

Torture in Counterterrorism: Agency Incentives and Slippery Slopes

Hugo M. Mialon, Sue H. Mialon, and Maxwell B. Stinchcombe¹

March 21, 2011

Abstract

We develop a counterterrorism model to analyze the effects of allowing a government agency to torture suspects when evidence of terrorist involvement is strong. We find that legalizing torture in strong-evidence cases has offsetting effects on agency incentives to counter terrorism by means other than torture. It increases these incentives because other efforts may increase the probability of having strong enough evidence to warrant the use of torture if other efforts fail—a complementarity effect. However, it also lowers these incentives because the agency might come to rely on torture to avert attacks—a decommitment effect. The decommitment effect is more likely to dominate if the agency’s non-torture efforts are good at stopping attacks. Moreover, legalizing torture in strong-evidence cases is likely to reduce security if the effectiveness of torture is low while that of non-torture efforts is high. Lastly, we find that legalizing torture in strong-evidence cases can increase agency incentives to torture even in weak-evidence cases—a slippery slope. (JEL H1, K4)

¹ Hugo Mialon and Sue Mialon, Department of Economics, Emory University, Atlanta GA 30322-2240 (hmialon@emory.edu and smialon@emory.edu). Maxwell Stinchcombe, Department of Economics, University of Texas at Austin, Austin TX 78712-1173 (maxwell@eco.utexas.edu). We are extremely grateful to the Editor, two anonymous referees, Satyajit Chatterjee, Martin Dumav, Amy Farmer, Andrew Francis, Louis Kaplow, Tilman Klumpp, Preston McAfee, Tim Mathews, Erik Nesson, Josh Robinson, Paul Rubin, Louis Seidman, Steve Shavell, Jeroen Swinkels, Kevin Tsui, Adam Winship, Kathy Zeiler, and seminar participants at Clemson, Emory, Harvard, University of Texas, and the 2009 SEA Meeting for helpful comments.

1. Introduction

Since the September 11, 2001 terrorist attacks on the U.S., a number of legal authorities, political authorities, and poll results have favored the use of torture in counterterrorism. Many arguments in favor of the use of torture begin with, or prominently feature, some variant of a “ticking bomb” scenario in which torturing one suspect leads, with near certainty, to saving many lives. However, if a ticking bomb scenario arises, it is because other preventive efforts have failed. In this paper, we analyze the effects of legalizing torture that arise through substitution within the portfolio of counterterrorism efforts.

The Push for Torture. In 2002, as Assistant Attorney General in the U.S. Department of Justice, John Yoo recommended that the White House withdraw its recognition of the rules prohibiting torture imposed by the Geneva Conventions.² In an influential book, Dershowitz (2002) proposed controlling torture through a system of judicial warrants for the use of torture in limited circumstances. In 2004, Senator Charles Schumer publicly rejected the idea that torture should never be used.³ In 2007, Attorney General Michael Mukasey refused to state that waterboarding is illegal.⁴ In 1987, the Supreme Court of Israel consented to the use of torture to stop attacks by Palestinian terrorists (Imseis, 2001, and Strauss, 2003). In a 1999 decision, the Court returned to a ban on torture, but the ban is not absolute and has not consistently been enforced as discussed by Imseis (2001). The Economist (2006) reported results of public opinion polls about torture. In a BBC poll of 27,000 people in 25 countries, 33 percent of people polled, 36 percent of Americans polled, and 46 percent of Israelis polled said that “some degree of torture is permissible.” In a recent

² Memorandum from John Yoo, Deputy Assistant Attorney General, and Robert J. Delahunty, Special Counsel, to William J. Haynes II, General Counsel, Department of Defense (Jan. 9, 2002).

³ Federal Government’s Counterterrorism Efforts: Hearing Before the Sen. Judiciary Subcommittee, 108th Cong. (2004) (statement of Sen. Charles Schumer, Member, S. Judiciary Committee).

⁴ Executive Nomination: Hearing on the Nomination of Michael B. Mukasey To Be Attorney General of the United States Before the S. Comm. on the Judiciary, 110th Cong. (2007). See Egan (2007)

poll of 2,000 people living in the U.S., which was conducted by The Pew Research Foundation (2009) between October and November of 2009, 54 percent of people polled thought that “torturing terrorist suspects is often or sometimes justified.”

The Ticking Bomb Scenario. Most arguments in favor of legalizing torture are presented as variants of the following mass terrorism scenario.⁵ Suppose the government learns that terrorists are planning an attack in a populated area. If the attack succeeds, many people could die. The government has in custody a suspect who may know about the attack but refuses to cooperate. The government can perhaps force the suspect to reveal what he knows through torture. The suspect could survive the torture, and the information extracted could save many lives. The government does not, at that time, have other means of stopping the attack. In some versions of this scenario, cost-benefit analysis suggests that allowing the government to torture the suspect is the socially efficient policy.

In this paper, we examine this efficiency argument more carefully. There are unintended consequences to allowing torture, and those consequences vary with the degree of evidence required to make torture legal. Specifically, if a ticking bomb scenario arises, it reflects a failure of other preventive efforts.⁶ Assuming that situations can arise in which the cost-benefit calculations described above favor torture, we develop a model of counterterrorism to evaluate the overall effects of legalizing torture. Throughout, we maintain the assumption that there is a fundamental agency problem: the counterterrorism agency puts, relative to society, more weight on preventing terrorist attacks than on protecting individual rights.

When The Agency Obeys Directives. In the first part of our analysis, we assume the agency obeys directives on torture. In this case, allowing the agency to torture when

⁵ See e.g. Dershowitz (2002), Strauss (2003), Luban (2005), and Bagaric and Clarke (2006).

⁶ For domestically based attacks, such efforts include tracking of materials used in bomb-making, restrictions on bomb-making activities, increased security at likely targets, and baggage and cargo screening at airports. For extraterritorial sources, such efforts include hiring, training and paying attention to analysts fluent in language and culture, cultivating allies, and bilateral or multilateral cooperative international policing.

evidence of guilt is strong has two opposing effects on the agency's incentives to counter terrorism by means other than torture; first, it tends to reduce these incentives because it ameliorates a situation in which other efforts have failed—a decommitment effect; second, it tends to increase these incentives because other efforts improve the chances of gaining enough evidence to warrant the use of torture if other efforts fail—a complementarity effect.

The decommitment effect is more likely to dominate if the agency's non-torture efforts are good at stopping attacks and the agency is accountable for damages from successful attacks. Furthermore, allowing torture in a broader range of cases lowers the complementarity effect. When the decommitment effect dominates, legalizing torture reduces the agency's other preventive efforts. In this case, we have a formalization of the observation in Rejali (2007) that reliance on torture typically makes an agency sloppier in its other preventive work and leads to agency "deskilling." In the longer run, it might reduce investment in the development of alternative technologies and prevention techniques.⁷

If legalizing torture reduces the agency's preventive effort, it can reduce security. In particular, it is likely to reduce security if the effectiveness of torture is low while that of non-torture efforts is high. Moreover, it reduces welfare if torture is sufficiently ineffective

⁷ Rejali (*op. cit.* Chap 21 and 22) documents the deskilling effect across a large number of 19'th and 20'th century instances. For example, there is evidence that the Gestapo's suppression of the Resistance in World War II was far more effective when it relied on informers and careful interrogation before it turned extensively and "unprofessionally" to torture. In the French-Algerian war, French army units that tortured became insubordinate to central army authority and abandoned basic police techniques – in one instance, going directly to torture rather than checking the personal effects of an apprehended suspect allowed an Algerian resistance bomb factory to be safely relocated. The radical part of the Algerian resistance movement followed a policy of identifying members of the moderate opposition when tortured, and because the French army did not check the veracity of what was revealed under torture, it wiped out the moderate opposition.

For a more recent example, when Abdul Hakim Murad was arrested by Filipino police in 1995 with fake passports, bomb-making materials, and an encrypted computer, police tortured him instead of trying to decrypt the computer. He revealed little specific information under torture, but when the CIA decrypted his computer years later, it revealed detailed information about Al Qaeda plots to blow up planes in the US, down to specific procedures and flight schedules.

The CIA's unedited *Human Resources Exploitation Training* manual summarizes the deskilling effect of torture with "The routine use of torture ... corrupts those that rely on it as the quick and easy way out" (See the National Security Archives at website <http://www2.gwu.edu/~nsarchiv/NSAEBB/NSAEBB122/>).

and the costs of torturing the innocent are sufficiently high.

The Enforcement Problem. The core of the Dershowitz (2002) argument is that agencies do not obey directives, that torture happens even though it is illegal, and that an enforced system of judicial warrants could bring this under control, resulting in less risk of torturing the innocent. In the second part of our analysis, we extend our model to encompass the possibility that the agency is willing to disobey directives on torture at the risk of legal sanction.⁸ In the extended model, the agency can choose whether to use torture even when torture is not allowed, and it faces potential punishment if it uses torture illegally. If torture is legal in strong-evidence cases, the agency only faces potential punishment if it uses torture in cases where evidence is weak.

