Access Econ

Submission Number: PET11-11-00222

The socially optimal level of altruism

Richard Adam Povey New College and Wadham College, Oxford University

Abstract

It is already recognized by some specific models in the existing literature that altruism may have socially counterproductive effects. Economic theory also shows that selfinterest often produces socially efficient outcomes. This paper explores the relationship between altruistic preferences and systems of punishment. The central argument is that altruism has detrimental effects on the efficacy of punishment and the resultant incentives of agents to co-operate with socially efficient equilibria. The sequential punishment model is presented - akin to an infinitely-repeated stage game, but sufficiently simple to allow determinate optimal penal codes to be derived - and the impact of different levels of altruism fully analysed. It is shown that high levels of altruistic motivation - close to but slightly less than full altruism - cause the socially efficient equilibrium to break down. Although the model is a stylized representation of social interaction, the key effects that drive these results should appear in many more specific examples. In summary, these are the temptation effect (more altruistic individuals are less tempted to do harm to others), the willingness effect (more altruistic individuals are less willing to inflict punishment), and the severity effect (punishments, such as a fine where the revenue is redistributed, are less severe for more altruistic individuals, because they place a higher value on the contribution of the revenue to the welfare of others).

Submitted: March 13, 2011.

The Socially Optimal Level of Altruism^{*}

Richard Povey, New College, Oxford University

December 30, 2010

[W]hen altruism improves static non-cooperative outcomes, it lessens the severity of credible punishments. An altruist may well be perceived as a "softy" and his threats may not be taken seriously. [Bernheim & Stark, 1988]

[T]he most efficient way to provide low payoffs, in terms of incentives to cheat, is to combine a grim present with a credibly rosy future. [Abreu, 1986]

1 Introduction - A Parable

The central problem of economic and social policy, indeed the essential prerequisite of the social order itself, is that of bestowing upon the individual agent the inventive to act in a manner which is beneficial for society as a whole. Such incentives can be *intrinsic* to the individual (altruistic preferences) or *extrinsic* (threats of punishment if the individual agent does not comply with the socially prescribed action). This paper analyses the interaction between these two alternative "technologies". We see that the two methods of achieving social order cannot be freely mixed at will, and that, in order for extrinsic incentives to work most effectively, it is necessary to limit the operation of intrinsic incentives, so as to avoid counterproductive interference between the two. An important implication is that the heavy (perhaps predominant) reliance upon extrinsic incentives in sophisticated human societies may in fact be a socially optimal "policy mix", rather than a mere second-best correction for the inadequate intrinsic motivation to "do the right thing".

Altruistic behaviour has been fruitfully modelled in economic theory using infinitely-repeated stage games [Fudenberg & Maskin, 1986] [Abreu, 1986], and infinite dynamic sequential games such as models with overlapping generations [Samuelson, 1958] [Hammond, 1975] [Cremer, 1986]. The traditional view is that apparent altruistic co-operation can occur in such models with selfinterested agents through the use of punishment equilibria which can deter actions which would be individually optimal but socially sub-optimal in a nonrepeated or non-sequential game. This paper provides a general workhorse model in which the results of such models can be extended in order to

^{*}I would like to express my gratitude to the United Kingdom Economic and Social Research Council and the Royal Economic Society for funding this research. I would also like to thank Kevin Roberts, Peter Hammond, Chris Wallace, Godfrey Keller, Alan Beggs, Scot Peterson and Michael Griebe for their helpful advice and comments on various drafts of this paper.

accommodate bona fide altruistic motivation. In common with the infinitely-repeated stage game model, agents in the model are infinitely-lived and discount the future. In common with overlapping generations models, on the other hand, players move sequentially, and each player only gets to move once during the entire game.¹

The sequential punishment model presented in this paper is intended to capture an abstract essential feature of the social and economic world in a simple but general manner. Other models with a similar idiom include the "Robinson Crusoe" economy [Ruffin, 1972], Samuelson's "pension game" [Samuelson, 1958] and Diamond's model of fiat money in a "coconut economy" [Diamond, 1984]. Each of these models can be illustrated intuitively with the help of a simple "parable", as can the sequential punishment model.

Consider a desert island where individuals are sufficiently settled to have established their own "back gardens". Each individual is sitting in their garden drinking a cold beer. One by one, at regular discrete intervals, one inhabitant finishes their drink, and must decide whether or not to walk to the bin, or to throw their bottle into one of their neighbour's gardens. All gardens are adjacent, and each person's bin is a variable distance away.² It is possible for each inhabitant to be threatened that, if they throw their bottle, everyone who subsequently finishes *their* beer will throw a bottle into the malefactor's garden. Sometimes this threat will be enough to make every inhabitant walk to the bin every time, leading to a socially efficient litter-free island. Sometimes the threat will not be enough, and some or all of the bottles will be thrown, leading to a socially inefficient outcome.

The central feature encapsulated in this model is the fundamental vicariousness of human social interaction. In any society, individuals are able to impose negative externalities upon one another for personal gain, in myriad ways. However, the very existence of this problem also offers a potential solution, in that it creates the possibility of punishing miscreants who take such opportunities, by threatening *them* with harm in the future.³ Another potential solution to the problem, however, is altruism; if people care about others, they will refrain from harming them.

This paper aims to show that these two alternative incentive technologies can interact with one another in a perverse manner. The question we will set out to answer is whether greater altruism on the part of the inhabitants of the island will always make it easier to achieve a litter-free island. The answer is a resounding and conclusive "no". In particular, it is shown that too high a level of altruism will in general lead to a worsening of the societal outcome. Altruism "dents" social incentive systems based upon extrinsic rewards or punishments.

2 Overview

It is commonly recognised that the repetition of stage games such as the prisoners' dilemma, or the Cournot and Bertand oligopoly games, allows

¹This may seem an unrealistic assumption, but it can be argued that we only need split a finite number of players into an infinite number of "egos" [Hammond, 1975].

 $^{^{2}}$ Imagine a giant pie-shaped island, each garden being a "wedge" - everyone is sitting at the middle of the island.

³Throughout the paper, we will use "harm" to refer to the inflicting of a negative externality and "punish" to the specific use of such harm opportunities to construct punishment equilibria.

apparently altruistic behaviour to be incentivized, resulting in a Pareto superior outcome for the players.⁴ However, such apparent altruism only reflects "enlightened self-interest". This paper introduces genuine altruistic motivation, exploring its effect on the achievability of a socially efficient equilibrum.

The well-known "Folk Theorem" establishes that, if players are sufficiently patient, any equilibrium which Pareto dominates the min-max payoff in the stage game can be supported as a subgame perfect equilibrium in the infinitely repeated game [Aumann & Shapley, 1992] [Rubinstein, 1979] [Fudenberg & Maskin, 1986]. The key implication is that, with imperfect altruism leading to a Pareto-inefficient equilibrium in the stage game, if there is sufficiently low discounting of the future then a Pareto efficient outcome in the infinitely-repeated game can be achieved. This paper asks the reverse question. Given a certain positive level of impatience, how high or low can the level of altruism be in order for a socially efficient outcome to be attained? It is shown that, in general, there must be a "Goldilocks" level of altruism, which is neither too high or too low.

In a game in which individuals make sequential moves (or in a repeated simultaneous-move game structure), they are able to punish one another based upon their previously observed behaviour. Individuals who are less altruistic are more willing to harm others because they place a lower value on the cost to the person being harmed. We shall therefore call this the **willingness effect**. Agents are also more afraid of being harmed because they value their own welfare more relative to that of others. For example, if a criminal is fined a certain amount and the revenue spent on other individuals, this is a more severe punishment for a less altruistic criminal because the fact that the fine revenue is spent on others mitigates the effect of the fine on the criminal's utility by a smaller amount than if the criminal were more altruistic.⁵ We shall therefore refer to this as the **severity effect**. Together, these two effects create a potential social benefit from individuals not being too altruistic. However, this must be balanced against the greater temptation towards wrongdoing by a less altruistic individual. Hence there is also a **temptation effect** from greater altruism.

Stark and Bernheim have observed that altruism can reduce the credibility of punishment if the potential punisher is perceived as a "softy" [Bernheim & Stark, 1988]. They observe that this can lead greater altruism to have a negative impact⁶ but argue that this must be analysed on a case-bycase basis. We establish a canonical framework in which to study the interaction of this willingness effect with the severity and temptation effects. This enables us to establish the stronger result that a high level of altruism will *in general* be socially detrimental. The central result of this paper is that, under certain fairly non-restrictive assumptions, the three effects conspire to render a socially efficient outcome impossible if the level of altruism becomes high enough.⁷

 $^{^{4}}$ This framework has important implications for such varied theoretical areas as environmental economics (international co-operation in pollution abatement can be seen as a repeated game) and competition policy (firms are sometimes able to collude in a repeated oligopoly game, to the detriment of overall social welfare) [Rees, 1993].

⁵The revenue from the fine could, of course, be "thrown away" in order to avoid this adverse effect, but this would be wasteful in that it would create a deadweight loss to punishment.

 $^{^6\}mathrm{This}$ is precisely the phenomenon we term the willingness effect.

 $^{^{7}}$ In terms of the island parable, the assumptions we must make are that the cost of having a bottle land is weakly greater than the cost of walking to a bin at the edge of the island, that

3 The Model

Suppose there are an infinite number of players and that distinct players, referenced by the period in which they move, each get a chance in sequence to impose damage upon another player. In period t, player t receives a **harm opportunity**, and must decide whether to accept or reject it, and a "target" for the harm opportunity, player A_t . If player t chooses to **accept** their harm opportunity , then player A_t suffers a **cost** in **felicity** of 1 unit. If player t accepts, then he gains felicity equal to the **benefit**, π_t ,⁸ which is drawn randomly and independently from a distribution defined by the density function $g(\pi)$, with support $[\hat{\pi}, 1]$, where $0 \leq \hat{\pi} < 1$.⁹ Therefore, $\int_{\hat{\pi}}^1 g(\pi) d\pi = 1$. All players publicly observe the value of π_t before player t moves.

The move made by player t in period t can be conceived to consist of choosing a **trigger level**, T_t , for the realized value of the benefit above which they will accept the harm opportunity¹⁰, and the individual whom they harm, A_t . A player's strategy maps the observed past history at period t (including the observed value of π_t) to the move they make.

We will derive results assuming a continuous distribution of the benefit with bounded support. We assume throughout that $g(\pi)$ is twice continuously differentiable. We will frequently use as an exemplar the case where the benefit has a continuous uniform distribution. Here the probability density function for the distribution of π_t will take the value $\frac{1}{1-\hat{\pi}}$ between $\hat{\pi}$ and 1. Most straightforwardly, with $\hat{\pi} = 0$, $g(\pi) = 1$ between 0 and 1.

4 Players' Preferences

Players do not act to maximise their "private utility", or **felicity**. Instead, players act to maximise their **social utility function**¹¹, which is a weighted sum of the felicities of all players.¹² We assume all players are risk-neutral and share the same discount rate δ . We let θ represent the weighting placed upon

players move in sequence, with perfect information about past play, and that all individuals share the same coefficient of altruism θ (the weighting they place on the welfare of others relative to their own welfare) and discount the future at rate δ .

 $^{^{8}\}pi$ is used indicate the vector of random benefit values.

⁹Intuitively, this assumption is sufficient to ensure that no partially altruistic individual ever wants to be harmed. In terms of the parable from section 1, the cost of having a bottle land in one's garden is equal to the cost of walking to a bin at the edge of the island. Although this might seem a restrictive assumption, we can always, by taking the limit as the expected value of π , $\bar{\pi} = \int_{\pi}^{1} \pi g(\pi) d\pi$, goes either to $\hat{\pi}$ or 1, examine the cases where the cost of a bottle landing is almost identical to the cost of walking to the bin ($\bar{\pi} \longrightarrow 1$), or almost always greater than the cost of walking to the bin ($\bar{\pi} \longrightarrow \hat{\pi}$). (See subsection 8.1.)

 $^{^{10}}$ Note that they do so after they observe the benefit value, so this is tantamount to choosing whether or not to punish for any given benefit value. It nonetheless proves useful later on to think in terms of trigger levels.

¹¹We will frequently use the term "social utility" in order to emphasise the inclusion of altruistic preferences. However, the modelling role of the social utility function is similar to any standard utility function.

 $^{^{12}}$ Effectively, every individual's social utility function is a social welfare functional which aggregates the orderings represented by all players' felicities, and which satisfies the Pareto principle, independence of irrelevant alternatives and unrestricted domain. Ratio scale comparability [Roberts, 1980] must be assumed, with all individuals gaining 0 felicity in the state of the world where no harm opportunities at all are taken.

the felicities of others in each player's social utility function. We assume that θ is identical for all players and is always strictly less than 1.

Definition 1. Let $f_{i,t}$ be the **felicity** of player *i* in period *t*:

$$f_{i,t} = \left\{ \begin{array}{ccc} -1 & if \quad T_t < \pi_t \text{ and } t \neq i \text{ and } A_t = i \\ \pi_t & if \quad T_t < \pi_t \text{ and } t = i \text{ and } A_t \neq i \\ \pi_t - 1 & if \quad T_t < \pi_t \text{ and } t = i \text{ and } A_t = i \\ 0 & otherwise \end{array} \right\}$$
(1)

Let δ be the discount rate¹³ and let θ be the coefficient of altruism:

$$0 \le \delta < 1 \tag{2}$$

$$\theta < 1$$
 (3)

Let u_t be the expected utility of player t looking forward from period t.¹⁴

$$u_{t} = E_{\pi} \left[\sum_{j=t+1}^{\infty} \left(\delta^{j-t} \left(f_{t,j} + \theta \sum_{k \neq t}^{\infty} f_{k,j} \right) \right) \right] \Big|_{\pi_{1}...\pi_{t}}$$
$$= \sum_{j=t+1}^{\infty} \left(\delta^{j-t} \left(E_{\pi} \left[f_{t,j} \right] \Big|_{\pi_{1}...\pi_{t}} + \theta \sum_{k \neq t}^{\infty} E_{\pi} \left[f_{k,j} \right] \Big|_{\pi_{1}...\pi_{t}} \right) \right)$$

5 The Single-Move Game

Consider first a single-move sequential punishment game in which a single individual (individual 1) has an opportunity to harm another (this can be thought of as a special case of the infinite-move game in which $\delta = 0$ so that there is no future). It is obvious that the individual's altruism level must be sufficiently high in order to prevent him or her from yielding to the temptation to inflict harm socially inefficiently, and so here the deleterious willingness and severity effects of greater altruism do not apply. In this simple case there is therefore no sense in which too much altruism is bad for society. It is clearly socially efficient for a harm opportunity to be taken if and only if $\pi_1 > 1$. The individual receiving the harm opportunity (individual 1), meanwhile, will choose to inflict harm if and only $\pi_1 > \theta$, since he values 1 unit of harm done to another individual at θ .¹⁵ The outcome can therefore only be socially efficient, meaning socially efficient for all possible revealed benefit values, if $\theta = 1$.

 $^{^{13}}$ The role of the assumption of discrete time periods with discounting of the future can be justified as the simplest way of capturing the idea that the technology used to detect deviation is imperfect and thus takes time [Cremer, 1986].

¹⁴The assumption that players are infinitely lived may appear restrictive, but its primary role is to simplify the model. Versions of the Folk Theorem have been proved for games with finitely-lived players and overlapping generations [Kotlikoff et al., 1988] [Kandori, 1992] [Messner & Polborn, 2003], and the general result is that having finitely-lived agents reduces, but does not eliminate, the possibility of supporting mutually beneficial equilibria in an infinitely-repeated stage game framework. It therefore seems reasonable to focus on the role of altruism by assuming away the issue of finite lifespans.

¹⁵We can assume that the individual receiving the harm opportunity will definitely harm another individual rather than himself because $\theta < 1$. We assume that if he is indifferent, he will not inflict harm.

6 The Infinite-Move Game

Once we introduce an infinite series of sequential moves, where different individuals receive an opportunity to inflict harm one after the other, the effect that altruism has on the willingness to punish and on the severity of punishment becomes important. The socially efficient outcome can now be achieved when individuals are less altruistic than the level required in the single-move game, through the use of punishment equilibria.

A socially efficient outcome can be achieved in the infinite-move game by using the information available, and the fact that players are able to choose whom they harm, to enforce credible threats of future punishment upon players who are tempted to socially inefficiently inflict harm in the current period. The ability of a player to inflict harm then plays the dual role of a temptation to impose a deadweight loss upon society at benefit to oneself, but also the opportunity for society to credibly threaten to punish those who do so. This means that there can then be some advantages to players being less than fully altruistic, since the threat of punishment is more severe the less altruistic players are, both because less altruistic players are willing to inflict harm more often (the willingness effect), but also because the loss of social utility from being harmed in place of another is greater for a less altruistic player (the severity effect).

It turns out that, provided there is sufficiently low discounting, the severity effect dominates and the lower constraint on the required level of altruism drops away because decreasing the level of altruism beyond a certain point always increases the severity of punishment more than enough to outweigh the increased temptation to deviate from the socially efficient equilibrium by inflicting harm (see Theorem 3). Most significantly, however, it transpires that too much altruism will, for any value of the discount rate, *prevent* the socially efficient outcome from being achieved (see Theorem 2).

The possibility of achieving this kind of socially efficient subgame perfect Nash equilibrium begs the question of how it is that the players are to be coordinated upon playing it. We might think of a social planner declaring the equilibrium that will be played, and then each player behaving unilaterally with no ability to communicate with the others. On the other hand, we could imagine a kind of "original position" [Rawls, 1999], where the players together agree to the planned equilibrium that will maximise their collective expected social utility. This is relevant since generally speaking although the equilibrium outcome will be socially efficient, the off-equilibrium behaviour prescribed by the punishment strategies will not be, and so there would be the temptation for a social planner to intervene in order to achieve a socially efficient outcome in the subgame beginning once someone has actually deviated. The social planner could, if able to alter the expectations of all players about the strategies being played by all the others, improve social welfare once someone has actually deviated by "letting bygones be bygones" and re-coordinating all players upon a new socially efficient equilibrium in the subgame starting in the current period.

The most stringent requirement we could put on punishment equilibria is that they be renegotiation-proof. This equilibrium concept refinement has been applied to other repeated games [Farrell & Maskin, 1989] [Benoit & Krishna, 1993], and shown to reduce the number of supportable equilibria. Sometimes, depending on the context, efficient subgame-perfect



Figure 1: Socially efficient equilibria and the socially optimal level of altruism

equilibria can be rendered unsupportable. The reason is that often the most severe punishments are not renegotiation proof because everyone would prefer to "let bygones be bygones" and renegotiate to a Pareto-superior path. In this paper, however, we stick with the requirement that punishment equilibria be subgame-perfect rather than renegotiation-proof. Intuitively, we assume that the social planner (or the community's decision-making process) is able to avoid the temptation to let malefactors "off the hook".

Figure 1 provides a preliminary schematic for the possible subgame-perfect socially efficient equilibria which are supportable for different values of θ and δ in the sequential punishment model. The most lightly shaded area A shows values of (θ, δ) where, by using "Nash reversion" (which requires each individual along an equilibrium path where a previous deviator is being punished to take their harm opportunity whenever $\pi_t > \theta$, as they would in a single-move game) a socially efficient equilibrium can be constructed. Individuals are incentivized to co-operate with the equilibrium due to the threat of focusing future punishment onto any deviator from the initial path.

