# The Complexity of Agreement

Scott Aaronson*

## Abstract

A celebrated 1976 theorem of Aumann asserts that honest, rational Bayesian agents with common priors will never "agree to disagree": if their opinions about any topic are common knowledge, then those opinions must be equal. Economists have written numerous papers examining the assumptions behind this theorem. But two key questions went unaddressed: first, can the agents reach agreement after a conversation of reasonable length? Second, can the computations needed for that conversation be performed efficiently? This paper answers both questions in the affirmative, thereby strengthening Aumann's original conclusion.

We first show that, for two agents with a common prior to agree within $\varepsilon$ about the expectation of a $[0, 1]$ variable with high probability over their prior, it suffices for them to exchange order $1/\varepsilon^2$ bits. This bound is completely independent of the number of bits $n$ of relevant knowledge that the agents have. We then extend the bound to three or more agents; and we give an example where the economists' "standard protocol" (which consists of repeatedly announcing one's current expectation) nearly saturates the bound, while a new "attenuated protocol" does better. Finally, we give a protocol that would cause two Bayesians to agree within $\varepsilon$ after exchanging order $1/\varepsilon^2$ messages, and that can be *simulated* by agents with limited computational resources. By this we mean that, after examining the agents' knowledge and a transcript of their conversation, no one would be able to distinguish the agents from perfect Bayesians. The time used by the simulation procedure is exponential in $1/\varepsilon^6$ but not in $n$.

## 1   Introduction

A vast body of work in AI, economics, philosophy, and other fields seeks to model human beings as *Bayesian agents*—agents that start out with some prior probability distribution over possible states of the world, then update the distribution as they gather new information [16]. Because of its simplicity, the "humans-as-roughly-Bayesians" thesis has remained popular, despite the work of Allais [1], Tversky and Kahneman [19], and others; and despite well-known problems such as old evidence [7]. But one aspect of human experience seems especially hard to reconcile with the thesis.

Pick any two people, and there will be some topic they disagree about: capitalism versus socialism, the Israeli-Palestinian conflict, the interpretation of quantum mechanics, etc.[1] The more intelligent the people, the easier it will be to find such a topic. If they discuss the topic, chances are excellent that they will not reach agreement, but will instead become more confirmed in their previous beliefs. This is so even if the people respect each other's intelligence and honesty.

The above facts are known to everyone, yet as Aumann [2] observed in 1976, they constitute a serious challenge to Bayesian accounts of human reasoning. For suppose Alice and Bob are Bayesians, who have the same prior probabilities for all states of the world, but who have since gained different knowledge and thus have different posterior probabilities. Suppose further that, conditioned on everything she knows, Alice assigns a posterior probability $p$ to (say) extraterrestrial life existing. Bob likewise assigns a posterior probability $q$. Then provided both agents know $p$ and $q$ (and know that they know them, etc.), Aumann showed that $p$ and $q$ must be equal. This is true even if neither agent has any idea on what sort of evidence the other's estimate is based. For the sort of evidence can itself be considered a random variable, which is ultimately governed by a prior probability distribution that is the same for both agents.

Admittedly, the agents are unlikely to agree *immediately* after exchanging $p$ and $q$. For conditioned on Alice's estimate being $p$, Bob will revise his estimate $q$, and similarly Alice will revise $p$ conditioned on Bob's

---

[1]If you disagree with this assertion, you are simply providing further evidence for it!

estimate $q$. The agents will then have to exchange their *new* estimates $p'$ and $q'$, and so on iteratively. But provided the set of possible states is finite, it is easy to show that this iterative process must terminate eventually, with both agents having the same estimate [8]. In conclusion, then, there is no reason for the agents ever to disagree about anything!

On hearing this theorem for the first time, all of us come up with plausible ways in which actual human beings might evade its conditions. People have self-serving biases; they often discard or distort evidence that goes against what they want to believe [9]. (According to an often-cited study [5], 94% of professors consider themselves better than their average colleagues.[2]) People might interpret the same assertion differently. Or the assertion might be inherently ambiguous, if it deals with beauty or morality for instance. People might weigh the same evidence by different criteria. They might not understand the evidence. They might defend their opinions as high-school debaters do, out of sport rather than a desire for truth. They might not report their opinions with candor; or if they do, they might not trust others to do likewise.

In our view, the real challenge is not to list such caveats, but to sift through them and to discover which ones are fundamental. As an illustration, several of the caveats listed above disappear once we assume that all people have a common prior. For among other things, such a prior would assign common probabilities to all possible ways of parsing an ambiguous sentence, and to all possible ways of weighing evidence. Understandably, then, much of the criticism of Aumann's theorem has focused on the common prior assumption (see [3, 4, 10] for a discussion of that assumption).

But suppose we accept that two people have different priors. The obvious question is, *what caused their priors to differ?* Different career choices? Different friends? Different kindergarten teachers? Whatever is named as the first influence, we need merely go back in time to before that influence took effect. At the earlier time, the two people had the same prior by assumption. So at later times, they would not really have different priors, just different posteriors obtained by starting from the same prior and then conditioning on different life experiences. If we push this reasoning to its limit, as Cowen and Hanson [4] do, we are left wondering whether prior differences could be encoded in DNA at conception. Even then, how much confidence should you place in an opinion, if you know that were your genes different, you would have the opposite opinion? More generally, on what grounds can you favor your own prior over another's? For all you know, your prior was "switched by accident" with someone else's at birth!

