

Deterring Transgressions at the Apex of Government*

Weijia Li[†] Yang Xie[‡]

April 23, 2026

Abstract

We analyze a model of deterring transgressions by the head of government, in which the conditions determining whether a transgression will be punished only become clear after the transgression occurs. Under a general assumption that these conditions evolve stochastically and path-dependently, we show that, since effective deterrence prevents new information about the underlying conditions from being revealed, the threat of punishment may over time become not credible enough to deter any transgressions that are sufficiently profitable for the leader. We demonstrate that several institutional or cultural solutions either fail to fully address this problem or face inherent limitations. The same mechanism applies to deterring other undesirable behaviors when the precise conditions for punishment only become clear after the fact, such as international aggression, and to maintaining other preventive policies whose success could undermine their perceived necessity, such as preemptive measures against discrimination, and vaccination efforts.

*We thank Yunus Aybas, Scott Ashworth, Cuimin Ba, Matheus Bandeira, Chris Bidner, Alberto Bisin, Patrick Bolton, Michael Callen, Kyle Chauvin, Liuchun Deng, Gabriele Gratton, Guojun He, Da Huang, Aziz Huq, Alexander Jakobsen, Ye Jin, Matthew Kahn, Ethan Kaplan, Urme Khan, Andrew Little, Zhao Liu, Zhaotian Luo, Aprajit Mahajan, Roger Myerson, Paulo Natenzon, Mariann Ollár, Marcin Peški, Ariel Roginsky, Gérard Roland, Konstantin Sonin, Yufeng Sun, Michael Ting, Scott Tyson, Fan Wang, Shaoda Wang, Melanie Wasserman, Yinxi Xie, Siyang Xiong, Dingxin Zhao, Congyi Zhou, and participants in seminars at Monash, NYU Shanghai, Seikei, and THU and 2023 SIOE and 2024 GAMES Meetings for their valuable comments. Earlier versions of the paper were titled “Forever Deterred, the Highest Crimes? A Difficulty and Potential Solutions” or “Can High Crimes Be Perpetually Deterred? A Difficulty and Potential Solutions.”

[†]Department of Economics, Monash University; weijia.li@monash.edu.

[‡]Department of Economics, University of California, Riverside; yang.xie@ucr.edu.

1 Introduction

Transgressions by leaders at the apex of government against legal or political norms can impose “truly massive costs” on society (Huq, 2018, p. 1510; Rivera et al., 2025). Deterring such acts often relies on political accountability mechanisms, such as elections, impeachment, or backlash from key supporters (e.g., Madison, 1845, p. 528; Manin et al., 1999; Myerson, 2008; Bidner and Francois, 2013; Myerson, 2019; Ginsburg et al., 2021). On political accountability, a vast literature has enlightened us on the problems imposed by the hidden type and actions of the office holder – thus the issues of selection and moral hazard – and coordination problems in implementing these accountability mechanisms.¹ These problems are of central importance, and the literature has provided many valuable insights on potential solutions or remedies to them.

In this paper, we propose a mechanism by which the deterrence of apex transgressions through these political accountability mechanisms may still fail, even when the aforementioned problems are absent or perfectly solved. Instead, the mechanism we propose arises when, heuristically speaking, the political threat of accountability fails to remain credible over time, since constant deterrence achieved by the credible threat also means that the threat has not been carried out for a possibly very long period of time.

We formalize and analyze this mechanism in a simple model of presidential transgression and impeachment. In the model, a President (P) in each period chooses whether to transgress, and a Congress (C) decides whether to punish the transgressing P through impeachment. Main insights from the model also apply to deterring transgressions by top leaders in non-presidential or non-democratic regimes, as P can represent the leader in question – a prime minister, general secretary of the ruling party, or monarch; impeachment thus stands for the relevant political accountability mechanism, which C implements.

In this model, we abstract away from selection issues, moral hazard, and coordination problems, by assuming that P has only one type, all actions of the players are observable by

¹For examples of influential surveys on political accountability with selection issues, moral hazard, or both, see Persson and Tabellini (2000, Ch. 4), Besley (2005, 2006), Ashworth (2012), Gehlbach (2013, Ch. 7), Gailmard (2014), Duggan and Martinelli (2017), and Dal Bó and Finan (2018); examples of individual studies are not limited to Ferejohn (1986), McCubbins et al. (1987), Austen-Smith and Banks (1989), Banks and Sundaram (1993, 1998), Besley and Case (1995), Coate and Morris (1995), Lohmann (1998), Fearon (1999), Berganza (2000), Duggan (2000), Canes-Wrone et al. (2001), Gailmard (2002, 2009), Bernhardt et al. (2004), Maskin and Tirole (2004), Ashworth (2005), Besley and Smart (2007), Fox (2007), Banks and Duggan (2008), Snyder and Ting (2008), Fox and Shotts (2009), Bonfiglioli and Gancia (2013), Smart and Sturm (2013), Ashworth and Bueno de Mesquita (2014), Myerson (2015), Aghion and Jackson (2016), Ashworth et al. (2017), Aytimur and Bruns (2019), Duggan and Martinelli (2020), Gratton and Lee (2024), Kasamatsu and Kishishita (2024a,b). For examples addressing coordination problems, see Kuran (1989), Lohmann (1994a,b), Weingast (1997), Aragonès et al. (2007), Myerson (2008), Persson and Tabellini (2009), Fearon (2011), Shadmehr and Bernhardt (2011), Gitmez et al. (2023), and Francois and Bidner (2024).

everyone, and C is a single player. Instead, the key feature of our model is that the incentive for C to impeach the transgressing P – referred to as the state of the world – follows a stochastic process, and the current state of the world is only known to the players after P transgresses. This feature is meant to capture a reality that arises from the political nature of impeachment, or more broadly, punishment for transgressions by the head of government: among all the factors behind whether such punishment will be carried out, many may be known or predictable before the transgression occurs; yet there is always a component, however small it may be, that is difficult to fully predict *ex ante* – an ambiguity that is “intrinsic to democratic politics and human existence” (Olsen, 2014, p. 107). We provide more discussion on this point when we introduce our model in Section 2.

This feature implies that, in our model, successful deterrence prevents new empirical knowledge about the state of the world from emerging. As a result, over time, in a generally uncertain yet path-dependent world, this lack of new information renders historical knowledge from the times of past transgressions increasingly irrelevant to the current state of the world. Consequently, the credibility of the threat of impeachment becomes increasingly akin to a coin toss with a fixed probability. Therefore, as long as transgression is sufficiently profitable for P, there will always come a time when the threat becomes not credible enough to deter transgression. Deterrence is thus dynamically self-undermining. This formalizes the mechanism we propose.

Further analysis reveals that this mechanism arises distinctively from our model: it is transgressions by the head of government, not misconduct or crimes in ordinary settings, that are in question; the state of the world evolves in an uncertain yet path-dependent manner; and the uncertainty centers on whether C will impeach a transgressing P, rather than P’s payoffs from transgressing. We also show that this mechanism remains even if part of the state of the world is perfectly predictable *ex ante*, provided that it is bounded; the problem the mechanism induces in deterring apex transgressions also applies more strongly when the transgression in question is truly exceptional and fundamentally distinct from those minor ones.

We then examine several potential institutional or cultural solutions to the problem within our framework, finding that each of them faces limitations:

- First, while equipping impeachment with a sufficiently severe punishment can make perpetual deterrence possible, the required severity may incentivize P to avoid punishment, to the extent that it may encourage the transgressing P to execute a self-coup. The required severity of punishment may also be too harsh to be feasible in a civilized society, as P’s gain from apex transgressions can be enormous. Our analysis also suggests that it could be difficult for outsiders of the model to tell whether or not

the punishment is indeed sufficiently severe by only observing an ongoing period of successful deterrence.

- Second, admitting many non-partisan members to Congress with a super-majority rule for conviction, which is common in practice, may exacerbate the problem, making perpetual deterrence impossible even with an arbitrarily harsh punishment.
- Third, establishing an additional punitive authority, such as an independent judiciary, can only partially alleviate the problem.
- Fourth, subjecting the transgressions in question to many overlapping punitive authorities, as in federalism, could be a solution to the problem, but this requires all these authorities to be independent of each other – a condition that may be unrealistic – and risks creating political conflict due to inconsistent decisions or procedures.
- Finally, fostering a political culture against transgressions by the head of government cannot fully address the problem, as long as the influence of such efforts decays, however slightly, over time.

Besides selection issues, moral hazard, and coordination problems, there are another two classic obstacles to achieving political accountability, both contained in Gehlbach (2013, p. 158–161), which captures the insight of Barro (1973), one of the earliest models in the literature. The first is that full accountability will be compromised only if the leader in question prefers transgressing and subsequently being held accountable for it to not transgressing in the first place. Along this line, Ferejohn (1986) shows that even partial accountability is in danger as the leader can play the electorates against each other, while Persson et al. (1997) suggest that full accountability is still achievable as separation of powers can play the high offices against each other. Different from the literature, in the mechanism we analyze, even if the leader prefers not transgressing to being held accountable for transgressing, deterrence with the accountability mechanism will still fail eventually, as long as the reward from not being held accountable for transgressing is large enough; also, separation of powers in the form of having an independent judiciary can only partially alleviate the problem.

The other is the tool for accountability being quite “lumpy,” in that the leader can only be either reelected or removed. Along this line, the literature has shown that adopting or committing to a mixed strategy may help to solve the problem (e.g., Meirowitz, 2007; Aghion and Jackson, 2016). That said, in our model, C randomizing impeachment is generally not credible, since everyone, including C, would come to know the state of the world when C chooses whether to impeach the transgressing P, such that C would generally strictly prefer

either of the two options, and everyone knows this. The mixed-strategy approach would thus not help to solve the problem we characterize.

The literature on political accountability, especially in a dynamic electoral framework, has also called for more effort to allow for a state variable evolving over time (e.g., survey by Duggan and Martinelli, 2017, p. 980), and recent contributions have allowed current policy choices to influence the future state of the world (e.g., Duggan, 2012; Battaglini, 2014; Duggan and Forand, 2025). Compared with this thread of literature, we abstract away from such influence, but allow the current policy choice, i.e., whether or not P transgresses, to affect whether the state of the world is revealed, and analyze how such effect would present a problem for constant deterrence of apex transgressions over time.

Recent history shows that even long-established, advanced democracies are not immune to transgressions by their leaders against democratic institutions and norms (e.g., Bowman, 2024; Stokes, 2025). This challenges the conventional wisdom that prolonged, successful democratic governance fosters immunity to such transgressions.² We reconcile this observation with the literature by showing that constant deterrence of such transgressions can become its own enemy, a possibility rooted in their unique setting, i.e., the apex of government, and the political nature of their punishment.

Interpreting transgressions in our model as attacks on democratic norms and executive constraints would also link our paper to the growing literature on democratic backsliding and executive transgressions (e.g., survey by Grillo et al., 2024). The literature has focused on the roles played by the informational or political disadvantages of the restrainer or electorate, coordination failure among the restrainers, permanent effects of holding the transgressor accountable, and current transgressions' electoral implications for the future.³ Compared with the literature, the mechanism we analyze does not depend on any informational or political disadvantage of the restrainer, is not caused by potential coordination failure among the restrainers, and can still happen even when the punishment for transgression is transient and everyone is myopic. Moreover, the mechanism can, by itself, generate stochastic cycles between lengthy periods of constant deterrence and consecutive transgressions, i.e., backsliding, and the revealed state of the world at each transgression indicates for how long the

²Examples are not limited to Przeworski (2006), Persson and Tabellini (2009), Fearon (2011), Bidner and Francois (2013), Acemoglu and Jackson (2015), Adriani and Sonderegger (2018), Besley and Persson (2019), Acemoglu and Robinson (2022), Acemoglu et al. (2025), Kasamatsu and Kishishita (2024b), Weyland (2024), and Fuchs and Fukuda (2025).

³For examples of studies on the disadvantages of the restrainers, see Acemoglu et al. (2013), Svobik (2020), Miller (2021), Grillo and Prato (2023), Luo and Przeworski (2023), Gratton and Lee (2024), and Chiopris et al. (2025); on coordination failure, Francois and Bidner (2024); on permanent effects of accountability, Howell et al. (2023); on future electoral implications, Howell and Wolton (2018), Helmke et al. (2022), and Hollyer et al. (Forthcoming).

threat of punishment can remain credible to prevent backsliding. This makes each transgression a critical juncture, i.e., a brief moment when deep uncertainty about future political trajectories may be resolved (e.g., Capoccia and Kelemen, 2007; Acemoglu and Robinson, 2012; Callen et al., 2024).

In the mechanism we analyze, the state of the world is hidden by a desirable outcome (successful deterrence), and a prolonged period of this desirable outcome may eventually lead to its opposite (transgression). This mechanism also applies to deterring other undesirable behaviors when the precise conditions for punishment only become clear after the fact, and to other contexts where the success of a preventive policy may hinder updated understanding of its necessity, thus threatening its continuation. We discuss three examples of such applications at the end of the paper: deterring international aggression, upholding or not the United States Voting Rights Act of 1965, and vaccine hesitancy.

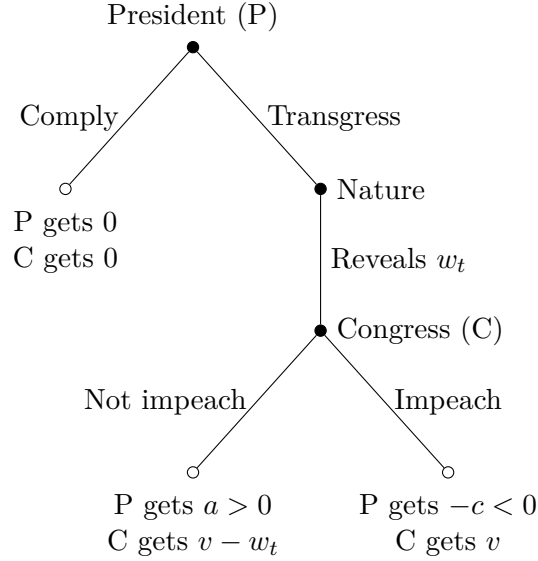
The paper proceeds as follows. Section 2 presents the baseline model. Section 3 analyzes it and delivers the main result. Section 4 discusses alternative setups of the model. Section 5 examines proposed institutional and cultural solutions. Section 6 concludes by discussing applications to other contexts. Proofs, technical details, and extensions are provided in the Appendix.

2 Baseline Model

In each period t , a President (P) and a Congress (C) interact in a game as in Figure 1. P first decides whether to comply with or transgress the law and norms governing the presidency. If P complies, both P and C receive a payoff of zero, and the period ends, transitioning to period $t + 1$.

If P transgresses, C will subsequently decide whether to punish P through impeachment. Modeling impeachment as punishment aligns with the constitutional law literature, which views punishment as a central purpose of impeachment (e.g., Tribe and Matz, 2018; Crespo, 2019). If C impeaches P, P incurs an exogenous cost, $c > 0$, while C receives an exogenous payoff, v . If C chooses not to impeach P, P receives an exogenous, positive gain, $a > 0$, while C receives $v - w_t$, where w_t represents the net incentive for C to impeach P. The period then ends, transitioning to period $t + 1$.

State of the world. We interpret w_t as the state of the world at period t . This approach aligns with insights from comparative and constitutional law. Impeachment, or more broadly, punishment for transgressions by the head of government, is ultimately a political solution to political crises where the transgressions in question have significantly undermined



Players know time, $s < t$, and state of the world, w_s , at last transgression.

Figure 1: Baseline model, period t

public confidence in the incumbent or threatened constitutional governance. Whether such punishment will be carried out is thus primarily a political question, and its answer depends on many specific circumstances at the time, whether political, social, economic, or moral. This holds true even when supposedly apolitical experts, such as prosecutors and judges, are involved (e.g., Story, 1833; Baumgartner and Kada, 2003; Tribe and Matz, 2018; Crespo, 2019; Ohnesorge, 2020; Ginsburg et al., 2021; Congressional Research Service, 2023; Bowman, 2024; Monaghan et al., 2024).