The darker grey area B shows those values of (θ, δ) for which social efficiency can only be supported by using a punishment path more severe than Nash reversion. This requires that the individuals doing the punishing be required to go "beyond their comfort zone" by inflicting harm for values of $\pi_t \leq \theta$ for which they would not do so in a single-move game, due to their partial altruism. The main analytic task of this paper is to characterize the nature of the socially efficient equilibria that can be supported using the most severe available path. The black region shows those values of (θ, δ) for which social efficiency is not supportable, even with the use of such an optimal path. The central result of the paper is that this region is "thinnest" at a **socially optimal level of altruism**, θ^* , and it will be shown that, under the fairly general assumptions made regarding the distribution of the benefit and the shared value of θ and δ , it is always the case that $0 < \theta^* < 1$ (see Theorem 4).

7 Punishment Paths

The sequential punishment model has close parallels with the traditional framework of infinitely-repeated games with discounting. Seminal results for the nature of the optimal penal codes in these types of game were provided by Abreu [Abreu, 1988], who showed that optimal punishment can be exhaustively described using **punishment paths**. These will in general have a **carrot-and-stick** structure, with players incentivized to co-operate with the more unpleasant early stages of the path by the "carrot" offered by the return to more pleasant co-operative behaviour in the later part of the path.¹⁶ The introduction of non-stationary carrot-and-stick punishments is particularly interesting in the sequential punishment model because partially altruistic individuals must themselves be threatened with harm if they refuse to co-operate with the punishment of others. This feature of the model generates a rich interaction between the altruistic preferences of the players and the structure of optimal punishment paths.

Strategy profiles and the corresponding equilibria in the sequential punishment model can be described in terms of an **initial path** and a **punishment path**. Along the initial path, no harm opportunities are permitted to be taken. If a player deviates from the initial path, then a punishment path tailored for that player is initiated. If a player deviates from an ongoing punishment path, then a new punishment path tailored for the most recent deviator is initiated.

A punishment path, denoted by ψ , is a vector of trigger levels for π above which harm opportunities are taken in a punishment equilibrium. Punishment paths provide a natural way to conceive of punishment equilibria in the sequential punishment model. If a punishment path, which was initiated in period j through a deviation by player j, is being followed in period t, then player t sets their trigger level T_t equal to ψ_{t-j} (so that player t takes the harm opportunity when $\pi_t > \psi_{t-j}$) and punishes player j by setting $A_t = j$.

[Abreu, 1988]

 $^{^{16}}$ Abreu also foresaw that his method would have far-reaching applications in other models:

Analogues to the theorems established here ought to appear in any model with discounting and a "repeated" structure. Finally, the conceptualization of punishment in terms of paths and deviations from prescribed paths should prove useful in other contexts.

The sequential punishment model analysed here is one such context. Although the sequential punishment model is not strictly-speaking a repeated stage game, the ability of individuals to condition their behaviour on the past, with deviations immediately observable next period, gives it an essentially analogous structure.

Definition 2. A punishment path, denoted ψ , is a vector of trigger levels, subscripted by the point reached along the path.¹⁷ Trigger levels must lie within the support for π , therefore $\forall_k : \psi_k \in [\hat{\pi}, 1]$. The set of possible punishment paths is Ψ , so that $\forall_{\psi} : \psi \in \Psi$. A flat punishment path, $\tilde{\psi} \in \tilde{\Psi}$, has the property that $\forall_k : \tilde{\psi}_k = \tilde{\psi}$. (We use $\tilde{\psi}$ to denote both a flat path and the trigger level that defines it.) The set of flat paths is $\tilde{\Psi}$.

Following Abreu's argument, in order to find out if the socially efficient outcome is supportable for any given θ and δ , it is in general necessary to derive the **optimal punishment path**. Along a punishment path, it will clearly be desirable to harm the most recent deviator as much as possible. Since players are indifferent as to whom they harm, any harm opportunities taken along an optimal punishment path will therefore be "focussed" upon the most recent deviator.

We can imagine choosing a fixed punishment, and then finding out the most severe path we can support given the use of that fixed punishment for any deviation. However, as argued by Abreu, we will only have found the most severe path we can support if we are in fact using that path to punish any deviation from any ongoing punishment path. Hence the optimal punishment path must be used to punish any deviation from itself. This is a useful recursive symmetry which we exploit in constructing the conditions for supportability below.

There are two constraints at each point along a punishment path. The first concerns the "sqeamishness" of partially altruistic individuals in implementing the "stick". Individual t is only willing to take a harm opportunity when $\pi_t \leq \theta$ if they are themselves threatened with punishment, in order to give them an incentive to inflict harm when it is unpleasant for them to do so. The second constraint concerns the "carrot" part of the path. In order to provide a carrot, it is necessary that trigger levels be higher later in the path (so that harm is inflicted only for high benefit values). This may involve individuals being required to *abstain* from taking a harm opportunity when $\pi_t > \theta$, for which they will also need to be given an incentive via carrot-and-stick punishment.

The second constraint turns out to be more difficult to deal with, but we are able to prove that, as $\theta \longrightarrow 1^-$, this constraint becomes insignificant, because even when it is not imposed, the socially efficient outcome becomes unsupportable using the optimal path anyway. Also, in many cases the second ("upper") constraint will not bind at any point along the path, whereas the first ("lower") constraint must always bind at the beginning of the optimal path. It is therefore the first constraint which primarily drives the shape of optimal punishment paths in the sequential punishment model.

Ignoring the upper constraint, optimal paths will be shown to have a **quasi-flat** structure, meaning that the trigger level is identical following the second point along the path. This is a surprising result, since optimal penal codes in infinitely-repeated stage game models, such as the Cournot and Bertrand oligopoly models, involve a finite number of periods of punishment followed by a return to full co-operation, where the Pareto efficient outcome in the stage game is restored [Abreu, 1986] [Lambson, 1987]. The different result in the

 $^{^{17}{\}rm We}$ use the term "period" to refer to "game time" and "point" to refer to the current position along an ongoing path.

sequential punishment model is driven by the presence of altruistic preferences, which cause "neutral observers" who are not being punished (but who are still affected by the carrot created by the remainder of the punishment path) to be more sensitive to variation in the trigger levels than the individual being punished (the first have a more concave inter-temporal utility function than the second). Intuitively, with partial altruism ($\theta < 1$), the individual being punished is hurt partly or primarily simply because *they* are being punished, whereas the values of the benefit for which harm opportunities are taken makes more relative difference to a "neutral observer".

The socially efficient outcome is **supportable** for given values of δ and θ if and only if there exists a ψ such that the corresponding strategy profile forms a subgame perfect Nash equilibrium. Checking for supportability involves two conditions. Firstly, the punishment path ψ must be **sustainable**. This requires that individuals be incentivized to co-operate with the punishment path, either by punishing when they would prefer not to in a single-move game, or by refraining from punishing when they would prefer to. Secondly, given a sustainable path, it must be of sufficient **severity** to incentivize all players to co-operate with the initial path, so that the socially efficient outcome occurs in equilibrium.

Definition 3. Let $U_k : \Psi \longrightarrow \mathcal{R}$ be the per-period average discounted expected utility of the individual being punished along path ψ , looking forward from point k. Let $V_k : \Psi \longrightarrow \mathcal{R}$ be the per-period average discounted expected utility of a "neutral observer" who is not being punished along path ψ .¹⁸

$$U_k(\psi) \equiv \left(\frac{1-\delta}{\delta}\right) \sum_{i=k+1}^{\infty} \left[\delta^{i-k} \int_{\psi_i}^1 (\theta\pi - 1)g(\pi)d\pi\right]$$
(4)

$$V_k(\psi) \equiv \left(\frac{1-\delta}{\delta}\right) \sum_{i=k+1}^{\infty} \left[\delta^{i-k} \int_{\psi_i}^1 (\theta\pi - \theta)g(\pi)d\pi\right]$$
(5)

Note that, for a flat path $\tilde{\psi}$, these functions simplify to give $\forall_k : U_k(\tilde{\psi}) \equiv U(\tilde{\psi}) \equiv \int_{\tilde{\psi}}^1 (\theta \pi - 1)g(\pi)d\pi$ and $\forall_k : V_k(\tilde{\psi}) \equiv V(\tilde{\psi}) \equiv \int_{\tilde{\psi}}^1 (\theta \pi - \theta)g(\pi)d\pi$. (Note also that the subscript to indicate the point reached along the path can be suppressed for a flat path.)

The supportability constraints are as follows. $\lambda_k : \Psi \longrightarrow \mathcal{R}$ is the lowest possible net loss of utility from refusing to punish when required to at point k along punishment path ψ (this only "bites" when $\psi_k < \theta$) and $\mu_k : \Psi \longrightarrow \mathcal{R}$ is the lowest possible net loss of utility from punishing when required not to along punishment path ψ (this only "bites" when $\psi_k > \theta$). $\kappa : \Psi \longrightarrow \mathcal{R}$, meanwhile, is the lowest possible net loss in utility from defecting from the initial path, given that path ψ is used to punish such a deviation.

In order for punishment path ψ to support a socially efficient equilibrium, it must be the case that $\forall_k : \lambda_k(\psi) \ge 0, \forall_k : \mu_k(\psi) \ge 0$ and that $\kappa(\psi) \ge 0$. The optimal punishment path is the one which minimises $U_0(\psi)$ subject to these constraints, which is the same as maximizing the severity of the punishment path for the punishee, denoted by $\phi : \Psi \longrightarrow \mathcal{R}$.

 $^{^{18}}$ R denotes the set of real numbers.

$$\lambda_k(\psi) \equiv \left(\frac{\delta}{1-\delta}\right) V_k(\psi) - \left(\frac{\delta}{1-\delta}\right) U_0(\psi) + \psi_k - \theta \tag{6}$$

$$\mu_k(\psi) \equiv \left(\frac{\delta}{1-\delta}\right) V_k(\psi) - \left(\frac{\delta}{1-\delta}\right) U_0(\psi) - \psi_k + \theta \tag{7}$$

$$\kappa(\psi) \equiv -\left(\frac{\delta}{1-\delta}\right) U_0(\psi) + \theta - 1 \tag{8}$$

$$\phi(\psi) \equiv -\left(\frac{\delta}{1-\delta}\right) U_0(\psi) \tag{9}$$

The **optimal path**, ψ^* is therefore the path that maximizes $\phi(\psi)$ whilst satisfying all the supportability constraints. The **optimal flat path**, $\tilde{\psi}^*$ is defined in an analogous manner.

8 The Optimal Flat Path

In this section, we proceed to exhaustively derive the nature of the socially efficient equilibria which can be supported using the optimal flat punishment path. The uniform distribution is used as an illustrative example, but Theorems 1, 3 and 4 apply for any continuous benefit distribution. Theorem 2 only applies to equilibria supported by flat paths. It will be generalized later on.

Intuitively, the impact of the level of altruism on the severity of the optimal flat punishment path, and the resultant supportability of the socially efficient equilibrium, depends upon the interaction of the temptation, severity and willingness effects laid out in section 2. The temptation effect means that lower altruism makes a socially efficient outcome harder to support, *ceteris paribus*. The severity effect, by contrast, makes a given punishment path more severe with a lower coefficient of altruism, and so increases the supportability of the socially efficient equilibrium, *ceteris paribus*. The willingness effect makes individuals more willing to punish with lower altruism, thus also rendering sustainable punishment paths more severe, and social efficiency easier to support, *ceteris paribus*.

The key result to be established is that, as $\theta \to 1^-$ and individuals become perfectly altruistic, the interaction of the three effects leads to a breakdown of the socially efficient equilibrium. The intuition is, firstly, that when $\theta = 1$, the severity of any punishment path will be 0, and the optimal path will not involve any harm being inflicted. This is because perfectly altruistic individuals do not mind harm being focused from other people onto them, and so there is no loss of utility from defecting from the punishment path, and therefore individuals cannot be incentivized to do any punishing at all. The constraint for supportability of the initial path must therefore be just fulfilled with equality at this point (because there is also no temptation to defect).

If the coefficient of altruism is reduced slightly below θ then, since very little punishment can be sustained with such a high coefficient of altruism, the willingness and severity effects must be negligible. Hence the temptation effect must dominate, and social efficiency must be rendered unsupportable. However, as θ is further reduced, provided δ is sufficiently high, the combined willingness and severity effects will eventually become large enough to offset the temptation effect and lead the condition for social efficiency to be supported again. For high enough values of δ , this can continue to occur as $\theta \longrightarrow -\infty$, meaning that social efficiency can be supported even with infinite malevolence. (This phenomenon is driven by the severity effect - see Theorem 3.)

There will in general exist an optimal level of altruism θ^* , in the sense that it allows the socially efficient outcome to be supported for the widest range of δ , δ^* . θ^* has the general property that it is low enough for the punishment path to involve inflicting harm for the widest possible range of benefit values, but not any lower, in order to avoid the temptation effect outweighing the severity effect as θ is further reduced. As we shall see later, for a "sufficiently flat" benefit distribution, the optimal punishment path is always flat. Even when the optimal path is not flat, however, Theorem 2 provides an important intermediate step in proving the core result of the paper - that too high a level of altruism is socially detrimental - for these more general cases.

In this section, we will derive the properties of the optimal flat punishment path denoted by $\tilde{\psi}^*$. Along such a path, the co-operation constraints at every point are identical, and so $\forall_k : \lambda_k(\tilde{\psi}) \equiv \lambda(\tilde{\psi})$ and $\forall_k : \mu_k(\tilde{\psi}) \equiv \mu(\tilde{\psi})$. Assuming that $\tilde{\psi}^*$ is greater than $\hat{\pi}$ (i.e. that individuals cannot be incentivized to punish for all possible values of the benefit), then an expression for the optimal flat trigger level can be found by setting $\lambda(\tilde{\psi}^*) = 0$ so that all individuals are being pushed right up against the limit of their willingness to punish given that they themselves are threatened with punishment if they refuse. Substituting in (4) and (5) into (6), equating with 0 and rearranging gives us the following equation for $\tilde{\psi}^*$.

$$\widetilde{\psi}^* = \theta - (1 - \theta) \frac{\delta}{1 - \delta} \int_{\widetilde{\psi}^*}^1 g(\pi) d\pi$$
(10)

Although this only implicitly defines $\tilde{\psi}^*$, and cannot be solved without making specific assumptions about the functional form of $g(\pi)$, it can be totally differentiated and rearranged to yield the following expression for the derivative $\frac{d\tilde{\psi}^*}{d\theta}$. The shows the willingness effect - the impact of a change in the coefficient of altruism upon the optimal flat trigger level. We should note at this point that as $\theta \longrightarrow 1^-$, this expression becomes unambiguously positive; the willingness effect negatively affects the severity of punishment as θ increases at an already high level of altruism (and thus positively affects the severity of punishment as θ is decreased).

$$\frac{d\widetilde{\psi}^*}{d\theta} = \frac{(1-\delta) + \delta \int_{\widetilde{\psi}^*}^1 g(\pi) d\pi}{(1-\delta) - \delta(1-\theta)g(\widetilde{\psi}^*)}$$
(11)

If θ is low enough, individuals will be willing to punish for all possible vales of the benefit. This will imply that $\lambda(\tilde{\psi}^*) \geq 0$ when $\tilde{\psi}^* = \hat{\pi}$. The following theorem derives the relevant value of θ .

Theorem 1. If $\theta \leq \delta + (1-\delta)\hat{\pi}$ then punishment will occur for all benefit values along the optimal path.

$$(If \ \theta \le \delta + (1-\delta)\hat{\pi} \ then \ \hat{\psi}^* = \hat{\pi}, \ otherwise \ \hat{\psi}^* = \theta - (1-\theta)\frac{\delta}{1-\delta}\int_{\tilde{\psi}^*}^1 g(\pi)d\pi.)$$

Proof. Substituting expressions (4) and (5) into (6) and making this greater than or equal to 0 where $\tilde{\psi}^* = \hat{\pi}$ gives us $(1 - \theta) \frac{\delta}{1 - \delta} \int_{\hat{\pi}}^1 g(\pi) d\pi + \hat{\pi} - \theta \ge 0$.

Since $\int_{\hat{\pi}}^{1} g(\pi) d\pi = 1$, rearranging yields the stated inequality. If $\theta > \delta + (1-\delta)\hat{\pi}$, on the other hand, then $\tilde{\psi}^*$ must be where $\lambda(\tilde{\psi}^*) = 0$ in an interior solution as described by equation (10).

Having derived the optimal flat punishment path, we are now in a position to characterise the socially efficient equilibria which can be supported using it. Substituting in (4) into (8) gives us the following for $\kappa(\tilde{\psi}^*)$, (along with its total derivative with respect to θ):

$$\kappa\left(\widetilde{\psi}^*\right) = -\frac{\delta}{1-\delta} \int_{\widetilde{\psi}^*}^1 (\theta\pi - 1)g(\pi)d\pi + \theta - 1 \tag{12}$$

$$\frac{d\kappa}{d\theta} = 1 - \frac{\delta}{1 - \delta} \left(\int_{\widetilde{\psi}^*}^1 \pi g(\pi) d\pi + \left(1 - \theta \widetilde{\psi}^*\right) g(\widetilde{\psi}^*) \frac{d\widetilde{\psi}^*}{d\theta} \right)$$
(13)

We can now prove that, for any functional form for $g(\pi)$, there will exist values of θ close to but below 1 for which social efficiency will not be supportable (i.e. for which $\kappa(\tilde{\psi}^*) < 0$).¹⁹

Theorem 2. As altruism becomes perfect, the optimal flat punishment path cannot support the socially efficient equilibrium, for any value of the discount rate.

$$(As \ \theta \longrightarrow 1^{-}, \ \widetilde{\psi}^{*} \longrightarrow 1, \ \kappa \left(\widetilde{\psi}^{*}\right) \longrightarrow 0 \ and \ \frac{d\kappa}{d\theta} \longrightarrow 1, \ therefore \ as \ \theta \longrightarrow 1^{-}, \\ \kappa \left(\widetilde{\psi}^{*}\right) \longrightarrow 0^{-}.)$$

Proof. As $\theta \longrightarrow 1$, it can be seen from expression (10) that $\tilde{\psi}^* \longrightarrow 1$. The RHS of (12) thus goes to 0. Meanwhile, the RHS of (13) goes to 1. Since $\kappa(\psi)$ is a continuously differentiable function, it must therefore be the case that $\kappa(\psi)$ falls below 0 for some values of θ close to but less than 1.

We will now show that, if δ is high enough, then, once $\tilde{\psi}^* = \hat{\pi}$, so that punishment is occurring for all possible values of the benefit, the severity effect will dominate. This means that as $\theta \longrightarrow -\infty$, $\kappa(\tilde{\psi}^*) \longrightarrow \infty$ and so social efficiency becomes unambiguously supportable. The following theorem derives the required condition on δ .

Theorem 3. If $\delta > \frac{1}{1+\pi}$ then the socially efficient equilibrium can be supported with infinite malevolence.

$$(If \, \delta > \frac{1}{1+\bar{\pi}} \, then \, \kappa \left(\widetilde{\psi}^* \right) \longrightarrow \infty \, as \, \theta \longrightarrow -\infty.)$$

Proof. By Theorem 1, when $\theta \leq \delta + (1 - \delta)\hat{\pi}$ and so $\tilde{\psi}^* = \hat{\pi}$, there is no further willingness effect and so $\frac{d\tilde{\psi}^*}{d\theta} = 0$. As $\theta \longrightarrow -\infty$, this must occur. Therefore, as can be seen from (12), as $\theta \longrightarrow -\infty$, $\kappa(\tilde{\psi}^*) \longrightarrow \infty$ provided that $\frac{\delta}{1-\delta} > \frac{1}{\int_{\pi}^{1} \pi g(\pi) d\pi}$. (This is because, as $\theta \longrightarrow -\infty$, any part of (12) not containing θ becomes negligible.) It can similarly be seen from (13) that these same conditions will ensure that $\frac{d\kappa}{d\theta} < 0$ once $\theta \leq \delta + (1 - \delta)\hat{\pi}$. Letting $\bar{\pi} \equiv \int_{\pi}^{1} \pi g(\pi) d\pi$ and rearranging yields the stated result.

