

Catastrophic climate change under time-inconsistent preferences

Thomas Michielsen*
Tilburg University

July 31, 2012

Abstract

We analyze optimal fossil fuel use in an intergenerational model with the risk of a climate catastrophe. Each generation maximizes a weighted sum of discounted utility (positive) and the probability that a climate catastrophe will occur at any point in the future (negative). The model generates time-inconsistency as generations disagree on the relative weights on consumption and catastrophe prevention. As a consequence, future generations emit too much from the current generation's perspective and a dynamic game ensues. We consider a sequence of models. If the fossil fuel stock is finite, early generations are more conservationist in Markov equilibrium than under commitment in order to smooth fossil fuel use over time. When fossil fuels are expected to become redundant in the near future, early generations may increase or decrease their fossil fuel consumption in Markov equilibrium depending on the utility and threshold distribution functions. When an abundant fossil fuel will remain essential as a production factor, the catastrophe becomes a self-fulfilling prophecy.

JEL-Classification: C73, D83, Q54

Keywords: catastrophic events, decision theory, uncertainty, time consistency

1 Introduction

A wide body of research suggests that the impact of man-made environmental pressures on the flow of ecosystem services is uncertain and highly non-linear.

*CentER, Department of Economics and Tilburg Sustainability Center, Tilburg University, P.O. Box 90153, 5000 LE Tilburg, The Netherlands. E-mail: t.o.michielsen@uvt.nl

Estimates of the climate sensitivity - the temperature increase as a result of a doubling of the CO₂ stock in the atmosphere - range from one to more than ten degrees [14]. Almost all of the uncertainty stems from possibility that at some threshold level of emissions, positive feedback mechanisms are set in motion: the melting of polar ice caps will increase solar radiation absorption and permafrost melting in the Arctic could cause large methane releases [10].¹

Whereas the impacts of modest temperature increases are potentially manageable, a rise in the high single digits and upwards will likely have catastrophic and irreversible consequences, including large permanent loss of biodiversity, sea level rise and increased prevalence of extreme weather events. The possibility of a climate catastrophe has important implications for intergenerational welfare analysis [9, 19, 20].

We analyze optimal carbon emissions using a welfare function in the spirit of [4, 5, 1] that consists of a weighted sum of expected discounted utility (positive) and the probability that a catastrophe will occur at any point in the future (negative). We consider uncertainty of the type in [17, 18, 11]. The catastrophe occurs when an unknown emission threshold is breached. The threshold must lie above the maximum concentration level that has been reached in the past; otherwise, the catastrophe would have occurred already. Keeping the carbon stock constant eliminates the hazard. Increasing the stock increases consumption temporarily or permanently if the threshold is not breached, but may trigger the catastrophe.

Our welfare function has positive and normative appeal. Firstly, it accommodates both impatience from observed behaviour and concerns about the far-distant future. Secondly, the intrinsic welfare loss from a catastrophe reflects the notion that while some environmental risks are justifiable from an expected cost-benefit point of view, it may be ethically unpalatable to gamble with future generations' welfare. Thirdly, the intrinsic loss term captures the existence value of the Earth's environment and intrinsic aversion to man-made disaster. As the value of species preservation is only partly tied to current and future use, it is natural to suppose that the welfare loss from future extinctions only partly depends on the time of occurrence. Alternative welfare criteria, such as discounted utility [13, 15] or hyperbolic discounting [8, 6] cannot incorporate some or all of these concerns.

¹Although we will focus on climate change, non-linearities are also present in other applications. Ecosystems can rapidly collapse if biodiversity drops below a critical value [16] and tipping points for deforestation and temperature increases may turn rainforests into savannahs [12].

Table 1: Time inconsistency

	current consumption	welfare weights	
		future consumption	catastrophe prevention
current generation	1	$\rho < 1$	ξ
future generation		1	ξ

When the welfare loss from a catastrophe does not depend on the time of occurrence, the preference structure becomes time-inconsistent: current generations discount future consumption relative to their intrinsic welfare loss from a catastrophe, but future generations do not discount their own consumption relative to the intrinsic catastrophe loss (see Table 1). As a result, future generations emit too much from the current generation's perspective and a dynamic game ensues. Early generations have a strategic motive to distort their fossil fuel use in order to influence future emissions.

We consider a sequence of models. Firstly, we introduce a two-period model with an abundant fossil fuel (e.g. coal). This model represents a setting in which a clean alternative to fossil fuels is expected to be developed in the near future (at the end of the second period). Secondly, we analyze an infinite-horizon model with emission decay and an abundant fossil fuel. Lastly, we consider an infinite-horizon model with emission decay and a finite stock of fossil fuel (e.g. oil). In each model, we compare fossil fuel use and the probability of a catastrophe in three cases: (a) when the first generation can commit all current and future fossil fuel use (the commitment solution), (b) when current generations do not anticipate that future generations have different preferences (the naive solution) and (c) when current generations take into account the reaction of future generations (the Markov equilibrium).

In the two-period model, first-period emissions are higher in the commitment solution than in Markov equilibrium for uniformly distributed thresholds and iso-elastic utility. The benefit of mitigating the imbalance between discounted marginal utilities in the two periods outweighs the cost of increasing the probability of a catastrophe. In the infinite-horizon model with an abundant fossil fuel, the catastrophe becomes a self-fulfilling prophecy. In Markov equilibrium, early generations realize that they cannot influence the steady-state stock and that future generations will undo their mitigation efforts, so they opt for higher emission flows than in the commitment solution. When the fossil fuel stock is finite and sufficiently small, the ability of future generations to drive

up the carbon concentration is limited. Early generations reduce emissions in Markov equilibrium compared to the commitment and naive solutions, in order to smooth fossil fuel use over time and so reduce the maximum carbon stock.

2 Two-period model

We consider a model with two generations, living in periods $t = 1, 2$. A representative agent in each generation derives utility $u_t(z_t)$ from consuming an abundant fossil fuel z , the economy's single commodity. The utility functions satisfy $u'_t \geq 0$, $u''_t \leq 0$, $\exists \bar{u}_t : u_t(z) < \bar{u}_t \forall z$. Fossil fuel use z_t contributes to the stock of carbon in the atmosphere D_t . We abstract from natural decay of carbon, so that

$$D_t = D_{t-1} + z_t, \quad D_0 \text{ given}$$

A climate catastrophe occurs when the stock of carbon reaches an unknown threshold \hat{D} at the end of a period. The threshold is randomly distributed on the interval $[0, \bar{D}]$. We express the probability of a catastrophe as a function of cumulative emissions through pdf $f(D)$ and cdf $F(D)$. When $z_1 > \hat{D}$, the second generation gets utility \underline{u} .² We may think of $u_2(z_2) - \underline{u}$ as the extrinsic welfare loss from the catastrophe.

The welfare function of each generation is given by a weighted sum of discounted utility (positive) and the probability that a catastrophe will occur in either period (negative). The first generation discounts utility of the second generation by a factor $\rho < 1$, but its welfare loss from the catastrophe does not depend on the time of occurrence. The second generation observes whether the first generation's emissions have triggered the catastrophe or not. We relax this assumption in the Appendix. When the catastrophe is only observed at the end of the second period, the second generation chooses a higher z_2 because there is a probability that the first generation has already triggered the catastrophe, in which case second-period mitigation is fruitless. The welfare function for the two generations read

$$w_1 = u_1(z_1) + (1 - F(z_1)) \rho u_2(z_2) + F(z_1) \rho \underline{u} - \xi F(z_1 + z_2) \quad (1a)$$

$$w_2 = \begin{cases} u_2(z_2) - \xi \frac{F(z_1 + z_2)}{1 - F(z_1)} & \text{if } z_1 < \hat{D} \\ \underline{u} - \xi & \text{if } z_1 \geq \hat{D} \end{cases} \quad (1b)$$

²In the remainder of this paper, we will regard \underline{u} to be sufficiently small such that the catastrophe is also undesirable from a point of view of material consumption. This is not necessary for the formal analysis however. If the catastrophe does not affect utility, all post-catastrophe generations choose $z_t \rightarrow \infty$ and $\underline{u} = \bar{u}_t$.

where $\xi > 0$ indicates the welfare weight on catastrophe prevention. Catastrophe risk in the second period is evaluated using the conditional cdf $\frac{F(z_1^C + z_2^C)}{1 - F(z_1^C)}$. The discount factor generates time-inconsistency in the preference structure: the second generation places a higher weight on second-period consumption $u_2(z_2)$ relative to the probability of a catastrophe $F(D_2)$ than the first generation does.

We distinguish between three solutions. Firstly, the commitment solution (superscript C), in which the first generation commits all current and future fossil fuel use. Secondly, the 'naive' solution (superscript N), in which the first generation does not anticipate that future generations will make a different trade-off between $u_2(z_2)$ and $F(D_2)$. Lastly, we consider the Markov solution (superscript M), in which the first generation foresees the preference reversal of the second generation and selects z_1 by backward induction, maximizing its welfare given the optimal response of the second generation.

2.1 Commitment solution

When the first generation can commit second-period fossil fuel use conditional on whether the catastrophe occurs at the end of the first period and the solution is interior, z_1^C and z_2^C immediately follow from (1b)

$$u'_1(z_1^C) - \rho f(z_1^C) [u_2(z_2^C) - \underline{u}] - \xi f(z_1^C + z_2^C) = 0 \quad (2a)$$

$$\rho u'_2(z_2^C) - \xi \frac{f(z_1^C + z_2^C)}{1 - F(z_1^C)} = 0 \quad \text{if } z_1^C < \hat{D} \quad (2b)$$

The first generation equates discounted marginal utility in both periods with the marginal welfare loss from catastrophe risk. In case of a corner solution, $(z_1^C, z_2^C) \rightarrow (\infty, \infty)$.

2.2 Naive solution

In the naive solution, the first generation behaves as if it could commit both z_1 and z_2 . The second generation however selects z_2^N to maximize (1b) rather than (1a), yielding

$$u'_1(z_1^N) - \rho f(z_1^N) [u_2(z_2^C) - \underline{u}] - \xi f(z_1^N + z_2^C) = 0 \quad (3a)$$

$$u'_2(z_2^N) - \xi \frac{f(z_1^N + z_2^N)}{1 - F(z_1^N)} = 0 \quad \text{if } z_1^N < \hat{D} \quad (3b)$$

By definition, z_1 is the same in the naive solution as in the commitment solution. Substituting $z_1^N = z_1^C$ in (3a) and comparing with (2a), we obtain $z_2^N > z_2^C$. The

second generation may give up on catastrophe prevention altogether ($z_2^N \rightarrow \infty$), voiding the first generation's mitigation efforts when $z_1^N < \hat{D}$.

