

I Don't Know*

Matthew Backus[†] Andrew T. Little[‡]

March 15, 2020

Abstract

Political decision-makers make choices in a complex and uncertain world, where even the most qualified experts may not know what policies will succeed. Worse, if these experts care about their reputation for competence, they may be averse to admitting what they don't know. We model the strategic communication of uncertainty, allowing for the salient reality that sometimes the effects of proposed policies are impossible to know. Our model highlights the challenge of getting experts to admit uncertainty, even when it is possible to check predictive success. Moreover, we identify a novel solution: checking features of the question that only good experts will infer – in particular, whether the effect of policies is knowable – allows for equilibria where uninformed experts do say “I Don't Know.”

*Thanks to Charles Angelucci, Jonathan Bendor, Sylvain Chassang, Wouter Dessein, Sean Gailmard, James Hollyer, Ryan Hübert, Navin Kartik, Greg Martin, Mallesh Pai, Matthew Mitchell, Andrea Prat, Michael Raith, Daniel Rappaport, Maher Said, Jim Snyder, Joel Sobel, Philipp Strack, and audiences at MPSA 2015, EARIE 2017, The 28th International Game Theory Conference, QPEC 2017, Petralia Workshop 2017, SAET 2017, ESSET 2018, Columbia, Harvard, the Higher School of Economics, the University of Hong Kong, Peking University, and Stanford for thoughtful comments and suggestions. We thank Jett Pettus, Alphonse Simon, and Brenden Eum for excellent research assistance. All remaining errors are our own.

[†]Columbia University, NBER, and CEPR, matthew.backus@columbia.edu

[‡]Corresponding author. UC Berkeley, andrew.little@berkeley.edu

[...] it is in the admission of ignorance and the admission of uncertainty that there is a hope for the continuous motion of human beings in some direction that doesn't get confined, permanently blocked, as it has so many times before in various periods in the history of man.

— Richard Feynman, John Danz Lecture, 1963

Policy-making is hard.

— Callander (2011)

1 Introduction

Political decision-makers frequently make disastrous choices. They waste blood and treasure on unwinnable wars, destroy economies with poor monetary policy, and underestimate the threat of coups or revolutions before their opponents show up at the gates. Sometimes poor decisions are made despite the availability of information about how they will turn out. Decision-makers may not consult the right experts, or ignore their advice. Other times the best course of action isn't even knowable, and the real danger is being persuaded to take risky action by “experts” who pretend to know what policies will work.

Most work on strategic communication in political science focuses on problems driven by differences of preference (“bias”) among experts and decision-makers (e.g., Gilligan and Krehbiel, 1987; Patty, 2009; Gailmard and Patty, 2012). However, even if experts have the same *policy* preferences as decision-makers, some are more competent than others at assessing the available evidence required to give good advice. And, as a literature in eco-

nomics and finance on career concerns following Holmström (1999) makes clear, experts' desire to appear competent ("reputational concerns") can distort the actions of agents, including the advice they give to principals (e.g., Ottaviani and Sørensen, 2006).

We bring this style of communication model to a political context, and also place a novel focus on the *difficulty* of policy questions. As is familiar to anyone who has tried to study the causal effect of policies (empirically or theoretically), some questions are harder to answer than others. Uncertainty may be driven by expert incompetence or by the difficulty of the policy question. However, as long as knowledge about the effects of policies is correlated with competence, uninformed experts (competent or not) risk a reputational hit for admitting uncertainty. As a result, experts who care about perceptions of their competence will be reluctant to say "I don't know".

Can this problem be solved by validating experts' claims, either by asking other experts, checking other sources, or waiting to see if their predictions come true? Our core contention is that the answer depends on what exactly gets validated. Perhaps the most intuitive kind of validation is what we call *state validation*, or checking whether the expert claims are correct. We find that—at least by itself—this is not effective at getting uninformed experts to report their ignorance. On the other hand, *difficulty validation*, which means checking whether the question was answerable in the first place, tends to be much more effective at inducing experts to admit uncertainty.

We develop a formal model which highlights this problem and our proposed solution in a clear fashion. A decision-maker (DM, henceforth "she") consults an expert (henceforth "he") before making a policy decision. The DM is uncertain about a state of the world that dictates which policy choice is ideal. The DM is also uncertain about the quality of the expert and whether the optimal policy is knowable. Experts are either competent ("good") or incompetent ("bad"), and the question of which policy is better is either answerable

(“easy”) or not (“hard”). Bad experts learn nothing about the state of the world. Good experts learn the state if and only if the question is answerable. This means there are two kinds of uninformed experts: bad ones who never learn truth, and good ones faced with an unanswerable question.

The expert then communicates a message to the decision-maker, who chooses a policy. Finally, the decision-maker – potentially endowed, *ex post*, with information about the state or question difficulty – forms posterior beliefs about the expert quality. The best decisions are made in an *honest* equilibrium, where experts who know which policy is better reveal this information, and the uninformed types all say “I don’t know.”

Our main analysis studies the scenario where an expert primarily cares about his reputation. That is, the expert has relatively little concern for the quality of the policy made, though he is unbiased in the sense that the expert and DM agree on which policy is best (conditional on the state). If the DM gets no *ex post* information about the state or difficulty – the *no validation* case – our answer is bleak. Since they face no chance of being caught guessing, uninformed types could claim to know whether the policy would succeed and appear competent. As a result, honesty is impossible.

What if the DM learns the truth about the ideal policy (state validation) before evaluating the expert? One might expect that state validation can induce honesty, since it is possible to “catch” uninformed experts guessing incorrectly. But what should the DM infer when seeing an incorrect guess: that the expert is incompetent, or just that he is uninformed? Under a restriction to strategies and beliefs related to the Markov refinement common to repeated games, we show that the competent uninformed experts and the incompetent uninformed experts must play the same strategy. Further, upon observing an incorrect guess, the DM should infer that the expert is uninformed but still possibly competent; i.e., the exact same inference if the expert said “I don’t know”. Since the expert might get away with guessing,

and being caught is no worse than admitting uncertainty, an honest equilibrium is still not possible.

The limits of state validation echo past pessimistic results about how reputational concerns limit communication (e.g., Ottaviani and Sørensen, 2006). However, our focus on the importance of problem difficulty also suggests a novel path forward. The key barrier to honest communication with state validation is that it does not allow the competent experts asked an unanswerable policy question to differentiate themselves from the incompetent experts. Consider this from the perspective of a competent but uninformed expert: he knows that he is uninformed because it is impossible to know which policy is better, but precisely for this reason he can't do a better job of guessing the state than an incompetent expert. Where the competent expert does have a definitive advantage over the incompetent expert is not in knowing which policy is better, but in knowing whether the ideal policy is knowable in the first place.

We build on this insight to reach our key positive result: if the DM learns *ex post* whether the question was answerable (difficulty validation), partial if not complete admission of uncertainty is possible.¹ The good uninformed experts admit uncertainty, confident that the DM will learn the ideal policy wasn't knowable. Bad experts may admit uncertainty as well, if this is safer than guessing and potentially getting caught making a prediction when the validation reveals that the question was unanswerable.

These results have practical implications for how political decision-makers should structure their interactions with experts. Consulting experts with an interest in good policy being made, or investing in methods to check if their predictions are correct (e.g., running pilot studies) is useful for some purposes, but not for eliciting admission of uncertainty. Rather, it is important for decision-makers to be able to eventually learn whether the questions they

¹As elaborated in Section 5, this also requires either non-zero policy concerns or state validation.

ask are answerable. We discuss several ways this can be accomplished, such as consulting multiple experts or ensuring decision-makers (or someone who evaluates experts) have some general expertise in research methods in order to be able to evaluate what claims are credible.

In addition to highlighting the importance of difficulty validation, we derive several comparative static results. First, incentives to guess are weaker when experts are generally competent. This implies that admission of uncertainty can be more frequent in environments with more qualified experts. Second, when questions are *ex ante* likely to be answerable, experts face a stronger incentive to guess. So, the quality of policies can be *lower* in environments where the ideal policy is more frequently knowable, because this is precisely when bad experts guess the most, diluting the informative value of any message.

The reader will find a discussion of related literature in Section 2, the setup of the model in Section 3, and our solution concept in Section 4. In Section 5, we show how this solution concept quickly rules out admission of uncertainty for several cases of the model. Then, in Section 6, we focus on the “minimal” case where admission of uncertainty is possible: when the DM exercises difficulty validation and the expert has small policy concerns. In this setting we explicitly characterize the conditions for honest equilibria and we explore several comparative statics of interest. Section 7 concludes with a summary and discussion of some of the broader implications of the model.

2 Related Work

Uncertainty about ideal policy has at least two causes. First, different people want different things. Even if the effects of policy choices are well-known, it may be hard to determine

what is best collectively. Second, consensus about the effects of different policies is rare. The world is complicated, and frequently the most credible research gives limited if any guidance about the effects of political decisions.

Decision-makers can try to learn about what policies will work in several ways. They can hold debates about policy (Austen-Smith, 1990), try to learn from the experience of other polities, or experiment with new policies on a tentative basis (Callander, 2011). Either as a part of these processes or separately, they can consult experts employed by the government (staffers, bureaucrats, other politicians) or elsewhere (think tank employees, academics, pundits) (Calvert, 1985),² or delegate to them where optimal (Dessein, 2002).³

Even if there are experts who have a solid grasp on the wisdom of proposed policies, there are always less competent “experts” who lack valuable information but may still pretend to be informed. And with heterogenous preferences, even good experts may disagree with decision-makers about how to best use this information. As a result, the challenges to knowing the ideal course of action in the first place (ignorance and preference bias) generate parallel problems of getting this information into the hands of actual policy-makers. Further, if the policy-makers themselves have more information than voters and heterogenous competence, the signaling implications of their choices may also lead to sub-optimal policies.

The main literatures we draw on study these problems, and Table 1 presents a classification of the most related work to clarify where our model fits in. The rows correspond to the identity of the “sender”, and the columns to whether the main conflict of interest is differing policy preferences or the sender’s desire to appear competent.

²See also Manski (2019) for a broader discussions about communication of scientific uncertainty in policy analysis.

³See Fehrler and Janas (2019) for a model and experiment on delegation with a focus on competence.

Table 1: Classification of Related Literature

	Preference Bias	Reputation for Competence
Expert/Advisor	Crawford and Sobel (1982), Gilligan and Krehbiel (1987), Gailmard and Patty (2013)	Ottaviani and Sørensen (2006), Rappaport (2015), <i>This Paper</i>
Decision-maker	Fearon (1999), Fox and Shotts (2009)	Holmström (1999), Ashworth (2012), Canes-Wrone et al. (2001)

Notes: Here we classify related literature by whether the informed party/sender is an advisor (top row) or the actual decision-maker (bottom row), and whether the main conflict of interest is different preferences over ideal policy (left column) or reputation for competence (right column).

Most of the political science literature on decision-making under uncertainty highlights the problems which arise when actors (experts, policy-makers, bureaucrats, voters) have different preferences or ideologies. The top left cell of Table 1 contains examples where, following Crawford and Sobel (1982), an informed advisor or expert (not the decision-maker) may not communicate honestly because of a difference in ideal policy. Much of this work has studied how different institutional features like committee rules in congress (Gilligan and Krehbiel, 1987), bureaucratic hierarchies (Gailmard and Patty, 2012), alternative sources of information (Gailmard and Patty, 2013), and voting rules (Schnakenberg, 2015) either solve or exacerbate communication problems. Formal models of communication in other political settings such as campaigns (Banks, 1990), lobbying (Schnakenberg, 2017), and international negotiations (Kydd, 2003) also focus on how different preferences affect equilibrium communication.

Our main innovation with respect to the formal theories of expert communication in political science is to bring focus to the other main barrier to communication: reputation concerns (right column of Table 1). Most related work on reputation for competence focuses not on experts but politicians themselves (bottom right cell of Table 1). In some of these models (following Holmström 1999), politicians exert effort (e.g., Ashworth, 2012)

or otherwise manipulate the information environment (e.g., Little, 2017) to make signals of their ability or performance more favorable. Closer to our setting, others model the competence of decision-makers as affecting their ability to discern the best policy. In these models, concern about reputation can lead to suboptimal policy choices if, for example, voters think certain policies are *ex ante* more likely to be favored by competent politicians (Canes-Wrone et al., 2001). The bottom left cell of Table 1 contains examples of models where politicians also want to develop a reputation for being ideologically aligned with citizens (Fearon, 1999), which sometimes creates tradeoffs with concerns for competence (Fox and Shotts, 2009).⁴

We argue that studying the reputation concerns of experts (top right cell of Table 1) is fundamental in political settings. By definition, experts typically have more policy-relevant information than politicians. Further, particularly in a cheap talk environment where experts have no direct control over the policy being made (and frequently have a relatively small impact on what choice is made), they sometimes, if not usually, care more about their reputation for competence than the impact of their advice on policy. To our knowledge, there are no other models of “reputational cheap talk” in political science. Related work in other disciplines has shown that reputation concerns lead experts to bias and overstate their reports in order to convince a decision-maker that they are the “good” type (e.g., Ottaviani and Sørensen, 2006), though the exact kind of lie this induces may depend on the career stage of the agent (Prendergast and Stole, 1996) or the precise information structure (Rappaport, 2015).

In addition to bringing this style of model to a political context, we contribute to the reputational cheap talk literature by focusing on heterogeneity in the difficulty of questions, and

⁴Not all work in the bottom row involves choices made by politicians. For example, Leaver (2009) studies bureaucrats with reputation concerns (and a desire to avoid public criticism) may make suboptimal choices (bottom right cell), and judges fear having decisions overturned by higher courts who may have different preferences (e.g., Hübert, 2019) (bottom right).

introducing the notion of *ex post* validation of question difficulty rather than whether experts' claims were "correct". In one sense this is related to the precision of experts' private information in models such as Ottaviani and Sørensen (2006), or their accuracy in related and recent work on screening forecasters by Deb et al. (2018) – it creates variation in the quality of signals. However, there is an important feature that differentiates difficulty as we have formulated it: it is a property of the problem itself, not the expert or the expert's signal. This drives our results. Validating the difficulty of problems generates the key informational wedge between good uninformed experts (who know the problem difficulty) and bad experts (who do not).⁵

3 The Model

Our model is motivated by the following scenario: a decision-maker (abbreviated DM, pronoun "she") is making a policy choice. The DM could be the chief executive of a country, a local executive (governor, mayor), or the head of a bureaucracy. There is an unknown state of the world which affects the optimal policy. However, the DM does not observe this state; to this end she employs an expert (pronoun "he"). Experts may be competent or incompetent. Competent experts sometimes know the state of the world, but other times the state is unknowable. Incompetent experts know nothing.

For concreteness, we will use a running example where the policy in question is how much to restrict the emissions of a chemical, and the DM is the head of an environmental agency with statutory authority to set this regulation level. The expert could be a scientist within the agency or hired as an external consultant. A natural source of uncertainty is whether

⁵The closest to what we are calling difficulty in the prior literature that we are aware of is the information endowment of managers in Dye (1985) and Jung and Kwon (1988), in an altogether different setting where shareholders are uncertain as to the informational endowment of managers.

the chemical in question is harmful to humans. If it is harmful, the DM will want to choose more stringent regulations; if it is not harmful there is no reason to regulate emissions. The expert will learn whether the chemical is harmful if two things are true: (1) he is competent enough to locate and digest the relevant scientific literature, and (2) this literature contains an answer to whether the chemical is harmful. If the expert is not competent or there is no scientific consensus on the effect of the chemical, the expert will not learn anything useful.

3.1 The Information Environment

We formalize this information structure as follows.

State of the World. Let the state of the world be $\omega \in \Omega \equiv \{0, 1\}$. The state of the world encodes the decision-relevant information for the DM (e.g., $\omega = 1$ if the chemical in question is harmful to humans and $\omega = 0$ if not). It is unknown to the DM at the outset, which is why she consults an expert.

Let p_1 represent the common knowledge probability that the state is 1, and so it is equal to 0 with probability $1 - p_1$. To reduce the cases to consider, we also assume that $\omega = 1$ is the *ex ante* more likely state, so $p_1 \geq 1/2$.⁶

Expert Types. The expert has a type $\theta \in \Theta \equiv \{g, b\}$, which indicates whether he is *good* (able to digest the relevant scientific literature) or *bad* (not able to digest the scientific literature). For linguistic variety, we often call good experts “competent” and bad experts “incompetent”. Let $p_g \in (0, 1)$ be the probability than an expert is good, and so $1 - p_g$

⁶By the symmetry of the payoffs introduced below, substantively identical results hold if state 0 is more likely.

represents the probability of a bad expert. This probability is common knowledge, and the expert knows his type.

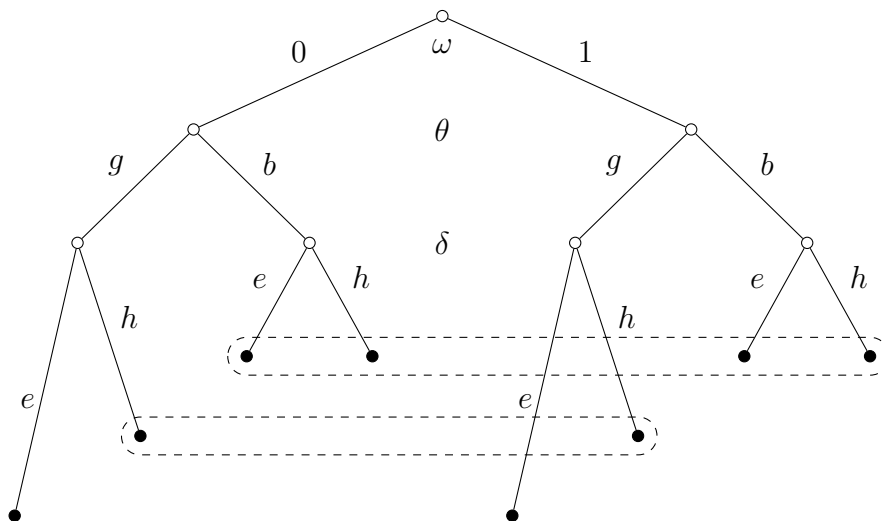
Question Difficulty. The difficulty of the question is captured by another random variable $\delta \in \Delta \equiv \{e, h\}$. That is, the question may be *easy* (the scientific literature is informative about whether the chemical is harmful), or *hard*, (the best research does not indicate whether the chemical is harmful). Let p_e be the common knowledge probability that $\delta = e$, and so $\delta = h$ with probability $1 - p_e$. The difficulty of the question is not directly revealed to either actor at the outset.

Experts Signal. The expert's type and the question difficulty determine what he learns about the state of the world. In our main analysis, we assume that the expert receives a completely informative signal if and only if he is good *and* the question is easy. If not, he learns nothing about the state. (In Appendix E, we analyze a more general information structure which only assumes that a signal is more likely to be informative when the expert is good and the question is easy.) Formally, let the signal be:

$$s = \begin{cases} s_1 & \omega = 1, \theta = g \text{ and } \delta = e \\ s_0 & \omega = 0, \theta = g \text{ and } \delta = e \\ s_\emptyset & \text{otherwise.} \end{cases} \quad (1)$$

In what follows, we will often refer to an expert who observes s_0 or s_1 as *informed*, and an expert who observes s_\emptyset as *uninformed*. Importantly, this distinction is not the same as good and bad. If an expert is informed he must be good, because bad experts always observe s_\emptyset . However, an uninformed expert may be good (if $\delta = h$) or bad.

Figure 1: Nature's Play and Experts' Information Sets



An important implication of our signal structure is that the good expert infers δ from his signal, since he observes s_0 or s_1 when $\delta = e$ and s_0 if $\delta = h$. More concretely, a competent expert can review the relevant literature and always determine whether the harmfulness of a chemical is known. On the other hand, bad experts, who always observe s_0 , learn nothing about the question difficulty.

3.2 Sequence of Play and Payoffs

The game proceeds in the following sequence: first, nature picks the random variables (ω, θ, δ) , according to independent binary draws with the probabilities p_1, p_g , and p_e specified above. Second, the expert observes his competence and signal and chooses a message from a infinite message space \mathcal{M} . The information sets of the expert are summarized in Figure 1. There are four: first, the expert may be bad; second, the expert may be good and the question hard; third, the expert may be good, the question easy, and the state 0, and finally, the expert may be good, the question easy, and the state 1.

Next the DM observes m and takes an action $a \in [0, 1]$, the *policy* choice. Her information set in this stage consists only of the expert report, i.e. $\mathcal{I}_{DM1} = (m)$.

In the running example, we can interpret a as the stringency of regulations applied to the chemical in question (restrictions on emissions, taxes, resources to spend on enforcement). To map cleanly to the formalization, $a = 0$ corresponds to the optimal policy if the chemical is not harmful. Policy $a = 1$ corresponds to the optimal level of regulation if the chemical is harmful.

Formally, let $v(a, \omega)$ be the *value* of choosing policy a in state ω . We assume the policy value is given by $v(a, \omega) \equiv 1 - (a - \omega)^2$. If taking a decisive action of $a = 0$ or $a = 1$, the value of the policy is equal to 1 for making the correct choice ($a = \omega$) and 0 for making the wrong choice ($a = 1 - \omega$). Taking an interior action ($0 < a < 1$) gives an intermediate policy value, where the v function implies an increasing marginal cost the further the action is from the true state. This function is common knowledge, but since the DM may be uncertain about ω , he may be uncertain about how the policy will turn out and hence the ideal policy. Let $\pi_1 = \mathbb{P}(\omega = 1 | \mathcal{I}_{DM1})$ denote the DM's belief that $\omega = 1$ when he sets the policy. Then, the expected value of taking action a is

$$1 - [\pi_1(1 - a)^2 + (1 - \pi_1)a^2], \quad (2)$$

which is maximized at $a = \pi_1$.