¹⁹We will prove this result more generally for any generic optimal punishment path in section 10, Theorem 7.

The above theorem shows that as the coefficient of altruism becomes infinitely negative, the severity effect will dominate if $\delta > \frac{1}{1+\pi}$. Since this lower bound for δ is less than 1, there will be a range where too high a level of altruism renders the socially efficient equilibrium unsupportable but, once θ is below the upper limit, no arbitrarily high degree of malevolence will do so. If, however, $\delta^* < \delta < \frac{1}{1+\pi}$ then both too high and too low a level of altruism will cause a breakdown of efficiency.

There are a number of approaches which we could take in defining the socially optimal level of altruism in the sequential punishment model. In a world where we were unable to achieve the first-best solution, we could ask what the impact of a change in the coefficient of altruism is upon the efficiency of the second-best equilibrium. This we do in section 11. In this section, we concentrate on worlds where the first-best solution is available, and ask what value of θ allows the socially efficient outcome to be supportable for the widest range of δ . We thus not only consider the best we can do in each possible world, but begin by considering the broader and more "philosophical" issue of which is the best of all possible worlds to be in.

The following theorem defines the socially optimal level of altruism, θ^* and corresponding minimum δ , δ^* . It also establishes that $\delta^* \leq \frac{1}{1+\pi}$ for any benefit distribution, so that both too high and too low a level of altruism relative to θ^* will cause a break-down of social efficiency, for values of δ close to, but above, δ^* . The optimal coefficient of altruism has a number of key features. Firstly, it must be "knife-edge" socially efficient equilibrium so that $\kappa(\tilde{\psi}^*) = 0$. Secondly, it must be the case that δ is just high enough so that punishment can occur for all values of π , in order that punishment paths are maximally severe for the individual being punished.

Theorem 4. The socially optimal level of altruism is always strictly positive and strictly less than 1.

(The socially optimal level of altruism is $\theta^* = \frac{3+\bar{\pi}\hat{\pi}-\sqrt{5+2\,\bar{\pi}\hat{\pi}+\bar{\pi}^2\hat{\pi}^2-4\,\bar{\pi}-4\,\bar{\pi}}}{2(1+\bar{\pi})}$, where $0 < \theta^* < 1$.)

Proof. The socially optimal level of altruism is where both $\kappa(\tilde{\psi}^*) = 0$ and $\theta = \delta + (1-\delta)\hat{\pi}$. The following two equations must therefore hold simultaneously: (Equation (15) is derived from Theorem 1. Equation (14) is derived from setting equation (12) equal to 0 and plugging in $\tilde{\psi}^* = \hat{\pi}$ and $\bar{\pi} \equiv \int_{\hat{\pi}}^1 \pi g(\pi) d\pi$.)

$$\frac{\delta}{1-\delta} = \frac{1-\theta}{1-\theta\bar{\pi}} \tag{14}$$

$$\theta = \delta + (1 - \delta)\hat{\pi} \tag{15}$$

Equations (14) and (15) together form a quadratic equation system, yielding the following solution: (Note that the second solution to the quadratic can be discounted since we require that $\theta^* < 1$ in order for (15) to be satisfied with $\delta^* < 1$.)

$$\delta^* = \frac{3 - 2\,\hat{\pi} - \bar{\pi}\hat{\pi} - \sqrt{5 + 2\,\bar{\pi}\hat{\pi} + \bar{\pi}^2\hat{\pi}^2 - 4\,\bar{\pi} - 4\,\hat{\pi}}}{2\,(1 - \hat{\pi})\,(1 + \bar{\pi})}\tag{16}$$

$$\theta^* = \frac{3 + \bar{\pi}\hat{\pi} - \sqrt{5 + 2\,\bar{\pi}\hat{\pi} + \bar{\pi}^2\hat{\pi}^2 - 4\,\bar{\pi} - 4\,\hat{\pi}}}{2\,(1 + \bar{\pi})} \tag{17}$$

To interpret (16) and (17), we first need to observe that $\sqrt{5+2\pi\hat{\pi}\hat{\pi}+\pi^2\hat{\pi}^2-4\pi-4\hat{\pi}}$ is decreasing in $\hat{\pi}$ and $\bar{\pi}$ and so its value will lie between $\sqrt{5}$ (when $\hat{\pi}=0$ and $\bar{\pi}=0$) and 0 (when $\hat{\pi}=1$ and $\bar{\pi}=1$). Therefore, it can immediately be seen that, since $\sqrt{5} < 3$, the socially optimal level of altruism defined by (17) will always be positive. Also, since θ^* is increasing in $\hat{\pi}$, its upper limiting value will be 1, as $\hat{\pi} \longrightarrow 1^-$. Since $0 < \theta^* < 1$ and $0 \leq \pi < 1$, it can therefore also be seen from (14) that $0 < \delta^* < 1$.

In order to be certain that (17) defines a point where a further decrease in θ will render the socially efficient initial path unsupportable, we need to show that δ^* from (16) lies weakly below $\frac{1}{1+\bar{\pi}}$, derived in Theorem 3. Dividing the RHS of (16) by $\frac{1}{1+\bar{\pi}}$ yields the following ratio, which we need to show is always weakly less than 1:

$$\frac{3 - 2\,\hat{\pi} - \bar{\pi}\hat{\pi} - \sqrt{5 + 2\,\bar{\pi}\hat{\pi} + \bar{\pi}^2\hat{\pi}^2 - 4\,\bar{\pi} - 4\,\hat{\pi}}}{2\,(1 - \hat{\pi})}\tag{18}$$

Note first that when $\bar{\pi} = 1$, (18) also equals 1. If the derivative of (18) with respect to $\bar{\pi}$ can be shown to be always positive when (18) is positive, then this will be sufficient to establish that (18) is always less than 1. Differentiating (18) with respect to $\bar{\pi}$ gives us:

$$\frac{2 - \hat{\pi} - \bar{\pi}\hat{\pi}^2 - \hat{\pi}\sqrt{5 + 2\,\bar{\pi}\hat{\pi} + \bar{\pi}^2\hat{\pi}^2 - 4\,\bar{\pi} - 4\,\hat{\pi}}}{2\sqrt{5 + 2\,\bar{\pi}\hat{\pi} + \bar{\pi}^2\hat{\pi}^2 - 4\,\bar{\pi} - 4\,\hat{\pi}\,(1 - \hat{\pi})}}\tag{19}$$

Denoting $\sqrt{5+2\pi\hat{\pi}\hat{\pi}+\pi^2\hat{\pi}^2-4\pi-4\pi}$ by β , (18) and (19) become respectively:

$$\frac{3 - 2\,\hat{\pi} - \bar{\pi}\hat{\pi} - \beta}{2\,(1 - \hat{\pi})}\tag{20}$$

$$\frac{2 - \hat{\pi} - \bar{\pi}\hat{\pi}^2 - \hat{\pi}\beta}{2\beta \ (1 - \hat{\pi})} \tag{21}$$

Both (20) and (21) are decreasing in β . Therefore they can only be positive when β is low enough. Setting (20) equal to 0 and solving for β yields the following:

$$\beta = 3 - 2\,\hat{\pi} - \bar{\pi}\hat{\pi} \tag{22}$$

Since (18) must always be positive, β must be weakly lower than this.

Plugging (22) into (21) yields the following expression, which is always positive, showing that β can never be high enough to make (21) negative.

$$\frac{1-\hat{\pi}}{3-2\,\hat{\pi}-\bar{\pi}\hat{\pi}}\tag{23}$$

We have therefore established that $\delta^* \leq \frac{1}{1+\pi}$ for all the relevant values of $\hat{\pi}$ and $\bar{\pi}$. Thus, if θ is further reduced below θ^* with δ unchanged at δ^* , the supportability constraint on the socially efficient initial path will be broken. θ^* is therefore optimal in the sense that a lower level of altruism could, for some δ values, render the socially efficient outcome unsupportable. Also, by Theorem 2, a big enough *increase* in θ will cause the socially efficient equilibrium to break down.

Note that if we take the case where $\hat{\pi} = 0$ for simplicity, the solution to (14) and (15) becomes the following:

$$\theta^* = \delta^* = \frac{3 - \sqrt{5 - 4\bar{\pi}}}{2(1 + \bar{\pi})} \tag{24}$$

8.1 Illustration: unitary distribution cases

In this subsection, we produce some graphics to compare the socially optimal level of altruism in a number of key cases. For one of them, we assume that $\hat{\pi} = 0$ and we vary $\bar{\pi}$ (this case was solved in equation (24)). We compare this to the two extreme **unitary distribution** cases (where we vary $\hat{\pi}$). The first is where $\bar{\pi} \longrightarrow 1$, so that the probability density is infinitely concentrated around the point where the benefit of punishing to the punisher is equal to the cost of punishing to the punishee. The second is where $\bar{\pi} \longrightarrow \hat{\pi}$, so that the probability density is infinitely concentrated around $\hat{\pi}$, and thus the benefit is always less than the cost. By taking the limit of (14) and (15) as $\bar{\pi} \longrightarrow \hat{\pi}$, we get the following solution for the case where the benefit is always $\hat{\pi}$:

$$\theta^* = \frac{3 + \hat{\pi}^2 - \sqrt{(\hat{\pi}^2 + 2\hat{\pi} + 5)(1 - \hat{\pi})^2}}{2(\hat{\pi} + 1)} \qquad \delta^* = \frac{3 + \hat{\pi} - \sqrt{\hat{\pi}^2 + 2\hat{\pi} + 5}}{2(\hat{\pi} + 1)}$$

For the case where $\bar{\pi} \longrightarrow 1$ and so the benefit is always 1, we get the following:

$$\theta^* = \frac{\hat{\pi}}{2} + \frac{1}{2} \qquad \delta^* = \frac{1}{2}$$

Finally, if we take a "symmetric" distribution where $\bar{\pi} = \frac{\hat{\pi}}{2} + \frac{1}{2}$ then we get the following:

$$\theta^* = \frac{6 + \hat{\pi}^2 + \hat{\pi} - \sqrt{(\hat{\pi}^2 + 4\hat{\pi} + 12)(1 - \hat{\pi})^2}}{2(3 + \hat{\pi})}$$
$$\delta^* = \frac{6 - \hat{\pi}^2 - 5\hat{\pi} - \sqrt{(\hat{\pi}^2 + 4\hat{\pi} + 12)(1 - \hat{\pi})^2}}{2(3 + \hat{\pi})(1 - \hat{\pi})}$$

Figures 2 through 4 illustrate the value of (24) as $\bar{\pi}$ changes and the values of the above expressions as $\hat{\pi}$ changes. The two unitary distribution cases represent the two extreme values of θ^* given a particular value of $\hat{\pi}$. The "symmetric" distribution case can be seen to lie somewhere between them. A number of observations are worth noting. Firstly, θ^* is increasing in both $\hat{\pi}$ and $\bar{\pi}$. This can be related to the impact of the willingness, severity and temptation effects. When $\hat{\pi}$ is high, the point where optimal paths become maximal, and the willingness effect becomes 0, is reached more quickly as θ is reduced, and so the socially optimal level of altruism is higher.

When $\hat{\pi}$ is fixed and $\bar{\pi}$ is increased, on the other hand, then the severity effect is reduced at all values of θ so that, when the willingness effect becomes 0, the value of δ must be higher in order to be at a "knife-edge" point where the social efficient outcome is just supportable. Again, the balance between the severity effect and the temptation effect is effectively tilted in favour of the temptation effect, leading to a higher socially optimal level of altruism.



Figure 2: Range of θ^* given $\hat{\pi}$



Figure 3: Range of δ^* given $\hat{\pi}$



Figure 4: Range of θ^*, δ^* given $\bar{\pi}$

The second key observation to make is that the range of θ^* is quite narrow for any given $\hat{\pi}$. Again, this is due to the interaction of the three effects. This can be most clearly seen by considering the unitary distribution where the benefit is always close to 1. This means that the willingness effect is always 0, because as soon as the optimal flat trigger level $\tilde{\psi}^*$ falls below 1, all of the probability mass lies above the trigger level, and so any further reduction in θ and $\tilde{\psi}^*$ has no impact. This explains why δ^* is always $\frac{1}{2}$, because this is the value of δ that will lead the temptation and severity effects to exactly cancel out as θ is reduced. It is, however, still necessary that $\tilde{\psi}^* = \hat{\pi}$ in order to reach θ^* .²⁰

8.2 Illustration: continuous uniform distribution

Figure 5 illustrates the application of Theorems 2, 3 and 4 to the case of a uniform distribution with support between 0 and 1 and therefore where $\forall_{\pi} : g(\pi) = 1$. The key features are the socially optimal level of altruism θ^* and corresponding minimum δ^* where, by substituting in $\bar{\pi} = \frac{1}{2}$, expression (24) tells us is at $\theta^* = \delta^* = 1 - \frac{1}{\sqrt{3}}$, and the value of $\delta = \frac{1}{1+\bar{\pi}} = \frac{2}{3}$ above which the severity effect dominates as $\theta \longrightarrow -\infty$, resulting in a horizontal asymptote for the black region where social efficiency is not supportable.

Figure 6 illustrates how the value of $\kappa(\tilde{\psi}^*)$ changes (on the y-axis) for a "cross section" taken through figure 5 where $\delta = 1 - \frac{1}{\sqrt{3}}$ and θ is allowed to

 $^{^{20}}$ Intuitively, if $\bar{\pi}$ is high then the willingness effect becomes negligible but, since the temptation and severity effects are then roughly equal, it still takes quite a low value of θ before maximal punishment occurs.



Figure 5: Socially efficient equilibria where $\hat{\pi} = 0$ and $g(\pi) = 1$

vary along the x-axis.²¹ It can be seen that $\kappa(\tilde{\psi}^*)$ lies below 0 for any value of θ apart from $\theta^* = \delta^* = 1 - \frac{1}{\sqrt{3}}$ and $\theta = 1$.²²

Figure 7 shows a "cross-section" at a second key value of $\delta = 0.5$.²³ The significance of this point is that it is where the co-operation constraint for the Nash reversion punishment path and the co-operation constraint for the maximal flat punishment path both bind at the boundary of the black region (which is, for this point only, also on the boundary of the dark grey region).

A third important value of δ is that above which the temptation effect never outweighs the severity effect as $\delta \longrightarrow -\infty$. In the case of this model, this is where $\frac{\delta}{1-\delta} = 2 \Longrightarrow \delta = \frac{2}{3}$.²⁴ This is illustrated graphically in figure 8. The property of this particular value of δ that the severity and temptation effects exactly cancel is the fact that the value of $\kappa(\tilde{\psi}^*)$ is constant for any $\theta < \delta$. It is instructive to compare this diagram to figures 11 and 12, which illustrate values of δ slightly above and below $\frac{2}{3}$ respectively.²⁵ Here we notice that the value of $\kappa(\tilde{\psi}^*)$ increases as θ is reduced below δ when $\delta > \frac{2}{3}$, showing that the severity

 $^{^{21}\}mathrm{This}$ corresponds to the line labelled f in figure 5.

 $^{^{22}}$ When $\theta=1$ there is no temptation to defect and so social efficiency can always be supported.

 $^{^{23}}$ This corresponds to the line labelled d in figure 5.

 $^{^{24}\}mathrm{This}$ corresponds to the line labelled b in figure 5.

 $^{^{25}}$ These correspond to the lines labelled a and c in figure 5.



Figure 6: Values of $\tilde{\psi}^*$ and $\kappa \left(\tilde{\psi}^*\right)$ when $\delta = 1 - \frac{1}{\sqrt{3}}$



Figure 7: Values of $\tilde{\psi}^*$ and $\kappa\left(\tilde{\psi}^*\right)$ when $\delta = \frac{1}{2}$



Figure 8: Values of $\tilde{\psi}^*$ and $\kappa\left(\tilde{\psi}^*\right)$ when $\delta = \frac{2}{3}$

effect outweighs the temptation effect, and the opposite occurs when $\delta < \frac{2}{3}$. Figures 9 and 10 show cross-sections for values of δ slightly above and below $1 - \frac{1}{\sqrt{3}}$.²⁶ The key feature is that when δ is below $1 - \frac{1}{\sqrt{3}}$, the only value of θ for which $\kappa(\tilde{\psi}^*)$ is not negative is 1.

9 Quasi-Flat Paths

The next two sections will primarily be concerned with proving the result from Theorem 2 for the general case where the optimal punishment path is not flat. We begin by showing that, as $\theta \longrightarrow 1^-$, the socially efficient equilibrium becomes unsupportable using the optimal **quasi-flat** punishment path. We then proceed, in section 10, to establish that the same result holds even when the optimal generic punishment path is used.

Definition 4. A quasi-flat path is one which is flat from point 2 onwards, and is denoted by $\ddot{\psi}$. The point 1 trigger level is $\ddot{\psi}_1$. The point 2 and after trigger level is $\ddot{\psi}_2$.²⁷

An important concept that will be used repeatedly in the lemmas and theorems to follow is the definition of an average trigger level which defines a flat path which is equivalent in terms of per-period average discounted utility

 $^{^{26}}$ These correspond to the lines labelled e and g in figure 5.

 $^{^{27}}$ This is the simplest punishment path structure enabling carrot-and-stick punishment, because the individual required to punish at point 1 will take into account the future they face if they co-operate, where the path continues to the less severe "carrot" part, whereas if they defect the path will reset and the "stick" at point 1 will be repeated.



Figure 9: Values of $\tilde{\psi}^*$ and $\kappa \left(\tilde{\psi}^*\right)$ when $\delta = 1 - \frac{1}{\sqrt{3}} - 0.03$



Figure 10: Values of $\tilde{\psi}^*$ and $\kappa \left(\tilde{\psi}^*\right)$ when $\delta = 1 - \frac{1}{\sqrt{3}} + 0.03$



Figure 11: Values of $\widetilde{\psi}^*$ and $\kappa \left(\widetilde{\psi}^*\right)$ when $\delta = \frac{2}{3} - 0.03$



Figure 12: Values of $\widetilde{\psi}^*$ and $\kappa \left(\widetilde{\psi}^*\right)$ when $\delta = \frac{2}{3} + 0.03$

to a given non-flat path looking forward from a particular point. This average will in general be different for the punishee and for a "neutral observer".