Conditional revelation. The dependence on these specific circumstances suggests that, while many factors, such as the alignment of electoral interests between P and C, may be known or predictable before a transgression occurs, there is always a component of these circumstances that is difficult to fully predict *ex ante*. Additionally, constitutional law is inherently general and often leaves punishable transgressions by the head of government incompletely defined, making it difficult to predict beforehand how a transgression will be addressed (e.g., Story, 1833, p. 287; Huq, 2018, p. 1515; Tribe and Matz, 2018, p. 45; Bowman, 2024, p. 99).⁴ As a result, one today has to estimate whether a transgression will

⁴A notable example is Section 4 of Article II of the Constitution of the United States, which states: “[t]he President, Vice President and all civil Officers of the United States, shall be removed from Office on Impeachment for, and Conviction of, Treason, Bribery, or other high Crimes and Misdemeanors.” The term “high Crimes and Misdemeanors” is a “catchall phrase” and “not defined in the Constitution or statutes,” creating “considerable vagueness” (Black and Bobbitt, 2018, p. 26; Congressional Research Service, 2023, p. 839). For further literature, see Tribe (1998), Tribe and Matz (2018), and Bowman (2024). For theoretical

be punished, while relying on her understanding of the circumstances of past transgressions (e.g., Tribe and Matz, 2018; Congressional Research Service, 2023, p. 839; Bowman, 2024).⁵

This leads to the key feature of our model: the conditional revelation of the state of the world. As in Figure 1, at the start of period t , the players do not know the exact value of w_t ; once P chooses to transgress, nature reveals w_t to them. In Section 4.2, we discuss the robustness of our analysis to including an additional known component of the state of the world at the start of each period.

Apex of government vs. ordinary settings. Such conditional revelation distinguishes transgressions at the apex of government from misconduct and crimes in ordinary settings. For such misconduct or crimes, players typically have abundant ex ante knowledge about the state of the world, as there are many alternative settings to the setting in question to learn from. In Section 4.1, we explore the dynamics of transgression and punishment when the transgression occurred in an ordinary setting, not at the apex of government.

Evolution of the state of the world. At the start of period t , the players know when past transgressions occurred and the corresponding states of the world. How they use this information depends on how w_t evolves. We assume that w_t follows a Markov process,

$$w_t = \rho w_{t-1} + \epsilon_t, \text{ where } 0 < \rho \leq 1, \text{ and } \epsilon_t \sim \mathcal{N}(0, \sigma^2), \text{ i.i.d., with } \sigma > 0. \quad (1)$$

Here, ϵ_t s represent contemporary shocks, which are mutually independent and identically and normally distributed, with a mean of zero. This reflects our approach that w_t captures the component of the state of the world that is unknown and difficult to predict ex ante. We assume $\sigma > 0$ and $0 < \rho \leq 1$ because the state of the world relevant to the transgressions in question is most reasonably uncertain yet path-dependent, but not oscillatory or explosive. In particular, we allow the influence of historical shocks to either decay ($0 < \rho < 1$) or persist ($\rho = 1$) over time. In Sections 4.4 and 5.5, we discuss alternative assumptions about the evolution of w_t , such as alternative ranges of ρ or σ , or including a drift term.

Complete the setup. For simplicity, we assume that players are myopic, caring only about their payoffs within the current period. This is justifiable if each period represents

arguments for constitutional law to remain silent on this issue, see Huq (2018) and Crespo (2019).

⁵In the context of the United States, the *Constitution Annotated* attests: “impeachment is essentially a political process that is largely unreviewable by the Judicial Branch. As such, the historical practice of impeachment proceedings ...informs our understanding of the Constitution’s meaning in this area. In this vein, the meaning of ‘high crimes and misdemeanors’ is informed ...by the history of congressional impeachments” (Congressional Research Service, 2023, p. 839). The statement is followed by a 24-page review of that history (Congressional Research Service, 2023, p. 840–863).

the finite term of office for the players. We discuss the robustness of our results to forward-looking players in Section 4.6.

We solve the model using backward induction. For simplicity, we assume that C will not impeach a transgressing P if indifferent between impeaching and not impeaching, and P will comply if indifferent between transgressing and complying.

3 Analysis and Main Result

We now solve the game at period t . C will impeach a transgressing P if and only if $w_t > 0$. Since w_t follows a Markov process, the time of the last transgression, $s < t$, and the state of the world at that time, w_s , provide all the information P needs to estimate w_t at the beginning of period t . The estimated probability of punishment is $\mathbf{P}[w_t > 0 \mid w_s, s]$. Thus, P transgresses if and only if

$$a \cdot (1 - \mathbf{P}[w_t > 0 \mid w_s, s]) - c \cdot \mathbf{P}[w_t > 0 \mid w_s, s] > 0, \quad (2)$$

or equivalently, if the conditional probability of punishment is sufficiently low:

$$\mathbf{P}[w_t > 0 \mid w_s, s] < a/(a + c). \quad (3)$$

We first examine how the conditional probability of punishment, $\mathbf{P}[w_t > 0 \mid w_s, s]$, depends on both the time elapsed since the last transgression, $t - s$, and the state of the world revealed at that time, w_s . The evolution of w_t implies that

$$w_t = \rho^{t-s} w_s + \sum_{\tau=s+1}^t \rho^{t-\tau} \epsilon_\tau. \quad (4)$$

The conditional distribution of w_t is thus

$$(w_t \mid w_s, s) \sim \begin{cases} \mathcal{N}(\rho^{t-s} w_s, (1 - \rho^{2(t-s)}) \cdot \sigma^2 / (1 - \rho^2)) & \text{if } \rho \in (0, 1); \\ \mathcal{N}(w_s, (t - s) \cdot \sigma^2) & \text{if } \rho = 1. \end{cases} \quad (5)$$

This implies that the conditional probability of punishment is

$$\mathbf{P}[w_t > 0 \mid w_s, s] = \begin{cases} \Phi\left(\rho^{t-s} w_s / \left(\sigma \cdot \sqrt{(1 - \rho^{2(t-s)}) / (1 - \rho^2)}\right)\right) & \text{if } \rho \in (0, 1); \\ \Phi(w_s / (\sigma \cdot \sqrt{t - s})) & \text{if } \rho = 1, \end{cases} \quad (6)$$

where $\Phi(\cdot)$ is the cumulative distribution function of the standard normal distribution. This

conditional probability has the following properties:

Lemma 1 (Credibility of impeachment threat). *It is more likely than not for C to impeach a transgressing P if and only if C impeached P at the last transgression. As time passes with constant deterrence, the threat to punish increasingly resembles a coin toss with a fixed probability, i.e., one half. Mathematically, for any $t > s$,*

- *if $w_s \leq 0$, then $\mathbf{P}[w_t > 0 \mid w_s, s] \leq 0.5$, increases with $t - s$, and converges to 0.5 as $t - s$ increases;*
- *if $w_s > 0$, then $\mathbf{P}[w_t > 0 \mid w_s, s] > 0.5$, decreases with $t - s$, and converges to 0.5 as $t - s$ increases.*

The intuition behind Lemma 1 is as follows. First, why does whether C impeached P at the last transgression, i.e., whether $w_s > 0$, indicate whether it is more likely than not for C to impeach a transgressing P now? This is because, since the world is path-dependent and no new empirical knowledge about the state of the world has been revealed since the last transgression, the current estimate of the state of the world must be, adjusted with the degree of path-dependence, centered on the last revealed state.

Second, why is the monotonic convergence to a coin toss with a fixed probability? In the case of $\rho \in (0, 1)$, this is because the influence of historical shocks decays over time, which is set mechanically by $\rho \in (0, 1)$. As time passes with constant deterrence, because of the conditional revelation of the state of the world, no new information about the state of the world has been revealed. The conditional distribution of the current state of the world thus converges to a normal distribution with a mean of zero and a finite variance, i.e.,

$$\text{if } \rho \in (0, 1), \quad \text{then } (w_t \mid w_s, s) \xrightarrow{d} \mathcal{N}(0, \sigma^2 / (1 - \rho^2)) \text{ as } t - s \rightarrow \infty. \quad (7)$$

This makes the threat of impeachment increasingly resemble a 50/50 coin toss.

In the case of $\rho = 1$, even though the influence of historical shocks would remain constant, its relevance would dwindle, still. This is because, without new revelation, the uncertainty in the state of the world piles up over time and is not bounded from above. The conditional distribution of the state of the world thus converges to a normal distribution with an infinite variance, i.e.,

$$\text{if } \rho = 1, \quad \text{then } (w_t \mid w_s, s) \sim \mathcal{N}(w_s, (t - s)\sigma^2), \text{ where } (t - s)\sigma^2 \rightarrow \infty \text{ as } t - s \rightarrow \infty. \quad (8)$$

This makes the threat of impeachment increasingly resemble a 50/50 coin toss, too.

Lemma 1 implies that, to understand the dynamics of impeachment and transgression under constant deterrence, we need to compare the critical threshold, $a/(a + c)$, with the conditional probability of punishment in the limit over time with constant deterrence, which is one-half. On that, $a/(a + c)$ has the following property:

Lemma 2 (Payoff threshold). *A threat of impeachment equivalent to a 50/50 coin toss cannot deter P from transgressing, if and only if the potential gain from transgression exceeds the punishment from impeachment. Mathematically,*

$$0.5 < a/(a + c) \text{ if and only if } a > c. \quad (9)$$

The intuition behind Lemma 2 is straightforward. If P perceives the threat of impeachment as a 50/50 coin toss, he will simply compare the gain from transgression with the punishment, and will transgress if and only if the gain is greater, however slightly.

Combining Lemmas 1 and 2 leads to our main result:

Proposition 1 (A problem in deterring apex transgressions). *For any $a > c$, define $\bar{w} \equiv (\sigma/\rho) \cdot \Phi^{-1}(a/(a + c)) > 0$. Then,*

1. *if $w_s \leq 0$, then C will not impeach P at period s, and P will transgress at s + 1;*
2. *if $0 < w_s < \bar{w}$, then C will impeach P at s, but P will still transgress at s + 1;*
3. *if $w_s \geq \bar{w}$, then C will impeach P at s, and there exists a period in the finite future, $T \in [s + 2, \infty)$, such that P will comply from s + 1 to T - 1 but will transgress at T.*

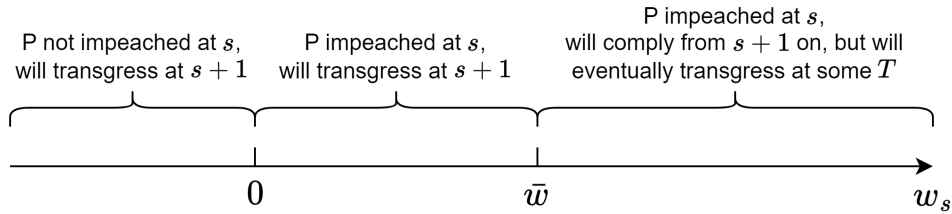


Figure 2: Proposition 1: Perpetual deterrence impossible for any $a > c$

Figure 2 visualizes Proposition 1. We prove the proposition in Appendix A. The intuition is as follows. About Claim 1, by Lemma 1, an unpunished transgression ($w_s \leq 0$) suggests that the impeachment threat in the following period is no more credible than a 50/50 coin toss. By Lemma 2, this is not enough to deter another transgression, whenever the potential gain from transgression exceeds the punishment ($a > c$).

About Claim 2, having the last transgression punished ($w_s > 0$) may not be enough to deter transgression in the following period, since, by Lemma 1, the credibility of the threat of impeachment may drop too much even in just one period. Therefore, the last revealed state of the world has to be sufficiently against transgression ($w_s \geq \bar{w}$) to deter another transgression in the period immediately after the last transgression.

About Claim 3, by Lemma 1, even if the threat of impeachment deters transgression right after the last transgression ($w_s \geq \bar{w}$), as time passes with constant deterrence, the threat increasingly resembles a 50/50 coin toss. By Lemma 2, a threat equivalent to a 50/50 coin toss cannot deter transgression, as long as the potential gain from transgression exceeds the punishment. Therefore, there must exist a period when the threat will eventually become not credible enough to deter transgression.

Proposition 1 thus formalizes the mechanism we propose, presenting a problem in deterring transgressions at the apex of government: for any transgression that is sufficiently profitable, constant deterrence would over time make the threat of punishment as non-credible as a coin toss with a fixed probability, i.e., one half in our baseline model, eventually leading to a failure of deterrence.

Stochastic cycles and critical junctures. Proposition 1 also implies that stochastic cycles between constant deterrence and consecutive transgressions can emerge. At each transgression, if the revelation is strongly against transgression, i.e., $w_s \geq \bar{w}$, where s is the current period, not only will the current transgression be punished, but deterrence of future transgressions can also be secured for potentially many subsequent periods. Indeed, we show in Appendix A that T in Proposition 1 increases with w_s .

If, instead, the revealed state of the world is not so much against transgression, i.e., $w_s < \bar{w}$ – in particular, if it prevents C from impeaching P, i.e., $w_s \leq 0$ – then another transgression will occur right after the current one. As the world is unlikely to have moved away so much within such a short time, the revealed state of the world at this following transgression may still be not so much against transgression, i.e., $w_{s+1} < \bar{w}$. In this way, a prolonged period of consecutive transgressions may follow. Each transgression can thus be viewed as a critical juncture, indicating future dynamics of transgressions and punishment.

Numerical illustration. To illustrate such dynamics, we provide a numerical simulation of our model. Figure 3 plots the typical dynamics over 500 periods, with the payoffs set to satisfy $a > c$. In this simulation, the initial state of the world w_1 is revealed to everyone before period 2 and is set exogenously to be sufficiently against transgression, i.e., $w_1 \geq \bar{w} > 0$, ensuring that P does not transgress at period 2. Consistent with Lemma 1, the conditional

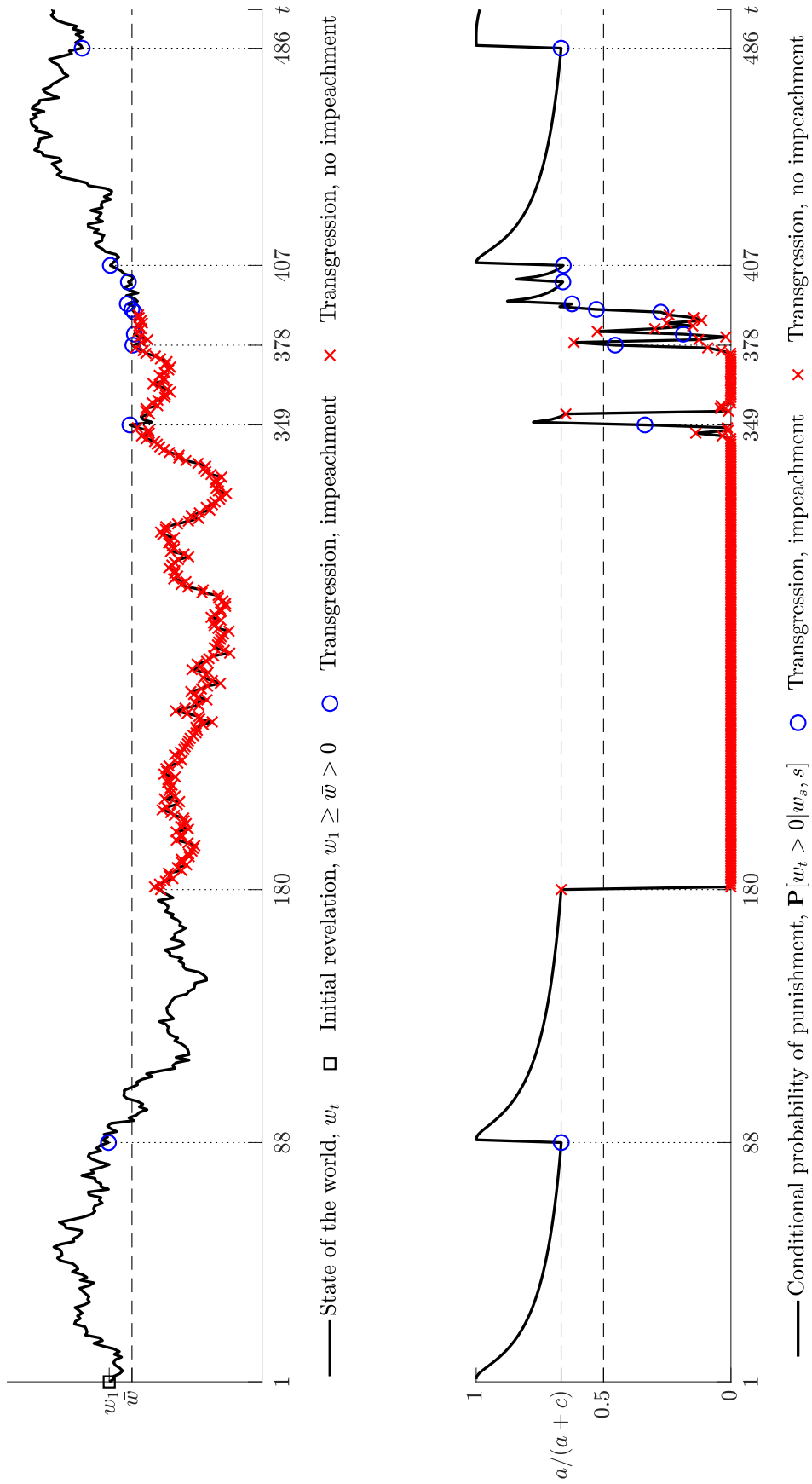


Figure 3: Typical dynamics with transgression gain exceeding punishment ($a > c$)

probability of punishment declines over time with constant deterrence. Consistent with Proposition 1, at period 88, the threat of impeachment becomes not credible enough to deter transgression, i.e., $\mathbf{P}[w_{88} > 0 \mid w_1, 1] < a/(a + c)$. At this point, the revealed state of the world is strongly against transgression, i.e., $w_{88} \geq \bar{w}$, leading to another lengthy period of deterrence. Yet, consistent with Lemma 1, the conditional probability of punishment still declines over time due to constant deterrence.