The quadratic loss formulation conveniently captures the realistic notion that when the expert does not learn the state, the decision-maker makes a better policy choice (on average) when learning this rather than being misled into thinking the state is zero or one. That is, in our formalization it is best to pick an intermediate level of regulation when it is unclear whether the chemical is harmful (e.g., modest emissions restrictions, labeling

requirements). Formally, if the question is unsolvable, the optimal action is $a = p_1$, giving an average payoff of $1 - p_1(1 - p_1)$, which is strictly higher than the average value of the policy for any other action.

After the policy choice is made, the DM may receive some additional information (validation), and then forms an inference about the expert competence $\pi_g \equiv \mathbb{P}(\theta = g | \mathcal{I}_{DM2})$. Different validation regimes (described below) affect what information the DM has at this stage, \mathcal{I}_{DM2} .

The decision-maker only cares about the quality of the policy:

$$u_{DM} = v(a, \omega). \tag{3}$$

So, as derived above, the DM will always pick a policy equal to his posterior belief that the state is 1, π_1 .

The expert cares about the belief that he is competent (π_g), and potentially also about the quality of the policy choice. We parameterize his degree of *policy concerns* by $\gamma \geq 0$ and write his payoff

$$u_E = \pi_g + \gamma v(a, \omega). \tag{4}$$

We first consider the case where $\gamma = 0$, i.e., the expert only cares about his reputation. We then analyze the case where $\gamma > 0$, focusing attention on the case where policy concerns are small ($\gamma \rightarrow 0$) in the main text.

3.3 Validation

Finally, we formalize how different kinds of *ex post* validation affect what the DM knows (\mathcal{I}_{DM2}) when forming a belief about the expert competence. For our regulation example, there are several ways the DM might later learn things about the state or difficulty of the question.

First, there may be studies in progress which will later clearly indicate whether the chemical is harmful. Alternatively, if the question is whether the chemical has medium or long term health consequences, this may not become clear until after the initial regulation decisions are made. To pick a prominent contemporary example, the degree to which we should regulate “vaping” of nicotine and cannabis depends on the long-term health consequences of this relatively new technology — perhaps relative to smoking — which are not currently known.

Second, the DM could consult other experts (who may themselves have strong or weak career concerns). These other experts could be asked about the state of the world, or perhaps the quality of the evidence on dimensions the DM may not be able to assess (Anecdotal evidence or scientific? Randomized control trials or observational studies? Have they been replicated? Human or animal subjects?).

In other contexts, validation about the state of the world could naturally arise if it is about an event that will later be realized (the winner of an election, the direction of a stock’s movement), but where the DM must make a choice before this realization. Alternatively, in many policy or business settings, the decision-maker may be able to validate the expert’s message directly, whether through experimental A/B testing or observational program evaluation. Closer to difficulty validation is the familiar notion of “peer review,” whereby other experts evaluate the feasibility of an expert’s design without attempting the question them-

selves. Another possibility is if the decision-maker has expertise (substantive or methodological) in the general domain of the question but does not have the time to investigate herself, and so she may be able to validate whether the question was answerable only after seeing what the expert comes up with.

Alternatively, subsequent events may reveal auxiliary information about whether the state should have been knowable, such as an extremely close election swayed by factors which should not have been *ex ante* predictable (i.e., this reveals that forecasting the winner of the election was a difficult question).

Formally, we consider several variations on \mathcal{I}_{DM2} . In all cases, the structure of \mathcal{I}_{DM2} is common knowledge.

In the *no validation* case, $\mathcal{I}_{DM2} = (m)$. This is meant to reflect scenarios where it is difficult or impossible to know the counterfactual outcome had the decision-maker acted differently. The *state validation* case, $\mathcal{I}_{DM2} = (m, \omega)$, reflects a scenario in which the DM can check the expert's advice against the true state of the world. In the *difficulty validation* case, $\mathcal{I}_{DM2} = (m, \delta)$, meaning that the DM learns whether the question was hard, i.e. whether the answer could have been learned by a good expert. In the *full validation* case, $\mathcal{I}_{DM2} = (m, \omega, \delta)$, the DM learns both the state and the difficulty of the question.

To be clear, many of the motivating examples blend validation about the state and difficulty. We consider the cases of “pure” difficulty or state validation to highlight what aspects of checking expert claims are most important for getting them to admit uncertainty.

4 Equilibrium Definition and Properties

The standard solution concept for a model like ours (with sequential moves and incomplete information) is Perfect Bayesian Equilibrium (PBE). While our central conclusions hold when using this solution concept, they become much clearer when we add a refinement which builds on the Markov restriction common to dynamic games of complete information (Maskin and Tirole, 2001). In this section we define the equilibrium concept in the context of our game; Appendix A contains a more general definition and detailed discussion of how the results change when using PBE. Appendix A also contains a discussion of past usage of related refinements; in short, there are several applied papers which use the same restriction of strategies that we employ, but to our knowledge this is the first paper to make use of the refinement to beliefs that the Markov strategies restriction implies.

4.1 Markov Sequential Equilibrium

The Markov restriction in repeated games requires that if two histories of play (h_1 and h_2), result in a strategically equivalent scenario starting at both histories (meaning that the players have the same expected utilities over the actions taken starting at both h_1 and h_2), then the players must use the same strategies starting in both histories. The analogous restriction in our setting has implications for strategies and beliefs.

In terms of strategies, we require that if two *types* of expert face a strategically equivalent scenario, they must play the same strategy. Most important to our model, in some cases a competent but uninformed expert and an incompetent expert have the exact same expected utility (for any DM strategy). PBE would allow these two types to play different strategies despite the fact that they face identical incentives; the Markov restriction requires them

to play the same strategy. Our restriction to beliefs essentially assumes that, even when observing an off-path message, the DM still believes that expert still plays some Markov strategy.

Formally, let $\sigma_{\theta,s}(m)$ be the probability of sending message m as a function of the sender type, $\pi_1(m)$ be the posterior belief that the state is 1 given message m , $\pi_g(m; \mathcal{I}_{DM2})$ be the posterior belief about the expert competence, and $a(m)$ be the policy action given message m . Let U_E and U_{DM} be the expected utilities for each player.

PBE requires that both players maximize their utility given the other's strategy and their beliefs, and that these beliefs are formed by Bayes' rule when possible. To these we add two requirements:

Definition 1. *A Markov Sequential Equilibrium to the model is a PBE which also meets the following requirements:*

- *(Expert Markov Strategies): If there are two types (θ', s') and (θ'', s'') such that $U_E(m; \theta', s') = U_E(m; \theta'', s'')$ for all m given the DM strategies and beliefs, then $\sigma_{\theta',s'}(m) = \sigma_{\theta'',s''}(m)$ for all m .*
- *(DM Markov Consistency): For all m and \mathcal{I}_{DM2} , there exists a sequence of non-degenerate Markov strategies for the expert σ^k , with corresponding beliefs formed by Bayes Rule $\pi_1^k(m)$ and $\pi_g^k(m, \mathcal{I}_{DM2})$, such that $\pi_1(m) = \lim_{k \rightarrow \infty} \pi_1^k(m)$ and $\pi_g(m, \mathcal{I}_{DM2}) = \lim_{k \rightarrow \infty} \pi_g^k(m, \mathcal{I}_{DM2})$.*

The key implication of Markov consistency is to rule out off-path inferences about payoff-irrelevant information, because off-path beliefs which condition on payoff-irrelevant information can not be reached by a sequence of Markov strategies. Our restriction is related to that implied by D1 and the intuitive criterion (Cho and Kreps, 1987). However, where

these refinements require players to make inferences about off-path play in the presence of strict differences of incentives between types, our restriction rules out inference about types in the absence of strict differences of incentives.

We use this solution concept for two reasons. The first is theoretical. People have endless “private information” which they could, in principle condition their behavior on. In a more complex but realistic version of our model, the expert may have private information about not only what he learns about the chemical in question, but his personal policy preferences, views on what kinds of scientific evidence are credible, and what he had for breakfast in the morning. When observing an off-path message, PBE would allow for updating (or not) on any of these dimensions, regardless of their relevance to the interaction at hand. The Markov strategies restriction provides a principled and precise justification for which kinds of private information experts can condition their strategies on.⁷ (Beliefs about the effects of the chemical in question? Usually. What evidence is credible? Sometimes. Breakfast? Rarely.⁸) The Markov beliefs restriction is then a logical implication of what observers can update on when observing off-path messages.

The second is more practical. As we show in Section 5, our main result about the importance of difficulty validation for getting experts to admit uncertainty is nearly immediate when using this restriction. As demonstrated in Appendix A, similar results hold when analyzing PBE, but since the set of PBE is much larger, the comparisons are not as clean.

⁷Of course, a classic interpretation of mixed strategies is that the sender conditions on some random (and “irrelevant”) information like a coin flip. However, as discussed in Appendix A, if mixed strategies are viewed as the limit of pure strategies with payoff perturbations a la Harsanyi (1973), then only Markov strategies are possible.

⁸Though see Cho and Kreps (1987, section II).

4.2 Properties of Equilibria

Since we allow for a generic message space, there will always be many equilibria to the model even with the Markov restrictions. To organize the discussion, we will focus on how much information can be conveyed (about ω and θ). On one extreme, we have babbling equilibria, in which all types employ the same strategy, and the DM learns nothing.

On the other extreme, there is never an equilibrium with full separation of types. To see why, suppose there is a message that is *only* sent by the good but uninformed types $m_{g,\emptyset}$ (“I don’t know because the optimal policy isn’t clear”) and a different message only sent by the bad uninformed types $m_{b,\emptyset}$ (“I don’t know because I am incompetent”). If so, the policy choice upon observing these messages would be the same. However, the reputation payoff for sending $m_{g,\emptyset}$ is strictly higher, and so the bad types have an incentive to deviate.

Still, such full separation is not required for the expert to communicate what he knows about the *state*. That is, there may still be equilibria where the uninformed types say “I don’t know” (if not why), and the informed types report the state of the world. We call these *honest* equilibria. In the main text we focus on the simplest version of an honest equilibrium, where there is a unique message m_x sent by each type observing s_x with probability 1, $x \in \{0, 1, \emptyset\}$; see Appendix B for a more general definition. This is a particularly important class of equilibria in our model as it conveys the most information about the state:

Proposition 1. *The expected value of the policy in an honest equilibrium is $p_g p_e + (1 - p_g p_e)(1 - p_1(1 - p_1)) \equiv \bar{v}$, which is strictly greater than the expected value of the policy in any equilibrium which is not honest.*

Proof. Unless otherwise noted, all proofs are in Appendix B. □

This result formalizes our intuition that it is valuable for the DM to learn when the expert is uninformed, and follows directly from the convexity of the loss function for bad policies. For example, if the expert on environmental policy sometimes says that a chemical is harmful (or is not harmful) when the truth is that he isn't sure, the regulator will never be entirely sure what the optimal policy is. Her best response to this garbled message is to not regulate as aggressively as she would if knowing the chemical is harmful for sure, even when the expert fully knows that this is the case.

As we will see, honest equilibria tend to fail when one of the uninformed types would prefer to “guess”, mimicking the message of the informed types. So, the other equilibria we consider in the main text are ones where the informed types still send their respective messages m_0 and m_1 , but the uninformed types at least sometimes send one of these messages. We refer to sending one of these messages as “guessing”, and sending m_0 as “admitting uncertainty”. See Appendix B for a definition of what it means to admit uncertainty with more general messaging strategies, and Appendix D for an extensive discussion of why focusing on this class of messaging strategies sacrifices no meaningful generality for our results.

Combined with proposition 1, these definitions highlight why admission of uncertainty is important for good decision-making. In any equilibrium with guessing, the fact that the uninformed types send messages associated with informed types leads to worse policies than in an honest equilibrium. This is for the two reasons highlighted in our opening paragraph. First, when the expert actually is informed, his advice will be partially discounted by the fact that the DM knows some uninformed types claim to know which policy is best. That is, decision-makers may ignore good advice. Second, when the expert is uninformed, he will induce the DM to take more decisive action than the expert's knowledge warrants; i.e., decision-makers may *take* bad advice.

5 When is Admission of Uncertainty Possible?

Between the four possible validation regimes and whether the expert exhibits policy concerns, there are many cases of the model to analyze. In this section we show that the Markov restrictions and some simple incentive compatibility constraints quickly allow us to see when admission of uncertainty is impossible, regardless of the particular parameter values. We then provide a complete analysis of one case where admission of uncertainty is possible but not guaranteed for more subtle comparative static results. In other words, we first highlight our negative results about state validation (and discuss why they do not afflict difficulty validation), and then derive concrete positive results about difficulty validation in the next section. Appendix C contains a full analysis of the remaining cases.

Markov sequential equilibrium has two main implications: Markov strategies, which requires that payoff-irrelevant information cannot affect equilibrium play, and Markov consistency, which requires that off-path beliefs cannot update on payoff-irrelevant information. To see the immediate implications of these restrictions, it is helpful to construct the classes of payoff-equivalent information sets. We put information sets in the same payoff equivalence class if experts at those decision nodes are payoff-equivalent *for any DM strategy*.⁹ Figure 2 illustrates for the case with no policy concerns. Each row represents an assumption on the DM's information set at the end of the game, \mathcal{I}_{DM2} . Each column represents one of the four information sets depicted in Figure 1.

No Validation, $\gamma = 0$. First, in the no validation (NV) case, $\mathcal{I}_{DM2} = (m)$. Since the DM only sees the message when evaluating the expert, if the expert has no policy concerns his private signal is payoff-irrelevant. Therefore there is a single payoff equivalence class

⁹Any two information sets can be payoff equivalent for *some* DM strategy: e.g., if she always picks the same policy and competence assessment for all messages.

Figure 2: Payoff Equivalence Classes With No Policy Concerns

NV	b, \cdot, \cdot	g, h, \cdot	$g, e, 0$	$g, e, 1$
SV	b, \cdot, \cdot	g, h, \cdot	$g, e, 0$	$g, e, 1$
DV	b, \cdot, \cdot	g, h, \cdot	$g, e, 0$	$g, e, 1$
FV	b, \cdot, \cdot	g, h, \cdot	$g, e, 0$	$g, e, 1$

Notes: This figure depicts equivalence classes under each validation regime for the case with no policy concerns. Each row represents a validation regime: respectively, no validation, state validation, difficulty validation, and full validation. Each column represents an expert information set.

comprised of all four information sets, as depicted in the first row of Figure 2, and the Markov strategies restriction implies that all experts play the same strategy. In this case, all that is left is a babbling equilibrium:

Proposition 2. *With no validation and no policy concerns (i.e., $\gamma = 0$), any MSE is babbling, and there is no admission of uncertainty.*

This highlights the importance of the Markov strategies restriction. In our main example, if the environmental expert does not care at all about whether a good policy choice is made (unlikely if he is a citizen who is affected by the policy) *and* there is no check on his claim about whether the chemical is harmful (again, unlikely), the information he has is not payoff relevant. We should not be surprised that this extreme premise leads to an extreme prediction that such an expert will not convey any useful information.

State Validation, $\gamma = 0$. A promising direction is to allow the DM to observe ω when forming beliefs about θ . Doing so breaks payoff equivalence between types with different information about the state, as we illustrate in the second row of Figure 2. This partition

of payoff equivalence is rich enough to support honesty in Markov strategies. Bad experts and uninformed good experts can pool on a message interpreted as “I don’t know”, and informed experts can send messages interpreted as “the optimal policy is zero” and “the optimal policy is one.”

To see this, consider the expected payoff for the expert with type and signal (θ, s) sending message m :

$$\sum_{\omega \in \{0,1\}} Pr(\omega|s, \theta) \pi_g(m, \omega).$$

Now, the informed types with different information about the state are not generally payoff-equivalent since $Pr(\omega|\theta, s)$ depends on the signal. However, good and bad uninformed experts, i.e. (g, s_\emptyset) and (b, s_\emptyset) , are always payoff equivalent since $Pr(\omega|g, s_\emptyset) = Pr(\omega|b, s_\emptyset)$. So, the Markov *strategies* restriction does not preclude admission of uncertainty, or even an honest equilibrium. However, the Markov consistency restriction does.

To see why, consider the simplest case of an honest equilibrium where types observing signal s_x send message m_x . In such an equilibrium, the decision-maker picks a stringent regulation when the expert says the chemical is harmful ($a = 1$ when observing m_1), no regulation when the chemical is not harmful ($a = 0$ when observing m_0), and intermediate regulation when the expert admits not knowing whether the chemical is harmful ($a = p_1$ when observing m_\emptyset). The on-path information sets include cases where a good expert makes an accurate recommendation, $(m_0, 0)$ and $(m_1, 1)$, and cases where an uninformed expert (good or bad) says “I don’t know” along with either validation result: $(m_\emptyset, 0)$, and $(m_\emptyset, 1)$. When validation indicates the expert was right about whether the chemical is harmful (i.e., $(m_0, 0)$ or $(m_1, 1)$), the DM knows the expert is good, i.e., $\pi_g(m_i, i) = 1$ for $i \in \{0, 1\}$. When observing m_\emptyset and either validation result, the belief about the expert

competence is:

$$\begin{aligned}
\pi_g(m_\emptyset, \omega) &= \frac{Pr(\theta = g, \delta = h)}{Pr(\theta = g, \delta = h) + Pr(\theta = b)} \\
&= \frac{p_g(1 - p_e)}{p_g(1 - p_e) + 1 - p_g} \\
&\equiv \pi_g^\emptyset.
\end{aligned}$$

The term π_g^\emptyset , which recurs frequently throughout the analysis, represents the share of uninformed types who are competent but asked a hard question.

For the uninformed types, the expected payoff for sending m_\emptyset in the honest equilibrium is π_g^\emptyset . Since $0 < \pi_g^\emptyset < p_g$, the expert revealing himself as uninformed leads to a lower belief about competence than the prior, but is not zero, since there are always competent but uninformed types.

Consider a deviation to m_1 , i.e., claiming to know that the chemical is harmful. When validation reveals this to be true, the DM observes $(m_1, 1)$, and believes the expert to be competent with probability 1. When validation reveals the chemical is not harmful, the DM observes $(m_1, 0)$, which is off-path.

However, MSE places some restriction on this belief, as it must be the limit of a sequence of beliefs consistent with a sequence of Markov strategies. Since the good and bad uninformed types are payoff equivalent and play the same strategy, the worst inference that the DM can make about the expert when observing an off-path message/validation combination is that the expert was uninformed, i.e., $\pi_g(m_1, 0) \geq \pi_g^\emptyset$ (see the proof of proposition 3). Given this restriction on the off-path belief, in any honest MSE the payoff to sending m_1 must be

at least:

$$p_1 + (1 - p_1)\pi_g^\theta > \pi_g^\theta.$$

The expert can look no worse from guessing that the chemical is harmful and being incorrect than he would when just admitting he is uncertain. Since there is a chance to look competent when guessing and being correct, the expert will always do so. This means there is always an incentive to deviate to m_1 (or, by an analogous argument, m_0), and hence no honest equilibrium.

A related argument implies that there is no MSE where the uninformed types *sometimes* admit uncertainty (i.e., $\sigma_\emptyset(m_\emptyset) \in (0, 1)$) and sometimes guess m_0 or m_1 : guessing and being correct always gives a higher competence evaluation than incorrectly guessing, which gives the same competence evaluation as sending m_\emptyset . So, guessing gives a strictly higher payoff than admitting uncertainty.

Proposition 3. *With state validation and no policy concerns, there is no MSE where an expert admits uncertainty.*

As shown in the proof of the proposition, there is an MSE where the informed types reveal their information: e.g., the experts who know the chemical is harmful report this, and those who know it is not harmful say so. However, all of the uninformed experts – competent or not – will make a guess at whether the chemical is harmful. The equilibrium condition is that their guessing probabilities are such that observing a claim that the chemical is harmful or not leads to the same belief about the expert competence. So, state validation does improve communication in general relative to no validation, but not in terms of admitting uncertainty.

Figure 3: Payoff Equivalence Classes With Policy Concerns

NV	b, \cdot, \cdot	g, h, \cdot	$g, e, 0$	$g, e, 1$
SV	b, \cdot, \cdot	g, h, \cdot	$g, e, 0$	$g, e, 1$
DV	b, \cdot, \cdot	g, h, \cdot	$g, e, 0$	$g, e, 1$
FV	b, \cdot, \cdot	g, h, \cdot	$g, e, 0$	$g, e, 1$

Notes: This figure depicts equivalence classes under each validation regime for the case with policy concerns. Each row represents a validation regime: respectively, no validation, state validation, difficulty validation, and full validation. Each column represents an expert information set, as derived in Figure 1.

Further, as shown in Appendix C, adding policy concerns (to the no validation or state validation case) will not solve this problem unless the expert cares so much about the policy as to accept the hit to his reputation from admitting uncertainty.

Difficulty, Full Validation, and Policy Concerns. For the DM to effectively threaten punitive off-path beliefs, we need to break the payoff equivalence of bad types and good but uninformed types, and this is precisely what difficulty validation (DV) does, depicted in the third row of Figure 2. However, difficulty validation is not enough to sustain honesty, since (unlike state validation) it does not break the payoff equivalence between the informed experts who actually know whether the chemical is harmful.

