

# PLEDGE-AND-REVIEW BARGAINING: FROM KYOTO TO PARIS\*

BÅRD HARSTAD

bard.harstad@econ.uio.no

June 2020

## Abstract

To improve on global climate policies, we must better understand past and present treaties. The agreements signed in Kyoto (1997) and Paris (2015) are associated with different bargaining procedures, participation levels, and enforcement mechanisms. The choice of procedure changed over time, and yet the commitment period length was and is five years. To shed light on these five facts, this paper presents a dynamic game with endogenous technology, participation, enforcement, contract terms, and bargaining procedure. The theoretical predictions are consistent with all five facts.

Keywords: Dynamic games, bargaining games, coalitions, climate change, Paris Agreement, Kyoto Protocol.

JEL codes: F55, H87, Q54.

---

\*I thank Scott Barrett, Ernesto Dal Bo, Faruk Gul, Jon Hovi, Steffen Lippert, Paolo Piacquadio, Santiago Rubio, Leo Simon, Håkon Sælen, David Victor, Christina Voigt, Joel Watson, and audiences in Adelaide (AARES pre-conference), U. Autònoma de Barcelona, U. of Barcelona, the BEET workshop at BI, UC Berkeley, UC3M, CEU, UC San Diego, U. of Chicago, CREST-Ecole Polytechnique, EIEF, ESEM 2018, University of Essex, HEC Paris, Hong Kong Baptist University, Ifo institute, London School of Economics, Manchester University, University of Melbourne, MIT, National Taiwan University, National University of Singapore, Northwestern University, University of Oslo, UPF, Princeton University, Queen Mary University, Singapore Management University, Stanford GSB, SURED 2018, Toulouse School of Economics, the 2019 Wallis conference, and WCERE 2018. Marie Karlsen and Johannes Hveem Alsvik provided excellent research assistance and Frank Azevedo helped with the editing.

- *The pledge-and-review strategy is completely inadequate.*

Christian Gollier and Jean Tirole  
The Economist, guest blog, 2015

Pledge-and-review (P&R) refers to the structure of the negotiations associated with the Paris Agreement, December 2015. Before the countries were expected to sign the climate agreement, each party was asked to submit an intended nationally determined contribution (NDC). For most developed countries, the NDC specified unconditional cuts in the emissions of greenhouse gases being effective from 2020 to 2025 (or to 2030). As an illustration, Table 1 presents the pledges for a sample of developed countries.<sup>1</sup>

<b>Party</b>	Australia	Canada	EU	New Zealand	Norway	Russia	USA
<b>Pledge</b>	26-28%	30%	40%	30%	40%	25-30%	26-28%

*Table 1. The pledges specify emission cuts relative to nationally chosen baselines.*

Every five years the parties shall review and make new pledges for another five-year period (Paris Agreement Art. 4.9).

This procedure is remarkably different from the one used for the Kyoto Protocol of 1997. There, a "top-down" approach was used to pressure governments to cut emissions by (on average) five percent relative to the 1990 levels. Bodansky and Rajamani (2018:23) find that: "In essence, the Kyoto Protocol was the product of mutual concessions [and, as a consequence] USA accepted a much stronger target (minus 7% from 1990 levels) than it had wanted..." The conditionality of concessions in Kyoto makes it unsurprising that most papers have associated the bargaining outcome with the Nash Bargaining Solution (NBS); see the literature review below.

By comparison, P&R has been referred to as a "bottom-up" approach since countries themselves determine how much to cut nationally, without making these cuts conditional on other countries' emissions cuts. According to the Paris Agreement (Art. 4.2): "Each

---

<sup>1</sup>The baseline year is 1990 for the European Union, Russia, and Switzerland, while it is 2005 for Australia, Canada, New Zealand, and the United States. Article 4.4 of the Paris Agreement encourages "economy-wide absolute emission reduction targets," although several developing countries state pledges in terms of emission per GDP and some of these are conditional on receiving transfers. The official list is at <http://www4.unfccc.int/ndcregistry> but for an overview see <http://cait.wri.org/indc/#/>.

Party shall prepare, communicate and maintain successive nationally determined contributions that it intends to achieve." Based on this, Victor (2015) observes: "Now, instead of setting commitments through centralized bargaining, the Paris approach sets countries free to make their own commitments" (Victor, 2015). *The New York Times* (Nov. 28, 2015) wrote that: "Instead of pursuing a [Kyoto-style] top-down agreement with mandated targets, [the organizers] have asked every country to submit a national plan that lays out how and by how much they plan to reduce emissions in the years ahead."

Several observers and scholars fear that P&R is unable to deliver as ambitious targets as would a top-down approach. Keohane and Oppenheimer (2016:142) predict that: "Many governments will be tempted to use the vagueness of the Paris Agreement, the discretion that it permits, to limit the scope or intensity of their proposed actions." Tirole (2017:209) adds: "The strategy of voluntary commitments has several significant defects, and is an inadequate response to the climate change challenge." Scholars and observers naturally expect a party's contribution to reflect that party's own interests to a larger extent, and the interests of other parties to a lower extent. This conjecture is confirmed in the formalization of P&R analyzed in the Harstad (2020a), which shows that each contribution maximizes an *asymmetric* Nash product where the weight on others' payoffs is smaller than in the NBS.

*Facts on Paris vs. Kyoto.*—The difference in bargaining procedure will be referred to as "Fact 1" on how Paris and Kyoto differ. There are several other remarkable differences. As a "Fact 2," the emission targets under Kyoto were binding for only 37 countries, while nearly every country in the world contributes to the Paris Agreement. For reasons that need to be understood, the Kyoto Protocol's procedure was endogenously chosen in the 1990s, but the participants opted for the P&R procedure in the 2010s: The different choice will be referred to as "Fact 3." Despite all the differences, "Fact 4" is that the commitment period length agreed on is the same – five years – for both treaties (see Section 4). Finally, "Fact 5" is that the commitments under the Kyoto Protocol were referred to as being "legally binding," but the NDCs are, explicitly, not (see Section 5).

It is, of course, possible to find alternative theories for one or two facts in isolation (see the below literature review). The challenge is to develop a coherent framework that is consistent with *all five facts* at the same time. This is extremely important, since we

need a deep understanding of existing treaties before we can outline improvements.

The ambition of this paper is to provide a framework that is consistent with all five facts. The game below is dynamic and every country emits as well as invests in renewables in every period. The goal is to explain, rather than to be as general as possible, and for this purpose a linear-quadratic model turns out to be sufficient. Furthermore, it is not necessary to take advantage of the vast number of subgame-perfect equilibria permitted in dynamic games; it will be sufficient to limit attention to Markov-perfect equilibria (MPEs). The framework accounts for the five facts in the following way.

(1) *The Bargaining Game*.—When the alternative bargaining outcomes are embedded in the dynamic climate policy game, I find that with P&R, relative to the NBS, the pledged emission cuts are smaller, investments in renewables are smaller, and everyone’s payoff is therefore lower. This finding provides a foundation for the skepticism and the criticism of the P&R procedure, discussed above.

(2) *Participation*.—The negative result is reversed, however, when we account for free riding and let the decision to participate in the bargaining game be voluntary. Since the P&R procedure permits a party to place a lower weight on the interest of others, it is not that costly for a party to participate, and this explains why the equilibrium coalition size is larger with P&R than with the NBS (i.e., Fact 2, discussed above.) The larger coalition size means, in equilibrium, that the *sum* of contributions is larger with P&R, the aggregate investments are larger, and so is welfare.

(3) *Choice of Procedure*.—The comparison of bargaining procedures is more interesting when we take into account that there is an upper boundary ( $\bar{n}$ ) for the number of potential members and that this constraint might bind. Furthermore, when the parties are heterogeneous, or with minimum participation constraints, a number ( $\underline{n}$ ) of them may participate regardless of the procedure. Then, the narrow-but-deep agreement under the NBS can be more attractive than the shallow-but-broad outcome under P&R. When these constraints are accounted for, P&R is preferred if and only if  $\bar{n}$  is large and/or  $\underline{n}$  is small, I show.

This result is in line with the development from Kyoto to Paris: In the 1990s, there were a large number of developing countries that could not be expected to contribute much to a global climate policy. Over the last twenty years, some of these have become

emerging economies that potentially have important roles to play. The number of relevant potential parties,  $\bar{n}$ , has therefore increased. During the same period, seven countries that initially signed the Kyoto Protocol declared that they did not intend to contribute to Kyoto's second commitment period (IPCC, 2014:1025). This can be interpreted as a smaller  $\underline{n}$ . Either (or both) of these developments makes P&R relatively more attractive for every participant. Thus, the theory is consistent with Fact 3 that the parties preferred the Kyoto Protocol in the 1990s, but P&R in the 2010s. With heterogeneous countries, Section 3 concludes by explaining the disagreement between the North and the South when it came to the choice of procedure.

(4) *Terms of Contract.*—The optimal contract duration in this model results from a novel trade-off: A long-term contract is unattractive because, after the parties have invested in new capacity, it becomes optimal to negotiate still more ambitious pledges. A short-term contract, however, creates a hold-up problem when the parties anticipate how their investments will influence the next bargaining outcome. The optimal term trades off these two concerns, but this trade-off is the same under P&R as under NBS, and it is independent of the number of participants in this model. Thus, if a five-year commitment period was optimal under the Kyoto Protocol, it is indeed optimal also for the Paris Agreement, according to this result. In other words, the theory rationalizes the similarity (Fact 4) as well as the differences between the agreements.

(5) *Enforcement.*—While the analysis of (1)-(4) restricts attention to MPEs, I end the analysis by discussing how this restriction can be relaxed. At the least, we should check whether the parties have incentives to defect or comply to the equilibrium pledges. As in the repeated games literature, we may require a party to be willing to comply under the threat that others can retaliate. The P&R bargaining outcome is more likely to be self-enforcing than the NBS, I show. Intuitively, when the pledges are less demanding, or the number of cooperators is large, then the temptation to defect is small. If the bargaining outcome is characterized by the NBS, in contrast, the parties might find it necessary to motivate compliance by raising the political cost of defection. As explained in Section 5, the political cost can be raised by requiring the emissions cuts to be "legally binding" or enforced by other punitive measures – and both these methods are indeed employed by the Kyoto Protocol, but not by the Paris Agreement (Fact 5).

Since the primary goal of the analysis is to present a coherent framework consistent with Facts 1-5, I do not force the model to be more complicated than is necessary. However, some readers may react to assumptions regarding the timing of the game, the choice of policy instruments, and the tools available at the bargaining stage. To deal with these concerns, the robustness section shows that the model can be generalized in ten such directions and it explains why the results survive in all of them.

*Literature.*—Although I am unaware of other frameworks consistent with all five facts, I draw on several contributions consistent with a *subset* of the facts. The dynamic climate change game below draws on standard assumptions introduced by Dutta and Radner (2004; 2006; 2020), Harstad (2012; 2016), and Battaglini and Harstad (2016), although the chosen functional forms here are different, making the analysis below especially tractable. For instance, Battaglini and Harstad (2016) assumed that technology and green infrastructure depreciated immediately and studied neither pledges nor enforcement. More importantly, *none* of the above papers compare or derive the consequences of alternative bargaining outcomes.

Dutta and Radner (2020) and Caparrós (2020) justify an aspect of the Paris Agreement by showing how its Green Climate Fund can lead to efficiency. Thus, they complement the present paper, which instead focuses on the game, between developed countries, with neither transfers nor conditionality. Dutta and Radner rationalize several of the characteristics of the Paris Agreement, but not Facts 3-5 on the change of procedure, the commitment period length, or the change in enforcement.

The coalition formation game below is the standard one when modeling collusion (d'Aspremont et al., 1983; Bloch, 2018) and environmental coalitions (Hoel, 1992; Carraro and Siniscalco, 1993; Barrett, 1994). The typical prediction is that the coalition size is very small. This prediction is inconsistent with Fact 3 and that real-world coalitions can be quite large. This inconsistency is referred to as a "paradox" by Kolstad and Toman (2005) and Nordhaus (2015). My formalization of P&R, however, is consistent with larger coalitions.

Dynamic models can explain large coalitions (see the surveys by Calvo and Rubio, 2012; Caparrós, 2016) depending on the contractual environment (Battaglini and Harstad, 2016), the beliefs (Karp and Sakamoto, 2019), or the willingness to avoid delays (Kováč

and Schmidt, 2019). But these papers predict that agreements will be long-lasting if and only if the coalition is large, so they do not rationalize Facts 4 and 5.

The trade-off between treaties that are narrow-but-deep vs. broad-but-shallow is real and therefore well known; see Schmalensee (1998), Barrett (2002), or Aldy et al. (2003). The model by Finus and Maus (2008) is especially relevant but it is static, and thus Finus and Maus cannot rationalize Facts 4 and 5. Furthermore, the assumptions in this literature have been criticized by political scientists such as Gilligan (2004) and Bernauer et al. (2013), but the present paper shows how the trade-off arises naturally from differences in bargaining procedures. In contrast to Schmalensee (1998), who recommended "broad, then deep," the model in this paper can rationalize the reverse, factual development from the deep-but-narrow Kyoto Protocol to the broad-but-shallow Paris Agreement.

Battaglini and Harstad (2020) show that participation (Fact 2) and weak enforcement (Fact 5) can be explained by policy makers' *reelection* concerns. Again, that framework does not compare procedures (so, it cannot explain Facts 1 and 3) and it is static (so, it cannot explain Fact 4). In dynamic games, Gerlagh and Liski (2018b) rationalize participation and significant contributions, while Harstad (2020b) rationalizes technology investments—even without conditionality—when parties are troubled by time inconsistency, but these frameworks cannot rationalize Facts 3-5. More importantly, and in line with the other papers above, they do not compare the consequences of different bargaining outcomes.

The more detailed strands of literatures on contract durations and compliance are discussed in Sections 4 and 5, respectively.

*Outline.*—The next section presents a dynamic climate policy game. After the equilibrium investment levels are derived, as a function of the pledged emission cuts, I show how the continuation values can be expressed as functions that are linear-quadratic in the pledges. When these payoffs are combined with the alternative bargaining solutions, it is shown that P&R leads to lower contributions, investments, and welfare than does the NBS. However, this negative finding is reversed when participation is endogenous in Section 2. Section 3 permits heterogeneous countries and argues that the theory can explain why P&R is preferred in the 2010s although the NBS was preferred in the 1990s. While these sections take the commitment period length as given, Section 4 proves that

the optimal length is indeed the same for the two procedures. Section 5 discusses how we can relax the assumption that parties can commit, and the restriction to MPEs, and shows that the P&R pledges are more likely to be self-enforcing, while the NBS outcome (characterizing the Kyoto agreement) may need to be "legally binding". Ten generalizations are analyzed in Section 6, before Section 7 concludes. The Appendix contains all proofs.

