

# Online Appendix for “Political Kludges”

Keiichi Kawai, Ruitian Lang, Hongyi Li

October 31, 2017

## Appendix B Proofs

### Short-Run Dynamics

**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

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

where  $\text{Conv}(S)$  is the convex closure of set  $S$ . 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  $v_\delta$  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

**Notation** Identify the state space as  $\mathcal{X} = [p_{-1}^*, p_{+1}^*] \times \{+1, -1\}$ , with generic element  $x = (p, i) \in \mathcal{X}$ .

We denote the sequence of random *transition 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,  $0 < t_1^{-1} < t_1^{+1} < t_2^{-1} < t_2^{+1} < \dots$ . A *transition history* is a sequence of transition times  $\{t_1^{-1}, t_1^{+1}, t_2^{-1}, t_2^{+1}, \dots\}$ . Notice that given a starting position  $\mathbf{p}(0)$ , a transition history 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 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.  $\blacksquare$

**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  $p_{+1}^*/\gamma$ , 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. For any set of states  $Y \subseteq \mathcal{X}$  and any starting state  $x \in \mathcal{X}$ , define  $\tau_x(Y) = \inf\{t \geq 0 : (q(t), i(t)) \in Y\}$  to be the first hitting time for  $Y$ , given starting state  $(q(0), i(0)) = x$ . A Markov process is *Harris recurrent* if, for some (finite or  $\sigma$ -finite) measure  $\varphi$ ,  $\Pr[\tau_x(Y) < \infty] = 1$  for all  $x \in \mathcal{X}$  and all  $Y \subseteq \mathcal{X}$  with  $\varphi(Y) > 0$ ; see, e.g., Meyn and Tweedie (1993) p. 490, or Theorem 1 of Kaspi and Mandelbaum (1994). An invariant probability measure for a Markov process is *ergodic* if every invariant subset  $Y \subseteq \mathcal{X}$  has mass of either 0 or 1; see, e.g., Definition 3.4 of Hairer (2008).

**Lemma B.1a.** *The process  $(q(t), i(t))$  is Harris recurrent.*

*Proof.* Consider a finite measure  $\varphi$  which puts all mass on the state  $(p_{-1}^*, -1)$ , so that  $\varphi(Y) > 0$  iff  $(p_{-1}^*, -1) \in Y$ . It suffices to show that  $\Pr[\tau_x(\{(p_{-1}^*, -1)\}) < \infty] = 1$  for all  $x \in \mathcal{X}$ . For a given point in sample space,  $\tau_\omega(\{(p_{-1}^*, -1)\}) = \infty$  only if  $\Delta t_k^{-1} < \frac{p_{+1}^* - p_{-1}^*}{\gamma}$  for all  $k \geq 1$ ; that is, if Party  $-1$  never remains in control long enough to move to position  $q = p_{-1}^*$ . But this is a probability-zero event because each  $\Delta t_k^{-1}$  is i.i.d. exponentially distributed. ■

**Proof of Lemma 1**

From Lemma B.1a, the process  $(q(t), i(t))$  is Harris recurrent. Any Harris recurrent process has a unique invariant measure (Azema, Kaplan-Duflo, and Revuz, 1967; see also the

discussions in Meyn and Tweedie, 1993, p. 491, and in Kaspi and Mandelbaum, 1994, p. 212). Our state space is compact, so this measure is finite and can be normalized to a (unique invariant) probability measure. Finally, if a Markov process has a unique invariant probability measure, then this measure is (uniquely) ergodic; see Corollary 5.6 of Hairer (2008). ■

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

$$(B.1a) \quad f(q, +1) \equiv f(q, -1) \equiv Ae^{\frac{\lambda_{+1}-\lambda_{-1}}{\gamma}q}$$

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

$$(B.1b) \quad \Delta F(p_{+1}^*, +1) = \frac{\gamma}{\lambda_{+1}} f(p_{+1}^*, +1) \text{ and } \Delta F(p_{-1}^*, -1) = \frac{\gamma}{\lambda_{-1}} f(p_{-1}^*, -1)$$

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  $\Delta t$ , the net change in the probability mass of  $[q, q + \Delta q] \times \{+1\}$  must be zero; that is,

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

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

$$(B.2) \quad \gamma f_q(q, +1) = \lambda_{-1} f(q, -1) - \lambda_{+1} f(q, +1) \text{ and}$$

$$(B.3) \quad \gamma f_q(q, -1) = \lambda_{-1} f(q, -1) - \lambda_{+1} f(q, +1)$$

for  $q \in [p_{-1}^*, p_{+1}^*]$ , where (B.3) holds by a symmetric argument. Solving the differential equations (B.2) and (B.3) simultaneously reveals that

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

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

$$\lambda_{+1} \Delta F(p_{+1}^*, +1) \Delta t - \gamma f(p_{+1}^*, +1) \Delta t \approx 0,$$

$$\lambda_{-1} \Delta F(p_{-1}^*, -1) \Delta t - \gamma f(p_{+1}^*, +1) \Delta t \approx 0$$

or, more compactly,

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

■

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  $\varepsilon$  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  $\|\mathbf{p}\| < \underline{c}$  for all  $\mathbf{p} \in \mathcal{B}$ . Suppose that the true and simulacrum process share the same transition history, as well as identical initial conditions:  $\|\mathbf{p}(0)\| = c(0) \geq \underline{c}$  and  $p(0) = q(0)$ . Define  $T = \inf\{t : c(t) < \underline{c}\}$ . Then  $q(t) = p(t)$  and  $c(t) \leq \|\mathbf{p}(t)\|$  for all  $t \leq T$ .*

*Proof.* This result requires only a straightforward inspection of the laws of motion of  $\mathbf{p}$  (outside  $\mathcal{B}$ ) and  $c, q$ . Specifically,  $c(t) \geq \underline{c}$  for all  $t < T$ , so  $\frac{d}{dt}p(t) = \frac{d}{dt}q(t)$  and  $\frac{d}{dt}\|\mathbf{p}(t)\| \leq \frac{d}{dt}c(t)$ , and thus  $p(t) \equiv q(t)$  and  $\|\mathbf{p}(t)\| \leq c(t)$ . ■

**Lemma B.2b.** *Suppose  $z_{+1} > 2$  and  $z_{-1} > 2$ . Select sufficiently large  $\underline{c}$  so that  $\|\mathbf{p}\| < \underline{c}$  for all  $\mathbf{p} \in \mathcal{B}$ , and select sufficiently small  $\varepsilon$  so that  $\varepsilon < 1 - \frac{\underline{c}-|p|}{\underline{c}+|p|}$  for  $p \in \{p_{-1}^*, p_{+1}^*\}$ . Suppose that the true and simulacrum process share the same transition history, as well as identical initial conditions:  $\|\mathbf{p}(0)\| \equiv c_\varepsilon(0) > \underline{c}$  and  $p(0) = q(0)$ . Suppose that  $c_\varepsilon(T) \leq \underline{c}$  at some time  $T > 0$ . Then there exists  $\tau \leq T$  such that  $\|\mathbf{p}(\tau)\| = \underline{c}$ .*

*Proof.* Suppose, towards a contradiction, that  $\|\mathbf{p}(t)\| > \underline{c}$  for all  $t \leq T$ . Then throughout this time interval,  $\frac{d}{dt}p(t) = \frac{d}{dt}q(t)$  and  $\frac{d}{dt}c_\varepsilon(t) \geq \frac{d}{dt}\|\mathbf{p}(t)\|$ , so  $p(t) \equiv q(t)$  and  $c_\varepsilon(t) \geq \|\mathbf{p}(t)\| > \underline{c}$ . This contradicts the assumption that  $c_\varepsilon(T) \leq \underline{c}$ . ■