This we view as a more minor problem, which can be solved by either combining state and difficulty validation (FV), as in the fourth row of Figure 2, or by adding small policy concerns, which yields payoff equivalence classes represented in Figure 3. We formally prove results about when (complete) communication of uncertainty is possible in the next section.

Summary. From the payoff equivalence classes alone, we now know what cases at least allow for the possibility of an honest equilibrium without strong policy concerns. First, we need to break the payoff equivalence between the types with different information about whether the chemical is harmful, which can be accomplished with either state validation or any policy concerns. Second, we need to break the payoff equivalence between good and bad uninformed experts, which, given the natural way we have set up the problem, can *only* be accomplished with difficulty validation (or full validation, which includes difficulty validation).

What remains is to check when honesty is in fact possible. In the next section, we do this for the “hardest” case, with only difficulty validation and with small policy concerns. As shown in Appendix C, the insights from this analysis are similar to what we obtain with larger policy concerns and/or full validation.

6 Analysis of Our Main Case

While the previous section shows that difficulty validation is necessary for the admission of uncertainty (absent large policy concerns), we have not yet shown when it is sufficient. This section explores when difficulty validation, combined with small policy concerns, is sufficient for honesty, or at least some admission of uncertainty.

We analyze difficulty validation rather than full validation in the main text not because we believe the former is necessarily more common, but to show circumstances under which the minimal validation regime is sufficient to induce admission of uncertainty. Again, in Appendix C we show similar results arise with full validation.

The focus on small policy concerns is for both technical and substantive reasons. If the

expert has exactly no policy concerns, then this renders the two informed types (i.e., those who know the chemical is harmful or not) payoff equivalent, which greatly undermines the amount of information which can be transmitted in an MSE. However, this is fragile to the small and entirely realistic perturbation where the expert has any non-zero concern about the quality of the policy. For example, even if the chemical in question only has a potentially small impact on the environment, the expert himself could be harmed by exposure. So, having small policy concerns allows the informed types to reveal their information honestly, while simplifying the analysis of the potential deviations for uninformed types since they (in the limit as $\gamma \rightarrow 0$) will send the message which maximizes their reputation for competence.

This also hints at the more substantive justification. In many, if not most policy-making domains, the effect of the policy change that experts will feel in their own personal life likely pales in comparison to their concern about perceptions of their competence, which can affect whether they keep their job or will be hired in the future. In general, we expect that our analysis applies more to “big” policy questions that affect many people, and where experts are very specialized and care about perceptions of their competence in that particular domain.

6.1 Equilibrium

We continue to study the most straightforward class of messaging strategies where the informed types send one message each which we give the natural labels m_0 and m_1 , and the uninformed types either send one of these messages or a third message labeled m_\emptyset (“I Don’t Know”).¹⁰

¹⁰ See Appendix D for an extensive discussion of the sense in which this is without loss of generality.

In an honest equilibrium, messages m_0 and m_1 are only sent by informed types. So, when observing these messages, the DM picks actions $a = 0$ and $a = 1$, respectively. At the evaluation stage, when observing (m_0, e) or (m_1, e) , the DM knows the expert is competent with certainty.

Upon observing m_\emptyset , the DM picks policy p_1 . The competence assessment when the expert says “I don’t know” depends on the result of the validation. When the problem is easy, the DM knows that the expert is incompetent since this fact means that a competent expert would have learned whether the chemical is harmful and sent an informative message. So, $\pi_g(m_\emptyset, e) = 0$. Upon observing (m_\emptyset, h) , the DM learns nothing about the expert competence since no expert gets an informative signal when the scientific literature is uninformative.

Combining these observations, the payoff to the good but uninformed type for sending m_\emptyset (who knows the validation will reveal $\delta = h$) is

$$p_g + \gamma(1 - p_1(1 - p_1)).$$

The bad uninformed type does not know if the validation will reveal the problem is hard, and so receives a lower expected competence evaluation and hence payoff for sending m_\emptyset :

$$(1 - p_e)p_g + \gamma(1 - p_1(1 - p_1)).$$

Since no types are payoff-equivalent, the Markov consistency requirement places no restrictions on off-path competence evaluations, and we can set these to zero.¹¹ In the case with only difficulty validation, these are the information sets (m_0, h) and (m_1, h) , i.e., get-

¹¹This belief can be reached by assuming a sequence of strategies where only the bad type sends a particular message, but with a probability which approaches zero.

ting caught guessing when validation reveals that the scientific literature is not informative. Importantly, the expert is not caught because the DM realizes his claim is wrong, but because she comes to learn that the question was not answerable. As formalized below, such an off-path belief is also “reasonable” in the sense that the bad uninformed types face a strictly stronger incentive to guess than the good uninformed types.

If the DM believes that the expert is bad with probability one upon observing (m_0, h) or (m_1, h) , then a good but uninformed type knows he will get a reputation payoff of zero if sending either of these messages. Further, if claiming he knows whether the chemical is harmful, the policy is worse as well, so he has no incentive to deviate. A bad type guessing that the chemical is harmful (m_1) gets expected payoff:

$$p_e + \gamma p_1.$$

which is strictly higher than the payoff for sending m_0 . So, the constraint for an honest equilibrium is:

$$(1 - p_e)p_g + \gamma(1 - p_1(1 - p_1)) \geq p_e + \gamma p_1$$

$$\gamma \geq \frac{p_e(1 + p_g) - p_g}{(1 - p_1)^2}.$$

Not surprisingly, when policy concerns are very strong, uninformed experts will admit uncertainty since it leads to better policy choices. In fact, this holds for any validation regime; see Appendix C for a comparison of “how strong” policy concerns must be to induce honesty for each case. However, this analysis also highlights that when policy concerns are small relative to reputation concerns, there is never an honest equilibrium with no validation or just state validation. In contrast, with just difficulty validation,¹² as

¹²Unsurprisingly, this constraint is even easier to meet with full validation; the key point is that difficulty validation is *necessary* for honesty with small policy concerns, and sometimes sufficient.

$\gamma \rightarrow 0$ there is an honest equilibrium when:

$$p_e \leq \frac{p_g}{1 + p_g}. \quad (5)$$

This inequality implies that an honest equilibrium is easy to sustain when p_e is low, and p_g is high. The former holds for two reasons: a prior belief that the literature is likely not informative about the chemical means that uninformed experts are likely to be competent, and the uninformed expert is more likely to be “caught” claiming to have an answer to an impossible problem. More competent experts (high p_g) makes honesty easier as it means the uninformed types are frequently competent, making admitting uncertainty look less bad.

What if (5) does not hold? Of course, there is always a babbling equilibrium, and there is also an always guessing equilibrium where the good and bad uninformed types always send m_0 or m_1 . More interesting for our purposes, there can also be an MSE where all of the good types report honestly and the bad types play a mixed strategy over (m_0, m_1, m_\emptyset) .¹³ There is a more subtle incentive compatibility constraint that must be met for this equilibrium to hold: if the bad types are indifferent between sending m_0 and m_1 , it can be the case that the *informed* types prefer to deviate to sending the other informed messages (i.e., the s_1 type prefers to send m_0) when policy concerns get very small. See the proof in Appendix C for details; in short, if the probability of a solvable problem is not too low or the probability of the state being one is not too high, then this constraint is not violated and there is an MSE where all of the good types send their honest message.¹⁴

¹³ This is the only MSE where the good types report honestly, and we conjecture that the equilibrium we check for (when it exists) maximizes both the probability that the expert admits uncertainty and the expected value of the decision. It is hypothetically possible, though we believe unlikely, that there could be an MSE where the good types do not report honestly and this induces the bad types to admit uncertainty more often in a manner which outweighs the loss of information from the good types.

¹⁴The proof of the proposition discusses how the intuition behind this constraint.

Proposition 4. *As $\gamma \rightarrow 0$ with difficulty validation, there is an honest MSE if and only if $p_e \leq \frac{p_g}{1+p_g}$. If not, and $p_e \geq 2p_1 - 1$, then there is an MSE where the good types send their honest message and the bad types use the following strategy:*

$$\sigma_b^*(m_\emptyset) = \begin{cases} \frac{1-p_e(1+p_g)}{1-p_g} & p_e \in \left(\frac{p_g}{1+p_g}, \frac{1}{1+p_g} \right), \\ 0 & p_e > \frac{1}{1+p_g}, \end{cases}, \quad (6)$$

$$\sigma_b^*(m_0) = (1 - p_1)(1 - \sigma_b^*(m_\emptyset)),$$

$$\sigma_b^*(m_1) = p_1(1 - \sigma_b^*(m_\emptyset)).$$

Proof. This is a special case of the more complete characterization of a class of MSE with difficulty validation and policy concerns stated and proven as proposition 12 in Appendix C. □

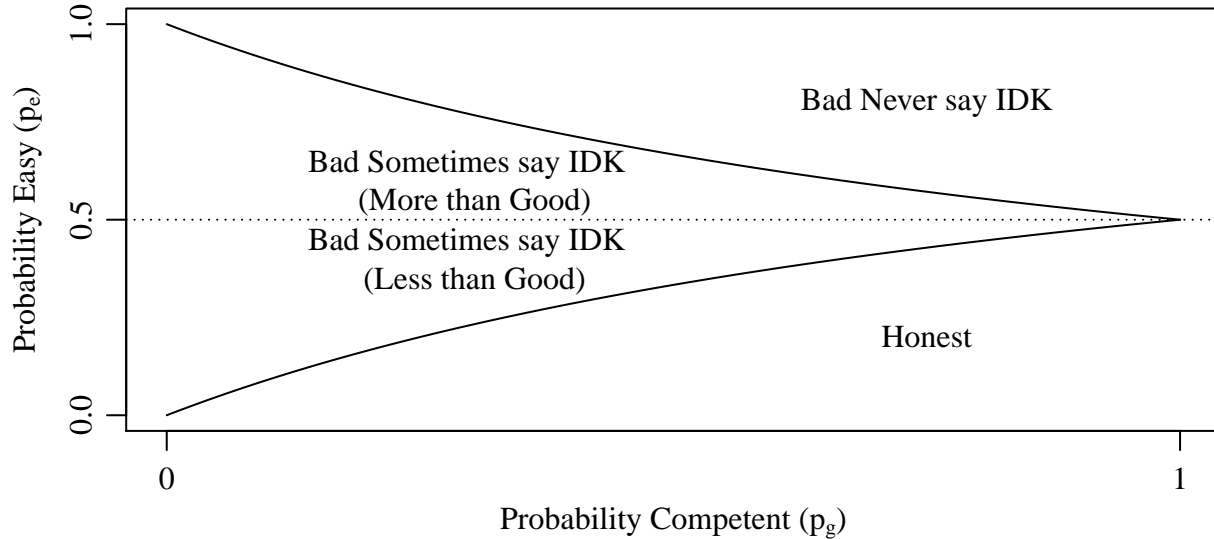
In sum, other than in a restrictive corner of the parameter space,¹⁵ small policy concerns and difficulty validation are sufficient to induce at least good and uninformed experts to admit uncertainty, and potentially full honesty among all experts.

6.2 When do we observe admission of uncertainty, and from whom?

Figure 4 illustrates the equilibrium described by proposition 4, which helps provide some more concrete insights into when we should expect to see admission of uncertainty and good decisions. All claims that follow are comparisons within this equilibrium, though again this is the only MSE where the good types report honestly (see footnote 13).

¹⁵The only time this does not hold is if $\frac{p_g}{1+p_g} < p_e$, and $p_e < 2p_1 - 1$. So, three sufficient conditions for this constraint to remain unviolated are (1) p_g is high, (2) p_1 is high, and (3) p_e is either low or high. Even if both of these inequalities hold, as shown in the proof of proposition 12 in Appendix C, it is a necessary but not sufficient condition for an incentive compatibility constraint to be violated.

Figure 4: Illustration of expert strategies in the equilibrium identified by proposition 4, as a function of p_e and p_g , with $p_1 = 1/2$.



In the bottom right corner (most experts are competent, most problems are hard), there is an honest equilibrium. Substantively, this corner of the parameter space plausibly corresponds to environments where the best experts are tackling questions on the research frontier, be it at conferences or in the pages of academic journals. In general, the admission of uncertainty is quite frequent in these settings: even good experts often don't know the answer but feel comfortable admitting this, and as a result bad experts can pool with them and admit uncertainty too. So, while decision-makers may not always get useful advice here – recall, most questions have no good answer — they will at least have no problems getting experts to honestly report what they know.

In the top right corner, most experts are competent and most problems are easy. We can think of this part of the parameter space as “business as usual” in bureaucracies, companies or other organizations where qualified experts are addressing relatively mundane problems that they generally know how to solve. In this case, there is little admission of uncertainty, both because experts typically get informative signals, and because admitting uncertainty

in this setting is too risky for bad experts to ever do it. While experts have the most information in this region, this comes at a cost from the perspective of the decision-maker, as bad types always guess, which dilutes the value of the informative messages.

Comparing across these two cases also hints at the question of when good or bad experts are more likely to admit uncertainty. Our model contains a straightforward assumption about which experts are likely to *be* uncertain about the state of the world: good experts sometimes, and bad experts always. And, in an honest equilibrium (when $p_e < p_g/(1+p_g)$), this is what they will report, and so bad experts admit uncertainty more often.

However, on the other extreme, when $p_e > 1/(1+p_g)$, it is *only* the good experts who will ever admit uncertainty. Put another way, if an outside observer (who the expert does not care to impress) were to see an expert admit uncertainty in this part of the parameter space (and in this equilibrium), she would know for sure that anyone who says “I don’t know” is in fact competent.

More generally, good experts admit uncertainty whenever the problem is hard, so with probability $1 - p_e$. For parameters where good experts sometimes admit uncertainty (the triangle in the left of Figure 4), bad experts send m_\emptyset with probability $\frac{1-p_e(1+p_g)}{1-p_g}$, which is less than $1 - p_e$ if and only if $p_e > 1/2$. So, in domains of difficult problems, bad experts are more likely to admit uncertainty, and in domains with easier questions good experts are more likely to admit uncertainty.

6.3 Comparative Statics

We now ask how changing the probability parameters of the model affects the communication of uncertainty by uninformed experts in the MSE identified by proposition 4.

In general, one might expect that adding more competent experts will lead to less admission of uncertainty and better decisions, and that making problems easier will also lead to less admission of uncertainty and better decisions. Both of these always hold within an honest equilibrium, but can break down elsewhere. Here we highlight cases where these intuitive results do not hold, and then summarize with a complete description of the comparative statics.

More competent experts Adding more competent experts always leads to better decisions (on average), but there are two scenarios where adding more competent experts (increasing p_g) can lead to more admission of uncertainty.

First, in the region where the bad types always guess, the good types are the only experts who send m_\emptyset , so increasing p_g leads to more admission of uncertainty.

More subtly, when the bad types sometimes send m_\emptyset and $p_e < 1/2$ (the bottom left triangle), adding more competent experts leads the bad types to admit uncertainty more often. And since this is the part of the parameter space where good types usually admit uncertainty as well, adding more competent types leads to more admission of uncertainty overall:

Proposition 5. *In the equilibrium identified by proposition 4:*

- (i) *The expected value of the decision is strictly increasing in p_g , and*
- (ii) *The unconditional probability that the expert admits uncertainty is strictly decreasing in p_g if $p_e \leq p_g/(1 + p_g)$ (honest equilibrium) or if $p_e \in (1/2, 1/(1 + p_g))$, and strictly increasing in p_g otherwise.*

More easy questions When the probability of an easy problem increases, this always decreases the admission of uncertainty as good types are more likely to be informed and bad types are more apt to guess. The potentially counterintuitive result here is that more

easy problems can sometimes lead to worse decisions. This is possible because when problems are more likely to be easy, the bad types are more tempted to guess. If the bad types never actually guess (bottom right corner) or always guess (top right corner), this does not matter. However, when making problems easier actually makes the bad types guess more, the messages m_0 and m_1 become less informative. As shown in the following proposition, this can sometimes lead to worse decisions:

Proposition 6. *In the equilibrium identified by proposition 4,*

(i) The unconditional probability that an expert admits uncertainty is strictly decreasing in p_e , and

(ii) For any p_g , there exists a $\tilde{p}_e \in (p_g/(1+p_g), 1/(1+p_g)]$ such that v^ is strictly decreasing in p_e for $p_e \in (p_g/(1+p_g), \tilde{p}_e)$.*

7 Discussion

This paper has studied the strategic communication of uncertainty by experts with reputation concerns. Our analysis is built on two theoretical innovations: first, in our setup, the decision-maker is uncertain not only about the state of the world, but also about whether the state is even knowable for qualified experts. This feature is related to prior work on classification of uncertainty; in particular, “aleatory” and “epistemic” uncertainty, language dating to Hacking (1975). Aleatory uncertainty characterizes what we cannot know, i.e. “difficult” questions (e.g., the roll of dice), where epistemic uncertainty is what we do not know, but could if we were more informed.¹⁶ We show that these properties of uncertainty have real implications, both in understanding why communication about uncertainty is hard and how to overcome that challenge. The way we formalize the distinction between easy and hard

¹⁶The language of aleatory and epistemic uncertainty has taken particular hold in structural engineering, where it is important in the classification of risk, see Kiureghian and Ditlevsen (2009).

problems highlights the idea that part of being a domain expert is not merely knowing the answers to questions, but knowing the limits of what questions are answerable.

A second innovation concerns the notion of “credible beliefs,” which is closely tied to structural consistency of beliefs (Kreps and Wilson, 1982). Honest communication in our model is disciplined by experts’ reputational concerns — off-path, they are punished by the low opinion of the decision-maker. But what can she credibly threaten to believe? Our use of Markov sequential equilibrium rules out non-credible beliefs that stipulate updating on payoff-irrelevant information.

A pragmatic way to frame our inquiry is that we ask what would the decision maker want to learn, *ex post*, in order to induce the experts, *ex ante*, to communicate their information honestly? We found that the intuitive answer – checking experts’ reports against the true state of the world – is insufficient. Even when decision-makers catch an expert red-handed in a lie, the severity of their beliefs is curtailed by the fact that good experts facing unanswerable questions are in the same conundrum as bad experts. Therefore, we show, state validation alone never induces honesty. In order to elicit honest reports from experts, it is necessary that the decision-maker also learns whether the problem is difficult. Indeed, in environments where the expert has even very small policy concerns, difficulty validation alone may be sufficient.

Is such difficulty validation common in the real world? As discussed throughout, we believe sometimes information about difficulty naturally comes *ex post*, and it can sometimes be accomplished by methods like peer review. We conclude with a practical suggestion: that when consulting multiple experts, decision-makers may want to give heterogeneous incentives and ask different questions. Rewarding some for “getting things right” gives good incentives to avoid letting personal biases contaminate advice, and potentially for collecting information in order to be informed in the first place. However, as emphasized here,

these kinds of reward schemes may exacerbate the problem of getting experts to admit uncertainty. On the other hand, paying a flat fee to an expert who will likely not be consulted again for their services may have drawbacks, but will make the expert much more comfortable admitting uncertainty. Further, some experts can simply be asked “do you think the evidence about this question is solid” rather than emphasizing what the expert thinks the truth is. Finding other ways to achieve difficulty validation could be a path to improving communication in politics and organizations more generally.

References

- Ashworth, S. (2012). Electoral accountability: recent theoretical and empirical work. *Annual Review of Political Science*, 15:183–201.
- Austen-Smith, D. (1990). Information transmission in debate. *American Journal of Political Science*, 34(1):124–152.
- Banks, J. S. (1990). A model of electoral competition with incomplete information. *Journal of Economic Theory*, 50(2):309–325.
- Bergemann, D. and Hege, U. (2005). The financing of innovation: Learning and stopping. *RAND Journal of Economics*, 36(4):719–752.
- Bergemann, D. and Hörner, J. (2010). Should auctions be transparent? Cowles Foundation Discussion Paper No. 1764.
- Bhaskar, V., Mailath, G. J., and Morris, S. (2013). A foundation for markov equilibria in sequential games with finite social memory. *Review of Economic Studies*, 80(3):925–948.
- Callander, S. (2011). Searching for good policies. *American Political Science Review*, 105(4):643–662.
- Calvert, R. L. (1985). The value of biased information: A rational choice model of political advice. *The Journal of Politics*, 47(2):530–555.
- Canes-Wrone, B., Herron, M. C., and Shotts, K. W. (2001). Leadership and pandering: A theory of executive policymaking. *American Journal of Political Science*, 45(3):532–550.
- Cho, I.-K. and Kreps, D. M. (1987). Signaling games and stable equilibria. *The Quarterly Journal of Economics*, pages 179–221.
- Crawford, V. P. and Sobel, J. (1982). Strategic information transmission. *Econometrica: Journal of the Econometric Society*, pages 1431–1451.
- Deb, R., Pai, M., and Said, M. (2018). Evaluating strategic forecasters. forthcoming, *American Economic Review*.
- Dessein, W. (2002). Authority and communication in organizations. *The Review of Economic Studies*, 69(4):811–838.
- Dye, R. A. (1985). Disclosure of nonproprietary information. *Journal of Accounting Research*, 23(1):123–145.
- Ericson, R. and Pakes, A. (1995). Markov-perfect industry dynamics: A framework for empirical work. *Review of Economic Studies*, 62(1):53–82.

