

# Disliking to disagree\*

Florian Hoffmann<sup>†</sup>   Kiryl Khalmetski<sup>‡</sup>   Mark T. Le Quement<sup>§</sup>

November 4, 2019

## Abstract

Abundant experimental and field evidence suggests that people tend to dislike open disagreement. We propose a formalization of perceived disagreement and study the implications of perceived disagreement aversion in a disclosure game involving agents with different priors. Across a variety of settings, the ideal conditions for disclosure involve identical prior variances and differing prior means. When equilibrium disclosure is partial, it is biased towards evidence that is congruent with the most confident agent's prior bias. Perceived disagreement aversion leads to assortative matching in prior beliefs that provides a theoretical basis for echo chambers. Finally, equilibria may feature higher average perceived or actual disagreement than a hypothetical full disclosure scenario.

**Keywords:** strategic disclosure, psychological games, disagreement aversion

**JEL classification:** D81, D83, D91

## 1 Introduction

Decentralized information exchange within social networks is an important channel shaping public opinion, which is ever more important in the digital era (Internet, social media).<sup>1</sup> While avoiding some of the distortions that are particularly relevant for centralized information flows, this source of information itself exhibits many forms of bias. In particular,

---

\*We thank Pierpaolo Battigalli, Oana Borcan, Gary Charness, Martin Dufwenberg, Ángel Hernando-Veciana, Martin Kocher, Matias Nunez, Axel Ockenfels, Amrish Patel, Karl Schlag, Rajiv Sethi, Joel Sobel, Peter Norman Sørensen, Robert Sugden, Joel van der Weele and Jörgen Weibull for helpful comments and suggestions. We thank participants at the following workshops and seminars: Southampton research seminar 2019, ESEM 2018, Toulouse IAS 2018, Paris-Cergy IAS 2018, Warwick Dr@w seminar 2018, 2nd workshop on Psychological Game Theory 2017, UEA Theory workshop 2017, Political Economy workshop in Konstanz 2017 and seminar of the DFG Research Unit "Design and Behavior" 2017. Khalmetski gratefully acknowledges financial support of the German Research Foundation (DFG) through the Research Unit "Design and Behavior" (FOR 1371).

<sup>†</sup>Erasmus University Rotterdam. E-mail: hoffmann@ese.eur.nl.

<sup>‡</sup>University of Cologne. E-mail: kiryl.khalmetski@uni-koeln.de.

<sup>§</sup>University of East Anglia. E-mail: m.le-quement@uea.ac.uk.

<sup>1</sup>See Sunstein (2007), p. 52.: *"In contrast to television, many of the emerging technologies are extraordinarily social, increasing people's capacity to form bonds with individuals and groups that would otherwise have been entirely inaccessible. Email, instant messaging, texting and Internet discussion groups provide increasingly remarkable opportunities, not for isolation, but for the creation of new groups and connections."*

people do not talk equally easily about all topics, are not equally willing to disclose all facts or opinions, and are not equally likely to talk to everyone. A 2016 poll by the online employment website *CareerBuilder* finds that 42 percent of respondents avoid talking politics at the office while 44 percent may talk about it but interrupt the conversation if it becomes heated.<sup>2,3</sup> Social-psychologists have developed a wide repertoire of concepts to describe informational biases arising in social networks, e.g. *Taboos*, *Overton windows*, *opinion corridors*, *political correctness*, *conversational minefields*, *echo chambers*, *confirmation bias*, *pluralistic ignorance*, *information avoidance*.

An important role in generating these biases can be attributed to the tendency to avoid conflict in opinions (i.e. to the desire to be perceived as having similar beliefs as the counterparty), or *perceived disagreement aversion*. A large body of experimental and field evidence documents that individuals tend to state opinions that conform to what they believe others think. In the seminal experiments conducted by Asch (1955), subjects wrongly evaluated the length of a line in public after being exposed to other participants' (artificially induced) wrong assessment. Deutsch and Gerard (1955) showed that this effect is weaker if subjects report their judgment privately, so that others' *perceived* disagreement is unaffected. Mutz (2006) reviews a number of studies showing that Americans avoid discussing politics with non like-minded people for fear of creating tensions.<sup>4</sup> Bursztyn et al. (2017) found that subjects were more likely to publicly reveal immigration-critical views two weeks after Donald Trump's victory than two weeks before it (i.e. before it became apparent that such views might be shared by a large fraction of the population). Prentice and Miller (1993) established that a large fraction of students refrained from expressing dissent with campus alcohol practices for fear of stigma, vastly underestimating the share of people sharing their opinion.<sup>5</sup>

---

<sup>2</sup>*Political Talk Heats Up the Workplace, According to New CareerBuilder Survey*, CareerBuilder.com, Press Releases, July 2016.

<sup>3</sup>See also for example the following recommendation from the gentleman's manual "*Hills Manual of Social and Business Forms*" from (1879): "*Do not discuss politics or religion in general company. (...) To discuss those topics is to arouse feeling without any good result.*"

<sup>4</sup>See Mutz (2006), p. 107: "*There is already ample evidence in support of the idea that people avoid politics as a means of maintaining interpersonal harmony. For example, in the mid 1950s, Rosenberg noted in his in-depth interviews that the threat to interpersonal harmony was a significant deterrent to political activity. More recent case studies have provided further support for this thesis. Still others have described in great detail the lengths to which people will go in order to maintain an uncontroversial atmosphere. Likewise, in focus group discussions of political topics, people report being aware of, and wary of, the risks of political discussion for interpersonal relationships. As one focus group participant put it, "It s not worth it to try and have an open discussion if it gets them [other citizens] upset."*

<sup>5</sup>Disagreement aversion has many potential causes (see Golman et al., 2016, for a general review of what the authors term a preference for *belief consonance*). Individuals might experience an intrinsic psychological discomfort from being explicitly confronted with disagreement in views (Festinger, 1957; Domínguez et al., 2016). The aversion may instead be driven by the anticipation of adverse consequences stemming from disagreement. For instance, political practice in north-western Europe (e.g. Netherlands and the so-called Polder model, Scandinavia) puts a strong emphasis on reaching consensus, in particular in negotiations between different labor market organizations.

While there is ample evidence of the relevance of perceived disagreement aversion, to the best of our knowledge it has not been formally modeled. This paper is a first step towards filling this gap. We suggest a formalization of perceived disagreement aversion and analyze its consequences for the incentives to share hard (verifiable) information with other agents which may have different prior beliefs.<sup>6</sup> Given mutually known priors, our baseline specification measures  $i$ 's perceived disagreement between herself and another agent  $j$  simply as the (absolute) distance between  $i$ 's expected value of the state ( $i$ 's first order belief) and  $i$ 's expectation of  $j$ 's expected value of the state ( $i$ 's second order belief), the utility of a perceived disagreement averse agent being strictly decreasing in this distance.

A main source of tension for information sharing originates in the mechanics of Bayesian updating: Though agents update their prior expectation in the same direction whatever the observed signal, the magnitude of belief adjustment depends on the prior belief distribution. It follows that disagreement may well increase following the disclosure of particular signal realizations.<sup>7</sup> As a consequence, a perceived disagreement averse agent has incentives to selectively reveal or hide his private information. The same robust intuition also implies that the benefits of generating costly information, as regards the reduction in (perceived) disagreement, also depends on agents' prior distributions. One of our contributions is to characterize how different specifications of prior heterogeneity may facilitate or hinder information sharing within a given group of disagreement averse agents, as well as affect the choice of conversation partners.

Our baseline model is a simple game of strategic disclosure by a potentially informed sender ( $S$ ) who is averse to disagreement as perceived by an uninformed receiver ( $R$ ). The sender privately observes, with some commonly known probability, an informative signal drawn from a known distribution, and can decide whether to disclose it to the receiver or not.<sup>8</sup>

Information transmission in equilibrium can be characterized based on the differences in means and variances of the heterogeneous prior distributions. This is interesting because these quantities have a natural interpretation. The prior mean represents an agent's prior stance. The prior variance represents his confidence in his prior stance and his willingness to revise his stance as new information becomes available. While we consider various specifications of the state space, the prior distributions and the signal structure, the basic intuition

---

<sup>6</sup>Heterogeneous priors are an integral part of many social situations. Instances range from views on general questions (climate change, immigration, free trade, religion, minority rights) to how to manage a firm or optimize an investment portfolio. A key underlying source is that people have different personal histories (experiences, socialization, education). See Morris (1995) for an early general discussion, and Acemoglu et al. (2016), Banerjee and Somanathan (2001), Gentzkow and Shapiro (2006), and Dixit and Weibull (2007) for modeling applications.

<sup>7</sup>This is most easily seen by comparing the belief adjustment with a degenerate prior (which is zero) to the one with, e.g., a uniform prior (which is strictly positive).

<sup>8</sup>As is standard (see, e.g., Jung and Kwon, 1988), scope for selective disclosure emerges when the probability of being informed is interior, preventing full unraveling.

is most apparent within the simple setup in which the state of the world is either 0 or 1 and  $S$ 's signal is binary and of known precision (call this the binary-binary setting). In this setup, we denote the commonly known prior beliefs that the state is 1 by  $\beta_i$ ,  $i \in \{S, R\}$ .

With heterogeneous priors, this game has (almost always) a unique pure-strategy equilibrium that always features some information transmission, as at least one signal realization is disclosed. Whether full disclosure is feasible however crucially depends on the prior profile. For any signal precision, full disclosure is feasible if  $\beta_S$  is close enough to  $1 - \beta_R$  while for low enough precision, full disclosure is not feasible if  $\beta_S$  is close enough (but not identical) to  $\beta_R$ . The profile of priors that makes it easiest to achieve full disclosure thus features similar prior variances and a potentially large difference in prior means. In such a profile, agents' willingness to revise their stance, and hence the magnitude of their belief adjustments, in the face of confirming and contradictory evidence is similar.<sup>9</sup> In contrast, given a small difference in prior variances, a potentially significant difference in prior means ( $\beta_S \approx 1 - \beta_R$ ) is better than almost none ( $\beta_S \approx \beta_R$ ).<sup>10</sup> The reason is that sufficiently different means ensure that a player with a higher (lower) mean will be relatively less affected by a higher (lower) signal, which in turn implies convergence in posterior beliefs whatever signal is disclosed.

We find that if disclosure is partial, the information selectively revealed by  $S$  is biased towards evidence that is congruent with the more confident player's prior belief. As an illustration, in the binary-binary setting, consider the case in which the most confident player assigns higher probability to state 1. Then, if equilibrium involves partial disclosure, only 1-signals are shown.<sup>11</sup>

We demonstrate in the binary-binary setting that perceived disagreement aversion generates echo chamber-like dynamics in simple matching scenarios. If receivers are randomly matched with senders and priors are publicly observed, a more confident receiver is less likely to encounter contradicting information, this probability tending to zero as his prior variance tends to 0. Allowing for repeated random pairwise encounters, this leads to inertia in learning dynamics. We then show that confirmatory information bias is further strengthened if disagreement averse senders can choose whom to be matched with, while society exhibits a sufficiently high degree of polarization in priors. Senders, rationally anticipating the nature of equilibrium disclosure, only interact with receivers whose prior mean is similar to their own (assortative matching). Our equilibrium characterization then implies that in the exclusively like-minded matches that are formed, only information congruent with the shared bias will be disclosed.

---

<sup>9</sup>This result extends beyond the binary world: In the canonical normal priors - normal signals setting, full disclosure is possible if and only if prior variances are identical, and full disclosure is the only equilibrium outcome if and only if prior means furthermore differ.

<sup>10</sup>In the normal-normal setting, when prior variances differ, the set of disclosed signals has zero measure under identical prior means and instead positive measure when prior means differ.

<sup>11</sup>Within the normal-normal setting, consider a situation in which both prior means and variances differ across players. Then, only signals within a bounded interval are disclosed, and this interval is closer to the prior mean of the more confident player in terms of Hausdorff distance.

Our theory of perceived disagreement aversion hence offers a putative explanation of the following two stylized facts. First, many citizens are exposed disproportionately to information that confirms their worldview (echo chambers). Second, there is very significant positive assortative matching in communicative behavior on the basis of worldviews (worldview homophily), partially as a result of the Internet. These stylized facts are often presented and discussed together.<sup>12</sup> Our tentative explanation of these facts rests on rationality, heterogeneous priors and aversion to (perceived) disagreement.<sup>13</sup>

We extend our analysis in various directions.<sup>14</sup> First, we adapt our baseline measure of perceived disagreement, which presumes commonly known priors, to allow for uncertainty about priors. Within our baseline bilateral disclosure game, (expected) prior heterogeneity in means continues to be conducive to information transmission, echoing the results obtained under known priors. Furthermore, we show that uncertainty about priors may actually be beneficial for information sharing in equilibrium.

In our baseline model of equilibrium disclosure the sender aims at being "politically correct," in the sense of selectively disclosing only signals that reduce *perceived* disagreement. There is an ongoing debate about the value of such self-censorship. In particular, critics of this view of political correctness often point to the benefits of encouraging people to freely speak their minds. Linking to this debate, we evaluate the value of commitment to a full disclosure strategy - or, respectively, the hidden cost of political correctness - and find that equilibrium disclosure by a perceived disagreement averse sender induces higher perceived disagreement in expectation than commitment to a full disclosure strategy if and only if the sender is more confident. Interestingly, equilibrium disclosure might also dominate full

---

<sup>12</sup>See Mutz (2006), p. 9: "*Social network studies have long suggested that likes talks to likes; in other words, people tend to selectively expose themselves to people who do not challenge their view of the world. Network survey after network survey has shown that people talk more to those who are like them than to those who are not, and political agreement is no exception to this general pattern.*" See also Sunstein (2007), p. 145: "*because of self-sorting, people are often reading like-minded points of view, in a way that can breed greater confidence, more uniformity within groups, and more extremism. Note in this regard that shared identities are often salient on the blogosphere, in a way that makes polarization both more likely and more likely to be large.*" See also Sunstein (2007), p. 63: "*The phenomenon of group polarization has conspicuous importance for the communications market, where groups with distinctive identities increasingly engage in within-group discussion. (...) New technologies, emphatically including the Internet, make it easier for people to surround themselves (virtually of course) with the opinions of like-minded but otherwise isolated others, and to insulate themselves from competing views. For this reason alone, they are breeding ground for polarization, and potentially dangerous for both democracy and social peace.*"

<sup>13</sup>An alternative theory explaining these stylized facts is that people talk in order to make the right decisions (say, match the state) and induce others to do the same. This theory predicts that more similar worldviews lead to better information transmission but thereby fails to explain why homogeneity in groups seems to correlate with (confirmation) biased learning. One might assume that individuals underestimate the extent to which peer group members' information correlates with their own, and as a consequence overweight others' information. This theory of so-called correlation neglect has been explored in Levy and Razin (2015) and Glaeser and Sunstein (2009). The theory assumes an element of bounded rationality, which is not the case of ours.

<sup>14</sup>For expository reasons we chose to present these extensions within the binary state - binary signal model.

disclosure with respect to expected *actual* disagreement. To see this, we take the perspective of a third party who cares about minimizing the expected ex post *actual* disagreement between  $S$  and  $R$ . Then, while it is immediate that full disclosure is the optimal commitment strategy whenever the third party shares the prior of either  $S$  or  $R$ , this need no longer be the case when the third party has a different prior, highlighting once again the role of prior heterogeneity.

**Literature review** In its foundations, our paper relates to a literature studying how (public) information relates to disagreement in beliefs. The literature so far focused on actual (instead of perceived) disagreement, trying to explain phenomena such as polarization, which refers to situations where individuals update in opposite directions on the basis of the same information. This may result from different prior beliefs (Dixit and Weibull, 2007; Acemoglu et al., 2007; Sethi and Yildiz, 2012), different privately observed prior signals (Andreoni and Mylovanov, 2012) as well as ambiguity (Baliga et al., 2013).<sup>15</sup> Under certain conditions, disagreement in beliefs may persist in the long run, i.e. asymptotically (Acemoglu et al., 2016; Andreoni and Mylovanov, 2012).<sup>16</sup>

An extensive body of research dating back to Grossman (1981), Milgrom (1981), and Milgrom and Roberts (1986) studies strategic disclosure of verifiable signals by a privately informed sender.<sup>17</sup> These models typically involve a difference in players' preferences over the receiver's action conditional on the state. Newer papers study the case of different prior beliefs, often featuring identical material preferences given the state, such as Banerjee and Somanathan (2001) and Kartik et al. (2015). While these papers, thus, share important elements of our analysis, perceived disagreement does not play a role in shaping equilibrium incentives for disclosure. Relatedly, Che and Kartik (2009) examine the effect of prior belief misalignment on the sender's incentives to privately acquire costly information. Prior misalignment hurts disclosure but increases information acquisition, so that the receiver may ultimately benefit from more misalignment. While potential benefits of prior misalignment also feature prominently in our analysis, this results from a different mechanism.<sup>18</sup> In particular, when information transmission is driven by perceived disagreement aversion, (some) prior misalignment *encourages* disclosure independently of whether information is given exogenously or acquired at some cost.

A strand of the literature on strategic information transmission features an endogenous preference for belief conformity arising from reputational concerns. Morris (2001) (see also

---

<sup>15</sup>Several papers in network economics consider the effect of individual conformity to the beliefs or opinions of others on belief polarization (Dandekar et al., 2013; Buechel et al., 2015; Golub and Jackson, 2012).

<sup>16</sup>Sethi and Yildiz (2016) focus on the fact that observing others' opinion over time, an observer learns both about their subjective prior and about their private information concerning some objective state, thereby triggering non-trivial dynamics in belief updating.

<sup>17</sup>See in particular also Jung and Kwon (1988) for the baseline model of random disclosure as well as Shin (1994a,b). See Sobel (2013) for a general review of the literature on strategic information transmission.

<sup>18</sup>This is apparent by noting that their result rests on strictly positive costs of information acquisition.

Sobel, 1985; Benabou and Laroque, 1992; Ely and Välimäki, 2003) studies a sender-receiver game with an endogenous reputational concern of the sender for being perceived as unbiased, which leads to distorted communication.<sup>19</sup> In Gentzkow and Shapiro (2006), the sender wishes to signal a high quality of her information to the receiver who may remain uninformed about the actual state. This leads her to bias her message towards the receiver’s prior belief.<sup>20</sup> Similarly, in our setup if the sender is less confident, she omits signals contradicting the receiver’s prior. The motivation is however very different: In our model the sender wants to mitigate perceived ex post disagreement with the quality of her information being known. This same objective will as a matter of fact lead the sender to omit signals that confirm the receiver’s prior if the latter is less confident.

Our study also contributes to the growing body of literature on psychological game theory, which posits preferences that directly incorporate beliefs (of arbitrary order) about others’ strategies or beliefs (Geanakoplos et al., 1989; Battigalli and Dufwenberg, 2009). Here, our analysis is related to Ely et al. (2015) who consider the behavior of a principal wishing the beliefs of an agent to follow a specific time path exhibiting suspense or surprises. While this paper as well as our baseline specification focus on pure belief-based preferences, several more applied models allow preferences to depend on the interplay between beliefs and material payoffs, see, for instance, the models of reciprocity (Rabin, 1993; Dufwenberg and Kirchsteiger, 2004) or guilt aversion (Battigalli and Dufwenberg, 2007).

Our paper also relates to a rich theoretical and empirical literature in social psychology on biases in network formation, communication and norm adoption, dating back to the 1950s, 1960s and 1970s (see Newcomb, 1961; Homans, 1961; Asch, 1955; Lazarsfeld and Merton, 1954; Festinger, 1950; Rosenberg, 1954; Huston and Levinger, 1978; Goffman, 1959). Finally, our paper also links to a research agenda in political economy and political theory on deliberative and so-called epistemic democracy (see Estlund, 2009; Landemore and Elster, 2012; Sunstein, 2007; Mutz, 2006; Huckfeldt et al., 2004; Feddersen and Pesendorfer, 1998; Coughlan, 2000; Austen-Smith and Feddersen, 2006) originating in Condorcet’s seminal work on majority voting. The agenda evaluates democratic institutions and practices in terms of their ability to aggregate information (their so-called truth-tracking properties), which ultimately rests on citizens’ incentive or ability to use their private information as well as

---

<sup>19</sup>Lourey (1994) offers a stimulating discussion of self-censorship and political correctness in public discourse stemming from such concerns.

<sup>20</sup>The models in Ottaviani and Sørensen (2006a) and Ottaviani and Sørensen (2006b) embed a similar setting with ex post verifiable reports, resulting in  $S$ ’s reporting conforming to his own prior. Visser and Swank (2007) study deliberative committees whose members want to signal high expertise. This gives them an incentive to pretend to have similar signals (i.e. to agree) and to decide against the prior. Within a similar setup Levy (2007) focuses on the impact of transparency rules on decision making. In a principal-agent setting, Prendergast (1993) examines the agent’s incentive to match the (noisy) information of the principal in his report. Bursztyn et al. (2017) consider a setting where a sender has to communicate his type to a receiver and has an incentive to appear of the same type as the receiver. Bénabou (2012) shows that agents with anticipatory utility may converge to each other’s wrong beliefs due to the dependence of one’s payoffs on the actions of the others.

share it with each other.

The remainder of the paper is organized as follows. Section 2 presents the benchmark model and the main theoretical results. Section 3 considers extensions of the model. Section 4 concludes. All proofs, unless explicitly stated otherwise, are relegated to the online Technical Appendix.

## 2 Main analysis

### 2.1 The disclosure game

There are two agents - the sender ( $S$ , he) and the receiver ( $R$ , she) and a state of Nature  $\omega \in \{0, 1\}$ . Player  $i \in S, R$  assigns prior probability  $\beta_i \in (0, 1)$  to  $\omega = 1$ . Priors are common knowledge.<sup>21</sup>  $S$  holds with probability  $\varphi \in (0, 1)$  a privately observed informative signal which has a value of either 0 or 1. Thus,  $S$  holds information  $\sigma \in \{0, 1, \emptyset\}$ , where  $\emptyset$  stands for no signal. If  $S$  obtains a signal, it is identical to the state with probability  $p \in (\frac{1}{2}, 1]$ , i.e.  $P(\sigma = \omega) = p$  for  $\sigma \neq \emptyset$ . If  $S$  obtains a signal, he can disclose it to  $R$  or not. Denote  $S$ 's disclosed information by  $d$ , where  $d \in \{0, 1, \emptyset\}$  and where  $\emptyset$  stands for no disclosure.  $R$  simply observes  $S$ 's signal if disclosed and subsequently updates beliefs. Let  $\tilde{\beta}_S(\sigma)$  and  $\tilde{\beta}_R(d)$  denote the posterior probability assigned by  $S$  and  $R$ , respectively, to  $\omega = 1$  given obtained information ( $\sigma$  or  $d$ ) and respective priors  $\beta_S$  and  $\beta_R$ .

After the disclosure stage,  $R$  evaluates how much  $S$ 's expected posterior belief is different from her own. In particular,  $R$ 's perceived disagreement is

$$\begin{aligned} \Delta(d, \beta_S, \beta_R) &= |E_R[E_S[\omega|\sigma, \beta_S] | d] - E_R[\omega|d, \beta_R]| \\ &= \left| E_R[\tilde{\beta}_S(\sigma) | d] - \tilde{\beta}_R(d) \right|. \end{aligned} \quad (1)$$

The expression  $\Delta(d, \beta_S, \beta_R)$  captures the extent to which  $R$  thinks that  $S$  beliefs are ex ante biased in a specific direction relative to her own, conditional on disclosure  $d$ .<sup>22</sup>

$S$  is averse to perceived disagreement on the part of  $R$ , i.e. wants to minimize  $R$ 's ex post perception of disagreement. Hence,  $S$ 's utility function for given priors  $\beta_S$  and  $\beta_R$  is given as

$$U_S(\beta_S, \beta_R, d) = -\Delta(d, \beta_S, \beta_R). \quad (2)$$

In other words,  $S$ 's utility is maximized if  $R$  *thinks* that  $S$  holds the same posterior belief as she. Note also that  $S$ 's *actual* posterior belief does not enter  $S$ 's utility function.  $R$ 's

<sup>21</sup>We consider the case of privately known priors in section 2.5.

<sup>22</sup>Note that the values of  $E_R[\tilde{\beta}_S(\sigma) | \emptyset]$  and  $\tilde{\beta}_R(d)$  depend on the assumed disclosure strategy of  $S$ , whereas  $\Delta(1, \beta_S, \beta_R)$  and  $\Delta(0, \beta_S, \beta_R)$  do not. To avoid any ambiguities, we often explicitly write  $\Delta^X(d, \beta_S, \beta_R)$ , where  $X$  is the putative equilibrium disclosure strategy of  $S$ . In the Appendix, we similarly introduce  $E_R^X[\tilde{\beta}_S(\sigma) | \emptyset]$  and  $\tilde{\beta}_R^X(\emptyset)$ .



preferences are left unspecified, this player being entirely passive.

Our equilibrium concept throughout is Perfect Bayesian equilibrium: Players' strategies are sequentially rational given their beliefs and others' equilibrium strategies. Second, beliefs are derived via Bayes' rule whenever possible.

A disclosure strategy of  $S$  specifies a probability of disclosing at each information set of  $S$ , and a disclosure strategy is informative if  $S$  discloses with positive ex ante probability. The three informative and pure disclosure strategies are respectively full disclosure (called FD), disclosure of only 1-signals or only 0-signals (called D1 and D0, respectively). We denote by ND the strategy of never disclosing. An equilibrium featuring disclosure strategy  $X \in \{FD, D1, D0, ND\}$  is called an  $X$ -equilibrium. An equilibrium featuring an informative disclosure strategy is called informative. If  $\beta_i > (<)\frac{1}{2}$ , we say that  $i$ 's prior is *biased towards* state 1 (0). If  $\beta_i > \frac{1}{2}$ , a 1-signal is *congruent with*  $i$ 's prior bias and a 0-signal *contradicts* it (vice versa if  $\beta_i < \frac{1}{2}$ ). If  $\beta_i$  is strictly closer to the boundary than  $\beta_j$ , then  $i$  is said to be more *confident* than  $j$  (or  $i$  holds a more *confident* prior than  $j$ ).

## 2.2 Equilibrium characterization

As our next proposition shows,  $S$ 's optimal disclosure strategy depends on the relation between players' prior beliefs, i.e. on the position of  $\beta_S$  relative to the following thresholds:

$$\beta_S^*(\beta_R, p) = \frac{(1-p)(1-\beta_R)}{1-p+\beta_R(2p-1)},$$

$$\beta_S^{**}(\beta_R, p) = \frac{p(1-\beta_R)}{\beta_R+p(1-2\beta_R)}.$$

The above two functions have the following properties. For  $\beta_R \in (0, 1)$  and  $p \in (\frac{1}{2}, 1]$ , it always holds that  $0 \leq \beta_S^*(\beta_R, p) < \beta_S^{**}(\beta_R, p) \leq 1$ . Also,  $\beta_S^*(\beta_R, p)$  is decreasing in  $p$  while  $\beta_S^{**}(\beta_R, p)$  is increasing in  $p$ . Finally,  $\beta_S^*(\beta_R, \frac{1}{2}) = \beta_S^{**}(\beta_R, \frac{1}{2}) = 1 - \beta_R$  while  $\beta_S^*(\beta_R, 1) = 0$  and  $\beta_S^{**}(\beta_R, 1) = 1$ . As we shall see, for any given  $\beta_R$  these two functions divide the parameter space into three regions, each of which features a unique equilibrium prediction.

**Proposition 1** 1. *If  $\beta_S = \beta_R$ , then any disclosure strategy of  $S$  is an equilibrium disclosure strategy.*

2. *Given  $\beta_S \neq \beta_R$ :*

a) *There exists no ND-equilibrium.*

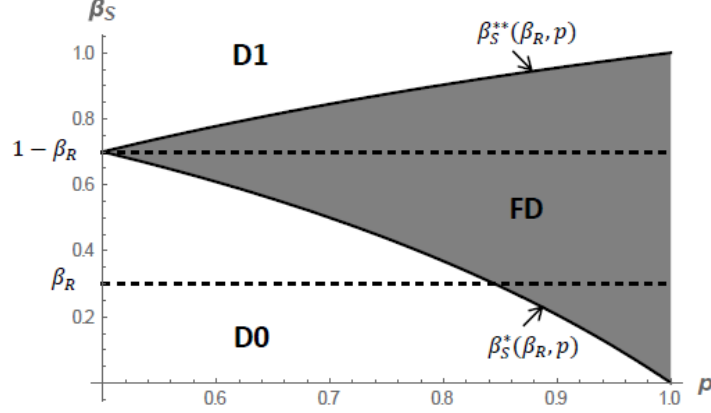
b) *The D0-equilibrium exists if and only if  $\beta_S \in (0, \beta_S^*(\beta_R, p)]$ .*

c) *The FD-equilibrium exists if and only if  $\beta_S \in [\beta_S^*(\beta_R, p), \beta_S^{**}(\beta_R, p)]$ .*

d) *The D1-equilibrium exists if and only if  $\beta_S \in [\beta_S^{**}(\beta_R, p), 1)$ .*

e) *Equilibria in mixed disclosure strategies exist if and only if*

*$\beta_S \in \{\beta_S^*(\beta_R, p), \beta_S^{**}(\beta_R, p)\}$ .*



**Figure 1:** Equilibrium characterization in the baseline model.

Figure 1 provides an illustration of our characterization for  $\beta_R = 0.3$ . The thick curves correspond to  $\beta_S^*(0.3, p)$  and  $\beta_S^{**}(0.3, p)$ . Strictly between the two thick curves (in the gray area), only the FD-equilibrium exists. Instead, strictly above (below) of the upward (downward) sloping thick curve, only the D1 (D0) equilibrium exists. Finally, for  $\beta_S = \beta_R$ , the FD-, D0-, D1- and ND-equilibria exist for any  $p \geq \frac{1}{2}$ . Note that  $\varphi$  does not affect the parameter values for which the different types of equilibrium exist, and it is thus left unspecified for this figure.