In this extended context, there are conditions under which the agency's optimal torture policy is to use torture in strong-evidence cases whether or not torture is legal in such cases. We find that, if these conditions hold, then legalizing torture in strong-evidence cases has the previous two effects, decommitment and complementarity, as well as a third, more subtle, decomplementarity effect, on the agency's non-torture efforts. If the agency's punishment for using torture in strong-evidence cases is removed, it has less incentive to invest in preventive effort in order to avoid having to use torture and thus subjecting itself to punishment (decommitment). However, it also has greater incentive to have strong evidence when it uses torture, so it can avoid the punishment. If attacks are frequent so any given suspect is more likely to be guilty, then the agency is more likely to get strong evidence by increasing its preventive effort (complementarity). But if attacks are infrequent, the agency is more likely to get strong evidence by reducing its effort (decomplementarity).

If the complementarity effect dominates, legalizing torture in strong-evidence cases in-

⁸ In wartime and on foreign grounds, the risk that torture will be used illegally may be high, as evidenced by the documented reports of sadistic torture at the Abu Ghraib and Guantanamo Bay detention facilities of the U.S. Military (Fay, 2004).

creases the agency's other efforts, which increases the accuracy of the agency's evidence, thereby reducing the probability that an innocent person is tortured. This case supports the Dershowitz argument that an open warrant system might actually increase incentives to obey the law and reduce torture of the innocent. However, if the decommitment and decomplementarity effects dominate, legalizing torture in strong-evidence cases reduces the agency's other efforts and increases the probability that an innocent person is tortured. The decommitment and decomplementarity effects dominate if the agency's non-torture efforts are good at stopping attacks and attacks are infrequent.

Slippery Slopes. We also find that legalizing torture in strong-evidence cases can lead to an increase in its use in other cases—a slippery slope. Intuitively, this involves the endogeneity of the quality of information. If the decommitment and decomplementarity effects dominate, legalizing torture in strong-evidence cases reduces the agency's efforts to counter terrorism by means other than torture, which in turn reduces the quality of the information on which the agency bases its decision to use torture if the other efforts fail, increasing agency incentives to use torture in other cases.

This mechanism differs from the three basic variants of the slippery slope arguments that we have found in the literature: utility change, cost change, and somewhat related bureaucratic structure arguments. Volokh (2003, p. 1077) elegantly summarizes the utility change arguments as “the normative power of the actual.” In more pedestrian language, a society that allows torture will perhaps come to see nothing wrong with it. Volokh (2003) and Rizzo and Whitman (2004) provide a detailed examination of cost-based slippery slope mechanisms in legal policymaking. These involve one decision lowering the cost to make another linked decision. In our context, if society pays the cost of training and supporting professional torturers, then the lower marginal cost of torture can lead to more frequent

torture. Posner (2002, p. 30) argues that if "... rules are promulgated permitting torture in defined circumstances, some officials are bound to want to explore the outer bounds of the rules. Having been regularized, the practice will become regular." Sobel (2000) provides a model of declining standards that may bear on the worry that any chosen evidence standard for torture may be prone to slip over time, perhaps by the accretion of precedents set by judges who are more sympathetic to agency arguments for torture.

Nonetheless, it is possible that a legal standard for torture would not slip if torture were legalized, just as the legal standard for capital punishment does not seem to have slipped after capital punishment was legalized (Bagaric and Clarke, 2006). However, according to the slippery slope mechanism that we identify here, even if the legal standard for torture were not to slip, legalizing torture in certain circumstances could still entail a slippery slope in that it could increase illegal torture in other circumstances.

Models of Torture in the Literature. There is a growing literature on the economics of terrorism (see Enders and Sandler, 2004 and 2005, Sandler and Siqueira, 2006, Siqueira and Sandler, 2007, Garoupa, Klick, and Parisi, 2006, and Berman and Laitin, 2008). However, this literature has not considered the use of torture in counterterrorism. If one considers avoidance of torture an individual right, the small but growing literature on the economics of individual rights is relevant (see Mialon and Rubin, 2008, for a summary and synthesis). Seidmann and Stein (2000), Mialon (2005), Leshem (2010), and Wickelgren (2010) examine the right to silence. These papers analyze the effects of preventing adverse inferences from a suspect's silence during interrogation, but they do not analyze the more fundamental right against torture in interrogation.

The only formal analyses of torture that we have found in the literature are Wantchekon and Healy (1999), Chen, Tsai, and Leung (2009), Chen, Chou, and Tsai (2009), and Baliga

and Ely (2010). Wantchekon and Healy analyze torture as a game of incomplete information between a torturer and a victim. They show that torture is carried out with positive probability in equilibrium because even a weak victim might hold out to test whether the torturer is professional or sadistic and even a professional might torture to test whether the victim is weak. Chen, Chou, and Tsai model torture as a war of attrition. They show that torture may occur because of the torturer's uncertainty about the suspect's limit of pain endurance and the suspect's uncertainty about the torturer's limit on pain infliction. In their model, torture can be effective because the innocent cannot verifiably confess while the guilty can. Thus, if torture is applied long enough, guilty suspects will confess and only innocent suspects will be tortured. Baliga and Ely show that torture may be largely ineffective if the torturer cannot commit to continue torturing a suspect he knows to be uninformed and to stop torturing a suspect he knows to be informed. In our paper, we do not focus on the strategic interaction between torturer and suspect but instead focus on the effects of legalizing and regulating torture on the behavior of the counterterrorism agency. Moreover, unlike the above papers, we consider the implications for security and torture of the innocent of agency problems and problems enforcing directives on torture.

Organization of the Paper. Section 2 presents our basic model of counterterrorism when the agency obeys directives on torture. Section 3 extends the model to consider the enforcement problems that arise from the agency being willing to run the risk of legal sanction. Section 4 summarizes.

2. When the Agency Obeys Directives

We begin with a description of the basic model of a government agency (e.g., the CIA, FBI, or DOD) and an individual who may have initiated a terrorist action. Under the assumption

that the agency obeys directives about torture, we compare outcomes, in terms of agency behavior and social welfare, under three scenarios: torture is illegal; torture is illegal except in the face of strong evidence of suspect guilt; and torture is legal. The next section extends the model to study the enforcement problem that arises if the agency may choose, at the cost of legal sanctions, to undertake torture.

2.1 The Basic Model

At time 1, an individual or group of individuals in a large population initiates a terrorist action with probability a , and the agency apprehends an individual. If a terrorist action is not initiated, the apprehended person is necessarily innocent. If a terrorist action is initiated, there is a probability b that the apprehended person is guilty. Thus, $\alpha := ab$ is the probability that the apprehended person has guilty knowledge, and $1 - \alpha$ is the probability that the person has no knowledge. Because we are interested in the logic of the ticking bomb scenario, we set $b = 1$, and therefore $\alpha = a$. We are assuming, in other words, that we are in the case in which the benefits from effective torture would be the largest.

At time 2, not knowing whether or not a terrorist action was initiated, the agency chooses effort $x \geq 0$ to stop a terrorist action by means other than torture. The effort might involve various forms of intelligence gathering and security checks. The cost of x is $c(x)$, where $c' > 0$ and $c'' > 0$. At time 3, if a terrorist action was initiated, the agency's effort stops it with probability $\varphi(x)$, where $\varphi' > 0$ and $\varphi'' \leq 0$. Thus, with probability $\alpha\varphi(x)$, the agency stops a terrorist action, and with probability $1 - \alpha\varphi(x)$, the agency does not stop a terrorist action and infers that either no terrorist action was initiated or one was initiated and the agency's effort did not stop it.

If the agency does not stop a terrorist action, at time 4, it receives evidence ε about whether the apprehended individual initiated a terrorist action. The evidence can be high,

ε_H , or low, ε_L . The probability of high evidence is $q_1(x)$ if a terrorist action was initiated and $q_2(x)$ if one was not initiated, where $q_1(x) \geq q_2(x)$ for any $x > 0$, $q_1(0) = q_2(0) = y \in (0, 1)$, $q'_1 > 0$, $q''_1 \leq 0$, $q'_2 < 0$, and $q''_2 \geq 0$. If the agency increases x , it obtains a more accurate signal, but it still gets a signal even if it chooses $x = 0$.⁹

At time 5, if torture is legal, the agency chooses whether or not to torture the apprehended individual given the available evidence (T or $\neg T$). We initially assume that the agency never uses torture illegally. If torture is not used and a terrorist action was initiated, the terrorist action succeeds, causing social damages D . If a terrorist action was initiated and torture is used, then with probability θ , torture is effective and the agency extracts the information from the guilty individual and stops the terrorist action, and with probability $1 - \theta$, torture is ineffective and the terrorist action succeeds.¹⁰ If torture is used and no terrorist action was initiated, society incurs damages t from torturing an innocent individual.

We assume that the agency internalizes a fraction, $\delta \in [0, 1]$, of social damages, D , and does not internalize the costs of torturing innocent people, t , which is part of the agency problem. This assumption motivates possible constraints on the agency's behavior. If the agency is not as concerned as society about protecting safety or individual rights, society may want to prevent the agency from using torture rather than leave the decision to use torture at the agency's discretion.¹¹ If a terrorist action is initiated but the agency stops it (either through torture or other means) or if a terrorist action is not initiated, then the

⁹ We assume that the same effort x is determining the prevention probability, $\varphi(x)$, and the quality of the evidence about suspects, as represented by $q_1(x)$ and $q_2(x)$. If instead we had $x = (x_1, x_2)$, where x_1 is about prevention while x_2 is about increasing the quality of evidence on suspects, then allowing torture could only reduce prevention effort in the model (see footnote 12 below).