At period 180, the threat of punishment again becomes not credible enough to deter transgression. This time, the revealed state of the world is too low even for C to impeach P, i.e., $w_{180} \leq 0$. Consistent with Proposition 1, this leads to another transgression at period 181. The revealed state of the world again prevents C from impeaching P. This pattern of consecutive transgressions without being punished continues until period 349.

From period 349 onward, transgressions still occur frequently. There are periods, such as period 378, when the transgressing P is impeached, but the revealed state of the world cannot sustain deterrence. It is only at period 407 that the world has become sufficiently against transgression, leading to a lengthy period of deterrence until period 486. The revelation at period 486 is even more strongly against transgression, starting a longer period of deterrence. That said, as Proposition 1 predicts, another transgression will eventually occur.

4 Alternative Setups

To better understand the main result, we discuss several alternative setups of the model:

4.1 Not at the Apex of Government

What if the transgression in question is not committed by the head of government, but in an ordinary setting? As discussed in Section 2, this would correspond to a setup where the players always know the value of w_t at the start of period t , which we formalize in Appendix B. In this case, constant deterrence would not imply a lack of new empirical knowledge about the state of the world. The potential transgressor's decision would only depend on w_t ; conditional on w_t , past transgression and punishment decisions would have no effect on the current transgression decision. Thus, deterrence would not be self-undermining over time as in Lemma 1 and Proposition 1.

4.2 An Additional Known Component

What if part of the state of the world is known or perfectly predictable ex ante? In Appendix C, we introduce an additional component to the net incentive for C to impeach P that is

independent of w_t and publicly known at the beginning of each period. We show that, as long as this component is bounded, i.e., if social knowledge about the state of the world is never too wild ex ante, over time with constant deterrence, the credibility of the threat of impeachment will still converge to a probability strictly lower than one. Therefore, the problem we characterize will persist, though the required payoff threshold, i.e., $a > c$ in Proposition 1, may change.

4.3 Testing the State of the World by Minor Transgressions

What if P could test the state of the world by committing a minor transgression in each period before deciding whether to proceed with the “real” transgression? We formalize this idea as an extension of the baseline model in Appendix D. In this extension, in addition to the baseline setup, P commits a costless, minor “test” transgression at the start of each period, which generates an additional public signal about the state of the world. We assume that the precision of this signal depends on how similar the real and test transgressions are in nature. Both players then process these signals using Bayesian updating.

We show that, at one extreme, if the real transgression is as minor as the test transgression, the additional signal in each period would perfectly reveal the net incentive for C to impeach P if the real transgression were committed. In this case, the game’s dynamics converge to those of deterring a misconduct or crime in an ordinary setting. At the other extreme, if the real transgression is egregious to an unthinkable degree, the signal from the minor test transgression would reveal very little about the current state of the world. Here, the game’s dynamics converge to those of our baseline model. Therefore, the problem presented in Proposition 1 applies more strongly to scenarios where the transgression in question is truly exceptional and fundamentally distinct from those minor ones.

4.4 Alternative Ways of Evolution of the State of the World

We have assumed that the state of the world is uncertain yet path-dependent, but not oscillatory or explosive, i.e., $\sigma > 0$ and $0 < \rho \leq 1$. What if it evolves differently? We discuss these alternatives below, with technical details provided in Appendix E.

Deterministic world, i.e., $\sigma = 0$. In this case, although constant deterrence would still not bring new empirical knowledge about the state of the world, the revelation at the last transgression would perfectly indicate the current state. Once a transgression is punished, P would understand that any future transgression will also be punished for sure, leading to perpetual deterrence. Thus, the problem presented in Proposition 1 would not exist.

Path-independent world, i.e., $\rho = 0$. In this case, although constant deterrence would still not bring new empirical knowledge about the state of the world, the lack of new empirical knowledge would not become more consequential over time, and P's decision to transgress would be independent of history. Thus, deterrence would not be self-undermining over time as in Lemma 1 and Proposition 1.

Oscillatory world, i.e., $\rho < 0$. In an oscillatory world, absent new empirical knowledge about the state of the world, a threat of impeachment that is more credible than a 50/50 coin toss in one period will become less credible than a 50/50 coin toss in the following period. Therefore, if the gain from transgression exceeds the punishment ($a > c$), deterrence would never last for more than one period. Thus, the problem presented in Proposition 1 would have no opportunity to emerge.

Explosive world, i.e., $\rho > 1$. In this case, as time passes with constant deterrence, the expectation and standard deviation of the state of the world, conditional on its value and time at the last transgression, explode at about the same rate, ρ . Therefore, once the state of the world is revealed to be sufficiently against transgression, deterrence can last forever. Thus, the problem presented in Proposition 1 would not arise.

4.5 Payoff Uncertainty

What if the conditional revelation were about P's payoffs from transgressing, rather than C's strategy to punish? In Appendix F, we examine this alternative setup. We assume that C impeaches P with an exogenous probability, while P's gain and punishment for transgression evolve (with the same degree of path-dependence) and are subject to conditional revelation, similar to w_t in the baseline model.

We show that perpetual deterrence is possible in this case. Although the payoff uncertainty accumulates over time with constant deterrence, it does not affect the expectation of the relative gain from transgression to punishment. Therefore, it will not affect whether P will comply or transgress. As a result, deterrence can last forever. The mechanism we propose thus arises from strategic uncertainty, not payoff uncertainty.

4.6 Forward-looking Players

What if the players in our model cared not only about their current payoffs, but also their future payoffs? In Appendix G, we analyze a setting in which any impeached P will be replaced by a new P, while C and all Ps are infinitely forward-looking. In the analysis, we

focus on symmetric, pure-strategy Markov-perfect equilibria, where “symmetric” means that all Ps adopt the same strategy.

We show that the problem presented in Proposition 1 still emerges whenever the one-time gain from transgression exceeds the net present value of the punishment from impeachment. This is because, first, about C’s strategy, the state of the world evolves exogenously, and the current state is determined and revealed before C decides whether to impeach a transgressing P. Thus, C’s decision about impeachment does not affect the future state of the world, so a forward-looking C behaves just like a myopic one in our baseline model.

Given that, for each forward-looking P, the intuition behind Proposition 1 still applies. Not only that; in equilibrium, P’s continuation value following an unpunished transgression must be weakly positive, as P can always choose to unconditionally comply from then on, receiving a zero continuation value. This weakly positive continuation value creates an additional incentive for P to transgress now, making constant deterrence even harder to sustain. Therefore, having forward-looking players makes the problem presented in Proposition 1 even more significant.

Summary. These discussions clarify that the mechanism we propose in this paper arises from 1) a lack of new empirical knowledge about the world due to constant deterrence, and the world being 2) uncertain and 3) path-dependent. Having significant transgressions at the apex of government in question that are fundamentally different from those minor ones implies a lack of new empirical knowledge under constant deterrence; uncertainty makes this lack of knowledge consequential, and path-dependence amplifies its impact over time, eventually making the threat of punishment not credible enough to deter transgression. This self-undermining nature of deterrence is rooted in strategic uncertainty, not payoff uncertainty, and remains robust even when the state of the world includes a known, bounded component, and when the players are forward-looking.

5 Potential Institutional or Cultural Solutions

Having formalized and analyzed the mechanism we propose and the problem it induces in deterring transgressions by the head of government, we now explore potential institutional and cultural solutions. As we will see, each of them faces limitations.

5.1 Severe Punishment

Proposition 1 addresses the case where P’s gain from transgression exceeds the punishment ($a > c$). Could a sufficiently severe punishment solve the problem? In Appendix H, we show that perpetual deterrence is possible if the punishment outweighs the gain from transgression ($c \geq a$). This is because, even though the threat of impeachment converges to a 50/50 coin toss over time with constant deterrence, it would still be sufficient for deterrence if the punishment were severe enough.

But relying on severe punishment as a solution has limitations. First, a classic idea in law and criminology argues that, as famously put by Beccaria (2008, p. 51), “impunity itself arises from the atrocity of punishments.” The more severe the punishment, the harder the offender will try to avoid it, and the more reluctant the punitive authority may become when considering whether to implement the punishment. Consequently, the punishment will be applied less consistently (e.g., Beccaria, 2008, p. 49–57; Sanchirico, 2006; Howell et al., 2023). Not only that, to preempt the punishment, a transgressing P may even choose to execute a self-coup, committing another transgression (e.g., Cameron, 1998, p. 134). Thus, independent of Proposition 1, having a severe punishment can make the threat of punishment less credible, or even encourage further transgressions.

Second, the gain from transgressions by leaders at the apex of government is often enormous in political value, especially when the objective is to control the highest office of government, which wields the power of legitimate violence or proclaiming states of exception (e.g., Weber, 2004; Schmitt, 1985; Agamben, 2005). For the punishment to outweigh such gains, it may need to be extremely harsh, making it infeasible in a civilized society (e.g., Beccaria, 2008, p. 85; Durkheim, 1973; Benabou and Tirole, 2026).⁶

A third limitation is that, when deterrence has been successful in the past, it may be difficult for outsiders of the game of transgression and punishment, such as social scientists, journalists, or concerned citizens, to determine whether the punishment is indeed sufficiently severe. Our model helps to elaborate on this point. As shown in Appendix H, an ongoing period of constant deterrence can occur in both cases ($a > c$ and $a \leq c$). Thus, observing an ongoing period of constant deterrence alone cannot reveal whether the punishment is sufficiently harsh to ensure perpetual deterrence.

⁶As a famous example, Priscus Attalus, a leading Roman senator, had a thumb and a forefinger maimed in 416 for collaborating with the Visigoths and usurping the Roman emperorship in 409 and 414 (e.g., Bury, 1889, p. 150). Such a punishment, combined with public humiliation, “was intended as a lesson for other Roman aristocrats of the dangers of independent action” (Salzman, 2021, p. 108). As a contemporary exception, in 2026, prosecutors sought the death penalty for Yoon Suk Yeol, a former president of South Korea, on charges of insurrection, after the 2024 South Korean martial law crisis; he was eventually found guilty and sentenced to life imprisonment (Rashid, 2026).

5.2 Wisdom of the Crowd

In our baseline model, C is a single player. As “the many are better judges ...when they come together” (Aristotle, 1998, p. 83), could the wisdom of the crowd, i.e., letting C consist of many non-partisan members, overcome the mechanism we propose and analyze? In Appendix I, we analyze this setting by incorporating multiple members into Congress, each with their own net incentive to convict a transgressing P, evolving and subject to conditional revelation as w_t in the baseline model. To ensure that these members are non-partisan, we assume that these net incentives are mutually independent. We also assume that a transgressing P will be punished if and only if he is convicted by a vote in Congress, given a voting rule for conviction, and that each member votes sincerely.

We show that, given any super-majority rule for conviction, which is often required by the seriousness of the matter in question, having many non-partisan members in Congress exacerbates the problem presented in Proposition 1. This is because, as in the baseline model, over time with constant deterrence, each Congress member will vote to convict a transgressing P almost like a 50/50 coin toss. Since these members are non-partisan, the Law of Large Numbers applies. Thus, in the limit over time, as the size of Congress increases, it becomes almost certain that only about half of the members will vote to convict, failing to clear the super-majority threshold.⁷ Therefore, even if the punishment is arbitrarily harsh and exceeds the gain from transgression, this almost entirely non-credible threat of punishment will still fail to deter P from transgressing.⁸ In this sense, for a wide range of reasonable voting rules for conviction, such as any super-majority rule, the wisdom of the crowd does not help overcome the mechanism we propose and analyze.⁹

⁷For example, in this limit, the probability of having at least two thirds of a Congress of 50 non-partisan members voting to convict a transgressing P is 0.0077; the probability for a 100-member Congress is 0.0004. One can also read this result as applying a special case of Condorcet’s jury theorem (Condorcet, 1785; Grofman et al., 1983, p. 264), in which each voter votes for the correct decision with probability 0.5, to super-majority voting rules.

⁸Bolton and Rosenthal (2002) recognize that a democratic process (majority or super-majority voting) can certify occasions (widespread financial distress) about which economic contracts (debts) are incomplete and warrant intervention (moratoria), and they show that such an ex post remedy, when anticipated, may improve efficiency of resource allocation ex ante. Along this line, we recognize that, on occasions (transgressions by the head of government) about which social contracts (constitutions) can be incomplete, a democratic process (super-majority voting) can be used to certify a proper intervention (impeachment or not), and we show that such an ex post remedy, when anticipated, can make such occasions inevitable. This provides an interpretation in the framework of incomplete contracts for this result.

⁹One may entertain the idea of having a voting rule that requires a super-majority for acquittal, i.e., a minority rule for conviction. Indeed, we show in Appendix I that, in the limit over time with constant deterrence, under such minority rule, as the size of Congress increases, it becomes almost certain that any transgressing P will be convicted, thus deterring transgression. Yet, in reality, such voting rules may encourage wrongful conviction of complying Ps.

5.3 Additional Independent Punitive Authority

In our baseline model, C is the only institutional authority that can hold P accountable. What if there were another independent punitive authority, such as an independent judiciary, to punish the transgressing P if impeachment fails? In Appendix J, we analyze this setting by introducing another player, the Judiciary (J), who decides whether to rule against a transgressing P with the same punishment, c . We assume that J has her own net incentive to rule against P, which evolves and is subject to conditional revelation, just like C's net incentive to impeach P. We also assume that J is independent, meaning her net incentive to rule against P and C's net incentive to impeach P are mutually independent.

We show that having such J can only partially alleviate the problem presented in Proposition 1. This is because, over time with constant deterrence, the additional threat of punishment from J also becomes as non-credible as a 50/50 coin toss. Thus, in the limit, the overall probability of punishment increases to 0.75, which is higher than 0.5 but still strictly less than one. As a result, although the required threshold for the gain from transgression increases from $a > c$ to $a > 3c$, the problem presented in Proposition 1 still exists.

5.4 Many Overlapping Punitive Authorities

As having an additional independent punitive authority can partially address the problem presented in Proposition 1, what if there were many such authorities, so that a transgressing P would be punished if at least one of them decides to act against him? This scenario is exemplified by federalism, where the federal government and each state government are considered separate sovereigns. Under the "separate sovereigns" doctrine, even if the transgression is committed only once, it can be addressed multiple times, each time by the judiciary of one of these separate sovereigns.

In Appendix K, we analyze this setting, where a transgressing P is punished if C impeaches him or at least one of many separate judiciaries rules against him. We assume that the net incentive for each judiciary to rule against P evolves and is subject to conditional revelation, just like C's net incentive to impeach P.

We show that, under certain conditions, the problem presented in Proposition 1 can be fully addressed in this setting. If all judiciaries and C are independent, meaning that their net incentives to act against P are mutually independent, then, in the limit over time with constant deterrence, the probability of eventual punishment is $1 - 0.5^{N+1}$, where N is the total number of these separate judiciaries. As this number grows, this probability approaches certainty, ensuring that P will comply regardless of the relative size of the punishment and the gain from transgression. In this way, the problem presented in Proposition 1 can be

overcome with many overlapping punitive authorities.

That said, this solution depends critically on the mutual independence of all overlapping punitive authorities, which is difficult to hold in a polarized political environment. For example, if all these punitive authorities can be categorized into only three independent groups, the probability of punishing a transgressing P in the limit over time with constant deterrence would be $1 - 0.5^3 = 7/8$. Thus, if the gain from transgression is sufficiently large, i.e., for any $a > 7c$, the problem presented in Proposition 1 would still exist.

Even if mutual independence were achievable, another concern can arise from having many overlapping punitive authorities. As seen in federalism, the lack of uniformity in state procedures and decisions may lead to political chaos at the federal level, as noted by the Supreme Court of the United States in *Trump v. Anderson* (2024). These limitations highlight the challenges of implementing this potential solution in practice.