Proposition 1 leads to the following corollary.

**Corollary 1** a) If  $\beta_S$  is sufficiently close to  $1 - \beta_R$ , then FD is the unique equilibrium.

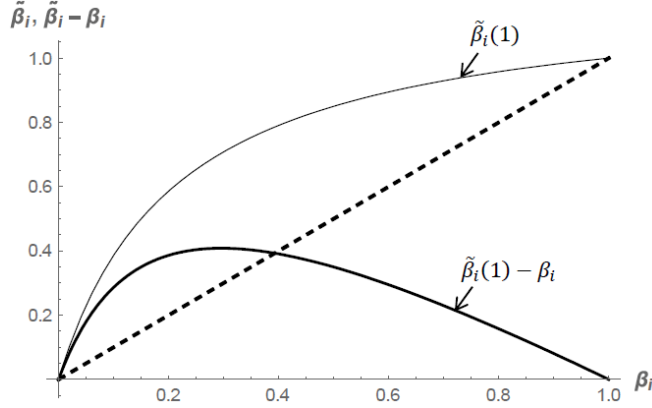
b) Given  $p < \max \left\{ \frac{(1-\beta_R)^2}{(1-\beta_R)^2 + (\beta_R)^2}, \frac{(\beta_R)^2}{(1-\beta_R)^2 + (\beta_R)^2} \right\}$  and  $\beta_R \neq 1/2$ , if  $\beta_S$  is sufficiently close (but not equal to)  $\beta_R$ , then FD is not an equilibrium.

c) For given  $\beta_i$ , the set of  $\beta_j$  for which FD exists is increasing in  $p$ . It is  $(0, 1)$  if  $p = 1$ .

d) If equilibrium features partial disclosure of the D0 or D1 type, the signal that is disclosed is the one that is congruent with the bias of the more confident player.

Summarizing, our characterization exhibits the following key properties:

1. Except under knife-edge conditions, there is a unique equilibrium.
2. Unless  $\beta_S = \beta_R$ , there exists no ND-equilibrium. The reason is that for any  $p$  and  $\beta_S \neq \beta_R$ , the disclosure of at least one type of signal (either 0 or 1) leads to a strict decrease in disagreement with respect to prior disagreement. This follows from the statistical property that (in this binary setup) from  $S$ 's *ex ante* perspective an informative signal always reduces disagreement by moving  $R$ 's belief towards his own prior in expectation.
3. Full disclosure is not always feasible. The intuition comes from contemplating the fact that belief updating has two dimensions: Direction and intensity. In our setup,



**Figure 2:** Intensity of belief updating given a 1-signal as a function of  $\beta_i$ .

players both update in the same direction after any signal (no polarization), but they update with different intensities. In Figure 2, the thin continuous curve shows  $\tilde{\beta}_i(1)$  as a function of  $\beta_i$  for  $p = 0.85$  and the thick curve plots  $\tilde{\beta}_i(1) - \beta_i$ , which is single peaked and concave in  $\beta_i$ . We see that very confident types update very little no matter the signal, while maximum updating arises for a prior moderately biased against the observed signal. Disagreement will increase after a signal if the player assigning the largest prior probability to the state indicated by the signal is also the player who updates the most. A 1-signal, for example, will increase disagreement if  $\beta_i < \beta_j$  and  $\beta_j$  shifts upward more than  $\beta_i$ .

4. Point a) of Corollary 1 states that for any  $p$ , FD is possible if  $\beta_S$  is close enough to  $1 - \beta_R$ . Such  $S$ -prior can have a very different mean than  $R$ 's prior but it has the same variance (i.e. exhibits the same confidence). A technical intuition is as follows. As noted above, for any  $\beta_S, \beta_R, p$  at least one signal (0 or 1) strictly decreases disagreement with respect to the status quo. Next, note that if  $\beta_S = 1 - \beta_R$ , both signals yields the same posterior disagreement, i.e.

$$\tilde{\beta}_S(1) - \tilde{\beta}_R(1) = \tilde{\beta}_S(0) - \tilde{\beta}_R(0).$$

Since at least one signal strictly reduces disagreement, the other must achieve the same. Hence, FD is achievable for any  $p$  for  $\beta_S = 1 - \beta_R$ .<sup>23</sup> A more concrete intuition (see Figure 2) is that when priors are symmetric around  $\frac{1}{2}$ , the prior with the highest (lowest) prior mean moves strictly less than the other prior after a 1-signal (0-signal). Thus the difference in posterior means is always smaller than the difference in prior means.

<sup>23</sup>Note furthermore that updating prior  $\beta_S^*$  with a 1-signal or instead  $\beta_S^{**}$  with a 0-signal yields  $1 - \beta_R$ .

5. Point b) of Corollary 1 states that for  $p$  low enough, FD is impossible if  $\beta_S$  is close enough (but not identical) to  $\beta_R$ . Note that prior variances are very similar if either prior means are approximately symmetric around  $\frac{1}{2}$ , or if these are approximately identical. Yet, full disclosure is significantly less robust in the latter case. The result shows that prior variances are not the only factor affecting disclosure and that prior means also play an important role. For an intuition, let both priors be strongly biased towards 0 and very close to each other. Given that  $\tilde{\beta}_i(1) - \beta_i$  is single peaked in  $\beta_i$  (see Figure 2), we see that after a 1-signal the player with the highest prior updates more intensely. In consequence, the spread between beliefs will increase after this signal.
6. Point c) of Corollary 1 means that a sufficiently precise signal allows for full disclosure. For an intuition, note that in the limit case of  $p = 1$  any signal trivially reduces disagreement to 0. Low signal quality thus triggers two types of costs for  $R$ ; exogenous and endogenous (i.e. strategic). The first is the lower informativeness of  $S$ 's signals and the second is the lower informativeness of  $S$ 's disclosure strategy.
7. For the intuition behind Point d) of Corollary 1, consider the case where the two players have opposite prior biases and let the most confident player be very confident and the other player's prior be close to  $\frac{1}{2}$ . The first player updates very little no matter the signal, so that her posterior is virtually identical to her prior no matter the signal observed. The moderate player instead updates significantly. Now, note that a signal congruent with (in contradiction with) the confident player's bias moves the belief of the moderate player closer to (away from) the confident player's prior.

## 2.3 Matching

### 2.3.1 Non-selective matching

Within a simple random matching setup, Point d) of Corollary 1 naturally implies that the more  $R$ 's prior is biased towards a given state, the less likely she is to be exposed to information contradicting her prior. Assume that a given receiver  $\tilde{R}$ , whose prior  $\beta_{\tilde{R}}$  is publicly observed, is randomly matched with a perceived disagreement averse sender whose publicly observed prior is drawn from the uniform distribution on  $[0, 1]$ . Given  $\beta_S, \beta_{\tilde{R}}$  and  $p$ , the standard disclosure game ensues. We call this game *non-selective matching*.

The following result characterizes the confirmatory bias arising under non-selective matching.

**Remark 1** *For any  $\omega, p$ , under non-selective matching, the ex ante probability that  $\tilde{R}$  observes a 0-signal (1-signal) is decreasing (increasing) in  $\beta_{\tilde{R}}$ .*

For instance, consider the case of a 0-signal. By Proposition 1 the ex ante probability that  $\tilde{R}$  is exposed to a 0-signal is the probability that  $S$ 's signal is 0 and that  $\beta_S$  is such that

the equilibrium is D0 or FD. This equals

$$\Pr[\sigma = 0 | \omega] \Pr[\beta_S < \beta_S^{**}(\beta_{\tilde{R}}, p)] = \Pr[\sigma = 0 | \omega] \beta_S^{**}(\beta_{\tilde{R}}, p),$$

which is strictly decreasing in  $\beta_{\tilde{R}}$ .

Within a dynamic version of the above non-selective matching scenario where  $\tilde{R}$  repeatedly plays the same one-shot disclosure game against short-sighted senders (i.e. who do not update from observing  $\tilde{R}$ 's period- $t$  prior), perceived disagreement aversion on the part of senders thus slows down  $\tilde{R}$ 's learning of the true state (i.e. causes inertia in beliefs) if the state is not congruent with  $\tilde{R}$ 's extreme prior bias. Note that  $\tilde{R}$ 's learning is only slowed down as opposed to entirely impeded, as  $\tilde{R}$  acknowledges that no disclosure by  $S$  does not necessarily imply that he holds no information.

### 2.3.2 Selective matching

In reality, individuals often choose their conversation partners and we now explore this possibility within the context of our model. We find that selective matching further increases the prospect of echo chambers: Individuals select matching partners with similar priors, which in turn induces partial and confirmatory disclosure.

We define the game of *selective matching* as follows. Suppose a large population of senders and receivers, all senders being (perceived) disagreement averse. As before, senders and receivers are randomly matched and observe each others' priors. Yet, in contrast to non-selective matching, a match becomes *active* if and only if the sender accepts it. We focus on the sender's payoffs. If the match does not become active, the sender obtains a payoff of 0, which represents his outside option. If the match becomes active, the standard disclosure game introduced in section 2.1 ensues and the sender's final payoff equals  $W > 0$  minus the ex post perceived disagreement after the disclosure stage. The interpretation of the payoffs is that psychological payoffs only arise once a sender explicitly decides to become involved in conversation. For simplicity, assume that the sender's match acceptance decision is made before his information is realized.

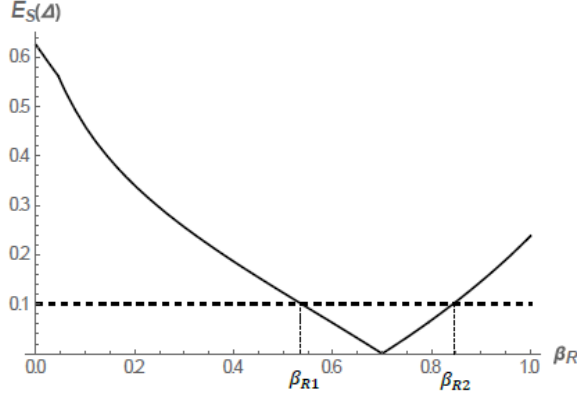
Let  $E_S[\Delta|\beta_S, \beta_R]$  denote the sender's ex ante expectation of the receiver's ex post perceived disagreement given  $\beta_S, \beta_R$  conditional on the match becoming active.<sup>24</sup> It follows that  $S$  will accept a match with  $R$  if and only if

$$W \geq E_S[\Delta|\beta_S, \beta_R]. \tag{3}$$

Furthermore,  $E_S[\Delta|\beta_S, \beta_R]$  satisfies the following description.

---

<sup>24</sup>Note that  $E_S[\Delta|\beta_S, \beta_R]$  is uniquely defined. In particular, if  $X$  and  $X'$  are two equilibrium disclosure rules given  $\beta_S, \beta_R$ , then  $E_S[\Delta^X|\beta_S, \beta_R] = E_S[\Delta^{X'}|\beta_S, \beta_R]$  (as is shown in the proof of Proposition 2). Recall furthermore that by Proposition 1 there is a unique equilibrium disclosure rule except under knife-edge conditions.



**Figure 3:** Perceived disagreement expected by  $S$  in equilibrium as a function of  $\beta_R$ .

**Proposition 2**  $E_S[\Delta|\beta_S, \beta_R]$  is continuous and V-shaped with respect to  $\beta_R$ , and it reaches its minimum of 0 at  $\beta_R = \beta_S$ .

Figure 3 illustrates the above proposition. We assume  $p = 0.9$ ,  $\varphi = 0.6$ ,  $\beta_S = 0.7$ . The thick curve shows  $E_S[\Delta|\beta_S, \beta_R]$  as a function of  $\beta_R$ . For  $W = 0.1$  (represented by the horizontal dotted line), only values of  $\beta_R$  situated between  $\beta_{R1}$  and  $\beta_{R2}$  satisfy (3). This is formalized in the following corollary.

**Corollary 2** Given  $W, p$ , there are thresholds  $\underline{\beta}_R < \beta_S < \bar{\beta}_R$  such that  $S$  accepts a match with  $R$  if and only if  $\beta_R \in [\underline{\beta}_R, \bar{\beta}_R]$ , where  $\lim_{W \rightarrow 0} \{\underline{\beta}_R, \bar{\beta}_R\} = \beta_S$ .

The above corollary states that the only matches that become active are those involving players whose priors are sufficiently similar.

Proposition 1 and this corollary imply that the prospect of confirmatory information bias is strengthened under selective matching in societies that are sufficiently polarized. To see this, consider the following scenario. Priors belong either to  $(0, \underline{\beta})$  or  $(\bar{\beta}, 1)$ , where  $\underline{\beta} < \frac{1}{2} < \bar{\beta} = 1 - \underline{\beta}$  and the conditional distribution of priors on each of the two intervals is uniform. Call this society  $(\underline{\beta}, \bar{\beta})$ . We consider throughout a receiver  $\tilde{R}$  with  $\beta_{\tilde{R}} > \bar{\beta}$  (i.e. biased towards 1) and compare outcomes under respectively non-selective and selective matching.

**Remark 2** Consider a society  $(\underline{\beta}, \bar{\beta})$ . Let  $\bar{\beta} > \frac{p - \sqrt{p(1-p)}}{2p-1}$  and  $\beta_{\tilde{R}} \geq \bar{\beta}$ . Given  $\omega$ , for  $W$  small enough  $\tilde{R}$  observes a 0-signal with probability weakly larger than  $\frac{1}{2} \left( \frac{1 - \beta_{\tilde{R}}}{\underline{\beta}} \right) P[\sigma = 0 | \omega]$  under non-selective matching and instead with probability zero under selective matching.

Under non-selective matching, with a probability bounded below by  $\frac{1}{2} \left( \frac{1 - \beta_{\tilde{R}}}{\underline{\beta}} \right)$ ,  $\tilde{R}$ 's match is such that the implied equilibrium is either FD or D0. To see this, note first that  $\tilde{R}$  is

matched half of the time with a sender satisfying  $\beta_S \in (0, \underline{\beta})$ . Second, conditional on  $\beta_S \in (0, \underline{\beta})$  the probability that  $\beta_S \leq 1 - \beta_{\tilde{R}}$  (yielding FD or D0 by Corollary 1) is  $\frac{1 - \beta_{\tilde{R}}}{\underline{\beta}}$ . As long as  $\beta_{\tilde{R}}$  is not extremely high, there is thus a significant probability, weakly larger than  $P[\sigma = 0 | \omega]_{\frac{1}{2}} \left( \frac{1 - \beta_{\tilde{R}}}{\underline{\beta}} \right)$ , that  $\tilde{R}$  encounters a 0-signal.

Selective matching yields a very different picture. For  $W$  very small, by Corollary 2 a sender  $S$  will accept a match with a receiver  $R$  if and only if  $\beta_R \approx \beta_S$ . As a result, any active match in which  $\tilde{R}$  (recall  $\beta_{\tilde{R}} > \bar{\beta}$ ) participates will involve an  $S$  satisfying  $\beta_S > \bar{\beta}$ . If furthermore  $\bar{\beta} > \frac{1}{2p-1} \left( p - \sqrt{p(1-p)} \right)$ , it holds true that  $\bar{\beta} > \beta_S^{**}(\beta_{\tilde{R}}, p)$  so that by transitivity, any active match involving  $\tilde{R}$  satisfies  $\beta_S > \beta_S^{**}(\beta_{\tilde{R}}, p)$  and thus yields the D1 equilibrium by Proposition 1. Hence,  $\tilde{R}$  never observes a 0-signal no matter the state. Finally, note that  $\frac{1}{2p-1} \left( p - \sqrt{p(1-p)} \right)$  is increasing in  $p$  and not very large as long as  $p$  is not very high (for  $p = \frac{3}{4}$  it equals 0.634), which shows that weak societal polarization suffices to create strong echo chamber dynamics under selective matching.

## 2.4 The hidden cost of political correctness

Can  $S$ 's attempt to minimize perceived disagreement be counter-productive from an ex ante perspective, thereby revealing a hidden cost of political correctness (relative to a hypothetical case of full transparency)? We address this question in two different ways: first, from  $S$ 's own perspective in terms of perceived disagreement, and then from the perspective of a third party (e.g., a social planner) who cares about actual disagreement (i.e. would like to reduce social polarization).

First, from  $S$ 's ex ante perspective, can the expected value of ex post *perceived* disagreement be higher in a (partial disclosure) equilibrium than it would be under (non-equilibrium) full disclosure? In such a case,  $S$  would prefer to commit to full disclosure if he could. This question is answered in the next proposition.

**Proposition 3** 1. *Let parameters be such that D1 is the unique equilibrium. If  $\beta_S > \beta_R$ , then  $S$  ex ante strictly prefers full disclosure over the D1-equilibrium. Vice versa if  $\beta_S < \beta_R$ .*  
 2. *Let parameters be such that D0 is the unique equilibrium. If  $\beta_S < \beta_R$ , then  $S$  ex ante strictly prefers full disclosure over the D0-equilibrium. Vice versa if  $\beta_S > \beta_R$ .*

In a partial disclosure equilibrium,  $S$  would thus ex ante prefer to instead commit to full disclosure if and only if he is the most confident player, which always holds true in D1 (D0) when  $\beta_S > \beta_R$  ( $\beta_S < \beta_R$ ). The intuition is as follows. In a partial disclosure equilibrium, e.g. D0, the omission of 1-signals has two countervailing effects. The upside is that  $S$  benefits from hiding a 1-signal once he holds it. The downside is that when  $S$  holds no signal,  $R$  interprets silence as a possible concealment of a 1-signal, which increases perceived disagreement relative to prior disagreement. The negative effect of equilibrium

concealment overweighs its positive effect if  $S$  is the most confident party. Recall that in this case,  $S$  omits signals contradicting his bias in a partial disclosure equilibrium (see Corollary 1.d). But  $R$  places a higher weight on the state corresponding to the omitted signal that  $S$  does, which leads  $R$  to overweight (in  $S$ 's eyes) the probability that such a signal is held (and omitted) by  $S$ , thereby inflating perceived disagreement after no disclosure. Instead, under full disclosure,  $R$ 's prior does not affect her ex post perception of  $S$ 's posterior (which is always common knowledge).

A second key question is whether from the ex ante perspective of a third party (TP) endowed with a prior  $\hat{\beta}$ , the expected value of ex post *actual* disagreement can be higher in equilibrium than it would be under FD. I.e. would TP prefer a truthful sender or a perceived disagreement averse sender if aiming at minimizing the expected ex post actual disagreement? Note that actual disagreement is different from perceived disagreement. The actual disagreement given that  $S$  holds signal  $\sigma$  and discloses  $d$  is  $|\tilde{\beta}_S(\sigma) - \tilde{\beta}_R(d)|$ , where  $\tilde{\beta}_R(d)$  is pinned down by  $R$ 's beliefs concerning  $S$ 's disclosure rule. In what follows, if  $\beta_i < \hat{\beta} < \beta_j$ , we say that  $S$  and  $R$ 's priors are on different sides of  $\hat{\beta}$ .

**Proposition 4** *Let parameters be such that there exists no FD-equilibrium. In the eyes of a third party with prior  $\hat{\beta}$  the expected actual disagreement:*

1. *is strictly larger in equilibrium than under FD if at least one of the following conditions holds:*

- a)  *$S$ 's and  $R$ 's priors are on different sides of  $\hat{\beta}$ ,*
- b)  *$R$ 's prior is further away from  $\hat{\beta}$  than is  $S$ 's prior.*

2. *is strictly smaller in equilibrium than under FD if the following two conditions hold simultaneously:*

- a)  *$S$ 's and  $R$ 's priors are either both strictly smaller or both strictly larger than  $\hat{\beta}$ ,*
- b)  *$S$ 's prior is further away from  $\hat{\beta}$  than  $R$ 's prior and is sufficiently close to the boundary.*

Part 1 of the proposition finds that equilibrium information omission can indeed be counterproductive while Part 2 identifies conditions under which it is helpful. A general intuition behind our results is that TP expects new information to lead  $S$ 's and  $R$ 's beliefs to converge to her prior. The disclosure strategy of  $S$  affects only the speed of convergence of  $R$ 's beliefs, as  $S$ 's actual posterior beliefs are independent of his disclosure strategy.

In Point 1.a),  $S$ 's and  $R$ 's priors are on different sides of  $\hat{\beta}$ . Here, given that  $S$ 's and  $R$ 's beliefs move closer to  $\hat{\beta}$  in expectation, they must also be moving closer to each other. Hence TP would prefer that both  $S$  and  $R$  learn as fast as possible and would thus prefer FD over partial disclosure. The second case in Point 1 is when  $\beta_S$  and  $\beta_R$  are on the same side of  $\hat{\beta}$ , but  $R$  is further away. An instance of this is the case of  $\hat{\beta} < \beta_S < \beta_R$ . Again TP expects  $S$  and  $R$  to converge to her prior  $\hat{\beta}$ , i.e. to both decrease.  $R$  will move towards  $S$  (since  $R$ 's prior decreases) but  $S$  will simultaneously move away from  $R$  (since  $S$ 's prior also



decreases). In consequence, TP would prefer to speed up  $R$ 's convergence by giving her full information.

Point 2 describes the case when  $\beta_S$  and  $\beta_R$  are on the same side of  $\widehat{\beta}$ , but  $S$  is further away from  $\widehat{\beta}$  and is close to the boundary. An instance of this is the case of  $\widehat{\beta} < \beta_R < \beta_S \approx 1$ . Here, both players' beliefs decrease. But decreasing  $R$ 's belief moves it away from  $S$ 's. So TP would prefer to slow down  $R$ 's learning and thus would choose partial disclosure.

## 2.5 Strangers' talk

Conversations often take place between parties who do not exactly know each others' priors but who might hold some relevant information concerning these priors (for example, by observing each other's accent, dressing style, profession, social networks). We now characterize equilibrium outcomes for a set of stylized scenarios featuring privately known priors.

Technically, since the measure of disagreement (1) is defined for given priors, the expected disagreement perceived by  $R$  under unknown  $\beta_S$ , given disclosure  $d$ , is

$$E_{R,\{\beta_S\}}[\Delta(d, \beta_S, \beta_R)] = E_{R,\{\beta_S\}} \left[ \left| E_{R,\{\sigma\}}[\widetilde{\beta}_S(\sigma) | d] - \widetilde{\beta}_R(d) \right| \right],$$

where  $E_{i,\{z\}}$  stands for an expectation over different possible realized values of the random variable  $z$ , as computed by  $i$ . Note that  $R$  takes the expectations sequentially: first, over all possible signal realizations to compute her second-order belief for given  $\beta_S$ , and only then over all possible realizations of  $\beta_S$  to compute the expected disagreement. In other words, she treats different cases of  $\beta_S$  as separate instances of disagreement. For example, an uninformed receiver with prior equal to 0.5 would not consider a sender with the same prior as disagreeing with her if  $S$  is known to hold either a 0- or a 1-signal (with 50% chance each), yet would treat their disagreement as having a positive value in case if  $S$  is known to hold the prior belief of either 0 or 1. Besides being more tractable, this approach reflects the intuition that disagreement caused by different priors is essentially more severe than the one caused by different private signals. Indeed, while the latter type of disagreement can generally be resolved by information exchange (Aumann, 1976), the disagreement caused by different prior beliefs cannot, since it is driven by a difference in initial worldviews which are beyond discussion.

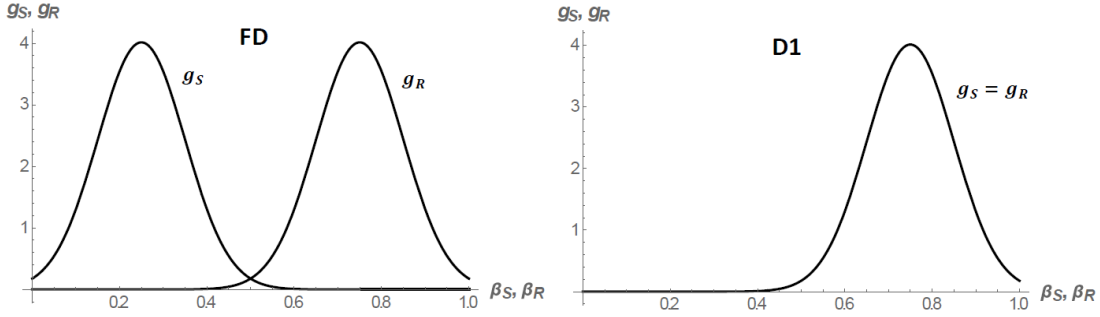
In turn, the expected utility of  $S$  under privately known priors becomes

$$-E_{S,\{\beta_R\}}[E_{R,\{\beta_S\}}[\Delta(d, \beta_S, \beta_R)]] = -E_{S,\{\beta_R\}}E_{R,\{\beta_S\}} \left[ \left| E_{R,\{\sigma\}}[\widetilde{\beta}_S(\sigma) | d] - \widetilde{\beta}_R(d) \right| \right].$$

These preferences give rise to the following equilibrium characterization.

**Proposition 5** *Let priors be privately observed and drawn from publicly known distributions  $G_S$  and  $G_R$ , endowed with respective probability density functions  $g_S$  and  $g_R$ .*

a) *If  $g_S$  and  $g_R$  are both symmetric around 1/2, then there exists an FD-equilibrium.*



**Figure 4:** Equilibrium characterization under uncertainty about priors.

b) If  $g_S$  and  $g_R$  are s.t.  $g_S(x) = g_R(1-x)$  for all  $x \in [0, 1]$  (i.e. they are symmetric w.r.t. each other around  $\frac{1}{2}$ ) and  $\frac{g_S(x)}{g_R(x)}$  is monotone in  $x$ , then there exists an FD-equilibrium.

c) If  $g_S$  and  $g_R$  are identical and sufficiently skewed to the right (left), then there exists a D1 (D0) equilibrium, but no FD and D0 (D1) equilibrium.

d) If  $S$ 's prior is commonly known and sufficiently close to  $1/2$  while  $g_R$  is symmetric around  $1/2$ , then there exists an FD-equilibrium.

Point a) shows that two-sided uncertainty about priors is beneficial to disclosure if none of the two players is a priori biased in one or the other direction. If this condition is satisfied, this provides an argument for not encouraging revelation of information about respective biases (e.g., disclosing one's own prior political stance in a conversation). For an intuition, note that if  $g_S$  and  $g_R$  are both symmetric around  $1/2$ , in a putative FD-equilibrium the payoff from disclosing is the same no matter the signal held by  $S$ . Since it is impossible that *both* signals increase disagreement under FD (i.e. that an informative experiment increases expected disagreement from the ex ante perspective), this implies that both should at worst leave disagreement unchanged. Full disclosure is thus incentive compatible for  $S$ .

Point b) shows that if players are both a priori biased in different directions but in an equivalent (i.e. mirror-image like) fashion, then an FD-equilibrium exists. This contrasts with point c), which states that FD may be infeasible if both priors are drawn from the same biased distribution. The findings of Points b) and c) echo Proposition 1, which highlights the benefit of diversity. For example, for  $p = 0.7$ , if  $\beta_S$  and  $\beta_R$  are both distributed according to the same truncated normal distribution with mean  $3/4$  and standard deviation  $\sigma = 0.3$ , the only (pure strategy) equilibrium is D1. The FD-equilibrium instead exists if the distribution of  $\beta_R$  stays the same while the distribution of  $\beta_S$  is reflected around  $1/2$ , i.e. changed to a truncated normal with mean  $1/4$  and standard deviation  $\sigma = 0.3$ . Figure 4 provides examples of profiles of distributions of prior beliefs, complemented by a description (in bold) of the implied equilibrium disclosure.

Finally, Point d) shows that two-sided uncertainty is not strictly necessary to ensure FD.

The latter is feasible if  $S$ 's prior is known and close to  $\frac{1}{2}$  while  $R$ 's prior is symmetrically distributed around  $\frac{1}{2}$ .

