# Bayesian Persuasion under Partial Commitment

Daehong Min*
Department of Economics
University of Arizona

April 14, 2016

## Abstract

This paper studies a variation of the Bayesian persuasion model in which the sender's commitment to a signaling device binds with probability less than one. The receiver knows the commitment probability but cannot tell whether the commitment is binding or not. We focus on the welfare implication of partial commitment: Are the sender and receiver better off or worse off as the commitment probability increases (or decreases)? We first show that the sender is *weakly* better off as the commitment probability increases, which does not depend on the assumptions on players' preferences and the common prior distribution over the state space. Then, we study the model in a specific environment: the uniform-quadratic case. In the uniform-quadratic case, we show that for any level of the sender's bias (even when the bias is arbitrarily high), both players are *strictly* better off as they move from the no-commitment through the partial-commitment to the full-commitment case. To establish the strict welfare improvement from the no-commitment to the partial-commitment case, it suffices to consider only three types of signaling devices. Interestingly, one of them can achieve the best outcome in Blume et al. (2007) and Goltman et al. (2009).

*Keywords*: Communication, Bayesian Persuasion, Cheap Talk

# 1    Introduction

The Bayesian persuasion literature assumes that a Sender's commitment to a signaling device is fully binding. That is, once a Sender chose a signaling device, the Sender should truthfully report signals from the chosen signaling device to a Receiver. For example, consider a drug company (a Sender) seeking the approval of a newly developed drug from the FDA (a Receiver). The drug approval procedure requires the company to conduct a drug experiment and report data from it. The drug company can choose a design for the drug experiment (a signaling device). And the company's commitment to the chosen experiment is binding by a law: it should report all data (signals) from the experiment as they are; hiding or falsifying data are legally prohibited.

However, we can find situations in which the Sender's commitment is compromised. The story of Dr. Robert Fiddes is a good example. Dr. Fiddes had been an outstanding figure in the field of drug experiments: he had led more than 200 experiments that drew approval decisions from the FDA. However, it turned out that he manipulated data from his study by falsifying data that are against the approval and even creating fake data in favor of the study. A surprising fact is that the FDA could not detect his misconduct and Dr. Fiddes knew this beforehand. Dr. Fiddes said that he would not be caught unless his former employee blew the whistle.[1] He knew that the FDA's monitoring system is not good enough to detect his misconduct. Simply, he knew that his commitment to an experimental design is not actually binding and also knew that he could report data in his discretion without being caught.

The short anecdote above demonstrates a possibility of imperfect commitment. We incorporate this possibility by relaxing the full-commitment assumption. Namely, we study a model of partial commitment in which the Sender's commitment is binding with a probability less than one. In the model, there are two players, the Sender and Receiver. The Sender moves first and he chooses a signaling device. Then, with a exogenously given probability, the Sender learns that his commitment is *not* binding (as Dr. Fiddes did). If this event occurs, the Sender can choose what signals he will send to the Receiver. With the other probability, the Sender learns that his commitment is binding (as the other "honest" drug experimenter believed the FDA's monitoring is austere enough to detect their misconduct). If this event occurs, the Sender truthfully reports signals from the previously chosen signaling device to the Receiver. The Receiver knows the probability that the Sender's commitment is binding *but* cannot tell whether the signal she observes comes from the Sender whose commitment is

---

[1]See the Newyork times article: http://www.nytimes.com/1999/05/17/business/a-doctor-s-drug-trials-turn-into-fraud.html?pagewanted=all or Beach, Judith E. "Clinical trials integrity: A CRO perspective*." Accountability in research 8.3 (2001): 245-260.

binding or not (as the FDA in the story above). Finally, the Receiver makes inference on the true state based on the signal she observes and takes an action that affects both the Sender's and Receiver's payoffs.[2]

In this paper, our focus is mainly on the welfare implication of partial commitment: Are the Sender and Receiver better off or worse off as the commitment probability increases (or decreases)? We first show that the Sender cannot be worse off as the commitment probability increases, which does not depend on the assumptions on players' preferences and the common prior distribution over the state space. Then, we study the model in a specific environment: the uniform-quadratic case introduced by Crawford and Sobel (1982). In the uniform-quadratic case, we show that for any level of the Sender's bias (even when the bias is arbitrarily high), both players are *strictly* better off as they move from the no-commitment through the partial-commitment to the full-commitment case. To establish the strict welfare improvement from the no-commitment to the partial-commitment case, it suffices to consider only three types of signaling devices. Interestingly, one of them can achieve the best outcome in Blume, Board, and Kawamura (2007) and Goltsman, Hörner, Pavlov, and Squintani (2009).

In our model the Sender can send signals in his discretion with a positive probability unlike in the Bayesian Persuasion model. Note that as the commitment probability decreases, the Sender is more likely to be free from his commitment and use the "freedom" to favor himself as Dr. Fiddes did (i.e. the less commitment probability, the more chance to behave like Dr. Fiddes). Hence, it may be tempting to say that the Sender is better off as the commitment probability decreases. However, our first finding predicts the exact opposite: the Sender is *weakly* better off as the commitment probability *increases*. In our model, the Sender chooses a signaling device and publicly announces it. Hence, a chosen signaling device induces a subgame. Then, the Sender can be considered as a *Stackellberg leader* who can choose which subgame he wants to be in. At an equilibrium, the Sender looks over all equilibria in all subgames and chooses to be in the subgame with the equilibrium that gives him the highest payoffs. Not surprisingly, the set of equilibrium outcomes that the Sender can choose *weakly* increases as the commitment probability increases. In other words, as the commitment probability increases, the Sender can always choose the equilibrium outcome that he could have chosen before the commitment probability has increased. Thus, the Sender cannot be harmed by having higher commitment probabilities.

Then, we study this model of partial commitment with the well-known uniform-quadratic framework in Crawford and Sobel (1982). That is, we assume that the state space and the

---

[2]One may think of the partial commitment model as a model lying between the cheap-talk and the Bayesian persuasion model: in one extreme case that the commitment probability is zero, the model collapses to the cheap-talk model; in the other extreme case, the model collapses to the Bayesian persuasion model.

common prior over the state space are the unit interval, $[0, 1]$, and the uniform distribution over this interval, respectively. In addition, both the Sender's and Receiver's preferences are represented by quadratic loss functions and the Sender's preference differs from the Receiver's by a positive bias. In this uniform-quadratic case, we show that the *ex ante* payoffs of both players are *strictly* improved as they move from the no-commitment through the partial-commitment to the full-commitment case.

We first compare the welfare under full commitment with that under partial commitment. We show that there exists the unique equilibrium outcome under full commitment: the Sender commits to a fully revealing signaling device that tells the exact true state and the Receiver learns the exact true state, which is a first-best outcome. Then, we show that this outcome is not attainable under partial commitment. Once we allow the Sender to have a chance to send signals in his discretion even with an arbitrarily small probability, we have the Sender's Incentive Compatibility constraints (henceforth IC constraints) that should be satisfied at an equilibrium. Apparently, the unique equilibrium outcome we get under full commitment cannot be achieved under partial commitment since the equilibrium outcome under full commitment is not compatible with the Sender IC constraints we necessarily have in the partial-commitment case: when the Sender's commitment is not binding, the Sender will not send a signal that tells the true state but send a signal that exaggerates the true state by the amount of a positive bias.

To establish the welfare improvment from the no-commitment to the partial-commitment case, we consider three types of signaling devices: a semi-fully revealing, a fully concealing, and a convexified signaling devices. We study each subgame after the Sender commits to one of these signaling devices. Then, we show that, in each subgame, we can construct a strictly better equilibrium than the best equilibrium we can have under the no-commitment case (or in the cheap-talk model). Those three types of signaling devices are sufficient to cover the entire parameter space with a strictly better equilibrium than the best equilibrium in the cheap-talk model. That is, for any commitment probability and any level of the bias (even arbitrarily large), the Sender's partial commitment to one of three types of signaling devices allows us to construct a strictly better equilibrium.

The semi-fully revealing signaling device tells the exact true states at the low states but pools high states into a single interval. The fully concealing signaling device merely generates random noise and the convexified signaling device is a convex combination of the fully concealing signaling device and a signaling rule called *the front-loading strategy*[3]. For high biases, the Sender's partial commitment to the semi-fully revealing signaling device

---

[3]We follow the terminology that is in Blume, Board, and Kawamura (2007)

3

strictly improves the *ex ante* payoffs of both players by guaranteeing a certain degree of the information transmission. When the bias is low, both players are strictly better off by the Sender's partial commitment to signaling devices that involve noise, the fully concealing and convexified signaling devices. Interestingly, a convexified signaling device can achieve the efficiency bound established in the mediation model by Goltsman et al. (2009) and in the noise model by Blume et al. (2007).

This paper is related to two strands of communication literature, the cheap talk and the Bayesian persuasion literature. In the Bayesian Persuasion literature, there have been papers that study many variations of the Bayesian persuasion model in Kamenica and Gentzkow (2011a). But these variations assume the Sender's commitment is always binding.[4] To my knowledge, this paper is the first one that considers a model with partial commitment.

In the cheap talk literature, there have been papers studying various ways to improve communication and welfare. Goltsman et al. (2009) show that a mediation scheme can improve welfare and also establish the efficiency bound for it. Blume et al. (2007) show that adding a small amount of noise in the communication process can improve welfare; furthermore, they show how the efficiency bound in Golsman et al. (2009) can be attained with simple noise. The main finding in this paper adds an observation to this strand of papers: the Sender's partial commitment can improve the welfare even when the probability that the commitment is binding is arbitrarily small and the Sender's bias is arbitrarily large.

This paper closely relates to Blume et al. (2007) and a result in this paper resembles that in Kim and Pogach (2014). On the one hand, when the Sender commits to the fully concealing signaling device, the subsequent game can be thought of as the noise model in Blume et al. (2007). We show that the welfare improvement is possible by the Sender's commitment to the fully concealing or the convexified signaling devices based on results in Blume et al. (2007). On the other hand, the equilibrium structure in the subgame after the semi-fully revealing signaling device resembles that of so-called type 1 equilibria in the honest model by Kim and Pogach (2014).

The remainder of this paper is structured as follows. In section 2, we formally introduce the model and provide a weak comparative statics of the model. In section 3, we focus on the uniform-quadratic case and present the welfare results. Section 4 concludes this paper.

---

[4]For example, Kamenica and Gentzkow (2011b) studies the setting where there are multiple Senders and Wang (2012) studies Bayesian persuasion with multiple Receivers. For other variations, Alonso and Câmara (2013) studies the Bayesian persuasion under the situation where a Sender and a Receiver have a different prior rather than common prior. Recently, Kolotilin (2014) studies the Bayesian persuaion model with a privately informed Receiver.

# 2 The Model

There are two players: a Sender (S) and a Receiver (R). Two players have a common prior distribution, $F(\omega)$, on the state space, $\Omega$. Assume that $\Omega$ is compact and $F(\omega)$ has a density which is positive everywhere.

The game proceeds as follows. The game starts with S's commitment to a signaling device. A signaling device, $\{\pi(\cdot|\omega)\}_{\omega \in \Omega}$, is a set of conditional distributions over a signal realization space, $M$, where $M$ is a set containing $\Omega$ (i.e. $\Omega \subset M$). We denote a signaling device by $\pi$. Importantly, S commits to a $\pi$ *before* he learns the realization of the true state. Once S chooses a $\pi$, he publicly announces his commitment to the $\pi$.

After his commitment to a $\pi$ is announced, S learns the true state and also if his commitment is binding or not. With probability $\alpha > 0$, S learns that his commitment is binding. We denote S whose commitment is binding by $S_b$. Since his commitment is binding, $S_b$ plays a simple role as a "truthful mediator" between $\pi$ and R: if the true state is $\omega'$, he observes a signal realization, $m \in M$, from the $\pi(m|\omega')$ and truthfully reports it to R.

With the other probability, $1 - \alpha > 0$, S learns that his commitment is not binding. We denote S whose commitment is not binding by $S_c$. Since his commitment is not binding, $S_c$ does not need to be a truthful mediator as $S_b$ does. Instead, $S_c$ can send any $m \in M$ in his discretion: he can discard a signal realization, $m \in M$, from a $\pi$ and report any $m \in M$ that he wants.[5] We denote a strategy of $S_c$ by $\sigma(\cdot) : \Pi \times \Omega \to \Delta M$. Especially, when we consider a specific subgame after S's commitment to a $\pi'$, $S_c$'s strategy in that subgame is reduced to a mapping from $\Omega$ to $\Delta M$. That is, in that subgame, if the true state is $\omega'$, $S_c$ who employs $\sigma'(\cdot)$ reports $m \in M$ according to his discretion, $\sigma'(m|\omega'; \pi')$, where $\sigma'(m|\omega'; \pi')$ is a distribution over $M$ given that the true state is $\omega'$.

We refer reports, $m \in M$, sent by $S_c$ to cheap-talk messages and reports, $m \in M$, sent by $S_b$ (or a $\pi$) to signals. We also assume that the value of $\alpha$ (the probability that S's commitment is binding) is common knowledge for both S and R.

