

Why Do Authoritarian Regimes Sign the Convention Against Torture?

Signaling, Domestic Politics and Non-Compliance*

James R. Hollyer[†]
Yale University

B. Peter Rosendorff[‡]
New York University

Current Version: June 2011

Abstract

Traditional international relations theory holds that states will join only those international institutions with which they generally intend to comply. Here we show when this claim might not hold. We construct a model of an authoritarian government's decision to sign the UN Convention Against Torture (CAT). Authoritarian governments use the signing of this treaty – followed by the willful violation of its provisions – as a costly signal to domestic opposition groups of their willingness to employ repressive tactics to remain in power. In equilibrium, authoritarian governments that torture heavily are more likely to sign the treaty than those that torture less. We further predict that signatory regimes survive longer in office than non-signatories, and enjoy less domestic opposition – and we provide empirical support for these predictions.

*We would like to thank James Vreeland and Jennifer Gandhi for their generosity in providing access to their data, and Leslie Johns and two anonymous reviewers for their detailed comments and suggestions. We would also like to thank Bruce Bueno de Mesquita, David Stasavage, Jon Eguia, Fernando Martel Garcia, Art Stein, Johannes Urpelainen, Joanne Gowa, the participants in seminars at Claremont, Columbia, Georgetown, NYU, UCLA, UCSD and USC; the 2009 MPSA Panel on International Human Rights Agreements, the 2009 APSA Panel on the Political Economy of International Regimes and the 3rd annual PEIO conference for helpful comments and suggestions. All remaining errors are our own.

[†]Corresponding author: Leitner Program in International and Comparative Political Economy, the MacMillan Center for International and Area Studies, Yale University. 115 Prospect Place, New Haven, CT 06520.

[‡]New York University, Department of Politics. 19 W. 4th St., 2nd Floor. New York, NY 10012.

Sovereign states that sign international treaties, we are told, intend to broadly comply with the obligations imposed by these treaties. The reasons for this claim are varied: International law (the Vienna Convention on the Law of Treaties in particular) declares that “every treaty ... is binding upon the parties.” This declaration follows from the basic principle of international law *pacta sunt servanda* – treaties are to be obeyed. Downs, Rocke & Barsoom (1996) establish that those countries that are most likely to abide by the rules promulgated by an international institution are also those countries that are most likely to join in the first place. Failure to comply with treaty provisions is described as a ‘managerial problem’ (Chayes & Chayes 1993), or as a temporary aberration to be remedied by (re)negotiation (Koremenos 2005). Tolerated temporary escape (Bagwell & Staiger 2005, Rosendorff & Milner 2001), exchanges of information, and dispute resolution mechanisms are designed to complete gaps in treaty language or to generate better information about signatories’ behavior (Rosendorff 2005) and to thus bring about treaty compliance. Where international treaties address issues of international externalities – such as trade, security, or the environment – the intent to comply is strengthened by the mutual gains associated with a predictable, stable, and cooperative international order. Similarly, states facing collective action problems may be inclined to forgo the temporary benefits of defection in order to remain within the society of cooperative nations, especially when future (relative to current) consumption is highly valued (Downs & Rocke 1995).

Recent scholarship has explored if these general findings also apply to the specific case of human rights treaties. Simmons (2009) argues that the major human rights treaties have been successful in reducing the prevalence of torture worldwide. She claims that countries accede to and ratify these treaties because they intend to comply with treaty provisions (p.42). She acknowledges that there are some states that ratify, but do not greatly adjust their behavior. She describes these states as the “false positives,” the countries that sign human rights treaties and continue to torture. And she observes that these states tend to have relatively authoritarian regimes.

Simmons’ observation reinforces the conclusions of Hathaway (2007), who finds a positive association between the practice of torture and the signing of human rights treaties by highly authoritarian regimes. Hafner-Burton & Tsutsui (2005, 2007) confirm that signing human rights

treaties has little or no effect on the behavior of the world's worst repressors. As they put it, there is a "rising gap between states' propensity to join the international human rights regime and to bring their human rights practice into compliance," and this gap brings the efficacy of international law into fundamental question.

We thus have a puzzle: if states join agreements because they intend to comply with them, why do some states, particularly authoritarian states, sign and fail to comply with human rights treaties?

We argue that authoritarian states sign human rights treaties explicitly because they do *not* intend to comply. And it is important to those signatories that all observers understand that they have no intention of complying at the time of signing. The logic, while counterintuitive, is straightforward: An elite facing threats from a domestic opposition can mitigate these threats by engaging in torture. Human rights agreements inflict penalties on the elites in signatory states for engaging in torture, in the event that these elites are removed from office.¹ The signing of a human rights treaty is, therefore, a signal to the opposition of the high value the elite places on holding onto power and its willingness to use torture if necessary. On observing the government's actions, the opposition – now better informed about the value the elite places on holding power – will rationally reduce its anti-regime activities. On the other hand a regime that doesn't sign shows itself to be vulnerable to the added costs associated with the use of torture. Thus, the opposition will increase efforts to remove the regime on seeing that the government does not sign.

This logic leads to two conclusions: First, more repressive regimes (regimes with elites more willing to use force to hold onto power) will sign more frequently than less (or non-) repressive governments. Second, opposition political action falls in signatory states, yielding reductions in the likelihood of regime collapse or transition. In the non-signatory states, opposition response actually rises, leading to more frequent regime failure.

The first finding is consistent with Hathaway's (2007) empirical results, and offers a theoretical explanation for the puzzle above. Authoritarian governments that torture more heavily in time t are more likely to sign the CAT in time $t + 1$. The magnitude of this difference is non-trivial – in our dataset, autocrats that go on to sign the CAT score roughly $\frac{1}{3}$ of a standard deviation higher on

¹Similar results follow if human rights agreements raise the marginal cost to engaging in torture. See the online appendix for this alternative variant of the model.

common torture measures than autocrats who never sign.

In order to check the veracity of the model, we test the second (novel) prediction: authoritarian regimes that sign the treaty will enjoy longer tenures in office than those that do not. This is true for three reasons: (1) A selection effect implies that those regimes that will fight most strongly to remain in power are the same regimes that sign the treaty; (2) An information effect implies that domestic opposition groups will engage in fewer activities designed to overthrow a signatory government; and (3) A commitment effect implies that regimes that sign the treaty will cling more tightly to power to avoid potential legal punishments after stepping down from office. We test this claim using data on the signing of the UN Convention Against Torture (CAT) and find that it enjoys robust empirical support. Signatory regimes face a lower hazard rate than observationally similar non-signatories across a wide variety of empirical specifications.

The model further predicts that opposition groups will reduce their effort to unseat signatory governments on witnessing the signing of the CAT. This claim follows from the *informational effect* of signing the CAT. Domestic opposition groups, on witnessing the government sign the CAT, conclude that it is a ‘strong’ type, likely to prevail in the contest for power. As a result, they reduce costly activities aimed at the government’s overthrow.

The CAT has contradictory effects on levels of torture. On the one hand, governments will be tempted to reduce their level of repression as oppositional efforts fall. This tendency follows from the *informational effect* of signing the CAT. However, signatories also face higher costs from relinquishing office after signing the CAT. Once removed from power, former government officials face the risk of potential prosecution and extradition for repressive acts. As a result, such leaders will be less likely to give up power and will employ higher levels of repression. This is the *commitment effect* of the CAT. The informational and commitment effects weigh against one another – and the precise direction of the net effect is sensitive to model assumptions.

We test these propositions in our empirical analysis below, and find suggestive evidence for a decline in opposition activity following the signing of the CAT. Our results also suggest that torture levels decline in signatory governments – and thus that the informational effect outweighs the commitment effect of the treaty.

Key to the causal logic of the argument is the notion that the CAT affects the costs to a member-state's elite of engaging in torture. We will argue that, aside from international opprobrium and withdrawals of concessions (or active sanctions) along other dimensions by the international community (Hafner-Burton 2005), signatories of the CAT must consider the role of "universal jurisdiction" and the extradition clauses of the CAT when determining whether or not to employ torture. These additional considerations serve to make torture more costly given accession to the CAT than not. However, these costs do not directly translate into higher levels of compliance by signatory states. Rather, they allow signing to act as a costly signaling mechanism, such that the states that sign are those that are most likely to defy their treaty obligations.

This paper makes contributions to three literatures. It first speaks to the literature on selection effects and international institutions. While it may generally be the case that governments join treaties by whose provisions they intend to abide; there may exist circumstances in which governments benefit by acceding to treaties whose provisions they intend to defy. Our model offers one instance in which this may take place.

Secondly, the paper speaks to the interaction of international regimes and domestic politics via an informational pathway. While the role of international institutions in generating information that facilitates cooperation among states is well recognized, here we identify an informational mechanism that affects domestic political conflict in significant and unexpected ways. The information generated by signing the CAT leads to less domestic opposition and the preservation of torturing regimes in power.

This paper also contributes to the literature on human rights law. We provide a theory of when and why authoritarian governments are likely to join human rights treaties, and provide empirical evidence in support of this theory. We also explore an "unintended consequence" of increased legalization of the human rights agenda – these legal instruments provide signaling opportunities to domestic oppositions of the elite's intent *not* to abide by its obligations, and result in the increased survival in office of torturing regimes.

Autocracies and the CAT

The United Nations Convention Against Torture and Other Cruel, Inhuman and Degrading Treatment or Punishment (CAT) was adopted in December 1984, went into effect in June 1987. It has been ratified by 139 states. It forbids

“any act by which severe pain or suffering, whether physical or mental, is intentionally inflicted on a person for such purposes as obtaining from him or a third person information or a confession, punishing him for an act he or a third person has committed or is suspected of having committed, or intimidating or coercing him or a third person, or for any reason based on discrimination of any kind, when such pain or suffering is inflicted by or at the instigation of or with the consent or acquiescence of a public official or other person acting in an official capacity.”(CAT Article 1)

The CAT requires that each member-state passes appropriate domestic laws making torture a crime, and requires that each state asserts jurisdiction when the crime occurs within its own territory, or the offender or the victim is a national of that state, or if the offender is present in its territory (if the member-state does not for some reason extradite the offender).

The emergence of a set of human rights treaties has been heralded as a major shift in the international system,² and a measure of the success and efficacy of international law. While these agreements have been ratified by most states in the world, repressive state behavior has continued to rise over time. Hafner-Burton & Tsutsui (2005) report that in 2000, while the average state had ratified 80% of all available human rights treaties, 35% of states are reported as having violated these agreements. States then are clearly willing to sign human rights agreements and continue to violate their treaty commitments.

The lack of compliance with human rights treaties is often viewed as stemming from a failure of enforcement. Enforcement of an international obligation has a number of prerequisites. First,

²The major human rights treaties (in addition to the CAT) are the International Convention on the Elimination of All Forms of Racial Discrimination (adopted 1965), the International Convention on Economic, Social, and Cultural Rights (1966), The International Convention on Civil and Political Rights (1966), the Convention on the Elimination of All Forms of Discrimination Against Women (1979) and the Convention on the Rights of the Child (1989). In addition other treaties, such as many Preferential Trading Agreements have both soft and hard prohibitions against human rights abuses (Hafner-Burton 2005).

failure to abide by the agreement must be observable. If violations are obscure, mixed in with noise or are otherwise difficult to observe or prove, enforcement is difficult and compliance unlikely. Second, there must exist a system of punishments to be imposed on a state or its elite in the event of a treaty violation to deter non-compliance (or rewards and incentives for compliance). And third, there must be some mechanism or process by which these costs are actually applied (or the benefits accrued). At the international level, this may be the withdrawal of concessions by a trading partner, or the application of sanctions (or enhancements to a state's trading or investment opportunities). At the domestic level, failure to abide by a ratified and implemented international agreement is likely to be a violation of domestic law and subject to sanction by domestic authorities currently or in the future.

Human rights treaties are viewed as being weak on all three dimensions. Violations are difficult to observe.³ The costs of non-compliance are low, and any potential benefits stemming from compliance such as enhanced trade or investment flows are absent (Nielsen & Simmons 2009). There are few mechanisms for enforcing the agreement.⁴ If non-compliance costs are in fact quite low, following Downs, Rocke & Barsoom (1996), we would expect most or all states to sign the CAT, and that state behavior on signing would be little (or un-)changed. The pattern of accession and compliance is somewhat different however. Many governments do not accede to the CAT (in our sample of 129 authoritarian regimes between 1985 and 1996, 74 regimes were never signatories); others sign and reduce torture levels, while still others sign and continue to torture.⁵

³The CAT does establish a monitoring committee, but it can only investigate and file a report of torture if the torturing state has explicitly accepted Article 21 and or Article 22; otherwise such allegations must be ignored, even if the state is a signatory to the rest of the CAT.

⁴A number of scholars have argued that even weak enforcement regimes can influence state behavior – by socialization into norms of appropriateness (Finnemore 1996) or cascades, where states feel pressured to conform (Keck & Sikkink 1998). Others have argued that the international regimes create openings for non-governmental actors (NGOs) to engage in information gathering, political action, legal maneuvering etc. that influence state behavior (Neumayer 2005, Simmons 2009). Moravcsik (2000) suggests that unstable democracies can “lock-in” human rights norms by treaty accession. Gilligan & Nesbitt (2007) argue that these norm-based arguments for the adoption of the CAT have not had any noticeable effect on torture levels. Nielsen & Simmons (2009) further find no evidence that signatory governments receive even praise from the US State Department on signing the CAT. If states exert pressure to sign the CAT and conform to its provisions, there is little evidence of this pressure in press statements.

⁵See Vreeland (2008) for a description of the variation in torture levels among dictatorships. Neumayer (2005) shows that torture levels fall in democracies and in other polities with richer civil society. Simmons (2009) argues that the CAT reduces torture in all but the most stable democracies and autocracies, due to the presence of NGOs and other civil society actors. Powell & Stanton (2009) demonstrate that an average of 83 percent of CAT signatories violate some CAT provisions each year, and 42 percent of signatories systematically violate CAT provisions.

Scholars have focused on domestic enforcement mechanisms to explain the observed variation in accession patterns and torture behavior. Hathaway contends that since domestic mechanisms to enforce compliance – such as an independent judiciary or an opposition party – are absent in autocracies, they find accession to human rights treaties essentially costless. In democracies, on the other hand, treaty violations are likely to impose costs on the incumbent government in the form of legal penalties or opposition attacks. Therefore, democracies are only likely to accede to human rights treaties if they are in compliance with these treaties' provisions *before* signing. Autocracies, however, should be willing to enter such treaties regardless of prior compliance. Both autocratic torturers and non-torturers will accede to the CAT. Contrary to this claim, however, Hathaway's (2007) empirical findings indicate that, amongst autocracies, there is a positive association between torture and the signing of the CAT.

