

# Flip-flopping and Electoral Concerns \*

Giovanni Andreottola \*

December 2018

## Abstract

Politicians who change their mind on a policy issue are often confronted with the accusation of being flip-floppers. However, a changing environment sometimes makes policy revisions necessary. The present analysis suggests that flip-flopping signals that politicians are poorly informed and is therefore detrimental to their reputation. As a result, electorally concerned politicians can have an incentive to keep their initial policy choice unchanged, despite it being inefficient, in order to avoid the stigma of flip-flopping. This distorted behaviour is not only damaging in terms of policy welfare, but also in terms of a worse selection of competent politicians through elections. Variations of the baseline model are used to provide an in-depth discussion of several possible ways to address the unwillingness of politicians to respond to information: these include term limits, the presence of media and the partial delegation of actions to independent agents.

*Keywords:* flip-flopping; elections; political agency; accountability; reputation; media; transparency; delegation

---

\*I am very grateful to my advisors Andrea Mattozzi and David Levine for their great supervision during my PhD at the EUI, and to Larry Samuelson for hosting me as a postdoc at the Cowles Foundation in Yale. The paper also benefitted from useful discussions with Ran Eilat, Piero Gottardi, Andrea Galeotti, Ronnie Razin, Alessandro Riboni and several others. I would also like to thank audiences at the Leitner Political Economy Seminar in Yale, the EUI Microeconomics Working Group, the EUI Workshop on Media Effects, the University of Bologna, the University of Naples, the Max Planck Institute in Munich, the NICEP Workshop in Nottingham, the IOEA in Cargese and the SIOE conference in Montréal for their questions and comments.

\*Cowles Foundation, Yale University.  
Email: Giovanni.Andreottola@yale.edu Website: [www.giovanniandreottola.com](http://www.giovanniandreottola.com)

*But with Kerry the charge isn't that he's inconstant. It's that in his inconstancy he flips wrong – the far more serious charge of bad judgment.*

Mickey Kaus (Slate)

*Bush's decision-making style was based on his gut instincts [...] Bush was quick to reach decisions, and, once reached, he saw change as a sign of weakness.*

(Newsweek)

## 1 Introduction

Consistency is one of the qualities that voters value the most in a politician. For example, the political scientist James Fearon (1999) writes: “*If I think of elections as a problem of choosing a competent, like-minded type not easily bought by special interests, then it makes perfect sense to be highly concerned with principledness and consistency*”. Conversely, voters tend to dislike politicians who change their mind on a policy issue, a practice which is often disparagingly denoted as *flip-flopping*.<sup>1</sup> The allegation of being a flip-flopper is, as a matter of fact, one of the most frequently used attacks in electoral races. In the recent past, two famous cases of American presidential candidates that have considerably suffered from being viewed as flip-floppers are John Kerry and Mitt Romney. The fact that voters tend to punish politicians who flip-flop is also to a large extent confirmed by empirical literature in political science.<sup>2</sup>

In a changing world in which politicians are constantly exposed to new information, however, changing one's mind on an issue is natural and might be the optimal thing to do in many situations: as Keynes put it, “*When the facts change, I change my mind*”. In this respect, the reputational stigma associated with flip-flopping can seem puzzling. In this paper I show how flip-flopping can be detrimental to a politician's reputation even in the absence of concerns about ideology or congruence, i.e. in situations in which changing one's mind simply reflects a change in the information available to the politician. In particular, flip-flopping is rationally punished by voters if i) the optimal policy choice is persistent ii) politicians have private information which cannot be credibly revealed to the public and iii) voters are not (fully) capable of judging the validity of a policy choice. The reason for this penalization is that policy shifts are more likely to be performed by

---

<sup>1</sup>The use of this expression dates back to at least 1890 according to the archives of The New York Times. Other terms used to shed negative light on the change of course of a politician are *u-turn* and *backflip*.

<sup>2</sup>See Tomz and Houweling (2012) and Doherty, Dowling, and Miller (2015): I will elaborate more about these and other studies in the related literature section.

poorly informed politicians; therefore, voters trying to select well-informed (competent) politicians assign a better reputation to politicians who do not flip-flop<sup>3</sup>.

If voters perceive policy shifts as a sign of incompetence, the behaviour of strategic politicians will be affected: in particular, office-motivated politicians will have the incentive to distort their behaviour in order to avoid flip-flops. My model therefore describes a form of electoral pandering which is endogenously induced by the previous action of the politician and which results in an excess of conservatism, as an effort to posture consistency.

A recent historical example in which such a logic seems to have played an important role is the decision to start the Iraq war: George W. Bush and Tony Blair believed Saddam Hussein had (or was in the process of developing) weapons of mass destruction, and that waging war on Iraq was necessary to stop him. As a result, the crisis escalated and the countries prepared for war. When, however, evidence pointing in the opposite direction was revealed (including for example confidential intelligence reports<sup>4</sup>), the two leaders should have probably revised the plan to invade Iraq. This flip-flop was unfortunately never performed. As The Guardian wrote concerning the Chilcot report: *“That was the point at which the UK government could and should have said the US must count the UK out. Blair should have admitted that this was a line in the sand. But he didn’t call a halt”*.

After establishing that reputational concerns lead to an insufficient amount of flip-flopping, the second part of the paper discusses some institutional design approaches to tackle the problem: the first is a single term limit policy. By prohibiting the re-election of incumbents, such a policy eliminates all policy distortions, since it shuts down the only source of misalignment between politicians and voters featured in the model. At the same time, however, the single term limit forces voters to forgo all the potential gains of learning about a politician’s type from his track-record. In addition, such policies require commitment, for example through constitutions.

The other institutional feature I discuss is the media environment and the way it relates to the issue of distorted consistency highlighted in my model. I focus on the

---

<sup>3</sup>There is indeed evidence of such a rationale in political commentary. Jack Shafer of the media outlet Politico, for example, wrote the following referring to Hillary Clinton: *“So if new or better information has been the impetus for her policy shifts, she must concede that she has a fat history of taking the wrong position in the early going and then requiring a re-do”*. The full article can be found at this link: <http://www.politico.com/magazine/story/2015/10/democratic-debate-hillary-clinton-flip-flop-213247>

<sup>4</sup>For example, the publication of the Chilcot report in the United Kingdom provided evidence that Tony Blair’s decision to maintain his support of the invasion of Iraq was not in accordance with the information he had received from intelligence sources.

two main roles played by the media, i.e. reporting policy choices (reporting media) and evaluating them (commentator media). One of the results standing out from the analysis of the media model is that fully accurate reporting of policy choices is never optimal. A noisy media can insulate politicians from the reputational stigma of flip-flopping: lies are crowded out by noise, with a positive net effect also in terms of selection of competent politicians through elections. If the noise is too large, however, selection eventually worsens (no learning is possible when the media is fully noisy). This result suggests that the faster and broader access to politicians' track records made possible by improvements in technology and the rise of social media might have led to an increase in policy distortions. This result is related to the idea that transparency can be damaging for political accountability: adding a noisy reporting media to the model is in fact similar to relaxing the assumption of full action transparency.

The other type of media that I analyze is what I define commentator media: instead of reporting on the actions of the politician, which are now assumed to be commonly known like in the benchmark model, the commentator media sends a signal on the state of the world. Whereas this usually acts towards disciplining politicians to follow their signal, there are situations in which increasing the informativeness of the commentator media results in politicians distorting their actions to a larger extent.<sup>5</sup> The reason is that in addition to increasing the chances of being contradicted by the media following postured consistency, increased informativeness makes media endorsements more valuable. This force can lead politicians to gamble on receiving a higher-powered endorsement (as consistent policy-makers) less frequently rather than a less-powered endorsement more frequently, thus increasing distortions. The main takeaway from this result is that more accurate commentator media might be bad for political accountability. This could have relevant implications, for example with respect to the public subsidization of media outlets.

Relatedly, I study a variation of the model in which a public signal about the state of the world is observed *before* the politician chooses the second action. In this setup I show that, unlike in the benchmark model, flip-flopping (and not the avoidance of it) can be motivated by electoral opportunism and that the equilibrium level of flip-flopping can be larger than the socially optimal one. The reason is that, when a public signal is available, there are conditions under which flip-flopping only occurs to match the public signal/opinion poll: therefore, politicians have an incentive to flip-flop and match the poll even if their information does not support such a move.

---

<sup>5</sup>This can happen if i) persistence is high ii) incompetent types have a sufficiently informative signal and iii) most politicians are competent.

Finally, in another extension I show that delegating the first action to an independent agent (for example a bureaucrat or a committee) can be beneficial even if this agent is incompetent: by increases the amount of flip-flops required from both competent and incompetent politicians, the signal of incompetence associated with flip-flopping is mitigated: this improves both accountability and political selection. On the other hand, delegating the action to a competent agent makes increases the consistency of incompetent politicians, but also makes it easier for incompetent incumbents to ignore their signals. The net effect on accountability is ambiguous, but the effect on political selection is always negative.

## 2 Related Literature

This paper is related to several streams of literature. The electoral concerns model that I consider builds on Canes-Wrone, Herron, and Shotts (2001), Prat (2005) and Ashworth and Shotts (2010). Besides some differences in how the model is built, of secondary importance for the results, the main novelty of my model is to introduce an additional period in which the incumbent takes an action. This feature is central for the contribution of the paper, since it allows me to delve into the intrinsically dynamic nature of flip-flopping. Moreover, by adding a stage to a model of electoral concerns, my analysis also allows for a simple endogenous interpretation of the concept of pandering towards a popular action described in electoral concerns models. In my model, as a matter of fact, the popular action is simply the previous policy choice. In this sense, my work contributes to the political economy literature on conformity and pandering: along with the aforementioned Canes-Wrone, Herron, and Shotts (2001) and Prat (2005), another seminal contributions is Maskin and Tirole (2004) - who compare the welfare properties of representative democracy, direct democracy and judicial power; other contributions include Levy (2007), who considers a committee of career concerned decision makers and Frisell (2009), who shows that voters' beliefs about the politician can have self-fulfilling consequences. Levy (2004), on the other hand, shows that in a similar setting also anti-herding (anti-conformism) can take place.

Connected to the idea of pandering towards a popular action is that of status-quo bias and propensity to reform: examples of models with endogenously status-quo biased politicians include Fu and Li (2014), where career-concerned policy makers undertake a reform with lower than optimal probability, and Dewan and Hortala-Vallve (2014), in whose model voters learn through either the success of a reform or the information

provided by a rival candidate. My work, therefore, links the concepts of pandering and status quo-bias: the past actions of the politician influence voters' beliefs to which the politician has an incentive to conform.

As far as the discussion on media is concerned, Prat (2005) and Ashworth and Shotts (2010) are the closest references. The introduction of a commentator media sector in the model draws mostly on Ashworth and Shotts (2010): my contribution relative to their work lies in the characterization of partially truthful equilibria and in the comparative statics analysis which shows how increasing the accuracy of the commentator media need not decrease the distortion of politicians' behaviour. These comparative statics also relate to one of the results of Gentzkow and Shapiro (2006): in their paper, the decision maker is a media outlet, which panders towards the prior belief of citizens in order to form a reputation of accuracy. Whereas Gentzkow and Shapiro (2006) show that the more likely it is for voters to learn the true state of the world, the lower the pandering by the media, I show that introducing an informative media might actually lead politicians to act in a more distorted way.

The analysis of the reporting media, on the other hand, is related to the discussion of action versus consequence transparency in Prat (2005). Since for many policy choices the assumption of action secrecy is not a realistic alternative, I show that the accuracy of the reporting media plays a similar role, and I demonstrate that the optimal arrangement is to have some but not perfect information on the action taken by the politician.

The single term limit rule I mention in the institutional design section, on the other hand, is related to the comparison between representative democracy and judicial power carried out by Maskin and Tirole (2004). A single term limit rule has similar effects to those of delegating decision making to a judicial power not subject to elections.

The consequences of reputation concerns on expert behaviour have also been studied outside the electoral environment. In particular, repeated action by experts has been analyzed by Prendergast and Stole (1996), Li (2007) and in Aghion and Jackson (2016): in addition to many differences in the modelling strategy, the main conceptual difference between my analysis and that of Prendergast and Stole (1996) has to do with the fact agents are forward-looking and the fact that I consider a changing environment rather than a fixed state. Aghion and Jackson (2016), on the other hand, consider repeated action over a sequence of independent and identically distributed states (whereas I allow for correlation) and with action consequences being observable to the principal (whereas they are unobservable in my baseline model). In Li (2007), an agent has to deliver two reports before being evaluated and paid a wage: the state of the world is fixed, but

decision makers' information becomes more accurate in the second period; similarly to what happens in my commentator media model, in her model agents sometimes gamble on being proved right. The idea of gambling on a policy likely to be proved wrong is also at the heart of the paper by Majumdar and Mukand (2004): theirs is a model of experimentation, in which an incumbent can choose to implement a risky policy and has to decide whether to continue the project after a potentially unsuccessful trial. In their model, low type governments are inefficiently reluctant to abandon bad projects, gambling on the small probability of success that would boost their reputation.

Whereas my paper is the first, to the best of my knowledge, to provide a theory of flip-flopping by a politician taking repeated decisions over an issue, there exists some theoretical work by Agranov (2016) and Hummel (2010) dealing with flip-flops between primary and general elections, hence with a completely different objective than that of my analysis.

On the empirical side, finally, there are several papers assessing how voters react to politicians flip-flopping. Doherty, Dowling, and Miller (2015), for example, find that flip-flopping affects the perception voters have of a politician. They show that voters are more forgiving of flip-flops on complex issues or issues which are far away in time. These predictions are in line with my model, in which the more persistent the state of the world, the worse is the reputation from flip-flopping (we can think that persistence is lower the longer the period of time one looks at and the more complex and multi-dimensional the subject is). In the same paper, the authors also show that the reputational cost of a flip-flop is compensated by the fact that the new position taken by the politician will be positively received by some voters, creating a trade-off for the politician. Similarly, Tomz and Van Houweling (2012) conduct survey experiments showing that candidates repositioning affects their support not only in terms of commitment to an ideological issue, but also in terms of perceived valence. The valence aspect is important in particular for issues that are not too salient for voters.

Another issue that has been studied is whether the effects of a flip-flop differ between issues of principle (as abortion or gay marriage) versus pragmatic issues related to a specific policy. Tavits (2007) shows that flip-flopping on pragmatic issues is seen less badly than flip-flopping on ideological issues. On the other hand, Tomz and Houweling (2012) do not find any difference across issues: independently of the issue, candidates who reposition themselves perform worse. My model focuses on issues where competence is key and on which all citizens would agree if perfectly informed. Tomz and Houweling (2012) also argue, without building a model, that the bad perception of flip-flopping

deters candidates away from it: this is exactly what happens in the formal model I present, where flip-flopping is not bad per se, but politicians avoid it because it carries a bad signal for their reputations.

Finally, Levendusky and Horowitz (2012) show experimentally that in the context of international relations, leaders who make threats and subsequently back down pay a cost in terms of electoral support and reputation (which in the international relations literature is called *audience cost*): one of the main reasons for this effect is that a leader changing his mind is seen as less competent than one who stays coherent. What is more, they show that partisanship does not play a significant role in the determination of audience costs. This evidence seems to capture a mechanism close to the one I present in my model. As a matter of fact, given the information asymmetry between politicians and voters and the fact that politicians are often evaluated for their foreign policy conduct before the consequences of their actions are fully known, foreign policy issues are among those where the theory I develop has more bite.

### 3 The Model

There are two periods  $t \in \{1, 2\}$ ; at the beginning of the first period, an incumbent politician is randomly drawn to enter the game and take an action  $a_t \in \{0, 1\}$  in each period on behalf of a representative voter (or a set of identical voters). The action's payoff depends on the underlying state of the world, which can take two values  $\omega_t \in \{0, 1\}$ . The commonly known prior probability that  $\omega_1$  takes value 0 or 1 is equal to  $\frac{1}{2}$ ; moreover, the state is persistent so that for  $j \in \{0, 1\}$ ,  $Pr(\omega_2 = j | \omega_1 = j) = \gamma > \frac{1}{2}$ . Incumbents can be of two types  $\theta \in \{H, L\}$ , i.e. competent and incompetent. In each period, both types receive an informative signal  $s_t \in \{0, 1\}$  on the state of the world, but the accuracy of the signal,  $Pr(s_t = j | \omega_t = j) = q_\theta$  depends on the politician's type:  $\frac{1}{2} < q_L < q_H \leq 1$ . To simplify exposition I fix  $q_H = 1$  (competent politicians perfectly observe the state) and therefore I drop the subscript from  $q_L$ , which will be simply denoted by  $q$ . Notice that since the prior is  $\frac{1}{2}$ , for all  $\gamma < 1$  the signal received by any politician is always decision-relevant, meaning that the probability of matching the action to the state is maximized by  $a_t = s_t$ . I denote by  $\rho_t((s_1, \dots, s_t), \theta) = Pr(\omega_t = s_t | (s_1, \dots, s_t), \theta)$  the posterior probability that the state is equal to the signal after observing realization  $s_t$ . Since the posterior of the perfectly informed competent politician is always equal to 1,  $\rho_t$  will denote, when not further specified, the posterior of the incompetent politician. Moreover, since most of the analysis revolves around that, when not further specified  $\rho_2$  will denote

$\Pr(\omega_2 = s_2 | s_1, s_2 \neq s_1, \theta = L)$ , whereas  $\bar{p}_2$  will indicate  $\Pr(\omega_2 = s_2 | s_1, s_2 = s_1, \theta = L)$ .