## 1. THE MODEL AND BENCHMARK RESULTS

### 1.1. *A Dynamic Game with Contributions and Investments*

The model describes a situation in which a set of parties can contribute to a public good as well as invest in their future capacities to contribute. In equilibrium, the negotiated contribution levels will influence how much the parties will invest, but past investments will also influence the future contribution levels. Although the model can be applied to other public good settings, it fits especially well to analyze climate policies. As required by the Paris Agreement (Art. 4.9): "Each Party shall communicate a nationally determined contribution every five years." Apparently, "The idea is that this short time frame would give countries the opportunity to regularly capture scientific and technological developments in their official targets."<sup>2</sup> The Stern Review (2006) also pointed out that new technology would be crucial to mitigate climate change. However, the treaties establish that "technology needs must be nationally determined, based on national circumstance and priorities" (§114 of the 2010 Cancun Agreement). For the model to be consistent with this practice, emissions cuts are assumed to be negotiable and contractible, while technology investments are not. (The assumption is relaxed in Section 4.)

In each period  $t$ , the utility for a party is the sum of three parts. First, if each party  $i$  contributes or abates the quantity  $q_{i,t}$ , the sum of abatements has the value  $a \sum_{i \in N} q_{i,t}$  to each party. This linearity assumption is made for simplicity, but it is common also because it is a reasonable approximation when it comes to climate change. As Golosov et al. (2014:78) write, for example: "Linearity is arguably not too extreme a simplification,

---

<sup>2</sup><https://www.carbonbrief.org/explainer-the-ratchet-mechanism-within-the-paris-climate-deal>



since the composition of a concave S-to-temperature mapping with a convex temperature-to-damage function may be close to linear."

An additional benefit of this linearity is that we can easily allow for a stock of greenhouse gases that accumulates over time, without changing the analysis, since  $a$  can be interpreted as the present discounted cost of emitting another unit of emission into the atmosphere, when we anticipate that this unit may contribute to climate change for decades. To see this, suppose party  $i$  emits  $g_{i,t}$  and the pollution stock is  $G_t = \sigma G_{t-1} + \sum_{i \in N} g_{i,t}$ , where  $\sigma \in [0, 1]$  measures the fraction of the past stock that survives to the next period. If parameter  $h > 0$  measures each party's per-period marginal environmental harm from the stock  $G_t$ , then the present discounted harm of another unit of emission is  $h / (1 - \sigma\delta)$  for each party. Consequently,  $a \equiv h / (1 - \sigma\delta)$  measures the present discounted benefit from abating one unit.

The second term in the utility function specifies the cost of contributing to the public good. For example, suppose a country can consume energy from both fossil fuels ( $g_{i,t}$ ) and renewables ( $R_{i,t}$ ). If the total consumption of energy is less than  $i$ 's bliss point,  $B_{i,t}$ , then  $i$  may experience a disutility that is quadratic in the difference:  $\frac{b}{2} (B_{i,t} - [g_{i,t} + R_{i,t}])^2$ . This disutility can be written as  $\frac{b}{2} (q_{i,t} - R_{i,t})^2$ , when  $q_{i,t}$  represents a cut in emissions relative to  $i$ 's bliss point (i.e., when  $q_{i,t} \equiv B_{i,t} - g_{i,t}$ ).<sup>3</sup>

Of course, also for other public good situations, it can be especially costly for  $i$  to contribute a lot relative to  $i$ 's capacity level, as represented by the stock  $R_{i,t}$ .

Each party can over time add to the capacity  $R_{i,t}$  by investing  $r_{i,t}$ . The investment cost is assumed to be convex and quadratic and it constitutes the third term in the per-period utility function:

$$u_{i,t} = a \sum_{j \in N} q_{j,t} - \frac{b}{2} (q_{i,t} - R_{i,t})^2 - \frac{c}{2} r_{i,t}^2, \text{ where}$$

$$R_{i,t+1} = R_{i,t} + r_{i,t}, \tag{1}$$

---

<sup>3</sup>For this particular interpretation of the model,  $g_{i,t} = B_{i,t} - q_{i,t}$  could be negative if  $q_{i,t}$  is very high compared to  $B_{i,t}$ . I do not impose any constraint  $g_{i,t} \geq 0$  because (a) of simplicity, (b)  $g_{i,t} \geq 0$  will not bind if  $B_{i,t}$  is growing sufficiently fast over time, (c)  $g_{i,t} < 0$  is in reality feasible with carbon-capture and storage technologies, and (d) it should be possible to interpret  $q_{i,t}$  as (unbounded) contributions to a public good, more generally. See Harstad (2012) for how one may deal with the constraint  $g_{i,t} \geq 0$  in a similar (although somewhat different) model without affecting the results qualitatively.

and where  $a$ ,  $b$ , and  $c$  are positive constants. The parties can have heterogeneous bliss points for consumption and initial technology levels ( $R_{i,1}$ ) but, for simplicity, the parties are assumed to be identical in other respects.

When  $\delta$  is the common discount factor, party  $i$  intends to maximize the following continuation value at each time  $t$ :

$$V_{i,t} = u_{i,t} + \delta V_{i,t+1} = \sum_{\tau=t}^{\infty} \delta^{\tau-t} u_{i,\tau},$$

*BAU.*—As a benchmark, consider the noncooperative MPE without any treaty, i.e., the "business as usual" (BAU) equilibrium. At every time  $t$ , when  $i$  takes as given  $R_{i,t}$ , the marginal abatement cost equals the marginal benefit for party  $i$ :

$$b(q_{i,t}^{BAU} - R_{i,t}) = a \Leftrightarrow q_{i,t}^{BAU} = R_{i,t} + \frac{a}{b}.$$

Consequently, the investment level does not influence  $i$ 's contribution cost, but only  $i$ 's contribution level in every future period. Party  $i$ 's preferred investment level is thus:

$$r_{i,t}^{BAU} = \frac{\delta}{1-\delta} \frac{a}{c}.$$

With this, it is straightforward to derive party  $i$ 's continuation value in BAU,  $V_{i,t}^{BAU}$ .<sup>4</sup>

The first-best outcome is given by the exact same equations if just  $a$  is replaced by  $na$ . In both cases, the second-order conditions trivially hold.

*Pledges.*—Now, consider the situation that arises after the parties have agreed to contribute *more* than the BAU levels. In particular, suppose  $i$  has agreed to contribute  $x_i \geq 0$  units, beyond  $i$ 's BAU level, for each of the next  $T$  periods. Since there is a large number of subgame-perfect equilibria in dynamic games, it is common to restrict attention to Markov-perfect strategies when there are stocks in the game. This will pin down the unique equilibrium that is the limit of the unique SPE if the time horizon

---

<sup>4</sup>As proven in the Appendix:

$$V_{i,t}^{BAU} = \frac{a}{1-\delta} \sum_{j \in N} R_{j,t} + \frac{a^2}{1-\delta} \left( n - \frac{1}{2} \right) \left( \frac{1}{b} + \frac{1}{c} \left[ \frac{\delta}{1-\delta} \right]^2 \right).$$

were finite but approached infinity. Clearly, the commitment  $x_i$  is payoff-relevant and it might motivate  $i$  to invest  $y_{i,t}$  units in addition to the BAU level. Total contributions and investments can then be written as:

$$q_{i,t} \equiv q_{i,t}^{BAU} + x_i \text{ and } r_{i,t} \equiv r_{i,t}^{BAU} + y_{i,t}. \quad (2)$$

*Timing.*—In each period, the parties simultaneously set the  $q_{i,t}$ 's and the  $r_{i,t}$ 's. (The outcome would be the same if these decisions were made sequentially.) If no agreement has been made regarding the contributions, party  $i$  is free to set any  $q_{i,t}$  and  $r_{i,t}$ . In each of the  $T$  periods after an agreement has been made, the pledge  $x_i$  pins down  $q_{i,t}$  but party  $i$  is free to set  $r_{i,t}$  or, equivalently,  $y_{i,t}$ . It will turn out to be most convenient to focus on the choices of  $x_i$  and  $y_{i,t}$  (of course, these decisions pin down  $q_{i,t}$  and  $r_{i,t}$ , given BAU).

The results below hold whether or not the parties negotiate new pledges after the present  $T$ -period commitment period. To distinguish the two cases, the index  $\iota \in \{0, 1\}$  takes the value of 1 if a new commitment period will be negotiated every  $T$  period, but  $\iota = 0$  if one returns to BAU after the current  $T$ -period commitment period. (Note that we have  $\iota = 1$  for the Paris Agreement since new pledges are supposed to be decided on every five years. This fact will be rationalized in Section 5, where I also relax the assumption that parties can commit.)

## 1.2. *Equilibrium Investments*

Given the above equations, party  $i$ 's continuation value can be written as a function of the  $x_i$ 's and the  $y_{i,t}$ 's. After the pledges have been agreed on, party  $i$ 's problem is to choose the investment levels over the next  $T$  periods. This boils down to a standard optimal control problem that is solved in the Appendix. The exact solution for the investment levels is presented here:

**Lemma 1.** For each  $i \in N$ ,  $t \in \{1, \dots, T\}$ , and  $\iota \in \{0, 1\}$ , equilibrium investments are linear in  $x_i$ :

$$\begin{aligned}
y_{i,t} &= x_i \left( l_1 m_1^{t-1} [1 - m_1] - l_2 m_2^{t-1} [m_2 - 1] \right), \text{ where} \\
m_1 &\equiv \frac{1}{2} \left( \frac{1}{\delta} + 1 + \frac{b}{c} \right) - \frac{1}{2} \sqrt{\left( \frac{1}{\delta} + 1 + \frac{b}{c} \right)^2 - \frac{4}{\delta}} \in (0, 1), \\
m_2 &\equiv \frac{1}{2} \left( \frac{1}{\delta} + 1 + \frac{b}{c} \right) + \frac{1}{2} \sqrt{\left( \frac{1}{\delta} + 1 + \frac{b}{c} \right)^2 - \frac{4}{\delta}} > 1, \\
l_1 &\equiv \frac{m_2^{T-1} (m_2 - 1)}{m_1^{T-1} (1 - m_1) + m_2^{T-1} (m_2 - 1)} \in (0, 1), \text{ and} \\
l_2 &\equiv \frac{m_1^{T-1} (1 - m_1)}{m_1^{T-1} (1 - m_1) + m_2^{T-1} (m_2 - 1)} = 1 - l_1 \in (0, 1).
\end{aligned}$$

Naturally, if  $i$  is committed to contribute a lot, in that  $x_i$  is large, then  $i$  invests more. It is easy to check that  $y_{i,t}$  increases in  $T$ , decreases in  $t$ , and reaches zero when  $t = T$ :

$$y_{i,T} = 0.$$

In the final period, a party invests exactly the same amount as in BAU. (This holds whether one expects to negotiate new pledges in the next period ( $\iota = 1$ ) or not ( $\iota = 0$ )). The intuition for all this is related to the hold-up problem: One more technology unit in the next period can – without any other change in investment or contribution cost – raise the total contribution level by one unit then and forever after. The party that invested captures  $1/n$  of this benefit, just as in BAU (see Harstad, 2016; Battaglini and Harstad, 2016). This intuition also explains why the equilibrium investment at any point in time is the same whether future agreements are expected (i.e.,  $\iota = 1$ ) or not ( $\iota = 0$ ). An important implication is that in every period in which the parties have not yet agreed to any pledge, contribution and investment levels are just as in BAU:  $x_i = 0$  and  $y_{i,t} = 0$ .

Conveniently, Lemma 1 states that the level of investments and thus technologies will be linear functions of  $x_i$ . We can substitute these functions into  $i$ 's utility function and write party  $i$ 's continuation value (i.e., the present discounted value of the future utility levels) as a linear-quadratic function,  $V_{i,1}(\mathbf{x})$ , where  $\mathbf{x} \equiv (x_1, \dots, x_n)$ . Given the benchmark continuation value without an agreement,  $V_{i,1}^{BAU}$ , we are especially interested

in the additional payoff with the pledges:  $U_i(\mathbf{x}) \equiv V_{i,1}(\mathbf{x}) - V_{i,1}^{BAU}$ . The additional payoff  $U_i(\mathbf{x})$  simplifies to a simple linear quadratic function, although with parameters  $\alpha$ ,  $\beta$ , and  $\gamma$  being complicated functions of  $a$ ,  $b$ ,  $c$ ,  $\delta$ , and  $T$ , as proven in Appendix.

**Lemma 2.** *Party  $i$ 's continuation value, relative to BAU, can be written as:*

$$\begin{aligned}
U_i(\mathbf{x}) &= \alpha \sum_{j \neq i} x_j - \frac{\beta}{2} x_i^2 + \gamma, \quad \text{where} & (3) \\
\alpha &\equiv \frac{a}{1 - \delta} [1 - \delta^T (l_1 m_1^{T-1} + l_2 m_2^{T-1})], \\
\beta &\equiv \sum_{t=1}^T \delta^{t-1} \left[ b (l_1 m_1^{t-1} + l_2 m_2^{t-1})^2 + c (l_1 m_1^{t-1} [1 - m_1] - l_2 m_2^{t-1} [m_2 - 1])^2 \right], \text{ and} \\
\gamma &\equiv \delta^T U_i(\mathbf{x}^*).
\end{aligned}$$

### 1.3. Equilibrium Pledges

As soon as the parties have agreed on  $\mathbf{x}$ ,  $i$  faces as continuation value  $U_i(\mathbf{x})$  in addition to the default payoff  $V_{i,t}^{BAU}$ . The bargaining surplus at stake for party  $i$  is thus  $U_i(\mathbf{x})$ .

*The Nash Bargaining Solution.*—The most used bargaining solutions in applied theory is the Nash Bargaining Solution (NBS), predicting that every  $x_i^*$  is the outcome of the following maximization problem:

$$x_i^* = \arg \max_{x_i} \prod_{j \in N} U_j(x_i, \mathbf{x}_{-i}^*) \quad (4)$$

It is easy to check that when all utility levels are the same in equilibrium (i.e., when the situation is symmetric), then the first-order condition (f.o.c.) of (4) is equivalent to the f.o.c. of:<sup>5</sup>

$$x_i^* = \max_{x_i} \sum_{j \in N} U_j(x_i, \mathbf{x}_{-i}^*).$$

---

<sup>5</sup>The two first-order conditions are identical because, when  $U_k(\mathbf{x}^*)$  is the same for every  $k \in N$ , then:

$$\frac{\partial}{\partial x_i} \prod_{j \in N} U_j(x_i, \mathbf{x}_{-i}^*) = \sum_{j \in N} \frac{\partial U_j(x_i, \mathbf{x}_{-i}^*)}{\partial x_i} \prod_{k \in N \setminus j} U_k(x_i, \mathbf{x}_{-i}^*) = \left( \sum_{j \in N} \frac{\partial U_j(x_i, \mathbf{x}_{-i}^*)}{\partial x_i} \right) (U_k(\mathbf{x}^*))^{n-1} = 0.$$

This simplicity is one reason for why the NBS is popular and reasonable in applications. In addition, the NBS relies on an appealing set of axioms (Nash, 1950) and it is the outcome of noncooperative bargaining games, such as the Nash demand game (Nash, 1953). The Rubinstein (1982) alternating-offer bargaining game also implements the NBS, even when there are multiple negotiators, under some consistency conditions (Binmore et al., 1986; Krishna and Serrano, 1996). In the Rubinstein bargaining game, a proposer suggests an outcome  $(x_1, \dots, x_n)$ , implying that when one  $x_i$  is accepted, it is conditional on  $x_j$  being accepted. This conditionality is a well-known characteristic of actual international negotiations, as in those leading up to the Kyoto Protocol (as discussed in the introduction). It is therefore not unreasonable to approximate the outcome of the Kyoto Protocol with the NBS for the participating parties.

