

The Dynamics of Policy Complexity

Keiichi Kawai, Ruitian Lang, Hongyi Li*

First Version: May 10, 2014
Current Version: April 25, 2017

Abstract

This paper explores the origins of bureaucratic complexity in public policy. In a model of incremental policymaking where entanglements between policy elements complicate attempts to undo existing policy, policymakers are biased towards increasing policy complexity – especially when policy is already complex. Policy evolution is thus path-dependent: simple policies remain simple, whereas complex policies become more complex. Complexity emerges and persists under political conflict between ideologically-extreme policymakers, and when legislative frictions impede policymaking. Patience is not a virtue: farsighted policymakers deliberately engage in obstructionism, introducing complex policies to hinder future opponents.

JEL Classification: C73, D72, D73

Keywords: kludges, bureaucracy, extremism, obstructionism, organizational change, kludgeocracy

1 Introduction

The complexity of public policy imposes significant costs on society. The United States Internal Revenue Service estimated that the various costs of tax compliance exceeded \$168 billion in 2010, which was fifteen percent of total tax receipts for that year.¹ In many areas of policy ranging from the tax code to education to healthcare, such complexity is pervasive and persistent.

This paper studies the evolution of policy complexity. It develops a theory where policy complexity emerges in the course of political conflict. Successive policymakers modify policy in pursuit of their own policy goals; they do so by layering new rules upon existing policy. As layers of rules accumulate, so does policy complexity.

A key aspect of our theory is that new rules take the form of *kludges*: piecemeal modifications that patch over old programs rather than replace them. Kludges serve to remedy

*Kawai: k.kawai@unsw.edu.au; Lang: ruitian.lang@anu.edu.au; Li: hongyi@unsw.edu.au. We thank Robert Akerlof, Alessandro Bonatti, Steve Callander, Heng Chen, Sven Feldmann, Robert Gibbons, Richard Holden, Anton Kolotilin, Jin Li, Hodaka Morita, Carlos Pimienta, Eric van den Steen, Peter Straka, Birger Wernerfelt, the MIT Organizational Economics Lunch, and the Hitotsubashi Theory Workshop for comments and suggestions; and Adam Solomon for excellent research assistance.

¹This may not be too surprising, given that the U.S. tax code contains more than four million words.

flaws in the implementation of existing policy, or even cancel out their impact. As such, they improve on existing policy, but do so in an inelegant and inefficient fashion relative to the alternative – to completely rewrite existing policy, unburdened by legacy concerns.²

We model a setting where policy control shifts intermittently between two rival policymakers. While in control, each policymaker may add or delete policy rules to achieve his ideological goal. By *complexity* we mean the measure of rules that make up policy (e.g., the number of lines in the tax code). In this setting, kludges are rules that are added to cancel out the ideological impact of old rules without having to delete those old rules. In other words, kludges allow policymakers to avoid elaborate policy overhauls, at the cost of excessive policy complexity.

This narrative has, so far, not yet addressed why policymakers may favor adding new rules as kludges, rather than deleting and replacing existing rules. Our theory incorporates features of the legislative process that are conducive to kludges. First, policymaking is incremental: rules may only be added or deleted gradually. Interest groups may oppose adding new rules that they dislike or deleting old rules that they like. Such resistance constrains policymakers with limited political capital from undertaking radical overhauls; instead, they make incremental changes.³

Further, rules are entangled with one another. Each rule is designed to fit well with the rest of policy – either by legislative intent, or through subsequent administrative and judicial interpretation of enacted legislation. Rules may rely on features of other rules, or fill gaps in other rules, or build upon other rules to modify their effect. Such interdependencies between rules create entanglements that hinder the undoing of existing policy. Deletion of a rule may cripple other dependent rules which rely on the functionality of the deleted rule. Consider the Alternative Minimum Tax (AMT) in the U.S. Tax Code. Many observers deem the AMT to be an unsatisfactory solution to the problems it was intended to solve, but also believe that it will be difficult to undo or significantly edit the AMT because many other aspects of the federal tax system have come to rely on the AMT.

In this paper, such entanglements are modeled as an exogenous constraint on the ability of policymakers to precisely undo existing policy. Essentially, an existing rule cannot be deleted without also deleting other rules that are (randomly) entangled with the targeted rule.

Consequently, each policymaker faces a trade-off between improving his ideological position and reducing policy complexity. He may add new rules that favor his ideological position. Or, he may delete rules that detract from his position. Deletion has the benefit of reducing complexity. However, the deletion process is stymied by entanglement. Unfavorable rules cannot be deleted surgically; they may be entangled with favorable rules which have to be deleted as well, thus slowing progress towards – or even moving policy away from – the policymaker's position. Entanglements thus induce a bias towards adding

²One example of policy kludge is the U.S. Affordable Care Act (ACA) of 2010, which introduced mechanisms (including mandates, subsidies and insurance exchanges) to fill gaps in the existing patchwork of private and public insurance options. A common view of both proponents and opponents was that the ACA is excessively complex compared to alternatives such as a single-payer healthcare system. These alternatives, however, would have required a politically infeasible complete overhaul of healthcare policy.

³For recent discussion, see As Levy and Razin (2013) and Teles (2013). Further, besides political constraints, cognitive limitations may introduce uncertainty about the impact of large-scale policy changes and thus force policymakers to focus on making small 'local' changes to policy (see, e.g., Lindblom 1959, Bendor 1995, Callander 2011a and Callander 2011b).

rather than deleting rules.

In this setting, we analyze the long-run evolution of policy under political conflict. The main dynamic effects are not driven by strategic interactions between policymakers. In fact, for most of this paper, we completely ignore strategic interactions by focusing on myopic policymakers. We demonstrate how, and under what circumstances, complexity may emerge and persist from such myopic dynamics.

We show that when complexity is high, all parties are particularly prone to adding further complexity in the form of kludges. Consequently, policy outcomes exhibit path dependence: simple policies remain simple, whereas complex policies may become more complex.

When complexity is high, policy evolution switches intermittently between two phases. In one phase, the policymaker in control adds rules to shift policy towards his favored position, increasing complexity as he does so. In the other phase, the policymaker in control has attained his policy goals, and deletes rules to reduce complexity. So, long-run outcomes are determined by the tug-of-war between these two phases, which exert opposite forces on complexity. If the first force is more powerful than the second, then complexity may accumulate and become unbounded over the long-run – in which case we say that policy becomes *kludged*.

Having identified this tension, we characterize conditions under which policy may become kludged. One set of comparative static results relates to political institutions:

- Policy is likely to become kludged if parties hold control for relatively equal periods, i.e., political power is relatively balanced.
- Policy is likely to become kludged if power transitions between parties occur frequently; for example, if electoral terms are short.
- Policy is likely to become kludged if legislative *friction* is high; that is, if the legislative process is slowed by procedural hurdles such as veto points.

The same institutional features that generate kludged policies also serve to reduce ideological polarization in policy outcomes. Indeed, our results highlight a tradeoff in the design of political institutions between policy outcomes that are simple but ideologically polarized, and outcomes that are complex but ideologically moderate.

Thus, viewed through the lens of our theory, the American political system is geared towards generating ideologically moderate political outcomes (relative to the preferences of its political parties) – but with the downside of complex, kludged public policy. After all, American electoral competition has historically been relatively balanced and volatile, with control over the presidency and congress switching regularly between the two major parties over the past decades. Further, many hurdles in the American legislative process, such as supermajority voting requirements, a proliferation of veto points, and filibuster rules, hinder the creation of new laws and the undoing of existing laws. Conversely, a planner who prioritizes complexity reduction over ideological moderation should design political institutions that weaken (or even eliminate) political competition and minimize legislative frictions.

Another set of comparative statics relates to political preferences. We show that policy is likely to become kludged as parties' preferences become more polarized – that is,

as the ideological distance between parties' favoured positions increases, and also as parties' preferences over positions become more intense. These results suggest the following connection between two trends in American politics: an inexorable increase in policy complexity over recent decades (Teles 2013) may have been driven by increasing polarization in political preferences over the same period.⁴

Later in the paper, we move to a setting with forward-looking parties, so that strategic interactions come into play. Here, the main lesson is that patience is not a virtue. A forward-looking party may engage in obstructionism: he makes policy changes not to achieve his own policy goals, but rather to hinder his opponent's future policy moves. Specifically, in a conflict between ideologically zealous policymakers, parties exhibit *strategic extremism*. Each policymaker pursues policy positions that are even more ideologically extreme than his preferences would naively dictate. This serves to "shift the goalposts" against his opponent, which ensures that policy remains relatively close to the policymaker's preferred positions in the medium run.⁵ Such strategic extremism has long-run consequences for policy complexity. In particular, we show that relative to the myopic case, strategic extremism may increase the probability that policy becomes kludged.

Literature Review Most closely related is Ely (2011), who studies how inefficiency – in the form of kludged designs – may arise and persist in single-player adaptive processes. Ely considers a genetic code – a set of genes – that performs well if each gene aligns appropriately with the environment, as mediated by the alignment of a particular 'master gene'. A code is kludged if the 'master gene' is poorly-aligned with the external environment. In an evolutionary setting where the genetic code grows increasingly long while being subject to fitness selection over random mutations, Ely focuses on showing that kludge may persist indefinitely: even though mutations may be arbitrarily large and thus a mutation that unkludges the code while maintaining internal alignment will eventually occur for any fixed code length, the increase in code length over time ensures that such mutations grow increasingly rare, and in fact may never occur. Ely (2011) shares with our paper the central conceit of kludges: that interdependencies between elements make kludges difficult to undo, and that increases in complexity lengthen the process of undoing kludges. But our approach differs in various ways; we highlight three clear distinctions here. First, Ely introduces a 'mechanical' evolutionary force that increases complexity over time, whereas in our model, players endogenously choose whether to increase complexity. This difference reflects our model's central focus on understanding the origins of complexity, whereas complexity in Ely (2011) is principally a device to inhibit unkludging.⁶ Second, we consider a two-player game between policymakers with conflicting objectives to highlight the role of political competition in producing kludges; in a one-player version of our model, kludges would never persist. Third, moving beyond the focus in Ely (2011) on a single myopic player, we discuss how strategic motives may lead to kludges.

⁴McCarty, Poole, and Rosenthal (2016) find that American political parties have become increasingly extreme since the 1970s; Azzimonti (2015) finds that political disagreement between parties has intensified in the same period.

⁵Glaeser, Ponzetto, and Shapiro (2005) present a voting model where politicians may *declare* extreme positions (relative to the voting public) to pander to their base. In contrast, in our model, politicians may *implement* extreme policies (relative to their own preferences).

⁶Relatedly, whereas unkludged codes can be arbitrarily complex in Ely (2011), kludge and complexity are essentially synonymous in our model.

Like our paper, Gratton, Guiso, Michelacci, and Morelli (2015) study the dynamics of policy complexity. In their model, policymakers enact legislation purely to bolster their reputation with the public. Reputational incentives to avoid bad legislation are muted if existing policy is already complex, potentially leading to a ‘complexity trap’ that is superficially reminiscent of our model’s path dependence result – albeit via a different mechanism.

A number of other papers from various literatures explore the idea that incremental rule development may be path-dependent. Callander and Hummel (2014) consider a model where successive policymakers with conflicting preferences strategically experiment to find their preferred policy. The first policymaker benefits from a ‘surprising’ experiment outcome, because it deters experimentation by the second policymaker and thus preserves any policy gains by the first policymaker. Ellison and Holden (2013) study a model of endogenous rule development where there are exogenous constraints on the extent to which new rules may ‘overwrite’ old rules. Compared to these models, our paper introduces path dependence through a distinct mechanism – entanglement – and thus produces very different implications.

Our results on strategic extremism are also related to the literature on agenda-setting in politics. Chen and Eraslan (2017) consider a model where competing policymakers take turns to address outstanding policy issues; their key assumption is that an issue that has previously been addressed cannot be revisited by succeeding policymakers. Buisseret and Bernhardt (2015) consider a related setting where this period’s policy outcome is determined by the interaction between an agenda-setter and a decision-maker, and serves as an endogenous status quo for the next period.⁷ A number of interesting strategic effects arise in these settings; for example, in opposition to the strategic extremism effect of the present paper, dynamic concerns in Buisseret and Bernhardt (2015) may restrain the agenda-setter from aggressive policy-setting. These papers do not address policy complexity, which is of course the central focus of the present paper. Further, unlike these other papers, the status-quo effect in the present paper is technological; any changes to policy take time for future policymakers to undo, which drives the dynamics of policy position and complexity.

2 Model

Policy The policy is a set of infinitesimal rules. Each rule’s ideological *direction* is either positive (+) or negative (−). The policy is summarized as a pair of numbers $\mathbf{p} = (p_+, p_-)$, where $p_j \geq 0$ is the mass of rules with direction $j \in \{+, -\}$. The policy’s *position* is the difference between the masses of positive and negative rules,

$$p = p_+ - p_-;$$

and the policy’s *complexity* is its total mass, denoted as

$$\|\mathbf{p}\| = p_+ + p_-.$$

Policy evolves in continuous time, $t \geq 0$. So we write, for example, $\mathbf{p}(t) = (p_+(t), p_-(t))$; but we will often conveniently suppress the time-dependence of policy variables. We take the initial policy $\mathbf{p}(0)$ as given.

⁷Other papers in this literature include Bernheim, Rangel, and Rayo (2006) and Messner and Polborn (2012), and Levy and Razin (2013).

Players and Preferences There are two parties, +1 and -1, generically identified as i . The instantaneous payoff of party $i \in \{+1, -1\}$ is a function of policy position and complexity:

$$u_i(\mathbf{p}) = -z_i |p - p_i^*| - \|\mathbf{p}\| \quad (1)$$

where $p_i^* \in \mathbb{R}$ is his positional *ideal*, $|p - p_i^*|$ is the absolute value of $p - p_i^*$, and $z_i > 1$ is his ideological *zeal*. That is, parties dislike policy positions that are distant from their ideal, and dislike complex policies. We assume that $p_{+1}^* > 0, p_{-1}^* < 0$; and that $z_{+1} > 1, z_{-1} > 1$.⁸

Some descriptive terminology: parties with small (large) $|p_i^*|$ are called *moderates* (*extremists*). Parties with high z_i are *zealous*. A policy with position $p = p_i^*$ is *i-ideal*.

Each party i maximizes his discounted payoff,

$$\max \mathbb{E} \left[\int_0^\infty u_i(\mathbf{p}(t)) e^{-r_i t} dt \right].$$

Most of this paper considers myopic parties: $r_{+1}, r_{-1} \rightarrow \infty$. Two features of myopic behavior are convenient. First, strategic interactions vanish: the myopic party i is unconcerned about what his opponent $-i$ does after i loses control. Second, only the neighbourhood of the current policy $\mathbf{p}(t)$ is relevant, because only nearby policies can be attained in the near future. In particular, as $r_i \rightarrow \infty$, party i 's problem reduces to maximizing the rate of change of his payoff,

$$\max \left\{ \frac{d}{dt} u_i(\mathbf{p}(t)) \right\}. \quad (2)$$

Policymaking Technology At any time t , one party $i(t) \in \{+1, -1\}$ is in control. Control transitions from i to $-i$ are random and arrive at rate $\lambda_i > 0$. Without loss of generality, party +1 starts the game in control: $i(0) = +1$. We interpret λ_i as i 's political *vulnerability*.

At each time t , party $i(t)$ chooses non-negative addition rates $\alpha_+(t), \alpha_-(t)$ and deletion rates $\delta_+(t), \delta_-(t)$ which move (p_+, p_-) :

$$\frac{d}{dt} p_j(t) = \alpha_j(t) - \delta_j(t) \text{ for each } j \in \{+, -\}, \quad (3)$$

subject to a *flow constraint*, reflecting the party's limited capacity to make policy changes, where γ^{-1} parametrizes the degree of legislative *friction*:

$$\begin{aligned} \alpha_+(t) + \alpha_-(t) + \delta_+(t) + \delta_-(t) &\leq \gamma, \\ \delta_j(t) &= 0 \text{ if } p_j(t) = 0 \text{ for } j \in \{+, -\}, \end{aligned} \quad (4)$$

and an *entanglement constraint* on the direction of rule deletion:

$$\frac{\delta_+(t)}{\delta_-(t)} = \frac{p_+(t)}{p_-(t)}. \quad (5)$$

⁸ The assumption $z_i > 1$ ensures that parties have sufficiently intense preferences over ideological position, and thus face a nontrivial tradeoff between adding and deleting rules. Alternatively, one might posit that payoffs are quadratic in position, $u_i(\mathbf{p}) = -z_i (p_i^* - p)^2 - \|\mathbf{p}\|$, so that party i 's positional preferences intensify as policy strays from p_i^* . This alternative formulation is analytically and expositionally less convenient, but produces qualitatively similar results.

In words, the entanglement constraint states that deleted rules must have the same proportions, by direction, as existing rules. This specification of the entanglement constraint is quite tight: given (total) deletion rate

$$\delta(t) = \delta_+(t) + \delta_-(t),$$

each of $\delta_+(t)$ and $\delta_-(t)$ are fully determined from (5). Consequently, we may use $\delta(t)$ to summarize the pair of deletion rates $(\delta_+(t), \delta_-(t))$.

Discussion of the Entanglement Constraint The entanglement constraint (5) captures, in reduced form, the notion of dependencies between rules. The premise is that parties cannot surgically target specific rules for deletion: a party who targets rule π for deletion has to also delete other rules that are entangled with π .

The specific form of (5) is a tractable depiction of severe entanglement, whereby each rule is entangled with many other rules. If policy is severely entangled, then deletions will mostly be “indirect” (i.e., of rules entangled with targeted rules) rather than “direct” (i.e., of targeted rules). Consequently, parties will have little control over the directions of deleted rules, especially if they have limited knowledge about which rules are entangled with each other. The overall composition of deleted rules will match the composition of the policy as a whole, rather than the direction of those rules targeted for deletion. This notion is captured succinctly by our entanglement constraint (5).

In Appendix A, we argue that our formulation of the entanglement constraint is quite natural. We present two alternative approaches to model the notion of policy entanglements, and show that both models generate our entanglement constraint under assumptions that reflect severe entanglement.

Appendix A.1 considers a linear network. Rules are totally ordered along a line. A party who seeks to delete a rule π first has to delete all the rules above π in the ordering. This model produces (5) as a limiting outcome when the number of rules is large.

Appendix A.2 considers a random network where any two rules are connected with some small probability. Dependencies are captured by the network structure: when a policymaker targets a rule π for deletion, he also has to delete all of π 's neighbours. This model produces (5) at the limit where each rule has infinitely many neighbours. Away from this limit, so that entanglement is not severe, a looser version of the entanglement constraint is obtained; we show that our results continue to hold there as well.

Policy Simplicity and Efficiency The following terminology will be helpful. Let the policy's *positive-simplicity* be the ratio of position to complexity, $\frac{p}{\|\mathbf{p}\|}$. (Conversely, *negative-simplicity* is defined as $-\frac{p}{\|\mathbf{p}\|}$.) So, a policy that is very j -simple (j -simplicity close to one) consists mostly of direction- j rules.⁹

Correspondingly, let the policy's *simplicity* be $\frac{|p|}{\|\mathbf{p}\|}$, where $|p| = |p_+ - p_-|$ is the absolute value of position. That is, policy is simple if most rules have the same direction. Restated slightly, policy is simple if complexity $\|\mathbf{p}\|$ is low relative to $|p|$. At the extreme, if all rules have the same direction, then $|p| = \|\mathbf{p}\|$, and we say that policy is *perfectly simple*.

⁹Indeed, j -simplicity $(j \frac{p}{\|\mathbf{p}\|} = \frac{p_j - p_{-j}}{p_j + p_{-j}})$ is just the difference between the proportion of direction- j and the proportion of direction- $(-j)$ rules in the policy.

Notice that any policy \mathbf{p} that is not perfectly simple, so that $|p| < \|\mathbf{p}\|$, is inefficient in the following sense: an alternative policy that achieves the same position p – but has lower complexity – can be constructed by deleting equal masses of positive and negative rules from \mathbf{p} . Indeed, both parties dislike complexity and thus are strictly better off under this alternative policy than under \mathbf{p} .

3 Myopic Dynamics

This section considers myopic parties, who maximize their objective (2) subject to the flow and entanglement constraints: (4) and (5).

3.1 Short-Run Dynamics

We start by characterizing each parties’ optimal addition and deletion choices at each instant, which determine how policy \mathbf{p} evolves in the short run. This sets the stage for Section 3.2 to discuss long-run outcomes.

It is also instructive to rewrite the law of motion (3) in terms of complexity $\|\mathbf{p}\|$ and position p .

$$\frac{d}{dt}\|\mathbf{p}(t)\| = \alpha_+(t) + \alpha_-(t) - \delta(t), \quad (6a)$$

$$\frac{d}{dt}p(t) = \alpha_+(t) - \alpha_-(t) - \delta(t)\frac{p(t)}{\|\mathbf{p}(t)\|}. \quad (6b)$$

Figure 1 illustrates. Complexity $\|\mathbf{p}\|$ increases at unit rate when adding rules in either direction, and decreases at unit rate when deleting rules. Position p increases (decreases) at unit rate when adding positive rules (negative rules). The effect of deletion on position p is more subtle. Under deletion, position shifts at rate $-\frac{p(t)}{\|\mathbf{p}(t)\|}$: equal in magnitude to the policy’s simplicity, and with *opposite* sign. For example, with a relatively positive-simple policy, deletion would shift position “downwards” relatively quickly, as many more positive rules than negative rules are deleted. This has a straightforward geometric interpretation, highlighted in Figure 1: deletion moves \mathbf{p} towards the empty policy $(0, 0)$. Note that $\frac{|p(t)|}{\|\mathbf{p}(t)\|} \leq 1$: the rate at which position shifts under deletion is (weakly) lower than under addition. Only for perfectly simple policies does deletion shift position as rapidly as addition (in the corresponding direction).

For concrete exposition, focus on party $i = +1$.¹⁰ Figure 2 illustrates party +1’s optimal strategy by depicting, as a function of complexity $\|\mathbf{p}\|$ and position p , the direction in which policy evolves.

Start with policies that lie “below” of +1’s ideal ($p < p_{+1}^*$). Here, party +1’s payoff function (1) simplifies to

$$u_{+1}(\mathbf{p}) = z_{+1}(p - p_{+1}^*) - \|\mathbf{p}\|;$$

Player +1’s payoff improves as complexity $\|\mathbf{p}\|$ decreases, and as position p increases towards p_{+1}^* . Combining this observation with the laws of motion (6a) and (6b) yields

$$\frac{d}{dt}u_{+1}(\mathbf{p}) = \alpha_+(z_{+1} - 1) + \alpha_-(-z_{+1} - 1) + \delta\left(-z_{+1}\frac{p}{\|\mathbf{p}\|} + 1\right). \quad (7)$$

¹⁰The focus on party +1 is without loss of generality; by symmetry, party +1’s and party –1’s optimal strategies are identical, up to a reversal of directions + and –.