5.5 Cultivating Political Culture

In light of these findings, a cultural approach is worth considering, and the literature of constitutional law has been advocating for persistent efforts promoting a political culture against the transgression in question (e.g., Huq, 2018). In our framework, we can model the effect of such efforts as an additional known, positive drift every period, $\mu > 0$, in how the state of the world evolves:

$$w_t = \mu + \rho w_{t-1} + \epsilon_t, \text{ where } \mu > 0, 0 < \rho \leq 1, \text{ and } \epsilon_t \sim \mathcal{N}(0, \sigma^2), \text{ i.i.d., with } \sigma > 0. \quad (10)$$

In Appendix L, we show that the cumulative effect of such efforts in the long run would still be too limited to fully address the problem presented in Proposition 1, as long as the state of the world is subject to a mean-reverting force ($0 < \rho < 1$), such as a natural decay, however slightly, over time. This is because, as the effect of the effort in each period decays over time, in the limit over time with constant deterrence, the cumulative effect of such efforts would shift the conditional distribution of w_t only finitely, i.e., by $\mu/(1 - \rho)$. Thus, the threat of punishment would still resemble a coin toss with a fixed probability, and the problem presented in Proposition 1 would still persist.

6 Concluding Remarks: Other Applications

In this paper, we model the political nature of punishing transgressions at the apex of government through the conditional revelation of the conditions that determine whether a

transgression will be punished. Our analysis highlights a problem in deterring such transgressions, even when selection issues, moral hazard, or coordination problems are absent or perfectly solved: the very success of deterrence can lead to a lack of new empirical knowledge about the underlying conditions, gradually eroding the credibility of punishment over time. We demonstrate that several institutional or cultural solutions either fail to fully solve the problem or face inherent limitations.

Beyond transgressions by the head of government, our model also applies to deterring other undesirable behaviors when the exact conditions determining whether such behaviors will get punished only become clear after such behaviors have been committed. For instance, in the context of international aggression, an aggressor may be poised to invade another country, but can be deterred by potential intervention from the international community. That said, the conditions determining whether such intervention will be carried out may only become clear after the invasion occurs. Our analysis suggests that, if these conditions evolve path-dependently yet stochastically and the prize of invasion is sufficiently big for the aggressor, then after a period of successful deterrence, the aggressor will eventually invade as the threat of international intervention gradually loses credibility, and consecutive aggressions may follow once the international community does not intervene resolutely.

Beyond deterring undesirable behaviors, the insight from our analysis can shed light on other contexts where the success of a preventive measure may obscure empirical evidence of its necessity, jeopardizing its continuation. For example, Section 4 of the United States Voting Rights Act of 1965 provided a formula to identify jurisdictions requiring federal preclearance for changes to voting laws – a measure designed to combat racial discrimination in voting (e.g., Cascio and Washington, 2014). This means that the latent tendency of jurisdictions to enact discriminatory voting laws would only become observable if preclearance were removed or rendered ineffective. If this tendency evolves path-dependently yet stochastically, our analysis suggests that the Supreme Court of the United States might eventually perceive preclearance as unnecessary and move to strike it down or significantly weaken it.

Indeed, in *Shelby County v. Holder* (2013), the Court’s majority held that Section 4 of the Voting Rights Act is unconstitutional, arguing that “the conditions that originally justified these measures [of preclearance] no longer characterize voting in the covered jurisdictions.” In her dissent, Justice Ginsburg countered that voting equality improvements were observed in these jurisdictions precisely because “the preclearance remedy [has been] in place in the covered jurisdictions [and] has worked and is continuing to work to stop discriminatory changes.” Supporting this view, recent studies have documented negative relative impacts of this ruling on voter registration, turnout, and city-council representation among racial or ethnic minorities (e.g., De Rienzo, 2022; Ricca and Trebbi, 2022; Billings et al., 2024).

As another example, the public health literature has long observed that vaccine hesitancy – “delay in acceptance or refusal of vaccination” – can result from vaccination complacency, which occurs when “perceived risks of vaccine-preventable diseases are low and vaccination is not deemed a necessary preventive action” (MacDonald et al., 2015, p. 4161–4162). In particular, “[i]mmunization programme success may, paradoxically, result in complacency and ultimately, hesitancy, as individuals [perceive] the disease the vaccine prevents [as] no longer common” (MacDonald et al., 2015, p. 4163).

In this context, the conditions determining whether an outbreak would occur in an under-vaccinated population might not be fully predictable until the population has indeed become under-vaccinated. When these conditions evolve path-dependently but remain uncertain, the mechanism we analyze could apply: prolonged success of vaccination may eventually lead some people to forgo vaccination as its perceived necessity diminishes. This dynamic, potentially exacerbated by misinformation and conspiracy theories (e.g., Vergano, 2025), can have severe public health and economic consequences, as simulations show in the context of the United States (e.g., Lo and Hotez, 2017; Kiang et al., 2025).¹⁰

References

- Acemoglu, Kamer Daron, Nicolás Ajzenman, Cevat Giray Aksoy, Martin Fiszbein, and Carlos Molina. 2025. (Successful) democracies breed their own support. *Review of Economic Studies* **92**: 621–655.
- Acemoglu, Kamer Daron, and Matthew O. Jackson. 2015. History, expectations, and leadership in the evolution of social norms. *The Review of Economic Studies* **82**: 423–456.
- Acemoglu, Kamer Daron, and James A. Robinson. 2012. *Why Nations Fail: The Origins of Power, Prosperity and Poverty*. New York: Crown Publishing Group.
- Acemoglu, Kamer Daron, and James A. Robinson. 2022. Non-modernization: Power–culture trajectories and the dynamics of political institutions. *Annual Review of Political Science* **25**: 323–339.
- Acemoglu, Kamer Daron, James A. Robinson, and Ragnar Torvik. 2013. Why do voters dismantle checks and balances? *Review of Economic Studies* **80**: 845–875.

¹⁰Lo and Hotez (2017, p. 887) warn that “[a] 5% decline in MMR vaccine coverage in the United States would result in an estimated 3-fold increase in measles cases for children aged 2 to 11 years nationally every year, with an additional \$2.1 million in public sector costs”; “[u]nder a 50% decline in childhood vaccination in the US,” Kiang et al. (2025, p. 2177) predict “51.2 million measles cases over a 25-year period, 9.9 million rubella cases, 4.3 million poliomyelitis cases, 197 diphtheria cases, 10.3 million hospitalizations, and 159200 deaths.”

- Adriani, Fabrizio, and Silvia Sonderegger. 2018. The signaling value of punishing norm-breakers and rewarding norm-followers. *Games* **9**: 102.
- Agamben, Giorgio. 2005. *State of Exception*. Chicago: University of Chicago Press.
- Aghion, Philippe, and Matthew O. Jackson. 2016. Inducing leaders to take risky decisions: Dismissal, tenure, and term limits. *American Economic Journal: Microeconomics* **8**: 1–38.
- Aragonès, Enriqueta, Andrew Postlewaite, and Thomas Palfrey. 2007. Political reputations and campaign promises. *Journal of the European Economic Association* **5**: 846–884.
- Aristotle. 1998. *Politics*. Indianapolis: Hackett Publishing Company.
- Ashworth, Scott. 2005. Reputational dynamics and political careers. *Journal of Law, Economics, and Organization* **21**: 441–466.
- Ashworth, Scott. 2012. Electoral accountability: Recent theoretical and empirical work. *Annual Review of Political Science* **15**: 183–201.
- Ashworth, Scott, and Ethan Bueno de Mesquita. 2014. Is voter competence good for voters? Information, rationality, and democratic performance. *American Political Science Review* **108**: 565–587.
- Ashworth, Scott, Ethan Bueno de Mesquita, and Amanda Friedenberg. 2017. Accountability and information in elections. *American Economic Journal: Microeconomics* **9**: 95–138.
- Austen-Smith, David, and Jeffrey Banks. 1989. Electoral accountability and incumbency. In Ordeshook, Peter C. (ed.) *Models of Strategic Choice in Politics*. Ann Arbor: University of Michigan Press, 121–148.
- Aytimur, R. Emre, and Christian Bruns. 2019. Accountability with large electorates. *Economic Journal* **129**: 1529–1560.
- Banks, Jeffrey S., and John Duggan. 2008. A dynamic model of democratic elections in multidimensional policy spaces. *Quarterly Journal of Political Science* **3**: 269–299.
- Banks, Jeffrey S., and Rangarajan K. Sundaram. 1993. Adverse selection and moral hazard in a repeated elections model. In Barnett, William A., Melvin J. Hinich, and Norman J. Schofield (eds.) *Political Economy: Institutions, Competition and Representation*. New York: Cambridge University Press, 295–312.
- Banks, Jeffrey S., and Rangarajan K. Sundaram. 1998. Optimal retention in agency problems. *Journal of Economic Theory* **82**: 293–323.
- Barro, Robert J. 1973. The control of politicians: An economic model. *Public Choice* **14**: 19–42.
- Battaglini, Marco. 2014. A dynamic theory of electoral competition. *Theoretical Economics* **9**: 515–554.

- Baumgartner, Jody C., and Naoko Kada (eds.). 2003. *Checking Executive Power: Presidential Impeachment in Comparative Perspective*. Westport: Praeger Publishers.
- Beccaria, Cesare. 2008. *On Crimes and Punishments and Other Writings*. Toronto: University of Toronto Press.
- Benabou, Roland, and Jean Tirole. 2026. Laws and norms. *Journal of Political Economy* **134**: 731–772.
- Berganza, Juan Carlos. 2000. Two roles for elections: Disciplining the incumbent and selecting a competent candidate. *Public Choice* **105**: 165–194.
- Bernhardt, Dan, Sangita Dubey, and Eric Hughson. 2004. Term limits and pork barrel politics. *Journal of Public Economics* **88**: 2383–2422.
- Besley, Timothy. 2005. Political selection. *Journal of Economic Perspectives* **19**: 43–60.
- Besley, Timothy. 2006. *Principled Agents? The Political Economy of Good Government*. Oxford: Oxford University Press.
- Besley, Timothy, and Anne Case. 1995. Incumbent behavior: Vote seeking, tax setting, and yardstick competition. *American Economic Review* **85**: 25–45.
- Besley, Timothy, and Torsten Persson. 2019. Democratic values and institutions. *American Economic Review: Insights* **1**: 59–76.
- Besley, Timothy, and Michael Smart. 2007. Fiscal restraints and voter welfare. *Journal of Public Economics* **91**: 755–773.
- Bidner, Chris, and Patrick Francois. 2013. The emergence of political accountability. *Quarterly Journal of Economics* **128**: 1397–1448.
- Billings, Stephen B., Noah Braun, Daniel B. Jones, and Ying Shi. 2024. Disparate racial impacts of *Shelby County v. Holder* on voter turnout. *Journal of Public Economics* **230**: 105047.
- Black, Charles L., Jr., and Philip Bobbitt. 2018. *Impeachment: A Handbook, New Edition*. New Haven: Yale University Press.
- Bolton, Patrick, and Howard Rosenthal. 2002. Political intervention in debt contracts. *Journal of Political Economy* **110**: 1103–1134.
- Bonfiglioli, Alessandra, and Gino Gancia. 2013. Uncertainty, electoral incentives and political myopia. *Economic Journal* **123**: 373–400.
- Bowman, Frank O., III. 2024. *High Crimes and Misdemeanors: A History of Impeachment for the Age of Trump*. Cambridge: Cambridge University Press, 2nd edition.
- Bury, John Bagnell. 1889. *A History of the Later Roman Empire: From Arcadius to Irene (395 A.D. to 800 A.D.)*, Vol. I. London: Macmillan.

- Callen, Michael, Jonathan L. Weigel, and Noam Yuchtman. 2024. Experiments about institutions. *Annual Review of Economics* **16**: 105–131.
- Cameron, Maxwell A. 1998. Self-coups: Peru, Guatemala, and Russia. *Journal of Democracy* **9**: 125–139.
- Canes-Wrone, Brandice, Michael C. Herron, and Kenneth W. Shotts. 2001. Leadership and pandering: A theory of executive policymaking. *American Journal of Political Science* **45**: 532–550.
- Capoccia, Giovanni, and R. Daniel Kelemen. 2007. The study of critical junctures: Theory, narrative, and counterfactuals in historical institutionalism. *World Politics* **59**: 341–369.
- Cascio, Elizabeth U., and Ebonya Washington. 2014. Valuing the vote: The redistribution of voting rights and state funds following the Voting Rights Act of 1965. *Quarterly Journal of Economics* **129**: 379–433.
- Chiopris, Caterina, Monika Nalepa, and Georg Vanberg. 2025. A wolf in sheep’s clothing: Citizen uncertainty and democratic backsliding. *Journal of Politics* **87**: 1272–1287.
- Coate, Stephen, and Stephen Morris. 1995. On the form of transfers to special interests. *Journal of Political Economy* **103**: 1210–1235.
- Condorcet, Marquis de. 1785. *Essai sur l’Application de l’Analyse à la Probabilité des Décisions Rendus à la Pluralité des Voix*. Paris: Imprimerie Royale.
- Congressional Research Service. 2023. *The Constitution of the United States of America: Analysis and Interpretation: Analysis of Cases Decided by the Supreme Court of the United States to June 30, 2022*. Washington: United States Government Publishing Office.
- Crespo, Andrew Manuel. 2019. Impeachment as punishment. *Harvard Law and Policy Review* **13**: 579–592.
- Dal Bó, Ernesto, and Frederico Finan. 2018. Progress and perspectives in the study of political selection. *Annual Review of Economics* **10**: 541–575.
- De Rienzo, Salvatore M., Jr. 2022. *Shelby County v. Holder* and changes in voting behavior. *American Economist* **67**: 195–210.
- Duggan, John. 2000. Repeated elections with asymmetric information. *Economics and Politics* **12**: 109–135.
- Duggan, John. 2012. Noisy stochastic games. *Econometrica* **80**: 2017–2045.
- Duggan, John, and Jean Guillaume Forand. 2025. Accountability in Markovian elections. *Games and Economic Behavior* **151**: 183–217.
- Duggan, John, and César Martinelli. 2017. The political economy of dynamic elections: Accountability, commitment, and responsiveness. *Journal of Economic Literature* **55**:

916–984.

- Duggan, John, and César Martinelli. 2020. Electoral accountability and responsive democracy. *Economic Journal* **130**: 675–715.
- Durkheim, Émile. 1973. Two laws of penal evolution. *Economy and Society* **2**: 285–308.
- Fearon, James D. 1999. Electoral accountability and the control of politicians: Selecting good types versus sanctioning poor performance. In Przeworski, Adam, Susan C. Stokes, and Bernard Manin (eds.) *Democracy, Accountability, and Representation*. Cambridge: Cambridge University Press, 55–97.
- Fearon, James D. 2011. Self-enforcing democracy. *Quarterly Journal of Economics* **126**: 1661–1708.
- Ferejohn, John. 1986. Incumbent performance and electoral control. *Public choice* **50**: 5–25.
- Fox, Justin. 2007. Government transparency and policymaking. *Public choice* **131**: 23–44.
- Fox, Justin, and Kenneth W. Shotts. 2009. Delegates or trustees? A theory of political accountability. *Journal of Politics* **71**: 1225–1237.
- Francois, Patrick, and Chris Bidner. 2024. The problem with authoritarian populists. *Studies in Microeconomics* **12**: 59–73.
- Fuchs, William, and Satoshi Fukuda. 2025. Shaping institutions. Working paper, Santa Clara University.
- Gailmard, Sean. 2002. Expertise, subversion, and bureaucratic discretion. *Journal of Law, Economics, and Organization* **18**: 536–555.
- Gailmard, Sean. 2009. Multiple principals and oversight of bureaucratic policy-making. *Journal of Theoretical Politics* **21**: 161–186.
- Gailmard, Sean. 2014. Accountability and principal-agent theory. In Bovens, Mark, Robert E. Goodin, and Thomas Schillemans (eds.) *The Oxford Handbook of Public Accountability*. Oxford: Oxford University Press, 90–105.
- Gehlbach, Scott. 2013. *Formal Models of Domestic Politics*. New York: Cambridge University Press.
- Ginsburg, Tom, Aziz Z. Huq, and David Landau. 2021. The comparative constitutional law of presidential impeachment. *University of Chicago Law Review* **88**: 81–164.
- Gitmez, A. Arda, James A. Robinson, and Mehdi Shadmehr. 2023. Missing discussions: Institutional constraints in the Islamic political tradition. National Bureau of Economic Research Working Paper 30916.
- Gratton, Gabriele, and Barton E. Lee. 2024. Liberty, security, and accountability: The rise and fall of illiberal democracies. *Review of Economic Studies* **91**: 340–371.