*Pledge-and-Review.*—However, the procedure adopted in Paris was quite different. As explained in the introduction, the parties were themselves free to make any nationally determined contribution (NDC) they wanted. This weakened the conditionality aspect that characterized earlier negotiations. Observers have feared that the weight on others' payoffs may be less with P&R than with the NBS, since a party has more discretion when deciding on  $x_i$  under P&R. It is consistent with this logic to let  $x_i^*$  maximize an *asymmetric* Nash product, where the weight on others' payoffs is smaller:

$$x_i^* = \arg \max_{x_i} \prod_{j \in N} U_j(x_i, \mathbf{x}_{-i}^*)^{w_j^i}, \text{ where } w_j^j/w_j^i = w \in [0, 1) \text{ for } j \neq i, \quad (5)$$

When all utility levels are the same in equilibrium (i.e., when the situation is symmetric), then the f.o.c. of (5) is equivalent to the f.o.c. of:

$$x_i^* = \arg \max_{x_i} \left[ U_i(x_i, \mathbf{x}_{-i}^*) + w \sum_{j \in N \setminus i} U_j(x_i, \mathbf{x}_{-i}^*) \right]. \quad (6)$$

*A Micro-foundation for Pledge-and-Review.*—The following game is detailed and analyzed in a working paper (Harstad, 2020a). Suppose all parties announce their pledges simultaneously and independently, but that disapprovals may lead to some delay,  $\Delta$ , before the parties can revise the pledges. If  $\mathbf{x}^*$  is expected in the next round, party  $i$  approves  $\mathbf{x} \neq \mathbf{x}^*$  with probability 1 if  $U_i(\mathbf{x}) \geq \delta^\Delta U_i(\mathbf{x}^*)$ , and with probability 0 otherwise.

This probability is between 0 and 1 if the willingness to accept is uncertain (due to, for example, trembles in the willingness to wait). With such uncertainty, the equilibrium  $\mathbf{x}^*$  is characterized by (5) and the weights are functions of the underlying uncertainty (Theorem 3, Harstad, 2020a). The other parts of the model are left unchanged when the uncertainty is based on temporary shocks. Under reasonable conditions,  $w \in (0, \frac{1}{2})$  (Corollary 1, Harstad, 2020a).

Thus, the viewpoints of commentators as well as bargaining theory support the intuition that the weight  $w$  (on others' payoffs) is larger (say,  $\bar{w}$ ) under the Kyoto Agreement than it is (say,  $\underline{w}$ ) under the pledge-under-review procedure used for the Paris Agreement. The subsequent analysis will investigate and compare the consequences of  $\underline{w} \in (0, 1)$  vs.  $\bar{w} \in (\underline{w}, 1]$ . It is permitted but not necessary to require  $\underline{w} < \frac{1}{2}$  and  $\bar{w} = 1$ .

When we combine eq. (6) with Lemma 2, we obtain:

**Lemma 3.** *For both P&R ( $w = \underline{w}$ ) and the NBS ( $w = \bar{w}$ ), the equilibrium pledge is:*

$$x_i^* = w(n-1)\alpha/\beta. \quad (7)$$

The smaller  $w$  is, the smaller are the  $x_i^*$ 's, and the smaller are investment levels, according to Lemma 1. Both effects make parties worse off than in the situation in which  $w = 1$ . By combining (3) and (7), we see that  $U_i(\mathbf{x}^*)$  increases in  $w$  when  $w \in (0, 1)$ .

**Proposition 1.** *With a smaller  $w$  (as with  $\underline{w} < \bar{w}$ ), contributions are lower, investments are lower, and so is welfare:*

$$U_i(\mathbf{x}^*) = \frac{\alpha^2(n-1)^2}{\beta(1-\delta^T l)} w \left(1 - \frac{w}{2}\right). \quad (8)$$

*Fact 1.*—As explained in the introduction, the Paris Agreement on climate change calls for P&R, while the top-down negotiations associated with the Kyoto Protocol can be approximated by the NBS. Given this factual difference, Proposition 1 is consistent with the criticism mentioned in the introduction as well as with negative experimental evidence on P&R (Barrett and Dannenberg, 2016). The following sections show that the picture will be more nuanced when we endogenize participation.

## 2. PARTICIPATION

This section endogenizes the coalition size and studies how it depends on the bargaining procedure. There are alternative ways of allowing for participation, but the standard approach is simple and it isolates my contribution to the literature. As discussed in the introduction, and according to Nordhaus (2015:1344), "the standard approach in environmental economics" when modelling coalitions begins with a participation stage at which every potential party,  $i \in \{1, \dots, \bar{n}\}$ , decides whether to participate in the coalition. These decisions are made simultaneously and everyone expects that participants will continue by playing the game analyzed in Section 1. Given the restriction to MPEs, free riders will simply follow their dominant BAU strategy and set  $x_i = 0$ .

It is most natural (and common) to focus on pure-strategy equilibria at the participation stage, and doing so pins down the equilibrium coalition size,  $n$ . I start by ignoring the constraint  $n \leq \bar{n}$  as well as a possible minimum participation threshold,  $\underline{n}$ , but these constraints are extensively discussed in Section 3. I also begin assuming that the participation decision is made once and for all, but Section 4 explains why the results continue to hold when this assumption is relaxed.

### 2.1. *Equilibrium Participation*

Since coalition members ends up contributing more than the level that would maximize their own utility, there is a cost of participating in the coalition. For a member to be willing to participate, the benefit of participating must outweigh this cost. The benefit of participating is that other participants will internalize (a fraction  $w$  of) the utility of one additional coalition member.

For each of the  $n$  participants, the equilibrium payoff (in addition to the BAU-payoff) is given by (8). If one of these parties instead free rides, the free rider's additional payoff will be  $\alpha(n-1)w(n-2)\alpha/\beta(1-\delta^T\iota)$ , since each of the other  $n-1$  parties will now commit to contribute  $w(n-2)\alpha/\beta$ , again and again, every  $T$  period (if  $\iota = 1$ ). By



comparison, participation is beneficial if:

$$U_i(\mathbf{x}^*) = \frac{\alpha^2 (n-1)^2}{\beta (1 - \delta^T \iota)} w \left(1 - \frac{w}{2}\right) \geq \frac{\alpha^2 (n-1)(n-2)}{\beta (1 - \delta^T \iota)} w \Rightarrow n \leq 1 + \frac{2}{w}. \quad (9)$$

The size  $n$  cannot be too great since then individual contributions would be so large and so costly that free riding would be preferable. For a coalition to be stable, (9) must hold for the equilibrium  $n$ . At the same time, we must have that  $n' > 1 + 2/w$  for every  $n' > n$ , since, otherwise, nonmembers would also like to participate. To characterize the equilibrium  $n$ , it is useful to employ the function  $\lfloor \cdot \rfloor$ , mapping its argument to the largest weakly smaller integer.

**Proposition 2.** *With a smaller  $w$  (as with  $\underline{w} < \bar{w}$ ), the coalition size is larger:*

$$n = \lfloor 1 + 2/w \rfloor.$$

Note that  $n = 3$  if  $w = 1$ , as when applying the NBS. This "small-coalition paradox" is well known in the literature, which also discusses the trade-off between "narrow-but-deep" vs. "broad-but-shallow" coalitions (see the literature review in the introduction). With P&R,  $w$  is small and a coalition member is not expected to contribute a lot. Lower contributions reduces both the cost and the benefit of participating. The first effect dominates because the cost is a strictly convex function of the contribution, while the benefit function is concave (linear). Thus, when  $w$  is small, participation is attractive for a larger set of  $n$ 's.

Since the number of participants must be an integer,  $n$  is a step function of  $w$ . When comparing bargaining procedures, we are interested in large rather than small differences in  $w$ . Thus, it is not unreasonable to abstract from the fact that  $n$  must be an integer and to use the approximation

$$n \approx n(w) \equiv 1 + 2/w. \quad (10)$$

With this approximation, the product  $(n-1)w$  is a constant that is pinned down when a (marginal) member must be indifferent between free-riding and participating.

## 2.2. Equilibrium Pledges - Revisited

When  $(n - 1)w$  stays constant as  $w$  is reduced,  $x_i^* = (n - 1)w\alpha/\beta$  also remains constant, and so does every investment level  $y_{i,t}$ . Since the individual contributions are invariant in  $w$ , while  $n$  is decreasing in  $w$ , the sum of payoffs will be larger when  $w$  is small. A participant's payoff is also larger when  $w$  is small: this is evident when the endogenous  $n$ , as described by (10), is combined with the utility (8). This gives:

$$U_i^* = \frac{4\alpha^2}{\beta(1 - \delta^T L)} \left( \frac{1}{w} - \frac{1}{2} \right). \quad (11)$$

**Corollary 1.** *When participation is endogenous, and approximated by  $n(w)$ , Proposition 1 is reversed: With a smaller  $w$  (as with  $\underline{w} < \bar{w}$ ), aggregate contributions are larger, aggregate investments are larger, and so is welfare.*

*Fact 2.*—While only 37 countries promised emission cuts for the Kyoto Protocol's first commitment period, 195 countries have pledged to contribute to the Paris Agreement. This fact is consistent with Proposition 2 since the Paris Agreement is associated with P&R and thus a smaller  $w$ . The result that  $x_i$  continues to be large even if  $w$  is small (because  $n$  increases) might shed some light on why the pledges in Table 1 are substantial, despite the P&R procedure.

*Remarks on the optimal  $w$ , Integers, and Robustness.*—Corollary 1 suggests that welfare increases when  $w \rightarrow 0$ . This result holds also when  $n$  is an integer step function, as in Proposition 2, but the result does *not* hold when we take into account that there is a finite number of countries in the world, as in the next section.

The result that  $x_i$  stays unchanged when  $n$  varies endogenously with  $w$  hinges on the functional forms. If, for example, the continuation value had ended up being:

$$\alpha \sum_{j \neq i} x_j - \frac{\beta}{2} x_i^\varphi, \quad \text{with } \varphi > 1,$$

then one can show that:  $n$  always decreases in  $w$ ;  $x_i$  decreases in  $w$  if  $\varphi < 2$  but increases in  $w$  if  $\varphi > 2$ ;  $U_i^*$  always decreases in  $w$  when  $w \in [0, 1]$ . Thus, the rationalization of Fact 2 survives this generalization.

### 3. WHEN TO CHOOSE PLEDGE-AND-REVIEW

With reasonable modifications, the preference concerning bargaining procedure involves a trade-off. This section discusses realistic constraints on  $n$  in order to compare participants' payoffs under  $\underline{w}$  and  $\bar{w} > \underline{w}$ .

#### 3.1. *The Grand Coalition and Maximum Participation*

There is a limit,  $\bar{n}$ , for how large the coalition can be. One may interpret  $\bar{n}$  as the number of countries in the world or, alternatively, as the number of countries that are of significance. If  $\bar{n} < n(\bar{w}) < n(\underline{w})$ , where  $n(\cdot)$  is defined by (10), then both bargaining games (characterized by  $\underline{w}$  or  $\bar{w}$ ) induce full participation. In this case,  $\bar{w}$  is preferable, according to Proposition 1. If, instead,  $n(\bar{w}) < n(\underline{w}) < \bar{n}$ , the upper boundary on  $n$  is nonbinding and  $\underline{w}$  is preferable, according to Corollary 1. A trade-off arises when  $n(\bar{w}) < \bar{n} < n(\underline{w})$ , since then participation is larger, but individual contributions smaller, when  $w$  is small. In this case, a sufficiently large  $\bar{n}$  is necessary to ensure that a participant's payoff is larger under  $\underline{w}$ .

The exact condition follows when comparing a participant's utility, as given by equation (8), for the two cases. The payoff is larger when  $w = \underline{w}$  than when  $w = \bar{w}$  if:

$$\frac{\alpha^2 (\bar{n} - 1)^2}{\beta (1 - \delta^T \iota)} \underline{w}^2 \left( \frac{1}{\underline{w}} - \frac{1}{2} \right) > \frac{\alpha^2 (n(\bar{w}) - 1)^2}{\beta (1 - \delta^T \iota)} \bar{w}^2 \left( \frac{1}{\bar{w}} - \frac{1}{2} \right) \Rightarrow \frac{\bar{n} - 1}{n(\bar{w}) - 1} > \Omega, \text{ where}$$

$$\Omega \equiv \sqrt{\frac{\bar{w} (1 - \bar{w}/2)}{\underline{w} (1 - \underline{w}/2)}} \in \left( 1, \frac{\bar{w}}{\underline{w}} \right).$$

#### 3.2. *Heterogeneity and Minimum Participation*

Countries are more heterogeneous in reality than permitted in the model above. If the willingness to participate varied across countries,  $n$  would not decline as fast as predicted by Proposition 2. A simple way of capturing this heterogeneity is to permit  $\underline{n}$  parties to be committed in that they participate regardless of  $w$ . The reason these parties are committed can be outside the model, but one may think of existing international treaties on non-climate issues such as international trade or regulatory politics. To be specific,

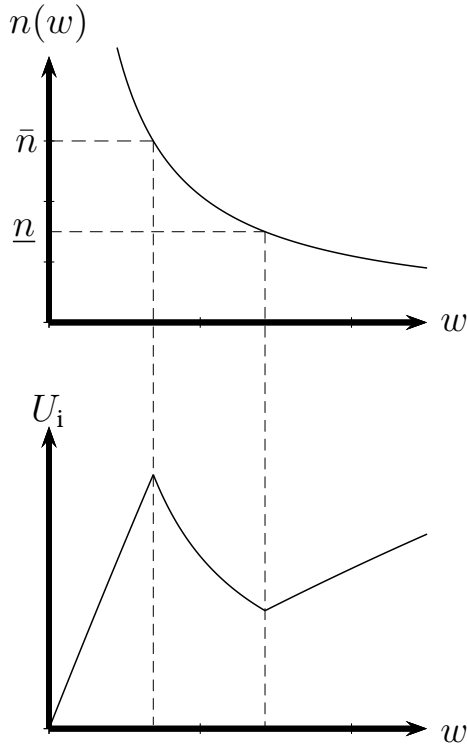


Figure 1: *Participation and participants' payoffs are strictly decreasing in  $w$  only when  $n(w) \in (\underline{n}, \bar{n})$ .*

European Union member countries cannot easily opt out of an environmental agreement unilaterally.

