1 =$
 21 (D, L) and $\underline{a}^2 = (T, R)$. Concluding that (TT, LL) can be sustained in $2^N \setminus \{\emptyset\}$ -Negotiated
 22 Binding Agreement. ▼

23 The sufficient conditions for outcomes of \mathcal{C} -Negotiated Binding Agreements presented
 24 in theorem 7 can be seen as a further refinement of the β -core of Aumann (1961), where
 25 within the β -core any constant action profile in G of those outside of a coalition may be used
 26 in order to prevent deviations, whereas in this paper we must satisfy additional conditions
 27 to ensure such a profile in G can be mutually justified by all coalitions. Note that this is
 28 not necessarily true in the notion of the β -core with respect to $ICIR(\mathcal{C})$, as some profiles
 29 within $ICIR(\mathcal{C})$ do not satisfy this notion of mutual coalitional rationality.

²⁷Shubik (2012) examines the 78 2x2 games which can be induced by strict ordinal preferences, of these 78,
 67 allow for the sufficient conditions for outcomes of a \mathcal{C} -Negotiated Binding Agreement to be applied. Note
 that is only 2 less than the existence of Nash equilibrium in pure strategies. In this sense, these sufficient
 conditions apply to more scenarios than initial inspection may suggest.

1 If there is only a single action profile consistent with $ICIR(\mathcal{C})$ then this must be sup-
 2 ported by a \mathcal{C} -Negotiated Binding Agreement, and further to this is the only profile that
 3 can be the outcome of \mathcal{C} -Negotiated Binding Agreement.

4 **Corollary 3.** *If G is such that u_i is continuous and A_i is compact for all agents, if*
 5 *$ICIR(\mathcal{C}) = \{a^*\}$, then a^* , then s^* is a \mathcal{C} -Negotiated Binding Agreement if and only if*
 6 *$s^*(h) = a^*$ for all $h \in H$.*

7 Note that uniqueness can occur more often than when $\{i\}_{i \in N} \subset \mathcal{C}$, as $ICIR(\mathcal{C})$ may
 8 involve more deletion in the underlying game G . However, as $ICIR(\mathcal{C})$ may be empty and
 9 leave us with no \mathcal{C} -Negotiated Binding Agreement.

10 I now turn to an application.

11 5.2.3. An Application of \mathcal{C} -Negotiated Binding Agreements

12 As with strong Nash equilibrium, conditions for existence of a \mathcal{C} -Negotiated Binding Agree-
 13 ment are not generically satisfied. Nonetheless, there exist interesting applications for which
 14 \mathcal{C} -Negotiated Binding Agreements exist. Consider the following Cournot game.

15 Application 1. (Symmetric Cournot with Fixed Cost)

16 Consider a simple model of Cournot with fixed costs as the underlying game G . I will
 17 take these fixed costs to depend on the total number of firms that enter the market. This
 18 captures a situation where the fixed cost is due to the purchase of equipment. The cost of
 19 the equipment itself is dictated by the law of supply and demand and therefore this cost
 20 increases with the number of firms purchasing this.

21 I model this in the following way. Let there be $n = 4$ firms. Let each firm choose the
 22 quantity that they will sell, $q_i \geq 0$. Let inverse demand, as a function of the total quantity,
 23 be given by $\max\{b - \sum_{j=1}^4 q_j, 0\}$, where $b > 0$. I assume that the marginal cost is constant
 24 and symmetric, therefore it is without loss to set it to 0. Therefore gross profits for player
 25 $i \in \{1, 2, 3, 4\}$ are given by $\max\{(b - \sum_{j=1}^n q_j), 0\}q_i$. Let fixed costs take the following
 26 form: $\left(\frac{3}{32}b \sum_{j \neq i} \mathbf{1}_{q_j > 0}\right)^2 \mathbf{1}_{q_i > 0}$. Notice that this cost increases with the number of firms
 27 entering, and the first firm to enter the market may do so for free. Therefore utility takes
 28 the following form:

$$u_i(q) = \max \left\{ \left(b - \sum_{j=1}^4 q_j \right), 0 \right\} q_i - \left(\frac{3}{32} b \sum_{j \neq i} \mathbf{1}_{q_j > 0} \right)^2 \mathbf{1}_{q_i > 0}$$

29 Notice the individual best responses are given by:

$$q_i^*(q_{-i}) = \begin{cases} \left\{ \frac{b - \sum_{j \neq i} q_j}{2} \right\} & \text{if } \sum_{j \neq i} q_j < b - \frac{3}{16}b \sum_{j \neq i} \mathbf{1}_{q_j > 0} \\ \left\{ 0, \frac{b - \sum_{j \neq i} q_j}{2} \right\} & \text{if } \sum_{j \neq i} q_j = b - \frac{3}{16}b \sum_{j \neq i} \mathbf{1}_{q_j > 0} \\ \{0\} & \text{if } \sum_{j \neq i} q_j > b - \frac{3}{16}b \sum_{j \neq i} \mathbf{1}_{q_j > 0} \end{cases}$$

1 All Nash equilibria of this game involve 2 firms enters the market with quantities of
 2 $q_i^* = \frac{b}{3}$. This leads to a payoff of $\frac{943}{9216}b^2$ for the firm who enters and 0 for those who do not.
 3 Further, note that this is not a Strong Nash equilibrium but all such equilibria are coalition-
 4 proof Nash equilibria. There are many Pareto efficient outcomes. For example, when sellers
 5 in aggregate sell the monopoly quantity, $\sum_{i=1}^4 q_i = \frac{b}{2}$, while profits are strictly positive for
 6 all those who produce strictly positive quantities, is Pareto efficient. Note that there exists
 7 such a profile for any number of firms entering. For instance, all firms producing $\frac{b}{8}$ leads
 8 to profits of $u_i(q^{p,all}) = \frac{55}{1024}b^2 > 0$. Due to the fixed costs, the weakly Pareto efficient
 9 outcomes do not require the monopoly quantity unless all firms enter, but all involve all
 10 firms receiving weakly positive payoff. Any weakly Pareto efficient profile is in the α -core
 11 and the β -core.

12 Consider $\mathcal{C} = 2^N \setminus \{\emptyset\}$. Let $q^* = (q_1^*, q_2^*, q_3^*, q_4^*)$ be the quantity that is trying to be
 13 sustained. I will argue that it is possible to sustain an efficient outcome where all agents
 14 produce in strong Negotiated Binding Agreement. That is an agreement such that $q_i^* = \frac{b}{8}$
 15 can be supported. Consider the following strategies.

16 1. [Punishment for a deviation of coalition C] If $h = (q^1, q^2, \dots, q^k)$ is such that $q_{-C}^{k-1} =$
 17 $s_{-C}^*((q^1, q^2, \dots, q^{k-2}))$ and either

- 18 (a) $q_l^k = s_l^*(q^1, q^2, \dots, q^{k-1})$ for all $l \notin C$ and $q_j^k \neq s_j^*(q^1, q^2, \dots, q^{k-1})$ for all $j \in C$
 19 (b) or $q_{-C}^k = \frac{16+\sqrt{137}}{64}b$ if $|C| = 3$, $q_{-C}^k = \left(\frac{16+\sqrt{137}}{64}b, \frac{16+\sqrt{137}}{64}b \right)$ if $|C| = 2$ and
 20 $q_{-C}^k = \left(\frac{16+\sqrt{137}}{64}b, \frac{16+\sqrt{137}}{64}b, 0 \right)$ if $|C| = 1$

21 Then:

- 22 • $s_i^*(h) = \frac{16+\sqrt{137}}{64}b$ for $i = \min_{j \notin C} j$ if $|C| \leq 3$ or $i = 1$ if $|C| = 4$.
 23 • $s_i^*(h) = \frac{16+\sqrt{137}}{64}b$ for $i = \min_{j \notin C \setminus \{\min_{j \notin C} j\}} j$ if $|C| \leq 2$ or $i = \text{mod}(j+2, 4)$,
 24 $j \in C, j \geq k, k \in C$ otherwise.
 25 • $s_i^*(h) = 0$ for all other $i \in N$.

26 2. [No deviation / deviation of N] For all other histories, let $s_i^*(h) = \frac{b}{8}$

27 The logic of this strategy is as follows. Suppose we are at a history only one coalition
 28 has deviated in the penultimate period of the history, while in the period before that either

all firms have made proposals in accordance with their assigned strategy or only those within the deviating coalition have deviated. Note that this may involve a smaller coalition deviating in the penultimate period, while in the next a larger coalition deviates. If this is the case, assign proposals such that two firms each receive exactly half the aggregate payoff of all firms entering and producing the efficient quantity. All other firms produce 0. At least one of these firms is not within the deviating coalition if the cardinality of that coalition is 3 or less. At all other histories, all agents propose their share of the equal division of the monopoly quantity.

Now I will show that this does indeed constitute a strong Negotiated Binding Agreement. First consider a coalition deviating from a history that does not fall into case 2, where no deviation leads to the agreement that all firms enter and divide the monopoly quantity. It cannot be that the grand coalition deviates to improve the utility of all members. Therefore it must be that deviation does not involve one firm. By the structure of case 1, which any deviation must lead to, it is then the case that those outside the coalition are proposing, in aggregate, at least $\frac{16+\sqrt{137}}{64}b$ in every period. As $\frac{16+\sqrt{137}}{64}b > \frac{1}{8}b$, there is less total demand left for the three deviating firms, and therefore it cannot be that all firms who deviate are producing and improving the utility of all members.²⁸ Therefore it must be that the deviation only involves a coalition of at most two firms. It cannot be that they are both assigned to not produce in all periods, as this implies that the profits are bounded above by $\frac{249-32\sqrt{137}}{4096}b < 0$ if producing and 0 if not. This bound is the same if only one firm deviates. Therefore no profitable deviation can exist from case 2. Now suppose that a profitable deviation exists from case 1. Similarly, it cannot be that two firms deviate and improve their utility. This is because a “punishing” firm does not wish to deviate, as they would be punished for this. A “punished” firm also does not wish to deviate, as the punishment is sufficiently high to ensure that they do not wish to enter the market. It cannot be that all agents jointly deviate, as there are two punishing firms, who, in aggregate, receive the utility that is the maximum that can be achieved for all firms entering. Therefore they have no incentive to do so.

Notice that the punishments can be tailored such that any punished firms do not produce, leading to a profit of 0. Given for any deviating coalition such a punishment could be used, we conclude that any weakly Pareto efficient profile could be sustained, that is, only the conditions of efficiency for the grand coalition need to be respected. This leads to an equivalence to the β -core, while being fully justified via the use of \mathcal{C} -Negotiated Binding

²⁸Notice that this could be strengthened by assigning the one outside of the deviating coalition a quantity that would provide them with the aggregate monopoly profit when all enter: that is they receive a profit of $\frac{220}{1024}b^2$, specifically by this firm producing $\frac{44}{64}b$ and others producing 0. However, for the sake of a simpler strategy, I use it as is. As $\frac{44}{64}b > b - \frac{3}{16}b \times 3$ it follows that those three deviating firms cannot make a positive profit.

1 Agreements. ▼

2 6 Literature Review

3 A number of papers have approached the question of binding agreements that can be made
 4 for normal form games using an approach close to or inspired by the farsighted stable set of
 5 [Harsanyi \(1974\)](#). I instead take a more non-cooperative game theoretic approach, exploring
 6 a refinement of SPE in a fully specified negotiation game. Within this strand of literature,
 7 [Mariotti \(1997\)](#) has the closest model and also considers an explicit negotiation protocol.
 8 The extensive form of the negotiation protocol is similar, but the payoff of perpetual dis-
 9 agreement is set to $-\infty$. In this work, [Mariotti \(1997\)](#) takes an approach close to the
 10 strong Subgame Perfect Equilibrium of [Rubinstein \(1980\)](#). He also imposes a refinement
 11 on this subgame perfect type concept based on the farsighted stable set. [Mariotti \(1997\)](#)
 12 does not provide general conditions for his solution concept, due to the complexity that the
 13 history-dependent negotiation entails. He instead proposes a history-independent version
 14 of his solution concept, in line with [Harsanyi \(1974\)](#), where agents strategies only map from
 15 the current proposal to the next proposal, rather than all possible previous proposals being
 16 considered. In this history independent version, [Mariotti \(1997\)](#) provides some necessary
 17 conditions for agreement outcomes similar to those provided in this paper for both Negoti-
 18 ated Binding Agreements and \mathcal{C} -Negotiated Binding Agreements. He also provides sufficient
 19 conditions for agreement outcomes for a class of two-player games with conditions on the
 20 Pareto Frontier, similarly using a notion of individual punishments.

