

# Curbing fossil fuels through global reward payment funds inducing countries to reduce supply, reduce demand and expand substitutes

Lennart Stern \*

August 22, 2022

The latest version can be found [here](#).

## Abstract

Existing international environmental institutions curb fossil fuels by rewarding countries for reducing demand and expanding substitutes. This paper argues that it would be beneficial to create new institutions that would reward countries for reducing their fossil fuel supply. Assuming complete information, I prove a Price Preservation Lemma: For any given budget, the optimal way to split the budget between the three approaches (rewarding supply reduction, demand reduction and substitute expansion) leaves the world market price of the fossil fuel unchanged. In a dynamic setting, this result holds under full commitment for any given intertemporal budget at the global institution's disposal. Using this Lemma, I show that the optimum cannot be implemented by relying on the supply side entirely on a deposit purchase fund. The results suggest that it would be valuable to also create funds rewarding countries for taxing fossil fuel extraction.

## 1 Introduction

This project studies goods with global externalities. In applying the model I focus throughout the paper on the case of fossil fuels and particularly on coal and oil. However, the results of the static model that I analyze in section 3 are arguably relevant for most goods with global externalities.

Specifically, I study the problem faced by a global institution having an exogenous budget which it can split between the following three approaches to curbing coal: It can pay countries to reduce coal extraction (supply reduction),

---

\*Paris School of Economics and EHESS, e-mail: [lennart.stern@psemail.eu](mailto:lennart.stern@psemail.eu)  
For helpful comments and suggestions I thank Achim Hagen, Pierre Fleckinger, Antony Millner and the participants at the 2021 IAEE conference and the 2021 IIPF conference.

to reduce energy use (demand reduction) and to expand renewables (substitute expansion).

The question arises as to how to split the budget between these three approaches. An analogous question arises for any good with global externalities. The following diagram summarizes how the world is currently answering this question for several important goods with negative global externalities (left column in brown) and goods with positive global externalities (left column in green):

good with global externalities	substitute	international institutions focused on supply	international institutions focused on demand	international institutions focused on substitute
fossil fuels	renewables, nuclear		Clean Development Mechanism, Green Climate Fund, Global Environment Facility, Climate Investment Fund, UN Environment Program	Clean Development Mechanism, Green Climate Fund, Global Environment Facility, Climate Investment Fund, UN Environment Program, IAEA's Nuclear Fuel Bank
goods produced on previously forested land (soy, palm oil...)	same goods produced in Savannah	REDD+ UN-REDD		
fish	bivalve aquaculture	WTO negotiations on reducing fishery subsidies		
drugs and vaccines for infectious disease control	single product treatments increasing risk of drug resistance		Global Alliance for Vaccines and Immunization Global Fund for AIDS, TB and Malaria	
goods for pandemic surveillance			International Health Regulations	

To analyze this question, this paper starts by studying a static model with complete information. A global institution announces reward payment schemes for each country conditioning a positive transfer on the country's coal extraction, energy use and renewable energy production. Each country takes these reward payment schemes and world market prices as given. Assuming that all demand and supply elasticities are finite, I prove that the optimal amount of funding allocated to each of the three approaches is always strictly positive and an increasing function of the total available budget (see corollary 1).

For the case of coal, I find based on middle of the road elasticity estimates taken from the literature, that for an exogenous budget it is optimal for the

global institution to spend 43% on paying countries to reduce coal supply. This contrasts with the current way that global institutions try to curb fossil fuels. So far, all of the money has been spent on demand reduction and substitute expansion. I find that for a given (not very large) available amount of money, 22% more welfare gains can be achieved if the money is split optimally between the three approaches than if the world deprives itself of the supply side approach. This provides a case for establishing a new global fund rewarding countries for reducing fossil fuel supply.

In practice, such a new global fund could define for each country reference levels of stocks of cumulative coal extraction based on business as usual scenario projections and then reward countries each year to the extent that their actual cumulative coal extraction is below the reference level for that year. This kind of scheme is already being used for rewarding countries for preserving tropical forests (Seymour and Busch (2016)).<sup>1</sup>

I show that under full commitment and assuming that such a global fund can freely save and borrow, such a global fund does not lose anything by restricting itself to this form of reward payment schemes rather than conditioning its reward payments in a more general way on the countries' extraction paths. However, a widely recognized problem with this type of scheme is the difficulty of fixing the right paths of reference levels (Mertz et al. (2018)).

One approach to avoiding this difficulty is to create instead a global fund that buys up appropriately chosen fossil fuel deposits. However, I show that the optimal mechanism cannot be implemented by relying entirely on such a deposit purchase fund on the supply side. In fact, such an approach would amount to only rewarding countries on the supply side on the basis of their eventual cumulative coal extraction. However, I prove that the optimal mechanism always involves rewarding in all period countries for having low cumulative coal extraction (see Corollary 4).

A potentially promising approach to remedy this shortcoming of deposit purchase funds could be to complement it with carbon pricing reward funds on the supply side. Countries could in each period be rewarded on the basis of the carbon taxes that they levy on fossil fuel extraction on their territories. Such an approach would also obviate the need for defining reference levels for cumulative coal extraction and with it the above-mentioned difficulties.

---

<sup>1</sup>The above diagram shows other goods with global externalities where supply-side approaches are currently absent, for example in the case of drugs for infectious disease control as displayed in the diagram. However, in that case, marginal costs of production are presumably approximately constant in the long run and so the price elasticity of supply is arguably very large in the long run. It turns out that this implies that the optimal amount of spending on rewarding supply reduction is very small. Thus the model can rationalize the fact that the world is focusing on demand side contracting in this case. Some of the other empty boxes in the above diagram appear to not be rationalizable within the model I will present. But I leave it for future work to analyze the specifics of these cases.

## 2 Related literature and contribution

Harstad (2012) studies a model where the countries adversely affected by climate change act in a coordinated way. He finds that the coalition's best policy is to simply buy foreign deposits and conserve them.

The model presented in the current paper differs in that in it climate change mitigation happens only due to the countries' responses to the global institution's reward payment schemes. I find that for exogenous funding it is optimal for the global institution to use strictly positive amounts of money on contracts rewarding supply reduction, demand reduction and substitute expansion. The Coasian approach of simply buying up fossil fuel deposits is never optimal in the model.

The current study tries to complement the literature on carbon leakage by drawing out its implications for the design of global institutions. Fæhn et al. (2017) analyze the problem of a country that tries to cause a given reduction in global emissions at a minimal cost for itself. For the case of Norway they find that two thirds of the emissions reductions should optimally come from supply reduction. This result bears some similarity to my result that under exogenous funding a global institution should optimally use 43% of its budget on spending on rewarding supply reduction. Collier and Venables (2015) provide further considerations in favor of focusing on supply side approaches.

I also contribute to the literature on the optimal roles of deposit purchase contracts and leasing contracts as instruments of supply side climate policy. I find that restricting supply side approaches to deposit purchase comes at a welfare cost, which echoes the results from Eichner et al. (2020), despite the difference in the settings. In the setting of the current paper, this result holds even though I assume climate change damages to only depend of eventual cumulative emissions.

A major limitation of the current study is that it only models a single fossil fuel (interpreted to be coal) and a clean substitute. Daubanes et al. (2020) highlight the importance of taking into account the substitution between coal and gas. Extending the model presented here to simultaneously include coal, oil and gas is left for future research.

## 3 The model

The set of countries is denoted  $I$ . Each country is assumed to be of negligible size so that it acts as a price taker on the world market.  $z_i$  is the amount of energy from renewables that country  $i$  produces.  $x_i$  denotes the amount of coal that country  $i$  extracts. Coal is measured so that one unit of coal generates one unit energy via combustion. The energy generated from coal is assumed to be a perfect substitute to the energy generated from renewables. We denote by  $y_i$  the amount of energy that country  $i$  uses. All other energy sources are assumed away.

There is a common numeraire good. Its price is normalized to 1. There are

no trade costs and there is a global market for coal. The world market price of coal is denoted  $p$ . Each country  $i \in I$  takes world market prices as given. Each country  $i$  has energy  $x_i$  from coal and energy  $z_i$  from renewables. If the sum  $x_i + z_i$  exceeds its energy use  $y_i$  then the country exports the excess amount of energy,  $x_i - y_i + z_i$ , in the form of coal. If the sum  $x_i + z_i$  is less than its energy use  $y_i$  then the country imports the shortfall of energy,  $x_i - y_i + z_i$ , in the form of coal. In either case, the net revenue that the country gets is  $p(x_i - y_i + z_i)$ .<sup>2</sup>

Country  $i$ 's utility is quasilinear in the numeraire:

$$U_i(x_i, y_i, z_i, p) = B_i(y_i) - C_i(x_i) - G_i(z_i) + p(x_i - y_i + z_i) + f_i(x_i, y_i, z_i)$$

Here  $B_i(y_i)$  is the benefit that country  $i$  derives from energy use.<sup>3</sup>  $C_i(x_i)$  is country  $i$ 's cost of extracting  $x_i$  of coal.  $G_i(z_i)$  is country  $i$ 's cost of producing  $z_i$  of energy from renewables. Moreover,  $f_i(x_i, y_i, z_i)$  denotes the transfer that country  $i$  gets from the global institution, as explained further below<sup>4</sup>.

It will be convenient to impose that costs are strictly convex and benefits strictly concave<sup>5</sup>:

**Assumption 1 (strict convexity of costs and strict concavity of benefits).**  $C_i'(x_i) > 0, C_i''(x_i) > 0 \forall i \forall x_i, G_i'(z_i) > 0, G_i''(z_i) > 0 \forall i \forall z_i, B_i'(y_i) > 0, B_i''(y_i) < 0 \forall i \forall y_i$ .

There is a global institution which evaluates global welfare as follows:

$$W = \sum_{i \in I} U_i - \eta \left( \sum_{j \in I} x_j \right)$$

The interpretation is as follows:  $\eta$  is a positive and strictly increasing function.  $\eta(\sum_{j \in I} x_j)$  is the aggregate value of the global climate change damages due to the aggregate amount  $\sum_{j \in I} x_j$  of coal combusted.  $U_i$  is the utility that country  $i$  tries to optimize. Strictly, this should include the damage due to climate change that country  $i$  suffers due to its own coal use. However, I make the simplifying assumption that country  $i$  neglects this, in line with our assumption that all countries are of negligible size.

The global institution is endowed with an exogenous budget  $F$ . It offers reward payments to countries to induce them to reduce their coal supply, their

<sup>2</sup>Currently, 20% of all coal is traded internationally. Consistent with our assumption of a globally integrated coal market, Steckel et al. (2015) find that "in the increasingly integrated global coal market the availability of a domestic coal resource does not have a statistically significant impact on the use of coal and related emissions".

<sup>3</sup>This should be interpreted to be the entire surplus that the country reaps from energy use, both in the form of consumer surplus accruing to end users and producer surplus from production using energy as an input.

<sup>4</sup>It turns out that the global institution does not lose anything by restricting itself to using additively separable reward payment functions, i.e.  $f_i(x_i, y_i, z_i) = f_{x_i}(x_i) + f_{y_i}(y_i) + f_{z_i}(z_i)$ . I will therefore restrict attention to such additively separable reward payment functions later on. For now I keep the notation more general, to avoid the impression that the additively separable form is a restrictive assumption.

<sup>5</sup>This assumption excludes the case of constant marginal cost, which is of interest for applications like the global health examples mentioned in the introduction. However, instead of treating this case separately, we will discuss this as a limiting case as price elasticities of supply go to infinity.

coal demand and to expand their renewables supply. The timing is as follows: First, the global institution announces transfers that it will pay to countries conditional on their choices.  $f_i(x_i, y_i, z_i)$  denotes the transfer that country  $i$  will receive if it chooses  $(x_i, y_i, z_i)$ . Since countries are sovereign, the global institution cannot ask countries to pay it money, which means that the transfers  $f_i(x_i, y_i, z_i)$  are constrained to be non-negative. Each country  $i$  takes the reward payment scheme  $f_i(x_i, y_i, z_i)$  and the price  $p$  as given and chooses  $(x_i, y_i, z_i)$  so as to maximize its utility  $U_i(x_i, y_i, z_i, p)$ .

**Definition 1.** A **reward payment scheme** offered to country  $i$  is a map  $f_i(x_i, y_i, z_i)$  assigning a nonnegative transfer to country  $i$ . A **world market equilibrium under a given set of reward payment schemes**  $(f_i)_{i \in I}$  is a combination of an allocation  $(x_i, y_i, z_i)_{i \in I}$  and world market price  $p$  such that the following 2 conditions hold:

- 1) market clearing:  $\sum_{i \in I} x_i - y_i + z_i = 0$
- 2) individual rationality:  $(x_i, y_i, z_i) = \operatorname{argmax}_{(x, y, z)} -C_i(x) + B_i(y) - G_i(z) + p(x - y + z) + f_i(x, y, z) \forall i \in I$

**Definition 2.** A set  $(f_i)_{i \in I}$  of reward payment schemes **implements** the allocation-price pair  $((x_i, y_i, z_i)_{i \in I}, p)$  with a budget  $F$  if  $((x_i, y_i, z_i)_{i \in I}, p)$  is a world market equilibrium under  $(f_i)_{i \in I}$  and  $\sum_{i \in I} f_i(x_i, y_i, z_i) = F$ .

**Definition 3.** A reward payment scheme  $f_i(x_i, y_i, z_i)$  is called **additively separable** if it can be written as  $f_i(x_i, y_i, z_i) = f_{ix}(x_i) + f_{iy}(y_i) + f_{iz}(z_i)$ . A reward payment scheme is called a “**positive affine linear scheme**” if it can be written as  $f_i(x_i, y_i, z_i) = \max(0, \theta_{ix}(\tilde{x}_i - x_i)) + \max(0, \theta_{iy}(\tilde{y}_i - y_i)) + \max(0, \theta_{iz}(z_i - \tilde{z}_i))$ .

We will now prove that (under our assumption of convex cost functions and concave benefit functions) nothing is lost by restricting attention to the positive affine linear schemes for the reward payment schemes. We will also show that we can view the global institution as if it was choosing the allocation and the world market prices.

**Lemma 1 (The Surjectivity Lemma).** *Consider a combination of an allocation  $(x_i, y_i, z_i)_{i \in I}$  satisfying  $\sum_{i \in I} x_i - y_i + z_i = 0$  and world market price  $p$ . There exists a set  $(f_i)_{i \in I}$  of positive affine linear schemes implementing  $((x_i, y_i, z_i)_{i \in I}, p)$ . Moreover, the minimal transfers required to implement  $((x_i, y_i, z_i)_{i \in I}, p)$  under affine linear schemes are  $F_{ix}$  for rewarding country  $i$  for supply reduction,  $F_{iy}$  for rewarding country  $i$  for demand reduction and  $F_{iz}$  for rewarding country  $i$  for substitute expansion with*

$$\begin{aligned} F_{ix} &:= \sup_x px - C_i(x) - (px_i - C_i(x_i)) \\ F_{iy} &:= \sup_y B_i(y) - py - (B_i(y) - py) \\ F_{iz} &:= \sup_z pz - G_i(z) - (pz - G_i(y)) \end{aligned}$$

Moreover, there does not exist any set of reward payment schemes implementing  $((x_i, y_i, z_i)_{i \in I}, p)$  with a strictly smaller budget, i.e. with a budget strictly less than  $\sum_{i \in I} F_{ix} + F_{iy} + F_{iz}$ . Furthermore, these minimal required transfers are the same if we allow for any reward payment schemes (instead of restricting them to be positive affine linear).

*Proof.* See Appendix A.1 □

Let us now summarize the Surjectivity Lemma using the diagram below. In practice, the global institution chooses reward payment schemes. The notion of market equilibrium defined above yields a (potentially multivalued) mapping assigning a (or several) combination(s) of a world market price and an allocation to each set of reward payment schemes. By the surjectivity Lemma 1, this map is surjective.

The surjectivity of this market equilibrium map allows us to view the global institution as if it was choosing a combination of a world market price and an allocation. Given a world market price and an allocation, there are many reward payment schemes inducing them via the market equilibrium map. The transfers that end up being paid are  $f_{xi}(x_i), f_{yi}(y_i), f_{zi}(z_i)$ . The minimal required transfers  $F_{ix}, F_{iy}, F_{iz}$  are given by the formulae shown below:

$$\begin{array}{ccc}
 & \text{market equilibrium} & \\
 & \text{(surjective map)} & \\
 (f_{xi}, f_{yi}, f_{zi})_{i \in I} & \longrightarrow & (p, (x_i, y_i, z_i)_{i \in I}) \\
 & & \downarrow \text{minimal required} \\
 & & \text{transfers} \\
 & & (F_{ix}, F_{iy}, F_{iz}) \\
 & & F_{ix} = \sup_x px - C_i(x) - (px_i - C_i(x_i)) \\
 & & F_{iy} = \sup_y B_i(y) - py - (B_i(x_i) - px_i) \\
 & & F_{iz} = \sup_z pz - G_i(z) - (pz_i - G_i(z_i))
 \end{array}$$

In this diagram we are restricting attention to additively separable reward payment schemes. This is justified by the Surjectivity Lemma : This restriction does not affect the minimal transfers that are required.<sup>6</sup>

Hence the global institution's problem can be written as:

$\max_{(p, (x_i, y_i, z_i)_{i \in I})} \sum_{j \in I} U_j - \eta(\sum_{j \in I} x_j)$  subject to the market clearing constraint,  $\sum_{i \in I} x_i - y_i + z_i = 0$ , and the budget balance constraint that  $\sum_{j \in I} F_{jx} + F_{jy} + F_{jz} \leq F$ , where  $F$  is the exogenous budget at the global institution's disposal.

**Lemma 2 (The First Best).** *For a sufficiently large budget  $F$ , the first best is characterized by the following conditions plus the feasibility condition  $\sum_{i \in I} (x_i - y_i + z_i) = 0$ :*

<sup>6</sup>Interestingly, this no longer holds if we were to depart from the assumption of complete information. In fact, one can deduce from the results in Armstrong and Rochet (1999) that even if types are independent across dimensions the optimal mechanism will not be additively separable.

$$B'_i(y_i) = B'_j(y_j) = G'_i(z_i) = G'_j(z_j) = C'_i(x_i) + \eta = C'_j(x_j) + \eta \forall i, j \in I$$

There exists a continuum of ways to split the budget between supply reduction reward payments schemes on the one hand and demand reduction and substitute expansion reward payment schemes on the other, all of which achieve the optimal global welfare.