There can be a minimum participation level also for other reasons. Most international treaties specify minimum participation thresholds that must be met for the treaty to enter into force. This threshold was the same for the Kyoto Protocol and the Paris Agreement. In isolation, the effect of such a threshold, if we refer to it as  $\underline{n}$ , is that  $n = \max\{\underline{n}, n(w)\}$ . After all, when the threshold binds, none of the  $\underline{n}$  participants prefer to free-ride given that the consequence will be the BAU outcome.

The minimum participation threshold  $\underline{n}$  is relevant only if  $\underline{n} > n(\bar{w})$ . If also  $\underline{n} > n(\underline{w})$ , the number of participants is always  $\underline{n}$  and then the larger  $w$  is optimal, according to Proposition 1. To isolate the trade-off associated with  $\underline{n}$ , suppose  $n(\bar{w}) < \underline{n} < n(\underline{w}) < \bar{n}$ . In this case, only  $\underline{n}$  parties participate under  $\bar{w}$ , while participation under  $\underline{w}$  is given by  $n(\underline{w})$ . By comparison, a participant's payoff can be larger under  $\underline{w}$  if and only if  $\underline{n}$  is sufficiently small.

The exact condition follows when we use the utility function (8) to compare the two

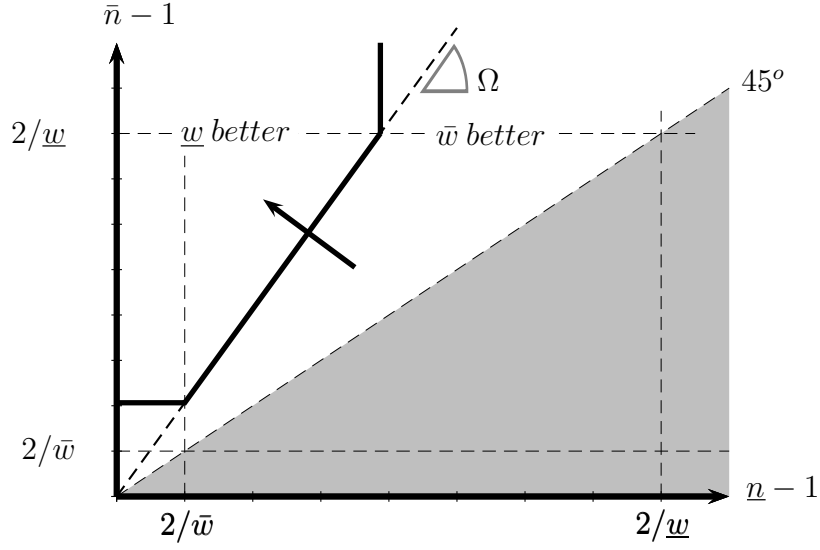


Figure 2: Participants prefer to switch to pledge-and-review ( $\underline{w}$ ) above the solid line.

cases. The payoff is larger when  $w = \underline{w}$  than when  $w = \bar{w}$  if:

$$\frac{\alpha^2 (n(\underline{w}) - 1)^2}{\beta (1 - \delta^T \iota)} \underline{w}^2 \left( \frac{1}{\underline{w}} - \frac{1}{2} \right) > \frac{\alpha^2 (\underline{n} - 1)^2}{\beta (1 - \delta^T \iota)} \bar{w}^2 \left( \frac{1}{\bar{w}} - \frac{1}{2} \right) \Rightarrow \frac{n(\underline{w}) - 1}{\underline{n} - 1} > \Omega.$$

### 3.3. The Preferred Bargaining Game

Figure 1 illustrates that payoffs are non-monotonic in  $w$ . Clearly, it is possible that both the minimum and the maximum participation levels bind at the same time. This happens if  $n(\bar{w}) < \underline{n} < \bar{n} < n(\underline{w})$ . In this case, there is full participation under  $\underline{w}$ , but only  $\underline{n}$  parties participate under  $\bar{w}$ . In this situation,  $\underline{w}$  is preferred when  $\bar{n}$  is large and  $\underline{n}$  is small. From (8), we get that the exact condition is:

$$\frac{\alpha^2 (\bar{n} - 1)^2}{\beta (1 - \delta^T \iota)} \underline{w}^2 \left( \frac{1}{\underline{w}} - \frac{1}{2} \right) > \frac{\alpha^2 (\underline{n} - 1)^2}{\beta (1 - \delta^T \iota)} \bar{w}^2 \left( \frac{1}{\bar{w}} - \frac{1}{2} \right) \Rightarrow \frac{\bar{n} - 1}{\underline{n} - 1} > \Omega.$$

The three conditions can be combined in the following way.

**Proposition 3.** *Everyone who participates when  $w = \bar{w}$  prefers to switch to P&R and  $w = \underline{w}$  if  $\bar{n}$  is large while  $\underline{n}$  is small. The exact condition is:*

$$\frac{\min \{ \bar{n} - 1, 2/\underline{w} \}}{\max \{ \underline{n} - 1, 2/\bar{w} \}} > \Omega. \quad (12)$$

This condition is drawn as the solid line in Figure 2. If there is a larger number of potential parties, or if fewer of them are committed to participate, we move in the direction of the arrow in the figure. Then, the "shallow" agreement ( $\underline{w}$ ) becomes preferred by all participants though the "deep" agreement was preferred given a smaller number of potential parties or a larger number of committed parties.

*Fact 3.—From Kyoto to Paris:* One may argue that both these developments (i.e., a larger  $\bar{n}$  and/or a smaller  $\underline{n}$ ) are in line with changes in world politics over the last couple of decades. Today we have a large number of emerging economies that in the 1990s were developing countries that could not be expected to contribute much to an international climate change treaty. For the model, this development implies that the number of relevant parties,  $\bar{n}$ , has increased. When  $\bar{n}$  is large, it becomes more important to select a procedure that is acceptable even when the set of participants is large.

During the same period, seven of the original Annex I countries, which initially signed the Kyoto Protocol, announced that they would not contribute to the Kyoto Protocol's second commitment period.<sup>6</sup> These withdrawals may be interpreted as a reduction in the number of committed countries,  $\underline{n}$ . Each of these changes (or a combination of the two) can motivate the switch to P&R, according to Proposition 3.

To get a sense for exactly when P&R is preferred, note that the Paris Agreement succeeded in motivating nearly every country in the world to participate. Given that the United States has announced that it will withdraw, it may be reasonable to assume that the weight associated with P&R satisfied:

$$n(\underline{w}) \approx 195 \Leftrightarrow \underline{w} \approx \frac{1}{97}.$$

For the sake of illustration, suppose that the Kyoto Protocol can be approximated by the NBS, with  $\bar{w} = 1$ . When we substitute these numbers into inequality (12), Proposition 3 implies that the Paris Agreement ( $\underline{w}$ ) is preferred to the Kyoto Protocol ( $\bar{w}$ ) by every

---

<sup>6</sup>According to the IPCC (2014:1025), "a number of Annex I countries (Belarus, Canada, Japan, New Zealand, Russia, the United States, and Ukraine) decided not to participate in the second commitment period."

committed country if and only if  $\underline{n}$  is weakly smaller than 28:

$$\text{P\&R} \succ \text{NBS} \Leftrightarrow \underline{n} \leq 28.$$

So, with 37 committed parties, as in the 1990s, the participants would unanimously vote for the Kyoto Protocol. If the number of committed parties is reduced to 28, these participants prefer P&R in order to motivate the larger set of countries to participate. Coincidentally, the European Union, evidently including the most committed set of countries, consists of exactly 28 countries.

Interestingly, there may be a disagreement between the North and the South regarding what procedure to choose. By using the same methodology as above, one can show that the *uncommitted* countries prefer the broad agreement ( $\underline{w}$ ) only when:

$$\min \{\bar{n} - 1, 2/\underline{w}\} > \sqrt{\frac{\max \{\bar{w} (\underline{n} - 1) \underline{n}, 4/\bar{w} + 2\}}{\underline{w} (1 - \underline{w}/2)}},$$

which is stronger than (12), although the comparative statics are similar. With this, the theory can rationalize why developing countries preferred to continue with the Kyoto Protocol, even when the committed countries preferred P&R.<sup>7</sup> In other words, the original set of (committed) participants prefer to switch to P&R too soon, that is, for a larger set of parameters than the set under which such a switch increases global welfare. Analogously, if the new potential members were pivotal in the decision on treaty design, they would accept P&R too late or too seldom, relative to the decision that is optimal when the original members' payoffs are taken into account.

#### 4. THE COMMITMENT PERIOD LENGTH

The results above hold for any commitment period length. The optimal  $T$ , from the participants' point of view, trades off two effects. On the one hand, the shorter the

---

<sup>7</sup>Bodansky et al. (2017:202) write: "Developing countries, for which the Kyoto model has obvious attractions because they are exempt from emissions targets, were keen to extend the protocol for a second and future commitment periods. Kyoto Annex B parties, in contrast, were reluctant to do so, for some countries because of Kyoto's prescriptive architecture, and for others because they did not want to be subject to emissions targets if the US, China, and other large emitters were not."

length of the commitment period is, the lower the equilibrium investments are at every point in time. This comparative static can be seen from Lemma 1 and it was above explained by the classic hold-up problem. On the other hand, with a large  $T$ , the  $x_i$ 's will soon be outdated since it will be optimal to deepen the cuts and the contributions after investments have accumulated.<sup>8</sup> The combined trade-off is new to the literature.<sup>9</sup>

The trade-off when it comes to deciding on  $T$  is independent of  $w$  and  $n$  in the model above. When  $n$  is exogenous, a party's payoff is given by equation (8). When  $n$  is instead endogenous, a party's payoff is given by (11). Every participant's preferred  $T$  is the same in either case:

$$T^* = \arg \max_T U_i(\mathbf{x}^*) = \arg \max_T U_i^* = \arg \max_T \frac{\alpha^2}{\beta(1 - \delta^T \iota)},$$

where  $\alpha$  and  $\beta$  are functions of  $T$ , as described by Lemma 2. This expression is complicated and depends on many of the model's parameters. For example, the optimal  $T$  is larger if  $\iota = 0$  than if  $\iota = 1$  since, in the first case, when no new commitment period will replace the current one, then the duration of cooperation is essentially given by  $T$ . At the same time, the optimal  $T$  does not depend on parameter  $a$ . The intuition is that a larger  $a$  increases *all* benefits and costs of abatements, without altering the trade-off, described above.

More interestingly,  $T^*$  does not depend on  $w$  and  $n$ .

**Proposition 4.** *The optimal commitment period length,  $T^*$ , is invariant in  $w$  and thus the same for  $\underline{w}$  and  $\bar{w}$ , regardless of whether  $n$  and  $w$  are endogenous or exogenous.*

*Fact 4.*—Given the many differences between the Kyoto Protocol and the Paris Agreement, the two are surprisingly similar regarding how frequently commitments must be updated. Pledges under the Paris Agreement must be updated every five years, and the

---

<sup>8</sup>The pledges do permit decreasing pollution levels, since the  $x_i$ 's are defined as cuts relative to the BAU outcome. However, countries invest more, given the pledges, than they do in the BAU outcome, and the importance of these additional investments accumulate over time, implying that the initial pledges can be improved upon. Section 6 explains that the results continue to hold if the pledges are functions of time or of investments.

<sup>9</sup>On the one hand, Harris and Holmstrom (1987) observed that a small  $T$  is beneficial since it permits a rigid contract to be updated when the external environment changes. On the other hand, the hold-up problem associated with a small  $T$  is recognized by, for example, Beccherle and Tirole (2011) and Harstad (2016). These two effects have never been combined, as far as I know.



Kyoto Protocol's first commitment period was also five years (2007-2012). This similarity is consistent with Proposition 4, stating that the optimal commitment period length is the same, despite the many other differences between the two treaties.<sup>10</sup>

Not only this result, but also the mechanism driving it seems to match well with reality. OECD's (2018:5) first argument for a five- rather than a ten-year commitment period is: "More regular opportunities to make technical and fundamental adjustments to NDCs as well as to incorporate effects of technology..."

## 5. COMPLIANCE AND ENFORCEMENT

So far, it has been assumed that the parties are able to commit to the pledges for  $T$  periods. However, given the incentive to free ride, discussed in Sections 2 and 3, it is reasonable to also be concerned with the temptation to contribute less at the time when other participants are expected to deliver on their promises. This section contributes to the literature on self-enforcing agreements (see, for example, Barrett (1994; 2002); Dutta and Radner (2004; 2006); Harstad et al. (2019), and the references therein) by showing when and why certain procedures, such as pledge-and-review, are more likely than others to be self-enforcing. To study self-enforcement, I relax the assumption that the parties *commit* to pledges, and also the restriction to MPEs by permitting history-dependent strategies such as trigger strategies.<sup>11</sup>

Since decisions are made simultaneously, a party that "defects" by not contributing will be able to enjoy the benefit from other participants' contributions in that period. When  $n$  is exogenously given, individual contributions are relatively small when  $w$  is small. The temptation to defect may thus also be small when  $w$  is small. When  $n$  is

---

<sup>10</sup>It is reasonable that Kyoto's second commitment period would also have been five years, if the parties had not anticipated that a new global treaty would be effective from 2020. According to Bodansky et al. (2017:203), in 2011, "Parties disagreed on several issues including: the length of the commitment period—whether it should be five years (like the first commitment period) or eight years (to coincide with the scheduled launch of the 2015 agreement)." In 2012, "the eight-year duration of the second commitment period was chosen so as to end when the Paris Agreement's NDCs were expected to take effect, and thus to avoid a commitment gap" (p. 205).

<sup>11</sup>Although space constraints prevent a full characterization of SPEs, note that if  $\delta \rightarrow 1$  then "folk theorems" imply that a large set of outcomes, including the first best (characterized in Section 1.1), can be supported as SPEs. For a smaller  $\delta$ , the best SPE is likely to be distorted in ways already investigated by Harstad et al. (2019).

endogenous, a smaller  $w$  motivates a larger  $n$  and that, in turn, implies that it is more important for a party that cooperation continues. In both cases, it is intuitive that the incentive constraint is more likely to hold when  $w$  is small.

To illustrate this intuition, suppose the parties revert to BAU (i.e., the noncooperative MPE) forever as soon as one party has defected by contributing less than pledged.<sup>12</sup> Section 6 and the Appendix permit the punishment to last for any length  $L \leq \infty$  of periods and to be triggered by any probability  $\phi \in (0, 1]$ . The result below holds, qualitatively, for every  $L > 0$  and  $\phi > 0$ .

