

# PARTICIPATION GAMES AND INTERNATIONAL ENVIRONMENTAL AGREEMENTS: A NON-PARAMETRIC MODEL

LARRY KARP<sup>†</sup>      AND      LEO SIMON<sup>††</sup>

FEBRUARY 3, 2011

**ABSTRACT.** We examine the size of stable coalitions in a participation game that has been used to model international environmental agreements, cartel formation, R&D spillovers, and monetary policy. The literature to date has relied on parametric examples; based on these examples, a consensus has emerged that in this kind of game, the equilibrium coalition size is small, except possibly when the potential benefits of cooperation are also small. In this paper, we develop a non-parametric approach to the problem, and demonstrate that the conventional wisdom is not robust. In a general setting, we identify conditions under which the equilibrium coalition size can be large even when potential gains are large. Contrary to previously examined leading special cases, we show that reductions in abatement costs in an international environmental game can increase equilibrium membership, and we provide a measure of the smallest reduction in costs needed to support a coalition of arbitrary size.

**KEYWORDS:** Stable coalitions, participation game, International Environmental Agreement, climate agreement, trans-boundary pollution, investment spillovers.

**JEL CLASSIFICATION NUMBERS:** C72, H4, Q54

---

<sup>†</sup>Department of Agricultural and Resource Economics, University of California, Berkeley, and the Ragnar Frisch Center for Economic Research, email: karp@berkeley.edu

<sup>††</sup>Department of Agricultural and Resource Economics, University of California, Berkeley, and the Giannini Foundation of Agricultural Economics, email: leosimon@berkeley.edu .

## 1. INTRODUCTION

There is a rapidly expanding literature on the prospects for international environmental agreements (IEAs) to combat trans-boundary pollution problems, including the depletion of the ozone layer, the proliferation of greenhouse gases and, more generally, global warming. Much of this literature is based on a particular participation game. Due to the difficulty of analyzing this game, virtually all research on the topic specifies particular functional forms for the costs and benefits of abatement. Based on this research, the unchallenged consensus opinion among researchers in this field is that IEAs are ineffective precisely when the potential gains from cooperation are large (Ioannidis et al. (2000), Finus (2001), Finus (2003), Barrett (2003)). Moreover, in leading special cases, an innovation—for example, a reduction in abatement costs—that increases the potential gains from cooperation will leave unchanged or reduce equilibrium participation. In these cases, a cost reduction that one would expect to be welfare enhancing has no effect on, or even reduces the fraction of the welfare gains from cooperation that could potentially be realized through the formation of an IEA. The existing theory on IEAs is thus pessimistic: sovereign nations cannot easily be induced to provide public goods; voluntary, self-enforcing agreements are not effective in providing public goods, and plausible remedies, such as investments that reduce abatement costs, will not help and might even be counterproductive. This pessimism spills into policy advice: it has been suggested that efforts devoted to promoting a climate agreement requiring signatories to reduce greenhouse gas emissions are bound to be unsuccessful. A corollary, discussed by Barrett (2003, Ch. 15) and Stiglitz (2006), is that an IEA cannot be successful in the absence of some kind of external punishment.

We show that the basis for this pessimism is fragile, in that it relies on specific functional forms. All of the forms considered in the IEA literature share a common property: the marginal abatement cost function is convex. This property is not required by classical micro-theory, which assumes only that the cost function itself is convex. When it is relaxed, the pessimistic results in the literature can be easily overturned. In this paper, we avoid the straightjacket of specific functional forms, focusing instead on a non-parametric expression for the function—we call it the *joiner's gain function*—that determines the equilibrium level of participation. When the marginal benefits from abatement are linear, there is a simple decomposition of this function that leads to a quite general theory of IEAs, and provides intuition that helps in analyzing the case where marginal benefits of abatement are concave. We begin by showing that if marginal abatement costs are weakly convex and benefits are linear, then the equilibrium size of the IEA never exceeds three countries. To this extent, the basis for the literature's pessimism is more general than previously thought. However, when the restriction that marginal abatement costs be convex is relaxed (maintaining the classical assumption that total costs are convex) the equilibrium IEA can be arbitrarily large. Thus, for some classically admissible cost functions—although not the ones previously analyzed in the literature—the pessimistic conclusion is overturned.

When marginal benefits are linear, our decomposition of the joiner's gain function is a powerful analytic tool. We use it to determine, for any integer up to the total number of potential IEA signatories, the type and minimal size of the cost reduction needed to support an equilibrium with at least this many members, in circumstances where under the original

cost function an IEA of this size would not be stable. It is thus possible to capture up to 100% of the potential gains from cooperation that are available under the appropriately specified, reduced cost function. We then extend in part these results to the general case of decreasing marginal benefits: we show that an agreement of any size can be supported, but do not attempt to measure the the minimum required cost reduction. Our results have important policy implications. The existing literature is pessimistic about the possibility that lower abatement costs can ameliorate the collective action problem. We show how to design cost reductions that do exactly that, and, when marginal benefits are linear, we provide a measure of the magnitude of the reduction required for them to be effective.

The literature on IEAs to which we contribute is based on a two-stage game developed in D'Aspremont et al. (1983).<sup>1</sup> In the first stage agents have a binary decision, whether to join or stay out of the coalition; in the second stage, agents take an action (abatement, sales, etc.). Members of the coalition internalize the effect of their actions on other coalition members. The outcome at the participation stage is a noncooperative Nash equilibrium, and the outcome at the second ("abatement") stage is either a Nash or Stackelberg equilibrium.

While our research was motivated by the difficulty of sustaining IEA's, this two-stage game has been used to study a wide variety of problems. Since the essential features of the model arise in several fields, the relevance of our results extends well beyond the setting of this paper. In an IO setting, Donsimoni et al. (1986) provide conditions under which the size of the stable coalition is uniquely determined. Katz (1986) and De Bondt (1997) use the D'Aspremont et al. (1983) definition of coalition stability to study collusion in R&D games with spillovers. Thoron (1998) shows that stable coalitions and Nash equilibrium outcomes to the binary action participation game are equivalent. Kohler (2002) uses the definition of stability to study cooperation in monetary policy. Escriva-Villar (2009) and Bos & Harrington Jr (2010) use the same definition of stability in a setting where the post-participation outcome is a non-cooperative equilibrium to a repeated game, rather than the result of optimization by the coalition. Bloch & Dutta (2008) explain the relation between the particular definition of coalition stability used in these papers, and other prominent definitions.

Early applications of this model to IEA formation include Hoel (1992), Carraro & Siniscalco (1993), and Barrett (1994). Barrett (2002) and Finus & Maus (2008) show that coalitions with modest ambitions may be more successful than those that try to fully internalize damages. Using functional forms that are standard for this literature, Barrett (2006) presents a model in which cost-reducing investments can improve the outcome under a stable coalition if there are increasing returns to scale (IRTS) to the adoption decision. With IRTS, the participation game becomes a coordination game, and there may exist a non-cooperative equilibrium in which all nations participate. Following up on this result, Hoel & De Zeeuw (2010) show that cost-reducing investments (without IRTS) can improve the performance of an IEA if investment can reduce abatement costs sufficiently that the a dominant strategy for countries is to abate at the socially optimal level. From our perspective, it is the pessimistic converse of these results that is significant: under the usual functional forms, investments that reduce abatement costs do little or nothing to improve the performance of IEAs, except in the special cases of IRTS and the "dominant strategy scenario;" we show that this converse

is not robust, once the possibility of non-convex marginal costs is introduced. Dixit & Olson (2000) and Hong & Karp (2010) study the game when agents use mixed strategies at the participation stage. Ulph (2004), Kolstad (2007), and Kolstad & Ulph (2008) examine the effect on participation and welfare when agents anticipate learning about costs or damages after the participation stage, before they choose abatement. Kosfeld, Okada & Riedl (2009), Burger & Kolstad (2009) and Dannenberg, Lange & Sturm (2009) test in laboratory settings the predictions of the simplest participation game and of some of its extensions.

Section 2 reviews the standard model, discusses the main conclusions of the literature, and illustrates that these conclusions are not robust. In Section 3, we consider arbitrary convex abatement costs, but require that marginal benefits from abatement are linear. In Section 4, we assume that the abatement benefit function is strictly concave. Section 5 concludes.

## 2. PRELIMINARIES

In this section, we review the standard model and discuss two specifications that have been frequently used. In the first, the marginal cost of abatement is constant up to a capacity constraint; in the second, marginal costs are linear. These special cases serve as a tutorial for readers unfamiliar with the topic, and illustrate the broad conclusions of the literature reviewed above. We show by example that these conclusions are not robust.

In the model developed below, we use “s” (for signatories) subscripts to denote signatories to the agreement and “f” (for freeloaders) subscripts to denote non-signatories. Superscripts indicate the number of signatories in the agreement. Signatories delegate their abatement decisions to the coalition, which chooses abatement to maximize coalition welfare. Each non-signatory ignores the public good nature of abatement and maximizes its individual welfare. Each of the  $N$  countries is identical, so the model predicts equilibrium participation and abatement, but not the identity of participants. There are two variants of the second stage, a Cournot variant in which signatories and non-signatories choose abatement levels simultaneously, and a Stackelberg variant in which the signatories choose first. For most of this paper, we will focus on Cournot. We consider Stackelberg in subsection 4.2.

Abatement is a public good. Each country derives a benefit  $B(Q)$  from the global level of abatement,  $Q$ , and incurs a cost based on its individual abatement level,  $q$ . The individual abatement cost function is convex, denoted by  $C(q)$ . Since  $B'(0) = 1$ , non-signatories choose an abatement level of zero whenever  $C'(0) > 1$ ; we examine a simple example of this case in subsection 2.1. In the remainder of this paper, we assume that  $C'(0) \leq 1$ .

Although the number of signatories,  $n$ , has a natural interpretation only when  $n$  is an integer, all of the variables in the paper are mathematically well-defined for any nonnegative real value of  $n$ . Accordingly, we define all functions below on the non-negative reals. For  $r \in \mathbb{R}_+$ , we denote by  $\lceil r \rceil$  the smallest integer weakly greater than  $r$  and by  $\lfloor r \rfloor$  the greatest integer weakly less than  $r$ .

Suppose that in the first stage of the participation game,  $r \in \mathbb{R}_+$  countries choose to participate in an IEA. In the Cournot variant of the abatement stage, the abatement levels  $q_s^r$

and  $q_f^r$  for signatories and non-signatories respectively are simultaneously determined as the solutions to the first order conditions (1a) and (1b) below:

$$0 = r \frac{\partial B(Q^r)}{\partial q_s^r} - C'(q_s^r) \quad (1a)$$

$$0 = \frac{\partial B(Q^r)}{\partial q_f^r} - C'(q_f^r), \quad (1b)$$

where  $Q^r = r q_s^r + (N - r) q_f^r$  denotes aggregate abatement when there are  $r$  members.

For  $r \in \mathbb{R}_+$ , we let  $g(r)$  denote the “joiner’s gain” function, representing the increment to a non-signatory’s payoff if it becomes the  $r$ ’th member of an IEA:

$$g(r) = \begin{cases} [B(Q^r) - C(q_s^r)] - [B(Q^{r-1}) - C(q_f^{r-1})] & \text{if } r \leq N \\ -1 & \text{if } r > N \end{cases}. \quad (2)$$

For  $r \leq N$ , the first term in square brackets is the net (private) benefit obtained by a signatory to an  $r$ -member IEA; the second term is the net benefit to not joining: if it did not join, the non-signatory would benefit from aggregate abatement  $Q^{r-1}$  and incur a cost  $C(q_f^{r-1})$ . Therefore, the difference between the two terms in square brackets is indeed the net gain to a potential signatory of becoming the  $r$ ’th member of an IEA.

We assume that a country indifferent between joining an IEA or not chooses to join. Breaking ties in this way, an IEA with an integer  $n \geq 2$  members<sup>2</sup> is a *stable equilibrium* iff

$$g(n) \geq 0 > g(n+1). \quad (3)$$

The first inequality is known as the “internal stability condition”—a member of an IEA with  $n$  members has no incentive to leave—while the second is called the “external stability condition”—a non-member strictly prefers not to join an IEA that has  $n$  members. Note that the external stability condition is vacuously satisfied if  $n = N$  since in this case there are no non-members; for this reason we have set  $g(n) = -1$  for  $n > N$ .

To solve the model, the natural approach is to compute the real-valued roots of the equation  $g(\cdot) = 0$ , and then check each root to see if the integers on either side of it satisfy (3). For a *non-integer* value<sup>3</sup> of  $r \in \mathbb{R}_+$  s.t.  $g(r) = 0$ , we say that  $r$  is a *stable root of  $g$*  if

$$g(\lfloor r \rfloor) \geq 0 > g(\lceil r \rceil). \quad (3')$$

Having made this point, we shall for the remainder of the paper focus on integer-valued solutions to the model, and index solutions with  $n$  rather than  $r$ .

In sections 2-3, we assume that benefits are linear in abatement, and choose units so that each country’s benefit  $B(Q)$  is equal to  $Q$ . Linearity simplifies the analysis considerably, since each non-signatory has a dominant abatement strategy: in particular, its abatement choice is independent of both the size of, and the collective abatement decision made by the

IEA; hence under this assumption the Cournot and Stackelberg variants are equivalent. In section 4, we generalize the model to allow for benefits that are concave in  $Q$ .

When there are  $n$  signatories and  $B(Q) = Q$ , each signatory in an interior equilibrium abates at the level  $q_s^n > 0$ , where  $q_s^n$  solves  $n = C'(q)$ , while each non-signatory abates at  $q_f > 0$ , where  $q_f$  solves  $1 = C'(q)$ , i.e.,  $q_f$  is independent of  $n$ . In this case, we can rewrite the joiner's gain function (2) as:

$$\begin{aligned} g(n) &= nq_s^n + (N - n)q_f - C(q_s^n) - [(n - 1)q_s^{n-1} + (N - n + 1)q_f - C(q_f)] \\ &= n(q_s^n - q_s^{n-1}) + (q_s^{n-1} - q_f) - [C(q_s^n) - C(q_s^{n-1}) + C(q_s^{n-1}) - C(q_f)]. \end{aligned} \quad (4)$$

For heuristic reasons, it is helpful in the linear case to write

$$g(n) = CR(n) - UA(n), \quad \text{where} \quad (5a)$$

$$UA(n) = C(q_s^{n-1}) - C(q_f) - (q_s^{n-1} - q_f) > 0; \quad (5b)$$

$$CR(n) = n(q_s^n - q_s^{n-1}) - [C(q_s^n) - C(q_s^{n-1})] > 0. \quad (5c)$$

$UA$  (unilateral action) is the utility loss to a non-signatory if it increases its abatement from  $q_f$  to the level  $q_s^{n-1}$  produced by the signatories to an  $(n - 1)$ -member IEA, while other countries maintain the same abatement levels; because the non-signatory's abatement level  $q_f$  is individually optimal, the cost increment of this unilateral step is greater than the benefit increment.  $CR$  (collective response) is the utility gain to the signatory when all signatories to the augmented coalition respond to the additional member by increasing abatement from  $q_s^{n-1}$  to  $q_s^n$ ; because  $q_s^n$  is collectively optimal for the augmented coalition, this second increment in aggregate abatement is greater than the second increment in cost.