Politicians privately observe their competence, and the signals they receive are also private information: competent politicians represent a fraction  $\lambda$  of incumbents.<sup>6</sup> At the end of period  $t = 2$  there is an election, in which the representative voter takes action  $e \in \{r, f\}$ , i.e. decides whether to retain ( $r$ ) or fire ( $f$ ) the incumbent politician. Voters know the statistical process governing the economy and they observe the track-record of the incumbent politician, i.e. the actions taken over the two periods, denoted by  $\tau = (a_1, a_2)$ . Track-records are used to form beliefs  $\mu(a_1, a_2) = \Pr(\theta = H | a_1, a_2)$  over an incumbent's competence, which I call reputation. Before the election takes place, a challenger appears. The challenger is competent with probability  $\lambda_O$ . The representative voter's utility depends on whether the incumbent's actions match the state of the world and on the type of politician winning the election. Selecting a competent politician (denote by  $\theta_e$  the type of the election winner) gives the voter a benefit of  $b$ , whereas  $v$ , which is drawn from a uniform distribution in  $[-b, b]$  after the incumbent has taken the second action, represents the relative valence of the challenger versus the incumbent. Hence, the representative voter's utility, denoted by  $U_c$ , reads:

$$U_c = \sum_{t=1}^2 \mathbb{1}_{a_t=\omega_t} + \mathbb{1}_{e=f}v + \mathbb{1}_{\theta_e=H}b.$$

Politicians derive utility both from taking the right action while in office and winning the election, with  $2\phi$  denoting the additional utility received if re-elected. Formally:

$$U_p = \sum_{t=1}^2 \mathbb{1}_{a_t=\omega_t} + \mathbb{1}_{e=r}2\phi.$$

A politician's information set, or history, denoted by  $h_t$ , includes all actions up to  $t-1$  and signals up to  $t$ . The strategy of the incumbent is a mapping  $\Psi$  from any history  $h_t$  to a probability of playing each action  $a_t \in \{0, 1\}$ . In particular, it is useful to express the incumbent's strategy as  $\sigma_t(h_t, \theta)$ , where  $\sigma(h_t, \theta) = \Pr(a_t = s_t | (s_1, \dots, s_t), (a_1, \dots, a_t), \theta)$  is the probability of choosing a policy  $a_t$  in accordance with the realization of the signal  $s_t$  at time  $t$ . Since that will be crucial to the equilibrium characterization, when not further specified  $\sigma$  will denote the probability of choosing an action matching the signal at time  $t = 2$  given  $s_2 \neq s_1$ .

The voter's strategy is a mapping between the track record  $\tau$  and the valence realization  $v$  and a probability distribution over decisions  $e \in \{r, f\}$ . Voters choose each

---

<sup>6</sup>The fact that politicians know their competence is not necessary for the core results of the paper.

politician with probability  $1/2$  when indifferent.<sup>7</sup>

The equilibrium concept I use is Perfect Bayesian Equilibrium (PBE). In order to rule out pooling equilibria, that is equilibria where at each time  $t$ , both types play the same action with probability 1, I use a trembling-hand perfection refinement.<sup>8</sup> Finally, for the uniqueness result I focus on symmetric equilibria, that is equilibria that are robust to relabelings of states of the world, signals and actions that keep the information structure identical (e.g. switching the name of state/action/signal 1 to 0 and viceversa in a symmetric information environment with symmetric priors and symmetric signals). In other words, in situations that are informationally the same, players play in the same way. As I will show in the appendix, there is a unique non-pooling and symmetric PBE, which coincides with the most informative PBE.<sup>9</sup>

## 4 Results

I start the analysis from the election in which the voter chooses between the incumbent and the challenger. Just before the election, the valence draw  $v$  is realized. As a result, voting for the challenger gives the voter a utility of  $v + \lambda_O b$ . Voting for the incumbent gives instead the voter a utility of  $\mu(a_1, a_2)b$ . Therefore, the incumbent is re-elected if  $v \leq (\mu(a_1, a_2) - \lambda_O)b$  and

$$Pr(v \leq (\mu(a_1, a_2) - \lambda_O)b) = \frac{(\mu(a_1, a_2) - \lambda_O)b - (-b)}{2b} = \frac{1}{2} + \frac{\mu(a_1, a_2) - \lambda_O}{2} \equiv r(\mu(a_1, a_2)).$$

As we can see, the re-election probability is linearly increasing in reputation, from which it follows that at  $t = 2$ , the difference in the probability of winning the election conditional on taking action 0 or 1 is  $\frac{\mu(a_1, 0) - \mu(a_1, 1)}{2}$ .

Having described the election stage, consider the decisions of the incumbent in periods  $t = 1$  and  $t = 2$ . Signals are always decision-relevant, so maximizing the probability to match the action to the state requires that politicians follow their signals. If that happens in equilibrium, the resulting equilibrium is said to be truthful.

**Definition 1.** *An equilibrium is truthful if and only if  $\sigma(s_t) = 1$  for each  $s_t$  at any  $t$  and*

---

<sup>7</sup>Notice, however, that this is a zero probability event, since challengers are distributed according to an atomless distribution.

<sup>8</sup>For additional details on how the trembling-hand refinement eliminates these equilibria, I refer to the Appendix.

<sup>9</sup>I have not been able to construct examples of non-symmetric equilibria, so I believe uniqueness does not rely on the symmetry requirement, but proving uniqueness relaxing the assumption of symmetry is substantially more complicated and not informative for the sake of the results presented in the paper.

for each type  $\theta$ .

Denote the reputation of different track records under truthful play by  $\mu^T(a_1, a_2)$ . There are four possible track records:  $\tau \in \{(0, 0), (0, 1), (1, 0), (1, 1)\}$ . Given the symmetric initial prior, the probability of obtaining each of the two consistent and flip-flopping signal sequences is the same: therefore, as the next claim will prove,  $\mu^T(0, 0) = \mu^T(1, 1)$  and  $\mu^T(0, 1) = \mu^T(1, 0)$ . In other words, all that matters is whether or not the politician changed his mind. I will call a track-record in which the politician did not change his mind *consistent* and one in which he did change his mind *flip-flopping*;  $\tau = C$  and  $\tau = F$  indicate consistent and flip-flopping track-records respectively. Similarly,  $\mu_C$  ( $\mu_C^T$  for truthful play) and  $\mu_F$  ( $\mu_F^T$  for truthful play) denote the reputation from consistent and flip-flopping track-records.

**Lemma 1.** *Under truthful play, the reputation of a consistent track-record is strictly larger than that of a flip-flopping track-record:  $\mu_C^T > \mu_F^T$ .*

*Proof.* Take any flip-flopping track record. Under truthful play,  $a_t = s_t$ , so that  $Pr(\tau|\theta) = Pr(s_2, s_1|\theta)$ . Denote by  $A(q)$  the following:

$$A(q) = (1 - \gamma)(q^2 + (1 - q)^2) + \gamma 2q(1 - q).$$

It can be easily verified that  $\frac{A(q)}{2}$  represents the probability that a politician with a signal of accuracy  $q$  receives a flip-flopping sequence of signals (the fact that the same probability holds for both sequences  $(0, 1)$  and  $(1, 0)$  follows from the prior being equal to  $\frac{1}{2}$ ). The expression for  $A(q)$  can be rearranged to  $A(q) = (1 - \gamma) + (2\gamma - 1)2q(1 - q)$  and since  $\gamma > \frac{1}{2}$ ,  $2\gamma - 1 > 0$  and  $A(q)$  is decreasing in  $q$ . Moreover,  $A(1) = 1 - \gamma$ . Finally, I construct the reputations from the 4 possible track-records. Since  $A(q)$  only depends on whether the politician received a consistent or flip-flopping sequence of signals, there are only two possible levels of reputation: one from flip-flopping and one from being consistent. To see that the reputation from flip-flopping is lower than that from being consistent, notice that  $A(q) > A(1) \Leftrightarrow \frac{A(q)}{A(1)} > \frac{1 - A(q)}{1 - A(1)}$  and therefore:

$$\mu_F^T = \frac{\lambda A(1)}{\lambda A(1) + (1 - \lambda)A(q)} < \frac{\lambda(1 - A(1))}{\lambda(1 - A(1)) + (1 - \lambda)(1 - A(q))} = \mu_C^T$$

From now on, I will simply write  $1 - \gamma$  for  $A(1)$  and write  $A$  for  $A(q)$ . □

The result I just stated is very important for the development of the whole paper since it highlights the reason that leads incumbents to distort their actions: the bad reputation

associated with flip-flopping. The stigma of flip-flopping, as a matter of fact, puts an office motivated incumbent in front of a trade off. Whenever receiving a second signal that contradicts the first one, i.e.  $s_2 \neq s_1$ , a politician knows that doing the optimal thing for society will result in a worse reputation. If the flip-flopping stigma is large enough compared to the benefit from following his signal, the politician will therefore choose to act against his private information.

In light of this, it can be shown that a necessary and sufficient condition for the existence of a truthful equilibrium is that incumbents have an incentive to follow their signal after  $s_2 \neq s_1$  (we will see that politicians always follow their signal at  $t = 1$ ). Moreover, a single-crossing property makes it sufficient to simply look at the incentives of incompetent incumbents. The reason is that whereas both competent and incompetent politicians enjoy the same benefit from avoiding a flip-flop, they do not sustain the same costs: having a less accurate signal, as a matter of fact, means having a larger probability of matching the state of the world when playing  $a_t \neq s_t$ . Lying is therefore cheaper for incompetent politicians. This property is very useful for the characterization of the equilibrium.

**Lemma 2 (Single Crossing Property).** *The cost of acting against one's signal is  $\rho_t(\theta) - (1 - \rho_t(\theta)) = 2\rho_t(\theta) - 1$ , and since  $\rho_t(\theta = L) < \rho_t(\theta = H) = 1$ , contradicting the signal is more costly for the competent politician.*

It follows that a truthful equilibrium is only sustainable as long as incompetent politicians are willing to follow their signal at  $t = 2$  after receiving two conflicting signals. This in turn requires the office motivation parameter to be low enough, because as  $\phi$  grows larger, the benefits from holding office progressively dwarf the utility from matching the action to the state of the world.

**Proposition 1.** *A truthful equilibrium is sustainable as long as  $\phi \leq \bar{\phi}$ , where*

$$\bar{\phi} = \frac{2\rho_2 - 1}{\mu_C^T - \mu_F^T}$$

*Proof.* See Appendix. □

The parameter  $\bar{\phi}$  is therefore the upper-bound on electoral rents under which a truthful equilibrium is sustainable.<sup>10</sup> The following theorem extends the characterization to the case of electoral rents exceeding the upper bound  $\bar{\phi}$ .

---

<sup>10</sup>The actual rents for the politician are  $2\phi$ , but that only serves the purpose of simplifying calculations and has no other consequence on the model.

**Theorem 1.** *The game has a unique non-pooling symmetric Perfect Bayesian Equilibrium. For  $\phi \leq \bar{\phi}$ , the unique equilibrium is the truthful equilibrium. For  $\phi > \bar{\phi}$ , the unique equilibrium is partially truthful, meaning that:*

$$\begin{aligned}\sigma(s_1) &= \sigma(s_2|\theta = H) = \sigma(s_2 = s_1 = a_1|\theta = L) = 1 \\ \sigma(s_2 \neq a_1|\theta = L) &= \sigma^* < 1.\end{aligned}$$

*Proof.* See Appendix. □

This theorem proves that the game always has a unique symmetric non-pooling equilibrium. It must therefore be the case that when  $\phi \leq \bar{\phi}$ , the unique such equilibrium is truthful. When  $\phi > \bar{\phi}$ , on the other hand, the unique such equilibrium is partially truthful, since whenever the signal received by the incumbent in the second period prescribes to flip-flopping, incompetent politicians mix between following their signal and repeating their previous action.

Theorem 1 has two interesting implications: the first is that for any level of  $\phi$ , flip-flopping always decreases the incumbent's reputation (since the two reputations average at  $\lambda$ , a bad reputation is always below  $\lambda$ ). However, the larger the flip-flopping avoidance distortion caused by incompetent politicians not following their signal, the smaller is the reputation gap between consistent play and flip-flopping, i.e. the less bad is the reputation from flip-flopping. As  $\phi$  tends to infinity and politicians only care about re-election, the reputation gap between consistency and flip-flopping approaches zero.

The second (and related) implication is that when the equilibrium is partially truthful, there is an insufficient amount of flip-flopping compared to the truthful equilibrium. In other words, voters stigmatize flip-flopping but at the same time they would be better off if more flip-flopping took place. As a matter of fact, it could even happen that voters are more confident about the policy being correct after seeing a flip-flop rather than a consistent policy: a flip-flop sends a bad signal on the incumbent's type but is always earnest, whereas a consistent policy is a good sign on the politician's type but not necessarily earnest. I summarize these insights in the following corollary:

**Corollary 1.** *In equilibrium, flip-flopping gives a bad reputation compared to truthful play:  $\mu_C^* - \mu_F^* > 0$  for any  $\phi$ . Equivalently,  $\mu_C^* > \lambda > \mu_F^*$ . In a partially truthful equilibrium, there is insufficient flip-flopping compared to the truthful equilibrium.*

Notice that another interesting implication of the model is that change hurts incumbents: when the state of the world changes between period 1 and 2, the reputation of incumbents is likely to fall (because of flip-flopping) and this means the incumbent is

more likely to be replaced. Notice that conditional on the state changing, good leaders are more likely to be replaced than bad leaders. However, conditional on having a bad leader in office, a change in the state increases the chances of having a better leader in the following period. In other words, improvements in leadership are more likely after a change in the state of the world.

**Corollary 2.** *The following properties hold:*

- i A leadership change is more likely after a change in the state of the world.*
- ii If  $\lambda_O$  is high enough, the probability of having a competent leader in office after elections is higher when the state changes than when the state does not change.*
- iii An incompetent incumbent is more likely to remain in office after a change in the state of the world compared to a competent incumbent; a competent incumbent is more likely to remain in office when the state does not change with respect to an incompetent incumbent.*

*Proof.* See Appendix. □

## 4.1 Comparative Statics and Welfare

I begin with the characterization of how the threshold  $\bar{\phi}$  moves with the parameters of the model. This characterization tells us what level of electoral incentives is sustainable without distorting the behaviour of politicians.

**Comparative Statics 1.**  *$\bar{\phi}$ , that is the maximum level of  $\phi$  such that a truthful equilibrium is sustainable, is increasing in  $q$ , decreasing in  $\gamma$  and can be both increasing or decreasing in  $\lambda$ .*

*Proof.* See Appendix. □

I next analyze the welfare implications of parameter changes. The definition of welfare in this model is based on the expected utility of voters, calculated as of time  $t = 0$  (i.e. before randomly picking the incumbent); as such, welfare does not account for the utility of politicians.<sup>11</sup>

---

<sup>11</sup>One can think that politicians are too small a fraction of the population for their welfare to matter in the aggregate; moreover, whereas the policy component of politicians' utility is the same as voters' (and therefore it is captured in welfare), the electoral rent part is constant in aggregate and simply transferred between politicians, hence neutral with respect to total welfare.

**Definition 2.** *Social welfare  $W$  is defined as:*

$$W = \underbrace{\mathbb{E}_0 \left[ \sum_{t=1}^2 \mathbb{1}_{a_t=\omega_t} \right]}_{\text{Accountability}} + \underbrace{\mathbb{E}_0 [\mathbb{1}_{e=f}v + \mathbb{1}_{\theta_e=H}b]}_{\text{Selection}}$$

Welfare can be decomposed in two parts, accountability and selection. Accountability indicates whether the incumbent acts in the best interest of society, which in this case means using information efficiently by following the signal: accountability welfare is therefore the probability that the incumbent chooses the optimal policy, formally  $\mathbb{E}_0 \sum_{t=1}^2 \mathbb{1}_{a_t=\omega_t}$ . Selection welfare, on the other hand, indicates the utility the voter derives from the election winner, also accounting for the valence shock: formally  $\mathbb{E}_0 [\mathbb{1}_{e=f}v + \mathbb{1}_{\theta_e=H}b]$ .

Let's denote by  $\tilde{q} \equiv q - A(1 - \sigma^*)(2\rho_2 - 1)$  the accuracy of the incompetent politician's policy choice, taking the flip-flopping avoidance distortion into account.

**Lemma 3.** *The expression for welfare can be rewritten in the following way:*

$$W = \underbrace{[\lambda + (1 - \lambda)q] + [\lambda + (1 - \lambda)\tilde{q}]}_{\text{Accountability}} + b \underbrace{\left[ \frac{1}{4} + \frac{\lambda_O^2}{4} - \frac{\lambda\lambda_O}{2} + \frac{\lambda + \lambda_O}{2} + \frac{\mathbb{E}\mu^2}{4} \right]}_{\text{Selection}}$$

*Proof.* See Appendix. □

From Lemma 3 we can see that the strategic behaviour of politicians affects accountability welfare through  $\tilde{q}$  and it affects selection welfare through  $\mathbb{E}\mu^2$ . Concerning the effect on  $\tilde{q}$ , it is straightforward to see that the higher the probability of not following the signal, the lower accountability welfare. In terms of the effect of  $\sigma^*$  on selection welfare, what matters is the variance of the reputation generated by the equilibrium behaviour. The lower the  $\sigma^*$ , the lower  $\mu_C$  and the higher  $\mu_F$ , which tend to converge as the probability of following the signal gets lower. However, as  $\sigma^*$  decreases, at the same time the probability of the realization  $\mu_C$  increases and the probability of  $\mu_F$  decreases, since politicians increasingly play consistent track records. This makes the result on the change in  $\mathbb{E}\mu^2$  potentially ambiguous. The following lemma tackles this issue by providing a useful connection between selection welfare and the equilibrium condition, which is written in terms of  $\mu_C - \mu_F$ .