Note that a finitely long agreement cannot be self-enforcing if  $\iota = 0$ , i.e., if one returns to the BAU outcome after period  $T$ . In such a situation, there would be no incentive to comply in period  $T$ , and thus not in period  $T - 1$ , etc. This observation can rationalize why the Paris Agreement specifies that new pledges must be set every five year (i.e.,  $\iota = 1$ ) and this is also why I henceforth restrict attention to  $\iota = 1$ .

The dynamic game in Section 1 is different from a repeated game because past investments influence BAU and thus all future contributions. When party  $i$  invests  $y_{i,t}$ , then  $i$ 's contribution will increase by  $y_{i,t}$  in every future period if the parties revert to BAU. Since every  $y_{i,t}$  is largest at the beginning of the commitment period, the temptation to defect is also largest in the beginning. Then, the payoff if  $i$  defects (by not contributing) is as expressed on the right-hand side in the following inequality. This payoff must be smaller than  $i$ 's equilibrium payoff, on the left-hand side:

$$U_i(\mathbf{x}^*) = \frac{\alpha^2 (n-1)^2}{\beta (1-\delta^T)} w \left(1 - \frac{w}{2}\right) \geq a \left( \sum_{j \neq i} x_j + \frac{\delta}{1-\delta} \sum_{j \neq i} y_{j,1} \right) \Leftrightarrow$$

$$w \leq \hat{w} \equiv 2 - 2 [1 - \delta (l_1 m_1 + l_2 m_2)] \frac{a (1 - \delta^T)}{\alpha (1 - \delta)}. \quad (13)$$

The implication follows when we substitute for the equilibrium  $y_{i,1}$ ,  $x_i$ , and  $\alpha$  and rewrite. The condition is easier to satisfy when  $w$  is small, i.e., if the bargaining procedure is characterized by P&R rather than by the NBS, for example.

Interestingly,  $n$  drops out from the inequality, and thus  $n$  does not influence whether

---

<sup>12</sup>On the one hand, it is possible to sustain as SPEs harsher punishments than the reversion to BAU. With harsher punishments, a treaty would be self-enforcing under a larger set of circumstances than those derived here. On the other hand, if parties could renegotiate punishments, then a treaty would be self-enforcing for a smaller set of parameters.

the bargaining outcome will be self-enforcing. It follows that condition (13) is robust to whether  $n$  is exogenous (as in Section 1) or endogenous (as in Section 2). Technically, the invariance follows because both the cost of the individual contribution and the benefit from the others' contributions are proportional to  $(n - 1)^2$ .

**Proposition 5.** *Regardless of whether participation is exogenous or endogenous, the bargaining outcome is self-enforcing if and only if  $w \leq \hat{w}$ , defined by (13).*

If  $w$  is so large that the incentive constraint is violated, then the parties must find additional ways of raising the cost of noncompliance. In reality, there are several ways of increasing these costs, since the exact wording in an international treaty influences the political and reputational costs if one later defects. Although there exists no world government ready to enforce contracts, it is not irrelevant whether a treaty is called "legally binding." IPCC (2014:1020) explains that "a more legally binding commitment ... signals a greater seriousness by states ... These factors increase the costs of violation (through enforcement and sanctions at international and domestic scales, the loss of mutual cooperation by others, and the loss of reputation and credibility in future negotiations)."

*Fact 5.*—The pledges are not legally binding under the Paris Agreement, but: "the Kyoto Protocol represents a much harder, more prescriptive approach, including legally binding, quantified emissions limitation targets" (Bodansky and Rajamani, 2018:32). This difference between the two agreements is consistent with Proposition 5. Since the Paris Agreement applies P&R bargaining, where  $w$  is smaller, it is possible that the incentive constraint holds for this agreement without making it legally binding. In this case, the parties would strictly prefer non-binding commitments if there were tiny costs associated with legally bindingness (e.g., the observed emission level might be only partly under the government's control, etc.).

Of course, when one can raise the cost of noncompliance by modifying the legal status of the agreement, then countries will comply on the equilibrium path regardless of the bargaining procedure. In line with this prediction, the remaining "thirty-six Kyoto parties [after Canada pulled out] were in full compliance with their first commitment period targets" (Bodansky and Rajamani, 2018:42).

## 6. ROBUSTNESS

The model above is simple and stylized yet able to rationalize Facts 1-5, discussed above. This rationalization turns out to be quite robust in that it continues to hold for a number of model modifications. This relatively technical section explains (and the Appendix proves) that Propositions 1-5 hold even if the parties negotiate investment levels or emission taxes (or both) instead of (or in addition to) the  $x_i$ 's. The  $x_i$ 's can also be time dependent, and the investment levels might be decided by profit-maximizing firms, without changing the propositions. The results are also quite robust to changes in timing.

(i) *Pledging to Invest.*—Some of the NDCs in the Paris Agreement specify national targets for renewable energy.<sup>13</sup> This possibility can be captured by letting parties decide on the  $y_{i,t}$ 's instead of on the  $x_i$ 's. As discussed in the Appendix, it is straightforward to analyze this scenario: when the  $y_{i,t}$ 's, but not the  $x_{i,t}$ 's, are pinned down, then  $i$ 's choice of  $x_{i,t}$  will satisfy  $b(x_{i,t} - Y_{i,t}) = 0$ , just as in BAU, where

$$Y_{i,t+1} \equiv Y_{i,t} + y_{i,t} \text{ and } Y_{i,1} \equiv 0.$$

. If the investment pledge must be time independent ( $y_i$ ) throughout a commitment period, then  $i$ 's continuation value can be written as in Lemma 2, where  $x_i$  is replaced by  $y_i$ , although the definitions of  $\alpha$  and  $\beta$  will be different. The proofs of Propositions 1-5 are thus similar to earlier proofs.

In fact,  $i$ 's continuation value will be separable in the  $\mathbf{y}_t$ 's, where  $\mathbf{y}_t = (y_{1,t}, \dots, y_{n,t})$ . Consequently, we can apply (6) when parties negotiate  $\mathbf{y}_t$ , while keeping fixed the investment levels for other periods. (6) will imply that the P&R outcome for  $y_{i,t}$  is:

$$y_{i,t}^* = (n - 1) w \frac{\delta a/c}{1 - \delta}. \tag{14}$$

Since this  $y_{i,t}$  is time independent, there is no loss for the parties if they restrict attention to time-independent investment levels. For these reasons, the length of the commitment

---

<sup>13</sup>For example, China pledges to increase the share of non-fossil fuels in its primary energy consumption to around 20 percent, while India pledges to produce about 40 percent of its electric power from non-fossil-fuel-based energy resources by 2030. For a recent overview, see <http://cait.wri.org/indc/#/>.

period ( $T$ ) will not influence payoffs, and any  $T$  is here equally good, regardless of the levels of  $n$  and  $w$ .

(ii) *Pledging on Emission Taxes.*—It is also straightforward to allow parties to pledge on domestic emission taxes, instead of on emission cuts. With an emission tax  $z_{i,t}$ , it is natural that consumption of fossil fuel be given by the condition in which the marginal benefit of consuming (or the marginal cost of abating) equals the tax:  $b(x_{i,t} - Y_{i,t}) = z_{i,t}$ . When parties are free to decide their investment levels, they will invest just as in BAU, so  $y_{i,t} = 0$ . If the emission tax level must be time independent ( $z_i$ ) throughout the commitment period, then  $i$ 's continuation value can be written as in Lemma 2, where  $x_i$  is replaced by  $z_i$ , although the definitions of  $\alpha$  and  $\beta$  differ. Again, the proofs of Propositions 1-5 are similar to earlier proofs.

In fact,  $i$ 's continuation value will be separable in the  $z_t$ 's, where  $z_t = (z_{1,t}, \dots, z_{n,t})$ . Consequently, we can apply (6) when parties negotiate  $z_t$ , while keeping fixed the emission taxes for other periods. Equation (6) will imply that the P&R outcome for  $z_{i,t}$  is:

$$z_{i,t}^* = (n - 1) wa. \quad (15)$$

Since this  $z_{i,t}$  is time independent, there is no loss for parties if they restrict attention to time-independent emission taxes. For these reasons, the length of the commitment period ( $T$ ) will not influence payoffs, and any  $T$  is equally good, regardless of the levels of  $n$  and  $w$ .

As a side remark, it is worth noting that the choice of instrument (i.e., whether parties should negotiate  $x_i$ 's,  $y_i$ 's, or  $z_i$ 's) is also independent of  $n$  and  $w$ . As proven in the Appendix, negotiating investment levels is better for all parties than negotiating emission taxes if and only if investments are inexpensive and the future is important.<sup>14</sup>

$$\frac{1}{\delta} < 1 + \sqrt{\frac{b}{c}}.$$

(iii) *Pledging on Investment Levels and Emission Taxes.*—Party  $i$ 's continuation value is separable in the  $y_t$ 's and the  $z_t$ 's, it can be shown. Thus, (6) can be applied for each

---

<sup>14</sup>The comparison to the situation in which the  $x_i$ 's are negotiated is more complex, however.

instrument separately, while keeping the other fixed. With this procedure, the outcome is given by the combination of (14) and (15). In this case, we have a "complete contract" since, given the negotiated investment levels (and thus the  $Y_{i,t}$ 's), the emission taxes pin down the contribution levels.

(iv) *Pledging on Investment Levels and Contribution Levels.*—Once the investment levels (and thus the  $Y_{i,t}$ 's) are pinned down, negotiating  $z_{i,t} = b(x_{i,t} - Y_{i,t})$  is equivalent to negotiating  $x_{i,t}$ . Thus, Scenario (iii) leads to the same outcome as that which occurs when parties can negotiate every investment level and every contribution level. As before, the choice of  $T$  is irrelevant, regardless of the  $n$  and  $w$  levels.<sup>15</sup>

(v) *Time-dependent Contribution Levels.*—In Scenario (iv), one best choice of  $T$  is  $T = \infty$ . With  $T = \infty$ , it is actually irrelevant that parties have negotiated investment levels in addition to contribution levels. The irrelevance follows because, once the  $x_{i,t}$ 's are given for every time, there is no externality associated with the  $y_{i,t}$ 's and, hence, every party will have incentives to invest optimally, without any need to negotiate  $y_{i,t}$ . As is shown in the Appendix, the equilibrium time-dependent contribution level is:

$$x_{i,t}^* = (n-1)w\frac{a}{b} + (n-1)w\frac{a}{c}\frac{\delta}{1-\delta}t.$$

Given this pledge, party  $i$  prefers to invest as in (14), ensuring that the marginal benefit from consuming (and from cutting emissions) is  $b(x_{i,t}^* - ty_{i,t}^*) = (n-1)wa$ , which coincides with  $z_{i,t}^*$  in (15).

In this situation, it is clear that parties are strictly better off with  $T = \infty$  than with  $T < \infty$ , since, with any finite  $T$ , the equilibrium  $y_{i,t}$  is lower (and less efficient) than the  $y_{i,t}$  that would follow in Scenario (iv), which coincides with the equilibrium  $y_{i,t}$  when  $T = \infty$ . When referring to the trade-off discussed in Section 4, there is here no reason to reduce  $T$  in order to update the pledges when the pledges can be time dependent. It is thus optimal with  $T = \infty$  to mitigate the hold-up problem.

---

<sup>15</sup>If parties can negotiate time-independent  $x_j$ 's and  $y_j$ 's, which must stay constant throughout the commitment period, then the parties would strictly prefer  $T = 1$ . With  $T = 1$ , the outcome will be the same as with time-dependent policies (Scenario (iv) and Scenario (iii)), while  $T > 1$  would be less efficient. In contrast to the discussion on the optimal  $T$ , in Section IV, there is no need to have a large  $T$  when the first-period investment level can be negotiated, since agreeing on  $y_{i,1}$  circumvents the hold-up problem.

Of course,  $T < \infty$  continues to be optimal if we introduce e.g. uncertainty (see below).

(vi) *Firms Invest.*—All three Scenarios (iii)-(v) implement the complete contract outcome, i.e., as when all  $y_{i,t}$  and  $x_{i,t}$  are negotiated according to P&R. The same outcome can be achieved if parties negotiate  $x_{i,t}$  at time  $t$ , for  $T = 1$ , while letting firms invest. The equilibrium pledge  $x_{i,t}$  will satisfy  $b(x_{i,t} - Y_{i,t}) = (n - 1)wa$ , which thus also characterizes the marginal willingness to pay for another unit of  $Y_{i,t}$  at time  $t$ . Consequently, the present discounted value of a unit invested today is  $\delta(n - 1)wa / (1 - \delta)$ , while the marginal investment cost is  $cy_{i,t}$ . The two are equalized when profit-maximizing price-taking firms decide on  $y_{i,t}$  and, then, the result is (14), just as when the parties negotiate the investment levels directly. In this situation, it is clear that parties are strictly better off with  $T = 1$  than with  $T > 1$  (unless the contribution levels are time dependent). Firms, unlike governments, are not discouraged by the nations' hold-up problem when new pledges are negotiated.<sup>16</sup>

(vii) *The Timing of T.*—Proposition 4 showed that every participant agreed on the choice of  $T$  and that this choice was independent of  $n$  and  $w$ . Thus, the choice of  $T$  remains the same whether participants decide on  $T$  after the participation stage, before the bargaining-choice stage, or in between the two. The timing of  $T$  influences neither the equilibrium level of  $n$  nor the preference regarding  $w$ .

(viii) *Multiple Participation Stages.*—Propositions 2-5 continue to hold if there is a participation stage before pledges are negotiated at the beginning of every commitment period (i.e., every  $T$  period). Participation is then attractive if and only if (9) holds, just as before. The identity of the  $n$  participants is also the same in every commitment period in an MPE, implying that every participant's continuation value is given by (11). Thus, the proofs of Propositions 2-5 continue to hold.