21 [Chwe \(1994\)](#); [Xue \(1998\)](#); [Ray and Vohra \(2015, 2019\)](#) also consider versions of the
 22 farsighted stable set. The closest with respect to my paper is [Ray and Vohra \(2019\)](#), which
 23 games with transferable utility, and defines the notion of the maximal farsighted stable set,
 24 which additionally requires a subgame perfect-like condition, imposing optimality given
 25 others' strategies at all histories of the negotiation. They provide general conditions linking
 26 the farsighted stable set as defined in [Ray and Vohra \(2015\)](#) to this concept. I instead take
 27 an approach that looks at general games, rather than a game with transferable utility, and
 28 instead link the concept of \mathcal{C} -Negotiated Binding Agreements to an alternative cooperative
 29 game theoretic Notion of the β -core of [Aumann \(1959, 1961\)](#). Finding the farsighted stable
 30 set is challenging and some papers have looked at finding the farsighted stable set for a
 31 specific underlying game ([Suzuki and Muto, 2005](#); [Nakanishi, 2009](#)).

32 Other papers have also proposed fully non-cooperative models of negotiation over bind-
 33 ing agreements for normal form games, based on a dynamic game of negotiation. [Kalai](#)
 34 [\(1981\)](#) looks at a fully specified model of negotiation by proposing a non-cooperative ex-
 35 tensive form game. In that model, agents propose an individual action in the underlying

1 game. If an agent changes their proposal within a period then they are no longer permitted
 2 to change their proposal again. The process ends at time t with the proposal profile pro-
 3 posed in that period. Kalai (1981) looks at the perfect equilibria of Selten (1988) and shows
 4 that only cooperation can be sustained in the 2-player prisoners' dilemma game. Nishihara
 5 (2022) has extended this to an n -player prisoners' dilemma, maintaining Kalai's negotiation
 6 protocol. The philosophy of Kalai's approach is similar to that of this paper, where agents
 7 negotiate over the agreement and can do so by proposing their own action. Bhaskar (1989)
 8 examines a model of pre-play agreement over a symmetric two-player Bertrand game. In a
 9 similar sense to this model, agents make proposals of the prices they will take, and have the
 10 opportunity to revise their proposals sequentially. Confirmation requires one agent not to
 11 change their proposal after seeing the others. Bhaskar (1989) looks at the perfect equilibria
 12 of such an agreement game and concludes that only the monopoly price can be sustained.
 13 The closest model in the non-cooperative literature is that of Harstad (2022), who pro-
 14 poses a "pledge-and-review" bargaining protocol, similar to the one in this paper, for public
 15 goods games. In his model, Harstad (2022) shows that when agents confirm by default,
 16 and discounting is hyperbolic, a folk theorem remains for the subgame perfect equilibria
 17 outcomes of this game. When considering typical refinements and variations of subgame
 18 perfect equilibrium (stationary subgame perfection and local perfection / trembling hand),
 19 each agent's pledge must be the result of maximising *some* weighted Nash product. How-
 20 ever, the weights used by each agent may differ, and therefore many inefficient equilibria
 21 arise despite this. In my work, I instead consider a more general class of games and consider
 22 and alternative refinement of SPE.

23 A number of papers have provided a more cooperative game theoretic approach for the
 24 agreements that can be made for games, for instance Strong Nash equilibrium (Aumann,
 25 1959) and the β -core (Aumann, 1959, 1961). In my paper, \mathcal{C} -Negotiated Binding Agreement
 26 outcomes lie somewhere between the β -core and Strong Nash equilibrium, as agents are
 27 permitted to change their proposals when they observe a proposal of others change, but can
 28 only do so in a way pinned down by an optimal strategy in the sense of equilibrium. Given
 29 this, my paper can also be seen in the light of the Nash program pointed to in Nash (1953),
 30 as the necessary and sufficient conditions \mathcal{C} -Negotiated Binding Agreement outcomes can
 31 be seen as a perturbed version of the β -core.

32 There are a number of other related papers that take the cooperative game theoretic
 33 approach. Notably, the γ -core (Chander and Tulkens, 1997). Chander (2007) provides
 34 further justification for the γ -core by showing it is *an* equilibrium to an infinitely repeated
 35 game where agents decide whether to cooperate or not in each round. Chander and Wooders
 36 (2020) define a notion of coalitional Subgame Perfect Equilibrium for underlying games with
 37 transferable utility, where a coalition's deviation payoff is with respect to the best Subgame

Perfect Equilibrium assuming all other players act without cooperation. A number of papers have also proposed notions of rationalizability for coalitions in a cooperative sense, for instance Herings et al. (2004); Ambrus (2006, 2009); Grandjean et al. (2017), which iterative elimination of coalitionally irrational actions can be seen as, but are all distinct. A strand of literature abstracts from the negotiation process *within* a group and takes a cooperative perspective, focusing on Pareto undominated actions that prevent new groups from breaking and forming (Ray and Vohra, 1997; Diamantoudi and Xue, 2007).

A number of papers consider a form of communication for equilibrium selection (Bernheim et al. (1987); Farrell and Maskin (1989); Bernheim and Ray (1989); Rabin (1994), etc.). My paper is related in the sense that agents can communicate via the negotiation procedure to select the outcome of the underlying game that will be played. However, the perspective is different, as these concepts are about refining a given set of non-binding agreements represented by the (potential mix over) SPE or Nash Equilibria of an underlying game, whereas I allow agents to make a binding agreement of potentially any outcome.

Negotiated Binding Agreements is also related to another literature on binding agreements, contract theory. Most closely related are the works of Jackson and Wilkie (2005); Yamada (2003); Ellingsen and Paltseva (2016) who all propose model allowing agents all have a strategic input on the *structure* of the contract over an underlying strategic environment, rather than allowing one agent or a mediator to completely define the structure and then make a take-it-or-leave-it offer on the contract. In a similar way, Negotiated Binding Agreements allows for all agents to have a strategic input on the action they will agree to in the underlying game. On the other hand, Kalai et al. (2010), Peters and Szentes (2012) and Tennenholtz (2004) all consider the possibility of all agents proposing contracts surrounding their own play in an underlying game, where these contracts can be a function of the contracts of others. This allows agents to specify reactions to deviations in full, and can allow for these to be fully specified at a higher level also. In contrast to these, my paper is requires that agents are required to only propose actions they could agree to, whereas these papers allow contracts to specify actions in the contract that would never be the result of equilibrium.

The way payoffs are defined for perpetual disagreement can be seen as similar to the literature of infinitely repeated games with no discounting. When well defined, the limit of means criteria of Aumann and Shapley (1994); Rubinstein (1994) can be used. The sufficient conditions within the paper are also similar to the sufficient conditions of player-specific punishment is used in infinitely repeated games (Fudenberg and Maskin, 1986; Abreu et al., 1994). The sufficient conditions I use are more restrictive as player-specific punishment only requires that their punishments' provide them an individually rational payoff and they prefer to punish rather than be punished. In contrast, I also require that

1 individuals are best responding to their punishment in the underlying game. These are
 2 used as there are no further rewards from following their punishments, which are held
 3 in the continuation of an infinitely repeated game. Therefore it must be the case that
 4 agents cannot improve the utility they would get facing the constant punishment of others,
 5 requiring that they best respond.

6 7 Conclusion

7 I propose a model of negotiated binding agreements over agents' play in an underlying
 8 normal form game. I study the outcomes of the underlying game that be supported using
 9 a refinement of Subgame Perfect Equilibrium, where agents only propose actions that they
 10 could agree to. I refer to this concept as Negotiated Binding Agreements. I show that
 11 the outcomes of the underlying game that can be agreed upon must satisfy a condition
 12 of *iterative* individual rationality. Further, any outcome in the underlying game, where
 13 appropriate individual punishments can be found, can be agreed to. These individual
 14 punishments are defined on the underlying game, where agents must be prescribed the
 15 action that best responds to their punishment in the baseline game. The sufficient condition
 16 for agreement outcomes is also shown to be necessary for two-player games, leading to a
 17 full characterisation within this class. By providing conditions for outcomes that can be
 18 agreed upon that are solely based on characteristics of the underlying game, I reconcile the
 19 rigour of the solution of a fully specified model of negotiation with easy-to-use conditions
 20 for agreement outcomes for the underlying game.

21 To display the ease of use of these conditions, I explore two key applications.²⁹ In a
 22 Cournot Duopoly, I show that when marginal costs are the same, any profile of payoffs such
 23 that each player receives positive profits is sustainable. In contrast, when marginal costs are
 24 very different only the firm with the lowest marginal cost receiving their monopoly profit is
 25 supported. In a simple First Price Auction, I show that these conditions lead to intuitive
 26 results about what can be agreed upon, where a minimal bound is put on the payoff the
 27 highest valuation bidder. In these applications, I fully characterise the Negotiated Binding
 28 Agreement outcomes.

29 I show how the necessary and sufficient conditions for the outcomes of the Negotiated
 30 Binding Agreements within this negotiation game naturally generalise to the case where
 31 agents may agree upon *how* to negotiate. I show these generalised conditions are linked
 32 to a perturbed version of the cooperative game theoretic notion of the β -core of [Aumann](#)
 33 (1961), while having the full backing of a fully specified negotiation procedure. I apply this
 34 to a Cournot model, with a fixed cost that depends on the number of entrants. Within this

²⁹An additional application of a public goods game is provided in the online appendix.

1 setting, the outcomes coincide with the β -core, but are fully justified by a fully specified
2 negotiation protocol.

3 A number of questions remain open. Firstly, there are a number of applied theory
4 questions that can use the results of this paper. A number of applied theory papers have
5 made use of cooperative solutions, for example in environmental agreements (Chander and
6 Tulkens, 1997; Carraro, 1998; Carraro et al., 2006) and trade agreements (Aghion et al.,
7 2007; Conconi and Perroni, 2002). Due to the easy-to-use conditions, my results may also
8 provide some interesting insights in some applied theoretical settings, while having the
9 backing of a fully specified negotiation protocol.

10 Additionally, the results of this paper may shed light on which environments should
11 be negotiated jointly, that is, when is bundling issues or games in negotiation beneficial.³⁰
12 This is particularly interesting from the applied theory perspective. For instance, inter-
13 national trade agreements involve simultaneously negotiating tariffs for multiple markets
14 and, for instance, environmental policy.³¹ However this is not always the case and therefore
15 understanding when it is theoretically beneficial is an interesting line to follow.³² Further,
16 the results of this paper may provide an understanding of when there is a benefit from
17 unilaterally giving up some actions in the underlying game, essentially allowing agents to
18 “take chips off the table”. The results of this paper show that unilaterally giving up an
19 action *can* be beneficial.³³ Nonetheless, understanding the removal of *which* actions leads
20 to this improvement is an open question. I leave these questions for future work.

21 Finally, allowing coalitions to overlap provides a direction for interesting insights. Such
22 arrangements of groups frequently occur in economic environments, such as international
23 relations and trade, but are not typically considered in the literature. In this work, I allow
24 for groups to overlap, but take the set of permissible coalitions to be exogenously set. An
25 interesting question is *which* coalitions would form when they are permitted to overlap,
26 allowing for an endogenous formation of coalitions. This would build on the literature of
27 endogenous coalition formation, for instance Ray and Vohra (1997); Diamantoudi and Xue

³⁰Bloch and De Clippel (2010) studies the core of cooperative games and characterise for which cooperative games is the core of the sum of those games the same as the sum of the cores.

³¹When these issues are independent, it is trivial to show that the set of payoffs sustainable in Negotiated Binding Agreements is weakly larger, however, when there is interdependence between these games the relation is unclear.

