# Information Manipulation, Coordination, and Regime Change*

CHRIS EDMOND

*University of Melbourne*

This article presents a model of information manipulation and political regime change. There is a regime that can be overthrown but only if enough citizens participate in an uprising. Citizens are imperfectly informed about the regime's ability to resist an uprising and the regime can engage in propaganda that, taken at face-value, makes the regime seem stronger than it truly is. This coordination game with endogenous information manipulation has a unique equilibrium and the article gives a complete analytic characterization of the equilibrium's comparative statics. Holding fixed the number of signals available to citizens, if the per-unit signal precision is sufficiently high then the regime is harder to overthrow. In contrast, if the number of signals increases, so that both total signal precision and the regime's costs of manipulation rise together, then the regime is easier to overthrow unless there are strong economies of scale in information control.

*Key words*: Global games, Signal-jamming, Hidden actions, Propaganda, Bias, Media

*JEL Codes*: C7, D7, D8

## 1. INTRODUCTION

Will improvements in information technologies help in overthrowing autocratic regimes? Optimists on this issue stress the role of new technologies in facilitating coordination and in improving information about a regime's intentions and vulnerabilities. The "Arab Spring" of uprisings against regimes in Tunisia, Egypt, Libya, and elsewhere that began in December 2010 has led to widespread discussion of the role of modern social media such as Facebook, Twitter, Skype, and YouTube in facilitating regime change. Similar discussion followed the use of such technologies during the demonstrations against the Iranian regime in June 2009 (*e.g.* Musgrove, 2009; Kirkpatrick, 2011).

But optimism about the use of new technologies in putting autocratic regimes under sustained pressure is hardly new; social media is only the latest technology to be viewed as a catalyst for regime change. Simple internet access, cell phones, satellite television, radio, and newspaper

---

have all been viewed as potential catalysts too. And while information optimism has a long and somewhat mixed history, it is also worth bearing in mind that the relationship between information technologies and autocratic regimes has a prominent dark side. Perhaps the most well-known examples are the use of mass media propaganda by Nazi Germany and the Soviet Union. Moreover, it has become increasingly clear that recent breakthroughs in information technology also provide opportunities for autocracies. During the Iranian demonstrations, technologies like Twitter allowed the regime to spread rumors and disinformation (Esfandiari, 2010). Similarly, the Chinese regime's efforts to counter online organization make use of the exact same technologies that optimists hope will help in bringing regime change (Kalathil and Boas, 2003; Fallows, 2008; Morozov, 2011).

So, should we be optimistic that recent breakthroughs in information technology will lead to the collapse of present-day autocratic regimes? To help address this question, I develop a simple model of information and regime change. While stylized, the model provides a number of insights into the ways in which a regime's chances of survival are affected by changes in information technology. In particular, the model implies the degree of *scale economies*, if any, in the regime's control over multiple sources of information plays a crucial role in determining whether an information revolution is likely to bring about regime change. For example, suppose an information revolution consists of technologies such as newsprint, radio, and cinema that are relatively easily *centralized* in the sense that, by investing in a large fixed propaganda establishment, the regime can then exert influence over additional newspapers, radio stations, or cinemas at low marginal cost. Then the regime may have sufficiently strong economies of scale in information control that its chances of survival will increase. In contrast, if the technology is more *decentralized*, so that there are diseconomies of scale in information control, then the model predicts the regime will become easier to overthrow as the number of sources of information increases.

Section 2 outlines the model. There is a single regime and a large number of citizens with heterogeneous information. Citizens can participate in an attack on the regime or not. If enough do, the regime is overthrown. The regime's ability to withstand an attack is given by a single parameter, the regime's type. Citizens are imperfectly informed about the regime's type and may coordinate either on overturning the regime or not. The regime is informed about its type and seeks to induce coordination on the status quo. It does this by taking a costly *hidden action* which influences the distribution of signals so that citizens receive *biased* information that, taken at face-value, suggests the regime is hard to overthrow. Citizens are rational and internalize the regime's incentives when forming their beliefs.

Section 3 gives the article's first main result: this coordination game with endogenous information manipulation has a unique perfect Bayesian equilibrium. Section 4 then turns to the implications of changes in information precision, holding the number of sources of information fixed. The second main result of the article is that the regime's information manipulation increases the regime's ex ante chances of surviving when the intrinsic *per-unit* signal precision is sufficiently high. If this per-unit precision is sufficiently high, then the regime survives in all circumstances where it is possible for the regime to survive. Section 5 then allows for an increase in the *number* of signals, which simultaneously (i) increases the total signal precision and (ii) increases the regime's costs of information manipulation. These two effects have offsetting implications for the regime's chances of surviving. The third main result of the article is that the net effect of an increase in the number of signals is to reduce the regime's chances of surviving—unless there are strong *economies of scale* in information control. This section also discusses the model's predictions in light of the historical relationship between autocracies and information control. Section 6 considers two extensions. Section 7 concludes. All proofs are in the Appendix. A supplementary online appendix contains further extensions and discussion.

## 1.1. *Effective information manipulation*

Perhaps the most interesting feature of the model is that the regime's information manipulation—a form of *signal-jamming*—can be effective in equilibrium. The regime is informed and takes a type-specific hidden action that cannot be directly observed by the individual citizens receiving the information. While the citizens know, in equilibrium, the *function* mapping the regime's type into its hidden action, they must infer the specific value of the action jointly with making inferences about the type itself. This distinguishes the setup from standard *career concerns* models, such as Holmström (1999), where the information sender is uninformed and chooses a single action to maximize its expected payoff, thereby making it easy for receivers to deduce any bias in their signals and respond accordingly. In my model, the extent to which signal-jamming is effective depends on the signal precision. If signals are precise, clustered tightly about the mean, then it takes only a small amount of bias to deliver a large shift in aggregate behaviour, and, consequently, signal-jamming can be especially effective.

Moreover, the signal-jamming effect is *amplified* by the fact that information receivers are playing a coordination game. It is common knowledge that the regime has the ability to manipulate information and that *any* equilibrium manipulation tends to reduce the mass of citizens who participate in an attack. Since individual actions are strategic complements, at the margin this reduces *every* individual's incentives to participate in an attack. In short, the strategic complementarities amplify the effectiveness of any bias the regime is able to impart. These coordination aspects of the model also distinguish the setup from standard models of strategic information transmission, such as Crawford and Sobel (1982), that focus on a single sender and single receiver. As in such models, different types of senders are able to partially pool in equilibrium. But whereas in Crawford and Sobel such pooling takes the form of different sender types choosing *the same* observable (and payoff irrelevant) message, here it takes the form of different regimes choosing different but unobservable actions in a way that allows lower types to mimic the strength of higher types and thereby reduce the size of the attack against them.

## 1.2. *Uniqueness*

The model is also closely related to Angeletos *et al.* (2006). Both my model and theirs start with a coordination game with imperfect information—often referred to as a *global game*[1]—and to that add an informed policy-maker whose actions make information endogenous. Angeletos *et al.* show that an endogenous information structure can lead to equilibrium multiplicity even if the underlying global game has a unique equilibrium. My model features a unique equilibrium despite its endogenous information structure and in that sense is more reminiscent of standard global games. The key reason for the uniqueness in my model relative to theirs is that in my setting the endogenous signals are a *monotone* function of the regime's type whereas in Angeletos *et al.* the endogenous signals are a non-monotone function that gives information receivers more common belief as to whether the status quo is likely to prevail or not. The greater common belief in their setting makes it easier to sustain multiple equilibria.

In short, this article provides a framework with a form of costly *noisy signalling* that, in contrast with the indeterminacies of many signalling models, instead features determinate outcomes of the kind familiar from standard global games.

---

1. See Carlsson and van Damme (1993) and Morris and Shin (1998, 2003). Other political economy applications of global games include Boix and Svolik (2013), Bueno de Mesquita (2010), Chassang and Padro-i-Miquel (2010) and Shadmehr and Bernhardt (2011).

## 2. MODEL OF INFORMATION MANIPULATION AND REGIME CHANGE

There is a unit mass of ex ante identical citizens, indexed by $i \in [0,1]$. The citizens face a regime that seeks to preserve the status quo. Each citizen decides whether to participate in an attack on the regime, $s_i = 1$, or not, $s_i = 0$. The aggregate attack is $S := \int_0^1 s_i di$. The type of a regime $\theta$ is its private information. The regime is overthrown if and only if $\theta < S$.

### 2.1.  *Citizen payoffs*

The payoff to a citizen is given by

$$u(s_i, S, \theta) = (\mathbb{1}\{\theta < S\} - p)s_i \tag{1}$$

where $\mathbb{1}\{\cdot\}$ denotes the indicator function and where $p > 0$ is the opportunity cost of participating in an attack. A citizen will participate, $s_i = 1$, whenever they expect regime change to occur with more than probability $p$. To focus on scenarios where individual decisions are not trivial, I also assume that $p < 1$. Otherwise, if $p \geq 1$, it is optimal for an individual to not participate independent of $\theta$ and $S$. With $p < 1$ the individual $s_i$ and aggregate $S$ are *strategic complements*: the more citizens participate in an attack, the more likely it is that the regime is overthrown and so the more likely it is that any individual also participates.

### 2.2.  *Citizen information and media outlets*

Citizens begin with common priors for $\theta$, specifically the improper uniform on $\mathbb{R}$. Citizens obtain information about the regime from $n \geq 1$ identical *media outlets*. Each outlet $j = 1, ..., n$ releases a report $y$ and each citizen observes that report with idiosyncratic noise. In particular, citizen $i$ receives $n$ signals of the form $x_{ij} = y + \varepsilon_{ij}$ where the noise $\varepsilon_{ij}$ is jointly IID normal across citizens and media outlets with mean zero and precision $\hat{\alpha} > 0$. Crucially, media outlets report the regime's *apparent strength* $y = \theta + a$, the sum of the regime's true type $\theta$ plus an action $a \geq 0$ chosen by the regime in an effort to convey a biased impression of how difficult it will be to overthrow.[2] Citizens cannot directly observe this action $a$ and must infer it jointly with inferring $\theta$ itself.

The *average signal* of citizen $i$ is $x_i := \frac{1}{n}\sum_{j=1}^n x_{ij}$, which is normal with mean $y = \theta + a$ and precision $\alpha := n\hat{\alpha}$, proportional to $n$. Given the regime's apparent strength $y$, the density of average signals is

$$f(x_i | y) := \sqrt{\alpha}\phi(\sqrt{\alpha}(x_i - y)), \qquad y = \theta + a \tag{2}$$

where $\phi(\cdot)$ denotes the standard normal PDF. Since the average signal is a sufficient statistic, I simply refer to this as the citizen's signal for short.

### 2.3.  *Regime's costs of information manipulation*

The regime's choice of hidden action $a$ is made after observing $\theta$, *i.e.* the regime is *informed*. The action $a$ enters the citizens' $n$ signals symmetrically with cost $C(a, n)$ to the regime. If $a = 0$, the regime pays no cost, $C(0, n) = 0$ for all $n$. But for $a > 0$, the cost is strictly increasing in both arguments, $C_a > 0$ and $C_n > 0$, and convex in the hidden action, $C_{aa} \geq 0$. An important special case is $C(a, n) = c(n)a$, *i.e.* where the marginal cost of taking action $a$ is a constant $c(n) > 0$, with $c(n)$ increasing in $n$. To simplify notation, I suppress the exogenous parameter $n$ and simply write the cost function $C(a)$ whenever no confusion ought to arise.

---

2.  The supplementary appendix provides simple micro-foundations for this form of media report and also considers *heterogeneous* media outlets, some of which pass on the manipulation and some of which do not.

### 2.4. Regime's payoffs

The gross benefit to the regime $B(\theta, S)$ depends on whether it survives or not. If $\theta < S$ the regime is overthrown and obtains an outside option normalized to zero. Otherwise, if $\theta \geq S$, the regime obtains a gross benefit $\theta - S$ from remaining in power. The key assumption here is that the benefit is *separable* in $\theta$ and $S$. Given separability, the *linearity* in $\theta$ and $S$ is without further loss of generality and simply represents a normalization of the type. The benefit is strictly decreasing in $S$, *i.e.* the regime wants to minimize the costs of unrest and so wants $S$ small even when it survives.

The net payoff to a regime can therefore be written

$$B(\theta, S) - C(a) = \max[0, \theta - S] - C(a). \tag{3}$$

Observe that if $\theta < S$, the regime will choose $a = 0$ and thereby obtain a net payoff of zero.

### 2.5. Equilibrium

A symmetric *perfect Bayesian equilibrium* of this model consists of individual beliefs $\pi(\theta | x_i)$ and participation decisions $s(x_i)$, an aggregate attack $S(y)$, and regime actions $a(\theta)$ such that: (i) a citizen with information $x_i$ rationally takes into account the manipulation $y(\theta) = \theta + a(\theta)$ when forming their posterior beliefs, (ii) given these beliefs, $s(x_i)$ maximizes individual expected utility, (iii) the aggregate attack is consistent with the individual decisions, and (iv) the actions $a(\theta)$ maximize the regime's payoff given the aggregate attack. In equilibrium, the regime is overthrown if $\theta < S(y(\theta))$ and otherwise survives.

### 2.6. Standard global game benchmark

If there are no hidden actions, the analysis reduces to that of a *standard global game*. In particular, if $a(\theta) = 0$ for all $\theta$, then each citizen has exogenous signal $x_i = \theta + \varepsilon_i$ and there is a unique equilibrium, as in Carlsson and van Damme (1993) and Morris and Shin (1998). In this equilibrium, strategies are threshold rules: there is a unique $\theta^*_{MS}$ such that the regime is overthrown for $\theta < \theta^*_{MS}$ and a unique signal $x^*_{MS}$ such that a citizen participates for $x_i < x^*_{MS}$. These thresholds are given by:

---

MORRIS–SHIN BENCHMARK. The equilibrium thresholds $x^*_{MS}, \theta^*_{MS}$ simultaneously solve

$$\text{Prob}[\theta < \theta^*_{MS} | x^*_{MS}] = \Phi(\sqrt{\alpha}(\theta^*_{MS} - x^*_{MS})) = p \tag{4}$$

$$\text{Prob}[x_i < x^*_{MS} | \theta^*_{MS}] = \Phi(\sqrt{\alpha}(x^*_{MS} - \theta^*_{MS})) = \theta^*_{MS} \tag{5}$$

where $\Phi(\cdot)$ denotes the standard normal CDF and $p$ is the individual opportunity cost of attacking. In particular, $\theta^*_{MS} = 1 - p$, independent of $\alpha$, and $x^*_{MS} = 1 - p - \Phi^{-1}(p)/\sqrt{\alpha}$.

---

The first condition says that if the regime's threshold is $\theta^*_{MS}$, the marginal citizen with signal $x_i = x^*_{MS}$ expects regime change with exactly probability $p$. The second says that if the signal threshold is $x^*_{MS}$, a regime with $\theta = \theta^*_{MS}$ will be indifferent to abandoning its position.

## 3. UNIQUE EQUILIBRIUM WITH INFORMATION MANIPULATION

The first main result of this article is that there is a unique equilibrium.

**Proposition 1.** *There is a unique perfect Bayesian equilibrium. The equilibrium is monotone in the sense that there exist thresholds $x^*$ and $\theta^*$ such that $s(x_i) = 1$ for $x_i < x^*$ and zero otherwise, while the regime is overthrown for $\theta < \theta^*$ and not otherwise.*

A detailed proof is given in Appendix A. Briefly, the proof involves first showing (i) that there is a unique equilibrium in monotone strategies, and (ii) that the unique monotone equilibrium is the only equilibrium which survives the iterative elimination of interim strictly dominated strategies. Here in the main text I briefly characterize the equilibrium.

### 3.1. *Equilibrium characterization*

Let $\hat{x}$ denote a candidate for the citizens' threshold and let $\Theta(\hat{x})$ and $a(\theta, \hat{x})$ denote candidates for the regime's threshold and hidden actions given $\hat{x}$.

**Regime problem.** Since citizens participate $s(x_i) = 1$ for $x_i < \hat{x}$, for any given $\hat{x}$ the aggregate attack facing a regime of apparent strength $y$ is simply

$$S(y) = \Phi(\sqrt{\alpha}(\hat{x} - y)), \qquad y = \theta + a \tag{6}$$

(using the form of the signal density given in (2) above). Since the regime is overthrown for $\theta < \Theta(\hat{x})$, hidden actions are $a(\theta, \hat{x}) = 0$ for all $\theta < \Theta(\hat{x})$, otherwise the regime would be incurring a cost but receiving no benefit. For all $\theta \geq \Theta(\hat{x})$, the regime chooses

$$a(\theta, \hat{x}) \in \underset{a \geq 0}{\operatorname{argmin}} \left[ \Phi(\sqrt{\alpha}(\hat{x} - \theta - a)) + C(a) \right]. \tag{7}$$

A key step in proving equilibrium uniqueness is to recognize that hidden actions are given by $a(\theta, \hat{x}) = A(\theta - \hat{x})$, where the auxiliary function $A : \mathbb{R} \to \mathbb{R}_+$ is *exogenous* and in particular does *not* depend on the citizen threshold $\hat{x}$. This function is defined by

$$A(t) := \underset{a \geq 0}{\operatorname{argmin}} \left[ \Phi(\sqrt{\alpha}(-t - a)) + C(a) \right]. \tag{8}$$

The regime threshold $\Theta(\hat{x})$ is then found from the indifference condition

$$\Theta(\hat{x}) = \Phi(\sqrt{\alpha}(\hat{x} - \Theta(\hat{x}) - A(\Theta(\hat{x}) - \hat{x}))) + C(A(\Theta(\hat{x}) - \hat{x})). \tag{9}$$

This condition requires that total costs equal total benefits at the extensive margin. For any given candidate citizen threshold $\hat{x}$, equations (8)–(9) determine the regime threshold $\Theta(\hat{x})$ and hidden actions $a(\theta, \hat{x}) = A(\theta - \hat{x})$ solving the regime's problem.

**Citizen problem.** An individual participates only if they believe the regime will be overthrown with probability at least $p$. Given $\hat{x}$ and the solution to the regime's problem, the posterior probability of regime change for an individual with arbitrary signal $x_i$ is

$$\operatorname{Prob}[\theta < \Theta(\hat{x}) \mid x_i, a(\cdot, \hat{x})] := \frac{\int_{-\infty}^{\Theta(\hat{x})} \sqrt{\alpha} \phi(\sqrt{\alpha}(x_i - \theta)) d\theta}{\int_{-\infty}^{\infty} \sqrt{\alpha} \phi(\sqrt{\alpha}(x_i - \theta' - a(\theta', \hat{x}))) d\theta'}$$

(using $a(\theta, \hat{x}) = 0$ for all $\theta < \Theta(\hat{x})$ in the numerator). Writing the hidden actions in terms of the auxiliary function $a(\theta, \hat{x}) = A(\theta - \hat{x})$, evaluating at $x_i = \hat{x}$, and then equating the result to the

opportunity cost $p$ gives the indifference condition characterizing the citizen threshold

$$\frac{\int_{-\infty}^{\Theta(\hat{x})} \sqrt{\alpha}\phi(\sqrt{\alpha}(\hat{x}-\theta))d\theta}{\int_{-\infty}^{\infty} \sqrt{\alpha}\phi(\sqrt{\alpha}(\hat{x}-\theta'-A(\theta'-\hat{x})))d\theta'} = p. \tag{10}$$

**Monotone equilibrium.**    As shown in Appendix A, there is a unique monotone equilibrium with thresholds $x^*$ and $\theta^*$ simultaneously solving conditions (9) and (10). The regime's equilibrium hidden actions are then given by $a(\theta)=A(\theta-x^*)$ using the auxiliary function from (8). A key step in the proof is showing that the posterior probability on the left-hand side of equation (10) depends only on the difference $\Theta(\hat{x})-\hat{x}$ and is monotone increasing in this argument so that (10) can be solved for a unique difference $\theta^*-x^*$. Similarly, the right-hand side of (9) only depends on the difference $\Theta(\hat{x})-\hat{x}$ so we can take the unique solution $\theta^*-x^*$ from (10) and plug it into the right-hand side of (9) to determine $\theta^*$ separately.