Vreeland (2008) explores the domestic political and institutional dynamics of autocracies, and offers an explanation for Hathaway's puzzling finding. He contends that the positive association between levels of torture and accession to the CAT stems from omitted variable bias. More precisely, Vreeland argues that the presence of domestic opposition parties both causes autocrats to torture more heavily and forces these governments to sign human rights treaties. When opposition parties exist, there must be some freedom to engage in speech and activities that contradict the will of the incumbent government. In such a situation, opposition activists are likely to "cross the line" in their criticisms, leading the government to employ torture to maintain its control. These opposition parties will also pressure the government to enter into human rights agreements. Since Hathaway's regressions do not control for the presence of opposition parties, she finds a spurious association between torture and accession. When the presence of such parties is controlled for, the association between torture and the signing of the CAT drops to insignificance.

Vreeland's theory appears to rely on out-of-equilibrium behavior. If, as Hathaway claims, human rights treaties do not constrain autocratic governments, what is motivating the domestic opposition to push for treaty accession in the first place? Opposition groups are acting on out-of-equilibrium beliefs. If, on the other hand, human rights treaties do constrain autocratic governments, Vreeland does not articulate how these treaties do so. Nor is it clear why, if autocratic governments are

willing to so tie their hands, a treaty is necessary to enforce cooperation between the government and opposition.⁶

It may be argued that the CAT acts as a commitment mechanism that constrains the government from acting against the opposition following an agreement exchanging reduced levels of repression by the government for reduced anti-regime activities by the opposition. Simmons & Danner (2010) make just such an argument with regards to the International Criminal Court (ICC). They contend that human rights enforcement bodies allow governments to commit to refrain from aggressive acts against opposition groups.⁷ However, commitment problems between a government and an opposition are two-sided. If the government is tying its hands by signing the CAT, how does the opposition similarly commit to refrain from future anti-regime activities once the CAT is signed? Presumably, if the government has tied its hands by signing the CAT – or any similar human rights agreement – it opens itself to future oppositional efforts. Any theory postulating such a mechanism should fully specify the means by which the opposition can commit to a compromise agreement with the government, or else rely on *ad hoc* assumptions about the the credibility of opposition promises.

Such commitment-based arguments suffer from a further difficulty: Punishments meted out by human rights bodies to authoritarian elites are typically only applied *after* these elites are removed from power. With regards to the CAT, domestic courts are unlikely to indict a sitting leader in the event of human rights violations. And foreign governments are unlikely to detain prominent officials, even in the event these governments are in a position to do so. Consequently, authoritarian rulers who expect that their efforts to crush the opposition will be successful are unlikely to be deterred from human rights violations by the threat of prosecution. The CAT, or other such treaties, cannot serve as credible commitment devices for such powerful or determined regimes. Witness

⁶Empirically, the inclusion of a control for the presence of opposition parties causes the association between torture and the signing of the CAT to drop to insignificance only when a broad spectrum of other controls are also included in the Vreeland regressions. When only ‘opposition parties’ and ‘torture levels’ are used to predict signing, both are significant. Moreover, the inclusion of additional variables does not significantly reduce the magnitude of the coefficient on torture. Since there is a substantial amount of multi-collinearity between these ‘torture’ and ‘opposition’ measures, one cannot determine whether the newfound insignificance of ‘torture’ is simply due to problems of estimation. Without a more convincing theory of why the presence of opposition parties leads an autocrat to sign human rights treaties, there seems little reason to suppose otherwise.

⁷Simmons & Danner (2010) specifically refer to the negotiation of peace agreements in the wake of civil wars.

the recent behavior of Omar Bashir of Sudan, who remains free to travel and – by many reports – continues to commit rampant human rights violations, despite an indictment by the ICC.

Hafner-Burton & Tsutsui (2007) also explore the link between autocratic accession to the CAT and torture levels. They argue (as in Hathaway 2007) that while there are vague political benefits from CAT membership (“window-dressing”); these treaties lack coercion and enforcement mechanisms, fail to make states internalize or acculturate international norms, and do not cause a domestic human rights institutional capacity to emerge. Autocracies are therefore unlikely to show any evidence of improvement in repressive behavior after accession.

Yet explanations that stress the CAT’s lack of enforcement would seemingly suggest that *all* authoritarian regimes will sign. As mentioned above, this is not empirically the case. Moreover, it is unclear from either Hafner-Burton & Tsutsui (2007) or Hathaway (2007) why there exists a positive association between torture levels and the signing of the CAT by authoritarian regimes. It therefore remains an open question as to why those states with the worst human rights records sign these agreements with the greatest frequency and then ignore their obligations.

In the theory developed below, we concur with Vreeland’s focus on the role of the interaction between autocratic governments and their domestic opposition as it affects the signing of human rights treaties. But we view this interaction quite differently. We assume a game is played between an office-seeking government and an opposition party. Their interaction is characterized by attempts to maintain (seize) power: The government can undertake costly measures to repress the opposition even as the opposition can take costly actions to remove the government. We assume that the opposition is imperfectly informed as to the willingness of the government to employ (costly) repressive tactics. We demonstrate that, in such a game, the government may use the signing of human rights treaties as a signal to the domestic opposition that it is willing to take drastic action to retain power.⁸ In such an equilibrium, those governments that sign the treaty would torture more heavily *ex ante* than those that do not. Moreover, we find that signatory governments are likely to survive longer in office than non-signatories.

⁸Our theory is, in some ways, analogous to the literature on audience costs (see, for instance, Fearon 1994, Smith 1998). Whereas audience cost theories often presume that failure to comply by an international agreement reveals negative information about a government’s type (e.g. a lack of ability), we demonstrate that the willful defiance of an international agreement may be used to signal the government’s ‘strength.’

This logic may at first appear highly counterintuitive. However, it is in keeping with common perceptions of governments' defiance of international actors in situations not pertaining to human rights. For instance, it is widely believed that North Korea's 2009 nuclear test – despite international disapprobation – was meant to reinforce the regime's control following Kim Jung-il's ill-health and designation of a successor.⁹ Since a 'weak' regime would be unable to face the international pressures stemming from such a test, this action is a credible signal of the regime's strength to a domestic audience. Similarly, it is often argued that the Castro regime in Cuba enhanced its domestic stability by provoking the United States. One could view these actions as signals meant to intimidate domestic political opponents. Our theory suggests that the signing, and willful defiance, of human rights treaties might play a similar role.

Theory

Article 4 of the CAT states that "Each State Party shall ensure that all acts of torture are offences under its criminal law." Moreover "[e]ach State Party shall make these offences punishable by appropriate penalties." Article 5 requires that any State Party to the CAT take into custody any alleged offender that is present in its territory. And Article 6 requires that, if requested to do so, any State Party must extradite the alleged offender to any state with jurisdiction over the case, which may be defined by the nationality of the perpetrator or the victim. If no such extradition occurs, the State Party must try the offender domestically.¹⁰ Finally Article 8 further requires signatories to treat violations of the prohibition on torture as extraditable offenses.¹¹

The CAT does, therefore, make torture a more serious offense.¹² Consider an autocrat inclined

⁹Fackler, Martin. "Test Delivers a Message for Domestic Audience." *The New York Times*. May 25, 2009. <http://www.nytimes.com/2009/05/26/world/asia/26northk.html> – accessed November 7, 2009.

¹⁰This requirement is often referred to as establishing 'universal jurisdiction' for human rights offenses.

¹¹United Nations Convention Against Torture and Other Cruel, Inhuman or Degrading Treatment or Punishment. <http://www.hrweb.org/legal/cat.html>

¹²Some scholars question whether torture and other human rights offenses are currently covered under customary international law, and perhaps even enforceable in domestic courts (see e.g. Klein (1988)). Our point is that the CAT increases the potential costs of engaging in torture over and above that which might be expected under customary international law. Moreover, of the seven "core" human rights treaties, it may be argued that the CAT possesses the most serious enforcement mechanism. Goodliffe & Hawkins (2006) argue the CAT was the first treaty to apply the principle of universal jurisdiction to human rights law – jurisdiction is based "on the nature of the crime rather than .. where the crime occurred or the nationality of the alleged perpetrator or victim" (p.2) . As such, they suggest its enforcement mechanisms are more coercive than those of other human rights treaties or customary international law alone.

to torture in order extract information from or to punish a domestic political opponent. Should the autocrat, at some point in the future, find himself (and always it is himself, not herself) out of power, deposed or otherwise overthrown, the consequences will differ depending on whether the state was a signatory to the CAT. The usual act of a falling autocrat is to abscond to another country, if he manages to remain alive or out of jail. Assume that the autocrat's country were a Party to the CAT. If the country to which he escaped were also a CAT signatory, the autocrat's successor can demand the autocrat's extradition for trial for human rights violations. No such obligation would necessarily exist if the country is not a signatory to the CAT. On this basis, we argue that signing the CAT will at least weakly increase the penalties an autocrat would suffer after being evicted from office.

If an autocrat flees into exile – and if his state is unable or unwilling to try him domestically – the now host nation, if it is a CAT signatory, has an obligation to try the ex-dictator for human rights offenses. It is reasonable to think that, if the state of the offending dictator had signed the CAT too, the pressures for arrest and indictment would be higher than if it were not a CAT signatory. And finally, a third state can demand extradition from a fellow CAT member if the now host country fails to try the alleged perpetrator.

These provisions increase the expected costs of torture substantially. In the event an autocrat is removed from office, the danger of extradition may significantly limit his possible destinations for exile. The long term costs of this restriction on his movement would be considerable. Clearly therefore, and contrary to much of the scholarship on the CAT, there are post-tenure liabilities associated with engaging in torture. While these might be perceived to be unlikely to occur, or might happen only in the distant future, the costs from treaty violation are non-trivial in expected value. Moreover the probability that these liabilities will be applied increases as more countries apply the principle of universal jurisdiction – Goodliffe, Hawkins & Vreeland (2009) find that 109 states have incorporated universal jurisdiction in domestic law, 14 have tried cases under the principle, and the courts have enforced the law in 12 of them. What matters for our argument is not that universal jurisdiction is always applied, but rather that autocratic governments anticipate that it might be,

effectively raising the expected costs of engaging in torture.¹³

The costs imposed by the CAT have most vividly been illustrated in the extradition proceedings in the British House of Lords against Augusto Pinochet in 1998. Famously, the Law Lords ruled that Pinochet may be extradited to face criminal charges in Spain. Offenses after 1988 were ruled as extraditable, as 1988 marked the year that the UK ratified the CAT and passed domestic implementing legislation (Roht-Arriaza 2001). This finding allowed the prosecution of Pinochet to proceed despite a negotiated amnesty with his successor regime. Moreover, the Spanish prosecution (with the consent of the UK Law Lords) catalyzed Chilean courts to permit filings and as many as 170 complaints were subsequently brought in Chilean courts (Jonas 2004, Roht-Arriaza 2001). Similarly, former Chadian dictator Hissène Habré is under house arrest in Senegal for CAT violations. Despite findings from both the UN Commission on Torture and the African Union that Senegal is obliged under CAT provisions to either extradite or try Habré for torture that took place while he was in office, Senegal appears to be dragging its feet. Belgium has sued Senegal at the International Court of Justice arguing that Senegal is violating the CAT by neither prosecuting nor extraditing him, effectively giving Habré asylum. There is no ruling from the ICJ as yet, but this provides another example that the CAT is preventing Habré from escaping to a villa on the French Riveira, and raising the personal costs of exile (ICJ 2009).

Paradoxically, the increased cost signing the CAT places on repression ensures that those countries that torture heavily are most likely to sign.¹⁴ Assume that autocrats vary in the value they place on office and further assume that opposition groups are unable to perfectly observe this value.¹⁵ Those governments that most desire to retain office are more willing to employ torture

¹³In the costly signaling model developed below, as the costs of treaty violation go to zero, all governments pool on signing. As the costs increase, only the more repressive governments are likely to sign. As noted above, we do not witness all governments pooling on signing the CAT. Moreover, Hathaway's (2007) empirical findings are consistent with punitive costs that exceed the minimum threshold for separation. To the extent that Goodliffe & Hawkins (2006) are correct regarding the relatively punitive enforcement mechanisms of the CAT, we would be more likely to see this pattern of behavior in CAT accession than in the accession to other human rights treaties.

¹⁴Here – and throughout – we concentrate on the signing rather than the ratification of human rights treaties. We do so for two reasons: (1) The signing of treaties is the prerogative of the executive. Ratification may or may not (depending on the authoritarian system in question) be subject to the approval of other actors. (2) Ratification of a treaty follows its signing. Hence we argue that the act of signing a treaty likely carries the most informational content about the executive's costs of torture, rather than its ratification.

¹⁵It may be objected that opposition groups are aware of the costs a regime faces from repression and its willingness to employ draconian methods to remain in office. While such groups no doubt have some information in this regard, this information is not always perfect. The veterans of many successful opposition movements often express surprise at their

than those that value office less highly. However, since all governments would like to intimidate the opposition, no government can effectively communicate whether it is truly a ‘strong’ or ‘weak’ type.

Signing a human rights treaty imposes costs on autocrats who torture and are subsequently removed from office. Since ‘strong’ types are less likely to be so-removed, they are similarly less likely to bear any such cost. ‘Strong’ types are therefore more likely to sign such agreements than ‘weak’ types – i.e., human rights agreements serve as a costly signaling mechanism that conveys a government’s type to the opposition. Those governments that value office most highly – and are thus most inclined to employ torture – can sign human rights agreements at low expected cost. Such behavior would be consistent with existing empirical findings.

Moreover, those autocrats who sign such treaties will survive in office longer than non-signatories. A *selection effect* implies that those regimes that will fight most strongly to remain in power are the same regimes that sign the treaty. An *information effect* implies that opposition groups – on learning that the state has signed the treaty and is therefore a strong state – will engage in fewer activities designed to overthrow a signatory government. And a *commitment effect* implies that governments faced with the potential threat of prosecution on relinquishing office will be more willing to employ repressive tactics to remain in power.

Model

We model the signing of the CAT as the outcome of an interaction between an autocratic government G and its domestic opposition D . Both are assumed to be office-seeking: i.e. each derives some value from holding power. The opposition derives benefits $C > 1$ from office. The government derives benefits $R + C$ where $R \sim U[0, 1]$. The realization of R is known to the government, but unknown to the opposition. In keeping with our informal theoretical discussion above, the opposition is uninformed about the value the government places on holding power – and hence on the level of torture it will apply in equilibrium.

successes. And many failed opposition groups undertake costly activities in the vain hope of removing the regime. These actions are most readily explained by imperfect information. Theoretically, governments would only be able to perfectly reveal their willingness to employ repressive tactics if they had a continuous array of credible signals at their disposal. So long as some uncertainty exists, there remains an incentive for low cost governments to signal their type.