Finally, R observes a report, $m \in M$, from S (either $S_c$ or $S_b$). R is aware of the possibility that the report is sent by $S_c$ and knows the probability of it, $1 - \alpha$. But R cannot tell if the report is sent by $S_c$ or $S_b$ (or putting it differently, she cannot distinguish the cheap-talk messages from the signals generated from a $\pi$.) Given a report from S and the value of $\alpha > 0$, R makes an inference on the true state, $\omega \in \Omega$. Then, R takes an action, $a \in A$, based on her updated beliefs, where $A$ is any set that contains $\Omega$. Once R takes an action, $a \in A$, both

---

[5]Note that this is a restriction. We restrict the set of signals that $S_c$ can use to be $M$ on which a signaling device is constructed. However, one can easily imagine a situation where $S_c$ can send a signal that is not in $M$. In that case, that signal will immediately tell R that she is hearing from $S_c$.

players' payoffs are realized and the game ends.

Both players' payoffs depends on both the action taken by R, $a \in A$, and the true state, $\omega \in \Omega$. We denote S's and R's payoffs functions by $U^S(a, \omega)$ and $U^R(a, \omega)$, respectively.

## 2.1  Equilibrium

We use the Perfect Bayesian Equilibrium as the solution concept. (Henceforth, an equilibrium means a Perfect Bayesian Equilibrium). In the model, S chooses a $\pi \in \Pi$ and publicly announces his choice. Hence, a $\pi$ chosen by S induces a subgame. We denote a subgame after S's commitment to a $\pi'$ by $\Gamma(\pi')$. For the sake of sequential rationality, we work backward: we start by defining an equilibrium in a subgame and then define an equilibrium of the model.

Consider a subgame after a $\pi'$, $\Gamma(\pi')$. An equilibrium in $\Gamma(\pi')$ consists of $S_c$'s strategy, R's strategy, and the set of posteriors of R. Note that we do not need to consider $S_b$'s strategy in $\Gamma(\pi')$ since $S_b$ behaves according to his commitment. In $\Gamma(\pi')$, we denote $S_c$'s strategy by $\sigma(\pi') := \{\sigma(m|\omega, \pi')\}_{\omega \in \Omega}$, R's strategy by $a(m; \pi')$, and the set of R's posterior by $\mathcal{M}(\sigma(\pi'), \pi')$. Formally, an equilibrium in $\Gamma(\pi')$ is a triple, $(\sigma^*(\pi'), a^*(m; \pi'), \mathcal{M}(\pi', \sigma^*(\pi')))$ such that

(1) for any posterior $\mu(\omega|m; \pi') \in \mathcal{M}(\pi', \sigma^*(\pi'))$,

$$\mu(\omega|m; \pi') = \frac{(\alpha\pi'(m|\omega) + (1-\alpha)\sigma^*(m|\omega, \pi'))f(\omega)}{\int_\Omega (\alpha\pi'(m|\omega) + (1-\alpha)\sigma^*(m|\omega, \pi'))f(\omega)d\omega},$$

(2) given any posterior, $\mu(\omega|m; \pi') \in \mathcal{M}(\cdot)$, $a^*(m; \pi') = \underset{a \in A}{\arg\max} \int_\Omega U^R(a, \omega)d\mu(\omega|m; \pi')$,

(3) for all $\omega \in \Omega$, if $\sigma^*(m'|\omega, \pi') > 0$ for a $m'$, $m' = \underset{m \in M}{\arg\max} U^S(a^*(m; \pi'), \omega)$.

In words, at an equilibrium in $\Gamma(\pi')$, R needs to form a posterior according to the Bayes' rule, R's strategy should maximize R's expected payoffs given posteriors, and, finally, $S_c$'s strategy should be a best response to R's strategy. Now we define an equilibrium of the model as follows:

An equilibrium of the model is a quadruple, $(\pi^*, \sigma^*(\pi), a^*(m; \pi), \mathcal{M}(\pi, \sigma^*(\pi)))$, such that

(1) for any $\pi \in \Pi$, $(\sigma^*(\pi), a^*(m; \pi), \mathcal{M}(\pi, \sigma^*(\pi))$ constitutes an equilibrium in $\Gamma(\pi)$,

(2) $\pi^*$ is a maximizer of ex ante payoffs of S: $\pi^* = \underset{\pi \in \Pi}{\arg\max} EU^S(\pi)$, where

$$EU^S(\pi) \;\; = \;\; \alpha \left[ \int_\Omega \left( \int_M U^S(a^*(m, \pi), \omega)d\pi(m|\omega) \right) dF(\omega) \right]$$

$$+(1 - \alpha) \left[ \int_\Omega \left( \int_M U^S(a^*(m, \pi), \omega) d\sigma^*(m|\omega, \pi) \right) dF(\omega) \right].$$

The *ex ante* payoffs that S can get by choosing a $\pi$ is simply a weighted average of the expected payoffs of $S_c$ and $S_b$ at an equilibrium of $\Gamma(\pi)$: the term after $\alpha$ is $S_b$'s expected payoffs and the term after $1 - \alpha$ is $S_c$'s expected payoff. At an equilibrium of the model, by choosing $\pi^*$, S decides to be in the subgame, $\Gamma(\pi^*)$, that gives him the highest payoffs.

## 2.2 Weak Comparative Statics

Now we state the first main finding which is the weak comparative statics analysis on S's welfare.

**Proposition 1.** *The Sender is **weakly** better off as the commitment probability, $\alpha$, increases: if there is an equilibrium in the model with $\alpha_0$, the Sender can always reproduce that equilibrium outcome in the model with $\alpha > \alpha_0$.*

*Proof.* Consider the model with a commitment probability, $\alpha_0$. Suppose that there is an equilibrium and denote the equilibrium signaling device by $\pi^*$. Now consider the equilibrium path. After the Sender commits to $\pi^*$, $S_c$ is in the subgame after $\pi^*$. Denote $S_c$'s equilibrium strategy in the subgame after $\pi^*$ by $\sigma^*(\pi^*) := \{\sigma^*(m|\omega; \pi^*)\}_{\omega \in \Omega}$. In addition, denote the Receiver's equilibrium strategy in the subgame after $\pi^*$ by $a^*(m; \pi^*) : M \to A$, where $M$ is the signal realization space. The triple, $(\pi^*, \sigma^*(\pi^*), a^*(m; \pi^*))$, determines the equilibrium outcome of the model with $\alpha_0$. At this equilibrium, any $m \in M$ such that $\pi^*(m|\omega) > 0$ or $\sigma^*(m|\omega) > 0$ induces the Receiver's posterior belief according to Bayes' Rule:

$$\mu_0(\omega|m) = \frac{(\alpha_0 \pi^*(m|\omega) + (1 - \alpha_0)\sigma^*(m|\omega))\, dF(\omega)}{\int_\Omega (\alpha_0 \pi^*(m|\omega) + (1 - \alpha_0)\sigma^*(m|\omega))\, dF(\omega)}.$$

Then, $a^*(m; \pi^*)$ is the maximizer of $\int_\Omega U^R(a, \omega) d\mu_0(\omega|m)$. Finally, $\sigma^*(\pi^*)$ satisfies $S_c$'s incentive compatibility condition. That is, for all $\omega \in \Omega$, if $\sigma^*(m'|\omega) > 0$,

$$m' = \arg\max_{m \in M} U^S(a^*(m; \pi^*), \omega).$$

Now consider the model with a higher commitment probability, $\alpha > \alpha_0$. Consider the following signaling device, $\pi^c = \{\pi^c(m|\omega)\}_{\omega \in \Omega}$: for all $\omega \in \Omega$,

$$\pi^c(m|\omega) := \frac{\alpha_0}{\alpha}\pi^*(m|\omega) + (1 - \frac{\alpha_0}{\alpha})\sigma^*(m|\omega; \pi^*).$$

Suppose that the Sender commits to $\pi^c$. The Sender's commitment to $\pi^c$ induces the subgame after $\pi^c$. In the subgame after $\pi^c$, let $S_c$ employs $\sigma^*(\pi^*)$. Then, the triple, $(\pi^c, \sigma^*(\pi^*), a^*(m; \pi^*))$, constitutes an equilibrium in the subgame after $\pi^c$.

To see this, consider the Receiver's incentive compatibility first. Given $(\pi^c, \sigma^*(\pi^*))$, the Receiver's posterior belief about $\omega$ conditional on observing $m \in M$ is

$$
\begin{aligned}
\mu(\omega|m) &= \frac{(\alpha\pi^c(m|\omega) + (1-\alpha)\sigma^*(m|\omega))\, dF(\omega)}{\int_\Omega (\alpha\pi^c(m|\omega) + (1-\alpha)\sigma^*(m|\omega))\, dF(\omega)} \\
&= \frac{\left(\alpha(\frac{\alpha_0}{\alpha}\pi^*(m|\omega) + (1-\frac{\alpha_0}{\alpha})\sigma^*(m|\omega)) + (1-\alpha)\sigma^*(m|\omega)\right) dF(\omega)}{\int_\Omega \left(\alpha(\frac{\alpha_0}{\alpha}\pi^*(m|\omega) + (1-\frac{\alpha_0}{\alpha})\sigma^*(m|\omega)) + (1-\alpha)\sigma^*(m|\omega)\right) dF(\omega)} \\
&= \mu_0(\omega|m).
\end{aligned}
$$

Thus, the Receiver's inference about $\omega$ after observing $m \in M$ in this subgame is exactly same as her inference at the equilibrium of the model with $\alpha_0$. Accordingly, the Receiver's strategy, $a^*(m; \pi^*)$, is optimal given $(\pi^c, \sigma^*(\pi^*))$.

Now only thing left to check is $S_c$'s incentive compatibility. Since the Receiver employs $a^*(m; \pi^*)$ which is the same strategy in the equilibrium of the model with $\alpha_0$, $\sigma^*(\pi^*)$ is a best response of $S_c$ in the subgame after $\pi^c$ as it is at the equilibrium of the model with $\alpha_0$.

Hence, the triple, $(\pi^c, \sigma^*(\pi^*), a^*(m; \pi^*))$, indeed constitutes an equilibrium in the subgame after $\pi^c$ in the model with the higher commitment probability, $\alpha > \alpha_0$. The Sender can achieve this outcome by committing to $\pi^c$ and employing $\sigma^*(\pi^*)$ when the Sender learns his commitment is not binding.

Note that $(\pi^*, \sigma^*(\pi^*), a^*(m; \pi^*))$ and $(\pi^c, \sigma^*(\pi^*), a^*(m; \pi^*))$ yield the same outcome: Both strategy profiles yield the same set of posteriors, the same distribution over the set of posteriors, and the same distribution over equilibrium actions. Accordingly, both strategy profiles yield the same *ex ante* payoffs for the Sender.

The Sender may be better off but cannot be worse off with a higher commitment probability: given a commitment probability, the Sender can always obtain an *ex ante* payoffs that he could have obtained at any equilibrium of the model with any lower commitment probability. □

Proposition 1 simply states that the *ex ante* payoffs of S *cannot* decrease as the commitment probability, $\alpha$, increases. In the model with a $\alpha$, S can achieve any equilibrium outcome in models with any $\alpha_0 < \alpha$.

Note that this weak comparative statics result does not necessarily hold for R. Consider two models with two different commitment probabilities, $\alpha$ and $\alpha_0$, where $\alpha > \alpha_0$. If S can achieve a strictly better equilibrium outcome in the model with $\alpha$ than with $\alpha_0$, S will choose that outcome in the model with $\alpha$. But it is unclear if R is also better off at the outcome chosen by S in the model with $\alpha$ compared to the outcome that S would choose in the model with $\alpha_0$.

However, there are cases that the weak comparative statics result also holds for R.

**Corollary 1.** *If $U^S(a,\omega) = -(\omega + b - a)^2$ and $U^R(a,\omega) = -(\omega - a)^2$, both the Sender and Receiver are **weakly** better off as the commitment probability, $\alpha$, increases.*

*Proof.* One can easily show that $EU^S = EU^R - b^2$ should hold at any equilibrium when $U^S(a,\omega) = -(\omega + b - a)^2$ and $U^R(a,\omega) = -(\omega - a)^2$, where $EU^i$ is $i$'s *ex ante* payoffs for $i = S, R$. By Proposition 1, $EU^S$ is weakly increasing in $\alpha$. Hence, $EU^R = EU^s + b^2$ is also weakly increasing in $\alpha$. □

It is a well-known fact that both players' *ex ante* interests are perfectly aligned if both players' preferences are represented by a quadratic loss function.[6] Corollary 1 is an immediate consequence of this well-known fact and Proposition 1.

## 3   The Welfare under Partial Commitment in the Uniform-Quadratic Case

In this section, we study the partial-commitment model in a specific setup, the uniform-quadratic framework, which is widely used in the cheap-talk literature. From now on, we assume that the state space, $\Omega$, is the unit interval, $[0, 1]$, and the common prior on $\Omega = [0, 1]$ is the uniform distribution. Furthermore, we assume that both players' preferences are represented by a quadratic loss function: $U^S(a, \omega; b) = -(\omega + b - a)^2$ and $U^R(\omega, a) = -(\omega - a)^2$, where $b > 0$ is term for S's bias. Thus, while R wants to match her action to the true state, S wants R to take the action that is equal to the true state plus his bias.