The Appendix goes on to show that this monotone equilibrium is the *only* equilibrium.

## 4. EQUILIBRIUM INFORMATION MANIPULATION

The most interesting implication of this model is that the regime's information manipulation—or *signal-jamming*—can be an effective tool for increasing the regime's ex ante chances of surviving. Section 4.1 begins the analysis by showing that the regime's policy $y(\theta)$ is monotone in $\theta$, a general result that obtains independent of precise details of the equilibrium. Section 4.2 then specializes to the case of a linear cost function. This permits a simple characterization of the entire equilibrium, not just the signal-jamming policy. Section 4.3 then shows that, for a given number of signals $n$, the regime's signal-jamming will be effective in equilibrium if the per-unit signal precision $\hat{\alpha}$ is sufficiently high. Section 4.4 provides further intuition and explains how these results relate to standard models in the literature.

### 4.1.  *Signal-jamming policy*

To understand the regime's signal-jamming problem intuitively, it is helpful to recast the problem in slightly more general terms. In particular, let $B(\theta,S)$ denote the regime's gross benefit from surviving and let $S(y)$ denote the aggregate attack facing a regime that chooses apparent strength $y \geq \theta$. For the moment, $S(y)$ can be *anything*—no equilibrium properties need to be imposed at all. In these general terms, the regime's policy is

$$y(\theta) \in \underset{y \geq \theta}{\operatorname{argmax}} \, [B(\theta,S(y))-C(y-\theta)] \tag{11}$$

with $B(\theta,S)=0$ whenever the regime is overthrown. We then have:

**Proposition 2.**    *If $B(\theta,S)$ satisfies the Spence-Mirrlees sorting condition $B_{\theta S} \leq 0$, then for any aggregate attack $S(y)$ the solution $y(\theta)$ to (11) is increasing in $\theta$.*

The proof is given in Appendix A. The key steps are to observe (i) that for given $\theta$, at any interior solution higher values of $y$ cost more and so will be chosen only if they reduce $S(y)$, and (ii) since $B(\theta,S)$ satisfies a sorting condition in $\theta$ and $S$, stronger regimes are at least weakly better

off from a smaller attack than are weaker regimes,[3] hence stronger regimes are more willing, at the margin, to pay the cost of choosing a larger $y$.

This monotonicity result is intuitive, but, since the form of $S(y)$ is left open, we do not yet know if this policy will be effective in equilibrium.

## 4.2. *Equilibrium with linear costs*

To study the equilibrium problem, I first specialize to the case of linear costs, $C(a) := ca$ for some constant $c > 0$. This gives rise to a simple form of signal-jamming and permits an exact characterization of the conditions under which the signal-jamming is effective.

**Overview of the regime's problem.** With linear costs, if the choice of $y$ is to induce survival it must minimize $S(y) + cy$ subject to the constraint $y \geq \theta$. Since the first-order condition for interior solutions is $-S'(y) = c$, all regimes with interior solutions choose *the same* apparent strength $y^*$. Corner solutions can arise for two reasons: (i) for all $\theta$ such that $\theta < S(y^*) + c(y^* - \theta)$ the regime is too weak to sustain the cost of taking $y^*$ and hence is overthrown, and (ii) for all $\theta \geq y^*$ the innate strength of the regime is sufficiently high that they do not need to undertake any costly manipulation (*i.e.* the non-negativity constraint $y \geq \theta$ is binding). All regimes at a corner simply have apparent strength equal to actual strength, $y(\theta) = \theta$. All regimes that are at an interior solution have $y(\theta) = y^*$. In short,

$$y(\theta) = \begin{cases} \theta & \text{if } \theta \notin [\theta^*, \theta^{**}) \\ y^* & \text{if } \theta \in [\theta^*, \theta^{**}) \end{cases}. \tag{12}$$

The cutoff $\theta^*$ is just the threshold above which the regime survives and is implicitly determined by the indifference condition $\theta^* = S(y^*) + c(y^* - \theta^*)$. The higher cutoff $\theta^{**}$ is simply equal to the apparent strength $y^*$ since it is the smallest regime type for which the non-negativity constraint $y^* \geq \theta$ is binding. Figure 1 illustrates. Observe that the lower cutoff $\theta^*$ necessarily lies in $[0, 1]$ since no regime with $\theta < 0$ can survive and any regime with $\theta > 1$ must survive. The upper cutoff $\theta^{**}$ may be greater than 1 though. Even when the regime survives, it prefers $S$ to be small (to reduce the costs associated with insurrections).

To summarize, regimes that are sufficiently weak and that will be overthrown, $\theta < \theta^*$, do not engage in any manipulation. Regimes that are sufficiently strong, $\theta \geq \theta^{**}$, have enough innate ability to resist an attack that they do not engage in any costly manipulation to shore up their position. Only intermediate regimes $\theta \in [\theta^*, \theta^{**})$ find it worthwhile to pay the costs of manipulation in equilibrium and all these intermediate regimes mimic the strength of a strong regime, $\theta^{**}$.

**No manipulation, Morris–Shin outcome.** Intuitively, when the cost of manipulation is sufficiently high, no regime will find it worthwhile to manipulate and the unique equilibrium will coincide with the Morris–Shin outcome. Since the aggregate attack is $S(y) = \Phi(\sqrt{\alpha}(x^* - y))$, the first-order necessary condition $-S'(y) = c$ can be written

$$\sqrt{\alpha}\phi(\sqrt{\alpha}(x^* - y^*)) = c. \tag{13}$$

---

3. I thank an anonymous referee for suggesting this line of argument. Observe also that since the regime's payoff $B(\theta, S) - C(a)$ is separable in $a$, the condition $B_{\theta S} \leq 0$ implies that the regime's indifference curves in $(S, a)$ satisfy a standard *single-crossing* property.
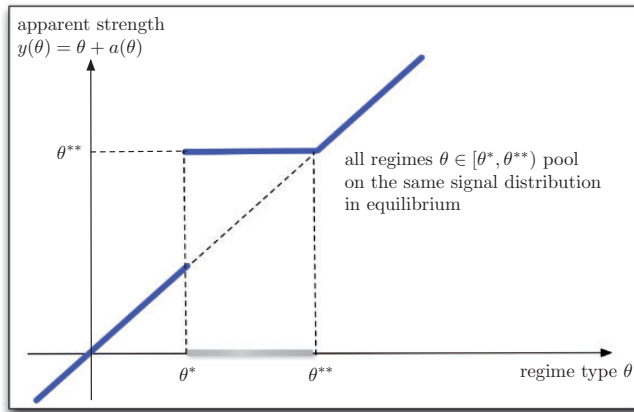
FIGURE 1

Signal-jamming with linear costs

*Notes*: Apparent strength $y(\theta) = \theta + a(\theta)$ when the regime has linear costs of manipulation. All regimes with $\theta < \theta^*$ are overthrown. All regimes with $\theta \in [\theta^*, \theta^{**})$ choose the same apparent strength $y^* = \theta^{**}$ and thus generate the same signal distribution in equilibrium. They mimic a stronger regime $\theta^{**}$ and generate signals for the citizens that are (locally) uninformative about $\theta$.

Observe that $c/\sqrt{\alpha}$ must be sufficiently small for this condition to have a solution.[4] In particular, a necessary condition for the existence of an interior solution is that $c/\sqrt{\alpha} < \phi(0)$ where $\phi(0) = 1/\sqrt{2\pi} \approx 0.399$ is the maximum value of the standard normal density. If $c/\sqrt{\alpha} \geq \phi(0)$, then the signal precision $\alpha$ is too low relative to the cost of manipulation and *all regimes* are at a corner with $y(\theta) = \theta$. In this case, $\theta^* = \theta^*_{MS} = 1 - p$ and $\theta^{**} = \theta^*$. In contrast, if the signal precision is high, *some* regimes will engage in information manipulation.

**Active manipulation.** In particular, let $\underline{\alpha}(c) := (c/\phi(0))^2$ denote the smallest precision such that (13) has a solution. Then all regimes that manipulate choose apparent strength

$$y^* = x^* + \gamma \tag{14}$$

where the parameter $\gamma$ solves $\sqrt{\alpha}\phi(\sqrt{\alpha}\gamma) = c$, that is

$$\gamma = \sqrt{\frac{1}{\alpha}\log\left(\frac{\alpha}{\underline{\alpha}(c)}\right)} > 0, \qquad \alpha > \underline{\alpha}(c). \tag{15}$$

The signal-jamming is *acute*. All regimes that manipulate information pool on the same distribution of signals, *i.e.* all regimes $\theta \in [\theta^*, \theta^{**})$ generate a mean of $y^* = \theta^{**}$ with signals $x_i = y^* + \varepsilon_i$ that are *locally completely uninformative* about $\theta$.

**Size of the attack.** Now let $S(y(\theta))$ denote the size of the aggregate attack facing a regime of type $\theta$ when the signal-jamming has this form. Using (12) and (14), the attack is

$$S(y(\theta)) = \begin{cases} \Phi(\sqrt{\alpha}(x^* - \theta)) & \text{if } \theta \notin [\theta^*, \theta^{**}) \\ \Phi(-\sqrt{\alpha}\gamma) & \text{if } \theta \in [\theta^*, \theta^{**}) \end{cases} . \tag{16}$$

---

4. This first-order condition may have zero or two solutions. If there are two solutions, it is straightforward to show that only the larger of these satisfies the second-order condition.
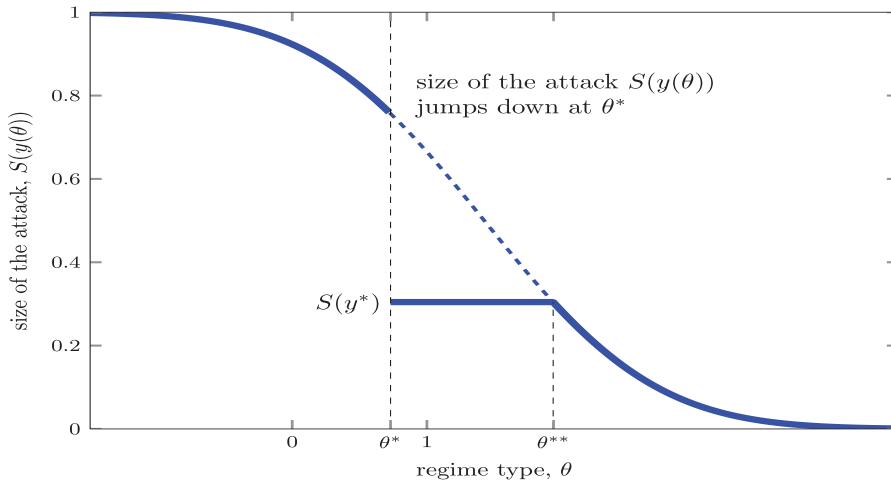
FIGURE 2

Manipulation leads to discrete fall in size of the attack

*Notes*: All regimes with $\theta \in [\theta^*, \theta^{**})$ choose the same apparent strength $y^* = \theta^{**} = x^* + \gamma$ and consequently face the same sized attack $S(y^*) = \Phi(-\sqrt{\alpha}\gamma)$. All regimes that do not manipulate have $y(\theta) = \theta$ and face the attack $S(y(\theta)) = \Phi(\sqrt{\alpha}(x^* - \theta))$. Observe that the attack is continuous at the upper boundary $\theta^{**}$.

Manipulation to convey a higher apparent strength causes the size of the attack to drop discontinuously at the regime threshold $\theta^*$. As shown in Figure 2, at $\theta^*$ the size of the attack jumps discretely from $\Phi(\sqrt{\alpha}(x^* - \theta^*))$ down to the lower value $\Phi(\sqrt{\alpha}(x^* - y^*)) = \Phi(-\sqrt{\alpha}\gamma)$. All regimes that manipulate face *the same-sized attack*, $\Phi(-\sqrt{\alpha}\gamma)$.

**Is manipulation effective in equilibrium?**    Thus, as in classic signalling games, the regime is able to send a (noisy) signal in equilibrium and this enables some weaker regime types to pool with stronger regime types. But, at the risk of repetition, this pooling tells us nothing about whether the signal-jamming is effective in equilibrium. It may be that the only regimes that are able to imitate stronger regime types are those regimes that would have survived in the absence of the signal-jamming technology. Moreover, it may be that some weak regimes that would survive if they could commit to not use manipulation are overthrown because they cannot make that commitment. Put differently, the discussion so far has taken the thresholds $\theta^*$ and $\theta^{**}$ as given. But these thresholds are endogenous and, in principle, may shift strongly against the regime so that in equilibrium information manipulation is ultimately ineffective.

It turns out that manipulation *is* effective when the per-unit signal precision is sufficiently high. But to see this we have to solve the rest of the model.

**Solving the model with linear costs.**    To complete the solution, recall that an individual citizen participates if they believe the regime will be overthrown with posterior probability at least $p$, that is, if $\text{Prob}[\theta < \theta^* | x_i, y(\cdot)] \geq p$. For the marginal citizen with signal $x_i = x^*$ this condition holds with equality so that

$$\text{Prob}[\theta < \theta^* | x^*, y(\cdot)] = \frac{\int_{-\infty}^{\theta^*} \sqrt{\alpha}\phi(\sqrt{\alpha}(x^* - \theta))d\theta}{\int_{-\infty}^{\infty} \sqrt{\alpha}\phi(\sqrt{\alpha}(x^* - y(\theta)))d\theta} = p$$

(since $y(\theta)=\theta$ for $\theta<\theta^*$). After collecting terms

$$\Phi[\sqrt{\alpha}(\theta^*-x^*)]=\frac{p}{1-p}\int_{\theta^*}^{\infty}\sqrt{\alpha}\phi(\sqrt{\alpha}(x^*-y(\theta)))d\theta. \qquad (17)$$

For the special case of linear costs, the integral on the right-hand side can be calculated using the first-order condition (13), giving

$$\int_{\theta^*}^{\infty}\sqrt{\alpha}\phi(\sqrt{\alpha}(x^*-y(\theta)))d\theta=(x^*-\theta^*+\gamma)c+\Phi(-\sqrt{\alpha}\gamma) \qquad (18)$$

(since $y(\theta)=\theta$ for $\theta\geq\theta^{**}$ and $\theta^{**}=y^*=x^*+\gamma$). Plugging this back into (17) gives

$$\Phi[\sqrt{\alpha}(\theta^*-x^*)]=\frac{p}{1-p}\big[(x^*-\theta^*+\gamma)c+\Phi(-\sqrt{\alpha}\gamma)\big]. \qquad (19)$$

It is straightforward to show that there is a unique threshold difference $\theta^*-x^*$ that solves this equation. The regime threshold $\theta^*$ is then determined using the indifference condition (9) which, with linear costs, can likewise be written

$$\theta^*=(x^*-\theta^*+\gamma)c+\Phi(-\sqrt{\alpha}\gamma) \qquad (20)$$

where the difference $\theta^*-x^*$ on the right-hand side is uniquely determined by (19) above.

    With the equilibrium thresholds $x^*,\theta^*$ determined by (19)–(20), we can now derive the conditions under which information manipulation is effective for the regime.

### 4.3. *Effectiveness of information manipulation*

I measure the effectiveness of signal-jamming by its ability to reduce the regime's threshold $\theta^*$ below the Morris–Shin benchmark $\theta^*_{\text{MS}}=1-p$. A lower $\theta^*$ increases the regime's ex ante survival probability by making it more likely that nature draws a $\theta\geq\theta^*$. In principle, it might be the case that lower $\theta^*$ is achieved through large, costly, actions that give the regime a lower net payoff than they would achieve in the Morris–Shin world. But it turns out that, for sufficiently high signal precision, the fall in $\theta^*$ also represents a genuine increase in payoffs.

    In this section, I focus on results *holding fixed* the number of signals $n$ (and hence holding fixed the regime's cost function) so that an increase in the *per-unit* precision $\hat{\alpha}$ is identical to an increase in the total precision $\alpha=n\hat{\alpha}$. Since these results hold $n$ fixed but allow for any marginal cost $c>0$ there is no loss of generality in setting $n=1$ so that $\alpha=\hat{\alpha}$. In Section 5, I return to the case of a change in the number of signals $n$ that simultaneously changes both the total precision of information *and* the regime's cost function.

        **Effective manipulation when signal precision is high.**    The main result here is:

**Proposition 3.** *For signal precision $\alpha$ sufficiently high, the regime threshold $\theta^*$ is strictly less than the Morris–Shin benchmark $\theta^*_{\text{MS}}=1-p$. In particular, for each c there is an $\underline{\alpha}(c)$ such that: (i) for $\alpha\leq\underline{\alpha}(c)$ all regimes are at a corner with $y(\theta)=\theta$ for all $\theta$ and $\theta^*=\theta^*_{\text{MS}}$; otherwise (ii) for $\alpha>\underline{\alpha}(c)$, regimes $\theta\in[\theta^*,\theta^{**})$ are at an interior solution and there is a critical precision $\alpha^*(c,p)\geq\underline{\alpha}(c)$ given by*

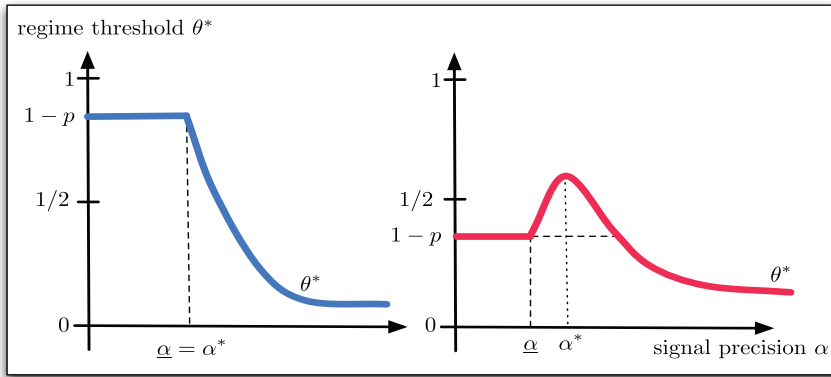$$\alpha^*(c,p):=\underline{\alpha}(c)\exp\left(\max\left[0,\Phi^{-1}(p)\right]^2\right) \qquad (21)$$

Information manipulation is effective when signal precision $\alpha$ is sufficiently high

*Notes*: The left panel shows the case when $p < 1/2$ and the regime threshold $\theta^*$ is *monotone* decreasing in the precision $\alpha$. The right panel shows the case when $p > 1/2$ so that $\theta^*$ is non-monotone in $\alpha$. In both cases, for low enough precision ($\alpha < \underline{\alpha}$) the regime threshold coincides with the Morris–Shin benchmark $\theta^*_{\mathrm{MS}} = 1 - p$ while for high enough $\alpha$ the threshold is less than $1 - p$.

*such that*

$$\frac{\partial}{\partial \alpha} \theta^* < 0 \qquad for\ all \qquad \alpha > \alpha^*(c, p) \tag{22}$$

*with* $\lim_{\alpha \to \infty} \theta^* = 0$, *less than the Morris–Shin benchmark.*

There are two cases to consider. First, when the opportunity cost of attacking the regime is *low*, $p < 1/2$, then all that matters is whether the signal precision $\alpha$ is high enough to induce *any* regime to manipulate. If so, the regime threshold is monotonically declining in $\alpha$ and so for all $\alpha > \underline{\alpha}$ the regime threshold is less than the Morris–Shin benchmark. This is shown in the left panel of Figure 3. Second, if the opportunity cost of attacking the regime is *high*, $p > 1/2$, then the regime threshold is "single-peaked" in $\alpha$, reaching a maximum at $\alpha^* > \underline{\alpha}$ before declining thereafter. This is shown in the right panel of Figure 3. In either case, in the limit as $\alpha \to \infty$, the regime threshold $\theta^* \to 0$ for any fixed $c$ and $p$.