³²Conconi and Perroni (2002) considers this question for international trade, but do so via a coalition formation procedure, a la Ray and Vohra (1997). Such procedures are fragile to changes in the games and definitions (see, example 1 of Gavan (2022)) and therefore it is difficult to make broad conclusions, whereas the results of my paper may allow for a better understanding within classes of games.

³³To see this, consider the two-player case. Notice that ruling out an action for player 1 in the underlying game makes the harshest punishment for player 2 weakly better (from player 2’s perspective). However, as player 2 must prefer punishing player 1 than being punished, this can limit the punishments that player 2 can use against player 1. With this, it may lead to an indirect improvement in the agreements that can occur for player 1.

(2007). Exploring overlapping coalition formation is especially interesting when transfers are not permitted, as then it is not possible to treat those coalitions that overlap equally as one larger coalition, by moving transfers between those agents to align incentives.

A Proofs

Proof of lemma 1: Notice that $\liminf_{k \rightarrow \infty} u_i(a^k) = (1 - \delta) \sum_{t=1}^{\infty} \delta^{t-1} \liminf_{k \rightarrow \infty} u_i(a^k)$. Therefore by continuity of subtraction we have that:

$$\lim_{\delta \rightarrow 1} (1 - \delta) \sum_{t=1}^{\infty} \delta^{t-1} u_i(a^t) - \liminf_{k \rightarrow \infty} u_i(a^k) = \lim_{\delta \rightarrow 1} (1 - \delta) \sum_{t=1}^{\infty} \delta^{t-1} \left(u_i(a^t) - \liminf_{k \rightarrow \infty} u_i(a^k) \right)$$

Note by definition of the \liminf , for all $\epsilon > 0 \exists T \in \mathbb{N}$ such that $\forall t > T$ we have that $u_i(a^t) - \liminf_{k \rightarrow \infty} u_i(a^k) > -\epsilon$. Therefore, for any such T , we may decompose the expression as follows.

$$\begin{aligned} \lim_{\delta \rightarrow 1} (1 - \delta) \sum_{t=1}^{\infty} \delta^{t-1} \left(u_i(a^t) - \liminf_{k \rightarrow \infty} u_i(a^k) \right) &= \lim_{\delta \rightarrow 1} (1 - \delta) \sum_{t=1}^T \delta^{t-1} \left(u_i(a^t) - \liminf_{k \rightarrow \infty} u_i(a^k) \right) + \dots \\ &\quad \dots + \lim_{\delta \rightarrow 1} (1 - \delta) \sum_{t=T+1}^{\infty} \delta^{t-1} \left(u_i(a^t) - \liminf_{k \rightarrow \infty} u_i(a^k) \right) \\ &= \lim_{\delta \rightarrow 1} (1 - \delta) \sum_{t=T+1}^{\infty} \delta^{t-1} \left(u_i(a^t) - \liminf_{k \rightarrow \infty} u_i(a^k) \right) \\ &> \lim_{\delta \rightarrow 1} (1 - \delta) \sum_{t=T+1}^{\infty} \delta^{t-1} (-\epsilon) \\ &= \lim_{\delta \rightarrow 1} -\delta^{T+1} \epsilon \\ &= -\epsilon \end{aligned}$$

Therefore $\lim_{\delta \rightarrow 1} (1 - \delta) \sum_{t=1}^{\infty} \delta^{t-1} \left(u_i(a^t) - \liminf_{k \rightarrow \infty} u_i(a^k) \right) > -\epsilon \forall \epsilon > 0$, concluding that $\lim_{\delta \rightarrow 1} (1 - \delta) \sum_{t=1}^{\infty} \delta^{t-1} \left(u_i(a^t) - \liminf_{k \rightarrow \infty} u_i(a^k) \right) \geq 0$ and therefore

$$\lim_{\delta \rightarrow 1} (1 - \delta) \sum_{t=1}^{\infty} \delta^{t-1} u_i(a^t) \geq \liminf_{k \rightarrow \infty} u_i(a^k)$$

By analogy $\lim_{\delta \rightarrow 1} (1 - \delta) \sum_{t=1}^{\infty} \delta^{t-1} u_i(a^t) \leq \limsup_{k \rightarrow \infty} u_i(a^k)$. ■

Proof of lemma 2: Suppose not, $U_i(s^*|h) < \inf_{a_{-i} \in A_{-i}} \sup_{a_i \in A_i} u_i(a_i, a_{-i})$. For any $\epsilon > 0$, let $\tilde{a}_i : A_{-i} \rightarrow A_i$ be such that $u_i(\tilde{a}_i(a_{-i}), a_{-i}) > \sup_{a_i \in A_i} u_i(a_i, a_{-i}) - \epsilon$. Note such a function exists for any $\epsilon > 0$. Let $s'_i(h) = (\tilde{a}_i(s_{-i}^*(h')), s_{-i}^*(h'))$ for all $h' \in H$. It follows that $U_i(s'_i, s_{-i}^*|h)$ is either such that it ends in agreement, in which case $U_i(s'_i, s_{-i}^*|h) > \inf_{a_{-i} \in A_{-i}} \sup_{a_i \in A_i} u_i(a_i, a_{-i}) - \epsilon$ and therefore, as we can construct such a function for any $\epsilon > 0$, we conclude that $U_i(s'_i, s_{-i}^*|h) \geq \inf_{a_{-i} \in A_{-i}} \sup_{a_i \in A_i} u_i(a_i, a_{-i})$. On the other hand, it may be that $U_i(s'_i, s_{-i}^*|h)$ ends in perpetual disagreement. In which case $(s'_i, s_{-i}^*|h) =$

1 $(a^1, a^2, \dots, a^T, \dots)$, where $a_i^t = \tilde{a}_i(a_{-i}^t)$. Therefore:

$$U_i(s'_i, s_{-i}^*|h) \geq \liminf_{t \rightarrow \infty} u_i(\tilde{a}_i(a_{-i}^t), a_{-i}^t) \geq \inf_{a_{-i} \in A_{-i}} u_i(\tilde{a}_i(a_{-i}), a_{-i}) > \inf_{a_{-i} \in A_{-i}} \sup_{a_i \in A_{-i}} u_i(a_i, a_{-i}) - \epsilon$$

2 $\Rightarrow U_i(s'_i, s_{-i}^*|h) \geq \inf_{a_{-i} \in A_{-i}} \sup_{a_i \in A_{-i}} u_i(a_i, a_{-i})$. A contradiction. ■

3 **Proof of theorem 1 and proposition 1:** By theorem 2 such a^* can be supported
4 and therefore that section of the proof is omitted.

5 To see that only such a^* can be sustained, take any a^* such that it is supported by a no
6 delay Negotiated Binding Agreement in the n -player case and a Negotiated Binding Agree-
7 ment in the two-player case given by the SPE s^* . Denote $\tilde{A} = \{a \in A | \exists h \in H \text{ s.t. } s^*(h) =$
8 $a\}$. Note by no delay these completely define the set of actions that can be agreed upon in
9 the n -player case by no delay. Further to this, note that $s_{-i}^*(h) \in \tilde{A}_{-i}$ for all $h \in H$ by no
10 delay and in the two-player case by no babbling. As s^* is an SPE it must be that there is no
11 profitable deviation. Notice that $U_i(s^*|h) \geq \max_{a_i \in A_i} \inf_{a_{-i} \in \tilde{A}_{-i}} u_i(a_i, a_{-i})$. Suppose not
12 $U_i(s^*|h) < \max_{a_i \in A_i} \inf_{a_{-i} \in \tilde{A}_{-i}} u_i(a_i, a_{-i})$. It follows that $\max_{a_i \in A_i} \inf_{a_{-i} \in \tilde{A}_{-i}} u_i(a_i, a_{-i}) -$
13 $U_i(s^*|h) > 0$. Consider a deviation to s'_i such that $s'_i(h') = s_i^*(h')$ for all h' such that
14 $h = (h', h'')$ while $s'_i(h')$ is such that $u_i((s'_i, s_{-i}^*)(h')) = \max_{a_i \in A_i} u_i(a_i, s_{-i}^*(h'))$ for all
15 other histories. Suppose such a deviation leads to perpetual disagreement. Denote the
16 sequence induced by such a strategy by $z' = (a^1, a^2, \dots, a^t, \dots)$. Notice that $u_i(a_i^t, a_{-i}^t) =$
17 $\max_{a_i \in A_i} u_i(a_i, a_{-i}^t)$. Note that therefore $u_i(a_i^t, a_{-i}^t) \geq \max_{a_i \in A_i} \inf_{a_{-i} \in \{a'_{-i} \in A_{-i} | a'_{-i} = a_{-i}^k\}} u_i(a_i, a_{-i})$.
18 By definition:

$$\begin{aligned} U_i(s_i, s_{-i}^*|h) &\geq \liminf_{t \rightarrow \infty} u_i(a_i^t) \\ &\geq \liminf_{t \rightarrow \infty} \max_{a_i \in A_i} \inf_{a_{-i} \in \{a'_{-i} \in A_{-i} | a'_{-i} = a_{-i}^k\}} u_i(a_i, a_{-i}) \\ &= \max_{a_i \in A_i} \inf_{a_{-i} \in \{a'_{-i} \in A_{-i} | a'_{-i} = a_{-i}^k\}} u_i(a_i, a_{-i}) \\ &\geq \max_{a_i \in A_i} \inf_{a_{-i} \in \tilde{A}_{-i}} u_i(a_i, a_{-i}) \Rightarrow U_i(s_i, s_{-i}^*|h) > U_i(s^*|h) \end{aligned}$$

19 therefore it cannot be that s^* is an SPE if the deviation ends in perpetual disagreement.
20 The argument for agreement is direct from the definition.

21 Therefore it must be that $U_i(s^*|h) \geq \max_{a_i \in A_i} \inf_{a_{-i} \in \tilde{A}_{-i}} u_i(a_i, a_{-i})$. As \tilde{A} are agreed
22 upon, therefore be $\forall \tilde{a} \in \tilde{A} u_i(\tilde{a}) \geq \max_{a_i \in A_i} \inf_{a_{-i} \in \tilde{A}_{-i}} u_i(a_i, a_{-i})$. Therefore $\exists a'_{-i} \in \tilde{\tilde{A}}_{-i}$,
23 where $\tilde{\tilde{A}}_{-i}$ is the limit points of \tilde{A}_{-i} such that $u_i(\tilde{a}) \geq \max_{a_i \in A_i} u_i(a_i, a'_{-i})$. As this holds for
24 all $\tilde{a} \in \tilde{A}$ it follows that $u_i(a') \geq \max_{a_i \in A_i} u_i(a_i, a'_{-i})$ therefore $u_i(a') = \max_{a_i \in A_i} u_i(a_i, a'_{-i})$.
25 therefore $\exists a^i \in \tilde{\tilde{A}}$ such that $u_i(\tilde{a}) \geq u_i(a^i) = \max_{a_i \in A_i} u_i(a_i, a^i_{-i})$. Notice that: $u_i(\tilde{a}) \geq$
26 $u_i(a^i)$ for all \tilde{A} and therefore $u_i(a^j) \geq u_i(a^i)$ and $u_i(a^*) \geq u_i(a^i)$. Therefore such a profile
27 of action profiles must exist for a^* to be supported. ■

Proof of theorem 2: Note within this proof I maintain the notation a^k to refer to the k^{th} period proposal in a history h , while I use \underline{a}^j to denote the action profile used in equilibrium as a punishment for j . Let s^* be as follows:

1. if $h = (a^1, \dots, a^k)$ is such that there is some $j \in N$, such that $a_{-j}^{k-1} = s_{-j}^*((a^1, \dots, a^{k-2}))$ and either

- (a) $a_l^k = s_l^*(h \setminus a^{k-1}) \quad \forall l \neq j$ while $a_j^k \neq s_j^*(h \setminus a^{k-1})$.

- (b) or $a_{-j}^k = \underline{a}_{-j}^j$.

then $s_i^*(h) = \underline{a}_i^j$.

2. $s_i^*(h) = a_i^*$ otherwise.