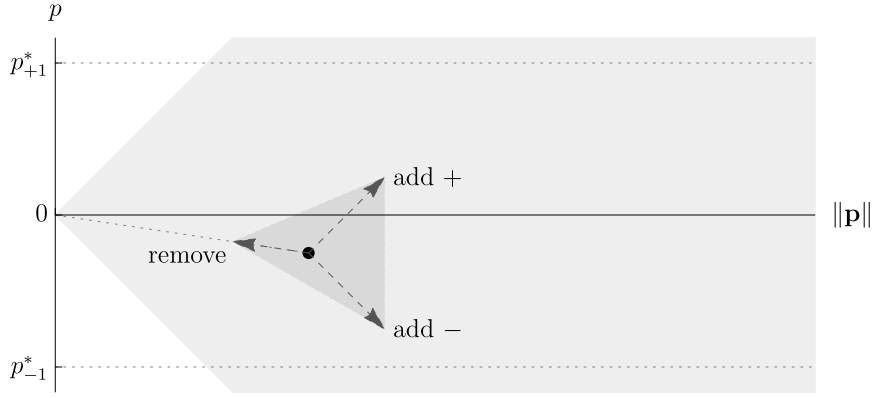


Figure 1: $(\frac{d}{dt} \|\mathbf{p}\|, \frac{d}{dt} p)$ under addition and deletion.

This representation clarifies the pros and cons of adding versus deleting elements. Party +1 has two partially conflicting goals: to increase position p towards his ideal, and to reduce complexity. Clearly, adding negative rules is never optimal for +1: complexity increases and position moves “downward”, away from p^* . The relevant tradeoff for +1 is between adding (positive) rules and deleting rules. Deletion reduces complexity. But, relative to positive addition, deletion slows or even reverses the shift in position towards +1’s ideal, especially if policy is highly positive-simple (so that deleted rules are mostly positive).

Given this tradeoff, party +1 optimally deletes rules iff policy is sufficiently negative-simple, so that deleted rules are mostly negative (and thus are “bad” for +1); specifically, iff $-\frac{p}{\|\mathbf{p}\|} > 1 - \frac{2}{z_{+1}}$. This deletion region, shaded grey in Figures 2a and 2b, shrinks as z_{+1} increases: a zealous party prioritizes positional gains over complexity reduction, and thus favors positive addition over deletion. On the other hand, wherever policy is sufficiently positive-simple (and lies below p^*), party +1 optimally adds positive rules.

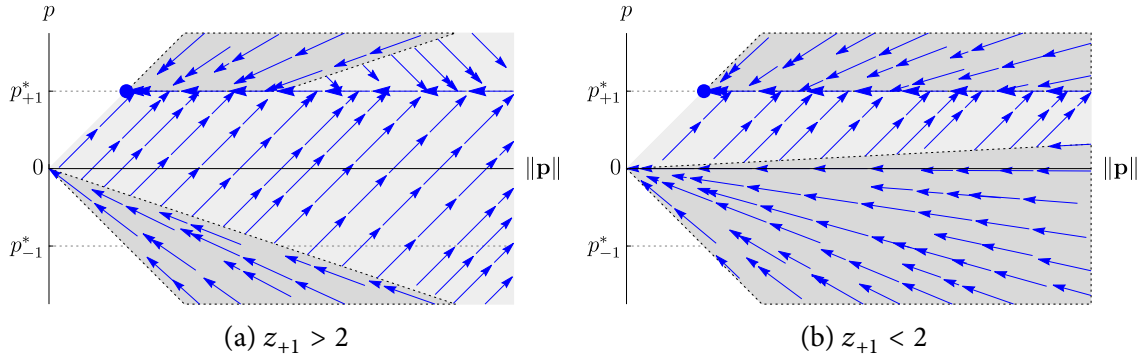


Figure 2: Party +1’s optimal strategy

The case where position lies above +1’s ideal ($p > p^*$) is identical, except that directions are reversed: +1 has to move “downward” to get closer to his ideal. Here, party +1 optimally deletes rules iff \mathbf{p} is sufficiently positive-simple ($\frac{p}{\|\mathbf{p}\|} > 1 - \frac{2}{z_{+1}}$), and optimally adds negative rules otherwise.

To summarize our discussion above, the following proposition specifies party i ’s optimal choice at non-ideal positions ($p \neq p_i^*$). It states that i deletes rules if a sufficiently large proportion of rules are “bad” for him, and adds rules otherwise. Given party i and policy

\mathbf{p} , let $j = \text{sgn}(p - p_i^*)$ be the direction *from* party i 's ideal p_i^* to the policy's position p ; that is, the direction of "bad" rules.

Proposition 1a. *Suppose that party i is in control and that policy \mathbf{p} is not i -ideal ($p \neq p_i^*$).*

1. *If policy is sufficiently j -simple ($j \frac{p}{\|\mathbf{p}\|} > 1 - \frac{2}{z_i}$), then party i deletes rules: $\delta = \gamma$.*
2. *Otherwise, if $j \frac{p}{\|\mathbf{p}\|} < 1 - \frac{2}{z_i}$, then party i adds direction- j rules: $\alpha_j = \gamma$.*

The final case consists of policies at party $+1$'s ideal ($p = p_{+1}^*$). Here, $+1$ has achieved his ideal position, and thus seeks to reduce complexity as *quickly* as possible while shifting position as *slowly* as possible. Thus, he optimally deletes rules if he is not too zealous and if policy is not too simple, so that deletion does not shift position away from his ideal too quickly. Otherwise, he instead chooses an appropriate combination of addition and deletion to maintain position at his ideal while reducing complexity.¹¹

Proposition 1b. *Suppose party i is in control and policy \mathbf{p} is i -ideal, $p = p_i^*$.*

1. *If policy is sufficiently simple ($\frac{|p|}{\|\mathbf{p}\|} > 1 - \frac{2}{z_i}$), then party i reduces complexity while staying on his ideal:*

$$(\alpha_j, \alpha_{-j}, \delta) = \gamma \cdot \left(\frac{|p|}{\|\mathbf{p}\| + |p|}, 0, \frac{\|\mathbf{p}\|}{\|\mathbf{p}\| + |p|} \right), \text{ so that}$$

$$\frac{d}{dt} \|\mathbf{p}(t)\| = -\gamma \frac{\|\mathbf{p}\| - p}{\|\mathbf{p}\| + p} \quad \text{and} \quad \frac{d}{dt} p(t) = 0.$$

2. *Otherwise, if $\frac{|p|}{\|\mathbf{p}\|} < 1 - \frac{2}{z_i}$, then party i deletes rules: $\delta = \gamma$.*

3.2 Path Dependence and Kludge

We now consider long-run policy dynamics. To move beyond the short run, we have to account for how political competition – as captured by the (random) switches of control between parties – affects the evolution of policy.

We will make two points about long-run outcomes. First, policy complexity is path-dependent: the starting point of policy strongly influences the long-run distribution of complexity. Second, there is a tight long-run relationship between complexity and the distribution of policy positions.

The following terminology hints at the outcomes we will be analyzing. We say that policy becomes *kludged* if $\lim_{t \rightarrow \infty} \|\mathbf{p}(t)\| = \infty$. We will focus on two statistics for the long-run distribution of complexity: the probability κ that policy becomes kludged, and a binary indicator K for the possibility of kludge,

$$\kappa = \Pr \left[\lim_{t \rightarrow \infty} \|\mathbf{p}(t)\| = \infty \right] \quad \text{and} \quad K = \begin{cases} 0 & : \quad \kappa = 0 \\ 1 & : \quad \kappa > 0 \end{cases}.$$

As a preliminary, observe that any starting policy eventually becomes *regular*, i.e., positioned at or between the parties' ideals: $p \in [p_{-1}^*, p_{+1}^*]$. This is unsurprising. If policy

¹¹At the perfectly simple i -ideal policy, where $p_{\text{sgn}(i)} = p_i^*$, party i cannot reduce complexity any further without moving away from his ideal, and policy stagnates: $\frac{d}{dt} p = \frac{d}{dt} \|\mathbf{p}\| = 0$.

position lies outside the ideals, then both parties will act to shift policy position in the same direction – towards the positional interval $[p_{-1}^*, p_{+1}^*]$, where both ideals lie. Only for regular policies does positional conflict arise between the parties' preferences, leading to interesting dynamics.

Remark 1.

1. *Suppose that policy $\mathbf{p}(t)$ is regular, i.e., $p(t) \in [p_{-1}^*, p_{+1}^*]$. Then policy (surely) remains regular forever.*
2. *Suppose that policy $\mathbf{p}(t)$ is not regular. Then policy (surely) becomes regular at some random time $\tau > t$, and remains regular thereafter.*

Remark 1 permits us to restrict attention to regular policies. We do so henceforth.

Our first result highlights one aspect of path dependence: *simple policies remain simple*. Let the *basin* \mathcal{B} be the set of regular policies at which at least one party chooses to delete rules (c.f. Propositions 1a and 1b):

$$\mathcal{B} = \left\{ \mathbf{p} : \left(-\frac{p}{\|\mathbf{p}\|} \geq 1 - \frac{2}{z_{+1}} \text{ or } \frac{p}{\|\mathbf{p}\|} \geq 1 - \frac{2}{z_{-1}} \right) \text{ and } p \in [p_{-1}^*, p_{+1}^*] \right\}. \quad (8)$$

Policies that are trapped within the basin can never escape, and tend to grow simpler over time. (See Figure 3.)

Proposition 2. *Suppose that policy lies within the basin: $\mathbf{p}(t) \in \mathcal{B}$.*

1. *Policy (surely) remains within the basin forever: $\mathbf{p}(t') \in \mathcal{B}$ for all $t' \geq t$.*
2. *Policy (almost surely) becomes perfectly simple at some random time $\tau \geq t$.*
3. *A perfectly simple policy (surely) remains perfectly simple forever.*

If both parties are sufficiently zealous ($z_{+1} > 2$ and $z_{-1} > 2$), then the basin \mathcal{B} consists of relatively simple policies. In this case, the intuition for Proposition 2 can be cleanly stated: *sufficiently simple policies grow (weakly) monotonically simpler*. This is because neither party benefits from reducing simplicity by ‘contaminating’ a sufficiently simple policy with new rules of the minority type. (See Figure 3a.)

To fix ideas, consider a (regular) policy that is highly positive-simple. From Party +1’s perspective, the policy’s position is either at or below his ideal ($p \leq p_{+1}^*$). He adds positive rules if $p < p_{+1}^*$. He reduces complexity, while maintaining position, if $p = p_{+1}^*$. In either case, simplicity increases. From Party –1’s perspective, the policy’s position is above his ideal ($p > p_{-1}^*$). He faces a tradeoff between adding negative rules and deleting rules. But, because policy is mostly positive-simple, deletion (of mostly positive rules) is optimal for –1. This leaves policy simplicity unchanged.

Given these incentives, simple policies grow progressively simpler as control changes hands between the two parties. In fact, any policy within \mathcal{B} eventually becomes perfectly simple, and remains so. In other words, \mathcal{B} serves as a basin of attraction for the set of perfectly simple policies.¹²

¹²As we will see shortly, this statement is somewhat imprecise. Depending on parameter values, the basin of attraction for the set of perfectly simple policies is either \mathcal{B} or the larger set of all regular policies.

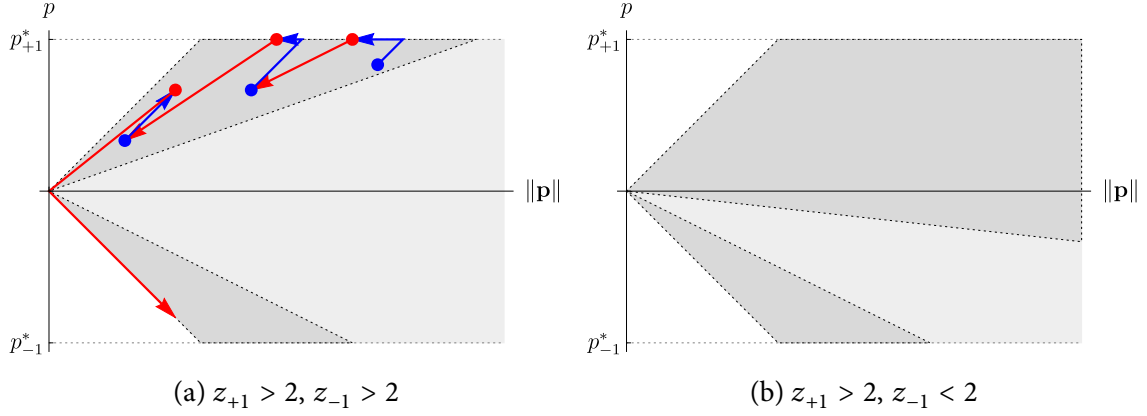


Figure 3: Basin \mathcal{B} (shaded grey region)

As parties become less zealous, the basin expands to include less-simple policies. If either party is insufficiently zealous ($z_{+1} \leq 2$ or $z_{-1} \leq 2$), the basin even contains all policies below or above the $\|\mathbf{p}\|$ -axis ($p \leq 0$ or $p \geq 0$). (See Figure 3b.) In this case, the basin becomes infinite in extent. Consequently, any starting policy inevitably becomes captured within \mathcal{B} . Long-run dynamics are mundane in this case.

Remark 2. *Suppose that $z_{+1} \leq 2$ or $z_{-1} \leq 2$. Then any policy $\mathbf{p}(t)$ (almost surely) becomes perfectly simple at some random time $\tau \geq t$, and remains perfectly simple thereafter.*

Hereafter, our analysis will focus on the case where both parties are sufficiently zealous ($z_{+1} > 2$ and $z_{-1} > 2$).

Outside the basin \mathcal{B} , how does policy complexity evolve? In particular, does policy always move into the basin and remain perfectly simple forever? Or, conversely, does complexity increase unboundedly?

Propositions 1a and 1b tell us that outside \mathcal{B} , each party i adds rules towards his ideal, and focuses on reducing complexity when at his ideal. This leads, in equilibrium, to the following laws of motion for position and complexity. For regular $\mathbf{p}(t) \notin \mathcal{B}$, position p moves towards (and stops at) p_i^* while party i is in control:

$$\frac{d}{dt} p(t) = \gamma \cdot \begin{cases} 1 : & p(t) \in [p_{-1}^*, p_{+1}^*) \text{ and } i(t) = +1 \\ -1 : & p(t) \in (p_{-1}^*, p_{+1}^*] \text{ and } i(t) = -1 \\ 0 : & p(t) = p_{i(t)}^* \end{cases} \quad (9a)$$

Whereas, complexity $\|\mathbf{p}\|$ decreases when position is *at* either ideal, and increases when position is *between* ideals;¹³ see Figure 4.

$$\frac{d}{dt} \|\mathbf{p}(t)\| = \gamma \cdot \begin{cases} 1 & : p(t) \in (p_{-1}^*, p_{+1}^*) \text{ or } p(t) = p_{-i(t)}^* \\ -\frac{\|\mathbf{p}(t)\| - |p(t)|}{\|\mathbf{p}(t)\| + |p(t)|} & : p(t) = p_{i(t)}^* \end{cases} \quad (9b)$$

¹³ The exception to this rule is the case where party i is in control and policy is at $-i$'s ideal: $p = p_{-i}^*$, in which case complexity increases ($\frac{d}{dt} \|\mathbf{p}\| = 1$). But this case may essentially be ignored, because policy spends zero time in this region of the state space: if party i takes control at $-i$'s ideal p_{-i}^* outside \mathcal{B} , he adds j -rules and instantaneously moves policy away from p_{-i}^* .

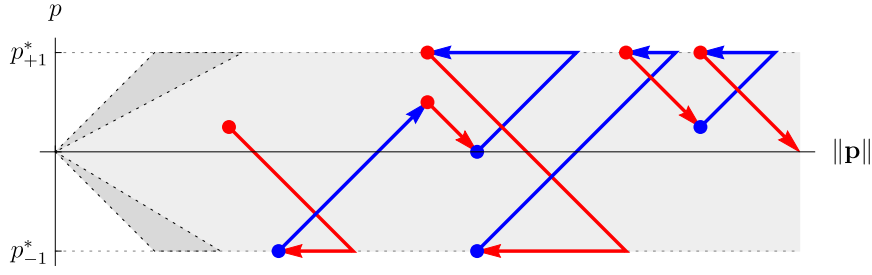


Figure 4: Increasing complexity outside the basin.

This short-run relationship between complexity and position, as expressed by (9b), extends naturally into the long-run. To state the long-run result precisely, first note that outside the basin \mathcal{B} , the long-run behavior of position p can be described in terms of its steady-state distribution.

Lemma 1. *Let $q(t) \in [p_{-1}^*, p_{+1}^*]$ be the random process that obeys, for all $t \geq 0$, the law of motion specified by (9a). Then the Markov process $(q(t), i(t))$ is uniquely ergodic, i.e., has a unique invariant (steady-state) distribution.*

Let $F(\cdot)$ be the steady-state marginal distribution of $q(t)$ from Lemma 1. Define

$$\mu = \int_{[p_{-1}^*, p_{+1}^*]} v(q) dF(q) \quad \text{where} \quad v(q) \equiv \begin{cases} 1 & : p_{-1}^* < q < p_{+1}^* \\ -1 & : q = p_{+1}^* \text{ or } p_{-1}^* \end{cases}.$$

We may interpret μ as the long-run average drift of complexity $\|\mathbf{p}\|$ outside the basin \mathcal{B} (given the normalization $\gamma = 1$). Alternatively, and equivalently, $-\mu$ represents the long-run average frequency of ideal positions: it captures how much time p spends *at* ideals instead of *between* ideals. Our next result builds on this equivalence and points out that kludge is possible if and only if ideal positions are achieved infrequently, so that complexity drifts upward in the long-run.¹⁴

Proposition 3. *Suppose that both parties are sufficiently zealous ($z_{+1} > 2$ and $z_{-1} > 2$), and that the starting policy $\mathbf{p}(0)$ is regular and not in the basin \mathcal{B} .*

1. *If $\mu > 0$, then $K = 1$ and policy (almost surely) becomes kludged or perfectly simple.*
2. *If $\mu < 0$, then $K = 0$ and policy (almost surely) becomes perfectly simple.*

Proposition 3 is a limited result in some respects. It does not fully characterize the probability κ that kludge occurs, and solves instead for the less-informative binary statistic K . This limitation arises because the policy process $(\mathbf{p}(t), i(t))$ is non-ergodic (i.e., path-dependent), so that long-run distributional outcomes such as κ are difficult to directly characterize.

¹⁴One might wonder whether our interpretation of μ as the long-run average drift of complexity is inaccurate, given that the rate at which complexity $\|\mathbf{p}\|$ decreases at ideals, $\left| \frac{d}{dt} \|\mathbf{p}(t)\| \right| = \frac{\|\mathbf{p}\| - |p|}{\|\mathbf{p}\| + |p|}$, is smaller than the rate at which complexity increases between ideals, $\left| \frac{d}{dt} \|\mathbf{p}(t)\| \right| = 1$. However, this difference vanishes at the high-complexity limit: $\frac{\|\mathbf{p}\| - |p|}{\|\mathbf{p}\| + |p|} \rightarrow 1$ as $\|\mathbf{p}\| \rightarrow \infty$. Indeed, the long-run statistics that we calculate are determined by the dynamics of policy at this high-complexity limit. Our interpretation of μ as the drift of complexity reflects this insight.

In other respects, Proposition 3 is quite powerful. It links K to the long-run properties of the (modified) position process $(q(t), i(t))$, which is ergodic and which permits a closed-form solution for the long-run distribution – derived in Lemma B.1b in the Appendix. This enables a rich set of sharp comparative statics results about how K changes with model primitives, which we explore in Section 3.3.

Let's return to our discussion of path dependence. Proposition 2 showed that simple policies grow simpler over time. If $\mu > 0$, then the (loose) converse holds as well: complex policies tend to grow more complex. This is illustrated crudely by Proposition 3, which shows that any policy outside the basin \mathcal{B} of relatively simple policies may become kludged. The following result is a cleaner formulation of the same basic points. If the initial policy has low (high) complexity, then the probability that policy eventually becomes kludged is low (high).

Proposition 4. *Fix a starting policy that is neutral and not perfectly simple: $p(0) = 0$ and $\|\mathbf{p}(0)\| > 0$. Suppose that $z_{+1} > 2$ and $z_{-1} > 2$, and that $\mu > 0$, so that $K = 1$. Then*

$$\begin{aligned} \kappa &\rightarrow 0 & \text{as } \|\mathbf{p}(0)\| &\rightarrow 0, \\ \kappa &\rightarrow 1 & \text{as } \|\mathbf{p}(0)\| &\rightarrow \infty. \end{aligned}$$

Such path dependence has policy implications. Because complexity begets further complexity, one-time interventions to reduce complexity may produce long-run gains that are underestimated in static analyses. As a concrete example, a simplification of the tax code obviously reduces the costs of tax compliance, and has an additional potential benefit: it may potentially prevent the tax code from growing ever more complex, or at least may slow the growth of said complexity.

3.3 Comparative Statics: The Politics of Kludges

For convenient exposition, relabel some of the model's primitives as follows. Define

$$\begin{aligned} \text{volatility: } \lambda &= \sqrt{\lambda_{+1} \lambda_{-1}}, \\ \text{imbalance: } \Lambda &= \max \left\{ \frac{\lambda_{+1}}{\lambda_{-1}}, \frac{\lambda_{-1}}{\lambda_{+1}} \right\}, \\ \text{distance: } \Delta_p^* &= p_{+1}^* - p_{-1}^*, \end{aligned}$$

where λ , the (geometric) average of control change arrival rates, captures the volatility of political control; where Δ_p^* , the difference between parties' ideals, captures the ideological distance between parties; and where Λ captures the degree of power imbalance between parties. Further, label the following increasing function of Λ as

$$\tilde{\Lambda} = \begin{cases} 1 & \text{if } \Lambda = 1 \\ \frac{\log \frac{3-\Lambda^{-1}}{3-\Lambda}}{\sqrt{\Lambda} - \sqrt{\Lambda^{-1}}} & \text{if } 1 < \Lambda < 3. \\ +\infty & \text{if } \Lambda \geq 3 \end{cases}$$

Propositions 5a and 5b encapsulate our comparative statics for complexity. We first state the results, then discuss the intuition.

Proposition 5a. Suppose that both parties are sufficiently zealous ($z_{+1} > 2$ and $z_{-1} > 2$), and that the starting policy $\mathbf{p}(0)$ is regular and not in the basin \mathcal{B} .

1. If $\Delta_p^* \lambda > \gamma \tilde{\Lambda}$, then $K = 1$ and policy (almost surely) becomes kludged or perfectly simple.
2. If $\Delta_p^* \lambda < \gamma \tilde{\Lambda}$, then $K = 0$ and policy (almost surely) becomes perfectly simple.