After staring into the metaphysical abyss of prior differences, the natural reaction of a computer scientist is to step back, and ask if there is some simpler explanation for why Aumann's theorem fails to describe the real world. Recall that in the theorem, Alice's and Bob's opinions only became equal by the end of a hypothetical conversation. Might that conversation last an absurdly long time? After all, if Alice and Bob exchanged everything they knew, then clearly they would agree about everything! But presumably they are not Siamese twins, and do not have their entire lives to talk to each other. Thus *communication complexity* might provide a fundamental reason for why even honest, rational people could agree to disagree. Indeed, this was our conjecture when we began studying the topic.

*Computational complexity* provides a second promising reason. If a "state of the world" consists of $n$ bits, then Aumann's theorem requires Alice and Bob to represent a prior probability distribution over $2^n$ possible states. Even worse, it requires them to calculate expectations over that distribution, and update it conditioned on new information. If $n$ is (say) 10000, then this is obviously too much to ask.

## 1.1 Summary of Results

This paper initiates the study of the communication complexity and computational complexity of agreement protocols. Its surprising conclusion is that complexity is *not* a major barrier to agreement—at least, not nearly as major as it seems from the above arguments. In our view, this conclusion strengthens Aumann's original theorem substantially, by forcing our attention back to the origin of prior differences.

For economists, the main novelty of the paper will be our relentless use of asymptotic analysis. We will never be satisfied to show that a protocol terminates eventually. Instead we will always ask: do the resources needed for the protocol scale 'reasonably' with the parameters of the problem being solved? Here 'resources' include the number of messages, the number of bits per message, and the number of computational steps; while 'parameters' include the number of agents, the number of bits each agent is given, and the desired

---

[2]This is logically possible, but one assumes the response would be similar were the professors asked about their *median* colleagues.

accuracy and probability of success. This approach will let us model the limitations of real-world agents without sacrificing simplicity and elegance.

For computer scientists, the main novelty will be that, when we analyze the communication complexity of a function $f$, we care only about how long it takes some set of agents to agree *among themselves* about the expectation of $f$. Whether the agents' expectations agree with external reality is irrelevant.

After introducing notation in Section 2, in Section 3 we present our first set of results, which concern the communication complexity of agreement.

Section 3.1 studies the "economists' standard protocol," introduced by Geanakoplos and Polemarchakis [8] and alluded to earlier. In that protocol, Alice and Bob repeatedly announce their current expectations of a $[0,1]$ random variable, conditioned on all previous announcements. The question we ask is how many messages are needed before the agents' expectations agree within $\varepsilon$ with probability at least $1 - \delta$ over their prior, given parameters $\varepsilon$ and $\delta$. We show that order $1/\left(\delta \varepsilon^2\right)$ messages suffice. We then show that order $1/\left(\delta \varepsilon^2\right)$ messages still suffice, if instead of sending their whole expectations (which are real numbers), the agents send "summary" messages consisting of only 2 bits each. What makes these upper bounds surprising is that they are completely independent of $n$, the number of bits needed to represent the agents' knowledge. By contrast, in ordinary communication complexity (see [14]), it is easy to show that given a random function $f : \{0,1\}^n \times \{0,1\}^n \to [0,1]$, Alice and Bob would need to exchange order $n$ bits to approximate $f$ to within (say) $1/10$ with high probability.

Given the results of Section 3.1, several questions demand our attention. Is the upper bound of $1/\left(\delta \varepsilon^2\right)$ bits tight, or can it be improved even further? Also, is the economists' standard protocol always optimal, or do other protocols sometimes need even less communication? Section 3.2 addresses these questions. Though we are unable to show any lower bound better than $\log 1/\varepsilon$ that applies to *all* protocols, we do give examples where the standard protocol needs almost $1/\varepsilon^2$ bits. We also show that the standard protocol is not optimal: there exist cases where the standard protocol uses almost $1/\varepsilon^2$ bits, while a new protocol (which we call the *attenuated protocol*) uses fewer bits.

In earlier work, Parikh and Krasucki [15] extended Aumann's agreement theorem to three or more agents, who send messages along the edges of a directed graph. Thus, it is natural to ask whether our *efficient* agreement theorem extends to this setting as well. Section 3.3 shows that it does: given $N$ agents with a common prior, who send messages along a strongly connected graph of diameter $d$, order $Nd^2/\left(\delta \varepsilon^2\right)$ messages suffice for every pair of agents to agree within $\varepsilon$ about the expectation of a $[0,1]$ random variable with probability at least $1 - \delta$ over their prior.

In Section 4 we shift attention to the *computational* complexity of agreement, the subject of our technically most interesting result. What we want to show is that, even if two agents are computationally bounded, after a conversation of reasonable length they can still probably approximately agree about the expectation of a $[0,1]$ random variable. A large part of the problem is to say what this even means. After all, if the agents both ignored their evidence and estimated (say) $1/2$, then they would agree before exchanging a single message! So agreement is only interesting if the agents have made some sort of "good-faith effort" to emulate Bayesian rationality.

Although we leave unspecified exactly what effort is necessary, we do propose a criterion that we think is certainly *sufficient*. This is that the agents be able to *simulate* a Bayesian agreement protocol, in such a way that a computationally-unbounded referee, given the agents' knowledge together with a transcript of their conversation, be unable to decide (with non-negligible bias) whether the agents are computationally bounded or not. The justification for this criterion is that, just as Turing [18] argued that a perfect simulation of thinking *is* thinking, so it seems to us that a statistically perfect simulation of Bayesian rationality *is* Bayesian rationality.