¹⁰As specified, x does not affect the probability θ that torturing the guilty individual yields the information necessary to stop an attack. Making θ an increasing function of x would not affect the main qualitative results of the paper but would add an additional source of complementarity between the agency's use of torture and its other efforts (see footnote 12 below).

¹¹We think of this as an extreme version of a reduced form of the difference between agency incentives and society's preferences. For a general analysis of optimal agency discretion when agency objectives are different from those of society, see Shavell (2007).

agency's payoff is $-c(x)$. If a terrorist action is initiated and the agency does not stop the terrorist action, then the agency's payoff is $-\delta D - c(x)$.

Assuming that the agency obeys directives, we compare outcomes when (1) torture is illegal whether the evidence is low or high, regime B ; (2) torture is legal only when evidence is high, regime H ; and (3) torture is legal whether the evidence is high or low, regime LH . The three regimes have effects on agency behavior, measured by the efforts put into non-torture activities, and on welfare, measured by public safety, the probability of torturing the innocent, and the cost of efforts.

We perform this comparison both in generality, when the functions $q_1(x)$, $q_2(x)$, $\varphi(x)$, and $c(x)$ satisfy only the derivative/inequality conditions listed above, and for the following specific parameterization: $x \geq 0$, $q_1(x) = 1 - (1/2)e^{-\gamma_1 x}$, $q_2(x) = (1/2)e^{-\gamma_2 x}$, $\varphi(x) = 1 - e^{-\lambda x}$, and $c(x) = c(e^x - 1)$, where $\gamma_1, \gamma_2, \lambda, c > 0$. The more general analysis yields results phrased in terms of derived quantities and a tradeoff between a decommitment and a complementarity effect. The parametric analysis is complementary, giving additional understanding of when one might expect one effect or the other to be larger.

2.2 Agency Behavior

Let $\mathbb{D}(x) = \alpha(1 - \varphi(x))[-\delta D]$ denote expected damages under regime B . Under regime B , the agency does not use torture whether it has low or high evidence, strategy $(\neg T, \neg T)$. Thus, its optimal action, x_B^* , solves

$$\max_{x \geq 0} EU_{Agency}^B(x | \neg T, \neg T) = \mathbb{D}(x) - c(x). \quad (1)$$

Under regime H , the agency uses tortures only when it has high evidence, strategy $(\neg T, T)$. Thus, its optimal action, x_H^* , solves

$$\max_{x \geq 0} EU_{Agency}^H(x | \neg T, T) = \psi(x)\mathbb{D}(x) - c(x), \quad (2)$$

where $\psi(x) = (1 - q_1(x)\theta)$ is the probability that strong evidence and torture fail. Under regime LH , the agency uses torture whether it has low or high evidence, strategy (T, T) .

Thus, its optimal action, x_{LH}^* , solves

$$\max_{x \geq 0} EU_{Agency}^{LH}(x|T, T) = (1 - \theta)\mathbb{D}(x) - c(x), \quad (3)$$

where $(1 - \theta)$ is the probability that torture of a guilty person fails.

We now characterize conditions under which a total ban elicits higher effort than a partial ban and show that having no ban on torture unambiguously reduces other agency effort. (Proofs are in the Appendix.)

Proposition 1 *If the optimal non-torture effort levels x_B^* , x_H^* , and x_{LH}^* , are strictly positive, then (a) $x_B^* > x_H^*$ iff $\mathbb{D}'(x_B^*) > [\psi(x_B^*)\mathbb{D}(x_B^*)]'$, and (b) $\min(x_B^*, x_H^*) > x_{LH}^*$.*

Proposition 1(a) shows that the partial ban on torture has two effects on agency efforts to counter terrorism by means other than torture, a decommitment effect and a complementarity effect. Intuitively, allowing the agency to torture when evidence of terrorist action is high reduces the incentives to avoid these situations, a decommitment effect that reduces other efforts. On the other hand, if other efforts raise the chances of having high enough evidence to torture and torture is effective, then torture and other efforts are complementary.

Formally, we see these two effects by noting that

$$\mathbb{D}'(x_B^*) > [\psi(x_B^*)\mathbb{D}(x_B^*)]' \text{ iff } q_1(x_B^*)\varphi'(x_B^*) - q_1'(x_B^*)(1 - \varphi(x_B^*)) > 0. \quad (4)$$

The term $q_1(x_H^*)\varphi'(x_H^*)$ measures the decommitment effect of agency torture on other agency efforts. To the extent that other agency efforts increase the probability of stopping an attack without using torture ($\varphi' > 0$) when they generate enough evidence to use torture (q_1), allowing torture reduces other agency efforts. The term $-q_1'(x_H^*)(1 - \varphi(x_H^*))$ measures the

complementarity effect of other agency efforts and agency torture. To the extent that other agency efforts increase the probability of having enough evidence to use torture ($q'_1 > 0$) when they fail to stop an attack without using torture ($1 - \varphi$), allowing torture increases other agency efforts.¹²

According to Proposition 1(b), allowing torture even when evidence is low unambiguously reduces the agency's other efforts. Allowing torture even when evidence is low does not increase the chances of having enough evidence to warrant the use of torture in the event that other efforts fail, since the agency always has enough evidence to warrant torture if torture is allowed even when evidence is low. Thus, allowing torture even when evidence is low has no complementarity effect and only a decommitment effect on other efforts and therefore unambiguously reduces other efforts.

For our specific parameterization (given at the end of Section 2.1), note that since x_{LH}^* is minimal, strictly positive optima are guaranteed by $(1 - \theta)\mathbb{D}'(0) > c'(0)$, i.e. by $(1 - \theta)Z\lambda > c$, where $Z = \alpha\delta D$. Solving $\mathbb{D}'(x_B) = c'(x_B)$, $(1 - \theta)\mathbb{D}'(x_{LH}) = c'(x_{LH})$, and $[\psi(x_H)\mathbb{D}(x_H)]' = c'(x_H)$, we find $x_B^* = (1/(1 + \lambda))\ln(Z\lambda/c)$, $x_{LH}^* = (1/(1 + \lambda))\ln((1 - \theta)(Z\lambda/c))$, and x_H^* is the solution to $(Z\lambda/c)[1 - \theta + (\theta/2)[1 + (\gamma_1/\lambda)]]e^{-\gamma_1 x} = e^{(1+\lambda)x}$. Clearly, $x_B^* > x_{LH}^*$. Moreover, from Proposition 1(a) and equation 4, we know that $x_B^* > x_H^*$ iff $q_1(x_B^*)\varphi'(x_B^*) > q'_1(x_B^*)(1 - \varphi(x_B^*))$. Plugging in $x_B^* = (1/(1 + \lambda))\ln(Z\lambda/c)$, we get that $x_B^* > x_H^*$ iff

$$\frac{1 + \gamma_1/\lambda}{2}\rho^{\gamma_1} < 1,$$

¹²If we were to also make the effectiveness of torture, θ , an increasing function of non-torture effort, x , then condition (4) would become

$$\theta(x_H^*)\{q_1(x_H^*)\varphi'(x_H^*) - q'_1(x_H^*)(1 - \varphi(x_H^*))\} - \theta'(x_H^*)q_1(x_H^*)(1 - \varphi(x_H^*)) > 0.$$

The complementarity effect would then have the additional term $-\theta'(x_H^*)q_1(x_H^*)(1 - \varphi(x_H^*))$.

If we were to have $x = (x_1, x_2)$, where x_1 is prevention effort and x_2 is effort to improve the quality of evidence on suspects if x_1 fails, so that $\varphi(x_1)$, $q_1(x_2)$, and $c(x_1, x_2) = f(x_1) + g(x_2)$, then there would be no complementarity effect at all and legalizing torture could only reduce prevention effort.

where $\rho = (c/Z\lambda)^{1/(1+\lambda)} < 1$. This condition is more likely to be satisfied if λ is higher, γ_1 is lower, c is lower, or $Z = \alpha\delta D$ is higher. Thus, legalizing torture in strong-evidence cases is more likely to reduce non-torture efforts if non-torture efforts are more effective at stopping attacks or less effective at turning up strong evidence when the suspect is guilty, the costs of non-torture efforts are lower, or the attack threat is higher and the agency is more accountable for realized damages.

2.3 Welfare Effects

We first analyze two central components of welfare, the probability of being safe from a terrorist action and the probability of torturing the innocent. If the agency does not torture, then the probability of safety and the probability of torturing the innocent are, respectively, $S_B^* = 1 - \alpha(1 - \varphi(x_B^*))$ and $Q_B^* = 0$. If the agency tortures only when it has high evidence, then the probabilities are $S_H^* = 1 - \alpha(1 - \varphi(x_H^*))\psi(x_H^*)$ and $Q_H^* = (1 - \alpha)q_2(x_H^*)$. If the agency tortures on the basis of any evidence at all, then the probabilities are $S_{LH}^* = 1 - \alpha(1 - \varphi(x_{LH}^*))(1 - \theta)$ and $Q_{LH}^* = (1 - \alpha)$.

The three safety probabilities are strictly increasing in agency effort, and the three probabilities of torturing the innocent are either flat or strictly decreasing in agency effort. Changes in agency behavior cannot change the part of the welfare ranking due to torturing the innocent—restricting the circumstances under which torture is allowed unambiguously reduces the likelihood of torturing the innocent.