- Fearon, J. D. (1999). Electoral accountability and the control of politicians: selecting good types versus sanctioning poor performance. In Przeworski, A., Stokes, S. C., and Manin, B., editors, *Democracy, accountability, and representation*. Cambridge University Press.
- Fehrler, S. and Janas, M. (2019). Delegation to a group. Manuscript.
- Fox, J. and Shotts, K. W. (2009). Delegates or trustees? a theory of political accountability. *The Journal of Politics*, 71(4):1225–1237.
- Fudenberg, D. and Tirole, J. (1991). *Game Theory*. MIT Press.
- Gailmard, S. and Patty, J. W. (2012). Formal models of bureaucracy. *Annual Review of Political Science*, 15:353–377.
- Gailmard, S. and Patty, J. W. (2013). Stovepiping. *Journal of Theoretical Politics*, 25(3):388–411.
- Gilligan, T. W. and Krehbiel, K. (1987). Collective decisionmaking and standing committees: An informational rationale for restrictive amendment procedures. *Journal of Law, Economics, & Organization*, 3(2):287–335.
- Hacking, I. (1975). *The Emergence of Probability: A Philosophical Study of Early Ideas about Probability Induction and Statistical Inference*. Cambridge University Press.
- Harsanyi, J. C. (1973). Games with randomly disturbed payoffs: A new rationale for mixed-strategy. *International Journal of Game Theory*, 2(1):1–23.
- Harsanyi, J. C. and Selten, R. (1988). *A General Theory of Equilibrium Selection in Games*. MIT Press, Cambridge, MA.
- Holmström, B. (1999). Managerial incentive problems: A dynamic perspective. *The Review of Economic Studies*, 66(1):169–182.
- Hübner, R. (2019). Getting their way: Bias and deference to trial courts. *American Journal of Political Science*, 63(3):706–718.
- Jung, W.-O. and Kwon, Y. K. (1988). Disclosure when the market is unsure of the information endowment of managers. *Journal of Accounting Research*, 26(1):146–153.
- Kartik, N. (2009). Strategic communication with lying costs. *Review of Economic Studies*, 76:1359–1395.
- Kiureghian, A. D. and Ditlevsen, O. (2009). Aleatory or epistemic? does it matter? *Structural Safety*, 31:105–112.
- Kreps, D. M. and Ramey, G. (1987). Structural consistency, consistency, and sequential rationality. *Econometrica*, 55(6):1331–1348.
- Kreps, D. M. and Wilson, R. (1982). Sequential equilibria. *Econometrica*, 50(4):863–894.

- Kydd, A. (2003). Which side are you on? bias, credibility, and mediation. *American Journal of Political Science*, 47(4):597–611.
- Leaver, C. (2009). Bureaucratic minimal squawk behavior: Theory and evidence from regulatory agencies. *American Economic Review*, 99(3):572–607.
- Little, A. T. (2017). Propaganda and credulity. *Games and Economic Behavior*, 102:224–232.
- Manski, C. F. (2019). Communicating uncertainty in policy analysis. *Proceedings of the National Academy of Sciences*, 116(16):7634–7641.
- Maskin, E. and Tirole, J. (1988a). A theory of dynamic oligopoly, i: Overview and quantity competition with large fixed costs. *Econometrica*, 56(3):549–569.
- Maskin, E. and Tirole, J. (1988b). A theory of dynamic oligopoly, ii: Price competition, kinked demand curves, and edgeworth cycles. *Econometrica*, 56(3):571–599.
- Maskin, E. and Tirole, J. (2001). Markov perfect equilibrium: 1. observable actions. *Journal of Economic Theory*, 100:191–219.
- Nash, J. F. (1950). The bargaining problem. *Econometrica*, 18(2):155–162.
- Ottaviani, M. and Sørensen, P. N. (2006). Reputational cheap talk. *RAND Journal of Economics*, 37(1).
- Patty, J. W. (2009). The politics of biased information. *The Journal of Politics*, 71(2):385–397.
- Prat, A. (2005). The wrong kind of transparency. *American Economic Review*, 95(3):862–877.
- Prendergast, C. and Stole, L. (1996). Impetuous youngsters and jaded old-timers: Acquiring a reputation for learning. *Journal of Political Economy*, 104(6):1105–1134.
- Rappaport, D. (2015). Humility in experts and the revelation of uncertainty. Manuscript.
- Schnakenberg, K. E. (2015). Expert advice to a voting body. *Journal of Economic Theory*, 160:102–113.
- Schnakenberg, K. E. (2017). Informational lobbying and legislative voting. *American Journal of Political Science*, 61(1):129–145.
- Tadelis, S. (2013). *Game Theory: An Introduction*. Princeton University Press, Princeton, NJ.

A Markov Sequential Equilibrium

A.1 General Definition of MSE

Take a general sequential game of incomplete information, and let each node (history) be associated with information set I and an action set A_I . Beliefs μ map information sets into a probability distribution over their constituent nodes. A strategy profile for the entire game σ maps each information set into a probability distribution over A_I . Write the probability (or density) of action a at information set I as $\sigma_I(a)$. Let the function $u_I(a, \sigma)$ denote the von Neumann-Morgenstern expected utility from taking action $a \in A_I$ at an information set I when all subsequent play, by all players, is according to σ . In our setting the payoff-relevant state depends on the information set of the DM, \mathcal{I}_{DM2} , and so in order to define it, we look to affine payoff equivalence following Harsanyi and Selten (1988) and Fudenberg and Tirole (1991).

Definition 2. A strategy σ is a **Markov strategy** if whenever, for any pair of information sets I and I' with associated action sets A_I and $A_{I'}$, and for some constants $\alpha > 0$ and β , there exists a bijection $f : A_I \rightarrow A_{I'}$ such that $u_I(a, \sigma) = \alpha u_{I'}(f(a), \sigma) + \beta, \forall a \in A_I$, then $\sigma_I(a) = \sigma_{I'}(f(a))$.

The extension of equilibrium in Markov strategies to a setting with incomplete information requires some additional language. Our notation and terminology parallels the treatment of sequential equilibrium in Tadelis (2013). As consistency is to sequential equilibrium, so Markov consistency is to Markov sequential equilibrium.

Definition 3. A profile of strategies σ and a system of beliefs μ is **Markov consistent** if there exists a sequence of non-degenerate, Markov mixed strategies $\{\sigma^k\}_{k=1}^{\infty}$ and a sequence of beliefs $\{\mu^k\}_{k=1}^{\infty}$ that are derived from Bayes' Rule, such that $\lim_{k \rightarrow \infty} (\sigma^k, \mu^k) \rightarrow (\sigma, \mu)$.

With this in hand, a notion of Markov sequential equilibrium follows directly.

Definition 4. *A profile of strategies σ , together with a set of beliefs μ , is a **Markov sequential equilibrium** if (σ^*, μ^*) is a Markov consistent perfect Bayesian equilibrium.*

Behavioral Motivation Markov strategies have axiomatic foundations (Harsanyi and Selten, 1988), and can be motivated by purification arguments as well as finite memory in forecasting (Maskin and Tirole, 2001; Bhaskar et al., 2013). In the complete information settings to which it is commonly applied, the Markovian restriction prevents the players from conditioning their behavior on payoff-irrelevant aspects of the (common knowledge) history.¹⁷ The natural extension of this idea to asymmetric information games is to prevent players from conditioning their strategies on payoff-irrelevant private information.

Our restriction on beliefs is also related to the notion of structural consistency proposed by Kreps and Wilson (1982).¹⁸ In that spirit, Markov consistency formalizes a notion of “credible” beliefs, analogous to the notion of credible threats in subgame perfect equilibrium. Instead of using arbitrarily punitive off-path beliefs to discipline on-path behavior, we require that off-path beliefs are credible in the sense that, *ex post*, on arriving at such an off-path node, the relevant agent could construct a convergent sequence of Markov strategies to rationalize them.

Robustness of MSE. Though the restriction to Markov strategies itself enforces a notion of robustness, there is a trivial sense in which the restriction of Markov equilibrium – whether in a complete information setting or an incomplete information setting – is non-robust. In particular, because it imposes symmetric strategies only when incentives are

¹⁷Applications of Markov equilibrium have been similarly focused on the infinitely-repeated, complete information setting. See, e.g. Maskin and Tirole (1988a,b); Ericson and Pakes (1995).

¹⁸Kreps and Ramey (1987) demonstrated that consistency may not imply structural consistency, as conjectured by Kreps and Wilson (1982). We observe that as the Markov property is preserved by limits, Markov consistency does not introduce any further interpretive difficulty.

exactly symmetric, small perturbations of a model may permit much larger sets of equilibria. In the standard applications of the Markov restriction, this could be driven by future payoffs being slightly different depending on the history of play. In our setting, good and bad uninformed experts could have marginally different expected payoffs. Either way, we maintain that this is a red herring. The point of the refinement, like the symmetry condition of Nash (1950), is to hold the theorist to a simple standard: that we be precise about exactly what *kind* of asymmetry in the model construction explains asymmetries in the predicted behavior. From this perspective, another interpretation of our work is that we are reflecting on exactly what informational structures introduce the asymmetry we need to obtain honesty in equilibrium.¹⁹

Note on a dynamic interpretation We have formulated our game as a one-shot sequential game with reputational concerns, so the payoff equivalence holds without the need for affine transformations. In the repeated game implied by the reputational concerns, this will not generally be the case: depending on the formulation of payoffs, good experts will most likely have a higher continuation value than bad ones in all but a babbling equilibrium. This is where the potential for affine transformations in our definition of Markov strategies is useful – if the prior beliefs of the DM at the beginning of a period are sufficient for the history of the game, then setting β equal to the difference in continuation values means that the Markov restriction still binds.

¹⁹We thank an anonymous referee for pointing out that one could also develop a notion of ε -Markov equilibrium to make this point. This is beyond the theoretical ambition of the current work, but an interesting direction for future work.

A.2 MSE and SE

Here we offer a brief discussion of the prior use of the Markov sequential equilibrium (MSE) solution concept as well as an illustration of its implications as a refinement on off-path beliefs.

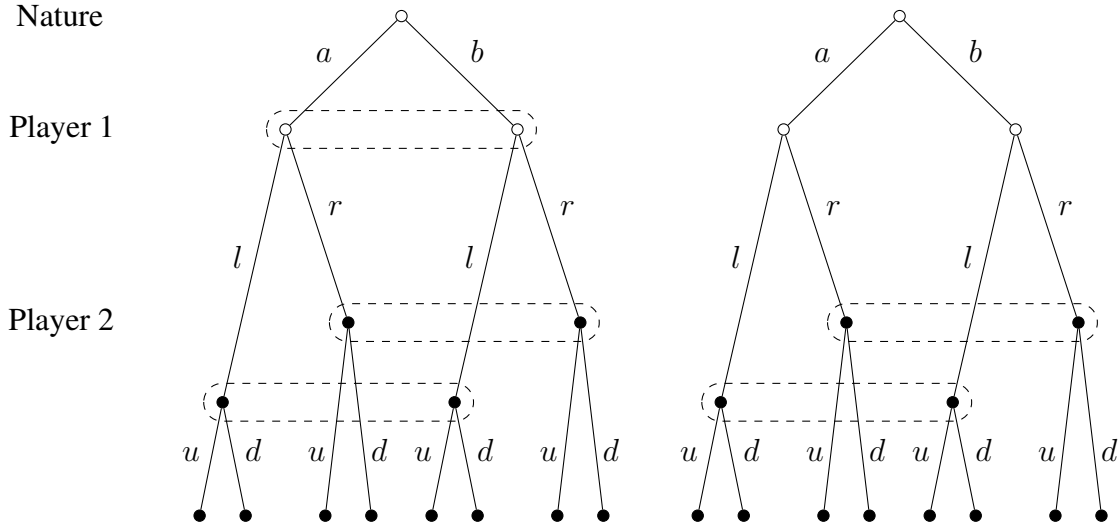
MSE is the natural extension of Markov Perfect Equilibrium to incomplete information games. However, its usage is infrequent and sometimes informal. To our knowledge, there is no general treatment nor general guidance to the construction of the maximally coarse (Markov) partition of the action space, unlike the case of MPE (Maskin and Tirole, 2001). Bergemann and Hege (2005) and Bergemann and Hörner (2010) employ the solution concept, defining it as a perfect Bayesian equilibrium in Markovian strategies. In other words, they impose the Markov restriction only on the sequential rationality condition. This is different and weaker than our construction; our definition of MSE imposes the Markov assumption on both sequential rationality as well as consistency. While they do not use the Markov restriction to refine off-path beliefs, this is of no consequence for their applications.

To see the relevance of MSE to off-path beliefs, consider the game illustrated in Figure A.1, which is constructed to mirror an example from Kreps and Wilson (1982).²⁰ First, nature chooses Player 1's type, a or b . Next, Player 1 chooses l or r . Finally, Player 2 chooses u or d . Player 2 is never informed of Player 1's type. Whether Player 1 knows their own type is the key difference between the two games.

In the first game, the player does not know their type. Posit an equilibrium in which Player 1 always chooses l . What must Player 2 believe at a node following r ? If the theorist is studying perfect Bayesian equilibrium (PBE), they may specify any beliefs they wish. Alternatively, if they are studying sequential equilibrium (SE), Player 2 must believe that

²⁰See, in particular, their Figure 5 (p.873).

Figure A.1: Consistency, Markov Consistency, and Off-Path Beliefs



Notes: This figure depicts two games, which differ in whether Player 1 knows their own type. Their type, a or b , is chosen by Nature with $Pr(a) = p$ and $Pr(b) = 1 - p$. Player 1 chooses l or r , and Player 2 sees this and reacts with u or d . Payoffs are omitted, but can be written $u_i(\cdot, \cdot, \cdot)$.

Player 1 is of type a with probability p .

In the second game depicted, SE imposes no restriction on Player 2's off-path beliefs. However, MSE may. If $u_1(a, l, \cdot) = u_1(b, l, \cdot)$ and $u_1(a, r, \cdot) = u_1(b, r, \cdot)$ (or, more generally, the expected utilities are equal up to an affine transformation) then we say that Player 1's type is *payoff irrelevant*. The restriction to Markov strategies implies that Player 1's strategy does not depend upon their type. Markov consistency implies that, further, Player 2 cannot update about payoff irrelevant information. Therefore Player 2 must believe that Player 1 is of type a with probability p .

A.3 Non-Markovian PBE

Here we briefly discuss PBE that fail the Markov consistency requirement of MSE, and argue why we believe these equilibria are less sensible.

In particular, we demonstrate that the most informative equilibrium under no policy concerns can involve more transmission of uncertainty and also information about the state. However, these equilibria are not robust to minor perturbations, such as introducing a vanishingly small random cost of lying.

Example 1: Admission of Uncertainty with No Validation. Even without the Markov restriction, it is immediate that there can be no fully honest equilibrium with no validation. In such an equilibrium, the competence assessment for sending either m_0 or m_1 is 1, and the competence assessment for sending m_\emptyset is strictly less than one. So the uninformed types have a strict incentive to deviate to m_0 or m_1 .

However, unlike the case with the Markov restriction which leads to babbling, there is an always guessing equilibrium: If informed types always send $m_x = s_x$, and all uninformed types send m_1 with probability p_1 and m_0 otherwise, the competence assessment upon observing either message is p_g . So no type has an incentive to deviate.

Further, it is possible to get admission of uncertainty if the good and bad uninformed types play different strategies. In the extreme, suppose the good types always send their honest message, including the uninformed sending m_\emptyset . If the bad types were to always send m_0 or m_1 , then the competence assessment upon sending m_\emptyset would be 1. In this case, saying “I don’t know” would lead to the highest possible competence evaluation, giving an incentive for all to admit uncertainty even if they know the state.

It is straightforward to check that if the bad types mix over messages (m_0, m_1, m_\emptyset) with probabilities $(p_e(1 - p_1), p_e p_1, 1 - p_e)$, then the competence assessment upon observing all messages is p_g , and so no expert has an incentive to deviate.

A common element of these equilibria is that the competence assessment for any on-path message is equal to the prior. In fact, a messaging strategy can be part of a PBE if and only if this property holds: the competence assessments must be the same to prevent deviation, and if they are the same then by the law of iterated expectations they must equal the prior. So, there is a range of informative equilibria, but they depend on types at payoff-equivalent information sets taking different actions, a violation of Markov strategies that reflects their sensitivity to small perturbations of the payoffs.

Example 2: Honesty with State Validation or Difficulty Validation. Now return to the state validation case, and the conditions for an honest equilibrium. Without the Markov restriction on beliefs, it is possible to set the off-path belief upon observing an incorrect guess to 0. With this off-path belief, the incentive compatibility constraint to prevent sending m_1 becomes $\pi_g^\emptyset \geq p_1$. Since π_g^\emptyset is a function of p_g and p_e (but not p_1), this inequality holds for a range of the parameter space. However, this requires beliefs that are not Markov consistent – the DM who reaches that off-path node cannot construct a Markov strategy to rationalize their beliefs. So we argue that the threat of these beliefs not credible.

Similarly, without the Markov restriction it is possible to get honesty with just difficulty validation. The binding constraint is that if any off-path message leads to a zero competence evaluation, the bad type gets a higher payoff from sending m_\emptyset (as will the case with $\gamma \rightarrow 0$ case, $(1 - p_e)p_g$) than from sending m_1 (now p_e). So, honesty is possible if $(1 - p_e)p_g > p_e$, i.e., the same condition as when $\gamma \rightarrow 0$.

The Fragility of These Examples. A standard defense of Markov strategies in repeated games is that they represent the simplest possible rational strategies (Maskin and Tirole, 2001). A similar principle applies here: rather than allowing for types with the same (effective) information to use different mixed strategies sustained by indifference, MSE focuses on the simpler case where those with the same incentives play the same strategy.

Further, as shown by Bhaskar et al. (2013) for the case of finite social memory, taking limits of vanishing, independent perturbations to the payoffs – in the spirit of Harsanyi and Selten (1988) “purification” – results in Markov strategies as well. Intuitively, suppose the expert receives a small perturbation to his payoff for sending each message which is independent of type and drawn from a continuous distribution, so he has a strict preference for sending one message over the others with probability one. Payoff-indifferent types must use the same mapping between the perturbations and messages, analogous to Markovian strategies. Further, if these perturbations put all messages on path, then all beliefs are generated by Markovian strategies.²¹

Summary It would be possible to construct additional informative equilibria if we allowed different types to play different actions, even when they are payoff equivalent. We view this as a modeling contrivance, and this is precisely what the Markov consistency restriction, above and beyond standard consistency, restricts. This point was previously made by Harsanyi and Selten (1988), who contend that the property of “invariance with respect to isomorphisms,” on which our definition of Markov strategies is based, is “an indispensable requirement for any rational theory of equilibrium point selection that is based on strategic considerations exclusively.” Or, in the appeal of Maskin and Tirole (2001) to payoff perturbations, “minor causes should have minor effects.”

²¹A related refinement more specific to our setting is to allow for a small random “lying cost” for sending a message not corresponding to the signal, which is independent of the type (Kartik, 2009).

B Proofs of Results in the Main Text

More general definitions Some of our results in this section rely on more general definitions of messaging strategies. Starting with “honesty”:

Definition 5. Let $\pi_s(m)$ be the DM posterior belief that the expert observed signal s upon sending message m . An equilibrium is **honest** if $\pi_s(m) \in \{0, 1\} \forall s \in \mathcal{S}$ and all on-path m .

As in all cheap-talk games, the messages sent only convey meaning by which types send them in equilibrium. We define admitting uncertainty as sending a message which is never sent by either informed type:

Definition 6. Let M_0 be the set of messages sent by the s_0 types and M_1 be the set of message sent by the s_1 types. Then an expert **admits uncertainty** if he sends a message $m \notin M_0 \cup M_1$

Finally, an important class of equilibria will be one in which the informed types send distinct message from each other, but the uninformed types sometimes if not always mimic these messages:

Definition 7. A **guessing equilibrium** is one where $M_0 \cap M_1 = \emptyset$, and $Pr(m \in M_0 \cup M_1 | \theta, s_\theta) > 0$ for at least one $\theta \in \{g, b\}$. In an **always guessing equilibrium**, $Pr(m \in M_0 \cup M_1 | \theta, s_\theta) = 1$ for both $\theta \in \{g, b\}$.

That is, an always guessing equilibrium is one where the informed types report their signal honestly, but the uninformed types never admit uncertainty.

Proof of Proposition 1: For convenience, we extend the definition of v so $v(a, \pi_1)$ represents the expected quality of policy a under the belief that the state is 1 with probability π_1 .

The DM's expected payoff from the game can be written as the sum over the (expected) payoff as a function of the expert signal:

$$\sum_{s \in \{s_0, s_1, s_\emptyset\}} Pr(s) \sum_m Pr(m|s) v(a^*(m), Pr(\omega|s)). \quad (7)$$

In the honest equilibrium, when the expert observes s_0 or s_1 , the DM takes an action equal to the state with probability 1, giving payoff 1. When the expert observes s_\emptyset , the equilibrium action is p_1 giving payoff $v(p_1, p_1) = 1 - p_1(1 - p_1)$. So, the average payoff is:

$$p_g p_e 1 + (1 - p_g p_e) p_1 (1 - p_1) = \bar{v}.$$

This payoff as expressed in (7) is additively separable in the signals, and v is and strictly concave in a for each s . So, for each $s \in \{s_0, s_1, s_\emptyset\}$, this component of the sum is maximized if and only if $a^*(m)$ is equal to the action taken upon observing the honest message is with probability 1. That is, it must be the case that:

$$a^*(m) = \begin{cases} 1 & m : Pr(m|s_1) > 0 \\ p_1 & m : Pr(m|s_\emptyset) > 0 \\ 0 & m : Pr(m|s_0) > 0 \end{cases} \quad (8)$$

If the equilibrium is not honest, then there must exist a message m' such that $Pr(s|m') < 1$ for all s . At least one of the informed types must send m' with positive probability; if not, $Pr(s_\emptyset|m') = 1$. Suppose the type observing s_0 sends m' with positive probability.