In the contest for power, the government may – at positive cost – engage in repressive measures entailing human rights violations against the opposition. Similarly, the opposition may undertake costly efforts to remove the government. The outcome of the contest for power will be determined, in part, by each party’s respective choice of repression and effort level.

If the government signs a human rights treaty, it faces increased costs ($P > 0$) on being removed from office. These costs are meant to reflect the dangers of extradition and prosecution faced by deposed government officials after signing the CAT. As noted above, sitting autocratic leaders need fear few legal repercussions from signing the CAT. But, those that are removed from office face the non-trivial dangers imposed by its imposition of extradition requirements and universal jurisdiction for human rights violations.

The sequence of the game is as follows: First, nature chooses the type of government $R \in [0, 1]$. This variable is observed by the government, but not by the opposition. Second, the government chooses whether or not to sign a human rights treaty $s \in \{0, 1\}$. Third, the opposition and the government simultaneously choose $e \geq 0$ – the level of effort put into deposing the government, and $t \geq 0$ – the level of repression. For simplicity, we assume that the choices of both government repression and opposition effort are made at a constant marginal cost of 1. Fourth, nature determines whether the government survives – with probability $\pi(t, e)$ – or not. All payoffs are realized and the game ends.

$\pi(t, e)$ is a standard contest success function (Hirschleifer 1991, Skaperdas 1996): $\pi(t, e) = \frac{t}{t+e}$.¹⁶ For simplicity, we let the distribution of government types be defined by the uniform distribution $f(\cdot)$ with support over the unit interval. This distribution is common knowledge.

Player utilities are defined by their expectation of holding office and the choice of s , t and e by the autocrat and the opposition respectively. The autocrat’s expected utility function is:

$$U_G(t, e, s; R) = \pi(t, e)[R + C] - t - s(1 - \pi(t, e))P$$

¹⁶The assumption that the probability of government survival is increasing in repression and decreasing in opposition activities is central to our results. This assumption is unlikely to be appropriate for democracies. We thus restrict our analysis to autocracies – for which this assumption is far more reasonable.

while the opposition's is defined as

$$U_D(t, e, s) = [1 - \pi(t, e)]C - e.$$

The government enjoys the rents from office $R + C > 0$ with probability $\pi(t, e)$ and pays a cost for repression equal to t . In the event that it is removed from office (with probability $[1 - \pi(t, e)]$), the government suffers a cost $P > 0$ if $s = 1$ and suffers no post-tenure punishment otherwise. The opposition, on the other hand, obtains C with probability $1 - \pi(t, e)$ and pays a cost for its anti-regime efforts of e .

Equilibrium

The game is solved through backwards induction using the perfect Bayesian equilibrium concept. Our first proposition establishes that there is an equilibrium in which governments that value office greatly sign the agreement, and those that value it little do not. In Appendix A we define two thresholds, $0 < \underline{P} < \bar{P}$.

Proposition 1. *If $\underline{P} < P < \bar{P}$, then there exists a unique semi-separating equilibrium where for $\tilde{R} \in (0, 1)$, the government will sign ($s = 1$) if $R \geq \tilde{R}$ and will not sign ($s = 0$) if $R < \tilde{R}$.*

All proofs can be found in Appendix A. Intuitively, this equilibrium implies that – for intermediate levels of post-tenure punishments – it is those governments that value office greatly that sign the CAT. Those that value office less will not sign. Note further that, *contra* standard selection arguments, this equilibrium posits that it is precisely those regimes that are *least* likely to follow the CAT's provisions that choose to sign. *Ceteris paribus* governments that value office greatly are more likely to apply repressive tactics to hold onto power. This finding is consistent with the puzzling empirical findings documented above – that authoritarian regimes that torture more heavily are more likely to sign the CAT than those that torture less.

The logic for this finding is straightforward. All autocratic governments seek to convince their opposition that they value office highly and are thus willing to employ high levels of repression, as this will serve to reduce the level of effort the opposition will put into removing the autocrat.

Signing a human rights treaty acts as a costly signal to the opposition of the value the government places on office. In equilibrium, the opposition learns that the government is a high value type, and that the marginal benefits of its anti-regime activities are lower than it thought. The opposition will thus reduce the level of its anti-regime effort e (since effort has declining marginal benefit). This benefits the government – as it faces a lower probability of removal from office $\pi(t, e)$.

The government benefits from the reduction in opposition effort – both because it increases its chances of survival in office and because it may respond by reducing levels of (costly) repression. However, it faces a higher stakes gamble insofar as it now risks the post-tenure punishments implied by the CAT. Post-tenure punishments may drive the government to exert greater efforts at repression such that they need not run the risk of post-tenure punishment.

The government must, therefore, weigh the costs of treaty accession against the benefits of lower opposition effort. If the level of post-tenure punishments is sufficiently low ($P < \underline{P}$), all autocrats pool on signing. Contrastingly, the level of post-tenure punishments is high $P > \bar{P}$, no government will sign (the threshold \tilde{R} goes to one making all types non-signers). For intermediate values of P , however, some governments choose to sign and others do not. Those that value office greatly suffer low expected costs from signing, as they will exert relatively high levels of repression and are thus quite likely to hold onto office. They also enjoy large marginal benefits from a reduction in opposition effort. Those that value office less highly are more likely to suffer the costs of post-tenure punishments, and derive a smaller marginal benefit from reduced opposition effort. Thus, it is those that value office most highly – and those that are most likely to employ repression – who sign the CAT.

Proposition 2. *In the semi-separating equilibrium, signatories will survive (weakly) longer in office than non-signatories.*

As we can derive the levels of government repression and opposition effort in equilibrium, we can also derive the probability of regime survival $\pi(t, e)$.¹⁷ This leads to the unambiguous conclusion that signatories survive in office with higher probability than non-signatories.

The survival effect stems from three causes. First, a selection effect is evident from Proposition

¹⁷See the appendix for the full specification of the semi-separating equilibrium and for the derivation of this result.

1 – autocrats value office greatly are more likely to sign than those that value it less. Signatory regimes are thus more willing to employ repressive tactics to cling to power than non-signatories. So the treaty selects those autocrats who would survive longer even in a world absent the CAT. But there are two additional causal effects of the CAT on regime survival: (1) Domestic opposition declines on signing the CAT, enhancing leader survival; (2) The CAT enhances the relative value of holding onto power, due to the threat of post-tenure punishment. Signatory governments may therefore become more willing to employ repression to remain in office.

In the two results that follow, we compare the opposition's effort levels, and the government's torture levels in two states of the world: the counterfactual – a world in which the CAT is not available as a signaling device, and the world with the CAT.

If there is no CAT, there is no opportunity for signaling. The domestic opposition makes its best guess about the type of government it is facing; given this level of domestic opposition, the more valuable is torture (the higher the value of office), the more torture the government will undertake. Combining this insight with that of Proposition 1, it is clearly evident that those authoritarian regimes that torture most heavily in the absence of the treaty are precisely those that are most likely to sign. We therefore now have a theoretical foundation for Hathaway's (2007) unexplained observation that the worst torturers are more likely to sign the CAT. Simmons (2009) makes a similar empirical finding of a positive association between torture and CAT signing. As a first test then, the model makes a prediction that is consistent with *extant* empirical work.

We can further compare equilibrium levels of torture absent the CAT with levels of torture in its presence. These results provide analytical insight into the expected effect of the CAT on the practice of torture.

Proposition 3. *Relative to a world where the CAT is absent, levels of government repression and opposition effort rise in non-signatory states.*

When the CAT is present, non-signatory governments are signaling their weakness (the low value they place on office) to domestic opposition groups. As the opposition now realizes that its efforts are relatively likely to be effective against a non-signatory government, it increases its anti-regime activities. As a result, non-signatory governments must similarly increase their levels

of repression to offset their greater risk of losing office.

Proposition 4. *Relative to a world where the CAT is absent, levels of opposition effort decline in all signatory states.*

The effect of the CAT on levels of repression in signatory states is more complicated. The *information effect* of the CAT implies that signing sends a signal of government strength to the opposition. Aware that their efforts are less likely to prove effective in unseating the government, the opposition reduces its anti-regime activities. *Ceteris paribus*, this will lead to a decline in the government's level of repression. However, signing the CAT increases the government's expected costs to losing office – the *commitment effect* of the CAT. The danger of post-tenure punishments causes governments to increase their levels of repression.

Our analytical results are silent on which effect dominates. This ambiguity is sensitive to modeling assumptions. In alternative versions of the model, we examine the effect of CAT signatory status when the CAT increases the marginal cost of torture, and derive a prediction in which torture declines in signatory states: the information effect outweighs the commitment effect.¹⁸ All other predictions remain unchanged. We explore this issue further in the empirical section below.

Examples

The equilibrium described above predicts that (1) those authoritarian regimes that torture most heavily will be most likely to sign the CAT, (2) signatory regimes will survive longer in office than observationally similar non-signatories, and (3) domestic opposition in authoritarian states declines on CAT signing.

These predictions – and the informational logic behind CAT accession – run counter to most *prima facie* expectations. However, several examples of authoritarian regimes that sign the CAT fit this logic rather well. For instance, Chad became a CAT signatory on June 9, 1995. The Chadian regime – headed by Idriss Déby – faced extensive armed opposition at the time, which it repressed through the extensive use of torture.¹⁹ The following year, the Déby regime unveiled a new con-

¹⁸For this version of the model, see Appendix B at homepages.nyu.edu/~jrh343.

¹⁹In 1995, Chad had a value of 4 on Hathaway's (2007) 5 point torture scale, and a 3 on CIRI's (2007) 3 point scale.

stitution, which controversially granted sweeping powers to the presidency. This constitution was adopted on March 31, 1996 and presidential elections – in which there were reports of extensive irregularities – followed soon after.²⁰ According to our theory, Déby's decision to sign of the CAT acted as a signal to opposition forces of his intention to cling to power. Following Propositions 2 and 4, Déby would be predicted to survive in office and the opposition would be expected to reduce its efforts at bringing about his ouster.

In fact, Déby did remain in power following elections in 1996 and remains in power currently. And, in 1997, several armed opposition groups ended their insurgency through negotiations with the government, consistent with Proposition 4.²¹

Following a similar logic, many authoritarian regimes signed the CAT immediately following or preceding a transition of power.²² For instance, the Museveni regime in Uganda signed the CAT in the year it assumed power.²³ Similarly, the Stevens government in Sierra Leone signed the CAT on March 18, 1985, immediately before handing power over to Stevens' chosen successor – Joseph Saidu Momoh – on November 28th of that year.²⁴ Our theory predicts that signing the CAT sends an informative signal of a regime's willingness to cling to office. Logically, the period immediately surrounding a change in the head of a regime would be a period of great uncertainty regarding the incoming elite's willingness and ability to cling to power. This is also likely to be the period during which the domestic opposition is particularly determined. As such, the informational value of signing the CAT is particularly great during transitional periods.

Of course, such results hardly constitute definitive evidence of the informational value of signing the CAT, though they are suggestive. Indeed, case studies are unlikely to provide strong support for our theoretical claims. We, in essence, argue that authoritarian regimes sign the CAT as part of

(Here, and throughout, the CIRI scale is inverted such that higher values on the CIRI index correspond to the more widespread use of torture.) See also: James, Odhiambo. "Human Rights: Pattern Changes but Violations Continue. *Africa News*. August, 1995.

²⁰Background Note: Chad.' US Department of State. Feb. 2009. <http://www.state.gov/r/pa/ei/bgn/37992.htm>

²¹Background Note: Chad.' US Department of State. Feb.2009.
<http://www.state.gov/r/pa/ei/bgn/37992.htm>

²²In our sample, just under 15 percent of authoritarian regimes that joined the CAT did so in a year of transition. This was the most commonly observed time of CAT signing in our sample.

²³Museveni assumed power on January 29, 1986 and the Museveni regime signed the CAT on November 3, 1986. Museveni currently remains in power.

²⁴Momoh remained in power until April of 1992.

an effort to deter opposition groups from undertaking anti-regime activities. As is argued by Achen & Snidal (1989), the use of case study evidence is problematic for assessing the effectiveness of deterrence. Our preferred tests of our theory therefore consists of a large-N analysis examining the claims of Propositions 1, 2, and 4.

Empirics

Since case study evidence is unlikely to provide conclusive in either supporting or falsifying our theory, we instead test the implications of the theoretical model above using large-N analyses. The model makes three testable predictions: (1) It is the most severe torturers that sign the CAT (Proposition 1); (2) Those authoritarian signatories that sign will survive in office with higher probability than observationally similar non-signatories (Proposition 2); and (3) Opposition effort declines on CAT signing (Proposition 4).

We further examine the association between signing the CAT and subsequent levels of torture. As noted above, our predictions with regard to levels of torture after signing the CAT are ambiguous and subject to modeling assumptions. The *information* and *commitment* effects of the CAT produce conflicting incentives for signatory governments with respect to torture, and our results do not clearly predict which effect dominates. While our theoretical expectations are mixed; these results are crucial to assessing the welfare implications of the CAT – so we include a preliminary empirical analysis below.

The first prediction – the most severe autocratic²⁵ torturers are more likely to sign the CAT – has been substantiated by previous work (Hathaway 2007, Simmons 2009 and Hafner-Burton and Tsutui 2007). We do not replicate these findings in detail here. However, we do provide evidence of the bivariate association between the eventual signing of the CAT and torture levels, to better demonstrate the empirical motivation at the root of this paper.

Simmons & Danner (2010) find suggestive evidence for our third prediction when examining

²⁵Autocratic regimes are identified following the definition advanced in Przeworski et al. (2000). Thus, autocracies are states which either lack executive or legislative elections, in which there exists either one party or no parties, or in which the ruling party has never been removed from power.

the effect of ratification of the Rome Statute, which grants jurisdiction to the ICC.²⁶ They find that ICC ratifications increase the probability of civil war terminations (which may reflect a reduction of opposition effort, government repression, or both). We examine a similar relationship with respect to civil war termination and the signing of the CAT – as well as the relationship between CAT signings and a host of other measures reflective of opposition effort.

Data

In all sets of regressions, we employ data from Vreeland's (2008) dataset on CAT accession, the Archigos database on political leaders and regime survival (Goemans 2006), and Gandhi and Przeworski's (2007) data on the longevity of authoritarian regimes.²⁷

Our measures of domestic unrest are derived from UCDP/PRIO Battle Deaths Dataset Version 3.0 (Lacina & Gleditsch 2005). We make use of annually observed observations of battle deaths resulting from civil wars, which can be viewed as a function both of opposition efforts to remove the government and of government repression of the opposition. We also rely on data from Banks' Cross-National Time-Series Data Archive (Banks 1979) (drawn from the dataset made available by Bueno de Mesquita et al. (2003)) to measure riots, strikes, revolutions, demonstrations, and other anti-government activities. We view these variables as measures of oppositional activity aimed at displacing the government.