- Grillo, Edoardo, Zhaotian Luo, Monika Nalepa, and Carlo Prato. 2024. Theories of democratic backsliding. *Annual Review of Political Science* **27**: 381–400.
- Grillo, Edoardo, and Carlo Prato. 2023. Reference points and democratic backsliding. *American Journal of Political Science* **67**: 71–88.
- Grofman, Bernard, Guillermo Owen, and Scott L. Feld. 1983. Thirteen theorems in search of the truth. *Theory and Decision* **15**: 261–278.
- Helmke, Gretchen, Mary Kroeger, and Jack Paine. 2022. Democracy by deterrence: Norms, constitutions, and electoral tilting. *American Journal of Political Science* **66**: 434–450.
- Hollyer, James, Marko Klačnja, and Rocío Titiunik. Forthcoming. Charismatic leaders and democratic backsliding. *Journal of Politics* .
- Howell, William G., Kenneth A. Shepsle, and Stephane Wolton. 2023. Executive absolutism: The dynamics of authority acquisition in a system of separated powers. *Quarterly Journal of Political Science* **18**: 243–275.
- Howell, William G., and Stephane Wolton. 2018. The politician’s province. *Quarterly Journal of Political Science* **13**: 119–146.
- Huq, Aziz Z. 2018. Legal or political checks on apex criminality: An essay on constitutional design. *UCLA Law Review* **65**: 1506–1530.
- Kasamatsu, Satoshi, and Daiki Kishishita. 2024a. Does informative opposition influence electoral accountability? *Quarterly Journal of Political Science* **19**: 459–498.
- Kasamatsu, Satoshi, and Daiki Kishishita. 2024b. Endogenous political trust and electoral accountability. *Journal of Politics* **86**: 358–363.
- Kiang, Mathew V., Kate M. Bubar, Yvonne Maldonado, Peter J. Hotez, and Nathan C. Lo. 2025. Modeling reemergence of vaccine-eliminated infectious diseases under declining vaccination in the US. *JAMA: The Journal of the American Medical Association* **333**: 2176–2187.
- Kuran, Timur. 1989. Sparks and prairie fires: A theory of unanticipated political revolution. *Public Choice* **61**: 41–74.
- Lo, Nathan C., and Peter J. Hotez. 2017. Public health and economic consequences of vaccine hesitancy for measles in the United States. *JAMA Pediatrics* **171**: 887–892.
- Lohmann, Susanne. 1994a. The dynamics of informational cascades: The Monday demonstrations in Leipzig, East Germany, 1989–91. *World Politics* **47**: 42–101.
- Lohmann, Susanne. 1994b. Information aggregation through costly political action. *American Economic Review* **84**: 518–530.
- Lohmann, Susanne. 1998. Rationalizing the political business cycle: A workhorse model.

Economics and Politics **10**: 1–17.

- Luo, Zhaotian, and Adam Przeworski. 2023. Democracy and its vulnerabilities: Dynamics of democratic backsliding. *Quarterly Journal of Political Science* **18**: 105–130.
- MacDonald, Noni E., Juhani Eskola, Xiaofeng Liang, Mohuya Chaudhuri, Eve Dubé, Bruce Gellin, Susan Goldstein, Heidi Larson, Mahamane Laouali Manzo, Arthur Reingold, Kinzang Tshering, Yuqing Zhou, Robb Butler, Philippe Duclos, Sherine Guirguis, Ben Hickler, and Melanie Schuster. 2015. Vaccine hesitancy: Definition, scope and determinants. *Vaccine* **33**: 4161–4164.
- Madison, James, Jr. 1845. *Debates on the Adoption of the Federal Constitution, with a Diary of the Debates of the Congress of the Confederation, Volume V*. Washington: United States Congress. Revised and arranged by Jonathan Elliot.
- Manin, Bernard, Adam Przeworski, and Susan C. Stokes. 1999. Elections and representation. In Przeworski, Adam, Susan C. Stokes, and Bernard Manin (eds.) *Democracy, Accountability, and Representation*. Cambridge: Cambridge University Press, 29–54.
- Maskin, Eric, and Jean Tirole. 2004. The politician and the judge: Accountability in government. *American Economic Review* **94**: 1034–1054.
- McCubbins, Mathew D., Roger G. Noll, and Barry R. Weingast. 1987. Administrative procedures as instruments of political control. *Journal of Law, Economics, and Organization* **3**: 243–277.
- Meirowitz, Adam. 2007. Probabilistic voting and accountability in elections with uncertain policy constraints. *Journal of Public Economic Theory* **9**: 41–68.
- Miller, Michael K. 2021. A republic, if you can keep it: Breakdown and erosion in modern democracies. *Journal of Politics* **83**: 198–213.
- Monaghan, Chris, Matthew Flinders, and Aziz Z. Huq (eds.). 2024. *Impeachment in a Global Context: Law, Politics, and Comparative Practice*. New York: Routledge.
- Myerson, Roger B. 2008. The autocrat’s credibility problem and foundations of the constitutional state. *American Political Science Review* **102**: 125–139.
- Myerson, Roger B. 2015. Moral hazard in high office and the dynamics of aristocracy. *Econometrica* **83**: 2083–2126.
- Myerson, Roger B. 2019. The dilemma of presidential impeachment. *Perspectives on Economics and Civilization*, December 15, 2019.
- Ohnesorge, John. 2020. Comparing impeachment regimes. *Duke Journal of Comparative and International Law* **31**: 259–299.
- Olsen, Johan P. 2014. Accountability and ambiguity. In Bovens, Mark, Robert E. Goodin,

- and Thomas Schillemans (eds.) *The Oxford Handbook of Public Accountability*. Oxford: Oxford University Press, 106–123.
- Persson, Torsten, Gérard Roland, and Guido Tabellini. 1997. Separation of powers and political accountability. *Quarterly Journal of Economics* **112**: 1163–1202.
- Persson, Torsten, and Guido Tabellini. 2000. *Political Economics: Explaining Economic Policy*. Cambridge: MIT Press.
- Persson, Torsten, and Guido Tabellini. 2009. Democratic capital: The nexus of political and economic change. *American Economic Journal: Macroeconomics* **1**: 88–126.
- Przeworski, Adam. 2006. Self-enforcing democracy. In Wittman, Donald A., and Barry R. Weingast (eds.) *The Oxford Handbook of Political Economy*. New York: Oxford University Press, 312–328.
- Rashid, Raphael. 2026. South Korea’s former president Yoon Suk Yeol jailed for life for leading insurrection. *The Guardian*, February 19, 2026.
- Ricca, Federico, and Francesco Trebbi. 2022. Minority underrepresentation in U.S. cities. National Bureau of Economic Research Working Paper 29738.
- Rivera, Eduardo, Enrique Seira, and Saumitra Jha. 2025. Apex corruption erodes democratic values. Working paper, Stanford University.
- Salzman, Michele Renee. 2021. *The Falls of Rome: Crises, Resilience, and Resurgence in Late Antiquity*. Cambridge: Cambridge University Press.
- Sanchirico, Chris William. 2006. Detection avoidance. *New York University Law Review* **81**: 1331–1399.
- Schmitt, Carl. 1985. *Political Theology: Four Chapters on the Concept of Sovereignty*. Cambridge: MIT Press.
- Shadmehr, Mehdi, and Dan Bernhardt. 2011. Collective action with uncertain payoffs: Coordination, public signals, and punishment dilemmas. *American Political Science Review* **105**: 829–851.
- Smart, Michael, and Daniel M. Sturm. 2013. Term limits and electoral accountability. *Journal of Public Economics* **107**: 93–102.
- Snyder, James M., Jr., and Michael M. Ting. 2008. Interest groups and the electoral control of politicians. *Journal of Public Economics* **92**: 482–500.
- Stokes, Susan C. 2025. *The Backsliders: Why Leaders Undermine Their Own Democracies*. Princeton: Princeton University Press.
- Story, Joseph. 1833. *Commentaries on the Constitution of the United States*. Cambridge: Hilliard, Gray, and Company.

- Supreme Court of the United States. 2013. *Shelby County v. Holder*. In *United States Reports*, volume 570. Washington: Supreme Court of the United States, 529–594.
- Supreme Court of the United States. 2024. *Trump v. Anderson*. In *United States Reports*, volume 601. Washington: Supreme Court of the United States.
- Svolik, Milan W. 2020. When polarization trumps civic virtue: Partisan conflict and the subversion of democracy by incumbents. *Quarterly Journal of Political Science* **15**: 3–31.
- Tribe, Laurence H. 1998. Defining “high Crimes and Misdemeanors”: Basic principles. *George Washington Law Review* **67**: 712–734.
- Tribe, Laurence H., and Joshua Matz. 2018. *To End a Presidency: The Power of Impeachment*. New York: Basic Books.
- Vergano, Dan. 2025. The brainwashing campaign that is measles misinformation. *Scientific American*, April 30, 2025.
- Weber, Maximilian K. E. 2004. *The Vocation Lectures*. Indianapolis: Hackett Publishing Company.
- Weingast, Barry R. 1997. The political foundations of democracy and the rule of the law. *American Political Science Review* **91**: 245–263.
- Weyland, Kurt. 2024. *Democracy’s Resilience to Populism’s Threat: Countering Global Alarmism*. Cambridge: Cambridge University Press.

Appendix

A Proof of Proposition 1

Proof. To prove Claim 1, suppose $w_s \leq 0$. C will thus not impeach P at period s . By $w_s \leq 0$, $a > c$, and Lemmas 1 and 2, we have

$$\mathbf{P}[w_{s+1} > 0 \mid w_s, s] \leq 0.5 < a/(a+c). \quad (\text{A.1})$$

P will thus transgress at period $s+1$. Claim 1 is proved.

To prove Claim 2, suppose $w_s > 0$. First, C will thus impeach P at period s . Second, at period $s+1$, by the evolution of w_t in Equation (1), P will transgress if and only if

$$a/(a+c) > \mathbf{P}[w_{s+1} > 0 \mid w_s, s] = \Phi(\rho w_s/\sigma), \quad (\text{A.2})$$

or equivalently, if

$$w_s < (\sigma/\rho) \cdot \Phi^{-1}(a/(a+c)) \equiv \bar{w}. \quad (\text{A.3})$$

Note $\bar{w} > 0$ because of $a > c$. Claim 2 is thus proved.

To prove Claim 3, suppose $w_s \geq \bar{w}$. By $\bar{w} > 0$, we have $w_s > 0$, and C will thus impeach P at period s . By Claim 2, which is proved above, P will not transgress at period $s+1$. Now consider the periods onward. By $w_s > 0$ and Lemma 1, as time passes with constant deterrence, the conditional probability of punishment decreases and converges to 0.5; by Lemma 2, $0.5 < a/(a+c)$. Therefore, there must exist a period T , when $\mathbf{P}[w_t > 0 \mid w_s, s]$ is so close to 0.5 for the first time that $\mathbf{P}[w_t > 0 \mid w_s, s] < a/(a+c)$ will come to hold. Indeed, in the case of $\rho = 1$, we have

$$T = s + \left\lfloor w_s^2 / (\sigma \cdot \Phi^{-1}(a/(a+c)))^2 \right\rfloor + 1, \quad (\text{A.4})$$

where $\lfloor \cdot \rfloor$ is the floor function, and, by $w_s \geq \bar{w}$, we have $T \geq s+2$; in the case of $\rho \in (0, 1)$, we have

$$T = s + \left\lfloor \ln \left(\frac{(\sigma \cdot \Phi^{-1}(a/(a+c)))^2}{w_s^2(1-\rho^2) + (\sigma \cdot \Phi^{-1}(a/(a+c)))^2} \right) / 2 \ln \rho \right\rfloor + 1, \quad (\text{A.5})$$

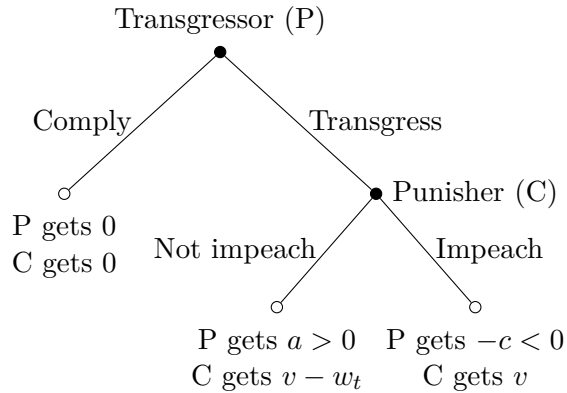
and, by $w_s \geq \bar{w}$, we have $T \geq s+2$, too. P will thus not transgress until period $T-1$ but will transgress at period T . Claim 3 is thus proved, and so is the proposition. \square

The proof also implies simple comparative statics of the length of constant deterrence:

Corollary A.1. *T in Proposition 1 increases with w_s and decreases with σ and a/c .*

B Not at the Apex of Government

Figure B.1 shows the structure of the game if the transgression in question were not at the apex of government, but a misconduct or crime in an ordinary setting. Compared with Figure 1, the only difference is that here, without the conditional revelation, the players already know w_t at the beginning of period t . By backward induction, the punisher will punish a transgression if and only if $w_t > 0$; knowing that, the potential transgressor will transgress if and only if $w_t > 0$, independent of the history of transgression and punishment. Therefore, deterrence would not be self-undermining in this setting.



Players know current state of the world, w_t .

Figure B.1: Not at the apex of government, period t

C An Additional Known Component

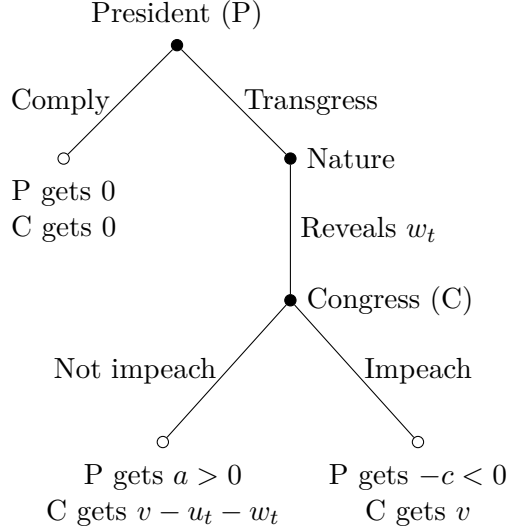
Now consider an alternative setup as in Figure C.1, where the net incentive for C to impeach a transgressing P in period t contains an additional component, u_t , known to the players at the beginning of period t . Since this component captures the “known,” i.e., not too wild, side of the state of the world, we further assume that this component is bounded from above, i.e., $u_t \leq \bar{u} < \infty$. Since w_t is still the unknown component, we assume that w_t still follows

$$w_t = \rho w_{t-1} + \epsilon_t, \text{ where } 0 < \rho \leq 1, \text{ and } \epsilon_t \sim \mathcal{N}(0, \sigma^2), \text{ i.i.d., with } \sigma > 0, \quad (\text{C.1})$$

and is independent of u_t .

Given this setup, by backward induction, C will impeach a transgressing P if and only if $w_t > -u_t$. Given this strategy of C, P will transgress if and only if

$$-c \cdot \mathbf{P}[w_t > -u_t | w_s, s, u_t] + a \cdot (1 - \mathbf{P}[w_t > -u_t | w_s, s, u_t]) > 0, \quad (\text{C.2})$$



Players know time, $s < t$, and state of the world, w_s , at last transgression, and u_t .

Figure C.1: An additional known component, period t

or equivalently, if

$$\mathbf{P}[w_t > -u_t | w_s, s, u_t] < a/(a + c). \quad (\text{C.3})$$

Given the evolution of w_t and its independence from u_t , we have

$$(w_t | w_s, s, u_t) \sim \begin{cases} \mathcal{N}(\rho^{t-s}w_s, (1 - \rho^{2(t-s)}) \cdot \sigma^2 / (1 - \rho^2)) & \text{if } \rho \in (0, 1), \\ \mathcal{N}(w_s, (t - s) \cdot \sigma^2) & \text{if } \rho = 1 \end{cases} \quad (\text{C.4})$$

Therefore, we have

$$\mathbf{P}[w_t > -u_t | w_s, s, u_t] = \begin{cases} \Phi\left(\frac{(\rho^{t-s}w_s + u_t)}{\left(\sigma \cdot \sqrt{\frac{1-\rho^{2(t-s)}}{1-\rho^2}}\right)}\right) & \text{if } \rho \in (0, 1), \\ \Phi\left(\frac{(w_s + u_t)}{(\sigma \cdot \sqrt{t-s})}\right) & \text{if } \rho = 1. \end{cases} \quad (\text{C.5})$$

By $u_t \leq \bar{u} < \infty$, we have, as $t - s \rightarrow \infty$,

$$\mathbf{P}[w_t > -u_t | w_s, s, u_t] \leq \begin{cases} \Phi\left(\frac{\rho^{t-s}w_s + \bar{u}}{\sigma \cdot \sqrt{\frac{1-\rho^{2(t-s)}}{1-\rho^2}}}\right) \rightarrow \Phi\left(\frac{\bar{u} \cdot \sqrt{1-\rho^2}}{\sigma}\right) < 1 & \text{if } \rho \in (0, 1), \\ \Phi\left(\frac{(w_s + \bar{u})}{(\sigma \cdot \sqrt{t-s})}\right) \rightarrow 0.5 & \text{if } \rho = 1. \end{cases} \quad (\text{C.6})$$

Therefore, if $\rho = 1$, whenever $a > c$, perpetual deterrence will be impossible; if $\rho \in (0, 1)$, whenever

$$a > c \cdot \Phi\left(\frac{\bar{u} \cdot \sqrt{1-\rho^2}}{\sigma}\right) / \left(1 - \Phi\left(\frac{\bar{u} \cdot \sqrt{1-\rho^2}}{\sigma}\right)\right) \quad (\text{C.7})$$

perpetual deterrence will be impossible, still. Therefore, the problem presented in Proposition 1 will still emerge even if we take into consideration an additional known component of the net incentive for C to impeach a transgressing P.