*Proof.* See appendix A.2. □

Lemma 2 asserts that if the global institution's budget constraint is not binding, then it does not matter how exactly the budget is split between supply reduction on the one hand and demand reduction and substitute expansion on the other, as long as the resulting required budget does not exceed the available budget. However, from now on we will assume that the global institution's budget constraint *is* binding:

**Assumption 2.** *The global institution's budget constraint is binding. In other words: its budget is insufficient to fully correct the global externalities from coal.*

This assumption will mean that it *will* matter for global welfare how the global institution's budget is split, thereby overturning the conclusion from Lemma 2. This is because the world market price of coal affects the sizes of the transfers required to make countries change their actions instead of just ignoring the reward payments. The world market price of coal, in turn, is affected by how countries are rewarded: The stronger the reward payments for supply reduction, the weaker the supply of coal on the world market and therefore the higher the resulting world market price. On the other hand, the stronger the reward payments for demand reduction and substitute expansion, the lower the demand for coal on the world market and thus the lower the resulting world market price for coal.

The following Lemma states that it is always optimal to choose a mixture of these two kinds of approaches, balanced precisely such that the net effect of the world market price of coal is neutral. It is important to emphasize that this "Price Preservation Lemma" refers to the *world market price* and not to the *net prices* that actors will face. Within a given country, the price that actors will face is the sum of the world market price and any taxes (or regulation-induced carbon prices, etc.) that the government will set. When the global institution rewards countries for supply reduction then this effectively means that it will pay countries for setting a carbon price on the coal extracted on its territory. When the global institution rewards countries for demand reduction then this effectively means that it will pay countries for taxing energy use (by households and firms). When the global institution rewards countries for expanding renewables it pays countries for exempting renewables from the tax on energy use or even for subsidizing renewables.

The result of all this will always be that the *net price* of coal combustion (including taxes and other implicit or explicit carbon prices) will increase. What will optimally be preserved is the *world market price* of coal:



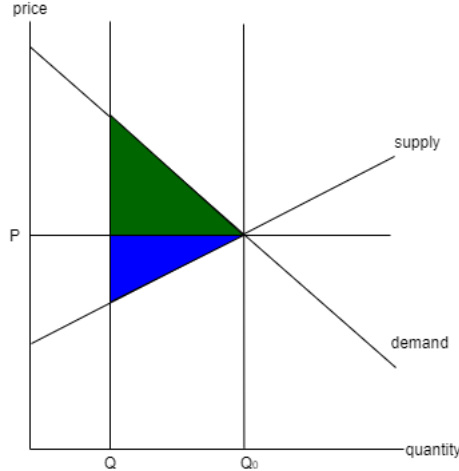
**Lemma 3 (The Price Preservation Lemma).** “If a set of reward payment schemes achieves a given allocation with minimal aggregate transfer payments then it must preserve the world market price  $p$  of coal.”

Formally: Consider a fixed allocation  $(x_i, y_i, z_i)_{i \in I}$ . Consider a set of reward payment schemes  $(f_i(x_i, y_i, z_i))_{i \in I}$  implementing  $((x_i, y_i, z_i)_{i \in I}, p)$  for some price  $p$  under a budget  $F$ . Then if there is no other set of reward payment schemes  $(\tilde{f}_i(x_i, y_i, z_i))_{i \in I}$  implementing  $((x_i, y_i, z_i)_{i \in I}, \tilde{p})$  for some price  $\tilde{p}$  under a budget  $\tilde{F}$  with  $\tilde{F} < F$  then  $p$  must equal the price of energy in the absence of any reward payment schemes.

In particular, the optimal reward payment schemes must leave the world market price  $p$  of coal at the same level as when there are no reward payment schemes.<sup>7</sup>

*Proof.* See Appendix A.3 □

For the case where there is no substitute, we can illustrate the Price Preservation Lemma graphically. Suppose for simplicity that all countries have identical demand and supply functions shown in the following diagram:

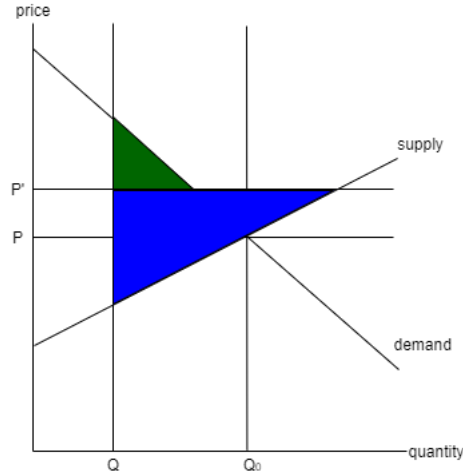


In the absence of any reward payment schemes, the world market price for coal is  $p$  and the quantity of coal produced and used by each country is  $Q_0$ . Now let us compare different ways of reducing the quantity produced (and used) to  $Q < Q_0$ . Suppose first the global institution achieves this using a reward payment scheme that leaves the price unchanged. In that case, the minimal

<sup>7</sup>It is straightforward to generalize both the Surjectivity Lemma and the Price Preservation Lemma to the case where there are intermediate inputs used only for the good in question. This is relevant in other applications. For example, consider the problem of how to best cause the production of vaccines in normal times to increase so that the world is better prepared for the next pandemic. Consider all the intermediate inputs to vaccines that are only used for them. The Price Preservation Lemma implies that it is optimal for the global institution to use a part of its budget for paying countries to expand production of these intermediate inputs. This is because if it does not then the prices for these intermediate inputs would increase as a result of the increased demand, in contradiction to the Price Preservation Lemma.

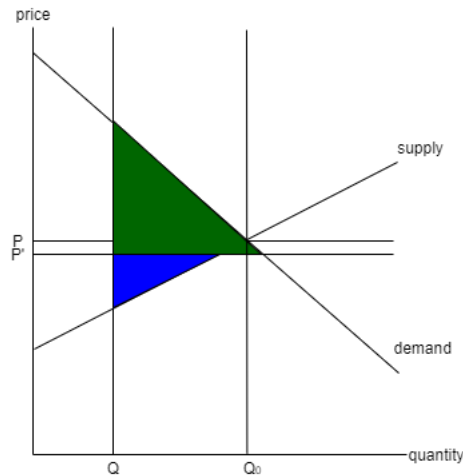
transfer that it has to pay each country for reducing their coal use from  $Q_0$  to  $Q$  is given by the green area. The minimal transfer that it has to pay each country for reducing their coal extraction is given by the blue area.

Now suppose the global institution were to implement  $Q$  with a higher world market price  $P' > P$ . This corresponds to a higher spending on rewarding supply reduction:



We see that the total size of the green and the blue areas together is larger now than when the price was preserved at  $P$ .

Similarly, greater demand side emphasis, corresponding to a smaller price  $P' < P$ , would require larger overall transfers:



Thus we have graphically recovered the Price Preservation Lemma: If a reward payment scheme is to achieve a given allocation with minimal aggregate transfers then the world market price  $p$  of coal must be the same as in the absence of any reward payment scheme.

Having established the Price Preservation Lemma by holding the allocation

constant and finding the price that minimizes the required transfers, let us now hold the price constant and find the optimal allocation for the given price.

**Lemma 4 (The Constrained Efficiency Lemma).** *At the optimal reward payment scheme we have: The allocation  $(x_i, y_i, z_i)$  achieves maximal welfare among all allocations satisfying market clearing and having the same value of  $\sum_{i \in I} x_i$ . Moreover, this constrained efficiency result even holds if we add the constraint that the world market price  $p$  of coal be any fixed value.*

*Proof.* See Appendix A.4 □

Given the Price Preservation Lemma and the Constrained Efficiency Lemma, it is intuitively clear that as the available budget  $F$  increases, so will the amounts spent on each of the three approaches at the optimal reward payment scheme. To see why, we note that if we were to only expand the budget for rewarding supply reduction, then the world market price  $p$  of coal would increase, in contradiction to the Price Preservation Lemma. Similarly, if we were to only expand the budget for rewarding demand reduction and the budget for rewarding substitute expansion, then the world market price of coal would fall. Moreover, from the Constrained Efficiency Lemma it is intuitively clear that the demand side budget and the substitute side budget must both expand: restricting marginal abatement to demand reduction or substitute expansion would come at an efficiency loss. I will formally validate this conclusion by proving the following:

**Corollary 1 (The Interior Solution Corollary).** *For the optimal reward payment scheme given the budget  $F$ , let  $F_x(F)$  denote the amount optimally used for supply side payments and similarly  $F_y(F)$  the optimal demand side budget and  $F_z(F)$  the optimal substitute side budget. We have:  $\frac{dF_x}{dF} > 0$ ,  $\frac{dF_y}{dF} > 0$ ,  $\frac{dF_z}{dF} > 0 \forall F$ .*

*Proof.* See Appendix A.5 □

**Corollary 2 (The Optimal Budget Split Corollary for Small Budgets).** *For the optimal reward payment scheme given the budget  $F$ , let  $F_x(F)$  denote the amount used for supply side reward payments and similarly  $F_y(F)$  the demand side budget and  $F_z(F)$  the substitute side budget. Moreover, let us denote by  $X(F)$  the resulting aggregate coal extraction and by  $\epsilon_x$  the aggregate price elasticity of supply of coal,  $Y(F)$  the resulting aggregate energy consumption and and by  $\epsilon_y$  the aggregate price elasticity of energy demand, by  $Z(F)$  the resulting aggregate renewable energy production and by  $\epsilon_z$  the aggregate price elasticity of supply of renewable energy. We “generically” have:*

$$\lim_{F \rightarrow 0} \frac{dF_x}{dF} = \frac{\epsilon_z \frac{Z(0)}{Y(0)} + \epsilon_y}{\epsilon_z \frac{Z(0)}{Y(0)} + \epsilon_x \frac{X(0)}{Y(0)} + \epsilon_y}$$

$$\lim_{F \rightarrow 0} \frac{dF_y}{dF} = \frac{\epsilon_y}{\epsilon_z + \epsilon_y} \frac{\epsilon_x \frac{X(0)}{Y(0)}}{\epsilon_z \frac{Z(0)}{Y(0)} + \epsilon_x \frac{X(0)}{Y(0)} + \epsilon_y}$$

$$\lim_{F \rightarrow 0} \frac{dF_z}{dF} = \frac{\epsilon_z}{\epsilon_z + \epsilon_y} \frac{\epsilon_x \frac{X(0)}{Y(0)}}{\epsilon_z \frac{Z(0)}{Y(0)} + \epsilon_x \frac{X(0)}{Y(0)} + \epsilon_y}$$

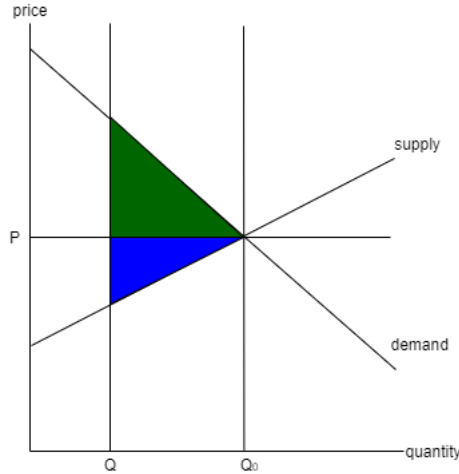
*Proof.* See Appendix A.5 □

We see that the more elastic the supply of coal, the smaller the proportion of money that will optimally be used to pay countries for reducing coal extraction. We even get the following:

**Corollary 3.**  $\lim_{\epsilon_x \rightarrow \infty} \lim_{F \rightarrow 0} \frac{dF_x}{dF} = 0$

*Proof.* This follows directly from Corollary 2 □

For the case without any substitute we can again illustrate this result diagrammatically:



By the Price Preservation Lemma we know that the optimal reward payment scheme consists of spending the amount corresponding to the blue area on rewarding supply reduction. Making supply more and more elastic corresponds to making the supply curve flatter and flatter. This decreases the blue area, so the optimal proportion of spending on rewarding supply reduction decreases. Thus in the limiting case where there are constant returns to scale on the supply side, it is optimal to focus all reward payments on the demand (and substitute) side.

Intuitively, this can be explained as follows: If there are constant marginal costs on the supply side then each ever so tiny country could reap large profits if the world market price was increased as result of reward payments inducing the other countries to reduce their supply. Thus the global institution would need to pay them large transfers to prevent them from seizing this opportunity to increase profits.

We conclude this section with some concavity results:

**Lemma 5.** *Let  $W(F)$  be the maximal welfare achievable with a given budget  $F$ . Suppose that  $\eta(\sum_i x_i)$  is linear. Then  $W(F)$  is strictly concave.*

*Proof.* See Appendix A.6 □

**Lemma 6.** *Let  $F(F_x, W)$  denote the budget required to achieve welfare  $W$  under the further constraint that the budget spent on reducing coal supply is  $F_x$ . Then  $F(F_x, W)$  is convex in  $F_x$ .*

*Proof.* See Appendix A.7. □

Lemmas 5 and 6 make the following conjecture plausible:

**Conjecture 1.** *Let  $W(F_x, F_y, F_z)$  be the maximal welfare achievable under the further constraint that the amount  $F_x$  be spent on coal supply reduction,  $F_y$  be spent on energy demand reduction and  $F_z$  be spent on renewable supply expansion.  $W(F_x, F_y, F_z)$  is concave.*

In the special case where all supply and demand functions have constant elasticities, conjecture 1 does seem to hold, as suggested by the numerical results shown in the next section.

## 4 Numerical results for the model under constant elasticity specification

All the numerical calibrations whose results are summarized here are fully documented in the accompanying Mathematica notebook that can be downloaded [here](#).

Drawing on the literature, I use the following middle-of-the-road for the parameters <sup>8</sup>:

$\epsilon_D = 0.85$  based on Freehan (2018) and Espey, J. A., & Espey, M. (2004).

$\epsilon_{S_G} = 2.7$  based on Johnson (2011)

$\epsilon_{S_C} = 1.3$  based on Dahl (2009)

$\eta = 0.4$  based on a social cost of carbon of \$36 per ton of CO<sub>2</sub> (based on EPA (2015))

$$\frac{X(0)}{Y(0)} = \frac{0.4}{0.26+0.4}$$

---

<sup>8</sup>For the numerical calibrations that follow I use a slightly more complicated (but formally isomorphic model) that is detailed in the accompanying Mathematica notebook. The model takes into account that energy is required as an input to produce renewable energy.

<sup>9</sup>In the accompanying Mathematica notebook I take into account that there are costs for generating coal-powered electricity other than the coal itself. The model is arguably most relevantly applied to the non-Annex 1 countries, given that existing global environmental institutions limit their reward payments to these countries. This is why I take India for calibrating the cost parameters for coal powered electricity:

Electricity prices are around \$0.08 per kwh in India.

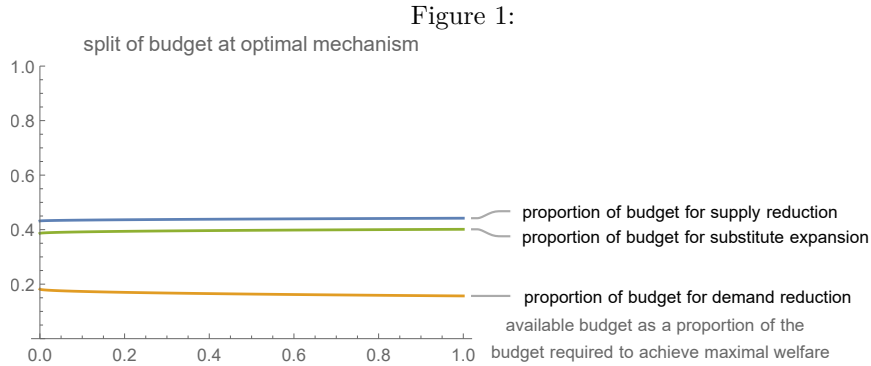
Per kwh of electricity from coal 900g of CO<sub>2</sub> gets emitted.

Assume a social cost of carbon of \$36 per ton of CO<sub>2</sub> (based on EPA (2015)).

Thus we have:  $p = 0.08$  per kwh and  $\eta = 0.9 \times 36 / 1000$  per kwh

$$\frac{Z(0)}{Y(0)} = \frac{0.26}{0.26+0.4} = 10$$

By Lemma 2, the optimal budget split for an infinitesimally small budget only depends on the elasticities at the situation where no reward payment scheme is in place. Shown in the following figure is the optimal budget split as a function of the available budget:



The plot shows the results starting from an infinitesimal budget all the way to the minimal budget allowing the global institution to fully correct the global externality. Interestingly, the result for the optimal budget split hardly depends on the budget.

In the following figure I show how global welfare depends on how the budget is split between the three approaches:

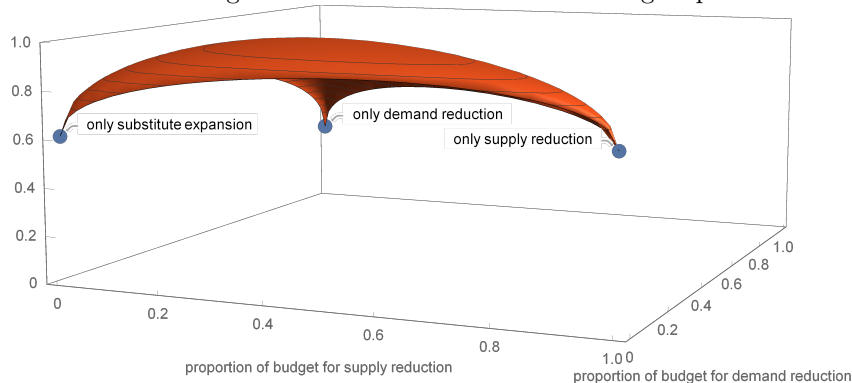
---

Hence  $\frac{\eta}{p} = 0.4$ . Thus with our normalization of  $p = 1$  we get  $\eta = 0.4$ .

The amount of coal used to generate 1 kWh is 0.00052 short tons by the US DOI. The price per short ton of bituminous coal is  $p = \$58.93$  (based on EIA). This gives a cost on coal inputs of \$0.0306436 per kWh of coal generated electricity. Assuming that the cost of generating coal-powered electricity equals its price in India (\$0.08 per kWh), this means that 38% of the cost is due to the coal itself. I use this figure in the accompanying Mathematica notebook to compute the plots that follow.

<sup>10</sup>As of 2018, renewables generates 26 % of global electricity (IEA (2019)), whilst coal generates 40 % of global electricity (World Bank). Since the model does not take into account the other electricity sources, I ignore them for the purposes of this illustrative calibration.

Figure 2: welfare as a function of budget split



On the x-axis and the y-axis are the proportions of the budget used for rewarding supply reduction and rewarding demand reduction, respectively. By definition, the remaining proportion of the budget is used for rewarding substitute expansion. On the z-axis is the ratio of the global welfare gains achieved divided by the maximal global welfare gains achievable with the given budget. The figure depicts the case where the budget is small. It turns out that the results change very little when one chooses any other budget value between 0 and half the minimal value required to fully correct the global externality.