**Lemma B.2c.** *Suppose  $z_{+1} > 2$  and  $z_{-1} > 2$ . For any complexity bound  $\underline{c} > 0$ , there exists some  $0 < v_\underline{c} < 1$  such that the following holds. Suppose that at some transition time  $t_k^i$ , complexity lies below this bound, i.e.,  $\|\mathbf{p}(t_k^i)\| \leq \underline{c}$ , and policy lies outside the basin, i.e.,  $\mathbf{p}(t_k^i) \notin \mathcal{B}$ . (i) Then with probability of at least  $v_\underline{c}$ , policy lies within the basin at the very next transition time:  $\|\mathbf{p}(t_{k'}^{-i})\| \in \mathcal{B}$ . (ii) Further, a.s., at some future transition time  $t_{k''}^{i''}$ , policy either exceeds the complexity bound, i.e.,  $\|\mathbf{p}(t_{k''}^{i''})\| > \underline{c}$ , or lies within the basin  $\mathcal{B}$ .*

*Proof.* Let  $\Delta\hat{t} = \frac{p_{+1}^* - p_{-1}^*}{\gamma}$  be the amount of time taken for policy to move from ideal  $p_i^*$  to  $p_{-i}^*$  by adding  $(-i)$ -rules. So, if player  $-i$  remains in control for a period of at least  $\Delta\hat{t}$  after taking control at time  $t_k^i$ , then he will reach position  $p_{-i}^*$  at some time  $t' \leq t_k^i + \Delta\hat{t}$  within this period – at which point  $\|\mathbf{p}(t')\| \leq \underline{c} + p_{+1}^* - p_{-1}^*$ .

Let  $\Delta\hat{t}_{-i} < \infty$  be the time taken for Party  $-i$  to reduce complexity along his ideal from  $\mathbf{p} = (\underline{c} + p_{+1}^* - p_{-1}^*, p_{-i}^*)$  to reach the basin  $\mathcal{B}$ . Let  $\Delta\hat{t} = \max\{\Delta\hat{t}_{-1}, \Delta\hat{t}_{+1}\}$ . So, if player  $-i$

remains in control for a period of at least  $\Delta\hat{t} + \hat{\Delta}t$  after taking control at time  $t_k^i$ , then policy will enter the basin  $\mathcal{B}$  at some time  $t'' \leq t_k^i + \Delta\hat{t} + \hat{\Delta}t$  within this period. This event occurs with probability of at least  $e^{-\max\{\lambda_{+1}, \lambda_{-1}\}(\Delta\hat{t} + \hat{\Delta}t)} > 0$ . Thus, part (i) holds.

Further, for (ii) not to occur, it must be that  $\|\mathbf{p}\| \leq \underline{c}$  at each future transition time after  $t_k^i$ . By part (i), at each transition time, policy enters the basin with probability of at least  $v_{\underline{c}}$ . It follows that (ii) occurs a.s.  $\blacksquare$

**Lemma B.2d.** *Define the random variable  $\|\mathbf{p}_\infty\|$  to take values on  $\{0_-\} \cup [0, \infty]$ , as follows:*

$$\|\mathbf{p}_\infty\| = \begin{cases} 0_- & \text{if } \mathbf{p}(t) \in \mathcal{B} \text{ for some } t \geq 0, \\ \liminf \{ \|\mathbf{p}_1^{-1}\|, \|\mathbf{p}_1^1\|, \|\mathbf{p}_2^{-1}\|, \dots \} & \text{otherwise} \end{cases}$$

where  $\|\mathbf{p}_k^i\| = \|\mathbf{p}(t_k^i)\|$  denotes complexity at the transition time  $t_k^i$ . Suppose that with positive probability, policy becomes neither perfectly simple nor kludged. Then there exists  $\underline{c} \in (0, \infty)$  such that  $\|\mathbf{p}_\infty\| \in [0, \underline{c})$  with positive probability.

*Proof.* While policy remains outside the basin, troughs in complexity coincide with transition times when: one party loses control after reducing complexity at his own ideal and the other party immediately starts increasing complexity. Consequently, if policy never enters the basin, then

$$\liminf_{t \rightarrow \infty} \|\mathbf{p}(t)\| = \liminf \{ \|\mathbf{p}_1^{-1}\|, \|\mathbf{p}_1^1\|, \|\mathbf{p}_2^{-1}\|, \dots \}.$$

Next, observe that  $\|\mathbf{p}_\infty\| \in [0, \infty)$  iff policy never becomes simple or kludged. Thus, by our supposition, the distribution of  $\|\mathbf{p}_\infty\|$  must have nonzero probability mass on  $[0, \infty)$ . The result follows.  $\blacksquare$

**Lemma B.2e.** *Fix  $\underline{c} > 0$ . The number of transition times  $t$  in the sequence  $\{t_1^{-1}, t_1^{+1}, t_2^{-1}, t_2^{+1}, \dots\}$  whereby  $\mathbf{p}(t) \notin \mathcal{B}$  and  $\|\mathbf{p}(t)\| \leq \underline{c}$  is a.s. finite.*

*Proof.* By Lemma B.2c, at any transition time  $t_k^i$  when  $\mathbf{p} \notin \mathcal{B}$  and  $\|\mathbf{p}\| \leq \underline{c}$ , policy enters the basin by the next transition time  $t_k^{-i}$  or  $t_{k+1}^{-i}$  – in which case the subsequence terminates – with probability at least  $v_{\underline{c}} > 0$ . It follows immediately that the existence of an infinite subsequence of such transition times is a probability-zero event.  $\blacksquare$

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  $\tau_k^i$  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  $\Delta \tau_1^i, \Delta \tau_2^i, \dots$ .*

*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  $q(t)$ .  $\blacksquare$

**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  $\tau_k^i$  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.,

$$(B.4) \quad \lim_{T \rightarrow \infty} \frac{1}{T - T_0} \int_{T_0}^T v_\varepsilon(q(t)) dt = \int v_\varepsilon(q) dF(q).$$

Now, write

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

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  $\mathbb{E}[\Delta c_k^{i,\varepsilon}] / \mathbb{E}[\Delta \tau_k^i]$ , which is finite by Lemma B.3c. 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  $\lim_{k \rightarrow \infty} c_\varepsilon(\tau_k^i) = \infty$  a.s.. Further, for any  $\underline{c} < c_1^{i,\varepsilon}$ ,  $\inf\{c_\varepsilon(\tau_k^i)\} \geq \underline{c}$  with positive probability, and  $\lim_{c_\varepsilon(\tau_1^i) - \underline{c} \rightarrow \infty} \Pr[\inf\{c_\varepsilon(\tau_k^i)\} \geq \underline{c}] = 1$ .
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. Accordingly, pick  $\underline{c} < \infty$  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: there exists some transition time  $t_k^i$  where  $\|\mathbf{p}(t_k^i)\| > \underline{c}$ . Conditioning on this event, specify initializations  $c(t_k^i) = \|\mathbf{p}(t_k^i)\|$  and  $q(t_k^i) = p(t_k^i)$ . From Lemma B.4, with positive probability,  $\lim_{t \rightarrow \infty} c(t) = \infty$  and  $c(t) \geq c(t_k^i) > \underline{c}$  for all  $t \geq t_k^i$ . Consequently, applying Lemma B.2a: with positive probability,  $\lim_{t \rightarrow \infty} \|\mathbf{p}\|(t) = \infty$ . In other words,  $\kappa > 0$ .