2.3 Markov solution

In the Markov solution, the first generation correctly takes into account the second generation's reaction. Condition (3a) implicitly defines the second generation's reaction function $r(z_1)$

$$u'_2(r(z_1^M)) = \xi \frac{f(z_1^M + r(z_1^M))}{1 - F(z_1^M)} \quad (4)$$

To avoid clutter, we omit the superscript M . Differentiating with respect to z_1 , we obtain

$$\begin{aligned} u''_2(r(z_1)) r'(z_1) &= \xi \left(\frac{f'(z_1 + r(z_1)) [1 - F(z_1)] + f(z_1) f(z_1 + r(z_1))}{[1 - F(z_1)]^2} + \right. \\ &\quad \left. r'(z_1) \frac{f'(z_1 + r(z_1))}{1 - F(z_1)} \right) \\ r'(z_1) &= \xi \frac{f'(z_1 + r(z_1)) [1 - F(z_1)] + f(z_1) f(z_1 + r(z_1))}{[1 - F(z_1)] [u''_2(r(z_1)) [1 - F(z_1)] - \xi f'(z_1 + r(z_1))]} \end{aligned} \quad (5)$$

The numerator in (5) is positive when F has an increasing hazard function. The sign of the denominator depends on the curvature of f . A sufficient condition for the second generation's reaction function to be downward-sloping is $f'(z_1 + r(z_1)) \geq 0$. When $f'(z_1 + r(z_1))$ is sufficiently negative, an increase in z_1 lowers the marginal probability of a catastrophe to such an extent that it becomes attractive for the second generation to choose a higher emission level.

The first-order condition for the first generation is

$$\begin{aligned} u'_1(z_1) - \rho f(z_1) u_2(r(z_1)) + \rho [1 - F(z_1)] u'_2(r(z_1)) r'(z_1) + \rho f(z_1) \underline{u} - \\ \xi f(z_1 + r(z_1)) (1 + r'(z_1)) &= 0 \\ \Leftrightarrow \underbrace{u'_1(z_1)}_I - \underbrace{\rho f(z_1) [u_2(r(z_1)) - \underline{u}]}_{II} - \underbrace{\xi (1 - \rho) f(z_1 + r(z_1)) r'(z_1)}_{III} - \underbrace{\xi f(z_1 + r(z_1))}_{IV} &= 0 \end{aligned} \quad (6)$$

The four components of (6) represent the first generation's considerations. The first term is the first generation's marginal utility. The second term reflects that higher first-period fossil fuel use increases the probability that the second generation's utility is reduced to \underline{u} . The third term represents the strategic

motive to influence the second generation's fossil fuel use. From the point of view of the first generation, the second generation consumes too much fossil fuel. When $r'(z_1)$ is negative (positive), the first generation can reduce z_2 by increasing (decreasing) its own fossil fuel use. This term is not present in (2a). The fourth term indicates the first generation's intrinsic desire to prevent a catastrophe.

Comparing (6) and (2a), it is not possible to say whether first-period emissions are higher in the Markov or in the commitment solution without assuming functional forms for u_t and F . In Lemma 1 we demonstrate that with iso-elastic utility and a uniformly distributed catastrophe threshold, first-period emissions are *higher* in the Markov solution than in the commitment outcome. Conversely, with quadratic utility, the first generation is more cautious in the Markov solution than under commitment.

Lemma 1. *Let \hat{D} be uniformly distributed ($F(D) = D/\bar{D}$ and $f(D) = 1/\bar{D}$) and ξ sufficiently large so that $z_1^C < \infty$, $z_1^M < \infty$. For iso-elastic utility $u_t(z_t) = \frac{z_t^{1-\eta}}{1-\eta}$, $z_1^M > z_1^C$. For quadratic utility $u_t(z_t) = az_t - \frac{1}{2}bz_t^2$, $z_1^M < z_1^C$.*

When \hat{D} is uniformly distributed, the last terms in the first-order conditions (6) and (2a) are equal. The FOCs differ in two respects. In the Markov solution, the sufficient condition $f'(z_1 + r(z_1)) \geq 0$ for the second generation's reaction function to be downward sloping is satisfied. This gives the first generation an incentive to consume more fossil fuel in the Markov than in the commitment solution. On the other hand, for the same level of z_1 the second generation consumes more fossil fuel in the Markov solution. This increases the utility loss to the second generation $u_2(z_2) - \underline{u}$ in case of a first-period catastrophe, which makes the first generation more prudent in the Markov solution. The first effect dominates for iso-elastic utility; the second effect for quadratic utility.

3 Infinite horizon, abundant fossil fuel

In this section, we consider an infinite-horizon model with a continuum of non-overlapping generations. As in the previous section, each generation derives utility from its own material consumption $u(z_t)$, cares about future material consumption (discounted at rate r) and the possibility of a climate catastrophe occurring at some point the future. A constant fraction of the emissions stock decays in each period. In the Appendix, we outline the necessary conditions for dynamic optimization problems with uncertain thresholds, as derived in [11].

3.1 Optimal fossil fuel use

In each period, a fraction α of the carbon stock D decays naturally so that

$$\dot{D} = z - \alpha D \quad (7)$$

Utility is concave and bounded and the hazard rate is increasing.

Assumption 1. $u(D, z) = u(z)$, $u'(z) > 0$, $u''(z) < 0 \forall z$ and $\lim_{z \rightarrow \infty} u(z) = \bar{u}$.

Assumption 2. $\psi(x) = \psi(D)$, $\psi'(D) \geq 0$.

When the catastrophe occurs, all subsequent generations receive utility $\underline{u} < \bar{u}$. We model this by creating a state variable γ that takes the value 0 when the catastrophe has not yet occurred, and is equal to $\bar{u} - \underline{u}$ after the threshold has been passed. Let τ denote the occurrence time of the catastrophe. For a given admissible trajectory $D(s)$, the welfare of generation t is³

$$V(t, D) = \begin{cases} \mathbb{E} \left(\int_t^\infty (-\gamma(s) + u(z(s))) e^{-r(s-t)} ds \right) - \xi \frac{P[\tau \in [t, \infty)]}{1 - P[\tau \in [0, t]]} & \text{for } \tau \notin [0, t] \\ s.t. \dot{D} = z - \alpha D, D(0) = D_0, \gamma(\tau^+) = \gamma(\tau^-) + \bar{u} - \underline{u} \\ \frac{\underline{u}}{r} - \xi & \text{for } \tau \in [0, t] \end{cases} \quad (8)$$

We analyze the commitment, naive and Markov solutions in turn. We assume existence of optimal solutions throughout. $D(t)$ is non-decreasing along the optimal path [18], so admissible paths satisfy Definition 1 in the Appendix.

3.2 Commitment solution

If the first generation can commit all current and future fossil fuel use, it maximizes (8) for $t = 0$. Its problem is

$$\begin{aligned} \max_z V^C(0, x(0)) &= \mathbb{E} \left(\int_0^\infty (-\gamma(s) + u(z(s))) e^{-rs} ds \right) - \xi P[\tau \in [0, \infty)] \\ s.t. \dot{D} &= z - \alpha D, D(0) = D_0, \gamma(\tau^+) = \gamma(\tau^-) + \bar{u} - \underline{u} \end{aligned} \quad (9)$$

We may rewrite the problem by including the intrinsic welfare loss from the catastrophe into the per-period post-catastrophe penalty γ . Since γ is subject to the discount rate, the per-period penalty $\hat{\gamma}$ in the alternative representation is

³ τ is distributed as a Poisson process, as we describe in the Appendix. For brevity, we omit the distribution of τ in the main text.

increasing in the time of occurrence τ to preserve the property that the intrinsic loss is independent of τ :

$$\begin{aligned} \max_z V^C(0, D(0)) &= \mathbb{E} \left(\int_0^\infty (-\hat{\gamma}(s) + u(z(s))) e^{-rs} ds \right) \\ \text{s.t. } \dot{D} &= z - \alpha D, \quad D(0) = D_0, \quad \hat{\gamma}(\tau^+) = \hat{\gamma}(\tau^-) + \bar{u} - \underline{u} + r\xi e^{r\tau} \end{aligned} \quad (10)$$

As time passes, it becomes prohibitively costly from the first generation's point of view to risk a catastrophe. It will therefore stabilize the emissions stock at some finite date t' such that the marginal benefit of increasing the emissions stock (higher current utility and higher steady-state utility if the threshold is not breached) equals the expected marginal cost (a permanent decrease in consumption and the intrinsic welfare loss evaluated at $\tau = t'$ if the catastrophe does occur).

Proposition 1. *The commitment solution is characterized by a steady-state emissions stock D^C . There exists a $t' < \infty$ such that $D^C(t') = D^C$ and $z^C(t) = \alpha D^C \forall t \geq t'$. D^C and t' satisfy*

$$u'(\alpha D^C) = \frac{\psi(\alpha D^C)}{r + \alpha} \left[u(\alpha D^C) - \underline{u} + r\xi e^{rt'} \right] \quad (11)$$

3.3 Naive solution

In the naive solution, each generation solves a problem that is similar to (9), with the initial carbon stock determined by previous generations. Each generation t envisions a preferred steady-state stock $D^{t,N}$, but as every subsequent generation places a lower relative welfare weight on catastrophe prevention, the emission targets $D^{t,N}$ increase over time. The emission targets converge to a unique level D^N that even the most distant generations do not want to exceed, as the marginal utility of higher steady-state consumption falls short of the permanent utility reduction and the welfare loss associated with a catastrophe.

Proposition 2. *The solution to generation t 's problem*

$$\begin{aligned} \max_z V^{t,N}(t, D(t)) &= \mathbb{E} \left(\int_t^\infty (-\hat{\gamma}(s) + u(z(s))) e^{-r(s-t)} ds \right) \\ \text{s.t. } \dot{D} &= z - \alpha D, \quad D(t) = D_t, \quad \hat{\gamma}(\tau^+) = \hat{\gamma}(\tau^-) + \bar{u} - \underline{u} + r\xi e^{r(\tau-t)} \end{aligned} \quad (12)$$

is characterized by a steady-state emissions stock $D^{t,N}$. Let D^N be given by

$$u'(\alpha D^N) = \frac{\psi(\alpha D^N)}{r + \alpha} \left[u(\alpha D^N) - \underline{u} + r\xi \right] \quad (13)$$

Then $D^{t,N} < D^N \forall t$ and $\lim_{t \rightarrow \infty} D^{t,N} = D^N$.