**Lemma 4.** *A sufficient condition for selection welfare to increase is that both  $\mu_C$  and  $\mu_C - \mu_F$  increase. A necessary condition selection welfare to increase if  $\mu_C - \mu_F$  decreases is that both  $\mu_C$  and  $\mu_F$  increase.*

*Proof.* See Appendix. □

A straightforward consequence of Lemma 4 is that decreasing  $\sigma^*$  decreases selection welfare, everything else equal. Clearly,  $\sigma^*$  is an endogenous object, and therefore what matters is the comparative statics results of changes in the exogenous parameters of the model. The first result concerns changes in  $\phi$ : since it enters welfare only through  $\sigma^*$ , a direct consequence of Lemma 4 is that increasing  $\phi$  weakly decreases welfare.

**Welfare Result 1.** *Increasing  $\phi$  weakly decreases welfare.*

*Proof.* See Appendix. □

It has to be kept in mind, of course, that this model abstracts from all those reasons why it might be a good idea to offer electoral incentives to politicians (for example to improve the share of competent politicians in the pool); since in this model politicians' interests are aligned with those of citizens except for electoral incentives, it is not surprising that providing electoral incentives can only make things worse.

Another conclusion from the analysis is that having a wide competence gap between the two politicians' types is bad for accountability: as a result, increasing  $q$  always increases welfare, both directly because incompetent politicians have better information when they choose a policy and indirectly because politicians become more truthful (i.e.  $\sigma^*$  increases). On the other hand, if there were no distortions, i.e.  $\sigma^* = 1$ , then the direct effect is coupled with a negative indirect effect, since it becomes more difficult to tell apart good politicians from bad ones. This is an example in which understanding the strategic behaviour of politicians can lead to opposite policy implications compared to a model ignoring them.

**Welfare Result 2.** *Within a partially truthful equilibrium, increasing  $q$  strictly increases welfare. Within a truthful equilibrium, on the other hand, increasing  $q$  strictly increases accountability welfare but has ambiguous selection welfare effects.*

*Proof.* See Appendix. □

Notice that when  $q$  is high, i.e. politicians are in general competent, then the reputation  $\mu_F$  is lower, i.e. flip-flopping hurts more. This might be one of the reasons why flip-flopping can hurt candidates in a race such as the US presidential election, despite the fact that electoral incentives are high. In other words, flip-flopping is worse for a candidate's reputation when there is more homogeneity between competent and incompetent candidates.<sup>12</sup>

---

<sup>12</sup>Notice that a bad reputation in this setup is just the probability of being of the incompetent type, no matter how bad the incompetence is.

The next comparative statics concern the fraction of competent politicians  $\lambda$ . As this fraction increases, both  $\mu_C$  and  $\mu_F$  increase. However,  $\mu_C - \mu_F$  can either increase or decrease. As a result,  $\sigma^*$  can move in either direction. This means that there exist cases in which having a better pool of politicians increases the distortion generated by incompetent politicians avoiding flip-flops. In terms of welfare from the selection of politicians, however, an increase in  $\lambda$  is always beneficial, at least whenever the starting point is a partially truthful equilibrium. The reason is that the equilibrium level of  $\mu_C - \mu_F$  remains constant, and therefore both  $\mu_C$  and  $\mu_F$  increase. From Lemma 1, this means that the selection of politicians improves.

**Welfare Result 3.** *Increasing  $\lambda$  in a partially truthful equilibrium improves selection welfare, but it can either increase or decrease accountability welfare. In a truthful equilibrium, on the other hand, accountability welfare increases but selection welfare might increase or decrease.*

*Proof.* See Appendix. □

Things get more complicated when evaluating a change in the persistence parameter  $\gamma$ . With more persistence there are several possible cases: first of all,  $\sigma^*$  can be either increasing or decreasing in  $\gamma$ . In the former case, accountability improves with higher persistence, whereas in the latter it could go either way. In terms of selection, however, it can be shown that a more persistent state of the world decreases the effectiveness of elections in selecting competent politicians:

**Welfare Result 4.** *In a partially truthful equilibrium, selection welfare decreases as  $\gamma$  increases. In a truthful equilibrium, on the other hand, selection welfare increases.*

*Proof.* See Appendix. □

This property is interesting since it tells us that when the game is in a partially truthful equilibrium, elections perform worse in a less variable world. This can seem surprising, especially given Corollary 2: however, notice that these two results are different, since Corollary 2 considers what happens within an equilibrium, whereas this comparative statics result compares different equilibria, with different strategies and reputations.

## 5 Institutional Design

So far I have shown how the bad reputation from flip-flopping can give (incompetent) politicians an incentive to distort their actions. In the baseline model I have presented,

there are three fundamental ingredients leading to the result: the first is the fact that politicians face an election at the end of  $t = 2$ , because being re-elected gives them a utility of  $2\phi$ ; the second is the fact that voters have no information on what constitutes the optimal policy in each period and are therefore only able to evaluate incumbents based on their track record in office; the third is that voters are perfectly able to observe the action taken by the incumbent, in such a way that a flip-flop is immediately caught and used to form reputations. My aim in this section is to relax these assumptions by introducing new institutional features to the baseline model.

## 5.1 Single Term Limit

In the baseline model I analyzed above, the incumbent can be re-elected at the end of period  $t = 2$ . As a result, it is interesting to see what would happen under a single term limit rule, in which all politicians can serve only one term in office. The single term limit rule for the President is an institution in several Latin American countries, Armenia, Israel, South Korea, The Philippines, as well as the European Central Bank, among others.<sup>13</sup> While in many instances the purpose of term limits is to prevent the entrenchment of incumbents and the manipulation of the democratic process, this model suggests an additional possible rationale for it, similar to the one proposed by Smart and Sturm (2013). As a matter of fact, in the setup I describe in this model, the single term limit is a blunt yet effective instrument to eliminate all distortions due to politicians' willingness to avoid flip-flopping. At the same time, however, having a single term limit also means forgoing the possibility to condition the reelection decision on the beliefs about the incumbent's type. It follows that banning reelections is only welfare improving if the accountability distortion from the flip-flopping avoidance is large and the upside from retaining incumbents with a good reputation is low.

**Proposition 2.** *A single term limit is welfare improving if and only if:*

$$b \leq 2 \frac{(1 - \lambda)A(1 - \sigma^*)(2\rho_2 - 1)}{\frac{1 + \lambda_O^2 + \mathbb{E}\mu^2}{2} + \lambda - \lambda_O - \lambda\lambda_O}$$

Intuitively, the condition states that the single-term limit rule is beneficial if the benefit from selecting a competent politician is low enough. Moreover, notice that if  $\lambda \rightarrow 1$ , the right-hand side goes to 0, meaning that if the incumbent is competent with

---

<sup>13</sup>In Latin America, countries that adopted a single term limit include Colombia, El Salvador, Guatemala, Honduras, Mexico, Paraguay.

a sufficiently high probability, even a very small benefit  $b$  is sufficient to prefer the re-electability of the incumbent. The same happens if  $q \rightarrow \frac{1}{2}$ , since the gain from better accountability converges to zero as incompetent politicians become less and less informed.

Finally, notice that implementing a single term limit rule requires the ability to commit (for example through a constitution) not to re-elect an incumbent thought to be competent with a high probability.<sup>14</sup>

## 5.2 Transparency of Actions: Reporting Media

Voters usually rely on the media to learn about the policies chosen by politicians. In some circumstances, for example when bills containing multiple prescriptions are voted, it is not so straightforward to understand whether the incumbent politician flip-flopped or played consistently.<sup>15</sup> The same might happen when voting in a committee is secret and knowing the result only enables to make probabilistic statements on whether a member voted in a certain direction.

In this section I therefore relax the assumption of full observability of the incumbent's actions and evaluate its impact on social welfare. The fact that action transparency is not always beneficial is well known in the literature.<sup>16</sup> My contribution is to characterize the level of partial transparency (or media accuracy) that is optimal when the starting equilibrium is partially truthful and to show that it satisfies the intuitive property of being the maximal level of accuracy consistent with politicians using their information efficiently.

Assume that the incumbent's track record is only observable through the report of a media company; the media company's reporting technology, however, is not perfectly accurate: whereas the first action taken by the incumbent is reported accurately, the second one might be misreported.<sup>17</sup> This means that given a true incumbent track record

---

<sup>14</sup>The idea of commitment to a single term limit is developed in a rather hyperbolic form by the Italian writer Italo Calvino (1969) in a short story in which he describes a hypothetical society in which leaders are beheaded at the end of their term in office.

<sup>15</sup>An example of a situation in which it was not simple to label a policy choice as a flip-flop is the vote by Bernie Sanders against the automobile industry bailout in January 2009: Sanders had actually supported the bailout previously and supported it afterwards, but after having voted in favour of it, he voted against the release of that tranche of aid since it also contained financial aid for the banking sector, which Sanders was not in favour of bailing out. In the recent presidential primary election, Hillary Clinton used this alleged flip-flop to attack Sanders.

<sup>16</sup>See Prat (2005) as one of the main examples.

<sup>17</sup>There are several reasons why it makes sense to think of only the second action being reported with noise: on one hand, if the first action were misreported the politician would probably have time to realign the public opinion before taking the second action, whereas that might not be possible after taking the second action, with the urgency of the election. Moreover, as I will discuss later in the paragraph, as long as the misreporting is realized ex-post with respect to the action, it does not matter which action

$\tau$  that displays flip-flopping, with probability  $1 - g$  the media company sends out a report indicating that the politician acted consistently, and the same happens when the true track record is consistent. Therefore, a voter observing a consistent track-record infers that the actual track-record of the incumbent is consistent with probability:

$$p_C \equiv Pr(\tau = C | \tilde{\tau} = C) = \frac{gPr(\tau = C)}{gPr(\tau = C) + (1 - g)Pr(\tau = F)}$$

while an analogous expression, denoted by  $p_F$ , represents the probability that the true track record is flip flopping given an observed flip-flopping record:

$$p_F \equiv Pr(\tau = F | \tilde{\tau} = F) = \frac{gPr(\tau = F)}{gPr(\tau = F) + (1 - g)Pr(\tau = C)}.$$

The voter updates her beliefs and assigns a reputation to each true track-record just like in the baseline game. However, since she does not observe the true track record, the actual reputation that she assigns to the politician is simply the weighted average of the reputations following each track-record, weighted by the probability that the observed track-record is of each type conditional on the observed media report. As a result, if  $g = 1$  there is no noise (full transparency) in the reporting media and we recover the baseline model, whereas if  $g = \frac{1}{2}$  there is no transparency and the reputation associated to each track-record is simply  $\lambda$ .

Since the lower the parameter  $g$ , the less voters can learn about the incumbent independently of his behaviour, it is intuitive that truthful play can be restored for  $g$  lower than threshold  $g^*$ . Moreover,  $g^*$  is always larger than  $1/2$ , since politicians are always truthful when  $g = 1/2$  and by continuity there can always be found a  $g > 1/2$  such that politicians are still truthful. Setting  $g < g^*$  is therefore never optimal, because once truthful play has been restored, which happens at  $g^*$ , further lowering  $g$  only worsens the learning of the voter. However, it can also be shown that increasing  $g$  starting from  $g^*$  is never optimal, making  $g^*$  the optimal level of media accuracy.

**Proposition 3.** *The optimal media reporting accuracy is  $g^* \in (1/2, 1]$ . The value of  $g^*$  is the largest possible such that incumbents play truthfully. Therefore,  $g^* < 1$  whenever a truthful equilibrium is not sustainable in the baseline model.*

*Proof.* See Appendix. □

Finally, notice that as long as reporting on the action takes place after both actions have been taken, it doesn't matter whether misreporting concerns the first, the second 

---

is misreported (it could even be that both are).

or both actions, since, thanks to the symmetry of the game, all that matters is whether a consistent track-record is reported as flip-flopping or viceversa. If instead reporting were noisy on the first action and the politician knew whether the report was correct or wrong before taking the second action, things could potentially differ: noise on the first action increases the cost of not following the signal after a wrong report, which leads to more truthful behaviour, but at the same time it leads to more distorted behaviour after a correct report.

### 5.3 Delegating the First Action

An alternative way of thinking about the transparency of the politician's action is to consider the possibility of delegating the first action to an independent agent. In many contexts in which multiple decisions have to be taken on an issue, as a matter of fact, the politician or top level official does not take care of all the steps in the process, but only acts at some stages. One can for instance think of committees doing preliminary work before a bill is voted, of bureaucrats drafting reforms, of a local government handling an issue before the central government steps in.

In order to analyze this feature, I twist the baseline model in the following way: the first action is taken by an independent and non-strategic agent, whereas the second action is taken by the incumbent as in the baseline model. As a result of delegating the first action, the incumbent observes the action taken by the agent but only receives a private signal before taking the second action. This variation of the model also allows me to contribute to the topic of delegation and political accountability considered, among others, by Fox and Jordan (2011).

The main result of this section is that delegating the first action to an incompetent agent increases both accountability at  $t = 2$  and the selection through elections, at the cost of a worse decision at  $t = 1$ . On the other hand, delegating to a competent agent improves the quality of the first action but worsens selection, whereas it has ambiguous effects on accountability at  $t = 2$ .

Finally, notice that given the timing of the game, it would be optimal to delegate the second action to the agent and let the politician take the first one, since that would give incumbents the incentive to always following the signal in the first period, eliminating the accountability distortion and improving selection. The game with the second action delegated resembles a policy oversight game (see for example Fox and Stephenson (2011)), with the peculiarity that oversight works through the persistence of the state of the

world.<sup>18</sup>

### 5.3.1 Delegating the First Action to an Incompetent Agent

Suppose that the utility citizens receive from the first action matching the state is now  $\alpha > 0$  (which could be larger or smaller than one) instead of 1 as in the benchmark model. Moreover, throughout this discussion I assume that the first decision maker (which I'll also call bureaucrat) is non-strategic and incompetent, i.e. has a signal with accuracy  $q$ . Thanks to this assumption, the information an incompetent incumbent has before taking the second action is exactly the same as in the benchmark model: this also means that the cost of ignoring the signal is the same. The competent type, on the other hand, has less initial information than in the benchmark model, since the first signal is less informative than the one a competent agent would have. As a result, the second period signal of the competent incumbent contradicts the action of the bureaucrat more often than it would contradict his own first action, had the first action not been delegated. Therefore, flip-flopping on the action of the bureaucrat comes at a smaller reputational cost than flip-flopping in the benchmark model. In a partially truthful equilibrium, the reputational wedge between consistency and flip-flopping is the same as in the benchmark model, i.e. it is equal to  $\frac{2\rho_2-1}{\phi}$ . However, achieving that same wedge requires a smaller distortion in the game with delegation. Therefore, delegating leads to a gain in terms of second action welfare, as well as selection, but to a loss in terms of first action welfare. It follows that if  $\alpha$  is small enough, delegating to an incompetent agent is beneficial.

**Proposition 4.** *Suppose that the first decision is delegated to a bureaucrat with signal accuracy  $q$ , i.e. just like the incompetent incumbent. Suppose further that  $\phi$  is larger than the threshold level  $\phi_D$ , so that truthful equilibria are not sustainable independently of whether delegation occurs or not. Then, delegating the first action leads to an increase in accountability welfare at  $t = 2$  as well as in selection welfare. Given that the accuracy of the first action decreases, overall welfare increases if and only if  $\alpha$  is sufficiently small.*

Notice that compared to a single term limit, which provides accountability at the expense of selection, delegation to an incompetent bureaucrat improves selection and (partially) accountability at  $t = 2$  at the expense of a worse decision at  $t = 1$ . Another consequence of this result is that delegation to a less capable agent is more likely to be beneficial when electoral concerns are high rather than when they are low, since with no distortions it is always better not to delegate to a bureaucrat of inferior quality.

---

<sup>18</sup>The statement and proof of this result are very simple and omitted from the current version of the paper for reasons of space.

### 5.3.2 Delegating the First Action to a Competent Agent

A natural question arising from the previous analysis is whether the first action should be delegated to a competent rather than an incompetent agent.<sup>19</sup> As Proposition 4 shows, delegating to an incompetent agent always increases the accountability welfare at  $t = 2$  and selection welfare, increasing total welfare if the weight  $\alpha$  of the first action is low enough. Delegating to a competent agent clearly increases the quality of the first action, but how does it change accountability and selection welfare?

First of all, delegating the first action to a competent agent means increasing the strength of the prior on  $\omega_2$  for the incompetent agent entering the game at time  $t = 2$ . In order to keep the game as comparable as possible with the baseline model, therefore, I assume that  $q > \gamma$ , i.e. even an incompetent agent receiving  $s_2 \neq a_1$  should choose  $a_2 = s_2$  in order to maximize the probability of matching the state of the world  $\omega_2$ .

In terms of flip-flopping, delegating to a competent agent does not change the number of flip-flops of the competent agent, compared to the baseline game, but it decreases the number of flip-flops of the incompetent agent. As a result, under truthful play flip-flopping gives a less bad reputation than in the baseline game. Symmetrically, the good reputation from avoiding flip-flops is smaller than in the baseline game. This result is the mirror image of what happened in the game with delegation to the incompetent agent, and it pushes for less distortions in equilibrium. However, when delegating the action to a competent agent, there is a crucial difference: after receiving  $s_2 \neq a_1$ , the incompetent politician's posterior is closer to  $1/2$  than under delegation to an incompetent agent. This makes lying cheaper and it pushes for more distortions. The presence of these two contrasting forces leads to an ambiguous net effect on accountability welfare at  $t = 2$ .