(An identical argument works if it is s_1 .) To prevent $Pr(s_0|m') = 1$ another type must send this message as well, and so in response the DM chooses an action strictly greater than 0, contradicting condition (8) and hence the expected quality of the decision in any equilibrium which is not honest is strictly less than \bar{v} . \square

Proof of Proposition 2: For any messaging strategy, the DM must form a belief about the expert competence for any message (on- or off-path); since it only depends on this message (and not the validation as in other cases) write these $\pi_g(m)$. So, for any type θ , the expected utility for sending message m is just $\pi_g(m)$. All types are payoff-equivalent in any equilibrium, and therefore in any MSE they must use the same strategy. Since all messages are sent by both informed and uninformed types, there is no admission of uncertainty. \square

Proof of Proposition 3: A more complete description of the MSE with no policy concerns and state validation is:

Proposition 7. *With state validation and no policy concerns:*

- i. In any MSE, there is no admission of uncertainty, and*
- ii. any non-babbling MSE is equivalent, subject to relabeling, to an MSE where both uninformed types send m_1 with probability*

$$\sigma_\theta^*(m_1) = \begin{cases} \frac{p_1(1+p_g p_e) - p_g p_e}{1 - p_g p_e} & \text{if } p_1 < 1/(1 + p_g p_e) \\ 1 & \text{otherwise} \end{cases}$$

—and m_0 with probability $\sigma_\theta^(m_0) = 1 - \sigma_\theta^*(m_1)$.*

Proof. Part i is immediate in a babbling equilibrium: there is no admission of uncertainty since there are no messages only sent by the uninformed types. Propositions 14 and 15 in

Appendix D, shows that all MSE are equivalent, subject to a relabeling of the messages, to one where the the s_0 types send m_0 , the s_1 types send m_1 , and there is only one other potential message m_\emptyset . What remains to be shown is that in the MSE of this form, the uninformed types never send m_\emptyset and send m_0 and m_1 with the probabilities in the statement of the proposition.

Recall the Markov strategy restriction implies the good and bad uninformed types use the same strategy. As shown in the main text, in a conjectured honest equilibrium, the payoff for an uninformed type to send m_\emptyset is π_g^\emptyset . To formally show a deviation to m_1 is profitable, recall the payoff to sending this message when $\omega = 1$ is 1. When $\omega = 0$, the competence assessment is off path. Markov consistency requires that this belief be formed as the limit to a sequence of well-defined beliefs which are consistent with a corresponding sequence of Markov strategies. Take any sequence of Markov strategies σ^k and resulting beliefs:

$$\pi_g^k(m_1, 0) = \frac{Pr(m_1, 0, \theta = g)}{Pr(m_1, 0)} = \frac{(1 - p_1)p_g p_e \sigma_0^k(m_1) + (1 - p_1)p_g(1 - p_e)\sigma_\emptyset^k(m_1)}{(1 - p_1)p_g p_e \sigma_0^k(m_1) + (1 - p_1)(1 - p_g p_e)\sigma_\emptyset^k(m_1)}.$$

This belief is increasing in $\sigma_0^k(m_1)$ and decreasing in $\sigma_\emptyset^k(m_1)$, and can range from π_g^\emptyset (when $\sigma_0^k(m_1) = 0$ and $\sigma_\emptyset^k(m_1) > 0$) to 1 (when $\sigma_0^k(m_1) > 0$ and $\sigma_\emptyset^k(m_1) = 0$). So, $\pi_g(m_1, 0)$ must be the limit of a sequence which lies in $[\pi_g^\emptyset, 1]$, and the off path belief must lie on this interval as well. Hence the payoff to deviating to m_1 is at least $p_1 + (1 - p_1)\pi_g^\emptyset > \pi_g^\emptyset$, completing the proof that there is no honest equilibrium.

Now suppose the uninformed types send m_\emptyset with a probability strictly between 0 and 1. The competence assessment for sending m_\emptyset is π_g^\emptyset . Writing the probability the uninformed types send m_1 with $\sigma_\emptyset(m_1)$, the competence assessment for sending m_1 and observing

validation that $w = 1$ is:

$$\begin{aligned}
Pr(\theta = g|m_1; \sigma_\emptyset(m_1)) &= \frac{p_g p_e p_1 + p_g(1 - p_e)\sigma_\emptyset(m_1)}{p_g p_e p_1 + (p_g(1 - p_e) + (1 - p_g))\sigma_\emptyset(m_1)} \\
&\geq \frac{p_g p_e p_1 + p_g(1 - p_e)}{p_g p_e p_1 + (p_g(1 - p_e) + (1 - p_g))} \\
&> \frac{p_g(1 - p_e)}{p_g(1 - p_e) + (1 - p_g)} = \pi_g^\emptyset.
\end{aligned}$$

Since the competence assessment for sending m_1 is strictly higher than for sending m_0 , there can be no MSE where the uninformed types admit uncertainty, completing part i.

For part ii, first consider the condition for an equilibrium where both m_0 and m_1 are sent by the uninformed types. The uninformed types must be indifferent between guessing m_0 and m_1 . This requires:

$$p_1 \pi_g(m_1, \omega = 1) + (1 - p_1) \pi_g^\emptyset = (1 - p_1) \pi_g(m_0, \omega = 0) + p_1 \pi_g^\emptyset \quad (9)$$

where the posterior beliefs about competence when “guessing wrong” are π_g^\emptyset and when “guessing right” are given by Bayes’ rule:

$$\begin{aligned}
\pi_g(m_1, \omega = 1) &= \frac{Pr(\theta = g, \omega = 1, m_1)}{Pr(m_1, \omega = 1)} = \frac{p_1 p_g (p_e + (1 - p_e) \sigma_\emptyset(m_1))}{p_1 (p_g p_e + (1 - p_g p_e) \sigma_\emptyset(m_1))} \\
\pi_g(m_0, \omega = 0) &= \frac{Pr(\theta = g, \omega = 0, m_0)}{Pr(m_0, \omega = 0)} = \frac{(1 - p_1) p_g (p_e + (1 - p_e) \sigma_\emptyset(m_0))}{(1 - p_1) (p_g p_e + (1 - p_g p_e) \sigma_\emptyset(m_0))}.
\end{aligned}$$

Plugging these into (9) and solving for the strategies with the additional constraint that $\sigma_\emptyset(m_0) + \sigma_\emptyset(m_1) = 1$ gives:

$$\begin{aligned}
\sigma_\emptyset(m_0) &= \frac{1 - p_1(1 + p_g p_e)}{1 - p_g p_e} \\
\sigma_\emptyset(m_1) &= \frac{p_1(1 + p_g p_e) - p_g p_e}{1 - p_g p_e}.
\end{aligned}$$

For this to be a valid mixed strategy, it must be the case that both of these expressions are between zero and one, which is true if and only if $p_1 < 1/(1 + p_g p_e) \in (1/2, 1)$. So, if this inequality holds and the off-path beliefs upon observing m_0 are sufficiently low, there is an MSE where both messages are sent by the uninformed types. And the competence assessment for any off-path message/validation can be set to π_g^\emptyset (i.e., the lowest belief possible with Markov consistency), which is less than the expected competence payoff for sending either m_0 or m_1 .

Now consider an equilibrium where uninformed types always send m_1 . The on-path message/validation combinations are then $(m_1, \omega = 0)$, $(m_1, \omega = 1)$, and $(m_0, \omega = 0)$, with the following beliefs about the expert competence:

$$\begin{aligned}\pi_g(m_1, \omega = 0) &= \frac{p_g(1 - p_e)}{p_g(1 - p_e) + 1 - p_g}; \\ \pi_g(m_1, \omega = 1) &= \frac{p_g p_e + p_g(1 - p_e)}{p_g p_e + p_g(1 - p_e) + (1 - p_g)} = p_g, \quad \text{and} \\ \pi_g(m_0, \omega = 0) &= 1.\end{aligned}$$

Preventing the uninformed types from sending m_0 requires:

$$p_1 p_g + (1 - p_1) \frac{p_g(1 - p_e)}{p_g(1 - p_e) + 1 - p_g} \geq p_1 \pi_g(m_0, \omega = 1) + (1 - p_1).$$

This inequality is easiest to maintain when $\pi_g(m_0, \omega = 1)$ is small, and by the argument in the main text, in an MSE it must be at least π_g^\emptyset . Setting $\pi_g(m_0, \omega = 1) = \pi_g^\emptyset$ and simplifying gives $p_1 \geq 1/(1 + p_g p_e)$, i.e., the reverse of the inequality required for an MSE where both m_0 and m_1 are sent. Again, setting the competence assessment for an off-path message to π_g^\emptyset prevents this deviation.

So, if $p_1 \leq 1/(1 + p_g p_e)$ there is an MSE where both messages are sent, and if not there is

an MSE where only m_1 is sent.

Finally, it is easy to verify there is never an MSE where only m_0 is sent, as the uninformed types have an incentive to switch to m_1 . \square

Proof of Proposition 4 See the proof of proposition 12 in Appendix C.

Proof of Proposition 5: For part i, the result is immediate in the range of p_e where p_g does not change the bad type strategy. For the range where the bad type strategy is a function of p_g , plugging in the strategies identified in (6) and simplifying gives the expected quality of the decision is:

$$1 - p_1(1 - p_1) + \frac{(p_e p_g)^2 p_1(1 - p_1)}{p_e - p_g(1 - 2p_e)}. \quad (10)$$

The derivative of (10) with respect to p_g is:

$$\frac{p_1(1 - p_1)p_e^2 p_g(2p_e(1 + p_g) - p_g)}{(p_e - p_g(1 - 2p_e))^2}.$$

which is strictly positive if $p_e > \frac{p_g}{2(1+p_g)}$. Since the range of p_e where the bad type plays a mixed strategy is $p_e \in (p_g/(1 + p_g), 1/(1 + p_g))$, this always holds.

For part ii, the unconditional probability of admitting uncertainty in this equilibrium is:

$$p_g(1 - p_e) + (1 - p_g)\sigma_b(m_\emptyset). \quad (11)$$

Within the range for an honest equilibrium (where $\sigma_b(m_\emptyset) = 1$) the derivative of (11) with respect to p_g is $-p_e < 0$. In the range where the bad type always guesses ($\sigma_b(m_\emptyset) = 0$) the derivative is $(1 - p_e) > 0$. Plugging in the equilibrium strategy when interior and

differentiating with respect to p_g gives $1 - 2p_e$, which is positive if and only if $p_e < 1/2$.

□

Proof of Proposition 6: For part i, $\sigma_b(m_\emptyset)$ is weakly decreasing in p_e , so (11) is strictly decreasing in p_e .

For part ii, in the range $p_e \in (p_g/(1+p_g), 1/(1+p_g))$, the expected quality of the decision is (10). Differentiating this with respect to p_e gives:

$$\frac{p_1(1-p_1)p_e p_g^2 (p_e - 2p_g + 2p_e p_g)}{(p_e - p_g + 2p_e p_g)^2}$$

which, evaluated at $p_e = p_g/(1+p_g)$ simplifies to $-p_1(1-p_1)$. By continuity, this derivative must be negative on some nonempty interval $(p_g/(1+p_g), \tilde{p}_e)$. So, the value of the decision must be locally decreasing at $p_e = p_g/(1+p_g)$, and by continuity, for an open interval $p_e \in (p_g/(1+p_g), \tilde{p}_\delta)$. □

C Analysis of Other Cases

Section 5 (and associated proofs) contains a full analysis of the cases with no policy concerns and no validation/state validation, and Section 6 analyzes the case with small policy concerns and difficulty validation.

Here we first tie up the cases of difficulty and full validation with no policy concerns, and then the other validation cases with policy concerns.

Difficulty and Full Validation with No Policy Concerns

Difficulty validation A complete description of the MSE with difficulty validation proves challenging. However, with no policy concerns, we can show that there is never any information communicated about the state, though there can be information communicated about the difficulty of the problem:

Proposition 8. *With no policy concerns and difficulty validation,*

- i. in any MSE, $a^*(m) = p_1$ for all on-path m , and*
- ii. there is an MSE where the good uninformed types always admit uncertainty.*

Proof. Given the payoff equivalence classes, the good and informed types must use the same mixed strategy. In any MSE, the posterior belief about the state upon observing an on-path message m can be written as a weighted average of the belief about the state conditional on being in each equivalence class, weighted by the probability of being in the class:

$$\begin{aligned}
 Pr(\omega = 1|m) &= Pr(\omega = 1|m, \theta = g, s \in \{s_0, s_1\})Pr(\theta = g, s \in \{s_0, s_1\}|m) \\
 &\quad + Pr(\omega = 1|m, \theta = g, s = s_\emptyset)Pr(\theta = g, s = s_\emptyset|m) \\
 &\quad + Pr(\omega = 1|m, \theta = b)Pr(\theta = b|m) \\
 &= p_1Pr(\theta = g, s \in \{s_0, s_1\}|m) + p_1Pr(\theta = g, s = s_\emptyset|m) + p_1Pr(\theta = b|m) = p_1.
 \end{aligned}$$

For each equivalence class there is no information conveyed about the state, so these conditional probabilities are all p_1 , and hence sum to this as well.

For part ii, we construct an equilibrium where the informed types always send m_e (“the problem is easy”), the good but uninformed types send m_h (“the problem is hard”), and the bad types mix over these two messages with probability $(\sigma_b(m_e), \sigma_b(m_h))$. Since m_h

is never sent by the informed types, sending this message admits uncertainty.

There can be an equilibrium where both of these messages are sent by the bad types if and only if they give the same expected payoff. Writing the probability of sending m_e as $\sigma_b(m_e)$, this is possible if:

$$p_e \pi_g(m_e, e) + (1 - p_e) \pi_g(m_e, h) = p_e \pi_g(m_h, e) + (1 - p_e) \pi_g(m_h, h),$$

– or, rearranged:

$$p_e \frac{p_g p_e}{p_g p_e + (1 - p_g) \sigma_b(m_e)} = (1 - p_e) \frac{p_g (1 - p_e)}{p_g (1 - p_e) + (1 - p_g) (1 - \sigma_b(m_e))}. \quad (12)$$

The left-hand side of this equation (i.e., the payoff to guessing the problem is easy) is decreasing in $\sigma_b(m_e)$, ranging from p_e to $p_e \frac{p_g p_e}{p_g p_e + (1 - p_g)}$. The right-hand side is increasing in $\sigma_b(m_e)$, ranging from $(1 - p_e) \frac{p_g (1 - p_e)}{p_g (1 - p_e) + (1 - p_g)}$ to $1 - p_e$. So, if

$$p_e \frac{p_g p_e}{p_g p_e + (1 - p_g)} - (1 - p_e) \geq 0, \quad (13)$$

then payoff to sending m_e is always higher. After multiplying through by $p_g p_e + (1 - p_g)$, the left-hand side of (13) is quadratic in p_e (with a positive p_e term), and has a root at $\frac{2p_g - 1 + \sqrt{1 + 4p_g - 4p_g^2}}{4p_g}$ which is always on $(1/2, 1)$, and a negative root.²² So, when p_e is above this root, the payoff to sending m_e is always higher, and hence there is a MSE where the uninformed types always send this message.

On the other hand, if

$$(1 - p_e) \frac{p_g (1 - p_e)}{p_g (1 - p_e) + (1 - p_g)} - p_e \geq 0,$$

²²All of these observations follow from the fact that $1 + 4p_g - 4p_g^2 \in (1, (2p_g + 1)^2)$.

then the payoff for sending m_h is always higher, which by a similar argument holds if $p_e \leq \frac{2p_g+1-\sqrt{1+4p_g-4p_g^2}}{4p_g}$. However, if neither of these inequalities hold, then there is a $\sigma_b(m_e) \in (0, 1)$ which solves (12), and hence there is an MSE where m_e is sent with this probability and m_h with complementary probability. Summarizing, there is an MSE where the bad type sends message m_e with probability:

$$\sigma_b^*(m_e) = \begin{cases} 0 & p_e \leq \frac{2p_g+1-\sqrt{1+4p_g-4p_g^2}}{4p_g} \\ \frac{p_e(p_e-p_g+2p_e p_g-2p_e^2 p_g)}{(1-p_g)(1-2p_e(1-p_e))} & p_e \in \left(\frac{2p_g+1-\sqrt{1+4p_g-4p_g^2}}{4p_g}, \frac{2p_g-1+\sqrt{1+4p_g-4p_g^2}}{4p_g} \right) \\ 1 & p_e \geq \frac{2p_g-1+\sqrt{1+4p_g-4p_g^2}}{4p_g} \end{cases}$$

and message m_h with probability $\sigma_b^*(m_h) = 1 - \sigma_b^*(m_e)$. \square

This implies that, while we can learn something about the question difficulty with difficulty validation alone, learning about the state (and attaining an honest equilibrium) require either nonzero policy concerns or state validation (or both).

Full Validation with No Policy Concerns With full validation, there are four possible validation results for each message. The expected payoff to sending message m given one's type and message is:

$$\sum_{\delta \in \{e, h\}} \sum_{\omega \in \{0, 1\}} Pr(\delta|s, \theta) Pr(\omega|s, \theta) \pi_g(m, \omega, \delta).$$

No pair of types share the same $Pr(\omega|s, \theta)$ and $Pr(\delta|s, \theta)$, so none must be payoff equivalent. As a result, all types can use distinct strategies, and off-path beliefs are unrestricted.

In an honest equilibrium, upon observing $(m_0, 0, e)$ or $(m_1, 1, e)$, the DM knows that the expert is competent. Upon observing (m_\emptyset, ω, e) the DM knows that the expert is not compe-

tent, as a competent expert would have received and sent an informative message since the problem is easy. Upon observing (m_\emptyset, ω, h) , the DM belief about the expert competence is the same as the prior, since if the problem is hard no one gets an informative message (and all send m_\emptyset).²³ So, the competence evaluations for the on-path messages are:

$$\begin{aligned}\pi_g(m_0, 0, e) &= 1 & \pi_g(m_1, 1, e) &= 1 \\ \pi_g(m_\emptyset, \omega, e) &= 0 & \pi_g(m_\emptyset, \omega, h) &= p_g\end{aligned}$$

To make honesty as easy as possible to sustain, suppose that for any off-path message (“guessing wrong”), the competence evaluation is zero.

The informed types get a competence evaluation of 1 for sending their honest message, so face no incentive to deviate.

A good but uninformed type knows the difficulty validation will reveal $\delta = h$, but does not know ω . Sending the honest message m_\emptyset gives a competence payoff of p_g . However, sending either m_0 or m_1 will lead to an off-path message/validation combination, and hence a payoff of zero. So, these types face no incentive to deviate.

Finally, consider the bad uninformed types, who do not know what either the state or difficulty validation will reveal. If they send m_\emptyset , they will be caught as uninformed if the problem was in fact easy (probability p_e). However, if the problem is hard, the DM does not update about their competence for either state validation result. So, the expected payoff to sending m_\emptyset is $(1 - p_e)p_g$.

If guessing m_1 , the expert will be “caught” if either the problem is hard *or* the state is 0. However, if guessing correctly, the competence evaluation will be 1. So, the expected

²³Formally, applying Bayes’ rule gives $Pr(\theta = g|m_\emptyset, \omega, h) = \frac{p_g(1-p_e)}{p_g(1-p_e)+(1-p_g)(1-p_e)} = p_g$.

payoff to this deviation is $p_e p_1$. Similarly, the expected payoff to guessing m_0 is $p_e(1 - p_1) < p_e p_1$, so m_1 is the best deviation.

Honesty is possible if admitting uncertainty leads to a higher competence evaluation than guessing m_1 , or:

$$(1 - p_e)p_g \geq p_e p_1 \implies p_e \leq \frac{p_g}{p_g + p_1}.$$

If this inequality does not hold, a fully honest MSE is not possible. However, there is always an MSE where the good but uninformed types always send m_\emptyset . In such an equilibrium, the bad types pick a mixed strategy over m_0 , m_1 , and m_\emptyset . Whenever the DM observes an “incorrect guess” she assigns a competence evaluation of zero. So, the informed types never deviate (as this ensures an incorrect guess), and good uninformed types have no reason to guess since they know the problem is hard. Returning to the derivation of the honest equilibrium, the off-path beliefs in this MSE are justified, in the sense that the good types all have strict incentives to report their honest message, and the bad types are the only ones who potentially face an incentive to send m_0 or m_1 when the problem is hard or m_\emptyset when the problem is easy.

Summarizing:

Proposition 9. *With no policy concerns and full validation, there is an MSE where the informed types send distinct messages and the good but uninformed types always admit uncertainty. If $p_e \leq \frac{p_g}{p_g + p_1}$, there is an honest MSE.*

Proof. The condition for the honest equilibrium is derived above. So what remains is to show there is always an MSE where the good but uninformed type always sends m_\emptyset .

In such an equilibrium, message/validation combinations $(m_0, 0, e)$, $(m_1, 1, e)$ and $(m_\emptyset, 0, h)$ and $(m_\emptyset, 1, h)$ are the only ones observed when the expert is competent. So, any other message/validation combination is either on-path and only sent by the bad types, in which case the competence assessment must be 0, or is off-path and can be set to 0.