The left and right hand side of (13) represent the marginal benefit and cost of increasing the steady-state emission level, respectively. Because of Assumptions 1 and 2, the left hand side is decreasing and the right hand side is increasing. Therefore, it cannot be optimal for any generation t that inherits carbon stock D^N to choose $z^{t,N}(t) > \alpha D^N$. As a consequence, the carbon stock never exceeds D^N .

Corollary 1. $D^C < D^N$

The steady-state carbon stock is higher in the naive solution than in the commitment solution, because later generations with lower present-value welfare weights on catastrophe prevention reoptimize towards higher steady-state carbon stocks.

3.4 Markov solution

When all generations are symmetric, the existence of a Markov equilibrium implies the existence of a policy function $\zeta^M(D)$ such that $z^M(t, D) = \zeta^M(D) \forall t$. In Proposition 3, we show that there exists a continuum of Markov equilibria which can be ranked by their steady-state carbon stocks. The actions of early generations depend on their beliefs about future generations' actions. When early generations believe future generations will increase the carbon stock up to a certain level D^M , the formers' choice of $z^M(t, D)$ has no effect on the second component of (8). Each generation then maximizes the integral of expected discounted utility subject to the carbon stock not exceeding the perceived maximum. The range of equilibria is bounded because of two considerations. When the perceived maximum is very high, it will not be reached because it would interfere with maximizing expected discounted utility (the first component of (8)). When the perceived maximum is too low (below D^N), future generations can increase their welfare by raising the carbon stock above this maximum.

Proposition 3. Let $D_1^M = D^N$ given by (13) and D_2^M be given by

$$u'(\alpha D_2^M) = \frac{\psi(D_2^M)}{r + \alpha} (u(\alpha D_2^M) - \underline{u}) \quad (14)$$

Assume a solution to

$$V^M(t, D) = \max_z \mathbb{E} \left(\int_t^\infty (-\gamma(s) + u(z(s))) e^{-r(s-t)} ds \right) \quad (15)$$

$$s.t. \dot{D} = z - \alpha D, \quad D(t) = D, \quad D(s) \leq D^M \quad \forall s \geq t, \quad \gamma(\tau^+) = \gamma(\tau^-) + \bar{u} - \underline{u} \quad (16)$$

exists $\forall x : D \leq D^M$, $D^M \in [D_1^M, D_2^M]$. Then there exists a continuum of Markov equilibria indexed by $D^M \in [D_1^M, D_2^M]$ such that

$$\zeta^M(D) = \begin{cases} \operatorname{argmax}_z V^M(t, D) & \text{if } D < D^M \\ \alpha D & \text{if } D \geq D^M \end{cases} \quad (17)$$

When all generations hold the same beliefs about the steady-state stock, those beliefs become self-fulfilling, even if they result in an inefficient equilibrium $D^M > D^N$. The $D^M = D^N$ equilibrium yields the highest welfare for all generations as it comes closest to internalizing the intrinsic welfare loss from a catastrophe. By contrast, in the naive solution each generation believes it decides the steady-state stock. Since it is in no generation's interest to exceed D^N , $D(t) > D^N$ is ruled out.

Corollary 2. *The first generation's welfare in the naive solution is lower than in the Markov solution when $D^M = D^N$.*

The naive solution suffers from a different inefficiency. Generations mistakenly perceive the steady-state carbon stock to be $D^{t,N} < D^N$, so their fossil fuel consumption does not maximize expected discounted utility under the correct belief D^N . In the Markov solution with $D^M = D^N$, all generations have consistent beliefs, so the $z(s)$ path does maximize $\mathbb{E} \left(\int_t^\infty (-\gamma(s) + u(z(s))) e^{-r(s-t)} ds \right)$ subject to $D(s) \leq D^N \forall s \geq t$. Figure 1 and 2 illustrate fossil fuel use in the three scenarios. Emissions in the naive solution are initially close to those in the commitment solution, but increasingly diverge as future generations put more weight on the cost of mitigation. The Markov solution converges to the same maximum emissions stock as the naive solution, but the maximum is attained much earlier, resulting in higher welfare for early generations than in the naive solution.

Proposition 4. *Let $\zeta^C(D)$, $\zeta^N(D)$ and $\zeta^M(D)$ denote the optimal strategies conditional on the emissions stock D in the commitment, naive and Markov solutions, respectively.⁴ $\zeta^C(D) < \zeta^N(D) < \zeta^M(D) \forall D > D_0$ and $D^C(t) < D^N(t) < D^M(t) \forall t > 0$.*

By Propositions 1, 2 and 3, the solution of the commitment, naive and Markov problems is the same as that of a constrained optimization problem ala [2, 3] in which the current generation maximizes expected discounted utility subject to the emissions stock not exceeding an exogenous ceiling D^C , $D^{t,N}$ or

⁴ $\zeta^C(D)$ is only optimal along the equilibrium path.

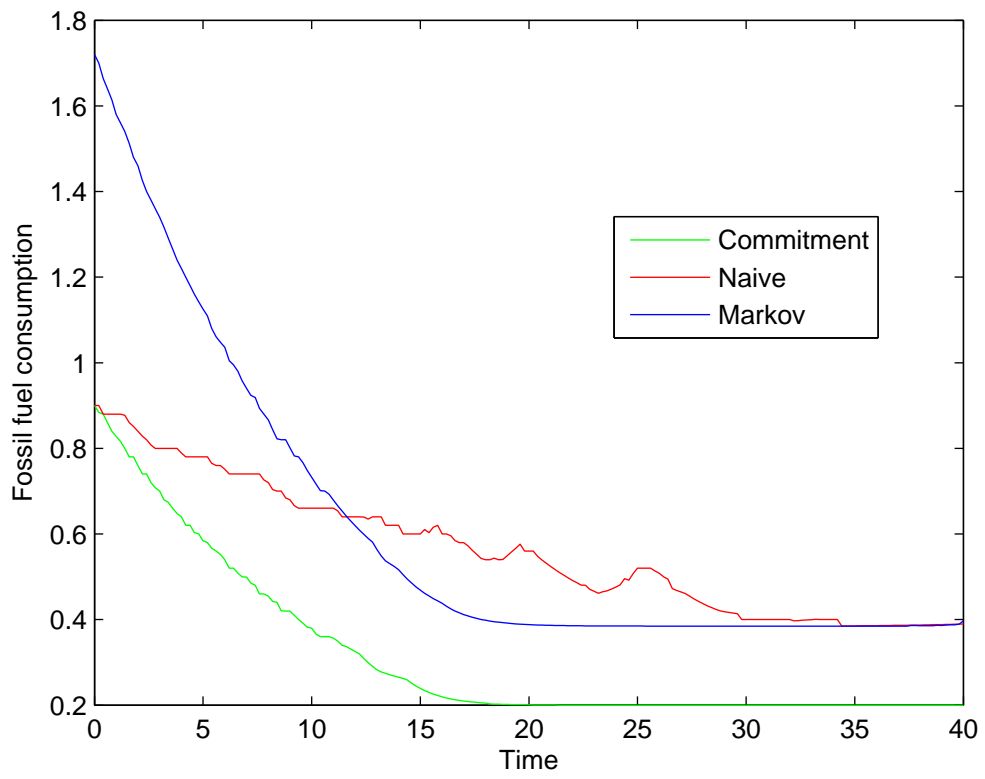


Figure 1: Emission flows $z(t)$ in commitment, naive and Markov solutions

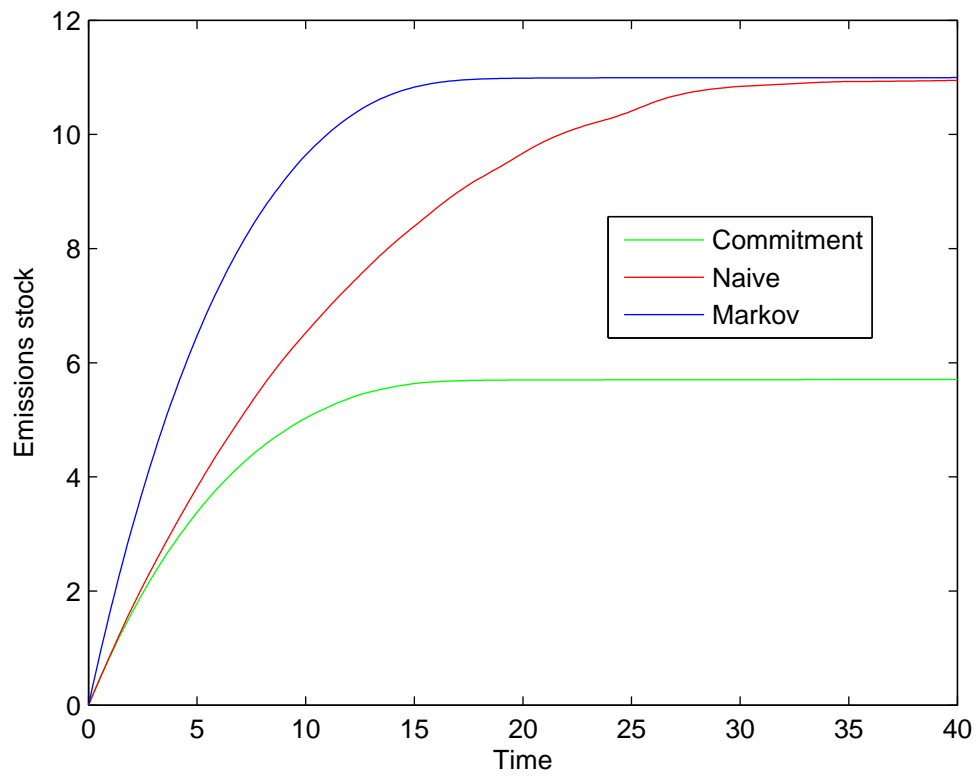


Figure 2: Carbon stocks $D(t)$ in commitment, naive and Markov solutions

D^M at any point in time. This 'carbon budget' is larger in the naive and Markov solutions, so conditional on the stock of carbon D the emission flows $\zeta^N(D)$ and $\zeta^M(D)$ are higher than in the commitment solution. Because emission flows can be ranked for any given stock, the carbon stocks can also be ranked unambiguously at each point in time.