Proposition 5b. Suppose that the conditions of Proposition 5a.1 are satisfied, so that $K = 1$. Then κ is increasing in z_{+1} and z_{-1} .

Proposition 5a follows closely from Proposition 3. It specifies conditions under which the long-run frequency of ideal positions is low (or high) enough that (outside the basin \mathcal{B}) complexity p drifts upward (or downward) in the long-run, $\mu > 0$ (or $\mu < 0$). This equivalence helps us to understand the comparative statics specified in Proposition 5a.

As imbalance Λ in political power increases, ideal positions are attained more frequently, so the complexity statistic K decreases. The intuition is clearest when $\frac{\lambda_{+1}}{\lambda_{-1}}$ is small and Λ is large, so that party +1 has much more influence over policy than party -1 on average. In that case, policy spends more time at +1's ideal than anywhere else. See Figure 5.

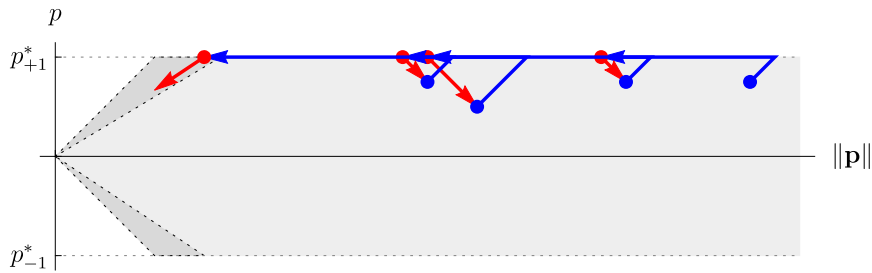


Figure 5: Large power differential leads to less kludge.

As legislative friction γ^{-1} increases, ideal positions are attained less frequently, and K increases. The intuition is clearest when γ^{-1} is low. In this case, each party i can rapidly change position while in power; he quickly reaches his ideal and spends most of his time there, reducing complexity. Conversely, when friction is high, each party spends more time moving (slowly) towards his ideal.

As political volatility λ increases, ideal positions are attained less frequently, and K increases. With high volatility, each party spends little time in control. He is thus unlikely to reach his ideal, and spends little time on average there, before he loses control.¹⁵

As ideological distance Δ_p^* increases, ideal positions are attained less frequently, and K increases. The intuition is clearest if Δ_p^* is small. In that case: whenever either party i takes control, for any (regular) starting position, he quickly reaches his ideal and stays there for most of his time in control. Conversely, if Δ_p^* is large, then each party starts off further from his ideal on average, and thus spends less time at his ideal.

Proposition 5b states that as parties becomes more zealous, kludge becomes less likely. The logic of this result differs from those of Proposition 5a: a change in zealousness has

¹⁵In fact, an increase in volatility Λ is equivalent to an increase in friction γ^{-1} (and a corresponding time dilation); both changes result in less policy change for each control interval.

no effect on the frequency of ideal positions and thus no effect on the drift of complexity. Instead, an increase in zealotry shrinks the basin \mathcal{B} . Consequently, policy becomes less likely to enter the basin and get trapped; rather, policy is more likely to ‘escape’ and become kludged.

A second set of comparative statics relate to policy polarization; that is, the long-run extent to which policy deviates from a ‘neutral’ position. To measure policy polarization, let’s adopt the normalization $p_{+1}^* = -p_{-1}^*$, so that the midpoint $p = 0$ between ideals is neutral. Let H be the steady-state marginal distribution of $|p(t)| \in [0, p_{+1}^*]$ under the law of motion (9a). In fact, H is a natural long-run measure of policy polarization: recall that policy eventually becomes either kludged or perfectly simple, and that $p(t)$ obeys the law of motion (9a) in either case. Accordingly, we say that polarization *increases* if $H(\cdot)$ increases in the sense of first-order stochastic dominance.

Proposition 6. *Suppose that $p_{+1}^* = -p_{-1}^*$. Then, policy polarization is:*

- (i) *increasing in imbalance Λ ,*
- (ii) *constant in zealotry z_{+1} and z_{-1} .*

Further fixing $\Lambda = 1$ (political power is balanced), policy polarization is:

- (iii) *decreasing in friction γ^{-1} ,*
- (iv) *decreasing in volatility λ ,*
- (v) *increasing in ideological distance $\Delta_p^* = 2p_{+1}^*$.*

Our comparative statics for complexity and our comparative statics for polarization are closely related. To start, consider those primitives corresponding to technological aspects of political and legislative processes $(\Lambda, \gamma^{-1}, \lambda)$. As recorded in Table 1, any change to one of these primitives that decreases polarization also increases complexity.

Because our measure of deviation $|p|$ is maximized at either ideal, an increase in polarization corresponds (informally speaking) to an increase in $-\mu$, the fraction of time spent at ideal positions relative to between-ideal positions – and thus in a decrease in long-run complexity, as measured by K . That is, holding the ideological distance Δ_p^* fixed, *polarization reduces complexity*.

The comparative static effects of changes to political preferences $(\Delta_p^*, z_{+1}, z_{-1})$ differ from those of changes to political processes. (i) A decrease in ideological distance Δ_p^* reduces both polarization and complexity. This is because Δ_p^* has opposite effects on $-\mu$ and polarization. A decrease in Δ_p^* , by reducing the time taken to travel between ideals, ensures that policy spends more time at than between ideals, thus increasing $-\mu$ and (consequently) decreasing complexity. But a decrease in Δ_p^* also forces policy to move within a narrower range of positions, and thus mechanically decreases polarization. (ii) A decrease in zealotry z_{+1}, z_{-1} , while decreasing complexity (Proposition 5b), has no effect on polarization. After all, changes in zealotry preserve the law of motion (9a) and thus also the long-run distribution of position p .

Policy Implications Our comparative statics provide some prescriptions for the design of political institutions. A patient planner who seeks to reduce long-run complexity should remove legislative impediments to rulemaking such as supermajority rules and vetoes (i.e.,

| | parameter | symbol | complexity | polarization |
|--------------|-------------|------------------|------------|---------------|
| institutions | volatility | λ | \nearrow | \searrow |
| | friction | γ^{-1} | \nearrow | \searrow |
| | imbalance | Λ | \searrow | \nearrow |
| preferences | distance | Δ_p^* | \nearrow | \nearrow |
| | zealousness | z_{+1}, z_{-1} | \nearrow | \rightarrow |

Table 1: Comparative Statics for Long-Run Outcomes

increase γ). She should reduce political volatility, perhaps by increasing the length of election cycles (i.e., reduce λ). Perhaps controversially, rather than balancing power between political parties, she should instead design political institutions that favour one party over others (i.e., increase Λ); stated crudely, she should support autocracies over democracies. However, our model suggests that such changes to the political process may not be costless, even in the long run: decreasing complexity in this fashion may come with an increase in polarization.

On the other hand, a planner can avoid the complexity-polarization tradeoff by manipulating political preferences. We prefer to think of such preference manipulation in terms of cultural change: by fostering a moderate political culture and curbing the extremist tendencies of political parties (reducing Δ_p^*), a polity may reduce both complexity and polarization in policy.

4 Strategic Extremism

This section considers strategic behavior by non-myopic parties. We will show that ideologically zealous (high z_i) parties may engage in *strategic extremism*: i.e., move towards extreme positions that lie beyond their ideals. Such strategic extremism may increase policy complexity.

We consider the limit where both parties are infinitely zealous, but nonetheless care infinitesimally about complexity. This simplification renders the problem particularly tractable by reducing the associated two-dimensional optimal control problem (over $\|\mathbf{p}\|$ and p) to a one-dimensional problem (over p). Importantly, as we will argue, this simplification preserves the essential dynamic forces in our model, and thus allows us to cleanly highlight the impact of strategic interactions on long-run outcomes.

Start by introducing purely positional preferences: $u_i(\mathbf{p}(t)) = -|p_i^* - p(t)|$, so each party i maximizes

$$\mathbb{E} \left[- \int_0^\infty |p_i^* - p(t)| e^{-r_i t} dt \right].$$

We restrict attention to strategies where each party i has a favoured position p_i^{**} called his *target* and – subject to the flow constraint (4) and entanglement constraint (5) – moves

towards it as quickly as possible:

$$\frac{d}{dt} p(t) = \gamma \cdot \begin{cases} 1 & : p(t) < p_i^{**} \\ -1 & : p(t) > p_i^{**} \\ 0 & : p(t) = p_i^{**} \end{cases} . \quad (10)$$

Note that (10) does not specify how complexity evolves, and thus does not uniquely define a strategy. For example, if $p_i^{**} > 0$ and $\mathbf{p}(t)$ contains only negative rules, then (10) may be satisfied either by adding positive rules or by deleting (negative) rules.

Let's mechanically introduce an infinitesimal distaste for complexity. Suppose that each party minimizes $\frac{d}{dt} \|\mathbf{p}\| = \alpha_+ + \alpha_- - \delta$ given the constraint (10). With this additional assumption, the strategy is uniquely defined. The strategy adds complexity everywhere – except at perfectly simple policies with opposite sign to p_i^{**} and at p_i^{**} itself, where the strategy reduces complexity as quickly as possible. That is,

$$\frac{d}{dt} \|\mathbf{p}(t)\| = \gamma \cdot \begin{cases} 1 & : p \neq p_i^{**} \text{ and } |p| < \|\mathbf{p}\| \\ \text{sgn}(p) \text{sgn}(p_i^{**} - p) & : p \neq p_i^{**} \text{ and } |p| = \|\mathbf{p}\| \\ -\frac{\|\mathbf{p}\| - p}{\|\mathbf{p}\| + p} & : p = p_i^{**} \text{ and } |p| \leq \|\mathbf{p}\| \\ 0 & : p = p_i^{**} \text{ and } |p| = \|\mathbf{p}\| \end{cases} . \quad (11)$$

We say that a strategy is *focused* on a target p_i^{**} if it obeys constraints (10) and (11). Figure 6 depicts focused strategies for party +1. We show in the Appendix that a (subgame perfect) equilibrium in focused strategies always exists (Lemma B.6f). Hereafter, we restrict attention to equilibria in focused strategies, and refer to them simply as *equilibria*.

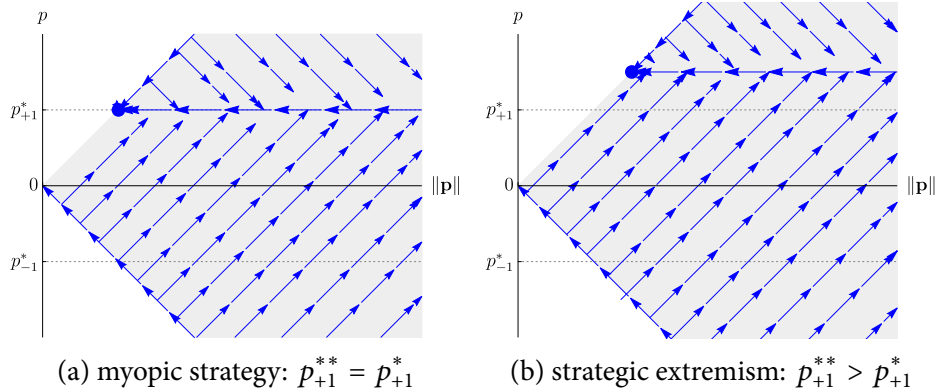


Figure 6: Focused Strategies for Party +1

Focused strategies are very similar to optimal myopic strategies, especially for highly zealous parties. Indeed, inspecting Figure 2, we see that a myopic strategy is identical to a focused strategy which targets the party's ideal ($p_i^{**} = p_i^*$) everywhere except the slivers of policies containing mostly “bad” rules ($j \frac{p}{\|\mathbf{p}\|} < 1 - \frac{2}{z_i}$ where $j = \text{sgn}(p - p_i^*)$). These slivers narrow as zealousness z_i increases; at the limit of infinite zealousness, $z_i \rightarrow \infty$, the difference between the optimal myopic strategy and the focused strategy vanishes. There, focused strategies generalize the myopic strategy by allowing the party to target a non-ideal position.

Consequently, with focused strategies, the dynamics of complexity from the myopic setting of Section 3 are essentially preserved, with the twist that parties' targets act as "endogenous ideals". Complexity increases when policy lies *between* targets, and decreases when policy is *at* either target. Thus a version of Proposition 5a holds in this setting, with targets taking the place of ideals: the average drift μ of complexity – and thus the long-run complexity statistic K – increases with the distance between players' targets, denoted as

$$\Delta_p^{**} = |p_{+1}^{**} - p_{-1}^{**}|.$$

Proposition 7. *Suppose that both parties play focused strategies, and that the starting policy is not perfectly simple ($|p(0)| < \|p(0)\|$).*

1. *If $\Delta_p^{**} \lambda > \gamma \tilde{\Lambda}$, then $K = 1$.*
2. *If $\Delta_p^{**} \lambda < \gamma \tilde{\Lambda}$, then $K = 0$.*

Given that target locations affect long-run complexity, we seek to understand how strategic considerations affect the parties' equilibrium target choices. Myopic play serves as a benchmark: strategic considerations are absent there, and parties choose their ideals as targets, i.e., $\Delta_p^{**} = \Delta_p^*$. (See Figure 6a.)

We say that party i engages in *strategic extremism* if he chooses a target that is more extreme than his ideal: that is, if his ideal p_i^* lies between his target and his opponent's target.

Proposition 8. *In any equilibrium, both parties always engage in (weak) strategic extremism: $p_{+1}^{**} \geq p_{+1}^*$ and $p_{-1}^{**} \leq p_{-1}^*$. (See Figure 6b.)*

Combined, Propositions 7 and 8 deliver the main lesson of this section. By taking extreme positions, parties push their targets apart:

$$p_{-1}^{**} \leq p_{-1}^* < 0 < p_{+1}^* \leq p_{+1}^{**},$$

so

$$\Delta_p^{**} \geq \Delta_p^*.$$

This in turn increases (weakly, and sometimes strictly) the average drift μ of complexity, and correspondingly the long-run complexity statistic K . To summarize: strategic behavior takes the form of strategic extremism, which may generate long-run policy complexity.

To understand the logic of strategic extremism, consider the tradeoff that party +1 faces between targeting his ideal position p_{+1}^* versus a more extreme position $p_{+1}^{**} > p_{+1}^*$. Intuitively, an extreme target "shifts the goalposts" upwards. This initial upward shift puts policy further from his ideal p_{+1}^* in the short run. But, after party -1 takes control and moves policy below p_{+1}^* , the initial shift becomes advantageous for +1 because policy is now *closer* to p_{+1}^* than it would have been under the counterfactual where no initial shift occurred. This advantage is maintained even after subsequent control changes, at least until policy reaches the other target p_{-1}^{**} or moves back above p_{+1}^* . Strategic extremism is optimal if the later benefits outweigh the earlier costs.¹⁶

¹⁶Under stronger assumptions, we may also quantify the extent of strategic extremism by each party, $\Delta_i^{**} = |p_i^{**} - p_i^*|$. For example, in the Appendix (Proposition B.1), we calculate Δ_i^{**} at the asymptotic limit where Δ_p^* and γ^{-1} are large. We show that a more vulnerable party (high λ_i) engages in more strategic extremism (high Δ_i^{**}).

Because strategic extremism entails short-run costs and medium-run benefits, it arises only if players are sufficiently patient. (Indeed, we already know from Section 3 that myopic parties do not engage in strategic extremism.) The following proposition presents a version of this intuition. Given targets p_{+1}^{**} and p_{-1}^{**} , define a binary indicator for strategic extremism:

$$\phi = \begin{cases} 1 & \text{if } \Delta_p^{**} > \Delta_p^* \\ 0 & \text{if } \Delta_p^{**} = \Delta_p^* \end{cases} .$$

Proposition 9. *Fix p_{+1}^* , p_{-1}^* , and $\lambda_{+1} \neq \lambda_{-1}$, so that parties have unequal durations. There exists $\bar{\gamma} < \infty$ such that the following hold if friction is high ($\gamma^{-1} > \bar{\gamma}^{-1}$):*

1. *There is a unique pair of equilibrium targets, $(p_{+1}^{**}, p_{-1}^{**})$.*
2. *ϕ is weakly decreasing in the parties' discount rates r_{+1} and r_{-1} .*
3. *For sufficiently small discount rates, strategic extremism occurs: $\phi = 1$.*

So, given that strategic extremism occurs only if players are sufficiently patient, kludge may occur with patient players despite being impossible under myopic players. From the perspective of a planner who dislikes complexity, *patience is not necessarily a virtue*.

5 Concluding Remarks

Throughout this paper, we have emphasized the applications of our model to public policy. However, we view our model as also being relevant to other settings where the design of complicated contracts or policies involves political or ideological disagreement; for example, in the politics of organizational design, or in the decentralized development of open-source software. In particular, the insights we derive in the model can be straightforwardly reinterpreted for an organizational context. For example, our results on long-run kludge suggest that political conflict between different factions within an organization may give rise to persistently inefficient bureaucratic routines and procedures within the organization.

In our model, the structure and density of entanglement – as captured by the entanglement constraint – is specified exogenously. This captures crudely the premise that entanglements between elements of complicated systems are difficult to anticipate, and arise inevitably during the design process.¹⁷ A more nuanced approach would be to partially endogenize entanglement; for example, by allowing parties to reduce or increase the entanglement of new rules from a ‘baseline’ level, perhaps at a cost. Such a setting may produce additional insights. For example, policymakers may deliberately enact highly-entangled rules, so as to obstruct their rivals from undoing those rules in the future.

¹⁷Readers who have written and debugged computer programs will surely sympathize with this premise.

A The Entanglement Constraint

This appendix presents two distinct microfoundations for the entanglement constraint (5). The second formulation (Appendix A.2) allows us to parametrize the degree of entanglement, and derive comparative statics.

A.1 Linear Network

Here, the policy is a set of *rules* endowed with a total order \succ . We say that π' depends on π if $\pi \succ \pi'$. We start by describing the policymaking technology in a discrete setting where each rule has small but positive mass ϵ , and time proceeds in discrete intervals of length ϵ/γ . (We will interpret γ later.) This discrete setting serves to build intuition for the role of dependencies in our analysis. We subsequently focus on the limit $\epsilon \rightarrow 0$, where the policymaking technology simplifies to a tractable continuous formulation. As before, p_j denotes the mass of j -rules in \mathbf{p} , and $\|\mathbf{p}\| = p_+ + p_-$.

A new rule π added at time t is uniformly randomly allocated a position in the order \succ . That is, at the moment of addition, π is equiprobably k -th in the ordering for all $k \in \{1, 2, \dots, |\mathbf{p}(t)|\}$, where $|\mathbf{p}(t)|$ is the number of rules in $\mathbf{p}(t)$ (including π). Ordering is pairwise persistent: if $\pi \succ \pi'$ at time t , then $\pi \succ \pi'$ for all future times τ that $\pi, \pi' \in \mathbf{p}(\tau)$. For simplicity of exposition, consider a single party who is always in control. At the start of each interval, the policymaker may choose any of the following actions, which is then realized at the end of the interval.

1. Add a new j -rule in either direction $j \in \{+, -\}$.
2. Delete the \succ -maximal rule.

At time t , the party observes the direction of each rule in $\mathbf{p}(t)$ and the history of all added and deleted rules up till time t , but does not observe the ordering between rules. So, if he chooses “delete” at the start of a time interval, then he observes which rule was \succ -maximal (and thus was deleted) only at the end of the interval.

In general, one might expect the party’s beliefs about the dependency ordering \succ to evolve in a complicated fashion. Conveniently, our technical assumptions allow us to abstract from the details of (beliefs about) the ordering.

Remark A.1. *At any time- t history, from the policymaker’s perspective, every permutation of the dependency ordering over $\mathbf{p}(t)$ is equally likely.*

So, rules are indistinguishable beyond their direction: all positive rules look alike, and all negative rules look alike. Consequently, the policymaker’s beliefs are summarized by the masses (p_+, p_-) of positive and negative rules.

Now, we calculate how p_+ and p_- change over a single time interval under addition and deletion. Let Δp_j denote the change in p_j over a single time interval.

If the party adds a j -rule, then (remembering that $\Delta t = \epsilon/\gamma$)

$$\Delta p_j = \gamma \Delta t \text{ and } \Delta p_{-j} = 0.$$

If the party deletes the \succ -maximal rule, it is equally likely to be any of the existing rules in $\mathbf{p}(t)$. So, deletion preserves (in expectation) the ratio of positive to negative rules in the policy:

$$\mathbb{E} [\Delta p_j] = -\gamma \frac{p_j}{\|\mathbf{p}\|} \Delta t \text{ for each } j \in \{+, -\}.$$

More generally, if party mixes over positive rule addition, negative rule addition, deletion, and doing nothing, then he can achieve (in expectation) any convex combination of addition and deletion outcomes:

$$\mathbb{E} [\Delta p_j] = \left(\alpha_j - \frac{p_j}{\|\mathbf{p}\|} \delta \right) \Delta t \text{ for each } j \in \{+, -\}, \quad (12a)$$

for any

$$\begin{aligned} \alpha_+ \geq 0, \alpha_- \geq 0, \delta \geq 0 \text{ such that} \\ \alpha_+ + \alpha_- + \delta \leq \gamma. \end{aligned} \quad (12b)$$

Now, focus on the limit $\epsilon \rightarrow 0$, so that each rule becomes infinitesimally small and time is continuous. Here, the laws of motion (12a)–(12b) can be expressed in differential form. The party chooses addition and deletion rates $\alpha_+(t) \geq 0, \alpha_-(t) \geq 0, \delta(t) \geq 0$ which determine the velocity of (p_+, p_-) :

$$\frac{d}{dt} p_j(t) = \alpha_j(t) - \frac{p_j(t)}{\|\mathbf{p}(t)\|} \delta(t) \text{ for each } j \in \{+, -\}, \quad (13a)$$

subject to a flow constraint

$$\alpha_+(t) + \alpha_-(t) + \delta(t) \leq \gamma. \quad (13b)$$

Together, Equations (13a) and (13b) are equivalent to the law of motion (3) and entanglement constraint (5).

A.2 Random Network

We start with an informal description of the model.

A policy is a continuum of infinitesimal rules. We adopt the convenient expositional convention that each rule has (infinitesimal) mass ϵ . Rules are linked to form an undirected network. Whenever a new rule π is created, it randomly forms a link to each existing rule with (infinitesimal) probability $\rho\epsilon$. So, each new rule π forms ρ links (in expectation) per unit mass of existing rules. We interpret ρ as the degree of *entangledness*. Once formed, links between pairs of rule persist until one or both rules in the pair are deleted.

As in Appendix A.1, consider a single party who is always in control. The party can add rules in either direction, but cannot precisely target a given rule for deletion. Specifically, if the party targets a rule π to be deleted, the direct neighbours of π will also be simultaneously deleted. As before, the maximum rate of addition and deletion (of all rules, including the neighbours of rules targeted for deletion) is mass γ per unit time.