Around the optimum the surface is quite flat. In fact, as long as each of the three approaches to curbing coal is funded at at least 50% its optimal proportion, welfare losses relative to the optimal budget split are at most 10%. In E I show that this result is quite robust across the ranges of elasticity estimates found in the literature.

This result has important implications for the design of mechanisms to fund the global institution. It suggests that it might not be so important to get the budget split exactly right and thus weighs in favor of decentralized funding mechanisms that have no guarantee for allocative efficiency but that create strong participation incentives by giving participating countries the opportunity to influence the allocation of funding across the different approaches to curbing fossil fuels (see Stern (2020) for an example of such a mechanism).

## 5 Implementation via tax-based reward payment schemes in the static model

Let us now assume that all governments determine the amount of fossil fuel extraction, fossil fuel use and renewable energy production by setting taxes and subsidies. In this case, the global institution can reward countries on the basis of these tax and subsidy rates. The definitions 1 and 2 from the previous section can be straightforwardly adapted to this setting.

**Lemma 7.** *Let  $p$  denote the world market price of coal in the absence of any reward payment schemes. Let  $x_i^*(p_x)$  the supply of coal in country  $i$  when the after-tax producer price for extracted coal is  $p_x$ ,  $y_i^*(p_y)$  the demand for energy in country  $i$  when the after-tax user price for energy is  $p_y$  and  $z_i^*(p_z)$  the supply of renewable energy in country  $i$  when the after-subsidy producer price for renewable energy is  $p_z$ .*

*The optimal mechanism can be implemented as follows: The global institution makes each country a take-it-or-leave-it offer to set a tax  $\tau_x$  on the extraction of fossil fuel against a reward payment of  $t_x := \sup_x px - C_i(x) - (px_i^*(p - \tau_x) - C_i(x_i^*(p - \tau_x)))$ , a take-it-or-leave-it offer to set a tax  $\tau_y$  on the energy use against a reward payment of  $\sup_y B_i(y) - py - (B_i(y_i^*(p + \tau_y)) - py_i^*(p + \tau_y))$ , a take-it-or-leave-it offer to set a subsidy on renewable energy production of  $\tau_y$  against a reward payment of  $pz - G_i(z) - (pz_i^*(p + \tau_y) - G_i(p + \tau_y))$ . Here  $(\tau_x, \tau_y)$  is uniquely determined by the requirement of global market clearing under  $p$ ,  $\sum_{i \in I} x_i^*(p - \tau_x) - y_i^*(p + \tau_y) + z_i^*(p + \tau_y) = 0$  and the condition that the sum of the transfers equals the global institution's budget  $F$ .*

*Proof.* The Constrained Efficiency Lemma 4 implies that the tax and subsidy rates have to be equal across countries and also that the marginal benefit of energy use equals the marginal cost of renewable energy production. By the Price Preservation Lemma 3, optimality requires implementing the world market price  $p$ . Lemma 1 implies that at this world market price the transfers are the minimal transfers making the countries accept the offers.  $\square$

The take-it-or-leave-it offers of the previous Lemma are clearly not ideal for use in practice since inaccuracies in the calculation of the required reward payments could lead to some countries rejecting some offers. It would be more robust if we could use for example affine linear reward payment schemes. In the setting of the previous section, where countries are rewarded on the basis of their quantities (i.e.  $x_i, y_i, z_i$ ), Lemma 1 shows that restricting reward payments to positive affine linear reward payment schemes does not lead to any welfare loss. This is reassuring: Under such schemes, if a country ends up having for example unexpectedly large fossil fuel extraction then it might still have some incentives to reduce it under the affine linear scheme even in situations where under the take-it-or-leave-it offers it would be best off rejecting the offer.

In addition to this, the positive affine linear reward payment schemes also have the advantage of simplicity which can explain why they are used, for example in the case of Norway's REDD contracts Angelsen (2017).

Unfortunately, positive affine linear reward payment schemes are in general not sufficient to implement the optimal allocations when countries are rewarded on the basis of their prices (i.e. their tax and subsidy rates), the focus of the current section. However, I will now provide sufficient conditions under which we can use a truncated form of the positive affine linear reward payment schemes without loss of welfare:

**Definition 4.** A reward payment scheme is called an “**upwardly truncated positive affine linear scheme**” if it can be written as  $f_i(x_i, y_i, z_i) = \min(\gamma_{ix}, \max(0, \theta_{ix}(\tilde{x}_i -$



$x_i)) + \min(\gamma_{iy}, \max(0, \theta_{ix}(\tilde{y}_i - y_i))) + \min(\gamma_{iz}, \max(0, \theta_{iz}(z_i - \tilde{z}_i))).$

**Lemma 8.** *Let  $p$  denote the world market price of coal in the absence of any mechanism. Denote  $x_i^*(p) := \operatorname{argmax}_x px - C_i(x)$ ,  $y_i^*(p) := \operatorname{argmax}_y B_i(y) - py$ ,  $z_i^*(p) := \operatorname{argmax}_z pz - G_i(z)$ . Consider a budget  $F$ . Let  $(x_i, y_i, z_i)_{i \in I}$  denote the optimal allocation implementable under  $F$ . Then each of the 3 sets of conditions listed below is sufficient to ensure that the optimal allocation can be implemented as an interior solution under upwardly truncated positive affine linear reward payment schemes (in the sense that it can be implemented under a reward payment scheme of the form  $f_i(\tau_{ix}, \tau_{iy}, \tau_{iz}) = \min(\gamma_{ix}, \max(0, \theta_{ix}(\tilde{\tau}_{ix} - \tau_{ix})) + \min(\gamma_{iy}, \max(0, \theta_{iy}(\tilde{\tau}_{iy} - \tau_{iy}))) + \min(\gamma_{iz}, \max(0, \theta_{iz}(\tau_{iz} - \tilde{\tau}_{iz})))$  with  $\tau_{ix} < \gamma_{ix}, \tau_{iy} < \gamma_{iy}, \tau_{iz} < \gamma_{iz}$ ):*

1)

$$1 - \left( \frac{B'_i(y) - p}{B'_i(y)} \right) \frac{\epsilon_{B'_i}(y)}{\epsilon_{B'_i}(y)} > 0 \forall y \in [y_i, y_i^*(p)] \quad (1)$$

where  $\epsilon_{B'_i}(y) := \frac{y}{B'_i(y)} B'_i(y)$  and  $\epsilon_{B''_i}(y) := \frac{y}{B''_i(y)} B''_i(y)$ ,

$$1 - \left( \frac{C'_i(x) - p}{C'_i(x)} \right) \frac{\epsilon_{C''_i}(x)}{\epsilon_{C'_i}(x)} > 0 \forall x \in [x_i^*(p), x_i] \quad (2)$$

where  $\epsilon_{C'_i}(x) := \frac{x}{C'_i(x)} C'_i(x)$  and  $\epsilon_{C''_i}(x) := \frac{x}{C''_i(x)} C''_i(x)$ ,

$$1 - \left( \frac{G'_i(z) - p}{G'_i(z)} \right) \frac{\epsilon_{G''_i}(z)}{\epsilon_{G'_i}(z)} > 0 \forall z \in [z_i^*(p), z_i] \quad (3)$$

where  $\epsilon_{G'_i}(z) := \frac{z}{G'_i(z)} G'_i(z)$  and  $\epsilon_{G''_i}(z) := \frac{z}{G''_i(z)} G''_i(z)$

or

2) The budget  $F$  is sufficiently small.

or

3)  $C'_i$  is convex on  $[x_i^*(p), x_i]$ ,  $B'_i$  is concave on  $[y_i, y_i^*(p)]$ , and  $G'_i$  is concave on  $[z_i^*(p), z_i]$

*Proof.* See appendix A.8 □

## 6 The dynamic model

There are  $T$  discrete time periods, denoted by  $t = 1, \dots, T$ . Country  $i$ 's coal reserves are such that the order in which it is best to extract is always the same. Thus country  $i$ 's only relevant choice is how much to extract in each period. Country  $i$ 's cumulative extraction of coal until the end of period  $t$  is denoted by  $x_{it}$ .

Country  $i$ 's energy use in period  $t$  is denoted by  $y_{it}$  and its renewable energy production by  $z_{it}$ . Country  $i$ 's pursues the objective of maximizing  $U_i$  given by:

$$U_i = \sum_{t=1}^T \frac{1}{(1+r)^t} (b_{it}(y_{it}) - g_{it}(z_{it}) - (c_{it}(x_{it}) - c_{it}(x_{it-1})) + f_{it}(x_{it}, y_{it}, z_{it}) + p_t(x_{it} + z_{it} - y_{it}))$$

Here I am assuming that each country  $i$  discounts future benefits and costs using the market interest rate  $r$ .  $b_{it}(y_{it})$  denotes country  $i$ 's benefit in period  $t$  from using the amount  $y_{it}$  of energy in that period.  $g_{it}(z_{it})$  denotes country  $i$ 's cost of producing  $z_{it}$  of renewable energy in period  $t$ .

The term  $c_{it}(x_{it}) - c_{it}(x_{it-1})$  represents the incremental cost in of extracting in period  $t$  the amount  $x_{it} - x_{it-1}$ , given that the cumulative coal extraction by the end of period  $t - 1$  is  $x_{it-1}$ . The form of the coal extraction cost assumed here is a discrete-time analogue of the assumption of stock-dependent extraction cost common in the continuous time literature (see e.g. van der Ploeg and Withagen (2014)). In fact, it is more general in that it allows the stock-dependent extraction cost to also explicitly depend on time.

The transfer that country  $i$  receives from the global institution in period  $t$  is denoted by  $f_{it}(x_{it}, y_{it}, z_{it})$ . With this notation I am implicitly constraining the global institution to not condition its reward payments in more complicated ways on the past. However, with Lemma 11 I establish that this is in fact without any loss.

As in the static model, I assume that the benefit functions are concave and the cost functions convex:

**Assumption 3.**  $b_{it}$  is concave  $\forall i, t$ . Moreover,  $g_{it}$  and  $c_{ii}$  are convex  $\forall i, t$ .

Analogously to the static model, the global institution evaluates global welfare as follows:

$$W = \sum_{i \in I} U_i - \eta \left( \sum_{i \in I} x_{iT} \right)$$

where  $\eta$  is a positive and strictly increasing function representing the global aggregate climate change damages. I am thus assuming here that climate change damages are determined by the aggregate cumulative emissions by the end of the last period. I conjecture that the results that follow hold under much weaker assumptions on how the aggregate emissions path translates into climate change damages<sup>11</sup>.

As throughout all of the formal analysis in this paper, I assume that the funding available for the global institution is exogenous. Specifically, I now assume that there is an exogenous flow of funding  $F_t$  being given to the global institution in period  $t$ . Moreover, I assume that the global institution can freely save and borrow at the interest rate  $r$ , which is also the rate at which the countries discount future money flows. Thus the global institutions has an intertemporal budget constraint: The discounted value of the aggregate transfers paid cannot exceed  $F := \sum_{t \in \{1, \dots, T\}} \frac{F_t}{(1+r)^t}$ .

<sup>11</sup>In Stern (2021) I prove other results about the dynamic model that do in fact hold under more general assumptions.

The following definition extends to the dynamic model the market equilibrium notion that I introduced already for the static model:

**Definition 5.** A *path of reward payment schemes* offered to country  $i$  is a sequence of maps  $(f_{it}((x_{is})_{s=1,\dots,t}), (z_{is})_{s=1,\dots,t}, (z_{is})_{s=1,\dots,t}))_{t \in \{1,\dots,T\}}$ , each assigning a non-negative transfer to country  $i$  depending on the prior history of the country's choices. A *path of reward payment schemes* is called **period-by-period** if it can be written as  $(f_{it}(x_{it}, y_{it}, z_{it}))_{t \in \{1,\dots,T\}}$ , meaning that the reward payment in a period  $t$  only depends on the the country's cumulative coal extraction  $x_{it}$  by the end of the period  $t$ , its energy use  $y_{it}$  in period  $t$  and its renewable energy production  $z_{it}$  in period  $t$ .

A **world market equilibrium under a given set of paths of reward payment schemes**  $((f_{it})_{t \in \{1,\dots,T\}})_{i \in I}$  is a combination of an allocation  $(x_{it}, y_{it}, z_{it})_{i \in I, t \in \{1,\dots,T\}}$  and a world market price path  $p = (p_t)_{t \in \{1,\dots,T\}}$  such that:

- 1) market clearing:  $\sum_{i \in I} x_{it} - x_{it-1} - y_{it} + z_{it} = 0 \forall t \in \{1, \dots, T\}$
  - 2) individual rationality:  $(x_{it}, y_{it}, z_{it})_{t \in \{1,\dots,T\}} \in \operatorname{argmax}_{(x_{it}, y_{it}, z_{it})_{t \in \{1,\dots,T\}}} U_i \forall i$ ,
- where country  $i$ 's utility is:

$$U_i = \sum_{t=1}^T \frac{1}{(1+r)^t} (b_{it}(y_{it}) - g_{it}(z_{it}) - (c_{it}(x_{it}) - c_{it}(x_{it-1})) + f_{it}((x_{is})_{s=1,\dots,t}, (z_{is})_{s=1,\dots,t}, (z_{is})_{s=1,\dots,t}) + p_t(x_{it} + z_{it} - y_{it}))$$

It turns out that the global institution actually does not loose anything by restricting itself to using paths of period-by-period reward payment schemes, as long as it can freely borrow and save at the market interest rate  $r$ . I will show this in Lemma 11.

From now on we will use the following assumption:

**Assumption 4.**  $c''_{it+1}(x) < (1+r)c''_{it}(x) \forall i, t, x$ .

Assumption 4 has some plausibility: If extraction technology does not change much over time then assumption 4 clearly holds.<sup>12</sup>

**Lemma 9.** *Suppose assumption 4 holds. Then for any world market price path  $p$  there exists for each country a unique utility maximizing extraction path.*

*Proof.* See Appendix B.1. □

**Definition 6.** Consider a path of world market prices  $p = (p_t)_{t \in \{1,\dots,T\}}$  for coal. Let  $(x_{it}^*(p))_{t \in \{1,\dots,T\}}$  denote the unique (by Lemma 9) corresponding coal extraction path that maximizes country  $i$ 's utility in the absence of any reward payment schemes, assuming the world market price path  $p$ .

<sup>12</sup>In fact, assumption 4 is also plausible if extraction technology improves over time: The marginal extraction cost of coal increases as one extracts more and more of the reserves. This is because the remaining coal is harder to access. However, as time passes, technological improvements might reduce the associated increases in marginal extraction costs. This makes it plausible that we even have  $c''_{it+1}(x) \leq c''_{it}(x) \forall i, t, x$ .

**Lemma 10.** *Suppose assumption 4 holds. Then in the absence of any reward payment schemes there is a unique world market equilibrium.*

*Proof.* See Appendix B.2. □

**Definition 7. “Positive affine linear reward payment schemes”** are defined to be schemes where the global institution offers country  $i$  in period  $t$  the transfers  $f_{iyt}(y_{it}) = \max(0, \theta_{iyt}(\tilde{y}_{it} - y_{it}))$ ,  $f_{izt}(z_{it}) = \max(0, \theta_{izt}(z_{it} - \tilde{z}_{it}))$  and  $f_{ixt}(x_{it}) = \max(0, \theta_{ixt}(\tilde{x}_{it} - x_{it}))$ .

**Lemma 11 (The Surjectivity Lemma, dynamic version).** *Consider a combination of an allocation,  $(x_{it}, y_{it}, z_{it})_{i \in I, t \in \{1, \dots, T\}}$  with  $(x_{it})_{i \in \{1, \dots, T\}}$  non-decreasing satisfying  $\sum_{i \in I} x_{it} - y_{it} + z_{it} = 0 \forall t \in \{1, \dots, T\}$  and a world market price  $p$ . Then for any set of paths of period-by-period reward payment schemes implementing this combination, the global institution must end up paying at least the amount  $F_{iyt} := \sup_y b_{it}(y) - p_t y - (b_{it}(y_{it}) - p_t y_{it})$  on rewarding country  $i$  for demand reduction in period  $t$  and at least the amounts  $F_{izt} := \sup_z p_t z - g_{it}(z) - (p_t z_{it} - g_{it}(z_{it}))$  on rewarding country  $i$  for substitute expansion in period  $t$ .*

Moreover, the global institution must at least spend the amount  $F_{ix} := \sum_{t=1, \dots, T} \frac{1}{(1+r)^t} F_{ixt}$  in discounted money on rewarding country  $i$  for supply reduction, where  $F_{ixt}$  is defined as follows:

$$F_{ixt} := (p_t - \frac{1}{1+r} p_{t+1})(x_{it}^*(p) - x_{it}) + (c_{it}(x_{it}) - c_{it}(x_{it}^*(p))) + \frac{1}{1+r} (c_{it+1}(x_{it}^*(p)) - c_{it+1}(x_{it})) \text{ for } t \in \{1, \dots, T-1\}$$

$$F_{ixT} := p_T (x_{iT}^*(p) - x_{iT}) + (c_{iT}(x_{iT}) - c_{iT}(x_{iT}^*(p)))$$

Furthermore, there does indeed exist a set of paths of reward payment schemes implementing the allocation  $(x_{it}, y_{it}, z_{it})_{i \in I, t \in \{1, \dots, T\}}$  and the world market price path  $p$  that ends up paying exactly the amounts  $F_{iyt}$  on rewarding country  $i$  for demand reduction in period  $t$ , the amounts  $F_{izt}$  on rewarding country  $i$  for substitute substitute expansion in period  $t$  and the discounted amount  $F_{ixt}$  on rewarding country  $i$  for supply reduction in period  $t$ . Furthermore, there exists a set of positive affine linear period-by-period reward payment schemes achieving this.

Moreover, if the global institution can freely save and borrow then it does not anything by restricting itself to using only period-by-period reward payment schemes in the following sense: If a pair of an allocation and a price path can be implemented through some path of reward payment schemes  $f$  then it can also be implemented through a path  $f'$  of period-by-period reward payment schemes at which the total discounted value of transfers that end up getting paid out to countries is no more than under  $f$ .