D Additional Signals from Test Transgressions

Now consider an extension of the baseline model, where we denote the severity of the “real” transgression in question as $\bar{x} \in R$ and the state of the world as $w_t(\bar{x})$. With this notation, the evolution of the state of the world becomes

$$w_t(\bar{x}) = \rho w_{t-1}(\bar{x}) + \epsilon_t, \text{ where } 0 < \rho \leq 1, \text{ and } \epsilon_t \sim \mathcal{N}(0, \sigma^2), \text{ i.i.d., with } \sigma > 0. \quad (\text{D.1})$$

We assume that, at the start of each period t and before the game is played following the baseline setup, P first commits a minor, “test” transgression of severity $\hat{x} < \bar{x}$. This test transgression is costless and generates a public signal of the state of the world,

$$w_t(\hat{x}) = w_t(\bar{x}) + e_t, \text{ where } e_t \sim \mathcal{N}(0, (\bar{x} - \hat{x})\zeta^2), \text{ i.i.d., with } \zeta > 0, \quad (\text{D.2})$$

and all ϵ_t s and e_t s are mutually independent. This signal-generating process implies that the more different the test transgression and the real transgression in question are, the more coarse this signal would be. At one extreme, when $\bar{x} - \hat{x} \rightarrow 0$, the real transgression in question is as minor as the test transgression, and the signal from the test transgression would precisely reveal the state of the world; at the other extreme, when $\bar{x} - \hat{x} \rightarrow \infty$, the real transgression is extremely different from the test transgression in their nature, and the signal would not reveal much of the state of the world.

We assume that all players will process the signals using Bayesian updating, and we want to understand how these signals may change the dynamics of the game. First, consider the end of period s , when the last real transgression happened with the state of the world revealed as $w_s(\bar{x})$. P’s posterior about $w_s(\bar{x})$ is

$$(w_s(\bar{x}) \mid w_s(\bar{x}), s) \sim \mathcal{N}(\mu_s^{\text{post}}, 1/p_s^{\text{post}}), \quad (\text{D.3})$$

where the mean and precision are, respectively,

$$\mu_s^{\text{post}} = w_s(\bar{x}), \quad p_s^{\text{post}} = \infty. \quad (\text{D.4})$$

Second, at the start of period $s + 1$, before receiving the signal, $w_{s+1}(\hat{x})$, while knowing

the evolution of the state of the world, P's prior about $w_{s+1}(\bar{x})$ is

$$(w_{s+1}(\bar{x}) | w_s(\bar{x}), s) \sim \mathcal{N}(\mu_{s+1}^{\text{prior}}, 1/p_{s+1}^{\text{prior}}), \quad (\text{D.5})$$

where the mean and precision are, respectively,

$$\mu_{s+1}^{\text{prior}} = \rho \mu_s^{\text{post}} = \rho w_s(\bar{x}), \quad p_{s+1}^{\text{prior}} = 1/(\rho^2/p_s^{\text{post}} + \sigma^2) = 1/\sigma^2. \quad (\text{D.6})$$

Since the signal follows

$$(w_{s+1}(\hat{x}) | w_{s+1}(\bar{x})) \sim \mathcal{N}(w_{s+1}(\bar{x}), (\bar{x} - \hat{x})\zeta^2), \quad (\text{D.7})$$

by Bayesian updating, P's posterior about $w_{s+1}(\bar{x})$ is

$$(w_{s+1}(\bar{x}) | w_s(\bar{x}), s, w_{s+1}(\hat{x})) \sim \mathcal{N}(\mu_{s+1}^{\text{post}}, 1/p_{s+1}^{\text{post}}), \quad (\text{D.8})$$

where the mean and precision are, respectively,

$$\mu_{s+1}^{\text{post}} = \frac{p_{s+1}^{\text{prior}} \cdot \mu_{s+1}^{\text{prior}} + w_{s+1}(\hat{x}) / (\bar{x} - \hat{x})\zeta^2}{p_{s+1}^{\text{prior}} + 1/(\bar{x} - \hat{x})\zeta^2}, \quad p_{s+1}^{\text{post}} = p_{s+1}^{\text{prior}} + 1/(\bar{x} - \hat{x})\zeta^2. \quad (\text{D.9})$$

We can do this up to the start of period t , assuming that no real transgression has happened from period $s+1$ to $t-1$. At the start of period t , P's prior about $w_t(\bar{x})$ is

$$(w_t(\bar{x}) | w_s(\bar{x}), s, w_{s+1}(\hat{x}), \dots, w_{t-1}(\hat{x})) \sim \mathcal{N}(\mu_t^{\text{prior}}, 1/p_t^{\text{prior}}), \quad (\text{D.10})$$

where the mean and precision are, respectively,

$$\mu_t^{\text{prior}} = \rho \mu_{t-1}^{\text{post}}, \quad p_t^{\text{prior}} = 1/(\rho^2/p_{t-1}^{\text{post}} + \sigma^2). \quad (\text{D.11})$$

Since the signal follows

$$(w_t(\hat{x}) | w_t(\bar{x})) \sim \mathcal{N}(w_t(\bar{x}), (\bar{x} - \hat{x})\zeta^2), \quad (\text{D.12})$$

by Bayesian updating, P's posterior about $w_t(\bar{x})$ is

$$(w_t(\bar{x}) | w_s(\bar{x}), s, w_{s+1}(\hat{x}), \dots, w_t(\hat{x})) \sim \mathcal{N}(\mu_t^{\text{post}}, 1/p_t^{\text{post}}), \quad (\text{D.13})$$

where the mean and precision are, respectively,

$$\mu_t^{\text{post}} = \frac{p_t^{\text{prior}} \cdot \mu_t^{\text{prior}} + w_t(\hat{x}) / (\bar{x} - \hat{x}) \zeta^2}{p_t^{\text{prior}} + 1 / (\bar{x} - \hat{x}) \zeta^2}, \quad p_t^{\text{post}} = p_t^{\text{prior}} + 1 / (\bar{x} - \hat{x}) \zeta^2. \quad (\text{D.14})$$

It is not straightforward to write out the general formulas for μ_t^{post} and p_t^{post} , but we can analyze two extreme cases. If the real transgression in question is as minor as the test transgression, i.e., if $\bar{x} - \hat{x} \rightarrow 0$, then we have

$$\mu_t^{\text{prior}} \rightarrow \rho w_{t-1}(\hat{x}), \quad p_t^{\text{prior}} \rightarrow 1/\sigma^2, \quad \mu_t^{\text{post}} \rightarrow w_t(\hat{x}), \quad p_t^{\text{post}} \rightarrow \infty. \quad (\text{D.15})$$

The game thus converges to the model in Appendix B, which is for misconduct and crimes not at the apex of government, where the current state of the world is known at the start of each period, and the problem presented in Proposition 1 would not exist. If the real transgression in question is extremely different from the test transgression, instead, i.e., if $\bar{x} - \hat{x} \rightarrow \infty$, then we have

$$\mu_t^{\text{post}} \rightarrow \mu_t^{\text{prior}} \rightarrow \rho^{t-s} w_s(\bar{x}), \quad p_t^{\text{post}} \rightarrow p_t^{\text{prior}} \rightarrow 1 / \left(\sum_{\tau=0}^{t-s-1} \rho^{2\tau} \right) \sigma^2. \quad (\text{D.16})$$

The game thus converges to the baseline model, and the problem presented in Proposition 1 would remain. This demonstrates that the problem presented in Proposition 1 arises distinctively when the real transgression in question is not any minor ones, so that little can be learned from these minor ones about the state of the world relevant to the real transgression in question.

E Alternative Ways of Evolution of the World

Deterministic world. If $\sigma = 0$, the current state of the world follows

$$w_t = \rho w_{t-1} = \rho^{t-s} w_s, \quad \text{where } 0 < \rho \leq 1. \quad (\text{E.1})$$

This implies that the revelation of the state of the world at the last transgression can perfectly indicate the current state, i.e., $w_t > 0$ if and only if $w_s > 0$. Therefore, once a transgression is punished, at each future period, P will understand that any transgression will be punished, so will not transgress. Thus, the problem presented in Proposition 1 would not emerge.

Path-independent world. If $\rho = 0$, the current state of the world at each period is independent of past states of the world, i.e.,

$$w_t = \epsilon_t \sim \mathcal{N}(0, \sigma^2), \text{ i.i.d., with } \sigma > 0. \quad (\text{E.2})$$

The probability of punishment is thus independent of the state of the world at the last transgression, i.e.,

$$\mathbf{P}[w_t > 0 \mid w_s, s] = \mathbf{P}[w_t > 0] = 0.5, \quad (\text{E.3})$$

which is a constant. P will thus always transgress, or always not transgress, in each period, depending on whether or not $a > c$, while independent of the last revelation and how distant it has been in history. Thus, deterrence would not be self-undermining over time.

Oscillatory world. If $\rho < 0$, first observe that

$$(w_{s+1} \mid w_s, s) \sim \mathcal{N}(\rho w_s, \sigma^2), \text{ where } \sigma > 0 \text{ and } \rho < 0; \quad (\text{E.4})$$

$$(w_{s+2} \mid w_s, s) \sim \mathcal{N}(\rho^2 w_s, (1 + \rho^2)\sigma^2), \text{ where } \sigma > 0 \text{ and } \rho < 0. \quad (\text{E.5})$$

Now we consider two cases. First, for any $w_s > 0$, i.e., right after any impeachment, it must follow that $\rho w_s < 0$ and, thus, $\mathbf{P}[w_{s+1} > 0 \mid w_s, s] < 0.5$. By Lemma 2, this implies that, for any transgression gain that is greater than the punishment, i.e., $a > c$, deterrence will fail for sure at $s + 1$.

Second, for any $w_s \leq 0$, instead, i.e., right after any unpunished transgression, it must follow that $\rho^2 w_s \leq 0$ and, thus, $\mathbf{P}[w_{s+2} > 0 \mid w_s, s] \leq 0.5$. By Lemma 2, this implies that, for any transgression gain that is greater than the punishment, i.e., $a > c$, deterrence will fail for sure at $s + 2$, if it succeeded at $s + 1$. Thus, in an oscillatory world, deterrence would never last, and the problem presented in Proposition 1 would have no opportunity to emerge.

Explosive world. If $\rho > 1$, first observe that the conditional distribution is

$$(w_t \mid w_s, s) \sim \mathcal{N}(\rho^{t-s} w_s, (1 - \rho^{2(t-s)}) \cdot \sigma^2 / (1 - \rho^2)), \text{ where } \sigma > 0 \text{ but } \rho > 1. \quad (\text{E.6})$$

The conditional probability of punishment is thus

$$\mathbf{P}[w_t > 0 \mid w_s, s] = \Phi\left(\rho^{t-s} w_s / \left(\sigma \cdot \sqrt{(1 - \rho^{2(t-s)}) / (1 - \rho^2)}\right)\right). \quad (\text{E.7})$$

Observe that, for any $t - s \geq 1$ and $w_s > 0$, $\mathbf{P}[w_t > 0 \mid w_s, s]$ decreases with $t - s$ and

$$\mathbf{P}[w_t > 0 \mid w_s, s] \rightarrow \Phi\left(\sqrt{\rho^2 - 1} \cdot w_s / \sigma\right) \text{ as } t - s \rightarrow \infty. \quad (\text{E.8})$$

Therefore, for any $a > 0$ and $c > 0$, there exists

$$\hat{w} \equiv \left(\sigma / \sqrt{\rho^2 - 1}\right) \cdot \Phi^{-1}(a / (a + c)) \quad (\text{E.9})$$

such that, for any $w_s \geq \hat{w}$,

$$\mathbf{P}[w_t > 0 \mid w_s, s] \geq a / (a + c) \quad (\text{E.10})$$

will always hold for any $t - s \geq 1$, i.e., transgression is forever deterred. The problem presented in Proposition 1 would thus not exist.

F Payoff Uncertainty

Now consider an alternative setup as in Figure F.1, where strategic uncertainty, i.e., whether C will impeach a transgressing P, is replaced by an exogenous probability of punishment, $p \in (0, 1)$, and payoff uncertainty about the gain and punishment, a_t and c_t , with conditional revelation. About the evolution of (a_t, c_t) , we assume that

$$a_t = \rho a_{t-1} + \epsilon_t^a, \quad c_t = \rho c_{t-1} + \epsilon_t^c, \quad \text{where } \rho > 0; \quad (\text{F.1})$$

$$\epsilon_t^a \sim \mathcal{N}(0, \sigma_a^2), \text{ i.i.d., with } \sigma_a \geq 0; \quad \epsilon_t^c \sim \mathcal{N}(0, \sigma_c^2), \text{ i.i.d., with } \sigma_c \geq 0, \quad (\text{F.2})$$

and all ϵ_t^a s and ϵ_t^c s are mutually independent. The time of the last transgression is still denoted as $s < t$. We still assume a myopic P.

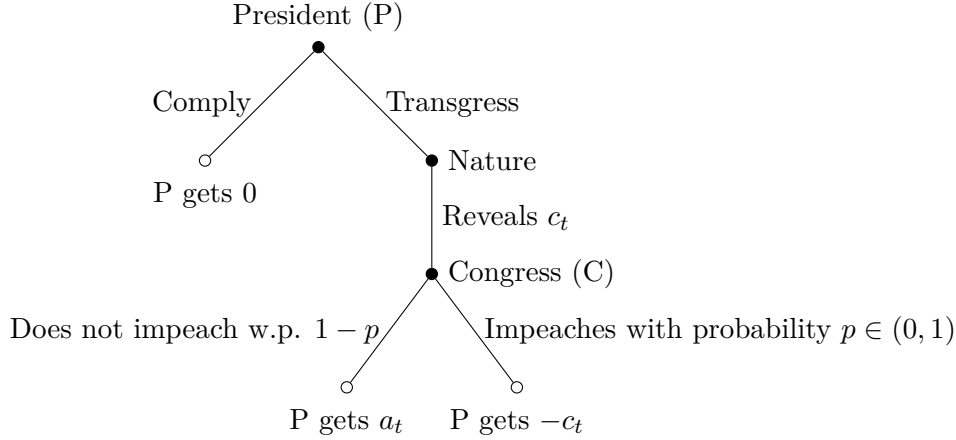
Given this setup, at the beginning of each period t , P will transgress if and only if

$$\mathbf{E}[(1 - p)a_t - pc_t \mid a_s, c_s, s] > 0, \quad \text{i.e.,} \quad (1 - p) \cdot \mathbf{E}[a_t \mid a_s, c_s, s] > p \cdot \mathbf{E}[c_t \mid a_s, c_s, s]. \quad (\text{F.3})$$

Given the evolution of (a_t, c_t) , this condition is equivalent to

$$(1 - p) \cdot \rho^{t-s} a_s > p \cdot \rho^{t-s} c_s, \quad \text{i.e.,} \quad (1 - p) \cdot a_s > p \cdot c_s \quad (\text{F.4})$$

which is independent of σ_a and σ_c , i.e., the payoff uncertainty, and $t - s$, i.e., how long transgression has been deterred. Therefore, once transgression is deterred, it is forever deterred. We see that the mechanism in Lemma 1 and Proposition 1 is about strategic uncertainty, not payoff uncertainty.



P knows time, $s < t$, and state of the world, (a_s, c_s) , at last transgression.

Figure F.1: Payoff uncertainty, period t

G Forward-looking Players

In the main text, we have assumed that the players are myopic. In this section, we show that the problem presented in Proposition 1 is robust with respect to forward-looking players. In particular, we assume that, once a P is impeached, a new P will replace him immediately, and the exiting P will suffer a permanent punishment of $c > 0$ every period afterward; all the players are forward-looking with an infinite horizon, and the discount factor is $\beta \in (0, 1)$.