But what do we mean by computationally-bounded agents? We discuss this question in detail in Section 4, but the basic point is that we assume two "subroutines": one that computes the $[0,1]$ variable of interest, given a state of the world $\omega$; and another that samples a state $\omega$ from any set in either agent's initial knowledge partition. The complexity of the simulation procedure is then expressed in terms of the number of calls to these subroutines.

Unfortunately, there is no way to simulate the economists' standard protocol—even our discretized version of it—using a small number of subroutine calls. The reason is that Alice's ideal estimate $p$ might lie on a "knife-edge" between the set of estimates that would cause her to send message $m_1$ to Bob, and the set that would cause her to send a different message $m_2$. In that case, it does not suffice for her to approximate $p$

using random sampling; she needs to determine it exactly. Our solution, which we develop in Section 4.1, is to have the agents "smooth" their messages by adding random noise to them. By hiding small errors in the agents' estimates, such noise makes the knife-edge problem disappear. On the other hand, we show that in the computationally-unbounded case, the noise does not prevent the agents from agreeing within $\varepsilon$ with probability $1 - \delta$ after order $1/\left(\delta\varepsilon^2\right)$ messages. In Sections 4.2 and 4.3 we prove the main result: that the smoothed standard protocol can be simulated using a number of subroutine calls that depends only on $\varepsilon$ and $\delta$, not on $n$. The dependence, alas, is exponential in $1/\left(\delta^3\varepsilon^6\right)$, so our simulation procedure is still not practical. However, we expect that both the procedure and its analysis can be considerably improved.

We conclude in Section 5 with some suggestions for future research, and some speculations about the causes of disagreement.

## 2    Preliminaries

Let $\Omega$ be a set of possible states of the world. Throughout this paper, $\Omega$ will be finite—both for simplicity of presentation, and because we do not believe that any physically realistic agent can ever have more than finitely many possible experiences. Let $\mathcal{D}$ be a prior probability distribution over $\Omega$ that is shared by some set of agents. We can assume $\mathcal{D}$ assigns nonzero probability to every $\omega \in \Omega$, for if not, we simply remove the probability-0 states from $\Omega$. Whenever we talk about a probability or expectation over a subset $S$ of $\Omega$, unless otherwise indicated we mean that we start from $\mathcal{D}$ and conditionalize on $\omega \in S$.

Throughout this paper, we will consider protocols in which agents send messages to each other in some order. Let $\Omega_{i,t}(\omega)$ be the set of states that agent $i$ considers possible immediately after the $t^{th}$ message has been sent, given that the true state of the world is $\omega$.[3] Then $\omega \in \Omega_{i,t}(\omega) \subseteq \Omega$, and indeed the set $\{\Omega_{i,t}(\omega)\}_{\omega \in \Omega}$ forms a partition of $\Omega$. Furthermore, since the agents never forget messages, we have $\Omega_{i,t}(\omega) \subseteq \Omega_{i,t-1}(\omega)$. Thus we say that the partition $\{\Omega_{i,t}\}_{\omega \in \Omega}$ *refines* $\{\Omega_{i,t-1}\}_{\omega \in \Omega}$, or equivalently that $\{\Omega_{i,t-1}\}_{\omega \in \Omega}$ *coarsens* $\{\Omega_{i,t}\}_{\omega \in \Omega}$. (As a convention, we freely omit arguments of $\omega$ when doing so will cause no confusion.) Notice also that if the $t^{th}$ message is not sent to agent $i$, then $\Omega_{i,t}(\omega) = \Omega_{i,t-1}(\omega)$.

Now let $f : \Omega \to [0,1]$ be a real-valued function that the agents are interested in estimating. The assumption $f(\omega) \in [0,1]$ is without loss of generality—for since $\Omega$ is finite, any function from $\Omega$ to $\mathbb{R}$ has a bounded range, which we can take to be $[0,1]$ by rescaling. We can think of $f(\omega)$ as the probability of some future event conditioned on $\omega$, but this is not necessary. Let $E_{i,t}(\omega) = \mathrm{EX}_{\omega' \in \Omega_{i,t}(\omega)}[f(\omega')]$ be agent $i$'s expectation of $f$ at step $t$, given that the true state of the world is $\omega$. Also, let $\Theta_{i,t}(\omega) = \{\omega' : E_{i,t}(\omega') = E_{i,t}(\omega)\}$ be the set of states for which agent $i$'s expectation of $f$ equals $E_{i,t}(\omega)$. Then the partition $\{\Theta_{i,t}\}_{\omega \in \Omega}$ coarsens $\{\Omega_{i,t}\}_{\omega \in \Omega}$, and $E_{i,t}(\omega) = \mathrm{EX}_{\omega' \in \Theta_{i,t}(\omega)}[f(\omega')]$.

The following simple but important fact is due to Hanson [11].