We focus on the welfare comparison of three different cases, the no-commitment, the partial-commitment, and the full-commitment cases. Then, we show that both S and R are *strictly* better off as they move from the no-commitment to the partial-commitment to the full-commitment cases, which is true for any level of S's bias (even when $b$ is arbitrarily high).

Finally, note that $EU^S = EU^R - b^2$ holds at any equilibrium in this uniform-quadratic case. Thus, when we discuss the welfare, it is enough to discuss S's *ex ante* payoffs. Furthermore, we will interchangeably use "the more (less) information transmission" and "the welfare improvement (loss)" since $EU^R$ increases in the amount of information transmitted and $EU^S$ is simply $EU^R - b^2$.

---

[6]It is worth mentioning that there is more general condition (called condition (M) in Crawford and Sobel (1982)) which ensures the perfect alignment of both players' *ex ante* preferences. Under the condition (M), S *ex ante* prefers an equilibrium to the other if and only if R *ex ante* prefers the former to the latter. Hence Corollary 1 also holds under the condition (M).

## 3.1 Benchmarks

Note that if $\alpha = 0$, the model is equivalent to the leading example in Crawford and Sobel (1982), and if $\alpha = 1$, the model corresponds to a special case of the Bayesian Persuasion model by Kamenica and Gentzkow (2011a). In this subsection, we summarize results in these two cases that will be used as benchmarks in later sections.

### 3.1.1 Cheap Talk: Case when $\alpha = 0$

First, suppose $\alpha = 0$. Then, all reports that R observes are sent by $S_c$. Hence, S's commitment to a $\pi$ dose not play any role and what matters is $S_c$'s strategy. Furthermore, R is aware of it. Thus, we return back to the uniform-quadratic case in Crawford and Sobel (1982) (henceforth, CS model). We summarize the results from the CS model and define the best equilibrium of it in the following.

We define the best equilibrium in the CS model (henceforth, BECS, the abbreviation for the Best Equilibrium in the CS model) as the equilibrium that gives the highest *ex ante* payoffs for both S and R. In the CS model, every equilibrium is characterized by a partition of the state space, where a partition reflects how much R learns about the true state at an equilibrium. The number of elements in an equilibrium partition is always finite and the maximum number of elements in an equilibrium partition is determined by S's bias $b > 0$. Given any $b > 0$, the equilibrium that induces the partition with the maximum number of elements (the finest partition) is the best equilibrium among all other equilibria

For any $b \geq 1/4$, the finest partition that can be induced at an equilibrium has only one element in it. Thus, the only equilibrium outcome is the pooling outcome. Accordingly, the pooling equilibrium is the BECS in this case. If $b < 1/4$, the maximum number of elements in an equilibrium partition is equal to $[-\frac{1}{2} + \frac{1}{2}\left(1 + \frac{2}{b}\right)^{1/2}]$, where $[I]$ means the smallest integer that is greater than or equal to $I$. Thus, the BECS is the equilibrium that induces the partition with $[-\frac{1}{2} + \frac{1}{2}\left(1 + \frac{2}{b}\right)^{1/2}]$ numbers of elements in it. For example, if $b = 1/8$, we have $[-\frac{1}{2} + \frac{1}{2}\left(17\right)^{1/2}] = [1.56...] = 2$. And given $b = 1/8$, the BECS is the equilibrium that induces the partition with 2 elements in it. We will use the BECS under a given value of $b > 0$ as our benchmark in the later section.

### 3.1.2 Bayesian Persuasion: Case when $\alpha = 1$

Now suppose $\alpha = 1$. Then, R knows that any $m \in M$ she observes is sent by $S_b$. Thus, the model collapses to the Bayesian Persuasion model by Kamenica and Gentzkow (2011a). In this case, a $\pi$ to which S commits fully determines an outcome since there is no interruption by $S_c$. Thus, at an equilibrium, S commits to a signaling device that maximizes his *ex*

10

*ante* payoffs among all possible signaling devices. It turns out that there exists the unique equilibrium outcome and the equilibrium outcome is independent of $b > 0$. Furthermore, the unique equilibrium outcome is also a first-best outcome in the uniform-quadratic framework. The following Remark 1 formally states the result.

**Remark 1.** *Suppose $\alpha = 1$. Then, for any $b > 0$, there exists the unique equilibrium outcome that is characterized by full revelation of the true states to the Receiver. Thus, at any equilibrium, the Sender's expected payoffs is $-b^2$ and the Receiver's expected payoffs is 0, which is a first-best outcome.*

The detailed proof of Remark 1 can be found in the Appendix. In the proof we simply use the fact that $EU^S = EU^R - b^2$ should hold at any equilibrium (or as far as R maximizes her expected payoffs given her posteriors). It is easy to see that $EU^R$ increases as there is less uncertainty about the true state; R can take better action as she is more certain about the true state. Since $EU^S = EU^R - b^2$, S can achieve the maximum *ex ante* payoffs by fully resolving R's uncertainty. Thus S chooses a $\pi$ that reveals true states to R. An example of such a signaling device is the one that directly tells the exact true state as follows: for all $\omega \in [0, 1]$,

$$\begin{aligned} \pi^f(m|\omega) &= 1 \quad \text{if} \ \ m = \omega, \\ &= 0 \quad \text{if} \ \ m \neq \omega. \end{aligned}$$

Given this signaling device, R learns the exact true state and takes the action that is equal to the true state, $a^*(m = \omega) = \omega$. Finally, R gets the payoffs of 0 and S gets the payoffs of $-b^2$. And it is a first best outcome since the sum of both players' payoffs cannot be higher than $-b^2$.

It is not hard to imagine other equilibria in which S chooses a signaling device that reveals the truth in a different way than the example above. But any equilibrium signaling device should be a signaling device that fully reveals the true state. Thus, all equilibrium outcomes should be the same as in the example above. Hence Remark 1 establishes the unique benchmark that will be used in the later section.

## 3.2  The Welfare Loss under Partial Commitment

In this subsection we compare the welfare under partial commitment with that under full commitment: both S and R obtain a *strictly* higher *ex ante* payoffs under full commitment than partial commitment.

From Remark 1, we see that the unique outcome is attained under the full-commitment assumption. The following Remark 2 shows that the unique outcome in Remark 1 cannot be

an equilibrium outcome if $\alpha < 1$ (i.e. if S has an opportunity to do cheap-talk with a positive probability).

**Remark 2.** *If $\alpha < 1$, the unique outcome under full commitment cannot be achieved at any equilibrium. i.e. if $\alpha < 1$, there is no equilibrium that has the outcome characterized by the full revelation of the true states to the Receiver.*

*Proof.* Suppose that the first-best outcome in the full-commitment case can be achieved at an equilibrium in the partial-commitment model. At the equilibrium, R knows the exact true state. Thus, for any signal R gets, R can precisely infer the true state and take the action that is equal to the true state. Consider two different states, $\omega'$ and $\omega' + b$. R's optimal actions at these states are $\omega'$ and $\omega' + b$ respectively. Denote signals that are sent in these state, $\omega'$ and $\omega' + b$, by $m'$ and $m' + b$. Then, $S_c$ who is in state $\omega'$ can be better off by sending $m' + b$. Thus, it cannot be an equilibrium. $\square$

Remark 2 partially characterizes equilibrium outcomes in a model with $\alpha < 1$. That is, at any equilibrium in a model with $\alpha < 1$, R should be uncertain about some states of the world. Recall the monotonicity between R's *ex ante* payoffs and the uncertainty about the true state; the more uncertainty is, the lower $EU^R$ is. Hence, Remark 2 implies that, at any equilibrium under partial commitment, R's payoffs should be strictly less than 0 that R would obtain under full commitment. Because $EU^S = EU^R - b^2$ should hold at any equilibrium, the loss in S's *ex ante* payoffs immediately follows. Hence, the welfare loss under partial commitment is established: both S and R are *strictly* worse off as S has an opportunity to do cheap-talk with any positive probability that is even arbitrarily small.

In a model with $\alpha < 1$, S can send a $m \in M$ in his discretion with a positive probability, which can be thought of as a data manipulation behavior.[7] Hence, it is intuitively clear that R will be suffered when S can manipulate data with a positive probability. However, why does S also suffer from having data manipulation chance though it seems that S can benefit from having such a chance?

The mathematical formulation of two different cases clarifies the reason why S is worse off as he has a chance of data manipulation. The main difference between the full-commitment and partial-commitment cases is whether $S_c$ is "active" or not in the model. If $\alpha = 1$, $S_c$ is not active. In other words, S has only one future-self, $S_b$, who will behave exactly as S meant. Hence, S who is the "twin" of $S_b$ solves a simple maximization problem given that R chooses an action that maximizes her expected payoffs. That is, S (or $S_b$)'s problem is as

---

[7]Recall that a signaling device, $\pi$, can be thought of as an experiment and signals generated by a $\pi$ can be considered as data from the experiment.

follows:

$$\max_{\pi \in \Pi} \; EU^S = EU^R - b^2.$$

Now suppose that $\alpha < 1$. Since S assumes R's incentive compatibility as in $\alpha = 1$ case, his objective function will be the same as before, $EU^S = EU^R - b^2$. But, since $S_c$ is active in this case, S should take into account how $S_c$ will behave. Especially, at an equilibrium, S rationally expects that $S_c$ will behave in the way that maximizes his payoffs. That is, S solves the same maximization problem above *but* assuming $S_c$'s IC behavior:

$$\max_{\pi \in \Pi} \; EU^S = EU^R - b^2, \; \text{s.t. } S_c \text{'s IC behavior.}$$

Hence, as soon as we decrease $\alpha$ from 1, we immediately bring some constraints into S's maximization problem. At the first glance, having data manipulation chance seems making S be more "free" *but* it actually makes S be more "restrictive".[8]

By looking at two different maximization problems, it is obvious that S cannot be better off as $\alpha$ decreases, which we already saw in Proposition 1. However, we cannot tell if S is *strictly* worse off or not: if the constraints that are brought into are slack, S can sustain the same payoffs as before. Remark 2 shows that the newly introduced constrains are not slack. Thus, S is *strictly* worse off and so is R as $\alpha$ is drifting away from 1.

## 3.3 The Welfare Improvement under Partial Commitment

In this subsection, we show that both S and R obtain *strictly* higher *ex ante* payoffs under partial commitment than under no commitment. This welfare improvement result holds for any pair of parameters, $(\alpha, b) \in (0, 1) \times (0, \infty)$. In other words, for any commitment probability, $\alpha \in (0, 1)$, and for any bias, $b > 0$, having S partially commit to some signaling devices *strictly* improves both S's and R's *ex ante* payoffs.

To establish this result, we will consider three types of signaling devices: a semi-fully revealing, a fully concealing, and a convexified signaling devices. Then, we show that both S and R can be strictly better off via S's partial commitment to one of these signaling device. Since S publicly announces the signaling device to which he commits, S's announcement induces a subgame. We will construct an equilibrium in each subgame following the announcement of S's commitment to each signaling device for all parameter pairs, $(\alpha, b) \in (0, 1) \times (0, \infty)$. Then, we show that, at the constructed equilibrium, both S's and R's *ex ante* payoffs are strictly higher than those in a BECS under no commitment (the best equilibrium in the cheap-talk model).

---

[8]Note that this is also the key intuition we saw in Proposition 1.

### 3.3.1 The Welfare Improvement via a Fully Concealing Signaling device

A fully concealing signaling device is a signaling device that fully conceals information about the state of world. There are many signaling devices that fully conceals information. Among those signaling devices, we choose a signaling device that is similar to the noise channel in Blume, Board, and Kawamura (2007). The fully concealing signaling device is defined as follows:

**Definition 1.** *The fully concealing signaling device is defined as* $\pi^n := \{\pi(m|\omega)\}_{\omega \in [0,1]}$, *where for each* $\omega \in [0,1]$,

$$\pi(m|\omega) \text{ is the uniform distribution on } [0,1].$$

Note that $\pi^n$ does not convey any information about $\omega \in [0,1]$ as we meant by its definition. One may think of $\pi^n$ as a simple noise generator; Observing a $m \in M$ generated by $\pi^n$ is same as observing a "noise" that is randomly drawn from the uniform distribution over $M$, where $M := [0,1]$.

Now suppose that S commits to $\pi^n$ at the start of the whole game. S publicly announces his commitment to $\pi^n$, which induces a subgame. In this subgame, when R observes a $m \in [0,1]$, she knows that, with probability $\alpha > 0$, $m$ is sent by $S_b$ (or generated by $\pi^n$). And $S_c$ also knows this: $S_c$ knows that, if he sends a $m$, R would interpret the $m$ as a noise from $\pi^n$ with probability $\alpha > 0$. Hence, this subgame coincides with the noise model in Blume et al. (2007).[9] In this subgame, a given commitment probability, $\alpha \in (0,1)$, can be thought of as the noise level, $\epsilon \in (0,1)$, in Blume et al. (2007).