Torture and the CAT

As noted above, many authors – starting with Hathaway (2007) – have noted a positive association between the torture levels employed by authoritarian regimes and their willingness to sign the CAT. Authoritarian governments that torture more heavily in time t are more likely to sign the

²⁶It is important to note that Simmons & Danner (2010) offer a quite different explanation for this relationship than that advanced here. They contend that the ICC acts as a commitment mechanism that enables governments to strike a bargain with opposition groups to reduce levels of violence. We do not find such an argument credible with respect to the CAT. As argued above, such a theory ignores the two-sided nature of commitment problems in domestic conflicts and fails to account for the post-tenure nature of punishments inflicted by the CAT. It is beyond the scope of this paper to test the extent to which our theory also applies to the ICC, but we note that the findings of Simmons & Danner (2010) are consistent with the mechanism postulated herein.

²⁷Throughout our analysis, we define a 'regime' as a single leader's tenure in the Archigos dataset. In this our analysis differs from both Hathaway's and Vreeland's. Both of these authors use countries as subjects and the country-year as the unit of observation. We use both the regime and the regime-year as our unit of observation in different analyses.

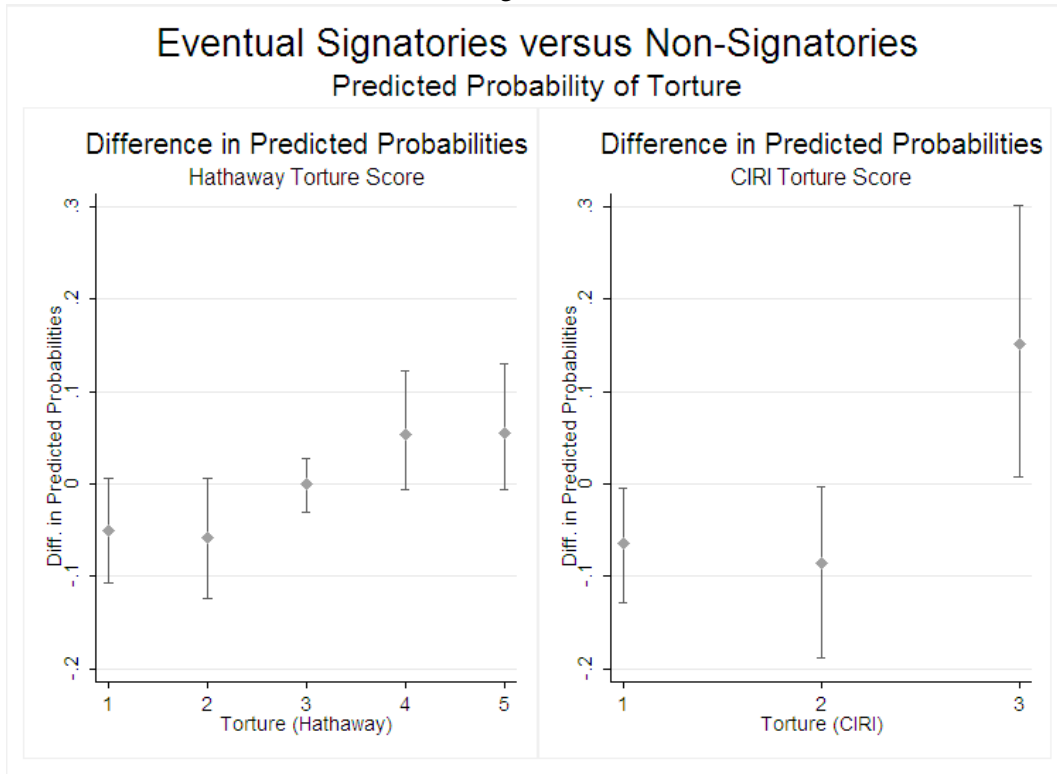
CAT in time $t + 1$. We argue that this paradoxical association arises because of the CAT's role in the interaction between authoritarian governments and their oppositions – those regimes that are more willing to employ repressive methods to hold on to power sign the CAT as a signal of their determination.

The empirical difference in the levels of torture employed by eventual signatories and non-signatories of the CAT is non-trivial. To assess this difference, we regress both the CIRI and Hathaway torture measures against an indicator variable that takes the value 1 if a regime eventually signs the CAT. Since this difference is expected to exist *prior* to signing the CAT, we drop all leader-years *after* signing takes place. These regressions are run only on authoritarian regimes.

Figure 1 demonstrates the difference in the predicted probability of witnessing different torture levels between eventual CAT signatories and non-signatories. Results are obtained from an ordered probit regression of torture levels on an indicator for eventual signatory status. As the regression results make clear, the differences in predicted probabilities are substantively meaningful – eventual signatories are roughly 5 percentage points more likely to score a 4 or 5 on the Hathaway torture measure and roughly 15 percentage points more likely to score a 3 on the CIRI torture measure than non-signatories. These differences border on significance at the 95 percent level of confidence. A simple difference in means test (not reported) reveals that this difference corresponds to a roughly $\frac{1}{3}$ of a standard deviation higher level of torture according to the Hathaway measures, and roughly $\frac{1}{2}$ of a standard deviation higher level of torture according to the CIRI measure.

The torture measures used here are no doubt subject to substantial measurement error. Moreover – as is widely noted in the human rights literature – these indexes are truncated at high levels of torture (Clark & Sikkink 2010, Hafner-Burton & Ron 2009, Wood & Gibney 2010). The results reported above therefore likely understate the true magnitude of the difference in torture levels between eventual signatories and non-signatory regimes. Both attenuation bias and the truncation of the indexes tend to understate the magnitude of the relationship between signing the CAT and torture.

Figure 1:



Differences between eventual CAT signatories and non-signatories in predicted torture scores. Predicted probabilities are generated from a bivariate ordered probit regression of torture levels on an ‘eventual signatory’ dummy. Dots indicate point estimates, lines indicate 95 percent confidence intervals. All standard errors are clustered at the country level. A score of 1 on the CIRI measure corresponds to zero reported instances of torture; a score of 2 to 1-49 reported instances of torture; and a score of 3 to 50+ reported instances of torture. A score of 1 on the Hathaway measure corresponds to no allegations of torture or individual security personnel are punished for instances of torture; a score of 2 to unsubstantiated instances of torture; a score of 3 to isolated or occasional instances of torture; a score of 4 to numerous or common instances of torture; and a score of 5 to prevalent or widespread torture. Predicted probabilities and confidence intervals are generated using the CLARIFY software package (Tomz, Wittenberg & King 2001) run from Stata 11.

Leader Survival

We next subject the claim that signatories of the CAT enjoy more secure tenures in office than non-signatories – derived in Proposition 2 above – to empirical scrutiny. Note that this claim is *not* equivalent to the statement that signing the CAT causes an authoritarian regime to survive longer in office. Our prediction is that signatories face a lower probability of removal due, in part, to a selection effect. Signatories face a systematically lower cost to repression than non-signatories.

Signatories also benefit from the revelation of information to opposition groups. Precisely because only low-cost repressors choose, in equilibrium, to sign the CAT, these regimes face a less restive domestic opposition. Domestic opposition groups are less willing to exert costly effort against regimes that reveal that they are low cost types by signing the CAT. Proposition 2 holds as a result of the cumulative force of these selection and information effects. Our tests of this proposition cannot differentiate between the two effects.

Single Record Cox Estimates

To test the association between the signing of the CAT and the survival time of regimes, we first run a Cox proportional hazards model of the probability of regime failure. The unit of observation in this model is the regime. Time is defined by *sumten*, the sum total number of days a given regime served in office, as taken from the Archigos dataset. The key explanatory variable is *Eversign*, a binary indicator variable that takes the value of 1 if a given regime is ever a signatory of the CAT.²⁸ 129 authoritarian regimes are observed during the 1985-1996 period.^{29,30}

The Cox model provides an estimate of the hazard rate of a given regime i (the probability regime i collapses at time t given that it has survived until time t) conditional upon observed covariates: $h_i(t) = h_o(t)e^{X_i\beta}$, where $h_o(t)$ is the baseline hazard function. The Cox model makes no assumptions about the parametric form of the baseline hazard function, which is estimated non-parametrically based on the exit times of regimes in the dataset. As our theory makes no predictions regarding the effect of time-in-office on regime survival, we treat time as – in effect – a nuisance parameter.³¹ The results of this regression are reported in Table 1.

²⁸The use of a similar variable *Everratify* produces estimates of similar magnitude and direction; though the coefficient is not significant at conventional levels.

²⁹These data are subject to a left-censoring problem. The *sumten* variable takes on values such that some regimes begin life before they come under observation. This censoring is unlikely to effect our results, as regimes that sign the CAT do not significantly vary in their age (at the time of signing) from those that do not. There is also a right-censoring issue, insofar as we do not observe all regimes that will sign the CAT. Since regimes that go on to sign the CAT after the sample ends are expected to have a lower hazard rate than those that will never sign the CAT, this would be expected to bias our results downwards.

³⁰The period under observation begins when the CAT comes into force. The 1985-1996 period is identical to that covered in the Vreeland (2008) dataset.

³¹The Cox model does assume hazard rates are proportional – that the shape of the hazard function does not differ across units. We test this assumption using the Grambsch-Therneau of the Schoenfeld residuals (Box-Seffensmeier & Jones 2004). We do not reject the null that the hazard functions are proportional. However, Harrell's rho statistics, testing individual covariates for violations of the proportional hazards assumption, suggest that the measures of torture

The controls include the variables identified as significant to the survival of authoritarian regimes by Gandhi & Przeworski (2007) – the number of changes in the executive during a given authoritarian spell, whether or not the executive has a civilian background, whether or not the government inherited one or more opposition parties, and an indicator for resource dependence.³² We also control for variables relevant to the signing the CAT identified by Vreeland (2008) and Hathaway (2007), notably the average level of torture employed by the regime,³³ an indicator variable that takes the value one if multiple parties are allowed to exist, and the population (in millions). Finally, we add a control for the military capabilities index compiled by the Correlates of War project, averaged over each regime (Singer 1987). As can be readily seen in Table 1, regimes that are signatories of the CAT have lower hazard rates than those that do not.³⁴ This difference is apparent even after controlling for other factors related to regime survival and for factors related to CAT accession.

It may be objected that our definition of the explanatory variable of interest in these models is problematic. The *Eversign* variable takes a value of 1 for *all* regimes that are signatories of the CAT, even those that inherit their signatory status from a predecessor government. It may reasonably be argued that such regimes are unlikely to withdraw from the CAT, regardless of their willingness to employ repression to stay in power. Choosing to remain in the CAT sends a very different signal than choosing to sign the CAT. If it is true that regimes that inherit their signatory status are unwilling to remove themselves from the treaty, then our results in Table 1 will be biased downwards. We therefore drop all regimes that inherit their signatory status and rerun the model above. Results are reported in Table 2.

The sign on the *Eversign* variable remains negative and increases in value and in significance.

used should be interacted with time. We include these interactions in our specifications below. Results are substantively unchanged if these interactions are not included.

³²This indicator takes the value 1 if primary commodity exports exceed 50 percent of total exports.

³³We employ both the Hathaway (Hathaway 2007) and CIRI (Cingranelli & Richards 2007) measures of torture. Both are based on data made available by the US Department of State and by Amnesty International. The Hathaway measure is an ordinal index running from 1 to 5 with higher values indicating the more extensive practice of torture. The CIRI index is an analogous measure with values running from 1 to 3. Both are drawn from the dataset used in Vreeland (2008).

³⁴The reported coefficients are not hazard ratios. A coefficient of zero implies that the variable in question has no effect on the hazard rate, negative coefficients imply that an increase in the variable reduces the hazard rate, and positive coefficients imply the reverse. We have also run models wherein the *Eversign* variable is interacted with the indicators for regime-type. These estimates do not indicate any significant difference in the relationship between CAT accession and survival between military and non-military regimes.

Table 1: Coefficient Estimates from a Single Record Cox Model

	Hath1	CIRI1	Hath2	CIRI2
Eversign	-.662 (.305)	-.715 (.305)	-.661 (.282)	-.699 (.293)
Torture	.232 (.233)	-.299 (.261)	.361 (.228)	.889 (.477)
Torture*time	-.0001 (.00004)	-.0002 (.00007)	-.0001 (.00003)	-.0002 (.00006)
No. Changed in Exe.	.146 (.032)	.136 (.031)	.133 (.034)	.117 (.037)
Civilian Exe.	.042 (.271)	.07 (.26)	-.033 (.289)	-.054 (.297)
Inherited Opposition Party	.303 (.23)	.321 (.224)	.4 (.25)	.387 (.242)
Multiple Parties	-.536 (.336)	-.662 (.33)	-.555 (.346)	-.633 (.334)
Resource Dependence	.416 (.422)	.462 (.441)	.203 (.478)	.158 (.463)
Population	.004 (.007)	.006 (.008)	.	.
COW Capabilities Index	40.132 (67.525)	12.679 (68.319)	.	.
GDP <i>per capita</i>	-.091 (.045)	-.074 (.04)	.	.
No. Subjects	129	129	129	129
No. Failures	71	71	71	71

Results from a single record Cox Survival analysis of regime tenure. Hazard rates estimates based on number of days in office as measured by the Archigos dataset. Estimates are constructed for a sample of 129 authoritarian regimes in the years 1985-1996. Coefficient estimates are of the form $h_i(t) = h_0(t)e^{\mathbf{x}\beta}$. **Hath** models use Hathaway's (2007) torture index; **CIRI** models use the CIRI (2007) index. All standard errors are clustered by country.

Table 2: Coefficient Estimates Dropping Regimes that Inherited Signatory Status

	Hath1	CIRI1	Hath2	CIRI2
Eversign	-1.004 (.372)	-1.076 (.352)	-.985 (.337)	-1.032 (.328)
Torture	.398 (.291)	1.31 (.586)	.481 (.268)	1.323 (.524)
Torture*time	-.0001 (.00003)	-.0002 (.00006)	-.0001 (.00003)	-.0002 (.00006)
No. Changes in Exe.	.142 (.039)	.122 (.039)	.125 (.045)	.103 (.046)
Civilian Exe.	-.028 (.314)	-.028 (.289)	-.122 (.343)	-.156 (.319)
Inherited Opp. Party	.238 (.261)	.236 (.261)	.313 (.277)	.329 (.269)
Multiple Parties	-.552 (.375)	-.774 (.382)	-.515 (.368)	-.658 (.359)
Resource Dependence	.309 (.454)	.227 (.463)	-.06 (.508)	-.127 (.498)
Population	.006 (.008)	.011 (.009)	.	.
COW Capabilities Index	18.847 (87.747)	-55.968 (96.877)	.	.
GDP <i>per capita</i>	-.111 (.054)	-.079 (.044)	.	.
No. Subjects	108	108	108	108
No. Failures	62	62	62	62

Results from a Cox Proportional Hazards estimate of regime tenure. Hazard rates estimates based on number of days in office as measured by the Archigos dataset. All regimes that inherited their status as a signatory government from a previous regime are dropped from the sample. Estimates are constructed for a sample of 108 authoritarian regimes in the years 1985-1996. Coefficient estimates are of the form $h_i(t) = h_0(t)e^{\mathbf{x}\beta}$. **Hath** models use Hathaway's (2007) torture index; **CIRI** models use the CIRI (2007) index. All standard errors are clustered by country.