However, the decrease in the posterior of the incompetent incumbent results in a negative effect on selection welfare. In other words, delegation to a competent bureaucrat hurts the ability of voters to select good politicians through elections. Interestingly, this happens entirely through the effect of the presence of bureaucrats on the reputation of flip-flopping versus consistent policy-making.

Overall, given the increase in the quality of the first action, it is not surprising that delegation to a competent bureaucrat is beneficial as long as  $\alpha$  is sufficiently high.

**Proposition 5.** *Suppose that the first decision is delegated to a bureaucrat with perfect signal accuracy. Suppose further that  $\phi$  is larger than the threshold level  $\max\{\bar{\phi}, \bar{\phi}_d\}$ , so*

---

<sup>19</sup>Clearly, if observably competent agents are available, it is not totally clear why this agent should not be in charge for both actions in the first place. However, one could imagine that there are other dimensions, absent from this model, for which an accountable politician is preferable to a bureaucrat.

that truthful equilibria are not sustainable independently of whether delegation occurs or not. Then, delegating the first action leads to a decrease in selection welfare. The effect on accountability is instead ambiguous, whereas the quality of the first action obviously increases. Therefore, a sufficient condition for delegation to a competent agent to be beneficial for overall welfare is that  $\alpha$  is sufficiently high.

## 5.4 Transparency of Consequences: Commentator Media

So far I assumed that voters have no feedback on the state of the world before elections. This serves the purpose of creating an environment in which all learning about the incumbent is done through his track-record. In many situations, however, voters have some information about the right policy choice. A fundamental role in this respect is played by the media.

In this section, therefore, I augment the baseline model with an additional player, which I call commentator media, in a similar fashion as what is done in the paper by Ashworth and Shotts (2010). To keep the analysis as simple as possible I assume that the media is only active at the end of  $t = 2$ , i.e. just before elections: this also reflects the fact that the coverage of politics in the vicinity of elections is particularly salient. I assume that the media is endowed with a signalling technology of accuracy  $q_M > \frac{1}{2}$  (conditionally independent of the signal that incumbents receive) and I abstract from strategic considerations on the part of the media assuming that before the election, the media truthfully reveals the realization of its signal  $s_M$ .

In this environment, voters have another piece of information to use when updating their beliefs on the incumbent type: as a result, reputation does not only depend on whether the incumbent flip-flopped or played consistently, but also on whether the media report matches the second period action (which I also call media endorsement). In such a setup it seems natural that having an informative signal on the state of the world will discipline the incumbent towards following their private signal. However, the following proposition demonstrates that, in some circumstances, increasing the accuracy of the commentator media can lead to more distorted behaviour by the incumbent.

**Proposition 6.** *Increasing the accuracy of the commentator media can increase flip-flopping avoidance, that is  $\sigma^*$  can be decreasing in  $q_M$ .*

*Proof.* See Appendix. □

It is interesting to describe in some detail the circumstances in which an increase in media accuracy increases the distortion to accountability. What is needed is a highly

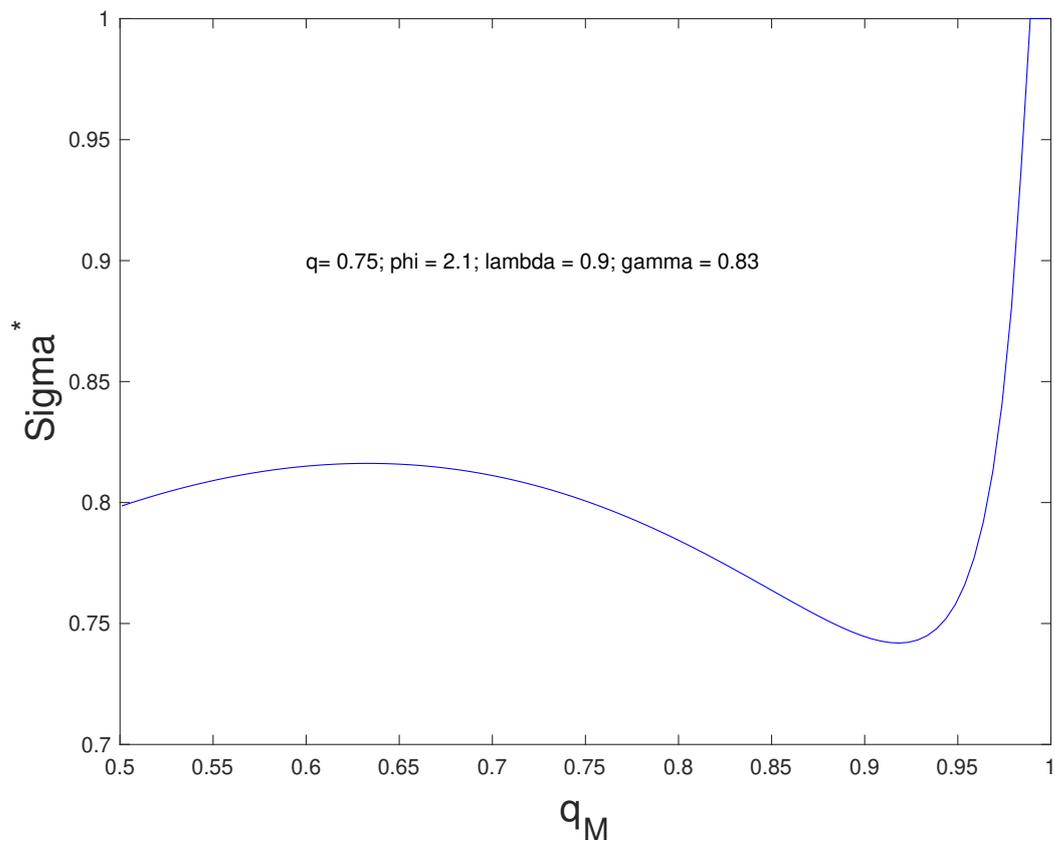


Figure 1: Example of Higher Quality Commentator Increasing Distortions

persistent environment (high  $\gamma$ ), a relatively high level of accuracy in the signal of incompetent politicians (high  $q$ ) and a large share of competent politicians (high  $\lambda$ ). In the example depicted in 1, increasing  $q_M$  from the lower bound of  $\frac{1}{2}$  initially increases  $\sigma^*$ , but as  $q_M$  increases further the effect reverses, to the extent that politicians distort their behaviour more when  $q_M = 0.9$  than when  $q_M = 0.5$ . Only when the media starts becoming extremely precise (with  $q_M$  close to 1) the distortion is eventually eliminated and  $\sigma^*$  converges to 1.

The intuition behind this result is what we might define as politicians gambling on the endorsement by the media. When  $\lambda$ ,  $\gamma$  and  $q$  are large, then an increase in media accuracy  $q_M$  sharply increases the payoff value of a successful gamble, which is the difference in reputation between being consistent and endorsed by the media and being a flip-flopper and opposed by the media. At the same time, the payoff from a failed gamble, which is the difference in reputation between a consistent track-record opposed by the media and a flip-flopping one endorsed by the media, decreases more slowly (but it suddenly drops when  $q_M$  is sufficiently high, meaning that when the media is very well informed, a non-endorsement is very costly in terms of reputation). Moreover, as long as the media is not too informative, the probability that the media signal matches the politician's flip-flopping signal remains close to  $\frac{1}{2}$ , making it likely for a politician avoiding a flip-flop to gamble successfully. This result relies heavily on the fraction of competent politicians  $\lambda$  being high.

The existence of potentially non-monotone responses of policy distortions to media accuracy can have interesting implications for issues such as the public subsidization of media outlets. In persistent environment with a prevalence of competent politicians, subsidizing media is only beneficial if a very precise commentary is achieved. This could for example suggest that concentrating resources into one or few high-quality outlets is better than subsidizing many average ones. In particular, once the media is sufficiently informative, returns to small increases in informativeness can be very large. Another possible application of this result concerns policy domains where most decision makers are highly competent, such as monetary policy. In such an environment, at times when persistence is high (for example periods of macroeconomic stability), even a very accurate commentator media outlet might lead to significant distortions, providing a further rationale for central bank independence.

Finally, notice that when an increased accuracy of the commentator media increases flip-flopping avoidance, accountability welfare decreases but selection welfare is still likely to improve, so I am not able to conclude that a more accurate commentator media is *tout*

court detrimental to welfare.

## 5.5 Transparency of Consequences: Poll-Matching

In this section I analyze the game in which a public signal concerning the state of the world in the second period is released *before* the politician takes the second action. In this modified game, the incumbent has one additional piece of information before choosing the second action  $a_2$ . In particular, I assume that the signal  $z$  (which I will call the poll) has accuracy  $Pr(z = \omega_2) = p = q$ , i.e. the same accuracy as the signal of the incompetent incumbent, and that it is conditionally independent with respect to the signal of the incumbent. Moreover, unlike in the benchmark model, the accuracy of the signal of the competent incumbent is now  $q_H = h < \bar{h}$ , where the upper bound  $\bar{h}$  will be defined later. In this setup, the optimal action in the second period does not only depend on the private signal of the incumbent, but also on the poll. In particular, the truthful equilibrium benchmark does not involve the politician always following his private signal: if the poll confirms the first action of the politician (that is to say  $z = s_1 = a_1$ ), then the optimal action for both types is to match the poll independently of the private signal.<sup>20</sup> If the poll instead calls for a flip-flop, the optimal action depends on the private signal  $s_2$ . In this setup, flip-flopping signals the fact that the private signal matches the state, i.e.  $s_2 = z$ ; therefore, flip-flopping to match the poll delivers a better reputation than being consistent but contradicting the poll. Hence, politicians have an incentive to flip-flop whenever the poll calls for it, and in equilibrium incompetent politicians mix by flip-flopping to match the poll even when their signal calls for sticking to the previous action. I summarize this result in the following proposition.

**Proposition 7.** *If  $h \in (p, \bar{h}]$ ,  $q = p \in [p, h]$  and  $\phi$  is higher than  $\max\{\bar{\phi}_z, \phi_{zz}\}$ , there exists a partially truthful equilibrium of the game such that  $a_1 = s_1$  for both types and:*

$$a_2 = \left\{ \begin{array}{ll} m, & \text{if } z = a_1 \\ s_2, & \text{if } z \neq a_1 \end{array} \right\}$$

*for the competent type, whereas when  $z \neq s_1$  and  $s_2 \neq z$ , the incompetent politician plays  $a_2 = z$  with probability  $1 - \sigma_z^*$ .*

*Proof.* See Appendix. □

---

<sup>20</sup>For this reason we need  $q_H$  to be sufficiently low, since if  $q_H$  is too large then the high type would choose  $a_2 = s_2$  independently of  $z$ .

Compared to the benchmark model, in this partially truthful equilibrium a flip-flopping track-record can be the result of opportunistic posturing carried out by an incompetent politician to match the public poll. Unlike in the benchmark model, however, flip-flopping is now beneficial for a politician's reputation, given that it only happens to match an informative poll. As a result, it still remains true that flip-flopping per se signals incompetence, but when the poll is informative enough, flip-flopping and matching the result of the poll is a signal of competence.

## 6 Conclusion

This paper describes the circumstances in which career concerned politicians have the incentive to inefficiently stick to their previous policy positions in order to avoid the stigma of flip-flopping. This happens because voters are aware that policy shifts are more likely to be performed by incompetent leaders. The incentive to avoid efficient policy shifts damages voters' welfare both in terms of policy effectiveness and selection of competent candidates through elections.

In other words, this paper rationalizes the conventional wisdom that flip-flopping is bad for the reputation of a politician; however, my results also suggest that the level of flip-flopping delivered by electoral competition might be excessively low and that democracy could benefit from politicians being more willing to change their mind.

An additional interesting implication of the mechanism described in the paper is that changes in the fundamentals driving policy choices are likely to bring to leadership change, but incompetent incumbents are more likely to remain in office after a change of state. Therefore, the probability of having a competent leader in office after a change in the state of the world might be lower.

To sum up, my analysis adds to the literature on the negative effects of electoral competition: elections enable voters to retain the politicians they believe to be competent, but at the same time they give potentially distortive incentives to politicians. In this context, therefore, it might be optimal to limit electoral incentives, by either setting up a single term limit rule or handing over some decision power to a judiciary.

Another implication of my model concerns the observability of politicians' actions: making some of the actions taken by the politician unobservable (or observable with noise) can eliminate the incentives for inefficient policy choices. One interpretation of this result is that making some parts of a policy-making process secret might be beneficial: one can think of secret voting on parliaments, closed-doors committee meetings, or non-disclosure of the early stages of a policy process. Another interpretation has to do with the media environment and the fact technological developments such as social media have made it very easy for voters to know what a politician previously did or stated: according to my model, this might lead to distortions and decrease the responsiveness of policies to information.

A similar trade-off between noise and distortion arises when the first action is observable but delegated to an independent agent: even if the agent is less competent than the incumbent, delegation might be beneficial by reducing the incentives to distort the second action. As a matter of fact, electoral selection is better under delegation to an

incompetent rather than a competent agent.

Finally, I also show that when media act as commentators of the quality of a policy choice, an increase in the accuracy of the media signal can in some circumstances incentivize rather than deter the distortive behaviour of politicians. In other words, giving voters more information to evaluate politicians' track records can backfire. In a similar context, if an informative poll is released before the politician takes the second action, there can be an incentive to flip-flop to match the poll.

An interesting direction for future research would in my opinion be to consider more in detail the role of flip-flopping as an electoral campaign tool used by a strategic challenger.

## References

- AGHION, P., AND M. O. JACKSON (2016): “Inducing Leaders to Take Risky Decisions: Dismissal, Tenure, and Term Limits,” *American Economic Journal: Microeconomics*, 8(3), 1–38.
- AGRANOV, M. (2016): “Flip-Flopping, Primary Visibility, and the Selection of Candidates,” *American Economic Journal: Microeconomics*, 8(2), 61–85.
- ASHWORTH, S., AND K. W. SHOTTS (2010): “Does informative media commentary reduce politicians’ incentives to pander?,” *Journal of Public Economics*, 94(11-12), 838–847.
- CALVINO, I. (1969): *Prima che tu dica pronto*. Mondadori.
- CANES-WRONE, B., M. C. HERRON, AND K. W. SHOTTS (2001): “Leadership and Pandering: A Theory of Executive Policymaking,” *American Journal of Political Science*, 45(3), 532–550.
- DEWAN, T., AND R. HORTALA-VALLE (2014): “Electoral Competition, Control and Learning,” *Working Paper*, pp. 1–31.
- DOHERTY, D., C. M. DOWLING, AND M. G. MILLER (2015): “When is Changing Policy Positions Costly for Politicians? Experimental Evidence,” *Political Behavior*, 38(2), 455–484.
- FEARON, J. D. (1999): “Electoral Accountability and the Control of Politicians: Selecting Good Types versus Sanctioning Poor Performance,” in *Democracy, Accountability and Representation*, ed. by A. Przeworski, and S. C. Stokes, chap. 2, pp. 55–97. Cambridge University Press.
- FOX, J., AND S. V. JORDAN (2011): “Delegation and accountability,” *The Journal of Politics*, 73(3), 831–844.
- FOX, J., AND M. C. STEPHENSON (2011): “Judicial review as a response to political posturing,” *American Political Science Review*, 105(2), 397–414.
- FRISELL, L. (2009): “A theory of self-fulfilling political expectations,” *Journal of Public Economics*, 93(5–6), 715 – 720.
- FU, Q., AND M. LI (2014): “Reputation-concerned policy makers and institutional status quo bias,” *Journal of Public Economics*, 110, 15–25.

- GENTZKOW, M., AND J. M. SHAPIRO (2006): “Media Bias and Reputation,” *Journal of Political Economy*, 114(2), 280–316.
- HUMMEL, P. (2010): “Flip-flopping from primaries to general elections,” *Journal of Public Economics*, 94(11–12), 1020 – 1027.
- LEVENDUSKY, M. S., AND M. HOROWITZ (2012): “When Backing Down Is the Right Decision: Partisanship, New Information and Audience Costs,” *Journal of Politics*, 74, 323–338.
- LEVY, G. (2004): “Anti-herding and strategic consultation,” *European Economic Review*, 48(3), 503–525.
- (2007): “Decision Making in Committees: Transparency, Reputation, and Voting Rules,” *The American Economic Review*, 97(1), 150–168.
- LI, W. (2007): “Changing One’s Mind When the Facts Change: Incentives of Experts and the Design of Reporting Protocols,” *The Review of Economic Studies*, 74(4), 1175–1194.
- MAJUMDAR, S., AND S. W. MUKAND (2004): “Policy Gambles,” *American Economic Review*, 94(4), 1207–1222.
- MASKIN, E. S., AND J. TIROLE (2004): “The Politician and the Judge: Accountability in Government,” *American Economic Review*, 94(4), 1034–1054.
- PRAT, A. (2005): “The Wrong Kind of Transparency,” *The American Economic Review*, 95(3), 862–877.
- PRENDERGAST, C., AND L. STOLE (1996): “Impetuous Youngsters and Jaded Old-Timers: Acquiring a Reputation for Learning,” *The Journal of Political Economy*, 104(6), 1105–1134.
- SMART, M., AND D. M. STURM (2013): “Term limits and electoral accountability,” *Journal of public economics*, 107, 93–102.
- TAVITS, M. (2007): “Principle vs. pragmatism: Policy shifts and political competition,” *American Journal of Political Science*, 51(1), 151–165.
- TOMZ, B. M., AND R. V. HOUWELING (2012): “Political Repositioning,” *Working Paper*, pp. 1–34.
- TOMZ, M., AND R. P. VAN HOUWELING (2012): “Candidate Repositioning,” pp. 1–43.