**Even the most fragile regimes can survive.**    This result is striking. As the precision becomes sufficiently high, *all* the regimes that can survive, do survive. To see an extreme example of this, consider an economy with opportunity cost $p \to 0$ so that it requires almost no individual sacrifice to participate in an attack on the regime. In the Morris–Shin benchmark, we would have $\theta^*_{\mathrm{MS}} \to 1$ and only the strongest regimes, those with $\theta \geq 1$, can survive. But with information manipulation, we have $\theta^* \to 0$ so all regimes $\theta \geq 0$ survive even though $p$ is very low. If information can be manipulated and signals are sufficiently precise, then even the very most fragile regimes can survive.[5]

**Signal precision and the size of the marginal attack.**    To understand this result further, observe that the size of the attack facing the marginal regime, $S(y^*)$, is always monotonically

---

5. Perhaps surprisingly, this result does not depend on a diffuse prior for $\theta$. The supplementary appendix shows how this limiting result still obtains even in the case of informative priors.

declining in the signal precision $\alpha$. This follows immediately on plugging the expression for $\gamma$ from (15) into (16). But there is another, perhaps more insightful, way to see this too. Recall the problem of a regime minimizing $S(y)+cy$ with $S(y)=\Phi(\sqrt{\alpha}(x^*-y))$. Since $S(y)$ is strictly monotone in $y$ we can equivalently choose $S$ to minimize

$$S-\frac{c}{\sqrt{\alpha}}\Phi^{-1}(S) \tag{23}$$

with first-order condition $\phi(\Phi^{-1}(S))=c/\sqrt{\alpha}$. Any solution is independent of the citizen threshold $x^*$ and depends only on the ratio $c/\sqrt{\alpha}$ (rather than each parameter separately). The ratio $c/\sqrt{\alpha}$ determines the *effective cost* of obtaining a given-sized attack. And, from the second-order condition, the size of the attack that solves this cost minimization problem must be at least locally increasing in $c/\sqrt{\alpha}$. From this point of view, we see clearly that an increase in the signal precision $\alpha$ reduces the effective cost of obtaining a smaller attack; the more probability density is located near the mean, the easier it is for the regime to achieve a relatively large change in the size of the attack from a relatively small perturbation to the signal mean. Other things equal, a regime can achieve a smaller attack when $\alpha$ is large.

Whenever $\alpha > \underline{\alpha}(c)$, we can solve (23) to obtain the marginal attack

$$S(y^*)=\Phi(-\sqrt{\alpha}\gamma)=\Phi\left(-\sqrt{\log\left(\frac{\alpha}{\underline{\alpha}(c)}\right)}\right)=\Phi\left(-\sqrt{2\log\left(\frac{\sqrt{\alpha}}{c}\phi(0)\right)}\right). \tag{24}$$

Otherwise, for $\alpha < \underline{\alpha}(c)$, the marginal attack is simply the Morris–Shin benchmark, $1-p$.

This discussion highlights the special role of the effective cost parameter $c/\sqrt{\alpha}$. When we return to the case of a variable number of signals $n$ in Section 5 below, this will have the form $c(n)/\sqrt{n\hat{\alpha}}$. Depending on the curvature of $c(n)$, an increase in the number of signals may either increase, decrease, or leave unchanged the effective cost of manipulation.

**Critical signal precision.**    If the regime did not have to pay to obtain $y^*$, then we would be done; the reduction in the size of the marginal attack would directly translate into a reduction in $\theta^*$. But since the manipulation is costly, the regime threshold satisfies $\theta^*=S(y^*)+c(y^*-\theta^*)$ so that a fall in $S(y^*)$ is not enough to guarantee that $\theta^*$ falls as the signal precision $\alpha$ increases. In particular, if the cost of $y^*$ increases sufficiently faster than the size of the attack $S(y^*)$ falls, then the threshold $\theta^*$ increases. When this happens, there is a non-monotone relationship, as in the right panel of Figure 3. Proposition 3 establishes that this non-monotonicity in $\alpha$ can only occur if the individual opportunity cost of attacking the regime is large ($p > 1/2$). The critical precision $\alpha^*$ is exactly the value of $\alpha$ such that the marginal regime faces an attack equal to the Morris–Shin benchmark $1-p$. Using (24), the critical precision $\alpha^*$ solves

$$\Phi\left(-\sqrt{2\log\left(\frac{\sqrt{\alpha^*}}{c}\phi(0)\right)}\right)=1-p \tag{25}$$

which can be inverted to obtain the expression in (21) above. When the signal precision is larger than $\alpha^*$ the regime threshold $\theta^*$ is decreasing in $\alpha$ and is eventually less than $1-p$.

**Asymptotic results with general convex costs.**    While the linear cost case is particularly tractable, the key result that the regime benefits from manipulation when the signal precision is high enough is more general. In particular:

**Proposition 4.** *For general convex cost functions* $\lim_{\alpha \to \infty} \theta^* = 0^+$. *In addition, if the cost function is strictly convex and* $C'(0) = 0$, *then* $\lim_{\alpha \to 0} \theta^* = 1^-$.

So, for high enough $\alpha$ (again, keeping $n$ fixed), the regime is indeed able to reduce the threshold $\theta^*$ below the Morris–Shin benchmark $1 - p$. Of course the threshold may rise significantly above the Morris–Shin benchmark when the signal precision is low enough, just as in the linear case. A striking example of this arises when the cost function is strictly convex with $C'(0) = 0$, for example a quadratic cost $C(a) = a^2/2$. In this case, as $\alpha \to 0$ the regime threshold not only rises above the Morris–Shin benchmark, it is driven all the way up to 1, the worst possible outcome for the regime.

With linear costs, as the signal precision becomes small enough the effective cost $c/\sqrt{\alpha}$ blows up, it is common knowledge that no regime will manipulate and so we have the Morris–Shin outcome. But with $C'(0) = 0$ the regime *always* has an incentive to manipulate even for very low levels of $\alpha$, since it can always take a very small action with very small marginal cost. For low values of $\alpha$, the regime would want to be able to commit to refrain from information manipulation. But because such commitments cannot be made, the regime is "trapped" into taking costly actions even as the precision falls. In this sense, the manipulation completely *backfires* on the regime when $\alpha$ is very low.

### 4.4. *Intuition and further discussion*

The result that signal-jamming is effective when the precision is sufficiently high depends on two key features of the model. First, manipulation is a *potentially* powerful tool when the precision is high, because then most citizens have signals near the mean and it takes only a small amount of bias to deliver a large change in the aggregate attack. Second, the regime is able to *use* this tool, because citizens face a difficult joint coordination-and-information-filtering problem that impedes their ability to infer the amount of bias in equilibrium.

**A powerful tool ….** To see why signal-jamming is potentially powerful, observe that in the Morris–Shin benchmark the size of the aggregate attack is $S_{\text{MS}}^*(\theta) = \Phi(\sqrt{\alpha}(x_{\text{MS}}^* - \theta))$, which approaches a step function $\mathbb{1}\{1 - p > \theta\}$ as $\alpha \to \infty$. A small "shock" that increased the signal mean from $\theta$ to $\theta + \tilde{a}$, say, would cause a shift of the step function so that the attack facing the marginal regime would fall from $1 - p$ to 0 and all regimes $\theta \geq 0$ would be able to survive. But this only delivers a reduction in the attack if $\tilde{a}$ is *unanticipated*. If the citizens could correctly anticipate such an increase in the signal mean, then they would discount their signals appropriately.

This illustrates two points: (i) the ability to shift the mean is potentially powerful when the precision is high, for then it takes only a small amount of bias to achieve a large reduction in the size of the attack and hence a large fall in $\theta^*$, but (ii) for this to be useful to the regime, it must be the case that the citizens are, somehow, impeded in their efforts to infer the bias.

**… and the regime can make use of it.** In turn, two features of the model account for the citizen's inability to infer the bias. First, different regimes have different incentives so that given imperfect information about $\theta$ there is also imperfect information about $a(\theta)$. Second, citizens are imperfectly coordinated. These two features are also key differences between the model of information manipulation here and canonical environments such as the careers concerns model of Holmström (1999) and the strategic information transmission model of Crawford and Sobel (1982).

**Contrast with career concerns models.** Citizens know the regime's incentives, so why can they not back out the regime's manipulation? The short answer is that the regime is *informed* and takes contingent actions $a(\theta)$ so that while citizens know the function $a(\cdot)$ in equilibrium, they do not know $\theta$ and hence do not know the exact value $a(\theta)$ to extract. Instead, they must infer $a(\theta)$ jointly with their inferences about $\theta$ itself and cannot perfectly decompose their signal into the true $\theta$ component and the endogenous bias component. In contrast, in standard *career concerns* models, such as Holmström (1999), the signal sender is itself uninformed and chooses a single action, $\tilde{a}$ say, to maximize its expected payoff. But since this action is common to all senders, it is then straightforward for signal receivers to infer that exact amount of bias and react accordingly.

**Contrast with strategic information transmission models.** As in traditional models of strategic information transmission, such as Crawford and Sobel (1982), the regime here makes use of a form of *partial pooling*. This is most clear in the case of linear costs where all regimes $\theta \in [\theta^*, \theta^{**})$ that intervene have the same apparent strength $y^*$. In this model, the partial pooling is achieved by different regimes taking different unobservable actions, namely $a(\theta) = y^* - \theta$. In contrast, in Crawford and Sobel this partial pooling takes the form of identical observable payoff-irrelevant messages.[6] In related work, Kartik *et al.* (2007) study a model of strategic information transmission where an informed sender faces a cost of sending biased messages (or where some fraction of receivers are naive and take messages at face-value) and show that there exist fully separating equilibria (*i.e.* in equilibrium receivers can infer the sender's type) and that these equilibria feature "inflated communication" in that the sender's message is greater than the sender's true type. The construction of these equilibria relies crucially on the assumption that the sender's type space is unbounded above. Kartik (2009) studies a related model with costs of sending biased messages but where the type space is instead bounded above and shows that the upper bound prevents full separation. While these models feature inflated communication, in contrast with my model the sender's information manipulation does not in fact deceive rational information receivers.

Another important difference is that information manipulation in my model is taken to influence the actions of heterogenous information receivers playing a game of *strategic complementarities*.

**Role of strategic complementarities.** To see the importance of strategic complementarities, recall that it is common knowledge that *any* manipulation serves to increase the regime's apparent strength and hence to reduce the aggregate attack. Since the individual $s_i$ and aggregate $S$ are strategic complements, this makes each individual less likely to participate. Consequently, the size of the attack facing the marginal regime falls. Moreover, when $\alpha$ is high, it takes only a small amount of bias to achieve a large reduction in $S$. Knowing this makes individuals even more reluctant to participate (a standard "multiplier effect" with strategic complementarities). As $\alpha \to \infty$, the regime threshold is driven to zero because in that limit it takes only an infinitesimal amount of bias to deter *all* citizens from attacking.

---

6. In my model the "message" $y(\theta)$ is observed with noise. See Blume *et al.* (2007) for a Crawford and Sobel model where information receivers sometimes see only noise rather than the sender's intended message. Blume *et al.* show that the presence of noise generally allows for more communication than the noiseless benchmark and that for low levels of noise this can be welfare-improving.

**Role of coordination.**    This argument crucially relies on the citizens being imperfectly coordinated, *i.e.* on individuals taking the aggregate $S$ as given. To see the importance of this, suppose to the contrary that citizens were *perfectly coordinated* and able to act as a single "large" agent who could force regime change for all $\theta < 1$. This agent receives a signal $x = \theta + a + \varepsilon$ with precision $\alpha \to \infty$. For simplicity, suppose also that costs are strictly convex. This implies $y(\theta)$ is *strictly* increasing and hence can be inverted to recover $\theta = y^{-1}(x)$. But knowing $\theta$ the single agent can deduce any manipulation $a(\theta)$ and discard it—so the regime has no incentive to undertake costly manipulation. Thus if citizens are perfectly coordinated, they know $x \to \theta$ and attack if and only if $x = \theta < 1$ and all regimes $\theta \in [0, 1)$ are wiped out. In contrast, if citizens are imperfectly coordinated, all regimes $\theta \in [0, 1)$ survive. In this sense, the imperfect coordination drives the equilibrium selection in the regime's favour (see Appendix B for more details).

**Contrast with Angeletos *et al*.**    The strategic interaction between an informed policy-maker and heterogeneous information receivers is a feature this model shares with Angeletos *et al.* (2006). A key difference, however, is that the equilibrium in my model is unique while their model features equilibrium multiplicity. The reason for this is a subtle difference in the information structure. In the version of their model closest to this article, Angeletos *et al.* let individuals receive two noisy signals, (i) a pure signal of $\theta$, and (ii) a separate signal on the endogenous action $a(\theta)$. In my model, in contrast, individuals get one noisy observation of the sum $y(\theta) = \theta + a(\theta)$, not separate signals for the two constituent parts. This makes an important difference.

In their setting individuals have two qualitatively different kinds of information; the pure signal is a monotone function of the policy-maker's type, but the signal on the endogenous action $a(\theta)$ is a non-monotone function of $\theta$. In particular, the policy-maker chooses $a(\theta) = a^*$ for some critical interval $\theta \in [\theta^*, \theta^{**})$ and $a(\theta) = \underline{a} < a^*$ otherwise. To see the importance of this, consider the benchmark version of their model where individuals observe this endogenous action without noise. In this case, the action creates *common certainty* that $\theta$ is either in the critical interval $[\theta^*, \theta^{**})$ or not. This common certainty then serves as the basis of multiple possible coordinated responses. If $a(\theta)$ is observed with unbounded idiosyncratic noise, then there is not common certainty but, so long as the noise on $a(\theta)$ is not too great, there is enough *common p-belief* (in the sense of Monderer and Samet, 1989) that individuals can likewise coordinate on multiple responses to their observations of the endogenous action.

In my setting, in contrast, while the *actions* $a(\theta)$ may be non-monotone in $\theta$, what matters for individual *signals* is the regime's apparent strength $y(\theta) = \theta + a(\theta)$ which, from Proposition 2, *is* monotone in $\theta$. Put differently, in their setting the actions $a(\theta)$ simply *are* the signal mean, whereas in my setting the actions are only part of the signal mean. The monotonicity of the signal is the key force driving equilibrium uniqueness in my setting; with monotonicity, there is not enough common p-belief to coordinate on multiple responses.

**Alternative models of information manipulation.**    In my model, the regime's information manipulation takes a particularly simple form, namely shifting the signal mean at a cost. Marinovic (2011) considers a related model where, with some fixed probability, an informed sender receives the opportunity to send a biased message (at no cost) but with complementary probability is instead compelled to report the truth. Information receivers are uninformed about *both* the sender's type and whether the sender has the opportunity to manipulate information. In equilibrium, a sender who has the opportunity to manipulate uses a mixed strategy (to avoid detection) which gives more likelihood to higher messages and consequently information receiver's find higher messages less believable (since it is more likely that they result from manipulation). Although different in

details, Marinovic (2011) shares with my benchmark model a focus on manipulation of the mean. Section 6.1 considers a quite distinct approach where the information sender strategically controls signal precision rather than the signal mean.


## 5. AN INFORMATION REVOLUTION

I now return to the general model where citizens have many signals and where the regime's cost of information manipulation and the number of signals are linked. For simplicity I focus on the case of linear costs, $C(a,n):=c(n)a$. Section 5.1 begins the analysis by showing that an increase in the cost of manipulation per se has the opposite effect on the regime's chances of survival to an increase in the total signal precision. Section 5.2 then considers an *information revolution*, with an increase in the number $n$ of signals simultaneously increasing the total signal precision $\alpha:=n\hat{\alpha}$ and increasing the regime's marginal cost of information manipulation $c(n)$, with offsetting implications for the regime's chances. The net effect is that the regime is made easier to overthrow unless there are relatively strong *economies of scale* in information control. Section 5.3 then interprets these implications of the model in light of the historical relationship between autocracies and information control.


### 5.1. *Changes in the cost of manipulation*

As a simple first step, I characterize the effect of a change in the marginal cost of information manipulation, continuing to hold fixed the number of signals $n$. As shown in Section 4.3 above, holding fixed the number of signals $n$ an increase in signal precision reduces a regime's cost of obtaining a given-sized reduction in the aggregate attack $S$, and, at least when $\alpha$ is high enough, this reduces the regime threshold $\theta^*$. This suggests there is a close link between changes in the signal precision and changes in the cost of manipulation. Indeed, from (24), the size of the aggregate attack facing the marginal regime $S(y^*)$ depends only on the ratio

$$\frac{c}{\sqrt{\alpha}} \tag{26}$$

not each of these parameters separately. The same is true for the regime threshold $\theta^*$ itself.

**Proposition 5.** *The marginal cost of manipulation $c$ and signal precision $\alpha$ only affect the regime threshold $\theta^*$ through the ratio $c/\sqrt{\alpha}$. Consequently, an increase in the marginal cost always has the opposite sign to the effect of an increase in the signal precision*

$$\frac{\partial \log\theta^*}{\partial \log c} = -2\frac{\partial \log\theta^*}{\partial \log\alpha}. \tag{27}$$

In particular, for high enough $\alpha$ we know that increasing the signal precision would reduce the threshold $\theta^*$ thereby increasing the regime's chances of surviving. But if so, it is also the case that increasing the marginal cost $c$ works in an offsetting way, decreasing the regime's chances of surviving. Moreover, this latter cost effect is *twice as large* as the signal precision effect. Of course the *level* of the regime threshold may well be $\theta^* < \theta^*_{MS} = 1-p$, so that overall the signal-jamming technology is effective for the regime—it is just that a marginal increase in $c$ makes the regime marginally easier to overthrow.

## 5.2. *Changes in the number of signals*

**Equilibrium with $n$ signals.**    With $n$ signals the analysis proceeds as before except that we need to explicitly write the total precision $\alpha := n\hat{\alpha}$ and marginal cost of information manipulation $c(n)$. The effective cost for a regime to obtain a given-sized attack $S$ is then

$$\frac{c(n)}{\sqrt{n\hat{\alpha}}} \tag{28}$$

namely the generalization of (26) to this setting with $n$ signals. If this ratio is sufficiently large, *i.e.* if $\hat{\alpha} < (c(n)/\phi(0))^2/n =: \underline{\alpha}_n$, then it is too expensive to manipulate and the equilibrium coincides with the Morris–Shin benchmark. Otherwise, there is an interval of regimes $[\theta_n^*, \theta_n^{**})$ that do manipulate and that choose apparent strength $y_n^* = x_n^* + \gamma_n$. These thresholds are determined as in (19)–(20) but with $c(n)$ replacing $c$ and $n\hat{\alpha}$ replacing $\alpha$, etc.

**Effect of an increase in the number of signals.**    The main result here is:

**Proposition 6.**    *The effect of an increase in the number of signals $n$ on the regime threshold $\theta_n^*$ is given by*

$$\frac{\partial \log \theta_n^*}{\partial \log n} = -2\left(\frac{c'(n)n}{c(n)} - \frac{1}{2}\right)\frac{\partial \log \theta_n^*}{\partial \log \hat{\alpha}}. \tag{29}$$

As before, if the per-unit signal precision is too low, $\hat{\alpha} < \underline{\alpha}_n$, then all regimes are at a corner and $\theta_n^* = 1 - p$, independent of $n$ (and $\hat{\alpha}$). In this case, both effects in (29) are zero. Otherwise, if $\hat{\alpha} > \underline{\alpha}_n$, the effect of the signal precision depends on whether $\hat{\alpha}$ is larger than a critical signal precision $\alpha_n^*$, analogous to (21) above. For $\hat{\alpha} > \alpha_n^*$ an increase in signal precision reduces the regime threshold $\theta_n^*$. For fixed $n$ this would help regimes survive, but now with $n$ increasing the overall effect also depends on the amount of curvature in $c(n)$. In short, the overall effect depends on the degree of *scale economies* in information control.