Following Blume et al. (2007), in this subgame, we can construct an equilibrium in which $S_c$ employs so-called *front-loading* strategy. We call such an equilibrium a front-loading equilibrium. A front-loading equilibrium is characterized by a partition of the state space as a BECS is in the no-commitment case (the cheap-talk model). Given a $\alpha \in (0,1)$ and a $b \in (\frac{1}{2N^2}, \frac{1}{2(N-1)2})$ for $N = 2, 3, ..$, we can construct a front-loading equilibrium in which the state space is partitioned in $N$-interval and $S_c$ employs a *front-loading* strategy as follows:

> if $\omega \in [0, \omega_1]$, randomize uniformly on $[0,1] \backslash \{m_2, m_3, ..., m_N\}$,
>
> if $\omega \in (\omega_{i-1}, \omega_i]$, send $m_i$ for $i = 2, ..., N$, where $\omega_N = 1$.

We denote $S_c$'s *front-loading* strategy by $\sigma_{fl}$. One can clearly see why it is called a front-loading strategy. Note that almost all messages, $[0,1] \backslash \{m_2, m_3, ...m_N\}$, are loaded at the "front" types, $[0, \omega_1]$.

---

[9]In their model, when the sender sends a massage, the message is delivered to a receiver only with a positive probability. With the other probability, the message sent by the sender is replaced by a noise from a given noise generating process and the receiver observes the noise.

At a front-loading equilibrium, R's inference is as follows. First, when R observes $m \in [0,1] \setminus \{m_1, m_2, ... m_N\}$, R believes that, with probability $\alpha$, $m$ is merely drawn from the uniform distribution on $[0,1]$ (or $m$ is sent by $S_b$) and, with probability $1 - \alpha$, $m$ is sent by $S_c$'s types who are in $[0, \omega_1]$. Thus, R's optimal action is a weighted average of two midpoints, $\frac{1}{2}$ of $[0,1]$ and $\frac{\omega_1}{2}$ of $[0, \omega_1]$.[10] Secondly, when R observes $m_i \in \{m_2, m_3, ..., m_N\}$, R believes that $m_i$ is sent by $S_c$'s types who are in $(\omega_{i-1}, \omega_i]$ with probability 1 since the event that $m_i$ is generated by $\pi^n$ (or sent by $S_b$) has measure zero. Thus, R's optimal action is merely the midpoint of the interval, $(\omega_{i-1}, \omega_i]$.

It is worth noting that a "meaningful" front-loading equilibrium exists only if $0 < b < 1/2$. If $b \geq 1/2$, the only incentive compatible $\sigma_{fl}$ for $S_c$ is the one that has no partition structure (i.e. all types uniformly randomizes on $[0,1]$). Consequently, a front-loading equilibrium when $b \geq 1/2$ is the same as the pooling equilibrium. Hence, it is immediate that, if $b \geq 1/2$, the front-loading equilibrium yields the same outcome as a BECS in the no-commitment case.

According to Blume et al. (2007), if $b < 1/2$, a front-loading equilibrium can yield an outcome in which both S and R obtain strictly higher *ex ante* payoffs than in a BECS outcome. But this welfare improvement is only possible if the noise level, $\epsilon \in (0,1)$, is less than an upper bound, $\bar{\epsilon}(b)$, where the upper bound, $\bar{\epsilon}(b)$, is a function of $b$ and its definition is relegated to Appendix for expositional convenience. This result can be interpreted in our frame work (by treating the noise level, $\epsilon$, as the commitment probability, $\alpha$). The following proposition states the interpretation of their result in our context.

**Proposition 2.** *If $\alpha \in (0, \bar{\epsilon}(b))$ and $b \in (0, 1/2) \setminus \{b = \frac{1}{2N^2}$ for $N = 2, 3, ...\}$, the Sender's partial commitment to $\pi^n$ **strictly** improves the ex ante payoffs for both the Sender and Receiver compared to the no-commitment case. Otherwise, the Sender's partial commitment to $\pi^n$ does not improve the ex ante payoffs for both players compared to the no-commitment case.*

For the detailed proof, refer to the proof for Proposition 9 in Blume et al. (2007).

Figure 1 summarizes Proposition 2. The shaded area in Figure 1 indicates the set of parameter pairs, $(\alpha, b)$, in which S's partial commitment to $\pi^n$ results in a welfare improvement. The unshaded area indicates the set of parameter pairs in which S's partial commitment to $\pi^n$ cannot improve the welfare along with the construction of a front-loading equilibrium.

---

[10] The following is R's optimal action, $a^*(m)$, that is calculated with the appropriate weight:

$$a^*(m) = \frac{\alpha(1-0)\frac{(1+0)}{2} + (1-\alpha)(\omega_1 - 0)\frac{\omega+0}{2}}{\alpha(1-0) + (1-\alpha)(\omega_1 - 0)}$$

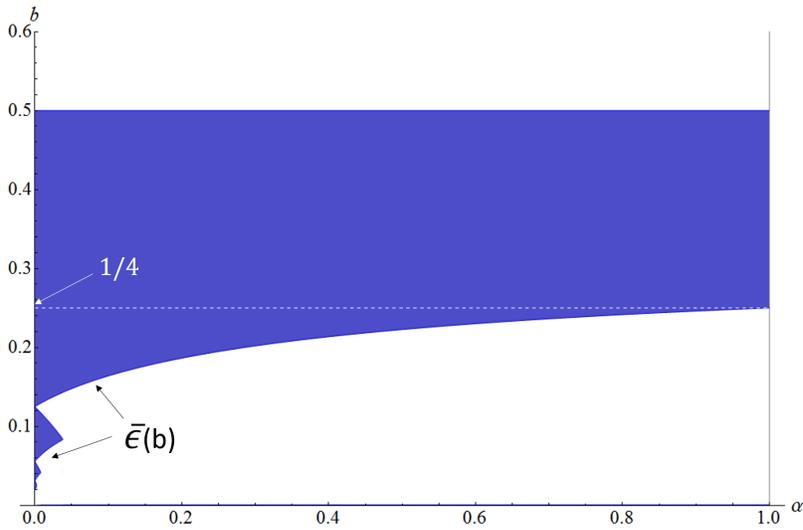$$= \frac{\alpha\frac{1}{2} + (1-\alpha)\omega_1\frac{\omega_1}{2}}{\alpha + (1-\alpha)\omega_1}.$$

Figure 1: Summary of Proposition 2

S's commitment to $\pi^n$ leads both S and R to the subgame in which noise is generated with a positive probability as in the noise model by Blume et al. (2007). As discussed in their paper, the presence of noise has two effects on the welfare. The presence of noise "loosens" $S_c$'s incentive compatibility constraints. Due to this loosening effect[11], we can construct a front-loading equilibrium that has either *finer* or *evener* partition than a BECS, which implies *potentially* more information transmission (thus a welfare improvement) at a front-loading equilibrium than at a BECS. However, the noise by itself hinders the information transmission, which implies there are some negative effects that may offset the positive effect above. Blume et al. (2007) show that, when the noise level is less than some upper bound (i.e. $0 < \alpha < \bar{\epsilon}(b)$) and $b \in (0, 1/2) \backslash \{b = \frac{1}{2N^2} \ for \ N = 2, 3, ...\}$, the positive effect of the noise dominates the negative effect. Hence, we can have welfare improvements in the shaded area.

However, once the noise level is too high (i.e. $\alpha > \bar{\epsilon}(b)$), the negative effect of the noise defeats the positive effect and we cannot have the welfare improvement results. This appears as the unshaded area under $b = 1/2$ in Figure 1. Furthermore, if $b \geq 1/2$, constructing front-loading equilibria is meaningless: for any $(\alpha, b) \in (0, 1) \times [1/2, \infty)$, any front-loading equilibrium is the pooling equilibrium. Thus, we cannot have welfare improvements via $\pi^n$ when $b \geq 1/2$ and it appears as the unshaded area above $b = 1/2$ in Figure 1.

---

[11]In their paper, they refer this effect to *strategic* effect.

### 3.3.2 The Welfare Improvement via Convexified Signaling Devices

In this section, we explicitly make use of the proof technique we used to prove Proposition 1. We start by constructing a new signaling device that we call a convexified signaling device. The construction (which is similar to that in the proof of Proposition 1) is as follows. First, given a parameter pair $(\alpha, b) \in (0, 1) \times (0, \infty)$, choose a $p \in (0, 1]$ and multiply $\alpha$ to the chosen $p$. Then, find the front-loading strategy that constitutes the front-loading equilibrium as in the previous section *but* with the noise level of $\alpha \cdot p$ (instead of just $\alpha > 0$). Denote this front-loading strategy by $\sigma_{fl}(m|\omega; \alpha p)$. Then, construct a new signaling device by making the convex combination of $\sigma_{fl}(m|\omega; \alpha p)$ and $\pi^n$ with the weight $p$ on $\pi^n$: for all $\omega \in [0, 1]$, $p\pi^n(m|\omega) + (1-p)\sigma_{fl}(m|\omega; \alpha p)$. Denote this convex combination by $\pi^c(p)$. We refer this convex combination, $\pi^c(p)$, to a convexified signaling device.

**Definition 2.** *A convexified signaling device is denoted by $\pi^c(p)$, where, for some $p \in (0, 1]$,*

$$\pi^c(p) := p\{\pi^n(m|\omega)\}_{\omega \in [0,1]} + (1-p)\{\sigma_{fl}(m|\omega; \alpha p)\}_{\omega \in [0,1]}.$$

Note that $\pi^c(p)$ is a valid signaling device since, given any $\omega \in [0, 1]$,

$$
\begin{aligned}
\int_M \pi^c(m|\omega; \alpha p)dm &= p\int_0^1 \pi^n(m|\omega)dm + (1-p)\int_0^1 \sigma_{fl}(m|\omega; \alpha p)dm \\
&= p \cdot 1 + (1-p) \cdot 1 = 1.^{12}
\end{aligned}
$$

A simple interpretation of $\pi^c(p)$ is as follows: with probability $p$, $\pi^c(p)$ generates signals according to $\pi^n$ and with probability $1-p$, it generates signals according to $\sigma_{fl}(m|\omega; \alpha p)$. Thus, $\pi^c(p)$ conceals the information with probability $p$ and, with probability, $1-p$, it conveys some information according to $\sigma_{fl}(m|\omega; \alpha p)$.

Note that, depending on which value we choose for $p \in (0, 1]$, we have a different $\pi^c(p)$. Thus, the construction above gives us a set of convexified signaling devices, $\{\pi^c(p) \ for \ p \in (0, 1]\}$. The other thing we need to notice is that, if $b \geq 1/2$, the above construction gives us the fully concealing signaling device, $\pi^n$. The reason is that, if $b \geq 1/2$, the only incentive compatible front-loading strategy is the one having no partition. Thus, for any chosen $p \in (0, 1]$, $\sigma_{fl}(m|\omega; \alpha p)$ is equal to $\pi^n$ and $\pi^c(p) = p\pi^n + (1-p)\pi^n = \pi^n$.

S's commitment to one of convexified signaling devices induces a subgame following the chosen signaling device. Suppose that S commits to $\pi^c(p)$ for some $p \in (0, 1]$. In this subgame, we can simply construct an equilibrium by making $S_c$ use $\sigma_{fl}(m|\omega, \alpha p)$ that is built in the chosen $\pi^c(p)$.

---

[12]For the sake of clarity, we abuse notation here: note that $\pi(m|\omega)$ is not a density but a distribution function.

Suppose $S_c$ uses $\sigma_{fl}(m|\omega;\alpha p)$ in the subgame following $\pi^c(p)$. Knowing the probability that $m \in M$ is sent according to $\pi^c(p)$ is equal to $\alpha$, R believes that signals are generated by $\pi^n$ with probability $\alpha p$. With the other probability, $1-\alpha p = \alpha\cdot(1-p)+1-\alpha$, R believes that signals are generated by $\sigma_{fl}(m|\omega;\alpha p)$.[13] Hence, R's optimal action is the same as the action chosen at the front-loading equilibrium in the noise model with the noise level of $\alpha p$. Now, the only thing we need to check is if $\sigma_{fl}(m|\omega;\alpha p)$ is incentive compatible for $S_c$. Note that $\sigma_{fl}(m|\omega;\alpha p)$ is chosen to be incentive compatible at the font-loading equilibrium in the noise model with the noise level of $\alpha p$ by the construction above. Hence, it is indeed an equilibrium in this subgame. We choose to use this equilibrium construction in every subgame after every convexified signaling device.

Before making the statement of the main finding in this section, we need to introduce the notion of optimal noise level that appears in Blume et al. (2007).

**Definition 3.** *The optimal noise level is the noise level that maximizes S's ex ante payoffs (and also R's ex ante payoffs) at a front-loading equilibrium in the noise model and is defined as follows:*

$$\epsilon^*(b) = \frac{(1-2b(N-1)^2)(1-2bN^2)}{4(N-1)N(b+b^2(N-1)N-1)} \quad if\ b \in (0,1/2)\backslash\{b = \frac{1}{2N^2}\ for\ N = 2,3,...\},$$

$$= 0 \quad if\ \{b = \frac{1}{2N^2}\ for\ N = 2,3,...\},$$

*where $N$ is the number of elements in the partition at a front-loading equilibrium given a $b \in (\frac{1}{2N^2}, \frac{1}{2(N-1)^2})$ for $N = 2,3,....$*

For example, suppose that $b \in (1/8, 1/2) \equiv (\frac{1}{2\cdot 2^2}, \frac{1}{2\cdot(2-1)^2})$. Then we can construct a front-loading equilibrium that partitions the state space, $[0,1]$, into 2-interval. Thus, $N = 2$. Then, the optimal noise level is the noise level that maximizes S's equilibrium payoffs at this front-loading equilibrium and is equal to $\epsilon^*(b) = \frac{(1-2b)(1-8b)}{8(b+2b^2+b-1)}$ for $1/8 < b < 1/2$. We refer the equilibrium outcome with the optimal noise level to the optimal noise outcome. Now we have Lemma 1.