Note that all results in Proposition 5 hold in approximation, i.e. as one density function (uniformly) converges to the other density, or instead to the symmetric reflection around  $\frac{1}{2}$  of the other density.<sup>25</sup>

### 3 Extensions

In what follows, we consider a set of key extensions of our main setup. We first consider a general information structure with continuous signals satisfying the marginal likelihood ratio property (MLRP), while maintaining the assumption of a binary state. In the second subsection, we assume a continuous state space, considering a Normal priors-Normal signals setup. The third subsection finds that perceived disagreement aversion arises endogenously in a variety of simple dynamic games. In our final subsection, we consider disagreement aversion within a game of costly collective acquisition of public signals.

#### 3.1 Binary state and continuous signals

We now show that our equilibrium characterization in the baseline setting carries over qualitatively to the case of an information structure with continuous signals satisfying MLRP. Assume that  $S$ 's signal  $s$  is drawn from an interval  $[\underline{s}, \bar{s}]$ . Given state  $\omega \in \{0, 1\}$ ,  $s$  is distributed according to  $F(s|\omega)$  with continuous and differentiable density  $f(s|\omega)$ . Assume that  $\frac{d}{ds} \frac{f(s|1)}{f(s|0)} > 0$  (MLRP), meaning that a higher signal implies a higher conditional probability of state 1. Upon learning  $s$ , the updated belief of  $i$  is

$$\tilde{\beta}_i(s) = \frac{\beta_i f(s|1)}{\beta_i f(s|1) + (1 - \beta_i) f(s|0)} = \frac{\beta_i}{\beta_i + (1 - \beta_i) \frac{f(s|0)}{f(s|1)}}$$

which is increasing in  $s$ . Assume furthermore that the extreme signals  $\underline{s}$  ( $\bar{s}$ ) are such that  $\lim_{s \rightarrow \underline{s}} \frac{f(s|1)}{f(s|0)} = 0$  and  $\lim_{s \rightarrow \bar{s}} \frac{f(s|1)}{f(s|0)} = \infty$ . Each of these two extreme signal realizations makes the observer (almost) sure that the state is 0 or 1, respectively. Note that there exists a threshold signal  $\tilde{s} \in (\underline{s}, \bar{s})$  such that whatever  $\beta_i \in (0, 1)$ , it holds true that  $\tilde{\beta}_i(s) \gtrless \beta_i$  for  $s \gtrless \tilde{s}$ . Signal  $\tilde{s}$  satisfies  $f(s|0) = f(s|1)$  and we call it the uninformative signal. We say that signal  $s > (<) \tilde{s}$  indicates state 1 (0). We say that signal  $s > (<) \tilde{s}$  is congruent with  $j$ 's prior bias if  $\beta_j > (<) \frac{1}{2}$ . We call the above setup the *binary state-continuous signals environment*. We call *simple disclosure equilibrium* (SD equilibrium) an equilibrium featuring two thresholds  $\underline{s} < s_1 < s_2 < \bar{s}$  such that  $S$  discloses  $s$  if and only if  $s \leq s_1$  or  $s \geq s_2$ . As with the binary signals environment, we call full disclosure (FD) an equilibrium

---

<sup>25</sup>A formal proof is available upon request.

where  $S$  discloses all signals. We obtain the following equilibrium characterization.

**Proposition 6** 1. If  $\beta_S \in \{\beta_R, 1 - \beta_R\}$  then there exists an FD-equilibrium. If  $\beta_S \notin \{\beta_R, 1 - \beta_R\}$ , then the unique equilibrium is an SD equilibrium.

2. In equilibrium, all signals congruent with the bias of the more confident player are disclosed. Signals contradicting the bias of the more confident player are partially disclosed.

The fundamental qualitative features of equilibrium echo those arising under binary signals. Except under knife-edged conditions, the equilibrium is unique. Only signals that are congruent with the prior of the more confident player are fully revealed. Furthermore, if  $\beta_S = 1 - \beta_R$ , a full disclosure equilibrium exists, implying that increasing prior misalignment can be helpful.

We now reexamine the issue of the hidden cost of political correctness already studied for the case of binary signals. Our original results (Propositions 2 and 3) carry over essentially identically to the continuous signals setup.

**Proposition 7** 1. Let parameters be such that in the unique equilibrium, the non-disclosure interval contains signals indicating state 0. If  $\beta_S > \beta_R$ , then  $S$  ex ante strictly prefers full disclosure over the SD equilibrium. Vice versa if  $\beta_S < \beta_R$ .

2. Let parameters be such that in the unique equilibrium, the non-disclosure interval contains signals indicating state 1. If  $\beta_S < \beta_R$ , then  $S$  ex ante strictly prefers full disclosure over the D0-equilibrium. Vice versa if  $\beta_S > \beta_R$ .

**Proposition 8** All the statements in Proposition 4 apply.

### 3.2 Continuous state space and continuous signals

Assume that the state space is  $\mathfrak{R}$ .  $S$  and  $R$ 's commonly known priors are normal and given by respectively  $N(\mu_S, \gamma_S^2)$  and  $N(\mu_R, \gamma_R^2)$ .  $S$  is known to hold a signal with probability  $\varphi \in (0, 1)$ . Given realized state  $\omega$ ,  $S$ 's signal equals  $\omega + \varepsilon$ , where  $\varepsilon \sim N(0, \gamma_\varepsilon)$ , this being commonly known.<sup>26</sup> We denote signal realizations by  $\sigma$ . Note the standard result that

$$E_i[\omega | \sigma] = \frac{\mu_i \frac{1}{\gamma_i^2} + \sigma \frac{1}{\gamma_\varepsilon^2}}{\frac{1}{\gamma_i^2} + \frac{1}{\gamma_\varepsilon^2}}.$$

We provide an equilibrium characterization for the same one-shot disclosure game.  $S$ , if he holds a signal, is free to either disclose it or omit it. We say that a signal  $\sigma$  increases (decreases) disagreement if and only if  $\Delta(\sigma) > (<) |\mu_S - \mu_R|$ , i.e. if the distance between

---

<sup>26</sup>A previous version of this paper contains an analysis of the case where priors are beta distributions and signals are drawn according to a state-dependent binomial distribution. Results (available upon request) echo those obtained in the binary and normal environments, in that they highlight the central role of differences in prior variances.

posterior means conditional on  $\sigma$  is larger (smaller) than that between prior means. We say that player  $i$  is more confident than  $j$  if and only if the variance of  $i$ 's prior belief is smaller, i.e.  $\gamma_i < \gamma_j$ . This is analogous to our previous definition of confidence for the binary state case, which also implies that a more confident player has a smaller variance of prior beliefs.<sup>27</sup>

We obtain the following equilibrium characterization.

**Proposition 9** 1. Let  $\gamma_S \neq \gamma_R$  and  $\mu_S \neq \mu_R$ .

a) Any equilibrium features a finite  $\eta > 0$  such that  $S$  discloses his signal if and only if  $\sigma \in I = [\tilde{\sigma} - \eta, \tilde{\sigma} + \eta]$ , where

$$\tilde{\sigma} = \frac{\mu_S(\gamma_R^2 + \gamma_\varepsilon^2) - \mu_R(\gamma_S^2 + \gamma_\varepsilon^2)}{\gamma_R^2 - \gamma_S^2}.$$

b) In any equilibrium, the interval of disclosed signals is closer to the prior mean of the more confident player in terms of Hausdorff distance. In particular,  $\tilde{\sigma} \notin (\mu_S, \mu_R)$  and  $\tilde{\sigma}$  is strictly closer to the prior mean of the more confident player.

2. Let  $\gamma_S \neq \gamma_R$  and  $\mu_S = \mu_R$ . Any signal other than  $\sigma = \mu$  strictly increases disagreement and there exists no equilibrium in which any signal other than  $\mu$  is disclosed. In any equilibrium,  $S$  discloses with ex ante probability zero.

3. Let  $\gamma_S^2 = \gamma_R^2$ . If  $\mu_S = \mu_R$ , then any signal leaves disagreement equal to the (zero) prior disagreement and any disclosure rule is an equilibrium disclosure rule. If  $\mu_S \neq \mu_R$ , then any signal strictly decreases disagreement and the only equilibrium is FD.

Point 1 states that if both priors have the same variance, then all signals weakly reduce disagreement, resulting in existence of the FD-equilibrium.<sup>28</sup> This is analogous to the binary state case, where FD exists when players are equally confident as they have identical prior variances (i.e.  $\beta_S = \beta_R$  or  $\beta_S = 1 - \beta_R$ ). Note that if and only if  $\mu_S \neq \mu_R$ , all signals strictly reduce disagreement and FD is the *unique* equilibrium, which indicates a positive role of differences in prior means, as in the binary case.

Points 2 and 3 consider the case of different prior variances. Point 2 assumes  $\gamma_S \neq \gamma_R$  and  $\mu_S = \mu_R = \mu$ . Here, any signal  $\sigma \neq \mu$  strictly increases disagreement, as posterior means always differ after disclosure. For any  $\sigma \neq \mu$ , the posterior mean of the more confident player is closer to  $\mu$  than that of the other player, as a lower prior variance causes higher inertia in belief updating. In equilibrium,  $S$  always conceals  $\sigma \neq \mu$  and thereby induces a perceived disagreement of 0. In consequence,  $S$  essentially never discloses (i.e. with ex ante probability 0).

Point 3.a states that given  $\gamma_S \neq \gamma_R$  and  $\mu_S \neq \mu_R$ , equilibrium communication features a non-degenerate interval of signals that are disclosed. A difference in means, conditional

<sup>27</sup>In the binary state case, the variance of prior beliefs of player  $i$  is given by  $\beta_i(1 - \beta_i)$ , which is strictly decreasing in the distance of  $\beta_i$  from  $1/2$ .

<sup>28</sup>The result that all signals reduce disagreement under  $\gamma_i = \gamma_j$  in the normal-learning setup has also been shown in Che and Kartik (2009).

on different (though potentially arbitrary close) variances, thus improves disclosure in comparison to the case of identical means. This echoes our finding for the binary model (cf. Corollary 1, points a) and b)). Qualitatively, the equilibrium exhibits the "opinion corridor" property. Only evidence that belongs to some predetermined interval  $I$  is disclosed. Again, the underlying mechanism is that the difference in belief inertias implies that sufficiently high and sufficiently low signals increase disagreement. Finally, Point 3.b is reminiscent of Corollary 1.d, obtained in the binary setting. The set of disclosed signals is biased towards the prior mean of the more confident player. Concluding, our main qualitative insights from the analysis of the discrete state space carry over to this continuous state space setup.

## 4 Conclusion

This paper introduces a new type of belief-dependent preferences reflecting an aversion to perceived disagreement. Our analysis has identified a range of implications for important instances of strategic communication and social learning. Disagreement aversion often leads to biases in information disclosure, in which case selective disclosure is biased towards the prior mean of the most confident player. Such disclosure bias may in turn be counterproductive from an ex ante perspective, in terms of minimizing ex post perceived disagreement or actual disagreement in beliefs. Generally, more similar prior variances beneficially affect disclosure while some heterogeneity in prior means is helpful. If matching of informed and uninformed parties is endogenous, informed parties unfortunately prefer to interact with parties whose prior is similar, leading to incomplete disclosure featuring confirmatory bias.

Our results provide a plausible explanation for stylized facts such as echo chambers and increasing social polarization. Further work building on the assumption of disagreement-aversion might provide more insight into the causes and consequences of contemporary societal patterns of belief heterogeneity.

## References

- Acemoglu, D., V. Chernozhukov, and M. Yildiz (2016). Fragility of asymptotic agreement under Bayesian learning. *Theoretical Economics* 11(1), 187–225.
- Acemoglu, D., V. Chernozhukov, M. Yildiz, et al. (2007). Learning and Disagreement in an Uncertain World. Technical report, Collegio Carlo Alberto.
- Andreoni, J. and T. Mylovanov (2012). Diverging opinions. *American Economic Journal: Microeconomics* 4(1), 209–232.
- Asch, S. E. (1955). Opinions and social pressure. *Readings about the social animal* 193, 17–26.

- Aumann, R. J. (1976). Agreeing to disagree. *The annals of statistics* 4(6), 1236–1239.
- Austen-Smith, D. and T. J. Feddersen (2006). Deliberation, preference uncertainty, and voting rules. *American political science review* 100(2), 209–217.
- Baliga, S., E. Hanany, and P. Klibanoff (2013). Polarization and ambiguity. *The American Economic Review* 103(7), 3071–3083.
- Banerjee, A. and R. Somanathan (2001). A simple model of voice. *The Quarterly Journal of Economics* 116(1), 189–227.
- Battigalli, P. and M. Dufwenberg (2007). Guilt in games. *The American economic review* 97(2), 170–176.
- Battigalli, P. and M. Dufwenberg (2009). Dynamic psychological games. *Journal of Economic Theory* 144(1), 1–35.
- Bénabou, R. (2012). Groupthink: Collective delusions in organizations and markets. *The Review of Economic Studies* 80, rds030.
- Benabou, R. and G. Laroque (1992). Using privileged information to manipulate markets: Insiders, gurus, and credibility. *The Quarterly Journal of Economics* 107(3), 921–958.
- Buechel, B., T. Hellmann, and S. Klößner (2015). Opinion dynamics and wisdom under conformity. *Journal of Economic Dynamics and Control* 52, 240–257.
- Bursztyń, L., G. Egorov, and S. Fiorin (2017). From extreme to mainstream: How social norms unravel. Technical report, National Bureau of Economic Research.
- Che, Y.-K. and N. Kartik (2009). Opinions as incentives. *Journal of Political Economy* 117(5), 815–860.
- Coughlan, P. J. (2000). In defense of unanimous jury verdicts: Mistrials, communication, and strategic voting. *American Political science review* 94(2), 375–393.
- Dandekar, P., A. Goel, and D. T. Lee (2013). Biased assimilation, homophily, and the dynamics of polarization. *Proceedings of the National Academy of Sciences* 110(15), 5791–5796.
- Deutsch, M. and H. B. Gerard (1955). A study of normative and informational social influences upon individual judgment. *The journal of abnormal and social psychology* 51(3), 629.
- Dixit, A. K. and J. W. Weibull (2007). Political polarization. *Proceedings of the National Academy of Sciences* 104(18), 7351–7356.

- Domínguez, D., F. Juan, S. A. Taing, and P. Molenberghs (2016). Why do some find it hard to disagree? An fMRI study. *Frontiers in human neuroscience* 9, 718.
- Dufwenberg, M. and G. Kirchsteiger (2004). A theory of sequential reciprocity. *Games and economic behavior* 47(2), 268–298.
- Ely, J., A. Frankel, and E. Kamenica (2015). Suspense and surprise. *Journal of Political Economy* 123(1), 215–260.
- Ely, J. C. and J. Välimäki (2003). Bad reputation. *The Quarterly Journal of Economics* 118(3), 785–814.
- Estlund, D. M. (2009). *Democratic authority: A philosophical framework*. Princeton University Press.
- Feddersen, T. and W. Pesendorfer (1998). Convicting the innocent: The inferiority of unanimous jury verdicts under strategic voting. *American Political science review* 92(1), 23–35.
- Festinger, L. (1950). Informal social communication. *Psychological review* 57(5), 271.
- Festinger, L. (1957). *A theory of cognitive dissonance*. Evanston, IL: Row, Peterson.
- Geanakoplos, J., D. Pearce, and E. Stacchetti (1989). Psychological Games and Sequential Rationality. *Games and Economic Behavior* 1, 60–79.
- Gentzkow, M. and J. M. Shapiro (2006). Media bias and reputation. *Journal of political Economy* 114(2), 280–316.
- Glaeser, E. L. and C. R. Sunstein (2009). Extremism and social learning. *Journal of Legal Analysis* 1(1), 263–324.
- Goffman, E. (1959). *The presentation of self in everyday life*. Garden City, NY: Doubleday Anchor Books.
- Golman, R., G. Loewenstein, K. O. Moene, and L. Zarri (2016). The preference for belief consonance. *The Journal of Economic Perspectives* 30(3), 165–187.
- Golub, B. and M. O. Jackson (2012). How homophily affects the speed of learning and best-response dynamics. *The Quarterly Journal of Economics* 127(3), 1287–1338.
- Grossman, S. J. (1981). The Informational Role of Warranties and Private Disclosure about Product Quality. *The Journal of Law & Economics* 24(3), 461–483.
- Homans, G. C. (1961). *Human behavior: Its elementary forms*.



- Huckfeldt, R., P. E. Johnson, and J. Sprague (2004). *Political disagreement: The survival of diverse opinions within communication networks*. Cambridge University Press.
- Huston, T. L. and G. Levinger (1978). Interpersonal attraction and relationships. *Annual review of psychology* 29(1), 115–156.
- Jung, W.-O. and Y. K. Kwon (1988). Disclosure When the Market Is Unsure of Information Endowment of Managers. *Journal of Accounting Research* 26(1), 146–153.
- Kartik, N., F. X. Lee, and W. Suen (2015). Does competition promote disclosure. Technical report.
- Landemore, H. and J. Elster (2012). *Collective wisdom: Principles and mechanisms*. Cambridge University Press.
- Lazarsfeld, P. F. and R. K. Merton (1954). Friendship as a social process: A substantive and methodological analysis. *Freedom and control in modern society* 18(1), 18–66.
- Levy, G. (2007). Decision making in committees: Transparency, reputation, and voting rules. *American economic review* 97(1), 150–168.
- Levy, G. and R. Razin (2015). Correlation neglect, voting behavior, and information aggregation. *American Economic Review* 105(4), 1634–45.
- Loury, G. C. (1994). Self-censorship in public discourse: a theory of "political correctness" and related phenomena. *Rationality and Society* 6(4), 428–461.
- Milgrom, P. and J. Roberts (1986). Relying on the Information of Interested Parties. *The RAND Journal of Economics* 17(1), 18–32.
- Milgrom, P. R. (1981). Good news and bad news: Representation theorems and applications. *The Bell Journal of Economics* 12(2), 380–391.
- Morris, S. (1995). The common prior assumption in economic theory. *Economics & Philosophy* 11(2), 227–253.
- Morris, S. (2001). Political correctness. *Journal of political Economy* 109(2), 231–265.
- Mutz, D. C. (2006). *Hearing the other side: Deliberative versus participatory democracy*. Cambridge University Press.
- Newcomb, T. M. (1961). *The acquaintance process*. Holt, Rinehart & Winston.
- Ottaviani, M. and P. N. Sørensen (2006a). Reputational cheap talk. *The Rand journal of economics* 37(1), 155–175.

- Ottaviani, M. and P. N. Sørensen (2006b). The strategy of professional forecasting. *Journal of Financial Economics* 81(2), 441–466.
- Prendergast, C. (1993). A theory of "yes men". *The American Economic Review* 83(4), 757–770.
- Prentice, D. A. and D. T. Miller (1993). Pluralistic ignorance and alcohol use on campus: some consequences of misperceiving the social norm. *Journal of personality and social psychology* 64(2), 243.
- Rabin, M. (1993). Incorporating Fairness Into Game Theory and Economics. *American Economic Review* 83, 1281–1302.
- Rosenberg, M. (1954). Some determinants of political apathy. *Public Opinion Quarterly* 18(4), 349–366.
- Sethi, R. and M. Yildiz (2012). Public Disagreement. *American Economic Journal. Microeconomics* 4(3), 57.
- Sethi, R. and M. Yildiz (2016). Communication with unknown perspectives. *Econometrica* 84(6), 2029–2069.
- Shin, H. S. (1994a). The burden of proof in a game of persuasion. *Journal of Economic Theory* 64(1), 253–264.
- Shin, H. S. (1994b). News management and the value of firms. *The RAND Journal of Economics* 25(1), 58–71.
- Sobel, J. (1985). A Theory of Credibility. *Review of Economic Studies* 52, 557–573.
- Sobel, J. (2013). Giving and receiving advice. In M. Acemoglu, D. Arellano and E. Dekel (Eds.), *Advances in Economics and Econometrics: Tenth World Congress*, pp. 305–341. New York: Cambridge University Press.
- Sunstein, C. R. (2007). *Republic. com 2.0*. Princeton University Press.
- Visser, B. and O. H. Swank (2007). On committees of experts. *The Quarterly Journal of Economics* 122(1), 337–372.
- Vives, X. (2010). *Information and learning in markets: the impact of market microstructure*. Princeton University Press.

# Disliking to disagree

Florian Hoffmann, Kiryl Khalmetski, Mark T. Le Quement

## Technical Appendix

### Appendix I: Preliminaries

Throughout the proofs we use the following notation for the perceived disagreement under equilibrium of type  $X = \{D0, D1, FD, ND\}$  given the disclosed information  $d = \{0, 1, \emptyset\}$ :

$$\Delta^X(d) = \left| E_R[\tilde{\beta}_S(\sigma)|d] - \tilde{\beta}_R(d) \right|.$$

Note that trivially,  $\Delta^X(1)$  and  $\Delta^X(0)$  are actually independent of  $X$  while  $\Delta^X(\emptyset)$  is not, so that we typically omit the superscript in the first two cases. Note also that the notation is slightly abusive in the sense that it does not make explicit that  $\Delta^X(d)$  depends on  $\beta_S, \beta_R$ . Besides, it is convenient to denote the highest and the lowest prior belief as, respectively

$$\begin{aligned} \bar{\beta} &= \max\{\beta_S, \beta_R\}, \\ \underline{\beta} &= \min\{\beta_S, \beta_R\}. \end{aligned}$$

We now characterize equilibrium posterior beliefs (obtained by applying Bayes' rule) that shall be used in checking incentives in different putative equilibria. Note that in any equilibrium

$$\begin{aligned} \tilde{\beta}_i(1) &= \frac{\Pr[\sigma = 1|\omega = 1]\beta_i}{\Pr[\sigma = 1|\omega = 1]\beta_i + \Pr[\sigma = 1|\omega = 0](1 - \beta_i)} = \frac{p\beta_i}{p\beta_i + (1 - p)(1 - \beta_i)}, \\ \tilde{\beta}_i(0) &= \frac{\Pr[\sigma = 0|\omega = 1]\beta_i}{\Pr[\sigma = 0|\omega = 1]\beta_i + \Pr[\sigma = 0|\omega = 0](1 - \beta_i)} = \frac{(1 - p)\beta_i}{(1 - p)\beta_i + p(1 - \beta_i)}. \end{aligned}$$

In a ND-equilibrium,  $\tilde{\beta}_R^{ND}(\emptyset) = \beta_R$  and

$$\begin{aligned} E_R^{ND}[\tilde{\beta}_S|\emptyset] &= \Pr[\sigma = 0|d = \emptyset, ND]\tilde{\beta}_S(0) + \Pr[\sigma = 1|d = \emptyset, ND]\tilde{\beta}_S(1) \\ &\quad + \Pr[\sigma = \emptyset|d = \emptyset, ND]\beta_S \\ &= \varphi((1 - p)\beta_R + p(1 - \beta_R))\tilde{\beta}_S(0) \\ &\quad + \varphi(p\beta_R + (1 - p)(1 - \beta_R))\tilde{\beta}_S(1) \\ &\quad + (1 - \varphi)\beta_S. \end{aligned}$$

In a D1-equilibrium:

$$\begin{aligned}
\tilde{\beta}_R^{D1}(\emptyset) &= \Pr[\sigma = 0|d = \emptyset, D1]\tilde{\beta}_R(0) + \Pr[\sigma = \emptyset|d = \emptyset, D1]\beta_R \\
&= \frac{\Pr[\sigma = 0]}{\Pr[\sigma = 0] + \Pr[\sigma = \emptyset]}\tilde{\beta}_R(0) + \frac{\Pr[\sigma = \emptyset]}{\Pr[\sigma = 0] + \Pr[\sigma = \emptyset]}\beta_R \\
&= \frac{\varphi((1-p)\beta_R + p(1-\beta_R))}{\varphi((1-p)\beta_R + p(1-\beta_R)) + (1-\varphi)}\tilde{\beta}_R(0) \\
&\quad + \frac{1-\varphi}{\varphi((1-p)\beta_R + p(1-\beta_R)) + (1-\varphi)}\beta_R, \\
E_R^{D1}[\tilde{\beta}_S|\emptyset] &= \Pr[\sigma = 0|d = \emptyset, D1]\tilde{\beta}_S(0) + \Pr[\sigma = \emptyset|d = \emptyset, D1]\beta_S \\
&= \frac{\varphi((1-p)\beta_R + p(1-\beta_R))}{\varphi((1-p)\beta_R + p(1-\beta_R)) + (1-\varphi)}\tilde{\beta}_S(0) \\
&\quad + \frac{1-\varphi}{\varphi((1-p)\beta_R + p(1-\beta_R)) + (1-\varphi)}\beta_S.
\end{aligned}$$

In a D0-equilibrium:

$$\begin{aligned}
\tilde{\beta}_R^{D0}(\emptyset) &= \Pr[\sigma = 1|d = \emptyset, D0]\tilde{\beta}_R(1) + \Pr[\sigma = \emptyset|d = \emptyset, D0]\beta_R \\
&= \frac{\Pr[\sigma = 1]}{\Pr[\sigma = 1] + \Pr[\sigma = \emptyset]}\tilde{\beta}_R(1) + \frac{\Pr[\sigma = \emptyset]}{\Pr[\sigma = 1] + \Pr[\sigma = \emptyset]}\beta_R \\
&= \frac{\varphi(p\beta_R + (1-p)(1-\beta_R))}{\varphi(p\beta_R + (1-p)(1-\beta_R)) + (1-\varphi)}\tilde{\beta}_R(1) \\
&\quad + \frac{1-\varphi}{\varphi(p\beta_R + (1-p)(1-\beta_R)) + (1-\varphi)}\beta_R, \\
E_R^{D0}[\tilde{\beta}_S|\emptyset] &= \Pr[\sigma = 1|d = \emptyset, D0]\tilde{\beta}_S(1) + \Pr[\sigma = \emptyset|d = \emptyset, D0]\beta_S \\
&= \frac{\varphi(p\beta_R + (1-p)(1-\beta_R))}{\varphi(p\beta_R + (1-p)(1-\beta_R)) + (1-\varphi)}\tilde{\beta}_S(1) \\
&\quad + \frac{1-\varphi}{\varphi(p\beta_R + (1-p)(1-\beta_R)) + (1-\varphi)}\beta_S.
\end{aligned}$$

## Appendix II: Proposition 1 and Corollary 1

### Proof of Proposition 1

Proposition 1 follows from a set of Lemmas, which are stated and proved in what follows.

**Lemma II.A** *Assume  $\beta_S \neq \beta_R$ . Then,  $E_S[\Delta^{FD}] < |\beta_S - \beta_R|$ . I.e. under full disclosure, from  $S$ 's ex ante perspective, the expected value of  $R$ 's ex post perceived disagreement is strictly reduced relative to prior disagreement.*

**Proof.** Assume without loss of generality that  $\beta_S > \beta_R$ . Then, the difference between

prior disagreement and  $S$ 's ex ante expected value of ex post perceived disagreement under full disclosure is

$$\begin{aligned}
(\beta_S - \beta_R) - E_S[\Delta^{FD}] &= (\beta_S - \beta_R) \\
&\quad - \varphi(\beta_S p + (1 - \beta_S)(1 - p))\Delta(1) \\
&\quad - \varphi(\beta_S(1 - p) + (1 - \beta_S)p)\Delta(0) \\
&\quad - (1 - \varphi)(\beta_S - \beta_R) \\
&= (\beta_S - \beta_R) - \varphi(\beta_S p + (1 - \beta_S)(1 - p)) \\
&\quad \times \left( \frac{\beta_S p}{\beta_S p + (1 - \beta_S)(1 - p)} - \frac{\beta_R p}{\beta_R p + (1 - \beta_R)(1 - p)} \right) \\
&\quad - \varphi(\beta_S(1 - p) + (1 - \beta_S)p) \\
&\quad \times \left( \frac{\beta_S(1 - p)}{\beta_S(1 - p) + (1 - \beta_S)p} - \frac{\beta_R(1 - p)}{\beta_R(1 - p) + (1 - \beta_R)p} \right) \\
&\quad - (1 - \varphi)(\beta_S - \beta_R) \\
&= \varphi \frac{(\beta_S - \beta_R)(1 - \beta_R)\beta_R(2p - 1)^2}{(1 - p + \beta_R(2p - 1))(\beta_R + p(1 - 2\beta_R))} > 0.
\end{aligned}$$

Hence, from  $S$ 's ex ante perspective, the full disclosure strategy will on average reduce perceived disagreement relative to prior disagreement. ■

**Lemma II.B** *Unless  $\beta_S = \beta_R$ , there exists no ND-equilibrium (where  $S$  omits to disclose all signals).*

**Proof.**

Consider  $\beta_S > \beta_R$ . Assume by contradiction that there exists an ND-equilibrium. Then,  $R$ 's perceived disagreement conditional on no disclosure is

$$\begin{aligned}
\Delta^{ND}(\emptyset) &= \left| E_R^{ND}[\tilde{\beta}_S(\sigma)|\emptyset] - \tilde{\beta}_R^{ND}(\emptyset) \right| \\
&= \sum_{\sigma \in \{0,1,\emptyset\}} \left( P(\sigma|\beta_R) \left( \tilde{\beta}_S(\sigma) - \tilde{\beta}_R(\sigma) \right) \right),
\end{aligned}$$

where  $P(\sigma|\beta_R)$  is the ex ante probability attributed by  $R$  to signal  $\sigma \in \{0, 1, \emptyset\}$ , where  $\emptyset$  stands for no signal. If  $\beta_S \neq \beta_R$ , then Lemma II.A implies

$$\min\{\tilde{\beta}_S(0) - \tilde{\beta}_R(0), \tilde{\beta}_S(1) - \tilde{\beta}_R(1)\} < \beta_S - \beta_R.$$

Consequently, for  $\beta_S \neq \beta_R$  we have (given that  $\tilde{\beta}_S(\emptyset) - \tilde{\beta}_R(\emptyset) = \beta_S - \beta_R$ )

$$\sum_{\sigma \in \{0,1,\emptyset\}} \left( P(\sigma|\beta_R) \left( \tilde{\beta}_S(\sigma) - \tilde{\beta}_R(\sigma) \right) \right) > \min\{\tilde{\beta}_S(0) - \tilde{\beta}_R(0), \tilde{\beta}_S(1) - \tilde{\beta}_R(1)\}.$$

Hence, in a putative ND-equilibrium, for some  $\sigma \in \{0, 1\}$   $S$  would have a strict incentive to deviate by disclosing  $\sigma$ . The case of  $\beta_R > \beta_S$  proceeds analogously.