The informed type observing s_0 knows the validation will be $(0, e)$, and $(m, 0, e)$ leads to competence assessment zero for $m \neq m_0$. So, this type has no incentive to deviate, nor does the s_1 type by an analogous argument. The good but uninformed type knows the validation will reveal h , and the DM observing (m_i, ω, h) for $i \in \{0, 1\}$ and $\omega \in \{0, 1\}$ will lead to a competence assessment of zero. So this type faces no incentive to deviate.

Now consider the bad type strategy. While explicitly deriving the equilibrium strategies here is tedious, a simple fixed point argument can be used to show existence. Write the whole strategy with $\sigma_b = (\sigma_b(m_0), \sigma_b(m_1), \sigma_b(m_\emptyset))$, and the bad type's expected competence assessment for sending each message when the DM expects strategy σ (averaging over the validation result) as:

$$\begin{aligned} U_\theta(m_\emptyset, b, \sigma) &\equiv p_e 0 + (1 - p_e) \frac{p_g}{p_g + (1 - p_g)\sigma_b(m_\emptyset)}, \\ U_\theta(m_0, b, \sigma) &\equiv p_e(1 - p_1) \frac{p_g}{p_g + (1 - p_g)\sigma_b(m_0)} + (1 - p_e)0, \text{ and} \\ U_\theta(m_1, b, \sigma) &\equiv p_e p_1 \frac{p_g}{p_g + (1 - p_g)\sigma_b(m_1)} + (1 - p_e)0. \end{aligned}$$

Write the expected payoff to the bad expert choosing mixed strategy σ when the decision-maker expects mixed strategy $\hat{\sigma}_b$ as $U(\sigma, \hat{\sigma}) = \sum_{i \in \{0, 1, \emptyset\}} \sigma_b(m_i) U_\theta(m_i; \hat{\sigma})$, which is continuous in all $\sigma_b(m_i)$, so optimizing this objective function over the (compact) unit simplex must have a solution. So, $BR(\hat{\sigma}_b) = \arg \max_\sigma U(\sigma; \hat{\sigma})$ is a continuous mapping from the unit simplex to itself, which by the Kakutani fixed point theorem must have a solution. So, the strategy (or strategies) given by such a fixed point are a best response for the bad type

when the decision-maker forms correct beliefs given this strategy.

□

Nonzero Policy Concerns

Now we the case where the expert cares about the policy made, $\gamma > 0$. Not surprisingly, when policy concerns are “large”, there is always an honest MSE since the expert primarily wants the DM to take the best possible action. Here we analyze how high policy concerns have to be in order to attain this honest equilibrium, and provide some results about what happens when policy concerns are not small but not large enough to induce honesty.

In Appendix D, we show that with no validation and state validation, in any MSE which is not babbling, the types observing s_0 and s_1 cannot send any common messages. Combined with a relabeling argument, for all of the analysis (of these validation regimes) with policy concerns we can again restrict attention to MSE where the informed types always send m_0 and m_1 , respectively, and uninformed types send at most one other message m_\emptyset .

While we do not rely on this result for the difficulty and full validation cases (where we only focus on the existence of equilibria with certain properties), we analyze the analogous messaging strategies in these cases to facilitate comparison.

No Validation Informed types never face an incentive to deviate from the honest equilibrium: upon observing s_x for $x \in \{0, 1\}$, the DM chooses policy $a^*(s_x) = x$, and knows the expert is competent, giving the highest possible expert payoff.

Uninformed types, however, may wish to deviate. Upon observing m_\emptyset , the DM takes action

$a = \pi_1 = p_1$, which gives expected policy value $1 - p_1(1 - p_1)$, and the belief about the competence is π_g^\emptyset . So, for the uninformed experts of either competence type, the payoff for reporting honestly and sending signal m_\emptyset is:

$$\pi_g^\emptyset + \gamma(1 - p_1(1 - p_1)). \quad (14)$$

If the expert deviates to $m \in \{m_0, m_1\}$, his payoff changes in two ways: he looks competent with probability 1 (as only competent analysts send these messages in an honest equilibrium, and without validation this is always on path), and the policy payoff gets worse on average. So, the payoff to choosing m_1 is:

$$1 + \gamma p_1. \quad (15)$$

It is easy to check that the payoff to deviating to m_0 is weakly lower, and so m_1 is the binding deviation to check. Preventing the uninformed type from guessing m_1 requires

$$\pi_g^\emptyset + \gamma(1 - p_1(1 - p_1)) \geq 1 + \gamma p_1.$$

Rearranging, define the threshold degree of policy concerns γ_{NV}^H required to sustain honesty by

$$\begin{aligned} \gamma &\geq \frac{1 - \pi_g^\emptyset}{(1 - p_1)^2} \\ &= \frac{(1 - p_g)}{(1 - p_g p_e)(1 - p_1)^2} \\ &\equiv \gamma_{NV}^H. \end{aligned} \quad (16)$$

If $\gamma < \gamma_{NV}^H$, the uninformed types strictly prefer sending m_1 to m_\emptyset *if the DM expects honesty*. Given our concern with admission of uncertainty, it is possible that there is a

mixed strategy equilibrium where the uninformed types sometimes send m_\emptyset and sometimes send m_0 or m_1 . However, as the following result shows, when policy concerns are too small to induce full honesty, the payoff for sending m_1 is always higher than the payoff for admitting uncertainty. Moreover, since γ_{NV}^H is strictly greater than zero, when policy concerns are sufficiently small some form of validation is required to elicit any admission of uncertainty.

Proposition 10. *When $\gamma > 0$ and no validation:*

i. If $\gamma \geq \gamma_{NV}^H$, then there is an honest MSE,

ii. If $\gamma \in (0, \gamma_{NV}^H)$, then all non-babbling MSE are always guessing (i.e., $\sigma_\emptyset^(m_\emptyset) = 0$)*

Proof. Part i is shown above.

For part ii, it is sufficient to show that if $\gamma < \gamma_{NV}^H$, then in any proposed equilibrium where $\sigma_\emptyset(m_\emptyset) > 0$, the payoff for an expert to send m_1 is always strictly higher than the payoff to sending m_\emptyset .

The competence evaluation upon observing m_1 as a function of the uninformed expert mixed strategy is:

$$\pi_g(m_1; \sigma_\emptyset(m_1)) = \frac{Pr(\theta = g, m_1)}{Pr(m_1)} = \frac{p_g p_1 p_e + p_g(1 - p_e)\sigma_\emptyset(m_1)}{p_g p_1 p_e + (p_g(1 - p_e) + (1 - p_g))\sigma_\emptyset(m_1)}$$

– and the belief about the state is:

$$\pi_1(m_1; \sigma_\emptyset(m_1)) = \frac{Pr(\omega = 1, m_1)}{Pr(m_1)} = \frac{p_1(p_g p_e + (1 - p_g p_e)\sigma_\emptyset(m_1))}{p_1 p_g p_e + (p_g(1 - p_e) + (1 - p_g))\sigma_\emptyset(m_1)}.$$

When observing m_\emptyset , the DM knows with certainty that the expert is uninformed, so $\pi_g(m_\emptyset) = \pi_g^\emptyset$ and $\pi_1(m_\emptyset) = p_1$.

Combining, the expected payoff for an uninformed type to send each message is:

$$U(m_1; s_\emptyset, \sigma_\emptyset(m_1)) = \pi_g(m_1; \sigma_\emptyset(m_1)) \\ + \gamma(1 - [p_1(1 - \pi_1(m_1; \sigma_\emptyset(m_1)))^2 + (1 - p_1)\pi_1(m_1; \sigma_\emptyset(m_1))^2])$$

and

$$U(m_\emptyset) = \pi_g^\emptyset + \gamma(1 - p_1(1 - p_1)).$$

Conveniently, $U(m_\emptyset)$ is not a function of the mixed strategy.

If $\gamma = 0$, then $U(m_i; \sigma_i) > U(m_\emptyset)$ for both $i \in \{0, 1\}$, because $\pi_g(m_i; \sigma_i) > \pi_g^\emptyset$. Further, by the continuity of the utility functions in γ and $\sigma_\emptyset(m_1)$, there exists a $\gamma^* > 0$ such that message m_1 will give a strictly higher payoff than m_\emptyset for an open interval $(0, \gamma^*)$. The final step of the proof is to show that this γ^* is exactly γ_{NV}^H .

To show this, let $\sigma^{\text{cand}}(\gamma)$ be the candidate value of $\sigma_\emptyset(m_1)$ that solves $U(m_1; s_\emptyset, \sigma_\emptyset(m_1)) = U(m_\emptyset)$. Rearranging, and simplifying this equality gives:

$$\sigma^{\text{cand}}(\gamma) = -\frac{p_1 p_g p_e}{1 - p_g p_e} + \gamma \frac{p_1 p_g p_e (1 - p_1)^2}{1 - p_g}$$

which is linear in γ . When $\gamma = 0$, $\sigma^{\text{cand}}(\gamma)$ is negative, which re-demonstrates that with no policy concerns the payoff to sending m_1 is always higher than m_\emptyset . More generally, whenever $\sigma^{\text{cand}}(\gamma) < 0$, the payoff to sending m_1 is always higher than m_\emptyset so there can be

no admission of uncertainty. Rearranging this inequality gives:

$$-\frac{p_1 p_g p_e}{1 - p_g p_e} + \gamma \frac{p_1 p_g p_e (1 - p_1)^2}{1 - p_g} < 0$$

$$\Leftrightarrow \gamma < \frac{1 - p_g}{(1 - p_g p_e)(1 - p_1)^2} = \gamma_{NV}^H,$$

completing part ii.

Now that we have demonstrated any MSE is always guessing, we can prove proposition 10. As $\gamma \rightarrow 0$, the condition for an equilibrium where the uninformed types send both m_0 and m_1 is that the competence assessments are the same. Writing these out gives:

$$\pi_g(m_0; \sigma_\emptyset) = \pi_g(m_1; \sigma_\emptyset)$$

$$\frac{p_g(1 - p_1)p_e + p_g(1 - p_e)\sigma_\emptyset(m_0)}{p_g(1 - p_1)p_e + (p_g(1 - p_e) + (1 - p_g))\sigma_\emptyset(m_0)} = \frac{p_g p_1 p_e + p_g(1 - p_e)\sigma_\emptyset(m_1)}{p_g p_1 p_e + (p_g(1 - p_e) + (1 - p_g))\sigma_\emptyset(m_1)}$$

which, combined with the fact that $\sigma_\emptyset(m_1) = 1 - \sigma_\emptyset(m_0)$ (by part ii) is true if and only if $\sigma_\emptyset(m_0) = 1 - p_1$ and $\sigma_\emptyset(m_1) = p_1$. There is no equilibrium where $\sigma_\emptyset(m_0) = 0$; if so, $\pi_g(m_0; \sigma_\emptyset) = 1 > \pi_g(m_1; \sigma_\emptyset)$. Similarly, there is no equilibrium where $\sigma_\emptyset(m_1) = 0$. \square

An intuition for this result is as follows. As uninformed types send m_1 more often, this has two affects on the appeal of sending this message. First, there is a complementarity where sending m_1 more often makes the policy response to this message less extreme, which the uninformed types like. Second, there is a substitution effect where it makes those sending m_1 look less competent. While these effects go in the opposite direction, the substitution effect that makes sending m_1 less appealing when other uninformed types do so is only strong when policy concerns are weak, which is precisely when sending m_1 is generally preferable to m_0 regardless of the uninformed type strategy.

State Validation. Suppose there is an honest equilibrium with state validation.

As in the case with no policy concerns, upon observing message $(m_0, 0)$ or $(m_1, 1)$ the DM knows the expert is competent and takes an action equal to the message, and upon $(m_0, 0)$ or $(m_0, 1)$ takes action p_1 and knows the expert is uninformed, giving competence evaluation π_g^\emptyset . So, the payoff for an uninformed type to send the equilibrium message is:

$$\pi_g^\emptyset + \gamma(1 - p_1(1 - p_1)). \quad (17)$$

By an identical argument to that made with no policy concerns, upon observing an off-path message, the payoff equivalence of the good and bad uninformed types implies the belief about competence in an MSE must be greater than or equal to π_g^\emptyset . So, the payoff to deviating to m_1 must be at least

$$p_1 + (1 - p_1)\pi_g^\emptyset + \gamma p_1$$

–and the corresponding policy concerns threshold to prevent this deviation is:

$$\pi_g^\emptyset + \gamma(1 - p_1(1 - p_1)) \geq p_1 + (1 - p_1)\pi_g^\emptyset + \gamma p_1$$

– which reduces to

$$\begin{aligned} \gamma &\geq p_1 \gamma_{NV}^H \\ &\equiv \gamma_{SV}^H \end{aligned} \quad (18)$$

Adding state validation weakens the condition required for an honest equilibrium, particularly when p_1 is close to $1/2$. However, this threshold is always strictly positive, so for

small policy concerns there can be no honesty even with state validation.

As shown in the proof of the following, if this condition is not met, then as with the no validation case there can be no admission of uncertainty. Further, since adding policy concerns does not change the classes of payoff equivalence, the case as $\gamma \rightarrow 0$ is the same as $\gamma = 0$.

Proposition 11. *With policy concerns and state validation:*

- i. *If $\gamma \geq \gamma_{SV}^H = p_1 \gamma_{NV}^H$, then there is an honest MSE,*
- ii. *If $\gamma \in (0, \gamma_{SV}^H)$, then all non-babbling MSE are always guessing (i.e., $\sigma_\emptyset^*(m_\emptyset) = 0$).*

Proof. Part i is demonstrated above

For part ii, our strategy mirrors the proof with no validation – that is, by way of contradiction, if the constraint for honesty is not met, then the payoff to sending m_1 is always strictly higher than m_\emptyset . As above, in any MSE where $\sigma_\emptyset(m_1) > 0$, the payoff for sending m_\emptyset is $\pi_g^\emptyset + \gamma(1 - p_1(1 - p_1))$. The payoff to sending m_1 is:

$$p_1 \pi_g(m_1, 1) + (1 - p_1) \pi_g^\emptyset + \gamma(1 - p_1(1 - \pi_1(m_1, \sigma_\emptyset(m_1))))^2 + (1 - p_1) \pi_1(m_1, \sigma_\emptyset(m_1))^2.$$

Next, the posterior beliefs of the decision-maker are the same as in the no validation case except:

$$\pi_g(m_1, 1) = \frac{Pr(\theta = g, m_1, \omega = 1)}{Pr(m_1, \omega = 1)} = \frac{p_1 p_g p_e + p_1 p_g (1 - p_e) \sigma_\emptyset(m_1)}{p_1 p_g p_e + p_1 (1 - p_g p_e) \sigma_\emptyset(m_1)} = \frac{p_g p_e + p_g (1 - p_e) \sigma_\emptyset(m_1)}{p_g p_e + (1 - p_g p_e) \sigma_\emptyset(m_1)}.$$

The difference between the payoffs for sending m_1 and m_\emptyset can be written:

$$p_e p_g p_1 \frac{z(\sigma_\emptyset(m_1); \gamma)}{(1 - p_e p_g)(p_e p_g (1 - \sigma_\emptyset(m_1)) - \sigma_\emptyset(m_1))(p_e p_g (p_1 - \sigma_\emptyset(m_1)) + \sigma_\emptyset(m_1))^2}$$

– where

$$z(\sigma_\emptyset(m_1); \gamma) = \gamma p_e p_g (-1 + p_e p_g) (-1 + p_1)^2 p_1 (p_e p_g (-1 + \sigma_\emptyset(m_1)) - \sigma_\emptyset(m_1)) \\ + (-1 + p_g) (p_e p_g (p_1 - \sigma_\emptyset(m_1)) + \sigma_\emptyset(m_1))^2.$$

So any equilibrium where both m_1 and m_\emptyset are sent is characterized by $z(\sigma_\emptyset(m_1); \gamma) = 0$. It is then sufficient to show that for $\gamma < \gamma_{SV}^H$, there is no $\sigma_\emptyset(m_1) \in [0, 1]$ such that $z(\sigma_\emptyset(m_1); \gamma) = 0$.

Formally, it is easy to check that z is strictly decreasing in γ and that $z(0, \gamma_{SV}^H) = 0$. So, $z(0, \gamma) > 0$ for $\gamma < \gamma_{SV}^H$. To show z is strictly positive for $\sigma_\emptyset(m_1) > 0$, first observe that:

$$\left. \frac{\partial z}{\partial \sigma_\emptyset(m_1)} \right|_{\gamma=\gamma_{SV}^H} = (1 - p_g)(1 - p_e p_g)(p_e p_g(2 - p_1)p_1 + (2 - 2p_e p_g)\sigma_\emptyset(m_1)) > 0$$

– and

$$\frac{\partial^2 z}{\partial \sigma_\emptyset(m_1) \partial \gamma} = -p_e p_g (1 - p_e p_g)^2 (1 - p_1)^2 p_1 < 0.$$

Combined, these inequalities imply $\frac{\partial z}{\partial \sigma_\emptyset(m_1)} > 0$ when $\gamma < \gamma_{SV}^H$. So, $z(\sigma_\emptyset(m_1), \gamma) > 0$ for any $\sigma_\emptyset(m_1)$ when $\gamma < \gamma_{SV}^H$, completing part ii. \square

Difficulty Validation. As shown in the main text, the condition for an honest equilibrium with difficulty validation and policy concerns is.

$$(1 - p_e)p_g + \gamma(1 - p_1(1 - p_1)) \geq p_e + \gamma p_1 \\ \gamma \geq \frac{p_e(1 + p_g) - p_g}{(1 - p_1)^2} \equiv \gamma_{DV}^H.$$

As discussed in the main text, γ_{DV}^H can be negative, meaning that there is an honest equilibrium even with no policy concerns.

Proposition 12. *With policy concerns and difficulty validation:*

- i. *If $\gamma \geq \gamma_{DV}^H$, then there is an honest MSE.*
- ii. *If $\gamma \leq \gamma_{DV}^H$, then there is an MSE where the uninformed good types admit uncertainty, and if $p_e \geq \frac{p_1}{2-p_1}$ there is an MSE where all of the good types send their honest message.*

Proof. Part i is shown above. For part ii, first note the equilibrium constructed in proposition 8 also holds with policy concerns: the policy choice upon observing both equilibrium messages is p_1 , so each type's relative payoff in this equilibrium is unaffected by the value of γ . Since the good uninformed types always admit uncertainty in this equilibrium, this demonstrates the first claim.

Now suppose the good types all send their honest message. By the same fixed point argument as proposition 9, the bad types must have at least one mixed strategy which is a best response given the good types strategy and DM strategy. What remains is to show the good types have no incentive to deviate from the honest message.

The message/validation combinations (m_0, e) , (m_1, e) , and (m_\emptyset, h) are on-path and yield competence evaluations which are all strictly greater than zero.

Message/validation combinations (m_0, h) , (m_1, h) , and (m_\emptyset, e) are never reached with a good type. So, if the bad types send those respective messages, they are on-path and the competence assessment must be zero. If these information sets are off-path the competence assessment can be set to zero.

Since only uninformed types send m_\emptyset , the policy choice upon observing m_\emptyset must be $a^*(m_\emptyset) = p_1$. The m_0 message is sent by the informed type who knows $\omega = 0$, and

potentially also by uninformed bad types, so $a^*(m_0) \in [0, p_1)$. Similarly, $a^*(m_1) \in (p_1, 1]$. So $a^*(m_0) < a^*(m_\emptyset) < a^*(m_1)$.

The good and uninformed type has no incentive to deviate from sending message m_\emptyset because for $m \in \{m_0, m_1\}$, $\pi_g(m_\emptyset, h) > \pi_g(m, h)$ and $v(a^*(m_\emptyset), p_1) > v(a^*(m), p_1)$.

The s_0 type has no incentive to deviate to m_\emptyset since $\pi_g(m_0, e) > \pi_g(m_\emptyset, e) = 0$ and $v(a^*(m_0), 0) > v(a^*(m_\emptyset), 0)$. Similarly, the s_1 type has no incentive to deviate to m_\emptyset .

So, the final deviations to check are for the informed types switching to the message associated with the other state; i.e., the s_0 types sending m_1 and the s_1 types sending m_0 . Preventing a deviation to m_1 requires:

$$\begin{aligned} \pi_g(m_0, e) + \gamma v(a^*(m_0), 0) &\geq \pi_g(m_1, e) + \gamma v(a^*(m_1), 0) \\ \Delta_\pi + \gamma \Delta_v(0) &\leq 0, \end{aligned} \tag{19}$$

where $\Delta_\pi \equiv \pi_g(m_1, e) - \pi_g(m_0, e)$ is the difference in competence assessments from sending m_1 versus m_0 (when the problem is easy), and $\Delta_v(p) \equiv v(a^*(m_1), p) - v(a^*(m_0), p)$ is the difference in the expected quality of the policy when sending m_1 vs m_0 for an expert who believes $\omega = 1$ with probability p . This simplifies to:

$$\Delta_v(p) = (a^*(m_1) - a^*(m_0))(2p - a^*(m_1) - a^*(m_0)).$$

Since $a^*(m_1) > a^*(m_0)$, $\Delta_v(p)$ is strictly increasing in p , and $\Delta_v(0) < 0 < \Delta_v(1)$.

The analogous incentive compatibility constraint for the s_1 types is:

$$\Delta_\pi + \gamma \Delta_v(1) \geq 0 \tag{20}$$

If the bad types never send m_0 or m_1 , then $\Delta_\pi = 0$, and (19)-(20) both hold. So, while not explicitly shown in the main text, in the honest equilibrium such a deviation is never profitable.