In the Appendix, we analyze an extension in which the catastrophe occurs with an exogenous delay after the threshold is breached. Steady-state carbon stocks are higher in this case: because the catastrophe no longer reduces utility immediately, the damage from a catastrophe is not as severe as in the benchmark model.

4 Infinite horizon, scarce fossil fuel

In this section, we consider optimal policies when the fossil fuel is scarce. Cumulative extraction cannot exceed the initial stock S_0 . Because the emissions stock can no longer be monotonically increasing, we append the framework from the previous section by using a discrete approximation. We consider N generations; each generation is alive for a period of length ϵ . Let N be very large and ϵ very small. Let $D_{\max}(\epsilon i) \equiv \max_{j < i} D(\epsilon j)$ denote the maximum carbon stock that has been reached until time ϵi and $\tau \equiv \epsilon \operatorname{argmin}_i \{D_{\max}(\epsilon i) \geq \hat{D}\}$ be the occurrence time of the catastrophe. For simplicity, and because the scarcity of fossil fuel already limits post-catastrophe utility, we abstract from direct utility damages after a catastrophe. The welfare of generation i , given a path $z(j)$ is then

$$\begin{aligned} V^i(S, D, D_{\max}) &= \sum_{j=i}^N \epsilon u(z(j)) e^{-r\epsilon(j-i)} - \xi \frac{P(\tau \leq \epsilon N)}{1 - P(\tau \leq \epsilon i)} \quad (18) \\ &+ \Gamma(S(N), D(N), D_{\max}(N)) \\ &\quad \sum_{j=1}^N \epsilon z(j) \leq S(i) \\ \text{s.t. } D(j+1) &= (1 - \epsilon\alpha) D(j) + \epsilon z(j) \\ D_{\max}(j+1) &= \max\{D(j+1), D_{\max}(j)\} \end{aligned}$$

where the terminal payoff

$$\Gamma : \left\{ S, D, D_{\max} \mid \operatorname{argmax}_{z(i)} \left(\sum_{j=i}^{\infty} \epsilon u(z(j)) e^{-r\epsilon(j-i)} \right) \leq \alpha D \right\} \rightarrow \mathbb{R}$$

is defined as

$$\max_{z(j)} \sum_{j=0}^{\infty} \epsilon u(z(j)) e^{-r\epsilon(N+j)} \text{ s.t. } \sum_{j=0}^{\infty} \epsilon z(j) \leq S$$

When N is sufficiently large, $S(N)$ is sufficiently small so that the catastrophe will not occur after period N with certainty⁵, and the scrap value is simply given by a cake-eating problem with remaining stock $S(N)$.

Proposition 5. *A maximum to $V^i(S, D, D_{max})$ exists, and the intergenerational game has at least one Markov equilibrium. The equilibrium paths are characterized by a terminal phase in which there is no catastrophe hazard.*

Proof. By our definition of the terminal payoff Γ and the observation in footnote 3, the problem is a finite extensive-form game, which has at least one subgame-perfect equilibrium. \square

We simulate optimal fossil fuel use in the commitment and Markov solution for a quadratic utility function using a discrete grid for (S, D, D_{max}) . The results are depicted in Figures 3 and 4. Initial emissions are lower in the Markov equilibrium than under commitment. If the first generation decreases its fossil fuel use by say 10 units, only a fraction of this reduction will be undone by the subsequent generation. Though the fossil fuel stock is still exhausted eventually, fossil fuel use is spread more evenly over time. Because of the natural decay of carbon in the atmosphere, this reduces the maximum carbon stock, and hence the probability of a catastrophe.

5 Conclusion

When generations care about the near term and the far-distant future, preferences become time-inconsistent. Current generations would like their descendants to restrain their consumption in order to reduce the risk of a catastrophe, but future generations attach a higher importance to the costs they would have to bear in such mitigation efforts. If polluting inputs are expected to remain an essential production factor and future generations' ability to use these inputs is unrestrained, current generations recognize that they cannot prevent pollution stocks from reaching undesirably high levels and hence reduce their own mitigation efforts. Given the large reserves of coal and unconventional oil and gas, this is a disconcerting message. Likewise, the strain on renewable resources such as arable land, forests and fish stocks is unlikely to let up for a time to come, which could trigger regime shifts that cause irreversible biodiversity loss, land cover change and soil degradation. Current generations only have an incentive to reduce pollution in the face of future preference reversals if the polluting inputs

⁵If not, we must have $\lim_{N \rightarrow \infty} S(N) \gg 0$, which cannot be optimal in either the commitment, naive or the Markov solution.

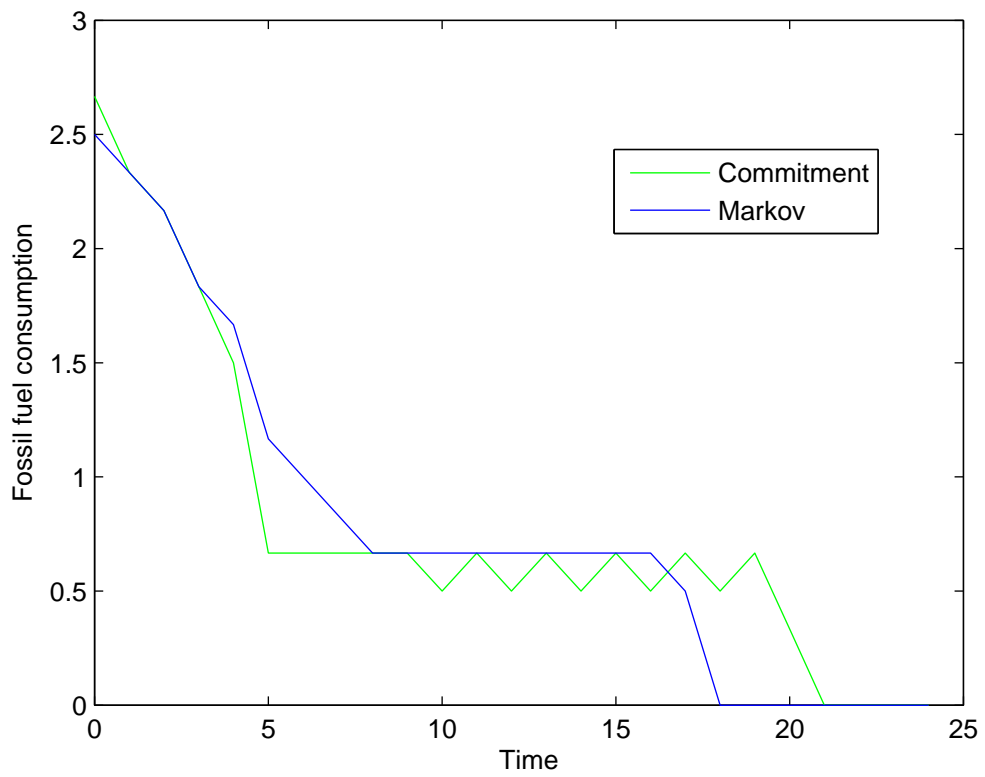


Figure 3: Emission flows $z(j)$ in commitment and Markov solutions

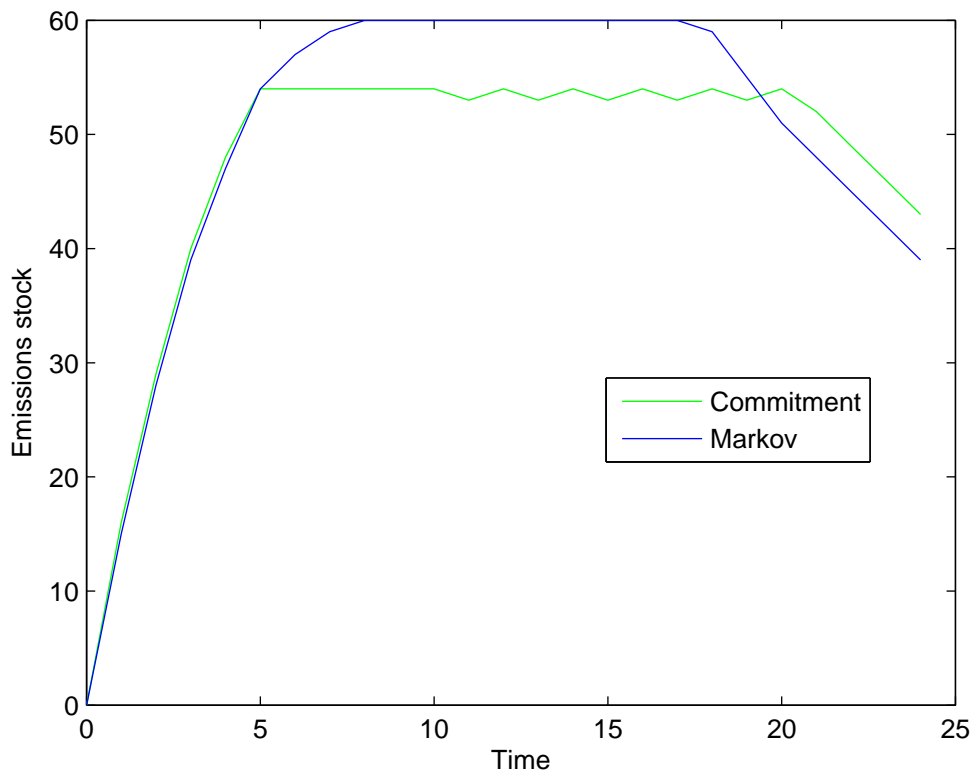


Figure 4: Carbon stocks $D(j)$ in commitment and Markov solutions

are expected to become obsolete or exhausted. In order to enact current preferences, regulators can reduce fossil fuel exploration and increase investments in clean alternatives.