**Lemma 1.** *If $(\alpha, b) \in [\epsilon^*(b), 1) \times (0, 1/2)\backslash\{b = \frac{1}{2N^2}\ for\ N = 2,3,...\}$, there exists a signaling device that implements the optimal noise outcome under partial commitment environment, and the signaling device is $\pi^c(p^*)$, where $p^* = \frac{\epsilon^*(b)}{\alpha}$. The optimal noise outcome can be achieved when the Sender commits to $\pi^c(p^*)$ and $S_c$ employs $\sigma_{fl}(m|\omega;\alpha p^*)$ in the subgame following $\pi^c(p^*)$. Finally, the Sender's ex ante payoffs at the optimal noise outcome is equal to $-\frac{1}{3}b(1-b) - b^2$.*

---

[13]With probability $\alpha(1-p)$ signals are generated by $\sigma_{fl}(m|\omega;\alpha p)$ that is built in $\pi^c(p)$ and with probability $1-\alpha$ signals are also generated by $\sigma_{fl}(m|\omega;\alpha p)$ since $S_c$ employs the strategy.

*Proof.* Suppose that $\alpha \geq \epsilon^*(b)$ and $b \in (0, 1/2)\backslash\{b = \frac{1}{2N^2} \; for \; N = 2, 3, ...\}$. Then S can have two cases: (1) $\alpha = \epsilon^*(b)$ or (2) $\alpha > \epsilon^*(b)$. In case (1), $p^* = 1$ and the optimal noise outcome can be achieved when S commits to $\pi^c(p^* = 1) = \pi^n$ and $S_c$ uses $\sigma_{fl}(m|\omega; \alpha)$ since $\alpha$ exactly coincides with the optimal noise level. In case (2), S can choose a $\pi^c(p)$ such that $\alpha p = \epsilon^*(b)$. Then, along with $S_c$'s use of $\sigma_{fl}(m|\omega; \alpha p)$, $\pi^c(p)$ can achieve the optimal noise outcome. For the Sender's payoffs from the optimal noise outcome, refer to Blume et al. (2007). $\qquad\square$

Lemma 1 provides a sufficient condition ($\alpha \geq \epsilon^*(b)$) for the existence of a convexified signaling device that can achieve the optimal noise outcome: the commitment probability should be greater or equal to $\epsilon^*(b)$ given a $b \in (0, 1/2)$.

With the set of convexified signaling devices, $\{\pi^c(p) \; for \; p \in (0, 1]\}$, S can control the frequency of noisy signals generated by $\pi^n$ that is built in a $\pi^c(p)$. For example, if S commits to $\pi^c(p = 0.2)$, the probability that noisy signals are generated by $\pi^n$ is $0.2 \times \alpha$ in the subgame following $\pi^c(p = 0.2)$. Namely, S can reduce a given frequency of noisy signals, $\alpha > 0$, to any level below $\alpha$, $0 < \alpha p \leq \alpha$. Hence, if $\alpha \geq \epsilon^*(b)$, S can choose to be in the "appropriate" subgame with the optimal noise level as he reduces $\alpha$ to $\epsilon^*(b)$ by committing to $\pi^c(p)$ with $p = \epsilon^*(b)/\alpha$.

However, S cannot increase the frequency of noisy signals: S's control on the noisy signals only works in one direction. Thus, if $\alpha < \epsilon^*(b)$, there does not exist a $\pi^c(p)$ that can achieve the optimal noise outcome; $S_c$ cannot magnify a given frequency of noisy signals, $\alpha$, to have the optimal noise level, $\epsilon^*(b)$.

Finally, as Blume et al. (2007) pointed out, S's *ex ante* payoff at the optimal noise equilibrium, $-\frac{1}{3}b(1 - b) - b^2$, is equal to the upper bound for S's *ex ante* payoffs that can be obtained at any equilibrium in the mediation model in Goltsman, Hörner, Pavlov, and Squintani (2009). It implies that, if the commitment probability is high enough (i.e. $\alpha \geq \epsilon^*(b)$), S's partial commitment to $\pi^c(p^*)$ with $p = \epsilon^*(b)/\alpha$ can achieve the best outcome in the mediation model in Goltsman et al. (2009).

With Lemma 1, we finally have the following Proposition 3.

**Proposition 3.** *For any $\alpha \in (0, 1)$ and $b \in (0, 1/2)\backslash\{b = \frac{1}{2N^2} \; for \; N = 2, 3, ...\}$, the Sender's partial commitment to $\pi^c(p)$ **strictly** improves the ex ante payoffs of both the Sender and Receiver compared to BECS in the no-commitment case:*

(a) *if $\alpha \geq \epsilon^*(b)$, the Sender's partial commitment to $\pi^c(p^*)$ can yield the optimal noise outcome that is strictly better than BECS outcome in the no-commitment case, where $p^* = \frac{\epsilon^*(b)}{\alpha}$,*

(b) if $\alpha < \epsilon^*(b)$, the Sender's partial commitment to any $\pi^c(p)$ for $p \in (0,1]$ can yield the front-loading equilibrium outcome that is strictly better than BECS outcome in the no-commitment case,
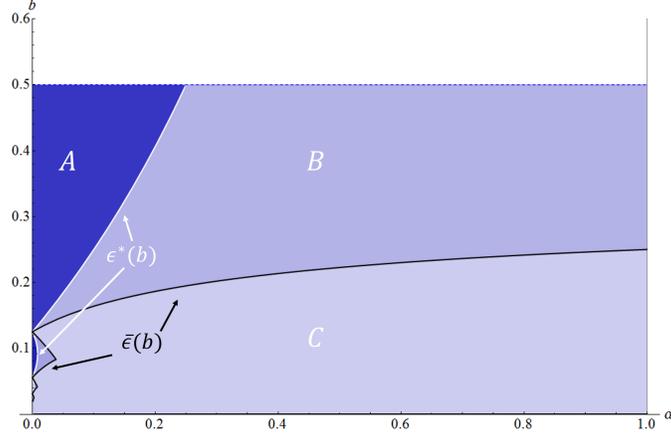


Figure 2: Summary of Proposition 3

Proposition 3 states that, as far as $b < 1/2$, there always exist some convexified signaling devices that *strictly* improve the both players' welfare for any given commitment probability, $\alpha \in (0,1)$.

Figure 2 summarizes Proposition 3. In Figure 2, we depict two functions, $\bar{\epsilon}(b)$ and $\epsilon^*(b)$. These two functions divide the relevant parameter space, $(\alpha, b) \in (0,1) \times (0,1/2)$, into 3 pieces. Depending on the value of $\alpha$, we have 3 cases: (1) if $\alpha < \epsilon^*(b)$, we are in area A, (2) if $\epsilon^*(b) \le \alpha < \bar{\epsilon}(b)$, we are in area B, and (3) if $\alpha \ge \bar{\epsilon}(b)$, we are in area C. First, in either area B or C, S's partial commitment to $\pi^c(p^*)$ can yield the optimal noise outcome as we see in Lemma 1. Both players' *ex ante* payoffs at the optimal noise outcome is unambiguously greater than those at any BECS that can be constructed in area B or C[14]. Thus, in area B or C, the welfare improvements can be achieved via S's partial commitment to $\pi^c(p^*)$ .

Secondly, in area A, S's partial commitment to any of $\pi^c(p)$ for $p \in (0,1]$ cannot yield the optimal noise outcome since $\alpha < \epsilon^*(b)$. However, note that $\epsilon^*(b) \le \bar{\epsilon}(b)$ in Figure 2. Thus, $\alpha < \bar{\epsilon}(b)$ in area A. As we saw in the previous section, if $\alpha < \bar{\epsilon}(b)$, we can construct a front-loading equilibrium that is Pareto-superior than BECS. Hence, S's partial commitment to a $\pi^c(p)$ for any $p \in (0,1]$ can yield a front-loading equilibrium that is Pareto-superior than any BECS in area A. Thus, even in region A, a welfare improvement is also possible via S's commitment to any convexified signaling device.

Finally, note that we cannot have a welfare improvement in the unshaded area in Figure

---

[14]To see this, refer to the appendix

20

2, $(\alpha, b) \in (0, 1) \times [1/2, \infty)$. Recall that any convexified signaling devices that are constructed when $b \geq 1/2$ is same as the fully concealing signaling device, $\pi^n$. Hence, when $b \geq 1/2$, S's commitment to any convexified signaling device is the same as the commitment to $\pi^n$. Again, recall that S's commitment to $\pi^n$ cannot improve the welfare when $b \geq 1/2$. Accordingly, S's partial commitment to any of $\pi^c(p)$ cannot improve the welfare in the unshaded area in Figure 2.

### 3.3.3 The Welfare Improvement via Semi-fully Revealing Signaling Devices

So far we have showed that, at the low-bias area, $(\alpha, b) \in (0, 1) \times (0, 1/2)$, S's partial commitment to either $\pi^n$ or $\pi^c(p)$ can strictly improve both players' *ex ante* payoffs compared to the no-commitment case. In this section, we provide the result that a welfare improvement is also possible even in the high-bias area, $(\alpha, b) \in (0, 1) \times [1/2, \infty)$.

It is helpful to remind us of two facts: (1) if $b \in [1/2, \infty)$, the only equilibrium that can be sustained in the no-commitment case is the pooling equilibrium and (2) both players' *ex ante* payoffs are strictly increasing in the amount of information transmitted. These two facts make our job easier: for a given $b \in [1/2, \infty)$ and $\alpha \in (0, 1)$, we only need to find a signaling device that induces a subgame having an equilibrium with a certain degree of information transmission. If there is such a signaling device for all $\alpha \in (0, 1)$ and all $b \in [1/2, \infty)$, we can establish the result that we desire.

It turns out that there are such signaling devices and those signaling devices belong to a specific class that we call semi-fully revealing signaling devices. A semi-fully revealing signaling device is defined as follows:

**Definition 4.** *A semi-fully signaling device is defined as $\pi^s(t) := \{\pi(m|\omega)\}$, where given some $t \in (0, 1)$,*

*(i) if $\omega \in [0, t]$,*

$$\pi(m|\omega) = 1 \quad \text{if } m = \omega,$$
$$= 0 \quad \text{if } m \neq \omega,$$

*(ii) if $\omega \in (t, 1]$,*

$$\pi(m|\omega) \text{ is the uniform distribution over } (t, 1] \subset M.$$

In words, if the true state, $\omega$, is in $(0, t]$, $\pi^s(t)$ reveals the true states by generating $m$ that is equal to the true states; if the true state, $\omega$, is in $(t, 1]$, $\pi^s(t)$ does not reveal the exact true state, for some value $t \in (0, 1)$. Note that $\pi^s(t)$ consumes all signals in $M := [0, 1]$. It is also worth noting that there are possibly many $\pi^s(t)$s depending on the value of $t$; there is a set of $\pi^s$, $\{\pi^s(t) \ for \ t \in (0, 1)\}$.

Now consider a subgame after S commits to a $\pi^s(t)$ with some $t \in (0,1)$. Consider a *hypothetical* equilibrium in this subgame at which $S_c$ employs the following strategy:

For all $\omega \in [0,1]$,

$$\sigma(m|\omega) \text{ is the uniform distribution over } (t,1] \subset M = [0,1].$$

In words, every type of $S_c$ employs the same mixed strategy that uniformly randomizes over $(t,1] \subset M$. Denote $S_c$'s strategy above by $\sigma_s := \{\sigma(m|\omega)\}_{\omega \in [0,1]}$, where $\sigma(m|\omega)$ is defined as above. For a while we will not discuss whether $\sigma_s$ is incentive compatible for $S_c$ and we will only focus on R's IC strategy.

At this hypothetical equilibrium at which $S_c$ employs $\sigma_s$, R's inference on $\omega$ and the corresponding optimal action are as follows. First, note $S_c$ only uses a part of $M$ which is $(t,1]$; $S_c$ does not "contaminate" the meaning of $m \in [0,t]$ which are generated by $\pi^s(t)$. Knowing this, if R observes a $m \in [0,t]$, R will believe that $m \in [0,t]$ is generated by $\pi^s(t)$ (or sent by $S_b$) with probability 1. Thus, when R observes a $m \in [0,t]$, she learns the true states and her optimal action is $a^*(m) = m$.