**2.1. Constant marginal costs with a capacity constraint.** To motivate the research orientation of this paper, we begin with a simple example that has been widely discussed in the literature cited above: benefits are linear in aggregate abatement, while marginal costs are constant up to a capacity constraint, normalized to 1. A country's cost of abatement level  $q$  is

$$C(q) = \begin{cases} cq & \text{for } 0 \leq q \leq 1 \\ \infty & \text{for } q > 1 \end{cases}. \quad (6)$$

We choose units so that each country's benefit is  $Q$ , where  $Q = \sum q$  is aggregate abatement. As the marginal benefit of abatement is constant at 1, this model is interesting only when  $c > 1$ , so that abatement at capacity is not a dominant strategy in the abatement stage.

In this model, signatories to an IEA of size  $n$  will abate to the capacity level of 1 iff  $n \geq \lceil c \rceil$ ; if  $n < \lceil c \rceil$ , signatories abate zero. Non-signatories always abate zero. To obtain the joiner's

gain function for this model, we substitute these properties into (4):

$$g(n) = \begin{cases} 0 & \text{for } n < \lceil c \rceil \\ n - c & \text{for } n = \lceil c \rceil . \\ 1 - c & \text{for } n > \lceil c \rceil \end{cases} \quad (7)$$

Note in particular that if  $n > \lceil c \rceil$ , so that signatories to an IEA with  $n-1$  members would abate at capacity, the  $n$ 'th signatory would incur a cost of  $c > 1$  to obtain only one more unit of abatement. Given our tie-breaking assumption (see p. 4 and eq. (3)), it follows from (7) that the model under this specification has a unique stable equilibrium which results in a positive level of abatement: in this equilibrium,  $n = \lceil c \rceil$ .

A striking property of this specification is that equilibrium coalition size weakly increases with abatement costs. Moreover, equilibrium global welfare equals<sup>4</sup>  $(N - c) \lceil c \rceil$ , which is a “saw-toothed” function of  $c$ : when  $c$  reaches any integer  $n$ , participation increases discontinuously while costs increase continuously, leading to a jump in equilibrium welfare; as  $c$  increases through the interval  $(n, n+1)$ , participation remains unchanged but costs increase, so that global welfare declines.

One might expect technological innovations that reduce the cost of abatement to increase global welfare, but the model with cost specification (6) delivers the opposite result: it implies that IEAs with large numbers of signatories are sustainable only if abatement is very costly, and that efforts to increase abatement efficiency may be counter-productive. This conclusion is reinforced by other papers in the literature utilizing different cost specifications.<sup>5</sup> The present paper is written with this intellectual background in mind: in particular, we focus our attention on how cost-reducing technological enhancements *can* in fact increase both equilibrium participation and aggregate welfare.

**2.2. Linear marginal costs.** We begin this subsection with a review of the game with linear marginal costs. We maintain the assumption that the benefit function is linear. Using a simple diagrammatic argument, we confirm the result in Barrett (1994) that under this specification, the unique stable equilibrium is  $n = 3$ , regardless of the slope of the marginal cost function. (Footnote 6 provides a formal proof.) A change in the slope of the marginal cost function changes the equilibrium level of abatement, leaving unchanged both the equilibrium number of signatories and the fraction of potential global welfare gains realized in the game. We then show that by introducing a kink in the marginal cost function, and flattening it north-east of the kink, we can support any integer  $3 < n \leq N$  as a second stable equilibrium ( $n = 3$  remains as an equilibrium). Thus we show by example that technological innovations that reduce costs *are* compatible with increased, indeed even 100%, participation. Fig. 1 is a diagrammatic representation of our decomposition of the joiner's gain function  $g(n)$  into two components,  $CR(n)$  and  $UA(n)$  (See (4)-(5b)), when the cost function is  $C(q) = \frac{q^2}{2M}$ , for some  $M \in \mathbb{R}$ , so that  $C'(q) = q/M$ . Recall that  $q_s^n$  solves  $n = C'(q)$ , while  $q_f$  solves  $1 = C'(q)$ . From (5c),  $UA(n)$  is the area under the marginal cost curve between  $q_f$  and  $q_s^{n-1}$ , minus the area of the rectangle with boundaries  $q_f$  and  $q_s^{n-1}$  and height 1. The difference is the area of the cross-hatched triangle labeled  $UA(n)$ . From (5b),  $CR(n)$  is the area of

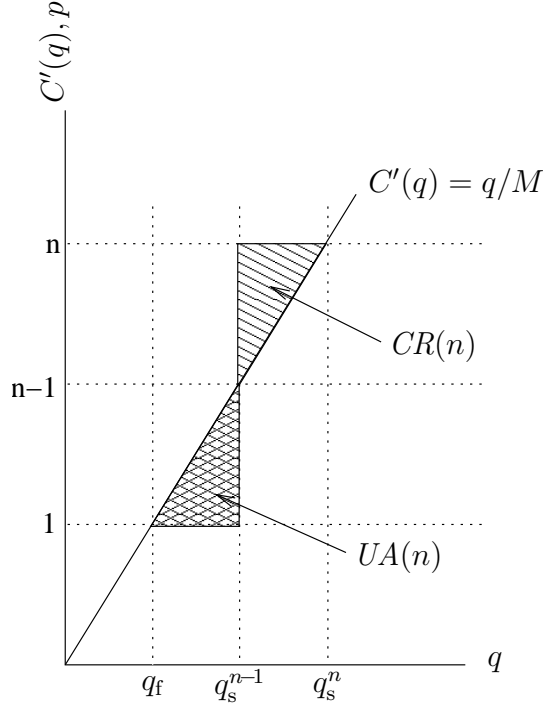


FIGURE 1. Linear benefits and linear marginal costs.

the rectangle with boundaries  $q_f$  and  $q_s^{n-1}$  and height  $n$  minus the area under the marginal cost curve between  $q_s^{n-1}$  and  $q_s^n$ . The difference is the area of the hatched triangle labeled  $CR(n)$ . This second triangle has unit height, regardless of  $n$ ; the height of the first triangle is  $n - 2$ . Since the two triangles are similar, the area  $CR(n) \geq UA(n)$  iff  $n \leq 3$ . The equilibrium condition (3) is that  $CR(n) \geq UA(n)$  while  $CR(n+1) < UA(n+1)$ . We have thus established that this condition is satisfied iff  $n = 3$ .<sup>6</sup>

When three countries form an IEA, the equilibrium global gain is  $6M(N - 2)$ ; this is a fraction  $\frac{12(N-2)}{N(N^2-2N+1)}$  of the total potential gain that could be achieved if all  $N$  countries joined the agreement.<sup>7</sup> This fraction is independent of the cost parameter  $M$  and decreases with  $N$  for  $N > 3$ . Although an increase in  $M$  lowers total and marginal costs and increases global welfare, technology enhancements that reduce  $M$  have no effect on the fraction of the potential welfare gain that is achieved in equilibrium.

We prove in proposition 1 in section 3 that an equilibrium size exceeding  $n = 3$  cannot be sustained if the *marginal* cost function is convex, i.e., if  $C'''(\cdot) \geq 0$ . Membership can exceed  $n = 3$  if this restriction is relaxed, allowing the marginal cost function to be locally concave beyond  $q = q_s^{\bar{n}-1}$ . (Note that conventional economic theory requires only that the *cost* function is convex; no restrictions are imposed on marginal costs except that they be non-decreasing.) Fig. 2 illustrates that for any  $3 < \bar{n} \leq N$ , a stable equilibrium of size  $\bar{n}$  can be implemented by a cost-reducing innovation. In the figure, the modified cost function is,  $\tilde{C}(\cdot)$ , whose derivative is piecewise linear and agrees with  $C'(\cdot)$  for  $q \leq q_s^{\bar{n}-1} = (\bar{n} - 1)M$ , but

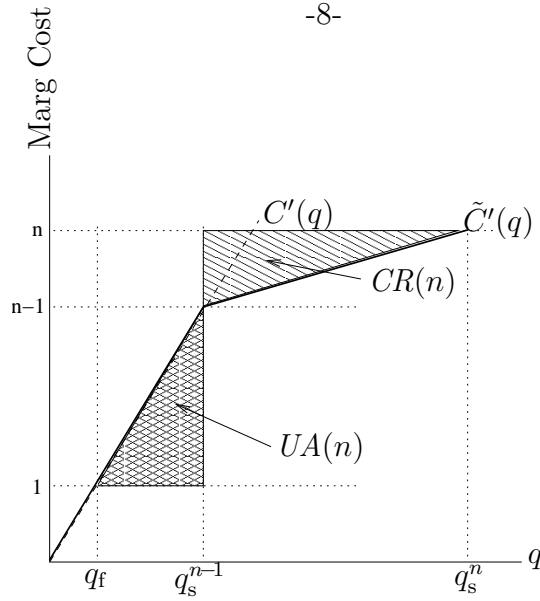


FIGURE 2. Linear benefits and kinked marginal costs.

is flatter than  $C'(\cdot)$  for larger values of  $q$ . Specifically, let

$$\tilde{C}'(q) = \begin{cases} q/M & \text{if } q \leq (\bar{n} - 1)M \\ c_0 + q/K & \text{if } q > (\bar{n} - 1)M \end{cases},$$

where  $K > M$  and  $c_0 = (\bar{n} - 1)(1 - M/K)$ . ( $c_0$  is chosen so that  $\tilde{C}'(\cdot)$  is continuous at  $q = q_s^{\bar{n}-1} = (\bar{n} - 1)M$ .) The new marginal cost curve  $\tilde{C}'(\cdot)$  is depicted in Fig. 2 by a heavy solid line; the original one by a dashed line.<sup>8</sup> Note that the area of  $CR(\bar{n})$  just exceeds the area of  $UA(\bar{n})$ : the shorter height of the first triangle is now more than offset by its longer base. The cost reduction from  $C(\cdot)$  to  $\tilde{C}(\cdot)$  leaves individual abatement unchanged at  $M$  if no IEA is implemented; in this case the reduction has no effect on global welfare. The cost reduction can however increase participation in an IEA. Indeed, if the reduction is sufficiently great—specifically, if  $K \geq M(N - 2)^2$  (see footnote 8)—an agreement to which all  $N$  countries are signatories is stable, and all of the potential gains from cooperation under this cost structure are realized.

### 3. CONVEX COSTS AND LINEAR BENEFITS

In this section we generalize the model to incorporate arbitrary convex cost functions, while maintaining the restriction that benefits are linear in abatement. We begin by showing that if the *marginal* cost function is initially strictly convex in abatement, then an IEA with more than two members is unsustainable. Once this restriction is relaxed (maintaining convexity of the *cost* function), an IEA of any size can be sustained. Our approach is to fix  $n \leq N$  and a given strictly convex cost function for which  $n$  is not a stable equilibrium, then lower this function to the point at which a stable equilibrium of size *at least*  $n$  exists. When benefits are linear in abatement, there are countless ways of attaining this goal: we identify the “smallest possible” reduction in the cost function that can accomplish it. This result is formalized in Prop. 2. Our motivation for this exercise is that cost-reducing technological progress are more

difficult to achieve, the greater is the extent to which costs must be reduced. Our “smallest possible” reduction can be interpreted as the least resource-intensive way to implement a coalition with at least  $n$  members.

The cost reduction we construct on p. 15 below (see (23) and (24)) implements a coalition with exactly  $n$  members iff under the initial cost function, the joiner’s gain function  $g(\cdot)$  is more negative at  $n+1$  than at  $n$ . If  $g(n+1) \geq g(n)$ , this reduced cost function will satisfy only the internal, but not the external stability condition at  $n$  (see (3)). However,

$$\text{if } n \text{ is internally stable, there exists } n' \geq n \text{ such that } n' \text{ is a stable equilibrium.} \quad (8)$$

If  $n$  is also externally stable, then by definition  $n$  is stable. Assume therefore, that  $n$  is not externally stable and let  $\underline{N} = \{n' \in (n+1, N] \cap \mathbb{N} : \tilde{g}(n') < 0\}$ , where  $\tilde{g}(\cdot)$  is the joiner’s gain function corresponding to the reduced cost function. If  $\underline{N}$  is non-empty, let  $\underline{n}$  be its smallest element. In this case, since  $\tilde{g}(\underline{n} - 1) > 0 > \tilde{g}(\underline{n})$ ,  $n' = \underline{n} - 1$  is both internally and externally stable. On the other hand, if  $\underline{N}$  is empty, then  $\tilde{g}(N) \geq 0$  and  $N$  will be both internally and externally stable, because the external stability condition is satisfied vacuously at  $N$ .

A straightforward extension of the methodology we develop in this section shows that under the same conditions on costs and benefits, it is possible to implement a coalition with exactly  $n$  members even when  $g(n+1) > g(n)$ . We establish this in Prop. 3, by reducing costs even further than in the construction (23) that we use to prove Prop. 2. The additional reduction differs in one significant respect from the original one: the construction (23) reduces costs while leaving unchanged the three abatement levels,  $q_s^n$ ,  $q_s^{n-1}$  and  $q_f$ , that are arguments of  $g(n)$ . This property enables us to obtain a simple expression for the “smallest possible” cost reduction such that internal stability is satisfied at  $n$ . When  $g(n+1) > g(n)$ , however, it is necessary to increase  $q_s^n$  in order to satisfy external stability at  $n$ . In this case, it no longer possible to obtain a closed-form expression for the smallest possible cost reduction.

For the rest of this section, we fix a differentiable, strictly convex cost function  $C$ , and let  $S(\cdot)$  denote the inverse of  $MC(\cdot)$ . When the benefit function is linear, the functions  $UA(n)$  and  $CR(n)$  defined above ((5c) and (5b)) can be rewritten in a particularly convenient form:

$$UA(n) = C(q_s^{n-1}) - C(q_f) - (q_s^{n-1} - q_f) = (n-2)q_s^{n-1} - \int_1^{n-1} S(p)dp \quad (9a)$$

$$CR(n) = n(q_s^n - q_s^{n-1}) - [C(q_s^n) - C(q_s^{n-1})] = \int_{n-1}^n S(p)dp - q_s^{n-1}. \quad (9b)$$

Geometric intuition for this reformulation is provided below.