# A Proofs

## Proposition 1

*Proof.* Let's consider the choice of action at  $t = 2$ , politicians have already taken action  $a_1$  and they know that if re-elected, which happens with probability  $\frac{1}{2} + \frac{\mu(a_1, a_2) - \lambda_O}{2}$ , they will get  $2\phi$ . On top of that, politicians get utility of 1 whenever they match the state of the world, the probability of which depends on the posterior belief, denoted by  $Pr(\omega_2 = s_2 | s_1, q_\theta) = \rho_2(s_2, s_1, q_\theta)$ . In order to slightly simplify notation, denote the probability of winning given reputation  $\mu$  by  $r(\mu)$ . So  $r(\mu) = \frac{1}{2} + \frac{\mu - \lambda_O}{2}$  as  $r(\mu)$ . Given that there are only two actions available, the politicians will calculate the reputation associated to each action and follow his signal if and only if:

$$\rho_2(s_2, s_1, q_\theta) + r(\mu(a_2 = s_2, s_1))2\phi \geq [1 - \rho_2(s_2, s_1, q_\theta)] + r(\mu(a_2 \neq s_2, s_1))2\phi.$$

It follows that in order to have  $\sigma(s_2) = 1$ , the above condition needs to hold for both types and both signal realizations given any of the two possible choices  $a_1$ , which results in a set of 8 inequalities.

However, it is immediate to notice that whenever  $a_2 = s_2$  is the most reputable action, i.e.  $\mu(a_2 = a_1, a_1) \geq \mu(a_2 \neq a_1, a_1)$ , following the signal is unquestionably optimal. Since consistency has a better reputation than flip-flopping, it follows that whenever  $s_2 = a_1$ ,  $a_2 = s_2$ , for each  $a_1$  and each type. This means that we are left with 4 conditions. Moreover, since  $\mu^T(0, 0) = \mu^T(1, 1)$  and  $\mu^T(1, 0) = \mu^T(0, 1)$  and the same holds for  $\rho_2(1, 0) = \rho_2(0, 1)$ , we are left with only two conditions, one for each type:

$$\rho_2(s_2 \neq s_1, q_\theta) + r(\mu_F^T)2\phi \geq [1 - \rho_2(s_2 \neq s_1, q_\theta)] + r(\mu_C^T)2\phi,$$

which can be rearranged to:

$$\frac{2\rho_2(s_2 \neq s_1, q_\theta) - 1}{\phi} \geq \mu_C^T - \mu_F^T.$$

Thanks to the single-crossing property, the binding constraint for a truthful equilibrium is the condition concerning the incompetent politician: whenever the incompetent politician follows his signal, or is at least indifferent, the competent politician does, too. Conversely, if the competent politician is indifferent or he doesn't follow his signal, the incompetent will also not follow it.

As a result, we know that if  $\frac{2\rho_2(s_2 \neq s_1, q) - 1}{\phi} \geq \mu_C^T - \mu_F^T$ , politicians will follow their signal at  $t = 2$ . Rearranging we get the condition on  $\phi$ ,  $\phi \leq \frac{2\rho_2(s_2 \neq s_1, q) - 1}{\mu_C^T - \mu_F^T}$ .

Let's now see what happens at  $t = 1$ . If politicians know that they are going to follow their signal at  $t = 2$ , then the dominant strategy at  $t = 1$  is to follow their signal. Since all that matters is being consistent versus flip-flopping, then given the persistence of the state, it is more likely to end up in the favourable situation of playing consistently by following one's signal in the first period. □

### **Trembling-hand perfection refinement eliminates pooling equilibria**

*Proof.* In a pooling equilibrium, both types of politicians pool on the same action (that is they play said action with probability one) in either one of the two periods or in both, independently of the signals received. In other words, in these equilibria only one track-record is played in equilibrium. For the off-equilibrium track-records, Bayes rule does not offer any restriction in beliefs. It follows that if the reputation attached to any off-equilibrium track-record is sufficiently bad and electoral concerns are sufficiently high, politicians will not have any incentive to deviate from the pooling track-record. Sufficient conditions for any pooling equilibrium to be sustainable read:

$$2\rho_1 - 1 + Pr(s_2 = s_1 | s_1)(2\bar{\rho}_2 - 1) < (r(\mu = \lambda) - r(\mu = 0))2\phi$$

and

$$2\bar{\rho}_2 - 1 < r(\mu = \lambda) - r(\mu = 0))2\phi,$$

where  $\bar{\rho}_2 = Pr(\omega_2 = s_2 | s_1, s_2 = s_1, \theta = L)$ . Let's now introduce the trembling-hand perfection requirement. Assume that with probability  $\epsilon > 0$  close to zero, a politician willing to play action  $a$  plays action  $a'$  instead. Assume that  $\epsilon$  is the same for both types. Take a pooling equilibrium. In any period, with probability  $\epsilon$  the voter observes an action different from that on the pooling track-record. Since both politicians have the same strategy and  $\epsilon$  is the same for both types, then the reputation the voter must attach to actions outside the pooling track-record has to be  $\lambda$ , the same as the reputation of the pooling track-record. However, this cannot happen in equilibrium, because if the reputation of any track-record is the same, then incumbents have an incentive to always follow their signal. In other words, no pooling equilibrium can survive the trembling hand perfection requirement. □

## Proof of Theorem 1

*Proof.* I split the proof in two parts. First of all I characterize the symmetric partially truthful equilibrium.

**Claim 1:** A partially truthful equilibrium exists.

In this equilibrium,  $\mu(0, 0) = \mu(1, 1) \equiv \mu_C$  and  $\mu(0, 1) = \mu(1, 0) \equiv \mu_F$  by symmetry; moreover,  $\mu_C > \mu_F$ . Assume for now that incumbents always follow their signal at  $t = 1$ . Consider period  $t = 2$ : after following their signal in period 1, in period 2 the incumbent has to decide whether to follow his signal or not. When  $s_2 = a_1$ , the incumbent always follows his signal, because:

$$\rho_2(s_1, s_2 = s_1, \theta) + r(\mu_C)2\phi > (1 - \rho_2(s_1, s_2 = s_1, \theta)) + r(\mu_F)2\phi$$

since  $\mu_C > \mu_F$  insures that  $r(\mu_C) > r(\mu_F)$  and  $\rho_2 > 1 - \rho_2$  by the decision relevance of signals. However, when  $s_2 \neq a_1$ , the incumbent has a tradeoff. From Proposition 1, a truthful equilibrium requires that:

$$\rho_2(s_1, s_2 \neq s_1, L) + r(\mu_F)2\phi \geq (1 - \rho_2(s_1, s_2 \neq s_1, L)) + r(\mu_C)2\phi,$$

from which the upper bound  $\bar{\phi}$  was derived. Moreover, we know from Lemma 2 (single-crossing property), that the following holds:

$$\begin{aligned} \rho_2(s_1, s_2 \neq s_1, L) + r(\mu_F)2\phi &\geq 1 - \rho_2(s_1, s_2 \neq s_1, L) + r(\mu_C)2\phi \Rightarrow \\ \rho_2(s_1, s_2 \neq s_1, H) + r(\mu_F)2\phi &> 1 - \rho_2(s_1, s_2 \neq s_1, H) + r(\mu_C)2\phi \end{aligned}$$

and

$$\begin{aligned} \rho_2(s_1, s_2 \neq s_1, H) + r(\mu_F)2\phi &\leq 1 - \rho_2(s_1, s_2 \neq s_1, H) + r(\mu_C)2\phi \Rightarrow \\ \rho_2(s_1, s_2 \neq s_1, L) + r(\mu_F)2\phi &< 1 - \rho_2(s_1, s_2 \neq s_1, L) + r(\mu_C)2\phi \end{aligned}$$

This means that we are left with three possibilities: both politicians play  $a_2 \neq s_2$  when that involves flip-flopping, or the high type mixes between  $a_2 = s_2$  and  $a_2 \neq s_2$ , or the high type always plays  $a_2 = s_2$  and the low type mixes. I will now prove that the first two cannot be part of an equilibrium. Assume that both politicians play  $a_2 \neq s_2$  when  $s_2 \neq a_1$ . Then, nobody would flip-flop and therefore, in a candidate equilibrium,  $\mu_C = \mu_F$  would hold; in this case, however, the optimal strategy for the incumbent is to follow his

signal. Consider now the other case, i.e. that when  $s_2 \neq a_1$ , the high type mixes between  $a_2 = s_2$  and  $a_2 \neq s_2$  while the low type always plays  $a_2 \neq s_2$ . If this were an equilibrium, then flip-flopping would reveal the high type, and therefore  $\mu_F = 1 > \mu_C$ . This cannot be part of an equilibrium, since the low type incumbent would profitably deviate by flip-flopping. It follows that the only possibility when  $s_2 \neq a_1$  is that the high type follows his signal whereas the low type mixes between the two actions. The low type mixes when the following holds:

$$\frac{2\rho_2(s_1, s_2 \neq s_1, L) - 1}{\phi} = \underbrace{\frac{\lambda\gamma}{\lambda\gamma + (1-\lambda)(1-A\sigma^*)}}_{\mu_C} - \underbrace{\frac{\lambda(1-\gamma)}{\lambda(1-\gamma) + (1-\lambda)A\sigma^*}}_{\mu_F}$$

In order to show existence and uniqueness of such an equilibrium, consider  $\sigma^* \in [0, 1]$ . It holds from Proposition 1 that at  $\sigma^* = 1$ ,  $\frac{2\rho_2(s_1, s_2 \neq s_1, L) - 1}{\phi} < \mu_C - \mu_F$ . At the same time, at  $\sigma^* = 0$  it has to be that  $\frac{2\rho_2(s_1, s_2 \neq s_1, L) - 1}{\phi} > \mu_C - \mu_F$ , since in that case  $\mu_F = 1$ . By continuity of  $\mu_C - \mu_F$ , an equilibrium with  $\sigma^* \in (0, 1)$  exists. Moreover, notice that  $\mu_C$  is strictly increasing in  $\sigma^*$  while  $\mu_F$  is strictly decreasing in  $\sigma^*$ , and therefore there is a unique equilibrium mixing probability  $\sigma^*$ . Finally, consider the action choice at  $t = 1$ . It holds that  $\mu_C > \mu_F$  and the incumbent knows he is going to follow his signal when that implies receiving a consistent reputation  $\mu_C$ . Moreover, the incumbent knows that, because of the persistence of the state of the world, given the realization of  $s_1$  it is more likely for him to receive  $s_2 = s_1$ . As a result, following the signal at  $t = 1$  is optimal both in terms of instantaneous payoff (it is more likely to match the state) and in terms of future payoff (it allows the politician to receive the higher reputation more often without having to distort his action). Mathematically, denote by  $\pi = Pr(s_2 = s_1 | s_1, \theta)$ . Notice that since  $\gamma > \frac{1}{2}$ ,  $\pi > \frac{1}{2}$  and notice that

$$\pi\rho_2(s_1, s_2 = s_1, \theta) + (1 - \pi)\rho_2(s_1, s_2 \neq s_1, \theta) = \rho_1(s_1, \theta)\gamma + (1 - \rho_1(s_1, \theta))(1 - \gamma).$$

Denote also by  $\rho_1 = Pr(\omega_2 = s_1 | s_1, \theta)$ ,  $\bar{\rho}_2 = Pr(\omega_2 = s_1 | s_2 = s_1, \theta)$  and by  $\rho_2 = Pr(\omega_2 = s_1 | s_2 \neq s_1, \theta)$ . With this notation, following one's signal at  $t_1$  requires the following condition to hold:

$$\rho_1 + \pi(\bar{\rho}_2 + r(\mu_C)2\phi) + (1 - \pi)(\rho_2 + r(\mu_F)2\phi) \geq (1 - \rho_1) + \pi(\bar{\rho}_2 + r(\mu_F)2\phi) + (1 - \pi)(\rho_2 + r(\mu_C)2\phi)$$

which can be rearranged to

$$(2\rho_1 - 1) \geq -(2\pi - 1)(r(\mu_C) - r(\mu_F))2\phi$$

which always holds. Hence both types follow their signal at  $t = 1$ .

To sum up, I have showed the existence of a unique symmetric partially truthful equilibrium (it would be truthful if  $\phi \leq \bar{\phi}$ ) with  $\mu_C > \mu_F$ . Both politicians follow their signal at  $t = 1$ . At  $t = 2$ , the high type always follows his signal, whereas the low type follows his signal with probability 1 when if  $s_2 = s_1$  and mixes playing  $a_2 = s_2$  with probability  $\sigma^* \in (0, 1)$  when  $s_2 \neq s_1$ .

**Claim 2:** The symmetric non-pooling equilibrium is unique.

First of all I elaborate on the concept of symmetric equilibrium I use. In this context, symmetry means that if the information structure is symmetric across signal realizations, the equilibrium play across signal realizations needs to also be symmetric. That means that taken two signal realizations after which the information available to players (including the expectations about future realizations of signals) is the same up to the labeling of the states of the world, the equilibrium play must be the same. In other words, symmetry rules out situations in which the strategies played after one signal realization depend arbitrarily on the labeling of the state of the world as 0 or 1. Specifically, in this game the only additional restriction added by the symmetry requirement is that  $Pr(a_1 = 0|s_1 = 0) = Pr(a_1 = 1|s_1 = 1)$ ; this condition, together with the no-pooling condition, implies symmetry also in the second period<sup>21</sup>, which ultimately results in reputations being symmetric, that is  $\mu(0, 0) = \mu(1, 1) \equiv \mu_C$  and  $\mu(0, 1) = \mu(1, 0) \equiv \mu_F$ .

In Claim 1 I have shown that when reputations are symmetric and  $\mu_C > \mu_F$ , then the equilibrium is the partially truthful one I characterized (or truthful, depending on  $\phi$ ). I will now relax the assumption that  $\mu_C > \mu_F$  and prove that remaining in the realm of symmetric equilibria, there exists no equilibrium where  $\mu_F > \mu_C$ . Suppose by contradiction that  $\mu_F > \mu_C$ . There are two cases: first, suppose that  $a_1 = s_1$  for both types. In this case, for analogous reasons as those explained in Claim 1, the only candidate equilibrium is one where competent politicians play truthfully and incompetent ones either play truthfully or mix, in this case when  $s_2 = s_1$  instead of  $s_2 \neq s_1$ . If incompetent politicians mix, then reputations are just like in the truthful equilibrium and  $\mu_C > \mu_F$ . If they mix, incompetent politicians play  $\tau_F$  even more often than in a truthful equilibrium and  $\mu_F$  decreases even further, while  $\mu_C$  increases. Therefore,  $\mu_C > \mu_F$ . It follows that such an equilibrium cannot exist. In other words, if  $a_1 = s_1$  for both politicians the unique equilibrium is the (partially) truthful one.

<sup>21</sup>That is that for a given  $i, j, k \in \{0, 1\}$ , with  $i' \neq i$ ,  $j' \neq j$  and  $k' \neq k$ ,  $Pr(a_2 = i|s_2 = j, s_1 = k) = Pr(a_2 = i'|s_2 = j', s_1 = k')$

Let's consider therefore the other possible case, i.e. the one where  $a_1 \neq s_1$  for some type and some signal realization. In particular, for  $a_1 \neq s_1$  to be an optimal action it has to be that  $(2\rho_1(\theta) - 1) \leq (2\pi(\theta) - 1)(\mu_F - \mu_C)\phi$  for some incumbent type, where as before  $\pi(\theta) = Pr(s_2 = s_1 | s_1, \theta) = [\gamma q(\theta) + (1 - \gamma)(1 - q(\theta))]q(\theta) + [\gamma(1 - q(\theta)) + (1 - \gamma)q(\theta)](1 - q(\theta)) < q(\theta)$ . Consider the subgame after the incumbent played  $a_1 \neq s_1$ : since  $\mu_F > \mu_C$ , at  $t = 2$  the optimal choice is to follow the signal when  $s_2 = s_1$  whereas a trade-off arises when  $s_2 \neq s_1$ . In equilibrium,  $\phi(\mu_F - \mu_C) \leq 2\rho_2(\theta) - 1$ , and therefore not following the signal in the first period requires:

$$(2\rho_1(\theta) - 1) \leq (2\pi(\theta) - 1)(2\rho_2(\theta) - 1)$$

However, this is impossible, since  $\frac{2\rho_1(\theta) - 1}{2\pi(\theta) - 1} > 2\rho_2(\theta) - 1$ , for each  $q(\theta) > 1/2$ , given that  $\pi(\theta) < q(\theta)$ . □

## Proof of Corollary 2

*Proof.* I prove the three results in sequence:

- i. Since  $\mu_C > \mu_F$ , a leadership change is more likely when the politician has a flip-flopping rather than a consistent reputation. In order to see the first result, compare the probability of a  $\mu_F$  reputation conditional on the state changing versus the state not changing:

$$\underbrace{\lambda + (1 - \lambda)(q^2 + (1 - q)^2)\sigma^*}_{Pr(\mu_F | \text{change})} > \underbrace{(1 - \lambda)2q(1 - q)\sigma^*}_{Pr(\mu_F | \text{no change})}$$

This clearly holds since  $q^2 + (1 - q)^2 > 2q(1 - q)$ .