Finally, if  $\beta_S = \beta_R$ , we trivially have

$$\sum_{\sigma \in \{0, 1, \emptyset\}} \left( P(\sigma | \beta_R) \left( \tilde{\beta}_S(\sigma) - \tilde{\beta}_R(\sigma) \right) \right) = 0 = \tilde{\beta}_S(0) - \tilde{\beta}_R(0) = \tilde{\beta}_S(1) - \tilde{\beta}_R(1),$$

so that  $S$  has no strict incentive to deviate from the equilibrium strategy given  $\sigma \in \{0, 1\}$ .

■

**Lemma II.C** *If  $\beta_S = \beta_R$ , then any disclosure strategy is an equilibrium disclosure strategy.*

**Proof.** Let  $\beta_S = \beta_R$  and fix any disclosure strategy  $\tilde{D}$  of  $S$ . Denote by  $\delta_R^\eta$  the probability assigned by  $R$  to  $\sigma = \eta$  given  $d = \emptyset$ ,  $R$ 's prior being  $\beta_R$ . I.e., formally,  $\delta_R^\eta = \Pr[\sigma = \eta | d = \emptyset, \beta_R]$ . We have

$$\begin{aligned} \Delta^{\tilde{D}}(\emptyset) &= \left| E_R^{\tilde{D}}[\tilde{\beta}_S | \emptyset] - \tilde{\beta}_R(\emptyset) \right| \\ &= \left| \begin{array}{l} \delta_R^0 \tilde{\beta}_S(0) + \delta_R^1 \tilde{\beta}_S(1) + (1 - \delta_R^0 - \delta_R^1) \beta_S \\ - \delta_R^0 \tilde{\beta}_R(0) - \delta_R^1 \tilde{\beta}_R(1) - (1 - \delta_R^0 - \delta_R^1) \beta_R \end{array} \right| \\ &= \left| \begin{array}{l} \delta_R^0 (\tilde{\beta}_S(0) - \tilde{\beta}_R(0)) + \delta_R^1 (\tilde{\beta}_S(1) - \tilde{\beta}_R(1)) \\ + (1 - \delta_R^0 - \delta_R^1) (\beta_S - \beta_R) \end{array} \right| = 0, \end{aligned} \quad (4)$$

where the last equality is due to  $\beta_S = \beta_R$ . Hence,  $S$  will be indifferent between disclosure of any signal (leading to 0 posterior disagreement) and non-disclosure. Consequently, the specified disclosure strategy constitutes an equilibrium. ■

**Lemma II.D** *Let  $\beta_S \neq \beta_R$ . D0 exists if and only if  $\beta_S \leq \beta_S^*(\beta_R)$ .*

**Proof.** The D0-equilibrium exists if and only if the following  $S$ 's incentive constraints are satisfied:

$$\Delta^{D0}(0) \leq \Delta^{D0}(\emptyset) \leq \Delta^{D0}(1). \quad (5)$$

Using (4), the second incentive constraint simplifies to (denoting again  $\delta_R^\eta = \Pr[\sigma = \eta | d = \emptyset, \beta_R]$ ):

$$\Delta^{D0}(\emptyset) - \Delta^{D0}(1) \leq 0 \Leftrightarrow \quad (6)$$

$$\left| \delta_R^1 (\tilde{\beta}_S(1) - \tilde{\beta}_R(1)) + (1 - \delta_R^1) (\beta_S - \beta_R) \right| - \left| \tilde{\beta}_S(1) - \tilde{\beta}_R(1) \right| \leq 0 \Leftrightarrow \quad (7)$$

$$(1 - \delta_R^1) \left( \bar{\beta} - \underline{\beta} - \frac{\bar{\beta}p}{\bar{\beta}p + (1 - \bar{\beta})(1 - p)} + \frac{\underline{\beta}p}{\underline{\beta}p + (1 - \underline{\beta})(1 - p)} \right) \leq 0 \Leftrightarrow \quad (8)$$

$$\left(1 - \frac{\varphi(\underline{\beta}p + (1 - \underline{\beta})(1 - p))}{\varphi(\underline{\beta}p + (1 - \underline{\beta})(1 - p)) + (1 - \varphi)}\right) \times \left(\bar{\beta} - \underline{\beta} - \frac{\bar{\beta}p}{\bar{\beta}p + (1 - \bar{\beta})(1 - p)} + \frac{\underline{\beta}p}{\underline{\beta}p + (1 - \underline{\beta})(1 - p)}\right) \leq 0 \Leftrightarrow \quad (9)$$

$$\frac{(\bar{\beta} - \underline{\beta})(2p - 1)(1 - \varphi)}{\bar{\beta}(1 - p + \underline{\beta}(2p - 1)) - (1 - p)(1 - \underline{\beta})} \leq 0. \quad (10)$$

On the left-hand side of the last inequality, all terms are always positive except for the numerator, which is increasing in both  $\bar{\beta}$  and  $\underline{\beta}$  and is equal to 0 if and only if

$$\bar{\beta} = \frac{(1 - \underline{\beta})(1 - p)}{1 - p + \underline{\beta}(2p - 1)} \Leftrightarrow \underline{\beta} = \frac{(1 - \bar{\beta})(1 - p)}{1 - p + \bar{\beta}(2p - 1)}.$$

Thus, independently of whether  $\beta_S = \bar{\beta}$  or  $\beta_S = \underline{\beta}$  (i.e. of whether  $\beta_S > \beta_R$  or  $\beta_S < \beta_R$ ) we have

$$\Delta^{D0}(\emptyset) - \Delta^{D0}(1) \leq 0 \text{ if and only if } \beta_S \leq \frac{(1 - \beta_R)(1 - p)}{1 - p + \beta_R(2p - 1)} = \beta_S^*(\beta_R). \quad (11)$$

Note further that  $\Delta^{D0}(\emptyset) - \Delta^{D0}(1) \leq 0$  implies  $\Delta^{D0}(0) \leq \Delta^{D0}(\emptyset)$ . Indeed, otherwise we would have  $\Delta^{D0}(\emptyset) \leq \min\{\Delta^{D0}(0), \Delta^{D0}(1)\}$ , i.e. it would hold true that

$$\delta_R^1 |\tilde{\beta}_S(1) - \tilde{\beta}_R(1)| + (1 - \delta_R^1) |\beta_S - \beta_R| \leq \min \left\{ \left| \tilde{\beta}_S(0) - \tilde{\beta}_R(0) \right|, \left| \tilde{\beta}_S(1) - \tilde{\beta}_R(1) \right| \right\}.$$

This, in turn, yields

$$|\beta_S - \beta_R| \leq |\tilde{\beta}_S(\eta) - \tilde{\beta}_R(\eta)|, \forall \eta \in \{0, 1\}.$$

The latter would imply that (a putative) full disclosure equilibrium does not yield an expected value of ex post perceived disagreement that is strictly smaller than prior disagreement. But this contradicts Lemma II.A. Thus, (5) holds if and only if  $\beta_S \leq \beta_S^*(\beta_R)$ . ■

**Lemma II.E** *Let  $\beta_S \neq \beta_R$ . D1 exists if and only if  $\beta_S \geq \beta_S^{**}(\beta_R)$ .*

**Proof.** The D1-equilibrium exists if and only if the following  $S$ 's incentive constraints are satisfied:

$$\Delta^{D1}(1) \leq \Delta^{D1}(\emptyset) \leq \Delta^{D1}(0). \quad (12)$$

Using (4), the second incentive constraint simplifies to (denoting again  $\delta_R^\eta = \Pr[\sigma =$

$\eta|d = \emptyset, \beta_R]$ :

$$\Delta^{D1}(\emptyset) - \Delta^{D1}(0) \leq 0 \Leftrightarrow \quad (13)$$

$$\left| \delta_R^0(\tilde{\beta}_S(0) - \tilde{\beta}_R(0)) + (1 - \delta_R^0)(\beta_S - \beta_R) \right| - \left| \tilde{\beta}_S(0) - \tilde{\beta}_R(0) \right| \leq 0 \Leftrightarrow \quad (14)$$

$$(1 - \delta_R^0) \left( \bar{\beta} - \underline{\beta} - \frac{\bar{\beta}(1-p)}{\bar{\beta}(1-p) + (1-\bar{\beta})p} + \frac{\underline{\beta}(1-p)}{\underline{\beta}(1-p) + (1-\underline{\beta})p} \right) \leq 0 \Leftrightarrow \quad (15)$$

$$\begin{aligned} & \left( 1 - \frac{\varphi(\underline{\beta}(1-p) + (1-\underline{\beta})p)}{\varphi(\underline{\beta}(1-p) + (1-\underline{\beta})p) + (1-\varphi)} \right) \\ & \times \left( \bar{\beta} - \underline{\beta} - \frac{\bar{\beta}(1-p)}{\bar{\beta}(1-p) + (1-\bar{\beta})p} + \frac{\underline{\beta}(1-p)}{\underline{\beta}(1-p) + (1-\underline{\beta})p} \right) \leq 0 \Leftrightarrow \quad (16) \\ & (\bar{\beta} - \underline{\beta})(2p-1)(1-\varphi) \end{aligned}$$

$$\times \frac{\bar{\beta}(\underline{\beta}(2p-1) - p) + p(1-\underline{\beta})}{(p - \bar{\beta}(2p-1))(p - \underline{\beta}(2p-1))(1 - \varphi(1-p + \underline{\beta}(2p-1)))} \leq 0. \quad (17)$$

On the left-hand side of the last inequality, all terms are always positive except for the numerator, which is decreasing in both  $\bar{\beta}$  and  $\underline{\beta}$  and is equal to 0 if and only if

$$\bar{\beta} = \frac{p(1-\underline{\beta})}{\underline{\beta} + p(1-2\underline{\beta})} \Leftrightarrow \underline{\beta} = \frac{p(1-\bar{\beta})}{\bar{\beta} + p(1-2\bar{\beta})}.$$

Thus, independently of whether  $\beta_S = \bar{\beta}$  or  $\beta_S = \underline{\beta}$  (i.e. of whether  $\beta_S > \beta_R$  or  $\beta_S < \beta_R$ ) we have

$$\Delta^{D1}(\emptyset) - \Delta^{D1}(0) \leq 0 \text{ if and only if } \beta_S \geq \frac{p(1-\beta_R)}{\beta_R + p(1-2\beta_R)} = \beta_S^{**}(\beta_R). \quad (18)$$

Note further that  $\Delta^{D1}(\emptyset) - \Delta^{D1}(0) \leq 0$  implies  $\Delta^{D1}(1) \leq \Delta^{D1}(\emptyset)$  by the same argument as in the proof of Lemma II.D. Thus, (12) holds if and only if  $\beta_S \geq \beta_S^{**}(\beta_R)$ . ■

**Lemma II.F** Let  $\beta_S \neq \beta_R$ . *FD exists if and only if  $\beta_S \in [\beta_S^*(\beta_R), \beta_S^{**}(\beta_R)]$ .*

**Proof.** The FD-equilibrium exists if and only if the following  $S$ 's incentive constraints are satisfied:

$$|\beta_S - \beta_R| \geq \Delta^{FD}(1), \quad (19)$$

$$|\beta_S - \beta_R| \geq \Delta^{FD}(0). \quad (20)$$

Note that the reverse inequality to (19) holds under the same conditions as (8), which in turn is equivalent to (6). Hence, by the proof of Lemma II.D  $|\beta_S - \beta_R| \leq \Delta^{FD}(1)$  iff  $\beta_S \leq \beta_S^*(\beta_R)$  (with  $|\beta_S - \beta_R| = \Delta^{FD}(1)$  iff  $\beta_S = \beta_S^*(\beta_R)$ ). Consequently, (19) holds if and only if  $\beta_S \geq \beta_S^*(\beta_R)$ . Analogously, from the proof of Lemma II.E we obtain that (20) holds if and only if  $\beta_S \leq \beta_S^{**}(\beta_R)$ . Hence, both constraints hold simultaneously if and only if



$\beta_S \in [\beta_S^*(\beta_R), \beta_S^{**}(\beta_R)]$ . ■

**Lemma II.G** Let  $\beta_S \neq \beta_R$ . Mixed strategy equilibria exist if and only if  $\beta_S \in \{\beta_S^*(\beta_R), \beta_S^{**}(\beta_R)\}$ .

**Proof.** First, if  $\beta_S \neq \beta_R$ , there cannot be an equilibrium in which  $S$ 's disclosure strategy (call this strategy  $M$ ) specifies omitting to disclose with a non-degenerate probability after both signals 0 and 1. Indeed, by the same arguments as in the proof of Lemma II.B (using the same notation), one can show that in such an equilibrium (call it an  $M$ -equilibrium), it will be true that

$$\Delta^M(\emptyset) = \sum_{\sigma} P(\sigma | \beta_R) |\beta_S(\sigma) - \beta_R(\sigma)| > \min\{|\tilde{\beta}_S(0) - \tilde{\beta}_R(0)|, |\tilde{\beta}_S(1) - \tilde{\beta}_R(1)|\}.$$

Hence, after some signal  $\tilde{\sigma} \in \{0, 1\}$ , ex post perceived disagreement will be strictly smaller than  $\Delta^M(\emptyset)$ . But given this,  $S$  would deviate to disclosing for sure when holding signal  $\tilde{\sigma}$ .

Consider now the remaining case of an equilibrium in which  $S$ 's disclosure strategy (call this strategy  $\tilde{M}$ ) specifies mixing between disclosure and non-disclosure only for one signal  $\sigma^* \in \{0, 1\}$ .

For the case of  $\sigma^* = 1$ , such an equilibrium requires the indifference condition  $\Delta^{\tilde{M}}(\emptyset) - \Delta(1) = 0$ . Letting  $\delta_R^\eta$  denote the probability, in the eyes of  $R$ , of  $S$  holding a signal  $\eta$  conditional on no disclosure, this is equivalent to:

$$\left| \delta_R^1(\tilde{\beta}_S(1) - \tilde{\beta}_R(1)) + (1 - \delta_R^1)(\beta_S - \beta_R) \right| - \left| \tilde{\beta}_S(1) - \tilde{\beta}_R(1) \right| = 0 \Leftrightarrow \quad (21)$$

$$(1 - \delta_R^1) \left( \bar{\beta} - \underline{\beta} - \frac{\bar{\beta}p}{\bar{\beta}p + (1 - \bar{\beta})(1 - p)} + \frac{\underline{\beta}p}{\underline{\beta}p + (1 - \underline{\beta})(1 - p)} \right) = 0 \Leftrightarrow \quad (22)$$

$$\bar{\beta} - \underline{\beta} - \frac{\bar{\beta}p}{\bar{\beta}p + (1 - \bar{\beta})(1 - p)} + \frac{\underline{\beta}p}{\underline{\beta}p + (1 - \underline{\beta})(1 - p)} = 0 \Leftrightarrow \quad (23)$$

$$\left\{ \underline{\beta}, \frac{(1 - \underline{\beta})(1 - p)}{1 - p + \underline{\beta}(2p - 1)} \right\} = \bar{\beta} \Leftrightarrow \quad (24)$$

$$\left\{ \bar{\beta}, \frac{(1 - \bar{\beta})(1 - p)}{1 - p + \bar{\beta}(2p - 1)} \right\} = \underline{\beta}. \quad (25)$$

Thus, given  $\beta_S \neq \beta_R$ ,  $S$  is indifferent between sending 1 and no disclosure if and only if  $\beta_S = \beta_S^*(\beta_R)$ . It then holds true that  $\Delta^{\tilde{M}}(\emptyset) > \Delta(0)$  since

$$\min\{|\tilde{\beta}_S(0) - \tilde{\beta}_R(0)|, |\tilde{\beta}_S(1) - \tilde{\beta}_R(1)|\} < |\beta_S - \beta_R|$$

by Lemma II.A.

For the case of  $\sigma^* = 0$ , analogously, such an equilibrium requires the indifference condition

$\Delta^{\widetilde{M}}(\emptyset) - \Delta^{\widetilde{M}}(0) = 0$ . This is further equivalent to:

$$\left| \delta_R^0(\widetilde{\beta}_S(0) - \widetilde{\beta}_R(0)) + (1 - \delta_R^0)(\beta_S - \beta_R) \right| - \left| \widetilde{\beta}_S(0) - \widetilde{\beta}_R(0) \right| = 0 \Leftrightarrow \quad (26)$$

$$(1 - \delta_R^0) \left( \overline{\beta} - \underline{\beta} - \frac{\overline{\beta}(1-p)}{\overline{\beta}(1-p) + (1-\overline{\beta})p} + \frac{\underline{\beta}(1-p)}{\underline{\beta}(1-p) + (1-\underline{\beta})p} \right) = 0 \Leftrightarrow \quad (27)$$

$$\left( \overline{\beta} - \underline{\beta} - \frac{\overline{\beta}(1-p)}{\overline{\beta}(1-p) + (1-\overline{\beta})p} + \frac{\underline{\beta}(1-p)}{\underline{\beta}(1-p) + (1-\underline{\beta})p} \right) = 0 \Leftrightarrow \quad (28)$$

$$\left\{ \underline{\beta}, \frac{p(1-\underline{\beta})}{\underline{\beta} + p(1-2\underline{\beta})} \right\} = \overline{\beta} \Leftrightarrow \quad (29)$$

$$\left\{ \overline{\beta}, \frac{p(1-\overline{\beta})}{\overline{\beta} + p(1-2\overline{\beta})} \right\} = \underline{\beta}. \quad (30)$$

This similarly leads to  $\beta_S = \beta_S^{**}(\beta_R)$ . ■

### Proof of Corollary 1.

**Step 1.** Point a) follows from the fact that  $\beta_S^*(\beta_R) < 1 - \beta_R < \beta_S^{**}(\beta_R)$ . By Proposition 1, this means that for  $\varepsilon$  small enough,  $\beta_S \in (1 - \beta_R - \varepsilon, 1 - \beta_R + \varepsilon)$  satisfies conditions such that the FD-equilibrium is the unique equilibrium.

**Step 2.** This proves Point b). Consider first  $\beta_R < 1/2$  and  $\beta_S$  sufficiently close to  $\beta_R$ . Then, by Proposition 1, given that  $\beta_S^*(\beta_R) < 1 - \beta_R < \beta_S^{**}(\beta_R)$ , the equilibrium features no full disclosure if and only if  $\beta_R < \beta_S^*(\beta_R, p)$ . In turn, for  $\beta_R < 1/2$  we have

$$\begin{aligned} \beta_R &< \beta_S^*(\beta_R, p) \Leftrightarrow \\ p &< \frac{(1 - \beta_R)^2}{(1 - \beta_R)^2 + (\beta_R)^2}. \end{aligned} \quad (31)$$

Hence, for any  $\beta_S$  sufficiently close to  $\beta_R < 1/2$  under the above condition we obtain  $\beta_S < \beta_S^*(\beta_R, p)$ , in which case D0 is the unique equilibrium.

Consider  $\beta_R > 1/2$  and  $\beta_S$  sufficiently close to  $\beta_R$ . Then, by Proposition 1, given that  $\beta_S^*(\beta_R) < 1 - \beta_R < \beta_S^{**}(\beta_R)$ , the equilibrium features no full disclosure if and only if  $\beta_R > \beta_S^{**}(\beta_R, p)$ . In turn, for  $\beta_R > 1/2$  we have

$$\begin{aligned} \beta_R &> \beta_S^{**}(\beta_R, p) \Leftrightarrow \\ p &< \frac{(\beta_R)^2}{(1 - \beta_R)^2 + (\beta_R)^2}. \end{aligned} \quad (32)$$

Hence, for any  $\beta_S$  sufficiently close to  $\beta_R > 1/2$  under the above condition we obtain  $\beta_S > \beta_S^{**}(\beta_R, p)$ , in which case D1 is the unique equilibrium.

Finally, note that (31) and (32) combine into

$$\beta_R \notin [\beta_S^*(\beta_R, p), \beta_S^{**}(\beta_R, p)] \Leftrightarrow$$

$$p < \max \left\{ \frac{(1 - \beta_R)^2}{(1 - \beta_R)^2 + (\beta_R)^2}, \frac{(\beta_R)^2}{(1 - \beta_R)^2 + (\beta_R)^2} \right\}.$$

This together with Proposition 1 leads to the claim.

**Step 3.** Point c) follows due to  $\beta_S^*(\beta_R, p)$  (resp.  $\beta_S^{**}(\beta_R, p)$ ) being continuously decreasing (resp. increasing) in  $p$  and being equal to 0 (1) if  $p = 1$ .

**Step 4.** This proves point d). Let  $\beta_R < 1/2$ , i.e.  $R$  is biased towards 0.

By Proposition 1, a D1-equilibrium exists if and only if  $\beta_S \geq \beta_S^{**}(\beta_R) > 1 - \beta_R$ . This implies that  $\beta_S$  is closer to 1 than  $\beta_R$  is close to 0, meaning that  $\beta_S$  is biased towards 1 and  $S$  is more confident than  $R$ .

By Proposition 1, a D0-equilibrium exists if and only if  $\beta_S \leq \beta_S^*(\beta_R) < 1 - \beta_R$ . This in turn is compatible with two cases: Either  $\beta_R$  is the most confident prior or  $\beta_S$  is the most confident prior, in which case it also holds true that  $\beta_S \leq \beta_R$ . In both cases, note that the most confident of the two priors is smaller than  $\frac{1}{2}$ , i.e. the more confident player is biased towards 0.

Let  $R$  be biased towards 1 ( $\beta_R \geq 1/2$ ). The symmetric argument as given for the case of  $\beta_R < 1/2$  applies. ■

## Appendix III: Propositions 3 and 4

### Proof of Proposition 3

**Step 1.** Consider the case  $\beta_S > \beta_R$  in D0-equilibrium. From  $S$ 's ex ante perspective, the expected ex post perceived disagreement is

$$E_S[\Delta^{D0}] = (\Pr[\sigma = 1 | \beta_S] + \Pr[\sigma = \emptyset | \beta_S])(E_R^{D0}[\tilde{\beta}_S | \emptyset] - \tilde{\beta}_R^{D0}(\emptyset))$$

$$+ \Pr[\sigma = 0 | \beta_S](\tilde{\beta}_S(0) - \tilde{\beta}_R(0)).$$

At the same time, under full disclosure

$$E_S[\Delta^{FD}] = \Pr[\sigma = 1 | \beta_S](\tilde{\beta}_S(1) - \tilde{\beta}_R(1)) + \Pr[\sigma = 0 | \beta_S](\tilde{\beta}_S(0) - \tilde{\beta}_R(0))$$

$$+ \Pr[\sigma = \emptyset | \beta_S](\beta_S - \beta_R).$$

Using the expressions obtained in Appendix I, it follows that:

$$\begin{aligned}
& E_S[\Delta^{D0}] - E_S[\Delta^{FD}] \\
&= \Pr[\sigma = 1 | \beta_S](E_R^{D0}[\tilde{\beta}_S | \emptyset] - \tilde{\beta}_R^{D0}(\emptyset) - (\tilde{\beta}_S(1) - \tilde{\beta}_R(1))) \\
&\quad + \Pr[\sigma = \emptyset | \beta_S](E_R^{D0}[\tilde{\beta}_S | \emptyset] - \tilde{\beta}_R^{D0}(\emptyset) \\
&\quad - (\beta_S - \beta_R)) \\
&= \varphi(\beta_S p + (1 - \beta_S)(1 - p)) \\
&\quad \times \left( \left( \frac{\varphi(\beta_R p + (1 - \beta_R)(1 - p))}{\varphi(\beta_R p + (1 - \beta_R)(1 - p)) + (1 - \varphi)} - 1 \right) (\tilde{\beta}_S(1) - \tilde{\beta}_R(1)) \right) \\
&\quad \quad + \left( \frac{(1 - \varphi)}{\beta_R \varphi p + (1 - \beta_R)\varphi(1 - p) + (1 - \varphi)} \right) (\beta_S - \beta_R) \\
&\quad + (1 - \varphi) \left( \left( \frac{\varphi(\beta_R p + (1 - \beta_R)(1 - p))}{\varphi(\beta_R p + (1 - \beta_R)(1 - p)) + (1 - \varphi)} \right) (\tilde{\beta}_S(1) - \tilde{\beta}_R(1)) \right) \\
&\quad \quad + \left( \frac{(1 - \varphi)}{\beta_R \varphi p + (1 - \beta_R)\varphi(1 - p) + (1 - \varphi)} - 1 \right) (\beta_S - \beta_R) \\
&= \Phi_1 \Phi_2
\end{aligned}$$

where

$$\begin{aligned}
\Phi_1 &= \frac{(\beta_S - \beta_R)^2 (1 - 2p)^2 (1 - \varphi) \varphi}{(\beta_R p + (1 - \beta_R)(1 - p))(\beta_S p + (1 - \beta_S)(1 - p))(1 - p\varphi + \beta_R \varphi(2p - 1))} > 0, \\
\Phi_2 &= (\beta_R + \beta_S - 1)(1 - p) + \beta_R \beta_S (2p - 1).
\end{aligned}$$

Note that  $\Phi_2$  is an increasing function of  $\beta_S$ . At the same time, by Proposition 1, it must be true that  $\beta_S < \beta_S^*$  if the D0-equilibrium is the unique equilibrium. Consequently,

$$\begin{aligned}
\Phi_2(\beta_S) &< \Phi_2(\beta_S^*) = \left( \beta_R + \frac{(1 - \beta_R)(1 - p)}{1 - p + \beta_R(2p - 1)} - 1 \right) (1 - p) \\
&\quad + \beta_R \frac{(1 - \beta_R)(1 - p)}{1 - p + \beta_R(2p - 1)} (2p - 1) \\
&= 0.
\end{aligned}$$

Hence,  $\Phi_1 \Phi_2 < 0$  so that

$$E_S[\Delta^{D0}] - E_S[\Delta^{FD}] < 0,$$

i.e. the sender would ex ante prefer D0 over FD.

**Step 2.** Consider the case  $\beta_S > \beta_R$  in D1-equilibrium. From  $S$ 's perspective, the ex ante expected ex post perceived disagreement is

$$\begin{aligned}
E_S[\Delta^{D1}] &= (\Pr[\sigma = 0 | \beta_S] + \Pr[\sigma = \emptyset | \beta_S])(E_R^{D1}[\tilde{\beta}_S | \emptyset] - \tilde{\beta}_R^{D1}(\emptyset)) \\
&\quad + \Pr[\sigma = 1 | \beta_S](\tilde{\beta}_S(1) - \tilde{\beta}_R(1))
\end{aligned}$$

It follows that

$$\begin{aligned}
& E_S[\Delta^{D1}] - E_S[\Delta^{FD}] \\
&= \varphi(\beta_S(1-p) + (1-\beta_S)p) \\
&\quad \times \left( \left( \frac{\varphi(\beta_R(1-p) + (1-\beta_R)p)}{\varphi(\beta_R(1-p) + (1-\beta_R)p) + (1-\varphi)} - 1 \right) (\tilde{\beta}_S(0) - \tilde{\beta}_R(0)) \right) \\
&\quad \quad + \left( \frac{(1-\varphi)}{\varphi(\beta_R(1-p) + (1-\beta_R)p) + (1-\varphi)} \right) (\beta_S - \beta_R) \\
&+ (1-\varphi) \left( \left( \frac{\varphi(\beta_R(1-p) + (1-\beta_R)p)}{\varphi(\beta_R(1-p) + (1-\beta_R)p) + (1-\varphi)} \right) (\tilde{\beta}_S(0) - \tilde{\beta}_R(0)) \right) \\
&\quad \quad + \left( \frac{(1-\varphi)}{\varphi(\beta_R(1-p) + (1-\beta_R)p) + (1-\varphi)} - 1 \right) (\beta_S - \beta_R) \\
&= \Phi_3 \Phi_4,
\end{aligned}$$

where

$$\begin{aligned}
\Phi_3 &= -\frac{(\beta_S - \beta_R)^2(1-2p)^2(1-\varphi)\varphi}{(\beta_R(1-p) + (1-\beta_R)p)(\beta_S(1-p) + (1-\beta_S)p)} \\
&\quad \times \frac{1}{1 - \varphi((1-\beta_R)(1-p) + \beta_R p)} \\
&< 0, \\
\Phi_4 &= p(1-\beta_R) - \beta_S(p(1-\beta_R) + \beta_R(1-p)).
\end{aligned}$$

Function  $\Phi_4$  is decreasing in  $\beta_S$ . At the same time, by Proposition 1 it must be true that  $\beta_S > \beta_S^{**}$  if the D1-equilibrium is the unique equilibrium. Consequently,

$$\Phi_4(\beta_S) < \Phi_4(\beta_S^{**}) = p(1-\beta_R) - \frac{p(1-\beta_R)}{\beta_R + p(1-2\beta_R)}(p(1-\beta_R) + \beta_R(1-p)) = 0.$$

Hence,  $\Phi_3 \Phi_4 > 0$ , i.e.

$$E_S[\Delta^{D1}] - E_S[\Delta^{FD}] > 0,$$

i.e. the sender would ex ante prefer FD over D1.