**Proposition 1 ([11])** *Suppose the partition* $\{\Omega_{i,t}\}_{\omega \in \Omega}$ *refines* $\{\Theta_{j,u}\}_{\omega \in \Omega}$. *Then*

$$\mathop{\mathrm{EX}}_{\omega' \in \Omega_{j,u}(\omega)}[E_{i,t}(\omega')] = \mathop{\mathrm{EX}}_{\omega' \in \Theta_{j,u}(\omega)}[E_{i,t}(\omega')] = E_{j,u}(\omega)$$

*for all* $\omega \in \Omega$. *As a consequence, an agent's expectation of its future expectation of $f$ always equals its current expectation. As another consequence, if Alice has just communicated her expectation of $f$ to Bob, then Alice's expectation of Bob's expectation of $f$ equals Alice's expectation.*

**Proof.** In each case, we are taking the expectation of $f$ over a subset $S \subseteq \Omega$ (either $\Omega_{j,u}(\omega)$ or $\Theta_{j,u}(\omega)$) for which $\mathrm{EX}_{\omega' \in S}[f(\omega')] = E_{j,u}(\omega)$. How $S$ is "sliced up" has no effect on the result. ∎

Proposition 1 already demonstrates a dramatic difference between Bayesian agreement protocols and actual human conversations. Suppose Alice and Bob are discussing whether useful quantum computers will be built by the year 2050. Bob says that, in his opinion, the chance of this happening is only 5%. Alice says she disagrees: she thinks the chance is 90%. How much should Alice expect her reply to influence Bob's estimate? Should she expect him to raise it to 10%, or even 15%, out of deference to his friend Alice's judgment? According to Proposition 1, she should expect him to raise it to 90%! That is, depending on what else Bob knows, his new estimate might be 85% or 95%, but its *expectation* from Alice's point of view is 90%.

---

[3]We assume for now that messages are "noise-free"; that is, they partition the state space sharply. Later we will remove this assumption.
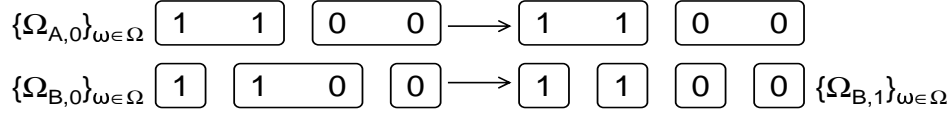
Figure 1: After Alice tells Bob whether $E_{A,0}$ is 1 or 0, Bob's partition $\{\Omega_{B,0}\}_{\omega \in \Omega}$ is refined to $\{\Omega_{B,1}\}_{\omega \in \Omega}$.

## 2.1 Miscellany

Asymptotic notation is standard: $F(n) = O(G(n))$ means there exist positive constants $a, b$ such that $F(n) \leq a + bG(n)$ for all $n \geq 0$; $F(n) = \Omega(G(n))$ means the same but with $F(n) \geq a + bG(n)$; $F(n) = \Theta(G(n))$ means $F(n) = O(G(n))$ and $F(n) = \Omega(G(n))$; and $F(n) = o(G(n))$ means $F(n) = O(G(n))$ and not $F(n) = \Omega(G(n))$.

We will have several occasions to use the following well-known bound.

**Theorem 2 (Chernoff, Hoeffding)** *Let $x_1, \ldots, x_K$ be $K$ independent samples of a $[0, 1]$ random variable with mean $\mu$. Then for all $\alpha \in (0, 1)$,*

$$\Pr[x_1 + \cdots + x_K \leq (1 - \alpha)\mu K] \leq e^{-\mu \alpha^2 K / 2},$$
$$\Pr[|x_1 + \cdots + x_K - \mu K| > \alpha K] \leq 2e^{-2\alpha^2 K}.$$

# 3 Communication Complexity

We now introduce and justify the communication complexity model. Assume for the moment that there are two agents, Alice ($A$) and Bob ($B$); Section 3.3 will generalize the model to three or more agents. We can imagine if we like that Alice and Bob are given $n$-bit strings $x$ and $y$ respectively, so that $\Omega \subseteq \{0, 1\}^n \times \{0, 1\}^n$. Letting $\omega = (x, y)$, we then have $\Omega_{A,0}(\omega) \subseteq x \times \{0, 1\}^n$ and $\Omega_{B,0}(\omega) \subseteq \{0, 1\}^n \times y$.

In an *agreement protocol*, Alice and Bob take turns sending messages to each other. Any such protocol is characterized by a sequence of functions $m_1, m_2, \ldots : 2^\Omega \to \mathcal{M}$, known to both agents, which map subsets of $\Omega$ to elements of a message space $\mathcal{M}$. Possibilities for $\mathcal{M}$ include $[0, 1]$ in a continuous protocol, or $\{0, 1\}$ in a discretized protocol. In all protocols considered in this paper, the $m_t$'s will be extremely simple; for example, we might have $m_t(S) = \mathrm{EX}_{\omega' \in S}[f(\omega')]$ be the agent's current expectation of $f$.

The protocol proceeds as follows: first Alice computes $m_1(\Omega_{A,0}(\omega))$ and sends it to Bob. After seeing Alice's message, and assuming the true state of the world is $\omega$, Bob's new set of possible states becomes

$$\Omega_{B,1}(\omega) = \Omega_{B,0}(\omega) \cap \{\omega' : m_1(\Omega_{A,0}(\omega')) = m_1(\Omega_{A,0}(\omega))\}$$

as in Figure 1. Then Bob computes $m_2(\Omega_{B,1}(\omega))$ and sends it to Alice, whereupon Alice's set of possible states becomes

$$\Omega_{A,2}(\omega) = \Omega_{A,0}(\omega) \cap \{\omega' : m_2(\Omega_{B,1}(\omega')) = m_2(\Omega_{B,1}(\omega))\}.$$

Then Alice computes $m_3(\Omega_{A,2}(\omega))$ and sends it to Bob, and so on.