*Proof.* See Appendix B.3 □

A natural question is whether the dynamic version of the Surjectivity Lemma could be strengthened. Are the  $F_{ixt}$  actually the minimal amounts that a global institution relying on period-by-period reward payment schemes has to end up paying in period  $t$  to induce country  $i$  to choose  $x_{it}$ ? It turns out that this

is not true in general. To see why, we note that if in some period the global institution were to pay out sufficiently large rents to country  $i$  for sufficiently low cumulative extraction, then this could even create incentives to conserve in previous periods and thereby decrease the required transfers in those previous periods.<sup>13</sup>

**Lemma 12 (The Price Preservation Lemma, dynamic version).** *Suppose  $(f_{it})_{i \in I, t \in \{1, \dots, T\}}$  is a set of reward payment schemes implementing a given allocation  $(x_{it}, y_{it}, z_{it})_{i \in I, t \in \{1, \dots, T\}}$  such that the discounted value of the aggregate transfers paid to countries is minimal. Then we must have: The entire world market price path  $(p_t)_{t \geq 0}$  is identical to when there is no reward payment scheme.*

*In particular, under a binding intertemporal budget constraint the global institution's optimal path of reward payment schemes leaves the world market price path of coal unchanged.*

*Proof.* See Appendix B.4. □

**Lemma 13 (The Constrained Efficiency Lemma, dynamic version).** *Suppose the global institution has an intertemporal budget of  $F$  and it can fully commit to any path of reward payment schemes. Consider the following set of allocations*

$$S_{F,p} := \{(x_{it}, y_{it}, z_{it})_{i \in I, t \in \{1, \dots, T\}} : \exists (f_{it})_{i,t}, p : (f_{it})_{i,t} \text{ implements } ((x_{it}, y_{it}, z_{it})_{i \in I, t \in \{1, \dots, T\}}, p) \text{ with spending } F\}$$

*We have: If  $(x_{it}, y_{it}, z_{it})_{i \in I, t \in \{1, \dots, T\}}$  maximizes global welfare  $W$  on  $S_{F,p}$  it maximizes global welfare amongst all allocations having the same value for the global damages  $\eta(\sum x_{iT})$  due to climate change.*

*In particular we have: The allocation implemented by the optimal reward payment scheme for a given discounted budget  $F$  is the unique allocation maximizing global welfare  $W$  under the constraint that the the global damages  $\eta(\sum x_{iT})$  due to climate change be a given constant.*

*Proof.* See Appendix B.5 □

Intuitively, the Constrained Efficiency Lemma should imply that it is optimal to reduce coal combustion in all periods so as to efficiently spread the mitigation effort. The following Lemma confirms this:

**Lemma 14 (The Monotone Mitigation Lemma).** *Suppose the global institution has an intertemporal budget of  $F = \sum_{t=1}^T \frac{F_t}{(1+r)^t}$ .*

*Denoting by  $x_{it}(F)$  the cumulative coal extraction of  $i$  by the end of period  $t$  at the optimal mechanism, we have:*

$$0 > \frac{dx_{i1}}{dF} > \dots > \frac{dx_{iT}}{dF} \forall i, t.$$

*In particular, increasing the budget  $F$  reduces coal extraction (and use),  $x_{it}(F) - x_{it-1}(F)$ , of each country in each period.*

<sup>13</sup>If we suppose that the global institution cannot commit to a path of reward payment schemes then this possibility vanishes: Even if future rents *could* contribute to participation incentives now, once the future arrives the global institution will no longer have incentives to pay rents. In fact, if the global institution cannot commit to a path of reward payment schemes, the Surjectivity Lemma can be strengthened, as shown in Lemma Stern (2021).

*Proof.* See Appendix B.6 □

**Corollary 4 (The Monotone Optimal Spending Corollary).** *Suppose the global institution has an intertemporal budget of  $F = \sum_{t=1}^T \frac{F_t}{(1+r)^t}$ .*

*Denoting by  $F_{tx}(F)$ ,  $F_{ty}(F)$ ,  $F_{tz}(F)$  the global institution's optimal spending on supply reduction, demand reduction and substitute expansion in period  $t$ , we have:*

$$\frac{dF_{tx}}{dF} > 0, \frac{dF_{ty}}{dF} > 0, \frac{dF_{tz}}{dF} > 0 \forall t$$

*Proof.* By the Dynamic Price Preservation Lemma 12, world market prices must under the optimal reward payment scheme be just as in the absence of any reward payment scheme. By the Monotone Mitigation Lemma 14, both cumulative coal extraction and marginal coal extraction are strictly decreasing functions of  $F$ . The former implies that the global institution's spending on rewarding supply reduction strictly increases and the latter implies (invoking the Constrained Efficiency Lemma) that the spending on demand reduction and substitute expansion must strictly increase. □

**Corollary 5 (The Coase-Does-Not-Hold Corollary).** *Suppose that the global institutions restricts its supply side spending to the last period. (This would effectively be the case if the global institution were to restrict itself to buying up coal deposits.) Then the global institution cannot achieve optimal welfare.*

*Thus we have: The "Coase Theorem" does not hold when the global institution's budget is insufficient to achieve a full correction of the global externalities.*

*Proof.* This follows directly from the Monotone Optimal Spending Corollary. □

**Corollary 6 (The Time Inconsistency Corollary).** *Suppose the global institution announces at time 1 the optimal reward payment scheme assuming it fully commits to it. Suppose that at time  $T$  the global institution reneges, to all countries' surprise, on its promised reward payment schemes and instead offers a set of reward payment schemes that maximize global welfare, given the state  $(x_{iT-1})_{i \in I}$  of the world then. Then this new reward payment scheme used for the last period involves less spending on rewarding supply reduction than the originally announced reward payment scheme.*

*Proof.* See Appendix B.7. □

Intuitively, corollary 6 makes sense: Part of the benefits of future spending on rewarding supply reduction are that they will by raising world market prices then increase the incentives to conserve in the prior periods. However, once the future arrives, this consideration disappears and thus the global institution is better off spending less on rewarding supply reduction.<sup>14</sup> This result suggests a

<sup>14</sup>This argument also suggest that the Time Inconsistency Corollary 6 could be strengthened to a claim that if the global institution reoptimizes from a period  $t > 1$  onwards (instead of just for period  $T$  as in the statement of the corollary) then it will spend less on rewarding supply reduction than under the initially announced reward payment scheme.

valuable role for deposit purchase funds because they could mitigate this time consistency problem: The global institution could early on purchase appropriate coal deposits in a way that does not undermine future incentives to spend on supply reduction: Buying up these appropriately chosen coal deposits early on amounts in the model to effectively increasing spending on supply reduction in the last period.

## 7 Implementation

We have seen that for any given budget at the global institution's disposal the optimal (in terms of global welfare and also in terms of reductions in eventual cumulative emissions) allocation can be implemented by having the global institution announce and commit to a set of reward payment schemes of the following form: In each period, each country is rewarded on the basis of its energy use and renewable energy production during the period and also on the basis of its cumulative coal extraction by the end of the period. Moreover, we have seen that positive affine linear reward payment schemes can be used for this (Lemma 1).

On the supply side, this looks as follows: For each period  $t$ , the global institution can define a reference level for country  $i$ 's cumulative coal extraction and reward it proportionally to the amount by which its actual cumulative coal extraction by the end of the period is below the reference level.

A well-recognized problem with these kinds of schemes in practice is that errors in the specification of the reference levels can drastically undermine its performance (Angelsen (2017)). For example, suppose that the global institution's budget is such that the optimal reward payment scheme (under complete information) reduces cumulative coal extraction in a given country by 10% by a given year relative to the situation in the absence of the global institution. Now suppose that the global institution underestimates for that year by 10% the cumulative emissions occurring in the country in the absence of any reward payment schemes. Then it would end up not rewarding the country at all even if the country makes the substantial effort of reducing its cumulative coal extraction by 10%. Such an estimation mistake by 10% can arguably happen quite easily, given that the global institution only has incomplete information about factors such as the future discovery of fossil fuel deposits.

This problem can to varying degrees be mitigated through alternative schemes on the supply side. One alternative scheme is to buy up coal deposits and thereby prevent them from being exploited. The informational requirements for this approach are arguably much less demanding. Fossil fuel deposits that are traded on the world market are expected by the market participants to be exploited eventually. Thus by buying up such deposits to conserve them a global institution could be quite confident that it is causing emission reductions.

Given these potential informational advantages of the deposit purchase approach for supply reduction, the question arises as to whether a global institution with a fixed intertemporal budget as analyzed in this paper should do all its

supply side intervention via deposit purchase. First best reasoning would suggest that nothing would be lost by restricting supply side intervention to deposit purchase, at least in this paper’s model where global externalities are determined purely by the eventual amount of coal that is combusted. However, by the Coase-Does-Not-Hold-Corollary 5 this is not the case in the realistic case where the global institution’s budget is insufficient to fully correct the externalities from fossil fuels.

Thus restricting the supply side approaches to deposit purchase would lead to a loss in global welfare. Given this result, it is important to explore alternative ways to reward supply reduction. One such approach could be to reward countries in each period for taxing coal extraction in that period.

Specifically, the reward payment for a given country in a given period could be defined to be proportional to some function of the country’s tax rate on coal extraction. The actual reward payment could be defined to be this function times an estimate of the country’s coal extraction in the period in the absence of any reward payments. This would ensure that countries of different “sizes” would get commensurate reward payments. Thus the incentive power would be distributed evenly across countries of different sizes which is required for efficiency given the convexity of abatement costs. Analogous schemes could be used for rewarding countries for taxing fossil fuel combustion which induce them to reduce energy use and to expand renewables (see Stern (2021) for an operational version of such a proposal).

Such schemes would be much less sensitive to errors in the estimations of the countries’ business as usual (BAU) extraction paths. To see why, consider the example discussed above where the global institution could under complete information induce countries to reduce coal extraction by 10% relative to what they would do otherwise, but where it underestimates a country’s BAU coal extraction by 10%. Under the tax-based reward payment scheme just discussed this estimation mistake would lead to the reward payment being 10% too low. The associated welfare losses would likely be small.

This suggests that rewarding countries each period for taxing coal extraction is a promising approach to inducing countries to reduce coal supply. Intuitively, it is clear that in this paper’s model a version of this approach can achieve the optimum for the global institution, for any given value of its budget. To see why, we first note that by Lemma 7, the global institution can in the static model restrict itself to rewarding countries on the basis of tax and subsidy rates. Now in the dynamic setting, we know by the Monotone Optimal Spending Corollary 4 that the global institution needs to spend in each period strictly positive amounts on rewarding supply reduction, demand reduction and substitute expansion. Intuitively, it is clear that the global institution can use tax/subsidy based reward payment schemes for this instead of the quantity based reward payment schemes in the original Corollary 4. It turns out that this argument goes through formally with some minor qualifications, as I show in appendix C.

Given this result about the sufficiency of tax-based reward payment schemes, the question arises as to whether this approach should be used exclusively or whether deposit purchase should be used in parallel on the supply side. The



Time Inconsistency Corollary 6 provides an argument for relying at least partially on deposit purchase: The ex ante optimal path of reward payment schemes involves spending more on rewarding supply reduction in the future than what will be optimal ex post once the future arrives. In the model, the global institution can partially mitigate this time inconsistency problem through deposit purchase, as I explained at the end of section 6.

## 8 Conclusion and limitations

This paper has analyzed the problem faced by global institutions such as the Green Climate Fund, the Global Environment Facility and the Climate Investment Fund. With the part of the budgets allocated to climate change mitigation, these institutions can be viewed as being able to reward countries for reducing fossil fuel supply, reducing energy demand and expanding renewable energy. So far, these institutions only pursue the demand reduction and the substitute expansion approaches. The results in this paper suggest that it would be optimal for them to pursue a mixture of supply, demand and substitute based approaches to curbing fossil fuels.

In fact, for a given intertemporal budget and under full commitment, it is optimal to reward each country in each period for demand reduction and substitute expansion as well as for having extracted less of the fossil fuel than what would have been optimal for it in the absence of such a reward payment. On the demand and substitute side, this can be achieved by rewarding countries on the basis of the rate at which they tax the combustion of the fossil fuel. Fossil fuel tax rates are a good measure of a country's effort to reduce demand for fossil fuels and to expand substitutes. Thus they are plausibly superior to quantities as variables to contract on, since the latter depend on many other factors such as economic growth.

The results of this paper suggest that it could be valuable to explore how carbon pricing reward funds could best be used on the supply side, where countries could be rewarded for taxing the extraction of fossil fuel. Indeed, I find that the commonly proposed deposit purchase approach to supply reduction is on its own unable to implement the optimal mechanism (corollary 5). Deposit purchase funds could play a valuable role in mitigating a time consistency problem identified in corollary 6. But it should plausibly be complemented by carbon pricing reward funds on the supply side.

An important limitation of this paper's model is that it abstracts away from informational asymmetries. Since in reality countries have private information about their costs and benefits of extracting coal, using energy and producing renewable energy, most of them will reap informational rents at the optimal mechanism. In Stern (2021) I develop a model taking this into account. Specifically, I study the optimal mechanism for reducing fossil fuel demand if the global institution can only condition its reward payments on each country's tax/subsidy rate on the fossil fuels. The model could be applied to the supply side and the substitute side and then integrated into the model presented in the current

paper.

A further limitation of the model used in this paper is that does not take into account the adverse effects of fossil fuel rents on global welfare via the natural resource curse (Ross (2015)). This consideration weighs in favor of focusing on demand reduction and substitute expansion instead of supply reduction. Possibly, it could rationalize the current absence of international institutions paying countries for reducing fossil fuel supply.

## A Proofs for section 3

### A.1 Proof of Surjectivity Lemma 1

Suppose first, hypothetically, the global institution could make countries pay transfers and suppose it were to impose the reward payment scheme  $\theta_{ixt}(\tilde{x}_{it} - x_{it})$ . Then since  $x \mapsto px - C_i(x) + \theta_{ix}(\tilde{x}_i - x)$  is concave, its global optimum is  $x_i$  iff  $\theta_{ix} = p - C'_i(x_i)$ . Now suppose the global institution offers instead the reward payment scheme  $f_{ix}(x_i) = \max(0, \theta_{ixt}(\tilde{x}_{it} - x_{it}))$ .  $x_{it}$  is still an optimal choice for  $i$  iff  $px_i - C_i(x_i) + \theta_{ixt}(\tilde{x}_{it} - x_{it}) \geq \sup_x px - C_i(x)$ , or equivalently:

$$\theta_{ixt}(\tilde{x}_{it} - x_{it}) \geq F_{ix} = \sup_x px - C_i(x) - (px_i - C_i(x_i))$$

Thus we have shown that the minimal transfer required to pay the country  $i$  to induce it to choose  $x_i$  under some affine linear reward payment scheme is indeed  $F_{ix} = \sup_x px - C_i(x) - (px_i - C_i(x_i))$ . Analogously, one can prove the claims for the other variables.

What is left to be proved is that there does not exist another set  $(f_i)_{i \in I}$  of reward payment schemes implementing  $((x_i, y_i, z_i)_{i \in I}, p)$  with a strictly smaller budget. To establish this, we note that the individual rationality condition implies:

$$B_i(y_i) - C_i(x_i) - G_i(z_i) + p(x_i - y_i + z_i) + f_i(x_i, y_i, z_i) \geq \sup_{(x,y,z)} B_i(y) - C_i(x) - G_i(z) + p(x - y + z) = \sup_y B_i(y) - py - \sup_x (C_i(x) - px) - \sup_z (G_i(z) - gz),$$

so we have:

$$f_i(x_i, y_i, z_i) \geq \sup_x px - C_i(x) - (px_i - C_i(x_i)) + \sup_y B_i(y) - py - (B_i(y_i) - py_i) + \sup_z pz - G_i(z) - (pz_i - G_i(z_i)) = F_{ix} + F_{iy} + F_{iz}.$$

### A.2 Proof of Lemma 2

*Proof.* If the budget is sufficiently large, the global institution's problem's Lagrangian becomes equal to:

$$L = \sum_{j \in I} U_j - \eta(\sum_{j \in I} x_j) + \mu \sum_{i \in I} (x_i - y_i + z_i)$$

where  $\mu$  is the Lagrange multiplier associated with the feasibility constraint.

The first order conditions are:

$$\frac{\partial L}{\partial x_i} = \frac{\partial U_i}{\partial x_i} - \eta + \mu = -C'_i(x_i) - \eta + \mu = 0$$

$$\frac{\partial L}{\partial y_i} = \frac{\partial U_i}{\partial y_i} - \eta + \mu = B'_i(y_i) - \mu = 0$$

$$\frac{\partial L}{\partial z_i} = \frac{\partial U_i}{\partial z_i} - \eta + \mu = -G'_i(z_i) + \mu$$

Thus the first best is characterized by the following conditions plus the feasibility condition  $\sum_{i \in I} (x_i - y_i + z_i)$ :

$$B'_i(y_i) = B'_j(y_j) = G'_i(z_i) = G'_j(z_j) = C'_i(x_i) + \eta = C'_j(x_j) + \eta \forall i, j \in I$$

In particular, the world market price  $p$  of energy does not appear in this characterization. The greater the emphasis on rewarding countries for supply reduction, the larger the resulting price  $p$  will be. But as long as the required overall budget does not exceed the available budget, this does not matter for global welfare.  $\square$

### A.3 Proof of Price Preservation Lemma 3

*Proof.* Based on the Surjectivity Lemma (1), we view the global institution as choosing the price  $p$  and the allocation  $(x_i, y_i, z_i)_{i \in I}$ . Global welfare is determined by the allocation,  $(x_i, y_i, z_i)_{i \in I}$ , and the price  $p$  only is relevant because it affects the aggregate transfers required to get all countries to participate. The minimal required transfers  $(F_{ix}, F_{iy}, F_{iz})_{i \in I}$  are given by the Surjectivity Lemma. In particular, the minimal aggregate required transfer is given by  $\sum_{i \in I} F_{ix} + F_{iy} + F_{iz}$ .

We now show that  $\sum_{i \in I} F_{ix} + F_{iy} + F_{iz}$  is convex when viewed as a function of  $p$ . To do so, we use that  $F_{ix} := \sup_x px - C_i(x) - (px_i - C_i(x_i))$ ,  $F_{iy} := \sup_y B_i(y) - py - (B_i(y) - py)$ ,  $F_{iz} := \sup_z pz - G_i(z) - (pz_i - G_i(z_i))$  and we compute:

$$\frac{d}{dp}(\sum_{i \in I} F_{ix} + F_{iy} + F_{iz}) = \sum_{i \in I} x_i^*(p) - x_i - (y_i^*(p) - y_i) + z_i^*(p) - z_i$$

where  $x_i^*(p)$  denotes country  $i$ 's supply function in the absence of any mechanism, i.e.  $px - C_i(x) = \operatorname{argmax}_x px - C_i(x)$  and analogously for  $y_i^*(p)$  and  $z_i^*(p)$ . By market clearing, we have:

$$\frac{d}{dp}(\sum_{i \in I} F_{ix} + F_{iy} + F_{iz}) = \sum_{i \in I} x_i^*(p) - y_i^*(p) + z_i^*(p)$$

But this is just the excess supply function, which is strictly increasing in  $p$  as can be deduced directly from assumption 1. Thus the optimal  $p$  is characterized by the condition

$$\sum_{i \in I} x_i^*(p) - y_i^*(p) + z_i^*(p) = 0$$

The price obtaining at the market equilibrium in the absence of any mechanism satisfies this condition by market clearing. It is the unique price satisfying this condition.  $\square$

### A.4 Proof of Constrained Efficiency Lemma 4

*Proof.* The basic reason for this result is as follows: The global institution has two considerations to take into account: it cares intrinsically about the countries' aggregate utility and it wants to minimize the required transfers. But the outside options are determined by the price and so the required transfers decrease in the countries' aggregate utility. Thus the two considerations perfectly align. I will now flesh out this argument in full formal detail.