These results are consistent both with Proposition 2 and with the claim that regimes that inherit their signatory status are unlikely to remove themselves from the CAT.³⁵

These results *cannot* be interpreted causally. It is very possible that our model results are being driven by a selection effect – particularly selection on unobserved covariates. However, Proposition 2 posits selection on covariates unobserved to the opposition party. The simple binary relationship – that those regimes that sign the CAT survive in office longer than those that do not – is consistent with model predictions.

Multiple Record Survival Analysis

In addition to the single-record model analyzed above, we run a multiple-record Cox survival analysis of the relationship between CAT signing and authoritarian regime survival. The unit of analysis in this model is the regime year. Time is defined by the difference between the year in which a regime took office and the current year.³⁶ The explanatory variable of interest is Lagged CAT Signing, a binary indicator variable that takes the value 1 in the year following a CAT signing. The model thus tests if signing the CAT in time t increases regime stability in year $t + 1$. 116 authoritarian regimes are observed over the 1985-1996 period.

As in the previous section, we run a Cox proportional hazard model of the form $h_i(t) = h_0(t)e^{\mathbf{X}\beta}$. We have tested the proportional hazards assumption using Grambsch-Therneau tests of the Schoenfeld residuals. The tests do not reject the null hypothesis that the proportional hazards assumption holds.³⁷

Table 3 reports coefficient estimates from this model. Columns marked **Hath** control for the Hathaway (2007) measure of torture, whereas those marked **CIRI** control for the CIRI (2007) measure. In all models the coefficient on Lagged CAT Signing is negative, implying that signing the

³⁵We have also run robustness checks on the single record model by (1) dropping all regimes that survive for less than either one or two years from the dataset, and (2) running a propensity score matching algorithm to ensure that the data are balance with respect to CAT signing. These results do not change the direction of the coefficient estimates reported above, and all results are significant at the 10 percent level or above. Results are available from the authors on request.

³⁶Data is `stset` in Stata to adjust for problems of both left and right-censoring. See Box-Seffensmeier & Jones (2004) for more details on censoring.

³⁷Harrell's rho statistics on parameter specific violations of the proportional hazards assumption are bordering on significant for `cinc`, `war` and `population` across both specifications. Including time interactions with these terms does not substantively affect the coefficients of interest. Results are available from the authors on request.

Table 3: Coefficient Estimates from a Multiple Record Cox Model

	Hath 1	Hath 2	CIRI 1	CIRI 2
Lagged CAT Signing	-1.328 (.696)	-1.418 (.667)	-1.142 (.753)	-1.248 (.742)
Torture	-.385 (.221)	-.349 (.194)	-.493 (.266)	-.468 (.25)
No. of Changes in Exe.	.159 (.045)	.169 (.042)	.139 (.041)	.149 (.042)
Civilian Exe.	-.327 (.369)	-.073 (.317)	-.386 (.389)	-.151 (.325)
Inherited Opp. Party	.264 (.28)	.	.269 (.29)	.
Multiple Parties	-.272 (.383)	.	-.271 (.377)	.
Resource Dependence	.016 (.484)	.	.101 (.478)	.
Growth	-.018 (.008)	-.018 (.008)	-.019 (.009)	-.019 (.008)
GDP per capita	-.007 (.04)	.	-.023 (.04)	.
COW Capabilities Index	79.097 (75.767)	34.521 (44.924)	71.092 (75.674)	25.545 (47.662)
War	.869 (.398)	.758 (.329)	.987 (.454)	.796 (.376)
Population	-.006 (.009)	.	-.007 (.009)	.
No. of Subjects	113	116	112	116
No. Failures	41	43	41	42

Results from a Cox Proportional Regime survival function estimates. Survival function estimates estimates based on number of years in office as measured by the Archigos dataset. Estimates are constructed for a sample of 90 authoritarian regimes in the years 1985-1996. Coefficient estimates are of the form $h_i(t) = h_0(t)e^{\mathbf{x}_i\beta}$. **Hath** models use Hathaway's (2007) torture index; **CIRI** models use the CIRI (2007) index. All standard errors are clustered by country.

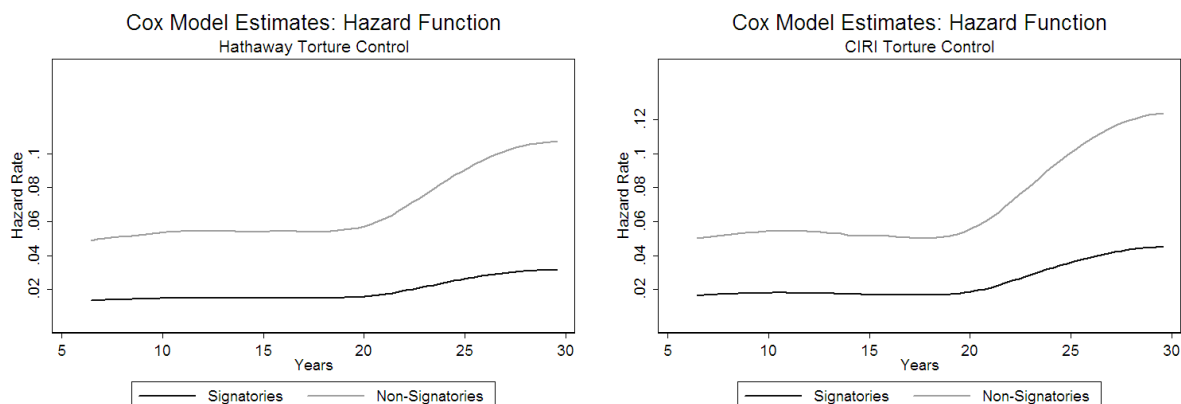
CAT in year t reduces a regime's hazard rate in year $t + 1$. These estimates are significant in all but one specification.³⁸ These results offer robust support for Proposition 2.

Since the Cox model is non-linear, the coefficients reported in Table 3 can be difficult to interpret. To give a better sense of the strength of the association implied by these estimates, we plot the baseline hazard function when Lagged CAT Signing is equal to 0 and when it is equal to 1 in Figure 2.

It is also worth noting that the coefficient on Torture is negative and significant at the 10 percent

³⁸The level of significance in this specification $p = 0.13$ borders on significance.

Figure 2: Hazard Function Estimates



Hazard function plots based on the models reported in Table 3 above. Plots using the Hathaway measure of torture as a control variable are plotted to the left, those using the Ciri measure of torture are plotted to the right. Hazard rates are depicted on the y-axis, while time in office (in years) is plotted on the x-axis. Signatories are depicted by the more darkly shaded line, non-signatories by the more lightly shaded line. All other covariates are set to their mean levels to produce this estimate.

level in all multiple record specifications. These results are consistent with the assumptions of the model – that regimes employ repressive tactics to maintain their hold on power.

Opposition Effort and Government Repression

While the above results offer strong support for a novel prediction generated by our model, they do not constitute direct tests of the mechanisms driving our results. We claim that an authoritarian government's decision of whether or not to sign the CAT can be modeled as a costly-signaling game. Since regimes that need to fear losing their grip on power (i.e. those regimes that value power slightly) face the greatest expected punishment under the CAT, these regimes will opt not to sign. Those that value office highly, on the other, on the other hand, will be willing to sign. As a result, the opposition will cue off of the government's signatory status, and will reduce their costly anti-regime efforts on witnessing a government sign the CAT.

It therefore follows that signing the CAT should be associated with a reduction in opposition efforts to remove the government. It *may* also be associated with reductions in government levels of repression. However, both of these terms are extremely difficult to accurately measure. As proxies, we turn both to the UCDP/PRIO estimates of battle deaths suffered during civil wars and

to the Banks measures of regime instability.

We employ both sets of measures to ensure the robustness of our results. Variation in results across these different indexes may also reflect variation in types of opposition movements. If the signing of the CAT leads to reductions in battle deaths experienced in civil wars, this suggests that the CAT acts as a useful signal to organized and armed domestic opposition groups. Similar results in the Banks dataset indicate the effect of the CAT in signaling to mass political movements. These different forms of domestic opposition may plausibly be expected to possess different levels of information regarding the willingness of the regime to employ torture to remain in power. They may also plausibly differ in their willingness to endure great costs to engaging in anti-regime activities.

The UCDP/PRIO dataset contains estimates of the number of battle deaths suffered in all wars with at least 25 fatalities between 1946 and 2005. Our data contain battle deaths estimates from wars UCDP/PRIO classifies as civil wars (types 3 and 4 in the PRIO data). We rely only on observations for which annual battle deaths estimates are available.³⁹ And we code all observations wherein no civil war was active (i.e. fewer than 25 battle deaths) as having zero battle deaths.

Our causal claim is that the signing of the CAT should lead to a decline in the number of battle deaths experienced in the following year. To test this claim, we employ a difference-in-differences specification $\Delta \text{battleddeaths}_{i,t} = \gamma \Delta \text{CATsignatory}_{i,t-1} + \Delta \mathbf{X}_{i,t} \beta + \epsilon_{i,t}$ where i denotes country, t denotes year t and Δ is the difference operator. The difference operator ensures that all unit specific effects that are constant between year t and $t - 1$ are controlled for in the regression estimates.⁴⁰ All specifications include controls for military capabilities, GDP *per capita*, and economic openness as well as for a cubic polynomial of time. We estimate this model using both the high and low estimates for battle deaths from the UNCP/PRIO dataset, and employ a seemingly unrelated regressions model, given the likely correlation of the error term across these estimates.

In addition to estimating this model on the full time-series cross-section of country-years in our dataset, we also pre-process our data using propensity score matching. Pre-processing the data in this manner is recommended as a measure to ensure covariate balance and reduce model

³⁹In some instances UCDP/PRIO estimates the total battle deaths from a given war and divides these deaths evenly over the period during which the war took place. Such observations are treated as missing values in our dataset.

⁴⁰For this reason, we do not include the Gandhi & Przeworski (2007) controls in these specifications, as these measures are largely constant over time.

dependence, particularly when working with a binary regressor (Ho et al. 2007). Since we are attempting to match panels – rather than individual observations – we employ the method of Simmons & Hopkins (2005). For countries that eventually sign the CAT, we collapse all pre-signing observations. For countries that do not, we collapse all observations over the full 1985-1996 period. We then estimate the probability that a given country ever signs the CAT based upon the covariates in the collapsed dataset. Panels are matched to one another using the MatchIt package (Ho et al. 2004) run from R 2.7.1, according to the genetic matching algorithm of Diamond & Sekhon (2006). The resultant weights are then merged with the full dataset.⁴¹

The results of these models – run on both the matched and unmatched datasets – are reported in Table 4. The coefficient on the lagged change in CAT signatory status is negative and highly significant (either at the 1 percent or 5 percent level) in all regressions. The association is also substantively large. Changes in the low estimates of battle deaths in the full sample had a mean of approximately -32 and a standard deviation of approximately 856. Thus, the signing of the CAT is associated with a roughly 0.7 standard deviation decline in the low estimates of battle deaths.⁴²

As a robustness check of our estimates reported in Table 4, we run a similar model using the Banks (1979) measures of regime instability. Included in these measures are the number of assassinations, strikes, government crises, riots, revolutions, opposition demonstrations, and guerrilla movements in a given country in a given year. The models are identical to the above in all but the regressand.

The Banks data, it is important to note, are derived from reports in the *New York Times*. As a result, they are likely subject to a substantial degree of measurement error of a particular kind. Since the *Times* is unlikely to report on a strike or demonstration that did not in fact take place, it is reasonable to assume that any measurement error is likely to result in an *underreporting* of events.

⁴¹Matching diagnostics are available from the authors on request.

⁴²The identification of these effects relies on a very small pool of countries that were experiencing a civil war at the time they signed the CAT. The results are thus sensitive to the exclusion of outliers. To adjust for this problem, we ran a fixed-effects OLS specification, with the regressor of interest a dummy that takes the value of 1 if the country-year is a CAT signatory. These results are consistent with the difference-in-differences estimate and indicate that CAT signatory status is associated with a substantial and significant decline in battle deaths. Moreover, these results are not sensitive to outlier observations (dropping a single observation never causes the results to decline in significance below the 95 percent level).

Table 4: CAT Signing and Civil War Battle Deaths: Seemingly Unrelated Regression Model

	Δ PRIO Low Battledeaths Estimate			
	Matched Dataset		Full Dataset	
Lagged CAT Signing	-623.486 (183.765)	-610.555 (183.84)	-582.804 (142.733)	-479.823 (124.52)
change Military Capabilities	39.279 (149.752)	-35.824 (140.141)	10.625 (87.904)	-20.987 (74.277)
change GDP <i>per capita</i>	-171808 (208743)	-150511.5 (208494.4)	-34204.07 (54874.4)	-2203.521 (17050.05)
change growth	-6.617 (4.753)	.	-3.419 (2.805)	.
change openness	1.161 (3.86)	.	.847 (2.212)	.
Cubic Time Polynomial	✓	✓	✓	✓
	Δ PRIO High Battledeaths Estimate			
	Matched Dataset		Full Dataset	
Lagged CAT Signing	-2092.317 (864.664)	-2038.021 (863.805)	-1872.469 (715.704)	-1470.178 (645.406)
change Military Capabilities	97.228 (704.625)	-167.57 (658.477)	-189.672 (440.777)	-347.977 (384.989)
change GDP <i>per capita</i>	-517324.4 (982191.6)	-432419.2 (979650)	-107367.3 (275156.3)	-2685.523 (88373.34)
change growth	-23.354 (22.365)	.	-19.158 (14.067)	.
change openness	-.298 (18.162)	.	-9.189 (11.094)	.
Cubic Time Polynomial	✓	✓	✓	✓
N	333	333	493	548

Results from a difference-in-differences model run using the UNCDP/PRIO measures of battle deaths. Estimates are from a seemingly unrelated regressions model. Coefficients to the left are estimates from models run on the matched dataset; those on the right are estimates from models run on the full dataset.