Now, assume towards a contradiction that with positive probability, policy neither becomes simple nor kludged. By Lemma B.2d, there exists a complexity bound  $\underline{c} > 0$  such that with positive probability, there exists some infinite subsequence of transition times where  $\mathbf{p} \notin \mathcal{B}$  and  $\|\mathbf{p}\| \leq \underline{c}$ . But this contradicts Lemma B.2e.

$\int v(q)dF(q) < 0$ : Select sufficiently small  $\varepsilon$  and sufficiently large  $\underline{c}$  so that  $\mathcal{B} \subset \{\mathbf{p} : \|\mathbf{p}\| < \underline{c}\}$  and so that  $\varepsilon < 1 - \frac{\underline{c}-|p|}{\underline{c}+|p|}$  for  $p = p_{+1}^*$  and  $p = p_{-1}^*$ . Lemmas B.2b and B.4 together imply that if policy is above the complexity bound  $\underline{c}$  at some time  $t$ ,  $\|\mathbf{p}(t)\| > \underline{c}$ , then (a.s.)  $\|\mathbf{p}(t')\| \leq \underline{c}$  at some future time  $t' > t$ . Further, we may assume without loss that  $t'$  is a transition time. This then implies that (a.s.) there exists some infinite subsequence of transition times whereby for each time  $t$  in this subsequence,  $\|\mathbf{p}(t)\| \leq \underline{c}$ . Combined with Lemma B.2e, we conclude that (a.s.) policy is within the basin (and thus eventually becomes simple) during some transition time in this subsequence. ■

**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  $\nu \in (0, 1)$ ,

$$(B.5) \quad \lim_{c(t_0) \rightarrow \infty} \Pr \left[ c(\tau_1^{+1}) \geq (1 - \nu) c(t_0) \right] = 1,$$

$$(B.6) \quad \lim_{c(t_0) \rightarrow \infty} \Pr \left[ c(\tau_1^{-1}) \geq (1 - \nu) c(t_0) \right] = 1$$

WLOG suppose that policy hits +1's ideal first, at time  $\tau_1^i$ ; note that  $c(\tau_1^i) \geq c(t_0)$ . Notice that, subsequent to  $\tau_1^i$ , Party  $i$  loses control with arrival rate  $\lambda_i$ ; so  $c(\tau_1^i) - c(\tau_1^i)$  is exponentially distributed with parameter  $\lambda_i$ . Consequently, as  $c(t_0) \rightarrow \infty$ , the probability that  $c(\tau_1^i) - c(\tau_1^i) \geq \nu c(t_0)$  vanishes. Our claim (B.5) follows immediately. The demonstration of the claim (B.6) 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)  $\lim_{c(t_0) \rightarrow \infty} \Pr \left[ \inf_{t \geq t_0} c(t) \geq \underline{c} \right] = 1$ .  $\blacksquare$

#### 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 as  $\|\mathbf{p}(0)\| \rightarrow 0$ , 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  $\mathcal{B} \subset \{\mathbf{p} : \|\mathbf{p}\| < \underline{c}\}$ .  $\blacksquare$

## 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,

$$\begin{aligned} \int v(q, i) dF(q)(q, i) &= -(\Delta F(p_{+1}^*, +1) + \Delta F(p_{-1}^*, -1)) + \int_{p_{-1}^*}^{p_{+1}^*} (f(q, +1) + f(q, -1)) dq \\ (B.7) \quad &= \frac{-\gamma \left( \frac{1}{\lambda_{-1}} e^{\frac{\lambda_{+1}-\lambda_{-1}}{\gamma} p_{-1}^*} + \frac{1}{\lambda_{+1}} e^{\frac{\lambda_{+1}-\lambda_{-1}}{\gamma} p_{+1}^*} \right) + \int_{p_{-1}^*}^{p_{+1}^*} \left( 2e^{\frac{\lambda_{+1}-\lambda_{-1}}{\gamma} q} \right) dq}{\gamma \left( \frac{1}{\lambda_{-1}} e^{\frac{\lambda_{+1}-\lambda_{-1}}{\gamma} p_{-1}^*} + \frac{1}{\lambda_{+1}} e^{\frac{\lambda_{+1}-\lambda_{-1}}{\gamma} p_{+1}^*} \right) + \int_{p_{-1}^*}^{p_{+1}^*} \left( 2e^{\frac{\lambda_{+1}-\lambda_{-1}}{\gamma} q} \right) dq}. \end{aligned}$$

The denominator of the last expression (B.7) 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 (B.7) 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.  $\blacksquare$

### 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  $\kappa$  is weakly increasing, in zealouslyness  $\mathbf{z}$ . It remains to show that  $\kappa$  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  $\kappa$  of kludge is at least  $\rho$ . (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,  $\kappa$  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  $0 \leq q \leq q' < p_{+1}^*$ ,

$$(B.8) \quad \frac{f(q')}{f(q)} = \frac{e^{\lambda \frac{\Delta-1/\Delta}{\gamma} q'} + e^{-\lambda \frac{\Delta-1/\Delta}{\gamma} q'}}{e^{\lambda \frac{\Delta-1/\Delta}{\gamma} q} + e^{-\lambda \frac{\Delta-1/\Delta}{\gamma} q}} \text{ and } \frac{\Delta F(\Delta_p^*)}{\lim_{q \rightarrow \Delta_p^*} f(q)} = \frac{\frac{\gamma}{\lambda} \left( \Lambda e^{\lambda \frac{\Delta-1/\Delta}{\gamma} q'} + \frac{1}{\Lambda} e^{-\lambda \frac{\Delta-1/\Delta}{\gamma} q'} \right)}{\Lambda e^{\lambda \frac{\Delta-1/\Delta}{\gamma} q} + \frac{1}{\Lambda} e^{-\lambda \frac{\Delta-1/\Delta}{\gamma} q}}$$

are both increasing in  $\Lambda$ . That is,  $F$  satisfies the monotone-likelihood ratio property in  $\Lambda$ . Thus,  $F$  increases in the sense of first-order stochastic-dominance as  $\Lambda$  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 (B.8) simplifies further: for all  $0 \leq q \leq q' < p_{+1}^*$ ,

$$(B.9) \quad \frac{f(q')}{f(q)} = 1 \text{ and } \frac{\Delta F(p_{+1}^*)}{\lim_{q \rightarrow p_{+1}^*} f(q)} = \frac{\gamma}{\lambda}, \text{ so that}$$

$$(B.10) \quad F(q) = \begin{cases} \frac{q}{\Delta_p^* + \frac{\gamma}{\lambda}} & : p < \Delta_p^* \\ 1 & : p = \Delta_p^* \end{cases}.$$

By inspection,  $F$  increases in the sense of first-order stochastic-dominance as  $\gamma$  increases, as  $\lambda$  decreases, and as  $\Delta_p^*$  increases.  $\blacksquare$

## 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. We'll use both  $i$  and  $\ell$  to generically identify a Party.