When formalizing this description, we distinguish between rules added at different times when describing the policy. Say that a rule has vintage- τ if it was added at time τ .

Define $m_j(t, \tau)$ to be the “quantity” of vintage- τ j -rules that remain at time $t \geq \tau$. Specify the law of motion of $m_j(t, \tau)$ to be

$$p_j(t, \tau) = \alpha_j(\tau) - \int_{\tau}^t \delta_j(\tilde{t}, \tau) d\tilde{t}, \quad (14)$$

where $\alpha_j(\tau)$ is the time- τ addition rate for j -rules and $\delta_j(\tilde{t}, \tau)$ is the time- \tilde{t} deletion rate for j -rules of vintage τ . That is, the time- t quantity of vintage- τ rules equals the time- τ addition rate, less the total quantity of vintage- τ rules deleted up till time t .

Let $\hat{\delta}_+(t, \tau)$ and $\hat{\delta}_-(t, \tau)$ be the time- t rate at which the party targets vintage- τ rules for deletion. Let the network structure be characterized by $\rho(j, \tau, \tilde{j}, \tilde{\tau})$, which represents the density of connections between j -rules of vintage- τ and \tilde{j} -rules of vintage- $\tilde{\tau}$. To capture the idea that immediate neighbours of deleted rules must also be deleted, we specify that the vintage- τ deletion rate accounts both for directly targeted rules, and for neighbours of targeted rules from other vintages:

$$\delta_j(t, \tau) = \hat{\delta}_j(t, \tau) + p_j(t, \tau) \sum_{\tilde{j}} \int_0^t \rho(j, \tau, \tilde{j}, \tilde{\tau}) \hat{\delta}_{\tilde{j}}(t, \tilde{\tau}) d\tilde{\tau}.$$

Define the mass of j -rules at time t to be the total quantity of rules, integrated over all vintages: $m_j(t) = \int_0^t m_j(t, \tau) d\tau$. Applying (14), we get

$$m_j(t) = \int_0^t \alpha_j(\tau) d\tau - \int_0^t \delta_j(\tilde{t}) d\tilde{t}, \text{ where } \delta_j(t) = \int_0^t \delta_j(t, \tau) d\tau.$$

In differential form, this replicates (3) from Section 2:

$$\frac{d}{dt} p_j(t) = \alpha_j(t) - \delta_j(t).$$

Naturally, we interpret $\delta_j(t)$ to be the rate at which j -rules are being deleted.

We specify that at each time t , the party chooses addition rates $\alpha_+(t)$ and $\alpha_-(t)$, and vintage-specific deletion rates $\hat{\delta}_+(t, \tau)$ and $\hat{\delta}_-(t, \tau)$, subject to the familiar flow constraint (4) on the overall rate of addition and deletion:

$$\alpha_+(t) + \alpha_-(t) + \delta_+(t) + \delta_-(t) \leq \gamma. \quad (15)$$

Our key simplifying assumption is that the density of links across vintages and directions of rules is completely homogenous: $\rho(\tau, \tilde{\tau}, j, \tilde{j}) \equiv \rho > 0$. In that case, some algebra reveals that the deletion rate is

$$\delta_j(t) = \hat{\delta}_j(t) + \rho m_j(t) (\hat{\delta}_+(t) + \hat{\delta}_-(t)), \text{ where } \hat{\delta}_j(t) = \int_0^t \hat{\delta}_j(t, \tau) d\tau. \quad (16)$$

Thus, the deletion rates $\delta_+(t)$ and $\delta_-(t)$ are determined entirely by the targeted deletion rates $\hat{\delta}_+(t)$ and $\hat{\delta}_-(t)$. In other words, it does not matter which rules to target; the only relevant decision is how many rules to delete.

Further, inspection of (16) indicates that the party can, by appropriately choosing deletion rates $\hat{\delta}_+(t)$ and $\hat{\delta}_-(t)$, achieve any combination of deletion rates $\delta_+(t)$ and $\delta_-(t)$ satisfying

$$\frac{\delta_+(t)}{\delta_-(t)} \in \left[\frac{1 + \rho p_+(t)}{\rho p_-(t)}, \frac{\rho p_+(t)}{1 + \rho p_-(t)} \right]. \quad (17)$$

Accordingly, we may restate the laws of motion as follows.

At each time t , the party chooses addition rates $\alpha_+(t)$ and $\alpha_-(t)$ and deletion rates $\delta_+(t)$ and $\delta_-(t)$, subject to the flow constraint (15) and the entanglement constraint (17).

Notice that (17) is a relaxed version of Section 2's entanglement constraint, (5). At the limit $\rho \rightarrow \infty$, where the network density becomes large, (17) tightens into (5).

Our results from Section 3 continue to hold in this setting, even with finite entangledness ρ . As before, suppose that both parties are myopic. Define the basin \mathcal{B} , as before, to be the set of regular policies where at least one party deletes rules. Here :

$$\mathcal{B} = \left\{ \mathbf{p} : \left(\frac{\rho p_+}{1 + \rho \|\mathbf{p}\|} < \frac{1}{z_{+1}} \text{ or } \frac{\rho p_-}{1 + \rho \|\mathbf{p}\|} > \frac{1}{z_{-1}} \right) \text{ and } p \in [p_{-1}^*, p_{+1}^*] \right\}$$

The basin \mathcal{B} expands as entangledness ρ decreases. This is intuitive: as the entanglement constraint loosens, the ability of each party to target rules for deletion improves, and thus deletion becomes optimal over a larger range of policies. Proposition 2 continues to hold in this setting: any policy in \mathcal{B} remains forever in \mathcal{B} .

Outside the basin \mathcal{B} , the laws of motion (9a) and (9b) continue to hold. Consequently, all of our results about kludge – Propositions 3, Proposition 4, 5a, 5b, and 6 – are preserved. Further, we may show that kludge increases with entangledness ρ :

Proposition A.1. *Suppose $K = 1$. Then κ is decreasing in ρ .*

The proof is almost identical to that of Proposition 5b, and thus is omitted. An increase in entangledness ρ shrinks the basin \mathcal{B} . Consequently, policy becomes less likely to enter the basin and get trapped; rather, policy is more likely to ‘escape’ and become kludged.

B Proofs

Short-Run Dynamics

Notation Denote the convex closure of a set X as $\text{Conv}(X)$. Denote the sequence of times at which control changes hands from party i to party $-i$ as t_1^i, t_2^i, \dots . Throughout, we assume WLOG that party +1 has control at $t = 0$; so, $0 < t_1^{+1} < t_1^{-1} < t_2^{+1} < t_2^{-1} \dots$; notice that the sequence of transition times $\{t_1^{+1}, t_1^{-1}, t_2^{+1}, t_2^{-1} \dots\}$ fully determines the equilibrium path $(\mathbf{p}(t), i(t))$. Define $\Delta t_k^{+1} \equiv t_k^{+1} - t_{k-1}^{-1}$ and $\Delta t_k^{-1} \equiv t_k^{-1} - t_k^{+1}$ to be the sequences of durations for which each party was in control.

Proof of Propositions 1a and 1b Focus on party +1; the calculation for party –1 is similar. Start with the case $p \in [p_{-1}^*, p_{+1}^*]$. There, party +1's problem is to maximize the (linear) objective

$$\frac{\partial}{\partial t} u_{+1}(\mathbf{p}(t)) = z_{+1} \frac{d}{dt} p - \frac{d}{dt} \|\mathbf{p}\|$$

subject to the constraint (2), which corresponds in $(\frac{d}{dt} p, \frac{d}{dt} \|\mathbf{p}\|)$ -space to a triangle

$\text{Conv}(\{v_+, v_-, v_\delta\})$ with vertices

$$v_+ = \gamma \cdot (1, 1), v_- = \gamma \cdot (-1, 1), v_\delta = \gamma \cdot (p/\|\mathbf{p}\|, -1).$$

A linear objective over a simplex is, of course, maximized at one of the vertices of the simplex. Some algebra reveals that vertex v_+ is optimal (maximizes the objective) when $-\frac{p}{\|\mathbf{p}\|} > 1 - \frac{2}{z_{+1}}$; otherwise, vertex v_- is optimal.

The case $p = p_{+1}^*$ is slightly more involved. Here, the objective is no longer linear in $(\frac{d}{dt}p, \frac{d}{dt}\|\mathbf{p}\|)$; specifically,

$$\frac{\partial}{\partial t} u_{+1}(\mathbf{p}(t)) = z_{+1} \left| \frac{d}{dt} p \right| - \frac{d}{dt} \|\mathbf{p}\|.$$

Notice, however, that this objective is linear on each of the half-planes $\frac{d}{dt}\|\mathbf{p}\| \leq 0$ and on $\frac{d}{dt}\|\mathbf{p}\| \geq 0$. The intersection of each half-plane with the triangle $\text{Conv}(\{v_+, v_-, v_\delta\})$ defines two simplices in $(\frac{d}{dt}p, \frac{d}{dt}\|\mathbf{p}\|)$ -space over which the objective function is linear:

$$\begin{aligned} \frac{\partial}{\partial t} u_{+1}(\mathbf{p}(t)) &= -z_{+1} \frac{d}{dt} p - \frac{d}{dt} \|\mathbf{p}\| \text{ over } \text{Conv}(\{v_+, v_{m+}, v_{m-}\}) \text{ and} \\ \frac{\partial}{\partial t} u_{+1}(\mathbf{p}(t)) &= z_{+1} \frac{d}{dt} p - \frac{d}{dt} \|\mathbf{p}\| \text{ over } \text{Conv}(\{v_-, v_\delta, v_{0-}, v_{0+}\}) \text{ where} \\ v_{0-} &= \gamma \cdot (0, -\frac{\|\mathbf{p}\| - p}{\|\mathbf{p}\| + p}) \text{ and } v_{0+} = \gamma \cdot (0, 1). \end{aligned}$$

Consequently, the objective function is maximized on one of the vertices of the two simplices. Some further algebra reveals that vertex v_δ is optimal if $-\frac{p}{\|\mathbf{p}\|} < 1 - \frac{2}{z_{+1}}$; otherwise, vertex v_{0-} is optimal.

One final point: when $p = \|\mathbf{p}\| = p_{+1}^*$, vertex v_{0-} results in $\frac{d}{dt}p = \frac{d}{dt}\|\mathbf{p}\| = 0$, and thus is equivalent to stagnation: $\alpha_j = \alpha_{-j} = \delta = 0$. ■

Path Dependence and Kludge

Proof of Remark 1 Remark 1.1 follows immediately from Propositions 1a and 1b, so we only prove Remark 1.2 here. WLOG, consider the case where $p > p_{+1}^*$. Similarly to the proof of Propositions 1a and 1b, we may characterize each parties optimal strategy.

- If $\frac{p_-}{\|\mathbf{p}\|} < \frac{1}{z_{+1}}$, then party +1 deletes rules, thus moving towards his ideal: $(\alpha_+, \alpha_-, \delta) = \gamma \cdot (0, 0, 1)$, so $\frac{d}{dt}p = -\frac{\|\mathbf{p}\| - p}{\|\mathbf{p}\| + p}$.
- If $\frac{p_-}{\|\mathbf{p}\|} > \frac{1}{z_{+1}}$, then party +1 adds negative rules: $(\alpha_+, \alpha_-, \delta) = \gamma \cdot (0, 1, 0)$, so $\frac{d}{dt}p = -1$.
- -1 always adds negative rules: $(\alpha_+, \alpha_-, \delta) = \gamma \cdot (0, 1, 0)$, so $\frac{d}{dt}p = -1$.

The take-away point is that policy position always shifts negatively: $\frac{d}{dt}p \leq -\frac{\|\mathbf{p}\| - p}{\|\mathbf{p}\| + p}$. In fact, we may show by induction that $\|\mathbf{p}(t)\| \leq \|\mathbf{p}(0)\| + p(0) - p_{+1}^*$; consequently, $\frac{d}{dt}p(t) \leq -\frac{\|\mathbf{p}(0)\| + p(0) - p_{+1}^* - p(0)}{\|\mathbf{p}(0)\| + p(0)} = -\frac{\|\mathbf{p}(0)\| - p_{+1}^*}{\|\mathbf{p}(0)\| + p(0)}$ for all $t \geq 0$. We conclude that policy reaches the +1-ideal position $p = p_{+1}^*$ in finite time. ■

Proof of Proposition 2 Proposition 2.3 follows directly from Proposition 1a: consider a perfectly simple policy consisting purely of j -rules, $m_{-j} = 0$, and let $j = \text{sgn } i$. Then party i adds j -rules, whereas party $-i$ deletes j -rules. In either case, policy remains perfectly simple.

Next, consider Proposition 2.1. Given $j = \text{sgn } i$, let \mathcal{B}_i be the set of policies at which party i deletes rules,

$$\mathcal{B}_i = \left\{ \mathbf{p} : \frac{p_j}{\|\mathbf{p}\|} < \frac{1}{z_i} \text{ and } p \in [p_{-1}^*, p_{+1}^*] \right\};$$

note that $\mathcal{B} = \mathcal{B}_{+1} \cup \mathcal{B}_{-1}$. Consider, WLOG, \mathcal{B}_{+1} . Assume for now that $\frac{1}{z_{+1}} + \frac{1}{z_{-1}} < 1$. Here, \mathcal{B}_{+1} and \mathcal{B}_{-1} do not intersect, except at the empty policy $\mathbf{p} = (0, 0)$. Consequently, policy dynamics within \mathcal{B}_{+1} , other than at the empty policy, take the following form:

- If $p > p_{-1}^*$, then party -1 adds negative rules, so $\frac{d}{dt} p_+ = 0$, $\frac{d}{dt} p_- = 1$, $\frac{d}{dt} \|\mathbf{p}\| = 1$. Calculations reveal $\frac{d}{dt} \frac{p_+(t)}{\|\mathbf{p}(t)\|} = \frac{-p_+}{\|\mathbf{p}\|^2} \leq 0$.
- If $p = p_{-1}^* < 0$, then party -1 reduces complexity, so $\frac{d}{dt} p = 0$, $\frac{d}{dt} \|\mathbf{p}\| < 0$. We immediately see that $\frac{d}{dt} \frac{p_+(t)}{\|\mathbf{p}(t)\|} = \frac{1}{2} \frac{d}{dt} \left(\frac{p(t)}{\|\mathbf{p}(t)\|} + 1 \right) < 0$.
- Party $+1$ always deletes rules, so $\frac{d}{dt} p = -p/\|\mathbf{p}\|$, $\frac{d}{dt} \|\mathbf{p}\| = -1$. Clearly, $\frac{d}{dt} \frac{p_+(t)}{\|\mathbf{p}(t)\|} = \frac{1}{2} \frac{d}{dt} \left(\frac{p(t)}{\|\mathbf{p}(t)\|} + 1 \right) = 0$.

In all cases (except the empty policy), $\frac{p_+(t)}{\|\mathbf{p}(t)\|}$ is weakly decreasing; so policy remains within \mathcal{B}_{+1} .

Now, relax the assumption that $\frac{1}{z_{+1}} + \frac{1}{z_{-1}} < 1$. Policy dynamics remain the same as above, except that at the intersection of \mathcal{B}_{+1} and \mathcal{B}_{-1} , party -1 deletes rules (instead of adding rules or reducing complexity), so that $\frac{d}{dt} \frac{p_+(t)}{\|\mathbf{p}(t)\|} = 0$. Clearly, this does not change our conclusion, as policy remains within \mathcal{B}_{+1} .

Our argument so far for Proposition 2.1 has neglected the empty policy; but this case is covered by Proposition 2.1. Both parties add rules at the empty policy, so policy remains perfectly simple ($p/\|\mathbf{p}\| = 1$) and thus remains in \mathcal{B} .

Finally, consider Proposition 2.2. Note that the complexity of any policy in \mathcal{B} is bounded above by some \bar{c} . Note, also, that if policy is initially in \mathcal{B}_i , then it always remains within \mathcal{B}_i unless policy becomes perfectly simple. Because the time periods between changes of control are i.i.d. and exponentially distributed, almost surely, the following event will eventually occur: (i) party i is in control at time t , (ii) policy $\mathbf{p}(t)$ is in \mathcal{B} , and (iii) i retains control for a period of at least \bar{c} . But because party i deletes rules from policy until he loses control, at time $t + \|\mathbf{p}(t)\|$, he reaches the empty policy (which is perfectly simple). ■

Proof of Remark 2 Consider the case where $z_{+1} = 2$. We will generalize to the case where $z_{+1} < 2$ later. The focus on z_{+1} is WLOG. Given $z_{+1} = 2$, \mathcal{B}_{+1} takes the form $\{\mathbf{p} : 0 \geq p \geq p_{-1}^*\}$. As a result, policy avoids the basin only if position remains forever within the interval $0 < p \leq p_{+1}^*$.

Outside the basin, party -1 adds negative rules. If -1 is ever in control for a contiguous period of longer than p_{+1}^*/γ , then he will decrease policy position by at least p_{+1}^* , and thus move policy into the basin. Such an event occurs with probability $e^{-\lambda_{-1}^{-1} p_{+1}^*/\gamma} > 0$ each time that -1 regains control. Since -1 regains control an infinite number of times almost surely, it follows that policy will almost surely enter the basin.

For the case $z_{+1} < 2$, notice that \mathcal{B}_{+1} expands as z_{+1} decreases; so the argument above continues to hold. ■

To prove Lemma 1 and Proposition 3, let's introduce some tools.

Lemma B.1a. *There exists $\tau > 0$ and $\xi > 0$ such that for all $(q(0), i(0))$,*

$$\Pr [(q(\tau), i(\tau)) = (p_{+1}^*, +1)] \geq \xi.$$

Proof. Choose $\tau = (p_{+1}^* - p_{-1}^*)/\gamma + \epsilon$ for some $\epsilon > 0$. Consider a history such that for some $t_0 < \epsilon$, $i(t) = +1$ for all $t \in [t_0, \tau]$. Then, obviously, $(q(\tau), i(\tau)) = (p_{+1}^*, +1)$. Further, the probability of such a history is at least $\xi = e^{-\lambda_{+1}\tau} (1 - e^{-\lambda_{-1}\epsilon})$. ■

Proof of Lemma 1

Lemma B.1a ensures that, in the language of ergodic theory, the entire state space $[p_{-1}^*, p_{+1}^*] \times \{+1, -1\}$ is *small* with respect to the process $(q(t), i(t))$. Lemma 1 then follows immediately from standard results in ergodic theory – see, for example, Bhattacharya and Majumdar (2003, Corollary 3.3). ■

Lemma B.1b. *The unique invariant (steady-state) distribution G of the process $(q(t), i(t))$ on $[p_{-1}^*, p_{+1}^*] \times \{+1, -1\}$ has density*

$$g(q, +1) \equiv g(q, -1) \equiv Ae^{\frac{\lambda_{+1}-\lambda_{-1}}{\gamma}q} \quad (18a)$$

for $p_{-1}^* \leq q \leq p_{+1}^*$, where A is a normalizing constant, and has atoms

$$\Delta G(p_{+1}^*, +1) = \frac{\gamma}{\lambda_{+1}} g(p_{+1}^*, +1) \text{ and } \Delta G(p_{-1}^*, -1) = \frac{\gamma}{\lambda_{-1}} g(p_{-1}^*, -1) \quad (18b)$$

at the each party's ideal, $(p_{-1}^*, -1)$ and $(p_{+1}^*, +1)$.

Proof. The steady-state distribution of (q, i) is invariant to the law of motion (9a) of $q(t)$ and of $i(t)$. For $q < p_{+1}^*$, over a small time interval Δt , the net change in the probability mass of $[q, q + \Delta q] \times \{+1\}$ must be zero; that is,

$$[\gamma g(q, +1) \Delta t - \gamma g(q + \Delta p, +1) \Delta t] + [\lambda_{-1} g(q, -1) \Delta q \Delta t - \lambda_{+1} g(q, +1) \Delta q \Delta t] \approx 0.$$

Taking the limit $\Delta q, \Delta t \rightarrow 0$, we get

$$\gamma g_q(q, +1) = \lambda_{-1} g(q, -1) - \lambda_{+1} g(q, +1) \text{ and} \quad (19)$$

$$\gamma g_q(q, -1) = \lambda_{-1} g(q, -1) - \lambda_{+1} g(q, +1) \quad (20)$$

for $q \in [p_{-1}^*, p_{+1}^*]$, where (20) holds by a symmetric argument. Solving the differential equations (20) and (19) simultaneously reveals that

$$g(q, +1) \equiv g(q, -1) \equiv Ae^{\frac{\lambda_{-1}-\lambda_{+1}}{\gamma}q}$$

for some constant A .

Notice that we have implicitly assumed that there are no atoms on $[p_{-1}^*, p_{+1}^*] \times \{+1\}$ or (symmetrically) on $(p_{-1}^*, p_{+1}^*) \times \{-1\}$. This holds because, if $(q, +1)$ were an atom, then the

law of motion (9a) dictates (impossibly) that $(q', +1)$ would also be an atom for all q' in some right-neighbourhood of q .

Finally, consider the (potential) atoms $\Delta G(p_{+1}^*, +1)$ and $\Delta G(p_{-1}^*, -1)$. Over a small time interval Δt , the net change in the probability mass of each atom must be zero; that is,

$$\begin{aligned}\lambda_{+1}\Delta G(p_{+1}^*, +1)\Delta t - \gamma g(p_{+1}^*, +1)\Delta t &\approx 0, \\ \lambda_{-1}\Delta G(p_{-1}^*, -1)\Delta t - \gamma g(p_{+1}^*, +1)\Delta t &\approx 0\end{aligned}$$

or, more compactly,

$$\Delta G(p_{+1}^*, +1) = \frac{\gamma}{\lambda_{+1}}g(p_{+1}^*, +1) \text{ and } \Delta G(p_{-1}^*, -1) = \frac{\gamma}{\lambda_{-1}}g(p_{-1}^*, -1).$$

■

Define a class of simulacra $c_\varepsilon(t)$ of the ‘true’ complexity process $\|\mathbf{p}(t)\|$, each of which is coupled to the position simulacrum $q(t)$: for $\varepsilon \geq 0$,

$$\frac{d}{dt}c_\varepsilon(t) \equiv v_\varepsilon(q(t))$$

where

$$v_\varepsilon(q) \equiv \gamma \cdot \begin{cases} -(1 - \varepsilon) & : q \in \{p_{-1}^*, p_{+1}^*\} \\ 1 & : q \in (p_{-1}^*, p_{+1}^*) \end{cases}.$$

The parameter ε captures how quickly the complexity simulacrum c decreases whenever the position simulacrum q is at either ideal. Conveniently, denote $c(t) \equiv c_0(t)$. Notice that at the extreme $\varepsilon = 0$, $v_0(q) \equiv v(q)$: the complexity simulacrum behaves as true complexity does at the limit $\|\mathbf{p}\| \rightarrow \infty$.