(ix) *Multiple Bargaining-choice Stages.*—Propositions 2-5 also hold if  $w$ , as well as  $n$ , are endogenously chosen at the beginning of every commitment period, for the same reasons as in Scenario (viii). In fact, if some parameters (such as  $\underline{n}$  and/or  $\bar{n}$ ) change every  $T$  period, then Propositions 2-5 characterize the outcome, and Proposition 3 characterizes the best bargaining procedure, for every commitment period, regardless of the parameter

<sup>16</sup>If each government can subsidize/tax the firms' investments, it can implement its preferred choice of investment, as described in the previous sections. Then, even the exact equations in Sections I-V stay unchanged, one can argue.

values after period  $T$ . This generalization implies that Proposition 3 can indeed rationalize a change from one procedure to another, if  $\underline{n}$  and/or  $\bar{n}$  has changed.<sup>17</sup>

(x) *Limited Punishments*.—When the self-enforcement constraint was discussed, Proposition 5 relied on the assumption that if one party defected, then all parties would play BAU forever after. On the one hand, one may argue that it is more realistic to assume that a defection can be observed with probability  $\phi < 1$ . On the other hand, one may also argue that, if cooperation has broken down, then parties might renegotiate to start cooperating again. To capture these concerns to some extent, the proof of Proposition 5 permits defection to be punished with a reversion to BAU for  $L \leq \infty$  periods with probability  $\phi \leq 1$  (while, with probability  $1 - \phi$ , there is no punishment). The incentive constraint becomes:

$$w \leq 2 - 2 \left[ \frac{1 - \delta (l_1 m_1 + l_2 m_2)}{(1 - \delta) (1 - \delta ([1 - \phi + \phi \delta^L]))} \right] \frac{a (1 - \delta^T)}{\alpha}.$$

A smaller  $\phi$  or  $L$  strengthens the incentive constraint, but note that Proposition 5 holds for all  $\phi \in (0, 1]$  and  $L > 0$ .

These generalizations can be summarized in the following proposition (proven in the Appendix).

**Proposition 6.** *Propositions 1-5 continue to hold if parties pledge-and-review bargain:*

- (i) *investment levels instead of  $\mathbf{x}$ ;*
- (ii) *emission taxes instead of  $\mathbf{x}$ ;*
- (iii) *investment levels and emission taxes instead of  $\mathbf{x}$ ;*
- (iv) *investment levels and  $\mathbf{x}$  instead of only  $\mathbf{x}$ ;*
- (v) *a time profile  $\{\mathbf{x}_t\}_{t=1}^{\infty}$  instead of a time-independent  $\mathbf{x}$ ;*
- (vi)  *$\mathbf{x}$ , while profit-maximizing price-taking firms invest;*
- (vii)  *$T$  after the  $n$  stage, or before  $n$  but after the  $w$  stage.*
- (viii) *Propositions 2-5 continue to hold if there is a participation stage before every commitment period.*

---

<sup>17</sup>The analysis would have been more complicated, however, if parameters changed also within commitment periods.



(ix) Propositions 3-5 continue to hold if both  $w$  and  $n$  are decided on every commitment period.

(x) Proposition 5 holds if defection leads to BAU for  $L \in (0, \infty]$  periods with probability  $\phi \in (0, 1]$ .

The optimal level of  $T$  varies across the scenarios, but for every scenario the optimal  $T$  is independent of the bargaining procedure. Obviously, the optimal  $T$ , as well as the other results, may depend on many things that are outside of this model, such as policy makers' ability to commit or predict the optimal level of contributions in the distant future. Propositions 1-3 and 5 have thus been derived for any fixed  $T$ , and they hold for every  $T$ .

It is possible to extend the model in many other directions as well.<sup>18</sup> This paper has simplified tremendously partly because the additional insight of some of the generalizations would overlap with results in earlier papers and partly because these extensions are evidently not necessary to rationalize the five facts on how the Paris Agreement compares with the Kyoto Protocol. On the contrary, if the modified model did predict that  $T$  should be a function of  $w$  or  $n$ , or that any of the other propositions would change, then it would not be supported by Facts 1-5.

## 7. CONCLUSIONS

The world community must develop better ways of dealing with climate change, but improvements require a deeper understanding of past and present treaties. This paper has shown that (1) the novelty of the pledge-and-review procedure can rationalize four other facts regarding how the Paris Agreement differs from the Kyoto Protocol: (2) Since

---

<sup>18</sup>Consider, for example, the uncertainty in Gerlagh and Liski (2018a). In Harstad (2016), relying on the NBS, both pollution and shocks on the marginal environmental harm accumulate over time. The shocks make it hard to predict optimal pledges and they motivate a small  $T$ , while the hold-up problem motivates a large  $T$ , especially when there are large technological spillovers. In Battaglini and Harstad (2016),  $n$  and subsequently  $T$  are set endogenously before every commitment period. Then, participants may prefer a small  $T$  if  $n$  is small, since a small  $T$  facilitates the admission of new members sooner. Since a small  $T$  also leads to hold-up problems, countries are motivated to participate to ensure a large  $T$ . (Note that their result cannot explain Facts 2 and 4.) Acemoglu et al. (2012) permit investments in dirty as well as in green technology, Dutta and Radner (2020) transfers from the North to the South, Karp (2017) altruism, Rubio (2018) adaptation, and Martimort and Sand-Zantman (2016) a mechanism-design approach.

P&R permits less ambitious pledges, it attracts a larger number of participants. This result can explain why many more countries took on emission cuts in Paris than in Kyoto. (3) Since raising participation is the main benefit of P&R, this procedure is preferable if and only if there is a large number of relevant players. This logic can explain why P&R was preferred in the 2010s, after several developing economies had become emerging economies, whereas the Kyoto Protocol's top-down procedure was chosen in the 1990s. (4) Despite the differences between the two treaties, the theory is consistent with the fact that the commitment period's length is the same for both. (5) Since P&R is not very demanding, and because it attracts a larger number of participants, the equilibrium pledges are more likely to be self-enforcing than is the NBS. This result is consistent with the fact that the Kyoto Protocol's emission cuts were legally binding, whereas they are not for the Paris Agreement.

Alternative theories may successfully explain one or two of the facts, in isolation. The view that the NDCs are simply statements of what the countries do anyway can explain participation but not the move towards the P&R procedure (Fact 3). The claim that only large coalitions sign long-term commitments (Battaglini and Harstad, 2016) rationalizes Fact 2, but not Facts 4 and 5. A coherent framework, consistent with all facts, is necessary for a complete understanding of climate treaties.

Although the paper focuses on a positive analysis, the reader may instinctively search for normative lessons. One lesson is that pledge-and-review might not be as inadequate as it at first appears to be; it can actually be preferable to the alternative when participation is endogenous. However, if participation can be encouraged by other means, then a more demanding conditional-offer bargaining game becomes preferable. Consequently, the benefit of offering "club benefits" (such as the lower tariffs in Nordhaus, 2015) is not, ultimately, that participation will increase, but that parties can choose a more ambitious bargaining procedure without fearing that participation will fall by too much.

## REFERENCES

- Acemoglu, D., P. Aghion, L. Bursztyn, and D. Hemous (2012): "The Environment and Directed Technical Change," *American Economic Review* 102(1):131-66.
- Aldy, J. E., S. Barrett, and R. N. Stavins (2003): "Thirteen plus one: a comparison of global climate policy architectures," *Climate Policy* 3(4): 373-97.
- Barrett, S. (1994): "Self-enforcing international environmental agreements," *Oxford Economic Papers*, 46: 878-94.
- Barrett, S. (2002): "Consensus Treaties," *Journal of Institutional and Theoretical Economics* 158(4): 529-47.
- Barrett, S., and A. Dannenberg (2016): "An experimental investigation into 'pledge and review' in climate negotiations," *Climatic Change* 138(1): 339-51.
- Battaglini, M., and B. Harstad (2016): "Participation and Duration of Environmental Agreements," *Journal of Political Economy* 124(1): 160-204.
- Battaglini, M., and B. Harstad (2020): "The Political Economy of Weak Treaties," *Journal of Political Economy* 128(2): 544-90
- Beccherle, J., and J. Tirole (2011): "Regional Initiatives and the Cost of Delaying Binding Climate Change Agreements," *Journal of Public Economics* 95: 1339-48.
- Bernauer, T., A. Kalbhenn, V. Koubi, and G. Spilker (2013): "Is there a 'Depth versus Participation' Dilemma in International Cooperation?" *Review of International Organization* 8 (4): 477-97.
- Binmore, K., A. Rubinstein, and A. Wolinsky (1986): "The Nash Bargaining Solution in Economic Modelling," *The RAND Journal of Economics* 17(2): 176-88.
- Bloch, F. (2018): "Coalitions and networks in oligopolies," *Handbook of Game Theory and Industrial Organization*, ed. by L. Corchon and M. Marini, Edward Elgar.
- Bodansky, D., J. Brunnee, and L. Rajamani (2017): *International Climate Change Law*, Oxford University Press.
- Bodansky, D., and L. Rajamani (2018): "The Evolution and Governance Architecture of the United Nations Climate Change Regime," in *Global Climate Policy: Actors, Concepts, and Enduring Challenges*, ed. by U. Luterbacher and D. Sprinz, MIT Press.
- Calvo, E., and S. J. Rubio (2012): "Dynamic Models of International Environmental Agreements: A Differential Game Approach," *International Review of Environmental and Resource Economics* 6: 289-339.
- Caparrós, A. (2016): "Bargaining and International Environmental Agreements," *Environmental and Resource Economics* 65(1): 5-31.
- Caparrós, A. (2020): "Pledge and implement bargaining in the Paris Agreement on climate change," mimeo.
- Carraro, C., and D. Siniscalco (1993): "Strategies for the International Protection of the Environment." *Journal of Public Economics* 52(3): 309-28.
- d'Aspremont, C., A. Jacquemin, J. J. Gabszewicz, and J. A. Weymark (1983): "On the stability of collusive price leadership," *The Canadian Journal of Economics* 16(1): 17-25.
- Dutta, P. K., and R. Radner (2004): "Self-enforcing climate-change treaties," *Proceedings of the National Academy of Science* 101: 4746-51.
- Dutta, P. K., and R. Radner (2006): "A Game-Theoretic Approach to Global Warming," *Advances in Mathematical Economics* 8: 135-53.

- Dutta, P. K., and R. Radner (2020): "The Paris Accord Can Be Effective if the Green Climate Fund is Effective," mimeo, Columbia University.
- Finus, M., and S. Maus (2008): "Modesty may pay!" *Journal of Public Economic Theory* 10: 801–26.
- Gerlagh, R., and M. Liski (2018a): "Carbon Prices for The Next Hundred Years," *The Economic Journal* 128(609): 728–57.
- Gerlagh, R., and M. Liski (2018b): "Consistent climate policies," *Journal of the European Economic Association* 16(1): 1–44.
- Gilligan, M. J. (2004): "Is There a Broader-Deeper Trade-off in International Multilateral Agreements?" *International Organization* 58 (3): 459-84.
- Gollier, C., and J. Tirole (2015): "Making Climate Agreements Work," *The Economist*, guest blog, June 1st: <https://www.economist.com/free-exchange/2015/06/01/making-climate-agreements-work>
- Golosov, M., J. Hassler, P. Krusell, and A. Tsyvinski (2014): "Optimal Taxes on Fossil Fuel in General Equilibrium," *Econometrica* 82(1): 41-88.
- Harris, M., and B. Holmstrom (1987): "On The Duration of Agreements," *International Economic Review* 28(2): 389-406.
- Harstad, B. (2012): "Climate Contracts: A Game of Emissions, Investments, Negotiations, and Renegotiations," *Review of Economic Studies* 79(4): 1527-57.
- Harstad, B. (2016): "The Dynamics of Climate Agreements," *Journal of the European Economic Association* 14(3): 719-52.
- Harstad, B. (2020a): "A Theory of Pledge-and-Review Bargaining," mimeo.
- Harstad, B. (2020b): "Technology and Time Inconsistency," forthcoming, *Journal of Political Economy*.
- Harstad, B., F. Lancia, and A. Russo (2019): "Compliance Technology and Self-Enforcing Agreements," *Journal of the European Economic Association* 17(1): 1-30.
- Hoel, M. (1992): "International environmental conventions: the case of uniform reductions of emissions," *Environmental and Resource Economics* 2(2): 141-59.
- IPCC (2014): *Climate Change 2014: Mitigation of Climate Change. Contribution of Working Group III to the Fifth Assessment Report of the Intergovernmental Panel on Climate Change*. Cambridge University Press.
- Karp, L. (2017): "Provision of a Public Good with Multiple Dynasties," *The Economic Journal* 127(607): 2641–64.
- Karp, L., and H. Sakamoto (2019): "Sober optimism and the formation of international environmental agreements," mimeo, UC Berkeley.
- Keohane, R. O., and M. Oppenheimer (2016): "Paris: Beyond the Climate Dead End through Pledge and Review," *Politics and Governance* 4(3): 42-51.
- Kováč, E., and R. C. Schmidt (2019): "A simple dynamic climate cooperation model," mimeo, University of Hagen.
- Krishna, V., and R. Serrano (1996): "Multilateral Bargaining", *Review of Economic Studies* 63(1): 61-80.
- Kolstad, C. D., and M. Toman (2005): "The Economics of Climate Policy," *Handbook of Environmental Economics* 3: 1562-93.
- Leonard, D., and N. Van Long (1992): *Optimal Control Theory and Static Optimization in Economics*, Cambridge University Press.

- Martimort, D., and W. Sand-Zantman (2016): "A Mechanism Design Approach to Climate-change Agreements," *Journal of the European Economic Association* 14(3): 669-718.
- Nash, J. (1950): "The Bargaining Problem," *Econometrica* 18: 155-62.
- Nash, J. (1953): "Two-Person Cooperative Games," *Econometrica* 21(1): 128-40.
- Nordhaus, W. D. (2015): "Climate Clubs: Overcoming Free-riding in International Climate Policy," *American Economic Review* 105(4): 1339-70.
- OECD (2018): "Common time frames: Summary of discussions at the March 2018 Climate Change Expert Group Global Forum," Note prepared by the OECD/IEA Climate Change Expert Group.
- Rubinstein, A. (1982): "Perfect Equilibrium in a Bargaining Model," *Econometrica* 50(1): 97-109.
- Rubio, S. J. (2018): "Self-Enforcing International Environmental Agreements: Adaptation and Complementarity," FEEM Working Paper 29.2018.
- Schmalensee, R. (1998): "Greenhouse policy architectures and institutions," *Economics and Policy Issues in Climate Change*, ed. by W. D. Nordhaus, Resources for the Future Press, Washington, D. C.
- Stern, N. (2006): *The Economics of Climate Change: The Stern Review*, Cambridge University Press
- Sydsaeter, K., and P. J. Hammond (1995): *Mathematics for Economic Analysis*, Prentice Hall.
- Tirole, J. (2017): *Economics for the Common Good*, Princeton University Press.
- Victor, D. (2015): "Why Paris Worked: A Different Approach to Climate Diplomacy," *Yale Envir.* 360:  
[https://e360.yale.edu/features/why\\_paris\\_worked\\_a\\_different\\_approach\\_to\\_climate\\_diplomacy](https://e360.yale.edu/features/why_paris_worked_a_different_approach_to_climate_diplomacy).

## APPENDIX: PROOFS

I will start by reformulating the optimal control problem described in Section 1.