**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  $i$  loses 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  $V_{i\ell}(p_0)$  be Party  $i$ 's expected payoff when Party  $\ell$  is in power and position equals  $p_0$ . Let  $T_\ell$  be the first time when Party  $\ell$  loses power to his opponent  $-\ell$ . Then

$$V_{i\ell}(p) = \mathbb{E} \left[ - \int_0^{T_\ell} e^{-r_\ell t} |g_\ell(t, p) - p_i^*| dt + e^{-r_\ell T_\ell} V_{i,-\ell}(g_\ell(T_\ell, p)) \right],$$

where  $g_\ell(t, p)$  evolves according to the law of motion

$$\frac{dg_\ell}{dt}(t, p) = a_\ell(g_\ell(t, p)),$$

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

$$(B.11) \quad V_{i\ell}(p) = \int_0^\infty [-|g_\ell(t, p) - p_i^*| + \lambda_\ell V_{i,-\ell}(g_\ell(t, p))] e^{-(r_\ell + \lambda_\ell)t} dt, \text{ for every } p_0 \in \mathbb{R}.$$

The Bellman equation associated with this integral is<sup>1</sup>

$$(B.12) \quad -|p_0 - p_i^*| + \lambda_\ell V_{i,-\ell}(p_0) - (r_i + \lambda_\ell)V_{i\ell}(p_0) + V_{i\ell}'(p_0)a_\ell(p_0) = 0, \text{ for every } p_0 \in \mathbb{R}.$$

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 a Markov Perfect Equilibrium. Then  $a_{+1}(p) = \gamma$  when  $p < p_{+1}^{**}$  and  $a_{-1}(p) = -\gamma$  when  $p > p_{-1}^{**}$ . Therefore, Eq. (B.12) implies that

$$(B.13) \quad \gamma V_{i,+1}'(p) = |p - p_i^*| - \lambda_{+1}V_{i,-1}(p) + (r_i + \lambda_{+1})V_{i,+1}(p), \text{ for } p < p_{+1}^{**};$$

$$(B.14) \quad \gamma V_{i,-1}'(p) = -|p - p_i^*| + \lambda_{-1}V_{i,-1}(p) - (r_i + \lambda_{-1})V_{i,-1}(p), \text{ for } p > p_{-1}^{**};$$

$$(B.15) \quad 0 = |p_\ell^{**} - p_i^*| + \lambda_\ell V_{i,-\ell}(p_\ell^{**}) - (r_i + \lambda_\ell)V_{i\ell}(p_\ell^{**}).$$

In equilibrium,  $V_{ii}'(p)a_i(p) \geq 0$  for every  $p$ . Therefore, Eq. (B.12) implies that

$$(B.16) \quad |p - p_i^*| - \lambda_i V_{i,-i}(p) + (r_i + \lambda_i)V_{ii}(p) = V_{ii}'(p)a_i(p) \geq 0 \text{ for every } p \in \mathbb{R}.$$

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  $\gamma$ ; assuming that  $V_{ii}'(p_i^{**}) < 0$  leads to a similar contradiction.) Differentiating the left hand side of Eq. (B.16) at  $p_i^{**}$  yields that

$$(B.17) \quad V_{i,-i}'(p_i^{**}) \begin{cases} = -\lambda_i^{-1}, & \text{if } p_i^{**} < p_i^*; \\ \in [-\lambda_i^{-1}, \lambda_i^{-1}], & \text{if } p_i^{**} = p_i^*; \\ = \lambda_i^{-1}, & \text{if } p_i^{**} > p_i^*. \end{cases}$$

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 Markov 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

<sup>1</sup>Formally, the Bellman equation can be derived as follows: replace  $p$  in Eq. (B.11) with  $g_\ell(s, p)$  and  $g_\ell(t, p)$  with  $g_\ell(s + t, p)$  and rewrite Eq. (B.11) as  $V_{i\ell}(g_\ell(s, p))e^{-(r_i + \lambda_\ell)s} = \int_s^\infty [-|g_\ell(\tau, p) - p_i^*| + \lambda_\ell V_{i,-\ell}(g_\ell(\tau, p))]e^{-(r_i + \lambda_\ell)\tau} d\tau$  where  $\tau = s + t$ . Differentiating both sides with respect to  $s$  at  $s = 0$  yields the Bellman equation.

be introduced shortly. Define

(B.18)

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

(B.19)

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

(B.20)

$$1_{+1} = \begin{pmatrix} 1 \\ 0 \end{pmatrix};$$

(B.21)

$$1_{-1} = \begin{pmatrix} 0 \\ 1 \end{pmatrix};$$

(B.22)

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

(B.23)

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

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

**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  $\eta$ . 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  $\eta$ . 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

$$(B.24) \quad \mu_{i\pm} = \frac{1}{2} \left[ \lambda_{+1} - \lambda_{-1} \pm \sqrt{(\lambda_{-1} - \lambda_{+1})^2 + 4r_i^2 + 4(\lambda_{-1} + \lambda_{+1})r_i} \right]$$

are the eigenvalues of  $A_i$ . Note that  $\mu_{i+} > 0 > \mu_{i-}$  for  $i \in \{+1, -1\}$ . To avoid confusion, the eigenvalue  $\mu_{i+}$  when  $i = +1$  will be written as  $\mu_{++}$  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

$$(B.25) \quad \xi_{+\pm}(p, p') = \gamma\mu_{+1,\pm}^{-1} \left[ 1 + e^{\gamma^{-1}(p-p')\mu_{+\pm}} - 2e^{\gamma^{-1}(p-p_{+1}^*)\mu_{+\pm}} \right];$$

$$(B.26) \quad \xi_{-\pm}(p, p') = -\gamma\mu_{-1,\pm}^{-1} \left[ 1 + e^{\gamma^{-1}(p-p')\mu_{-\pm}} - 2e^{\gamma^{-1}(p-p_{-1}^*)\mu_{-\pm}} \right].$$

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

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

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

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. (B.25) 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  $p' \geq p_{+1}^*$ ,

$$(B.27) \quad H_{+1,2}(p, p', 1) = (\mu_{++} - \mu_{+-})^{-1} \left( \zeta_{+-} e^{-\gamma^{-1}(p'-p)\mu_{+-}} - \zeta_{++} e^{-\gamma^{-1}(p'-p)\mu_{++}} \right),$$

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.  $\blacksquare$

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. (B.13) and (B.14) 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. (B.18). If  $\vec{V}_i(p')$  is known, the solution to this differential equation is

$$(B.28) \quad \vec{V}_i(p) = e^{\gamma^{-1}(p-p')A_i} \vec{V}_i(p') + e^{\gamma^{-1}pA_i} [L_i(p) - L_i(p')], \text{ for every } p \in [p_{-1}^{**}, p_{+1}^{**}],$$

where  $L_i$  is defined in Eq. (B.19). Eqs. (B.15) (for the case where  $\ell = i$ ) and (B.17) 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. (B.28) 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. (B.15) for the case where  $i \neq \ell$  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. (B.22) and (B.23). 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. (B.13)-(B.17) will be satisfied, implying that the focused strategy profile is a Markov Perfect Equilibrium.  $\blacksquare$

**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. (B.25), 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  $\infty$  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  $\eta$  such that  $H_{+1}(p, p_{+1}^*, \eta) = 0$  when  $BR_{+1}(p) = p_{+1}^*$ . Then

$$(B.29) \quad H_{+1}(p, BR_{+1}(p), \eta_{+1}(p)) = 0, \text{ for every } p \leq p_{-1}^*.$$