Consider a fixed price  $p$ . Consider the set  $S$  of all allocations satisfying the market clearing condition and having a given value  $X$  for  $\sum x_i$ .

$$S := \{(x_i, y_i, z_i) : \sum_{i \in I} (x_i - y_i + z_i) = 0, \sum_{i \in I} x_i = X\}$$

Let  $\mu$  denote the Lagrange multiplier associated with the market clearing constraint. Let  $\beta$  denote the Lagrange multiplier associated with the global institution's budget constraint. The global institution's Lagrangian is:

$$L = \sum_{i \in I} U_i - \eta(\sum_{j \in I} x_j) + \mu \sum_{i \in I} (x_i - y_i + z_i) - \beta(\sum_{i \in I} F_{ix} + F_{iy} + F_{iz} - F)$$

When choosing amongst allocations in this set  $S$ , there are only two terms in the Lagrangian that are affected, namely  $\sum_{i \in I} U_i$  and  $-\beta(\sum_{i \in I} F_{ix} + F_{iy} + F_{iz})$ :

$$L = \sum_{i \in I} U_i - \eta(\sum_{j \in I} x_j) + \mu \sum_{i \in I} (x_i - y_i + z_i) - \beta(\sum_{i \in I} F_{ix} + F_{iy} + F_{iz} - F)$$

So we can write:

$$L = \sum_{i \in I} U_i - \beta(\sum_{i \in I} F_{ix} + F_{iy} + F_{iz}) + \phi$$

where  $\phi$  does not depend on the allocation (as long as the allocation is chosen from the set  $S$ ).

Using the expressions for the minimal transfers, we obtain:

$$L = \sum_{i \in I} U_i - \beta(\sum_{i \in I} F_{ix} + F_{iy} + F_{iz}) + \phi$$

We have:

$$\sum_{i \in I} F_{ix} + F_{iy} + F_{iz} = \sup_x px - C_i(x) - (px_i - C_i(x_i)) + \sup_y B_i(y) - py - (B_i(y) - py) + \sup_z pz - G_i(z) - (py - G_i(y))$$

$$\sum_{i \in I} F_{ix} + F_{iy} + F_{iz} = \sum_{i \in I} \sup_{x,y,z} p(x - y + z) - C_i(x) + B_i(y) - G_i(z) + C_i(x_i) - B_i(y_i) + G_i(z_i)$$

$$\sum_{i \in I} F_{ix} + F_{iy} + F_{iz} = F + \sum_{i \in I} \sup_{x,y,z} p(x - y + z) - C_i(x) + B_i(y) - G_i(z) - U_i$$

Now we can write the Lagrangian as follows:

$$L = (1 + \beta) \sum_{i \in I} U_i + \phi^\#$$

where  $\phi^\#$  does not depend on the allocation (as long as the allocation is chosen from the set  $S$ ). Thus for any given  $p$  and a given value  $X$  for  $\sum x_i$ , the global institution's optimization problem for the allocation  $(x_i, y_i, z_i)_{i \in I}$  boils down to simply maximizing  $\sum_{i \in I} U_i$  under the constraint that  $\sum x_i = X$ .  $\square$

## A.5 Proof of the Interior Solution Corollary 1 and the Optimal Budget Split Corollary for Small Budgets 2

*Proof.* (of the Interior Solution Corollary and the Optimal Budget Split Corollary for Small Budgets) By the Price Preservation Lemma, the price at the optimal mechanism is the same as at the market equilibrium in the absence of any mechanism. We denote this price simply by  $p$ .

As before, let  $\mu$  denote the Lagrange multiplier associated with the market clearing constraint. Let  $\beta$  denote the Lagrange multiplier associated with the global institution's budget constraint. The global institution's Lagrangian is:

$$L = \sum_{i \in I} U_i - \eta(\sum_{j \in I} x_j) + \mu \sum_{i \in I} (x_i - y_i + z_i) - \beta(\sum_{i \in I} F_{ix} + F_{iy} + F_{iz} - F)$$

By differentiating the Lagrangian with respect to the allocation, we obtain the following optimality conditions:

$$h_{xi}(x, y, z, \mu) := (1 + \beta)(p - C'_i(x_i)) + \mu - \eta' = 0$$

$$h_{yi}(x, y, z, \mu) := (1 + \beta)B'_i(y_i) - p - \mu = 0$$

$$h_{zi}(x, y, z, \mu) := (1 + \beta)(p - G'_i(z_i)) - p + \mu = 0$$

We also have the market clearing condition:

$$f_\mu(x, y, z, \mu) := \sum_i x_i - y_i + z_i = 0 \quad (4)$$

Define  $\sigma := \frac{1}{1+\beta}$ , so  $\sigma = 0$  corresponds to the case where the budget  $F = 0$  and  $\sigma = 1$  corresponds to the case where the budget  $F$  is the minimal amount sufficient to implement the global optimum. With this, we can rewrite the optimality conditions as follows:

$$h_{x_i}(x, y, z, \mu) := p - C'_i(x_i) + \sigma(\mu - \eta') = 0 \quad (5)$$

$$h_{y_i}(x, y, z, \mu) := B'_i(y_i) - p - \sigma\mu = 0 \quad (6)$$

$$h_{z_i}(x, y, z, \mu) := (1 + \beta)(p - G'_i(z_i)) - p + \sigma\mu = 0 \quad (7)$$

Now we can study what happens as we relax the budget constraint, which corresponds to increasing  $\sigma$ . Let us denote by  $h := ((h_{x_i})_{i \in I}, (h_{y_i})_{i \in I}, (h_{z_i})_{i \in I}, h_\mu)$ , where  $h$  is a vector function whose components are defined in equations 5, 6, 7 and 4. The  $\sigma$  determines  $(x, y, z, \mu)$  via the condition  $h(x, y, z, \mu) = 0$ .

For  $\sigma \in (0, \infty)$  we know that this system has a unique solution. To see this, consider a fixed  $\sigma > 0$ . We can think of a choice of  $\mu$  as determining the  $(x_i, y_i, z_i)_{i \in I}$ . Define  $g(\mu) := h_\mu(x(\mu), y(\mu), z(\mu), \mu)$ , where  $x(\mu)$  denotes the vector of the  $x_i$  determined via equation 5 etc.. Each of the  $x_i$  and  $z_i$  are strictly increasing in  $\mu$ , whilst all of the  $y_i$  are strictly decreasing in  $\mu$ . Hence  $g$  is increasing. Moreover,  $g(0) < 0$ . To see why, we note that for  $\mu = 0$  the  $y_i$  and the  $z_i$  are as in the absence of any mechanism whilst the  $x_i$  are strictly smaller. We also have that  $g(\eta) > 0$ . To see why, we note that for  $\mu = \eta$  the  $x_i$  are as in the absence of any mechanism, whilst the  $y_i$  are strictly smaller and the  $z_i$  are strictly larger.

Given that thus  $g(0) < 0$ ,  $g(\eta) > 0$  and that  $g(\mu)$  is increasing, we can apply the intermediate value theorem as long as  $g$  is continuous. But all the  $x_i(\mu), y_i(\mu), z_i(\mu)$  are continuous which implies that  $g$  is continuous. Hence, by the intermediate value theorem, there exists a unique  $\mu$  such that equation 4 holds and this  $\mu$  is in  $(0, \eta)$ . This shows that there is a unique solution to equations 5, 5, 7 and 4. We now denote the unique solution by  $(x, y, z, \mu)(\sigma)$ .

Now we will show that  $(x, y, z, \mu)(\sigma)$  is continuously differentiable at all  $\sigma \in (0, 1)$ . For this it suffices by the implicit function theorem to show that the Jacobian of  $h$  is nonsingular on  $(0, 1)$ , since  $f$  is continuously differentiable. The Jacobian  $J$  is as follows:

$$J = \begin{pmatrix} -C''(x) & 0 & 0 & \sigma \\ 0 & B''(y) & 0 & -\sigma \\ 0 & 0 & -G''(z) & \sigma \\ 1 & -1 & 1 & 0 \end{pmatrix}$$

where  $C''(x)$  is a diagonal matrix with entries  $C''_i(x_i)$  etc. and in a slight abuse of notation the 6 zeros in the upper left denote  $|I|$  by  $|I|$  matrices with

all entries 0 and the  $\sigma$  denote column vectors of length  $|I|$ . For  $\sigma > 0$ ,  $J$  is non-singular, since by assumption  $C_i'' > 0, B_i'' < 0, G_i'' > 0$ .

To see why, we note that if we want to write the bottom row vector as a linear combination of the other rows, then the weight given to each of the first  $|I|$  rows (corresponding to the  $(x_i)_{i \in I}$ ) must be strictly positive, whilst the weight given to the rows from row  $|I| + 1$  to row  $2|I|$  must be negative, whilst the weight given to the rows from row  $2|I| + 1$  to row  $3|I|$  must be negative. But all this together implies that the last component of this linear combination will be strictly positive, in contradiction to the fact that the last component of the last row is 0.

Hence we have established by the implicit function theorem that  $(x, y, z, \mu)(\sigma)$  is continuously differentiable at all  $\sigma \in (0, \infty)$ . We also note that  $\lim_{\sigma \rightarrow 0^+} J$  is singular. This explains why we will now need to do some more work to establish a fact that we will later need, namely that “generically”  $\lim_{\sigma \rightarrow 0^+} \sigma \frac{d\mu}{d\sigma} = 0$ .

Differentiation of the first order conditions (5,6,7) with respect to  $\sigma$  yields:

$$C_i''(x_i) \frac{dx_i}{d\sigma} = \mu - \eta + \sigma \frac{d\mu}{d\sigma} \quad (8)$$

$$B_i''(y_i) \frac{dy_i}{d\sigma} = \mu + \sigma \frac{d\mu}{d\sigma} \quad (9)$$

$$G_i''(z_i) \frac{dz_i}{d\sigma} = \mu + \sigma \frac{d\mu}{d\sigma} \quad (10)$$

Above we showed that  $\mu(\sigma) \in (0, \eta) \forall \sigma > 0$ . From this it follows that  $\sigma \frac{d\mu}{d\sigma} \in (-\mu, \eta) \forall \sigma > 0$ . This is because for  $\sigma$  such that  $\sigma \frac{d\mu}{d\sigma} > \eta - \mu$  we would have  $\frac{dx_i}{d\sigma} > 0 \forall i$  (by equation 8),  $\frac{dy_i}{d\sigma} < 0 \forall i$  (by equation 9) and  $\frac{dz_i}{d\sigma} > 0 \forall i$  (by equation 10), so that  $\frac{d}{d\sigma}(x_i - y_i + z_i) > 0 \forall i$  which would contradict the market clearing condition.

Similarly, for  $\sigma \frac{d\mu}{d\sigma} < -\mu$  we would have  $\frac{dx_i}{d\sigma} < 0 \forall i$  (by equation 8),  $\frac{dy_i}{d\sigma} > 0 \forall i$  (by equation 9) and  $\frac{dz_i}{d\sigma} < 0 \forall i$  (by equation 10) so that  $\frac{d}{d\sigma}(x_i - y_i + z_i) > 0 \forall i$ , which would contradict the market clearing condition.

We are now ready to show that “generically” we must have  $\lim_{\sigma \rightarrow 0} \sigma \frac{d\mu}{d\sigma} = 0$ . To establish this, let us first suppose that  $\lim_{\sigma \rightarrow 0} \sigma \frac{d\mu}{d\sigma}$  exists. Suppose we have  $\lim_{\sigma \rightarrow 0} \sigma \frac{d\mu}{d\sigma} = K \in (0, \infty]$ . Then there exists some  $\tilde{\sigma} > 0$  such that for all  $\sigma \in (0, \tilde{\sigma}]$  we have  $\frac{d\mu}{d\sigma} > \frac{K}{2\sigma}$ . Integrating this yields for all  $\sigma \in (0, \tilde{\sigma}]$ :  $\mu(\tilde{\sigma}) - \mu(\sigma) = \int_{s=\sigma}^{\tilde{\sigma}} \frac{d\mu}{ds} > \int_{s=\sigma}^{\tilde{\sigma}} \frac{K}{s} ds = K(\log(\tilde{\sigma}) - \log(\sigma))$ . Rearranging yields:  $\mu(\sigma) < \mu(\tilde{\sigma}) - K(\log(\tilde{\sigma}) - \log(\sigma))$ . But this would imply that  $\lim_{\sigma \rightarrow 0} \mu(\sigma) = -\infty$ , in contradiction to the fact that, as shown above, we always have  $\mu(\sigma) \in (0, \eta) \forall \sigma$ . Similarly, we can show that  $\lim_{\sigma \rightarrow 0} \sigma \frac{d\mu}{d\sigma} = K < 0$  leads to a contradiction.

Now in our quest to show that  $\lim_{\sigma \rightarrow 0} \sigma \frac{d\mu}{d\sigma} = 0$  we only have one more case to show to be impossible, namely the case where  $\lim_{\sigma \rightarrow 0} \sigma \frac{d\mu}{d\sigma}$  does not exist. But in this case  $\frac{d\mu}{d\sigma}$  has to fluctuate infinitely often by unbounded amounts as  $\sigma$  approaches 0. Intuitively, this will “generically” never happen.

Now using the surjectivity lemma 1 yields:

$$\frac{dF_x}{d\sigma} = \frac{d(\sum_i F_{ix})}{d\sigma} = \sum_i (C'_i(x_i) - p) \frac{\mu - \eta + \sigma \frac{d\mu}{d\sigma}}{C''_i(x_i)}$$

$$\frac{dF_y}{d\sigma} = \frac{d(\sum_i F_{iy})}{d\sigma} = \sum_i -(B'_i(y_i) - p) \frac{\mu + \sigma \frac{d\mu}{d\sigma}}{B''_i(y_i)}$$

$$\frac{dF_z}{d\sigma} = \frac{d(\sum_i F_{iz})}{d\sigma} = \sum_i (G'_i(x_i) - p) \frac{\mu + \sigma \frac{d\mu}{d\sigma}}{G''_i(x_i)}$$

Using the optimality conditions yields:

$$\frac{dF_x}{d\sigma} = \sum_i \sigma (\mu - \eta) \frac{\mu - \eta + \sigma \frac{d\mu}{d\sigma}}{C''_i(x_i)}$$

$$\frac{dF_y}{d\sigma} = \sum_i \sigma \mu \frac{\mu + \sigma \frac{d\mu}{d\sigma}}{-B''_i(y_i)}$$

$$\frac{dF_{i,z}}{d\sigma} = \sum_i \sigma \mu \frac{\mu + \sigma \frac{d\mu}{d\sigma}}{G''_i(x_i)}$$

$$\frac{\frac{dF_x}{d\sigma}}{\frac{dF_x}{d\sigma} + \frac{dF_y}{d\sigma} + \frac{dF_z}{d\sigma}} = \frac{(\mu - \eta + \sigma \frac{d\mu}{d\sigma})(\mu - \eta) \sum \frac{1}{C''_i(x_i)}}{(\mu - \eta + \sigma \frac{d\mu}{d\sigma})(\mu - \eta) \sum \frac{1}{C''_i(x_i)} + (\mu + \sigma \frac{d\mu}{d\sigma}) \mu \sum \frac{1}{-B''_i(y_i)} + (\mu + \sigma \frac{d\mu}{d\sigma}) \mu \sum \frac{1}{G''_i(x_i)}}$$

Now we differentiate the market clearing condition, getting

$$\sum \frac{dx_i}{d\sigma} + \frac{dz_i}{d\sigma} = \sum \frac{dy_i}{d\sigma}$$

Substituting into this gives:

$$(\mu - \eta + \sigma \frac{d\mu}{d\sigma}) \sum \frac{1}{C''_i(x_i)} + (\mu + \sigma \frac{d\mu}{d\sigma}) \sum \frac{1}{G''_i(x_i)} = (\mu + \sigma \frac{d\mu}{d\sigma}) \sum \frac{1}{B''_i(y_i)}$$

$$\text{so } (\mu - \eta + \sigma \frac{d\mu}{d\sigma}) \sum \frac{1}{C''_i(x_i)} = -(\mu + \sigma \frac{d\mu}{d\sigma}) (\sum \frac{1}{-B''_i(y_i)} + \sum \frac{1}{G''_i(x_i)})$$

which we can plug in to get:

$$\begin{aligned} \frac{\frac{dF_x}{d\sigma}}{\frac{dF_x}{d\sigma} + \frac{dF_y}{d\sigma} + \frac{dF_z}{d\sigma}} &= \frac{-(\mu + \sigma \frac{d\mu}{d\sigma}) (\sum \frac{1}{-B''_i(y_i)} + \sum \frac{1}{G''_i(x_i)}) (\mu - \eta)}{-(\mu + \sigma \frac{d\mu}{d\sigma}) (\sum \frac{1}{-B''_i(y_i)} + \sum \frac{1}{G''_i(x_i)}) (\mu - \eta) + (\mu + \sigma \frac{d\mu}{d\sigma}) \mu \sum \frac{1}{-B''_i(y_i)} + (\mu + \sigma \frac{d\mu}{d\sigma}) \mu \sum \frac{1}{G''_i(x_i)}} \\ &= \frac{-(\mu + \sigma \frac{d\mu}{d\sigma}) (\sum \frac{1}{-B''_i(y_i)} + \sum \frac{1}{G''_i(x_i)}) (\mu - \eta)}{-(\mu + \sigma \frac{d\mu}{d\sigma}) (\sum \frac{1}{-B''_i(y_i)} + \sum \frac{1}{G''_i(x_i)}) (\mu - \eta) + (\mu + \sigma \frac{d\mu}{d\sigma}) \mu \sum \frac{1}{-B''_i(y_i)} + (\mu + \sigma \frac{d\mu}{d\sigma}) \mu \sum \frac{1}{G''_i(x_i)}} \\ &= \frac{-\sum \frac{1}{-B''_i(y_i)} + \sum \frac{1}{G''_i(x_i)} (\mu - \eta)}{-\sum \frac{1}{-B''_i(y_i)} + \sum \frac{1}{G''_i(x_i)} (\mu - \eta) + \mu \sum \frac{1}{-B''_i(y_i)} + \mu \sum \frac{1}{G''_i(x_i)}} \\ &= \frac{-(\mu - \eta)}{-(\mu - \eta) + \mu} \\ &= \frac{\eta - \mu}{\eta} \end{aligned}$$

Increasing  $\sigma$  corresponds to increasing  $F$ . Therefore, we must have  $\frac{dF_x}{d\sigma} + \frac{dF_y}{d\sigma} + \frac{dF_z}{d\sigma} > 0$ . But we already established that  $\mu \in (0, \eta)$ , so it follows that  $\frac{dF_x}{d\sigma} > 0$ .