Corollary 1.1 *The probabilities of torturing the innocent, Q_B^* , Q_H^* , and Q_{LH}^* , satisfy $0 = Q_B^* < Q_H^* < Q_{LH}^*$.*

By contrast, changes in agency behavior can change the part of the welfare ranking due to safety. Legalizing torture has a direct effect on security—it stops terrorist attacks when

they have been initiated and a torture-susceptible guilty person is available. By altering agency behavior, it also has an indirect effect on security. The effect on safety depends on which of these effects dominates.

Corollary 1.2 *If the optimal non-torture efforts are strictly positive, then safety probabilities, S_B^* , S_H^* , and S_{LH}^* , satisfy (a) $S_B^* > S_H^*$ iff $[\varphi(x_B^*) - \varphi(x_H^*)] > \theta q_1(x_H^*)[1 - \varphi(x_H^*)]$, (b) $S_B^* > S_{LH}^*$ iff $[\varphi(x_B^*) - \varphi(x_{LH}^*)] > \theta[1 - \varphi(x_{LH}^*)]$, and (c) $S_H^* > S_{LH}^*$ iff $[\varphi(x_H^*) - \varphi(x_{LH}^*)] > \theta[1 - \varphi(x_{LH}^*) - q_1(x_H^*)(1 - \varphi(x_H^*))]$.*

If the availability of torture sufficiently reduces other agency efforts, it reduces safety and increases torture of the innocent; otherwise, it increases safety and torture of the innocent. In more detail, in Corollary 1.2(a), the left-hand side is positive if the decommitment effect of legalizing torture when evidence is high dominates its complementarity effect as identified in Proposition 1(a), so that $x_B^* > x_H^*$. If x_B^* is sufficiently greater than x_H^* , then $S_B^* > S_H^*$. However, if $x_B^* < x_H^*$, the left hand side of the inequality is always negative, implying that if the complementarity effect dominates the decommitment effect, public safety is always higher when torture is legalized than when torture is banned.

For our specific parameterization, $S_B^* > S_H^*$ iff $1 - \theta(1 - (1/2)e^{-\gamma_1 x_H^*}) > e^{-\lambda(x_B^* - x_H^*)}$, where $x_B^* = (1/(1+\lambda)) \ln(Z\lambda/c)$ and x_H^* is the solution to $(Z\lambda/c)[1 - \theta + (\theta/2)[1 + (\gamma_1/\lambda)]e^{-\gamma_1 x}] = e^{(1+\lambda)x}$. A sufficient condition for this is $1 - \theta > e^{-\lambda(x_B^* - x_H^*)}$. If λ is large and θ is small but positive, this condition is likely to be satisfied. Thus, legalizing torture in high-evidence cases is likely to reduce safety if the effectiveness of non-torture efforts is high while the effectiveness of torture is low.

The full welfare levels under each of the regimes are $W_B^* = -(1 - S_B^*)D - c(x_B^*)$, $W_H^* = -(1 - S_H^*)D - Q_H^*t - c(x_H^*)$, and $W_{LH}^* = -(1 - S_{LH}^*)D - Q_{LH}^*t - c(x_{LH}^*)$. The agency's payoffs under each of the regimes are the same as society's payoffs except with D replaced by δD and with $t = 0$. The next corollary compares welfare across regimes.

Corollary 1.3 *If the optimal non-torture efforts are strictly positive, then the welfare levels, W_B^* , W_H^* , and W_{LH}^* , satisfy (a) $W_B^* > W_H^*$ iff $[S_B^* - S_H^*]D + Q_H^*t > [c(x_B^*) - c(x_H^*)]$, (b) $W_B^* > W_{LH}^*$ iff $[S_B^* - S_{LH}^*]D + Q_{LH}^*t > [c(x_B^*) - c(x_{LH}^*)]$, and (c) $W_H^* > W_{LH}^*$ iff $[S_H^* - S_{LH}^*]D + [Q_{LH}^* - Q_H^*]t > [c(x_H^*) - c(x_{LH}^*)]$.*

The availability of torture reduces welfare if the social costs of torturing the innocent, t , are sufficiently high (since t only enters welfare directly, by assumption).

For our specific parameterization, $W_B^* > W_H^*$ iff $\alpha D e^{-\lambda x_H^*} [1 - \theta(1 - (1/2)e^{-\gamma_1 x_H^*}) - e^{-\lambda(x_B^* - x_H^*)}] + (1/2)(1 - \alpha)t e^{-\gamma_2 x_H^*} > c(e^{x_B^*} - e^{x_H^*})$. As $\theta \rightarrow 0$, the *RHS* $\rightarrow 0$, but the *LHS* $\rightarrow (1/2)(1 - \alpha)t e^{-\gamma_2 x_H^*} > 0$. Therefore, if θ is sufficiently small, $W_B^* > W_H^*$. Intuitively, if torture is ineffective, the differences in non-torture effort costs and safety levels between regimes B and H are very small. The only significant difference between B and H is then that the agency uses torture in H , so the probability of torturing the innocent is significantly higher in H , which implies that welfare is lower in H .

Thus, if torture is sufficiently ineffective and the costs of torturing the innocent are sufficiently high, legalizing torture in strong-evidence cases reduces welfare. However, if the extent δ to which the agency internalizes D is not too high, then the availability of torture can increase the agency's payoff although it reduces welfare. In such a case, the agency has incentives to disobey directives on torture, which leads us to the enforcement problem.

3. The Enforcement Problem

Our analysis of the agency problem has focused on the choice of non-torture efforts assuming that the agency obeys directives on torture. We now suppose that the agency faces a penalty, p , if it tortures when not allowed to.¹³ We examine and compare what happens when (1)

¹³The penalty p might correspond to the utility cost of the possibility of being prosecuted within the domestic judiciary system or being declared a war criminal by international courts after choosing to use torture. *Ex post* enforcement creates a "liability rule" against torture. For an economic analysis of liability rules in the protection of individual rights, see Kontorovich (2004) and Kaplow and Shavell (1996).

torture is banned by penalties, regime $\mathcal{B}(p, p)$ with penalties $(p_L, p_H) = (p, p)$, and (2) torture is regulated by a warrant system, regime $\mathcal{W}(p, 0)$, in which torture does not carry a penalty if evidence is high, i.e., the penalties are $(p_L, p_H) = (p, 0)$.

Regime B of the previous section achieves a complete torture ban because the agency obeys directives, while regime $\mathcal{B}(p, p)$ penalizes any torture, and only achieves the ban if the penalty is sufficiently high. The comparison between the regimes H and $\mathcal{W}(p, 0)$ is similar. Though we analyze regime $\mathcal{W}(p, 0)$ as a system of torture warrants, it is formally identical to nonprosecution of torture in high-evidence cases, e.g., if the “necessity” defense successfully avoids prosecution of torturers in cases where evidence of the suspect’s guilt is compelling.

We first gather results that allow us to study agency behavior under the two regimes. With these in place, we then analyze, within the context of the model, the Dershowitz argument for reductions in the frequency of torture arising from a torture warrant system. We then turn to the existence of slippery slopes.

3.1 Agency Behavior

Let $P(\varepsilon_L|x)$ and $P(\varepsilon_H|x)$ be the likelihoods of a low and a high evidence suspect, and let $\beta_L(x)$ and $\beta_H(x)$ be the likelihoods of a low and a high evidence suspect being guilty, respectively. From Bayes’ rule, we have

$$\beta_L(x) = \frac{\alpha(1 - \varphi(x))(1 - q_1(x))}{P(\varepsilon_L|x)} \text{ and } \beta_H(x) = \frac{\alpha(1 - \varphi(x))q_1(x)}{P(\varepsilon_H|x)}, \text{ where} \quad (5)$$

$$P(\varepsilon_L|x) = \alpha(1 - \varphi(x))(1 - q_1(x)) + (1 - \alpha)(1 - q_2(x)) \text{ and} \quad (6)$$

$$P(\varepsilon_H|x) = \alpha(1 - \varphi(x))q_1(x) + (1 - \alpha)q_2(x). \quad (7)$$

Lemma 1 *For all x , (a) $\beta_L(x) < \beta_H(x)$ and $\beta'_L(x) < 0$, (b) the agency’s optimal choice is T at ε_L if $\beta_L(x) > \frac{p_L}{\theta\delta D}$ and T at ε_H if $\beta_H(x) > \frac{p_H}{\theta\delta D}$, and (c) $[P(\varepsilon_L|x) + P(\varepsilon_H|x)]' < 0$, which implies that at least one of the derivatives, $[P(\varepsilon_L|x)]'$ and $[P(\varepsilon_H|x)]'$, is negative.*

$\beta_L(x)$ is decreasing in x because increases in x reduce the probability of a terrorist action eluding the agency's non-torture efforts; reduce the probability of false exculpatory evidence; and increase the probability of valid exculpatory evidence.

The agency's behavior under the different regimes can be understood in two parts: (A) for a given torture policy, finding the agency's optimal effort as a function of the penalty p and calculating the associated payoffs, i.e. the value function; and (B) comparing the value functions to pick the optimal torture policy as a function of p . From the first two parts of Lemma 1, we know that under regime $\mathcal{B}(p, p)$, the agency's optimal torture policy is either (T, T) , $(-T, T)$, or $(-T, -T)$, and that under regime $\mathcal{W}(p, 0)$, it is either (T, T) or $(-T, T)$.