Table 5: CAT Signing and Levels of Unrest: Seemingly Unrelated Regression Model

	Matched Dataset		Full Dataset	
	Controls	No Controls	Controls	No Controls
Δ Assassinations	.067 (.205)	-.002 (.193)	.044 (.168)	.012 (.12)
Δ Strikes	-.275 (.157)	-.219 (.148)	-.204 (.126)	-.157 (.085)
Δ Gov't Crises	-.048 (.098)	-.035 (.097)	-.027 (.088)	-.016 (.071)
Δ Riots	-.112 (.381)	-.141 (.366)	.044 (.318)	-.01 (.222)
Δ Revolutions	-.15 (.205)	-.162 (.19)	-.146 (.205)	-.147 (.146)
Δ Demonstrations	-.023 (.414)	-.045 (.405)	.094 (.353)	-.134 (.281)
Δ Guerrilla Movements	-.077	-.065	-.063	-.064
N	328	365	497	743

Coefficient estimates on a lagged CAT signing indicator from a difference-in-differences model run using the Banks measures of domestic unrest. Estimates are from a seemingly unrelated regressions model. Regressands are reported in the left column, coefficient estimates and standard errors in the four other columns. Coefficients to the left are estimates from models run on the matched dataset; those on the right are estimates from models run on the full dataset. Columns marked 'Controls' include first-differenced measures of the growth rate of GDP, the level of GDP, the level of military capabilities, and the degree of trade openness. All estimates include controls for a cubic polynomial of time.

Such underreporting is likely to bias our difference-in-differences estimator towards zero.⁴³ As such is the case, it would be particularly surprising to find any systematic relationship between CAT signing and opposition activity using this data.

Our estimates from these difference-in-differences models are reported in Table 5. As was true with the battle deaths data, we employ a seemingly unrelated regressions model to adjust for the likely correlation in the error terms across these specifications.

As might be expected, our coefficient estimates are rarely significant (the lone exceptions are for two models testing the association between CAT signings and changes in the number of strikes). However, the coefficient estimates are nearly uniformly (23 of 28) negative – i.e. in the direction

⁴³Assume, by way of example, that every event has a 95 percent chance of appearing in the *Times* data, and that the probability of observation is independent across events. Let us say that there are 5 riots in a signatory country in year $t - 1$ and this number is reduced to 3 riots in the year after signing the CAT. The expected observed difference is given by $.95 * 5 - .95 * 3 = 1.9$, which is strictly less than the true difference of 2. Conversely, a non-signatory government that experience the same number of riots in year $t - 1$ and did not experience a reduction in coups stemming from signing the CAT would have an expected observed difference of 0.

postulated by our theory. Moreover, the coefficients are often substantively large. For instance, the coefficient on strikes suggests that CAT signing is associated with a 0.5 standard deviation decline in the number of strikes from one year to the next. The coefficient on revolutions suggests a 0.2 standard deviation decline in the change in the number of revolutions. Our interpretation of these results, therefore, is that the data weakly support our theory. Given the substantial problems with these data, consistently negative sign and non-trivial magnitude of our coefficient estimates are unlikely to be purely the result of chance.

Change in Torture Levels

Our theoretical model suggest that the signing of the CAT could cause either an increase or a decrease in realized torture levels. This depends on the net result of two contrasting effects: the *information effect* of the CAT – by which opposition groups reduce their anti-regime efforts following signing – and the *commitment effect* – by which incumbent regimes cling more ferociously to power to avoid the post-tenure punishments associated with torture. While our model is silent on which dominates, we can see if the data can provide any insight.

We test which effect dominates below. To conduct these tests, we rely on the torture indexes created by Hathaway (2007) and CIRI (2007). The former is an ordinal index ranging from 1 to 5, with higher values corresponding to the more widespread use of torture. The latter is a similar index with values ranging from 1 to 3. Both indexes rely on information from Amnesty International and the US Department of State to derive their torture scales.

Of course, the extent of torture is inherently difficult to measure. The practice of torture is unlikely to be fully documented and highly repressive regimes are unlikely to be particularly transparent to either Amnesty International or the US Department of State. Both the Hathaway and CIRI data thus no doubt suffer from a good deal of measurement error.

Since these indexes have minimal and maximal values – and vary only over a restricted range – this measurement error *cannot* be normally distributed with mean zero. A country-year that should be classified as a 1 on the CIRI index can *only* be misclassified as having a higher level of torture than actually practiced and, conversely, a country-year that should be classified as a 3 can only be

misclassified as having a lower level of torture than actually practiced. As was true of measurement error in the Banks dataset, errors in the Hathaway and CIRI measures will bias any results towards zero.

An alternative source of bias may also exist in the data. The coding of the Hathaway and CIRI torture measures relies on reports assembled by Amnesty International and the US Department of State. If the expectations of the authors of these reports are affected by CAT signatory status – such that they are inclined to believe that signatories employ low levels of torture – then our results may be biased in favor of concluding that the CAT reduces torture levels.

As with the analyses above, we run a difference-in-differences estimation to assess whether CAT signings are associated with a decline in the prevalence of torture relative to prevailing existing practices $\Delta torture_{i,t} = Probit(\gamma \Delta CATsignatory_{i,t-1} + \Delta \mathbf{X}_{i,t} \beta + \epsilon_{i,t})$ where i denotes country i , t denotes year t , and Δ is the difference operator.⁴⁴ Controls include the change in the military capabilities index, the change in GDP *per capita*, the change in the rate of economic growth, the change in levels of economic openness, and a cubic polynomial of time. These models are run both on the full panel of observations and on a panel consisting only of countries that eventually sign the CAT. The latter estimates rely only on time-series variation amongst CAT signatories – signing the CAT is assumed to affect torture levels in the year after signing and have no effect on the rate of change thereafter. Reports from all specifications are reported in Table 6.

The sign on the lagged change in CAT signatory status is negative in three of the four specifications, and is significant at the 90 percent level when the more fine-grained Hathaway measure is used. The magnitude of the coefficients is fairly consistent across both the panel of signatories and the full dataset. These findings (weakly) suggest that signing the CAT leads to a reduction in levels of torture. Given that the Hathaway index contains more information than the CIRI, it is not surprising that the former measure produces more significant results than the latter.

⁴⁴ $\Delta torture_{i,t}$ is a trichotomous variable coded as 1 if torture levels increase, 0 if they remain the same, and -1 if they decrease.

Table 6: CAT Signing and Torture Levels

	Signatories		Full Dataset	
	Δ Tort. Hathaway	Δ Tort. Ciri	Δ Tort. Hathaway	Δ Tort. Ciri
Lagged CAT Signing	-.503 (.301)	.14 (.277)	-.5 (.275)	-.034 (.257)
change Military Capabilities	-.096 (.205)	.11 (.194)	-.073 (.102)	.111 (.115)
change GDP <i>per capita</i>	51.513 (52.163)	35.232 (61.851)	16.142 (60.291)	58.678 (63.075)
change growth	-.003 (.006)	.0007 (.007)	-.001 (.006)	-.002 (.004)
change openness	-.008 (.004)	-.018 (.005)	-.002 (.005)	-.008 (.004)
Cubic Time Polynomial	✓	✓	✓	✓
N	262	259	514	493

Coefficient estimates from an ordered probit regression of changes in torture levels on changes in CAT signatory status (a difference-in-differences specification). Models to the left are run on a panel of country years containing countries that eventually sign the CAT. Models to the right are run on the full panel of country-years.

Alternative Explanations for the Findings

If the CAT does indeed raise the costs of engaging in torture, consider an alternative story: The treaty acts instead as a mechanism for tying the hands of a repressive government – the signing of the CAT represents a credible commitment by the government to reduce torture in exchange for reduced opposition efforts by the opposition. The government trades away torture in exchange for increased survival in office.

There are two problems with this interpretation. First, as mentioned above, there is no mechanism for the domestic opposition to commit to reducing its anti-government efforts when the government ties its hands with the treaty. What is to prevent the opposition from pressing the newly constrained government in its attempt to oust the incumbent leadership from power? It is improbable that a determined opposition would reduce its efforts to remove the sitting government simply by virtue of the fact the government has agreed to refrain from torture. Second, the hands-tying story does not predict that those governments that torture more heavily in time t are more likely to sign the CAT in time $t + 1$. In such a model, it is the weaker types, the types with the highest costs of torture that have the most to gain from such an arrangement. So

this story would require the weaker types to be most likely to sign the treaty. Instead we find that it is the most severe torturers that sign, the states with the lowest costs of engaging in torture. As noted above, our theoretical model is consistent with the empirical evidence produced by Hathaway (2007) and Vreeland (2008), and with the empirical evidence presented above.

Notice that the argument made here is also a “tying-of-the-hands” story. The treaty raises the costs of violations, and hence is a credible commitment of the tough type’s willingness to hold onto office. In our model, it is the strong types who are willing to bind their hands.

A second alternative explanation questions the utility of the CAT as a signaling device. For instance, do we need the CAT in order for the elite to successfully signal its toughness? A tough autocrat could torture in excessive amounts in early periods, attempting to communicate its toughness to all observers; the domestic opposition might learn it is facing a tough opponent and reduce its efforts accordingly. Hence the torture itself is a credible signal of the government’s type. This would be consistent with the observation that it is the tough types that survive in office longer.

This claim rests on the notion that early period torture levels can separate out the types; but there will be incentives for even weak types to try mimic the strong types in the early periods in order to try to convince the domestic opposition that they are indeed tough. Hence, as in the canonical entry deterrence game, all types would pool on the same signal. Early period torture levels are unlikely to act as a signal that would permit separation of types.

Undoubtedly alternative signaling mechanisms to exist in order for governments to credibly signal their toughness. We do not argue that the CAT is unique in this regard; we do argue that the CAT play this role.

Conclusion

Autocracies that torture more are more likely to sign the CAT than those that torture less and autocracies that sign the CAT continue to torture; though it does appear that the CAT slightly reduces average torture levels in signatory states. Furthermore, those autocrats who sign the CAT survive longer in office than those that do not, and that oppositional activities in signatory states fall when the CAT is signed.

If authoritarian governments use the signing – and violation – of human rights agreements as a signal of their willingness to repress domestic opponents, those regimes that practice repression *ex ante* are most likely to sign. Moreover, in equilibrium, those states that sign continue to torture after signing. The informational effect acts as a threat: signing the treaty signals strength and a willingness to torture if necessary. A rational opposition reduces its political activity in response. In signatory states, declining levels of opposition effort create incentives to reduce the extent of repression. However, the CAT also creates a contrasting incentive to increase torture levels so that the regime is protected from the threat of post-tenure enforcement. Empirically, we find that the former effect dominates. So while the treaty is signed with an intent to defy its provisions, torture levels do fall relative to what they would have been absent signing. They do not, however, go to zero.

While the level of torture in CAT signatories declines, the presence of the CAT causes torture in non-signatories to rise. Moreover, signing the CAT prolongs the tenure in office of the worst torturers relative to the non-signers who are lesser torturers *ex ante*.

What then to make of the CAT? While the CAT may reduce torture in the most autocratic of states, those states sign precisely because they intend to continue torturing. Moreover, those regimes that sign become more secure. The good intentions of the international community may have the unintended consequence of strengthening undemocratic regimes around the world. Additionally, levels of torture rise in non-signatory states. The CAT thus creates a redistribution of torture across authoritarian regimes.

The aggregate welfare implications of the CAT remain to be established. The tenure in office of the worst torturers rises, but the tenure in office falls for less oppressive autocrats; the overall tenure in office of autocrats may rise or fall with the introduction of the CAT. Moreover, we do not explore the implications of the CAT for democratic governments, which may offset some of the negative implications of the behavior of autocrats. We leave further investigation of the aggregate effect of the CAT to future work.

Our model indicates that the CAT presents an opportunity for authoritarian regimes that value office very highly to signal their intent to hold firmly onto that office. The signal associated with

signing the agreement is informative in the sense that signing acts as a credible threat of a willingness to exert effort and incur high costs in order to remain in power. While the treaty designers probably had no intention for the treaty to play this role, hard-core authoritarian regimes appear to have taken advantage of this mechanism. This suggests that these autocrats have an incentive to try to find other mechanisms to signal their type in a credible fashion. While such signaling devices may exist, few are credible in their capacity to separate out types. So long as some degree of uncertainty over governments' willingness to employ repression exists, low-cost types will have an incentive to signal their status. Thus, they will resort to a variety of signaling mechanisms, including – it seems – the signing of the CAT.

The results here offer a response to those that argue for a strengthening of the international human rights regime. A stronger regime means greater penalties for engaging in torture, raising the costs of compliance. Standard accounts suggest that the lower a country's cost of compliance, the greater its probability that it joins an international organization. In this model, the prediction is quite the opposite. As compliance costs rise, the pool of signatories becomes increasingly dominated by those that intend to defy the treaty. For the higher are the costs, the greater the separation of high- from low-cost types, and the more potent the signal sent to the domestic opposition. Higher compliance costs in this case may be more effective at protecting a high-torturing regime.

The findings of this research suggest an under-appreciated element of international institutional design. Agreements that focus on the nation state as a unitary actor, and ignore the effect of the institution on domestic politics – and in particular domestic conflict – may generate unanticipated and adverse effects. Policymakers, engaged in negotiations at the international level over the design of international institutions need to anticipate the effect of these agreements on the domestic polity. Agreements may come into effect exactly because they bolster the political survival of those leaders that sign them. When these leaders are autocratic, it is likely that they will use participation in international agreements to help prevent democratic reform.

How is possible that such a simple model yields such counterintuitive results? We believe this is a consequence of taking two aspects of international politics more seriously. First, domestic political contestation matters when it comes to a state's decision to accede or not to an international obliga-

tion. Second, international institutions generate information (if only by states' accession) that will affect the political calculus of the domestic groups engaged in political competition. By combining the information generated by the international institution and an explicit political contest at the domestic level, we generate results that depart somewhat from the standard canon: countries may accede to treaties they intend to defy.