**3.1. A pessimistic result when marginal cost is convex.** Figure 3 provides intuition for why the formulations in (5) and (9) are equivalent provided benefits are linear. The upward sloping curve in Fig. 3 represents both the marginal cost curve  $MC'(q)$  and its inverse,  $S(p)$ . In our discussion of Fig. 1 on p. 7, the cross-hatched area  $UA(n)$  was identified as the area under the marginal cost curve between  $q_f$  and  $q_s^{n-1}$ , minus the area of the rectangle with

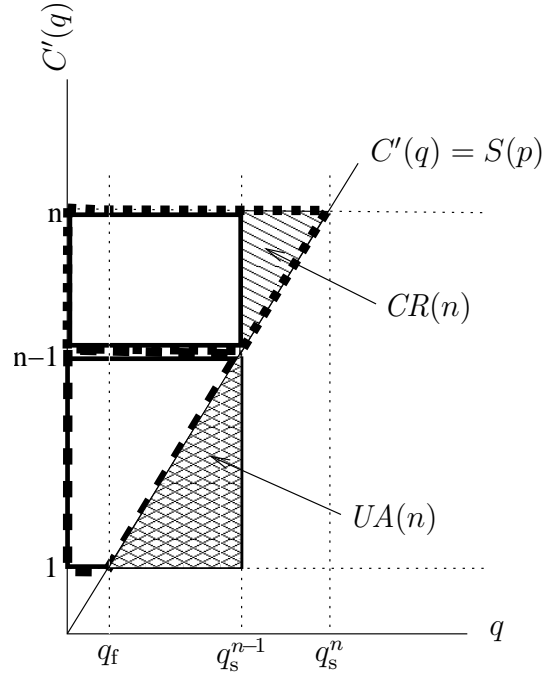


FIGURE 3. Linear benefits and linear marginal costs.

vertical boundaries  $q_f$  and  $q_s^{n-1}$  and height 1. As Fig 3 illustrates, this area is also equal to the area of the rectangle bounded vertically by 0 and  $q_s^{n-1}$ , and horizontally by 1 and  $n-1$ , minus the trapezoid—marked in the figure by widely spaced dashes—to the left of  $S(p)$  between 1 and  $n-1$ . The area of the difference between the trapezoid and the rectangle is given by the right-hand side of (9a). The hatched area  $CR(n)$  was identified on p. 7 as equal to the area of the rectangle with vertical boundaries  $q_s^{n-1}$  and  $q_s^n$  and height  $n$  minus the area under the marginal cost curve between  $q_s^{n-1}$  and  $q_s^n$ . It is also equal to the trapezoid—marked in the figure by closely spaced dashes—to the left of the curve  $S(p)$  between  $n-1$  and  $n$ , minus the area of the rectangle between 0 and  $q_s^{n-1}$  with height one. The area of the difference between this second rectangle and smaller trapezoid is equal to the right-hand side of (9b).

We use (9) to simplify expression (4) for the joiner's gain function  $g(n) = CR(n) - UA(n)$ :

$$\begin{aligned}
 g(n) &= \int_{n-1}^n S(p)dp - q_{in}^{n-1} - \left[ (n-1-1)q_{in}^{n-1} - \int_1^{n-1} S(p)dp \right] \\
 &= \int_1^n S(p)dp - (n-1)q_{in}^{n-1} = \int_1^n S(p)dp - (n-1)S(n-1). \quad (10)
 \end{aligned}$$

Using expression (10), Prop. 1 establishes that when marginal costs are strictly convex, there cannot exist a stable equilibrium with more than two signatories.<sup>9</sup>

**Proposition 1.** *If the marginal cost function is strictly convex, so that  $S(p)$  is strictly concave, then the largest stable equilibrium is less than or equal to 2.*

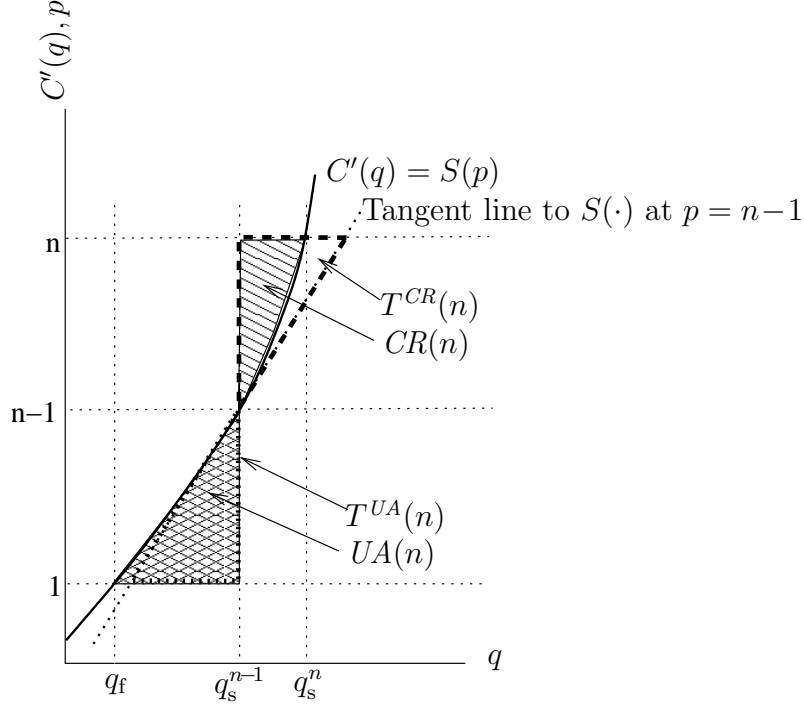


FIGURE 4. Intuition for Prop. 1

The Appendix contains the formal proof, and Fig. 4 provides the intuition. Since  $MC'(\cdot)$  is strictly convex (i.e.,  $S(\cdot)$  is strictly concave), its graph lies everywhere except at  $p = n-1$  strictly to the left of the line that is tangent to  $S(\cdot)$  at  $p = n-1$ . Therefore the hatched region  $CR(n)$  is strictly contained in the triangle  $T^{CR}(n)$ , indicated by dashed lines, which is bounded below by  $n-1$  and above by  $n$ , and so has height 1. Since  $T^{CR}(n)$  is similar to the triangle  $T^{UA}(n)$ , marked by dotted lines, with height  $n-2$ , the area of  $T^{UA}(n)$  weakly exceeds that of  $T^{CR}(n)$  for all  $n \geq 3$ . Moreover, the cross-hatched region  $UA(n)$  strictly contains  $T^{UA}(n)$ . It follows that for all  $n \geq 3$ ,  $UA(n) > CR(n)$ , verifying that  $g(n) < 0$ , for every integer  $n > 2$ . Prop. 1 now follows from (3).

The intuition underlying Prop. 1 can be extended to provide a helpful insight into the relationship between the internal and external stability conditions. Say that a function  $f : \mathbb{R}_+ \rightarrow \mathbb{R}_+$  is *convex beyond*  $x$  if  $f(\cdot)$  is convex on the interval  $[x, \infty]$ . It will be clear by inspection of Fig. 4 that if for some  $\bar{n}$ , the marginal cost curve is convex beyond  $MC^{-1}(\bar{n})$ , then as  $n$  increases beyond  $\bar{n}$ , the breadth of  $CR(n)$  at its base will either remain constant or decrease. Its height, however, will remain constant at 1. On the other hand, both the height and breadth of  $UA(n)$  will increase with  $n$ . It follows immediately that

$$\begin{aligned} \text{for any } \bar{n} \in \mathbb{N}, \text{ if marginal cost is convex beyond } MC^{-1}(\bar{n}), \text{ then} \\ \text{the joiner's gain function } g(\cdot) \text{ is strictly decreasing beyond } \bar{n}. \end{aligned} \quad (11)$$

An immediate consequence of (11) is that if marginal cost is convex beyond  $MC^{-1}(\bar{n})$ , then the external stability requirement is satisfied beyond  $\bar{n}$  whenever the internal stability

requirement is satisfied with equality (i.e., for  $n \geq \bar{n}$ ,  $g(n) = 0$  implies  $g(n + 1) < 0$ ). Summarizing

$$\text{if MC is convex beyond } \text{MC}^{-1}(\bar{n}), \text{ then } n \geq \bar{n} \text{ is stable if } g(n) = 0. \quad (12)$$

In Prop. 2 below, we reduce a given initial cost function in a way that leaves the original cost function unchanged beyond  $\bar{n}$ , while the modified joiner's gain function is zero at  $\bar{n}$ . It follows from (12) that if the marginal cost were originally convex, then our construction is sufficient to ensure that  $\bar{n}$  is stable. This fact is significant because in the literature on participation games to date, costs are almost invariably parameterized so that marginal costs are convex. Indeed, we are unaware of *any* paper in the related literature in which this property is violated.

**3.2. Admissible cost reductions.** We now consider the effect on the equilibrium size of an IEA of a technological improvement that lowers the abatement cost function. We model this kind of improvement by adding to an initially specified inverse marginal cost function  $S(\cdot)$  a non-negative function that is positive on an open subset of its domain. Observe that a (weak) increase in  $S(\cdot)$ , for all  $p$ , implies a (weak) decrease in marginal costs, for all  $q$ .

Let  $\varepsilon(p) \geq 0$  denote a *cost reduction function*. The new (lower) inverse marginal cost function is denoted by  $\tilde{S}(\cdot) = S(\cdot) + \varepsilon(\cdot)$ . In an  $n$ -member IEA, the per-member equilibrium level of abatement is  $S(n)$  before the cost reduction, and  $\tilde{S}(n)$  after it. Thus,  $\varepsilon(n)$  is the amount by which the cost reduction increases per-member abatement for an  $n$ -member IEA. Let  $\tilde{g}(\cdot)$  denote the joiner's gain function under the modified cost structure  $\tilde{S}(\cdot)$ .

We say that a cost reduction function  $\varepsilon(\cdot)$  is *admissible* if it satisfies certain properties. To streamline the exposition, we allow  $\varepsilon(\cdot)$  to be discontinuous, in which case there will be  $p$ -values at which the optimal abatement choice is not uniquely defined.<sup>10</sup> To resolve indeterminacies, we assume:

$$\begin{aligned} &\text{If multiple abatement levels satisfy the first order condition (1)} \\ &\quad \text{each non-signatory chooses the } \textit{lowest} \text{ optimal abatement level} \end{aligned} \quad (13a)$$

$$\quad \text{each signatory chooses the } \textit{highest} \text{ optimal abatement level.} \quad (13b)$$

From the exposition which follows, it will be clear that we could have restricted  $\varepsilon(\cdot)$  to be continuous, and thus eliminated all indeterminacies. This would have required more notation, increased the complexity of proofs, and made our heuristic explanations more cumbersome, without changing any of our formal results.

We now define the conditions for admissibility:

**Definition 1.** *For given  $S(\cdot)$ , a cost reduction function  $\varepsilon(p)$  is admissible if*

$$\varepsilon(p) \geq 0, \text{ with strict inequality on some open subset of } \mathbb{R} \quad (14a)$$

$$S(\cdot) + \varepsilon(\cdot) \text{ is a non-decreasing function of } p. \quad (14b)$$

Inequality (14a) implies that the cost reduction increases abatement for some levels of membership; inequality (14b) states that marginal cost after the reduction is nondecreasing, so that costs remain convex after the cost reduction. For convenience, we add a third restriction:

$$\varepsilon(p) = 0 \text{ for all } p \in (N + 1, \infty). \quad (14c)$$

Because the largest possible IEA has  $N$  members,  $q_s^n$  never exceeds  $S(N)$ . Hence marginal cost reductions at levels of  $q > S(N)$  can have no impact on equilibrium outcomes, so restriction (14c) is without loss of generality.

Condition (14b) has an implication that is significant in what follows.<sup>11</sup>

$$\text{if } \varepsilon(\cdot) \text{ satisfies (14) and } \varepsilon(p) > 0 \text{ then there exists } \delta > 0 \text{ s.t. } \varepsilon(\cdot) > 0 \text{ on } (p, p + \delta). \quad (15)$$

We will invoke property (15) in the proof of Prop. 2.

Replacing  $S(\cdot)$  with  $\tilde{S}(\cdot)$  in (10), we see that the change in the joiner's gain function at  $n$  due to the cost reduction  $\varepsilon(\cdot)$  is

$$\Delta^\varepsilon g(n) = \tilde{g}(n) - g(n) = \int_1^n \varepsilon(p) dp - (n - 1)\varepsilon(n - 1). \quad (16a)$$

Similarly, the change in  $g(\cdot)$  at  $n+1$  is

$$\Delta^\varepsilon g(n+1) = \int_1^{n+1} \varepsilon(p) dp - n\varepsilon(n) \quad (16b)$$

$$= \Delta^\varepsilon g(n) - n\varepsilon(n) + (n - 1)\varepsilon(n - 1) + \int_n^{n+1} \varepsilon(p) dp. \quad (16c)$$

For any  $3 \leq n \leq N$ , we identify in subsection 3.3 the “smallest possible” or “most efficient” reduction in abatement costs such that an IEA with at least  $n$  members is stable. Because we consider cost reductions over an interval, there are many alternative definitions of “smallest possible.” We adopt a natural specification, evaluating  $\varepsilon(\cdot)$ 's that satisfy conditions (14) according to the following norm:

$$\|\varepsilon\| = \int_0^{N+1} \varepsilon(p) dp; \quad (17)$$

that is, we take the integral over the change in inverse marginal costs from 0 to  $N + 1$ . For  $\varepsilon(\cdot)$ 's satisfying restriction (14c), this norm is equivalent to one that computes the integral of marginal costs (cf. inverse marginal costs) over an interval of  $q$ 's containing the set on which the new and the old marginal cost curves (cf. the inverse marginal cost curves) are permitted to differ.

**3.3. Cost decreases that increase equilibrium membership.** Fix  $\bar{n} > 2$  and a marginal cost function whose inverse is  $S$ . The integer  $\bar{n}$  may be unstable under  $S$  because it fails the test either for internal stability—i.e.,  $g(\bar{n}) < 0$ —or for external stability—i.e.,  $g(\bar{n}+1) \geq 0$ . The pessimistic results in the literature to date are all consequences of the fact that for the functional forms which that literature considers, it is the internal rather than the external condition that is violated for large  $n$ . Accordingly, to limit the number of cases we need to consider we assume that under the original cost structure, the internal stability condition fails at  $\bar{n}$ , i.e.,  $g(\bar{n})$  is negative.