Lemma B.2a. *Consider the simulacrum process with $\varepsilon = 0$. Suppose $z_{+1} > 2$ and $z_{-1} > 2$, which ensures that \mathcal{B} is finite in extent. Select sufficiently large \underline{c} so that $\mathcal{B} \subset \{\mathbf{p} : \|\mathbf{p}\| < \underline{c}\}$. Fix a start time $t_0 \geq 0$, and suppose initial conditions are identical for the true and simulacrum process, as follows: $\|\mathbf{p}(t_0)\| \equiv c(t_0)$ and $p(t_0) = q(t_0) > \underline{p}$. Consider any history up to time $T \leq \infty$ such that $c(t) \geq \underline{c}$ for all $t \in [t_0, T)$. Then $q(t) = \underline{p}(t)$ and $c(t) \leq \|\mathbf{p}(t)\|$ for all $t \in [t_0, T)$.*

Proof. This result requires only a straightforward inspection of the laws of motion of \mathbf{p} (outside \mathcal{B}) and c, q : we have that $\frac{d}{dt}p(t) = \frac{d}{dt}q(t)$ and $\frac{d}{dt}\|\mathbf{p}(t)\| \leq \frac{d}{dt}c(t)$, so $p(t) \equiv q(t)$ and $\|\mathbf{p}(t)\| \leq c(t)$. ■

Lemma B.2b. *Suppose $z_{+1} > 2$ and $z_{-1} > 2$. Select sufficiently small ε and sufficiently large \underline{c} so that $\mathcal{B} \subset \{\mathbf{p} : \|\mathbf{p}\| < \underline{c}\}$ and so that $1 - \varepsilon < \frac{\underline{c} - \underline{p}}{\underline{c} + \underline{p}}$ for $p = p_{+1}^*$ and $p = p_{-1}^*$. Fix a start time $t_0 \geq 0$, and suppose initial conditions are identical for the true and simulacrum process, as follows: $\|\mathbf{p}(t_0)\| \equiv c_\varepsilon(t_0) \geq \underline{c}$ and $p(t_0) = q(t_0)$. Suppose that at some time $T > t_0$, $c_\varepsilon(T) \leq \underline{c}$. Then at some time $\tau \in (t_0, T]$, $\|\mathbf{p}(\tau)\| = \underline{c}$ and $p(\tau) = p_{i(\tau)}$.*

Proof. Suppose, towards a contradiction, that $\|\mathbf{p}(t)\| > \underline{c}$ for all $t \in (t_0, T]$. Then throughout this time interval, $\frac{d}{dt}p(t) = \frac{d}{dt}q(t)$ and $\frac{d}{dt}\|\mathbf{p}(t)\| \geq \frac{d}{dt}c_\varepsilon(t)$, so $p(t) \equiv q(t)$ and $c_\varepsilon(t) \geq \|\mathbf{p}(t)\| > \underline{c}$. This contradicts the assumption that $c(T) \leq \underline{c}$. Finally, note that for $\|\mathbf{p}\| \geq \underline{c}$, $\|\mathbf{p}\|$ decreases only while position is some ideal p_i^* and i is in control; thus whenever complexity $\|\mathbf{p}\|$ attains \underline{c} from above at some time τ , we must have $p(\tau) = p_{i(\tau)}$. ■

Lemma B.2c. *Suppose $z_{+1} > 2$ and $z_{-1} > 2$. Select sufficiently large \underline{c} so that $\mathcal{B} \subset \{\mathbf{p} : \|\mathbf{p}\| < \underline{c}\}$. Define $\mathcal{B}' = \{\mathbf{p} : \|\mathbf{p}\| \leq \underline{c} \text{ and } p \in [p_{-1}^*, p_{+1}^*] \text{ and } \mathbf{p} \notin \mathcal{B}\}$. If $\mathbf{p}(0) \in \mathcal{B}'$, then $\mathbf{p}(t)$ eventually leaves (and possibly returns to) \mathcal{B}' in finite time (a.s.). Further, there exists $\nu > 0$ so that if $p(0) = p_{i(0)}$, then with probability of at least ν , $\mathbf{p}(t)$ exits \mathcal{B}' by entering the basin \mathcal{B} .*

Proof. Let $\tau_1 = \frac{p_{+1}^* - p_{-1}^*}{\gamma}$ be the amount of time taken for policy to move from ideal p_{-i}^* to p_i^* if j -rules are added at the maximum rate γ . Let $\tau_i < \infty$ be the time taken for policy to move from $\mathbf{p} = (\underline{c}, p_i^*)$ to reach \mathcal{B} if parties reduce complexity along p_i^* at the maximum rate, $\frac{d}{dt} \|\mathbf{p}\| = -\frac{\|\mathbf{p}\| - |p_i^*|}{\|\mathbf{p}\| + |p_i^*|}$. Let $\tau_2 = \max\{\tau_{+1}, \tau_{-1}\}$. If party i takes control at time t_0 with $\mathbf{p}(t_0) \in \mathcal{B}'$ and retains control for duration $\Delta t^i > \tau_1 + \tau_2$, then policy must exit \mathcal{B}' in one of two ways at some $t \in (t_0, t_0 + \Delta t^i)$, while i is in control: either (i) policy complexity $\|\mathbf{p}(t)\|$ exceeds \underline{c} , or (ii) policy enters the basin, $\|\mathbf{p}(t)\| \in \mathcal{B}$. In other words, $\mathbf{p}(t)$ remains in \mathcal{B}' forever only if for all control durations, $\Delta t_k^i \leq \tau_1 + \tau_2$. But this almost never occurs; so, policy almost surely leaves \mathcal{B}' .

Further, suppose $p(0) = p_{i(0)}$. If $i(0)$ retains control for duration for at least τ_2 , then he will reduce complexity along $p_{i(0)}$; by definition of τ_2 , policy will thus enter the basin \mathcal{B} while $i(0)$ is in control. This occurs with probability of at least $\nu = e^{-\max\{\lambda_{+1}, \lambda_{-1}\}\tau_2} > 0$. ■

Let's introduce some further notation. For each $i \in \{+1, -1\}$, define $\{\tau_1^i < \tau_2^i < \dots\}$ to be the subsequence of $\{t_1^i, t_2^i, \dots\}$ corresponding to the times where i loses control to $-i$ while the position simulacrum is at i 's ideal (i.e., $q(t_k^i) = p_i^*$). Note that each τ_k^i is a stopping time relative to the filtration generated by $(q(t), i(t))$. For $k = 1, 2, \dots$, define $\Delta c_k^{i,\varepsilon} \equiv c_\varepsilon(\tau_{k+1}^i) - c_\varepsilon(\tau_k^i)$ to be the change in the complexity simulacrum between the k -th and $(k+1)$ -th times i loses control while at his ideal. Analogously, define $\Delta \tau_k^i \equiv \tau_{k+1}^i - \tau_k^i$.

The sequences $\{\Delta c_1^{+1,\varepsilon}, \Delta c_2^{+1,\varepsilon}, \dots\}$ and $\{\Delta c_1^{-1,\varepsilon}, \Delta c_2^{-1,\varepsilon}, \dots\}$ have the following useful properties.

Lemma B.3a. *For each $i \in \{+1, -1\}$, the random variables $\Delta c_1^{i,\varepsilon}, \Delta c_2^{i,\varepsilon}, \dots$ are i.i.d., as are the random variables $\Delta t_1^i, \Delta t_2^i, \dots$*

Proof. Follows immediately from the fact that $(q(t), i(t))$ is a strong Markov process and $\frac{d}{dt} c_\varepsilon(t)$ depends only on $q(t)$. ■

Lemma B.3b. $\inf\{c(t) : t \geq 0\} = \inf\{c(\tau_1^{+1}), c(\tau_2^{+1}), \dots\} \cup \{c(\tau_1^{-1}), c(\tau_2^{-1}), \dots\}$

Proof. This follows immediately from the fact that the complexity simulacrum increases between ideals and decreases at ideals; and that each τ_k^i corresponds to a time at which the position simulacrum departs i 's ideal. Consequently, $\{c(\tau_1^{+1}), c(\tau_2^{+1}), \dots\} \cup \{c(\tau_1^{-1}), c(\tau_2^{-1}), \dots\}$ corresponds to the set of local minima of the complexity simulacrum process. ■

Lemma B.3c. $\mathbb{E}[\Delta \tau_k^i] < \infty$ and $|\mathbb{E}[\Delta c_k^{i,\varepsilon}]| < \infty$.

Proof. $|\mathbb{E}[\Delta c_k^{i,\varepsilon}]| \leq \gamma \mathbb{E}[\Delta \tau_k^i]$, so it is sufficient to prove that $\mathbb{E}[\Delta \tau_k^i] < \infty$. The proof of this last point involves showing that $\Delta \tau_k^i$ has exponentially-bounded tails; it is tedious and not very insightful, and thus is omitted. ■

Lemma B.3d. *For any $\varepsilon \geq 0$ and every k , the following statements are equivalent:*

1. $\mathbb{E} [\Delta c_k^{+1,\varepsilon}] \geq 0$.
2. $\mathbb{E} [\Delta c_k^{-1,\varepsilon}] \geq 0$.
3. $\int v_\varepsilon(q) dF(q) \geq 0$.

Proof. We show that 1 \iff 3; the argument that 2 \iff 3 is identical. From Lemma 1, $(q(t), i(t))$ is uniquely ergodic, so Birkhoff's ergodic theorem applies: a.s.,

$$\lim_{T \rightarrow \infty} \frac{1}{T - T_0} \int_{T_0}^T v_\varepsilon(q(t)) dt = \int v_\varepsilon(q) dF(q). \quad (21)$$

Now, write

$$\lim_{k \rightarrow \infty} \frac{1}{\tau_{k+1}^i - \tau_1^i} \int_{\tau_1^i}^{\tau_{k+1}^i} v_\varepsilon(q(t)) dt = \lim_{k \rightarrow \infty} \frac{c_\varepsilon(\tau_{k+1}^i) - c_\varepsilon(\tau_1^i)}{\tau_{k+1}^i - \tau_1^i} = \lim_{k \rightarrow \infty} \frac{\frac{1}{k} \sum_{j=1}^k \Delta c_j^{i,\varepsilon}}{\frac{1}{k} \sum_{j=1}^k \Delta \tau_j^i}.$$

Note that $\lim_{k \rightarrow \infty} \tau_{k+1}^i = \infty$ almost surely, so the LHS converges almost surely to $\int v_\varepsilon(q) dF(q)$. By the strong law of large numbers, the RHS converges almost surely to $\mathbb{E} [\Delta c_k^{i,\varepsilon}] / \mathbb{E} [\Delta \tau_k^i]$. So,

$$\int v_\varepsilon(q) dF(q) = \frac{\mathbb{E} [\Delta c_k^{i,\varepsilon}]}{\mathbb{E} [\Delta \tau_k^i]}.$$

The result follows. ■

Lemma B.3e.

1. Suppose $\mathbb{E} [\Delta c_1^{i,\varepsilon}] > 0$. Then $\lim_{k \rightarrow \infty} c(\tau_k^i) = \infty$ a.s.. Further, for any $\underline{c} < c_1^{i,\varepsilon}$, $\inf\{c(\tau_k^i)\} \geq \underline{c}$ with positive probability, and $\lim_{c(\tau_1^i) - \underline{c} \rightarrow \infty} \Pr [\inf\{c(\tau_k^i)\} \geq \underline{c}] = 1$.
2. Suppose $\mathbb{E} [\Delta c_1^{i,\varepsilon}] \leq 0$. Then $\inf_k \{c(\tau_k^i)\} = -\infty$ a.s..

Proof. This lemma is simply a restatement of classic results from large deviation theory. The cases where $\mathbb{E} [\Delta c_1^{i,\varepsilon}] \geq 0$ follow from the strong law of large numbers. The case where $\mathbb{E} [\Delta c_1^{i,\varepsilon}] = 0$ follows from the recurrence theorem. ■

Lemma B.4.

1. If $\int v_\varepsilon(q) dF(q) > 0$, then with positive probability, $\lim_{t \rightarrow \infty} c_\varepsilon(t) = \infty$ and $c_\varepsilon(t) \geq c_\varepsilon(0)$ for all $t \geq 0$.
2. If $\int v_\varepsilon(q) dF(q) < 0$, then $\inf_{t \geq 0} \{c_\varepsilon(t)\} = -\infty$ almost surely.

Proof. Follows immediately from Lemmas B.3b, B.3d, and B.3e. ■

Proof of Proposition 3

$\int v(q)dF(q) > 0$: The assumptions $z_{+1} > 2$ and $z_{-1} > 2$ ensure that the basin \mathcal{B} is finite in extent: there exists $\underline{c} < \infty$ such that $\mathcal{B} \subset \{\mathbf{p} : \|\mathbf{p}\| < \underline{c}\}$. A moment of reflection reveals that if $\mathbf{p}(0) \notin \mathcal{B}$, then the following event occurs with positive probability: at some time t_0 , party +1 (WLOG) loses control to -1, and does so at a policy with +1-ideal position $p(t_0) = p_{+1}^*$ and complexity $\|\mathbf{p}(t_0)\| \geq \underline{c}$. Conditioning on this event, fix $c(t_0) = \|\mathbf{p}(t_0)\|$ and $q(t_0) = p(t_0)$. From Lemma B.4, with positive probability, $\lim_{t \rightarrow \infty} c(t) = \infty$ and $c(t) \geq c(0)$ for all $t \geq t_0$. Consequently, applying Lemma B.2a twice: (i) with positive probability, $c(t) \geq \underline{c}$ for all $t \geq t_0$; so, (ii) with positive probability, $\lim_{t \rightarrow \infty} \|\mathbf{p}\|(t) = \infty$. In other words, $\kappa > 0$.

$\int v(q)dF(q) < 0$: Select sufficiently small ε and sufficiently large \underline{c} so that $\mathcal{B} \subset \{\mathbf{p} : \|\mathbf{p}\| < \underline{c}\}$ and so that $1 - \varepsilon < \frac{c - |p|}{c + |p|}$ for $p = p_{+1}^*$ and $p = p_{-1}^*$. Define \mathcal{B}' as in Lemma B.2c. We make the following two observations.

1. Lemmas B.4 and B.2b imply that if policy is above the complexity bound \underline{c} at some time t , $\|\mathbf{p}(t)\| > \underline{c}$, then it almost surely enters \mathcal{B}' , and does so via one of the ideals p_i^* .
2. Lemma B.2c implies that whenever policy enters \mathcal{B}' via one of the ideals p_i^* , then it almost surely exits \mathcal{B}' (either by entering the basin \mathcal{B} or by exceeding the complexity bound \underline{c}), and (with probability of at least $\nu > 0$, where ν is defined as in Lemma B.2c) does so by entering the basin \mathcal{B} .

Combining these two observations, we conclude that policy almost surely enters the basin eventually, and thus that $\kappa = 0$. ■

Lemma B.5. *Suppose $\mu > 0$. Consider the simulacrum process with $\varepsilon = 0$. Fix a start time $t_0 \geq 0$. For any \underline{c} ,*

$$\lim_{c(t_0) \rightarrow \infty} \Pr \left[\inf_{t \geq t_0} c(t) \geq \underline{c} \right] = 1.$$

Proof. Let $\tau_1^i \geq t_0$ be the first stopping time where Party i loses control at his ideal. We claim that for any $\nu \in (0, 1)$,

$$\lim_{c(t_0) \rightarrow \infty} \Pr \left[c(\tau_1^{+1}) \geq (1 - \nu) c(t_0) \right] = 1, \quad (22)$$

$$\lim_{c(t_0) \rightarrow \infty} \Pr \left[c(\tau_1^{-1}) \geq (1 - \nu) c(t_0) \right] = 1 \quad (23)$$

WLOG suppose that policy hits +1's ideal first, at time τ_1' ; note that $c(\tau_1') \geq c(t_0)$. Notice that, subsequent to τ_1' , Party i loses control with arrival rate λ_i ; so $c(\tau_1') - c(\tau_1^i)$ is exponentially distributed with parameter λ_i . Consequently, as $c(t_0) \rightarrow \infty$, the probability that $c(\tau_1') - c(\tau_1^i) \geq \nu c(t_0)$ vanishes. Our claim (22) follows immediately. The demonstration of the claim (23) is more involved, but proceeds similarly.

Condition on the event that $c(\tau_1^i) \geq (1 - \nu)c(t_0)$ for $i \in \{+1, -1\}$. As $c(t_0) \rightarrow \infty$, we have $(1 - \nu)c(t_0) - \underline{c} \rightarrow \infty$, so

$$\begin{aligned} \lim_{c(t_0) \rightarrow \infty} \Pr \left[\inf_{t \geq t_0} c(t) \geq \underline{c} \right] &= \lim_{c(t_0) \rightarrow \infty} \Pr \left[\inf_{i \in \{+1, -1\}; k \geq 1} c(\tau_k^i) \geq \underline{c} \right] \\ &\geq \lim_{c(t_0) \rightarrow \infty} \Pr \left[\inf_{i \in \{+1, -1\}; k \geq 1} c(\tau_k^i) \geq \underline{c} \right] = 1, \end{aligned}$$

where the last equality follows from Lemma B.3e.1. At the limit $c(t_0) \rightarrow \infty$, we conclude that (unconditionally) $\lim_{c(t_0) \rightarrow \infty} \Pr \left[\inf_{t \geq t_0} c(t) \geq \underline{c} \right] = 1$. ■

Proof of Proposition 4

$\|\mathbf{p}(0)\| \rightarrow 0$: Assume WLOG that Party +1 starts the game in control: $i(0) = +1$. We will argue that $\|\mathbf{p}(0)\| \rightarrow 0$, the distance of the starting policy $\mathbf{p}(0)$ from the basin \mathcal{B} vanishes. Note that the region where Party -1 deletes rules is bounded by the line $\frac{p}{\|\mathbf{p}\|} = 1 - \frac{2}{z_{-1}}$. While +1 remains in control, policy evolves along the line $(p, \|\mathbf{p}\|) = \gamma \cdot (t, t + \|\mathbf{p}(0)\|)$. The two aforementioned lines intersect where $p = \|\mathbf{p}(0)\| \frac{2}{z_{-1}-2}$. That is, if +1 remains in control for a time period longer than $\|\mathbf{p}(0)\| \gamma^{-1} \frac{2}{z_{-1}-2}$, then policy will enter the basin and eventually become perfectly simple. As $\|\mathbf{p}(0)\| \rightarrow 0$, the probability that this occurs converges to one.

$\|\mathbf{p}(0)\| \rightarrow \infty$: Consider the simulacrum process with $\varepsilon = 0$, and suppose that initial conditions are identical for the true and simulacrum process: $q(0) = 0$ and $c(0) = \|\mathbf{p}(0)\|$. The result then follows immediately from Lemma B.5 by choosing \underline{c} so that $\mathcal{B} \subset \{\mathbf{p} : \|\mathbf{p}\| < \underline{c}\}$. ■

Comparative Statics: The Politics of Kludges

Proof of Proposition 5a

From Proposition 3, the key object of interest is $\int v(q, i) dF(q)(q, i)$. We can rewrite, via some manipulations,

$$\begin{aligned} \int v(q, i) dF(q)(q, i) &= -(\Delta G(p_{+1}^*, +1) + \Delta G(p_{-1}^*, -1)) + \int_{p_{-1}^*}^{p_{+1}^*} (g(q, +1) + g(q, -1)) dq \\ &= \frac{-\gamma \left(\frac{1}{\lambda_{-1}} e^{\frac{\lambda_{+1}-\lambda_{-1}}{\gamma} p_{-1}^*} + \frac{1}{\lambda_{+1}} e^{\frac{\lambda_{+1}-\lambda_{-1}}{\gamma} p_{+1}^*} \right) + \int_{p_{-1}^*}^{p_{+1}^*} \left(2e^{\frac{\lambda_{+1}-\lambda_{-1}}{\gamma} q} \right) dq}{\gamma \left(\frac{1}{\lambda_{-1}} e^{\frac{\lambda_{+1}-\lambda_{-1}}{\gamma} p_{-1}^*} + \frac{1}{\lambda_{+1}} e^{\frac{\lambda_{+1}-\lambda_{-1}}{\gamma} p_{+1}^*} \right) + \int_{p_{-1}^*}^{p_{+1}^*} \left(2e^{\frac{\lambda_{+1}-\lambda_{-1}}{\gamma} q} \right) dq}. \end{aligned} \quad (24)$$

The denominator of the last expression (24) is positive; we may rewrite the numerator as

$$\frac{\gamma}{\lambda_{-1} - \lambda_{+1}} \left(e^{(p_{+1}^* - p_{-1}^*)(\lambda_{-1} - \lambda_{+1})} \left(3 - \frac{\lambda_{-1}}{\lambda_{+1}} \right) - \left(3 - \frac{\lambda_{+1}}{\lambda_{-1}} \right) \right),$$

so (24) has the same sign as

$$(p_{+1}^* - p_{-1}^*) - \frac{\log \frac{3 - \lambda_{+1}/\lambda_{-1}}{3 - \lambda_{-1}/\lambda_{+1}}}{\lambda_{-1} - \lambda_{+1}} = \Delta_p^* - \frac{\log \frac{3 - \Lambda^{-1}}{3 - \Lambda}}{\lambda(\sqrt{\Lambda} - \sqrt{\Lambda^{-1}})}.$$

The result then follows from Proposition 3. ■

Proof of Proposition 5b

Denote the parties' zealouslyness as $\mathbf{z} = (z_{+1}, z_{-1})$. We say that $\mathbf{z}' > \mathbf{z}$ if $z'_{+1} \geq z_{+1}$ and $z'_{-1} \geq z_{-1}$, with at least one strict inequality. Relabel the basin as $\mathcal{B}(\mathbf{z})$ to highlight its dependence on parties' zealouslyness. Our assumptions $z_{+1} > 2$ and $z_{-1} > 2$ ensure that $\mathcal{B}(\mathbf{z})$ is a compact set. Also, $\mathcal{B}(\mathbf{z})$ increases (strictly) in \mathbf{z} : if $\mathbf{z}' < \mathbf{z}$, then $\mathcal{B}(\mathbf{z}') \subset \mathcal{B}(\mathbf{z})$.

A history h is an infinite sequence of control durations $\{\Delta t_1^{+1}, \Delta t_1^{-1}, \Delta t_2^{+1}, \Delta t_2^{-1}, \dots\}$, whereas a k -truncated history h_k is characterized by the first $2k$ durations of control,

$$\{\Delta t_1^{+1}, \Delta t_1^{-1}, \dots, \Delta t_k^{+1}, \Delta t_k^{-1}\}.$$

Combined with the model's primitives, a history h determines the (equilibrium) path of policy for all time $t \geq 0$, whereas a truncated history h_k determines the path of policy up till time $t_k = \Delta t_1^{+1} + \dots + \Delta t_k^{-1}$. For $t \leq t_k$, we write $\mathbf{p}(t; h_k, \mathbf{z})$ to denote the time- t policy under truncated history h_k , given that parties have zealouslyness \mathbf{z} . Correspondingly, we write $\mathcal{P}(h_k; \mathbf{z}) = \cup_{t \leq t_k} \mathbf{p}(t; h_k, \mathbf{z})$ to denote the set of all policies attained under h_k up until (and including) time t_k .