Definition 5. Let the *U*-average and the *V*-average be respectively denoted as $U^{-1}(U_k(\psi))$ and $V^{-1}(V_k(\psi))$. These two averages are defined below, and total differentiation is also used to find their derivatives with respect to the trigger level at a particular point i + k (where i > 0 since the average is "forward looking"), and the implicit derivative of $V^{-1}(V_k(\psi))$ with respect to $U^{-1}(U_k(\psi))$.

$$U_{k}(\psi) \equiv \int_{U^{-1}(U_{k}(\psi))}^{1} (\theta\pi - 1)g(\pi)d\pi$$
$$\equiv \left(\frac{1-\delta}{\delta}\right) \left(\sum_{i=1}^{\infty} \left[\delta^{i} \left(\int_{\psi_{i+k}}^{1} (\theta\pi - 1)g(\pi)d\pi\right)\right]\right)$$
(25)

$$V_{k}(\psi) \equiv \int_{V^{-1}(V_{k}(\psi))}^{1} (\theta\pi - \theta)g(\pi)d\pi$$
$$\equiv \left(\frac{1-\delta}{\delta}\right) \left(\sum_{i=1}^{\infty} \left[\delta^{i} \left(\int_{\psi_{i+k}}^{1} (\theta\pi - \theta)g(\pi)d\pi\right)\right]\right)$$
(26)

$$\frac{d}{d\psi_{i+k}}U^{-1}(U_k(\psi)) = \left(\frac{1-\delta}{\delta}\right)\delta^i\left(\frac{g(\psi_{i+k})}{g(U^{-1}(U_k(\psi)))}\right)\left(\frac{1-\theta\psi_{i+k}}{1-\theta U^{-1}(U_k(\psi))}\right)$$
(27)

$$\frac{d}{d\psi_{i+k}}V^{-1}(V_k(\psi)) = \left(\frac{1-\delta}{\delta}\right)\delta^i\left(\frac{g(\psi_{i+k})}{g(V^{-1}(V_k(\psi)))}\right)\left(\frac{1-\psi_{i+k}}{1-V^{-1}(V_k(\psi))}\right)$$
(28)

$$\frac{\frac{d}{d\psi_{i+k}}V^{-1}(V_k(\psi))}{\frac{d}{d\psi_{i+k}}U^{-1}(U_k(\psi))} = \left(\frac{1-\psi_{i+k}}{1-\theta\,\psi_{i+k}}\right) \left(\frac{g\left(U^{-1}(U_k(\psi))\right)}{g(V^{-1}(V_k(\psi)))}\right) \left(\frac{1-\theta\,U^{-1}(U_k(\psi))}{1-V^{-1}(V_k(\psi))}\right)$$
(29)

The following lemma will prove useful in this and subsequent sections. It applies to all optimal punishment paths, not just quasi-flat ones, and states that U_k must we weakly minimised at point 0.

Lemma 1. The U-average must be weakly minimized at the beginning of an optimal punishment path.

((a) If a punishment path ψ^* is optimal then $\forall_k : U_k(\psi^*) \ge U_0(\psi^*)$. (b) If a punishment path ψ^* is optimal then $\psi_1^* \le U^{-1}(U_0(\psi^*)) \le U^{-1}(U_1(\psi^*))$.)

Proof. For the first claim, note that it would be possible to construct a new path ψ' identical to ψ^* except beginning at point k so that $\forall_i : \psi'_i = \psi^*_{k+i}$, resulting in the following sustainability constraints:

$$\lambda_i(\psi') \equiv \left(\frac{\delta}{1-\delta}\right) V_i(\psi') - \left(\frac{\delta}{1-\delta}\right) U_0(\psi') + \psi'_i - \theta$$
$$\mu_i(\psi') \equiv \left(\frac{\delta}{1-\delta}\right) V_i(\psi') - \left(\frac{\delta}{1-\delta}\right) U_0(\psi') - \psi'_i + \theta$$

These can be rewritten as:

$$\lambda_i(\psi') \equiv \left(\frac{\delta}{1-\delta}\right) V_{k+i}(\psi^*) - \left(\frac{\delta}{1-\delta}\right) U_k(\psi^*) + \psi_{k+i}^* - \theta$$
$$\mu_i(\psi') \equiv \left(\frac{\delta}{1-\delta}\right) V_{k+i}(\psi^*) - \left(\frac{\delta}{1-\delta}\right) U_k(\psi^*) - \psi_{k+i}^* + \theta$$

Now, since ψ^* must, by assumption, be sustainable, we know that, for any k and i:

$$\lambda_{k+i}(\psi^*) \equiv \left(\frac{\delta}{1-\delta}\right) V_{k+i}(\psi^*) - \left(\frac{\delta}{1-\delta}\right) U_0(\psi^*) + \psi^*_{k+i} - \theta \ge 0$$
$$\mu_{k+i}(\psi^*) \equiv \left(\frac{\delta}{1-\delta}\right) V_{k+i}(\psi^*) - \left(\frac{\delta}{1-\delta}\right) U_0(\psi^*) - \psi^*_{k+i} + \theta \ge 0$$

If we now suppose that there exists a k such that $U_k(\psi^*) < U_0(\psi^*)$, this would mean, by observation, that the supportability constraints for ψ' would unambiguously be fulfilled at every point. Also, this would mean that $\phi(\psi') > \phi(\psi^*)$. Therefore ψ' would be sustainable, and would be more severe than ψ^* . Hence ψ^* could not be optimal - a contradiction.

For the second claim, note that the following identity holds for any path ψ :

$$\frac{\delta}{1-\delta} \left(\int_{U^{-1}(U_0(\psi))}^{1} (\theta\pi - 1)g(\pi)d\pi \right) \equiv \delta \int_{\psi_1}^{1} (\theta\pi - 1)g(\pi)d\pi + \frac{\delta^2}{1-\delta} \int_{U^{-1}(U_1(\psi))}^{1} (\theta\pi - 1)g(\pi)d\pi$$
(30)

This can be rewritten as:

$$\left(\frac{\delta}{1-\delta}\right)U_0\left(\psi\right) \equiv \left(\frac{\delta}{1-\delta}\right)U_1\left(\psi\right) + \delta \int_{\psi_1}^{U^{-1}(U_1(\psi))} (\theta\pi - 1)g(\pi)d\pi$$

Since we know from the argument made above that $U_0(\psi^*) \leq U_1(\psi^*)$, we also know that $\int_{\psi_1^*}^{U^{-1}(U_1(\psi^*))} (1 - \theta \pi) g(\pi) d\pi \geq 0$, and therefore that $\psi_1^* \leq U^{-1}(U_1(\psi^*))$.

Finally, identity (30) can also be rewritten as:

$$0 \equiv \delta \int_{\psi_1}^{U^{-1}(U_0(\psi))} (\theta \pi - 1)g(\pi)d\pi + \frac{\delta^2}{1 - \delta} \int_{U^{-1}(U_1(\psi))}^{U^{-1}(U_0(\psi))} (\theta \pi - 1)g(\pi)d\pi$$

Since $U^{-1}(U_0(\psi^*)) \leq U^{-1}(U_1(\psi^*))$, in order for this to hold it must follow that $U^{-1}(U_0(\psi^*)) \geq \psi_1^*$. Intuitively, the U-average must be "dragged down" from below by the trigger level at point 1.

Lemma 1 implies that the optimal quasi-flat path must have a weakly lower trigger level in period 1 $(\ddot{\psi}_1^*)$ than in period 2 and after $(\ddot{\psi}_2^* \equiv U^{-1}(U_1(\psi^*)))$.²⁸ This fits the intuition that the "stick" of harsher punishment should come earlier in the punishment path so that the "carrot" of less harsh punishment later on will operate as an incentive to co-operate with the harsher punishment earlier on. The optimal quasi-flat path will clearly also satisfy $\lambda_1(\ddot{\psi}^*) = 0$.

The trade-off which drives the optimal quasi-flat path is between a "bigger stick" at point one and the resulting "nicer carrot" at point two and after. (See sections 9.5 and 9.6 for more detailed examples of this principle in action.) The gain in severity of the punishment path from a reduction in the point one trigger level depends upon the probability density at that benefit value,

²⁸Note that, for a quasi-flat path, $\forall_k : \ddot{\psi}_2^* \equiv V^{-1}(V_k(\ddot{\psi}^*)) \equiv U^{-1}(U_k(\ddot{\psi}^*)).$

whilst the effectiveness of the carrot in offsetting this to ensure sustainability also depends on the probability density at the point two and after trigger level. The *cost* of incentivizing co-operation with a lower trigger level at point one, however, does *not* depend upon the probability density at the point one trigger level, since it is "paid" in full if the value of the benefit turns out to be in the relevant range. Intuitively, therefore, if (and only if) the probability density function for the benefit is sufficiently flat will this cost will always outweigh the benefit of making the quasi-flat punishment path non-flat. Theorem 5 will establish this sufficient condition on the probability density function $g(\pi)$ to ensure that the optimal quasi-flat punishment path will be flat.

9.1 A taxonomy of optimal quasi-flat paths

An important lesson to draw from the discussion so far is that the framework of strategy profiles constructed from punishment paths does not, in and of itself, provide enough structure to allow a complete and comprehensive solution to the problem of finding the form of optimal punishment in a specific context such as that of the sequential punishment model. Although they all share a similar carrot-and-stick structure, each particular model requires its own toolkit of "tricks" to derive the precise shape of the optimal paths.

The concept of a quasi-flat punishment path turns out to be essential to analysing the equilibria supportable by optimal generic punishment paths in the sequential punishment model. This is because, as will be shown in section 10, it is always possible to construct a quasi-flat path whose severity (ϕ) value forms an upper bound for all sustainable generic paths, and to derive sufficient structure from this to extend Theorem 2 to all the necessary general cases. Also, quasi-flat paths themselves come in a variety of "flavours", the differences between them driven by the optimal structure of carrot-and-stick punishment in the sequential punishment model, and its interaction with the partially altruistic preferences of the players, along with the limits on the support for the distribution of the benefit π .

It is natural to ask at this point why quasi-flat paths emerge naturally from the structure of the sequential punishment model whilst, as we shall discuss later, optimal punishment in the infinitely-repeated Bertrand and Cournot games involves complete "front-loading", with as much punishment as possible packed into the early stages of the path. The answer is that it is the presence of partial altruism which drives the "flattening-out" of the tail of paths in the sequential punishment model.

The positioning of trigger levels at point 2 and onwards involves a trade-off, allowing greater punishment to be "bought" at point 1, but at the cost of less severe punishment later. The "sacrifice ratio" will be given by expression (29), which measures the increase in "carrot" (measured as a higher V-average) for a given reduction in "effectiveness" of the path (a higher U-average), brought about by a rise in a later trigger level, ψ_{k+i} . It can be seen that this ratio is more favourable when ψ_{k+i} is lower. (The benefits and costs are discounted, so the exchange ratio, given a particular trigger level, looks the same for all future periods.) It is therefore optimal to "spread out" the punishment evenly over the entire tail.

There are a number of possibilities for the precise form that the optimal quasi-flat path can take. By Lemma 1, it is impossible for the trigger level at point 1 to be higher than at point 2 onwards. This leaves seven possibilities, types A-G. Type A is the **maximal** path characterised in Theorem 1. Possibility B is a quasi-maximal path, where the trigger level is "maxed-out" at point 1 but not from point 2 onwards. Type C is a **flat** path. The next type, D, is a path where the amount of punishment at point 1 runs up against the constraint imposed by not being able to make the future "carrot" attractive enough to allow more severe punishment. This happens because we reach the top of the support of the distribution for π (i.e. $\ddot{\psi}_2$ reaches 1). We shall call this a quasi-minimal path. With type E paths, on the other hand, the amount of punishment at point 1 runs up against constraint (7) in that we cannot further increase ψ_2 without rendering the path unsustainable. We shall refer to this case as a **carrot-constrained** path. A sixth possibility, F, is that there is an optimal marginal trade-off between punishment in point 1 and punishment at point 2 and after. We shall call this a **carrot-maximized** path. The seventh and final case, G, is one where no punishment at all can be incentivized; this is a minimal path. Figures 13 through 19 illustrate the various possible structural types of optimal quasi-flat path.

We shall also find it essential in the lemmas and theorems to follow to distinguish between two different types of optimal path. **Fully-constrained paths** must satisfy all co-operation conditions defined by (6) and (7). **Semiconstrained paths** only need satisfy the conditions defined in (6). A fullyconstrained optimal quasi-flat path satisfies co-operation conditions $\lambda_1 \geq 0$, $\lambda_2 \geq 0$, $\mu_1 \geq 0$ and $\mu_2 \geq 0$. Given quasi-flatness, these form a sufficient condition for all the co-operation constraints to be fulfilled. A semi-constrained optimal quasi-flat path is one which is only constrained to satisfy $\lambda_1 \geq 0$ and $\lambda_2 \geq 0$ (though it may, by "chance", also satisfy $\mu_1 \geq 0$ and $\mu_2 \geq 0$, and therefore also be in the set of fully-constrained paths).

Definition 6. Let Ψ denote the set of unconstrained paths, in increasing order of ϕ . Let Ψ_{\circ} and Ψ_{\bullet} be the sets of sustainable semi-constrained and fullyconstrained paths respectively, so that $\Psi_{\bullet} \subset \Psi_{\circ} \subset \Psi$. These are similarly ordered by ϕ . The optimal generic path can be defined as $\psi^* = \sup \Psi_{\bullet} = \max_{\psi \in \Psi_{\bullet}} \{\phi(\psi)\}$. Let $\Psi \subset \Psi$ be the set of quasi-flat paths and $\widetilde{\Psi} \subset \Psi$ and be the set of flat paths. Let $\Psi_{\circ} = \Psi \cap \Psi_{\circ}, \ \widetilde{\Psi}_{\circ} = \widetilde{\Psi} \cap \Psi_{\circ}, \ \breve{\Psi}_{\bullet} = \breve{\Psi} \cap \Psi_{\bullet}$ and $\widetilde{\Psi}_{\bullet} = \widetilde{\Psi} \cap \Psi_{\bullet}$ denote analogous sets of semi-constrained and fully-constrained quasi-flat and flat paths. Other optimal paths can be defined using these sets. So, for example, $\ \psi^*_{\circ} = \sup \Psi_{\circ} = \max_{\psi \in \Psi_{\circ}} \{\phi(\psi)\}$ is the optimal semi-constrained quasi-flat path.

It should be noted at this point that the optimal semi-constrained quasi-flat path cannot be "carrot-constrained", since constraint (7) does not apply. Also, observe that the optimal semi-constrained path will be at least as effective as the optimal fully-constrained path. In other words, if all paths are ordered in ϕ , the optimal semi-constrained path will equal or beat the optimal fully-constrained path: $\phi(\sup \Psi_{\circ}) \ge \phi(\sup \Psi_{\bullet})$. Analogously, $\phi(\sup \Psi_{\circ}) \ge \phi(\sup \Psi_{\bullet})$. This observation is key in enabling Theorem 2 to be generalized to the case where the optimal generic path is used to punish deviations from the socially efficient initial path. Define ψ_k^{λ} and ψ_k^{μ} as the values of ψ_k that satisfy (6) and (7) with equality:²⁹



 $^{^{29}\}psi_k^{\lambda}$ and ψ_k^{μ} therefore represent the upper and lower limits for the trigger level that *could* be sustained at point k given the structure of the entire punishment path. Note also that, for a quasi-flat path, they will be the same for any k because $\forall_k : V_k(\ddot{\psi}) \equiv V_1(\ddot{\psi})$.











Theorems 1, 3 and 4 continue to hold for the equilibria supportable by generic optimal paths (and, therefore, quasi-flat paths) without alteration. In general, the quasi-maximal case only occurs when θ is close to the boundary established in Theorem 1 below which the optimal path is maximal. In section 9.5, we see an example of this with the triangular distribution. As we shall argue shortly, the minimal and quasi-minimal cases cannot possibly be optimal. (The proof of Lemma 5 in the appendix verifies this.) The flat, carrot-constrained and carrot-maximized paths illustrated in figures 15, 17 and 18 respectively represent the three possibilities for an interior solution.

9.2 Conditions for a flat path

To intuitively derive necessary and sufficient conditions for the optimal quasiflat path to be flat (type C), we can use a trick from Abreu by considering the optimal path we are able to construct using a fixed punishment for a deviation. If this can be shown to be flat, then the optimal path constructed using iself as a punishment will also be flat. Let $\bar{U} \equiv (\frac{\delta}{1-\delta})U_0(\bar{\psi})$ be the expected utility for the punishee along the fixed path $\bar{\psi}$. Let $\bar{\lambda}_k(\ddot{\psi}) \equiv (\frac{\delta}{1-\delta})V_k(\ddot{\psi}) - \bar{U} + \ddot{\psi}_k - \theta$ be the co-operation constraint at point k for quasi-flat path $\ddot{\psi}$ given the use of fixed path $\bar{\psi}$ to punish a deviation. Since the trigger level at point 1 should be set so that $\bar{\lambda}_1(\ddot{\psi}) = 0$, we know that the following condition must hold:

$$\ddot{\psi}_1 = \theta - \left(\frac{\delta}{1-\delta}\right) V_1(\ddot{\psi}) + \bar{U} = \theta + \left(\frac{\delta}{1-\delta}\right) \int_{\ddot{\psi}_2}^1 (\theta - \theta\pi) g(\pi) d\pi + \bar{U} \quad (31)$$

We are seeking to maximize the disutility of the person being punished along the quasi-flat punishment path $\ddot{\psi}$. This will be given by:

$$\phi = -\left(\frac{\delta}{1-\delta}\right)U_0\left(\ddot{\psi}\right) = \delta \int_{\ddot{\psi}_1}^1 \left(1-\theta\pi\right)g(\pi)d\pi + \left(\frac{\delta^2}{1-\delta}\right)\int_{\ddot{\psi}_2}^1 \left(1-\theta\pi\right)g(\pi)d\pi$$
(32)

Totally differentiating (31) gives us:

$$\frac{d\ddot{\psi}_1}{d\ddot{\psi}_2} = -\left(\frac{\delta}{1-\delta}\right)\left(\theta - \theta\ddot{\psi}_2\right)g\left(\ddot{\psi}_2\right) \tag{33}$$

Totally differentiating (32) with respect to $\ddot{\psi}_2$ gives us:

$$\frac{d\phi}{d\ddot{\psi}_2} = \delta\left(\theta\ddot{\psi}_1 - 1\right)g\left(\ddot{\psi}_1\right)\frac{d\ddot{\psi}_1}{d\ddot{\psi}_2} - \left(\frac{\delta^2}{1 - \delta}\right)\left(1 - \theta\ddot{\psi}_2\right)g\left(\ddot{\psi}_2\right)$$

Substituting in (33) and simplifying yields:

$$\frac{d\phi}{d\ddot{\psi}_2} = \left(\frac{\delta^2}{1-\delta}\right)g\left(\ddot{\psi}_2\right)\left(\left(1-\theta\ddot{\psi}_1\right)\left(\theta-\theta\ddot{\psi}_2\right)g\left(\ddot{\psi}_1\right) - \left(1-\theta\ddot{\psi}_2\right)\right)$$

This is unambiguously negative if the following condition holds:

$$g\left(\ddot{\psi}_{1}\right) < \frac{1 - \theta \ddot{\psi}_{2}}{\left(1 - \theta \ddot{\psi}_{1}\right) \theta \left(1 - \ddot{\psi}_{2}\right)}$$
(34)

The RHS of the above expression is increasing in $\ddot{\psi}_1$ and $\ddot{\psi}_2$. This means that the most stringent condition will be where $\dot{\psi}_1 = \dot{\psi}_2 = \hat{\pi}$. Requiring that the probability density function $g(\pi)$ always be less than this ensures that the above condition will always hold. This yields the following condition:

$$\forall_{\pi} : g(\pi) < \frac{1}{\theta(1-\hat{\pi})} \tag{35}$$

Provided this condition holds, increasing $\ddot{\psi}_2$ in order to further reduce $\ddot{\psi}_1$ always makes the punishment path less effective by reducing ϕ . It is therefore optimal to set $\ddot{\psi}_2 = \ddot{\psi}_1$ since to set $\ddot{\psi}_2 < \ddot{\psi}_1$ will result in a clearly non-optimal path, by Lemma 1. Condition (35) is therefore sufficient for the optimal quasi-flat path to be flat. (This result is verified in Theorem 5 below.)