We now proceed analogously for  $F_y$ :

$$\begin{aligned} \frac{\frac{dF_y}{d\sigma}}{\frac{dF_x}{d\sigma} + \frac{dF_y}{d\sigma} + \frac{dF_z}{d\sigma}} &= \frac{(\mu + \sigma \frac{d\mu}{d\sigma}) \mu \sum \frac{1}{-B''_i(y_i)}}{(\mu + \sigma \frac{d\mu}{d\sigma}) (\sum \frac{1}{-B''_i(y_i)} + \sum \frac{1}{G''_i(x_i)}) (\eta - \mu) + (\mu + \sigma \frac{d\mu}{d\sigma}) \mu \sum \frac{1}{-B''_i(y_i)} + (\mu + \sigma \frac{d\mu}{d\sigma}) \mu \sum \frac{1}{G''_i(x_i)}} \\ \frac{\frac{dF_y}{d\sigma}}{\frac{dF_x}{d\sigma} + \frac{dF_y}{d\sigma} + \frac{dF_z}{d\sigma}} &= \frac{\mu \sum \frac{1}{-B''_i(y_i)}}{\eta (\sum \frac{1}{-B''_i(y_i)} + \sum \frac{1}{G''_i(x_i)})} \end{aligned}$$

Hence  $\frac{dF_y}{d\sigma} > 0$ , since  $\mu \in (0, \eta)$ .

Now we proceed analogously to  $F_z$ :

$$\frac{\frac{dF_z}{d\sigma}}{\frac{dF_x}{d\sigma} + \frac{dF_y}{d\sigma} + \frac{dF_z}{d\sigma}} = \frac{\mu \sum \frac{1}{G''_i(x_i)}}{\eta (\sum \frac{1}{-B''_i(y_i)} + \sum \frac{1}{G''_i(x_i)})}$$

Hence  $\frac{dF_z}{d\sigma} > 0$ , since  $\mu \in (0, \eta)$ .

Rearranging the condition derived from market clearing yields:  $(\mu - \eta + \sigma \frac{d\mu}{d\sigma}) \sum \frac{1}{C_i''(x_i)} = -(\mu - \eta + \sigma \frac{d\mu}{d\sigma}) (\sum \frac{1}{-B_i''(y_i)} + \sum \frac{1}{G_i''(x_i)}) - \eta (\sum \frac{1}{-B_i''(y_i)} + \sum \frac{1}{G_i''(x_i)})$

$$\mu - \eta = -\eta \frac{\sum \frac{1}{-B_i''(y_i)} + \sum \frac{1}{G_i''(x_i)}}{\sum \frac{1}{C_i''(x_i)} + \sum \frac{1}{-B_i''(y_i)} + \sum \frac{1}{G_i''(x_i)}} - \sigma \frac{d\mu}{d\sigma} \frac{\sum \frac{1}{-B_i''(y_i)} + \sum \frac{1}{G_i''(x_i)}}{\sum \frac{1}{C_i''(x_i)} + \sum \frac{1}{-B_i''(y_i)} + \sum \frac{1}{G_i''(x_i)}} + \frac{\sigma}{\eta} \frac{d\mu}{d\sigma}$$

using that  $\mu_0 := \lim_{\sigma \rightarrow 0} \sigma \frac{d\mu}{d\sigma} = 0$  we get:

$$\lim_{\sigma \rightarrow 0} \frac{\frac{dt_x}{d\sigma} + \frac{dt_y}{d\sigma} + \frac{dt_z}{d\sigma}}{\sum \frac{1}{C_i''(x_i)} + \sum \frac{1}{-B_i''(y_i)} + \sum \frac{1}{G_i''(x_i)}} = \frac{\frac{dt_x}{d\sigma}}{\sum \frac{1}{C_i''(x_i)} + \sum \frac{1}{-B_i''(y_i)} + \sum \frac{1}{G_i''(x_i)}}$$

Letting  $y_i^*(p)$  denote as before the energy demand function for country  $i$  in the absence of any mechanism, we have:

$$B_i'(y_i^*(p)) = p, \text{ so } B_i''(p) \frac{dy_i^*}{dp}(p) = 1$$

Letting  $c$  denote the supply function for coal for country  $i$  in the absence of any mechanism, we have:

$$C_i'(x_i^*(p)) = p, \text{ so } C_i''(p) \frac{dx_i^*}{dp}(p) = 1$$

Letting  $z_i^*$  denote the supply function for coal for country  $i$  in the absence of any mechanism, we have:

$$G_i'(z_i^*(p)) = p(1 - q), \text{ so } G_i''(p) \frac{dz_i^*}{dp}(p) = 1$$

With this we obtain:

$\sum \frac{1}{-B_i''(y_i)} = \frac{dy_i^*}{dp} = \frac{y_i^*(p)}{p} \frac{p}{y_i^*(p)} \frac{dy_i^*}{dp} = \frac{y_i^*(p)}{p} \epsilon_y$ , where  $\epsilon_D$  denotes the price elasticity of demand for energy.

$\sum \frac{1}{C_i''(x_i)} = \frac{dx_i^*}{dp} = \frac{x_i^*(p)}{p} \frac{p}{x_i^*(p)} \frac{dx_i^*}{dp} = \frac{x_i^*(p)}{p} \epsilon_x$ , where  $\epsilon_x$  denotes the price elasticity of supply of coal.

$\sum \frac{1}{G_i''(z_i)} = \frac{dz_i^*}{dp} = \frac{z_i^*(p)}{p} \frac{p}{z_i^*(p)} \frac{dz_i^*}{dp} = \frac{z_i^*(p)}{p} \epsilon_{S_G}$ , where  $\epsilon_{S_G}$  denotes the price elasticity of supply of renewable energy.

Using these identities we get:

$$\lim_{\sigma \rightarrow 0} \frac{\frac{dF_x}{d\sigma} + \frac{dF_y}{d\sigma} + \frac{dF_z}{d\sigma}}{\frac{dF_x}{d\sigma} + \frac{dF_y}{d\sigma} + \frac{dF_z}{d\sigma}} = \frac{\epsilon_z \frac{Z(0)}{Y(0)} + \epsilon_D}{\epsilon_z \frac{Z(0)}{Y(0)} + \epsilon_x \frac{X(0)}{Y(0)} + \epsilon_y} = \frac{\epsilon_z \frac{Z(0)}{Y(0)} + \epsilon_y}{\epsilon_{S_G} \frac{Z(0)}{Y(0)} + \epsilon_x \frac{X(0)}{Y(0)} + \epsilon_y}$$

where  $X(0), Y(0), Z(0)$  denote the aggregate quantities for  $F = 0$ .

From this we deduce:

$$\lim_{\sigma \rightarrow 0} \frac{\frac{dF_y}{d\sigma}}{\frac{dF_x}{d\sigma} + \frac{dF_y}{d\sigma} + \frac{dF_z}{d\sigma}} = \frac{\epsilon_y}{\epsilon_z + \epsilon_y} \frac{\epsilon_x \frac{X(0)}{Y(0)}}{\epsilon_z \frac{Z(0)}{Y(0)} + \epsilon_x \frac{X(0)}{Y(0)} + \epsilon_y}$$

$$\lim_{\sigma \rightarrow 0} \frac{\frac{dF_z}{d\sigma}}{\frac{dF_x}{d\sigma} + \frac{dF_y}{d\sigma} + \frac{dF_z}{d\sigma}} = \frac{\epsilon_z}{\epsilon_z + \epsilon_y} \frac{\epsilon_x \frac{X(0)}{Y(0)}}{\epsilon_z \frac{Z(0)}{Y(0)} + \epsilon_x \frac{X(0)}{Y(0)} + \epsilon_y}$$

But we also have:

$$\frac{dF_x}{dF} = \frac{\frac{dF_x}{d\sigma}}{\frac{dF}{d\sigma}} = \frac{\frac{dF_x}{d\sigma}}{\frac{dF_x}{d\sigma} + \frac{dF_y}{d\sigma} + \frac{dF_z}{d\sigma}} \text{ so the claimed result follows. } \square$$

## A.6 Proof of Lemma 5

*Proof.* Let  $(x_i(F), y_i(F), z_i(F))$  be the optimal allocation under the budget  $F$ . Given any  $F_1, F_2 > 0$ , suppose we have a budget of  $\phi F_1 + (1 - \phi)F_2$  with  $\phi \in (0, 1)$ . Set  $p$  equal to the status quo value in the absence of any mechanism and set  $(x_i, y_i, z_i) = (\phi x_i(F_1) + (1 - \phi)x_i(F_2), \phi y_i(F_1) + (1 - \phi)y_i(F_2), \phi z_i(F_1) + (1 -$



$\phi)z_i(F_2))$ . The transfers required for this allocation to satisfy the participation constraints is lower than  $\phi F_1 + (1 - \phi)F_2$  by the assumed convexity of  $C_i$  and  $G_i$  and the assumed concavity of  $B_i$ . Suppose we set the transfer so that the participation constraints are satisfied with equality. Then, denoting by  $\tilde{U}_i$  the values of  $U_i$  under the status quo, we have:

$$\begin{aligned} W &= \sum \tilde{U}_i - \eta(\phi x_i(F_1) + (1 - \phi)x_i(F_2)) \\ &= \sum \phi(\tilde{U}_i - \eta x_i(F_1)) + (1 - \phi)(\tilde{U}_i - \eta x_i(F_2)) \\ &= \phi W(F_1) + (1 - \phi)W(F_2) \end{aligned}$$

But since we have not even used up all our budget, this shows that we can do strictly better than this.  $\square$

## A.7 Proof of Lemma 6

*Proof.* The first order conditions for the  $x_i$  require that  $C'_i(x_i) = C'_j(x_j) \forall i, j$ . Thus once we stipulate a value for  $\sum_i x_i$ , all the  $x_i$  are determined via  $C'_i(x_i) = C'_j(x_j) \forall i, j$ . Then all the  $y_i, z_i$  are determined via  $\sum_i x_i = \sum_i y_i - \sum_i z_i$ ,  $B'_i(y_i) = B'_j(y_j) = G'_i(z_i) = G'_j(z_j)$ . Thus in particular, once  $\sum_i x_i$  is fixed, the welfare  $W$  is determined and the transfers only depend on  $p$ . Conversely,  $W$  determines all the  $x_i, y_i, z_i$  and the  $F_x$  then only depends on  $p$ .

Specifically,  $F_x$  corresponds to  $p$  via  $F_x = \sum_i F_{ix} = \sum_i p x_i^*(p) - C_i(x_i^*(p)) - (p x_i - C_i(x_i))$

where, as usual, we denote by  $x_i^*(p)$   $i$ 's coal supply in the absence of the mechanism.

$$\frac{dF_x}{dp} = \sum_i x_i^* + (p - \alpha) \sum_i \frac{dx_i^*}{dp} - \sum_i \frac{dx_i^*}{dp} C'_i(x_i^*) = \sum_i x_i^*, \text{ so } \frac{dp}{dF_x} = \frac{1}{\sum_i x_i^*}$$

Similarly, we get:

$$\begin{aligned} \frac{d \sum_i F_{iy}}{dF_x} &= \frac{d \sum_i F_{iy}}{dp} \frac{dp}{dF_x} = \frac{-\sum_i y_i^*}{\sum_k x_k^*}, \frac{d \sum_i F_{iz}}{dF_x} = \frac{\sum_i z_i^*}{\sum_i x_i^*}, \text{ so } \frac{dF}{dF_x} = \frac{\sum_i x_i^* + \sum_i z_i^* - \sum_i y_i^*}{\sum_k x_k^*} \\ \frac{d^2 F}{dF_x^2} &= \frac{d}{dp} \frac{\sum_i x_i^* + z_i^* - y_i^*}{\sum_k x_k^*} \frac{dp}{dF_x} = \frac{-\sum_k \frac{dx_k^*}{dp} \sum_i (x_i^* + z_i^* - y_i^*) + \sum_k x_k^* (\sum_i \frac{dx_i^*}{dp} + \frac{dz_i^*}{dp} - \frac{dy_i^*}{dp})}{(\sum_k x_k^*)^2} \frac{dp}{dF_x} = \\ &= \frac{\sum_k \frac{dx_k^*}{dp} \sum_i (y_i^* - z_i^*) + \sum_k x_k^* (\frac{dz_i^*}{dp} - \frac{dy_i^*}{dp})}{(\sum_i x_i^*)^3} \end{aligned}$$

At the optimal mechanism, each country  $i$  will be paid to lower their energy use relative to what it would individually choose were it to ignore the mechanism. Hence we must have  $y_i \leq y_i^* \forall i$ . Similarly, we must have  $z_i \geq z_i^* \forall i$ . Moreover, market clearing implies that  $\sum_i y_i - z_i \geq 0$ , so  $\sum_i (y_i^* - z_i^*) \geq 0$ .

Since  $\frac{dz_i^*}{dp} \geq 0$ ,  $\frac{dx_i^*}{dp} \geq 0 \forall i$  by the law of supply and  $\frac{dy_i^*}{dp} \leq 0$  by the law of demand, it follows that  $\frac{d^2 F}{dF_x^2} \geq 0$   $\square$

## A.8 Proof of Lemma 8

*Proof.* Denote by  $u_{iy}$  the net payoff that country  $i$  gets from setting a tax rate  $\tau_{iy} > 0$  on energy use relative to setting  $\tau_{iy} = 0$  is.

$$u_{iy} = B_i(y_i^*(p + \tau_{iy})) - p y_i^*(p + \tau_{iy}) - (B_i(y_i^*(p)) - p y_i^*(p)) + \theta_{iy}(\tau_{iy} - \tilde{\tau}_{iy})$$

$$\frac{du_{iy}}{d\tau_{iy}} = (B'_i(y_i^*(p + \tau_{iy})) - p)y_i^{*'}(p + \tau_{iy}) + \theta_{iy}$$

Given that  $y_i^*$  is characterized by the optimality condition

$$B'_i(y_i) = p + \tau_{iy}$$

we have

$$y_i^{*'}(p + \tau_{iy}) = \frac{1}{B''_i(y_i^*(p + \tau_{iy}))}$$

Hence:

$$\frac{du_{iy}}{d\tau_{iy}} = (B'_i(y_i^*(p + \tau_{iy})) - p) \frac{1}{B''_i(y_i^*(p + \tau_{iy}))} + \theta_{iy}$$

$$\frac{d^2u_{iy}}{d\tau_{iy}^2} = B''_i(y_i^*(p + \tau_{iy})) \frac{1}{(B''_i(y_i^*(p + \tau_{iy})))^2} - (B'_i(y_i^*(p + \tau_{iy})) - p) B'''_i(y_i^*(p + \tau_{iy})) \frac{1}{(B''_i(y_i^*(p + \tau_{iy})))^3}$$

Given that  $B''_i(y) < 0 \forall y$ ,  $\frac{d^2u_{iy}}{d\tau_{iy}^2} < 0$  is equivalent to:

$$1 - (B'_i(y_i^*(p + \tau_{iy})) - p) B'''_i(y_i^*(p + \tau_{iy})) \frac{1}{(B''_i(y_i^*(p + \tau_{iy})))^2} > 0$$

Using the definition of the elasticities, this can be rewritten as:

$$1 - \frac{B'_i(y_i^*(p + \tau_{iy})) - p}{B'_i(y_i^*(p + \tau_{iy}))} \frac{\epsilon_{B'_i}(y_i^*(p + \tau_{iy}))}{\epsilon_{B'_i}(y_i^*(p + \tau_{iy}))} > 0$$

Now the sufficiency of conditions 2) follows from the fact that  $\lim_{F \rightarrow 0} \frac{B'_i(y) - p}{B'_i(y)} = 0$ .

To show the sufficiency of condition 3), we note that by assumption 1 we have  $\epsilon_{B'_i}(y) < 0$  and  $B''_i < 0$ . this implies that if  $B'''_i < 0$  then equation 1 is satisfied.

For the supply side, an analogous computation yields that the following is a sufficient condition:

$$1 - \frac{C'_i(x) - p}{C'_i(x)} \frac{\epsilon_{C'_i}(x)}{\epsilon_{C'_i}(x)} > 0$$

By assumption 1, we have  $\epsilon_{C'_i}(x) > 0$ . Moreover, we also have  $\frac{C'_i(x) - p}{C'_i(x)} < 0$ . Hence it is sufficient that  $\epsilon_{C''_i}(x) > 0$ , which is equivalent to  $C'''_i > 0$  on the relevant domain.

For the substitute side, an analogous computation yields that the following is a sufficient condition:

$$1 - \left( \frac{G'_i(z) - p}{G'_i(z)} \right) \frac{\epsilon_{G'_i}(z)}{\epsilon_{G'_i}(z)} > 0$$

By assumption 1, we have  $\epsilon_{G'_i}(z) > 0$ . Moreover, we also have  $\frac{G'_i(z) - p}{G'_i(z)} > 0$ . Hence it is sufficient that  $\epsilon_{G''_i}(z) < 0$ , which is equivalent to  $C'''_i < 0$  on the relevant domain.  $\square$

## B Proofs for section 6

### B.1 Proof of Lemma 9

*Proof.* Country  $i$ 's problem on the supply side is to choose a nondecreasing sequence of cumulative extraction  $(x_{it})_{t=1,\dots,T}$  so as to maximize the discounted payoffs:

$$\begin{aligned}\frac{dU_i}{dx_{it}} &= \frac{1}{(1+r)^t} (p_t - \frac{1}{1+r} p_{t+1} - c'_{it}(x_{it}) + \frac{1}{1+r} c'_{it+1}(x_{it})) \\ \frac{d^2U_i}{dx_{it}^2} &= \frac{1}{(1+r)^t} (-c''_{it}(x_{it}) + \frac{1}{1+r} c''_{it+1}(x_{it}))\end{aligned}$$

By assumption 4 we can from this conclude that  $\frac{d^2U_i}{dx_{it}^2} < 0 \forall t \in \{1, \dots, T - 1\}$ .

1). Since  $\frac{d^2U_i}{dx_{it}dx_{jt}} \forall i \neq j$ , we can conclude that the function  $U_i$  is concave in  $(x_{it})_{t \in \{1, \dots, T\}}$  on  $R^T$ . It follows in particular that it is concave on the subspace of  $R^T$  defined by the condition that  $(x_{it})_{t=1, \dots, T}$  be non-decreasing. Thus there exists a unique utility maximizing extraction path for country  $i$ .  $\square$

### B.2 Proof of Lemma 10

*Proof.* Formally, we can reinterpret the Walrasian World model of this paper as a classical Walrasian model with a numeraire good by viewing each country as being comprised of three parts:

- 1) a representative consumer with utility  $B_i$  that is quasilinear in the numeraire good<sup>15</sup>
- 2) a firm able to use the numeraire good to extract coal and turn it into energy
- 3) a firm able to use the numeraire good to produce renewable energy energy.