Note that $\mathcal{P}(h_k, \mathbf{z})$ is compact. Suppose that $\mathcal{P}(h_k, \mathbf{z})$ does not intersect with the basin $\mathcal{B}(\mathbf{z})$; i.e., policy does not enter the basin at any time $t \leq t_k$. Then $\mathcal{P}(h_k, \mathbf{z})$ is 'uniformly continuous' in h_k , in the following sense. For any neighbourhood of $\mathcal{P}(h_k, \mathbf{z})$, there exists a neighbourhood of h_k (with respect to the usual topology on \mathbb{R}^k) such that for every k -truncated history h'_k in this neighbourhood, $\mathcal{P}(h'_k, \mathbf{z})$ lies within the aforementioned neighbourhood of $\mathcal{P}(h_k, \mathbf{z})$. Similarly, $\mathcal{P}(h_k, \mathbf{z})$ is 'pointwise continuous' in h_k , in the following specific sense: for any $l \leq k$, treating t_l as a function of h_k , $\mathbf{p}(t_l; h_k, \mathbf{z})$ is continuous in h_k .

A preliminary observation is that fixing a history h , if policy ever enters the basin $\mathcal{B}(\mathbf{z})$ given zealouslyness \mathbf{z} , then it enters the (larger) basin $\mathcal{B}(\mathbf{z}')$ given zealouslyness $\mathbf{z}' \leq \mathbf{z}$. Thus the probability that policy ever enters the basin is weakly decreasing, and κ is weakly increasing, in zealouslyness \mathbf{z} . It remains to show that κ is *strictly* increasing in \mathbf{z} .

Choose \mathbf{z} and \mathbf{z}' such that $\mathbf{z}' < \mathbf{z}$. Choose $\rho > 0$ and $\underline{c} > 0$ such that $\kappa \geq \rho$ for any regular starting policy with $\|\mathbf{p}\| \geq \underline{c}$. Choose $k \geq 2$ and a k -truncated history h_k with the following properties. First, $\mathcal{P}(h_k; \mathbf{z})$ does not intersect with $\mathcal{B}(\mathbf{z})$. Second, for some $l < k$, $\mathbf{p}(t_l; h_k, \mathbf{z})$ lies within the interior of $\mathcal{B}(\mathbf{z}')$. Third, at time t_k , complexity strictly exceeds \underline{c} : that is, $\|\mathbf{p}(t_k; h_k, \mathbf{z})\| > \underline{c}$.

By continuity of $\mathcal{P}(h_k, \mathbf{z})$ in h_k (both uniform and pointwise), we can construct a neighbourhood H_k of h_k such that these three properties also hold for any truncated history $h'_k \in H_k$. These properties, in turn, imply the following additional properties. (i) Given that parties have zealouslyness \mathbf{z} , conditional on h'_k , the probability κ of kludge is at least ρ . (ii) Given that parties have zealouslyness \mathbf{z}' , conditional on h'_k , policy enters the basin $\mathcal{B}(\mathbf{z}')$ and thus (almost surely) becomes perfectly simple.

Since H_k is a neighbourhood in the usual \mathbb{R}^k -topology, there is a strictly positive probability mass of truncated histories $h'_k \in H_k$. Coupled with properties (i) and (ii), it follows that there is a strictly positive probability mass of (untruncated) histories where policy becomes kludged given zealouslyness \mathbf{z}' , but does not become kludged given zealouslyness \mathbf{z} . In other words, κ is strictly increasing in \mathbf{z} . ■

Proof of Proposition 6

(i) Let F and f be the marginal steady-state distribution and density of $|q|$. Applying

Lemma B.1b: for all $0 \leq q \leq q' < p_{+1}^*$,

$$\frac{f(q')}{f(q)} = \frac{e^{\lambda \frac{\Lambda-1}{\gamma} q'} + e^{-\lambda \frac{\Lambda-1}{\gamma} q'}}{e^{\lambda \frac{\Lambda-1}{\gamma} q} + e^{-\lambda \frac{\Lambda-1}{\gamma} q}} \text{ and } \frac{\Delta F(\Delta_p^*)}{\lim_{q \rightarrow \Delta_p^*} f(q)} = \frac{\frac{\gamma}{\lambda} \left(\Lambda e^{\lambda \frac{\Lambda-1}{\gamma} q'} + \frac{1}{\Lambda} e^{-\lambda \frac{\Lambda-1}{\gamma} q'} \right)}{\Lambda e^{\lambda \frac{\Lambda-1}{\gamma} q} + \frac{1}{\Lambda} e^{-\lambda \frac{\Lambda-1}{\gamma} q}} \quad (25)$$

are both increasing in Λ . That is, F satisfies the monotone-likelihood ratio property in Λ . Thus, F increases in the sense of first-order stochastic-dominance as Λ increases.

(ii) This follows from the observation that the dynamics of q are independent of z_{+1}, z_{-1} .

(iii)–(v) If $p_{-1}^* = -p_{+1}^*$ and $\Lambda = 1$, then (25) simplifies further: for all $0 \leq q \leq q' < p_{+1}^*$,

$$\frac{f(q')}{f(q)} = 1 \text{ and } \frac{\Delta F(p_{+1}^*)}{\lim_{q \rightarrow p_{+1}^*} f(q)} = \frac{\gamma}{\lambda}, \text{ so that} \quad (26)$$

$$F(q) = \begin{cases} \frac{q}{\Delta_p^* + \frac{\gamma}{\lambda}} & : p < \Delta_p^* \\ 1 & : p = \Delta_p^* \end{cases}. \quad (27)$$

By inspection, F increases in the sense of first-order stochastic-dominance as γ increases, as λ decreases, and as Δ_p^* increases. ■

Strategic Extremism

For this Appendix, we say that an equilibrium is Markov Perfect if the evolution of position $\frac{d}{dt}p(t)$ depends only on the payoff-relevant state variables $(p(t), i(t))$. In particular, equilibria in focused strategies are Markov Perfect.

Lemma B.6a. *If a focused strategy profile with targets $(p_{+1}^{**}, p_{-1}^{**})$ is a Markov Perfect Equilibrium, then $p_{+1}^{**} \geq p_{+1}^*$ and $p_{-1}^{**} \leq p_{-1}^*$.*

Proof. Let $a_i(p)$ be the rate at which party i shifts policy position when he is in power and the current policy position is p . A Markov strategy profile is described by two functions a_{+1} and a_{-1} . Let $V_{ij}(p_0)$ be party i 's expected payoff when party j is in power and the current policy bias is p_0 . Let T_j be the first time when j loses power to the other party. Then

$$V_{ij}(p) = E \left[- \int_0^{T_j} e^{-r_t} |g_j(t, p) - p_i^*| dt + e^{-r_{T_j}} V_{i,-j}(g_j(T_j, p)) \right],$$

where $g_j(t, p)$ evolves according to the law of motion

$$\frac{dg_j}{dt}(t, p) = a_j(g_j(t, p)),$$

with initial condition $g_j(0, p) = p$. The expectation in the expression of V_{ij} is taken over T_j . For notational simplicity, the dependence of V and g on a has been suppressed. Substituting in the probability density of T_j and performing a change of order of integral yields that

$$V_{ij}(p) = \int_0^\infty [-|g_j(t, p) - p_i^*| + \lambda_j V_{i,-j}(g_j(t, p))] e^{-(r_i + \lambda_j)t} dt, \text{ for every } p_0 \in \mathbb{R}. \quad (28)$$

The Bellman equation associated with this integral is¹⁸

$$-|p_0 - p_i^*| + \lambda_j V_{i,-j}(p_0) - (r_i + \lambda_j) V_{ij}(p_0) + V'_{ij}(p_0) a_j(p_0) = 0, \text{ for every } p_0 \in \mathbb{R}. \quad (29)$$

By the standard theory of optimal control, the optimal control satisfies the conditions that $a_i(p) = \gamma$ if $V'_{ii}(p) > 0$ and $a_i(p) = -\gamma$ if $V'_{ii}(p) < 0$. Now consider the special case where (a_{+1}, a_{-1}) is a focused strategy with targets $(p_{+1}^{**}, p_{-1}^{**})$ and is an MPE. Then $a_{+1}(p) = \gamma$ when $p < p_{+1}^{**}$ and $a_{-1}(p) = -\gamma$ when $p > p_{-1}^{**}$. Therefore, Eq. (29) implies that

$$\gamma V'_{i,+1}(p) = |p - p_i^*| - \lambda_{+1} V_{i,-1}(p) + (r_i + \lambda_{+1}) V_{i,+1}(p), \text{ for } p < p_{+1}^{**}; \quad (30)$$

$$\gamma V'_{i,-1}(p) = -|p - p_i^*| + \lambda_{-1} V_{i,-1}(p) - (r_i + \lambda_{-1}) V_{i,-1}(p), \text{ for } p > p_{-1}^{**}; \quad (31)$$

$$0 = |p_j^{**} - p_i^*| + \lambda_j V_{i,-j}(p_j^{**}) - (r_i + \lambda_j) V_{ij}(p_j^{**}). \quad (32)$$

In equilibrium, $V'_{ii}(p) a_i(p) \geq 0$ for every p . Therefore, Eq. (29) implies that

$$|p - p_i^*| - \lambda_i V_{i,-i}(p) + (r_i + \lambda_i) V_{ii}(p) = V'_{ii}(p) a_i(p) \geq 0 \text{ for every } p \in \mathbb{R}. \quad (33)$$

When $p = p_i^{**}$, the left hand side vanishes as $a_i(p_i^{**}) = 0$. Therefore, p_i^{**} is a global minimum of the left hand side (as a function of p). Moreover, $V'_{ii}(p_i^{**}) = 0$. (If $V'_{ii}(p_i^{**}) > 0$, then $a_i(p_i^{**})$ should be γ ; assuming that $V'_{ii}(p_i^{**}) < 0$ leads to a similar contradiction.) Differentiating the left hand side of Eq. (33) at p_i^{**} yields that

$$V'_{i,-i}(p_i^{**}) \begin{cases} = -\lambda_i^{-1}, & \text{if } p_i^{**} < p_i^*; \\ \in [-\lambda_i^{-1}, \lambda_i^{-1}], & \text{if } p_i^{**} = p_i^*; \\ = \lambda_i^{-1}, & \text{if } p_i^{**} > p_i^*. \end{cases} \quad (34)$$

Suppose that $p_{+1}^{**} < p_{+1}^*$. Then $g_{-1}(t, p) = \max\{p - \gamma t, p_{-1}^{**}\} < p_{+1}^*$ when $p < p_{+1}^{**}$. Therefore,

$$V_{+1,-1}(p) = \int_0^\infty [-(p_{+1}^* - g_{-1}(t, p)) + \lambda_{-1} V_{+1,+1}(g_{-1}(t, p))] e^{-(r_{+1} + \lambda_{-1})t} dt, \text{ for } p \in (p_{-1}^{**}, p_{+1}^{**}).$$

By assumption, $V'_{+1,+1}(p) \geq 0$ for every $p < p_{+1}^{**}$. Therefore, the terms in the bracket are increasing in $g_{-1}(t, p)$. Since $g_{-1}(t, p) = \max\{p - \gamma t, p_{-1}^{**}\}$, $g_{-1}(t, p)$ is non-decreasing in p . Therefore, $V_{+1,-1}(p)$ is non-decreasing in p , contradicting the result that $V'_{+1,-1}(p_{+1}^{**}) = -\lambda_{+1}^{-1}$. The assumption that $p_{-1}^{**} > p_{-1}^*$ leads to a similar contradiction. ■

It will be shown in Lemma B.6c that a focused strategy profile with targets $(p_{+1}^{**}, p_{-1}^{**})$ forms a Markov Perfect Equilibrium if and only if $p_i^{**} = BR_i(p_{-i}^{**})$ where the best response functions BR_{+1} and BR_{-1} will be defined from the functions H_{+1} and H_{-1} to be introduced

¹⁸Formally, the Bellman equation can be derived as follows: replace p in Eq. (28) with $g_j(s, p)$ and $g_j(t, p)$ with $g_j(s+t, p)$ and rewrite Eq. (28) as $V'_{ij}(g_j(s, p)) e^{-(r_i + \lambda_j)s} = \int_s^\infty [-|g_j(\tau, p) - p_i^*| + \lambda_j V_{i,-j}(g_j(\tau, p))] e^{-(r_i + \lambda_j)\tau} d\tau$ where $\tau = s + t$. Differentiating both sides with respect to s at $s = 0$ yields the Bellman equation.

shortly. Define

$$A_i = \begin{pmatrix} r_i + \lambda_{+1} & -\lambda_{+1} \\ \lambda_{-1} & -(r_i + \lambda_{-1}) \end{pmatrix}, \text{ for } i \in \{-1, +1\}; \quad (35)$$

$$L_i(p) = \int_0^p \gamma^{-1} |\tilde{p} - p_i^*| e^{-\gamma^{-1} \tilde{p} A_i} \begin{pmatrix} 1 \\ -1 \end{pmatrix} d\tilde{p}, \text{ for } i \in \{-1, +1\}; \quad (36)$$

$$\mathbf{1}_{+1} = \begin{pmatrix} 1 \\ 0 \end{pmatrix}; \quad (37)$$

$$\mathbf{1}_{-1} = \begin{pmatrix} 0 \\ 1 \end{pmatrix}; \quad (38)$$

$$H_{+1}(p, p', \eta) = \mathbf{1}_{-1}^\top e^{\gamma^{-1}(p-p')A_{+1}} \begin{pmatrix} -|p' - p_{+1}^*| \\ |p' - p_{+1}^*| + \gamma \lambda_{+1}^{-1} \eta \end{pmatrix} + \mathbf{1}_{-1}^\top e^{\gamma^{-1} p A_{+1}} A_{+1} [L_{+1}(p) - L_{+1}(p')] - |p - p_{+1}^*|; \quad (39)$$

$$H_{-1}(p, p', \eta) = \mathbf{1}_{+1}^\top e^{\gamma^{-1}(p-p')A_{-1}} \begin{pmatrix} -|p' - p_{-1}^*| - \gamma \lambda_{-1}^{-1} \eta \\ |p' - p_{-1}^*| \end{pmatrix} + \mathbf{1}_{+1}^\top e^{\gamma^{-1} p A_{-1}} A_{-1} [L_{-1}(p) - L_{-1}(p')] + |p - p_{-1}^*|. \quad (40)$$

In the last two equations, $\mathbf{1}_i^\top$ denotes the transpose of $\mathbf{1}_i$.

Lemma B.6b. *For every $p \leq p_{-1}^*$, $H_{+1}(p, p_{+1}^*, 0) < 0$ and $H_{+1}(p, p', 1)$ is strictly increasing in p' for $p' \geq p_{+1}^*$. For every $p \leq p_{-1}^*$ and $p' \geq p_{+1}^*$, $H_{+1}(p, p', \eta)$ is strictly increasing in η . Finally, $H_{+1}(p, p', 1) \rightarrow \infty$ as $p' \rightarrow \infty$. Similarly, for every $p \geq p_{+1}^*$, $H_{-1}(p, p_{-1}^*, 0) > 0$ and $H_{-1}(p, p', 1)$ is strictly increasing in p' for $p' \leq p_{-1}^*$. For every $p \geq p_{+1}^*$ and $p' \leq p_{-1}^*$, $H_{-1}(p, p', \eta)$ is strictly decreasing in η . Finally, $H_{-1}(p, p', 1) \rightarrow -\infty$ as $p' \rightarrow -\infty$.*

Proof. First perform the eigenvalue decomposition of A_i :

$$A_i = \frac{1}{\lambda_{-1}(\mu_{i+} - \mu_{i-})} \begin{pmatrix} \mu_{i+} + r_i + \lambda_{-1} & \mu_{i-} + r_i + \lambda_{-1} \\ \lambda_{-1} & \lambda_{-1} \end{pmatrix} \begin{pmatrix} \mu_{i+} & \\ & \mu_{i-} \end{pmatrix} \begin{pmatrix} \lambda_{-1} & -\mu_{i-} - r_i - \lambda_{-1} \\ -\lambda_{-1} & \mu_{i+} + r_i + \lambda_{-1} \end{pmatrix},$$

where

$$\mu_{i\pm} = \frac{1}{2} \left[\lambda_{+1} - \lambda_{-1} \pm \sqrt{(\lambda_{-1} - \lambda_{+1})^2 + 4r_i^2 + 4(\lambda_{-1} + \lambda_{+1})r_i} \right] \quad (41)$$

are the eigenvalues of A_i . Note that $\mu_{i+} > 0 > \mu_{i-}$ for $i \in \{+1, -1\}$. To avoid confusion, the eigenvalue μ_{i+} when $i = +1$ will be written as μ_{++} and the same rule applies to the other three eigenvalues as well as $\xi_{i\pm}$ and $\zeta_{i\pm}$ to be introduced below. Using this decomposition, H_i can be rewritten as

$$\begin{aligned} H_{+1}(p, p', \eta) &= (\mu_{++} - \mu_{+-})^{-1} \left[(r_{+1} + 2\lambda_{-1} + \mu_{+-}) \xi_{++}(p, p') - (r_{+1} + 2\lambda_{-1} + \mu_{++}) \xi_{+-}(p, p') \right] + \\ &\quad + (\mu_{++} - \mu_{+-})^{-1} \left[(r_{+1} + \lambda_{-1} + \mu_{++}) e^{\gamma^{-1}(p-p')\mu_{+-}} - (r_{+1} + \lambda_{-1} + \mu_{+-}) e^{\gamma^{-1}(p-p')\mu_{++}} \right] \gamma \lambda_{+1}^{-1} \eta; \\ H_{-1}(p, p', \eta) &= (\mu_{-+} - \mu_{--})^{-1} \left[(r_{-1} + 2\lambda_{+1} - \mu_{--}) \xi_{-+}(p, p') - (r_{-1} + 2\lambda_{+1} - \mu_{-+}) \xi_{--}(p, p') \right] + \\ &\quad + (\mu_{-+} - \mu_{--})^{-1} \left[(r_{-1} + \lambda_{+1} - \mu_{-+}) e^{\gamma^{-1}(p-p')\mu_{--}} - (r_{-1} + \lambda_{+1} - \mu_{--}) e^{\gamma^{-1}(p-p')\mu_{-+}} \right] \gamma \lambda_{-1}^{-1} \eta, \end{aligned}$$

where

$$\xi_{i\pm}(p, p') = \gamma^{-1} \mu_{i\pm} \int_{p'}^p e^{\gamma^{-1}(p-\tilde{p})\mu_{i\pm}} |\tilde{p} - p_i^*| d\tilde{p} + |p - p_i^*| - e^{\gamma^{-1}(p-p')\mu_{i\pm}} |p' - p_i^*|.$$

Splitting the first integral at p_i^* and integrating by parts yields that

$$\xi_{+\pm}(p, p') = \gamma \mu_{+1, \pm}^{-1} \left[1 + e^{\gamma^{-1}(p-p')\mu_{++}} - 2e^{\gamma^{-1}(p-p_{+1}^*)\mu_{++}} \right]; \quad (42)$$

$$\xi_{-\pm}(p, p') = -\gamma \mu_{-1, \pm}^{-1} \left[1 + e^{\gamma^{-1}(p-p')\mu_{+-}} - 2e^{\gamma^{-1}(p-p_{-1}^*)\mu_{+-}} \right]. \quad (43)$$

When $p \leq p_{-1}^*$ and $p' \geq p_{+1}^*$, $e^{\gamma^{-1}(p-p')\mu_{+-}} > e^{\gamma^{-1}(p-p')\mu_{++}}$, and $|r_{+1} + \lambda_{-1} + \mu_{++}| > |r_{+1} + \lambda_{-1} + \mu_{+-}|$, so

$$\frac{\partial H_{+1}}{\partial \eta}(p, p', \eta) = (\mu_{++} - \mu_{+-})^{-1} \left[(r_{+1} + \lambda_{-1} + \mu_{++})e^{\gamma^{-1}(p-p')\mu_{+-}} - (r_{+1} + \lambda_{-1} + \mu_{+-})e^{\gamma^{-1}(p-p')\mu_{++}} \right] \gamma \lambda_{+1}^{-1} > 0.$$

Therefore, $H_{+1}(p, p', \eta)$ is strictly increasing in η . A symmetric argument implies that $H_{-1}(p, p', \eta)$ is strictly decreasing in η when $p \geq p_{+1}^*$ and $p' \leq p_{-1}^*$.

In what follows, fix a $p \leq p_{-1}^*$. Then

$$\xi_{+\pm}(p, p_{+1}^*) = \int_p^{p_{+1}^*} e^{\gamma^{-1}(p-\bar{p})\mu_{+1\pm}} d\bar{p}.$$

Therefore, $0 < \xi_{++}(p, p_{+1}^*) < \xi_{+-}(p, p_{+1}^*)$, and thus

$$H_{+1}(p, p_{+1}^*, 0) = -(\mu_{++} - \mu_{+-})^{-1} (r_{+1} + 2\lambda_{-1} + \mu_{++}) [\xi_{+-}(p, p_{+1}^*) - \xi_{++}(p, p_{+1}^*)] - \xi_{++}(p, p_{+1}^*) < 0.$$

Moreover, as $p' \rightarrow \infty$, $\xi_{++}(p, p')$ remains bounded while $\xi_{+-}(p, p') \rightarrow \infty$. It follows immediately that

$$\lim_{p' \rightarrow \infty} H_{+1}(p, p', 1) = \infty.$$

Taking derivative with respect to p' on both sides of Eq. (42) yields that

$$\frac{\partial \xi_{+1\pm}}{\partial p'}(p, p') = -e^{-\gamma^{-1}(p'-p)\mu_{+1\pm}}.$$

Therefore, for $p \leq p_{-1}^*$ and $p' \geq p_{+1}^*$,

$$H_{+1,2}(p, p', 1) = (\mu_{++} - \mu_{+-})^{-1} \left(\zeta_{+-} e^{-\gamma^{-1}(p'-p)\mu_{+-}} - \zeta_{++} e^{-\gamma^{-1}(p'-p)\mu_{++}} \right), \quad (44)$$

where $H_{+1,2}$ denotes the partial derivative of H_{+1} with respect to its second argument, and

$$\begin{aligned} \zeta_{++} &= (r_{+1} + 2\lambda_{-1} + \mu_{+-}) - (r_{+1} + \lambda_{-1} + \mu_{+-})\lambda_{+1}^{-1}\mu_{++}; \\ \zeta_{+-} &= (r_{+1} + 2\lambda_{-1} + \mu_{++}) - (r_{+1} + \lambda_{-1} + \mu_{++})\lambda_{+1}^{-1}\mu_{+-}. \end{aligned}$$

Now $\zeta_{+-} > 0$ and $e^{-\gamma^{-1}(p'-p)\mu_{+-}} > e^{-\gamma^{-1}(p'-p)\mu_{++}}$ when $p' \geq p_{+1}^*$. Moreover,

$$\zeta_{+-} - \zeta_{++} = (\mu_{++} - \mu_{+-})[1 + \lambda_{+1}^{-1}(r_{+1} + \lambda_{-1})] > 0.$$

Therefore, $\frac{\partial H_{+1}}{\partial p'}(p, p', 1) > 0$ and thus $H_{+1}(p, p', 1)$ is strictly increasing in p' for $p' \geq p_{+1}^*$ and $p \leq p_{-1}^*$.