A necessary condition for the optimal quasi-flat path to be flat can be found by substituting $\ddot{\psi}_1 = \ddot{\psi}_2 = \tilde{\psi}^*$ into condition (34) to give the following:

$$g\left(\tilde{\psi}^*\right) < \frac{1}{\theta\left(1 - \tilde{\psi}^*\right)} \tag{36}$$

Theorem 5. If the benefit distribution is sufficiently flat, then the optimal quasi-flat path is flat. (If $\forall_{\pi} : g(\pi) < \frac{1}{\theta(1-\hat{\pi})}$ then $\phi\left(\tilde{\psi}^*\right) \ge \phi\left(\ddot{\psi}^*\right)$.) Proof. Lemma 5 from the appendix establishes that if $\forall_{\pi} : g(\pi) < \frac{1}{\theta(1-\hat{\pi})}$ then the optimal semi-constrained quasi-flat path must be flat. This imples that $\phi\left(\tilde{\psi}^*_{\circ}\right) \geq \phi\left(\ddot{\psi}^*_{\circ}\right)$. In order to establish the same result for the optimal fullyconstrained quasi-flat path, it is sufficient to show that if the optimal semiconstrained quasi-flat path is flat (i.e. $\sup(\ddot{\Psi}_{\circ}) \in \widetilde{\Psi}$), then it will be in the set of fully-constrained flat paths ($\sup(\ddot{\Psi}_{\circ}) \in \widetilde{\Psi}_{\bullet}$). This will mean that $\phi(\tilde{\psi}^*) \geq \phi(\ddot{\psi}^*)$. If this is the case then, since $\ddot{\Psi}_{\bullet} \subset \ddot{\Psi}_{\circ}$ and thus $\phi(\ddot{\psi}^*_{\circ}) \geq \phi(\ddot{\psi}^*)$, it must follow that $\phi(\tilde{\psi}^*) \geq \phi(\ddot{\psi}^*)$.

To see that $\sup(\ddot{\Psi}_{\circ}) \in \widetilde{\Psi}_{\bullet}$, firstly observe that if $\sup(\ddot{\Psi}_{\circ}) \in \widetilde{\Psi}$ then $\sup(\ddot{\Psi}_{\circ}) \in \widetilde{\Psi}_{\circ}$, since $\sup(\ddot{\Psi}_{\circ})$ is, by assumption, semi-constrained. Now note, from expressions (6) and (7), that, if $\tilde{\psi} \leq \theta$ and $\lambda(\tilde{\psi}) \geq 0$, then $\mu(\tilde{\psi}) \geq 0$. Hence $\sup(\ddot{\Psi}_{\circ})$ must be in the set of fully-constrained paths, and so $\sup(\ddot{\Psi}_{\circ}) \in \widetilde{\Psi}_{\bullet}$. (By observation of expression (10), it must be the case that $\tilde{\psi} \leq \theta$ for any conceivable optimal flat path.)

9.3 Conditions for non-flat paths

We now proceed to lay out the conditions which must hold for the various possibile configurations of a non-flat fully-constrained optimal quasi-flat path. Most obviously, the optimal path cannot possibly be minimal (type G) because it can be seen from condition (10) that there will always exist a sustainable, and more severe, flat path. Secondly, we already know the necessary and sufficient condition for a maximal (type A) path from Theorem 1. Thirdly, taking the case of a quasi-maximal (type B) path, we know that $\ddot{\psi}_1 = \hat{\pi}$. The value for $\ddot{\psi}_2$ can then be derived using the $\lambda_1(\ddot{\psi}) = 0$ condition. This should be checked as a candidate for the fully-constrained optimal quasi-flat path.

For a carrot-constrained (type E) path, constraint (6) binds at point one and constraint (7) at points two and after. In this case, the limit to how much "carrot" can be created is imposed by the difficulty in incentivizing individuals to refrain from punishing when they would like to along the "tail" of the punishment path.³⁰ The optimal carrot-constrained path is characterised by the property that both $\lambda_1(\ddot{\psi}) = 0$ and $\mu_2(\ddot{\psi}) = 0$. Solving $\mu_2(\ddot{\psi}) - \lambda_1(\ddot{\psi}) = 0$ (applying definitions (6) and (7)) gives us:

$$\ddot{\psi}_2 = 2\theta - \ddot{\psi}_1 \tag{37}$$

The fact that a carrot-maximized (type F) path is also a possibility can be seen by observation of condition (34). As $\ddot{\psi}_2 \longrightarrow 1$, the RHS of (34) goes to infinity. Therefore the inequality will definitely be fulfilled, and further increases in $\ddot{\psi}_2$ in order to decrease $\ddot{\psi}_1$ will no longer improve the severity of the path. If this happens before $\ddot{\psi}_2$ reaches ψ_j^{μ} then the optimal quasi-flat path will be "carrot-maximized". A carrot-maximized path must therefore have the property that condition (34) is satisfied with equality. Rearranging this gives us:

³⁰Although it might be felt intuitively that if the socially efficient initial path is to be sustainable using a particular path, then co-operation with the "tail" of the punishment path would automatically also be sustainable, this does not necessarily follow because there is still a less attractive future along the tail of the punishment path if punishers co-operate, rendering the severity of punishment lower and thus making co-operation with the "tail" more difficult to incentivize than co-operation with the initial path.

$$\ddot{\psi}_2 = \frac{g\left(\ddot{\psi}_1\right)\left(1-\theta\ddot{\psi}_1\right)\theta-1}{g\left(\ddot{\psi}_1\right)\left(1-\theta\ddot{\psi}_1\right)\theta-\theta}$$
(38)

The above argument from condition (34) also shows why a quasi-minimal (type D) path is impossible, because as $\ddot{\psi}_2 \longrightarrow 1^-$ the RHS of the inequality goes to ∞ . Therefore, when raising $\ddot{\psi}_2$ in search of the optimal quasi-flat path, a path would always become carrot-maximized before it become quasi-minimal.

9.4 The optimal quasi-flat path

We now have all the information we need to lay out the procedure for finding the optimal quasi-flat path. Assuming that the optimal path is not maximal, the only possibilities are a quasi-maximal, flat, carrot-constrained or carrotmaximized path. Solving conditions (37) and (38) respectively simultaneously with equation (40) below provides a shortcut in finding the optimal quasi-flat path since, once the various solutions have been compared with the optimal flat path defined by (10) and with the quasi-maximal candidate where $\lambda(\ddot{\psi}) = 0$ and $\ddot{\psi}_1 = \hat{\pi}$, and the most severe path among them found, we can be sure that it is optimal.³¹ We use this method in subsection 9.6 to analyse a specific example of a carrot-maximized and a carrot-constrained path.

Most importantly, we can now prove the analogue of Theorem 2 for the more general case of the optimal quasi-flat path. This requires the use of Lemma 5 from the appendix.³² Once we present Theorem 7, the result will be generalized to all globally optimal fully-constrained paths. This will allow us to substantiate generally, for any benefit distribution, the result that too high a level of altruism will cause the socially efficient equilibrium to break down.³³

Theorem 6. As altruism becomes perfect, the optimal quasi-flat punishment path cannot support the socially efficient equilibrium, for any value of the discount rate. $(As \ \theta \longrightarrow 1^-, \ \kappa \left(\ddot{\psi}^* \right) \longrightarrow 0^-.)$

Proof. Firstly, observe that the most severe path in $\ddot{\Psi}_{\circ}$, $\sup \ddot{\Psi}_{\circ}$ must be at least as severe as the most severe path in $\ddot{\Psi}_{\bullet}$, $\sup \ddot{\Psi}_{\bullet}$. Therefore $\phi(\sup \ddot{\Psi}_{\circ}) \geq \phi(\sup \ddot{\Psi}_{\bullet})$ and so $\phi(\ddot{\psi}_{\circ}^*) \geq \phi(\ddot{\psi}^*)$. Now, combining with expressions (8) and

(9), we know that:

$$\kappa\left(\ddot{\psi}^*\right) \le \phi\left(\ddot{\psi}^*_\circ\right) + \theta - 1 \tag{39}$$

Theorem 2, combined with Lemma 5, has already established that the RHS of expression (39) goes to 0^- as $\theta \longrightarrow 1^-$. Therefore the LHS of (39) must be strictly negative as $\theta \longrightarrow 1^-$.

 $^{^{31}\}rm Note that we must also check the sustainability of each "candidate" because these are necessary rather than sufficient conditions.$

 $^{^{32}}$ In Lemma 5, we impose only that the optimal path be semi-constrained, partly to simplify the proof but, more importantly, because we are later able to generate semi-constrained quasi-flat paths by "flattening-out" generic paths.

³³In the course of the proof for Lemma 5, we also exhaustively derive the conditions on $g(\pi)$ under which the optimal semi-constrained quasi-flat path can be maximal, quasi-maximal, carrot-maximized and flat.



Figure 20: Triangular probability density function

9.5 Illustration: quasi-maximal paths

Given assumption (3), condition (35) will definitely hold for a uniform distribution with support between $\hat{\pi}$ and 1. It is however, instructive to look at some examples where the optimal quasi-flat path is not flat. As we have seen, this requires a distribution with a high enough probability density around the optimal flat trigger level. The simplest distribution for the benefit which can produce this result is a triangular distribution. The triangular probability density function $g(\pi) = 50(1-\pi)$ with support between 0.8 and 1 and illustrated in figure 20 is one such example.

The optimal quasi-flat path can be computed by using the $\lambda_1(\ddot{\psi}) = 0$ condition to derive the following expression for $\ddot{\psi}_1$, and then numerically solving for $\ddot{\psi}_1$ given each particular $\ddot{\psi}_2$:

$$\ddot{\psi}_{1} = \theta - \frac{\delta}{1-\delta} \int_{\ddot{\psi}_{2}}^{1} (\theta\pi - \theta)g(\pi)d\pi + \delta \int_{\ddot{\psi}_{1}}^{1} (\theta\pi - 1)g(\pi)d\pi + \frac{\delta^{2}}{1-\delta} \int_{\ddot{\psi}_{2}}^{1} (\theta\pi - 1)g(\pi)d\pi + \delta \int_{\ddot{\psi}_{1}}^{1} (\theta\pi - 1)g(\pi)d\pi + \delta \int_{\dot{\psi}_{1}}^{1} (\theta\pi - 1)g($$

Figures 21 through 26 illustrate the results of this exercise for $\delta = 0.48$ and $\delta = 0.49$, with $\theta = 0.9$. Figures 21 and 24 show the value of $\ddot{\psi}_2$ imputed from (40) for each value of $\ddot{\psi}_1$. They also show the optimal flat trigger level $\tilde{\psi}^*$ (numerically solved using (10)), and the $\ddot{\psi}_2 = \ddot{\psi}_1$ line.

If $\delta = 0.5$, then solving equation (10) to find the optimal flat path for the assumed triangular benefit distribution yields $\tilde{\psi}^* = 0.8$. This would be a maximal (type A) path involving punishment for all benefit values all of the time, as $\hat{\pi} = 0.8$. By making δ slightly lower than this, we prevent the optimal flat path from being maximal because the future is no longer "important enough" to sustain it. For the triangular distribution, this results in the optimal path being quasi-maximal (type B).

Figures 23 and 26 show the values of the various co-operation constraints, and the severity of the punishment paths (ϕ) .³⁴ In both cases, all of the constraints are positive, showing that the punishment path is itself sustainable and supports the initial path. In both cases, the highest value of $\phi(\ddot{\psi})$ occurs where $\ddot{\psi}_1 = 0.8$. We have, therefore, two cases where the optimal quasi-flat punishment path is "maxed-out" at point 1, so that $\ddot{\psi}_1^* = \hat{\pi}$. The "tail" of the path is, however, not maxed-out.

A notable feature of the graph of the imputed $\ddot{\psi}_2$ in figure 21 is that it initially climbs as $\ddot{\psi}_1$ is reduced below $\widetilde{\psi}^*$. This means that $\ddot{\psi}_2$ must initially be raised to compensate for a reduction in $\ddot{\psi}_1$ if the point 1 co-operation constraint is to remain unbroken. Eventually, on the other hand, the imputed $\ddot{\psi}_2$ graph begins to fall as ψ_1 is reduced. Since this line continues to lie above the dashed $\ddot{\psi}_2 = \ddot{\psi}_1$ line, these more severe quasi-flat punishment paths can only be achieved by introducing some carrot-and-stick element into the punishment path. By doing this, however it is possible to reduce *both* trigger levels below the trigger level of the optimal flat punishment path.

Figures 24 and 26 show a situation where the graph of the imputed $\ddot{\psi}_2$ has a positive gradient for any change in $\ddot{\psi}_1$. This means that it is now possible to reach more effective paths which remain sustainable by simultaneously reducing both trigger levels below the optimal flat trigger level, provided that $\ddot{\psi}_2$ is reduced by less than $\ddot{\psi}_1$. Here we have a classic illustration of Abreu's principle that when being optimally punished, individuals can be persuaded to put up with an even more unpleasant present in exchange for a relatively less unpleasant future.

The carrot-and-stick structure of optimal quasi-flat punishment paths illustrated here is similar to the optimal penal codes derived for infinitely repeated oligopoly games. Abreu himself derived the properties of optimal penal codes in the infinitely repeated Cournot game [Abreu, 1986]. He found that, provided there is sufficiently low discounting, optimal punishment paths would involve one period of very high output and very negative profits followed by a return to fully collusive behaviour. Such paths cannot be any worse than the 0 discounted stream of profits that could be achieved by shutting down.

Lambson built on Abreu's work to derive optimal penal codes in the infinitely repeated homogeneous product Bertrand model with identical firms and capacity constraints [Lambson, 1987]. Optimal penal codes in this model involve a number of periods of low profits followed by a return to fully collusive behaviour. This is because the losses per firm in any one period cannot be infinite in the Bertrand model. The sequential punishment model is more like Bertrand in that the severity of the punishment path at point 1 is constrained by the fact that $\ddot{\psi}_1$ cannot be any lower than $\hat{\pi}$. This means that generally there will also be some punishment along the "tail" of an optimal quasi-flat path. The difference, however, is that in the sequential punishment model, the tail of the optimal path involves "flattening-out" rather than "front-loading" the remaining punishment.

³⁴Note that because $\ddot{\psi}_1 < \theta$ and $\ddot{\psi}_2 < \theta$, constraint (7) is fulfilled and so μ_2 does not need to be displayed.



Figure 21: Trigger levels for a quasi-maximal path - $\theta=0.9,\,\delta=0.48.$



Figure 22: Co-operation constraints for a quasi-maximal path - $\delta = 0.48$.



Figure 23: Severity of a quasi-maximal path - $\delta=0.48.$



Figure 24: Trigger levels for a quasi-maximal path - $\theta=0.9,\,\delta=0.49.$



Figure 25: Co-operation constraints for a quasi-maximal path - $\delta=0.49.$



Figure 26: Severity of a quasi-maximal path - $\delta=0.49.$

9.6 Illustration: carrot-constrained and carrotmaximized paths

In this subsection, we will present specific numerical examples of a flat, a carrot-constrained and a carrot-maximized path, to illustrate more of the general principles discussed in subsection 9.1. Figures 28, 29 and 30 respectively show the numerically calculated values for $\ddot{\psi}_2$ given $\ddot{\psi}_1$, the various co-operation constraints and the effectiveness of the punishment path ϕ for the triangular distribution shown in figure 20, with $\delta = \frac{1}{3}$ and $\theta = 0.9$. Figures 31, 32 and 33 show the results for the same exercise conducted using the probability density function given in (41) below, and shown in figure 27, with $\delta = \frac{1}{8}$ and $\theta = 0.9$, $\hat{\pi} = 0.8$ and support between 0.8 and 1. Figures 34, 35 and 36 show analysis for the same distribution but with $\delta = \frac{1}{16}$.

$$g(\pi) = 17.5 \left(1 - \left(\left(\frac{2\pi - \hat{\pi} - 1}{1 - \hat{\pi}} \right)^2 \right)^{\frac{1}{5}} \right)$$
(41)

The key features to note are, firstly, that in all three cases μ_2 reaches 0 before $\ddot{\psi}_1$ reaches $\hat{\pi} = 0.8$, showing that the optimal path cannot possibly be quasi-maximal. Secondly, whereas in figure 30, ϕ decreases as $\ddot{\psi}_1$ is reduced, showing that the optimal quasi-flat path is flat, in figure 33, ϕ reaches an interior maximum, showing that the optimal quasi-flat path is carrot-maximized. In figure 36, meanwhile, ϕ reaches a constrained maximum at the carrotconstrained value of $\ddot{\psi}_1$, where the x-axis begins. Note, finally, that the starting value on the x-axis, corresponding to the carrot-maximized "candidate" for the optimal path in each case, and the carrot-maximum denoted by the dotted line in figure 33 were derived using the shortcut method described earlier in section 9.4, thus verifying its applicability.



Figure 27: Probability density for a carrot-maximized or carrot-constrained path



Figure 28: Trigger levels for a flat path - $\theta = 0.9, \, \delta = \frac{1}{3}$.



Figure 29: Co-operation constraints for a flat path.



Figure 30: Severity of a flat path



Figure 31: Trigger levels for a carrot-maximized path - $\theta = 0.9$, $\delta = \frac{1}{8}$.



Figure 32: Co-operation constraints for a carrot-maximized path.



Figure 33: Severity of a carrot-maximized path



Figure 34: Trigger levels for a carrot-constrained path - $\theta = 0.9$, $\delta = \frac{1}{16}$.



Figure 35: Co-operation constraints for a carrot-constrained path.



Figure 36: Severity of a carrot-constrained path

10 The Optimal Generic Path

We are now ready to further generalize Theorems 2 and 6 to the case where the optimal generic punishment path is used to support the initial path. The result hinges upon three intuitive observations. Firstly, the optimal semiconstrained path is quasi-flat, because it is always possible to take any given optimal semi-constrained path and "flatten-out" the tail, producing an equally severe path without breaking any of the co-operation constraints. Secondly, as $\theta \longrightarrow 1^-$, it is impossible to support the socially efficient equilibrium using the optimal semi-constrained quasi-flat path, because it must become flat (this is proved in Lemma 5), and we already know (from Theorem 2) that the result holds for flat paths. Thirdly, since the optimal fully-constrained generic path must be weakly less severe than the optimal semi-constrained path, then it must also follow that, as $\theta \longrightarrow 1^-$, the socially efficient equilibrium cannot be supported by any sustainable path.

Theorem 7 will work by arguing, firstly, that any generic optimal punishment path can be replaced by a semi-constrained quasi-flat path constructed by "flattening out" to the point 1 U-average from point 2 onwards. This newly constructed quasi-flat path will continue to fulfil the point 1 and point 2 cooperation conditions³⁵ $\lambda_1 \geq 0$ and $\lambda_2 \geq 0$. It will therefore continue to be in the set of semi-constrained paths, and will be as severe as the path it was generated from.

Applying Lemma 5 from the appendix, it will therefore be shown that the supportability constraint on the socially efficient initial path, given the use of the generic optimal path to punish a deviation, is broken as $\theta \longrightarrow 1^-$. Whether or not the optimal path is flat in a particular case, it is thus established that intermediate values of the coefficient of altruism are best able to allow a socially efficient equilibrium to be supported using globally optimal punishment paths. Theorem 7 will, as a result, be crucial to establishing the general applicability of the key results from section 8, which form the core contribution of this paper.