A Proof of Lemma 1

Proof. First, consider iso-elastic utility $u_t(z_t) = \frac{z_t^{1-\eta}}{1-\eta}$. If the catastrophe has not occurred by the start of the second period, we have

$$w_2^M = \frac{(z_2^M)^{1-\eta}}{1-\eta} - \xi \frac{z_1^M + z_2^M}{\bar{D} - z_1^M} \Leftrightarrow (z_2^M)^{-\eta} = \frac{\xi}{\bar{D} - z_1^M} \Leftrightarrow r(z_1) = \left(\frac{\bar{D} - z_1}{\xi} \right)^{\frac{1}{\eta}}$$

Substituting in (1a), we obtain

$$w_1^M = \frac{(z_1^M)^{1-\eta}}{1-\eta} + \rho \frac{\left(\frac{\bar{D} - z_1}{\xi} \right)^{\frac{1-\eta}{\eta}}}{1-\eta} \left(1 - \frac{z_1^M}{\bar{D}} \right) + \frac{z_1^M \rho \underline{u}}{\bar{D}} - \xi \frac{z_1^M + \left(\frac{\bar{D} - z_1}{\xi} \right)^{\frac{1}{\eta}}}{\bar{D}}$$

The associated first-order condition is

$$(z_1^M)^{-\eta} - \frac{\rho}{\bar{D}} \left(\frac{\left(\frac{\bar{D} - z_1}{\xi} \right)^{\frac{1-\eta}{\eta}}}{1-\eta} - \underline{u} \right) + \frac{1-\rho}{\eta \bar{D}} \left(\frac{\bar{D} - z_1}{\xi} \right)^{\frac{1-\eta}{\eta}} - \frac{\xi}{\bar{D}} = 0 \quad (19)$$

Conversely, in the commitment outcome second-period emissions satisfy

$$w_2^C = \rho \frac{(z_2^C)^{1-\eta}}{1-\eta} - \xi \frac{z_1^C + z_2^C}{\bar{D} - z_1^C} \Leftrightarrow \rho (z_2^C)^{-\eta} = \frac{\xi}{\bar{D} - z_1^C} \Leftrightarrow z_2^C = \left(\frac{\rho(\bar{D} - z_1^C)}{\xi} \right)^{\frac{1}{\eta}}$$

This gives us

$$w_1^C = \frac{(z_1^C)^{1-\eta}}{1-\eta} + \rho \frac{\left(\frac{\rho(\bar{D} - z_1^C)}{\xi} \right)^{\frac{1-\eta}{\eta}}}{1-\eta} \left(1 - \frac{z_1^C}{\bar{D}} \right) + \frac{z_1^C \rho \underline{u}}{\bar{D}} - \xi \frac{z_1^C + \left(\frac{\rho(\bar{D} - z_1^C)}{\xi} \right)^{\frac{1}{\eta}}}{\bar{D}}$$

and FOC

$$(z_1^C)^{-\eta} - \frac{\rho}{\bar{D}} \left(\frac{\left(\frac{\rho(\bar{D} - z_1^C)}{\xi} \right)^{\frac{1-\eta}{\eta}}}{1-\eta} - \underline{u} \right) - \frac{\xi}{\bar{D}} = 0 \quad (20)$$

It can be shown that the left-hand side of (19) is larger than the left-hand side of (20) for all z_1 and $\rho \in (0, 1)$. Therefore, $z_1^M > z_1^C$.

Now consider quadratic utility $u_t(z_t) = az_t - \frac{1}{2}bz_t^2$. If the catastrophe has not occurred by the start of the second period, the second generation's welfare is

$$w_2^M = az_2^M - \frac{1}{2}b(z_2^M)^2 - \xi \frac{\frac{z_2^M}{\bar{D}}}{1 - \frac{z_1^M}{\bar{D}}} \Leftrightarrow a - bz_2^M - \frac{\xi}{\bar{D}\left(1 - \frac{z_1^M}{\bar{D}}\right)} = 0 \Leftrightarrow r(z_1) = \frac{a(\bar{D} - z_1) - \xi}{b(\bar{D} - z_1)}$$

Substituting in (1a), we obtain

$$\begin{aligned} w_1^M = & az_1^M - \frac{1}{2}b(z_1^M)^2 + \rho \left(a \left(\frac{a(\bar{D} - z_1) - \xi}{b(\bar{D} - z_1)} \right) - \frac{1}{2}b \left(\frac{a(\bar{D} - z_1) - \xi}{b(\bar{D} - z_1)} \right)^2 \right) \left(1 - \frac{z_1^M}{\bar{D}} \right) \\ & + \frac{z_1^M \rho \underline{u}}{\bar{D}} - \xi \frac{z_1 + \frac{a(\bar{D} - z_1) - \xi}{b(\bar{D} - z_1)}}{\bar{D}} \end{aligned}$$

The first-order condition is

$$a - bz_1^M - \frac{1}{2}\rho \frac{(a^2 - 2b\underline{u})(\bar{D} - z_1^M)^2 - \xi^2}{b\bar{D}(\bar{D} - z_1^M)^2} + \frac{\xi^2(1 - \rho)}{b\bar{D}(\bar{D} - z_1^M)^2} - \frac{\xi}{\bar{D}} = 0$$

In the commitment outcome, the first generation chooses z_2^C to maximize

$$w_2^C = \rho \left(az_2^C - \frac{1}{2}b(z_2^C)^2 \right) - \xi \frac{\frac{z_2^C}{\bar{D}}}{1 - \frac{z_1^C}{\bar{D}}} \Leftrightarrow \rho(a - bz_2^C) - \frac{\xi}{\bar{D}\left(1 - \frac{z_1^C}{\bar{D}}\right)} = 0 \Leftrightarrow z_2^C = \frac{\rho a(\bar{D} - z_1^C) - \xi}{\rho b(\bar{D} - z_1^C)}$$

The first generation's welfare is then

$$\begin{aligned} w_1^C = & az_1^C - \frac{1}{2}b(z_1^C)^2 + \rho \left(a \left(\frac{\rho a(\bar{D} - z_1^C) - \xi}{\rho b(\bar{D} - z_1^C)} \right) - \frac{1}{2}b \left(\frac{\rho a(\bar{D} - z_1^C) - \xi}{\rho b(\bar{D} - z_1^C)} \right)^2 \right) \left(1 - \frac{z_1^C}{\bar{D}} \right) \\ & + \frac{z_1^C \rho \underline{u}}{\bar{D}} - \xi \frac{z_1^C + \frac{\rho a(\bar{D} - z_1^C) - \xi}{\rho b(\bar{D} - z_1^C)}}{\bar{D}} \end{aligned}$$

giving rise to the following first-order condition

$$a - bz_1^C - \frac{1}{2} \frac{(\bar{D} - z_1^C)^2 (a^2 - 2b\underline{u}) \rho^2 - \xi^2}{\rho b \bar{D} (\bar{D} - z_1^C)^2} = 0 \quad (21)$$

Letting $z_1^C = z_1^M = z_1$, we have

$$\frac{\partial w_1^C}{\partial z_1} - \frac{\partial w_1^M}{\partial z_1} = \frac{1}{2} \frac{\xi^2(1 - \rho)^2}{\rho b \bar{D} (\bar{D} - z_1)^2} > 0$$

Therefore, for quadratic utility $z_1^C > z_1^M$. \square

B Proof of Proposition 1

Proof. From (10) it is apparent that if $z^C(t) = \alpha D(t)$ for some t , we must also have $z^C(s) = \alpha D(s) \forall s > t$. Otherwise, the first generation could improve its welfare by choosing $z^C(t) > \alpha D(t)$, as the current value cost of triggering a catastrophe is lower at t than at s . Moreover, the emissions stock must stabilize at some finite level because $\lim_{z \rightarrow \infty} u'(z) = 0$, $\lim_{D \rightarrow \infty} \psi(D) \gg 0$ and the monotonicity of $D(s)$ along the optimal path. Combining the above observations, there exists some t' such that $D^C(t') = D^C$ and $z^C(t) = \alpha D^C \forall t \geq t'$.

Now consider the alternative problem

$$\begin{aligned} \max_z \tilde{V}^C(t', D(t')) &= \mathbb{E} \left(\int_{t'}^{\infty} (-\tilde{\gamma}(s) + u(z(s))) e^{-r(s-t')} ds \right) \\ \text{s.t. } \dot{D} &= z - \alpha D, \quad D(t') = D_{t'}, \quad \tilde{\gamma}(\tau^+) = \tilde{\gamma}(\tau^-) + \bar{u} - \underline{u} + r\xi e^{rt'} \end{aligned} \quad (22)$$

The above problem has the same solution as (10) evaluated at $(t', D(t'))$ except that $\tilde{\gamma}(\tau^+) - \tilde{\gamma}(\tau^-)$ in (22) remains constant over time, whereas $\hat{\gamma}(\tau^+) - \hat{\gamma}(\tau^-)$ in (10) does not. The derivatives of $V^C(t', D(t'))$ and $\tilde{V}^C(t', D(t'))$ with respect to $z(t')$ have the same sign. Because (22) is stationary, we can analyze its steady state, assuming it is approached by a path that is monotonically increasing according to Definition 1. t' and $D(t') = D^C$ satisfy the conditions in the proposition text if and only if $z = \alpha D^C$ is the optimal steady-state policy in (22). Letting \tilde{v} denote the costate variable for $\tilde{V}(t, D(t))$, the steady-state conditions are

$$\dot{D} = z - \alpha D = 0 \quad (23a)$$

$$\dot{\tilde{\mu}} = (r + a)\tilde{\mu} + \psi(D)(z - \alpha D) - \psi(D)\alpha \left(\tilde{v} - \frac{\underline{u}}{r} + \xi e^{rt'} \right) = 0 \quad (23b)$$

$$\dot{\tilde{v}} = r\tilde{v} - u(z) + \psi(D)(z - \alpha D) \left(\tilde{v} - \frac{\underline{u}}{r} + \xi e^{rt'} \right) = 0 \quad (23c)$$

$$u'(z) + \tilde{\mu} - \psi(D) \left(\tilde{v} - \frac{\underline{u}}{r} + \xi e^{rt'} \right) = 0 \quad (23d)$$

Solving (23) for D , $\tilde{\mu}$, \tilde{v} and z yields

$$u'(\alpha D) = \frac{\psi(\alpha D)}{r + \alpha} \left[u(\alpha D) - \underline{u} + r\xi e^{rt'} \right]$$

Therefore, t' and D^C must satisfy (11). \square

C Proof of Proposition 2

Proof. By the argument in the main text, the steady-state carbon stock cannot exceed D^N . Consider a generation t that inherits carbon stock $D(t) < D^N$. Let

$D^{t,N}(t')$ and $z^{t,N}(t')$ denote the carbon stock and fossil fuel use respectively at time $t' > t$ in generation t 's preferred path. Suppose that $D^{t,N} = D^N$ and $D^{t,N}(t') = D^{t,N}$.⁶ Analogous to the proof of Proposition 1, it can only be optimal to choose $z^{t,N}(t') = \alpha D^N$ iff

$$u'(\alpha D^N) = \frac{\psi(\alpha D^N)}{r + \alpha} \left[u(\alpha D^N) - \underline{u} + r\xi e^{r(t'-t)} \right] \quad (24)$$

If (13) holds at D^N , the right hand side of (24) exceeds the left hand side at $D^{t,N} = D^N$ since $t' > t$. By Assumptions 1 and 2, we must therefore have $D^{t,N} < D^N$.