First note that from any history the continuation is terminal within two periods and therefore no babbling is satisfied. Now to show that s^* is a Subgame Perfect Equilibrium of the negotiation game. Suppose that a profitable deviation exists at a history $h \in H$ for $i \in N$. If the deviation does not include some different proposal within two periods of h it cannot be profitable, as the outcome remains the same. Therefore any deviation must occur within two periods. Any such deviation, denoted by s'_i , if it does not lead to the same terminal history and therefore cannot be profitable, of $i \in N$ must lead to \underline{a}_{-i}^i for all periods following. Let the terminal history following the deviation be denoted by $(s_{-i}^*, s'_i|h) = (h, a^k, a^{k+1}, \dots, a^t, \dots)$. When $(s_{-i}^*, s'_i|h) \in Z'$ let $(s_{-i}^*, s'_i|h) = (h, a'^1, a'^2, \dots, a((s_{-i}^*, s'_i|h)), a((s_{-i}^*, s'_i|h)), a((s_{-i}^*, s'_i|h)), \dots)$, i.e let the agreement that $(s_{-i}^*, s'_i|h)$ concludes in be infinitely repeated at the end of the sequence, with some abuse of notation. However, by construction, it must be that $\limsup_{t \rightarrow \infty} u(a^t) \leq u_i(\underline{a}^i)$ and therefore it must be at least weakly worse than any terminal history of the strategy s^* . Therefore no profitable deviation exists. ■

Proof of theorem 5: Suppose not, for some history $h' \in H$ we have that $s_i(h') = a_i$. By no babbling it follows that there exists some $h \in H$ such that $a_i(s|h) = a_i$. Therefore it must be that $U_i(s|h) = u_i(a(s|h)) \leq \sup_{a'_{-i} \in A_{-i}} u_i(a_i, a'_{-i})$. Take $\epsilon = \inf_{a'_{-i} \in A_{-i}} \sup_{a'_i \in A_i} u_i(a'_i, a'_{-i}) - u_i(a(s|h)) > 0$. Take a function $\tilde{a}_i : A_{-i} \rightarrow A_i$ such that $u_i(\tilde{a}_i(a_{-i}), a_{-i}) > \sup_{a'_i \in A_i} u_i(a'_i, a_{-i}) - \epsilon$. Consider a deviation $s'_i(h'') = \tilde{a}_i(s_{-i}(h''))$ for all $h'' \in H$. It follows that:

$$U_i(s'_i, s_{-i}|h) \geq \inf_{a_{-i} \in A_{-i}} u_i(\tilde{a}_i(a_{-i}), a_{-i}) > \inf_{a_{-i} \in A_{-i}} \sup_{a'_i \in A_i} u_i(a'_i, a_{-i}) - \epsilon$$

Therefore it follows that $U_i(s'_i, s_{-i}|h) > u_i(a(s|h)) = U_i(s|h)$, concluding that a profitable deviation exists and therefore it cannot be that s is a Subgame Perfect Equilibrium of the negotiation game. By no babbling, we conclude that $s_i(h) \notin D_i(A_{-i})$ for any $h \in H$.

Now suppose by contradiction that, for all $j \in N$ $s_j(h') \in \tilde{A}_j^k \quad \forall k < m$ and $h' \in H$ but for some $i \in N$ $s_j(h') = a_i \notin \tilde{A}_j^{m+1}$ for some $h' \in H$. By no babbling it must be

that a) $s_{-i}(h') \in \tilde{A}_i^m$ for all h' and b) by no babbling there is some $h \in H$ for which $a_i(s|h) = a_i$. Therefore it must be that $U_i(s|h) = u_i(a(s|h)) \leq \sup_{a'_{-i} \in \tilde{A}_{-i}^m} u_i(a_i, a'_{-i})$. Take $\epsilon = \inf_{a'_{-i} \in \tilde{A}_{-i}^m} \sup_{a'_i \in A_i} u_i(a'_i, a'_{-i}) - u_i(a(s|h)) > 0$. Take a function $\tilde{a}_i : \tilde{A}_{-i}^m \rightarrow A_i$ such that $u_i(\tilde{a}_i(a_{-i}), a_{-i}) > \sup_{a'_i \in A_i} u_i(a'_i, a_{-i}) - \epsilon$ for all $a_{-i} \in \tilde{A}_{-i}^m$. Consider a deviation $s'_i(h'') = \tilde{a}_i(s_{-i}(h''))$ for all $h'' \in H$. It follows that:

$$U_i(s'_i, s_{-i}|h) \geq \inf_{a_{-i} \in \tilde{A}_{-i}^m} u_i(\tilde{a}_i(a_{-i}), a_{-i}) > \inf_{a_{-i} \in \tilde{A}_{-i}^m} \sup_{a'_i \in A_i} u_i(a'_i, a_{-i}) - \epsilon$$

Therefore it follows that $U_i(s'_i, s_{-i}|h) > u_i(a(s|h)) = U_i(s|h)$, concluding that a profitable deviation exists and therefore it cannot be that s is a Subgame Perfect Equilibrium of the negotiation game. By no babbling, we conclude that $s_i(h) \notin D_i(\tilde{A}_{-i}^m)$ for any $h \in H$ and therefore $s_i(h) \in \tilde{A}_i^{k+1}$, a contradiction. ■

Proof of lemma 3: Note that $B^0 = \tilde{A}^0$. Now we will show that $B^k \subseteq \tilde{A}^k$ for all $k \geq 0$. By the inductive hypothesis suppose that $B^m \subseteq \tilde{A}^m$ for all $m < k$. Now notice that for any $a_i \in B_i^k$ we have that there is some $a_{-i} \in B_{-i}^{k-1} \subseteq \tilde{A}_{-i}^{k-1}$ such $u_i(a_i, a_{-i}) \geq u_i(a'_i, a_{-i})$ for all $a'_i \in A_i$. It follows that $u_i(a_i, a_{-i}) \geq \inf_{a'_{-i} \in B_{-i}^{k-1}} \sup_{a'_i \in A_i} u_i(a'_i, a'_{-i}) \geq \inf_{a'_{-i} \in \tilde{A}_{-i}^{k-1}} \sup_{a'_i \in A_i} u_i(a'_i, a'_{-i})$. Further,

$$u_i(a_i, a_{-i}) \leq \sup_{a''_{-i} \in B_{-i}^k} u_i(a_i, a''_{-i}) \leq \sup_{a''_{-i} \in \tilde{A}_{-i}^k} u_i(a_i, a''_{-i})$$

Therefore we conclude that if $a_i \in B_i^k$ then $a_i \in \tilde{A}_i^k$, concluding the proof. ■

Proof of theorem 4: Suppose not, then there is some $i \in N$ and $h \in H$ such that that $\inf_{a'_{-i} \in IIR_{-i}} \sup_{a'_i \in A_i} u_i(a'_i, a'_{-i}) > U_i(s^*|h)$. It must be that a) s^* is a Subgame Perfect Equilibrium of the negotiation game and b) by theorem 3 it must be that $s^*_{-i}(h) \in IIR_{-i}$ for all $h \in H$. Let $\epsilon = \inf_{a'_{-i} \in IIR_{-i}} \sup_{a'_i \in A_i} u_i(a'_i, a'_{-i}) - U_i(s^*|h) > 0$. Construct $\tilde{a}_i : IIR_{-i} \rightarrow A_i$ such that $u_i(\tilde{a}_i(a_{-i}), a_{-i}) \geq \sup_{a_i \in A_i} u_i(a_i, a_{-i}) - \frac{\epsilon}{2}$ for all $a_{-i} \in IIR_{-i}$. Consider a deviation to $s'_i(h')$ such that $s'_i(h') = \tilde{a}_i(s^*_{-i}(h'))$ for all $h' \in H$ at the history h . It follows that:

$$\begin{aligned} U_i(s'_i, s^*_{-i}|h) &\geq \inf_{a_{-i} \in IIR_{-i}} u_i(\tilde{a}_i(a_{-i}), a_{-i}) = \inf_{a_{-i} \in IIR_{-i}} \sup_{a_i \in A_i} u_i(a_i, a_{-i}) - \frac{\epsilon}{2} \\ &= \frac{\inf_{a_{-i} \in IIR_{-i}} \sup_{a_i \in A_i} u_i(a_i, a_{-i}) + U_i(s^*|h)}{2} > U_i(s^*|h) \end{aligned}$$

A contradiction, as therefore s^* is not a Subgame Perfect Equilibrium of the negotiation game and therefore not a Negotiated Binding Agreement. ■

Proof of lemma 4: As a^* is a strong Nash equilibrium, it follows that $\nexists C \in 2^N \setminus \{\emptyset\}, a'_C \in A_C$ such that $u_i(a'_C, a^*_{-C}) > u_i(a^*)$ for all $i \in C$. Therefore a^* is not coalitionally irrational. Now suppose that $a^* \in \tilde{A}^m(C)$ for all $m < k$. Notice that by the same statement this implies that $a^* \in \tilde{A}^{m+1}(C)$. This implies that $a^* \in ICIR(C)$ for all C . ■

Proof of theorem 5: Suppose not, for some history $h' \in H$ we have that $s_C(h') = a_C$. By C no babbling it follows that there exists some $h \in H$ such that $a_C(s|h) = a_C$. Therefore it must be that $U_i(s^*|h) = u_i(a(s^*|h)) \leq \sup_{a'_C \in A_{-C}} u_i(a_C, a'_C)$ for all $i \in C$. By definition of a_C being not coalitionally rational, there exists a function $a'_C : A_{-C} \rightarrow A_C$ such that $\inf_{a_{-C} \in A_{-C}} u_i(a'_C(a_{-C}), a_{-C}) > \sup_{a'_C \in A_{-C}} u_i(a_C, a'_C)$. Consider a deviation of C at history h such that $s_C(h') = a'_C(s_{-C}(h'))$ for all $h' \in H$. It follows that:

$$U_i(s'_C, s^*_{-C}|h) \geq \inf_{a_{-C} \in A_{-C}} u_i(a'_C(a_{-C}), a_{-C}) > \sup_{a'_C \in A_{-C}} u_i(a_C, a'_C) \geq U_i(s^*|h)$$

for all $i \in C$. Concluding that s^* is not a \mathcal{C} -Subgame Perfect Equilibrium.

Now suppose by contradiction that $s(h') \in \tilde{A}^k(\mathcal{C}) \forall k < m$ and $h' \in H$ but $s(h') = a \notin \tilde{A}^{m+1}(\mathcal{C})$ for some $h' \in H$. By definition, it must be that $a \in \bigcup_{C \in \mathcal{C}} [D_C(\tilde{A}^{m-1}(\mathcal{C})_{-C}) \times A_{-C}]$. Therefore it must be that $a_C \in D_C(\tilde{A}^{m-1}(\mathcal{C})_{-C})$ for some $C \in \mathcal{C}$. By \mathcal{C} -no babbling we have that $\exists h \in H$ such that $a_C = a^*_C(s^*|h)$. By definition of coalition rationality given $\tilde{A}^{m-1}(\mathcal{C})_{-C}$, as $a_C \in D_C(\tilde{A}^{m-1}(\mathcal{C})_{-C})$ there must be some that there is some $a'_C : \tilde{A}^{m-1}(\mathcal{C})_{-C}$ such that $\inf_{a_{-C} \in \tilde{A}^{m-1}(\mathcal{C})_{-C}} u_i(a'_C(a_{-C}), a_{-C}) > \sup_{a'_C \in \tilde{A}^{m-1}(\mathcal{C})_{-C}} u_i(a_C, a'_C)$. Consider a deviation of C at history h such that $s_C(h') = a'_C(s_{-C}(h'))$ for all $h' \in H$. It follows that:

$$U_i(s'_C, s^*_{-C}|h) \geq \inf_{a_{-C} \in \tilde{A}^{m-1}(\mathcal{C})_{-C}} u_i(a'_C(a_{-C}), a_{-C}) > \sup_{a'_C \in \tilde{A}^{m-1}(\mathcal{C})_{-C}} u_i(a_C, a'_C)$$

. Therefore $U_i(s'_C, s^*_{-C}|h) > U_i(s^*|h)$ for all $i \in C$. Concluding that s^* is not a \mathcal{C} -Subgame Perfect Equilibrium of the negotiation game. A contradiction. ■

Proof of theorem 6: Suppose this is not the case. There is some $C \in \mathcal{C}$ $a'_C : [ICIR(\mathcal{C})]_{-C} \rightarrow A_C$ such that $\inf_{a_{-C} \in [ICIR(\mathcal{C})]_{-C}} u_i(a'_C(a_{-C}), a_{-C}) > U_i(s^*|h)$ for all $i \in C$. It must be that s^* is a \mathcal{C} -Subgame Perfect Equilibrium of the negotiation game, and therefore there cannot exist a profitable deviation for C . Notice that $s^*_i(h) \in [ICIR(\mathcal{C})]_i$ for all $i \in N$.