Our goal, then, is to shift  $g(\cdot)$  up at  $\bar{n}$ . From (4),  $g(\cdot)|_{\bar{n}}$  decreases with  $q_s^{\bar{n}-1}$ ; moreover,  $q_s^{\bar{n}-1}$  increases as  $MC(\cdot)|_{q_s^{\bar{n}-1}}$  shifts down, so that a cost modification which lowers  $MC(\cdot)$  on a neighborhood of  $q_s^{\bar{n}-1}$  will, holding all else constant, shift  $g(\cdot)|_{\bar{n}}$  even further down. On the other hand, equation (4) also reveals that  $g(\cdot)$  shifts up at  $\bar{n}$  as the cost of producing  $q_s^{\bar{n}}$  decreases. Accordingly, the cost modification we construct below leaves  $MC(\cdot)$  unchanged at  $q_s^{\bar{n}-1}$  itself, while lowering it on an interval to the left of  $q_s^{\bar{n}-1}$ , thus lowering the cost of producing  $q_s^{\bar{n}}$ , while  $q_s^{\bar{n}-1}$  remains constant. As Fig 4 illustrates and Prop. 1 formalizes, this reduction must be sufficiently large that the modified marginal cost curve must on some interval to the left of  $q_s^{\bar{n}-1}$  lie below the tangent line to  $MC(\cdot)$  at  $q_s^{\bar{n}-1}$ —in other words the modified marginal cost curve must be locally concave on some interval—otherwise  $g(\cdot)$  will necessarily be negative for all  $n \geq 3$ . (As we note on p. 16, we could also increase  $g(\cdot)$  at  $\bar{n}$  to some extent by lowering marginal cost between  $q_s^{\bar{n}-1}$  and  $q_s^{\bar{n}}$ ; there is, however, no efficiency gain to this alternative approach.)

Before proceeding with our construction, we note some key inequalities and introduce some terminology. Given a cost reduction function  $\varepsilon(\cdot)$ , a necessary and sufficient condition for  $\bar{n}$  to be stable given the reduced cost structure  $\tilde{S}(\cdot) = S(\cdot) + \varepsilon(\cdot)$  is that

$$\tilde{g}(\bar{n}) = g(\bar{n}) + \Delta^\varepsilon g(\bar{n}) \geq 0 > g(\bar{n}+1) + \Delta^\varepsilon g(\bar{n}+1) = \tilde{g}(\bar{n}+1). \quad (18)$$

It follows from (18) that  $\bar{n}$  will be stable given a cost reduction  $\varepsilon(\cdot)$  iff

$$\Delta^\varepsilon g(\bar{n}) \geq -g(\bar{n}) > 0 \quad (19a)$$

$$g(\bar{n}+1) + \Delta^\varepsilon g(\bar{n}+1) < 0. \quad (19b)$$

We established at the beginning of section 3 that if (19a) is satisfied at  $\bar{n}$ , then there exists  $n' \geq \bar{n}$  such that  $n'$  is stable. Accordingly we say that

$$\varepsilon(\cdot) \text{ implements at least } \bar{n} \text{ if } \varepsilon(\cdot) \text{ satisfies (14) and if } \Delta^\varepsilon g(\bar{n}) \text{ satisfies (19a)}. \quad (20a)$$

$$\varepsilon(\cdot) \text{ implements exactly } \bar{n} \text{ if } \varepsilon(\cdot) \text{ satisfies both (19a) and (19b)} \quad (20b)$$

That is, to implement at least  $\bar{n}$ , it is necessary only to satisfy internal stability at  $\bar{n}$ ; to implement exactly  $\bar{n}$ , external stability must be satisfied as well. The cost reduction we construct below satisfies internal stability with equality, i.e., sets  $\Delta^\varepsilon g(\bar{n}) = -g(\bar{n})$ . As noted on p. 11 (see (12)), if the initial cost function belongs to a class that has been considered in the related literature to date and internal stability is satisfied with equality, then external

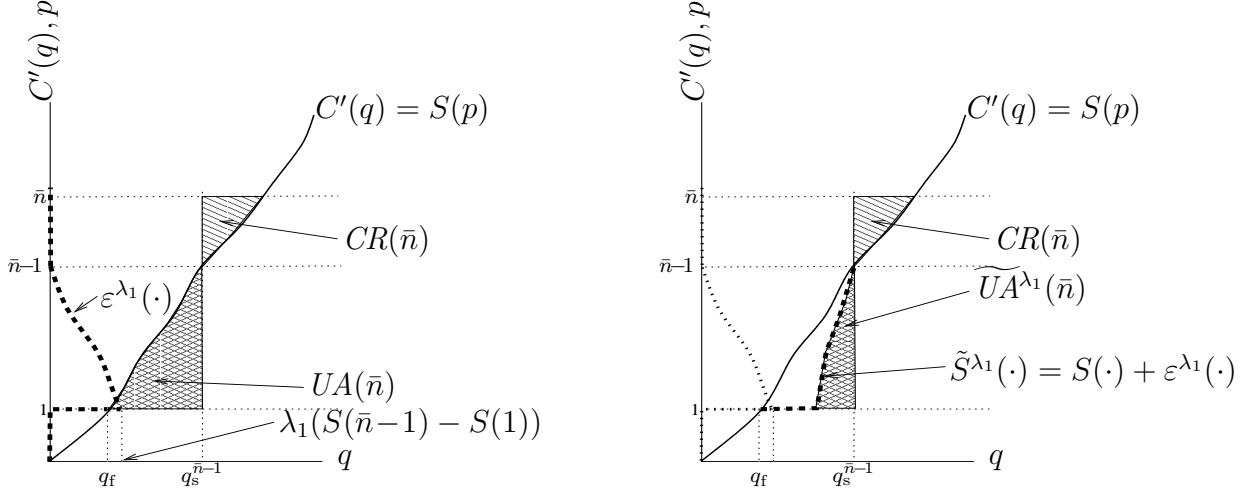


FIGURE 5. Shrinking  $UA(\bar{n})$  increases  $\tilde{g}(\bar{n})$ .

stability is satisfied as well. For such cost functions, the distinction between “at least” and “exactly” is moot.

There are countless admissible cost reductions that implement at least  $\bar{n}$ . We identify the minimal value of  $\|\varepsilon\|$ , taken over the set of all such reductions. When satisfying internal stability with equality does *not* imply external stability, there will not exist a minimal value of  $\|\varepsilon\|$  that implements exactly  $\bar{n}$ . Accordingly, we define the:

minimal cost reduction for at least  $\bar{n}$  as  $\epsilon_{\min}^{\bar{n}} = \min \{ \|\varepsilon\| : \varepsilon(\cdot) \text{ implements at least } \bar{n} \}$ . (21a)

infimal cost reduction for exactly  $\bar{n}$  as  $\epsilon_{\inf}^{\bar{n}} = \inf \{ \|\varepsilon\| : \varepsilon(\cdot) \text{ implements exactly } \bar{n} \}$ . (21b)

Finally, say that

$\varepsilon(\cdot)$  is *efficient for at least*  $\bar{n}$  if  $\varepsilon(\cdot)$  implements at least  $\bar{n}$  and  $\|\varepsilon\| = \epsilon_{\min}^{\bar{n}}$ . (22a)

$\varepsilon(\cdot)$  is *efficient for exactly*  $\bar{n}$  if  $\varepsilon(\cdot)$  implements exactly  $\bar{n}$  and  $\|\varepsilon\| = \epsilon_{\inf}^{\bar{n}}$ . (22b)

A cost reduction that is efficient for exactly  $\bar{n}$  exists only if satisfying internal stability with equality implies external stability.

The construction below identifies an  $\varepsilon(\cdot)$  that is efficient for at least  $\bar{n}$ . First, for each  $\lambda_1 \in [0, 1)$ , define

$$\varepsilon^{\lambda_1}(p) = \begin{cases} \lambda_1(S(\bar{n}-1) - S(p)) & \text{if } p \in [1, \bar{n}-1) \\ 0 & \text{otherwise} \end{cases}. \quad (23)$$

For  $\lambda_1 \approx 2/3$ , the graph of  $\varepsilon^{\lambda_1}(\cdot)$  is depicted in the left panel of Fig. 5 by heavy dashed lines. It is vertical until  $p = 1$ , at which point it jumps horizontally, then slopes negatively until  $p = (\bar{n}-1)$ , beyond which it is again vertical. In the right panel of the figure,  $\varepsilon^{\lambda_1}(\cdot)$  is added to  $S(\cdot)$  to form  $\tilde{S}^{\lambda_1}(\cdot)$ , which by construction agrees with  $S(\cdot)$  on  $[0, 1) \cup [\bar{n}-1, \infty)$ , is discontinuous at  $p = 1$ , then rises more steeply than  $S(\cdot)$  over the interval  $(1, \bar{n}-1]$ . Note

that as  $\lambda_1$  approaches 1,  $\lim_{\delta \searrow 0} \tilde{S}^{\lambda_1}(1 + \delta)$  approaches  $q_s^{\bar{n}-1}$ , while  $S^{\lambda_1}(\cdot)$  rises increasingly steeply on  $(1, \bar{n}-1]$ .

In the left panel,  $UA(\bar{n})$  is larger than  $CR(\bar{n})$ , so that  $g(\bar{n}) \equiv CR(\bar{n}) - UA(\bar{n}) < 0$ . In the right panel, the modified region  $\widetilde{UA}^{\lambda_1}(\bar{n})$  (the region above  $p = 1$ , below  $p = \bar{n}-1$ , right of  $\tilde{S}^{\lambda_1}(\cdot)$  and left of  $q_s^{\bar{n}-1}$ ) is considerably smaller. The location of  $q_f$  remains unchanged, in accordance with our assumption (13a) about how indeterminacies are resolved. It will be clear from the figure that when  $\lambda_1 = 0$ ,  $\tilde{S}(\cdot)$  is identical to  $S(\cdot)$ , so that the area of  $\widetilde{UA}^{\lambda_1}(\bar{n})$  exceeds that of  $CR(\bar{n})$ , while as  $\lambda_1$  approaches 1, the area of  $\widetilde{UA}^{\lambda_1}(\bar{n})$  shrinks to zero and is hence eventually smaller than that of  $CR(\bar{n})$ . Since the area  $\widetilde{UA}^{\lambda_1}(\bar{n})$  changes continuously with  $\lambda_1$  it follows from the intermediate value theorem that

$$\text{there exists } \lambda_1^* \in (0, 1) \text{ s.t. the areas of } \widetilde{UA}^{\lambda_1^*}(\bar{n}) \text{ and } CR(\bar{n}) \text{ are equal.} \quad (24)$$

Hence for the function  $\tilde{S}^{\lambda_1^*}$ ,  $\tilde{g}^{\lambda_1^*}(\bar{n}) = 0$ . From (16a),  $\Delta^{\lambda_1^*}g(\bar{n}) = \tilde{g}^{\lambda_1^*}(\bar{n}) - g(\bar{n})$ , so that  $\Delta^{\lambda_1^*}g(\bar{n}) = -g(\bar{n})$ . Moreover, by construction,  $\varepsilon^{\lambda_1^*}(\bar{n}-1) = 0$ . Hence from (16a) and (17),  $\|\varepsilon^{\lambda_1^*}\| = -g(\bar{n})$ , so that  $\varepsilon_{\min}^{\bar{n}}$ , the minimal cost reduction for at least  $\bar{n}$  is at most  $-g(\bar{n})$ . Part 1 of Prop. 2 establishes that in fact  $\varepsilon_{\min}^{\bar{n}} = -g(\bar{n})$ .

A property of the construction (23) is that  $\varepsilon^{\lambda_1}(p) = 0$ , for all  $p \geq \bar{n}-1$ . Hence from (16c),  $\Delta^{\lambda_1^*}g(\bar{n}) = \Delta^{\lambda_1^*}g(\bar{n}+1)$ . Since from (18),  $\tilde{g}^{\lambda_1^*}(\bar{n}+1) = g(\bar{n}+1) + \Delta^{\lambda_1^*}g(\bar{n}+1)$  while  $\tilde{g}^{\lambda_1^*}(\bar{n}) = 0 = g(\bar{n}) + \Delta^{\lambda_1^*}g(\bar{n})$ , we have that

$$\tilde{g}^{\lambda_1^*}(\bar{n}+1) = g(\bar{n}+1) - g(\bar{n}). \quad (25)$$

It now follows from (19b) and (25) that if  $g(\bar{n}+1) < g(\bar{n})$ , then  $\bar{n}$  is stable under  $\varepsilon^{\lambda_1^*}$ .

The construction above is by no means unique: if the heavy dashed line that bounds  $\widetilde{UA}^{\lambda_1}(\bar{n})$  on the left were replaced by *any* other non-decreasing line such that the area of the corresponding region  $UA(\bar{n})$  were equal to that of  $\widetilde{UA}^{\lambda_1}(\bar{n})$ , then the cost reduction function associated with that line would also be efficient for at least  $\bar{n}$ . There is also a second class of  $\varepsilon(\cdot)$ 's that have this property: functions in this class are positive on  $(\bar{n}-1, \bar{n})$  and thus increase the area of  $CR(\bar{n})$  rather than simply reducing the area of  $UA(\bar{n})$ . There is, however, a limit in general on how much can be accomplished by increasing  $CR(\bar{n})$  without also reducing  $UA(\bar{n})$ ; for example, if  $S(\cdot)$  were nearly affine on  $(\bar{n}-1, \bar{n})$ , then one could not do much better than double the area of  $CR(\bar{n})$ , which, if the initial difference between  $UA(\bar{n})$  and  $CR(\bar{n})$  were sufficiently large, would be alone insufficient to equalize the two areas. By contrast, regardless of the initial relative sizes of  $UA(\bar{n})$  and  $CR(\bar{n})$ , the area of  $UA(\bar{n})$  can always be reduce sufficiently to equalize them.

Prop 2 below summarizes and extends the preceding graphical analysis. The extensions are proved in the appendix. The result focuses exclusively on conditions under which the *internal* stability condition (19a) is satisfied. As we observed on p. 9 (statement (8)), if the internal stability condition is satisfied at  $\bar{n}$ , then both the internal and external stability conditions will be satisfied for some  $n' \geq \bar{n}$ . For this reason, the first three parts of Prop 2 relate to properties that hold for “at least  $\bar{n}$ .” The fourth part notes that when the external stability

constraint is slack—i.e., when  $g(\bar{n}+1) < g(\bar{n})$ —then a property will hold for at least  $\bar{n}$  if and only if it holds for exactly  $\bar{n}$ .

**Proposition 2.** *Let  $C$  be a twice differentiable, strictly convex cost function and let  $B$  be a linear benefit function. For  $1 < \bar{n} \leq N$ , assume that  $g(\bar{n}) < 0$  so that  $\bar{n}$  is not stable. Then*

- (1) *the minimal cost reduction for at least  $\bar{n}$  is  $-g(\bar{n})$ .*
- (2)  *$\exists \lambda_1^* \in (0, 1)$  s.t. the cost reduction  $\varepsilon^{\lambda_1^*}$  defined by (23) is efficient for at least  $\bar{n}$ .*
- (3)  *$\varepsilon(\cdot)$  is efficient for at least  $\bar{n}$  iff  $\int_1^{\bar{n}} \varepsilon(p) dp = -g(\bar{n})$  and  $\varepsilon\left([0, 1) \cup \{\bar{n}-1\} \cup [\bar{n}, \infty]\right) = 0$ .*
- (4) *if  $g(\bar{n}+1) < g(\bar{n})$  then  $\varepsilon(\cdot)$  is efficient for at least  $\bar{n}$  iff  $\varepsilon(\cdot)$  is efficient for exactly  $\bar{n}$ .*

An immediate corollary of part (1) of Prop. 2 is the intuitive result that if  $g(\cdot)$  is monotone decreasing in  $n$  over an interval, then larger cost reductions are required in order to implement larger coalitions in that interval. Note from part (3) that efficient cost reductions necessarily leave unchanged the original cost function beyond  $\bar{n}$ . Thus, if the initial *marginal* cost function is convex, then from (11) and part (4),  $\varepsilon(\cdot)$  is efficient for exactly  $\bar{n}$  iff it is efficient for at least  $\bar{n}$ . As we have emphasized above, marginal costs are indeed convex for the functions that have been considered in the literature on participation games to date, so that part (4) applies to these functions.