- ii. A change in the state of the world increases the chance that both types of incumbent are fired. The probability of having a competent leader in office after elections is:

$$\lambda Pr(e = r | \text{change}, \theta = H) + \lambda_0 [\lambda Pr(e = f | \text{change}, \theta = H) + (1 - \lambda) Pr(e = f | \text{change}, \theta = L)]$$

when the state changes and

$$\lambda Pr(e = r | \text{no ch.}, \theta = H) + \lambda_0 [\lambda Pr(e = f | \text{no ch.}, \theta = H) + (1 - \lambda) Pr(e = f | \text{no ch.}, \theta = L)]$$

when it doesn't change. Notice that

$$Pr(e = r|ch., \theta = H) = Pr(e = r|\mu_F) = 1 - Pr(e = f|\mu_F)$$

and thus  $Pr(e = f|ch., \theta = H) = Pr(e = f|\mu_F)$ , whereas

$$Pr(e = f|ch., \theta = L) = Pr(\mu_F|ch., \theta = L)Pr(e = f|\mu_F) + Pr(\mu_C|ch., \theta = L)Pr(e = f|\mu_C)$$

and analogously

$$Pr(e = f|no ch., \theta = L) = Pr(\mu_F|no ch., \theta = L)Pr(e = f|\mu_F) + Pr(\mu_C|no ch., \theta = L)Pr(e = f|\mu_C)$$

Denote  $\beta = Pr(\mu_F|change, \theta = L)$  and  $\beta' = Pr(\mu_F|not change, \theta = L)$  and the probability of having a competent leader in power after a state change is higher if and only if:

$$\lambda_O > \frac{\lambda}{(\beta - \beta')(1 - \lambda) + \lambda}.$$

The intuition is the following: conditional on the state of the world changing, competent politicians always get the correct signal, and since in equilibrium they play truthfully, they always end up with a reputation of  $\mu_F$  after a change in the state. Concerning the incompetent politicians, they sometimes receive a wrong signal and on top of that, they sometimes avoid flip-flopping. As a result, competent politicians are more likely to be fired after a change in the state whereas they are more likely to not be fired when the state did not change.

- iii. Finally, the third point is a direct consequence of the fact that when the state changes, a competent incumbent always gets the lower reputation  $\mu_F$ , whereas an incompetent incumbent sometimes gets the better reputation  $\mu_C$ . Conversely, when the state does not change competent incumbents always receive  $\mu_C$ , whereas incompetent incumbents sometimes receive  $\mu_F$ .

□

## Proof of Comparative Statics 1

*Proof.* Recall that

$$\bar{\phi} = \frac{2\rho_2 - 1}{\mu_C^T - \mu_F^T}$$

,

$$\rho_2 = \frac{[(1 - q)\gamma + q(1 - \gamma)]q}{[(1 - q)\gamma + q(1 - \gamma)]q + [q\gamma + (1 - q)(1 - \gamma)](1 - q)}$$

and

$$\mu_C^T - \mu_F^T = \frac{\lambda\gamma}{\lambda\gamma + (1-\lambda)(1-A(q))} - \frac{\lambda(1-\gamma)}{\lambda(1-\gamma) + (1-\lambda)A(q)},$$

where  $A(q) = (1-\gamma)[q^2 + (1-q)^2] + \gamma 2q(1-q)$ . Therefore, when  $q$  increases,  $\rho_2$  increases (the incompetent incumbent becomes better informed), whereas  $\mu_C^T - \mu_F^T$  decreases. This can be seen by noticing that  $A(q)$  decreases in  $q$ . The two types become more similar and therefore a flip-flop is less telling of the politician being of the incompetent type. When  $\gamma$  increases, it can be seen from the expression above that  $\rho_2$  decreases, since the opposite signal received in the first period matters more, whereas  $\mu_C^T - \mu_F^T$  increases: as a matter of fact,  $\frac{1-A(q)}{\gamma}$  decreases, so  $\mu_C$  increases, whereas  $\frac{A(q)}{1-\gamma}$  increases, so  $\mu_F$  decreases. In other words, flip-flopping becomes a more accurate signal of incompetence. Finally, when  $\lambda$  increases,  $\rho_2$  does not change but  $\mu_C^T - \mu_F^T$  can move in either direction, and therefore the effect on  $\bar{\phi}$  is ambiguous.  $\square$

### Proof of Lemma 3

*Proof.* As far as  $t = 1$  and  $t = 2$  are concerned, since the competent politician follows his signal, which is perfectly accurate, he always takes the right decision. The incompetent politician, instead, follows the signal in the first period, taking the right decision with probability  $q$ , but in the second period, if the signal indicates flip-flopping as the optimal action, he contradicts it with probability  $1 - \sigma^*$ . As a result, the accuracy of the incompetent incumbent's signal is  $(1-A)\bar{\rho}_2 + A(\sigma^*\rho_2 + (1-\sigma^*)(1-\rho_2)) = q - (1-\sigma^*)(2\rho_2 - 1)A$ , to get which I used the fact that  $(1-A)\bar{\rho}_2 + A\rho_2 = q$ , since  $q = Pr(\omega = s_1|s_1)$ . As far as the valence draw is concerned, the expression for expected welfare is the following:

$$\mathbb{E}[v|e = f] = \sum_{\tau \in \{C, F\}} Pr(\tau) \int_{b(\mu - \lambda_O)}^b \frac{v}{2b} dv.$$

Solving the integral yields:

$$\mathbb{E}[v|e = f] = \sum_{\tau \in \{C, F\}} Pr(\mu(\tau)) \frac{b}{4} [1 - (\mu - \lambda_O)^2]$$

which can be rewritten as:

$$\mathbb{E}[v|e = f] = \mathbb{E}_\mu \frac{b}{4} [1 - (\mu - \lambda_O)^2].$$

Consider now the term  $\mathbb{1}_{\theta_e=H}b$ . First of all,

$$Pr(e = f) = \mathbb{E}_\mu \left[ \frac{1}{2} - \frac{\mu - \lambda_O}{2} \right]$$

and

$$Pr(e = r|\theta = H) = \frac{1}{2} - \frac{\lambda_O}{2} + \frac{Pr(\tau = C|H)\mu_C + Pr(\tau = F|H)\mu_F}{2},$$

where  $Pr(\tau = C|H) = \gamma$  and  $Pr(\tau = F|H) = 1 - \gamma$ . Taking expectations with respect to  $\mu$  where necessary, and considering that

$$\mathbb{E}\mu = Pr(\tau = C)\mu_C + Pr(\tau = F)\mu_F = \lambda,$$

the welfare expression can be rewritten as

$$\begin{aligned} W &= [\lambda + (1 - \lambda)q] + [\lambda + (1 - \lambda)\tilde{q}] + \frac{b}{4}[1 - \mathbb{E}_\mu(\mu - \lambda_O)^2] \\ &\quad + b\lambda_O \left( \frac{1}{2} - \frac{\lambda - \lambda_O}{2} \right) + b\lambda \left[ \frac{1}{2} - \frac{\lambda_O}{2} + \frac{\gamma\mu_C + (1 - \gamma)\mu_F}{2} \right]. \end{aligned}$$

Consider now the expression  $\mathbb{E}\mu^2$ , where the expectation is taken over the distribution of  $\mu$ :

$$\mathbb{E}\mu^2 = Pr(\tau = C)\mu_C^2 + Pr(\tau = F)\mu_F^2$$

Using the fact that  $Pr(\tau = C) = \frac{\lambda\gamma}{\mu_C}$  and  $Pr(\tau = F) = \frac{\lambda(1-\gamma)}{\mu_F}$  I can rewrite the former expression to get:

$$\mathbb{E}\mu^2 = \lambda(\gamma\mu_C + (1 - \gamma)\mu_F).$$

Using this result and multiplying out all the terms, the welfare expression can finally be expressed as:

$$W = [\lambda + (1 - \lambda)q] + [\lambda + (1 - \lambda)\tilde{q}] + b \left[ \frac{1}{4} + \frac{\lambda_O^2}{4} - \frac{\lambda\lambda_O}{2} + \frac{\lambda + \lambda_O}{2} + \frac{\mathbb{E}\mu^2}{4} \right].$$

□

## Proof of Lemma 4

*Proof.* The average reputation is constant in equilibrium and equal to  $\lambda$ , hence  $Pr(\tau = C)\mu_C + Pr(\tau = F)\mu_F = \lambda$ . This allows us to rewrite  $Pr(\tau = C) = \frac{\lambda - \mu_F}{\mu_C - \mu_F}$ . From Lemma 3 we know that selection improves if and only if the second moment of the distribution of reputation, denoted by  $\mathbb{E}\mu^2 = Pr(\tau = C)\mu_C^2 + (1 - Pr(\tau = C))\mu_F^2$ , increases. Moreover, in equilibrium it holds that  $\mu_C \geq \mu_F$  and  $Pr(\tau = C) > Pr(\tau = F)$ . Taking the expression

for the second moment and substituting in the constant-mean constraint yields:

$$\mathbb{E}\mu^2 = \frac{\lambda - \mu_F}{\mu_C - \mu_F} \mu_C^2 + \left[ 1 - \frac{\lambda - \mu_F}{\mu_C - \mu_F} \right] \mu_F^2$$

Let's now write the expression of a contour line, in the  $(\mu_C, \mu_F)$  plane, along which the value of the second moment is constant; using implicit differentiation one gets:

$$\frac{\partial \mu_F}{\partial \mu_C} = \frac{\lambda - \mu_F}{\mu_C - \lambda}$$

It can be seen that this contour line is increasing in  $\mu_C$  and concave. This means that the second moment increases by increasing  $\mu_C$  and decreasing  $\mu_F$ . Moreover, notice that since in equilibrium  $\mu_C > \mu_F$ ,  $Pr(\mu_C) > Pr(\mu_F)$  and  $Pr(\mu_C)\mu_C + Pr(\mu_F)\mu_F = \lambda$ , then  $\mu_C - \lambda < \lambda - \mu_F$ , i.e.  $\lambda > \frac{\mu_C + \mu_F}{2}$ . So  $\frac{\partial \mu_F}{\partial \mu_C} \geq 1$  at any equilibrium point  $(\mu_C, \mu_F)$ . As a result, therefore, whenever  $\mu_C - \mu_F$  increases and  $\mu_C$  increases, the second moment  $\mathbb{E}\mu^2$  increases and thus selection welfare increases, too. The graphical intuition for the result is that the contour line of the second moment is increasing and concave with slope larger than 1, whereas the  $\mu_C - \mu_F$  isoquant is increasing with a slope of 1. Moreover, the second point at which the  $\mathbb{E}\mu^2$  isoquant crosses the  $\mu_C - \mu_F$  isoquant is never in the set of feasible  $(\mu_C, \mu_F)$ , given that at that point the slope of the  $\mathbb{E}\mu^2$  isoquant is less than 1.

A corollary of this result is that a necessary condition for  $\mathbb{E}\mu^2$  (and thus selection welfare) to increase when  $\mu_C - \mu_F$  decreases is that both  $\mu_C$  and  $\mu_F$  increase.  $\square$

### Proof of Welfare Result 1

*Proof.* When  $\phi$  increases, the benefits from office increase and this increases the accountability distortion. This happens because  $\mu_C - \mu_F$  is an increasing function of  $\sigma$ ; when  $\phi$  increases, therefore,  $\sigma^*$  decreases up to the point where  $\mu_C - \mu_F$  is equal to the new costs of deviating for the signal. This reasoning can be easily verified from the equilibrium condition, noting that the right-hand side increases as  $\sigma^*$  increases:

$$\frac{2\rho_2 - 1}{\phi} = \frac{\lambda\gamma}{\lambda\gamma + (1 - \lambda)(1 - A\sigma^*)} - \frac{\lambda(1 - \gamma)}{\lambda(1 - \gamma) + (1 - \lambda)A\sigma^*}$$

The decrease in  $\sigma^*$  negatively impacts welfare in the second period. In terms of selection of politicians, since a lower  $\sigma^*$  decreases  $\mu_C$  and increases  $\mu_F$ , then from Lemma 4 we know that this means that the second moment will decrease, hence also selection welfare worsens.  $\square$

### Proof of Welfare Result 2

*Proof.* Assume that  $q$  increases to  $q'$ . An increase in  $q$  moves  $A = (1 - \gamma)(q^2 + (1 - q)^2) + 2\gamma q(1 - q)$  down and  $\frac{2\rho_2 - 1}{\phi}$  up. As long as both  $\sigma^*(q)$  and  $\sigma^*(q')$  are strictly less than 1, then in equilibrium  $\frac{2\rho_2 - 1}{\phi} = \mu_C - \mu_F$  and therefore an increase of  $q$  to  $q'$  leads to a larger equilibrium value of  $\mu_C - \mu_F$ . The equilibrium value of  $\mu_C - \mu_F$  in response to a change in  $q$  is driven by  $A\sigma^*$ : as a result, a larger level of  $\mu_C - \mu_F$  can only occur when  $A\sigma^*$  increases. It follows that  $\mu_C$  increases and  $\mu_F$  decreases, and therefore by Lemma 4 we know that selection welfare improves. In terms of accountability welfare,  $\sigma^*(q') > \sigma^*(q)$  and  $q' > q$ , so not only incompetent politicians are better, but they also act in a less distorted way. Hence, accountability welfare increases, and therefore total welfare increases as well. Notice that once in a truthful equilibrium, an increase in  $q$  decreases  $\mu_C - \mu_F$ , since  $A$  keeps decreasing but  $\sigma$  cannot increase any further. As a result, the second moment decreases and the probability of having a competent politician in office in the second period decreases. However, bad politicians are better so the effect on selection welfare is ambiguous.  $\square$

### Proof of Welfare Result 3

*Proof.* Let's consider the case in which an increase in  $\lambda$  is such that the game remains in a partially truthful equilibrium. Accountability welfare can either increase or decrease. The reason is that, denoting by  $D_C \equiv \lambda\gamma + (1 - \lambda)(1 - A\sigma)$  and  $D_F = \lambda(1 - \gamma) + (1 - \lambda)A\sigma$ , the expression  $\frac{\partial(\mu_C - \mu_F)}{\partial\lambda} = \frac{A\sigma^*}{1 - \gamma} \frac{1}{D_F^2} - \frac{1 - A\sigma^*}{\gamma} \frac{1}{D_C^2}$  can move in both directions following an increase in  $\lambda$ . As a consequence,  $\sigma^*$  can either increase or decrease in order to move  $\mu_C - \mu_F$  back to the equilibrium level. In other words, despite more competent politicians being available, it is possible for the increasingly distorted behaviour of incompetent politicians to decrease accountability welfare.

In terms of selection welfare, on the other hand, the average  $\mu$  increases despite  $\mu_C - \mu_F$  remaining constant. This means that either  $\mu_C$  and  $\mu_F$  increase, or that at least  $Pr(\tau = C)$  has to increase. However, we can check using  $Pr(\tau = C) = \frac{\lambda\gamma}{\mu_C}$  that  $\mu_C$  has to increase, because otherwise both  $Pr(\tau = C)$  and  $Pr(\tau = F)$  would increase, which is a contradiction. It follows by Lemma 4 that selection welfare increases.

Suppose now that the game has a truthful equilibrium and that an increase in  $\lambda$  does not make it switch to a partially truthful equilibrium. It is straightforward that accountability welfare increases, since more incumbents are high types. However,  $\mu_C^T - \mu_F^T$  can move both up or down, and despite knowing that both  $\mu_C^T$  and  $\mu_F^T$  increase, selection welfare might increase or decrease.  $\square$

### Proof of Welfare Result 4

*Proof.* Consider a partially truthful equilibrium. When  $\gamma$  increases, the equilibrium value of  $\mu_C - \mu_F$  decreases. As a consequence of Lemma 4, we know that in such a situation selection can only improve if both  $\mu_C$  and  $\mu_F$  increase. However, since it has to be that  $Pr(\tau = C)\mu_C + Pr(\tau = F)\mu_F = \lambda$ , then if both  $\mu_C$  and  $\mu_F$  were to increase,  $Pr(\tau = F)$  would have to increase, too. However, this is not possible, because if  $\gamma$  increases and  $Pr(\tau = F)$  also increases, then  $\mu_F$  has to decrease, since  $Pr(\tau = F) = \frac{\lambda(1-\gamma)}{\mu_F}$ . This means that selection always worsens when  $\gamma$  increases.  $\square$

### Proof of Proposition 2

*Proof.* The welfare expression derived in Lemma 3 reads:

$$W = [\lambda + (1 - \lambda)q] + [\lambda + (1 - \lambda)\tilde{q}] + b \left[ \frac{1}{4} + \frac{\lambda_O^2}{4} - \frac{\lambda\lambda_O}{2} + \frac{\lambda + \lambda_O}{2} + \frac{\mathbb{E}\mu^2}{4} \right]$$

With a single term limit, politicians can do no better than following their signal, since there is no re-election possibility. From the voter's perspective, the utility gain is:

$$(1 - \lambda)A(1 - \sigma^*)(2\rho_2 - 1)$$

At the same time, however, a single term limit corresponds to a commitment to choosing the challenger no matter what the belief about the incumbent is. Therefore, society gets the benefit  $b$  with probability  $\lambda_O$ , plus  $\mathbb{E}v = 0$ . The utility from the selection of the right type of politician is therefore  $\lambda_O b$  rather than  $b \left[ \frac{1}{4} + \frac{\lambda_O^2}{4} - \frac{\lambda\lambda_O}{2} + \frac{\lambda + \lambda_O}{2} + \frac{\mathbb{E}\mu^2}{4} \right]$ . This generates a loss in terms of selection welfare expressed by:

$$b \left( \frac{1}{4} + \frac{\lambda_O^2}{4} - \frac{\lambda\lambda_O}{2} + \frac{\lambda - \lambda_O}{2} + \frac{\mathbb{E}\mu^2}{4} \right)$$

It follows that having a single term limit is beneficial if the following inequality holds:

$$(1 - \lambda)A(1 - \sigma^*)(2\rho_2 - 1) \geq b \left( \frac{1}{4} + \frac{\lambda_O^2}{4} - \frac{\lambda\lambda_O}{2} + \frac{\lambda - \lambda_O}{2} + \frac{\mathbb{E}\mu^2}{4} \right)$$

$\square$

### Proof of Proposition 3

*Proof.* Before proving the proposition, I provide an intermediate result that will be useful in the proof.