Consider a joint deviation from coalition C such that $s'_C(h) = a'_C(s^*_{-C}(h))$ for all $h \in H$. By the definition of the utilities that this can induce, it is clear that:

$$U_i(s'_C, s^*_{-C}|h) \geq \inf_{a_{-C} \in [ICIR(\mathcal{C})]_{-C}} u_i(a'_C(a_{-C}), a_{-C})$$

for all $i \in C$, and therefore $U_i(s'_C, s^*_{-C}|h) > U_i(s^*|h)$ for all $i \in C$. In conclusion, s^* cannot be a \mathcal{C} -Subgame Perfect Equilibrium. ■

Proof of theorem 7: Consider the following strategy:

1. if $h = (a^1, \dots, a^k)$ is such that there is some $C \in \mathcal{C}$, such that $a^{k-1}_{-C} = s^*_{-C}((a^1, \dots, a^{k-2}))$ and either $a^k_l = s^*_l(h \setminus a^{k-1}) \quad \forall l \notin C$ while $a^k_j \neq s^*_j(h \setminus a^{k-1})$ for all $j \in C$ or $a^k_{-C} = \underline{a}_{-C}$ then $s^*_i(h) = \underline{a}^C_i$.

2. $s_i^*(h) = a_i^*$ otherwise.

By definite, at no history can N deviate as a coalition to improve all their utilities if $N \in \mathcal{C}$. Now assume that some other coalition $C \in \mathcal{C}$ has a profitable deviation. If $a_j \neq s_j^*(h)$ for all $j \in C$, then it cannot be profitable as it leads to a history that induces the \underline{a}_C for all periods. If $a_j \neq s_j^*(h)$ for all $j \in B$, where $B \subset C$, while $a_j^* = s_j^*(h)$. Then it must induce a path such that either a member of B is worse off, or further deviations within C take place. Either way, it cannot be that this is a profitable deviation.

As all histories end within 2 periods we satisfy the condition of no babbling agreements and therefore we have a \mathcal{C} -Negotiated Binding Agreement. ■

References

- Abreu, D., Dutta, P. K., and Smith, L. (1994). The Folk Theorem for Repeated Games: A New Condition. *Econometrica*, 62(4):939–948.
- Aghion, P., Antràs, P., and Helpman, E. (2007). Negotiating Free Trade. *Journal of International Economics*, 73(1):1–30.
- Ambrus, A. (2006). Coalitional Rationalizability. *The Quarterly Journal of Economics*, 121(3):903–929.
- Ambrus, A. (2009). Theories of Coalitional Rationality. *Journal of Economic Theory*, 144(2):676–695.
- Aumann, R. J. (1959). Acceptable Points in General Cooperative n -person Games. *Contributions to the Theory of Games (AM-40)*, 4:287–324.
- Aumann, R. J. (1961). The Core of a Cooperative Game without Side Payments. *Transactions of the American Mathematical Society*, 98(3):539–552.
- Aumann, R. J. and Shapley, L. S. (1994). Long-Term Competition—a game-theoretic analysis. In *Essays in game theory*, pages 1–15. Springer.
- Baron, E. J. (2018). The Effect of Teachers’ Unions on Student Achievement in the Short Run: Evidence from Wisconsin’s Act 10. *Economics of Education Review*, 67:40–57.
- Bernheim, B. D. (1984). Rationalizable Strategic Behavior. *Econometrica: Journal of the Econometric Society*, pages 1007–1028.
- Bernheim, B. D., Peleg, B., and Whinston, M. D. (1987). Coalition-Proof Nash Equilibria i. Concepts. *Journal of Economic Theory*, 42(1):1–12.
- Bernheim, B. D. and Ray, D. (1989). Collective Dynamic Consistency in Repeated Games. *Games and Economic Behavior*, 1(4):295–326.
- Bhaskar, V. (1989). Quick Responses in Duopoly Ensure Monopoly Pricing. *Economics Letters*, 29(2):103–107.
- Biasi, B. (2021). The Labor Market for Teachers under Different Pay Schemes. *American Economic Journal: Economic Policy*, 13(3):63–102.

- 1 Biasi, B. and Sarsons, H. (2021). Information, Confidence, and the Gender Gap in Bargaining. *AEA Papers*
2 *and Proceedings*, 111:174–78.
- 3 Biasi, B. and Sarsons, H. (2022). Flexible Wages, Bargaining, and the Gender Gap. *The Quarterly Journal*
4 *of Economics*, 137(1):215–266.
- 5 Bloch, F. and De Clippel, G. (2010). Cores of combined games. *Journal of Economic Theory*, 145(6):2424–
6 2434.
- 7 Busch, L.-A. and Wen, Q. (1995). Perfect equilibria in a negotiation model. *Econometrica: Journal of the*
8 *Econometric Society*, pages 545–565.
- 9 Carraro, C. (1998). Beyond Kyoto: A Game-Theoretic Perspective. In *the Proceedings of the OECD*
10 *Workshop on “Climate Change and Economic Modelling. Background Analysis for the Kyoto Protocol”*,
11 *Paris*, pages 17–18. Citeseer.
- 12 Carraro, C., Eyckmans, J., and Finus, M. (2006). Optimal Transfers and Participation Decisions in Inter-
13 national Environmental Agreements. *The Review of International Organizations*, 1(4):379–396.
- 14 Chakrabarti, S. K. (1988). Refinements of the β -core and the strong equilibrium and the aumann proposition.
15 *International Journal of Game Theory*, 17:205–224.
- 16 Chander, P. (2007). The gamma-Core and Coalition Formation. *International Journal of Game Theory*,
17 35(4):539–556.
- 18 Chander, P. and Tulkens, H. (1997). The Core of an Economy with Multilateral Environmental Externalities.
19 *International Journal of Game Theory*, 26(3):379–401.
- 20 Chander, P. and Wooders, M. (2020). Subgame-Perfect Cooperation in an Extensive Game. *Journal of*
21 *Economic Theory*, page 105017.
- 22 Chatterjee, K., Dutta, B., Ray, D., and Sengupta, K. (1993). A Noncooperative Theory of Coalitional
23 Bargaining. *The Review of Economic Studies*, 60(2):463–477.
- 24 Chwe, M. S.-Y. (1994). Farsighted Coalitional Stability. *Journal of Economic Theory*, 63(2):299–325.
- 25 Conconi, P. and Perroni, C. (2002). Issue linkage and issue tie-in in multilateral negotiations. *Journal of*
26 *international Economics*, 57(2):423–447.
- 27 Currarini, S. and Marini, M. (2003). A Sequential Approach to the Characteristic Function and the Core in
28 Games with Externalities. In *Advances in Economic Design*, pages 233–249. Springer.
- 29 Diamantoudi, E. and Xue, L. (2007). Coalitions, Agreements and Efficiency. *Journal of Economic Theory*,
30 136(1):105–125.
- 31 Doval, L. and Ely, J. C. (2020). Sequential information design. *Econometrica*, 88(6):2575–2608.
- 32 Ellingsen, T. and Paltseva, E. (2016). Confining the Coase Theorem: contracting, ownership, and free-riding.
33 *The Review of Economic Studies*, 83(2):547–586.
- 34 Farrell, J. and Maskin, E. (1989). Renegotiation in Repeated Games. *Games and Economic Behavior*,
35 1(4):327–360.

- 1 Fudenberg, D. and Maskin, E. (1986). The Folk Theorem in Repeated Games with Discounting or with
2 Incomplete Information. *Econometrica*, 54(3):533–554.
- 3 Gavan, M. J. (2022). Weak Coalitional Equilibrium: Existence and Overlapping Coalitions. *Working Paper*.
- 4 Grandjean, G., Mauleon, A., and Vannetelbosch, V. (2017). Strongly Rational Sets for normal-form Games.
5 *Economic Theory Bulletin*, 5(1):35–46.
- 6 Greenberg, J. (1990). *The theory of social situations: an alternative game-theoretic approach*. Cambridge
7 University Press.
- 8 Grossman, G. M., McCalman, P., and Staiger, R. W. (2021). The “new” economics of trade agreements:
9 From trade liberalization to regulatory convergence? *Econometrica*, 89(1):215–249.
- 10 Halpern, J. Y. and Pass, R. (2018). Game Theory with Translucent Players. *International Journal of Game*
11 *Theory*, 47(3):949–976.
- 12 Harsanyi, J. C. (1974). An Equilibrium-Point Interpretation of Stable Sets and a Proposed Alternative
13 Definition. *Management Science*, 20(11):1472–1495.
- 14 Harstad, B. (2022). A theory of pledge-and-review bargaining. *Journal of Economic Theory*, page 105574.
- 15 Herings, P. J.-J., Mauleon, A., and Vannetelbosch, V. J. (2004). Rationalizability for Social Environments.
16 *Games and Economic Behavior*, 49(1):135–156.
- 17 Jackson, M. O. and Wilkie, S. (2005). Endogenous games and mechanisms: Side payments among players.
18 *The Review of Economic Studies*, 72(2):543–566.
- 19 Kalai, A. T., Kalai, E., Lehrer, E., and Samet, D. (2010). A Commitment Folk Theorem. *Games and*
20 *Economic Behavior*, 69(1):127–137. Special Issue In Honor of Robert Aumann.
- 21 Kalai, E. (1981). Preplay Negotiations and the Prisoner’s Dilemma. *Mathematical Social Sciences*, 1(4):375–
22 379.
- 23 Kimya, M. (2020). Equilibrium coalitional behavior. *Theoretical Economics*, 15(2):669–714.
- 24 Li, S. (2017). Obviously Strategy-Proof Mechanisms. *American Economic Review*, 107(11):3257–87.
- 25 Limao, N. (2016). Preferential Trade Agreements. In *Handbook of Commercial Policy*, volume 2, pages 281
26 – 360. Elsevier.
- 27 Mariotti, M. (1997). A Model of Agreements in Strategic Form Games. *Journal of Economic Theory*,
28 74(1):196–217.
- 29 Nakanishi, N. (2009). Noncooperative Farsighted Stable Set in an n-player Prisoners’ Dilemma. *International*
30 *Journal of Game Theory*, 38(2):249–261.
- 31 Nash, J. (1953). Two-person Cooperative Games. *Econometrica: Journal of the Econometric Society*, pages
32 128–140.
- 33 Nishihara, K. (2022). Resolution of the N-Person Prisoners’ Dilemma by Kalai’s Preplay Negotiation Pro-
34 cedure. Available at SSRN 4112007.