We complete the proof of $\lim_{t \rightarrow \infty} D^{t,N} = D^N$ by noting that whenever $D^{t,N} < D^N$ and $D^{t,N}(t') = D^{t,N}$, generation $t' > t$ prefers $D^{t',N} > D^{t,N}$. $D^{t,N}(t') = D^{t,N}$ implies

$$u'(\alpha D^{t,N}) = \frac{\psi(\alpha D^{t,N})}{r + \alpha} \left[u(\alpha D^{t,N}) - \underline{u} + r\xi e^{r(t'-t)} \right] \quad (25)$$

If $D^{t',N} = D^{t,N}$, we must also have

$$u'(\alpha D^{t,N}) = \frac{\psi(\alpha D^{t,N})}{r + \alpha} \left[u(\alpha D^{t,N}) - \underline{u} + r\xi \right] \quad (26)$$

Clearly, (25) and (26) cannot hold simultaneously. When (25) holds, the left hand side of (26) is larger than the right hand side at $D^{t,N}$. Generation t' will therefore choose $D^{t',N} > D^{t,N}$, so $z^{t',N}(t') > \alpha D^{t,N}$. As the carbon stock approaches D^N , the target levels $D^{t,N}$ must also approach D^N . \square

D Proof of Corollary 1

Suppose that D^C is reached for the first time at time $t' > 0$. Then it can only be optimal to choose $z(t') = \alpha D^C$ iff

$$u'(\alpha D^C) = \frac{\psi(\alpha D^C)}{r + \alpha} \left[u(\alpha D^C) - \underline{u} + r\xi e^{rt'} \right]$$

If (13) holds at D^N , the right hand side of the above equation exceeds the left hand side at $D^C = D^N$. By Assumptions 1 and 2, we must therefore have $D^C < D^N$.

⁶If $D^{t,N} = D^N$ but D_t does not reach D^N in finite time, a modified version of the below argument still applies: for t' arbitrarily large and ϵ arbitrarily small, the left hand side of (24) evaluated at $D^N - \epsilon$ is larger than the right hand side.

E Proof of Proposition 3

Proof. Recall that D_1^M and D_2^M are unique by Proposition 2. We verify that the equilibria in the proposition text satisfy the equilibrium conditions. Let t be sufficiently large and suppose that generation t believes that future generations will follow (17) and $D \geq D^M$. Then generation t believes that if it increases the emissions stock, future generations will keep the emissions stock constant.

First, consider the case in which $D^M < D^N$. By Proposition 2, generation t would prefer to reach a higher steady-state emissions stock in the naive solution, that is if it could commit all fossil fuel use from t onward. We show that this implies that in the Markov solution, generation t will choose $z > \alpha D$. When t is sufficiently large, $D^{t,N}$ is arbitrarily close to D^N . Furthermore, in generation t 's preferred path $z^{t,N}(s)$, $D^{t,N}$ is reached in finite time. This means there exists a $t' > t$ such that

$$(1 - \alpha) D^{t,N}(t') + z^{t,N}(t') = D^{t,N}$$

and

$$\left. \frac{\partial V^{t,N}(t', D^{t,N}(t'))}{\partial z^{t,N}(t')} \right|_{z^{t,N}(t') = D^{t,N} - (1 - \alpha) D^{t,N}(t')} = 0 \quad (27)$$

The interpretation of (27) is that, at $t' > t$ and $D^{t,N}(t') > D(t)$, generation t would choose to increase the emissions stock by $D^{t,N} - (1 - \alpha) D^{t,N}(t')$ if the emissions stock would remain constant in all subsequent periods. But then by Assumption 1 and Assumption 2, it must be welfare-improving to increase the emissions stock by the same amount at $(t, D(t))$, given that future generations keep the stock constant at the new level: the marginal utility of consumption is higher, the hazard rate is lower and the current-value cost of a catastrophe is lower. Therefore, $D^M < D^N$ cannot be an equilibrium.

Now turn to the decisions of early generations that inherit a carbon stock $D(t) < D^M$. If generation t believes that subsequent generations will follow (17), it realizes that its actions will not affect the maximum carbon stock D^M . When all future generations also believe the maximum stock equals D^M , the preferences of all generations that inherit $D(t) < D^M$ are no longer time-inconsistent. Then the problem of generation t reduces to maximizing the integral of expected discounted utility subject to $D(s) \leq D^M$, i.e.

$$\begin{aligned} & \max_z \mathbb{E} \left(\int_t^\infty (-\gamma(s) + u(z(s))) e^{-r(s-t)} ds \right) \\ & \text{s.t. } \dot{D} = z - \alpha D, \quad D(t) = D, \quad D(s) \leq D^M \quad \forall s \geq t, \quad \gamma(\tau^+) = \gamma(\tau^-) + \bar{u} - \underline{u} \end{aligned} \quad (28)$$

The solution to this optimal control problem coincides with the Markov solution. Analogous to Proposition 1, the steady-state of the unconstrained version of (28) satisfies

$$u'(\alpha D) = \frac{\psi(D)}{r + \alpha} (u(\alpha D) - \underline{u})$$

Therefore, carbon stocks larger than D_2^M are never visited in equilibrium. \square

F Proof of Proposition 4

Proof. Using the results from Propositions 1, 2 and 3, we can rewrite (10), (12) and (15) as constrained optimization problems

$$\begin{aligned} \max_z V^k(t, D(t)) &= \mathbb{E} \left(\int_t^\infty (-\hat{\gamma}(s) + u(z(s))) e^{-r(s-t)} ds \right) \\ \text{s.t. } \dot{D} &= z - \alpha D, \quad D(t) = D_t, \quad D(s) \leq D^k \quad \forall s \\ \hat{\gamma}(\tau^+) &= \hat{\gamma}(\tau^-) + \bar{u} - \underline{u}, \quad k \in \{C, \{t, N\}, M\} \end{aligned} \quad (29)$$

where $D^C < D^{t,N} < D^M$ for $0 < t < \infty$. We can represent the optimal strategy in each solution as $z = \zeta^k(D) = \zeta(D; D^k)$, $k \in \{C, \{t, N\}, M\}$, where $\zeta^C(D)$ is only optimal along the equilibrium path. $D^C < D^{t,N} < D^M$ implies $\zeta^C(D) < \zeta^{t,N}(D) < \zeta^M(D)$ if and only if $\frac{\partial \zeta(D; D^k)}{\partial D^k} > 0$. Let

$$\begin{aligned} V(D; D^k) &= \max_z \mathbb{E} \left(\int_t^\infty (-\hat{\gamma}(s) + u(z(s))) e^{-r(s-t)} ds \right) \\ \text{s.t. } \dot{D} &= z - \alpha D, \quad D(t) = D_t, \quad D(s) \leq D^k \quad \forall s \\ \hat{\gamma}(\tau^+) &= \hat{\gamma}(\tau^-) + \bar{u} - \underline{u} \end{aligned}$$

be the value of continuing optimally from carbon stock D subject to $D(s) \leq D^k \quad \forall s$. Writing $V = V(D; D^k)$, the HJB equation and the first order condition from the Hamiltonian stipulate

$$rV = \max_z \left\{ u(z) + V_D(z - \alpha D) - \psi(D)(z - \alpha D) \left(V - \frac{\underline{u}}{r} \right) \right\} \quad (30)$$

$$u'(z) + V_D - \psi(D) \left(V - \frac{\underline{u}}{r} \right) = 0 \quad (31)$$

By (30), along the optimal path

$$V_D = \frac{rV - u(z)}{z - \alpha D} + \psi(D) \left(V - \frac{\underline{u}}{r} \right) \quad (32)$$

Substituting (32) in (31), we obtain

$$\begin{aligned}
u'(z) + \frac{rV - u(z)}{z - \alpha D} &= 0 \\
\Leftrightarrow (z - \alpha D) u'(z) + rV - u(z) &= 0 \\
\Leftrightarrow \tilde{z} u'(\tilde{z} + \alpha D) + rV - u(\tilde{z} + \alpha D) &= 0
\end{aligned}$$

Totally differentiate with respect to D^k

$$\begin{aligned}
\frac{\partial \tilde{z}}{\partial D^k} u'(\tilde{z} + \alpha D) + \tilde{z} \frac{\partial \tilde{z}}{\partial D^k} u''(\tilde{z} + \alpha D) + r \frac{\partial V}{\partial D^k} - \frac{\partial \tilde{z}}{\partial D^k} u'(\tilde{z} + \alpha D) &= 0 \\
\Leftrightarrow \frac{\partial \tilde{z}}{\partial D^k} \underbrace{\tilde{z} u''(\tilde{z} + \alpha D)}_{<0} + \underbrace{r \frac{\partial V}{\partial D^k}}_{>0 \forall D^k < D_2^M} &= 0
\end{aligned}$$

By the above, we must have $\frac{\partial \tilde{z}}{\partial D^k} > 0$. Having established $\zeta^C(D) < \zeta^{t,N}(D) < \zeta^M(D) \forall D$, it automatically follows that $D^C(t) < D^N(t) < D^M(t)$. \square

G Time lags

In the analysis in the main text, we have assumed that the catastrophe occurs instantaneously when the carbon stock reaches the threshold \hat{D} . The reaction of temperatures to emissions is quite sluggish however: current emissions do not substantially affect global temperature levels until a few decades from now [7]. In this section, we incorporate this time lag by letting the catastrophe occur with an exogenous delay. As a consequence, generations only observe whether the lagged (rather than the current) carbon stock has triggered a catastrophe.