**Scale economies in information control.**    The curvature of $c(n)$ is a natural measure of the regime's ability to influence a changing informational environment. If it is harder to exert control over an ever-expanding array of media outlets, then that suggests the average cost $c(n)/n$ is increasing in $n$, *i.e.* diseconomies of scale in controlling media outlets. Alternatively, if there are complementarities in media control, with influence over one media outlet facilitating influence over others, then that suggests average cost $c(n)/n$ is decreasing in $n$, *i.e.* economies of scale in controlling media outlets.

Suppose we are in the more interesting scenario of a *high precision environment*, where an increase in the per-unit signal precision $\hat{\alpha}$ indeed reduces the regime threshold $\theta_n^*$. What then is the effect of an increase in $n$? Since the total precision $\alpha = n\hat{\alpha}$ is proportional to $n$ and the precision and cost effects work in offsetting ways, a priori one might reasonably conjecture that the case of constant average cost, $c(n) = n\hat{c}$ say, would be neutral for the regime. But it is not. In this case, the *effective cost* for the regime of obtaining a given-sized attack $S$ is $\hat{c}\sqrt{n/\hat{\alpha}}$, from (28), so doubling $n$ doubles the cost $n\hat{c}$, doubles the total precision $n\hat{\alpha}$ but *increases* the effective cost by a factor of $1/2$.

More generally, let $\epsilon(n)$ denote the elasticity of the average cost function $c(n)/n$, namely

$$\epsilon(n) := \frac{c'(n)n}{c(n)} - 1.$$

Then in a high precision environment, an increase in $n$ increases the regime threshold $\theta_n^*$ and makes the regime easier to overthrow unless the average cost $c(n)/n$ declines sufficiently fast—specifically, unless $\epsilon(n) < -1/2$. In short, not just scale economies but *sufficiently strong* scale economies.

This can be seen as a restriction on the number of media outlets $n$. For example, consider the cost function $c(n) = c_0 + c_1 n$ so that there is a "fixed cost" $c_0 > 0$ and a constant "marginal cost" of manipulating more signals $c_1 > 0$. Then the elasticity of average cost is $\epsilon(n) = -c_0/(c_0 + c_1 n)$ and the condition $\epsilon(n) < -1/2$ is satisfied only when the number of media outlets is relatively small, namely $n < (c_0/c_1)$. So in addition to influencing the information *content*, as measured by $a$, the regime may also have an incentive to directly seek control over the *number* of outlets $n$. Section 6.1 discusses an example of this.

### 5.3. *Implications of the model*

Thus whether a new information technology will be threatening to a regime will depend on those structural characteristics of the technology that determine the degree of scale economies in controlling it.

**Mass media and centralized technologies.** In the first half of the twentieth century, the relationship between then state-of-the-art technologies and autocratic regimes seemed particularly close, perhaps the most striking examples being the use of mass media propaganda by Nazi Germany and the Soviet Union (Friedrich and Brzezinski, 1965; Arendt, 1973). The information technologies underlying the immersive propaganda machinery of these regimes, most importantly mass print media, radio, and cinema, share the feature that they are relatively easy to *centralize*. By dominating the broadcast technology itself (*e.g.* radio stations, recording studios), and by permitting only a small number of approved outlets (*e.g.* the party newspaper), the regime can establish widespread control over the amount of information that individuals receive. Moreover, fear of punishment may also make it extremely risky for individuals to try to use these technologies to receive rival information (*e.g.* foreign radio broadcasts). This approach to information control has the feature that it involves a relatively large fixed cost—to establish and maintain the party broadcast machinery, to pursue covert attempts to circumvent official and quasi-official channels—but with that fixed establishment in place, the marginal cost of controlling additional media outlets is relatively small. In short, these are technologies most likely to feature large scale economies in information control and it is technologies of this kind that the model predicts ought to be most conducive to a regime's survival. So long as the fixed cost can be covered, an increase in the number of newspapers or radio stations or cinemas need not trouble the regime.

**Shared experiences and quasi-public information.** The mass media aspects of this kind of centralized propaganda machinery are important. In coordination settings, public information and related forms of shared experience typically have a disproportionate effect exactly because individuals want to coordinate their actions (*e.g.* Chwe, 2001). In this model, there is no public signal, but there are nonetheless important forms of shared information. In particular, the existence of the regime's propaganda arm is common knowledge, as is the direction of that propaganda, namely that it serves to suppress opposition. Given this, if individual signals are relatively tightly clustered around the regime's preferred message (either because the signals are innately technologically precise, *i.e.* the per-unit signal precision $\hat{\alpha}$ is high, or because there are strong enough scale economies and the total number $n$ of signals is high, or both) then propaganda will shift the mean belief of enough people that the regime's chances of surviving improve.

Thus although individuals have noisy information, the fact that their information comes from a channel of communication with a common component provides the regime with a potentially powerful tool. Put differently, what matters is not just the number of people who listen to a party broadcast on the radio or watch a mass parade in support of the regime, but also that fact that those people are keenly aware that their fellow citizens are sharing the experience, taking in the same information (hence also the importance of broadcasts in public or quasi-public locations, *e.g.* rallies in central locations, listening to radio broadcasts in the workplace). Given this, it is perhaps not surprising to learn that the Nazi regime actively subsidized the mass production of cheap radio sets; by 1939 some 70% of households owned a radio, the highest proportion in the world at the time (Zeman, 1973). Facilitating the diffusion of mass media technologies like this not only increases the mere fact of the regime's ability to communicate with many people simultaneously, it also increases the extent to which every individual knows that the information they are getting is similar to the information their fellow citizens are getting, *i.e.* increases shared experience.

**Social media and decentralized technologies.**    While cautionary tales about the role of then state-of-the-art technologies under autocracies in the twentieth century perhaps serves to dampen the most extreme forms of optimism about information technology and the prospects of regime change, does this have any relevance for thinking about the role of modern social media? The model does suggest that structural changes in the nature of technology may well make a difference. In particular, to the extent that social media technologies are *decentralized* and cannot easily be dominated by a single fixed propaganda establishment, the proliferation of new channels of communication may make it simply too costly for regimes to maintain their information control and thereby make it easier for them to be overthrown.

In practice, however, evidence on the effects of modern information technologies on autocratic regimes is mixed. On the one hand, the wave of uprisings against autocratic regimes in Tunisia, Egypt, Libya, Syria, and other Arab countries beginning in December 2010 has given renewed support to the idea that modern decentralized technologies can play a significant role in facilitating coordination against regimes. On the other hand, it is also clear that some regimes have managed to find effective ways to counter online organization and dissent. Kalathil and Boas (2003), for example, emphasize the Chinese regime's co-opting of modern information technologies against dissidents,[7] whereas Fallows (2008) provides a persuasive account of the surprising effectiveness of China's *national firewall* and the elaborate system of monitoring and censorship that the firewall interacts with. Similar examples of a regime's efforts to counteract modern technologies come from the 2009 presidential elections in Iran and subsequent demonstrations. In the run-up to the election, the regime suspended access to social media websites like Facebook. On election day, mobile phone communications were interrupted and foreign news providers such as the BBC experienced jamming designed to impede their broadcasts. During the subsequent demonstrations, services like Twitter were used by the regime to provide disinformation (Cohen, 2009; Esfandiari, 2010).

The model rationalizes this mixed evidence through two kinds of considerations. First, some kinds of social media should probably not be viewed as belonging to the "decentralized" category at all. To the extent that such mobilization can be prevented or significantly impeded by the use

---

7. Chase and Mulvenon (2002), in contrast, emphasize the Chinese regime's use of more traditional authoritarian methods in cracking down on online dissent. Kalathil and Boas (2003) also discuss related efforts by autocratic regimes to counteract the internet in Cuba, Saudi Arabia, and elsewhere. Soley (1987) discusses earlier efforts by the Cuban regime to counteract U.S. based satellite radio and television broadcasts.

of a national firewall or similar large-scale fixed investments in information control, then perhaps these technologies have more in common with older mass media technologies with an attendant risk that the scale economies in information control are sufficiently great so that, rather than being a risk to the regime, they instead increase its chances of survival. If so, such social media tools may facilitate online organization or dissent only to the extent that the regime permits a relatively porous firewall—perhaps because other costs, such a reduction in economic efficiency, are perceived as being too high to pay.

Second, not all extra sources of information should be viewed as an increase in the number of signals. A proliferation of Twitter accounts does not necessarily increase $n$. To the extent that these extra sources of information are highly correlated with existing sources—*e.g.* perhaps because an individual belongs to an online community and mostly obtains information from people with similar beliefs, creating an echo-chamber effect—these additional sources may not actually increase an individual's overall amount of information much. Moreover, to the extent that the regime is able to use these new technologies as a channel for strategic disinformation, then again the actual increase in information content may be less than the notional increase—either because of successful disinformation or simply because the latent possibility of disinformation reduces trust in the new technology in general.

## 6. EXTENSIONS

I now consider two alternative setups. Section 6.1 allows the regime to strategically control the *number* of signals (and hence the total signal precision), rather than the signal mean. In this setting, the regime engages in manipulation but cannot introduce any bias. Section 6.2 allows for a *struggle* over information as an opposition attempts to shift beliefs against the regime. The supplementary appendix contains several additional extensions.

### 6.1.  *Controlling the number of signals*

In the main model, manipulation shifts the signal mean so that individual $i$ has $j=1,...,n$ signals of the form $x_{ij}=y+\varepsilon_{ij}$ and the average signal $x_i=\frac{1}{n}\sum_{j=1}^{n}x_{ij}$ has total precision $\alpha=n\hat{\alpha}$. I now consider an alternative model where the regime strategically controls the *number* of signals and hence controls the total precision. To isolate these effects, I assume the regime *cannot introduce bias*, *i.e.* the signal mean is simply $y=\theta$.

Since what matters here is the total precision, I ignore integer constraints and assume there is an interval of *feasible* signals $[\underline{n},\overline{n}]$. Each of these signals is IID normal with per-unit precision normalized to $\hat{\alpha}=1$. By taking hidden action $a$, the regime chooses the fraction $\Phi(a)$ of these signals that are *operational*. The number of signals $n(a)$ available is then

$$n(a):=(1-\Phi(a))\underline{n}+\Phi(a)\overline{n}, \qquad 0<\underline{n}<\overline{n} \tag{30}$$

which is strictly increasing in $a$ with $n(-\infty)=\underline{n}$ and $n(+\infty)=\overline{n}$. When the regime is passive and takes no action, $a=0$, the number of signals available to a citizen is simply the midpoint $n(0)=(\underline{n}+\overline{n})/2=:\alpha$, which in this context should be thought of as the "natural" or "intrinsic" precision. By taking an action $a<0$, the regime pushes down the number of available signals below this natural level and so reduces signal precision (increases noise). By taking an action $a>0$, the regime pushes up the number of available signals thereby increasing signal precision. I assume the regime has a strictly convex cost function $C(a)$ that is symmetric about zero and increasing in the magnitude of $a$ with $C'(0)=0$.

For this model, I restrict attention to a monotone equilibrium where citizens participate if their signal is $x_i < x^*$ and the regime is overthrown if $\theta < \theta^*$, for thresholds $x^*, \theta^*$ to be determined. The main result here is:

**Proposition 7.** *A regime makes signals more noisy, $n(a(\theta)) < \alpha$, whenever $\theta \in [\theta^*, x^*)$ and makes signals more precise, $n(a(\theta)) > \alpha$, whenever $\theta > x^*$.*

Intuitively, if the regime has intermediate type $\theta \in [\theta^*, x^*)$ then it sets $n(a(\theta)) < \alpha$, i.e. it makes the average signal more noisy, muddying the signal so as to obscure its relative weakness. But if the regime has a high enough type, $\theta > x^*$, then it sets $n(a(\theta)) > \alpha$, i.e. it makes the average signal more precise so as to clarify its relative strength. Depending on parameters, it may be the case that $x^* < \theta^*$ in which case all regimes that survive, all $\theta \geq \theta^*$, choose $n(a(\theta)) \geq \alpha$ to clarify their strength. The left panel of Figure 4 illustrates with parameters chosen so that $x^* > \theta^*$, implying that the hidden actions jump down at $\theta^*$.

This version of the model is somewhat less tractable than the benchmark and it is difficult to analytically characterize the regime threshold $\theta^*$ as a function of the natural precision $\alpha$. But it is straightforward to solve for $\theta^*$ numerically. The right panel of Figure 4 shows $\theta^*$ as the natural precision $\alpha$ varies from 0 to 100. Notice that diminishing returns to information manipulation set in fast when the intrinsic $\alpha$ is high. A 5-fold increase in $\alpha$ from 0.5 to 2.5 hardly shifts the threshold $\theta^*$ at all. Despite the lack of bias in signals, the regime can still benefit from information manipulation in this setting. For high enough $\alpha$ we have $\theta^* < 1 - p$ so that again information manipulation is effective (see Appendix A for more details).

**Discussion.** This example shows that the regime can benefit from information manipulation *even if it is not able to introduce bias*. Edmond (2012) provides a related example, developing a model where citizens and the regime have quadratic preferences of the kind used in the cheap-talk literature and showing that while citizens in this setting can infer the bias in their information (and so there is no bias in equilibrium), their signals are endogenously noisier in a way that helps the regime prevent coordination against its interests.

The variation in the regime's preferences over the number of signals is reminiscent of results on information provision by a monopolist, as in Lewis and Sappington (1991, 1994) or Johnson and Myatt (2006). In Johnson and Myatt, for example, a monopolist's preferences over the posterior precision of consumer beliefs about a product depends on the amount of prior dispersion in beliefs; a monopolist facing a homogeneous population deploys a mass-market, low-markup strategy, and keeps information low, but a monopolist facing a more heterogeneous population deploys a more niche, high-markup strategy, and provides more information. In models of this kind, however, the private information lies on the consumers' side rather than the monopolist's and the ex post information provision (if any) is observable by consumers. In my model, in contrast, the ex post information provision varies with the regime's unobservable type and citizens have to make inferences about the endogenous precision of their information jointly with their inferences about the type itself.

## 6.2. *Struggles over information*

I now return to the setting of shifts in the signal mean with exogenous signal precision. I suppose, however, that there is an *opposition* and that if a regime is of type $\theta$ and takes action $a$ while the opposition takes action $e$, then citizens have signals $x_i = \theta + a - e + \varepsilon_i$ where $\varepsilon_i$ is IID normal with mean zero and precision $\alpha$. The regime's action $a$ increases the signal mean while the opposition's action $e$ decreases it. Both actions are unobserved by citizens.
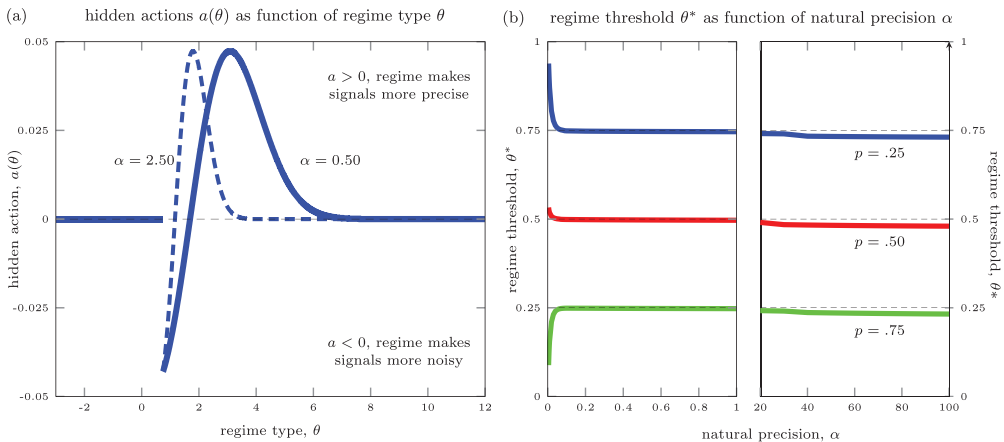
FIGURE 4

Hidden actions and threshold $\theta^*$ when regime can control the number of signals

*Notes*: Panel (a) shows the regime's actions to directly manipulate the number of signals $n(a)$. For intermediate $\theta$ it is optimal for $a(\theta) < 0$ so that the regime makes the signal more noisy than the natural precision, $n(a(\theta)) < \alpha$. For high $\theta$ it is optimal for $a(\theta) > 0$ so that the regime clarifies its strength by making the signal more precise, $n(a(\theta)) > \alpha$. In this example, the opportunity cost is $p = 0.25$. Panel (b) shows $\theta^*$ as a function of the $\alpha$ for various $p$. The regime benefits from information manipulation in that $\theta^* < 1 - p$ when $\alpha$ is high enough. All calculations use the bounds $\underline{n} = \alpha/2$, $\bar{n} = 3\alpha/2$ and cost function $C(a) = a^2/2$.

To highlight the struggle over manipulating information, I assume that the regime and the opposition *both* know the regime's type $\theta$. Along the equilibrium path citizens receive signals with mean $\theta + a(\theta) - e(\theta)$. If $a(\theta) = e(\theta)$, then the opposition simply undoes the efforts of the regime. I assume the regime has strictly convex cost $C(a)$ with $C'(0) = 0$ and the opposition has cost $C(\kappa e)/\kappa$ for some parameter $\kappa > 0$. If $\kappa > 1$ the regime has a cost advantage.

The payoff to the opposition is of the form $S - C(\kappa e)/\kappa$, so the opposition prefers the attack to be as large as possible, similar to the dissidents in Bueno de Mesquita (2010). Taking the aggregate attack $S(\theta, a, e)$ as given, an equilibrium in the subgame between the regime and the opposition consists of hidden actions $a(\theta), e(\theta)$ that are mutual best responses

$$a(\theta) \in \underset{a \geq 0}{\operatorname{argmax}} \{B(\theta, S(\theta, a, e(\theta))) - C(a)\} \tag{31}$$

$$e(\theta) \in \underset{e \geq 0}{\operatorname{argmax}} \{S(\theta, a(\theta), e) - C(\kappa e)/\kappa\}. \tag{32}$$

The regime's outside option introduces a key *asymmetry*. The regime does not care about the size of $S$ in those states where it is overthrown. In contrast, the opposition cares about $S$ both when the regime is overthrown and when it is not.

I again restrict attention to a monotone equilibrium where citizens participate if their signal is $x_i < x^*$ and the regime is overthrown if $\theta < \theta^*$, for thresholds $x^*, \theta^*$ to be determined. The main result here is:

**Proposition 8.** *Conditional on the regime's survival, $\theta \geq \theta^*$, the regime and opposition actions are proportional, $a(\theta) = \kappa e(\theta)$. Moreover, if the signal precision $\alpha$ is sufficiently high, then the regime threshold $\theta^*$ is strictly less than the Morris–Shin benchmark $\theta_{MS}^* = 1 - p$, despite the opposition.*
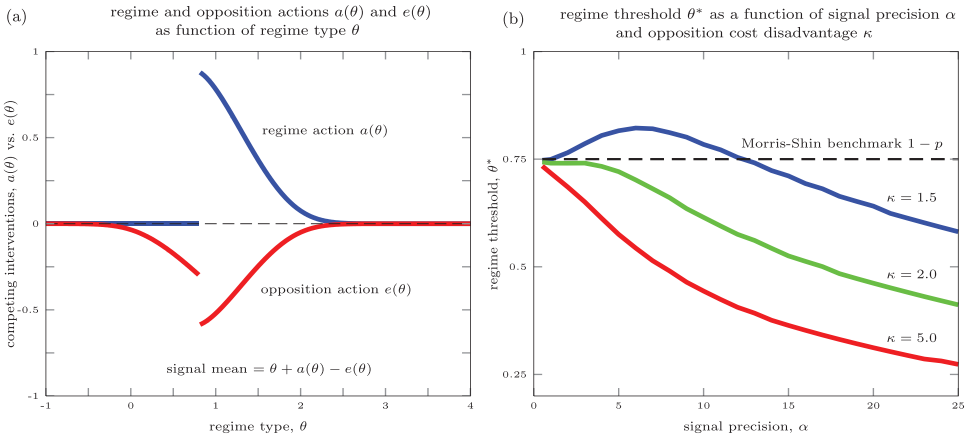
FIGURE 5

Hidden actions and regime threshold when there is an opposition

*Notes*: Panel (a) shows the regime's and opposition's actions $a(\theta)$ and $e(\theta)$ when there is a struggle over information. For clarity the opposition's $e(\theta)$ is plotted on a negative scale. For $\theta < \theta^*$, only the opposition takes an action. For $\theta \geq \theta^*$ the actions satisfy $a(\theta) = \kappa e(\theta)$, where $\kappa$ measures the costliness of the opposition's action. In this example, $\kappa = 1.5$ and the opposition's costs are greater than the regime's. Panel (b) shows $\theta^*$ as a function of the signal precision $\alpha$ for various $\kappa$. In these examples, the regime still benefits from information manipulation in that $\theta^* < 1 - p$ when $\alpha$ is high enough. In these calculations, the opportunity cost is $p = 0.25$ and the regime's cost function is $C(a) = a^2/2$ so that the opposition's cost function is $\kappa e^2/2$.