At this point we should address an obvious question: how do Alice and Bob know each other's initial partitions, $\{\Omega_{A,0}\}_{\omega \in \Omega}$ and $\{\Omega_{B,0}\}_{\omega \in \Omega}$? If the agents do not know each other's partitions, then messages between them are useless, since neither agent knows how to update its own partition based on the other's messages. This question is not specific to our setting; it can be asked about Aumann's original result as well as any of its extensions. The solution in each case is that the state of the world $\omega \in \Omega$ *includes the agents' mental states as part of it.* From this it follows that every agent has a uniquely defined partition known to every other agent. For suppose Alice calculates that if the state of the world is $\omega$, then Bob's knowledge is $\Omega_{B,0}(\omega)$, meaning that he knows (and knows only) that the state belongs to $\Omega_{B,0}(\omega)$. Then for all $\omega' \in \Omega_{B,0}(\omega)$, she must calculate that if the state is $\omega'$, then Bob's knowledge is $\Omega_{B,0}(\omega)$ as well. Otherwise one of her calculations was mistaken.

The reader might object on the following grounds. Suppose Alice and Bob are the only two agents, and let $\Omega^{(0)}$ be the set of possible states of the "external" world—meaning everything except Alice and Bob. Next let $\Omega^{(1)}$ be the set of possible states of the agents' knowledge regarding $\Omega^{(0)}$, let $\Omega^{(2)}$ be the set of possible states of their knowledge regarding $\Omega^{(1)}$, and so on. Then $\Omega = \Omega^{(0)} \times \Omega^{(1)} \times \Omega^{(2)} \times \cdots$, which contradicts the assumption that $\Omega$ is finite. The obvious response is that, since the agents' brains can store only finitely many bits, not all elements of $\Omega^{(0)} \times \Omega^{(1)} \times \Omega^{(2)} \times \cdots$ are actually possible.

However, the above response is open to a different objection, related to the diagonalization arguments of Gödel and Turing. Suppose Alice's and Bob's brains store $n$ bits each. Then in order to reason about the set of possible states of their brains, wouldn't they need brains that store more than $n$ bits? We leave this conundrum unresolved, confining ourselves to the following three remarks. First, only a tiny portion of the agents' brains is likely to be relevant to their topic of conversation, which means "plenty of room left over" for metareasoning about knowledge. Second, by reducing the number of brain states that the agents need to consider, our results in Section 4 will lessen the force of the self-reference argument, though not eliminate it. Third, the agents' "knowledge hierarchy" seems likely to collapse at a low level. That is, Alice might have little idea what sort of evidence shaped Bob's opinions about the external world. But Bob probably has some idea what sort of evidence shaped Alice's opinions about Bob's opinions, and Alice probably has a good idea what sort of evidence shaped Bob's opinions about Alice's opinions about Bob's opinions (assuming Bob even *has* nontrivial such opinions). The more indirect the knowledge, the fewer the ways of obtaining it.

Let us return to explaining the communication complexity model. After the $t^{th}$ message, we say Alice and Bob $(\varepsilon, \delta)$-*agree* if their expectations of $f$ agree to within $\varepsilon$ with probability at least $1 - \delta$; that is, if

$$\Pr_{\omega \in \mathcal{D}} \left[ |E_{A,t}(\omega) - E_{B,t}(\omega)| > \varepsilon \right] \leq \delta.$$

The goal will be to minimize the number of messages until the agents $(\varepsilon, \delta)$-agree.

In our view, $(\varepsilon, \delta)$-agreement is a much more fundamental notion than exact agreement. For suppose $f$ represents the probability that global warming, if left unchecked, will cause sea levels to rise at least 30 centimeters by the year 2100. If after an hour's conversation, any two people could agree within $1/4$ about $f$ with probability at least $3/4$, then the world would be a remarkably different place than it now is. That the agreement was inexact and uncertain would be less significant than the fact that it occurred at all.

But why do we calculate the success probability over $\mathcal{D}$, and not some other distribution? In other words, what if the agents' priors agree with each other, but not with external reality? Unfortunately, in that case it seems difficult to prove anything, since the "true" prior could be concentrated on a few states that the agents consider vanishingly unlikely. Furthermore, we conjecture that there exist $f, \mathcal{D}$ such that for all agreement protocols, the agents must exchange $\Omega(n)$ bits to agree within $\varepsilon$ on *every* state $\omega$ (that is, to $(\varepsilon, 0)$-agree). So given a protocol that causes Alice and Bob to $(\varepsilon, \delta)$-agree, what we should really say is that both agents enter the conversation *expecting* to agree within $\varepsilon$ with probability at least $1 - \delta$. This, of course, is profoundly unlike the situation in real life, where adversaries generally do not enter arguments expecting to convince or to be convinced.

Let us make two further remarks about the model. First, if the agents want to agree exactly (that is, $(0, 0)$-agree), it is clear that in the worst case they need $2n$ bits of communication, $n$ from Alice and $n$ from Bob. Note the contrast with ordinary communication complexity, where $n$ bits always suffice. Indeed, even to produce approximate agreement, two-way communication is necessary in general, as shown by the example $f(x, y) = (2x + y)/3$, where $x, y \in \{0, 1\}$ are uniformly distributed.