G.1 Two-period model

We modify the timing in section 2 as follows. When $D_1 \geq \hat{D}$, the catastrophe is not triggered at the end of period 1. Generation 2 does not observe whether $D_1 \geq \hat{D}$, so it uses the unconditional rather than the conditional probability distribution of \hat{D} . The catastrophe only occurs at the end of the second period if $D_2 \geq \hat{D}$. The welfare functions are as follows

$$w_1 = u_1(z_1) + \rho u_2(z_2) - \xi F(z_1 + z_2) \quad (33a)$$

$$w_2 = u_2(z_2) - \xi F(z_1 + z_2) \quad (33b)$$

giving rise to first-order conditions

$$u'_1(z_1^C) - \xi f(z_1^C + z_2^C) = 0 \quad (34a)$$

$$\rho u'_2(z_2^C) - \xi f(z_1^C + z_2^C) = 0 \quad (34b)$$

in the commitment solution and

$$u'_1(z_1^N) - \xi f(z_1^N + z_2^C) = 0 \quad (35a)$$

$$u'_2(z_2^N) - \xi f(z_1^N + z_2^N) = 0 \quad (35b)$$

in the naive solution. Comparing these FOCs to (2a) and (13), the fact that the second generation uses the unconditional pdf causes second-period emissions conditional on z_1 to be higher than if the catastrophe would occur instantaneously. The second generation must contend with the possibility that $z_1 \geq \hat{D}$, in which case it should choose a high level of consumption as even very cautious strategies cannot avert the catastrophe. The first generation in turn is no longer concerned that its fossil fuel consumption will reduce the second generation's utility to \underline{u} when $z_1 \geq \hat{D}$, so the $\rho f(z_1)[u_2(z_2) - \underline{u}]$ term is not present in (34) and (35). The Markov solution shows a similar picture. The second generation's reaction function satisfies

$$\begin{aligned} u'_2(r(z_1^M)) &= \xi f(z_1^M + r(z_1^M)) & (36) \\ u''_2(r(z_1^M)) r'(z_1) &= \xi f'(z_1 + r(z_1)) (1 + r'(z_1)) \\ r'(z_1) &= \frac{\xi f'(z_1 + r(z_1))}{u''_2(r(z_1)) - \xi f'(z_1 + r(z_1))} \end{aligned}$$

The second generation's reaction function is less likely to be downward sloping than in section 2. Higher emissions in the first period increase the probability that the catastrophe has already been triggered by z_1 , which makes the second generation less inclined to curb its emissions. The first generation's FOC is

$$\begin{aligned} u'_1(z_1) + \rho u'_2(r(z_1)) r'(z_1) - \xi f'(z_1 + r(z_1)) (1 + r'(z_1)) &= 0 \\ \Leftrightarrow u'_1(z_1) + \rho \xi f'(z_1 + r(z_1)) - \xi f'(z_1 + r(z_1)) (1 + r'(z_1)) &= 0 \\ \Leftrightarrow u'_1(z_1) + \xi (1 - \rho) f'(z_1 + r(z_1)) r'(z_1) - \xi f'(z_1 + r(z_1)) &= 0 \end{aligned}$$

The first transformation follows by substituting (36). The FOC is similar to (6) except for the absence of the second term.

G.2 Infinite horizon, abundant fossil fuel

We implement the the time lag in the infinite-horizon model in a similar way. The catastrophe occurs ℓ periods after $D = \hat{D}$, that is when $D_{t-\ell} = \hat{D}$. Generations between $t - \ell$ and t do not observe whether $D_{t-\ell} \geq \hat{D}$. The welfare

function for generation t is then

$$V_\ell(t, D) = \begin{cases} \mathbb{E} \left(\int_t^\infty (-\gamma(s) + u(z(s))) e^{-r(s-t)} ds \right) - \xi \frac{P[\tau \in [t, \infty)]}{1 - P[\tau \in [0, t]]} \\ \text{s.t. } \dot{D} = z - \alpha D, \quad D(0) = D_0, \quad \gamma(\tau^+) = \gamma(\tau^-) + \bar{u} - \underline{u} + r\xi e^{-r(\tau-t)} \\ \tau \sim \psi(x(\tau - \ell), z(\tau - \ell)) g(x(\tau - \ell), z(\tau - \ell)) \\ \exp\left(-\int_0^{\tau-\ell} \psi(x(s)) g(x(s), z(s)) ds\right) \\ \frac{\underline{u}}{r} - \xi \end{cases} \begin{array}{l} \text{for } \tau \notin [0, t] \\ \\ \\ \text{for } \tau \in [0, t] \end{array} \quad (37)$$

To have an interesting problem, we assume it is not optimal to choose $z_t \rightarrow \infty$ for the first ℓ periods and accept that the catastrophe occurs with certainty afterwards. This assumption holds when the delay ℓ is not too large.⁷

Assumption 3. *There exists a $\tilde{z} < \infty$ such that*

$$\sup \{V_\ell(t, D) : z(s) \in \mathbb{R}_+\} = \max_{z(s)} \{V_\ell(t, D) \text{ s.t. } z(s) < \tilde{z} \forall s\}$$

The time lag gives rise to higher steady-state carbon stocks because it mitigates the negative impact of a catastrophe: the utility of generations between $t - \ell$ and the time of occurrence t is not reduced to \underline{u} .

Proposition 6. *Assume a solution to (37) exists. Then the commitment solution is characterized by a steady-state emissions stock D_ℓ^C . There exists a $t' < \infty$ such that $D(t') = D_\ell^C$ and $z_\ell^C(t) = \alpha D_\ell^C \forall t \geq t'$. D_ℓ^C and t' satisfy*

$$u'(\alpha D_\ell^C) = \frac{\psi(\alpha D_\ell^C) e^{-r\ell}}{r + \alpha} \left[u(\alpha D_\ell^C) - \underline{u} + r\xi e^{r(t'+\ell)} \right] \quad (38)$$

Proof. The proof is analogous to Proposition 1. The action at time t only depends on the history of emissions through the current emissions stock $D(t)$ because the control at time t affects generations from $t + \ell$ onwards if and only if the catastrophe has not occurred by time $t + \ell$, which depends on the history only through the current stock. The steady-state conditions to the alternative problem

$$\begin{aligned} \max_z \tilde{V}_\ell^C(t', D(t')) &= \mathbb{E} \left(\int_{t'}^\infty (-\tilde{\gamma}(s) + u(z(s))) e^{-r(s-t')} ds \right) \\ \text{s.t. } \dot{D} &= z - \alpha D, \quad D(t') = D_{t'}, \quad \tilde{\gamma}(\tau^+) = \tilde{\gamma}(\tau^-) + \bar{u} - \underline{u} + r\xi e^{r(t'+\ell)} \end{aligned}$$

⁷A sufficient, albeit strong condition is $\frac{1-e^{-r\ell}}{r} [\bar{u} - u(0)] - \frac{e^{-r\ell}}{r} [u(0) - \underline{u}] - \xi < 0$.

which has the same solution as (37) at (t', D_ℓ^C) , are

$$z - \alpha D = 0 \quad (39a)$$

$$(r + a) \tilde{\mu} + \psi(D) (z - \alpha D) - \psi(D) \alpha \left(\tilde{v} - \frac{u}{r} + \xi e^{r(t'+\ell)} \right) e^{-r\ell} = 0 \quad (39b)$$

$$r\tilde{v} - u(z) + \psi(D) (z - \alpha D) \left(\tilde{v} - \frac{u}{r} + \xi e^{r(t'+\ell)} \right) e^{-r\ell} = 0 \quad (39c)$$

$$u'(z) + \tilde{\mu} - \psi(D) \left(\tilde{v} - \frac{u}{r} + \xi e^{r(t'+\ell)} \right) e^{-r\ell} = 0 \quad (39d)$$

Solving the above system, the steady-state emissions stock satisfies

$$u'(\alpha D_\ell^C) = \frac{\phi(D_\ell^C) e^{-r\ell}}{r + \alpha} \left[u(\alpha D_\ell^C) - \underline{u} + r\xi e^{r(t'+\ell)} \right]$$

□

The naive solution is described in Proposition 7.

Proposition 7. *Assume a solution to generation t 's problem*

$$\max_z V_\ell^{t,N}(t, D(t)) = \mathbb{E} \left(\int_t^\infty (-\hat{\gamma}(s) + u(z(s))) e^{-r(s-t)} ds \right)$$

s.t. $\dot{D} = z - \alpha D$, $D(t) = D_t$, $\hat{\gamma}(\tau^+) = \hat{\gamma}(\tau^-) + \bar{u} - \underline{u} + r\xi e^{r(\tau-t+\ell)}$

exists. Then this solution is characterized by a steady-state emissions stock $D_\ell^{t,N}$. Let D_ℓ^N be given by

$$u'(\alpha D_\ell^N) = \frac{\phi(D_\ell^N) e^{-r\ell}}{r + \alpha} \left[u(\alpha D_\ell^N) - \underline{u} + r\xi e^{r\ell} \right] \quad (40)$$

Then $D_\ell^{t,N} < D_\ell^N \forall t$ and $\lim_{t \rightarrow \infty} D_\ell^{t,N} = D_\ell^N$.

Proof. Analogous to Proposition 2. □

Lastly, we characterize the Markov solution.