In a monotone equilibrium, the aggregate attack is $S(\theta, a, e) = \Phi(\sqrt{\alpha}(x^* - \theta - a + e))$. The first-order condition for the regime is $C'(a) = \sqrt{\alpha}\phi(\sqrt{\alpha}(x^* - \theta - a + e(\theta)))$ for $\theta \geq \theta^*$ while the first-order condition for the opposition is $C'(\kappa e) = \sqrt{\alpha}\phi(\sqrt{\alpha}(x^* - \theta - a(\theta) + e))$. Thus if the regime survives, $\theta \geq \theta^*$, then along the equilibrium path the marginal benefit of manipulation is the same for both the regime and the opposition and the level of manipulation is simply determined by their respective marginal costs $C'(a(\theta)) = C'(\kappa e(\theta))$ so that $a(\theta) = \kappa e(\theta)$ (since $C(\cdot)$ is strictly convex). In short, when the regime survives its actions are larger than those of the opposition if and only if the regime has a cost advantage, $\kappa > 1$.

The left panel of Figure 5 illustrates, with $\kappa > 1$ so that the regime's actions $a(\theta)$ are larger than the opposition's actions $e(\theta)$ on $\theta \geq \theta^*$. Just as the regime's hidden actions jump discretely to $a(\theta^*) > 0$ at the threshold $\theta = \theta^*$, so too do the opposition's actions typically jump at the threshold (though their jump may be up *or* down, depending on parameters). For $\theta < \theta^*$ only the opposition takes any action.

Does the presence of an opposition change the effectiveness of a regime's manipulation? On the one hand, it is true that the presence of the opposition generally moves the threshold $\theta^*$ against the regime (it is higher than it would be in the model where there is no opposition, $\kappa = \infty$). On the other hand, the regime still manipulates information and for high enough signal precision $\alpha$ the regime threshold $\theta^*$ is less than the Morris–Shin benchmark $\theta^*_{MS} = 1 - p$, just as in the benchmark model. The right panel of Figure 5 illustrates, showing the effects of changing $\alpha$ for various levels of $\kappa$. Not surprisingly, when $\kappa$ is relatively high, so that the opposition is at a pronounced disadvantage, it takes only a small increase in $\alpha$ for the regime's manipulation to be effective.

These results suggest that while the presence of organized opposition is important for understanding how much manipulation takes place and for the equilibrium *level* of the regime

threshold $\theta^*$, it is less important for the result that the regime threshold can be decreasing in the signal precision so that more precise signals move the threshold in the regime's favour.

## 7. CONCLUSIONS

I develop a signal-jamming model of information and political regime change. In contrast with familiar signal-jamming models, the regime's manipulation presents citizens with a difficult signal-extraction problem and the manipulation is often payoff improving for the regime. The model predicts that (i) holding fixed the number of signals available to citizens, an increase in the *per-unit* precision of individual signals can make signal-jamming more effective and thereby make the regime harder to overthrow, but (ii) an increase in the *number* of signals, which simultaneously increases both the total signal precision and the regime's costs of information manipulation, instead makes the regime easier to overthrow—unless there are strong economies of scale in information control.

The model thus allows for two kinds of information revolutions. In the first kind, perhaps best associated with the role of mass media propaganda under the totalitarian regimes of the early twentieth century, an information revolution consists of many high precision sources of information but these technologies are more *centralized* and subject to strong economies of scale. By establishing a large fixed propaganda apparatus, the regime can then control additional newspapers, radio stations, cinemas, etc at low marginal cost. This kind of information revolution is favourable to a regime's survival prospects. In the second kind, perhaps best associated with the rise of social media and other more *decentralized* sources of information, however, there are less likely to be strong economies of scale in information control. This kind of information revolution is unfavourable to a regime's survival prospects.

The coordination game studied in this article is deliberately stylized so as to focus attention on the effectiveness of the regime's manipulation and its sensitivity to changes in the information environment. For example, the model takes as given the *degree of influence* the regime has over citizens' sources of information. It would clearly be interesting to develop a richer model where the degree of influence over the media is itself an equilibrium outcome, in the spirit of Besley and Prat (2006) or Gehlbach and Sonin (2008), that needs to be determined jointly with the regime's manipulation and survival probability.

At a more general level, this article concerns a policy-maker seeking to induce information receivers to coordinate on the policy-maker's preferred outcome. Similar issues are likely to be at work in many scenarios. For example, the managers of a bank may try to influence information to prevent a run on deposits. But in market-oriented scenarios (such as a bank run, or a currency crisis), *market prices* are also likely to aggregate information and may interact with the policy-maker's actions in complicated ways (*cf.,* Angeletos and Werning, 2006; Hellwig *et al.*, 2006). In that sense, the analysis in this article should be viewed as suggestive of a key element that would be operative in a more complete, but also more complex, account of information manipulation in market-oriented coordination games.

## APPENDIX

### A. PROOFS AND OMITTED DERIVATIONS

#### A.1. *Morris–Shin Benchmark*

Let $\hat{x}, \hat{\theta}$ denote candidates for the critical thresholds. The posterior beliefs of a citizen with $x_i$ facing $\hat{\theta}$ are given by $\text{Prob}[\theta < \hat{\theta} | x_i] = \Phi(\sqrt{\alpha}(\hat{\theta} - x_i))$. A citizen with $x_i$ will participate if and only if $\Phi(\sqrt{\alpha}(\hat{\theta} - x_i)) \geq p$. This probability is continuous and strictly decreasing in $x_i$, so for each $\hat{\theta}$ there is a unique signal for which a citizen is indifferent. Similarly,

if the regime faces threshold $\hat{x}$ the aggregate attack is $\text{Prob}[x_i < \hat{x}|\theta] = \Phi(\sqrt{\alpha}(\hat{x}-\theta))$. A regime $\theta$ will not be overthrown if and only if $\theta \geq \Phi(\sqrt{\alpha}(\hat{x}-\theta))$. The probability on the right-hand side is continuous and strictly decreasing in $\theta$, so for each $\hat{x}$ there is a unique regime type that is indifferent. The Morris–Shin thresholds $x^*_{\text{MS}}, \theta^*_{\text{MS}}$ simultaneously solve these best response conditions as equalities, as stated in equations (4)–(5) in the main text. It is then straightforward to verify that there is only one solution to these equations and that $\theta^*_{\text{MS}} = 1 - p$ independent of $\alpha$ and $x^*_{\text{MS}} = 1 - p - \Phi^{-1}(p)/\sqrt{\alpha}$.

## A.2. *Proof of Proposition 1*

The proof shows first that (i) there is a unique equilibrium in monotone strategies, and (ii) that the unique monotone equilibrium is the only equilibrium which survives the iterative elimination of interim strictly dominated strategies. For ease of exposition, the proof is broken down into separate lemmas.

### (i) Unique equilibrium in monotone strategies.

**Regime problem.**    Let $\hat{x} \in \mathbb{R}$ denote a candidate for the citizens' threshold.

**Lemma 1.**    *For each $\hat{x} \in \mathbb{R}$, the unique solution to the regime's decision problem is characterized by a pair of functions, $\Theta : \mathbb{R} \to [0,1)$ and $A : \mathbb{R} \to \mathbb{R}_+$ such that if citizens participate for all $x_i < \hat{x}$ then the best-response of the regime is to abandon if and only if its type is $\theta < \Theta(\hat{x})$ and to choose an action $a(\theta) = 0$ for $\theta < \Theta(\hat{x})$ and $a(\theta) = A(\theta - \hat{x})$ for $\theta \geq \Theta(\hat{x})$.*

*Proof of Lemma 1.* To begin, let

$$S(w) := \Phi(-\sqrt{\alpha}w). \tag{A.1}$$

The auxiliary function $S(w)$ is exogenous and does not depend on $\hat{x}$. In terms of this function, the aggregate attack facing the regime is

$$\int_{-\infty}^{\hat{x}} \sqrt{\alpha}\phi(\sqrt{\alpha}(x_i - \theta - a))dx_i = \Phi(\sqrt{\alpha}(\hat{x} - \theta - a)) = S(\theta + a - \hat{x}). \tag{A.2}$$

Since the regime has access to an outside option normalized to zero, its problem can be written

$$V(\theta, \hat{x}) := \max[0, W(\theta, \hat{x})] \tag{A.3}$$

where $W(\theta, \hat{x})$ is the best payoff regime $\theta$ can get if it is not overthrown

$$W(\theta, \hat{x}) := \max_{a \geq 0} \left[ \theta - S(\theta + a - \hat{x}) - C(a) \right]. \tag{A.4}$$

From the envelope theorem, the partial derivative $W_\theta(\theta, \hat{x}) = 1 - S'(\theta - \hat{x} + a) > 1$ since $S'(w) < 0$ for all $w \in \mathbb{R}$. Since $S(w) \geq 0$ and $C(a) \geq 0$ we know $W(\theta, \hat{x}) < 0$ for all $\theta < 0$ and all $\hat{x}$. Similarly, $W(1, \hat{x}) > 0$ for all $\hat{x}$. So by the intermediate value theorem there is a unique $\Theta(\hat{x}) \in [0,1)$ such that $W(\Theta(\hat{x}), \hat{x}) = 0$. And since $W_\theta(\theta, \hat{x}) > 1$ the regime is overthrown if and only if $\theta < \Theta(\hat{x})$. Since positive actions are costly, the regime takes no action for $\theta < \Theta(\hat{x})$. Otherwise, for $\theta \geq \Theta(\hat{x})$, the actions of the regime are given by

$$a(\theta) = A(\theta - \hat{x}) \tag{A.5}$$

where the function $A(t)$ is exogenous and does not depend on $\hat{x}$. This auxiliary function is defined by:

$$A(t) := \operatorname*{argmin}_{a \geq 0}[S(t + a) + C(a)] \tag{A.6}$$

The first-order necessary condition for interior solutions can be written $C'(a) = -S'(t + a)$, and, on using the formula for $S(\cdot)$ in equation (A.1) above, $C'(a) = \sqrt{\alpha}\phi(\sqrt{\alpha}(t + a))$ where $\phi(w) := \exp(-w^2/2)/\sqrt{2\pi}$ for all $w \in \mathbb{R}$. This first-order condition may have zero, one or two solutions for each $t$. If for a given $t$ there are zero (interior) solutions, then $A(t) = 0$. If for given $t$ there are two solutions, one of them can be ruled out by the second-order sufficient condition $\alpha\phi'(\sqrt{\alpha}(t + a)) + C''(a) > 0$. Using the property $\phi'(w) = -w\phi(w)$ for all $w \in \mathbb{R}$ shows that if there are two solutions to the first-order condition, only the "higher" of them satisfies the second-order condition. Therefore for each $t$ there is a single $A(t)$ that solves the regime's problem. Making the substitution $t = \theta - \hat{x}$, the regime's threshold $\Theta(\hat{x})$ is then found from the indifference condition $W(\Theta(\hat{x}), \hat{x}) = 0$, or more explicitly

$$\Theta(\hat{x}) = S(\Theta(\hat{x}) - \hat{x} + A(\Theta(\hat{x}) - \hat{x})) + C(A(\Theta(\hat{x}) - \hat{x})). \tag{A.7}$$

Taking $\hat{x}$ as given, equations (A.6) and (A.7) give the regime threshold $\Theta(\hat{x})$ and the hidden actions $a(\theta) = A(\theta - \hat{x})$ that solve the regime's problem.    ∥

**Citizen problem.**    Let $\hat{\theta} \in [0,1)$ and $a : \mathbb{R} \to \mathbb{R}_+$ denote, respectively, a candidate for the regime's threshold and a candidate for the regime's hidden actions with $a(\theta) = 0$ for $\theta < \hat{\theta}$.

**Lemma 2.**    *For each $\hat{\theta} \in [0,1)$ and $a : \mathbb{R} \to \mathbb{R}_+$*

   *(a) The unique solution to the problem of a citizen with signal $x_i$ is given by a mapping $P(\cdot | a(\cdot)) : \mathbb{R} \times \mathbb{R} \to [0,1]$ such that the citizen participates if and only if*

$$P(x_i, \hat{\theta} | a(\cdot)) := \text{Prob}[\theta < \hat{\theta} | x_i, a(\cdot)] \geq p \tag{A.8}$$

   *where $P$ is continuous and strictly decreasing in $x_i$ with limits $P(-\infty, \hat{\theta} | a(\cdot)) = 1$ and $P(+\infty, \hat{\theta} | a(\cdot))) = 0$ for any $\hat{\theta}$ and function $a(\cdot)$ satisfying $a(\theta) = 0$ for $\theta < \hat{\theta}$.*

   *(b) For any candidate citizen threshold $\hat{x}$, with implied regime threshold $\Theta(\hat{x})$ and hidden actions $A(\theta - \hat{x})$, an individual citizen with signal $x_i$ participates if and only if its signal is such that*

$$K(x_i, \hat{x}) := \text{Prob}[\theta < \Theta(\hat{x}) | x_i, A(\cdot)] \geq p \tag{A.9}$$

   *where $K : \mathbb{R} \times \mathbb{R} \to [0,1]$ is continuous, strictly increasing in $x_i$ with limits $K(-\infty, \hat{x}) = 0$ and $K(+\infty, \hat{x}) = 1$ for any $\hat{x}$. Moreover, $K(x_i, \hat{x}) = \text{Prob}[\theta < \Theta(\hat{x}) - \hat{x} | x_i - \hat{x}, A(\cdot)]$ for any $\hat{x}$.*

*Proof of Lemma 2.* (a) For notational simplicity, write $x$ for an individual's signal, $\theta$ for the regime threshold, and $P(x, \theta)$ for the probability an individual with $x$ assigns to the regime's type being less than $\theta$ when the actions are $a : \mathbb{R} \to \mathbb{R}_+$. That is,

$$P(x, \theta) = \frac{\int_{-\infty}^{\theta} \sqrt{\alpha} \phi(\sqrt{\alpha}(x_i - t)) dt}{\int_{-\infty}^{\infty} \sqrt{\alpha} \phi(\sqrt{\alpha}(x_i - t - a(t))) dt} \tag{A.10}$$

where the numerator uses $a(t) = 0$ for $t < \theta$. Hence $P : \mathbb{R} \times \mathbb{R} \to [0,1]$ is continuous in $x, \theta$. This probability can be written

$$P(x, \theta) = \frac{N(\theta - x)}{N(\theta - x) + D(x, \theta)} \tag{A.11}$$

where

$$N(\theta - x) := \Phi(\sqrt{\alpha}(\theta - x)), \quad \text{and} \quad D(x, \theta) := \int_{\theta}^{\infty} \sqrt{\alpha} \phi(\sqrt{\alpha}(x - \xi - a(\xi))) d\xi. \tag{A.12}$$

Differentiating (A.11) shows $P_x < 0$ if and only if $N'/N > -D_x/D$. Calculating the derivatives shows that this is equivalent to

$$H(\sqrt{\alpha}(x - \theta)) > -\frac{\int_{\theta}^{\infty} \phi'(\sqrt{\alpha}(x - y(\xi))) d\xi}{\int_{\theta}^{\infty} \phi(\sqrt{\alpha}(x - y(\xi))) d\xi} = \frac{\int_{\theta}^{\infty} \sqrt{\alpha}(x - y(\xi)) \phi(\sqrt{\alpha}(x - y(\xi))) d\xi}{\int_{\theta}^{\infty} \phi(\sqrt{\alpha}(x - y(\xi))) d\xi} \tag{A.13}$$

where $H(w) := \phi(w)/(1 - \Phi(w)) > 0$ denotes the standard normal *hazard function* for $w \in \mathbb{R}$, where $y(\xi) := \xi + a(\xi)$ is the mean of the signal distribution if $\xi \geq \theta$, and where the equality follows from $\phi'(w) = -w\phi(w)$ for all $w$. Now define a density $\varphi(\xi | x) > 0$ by

$$\varphi(\xi | x) := \frac{\phi(\sqrt{\alpha}(x - y(\xi)))}{\int_{\theta}^{\infty} \phi(\sqrt{\alpha}(x - y(\xi'))) d\xi'}, \qquad \xi \in [\theta, \infty). \tag{A.14}$$

Then after a slight rearrangement of terms in (A.13), $P_x < 0$ if and only if

$$H(\sqrt{\alpha}(x - \theta)) - \sqrt{\alpha}(x - \theta) > \sqrt{\alpha} \left[ \theta - \int_{\theta}^{\infty} y(\xi) \varphi(\xi | x) d\xi \right]. \tag{A.15}$$

Since the hazard function satisfies $H(w) > w$ for all $w \in \mathbb{R}$ and $\alpha > 0$, it is sufficient that

$$\int_{\theta}^{\infty} y(\xi) \varphi(\xi | x) d\xi \geq \theta. \tag{A.16}$$

But since $y(\xi) := \xi + a(\xi)$, $\xi \geq \theta$, and $a(\xi) \geq 0$, condition (A.16) is always satisfied. Therefore $P_x < 0$. Since $N' > 0$ and $D_\theta < 0$, $P_\theta > 0$ for all $x, \theta$. Moreover, since $N(-\infty) = 0$ and $D > 0$ we have $P(x, -\infty) = 0$ for all $x$. Similarly, since $a(\xi) = 0$ for all $\xi < \theta$ as $\theta \to \infty$ we have $D(x, \theta) \to 1 - N(\theta - x)$ and since $N(+\infty) = 1$ this means $D(x, +\infty) = 0$ for all $x$. Therefore $P(x, +\infty) = 1$ for all $x$. The limit properties in $x$ are established in parallel fashion.

   (b) Fix a $\hat{x} \in \mathbb{R}$ and let $A(\theta - \hat{x})$ denote the associated hidden actions. Analogous to (A.11), write $P(x, \theta, \hat{x}) = N(\theta - x)/[N(\theta - x) + D(x, \theta, \hat{x})$ where $N : \mathbb{R} \to [0,1]$ is defined as in (A.12) above and where

$$D(x, \theta, \hat{x}) := \int_{\theta}^{\infty} \sqrt{\alpha} \phi(\sqrt{\alpha}(x - t - A(t - \hat{x}))) dt. \tag{A.17}$$

Now define $K(x, \hat{x}) := P(x, \Theta(\hat{x}), \hat{x})$. That $K(x, \hat{x})$ is continuous and decreasing in $x$ is immediate from part (a) above. Finally, for $K(x, \hat{x}) = P(x - \hat{x}, \Theta(\hat{x}) - \hat{x}, 0)$ it is sufficient that $D(x, \Theta(\hat{x}), \hat{x}) = D(x - \hat{x}, \Theta(\hat{x}) - \hat{x}, 0)$. From (A.17) and using the change of variables $\xi := \theta - \hat{x}$ we have

$$D(x, \Theta(\hat{x}), \hat{x}) = \int_{\Theta(\hat{x}) - \hat{x}}^{\infty} \sqrt{\alpha} \phi(\sqrt{\alpha}(x - \hat{x} - \xi - a(\xi))) d\xi = D(x - \hat{x}, \Theta(\hat{x}) - \hat{x}, 0). \tag{A.18}$$

Therefore $K(x, \hat{x}) = P(x - \hat{x}, \Theta(\hat{x}) - \hat{x}, 0) = \text{Prob}[\theta < \Theta(\hat{x}) - \hat{x} | x - \hat{x}, A(\cdot)]$ as claimed.    ‖

**Fixed point.**    A citizen with signal $x_i$ will participate if and only if $K(x_i, \hat{x}) \geq p$. Since $K(x_i, \hat{x})$ is strictly increasing in $x_i$ with $K(-\infty, \hat{x}) < p$ and $K(+\infty, \hat{x}) > p$ for any $\hat{x} \in \mathbb{R}$, there is a unique signal $\psi(\hat{x})$ solving

$$K(\psi(\hat{x}), \hat{x}) = p \tag{A.19}$$

such that a citizen with signal $x_i$ participates if and only if $x_i < \psi(\hat{x})$.