**Step 3.** Consider the case  $\beta_S < \beta_R$ . Then, it can be shown that

$$\begin{aligned}
E_S[\Delta^{D0}] - E_S[\Delta^{FD}] &= -\Phi_1 \Phi_2 > 0, \\
E_S[\Delta^{D1}] - E_S[\Delta^{FD}] &= -\Phi_3 \Phi_4 < 0.
\end{aligned}$$

Thus, the sender would ex ante prefer FD over D0 and D1 over FD whenever D0 and D1 are the unique equilibria, respectively. ■

### Proof of Proposition 4

**Step 1.** In Steps 1-4 below, we consider the case that  $\beta_S > \beta_R$ . Define as  $\tilde{\Theta}(\text{Partial}, \hat{\beta})$  and  $\tilde{\Theta}(\text{Full}, \hat{\beta})$  the expected actual disagreement under partial and full disclosure respectively,

from the perspective of a third party endowed with prior  $\widehat{\beta}$ . Denote further by  $\widetilde{\beta}_i(\iota, \text{Partial})$  and  $\widetilde{\beta}_i(\iota, \text{Full})$  the posterior of player  $i$  conditional on obtained information  $\iota$  under partial and full disclosure respectively. We have:

$$\begin{aligned}
\widetilde{\Theta}(\text{Partial}, \widehat{\beta}) &= E_{\widehat{\beta}} \left[ \left| \widetilde{\beta}_S(\sigma, \text{Partial}) - \widetilde{\beta}_R(d, \text{Partial}) \right| \right] \\
&\geq E_{\widehat{\beta}} \left[ \widetilde{\beta}_S(\sigma, \text{Partial}) - \widetilde{\beta}_R(d, \text{Partial}) \right] \\
&= E_{\widehat{\beta}}[\widetilde{\beta}_S(\sigma, \text{Partial})] - E_{\widehat{\beta}}[\widetilde{\beta}_R(d, \text{Partial})] \\
&= E_{\widehat{\beta}} \left[ \widetilde{\beta}_S(\sigma, \text{Full}) \right] - E_{\widehat{\beta}} \left[ \widetilde{\beta}_R(d, \text{Partial}) \right]. \tag{33}
\end{aligned}$$

In the above, the equality  $E_{\widehat{\beta}}[\widetilde{\beta}_S(\sigma, \text{Partial})] = E_{\widehat{\beta}}[\widetilde{\beta}_S(\sigma, \text{Full})]$  follows from the fact that  $S$ 's expected posterior is independent of the disclosure rule. Note on the other hand that

$$\begin{aligned}
\widetilde{\Theta}(\text{Full}, \widehat{\beta}) &= E_{\widehat{\beta}} \left[ \left| \widetilde{\beta}_S(\sigma, \text{Full}) - \widetilde{\beta}_R(d, \text{Full}) \right| \right] \\
&= E_{\widehat{\beta}} \left[ \widetilde{\beta}_S(\sigma, \text{Full}) \right] - E_{\widehat{\beta}} \left[ \widetilde{\beta}_R(d, \text{Full}) \right]. \tag{34}
\end{aligned}$$

To see this, note that under FD, it always holds true that  $d = \sigma$ . Recall also that  $\widetilde{\beta}_S(\sigma) > \widetilde{\beta}_R(\sigma)$  for any  $\sigma$  given  $\beta_S > \beta_R$ .

It follows from the above analysis that

$$\widetilde{\Theta}(\text{Partial}, \widehat{\beta}) - \widetilde{\Theta}(\text{Full}, \widehat{\beta}) \geq E_{\widehat{\beta}} \left[ \widetilde{\beta}_R(d, \text{Full}) \right] - E_{\widehat{\beta}} \left[ \widetilde{\beta}_R(d, \text{Partial}) \right]. \tag{35}$$

**Step 2.** We now show that  $E_{\widehat{\beta}} \left[ \widetilde{\beta}_R(d, \text{Full}) \right] - E_{\widehat{\beta}} \left[ \widetilde{\beta}_R(d, \text{Partial}) \right] > 0$  if and only if  $\widehat{\beta} > \beta_R$ . Here we follow the steps of the analysis presented in Kartik et al. (2015). One can verify that

$$\widetilde{\beta}_R(d) = \frac{\widetilde{\beta}(d)^{\frac{\beta_R}{\beta}}}{\widetilde{\beta}(d)^{\frac{\beta_R}{\beta}} + (1 - \widetilde{\beta}(d))^{\frac{1 - \beta_R}{1 - \beta}}},$$

where  $\widetilde{\beta}(d)$  denotes the hypothetical posterior belief of  $R$  if she had a prior  $\widehat{\beta}$  and observed  $d$ . One can verify that the above function is concave in  $\widetilde{\beta}(d)$  if  $\beta < \beta_R$  and convex if the opposite inequality holds. Blackwell (1953) has shown that a garbling increases (resp. reduces) an individual's expectation of any concave (resp. convex) function of his posterior. Then, since partial disclosure is a garbling of full disclosure,<sup>29</sup> we obtain that

$$E_{\widehat{\beta}} \left[ \widetilde{\beta}_R(d, \text{Partial}) \right] < (>) E_{\widehat{\beta}} \left[ \widetilde{\beta}_R(d, \text{Full}) \right] \text{ if } \widehat{\beta} > (<) \beta_R \tag{36}$$

given that  $R$ 's posterior is a convex (concave) function of  $\widetilde{\beta}(\sigma)$  if  $\widehat{\beta} > (<) \beta_R$ .

<sup>29</sup>See Kartik et al. (2015) for a formal definition of garbling.

**Step 3.** (35) and (36) together imply

$$\tilde{\Theta}(\text{Partial}, \hat{\beta}) - \tilde{\Theta}(\text{Full}, \hat{\beta}) > 0 \text{ if } \hat{\beta} > \beta_R.$$

Thus, the third party would prefer full disclosure over partial disclosure whenever  $\hat{\beta} > \beta_R$ , i.e. whenever  $\beta_R < \hat{\beta} < \beta_S$  or  $\hat{\beta} \geq \beta_S > \beta_R$ .

**Step 4.** Consider  $\hat{\beta} < \beta_R < \beta_S$ . If  $\beta_S$  is sufficiently close to 1, then we have:

$$\begin{aligned} \tilde{\Theta}(\text{Partial}, \hat{\beta}) &= E_{\hat{\beta}} \left[ \left| \tilde{\beta}_S(\sigma, \text{Partial}) - \tilde{\beta}_R(d, \text{Partial}) \right| \right] \\ &= E_{\hat{\beta}} \left[ \tilde{\beta}_S(\sigma, \text{Partial}) - \tilde{\beta}_R(d, \text{Partial}) \right] \\ &= E_{\hat{\beta}}[\tilde{\beta}_S(\sigma, \text{Partial})] - E_{\hat{\beta}}[\tilde{\beta}_R(d, \text{Partial})] \\ &= E_{\hat{\beta}}[\tilde{\beta}_S(\sigma, \text{Full})] - E_{\hat{\beta}}[\tilde{\beta}_R(d, \text{Partial})]. \end{aligned}$$

Note in the above that we have equalities at all stages in contrast to (33). This together with (34) and (36) implies

$$\tilde{\Theta}(\text{Partial}, \hat{\beta}) - \tilde{\Theta}(\text{Full}, \hat{\beta}) = E_{\hat{\beta}}[\tilde{\beta}_R(d, \text{Full})] - E_{\hat{\beta}}[\tilde{\beta}_R(d, \text{Partial})] < 0.$$

Hence, in this case the third party would prefer partial disclosure over full disclosure in terms of minimizing expected actual disagreement.

**Step 5.** The proof for the remaining case of  $\beta_S < \beta_R$  is conceptually identical to what has been presented, and is hence omitted. We obtain the following counterparts of the statements proven above:

$$\begin{aligned} \tilde{\Theta}(\text{Partial}, \hat{\beta}) - \tilde{\Theta}(\text{Full}, \hat{\beta}) &> 0 \text{ if } \hat{\beta} < \beta_R, \\ \tilde{\Theta}(\text{Partial}, \hat{\beta}) - \tilde{\Theta}(\text{Full}, \hat{\beta}) &< 0 \text{ if } \beta_S < \beta_R < \hat{\beta} \text{ and } \beta_S \text{ is close to } 0. \end{aligned}$$

■

## Appendix IV: Proposition 5

### Proof of Proposition 5.a)

**Step 1.** Consider a putative FD-equilibrium. Let  $G_S(G_R)$  denote the (symmetric) cumulative distribution function of  $S$ 's ( $R$ 's) prior belief. Then, if the sender discloses 0-signal, the receiver with the prior  $\beta_R$  believes that the disagreement is

$$\Delta(0) = \int_{\beta_S=0}^1 \left| \tilde{\beta}_S(0) - \tilde{\beta}_R(0) \right| dG_S(\beta_S).$$

In turn, the sender expects that the receiver's perceived disagreement is

$$E_S[\Delta(0)] = \int_{\beta_R=0}^1 \int_{\beta_S=0}^1 \left| \tilde{\beta}_S(0) - \tilde{\beta}_R(0) \right| dG_S(\beta_S) dG_R(\beta_R).$$

If the sender does not disclose, the expected perceived disagreement is

$$E_S[\Delta^{FD}(\emptyset)] = \int_{\beta_R=0}^1 \int_{\beta_S=0}^1 |\beta_S - \beta_R| dG_S(\beta_S) dG_R(\beta_R).$$

In FD-equilibrium we must have  $E_S[\Delta(0)] - E_S[\Delta^{FD}(\emptyset)] < 0$ . We have

$$\begin{aligned} & E_S[\Delta(0)] - E_S[\Delta^{FD}(\emptyset)] \\ &= \int_{\beta_R=0}^1 \int_{\beta_S=0}^1 \left( \left| \tilde{\beta}_S(0) - \tilde{\beta}_R(0) \right| - |\beta_S - \beta_R| \right) dG_S(\beta_S) dG_R(\beta_R). \end{aligned}$$

Denote  $\tilde{\beta}(\sigma, \beta)$  the posterior belief given obtained/disclosed signal  $\sigma$  and prior belief  $\beta$ . Besides, denote  $\kappa(\beta_i, \beta_j) = \left| \tilde{\beta}(0, \beta_i) - \tilde{\beta}(0, \beta_j) \right| - |\beta_i - \beta_j|$ . Then,

$$\begin{aligned} & \int_{\beta_R=0}^1 \int_{\beta_S=0}^1 \left( \left| \tilde{\beta}_S(0) - \tilde{\beta}_R(0) \right| - |\beta_S - \beta_R| \right) dG_S(\beta_S) dG_R(\beta_R) \\ &= \int_{\beta_R=0}^1 \int_{\beta_S=0}^1 \kappa(\beta_S, \beta_R) dG_S(\beta_S) dG_R(\beta_R) \\ &= \int_{\beta_R=0}^{0.5} \int_{\beta_S=0}^1 \kappa(\beta_S, \beta_R) dG_S(\beta_S) dG_R(\beta_R) \\ & \quad + \int_{\beta_R=0}^{0.5} \int_{\beta_S=0}^1 \kappa(\beta_S, 1 - \beta_R) dG_S(\beta_S) dG_R(1 - \beta_R) \\ &= \int_{\beta_R=0}^{0.5} \int_{\beta_S=0}^1 \kappa(\beta_S, \beta_R) dG_S(\beta_S) dG_R(\beta_R) \\ & \quad + \int_{\beta_R=0}^{0.5} \int_{\beta_S=0}^1 \kappa(\beta_S, 1 - \beta_R) dG_S(\beta_S) dG_R(\beta_R) \\ &= \int_{\beta_R=0}^{0.5} \int_{\beta_S=0}^1 (\kappa(\beta_S, \beta_R) + \kappa(\beta_S, 1 - \beta_R)) dG_S(\beta_S) dG_R(\beta_R), \end{aligned}$$



where the third equality follows due to symmetry of  $G$ . Next, denote  $\lambda(\beta_S, \beta_R) = \kappa(\beta_S, \beta_R) + \kappa(\beta_S, 1 - \beta_R)$ . Then, similarly,

$$\begin{aligned}
& \int_{\beta_R=0}^{0.5} \int_{\beta_S=0}^1 (\kappa(\beta_S, \beta_R) + \kappa(\beta_S, 1 - \beta_R)) dG_S(\beta_S) dG_R(\beta_R) \\
&= \int_{\beta_R=0}^{0.5} \int_{\beta_S=0}^1 \lambda(\beta_S, \beta_R) dG_S(\beta_S) dG_R(\beta_R) \\
&= \int_{\beta_R=0}^{0.5} \left( \int_{\beta_S=0}^{0.5} \lambda(\beta_S, \beta_R) dG_S(\beta_S) + \int_{\beta_S=0}^{0.5} \lambda(1 - \beta_S, \beta_R) dG_S(1 - \beta_S) \right) dG_R(\beta_R) \\
&= \int_{\beta_R=0}^{0.5} \int_{\beta_S=0}^{0.5} (\lambda(\beta_S, \beta_R) + \lambda(1 - \beta_S, \beta_R)) dG_S(\beta_S) dG_R(\beta_R).
\end{aligned}$$

Let us now show that  $\lambda(\beta_S, \beta_R) + \lambda(1 - \beta_S, \beta_R) < 0$  for any  $\beta_S < 0.5$  and  $\beta_R < 0.5$  in which case the whole integral on the right-hand side is negative. Denote as before  $\bar{\beta} = \max\{\beta_S, \beta_R\}$  and  $\underline{\beta} = \min\{\beta_S, \beta_R\}$ . Then, (noting that  $1 - \underline{\beta} > 1 - \bar{\beta} > \bar{\beta} > \underline{\beta}$  due to both  $\bar{\beta} < 0.5$  and  $\underline{\beta} < 0.5$ )

$$\begin{aligned}
& \lambda(\beta_S, \beta_R) + \lambda(1 - \beta_S, \beta_R) \\
&= \kappa(\beta_S, \beta_R) + \kappa(\beta_S, 1 - \beta_R) + \kappa(1 - \beta_S, \beta_R) + \kappa(1 - \beta_S, 1 - \beta_R) \\
&= \left( \tilde{\beta}(0, \bar{\beta}) - \tilde{\beta}(0, \underline{\beta}) \right) - (\bar{\beta} - \underline{\beta}) \\
&\quad + \left( \tilde{\beta}(0, 1 - \bar{\beta}) - \tilde{\beta}(0, \underline{\beta}) \right) - (1 - \bar{\beta} - \underline{\beta}) \\
&\quad + \left( \tilde{\beta}(0, 1 - \underline{\beta}) - \tilde{\beta}(0, \bar{\beta}) \right) - (1 - \underline{\beta} - \bar{\beta}) \\
&\quad + \left( \tilde{\beta}(0, 1 - \underline{\beta}) - \tilde{\beta}(0, 1 - \bar{\beta}) \right) - (1 - \underline{\beta} - (1 - \bar{\beta})) \\
&= 2(\tilde{\beta}(0, 1 - \underline{\beta}) - \tilde{\beta}(0, \underline{\beta}) + 2\underline{\beta} - 1) \\
&= 2 \left( \frac{(1 - \underline{\beta})(1 - p)}{(1 - \underline{\beta})(1 - p) + \underline{\beta}p} - \frac{\underline{\beta}(1 - p)}{\underline{\beta}(1 - p) + (1 - \underline{\beta})p} + 2\underline{\beta} - 1 \right) \\
&= -\frac{2(1 - 2p)^2(1 - \underline{\beta})(1 - 2\underline{\beta})\underline{\beta}}{(1 - p + \underline{\beta}(2p - 1))(\underline{\beta} + p(1 - 2\underline{\beta}))} < 0,
\end{aligned}$$

where the inequality follows due to  $\underline{\beta} < 0.5$ .

**Step 2.** By symmetry considerations, the same property holds for 1-signals, i.e.  $E_S E_R[\Delta(1)] - E_S E_R[\Delta^{FD}(\emptyset)] < 0$ . Formally, the proof proceeds analogously redefining  $\kappa(\beta_i, \beta_j) = \left| \tilde{\beta}(1, \beta_i) - \tilde{\beta}(1, \beta_j) \right| - |\beta_i - \beta_j|$ . ■

**Proof of Proposition 5.b)**

In what follows, we assume without loss of generality that MLRP is satisfied as

$$\frac{\partial g_S(x)}{\partial x g_R(x)} > 0. \quad (37)$$

**Step 1.** Denote the difference in disagreement under disclosure and no disclosure in a putative FD-equilibrium as

$$\begin{aligned} \kappa_0(\beta_S, \beta_R) &= |\beta_S - \beta_R| - \left| \tilde{\beta}(0, \beta_S) - \tilde{\beta}(0, \beta_R) \right|, \\ \kappa_1(\beta_S, \beta_R) &= |\beta_S - \beta_R| - \left| \tilde{\beta}(1, \beta_S) - \tilde{\beta}(1, \beta_R) \right|. \end{aligned}$$

In FD, we have

$$\begin{aligned} \int_{\beta_R=0}^1 \int_{\beta_S=0}^1 \kappa_0(\beta_S, \beta_R) g_S(\beta_S) g_R(\beta_R) d\beta_S d\beta_R &\geq 0, \\ \int_{\beta_R=0}^1 \int_{\beta_S=0}^1 \kappa_1(\beta_S, \beta_R) g_S(\beta_S) g_R(\beta_R) d\beta_S d\beta_R &\geq 0. \end{aligned}$$

Since the joint distribution of priors is completely symmetric with respect to either boundary (0 or 1), the effect of 0-disclosure on the expected disagreement should be equivalent to the effect of 1-disclosure, i.e.

$$\begin{aligned} &\int_{\beta_R=0}^1 \int_{\beta_S=0}^1 \kappa_0(\beta_S, \beta_R) g_S(\beta_S) g_R(\beta_R) d\beta_S d\beta_R \\ &= \int_{\beta_R=0}^1 \int_{\beta_S=0}^1 \kappa_1(\beta_S, \beta_R) g_S(\beta_S) g_R(\beta_R) d\beta_S d\beta_R. \end{aligned}$$

This implies that for  $i = 0, 1$

$$\begin{aligned} \int_{\beta_R=0}^1 \int_{\beta_S=0}^1 \kappa_i(\beta_S, \beta_R) g_S(\beta_S) g_R(\beta_R) d\beta_S d\beta_R &\geq 0 \Leftrightarrow \\ \int_{\beta_R=0}^1 \int_{\beta_S=0}^1 \eta(\beta_S, \beta_R) g_S(\beta_S) g_R(\beta_R) d\beta_S d\beta_R &\geq 0. \end{aligned}$$

where  $\eta(\beta_S, \beta_R) = \kappa_0(\beta_S, \beta_R) + \kappa_1(\beta_S, \beta_R)$ .

**Step 2.** We have

$$\begin{aligned}
& \int_{\beta_R=0}^1 \int_{\beta_S=0}^1 \eta(\beta_S, \beta_R) g_S(\beta_S) g_R(\beta_R) d\beta_S d\beta_R \\
= & \int_{\beta_R=0}^{0.5} \int_{\beta_S=0}^1 \eta(\beta_S, \beta_R) g_S(\beta_S) g_R(\beta_R) d\beta_S d\beta_R \\
& + \int_{\beta_R=0}^{0.5} \int_{\beta_S=0}^1 \eta(\beta_S, 1 - \beta_R) g_S(\beta_S) g_R(1 - \beta_R) d\beta_S d\beta_R \\
= & \int_{\beta_R=0}^{0.5} \int_{\beta_S=0}^{0.5} \eta(\beta_S, \beta_R) g_S(\beta_S) g_R(\beta_R) d\beta_S d\beta_R \\
& + \int_{\beta_R=0}^{0.5} \int_{\beta_S=0}^{0.5} \eta(1 - \beta_S, \beta_R) g_S(1 - \beta_S) g_R(\beta_R) d\beta_S d\beta_R \\
& + \int_{\beta_R=0}^{0.5} \int_{\beta_S=0}^{0.5} \eta(\beta_S, 1 - \beta_R) g_S(\beta_S) g_R(1 - \beta_R) d\beta_S d\beta_R \\
& + \int_{\beta_R=0}^{0.5} \int_{\beta_S=0}^{0.5} \eta(1 - \beta_S, 1 - \beta_R) g_S(1 - \beta_S) g_R(1 - \beta_R) d\beta_S d\beta_R \\
= & \int_{\beta_R=0}^{0.5} \int_{\beta_S=0}^{0.5} \varsigma(\beta_S, \beta_R) d\beta_S d\beta_R,
\end{aligned}$$

where

$$\begin{aligned}
\varsigma(\beta_S, \beta_R) = & \eta(\beta_S, \beta_R) g_S(\beta_S) g_R(\beta_R) + \eta(1 - \beta_S, \beta_R) g_S(1 - \beta_S) g_R(\beta_R) \\
& + \eta(\beta_S, 1 - \beta_R) g_S(\beta_S) g_R(1 - \beta_R) + \eta(1 - \beta_S, 1 - \beta_R) g_S(1 - \beta_S) g_R(1 - \beta_R).
\end{aligned}$$

Hence, given Step 1, for the main claim it is sufficient to show that  $\varsigma(\beta_S, \beta_R) \geq 0$  for any  $\{\beta_S, \beta_R\} \in [0, 0.5]^2$ .

**Step 3.** Let us show that  $\varsigma(\beta_S, \beta_R)$  is increasing in  $p$  for  $p \in (1/2, 1)$  and  $\{\beta_S, \beta_R\} \in [0, 0.5]^2$ . To simplify the notation, let us denote  $g_S(\beta_S) \equiv g_{S1}$ ,  $g_R(\beta_R) \equiv g_{R1}$ ,  $g_S(1 - \beta_S) \equiv g_{S2}$ ,  $g_R(1 - \beta_R) \equiv g_{R2}$ .

Consider first  $0.5 \geq \beta_R > \beta_S$ . Substituting all expressions into  $\varsigma(\beta_S, \beta_R)$  and simplifying, we obtain

$$\begin{aligned}
\varsigma(\beta_S, \beta_R) = & \tau_1(\beta_R + \beta_S - 1)(g_{R2}g_{S1} + g_{R1}g_{S2}) \\
& + \tau_2(\beta_R - \beta_S)(g_{R1}g_{S1} + g_{R2}g_{S2}),
\end{aligned}$$

where

$$\begin{aligned}\tau_1 &= -2 - \frac{(1-p)p}{(\beta_R + p - 2\beta_{Rp})(\beta_S + p - 2\beta_{Sp} - 1)} \\ &\quad - \frac{(1-p)p}{(\beta_R + p - 2\beta_{Rp} - 1)(\beta_S + p - 2\beta_{Sp})}, \\ \tau_2 &= 2 - \frac{(1-p)p}{(\beta_R + p - 2\beta_{Rp} - 1)(\beta_S + p - 2\beta_{Sp} - 1)} \\ &\quad - \frac{(1-p)p}{(\beta_R + p - 2\beta_{Rp})(\beta_S + p - 2\beta_{Sp})}.\end{aligned}$$

Taking the derivative of  $\zeta(\beta_S, \beta_R)$  with respect to  $p$  and simplifying we obtain

$$\begin{aligned}\frac{\partial \zeta(\beta_S, \beta_R)}{\partial p} &= T_1 + T_2, \\ T_1 &= (1 - \beta_R)\beta_R \frac{(2p-1)(1-2\beta_R)}{(\beta_R + p - 2\beta_{Rp} - 1)^2(\beta_R + p - 2\beta_{Rp})^2} \\ &\quad \times (g_{R2} - g_{R1})(g_{S1} - g_{S2}), \\ T_2 &= (1 - \beta_S)\beta_S \frac{(2p-1)(1-2\beta_S)}{(\beta_S + p - 2\beta_{Sp} - 1)^2(\beta_S + p - 2\beta_{Sp})^2} \\ &\quad \times (g_{R2} + g_{R1})(g_{S1} + g_{S2}).\end{aligned}$$

Consider now the case  $\beta_R < \beta_S \leq 0.5$ . Substituting all expressions into  $\zeta(\beta_S, \beta_R)$  and simplifying, we obtain in this case

$$\begin{aligned}\zeta(\beta_S, \beta_R) &= \tau_1(\beta_R + \beta_S - 1)(g_{R2}g_{S1} + g_{R1}g_{S2}) \\ &\quad + \tau_2(\beta_S - \beta_R)(g_{R1}g_{S1} + g_{R2}g_{S2}),\end{aligned}$$

Taking the derivative of  $\zeta(\beta_S, \beta_R)$  with respect to  $p$  and simplifying we obtain

$$\frac{\partial \zeta(\beta_S, \beta_R)}{\partial p} = \widehat{T}_1 + \widehat{T}_2,$$

where

$$\begin{aligned}\widehat{T}_1 &= (1 - \beta_R)\beta_R \frac{(2p-1)(1-2\beta_R)}{(\beta_R + p - 2\beta_{Rp} - 1)^2(\beta_R + p - 2\beta_{Rp})^2} \\ &\quad \times (g_{R2} + g_{R1})(g_{S1} + g_{S2}), \\ \widehat{T}_2 &= (1 - \beta_R)\beta_R \frac{(2p-1)(1-2\beta_S)}{(\beta_S + p - 2\beta_{Sp} - 1)^2(\beta_S + p - 2\beta_{Sp})^2} \\ &\quad \times (g_{R2} - g_{R1})(g_{S1} - g_{S2}).\end{aligned}$$

Recall that  $\{\beta_S, \beta_R\} \in [0, 0.5]^2$  by assumption. Hence, to show that  $\frac{\partial \zeta(\beta_S, \beta_R)}{\partial p} \geq 0$  in

either case we need to show that

$$(g_{R2} - g_{R1})(g_{S1} - g_{S2}) > 0.$$

This is done in the next step.

**Step 4.** By initial assumption, we have that for any  $x$

$$g_R(x) = g_S(1 - x).$$

In particular, this implies

$$\frac{g_R(0.5)}{g_S(0.5)} = 1.$$

Note that then the MLRP in (37) implies that for any  $\beta_R < 0.5$  and  $\beta_S < 0.5$

$$\begin{aligned} \frac{g_S(\beta_S)}{g_R(\beta_S)} &< \frac{g_S(0.5)}{g_R(0.5)} < \frac{g_S(1 - \beta_R)}{g_R(1 - \beta_R)} \Leftrightarrow \\ \frac{g_S(\beta_S)}{g_R(\beta_S)} &< 1 < \frac{g_S(1 - \beta_R)}{g_R(1 - \beta_R)}. \end{aligned} \quad (38)$$

Since by initial assumption  $g_R(x) = g_S(1 - x)$ , (38) is equivalent to

$$\frac{g_S(\beta_S)}{g_S(1 - \beta_S)} < 1 < \frac{g_R(\beta_R)}{g_R(1 - \beta_R)}.$$

In terms of our previous notation, this is equivalent to

$$\begin{aligned} g_{S1} &< g_{S2}, \\ g_{R1} &> g_{R2}. \end{aligned}$$

Finally, this leads to

$$(g_{R2} - g_{R1})(g_{S1} - g_{S2}) > 0. \quad (39)$$

**Step 5.** Applying (39) to the expressions for  $\frac{\partial \varsigma(\beta_S, \beta_R)}{\partial p}$  from Step 3, we obtain

$$\frac{\partial \varsigma(\beta_S, \beta_R)}{\partial p} \geq 0.$$

At the same time, it is easy to verify that  $\varsigma(\beta_S, \beta_R) = 0$  for  $p = 1/2$ . Consequently,  $\varsigma(\beta_S, \beta_R) \geq 0$  for any  $p > 1/2$ . Then, by Step 2 this results in

$$\int_{\beta_R=0}^1 \int_{\beta_S=0}^1 \eta(\beta_S, \beta_R) g_S(\beta_S) g_R(\beta_R) d\beta_S d\beta_R \geq 0.$$

By Step 1, this implies that the incentive constraints for full disclosure are satisfied. ■

**Proof of Proposition 5.c)**

**Step 1.** Let us show that for sufficiently high  $x$ , it holds that  $\beta_S > \beta_S^{**}(\beta_R) = \frac{p(1-\beta_R)}{\beta_R+p(1-2\beta_R)}$  for any  $\{\beta_S, \beta_R\} \in [x, 1]^2$ . Indeed, it is easy to verify that  $x > \beta_S^{**}(x)$  if and only if  $x > \frac{p}{p+\sqrt{p(1-p)}}$ . Thus, we have that for  $x > \frac{p}{p+\sqrt{p(1-p)}}$  and any  $\{\beta_S, \beta_R\} \in [x, 1]^2$  it holds

$$\beta_S \geq x > \beta_S^{**}(x) \geq \beta_S^{**}(\beta_R),$$

where the last inequality is due to  $\beta_S^{**}(x)$  decreasing in  $x$ . Hence,  $\beta_S > \beta_S^{**}(\beta_R)$  for any  $\{\beta_S, \beta_R\} \in [x, 1]^2$ .