All the assertions about H_{-1} can be proved with a symmetric argument. ■

Fix a $p \leq p_{-1}^*$. If $H_{+1}(p, p_{+1}^*, 1) \geq 0$, then there exists a unique $\eta_{+1} \in (0, 1]$ such that $H_{+1}(p, p_{+1}^*, \eta_{+1}) = 0$. In this case, define $BR_{+1}(p) = p_{+1}^*$. If $H_{+1}(p, p_{+1}^*, 1) < 0$, then there exists a unique $p' \in (p_{+1}^*, \infty)$ such that $H_{+1}(p, p', 1) = 0$. Define $BR_{+1}(p) = p'$ in this case. Define BR_{-1} in a similar fashion.

Lemma B.6c. *The focused strategy with targets $(p_{+1}^{**}, p_{-1}^{**})$ is a Markov Perfect Equilibrium of the one-dimensional game if and only if $p_i^{**} = BR_i(p_{-i}^{**})$ for $i \in \{-1, +1\}$.*

Proof. Let

$$\vec{V}_i(p) = \begin{pmatrix} V_{i+1}(p) \\ V_{i-1}(p) \end{pmatrix}.$$

Then Eqs. (30) and (31) can be rewritten as

$$\vec{V}_i'(p) = \gamma^{-1} A_i \vec{V}_i(p) + \gamma^{-1} |p - p_i| \begin{pmatrix} 1 \\ -1 \end{pmatrix}, \text{ for every } p \in (p_{-1}^{**}, p_{+1}^{**}),$$

where A_i is defined in Eq. (35). If $\vec{V}_i(p')$ is known, the solution to this differential equation is

$$\vec{V}_i(p) = e^{\gamma^{-1}(p-p')A_i} \vec{V}_i(p') + e^{\gamma^{-1}pA_i} [L_i(p) - L_i(p')], \text{ for every } p \in [p_{-1}^{**}, p_{+1}^{**}], \quad (45)$$

where L_i is defined in Eq. (36). Eqs. (32) (for the case where $j = i$) and (34) can be rewritten as

$$\begin{aligned} A_{+1} \vec{V}_{+1}(p_{+1}^{**}) &= \begin{pmatrix} -|p_{+1}^{**} - p_{+1}^*| \\ |p_{+1}^{**} - p_{+1}^*| + \gamma \lambda_{+1}^{-1} \eta_{+1} \end{pmatrix}; \\ A_{-1} \vec{V}_{-1}(p_{-1}^{**}) &= \begin{pmatrix} -|p_{-1}^{**} - p_{-1}^*| - \gamma \lambda_{-1}^{-1} \eta_{-1} \\ |p_{-1}^{**} - p_{-1}^*| \end{pmatrix}, \end{aligned}$$

where $\eta_i \in [-1, 1]$ if $p_i^{**} = p_i^*$ and $\eta_i = 1$ if $p_i^{**} \neq p_i^*$. Substituting these two equations into Eq. (45) yields that

$$\begin{aligned} A_{+1} \vec{V}_{+1}(p_{-1}^{**}) &= e^{\gamma^{-1}(p_{-1}^{**}-p_{+1}^{**})A_{+1}} \begin{pmatrix} -|p_{+1}^{**} - p_{+1}^*| \\ |p_{+1}^{**} - p_{+1}^*| + \gamma \lambda_{+1}^{-1} \eta_{+1} \end{pmatrix} + e^{\gamma^{-1}p_{-1}^{**}A_{+1}} A_{+1} [L_{+1}(p_{-1}^{**}) - L_{+1}(p_{+1}^{**})]; \\ A_{-1} \vec{V}_{-1}(p_{+1}^{**}) &= e^{\gamma^{-1}(p_{+1}^{**}-p_{-1}^{**})A_{-1}} \begin{pmatrix} -|p_{-1}^{**} - p_{-1}^*| - \gamma \lambda_{-1}^{-1} \eta_{-1} \\ |p_{-1}^{**} - p_{-1}^*| \end{pmatrix} + e^{\gamma^{-1}p_{+1}^{**}A_{-1}} A_{-1} [L_{-1}(p_{+1}^{**}) - L_{-1}(p_{-1}^{**})] \end{aligned}$$

Now the remaining boundary condition Eq. (32) for the case where $i \neq j$ can be rewritten as $H_i(p_{-i}^{**}, p_i^{**}, \eta_i) = 0$ for $i \in \{-1, +1\}$ where H_i is defined in Eqs. (39) and (40). This proves the ‘‘only if’’ assertion of the lemma. Conversely, if $(p_{+1}^{**}, p_{-1}^{**})$ satisfies the system that $H_i(p_i^{**}, p_{-i}^{**}, \eta_i) = 0$ with $\eta_i \in [-1, 1]$ when $p_i^{**} = p_i^*$ and $\eta_i = 1$ when $p_i^{**} \neq p_i^*$, then Eqs. (30)-(34) will be satisfied, implying that the focused strategy profile is an MPE. ■

Lemma B.6d. $\lim_{P \rightarrow \infty} H_{+1}(-P, P, 1) = \infty$, and $\lim_{P \rightarrow \infty} H_{-1}(P, -P, 1) = -\infty$.

Proof. By Eq. (42), as $P \rightarrow \infty$,

$$\begin{aligned} \xi_{++}(-P, P) &\rightarrow \gamma \mu_{++}^{-1}; \\ \xi_{+-}(-P, P) &\sim \gamma \mu_{+-}^{-1} e^{-2\gamma^{-1}P\mu_{+-}}. \end{aligned}$$

Therefore,

$$H_{+1}(-P, P, 1) \sim (\mu_{++} - \mu_{+-})^{-1} [\gamma \lambda_{+1}^{-1} (r_{+1} + \lambda_{-1} + \mu_{++}) - \gamma \mu_{+-}^{-1} (r_{+1} + 2\lambda_{-1} + \mu_{++})] e^{-2\gamma^{-1} P \mu_{+-}}.$$

Clearly, the right hand side approaches ∞ as $P \rightarrow \infty$. A symmetric argument implies that $H_{-1}(P, -P, 1) \rightarrow -\infty$ as $P \rightarrow \infty$. \blacksquare

Lemma B.6e. *There exists a $p_{-1,c} \leq p_{-1}^*$ such that $BR_{+1}(p) = p_{+1}^*$ if and only if $p_{-1,c} \leq p \leq p_{-1}^*$, and $BR_{+1}'(p) < 0$ when $p < p_{-1,c}$. Similarly, there exists a $p_{+1,c} \geq p_{+1}^*$ such that $BR_{-1}(p) = p_{-1}^*$ if and only if $p_{+1}^* \leq p \leq p_{+1,c}$ and $BR_{-1}'(p) < 0$ when $p > p_{+1,c}$.*

Proof. We only prove the assertion on BR_{+1} , as the assertion on BR_{-1} follows from a symmetric argument. For every $p \leq p_{-1}^*$, define $\eta_{+1}(p) = 1$ if $BR_{+1}(p) > p_{+1}^*$ and $\eta_{+1}(p)$ be the unique η such that $H_{+1}(p, p_{+1}^*, \eta) = 0$ when $BR_{+1}(p) = p_{+1}^*$. Then

$$H_{+1}(p, BR_{+1}(p), \eta_{+1}(p)) = 0, \text{ for every } p \leq p_{-1}^*. \quad (46)$$

For every $\tilde{p} \in \mathbb{R}$, $p' \geq p_{+1}^*$ and $\eta \in [-1, 1]$, define

$$\tilde{U}_{+1}(\tilde{p}; p', \eta) = e^{\gamma^{-1}(\tilde{p}-p')A_{+1}} A_{+1}^{-1} \begin{pmatrix} -|p' - p_{+1}^*| \\ |p' - p_{+1}^*| + \gamma \lambda_{+1}^{-1} \eta \end{pmatrix} + e^{\gamma^{-1} \tilde{p} A_{+1}} [L_{+1}(\tilde{p}) - L_{+1}(p')].$$

Then

$$\tilde{U}_{+1}'(\tilde{p}; p', \eta) = \gamma^{-1} A_{+1} \tilde{U}_{+1}(\tilde{p}; p', \eta) + \gamma^{-1} |\tilde{p} - p_{+1}^*| \begin{pmatrix} 1 \\ -1 \end{pmatrix}, \text{ for every } \tilde{p} \in \mathbb{R}. \quad (47)$$

$$H_{+1}(p, p', \eta) = 1_{-1}^\top A_{+1} \tilde{U}_{+1}(p; p', \eta) - |p - p_{+1}^*|. \quad (48)$$

In the first equation, \tilde{U}_{+1}' is the derivative of \tilde{U}_{+1} with respect to its first argument (\tilde{p} in that equation). Therefore, for every $p < p_{+1}^*$ the partial derivative of H_{+1} with respect to its first argument is

$$\begin{aligned} H_{+1,1}(p, p', \eta) &= 1_{-1}^\top \gamma^{-1} A_{+1} \left[A_{+1} \tilde{U}_{+1}(p; p', \eta) + |p - p_{+1}^*| \begin{pmatrix} 1 \\ -1 \end{pmatrix} \right] + 1 \\ &= 1 + \lambda_{-1} 1_{+1}^\top \left[A_{+1} \tilde{U}_{+1}(p; p', \eta) + |p - p_{+1}^*| \begin{pmatrix} 1 \\ -1 \end{pmatrix} \right] - (r_{+1} + \lambda_{-1}) H_{+1}(p, p', \eta). \end{aligned}$$

Combining this result with Eq. (46) yields that

$$H_{+1,1}(p, BR_{+1}(p), \eta_{+1}(p)) = 1 + \lambda_{-1} \gamma^{-1} 1_{+1}^\top \left[A_{+1} \tilde{U}_{+1}(p; p', \eta) + |p - p_{+1}^*| \begin{pmatrix} 1 \\ -1 \end{pmatrix} \right]. \quad (49)$$

On the other hand, party +1's value function when party -1 has target p and party +1's target is $BR_{+1}(p)$ satisfies the same differential equation (Bellman equation) Eq. (47) and the same boundary conditions at p and $BR_{+1}(p)$. Therefore, on $[p, BR_{+1}(p)]$, $\tilde{U}_{+1}(\cdot; BR_{+1}(p), \eta_{+1}(p))$ coincides with party +1's value function. Because party +1's flow payoff is always non-positive,

$$U_{+1,+1}(p; BR_{+1}(p), \eta_{+1}(p)) \leq 0. \quad (50)$$

Furthermore, party +1 has the option to stay at p when receiving control at policy position p , by doing which he receives expected payoff $-\frac{1}{r_{+1}}|p - p_{+1}^*|$. (If party +1 does this, then both parties will remain stationary at p and the policy position will remain at p forever.) Therefore,

$$U_{+1,+1}(p; BR_{+1}(p), \eta_{+1}(p)) \geq -\frac{1}{r_{+1}}|p - p_{+1}^*|. \quad (51)$$

By Eq. (48), that $H_{+1}(p, BR_{+1}(p), \eta_{+1}(p)) = 0$ implies that

$$\lambda_{-1}U_{+1,+1}(p; BR_{+1}(p), \eta_{+1}(p)) - (r_{+1} + \lambda_{-1})U_{+1,-1}(p; BR_{+1}(p), \eta_{+1}(p)) - |p - p_{+1}^*| = 0. \quad (52)$$

Using Eq. (52) to eliminate $U_{+1,-1}(p; BR_{+1}(p), \eta_{+1}(p))$ from Eq. (49) yields that

$$H_{+1,1}(p, BR_{+1}(p), \eta_{+1}(p)) = 1 + \frac{\lambda_{-1}(r_{+1} + \lambda_{+1} + \lambda_{-1})}{\gamma(r_{+1} + \lambda_{-1})} [|p - p_{+1}^*| + r_{+1}U_{+1,+1}(p; BR_{+1}(p), \eta_{+1}(p))].$$

Combining this with Eqs. (50) and (51) yields that

$$1 \leq H_{+1,1}(p, BR_{+1}(p), \eta_{+1}(p)) \leq 1 + \frac{\lambda_{-1}(r_{+1} + \lambda_{-1} + \lambda_{+1})}{\gamma(r_{+1} + \lambda_{-1})} |p - p_{+1}^*|. \quad (53)$$

In particular, $H_{+1,1}(p, BR_{+1}(p), \eta_{+1}(p)) > 0$ for every $p \leq p_{-1}^*$. Lemma B.6b implies that the partial derivative of H_{+1} with respect to its third argument (η) is positive. By the Implicit Function Theorem,

$$\eta'_{+1}(p) = -\frac{H_{+1,1}(p, p_{+1}^*, \eta_{+1}(p))}{H_{+1,3}(p, p_{+1}^*, \eta_{+1}(p))} < 0, \text{ for every } p \text{ such that } BR_{+1}(p) = p_{+1}^*.$$

Therefore, $\eta_{+1}(p)$ is decreasing in p , and as long as $BR_{+1}(p) = p_{+1}^*$, $\eta_{+1}(\tilde{p})$ will remain below unit for every $\tilde{p} \geq p$. This proves the existence of $p_{-1,c}$.

Lemma B.6b also implies that $H_{+1,2}(p, BR_{+1}(p), \eta_{+1}(p)) > 0$. By the Implicit Function Theorem,

$$BR'_{+1}(p) = -\frac{H_{+1,1}(p, BR_{+1}(p), 1)}{H_{+1,2}(p, BR_{+1}(p), 1)} < 0, \text{ for every } p < p_{-1,c}. \quad \blacksquare$$

Lemma B.6f. *A Markov Perfect equilibrium in focused strategies exists. In any such equilibrium, targets are weakly extreme: $p_{+1}^{**} \geq p_{+1}^*$ and $p_{-1}^{**} \leq p_{-1}^*$.*

Proof. By Lemma B.6d, there exists a $P > \max\{|p_{+1}^*|, |p_{-1}^*|\}$ such that $H_{+1}(-P, P, 1) > 0$ and $H_{-1}(P, -P, 1) < 0$. Therefore, $BR_{+1}(-P) < P$ and $BR_{-1}(P) > -P$. By Lemma 7.4f, $BR_{+1}(p) < P$ for every $p \in [-P, p_{-1}^*]$ and $BR_{-1}(p) > -P$ for every $p \in [p_{+1}^*, P]$. Therefore, the map

$$BR(p_{+1}, p_{-1}) = (BR_{+1}(p_{-1}), BR_{-1}(p_{+1}))$$

is a continuous map of from $[p_{+1}^*, P] \times [-P, p_{-1}^*]$ into itself. The existence of equilibrium follows from Brouwer's fixed-point theorem. \blacksquare

Lemma B.6g. *There exists a $\gamma_1 > 0$ such that when $\gamma \leq \gamma_1$, $BR'_{+1}(p) > -1$ for every $p < p_{-1,c}$ and $BR'_{-1}(p) > -1$ for every $p > p_{+1,c}$.*

Proof. By Eq. (44), as $\gamma \rightarrow 0$,

$$H_{+1,2}(p, p', 1) \sim (\mu_{++} - \mu_{+-})^{-1} \zeta_{+-} e^{-\gamma^{-1}(p'-p)\mu_{+-}}.$$

Combining this result with Eq. (53) yields that when $p < p_{-1,c}$,

$$BR'_{+1}(p) = -\frac{H_{+1,1}(p, BR_{+1}(p), 1)}{H_{+1,2}(p, BR_{+1}(p), 1)} \geq -M\gamma^{-1}|p - p_{+1}^*|e^{\gamma^{-1}(p_{+1}^*-p)\mu_{+-}},$$

for some constant $M > 0$. (We have used the fact that $BR_{+1}(p) \geq p_{+1}^*$. The right hand side is strictly increasing in $|p - p_{+1}^*|$ when $|p - p_{+1}^*| > -\gamma\mu_{+-}^{-1}$. Therefore, when $\gamma < -\mu_{+-}(p_{+1}^* - p_{-1}^*)$,

$$BR'_{+1}(p) \geq -M\gamma^{-1}(p_{+1}^* - p_{-1}^*)e^{\gamma^{-1}(p_{+1}^*-p_{-1}^*)\mu_{+-}}, \text{ for every } p < p_{-1,c}.$$

The limit of the right hand side as $\gamma \rightarrow 0$ is zero, so $BR'_{+1}(p) > -1$ for every $p < p_{-1,c}$ when γ is below some threshold. A symmetric argument proves the assertion on BR_{-1} . \blacksquare

The following lemma is concerned with the dependence of $H_i(p, p', \eta)$ on r_i . To make the dependence explicit, the function will be written as $H_i(p, p', \eta; r_i)$ in the lemma and its proof.

Lemma B.6h. *Assume that $\lambda_{+1} \neq \lambda_{-1}$. There exists a $\gamma_2 > 0$ such that when $\gamma \leq \gamma_2$, the following hold:*

1. *There exists a $r_{+1,c} < \infty$ such that $H_{+1}(p, p_{+1}^*, 1; r_{+1}) > 0$ for every $p \leq p_{-1}^*$ and $r_{+1} > r_{+1,c}$; if $H_{+1}(p, p_{+1}^*, 1; r_{+1}) = 0$ for some $p \leq p_{-1}^*$ and $r_{+1} \leq r_{+1,c}$, then $\frac{\partial H_{+1}}{\partial r_{+1}}(p, p_{+1}^*, 1; r_{+1}) > 0$.*
2. *There exists a $r_{-1,c} < \infty$ such that $H_{-1}(p, p_{-1}^*, 1; r_{-1}) < 0$ for every $p \geq p_{+1}^*$ and $r_{-1} > r_{-1,c}$; if $H_{-1}(p, p_{-1}^*, 1; r_{-1}) = 0$ for some $p \geq p_{+1}^*$ and $r_{-1} \leq r_{-1,c}$, then $\frac{\partial H_{-1}}{\partial r_{-1}}(p, p_{-1}^*, 1; r_{-1}) < 0$.*

Proof. We only prove the assertion on H_{+1} . In this proof, the dependence of μ_{++} and μ_{+-} on r_{+1} will be made explicit. By Eq. (42), $\xi_{++}(p, p_{+1}^*; r_{+1})$ converges to a function of p while $\xi_{+-}(p, p_{+1}^*) \sim -\gamma\mu_{+-}^{-1}e^{\gamma^{-1}(p-p_{+1}^*)\mu_{+-}(r_{+1})}$. Therefore,

$$\begin{aligned} & (\mu_{++}(r_{+1}) - \mu_{+-}(r_{+1}))H_{+1}(p, p_{+1}^*, 1; r_{+1}) \\ & \sim \left[(r_{+1} + 2\lambda_{-1} + \mu_{++}(r_{+1}))\gamma\mu_{+-}(r_{+1})^{-1} + (r_{+1} + \lambda_{-1} + \mu_{++}(r_{+1}))\gamma\lambda_{+1}^{-1} \right] e^{-\gamma^{-1}(p_{+1}^*-p)\mu_{+-}(r_{+1})}. \end{aligned}$$

(As $r_{+1} \rightarrow \infty$, $e^{-\gamma^{-1}(p_{+1}^*-p)\mu_{++}(r_{+1})}$ approaches zero faster and $e^{-\gamma^{-1}(p_{+1}^*-p)\mu_{+-}(r_{+1})}$ approaches infinity faster. Therefore, the concern of $r_{+1} \rightarrow \infty$ does not jeopardize the above result.) Consequently, the sign of $H_{+1}(p, p_{+1}^*, 1; r_{+1})$ is the same as that of

$$(r_{+1} + 2\lambda_{-1} + \mu_{++}(r_{+1}))\mu_{+-}(r_{+1})^{-1} + (r_{+1} + \lambda_{-1} + \mu_{++}(r_{+1}))\lambda_{+1}^{-1}.$$

Note that $\mu_{+-}(r_{+1}) \sim \pm r_{+1}$ as $r_{+1} \rightarrow \infty$. Therefore, the first term in the above expression approaches -2 as $r_{+1} \rightarrow \infty$ and the second term approaches infinity as $r_{+1} \rightarrow \infty$. Therefore, $H_{+1}(p, p_{+1}^*, 1; r_{+1}) > 0$ when r_{+1} is above some threshold $r_{+1,c}$ that is independent of p .