**Lemma 3.**    *The function $\psi : \mathbb{R} \to \mathbb{R}$ is continuous and has a unique fixed point $x^* = \psi(x^*)$ with derivative $\psi'(x^*) \in (0, 1)$ at the fixed point. Moreover $\psi(x) \leq x^*$ for all $x < x^*$ and $\psi(x) \geq x^*$ for all $x > x^*$.*

*Proof of Lemma 3.* Since $K(x, \hat{x})$ is continuously differentiable in $x$, an application of the implicit function theorem to (A.19) shows that $\psi(\cdot)$ is continuous. Fixed points of $\psi(\cdot)$ satisfy $x^* = \psi(x^*)$. Equivalently, by part (b) of Lemma 2, they satisfy $K(x^*, x^*) = P(0, \Theta(x^*) - x^*, 0) = p$, where $\Theta(\hat{x})$ is the critical state in the regime's problem (A.6)–(A.7). By Lemma 2 and the intermediate value theorem there is a unique $z^* \in \mathbb{R}$ such that $P(0, z^*, 0) = p$. Then applying the implicit function theorem to (A.6)–(A.7) gives

$$\Theta'(\hat{x}) = \frac{\sqrt{\alpha}\phi(\sqrt{\alpha}(\hat{x} - \Theta(\hat{x}) - A(\Theta(\hat{x}) - \hat{x})))}{1 + \sqrt{\alpha}\phi(\sqrt{\alpha}(\hat{x} - \Theta(\hat{x}) - A(\Theta(\hat{x}) - \hat{x})))} \in (0, 1). \tag{A.20}$$

Since $\Theta(-\infty) = 0$ and $\Theta(+\infty) = 1$, there is a unique $x^* \in \mathbb{R}$ such that $\Theta(x^*) - x^* = z^*$, hence $\psi(\cdot)$ has a unique fixed point, the same $x^*$. Now using part (b) of Lemma 2 and implicitly differentiating (A.19) we have

$$\psi'(\hat{x}) = 1 + \frac{P_\theta(\psi(\hat{x}) - \hat{x}, \Theta(\hat{x}) - \hat{x}, 0)}{P_x(\psi(\hat{x}) - \hat{x}, \Theta(\hat{x}) - \hat{x}, 0)}(1 - \Theta'(\hat{x})). \tag{A.21}$$

By Lemma 2, $P_\theta > 0$ and $P_x < 0$ and $\Theta'(\hat{x}) \in (0, 1)$ from (A.20). Therefore $\psi'(\hat{x}) < 1$ for all $\hat{x}$. To see that $\psi'(x^*) > 0$, first notice that it is sufficient that $P_\theta/P_x \geq -1$ when evaluated at $\hat{x} = x^*$. Calculating the derivatives shows that this is true if and only if

$$\phi(\sqrt{\alpha}(y(\theta^*) - x^*)) + \int_{\theta^*}^{\infty} \sqrt{\alpha}\phi'(\sqrt{\alpha}(y(\theta) - x^*))d\theta \leq 0 \tag{A.22}$$

where $\theta^* := \Theta(x^*)$ and where $y(\theta) = \theta + a(\theta)$ is the mean of the signal distribution from which a citizen is sampling if the regime has type $\theta \geq \theta^*$. To show that this condition always holds, we need to consider the cases of linear costs and strictly convex costs separately. If costs are linear, $C(a) = ca$, then if $c \geq \bar{c} := \sqrt{\alpha}\phi(0)$ the result is trivial because $a(\theta) = 0$ for all $\theta \in \mathbb{R}$. So suppose $c < \bar{c}$. Then $a(\theta) = \max[0, x^* + \gamma - \theta]$ where $\gamma := \sqrt{2\log(\sqrt{\alpha}\phi(0)/c)/\alpha} > 0$. Calculating the integral and then simplifying shows that (A.22) holds if and only if $-\alpha\gamma\phi(\sqrt{\alpha}\gamma)a(\theta^*) \leq 0$ which is true because $a(\theta^*) \geq 0$. If costs are strictly convex, then from the optimality conditions for the regime's choice of action we have that $a(\theta) > 0$ for all $\theta \geq \theta^*$ and

$$\sqrt{\alpha}\phi(\sqrt{\alpha}(y(\theta) - x^*)) = C'(a(\theta)), \quad \theta \geq \theta^* . \tag{A.23}$$

Differentiating with respect to $\theta$ gives

$$\alpha\phi'(\sqrt{\alpha}(y(\theta) - x^*))y'(\theta) = C''(a(\theta))a'(\theta), \quad \theta \geq \theta^* . \tag{A.24}$$

Using the associated second-order condition shows that $y'(\theta) > 0$ for $\theta \geq \theta^*$. Since $y(\cdot)$ is invertible, a change of variables shows that (A.22) holds if and only if

$$\int_{\theta^*}^{\infty} \phi'(\sqrt{\alpha}(y(\theta) - x^*))\frac{a'(\theta)}{y'(\theta)}d\theta \geq 0 . \tag{A.25}$$

Using (A.24) we equivalently have the condition

$$\int_{\theta^*}^{\infty} \frac{\phi'(\sqrt{\alpha}(y(\theta) - x^*))^2}{C''(a(\theta))}d\theta \geq 0 \tag{A.26}$$

which is true since the integrand is non-negative. Therefore, $P_\theta/P_x \geq -1$ at $\hat{x} = x^*$ and $\psi'(x^*) > 0$.

Finally, $\psi(\hat{x}) \leq x^*$ for every $\hat{x} < x^*$ is proven by contradiction. Suppose not. Then by continuity of $\psi$ there exists $\tilde{x} < x^*$ such that $\psi(\tilde{x}) = x^*$. Moreover, since $\psi'(x^*) > 0$, we must have $\psi'(\tilde{x}) < 0$ for at least one such $\tilde{x}$. Since $\psi(\tilde{x}) = x^*$ and $K(x^*, x^*) = p$, under this hypothesis we can write $K(\psi(\tilde{x}), \psi(\tilde{x})) = p$ so by the implicit function theorem $\psi(\tilde{x})$ must satisfy

$$\psi'(\tilde{x})[K_1(x^*, x^*) + K_2(x^*, x^*)] = 0 \tag{A.27}$$

where the hypothesis $\psi(\tilde{x}) = x^*$ is used to evaluate the partial derivatives $K_1$ and $K_2$. Since $\psi'(\tilde{x}) < 0$, this can only be satisfied if $K_1(x^*, x^*) + K_2(x^*, x^*) = 0$. But for any $\hat{x} \in \mathbb{R}$, the value $\psi(\hat{x})$ is implicitly defined by $K(\psi(\hat{x}), \hat{x}) = p$ so that by the implicit function theorem $\psi'(\hat{x}) = -K_2(\psi(\hat{x}), \hat{x})/K_1(\psi(\hat{x}), \hat{x})$. From (A.21) we know $\psi'(\hat{x}) < 1$ for any $\hat{x}$ and since $K_1 < 0$ from Lemma 2 we conclude $K_1(\psi(\hat{x}), \hat{x}) + K_2(\psi(\hat{x}), \hat{x}) < 0$ for any $\hat{x}$. For $\hat{x} = x^*$ in particular, $K_1(x^*, x^*) + K_2(x^*, x^*) < 0$ so we have the needed contradiction. Therefore $\psi(\hat{x}) \leq x^*$ for every $\hat{x} < x^*$. A symmetric argument shows $\psi(\hat{x}) \geq x^*$ for every $\hat{x} > x^*$.    $\|$

**Concluding that there is a unique equilibrium in monotone strategies.** To conclude part (i) of the proof, we take an arbitrary $\hat{x} \in \mathbb{R}$ and solve the regime's problem to get $\Theta(\hat{x})$ and $a(\theta, \hat{x}) = A(\theta - \hat{x})$ using the auxiliary function from Lemma 1. We use these functions to construct $K(x_i, \hat{x})$ from (A.9) for each signal $x_i \in \mathbb{R}$ and use Lemma 2 to conclude that in particular $K(\hat{x}, \hat{x}) = P(0, \Theta(\hat{x}) - \hat{x}, 0)$. We then use the intermediate value theorem to deduce that there is a unique $z^* \in \mathbb{R}$ such that $P(0, z^*, 0) = p$. This gives a unique difference $z^* = \theta^* - x^*$ that can be plugged into the regime's indifference condition (A.7) to get the unique $\theta^* = \Theta(x^*) \in [0, 1)$ such that the regime is overthrown if and only if $\theta < \theta^*$. The unique signal threshold is then $x^* = \theta^* - z^*$ and the unique hidden action function is given by $a(\theta) := A(\theta - x^*)$.

*(ii) Iterative elimination of interim strictly dominated strategies.* We can now go on to show that there is no other equilibrium. The argument begins by showing that for sufficiently low signals it is a dominant strategy to participate in the attack on the regime and for sufficiently high signals it is a dominant strategy to not participate.

**Dominance regions.** If the regime has $\theta < 0$, any mass $S \geq 0$ can overthrow the regime. Similarly, if the regime has $\theta \geq 1$ it can never be overthrown. Any regime that is overthrown takes no action, since to do so would incur a cost for no gain. Similarly, any regime $\theta$ that is not overthrown takes an action no larger than the $a$ such that $\theta = C(a)$. Any larger action must result in a negative payoff which can be improved upon by taking the outside option. Given this:

**Lemma 4.** *There exists a pair of signals $\underline{x} < \overline{x}$, both finite, such that $s(x_i) = 1$ is strictly dominant for $x_i < \underline{x}$ and $s(x_i) = 0$ is strictly dominant for $x_i > \overline{x}$.*

*Proof of Lemma 4.* The most *pessimistic* scenario for any citizen is that regimes are overthrown only if $\theta < 0$ and that regimes take the largest hidden actions that could be rational $\underline{a}(\theta) := C^{-1}(\theta)$ for $\theta \geq 0$ and zero otherwise. Let $\underline{P}(x_i) := \text{Prob}[\theta < 0 \,|\, x_i, \underline{a}(\cdot)]$ denote the probability the regime is overthrown in this most pessimistic scenario. Part (a) of Lemma 2 holds for hidden actions of the form $\underline{a}(\theta)$ and implies $\underline{P}'(x_i) < 0$ for all $x_i$, and since $\underline{P}(-\infty) = 1$ and $\underline{P}(+\infty) = 0$ by the intermediate value theorem there is a unique value, $\underline{x}$, finite, such that $\underline{P}(\underline{x}) = p$. For $x_i < \underline{x}$ it is (iteratively) strictly dominant for $s(x_i) = 1$. Similarly, the most *optimistic* scenario for any citizen is that regimes are overthrown if $\theta < 1$ and that regimes take the smallest hidden actions that could be rational $\overline{a}(\theta) := 0$. Let $\overline{P}(x_i) := \text{Prob}[\theta < 1 \,|\, x_i, \overline{a}(\cdot)]$ denote the probability the regime is overthrown in this most optimistic scenario. A parallel argument establishes the existence of a unique value, $\overline{x}$, finite, such that $\overline{P}(\overline{x}) = p$. For $x_i > \overline{x}$ it is (iteratively) strictly dominant for $s(x_i) = 0$. ‖

**Iterative elimination.** Starting from the dominance regions implied by $\underline{x}$ and $\overline{x}$ it is then possible to iteratively eliminate (interim) strictly dominated strategies. Recall that $S(w) := \Phi(-\sqrt{\alpha}w)$ and $A(t) := \text{argmin}_{a \geq 0}[S(t + a) + C(a)]$. Again, these auxiliary functions do not depend on any endogenous variable and in particular do not depend on citizen thresholds.

**Lemma 5.** *Let $x_{n+1} = \psi(x_n)$ for $n = 0, 1, 2, \ldots$ where*

$$K(\psi(x_n), x_n) = p$$

(a) *If it is strictly dominant for $s(x_i) = 1$ for all $x_i < \underline{x}_n$, then the regime is overthrown for at least all $\theta < \underline{\theta}_n := \Theta(\underline{x}_n)$ where the function $\Theta : \mathbb{R} \to [0, 1)$ solves*

$$\Theta(x) = S(\Theta(x) - x + A(\Theta(x) - x)) + C((\Theta(x) - x)). \tag{A.28}$$

*Similarly, if it is strictly dominant for $s(x_i) = 0$ for all $x_i > \overline{x}_n$, then regime is not overthrown for at least all $\theta > \overline{\theta}_n := \Theta(\overline{x}_n)$.*

(b) *Moreover, if it is strictly dominant for $s(x_i) = 1$ for all $x_i < \underline{x}_n$, then it is strictly dominant for $s(x_i) = 1$ for all $x_i < \underline{x}_{n+1} = \psi(\underline{x}_n)$. Similarly, if it is strictly dominant for $s(x_i) = 0$ for all $x_i > \overline{x}_n$, then it is strictly dominant for $s(x_i) = 0$ for all $x_i > \overline{x}_{n+1} = \psi(\overline{x}_n)$.*

*Proof of Lemma 5.* (a) Fix an $\underline{x}_n$ and $\overline{x}_n$ such that citizens with signals $x_i < \underline{x}_n$ have $s(x_i) = 1$ and likewise citizens with signals $x_i > \overline{x}_n$ have $s(x_i) = 0$. From Lemma 4 this can be done at least for the signals $\underline{x}, \overline{x}$ that determine the bounds of the dominance regions. All citizens with signals $x_i < \underline{x}_n$ have $s(x_i) = 1$ so the mass of subversives is at least $\Phi(\sqrt{\alpha}(\underline{x}_n - \theta - a))$. To acknowledge this, write the total mass of subversives

$$\Phi(\sqrt{\alpha}(\underline{x}_n - \theta - a)) + \Delta(\theta + a) \tag{A.29}$$

for some function $\Delta : \mathbb{R} \to [0, 1]$. First consider the case $\Delta(\cdot) = 0$ where *only* citizens with $x_i < \underline{x}_n$ subvert the regime. From Lemma 1 there is a unique threshold $\underline{\theta}_n := \Theta(\underline{x}_n) \in [0, 1)$ sustained by hidden actions $a(\theta) = A(\theta - \underline{x}_n)$ solving (A.6)–(A.7)

such that the regime is overthrown if $\theta < \underline{\theta}_n = \Theta(\underline{x}_n)$. Now consider the case $\Delta(\cdot) > 0$ where *some* citizens with signals $x_i \geq \underline{x}_n$ also subvert the regime. The proof that the regime is overthrown for at least all $\theta < \Theta(\underline{x}_n)$ is by contradiction. Suppose that when $\Delta(\cdot) > 0$ regime change occurs for all $\theta < \tilde{\theta}_n$ for some $\tilde{\theta}_n \leq \Theta(\underline{x}_n)$. A marginal regime $\tilde{\theta}_n$ must be indifferent between being overthrown and taking the outside option, so this threshold satisfies $\tilde{\theta}_n = S(\tilde{\theta}_n + \tilde{a}_n - \underline{x}_n) + C(\tilde{a}_n)$ where $\tilde{a}_n \geq 0$ is the optimal action for the marginal regime $\tilde{\theta}_n$. Then observe

$$\begin{aligned}
\Theta(\underline{x}_n) &= \Phi[\sqrt{\alpha}(\underline{x}_n - \Theta(\underline{x}_n) - A(\Theta(\underline{x}_n) - \underline{x}_n))] + C[A(\Theta(\underline{x}_n) - \underline{x}_n)] \\
&\leq \Phi[\sqrt{\alpha}(\underline{x}_n - \Theta(\underline{x}_n) - a)] + C(a), \\
&< \Phi[\sqrt{\alpha}(\underline{x}_n - \Theta(\underline{x}_n) - a)] + \Delta(\tilde{\theta}_n + \tilde{a}_n) + C(a), \qquad \text{for any } a \geq 0
\end{aligned}$$

where the first inequality follows because $A(\cdot)$ minimizes $\Phi[\sqrt{\alpha}(\underline{x}_n - \theta - a)] + C(a)$ and where the second inequality follows from $\Delta(\cdot) > 0$. Taking $a = \tilde{a}_n \geq 0$ we then have

$$\begin{aligned}
\Theta(\underline{x}_n) &< \Phi[\sqrt{\alpha}(\underline{x}_n - \Theta(\underline{x}_n) - \tilde{a}_n)] + \Delta(\tilde{\theta}_n + \tilde{a}_n) + C(\tilde{a}_n) \\
&= \Phi[\sqrt{\alpha}(\underline{x}_n - \Theta(\underline{x}_n) - \tilde{a}_n)] + \Delta(\tilde{\theta}_n + \tilde{a}_n) + C(\tilde{a}_n) \\
&\quad + \Phi[\sqrt{\alpha}(\underline{x}_n - \tilde{\theta}_n - \tilde{a}_n)] - \Phi[\sqrt{\alpha}(\underline{x}_n - \tilde{\theta}_n - \tilde{a}_n)] \\
&= \tilde{\theta}_n + \Phi[\sqrt{\alpha}(\underline{x}_n - \Theta(\underline{x}_n) - \tilde{a}_n)] - \Phi[\sqrt{\alpha}(\underline{x}_n - \tilde{\theta}_n - \tilde{a}_n)] \\
&\leq \tilde{\theta}_n
\end{aligned}$$

where the last inequality follows because the hypothesis $\tilde{\theta}_n \leq \Theta(\underline{x}_n)$ implies $\Phi[\sqrt{\alpha}(\underline{x}_n - \tilde{\theta}_n - \tilde{a}_n)] \geq \Phi[\sqrt{\alpha}(\underline{x}_n - \Theta(\underline{x}_n) - \tilde{a}_n)]$. This is a contradiction, and so $\tilde{\theta}_n > \Theta(\underline{x}_n)$. Therefore, the regime is overthrown for at least all $\theta < \Theta(\underline{x}_n)$. A parallel argument shows that if it is strictly dominant for $s(x_i) = 0$ for all $x_i > \overline{x}_n$, then the regime is not overthrown for at least all $\theta > \overline{\theta}_n := \Theta(\overline{x}_n)$.