Now existence of a market equilibrium follows from proposition 17.BB.2 of Mas-Colell et. al 1995 (page 634), using assumption 3.

To establish uniqueness, we note that the first welfare theorem requires that any Walrasian equilibrium is Pareto optimal. However, given that all utilities are quasilinear in the same numeraire, this means that every Walrasian equilibrium must maximize the sum of all the utilities. But since all the utilities are concave functions of the allocations, this uniquely determines the allocation up to transfers in the numeraire good between consumers.  $\square$

### B.3 Proof of Surjectivity Lemma, Dynamic Version with full Commitment 11

*Proof.* For the demand side and the substitute side each country's optimization problem is separable across periods. Thus the claims follow by the same arguments as in the proof of Lemma 1.

<sup>15</sup>Of course, for the application of the model, the appropriate interpretation of the model is that  $B_i$  captures both the utility that consumers in  $i$  derive from energy consumption and the profits the benefits generated from energy for the production of other goods. The narrower interpretation adopted in this proof simply helps to import the results from standard Walrasian theory.

On the supply side, however, there is a complication to be dealt with: The constraint that the extraction sequence be non-decreasing. Country  $i$ 's problem on the supply side is to choose a nondecreasing sequence of cumulative extraction  $(x_{it})_{t=1,\dots,T}$  so as to maximize the discounted payoffs:

$$\sum_{t=1}^T \frac{1}{(1+r)^t} (p_t(x_{it} - x_{it-1}) + f_{ixt}(x_{it}) - (c_{it}(x_{it}) - c_{it}(x_{it-1})))$$

Let us rewrite this by grouping together the terms involving a given  $x_{it}$ :

$$\sum_{t=1}^{T-1} \frac{1}{(1+r)^t} \left( \left( p_t - \frac{1}{1+r} p_{t+1} \right) x_{it} + f_{ixt}(x_{it}) - c_{it}(x_{it}) + \frac{1}{1+r} c_{it+1}(x_{it}) \right) + \frac{1}{(1+r)^T} (p_T x_{iT} + f_{ixT}(x_{iT}) - c_{iT}(x_{iT})) \quad (11)$$

Suppose the global institution were to offer the following reward payment schemes: In period  $t$  pay the amount  $F_{ixt}$  to  $i$  if  $i$  chooses  $x_{it}$  and 0 otherwise, where  $F_{ixt}$  is defined as follows

$$F_{ixt} := \left( p_t - \frac{1}{1+r} p_{t+1} \right) (x_{it}^*(p) - x_{it}) + (c_{it}(x_{it}) - c_{it}(x_{it}^*(p))) + \frac{1}{1+r} (c_{it+1}(x_{it}^*(p)) - c_{it+1}(x_{it})) \text{ for } t \in \{1, \dots, T-1\}$$

$$F_{ixT} := p_T (x_{iT}^*(p) - x_{iT}) + (c_{iT}(x_{iT}) - c_{iT}(x_{iT}^*(p)))$$

Recall that by definition 6,  $(x_{it}^*(p))_{t \in \{1, \dots, T\}}$  denotes the unique utility maximizing extraction path for country  $i$ . By equation 11,  $x_{it}^*(p)$  is the unique value that maximizes the term of the profit in which it appears. Thus we have  $F_{ixt} \geq 0 \forall i, t$ . Thereby we have established that the reward payment schemes we have defined are positive.

Consider the relaxed problem for country  $i$  where it ignores the constraint that  $(x_{it})_{t=1,\dots,T}$  has to be non-decreasing. From equation 11 it follows that  $(x_{it})_{t=1,\dots,T}$  is an optimum for  $i$  under the reward payment schemes defined above. But by assumption, this  $(x_{it})_{t=1,\dots,T}$  is non-decreasing, so it is also an optimum of the actual problem that country  $i$  faces.

Now we need to show that every reward payment scheme inducing country  $i$  to choose  $(x_{it})_{t=1,\dots,T}$  must end up paying to  $i$  at least the amount  $F_{ix} := \sum_{t=1,\dots,T} \frac{1}{(1+r)^t} F_{ixt}$  in discounted money on rewarding supply reduction. But for this we simply note that country  $i$  always has the option of ignoring all of the supply side reward payment schemes altogether. From this and equation 11 the claimed result follows.

The fact that each country always has the option of ignoring the reward payment schemes also establishes that the global institution does not increase the minimal required discounted value of reward payment schemes on the demand or substitute side by restricting itself to using period-by-period reward payment schemes.

Last, let us show that positive affine linear reward payment schemes are sufficient to achieve any given supply side allocation with minimal transfers.

Let  $V_{ixt}(x_{it-1}, (p_s)_{s \geq t}, (f_{ixs})_{s \geq t})$  denote the maximal discounted profit that country  $i$  can reap from coal extraction from period  $t$  onwards, given that its cumulative extraction before that is  $x_{it-1}$  and given the continuation world market price path  $(p_s)_{s \geq t}$  and given the continuation path  $(f_{ixs})_{s \geq t}$  of reward payment schemes.

The individual rationality condition for  $i$  to accept  $x_{iT}$  conditional on having accepted  $(x_{i1}, \dots, x_{iT-1})$  is:

$$x_{it} = \operatorname{argmax}_x p_t(x - x_{it-1}) + f_{ixt}(x - x_{it-1}) - (c_{it}(x) - c_{it}(x_{it-1})) + \frac{1}{1+r} V_{ixt+1}(x, (p_s)_{s \geq t+1}, (f_{ixs})_{s \geq t+1})$$

The corresponding first order condition is:

$$p_t - \theta_{ixt} - c'_{it}(x_{it}) + \frac{1}{1+r} \frac{\partial V_{ixt+1}}{\partial x_{it}} = 0$$

Using the envelope theorem, we compute:

$$\frac{\partial V_{ixt+1}}{\partial x_{it}} = -p_{t+1} + c'_{it+1}(x_{it})$$

Hence the first order conditions are:

$$p_t - \theta_{ixt} - c'_{it}(x_{it}) + \frac{1}{1+r} (-p_{t+1} + c'_{it+1}(x_{it})) = 0 \text{ for } t \leq T - 1$$

$$p_T - \theta_{ixT} - c'_{iT}(x_{iT}) = 0 \text{ for } t = T$$

Thus in each period  $t$  there is a unique  $\theta_{it}$  such that  $x_{it}$  satisfies the first order condition in period  $t$ . What remains to be checked is that the second order conditions hold. For period  $T$ , this is true by since  $c_{iT}$  is convex by assumption. For periods  $t < T$ , the second order condition for a maximum is:

$$-c''_{it}(x_{it}) + \frac{1}{1+r} c''_{it+1}(x_{it}) \leq 0$$

But by assumption 4, this condition holds. In fact, assumption 4 implies that the function  $x \mapsto p_t(x - x_{it-1}) + f_{ixt}(x - x_{it-1}) - (c_{it}(x) - c_{it}(x_{it-1})) + \frac{1}{1+r} V_{ixt+1}(x, (p_s)_{s \geq t+1}, (f_{ixs})_{s \geq t+1})$  is concave and thus the first order condition gives the unique maximum.

Adjusting the reference levels of the positive affine linear reward payment schemes simply amounts to varying the transfer that a country gets conditional on accepting the reward payment scheme (instead of ignoring it and maximizing its payoff without the rewards payments instead). Hence there are unique reference levels such that the the global institution ends up paying exactly the minimal transfers.  $\square$

## B.4 Proof of Price Preservation Lemma, Dynamic Version with Full Commitment 12

*Proof.* By the Surjectivity Lemma 11, we can view the global institution as if it was choosing a world market price path  $(p_t)_{t \in \{1, \dots, T\}}$  and an allocation  $(x_{it}, y_{it}, z_{it})_{i \in I, t \in \{1, \dots, T\}}$ . Also by the Surjectivity Lemma the minimal discounted sum of transfers that the global institution has to pay for rewarding supply reduction is  $F_{ix} := \sum_{t=1, \dots, T} \frac{1}{(1+r)^t} F_{ixt}$  where  $F_{ixt}$  is defined as follows:

$$F_{ixt} := (p_t - \frac{1}{1+r} p_{t+1})(x_{it}^*(p) - x_{it}) + (c_{it}(x_{it}) - c_{it}(x_{it}^*(p))) + \frac{1}{1+r} (c_{it+1}(x_{it}^*(p)) - c_{it+1}(x_{it})) \text{ for } t \in \{1, \dots, T - 1\}$$

$$F_{ixT} := p_T(x_{iT}^*(p) - x_{iT}) + (c_{iT}(x_{iT}) - c_{iT}(x_{iT}^*(p)))$$

Moreover, the minimal discounted sum of transfers that the global institution has to pay for inducing the demand reduction and the substitute expansion are  $F_{iy} := \sum_{t=1, \dots, T} \frac{1}{(1+r)^t} F_{iyt}$  and  $F_{iz} := \sum_{t=1, \dots, T} \frac{1}{(1+r)^t} F_{izt}$  with the following definitions:

$$F_{iyt} := \sup_y b_{it}(y) - p_t y - (b_{it}(y_{it}) - p_t y_{it})$$

$$F_{izt} := \sup_z p_t z - g_{it}(z) - (p_t z_{it} - g_{it}(z_{it}))$$

Let us use the envelope theorem to compute how the minimal value of discounted aggregate transfer is affected by the world market prices:

$$\begin{aligned}
\frac{dF_{ix}}{dp_t} &= \frac{1}{(1+r)^t} (x_{it}^*(p) - x_{it}) - \frac{1}{(1+r)^{t-1}} \frac{1}{1+r} (x_{it-1}^*(p) - x_{it-1}) \\
\frac{dF_{ix}}{dp_t} &= \frac{1}{(1+r)^t} (x_{it}^*(p) - x_{it-1}^*(p) - (x_{it} - x_{it-1})) \\
\frac{d}{dp_t} (\sum_{i \in I, t \in \{1, \dots, T\}} F_{ix} + F_{iy} + F_{iz}) &= \sum_{i \in I} \frac{1}{(1+r)^t} (x_{it}^*(p) - x_{it-1}^*(p) - (x_{it} - \\
&x_{it-1}) - y_{it}^*(p) + y_{it} + z_{it}^*(p) - z_{it}) \\
\text{By market clearing in period } t \text{ we deduce from this:} \\
\frac{d}{dp_t} (\sum_{i \in I} F_{ix} + F_{iy} + F_{iz}) &= \sum_{i \in I, t \in \{1, \dots, T\}} \frac{1}{(1+r)^t} (x_{it}^*(p) - x_{it-1}^*(p) - y_{it}^*(p) + \\
&z_{it}^*(p))
\end{aligned}$$

Optimality of the world market price path for the problem of minimising discounted aggregate transfers implies that this expression has to be 0 for each  $t \in \{1, \dots, T\}$ . But this set of conditions is equivalent to the statement that  $((x_{it}, y_{it}, z_{it})_{i \in I, t \in \{1, \dots, T\}}, p)$  be a market equilibrium in the absence of any reward payment schemes. By Lemma 10 there is a unique such market equilibrium.  $\square$

## B.5 Proof of Constrained Efficiency Lemma, Dynamic Version 13

*Proof.* The global institution's objective is to maximize global welfare:

$$W = \sum_i U_i((x_{it}, y_{it}, z_{it})_{t=1, \dots, T}) - \eta(\sum_{i \in I} x_{i1}, \dots, \sum x_{iT})$$

under the constraint that the total amount of required transfers does not exceed the available intertemporal budget:

$$\sum_i \sum_{t=1, \dots, T} \frac{1}{(1+r)^t} F_{it} \leq F$$

Consider a fixed value for the total climate change damages,  $\eta(\sum_{i \in I} x_{i1}, \dots, \sum x_{iT}) = \tilde{\eta}$ . How should the global institution optimally choose the allocation  $(x_{it}, y_{it}, z_{it})_{t=1, \dots, T}$  under the constraint that  $\eta(\sum_{i \in I} x_{i1}, \dots, \sum x_{iT}) = \tilde{\eta}$ ? There are two considerations: The institution should try to reduce the transfers it needs to pay. For this it is best to maximize the countries' aggregate utility,  $\sum_i U_i((x_{it}, y_{it}, z_{it})_{t=1, \dots, T})$ , since the required transfers are determined so as to make up for the utility loss that countries incur relative to ignoring the reward payment schemes and optimizing given the market price vector  $p$ . The global institution also cares intrinsically about the countries' utilities. Hence the two considerations perfectly align and the global institution should choose the allocation so as to maximize  $\sum_i U_i((x_{it}, y_{it}, z_{it})_{t=1, \dots, T})$  under the constraint that  $\eta(\sum_{i \in I} x_{i1}, \dots, \sum x_{iT}) = \tilde{\eta}$ .  $\square$

## B.6 Proof of Monotone Mitigation Lemma 14

*Proof.* By the Constrained Efficiency Lemma 13 we know that the allocation at the optimal mechanism must maximize global welfare under the constraint that the eventual cumulative coal extraction be some fixed value  $X$ . For  $X = \sum_i x_{iT}^*(p)$  the optimum is  $x_{it} = x_{it}^*(p)$  by the first Welfare Theorem.

From now on we shall reason about how to optimally choose the  $x_t := \sum_i x_{it}$ , it being understood that the aggregate quantities are split up across countries so as to maximize global welfare. This is justified by the Constrained Efficiency Lemma 13. We use the definition of aggregate costs:  $c_t(x_t) :=$

$\min_{(x_{it})_i} \sum c_{it}(x_{it})$ . For better readability I will now assume that there is no substitute (i.e. no renewable energy). It is straightforward to generalise the proof that I will now give.

Let us denote  $(x_t(X), y_t(X), z_t(X))_{t \in \{1, \dots, T\}}$  the optimal allocation under the constraint that  $x_T = X$ . We shall now prove that  $0 < \frac{dx_1}{dX} < \dots < \frac{dx_{T-1}}{dX} < 1$ . For this, we note that by optimality we have:

$$\frac{\partial W}{\partial x_t} = 0 \forall t \leq T - 1$$

Differentiating this with respect to  $X$  and noting that  $x_t$  only affects coal combustion in period  $t$  and  $t + 1$  yields :

$$\frac{\partial^2 W}{\partial x_t \partial x_{t-1}} \frac{dx_{t-1}}{dX} + \frac{\partial^2 W}{\partial x_t^2} \frac{dx_t}{dX} + \frac{\partial^2 W}{\partial x_t \partial x_{t+1}} \frac{dx_{t+1}}{dX} = 0 \forall t \in \{2, \dots, T - 2\}$$

$$\frac{\partial^2 W}{\partial x_1^2} \frac{dx_1}{dX} + \frac{\partial^2 W}{\partial x_1 \partial x_2} \frac{dx_2}{dX} = 0$$

$$\frac{\partial^2 W}{\partial x_{T-1}^2} \frac{dx_{T-1}}{dX} + \frac{\partial^2 W}{\partial x_{T-1} \partial X} = 0$$

Rearranging yields:

$$\frac{dx_{t+1}}{dX} = -\frac{\frac{\partial^2 W}{\partial x_t \partial x_{t-1}}}{\frac{\partial^2 W}{\partial x_t \partial x_{t+1}}} \frac{dx_{t-1}}{dX} - \frac{\frac{\partial^2 W}{\partial x_t^2}}{\frac{\partial^2 W}{\partial x_t \partial x_{t+1}}} \frac{dx_t}{dX} \forall t \in \{2, \dots, T - 2\} \quad (12)$$

$$\frac{dx_2}{dX} = \frac{-\frac{\partial^2 W}{\partial x_1^2}}{\frac{\partial^2 W}{\partial x_1 \partial x_2}} \frac{dx_1}{dX} \quad (13)$$

$$\frac{dx_{T-1}}{dX} = \frac{\frac{\partial^2 W}{\partial x_{T-1} \partial X}}{-\frac{\partial^2 W}{\partial x_{T-1}^2}} \quad (14)$$

We have:

$$\frac{\partial W}{\partial x_t} = \frac{1}{(1+r)^t} (b'_t(x_t - x_{t-1}) - \frac{1}{1+r} b'_{t+1}(x_{t+1} - x_t) - c'_t(x_t) + \frac{1}{1+r} c'_{t+1}(x_t))$$

$$\frac{\partial^2 W}{\partial x_t^2} = \frac{1}{(1+r)^t} (b''_t(x_t - x_{t-1}) + \frac{1}{1+r} b''_{t+1}(x_{t+1} - x_t) - c''_t(x_t) + \frac{1}{1+r} c''_{t+1}(x_t))$$

$$\frac{\partial^2 W}{\partial x_t \partial x_{t+1}} = \frac{1}{(1+r)^t} (-\frac{1}{1+r} b''_{t+1}(x_{t+1} - x_t))$$

Since by assumption 4, we have  $c''_t(x_t) - \frac{1}{1+r} c''_{t+1}(x_t) > 0$ , so we obtain in particular:

$$0 < \frac{\partial^2 W}{\partial x_t \partial x_{t+1}} < -\frac{\partial^2 W}{\partial x_t^2} \quad (15)$$

From 15 and 14 we deduce that  $0 < \frac{dx_{T-1}}{dX} < 1$ .

Let us now assume that  $0 < \frac{dx_1}{dX}$ . From this we will now deduce that  $\frac{dx_t}{dX} < \frac{dx_{t+1}}{dX} \forall t \in \{1, \dots, T - 2\}$  by mathematical induction on  $t$ .