As we noted on p. 9, cost reductions that are efficient for at least  $\bar{n}$  will violate the external stability constraint whenever  $g(\bar{n}+1) \geq g(\bar{n})$ . If costs are reduced further, however, it is possible to implement exactly  $\bar{n}$  in this case. The construction required to achieve this is similar to our first reduction (23), although the analysis is less “clean” in several respects. First, we were able to satisfy the internal stability condition at  $\bar{n}$  while holding constant the abatement levels  $q_f$ ,  $q_s^{\bar{n}-1}$  and  $q_s^{\bar{n}}$ ; to satisfy external stability, on the other hand, it is necessary to increase  $q_s^{\bar{n}}$ . (We need a cost reduction  $\varepsilon(\cdot)$  s.t.  $\Delta^\varepsilon g(\bar{n}+1) < 0$ . From (16b), it is now necessary that  $\varepsilon(\bar{n}) > 0$ , since  $\Delta^\varepsilon g(\bar{n}+1)$  increases with  $\varepsilon(p)$  for all other values of  $p$ .) Second, because of the strict inequality in the external stability condition (19b), a minimal cost reduction no longer exists. Third, while it is possible to identify an infimal cost reduction, a closed-form expression for this infimum does not exist. For completeness, Prop 3 summarizes most of what can be said when for the original cost function,  $g(\bar{n}+1) \geq g(\bar{n})$ .

**Proposition 3.** *Let  $C$  be a twice differentiable, strictly convex cost function and let  $B$  be a linear benefit function. For  $1 < \bar{n} \leq N$ , assume that  $g(\bar{n}) \leq g(\bar{n}+1) < 0$ , so that  $\bar{n}$  is not stable. Then*

- (1) *there exists a cost reduction that implements exactly  $\bar{n}$ .*
- (2) *if  $g(\bar{n}+1) = g(\bar{n})$  (resp.  $g(\bar{n}+1) > g(\bar{n})$ ), then  $\epsilon_{inf}^{\bar{n}} = -g(\bar{n})$  (resp.  $\epsilon_{inf}^{\bar{n}} > -g(\bar{n})$ ).*

Props 2 and 3 highlight the importance of the relative magnitudes of  $g(\bar{n})$  and  $g(\bar{n}+1)$ , under the assumption that both are negative. If  $g(\bar{n}) > g(\bar{n}+1)$  the condition for external stability is slack, in the sense that internal stability guarantees external stability. If  $g(\bar{n}) < g(\bar{n}+1)$ , the external stability condition is binding. In the latter case, some cost reduction on an open interval above  $\bar{n}$  is required in order to satisfy the external stability condition.

We examined cost reductions that increase the size of a stable coalition. As illustrated by the example in subsection 2.1, some cost reductions actually reduce the equilibrium coalition size. We have proved analogs to Propositions 2 and 3 which identify the smallest cost reduction that renders unstable a previously stable coalition. Because there is less policy interest in these kinds of cost reductions, we have not included the analogs here.

#### 4. GENERAL COSTS, CONCAVE ABATEMENT BENEFITS

In this section we assume that the benefit function  $B(\cdot)$  is increasing and strictly concave in aggregate benefits. We continue to assume, as in section 3, that the abatement cost function  $C$  is initially strictly convex. To avoid dealing with nonnegativity constraints (see for example, Diamantoudi & Sartzetakis (2006), Rubio & Ulph (2006)) we assume

$$C'(0) < B'(Nq^*) \text{ where } q^* \text{ is the socially optimal level of abatement.} \quad (26)$$

Since signatory abatement levels are bounded above by  $q^*$  in any Cournot or Stackelberg equilibrium, this assumption guarantees that an interior solution to the non-signatory's first order condition (1b) exists.

When we relax the assumption of linear marginal benefits, our problem becomes more complex for two related reasons. First, we can no longer normalize a non-signatory's marginal benefit function to unity; nor can we normalize a signatory's marginal benefit function to the number of members of the IEA. A consequence is that our simplification (4) of expression (2) for the joiner's gain function is no longer valid; also invalid is our decomposition (5), derived from the simplification (4), of the joiner's gain function into  $UA$  and  $CR$ . Second, a non-signatory's optimization problem is no longer independent of the decisions made by signatories, so that for each  $n$ , the abatement choices made by signatories and non-signatories must be determined simultaneously. Moreover, if an additional member joins an IEA, so that total abatement by signatories increases, abatement levels by non-signatories will decrease in response. (Accordingly, we now denote non-signatory abatement when there are  $n$  signatories by  $q_f^n$  rather than  $q_f$ .)

Because of these complexities, the cost-reduction required to implement an IEA with at least  $n$  members is less straightforward than the one we constructed to prove Prop. 2. Moreover, in this context, it would be difficult to characterize the class of efficient cost reductions. Accordingly the scope of Prop. 4, presented in subsection 4.1 below, is narrower than its counterpart, Prop 2, for linear benefits. We have not attempted to identify the further cost-reduction that would be required to implement exactly  $n$ , so we have no counterpart to Prop 3. Another consequence of the additional complexity is that the Cournot and Stackelberg versions of our game no longer yield identical results. We will return this difference in subsection 4.2, where we establish that given a benefit and initial cost function, a cost-reduction that implements at least  $n$  as a Cournot equilibrium will also implement at least  $n$  as a Stackelberg equilibrium.

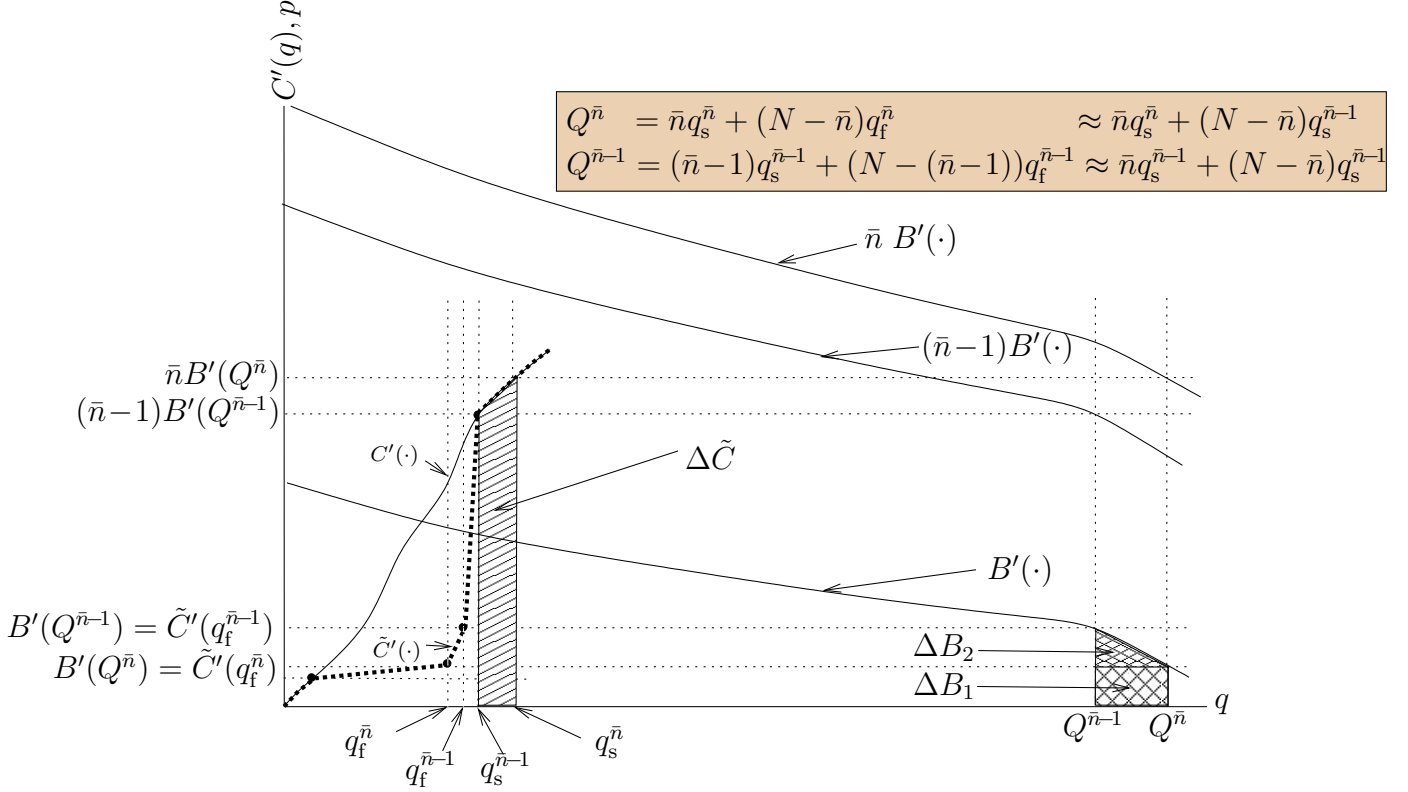


FIGURE 6. Implementing  $\bar{n}$  with concave benefits

4.1. **The Cournot game: a cost reduction that implements at least  $\bar{n}$ .** Fix a strictly increasing marginal cost function  $C'$ , a strictly decreasing marginal benefit function  $B'$ , and  $1 < \bar{n} \leq N$ . Fig. 6 illustrates a reduced marginal cost function  $\tilde{C}'$  which implements at least  $\bar{n}$  as a stable equilibrium. The original marginal cost curve  $C'$  is indicated by a thin solid line and the modified curve  $\tilde{C}'$  by a heavy dotted line. The two curves overlap near the origin, and beyond  $q_s^{\bar{n}-1}$ . The region where they differ consists of a segment that is arbitrarily close to horizontal, followed by one that is arbitrarily close to vertical.

We show that under the modified marginal cost function  $\tilde{C}'$ , the joiner's gain function  $\tilde{g}(\cdot)$  satisfies the internal stability condition (19a). From (8), both the internal and external stability conditions are satisfied for some  $n' \geq \bar{n}$ . The key to our construction is to ensure that for both  $\bar{n}-1$  and  $\bar{n}$ , the optimal levels of abatement for non-signatories, respectively  $q_f^{\bar{n}-1}$  and  $q_f^{\bar{n}}$ , are virtually the same as, but slightly smaller than  $q_s^{\bar{n}-1}$ ; we impose this by making modified marginal costs virtually vertical on the interval  $(q_f^{\bar{n}}, q_s^{\bar{n}-1})$  which contains  $q_f^{\bar{n}-1}$ .

Because  $q_s^{\bar{n}-1}$  and  $q_f^{\bar{n}-1}$  are arbitrarily close together under this construction, aggregate abatement with  $\bar{n}-1$  signatories is  $Q^{\bar{n}-1} \approx Nq_s^{\bar{n}-1}$ . Similarly, because  $q_s^{\bar{n}-1}$  and  $q_f^{\bar{n}}$  are arbitrarily close together, aggregate abatement with  $\bar{n}$  signatories is  $Q^{\bar{n}} \approx \bar{n}q_s^{\bar{n}} + (N - \bar{n})q_s^{\bar{n}-1}$ . Thus,  $Q^{\bar{n}} - Q^{\bar{n}-1} \approx \bar{n}(q_s^{\bar{n}} - q_s^{\bar{n}-1})$ . This property simplifies our analysis because, by construction, we have eliminated the complication that non-signatory abatement depend on  $\bar{n}$ .

A closely related consequence of our construction is that when there are  $\bar{n}-1$  signatories under the modified cost structure, it costs a signatory virtually the same to produce  $q_s^{\bar{n}-1}$  as it costs a non-signatory to produce  $q_f^{\bar{n}-1}$ . This implies that in the joiner's gain function  $\tilde{g}(\cdot)$  (see expression (2)), the term  $\tilde{C}(q_f^{\bar{n}-1})$  can be replaced by  $\tilde{C}(q_s^{\bar{n}-1})$ , so that  $\tilde{g}(\bar{n}) \approx \Delta B - \Delta \tilde{C}$ , where  $\Delta B$  is the per-signatory gain from increased abatements, and  $\Delta \tilde{C}$  is the per-signatory incremental cost, when the IEA acquires an additional signatory. To establish that  $\tilde{g}(\bar{n})$  is positive, we now only need to show that  $\Delta B$  exceeds  $\Delta \tilde{C}$ .

The magnitude of  $\Delta \tilde{C}$  (the tall thin hatched trapezoid in Fig. 6) is strictly less than the area of the rectangle that encloses it with height  $\bar{n}B'(Q^{\bar{n}})$  and width  $(q_s^{\bar{n}} - q_s^{\bar{n}-1})$ .  $\Delta B$  is the short wide cross-hatched trapezoid to the right of the figure, the union of the rectangle  $\Delta B_1$  and triangular-shaped  $\Delta B_2$ . The height of  $\Delta B_1$  is  $B'(Q^{\bar{n}})$ , i.e.,  $1/\bar{n}$  times the height of rectangle that contains  $\Delta \tilde{C}$ . The width of  $\Delta B_1$  is  $Q^{\bar{n}} - Q^{\bar{n}-1} \approx \bar{n}(q_s^{\bar{n}} - q_s^{\bar{n}-1})$  i.e.,  $\bar{n}$  times the width of the rectangle that contains  $\Delta \tilde{C}$ . Thus the area of  $\Delta B_1$  strictly exceeds that of  $\Delta \tilde{C}$ . Moreover, the trapezoid  $\Delta B$  also includes  $\Delta B_2$  which has positive area. Thus  $g(\bar{n}) = \Delta B - \Delta \tilde{C} > 0$ , verifying that  $\bar{n}$  is internally stable. This informal argument is made precise in Prop 4 below.