**Lemma A-1.** *Given the actual pledges,  $\mathbf{x}$ , and the future equilibrium pledges,  $\mathbf{x}^*$ , party  $i$ 's continuation value is  $V_{i,1}(\mathbf{x}) = V_{i,1}^{BAU} + U_i(\mathbf{x})$ , where:*

$$U_i(\mathbf{x}) \equiv \max_{\{y_{i,t}\}_{t=1}^T} \sum_{t=1}^T \delta^{t-1} \left[ a \sum_{j \neq i} x_j - \frac{b}{2} (x_i - Y_{i,t})^2 - \frac{c}{2} y_{i,t}^2 \right] + a \frac{\delta^T}{1-\delta} \sum_{j \neq i} Y_{j,T+1} + \delta^T U_i(\mathbf{x}^*),$$

$$Y_{i,t+1} \equiv Y_{i,t} + y_{i,t}, \text{ and } Y_{i,1} \equiv 0.$$

The lemma permits the current pledges ( $\mathbf{x}$ ) to be different from those expected in equilibrium in the subsequent commitment period (i.e.,  $\mathbf{x}^*$ ). Conveniently, the heterogeneous bliss points and initial technology levels drop out when utility is measured relative to BAU. It is also convenient that the investments' effects on  $Y_{j,T+1}$  are captured in terms that do not interact with the future continuation value,  $\delta^T U_i(\mathbf{x}^*)$ . The additional investments affect the future  $V_{i,1}^{BAU}$  but not  $U_i(\mathbf{x})$ .

*Proof.* I will first derive  $V_{i,t}^{BAU}$ . When we substitute in for  $u_{i,t}$ ,  $q_{i,t}^{BAU}$ , and  $r_{i,t}^{BAU}$  into  $U_{i,t}^{BAU} = \sum_{\tau=t}^{\infty} \delta^{\tau-t} u_{i,\tau}$ , we can rewrite  $V_{i,t}^{BAU}$  as:

$$\begin{aligned} V_{i,t}^{BAU} &= \sum_{\tau=t}^{\infty} \delta^{\tau-t} \left[ a \sum_{j \in N} \left( R_{j,\tau} + \frac{a}{b} \right) - \frac{b}{2} \left( \frac{a}{b} \right)^2 - \frac{c}{2} \left( \frac{\delta}{1-\delta} \frac{a}{c} \right)^2 \right] \\ &= \frac{a}{1-\delta} \sum_{j \in N} R_{j,t} + a \sum_{j \in N} \sum_{\tau=t}^{\infty} \frac{\delta^{\tau+1-t}}{1-\delta} y_{j,\tau} + \sum_{\tau=t}^{\infty} \delta^{\tau-t} \left[ \left( n - \frac{1}{2} \right) \frac{a^2}{b} - \frac{c}{2} \left( \frac{\delta}{1-\delta} \frac{a}{c} \right)^2 \right] \\ &= \frac{a}{1-\delta} \sum_{j \in N} R_{j,t} + \frac{a^2}{1-\delta} \left( n - \frac{1}{2} \right) \left( \frac{1}{b} + \frac{1}{c} \left[ \frac{\delta}{1-\delta} \right]^2 \right). \end{aligned}$$

Similarly, the BAU payoff at time  $T+1$  can be written as:

$$V_{i,T+1}^{BAU} = \frac{a}{1-\delta} \sum_{j \in N} (R_{j,T+1}^{BAU} + Y_{j,T+1}) + \frac{a^2}{1-\delta} \left( n - \frac{1}{2} \right) \left( \frac{1}{b} + \frac{1}{c} \left[ \frac{\delta}{1-\delta} \right]^2 \right),$$

where  $Y_{i,T+1}$  measures the additional investments, relative to BAU, thanks to the first commitment period. Each party's present-discounted value of  $Y_{i,T+1}$  is  $a \frac{\delta^T}{1-\delta} \sum_j Y_{i,T+1}$ , when evaluated in period 1. This term should be added when we derive the additional utility, relative to BAU, when the  $n$  parties commit to  $\mathbf{x}$  for  $T$  periods at time  $t=1$  (even

if the parties thereafter returned to BAU). The additional utility, relative to BAU, is thus:

$$\begin{aligned}
& \sum_{t=1}^T \delta^{t-1} \left[ a \sum_{j \in N} (q_{j,t}^{BAU} + x_j) - \frac{b}{2} (q_{i,t}^{BAU} + x_i - R_{i,t}^{BAU} - Y_{i,t})^2 - \frac{c}{2} (r_{i,t}^B + y_{i,t})^2 - u_{i,t}^{BAU} \right] \\
& + a \frac{\delta^T}{1-\delta} \sum_{j \in N} Y_{j,T+1} \\
& = \sum_{t=1}^T \delta^{t-1} \left[ a \sum_{j \neq i} x_j - \frac{b}{2} (x_i - Y_{i,t})^2 - \frac{c}{2} y_{i,t}^2 + a Y_{i,t} - a \delta \frac{Y_{i,t+1} - Y_{i,t}}{1-\delta} \right] + a \delta^T \frac{\sum_{j \in N} Y_{j,T+1}}{1-\delta} \\
& = \sum_{t=1}^T \delta^{t-1} \left[ a \sum_{j \neq i} x_j - \frac{b}{2} (x_i - Y_{i,t})^2 - \frac{c}{2} y_{i,t}^2 \right] + a \frac{\delta^T}{1-\delta} \sum_{j \neq i} Y_{j,T+1},
\end{aligned} \tag{16}$$

where the last equality follows because the three terms with  $Y_{i,\tau}$  in (16) sum to zero for each  $\tau = \{2, \dots, T+1\}$  and because  $Y_{i,1} = 0$ .

When the parties to *not* play BAU after the first commitment period, then, in order to obtain  $i$ 's total additional payoff relative to BAU, we must add the additional payoff  $\delta^T U_i(\mathbf{x}^*)$ , where  $U_i(\mathbf{x}^*)$  is the equilibrium additional utility relative to BAU, in order to get  $U_i(\mathbf{x})$  in Lemma A-1.  $\parallel$

#### *Proof of Lemma 1*

Lemma A-1 defines an optimal-control problem with control  $y_{i,t}$ . Note that the terminal value for  $Y_{i,T+1}$  is zero because  $U_i(\mathbf{x})$  is measured relative to  $V_{i,1}^{BAU}$ : this implies that  $y_{i,T} = 0$ , i.e., the investment level in the final period coincides with the equilibrium investment level in BAU. In other words, there is no *additional* investment in the final period.

When  $\lambda_t$  defines the shadow value of the stock  $Y_{i,t}$ , evaluated at time 1, the discrete-time Hamiltonian can be written as:<sup>19</sup>

$$H_t = \delta^{t-1} \left[ a \sum_{j \neq i} x_j - \frac{b}{2} (x_i - Y_{i,t})^2 - \frac{c}{2} y_{i,t}^2 \right] + \lambda_{i,t+1} (Y_{i,t} + y_{i,t}),$$

with first-order conditions

$$y_{i,t} = \arg \max_{y_{i,t}} H_t = \lambda_{i,t+1} / c \delta^{t-1},$$

adjoint equation

$$\lambda_{i,t+1} - \lambda_{i,t} = -\frac{\partial H_t}{\partial Y_{i,t}} = -\delta^{t-1} b (x_i - Y_{i,t}),$$

and terminal condition

$$\lambda_{i,T+1} = 0 \Leftrightarrow y_{i,T} = 0.$$

<sup>19</sup>I here apply Pontryagin's maximum principle for discrete time problems. For a general characterization and proof, see, for example, Leonard and van Long (1992:129-33).

Combining the first two conditions and (1), we get the second-order difference equation:

$$\begin{aligned} c\delta^{t-2}(Y_{i,t} - Y_{i,t-1}) - c\delta^{t-1}(Y_{i,t+1} - Y_{i,t}) &= \delta^{t-1}(x_i - Y_{i,t})b \Rightarrow \\ -Y_{i,t+1} + \left(\frac{1}{\delta} + 1 + \frac{b}{c}\right)Y_{i,t} - \frac{1}{\delta}Y_{i,t-1} &= x_i b/c, \end{aligned}$$

which has the solution (see, e.g., Sydsaeter and Hammond, 1995:751-53):

$$\begin{aligned} Y_{i,t} &= A_1 m_1^{t-1} + A_2 m_2^{t-1} + x_i, \text{ where} \tag{17} \\ m_1 &= \frac{1}{2} \left( \frac{1}{\delta} + 1 + \frac{b}{c} \right) - \frac{1}{2} \sqrt{\left( \frac{1}{\delta} + 1 + \frac{b}{c} \right)^2 - \frac{4}{\delta}} \in (0, 1), \\ m_2 &= \frac{1}{2} \left( \frac{1}{\delta} + 1 + \frac{b}{c} \right) + \frac{1}{2} \sqrt{\left( \frac{1}{\delta} + 1 + \frac{b}{c} \right)^2 - \frac{4}{\delta}} > 1. \end{aligned}$$

The constants  $A_1$  and  $A_2$  can be derived from the initial condition  $Y_{i,1} = 0$ , implying  $A_1 + A_2 = -x_i$ , and the terminal condition,  $y_{i,T} = 0$ , implying

$$\begin{aligned} y_{i,T} = Y_{i,T+1} - Y_{i,T} &= A_1 m_1^T \left( 1 - \frac{1}{m_1} \right) - (A_1 + x_i) m_2^T \left( 1 - \frac{1}{m_2} \right) = 0 \Rightarrow \\ A_1 &= -\frac{m_2^T \left( 1 - \frac{1}{m_2} \right)}{m_1^T \left( \frac{1}{m_1} - 1 \right) + m_2^T \left( 1 - \frac{1}{m_2} \right)} x_i, \text{ and} \\ A_2 &= -A_1 - x_i = \frac{m_2^T \left( 1 - \frac{1}{m_2} \right)}{m_1^T \left( \frac{1}{m_1} - 1 \right) + m_2^T \left( 1 - \frac{1}{m_2} \right)} x_i - x_i = -\frac{m_1^T \left( \frac{1}{m_1} - 1 \right)}{m_1^T \left( \frac{1}{m_1} - 1 \right) + m_2^T \left( 1 - \frac{1}{m_2} \right)} x_i. \end{aligned}$$

With the definitions  $l_1 = -A_1 x_i$  and  $l_2 = -A_2 x_i$ , (17) can be written as in Lemma 1. ||

*Proof of Lemma 2*

By substituting in for  $y_{i,t}$  and  $Y_{i,t}$  into  $U_{i,1}(\mathbf{x})$ , defined in Lemma A-1, we get:

$$\begin{aligned} U_i(\mathbf{x}) - \delta^T U_i(\mathbf{x}^*) &= \sum_{t=1}^T \delta^{t-1} \left[ a \sum_{j \neq i} x_j - \frac{b}{2} (x_i - Y_{i,t})^2 - \frac{c}{2} y_{i,t}^2 \right] + a \frac{\delta^T}{1 - \delta} \sum_{j \neq i} Y_{j,T+1} \\ &= \sum_{t=1}^T \delta^{t-1} \left[ \begin{aligned} &a \sum_{j \neq i} x_j - \frac{b}{2} x_i^2 (l_1 m_1^{t-1} + l_2 m_2^{t-1})^2 \\ &-\frac{c}{2} [x_i (l_1 m_1^{t-1} [1 - m_1] - l_2 m_2^{t-1} [m_2 - 1])]^2 \end{aligned} \right] \\ &\quad + a \frac{\delta^T}{1 - \delta} \sum_{j \neq i} Y_{j,T+1} \end{aligned}$$



$$\begin{aligned}
&= \alpha \sum_{j \neq i} x_j + \beta x_i^2/2, \text{ if just} \\
\alpha &\equiv \sum_{t=1}^T \delta^{t-1} a + a \frac{\delta^T}{1-\delta} \frac{Y_{j,T+1}}{x_j} = \frac{a}{1-\delta} [1 - \delta^T (l_1 m_1^{T-1} + l_2 m_2^{T-1})] \text{ and} \\
\beta &\equiv \sum_{t=1}^T \delta^{t-1} \left[ b (l_1 m_1^{t-1} + l_2 m_2^{t-1})^2 + c [(l_1 m_1^{t-1} [1 - m_1] - l_2 m_2^{t-1} [m_2 - 1])]^2 \right]. \quad \parallel
\end{aligned}$$

*Proofs of Propositions 1-4*

The proof of Proposition 1 follows from the earlier Lemmata, while Propositions 2 and 3 follow from the reasoning in the text. Proposition 4 follows straightforwardly from the equilibrium continuation values, derived above.

*Proof of Proposition 5*

If  $i$  defects by not contributing at time  $t$ , then  $i$  can still benefit  $a \sum_{j \neq i} x_j + \frac{a\delta}{1-\delta} \sum_{j \neq i} y_{j,t}$ , since  $j$ 's investments will raise  $j$ 's contribution in the future, even when the parties return to BAU. This benefit is largest at  $t = 1$ , since  $y_{j,t}$  is decreasing in  $t \in \{1, \dots, T\}$ , as noticed already.