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

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

Then

$$(B.30) \quad \vec{U}'_{+1}(\tilde{p}; p', \eta) = \gamma^{-1}A_{+1}\vec{U}_{+1}(\tilde{p}; p', \eta) + \gamma^{-1}|\tilde{p} - p_{+1}^*| \begin{pmatrix} 1 \\ -1 \end{pmatrix}, \text{ for every } \tilde{p} \in \mathbb{R}.$$

$$(B.31) \quad H_{+1}(p, p', \eta) = 1_{-1}^\top A_{+1}\vec{U}_{+1}(p; p', \eta) - |p - p_{+1}^*|.$$

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

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



Combining this result with Eq. (B.29) yields that

$$(B.32) \quad H_{+1,1}(p, BR_{+1}(p), \eta_{+1}(p)) = 1 + \lambda_{-1} \gamma^{-1} \mathbf{1}_{+1}^\top \left[ A_{+1} \vec{U}_{+1}(p; p', \eta) + |p - p_{+1}^*| \begin{pmatrix} 1 \\ -1 \end{pmatrix} \right].$$

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. (B.30) and the same boundary conditions at  $p$  and  $BR_{+1}(p)$ . Therefore, on  $[p, BR_{+1}(p)]$ ,  $\vec{U}_{+1}(\cdot; BR_{+1}(p), \eta_{+1}(p))$  coincides with Party +1's value function. Because Party +1's flow payoff is always non-positive,

$$(B.33) \quad U_{+1,+1}(p; BR_{+1}(p), \eta_{+1}(p)) \leq 0.$$

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,

$$(B.34) \quad U_{+1,+1}(p; BR_{+1}(p), \eta_{+1}(p)) \geq -\frac{1}{r_{+1}}|p - p_{+1}^*|.$$

By Eq. (B.31), that  $H_{+1}(p, BR_{+1}(p), \eta_{+1}(p)) = 0$  implies that

$$(B.35) \quad \lambda_{-1} U_{+1,+1}(p; BR_{+1}(p), \eta_{+1}(p)) - (r_{+1} + \lambda_{-1}) U_{+1,-1}(p; BR_{+1}(p), \eta_{+1}(p)) - |p - p_{+1}^*| = 0.$$

Using Eq. (B.35) to eliminate  $U_{+1,-1}(p; BR_{+1}(p), \eta_{+1}(p))$  from Eq. (B.32) 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. (B.33) and (B.34) yields that

$$(B.36) \quad 1 \leq H_{+1,1}(p, BR_{+1}(p), \eta_{+1}(p)) \leq 1 + \frac{\lambda_{-1}(r_{+1} + \lambda_{-1} + \lambda_{+1})}{\gamma(r_{+1} + \lambda_{-1})} |p - p_{+1}^*|.$$

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 ( $\eta$ ) 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,

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

■

**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. (B.27), 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. (B.36) 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  $\gamma$  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  $H_{-1}(p, p_{-1}^*, r_{-1}) < 0$  for every  $p \geq p_{+1}^*$  and  $r_{-1} > r_{-1,c}$ ; if  $H_{-1}(p, p_{-1}^*, r_{-1}) = 0$  for some  $p \geq p_{+1}^*$  and  $r_{-1} \leq r_{-1,c}$ , then  $\frac{\partial H_{-1}}{\partial r_{-1}}(p, p_{-1}^*, r_{-1}) < 0$ .*

*Proof.* We only prove the assertion on  $H_{+1}$ . In this proof, the dependence of  $\mu_{++}$  and  $\mu_{+-}$  on  $r_{+1}$  will be made explicit. By Eq. (B.25),  $\xi_{++}(p, p_{+1}^*; r_{+1})$  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  $\mu_{\pm\pm}(\tilde{r}_{+1})$ :

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

(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  $\gamma$  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  $\mu_{+-}$ . Its limit when  $\mu_{+-} \rightarrow 0$  is  $\frac{1}{2}\gamma^{-1}(p_{+1}^* - p)$ . Therefore,

$$(B.37) \quad c(r_{+1}) - b(r_{+1}) > \frac{(1 + m_{+1}(r_{+1}))\lambda_{-1}}{(r_{+1} + 2\lambda_{-1} + \mu_{++}(r_{+1}))(r_{+1} + \lambda_{-1} + \mu_{++}(r_{+1}))} + \frac{1}{2}\gamma^{-1}(p_{+1}^* - p).$$

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

$$(B.38) \quad c(r_{+1}) - a(r_{+1}) < c(r_{+1}) - d(r_{+1}) = \left[ \frac{1 + m_{+1}(r_{+1})}{r_{+1} + \lambda_{-1} + \mu_{++}(r_{+1})} - \frac{1 - m_{+1}(r_{+1})}{r_{+1} + \lambda_{-1} + \mu_{+-}(r_{+1})} \right] + 2\gamma^{-1}m_{+1}(r_{+1})(p_{+1}^* - p).$$

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

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

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

$$\begin{aligned} & \frac{(\lambda_{-1} - \lambda_{+1})(p_{+1}^* - p)}{\gamma \left( e^{\gamma^{-1}(p_{+1}^* - p)(\lambda_{-1} - \lambda_{+1})} - 1 \right)} \text{ if } \lambda_{-1} > \lambda_{+1}, \\ & \frac{\gamma \left( 1 - e^{-\gamma^{-1}(p_{+1}^* - p)(\lambda_{+1} - \lambda_{-1})} \right)}{(\lambda_{+1} - \lambda_{-1})(p_{+1}^* - p)} \text{ if } \lambda_{+1} > \lambda_{-1}. \end{aligned}$$

Either way, the ratio approaches zero at least as fast as  $\gamma$  as  $\gamma \rightarrow 0$ . Therefore, there exists a  $\eta > 0$  such that

$$(B.39) \quad \frac{A(r_{+1})}{B(r_{+1})} < \eta\gamma,$$

for all  $\gamma \leq \bar{\gamma}$  and  $r_{+1} \leq r_{+1,c}$ . Combining Eqs. (B.37)-(B.39) 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. (B.37) and  $E_2(r_{+1})$  is the term in the bracket on the right hand side of Eq. (B.38). 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  $\frac{\partial H_{+1}}{\partial r_{+1}}(p, p_{+1}^*, 1; r_{+1}) > 0$ . ■

### 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  $\gamma$  is sufficiently small, the sign of  $H_{+1}(p, p_{+1}^*, 1)$  is the same as the sign of

$$(B.40) \quad (r_{+1} + 2\lambda_{-1} + \mu_{++}(r_{+1}))\mu_{+-}(r_{+1})^{-1} + (r_{+1} + \lambda_{-1} + \mu_{++}(r_{+1}))\lambda_{-1}^{-1}.$$

According to Eq. (B.24), 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. (B.40) approaches  $-\infty$  as  $r_{+1} \rightarrow 0$ , and when  $\lambda_{+1} < \lambda_{-1}$ , the expression in Eq. (B.40) 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.  $\blacksquare$

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

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

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

*Proof.* By Eq. (B.24), 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  $r_{+1}, r_{-1} \rightarrow 0$ .  $\blacksquare$

## Reducing Complexity by Trigger Strategies

We retain the purely positional preferences from Section III: parties' flow payoffs are  $u_i(\mathbf{p}(t)) = -|p_i^* - p(t)|$ . Denote the sum of flow payoffs at position  $p$  as

$$w(p) = -|p - p_{+1}^*| - |p - p_{-1}^*| = \begin{cases} -\Delta_p^* - 2(p_{-1}^* - p), & \text{if } p < p_{-1}^*; \\ -\Delta_p^*, & \text{if } p_{-1}^* \leq p \leq p_{+1}^*; \\ -\Delta_p^* - 2(p - p_{+1}^*), & \text{if } p > p_{+1}^*. \end{cases}$$

We assume that both parties have common discount rate  $r$ . Given a focused Markov equilibrium, let  $W(p_0, i) = V_{+1,i}(p_0) + V_{-1,i}(p_0)$  be the sum of the two parties' value functions when the initial state is  $(p_0, i)$ . We maintain previous notation and use  $p_{-1}^{**}$  and  $p_{+1}^{**}$  to reference focused Markov equilibrium targets.