References

- Achen, Christopher H. & Duncan Snidal. 1989. "Rational Deterrence Theory and Comparative Case Studies." *World Politics* 41(2):143–169.
- Bagwell, Kyle & Robert Staiger. 2005. "Enforcement, Private Political Pressure and the GATT/WTO Escape Clause." *The Journal of Legal Studies* 34(2):471–513.
- Banks, Arthurs S. 1979. "Cross-National Time-Series Data Archive." Center for Social Analysis, State University of New York at Binghamton.
- Box-Seffensmeier, Janet M. & Bradford S. Jones. 2004. *Event History Modeling: A Guide for Social Scientists*. Cambridge University Press.
- Bueno de Mesquita, Bruce, Alastair Smith, Randolph M. Siverson & James D. Morrow. 2003. *The Logic of Political Survival*. Cambridge, MA: The MIT Press.
- Chayes, Abram & Antonia Handler Chayes. 1993. "On Compliance." *International Organization* 47(2):175–205.
- Cingranelli, David L. & David L. Richards. 2007. "The Cingranelli-Richards (CIRI) Human Rights Dataset."
- Clark, Ann Marie & Kathryn Sikkink. 2010. "Information Effects and Human Rights Data: Is the Good News about Increased Human Rights Information Bad News for Human Rights Measures."
- Diamond, Alexis & Jasjeet S. Sekhon. 2006. "Genetic Matching for Estimating Causal Effects: A General Multivariate Matching Method for Achieving Balance in Observational Studies."
- Downs, George W. & David M. Rocke. 1995. *Optimal Imperfection?* Princeton University Press.
- Downs, George W., David M. Rocke & Peter N. Barsoom. 1996. "Is the Good News About Compliance Good News About Cooperation?" *International Organization* 50(3):379–406.
- Fearon, James D. 1994. "Domestic Political Audiences and the Escalation of International Disputes." *The American Political Science Review* 88.3:577–592.
- Finnemore, Martha. 1996. *National Interests in International Society, Cornell Studies in Political Economy*. Cornell University Press.

- Fudenberg, Drew & Jean Tirole. 1991. *Game Theory*. The MIT Press.
- Gandhi, Jennifer & Adam Przeworski. 2007. "Authoritarian Institutions and the Survival of Autocrats." *Comparative Political Studies* 40:1279–1301.
- Gilligan, Michael J. & Nathaniel H. Nesbitt. 2007. "Do Norms Reduce Torture?"
- Goemans, Hein. 2006. Archigos: A Database on Political Leaders. Paper Presented at the Annual Meeting of the American Political Science Association.
- Goodliffe, Jay & Darren G. Hawkins. 2006. "Explaining Commitment: State and the Convention Against Torture." *Journal of Politics* 68:358–371.
- Goodliffe, Jay, Darren Hawkins & James Raymond Vreeland. 2009. "Identity and Norm Diffusion in the Convention Against Torture."
- Hafner-Burton, Emilie M. 2005. "Trading Human Rights: How Preferential Trade Agreements Influence Government Repression." *International Organization* 59(3):593–629.
- Hafner-Burton, Emilie M. & James Ron. 2009. "Seeing Double: Human Rights Impact through Qualitative and Quantitative Eyes." *World Politics* 61:360–401.
- Hafner-Burton, Emilie M. & Kiyoteru Tsutsui. 2005. "Human Rights in a Globalizing World: The Paradox of Empty Promises." *American Journal of Sociology* 110(5):1373–1411.
- Hafner-Burton, Emilie M. & Kiyoteru Tsutsui. 2007. "Justice Lost! The Failure of International Human Rights Law to Matter Where Its Needed Most." *Journal of Peace Research* 44(4):407–425.
- Hathaway, Oona. 2007. "Why Do Countries Commit to Human Rights Treaties?" *The Journal of Conflict Resolution* 51(4):588–621.
- Hirschleifer, Jack. 1991. "The Paradox of Power." *Economics and Politics* 3:177–200.
- Ho, Daniel E., Kosuke Imai, Gary King & Elizabeth A. Stuart. 2007. "Matching as Preprocessing for Reducing Model Dependence in Parametric Causal Inference." *Political Analysis* (forthcoming).
- Ho, Daniel, Kosuke Imai, Gary King & Elizabeth Stuart. 2004. "MatchIt: Matching as Nonparametric Preprocessing for Parametric Causal Inference." <http://gking.harvard.edu/matchit/>.
- ICJ. 2009. Questions Relating to the Obligation to Prosecute or Extradite (Belgium v. Senegal). Summary of Order 2009/3 International Court of Justice.
- Jonas, Stacie. 2004. "The Ripple Effect of the Pinochet Case." *Human Rights Brief* 11(3):36–38.
- Keck, Margaret E. & Kathryn Sikkink. 1998. *Activists Beyond Borders: Advocacy Networks in International Politics*. Cornell University Press.
- Klein, David F. 1988. "A Theory for the Application of the Customary International Law of Human Rights by Domestic Courts." *Yale Journal of International Law* 13:332–365.

- Koremenos, Barbara. 2005. "Contracting Around International Uncertainty." *The American Political Science Review* 99:549–565.
- Lacina, Bethany & Nils Petter Gleditsch. 2005. "Monitoring Trends in Global Combat: A New Dataset of Battle Deaths." *European Journal of Population* 21:145–166.
- Moravcsik, Andrew. 2000. "The Origins of Human Rights Regimes: Democratic Delegation in Postwar Europe." *International Organization* 54(2):217–252.
- Neumayer, Eric. 2005. "Do International Human Rights Treaties Improve Respect for Human Rights?" *The Journal of Conflict Resolution* 49(6):925–953.
- Nielsen, Richard & Beth A. Simmons. 2009. "Rewards for Rights Ratification? Testing for Tangible and Intangible Benefits of Human Rights Treaty Ratification."
- Powell, Emilia Justyna & Jeffrey K. Stanton. 2009. "Domestic Judicial Institutions and Human Rights Treaty Violation." *International Studies Quarterly* 53(1):149–174.
- Przeworski, Adam, Michael E. Alvarez, José Antonio Cheibub & Fernando Limongi. 2000. *Democracy and Development: Political Institutions and Well-Being in the World, 1950-1990*. Cambridge University Press.
- Roht-Arriaza, Naomi. 2001. "The Pinochet Precedent and Universal Jurisdiction." *New England Law Review* 35(2):311–320.
- Rosendorff, B. Peter. 2005. "Stability and Rigidity: Politics and the Design of the WTO's Dispute Resolution Procedure." *The American Political Science Review* 99(3):389–400.
- Rosendorff, B. Peter & Helen V. Milner. 2001. "The Optimal Design of International Trade Institutions: Uncertainty and Escape." *International Organization* 55(4):829–857.
- Simmons, Beth A. 2009. *Mobilizing for Human Rights: International Law in Domestic Politics*. Princeton University Press.
- Simmons, Beth A. & Allison Danner. 2010. "Credible Commitments and the International Criminal Court." *International Organization* 64:225–256.
- Simmons, Beth A. & Daniel J. Hopkins. 2005. "The Constraining Power of International Treaties: Theory and Methods." *The American Political Science Review* 99(4):623–631.
- Singer, J. David. 1987. "Reconstructing the Correlates of War Dataset on Material Capabilities of State, 1816-1985." *International Interactions* 14:115–132.
- Skaperdas, Stergios. 1996. "Contest Success Functions." *Economic Theory* 7:283–290.
- Smith, Alastair. 1998. "International Crises and Domestic Conflicts." *The American Political Science Review* 92.3:623–638.
- Tomz, Michael, Jason Wittenberg & Gary King. 2001. "CLARIFY: Software for Interpreting and Presenting Statistical Results. Version 2.0." <http://gking.harvard.edu>. Cambridge, MA: Harvard University.

Vreeland, James Raymond. 2008. "Political Institutions and Human Rights: Why Dictatorships Enter into the United Nations Convention Against Torture." *International Organization* 62:65–101.

Wood, Reed M. & Mark Gibney. 2010. "The Political Terror Scale (PTS): A Re-introduction and a Comparison to CIRI." *Human Rights Quarterly* 32:367–400.

Appendix: Proofs

A perfect Bayesian equilibrium of a signaling game consists of a strategy profile and a system of beliefs such that (1) the sender chooses her strategy to maximize her utility subject to the receiver's strategy; (2) the receiver chooses her strategy to maximize her utility subject both to the sender's strategy and to her beliefs conditional upon the sender's message; and (3) the receiver's beliefs are updated according to Bayes' Rule, whenever possible (Fudenberg & Tirole 1991).

Definitions:

1. Define a pair of strategies $\{(s, t), e\}$ where $s : [0, 1] \rightarrow \{0, 1\}$, $t : \{0, 1\} \times [0, 1] \rightarrow \mathbb{R}_+$, $e : \{0, 1\} \rightarrow \mathbb{R}_+$.
2. Define the functions $\Psi(\tilde{R}) = \frac{C(\sqrt{1+C+\bar{P}} - \sqrt{\tilde{R}+C+\bar{P}})}{1-\tilde{R}}$ and $\Upsilon(\tilde{R}) = \frac{C(\sqrt{\tilde{R}+C} - \sqrt{C})}{\tilde{R}}$.
3. Define $\underline{P} = \frac{8C+9}{16C+16}$ and $\bar{P} = \frac{C+4-2C^2}{4(C-1)} + \frac{1}{2}\sqrt{\frac{C^4+C^3}{(C-1)^2}}$

Proof of Proposition 1:

The proposition states that if $\underline{P} < P < \bar{P}$, then there exists a unique semi-separating equilibrium:

$$\text{There exists } \tilde{R} \in (0, 1) \text{ such that } s(R) = \begin{cases} 1 & \text{if } R \geq \tilde{R} \\ 0 & \text{if } R < \tilde{R} \end{cases} \text{ and } e(s) = \begin{cases} \Psi(\tilde{R})^2 & \text{if } s = 1 \\ \Upsilon(\tilde{R})^2 & \text{if } s = 0 \end{cases} \text{ and}$$

$$t(s, R) = \begin{cases} \Psi(\tilde{R})\sqrt{R+C+\bar{P}} - [\Psi(\tilde{R})]^2 & \text{if } R \geq \tilde{R} \\ \Upsilon(\tilde{R})\sqrt{R+C} - [\Upsilon(\tilde{R})]^2 & \text{if } R < \tilde{R} \end{cases}, \text{ with beliefs } f(1) = \begin{cases} \frac{1}{1-\tilde{R}} & \text{if } R \in [\tilde{R}, 1] \\ 0 & \text{otherwise} \end{cases}$$

$$\text{and } f(0) = \begin{cases} \frac{1}{\tilde{R}} & \text{if } R \in [0, \tilde{R}) \\ 0 & \text{otherwise} \end{cases}.$$

We prove this first by checking, given any signal, that each player is playing a best response.

Then we specify the conditions for the threshold type \tilde{R} to be interior to the type space and establish

that this cut point is unique. Finally, we demonstrate that no signatory government with type $R \geq \tilde{R}$ has any incentive to deviate from the equilibrium by not signing the CAT, nor does any government with type $R < \tilde{R}$ have any incentive to deviate from this equilibrium by signing the CAT.

First, suppose $s = 1$, $\frac{\partial U_G}{\partial t} = \pi_t(t, e)[R + C + P] - 1 = 0$ yields a reaction function $t(e, 1, R) = \sqrt{e[R + C + P]} - e \Rightarrow t(\frac{C^2(\sqrt{1+C+P}-\sqrt{\tilde{R}+C+P})^2}{(1-\tilde{R})^2}, 1, R) = \Psi(\tilde{R})\sqrt{R + C + P} - [\Psi(\tilde{R})]^2$. Therefore, autocratic signatory governments are playing a best response to their opposition when $t(1, R) = \Psi(\tilde{R})\sqrt{R + C + P} - [\Psi(\tilde{R})]^2$.

Suppose $s = 0$, $\frac{\partial U_G}{\partial t} = \frac{e}{(t+e)^2}[R + C] - 1 = 0$ yields a reaction function $t(e, 0, R) = \sqrt{e[R + C]} - e \Rightarrow t(\frac{C^2(\sqrt{\tilde{R}+C}-\sqrt{C})^2}{\tilde{R}^2}, 0, R) = \Upsilon(\tilde{R})\sqrt{R + C} - [\Upsilon(\tilde{R})]^2$. Therefore, non-signatory autocratic governments are playing a best response to their opposition when $t(0, R) = \Upsilon(\tilde{R})\sqrt{R + C} - [\Upsilon(\tilde{R})]^2$.

The opposition's problem, if the opposition observes $s = 1$ is to maximize $EU_D = \int_{\tilde{R}}^1 [(1 - \pi(t, e))C - e]f(1)dR$, where $f(1)$ is the posterior distribution of R over updated support $[\tilde{R}, 1]$, conditional on signal $s = 1$.

Substituting in from the response function above, and recalling that $\pi(t, e) = \frac{t}{t+e}$, one can derive the opposition's expected utility:

$$\begin{aligned} EU_D &= C \int_{\tilde{R}}^1 \frac{e}{t^1 + e} \frac{1}{1 - \tilde{R}} dR - e \\ &= \frac{2C\sqrt{e}(\sqrt{1+C+P} - \sqrt{\tilde{R}+C+P})}{1 - \tilde{R}} - e \end{aligned}$$

This can then be used to derive the opposition's response function:

$$\begin{aligned} \frac{\partial U_D}{\partial e} &= \frac{C(\sqrt{1+C+P} - \sqrt{\tilde{R}+C+P})}{(1 - \tilde{R})\sqrt{e}} - 1 = 0 \\ \Leftrightarrow e(1) &= \frac{C^2(\sqrt{1+C+P} - \sqrt{\tilde{R}+C+P})^2}{(1 - \tilde{R})^2} = [\Psi(\tilde{R})]^2 \end{aligned}$$

Similarly, if the opposition observes $s = 0$, $EU_D = \int_0^{\tilde{R}} [(1 - \pi(t(e, 0), e))C - e]f(0)dR$, where $f(0)$ is the posterior distribution of R over updated support $[0, \tilde{R}]$, conditional on signal $s = 0$.

Substituting in from the response function above, and recalling that $\pi(t, e) = \frac{t}{t+e}$, one can

derive the opposition's expected utility:

$$\begin{aligned} EU_D &= C \int_0^{\tilde{R}} \frac{e}{\sqrt{e^0[R+C]}} \frac{1}{\tilde{R}} dR - e \\ &= \frac{2C\sqrt{e}(\sqrt{\tilde{R}+C} - \sqrt{C})}{\tilde{R}} - e \end{aligned}$$

Therefore,

$$\begin{aligned} \frac{\partial U_D}{\partial t} &= \frac{C(\sqrt{\tilde{R}+C} - \sqrt{C})}{\tilde{R}\sqrt{e}} - 1 = 0 \\ \Leftrightarrow e(0) &= \frac{C^2(\sqrt{\tilde{R}+C} - \sqrt{C})^2}{\tilde{R}^2} = [\Upsilon(\tilde{R})]^2 \end{aligned}$$

In both cases the opposition is playing a best response to the government's action in the equilibrium specified.