For part (A) of the analysis of regime $\mathcal{B}(p, p)$, let $f_{\mathcal{B}}^{T,T}(p)$, $f_{\mathcal{B}}^{-T,T}(p)$, and $f_{\mathcal{B}}^{-T,-T}(p)$ be the agency's value functions at the optimal effort levels for a given p when the agency's policy is (T, T) , $(-T, T)$, and $(-T, -T)$, respectively. Specifically, these are

$$f_{\mathcal{B}}^{T,T}(p) = \max_{x \geq 0} [(1 - \theta)\mathbb{D}(x) - c(x)] - p \{P(\varepsilon_L|x) + P(\varepsilon_H|x)\}, \quad (8)$$

$$f_{\mathcal{B}}^{-T,T}(p) = \max_{x \geq 0} [\psi(x)\mathbb{D}(x) - c(x)] - p \{0 + P(\varepsilon_H|x)\}, \text{ and} \quad (9)$$

$$f_{\mathcal{B}}^{-T,-T}(p) = \max_{x \geq 0} [\mathbb{D}(x) - c(x)] - p \{0 + 0\}. \quad (10)$$

For part (B) of the analysis of regime $\mathcal{B}(p, p)$, agency behavior is the policy part of the solution to the problem

$$V_{\mathcal{B}}(p) = \max\{f_{\mathcal{B}}^{T,T}(p), f_{\mathcal{B}}^{-T,T}(p), f_{\mathcal{B}}^{-T,-T}(p)\}. \quad (11)$$

Similarly, for part (A) of the analysis of regime $\mathcal{W}(p, 0)$, we have

$$f_{\mathcal{W}}^{T,T}(p) = \max_{x \geq 0} [(1 - \theta)\mathbb{D}(x) - c(x)] - p \{P(\varepsilon_L|x) + 0\}, \text{ and} \quad (12)$$

$$f_{\mathcal{W}}^{-T,T}(p) = \max_{x \geq 0} [\psi(x)\mathbb{D}(x) - c(x)] - p \{0 + 0\}. \quad (13)$$

and part (B) asks for the policy part of the solution to

$$V_{\mathcal{W}}(p) = \max\{f_{\mathcal{W}}^{T,T}(p), f_{\mathcal{W}}^{-T,T}(p)\}. \quad (14)$$

Under both regimes $\mathcal{B}(p, p)$ and $\mathcal{W}(p, 0)$, as the agency tortures in fewer circumstances, i.e., as we move down through the three cases in (8), (9), and (10) under regime $\mathcal{B}(p, p)$, and down through the two cases in (12) and (13) under regime $\mathcal{W}(p, 0)$, for all fixed values of x , the terms that multiply the penalties p decrease. For example, under regime $\mathcal{B}(p, p)$, we have p multiplying the terms $\{P(\varepsilon_L|x) + P(\varepsilon_H|x)\}$, then $\{0 + P(\varepsilon_H|x)\}$, and then $\{0 + 0\}$. This suggests that the value functions in the two regimes have a “single-crossing from above in p ” property. This turns out to be true but is somewhat more subtle because we are evaluating the penalty terms at different values of x . The following lemma addresses this issue and gives additional useful properties of the value functions.

Lemma 2 *The value functions for part (A) of the analysis of regimes $\mathcal{B}(p, p)$ and $\mathcal{W}(p, 0)$ have the following properties: the value function for any policy that tortures is strictly decreasing in p ; convex in p ; and crosses the value function for a policy that tortures in fewer circumstances exactly once, from above, as p increases.*

Define $\underline{p}_{\mathcal{B}}$, $\overline{p}_{\mathcal{B}}$, and $p_{\mathcal{W}}$ to be the penalty levels at which we have the crossings,

$$f_{\mathcal{B}}^{T,T}(\underline{p}_{\mathcal{B}}) = f_{\mathcal{B}}^{-T,T}(\underline{p}_{\mathcal{B}}), f_{\mathcal{B}}^{-T,T}(\overline{p}_{\mathcal{B}}) = f_{\mathcal{B}}^{-T,-T}(\overline{p}_{\mathcal{B}}), \text{ and} \quad (15)$$

$$f_{\mathcal{W}}^{T,T}(p_{\mathcal{W}}) = f_{\mathcal{W}}^{-T,T}(p_{\mathcal{W}}). \quad (16)$$

3.2 The Dershowitz Argument

As noted above, the core of the Dershowitz (2002) argument is that torture happens although it is illegal and that an enforced system of judicial warrants could bring this under control, resulting in less torture. In our model, Dershowitz’s case corresponds to regime $\mathcal{B}(p^\circ, p^\circ)$ for any p° at which the agency’s torture policy is $(-T, T)$. The appropriate comparison is the change of regime to regime $\mathcal{W}(p^\circ, 0)$. From Lemma 1(b), the agency’s choice of torture policy in regime $\mathcal{W}(p^\circ, 0)$ is either $(-T, T)$, the case we analyze now, or (T, T) , which we analyze in Proposition 4 as part of our discussion of slippery slopes.

The patterns of the crossings of the value functions determine when $(-T, T)$ is simultaneously optimal under regime $\mathcal{B}(p, p)$ and regime $\mathcal{W}(p, 0)$.

Proposition 2 *The agency's optimal torture policy is $(-T, T)$ under both regimes $\mathcal{B}(p^\circ, p^\circ)$ and $\mathcal{W}(p^\circ, 0)$ if $\underline{p}_{\mathcal{B}} < \overline{p}_{\mathcal{B}}$, $p_{\mathcal{W}} < \overline{p}_{\mathcal{B}}$, and $p^\circ \in (\max\{p_{\mathcal{W}}, \underline{p}_{\mathcal{B}}\}, \overline{p}_{\mathcal{B}})$.*

We now compare agency effort, safety, and frequency of torture of the innocent under regime $\mathcal{B}(p, p)$ and regime $\mathcal{W}(p, 0)$ when the agency's optimal policy is $(-T, T)$ under both regimes $\mathcal{B}(p, p)$ and $\mathcal{W}(p, 0)$.

Proposition 3 *Suppose the agency's optimal torture policy is $(-T, T)$ under both regimes $\mathcal{B}(p^\circ, p^\circ)$ and $\mathcal{W}(p^\circ, 0)$. (a) If $[P(\varepsilon_H|x)]' < 0$ for all x , regime $\mathcal{W}(p^\circ, 0)$ has lower effort and thus lower safety and more frequent torture of the innocent. (b) If $[P(\varepsilon_H|x)]' > 0$ for all x , regime $\mathcal{W}(p^\circ, 0)$ has higher effort and safety and less frequent torture of the innocent.*

The effects of moving from $\mathcal{B}(p^\circ, p^\circ)$ to $\mathcal{W}(p^\circ, 0)$ depend crucially on the sign of the derivative $[P(\varepsilon_H|x)]'$. We have

$$[P(\varepsilon_H|x)]' = \alpha[q_1'(x)(1 - \varphi(x)) - \varphi'(x)q_1(x)] + (1 - \alpha)q_2'(x). \quad (17)$$

Thus, the sign of $[P(\varepsilon_H|x)]'$ depends on the relative sizes of the three terms, $-\alpha\varphi'(x)q_1(x)$, $\alpha q_1'(x)(1 - \varphi(x))$, and $(1 - \alpha)q_2'(x)$, which correspond to the following three effects. **De-commitment:** moving from $\mathcal{B}(p^\circ, p^\circ)$ to $\mathcal{W}(p^\circ, 0)$ tends to reduce agency effort x because it eliminates the agency's cost of using torture if the individual is guilty and the agency's evidence turns out to be high, which reduces the agency's commitment not to use torture. This effect is captured by the negative term $-\alpha\varphi'(x)q_1(x)$. **Complementarity:** moving from $\mathcal{B}(p^\circ, p^\circ)$ to $\mathcal{W}(p^\circ, 0)$ tends to increase x because increasing x increases the probability that the agency has high evidence if the individual is guilty, in which case the agency can use torture without punishment. This effect is captured by the positive term $\alpha q_1'(x)(1 - \varphi(x))$.

Decomplementarity: moving from $\mathcal{B}(p^\circ, p^\circ)$ to $\mathcal{W}(p^\circ, 0)$ tends to reduce x because reducing x increases the probability that the agency has high evidence if the individual is innocent, in which case the agency can escape punishment for using torture on an innocent individual. This effect is captured by the negative term $(1 - \alpha)q_2'(x)$.

If the complementarity effect dominates both the decommitment and decomplementarity effects, then moving from $\mathcal{B}(p^\circ, p^\circ)$ to $\mathcal{W}(p^\circ, 0)$ increases agency effort x , and thereby increases safety and reduces the probability of torturing the innocent. However, if the decommitment and decomplementarity effects together dominate the complementarity effect, then moving from $\mathcal{B}(p^\circ, p^\circ)$ to $\mathcal{W}(p^\circ, 0)$ reduces agency effort x , and thereby reduces safety and increases the probability of torturing the innocent.

For our parameterization, $[P(\varepsilon_H|x)]' < 0$ iff $e^{-\gamma_2 x}[-\alpha\lambda e^{(\gamma_2 - \lambda)x} + (1/2)\alpha(\lambda + \gamma_1)e^{(\gamma_2 - \gamma_1 - \lambda)x} - (1/2)(1 - \alpha)\gamma_2] < 0$. Sufficient conditions for this are (a) $\lambda > \gamma_2$ and (b) $\alpha\gamma_1 < \alpha\lambda + (1 - \alpha)\gamma_2$. Conditions (a) and (b) are both satisfied if λ is sufficiently high. Moreover, even if λ is only high enough to satisfy condition (a), condition (b) is still satisfied if α is low enough. Thus, moving from \mathcal{B} to \mathcal{W} is likely to reduce agency efforts and security if agency efforts are effective at stopping attacks and attacks are infrequent.