We maintain the assumption that  $\gamma < \bar{\gamma}$  as in Proposition 9. This assumption ensures that there exists  $\hat{r} > 0$  such that when  $r \leq \hat{r}$ , a focused Markov equilibrium exists (uniquely) and exhibits strategic extremism; in other words,  $p_{+1}^{**} > p_{+1}^*$  or  $p_{-1}^{**} < p_{-1}^*$ . (Uniqueness is convenient but not necessary for our results.)

**Lemma B.7a.** *Suppose  $\gamma < \bar{\gamma}$ . For some constant  $\varrho > 0$ , for all discount rates  $r \leq \hat{r}$ ,*

$$rW(p_0, i) \leq -\Delta_p^* - \varrho(\Delta_p^{**} - \Delta_p^*).$$

*Proof.* Suppose  $p(0) = p_0$  and  $i(0) = i$ . WLOG, assume that  $p_{+1}^{**} - p_{+1}^* \geq p_{-1}^* - p_{-1}^{**}$ . Let  $p_{-1}^{***} = p_{-1}^* - (p_{+1}^{**} - p_{+1}^*)$ ; note that  $p_{-1}^{***} \leq p_{-1}^{**}$  by assumption.

First, suppose  $p_0 < p_{-1}^{***}$ . No matter which party is in control, position  $p(t)$  increases from  $p_0$  to  $p_{-1}^{***}$  at some time  $t_0$ ; once position reaches  $p_{-1}^{***}$ , it stays forever within the interval  $[p_{-1}^{***}, p_{+1}^{**}]$ . Further, notice that  $w(p)$  is strictly increasing on  $[p_0, p_{-1}^{***}]$  and weakly increasing on  $[p_{-1}^{***}, p_{+1}^{**}]$ ; that is, the total flow payoff  $w(p(t))$  is weakly higher (strictly lower) at any time  $t \geq t_0$  ( $t < t_0$ ) than at  $t_0$ . Given that  $W(p(t), i(t))$  is a weighted mean of future flow payoffs, our observations imply that for  $p_0 < p_{-1}^{***}$ ,

$$W(p_0, i) < \mathbb{E}[W(p(t_0), i(t_0))] < \max\{W(p_{-1}^{***}, +1), W(p_{-1}^{***}, -1)\}.$$

Second, suppose  $p_0 \geq p_{-1}^{***}$ . Let  $t_{+1} = \gamma^{-1}(p_{+1}^{**} - p_{-1}^{***})$  be the amount of time taken for  $p(t)$  to travel from  $p_{-1}^{***}$  to  $p_{+1}^{**}$ . Fix any  $t \geq t_{+1} + 1$ . A moment's reflection reveals that there is some  $q_0 > 0$  (independent of  $r$  and  $t$ ) such that Party +1 is in control at time  $t - t_{+1}$  with probability of at least  $q_0$ . Conditional on this event, with probability  $e^{-\lambda_+ t_{+1}}$ , Party +1 remains in control until time  $t$ , in which case  $p(t) = p_{+1}^{**}$ . Combining these observations,  $p(t) = p_{+1}^{**}$  with probability  $\geq q_0 \cdot e^{-\lambda_+ t_{+1}}$ . Consequently,

$$\mathbb{E}[w(p(t))] \leq -\Delta_p^* - q_0 \cdot e^{-\lambda_+ t_{+1}} (p_{+1}^{**} - p_{+1}^*) \text{ for all } t \geq t_{+1} + 1.$$

Further,  $\mathbb{E}[w(p(t))] \leq -\Delta_p^*$  for  $t < t_{+1}$ . Combining these last two inequalities,

$$\begin{aligned} rW(p_0, +1) &\leq r \int_0^{t_{+1}+1} -\Delta_p^* e^{-rt} dt + r \int_{t_{+1}+1}^{\infty} [-\Delta_p^* - q_{+1}(p_{+1}^{**} - p_{+1}^*)] e^{-rt} dt \\ &= -\Delta_p^* - e^{-r(t_{+1}+1)} q_{+1}(p_{+1}^{**} - p_{+1}^*) \\ &= -\Delta_p^* - e^{-\hat{r}(t_{+1}+1)} q_{+1}(p_{+1}^{**} - p_{+1}^*) \\ &\leq -\Delta_p^* - e^{-\hat{r}(t_{+1}+1)} q_{+1} (\Delta_p^{**} - \Delta_p^*) / 2. \end{aligned}$$

for every  $p_0 \geq p_{-1}^{***}$ . In other words, the lemma holds with  $\varrho = \frac{1}{2} e^{-\hat{r}(t_{+1}+1)} q_{+1}$ .  $\blacksquare$

**Lemma B.7b.** Suppose  $\gamma < \bar{\gamma}$ . For some constants  $\tilde{q} > 0$  and  $\tilde{r} > 0$ , for all discount rates  $r \leq \tilde{r}$  and for all initial states  $(p_0, i)$ ,

$$rW(p_0, i) \leq -\Delta_p^* - \tilde{q}.$$

*Proof.* The proof of Proposition 9 implies that  $|p_{\pm 1}^{**} - p_{\pm 1}^*|$  is bounded away from zero as  $r \rightarrow 0$  if  $\lambda_{\pm 1} > \frac{1}{3}\lambda_{\mp 1}$ , so  $\Delta_p^*$  is bounded away from zero as  $r \rightarrow 0$  for any  $\lambda_{-1}$  and any  $\lambda_{+1}$ . This, combined with Lemma B.7a, proves our result.  $\blacksquare$

**Lemma B.8a.** Suppose  $\gamma < \bar{\gamma}$ , and choose  $\tilde{r}$  as in Lemma B.7b. There exists  $\Delta_p > 0$  such that  $\Delta^{**} \leq \Delta_p$  for all  $r \leq \tilde{r}$ .