Secondly, suppose R observes $m \in (t,1]$. Not as the case of $m \in [0,t]$, $S_c$ uses $m \in (t,1]$. Knowing this, if R observes $m \in (t,1]$, she takes $S_c$'s strategy into account. That is, R will believe that: (1) with probability $\alpha$, $m$ is generated by $\pi^s(t)$ and (2) with probability $1 - \alpha$, $m$ is sent by $S_c$. Then, on the one hand, R believes that with probability $\alpha$, the true state is between $(t,1]$ since this is the information contained in $\pi^s(t)$. On the other hand, R believes that with probability $1 - \alpha$, the true state is between $[0,1]$ since all types of $S_c$ employs the same strategy and it does not give any information about $\omega$. Accordingly, when R observes a $m \in (t,1]$, her optimal action is a simple weighted average of $\frac{1+t}{2}$ and $\frac{1}{2}$ with some appropriate weight as follows:

$$a^*(m \in (t,1]) = \frac{\alpha(1-t)}{\alpha(1-t) + (1-\alpha) \times 1} \left(\frac{1+t}{2}\right) + \frac{1-\alpha}{\alpha(1-t) + (1-\alpha) \times 1} \left(\frac{1}{2}\right) \qquad (1)$$

As a summary, R's optimal (or IC) strategy is as follows:

$$
\begin{aligned}
a^*(m) &= m \quad \text{if } m \in [0,t], \\
&= \bar{a} \quad \text{if } m \in (t,1],
\end{aligned}
$$

where $\bar{a}$ is the action specified in equation (1) above. Denote R's strategy above by $a_s = \{a^*(m)\}_{m \in [0,1]}$.

Note that, at this hypothetical equilibrium, there is a certain degree of information transmission: R learns the exact true states if she observes $m \in [0,t]$ and R "partially" learns

where $\omega$ is in when she observes $m \in (t, 1]$. Hence, if $(\sigma_s, a_s)$ actually constitutes an equilibrium in the subgame following a $\pi^s(t)$ with some $t \in (0, 1)$, S's partial commitment to the $\pi^s(t)$ can strictly improve both S's and R's *ex ante* payoffs compared to the BECS in the no-commitment case (the pooling equilibrium).

Now, to check if $(\sigma_s, a_s)$ actually constitutes an equilibrium, we only need to check if $\sigma_s$ is incentive compatible for $S_c$ in the subgame following S's commitment to a $\pi^s(t)$ for some $t \in (0, 1)$. The following Lemma 2 shows that there are a set of $\pi^s(t)$ that induces subgames at which $S_c$'s strategy, $\sigma_s$, is incentive compatible.

**Lemma 2.** *For any $(\alpha, b) \in (0, 1) \times [1/2, \infty)$, there exists a set of $\pi^s(t)$ that induces subgames in which $(\sigma_s, a_s)$ constitutes an equilibrium. Given any parameter pair, $(\alpha, b) \in (0, 1) \times [1/2, \infty)$, the set of $\pi^s(t)$ is determined as $\{\pi^s(t) \text{ for } t \in T(\alpha, b)\}$ which is not empty and a proper subset of $\{\pi^s(t) \text{ for } t \in (0, 1)\}$, where $T(\alpha, b)$ is a proper subset of $(0, 1)$.*

The detailed proof and the definition of $T(\alpha, b)$ is in Appendix. In the proof, we start by specifying conditions for $\sigma_s$ to be incentive compatible for $S_c$. These conditions put restrictions on $b$ and the set of semi-fully revealing signaling devices: (1) $b$ should be greater than $1/4$ and (2) $\sigma_s$ is incentive compatible only if S commits to a $\pi^s(t) \in \{\pi^s(t) \text{ for } t \in T(\alpha, b)\}$ that is a subset of $\{\pi^s(t) \text{ for } t \in (0, 1)\}$.

The first restriction intuitively makes sense. At the equilibrium, all types of $S_c$ obtains the highest action, $\bar{a}(m \in (t, 1])$. If $b$ is too low, some low types of $S_c$ may prefer some lower actions, $a^*(m \in [0, t]) = m < \bar{a}(m \in (t, 1])$. To prevent these lower types' deviations to any of lower actions, $b$ should be greater than $1/4$ and it is satisfied since we are considering the high-bias area, $b \in [1/2, \infty)$.

The second restriction on the set $\{\pi^s(t) \text{ for } t \in (0, 1)\}$ is nothing but a restriction on the values of $t \in (0, 1)$ that characterizes a $\pi^s(t)$. In the proof, we show that for any given parameter pair, $(\alpha, b) \in (0, 1) \times [1/2, \infty)$, the set, $\{\pi^s(t) \text{ for } t \in T(\alpha, b)\}$, is nonempty, which is equivalent to show that $T(\alpha, b)$ is nonempty.

Lastly, it is worth mentioning the relation between our equilibrium construction above and one of equilibria constructed in the honesty model by Kim and Pogach (2014). Kim and Pogach (2014) consider a model of honesty where S is "honest" with a positive probability. In their model, S tells the exact true states with a positive probability. Hence, their model coincides with a subgame after S commits to a fully revealing signaling device in our partial-commitment environment (which we do not consider in this paper). One of equilibria in their model exhibits the same outcome structure as our equilibrium outcome in this subsection.[15]

---

[15]Our equilibrium outcome coincides the outcome of the "simplest" Type 1 equilibrium in Kim and Pogach (2014).

However, their construction is less permissible than ours in the sense that their construction requires a lower bound for $\alpha$ and an upper bound for $b$ while ours only requires a lower bound for $b$.

Lemma 2 states that for any $\alpha \in (0,1)$ and any $b \in [1/2, \infty)$, there always exists a $\pi^s(t)$ that induces a subgame that has an equilibrium with higher *ex ante* payoffs for both S and R. It is worth noting that Lemma 2 holds even when S's bias, $b$, is arbitrarily large and the commitment probability $\alpha$ is arbitrarily small. For example, if $b$ is arbitrarily large, $T(\alpha, b)$ becomes independent of $b$ and turns to be $(0, \frac{1-\sqrt{1-\alpha}}{\alpha}]$. Note that $(0, \frac{1-\sqrt{1-\alpha}}{\alpha}]$ is nonempty for any $\alpha \in (0,1)$; as $\alpha \to 0$ (*or* 1), $\frac{1-\sqrt{1-\alpha}}{\alpha} \to \frac{1}{2}$ (*or* 1). The following Proposition 3 concludes this section.

**Proposition 4.** *For any $\alpha \in (0,1)$ and $b \in [1/2, \infty)$, the Sender's partial commitment to $\pi^s(t) \in \{\pi^s(t) \ for \ t \in T(\alpha, b)\}$ **strictly** improves the ex ante payoffs for both the Sender and Receiver compared to BECS in the no-commitment case.*

*Proof.* By Lemma 2, S's commitment to $\pi^s(t) \in \{\pi^s(t) \ for \ t \in T(\alpha, b)\}$ induces a subgame that has an equilibrium constituted by $(\sigma_s, a_s)$. At such an equilibrium, there is a certain degree of information transmission. Therefore R obtains strictly higher ex ante payoffs compared to the BECS (a pooling equilibrium) in the no-commitment case. Since S's ex ante payoffs at an equilibrium is $EU^S = EU^R - b^2$, S also obtains strictly higher ex ante payoffs compare to the no-commitment case. $\square$

### 3.3.4 The Welfare Improvement under Partial Commitment

We have considered three classes of signaling devices and how S's partial commitment to those signaling devices can strictly improve both players' *ex ante* payoffs. We put the results in the previous sections altogether and obtain one of the main findings in this paper: Both S and R are *strictly* better off compared to the no-commitment case as S has an opportunity to commit himself to signaling devices with a positive probability, $\alpha \in (0,1)$, regardless of S's bias, $b \in (0, \infty)$.

**Proposition 5.** *For any commitment probability, $\alpha \in (0,1)$, and almost all positive biases of the Sender (i.e. for any $b \in (0, \infty) \backslash \{\frac{1}{2N^2} \ for \ N = 2, 3, ....\}$), the Sender's partial commitment to either $\pi^n$, $\pi^c(p)$, or $\pi^s(t \in T(\alpha, b))$ **strictly** improves the ex ante payoffs for both the Sender and Receiver compared to the no-commitment case.*

*Proof.* Proof by Proposition 2,3, and 4. $\square$

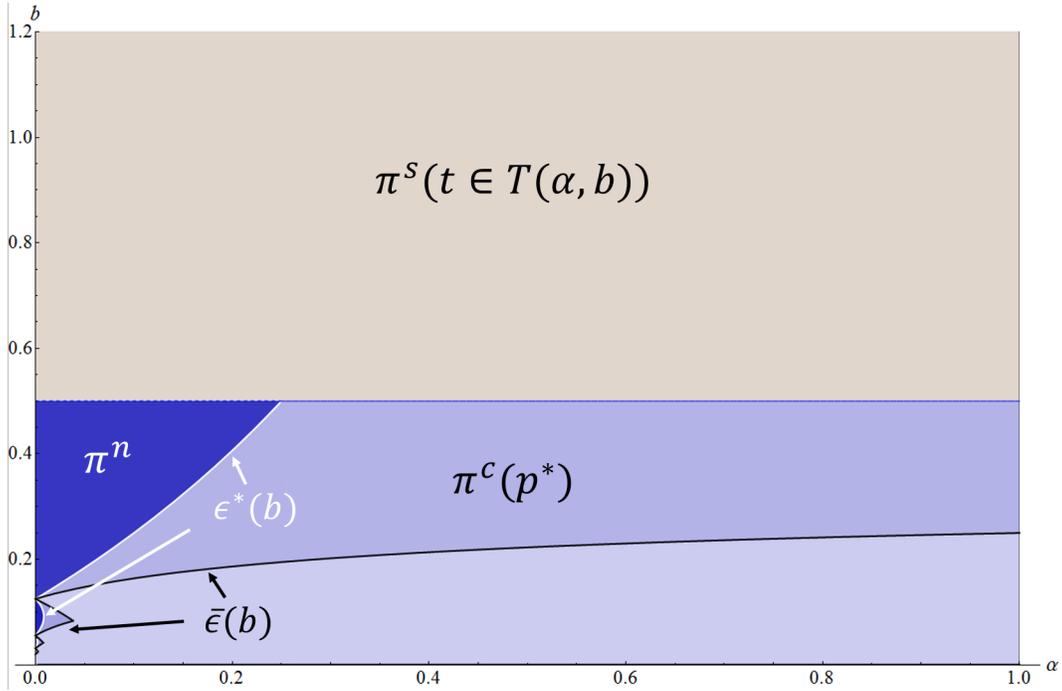The following Figure 3 summarizes Proposition 5.

24

Figure 3: Summary of Proposition 4

Given any value of commitment probability, $\alpha \in (0,1)$, if $b \in (0,1/2)\backslash\{\frac{1}{2N^2} \ for \ N = 2, 3, ...\}$, S's partial commitment to either $\pi^n$ or $\pi^c(p)$ induces a subgame that has an equilibrium at which both S and R obtain *strictly* higher *ex ante* payoffs compared to the no-commitment case. Especially, if $\alpha \in [\epsilon^*(b), 1)$, S's partial commitment to $\pi^c(p^* = \frac{\epsilon^*(b)}{\alpha})$ can lead both S and R to be in the equilibrium outcome that is the same as both the optimal noise outcome in Blume et al. (2007) and the best outcome in the mediation model by Goltsman et al. (2009).

Given any value of commitment probability, $\alpha \in (0,1)$, if $b \in [1/2, \infty)$, S's partial commitment to $\pi^s(t)$ for any $t \in T(\alpha, b)$ induces a subgame that has an equilibrium at which both S and R obtain *strictly* higher *ex ante* payoffs compared to the no-commitment case by guaranteeing a certain degree of communication at the equilibrium. Interestingly, this is true even when $\alpha$ is arbitrarily close to 0 and $b$ is arbitrarily large.

In summary, for high biases, $b \in [1/2, \infty)$, S's partial commitment to a semi-fully revealing signaling device, $\pi^s(t)$, can improve the welfare of both players; for almost all low biases, $b \in (0, 1/2)\backslash\{\frac{1}{2N^2} \ for \ N = 2, 3, ..\}$, S's partial commitment to a signaling device that involves noise can improve the welfare. Note that the only case that the welfare improvement fails via S's partial commitment is measure zero event; these appear as a set of horizonal lines, $\{(\alpha, b) \in (0,1) \times \{b = \frac{1}{2N^2} \ for \ N = 2, 3, ...\}\}$, in 2-dimensional parameter space.

Having S's commitment bind with a positive probability restricts S's ability to manipulate R's beliefs compared to the no-commitment case, which seems to only benefit R but not S. In fact, S cannot be worse off by having his commitment bind with a positive probability as we saw in Proposition 1.

Then, the question becomes "Can S be strictly better off compared to the no-commitment case?" and, if he can, "When and how S can be strictly better off?" The answer to the first question is "Yes" and the answer to the second question is "Almost always via commitment to three types of signaling devices: $\pi^n$, $\pi^c(p)$, and $\pi^s(t)$.

The intuition that we learned from Proposition 1 can be applied in the uniform-quadratic case. That is, S obtains more "freedom" by "tying" his hands: By having the chance to commit to signaling devices, S can have a strictly larger set of possible equilibria than the no-commitment case. On top of this, in that enlarged set of equilibria, there are equilibria that gives strictly higher *ex ante* payoffs for him in the uniform-quadratic case.

## 4  Conclusion

Departing from the full-commitment assumption commonly assumed in the Bayesian persuasion literature, this paper proposes a model of partial commitment in which the Sender's commitment is binding with a probability less than one. The model is in-between the cheap-talk (no commitment) and the Bayesian persuasion (full commitment) models, which we use as our benchmarks.