Second, we need to demonstrate that a unique \tilde{R} exists in the $[0, 1]$ interval that defines the type space. To do this, we first note that the utility for a signatory government is $[\sqrt{R+C+P} - \Psi(\tilde{R})]^2 - P$; while the utility for a non-signatory government is $[\sqrt{R+C} - \Upsilon(\tilde{R})]^2$. We then define $\Delta(R) = [\sqrt{R+C+P} - \Psi(\tilde{R})]^2 - P - [\sqrt{R+C} - \Upsilon(\tilde{R})]^2$. The function $\Delta(\cdot)$ returns the difference in the government's utility from signing the CAT and its utility from not signing the CAT. At the threshold value \tilde{R} , this function must equal zero: $\Delta(\tilde{R}) = 0$.

Substitute the expressions for $\Psi(\tilde{R})$ and $\Upsilon(\tilde{R})$ into the function $\Delta(R)$:

$$\begin{aligned} \Delta(R) &= \frac{C(\sqrt{1+C+P} - \sqrt{\tilde{R}+C+P})}{1 - \tilde{R}} \left[\frac{C(\sqrt{1+C+P} - \sqrt{\tilde{R}+C+P})}{1 - \tilde{R}} - 2\sqrt{R+C+P} \right] - \\ &\quad \frac{C(\sqrt{\tilde{R}+C} - \sqrt{C})}{\tilde{R}} \left[\frac{C(\sqrt{\tilde{R}+C} - \sqrt{C})}{\tilde{R}} - 2\sqrt{R+C} \right] \\ \Leftrightarrow \Delta(\tilde{R}) &= \frac{C^2(\sqrt{1+C+P} - \sqrt{\tilde{R}+C+P})^2}{(1 - \tilde{R})^2} - \frac{2C(\sqrt{\tilde{R}+C+P}\sqrt{1+C+P} - \tilde{R} - C - P)}{1 - \tilde{R}} - \\ &\quad \frac{C^2(\sqrt{\tilde{R}+C} - \sqrt{C})^2}{\tilde{R}^2} + \frac{2C(\tilde{R}+C - \sqrt{\tilde{R}+C}\sqrt{C})}{\tilde{R}} \end{aligned}$$

We can then take the limit of $\Delta(\tilde{R})$ as $\tilde{R} \rightarrow 1$.

$$\begin{aligned} \lim_{\tilde{R} \rightarrow 1} \Delta(\tilde{R}) &= \frac{C^2}{4(1+C+P)} - C - C^2(\sqrt{1+C} - \sqrt{C})^2 + 2C(1+C - \sqrt{1+C}\sqrt{C}) \\ &< 0 \text{ iff } C > 0 \text{ and } P > \frac{8C+9}{16C+16} = \underline{P} \end{aligned}$$

And, we take the limit of $\Delta(\tilde{R})$ as $\tilde{R} \rightarrow 0$.

$$\begin{aligned} \lim_{\tilde{R} \rightarrow 0} \Delta(\tilde{R}) &= C^2(\sqrt{1+C+P} - \sqrt{C+P})^2 - 2C(\sqrt{C+P}\sqrt{1+C+P} - C - P) - \frac{C}{4} + C \\ &> 0 \text{ iff } C > 1 \text{ and } -C - 1 < P < \frac{C+4-2C^2}{4(C-1)} + \frac{1}{2}\sqrt{\frac{C^4+C^3}{(C-1)^2}} = \bar{P} \end{aligned}$$

Thus, by the intermediate value theorem, there exists a value of $\tilde{R} \in [0, 1]$ such that $\Delta(\tilde{R}) = 0$ if $C > 1$ and $\underline{P} < P < \bar{P}$. We graph the C, P parameter space such that these two conditions hold below.

It remains to be shown that the value of \tilde{R} such that $\Delta(\tilde{R}) = 0$ is unique. Note that \tilde{R} will be unique if $\Delta(R)$ is monotonic in R .

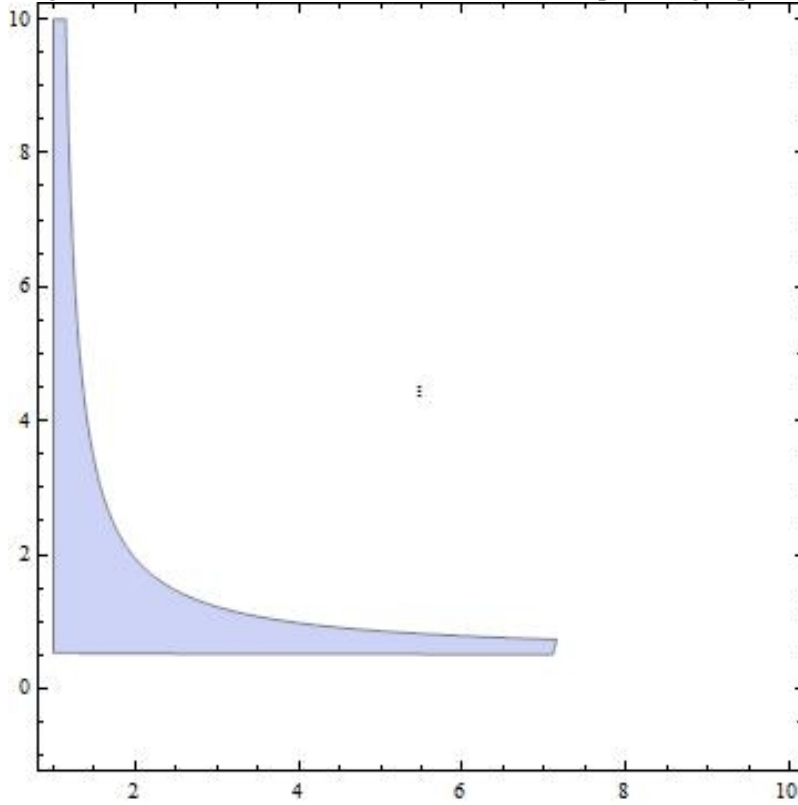
$$\begin{aligned} \frac{\partial \Delta(R)}{\partial R} &= \frac{[\sqrt{R+C+P} - \Psi(\tilde{R})]}{\sqrt{R+C+P}} - \frac{[\sqrt{R+C} - \Upsilon(\tilde{R})]}{\sqrt{R+C}} = \frac{\Upsilon(\tilde{R})}{\sqrt{R+C}} - \frac{\Psi(\tilde{R})}{\sqrt{R+C+P}} \\ \text{Therefore } \frac{\partial \Delta(R)}{\partial R} &> 0 \Leftrightarrow \frac{\sqrt{R+C+P}}{\sqrt{R+C}} > \frac{\Psi(\tilde{R})}{\Upsilon(\tilde{R})}. \end{aligned}$$

$$\text{Note that: } \Psi(\tilde{R}) = CE\left[\frac{1}{\sqrt{R+C+P}} \mid R \geq \tilde{R}\right] \text{ and: } \Upsilon(\tilde{R}) = CE\left[\frac{1}{\sqrt{R+C}} \mid R < \tilde{R}\right]$$

$$\text{Therefore, } \Upsilon(\tilde{R}) \geq \Psi(\tilde{R}) \text{ which implies } \frac{\sqrt{R+C+P}}{\sqrt{R+C}} > \frac{\Psi(\tilde{R})}{\Upsilon(\tilde{R})} \forall P > 0 \text{ and } \frac{\partial \Delta(R)}{\partial R} > 0 \forall R.$$

Third, we must demonstrate that no signatories in this equilibrium wish to deviate by not signing the CAT (setting $s = 0$) and that no non-signatories wish to deviate by signing (setting $s = 1$). To see this, recall that the function $\Delta(R)$ returns the difference in utility for a given government from signing versus no signing. As demonstrated above, this function is monotonic and increasing in R . Thus, for all governments with $R < \tilde{R}$ the value of $\Delta(R)$ is negative – implying that they receive greater utility from not signing than from signing. Thus no such government has

Figure 3: Range of Parameter Values Such that a Semi-Separating Equilibrium Exists



A graph of the parameter space in which the conditions $C > 1$ and $\frac{8C+9}{16C+16} < P < \frac{C+4-2C^2}{4(C-1)} + \frac{1}{2}\sqrt{\frac{C^4+C^3}{(C-1)^2}}$. Values of P are reported on the y-axis, values of C are reported on the x-axis. The range of parameter values where both conditions hold is shaded.

an incentive to deviate from the equilibrium.

Note further that all governments with $R \geq \tilde{R}$ have values of $\Delta(R) \geq 0$. This implies that they receive (weakly) greater utility from signing the CAT than from not signing. Thus, no government with $R \geq \tilde{R}$ has an incentive to deviate from the semi-separating equilibrium by not signing the treaty. \square

Proof of Proposition 2:

Recall, from the above, that the equilibrium level of torture exerted by a non-signatory government is given by

$$t\left(\frac{C^2(\sqrt{\tilde{R}+C}-\sqrt{C})^2}{\tilde{R}^2}, 0, R\right) = \Upsilon(\tilde{R})\sqrt{R+C} - [\Upsilon(\tilde{R})]^2.$$

And the equilibrium level of effort exerted by the opposition against a non-signatory government is given by $[\Upsilon(\tilde{R})]^2$. Substituting these two values into the contest success function $\pi(t, e) = \frac{t}{t+e}$ yields $1 - \frac{\Upsilon(\tilde{R})}{\sqrt{R+C}}$.

Similarly, the equilibrium level of torture exerted by a signatory government is given by

$$t\left(\frac{C^2(\sqrt{1+C+P} - \sqrt{\tilde{R}+C+P})^2}{(1-\tilde{R})^2}, 1, R\right) = \Psi(\tilde{R})\sqrt{R+C+P} - [\Psi(\tilde{R})]^2.$$

And the equilibrium level of effort exerted by the opposition against a signatory government is given by $[\Psi(\tilde{R})]^2$. Substituting these two values into the $\pi(t, e)$ yields $1 - \frac{\Psi(\tilde{R})}{\sqrt{R+C+P}}$. From above, $\Upsilon(\tilde{R}) > \Psi(\tilde{R})$ in any semi-separating equilibrium. $\frac{1}{\sqrt{R+C+P}} < \frac{1}{\sqrt{R+C}} \forall P > 0$. Therefore, $\frac{\Psi(\tilde{R})}{\sqrt{R+C+P}} < \frac{\Upsilon(\tilde{R})}{\sqrt{R+C}}$. Thus, signatories survive in office longer than non-signatories. \square

Proof of Proposition 3:

We first must determine the equilibrium behavior of the government and opposition absent the institution. Then the utility of the government is $U_G = \pi(t, e)[R+C] - t$, yielding the government's response function $t = \sqrt{e[R+C]} - e$.

The opposition does not receive any signal, and thus must make its best guess regarding government behavior based on its prior that $R \sim U[0, 1]$. $EU_D = E[1 - \pi(t, e)]C - e = C \int_0^1 \frac{e}{t+e} dR - e$. Substituting in the government's best response function for t yields $EU_D = 2C\sqrt{e}(\sqrt{1+C} - \sqrt{C}) - e$, which produces the opposition's best response function $e = C^2(\sqrt{1+C} - \sqrt{C})^2$.

Substituting the opposition's best response back into the government's response function yields $t = \sqrt{e[R+C]} - e = C(\sqrt{1+C} - \sqrt{C})\sqrt{R+C} - C^2(\sqrt{1+C} - \sqrt{C})$.

Using this value for the level of government torture absent the institution, and the value for the level of torture amongst non-signatories derived above $(\Upsilon(\tilde{R})\sqrt{R+C} - [\Upsilon(\tilde{R})]^2)$ we can state that torture increases in non-signatory states (after substituting for the definition of $\Upsilon(\tilde{R})$) iff $C(\sqrt{1+C} - \sqrt{C})\sqrt{R+C} - C^2(\sqrt{1+C} - \sqrt{C})^2 < \frac{C(\sqrt{\tilde{R}+C}-\sqrt{C})}{\tilde{R}}\sqrt{R+C} - \frac{C^2(\sqrt{\tilde{R}+C}-\sqrt{C})^2}{\tilde{R}^2}$.

Note that $\frac{C(\sqrt{\tilde{R}+C}-\sqrt{C})}{\tilde{R}} = CE[\frac{1}{\sqrt{R+C}} | R < \tilde{R}]$ and $C(\sqrt{1+C} - \sqrt{C}) = CE[\frac{1}{\sqrt{R+C}} | R \in [0, 1]]$ and that $CE[\frac{1}{\sqrt{R+C}} | R < \tilde{R}] > CE[\frac{1}{\sqrt{R+C}} | R \in [0, 1]]$. Note too that $C(\sqrt{1+C} - \sqrt{C})\sqrt{R+C} < \frac{C(\sqrt{\tilde{R}+C}-\sqrt{C})}{\tilde{R}}\sqrt{R+C} \forall C > 0$. Then $C(\sqrt{1+C} - \sqrt{C})\sqrt{R+C} - CE[\frac{1}{\sqrt{R+C}} | R \in [0, 1]] < \frac{C(\sqrt{\tilde{R}+C}-\sqrt{C})}{\tilde{R}}\sqrt{R+C} - CE[\frac{1}{\sqrt{R+C}} | R < \tilde{R}] \forall C > 0$ which establishes the result.

It remains to be shown that opposition effort levels increase in non-signatory states relative to a world without the CAT. Note from the response functions derived above this will be true iff $C^2(\sqrt{1+C} - \sqrt{C})^2 < [\Upsilon(\tilde{R})]^2$. Recall that $C^2(\sqrt{1+C} - \sqrt{C})^2 = (CE[\frac{1}{R+C}|R \in [0,1]])^2$ and $[\Upsilon(\tilde{R})]^2 = (CE[\frac{1}{\sqrt{R+C}}|R < \tilde{R}])^2$. Since $E[\frac{1}{R+C}|R \in [0,1]] < E[\frac{1}{\sqrt{R+C}}|R < \tilde{R}]$, it must be true that $C^2(\sqrt{1+C} - \sqrt{C})^2 < [\Upsilon(\tilde{R})]^2$. Thus, opposition effort levels rise in non-signatory states relative to a world with no CAT. \square

Proof of Proposition 4:

Absent the CAT, equilibrium levels of opposition effort are $C^2(\sqrt{1+C} - \sqrt{C})^2 = (E[\frac{1}{\sqrt{R+C}}|R \in [0,1]])^2$. In the presence of the CAT, opposition effort levels in signatory states are $[\Psi(\tilde{R})]^2 = (E[\frac{1}{\sqrt{R+C+P}}|R \geq \tilde{R}])^2$. Note that $E[\frac{1}{\sqrt{R+C}}|R \in [0,1]] > E[\frac{1}{\sqrt{R+C+P}}|R \geq \tilde{R}]$. Thus, $C^2(\sqrt{1+C} - \sqrt{C})^2 > [\Psi(\tilde{R})]^2$, opposition efforts decline on signing the CAT. \square