Now consider an equilibrium where the bad types send both m_0 and m_1 , in which case they must be indifferent between both messages:

$$\begin{aligned} p_e \pi_g(m_0, e) + \gamma v(a^*(m_0), p_1) &= p_e \pi_g(m_1, e) + \gamma v(a^*(m_1), p_1) \\ p_e \Delta_\pi + \gamma \Delta_v(p_1) &= 0 \end{aligned} \quad (21)$$

Substituting this constraint into (19) and (20) and simplifying gives:

$$p_e \Delta_v(0) - \Delta_v(p_1) \leq 0 \quad (22)$$

$$p_e \Delta_v(1) - \Delta_v(p_1) \geq 0. \quad (23)$$

If $\Delta_v(p_1) = 0$ the constraints are both met. If $\Delta_v(p_1) < 0$ then the second constraint is always met, and the first constraint can be written:

$$p_e \geq \frac{\Delta_v(p_1)}{\Delta_v(0)} = \frac{a^*(m_0) + a^*(m_1) - 2p_1}{a^*(m_0) + a^*(m_1)} \equiv \check{p}_\delta \quad (24)$$

This constraint is hardest to meet when \check{p}_δ is large, which is true when $a^*(m_0) + a^*(m_1)$ is high. The highest value this sum can take on is $p_1 + 1$, so $\check{p}_\delta \leq \frac{1-p_1}{1+p_1}$.

If $\Delta_v(p_1) > 0$, then the first constraint is always met, and the second constraint becomes:

$$p_e \geq \frac{\Delta_v(p_1)}{\Delta_v(1)} = \frac{2p_1 - (a^*(m_0) + a^*(m_1))}{2 - (a^*(m_0) + a^*(m_1))} \equiv \hat{p}_\delta \quad (25)$$

This is hardest to meet when $a^*(m_0) + a^*(m_1)$ is small, and the smallest value it can take on is p_1 . Plugging this in, $\hat{p}_\delta \geq \frac{p_1}{2-p_1} \geq \check{p}_\delta$.

For $p_1 \geq 1/2$, $\hat{p}_\delta \geq \check{p}_\delta$. Without placing any further restrictions on the value of $a^*(m_0) + a^*(m_1)$, this constraint ranges from $\hat{p}_\delta \in (1/3, 1)$. Still, if p_e is sufficiently high, the informed types never have an incentive to deviate when the bad types send both m_0 and m_1 .

If the bad types only send m_1 but not m_0 , then the s_0 types get the highest possible payoff, so the relevant deviation to check is the s_1 types switching to m_0 . The bad types sending weakly preferring m_1 implies $p_e \Delta_\pi + \gamma \Delta_v(p) \geq 0$, and substituting into equation 23 gives the same $p_e \geq \hat{p}_\delta$. Similarly, if the bad types only send m_0 but not m_1 , then the relevant constraint is the s_0 types sending m_1 , for which $p_e \geq \check{p}_\delta$ is sufficient.

Summarizing, a sufficient condition for the existence of a MSE where the good types report honestly (for any value of γ) is $p_e \leq p_g/(1 + p_g)$ (in which case $\gamma \leq \gamma_{DV}^H$), or $p_e \geq \frac{p_1}{2-p_1}$. This completes part ii.

Now to prove proposition 4 we first characterize the optimal strategy for the bad types as $\gamma \rightarrow 0$, assuming the good types send their honest message. If sending m_\emptyset , the expert will reveal his type if $\delta = e$, but appear partially competent if $\delta = h$, giving expected payoff

$$(1 - p_e) \frac{p_g}{p_g + (1 - p_g)\sigma_b(m_\emptyset)}.$$

When sending m_0 , the expert will reveal his type if $\delta = h$ (as only bad types guess when the problem is hard), but look partially competent if $\delta = e$:

$$p_e \frac{p_g(1 - p_1)}{p_g(1 - p_1) + (1 - p_g)\sigma_b(m_0)}.$$

and when sending m_1 the expect payoff is:

$$p_e \frac{p_g p_1}{p_g p_1 + (1 - p_g) \sigma_b(m_1)}.$$

setting these three equal subject to $\sigma_b(m_0) + \sigma_b(m_1) + \sigma_b(m_\emptyset) = 1$ gives:

$$\begin{aligned}\sigma_b(m_\emptyset) &= \frac{1 - p_e(1 + p_g)}{1 - p_g}; \\ \sigma_b(m_0) &= \frac{(1 - p_1)(p_e - p_g(1 - p_e))}{1 - p_g} \\ \sigma_b(m_1) &= \frac{p_1(p_e - p_g(1 - p_e))}{1 - p_g}.\end{aligned}$$

These are all interior if and only if:

$$0 < \frac{1 - p_e(1 + p_g)}{1 - p_g} < 1 \implies \frac{p_g}{1 + p_g} < p_e < \frac{1}{1 + p_g}.$$

If $p_e \leq \frac{p_g}{1 + p_g}$, then there can be no equilibrium where the bad expert uses a fully mixed strategy because he would always prefer to send m_\emptyset ; and recall this is exactly the condition for an honest equilibrium with no validation. If $p_e \geq \frac{1}{1 + p_g}$, then the bad type always guesses. Setting the payoff for a bad type sending m_0 and m_1 equal along with $\sigma_b(m_0) + \sigma_b(m_1) = 1$ gives the strategies in the statement of the proposition.

The final step is to ensure the informed types do not send the message associated with the other state. Recall the IC constraints depend on $a^*(m_0) + a^*(m_1)$, which we can now restrict to a narrower range given the bad type strategy:

$$\begin{aligned}a^*(m_0) + a^*(m_1) &= \frac{(1 - p_g)p_1(1 - p_1)(1 - \sigma_b(m_\emptyset))}{p_e p_g(1 - p_1 + (1 - p_g)(1 - p_1)(1 - \sigma_b(m_\emptyset)))} \\ &\quad + \frac{p_e p_g p_1 + (1 - p_g)p_1 p_1(1 - \sigma_b(m_\emptyset))}{p_e p_g p_1 + (1 - p_g)p_1(1 - \sigma_b(m_\emptyset))} \\ &= \frac{p_e p_g + (1 - \sigma_b(m_\emptyset))(1 - p_g)2p_1}{p_e p_g + (1 - \sigma_b(m_\emptyset))(1 - p_g)}.\end{aligned}$$

This can be interpreted as weighted average of 1 (with weight $p_e p_g$) and $2p_1 > 1$ (with weight $(1 - \sigma_b(m_0))(1 - p_g)$), and so must lie on $[1, 2p_1]$. So, (25) is always the binding constraint, and is hardest to satisfy when $a^*(m_0) + a^*(m_1) \rightarrow 1$, in which case the constraint becomes $\hat{p}_\delta = 2p_1 - 1$. So, $p_e \geq 2p_1 - 1$ is a sufficient condition for the informed types to never deviate. For any $p_e > 0$, this holds for p_1 sufficiently close to $1/2$, which completes the proof of proposition 4.

□

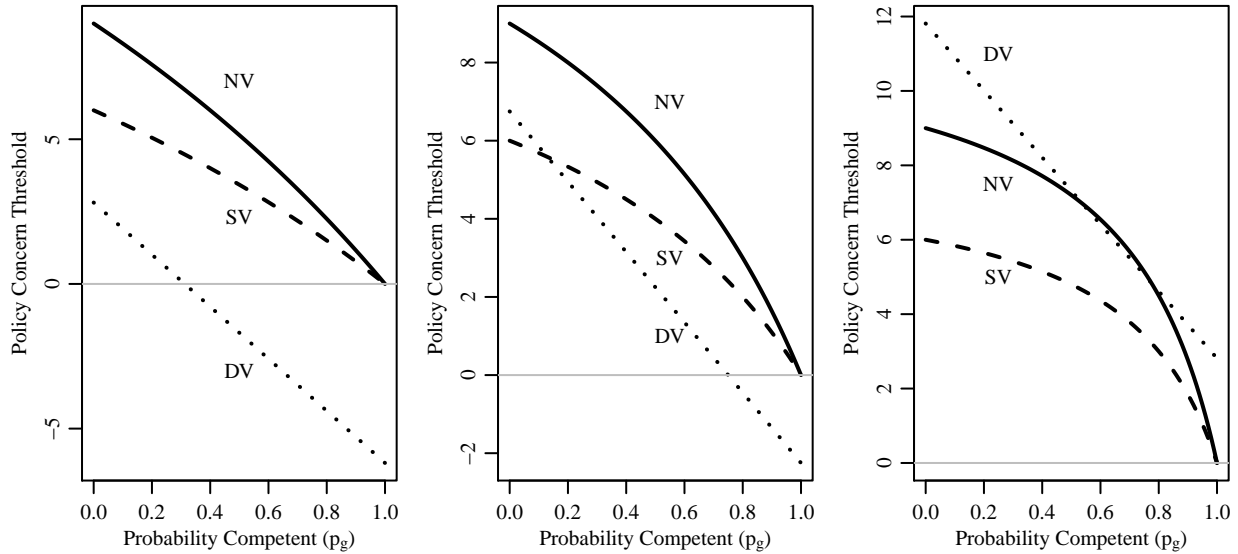
Here is an example where the constraint on the informed types is violated. Suppose p_1 is close to 1, and the bad types usually send m_1 , and rarely m_0 . Then the tradeoff they face is that sending m_1 leads to a better policy, but a lower competence payoff when the problem is easy (when the problem is hard, the competence payoff for either guess is zero). Now consider the good expert who observes signal s_1 . Compared to the bad expert, this type has a marginally stronger incentive to send m_1 (since p_1 is close to 1). However, this type *knows* that he will face a reputational loss for sending m_1 rather than m_0 , while the bad type only experiences this loss with probability p_e . So, the bad type being indifferent means the type who knows the state is 1 has a strict incentive to deviate to m_0 . In general, this deviation is tough to prevent when p_e is low and p_1 is close to 1, hence the condition in the proposition.

Comparative Statics: Difficulty Validation Can be the Wrong Kind of Transparency.

As long as policy concerns are strictly positive but small, difficulty validation is more effective at eliciting honesty than state validation.

For larger policy concerns the comparison becomes less straightforward. Figure C.2 shows the policy concern threshold for honesty under no validation (solid line), state validation

Figure C.2: Comparative Statics of Honesty Threshold



Notes: Comparison of threshold in policy concerns for full honesty under different validation regimes as a function of p_g . The panels vary in the likelihood the problem is solvable, which is 0.25 in the left panel, 0.5 in the middle panel, and 0.75 in the right panel.

(dashed line), and difficulty validation (dotted line) as a function of the prior on the expert competence, when the problem is usually hard ($p_e = 0.25$, left panel), equally likely to be easy or hard ($p_e = 0.5$, middle panel) and usually easy ($p_e = 0.75$, right panel). In all panels $p_1 = 0.67$; changing this parameter does not affect the conclusions that follow.²⁴ For intuition, difficulty validation makes it hard to compensate bad experts for saying “I don’t know,” as there are fewer good experts who don’t know. For very easy problems difficulty validation can be *worse* than no validation. This mirrors the result in Prat (2005), where transparency can eliminate incentives for bad types to pool with good types by exerting more effort.

This figure illustrates several results. First, in all cases, the policy concern threshold re-

²⁴In general, honesty is easier to sustain under all validation regimes when p_1 is lower, with state validation being particularly sensitive to this change.

quired is decreasing in p_g , which means it is easier to sustain honesty when the prior is that the expert is competent. This is because when most experts are competent in general, most uninformed experts are competent as well, and so there is less of a penalty for admitting uncertainty. Second, the threshold with state validation is always lower than the threshold with no validation, though these are always strictly positive as long as $p_g < 1$. Further, for most of the parameter space these thresholds are above two, indicating the expert must care twice as much about policy than about perceptions of his competence to elicit honesty. On the other hand, in the right and middle panels there are regions where the threshold with difficulty validation is below zero, indicating no policy concerns are necessary to induce admission of uncertainty (in fact, the expert could want the decision-maker to make a *bad* decision and still admit uncertainty).

Finally, consider how the relationship between the thresholds changes as the problem becomes easier. When problems are likely to be hard (left panel), difficulty validation is the best for eliciting honesty at all values of p_g . In the middle panel, difficulty validation is always better than no validation, but state validation is best for low values of p_g . When the problem is very likely to be easy, difficulty validation is always worse than state validation and is even worse than even no validation other than for a narrow range of p_g .

However, even in this case difficulty validation still can elicit honesty from good but uninformed experts when policy concerns are not high enough, while there is no admission of uncertainty at all when policy concerns are not high enough with no validation and state validation.

D Relabeling

Some of our formal results only rely on the existence of equilibria with certain properties. For these results the fact that we often restrict attention to the (m_0, m_1, m_\emptyset) message set poses no issues: it is sufficient to show that there is an equilibrium of this form with the claimed properties. However, propositions 2, 3, 10ii-iii, and 11, make claims that all (non-babbling) MSE have certain properties.²⁵ The proofs show that all equilibrium where the s_0 and s_1 types send distinct and unique messages (labelled m_0 and m_1) and there is at most one other message (labelled m_\emptyset) have these properties. Here we show this is WLOG in the sense that with no validation or state validation, any non-babbling equilibrium can be relabeled to an equilibrium of this form.

Consider a general messaging strategy where $M \subseteq \mathcal{M}$ is the set of messages sent with positive probability. Write the probability that the informed types observing s_0 and s_1 and $\sigma_0(m)$ and $\sigma_1(m)$. When the good and bad uninformed types are not necessarily payoff equivalent we write their strategies $\sigma_{\theta, \emptyset}(m)$. When these types are payoff equivalent and hence play the same strategy, we drop the θ : $\sigma_\emptyset(m)$. Similarly, let M_0 and M_1 be the set of messages sent by the respective informed types with strictly positive probability, and $M_{g, \emptyset}$, $M_{b, \emptyset}$, and M_\emptyset the respective sets for the uninformed types, divided when appropriate.

As is standard in cheap talk games, there is always a babbling equilibrium:

Proposition 13. *There is a class of babbling equilibria where $\sigma_0(m) = \sigma_1(m) = \sigma_{g, \emptyset}(m) = \sigma_{b, \emptyset}(m)$ for all $m \in M$.*

Proof. If all play the same mixed strategy, then $\pi_g(m, \mathcal{I}_{DM2}) = p_g$ and $a^*(m, \mathcal{I}_{DM}) = p_1$ for any $m \in M$ and \mathcal{I}_{DM} . Setting the beliefs for any off-path message to be the same as

²⁵Proposition 8 also makes a claim about all equilibria, but this is already proven in Appendix B.

the on-path messages, all types are indifferent between any $m \in \mathcal{M}$. □

The next result states that for all cases with either state validation or policy concerns, in any non-babbling equilibrium the informed types send no common message (note this result does *not* hold with difficulty validation; in fact, the proof of proposition 8 contains a counterexample):

Proposition 14. *With either no validation or state validation (and any level of policy concerns), any MSE where $M_0 \cap M_1 \neq \emptyset$ is babbling, i.e., $\sigma_0(m) = \sigma_1(m) = \sigma_{g,\emptyset}(m) = \sigma_{b,\emptyset}(m)$ for all $m \in M$.*

Proof. We first prove the result with state validation, and then briefly highlight the aspects of the argument that differ with no validation.

Recall that for this case the good and bad uninformed types are payoff equivalent, so we write their common message set and strategy M_\emptyset and $\sigma_\emptyset(m)$. The proof proceeds in three steps.

Step 1: If $M_0 \cap M_1 \neq \emptyset$, then $M_0 = M_1$. Let $m_c \in M_0 \cap M_1$ be a message sent by both informed types. Suppose there is another message sent only by the s_0 types: $m_0 \in M_0 \setminus M_1$. For the s_0 type to be indifferent between m_0 and m_c :

$$\pi_g(m_c, 0) + \gamma v(a^*(m_c), 0) = \pi_g(m_0, 0) + \gamma v(a^*(m_0), 0).$$

For this equation to hold, it must be the case that the uninformed types send m_0 with positive probability: if not, then $\pi_g(m_c, 0) \leq \pi_g(m_0, 0) = 1$, but $v(a^*(m_c), 0) < 1 = v(a^*(m_0), 0)$, contradicting the indifference condition.

For the uninformed types to send m_0 , it must also be the case that his expected payoff for sending this message, which can be written

$$p_1(\pi_g(m_0, 1) + \gamma v(a^*(m_0), 1)) + (1 - p_1)(\pi_g(m_0, 0) + v(a^*(m_0), 0))$$

– is at least his payoff for sending m_c :

$$p_1(\pi_g(m_c, 1) + \gamma v(a^*(m_c), 1)) + (1 - p_1)(\pi_g(m_c, 0) + v(a^*(m_c), 0)).$$

The second terms, which both start with $(1 - p_1)$, are equal by the indifference condition for s_0 types, so this requires:

$$\pi_g(m_0, 1) + \gamma v(a^*(m_0), 1) \geq \pi_g(m_c, 1) + \gamma v(a^*(m_c), 1).$$

Since m_0 is never sent by the s_1 types, $\pi_g(m_0, 1) = \pi_g^\emptyset$, while $\pi_g(m_c, 1) > \pi_g^\emptyset$. So, this inequality requires $v(a^*(m_0), 1) > v(a^*(m_c), 1)$, which implies $a^*(m_0) > a^*(m_c)$. A necessary condition for this inequality is $\frac{\sigma_\emptyset(m_0)}{\sigma_0(m_0)} > \frac{\sigma_\emptyset(m_c)}{\sigma_0(m_c)}$, which also implies $\pi_g(m_c, 0) > \pi_g(m_0, 0)$. But if $a^*(m_0) > a^*(m_c)$ and $\pi_g(m_c, 0) > \pi_g(m_0, 0)$, the s_0 types strictly prefer to send m_c rather than m_0 , a contradiction. By an identical argument, there can be no message in $M_1 \setminus M_0$, completing step 1.

Step 2: If $M_0 = M_1$, then $\sigma_0(m) = \sigma_1(m)$ for all m . If $M_0 = M_1$ is a singleton, the result is immediate. If there are multiple common messages and the informed types do not use the same mixed strategy, there must be a message m^0 such that $\sigma_0(m^0) > \sigma_1(m^0) > 0$ and another message m^1 such that $\sigma_1(m^1) > \sigma_0(m^1) > 0$. (We write the message “generally sent by type observing s_x ” with a superscript to differentiate between the subscript notation referring to messages always sent by type s_x .) The action taken by the DM upon observing m^0 must be strictly less than p_1 and upon observing m^1 must be strictly greater than p_1 ,²⁶

²⁶The action taken upon observing m can be written $\mathbb{P}(s_1|m) + p_1\mathbb{P}(s_0|m)$. Rearranging, this is greater

so $a^*(m^0) < a^*(m^1)$.

Both the s_1 and s_0 types must be indifferent between both messages, so:

$$\pi_g(m^0, 0) + \gamma v(a^*(m^0), 0) = \pi_g(m^1, 0) + \gamma v(a^*(m^1), 0)$$

$$\pi_g(m^0, 1) + \gamma v(a^*(m^0), 1) = \pi_g(m^1, 1) + \gamma v(a^*(m^1), 1)$$

Since $v(a^*(m^0), 0) > v(a^*(m^1), 0)$, for the s_0 to be indifferent it must be the case that $\pi_g(m^0, 0) < \pi_g(m^1, 0)$. Writing out this posterior belief:

$$\mathbb{P}(\theta = g|m, 0) = \frac{(1 - p_1)(p_g p_e \sigma_0(m) + (1 - p_e) \sigma_\emptyset(m))}{(1 - p_1)(p_g p_e \sigma_0(m) + (1 - p_g p_e) \sigma_\emptyset(m))}.$$

Rearranging, $\pi_g(m^0, 0) < \pi_g(m^1, 0)$ if and only if $\frac{\sigma_0(m^0)}{\sigma_0(m^1)} < \frac{\sigma_\emptyset(m^0)}{\sigma_\emptyset(m^1)}$. Similarly, it must be the case that $\pi_g(m^1, 1) < \pi_g(m^0, 1)$, which implies $\frac{\sigma_1(m^0)}{\sigma_1(m^1)} > \frac{\sigma_\emptyset(m^0)}{\sigma_\emptyset(m^1)}$. Combining, $\frac{\sigma_0(m^0)}{\sigma_0(m^1)} < \frac{\sigma_1(m^0)}{\sigma_1(m^1)}$, which contradicts the definition of these messages. So, $\sigma_0(m) = \sigma_1(m)$ for all m .

Step 3: If $M_0 = M_1$ and $\sigma_0(m) = \sigma_1(m)$, then $M_\emptyset = M_0 = M_1$ and $\sigma_\emptyset(m) = \sigma_0(m) = \sigma_1(m)$. By step 2, it must be the case that $a^*(m) = p_1$ for all messages sent by the informed types. So, the uninformed types can't send a message not sent by the informed types: if so, the payoff would be at most $\pi_g^\emptyset + \gamma v(p_1, p_1)$, which is strictly less than the payoff for sending a message sent by the informed types. If there is only one message in M then the proof is done. If there are multiple types, all must be indifferent between each message, and by step 2 they lead to the same policy choice. So, they must also lead to the same competence assessment for each revelation of ω , which is true if and only if $\sigma_\emptyset(m) = \sigma_0(m) = \sigma_1(m)$.

□

than p_1 if and only if $\frac{\mathbb{P}(s_1, m)}{\mathbb{P}(s_1, m) + \mathbb{P}(s_0, m)} > p_1$ which holds if and only if $\sigma_1(m) > \sigma_0(m)$.