**Lemma 5:**  $p_C + p_F$  is increasing in  $\sigma$

*Proof.* By definition,  $p_C = \frac{gPr(\tau=C)}{gPr(\tau=C)+(1-g)Pr(\tau=F)}$  and  $p_F = \frac{gPr(\tau=F)}{gPr(\tau=F)+(1-g)Pr(\tau=C)}$ . Now, remember that  $\mu_C = \frac{\lambda\gamma}{\lambda\gamma+(1-\lambda)(1-A\sigma)} = \frac{\lambda\gamma}{Pr(\tau=C)}$  and similarly,  $\mu_F = \frac{\lambda(1-\gamma)}{Pr(\tau=F)}$ . Substituting these into the expressions for  $p_C$  and  $p_F$  one gets the following expressions:

$$p_C = \frac{g\gamma\mu_F}{g\gamma\mu_F + (1-g)(1-\gamma)\mu_C}$$

and

$$p_F = \frac{g(1-\gamma)\mu_C}{g(1-\gamma)\mu_C + (1-g)\gamma\mu_F}.$$

Let's now denote by  $D_C$  and  $D_F$  the denominators of  $\mu_C$  and  $\mu_F$  respectively and by  $Den_C$  and  $Den_F$  the denominator of  $p_C$  and  $p_F$  respectively. First of all, let's establish that  $Den_F < Den_C$ . This clearly holds since  $g(1-\gamma)\mu_C + (1-g)\gamma\mu_F < g\gamma\mu_F + (1-g)(1-\gamma)\mu_C$  can be rearranged to yield  $D_C > D_F$ , which is always satisfied. Let's now denote as  $h(\sigma)$  the following ratio:

$$h(\sigma) = \frac{1-\gamma}{\gamma} \frac{\mu_C}{\mu_F}.$$

Notice that  $p_C = \frac{1}{1+\frac{1-g}{g}h(\sigma)}$  and  $p_F = \frac{1}{1+\frac{1-g}{g}\frac{1}{h(\sigma)}}$ . Therefore, differentiating  $p_C + p_F$  with respect to  $\sigma$  yields:

$$\frac{1-g}{g}h'(\sigma) \left[ \frac{1}{Den_F^2 h^2(\sigma)} - \frac{1}{Den_C^2} \right] > 0$$

In order to sign this quantity, notice that  $h'(\sigma) > 0$ , since  $\mu_C$  increases and  $\mu_F$  decreases as  $\sigma$  increases; moreover,  $h(\sigma) < 1$ , since  $(1-\gamma)\mu_C < \gamma\mu_F$  holds as a consequence of  $D_F < D_C$ . Remembering that it is also the case that  $Den_F < Den_C$ , we can conclude that the sign is positive. □

I will now proceed with the proof of the proposition. In terms of first period behaviour, nothing changes with respect to the baseline model: following the same arguments as in Theorem 1, all incumbents follow their signal at  $t = 1$ .

Consider now  $t = 2$ . Let's fix the strategy played by incumbents, and in particular let's assume that incompetent incumbents follow a flip-flopping signal with probability  $\sigma$ . If the voter were able to observe track-records perfectly, beliefs would be  $\mu_C(\sigma)$  and  $\mu_F(\sigma)$ . Given the noise, the belief on the incumbent being competent when observing a consistent track-record is a linear combination of the two, i.e.  $\tilde{\mu}_C = p_C\mu_C(\sigma) + (1-p_C)\mu_F(\sigma)$ , where, as before,  $p_C = \frac{gPr(\tau=C)}{gPr(\tau=C)+(1-g)Pr(\tau=F)}$ . Analogous expressions, denoted by  $\tilde{\mu}_F$  and  $p_F$ , hold for the case in which the voter observes a flip-flopping track-record and are derived by simply swapping  $C$  with  $F$ . When deciding whether to follow the signal, the incumbent knows that with probability  $1-g$  the voter will observe a flip-flopping track-record even

if he plays consistently, and viceversa. It follows that the incumbent prefers to follow a flip-flopping signal whenever the following inequality holds:

$$\rho_2 + gr(\tilde{\mu}_F)2\phi + (1-g)r(\tilde{\mu}_C)2\phi \geq 1 - \rho_2 + gr(\tilde{\mu}_C)2\phi + (1-g)r(\tilde{\mu}_F)2\phi$$

Using the definitions of  $\tilde{\mu}_C$  and  $\tilde{\mu}_F$ , this expression can be rearranged to yield the following:

$$\frac{2\rho_2 - 1}{\phi} \geq (2g - 1)(\tilde{\mu}_C - \tilde{\mu}_F)$$

and then further to get:

$$\frac{2\rho_2 - 1}{\phi} \geq (2g - 1)(p_C + p_F - 1)(\mu_C - \mu_F)$$

Notice that if  $g = 1$ , i.e. no noise, then one gets back to the expression from the baseline model, given also that  $p_C(g = 1) = 1 = p_F(g = 1)$ . Let's now consider the case in which a truthful equilibrium is not sustainable when  $g = 1$ , i.e.  $\frac{2\rho_2 - 1}{\phi} < (\mu_C^T - \mu_F^T)$ . When noise is introduced, the right hand side becomes  $(2g - 1)(p_C + p_F - 1)(\mu_C^T - \mu_F^T)$ . Thanks to Lemma 5, it is immediate to check that given  $\mu_C^T$  and  $\mu_F^T$ , this quantity strictly decreases as  $g$  decreases, and it reaches zero as  $g = \frac{1}{2}$ . This means that there is a level of noise  $g^* \in (\frac{1}{2}, 1)$  such that

$$\frac{2\rho_2 - 1}{\phi} = (2g^* - 1)(p_C(g^*) + p_F(g^*) - 1)(\mu_C^T - \mu_F^T).$$

Remember that  $\mu_C(\sigma) - \mu_F(\sigma)$  is strictly increasing in  $\sigma$  and by Lemma 5 it follows that  $p_C + p_F - 1$  is also strictly increasing in  $\sigma$ . As a result, then, decreasing  $g$  always increases the equilibrium level of  $\sigma$  in order to keep the equilibrium condition satisfied. This means that as long as  $\sigma < 1$ , decreasing  $g$  will alleviate the accountability distortion caused by incumbents avoiding flip-flops. Consider now selection welfare: as  $g$  decreases, as long as  $\sigma < 1$ , the equilibrium level of  $\tilde{\mu}_C - \tilde{\mu}_F$  increases as  $g$  decreases, since  $\frac{2\rho_2 - 1}{\phi} = (2g - 1)(\tilde{\mu}_C - \tilde{\mu}_F)$  has to hold. Concerning  $\tilde{\mu}_C$ , using the fact that  $\tilde{\mu}_C = p_C\mu_C + (1 - p_C)\mu_F$  and using the definitions of  $p_C$  and  $p_F$  yields the following expression:

$$\tilde{\mu}_C = \frac{[g\gamma + (1-g)(1-\gamma)]\mu_C\mu_F}{g\gamma\mu_F + (1-g)(1-\gamma)\mu_C} = \frac{[g\gamma + (1-g)(1-\gamma)]}{\frac{g\gamma}{\mu_C} + \frac{(1-g)(1-\gamma)}{\mu_F}}$$

Using the definition of  $\mu_C$  and  $\mu_F$  and differentiating the denominator of this expression with respect to  $\sigma$  yields  $\frac{\lambda}{1-\lambda}A(1-2g) < 0$  since  $g > \frac{1}{2}$ . It follows that  $\tilde{\mu}_C$  increases as  $\sigma$  increases. Hence, we can use Lemma 4 and conclude that also selection welfare

increases as  $g$  decreases and the equilibrium features  $\sigma < 1$ . As a result, decreasing  $g$  improves overall welfare as long as it crowds out the lies of politicians. Decreasing  $g$  below  $g^*$ , however,  $\mu_C - \mu_F$  cannot increase further and hence incompetent incumbents start having a strict preference towards following their signal when it prescribes a flip-flop. Therefore, decreasing  $g$  further will hurt learning and have no benefit on accountability. It follows that  $g^*$  is the optimal level of news informativeness.  $\square$

#### Proof of Proposition 4

*Proof.* I will first show that for each  $\sigma$ , the reputation spread in the benchmark model,  $\mu_C - \mu_F$ , is larger than its analogue with delegation, denoted by  $\mu_C^D - \mu_F^D$ . In other words:

$$\mu_C - \mu_F > \mu_C^D - \mu_F^D$$

In order to show this, write  $\mu_C - \mu_F$  as follows:

$$\frac{\lambda\gamma}{\lambda\gamma + (1-\lambda)(1-A\sigma)} - \frac{\lambda(1-\gamma)}{\lambda(1-\gamma) + (1-\lambda)A\sigma},$$

where, as usual,  $A = (1-\gamma)(q^2 + (1-q)^2) + 2\gamma q(1-q)$ . For the case of delegation,  $\mu_C^D - \mu_F^D$  can be written as:

$$\frac{\lambda\pi}{\lambda\pi + (1-\lambda)(1-A_D\sigma)} - \frac{\lambda(1-\pi)}{\lambda(1-\pi) + (1-\lambda)A_D\sigma},$$

where  $\pi = \gamma q + (1-\gamma)(1-q)$  and  $A_D = (1-\pi)q + \pi(1-q)$ . Let's start by comparing  $\mu_F$  and  $\mu_F^D$ . The latter is larger whenever  $\frac{A_D}{1-\pi}\sigma < \frac{A}{1-\gamma}\sigma$ . First of all, substituting immediately shows that  $A_D = A$ . This is due to the fact the first action is delegated to an agent that is identical to the incompetent incumbent. Since  $\pi < \gamma$ , we have that  $\mu_F^D > \mu_F$ .<sup>22</sup> Therefore,  $\mu_F^D > \mu_F$  for each value of  $\sigma$ . The fact that  $\mu_C > \mu_C^D$  holds can be shown analogously, using the fact that  $\frac{1}{\pi}(1-A_D\sigma) > \frac{1}{\gamma}(1-A\sigma)$ . Therefore, fixing a level of  $\sigma$ , the reputation spread between consistency and flip-flopping is always larger when the first action is not delegated. Since  $\mu_C - \mu_F$  and  $\mu_C^D - \mu_F^D$  are monotonically increasing in  $\sigma$ , this also means that, denoting by  $\sigma_D^*$  the  $\sigma^*$  the equilibrium levels of  $\sigma$  in the game with and without delegation respectively, the following holds:

$$\mu_C - \mu_F = \mu_C^D - \mu_F^D \Rightarrow \sigma_D^* > \sigma^*.$$

---

<sup>22</sup>If the action was delegated to a more competent agent, the number of flip-flops of the incompetent incumbent would decrease, improving even further the reputation associated to a flip-flop.

In other words, the same reputation spread implies more distortion in the game without delegation. Let's now denote by  $\mu_C^{D,T}$  and  $\mu_F^{D,T}$  the reputations obtained in the delegation game under  $\sigma_D = 1$  (i.e. truthful play): if  $\phi > \frac{2\rho_2-1}{\mu_C^{D,T}-\mu_F^{D,T}} \equiv \bar{\phi}_D$ , then a truthful equilibrium is not sustainable in the delegation game. Since  $\mu_C^T - \mu_F^T > \mu_C^{D,T} - \mu_F^{D,T}$ , a truthful equilibrium is also not sustainable in the game without delegation. It follows that in both games, the equilibrium value of the reputation spread is  $\frac{2\rho_2-1}{\phi}$ , since the information available to the incompetent incumbent when facing a trade-off is the same in both games. Therefore,  $\sigma^* < \sigma_D^*$ , meaning that accountability improves at  $t = 2$ . In order to show that selection welfare improves, too, consider, first notice that since  $\sigma_D^* > \sigma^*$  and  $A_D = A$ , then  $\mu_C > \mu_C^D$ . As a consequence of  $\mu_C - \mu_F = \mu_C^D - \mu_F^D$  and  $\mu_C^D > \mu_C$ , Lemma 4 allows us to conclude that the probability of electing a competent incumbent is larger in the delegation game, too.

To conclude, delegation improves accountability at  $t = 2$  and the selection of competent politicians through elections. However, delegation has a cost since the first period action is taken by an agent who is always incompetent. Therefore, a sufficient condition for delegation to be worthwhile in terms of total welfare is for  $\alpha$  to be sufficiently low.  $\square$

### Proof of Proposition 5

*Proof.* In some parts, this proof is analogous to the one of Proposition 4, which I follow as much as possible for consistency. I will first show that for each  $\sigma$ , the reputation spread in the benchmark model,  $\mu_C - \mu_F$ , is larger than its analogue with delegation, denoted by  $\mu_C^d - \mu_F^d$ . In other words:

$$\mu_C - \mu_F > \mu_C^d - \mu_F^d$$

In order to show this, write  $\mu_C - \mu_F$  as follows:

$$\frac{\lambda\gamma}{\lambda\gamma + (1-\lambda)(1-A\sigma)} - \frac{\lambda(1-\gamma)}{\lambda(1-\gamma) + (1-\lambda)A\sigma},$$

where, as usual,  $A = (1-\gamma)(q^2 + (1-q)^2) + 2\gamma q(1-q)$ . For the case of delegation,  $\mu_C^d - \mu_F^d$  can be written as:

$$\frac{\lambda\gamma}{\lambda\gamma + (1-\lambda)(1-A_d\sigma)} - \frac{\lambda(1-\gamma)}{\lambda(1-\gamma) + (1-\lambda)A_d\sigma},$$

where  $A_d = (1-\gamma)q + \gamma(1-q)$ . Let's start by comparing  $\mu_F$  and  $\mu_F^d$ . Recalling that  $A = \gamma 2q(1-q) + (1-\gamma)(q^2 + (1-q)^2)$ , it is straightforward to show that  $A_d < A$ .

Therefore,  $\mu_F^d > \mu_F$  given any value of  $\sigma$ . The fact that  $\mu_C > \mu_C^d$  holds can be shown analogously. Therefore, fixing a level of  $\sigma$ , the reputation spread between consistency and flip-flopping is always larger when the first action is not delegated. Since  $\mu_C - \mu_F$  and  $\mu_C^d - \mu_F^d$  are monotonically increasing in  $\sigma$ , this also means that, denoting by  $\sigma_d^*$  the  $\sigma^*$  the equilibrium levels of  $\sigma$  in the game with and without delegation respectively, the following holds:

$$\mu_C - \mu_F = \mu_C^d - \mu_F^d \Rightarrow \sigma_d^* > \sigma^*.$$

In other words, the same reputation spread implies more distortion in the game without delegation. Compared to the baseline game, however, the posterior of the incompetent incumbent when the signal at  $t = 2$  suggests to flip-flop is smaller. As a matter of fact,  $\rho_2^d = \frac{(1-\gamma)q}{(1-\gamma)q + \gamma(1-q)}$  is smaller than  $\rho_2$  since the prior on the opposite state of the world is now  $\gamma$  instead of  $\gamma q + (1-\gamma)(1-q)$ .

Let's now denote by  $\mu_C^{d,T}$  and  $\mu_F^{d,T}$  the reputations obtained in the delegation game under  $\sigma_d = 1$  (i.e. truthful play): if  $\phi > \frac{2\rho_2^d - 1}{\mu_C^{d,T} - \mu_F^{d,T}} \equiv \bar{\phi}_d$ , then a truthful equilibrium is not sustainable in the delegation game. Depending in whether the prior effect on  $\rho_2^d$  or the reputation effect dominates,  $\bar{\phi}_d$  can be either larger or smaller than  $\bar{\phi}$  (the threshold for the truthful equilibrium in the baseline game). However, suppose that  $\phi > \max\{\bar{\phi}, \bar{\phi}_d\}$ . In other words, we are in a partially truthful equilibrium no matter whether delegation occurs. In equilibrium,  $\mu_C^d - \mu_F^d < \mu_C - \mu_F$ , so whether  $\sigma_d^* > \sigma^*$  or the other way around is ambiguous. However, a conclusion that can be drawn from the analysis is that  $\mu_C^d - \mu_F^d < \mu_C - \mu_F \Leftrightarrow A_d \sigma_d^* < A \sigma^*$ .

In terms of selection welfare, I can use Lemma 4 to conclude that since  $\mu_C - \mu_F > \mu_C^d - \mu_F^d$  and  $\mu_C > \mu_C^d$ , selection welfare is higher in the baseline game.