When defection is punished by a reversion to BAU for  $L \leq \infty$  periods with probability  $\phi \in (0, 1]$ , then compliance (giving payoff  $U_i^*$ ) is better at time  $t = 1$  if:

$$a \sum_{j \neq i} x_j + \frac{a\delta}{1-\delta} \sum_{j \neq i} y_{j,1} + \delta (1 - \phi + \phi\delta^L) U_i^* \leq U_i^*.$$

When we substitute in for  $y_{j,1}$ ,  $x_j^*$ , and  $U_i^*$ , this inequality becomes:

$$\begin{aligned}
&a \left( \sum_{j \neq i} x_j^* + \frac{\delta}{1-\delta} \sum_{j \neq i} y_{j,1}^* \right) \leq [1 - \delta (1 - \phi + \phi\delta^L)] U_i^* \iff \\
&a \left[ 1 + \frac{\delta}{1-\delta} (1 - l_1 m_1 - l_2 m_2) \right] \sum_{j \neq i} x_j^* \leq [1 - \delta (1 - \phi + \phi\delta^L)] \frac{\alpha^2 (n-1)^2}{\beta (1-\delta^T)} w \left( 1 - \frac{w}{2} \right) \iff \\
&\frac{a (1 - \delta^T)}{\alpha [1 - \delta (1 - \phi + \phi\delta^L)]} \left[ \frac{1 - \delta (l_1 m_1 + l_2 m_2)}{1 - \delta} \right] \leq 1 - \frac{w}{2} \iff \\
&w \leq 2 - 2 \frac{1 - \delta (l_1 m_1 + l_2 m_2)}{(1 - \delta) [1 - \delta (1 - \phi + \phi\delta^L)]} \frac{a (1 - \delta^T)}{\alpha},
\end{aligned}$$

which equals (13) when  $\phi = 1$  and  $L \rightarrow \infty$ .  $\parallel$

*Proof of Proposition 6*

(i) *Contracts on investments:* I will first permit the negotiated  $\mathbf{y}_t = (y_{1,t}, \dots, y_{n,t})$  to be time-dependent, so that  $\mathbf{y} = (\mathbf{y}_1, \dots, \mathbf{y}_T)$  is a matrix. Lemma 1 presents a reformulation of the problem and (when we remove the max-operator) it holds regardless of how the  $x_{i,t}$ 's and the  $y_{i,t}$ 's are decided on. When  $y_{i,t}$  is committed to, but not  $x_{i,t}$ , the latter follows straightforwardly from  $i$ 's maximization problem and, just as in BAU,

$$q_{i,t} - R_{i,t} = a/b \Rightarrow x_i = Y_{i,t}.$$

The continuation value can thus be written as a function of the investments matrix  $\mathbf{y}$ :

$$U_i(\mathbf{y}) = \sum_{t=1}^T \delta^{t-1} \left[ a \sum_{j \neq i} \sum_{t'=1}^{t-1} y_{j,t'} - \frac{c}{2} y_{i,t}^2 \right] + a \frac{\delta^T}{1-\delta} \sum_{j \neq i} \sum_{t'=1}^T y_{j,t'} + \delta^T U_i(\mathbf{y}^*) \iff$$

$$U_i(\mathbf{y}) - \delta^T U_i(\mathbf{y}^*) = \sum_{t=1}^T \left[ \alpha_t \sum_{j \neq i} y_{j,t} - \frac{\beta_t}{2} y_{i,t}^2 \right], \text{ where } \alpha_t = \frac{a\delta^t}{1-\delta} \text{ and } \beta_t = \delta^{t-1}c.$$

If we require a time-independent  $y_{j,t} = y_j$ , we can write

$$U_i(\mathbf{y}) - \delta^T U_i(\mathbf{y}^*) = \alpha \sum_{j \neq i} y_j - \frac{\beta}{2} y_i^2, \text{ where } \alpha = \delta a \frac{1-\delta^T}{(1-\delta)^2} \text{ and } \beta = \sum_{t=1}^T \delta^{t-1}c = c \frac{1-\delta^T}{1-\delta}.$$

Just as before, we can write  $i$ 's payoff as in Example E. Consequently, the proofs for the other propositions follow the same steps as above. Proposition 1 gives, for example:

$$y_j^* = w(n-1)\alpha/\beta = w(n-1)\frac{\delta a/c}{1-\delta}.$$

*Time-dependent investment levels:* Since  $i$ 's payoff is separable in the  $y_{j,t}$ 's, we can apply (6) for each  $\mathbf{y}_t$ , if we fix the investment levels for the other periods, in order to get:

$$y_{j,t}^* = w(n-1)\alpha_t/\beta_t = w(n-1)\frac{\delta a/c}{1-\delta},$$

which equals  $y_j^*$ . Hence, the restriction to time-independent investment levels is nonbinding: the equilibrium is the same in both cases.

The choice of  $T$  is also irrelevant in both cases, since the equilibrium continuation value is:

$$U_i(\mathbf{y}^*) = \delta a \frac{1}{(1-\delta)^2} (n-1)^2 w \frac{\delta a/c}{1-\delta} - \frac{c/2}{1-\delta} \left[ (n-1) w \frac{\delta a/c}{1-\delta} \right]^2 = \frac{[\delta a(n-1)]^2}{c(1-\delta)^3} w \left( 1 - \frac{w}{2} \right).$$

(ii) *Contracts on carbon tax:* I will first permit  $z_t = (z_{1,t}, \dots, z_{n,t})$  to be time-dependent, so that  $z = (\mathbf{z}_1, \dots, \mathbf{z}_T)$  is a matrix.

With an emission tax equal to  $z_{i,t}$ , collected by the government in country  $i$ , the equilibrium ensures that the marginal benefit when consuming fossil fuel (or the marginal abatement cost) equals  $z_{i,t}$ . This implies:

$$x_{i,t} - Y_{i,t} = z_{i,t}/b,$$

and, therefore,  $i$ 's continuation value can be written as the function

$$U_i(\mathbf{z}) = \sum_{t=1}^T \delta^{t-1} \left[ a \sum_{j \neq i} (z_{j,t}/b + Y_{j,t}) - \frac{z_{i,t}^2}{2b} - \frac{c}{2} y_{i,t}^2 \right] + a \frac{\delta^T}{1-\delta} \sum_{j \neq i} Y_{j,T+1} + \delta^T U_i(\mathbf{z}^*)$$

so, there is no value for  $i$  to invest beyond the BAU-levels, and  $y_{i,t} = 0$ , so:

$$U_i(\mathbf{z}) - \delta^T U_i(\mathbf{z}^*) \equiv \sum_{t=1}^T \delta^{t-1} \left[ a \sum_{j \neq i} z_{j,t}/b - \frac{z_{i,t}^2}{2b} \right] = \sum_{t=1}^T \left[ \alpha_t \sum_{j \neq i} z_{j,t} - \frac{\beta_t}{2} z_{i,t}^2 \right], \text{ where}$$

$$\alpha_t = a\delta^{t-1}/b \text{ and } \beta_t = \delta^{t-1}/b.$$

If the emission tax is time-independent, we can write:

$$U_i(\mathbf{z}) - \delta^T U_i(\mathbf{z}^*) = \alpha \sum_{j \neq i} z_j - \frac{\beta}{2} z_i^2, \text{ where } \alpha = \frac{a}{b} \frac{1 - \delta^T}{1 - \delta} \text{ and } \beta = \frac{1}{b} \frac{1 - \delta^T}{1 - \delta}.$$

In this case, Lemma 2 implies:

$$z_i^* = w(n-1)\alpha/\beta = w(n-1)a.$$

*Time-dependent tax:* Since  $i$ 's payoff is separable in the  $z_{j,t}$ 's, we can apply Lemma 2 for each  $z_t$ , if we fix the emission tax levels for the other periods, in order to get:

$$z_{j,t}^* = w(n-1)\alpha_t/\beta_t = w(n-1)a,$$

which equals  $z_i^*$ . Hence, the restriction to time-independent emission tax levels is nonbinding: the equilibrium is the same in both cases.

The choice of  $T$  is also irrelevant in both cases, since the equilibrium continuation value is:

$$U_i(\mathbf{z}^*) = \frac{a}{b} \frac{1}{1 - \delta} (n-1)^2 wa - \frac{1}{2} \frac{1}{b} \frac{1}{1 - \delta} [(n-1)wa]^2 = \frac{[a(n-1)]^2}{b(1-\delta)} w \left(1 - \frac{w}{2}\right).$$

*By comparison,* a tax gives higher payoff than an investment agreement if:

$$\frac{[a(n-1)]^2}{b(1-\delta)} > \frac{[\delta a(n-1)]^2}{c(1-\delta)^3} \iff c(1-\delta)^2 > b\delta^2 \iff \frac{1}{\delta} > 1 + \sqrt{\frac{b}{c}}.$$

Clearly, the investment agreement is better if investments are inexpensive and the tax ineffective (because  $b$  is large). If  $\delta$  is large, investments are, in effect, less expensive, and thus the investment agreement is more attractive.

*(iii) Combining (i) and (ii):* When the parties face both a matrix of emission taxes and a matrix of investment levels,  $i$ 's payoff can be written as:

$$U_i(\mathbf{x}) - \delta^T U_i(\mathbf{x}^*) \equiv \sum_{t=1}^T \delta^{t-1} \left[ a \sum_{j \neq i} \left( \frac{z_{j,t}}{b} + Y_{j,t} \right) - \frac{z_{i,t}^2}{2b} - \frac{c}{2} y_{i,t}^2 \right] + a \frac{\delta^T}{1 - \delta} \sum_{j \neq i} Y_{j,T+1}$$

$$= \left[ \sum_{t=1}^T \delta^{t-1} \left( a \sum_{j \neq i} Y_{j,t} - \frac{c}{2} y_{i,t}^2 \right) + a \frac{\delta^T}{1 - \delta} \sum_{j \neq i} Y_{j,T+1} \right] + \left[ \sum_{t=1}^T \delta^{t-1} \left( a \sum_{j \neq i} \frac{z_{j,t}}{b} - \frac{z_{i,t}^2}{2b} \right) \right],$$

where the first (second) bracket can be recognized as  $i$ 's payoff in the situation when only the investment levels (the emission taxes) were negotiated. The two problems are thus

separable, and the results above continue to hold when the parties can negotiate both policy instruments. In this case, the additional payoff, relative to BAU, is also the sum of the two additional payoffs, derived above:

$$U_i(\mathbf{y}^*) + U_i(\mathbf{z}^*) = \left[ \frac{1}{c(1/\delta - 1)^2} + \frac{1}{b} \right] \frac{[a(n-1)]^2}{(1-\delta)} w \left( 1 - \frac{w}{2} \right).$$

(iv) *Complete contracts*: When the parties negotiate the investment levels, the  $z_{j,t}$ 's pin down the  $x_{j,t}$ 's, given the  $y_{j,t}$ 's, so negotiating the  $z_{j,t}$ 's is then equivalent to negotiating the  $x_{j,t}$ 's: Also when the  $y_{j,t}$ 's and the  $x_{j,t}$ 's are negotiated, the contract is complete and the choice of  $T$  is irrelevant. One optimal  $T$  is thus  $T = \infty$ .

(v) *Time-path for  $x$* : When the  $y_{j,t}$ 's and the  $x_{j,t}$ 's are negotiated, one optimal  $T$  is  $T = \infty$ . In this situation, pinning down the  $x_{j,t}$ 's is equivalent to pinning down both the  $y_{j,t}$ 's and the  $x_{j,t}$ 's, because there is no externality when it comes to the  $y_{j,t}$ 's (given every future  $x_{j,t}$ ) and, hence, every party will invest optimally, without any need to specify the investment levels.

This reasoning completes the proof but, to illustrate, consider the time profile for the contribution levels when the parties negotiate both the emission taxes and the investment levels:

$$x_{i,t} = (n-1)wa/b + t(n-1)w \frac{\delta a/c}{1-\delta}.$$

Given this path, optimal investments, from  $i$ 's point of view, are:

$$\begin{aligned} cy_{i,t-1} - \delta cy_{i,t} &= \delta b(x_{i,t} - Y_{i,t}) = \delta b \left( (n-1)wa/b + t(n-1)w \frac{\delta a/c}{1-\delta} - t(n-1)w \frac{\delta a/c}{1-\delta} \right) \\ &= \delta b((n-1)wa/b) \Rightarrow y_{i,t-1} = \frac{\delta(n-1)wa}{c(1-\delta)}, \end{aligned}$$

just as in the optimal contract. So, the combination of negotiating investment levels and emission taxes is indeed equivalent to pinning down the path of  $x_{i,t}$ .

(vi) *Firms*: It suffices to prove that when  $T = 1$  and the parties negotiate  $x_{i,t}$  at the start of every period  $t$ , and the firms invest to maximize profit, then the outcome coincides with the outcome when all the  $y_{i,t}$ 's and the  $x_{i,t}$ 's are negotiated at the very beginning.

When only this period's  $x_{i,t}$  are negotiated at the start of period  $t$ , then, when applying Lemma 2:

$$b(x_{i,t} - Y_{i,t}) = aw(n-1).$$

Firms invest such as to equalize the marginal investment cost to the present-discounted value of their investment, where the willingness to pay for more  $R_{i,t}$  equals  $b(q_{i,t} - R_{i,t})$  at time  $t$ . Thus:

$$\begin{aligned} cr_{i,t} &= \sum_{t=1}^{\infty} \delta^t b(q_{i,t} - R_{i,t}) = \sum_{t=1}^{\infty} \delta^t b(q_{i,t}^{BAU} - R_{i,t}^{BAU} + x_{i,t} - Y_{i,t}) \\ &= \sum_{t=1}^{\infty} \delta^t b \left( \frac{a}{b} + x_{i,t} - Y_{i,t} \right) = \sum_{t=1}^{\infty} \delta^t b \left( \frac{a}{b} + \frac{a}{b} w(n-1) \right) = \frac{\delta}{1-\delta} b \left( \frac{a}{b} + \frac{a}{b} w(n-1) \right). \end{aligned}$$

With  $r_{i,t} = r_{i,t}^{BAU} + y_{i,t}$  and  $r_{i,t}^{BAU} = \frac{\delta}{1-\delta} \frac{a}{c}$ , we get  $cy_{i,t} = \frac{\delta}{1-\delta} aw(n-1)$ , as with complete

contracts.

(vii)-(ix) are trivial and thus omitted.

(x) *Compliance*: In all the above situations, and also in the basic model if  $c \rightarrow \infty$ , we have that  $U_i^*$  is independent of  $T$  and it can, when the policy instrument is given by the matrix  $\psi = (\psi_1, \dots, \psi_K)$ , where  $\psi_k = (\psi_{1,k}, \dots, \psi_{n,k})$  for each  $k \in \{1, \dots, K\}$ , be written as the following (for some constants  $\alpha_k$  and  $\beta_k$ ):

$$U_i^* = \sum_{k \in \{1, \dots, K\}} \frac{1}{1 - \delta} \left[ \alpha'_k \sum_{j \neq i} \psi_{j,k} - \frac{\beta'_k}{2} \psi_{j,k}^2 \right], \text{ so } \psi_{j,k} = w(n-1) \alpha'_k / \beta'_k,$$

from Lemma 2. If defection is punished by reverting to BAU for  $L$  periods with probability  $\phi$ , then the incentive constraint is:

$$\begin{aligned} \sum_{k \in \{1, \dots, K\}} \alpha'_k \sum_{j \neq i} \psi_{j,k} + \delta(1 - \phi + \phi\delta^L) U_i^* \leq U_i^* &\iff \sum_{k \in \{1, \dots, K\}} \frac{(n-1)^2 (\alpha'_k)^2}{\beta'_k} w \leq \\ \sum_{k \in \{1, \dots, K\}} \frac{1 - \delta(1 - \phi + \phi\delta^L)}{1 - \delta} \left[ \frac{(n-1)^2 (\alpha'_k)^2}{\beta'_k} w - \frac{\beta'_k}{2} \left[ \frac{(n-1) \alpha'_k}{\beta'_k} w \right]^2 \right] &\iff \\ 1 \leq \frac{1 - \delta(1 - \phi + \phi\delta^L)}{1 - \delta} \left[ 1 - \frac{1}{2} w \right] &\iff \\ w \leq 2 - 2 \frac{1 - \delta}{1 - \delta(1 - \phi + \phi\delta^L)} = 2 \frac{1 - \delta(1 - \phi + \phi\delta^L) - 1 + \delta}{1 - \delta(1 - \phi + \phi\delta^L)} = 2\delta \frac{\phi(1 - \delta^L)}{1 - \delta(1 - \phi + \phi\delta^L)}, \end{aligned}$$

which simplifies to  $w \leq 2\delta$  if  $\phi = 1$  and  $L = \infty$ . ||