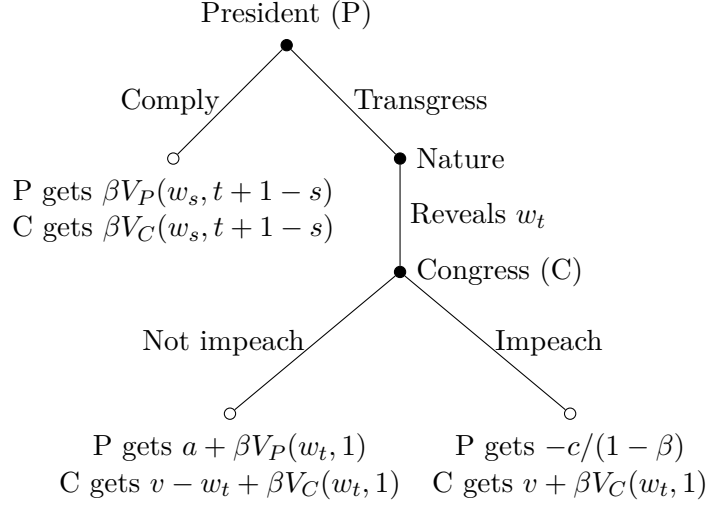
We consider symmetric, pure-strategy Markov-perfect equilibria (SMPE), where “symmetric” means that all Ps are restricted to adopting the same strategy. The payoff-relevant variables, i.e., state variables, are thus only the last revelation of the state of the world, i.e., at the last transgression, w_s , and how long it has been since the last transgression, $t - s$. We thus write the net present values that the players enjoy in equilibrium at the beginning of period t as $V_i(w_s, t - s)$, $i = C, P$.

Figure G.1 then illustrates the structure of the subgame starting from t , taking the continuation value in equilibrium for the future period as given. At the beginning of period t , the payoff-relevant variables, w_s and $t - s$, are known to all the players. As in Section 2, for simplicity, we still assume that C will not impeach P if indifferent between impeaching and not impeaching, and P will not transgress if indifferent between transgressing and complying.

We first establish the existence and uniqueness of the equilibrium. By Figure G.1, we can pin down C’s equilibrium strategy, since she cannot affect the future state of the world:

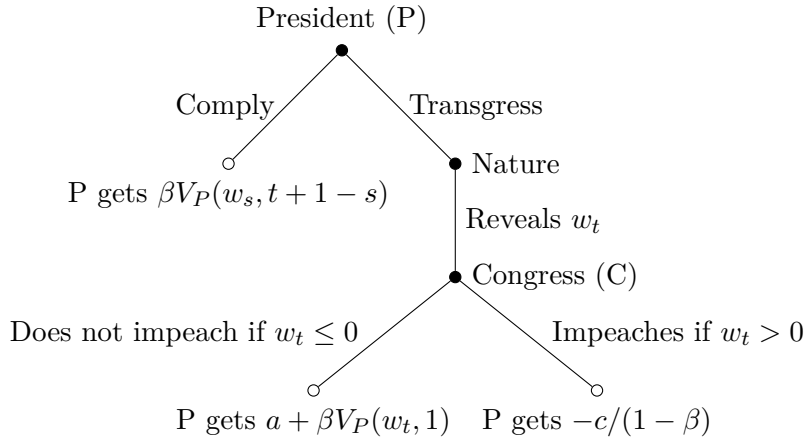
Lemma G.1. *In any SMPE, C’s strategy is to impeach any transgressing P at period t if and only if $w_t > 0$.*

Lemma G.1 helps us reduce the subgame at period t , as in Figure G.1, into a decision



Players know time, $s < t$, and state of the world, w_s , at last transgression.

Figure G.1: Forward-looking players and permanent punishment, subgame at period t



Players know time, $s < t$, and state of the world, w_s , at last transgression.

Figure G.2: Forward-looking players and permanent punishment, subgame at t reduced

problem for P, as in Figure G.2. P's strategy in any SMPE is thus to transgress at period t if and only if

$$\begin{aligned} & \mathbf{E}[a + \beta V_P(w_t, 1) \mid w_t \leq 0, w_s, s] \cdot \mathbf{P}[w_t \leq 0 \mid w_s, s] - (c/(1 - \beta)) \cdot \mathbf{P}[w_t > 0 \mid w_s, s] \\ & > \beta \cdot V_P(w_s, t + 1 - s), \end{aligned} \tag{G.1}$$

or equivalently, if

$$\begin{aligned} & a \cdot \mathbf{P}[w_t \leq 0 \mid w_s, s] - (c/(1 - \beta)) \cdot \mathbf{P}[w_t > 0 \mid w_s, s] \\ & + \beta \cdot \mathbf{E}[V_P(w_t, 1) \mid w_t \leq 0, w_s, s] \cdot \mathbf{P}[w_t \leq 0 \mid w_s, s] > \beta \cdot V_P(w_s, t + 1 - s), \end{aligned} \tag{G.2}$$

The Bellman equation for P at the beginning of period t is thus

$$\begin{aligned} V_P(w_s, t - s) = \max \{ & a \cdot \mathbf{P}[w_t \leq 0 \mid w_s, s] - (c/(1 - \beta)) \cdot \mathbf{P}[w_t > 0 \mid w_s, s] \\ & + \beta \cdot \mathbf{E}[V_P(w_t, 1) \mid w_t \leq 0, w_s, s] \cdot \mathbf{P}[w_t \leq 0 \mid w_s, s], \\ & \beta \cdot V_P(w_s, t + 1 - s) \}. \end{aligned} \quad (\text{G.3})$$

We can then establish the existence and uniqueness of the equilibrium:

Proposition G.1. *There exists a unique SMPE, in which C will impeach any transgressing P at period t if and only if $w_t > 0$, and P's strategy solves the Bellman equation,*

$$\begin{aligned} V_P(w_s, t - s) = \max \{ & a \cdot \mathbf{P}[w_t \leq 0 \mid w_s, s] - (c/(1 - \beta)) \cdot \mathbf{P}[w_t > 0 \mid w_s, s] \\ & + \beta \cdot \mathbf{E}[V_P(w_t, 1) \mid w_t \leq 0, w_s, s] \cdot \mathbf{P}[w_t \leq 0 \mid w_s, s], \\ & \beta \cdot V_P(w_s, t + 1 - s) \}. \end{aligned} \quad (\text{G.4})$$

Proof. First, by Lemma G.1, we have C's strategy in any SMPE, if exists. Second, observe that any $V_P(\cdot, \cdot)$ in equilibrium, if exists, must be bounded by $[0, a/(1 - \beta)]$. This is because, an always-complying strategy will generate a zero net present value for P, while the best possible scenario for P would not be better than an always-transgressing strategy being never punished, which would generate a net present value of $a/(1 - \beta)$ for P.

Now denote the set of functions from $\mathbb{R} \times \mathbb{N}^+$ to \mathbb{R} that are bounded by $[0, a/(1 - \beta)]$ as

$$S \equiv \{ \nu : \mathbb{R} \times \mathbb{N}^+ \rightarrow \mathbb{R}, \nu(\cdot, \cdot) \in [0, a/(1 - \beta)] \}, \quad (\text{G.5})$$

and endow it with the sup-norm metrics,

$$\rho(\nu, \mu) = \sup_{(x, y) \in \mathbb{R} \times \mathbb{N}^+} |\nu(x, y) - \mu(x, y)|. \quad (\text{G.6})$$

Also denote the operator from S to S that is based on the Bellman equation as

$$\begin{aligned} G(\nu)(w_s, t - s) \equiv \max \{ & a \cdot \mathbf{P}[w_t \leq 0 \mid w_s, s] - (c/(1 - \beta)) \cdot \mathbf{P}[w_t > 0 \mid w_s, s] \\ & + \beta \cdot \mathbf{E}[\nu(w_t, 1) \mid w_t \leq 0, w_s, s] \cdot \mathbf{P}[w_t \leq 0 \mid w_s, s], \\ & \beta \cdot \nu(w_s, t + 1 - s) \}. \end{aligned} \quad (\text{G.7})$$

It is easy to show that G satisfies Blackwell's sufficient conditions for a contraction of modulus $\beta \in (0, 1)$. By Banach's contraction mapping theorem, G thus admits a unique fixed point, which is the value function $V_P(\cdot, \cdot)$ in the unique SMPE. The proposition is thus proved. \square

We can then show a result parallel to Proposition 1:

Proposition G.2. *Given any $a > c/(1 - \beta)$, there does not exist a finite level of the state of the world at a transgression at period s , $w_s < \infty$, such that, in the unique SMPE, P will always comply from period $s + 1$ onward.*

Proof. Suppose that there exist such $w_s < \infty$ that, in the unique SMPE, P will comply in any period $t \geq s + 1$. Since P will comply in any period $t \geq s + 1$, then we have, for any $t \geq s + 1$,

$$\begin{aligned} & a \cdot \mathbf{P}[w_t \leq 0 \mid w_s, s] - (c/(1 - \beta)) \cdot \mathbf{P}[w_t > 0 \mid w_s, s] \\ & + \beta \cdot \mathbf{E}[V_P(w_t, 1) \mid w_t \leq 0, w_s, s] \cdot \mathbf{P}[w_t \leq 0 \mid w_s, s] \leq \beta \cdot V_P(w_s, t + 1 - s), \end{aligned} \quad (\text{G.8})$$

where

$$V_P(w_s, t + 1 - s) = 0. \quad (\text{G.9})$$

This condition thus becomes

$$\begin{aligned} & a \cdot \mathbf{P}[w_t \leq 0 \mid w_s, s] - (c/(1 - \beta)) \cdot \mathbf{P}[w_t > 0 \mid w_s, s] \\ & + \beta \cdot \mathbf{E}[V_P(w_t, 1) \mid w_t \leq 0, w_s, s] \cdot \mathbf{P}[w_t \leq 0 \mid w_s, s] \leq 0, \end{aligned} \quad (\text{G.10})$$

or equivalently,

$$\begin{aligned} & a - (a + (c/(1 - \beta))) \cdot \mathbf{P}[w_t > 0 \mid w_s, s] \\ & + \beta \cdot \mathbf{E}[V_P(w_t, 1) \mid w_t \leq 0, w_s, s] \cdot \mathbf{P}[w_t \leq 0 \mid w_s, s] \leq 0. \end{aligned} \quad (\text{G.11})$$

Note that, by Lemma 1, $\mathbf{P}[w_t > 0 \mid w_s, s] \rightarrow 0.5$ as $t - s \rightarrow \infty$; we also know $V_P(\cdot, \cdot) \in [0, a/(1 - \beta)]$. Therefore, by $a > c/(1 - \beta)$, as $t - s \rightarrow \infty$, the left-hand side of Inequality (G.11) is strictly positive. This inequality thus cannot hold as $t - s \rightarrow \infty$, contradicting what we have supposed. The proposition is thus proved by contradiction. \square

H Severe Punishment

Now we turn to the case where the punishment for transgression is sufficiently severe, i.e., $c \geq a$. Lemmas 1 and 2 imply that perpetual deterrence will become possible, in that once transgression is deterred, it will be forever deterred:

Proposition H.1. *For any $a \leq c$, we have $\bar{w} \equiv (\sigma/\rho) \cdot \Phi^{-1}(a/(a + c)) \leq 0$, and*

1. *if $w_s < \bar{w}$, then C will not impeach P at period s , and P will transgress at $s + 1$;*
2. *if $\bar{w} \leq w_s \leq 0$, then C will not impeach P at s , but P will never transgress from $s + 1$ onward;*

3. if $w_s > 0$, then C will impeach P at s , and P will never transgress from $s + 1$ onward.

Proof. About Claim 3, since $w_s > 0$, C will impeach P at period s . By Lemma 1, the conditional probability of punishment at period $s + 1$ is greater than 0.5; by $a \leq c$ and Lemma 2, even a 50/50 coin toss is able to deter transgression. Therefore, P will not transgress at period $s + 1$. Since no transgression has happened since period $s + 1$, period s was still the last time when the state of the world was revealed. The same argument applied to period $s + 1$ thus applies from period $s + 2$ onward, and P will thus not transgress from period $s + 2$ onward. Claim 3 is thus proved.

About Claim 2, by $w_s \leq 0$, C will not impeach P at period s . About period $s + 1$, as in the proof of Proposition 1, P will transgress if and only if

$$w_s < (\sigma/\rho) \cdot \Phi^{-1}(a/(a+c)) \equiv \bar{w}, \quad (\text{H.1})$$

though, by $a \leq c$, $\bar{w} \leq 0$. Therefore, if $\bar{w} \leq w_s \leq 0$, then the threat to punish is sufficiently credible to deter P from transgressing. About period $s + 2$, by the fact that P will not transgress at period $s + 1$, period s is still the last time when the state of the world was revealed. By $w_s \leq 0$ and Lemma 1, that implies that the threat of punishment at period $s + 2$ is even more credible than that at period $s + 1$, so it is sufficiently credible to deter P from transgressing at period $s + 2$. The same argument applies from period $s + 2$ onward. Claim 2 is thus proved.

About Claim 1, since $w_s < \bar{w}$, P will transgress at period $s + 1$; since $\bar{w} \leq 0$, $w_s < 0$ and C will thus not impeach P at period s . Claim 1 and the proposition are thus proved. \square

Proposition H.1 suggests that the self-undermining nature of deterrence of transgressions will not be a problem for perpetual deterrence, if the punishment is sufficiently severe. The relative size of P 's gain and punishment for transgression is thus a key parameter about the institution surrounding the players. A question follows: if someone cannot observe the institution directly, can she find it out from the dynamics of transgression and impeachment?

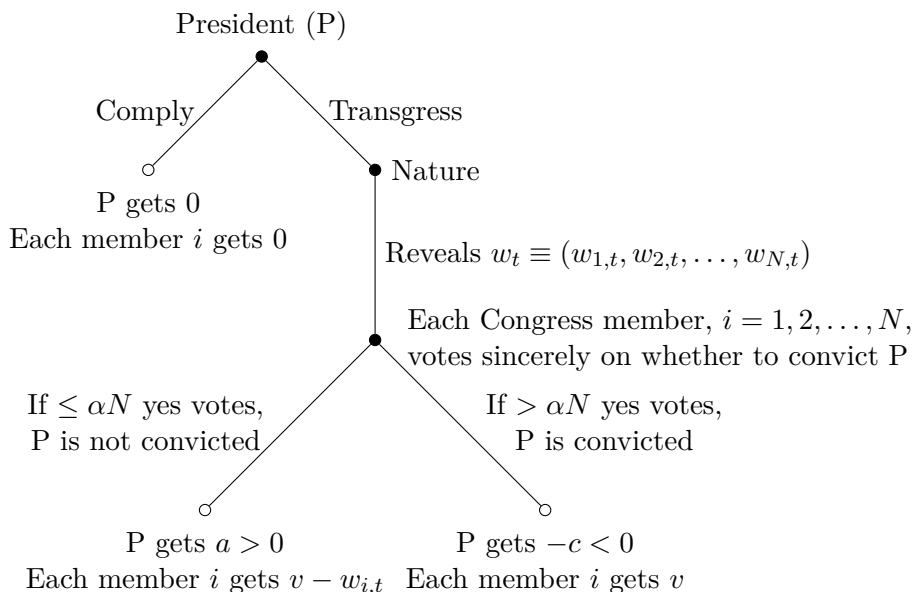
Corollary H.1. *Only observing a transgressing P being impeached at period $s < t$ and constant deterrence from then to the current period t , one cannot identify whether or not the institution will make the constant deterrence forever, i.e., whether $a \leq c$ or $a > c$; she can tell $a > c$ once another transgression happens.*

The intuition behind Corollary H.1 is simple: an ongoing period of constant deterrence can happen in both the cases of $a > c$ and $a \leq c$. In this sense, only observing an ongoing period of constant deterrence will not tell us much about the underlying institution. Therefore, not only does constant deterrence generate a lack of new empirical knowledge

about the state of the world, as contained in the setup of our model, but it cannot tell much about the institution by itself; in order to learn about the institution, one has to wait for the institution to be under stress, i.e., when a transgression happens.

I Congress of Many Non-partisan Members

Now consider an extension as in Figure I.1. The main difference from Section 2 is that now Congress consists of $N \geq 3$ members, each denoted by $i = 1, 2, \dots, N$. With this Congress of multiple members, after P transgresses and Nature reveals the state of the world, each Congress member will vote sincerely, i.e., considering oneself to be pivotal when voting, on whether to convict P. P will then be convicted if and only if there are greater than αN yes votes, where $\alpha \in [0, 1)$ is exogenous and represents the voting rule for conviction.