- 1 Pearce, D. G. (1984). Rationalizable Strategic Behavior and the Problem of Perfection. *Econometrica:*
2 *Journal of the Econometric Society*, pages 1029–1050.
- 3 Peters, M. and Szentes, B. (2012). Definable and Contractible Contracts. *Econometrica*, 80(1):363–411.
- 4 Rabin, M. (1994). A Model of pre-game Communication. *Journal of Economic Theory*, 63(2):370–391.
- 5 Ray, D. and Vohra, R. (1997). Equilibrium Binding Agreements. *Journal of Economic Theory*, 73(1):30–78.
- 6 Ray, D. and Vohra, R. (2015). The Farsighted Stable Set. *Econometrica*, 83(3):977–1011.
- 7 Ray, D. and Vohra, R. (2019). Maximality in the Farsighted Stable Set. *Econometrica*, 87(5):1763–1779.
- 8 Rubinstein, A. (1979). Equilibrium in Supergames with the Overtaking Criterion. *Journal of Economic*
9 *Theory*, 21(1):1–9.
- 10 Rubinstein, A. (1980). Strong perfect equilibrium in supergames. *International Journal of Game Theory*,
11 9(1):1–12.
- 12 Rubinstein, A. (1982). Perfect Equilibrium in a Bargaining Model. *Econometrica: Journal of the Econo-*
13 *metric Society*, pages 97–109.
- 14 Rubinstein, A. (1994). Equilibrium in Supergames. In *Essays in Game Theory*, pages 17–27. Springer.
- 15 Salcedo, B. (2017). Interdependent Choices. Technical report, University of Western Ontario.
- 16 Scarf, H. E. (1971). On the Existence of a Cooperative Solution for a General Class of N-person Games.
17 *Journal of Economic Theory*, 3(2):169–181.
- 18 Selten, R. (1965). Spieltheoretische behandlung eines oligopolmodells mit nachfrageträgheit: Teil i: Bes-
19 timmung des dynamischen preisgleichgewichts. *Zeitschrift für die gesamte Staatswissenschaft/Journal of*
20 *Institutional and Theoretical Economics*, (H. 2):301–324.
- 21 Selten, R. (1988). *Reexamination of the Perfectness Concept for Equilibrium Points in Extensive Games*,
22 pages 1–31. Springer Netherlands, Dordrecht.
- 23 Shubik, M. (2012). What is a Solution to a Matrix Game. *Cowles Foundation Discussion Paper N. 1866*,
24 Available at SSRN: https://papers.ssrn.com/sol3/papers.cfm?abstract_id=2220772.
- 25 Suzuki, A. and Muto, S. (2005). Farsighted Stability in an n-Person Prisoner’s Dilemma. *International*
26 *Journal of Game Theory*, 33(3):431–445.
- 27 Tennenholtz, M. (2004). Program Equilibrium. *Games and Economic Behavior*, 49(2):363–373.
- 28 Xue, L. (1998). Coalitional Stability under Perfect Foresight. *Economic Theory*, 11(3):603–627.
- 29 Yamada, A. (2003). Efficient Equilibrium Side Contracts. *Economics Bulletin*, 3(6):1–7.

1 B Online Appendix

2 B.1. Application of Public Good Provision

3 To further demonstrate the logic of the main proofs, I turn to a public goods game as the
4 underlying environment being negotiated over. In this setting, it is shown that contribu-
5 tion of an agent can only be supported if and only if sufficiently many other agents also
6 contribute. With this, full and no contribution can be supported.

7 **Application 2. (Public Goods Game)** Consider the underlying game, G , to be the
8 following Public Goods Game. $N = \{1, 2, \dots, n\}$. Let $A_i = \{c, d\}$ for each i . Let $u_i(a) =$
9 $1 + k \left[\sum_{j \in N} \mathbf{1}_{a_j=c} \right] - \mathbf{1}_{a_i=c}$ with $k \in (\frac{1}{n}, 1)$.

10 Firstly notice that for any player it is strictly dominant to choose d and hence the only
11 Nash equilibrium payoff is 1.

12 I will now construct a strategy that allows for any action profile that Pareto dominates
13 the Nash equilibrium to be supported by Negotiated Binding Agreement. Specifically, let
14 a^* denote an action profile such that $u_i(a^*) = 1 + k|\{i \in N : a_i^* = c\}| - \mathbf{1}_{a_i^*=c} \geq 1$.
15 Now construct s^* as follows. Let $s_i^*(\emptyset) = s_i^*(a^*) = a_i^*$. For all other partial histories let
16 $s_i^*(h) = d$. First, notice that for the partial histories \emptyset and (a^*) we have that $s^*(h) = a^*$ while
17 $a(s^*|h) = a^*$. Secondly, notice that for all other partial histories, we have $s^*(h) = (d)_{i \in N}$
18 while $a(s^*|h) = (d)_{i \in N}$. Concluding the condition for a no-babbling is satisfied. All that is
19 left to show is that s^* is a Subgame Perfect Equilibrium of the negotiation game. Suppose
20 not, there is some partial history $h \in H$ such that there is some other strategy $s_i \in S_i$ such
21 that $U_i(s_i, s_{-i}^*|h) > U_i(s^*|h)$. There are two possible cases.

22 1. The first possibility is that $h = \emptyset$ or (a^*) . Notice for a deviation to be profitable
23 it must be such that $a(s_i, s_{-i}^*|h) \neq a^*$, as otherwise, a strict inequality cannot hold.
24 Given this, the strategy s_i, s_{-i}^* must induce a history $h' \neq (a^*, a^*)$. Therefore, it must
25 be that all other players choose d for all periods other than the first. There are three
26 possibilities.

27 (a) Firstly, it may be that the strategy s_i, s_{-i}^* induces a terminal history, $z \in Z'$, with
28 the agreement $(d)_{i \in N}$. This induces a payoff of 1 for player i , while $U_i(s^*|\emptyset) =$
29 $u_i(a^*) \geq 1 = U_i(s_i, s_{-i}^*|\emptyset)$, therefore this cannot be profitable.

30 (b) It may be that the strategy induces a terminal history, $z \in Z'$, with the agreement
31 $((d)_{j \neq i}, c)$. However, this leads to a payoff of 0 for player i , while $U_i(s^*|\emptyset) =$
32 $u_i(a^*) \geq 1 > 0 = U_i(s_i, s_{-i}^*|\emptyset)$, therefore this cannot be profitable.

33 (c) It may be that s_i, s_{-i}^* induces a terminal history, z such that d is played by all
34 other players in all but the first period and no agreement is made, i.e. $z \in Z''$.

As no agreement is made, it must be that there are no two consecutive periods where the same action profile is played by all players it must be that s_i alternates between d and c . This implies that the lim sup of utilities induces by the proposals is 1. As $U_i(s^*|\emptyset) = u_i(a^*) \geq 1 \geq U_i(s_i, s_{-i}^*|\emptyset)$, this cannot be a profitable deviation.

2. Now suppose the history is partial and such that $h \neq \emptyset$ and $h \neq a^*$. No deviation leads to the agreement $(d)_{i \in N}$, with a payoff of 1. A deviation can only lead to the three cases examined above. Given this, the logic of the previous case remains true.

In conclusion, for any a^* such that $u_i(a^*) = 1 + k|\{i \in N : a_i^* = c\}| - \mathbf{1}_{a_i^* = c} \geq 1$ holds, we can provide a Negotiated Binding Agreement that supports such a profile. Further to this, it provides some intuition behind the sufficiency proof of theorem ?? which would imply this result.

To explore this further, notice that this implies that $a^* = (d)_{i \in N}$ may be supported. Further to this, a number of action profiles that maintain contribution can be supported by a Negotiated Binding Agreement. Specifically, for some a^* such that there exists some i such that $a_i^* = c$, we have that $1 + k|\{i \in N : a_i^* = c\}| > k|\{i \in N : a_i^* = c\}|$, i.e. the number of players contributing have a strictly lower utility than those who are not. It must be unprofitable for an agent i such that $a_i^* = c$ to instead propose d for all periods and receive a payoff of 1 via the agreement of $(d)_{i \in N}$ that would be induced by s_{-i}^* . With this, we can see that any a^* supported in a Negotiated Binding Agreement must be such that $k|\{i \in N : a_i^* = c\}| \geq 1$. More succinctly, when the number of contributors is above a lower bound, $|\{i \in N : a_i^* = c\}| \geq \frac{1}{k}$, the action profile can be supported by a Negotiated Binding Agreement. As $\frac{1}{k} < n$ this implies that full cooperation can be sustained.

Finally, to show that this fully characterises the Negotiated Binding Agreement, suppose that there is some equilibrium s^* that supports some a^* such that $u_i(a^*) < 1$ for some $i \in N$. For this to be the case it must be that $u_i(s^*|\emptyset) = a^*$. Now consider a deviation of $i \in N$ such that $s_i(h') = d$ for all histories $h' \in H$ at $h = \emptyset$. Such a deviation ensures that in any terminal history, the payoff is pinned down by $u_i(d, a_{-i})$, be that if the history ends in agreement or not. If it does not end in agreement, it is pinned down by between some $u_i(d, a_{-i})$ with $a_{-i} \in \{c, d\}^{n-1}$. However, $u_i(d, a_{-i}) \geq 1$ for all possible $a_{-i} \in \{c, d\}^{n-1}$. Therefore $U_i(s_i, s_{-i}^*|\emptyset) \geq 1 > U_i(s^*|\emptyset)$. Therefore it cannot be that s^* is a Subgame Perfect Equilibrium of the negotiation game and therefore cannot be a Negotiated Binding Agreement. ▼

1 B.2. Appendix: Robustness

2 In this section, I outline how the results of this paper are robust to changes in how the
 3 negotiation game is defined. I do so as follows. In subsection B.2.1. I show that necessarily
 4 proposals can only be made from actions that survive iterated elimination of absolutely
 5 dominated actions, which are tightly related to those that survive iterated elimination of
 6 individually irrational actions, and the sufficient conditions for agreement outcomes hold if
 7 agents make proposals sequentially rather than simultaneously in each period. In subsection
 8 B.2.2. I show that, if the payoffs of the infinite histories are appropriately defined, both
 9 the necessary and sufficient conditions for agreement outcomes hold if agents may make
 10 proposals of the joint action, rather than just their own, in each period. In subsection
 11 B.2.3., I show that the sufficient conditions for agreement outcomes remain to be true in a
 12 model where the payoff of the infinitely terminal histories are taken to be worse than the
 13 payoff of any finite terminal history.

14 In essence, these robustness checks show how the drivers of the results. Specifically, that
 15 agents cannot use a non-agreement outcome as a threat of deviating, whereas timing and
 16 the proposals used are not an important for driving the results.

17 B.2.1. Robustness to Order of Proposals

18 As in section 2, let G be an underlying game with bounded payoffs.

19 Define the negotiation game with order as follows.

20 Let $\mathcal{O} : N \rightarrow |N|$ be the order in which agents make proposals within a period. Note
 21 that this function may not be one-to-one, and therefore it may be that many agents make
 22 the proposals at the same time. Assume that if $\mathcal{O}(i) = k > 1$ then $\exists j \in N$ such that
 23 $\mathcal{O}(j) = k - 1$. That is, \mathcal{O} naturally defines an order: if I am not first, then there must be
 24 someone who proposes before me. I also assume that $\mathcal{O}(i) = 1$ for some $i \in N$ to ensure the
 25 first proposer is labelled as such. Let $\mathcal{O}^{-1}(k) = \{i \in N | \mathcal{O}(i) = k\}$, that is, define $\mathcal{O}^{-1}(k)$ is
 26 the set of agents who make the k^{th} proposal.

A history will be the empty set or a sequence of proposals for all agents followed by the
 first k proposals within the last period. That is,

$$h = (a^1, a^2, \dots, a^{k-1}, (a_{\mathcal{O}^{-1}(2)}^k, a_{\mathcal{O}^{-1}(1)}^k, \dots, a_{\mathcal{O}^{-1}(l)}^k))$$

27 , with $l \leq n$, i.e. there may be agents who are yet to make a proposal within the current
 28 period.

29 A history is terminal if, either:

a) Where the same action profile is proposed twice in consecutive periods, all agents have made a proposal within the last period, and no earlier occurrence of consecutive repetition is present. That is, $z = (a^1, \dots, a^{k-1}, a^k)$ is terminal if $a^k = a^{k-1}$ and $a^m \neq a^{m-1}$ for all $m < k$. Let the set of such histories be denoted by \tilde{Z}' and refer to these histories as ones where an *agreement* is made.

b) an infinite sequence where the same action profile is never proposed consecutively, and all agents have made a proposal within each period. Let the set of such histories be denoted by \tilde{Z}'' . I will again refer to these as perpetual disagreement histories.

Let the set of terminal histories be given by $\tilde{Z} = \tilde{Z}' \cup \tilde{Z}''$. The set of all possible histories is all terminal histories and all finite histories where there are no consecutive proposals that are the same action for all agents. Let the set of partial histories be denoted by \tilde{H} .

As before, whenever $z = (a^1, \dots, a^k) \in \tilde{Z}'$ let $U_i(z) = u_i(a^k)$.