Proposition 8. *Let $D_{\ell,1}^M = D_\ell^N$ and $D_{\ell,2}^M$ be given by*

$$u'(\alpha D_{\ell,2}^M) = \frac{\phi(D_{\ell,2}^M) e^{-r\ell}}{r + \alpha} \left[u(\alpha D_{\ell,2}^M) - \underline{u} \right] \quad (41a)$$

Assume a solution to

$$V_\ell^M(t, D) = \max_z \mathbb{E} \left(\int_t^\infty (-\gamma(s) + u(z(s))) e^{-r(s-t)} ds \right)$$

s.t. $\dot{D} = z - \alpha D$, $D(t) = D$, $D(s) \leq D_\ell^M \forall s \geq t$, $\gamma(\tau^+) = \gamma(\tau^-) + \bar{u} - \underline{u}$

exists $\forall x : D \leq D_\ell^M$, $D_\ell^M \in [D_{\ell,1}^M, D_{\ell,2}^M]$. Then there exists a continuum of Markov equilibria indexed by $D_\ell^M \in [D_{\ell,1}^M, D_{\ell,2}^M]$ such that

$$\zeta^M(D) = \begin{cases} \operatorname{argmax}_z V_\ell^M(t, D) & \text{if } D < D_\ell^M \\ \alpha D & \text{if } D \geq D_\ell^M \end{cases} \quad (42)$$

Proof. Assumption 3 rules out non-local equilibrium conditions⁸. The remainder of the proof follows that of Proposition 3. \square

H Piecewise deterministic optimal control

Consider a random variable ε with probability density function $f(\varepsilon)$ defined on $[0, \infty)$ and cumulative density function $F(\varepsilon)$. Denote the actual value of ε by $\tilde{\varepsilon}$. The hazard rate of ε is $\psi(\varepsilon) \equiv \frac{f(\varepsilon)}{1 - \int_0^\varepsilon f(\eta) d\eta}$. Let $x \in X \subseteq \mathbb{R}^n$ denote the vector of state variables and define a threshold function $\Phi(x, \varepsilon) = 0$. The catastrophe occurs when $\Phi(x, \tilde{\varepsilon}) = 0$. We assume $\frac{\partial \Phi}{\partial x_i} \geq 0$, $i = 1, \dots, n$ and $\frac{\partial \Phi}{\partial \varepsilon} \leq 0$: higher values of the state variables bring the system 'closer' to the threshold, and higher values of $\tilde{\varepsilon}$ imply a higher threshold. Define $\phi : X \rightarrow \mathbb{R}_+$ as $\{\varepsilon : \Phi(x, \varepsilon) = 0, x \in X\}$. $\phi(x)$ is the value of ε such that the threshold is reached when the state variables take on value x . Because of our assumptions on the partial derivatives of Φ , $\phi'(x) \geq 0$.

Definition 1. Let $x : \mathbb{R}_+ \rightarrow X$ be continuous and differentiable almost everywhere. $x(t)$ is monotonically increasing with respect to $\Phi(x(t), \varepsilon) = 0$ and ε if and only if for any t_0 and t_1 such that $t_0 < t_1$ it holds that

$$\Phi(x(t_0), \varepsilon_0) = \Phi(x(t_1), \varepsilon_1) \Leftrightarrow \varepsilon_0 \leq \varepsilon_1$$

For trajectories of the state variables $x(t)$ that are monotonically increasing with respect to $\Phi(x(t), \varepsilon) = 0$, $\phi(x(t))$ increases over time. From here on, we restrict attention to such trajectories, as trajectories with decreasing state variables will not be optimal. Then the occurrence time of the catastrophe τ is a Poisson process:

$$\tau \sim f(\varphi(x(\tau))) \varphi'(x(\tau)) x'(\tau)$$

Again following [11], we model the catastrophe as a discrete jump in the state variables. (**author?**) argues that this approach is more general than a discrete jump in instantaneous utility, because the latter can always be modeled as the

⁸TM: Non-trivial?

former, but not the other way around. When the catastrophe occurs at time τ , the jump in the state variables is given by

$$x(\tau^+) = Q(x(\tau^-)) = x(\tau^-) + q(x(\tau^-)) \quad (43)$$

where $x(\tau^-) = \lim_{t \uparrow \tau} x(t)$ and $x(\tau^+) = \lim_{t \downarrow \tau} x(t)$. [11] shows that expected discounted utility can be maximized by solving the following problem

$$\begin{aligned} W(t, x(t)) &= \max_z \mathbb{E} \left(\int_0^\infty f(x, z) e^{-rt} dt \right) \text{ s.t. } \dot{x} = g(x, z), x(0) = x_0 \\ x(\tau^+) &= x(\tau^-) + q(x(\tau^-)) \\ \tau &\sim \psi(x(\tau), z(\tau)) g(x(\tau), z(\tau)) \exp \left(- \int_0^\tau \psi(x(s)) g(x(s), z(s)) ds \right) \end{aligned} \quad (44)$$

where we write $g(x, z)$ for $x'(t)$. The risk-augmented Hamiltonian for this problem is

$$\begin{aligned} H(x, \mu, z) &= u(x, z) + \mu g(x, z) + \psi(\phi(x)) \phi'(x) g(x, z) \\ &\quad \times [W(t, x + q(x) | \tau = t) - W(t, x)] \end{aligned} \quad (45)$$

where

$$W(t, x | \tau = t) = \max_z \int_t^\infty u(y, z) e^{-r(s-t)} ds \text{ s.t. } \dot{y} = g(y, z), y(t) = x \quad (46)$$

is the value of continuing optimally when the catastrophe occurs at time t and results in state x . For brevity, we will write $(\cdot | \tau)$ as shorthand for $(\cdot | \tau = t)$. The post-catastrophe problem is a standard deterministic control problem with costate variables $\mu(s, t | \tau)$. Note that $\frac{\partial}{\partial x} W(t, x | \tau) = \mu(t, t | \tau)$ and $\frac{\partial}{\partial x} W(t, x + q(x) | \tau) = (I^n + q'(x)) \mu(t, t | \tau)$, where I^n is the n -dimensional identity matrix and $q'(x)$ is the Jacobian of $q(x)$. Lastly, $J(t, x)$ in (45) is

$$\begin{aligned} W(t, x) &= \max_z \mathbb{E} \left(\int_t^\infty u(y, z) e^{-r(s-t)} ds \right) \text{ s.t. } \dot{x} = g(y, z), y(0) = x \\ x(\tau^+) &= x(\tau^-) + q(x(\tau^-)) \\ \tau &\sim \psi(x(\tau), z(\tau)) g(x(\tau), z(\tau)) \exp \left(- \int_0^\tau \psi(x(s)) g(x(s), z(s)) ds \right) \end{aligned} \quad (47)$$

The differential equation for $w = W(t, x(t))$ is then (see the Appendix in [11])

$$\dot{w} = rw - u(x, z) + \psi(\phi(x)) \phi'(x) g(x, z) (w - W(t, x + q(x) | \tau)) \quad (48)$$

The Hamiltonian (45) gives rise to the following conditions

$$u = \operatorname{argmax}_v H(x, \mu, v) \quad (49)$$

$$\begin{aligned} \dot{\mu} = r\mu - \frac{\partial}{\partial x} f(x, z) - \mu \frac{\partial}{\partial x} g(x, z) - \lambda(x) (\mu(t|t, x + q(x)) (I^n + q'(x)) - \mu) \\ - \lambda'(x) (W(t, x + q(x) | \tau) - w) \end{aligned} \quad (50)$$

Lastly, define the transversality conditions. If x is the optimal path, then for all admissible y and $\dot{y} = g(y, u)$, we must have

$$\lim_{t \rightarrow \infty} \mu e^{-rt} (y(t) - x(t)) \geq 0 \quad \lim_{t \rightarrow \infty} z(t) e^{-rt} = 0 \quad (51)$$

References

- [1] F. Alvarez-Cuadrado and N. V. Long. A mixed Bentham-Rawls criterion for intergenerational equity: Theory and implications. Journal of Environmental Economics and Management, 58(2):154–168, 2009.
- [2] Ujjayant Chakravorty, Bertrand Magn, and Michel Moreaux. A Hotelling model with a ceiling on the stock of pollution. Journal of Economic Dynamics and Control, 30(12):2875–2904, 2006.
- [3] Ujjayant Chakravorty, Michel Moreaux, and Mabel Tidball. Ordering the extraction of polluting nonrenewable resources. American Economic Review, 98(3):1128–1144, 2008.
- [4] Graciela Chichilnisky. An axiomatic approach to sustainable development. Social Choice and Welfare, 13(2):231–257, 1996.
- [5] Graciela Chichilnisky. An axiomatic approach to choice under uncertainty with catastrophic risks. Resource and Energy Economics, 22:221–231, 2000.
- [6] Reyer Gerlagh and Matti Liski. Carbon prices for the next thousand years. Mimeographed, 2012.
- [7] J.T. Houghton, Y. Ding, D.J. Griggs, M. Noguer, P.J. van der Linden, X. Dai, K. Maskell, and C.A. Johnson, editors. Climate Change 2001, Working Group 1: The Scientific Basis. Intergovernmental Panel on Climate Change, 2001.

- [8] Larry Karp. Hyperbolic discounting and climate change. Journal of Public Economics, 89(2-3):261–282, 2005.
- [9] K. Keller, M. B. Bolker, and D. F. Bradford. Uncertain climate thresholds and optimal economic growth. Journal of Environmental Economics and Management, 48(1):723–741, 2004.
- [10] Timothy M. Lenton, Hermann Held, Elmar Kriegler, Jim W. Hall, Wolfgang Lucht, Stefan Rahmstorf, and Hans Joachim Schellnhuber. Tipping elements in the Earth’s climate system. Proceedings of the National Academy of Sciences of the United States of America, 105(6):1786–1793, 2008.
- [11] Eric Nævdal. Dynamic optimisation in the presence of threshold effects when the location of the threshold is uncertain - with an application to the possible disintegration of the Western Antarctic Ice Sheet. Journal of Economic Dynamics and Control, 30:1131–1158, 2006.
- [12] Daniel C. Nepstad, Claudia M. Stickler, Britaldo Soares-Filho, and Frank Merry. Interactions among Amazon land use, forests and climate: prospects for a near-term forest tipping point. Philosophical Transactions of the Royal Society B, 363(1498):1737–1746, 2008.
- [13] William D. Nordhaus. Rolling the ‘DICE’: an optimal transition path for controlling greenhouse gases. Resource and Energy Economics, 15(1):27–50, 1993.
- [14] Gerard H. Roe and Marcia B. Baker. Why is climate sensitivity so unpredictable? Science, 318:629–632, 2007.
- [15] N.H. Stern. The Economics of Climate Change: The Stern Review. Cambridge University Press, Cambridge, UK, 2007.
- [16] David Tilman and John A. Downing. Biodiversity and stability in grasslands. Nature, 367:363–365, 1994.
- [17] Yacov Tsur and Amos Zemel. Endangered species and natural resource exploitation: Extinction vs. coexistence. Natural Resource Modeling, 8(4):389–413, 1994.
- [18] Yacov Tsur and Amos Zemel. Accounting for global warming risks: Resource management under event uncertainty. Journal of Economic Dynamics and Control, 20:1289–1305, 1996.

- [19] M. Weitzman. On modeling and interpreting the economics of catastrophic climate change. Review of Economics and Statistics, 91(1):1–19, 2009.
- [20] M. Weitzman. GHG targets as insurance against catastrophic climate damages. Mimeographed, 2010.