3.3 Slippery Slopes

We have a slippery slope if legalizing torture in dire circumstances increases the set of circumstances in which torture occurs. This happens if either $(-T, T)$ or $(-T, -T)$ are optimal for the agency under $\mathcal{B}(p, p)$ while (T, T) is optimal under the torture warrant system $\mathcal{W}(p, 0)$. In either case, switching to a torture warrant system lowers agency effort and increases the frequency of torture of the innocent.

Proposition 4 *If p° is such that either (a) the agency's optimal torture policy is $(-T, -T)$ under regime $\mathcal{B}(p^\circ, p^\circ)$ and is (T, T) under regime $\mathcal{W}(p^\circ, 0)$, or (b) the agency's optimal*

torture policy is $(-T, T)$ under regime $\mathcal{B}(p^\circ, p^\circ)$ and is (T, T) under regime $\mathcal{W}(p^\circ, 0)$, then switching from regime $\mathcal{B}(p^\circ, p^\circ)$ to regime $\mathcal{W}(p^\circ, 0)$ reduces agency effort and increases the frequency of torture of the innocent.

If the torture policy changes in either of the ways indicated, we have direct evidence of lower agency effort—from Lemma 1, the agency only tortures at ε_L if $\beta_L(x)$ is sufficiently high, and higher values of β_L can only arise from lower values of x . More intuitively, the slippery slopes in Proposition 4(a) and (b) both arise from lower effort reducing the quality of exculpatory evidence. Moving from regime \mathcal{B} to regime \mathcal{W} eliminates the agency’s penalty for using torture in high-evidence cases. This reduces the agency’s commitment to non-torture efforts, whether or not the agency was already using torture in high-evidence cases under \mathcal{B} . A reduction in non-torture efforts reduces the quality of the agency’s evidence, which increases the agency’s incentives to adopt a strategy that uses torture even when evidence is low. Adoption of such a strategy further reduces agency efforts, further reinforcing the agency’s incentives to use torture even when evidence is low.

Whether or not moving down either type of slippery slope translates to lower safety depends, as in Corollary 1.2, on the effectiveness of torture. If torture is sufficiently ineffective relative to other efforts, then lowering other efforts and torturing more reduces safety. Whether or not each type of slippery slope arises depends on the relative locations of the crossing points defined in (15) and (16).

Proposition 5 *Suppose $\underline{p}_{\mathcal{B}} < \overline{p}_{\mathcal{B}}$. (a) The agency’s optimal torture policy is uniquely $(-T, -T)$ under regime $\mathcal{B}(p^\circ, p^\circ)$ and is uniquely (T, T) under regime $\mathcal{W}(p^\circ, 0)$ if and only if $\overline{p}_{\mathcal{B}} < p^\circ < p_{\mathcal{W}}$. (b) The agency’s optimal torture policy is uniquely $(-T, T)$ under regime $\mathcal{B}(p^\circ, p^\circ)$ and is uniquely (T, T) under regime $\mathcal{W}(p^\circ, 0)$ if and only if $\underline{p}_{\mathcal{B}} < p^\circ \min\{\overline{p}_{\mathcal{B}}, p_{\mathcal{W}}\}$.*

The crossing point conditions in Proposition 5(a) and (b) are more easily satisfied if the value functions for policies involving torture in regime \mathcal{B} are steeper functions of p while the

value function for the policy involving torture in regime \mathcal{W} is shallower. The value functions for \mathcal{B} are relatively steeper if the level of $P(\varepsilon_H|x)$ is higher, because this term is the difference between the penalties across \mathcal{B} and \mathcal{W} . If the probability of having high evidence suspects is higher, the agency is more tempted to rely on torture as a counterterrorism tool, and moving to \mathcal{W} is more likely to lower agency effort and increase torture of the innocent. Moreover, if $P(\varepsilon_H|x)$ does not vary much with x , then the value functions in \mathcal{B} are less convex, meaning that they stay steeper for a larger range of penalties. If the probability of having a high evidence suspect is less sensitive to effort, then there is less loss to lowering effort if torture is used more, which also makes a slippery slope more likely to arise.

4. Summary and Future Work

We developed a model of counterterrorism to analyze the effects of allowing the government to use torture when evidence of terrorist involvement is strong. We first examined the case in which the agency tasked with counterterrorism places a different weight on torture than society does, but in which it follows any directives. In this case, we showed that allowing the agency to use torture in strong-evidence cases may reduce its efforts to stop terrorism by means other than torture. This effect blunts any gain to safety that may arise through torture, and the net effect may be a reduction in security.

We then extended our analysis to encompass the possibility that there is an enforcement problem and the agency is willing to disobey torture directives at the risk of legal sanction. This extension allowed us to examine conditions under which the Dershowitz argument that a system of torture warrants could reduce torture holds or fails, showing that its validity depends on the agency's ineffectiveness at stopping attacks and finding exculpatory evidence when it exists. It also brought to light a slippery slope that works through the endogeneity

of the quality of information. Allowing torture in strong-evidence cases may reduce the agency's non-torture efforts. The resulting agency deskilling may then reduce the quality of exculpatory evidence, which may lead to torture even in weak-evidence cases. The main arguments we developed have a simple outline: loosening constraints on torture may induce changes in agency behavior that may compromise security and reduce the quality of the agency's evidence to such an extent that it motivates the use of torture even in the face of potentially exculpatory evidence.

In future work, it might be interesting to endogenize terrorism in our model and examine the effects of legalizing torture on the probability of attack. If legalizing torture has a sufficiently large decommitment effect, it might indirectly increase the probability of attack by reducing the agency's preventive effort and thereby increasing the probability that an attack would succeed. On the other hand, torture could directly reduce the attack probability if it is effective as a punishment or means of intimidation. However, if illegitimate regimes are more prone to use torture (e.g., for intimidation) than legitimate ones, legalizing torture might also signal that a regime is illegitimate, which may increase individuals' benefits of attacking it. Dreher, Gassebner, and Siemers (2010) find that terrorist attacks are positively associated with human rights violations (including torture) across regions. By engaging in torture, a government risks pooling with illegitimate regimes and inciting a terrorist backlash. In addition, the use of torture may affect the structure of terrorist organizations. Torture may be used to elicit information about other terrorist suspects or members of the terrorist network, which may induce the terrorist organization to become more decentralized. Lastly, we did not consider potential adverse selection problems with respect to torture. Once torture is allowed in extreme conditions, more sadistic individuals might be drawn to work in the agency and naturally "extreme conditions" may become less extreme.

A Mathematical Appendix

Proof of Proposition 1. As we are assuming smooth concave objective functions and interior solutions, analyses of the first order conditions are sufficient. For (a), note that the objective functions in (1), (2), and (3) are the sums of strictly concave functions. This means that optimal agency effort under regime B is characterized by (1') $\mathbb{D}'(x_B^*) = c'(x_B^*)$, and it is characterized by (2') $[\psi(x_H^*)\mathbb{D}(x_H^*)]' = c'(x_H^*)$ under regime H . Since $c''(x) > 0$, the derivative condition $\mathbb{D}'(x_B^*) > [\psi(x_B^*)\mathbb{D}(x_B^*)]'$ holds iff $x_B^* > x_H^*$. The first part of (b), $x_B^* > x_{LH}^*$, follows from $(1 - \theta) < 1$ and the supermodularity of $h(x, \theta) := (1 - \theta)\mathbb{D}(x) - c(x)$. For the second part of (b), note that under regime LH , optimal agency effort is characterized by (3') $(1 - \theta)\mathbb{D}'(x_{LH}^*) = c'(x_{LH}^*)$. Since $c''(x) > 0$, (2') and (3') deliver $[\psi(x_H^*)\mathbb{D}(x_H^*)]' > (1 - \theta)\mathbb{D}'(x_H^*)$ iff $x_H^* > x_{LH}^*$. Since $\psi(x) = (1 - q_1(x)\theta)$, rearrangement yields $x_H^* > x_{LH}^*$ iff $q_1'(x_H^*)\mathbb{D}(x_H^*) < (1 - q_1(x_H^*))\mathbb{D}'(x_H^*)$, and the left-hand side is negative while the right-hand side is positive. ■

Proof of Lemma 1. For the first part of (a), note that $\beta_L(x) < \beta_H(x)$ because $q_1(x) > q_2(x)$. For the second part, note that

$$\beta_L(x) = \frac{\alpha(1 - \varphi(x))(1 - q_1(x))}{\alpha(1 - \varphi(x))(1 - q_1(x)) + (1 - \alpha)(1 - q_2(x))} = \frac{a(x)}{a(x) + b(x)}, \quad (18)$$

where $a(\cdot)$ is decreasing in x and $b(x)$ is increasing. (b) is a direct implication of utility maximization. For (c), note that $P(\varepsilon_L|x) + P(\varepsilon_H|x) = \alpha(1 - \varphi(x)) + (1 - \alpha)$, so $[P(\varepsilon_L|x) + P(\varepsilon_H|x)]' < 0$ since $\varphi'(x) > 0$. ■