(b) Since cumulative distribution functions are non-decreasing, for any beliefs of the citizens, the posterior probability assigned by a citizen with signal $x_i$ to the regime's overthrow is at least as much as the probability they assign to $\theta < \Theta(\underline{x}_n)$. Equivalently, $K(x_i, \underline{x}_n) - p$ is the most conservative estimate of the expected gain to subverting. From Lemma 2 and the intermediate value theorem, there is a unique $\underline{x}_{n+1} = \psi(\underline{x}_n)$ solving $K(\psi(\underline{x}_n), \underline{x}_n) = p$ such that if it is strictly dominant for $s(x_i) = 1$ for all $x_i < \underline{x}_n$, then it is strictly dominant for $s(x_i) = 1$ for all $x_i < \underline{x}_n$. Similarly, there is a unique $\overline{x}_{n+1} = \psi(\overline{x}_n)$ solving $K(\psi(\overline{x}_n), \overline{x}_n) = p$ such that if it is strictly dominant for $s(x_i) = 0$ for all $x_i > \overline{x}_n$, then it is strictly dominant for $s(x_i) = 0$ for all $x_i > \overline{x}_{n+1}$. Applying the proof of part (a) at each step then completes the argument.    ∥

**Concluding that there is no other equilibrium.**    Let $\underline{x}_0 := \underline{x}$ and $\overline{x}_0 := \overline{x}$ and generate sequences $\{\underline{x}_n\}_{n=0}^{\infty}$ from $\underline{x}_{n+1} = \psi(\underline{x}_n)$ and $\{\overline{x}_n\}_{n=0}^{\infty}$ from $\overline{x}_{n+1} = \psi(\overline{x}_n)$ where

$$K(\psi(\overline{x}_n), \overline{x}_n) = p \tag{A.30}$$

and

$$K(\psi(\underline{x}_n), \underline{x}_n) = p. \tag{A.31}$$

Part (a) of Lemma 5 maps the sequences of citizen thresholds $\{\underline{x}_n\}_{n=0}^{\infty}$ and $\{\overline{x}_n\}_{n=0}^{\infty}$ into *monotone* sequences of regime thresholds, $\{\underline{\theta}_n\}_{n=0}^{\infty}$ from $\underline{\theta}_n := \Theta(\underline{x}_n)$ and $\{\overline{\theta}_n\}_{n=0}^{\infty}$ from $\overline{\theta}_n := \Theta(\overline{x}_n)$. Moreover, by Lemma 3 the function $\psi(\cdot)$ generating the sequences $x_{n+1} = \psi(x_n)$ is continuous, has a unique fixed point $x^* = \psi(x^*)$ with derivative $\psi'(x^*) \in (0, 1)$ at this fixed point and upper bound $\psi(\underline{x}_n) \leq x^*$ for all $\underline{x}_n < x^*$. From below, the sequence $\{\underline{x}_n\}_{n=0}^{\infty}$ is bounded above, strictly monotone increasing and so converges $\underline{x}_n \nearrow x^*$ as $n \to \infty$. Similarly the sequence $\{\underline{\theta}_n\}_{n=0}^{\infty}$ is bounded above, strictly monotone increasing and so converges $\underline{\theta}_n \nearrow \theta^* =: \Theta(x^*)$ as $n \to \infty$. From above, symmetrically, the sequence $\{\overline{x}_n\}_{n=0}^{\infty}$ is bounded below, strictly monotone decreasing and so converges $\overline{x}_n \searrow x^*$ as $n \to \infty$. Similarly the sequence $\{\overline{\theta}_n\}_{n=0}^{\infty}$ is bounded below, strictly monotone decreasing and so converges $\overline{\theta}_n \searrow \theta^* =: \Theta(x^*)$ as $n \to \infty$. After a finite $n$ iterations, the only candidates for a citizen's equilibrium strategy all have $s(x_i) = 1$ for $x_i < \underline{x}_n$ and $s(x_i) = 0$ for $x_i > \overline{x}_n$ with $s(x_i)$ arbitrary for $x_i \in [\underline{x}_n, \overline{x}_n]$. Similarly, the only candidate for the regime's strategy has the regime abandoning for all $\theta < \underline{\theta}_n$, not abandoning for $\theta \geq \underline{\theta}_n$ with arbitrary choices for $\theta \in [\underline{\theta}_n, \overline{\theta}_n]$. At each iteration, these regime thresholds are implicitly determined by hidden actions $\underline{a}_n(\theta) := A(\theta - \underline{x}_n)$ and $\overline{a}_n(\theta) := A(\theta - \overline{x}_n)$ respectively. In the limit as $n \to \infty$, the only strategy that survives the elimination of strictly dominated strategies is the one with $s(x_i) = 1$ for $x_i < x^*$ and $s(x_i) = 0$ otherwise for citizens, with the regime abandoning for $\theta < \theta^* = \Theta(x^*)$ and hidden actions given by $a(\theta) = A(\theta - x^*)$. Therefore the only equilibrium is the unique monotone equilibrium.    ∥

## A.3.    *Proof of Proposition 2*

Since the regime benefit function $B(\theta, S)$ satisfies the Spence-Mirrlees sorting condition we have, for any $\theta \geq \theta'$ and any $S \geq S'$, that

$$B(\theta, S') - B(\theta, S) \geq B(\theta', S') - B(\theta', S). \tag{A.32}$$

That is, stronger regimes benefit at least weakly more from a smaller aggregate attack. The proof that stronger regimes choose higher apparent strengths is by contradiction.[8] Suppose that regime $\theta_H \geq \theta_L$ chooses apparent strength $y_H = y(\theta_H) < y_L = y(\theta_L)$. Then the weaker regime must be paying a higher cost

$$C(y_L - \theta_L) > C(y_H - \theta_L) \geq C(y_H - \theta_H). \tag{A.33}$$

Moreover, since $y_L$ is optimal for $\theta_L$

$$B(\theta_L, S(y_L)) - B(\theta_L, S(y_H)) \geq C(y_L - \theta_L) - C(y_H - \theta_L) > 0. \tag{A.34}$$

That is, since the cost of choosing $y_L$ is greater, the benefit must be greater too. But then since $B(\theta, S)$ is strictly decreasing in $S$ it must be that $S_L = S(y_L) < S_H = S(y_H)$. In short, a higher apparent strength is chosen only if it induces a smaller aggregate attack. But since $\theta_H \geq \theta_L$ and $S_H \geq S_L$, from the sorting condition (A.32), we then have

$$B(\theta_H, S_L) - B(\theta_H, S_H) \geq B(\theta_L, S_L) - B(\theta_L, S_H) > 0. \tag{A.35}$$

Since higher types can achieve a desired apparent strength at lower cost, *i.e.* $C(y_L - \theta_H) < C(y_H - \theta_H)$, (A.35) can only be true if

$$B(\theta_H, S_L) - C(y_H - \theta_H) > B(\theta_H, S_H) - C(y_H - \theta_H). \tag{A.36}$$

But this contradicts the optimality of $y_H$ for $\theta_H$. Hence $y_H \geq y_L$, and hence $y(\theta)$ is increasing in $\theta$.  ‖

## A.4. *Proofs of Proposition 3, Proposition 5, and Proposition 6*

*Preliminaries.* The proofs of Propositions 3, 5, and 6 are similar. To begin, substitute the regime indifference condition (20) into the citizen indifference condition (19) to obtain

$$\Phi[\sqrt{\alpha}(\theta^* - x^*)] = \frac{p}{1-p}\theta^* \quad \Leftrightarrow \quad \theta^* - x^* = \frac{1}{\sqrt{\alpha}}\Phi^{-1}\left(\frac{p}{1-p}\theta^*\right). \tag{A.37}$$

And now substitute this expression for the difference $\theta^* - x^*$ back into the regime indifference condition (20) to get a single equation characterizing the critical regime threshold $\theta^*$, namely

$$\theta^* + \frac{c}{\sqrt{\alpha}}\Phi^{-1}\left(\frac{p}{1-p}\theta^*\right) = \sqrt{\alpha}\gamma\phi(\sqrt{\alpha}\gamma) + \Phi(-\sqrt{\alpha}\gamma) \tag{A.38}$$

(using the fact that $\gamma$ is implicitly defined by $c = \sqrt{\alpha}\phi(\sqrt{\alpha}\gamma)$). Now define the composite parameter $z := c/\sqrt{\alpha} > 0$ and observe that in terms of this parameter

$$\sqrt{\alpha}\gamma = \sqrt{2\log\left(\frac{\phi(0)}{z}\right)} =: \delta(z)$$

so that (A.38) can be written in terms of the composite parameter $z$ alone, namely

$$T(\theta^*) := \theta^* + z\Phi^{-1}\left(\frac{p}{1-p}\theta^*\right) = \delta(z)\phi(\delta(z)) + \Phi(-\delta(z)). \tag{A.39}$$

All of the results in Propositions 3, 5, and 6 follow straightforwardly from the comparative statics of $\theta^*$ with respect to $z$. Implicitly differentiating with respect to $z$ gives

$$T'(\theta^*)\frac{\partial\theta^*}{\partial z} + \Phi^{-1}\left(\frac{p}{1-p}\theta^*\right) = \delta(z)\phi'(\delta(z))\delta'(z) + \delta'(z)\phi(z) - \phi(-\delta(z))\delta'(z). \tag{A.40}$$

The right-hand side can be simplified using symmetry, $\phi(-\delta(z)) = \phi(\delta(z))$, and the property $\phi(\delta(z)) = z$ so that $\phi'(\delta(z))\delta'(z) = 1$. Thus

$$T'(\theta^*)\frac{\partial\theta^*}{\partial z} + \Phi^{-1}\left(\frac{p}{1-p}\theta^*\right) = \delta(z). \tag{A.41}$$

Then because $T'(\theta) > 0$ for all $\theta$ we have that

$$\frac{\partial}{\partial z}\theta^* > 0 \quad \Leftrightarrow \quad \theta^* < \theta_{\text{crit}} := \frac{1-p}{p}\Phi(\delta(z)). \tag{A.42}$$

And because $T'(\theta) > 0$ for all $\theta$ we have $\theta^* < \theta_{\text{crit}}$ if and only if $T(\theta^*) < T(\theta_{\text{crit}})$. Applying $T(\cdot)$ to both sides of equation (A.42) and simplifying we have that the regime threshold $\theta^*$ is increasing in $z$ if and only if

$$p < \Phi(\delta(z)). \tag{A.43}$$

Since $\delta(z) > 0$ and $\Phi^{-1}(p) < 0$ for any $p < 1/2$, this condition is necessarily satisfied if $p < 1/2$. Using the definition of $\delta(z)$ and rearranging then gives

$$\frac{\partial}{\partial z}\theta^* > 0 \quad \Leftrightarrow \quad z < z^* := \phi(0)\exp\left(-\frac{1}{2}\max\left[0, \Phi^{-1}(p)\right]^2\right). \tag{A.44}$$

This condition is the key to the proofs of each of Propositions 3, 5, and 6 below.

8. I thank an anonymous referee for suggesting this proof.

*Proof of Proposition 3.*    If $\alpha \leq \underline{\alpha}(c) := (c/\phi(0))^2$, any regime is at a corner solution and has hidden actions $a(\theta) = 0$. In this case, the regime threshold is the same as in the Morris–Shin benchmark economy, $\theta^* = 1 - p$ for all $\alpha \leq \underline{\alpha}$. Otherwise, for interior solutions, the comparative statics are found by using the definition of the composite parameter $z = c/\sqrt{\alpha}$ and the chain rule

$$\frac{\partial}{\partial \alpha} \theta^* = -\frac{1}{2} \frac{c}{\alpha \sqrt{\alpha}} \frac{\partial}{\partial z} \theta^* . \tag{A.45}$$

Hence

$$\frac{\partial}{\partial \alpha} \theta^* < 0 \quad \Leftrightarrow \quad \frac{\partial}{\partial z} \theta^* > 0 \quad \Leftrightarrow \quad z < z^* . \tag{A.46}$$

And then on using the definition of $z$ and rearranging

$$\frac{\partial}{\partial \alpha} \theta^* < 0 \quad \Leftrightarrow \quad \alpha > \alpha^*(c, p) := \underline{\alpha}(c) \exp\left( \max\left[0, \Phi^{-1}(p)\right]^2 \right). \tag{A.47}$$

Now to establish that $\lim_{\alpha \to \infty} \theta^* = 0$, observe that for any $w \in \mathbb{R}$ the cumulative density $\Phi(\sqrt{\alpha} w) \to \mathbb{1}\{w > 0\}$ as $\alpha \to \infty$, *i.e.* to the indicator function that equals one if $w > 0$ and zero otherwise. Moreover, as $\alpha \to \infty$ the parameter $\gamma = \sqrt{\log(\alpha/\underline{\alpha})/\alpha} \to 0$ and $\Phi(-\sqrt{\alpha} \gamma) \to 0$. Applying these to (19) we see that for large $\alpha$ solutions to the citizen's indifference condition are approximately the same as solutions to

$$\mathbb{1}\{\theta^* - x^* > 0\} = -\frac{p}{1-p}(\theta^* - x^*)c. \tag{A.48}$$

The only solution to (A.48) is $\theta^* - x^* = 0$. So as $\alpha \to \infty$, solutions to (19) approach zero too. Then from the regime's indifference condition (20), if $\theta^* - x^* \to 0$ it must also be the case that $\theta^* \to 0$ as claimed.    $\|$

*Proof of Proposition 5.*    Similarly from the chain rule

$$\frac{\partial}{\partial c} \theta^* = \frac{1}{\sqrt{\alpha}} \frac{\partial}{\partial z} \theta^* . \tag{A.49}$$

Hence eliminating the derivative with respect to $z$ between (A.45) and (A.49) we have

$$\frac{\partial}{\partial c} \theta^* = -2 \frac{\alpha}{c} \frac{\partial}{\partial z} \theta^* . \tag{A.50}$$

Multiplying both sides by $c/\theta^* > 0$ then gives the stated relationship between elasticities (27). Thus the cost effect is twice as large in magnitude as the signal precision effect. Observe that for corner solutions, $\alpha \leq \underline{\alpha}(c) := (c/\phi(0))^2$, both these effects are zero. Otherwise, if $\alpha > \underline{\alpha}(c)$ the signs of the two effects are determined by whether $\alpha$ is larger or smaller than $\alpha^*(c, p)$ but are always opposite to each other.    $\|$

*Proof of Proposition 6.*    Now let the composite parameter be

$$z = \frac{c(n)}{\sqrt{n\hat{\alpha}}}.$$

Following calculations identical to (A.37)–(A.44) above, the equilibrium threshold $\theta^*$ depends on $n$ and $\hat{\alpha}$ only through this composite parameter with the comparative statics with respect to $\hat{\alpha}$ being obtained exactly as in Proposition 3 above. Then using the chain rule

$$\frac{\partial \theta^*}{\partial n} = \frac{\partial \theta^*}{\partial z} \frac{\partial z}{\partial n}, \quad \text{and} \quad \frac{\partial \theta^*}{\partial \hat{\alpha}} = \frac{\partial \theta^*}{\partial z} \frac{\partial z}{\partial \hat{\alpha}}. \tag{A.51}$$

Eliminating the derivative with respect to $z$ then gives

$$\frac{\partial \theta^*}{\partial n} = \frac{\partial z}{\partial n} \left( \frac{\partial z}{\partial \hat{\alpha}} \right)^{-1} \frac{\partial \theta^*}{\partial \hat{\alpha}}. \tag{A.52}$$

Calculating the derivative of $z$ with respect to $\hat{\alpha}$ and then multiplying both sides by $n/\theta^* > 0$ then gives the relationship between elasticities

$$\frac{\partial \log \theta^*}{\partial \log n} = -2 \frac{\partial \log z}{\partial \log n} \frac{\partial \log \theta^*}{\partial \log \hat{\alpha}}. \tag{A.53}$$

And then using

$$\frac{\partial \log z}{\partial \log n} = \frac{c'(n)n}{c(n)} - \frac{1}{2}$$

gives the stated relationship between elasticities, equation (29).    $\|$

### A.5.   *Proofs of Proposition 4 and Proposition 8*

*Preliminaries.*    These results concern asymptotics with respect to the signal precision $\alpha$. The overall proof strategy is similar in that in both cases I adopt a guess-and-verify method, guessing the asymptotic behaviour of the equilibrium and then verifying that the relevant equilibrium conditions are satisfied. Proposition 4 holds for the benchmark model for which we know, from Proposition 1 above, that the equilibrium is unique. Proposition 8 refers to an extension of the benchmark model for which it is not known if the equilibrium is unique. Proposition 8 applies to any monotone equilibrium of this extended model but leaves open the asymptotics of other equilibria (if any).

### *Proof of Proposition 4.*

**Equilibrium conditions.**    For a general convex cost function $C(a)$, the key equilibrium conditions are the marginal citizen's indifference condition, which can be written

$$(1-p)\Phi(\sqrt{\alpha}(\theta^* - x^*)) = p\int_{\theta^*}^{\infty} \sqrt{\alpha}\phi(\sqrt{\alpha}(x^* - \theta - a(\theta)))d\theta \tag{A.54}$$

and the regime's indifference condition

$$\theta^* = \Phi(\sqrt{\alpha}(x^* - \theta^* - a(\theta^*))) + C(a(\theta^*)) \tag{A.55}$$

with the regime's hidden actions $a(\theta)$ satisfying the first-order condition

$$C'(a(\theta)) = \sqrt{\alpha}\phi(\sqrt{\alpha}(x^* - \theta - a(\theta))), \qquad \theta \geq \theta^* . \tag{A.56}$$

**Guess-and-verify.**    We now guess that

$$\lim_{\alpha\to\infty} \sqrt{\alpha}(x^* - \theta^* - a(\theta^*)) = \lim_{\alpha\to\infty} \sqrt{\alpha}(\theta^* - x^*) = -\infty. \tag{A.57}$$

For notational simplicity I suppress the dependence of the equilibrium thresholds and hidden actions on the signal precision $\alpha$. Note that the equilibrium hidden actions $a(\theta)$ depend on $\alpha$ both directly (via the auxiliary function $A(\cdot)$, as defined in (8) in the main text) and indirectly via the citizen threshold $x^*$. To verify this guess, first observe that if the first limit in (A.57) holds, then from (A.55) we have $\theta^* = C(a(\theta^*))$. Similarly, if the second limit in (A.57) holds, then $\Phi(\sqrt{\alpha}(\theta^* - x^*)) \to 0$ and the value of the integral on the right-hand side of (A.54) has to converge to zero in order for this guess to work. Plugging in the first-order condition (A.56) this requires

$$\lim_{\alpha\to\infty}\int_{\theta^*}^{\infty} C'(a(\theta))d\theta = 0. \tag{A.58}$$

Since $\theta^* \in [0,1]$ for any $\alpha$ and the composite function $C'(a(\theta))$ is non-negative and uniformly continuous in $\alpha$ for any $\theta$, this can only be true if the hidden actions $a(\theta) \to 0^+$ for all $\theta \geq \theta^*$. But if $a(\theta^*) \to 0^+$ for the marginal regime then $C(a(\theta^*)) \to 0^+$ and so $\theta^* \to 0^+$ too. Hence we have the result that for high enough signal precision, even the most fragile regime can survive.