*Proof.* It suffices to show that  $p_{+1}^{**} - p_{-1}^{**}$  remains bounded as  $r \rightarrow 0$ . Consider for now the case  $\lambda_{+1} > \lambda_{-1}$ . Suppose, towards a contradiction, that there exists a sequence  $r_n \rightarrow 0$  such that the corresponding sequence  $p_{+1,n}^{**} - p_{-1,n}^{**} \rightarrow \infty$ . Tedious but straightforward calculations reveal that as  $r_n \rightarrow 0$ , the corresponding sequences  $H_{i,n}$  behave as follows:

$$\begin{aligned} H_{+1,n}(p_{-1,n}^{**}, p_{+1,n}^{**}, 1) &\sim \frac{2\lambda_{-1}\gamma}{(\lambda_{+1} - \lambda_{-1})^2} \left[ 1 - 2e^{-\gamma^{-1}(p_{+1,n}^{**} - p_{-1,n}^{**})(\lambda_{+1} - \lambda_{-1})} \right] - \frac{\lambda_{+1} + \lambda_{-1}}{\lambda_{+1} - \lambda_{-1}} \gamma (2p_{+1,n}^{**} - p_{-1,n}^{**} - p_{+1,n}^{**}) + \\ &\quad + \frac{\gamma}{\lambda_{+1} - \lambda_{-1}} e^{-\gamma^{-1}(p_{+1,n}^{**} - p_{-1,n}^{**})\mu_{+-n}}, \\ H_{-1,n}(p_{-1,n}^{**}, p_{+1,n}^{**}, 1) &\sim -\frac{\gamma(\lambda_{+1} + \lambda_{-1})}{\lambda_{+1} - \lambda_{-1}} e^{\gamma^{-1}(p_{+1,n}^{**} - p_{-1,n}^{**})(\lambda_{+1} - \lambda_{-1})} + \frac{4\gamma\lambda_{+1}}{(\lambda_{+1} - \lambda_{-1})^2} e^{\gamma^{-1}(p_{+1,n}^{**} - p_{-1,n}^{**})(\lambda_{+1} - \lambda_{-1})}. \end{aligned}$$

In equilibrium (with strategic extremism),  $H_i(p_{-i}^{**}, p_i^{**}, 1) = 0$ . The asymptotic behavior of  $H_{-1,n}(p_{-1,n}^{**}, p_{+1,n}^{**}, 1)$  implies that  $p_{-1,n}^{**}$  remains bounded, as otherwise  $H_{-1,n}(p_{-1,n}^{**}, p_{+1,n}^{**}) \rightarrow -\infty$ . However, this implies that  $p_{+1,n}^{**} \rightarrow \infty$  and thus  $H_{+1,n}(p_{-1,n}^{**}, p_{+1,n}^{**}, 1) \rightarrow \infty$ , regardless of the asymptotic behavior of  $(p_{+1,n}^{**} - p_{-1,n}^{**})r_n$ .

A similar argument by contradiction holds in the case  $\lambda_{+1} < \lambda_{-1}$ .  $\blacksquare$

**Lemma B.8b.** Suppose  $\gamma < \bar{\gamma}$ , and choose  $\tilde{r}$  as in Lemma B.7b. There exist constants  $V_{+1}, V_{-1}$  and  $\Delta_V$  (with  $\Delta_V > 0$ ) such that for all  $r \leq \tilde{r}$ ,

$$|rV_{i\ell}(p) - V_i| < r\Delta_V, \text{ for every } i, \ell, \text{ and } p \in [p_{-1}^{**}, p_{+1}^{**}].$$

*Proof.* The equilibrium condition that  $H_{+1}(p_{-1}^{**}, p_{+1}^{**}, 1) = 0$  can be rewritten as

$$(B.41) \quad 1_{-1}^\top e^{\gamma^{-1}(p_{-1}^{**} - p_{+1}^{**})A_{+1}} A_{+1} r \vec{V}_{+1}(p_{+1}^{**}) + 1_{-1}^\top e^{\gamma^{-1}p_{-1}^{**}A_{+1}} A_{+1} r [L_{+1}(p_{-1}^{**}) - L_{+1}(p_{+1}^{**})] - r |p_{-1}^{**} - p_{+1}^{**}| = 0.$$

In equilibrium, Party +1's flow payoff is negative unless  $p = p_{+1}^*$ , but  $p \neq p_{+1}^*$  almost all the time, so  $r\vec{V}_{+1}(p_{+1}^{**})$  is bounded away from zero. On the other hand, Lemma B.8a ensures that  $p_{+1}^{**}$  and  $p_{-1}^{**}$  remain bounded as  $r \rightarrow 0$ . So, the second and third terms on the left hand side of Equation (B.41) are of order  $O(r)$ . Diagonalizing  $A_{+1}$  as

$$A_{+1} = \Lambda_{+1} \begin{pmatrix} \mu_{++} & \\ & \mu_{--} \end{pmatrix} (\Lambda_{+1})^{-1} \text{ where } \Lambda_{+1} = \begin{pmatrix} \mu_{++} + r + \lambda_{-1} & \mu_{+-} + r + \lambda_{-1} \\ & \lambda_{-1} \end{pmatrix},$$

we can reduce the first term on the left hand side of Equation (B.41) to

$$1_{-1}^\top e^{\gamma^{-1}(p_{-1}^{**} - p_{+1}^{**})\Lambda_{+1}} A_{+1} r \vec{V}_{+1}(p_{+1}^{**}) = 1_{-1}^\top \Lambda_{+1} \begin{pmatrix} \mu_{++} e^{\gamma^{-1}(p_{-1}^{**} - p_{+1}^{**})\mu_{++}} & \\ & \mu_{+-} e^{\gamma^{-1}(p_{-1}^{**} - p_{+1}^{**})\mu_{+-}} \end{pmatrix} (\Lambda_{+1})^{-1} r \vec{V}_{+1}(p_{+1}^{**}).$$



Therefore,

$$(B.42) \quad \mathbf{1}_{-1}^\top \Lambda_{+1} \begin{pmatrix} \mu_{++} e^{\gamma^{-1}(p_{-1}^{**} - p_{+1}^{**})\mu_{++}} \\ \mu_{+-} e^{\gamma^{-1}(p_{-1}^{**} - p_{+1}^{**})\mu_{+-}} \end{pmatrix} (\Lambda_{+1})^{-1} r \vec{V}_{+1}(p_{+1}^{**}) = O(r).$$

On the other hand,

$$(B.43) \quad r \vec{V}_{+1}(p) = e^{\gamma^{-1}(p - p_{+1}^{**})A_{+1}} r \vec{V}_{+1}(p_{+1}^{**}) + e^{\gamma^{-1}pA_{+1}} [L_{+1}(p) - L_{+1}(p_{+1}^{**})].$$

Diagonalization of  $A_{+1}$  reduces Equation (B.43) to the following:

$$(B.44) \quad r \vec{V}_{+1}(p) = \Lambda_{+1} \begin{pmatrix} e^{\gamma^{-1}(p - p_{+1}^{**})\mu_{++}} \\ e^{\gamma^{-1}(p - p_{+1}^{**})\mu_{+-}} \end{pmatrix} (\Lambda_{+1})^{-1} r \vec{V}_{+1}(p_{+1}^{**}) + O(r), \text{ for } p \in [p_{-1}^{**}, p_{+1}^{**}].$$

There are two cases:  $\lambda_{+1} > \lambda_{-1}$  and  $\lambda_{+1} < \lambda_{-1}$ . Consider the case  $\lambda_{+1} > \lambda_{-1}$ . In this case,  $\mu_{++} = \lambda_{+1} - \lambda_{-1} + O(r)$ , while  $\mu_{+-} = -\frac{\lambda_{+1} + \lambda_{-1}}{\lambda_{+1} - \lambda_{-1}} r + O(r^2)$ . It is straightforward to verify that all the entries of  $\Lambda_{+1}$  and  $(\Lambda_{+1})^{-1}$  converge to positive numbers as  $r \rightarrow 0$ . By Equation (B.42), the first component of  $(\Lambda_{+1})^{-1} r \vec{V}_{+1}(p_{+1}^{**})$  must be of the order  $O(r)$ . Substituting this fact into Eq. (B.44), we conclude that the amplitude of  $r \vec{V}_{+1}(p)$  on  $[p_{-1}^{**}, p_{+1}^{**}]$  is of the order  $O(r)$ . The case  $\lambda_{+1} < \lambda_{-1}$  proceeds similarly, albeit with the modifications  $\mu_{++} = \frac{\lambda_{+1} + \lambda_{-1}}{\lambda_{+1} - \lambda_{-1}} r + O(r^2)$  and  $\mu_{+-} = \lambda_{+1} - \lambda_{-1} + O(r)$ .