Definition 7. Let $\gamma : \Psi \longrightarrow \ddot{\Psi}$ be the function which constructs a quasi-flat path from a generic path by "flattening-out" the trigger levels from point 2 onwards to the point 1 U-average. This means that $\gamma_1(\psi) = \psi_1$ and $\gamma_2(\psi) = U^{-1}(U_1(\psi))$. Note also that $\phi(\gamma(\psi)) = \phi(\psi)$.³⁶

Lemma 2 below is an essential building block necessary for the proof of Theorem 7. It establishes that for any generic punishment path looking forward from any point k, the U-average, $U^{-1}(U_k(\psi))$, must lie equal to or above the Vaverage, $V^{-1}(V_k(\psi))$. It should be noted that the result from Lemma 2 is based on the assumption of risk neutral agents. Although mathematically analogous to differing levels of risk aversion, the result is in fact generated by differing attitudes towards the variability of benefit values under which punishment is inflicted. The sensitivity of a "neutral observer" to the "wastefulness" of punishment when watching others being punished is greater than the sensitivity of the individual being punished. This makes intuitive sense, since, to take the example of a fine, the person being fined is mainly affected in social utility terms by the fact that they are fined, whereas altruists who value the felicity of the

 $^{^{35}\}mathrm{The}$ co-operation conditions for point 3 and after are identical to those at point 2.

³⁶Note that $\gamma_k(\psi)$ is shorthand for $(\gamma(\psi))_k$ - the k^{th} trigger level in the path $\gamma(\psi)$.

person fined and the recipient of the revenue equally will be more sensitive to any deadweight loss from punishment. This result should therefore have wider implications in other models involving altruism and punishment.

Lemma 2. For any punishment path, at any point, the U-average is weakly greater than the V-average. $(\forall_{\psi}\forall_k: U^{-1}(U_k(\psi)) \geq V^{-1}(V_k(\psi)).)$

Proof. This result follows from an application of the principle of stochastic dominance, a key principle in the economics of risk [Rothschild & Stiglitz, 1970]. First, note from definition 3 that the per period average utility functions for the person being punished and a neutral observer can be rewritten as $U_k(\psi) \equiv \sum_{i=1}^{\infty} [p_i U(\psi_{i+k})]$ and $V_k(\psi) \equiv \sum_{i=1}^{\infty} [p_i V(\psi_{i+k})]$ where $p_i = (1-\delta)\delta^{i-1}$ and $\sum_{i=1}^{\infty} [p_i] = 1$. This means that $U(\psi)$ and $V(\psi)$ can be thought of as "expected utility functions", with the discount factor for each point along the path taking the role of probabilities for different outcomes of a lottery. The "expected value" of a particular path ψ looking forward from point k can then be defined as $E_k[\psi] = \sum_{i=1}^{\infty} [p_i\psi_{i+k}]$. Since this "expected value" is equal for both the punishee and the neutral observer, the discounted expected utility along a path is exactly analogous to the expected utility of a risky prospect. Therefore if we can show that the "neutral observer" is more "risk averse" than the punishee, the result of the lemma will follow.

The coefficient of absolute risk aversion $R_a = -\frac{U''}{U'}$ measures the degree of concavity of a utility function. If it is always higher for one function than another, then the corresponding agent is the more risk averse [Diamond & Stiglitz, 1974]. For the two types of agent under consideration (with utility functions $U(\psi)$ and $V(\psi)$ respectively), the CARA works out as the following.

$$R_a^U = \frac{\theta}{1 - \theta\psi} - \frac{g'(\psi)}{g(\psi)} \tag{42}$$

$$R_a^V = \frac{\theta}{\theta - \theta\psi} - \frac{g'(\psi)}{g(\psi)} \tag{43}$$

Since (43) is always unambigously greater than (42), the neutral observer is more "risk averse" than the punishee, and hence will always have a lower "certainty equivalent" from a given path looking forward, which, from definition 3, is precisely analogous to $V^{-1}(V_k(\psi))$ (as opposed to $U^{-1}(U_k(\psi))$ for the punishee).

The intuition for the result in Lemma 2 is also closely related to the "sacrifice ratio" derived in expression (29). Punishing within a certain "bracket" of values of π , with a fixed width, has a greater effect on $U^{-1}(U_k(\psi))$ (increasing the "stick") relative to $V^{-1}(V_k(\psi))$ (decreasing the "carrot") the higher the bracket. Given a particular value of $U^{-1}(U_k(\psi))$, $V^{-1}(V_k(\psi))$ is maximized when this is generated by a flat path looking forwards.

Lemma 3. If a punishment path is optimal, then the quasi-flat path constructed by "flattening it out" will be in the set of semi-constrained paths. (If a path ψ^* is optimal then $\forall_k : \lambda_k (\gamma(\psi^*)) \ge 0$, therefore $\gamma(\psi^*) \in \ddot{\Psi}_{\circ}$.) *Proof.* Firstly, observe that, given the result from Lemma 2, $U_0(\gamma(\psi^*)) = U_0(\psi^*)$ and $\forall_k : V_k(\gamma(\psi^*)) \ge V_1(\psi^*)$. To see this, note the following:

$$\left(\frac{\delta}{1-\delta}\right) U_0(\psi) \equiv \delta \int_{\psi_1}^1 (\theta\pi - 1)g(\pi)d\pi + \frac{\delta^2}{1-\delta} \int_{U^{-1}(U_1(\psi))}^1 (\theta\pi - 1)g(\pi)d\pi \right) d\pi$$

$$\left(\frac{\delta}{1-\delta}\right) U_0(\gamma(\psi)) \equiv \delta \int_{\gamma_1(\psi)}^1 (\theta\pi - 1)g(\pi)d\pi + \frac{\delta^2}{1-\delta} \int_{\gamma_2(\psi)}^1 (\theta\pi - 1)g(\pi)d\pi \right) d\pi$$

$$\left(\frac{\delta}{1-\delta}\right) V_1(\psi) \equiv \frac{\delta}{1-\delta} \int_{V^{-1}(V_1(\psi))}^1 (\theta\pi - \theta)g(\pi)d\pi$$

$$\forall_k : \left(\frac{\delta}{1-\delta}\right) V_k(\gamma(\psi)) \equiv \frac{\delta}{1-\delta} \int_{\gamma_2(\psi)}^1 (\theta\pi - \theta)g(\pi)d\pi$$

The claim is then straightforward to verify once it is noted that $\gamma_1(\psi) = \psi_1$ and $\gamma_2(\psi) = U^{-1}(U_1(\psi))$, since, by Lemma 2, $U^{-1}(U_1(\psi)) \ge V^{-1}(V_1(\psi))$.

Now we can proceed to note that:

$$\lambda_{1}(\gamma(\psi)) \equiv \left(\frac{\delta}{1-\delta}\right) V_{1}(\gamma(\psi)) - \left(\frac{\delta}{1-\delta}\right) U_{0}(\gamma(\psi)) + \gamma_{1}(\psi) - \theta \qquad (44)$$

$$\forall_{k\geq 2} : \lambda_k \left(\gamma \left(\psi \right) \right) \equiv \left(\frac{\delta}{1-\delta} \right) V_k \left(\gamma \left(\psi \right) \right) - \left(\frac{\delta}{1-\delta} \right) U_0 \left(\gamma \left(\psi \right) \right) + \gamma_2 \left(\psi \right) - \theta$$
(45)

Now take an optimal path ψ^* . By assumption, ψ^* is sustainable, and so:

$$\lambda_1\left(\psi^*\right) \equiv \left(\frac{\delta}{1-\delta}\right) V_1\left(\psi^*\right) - \left(\frac{\delta}{1-\delta}\right) U_0\left(\psi^*\right) + \psi_1^* - \theta \ge 0 \tag{46}$$

Given the result from Lemma 1 that $\psi_1^* \leq U^{-1}(U_1(\psi^*))$, along with all the observations noted so far, condition (46) is sufficient for (44) and (45) to be weakly positive for all relevant k.

Lemma 4. The optimal semi-constrained quasi-flat punishment path is at least as severe as the optimal fully-constrained generic punishment path. $(\phi \left(\ddot{\psi}_{\circ}^{*} \right) \geq \phi \left(\psi^{*} \right).)$

Proof. Consider the optimal generic path $\psi^* \in \Psi_{\bullet}$. By Lemma 3, $\gamma(\psi^*) \in \ddot{\Psi}_{\circ}$. Also, by definition, $\phi(\gamma(\psi^*)) = \phi(\psi^*)$. Therefore the most severe path in $\ddot{\Psi}_{\circ}$, $\sup \ddot{\Psi}_{\circ}$ must be at least as severe as the most severe path in Ψ_{\bullet} , $\sup \Psi_{\bullet}$. Thus $\phi(\sup \ddot{\Psi}_{\circ}) \ge \phi(\sup \Psi_{\bullet})$ and so $\phi\left(\ddot{\psi}_{\circ}^*\right) \ge \phi(\psi^*)$.

Theorem 7. As altruism becomes perfect, the optimal generic punishment path cannot support the socially efficient equilibrium, for any value of the discount rate.

$$(As \ \theta \longrightarrow 1^{-}, \ \kappa \left(\psi^*\right) < 0.)$$

Proof. First, note, from expressions (8) and (9) that:

$$\kappa\left(\psi^*\right) = \phi\left(\psi^*\right) + \theta - 1$$

By Lemma 4, this implies that:

$$\kappa\left(\psi^*\right) \le \phi\left(\ddot{\psi}^*_{\circ}\right) + \theta - 1 \tag{47}$$

Theorem 2, combined with Lemma 5 in the appendix, has already established that the RHS of expression (47) goes to 0^- as $\theta \longrightarrow 1^-$. Therefore the LHS of (47) must be strictly negative as $\theta \longrightarrow 1^-$.

Theorem 8. If the benefit distribution is sufficiently flat, then the optimal punishment path is flat.

$$(If \forall_{\pi} : g(\pi) < \frac{1}{\theta(1-\hat{\pi})} then \phi\left(\widetilde{\psi}^*\right) \ge \phi\left(\psi^*\right).)$$

Proof. From Lemma 5, we know that, if $\forall_{\pi} : g(\pi) < \frac{1}{\theta(1-\hat{\pi})}$, then $\ddot{\psi}^*_{\circ} \in \widetilde{\Psi}_{\bullet}$. Therefore, $\phi\left(\widetilde{\psi}^*\right) \geq \phi\left(\ddot{\psi}^*_{\circ}\right)$. Applying Lemma 4, the result follows straightforwardly.

11 Second-Best Equilibria

For any particular level of the coefficient of altruism θ , if the discount rate δ is low enough, so that players are sufficiently impatient, then the socially efficient initial path will not be supportable. There will still, however, exist a second-best optimal subgame-perfect Nash equilibrium, supported by the optimal punishment path, in the sense that the associated initial path maximizes efficiency by minimizing the range of benefit values for which punishment occurs. Given condition (35) then, as shown in Theorem 8, the optimal punishment path which supports such an optimal equilibrium will be flat. For the remainder of the paper we will, for the sake of simplicity, assume that this condition is satisfied.

Along the second-best optimal initial path, the punishment path will be initiated when an individual punishes for a value of the benefit below trigger level $\tilde{\psi}_h$. (Intuitively, the most attractive punishment opportunities will be the most difficult to deter along the initial path.) All players will then switch to a path where punishment is carried out above trigger level $\tilde{\psi}^*$, derived from the optimal flat punishment path (which is globally optimal, under the assumptions made). Assuming an interior solution for the optimal flat path where $\tilde{\psi}^* > \hat{\pi}, \tilde{\psi}^*$ and its total derivative with respect to θ will be respectively given by equations (10) and (11). The supportability constraint on the initial path will be:

$$\widetilde{\psi}_h - \theta \le \frac{\delta}{1-\delta} \int_{\widetilde{\psi}_h}^1 (\theta \pi - \theta) g(\pi) d\pi - \frac{\delta}{1-\delta} \int_{\widetilde{\psi}^*}^1 (\theta \pi - 1) g(\pi) d\pi \qquad (48)$$

When $\tilde{\psi}_h$ is set just high enough to make (48) bind, to give us the optimal second-best equilibrium, we will have the following implicit definition of $\tilde{\psi}_h^*$ in terms of itself, $\tilde{\psi}^*$, θ and δ :

$$\widetilde{\psi}_{h}^{*} = \theta + \frac{\delta}{1-\delta} \left(\int_{\widetilde{\psi}^{*}}^{1} g\left(\pi\right) d\pi - \theta \int_{\widetilde{\psi}^{*}}^{\widetilde{\psi}_{h}^{*}} \pi g\left(\pi\right) d\pi - \theta \int_{\widetilde{\psi}_{h}^{*}}^{1} g\left(\pi\right) d\pi \right)$$
(49)



Figure 37: Second-best equilibria where $\hat{\pi} = 0$ and $g(\pi) = 1$

Totally differentiating (49) with respect to θ and solving for $\frac{d\tilde{\psi}_{h}^{*}}{d\theta}$ gives us the following expression for the overall derivative of $\tilde{\psi}_{h}^{*}$ with respect to the coefficient of altruism:

$$\frac{d\tilde{\psi}_{h}^{*}}{d\theta} = \frac{(1-\delta) - \delta\left(\int_{\tilde{\psi}_{*}}^{\tilde{\psi}_{h}^{*}} g\left(\pi\right) \pi d\pi + \int_{\tilde{\psi}_{h}^{*}}^{1} g\left(\pi\right) d\pi\right) - \delta\left(1-\theta\tilde{\psi}^{*}\right) g\left(\tilde{\psi}^{*}\right) \frac{d\tilde{\psi}^{*}}{d\theta}}{(1-\delta) - \delta\theta\left(1-\tilde{\psi}_{h}^{*}\right) g\left(\tilde{\psi}_{h}^{*}\right)}$$
(50)

By observation of (10) and (49), as $\theta \longrightarrow 1^-$, $\tilde{\psi}^* \longrightarrow 1^-$ and therefore $\tilde{\psi}_h^* \longrightarrow 1^-$. Additionally, note that as $\theta \longrightarrow 1^-$, the RHS of (50) goes to 1. These observations imply that there must be a region where increasing the coefficient of altruism results in increasing $\tilde{\psi}_h^*$. Note also that as $\delta \longrightarrow 0$, the RHS of (50) again goes to 1. This implies that, for a low enough δ , increasing θ always increases $\tilde{\psi}_h^*$ (intuitively, this is where the future counts for sufficiently little that the severity effect, even in combination with the willingness effect, is never sufficient to outweigh the temptation effect).

Figure 37 shows the highest ψ_h^* which is supportable along the initial path given different values of θ (along the x-axis) and δ . It can be seen that the curve always has a slope of 1 as $\theta \longrightarrow 1$, and that the gradient is always positive for all θ when δ is low. Importantly, there is always a level of θ high enough but lower than 1 where $\tilde{\psi}_h^*$ falls below one and then back up to one as $\theta \longrightarrow 1$. This corresponds to the black region in figure 5 and derived analytically in Theorems 1 through 4. Finally, for high enough δ with a low enough θ , $\tilde{\psi}_h^*$ goes above one (i.e. the graph gets "cut off"). This is corresponds to the region where the first-best socially efficient equilibrium is supportable. An important conclusion to draw from figure 37 is that the efficiency loss from too high a level of altruism can be non-negligible. Although as $\theta \longrightarrow 1$, $\tilde{\psi}_h^* \longrightarrow 1$, there will exist intermediate levels of altruism where an increase in the coefficient of altruism to a higher intermediate level (which is still less than 1) could make the efficiency of the optimal second-best outcome significantly lower. Altruism is in many realistic cases a "double-edged sword" in the sequential punishment model, and too high a level of altruism will in general be socially detrimental in a significant manner.

12 Conclusion

This paper has taken two areas of economic theory, the modelling of altruistic preferences and the structure of optimal punishment paths, and shown how they can interact to produce interesting results in a new type of model, the sequential punishment model - a simple infinite-move sequential game with perfect information and discounting - where players move by choosing whether or not to take opportunities to inflict harm upon others with benefit to themselves. Essentially, the model is an abstract representation of the fundamentally vicarious nature of human interaction in any kind of society, whatever its organizational principles. The central implication of the analysis is that excessive altruism will interfere detrimentally with punishment systems, "denting" them in such a manner that social welfare is reduced compared to a situation with lower altruism.

The concepts of willingness, severity and temptation effects which were used to analyse the subgame perfect equilibria of the sequential punishment model, and to establish the existence of a socially optimal level of altruism, should present themselves in other contexts. They would, for instance, be highly relevant to the analysis of optimal taxation, the economic theory of criminal punishment, and to issues surrounding the evolution of altruistic preferences; see, for example, "Punishment and the Potency of Group Selection" [Povey, 2010].

Sections 1 through 7 set up the notation for the sequential punishment model, and related this to Abreu's framework of optimal penal codes, originally developed for repeated stage games. The main body of novel results for this paper is in sections 8, 9 and 10, which progressively generalize the core result of the paper - that as altruism becomes perfect the socially efficient equilibrium breaks down - to the equilibria supportable by the optimal flat punishment path, the optimal quasi-flat path and, finally, the optimal generic punishment path.

In the sequential punishment model, the interaction between the severity, willingness and temptation effects can be conclusively seen to lead to the result that an intermediate "Goldilocks" level of altruism is socially optimal. This result was established (initially for equilibria supportable by flat paths), in Theorem 2. The key intuition for this result is that, for a low enough value of the discount rate δ , the temptation effect must initially dominate the severity and willingness effects as θ is reduced from below 1. Since social efficiency is only barely supportable at $\theta = 1$, the constraint for supportability must be broken as $\theta \longrightarrow 1^-$.

If individuals are sufficiently impatient, excessive malevolence will, on the other hand, be socially damaging, due to the dominance of the temptation effect over the severity effect as $\theta \longrightarrow -\infty$. In contrast however, if the discount rate δ is high enough, then even infinite malevolence does not break the socially efficient equilibrium, because the severity effect will outweigh the temptation effect as $\theta \longrightarrow -\infty$ (Theorem 3). The sequential punishment model also demonstrates, therefore, an asymmetry, in that high altruism is in general more damaging than extreme malevolence. With sufficiently effective monitoring (leading to a δ close to 1) malevolent preferences are not of concern from a social welfare perspective, but excessively altruistic ones remain so.

The analysis in section 9 involved the investigation of quasi-flat paths, which maintain a simple structure of carrot-and-stick punishment, and are an interesting illustration of the general principles governing optimal penal codes in and of themselves. Quasi flat paths are a specific feature which emerges from the application of optimal penal codes to the sequential punishment model. They occur because the presence of partially altruistic preferences leads to a "flattening-out" of the tail of the optimal punishment path. This is a key difference between the nature of optimal penal codes in the sequential punishment model and those found in the infinitely-repeated Cournot and Bertrand simultaneous oligopoly stage games.

The optimality of quasi-flat paths is driven by the higher relative concavity of the "inter-temporal" utility function of a "neutral observer" relative to the punishee in a punishment equilibrium. This is also an interesting result which should have wider implications. The broader society is more sensitive to the inefficiencies brought about by random variations in punishment technology than the individual being punished. This is because the individual cares primarily about the fact that they are punished, and so places less relative weighting on the deadweight loss to society from punishment. A similar phenomenon should emerge in any model with altruistic agents and punishment technologies. There are therefore potential applications in optimal taxation theory and in rational choice models of criminality and punishment.

Theorem 7 completes the generalization of the results from section 8 so that they apply to any continuous distribution of the benefit with support between $\hat{\pi}$ and 1 (where $0 \leq \hat{\pi} < 1$), and when the optimal generic punishment path is used to support the socially efficient equilibrium by punishing any deviation from the socially efficient initial path. This is a satisfyingly general result, although it does require ruling out situations where individuals would actually *enjoy* being harmed. Another limitation of the treatment of the model in this paper is that the case where $\theta \geq 1$, so that individuals are "martyrs", who care about others more than themselves, has been excluded from the outset.