Whenever $z \in \tilde{Z}''$ let $U_i(z) \in [\liminf_{t \rightarrow \infty} u_i(a^t), \limsup_{t \rightarrow \infty} u_i(a^t)]$. Only take these definitions over well-defined action profiles.

Let \tilde{H}_i be the set of partial histories where $i \in N$ is active. That is $h \in \tilde{H}_i$ is such that $h = (a^1, a^2, \dots, a^{k-1}, (a_{\mathcal{O}^{-1}(1)}^k, \dots, a_{\mathcal{O}(i)-1}^k))$ when $\mathcal{O}(i) \neq 1$ and $h = (a^1, a^2, \dots, a^{k-1}, a^k)$. the strategy of $i \in N$ dictates the proposal i would make at any history for which they are active: $s_i : \tilde{H}_i \rightarrow A_i$. Let S_i be the space of all such mappings.

For a partial history $h \in \tilde{H}$, let $U_i(s|h)$ denote the payoff that would be received from the terminal history that the strategy s would induce, starting from the history $h \in \tilde{H}$. I will refer to such a history as $(s|h)$. When $z \in \tilde{Z}'$, i.e. an agreement is made, let $a(h)$ as the action profile that terminates z .

I define Subgame Perfect Equilibrium for this model here:

Definition (Subgame Perfect Equilibrium). s^* is Subgame Perfect Equilibrium, if for all $i \in N$, for all partial histories where $i \in N$ is active $h \in \tilde{H}_i$, $U_i(s^*|h) \geq U_i(s_i, s_{-i}^*|h)$, for all $s_i \in S_i$.

This leads to the natural definition of Negotiated Binding Agreement in this setting. To make the distinction clear, I refer to this as Negotiated Binding Agreement with order.

Definition 10 (Negotiated Binding Agreement with Order). s^* is a Negotiated Binding Agreement with order \mathcal{O} supporting $a^* = a * (s^*|\emptyset)$ if:

a) s^* is a Subgame Perfect Equilibrium.

b) For all $h \in \tilde{H}_i \exists h' \in \tilde{H}_i$ such that $s_i(h) = a_i(s^*|h')$.

Now I show that some necessary conditions related in section ?? remain to be true for this specification of the model. First, I will show that any action proposed at some history of Negotiated Binding Agreement must survive a procedure of iterated deletion of absolutely dominated actions, also known as interdependent choice rationalizability (Salcedo, 2017) and min-max rationalizability (Halpern and Pass, 2018).

Definition 11 (Absolute Domination given $C_{-i} \subseteq A_{-i}$). $a_i \in A_i$ is absolutely dominated given $C_{-i} \subseteq A_{-i}$ if $\exists a'_i \in A_i$ such that

$$\inf_{a_{-i} \in C_{-i}} u_i(a'_i, a_{-i}) > \sup_{a_{-i} \in C_{-i}} u_i(a_i, a_{-i})$$

Denote the set of absolutely dominated actions given C_{-i} by $D_i(C_{-i})$.

As I do not require that the utility functions are continuous and defined over a compact set, the minimum or maximum may not exist. With this, I take the supremum and infimum, which by the assumption that the utility function is bounded are always well-defined. Bar this change, the above definition is equivalent to that of Salcedo (2017). Note that, if in a normal form game, there is a single action that is not absolutely dominated given A_{-i} , then this action is an obviously dominant strategy as defined by Li (2017).

Definition 12 (Iterated Elimination of Absolutely Dominated Actions). Let $\tilde{A}_i^0 = A_i$ for all $i \in N$. Let $\tilde{A}_{-i}^0 = A_{-i}$. Then for all $m > 0$ let $\tilde{A}_i^m = \tilde{A}_i^{m-1} \setminus D_i(\tilde{A}_{-i}^{m-1})$ where $\tilde{A}_{-i}^{m-1} = \times_{j \neq i} \tilde{A}_j^{m-1}$.

The set of actions that survives Iterated Elimination of Absolutely Dominated Actions (IAD) for i is given by $IAD_i = \bigcap_{m \geq 0} \tilde{A}_i^m$. Let $IAD = \times_{i \in N} IAD_i$.

Note that if at each level of iteration if the min-max and max-min payoffs are the same then IAD coincides with IIR . Note that typically, the concept of iterated elimination of individually irrational actions and iterated elimination of absolutely dominated actions are different, for instance, consider the following example.

Example 6. Consider the following underlying two-player game.

$1 \setminus 2$	L	R
U	1, 2	-1, 0.5
M	-1, 1	1, 0.5
D	-0.7, 3	-0.7, 3

22

Here, in iterated elimination of absolutely dominated actions, all profiles survive. However, if we consider iterated elimination of individually irrational actions, we may remove D , as the min-max payoff for player 1 is 1. Given this, we may also eliminate R for player

25

2, as her min-max payoff is 0.5. Finally, we remove M , therefore we conclude that iterated elimination of individually rational actions leads to the unique prediction of (U, L) , while iterated elimination of absolutely dominated actions allows for any action profile. ▼

These definitions lead to the following proposition; any proposal made on or off the path in a Negotiated Binding Agreement with order must come from the set of actions that survives iterated elimination of absolutely dominated actions.

Proposition 2. *For any order \mathcal{O} , if s^* is a Negotiated Binding Agreement with order then, for all histories where i is active $h \in \tilde{H}_i$, $s_i(h) \in IAD_i$.*

I reserve this proof, and all other proofs within this appendix, for the appendix B.3..

Further to this, the following proposition shows that the sufficient conditions for agreement outcomes are relevant within this specification of the model. Indeed, further to this, any outcome that can be sustained with a Negotiated Binding Agreement can be sustained within a model of negotiation with order, no matter the order. This is highlighted by the following proposition.

Proposition 3. *Take any order \mathcal{O} . If a^* is supported in a Negotiated Binding Agreement then it is supported in Negotiated Binding Agreement with order \mathcal{O} .*

In essence, this shows that the qualitative results of having agreements be based on player specific punishments and all agreement outcomes having to satisfy some iterative individual rationality constraint are robust to having sequential proposals. Rather, the structure of the terminal histories, and the associated payoffs, as well as the ability of all agents to make some proposal, are the key features of the model. Within the next sub-appendix, I go on to show that when the payoffs of infinite histories are correctly specified, the robustness of these results also holds when agents propose the action profile, rather than only their action. This further highlights this point.

B.2.2. Robustness to Joint Proposals

As in section 2, let G be an underlying game with bounded payoffs.

Define the negotiation game with all proposals as follows.

A history will be the empty set or a sequence of proposals for all agents, where each agent may propose a joint action profile. That is,

$$h = ((a^{1,1}, a^{2,1}, \dots, a^{n,1}), (a^{1,2}, a^{2,2}, \dots, a^{n,2}), \dots, (a^{1,k}, a^{2,k}, \dots, a^{n,k}))$$

, where $a^{i,t} \in A$. With some abuse of notation, let $a^t = (a^{1,t}, a^{2,t}, \dots, a^{n,t})$.

1 A history is terminal if, either:

2 a) Where the same action profile is proposed twice in consecutive periods by all agents
 3 and no earlier occurrence of consecutive repetition is present. That is, $h = (a^1, \dots, a^{k-1}, a^k)$
 4 is terminal if $a^k = a^{k-1}$, $a^{i,k} = a^{j,k}$ for all $i, j \in N$, and either $a^m \neq a^{m-1}$ for all
 5 $m < k$ or $a^{i,m} \neq a^{j,m}$ for some $i, j \in N$. Let the set of such histories be denoted by
 6 \tilde{Z}' and refer to these histories as ones where an *agreement* is made.

7 b) an infinite sequence where the same action profile for all agents is never proposed
 8 consecutively. Let the set of such histories be denoted by \tilde{Z}'' . Refer to these as
 9 perpetual disagreement histories.

10 Let the set of terminal histories be given by $\tilde{Z} = \tilde{Z}' \cup \tilde{Z}''$. The set of all possible histories is
 11 all terminal histories and all finite histories where there are no consecutive proposals that
 12 are the same action profile for all agents. Let the set of all partial histories given by \tilde{H} .

Whenever $z = (a^1, \dots, a^k, \dots) \in \tilde{Z}'$, let

$$\tilde{h} = ((a_i^{i,1})_{i \in N}, (a_i^{i,2})_{i \in N}, \dots, (a_i^{i,k})_{i \in N}, \dots)$$

13 , i.e., take the proposals that each agent makes for themselves. Let this sequence be denoted
 14 by $\tilde{z} = (\tilde{a}^1, \tilde{a}^2, \dots, \tilde{a}^k, \dots)$. Let the lower bound of the utility of $z = (a^1, \dots, a^k, \dots) \in \tilde{Z}'$ be
 15 given by $\liminf_{t \rightarrow \infty} u_i(\tilde{a}^t)$ and an upper bound of the $\limsup_{t \rightarrow \infty} u_i(\tilde{a}^t)$. This implies that
 16 if no agreement is made, then only your own proposals matter, you cannot impact what
 17 others do in this case.

18 the strategy of $i \in N$ dictates the proposal i would make when they are active: $s_i :$
 19 $\tilde{H} \rightarrow A$. Let S_i be the space of all such mappings.

20 With some abuse of notation, for a partial history $h \in \tilde{H}$, let $U_i(s|h)$ denote the payoff
 21 that would be received from the terminal history that the strategy s would induce, starting
 22 from the history $h \in \tilde{H}$. I will again refer to such a history as $(s|h)$. As before, when $z \in \tilde{Z}'$,
 23 i.e. an agreement is made, let $a(h)$ as the action profile that terminates z .

24 **Definition** (Subgame Perfect Equilibrium). s^* is Subgame Perfect Equilibrium, if for all
 25 $i \in N$, for all partial histories $h \in \tilde{H}$, for all $i \in N$, $U_i(s^*|h) \geq U_i(s_i, s_{-i}^*|h)$, for all $s_i \in S_i$.

26 This leads to the definition of Negotiated Binding Agreement in this setting. To make
 27 the distinction clear, I refer to this as Negotiated Binding Agreement with all proposals.

28 **Definition 13** (Negotiated Binding Agreement with all Proposals). s^* is a Negotiated
 29 Binding Agreement with all proposals supporting $a^* = a(s|\emptyset)$ if:

30 a) s^* is a Subgame Perfect Equilibrium.

b) $\forall h \in \tilde{H} \exists h' \in \tilde{H}$ such that $s_i(h) = a(s^*|h)$.

As before, the following proposition shows that the necessary conditions previously shown for Negotiated Binding Agreement hold for this specification of the model.

Proposition 4. *If s^* is a Negotiated Binding Agreement with all proposals, for all histories $h \in \tilde{H}$, $s_i(h) \in IIR_i$.*

Further, for any negotiated with order s^* be such that, for any history $h \in \tilde{H}$, $U_i(s^*|h) \geq \underline{u}_i$ where

$$\underline{u}_i = \inf_{a_{-i} \in IIR_{-i}} \sup_{a_i \in A_i} u_i(a_i, a_{-i})$$

Further to this, the sufficient conditions for Negotiated Binding Agreement outcomes are also sufficient for Negotiated Binding Agreements with all Proposals. This is captured by the following proposition, which shows us that any Negotiated Binding Agreement can be replicated by a Negotiated Binding Agreement with all proposals.

Proposition 5. *a^* is supported by a Negotiated Binding Agreement with all proposals if a^* is supported by a Negotiated Binding Agreement.*

This again highlights the important features and drivers of the results of the model. In essence, it is the ability of agents to make a meaningful impact on their payoff via their proposals, while ensuring they do not force other agents to take some action. This is highlighted by the idea that the payoffs of infinite terminal histories, i.e. when there is no agreement, take the actions for individuals that they propose for themselves.

B.2.3. Robustness to Outside Options

Within this sub-appendix, I take the model to be exactly as in section ???. That is, agents simultaneously propose the action that they will take. The only caveat is that whenever a terminal history is infinite they receive a payoff that is worse than the payoff within the underlying game. That is, when $z \in Z''$ let $U_i(z) = \inf_{a \in A} u_i(a)$. Negotiated Binding Agreement can be defined as before. To distinguish between these cases I will refer to Negotiated Binding Agreement for the model in this sub-appendix as *constant outside option Negotiated Binding Agreement*. In this setting, it is no longer true that the necessary conditions for agreement outcomes hold. However, the sufficient conditions for agreement outcomes remain to be valid. This is highlighted by the following proposition.