**Proposition 4.** *Assume that the benefit function  $B$  and initial cost function  $C$  are both twice continuously differentiable, that  $B$  is strictly concave and  $C$  is strictly convex. Then for any  $1 < \bar{n} \leq N$ , there exists an admissible cost reduction such that under the modified cost function  $\tilde{C}$ , an IEA with at least  $\bar{n}$  members is stable.*

**4.2. The Stackelberg game.** As noted above (p. 3), the literature has modeled the abatement stage of the participation game as both a Cournot and a Stackelberg game. So far, we have focused on the Cournot variant; we now turn to Stackelberg. In a simulation model with quadratic benefits and costs, Barrett (1994) shows that the equilibrium IEA in the Stackelberg game can be arbitrarily large. Diamantoudi & Sartzetakis (2006) counter by showing that when abatement is restricted to not exceed business-as-usual emissions, equilibrium membership cannot exceed four. With this reasonable restriction on abatement, the Stackelberg variant is thus only slightly less pessimistic than the Cournot. By contrast, with the cost reductions that we construct, any level of equilibrium participation is possible.

In the Stackelberg game, signatories treat  $q_f^n$  as a function of  $q_s^n$ . The FOC for a signatory is now

$$0 = n \frac{dB(Q^n)}{dq_s^n} - C'(q_s^n) = nB'(Q^n) \left(1 + \frac{dq_f^n}{dq_s^n}\right) - C'(q_s^n). \quad (27)$$

The FOC for Stackelberg non-signatories is the same as for Cournot non-signatories, i.e., (1b). Our assumption (26) ensures that (1b) has an interior solution, so that we can apply

the implicit function theorem to this expression to obtain  $\frac{dq_f^n}{dq_s^n}$ :

$$\frac{dq_f^n}{dq_s^n} = \frac{n B''(Q^n)}{C''(q_f^n) - (N - n)B''(Q^n)} \leq 0 \text{ with strict inequality if } B'' < 0. \quad (28)$$

We distinguish between the solutions to the Cournot and Stackelberg variants with the additional superscripts ‘C’ and ‘S’; Thus, for example  $q_s^{nC}$  denotes the solution to the signatory’s Cournot first order condition (1a) when the number of signatories is  $n$ .

To study the relationship between the two variants of the participation game, we first observe that when  $Q^{nS}(q_s^n) = nq_s^n + (N - n)q_f^n(q_s^n)$  and  $B$  is strictly concave, then

$$\frac{dQ^{nS}(q_s^n)}{dq_s^n} = n + (N - n)\frac{dq_f^n}{dq_s^n} = \frac{n C''(q_f^n)}{C''(q_f^n) - (N - n)B''(Q^n)} > 0. \quad (29)$$

Now compare the signatory’s first order conditions (1a) and (27) for the Cournot and Stackelberg variants. From (28),  $\frac{dq_f^n}{dq_s^n}$  is negative when  $B$  is strictly concave; hence the the right-hand side of the Stackelberg FOC (27) is negative when evaluated at the solution to the Cournot FOC (1a), i.e.,

$$0 > n\frac{dB(Q^{nC})}{dq_s^n} - C'(q_s^{nC}). \quad (30)$$

Summarizing, we have

$$q_s^{nS} < q_s^{nC}, \quad q_f^{nS} > q_f^{nC} \quad \text{and} \quad Q^{nS} < Q^{nC}. \quad (31)$$

The first inequality follows from (30) and the strict concavity of the signatory’s payoff in  $q_s^n$ , the second from (28) and the third from (29).

We can now sign the difference  $g^S(\cdot) - g^C(\cdot)$  between the Stackelberg and Cournot joiner’s gain functions. Clearly,  $q_f^{nS}(q_s^{nC}) = q_f^{nC}$ , so that for signatories to an  $n$ -member IEA, the Cournot benefit  $B(Q^{nC}) = B(nq_s^{nC} + (N - n)q_f^{nC})$  is attainable at the Cournot cost of  $C(q_s^{nC})$ . Since  $q_s^{nS} \neq q_s^{nC}$  and the signatory’s payoff is strictly concave, it follows that the first part of  $g^S(n)$ —i.e.,  $(B(Q^{nS}) - C(q_s^{nS}))$ —strictly exceeds the corresponding part of  $g^C(n)$ —i.e.,  $(B(Q^{nC}) - C(q_s^{nC}))$ . On the other hand, since from (31),  $Q^{(n-1)S} < Q^{(n-1)C}$  while  $q_f^{(n-1)S} > q_f^{(n-1)C}$ , the second part of  $g^S(n)$ — i.e.,  $(B(Q^{(n-1)S}) - C(q_f^{(n-1)S}))$ —is strictly less than the corresponding part of  $g^C(n)$ —i.e.,  $(B(Q^{(n-1)C}) - C(q_f^{(n-1)C}))$ . Hence  $g^S(n) > g^C(n)$ .

It follows that if there exists a cost-reduction such that for the resulting cost function  $\tilde{S}$ ,  $\tilde{g}^C(n) > 0$ , then for the same cost function,  $\tilde{g}^S(n) > 0$ . The preceding discussion thus proves the following proposition.

**Proposition 5.** *Assume that the benefit function  $B$  and initial cost function  $C$  are twice continuously differentiable, that  $B$  is strictly concave and  $C$  is strictly convex. Then for*

any  $1 < n \leq N$ , there exists an admissible cost reduction such that under the modified cost function  $\tilde{C}$ , an IEA with at least  $n$  members is stable in the Stackelberg variant of the participation game.

## 5. CONCLUSION

An important strand of the theory of IEAs focuses on the noncooperative Nash equilibria of the participation game introduced by D'Aspremont et al. (1983). The literature to date, which has focused exclusively on parametric examples, has been highly pessimistic about the prospects for international environmental cooperation. Its conclusion is that the equilibrium size of a stable IEA is small except when the potential gains from cooperation are also small. Moreover, although reductions in abatement costs necessarily increase the potential gains from cooperation, analysis of leading special cases has shown that such reductions leave unchanged or reduce equilibrium membership.

This paper approaches the problem from a non-parametric perspective. We first analyze the participation game assuming linear marginal benefits and general convex abatement costs. We introduce a novel decomposition of the gains to participation, and use it to show that regardless of the potential gains from cooperation, if *marginal* abatement costs are convex then the equilibrium coalition size cannot exceed three. Under the same assumptions, reductions in abatement costs cannot increase membership. In this respect, the results obtained in the previous literature are robust.

In an important respect, however, these previous results are fragile. When marginal abatement costs can be locally concave over a critical range of abatement levels, there are no restrictions on the size of the equilibrium coalition. Indeed, the noncooperative Nash equilibrium can realize up to 100% of the potential gains from cooperation, even when these potential gains are large. Moreover, we show that cost reductions *can* increase equilibrium membership and, when marginal benefits are linear, we derive a lower bound on the magnitude of the cost reduction needed to induce an arbitrary level of cooperation. Finally, we show that when benefits are strictly concave in abatement, it is still possible to support agreements among an arbitrarily large fraction of potential member countries.

Our analysis highlights the danger of drawing sweeping conclusions, however plausible, from parametric examples. In particular, our results undermine the apparently robust basis for pessimism that has become almost conventional wisdom in the field. On the other hand, one should not interpret our results as sufficient grounds for optimism about the prospects for real-world global agreements, since the kinds of cost functions we show to be consistent with such agreements might be difficult to design in practice.

## REFERENCES

- Barrett, S. (1994). Self-enforcing international environmental agreements, *Oxford Economic Papers* **46**: 878–894.
- Barrett, S. (2002). Consensus treaties, *Journal of Institutional and Theoretical Politics* **158**: 519–41.

- Barrett, S. (2003). *Environment and Statecraft*, Oxford University Press.
- Barrett, S. (2006). Climate treaties and breakthrough technologies, *American Economic Review* **96**: 22–25.
- Bloch, F. & Dutta, B. (2008). Formation of Networks and Coalitions, *Handbook of Social Economics, J. Benhabib, A. Bisin and M. Jackson, eds* .
- Bos, I. & Harrington Jr, J. (2010). Endogenous cartel formation with heterogeneous firms, *The RAND Journal of Economics* **41**(1): 92–117.
- Burger, N. & Kolstad, C. (2009). Voluntary public goods provision: coalition formation and uncertainty. NBER Working Papers 15543.
- Carraro, C. & Siniscalco, D. (1993). Strategies for the international protection of the environment, *Journal of Public Economics* **52**: 309–28.
- Chwe, M. S. (1994). Farsighted coalitional stability, *Journal of Economic Theory* **63**: 299–325.
- Dannenber, A., Lange, A. & Sturm, B. (2009). On the formation of coalitions to provide public goods – experimental evidence from the lab. Working paper.
- D’Aspremont, C., Jacquemin, A., Gabszewicz, J. J. & Weymark, J. A. (1983). On the stability of collusive price leadership, *The Canadian Journal of Economics / Revue canadienne d’Economie* **16**(1): pp. 17–25.
- De Bondt, R. (1997). Spillovers and innovative activities\* 1, *International Journal of Industrial Organization* **15**(1): 1–28.
- de Zeeuw, A. (2008). Dynamic effects on the stability of international environmental agreements, *Journal of Environmental Economics and Management* **55**(2): 163–74.
- Diamantoudi, E. & Sartzetakis, E. (2002). International environmental agreements - the role of foresight. UNiversity of Aarhus working paper 2002-10.
- Diamantoudi, E. & Sartzetakis, E. (2006). Stable international environmental agreements: An analytical approach, *Journal of Public Economic Theory* **8**(2): 247–263.
- Dixit, A. & Olson, M. (2000). Does voluntary participation undermine the coase theorem, *Journal of Public Economics* **76**: 309 – 335.
- Donsimoni, M., Economides, N. & Polemarchakis, H. (1986). Stable cartels, *International Economic Review* **27**(2): 317–327.
- Ecchia, G. & Mariotti, M. (1998). Coalition formation in international environmental agreements and the role of institutions, *European Economic Review* **42**(3-5): 573–582.
- Escribuela-Villar, M. (2009). A note on cartel stability and endogenous sequencing with tacit collusion, *Journal of Economics* **96**(2): 137–147.
- Eyckmans, J. (2001). On the farsighted stability of the Kyoto Protocol. CLIMNEG Working Paper 40, CORE, Universite Catholique de Louvain.
- Finus, M. (2001). *Game Theory and International Environmental Cooperation*, Cheltenham:Edward Elgar.
- Finus, M. (2003). 3. Stability and design of international environmental agreements: the case of transboundary pollution, *The international yearbook of environmental and resource economics 2003/2004: a survey of current issues* p. 82.
- Finus, M. & Maus, S. (2008). Modesty may pay!, *Journal of Public Economic Theory* **10**: 801–826.
- Hoel, M. (1992). International environmental conventions: the case of uniform reductions of emissions, *Environmental and Resource Economics* **2**: 141–59.
- Hoel, M. & De Zeeuw, A. (2010). Can a Focus on Breakthrough Technologies Improve the Performance of International Environmental Agreements?, *Environmental and Resource Economics* pp. 1–12.
- Hong, F. & Karp, L. (2010). International environmental agreements with mixed strategies and investment. unpublished Working Paper.
- Ioannidis, A., Papandreou, A. & Sartzetakis, E. (2000). International environmental agreements: A literature review, *Cahiers de Recherche du GREEN 00-08. Québec: Université Laval* .
- Katz, M. (1986). An analysis of cooperative research and development, *The Rand Journal of Economics* **17**(4): 527–543.
- Kohler, M. (2002). Coalition formation in international monetary policy games, *Journal of International Economics* **56**(2): 371–385.
- Kolstad, C. D. (2007). Systematic uncertainty in self-enforcing international environmental agreements, *Journal of Environmental Economics and Management* **53**: 68–79.

- Kolstad, C. D. & Ulph, A. (2008). Learning and international environmental agreements, *Climatic Change* **89**: 125–41.
- Kosfeld, M., Okada, A. & Riedl, A. (2009). Institution formation in public goods games, *American Economic Review* **99**: 1335–1355.
- Osmani, D. & Tol, R. (2009). Toward farsightedly stable international environmental agreements, *Journal of Public Economic Theory* **11**: 455–92.
- Ray, D. & Vohra, R. (2001). Coalitional power and public goods, *Journal of Political Economy* **109**(6): 1355 – 1382.
- Rubio, S. J. & Ulph, A. (2006). Self-enforcing international environmental agreements revisited, *Oxford Economic Papers* **58**: 233–263.
- Stiglitz, J. E. (2006). A new agenda for global warming, *Economists’ Voice* pp. 1–4.
- Thoron, S. (1998). Formation of a coalition-proof stable cartel, *Canadian Journal of Economics* **31**(1): 63–76.
- Ulph, A. (2004). Stable international environmental agreements with a stock pollutant, uncertainty and learning, *Journal of Risk and Uncertainty* **29**: 53–73.
- Xue, L. (1998). Coalitional stability under perfect foresight, *Economic Theory* **11**: 603–627.

## NOTES

<sup>1</sup>An alternative to the participation game that we use is based on games of coalition formation, such as Chwe (1994), Ecchia & Mariotti (1998), Xue (1998) and Ray & Vohra (2001). These models use a more sophisticated interpretation of rationality, in which agents understand how their provisional decision to join or leave a coalition would affect other agents’ participation decisions; nations are “farsighted”. Diamantoudi & Sartzetakis (2002), Eyckmans (2001), de Zeeuw (2008) and Osmani & Tol (2009) apply this notion to IEA models. An entirely distinct line of research models climate policy using a cooperative game.

<sup>2</sup>For many specifications of the model, (3) can be satisfied for  $n = 1$ . Since in this case there is no distinction between signatories and non-signatories, such solutions are clearly uninteresting; we ignore them throughout the paper.

<sup>3</sup> Since  $\lfloor n \rfloor = \lceil n \rceil$  for any integer  $n$ , (3′) will be violated whenever the root of  $g$  is an integer. A necessary and sufficient condition for *any* real-valued root  $n$  to be stable is that  $g(\lfloor n \rfloor) \geq 0 > g(\lceil n + \varepsilon \rceil)$  for sufficiently small  $\varepsilon > 0$ .

<sup>4</sup>In equilibrium,  $\lceil c \rceil$  countries each abate one unit at a cost of  $c$ , so that each of  $N$  countries benefit from  $Q = \lceil c \rceil$  units of abatement. Hence, global welfare is  $N\lceil c \rceil - \lceil c \rceil c = (N - c)\lceil c \rceil$ .

<sup>5</sup>See in particular our discussion of Barrett (2006) and Hoel & De Zeeuw (2010) on p. 2.

<sup>6</sup>To make the above argument precise, let  $C'(q) = \frac{q}{M}$ . Since marginal costs are linear,  $q_s^n - q_s^{n-1} = M$  which implies the area of  $CR(n) = \frac{M}{2}$  and  $\frac{n-1-1}{q_s^{n-1}-q_i^n} = \frac{1}{M}$ , so  $q_s^{n-1} - q_i^n = M(n-2)$  which implies that the area of  $UA(n) = \frac{M}{2}(n-2)^2$ . The root of  $g(n) = 0$  satisfies  $g(n) = \frac{M}{2} - \frac{M}{2}(n-2)^2 = 0$  or  $1 - (n-2)^2 = 0$  so the solutions are 3, 1. Since we have  $g'(n) = -M(n-2) < 0$  for  $n > 2$  we know that the root  $n = 3$  is stable, i.e.