Section 11 demonstrated that there is a potentially significant loss of social efficiency from the detrimental effects of too high a level of altruism on the social incentive systems used to induce co-operation via the punishment of transgressors. This was an important final piece of the argument, as it is necessary not only to show that too high a level of altruism will break the supportability constraint on the socially efficient outcome, but also that the loss of social welfare in the resulting second-best world will often be non-negligible.

13 Appendix

Lemma 5. (a) As altruism becomes perfect, the optimal semi-constrained quasiflat path becomes flat. (b) If the benefit distribution is sufficiently flat, the optimal semi-constrained quasi-flat path becomes flat. ((a) As $\theta \longrightarrow 1^-$, $\ddot{\psi}^*_{\circ} \in \tilde{\Psi}$. (b) If $\forall_{\pi} : g(\pi) < \frac{1}{\theta(1-\hat{\pi})}$ then $\ddot{\psi}^*_{\circ} \in \tilde{\Psi}$.)

Proof. The optimal semi-constrained quasi flat path $\ddot{\psi}^*_{\circ}$ can be found by solving a constrained optimization problem [Lagrange, 1806] [Simon & Blume, 1994] using the following Lagrangean function:

$$\mathcal{L} = -\left(\frac{\delta}{1-\delta}\right) U\left(\ddot{\psi}\right) + \bar{\lambda}_1 \left(\left(\frac{\delta}{1-\delta}\right) V\left(\ddot{\psi}\right) - \left(\frac{\delta}{1-\delta}\right) U\left(\ddot{\psi}\right) + \ddot{\psi}_1 - \theta\right) \\ + \bar{\lambda}_2 \left(\left(\frac{\delta}{1-\delta}\right) V\left(\ddot{\psi}\right) - \left(\frac{\delta}{1-\delta}\right) U\left(\ddot{\psi}\right) + \ddot{\psi}_2 - \theta\right) \\ + \bar{\eta}_1 \left(1 - \ddot{\psi}_1\right) + \bar{\eta}_2 \left(1 - \ddot{\psi}_2\right) + \bar{\zeta}_1 \left(\ddot{\psi}_1 - \hat{\pi}\right) + \bar{\zeta}_2 \left(\ddot{\psi}_2 - \hat{\pi}\right)$$

We know that the following first order conditions must hold at the constrained maximum: $\frac{\partial \mathcal{L}}{\partial \dot{\psi}_1} = 0$, $\frac{\partial \mathcal{L}}{\partial \dot{\psi}_2} = 0$, $\frac{\partial \mathcal{L}}{\partial \lambda_1} \ge 0$, $\bar{\lambda}_1 \ge 0$, $\frac{\partial \mathcal{L}}{\partial \lambda_1} \bar{\lambda}_1 = 0$, $\frac{\partial \mathcal{L}}{\partial \lambda_2} \ge 0$, $\bar{\lambda}_2 \ge 0$, $\frac{\partial \mathcal{L}}{\partial \lambda_2} \bar{\lambda}_2 = 0$, $\frac{\partial \mathcal{L}}{\partial \bar{\eta}_1} \ge 0$, $\bar{\eta}_1 \ge 0$, $\frac{\partial \mathcal{L}}{\partial \bar{\eta}_1} \bar{\eta}_1 = 0$, $\frac{\partial \mathcal{L}}{\partial \bar{\eta}_2} \ge 0$, $\bar{\eta}_2 \ge 0$, $\frac{\partial \mathcal{L}}{\partial \bar{\chi}_1} \bar{\lambda}_2 = 0$, $\frac{\partial \mathcal{L}}{\partial \bar{\chi}_1} \bar{\chi}_1 = 0$, $\frac{\partial \mathcal{L}}{\partial \bar{\chi}_2} \bar{\chi}_2 = 0$, $\frac{\partial \mathcal{L}}{\partial \bar{\chi}_1} \bar{\chi}_1 = 0$, $\frac{\partial \mathcal{L}}{\partial \bar{\chi}_2} \bar{\chi}_2 = 0$, $\frac{\partial \mathcal{L}}{\partial \bar{\chi}_1} \bar{\chi}_2 = 0$, $\frac{\partial \mathcal{L}}{\partial \bar{\chi}_1} \bar{\chi}_1 = 0$, $\frac{\partial \mathcal{L}}{\partial \bar{\chi}_2} \bar{\chi}_2 = 0$, $\frac{\partial \mathcal{L}}{\partial \bar{\chi}_2} \bar{\chi}_2 = 0$. Substituting in to the Lagrangean from equations (4) and (5) and rearranging

$$\mathcal{L} = -\delta \left(1 + \bar{\lambda}_1 + \bar{\lambda}_2\right) \int_{\ddot{\psi}_1}^1 (\theta \, \pi - 1) \, g \left(\pi\right) d\pi - \frac{\delta^2 \left(1 + \bar{\lambda}_1 + \bar{\lambda}_2\right) \int_{\ddot{\psi}_2}^1 (\theta \, \pi - 1) \, g \left(\pi\right) d\pi}{1 - \delta} + \frac{\delta \left(\bar{\lambda}_1 + \bar{\lambda}_2\right) \int_{\ddot{\psi}_2}^1 (\theta \, \pi - \theta) \, g \left(\pi\right) d\pi}{1 - \delta} + \left(\bar{\zeta}_1 + \bar{\lambda}_1 - \bar{\eta}_1\right) \ddot{\psi}_1 + \left(\bar{\zeta}_2 + \bar{\lambda}_2 - \bar{\eta}_2\right) \ddot{\psi}_2 - \left(\bar{\zeta}_1 + \bar{\zeta}_2\right) \hat{\pi} - \left(\bar{\lambda}_1 + \bar{\lambda}_2\right) \theta + \bar{\eta}_1 + \bar{\eta}_2$$

The relevant derivatives for the first order conditions are:

$$\frac{d\mathcal{L}}{d\dot{\psi}_{1}} = \delta \left(1 + \bar{\lambda}_{1} + \bar{\lambda}_{2}\right) \left(\theta \ddot{\psi}_{1} - 1\right) g\left(\ddot{\psi}_{1}\right) + \bar{\zeta}_{1} + \bar{\lambda}_{1} - \bar{\eta}_{1}$$

$$\frac{d\mathcal{L}}{d\ddot{\psi}_{2}} = \frac{\delta^{2} (1 + \bar{\lambda}_{1} + \bar{\lambda}_{2}) (\theta \ddot{\psi}_{2} - 1) g(\ddot{\psi}_{2})}{1 - \delta} - \frac{\delta (\bar{\lambda}_{1} + \bar{\lambda}_{2}) (\theta \ddot{\psi}_{2} - \theta) g(\ddot{\psi}_{2})}{1 - \delta} + \bar{\lambda}_{2} - \bar{\eta}_{2} + \bar{\zeta}_{2}$$

$$\frac{d\mathcal{L}}{d\lambda_{1}} = -\delta \int_{\ddot{\psi}_{1}}^{1} (\theta \pi - 1) g(\pi) d\pi - \frac{\delta^{2} \int_{\dot{\psi}_{2}}^{1} (\theta \pi - 1) g(\pi) d\pi}{1 - \delta} + \frac{\delta \int_{\dot{\psi}_{2}}^{1} (\theta \pi - \theta) g(\pi) d\pi}{1 - \delta} + \frac{\psi_{1} - \theta}{1 - \delta}$$

$$\frac{d\mathcal{L}}{d\lambda_{2}} = -\delta \int_{\ddot{\psi}_{1}}^{1} (\theta \pi - 1) g(\pi) d\pi - \frac{\delta^{2} \int_{\dot{\psi}_{2}}^{1} (\theta \pi - 1) g(\pi) d\pi}{1 - \delta} + \frac{\delta \int_{\dot{\psi}_{2}}^{1} (\theta \pi - \theta) g(\pi) d\pi}{1 - \delta} + \frac{\psi_{2} - \theta}{1 - \delta}$$

In order to prove the result, we must exhaustively consider the various possibilities for the different constraints, and whether or not they bind. We also see that this procedure relates straightforwardly and directly to the taxonomy of quasi-flat paths laid out in section 9.1.

13.0.1 Maximal paths

Suppose, first of all, that $0 < \bar{\zeta}_1$ and $0 < \bar{\zeta}_2$. This means that both "lower constraints" on the trigger levels bind, so that $\ddot{\psi}_1 = \hat{\pi}$ and $\ddot{\psi}_2 = \hat{\pi}$. Therefore, the upper constraints cannot possibly bind, and so $\bar{\eta}_1 = 0$ and $\bar{\eta}_2 = 0$. The requirements that $\frac{\partial \mathcal{L}}{\partial \lambda_1} \ge 0$ and $\frac{\partial \mathcal{L}}{\partial \lambda_2} \ge 0$ therefore both become the following:

$$0 \le \int_{\hat{\pi}}^{1} \frac{g\left(\pi\right)\delta\left(1-\theta\right)}{1-\delta} d\pi - \theta + \hat{\pi}$$

This simplifies to give $\theta \leq \delta + (1 - \delta) \hat{\pi}$. So, only when the inequality condition from Theorem 1 is either not fulfilled, or just fulfilled with equality, can we have a maximal path. Rearranging this condition for δ , we know that the maximal path will be the constrained optimum if and only if:

$$\frac{\theta - \hat{\pi}}{1 - \hat{\pi}} \le \delta \tag{51}$$

If this condition is satisfied as a strict inequality, then, since $\frac{\partial \mathcal{L}}{\partial \bar{\lambda}_1} > 0$ and $\frac{\partial \mathcal{L}}{\partial \bar{\lambda}_2} > 0$, we know that $\bar{\lambda}_1 = 0$ and $\bar{\lambda}_2 = 0$. The first order conditions on $\ddot{\psi}_1$ and $\ddot{\psi}_2$ therefore yield the following solution for $\bar{\zeta}_1$ and $\bar{\zeta}_2$.

$$\frac{\partial \mathcal{L}}{\partial \ddot{\psi}_1} = \delta \left(\theta \,\hat{\pi} - 1\right) g\left(\hat{\pi}\right) + \bar{\zeta}_1 = 0$$
$$\frac{\partial \mathcal{L}}{\partial \ddot{\psi}_2} = \frac{\left(\theta \,\hat{\pi} - 1\right) \delta^2 g\left(\hat{\pi}\right)}{1 - \delta} + \bar{\zeta}_2 = 0$$
$$\bar{\zeta}_1 = \delta \left(1 - \theta \,\hat{\pi}\right) g\left(\hat{\pi}\right) \qquad \bar{\zeta}_2 = \frac{\left(1 - \theta \,\hat{\pi}\right) \delta^2 g\left(\hat{\pi}\right)}{1 - \delta}$$

The solutions for $\overline{\zeta}_1$ and $\overline{\zeta}_2$ are positive, which is consistent with our initial assumptions.

If we take the limit of the LHS of condition (51), we can see that a maximal path is not possible as $\theta \longrightarrow 1^-$, since, by assumption, $\delta < 1$.

$$\lim_{\theta \to 1^{-}} \left\{ \frac{\theta - \hat{\pi}}{1 - \hat{\pi}} \right\} = 1^{-}$$

Suppose instead that $\bar{\zeta}_1 = 0$ and $0 < \bar{\zeta}_2$. This implies that $\ddot{\psi}_2 = \hat{\pi}$ and that $\bar{\eta}_2 = 0$. The requirements that $\frac{\partial \mathcal{L}}{\partial \lambda_1} \ge 0$ and $\frac{\partial \mathcal{L}}{\partial \lambda_2} \ge 0$ become the following:

$$0 \leq -\int_{\hat{\pi}}^{1} \frac{g\left(\pi\right)\delta\left(-\theta\,\pi+\theta+\delta\,\theta\,\pi-\delta\right)}{1-\delta} d\pi - \int_{\ddot{\psi}_{1}}^{1} \delta\left(\theta\,\pi-1\right)g\left(\pi\right)d\pi + \ddot{\psi}_{1} - \theta \tag{52}$$

$$0 \leq -\int_{\hat{\pi}}^{1} \frac{g\left(\pi\right)\delta\left(-\theta\,\pi+\theta+\delta\,\theta\,\pi-\delta\right)}{1-\delta} d\pi - \int_{\ddot{\psi}_{1}}^{1} \delta\left(\theta\,\pi-1\right)g\left(\pi\right)d\pi + \hat{\pi} - \theta \tag{53}$$

Assuming that $\ddot{\psi}_1 > \hat{\pi}$ (otherwise we would have a maximal path, as above), then, if the second of these inequalities is fulfilled, so is the first. This means that $\bar{\lambda}_1 = 0$, since $\frac{dL}{d\lambda_1} > 0$. The second inequality can be rearranged to give the following:

$$\frac{\int_{\bar{\psi}_1}^1 \delta \left(1-\delta\right) \left(\theta \,\pi-1\right) g\left(\pi\right) d\pi - \delta \,\theta \,\bar{\pi} + \delta^2 \theta \,\bar{\pi} - \delta^2 + \theta}{1-\delta} \le \hat{\pi} \tag{54}$$

The following inequality, however, must be satisfied:

$$\delta (1-\delta) (\theta \bar{\pi} - 1) \leq \int_{\ddot{\psi}_1}^1 \delta (1-\delta) (\theta \pi - 1) g(\pi) d\pi$$

Therefore, condition (54) can only be fulfilled if $\theta \leq \delta + \hat{\pi}(1-\delta)$, which is the same sufficient condition already derived for the maximal path to be optimal. Hence this case collapses into the previous one.

13.0.2 Quasi-maximal paths

Suppose that $0 < \bar{\zeta}_1$ and $\bar{\zeta}_2 = 0$. This implies that $\ddot{\psi}_1 = \hat{\pi}$ and $\bar{\eta}_1 = 0$. The requirements that $\frac{\partial \mathcal{L}}{\partial \lambda_1} \ge 0$ and $\frac{\partial \mathcal{L}}{\partial \lambda_2} \ge 0$ become the following:

$$0 \leq -\int_{\ddot{\psi}_2}^{1} \frac{g\left(\pi\right)\delta\left(-\theta\,\pi+\theta+\delta\,\theta\,\pi-\delta\right)}{1-\delta} d\pi - \int_{\hat{\pi}}^{1}\delta\left(\theta\,\pi-1\right)g\left(\pi\right)d\pi + \hat{\pi} - \theta \tag{55}$$

$$0 \leq -\int_{\ddot{\psi}_2}^{1} \frac{g\left(\pi\right)\delta\left(-\theta\,\pi+\theta+\delta\,\theta\,\pi-\delta\right)}{1-\delta} d\pi - \int_{\hat{\pi}}^{1} \delta\left(\theta\,\pi-1\right)g\left(\pi\right)d\pi + \ddot{\psi}_2 - \theta \tag{56}$$

Assuming that $\ddot{\psi}_2 > \hat{\pi}$ (otherwise we would have a maximal path), then if the first of these inequalities is fulfilled, so is the second. Therefore, $\bar{\lambda}_2 = 0$. Condition (55) can be rearranged to give:

$$\int_{\ddot{\psi}_2}^1 \frac{\delta g(\pi) \left(-\theta \pi + \theta + \delta \theta \pi - \delta\right)}{1 - \delta} d\pi + \delta \theta \bar{\pi} - \delta + \theta \le \hat{\pi}$$
(57)

The first order conditions on $\ddot{\psi}_1$ and $\ddot{\psi}_2$ become the following:

$$\frac{\partial \mathcal{L}}{\partial \ddot{\psi}_1} = \delta \left(1 + \bar{\lambda}_1 \right) \left(\theta \, \hat{\pi} - 1 \right) g \left(\hat{\pi} \right) + \bar{\zeta}_1 + \bar{\lambda}_1 = 0$$
$$\frac{\partial \mathcal{L}}{\partial \ddot{\psi}_2} = \frac{\delta^2 \left(1 + \bar{\lambda}_1 \right) \left(\theta \, \ddot{\psi}_2 - 1 \right) g \left(\ddot{\psi}_2 \right)}{1 - \delta} - \frac{\delta \, \bar{\lambda}_1 \left(\theta \, \ddot{\psi}_2 - \theta \right) g \left(\ddot{\psi}_2 \right)}{1 - \delta} - \bar{\eta}_2 = 0$$

We now need to consider two possible cases, one where $\bar{\lambda}_1 = 0$ and one where $\bar{\lambda}_1 > 0$. Taking the first, the first order conditions on $\ddot{\psi}_1$ and $\ddot{\psi}_2$ become:

$$\delta \left(\theta \,\hat{\pi} - 1\right) g\left(\hat{\pi}\right) + \bar{\zeta}_1 = 0$$
$$\frac{\delta^2 \left(\theta \,\ddot{\psi}_2 - 1\right) g\left(\ddot{\psi}_2\right)}{1 - \delta} - \bar{\eta}_2 = 0$$

The second of these cannot possibly be fulfilled, since the LHS is unambiguously negative. Therefore, it must be the case that $0 < \overline{\lambda}_1$. The two first order conditions then yield the following solutions for $\overline{\lambda}_1$:

$$\bar{\lambda}_1 = -\frac{\delta \left(\theta \,\hat{\pi} - 1\right) g\left(\hat{\pi}\right) + \bar{\zeta}_1}{1 + \delta \left(\theta \,\hat{\pi} - 1\right) g\left(\hat{\pi}\right)} \tag{58}$$

$$\bar{\lambda}_1 = -\frac{\delta^2 \left(\theta \,\ddot{\psi}_2 - 1\right) g \left(\ddot{\psi}_2\right) - \bar{\eta}_2 (1 - \delta)}{\delta g \left(\ddot{\psi}_2\right) \left(\left(\theta \,\ddot{\psi}_2 - 1\right) \delta - \theta \left(\ddot{\psi}_2 - 1\right)\right)}$$
(59)

In order for the second of these expressions to be weakly positive, we must have that:

$$\delta < \frac{\theta - \theta \psi_2}{1 - \theta \, \ddot{\psi}_2} \tag{60}$$

Now consider whether $\bar{\eta}_2 = 0$ or $\bar{\eta}_2 > 0$. If $\bar{\eta}_2 > 0$ then $\hat{\psi}_2 = 1$ and so (60) cannot be fulfilled. Hence $\bar{\eta}_2 = 0$. In that case, solving (58) and (59) simultaneously for $\bar{\zeta}_1$ yields:

$$\bar{\zeta}_{1} = \frac{\delta\left(\theta \left(\theta \,\hat{\pi} - 1\right)\left(\ddot{\psi}_{2} - 1\right)g\left(\hat{\pi}\right) + \left(\theta \,\ddot{\psi}_{2} - 1\right)\right)}{\left(\theta \,\ddot{\psi}_{2} - 1\right)\delta - \theta \,\left(\ddot{\psi}_{2} - 1\right)}$$

Given condition (60), the denominator is positive. The numerator will be positive if and only if the following condition holds:

$$g\left(\hat{\pi}\right) \ge \frac{1 - \theta \,\tilde{\psi}_2}{\theta \,\left(1 - \ddot{\psi}_2\right) \left(1 - \theta \,\hat{\pi}\right)} \tag{61}$$

Since condition (55) must be satisfied with equality, we have that:

$$-\int_{\ddot{\psi}_2}^1 \frac{\delta g\left(\pi\right)\left(-\theta \pi + \theta + \delta \theta \pi - \delta\right)}{1 - \delta} d\pi = \delta \left(\theta \,\bar{\pi} - 1\right) - \hat{\pi} + \theta$$

Taking limits of both sides as $\theta \longrightarrow 1^-$ yields:

$$-\delta \int_{\ddot{\psi}_{2}}^{1} (1-\pi) g(\pi) d\pi = 1 - \hat{\pi} - \delta (1-\bar{\pi})$$

This is not possible, since the LHS is unambiguously negative whilst the RHS is unambiguously positive. Therefore, a quasi-maximal path is not possible as $\theta \longrightarrow 1^{-}$.