Proposition 6. *If s^* is a Negotiated Binding Agreement then s^* is a constant outside option Negotiated Binding Agreement.*

As Negotiated Binding Agreements do need not make use of the infinitely long terminal histories as part of equilibrium, this result shows us that they are important only for restricting deviations. That is, if we were to make such an option worse for each player, they have less incentive to deviate than before. Therefore Negotiated Binding Agreement captures a set of strategies and outcomes that work regardless of whether the outside option is specified as within this chapter or normalised to be worse than any outcome as typically assumed in bargaining games.

B.3. Proofs for Appendix B.2.

Proof of proposition 2: By induction. Firstly, note that $s_i(h) = a_i \notin D_i(A_{-i})$ for all $h \in \tilde{H}_i$. To see this suppose by contradiction it is not the case. Then $s_i^*(h) = a_i \in D_i(A_{-i})$ for some $i \in N$ and some history $h \in \tilde{H}_i$. It must be that $a_i(s^*|h) = a_i$ for some $h' \in H$. Given this, $U_i(s^*|h) \leq \sup_{a_{-i} \in A_{-i}} u_i(a_i, a_{-i})$. Further, s^* is a Subgame Perfect Equilibrium of the negotiation game, and therefore there is no profitable deviation for s_i^* at any history for which i is active, including h' . Notice that as $a_i \in D_i(A_{-i})$ then $\exists a'_i \in A_i$ such that $\inf_{a'_{-i} \in A_{-i}} u_i(a'_i, a'_{-i}) > \sup_{a_{-i} \in A_{-i}} u_i(a_i, a_{-i})$. Now consider a strategy s'_i such that $s'_i(h'') = a'_i$ for all h'' for which i is active. Notice that, by construction of s'_i , the history $(s'_i, s_{-i}^*|h')$ must either terminate in a'_i or be such that only action profiles with a'_i appear after h . In either case, we can conclude that $U_i(s'_i, s_{-i}^*|h') \geq \inf_{a'_{-i} \in A_{-i}} u_i(a'_i, a'_{-i}) > \sup_{a_{-i} \in A_{-i}} u_i(a_i, a_{-i}) \geq U_i(s^*|h')$. A contradiction to s^* be a Subgame Perfect Equilibrium of the negotiation game.

By the inductive hypothesis, suppose that $s_i^*(h) \in \tilde{A}_i^m$ for all $h \in \tilde{H}_i$ and $i \in N$. Now suppose by contradiction that $s_i^*(h) = a_i \in D_i(\tilde{A}_{-i}^m)$. It must be that $a_i(s^*|h') = a_i$. Given this, $U_i(s^*|h') \leq \sup_{a_{-i} \in \tilde{A}_{-i}^m} u_i(a_i, a_{-i})$, as $s_{-i}^*(h'') \in \tilde{A}_{-i}^m$ for all $h'' \in \tilde{H}_i$. Further, s^* is a Subgame Perfect Equilibrium of the negotiation game, and therefore there is no profitable deviation for s_i^* at any history, including h' . Notice that as $a_i \in D_i(\tilde{A}_{-i}^m)$ then $\exists a'_i \in A_i$ such that $\inf_{a'_{-i} \in \tilde{A}_{-i}^m} u_i(a'_i, a'_{-i}) > \sup_{a_{-i} \in \tilde{A}_{-i}^m} u_i(a_i, a_{-i})$. Now consider a strategy s'_i such that $s'_i(h'') = a'_i$ for all $h'' \in \tilde{H}_i$. Notice that, by definition and construction of s'_i $U_i(s'_i, s_{-i}|h')$ must only be constructed using the utility of $u_i(a'_i, \cdot)$, as either $(s'_i, s_{-i}|h') \in Z'$, in which case it must terminate in a'_i by definition, or $(s'_i, s_{-i}|h') \in Z''$, in which case all histories following h' use only a'_i . Further, as $s_{-i}^*(h'') \in \tilde{A}_{-i}^m$ that from this history on the only action profiles proposed are a'_i, a'_{-i} such that $a'_{-i} \in \tilde{A}_{-i}^m$. Given this, we can conclude that $U_i(s'_i, s_{-i}|h') \geq \inf_{a'_{-i} \in \tilde{A}_{-i}^m} u_i(a'_i, a'_{-i}) > \sup_{a_{-i} \in \tilde{A}_{-i}^m} u_i(a_i, a_{-i}) \geq U_i(s^*|h')$. A contradiction to s^* being a Subgame Perfect Equilibrium of the negotiation game. ■

Proof of proposition 3: We will show that if a^* is sustained in a Negotiated Binding Agreement then it can be sustained in a Negotiated Binding Agreement with order \mathcal{O} for

1 any order. Take any order \mathcal{O} . Take s^* that sustains a^* in a Negotiated Binding Agreement.
 2 Let $s'_i : \tilde{H}_i \rightarrow A_i$ such that, for all $h \in \tilde{H}_i$ such that $h = (h', (a_{\mathcal{O}^{-1}(1)}, \dots, a_{\mathcal{O}^{-1}(i)-1}))$ we
 3 have that $s'_i(h) = s_i^*(h')$. First note that $a(s'|\emptyset) = a^*$ and $a(s'|h') = a(s^*|h)$ whenever
 4 $h' = h$ while $h' \in \tilde{H}$ and $h \in H$. Next we will show that s' is subgame perfect of the
 5 negotiation game. Suppose not, there is some $i \in N$ for which there exists some $h \in H'_i$
 6 and some $s''_i \in S_i$ such that $U_i(s''_i, s'_{-i}|h) > U_i(s'|h)$. However, given agents are rational
 7 and the structure of s' , they can replicate any deviation from s'_i with a deviation from s_i^* .
 8 With this, we must conclude that s_i^* is not subgame perfect of the negotiation game. A
 9 contradiction. Concluding that s' is a Negotiated Binding Agreement with order \mathcal{O} , leading
 10 to the outcome a^* . ■

11 **Proof of proposition 4:** By induction. Firstly, note that $s_i(h) = a \notin D_i(A)$ for all
 12 $h \in \tilde{H}$. Suppose by contradiction it is the case. Then $s_i^*(h) = a \in D(A)$ for some $i \in N$ and
 13 some history $h \in \tilde{H}$. It must be that $a(s^*|h') = a$ for some history $h' \in H$. This implies
 14 that $[s_i^*(h')]_j = a_j \in D_j(A_{-j})$ for some j . Given this, $U_j(s^*|h) \leq \sup_{a_{-j} \in A_{-j}} u_j(a_j, a_{-j})$.
 15 Further, s^* is a Subgame Perfect Equilibrium of the negotiation game, and therefore there
 16 is no profitable deviation for s_j^* at any history, including h . Notice that as $a_j \in D_j(A_{-j})$
 17 then, for all $\epsilon > 0 \exists a'_j : A_{-j} \rightarrow A_j$ such that $u_i(a'_j(a_{-i}), a'_{-j}) > \sup_{a_{-j} \in A_{-j}} u_i(a_j, a_{-j}) - \epsilon$.
 18 Now consider a strategy s'_j such that $s'_j(h'') = (a'_j(s_{-j}(h'')), a''_{-j})$, for some $a''_{-j} \in A_{-j}$ for
 19 all h'' . Notice that, by construction of s'_i , the history $(s'_j, s_{-j}^*|h')$ must either terminate in a'_j
 20 or be such that only action profiles with a'_j appear after h' . In either case, we can conclude
 21 that $U_j(s'_j, s_{-j}^*|h') \geq \inf_{a'_{-j} \in A_{-j}} u_i(a'_j(a'_{-j}), a'_{-j}) > \sup_{a_{-j} \in A_{-j}} u_i(a_j, a_{-j}) \geq U_j(s^*|h')$. A
 22 contradiction to s^* be a Subgame Perfect Equilibrium of the negotiation game.

23 By the inductive hypothesis, suppose that $s_i^*(h) \in \tilde{A}^m$ for all $h \in \tilde{H}_i$ and $i \in N$. Now
 24 suppose by contradiction that $s_i^*(h) = a \in D(\tilde{A}^m)$. It must be that $a(s^*|h') = a$ for some
 25 history $h' \in H$. Further, for some $j \in N$ $a_j \in D(\tilde{A}^m)$. Without loss of generality let
 26 $j = i$. Given this, $U_i(s^*|h') \leq \sup_{a_{-i} \in \tilde{A}_{-i}^m} u_i(a_i, a_{-i})$, as $s_{-i}^*(h'') \in \times_{j \neq i} \tilde{A}^m$ for all $h'' \in \tilde{H}$.
 27 Further, s^* is a Subgame Perfect Equilibrium of the negotiation game, and therefore there is
 28 no profitable deviation for s_i^* at any history, including h . Notice that as $a_j \in D_j(A_{-j})$ then,
 29 for all $\epsilon > 0 \exists a'_j : \tilde{A}_{-i}^m \rightarrow A_j$ such that $u_i(a'_j(a_{-i}), a'_{-j}) > \sup_{a_{-j} \in \tilde{A}_{-j}^m} u_i(a_j, a_{-j}) - \epsilon$. Now
 30 consider a strategy s'_i such that $s'_i(h'') = (a'_i(s_{-i}^*(h'')), a''_{-i})$, with $a''_{-i} \notin A_{-i}$ for all $h'' \in \tilde{H}_i$.
 31 Notice that, by definition and construction of s'_i $U_i(s'_i, s_{-i}^*|h')$ must only be constructed
 32 using the utility of $u_i(a'_i, \cdot)$, as with the before logic, we can only terminate in histories that
 33 have a'_i infinitely repeated or an agreement is reached with a'_i . Given this, we can conclude
 34 that $U_i(s'_i, s_{-i}^*|h') \geq \inf_{a'_{-i} \in \tilde{A}_{-i}^m} u_i(a'_i(a'_{-i}), a'_{-i}) > \sup_{a_{-i} \in \tilde{A}_{-i}^m} u_i(a_i, a_{-i}) \geq U_i(s^*|h')$. A
 35 contradiction to s^* be a Subgame Perfect Equilibrium of the negotiation game.

As proposals are simultaneous, the logic of showing that for any negotiated with order

s^* be such that, for any history $h \in \tilde{H}$, $U_i(s^*|h) \geq \underline{u}_i$ where

$$\underline{u}_i = \inf_{a_{-i} \in IIR_{-i}} \sup_{a_i \in A_i} u_i(a_i, a_{-i})$$

1 is identical to theorem 4, where s'_i is selected to intentionally cause perpetual disagreement.
 2 ■

3 **Proof of proposition 5:** If a^* is supported by a Negotiated Binding Agreement then
 4 a^* is supported by a all proposal Negotiated Binding Agreement. Take s^* that supports a^*
 5 in a Negotiated Binding Agreement. Construct $s'_i : \tilde{H} \rightarrow A$ as follows. Let $s'_i(h'') = s^*(\tilde{h}'')$,
 6 where \tilde{h}'' is as defined to define payoffs of infinite histories. Clearly if s_i^* is optimal so is s'_i as
 7 a deviation to a partial infinite history leads to the same payoff that could be achieved under
 8 s_{-i}^* . A deviation to another terminal history must be such that it could not be achieved
 9 under a deviation from s_i^* . However, by definition of s'_i , this cannot be the case. ■

10 **Proof of proposition 6:** As s^* is a Negotiated Binding Agreement it must be that s^* is
 11 a Subgame Perfect Equilibrium of the negotiation game with the terminal infinite histories
 12 giving a payment as defined in section 2. As s^* never dictates that a history should be
 13 infinite and terminal, it follows that there is no profitable deviation where the outcome
 14 leads to a deterministic outcome. It follows that the payoff on the path remains the same
 15 when the model of a constant outside option is taken. Finally, as there is no profitable
 16 deviation when the deviation would induce a terminal infinite history when the payoff is
 17 defined as in section 2, there cannot be a profitable deviation when the constant outside
 18 option is taken. Therefore s^* is a constant outside option Negotiated Binding Agreement.
 19 ■