Analogously, one can show that for any sufficiently small  $y$  (in particular, for any  $y < \frac{p+\sqrt{p(1-p)-1}}{2p-1}$ ), it holds  $\beta_S < \beta_S^*(\beta_R)$  for any  $\{\beta_S, \beta_R\} \in [0, y]^2$ .

**Step 2.** Let us show that if the common distribution of priors  $g$  is shifted to the right, then D1-equilibrium always exists. The incentive constraints for D1 are (see Step 1 in the proof of Proposition 5.a)

$$\int_{\beta_R=0}^1 \int_{\beta_S=0}^1 \tilde{\kappa}_0(\beta_S, \beta_R) g(\beta_S) g(\beta_R) d\beta_S d\beta_R \leq 0, \quad (40)$$

$$\int_{\beta_R=0}^1 \int_{\beta_S=0}^1 \tilde{\kappa}_1(\beta_S, \beta_R) g(\beta_S) g(\beta_R) d\beta_S d\beta_R \geq 0, \quad (41)$$

where

$$\begin{aligned} \tilde{\kappa}_0(\beta_S, \beta_R) &= \Delta^{D1}(\emptyset; \beta_S, \beta_R) - \Delta(0; \beta_S, \beta_R), \\ \tilde{\kappa}_1(\beta_S, \beta_R) &= \Delta^{D1}(\emptyset; \beta_S, \beta_R) - \Delta(1; \beta_S, \beta_R). \end{aligned}$$

At the same time, for any constellation  $\{\beta_S, \beta_R\} \in [x, 1]^2$  and  $x$  sufficiently high we have  $\beta_S > \beta_S^{**}(\beta_R)$  by Step 1, which then implies by Proposition 1

$$\begin{aligned} \tilde{\kappa}_0(\beta_S, \beta_R) &\leq 0, \\ \tilde{\kappa}_1(\beta_S, \beta_R) &\geq 0. \end{aligned}$$

Consequently,

$$\int_{\beta_R=x}^1 \int_{\beta_S=x}^1 \kappa_0(\beta_S, \beta_R) g_S(\beta_S) g_R(\beta_R) d\beta_S d\beta_R \leq 0, \quad (42)$$

$$\int_{\beta_R=x}^1 \int_{\beta_S=x}^1 \kappa_1(\beta_S, \beta_R) g_S(\beta_S) g_R(\beta_R) d\beta_S d\beta_R \geq 0. \quad (43)$$

Finally, (42) and (43) result in (40) and (41) as far as  $g$  is sufficiently skewed to the right.

**Step 3.** The non-existence of other pure strategy equilibria (besides D1) if  $g$  is sufficiently shifted to the right follows by the analogous argument. In particular, by Step 1 and

Proposition 1 for any given constellation  $\{\beta_S, \beta_R\} \in [x, 1]^2$  with  $x$  sufficiently high the  $S$ 's incentive constraints for other equilibria (D0 and FD) are not satisfied. Consequently, they are still not satisfied once we integrate them over all possible constellations  $\{\beta_S, \beta_R\} \in [x, 1]^2$  like in Step 2. If the probability mass set on  $\{\beta_S, \beta_R\} \notin [x, 1]^2$  gets sufficiently small, the same applies to the integration over all possible constellations  $\{\beta_S, \beta_R\} \in [0, 1]^2$ .

**Step 4.** Consider the case when the distribution  $g$  is sufficiently skewed to the left, i.e. to values  $[0, y]^2$ . As before, Step 1 implies that for any given  $\{\beta_S, \beta_R\} \in [0, y]^2$  we have  $\beta_S < \beta_S^*(\beta_R)$ , i.e. the  $S$ 's incentive constraints for D0 are satisfied, while for D1 and FD they are not satisfied. Consequently, the same holds once we integrate them over all possible priors constellations in  $[0, y]^2$ , and hence in  $[0, 1]^2$  (under sufficiently skewed distribution).

### Proof of Proposition 5.d)

Suppose that  $S$ 's prior  $\beta_S$  is commonly known. That of  $R$  is drawn from a symmetric distribution  $G$  over  $[0, 1]$ . Then, by the same steps as in the proof of Proposition 5.a we obtain

$$\begin{aligned} E_S[\Delta(0)] - E_S[\Delta^{FD}(\emptyset)] &= \int_{\beta_R=0}^1 \left( \left| \tilde{\beta}_S(0) - \tilde{\beta}_R(0) \right| - |\beta_S - \beta_R| \right) dG_R(\beta_R) \\ &= \int_{\beta_R=0}^{0.5} (\kappa_0(\beta_S, \beta_R) + \kappa_0(\beta_S, 1 - \beta_R)) dG_R(\beta_R). \end{aligned} \quad (44)$$

Consider  $\beta_R < 0.5$  such that  $1 - \beta_R > \beta_S > \beta_R$ . For such  $\beta_R$  it holds

$$\begin{aligned} &\kappa_0(\beta_S, \beta_R) + \kappa_0(\beta_S, 1 - \beta_R) \\ &= \left( \tilde{\beta}(0, \beta_S) - \tilde{\beta}(0, \beta_R) \right) - (\beta_S - \beta_R) \\ &\quad + \left( \tilde{\beta}(0, 1 - \beta_R) - \tilde{\beta}(0, \beta_S) \right) - (1 - \beta_R - \beta_S) \\ &= \tilde{\beta}(0, 1 - \beta_R) - \tilde{\beta}(0, \beta_R) + 2\beta_R - 1 \\ &= \frac{(1 - \beta_R)(1 - p)}{(1 - \beta_R)(1 - p) + \beta_R p} - \frac{\beta_R(1 - p)}{\beta_R(1 - p) + (1 - \beta_R)p} + 2\beta_R - 1 \\ &= -\frac{(1 - 2p)^2(1 - \beta_R)(1 - 2\beta_R)\beta_R}{(1 - p + \beta_R(2p - 1))(\beta_R + p(1 - 2\beta_R))} < 0. \end{aligned}$$

Since the probability mass of  $\beta_R < 0.5$  such that the condition  $1 - \beta_R > \beta_S > \beta_R$  is satisfied is sufficiently large for  $\beta_S$  sufficiently close to 0.5, the right-hand side of (44) is negative as well. Hence, the sender would prefer to disclose 0-signal over no disclosure. The same claim for 1-signals follows by symmetry considerations. Consequently, the FD-equilibrium exists. ■

## Appendix V: Proposition 2

**Step 1.** First, note that  $\Delta(0, \beta_S, \beta_R)$  and  $\Delta(1, \beta_S, \beta_R)$  are V-shaped with respect to either  $\beta_S$  or  $\beta_R$  reaching its minimum at  $\beta_S = \beta_R$ . Indeed, since  $\tilde{\beta}_i(0)$  is increasing in  $\beta_i$ , it follows that  $\Delta(0, \beta_i, \beta_j)$  decreases in  $\beta_i$  if  $\beta_i < \beta_j$  and increases in  $\beta_i$  otherwise, being equal to 0 for  $\beta_i = \beta_j$ . The same argument applies to  $\Delta(1, \beta_S, \beta_R)$ .

**Step 2.** Let us show another auxiliary result that  $E_S[\Delta^{D0}(\emptyset, \beta_S, \beta_R)]$  and  $E_S[\Delta^{D1}(\emptyset, \beta_S, \beta_R)]$  are V-shaped with respect to  $\beta_R$  reaching its minimum at  $\beta_S = \beta_R$ . Consider  $E_S[\Delta^{D1}(\emptyset, \beta_S, \beta_R)]$ . Using the expressions from Appendix I, we get:

$$\begin{aligned} & E_S[\Delta^{D1}(\emptyset, \beta_S, \beta_R)] \\ = & \frac{p(1-p\varphi) - \beta_S(2p-1)(1-\varphi)}{p - \beta_S(2p-1)} \frac{|\beta_S - \beta_R|}{1 - \varphi(1-p + \beta_R(2p-1))}. \end{aligned} \quad (45)$$

Taking the derivative with respect to  $\beta_R$  and simplifying we obtain (for  $\beta_R \neq \beta_S$ )

$$\begin{aligned} & \frac{\partial E_S[\Delta^{D1}(\emptyset, \beta_S, \beta_R)]}{\partial \beta_R} \\ = & \text{sgn}[\beta_R - \beta_S] \frac{p(1-p\varphi) - \beta_S(2p-1)(1-\varphi)}{p - \beta_S(2p-1)} \frac{1 - \varphi(1-p + \beta_S(2p-1))}{(1 - \varphi(1-p + \beta_R(2p-1)))^2} \end{aligned}$$

It is easy to verify that all terms on the right-hand side following the sign function are always positive. Hence, the sign of the derivative is determined by  $\text{sgn}[\beta_R - \beta_S]$ , which implies that function  $E_S[\Delta^{D1}(\emptyset, \beta_S, \beta_R)]$  is V-shaped with respect to  $\beta_R$ , being kinked at  $\beta_S = \beta_R$  where it is equal to 0 (see 45).

Consider  $E_S[\Delta^{D0}(\emptyset, \beta_S, \beta_R)]$ . Using the expressions from Appendix I, we get:

$$\begin{aligned} & E_S[\Delta^{D0}(\emptyset, \beta_S, \beta_R)] \\ = & \frac{1-p + \beta_S(2p-1)(1-\varphi) - (1-p)^2\varphi}{1-p + \beta_S(2p-1)} \frac{|\beta_S - \beta_R|}{1 - \varphi(p - \beta_R(2p-1))}. \end{aligned} \quad (46)$$

Taking the derivative with respect to  $\beta_R$  and simplifying we obtain (for  $\beta_R \neq \beta_S$ )

$$\begin{aligned} & \frac{\partial E_S[\Delta^{D0}(\emptyset, \beta_S, \beta_R)]}{\partial \beta_R} \\ = & \text{sgn}[\beta_R - \beta_S] \frac{1-p + \beta_S(2p-1)(1-\varphi) - (1-p)^2\varphi}{1-p + \beta_S(2p-1)} \frac{1 - \varphi(p - \beta_S(2p-1))}{(1 - \varphi(p - \beta_R(2p-1)))^2} \end{aligned}$$

It is easy to verify that all terms on the right-hand side following the sign function are always positive. Hence, the sign of the derivative is determined by  $\text{sgn}[\beta_R - \beta_S]$ , which again implies that function  $E_S[\Delta^{D0}(\emptyset, \beta_S, \beta_R)]$  is V-shaped with respect to  $\beta_R$ , being kinked at  $\beta_S = \beta_R$  where it is equal to 0 (see 46).

**Step 3.** Let us show that  $E_S[\Delta^{FD}]$ ,  $E_S[\Delta^{D0}]$ , and  $E_S[\Delta^{D1}]$  are all V-shaped with respect



to  $\beta_R$  and reach their minimum at  $\beta_S = \beta_R$ . We have

$$\begin{aligned} E_S[\Delta^{FD}] &= \Pr[\sigma_S = 1|\beta_S]\Delta(1, \beta_S, \beta_R) \\ &+ \Pr[\sigma_S = 0|\beta_S]\Delta(0, \beta_S, \beta_R) + \Pr[\sigma_S = \emptyset|\beta_S]|\beta_S - \beta_R|, \end{aligned} \quad (47)$$

$$\begin{aligned} E_S[\Delta^{D0}] &= \Pr[\sigma_S = 0|\beta_S]\Delta(0, \beta_S, \beta_R) \\ &+ (1 - \Pr[\sigma_S = 0|\beta_S])\Delta^{D0}(\emptyset, \beta_S, \beta_R), \end{aligned} \quad (48)$$

$$\begin{aligned} E_S[\Delta^{D1}] &= \Pr[\sigma_S = 1|\beta_S]\Delta(1, \beta_S, \beta_R) \\ &+ (1 - \Pr[\sigma_S = 1|\beta_S])\Delta^{D1}(\emptyset, \beta_S, \beta_R). \end{aligned} \quad (49)$$

Note now that by Steps 1 and 2, it holds true that  $\Delta(0, \beta_S, \beta_R)$ ,  $\Delta(1, \beta_S, \beta_R)$ ,  $\Delta^{D0}(\emptyset, \beta_S, \beta_R)$ ,  $\Delta^{D1}(\emptyset, \beta_S, \beta_R)$  and  $|\beta_S - \beta_R|$  are all V-shaped and reach their minimum (which equals 0) for  $\beta_S = \beta_R$ . It follows immediately that  $E_S[\Delta^{FD}]$ ,  $E_S[\Delta^{D0}]$  and  $E_S[\Delta^{D1}]$  exhibit these same properties.

**Step 4.** Let us show that  $E_S[\Delta]$  is uniquely defined, i.e. that if  $X$  and  $X'$  are two equilibrium disclosure rules given  $\beta_S, \beta_R$ , then  $E_S[\Delta^X|\beta_S, \beta_R] = E_S[\Delta^{X'}|\beta_S, \beta_R]$ . First note that  $E_S[\Delta] = 0$  in any equilibrium if  $\beta_S = \beta_R$  (see the proof of Lemma II.C). Consider  $\beta_S \neq \beta_R$ . Then, by Proposition 1 the only instances where the equilibrium is not unique are when  $\beta_S = \beta_S^*(\beta_R)$  and  $\beta_S = \beta_S^{**}(\beta_R)$ . Consider  $\beta_S = \beta_S^*(\beta_R)$ , in which case by Proposition 1 there exist FD, D0 and mixed disclosure equilibria. By (21)-(25) in the proof of Lemma II.G, if  $\beta_S = \beta_S^*(\beta_R)$ , then  $S$  must be indifferent between disclosing  $\sigma = 1$  and non-disclosure for any  $\delta_R^1$ , i.e. in any disclosure equilibrium in which  $\sigma = 1$  is not disclosed with positive probability (this includes D0). I.e. it must be true that:

$$\Delta(1, \beta_S^*(\beta_R), \beta_R) = \Delta^{D0}(\emptyset, \beta_S^*(\beta_R), \beta_R). \quad (50)$$

Note that by (4)  $\Delta^{D0}(\emptyset, \beta_S^*(\beta_R), \beta_R)$  is a weighted average between  $\Delta(1, \beta_S^*(\beta_R), \beta_R)$  and  $|\beta_S^*(\beta_R) - \beta_R|$ . Together with (50), this implies

$$\Delta(1, \beta_S^*(\beta_R), \beta_R) = |\beta_S^*(\beta_R) - \beta_R|. \quad (51)$$

(47), (48), (50) and (51) jointly imply that  $E_S[\Delta^{FD}|\beta_S = \beta_S^*(\beta_R)]$  is equal to  $E_S[\Delta^{D0}|\beta_S = \beta_S^*(\beta_R)]$ , as well as to the corresponding value under any other equilibrium involving randomization between disclosure and non-disclosure when  $\sigma = 1$  (recall that this is the only possible mixed-disclosure strategy equilibrium if  $\beta_S = \beta_S^*(\beta_R)$  by the proof of Lemma II.G). Consequently,  $E_S[\Delta]$  is uniquely defined if  $\beta_S = \beta_S^*(\beta_R)$ . By an analogous argument,  $E_S[\Delta]$  is uniquely defined if  $\beta_S = \beta_S^{**}(\beta_R)$ .

**Step 5.** Let us show that  $E_S[\Delta]$  is continuous in  $\beta_R$ . Note first that  $E_S[\Delta^{D0}]$ ,  $E_S[\Delta^{FD}]$  and  $E_S[\Delta^{D1}]$  are all continuous in  $\beta_R$ . Besides, by Step 4,  $E_S[\Delta]$  is uniquely defined. This together with Proposition 1 implies that  $E_S[\Delta]$  is equal either to  $E_S[\Delta^{D0}]$ ,  $E_S[\Delta^{FD}]$  or  $E_S[\Delta^{D1}]$  depending on whether, respectively,  $\beta_S \in (0, \beta_S^*(\beta_R)]$ ,  $\beta_S \in [\beta_S^*(\beta_R), \beta_S^{**}(\beta_R)]$  and

$\beta_S \in [\beta_S^{**}(\beta_R), 1)$ , being continuous at  $\beta_R = (\beta_S^*)^{-1}(\beta_S)$  and  $\beta_R = (\beta_S^{**})^{-1}(\beta_S)$ . Consequently,  $E_S[\Delta]$  is also continuous in  $\beta_R$ .

**Step 6.** Consider finally the perceived disagreement from  $S$ 's perspective. By Step 5,  $E_S[\Delta]$  is continuous in  $\beta_R$  and is equal either to  $E_S[\Delta^{D0}]$ ,  $E_S[\Delta^{FD}]$  or  $E_S[\Delta^{D1}]$ . By Step 3, all these functions are V-shaped with respect to  $\beta_R$  reaching its minimum at  $\beta_S = \beta_R$ . Consequently, the same holds for  $E_S[\Delta]$ . ■

## Appendix VI: Proposition 6

Proposition 6 follows from a set of Lemmas (Lemmas V.A to V.D), which are stated and proved in what follows. In a given SDE featuring the non-disclosure interval  $(s_1, s_2)$ , we denote  $R$ 's perceived disagreement conditional on disclosure of a signal  $s$  by  $\Delta(s)$  and conditional on non-disclosure by  $\Delta^{(s_1, s_2)}(\emptyset)$ :

$$\begin{aligned}\Delta(s) &= \left| \tilde{\beta}_S(s) - \tilde{\beta}_R(s) \right| \text{ for } s \in [\underline{s}, \bar{s}], \\ \Delta^{(s_1, s_2)}(\emptyset) &= \left| E_R^{s_1, s_2}[\tilde{\beta}_S | \emptyset] - \tilde{\beta}_R^{s_1, s_2}(\emptyset) \right|.\end{aligned}$$

**Lemma V.A** *If  $\beta_S \neq \beta_R$ , then  $\Delta(s)$  satisfies the following:*

- i)  $\lim_{s \rightarrow \underline{s}} \Delta(s) = \lim_{s \rightarrow \bar{s}} \Delta(s) = 0$ .*
- ii) There exists  $\hat{s}$  such that  $\Delta(s)$  is increasing in  $s$  for all  $s < \hat{s}$  and decreasing in  $s$  for all  $s > \hat{s}$ .*
- iii)  $\tilde{s} > (<) \hat{s}$  if and only if the player with the lower prior is less (more) confident. Instead,  $\tilde{s} = \hat{s}$  if and only if  $\beta_S = 1 - \beta_R$ , i.e. if players are equally confident.*

**Proof.**

**Step 1.** i) is immediate. To show ii) we first prove that there is a unique  $\hat{s}$  such that

$$\frac{d}{ds} \left( \tilde{\beta}_S(\hat{s}) - \tilde{\beta}_R(\hat{s}) \right) = 0$$

Indeed,

$$\begin{aligned}& \frac{d}{ds} \left( \tilde{\beta}_S(s) - \tilde{\beta}_R(s) \right) \\ &= \frac{d}{ds} \left( \frac{\beta_S}{\beta_S + (1 - \beta_S) \frac{f(s|0)}{f(s|1)}} - \frac{\beta_R}{\beta_R + (1 - \beta_R) \frac{f(s|0)}{f(s|1)}} \right) \\ &= \frac{f(s|0)}{f(s|1)} \left( \frac{\beta_R(1 - \beta_R)}{\left( \beta_R + (1 - \beta_R) \frac{f(s|0)}{f(s|1)} \right)^2} - \frac{\beta_S(1 - \beta_S)}{\left( \beta_S + (1 - \beta_S) \frac{f(s|0)}{f(s|1)} \right)^2} \right) \frac{d}{ds}.\end{aligned}\tag{52}$$

Consider the solution to

$$\beta_R(1 - \beta_R) \left( \beta_S + (1 - \beta_S) \frac{f(s|0)}{f(s|1)} \right)^2 = \beta_S(1 - \beta_S) \left( \beta_R + (1 - \beta_R) \frac{f(s|0)}{f(s|1)} \right)^2.$$

Both sides are decreasing in  $s$ , but we claim that they increase at different rates. To see this, note that

$$\begin{aligned} & \frac{d}{ds} \beta_R(1 - \beta_R) \left( \beta_S + (1 - \beta_S) \frac{f(s|0)}{f(s|1)} \right)^2 \\ &= 2\beta_R(1 - \beta_R)(1 - \beta_S) \left( \beta_S + (1 - \beta_S) \frac{f(s|0)}{f(s|1)} \right) \frac{d}{ds} \frac{f(s|0)}{f(s|1)}, \\ & \frac{d}{ds} \beta_S(1 - \beta_S) \left( \beta_R + (1 - \beta_R) \frac{f(s|0)}{f(s|1)} \right)^2 \\ &= 2\beta_S(1 - \beta_R)(1 - \beta_S) \left( \beta_R + (1 - \beta_R) \frac{f(s|0)}{f(s|1)} \right) \frac{d}{ds} \frac{f(s|0)}{f(s|1)}. \end{aligned}$$

The result then follows from the fact that

$$\begin{aligned} & 2\beta_R(1 - \beta_R)(1 - \beta_S) \left( \beta_S + (1 - \beta_S) \frac{f(s|0)}{f(s|1)} \right) \frac{d}{ds} \frac{f(s|0)}{f(s|1)} \\ & \geq 2\beta_S(1 - \beta_R)(1 - \beta_S) \left( \beta_R + (1 - \beta_R) \frac{f(s|0)}{f(s|1)} \right) \frac{d}{ds} \frac{f(s|0)}{f(s|1)} \end{aligned}$$

is equivalent to

$$\beta_R\beta_S + \beta_R(1 - \beta_S) \frac{f(s|0)}{f(s|1)} \leq \beta_R\beta_S + \beta_S(1 - \beta_R) \frac{f(s|0)}{f(s|1)}$$

which, in turn, is equivalent to  $\beta_R \leq \beta_S$ . Hence,  $\hat{s}$  (where  $\Delta(s)$  reaches its extremum) must be unique. Then, claim ii) follows from continuity and (i) together with  $\Delta(\tilde{s}) = |\beta_S - \beta_R| > 0$ .

**Step 2.** To show (iii), define again  $\bar{\beta} = \max\{\beta_S, \beta_R\}$  and  $\underline{\beta} = \min\{\beta_S, \beta_R\}$  such that  $\Delta(s) = \tilde{\beta}(s, \bar{\beta}) - \tilde{\beta}(s, \underline{\beta})$ . From (52) we then have:

$$\begin{aligned} \frac{d}{ds} \Delta(\tilde{s}) &= (\underline{\beta}(1 - \underline{\beta}) - \bar{\beta}(1 - \bar{\beta})) \frac{d}{ds} \frac{f(\tilde{s}|0)}{f(\tilde{s}|1)} \geq 0 \\ &\iff \underline{\beta}(1 - \underline{\beta}) \leq \bar{\beta}(1 - \bar{\beta}), \end{aligned}$$

so that by claim ii)  $\tilde{s} > (<) \hat{s}$  if and only if  $\underline{\beta}$  is less (more) confident than  $\bar{\beta}$ , and  $\tilde{s} = \hat{s}$  if and only if players are equally confident. ■

**Lemma V.B** (i) If  $\beta_S = \{\beta_R, 1 - \beta_R\}$ , then there exists an FD-equilibrium.

(ii) If  $\beta_S \neq \{\beta_R, 1 - \beta_R\}$ , then in any equilibrium a positive measure of signals is not disclosed.

**Proof.**

**Step 1.** Let us show the existence of FD for  $\beta_S = \{\beta_R, 1 - \beta_R\}$ . If  $\beta_S = \beta_R$  then trivially  $\Delta(s) = 0$  for any  $s$ , so that  $S$  has always an incentive to disclose  $s$ . If  $\beta_S = 1 - \beta_R$ , then  $\tilde{s} = \hat{s}$  by Lemma V.A(iii). Consequently, for any  $s \in [\underline{s}, \bar{s}]$  we obtain

$$\Delta^{FD}(\emptyset) = |\beta_S - \beta_R| = \Delta(\tilde{s}) = \Delta(\hat{s}) \geq \Delta(s),$$

where the last inequality is by Lemma V.A (ii). Hence,  $S$  has an incentive to disclose all signals in equilibrium.

**Step 2.** Let us show that for  $\beta_S \neq \{\beta_R, 1 - \beta_R\}$  there exists no equilibrium where the set of non-disclosed signals has 0-measure. Assume by contradiction that this is the case. Consider thus a putative equilibrium featuring a disclosure rule  $\tilde{D}$  such that the set of non-disclosed signals has 0-measure. Then, the perceived disagreement upon non-disclosure is  $\Delta^{\tilde{D}}(\emptyset) = |\beta_S - \beta_R|$ , since  $R$  assigns probability 1 to the fact that  $S$  is uninformed. At the same time, since  $\Delta(s)$  is single peaked at  $\hat{s}$  by Lemma V.A(ii) and  $\tilde{s} \neq \hat{s}$  by Lemma V.A(iii), we have

$$\Delta(\hat{s}) > \Delta(\tilde{s}) = |\beta_S - \beta_R| = \Delta^{\tilde{D}}(\emptyset),$$

so that  $S$  has an incentive not to disclose all signals located sufficiently close to  $\hat{s}$ , which is a contradiction.

**Lemma V.C** *If  $\beta_S \neq \{\beta_R, 1 - \beta_R\}$ , then the unique equilibrium is an SDE.*

**Proof.**

**Step 0.** Steps 1-2 introduce key equilibrium conditions. In steps 3-4, we show that there exists a unique SDE. Step 5 proves that any equilibrium is an SDE.

In what follows, we assume  $\beta_S > \beta_R$ . The proof for the reverse case follows the same steps and is omitted.

**Step 1.** Consider a putative simple disclosure equilibrium with non-disclosure interval  $(s_1, s_2)$ . From  $R$ 's point of view,  $S$  does not disclose an observed signal with probability

$$\Pr_R(s \in (s_1, s_2)) = \beta_R \int_{s_1}^{s_2} f(s|1)ds + (1 - \beta_R) \int_{s_1}^{s_2} f(s|0)ds.$$

When  $S$  does not disclose,  $R$ 's posterior is

$$\begin{aligned} & \tilde{\beta}_R^{s_1, s_2}(\emptyset) \\ = & \frac{\varphi}{(1 - \varphi) + \varphi \Pr_R(s \in (s_1, s_2))} \int_{s_1}^{s_2} (\beta_R f(s|1) + (1 - \beta_R) f(s|0)) \tilde{\beta}_R(s) ds \\ & + \frac{(1 - \varphi)}{(1 - \varphi) + \varphi \Pr_R(s \in (s_1, s_2))} \beta_R \end{aligned}$$

Similarly,  $R$ 's belief about  $S$ 's posterior in this case is

$$\begin{aligned} & E_R^{s_1, s_2}[\tilde{\beta}_S | \emptyset] \\ = & \frac{\varphi}{(1 - \varphi) + \varphi \Pr_R(s \in (s_1, s_2))} \int_{s_1}^{s_2} (\beta_R f(s|1) + (1 - \beta_R) f(s|0)) \tilde{\beta}_S(s) ds \\ & + \frac{(1 - \varphi)}{(1 - \varphi) + \varphi \Pr_R(s \in (s_1, s_2))} \beta_S. \end{aligned}$$

**Step 2.** Given the definition of SDE and the fact that  $\Delta(s)$  is single peaked,  $S$  must be indifferent between disclosure and non-disclosure at  $s_1$  and  $s_2$ . Hence, we require

$$\Delta(s) = \Delta^{(s_1, s_2)}(\emptyset) \text{ for } s = s_1, s_2. \quad (53)$$

Next, implicitly define  $s_2^*(s_1)$  as a value of  $s_2 \neq s_1$  equalizing

$$\Delta(s_1) = \Delta(s_2)$$

for given  $s_1 < \hat{s}$ . There is a unique such value by Lemma V.A (ii). We additionally (abusively) define  $s_2^*(\hat{s}) = \hat{s}$ . Then, the equilibrium condition (53) holds if and only if

$$\Delta(s_1) = \Delta^{(s_1, s_2^*(s_1))}(\emptyset).$$

Define the function

$$\gamma(s_1) \equiv \Delta^{(s_1, s_2^*(s_1))}(\emptyset) - \Delta(s_1).$$

It follows that the SDE featuring the non-disclosure interval  $(s_1, s_2^*(s_1))$  exists if and only if

$$\gamma(s_1) = 0. \quad (54)$$

In the next steps, we show that there always exists a unique value of  $s_1$  satisfying the above condition, which implies that there always exists a unique SDE.