In conclusion, the effect of delegation to a competent agent on accountability at time  $t = 2$  is ambiguous, whereas the effect on selection welfare is negative. Since the effect on  $t = 1$  policy welfare is positive, for large enough  $\alpha$  delegation to a competent agent is preferred to no delegation. □

## Proof of Proposition 6

*Proof.* Let's start the analysis by writing the modified reputations. Given that there are now two signals (the track record and the media signal) we now have four different reputations. I denote by  $\mu_{C,E}$  and  $\mu_{F,O}$  the reputation from consistent play given that the media endorses (*E*) or opposes (*O*) the politician's decision. Given this notation, the reputation expressions can be written in the following way, denoting by  $p_2$  the probability

that  $\omega_2 = s_1$ , that is  $p_2 = \gamma q + (1 - \gamma)(1 - q)$ :

$$\begin{aligned}\mu_{C,E} &= \frac{\lambda \gamma q_M}{\lambda \gamma q_M + (1 - \lambda)[(p_2 q q_M + (1 - p_2)(1 - q)(1 - q_M)) + ((1 - p_2)q(1 - q_M) + p_2(1 - q)q_M)(1 - \sigma)]} \\ \mu_{F,E} &= \frac{\lambda(1 - \gamma)q_M}{\lambda(1 - \gamma)q_M + (1 - \lambda)((1 - p_2)q(1 - q_M) + p_2(1 - q)q_M)\sigma} \\ \mu_{C,O} &= \frac{\lambda \gamma(1 - q_M)}{\lambda \gamma(1 - q_M) + (1 - \lambda)[(p_2 q(1 - q_M) + (1 - p_2)(1 - q)q_M) + ((1 - p_2)q q_M + p_2(1 - q)(1 - q_M))(1 - \sigma)]} \\ \mu_{F,O} &= \frac{\lambda(1 - \gamma)(1 - q_M)}{\lambda(1 - \gamma)(1 - q_M) + (1 - \lambda)((1 - p_2)q(1 - q_M) + p_2(1 - q)q_M)\sigma}\end{aligned}$$

Notice that if  $q_M = \frac{1}{2}$ , then we get back to the expressions used in the baseline model. Moreover, compared to the reputations from the baseline model,  $\mu_{C,E} > \mu_C > \mu_{C,O}$  and the analogous inequality holds for the flip-flopping reputations. Now, denote by  $S = Pr(s_M = j | s_2 = j, s_1 = \neg j) = \rho_2 q_M + (1 - \rho_2)(1 - q_M)$  for  $j \in \{0, 1\}$  ( $\rho_2 = Pr(\omega_2 = s_2 | s_2 \neq s_1, \theta = L)$  is defined as usual), the probability that, given the incumbent's signal is flip-flopping, the media signal also endorses a flip-flop. It turns out that when they receive a flip-flopping signal, incompetent incumbents follow their signal if the following inequality holds:

$$\frac{2\rho_2 - 1}{\phi} \geq [S\mu_{C,O}(\sigma = 1) + (1 - S)\mu_{C,E}(\sigma = 1)] - [S\mu_{F,E}(\sigma = 1) + (1 - S)\mu_{F,O}(\sigma = 1)]$$

The expression can be rearranged to yield:

$$\frac{2\rho_2 - 1}{\phi} \geq S[\mu_{C,O}(\sigma = 1) - \mu_{F,E}(\sigma = 1)] + (1 - S)[\mu_{C,E}(\sigma = 1) - \mu_{F,O}(\sigma = 1)]$$

These inequalities deliver, as usual, an upper bound on  $\phi$  such that the equilibrium is truthful. When  $\phi$  exceeds this upper bound, the equilibrium value of  $\sigma$  becomes less than 1. First of all, it can be verified that the right-hand side of the above equation is increasing in  $\sigma$ : similarly to the baseline model, notice from the expressions in the previous page that  $\mu_{C,E}$  is increasing in  $\sigma$ , and so is  $\mu_{C,O}$ , whereas  $\mu_{F,E}$  and  $\mu_{F,O}$  are decreasing in  $\sigma$ . This ensures the existence and uniqueness of a partially truthful equilibrium. Having done that, the question is whether the right-hand side of the equation can be increasing in  $q_M$ . If that is the case, then  $\sigma^*$  needs to decrease for equilibrium to be restored. It can be numerically verified that there exist parameter values such that  $\sigma^*$  is decreasing in  $q_M$ . The Matlab code is available upon request. Finally, notice that in the first period incumbents have the incentive to follow their signal, just like in the baseline game: not following the signal in the first period leads to a gain conditional on the second signal being different from the first, whereas it leads to a loss when the second signal matches the first. In particular, conditional on  $s_2 = s_1$  and  $m = s_2$ , which happens with probability

$Pr(s_2 = m = s_1 | s_1) \equiv \pi_{C,E}$  the gain from following the first signal is  $2\phi(\mu_{C,E} - \mu_{F,E})$ ; conditional on  $s_2 = s_1 \neq m$ , which happens with probability  $Pr(s_2 = s_1 \neq m | s_1) \equiv \pi_{C,O}$  the gain is  $2\phi(\mu_{C,O} - \mu_{F,O})$ ; conditional on  $s_2 \neq s_1$  and  $m = s_2$ , which occurs with probability  $Pr(s_2 = m \neq s_1 | s_1) \equiv \pi_{F,E}$ , the loss is  $2\phi(\mu_{C,E} - \mu_{F,E})$  and finally conditional on  $s_2 \neq m = s_1$ , which happens with probability  $Pr(s_2 \neq s_1 = m | s_1) \equiv \pi_{F,O}$ , the loss is  $2\phi(\mu_{C,O} - \mu_{F,O})$ . Putting everything together, the gain from following the signal is:

$$[\pi_{C,E} - \pi_{F,E}](\mu_{C,E} - \mu_{F,E}) + [\pi_{C,O} - \pi_{F,O}](\mu_{C,O} - \mu_{F,O}) > 0$$

Notice that

$$\pi_{C,E} - \pi_{F,E} = qq_M[2(\gamma q + (1 - \gamma)(1 - q)) - 1] + (1 - q)(1 - q_M)[1 - 2(\gamma q + (1 - \gamma)(1 - q))]$$

which is equal to:

$$[2(\gamma q + (1 - \gamma)(1 - q)) - 1](qq_M - (1 - q)(1 - q_M)) > 0$$

and similarly:

$$\pi_{C,O} - \pi_{F,O} = [2(\gamma q + (1 - \gamma)(1 - q)) - 1][(q(1 - q_M) - (1 - q)q_M)] > 0$$

The key to this result is, similarly to the baseline model, that  $2(\gamma q + (1 - \gamma)(1 - q)) - 1 > 0$ . □

## Proof of Proposition 7

*Proof.* Let's start by analyzing the undistorted equilibrium (i.e. the analogue of the truthful equilibrium in this setup). In the second period, there are 4 possible combinations of politician's private signal and public poll: in particular, what is crucial is whether  $s_2$  matches the poll  $z$  or differs from it. In the former case, the optimal choice is always to play  $a_2 = s_2 = z$  independently of  $s_1$  and the type. If however the poll and the private signal are conflicting, the first signal becomes pivotal for the decision: for the low type, the private signal and the poll offset each other and, given the persistence of the state, the first signal dictates the optimal action. For the high type, the private signal is stronger, but not enough to drive the decision when the private signal suggests to flip-flop and the poll suggests not to, due to the assumption that  $h < \bar{h}$ . The value of  $\bar{h}$  is therefore such

that for  $h < \bar{h}$ :

$$p > \frac{h[h(1-\gamma) + (1-h)\gamma]}{h[h(1-\gamma) + (1-h)\gamma] + (1-h)[h\gamma + (1-h)(1-\gamma)]},$$

with equality at  $h = \bar{h}$ . As a result, when the private signal and the poll are conflicting, the optimal decision is to stick to the first decision. In this context, therefore, flip-flopping only occurs to match the poll, whereas the poll is sometimes not matched when the politician plays consistently.

Let's now consider reputations. The voter now has an additional signal  $z$ , therefore the reputation depends not only on flip-flopping versus consistency, but also on whether the action matches the poll or not: for  $\tau \in \{C, F\}$ ,  $\mu_{\tau,K}$  denotes the reputation after contradicting the poll, whereas  $\mu_{\tau,M}$  that after matching it. Moreover, notice that since no politician would flip-flop to contradict the poll,  $\mu_{F,K}$  is not defined by Bayes rule but is derived out-of-equilibrium. I will assume that out-of-equilibrium beliefs are pessimistic enough to discourage a flip-flopping track-record when it doesn't match the poll: for example,  $\mu_{F,K} = 0$ .

Suppose now that incumbents play using truthful strategies. If the poll signal matches the first action, the incumbent always plays  $a_2 = z$  and reputation can be written as:

$$\mu_{C,M} = \frac{\lambda C_H}{\lambda C_H + (1-\lambda)C_L}$$

where the expression for  $C_\theta$  denotes, for a player of type  $\theta$ , that is a player with signal accuracy  $q(\theta)$ :

$$C_\theta = \gamma[q(\theta)^2 p + (1-q(\theta))^2(1-p)] + (1-\gamma)[q(\theta)(1-q(\theta))]$$

If on the other hand the signal does not match the first action, maximizing the probability of matching the state requires the incumbent to follow his private signal, as a result of which reputations can be written as:

$$\mu_{F,M} = \frac{\lambda F_H}{\lambda F_H + (1-\lambda)F_L}$$

in the case of flip-flopping and matching the poll with the second action, and

$$\mu_{C,K} = \frac{\lambda K_H}{\lambda K_H + (1-\lambda)K_L}$$

in the case of consistent play and not matching the poll with the second action, where

for an agent of type  $\theta$  with signal accuracy  $q(\theta)$ ,  $F_\theta$  and  $K_\theta$  are defined as:

$$F_\theta = \gamma[q(\theta)(1 - q(\theta))] + (1 - \gamma)[q(\theta)^2p + (1 - q(\theta))^2(1 - p)]$$

and

$$K_\theta = \gamma[q(\theta)^2(1 - p) + (1 - q(\theta))^2p] + (1 - \gamma)[q(\theta)(1 - q(\theta))].$$

Let's focus on the case in which the poll does not match the first action. If the incumbent flip-flops, that signals incompetence but at the same time having a signal that matches the poll is a sign of competence. If  $p$  is large enough, therefore, the second effect outweighs the first, and flip-flopping to match the poll gives a higher reputation than not flip-flopping but not matching the poll. In particular, there exist values of  $p$  such that,  $\mu_{F,M} > \mu_{C,K}$ , which holds if and only if:

$$\frac{F_L}{F_H} < \frac{K_L}{K_H}$$

Given the assumption that  $p = q(L)$  the former inequality can be rewritten as:

$$\frac{\gamma p(1 - p) + (1 - \gamma)(p^3 + (1 - p)^3)}{\gamma(p^2(1 - p) + (1 - p)^2p) + (1 - \gamma)p(1 - p)} < \frac{\gamma h(1 - h) + (1 - \gamma)(h^2p + (1 - h)^2(1 - p))}{\gamma(h^2(1 - p) + (1 - h)^2p) + (1 - \gamma)h(1 - h)}$$

and it can be checked that there exists  $p \in [\underline{p}, \bar{p}]$  for which it is verified, with  $\underline{p} < \gamma$  and  $\bar{p} = h$ . In other words, for  $p$  sufficiently large, the reputation from flip-flopping and matching the poll becomes greater than the reputation from being consistent but contradicting the poll.

Therefore, when  $s_2 = z = s_1 = a_1$  or when  $s_2 \neq z = s_1 = a_1$ , matching the poll is optimal both in terms of matching the state and in terms of reputation, since  $\mu_{C,M} > \mu_{F,K}$ . When on the other hand  $z \neq s_1 = a_1$ , if  $s_2 = z$  the choice of  $a_2 = z$  is straightforward since it maximizes the probability of matching the state and  $\mu_{F,M} > \mu_{C,K}$ ; however, if  $z \neq s_1$  and  $s_2 \neq z$ , then following the private signal leads to the optimal state matching decision but leads to a worse reputation, whereas posturing to match the poll is costly in terms of policy performance but gives a reputation of  $\mu_{F,M} > \mu_{C,K}$ . As a result, the politician has a trade-off and the size of the electoral concerns  $\phi$  determines whether an undistorted equilibrium is feasible or not. In particular, the maximum value of  $\phi$  such that a truthful equilibrium is sustainable is:

$$\bar{\phi}_z = \frac{2\rho_z(L) - 1}{\mu_{F,M} - \mu_{C,K}}$$

In a partially truthful equilibrium, in which the reputations  $\mu_{F,M}$  and  $\mu_{C,K}$  become  $\mu_{F,M}^*$  and  $\mu_{C,K}^*$ , the probability that the incompetent incumbent follows his signal when  $z \neq a_1 = s_2$ , denoted by  $\sigma_z$ , takes the value  $\sigma_z^*$  which solves the following equation, in a very similar way to what happened in the benchmark model:

$$\frac{\lambda F_H}{\lambda F_H + (1 - \lambda)(F_L + K_L(1 - \sigma_z^*))} - \frac{\lambda K_H}{\lambda K_H + (1 - \lambda)K_L\sigma_z^*} = \frac{2\rho_z(L) - 1}{\phi},$$

where  $\rho_z = \frac{\bar{\rho}_2(1-p)}{\bar{\rho}_2(1-p)+(1-\bar{\rho}_2)p}$  and  $\bar{\rho}_2 = \frac{[\gamma q+(1-\gamma)(1-q)]q}{[\gamma q+(1-\gamma)(1-q)]q+[\gamma(1-q)+(1-\gamma)q](1-q)}$ .

Since  $\rho_z(\theta = L) < \rho_z(\theta = H)$ , in a partially truthful equilibrium incompetent politicians mix between following their signal and contradicting the poll and matching the poll, whereas competent politicians always follow their signal.

Finally, notice that in the first period it is still optimal for the politician to follow the signal. The reason is the following: suppose without loss of generality that the initial signal was  $s_1 = 0$ . Consider first the case of  $z = 0$  and the decision of the incompetent politician. Suppose the politician follows the first signal, i.e.  $a_1 = 0$ : the optimal action for the politician in the second period is  $a_2 = 0$  and the reputation  $\mu_{C,M}$ . Suppose instead that the politician were to deviate and not follow the first signal: the optimal action in the second period would still be  $a_2 = 0$ , resulting in reputation  $\mu_{F,M}$ , independently of the private signal. Since  $\mu_{C,M} > \mu_{F,M}^*$ , deviating in the first period results in a loss of  $2\phi(\mu_{C,M} - \mu_{F,M}^*)$  conditional on the poll matching the first action. Let's now do the same comparison for  $z = 1$  (we are still considering  $s_1 = 0$ ). In this case, if the politician played  $a_1 = 0$ , then we have the following: conditional on  $s_2 = 1$ , the optimal action is to match the poll (and signal) and receive  $\mu_{F,M}^*$ . If  $s_2 = 0$ , the politician is indifferent between following the signal and not matching the poll and not following and matching the poll. For our comparison, consider the utility of not following the signal and getting  $\mu_{F,M}^*$ . Suppose now that the politician did not follow the signal at  $t = 1$ : after  $s_2 = 1$ , the optimal action is to follow signal and poll and play  $a_2 = 1$ , resulting in a reputation of  $\mu_{C,M}$ . When  $s_2 = 0$ , instead, the incompetent politician now strictly prefers action  $a_2 = 0$  and reputation  $\mu_{C,M}$  rather than  $\mu_{F,K}$  and action  $a_2 = 1$ . In other words, in this case the gain of not following the signal in the first period is  $2\phi(\mu_{C,M} - \mu_{F,M}^*)$ , i.e. the same as the gain from following the signal in the first period conditional on  $z = 0$ . However, after receiving  $s_1 = 0$ , the probability of  $z = 0$  in the second period is larger than  $1/2$ , given the informativeness of the poll and the persistence of the state. Therefore, from  $t = 1$  perspective it is optimal to follow the signal and play  $a_1 = 0$  after  $s_1 = 0$  (and analogous conclusions can be drawn for the  $s_1 = 1$  case). Consider now the decision of

the competent politician: nothing changes following  $z = 0$ . Following  $z = 1$ , nothing changes if  $s_2 = 1$ , whereas if  $s_2 = 0$ , the competent politician would always follow the signal and get  $\mu_{C,K}^*$ , whereas after not following the first signal, the choice could be either  $a_2 = 0$  and  $\mu_{F,K}$  or  $a_2 = 1$  and  $\mu_{C,M}$ . In the former case, since in the equilibrium we are considering  $\mu_{F,K} = 0$ , the potential gain from not following the signal at  $t = 1$  is even smaller. In the latter case, the gain in terms of reputation would be larger but along with it there would be a loss from not following the signal after  $s_2 = 0$ . In this case, the gain conditional on  $z = 1$  is

$$\phi(\mu_{C,M} - \mu_{C,K}^*) - (2\rho_2(H) - 1)$$

Notice that this can be rewritten in the following way:

$$2\phi(\mu_{C,M} - \mu_{F,M}^*) + 2\phi(\mu_{F,M}^* - \mu_{C,K}^*) - (2\rho_2(H) - 1)$$

and further as

$$2\phi(\mu_{C,M} - \mu_{F,M}^*) + 2(\bar{\rho}_2(L) - \rho_2(H)) > 2\phi(\mu_{C,M} - \mu_{F,M}^*).$$

In other words, in this case the gain from not following the signal conditional on  $z = 1$  is larger than for the incompetent politician, which might induce him to deviate from the first signal unless the following holds:

$$\phi(\mu_{C,M} - \mu_{F,M}^*)(Pr(z = 0|s_1 = 0) - Pr(z = 1|s_1 = 0)) > Pr(z = 1|s_1 = 0)(\bar{\rho}_2(L) - \rho_2(H))$$

Notice that since  $\mu_{C,M}$  is not affected by  $\phi$  whereas  $\mu_{F,M}^*(\sigma_z^*) < \mu_{F,M}$ , there exists a value of  $\phi$  sufficiently high such that the condition is satisfied and which works as sufficient condition for the competent politician to follow the signal at  $t = 1$ . This is the value  $\phi_{zz}$  mentioned in the statement of the proposition. □