Next, consider the no validation case. For step 1, define m_0 and m_1 analogously. The uninformed types must send m_0 by the same logic, and these types at least weakly prefer sending this to m_c (while the s_0 types are indifferent) requires:

$$\pi_g(m_0) + \gamma v(a^*(m_0), 1) \geq \pi_g(m_c) + \gamma v(a^*(m_c), 1).$$

This can hold only weakly to prevent the s_1 types from sending m_0 (as required by the definition). Combined with the s_0 indifference condition:

$$\pi_g(m_0) - \pi_g(m_c) = \gamma v(a^*(m_c), 1) - \gamma v(a^*(m_0), 1) = \gamma v(a^*(m_c), 0) - \gamma v(a^*(m_0), 0),$$

which requires $a^*(m_0) = a^*(m_c)$. Since the s_1 types send m_c but not m_0 this requires $\frac{\sigma_\theta(m_0)}{\sigma_0(m_0)} > \frac{\sigma_\theta(m_c)}{\sigma_0(m_c)}$, which implies $\pi_g(m_0) < \pi_g(m_c)$, contradicting the s_0 types being indifferent between both messages.

Steps 2 and 3 follow the same logic.

□

Finally, we prove that any MSE where the messages sent by the s_0 and s_1 types do not overlap is equivalent to an MSE where there is only one message sent by each of these types and only one “other” message. This provides a formal statement of our claims about equilibria which are “equivalent subject to relabeling”:

Proposition 15. *Let $M_U = M_\emptyset \setminus (M_0 \cup M_1)$ (i.e., the messages only sent by the uninformed types). With no validation or state validation:*

- i. *In any MSE where $M_0 \cap M_1 = \emptyset$, for $j \in \{0, 1, U\}$, and any $m', m'' \in M_j$, $a^*(m') = a^*(m'')$ and $\pi_g(m', \mathcal{I}_{DM2}) = \pi_g(m'', \mathcal{I}_{DM2})$*
- ii. *Take an MSE where $|M_j| > 1$ for any $j \in \{0, 1, U\}$, and the equilibrium actions and pos-*

terior competence assessments for the messages in this set are $a^*(m_i)$ and $\pi_g(m_i, \mathcal{I}_{DM2})$ (which by part i are the same for all $m_i \in M_j$). Then there is another MSE where $M_j = \{m\}$, and equilibrium strategy and beliefs a_{new}^* and $\pi_{g,new}$ such that $a^*(m_i) = a_{new}^*(m)$, and $\pi_g(m_i, \mathcal{I}_{DM2}) = \pi_{g,new}(m, \mathcal{I}_{DM2})$

Proof. For part i, first consider the message in M_U . By construction the action taken upon observing any message in this set is p_1 . And since the good and bad uninformed types are payoff equivalent and use the same strategy, the competence assessment upon observing any message in this set must be π_g^θ .

For M_0 , first note that for any $m', m'' \in M_0$, it can't be the case that the uninformed types only send one message but not the other with positive probability; if so, the message not sent by the uninformed types would give a strictly higher payoff for the s_0 types, and hence they can't send both messages. So, either the uninformed types send neither m' nor m'' , in which case the result is immediate, or they send both, in which case they must be indifferent between both. As shown in the proof of proposition 14, this requires that the action and competence assessment are the same for both m' and m'' . An identical argument holds for M_1 , completing part i.

For part ii and M_θ , the result immediately follows from the same logic as part i.

For M_0 , if the uninformed types do not send any messages in M_0 , then the on-path response to any $m_0^j \in M_0$ are $a^*(m_0^j) = 0$ and $\pi_g(m_0^j, 0) = 1$. Keeping the rest of the equilibrium fixed, the responses in a proposed MSE where the s_0 types always send m_0 are also $a_{new}^*(m_0) = 0$ and $\pi_{g,new}(m_0^j, 0) = 1$. So there is an MSE where the s_0 types all send m_0 which is equivalent to the MSE where the s_0 types send multiple messages.

If the uninformed types do send the messages in M_0 , then part i implies all messages must

lead to the same competence evaluation, which implies for any $m'_0, m''_0 \in M_0$, $\frac{\sigma_\theta(m'_0)}{\sigma_0(m'_0)} = \frac{\sigma_\theta(m''_0)}{\sigma_0(m''_0)} \equiv r_0$. In the new proposed equilibrium where $M_0 = \{m_0\}$, set $\sigma_{0,\text{new}}(m_0) = 1$ and $\sigma_{\theta,\text{new}}(m_0) = r_0$. Since $\frac{\sigma_{\theta,\text{new}}(m_0)}{\sigma_{0,\text{new}}(m_0)} = \frac{\sigma_\theta(m'_0)}{\sigma_0(m'_0)}$, $a_{\text{new}}^*(m_0) = a^*(m'_0)$ and $\pi_{g,\text{new}}(m'_0, 0) = 1$, and all other aspects of the MSE are unchanged. \square

E Alternative Signal Structure

In this section, we consider two alternative signal specifications.

E.1 More general binary signal structure

Here is a more general formulation of the signal structure. We again assume a binary incumbent type $\theta \in \{g, b\}$ and problem difficulty $\delta \in \{e, h\}$. Now assume that the signal is given by:

$$s = \begin{cases} s_\omega & \text{with probability } P(\theta, \delta) \\ s_\emptyset & \text{o.w} \end{cases} \quad (26)$$

where $P(g, \delta) \geq P(b, \delta)$ (with the inequality strict for at least one δ) and $P(\theta, e) \geq P(\theta, h)$ (with the inequality strict for at least one θ). That is, more competent experts are (weakly) more likely to get an informative signal for either problem difficulty, and easy problems are (weakly) more likely to result in an informative signal. We assume that at least one of the inequalities is strict so that both variable “matter”.

All other aspects of the model are the same as in the main text.

The analysis in the main text is a special case of these assumptions where $P(\theta, \delta)$ is equal to 1 if $\theta = g$ and $\delta = e$ and zero otherwise. With the more general signal structure, there are now up to 6 potential types, as a function of the expert signal and competence. A type of competence θ who observes $s_x, x \in \{0, 1\}$ (if such a signal is possible for type θ ; recall this is not possible in the main formulation for $\theta = b$ types) knows that $\omega = x$, and his posterior belief about the problem difficulty is:

$$Pr(\delta = e | s_x, \theta) = \frac{p_e p_x p_\theta P(\theta, e)}{p_e p_x p_\theta P(\theta, e) + p_h p_x p_\theta P(\theta, h)} = \frac{p_e p_\theta P(\theta, e)}{p_e p_\theta P(\theta, e) + p_h p_\theta P(\theta, h)}. \quad (27)$$

Since $P(\theta, e) \geq P(\theta, h)$, $Pr(\delta = e | s_x, \theta) \geq p_e$, and if $P(\theta, e) > P(\theta, h)$ the inequality is strict. That is, since each type is (weakly) more likely to observe an informative signal when the problem is easy, they are (weakly) more likely to believe the problem is easy given an informative signal.

Also important for what comes, both types have an equal belief about the problem difficulty if and only if $\frac{P(g,e)}{P(g,h)} = \frac{P(b,e)}{P(b,h)}$. This condition might hold. For example, suppose $P(b, h) = 1/4$, $P(b, e) = 1/2$, $P(g, h) = 1/2$, and $P(g, e) = 1$. Then both types are twice as likely to receive an informative signal when the problem is easy, and hence learn the same about the problem difficulty from getting an informative signal. However, this is a knife-edged condition, and if $\frac{P(g,e)}{P(g,h)} > \frac{P(b,e)}{P(b,h)}$ the competent type will update about the easiness of the problem more sharply and if $\frac{P(g,e)}{P(g,h)} < \frac{P(b,e)}{P(b,h)}$ the bad type will update more in the positive direction.

A type of competence θ who observes $s = s_\theta$ maintains his prior belief about the state ($Pr(\omega = 1 | s_\theta) = p_1$) and his belief about the problem difficulty becomes:

$$Pr(\delta = e | s_\theta, \theta) = \frac{p_e p_\theta (1 - P(\theta, e))}{p_e p_\theta (1 - P(\theta, e)) + p_h p_\theta (1 - P(\theta, h))}. \quad (28)$$

Following a similar logic as the above, both sides will (weakly) come to believe the problem is less likely to be easy when getting an uninformative signal. These updates are equal if and only if $\frac{1-P(g,e)}{1-P(g,h)} = \frac{1-P(b,e)}{1-P(b,h)}$

As with the example in the main text, MSE helps quickly pin down which types can send different messages in equilibrium.

No policy concerns, no validation The analysis with no policy concerns or validation is identical: all types are payoff equivalent, and any equilibrium is babbling.

No policy concerns, state validation With state validation, types who observe different signals have different beliefs about ω , but types observing the same signal are still payoff equivalent regardless of their competence.²⁷ However, this will never induce honesty for a similar reason as the main model. In an honest equilibrium, the posterior belief about the expert competence when observing $(m_1, \omega = 1)$ is:

$$Pr(\theta = g|m_1, \omega = 1) = \frac{p_1 p_g (p_e P(g, e) + p_h P(g, h))}{p_1 p_g (p_e P(g, e) + p_h P(g, h)) + p_1 p_b (p_e P(b, e) + p_h P(b, h))} > p_g, \quad (29)$$

where the inequality follows from the assumption that our assumption that $P(g, \delta) \geq P(b, \delta)$ for both δ and one of the inequalities is strict, and hence $(p_e P(g, e) + p_h P(g, h)) > (p_e P(b, e) + p_h P(b, h))$.

Similarly, $Pr(\theta = g|m_1, \omega = 1) > p_g$ and $Pr(\theta = g|m_\emptyset) < p_g$. Given the Markov strategies requirement, Markov consistency implies that the worst inference the DM can

²⁷Depending on the P function they might have different views of the problem difficulty for some signals, but since there is no difficulty validation this does not affect their expected payoff for any possible DM strategy.

make about the expert competence upon observing something off path is $Pr(\theta = g|s = s_\emptyset) = Pr(\theta = g|m_\emptyset)$. So, the expected utility of sending m_1 is:

$$p_1 Pr(\theta = g|m_1, \omega = 1) + (1 - p_1) Pr(\theta = g|s_\emptyset) > Pr(\theta = g|m_\emptyset) \quad (30)$$

and hence this is a profitable deviation. So, there is no honest MSE with state validation alone.

Difficulty validation and small policy concerns With difficulty validation and small policy concerns, no two types with different beliefs about the difficulty of the problem are always payoff equivalent. So, as long as $\frac{P(g,e)}{P(g,h)} \neq \frac{P(b,e)}{P(b,h)}$ and $\frac{1-P(g,e)}{1-P(g,h)} \neq \frac{1-P(b,e)}{1-P(b,h)}$, the Markov strategies and Markov consistency requirements have no bite, making it possible to punish those who guess incorrectly with a belief that they are competent with probability zero.

However, an important aspect for this to make honesty possible is for sending an informative message when the problem is hard is actually off-path. In a proposed honest equilibrium with just difficulty validation, if $P(\theta, h) > 0$ for some θ , then not only are message/validation combinations (m_1, h) and (m_1, h) on path, but tend to be reached when the expert is competent.

So, an important assumption to make difficulty (and small policy) concerns effective at inducing honesty is that $P(\theta, h) = 0$ (as was true in the main model). In words, this implies that there are not only relatively hard or easy questions that the expert might be asked, but there are *impossible* questions. We think in most domains this is reasonable, particularly if we interpret informative messages as stating that the state is zero or one with certainty.

If so, the key constraint for sustaining an honest equilibrium is that good and bad uninformed experts face no incentive to guess. The payoff to admitting uncertainty (i.e., sending m_0) for type θ is:

$$Pr(\delta = e|s_0, \theta)Pr(\theta = g|s_0, \delta = e) + Pr(\delta = h|s_0, \theta)Pr(\theta = g|s_0, \delta = h). \quad (31)$$

The $Pr(\delta|s_0, \theta)$ are derived above, and the second halves are:

$$Pr(\theta = g|s_0, \delta) = \frac{p_g(1 - P(g, \delta))}{p_g(1 - P(g, \delta)) + (1 - p_g)(1 - P(b, \delta))}. \quad (32)$$

The payoff to guessing m_1 (which is a better deviation than m_0 as long as $p_1 \geq 1/2$.) is the probability that his guess is correct and the problem is easy (since an informative signal with a hard problem is off path), times the competence evaluation in this circumstance:

$$p_1 Pr(\delta = e|s_0, \theta)Pr(\theta = g|s_1, m_1, \delta = e), \quad (33)$$

where

$$Pr(\theta = g|s_1, m_1, \delta = e) = \frac{p_g P(g, e)}{p_g P(g, e) + (1 - p_g)P(b, e)}. \quad (34)$$

So, as $\gamma \rightarrow 0$, the condition for an honest equilibrium is that:

$$\begin{aligned} &Pr(\delta = h|s_0, \theta)Pr(\theta = g|s_0, \delta = h) \geq \\ &Pr(\delta = e|s_0, \theta)(p_1 Pr(\theta = g|s_1, m_1, \delta = e) - Pr(\theta = g|s_0, \delta = e)) \end{aligned} \quad (35)$$

for $\theta \in \{g, b\}$. While the algebra is messier, the core idea is just like in the main case. The trade-off here is that if the problem turns out to be easy it can be more profitable to guess, while when the problem turns out to be hard it is better to admit uncertainty.

Full Validation, no policy concerns Finally, consider the full validation case. Because of state validation it is possible for experts with different views about the state to send different messages, and as long as $\frac{P(g,e)}{P(g,h)} \neq \frac{P(b,e)}{P(b,h)}$ and $\frac{1-P(g,e)}{1-P(g,h)} \neq \frac{1-P(b,e)}{1-P(b,h)}$ it is possible to set any off path beliefs to zero.

This allows for the possibility of an honest equilibrium even without the assumption that some problems are impossible, since in an honest equilibrium incorrect guesses are never on path, and unlike the case with just state validation can be punished with an off-path belief that the expert must be the bad type, regardless of what the difficulty validation says.

The utility for admitting uncertainty in such an honest equilibrium is again given by equation (31). The expected utility for guessing 1 (assuming that sending an informative signal when validation reveals that $\delta = h$ is on-path) is:

$$p_1(p_e Pr(\theta = g|s_1, \delta = e) + p_h Pr(\theta = g|s_1, \delta = h)) \quad (36)$$

so an honest equilibrium is possible if:

$$\begin{aligned} & p_e(Pr(\theta = g|s_\emptyset, \delta = e) - p_1 Pr(\theta = g|s_1, \delta = e)) \\ & + p_h(Pr(\theta = g|s_\emptyset, \delta = h) - p_1 Pr(\theta = g|s_1, \delta = h)) \geq 0. \end{aligned} \quad (37)$$

As $p_1 \rightarrow 1$, this inequality never holds, but for smaller p_1 it is possible.

E.2 Continuous expert competence and question difficulty

Now let the competence of the expert θ and the difficulty of the problem δ both be uniform on $[0, 1]$. Let the private signal to the expert be:

$$s = \begin{cases} s_0 & \omega = 0, \theta > \beta\delta \\ s_1 & \omega = 1, \theta > \beta\delta \\ s_\emptyset & \text{o/w} \end{cases}$$

where $\beta > 0$.

If $\beta < 1$, then even the hardest problems ($\delta = 1$) are solvable by a strictly positive proportion of experts. If $\beta > 1$, then there are some problems which are so difficult that no expert can solve them. For reasons which will become apparent, we focus on the case where $\beta > 1$, and so $\bar{\delta} = 1/\beta < 1$ is the “least difficult unsolvable problem”.

The expert learns s and θ , which is always partially informative about δ . In particular, an expert who gets an informative signal knows that $\delta \in [0, \theta/\beta]$, and an expert who does not get an informative signal knows that $\delta \in [\theta/\beta, 1]$. An interesting contrast with the binary model is that better experts don’t always know more about the problem difficulty: when they learn the state, the range of possible values of δ is increasing in θ . However, when the expert learns the state (particularly with state validation) knowing the difficulty is not particularly relevant. On the other hand, when the expert is uninformed those who are more competent can restrict the difficulty of the problem to a smaller interval.

As with the binary model, we search for honest equilibria in the sense that the expert fully separates with respect to their signal (if not with respect to their competence). Here we only consider the case with no policy concerns.

No validation, State validation, Difficulty validation Even though the state space is much larger in this version of the model, with no validation (and no policy concerns) all types are still payoff equivalent. So any MSE must be babbling.

Similarly, with state validation there are now multiple types (differentiated by θ) who learn the state and multiple types who do not learn the state. However, since the knowledge of the state is the only payoff-relevant component of the type space, all types observing a particular s must play the same strategy in an MSE. So, by the same logic as the binary model, there is no honest MSE.

Difficulty validation alone (and again with no policy concerns; with small policy concerns honesty might be possible for some parameters) also hits the same problem as in the binary model. Among the informed types with competence θ , those observing s_0 and those observing s_1 are payoff equivalent, and so no information can be communicated about the state. It is possible that information about the difficulty of the problem can be conveyed.

Full Validation Now consider the full validation case. No pairs of types are payoff equivalent, since even those observing the same signal have different beliefs about what the difficulty validation will reveal for each value of θ . So, it is possible to use punitive off-path beliefs where those who guess incorrectly or when no expert could solve the problem, which is possible when $\beta > 1$.

We now show an honest equilibrium can be possible in this case. First, consider the on-path inferences by the DM. When seeing a correct message and difficulty δ , the DM knows the expert competence must be on $[\beta\delta, 1]$, and so the average competence assessment is:

$$\pi_g(m_1; \omega = 1, \delta) = \pi_g(m_0, \omega = 0, \delta) = \frac{\beta\delta + 1}{2}$$

which is at least $1/2$, and increasing in δ .

Upon observing m_\emptyset and δ , there are two possible cases. If $\delta > 1/\beta$, then no expert could have solved the problem, and so there is no information conveyed about the expert competence. If $\delta < 1/\beta$, then the DM learns that the expert competence is uniform on $[0, \beta\delta]$. Combining:

$$\pi_g(m_\emptyset, \omega, \delta) = \begin{cases} 1/2 & \delta > 1/\beta \\ \frac{\beta\delta}{2} & \delta \leq 1/\beta \end{cases}.$$

All other message and validation combinations are off-path, and can be set to zero.

Now consider the expert payoffs.

An informed expert (of any competence) gets a payoff of $\frac{\beta\delta+1}{2} > 1/2$ for sending the equilibrium message, 0 for sending the other informed message (i.e., m_1 rather than m_0 when $s = s_0$), and $\pi_g(m_\emptyset, \omega, \delta) \leq 1/2$ for sending m_\emptyset . So these types never deviate.

Uninformed experts know the difficulty – which again will be revealed to the DM – is uniform on $[\theta/\beta, 1]$. Note that for all but the (measure zero) $\theta = 1$ types, $1/\beta$ lies on this interval. So, all but the most competent experts don't know for sure if the problem is solvable by some experts, though very competent experts can become nearly certain the problem is unsolvable.

We can write the expected competence for admitting uncertainty to be the probability that $\delta \geq 1/\beta$ times $1/2$, plus the probability that $\delta < 1/\beta$ times the average competence assigned on this interval. Since the competence assessment is linear in δ on this interval, ranging from $\frac{\beta(\theta/\beta)}{2} = \frac{\theta}{2}$ to $1/2$, this average is $\frac{\theta+1}{4}$. Combining, the expected competence

for sending m_0 is:

$$\begin{aligned}
\mathbb{E}_\delta[\pi_g(m_0, \omega, \delta)] &= Pr(\delta \leq 1/\beta) \frac{\theta + 1}{4} + Pr(\delta > 1/\beta) \frac{1}{2} \\
&= \frac{1/\beta - \theta/\beta}{1 - \theta/\beta} \frac{\theta + 1}{4} + \frac{1 - 1/\beta}{1 - \theta/\beta} \frac{1}{2} \\
&= \frac{1 - \theta}{\beta - \theta} \frac{\theta + 1}{4} + \frac{\beta - 1}{\beta - \theta} \frac{1}{2}
\end{aligned}$$

which is increasing in θ .

Next consider the payoff for sending m_1 ; as before, this is the “better guess” since it is more likely to be matched by the state validation. This will lead to a competence evaluation of $\frac{\beta\delta+1}{2}$ if $\omega = 1$ (probability p_1) and if $\delta < 1/\beta$ (probability $\frac{1-\theta}{\beta-\theta}$), and 0 otherwise. Since the guessing correct payoff is linear in δ and the belief about δ conditional on a solvable problem is uniform on $[\theta/\beta, 1/\beta]$, the average competence assessment when getting away with a guess is:

$$\frac{\frac{\theta+1}{2} + 1}{2} = \frac{3 + \theta}{4}.$$

So the payoff to this deviation is:

$$p_1 \frac{1 - \theta}{\beta - \theta} \frac{3 + \theta}{4}$$

which is decreasing in θ .

So, the binding constraint is that the $\theta = 0$ prefers sending m_0 , which again reinforces the

assumption made about off-path beliefs. Honesty is possible when:

$$\frac{1}{\beta} \frac{1}{4} + \frac{\beta - 1}{\beta} \frac{1}{2} \geq p_1 \frac{1}{\beta} (3/4)$$

$$\beta \geq (3/2)p_1 + 1/2.$$

Since $p_1 \in [1/2, 1]$, this threshold ranges from $5/4$ to 2 . That the threshold in β is strictly greater than 1 means that there must be some possibility of getting caught answering an unanswerable question. The threshold is lower when p_1 is lower since this makes guessing less attractive as one is more likely to be caught guessing wrong.