Proof of Lemma 2. By the envelope theorem, the value functions in (8), (9), and (12) are strictly decreasing in the penalty p . For convexity, note that all three problems are of the form $f(p) = \max_{x \geq 0} n(x) - pm(x)$ and that, $x^*(p)$, the optimal effort x for a given value of p has the property that $\frac{d}{dp}x^*(p)$ has the opposite sign of $m'(x^*(p))$. Now,

$f(p) = n(x^*(p)) - pm'(x^*(p))$ where $x^*(p)$ satisfies $[n'(x^*(p)) - pm'(x^*(p))] \equiv 0$. Therefore, $f'(p) = [n'(x^*) - pm'(x^*)] \frac{d}{dp} x^*(p) - m(x^*(p))$ and $f''(p) = -m'(x^*) \frac{d}{dp} x^*(p)$. Since the two terms in this product have the opposite sign, we have $f''(p) \geq 0$. We now show that the value functions in the two regimes have the “single-crossing from above in p ” property. We treat the simpler case, regime \mathcal{W} , first. Since (T, T) is the agency’s optimal policy at $p = 0$, $f_{\mathcal{W}}^{T,T}(0) > f_{\mathcal{W}}^{-T,T}(0)$. By the envelope theorem, under regime \mathcal{W} with penalties $(p, 0)$, $\frac{d}{dp} f_{\mathcal{W}}^{T,T}(p) = -P(\varepsilon_L|x) < 0$, while $\frac{d}{dp} f_{\mathcal{W}}^{-T,T}(p) = 0$. For the same reasons, under regime $\mathcal{B}(p, p)$, the crossings of $f_{\mathcal{B}}^{-T,-T}(p)$ happen from above. All that is left to consider is the crossing of $f_{\mathcal{B}}^{T,T}(p)$ and $f_{\mathcal{B}}^{-T,T}(p)$. In order to show that $f_{\mathcal{B}}^{T,T}(p)$ crosses $f_{\mathcal{B}}^{-T,T}(p)$ at most once from above as p increases, it is sufficient to show that for at any p^\dagger where $f_{\mathcal{B}}^{T,T}(p^\dagger) = f_{\mathcal{B}}^{-T,T}(p^\dagger)$, $f_{\mathcal{B}}^{T,T}(\cdot)$ is steeper than $f_{\mathcal{B}}^{-T,T}(\cdot)$. By the envelope theorem again, we need to show that

$$[P(\varepsilon_L|x_{T,T}^*(p^\dagger)) + P(\varepsilon_H|x_{T,T}^*(p^\dagger))] > P(\varepsilon_H|x_{-T,T}^*(p^\dagger)), \quad (19)$$

where $x_{T,T}^* = x_{T,T}^*(p^\dagger)$ and $x_{-T,T}^* = x_{-T,T}^*(p^\dagger)$ are the agency’s corresponding optimal effort levels under regime $\mathcal{B}(p^\dagger, p^\dagger)$ when following the torture policies (T, T) and $(-T, T)$. Now, at p^\dagger , both torture policies $(-T, T)$ and (T, T) are optimal, which implies that $\beta_L(x_{-T,T}^*) \leq p^\dagger/\theta\delta D \leq \beta_L(x_{T,T}^*)$. Since $\beta'_L(x) < 0$, this implies that $x_{-T,T}^* \geq x_{T,T}^*$. Then,

$$\begin{aligned}
 & P(\varepsilon_L|x_{TT}^*) + P(\varepsilon_H|x_{TT}^*) - P(\varepsilon_H|x_{-TT}^*) \\
 &= \alpha(1 - \varphi(x_{TT}^*)) + (1 - \alpha) - \alpha(1 - \varphi(x_{-TT}^*))q_1(x_{-TT}^*) - (1 - \alpha)q_2(x_{-TT}^*) \\
 &= (1 - \alpha) [1 - q_2(x_{-TT}^*)] + \alpha \{1 - \varphi(x_{TT}^*) - q_1(x_{-TT}^*) + \varphi(x_{-TT}^*)q_1(x_{-TT}^*)\} \\
 &= (> 0) + \alpha \{1 - \varphi(x_{-TT}^*) + \varphi(x_{-TT}^*) - \varphi(x_{TT}^*) - q_1(x_{-TT}^*) + \varphi(x_{-TT}^*)q_1(x_{-TT}^*)\} \\
 &= (> 0) + \underbrace{\alpha\{1 - \varphi(x_{-TT}^*)\}}_{>0} \underbrace{\{1 - q_1(x_{-TT}^*)\}}_{\geq 0} + \underbrace{\alpha\{\varphi(x_{-TT}^*) - \varphi(x_{TT}^*)\}}_{\geq 0} > 0,
 \end{aligned} \quad (20)$$

where the weak inequality comes from $x_{-TT}^* \geq x_{TT}^*$ and $\varphi'(x) > 0$. ■

Proof of Proposition 2. From Lemma 2, $f_{\mathcal{B}}^{T,T}(p)$ crosses $f_{\mathcal{B}}^{-T,T}(p)$ and $f_{\mathcal{B}}^{-T,T}(p)$ crosses $f_{\mathcal{B}}^{-T,-T}(p)$ once from above as p increases. Let $\overline{p_{\mathcal{B}}}$ and $\underline{p_{\mathcal{B}}}$ be the points at which $f_{\mathcal{B}}^{T,T}(\overline{p_{\mathcal{B}}}) = f_{\mathcal{B}}^{-T,T}(\overline{p_{\mathcal{B}}})$ and $f_{\mathcal{B}}^{-T,T}(\underline{p_{\mathcal{B}}}) = f_{\mathcal{B}}^{-T,-T}(\underline{p_{\mathcal{B}}})$, respectively. If $\underline{p_{\mathcal{B}}} < \overline{p_{\mathcal{B}}}$, under regime $\mathcal{B}(p^\circ, p^\circ)$, for $p^\circ \in [0, \underline{p_{\mathcal{B}}})$, $V_{\mathcal{B}}(p^\circ) = f_{\mathcal{B}}^{T,T}(p^\circ)$, for $p^\circ \in (\underline{p_{\mathcal{B}}}, \overline{p_{\mathcal{B}}})$, $V_{\mathcal{B}}(p^\circ) = f_{\mathcal{B}}^{-T,T}(p^\circ)$, and for $p^\circ > \overline{p_{\mathcal{B}}}$, $V_{\mathcal{B}}(p^\circ) = f_{\mathcal{B}}^{-T,-T}(p^\circ)$. The agency chooses $(-T, T)$ under regime \mathcal{B} if $p^\circ \in (\underline{p_{\mathcal{B}}}, \overline{p_{\mathcal{B}}})$. Similarly, from Lemma 2, $f_{\mathcal{W}}^{T,T}(p)$ crosses $f_{\mathcal{W}}^{-T,T}(p)$ once from above as p increases. Let $p_{\mathcal{W}}$ be the point at which $f_{\mathcal{W}}^{T,T}(p_{\mathcal{W}}) = f_{\mathcal{W}}^{-T,T}(p_{\mathcal{W}})$. Under regime $\mathcal{W}(p^\circ, 0)$, for $p^\circ \in [0, p_{\mathcal{W}})$, $V_{\mathcal{W}}(p^\circ) = f_{\mathcal{W}}^{T,T}(p^\circ)$, and for $p^\circ > p_{\mathcal{W}}$, $V_{\mathcal{W}}(p^\circ) = f_{\mathcal{W}}^{-T,T}(p^\circ)$. The agency chooses $(-T, T)$ under regime \mathcal{W} if $p^\circ > p_{\mathcal{W}}$. Thus, if $p_{\mathcal{W}} < \overline{p_{\mathcal{B}}}$, then for $p^\circ \in (\max\{p_{\mathcal{W}}, \underline{p_{\mathcal{B}}}\}, \overline{p_{\mathcal{B}}})$, the agency chooses $(-T, T)$ under both regimes \mathcal{B} and \mathcal{W} . ■

Proof of Proposition 3. For $t \in \{0, 1\}$ and $x \geq 0$, define

$$h(x, t) = [\psi(x)\mathbb{D}(x) - c(x)] - p^\circ \{0 + t \cdot P(\varepsilon_H|x)\} \quad (21)$$

so that $\max_{x \geq 0} h(x, 1)$ is the agency's optimal effort problem in regime $\mathcal{B}(p^\circ, p^\circ)$ if following the policy $(-T, T)$, and $\max_{x \geq 0} h(x, 0)$ is the agency's optimal effort problem in regime $\mathcal{W}(p^\circ, 0)$ if following the policy $(-T, T)$. Simple increasing differences comparative statics (e.g., Corbae, Stinchcombe, and Zeeman, 2009, §2.8.b) show that $x^*(1) \geq x^*(0)$ if for all $x' > x$, $h(x', 1) - h(x, 1) > h(x', 0) - h(x, 0)$, and $x^*(1) \leq x^*(0)$ if for all $x' > x$, $h(x', 1) - h(x, 1) < h(x', 0) - h(x, 0)$. Applied to (21), $h(x', 1) - h(x, 1) > h(x', 0) - h(x, 0)$ iff $-p^\circ \cdot \{P(\varepsilon_H|x') - P(\varepsilon_H|x)\} > 0$. Thus, $[P(\varepsilon_H|x)]' < 0$ for all x implies that $\mathcal{B}(p^\circ, p^\circ)$ has higher optimal effort than $\mathcal{W}(p^\circ, 0)$, part (a) of the Proposition, and $[P(\varepsilon_H|x)]' > 0$ for all x implies that $\mathcal{B}(p^\circ, p^\circ)$ has lower optimal effort than $\mathcal{W}(p^\circ, 0)$, part (b) of the Proposition. ■