A symmetric argument applies to  $\vec{V}_{-1}$ . ■

### Proof of Proposition 10

Recall that  $W(p, i) = V_{+1,i}(p) + V_{-1,i}(p)$  is the sum of value functions under the unique focused Markov equilibrium. Continue to denote the parties' targets under the focused Markov equilibrium as  $p_{-1}^{**}$  and  $p_{+1}^{**}$ . Combining Lemmas B.7b and B.8b, we may ensure that for sufficiently small  $r$ , there exists a regular position  $\tilde{p}^{**} \in [p_{-1}^*, p_{+1}^*]$  and positive numbers  $\varrho_{+1}$  and  $\varrho_{-1}$  such that

$$(B.45) \quad r V_{i\ell}(p) \leq -|p_i^* - \tilde{p}^{**}| - \varrho_i, \text{ for } i, \ell = \pm 1, \text{ and } p \in [p_{-1}^*, p_{+1}^*].$$

We now show that a focused trigger-strategy profile with common target  $\tilde{p}^{**}$  and punishment targets  $\hat{p}_{-1}^{**} = p_{-1}^*$  and  $\hat{p}_{+1}^{**} = p_{+1}^*$  is an equilibrium. Let  $\tilde{V}_{i\ell}(\mathbf{p})$  be Party  $i$ 's value function under the focused trigger-strategy profile if Party  $\ell$  is in control initially, no party has previously deviated, and the initial policy is  $\mathbf{p}$ . By construction, following any deviation, the continuation equilibrium coincides with the unique focused Markov equilibrium, and thus the continuation value for each Party  $i$  given state  $(\ell, p)$  equals  $V_{i\ell}(p)$ .

On the equilibrium path, given any initial policy position  $p_0 \in [\hat{p}_{-1}^{**}, \hat{p}_{+1}^{**}]$ , policy reaches  $\tilde{p}^{**}$  within time  $T^{**} = \gamma \Delta^{**}$  and stays on  $\tilde{p}^{**}$  thereafter. Then

$$(B.46) \quad r \tilde{V}_{i\ell}(\mathbf{p}) \geq -\Delta^{**} (1 - e^{-rT^{**}}) - e^{-rT^{**}} |p_i^* - \tilde{p}^{**}|.$$

Combining Eqs. (B.45) and (B.46), we see that there exists an  $\underline{r} > 0$  such that for  $r \leq \underline{r}$ ,

$$r \tilde{V}_{i\ell}(\mathbf{p}) \geq r V_{i\ell}(p) + \frac{1}{2} \varrho_i, \text{ for every } i, \ell, \text{ and } \mathbf{p} \text{ such that } p \in [p_{-1}^*, p_{+1}^*];$$

that is, neither party prefers to deviate from the trigger strategy provided that no party deviated before and  $p \in [p_{-1}^{**}, p_{+1}^{**}]$ . If either party has deviated before, the parties are simply playing a focused strategy equilibrium, so neither party has incentive to deviate. Finally, when  $p \notin [p_{-1}^{**}, p_{+1}^{**}]$ , the focused trigger-strategy profile coincides with the focused Markov profile, and neither party has incentive to deviate. We conclude that the focused trigger-strategy profile is a subgame-perfect equilibrium. ■

**An Aside: Asymptotic Simplicity** As discussed in footnote 17, given our construction of the focused trigger-strategy equilibrium, policy asymptotically approaches – but never attains – perfect simplicity. Here, we modify this construction slightly to ensure that perfect simplicity is always attained within finite time on the equilibrium path. Define the *complexity threshold* to be  $\overline{\|\mathbf{p}\|} = \tilde{p}^{**} + 1$ .

Behavior following any deviation remains entirely unmodified: each Party  $i$  focuses on his punishment target  $\hat{p}_i^{**} = p_i^{**}$ . Prior to any deviation, behavior above the complexity threshold ( $\|\mathbf{p}\| > \overline{\|\mathbf{p}\|}$ ) also remains unmodified: both parties focus on the common target  $\tilde{p}^{**}$ .

Modify behavior below the complexity threshold ( $\|\mathbf{p}\| \leq \overline{\|\mathbf{p}\|}$ ), and prior to any deviation, as follows. If policy is perfectly simple and all existing rules have direction  $j = \text{sgn}(\tilde{p}^{**})$ , then each party adds or removes  $j$ -rules until he attains the common target position  $\tilde{p}^{**}$ , and subsequently stays there forever. Otherwise, each party reduces complexity as quickly as possible ( $\delta = \gamma$ ), until he attains the empty policy. He then adds  $j$ -rules ( $\alpha_j = \gamma$ ) until he attains the common target position  $\tilde{p}^{**}$ , and subsequently stays there forever. Figures B.1a and B.1b illustrate pre-deviation behavior in the modified equilibrium and the (unmodified) focused trigger-strategy equilibrium.

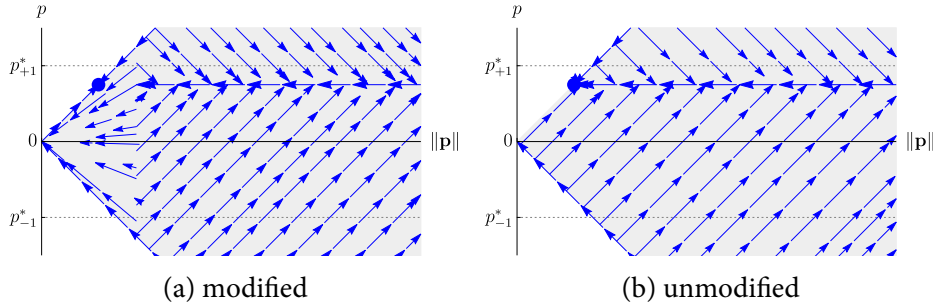


Figure B.1: Trigger-Strategy Equilibrium Path: Modified vs. Unmodified

With this modification, from any initial policy, the perfectly simple policy with common target position  $\tilde{p}^{**}$  is attained in finite time on the equilibrium path. Relative to the unmodified equilibrium, policy may spend an additional time period of up to  $\frac{\overline{\|\mathbf{p}\|} + \tilde{p}^{**}}{\gamma}$  away from the common target position (while respecting the positional bound  $p \in [p_{-1}^*, p_{+1}^*]$ ). It follows that the total time spent away from the common target position still remains bounded. Proposition 10 thus holds for the modified equilibrium as well.

## References

- Azema, J., Kaplan-Duflo, M., & Revuz, D. (1967). "Mesure invariante sur les classes récurrentes des processus de Markov." *Zeitschrift für Wahrscheinlichkeitstheorie und Verwandte Gebiete*, 8(3), 157–181.
- Hairer, M. (2008). "Ergodic theory for stochastic PDEs." *Preprint*.
- Kaspi, H. & Mandelbaum, A. (1994). "On Harris Recurrence in Continuous Time." *Mathematics of Operations Research*, 19(1), 211–222.
- Meyn, S. P. & Tweedie, R. L. (1993). "Stability of Markovian Processes II: Continuous-Time Processes and Sampled Chains." *Advances in Applied Probability*, 25(3), 487–517.