**Step 3.** This step proves existence of an SDE, i.e. show that there exists  $s_1$  such that  $\gamma(s_1) = 0$ . Denote  $s'_1$  the smallest value of  $s$  such that  $\Delta(s) = \Delta(\hat{s}) = \beta_S - \beta_R$ . Note that  $s'_1 < \hat{s}$  given  $\beta_S \neq \{\beta_R, 1 - \beta_R\}$ . Indeed, we know from Lemma V.B. that in this case  $\Delta(\hat{s}) > \beta_S - \beta_R$ , and we also know from Lemma V.A that  $\Delta(s)$  is single peaked in  $s$  with a maximum at  $\hat{s}$  and with  $\lim_{s \rightarrow \underline{s}} \Delta(s) = \lim_{s \rightarrow \bar{s}} \Delta(s) = 0$ .

Let us prove that  $\gamma(s'_1) > 0$ . By Step 1 we have

$$\begin{aligned}
\Delta^{(s'_1, s_2^*(s'_1))}(\emptyset) &= \frac{\varphi}{(1-\varphi) + \varphi \Pr_R[s \in (s'_1, s_2^*(s'_1))]} \int_{s'_1}^{s_2^*(s'_1)} (\beta_S(s) - \beta_R(s)) \tilde{f}(s) ds \\
&\quad + \left( 1 - \frac{\varphi \Pr_R[s \in (s'_1, s_2^*(s'_1))]}{(1-\varphi) + \varphi \Pr_R[s \in (s'_1, s_2^*(s'_1))]} \right) (\beta_S - \beta_R) \\
&= \frac{\varphi}{(1-\varphi) + \varphi \Pr_R[s \in (s'_1, s_2^*(s'_1))]} \int_{s'_1}^{s_2^*(s'_1)} (\beta_S(s) - \beta_R(s)) \tilde{f}(s) ds \\
&\quad + \left( 1 - \frac{\varphi \Pr_R[s \in (s'_1, s_2^*(s'_1))]}{(1-\varphi) + \varphi \Pr_R[s \in (s'_1, s_2^*(s'_1))]} \right) \Delta(s'_1) \\
&> \Delta(s'_1)
\end{aligned} \tag{55}$$

where the second equality is by construction of  $s'_1$ , and the strict inequality follows from the fact that  $\beta_S(s) - \beta_R(s) > \Delta(s'_1)$  for all  $s \in (s'_1, s_2^*(s'_1))$  since  $s'_1 < \hat{s}$  as noted above. This implies that  $\gamma(s'_1) = \Delta^{(s'_1, s_2^*(s'_1))}(\emptyset) - \Delta(s'_1) > 0$ .

Now let us show that  $\gamma(\hat{s}) < 0$ . Since  $\hat{s} = s_2^*(\hat{s})$  by construction, it holds that  $\Pr_R[s \in (\hat{s}, s_2^*(\hat{s}))] = 0$  so that

$$\Delta^{(\hat{s}, s_2^*(\hat{s}))}(\emptyset) = \beta_S - \beta_R = \Delta(s'_1) < \Delta(\hat{s}), \tag{56}$$

where the inequality is due to  $s'_1 < \hat{s}$ .

Thus, we have shown that  $\gamma(s'_1) > 0$  and  $\gamma(\hat{s}) < 0$ . From continuity of  $\gamma(s)$  on  $[s'_1, \hat{s}]$ , it then follows that there exists at least one  $s_1 \in (s'_1, \hat{s})$  such that  $\gamma(s_1) = 0$ .

**Step 4.** We now show that there exists a unique SDE. By Step 1 we have

$$\begin{aligned}
&\Delta^{(s_1, s_2^*(s_1))}(\emptyset) \\
&= \frac{\varphi}{(1-\varphi) + \varphi \Pr_R(s \in (s_1, s_2^*(s_1)))} \int_{s_1}^{s_2^*(s_1)} (\beta_S(s) - \beta_R(s)) \tilde{f}(s) ds \\
&\quad + \left( 1 - \frac{\varphi \Pr_R(s \in (s_1, s_2^*(s_1)))}{(1-\varphi) + \varphi \Pr_R(s \in (s_1, s_2^*(s_1)))} \right) (\beta_S - \beta_R) \\
&= \frac{\varphi}{(1-\varphi) + \varphi \Pr_R(s \in (s_1, s_2^*(s_1)))} \left( \begin{array}{l} \int_{s_1}^{s_2^*(s_1)} (\beta_S(s) - \beta_R(s)) \tilde{f}(s) ds \\ - \Pr_R[s \in (s_1, s_2^*(s_1))] (\beta_S - \beta_R) \end{array} \right) \\
&\quad + (\beta_S - \beta_R).
\end{aligned}$$

Denote  $\eta(s_1) = \frac{\varphi}{(1-\varphi) + \varphi \Pr_R[s \in (s_1, s_2^*(s_1))]}$  so that

$$\begin{aligned} \eta'(s_1) &= \frac{\partial \eta(s_1)}{\partial s_1} = - \left( \frac{\varphi}{(1-\varphi) + \varphi \Pr_R[s \in (s_1, s_2^*(s_1))]} \right)^2 \\ &\quad \times \left( \frac{\partial \Pr_R[s \in (s_1, s_2^*(s_1))]}{\partial s_1} + \frac{\partial \Pr_R[s \in (s_1, s_2^*(s_1))]}{\partial s_2^*} \frac{\partial s_2^*(s_1)}{\partial s_1} \right) \\ &= [\eta(s_1)]^2 \left( \tilde{f}(s_1) - \tilde{f}(s_2^*) \frac{\partial s_2^*}{\partial s_1} \right) > 0. \end{aligned}$$

Then, taking the derivative of  $\Delta^{(s_1, s_2^*(s_1))}(\emptyset)$  with respect to  $s_1$  we obtain

$$\begin{aligned} &\frac{\partial \Delta^{(s_1, s_2^*(s_1))}(\emptyset)}{\partial s_1} \\ &= \eta'(s_1) \left( \int_{s_1}^{s_2^*} (\beta_S(s) - \beta_R(s)) \tilde{f}(s) ds - \Pr_R[s \in (s_1, s_2^*(s_1))] (\beta_S - \beta_R) \right) \\ &\quad + \eta(s_1) (-(\beta_S(s_1) - \beta_R(s_1)) \tilde{f}(s_1) + (\beta_S(s_2^*) - \beta_R(s_2^*)) \tilde{f}(s_2^*) \frac{\partial s_2^*}{\partial s_1} \\ &\quad + (\beta_S - \beta_R) \left( \tilde{f}(s_1) - \tilde{f}(s_2^*) \frac{\partial s_2^*}{\partial s_1} \right)) \\ &= [\eta(s_1)]^2 \left( \tilde{f}(s_1) - \tilde{f}(s_2^*) \frac{\partial s_2^*}{\partial s_1} \right) \\ &\quad \times \left( \int_{s_1}^{s_2^*} (\beta_S(s) - \beta_R(s)) \tilde{f}(s) ds - \Pr_R[s \in (s_1, s_2^*(s_1))] (\beta_S - \beta_R) \right) \\ &\quad - \eta(s_1) ((\beta_S(s_1) - \beta_R(s_1)) \left( \tilde{f}(s_1) - \tilde{f}(s_2^*) \frac{\partial s_2^*}{\partial s_1} \right) \\ &\quad + (\beta_S - \beta_R) \left( \tilde{f}(s_1) - \tilde{f}(s_2^*) \frac{\partial s_2^*}{\partial s_1} \right)) \\ &= \eta(s_1) \left( \tilde{f}(s_1) - \tilde{f}(s_2^*) \frac{\partial s_2^*}{\partial s_1} \right) \left( \begin{array}{c} \eta(s_1) \int_{s_1}^{s_2^*} (\beta_S(s) - \beta_R(s)) \tilde{f}(s) ds \\ + (1 - \eta(s_1) \Pr_R[s \in (s_1, s_2^*(s_1))]) (\beta_S - \beta_R) \\ - (\beta_S(s_1) - \beta_R(s_1)) \end{array} \right) \\ &= \eta(s_1) \left( \tilde{f}(s_1) - \tilde{f}(s_2^*) \frac{\partial s_2^*}{\partial s_1} \right) (\Delta^{s_1, s_2^*(s_1)}(\emptyset) - \Delta(s_1)) \\ &= \eta(s_1) \left( \tilde{f}(s_1) - \tilde{f}(s_2^*) \frac{\partial s_2^*}{\partial s_1} \right) \gamma(s_1). \end{aligned}$$

Thus,

$$\frac{\partial \gamma(s_1)}{\partial s_1} = \eta(s_1) \left( \tilde{f}(s_1) - \tilde{f}(s_2^*) \frac{\partial s_2^*}{\partial s_1} \right) \gamma(s_1) - \Delta'(s_1). \quad (57)$$

Note that  $\Delta'(s_1) > 0$  and  $\frac{\partial s_2^*}{\partial s_1} < 0$  for any  $s_1 < \hat{s}$  by Lemma V.A (ii). Then, (57) implies that for any  $s_1 < \hat{s}$  such that  $\gamma(s_1) \leq 0$  it holds  $\gamma'(s_1) < 0$ . Consequently, if  $\gamma(s') = 0$  for some

$s' < \hat{s}$ , it is strictly decreasing for all  $s_1 \in [s', \hat{s})$ . Hence, the equilibrium condition  $\gamma(s_1) = 0$  can be satisfied for at most one value of  $s_1 < \hat{s}$ .

**Step 5.** We prove by contradiction that any equilibrium is a simple disclosure equilibrium. Assume thus an equilibrium which is not an SDE. By Lemma V.B, the set of non-disclosed signals has a positive measure. Upon non-disclosure, let the perceived disagreement be denoted by  $C > 0$ . Conditional on obtaining a signal,  $S$  wants to disclose if and only if the resulting disagreement  $\Delta(s)$  is smaller than  $C$ . Recall now that  $\Delta(s)$  is single peaked at  $\hat{s}$  by Lemma V.A(ii). Hence, given that a positive measure of signals is not disclosed, we must have  $C < \Delta(\hat{s})$ . Then, there are  $s_1, s_2$  satisfying  $\underline{s} < s_1 < s_2 < \bar{s}$  such that the actual disagreement is strictly higher than  $C$  after disclosing  $s \in (s_1, s_2)$  and strictly lower than  $C$  after disclosing  $s < s_1$  and  $s > s_2$ . In other words, this implies that for any putative equilibrium, there are  $s_1, s_2$  satisfying  $\underline{s} < s_1 < s_2 < \bar{s}$  such that  $S$  would strictly prefer not to disclose for  $\sigma \in (s_1, s_2)$  and strictly prefer to disclose if  $\sigma < s_1$  and  $\sigma > s_2$ . A putative equilibrium which is not an SDE thus gives rise to strict deviation incentives for  $S$ . ■

#### Lemma V.D

a) Assume that  $\beta_S > \beta_R$ . If  $\beta_R < 1 - \beta_S$ , i.e.  $R$  is more confident than  $S$ , then the equilibrium features  $\tilde{s} < s_1 < s_2$ , i.e. all signals congruent with  $R$ 's prior bias are disclosed. If  $\beta_R > 1 - \beta_S$ , i.e.  $R$  is less confident than  $S$ , then the equilibrium features  $s_1 < s_2 < \tilde{s}$ , i.e. all signals congruent with  $S$ 's prior bias are disclosed.

b) Assume that  $\beta_S < \beta_R$ . If  $\beta_R > 1 - \beta_S$ , i.e.  $R$  is more confident than  $S$ , then the equilibrium features  $s_1 < s_2 < \tilde{s}$ , i.e. all signals congruent with  $R$ 's prior bias are disclosed. If  $\beta_R < 1 - \beta_S$ , i.e.  $R$  is less confident than  $S$ , then the equilibrium features  $\tilde{s} < s_1 < s_2$ , i.e. all signals congruent with  $S$ 's prior bias are disclosed.

#### Proof.

We use the definitions of  $\gamma(s_1)$ ,  $s'_1$  and  $s_2^*(s_1)$  used in the proof of Lemma V.C (see Steps 2 and 3 there). By (56),  $\gamma(\hat{s}) < 0$ . At the same time, by (55) we have that  $\gamma(s'_1) > 0$ . Consequently, by the uniqueness of the SDE

$$s'_1 < s_1. \quad (58)$$

Given that  $\Delta(s)$  is single-peaked and the definition of  $s_2^*$ , this further implies

$$s_2 = s_2^*(s_1) < s_2^*(s'_1). \quad (59)$$

Now note that by construction  $s'_1 = \tilde{s}$  if  $\tilde{s} < \hat{s}$ , and  $s_2^*(s'_1) = \tilde{s}$  if  $\tilde{s} > \hat{s}$ . Then, the claims a) and b) follow by Lemma V.A(ii) together with (58) and (59). ■



## Appendix VII: Propositions 7 and 8

### Proof of Proposition 7

**Step 0.** We prove Point 1 in what follows. By assumption, it holds true that  $s_1 < s_2 < \tilde{s}$ . By Lemmas V.C and V.D it follows that  $\beta_R > 1 - \beta_S$ . We focus on proving that  $S$  would strictly prefer to commit to full disclosure if  $\beta_S > \beta_R$ . Note that combining  $\beta_R > 1 - \beta_S$  and  $\beta_S > \beta_R$  implies  $\beta_S > \frac{1}{2}$  and  $\beta_R \in (1 - \beta_S, \beta_S)$ . The proof that  $S$  instead prefers equilibrium disclosure given  $\beta_S < \beta_R$  and  $s_1 < s_2 < \tilde{s}$  is briefly outlined in our final step. The proof of Point 2 is conceptually identical to that of Point 1 and thus entirely omitted.

**Step 1.** Assume that  $\beta_S > \beta_R$ . From  $S$ 's perspective, the ex ante perceived disagreement in the SDE featuring thresholds  $\{s_1, s_2\}$  is given by:

$$\begin{aligned} & (1 - \varphi) \left[ E_R[\tilde{\beta}_S^{s_1, s_2} | \emptyset] - \tilde{\beta}_R^{s_1, s_2}(\emptyset) \right] \\ & + \varphi \int_{s_1}^{s_2} (\beta_S f(s|1) + (1 - \beta_S) f(s|0)) ds \left[ E_R[\tilde{\beta}_S^{s_1, s_2} | \emptyset] - \tilde{\beta}_R^{s_1, s_2}(\emptyset) \right] \\ & + \varphi \int_{s_1}^{s_2} (\beta_S f(s|1) + (1 - \beta_S) f(s|0)) \Delta(s) ds. \end{aligned}$$

Recall also that we know from Step 1 in the proof of Lemma V.C that

$$\begin{aligned} E_R[\tilde{\beta}_S^{s_1, s_2} | \emptyset] - \tilde{\beta}_R^{s_1, s_2}(\emptyset) &= \frac{\varphi}{(1 - \varphi) + \varphi \Pr_R(s \in (s_1, s_2))} \\ &\times \int_{s_1}^{s_2} (\beta_R f(s|1) + (1 - \beta_R) f(s|0)) \Delta(s) ds \\ &+ \frac{(1 - \varphi)}{(1 - \varphi) + \varphi \Pr_R(s \in (s_1, s_2))} (\beta_S - \beta_R). \end{aligned}$$

**Step 2.** We here consider a putative full disclosure equilibrium. From  $S$ 's perspective, the ex ante perceived disagreement in an equilibrium with full disclosure is simply

$$\begin{aligned} & \varphi \int_{s_1}^{s_2} (\beta_S f(s|1) + (1 - \beta_S) f(s|0)) \Delta(s) ds \\ & + \varphi \int_{s_1}^{s_2} (\beta_S f(s|1) + (1 - \beta_S) f(s|0)) \Delta(s) ds \\ & + (1 - \varphi) [\beta_S - \beta_R]. \end{aligned}$$

**Step 3.** We introduce two expressions which we shall call  $\Theta(\text{Partial})$  and  $\Theta(\text{Full})$ . These describe the expected perceived disagreement in  $S$ 's eyes under each of the two disclosure rules, when restricting ourselves to those events where either  $s \in [s_1, s_2]$  or  $S$  holds no signal

(as otherwise the perceived disagreement is identical under the two regimes). We have:

$$\begin{aligned}
& \Theta(\text{Partial}) \\
&= [\varphi \Pr_S(s \in (s_1, s_2)) + (1 - \varphi)] \left[ E_R[\tilde{\beta}_S^{s_1, s_2} | \emptyset] - \tilde{\beta}_R^{s_1, s_2}(\emptyset) \right] \\
&= [\varphi \Pr_S(s \in (s_1, s_2)) + (1 - \varphi)] \\
&\quad \times \left[ \frac{\varphi}{(1-\varphi) + \varphi \Pr_R(s \in (s_1, s_2))} \int_{s_1}^{s_2} (\beta_R f(s|1) + (1 - \beta_R) f(s|0)) \Delta(s) ds \right. \\
&\quad \quad \left. + \frac{(1-\varphi)}{(1-\varphi) + \varphi \Pr_R(s \in (s_1, s_2))} (\beta_S - \beta_R) \right]
\end{aligned}$$

and

$$\Theta(\text{Full}) = \varphi \int_{s_1}^{s_2} [\beta_S f(s|1) + (1 - \beta_S) f(s|0)] \Delta(s) ds + (1 - \varphi) (\beta_S - \beta_R).$$

Our objective is to identify conditions under which  $\Theta(\text{Partial}) > \Theta(\text{Full})$ , i.e.

$$\begin{aligned}
& [\varphi \Pr_S(s \in (s_1, s_2)) + (1 - \varphi)] \left[ E_R[\tilde{\beta}_S^{s_1, s_2} | \emptyset] - \tilde{\beta}_R^{s_1, s_2}(\emptyset) \right] \\
&> \varphi \int_{s_1}^{s_2} [\beta_S f(s|1) + (1 - \beta_S) f(s|0)] \Delta(s) ds + (1 - \varphi) (\beta_S - \beta_R).
\end{aligned}$$

**Step 4.** Define  $\Pr_{\hat{\beta}_R}(s \in (s_1, s_2))$  as the ex ante probability assigned  $s \in [s_1, s_2]$ , when using the prior  $\hat{\beta}_R$ . I.e. let:

$$\Pr_{\hat{\beta}_R}(s \in (s_1, s_2)) = \int_{s_1}^{s_2} (\hat{\beta}_R f(s|1) + (1 - \hat{\beta}_R) f(s|0)) ds.$$

We define  $\Delta^{(s_1, s_2)}(\emptyset, \hat{\beta}_R)$  as a slightly modified version of  $E_R[\tilde{\beta}_S^{s_1, s_2} | \emptyset] - \tilde{\beta}_R^{s_1, s_2}(\emptyset)$ , with the only difference that the distribution of signals is calculated based on the prior  $\hat{\beta}_R$ . We let

$$\begin{aligned}
& \Delta^{(s_1, s_2)}(\emptyset, \hat{\beta}_R) \\
&= \frac{\varphi}{(1 - \varphi) + \varphi \Pr_{\hat{\beta}_R}(s \in (s_1, s_2))} \\
&\quad \times \int_{s_1}^{s_2} (\hat{\beta}_R f(s|1) + (1 - \hat{\beta}_R) f(s|0)) \Delta(s) ds \\
&\quad + \frac{(1 - \varphi)}{(1 - \varphi) + \varphi \Pr_{\hat{\beta}_R}(s \in (s_1, s_2))} (\beta_S - \beta_R).
\end{aligned}$$

Let us finally define

$$\hat{\Theta}(\text{Partial}, \hat{\beta}_R) = [\varphi \Pr_S(s \in (s_1, s_2)) + (1 - \varphi)] \left[ \Delta^{(s_1, s_2)}(\emptyset, \hat{\beta}_R) \right]$$

and note that  $\hat{\Theta}(\text{Partial}, \beta_R) = \Theta(\text{Partial})$ .

In what follows, we shall consider the value of the above function for  $\widehat{\beta}_R = \beta_S$  and for  $\widehat{\beta}_R \in (1 - \beta_S, \beta_S)$ . We show in step 5 that  $\widehat{\Theta}(\text{Partial}, \beta_S) = \Theta(\text{Full})$ . We show in step 6 that for any  $\widehat{\beta}_R \in (1 - \beta_S, \beta_S)$ , we have  $\widehat{\Theta}(\text{Partial}, \widehat{\beta}_R) > \Theta(\text{Full})$ . Given that by assumption  $\beta_R \in (1 - \beta_S, \beta_S)$ , this implies that in particular  $\widehat{\Theta}(\text{Partial}, \beta_R) = \Theta(\text{Partial}) > \Theta(\text{Full})$ .

**Step 5.** Note that when setting  $\widehat{\beta}_R = \beta_S$ , we have:

$$\begin{aligned}
& \widehat{\Theta}(\text{Partial}, \beta_S) \\
&= [\varphi \Pr_S(s \in (s_1, s_2)) + (1 - \varphi)] [\Delta^{(s_1, s_2)}(\emptyset, \beta_S)] \\
&= [\varphi \Pr_S(s \in (s_1, s_2)) + (1 - \varphi)] \\
&\quad \times \left[ \frac{\varphi}{(1 - \varphi) + \varphi \Pr_S(s \in (s_1, s_2))} \int_{s_1}^{s_2} (\beta_S f(s|1) + (1 - \beta_S) f(s|0)) \Delta(s) ds \right. \\
&\quad \left. + \frac{(1 - \varphi)}{(1 - \varphi) + \varphi \Pr_S(s \in (s_1, s_2))} (\beta_S - \beta_R) \right] \\
&= \varphi \int_{s_1}^{s_2} [\beta_S f(s|1) + (1 - \beta_S) f(s|0)] \Delta(s) ds + (1 - \varphi) (\beta_S - \beta_R) \\
&= \Theta(\text{Full}).
\end{aligned}$$

**Step 6.** Here, we show that  $\Delta^{(s_1, s_2)}(\emptyset, \widehat{\beta}_R)$  increases (resp. decreases) as  $\widehat{\beta}_R$  decreases (resp. increases), for  $\widehat{\beta}_R \leq \beta_S$ . Note that we can rewrite  $\Delta^{(s_1, s_2)}(\emptyset, \widehat{\beta}_R)$  as follows:

$$\begin{aligned}
& \Delta^{(s_1, s_2)}(\emptyset, \widehat{\beta}_R) \\
&= \left[ \frac{\varphi \Pr_{\widehat{\beta}_R}(s \in (s_1, s_2))}{(1 - \varphi) + \varphi \Pr_{\widehat{\beta}_R}(s \in (s_1, s_2))} \int_{s_1}^{s_2} \frac{(\widehat{\beta}_R f(s|1) + (1 - \widehat{\beta}_R) f(s|0))}{\Pr_{\widehat{\beta}_R}(s \in (s_1, s_2))} \Delta(s) ds \right. \\
&\quad \left. + \frac{(1 - \varphi)}{(1 - \varphi) + \varphi \Pr_{\widehat{\beta}_R}(s \in (s_1, s_2))} (\beta_S - \beta_R) \right].
\end{aligned}$$

From the above expression, note that  $\Delta^{(s_1, s_2)}(\emptyset, \widehat{\beta}_R)$  is thus a weighted average of the expressions

$$\begin{aligned}
& E_{\widehat{\beta}_R} [\widetilde{\beta}_S(s) - \widetilde{\beta}_R(s) | s \in [s_1, s_2]] \\
&= \int_{s_1}^{s_2} \frac{(\widehat{\beta}_R f(s|1) + (1 - \widehat{\beta}_R) f(s|0))}{\Pr_{\widehat{\beta}_R}(s \in (s_1, s_2))} \Delta(s) ds
\end{aligned}$$

and  $(\beta_S - \beta_R)$ . The first expression is weighted by  $\frac{\varphi \Pr_{\widehat{\beta}_R}(s \in (s_1, s_2))}{(1 - \varphi) + \varphi \Pr_{\widehat{\beta}_R}(s \in (s_1, s_2))}$  and the second is weighted by  $\frac{(1 - \varphi)}{(1 - \varphi) + \varphi \Pr_{\widehat{\beta}_R}(s \in (s_1, s_2))}$ . In other words,  $\Delta^{(s_1, s_2)}(\emptyset, \widehat{\beta}_R)$  can be written as:

$$\Delta^{(s_1, s_2)}(\emptyset, \widehat{\beta}_R) = p(\widehat{\beta}_R) A(\widehat{\beta}_R) + (1 - p(\widehat{\beta}_R)) (\beta_S - \beta_R),$$

where we let

$$p(\widehat{\beta}_R) = \frac{\varphi \Pr_{\widehat{\beta}_R}(s \in (s_1, s_2))}{(1 - \varphi) + \varphi \Pr_{\widehat{\beta}_R}(s \in (s_1, s_2))}$$

and we let

$$A(\widehat{\beta}_R) = E_{\widehat{\beta}_R} \left[ \widetilde{\beta}_S(s) - \widetilde{\beta}_R(s) \mid s \in [s_1, s_2] \right].$$

The derivative of  $\Delta^{(s_1, s_2)}(\emptyset, \widehat{\beta}_R)$  w.r.t.  $\widehat{\beta}_R$  is thus given by

$$\begin{aligned} \frac{\partial \Delta^{(s_1, s_2)}(\emptyset, \widehat{\beta}_R)}{\partial \widehat{\beta}_R} &= \frac{\partial p(\widehat{\beta}_R)}{\partial \widehat{\beta}_R} A(\widehat{\beta}_R) + p(\widehat{\beta}_R) \frac{\partial A(\widehat{\beta}_R)}{\partial \widehat{\beta}_R} - \frac{\partial p(\widehat{\beta}_R)}{\partial \widehat{\beta}_R} (\beta_S - \beta_R) \\ &= p(\widehat{\beta}_R) \frac{\partial A(\widehat{\beta}_R)}{\partial \widehat{\beta}_R} + \frac{\partial p(\widehat{\beta}_R)}{\partial \widehat{\beta}_R} \left[ A(\widehat{\beta}_R) - (\beta_S - \beta_R) \right]. \end{aligned}$$

In order to prove that  $\frac{\partial \Delta^{(s_1, s_2)}(\emptyset, \widehat{\beta}_R)}{\partial \widehat{\beta}_R} < 0$  for  $\widehat{\beta}_R \in (1 - \beta_S, \beta_S)$ , it thus suffices to show that  $\frac{\partial A(\widehat{\beta}_R)}{\partial \widehat{\beta}_R} < 0$ ,

$$\left[ A(\widehat{\beta}_R) - (\beta_S - \beta_R) \right] > 0$$

and  $\frac{\partial p(\widehat{\beta}_R)}{\partial \widehat{\beta}_R} < 0$ . We show in what follows that these properties are indeed satisfied for  $\widehat{\beta}_R \in (1 - \beta_S, \beta_S)$ .

Note first that  $\frac{\partial \Pr_{\widehat{\beta}_R}(s \in (s_1, s_2))}{\partial \widehat{\beta}_R} = \int_{s_1}^{s_2} (f(s|1) - f(s|0)) ds$ , which is strictly negative given that we know that  $f(s|0) > f(s|1)$  for any  $s \in [s_1, s_2]$ , recalling that  $s_1 < s_2 < \widetilde{s}$  by assumption. It follows immediately that  $\frac{(1-\varphi)}{(1-\varphi) + \varphi \Pr_{\widehat{\beta}_R}(s \in (s_1, s_2))} = 1 - p(\widehat{\beta}_R)$  increases in  $\widehat{\beta}_R$  and that  $\frac{\varphi \Pr_{\widehat{\beta}_R}(s \in (s_1, s_2))}{(1-\varphi) + \varphi \Pr_{\widehat{\beta}_R}(s \in (s_1, s_2))} = p(\widehat{\beta}_R)$  decreases in  $\widehat{\beta}_R$ . Second, to show that  $A(\widehat{\beta}_R) - (\beta_S - \beta_R) > 0$  note that by the facts that  $s_1 < s_2 < \widetilde{s}$ , that  $\Delta(s_1) = \Delta(s_2)$  in SDE and that  $\Delta(s)$  is hump shaped in  $s$ , we obtain

$$\beta_S - \beta_R = \Delta(\widetilde{s}) < \Delta(s_1) < \Delta(s) \mid s \in (s_1, s_2).$$

Third, we now show that  $A(\widehat{\beta}_R) = E_{\widehat{\beta}_R} \left[ \left( \widetilde{\beta}_S(s) - \widetilde{\beta}_R(s) \right) \mid s \in [s_1, s_2] \right]$  decreases as  $\widehat{\beta}_R$  increases.