Now we prove the induction step. For this, suppose that the claim holds for  $t$ , i.e suppose that  $\frac{dx_t}{dX} > \frac{dx_{t-1}}{dX} > 0$ . By 12, we have:

$$\frac{dx_{t+1}}{dX} = -\frac{\frac{\partial^2 W}{\partial x_t \partial x_{t-1}}}{\frac{\partial^2 W}{\partial x_t \partial x_{t+1}}} \frac{dx_{t-1}}{dX} - \frac{\frac{\partial^2 W}{\partial x_t^2}}{\frac{\partial^2 W}{\partial x_t \partial x_{t+1}}} \frac{dx_t}{dX}$$

Now using the induction hypothesis according to which  $\frac{dx_t}{dX} > \frac{dx_{t-1}}{dX}$  and the fact  $\frac{\partial^2 W}{\partial x_t \partial x_{t+1}} = \frac{1}{(1+r)^t} (-\frac{1}{1+r} b''_{t+1}(x_{t+1}-x_t)) > 0$  and  $\frac{\partial^2 W}{\partial x_t \partial x_{t-1}} = \frac{1}{(1+r)^t} (-b''_t(x_t - x_{t-1})) > 0$ , we deduce:

$$\frac{dx_{t+1}}{dX} > -\frac{\frac{\partial^2 W}{\partial x_t \partial x_{t-1}}}{\frac{\partial^2 W}{\partial x_t \partial x_{t+1}}} \frac{dx_t}{dX} - \frac{\frac{\partial^2 W}{\partial x_t^2}}{\frac{\partial^2 W}{\partial x_t \partial x_{t+1}}} \frac{dx_t}{dX} = \frac{-\frac{\partial^2 W}{\partial x_t^2} - \frac{\partial^2 W}{\partial x_t \partial x_{t-1}}}{\frac{\partial^2 W}{\partial x_t \partial x_{t+1}}} \frac{dx_t}{dX}$$

Substituting in yields:

$$\frac{dx_{t+1}}{dX} > \frac{-\frac{1}{1+r} b''_{t+1}(x_{t+1}-x_t) + c''_t(x_t) - \frac{1}{1+r} c''_{t+1}(x_t)}{-\frac{1}{1+r} b''_{t+1}(x_{t+1}-x_t)} \frac{dx_t}{dX}$$

But by assumption 4, we have  $c''_t(x_t) - \frac{1}{1+r} c''_{t+1}(x_t) > 0$ , so this implies that that  $\frac{dx_{t+1}}{dX} > \frac{dx_t}{dX}$

Thus we have shown that if we assume that if  $0 < \frac{dx_1}{dX}$  then by induction it follows that  $\frac{dx_t}{dX} < \frac{dx_{t+1}}{dX} \forall t \in \{1, \dots, T-2\}$ .

Now we will show that  $0 \geq \frac{dx_1}{dX}$  would lead to a contradiction. In fact, the same inductive argument just provided would in this case imply that  $\frac{dx_{T-1}}{dX} \leq \dots \leq \frac{dx_1}{dX} \leq 0$ . But we already showed above that  $0 < \frac{dx_{T-1}}{dx_T} < 1$ .

Thus we have shown that  $0 < \frac{dx_1}{dX} < \dots < \frac{dx_1}{dX} < 1$ . In particular, we have  $x_t(X) < \sum_i x_{it}^*(p) \forall t \forall W$  as long as  $X < \sum_i x_{iT}^*(p)$ .  $\square$

## B.7 Proof of the Time Inconsistency Corollary 6

*Proof.* Consider the situation at time  $T$ , assuming that to all countries' surprise at that time the global institution announces the new reoptimized mechanism. By the first part of the Monotone Mitigation Lemma 14, until time  $T$  less coal has been extracted than what would have been extracted in the absence of any mechanism. Let  $p_T^\#$  denote the price that would occur if for the time period  $T$  no more mechanism was used and given that the cumulative coal extraction by the end of period  $t-1$  was  $(x_{iT-1})_{i \in I}$ . Let  $\hat{p}_T$  denote the price arising in period  $T$  if there was not any mechanism at any period. We have  $p_T^\# < \hat{p}_T$ , since otherwise there would be an excess coal supply under  $p_T^\#$ .

By the Dynamic Price Preservation Lemma 12,  $\hat{p}_T$  is also the price arising under the optimal mechanism under full commitment. Thus if the global institution reneges on the announced mechanism at period  $T$  and reoptimizes, it will spend less on rewarding supply reduction.  $\square$

## C Tax-based implementation in the dynamic model

Let us now study in the dynamic model how the global institution can best reward countries on the basis of their tax/subsidy rates. As throughout this paper, let us assume that the global institution has an exogenous intertemporal budget and that it can fully commit. Of course it is impossible to achieve more global welfare than using quantity-based reward payment schemes as studied in section 6, given any intertemporal budget at the global institution's disposal. Let us now consider this optimal allocation that the global institution can achieve



for a given budget using quantity-based reward payment schemes and let us characterize the corresponding tax and subsidy rates.

Let  $\tau_{ixt}$  denote the tax that country  $i$  charges coal producers per unit of coal extraction,  $\tau_{iyt}$  the tax that it charges coal users per unit of coal combustion and  $\tau_{izt}$  the subsidy that it pays renewable energy producers per unit of renewable energy produced. From the Constrained Efficiency Lemma 13 we obtain the following 3 conditions:

**Condition 1.**  $\tau_{ixt} = \tau_{jxt} \forall i, j, \tau_{iyt} = \tau_{jyt} \forall i, j, \tau_{izt} = \tau_{jzt} \forall i, j$

**Condition 2.**  $\tau_{izt} = \tau_{iyt} \forall i, t$

**Condition 3.**  $\tau_{ixt+1} + \tau_{iyt+1} = (1 + r)(\tau_{ixt} + \tau_{iyt}) \forall i, j, t$

Condition 2) states that the marginal cost of reducing fossil fuel demand via reduced energy use has to equal the marginal cost of reducing it via expanded renewable energy production.

Condition 3) states that the total wedge between the consumer and the producer price for coal has to increase at the rate  $r$  at which countries discount the future.

From the Price Preservation Lemma we obtain:

**Condition 4.** Market clearing under preserved price path:

$$\sum_{i \in I} x_{it}^*(p - \tau_{ix}) + z_{it}^*(p_t + \tau_{izt}) - y_{it}^*(p_t + \tau_{iyt}) = 0 \forall i, t$$

where  $p$  denotes the price path in the absence of any reward payment schemes and  $\tau_{ix}$  denotes the path of tax rates on coal extraction.

The conditions 1), 2), 3) leave  $T$  degrees of freedom, corresponding to the split of the total wedge  $\tau_{ixt} + \tau_{iyt}$  between extraction-level tax and combustion level tax on coal. These  $T$  degrees of freedom are pinned down by the  $T$  requirements provided by condition 4).

In the model, given our assumption of complete information, the global institution could simply compute these paths of tax rates and make countries take-it-or-leave-it offers for exactly adopting these paths. Given our assumption that the global institution only has an intertemporal budget constraint (i.e. that it can freely save from and borrow against the exogenous flow of funding that it receives), the global institution could offer these conditional reward payments in the last period.

With that, the global institution can clearly implement the above paths of tax rates. However, in practice it would presumably be much better for the global institution to make reward payments in each period. Ideally, these reward payments should only depend on the choices made by the country in that period, thus making the link between the choice and the reward payment more transparent.

For the case of quantity-based reward payment schemes, we saw in section 6 that nothing is lost by restricting the global institution to such period-by-period reward payment schemes. For the demand side and the substitute side, these results trivially also hold for the tax-based reward payment schemes.

However, for the supply side a potential complication arises if one restricts the global institution to using period-by-period tax-based reward payment schemes: The effect of a country's period  $t$  tax rate on coal extraction on its utility depends on its extraction tax rates in previous periods. For example, if a country subsidizes coal extraction in early periods then it might not mind taxing coal extraction at a high rate in later periods since it might not extract any coal any more in any case.

This complication suggests that the global institution should in each period only pay supply side reward payments to countries whose cumulative coal extraction by the end of the preceding period did not exceed the cumulative coal extraction that it would have had in the absence of any reward payment schemes. In the model, it turns out that with this qualification nothing is lost by restricting the global institution to paying on the supply side each period each country solely on the basis of the country's tax rate in that period. A formal proof of this result is available upon request.

Now in practice the global institution has only imperfect information about the countries' counterfactual extraction paths in the absence of any reward payment schemes. In fact, this predicament has been the main motivation for considering alternatives to the quantity-based reward payment schemes in the first place, as I explained in section 5. However, the approach I have proposed in the previous paragraphs does not require the estimates of the counterfactual extraction paths in the absence of any reward payment schemes to be very precise. In fact, these estimates only need to be low enough to avoid creating perverse incentives for countries to subsidize coal extraction early on to reap reward payments for taxing it later on and high enough to exceed the cumulative extraction paths that countries get when they implement the proposed tax rates. Between these two failure modes there is a margin of error.

## D Mathematica notebooks

A Mathematica notebook for the numerical computations under constant elasticity specifications can be downloaded [here](#).

A Mathematica notebook computing the spending paths on the three approaches by the global institution at the optimal mechanism with full commitment and no borrowing or saving constraints can be downloaded [here](#).

A Mathematica notebook computing the surfaces shown in section E about the loss from misallocation can be downloaded [here](#).

## E Robustness checks about the loss from misallocation

In section 4 I showed for a particular combination of elasticity estimates how welfare depends on the budget split. One takeaway was that the loss from misallocation is relatively small: as long as each of the three approaches gets at

least 50% its optimal proportion of the budget, welfare losses are at most 10%. It suggests that it might not be so important to get the budget split exactly right and thus weighs in favor of decentralized funding mechanisms that have no guarantee for allocative efficiency but that create strong participation incentives by giving participating countries the opportunity to influence the allocation of funding across the different approaches to curbing fossil fuels.

We should expect countries' allocation decisions in such mechanisms to be guided by a mixture of concern for global welfare and their own payoffs. Large fossil fuel exporters will strongly prefer money to go to the supply reduction approach since this will raise fossil fuel world market prices. Fossil fuel importers, on the other hand, will prefer money to go to the demand reduction and substitute expansion approaches. Those countries primarily concerned about climate change will prefer money to go at the margin to whatever of the three approaches is underfunded relative to the others. Thus we should expect the overall allocation to be somewhat responsive to what actually turns out to be good for global welfare. Based on this, I now assume for concreteness that for any given overall budget each of the three approaches gets at least half its optimal proportion.

Under this constraint, the worst outcome in terms of global welfare occurs when two of the three approaches each get only half their optimal proportion of the budget, with the third approach getting the rest. I refer to the three corresponding cases as "supply-side-heavy", "demand-side-heavy" and "substitute-side-heavy". I plot below the proportion of welfare realised under these three cases relative to the welfare that would be realised if the budget was split optimally across the three approaches. For these numerical simulations I assume constant-elasticity specifications for coal supply, energy demand and renewable energy supply. Throughout I assume that the overall budget is small. It turns out that all the results change little with the size of the budget. The plots explore the entire range of elasticity estimates that I have found in the literature, as I detail in the following subsections.

## E.1 Estimates of long run price elasticities of demand for energy

Espey and Espey (2004) carried out a meta-analysis about residential electricity demand. of price and income elasticity estimates from 36 studies published over the period 1947 to 1997. The 125 estimates of long-run price elasticity fell in the range from  $-2.25$  to  $-0.04$  with a mean of  $-0.85$ . All the more recent studies that I have seen have estimates falling in this range<sup>16</sup> I thus consider the range  $-2.25$  to  $-0.04$  in the plots shown below.

---

<sup>16</sup>E.g. Burke & Abayasekara (2018) find  $-1$ .

## E.2 Estimates of price elasticities of supply of renewable energy

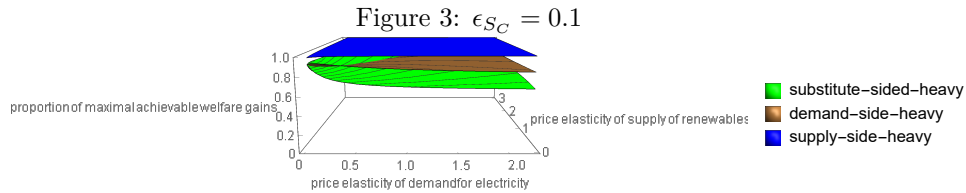
I have only found a single study, namely Johnson (2011), which gives an estimate of 2.7. In the plots shown below I consider the range from 0.1 to 3 for the price elasticity of supply of renewables.

## E.3 Estimates of the price elasticity of supply of coal

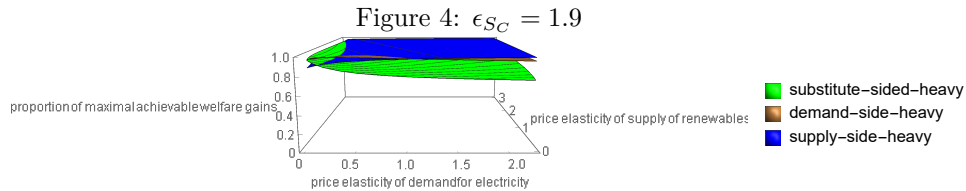
Daubanes, J., Henriot, F., & Schubert, K. (2020). note that the empirical literature on the price elasticity of coal supply—e.g., Labys et al. (1979), Beck et al. (1991), Light (1999), Light et al. (1999), and Dahl (2009)—finds estimates ranging from 0.1 and 1.9. Based on this, I consider the range from 0.1 to 1.9.

## E.4 Results

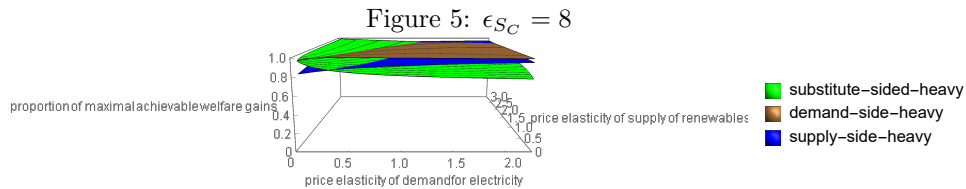
Here is the case where the price elasticity of supply of coal is 0.1:



Here is the case where the price elasticity of supply of coal is 1.9:



Whilst the estimates for the price elasticity of supply of coal range from 0.1 to 1.9, we presumably cannot rule out potentially much large value for it in the long term. For illustration, consider the case where the price elasticity of supply of coal is 8:



Overall, these results suggest that the conclusion that the welfare losses from misallocation are likely to be small is robust.

## References

- [1] Angelsen, A. (2017). REDD+ as result-based aid: General lessons and bilateral agreements of Norway. *Review of Development Economics*, 21(2), 237-264.
- [2] Armstrong, M., & Rochet, J. C. (1999). Multi-dimensional screening: A user's guide. *European Economic Review*, 43(4-6), 959-979.
- [3] Bayram, A. B., & Graham, E. R. (2017). Financing the United Nations: Explaining variation in how donors provide funding to the UN. *The Review of International Organizations*, 12(3), 421-459.
- [4] Bretschger and Pattakou (2018), Bretschger, L., & Pattakou, A. (2019). As bad as it gets: how climate damage functions affect growth and the social cost of carbon. *Environmental and resource economics*, 72(1), 5-26.
- [5] Burke, P. J., & Abayasekara, A. (2018). The price elasticity of electricity demand in the United States: A three-dimensional analysis. *The Energy Journal*, 39(2).
- [6] Cramton, P., & Stoft, S. (2012). Global climate games: How pricing and a green fund foster cooperation. *Economics of Energy & Environmental Policy*, 1(2), 125-136.
- [7] Collier, P., & Venables, A. J. (2014). Closing coal: economic and moral incentives. *Oxford Review of Economic Policy*, 30(3), 492-512.
- [8] Dahl, C. (2009). Energy demand and supply elasticities. *Energy Policy*, 72.
- [9] Daubanes, J., Henriot, F., & Schubert, K. (2020). Unilateral CO2 Reduction Policy with More Than One Carbon Energy Source.
- [10] Edenhofer, O., & Kalkuhl, M. (2011). When do increasing carbon taxes accelerate global warming? A note on the green paradox. *Energy Policy*, 39(4), 2208-2212.
- [11] Eichner, T., Kollenbach, G., & Schopf, M. (2020). Buying versus leasing fuel deposits for preservation. *The Scandinavian Journal of Economics*.
- [12] Espey, J. A., & Espey, M. (2004). Turning on the lights: A meta-analysis of residential electricity demand elasticities. *Journal of Agricultural and Applied Economics*, 36(1379-2016-112600), 65-81.
- [13] Fæhn, T., Hagem, C., Lindholt, L., Mæland, S., & Rosendahl, K. E. (2017). Climate policies in a fossil fuel producing country—demand versus supply side policies. *The Energy Journal*, 38(1).

- [14] Harstad, B. (2012). Buy coal! A case for supply-side environmental policy. *Journal of Political Economy*, 120(1), 77-115.
- [15] Hoel, M. (2014). Supply side climate policy and the green Paradox. *Climate policy and nonrenewable resources: The green paradox and beyond*, 21.
- [16] IPCC (2014). *Climate change, synthesis report, summary for policymakers*
- [17] Johnson, E. P. (2011). The price elasticity of supply of renewable electricity generation: Evidence from state renewable portfolio standards.
- [18] Kornek, U., & Edenhofer, O. (2020). The strategic dimension of financing global public goods. *European Economic Review*, 127, 103423.
- [19] Long, N. V. (2015). The green paradox in open economies: Lessons from static and dynamic models. *Review of Environmental Economics and Policy*, 9(2), 266-284.
- [20] Martimort, D., & Sand-Zantman, W. (2016). A mechanism design approach to climate-change agreements. *Journal of the European Economic Association*, 14(3), 669-718.
- [21] Mas-Colell, A., Whinston, M. D., & Green, J. R. (1995). *Microeconomic theory* (Vol. 1). New York: Oxford university press.
- [22] Mertz, O., Grogan, K., Pflugmacher, D., Lestrelin, G., Castella, J. C., Vongvisouk, T., ... & Müller, D. (2018). Uncertainty in establishing forest reference levels and predicting future forest-based carbon stocks for REDD+. *Journal of Land Use Science*, 13(1-2), 1-15.
- [23] Rezai, A., & Van Der Ploeg, F. (2017). Second-best renewable subsidies to de-carbonize the economy: commitment and the Green Paradox. *Environmental and Resource Economics*, 66(3), 409-434.
- [24] Ross, M. L. (2015). What have we learned about the resource curse?. *Annual Review of Political Science*, 18, 239-259.
- [25] Seymour, F., & Busch, J. (2016). *Why forests? Why now?: The science, economics, and politics of tropical forests and climate change*. Brookings Institution Press.
- [26] Steckel, J. C., Edenhofer, O., & Jakob, M. (2015). Drivers for the renaissance of coal. *Proceedings of the National Academy of Sciences*, 112(29), E3775-E3781.
- [27] Teng, M., Burke, P. J., & Liao, H. (2019). The demand for coal among China's rural households: Estimates of price and income elasticities. *Energy Economics*, 80, 928-936.
- [28] Van der Ploeg, F., & Withagen, C. (2014). Growth, renewables, and the optimal carbon tax. *International Economic Review*, 55(1), 283-311.

- [29] Van der Ploeg, F., & Withagen, C. (2015). Global warming and the green paradox: A review of adverse effects of climate policies. *Review of Environmental Economics and Policy*, 9(2), 285-303.
- [30] Van der Ploeg, F. (2016). Second-best carbon taxation in the global economy: the Green Paradox and carbon leakage revisited. *Journal of Environmental Economics and Management*, 78, 85-105.