Now for the limit as $\alpha \to 0^+$. Recall that for this part we assume strictly convex costs *and* that $C'(0) = 0$. Let $\alpha \to 0^+$ and guess that $\sqrt{\alpha}x^* \to \infty$ holds. Then $x^* \to \infty$. Since $\theta^* \in [0,1]$, we have $\sqrt{\alpha}(x^* - \theta^*) \to \infty$ and the integral on the right-hand side of (A.54) needs to converge to zero. Hence, by (A.56), $a(\theta) \to 0^+$ for all $\theta \geq \theta^*$ [Note: the strict convexity plus $C'(0) = 0$ is used here so that (A.56) continuously holds even as $\alpha \to 0^+$; with a jump in the marginal cost at $a = 0$, this would not generally be true]. But if $a(\theta^*) \to 0^+$, $\theta^* \in [0,1]$, and $\sqrt{\alpha}x^* \to \infty$, then from the indifference condition (A.55) the regime threshold $\theta^* \to 1^-$.    ‖

### *Proof of Proposition 8.*

**Equilibrium conditions when there is an opposition.**    Consider a monotone equilibrium with thresholds $x^*, \theta^*$ to be determined. The aggregate attack facing the regime is $S(\theta, a, e) = \Phi(\sqrt{\alpha}(x^* - \theta - a + e))$ and the first order necessary condition characterizing the regime's hidden action can be written

$$C'(a) = \sqrt{\alpha}\phi(\sqrt{\alpha}(x^* - \theta - a + e(\theta))), \qquad \theta \geq \theta^* \tag{A.59}$$

with $a(\theta) = 0$ for all $\theta < \theta^*$. Similarly, for the opposition

$$C'(\kappa e) = \sqrt{\alpha}\phi(\sqrt{\alpha}(x^* - \theta - a(\theta) + e)), \qquad \text{all } \theta. \tag{A.60}$$

Combining equations (A.59) and (A.60) and using the strict convexity of $C(\cdot)$ we have that for all $\theta \geq \theta^*$ where the regime survives

$$a(\theta) = \kappa e(\theta), \qquad \theta \geq \theta^* . \tag{A.61}$$

Thus if the regime survives the regime's actions are larger than those of the opposition if and only if $\kappa > 1$, *i.e.* when the regime has a cost advantage. Plugging in $e(\theta) = a(\theta)/\kappa$ into (A.59) gives

$$C'(a) = \sqrt{\alpha}\phi\left(\sqrt{\alpha}\left(x^* - \theta - \frac{\kappa - 1}{\kappa}a\right)\right), \qquad \theta \geq \theta^* \tag{A.62}$$

which can be implicitly solved for $a(\theta)$ for all $\theta \geq \theta^*$. For $\theta < \theta^*$ we have $a(\theta) = 0$ and we can solve for $e(\theta)$ from

$$C'(\kappa e) = \sqrt{\alpha}\phi(\sqrt{\alpha}(x^* - \theta + e)), \qquad \theta < \theta^* \tag{A.63}$$

Conditional on $x^*, \theta^*$ equations (A.61), (A.62), and (A.63) characterize the functions $a(\theta), e(\theta)$. The thresholds $x^*, \theta^*$ are then determined by the regime's indifference condition, which can now be written

$$\theta^* = \Phi\left(\sqrt{\alpha}\left(x^* - \theta^* - \frac{\kappa - 1}{\kappa}a(\theta^*)\right)\right) + C(a(\theta^*)) \tag{A.64}$$

and the citizen's indifference condition, which can likewise be written

$$(1 - p)\int_{-\infty}^{\theta^*}\sqrt{\alpha}\phi(\sqrt{\alpha}(x^* - \theta + e(\theta)))d\theta = p\int_{\theta^*}^{\infty}\sqrt{\alpha}\phi\left(\sqrt{\alpha}\left(x^* - \theta - \frac{\kappa - 1}{\kappa}a(\theta)\right)\right)d\theta \tag{A.65}$$

(analogous to (A.54) above but recognizing the opposition's action $e(\theta)$ on $\theta < \theta^*$ and $a(\theta) = \kappa e(\theta)$ on $\theta \geq \theta^*$).

**Guess-and-verify.**   Now guess that

$$\lim_{\alpha \to \infty}\sqrt{\alpha}a(\theta) = \lim_{\alpha \to \infty}\sqrt{\alpha}e(\theta) = 0 \qquad \text{and} \qquad \lim_{\alpha \to \infty}\sqrt{\alpha}(\theta^* - x^*) = -\infty. \tag{A.66}$$

To verify that this is an equilibrium, begin with the regime's indifference condition and observe that since $a(\theta) \to 0$ as $\alpha \to \infty$ for any $\theta$ we know $C(a(\theta^*)) \to 0$ and so by the last limit in (A.66) we have $\Phi(\sqrt{\alpha}(\theta^* - x^*)) \to 0^+$. Therefore from (A.64) the regime threshold $\theta^* \to 0^+$ too.

We now need to verify that this is also consistent with optimal behaviour by the citizens and by the opposition. Using the regime's and opposition's first-order conditions we can rewrite the citizen indifference condition (A.65) as

$$(1 - p)\int_{-\infty}^{\theta^*}C'(\kappa e(\theta))d\theta = p\int_{\theta^*}^{\infty}C'(a(\theta))d\theta.$$

Now under the conjectured equilibrium $e(\theta) \to 0$ and $a(\theta) \to 0$ for all $\theta$ so the marginal costs $C'(\kappa e(\theta)) \to 0$ and $C'(a(\theta)) \to 0$ for all $\theta$ too. Then using $\theta^* \in [0, 1]$ for all $\alpha$ and that both integrands are non-negative and uniformly continuous in $\alpha$ for each $\theta$ we can pass to the limit as $\alpha \to \infty$ to conclude that both integrals evaluate to zero also. So the citizen indifference condition holds in the conjectured equilibrium.

Finally, observe that optimal behaviour by the opposition requires the first-order condition $C'(\kappa e) = \sqrt{\alpha}\phi(\sqrt{\alpha}(x^* - \theta + e))$ to hold for all $\theta < \theta^*$. But if the thresholds $x^*, \theta^*$ for the citizens and the regime satisfy the limit $\sqrt{\alpha}(\theta^* - x^*) \to -\infty$ as conjectured in (A.66) then indeed the opposition's best response is for $e(\theta) \to 0$ for all $\theta < \theta^*$ [Note: we already know that $e(\theta) = a(\theta)/\kappa \to 0$ on $\theta \geq \theta^*$]. Hence this guess works, *i.e.* the conjectured equilibrium as $\alpha \to \infty$ is consistent with both the regime and citizen indifference conditions and with optimal behaviour by the opposition. Thus as $\alpha \to \infty$ we have $\theta^* \to 0^+$ and hence for signal precision $\alpha$ large enough we have regime threshold $\theta^*$ strictly less than the Morris–Shin benchmark $\theta_{\text{MS}}^* = 1 - p$ despite the oppostion.   $\|$

## A.6.   *Proof of Proposition 7 and related details*

Consider a monotone equilibrium with thresholds $x^*, \theta^*$ to be determined. The aggregate attack facing the regime is $S(\theta, a) = \Phi(\sqrt{n(a)}(x^* - \theta))$ and the first-order necessary condition characterizing the regime's hidden action can be written

$$C'(a) = (\theta - x^*)\phi\left(\sqrt{n(a)}(x^* - \theta)\right)\frac{\partial}{\partial a}\sqrt{n(a)}, \qquad \theta \geq \theta^*. \tag{A.67}$$

As usual, for $\theta < \theta^*$ we have $a(\theta) = 0$ before jumping discretely at $\theta^*$. For $\theta \geq \theta^*$ the hidden actions solve (A.67) and, as in the benchmark model, are determined by the difference $\theta - x^*$. Now however the *sign* of $a(\theta)$ is the sign of the difference $\theta - x^*$ (since the cost function $C(a)$ is symmetric about zero, the normal density is strictly positive, and $n(a)$ is strictly increasing in $a$). In particular, $a(\theta) < 0$ for $\theta < x^*$ and $a(\theta) > 0$ for $\theta > x^*$. Since $n(0) =: \alpha$ and $n(a)$ is strictly increasing, we then have $n(a(\theta)) < \alpha$ for $\theta \in [\theta^*, x^*)$ and $n(a(\theta)) > \alpha$ for $\theta > x^*$. If $x^* < \theta^*$ then $[\theta^*, x^*)$ is empty and we simply have $n(a(\theta)) > \alpha$ for all $\theta > \theta^*$. The knife-edge regime with type just equal to the citizen threshold $\theta = x^*$ chooses $a(x^*) = 0$ giving the neutral outcome $n(a(x^*)) = \alpha$.   $\|$

**Equilibrium conditions when regime controls number of signals.**    To solve this version of the model, observe that conditional on $x^*$, the first-order condition (A.67) characterizes the function $a(\theta)$ and hence $n(a(\theta))$. The thresholds $x^*, \theta^*$ are then determined by the regime's indifference condition, which in this version of the model can be written

$$\theta^* = \Phi\left(\sqrt{n(a(\theta^*))}(x^* - \theta^*)\right) + C(a(\theta^*)) \tag{A.68}$$

and the citizen's indifference condition, which can likewise be written

$$(1-p)\Phi(\sqrt{\alpha}(\theta^* - x^*)) = p \int_{\theta^*}^{\infty} \sqrt{n(a(\theta))}\phi\left(\sqrt{n(a(\theta))}(x^* - \theta)\right) d\theta \tag{A.69}$$

(analogous to (A.54) and using $n(a(\theta)) = n(0) =: \alpha$ for $\theta < \theta^*$). These two conditions are straightforward to solve numerically for any given parameters $\alpha, p$ and cost function $C(a)$.

# B.   FURTHER DISCUSSION

## B.1.   *Can we transform this to a standard global game?*

The signal-jamming monotonicity result Proposition 2 suggests an alternative approach to establishing equilibrium uniqueness. Namely, given monotonicity of $y(\theta)$, perhaps we can transform this game by taking the relevant notion of the "fundamental" to be the outcome $y$ (rather than $\theta$) so that the regime is overthrown if its $y$ is too low, with standard global games arguments then being invoked to obtain uniqueness. Unfortunately, there is a key difficulty with this approach. To ensure we are playing the same game, the citizens' prior over $y$ in the transformed game has to be consistent with their beliefs about $y(\theta)$ in the original game. We cannot simply take the prior over $y$ to be the improper uniform, for instance, since that would be completely inconsistent with the regime's manipulation. In fact, to be consistent with the regime's manipulation, the prior over $y$ will have (at least one) *gap*. Because $y(\theta)$ typically jumps from $y(\theta) = \theta$ for $\theta < \theta^*$ to some $y^* = y(\theta^*) > \theta^*$ at the critical threshold, the prior will have to give exactly zero probability density to the interval $(\theta^*, y^*)$. Standard global games arguments crucially assume, however, that the prior has strictly positive density over its support (*e.g.* Carlsson and van Damme (1993), Assumption A2) and that property would not hold in the transformed game. Worse yet, if, to take what will be a leading example, the function $y(\theta)$ is *constant* on an interval of $\theta$, then this will put an *atom* in the prior over $y$ (*cf.*, Edmond, 2013). Standard global games arguments generally require that the prior be sufficiently *diffuse* (Morris and Shin, 2000; Hellwig, 2002; Morris and Shin, 2003) and that property is harder to obtain if the prior has an atom. In short, we cannot simply transform the game and invoke standard global games arguments.

Proposition 1 shows that this game *does* have a unique equilibrium and that this equilibrium has many of the familiar global games properties. Rather than attempting to transform this game to a purportedly equivalent standard global game, the proof uses the usual global games *methodology* of establishing dominance regions and iteratively eliminating (interim) strictly dominated strategies. In employing this methodology, I make extensive use of the assumption that the benefit $B(\theta, S)$ is separable in $\theta$ and $S$, a standard payoff assumption in the global games literature. In my model, separability also allows the actions to be written $a(\theta, \hat{x}) = A(\theta - \hat{x})$ in terms of the exogenous function defined in (8). This greatly simplifies the equilibrium fixed point problem. Without separability, knowing $\hat{x}$ is not sufficient to determine the whole function.

## B.2.   *Role of coordination*

This appendix highlights the role of imperfect coordination in enabling the regime to survive even when signals are precise. Suppose to the contrary that citizens are perfectly coordinated and receive one $x = \theta + a + \varepsilon$. Collectively, they can overthrow the regime if $\theta < 1$. In a monotone equilibrium the mass attacks the regime, $S(x) = 1$, if and only if $x < x^*$ where $x^*$ solves $\text{Prob}[\theta < 1 | x^*] = p$.

The regime now faces aggregate uncertainty. It does not know what value of $x$ will be realized. The regime chooses its hidden action to maximize its expected payoff

$$a(\theta) \in \underset{a \geq 0}{\operatorname{argmax}} \left[ -C(a) + \int_{-\infty}^{\infty} \max[0, \theta - S(x)] \sqrt{\alpha}\phi(\sqrt{\alpha}(x - \theta - a)) dx \right]. \tag{B.1}$$

In a monotone equilibrium, the regime's objective simplifies to

$$-C(a) - \min[\theta, 1]\Phi(\sqrt{\alpha}(x^* - \theta - a)). \tag{B.2}$$

Regimes with $\theta < 0$ are overthrown and so never engage in costly manipulation.

**Example: strictly convex costs.**    Suppose that costs are strictly convex, $C''(a) > 0$. This implies all regimes $\theta > 0$ will choose some positive manipulation $a(\theta) > 0$, even regimes that are overthrown ex post. The key first order necessary condition for the regime's choice of action $a(\theta)$ is

$$\min[\theta, 1]\sqrt{\alpha}\phi(\sqrt{\alpha}(x^* - \theta - a)) = C'(a), \qquad \theta \geq 0 \tag{B.3}$$

As usual, there may be two solutions to this first-order condition; if so, the smaller is eliminated by the second-order condition. An equilibrium of this game is constructed by simultaneously determining $a(\theta)$ and the $x^*$ that solves $\text{Prob}[\theta < 1 | x^*] = p$.

The first-order condition implies that taking as given $x^*$ the regime's $a(\theta) \to 0$ as $\alpha \to \infty$. Given this, the probability of overthrowing the regime $\text{Prob}[\theta < 1 | x] \to \mathbb{1}\{x < 1\}$ as $\alpha \to \infty$. This implies $x^* \to 1$. With arbitrarily precise information, the regime takes no action and so $x$ is very close to $\theta$. The mass attacks only if it believes $\theta < 1$ and since $x$ is close to $\theta$ attacks only if $x < 1$. So if citizens are perfectly coordinated then for precise information regime change occurs for all $\theta < 1$. In contrast, if citizens are imperfectly coordinated then for precise information all regimes $\theta \geq 0$ survive.

Angeletos *et al.* (2006) provide a related analysis. In their model, if agents are imperfectly coordinated then for precise information $\theta^*$ can be any $\theta \in (0, \theta^*_{\text{MS}}]$ where $\theta^*_{\text{MS}} = 1 - p < 1$. But if agents are perfectly coordinated then for precise information regime change occurs for all $\theta < 1$. Thus, when information is precise the two models agree about the regime change outcome when agents are perfectly coordinated but come to different conclusions when agents are imperfectly coordinated.

## Supplementary Data

Supplementary data are available at *Review of Economic Studies* online.

## REFERENCES

ANGELETOS, G.-M. and WERNING, I. (2006), "Crises and Prices: Information Aggregation, Multiplicity and Volatility", *American Economic Review*, **96**, 1720–1736.

ANGELETOS, G.-M., HELLWIG, C. and PAVAN, A. (2006), "Signaling in a Global Game: Coordination and Policy Traps", *Journal of Political Economy*, **114**, 452–484.

ARENDT, H. (1973), *The Origins of Totalitarianism*, revised edn (London: André Deutsch).

BESLEY, T. and PRAT, A. (2006), "Handcuffs for the Grabbing Hand? The Role of the Media in Political Accountability", *American Economic Review*, **96**, 720–736.

BLUME, A., BOARD, O. J. and KAWAMURA, K. (2007), "Noisy Talk", *Theoretical Economics*, **2**, 395–440.

BOIX, C. and SVOLIK, M. (2013), "The Foundations of Limited Authoritarian Government: Institutions and Power-Sharing in Dictatorships", *Journal of Politics*, **75**, 300–316.

BUENO DE MESQUITA, E. (2010), "Regime Change and Revolutionary Entrepreneurs", *American Political Science Review*, **104**, 446–466.

CARLSSON, H. and VAN DAMME, E. (1993), "Global Games and Equilibrium Selection", *Econometrica*, **61**, 989–1018.

CHASE, M. S. and MULVENON, J. C. (2002), *You've Got Dissent: Chinese Dissident use of the Internet and Beijing's Counter-Strategies*. (RAND report MR-1543).

CHASSANG, S. and PADRO-I-MIQUEL, G. (2010), "Conflict and Deterrence under Strategic Risk", *Quarterly Journal of Economics*, **125**, 1821–1858.

CHWE, M. S.-Y. (2001), *Rational Ritual: Culture, Coordination, and Common Knowledge* (Princeton, NJ: Princeton University Press).

COHEN, N. (2009), "Twitter on the Barricades: Six Lessons Learned", *New York Times*.

CRAWFORD, V. P. and SOBEL, J. (1982), "Strategic Information Transmission", *Econometrica*, **50**, 1431–1451.

EDMOND, C. (2012), "Information Manipulation and Social Coordination", (University of Melbourne Working Paper).

EDMOND, C. (2013), "Non-Laplacian Beliefs in a Global Game with Noisy Signaling", (University of Melbourne Working Paper).

ESFANDIARI, G. (2010), "The Twitter Devolution", *Foreign Policy*.

FALLOWS, J. (2008), "The Connection Has Been Reset", *Atlantic Monthly*, **301**.

FRIEDRICH, C. J. and BRZEZINSKI, Z. K. (1965), *Totalitarian Dictatorship and Autocracy*, 2nd edn, (Cambridge, MA: Harvard University Press).

GEHLBACH, S. and SONIN, K. (2008), "Government Control of the Media", (University of Wisconsin Working Paper).

HELLWIG, C. (2002), "Public Information, Private Information, and the Multiplicity of Equilibria in Coordination Games", *Journal of Economic Theory*, **107**, 191–222.

HELLWIG, C., MUKHERJI, A. and TSYVINSKI, A. (2006), "Self-Fulfilling Currency Crises: The Role of Interest Rates", *American Economic Review*, **96**, 1769–1787.

HOLMSTRÖM, B. (1999), "Managerial Incentive Problems: A Dynamic Perspective", *Review of Economic Studies*, **66**, 169–182.

JOHNSON, J. P. and MYATT, D. P. (2006), "On the Simple Economics of Advertising, Marketing, and Product Design", *American Economic Review*, **96**, 756–784.

KALATHIL, S. and BOAS, T. C. (2003), *Open Networks, Closed Regimes: The Impact of the Internet on Authoritarian Rule*, (Washington, DC: Carnegie Endowment for International Peace).

KARTIK, N. (2009), "Strategic Communication with Lying Costs", *Review of Economic Studies*, **76**, 1359–1395.

KARTIK, N., OTTAVIANI, M. and SQUINTANI, F. (2007), "Credulity, Lies, and Costly Talk", *Journal of Economic Theory*, **134**, 93–116.

KIRKPATRICK, D. D. (2011), "Wired and Shrewd, Young Egyptians Guide Revolt", *New York Times*.

LEWIS, T. R. and SAPPINGTON, D. E. M. (1991), "All-or-Nothing Information Control", *Economics Letters*, **37**, 111–113.

LEWIS, T. R. and SAPPINGTON, D. E. M. (1994), "Supplying Information to Facilitate Price Discrimination", *International Economic Review*, **35**, 309–327.

MARINOVIC, I. (2011), "Internal Control System, Earnings Quality and the Dynamics of Financial Reporting", (Stanford GSB Working Paper).

MONDERER, D. and SAMET, D. (1989), "Approximating Common Knowledge with Common Beliefs", *Games and Economic Behavior*, **1**, 170–190.

MOROZOV, E. (2011), *The Net Delusion*. (London: Allen Lane).

MORRIS, S. and SHIN, H. S. (1998), "Unique Equilibrium in a Model of Self-Fulfilling Currency Attacks", *American Economic Review*, **88**, 587–597.

MORRIS, S. and SHIN, H. S. (2000), "Rethinking Multiple Equilibria in Macroeconomic Modeling", in Bernanke, B. S. and Rogoff, K. (eds) *NBER Macroeconomics Annual*, (MIT Press), pp. 139–161.

MORRIS, S. and SHIN, H. S. (2003), "Global Games: Theory and Applications", in Dewatripont, M., Hansen, L. P. and Turnovsky, S. J. (eds) *Advances in Economics and Econometrics: Theory and Applications*, (Cambridge University Press).

MUSGROVE, M. (2009), "Twitter Is a Player In Iran's Drama", *Washington Post*.

SHADMEHR, M. and BERNHARDT, D. (2011), "Collective Action with Uncertain Payoffs: Coordination, Public Signals and Punishment Dilemmas", *American Political Science Review*, **105**, 829–851.

SOLEY, L. C. (1987), *Clandestine Radio Broadcasting: A Study of Revolutionary and Counterrevolutionary Electronic Broadcasting* (New York, NY: Praeger).

ZEMAN, Z. A. B. (1973), *Nazi Propaganda*, 2nd edn, (London: Oxford University Press).