Note that:

$$\begin{aligned}
& \frac{\partial \left[ \int_{s_1}^{s_2} \frac{(\widehat{\beta}_R f(s|1) + (1 - \widehat{\beta}_R) f(s|0)) \Delta(s) ds}{\Pr_{\widehat{\beta}_R}(s \in (s_1, s_2))} \right]}{\partial \widehat{\beta}_R} \\
&= \int_{s_1}^{s_2} \frac{\left( \begin{aligned} & (f(s|1) - f(s|0)) \left[ \int_{s_1}^{s_2} \widehat{\beta}_R f(s|1) + (1 - \widehat{\beta}_R) f(s|0) ds \right] \\ & - \left[ \widehat{\beta}_R f(s|1) + (1 - \widehat{\beta}_R) f(s|0) \right] \left[ \int_{s_1}^{s_2} (f(s|1) - f(s|0)) ds \right] \end{aligned} \right)}{\left[ \Pr_{\widehat{\beta}_R}(s \in (s_1, s_2)) \right]^2} \Delta(s) ds \\
&= \frac{\left( \begin{aligned} & \left[ \int_{s_1}^{s_2} \widehat{\beta}_R f(s|1) + (1 - \widehat{\beta}_R) f(s|0) ds \right] \left[ \int_{s_1}^{s_2} (f(s|1) - f(s|0)) \Delta(s) ds \right] \\ & - \left[ \int_{s_1}^{s_2} (f(s|1) - f(s|0)) ds \right] \left[ \int_{s_1}^{s_2} (\widehat{\beta}_R f(s|1) + (1 - \widehat{\beta}_R) f(s|0)) \Delta(s) ds \right] \end{aligned} \right)}{\left[ \Pr_{\widehat{\beta}_R}(s \in (s_1, s_2)) \right]^2} \\
&= \frac{\left( \begin{aligned} & - \left[ \int_{s_1}^{s_2} (f(s|1) - f(s|0)) ds \right] \left[ \int_{s_1}^{s_2} (\widehat{\beta}_R f(s|1) + (1 - \widehat{\beta}_R) f(s|0)) \Delta(s) ds \right] \\ & + \left[ \int_{s_1}^{s_2} \widehat{\beta}_R f(s|1) + (1 - \widehat{\beta}_R) f(s|0) ds \right] \left[ \int_{s_1}^{s_2} (f(s|1) - f(s|0)) \Delta(s) ds \right] \end{aligned} \right)}{\left[ \Pr_{\widehat{\beta}_R}(s \in (s_1, s_2)) \right]^2} \\
&< \frac{\left( \begin{aligned} & - \left[ \int_{s_1}^{s_2} (f(s|1) - f(s|0)) ds \right] \left[ \int_{s_1}^{s_2} (\widehat{\beta}_R f(s|1) + (1 - \widehat{\beta}_R) f(s|0)) \Delta(s) ds \right] \\ & + \left[ \int_{s_1}^{s_2} \widehat{\beta}_R f(s|1) + (1 - \widehat{\beta}_R) f(s|0) ds \right] \left[ \int_{s_1}^{s_2} (f(s|1) - f(s|0)) \right] \left[ \int_{s_1}^{s_2} \Delta(s) ds \right] \end{aligned} \right)}{\left[ \Pr_{\widehat{\beta}_R}(s \in (s_1, s_2)) \right]^2} \\
&= \frac{- \left[ \int_{s_1}^{s_2} (f(s|1) - f(s|0)) ds \right] \left( \begin{aligned} & \left[ \int_{s_1}^{s_2} (\widehat{\beta}_R f(s|1) + (1 - \widehat{\beta}_R) f(s|0)) \Delta(s) ds \right] \\ & - \left[ \int_{s_1}^{s_2} \widehat{\beta}_R f(s|1) + (1 - \widehat{\beta}_R) f(s|0) ds \right] \left[ \int_{s_1}^{s_2} \Delta(s) ds \right] \end{aligned} \right)}{\left[ \Pr_{\widehat{\beta}_R}(s \in (s_1, s_2)) \right]^2} \\
&< 0.
\end{aligned}$$

Above, the first equality follows from the application of Leibniz' rule. The first and the second inequality follow from applying Hölder's inequality.

Thus, we have shown that  $\frac{\partial \Delta^{(s_1, s_2)}(\varnothing, \widehat{\beta}_R)}{\partial \widehat{\beta}_R} < 0$ . This implies that

$$\frac{\partial \widehat{\Theta}(\text{Partial}, \widehat{\beta}_R)}{\partial \widehat{\beta}_R} < 0.$$

In sum, we obtain that for  $\beta_S > \beta_R$  and  $s_1 < s_2 < \widetilde{s}$  it holds

$$\Theta(\text{Partial}) = \widehat{\Theta}(\text{Partial}, \beta_R) > \widehat{\Theta}(\text{Partial}, \beta_S) = \Theta(\text{Full}). \tag{60}$$

Here, the inequality follows from the previous inequality, while the second equality is by Step 5.

**Step 7.** Suppose now instead that  $\beta_S < \beta_R$  and  $s_1 < s_2 < \tilde{s}$ . Note that combining the assumptions  $\beta_S < \beta_R$  and  $s_1 < s_2 < \tilde{s}$  implies that  $\beta_R \in (\beta_S, 1)$  by Lemma V.D. The argument follows the same logic as above. It still holds true  $\widehat{\Theta}(\text{Partial}, \beta_S) = \Theta(\text{Full})$  and that  $\widehat{\Theta}(\text{Partial}, \beta_R) = \Theta(\text{Partial})$ . It also still holds true that  $\Theta(\text{Partial}, \widehat{\beta}_R)$  is decreasing in  $\widehat{\beta}_R$ . It follows that

$$\Theta(\text{Partial}) = \widehat{\Theta}(\text{Partial}, \beta_R) < \widehat{\Theta}(\text{Partial}, \beta_S) = \Theta(\text{Full}).$$

■

### Proof of Proposition 8

The argument here is exactly identical to the proof of the counterpart of this result for the case of binary signals (Proposition 4). ■

## Appendix VIII: Proposition 9

### Proof of Proposition 9.1.a).

**Step 1.** We here prove that there are finite  $\sigma' < \sigma''$  such that  $\sigma$  increases disagreement if and only if  $\sigma \notin [\sigma', \sigma'']$ . Recall that the perceived disagreement is given by

$$\Delta(d, \beta_S, \beta_R) = |E_R[E_S[\omega|\sigma] | d] - E_R[\omega | d]|. \quad (61)$$

Note that

$$E_i[\omega | \sigma] = \frac{\mu_i \frac{1}{\gamma_i^2} + \sigma \frac{1}{\gamma_\varepsilon^2}}{\frac{1}{\gamma_i^2} + \frac{1}{\gamma_\varepsilon^2}}.$$

Hence,

$$\Delta(\sigma) = |E_S[\omega | \sigma] - E_R[\omega | \sigma]| = |D(\sigma)|, \quad (62)$$

where

$$\begin{aligned} & D(\sigma) \\ &= \frac{\mu_S \frac{1}{\gamma_S^2} + \sigma \frac{1}{\gamma_\varepsilon^2}}{\frac{1}{\gamma_S^2} + \frac{1}{\gamma_\varepsilon^2}} - \frac{\mu_R \frac{1}{\gamma_R^2} + \sigma \frac{1}{\gamma_\varepsilon^2}}{\frac{1}{\gamma_R^2} + \frac{1}{\gamma_\varepsilon^2}} \\ &= -\frac{\gamma_\varepsilon^2}{(\gamma_R^2 + \gamma_\varepsilon^2)(\gamma_S^2 + \gamma_\varepsilon^2)} [\sigma (\gamma_R^2 - \gamma_S^2) - \gamma_R^2 \mu_S + \gamma_S^2 \mu_R + (\mu_R - \mu_S) \gamma_\varepsilon^2]. \end{aligned} \quad (63)$$

Note that  $D(\sigma)$  is a linear function, so that it has a unique root in  $\mathfrak{R}$ . Consequently,  $|D(\sigma)|$  is V-shaped in  $\sigma$ , with the minimum value of 0. It follows immediately that there exist  $\sigma' < \sigma''$

such that  $|D(\sigma)| > |\mu_S - \mu_R|$  (i.e.  $\sigma$  increases disagreement) if and only if  $\sigma \notin [\sigma', \sigma'']$ .

**Step 2.** Let us show that any equilibrium under  $\gamma_S \neq \gamma_R$  and  $\mu_S \neq \mu_R$  features a disclosure interval  $[\underline{\sigma}, \bar{\sigma}]$  such that  $\sigma$  is disclosed if and only if  $\sigma \in [\underline{\sigma}, \bar{\sigma}]$ .

First, note that FD never exists under  $\gamma_S \neq \gamma_R$  by Step 1, since otherwise  $S$  would have an incentive to deviate by concealing any signal  $\sigma$  such that  $|D(\sigma)| > |\mu_S - \mu_R|$ .

Second, note that an equilibrium with disclosure rule  $\tilde{D}$  must feature a positive and finite value of perceived disagreement conditional on no disclosure  $\Delta^{\tilde{D}}(\emptyset)$ . The fact that  $\Delta^{\tilde{D}}(\emptyset)$  should be finite can be shown by contradiction. Suppose indeed that  $\Delta^{\tilde{D}}(\emptyset)$  is not finite. Then there would exist an FD-equilibrium, as  $S$  would always favour disclosing over not disclosing. But we know that there exists no FD-equilibrium, as stated above. The fact that  $\Delta^{\tilde{D}}(\emptyset)$  must be positive can also be shown by contradiction. If this is not the case, then in equilibrium all signals must be concealed other than the unique signal  $\tilde{\sigma}$  such that  $D(\tilde{\sigma}) = 0$ . Then,  $R$ 's posterior belief distribution will not change after no disclosure, while  $R$  would expect that  $S$ 's posterior mean will be strictly between  $\mu_R$  and  $\mu_S$ .<sup>30</sup> Hence,  $\Delta^{\tilde{D}}(\emptyset) > 0$  which is a contradiction.

Consider thus an equilibrium with disclosure rule  $\tilde{D}$  featuring a positive and finite  $\Delta^{\tilde{D}}(\emptyset)$ . In this case, every signal  $\sigma$  such that  $\Delta(\sigma) \leq \Delta^{\tilde{D}}(\emptyset)$  will be disclosed, and every signal such that  $\Delta(\sigma) > \Delta^{\tilde{D}}(\emptyset)$  will not be disclosed. Given that  $\Delta^{\tilde{D}}(\sigma)$  is V-shaped in  $\sigma$ , the claim follows. Next, by (62) and linearity of  $D(\sigma)$  it follows that  $\Delta(\sigma)$  is symmetric around  $\tilde{\sigma}$ , where

$$\begin{aligned} \Delta(\tilde{\sigma}) &= 0 \Leftrightarrow \\ D(\tilde{\sigma}) &= 0 \Leftrightarrow \\ \tilde{\sigma} &= \frac{\mu_S(\gamma_R^2 + \gamma_\varepsilon^2) - \mu_R(\gamma_S^2 + \gamma_\varepsilon^2)}{\gamma_R^2 - \gamma_S^2}. \end{aligned} \tag{64}$$

Consequently, the disclosure interval (characterized by  $\{\sigma : \Delta(\sigma) \leq \Delta^{\tilde{D}}(\emptyset)\}$  for a given equilibrium  $\Delta^{\tilde{D}}(\emptyset)$ ) is also symmetric around  $\tilde{\sigma}$ .

**Step 3.** This shows existence of an equilibrium of the type characterized in Step 2. Denote  $R$ 's perceived disagreement conditional on disclosure in an equilibrium featuring

---

<sup>30</sup>See Chapter 10 of Technical Appendix in Vives (2010).

disclosure interval  $(\underline{\sigma}, \bar{\sigma})$  as  $\Delta^{(\underline{\sigma}, \bar{\sigma})}(\emptyset)$ . We have

$$\begin{aligned}
& \Delta^{(\underline{\sigma}, \bar{\sigma})}(\emptyset) \\
&= \left| \tau \mu_S + (1 - \tau) \left[ \int_{\sigma \leq \underline{\sigma}} E_S[\omega|\sigma] \tilde{f}(\sigma | d = \emptyset, \underline{\sigma}, \bar{\sigma}) d\sigma \right. \right. \\
&\quad \left. \left. + \int_{\sigma \geq \bar{\sigma}} E_S[\omega|\sigma] \tilde{f}(\sigma | d = \emptyset, \underline{\sigma}, \bar{\sigma}) d\sigma \right] \right. \\
&\quad \left. - \tau \mu_R - (1 - \tau) \left[ \int_{\sigma \leq \underline{\sigma}} E_R[\omega|\sigma] \tilde{f}(\sigma | d = \emptyset, \underline{\sigma}, \bar{\sigma}) d\sigma \right. \right. \\
&\quad \left. \left. + \int_{\sigma \geq \bar{\sigma}} E_R[\omega|\sigma] \tilde{f}(\sigma | d = \emptyset, \underline{\sigma}, \bar{\sigma}) d\sigma \right] \right| \\
&= \left| \tau(\mu_S - \mu_R) + (1 - \tau) \left[ \int_{\sigma \leq \underline{\sigma}} (E_S[\omega|\sigma] - E_R[\omega|\sigma]) \tilde{f}(\sigma | d = \emptyset, \underline{\sigma}, \bar{\sigma}) d\sigma \right. \right. \\
&\quad \left. \left. + \int_{\sigma \geq \bar{\sigma}} (E_S[\omega|\sigma] - E_R[\omega|\sigma]) \tilde{f}(\sigma | d = \emptyset, \underline{\sigma}, \bar{\sigma}) d\sigma \right] \right|,
\end{aligned}$$

where  $\tau = P(\sigma = \emptyset | d = \emptyset, \underline{\sigma}, \bar{\sigma})$ , and  $\tilde{f}(\sigma | d = \emptyset, \underline{\sigma}, \bar{\sigma})$  is the conditional distribution of  $\sigma$  given  $d = \emptyset$  from the perspective of  $R$ , in an equilibrium featuring disclosure interval  $[\underline{\sigma}, \bar{\sigma}]$ .

Given the V-shape of  $\Delta(\sigma)$ , a profile  $\underline{\sigma}, \bar{\sigma}$  constitutes an equilibrium if and only if:

$$\Delta^{(\underline{\sigma}, \bar{\sigma})}(\emptyset) = \Delta(\underline{\sigma}) = \Delta(\bar{\sigma}). \quad (65)$$

For any given  $x > 0$ , let  $\underline{\sigma}(x), \bar{\sigma}(x)$  denote the unique pair of signals satisfying

$$\Delta(\underline{\sigma}(x)) = \Delta(\bar{\sigma}(x)) = x,$$

which exists for any  $x > 0$  given that (63) is unbounded. Hence, the equilibrium condition (65) is equivalent to

$$\Delta^{(\underline{\sigma}(x), \bar{\sigma}(x))}(\emptyset) = x. \quad (66)$$

Note that for  $x = 0$ , i.e. if all signals are concealed, then it must be true that

$$\Delta^{(\underline{\sigma}(0), \bar{\sigma}(0))}(\emptyset) > 0 \quad (67)$$

(see Step 2).

Next, let us show that

$$\lim_{x \rightarrow \infty} \Delta^{(\underline{\sigma}(x), \bar{\sigma}(x))}(\emptyset) = |\mu_S - \mu_R| < x.$$

Note that

$$\Delta^{(\underline{\sigma}(x), \bar{\sigma}(x))}(\emptyset) = |P(\sigma = \emptyset | d = \emptyset, \underline{\sigma}(x), \bar{\sigma}(x))(\mu_S - \mu_R) + \varsigma(\underline{\sigma}(x), \bar{\sigma}(x))|,$$



where

$$\begin{aligned} & \zeta(\underline{\sigma}(x), \bar{\sigma}(x)) \\ &= P(\sigma \neq \emptyset | d = \emptyset, \underline{\sigma}(x), \bar{\sigma}(x)) \left[ \int_{\sigma \leq \underline{\sigma}(x)} D(\sigma) \tilde{f}(\sigma | d = \emptyset, \underline{\sigma}(x), \bar{\sigma}(x)) d\sigma \right. \\ & \quad \left. + \int_{\sigma \geq \bar{\sigma}(x)} D(\sigma) \tilde{f}(\sigma | d = \emptyset, \underline{\sigma}(x), \bar{\sigma}(x)) d\sigma \right]. \end{aligned}$$

Let us show that  $\lim_{x \rightarrow \infty} \zeta(\underline{\sigma}(x), \bar{\sigma}(x)) = 0$ . We have

$$\begin{aligned} \zeta(\underline{\sigma}(x), \bar{\sigma}(x)) &= \frac{\varphi P(\sigma \notin [\underline{\sigma}(x), \bar{\sigma}(x)])}{\varphi P(\sigma \notin [\underline{\sigma}(x), \bar{\sigma}(x)]) + (1 - \varphi)} \\ & \quad \times \left[ \int_{\sigma \leq \underline{\sigma}(x)} D(\sigma) \tilde{f}(\sigma | d = \emptyset, \underline{\sigma}(x), \bar{\sigma}(x)) d\sigma \right. \\ & \quad \left. + \int_{\sigma \geq \bar{\sigma}(x)} D(\sigma) \tilde{f}(\sigma | d = \emptyset, \underline{\sigma}(x), \bar{\sigma}(x)) d\sigma \right] \\ &= \frac{\varphi P(\sigma \notin [\underline{\sigma}(x), \bar{\sigma}(x)])}{\varphi P(\sigma \notin [\underline{\sigma}(x), \bar{\sigma}(x)]) + (1 - \varphi)} \left[ \int_{\sigma \leq \underline{\sigma}(x)} D(\sigma) \frac{f(\sigma)}{P(\sigma \notin [\underline{\sigma}(x), \bar{\sigma}(x)])} d\sigma \right. \\ & \quad \left. + \int_{\sigma \geq \bar{\sigma}(x)} D(\sigma) \frac{f(\sigma)}{P(\sigma \notin [\underline{\sigma}(x), \bar{\sigma}(x)])} d\sigma \right] \\ &= \frac{\varphi}{\varphi P(\sigma \notin [\underline{\sigma}(x), \bar{\sigma}(x)]) + (1 - \varphi)} \left[ \int_{\sigma \leq \underline{\sigma}(x)} D(\sigma) f(\sigma) d\sigma \right. \\ & \quad \left. + \int_{\sigma \geq \bar{\sigma}(x)} D(\sigma) f(\sigma) d\sigma \right]. \end{aligned} \quad (68)$$

Note that

$$\lim_{x \rightarrow \infty} \frac{\varphi}{\varphi P(\sigma \notin [\underline{\sigma}(x), \bar{\sigma}(x)]) + (1 - \varphi)} = \frac{\varphi}{1 - \phi}. \quad (69)$$

At the same time, given that  $E[D(\sigma)]$  must be finite (since  $\sigma$  is normally distributed and  $D(\sigma)$  is linear in  $\sigma$ ), we have

$$\begin{aligned} & \lim_{x \rightarrow \infty} \int_{\sigma \leq \underline{\sigma}(x)} D(\sigma) f(\sigma) d\sigma \\ &= \lim_{x \rightarrow \infty} \left( E[D(\sigma)] - \int_{\sigma > \underline{\sigma}(x)} D(\sigma) f(\sigma) d\sigma \right) \\ &= E[D(\sigma)] - \lim_{x \rightarrow \infty} \int_{\underline{\sigma}(x)}^{\infty} D(\sigma) f(\sigma) d\sigma \\ &= E[D(\sigma)] - E[D(\sigma)] = 0. \end{aligned} \quad (70)$$

where the third equality is due to the fact that  $\underline{\sigma}(x)$  must be linear in  $x$  (since its inverse function  $\Delta(\sigma)$  is linear in  $\sigma$ ). By the same argument,

$$\lim_{x \rightarrow \infty} \int_{\sigma \geq \bar{\sigma}(x)} D(\sigma) f(\sigma) d\sigma = 0. \quad (71)$$

(68), (69), (70) and (71) together imply

$$\lim_{x \rightarrow \infty} \zeta(\underline{\sigma}(x), \bar{\sigma}(x)) = 0. \quad (72)$$

We may conclude that

$$\lim_{x \rightarrow \infty} \Delta^{(\underline{\sigma}(x), \bar{\sigma}(x))}(\emptyset) = |\mu_S - \mu_R| < x. \quad (73)$$

Given the continuity of

$$\Delta^{(\underline{\sigma}(x), \bar{\sigma}(x))}(\emptyset)$$

in  $x$ , it follows from (67) and (73) that as  $x$  increases from  $\underline{x} = 0$  to  $+\infty$ , there is some  $x \in (0, +\infty)$  such that the equilibrium condition (66) is satisfied. ■

### Proof of Proposition 9.1.b).

Assume  $\mu_R > \mu_S$  (the proof for  $\mu_R < \mu_S$  proceeds symmetrically). We obtain

$$\begin{aligned} \mu_S - \tilde{\sigma} &= \frac{(\mu_R - \mu_S)(\gamma_S^2 + \gamma_\varepsilon^2)}{\gamma_R^2 - \gamma_S^2}, \\ \mu_R - \tilde{\sigma} &= \frac{(\mu_R - \mu_S)(\gamma_R^2 + \gamma_\varepsilon^2)}{\gamma_R^2 - \gamma_S^2}. \end{aligned}$$

Hence, if  $\gamma_R > \gamma_S$  we get  $\tilde{\sigma} < \mu_S < \mu_R$ , while if  $\gamma_R < \gamma_S$  we get  $\tilde{\sigma} > \mu_R > \mu_S$ . Thus,  $\tilde{\sigma} \notin [\mu_S, \mu_R]$  while  $\tilde{\sigma}$  is closer to the mean of the more confident player.

Finally, note that the Hausdorff distance between the disclosure interval  $D = [\tilde{\sigma} - \eta, \tilde{\sigma} + \eta]$  and a mean  $\mu_i$  is given by the largest distance between a point in  $D$  and  $\mu_i$ . In case if  $\gamma_R > \gamma_S$  so that  $\tilde{\sigma} < \mu_S < \mu_R$ , the furthest point from either  $\mu_S$  or  $\mu_R$  is  $\tilde{\sigma} - \eta$ , which is then closer to the prior of the more confident player  $\mu_S$ . In case if  $\gamma_R < \gamma_S$  so that  $\tilde{\sigma} > \mu_R > \mu_S$ , the furthest point from either  $\mu_S$  or  $\mu_R$  is  $\tilde{\sigma} + \eta$ , which is then closer to the prior of the more confident player  $\mu_R$ . ■

### Proof of Proposition 9.2.

Let  $\gamma_S^2 \neq \gamma_R^2$  and  $\mu_S = \mu_R = \mu$ . Recall from Step 1 in the proof of Proposition 9.1.b) that  $\Delta(\sigma)$  is symmetrically V-shaped around  $\tilde{\sigma}$  while  $\Delta(\tilde{\sigma}) = 0$ . This immediately implies that any signal weakly increases disagreement relative to the prior disagreement (of 0).

We first show that there exists an equilibrium  $\tilde{D}$  where all signals besides  $\tilde{\sigma} = \mu$  are not disclosed, while  $\tilde{\sigma}$  is disclosed with an arbitrary probability in  $[0, 1]$ . Indeed, in this case the posterior disagreement conditional on no disclosure  $\Delta^{\tilde{D}}(\emptyset)$  is equal to 0, since the posterior means of both  $S$ 's and  $R$ 's belief distributions in the eyes of  $R$  are then equal to the prior mean  $\mu$ . Hence, given the shape of  $\Delta(\sigma)$ ,  $S$  indeed strictly prefers non-disclosure over disclosure for any  $\sigma$  except for  $\tilde{\sigma}$ , where he is indifferent.

Let us show that no other equilibrium exists. Assume by contradiction that there exists an equilibrium disclosure rule  $D'$  such that some other signals besides  $\tilde{\sigma}$  are disclosed with positive probability. Then,  $\Delta^{D'}(\emptyset)$  must be strictly positive as otherwise  $S$  would strictly prefer to conceal all signals other than  $\tilde{\sigma}$ . But if  $\Delta^{D'}(\emptyset) > 0$ , then every signal  $\sigma$  such that

$\Delta(\sigma) \leq \Delta^{D'}(\emptyset)$  will be disclosed and every signal such that  $\Delta(\sigma) > \Delta^{D'}(\emptyset)$  will not be disclosed. Given the symmetric V-shape of  $\Delta(\sigma)$ ,  $S$  must disclose signals belonging to an interval  $[\underline{\sigma}, \bar{\sigma}]$  which is symmetric around  $\tilde{\sigma}$ . Then,

$$\Delta^{D'}(\emptyset) = \left| \tau(\mu_S - \mu_R) + (1 - \tau) \left[ \int_{\sigma \leq \underline{\sigma}} D(\sigma) \tilde{f}(\sigma | d = \emptyset, \underline{\sigma}, \bar{\sigma}) d\sigma + \int_{\sigma \geq \bar{\sigma}} D(\sigma) \tilde{f}(\sigma | d = \emptyset, \underline{\sigma}, \bar{\sigma}) d\sigma \right] \right|,$$

where  $D(\sigma)$  is given by (63),  $\tau = P(\sigma = \emptyset | d = \emptyset, \underline{\sigma}, \bar{\sigma})$ , and  $\tilde{f}(\sigma | d = \emptyset, \underline{\sigma}, \bar{\sigma})$  is the conditional distribution of  $\sigma$  given  $d = \emptyset$  from the perspective of  $R$ . Given  $\mu_S = \mu_R$  this further simplifies to

$$\Delta^{D'}(\emptyset) = \left| (1 - \tau) \left[ \int_{\sigma \leq \underline{\sigma}} Q(\sigma) d\sigma + \int_{\sigma \geq \bar{\sigma}} Q(\sigma) d\sigma \right] \right|, \quad (74)$$

where  $Q(\sigma) = D(\sigma) \frac{f(\sigma)}{P(\sigma \notin [\underline{\sigma}(x), \bar{\sigma}(x)])}$ . Note furthermore that given linearity of  $D(\sigma)$ , we have  $D(\tilde{\sigma} + z) = -D(\tilde{\sigma} - z)$ . Besides, by (64), we have  $\tilde{\sigma} = \mu$  and hence  $f(\tilde{\sigma} + z) = f(\tilde{\sigma} - z)$ . Thus,  $Q(\tilde{\sigma} + z) = -Q(\tilde{\sigma} - z)$ , which by the symmetry of disclosure interval around  $\tilde{\sigma}$  yields

$$\int_{\sigma \leq \underline{\sigma}} Q(\sigma) d\sigma + \int_{\sigma \geq \bar{\sigma}} Q(\sigma) d\sigma = 0.$$

This together with (74) implies  $\Delta^{D'}(\emptyset) = 0$ , which is a contradiction. Thus, there exists no equilibrium where any signal besides for  $\tilde{\sigma}$  is disclosed. ■

### Proof of Proposition 9.3.

**Step 1.** If  $\gamma_S^2 = \gamma_R^2$  and  $\mu_S \neq \mu_R$ , then

$$\begin{aligned} \Delta(\sigma) &= |E_S[\omega | \sigma] - E_R[\omega | \sigma]| = \left| \frac{\gamma_\varepsilon^2}{\gamma_R^2 + \gamma_\varepsilon^2} (\mu_S - \mu_R) \right| \\ &< |\mu_S - \mu_R|. \end{aligned} \quad (75)$$

Consequently, there exists an equilibrium with full disclosure. Let us show that no other equilibrium exists. Assume by contradiction that there exists an equilibrium featuring a

non-empty disclosure interval  $\tilde{I}$ . Then,

$$\begin{aligned}
& \Delta^{\tilde{I}}(\emptyset) \\
&= \left| \begin{array}{l} P(\sigma = \emptyset | d = \emptyset, \underline{\sigma}, \bar{\sigma})(\mu_S - \mu_R) \\ + P(\sigma \neq \emptyset | d = \emptyset, \underline{\sigma}, \bar{\sigma}) \int_{\sigma \in \tilde{I}} D(\sigma) \tilde{f}(\sigma | d = \emptyset) d\sigma \end{array} \right| \\
&= \left| \begin{array}{l} P(\sigma = \emptyset | d = \emptyset, \underline{\sigma}, \bar{\sigma})(\mu_S - \mu_R) \\ + P(\sigma \neq \emptyset | d = \emptyset, \underline{\sigma}, \bar{\sigma}) \left( \frac{\gamma_\varepsilon^2}{\gamma_R^2 + \gamma_\varepsilon^2} (\mu_S - \mu_R) \right) \end{array} \right| \\
&> \left| \frac{\gamma_\varepsilon^2}{\gamma_R^2 + \gamma_\varepsilon^2} (\mu_S - \mu_R) \right| = \Delta(\sigma)
\end{aligned}$$

for any  $\sigma$ . Hence,  $S$  would have incentive to deviate to disclosure for any  $\sigma \in \tilde{I}$ .

**Step 2.** Let  $\gamma_S^2 = \gamma_R^2$  and  $\mu_S = \mu_R$ , then by (63)  $D(\sigma) = \Delta(\sigma) = 0$  for any  $\sigma$ . It follows immediately that any disclosure rule is an equilibrium disclosure rule. ■