Second, our ending condition is simply that the agents $(\varepsilon, \delta)$-agree at *some* step $t$. We do not require them to fix this $t$ independently of $f$ and $\mathcal{D}$. The reason is that for any $t$, there might exist perverse $f, \mathcal{D}$ such that the agents nearly agree for the first $t - 1$ steps, then disagree violently at the $t^{th}$ step. However, it seems unfair to penalize the agents in such cases.

The following is the best lower bound we are able to show on agreement complexity.

**Proposition 3** *There exist $f, \mathcal{D}$ such that for all $\varepsilon \geq 2^{-n}$ and $\delta \geq 0$, Alice must send $\Omega\left(\log \frac{1-\delta}{\varepsilon}\right)$ bits to Bob and Bob must send $\Omega\left(\log \frac{1-\delta}{\varepsilon}\right)$ bits to Alice before the agents $(\varepsilon, \delta)$-agree. In particular, if $\delta$ is bounded away from 1 by a constant, then $\Omega(\log 1/\varepsilon)$ bits are needed.*

**Proof.** Let $\Omega = \{1, \ldots, 2^n\}^2$, let $\mathcal{D}$ be uniform over $\Omega$, and let $f(x, y) = (x + y)/2^{n+1}$ for all $(x, y) \in \Omega$. Thus if $\widehat{x}$ is Bob's expectation of $x$ at step $t$ and $\widehat{y}$ is Alice's expectation of $y$, then $E_{A,t} = (x + \widehat{y})/2^{n+1}$

6

and $E_{B,t} = (\widehat{x} + y) / 2^{n+1}$. Suppose one agent, say Alice, has sent only $t < \log_2 \left( \frac{1-\delta}{\varepsilon} \right) - 2$ bits to Bob. For each $i \in \{1, \dots, 2^t\}$, let $p_i$ be the probability of the $i^{th}$ message sequence from Alice. Conditioned on $i$, there are $2^n p_i$ values of $x$ still possible from Bob's point of view. So regardless of $E_{B,t}$, the probability of $|E_{A,t} - E_{B,t}| \leq \varepsilon$ can be at most $4\varepsilon / p_i$. Therefore the agents agree within $\varepsilon$ with total probability at most

$$\sum_{i=1}^{2^t} p_i \left( \frac{4\varepsilon}{p_i} \right) = 4\varepsilon 2^t < 1 - \delta.$$

∎

## 3.1 Convergence of the Standard Protocol

The two-player "standard protocol" is simply the following: first Alice sends $E_{A,0}$, her current expectation of $f$, to Bob. Then Bob sends his expectation $E_{B,1}$ to Alice, then Alice sends $E_{A,2}$ to Bob, and so on. Geanakoplos and Polemarchakis [8] observed that for any $f, \mathcal{D}$, if the agents use the standard protocol then after a finite number of messages $T$, they will reach *consensus*—meaning that $E_{A,T} = E_{B,T}$, both agents know this, both know that they know it, etc. In particular, in our terminology Alice and Bob $(0, 0)$-agree.

In this section we ask how many messages are needed before the agents $(\varepsilon, \delta)$-agree. The surprising and unexpected answer, in Theorem 5, is that $1/ \left( \delta \varepsilon^2 \right)$ messages always suffice, independently of $n$ and all other parameters of $f$ and $\mathcal{D}$. One might guess that, since the expectations $E_{A,0}, E_{B,1}, \dots$ are real numbers, the cost of communication must be hidden in the length of the messages. However, in Theorem 6 we show even if the agents send only 2-bit "summaries" of their expectations, $O \left( 1/ \left( \delta \varepsilon^2 \right) \right)$ messages still suffice for $(\varepsilon, \delta)$-agreement.

Given any function $F : \Omega \to [0, 1]$, let $\|F\|_2^2 = \text{EX}_{\omega \in \mathcal{D}} \left[ F(\omega)^2 \right]$. The following proposition will be used again and again in this paper.

**Proposition 4** *Suppose the partition $\{\Omega_{i,t}\}_{\omega \in \Omega}$ refines $\{\Theta_{j,u}\}_{\omega \in \Omega}$. Then*

$$\|E_{i,t}\|_2^2 - \|E_{j,u}\|_2^2 = \|E_{i,t} - E_{j,u}\|_2^2$$

*so in particular, $\|E_{i,t}\|_2^2 \geq \|E_{j,u}\|_2^2$. A special case is that $\|E_{i,t+1}\|_2^2 \geq \|E_{i,t}\|_2^2$ for all $i, t$.*

**Proof.** We have

$$\text{EX}\left[E_{i,t} E_{j,u}\right] = \underset{\omega \in \mathcal{D}}{\text{EX}} \left[ E_{j,u}(\omega) \underset{\omega' \in \Theta_{j,u}(\omega)}{\text{EX}} [E_{i,t}(\omega')] \right] = \underset{\omega \in \mathcal{D}}{\text{EX}} [E_{j,u}(\omega) \cdot E_{j,u}(\omega)] = \|E_{j,u}\|_2^2$$

by Proposition 1, and therefore

$$\|E_{i,t} - E_{j,u}\|_2^2 = \|E_{i,t}\|_2^2 + \|E_{j,u}\|_2^2 - 2 \, \text{EX} [E_{i,t} E_{j,u}] = \|E_{i,t}\|_2^2 - \|E_{j,u}\|_2^2.$$

∎

We can now prove an upper bound on the number of messages needed for agreement.