$g(2) > g(3) = 0 > g(4)$ . This confirms that for linear marginal costs,  $n = 3$  is the unique non-trivial IEA size.

<sup>7</sup>The equilibrium level of aggregate abatement with the coalition is  $3(3M) + (N - 3)M = M(6 + N)$ , so the global benefit of abatement is  $MN(6 + N)$ . The aggregate equilibrium cost of abatement is  $(N - 3)\frac{M}{2} + 3\frac{9M}{2} = \frac{1}{2}M(N + 24)$ , so global equilibrium net benefits are  $MN(6 + N) - \frac{1}{2}M(N + 24) = \frac{1}{2}M(11N + 2N^2 - 24)$ . In the absence of a coalition, each nation abates at  $M$ , so global net benefits are  $N(NM - \frac{M}{2})$ . The global gain due to the formation of the coalition is therefore  $\frac{1}{2}M(11N + 2N^2 - 24) - N(NM - \frac{M}{2}) = 6M(N - 2)$ . If every nation were to join the coalition, aggregate welfare would be  $N\left(N^2M - \frac{N^2M}{2}\right) = \frac{1}{2}N^3M$ . The actual gain relative to the potential gain, i.e. the fraction of the potential gain from cooperation actually achieved, is

$$\frac{6M(N - 2)}{\frac{1}{2}N^3M - N\left(NM - \frac{M}{2}\right)} = 12\frac{N - 2}{N(N^2 - 2N + 1)}.$$

<sup>8</sup>We solve for  $K$  such that  $\bar{n}$  is a stable equilibrium. As in fn. 6, the area of  $UA(\bar{n}) = \frac{M}{2}(\bar{n} - 2)^2$ , while the area of  $CR(\bar{n})$  is now  $(q_s^n - q_f^n)/2 = K/2$ . Hence  $CR(\bar{n}) = UA(\bar{n})$  if  $K = (\bar{n} - 2)^2M$ . Moreover, since for  $n > \bar{n} - 1$ ,  $g(n) = \frac{M}{2}((\bar{n} - 2)^2 - (n - 2)^2)$ ,  $g'(\cdot) < 0$  on  $(\bar{n} - 1, \infty)$ . Thus  $g(\bar{n}) \geq 0 > g(\bar{n} + 1)$ , verifying that condition (3) for  $\bar{n}$  to be a stable equilibrium is satisfied.

<sup>9</sup>Using a convex marginal cost function which had previously been used to estimate the costs of greenhouse gas abatement, Barrett (1994)'s Proposition 4 shows that the equilibrium IEA consists of 2 members.

<sup>10</sup>A function  $\varepsilon(p)$  such that for some  $p$ ,  $\lim_{\delta \searrow 0} \varepsilon(p - \delta) = \underline{q} < \bar{q} = \lim_{\delta \searrow 0} \varepsilon(p + \delta)$  is consistent with a cost function  $C(q)$  which is affine on the interval  $[\underline{q}, \bar{q}]$ , i.e., marginal cost is flat on this interval. In this case, all of the  $q$ 's in  $[\underline{q}, \bar{q}]$  will be equi-profitable. Strictly speaking,  $S(p) + \varepsilon(p)$  is not the inverse of any marginal cost function, since viewed as a function of  $q$  rather than  $p$ ,  $S(q) + \varepsilon(q)$  is not one-to-one.

<sup>11</sup>Suppose that (15) were false, i.e., that for some  $p$  and every  $\delta > 0$ ,

$$\varepsilon(p) > 0 \quad \text{and} \quad \varepsilon(\cdot) = 0 \quad \text{on} \quad (p, p + \delta). \quad (32a)$$

The admissibility condition (14b) requires that for every  $\gamma > 0$

$$S(p + \gamma) + \varepsilon(p + \gamma) \geq S(p) + \varepsilon(p). \quad (32b)$$

Properties (32a) and (32b) can hold simultaneously only if  $S(\cdot)$  fails to be right-continuous at  $p$ , i.e., if  $\liminf_{\gamma \searrow 0} S(p + \gamma) \geq S(p) + \varepsilon(p) > S(p)$ . But in this case, the marginal cost function MC would be undefined on the interval  $(S(p), \liminf_{\gamma \searrow 0} S(p + \gamma))$ , contradicting the assumption  $S$  is the inverse of MC.

APPENDIX: PROOFS

**Proof of Proposition 1:** Let  $n = 2$  and let the tangent to the marginal cost curve, evaluated at  $q_s^2$ , be  $MC'(q_s^2) = \frac{1}{M}$ . The area of triangles  $T^{CR}(3)$  and  $T^{UA}(3)$  both equal  $\frac{M}{2}$ . Since by strict convexity,  $MC(q)$  lies strictly above the tangent line through  $q_s^2$ , for all  $q \neq q_s^2$ ,  $CR(3) < \frac{M}{2} < UA(3)$ . Therefore,  $g(3) = CR(3) - UA(3) < 0$ . Now assume  $n > 2$ . From (10),

$$\frac{dg(n)}{dn} = S(n) - S(n-1) - (n-1) \frac{dS(n-1)}{dn}.$$

But since  $S(\cdot)$  is strictly concave,  $S(n) < S(n-1) + \frac{dS(n-1)}{dn}$ . Hence

$$\begin{aligned} \frac{dg(n)}{dn} &< S(n-1) + \frac{dS(n-1)}{dn} - S(n-1) - (n-1) \frac{dS(n-1)}{dn} \\ &= -(n-2) \frac{dS(n-1)}{dn}, \end{aligned}$$

which is negative for  $n > 2$ . Since  $g(3) < 0$  and  $g(n)$  is strictly decreasing for  $n \geq 3$ , there can be no integer  $n > 2$  such that  $g(n) \geq 0$ . Therefore, there are no stable equilibria with  $n$  greater than 2. ■

**Proof of Proposition 2:** Part (4) has been proved in the text (see the line immediately following (25)). We also proved in the text that  $\varepsilon^{\lambda^*}$  implements at least  $\bar{n}$  and that  $\|\varepsilon^{\lambda^*}\| = -g(\bar{n})$ . To prove that  $-g(\bar{n})$  is the minimal cost reduction for at least  $\bar{n}$  (part (1)), it suffices to prove that  $\|\varepsilon\| \geq -g(\bar{n})$ , for any admissible cost reduction  $\varepsilon(\cdot)$  that implements at least  $\bar{n}$ . It will then follow immediately that  $\varepsilon^{\lambda^*}$  is efficient for at least  $\bar{n}$  (part (2)).

Assume now that  $\varepsilon(\cdot)$  implements at least  $\bar{n}$ . Substituting (16a) into (17), we obtain

$$\begin{aligned} \|\varepsilon\| &= \Delta^\varepsilon g(\bar{n}) + (\bar{n}-1)\varepsilon(\bar{n}-1) + \int_0^1 \varepsilon(p)dp + \int_{\bar{n}}^{N+1} \varepsilon(p)dp \\ &\underset{\text{from (19a)}}{\geq} -g(\bar{n}) + (\bar{n}-1)\varepsilon(\bar{n}-1) + \int_0^1 \varepsilon(p)dp + \int_{\bar{n}}^{N+1} \varepsilon(p)dp. \end{aligned} \quad (33)$$

Since admissibility requires  $\varepsilon(\cdot) \geq 0$ , (33) implies that  $\|\varepsilon\| \geq -g(\bar{n})$ , which is what we needed to prove. This completes the proof of parts (1) and hence (2).

To prove part (3), we first establish sufficiency. Assume that  $\int_1^{\bar{n}} \varepsilon(p)dp = -g(\bar{n})$  and that  $\varepsilon([0, 1) \cup \{\bar{n}-1\} \cup [\bar{n}, \infty]) = 0$ . Since  $\varepsilon(\bar{n}-1) = 0$ , (16a) implies  $\Delta^\varepsilon g(\bar{n}) = -g(\bar{n})$ , so that from (19a),  $\varepsilon$  implements at least  $\bar{n}$ . Since  $\varepsilon([0, 1) \cup [\bar{n}, \infty]) = 0$ , (17) implies  $\|\varepsilon\| = -g(\bar{n})$ . It now follows that  $\varepsilon(\cdot)$  is efficient for at least  $\bar{n}$ .

We now prove necessity. Assume that  $\varepsilon(\cdot)$  is efficient for at least  $\bar{n}$ , so that from part (1) of the proposition,  $\|\varepsilon\| = -g(\bar{n})$ . This, together with (33) implies

$$0 \geq (\bar{n}-1)\varepsilon(\bar{n}-1) + \int_0^1 \varepsilon(p)dp + \int_{\bar{n}}^{N+1} \varepsilon(p)dp. \quad (34)$$

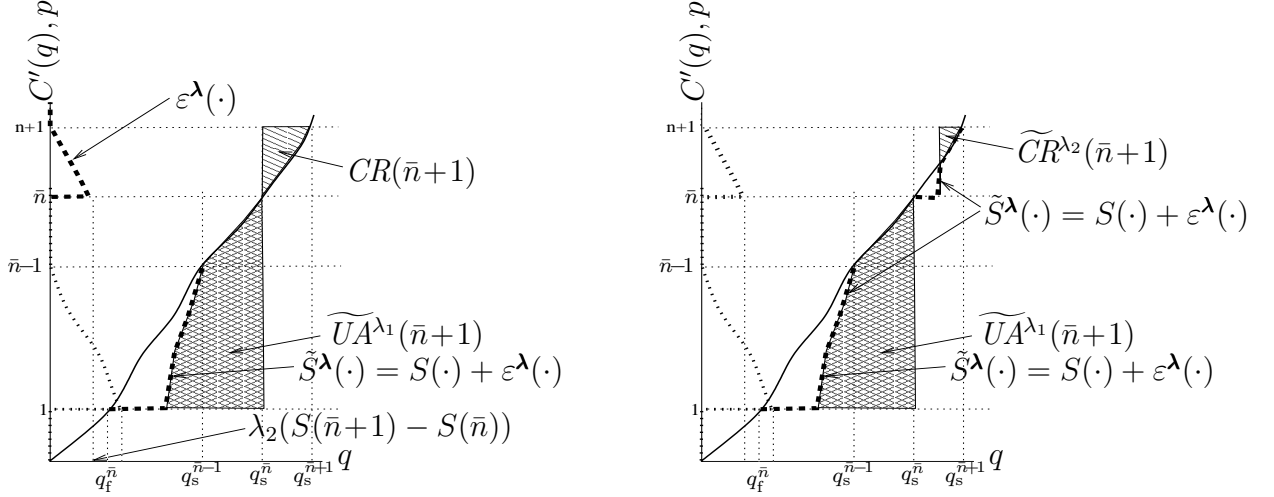


FIGURE 7. Shrinking  $CR$  reduces  $\tilde{g}(\bar{n}+1)$ .

But since  $\varepsilon(\cdot)$  is nonnegative, (34) implies

$$\varepsilon(\bar{n}-1) = \varepsilon((0, 1)) = \varepsilon((\bar{n}, N+1)) = 0. \quad (35)$$

(35) and (32a) imply  $\varepsilon(0) = \varepsilon(\bar{n}) = 0$ , verifying that  $\varepsilon([0, 1] \cup \{\bar{n}-1\} \cup [\bar{n}, \infty)) = 0$ . Hence  $-g(\bar{n}) = \|\varepsilon\| \equiv \int_0^{N+1} \varepsilon(p) dp = \int_1^{\bar{n}} \varepsilon(p) dp$ .  $\blacksquare$

**Proof of Proposition 3:** Let  $\lambda_1^*$  be defined as in display (24). The reduction  $\varepsilon^{\lambda_1^*}$  satisfies the internal stability condition (19a). From (25),  $\tilde{g}^{\lambda_1^*}(\bar{n}+1) = (g(\bar{n}+1) - g(\bar{n})) \geq 0$  so that a further cost reduction is required in order to satisfy the external stability condition (19b). The construction we use to obtain this reduction is similar of our previous one: earlier we shrank  $UA(\bar{n})$  until the area  $\widetilde{UA}^{\lambda_1^*}(\bar{n})$  equaled  $CR(\bar{n})$ ; now by assumption, the area of  $CR(\bar{n}+1)$  weakly exceeds that of  $\widetilde{UA}^{\lambda_1^*}(\bar{n}+1)$ , so we must shrink  $CR(\bar{n}+1)$ . Fig. 7 illustrates the second construction. To conserve space, we have graphed  $CR(\bar{n}+1)$  so that its area is actually *less than*  $\widetilde{UA}^{\lambda_1^*}(\bar{n}+1)$ —the reverse inequality could hold only if the initial cost curve flattened substantially between  $q_s^{\bar{n}}$  and  $q_s^{\bar{n}+1}$  (cf. Fig. 2). Nonetheless, the viewer of the figure is requested to imagine that  $CR(\bar{n}+1)$  is in fact weakly greater than  $\widetilde{UA}^{\lambda_1^*}(\bar{n}+1)$ .

Given  $\lambda_2 \in [0, 1)$ , let  $\lambda = (\lambda_1^*, \lambda_2)$  and extend the reduction  $\varepsilon^{\lambda_1^*}$  (see (23)) to  $\varepsilon^{\lambda}$  as follows:

$$\varepsilon^{\lambda}(p) = \begin{cases} \lambda_1^*(S(\bar{n}-1) - S(p)) & \text{if } p \in [1, \bar{n}-1) \\ \lambda_2(S(\bar{n}+1) - S(p)) & \text{if } p \in [\bar{n}, \bar{n}+1) \\ 0 & \text{otherwise} \end{cases} \quad (36)$$

The extension of  $\varepsilon^{\lambda_1^*}(\cdot)$  to  $\varepsilon^{\lambda}(\cdot)$  is depicted in the left panel of Fig. 7 by heavy dashed lines. It agrees with the graph of  $\varepsilon^{\lambda_1^*}$  for  $p < \bar{n}$ . At  $\bar{n}$  it jumps horizontally, then slopes negatively until  $p = \bar{n}+1$ , beyond which it is again vertical. In the right panel of the figure,  $\varepsilon^{\lambda}(\cdot)$  is

added to  $S(\cdot)$  to form  $\tilde{S}^\lambda(\cdot)$ , which by construction agrees with  $\tilde{S}^{\lambda_1}(\cdot)$  on  $[0, \bar{n}) \cup [\bar{n}+1, \infty)$ . when  $\lambda_2 = 0$ ,  $\tilde{S}^\lambda(\cdot)$  is identical to  $S^{\lambda_1}(\cdot)$ , so that the area of  $\widetilde{CR}^{\lambda_2}(\bar{n}+1)$  weakly exceeds that of  $\widetilde{UA}^{\lambda_1}(\bar{n}+1)$ . When  $\lambda_2 > 0$ , the tie-breaking rule (13b) implies that  $q_s^{\bar{n}+1}$  is discontinuous from below, i.e., jumps to the right. Since  $q_s^{\bar{n}+1}$  defines the vertical boundary dividing  $\widetilde{UA}^\lambda(\bar{n}+1)$  and  $\widetilde{CR}^{\lambda_2}(\bar{n}+1)$ , this boundary also moves to the right as  $\lambda_2$  increases. As  $\lambda_2$  approaches 1, the area of  $\widetilde{CR}^{\lambda_2}(\bar{n}+1)$  shrinks to zero, and the area of  $\widetilde{UA}^\lambda(\bar{n}+1)$  increases, so that the former is eventually strictly smaller than the latter. Since the two areas are continuous in  $\lambda_2$  it follows from the intermediate value theorem that there exists  $\lambda_2^* \in [0, 1)$  such that the areas of  $\widetilde{UA}^{\lambda_1}(\bar{n}+1)$  and  $\widetilde{CR}^{\lambda_2}(\bar{n}+1)$  are equal. Hence for  $\lambda_2 > \lambda_2^*$  and  $\boldsymbol{\lambda} = (\lambda_1^*, \lambda_2)$ ,  $\tilde{g}^\lambda(\bar{n}+1) < 0$  for the function  $\tilde{S}^\lambda$ . This proves part (1) of Prop. 3.