We first show that the Sender is *weakly* better off as the commitment probability increases: tightening a strap around his hands cannot make him worse off. When the preferences of the Sender and Receiver is *ex ante* well-aligned, the Receiver is also *weakly* better off as the commitment probability increases.

Then we study the model in a specific environment called the uniform-quadratic case which is widely used in the cheap-talk literature. We show that the welfare is *strictly* improved as we move from the no-commitment to the partial-commitment to the full-commitment cases, which holds for any level of the Sender's bias. Interestingly, the Sender can achieve the best outcome in Blume et al. (2007) and Goltsman et al. (2009) by committing himself to the convexified signaling device.

# 5   Appendix

Proof of Remark 1

*Proof.* We start by specifying S's payoffs from an arbitrary signaling device, $\pi$. Suppose S chooses and commits to a signaling device $\pi$. Since R knows that S's commitment is always binging, R will form the following posterior when she observes a $m \in M$,

$$\mu(\omega|m) = \frac{\pi(m|\omega)f(\omega)}{\int_0^1 \pi(m|\omega)f(\omega)d\omega}.$$

Given a $\mu(\omega|m')$, R will choose an optimal action $a^*(m') = \underset{a \in A}{argmax} \int_0^1 -(\omega - a)^2 d\mu(\omega|m')$. Note that $a^*(m') = -\int_0^1 \omega d\mu(\omega|m') = -E_{\mu_{m'}}[\omega]$. Then, S's expected payoffs from inducing the posterior, $\mu(\omega|m')$, (which is, in turn, equal to inducing the action, $a^*(m')$) is equal to

$$
\begin{aligned}
E_{\mu_{m'}}[U^S(\cdot)] &= \int_0^1 -(\omega + b - a^*(m'))^2 d\mu(\omega|m') \\
&= -E_{\mu_{m'}}[\omega^2] + \left(E_{\mu_{m'}}[\omega]\right)^2 - b^2 = -\sigma^2_{\mu_{m'}}[\omega] - b^2,
\end{aligned}
$$

where $-\sigma^2_{\mu_{m'}}[\omega]$ is the variance of R's posterior distribution, $\mu(\omega|m')$. Then, we can write S's *ex ante* payoffs as follows:

$$E_\tau[E_{\mu_{m'}}[U^S(\cdot)]] = E_\tau[-\sigma^2_{\mu'_m}[\omega] - b^2] = -E_\tau[\sigma^2_{\mu'_m}[\omega]] - b^2,$$

where $E_\tau[\cdot]$ means the expectation with respect to the probability measure, $\tau$, that measures the probability that a posterior, $\mu(\omega|m')$, is induced (or probability that $m$ is sent). Note that the term, $-E_\tau[\sigma^2_{\mu'_m}[\omega]]$, is the expected value of variances of R's posteriors, $\{\mu(\omega|m)\}_{m \in M}$, and $-b^2$ is a constant. Thus, at an equilibrium, S chooses to commit to a signaling device that minimizes $-E_\tau[\sigma^2_{\mu'_m}[\omega]]$. S can minimize $-E_\tau[\sigma^2_{\mu'_m}[\omega]]$ to zero by completely removing R's uncertainty about the states of the world. Thus, at any equilibrium, S should choose a $\pi$ that fully reveals the true states to R and S gets the maximized payoffs of $-b^2$.

A simple algebra shows that $-E_\tau[\sigma^2_{\mu'_m}[\omega]]$ is R's *ex ante* payoffs. Thus, at any equilibrium, R's *ex ante* payoffs is 0. In turn, this is a first-best outcome. Given the uniform-quadratic setting, this payoffs allocation hits the Pareto frontier: it is impossible to make S be better off while R is not worse off and vice versa. Hence, the equilibrium outcome is a first-best outcome. □

$\bar{\epsilon}(b)$ is defined as follows:

$$
\begin{aligned}
\bar{\epsilon}(b) &= \frac{2(1-2b(N-1)^2)(1+2b(N-1)N)}{N((2N-3)^2(1+2b(N-1))-8b^2(N-1)^3)} && \text{if } \frac{1}{2N(N-1)} \le b < \frac{1}{2(N-1)^2}, \\
&= \frac{2(1+2b(N-1)N)(2bN^2-1)}{(N-1)(1+2(1-b)N)(1+2N-4bN^2)} && \text{if } \frac{1}{2N^2} < b < \frac{1}{2N(N-1)}, \\
&= 0 && \text{if } b = \frac{1}{2N^2},
\end{aligned}
$$

where $N = 2, 3, 4, 5, \ldots$ and $N$ is the number of elements in an equilibrium partition at a front-loading equilibrium given a bias $b \in (\frac{1}{2N^2}, \frac{1}{2(N-1)^2})$.

For example, if $b = 1/10$, we can construct a front-loading equilibrium that partitions the state space into 3-interval since $b \in (1/18, 1/8)$. Then, $N = 3$ and note that $1/0 \in (\frac{1}{2\cdot3\cdot(3-1)}, \frac{1}{2\cdot(3-1)^2})$. Thus, $\bar{\epsilon}(b = 1/10) = \frac{2(1-8b)(1+12b)}{3(9(1+4b)-64b^2)} = 0.0245$. If $b = 1/15$, we can still construct a front-loading equilibrium with 3-interval partition. Thus, $N = 3$ but now $b = 1/15 \in (\frac{1}{2\cdot3^2}, \frac{1}{2\cdot3\cdot2})$. Hence, $\bar{\epsilon}(b = 1/15) = \frac{2(1+12b)(18b-1)}{2(1+6(1-b)(7-36b)} = 0.08088$. Finally, if $b = 1/18$, $\bar{\epsilon}(b = 1/18) = 0$.

When $(\alpha, b) \in (1/2, 1) \times (0, 1/4)$, S's payoffs at a BECS that partitions $[0,1]$ into $N$-interval is equal to $-\frac{1}{12} - \frac{b^2(N-1)^2}{3} - b^2$. This payoffs is always strictly less than $-\frac{1}{3}b(1-b) - b^2$ as the following show:

$$
\begin{aligned}
& -\frac{1}{3}(1-b)b - b^2 - \left( -\frac{1}{12} - \frac{b^2(N^2-1)}{3} - b^2 \right) \\
&= -\frac{1}{3}b + \frac{1}{3}b^2 + \frac{1}{12N^2} + \frac{b^2(N^2-1)}{3} = \frac{1}{12N^2} - \frac{1}{3}b + \frac{N^2b^2}{3} \\
&= \frac{1}{12N^2}\left(1 - 4N^2b + 4N^4b^2\right) = \frac{1}{12N^2}(1-2N^2b)^2 > 0.
\end{aligned}
$$

*Proof.* We need to show that there exists a signaling device that induces a subgame having $(\sigma_s, a_s)$ as an equilibrium. We consider a subgame following a $\pi^s(t)$ with some $t$. Conditions for $(\sigma_s, a_s)$ to be an equilibrium are identified as three conditions on $t$. Then, we show that, for any parameter pair, $(\alpha, b) \in (0,1) \times [1/2, \infty)$, there exist a set of $t$ that satisfies three conditions.

Consider a subgame following a $\pi^s(t)$ with some $t \in (0,1)$ and consider $(\sigma_s, a_s)$, where $\sigma_s$ and $a_s$ are defined as follows:

For all $\omega \in [0,1]$,

$$\sigma(m|\omega) \text{ is the uniform distribution over } (t,1] \subset M = [0,1],$$

and

$$
\begin{aligned}
a_s(m) &= m \quad \text{if } m \in [0,t], \\
&= \bar{a} \quad \text{if } m \in (t,1],
\end{aligned}
$$

where $\bar{a} = \frac{\alpha(1-t)}{\alpha(1-t)+(1-\alpha)\times 1}\left(\frac{1+t}{2}\right) + \frac{1-\alpha}{\alpha(1-t)+(1-\alpha)\times 1}\left(\frac{1}{2}\right)$.

Given $\sigma_s$, $a_s$ is optimal for R. Consider $S_c$'s IC conditions given $a_s$. For $\sigma_s$ to be IC for $S_c$, the following three conditions should be simultaneously satisfied:

(1) $t \in (0,b]$,

(2) $t \in (0,\bar{a}]$,

(3) $-(\omega + b - \bar{a})^2 \geq -(\omega + b - a_s(m))^2$, for all $a_s(m)$ and for all $\omega \in [0,1]$.

Condition (1) should be satisfied. Otherwise, $S_c$'s types near 0 can get their most preferred actions by deviating to messages that are near $t > b$. For example, suppose that $t > b$ and consider $S_c$'s type $\omega = 0$. Given R's strategy above, message $m = b$ will induce $a^*(m = b) = b$ since $t < m = b$. Thus, type $\omega = 0$ will deviate to $m = b$ instead of using $\sigma_s$.

Condition (2) simply means that $\bar{a}$ should be the highest action that $S_c$ can get. If $t > \bar{a}$, then there are actions that are greater than $\bar{a}$ and less than $t$. For example, $m' \in (\bar{a}, t)$ will induce action $a^*(m') = m'$ according to R's strategy above. One can easily see that some high types of $S_c$ would prefer $a^*(m') = m'$ to $\bar{a}$. For example, consider type $\omega = 1$ of $S_c$. This type will get the maximum payoffs of 0 when he induces the action, $1 + b$. Given the quadratic loss function, $-(\omega + b - a)^2$, his payoffs increases as $a \in A$ increases up to $1 + b$.

29

Hence, his most preferred action will be the highest action he can induce. If $t > \bar{a}$, he can induce higher actions than $\bar{a}$. Thus, he deviates to any $m \in (\bar{a}, t)$.

Condition (3) are usual IC constraints for $S_c$: all types of $S_c$ should (weakly) prefer the action they get, $\bar{a}$, to any other actions that they could get from a deviation.

Condition (2) and (3) can be simplified as follow. First, consider condition (2). A simple algebra yields $t \leq \bar{a} \Leftrightarrow t \leq \frac{1 - \sqrt{1-\alpha}}{\alpha}$.

Secondly, given condition (1) and (2), condition (3) can be simplified as follows:

$$|b - \bar{a}| \leq b - t$$

. Note that condition (3) by itself implies $|\omega + b - \bar{a}| \leq |\omega + b - a_s(m)|$, for all $\omega \in [0,1]$ and for all $a_s(m)$. Condition (2) implies $\bar{a} \geq a_s(m)$ for all $a_s(m)$.[16] In addition, since $\frac{\partial^2 U^S}{\partial \omega \partial a} > 0$, if $a'' > a'$ and $\omega'$ prefers $a''$ to $a'$, then any $\omega'' > \omega'$ also prefers $a''$ to $a'$. This implies that, if the inequality holds for $\omega = 0$, the inequality holds for all $\omega \in (0,1]$. Hence, it suffices to consider only the inequality with $\omega = 0$, $|b - \bar{a}| \leq |b - t|$. Finally, by condition (1), we can get rid of $|\cdot|$ on the right-hand side.

In summary, for $\sigma_s$ to be IC, there should exist some $t \in (0,1)$ that simultaneously satisfies the following:

(a) $t \in (0, b]$,

(b) $t \in (0, \frac{1 - \sqrt{1-\alpha}}{\alpha}]$ (from condition (2)),

(c) $|b - \bar{a}| \leq b - t$ (from condition (1), (2), and (3)).

Condition (c) can be satisfied in two ways depending on if $b \geq \bar{a}$ or $b < \bar{a}$.

If $b \geq \bar{a}$, condition (c) becomes $\bar{a} \geq t$,

If $b < \bar{a}$, condition (c) becomes $b \geq \frac{1}{2}(\bar{a} + t)$.

First note that either $b \geq t$ or $b \geq \frac{1}{2}(\bar{a} + t)$ reduces to conditions on $t$ since $\bar{a}$ is a function of $t$. Hence all three conditions are about $t$. If there exists a $t$ that satisfies three conditions, then there exists a $\pi^s(t)$ that we desire.

Secondly, note that, if $b \geq \bar{a}$, condition (c) becomes redundant since it is same as condition (b); But, if $b < \bar{a}$, it is different from condition (b). Thus, if $b < \bar{a}$, condition (c) will give an additional condition that $t$ should satisfy other than just condition (a) and (b).

One can easily see that

---

[16]Note that since any $m < t$ will induce action $a_s(m) = m$, any $a_s(m < t) < t$. Any $m \geq t$ will induce $\bar{a}$. By condition (2), $\bar{a} \geq t_1$. Hence, $\bar{a} \geq t \geq a_s(m)$ for all $a_s(m)$.

$b - \bar{a} \geq 0$ if and only if $\alpha t^2 - 2\alpha bt + 2b - 1 \geq 0$, and

$b - \bar{a} < 0$ if and only if $\alpha t^2 - 2\alpha bt + 2b - 1 < 0$

Denote $\alpha t^2 - 2\alpha bt + 2b - 1$ by $f(t)$. The discriminant of $f(t)$ is $4\alpha b^2 - 4\alpha(2b - 1)$. Note that the discriminant is less than or equal to zero if $b \in [\frac{1-\sqrt{1-\alpha}}{\alpha}, \frac{1+\sqrt{1-\alpha}}{\alpha}]$ and is strictly positive if either $b \in (1/2, \frac{1-\sqrt{1-\alpha}}{\alpha})$ or $b \in (\frac{1+\sqrt{1-\alpha}}{\alpha}, \infty)$. We have three cases depending on which set $b$ belongs to; if $b \in [\frac{1-\sqrt{1-\alpha}}{\alpha}, \frac{1+\sqrt{1-\alpha}}{\alpha}]$, $b \geq \bar{a}$ holds for all $t \in (0, 1)$; But, in the other two cases, either $b \geq \bar{a}$ or $b < \bar{a}$ can hold depending on the value of $t \in (0, 1)$.