**Theorem 5** *For all $f, \mathcal{D}$, the standard protocol causes Alice and Bob to $(\varepsilon, \delta)$-agree after at most $1/ \left( \delta \varepsilon^2 \right)$ messages.*

**Proof.** Intuitively, so long as the agents disagree by more than $\varepsilon$ with high probability, Alice's expectation $E_{A,1}, E_{A,2}, \dots$ follows an unbiased random walk with step size roughly $\varepsilon$. Furthermore, this walk has two absorbing barriers at 0 and 1, for the simple fact that $E_{A,t} \in [0, 1]$. And we expect a random walk with step size $\varepsilon$ to hit a barrier after about $1/\varepsilon^2$ steps.

To make this intuition precise, we need only track the expectation, not of $E_A$ and $E_B$, but of $E_A^2$ and $E_B^2$. Suppose Alice sends the $t^{th}$ message. Then Bob's partition $\{\Omega_{B,t}\}_{\omega \in \Omega}$ refines $\{\Theta_{A,t-1}\}_{\omega \in \Omega}$. It follows by Proposition 4 that

$$\|E_{B,t}\|_2^2 - \|E_{A,t-1}\|_2^2 = \|E_{B,t} - E_{A,t-1}\|_2^2.$$

7

Assuming $\Pr\left[|E_{B,t} - E_{A,t-1}| > \varepsilon\right] \geq \delta$, this implies that $\|E_{B,t}\|_2^2 > \|E_{A,t-1}\|_2^2 + \delta\varepsilon^2$. Similarly, after Bob sends Alice the $(t+1)^{st}$ message, we have $\|E_{A,t+1}\|_2^2 > \|E_{B,t}\|_2^2 + \delta\varepsilon^2$. So until the agents $(\varepsilon, \delta)$-agree, each message increases $\max\left\{\|E_{A,t}\|_2^2, \|E_{B,t}\|_2^2\right\}$ by more than $\delta\varepsilon^2$. But the maximum can never exceed 1 (since $E_{A,t}, E_{B,t} \in [0,1]$), which yields an upper bound of $1/\left(\delta\varepsilon^2\right)$ on the number of messages. $\blacksquare$

As mentioned previously, the trouble with the standard protocol is that sending one's expectation might require too many bits. A simple way to discretize the protocol is as follows. Imagine a "monkey in the middle," Charlie, who has the same prior distribution $\mathcal{D}$ as Alice and Bob and who sees all messages between them, but who does not know either of their inputs. In other words, letting $\Omega_{C,t}(\omega)$ be the set of states that Charlie considers possible after the first $t$ messages, we have $\Omega_{C,0}(\omega) = \Omega$ for all $\omega$. Then the partition $\{\Omega_{C,t}\}_{\omega \in \Omega}$ coarsens both $\{\Omega_{A,t}\}_{\omega \in \Omega}$ and $\{\Omega_{B,t}\}_{\omega \in \Omega}$; therefore both Alice and Bob can compute Charlie's expectation $E_{C,t}(\omega) = \mathrm{EX}_{\omega' \in \Omega_{C,t}(\omega)}[f(\omega')]$ of $f$.

Now whenever it is her turn to send a message to Bob, Alice sends the message "high" if $E_{A,t} > E_{C,t} + \varepsilon/4$, "low" if $E_{A,t} < E_{C,t} - \varepsilon/4$, and "medium" otherwise. This requires 2 bits. Likewise, Bob sends "high" if $E_{B,t} > E_{C,t} + \varepsilon/4$, "low" if $E_{B,t} < E_{C,t} - \varepsilon/4$, and "medium" otherwise.

**Theorem 6** *For all $f, \mathcal{D}$, the discretized protocol described above causes Alice and Bob to $(\varepsilon, \delta)$-agree after $O\left(1/\left(\delta\varepsilon^2\right)\right)$ messages.*

**Proof.** The plan is to show that either $\|E_{A,t}\|_2^2$, $\|E_{B,t}\|_2^2$, or $\|E_{C,t}\|_2^2$ increases by at least $\delta\varepsilon^2/512$ with every message of Alice's, until Alice and Bob $(\varepsilon, \delta)$-agree. Since $\|E_{i,t}\|_2^2 \leq 1$ for all $i$, this will imply an upper bound of $3072/\left(\delta\varepsilon^2\right)$ on the number of messages (we did not optimize the constant!).