Players know time, $s < t$, and state of the world, w_s , at last transgression.

Figure I.1: Congress of $N \geq 3$ members, with voting rule $\alpha \in [0, 1)$, period t

All the other features of the extension follow the model in Section 2, with accommodations necessary for a Congress of multiple members. About the state of the world, we assume that each Congress member faces her own net incentive to convict a transgressing P, $w_{i,t}$. The state of the world is thus a vector now, $w_t \equiv (w_{1,t}, w_{2,t}, \dots, w_{N,t})$. We assume that each component of the state of the world evolves as modeled in Section 2, i.e., for any i , $w_{i,t} = \rho_i w_{i,t-1} + \epsilon_{i,t}$, where $0 < \rho_i \leq 1$ and $\epsilon_{i,t} \sim \mathcal{N}(0, \sigma_i^2)$, independent and identically distributed, with $\sigma_i > 0$. We further assume that all $w_{i,t}$ s are mutually independent across all i , so that all the Congress members are non-partisan.

Given this setting, at each period t , by backward induction, each Congress member i will vote to convict a transgressing P if and only if $w_{i,t} > 0$, while P will transgress if and only if the threat of conviction is not credible enough, i.e.,

$$\mathbf{P}[W_t > \alpha N \mid w_s, s] < a/(a + c), \quad (\text{I.1})$$

where W_t is the number of convicting votes against a transgressing P.

Limit over time with constant deterrence. About this threat of conviction, since all $w_{i,t}$ s are mutually independent across all i , Lemma 1 still applies. In particular, as time passes with constant deterrence, the threat from each Congress member i to vote to convict P, conditional on the time and state of the world at the last transgression, s and w_s , will converge to a 50/50 coin toss, i.e.,

$$\mathbf{P}[w_{i,t} > 0 \mid w_s, s] \rightarrow 0.5 \text{ as } t - s \rightarrow \infty. \quad (\text{I.2})$$

Therefore, in this limit, the conditional distribution of the number of votes to convict a transgressing P will converge to a Binomial distribution, with the number of Bernoulli trials being N and the success probability being 0.5, i.e.,

$$(W_t \mid w_s, s) \xrightarrow{d} W \sim \mathcal{B}(N, 0.5) \text{ as } t - s \rightarrow \infty. \quad (\text{I.3})$$

By the Weak Law of Large Numbers, this implies that, in this limit over time with constant deterrence, as the size of Congress, N , increases, it will be almost certain that the proportion of the Congress members voting to convict a transgressing P will be close to one half, i.e.,

$$W/N \xrightarrow{p} 0.5 \text{ as } N \rightarrow \infty. \quad (\text{I.4})$$

Applying the voting rule to this limit, we have the following proposition:

Proposition I.1. *For any transgression gain and punishment for P, given any super-majority rule for conviction in impeachment, in the limit over time with constant deterrence, as the size of Congress increases, it will be almost certain that Congress will not convict any transgressing P and, therefore, P will transgress; given any minority rule, in the limit over time with constant deterrence, as the size of Congress increases, it will be almost certain that Congress will convict any transgressing P and, therefore, P will comply. Mathematically, for*

any $a > 0$ and $c > 0$,

$$\begin{aligned} & \text{if } 0.5 < \alpha < 1, \text{ then } \mathbf{P}[W > \alpha N] \rightarrow 0 < a/(a+c) \text{ as } N \rightarrow \infty; \\ & \text{if } 0 \leq \alpha < 0.5, \text{ then } \mathbf{P}[W > \alpha N] \rightarrow 1 \geq a/(a+c) \text{ as } N \rightarrow \infty. \end{aligned} \quad (\text{I.5})$$

Proposition I.1 suggests that admitting many non-partisan members to Congress, given any super-majority rule for conviction, will make the problem presented in Proposition 1 even worse: it is almost certain that transgression cannot be forever deterred, even if the transgression gain is smaller than the punishment, i.e., $a \leq c$.

J Another Call by an Independent Judiciary

Now consider another extension as in Figure J.1. The main difference from Section 2 is that now, in case that a transgressing P is not impeached by C, a Judiciary (J) can make another call on the matter. If J rules against P, then P will still be punished; if J rules in favor of P, then P will receive the gain from transgression. The net incentive for J to rule against P is w'_t , where we assume that

$$w'_t = \rho' w'_{t-1} + \epsilon'_t, \text{ where } 0 < \rho' \leq 1, \text{ and } \epsilon'_t \sim \mathcal{N}(0, \sigma'^2), \text{ i.i.d., with } \sigma' > 0, \quad (\text{J.1})$$

and is subject to the same conditional revelation as w_t . We also assume that the judiciary is independent, such that w'_t and w_t are mutually independent.

Given this setting, P will transgress if and only if

$$\mathbf{P}[w_t > 0 \text{ or } w'_t > 0 \mid w_s, w'_s, s] < a/(a+c). \quad (\text{J.2})$$

Note that, by the mutual independence between w'_t and w_t , we have

$$\mathbf{P}[w_t > 0 \text{ or } w'_t > 0 \mid w_s, w'_s, s] = 1 - \mathbf{P}[w_t \leq 0 \mid w_s, s] \cdot \mathbf{P}[w'_t \leq 0 \mid w'_s, s]. \quad (\text{J.3})$$

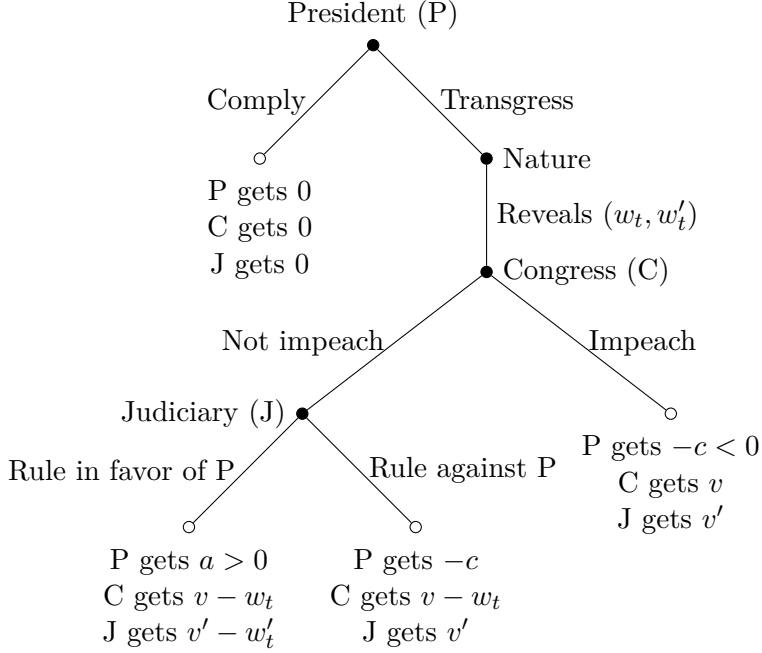
As $t - s \rightarrow \infty$, we have

$$\mathbf{P}[w_t > 0 \text{ or } w'_t > 0 \mid w_s, w'_s, s] \rightarrow 1 - 0.5 \times 0.5 = 0.75. \quad (\text{J.4})$$

Therefore, if

$$a/(a+c) > 0.75, \text{ i.e., } a > 3c, \quad (\text{J.5})$$

then constant deterrence cannot last forever. Therefore, the problem presented in Proposi-



Players know time, $s < t$, and state of the world, (w_s, w'_s) , at last transgression.

Figure J.1: Another call by judiciary, period t

tion 1 will still emerge, even though the required threshold for the transgression gain relative to the punishment is higher, i.e., $a > 3c$, instead of $a > c$.

K Federalism

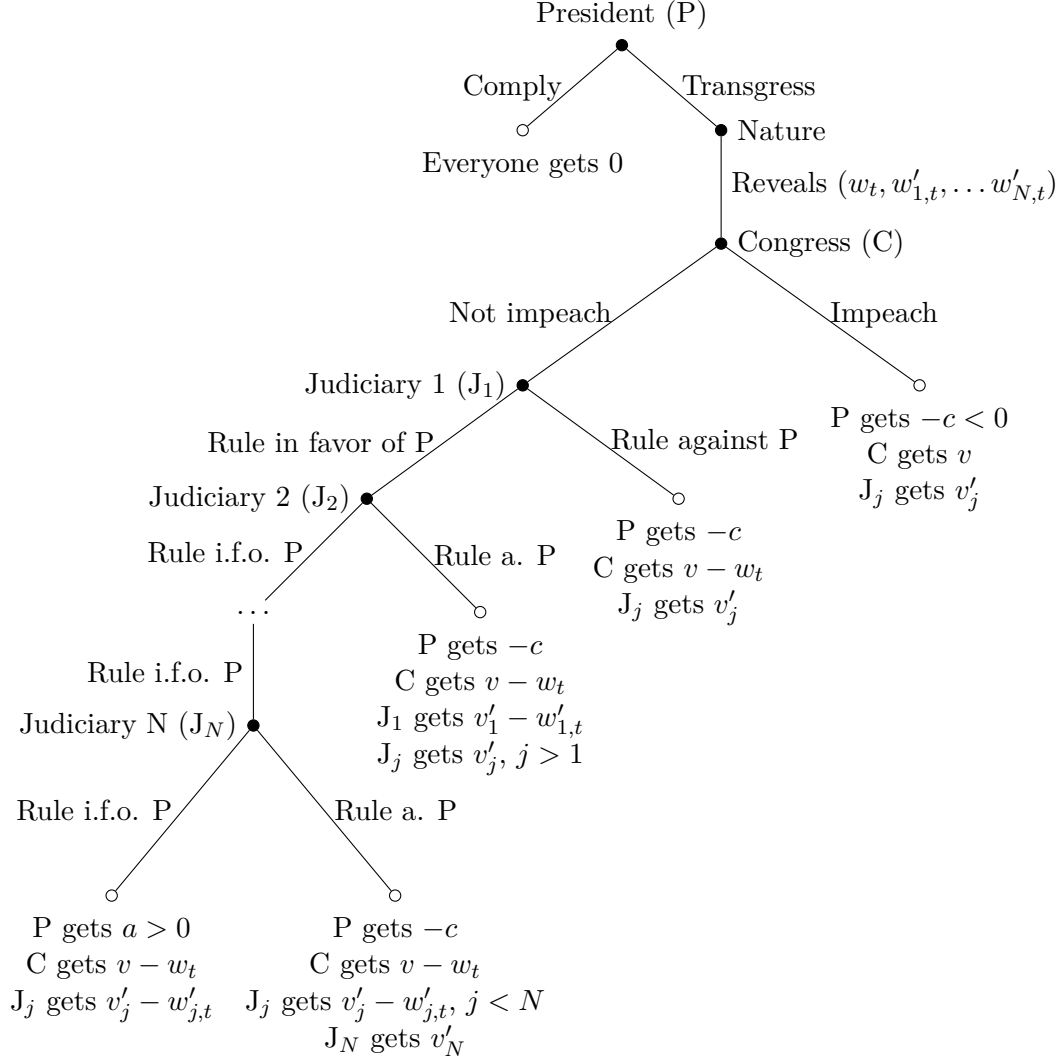
Now consider a further extension as in Figure K.1. The main difference from Section J is that now, each of N Judiciaries, J_i , $i = 1, \dots, N$, in sequence, can make her call if the transgressing P is not punished by C and the judiciaries before her, such that a transgressing P will be punished if C impeaches him, or if at least one J_i rules against him. The net incentive for each J_i to rule against the transgressing P is $w'_{i,t}$ and subject to conditional revelation.

Now suppose that all judiciaries are independent of C and each other, i.e., for each $i = 1, \dots, N$,

$$w'_{i,t} = \rho'_i w'_{i,t-1} + \epsilon'_{i,t}, \text{ where } 0 < \rho'_i \leq 1, \text{ and } \epsilon'_{i,t} \sim \mathcal{N}(0, \sigma_i'^2), \text{ i.i.d., with } \sigma'_i > 0, \quad (\text{K.1})$$

and $w'_{1,t}, \dots, w'_{N,t}$, and w_t are mutually independent. Thus, P will transgress if and only if

$$\mathbf{P} [w_t > 0 \text{ or } w'_{1,t} > 0 \text{ or } \dots \text{ or } w'_{N,t} > 0 \mid w_s, w'_{1,s}, \dots, w'_{N,s}, s] < a/(a + c). \quad (\text{K.2})$$



Players know time, $s < t$, and state of the world, $(w_s, w'_{1,s}, \dots, w'_{N,s})$, at last transgression.

Figure K.1: Federalism, period t

Note that, by the mutual independence among $w'_{1,t}, \dots, w'_{N,t}$, and w_t , we have

$$\begin{aligned}
 & \mathbf{P} [w_t > 0 \text{ or } w'_{1,t} > 0 \text{ or } \dots \text{ or } w'_{N,t} > 0 \mid w_s, w'_{1,s}, \dots, w'_{N,s}, s] \\
 &= 1 - \mathbf{P} [w_t \leq 0 \mid w_s, s] \cdot \prod_{i=1}^N \mathbf{P} [w'_{i,t} \leq 0 \mid w'_{i,s}, s].
 \end{aligned} \tag{K.3}$$

As $t - s \rightarrow \infty$, we have

$$\mathbf{P} [w_t > 0 \text{ or } w'_{1,t} > 0 \text{ or } \dots \text{ or } w'_{N,t} > 0 \mid w_s, w'_{1,s}, \dots, w'_{N,s}, s] \rightarrow 1 - 0.5^{N+1}. \tag{K.4}$$

In this limit, as $N \rightarrow \infty$, we have

$$\mathbf{P} [w_t > 0 \text{ or } w'_{1,t} > 0 \text{ or } \dots \text{ or } w'_{N,t} > 0 \mid w_s, w'_{1,s}, \dots, w'_{N,s}, s] \rightarrow 1. \quad (\text{K.5})$$

Therefore, for any given $a > 0$ and $c > 0$, in the limit over time with constant deterrence, as the number of independent judiciaries increases, it will be almost certain that a transgressing P will be punished and, therefore, P will comply. In this sense, the problem presented in Proposition 1 can be overcome asymptotically.

Note that this result is conditional on the mutual independence of all the judiciaries and C. To see this, suppose instead that all the judiciaries face the same incentive as C's, i.e., for any $i = 1, \dots, N$, $w'_{i,t} = w_t$. This setting is then reduced to the baseline model, and Proposition 1 holds.

L A Positive Drift

Now suppose that the state of the world evolves with a positive drift every period, i.e.,

$$w_t = \mu + \rho w_{t-1} + \epsilon_t, \text{ where } \mu > 0, 0 < \rho \leq 1, \text{ and } \epsilon_t \sim \mathcal{N}(0, \sigma^2), \text{ i.i.d., with } \sigma > 0. \quad (\text{L.1})$$

Its distribution, conditional on its value and the time at the last transgression, is thus

$$(w_t \mid w_s, s) \sim \begin{cases} \mathcal{N} \left(\frac{1-\rho^{t-s}}{1-\rho} \cdot \mu + \rho^{t-s} w_s, \frac{1-\rho^{2(t-s)}}{1-\rho^2} \cdot \sigma^2 \right) & \text{if } \rho \in (0, 1), \\ \mathcal{N}((t-s)\mu + w_s, (t-s) \cdot \sigma^2) & \text{if } \rho = 1. \end{cases} \quad (\text{L.2})$$

The conditional probability of punishment is thus

$$\mathbf{P}[w_t > 0 \mid w_s, s] = \begin{cases} \Phi \left(\left(\frac{1-\rho^{t-s}}{1-\rho} \cdot \mu + \rho^{t-s} w_s \right) / \left(\sigma \cdot \sqrt{\frac{1-\rho^{2(t-s)}}{1-\rho^2}} \right) \right) & \text{if } \rho \in (0, 1), \\ \Phi((\mu(t-s) + w_s) / (\sigma \cdot \sqrt{t-s})) & \text{if } \rho = 1. \end{cases} \quad (\text{L.3})$$

Therefore, as $t - s \rightarrow \infty$, we have

$$\mathbf{P}[w_t > 0 \mid w_s, s] \rightarrow \begin{cases} \Phi \left((\mu/\sigma) \cdot \sqrt{(1+\rho)/(1-\rho)} \right) \in (0.5, 1) & \text{if } \rho \in (0, 1), \\ 1 & \text{if } \rho = 1. \end{cases} \quad (\text{L.4})$$

Therefore, having a constant positive drift can only partially alleviate the problem presented in Proposition 1, if the effect of historical shocks decays over time; it cannot overcome the problem, unless there is no such decay.