We denote the set of $t$ that satisfies tree conditions by $T(\alpha, b)$ and show that for each case, there always exists a set, $T(\alpha, b)$, which is non-empty.

*Case 1: $b \in [\frac{1-\sqrt{1-\alpha}}{\alpha}, \frac{1+\sqrt{1-\alpha}}{\alpha}]$*

In this case, $b \geq \bar{a}$ for any value of $t \in (0, 1)$: $f(t) \geq 0$ for all $t \in (0, 1)$ since the discriminant of $f(t)$ is strictly negative. Hence condition (c) becomes $\bar{a} \geq t$ which is equivalent to condition (b). It makes condition (c) redundant and we only need to find $t$ that satisfies condition (a) and (b). Any $t \in (0, b] \cap (0, \frac{1-\sqrt{1-\alpha}}{\alpha}]$ satisfies condition (a) and (b). Thus, in this case, $T(\alpha, b) = (0, b] \cap (0, \frac{1-\sqrt{1-\alpha}}{\alpha}] = (0, \min\{\frac{1-\sqrt{1-\alpha}}{\alpha}, b\}]$. Note that we have $b \geq \frac{1-\sqrt{1-\alpha}}{\alpha}$ in this case. Thus the set, $T(\alpha, b)$, is namely $(0, \frac{1-\sqrt{1-\alpha}}{\alpha}]$. Finally, $T(\alpha, b)$ is nonempty since $0 < \frac{1-\sqrt{1-\alpha}}{\alpha} < 1$ for any $\alpha \in (0, 1)$.

*Case 2: $b \in [1/2, \frac{1-\sqrt{1-\alpha}}{\alpha})$*

In this case, the quadratic equation, $f(t) = 0$, has two real roots since the discriminant of $f(t)$ is strictly positive. Thus, either $f(t) \geq 0$ or $f(t) < 0$ can hold depending on the value of $t \in (0, 1)$, which, in turn, implies that either $b \geq \bar{a}$ or $b < \bar{a}$ can hold; there are two ways to have condition (3) satisfied. For simplicity, we only consider the case that $b < \bar{a} (\Leftrightarrow f(t) < 0)$. The two real roots of $f(t) = 0$ are $b \pm \frac{\sqrt{\alpha(\alpha b^2 - 2b + 1)}}{\alpha}$ and denoted by $r_1, r_2$ respectively.

Suppose that $b < \bar{a}$. Then, $t \in [r_1 r_2]$. To have condition (a) satisfied, $t \leq b$. And note that $b < r_2 = b + \frac{\sqrt{\alpha(\alpha b^2 - 2b + 1)}}{\alpha}$. Condition (a) narrows $[r_1, r_2]$ down to $[r_1, b]$.

Now to have condition (b) satisfied, $t \leq \frac{1-\sqrt{1-\alpha}}{\alpha}$. Note that we are considering $b \in [1/2, \frac{1-\sqrt{1-\alpha}}{\alpha})$. Hence, any value of $t \in [r_1, b]$ satisfies condition (b) since $t < b < \frac{1-\sqrt{1-\alpha}}{\alpha}$.

Finally, check if some $t \in [r_1, b]$ can satisfy condition (c). Since $t \in [r_1, b] \subset [r_1, r_2]$, $f(t) < 0$ and $b < \bar{a}$. Since $b < \bar{a}$, condition (c) becomes $b \geq \frac{1}{2}(\bar{a} + t)$. To have $b \geq \frac{1}{2}(\bar{a} + t)$, $t$ should be in $(0, k_1] \cup [k_2, 1)$, where $k_1 < k_2$ are two real roots of $b = \frac{1}{2}(\bar{a} + b)$ and are $\frac{1 + 2\alpha b \pm \sqrt{4\alpha^2 b^2 - 8\alpha b + 1 + \alpha}}{3\alpha}$. [17]

---

[17] The existence of $k_1$ and $k_2$ are guaranteed if $b < \frac{1-\sqrt{1-\alpha}}{\alpha}$ and it holds since we are in the case that $b \in [1/2, \frac{1-\sqrt{1-\alpha}}{\alpha})$.

In summary, in this case, $T(\alpha, b)$ is equal to $[r_1, b] \cap \{(0, k_1] \cup [k_2, 1)\}$: if $t \in [r_1, b]$, condition (a) and (b) are satisfied and, if $t \in (0, k_1] \cup [k_2, 1)$, condition (c) is satisfied. Now only thing left is to show that $[r_1, b] \cap \{(0, k_1] \cup [k_2, 1)\}$ is nonempty.

First, with a tedious algebra, one can show that $k_1 \leq b < k_2$ if $b < \frac{1-\sqrt{1-\alpha}}{\alpha}$. Hence, $T(\alpha, b) = [r_1, b] \cap \{(0, k_1] \cup [k_2, 1)\}$ becomes $[r_1, k_1]$.

Then, to $T(\alpha, b) = [r_1, k_1]$ be nonempty, $r_1 < k_1$ should hold. One can easily check that $r_1 \geq 0$ if $b \geq 1/2$. We have two cases: (1) $b = 1/2$ and $b \in (1/2, \frac{1-\sqrt{1-\alpha}}{\alpha})$. First, suppose that $b = 1/2$, then $r_1 = 0$. Thus, we need to check if $k_1 > 0$ or not. $k_1 > 0$ if and only if $b > /14$. $k_1 > 0$ if and only the numerator of $k_1$ is strictly positive:

$$1 + 2\alpha b - \sqrt{4\alpha^2 b^2 - 8\alpha b + 1 + 3\alpha} > 0,$$

and it reduces to be $b > 1/4$. Thus, if $b = 1/2$, $T(\alpha, b) = (0, k_1]$ is nonempty. Secondly, suppose $b > 1/2$. Then $r_1 > 0$ and we need to check if $k_1$ is strictly greater than $r_1$. $k_1 > r_1$ if and only if $\alpha < 1$. To see this,

$$k_1 - r_1 > 0$$
$$\Leftrightarrow \quad (1 - \alpha b) + \sqrt{\alpha(\alpha b^2 - 2b + 1)} > \sqrt{4\alpha^2 b^2 - 8\alpha b + 3\alpha + 1},$$
$$\text{Define } A := \alpha(\alpha b^2 - 2b + 1) \text{ and } B := 4\alpha^2 b^2 - 8\alpha b + 3\alpha + 1,$$
$$\Leftrightarrow \quad (1 - \alpha b)^2 + 2(1 - \alpha b)\sqrt{A} + A > B,$$
$$\Leftrightarrow \quad 2\sqrt{A}(1 - \alpha b) > B - A - (1 - \alpha b)^2,$$
$$\Leftrightarrow \quad 2\sqrt{A}(1 - \alpha b) > 2\alpha(\alpha b^2 - 2b + 1) \Leftrightarrow 1 > \alpha.$$

Hence, $T(\alpha, b) = [r_1, k_1]$ is nonempty too.

Thus, when $b \in [1/2, \frac{1-\sqrt{1-\alpha}}{\alpha})$, $T(\alpha, b) = [r_1, k_1]$ and it is nonempty; Any $\pi^s(t)$ with $t \in T(\alpha, b)$ can induce the subgame having $(\sigma_s, a_s)$ as an equilibrium.

*Case 3:* $b \in (\frac{1+\sqrt{1-\alpha}}{\alpha}, \infty)$

Since $b > \frac{1+\sqrt{1-\alpha}}{\alpha}$, $f(t) = 0$ has two real roots, $r_1$ and $r_2$. Consider $t \leq r_1$. Then, $f(t) \geq 0$ which implies $b \geq \bar{a}$. Thus, for any $t \in (0, r_1]$, condition (c) becomes redundant.

Now we only need to check which $t \in (0, r_1]$ can satisfy both condition (a) and (b). First any $t \in (0, r_1]$ satisfies condition (a), $t \in (0, b]$, since $b < r_1 = b - \frac{\sqrt{\alpha(\alpha b^2 - 2b + 1)}}{\alpha}$. To have condition (b) satisfied, $t \in (0, \frac{1-\sqrt{1-\alpha}}{\alpha}]$.

Hence, $T(\alpha, b)$ in this third case is $(0, r_1] \cap (0, \frac{1-\sqrt{1-\alpha}}{\alpha}] = (0, \min\{r_1, \frac{1-\sqrt{1-\alpha}}{\alpha}\}]$. It turns out that $\min\{r_1, \frac{1-\sqrt{1-\alpha}}{\alpha}\} = \frac{1-\sqrt{1-\alpha}}{\alpha}$ if $\alpha \in (0, 1)$ and $b \in (\frac{1+\sqrt{1-\alpha}}{\alpha}, \infty)$. To see this,

$$r_1 - \frac{1 - \sqrt{1 - \alpha}}{\alpha} > 0,$$

$$\Leftrightarrow \quad \alpha b - \sqrt{\alpha(\alpha b^2 - 2b + 1)} - 1 + \sqrt{1-\alpha} > 0,$$
$$\Leftrightarrow \quad \alpha b - 1 + \sqrt{1-\alpha} > \sqrt{\alpha(\alpha b^2 - 2b + 1)}$$

Note that both sides are positive. $\alpha b - 1 > 0$ since $b > \frac{1}{\alpha} > \frac{1+\sqrt{1-\alpha}}{\alpha}$, and $\alpha(\alpha b^2 - 2b + 1) > 0$ since $b > \frac{1+\sqrt{1-\alpha}}{\alpha}$ implies $f(t) = 0$ has two real roots. Square both sides. After rearranging terms, we have

$$2(\alpha b - 1)\sqrt{1-\alpha} > 2(\alpha - 1),$$

which always holds if $\alpha \in (0,1)$ and $b \in (\frac{1+\sqrt{1-\alpha}}{\alpha}$. $\alpha \in (0,1)$ implies $2(\alpha-1) < 0$ and we know $2(\alpha b - 1)\sqrt{1-\alpha} > 0$ if $b \in (\frac{1+\sqrt{1-\alpha}}{\alpha}, \infty)$. Hence, $T(\alpha, b)$ in this case is $(0, \frac{1-\sqrt{1-\alpha}}{\alpha}]$ and it is also nonempty.

We have considered three cases. For each case, we have identified $T(\alpha, b)$ and showed it is nonempty set. As a summary,

If $b \in [1/2, \frac{1-\sqrt{1-\alpha}}{\alpha})$, $T(\alpha, b) = [r_1, k_1]$,

If $b \in [\frac{1-\sqrt{1-\alpha}}{\alpha}, \infty)$, $T(\alpha, b) = (0, \frac{1-\sqrt{1-\alpha}}{\alpha}]$,

where $r_1 = b - \frac{\sqrt{\alpha(\alpha b^2 - 2b + 1)}}{\alpha}$ and $k_1 = \frac{1 + 2\alpha b - \sqrt{4\alpha b^2 - 8\alpha b + 1 + 3\alpha}}{\alpha}$.

In words, for any $(\alpha, b) \in (0,1) \times [1/2, \infty)$, there always exist $T(\alpha, b)$ that satisfies all three conditions for $(\sigma_s, a_s)$ to constitute an equilibrium in a subgame following $\pi^s(t \in T(\alpha, b))$. In other words, for any parameter pair, $(\alpha, b) \in (0,1) \times [1/2, \infty)$ there always exists a set of $\pi^s(t)$ that induces a subgame having $(\sigma_s, a_s)$ as an equilibrium. $\qquad \square$

# References

ALONSO, R., AND O. CÂMARA (2013): "Persuading skeptics and reaffirming believers," *Available at SSRN 2306820*.

BLUME, A., O. J. BOARD, AND K. KAWAMURA (2007): "Noisy talk," *Theoretical Economics*, 2(4), 395–440.

CRAWFORD, V. P., AND J. SOBEL (1982): "Strategic information transmission," *Econometrica: Journal of the Econometric Society*, pp. 1431–1451.

GOLTSMAN, M., J. HÖRNER, G. PAVLOV, AND F. SQUINTANI (2009): "Mediation, arbitration and negotiation," *Journal of Economic Theory*, 144(4), 1397–1420.

KAMENICA, E., AND M. GENTZKOW (2011a): "Bayesian persuasion," *American Economic Review*, 101(6), 2590–2615.

——— (2011b): "Competition in persuasion," Discussion paper, National Bureau of Economic Research.

KIM, K., AND J. POGACH (2014): "Honesty vs. advocacy," *Journal of Economic Behavior & Organization*, 105, 51–74.

KOLOTILIN, A. (2014): "Optimal Information Disclosure: Quantity vs. Quality," *UNSW Australian School of Business Research Paper*, (2013-19).

WANG, Y. (2012): "Bayesian persuasion with multiple receivers," *Job Market Paper, Department of Economics, University of Pittsburgh*, pp. 1–45.