Proof of Proposition 4. In both (a) and (b), the switch from $\mathcal{B}(p^\circ, p^\circ)$ to $\mathcal{W}(p^\circ, 0)$ involves changing the policy from not torturing when evidence is low to torturing when evidence is low. By Lemma 1(a), $\beta'_L(x) < 0$. Therefore, by Lemma 1(b), the change in

torture policy indicates a decrease in effort. The argument for the increasing frequency of torture of the innocent follows, as above, from $0 < q_2(x) < 1$. ■

Proof of Proposition 5. From Lemma 2, $f_{\mathcal{B}}^{T,T}(p)$ crosses $f_{\mathcal{B}}^{-T,T}(p)$ and $f_{\mathcal{B}}^{-T,T}(p)$ crosses $f_{\mathcal{B}}^{-T,-T}(p)$ once from above as p increases. Let $\overline{p_{\mathcal{B}}}$ and $\underline{p_{\mathcal{B}}}$ be the points at which $f_{\mathcal{B}}^{T,T}(\overline{p_{\mathcal{B}}}) = f_{\mathcal{B}}^{-T,T}(\overline{p_{\mathcal{B}}})$ and $f_{\mathcal{B}}^{-T,T}(\underline{p_{\mathcal{B}}}) = f_{\mathcal{B}}^{-T,-T}(\underline{p_{\mathcal{B}}})$, respectively. Under regime $\mathcal{B}(p^\circ, p^\circ)$, if $\underline{p_{\mathcal{B}}} < \overline{p_{\mathcal{B}}}$, then for $p^\circ \in (\underline{p_{\mathcal{B}}}, \overline{p_{\mathcal{B}}})$, $V_{\mathcal{B}}(p^\circ) = f_{\mathcal{B}}^{-T,T}(p^\circ)$, and thus, the agency chooses $(-T, T)$, and for $p > \overline{p_{\mathcal{B}}}$, $V_{\mathcal{B}}(p^\circ) = f_{\mathcal{B}}^{-T,-T}(p^\circ)$, and thus, the agency chooses $(-T, -T)$. Similarly, from Lemma 2, $f_{\mathcal{W}}^{T,T}(p)$ crosses $f_{\mathcal{W}}^{-T,T}(p)$ once from above as p increases. Let $p_{\mathcal{W}}$ be the point at which $f_{\mathcal{W}}^{T,T}(p_{\mathcal{W}}) = f_{\mathcal{W}}^{-T,T}(p_{\mathcal{W}})$. Under regime $\mathcal{W}(p^\circ, 0)$, for $p^\circ \in [0, p_{\mathcal{W}})$, $V_{\mathcal{W}}(p^\circ) = f_{\mathcal{W}}^{T,T}(p^\circ)$, and thus, the agency chooses (T, T) . Case (a): if $\overline{p_{\mathcal{B}}} < p_{\mathcal{W}}$, then the agency chooses $(-T, -T)$ under regime \mathcal{B} and (T, T) under regime \mathcal{W} if $p^\circ \in (\overline{p_{\mathcal{B}}}, p_{\mathcal{W}})$. Case (b): if $\underline{p_{\mathcal{B}}} < p_{\mathcal{W}}$, then the agency chooses $(-T, T)$ under regime \mathcal{B} and (T, T) under regime \mathcal{W} if $p^\circ \in (\underline{p_{\mathcal{B}}}, \min\{\overline{p_{\mathcal{B}}}, p_{\mathcal{W}}\})$. ■

REFERENCES

- Bagaric, Mirko and Clarke, Julie (2006). *Torture: When the Unthinkable is Morally Permissible*, State University of New York Press: New York.
- Baliga, Sandeep and Ely, Jeffrey C. (2010). "Torture," Northwestern University Working Paper.
- Berman, Eli and Laitin, David D. (2008). "Religion, Terrorism and Public Goods: Testing the Club Model," *Journal of Public Economics* 92, 1942-1967.
- Chen, Kong-Ping, Tsai, Tsung-Sheng, and Leung, Angela (2009). "Judicial Torture as a Screening Device," Academia Sinica Working Paper.
- Chen, Kong-Ping, Chou, Chien-Fu, and Tsai, Tsung-Sheng (2009). "Judicial Torture as a War of Attrition," Academia Sinica Working Paper.
- Corbae, Dean, Stinchcombe, Maxwell B., and Zeeman, Juraj (2009). *An Introduction to Mathematical Analysis for Economic Theory and Econometrics*, Princeton University Press: New Jersey.
- Dershowitz, Alan M. (2002). *Why Terrorism Works*, Yale University Press: New Haven.
- Dershowitz, Alan M. (2003). "The Torture Warrant: A Response to Professor Strauss," *New York Law School Law Review* 48, 275-294.
- Dreher, Axel, Gassebner, Martin, and Siemers, Lars-H. R. (2010). "Does Terror Threaten Human Rights? Evidence From Panel Data," *Journal of Law and Economics* 53, 65-93.
- Eggen, Dan (2007). "Torture Stance Raises Doubts on Mukasey," *Washington Post* October 27.
- Enders, Walter and Sandler, Todd (2004). "An Economic Perspective on Transnational Terrorism," *European Journal of Political Economy* 20, 301-316.
- Enders, Walter and Sandler, Todd (2005). "Transnational Terrorism 1968-2000: Thresholds, Persistence, and Forecasts," *Southern Economic Journal* 71, 467-483.
- Fay, Major General George R. (2004). "Fay Report: Investigation of Intelligence Activities At Abu Ghraib: Executive Summary," Available at Website <http://f11.findlaw.com/news.findlaw.com/hdocs/docs/dod/fay82504rpt.pdf>.
- Garoupa, Nuno, Klick, Jonathan, and Parisi, Francesco (2006). "A Law and Economics Perspective on Terrorism," *Public Choice* 128, 147-168.
- Imseis, Ardi (2001). "'Moderate' Torture On Trial: The Israeli Supreme Court Judgment Concerning the Legality of the General Security Service Interrogation Methods," *International Journal of Human Rights* 5, 71-96.
- Kaplow, Louis and Shavell, Steven (1996). "Property Rules versus Liability Rules: An Economic Analysis," *Harvard Law Review* 109, 713-790.

- Kontorovich, Eugene (2004). "Liability Rules for Constitutional Rights: The Case of Mass Detentions," *Stanford Law Review* 56, 755-833.
- Leshem, Shmuel (2010). "The Benefits of a Right to Silence for the Innocent," *RAND Journal of Economics* 41, 398-416.
- Luban, David (2005). "Liberalism, Torture, and the Ticking Bomb," *Virginia Law Review* 91, 1425-1461.
- Mialon, Hugo M. (2005). "An Economic Theory of the Fifth Amendment," *RAND Journal of Economics* 36, 834-849.
- Mialon, Hugo M. and Rubin, Paul H. (2008). "The Economics of the Bill of Rights," *American Law and Economics Review* 10, 1-60.
- Posner, Richard (2002). "The Best Offense," *The New Republic* September 2, 28-31.
- Rejali, Darius (2007). *Torture and Democracy*, New Jersey: Princeton University Press.
- Rizzo, Mario J. and Whitman, Douglas Glen (2004). "The Camel's Nose is in the Tent: Rules, Theories, and Slippery Slopes," *UCLA Law Review* 51, 539-592.
- Sandler, Todd and Siqueira, Kevin (2006). "Global Terrorism: Deterrence versus Preemption," *Canadian Journal of Economics* 50, 1370-1387.
- Seidmann, Daniel J. and Stein, Alex (2000). "The Right to Silence Helps the Innocent: A Game-Theoretic Analysis of the Fifth Amendment Privilege," *Harvard Law Review* 114, 431-510.
- Shavell, Steven (2007). "Optimal Discretion in the Application of Rules," *American Law and Economics Review* 9, 175-194.
- Siqueira, Kevin and Sandler, Todd (2007). "Terrorist Backlash, Terrorism Mitigation, and Policy Delegation," *Journal of Public Economics* 91, 1800-1815.
- Sobel, Joel (2000). "A Model of Declining Standards," *International Economic Review* 41, 295-303.
- Strauss, Marcy (2003). "Torture," *New York Law School Law Review* 48, 201-274.
- The Economist (2006), "Saying No to Torture," *The Economist*, October 20.
- The Pew Research Foundation (2009), "America's Place in the World 2009: An Investigation of Public and Leadership Opinion About International Affairs," December 2009, Available at Website <http://people-press.org/reports/pdf/569.pdf>.
- Volokh, Eugene (2003). "The Mechanisms of the Slippery Slope," *Harvard Law Review* 116, 1026-1134.
- Waldron, Jeremy (2005). "Torture and Positive Law: Jurisprudence for the White House," *Columbia Law Review* 105, 1681-1750.
- Wantchekon, Leonard and Healy, Andrew (1999). "The 'Game' of Torture," *Journal of Conflict Resolution* 43, 596-609.
- Wickelgren, Abraham L. (2010). "A Right to Silence for Civil Defendants?" *Journal of Law, Economics, and Organization* 26, 92-114.