To verify part (2), observe that if  $g(\bar{n}+1) = g(\bar{n})$ , then  $\lambda_2^* = 0$ , so that for any  $\lambda_2 > 0$  and  $\boldsymbol{\lambda} = (\lambda_1^*, \lambda_2)$ ,  $\tilde{g}^\lambda(\bar{n}+1) < 0$ , i.e., external stability can be implemented for an infinitesimally small incremental cost. Thus if  $\lambda_2^* = 0$ , then  $\epsilon_{\inf}^{\bar{n}} = \epsilon_{\min}^{\bar{n}} = -g(\bar{n})$ . Now suppose that  $g(\bar{n}+1) > g(\bar{n})$  and consider an arbitrary cost reduction  $\varepsilon(\cdot)$  which satisfies (19). In this case

$$\Delta^\varepsilon g(\bar{n}+1) \underbrace{\leq}_{\text{from (19b)}} -g(\bar{n}+1) \underbrace{\leq}_{\text{by assumption}} -g(\bar{n}) \quad (37)$$

while from (16c)

$$\Delta^\varepsilon g(\bar{n}+1) = \Delta^\varepsilon g(\bar{n}) - \bar{n}\varepsilon(\bar{n}) + (\bar{n}-1)\varepsilon(\bar{n}-1) + \int_{\bar{n}}^{\bar{n}+1} \varepsilon(p) dp.$$

But from (19a),  $\Delta^\varepsilon g(\bar{n}) \geq -g(\bar{n})$  so that

$$\Delta^\varepsilon g(\bar{n}+1) \geq -g(\bar{n}) - \bar{n}\varepsilon(\bar{n}) + (\bar{n}-1)\varepsilon(\bar{n}-1) + \int_{\bar{n}}^{\bar{n}+1} \varepsilon(p) dp.$$

Moreover, since  $\varepsilon(\cdot)$  is nonnegative, it follows that

$$\Delta^\varepsilon g(\bar{n}+1) \geq -g(\bar{n}) - \bar{n}\varepsilon(\bar{n}). \quad (38)$$

(37) and (38) together imply that  $\bar{n}\varepsilon(\bar{n}) > 0$ . It now follows from (15) that there exists  $\delta > 0$

such that  $\int_{\bar{n}}^{\bar{n}+\delta} \varepsilon(p) dp > 0$ . But since  $\epsilon_{\min}^{\bar{n}} = \int_1^{\bar{n}} \varepsilon(p) dp = -g(\bar{n})$ , we have:

$$\|\varepsilon\| \equiv \int_1^{N+1} \varepsilon(p) dp \geq \int_1^{\bar{n}} \varepsilon(p) dp + \int_{\bar{n}}^{\bar{n}+\delta} \varepsilon(p) dp > -g(\bar{n}),$$

completing the proof of part (2). ■

**Proof of Proposition 4:** Let  $C$  be an initial strictly convex cost function. We distinguish below between “script Q’s” (denoted  $\mathcal{Q}$ ) and regular  $Q$ ’s: The  $\mathcal{Q}$ ’s are solutions to certain equations relating  $B'$  and  $C'$ , while the  $Q$ ’s, as usual, denote aggregate abatements, i.e.,  $Q^r = rq_s^r + (N-r)q_f^r$ . Having defined our  $\mathcal{Q}$ ’s, we will construct a modified cost function  $\tilde{C}$  with the property that for each  $r$ , aggregate abatement  $Q^r$  with  $r$  signatories equals  $\mathcal{Q}^r$ .

Now fix  $\bar{n} \leq N$  and  $\delta \geq 0$ . Let  $\mathcal{Q}^{\bar{n}-1}(\delta) = Nq_s^{\bar{n}-1}(\delta) - (N - (\bar{n}-1))\delta$ , and let  $\mathcal{Q}^{\bar{n}}(\delta) = Nq_s^{\bar{n}}(\delta) - 2(N - \bar{n})\delta$ , where

$$q_s^{\bar{n}-1}(\delta) \text{ is defined by the condition } C'(q_s^{\bar{n}-1}(\delta)) = (\bar{n}-1)B'(\mathcal{Q}^{\bar{n}-1}(\delta)) \quad (39a)$$

$$q_s^{\bar{n}}(\delta) \text{ is defined by the condition } C'(q_s^{\bar{n}}(\delta)) = \bar{n}B'(\mathcal{Q}^{\bar{n}}(\delta)). \quad (39b)$$

Thus

$$\mathcal{Q}^{\bar{n}}(\delta) - \mathcal{Q}^{\bar{n}-1}(\delta) = N(q_s^{\bar{n}}(\delta) - q_s^{\bar{n}-1}(\delta)) - \delta(N - (\bar{n}-1)). \quad (40)$$

Note that for some  $\epsilon > 0$ ,  $q_s^{\bar{n}}(0) - q_s^{\bar{n}-1}(0) \geq 2\epsilon$ , since otherwise,  $\mathcal{Q}^{\bar{n}}(0) - \mathcal{Q}^{\bar{n}-1}(0) \leq 0$ , and, since  $B$  is strictly concave,  $\frac{\bar{n}B'(\mathcal{Q}^{\bar{n}}(0))}{(\bar{n}-1)B'(\mathcal{Q}^{\bar{n}-1}(0))} \geq \frac{\bar{n}}{\bar{n}-1}$  while, since  $C$  is strictly convex,  $\frac{C'(q_s^{\bar{n}}(0))}{C'(q_s^{\bar{n}-1}(0))} < 1$ , in which case equalities (39a) and (39b) could not be satisfied simultaneously. Hence by continuity, there exists  $\bar{\delta}$  s.t. if  $\delta < \bar{\delta}$ , then  $q_s^{\bar{n}}(\delta) - q_s^{\bar{n}-1}(\delta) \geq \epsilon$  and  $\mathcal{Q}^{\bar{n}}(\delta) - \mathcal{Q}^{\bar{n}-1}(\delta) > 0$ .

Fig 6 illustrates a modified cost function  $\tilde{C}(\cdot|\delta) \leq C(\cdot)$ , for  $0 < \delta < \bar{\delta}$ . We first identify two intervals on which  $C$  and  $\tilde{C}(\cdot|\delta)$  agree, and define  $\tilde{C}(\cdot|\delta)$  at some critical points in between.

$$\tilde{C}'(q|\delta) = \begin{cases} C'(q) & \text{if } q \leq (C')^{-1}(B'(\mathcal{Q}^{\bar{n}}(\delta)) - \delta) \\ B'(\mathcal{Q}^{\bar{n}}(\delta)) & \text{if } q = q_s^{\bar{n}-1}(\delta) - 2\delta \\ B'(\mathcal{Q}^{\bar{n}-1}(\delta)) & \text{if } q = q_s^{\bar{n}-1}(\delta) - \delta \\ C'(q) & \text{if } q \geq q_s^{\bar{n}-1}(\delta) \end{cases}. \quad (41)$$

For the remainder of the proof we will suppress the argument  $\delta$  when referring to the functions that depend on it. It is straightforward to verify that  $\tilde{C}'(\cdot)$  strictly increases with  $q$  over the restricted domain for which it is has so far been defined. For the remaining intervals, define  $\tilde{C}'(\cdot)$  so that the entire function is strictly increasing and differentiable. Observe also that since  $C'$  is strictly increasing, if  $\delta > 0$  is sufficiently close to zero, then  $\tilde{C}'(\cdot)$  will lie strictly below  $C(\cdot)$  on the interval  $((C')^{-1}(B'(\mathcal{Q}^{\bar{n}}) - \delta), q_s^{\bar{n}-1})$ .

By construction, under the cost function  $\tilde{C}$ , a signatory to an IEA with  $\bar{n}-1$  members chooses  $q_s^{\bar{n}-1}$  while a non-signatory chooses  $q_f^{\bar{n}-1} = q_s^{\bar{n}-1} - \delta$ . To verify this, note that

$$Q^{\bar{n}-1} \equiv (\bar{n}-1)q_s^{\bar{n}-1} + (N - (\bar{n}-1))q_f^{\bar{n}-1} = Nq_s^{\bar{n}-1} - (N - (\bar{n}-1))\delta \equiv \mathcal{Q}^{\bar{n}-1} \quad (42)$$

so that

$$\begin{aligned} \tilde{C}'(q_f^{\bar{n}-1}) &\equiv \tilde{C}'(q_s^{\bar{n}-1} - \delta) \underset{\text{from (41)}}{=} B'(\mathcal{Q}^{\bar{n}-1}) \\ \tilde{C}'(q_s^{\bar{n}-1}) &\underset{\text{from (41)}}{=} C'(q_s^{\bar{n}-1}) \underset{\text{from (39a)}}{=} (\bar{n}-1)B'(\mathcal{Q}^{\bar{n}-1}) \underset{\text{from (42)}}{=} (\bar{n}-1)B'(\mathcal{Q}^{\bar{n}-1}). \end{aligned}$$

Similarly, a signatory to an IEA with  $\bar{n}$  members chooses  $q_s^{\bar{n}}$  while a non-signatory chooses  $q_f^{\bar{n}} = q_s^{\bar{n}-1} - 2\delta$ , since

$$Q^{\bar{n}} = \bar{n}q_s^{\bar{n}} + (N - \bar{n})q_f^{\bar{n}} = Nq_s^{\bar{n}} - 2(N - \bar{n})\delta = \mathcal{Q}^{\bar{n}}, \quad (43)$$

so that

$$\begin{aligned} \tilde{C}'(q_f^{\bar{n}}) &\equiv \tilde{C}'(q_s^{\bar{n}-1} - 2\delta) \underbrace{=}_{\text{from (41)}} B'(Q^{\bar{n}}) \\ \tilde{C}'(q_s^{\bar{n}}) &\underbrace{=}_{\text{from (41)}} C'(q_s^{\bar{n}}) \underbrace{=}_{\text{from (39b)}} (\bar{n})B'(Q^{\bar{n}}) \underbrace{=}_{\text{from (42)}} \bar{n}B'(Q^{\bar{n}}). \end{aligned}$$

We now show that  $g(\bar{n}) > 0$ . In the chain of relationships below, we use approximation signs ( $\approx$ ) to indicate that the difference between the right and left hand sides of the approximation sign is  $O(\delta)$ .<sup>12</sup> Because  $\delta$  can be chosen to be arbitrarily small, the approximations can be treated for our purposes as effectively equalities. From (2),

$$\begin{aligned} g(\bar{n}) &= B(Q^{\bar{n}}) - \tilde{C}(q_s^{\bar{n}}) - \left[ B(Q^{\bar{n}-1}) - \tilde{C}(q_f^{\bar{n}-1}) \right] \\ &= B(Q^{\bar{n}}) - \left( \tilde{C}(q_s^{\bar{n}-1}) + \int_{q_s^{\bar{n}-1}}^{q_s^{\bar{n}}} \tilde{C}'(q) dq \right) - \left[ B(Q^{\bar{n}-1}) - \tilde{C}(q_f^{\bar{n}-1}) \right]. \end{aligned}$$

Because  $\tilde{C}'$  is strictly increasing

$$g(\bar{n}) > B(Q^{\bar{n}}) - \left( \tilde{C}(q_s^{\bar{n}-1}) + \tilde{C}'(q_s^{\bar{n}})(q_s^{\bar{n}} - q_s^{\bar{n}-1}) \right) - \left[ B(Q^{\bar{n}-1}) - \tilde{C}(q_f^{\bar{n}-1}) \right].$$

Since  $\tilde{C}(q_f^{\bar{n}-1}) \approx \tilde{C}(q_s^{\bar{n}-1})$

$$\begin{aligned} g(\bar{n}) &\approx B(Q^{\bar{n}}) - \tilde{C}'(q_s^{\bar{n}})(q_s^{\bar{n}} - q_s^{\bar{n}-1}) - B(Q^{\bar{n}-1}) \\ &= \int_{Q^{\bar{n}-1}}^{Q^{\bar{n}}} B'(q) dq - \tilde{C}'(q_s^{\bar{n}})(q_s^{\bar{n}} - q_s^{\bar{n}-1}). \end{aligned}$$

Since  $B' < 0$  and, from (40),  $Q^{\bar{n}} - Q^{\bar{n}-1} = N(q_s^{\bar{n}} - q_s^{\bar{n}-1}) - \delta(N - (\bar{n} - 1)) > 0$ ,

$$\begin{aligned} g(\bar{n}) &> B'(Q^{\bar{n}}) \left( Nq_s^{\bar{n}} - Nq_s^{\bar{n}-1} - \delta(N - (\bar{n} - 1)) \right) - \tilde{C}'(q_s^{\bar{n}})(q_s^{\bar{n}} - q_s^{\bar{n}-1}) \\ &\approx B'(Q^{\bar{n}}) \left( Nq_s^{\bar{n}} - Nq_s^{\bar{n}-1} \right) - \tilde{C}'(q_s^{\bar{n}})(q_s^{\bar{n}} - q_s^{\bar{n}-1}) \\ &= \left( NB'(Q^{\bar{n}}) - \tilde{C}'(q_s^{\bar{n}}) \right) (q_s^{\bar{n}} - q_s^{\bar{n}-1}) \\ &\geq \underbrace{\left( \bar{n}B'(Q^{\bar{n}}) - \tilde{C}'(q_s^{\bar{n}}) \right)}_{\text{from (39b)}} (q_s^{\bar{n}} - q_s^{\bar{n}-1}) \underbrace{=} 0. \end{aligned}$$

■