13.0.3 Minimal paths

Now take the case where $0 < \bar{\eta}_1$ and $0 < \bar{\eta}_2$. This implies that $\ddot{\psi}_1 = 1$, $\ddot{\psi}_2 = 1$, $\bar{\zeta}_1 = 0$ and $\bar{\zeta}_2 = 0$. Thus we have a minimal path, where no punishment occurs at all. $\frac{\partial \mathcal{L}}{\partial \lambda_1}$ and $\frac{\partial \mathcal{L}}{\partial \lambda_2}$ both simplify to give $1 - \theta$. Clearly, therefore, $\bar{\lambda}_1 = 0$ and $\bar{\lambda}_2 = 0$. The first order conditions on $\ddot{\psi}_1$ and $\ddot{\psi}_2$ become:

$$\frac{\partial \mathcal{L}}{\partial \ddot{\psi}_1} = -\delta \left(1 - \theta\right) g\left(1\right) - \bar{\eta}_1 = 0$$
$$\frac{\partial \mathcal{L}}{\partial \ddot{\psi}_2} = -\frac{\delta^2 \left(1 - \theta\right) g\left(1\right)}{1 - \delta} - \bar{\eta}_2 = 0$$

Neither of these can possibly be fulfilled for $\theta < 1$, so a minimal path can never be optimal.

13.0.4 Quasi-minimal paths

Now we consider the case where $\bar{\eta}_1 = 0$ and $0 < \bar{\eta}_2$. This implies that $\ddot{\psi}_2 = 1$ and $\bar{\zeta}_2 = 0$. The requirements that $\frac{\partial \mathcal{L}}{\partial \lambda_1} \ge 0$ and $\frac{\partial \mathcal{L}}{\partial \lambda_2} \ge 0$ become the following:

$$\begin{split} 0 &\leq \int_{\vec{\psi}_1}^1 -\delta \, \left(\theta \, \pi - 1\right) g \left(\pi\right) d\pi + \vec{\psi}_1 - \theta \\ 0 &\leq \int_{\vec{\psi}_1}^1 -\delta \, \left(\theta \, \pi - 1\right) g \left(\pi\right) d\pi + 1 - \theta \end{split}$$

Assuming $\ddot{\psi}_1 < 1$ (otherwise we would have a minimal path), then if the first of these is fulfilled, so is the second. Therefore, $\frac{\partial \mathcal{L}}{\partial \lambda_2} > 0$ and $\bar{\lambda}_2 = 0$. The first order conditions on $\ddot{\psi}_1$ and $\ddot{\psi}_2$ thus become:

$$\frac{\partial \mathcal{L}}{\partial \ddot{\psi}_1} = \delta \left(1 + \bar{\lambda}_1 \right) \left(\theta \, \ddot{\psi}_1 - 1 \right) g \left(\ddot{\psi}_1 \right) + \bar{\zeta}_1 + \bar{\lambda}_1 = 0 \tag{62}$$

$$\frac{\partial \mathcal{L}}{\partial \ddot{\psi}_2} = -\frac{\delta^2 \left(1 + \bar{\lambda}_1\right) \left(1 - \theta\right) g\left(1\right)}{1 - \delta} - \bar{\eta}_2 = 0 \tag{63}$$

The LHS of the second expression is unambiguously negative, so this shows that quasi-minimal paths cannot be optimal.

Now consider the case where $0 < \bar{\eta}_1$ and $\bar{\eta}_2 = 0$. This implies that $\ddot{\psi}_1 = 1$ and $\bar{\zeta}_1 = 0$. The requirements that $\frac{\partial \mathcal{L}}{\partial \lambda_1} \ge 0$ and $\frac{\partial \mathcal{L}}{\partial \lambda_2} \ge 0$ become the following:

$$0 \leq -\int_{\ddot{\psi}_2}^{1} \frac{\delta g(\pi) \left(-\theta \pi + \theta + \delta \theta \pi - \delta\right)}{1 - \delta} d\pi + 1 - \theta$$
$$0 \leq -\int_{\ddot{\psi}_2}^{1} \frac{\delta g(\pi) \left(-\theta \pi + \theta + \delta \theta \pi - \delta\right)}{1 - \delta} d\pi + \ddot{\psi}_2 - \theta$$

If the second is fulfilled, so is the first. Thus, $\frac{\partial \mathcal{L}}{\partial \bar{\lambda}_1} > 0$ and $\bar{\lambda}_1 = 0$. The first order conditions $\ddot{\psi}_1$ and $\ddot{\psi}_2$ then become:

$$\frac{\partial \mathcal{L}}{\partial \ddot{\psi}_1} = -\delta \left(1 + \bar{\lambda}_2\right) (1 - \theta) g (1) - \bar{\eta}_1 = 0$$
$$\frac{\partial \mathcal{L}}{\partial \ddot{\psi}_2} = \delta \left(\frac{\left(\delta \theta \, \ddot{\psi}_2 - \delta - \theta \, \ddot{\psi}_2 + \theta\right) \bar{\lambda}_2}{1 - \delta} + \frac{\delta \left(\theta \, \ddot{\psi}_2 - 1\right)}{1 - \delta}\right) g \left(\ddot{\psi}_2\right) + \bar{\lambda}_2 + \bar{\zeta}_2 = 0$$

The first condition cannot possibly be satisfied, so this case is also impossible.

0

13.0.5 Carrot-maximized paths

The case where $\bar{\zeta}_1 = 0$, $\bar{\zeta}_2 = 0$, $\bar{\eta}_1 = 0$, $\bar{\eta}_2 = 0$ is the most interesting one, where the possibilities of a carrot-maximized and flat optimal quasi-flat path can occur. The requirements that $\frac{\partial \mathcal{L}}{\partial \lambda_1} \geq 0$ and $\frac{\partial \mathcal{L}}{\partial \lambda_2} \geq 0$ in this case become the following:

$$0 \leq -\int_{\ddot{\psi}_{2}}^{1} \frac{\delta g\left(\pi\right)\left(-\theta \pi + \theta + \delta \theta \pi - \delta\right)}{1 - \delta} d\pi + \int_{\ddot{\psi}_{1}}^{1} -\delta \left(\theta \pi - 1\right) g\left(\pi\right) d\pi + \ddot{\psi}_{1} - \theta \tag{64}$$

$$0 \leq -\int_{\ddot{\psi}_{2}}^{1} \frac{\delta g(\pi) \left(-\theta \pi + \theta + \delta \theta \pi - \delta\right)}{1 - \delta} d\pi + \int_{\ddot{\psi}_{1}}^{1} -\delta \left(\theta \pi - 1\right) g(\pi) d\pi + \ddot{\psi}_{2} - \theta$$
(65)

By Lemma 1, we know that an optimal path must have the property that $\ddot{\psi}_1 \leq \ddot{\psi}_2$. There are therefore two cases to consider. First, suppose that $\ddot{\psi}_1 \neq \ddot{\psi}_2$ and so $\ddot{\psi}_1 < \ddot{\psi}_2$. In that case, only (64) can bind, and so $\frac{\partial \mathcal{L}}{\partial \lambda_2} > 0$ and $\bar{\lambda}_2 = 0$. The first order conditions on $\ddot{\psi}_1$ and $\ddot{\psi}_2$ then become the following:

$$\frac{\partial \mathcal{L}}{\partial \ddot{\psi}_1} = \left(\delta \left(\theta \,\ddot{\psi}_1 - 1\right)g\left(\ddot{\psi}_1\right) + 1\right)\bar{\lambda}_1 + \delta \left(\theta \,\ddot{\psi}_1 - 1\right)g\left(\ddot{\psi}_1\right) = 0$$
$$\frac{\partial \mathcal{L}}{\partial \ddot{\psi}_2} = \left(-\theta \,\ddot{\psi}_2 + \frac{(\theta - \delta)}{1 - \delta}\right)\delta g\left(\ddot{\psi}_2\right)\bar{\lambda}_1 + \frac{\delta^2}{1 - \delta}\left(\theta \,\ddot{\psi}_2 - 1\right)g\left(\ddot{\psi}_2\right) = 0$$

Clearly, $\bar{\lambda}_1$ must be positive, or the first of these could not possibly be fulfilled. Solving these two equations respectively for $\bar{\lambda}_1$ gives us:

$$\bar{\lambda}_1 = -\frac{\delta \left(\theta \ddot{\psi}_1 - 1\right) g \left(\ddot{\psi}_1\right)}{\delta \left(\theta \ddot{\psi}_1 - 1\right) g \left(\ddot{\psi}_1\right) + 1}$$
$$\bar{\lambda}_1 = -\frac{\left(\theta \ddot{\psi}_2 - 1\right) \delta}{-\theta \left(1 - \delta\right) \ddot{\psi}_2 - \delta + \theta}$$

Finally, solving these two equations simultaneously for $g(\ddot{\psi}_1)$ gives us the following necessary condition for a carrot-maximized path, which, once rearranged to make $\ddot{\psi}_2$ the subject, verifies condition (38) derived by a more intuitive argument in the main text in section 9.2:

$$g\left(\ddot{\psi}_{1}\right) = \frac{1-\theta\,\ddot{\psi}_{2}}{\theta\,\left(1-\ddot{\psi}_{2}\right)\left(1-\theta\,\ddot{\psi}_{1}\right)}\tag{66}$$

Since $\overline{\lambda}_1 > 0$, condition (64) must be satisfied with equality. This gives us:

$$0 = -\int_{\ddot{\psi}_2}^1 \frac{\delta g\left(\pi\right)\left(-\theta \pi + \theta + \delta \theta \pi - \delta\right)}{1 - \delta} d\pi - \int_{\ddot{\psi}_1}^1 \delta \left(\theta \pi - 1\right) g\left(\pi\right) d\pi + \ddot{\psi}_1 - \theta$$

Taking the limit of both sides as $\theta \longrightarrow 1^-$ yields:

$$0 = \delta \int_{\ddot{\psi}_1}^{\ddot{\psi}_2} (1 - \pi) g(\pi) \, d\pi - (1 - \ddot{\psi}_1)$$

This condition cannot possibly be fulfilled, so a carrot-maximized path is also impossible as $\theta \longrightarrow 1^-$.

13.0.6 Flat paths

Consider, finally, the case where $\bar{\zeta}_1 = 0$, $\bar{\zeta}_2 = 0$, $\bar{\eta}_1 = 0$, $\bar{\eta}_2 = 0$ and $\ddot{\psi}_1 = \ddot{\psi}_2 = \ddot{\psi}$. The requirements that $\frac{\partial \mathcal{L}}{\partial \lambda_1} \ge 0$ and $\frac{\partial \mathcal{L}}{\partial \lambda_2} \ge 0$ then both become:

$$0 \le \int_{\tilde{\psi}}^{1} \frac{\delta g(\pi) (1-\theta)}{1-\delta} d\pi - \theta + \tilde{\psi}$$

This must hold with equality, yielding condition (10) from the main text, which defines the optimal flat path $\tilde{\psi}^*$.

The first order conditions on $\ddot{\psi}_1$ and $\ddot{\psi}_2$ become the following:

$$\left(\delta\left(\theta\,\tilde{\psi}-1\right)g\left(\tilde{\psi}\right)+1\right)\bar{\lambda}_{1}+\delta\left(\theta\,\tilde{\psi}-1\right)g\left(\tilde{\psi}\right)\left(1+\bar{\lambda}_{2}\right)=0$$

$$\frac{\delta\left(\delta\,\theta\,\tilde{\psi}-\delta-\theta\,\tilde{\psi}+\theta\right)g\left(\tilde{\psi}\right)\bar{\lambda}_{1}}{1-\delta}+\left(1+\frac{\delta\left(\delta\,\theta\,\tilde{\psi}-\delta-\theta\,\tilde{\psi}+\theta\right)g\left(\tilde{\psi}\right)}{1-\delta}\right)\bar{\lambda}_{2}+\frac{\delta^{2}\left(\theta\,\tilde{\psi}-1\right)g\left(\tilde{\psi}\right)}{1-\delta}=0$$

Solving these simultaneously for $\overline{\lambda}_1$ and $\overline{\lambda}_2$ yields

$$\bar{\lambda}_{1} = \frac{\delta g\left(\tilde{\psi}\right) \left(\theta \,\delta \left(\tilde{\psi}-1\right) \left(\theta \,\tilde{\psi}-1\right) g\left(\tilde{\psi}\right) - (1-\delta) \left(\theta \,\tilde{\psi}-1\right)\right)}{\delta \left(\theta-1\right) g\left(\tilde{\psi}\right) + 1 - \delta}$$
$$\bar{\lambda}_{2} = -\frac{\left(\theta \left(\tilde{\psi}-1\right) \left(\theta \,\tilde{\psi}-1\right) g\left(\tilde{\psi}\right) + \theta \,\tilde{\psi}-1\right) \delta^{2} g\left(\tilde{\psi}\right)}{\delta \left(\theta-1\right) g\left(\tilde{\psi}\right) + 1 - \delta}$$

The following two conditions are together sufficient for both of these expressions to be finite and positive.

$$\delta < \frac{1}{1 + (1 - \theta)g\left(\tilde{\psi}\right)} \tag{67}$$

$$g\left(\tilde{\psi}\right) \le \frac{1}{\left(1 - \tilde{\psi}\right)\theta} \tag{68}$$

Note now that the RHS of (34) is increasing in $\ddot{\psi}_1$ and $\ddot{\psi}_2$, and is the same as the RHS of (61) when $\ddot{\psi}_1 = \hat{\pi}$ and (68) when $\ddot{\psi}_1 = \ddot{\psi}_2 = \tilde{\psi}$. Hence, condition (35) from the main text (repeated below as (69)) is sufficient for (68) to definitely hold, and for (61) and (66) to never hold. Also, condition (69) is sufficient for (67) to hold provided that (51) does not. (Subtracting the LHS of (51) from (67), requiring that this be positive and solving for $g(\tilde{\psi})$ yields $g(\tilde{\psi}) < \frac{1}{\theta - \hat{\pi}}$, which will definitely be fulfilled if (69) is.)

$$\forall_{\pi} : g(\pi) < \frac{1}{\theta \left(1 - \hat{\pi}\right)} \tag{69}$$

Condition (69) is therefore sufficient for the optimal semi-constrained quasi-flat path to be flat. $\hfill \Box$

References

[Abreu, 1986]	ABREU, DILIP (1986). "Extremal Equilibria of Oligopolistic Supergames". Journal of Economic Theory, 39, 191–225.
[Abreu, 1988]	ABREU, DILIP (1988). "On the Theory of Infinitely Repeated Games with Dis- counting". <i>Econometrica</i> , 56(2), 383–396.
[Aumann & Shapley, 1992]	AUMANN, ROBERT J. AND SHAPLEY, LLOYD S. (1992). "Long Term Competition - A Game Theoretic Analysis". UCLA Eco- nomics Working Papers 676, UCLA Department of Economics.
[Benoit & Krishna, 1993]	BENOIT, JEAN-PIERRE AND KRISHNA, VIJAY (1993). "Renegotiation in Finitely Repeated Games". Econometrica, 61(2), 303–23.
[Bernheim & Stark, 1988]	BERNHEIM, DOUGLAS B. AND STARK, ODED (1988). "Altruism within the Family Reconsidered: Do Nice Guys Finish Last?". The American Economic Review, 78(5), 1034–1045.
[Cremer, 1986]	CREMER, JACQUES (1986). "Cooperation in Ongoing Organizations". The Quarterly Journal of Economics, 101(1), 33–49.
[Diamond, 1984]	DIAMOND, PETER (1984). "Money in Search Equilibrium". <i>Econometrica</i> , 52(1), 1–20.
[Diamond & Stiglitz, 1974]	DIAMOND, PETER A. AND STIGLITZ, JOSEPH E. (1974). "Increases in Risk and in Risk Aversion". Journal of Economic Theory, 8(3), 337–360.
[Farrell & Maskin, 1989]	FARRELL, JOSEPH T. AND MASKIN, ERIC S. (1989). "Renegotiation in Repeated Games". Games and Economic Behaviour, 1(1), 327–360.
[Fudenberg & Maskin, 1986]	FUDENBERG, DREW AND MASKIN, ERIC (1986). "The Folk Theorem in Repeated Games with Discounting or with Incomplete Information". <i>Econometrica</i> , 54(3), 533–54.
[Hammond, 1975]	HAMMOND, PETER (1975). Charity: Altru- ism or Cooperative Egoism? In E. S. Phelps (Ed.), Altruism, Morality and Economic Theory

(pp. 115–131). Russell Sage Foundation, New York.

lem". The American Economic Review, 78(4),

[Kandori, 1992]	KANDORI, MICHIHIRO (1992). "Repeated
	Games Played by Overlapping Genera-
	tions of Players". The Review of Economic
	Studies, 59(1), 81-92.
[Kotlikoff et al., 1988]	Kotlikoff, Laurence J. and Persson,
	TORSTEN AND SVENSSON, LARS E. O. (1988).
	"Social Contracts as Assets: A Possible
	Solution to the Time-Consistency Prob-

[Lagrange, 1806] LAGRANGE, JOSEPH-LOUIS (1806). Leçons Sur le Calcul des Fonctions. Chez Courgier, Paris.

662 - 677.

- [Lambson, 1987] LAMBSON, VAL EUGENE (1987). "Optimal Penal Codes in Price-Setting Supergames with Capacity Constraints". Review of Economic Studies, 54(3), 385–97.
- [Messner & Polborn, 2003] MESSNER, MATTHIAS AND POLBORN, MAT-TIAS K. (2003). "Cooperation in Stochastic OLG games". Journal of Economic Theory, 108(1), 152–168.
- [Povey, 2010] POVEY, RICHARD (2010). "Punishment and the Potency of Group Selection". Available at users.ox.ac.uk/~sedm1375.
- [Rawls, 1999] RAWLS, JOHN (1999). A Theory of Justice. Oxford University Press, Oxford.
- [Rees, 1993] REES, RAY (1993). **"Tacit Collusion"**. Oxford Review of Economic Policy, 9(2), 27–40.
- [Roberts, 1980] ROBERTS, KEVIN W S (1980). "Interpersonal Comparability and Social Choice Theory". Review of Economic Studies, 47(2), 421–39.
- [Rothschild & Stiglitz, 1970] ROTHSCHILD, MICHAEL AND STIGLITZ, JOSEPH E. (1970). "Increasing Risk: I. A Definition". Journal of Economic Theory, 2(3), 225 – 243.
- [Rubinstein, 1979] RUBINSTEIN, ARIEL (1979). "Equilibrium in Supergames with the Overtaking Criterion". Journal of Economic Theory, 21(1), 1–9.

[Ruffin, 1972]	RUFFIN, ROY J. (1972). "Pollution in a Crusoe Economy". The Canadian Journal of Economics, 5(1), 110–118.
[Samuelson, 1958]	SAMUELSON, PAUL A. (1958). "An Exact Consumption-Loan Model of Interest with or without the Social Contrivance of Money". Journal of Political Economy, 66, 467.
[Simon & Blume, 1994]	SIMON, CARL P. AND BLUME, LAWRENCE (1994). Mathematics for Economists. W. W. Norton and Company, New York.