Now consider a finite r_{+1} . Note that

$$\begin{aligned} & (\mu_{++} - \mu_{+-})H(p, p_{+1}^*, 1; \tilde{r}_{+1}) \\ = & (\tilde{r}_{+1} + 2\lambda_{-1} + \mu_{+-}(\tilde{r}_{+1})) \int_p^{p_{+1}^*} e^{-\gamma^{-1}(\tilde{p}-p)\mu_{++}(\tilde{r}_{+1})} d\tilde{p} - (\tilde{r}_{+1} + 2\lambda_{-1} + \mu_{++}(\tilde{r}_{+1})) \int_p^{p_{+1}^*} e^{-\gamma^{-1}(\tilde{p}-p)\mu_{+-}(\tilde{r}_{+1})} d\tilde{p} \\ + & \gamma\lambda_{+1}^{-1}(\tilde{r}_{+1} + \lambda_{-1} + \mu_{++}(\tilde{r}_{+1}))e^{-\gamma^{-1}(p_{+1}^*-p)\mu_{+-}(\tilde{r}_{+1})} - \gamma\lambda_{+1}^{-1}(r_{+1} + \lambda_{-1} + \mu_{+-}(\tilde{r}_{+1}))e^{-\gamma^{-1}(p_{+1}^*-p)\mu_{++}(\tilde{r}_{+1})}. \end{aligned}$$

The dependence of the $\mu_{\pm\pm}$ on \tilde{r}_{+1} has been made explicit. Denote the four terms on the right hand side by $A(\tilde{r}_{+1})$, $-B(\tilde{r}_{+1})$, $C(\tilde{r}_{+1})$, $-D(\tilde{r}_{+1})$, respectively. The choice of signs ensures that all the four new functions are positive. First compute the derivative of $\mu_{\pm\pm}(\tilde{r}_{+1})$:

$$\mu'_{\pm\pm}(\tilde{r}_{+1}) = \pm m_{+1}(\tilde{r}_{+1}) := \pm [(\lambda_{+1} - \lambda_{-1})^2 + 4\tilde{r}_{+1}^2 + 4\tilde{r}_{+1}(\lambda_{+1} + \lambda_{-1})]^{-1/2} (2r_{+1} + \lambda_{+1} + \lambda_{-1}).$$

(The symbol “:=” means that the right hand side is the definition of the left hand side.) It is easy to see that $m_{+1}(\tilde{r}_{+1}) > 1$. Next compute the log-derivatives of the four terms:

$$\begin{aligned} a(\tilde{r}_{+1}) := A(\tilde{r}_{+1})^{-1}A'(\tilde{r}_{+1}) &= \frac{1 - m_{+1}(\tilde{r}_{+1})}{\tilde{r}_{+1} + 2\lambda_{-1} + \mu_{+-}(\tilde{r}_{+1})} - \frac{m_{+1}(\tilde{r}_{+1})}{\mu_{++}(\tilde{r}_{+1})} \\ &+ \gamma^{-1}(p_{+1}^* - p)m_{+1}(\tilde{r}_{+1}) \left(e^{\gamma^{-1}(p_{+1}^*-p)\mu_{++}(\tilde{r}_{+1})} - 1 \right)^{-1}; \\ b(\tilde{r}_{+1}) := B(\tilde{r}_{+1})^{-1}B'(\tilde{r}_{+1}) &= \frac{1 + m_{+1}(\tilde{r}_{+1})}{\tilde{r}_{+1} + 2\lambda_{-1} + \mu_{++}(\tilde{r}_{+1})} + \frac{m_{+1}(\tilde{r}_{+1})}{\mu_{+-}(\tilde{r}_{+1})} \\ &+ \gamma^{-1}(p_{+1}^* - p)m_{+1}(\tilde{r}_{+1}) \left(1 - e^{\gamma^{-1}(p_{+1}^*-p)\mu_{+-}(\tilde{r}_{+1})} \right)^{-1}; \\ c(\tilde{r}_{+1}) := C(\tilde{r}_{+1})^{-1}C'(\tilde{r}_{+1}) &= \frac{1 + m_{+1}(\tilde{r}_{+1})}{\tilde{r}_{+1} + \lambda_{-1} + \mu_{++}(\tilde{r}_{+1})} + \gamma^{-1}m_{+1}(\tilde{r}_{+1})(p_{+1}^* - p); \\ d(\tilde{r}_{+1}) := D(\tilde{r}_{+1})^{-1}D'(\tilde{r}_{+1}) &= \frac{1 - m_{+1}(\tilde{r}_{+1})}{\tilde{r}_{+1} + \lambda_{-1} + \mu_{+-}(\tilde{r}_{+1})} - \gamma^{-1}m_{+1}(\tilde{r}_{+1})(p_{+1}^* - p). \end{aligned}$$

It is easy to see that $d(\tilde{r}_{+1}) < 0$. By assumption, $A(r_{+1}) - B(r_{+1}) + C(r_{+1}) - D(r_{+1}) = 0$. Therefore, $C(r_{+1}) = B(r_{+1}) - A(r_{+1}) + D(r_{+1}) > B(r_{+1}) - A(r_{+1})$. It follows that

$$\begin{aligned} \frac{\partial H_{+1}}{\partial r_{+1}}(p, p_{+1}^*, 1; r_{+1}) &= a(r_{+1})A(r_{+1}) - b(r_{+1})B(r_{+1}) + c(r_{+1})C(r_{+1}) - d(r_{+1})D(r_{+1}) \\ &> (c(r_{+1}) - b(r_{+1}))B(r_{+1}) - (c(r_{+1}) - a(r_{+1}))A(r_{+1}). \end{aligned}$$

Note that

$$\begin{aligned} & c(r_{+1}) - b(r_{+1}) \\ = & \frac{(1 + m_{+1}(r_{+1}))\lambda_{-1}}{(r_{+1} + 2\lambda_{-1} + \mu_{++}(r_{+1}))(r_{+1} + \lambda_{-1} + \mu_{++}(r_{+1}))} + \\ & + m_{+1}(r_{+1}) \left[-\frac{1}{\mu_{+-}(r_{+1})} - \gamma^{-1}(p_{+1}^* - p) \left(e^{-\gamma^{-1}(p_{+1}^*-p)\mu_{+-}(r_{+1})} - 1 \right)^{-1} \right]. \end{aligned}$$

The first term is independent of γ and is bounded away from zero for $r_{+1} \in [0, r_{+1,c}]$ as its limit when $r_{+1} \rightarrow 0$ is positive. (Here the assumption that $\lambda_{+1} \neq \lambda_{-1}$ has been used.) The

term in the bracket is positive and strictly decreasing in μ_{+-} . Its limit when $\mu_{+-} \rightarrow 0$ is $\frac{1}{2}\gamma^{-1}(p_{+1}^* - p)$. Therefore,

$$c(r_{+1}) - b(r_{+1}) > \frac{(1 + m_{+1}(r_{+1}))\lambda_{-1}}{(r_{+1} + 2\lambda_{-1} + \mu_{++}(r_{+1}))(r_{+1} + \lambda_{-1} + \mu_{++}(r_{+1}))} + \frac{1}{2}\gamma^{-1}(p_{+1}^* - p). \quad (54)$$

A similar calculation shows that $a(r_{+1}) - d(r_{+1}) > 0$. Therefore,

$$c(r_{+1}) - a(r_{+1}) < c(r_{+1}) - d(r_{+1}) = \left[\frac{1 + m_{+1}(r_{+1})}{r_{+1} + \lambda_{-1} + \mu_{++}(r_{+1})} - \frac{1 - m_{+1}(r_{+1})}{r_{+1} + \lambda_{-1} + \mu_{+-}(r_{+1})} \right] + 2\gamma^{-1}m_{+1}(r_{+1})(p_{+1}^* - p). \quad (55)$$

The term in the bracket is independent of γ and is bounded when $r_{+1} \in [0, r_{+1,c}]$. Finally,

$$\frac{A(r_{+1})}{B(r_{+1})} = \frac{r_{+1} + 2\lambda_{-1} + \mu_{++}(r_{+1})}{r_{+1} + 2\lambda_{-1} + \mu_{+-}(r_{+1})} \frac{|\mu_{+-}(r_{+1})| \left(1 - e^{-\gamma^{-1}(p_{+1}^* - p)\mu_{++}(r_{+1})}\right)}{\mu_{++}(r_{+1}) \left(e^{-\gamma^{-1}(p_{+1}^* - p)\mu_{+-}(r_{+1})} - 1\right)}.$$

The first fraction is independent of γ and is bounded for $r_{+1} \in [0, r_{+1,c}]$. The second fraction is actually the ratio between two integrals:

$$\frac{\int_p^{p_{+1}^*} e^{-\gamma^{-1}(\tilde{p}-p)\mu_{++}(r_{+1})} d\tilde{p}}{\int_p^{p_{+1}^*} e^{-\gamma^{-1}(\tilde{p}-p)\mu_{+-}(r_{+1})} d\tilde{p}},$$

which is strictly decreasing in r_{+1} . The limit of the ratio as $r_{+1} \rightarrow 0$ is

$$\frac{(\lambda_{-1} - \lambda_{+1})(p_{+1}^* - p)}{\gamma \left(e^{\gamma^{-1}(p_{+1}^* - p)(\lambda_{-1} - \lambda_{+1})} - 1\right)},$$

if $\lambda_{-1} > \lambda_{+1}$, and

$$\frac{\gamma \left(1 - e^{-\gamma^{-1}(p_{+1}^* - p)(\lambda_{+1} - \lambda_{-1})}\right)}{(\lambda_{+1} - \lambda_{-1})(p_{+1}^* - p)},$$

if $\lambda_{+1} > \lambda_{-1}$. Either way, the ratio approaches zero at least as fast as γ as $\gamma \rightarrow 0$. Therefore, there exists a $\eta > 0$ such that

$$\frac{A(r_{+1})}{B(r_{+1})} < \eta\gamma, \quad (56)$$

for all $\gamma \leq \bar{\gamma}$ and $r_{+1} \leq r_{+1,c}$. Combining Eqs. (54)-(56) yields that

$$\frac{\partial H_{+1}}{\partial r_{+1}}(p, p_{+1}^*, 1; r_{+1}) > m_{+1}(r_{+1})B(r_{+1}) \left[E_1(r_{+1}) + \frac{1}{2}\gamma^{-1}(p_{+1}^* - p) - \eta\gamma(E_2(r_{+1}) + 2\gamma^{-1}(p_{+1}^* - p)) \right],$$

where $E_1(r_{+1})$ is the first term on the right hand side of Eq. (54) and $E_2(r_{+1})$ is the term in the bracket on the right hand side of Eq. (55). Both E_1 and E_2 are positive and bounded. The above inequality holds for all $r_{+1} \in [0, r_{+1,c}]$ and $p \leq p_{-1}^*$. The bracket on the right hand side of the inequality approaches infinity as $\gamma \rightarrow 0$. Therefore, there exists some $\bar{\gamma}_{+1} \leq \bar{\gamma}$ such that for $\gamma \leq \bar{\gamma}_{+1}$ and $p \leq p_{-1}^*$, that $H_{+1}(p, p_{+1}^*, 1; r_{+1}) = 0$ implies that $\frac{\partial H_{+1}}{\partial r_{+1}}(p, p_{+1}^*, 1; r_{+1}) > 0$. \blacksquare

Proof of Proposition 8

Proposition 8 is a corollary of Lemma B.6f. ■

Proof of Proposition 9

Let $\bar{\gamma} = \min\{\gamma_1, \gamma_2\}$. By Lemma B.6g, the map $BR : (-\infty, p_{-1}^*] \times [p_{+1}^*, \infty) \rightarrow (-\infty, p_{-1}^*] \times [p_{+1}^*, \infty)$ defined by $BR(p_{-1}, p_{+1}) = (BR_{+1}(p_{-1}), BR_{-1}(p_{+1}))$ is a contraction mapping. Therefore, it has a unique fixed point. By Lemma B.6c, the game has a unique Markov Perfect Equilibrium in focused strategies, with the unique fixed point of BR as the parties' targets. By Lemma B.6b, $BR_{+1}(p_{-1}) = p_{+1}^*$ if and only if $H_{+1}(p_{-1}, p_{+1}^*, 1) \geq 0$ and $BR_{-1}(p_{+1}) = p_{-1}^*$ if and only if $H_{-1}(p_{+1}, p_{-1}^*, 1) < 0$. By Lemma B.6h, if $r_{+1} > r_{+1,c}$ and $r_{-1} > r_{-1,c}$, then $BR_{+1}(p_{-1}) = p_{+1}^*$ for every $p_{-1} \leq p_{-1}^*$ and $BR_{-1}(p_{+1}) = p_{-1}^*$ for every $p_{+1} \geq p_{+1}^*$ and thus (p_{+1}^*, p_{-1}^*) is the unique equilibrium target.

Next, we show that if in the unique equilibrium $(p_{+1}^{**}, p_{-1}^{**})$, $p_{+1}^{**} = p_{+1}^*$ for some r_{+1} , then $p_{+1}^{**} = p_{+1}^*$ when r_{+1} increases to any $\tilde{r}_{+1} > r_{+1}$. By Lemma B.6h, $H_{+1}(p_{-1}^*, p_{+1}^*, 1; r_{+1}) \geq 0$. Suppose that $H_{+1}(p_{-1}^*, p_{+1}^*, 1; \tilde{r}_{+1}) < 0$. Then let

$$r_{+1,0} = \sup\{r \geq r_{+1} : H_{+1}(p_{-1}^*, p_{+1}^*, 1; \tilde{r}) \geq 0 \text{ for every } \tilde{r} \in [r_{+1}, r]\}.$$

Then since $H_{+1}(p_{-1}^*, p_{+1}^*, 1; r)$ is continuously differentiable in r ,

$$\begin{aligned} H_{+1}(p_{-1}^*, p_{+1}^*, 1; r_{+1,0}) &= 0, \text{ and} \\ \frac{\partial H_{+1}}{\partial r_{+1}}(p_{-1}^*, p_{+1}^*, 1; r_{+1,0}) &\leq 0, \end{aligned}$$

contradicting Lemma B.6h. Therefore, $H_{+1}(p_{-1}^*, p_{+1}^*, 1; \tilde{r}_{+1}) \geq 0$ and thus $BR_{+1}(p_{-1}^*; \tilde{r}_{+1}) = p_{+1}^*$. Since r_{+1} does not affect BR_{-1} , (p_{+1}^*, p_{-1}^*) remains the unique equilibrium. A symmetric argument shows that if $p_{-1}^{**} = p_{-1}^*$ and r_{-1} increases to some $\tilde{r}_{-1} \geq r_{-1}$, then (p_{+1}^*, p_{-1}^*) remains the unique equilibrium.

Finally, consider the behavior of $H_{+1}(p, p_{+1}^*, 1)$ as $r_{+1} \rightarrow 0$ for an arbitrary $p \leq p_{-1}^*$. As shown in Lemma B.6h, when γ is sufficiently small, the sign of $H_{+1}(p, p_{+1}^*, 1)$ is the same as the sign of

$$(r_{+1} + 2\lambda_{-1} + \mu_{++}(r_{+1}))\mu_{+-}(r_{+1})^{-1} + (r_{+1} + \lambda_{-1} + \mu_{++}(r_{+1}))\lambda_{+1}^{-1}. \quad (57)$$

According to Eq. (41), as $r_{+1} \rightarrow 0$,

$$(\mu_{++}(r_{+1}), \mu_{+-}(r_{+1})) \rightarrow \begin{cases} (\lambda_{+1} - \lambda_{-1}, 0) & , \text{ if } \lambda_{+1} > \lambda_{-1}; \\ (0, \lambda_{+1} - \lambda_{-1}) & , \text{ if } \lambda_{+1} < \lambda_{-1}. \end{cases}$$

Therefore, when $\lambda_{+1} > \lambda_{-1}$, the expression in Eq. (57) approaches $-\infty$ as $r_{+1} \rightarrow 0$, and when $\lambda_{+1} < \lambda_{-1}$, the expression in Eq. (57) approaches $-\frac{2\lambda_{-1}}{\lambda_{-1}-\lambda_{+1}} + \frac{\lambda_{-1}}{\lambda_{+1}}$ as $r_{+1} \rightarrow 0$. Therefore, for $H_{+1}(p, p_{+1}^*, 1) < 0$ and thus $BR_{+1}(p; r_{+1}) > p_{+1}^*$ as $r_{+1} \rightarrow 0$ if $\lambda_{+1} \neq \lambda_{-1}$ and $\lambda_{+1} > \frac{1}{3}\lambda_{-1}$. By a symmetric argument, $BR_{-1}(p; r_{-1}) < p_{-1}^*$ as $r_{-1} \rightarrow 0$ if $\lambda_{-1} \neq \lambda_{+1}$ and $\lambda_{-1} > \frac{1}{3}\lambda_{+1}$. To sum up, as long as $\lambda_{-1} \neq \lambda_{+1}$, at least one party exhibits strategic extremism when both r_{+1} and r_{-1} approach zero. ■

The following result calculates (asymptotically) the extent of strategic extremism by each party, $\Delta_i^{**} = |p_i^{**} - p_i^*|$.

Proposition B.1. *Suppose that $r_{+1} = 0$ and that $\lambda_{+1} > \lambda_{-1}$. Then the extent of strategic extremism by party +1 is, asymptotically for large γ^{-1} and Δ_p^* ,*

$$\begin{aligned}\Delta_{+1}^{**} &\rightarrow \max \left\{ 0, p_{+1}^* - p_{-1}^* + \Delta_{-1}^{**} - \frac{\gamma}{\lambda_{+1} - \lambda_{-1}} + O\left(e^{-\gamma^{-1}(p_{+1}^* - p_{-1}^*)(\lambda_{+1} - \lambda_{-1})}\right) \right\}; \\ \Delta_{-1}^{**} &\rightarrow \max \left\{ 0, \frac{\gamma}{\lambda_{+1} - \lambda_{-1}} \log \left[\frac{4\lambda_{-1}}{\lambda_{+1} + \lambda_{-1}} + O\left(\gamma^{-1} e^{-\gamma^{-1}(p_{+1}^* - p_{-1}^*)(\lambda_{+1} - \lambda_{-1})}\right) \right] \right\},\end{aligned}$$

Proof. By Eq. (41), as r_{+1} and r_{-1} approach zero,

$$(\mu_{i+}, \mu_{i-}) \rightarrow (\lambda_{+1} - \lambda_{-1}, 0), \text{ for } i \in \{-1, +1\}.$$

Substituting these into the expressions of H_{+1} and H_{-1} in the proof of Lemma B.6b yields

$$(\mu_{++} - \mu_{+-})H_{+1}(p_{-1}, p_{+1}, 1) \rightarrow \frac{\gamma(\lambda_{+1} + \lambda_{-1})}{\lambda_{+1} - \lambda_{-1}} - (\lambda_{+1} + \lambda_{-1})(2p_{+1}^* - p_{-1} - p_{+1}) + O\left(e^{-\gamma^{-1}(p_{+1}^* - p_{-1}^*)(\lambda_{+1} - \lambda_{-1})}\right)$$

and

$$\begin{aligned} &(\mu_{-+} - \mu_{--})H_{-1}(p_{+1}, p_{-1}, 1)e^{-\gamma^{-1}(p_{+1} - p_{-1}^*)(\lambda_{+1} - \lambda_{-1})} \\ \rightarrow &\frac{2\gamma\lambda_{+1}}{\lambda_{+1} - \lambda_{-1}} \left[2 - e^{\gamma^{-1}(p_{-1}^* - p_{-1})(\lambda_{+1} - \lambda_{-1})} \right] - \frac{\gamma\lambda_{+1}}{\lambda_{-1}} e^{\gamma^{-1}(p_{-1}^* - p_{-1})(\lambda_{+1} - \lambda_{-1})} + O\left(e^{-\gamma^{-1}(p_{+1} - p_{-1}^*)(\lambda_{+1} - \lambda_{-1})}\right),\end{aligned}$$

as $r_{+1}, r_{-1} \rightarrow 0$ and for $p_{+1} \geq p_{+1}^*$ and $p_{-1} \leq p_{-1}^*$. By construction of the best response functions,

$$\begin{aligned}BR_{+1}(p_{-1}) &\rightarrow \max \left\{ p_{+1}^*, 2p_{+1}^* - p_{-1} - \frac{\gamma}{\lambda_{+1} - \lambda_{-1}} + O\left(e^{-\gamma^{-1}(p_{+1}^* - p_{-1}^*)(\lambda_{+1} - \lambda_{-1})}\right) \right\}; \\ BR_{-1}(p_{+1}) &\rightarrow \min \left\{ p_{-1}^*, p_{-1}^* - \frac{\gamma}{\lambda_{+1} - \lambda_{-1}} \log \left[\frac{4\lambda_{-1}}{\lambda_{+1} + \lambda_{-1}} + O\left(\gamma^{-1} e^{-\gamma^{-1}(p_{+1}^* - p_{-1}^*)(\lambda_{+1} - \lambda_{-1})}\right) \right] \right\}.\end{aligned}$$

Therefore, in the unique equilibrium when $\gamma^{-1}(p_{+1}^* - p_{-1}^*)(\lambda_{+1} - \lambda_{-1})$ is sufficiently big,

$$\begin{aligned}\Delta_{+1}^{**} &\rightarrow \max \left\{ 0, p_{+1}^* - p_{-1}^* + \Delta_{-1}^{**} - \frac{\gamma}{\lambda_{+1} - \lambda_{-1}} + O\left(e^{-\gamma^{-1}(p_{+1}^* - p_{-1}^*)(\lambda_{+1} - \lambda_{-1})}\right) \right\}; \\ \Delta_{-1}^{**} &\rightarrow \max \left\{ 0, \frac{\gamma}{\lambda_{+1} - \lambda_{-1}} \log \left[\frac{4\lambda_{-1}}{\lambda_{+1} + \lambda_{-1}} + O\left(\gamma^{-1} e^{-\gamma^{-1}(p_{+1}^* - p_{-1}^*)(\lambda_{+1} - \lambda_{-1})}\right) \right] \right\},\end{aligned}$$

as $r_{+1}, r_{-1} \rightarrow 0$. ■

References

- AZZIMONTI, M. (2015): “Partisan Conflict and Private Investment,” *Working paper*.
- BENDOR, J. (1995): “A Model of Muddling Through,” *The American Political Science Review*, 89(4), 819–840.
- BERNHEIM, B. D., A. RANGEL, AND L. RAYO (2006): “The Power of the Last Word in Legislative Policy Making,” *Econometrica*, 74(5), 1161–1190.

- BHATTACHARYA, R., AND M. MAJUMDAR (2003): "Random dynamical systems: a review," *Economic Theory*, 23(1), 13–1.
- BUISSERET, P., AND D. BERNHARDT (2015): "Dynamics of Policymaking: Stepping Back to Leap Forward, Stepping Forward to Keep Back," *Working Paper*.
- CALLANDER, S. (2011a): "Searching and Learning by Trial and Error," *The American Economic Review*, 101(6), 2277–2308.
- (2011b): "Searching for Good Policies," *American Political Science Review*, 105(4), 643–662.
- CALLANDER, S., AND P. HUMMEL (2014): "Preemptive Policy Experimentation," *Econometrica*, 82(4), 1509–1528.
- CHEN, Y., AND H. ERASLAN (2017): *American Economic Journal: Microeconomics* 9(2), 1–32.
- ELLISON, G., AND R. HOLDEN (2013): "A Theory of Rule Development," *Journal of Law Economics and Organization*.
- ELY, J. C. (2011): "Kludged," *American Economic Journal: Microeconomics*, 3(3), pp. 210–231.
- GLAESER, E. L., G. A. M. PONZETTO, AND J. M. SHAPIRO (2005): "Strategic Extremism: Why Republicans and Democrats Divide on Religious Values," *The Quarterly Journal of Economics*, 120(4), 1283–1330.
- GRATTON, G., L. GUISO, C. MICHELACCI, AND M. MORELLI (2015): "From Weber to Kafka: Political Activism and the Emergence of an Inefficient Bureaucracy," *Working paper*.
- LEVY, G., AND R. RAZIN (2013): "Dynamic legislative decision making when interest groups control the agenda," *Journal of Economic Theory*, 148(5), 1862–1890.
- LINDBLOM, C. E. (1959): "The Science of "Muddling Through"," *Public Administration Review*, 19(2), 79–88.
- MCCARTY, N., K. T. POOLE, AND H. ROSENTHAL (2016): *Polarized America, The Dance of Ideology and Unequal Riches*. MIT Press.
- MESSNER, M., AND M. K. POLBORN (2012): "The option to wait in collective decisions and optimal majority rules," *Journal of Public Economics*, 96(5-6), 524–540.
- TELES, S. (2013): "Kludgeocracy in America," *National Affairs*, 17, 97–114.