Assume that $\Pr\left[|E_{A,t} - E_{B,t}| > \varepsilon\right] \geq \delta$ and it is Alice's turn to send the $(t+1)^{st}$ message. By the triangle inequality, either

$$\Pr\left[|E_{A,t} - E_{C,t}| > \frac{\varepsilon}{2}\right] \geq \frac{\delta}{2}$$

or

$$\Pr\left[|E_{B,t} - E_{C,t}| > \frac{\varepsilon}{2}\right] \geq \frac{\delta}{2}.$$

We analyze these two cases separately. In the first case, with probability at least $\delta/2$ Alice's message is either "high" or "low." If the message is "high," then $E_{C,t+1}$ becomes an average of numbers each greater than $E_{C,t} + \varepsilon/4$, so $E_{C,t+1} > E_{C,t} + \varepsilon/4$. If the message is "low," then likewise $E_{C,t+1} < E_{C,t} - \varepsilon/4$. Since $\{\Omega_{C,t+1}\}_{\omega \in \Omega}$ refines $\{\Omega_{C,t}\}_{\omega \in \Omega}$, Proposition 4 thereby gives

$$\|E_{C,t+1}\|_2^2 - \|E_{C,t}\|_2^2 = \|E_{C,t+1} - E_{C,t}\|_2^2 > \frac{\delta}{2}\left(\frac{\varepsilon}{4}\right)^2.$$

Now for the second case. If, after Alice sends the $(t+1)^{st}$ message, we still have

$$\Pr\left[|E_{B,t+1} - E_{C,t+1}| > \frac{\varepsilon}{4}\right] \geq \frac{\delta}{4},$$

then the previous argument applied to Bob implies that

$$\|E_{C,t+2}\|_2^2 - \|E_{C,t+1}\|_2^2 > \frac{\delta}{4}\left(\frac{\varepsilon}{4}\right)^2$$

and we are done. So suppose otherwise. Then the difference between Bob's and Charlie's expectations must have changed significantly:

$$\Pr\left[|E_{B,t} - E_{C,t}| - |E_{B,t+1} - E_{C,t+1}| > \frac{\varepsilon}{4}\right] > \frac{\delta}{4}.$$

Hence by another application of the triangle inequality, either

$$\Pr\left[|E_{B,t+1} - E_{B,t}| > \frac{\varepsilon}{8}\right] > \frac{\delta}{8}$$

8

or

$$\Pr\left[|E_{C,t+1} - E_{C,t}| > \frac{\varepsilon}{8}\right] > \frac{\delta}{8}.$$

In the former case, Proposition 4 yields

$$\|E_{B,t+1}\|_2^2 - \|E_{B,t}\|_2^2 = \|E_{B,t+1} - E_{B,t}\|_2^2 > \frac{\delta}{8}\left(\frac{\varepsilon}{8}\right)^2,$$

while in the latter case,

$$\|E_{C,t+1}\|_2^2 - \|E_{C,t}\|_2^2 > \frac{\delta}{8}\left(\frac{\varepsilon}{8}\right)^2.$$

∎

## 3.2 Attenuated Protocol

We have seen that two agents, using the standard protocol, will always $(\varepsilon, \delta)$-agree after exchanging only $O\left(1/\left(\delta\varepsilon^2\right)\right)$ messages. This result immediately raises three questions:

(1) Is there a scenario where the standard protocol *needs* about $1/\varepsilon^2$ messages to produce $(\varepsilon, \delta)$-agreement?

(2) Is the standard protocol always optimal, or do other protocols sometimes outperform it?

(3) Is there a scenario where *any* agreement protocol needs a number of communication bits polynomial in $1/\varepsilon$?

Although we leave question (3) open, in this section we resolve questions (1) and (2). In particular, assume for simplicity that $\delta = 1/2$. Then for all $\varepsilon > 0$, Theorem 7 gives a scenario where the standard protocol uses almost $1/\varepsilon^2$ messages, even if the messages are continuous rather discrete. By contrast, a new "attenuated protocol" uses only 2 messages, both consisting of a constant number of bits (independent of $\varepsilon$). Theorem 8 then gives a *fixed* scenario where for all $\varepsilon > 0$, the standard protocol uses almost $1/\varepsilon^2$ messages, while the attenuated protocol uses only 2 messages, both consisting of $O\left(1/\varepsilon\right)$ bits.

The attenuated protocol is interesting in its own right. The idea is to imagine that in the standard protocol, the communication channel between Alice and Bob becomes gradually more noisy as time goes on, so that each message conveys slightly less information than the one before. It turns out that in some cases, such noise would actually help! For intuitively, each time the message intensity decreases by $\epsilon$, the "price" the agents pay in terms of disagreement is proportional to $\epsilon^2$. So it is better for them to attenuate their conversation gradually, than to send a sequence of "maximum-intensity" messages followed by no message (which we can think of as intensity 0).[4] Even if the noise that produces this strange effect is missing from the channel, the agents can easily simulate it. Furthermore, the messages will turn out to be nonadaptive, so they can all be concatenated into one message from Alice and one from Bob.

But how do we ensure that the standard protocol needs almost $1/\varepsilon^2$ messages? Intuitively, by forcing the random walk behavior of Section 3.1 actually to occur. That is, at the beginning there will be a disagreement that can only be resolved by Alice sending a bit to Bob. But then that bit will cause a new disagreement even as it resolves the old one, and so on.

**Theorem 7** *For all $\varepsilon > 0$, there exist $f, \mathcal{D}$ such that for all $\delta > 0$:*

(i) *Using the standard protocol, Alice and Bob need to exchange $\Omega\left(\frac{1}{\varepsilon^2 \log \frac{2}{(1-\delta)\varepsilon}}\right)$ messages before they $(\varepsilon, \delta)$-agree.*

(ii) *Using a different protocol, they need only exchange 2 messages, both consisting of $O\left(\log 1/\delta\right)$ bits.*

*In particular, if $\delta = 1/2$ then the standard protocol needs $\Omega\left(\frac{1/\varepsilon^2}{\log 1/\varepsilon}\right)$ bits whereas the attenuated protocol needs $O\left(1\right)$ bits.*

---

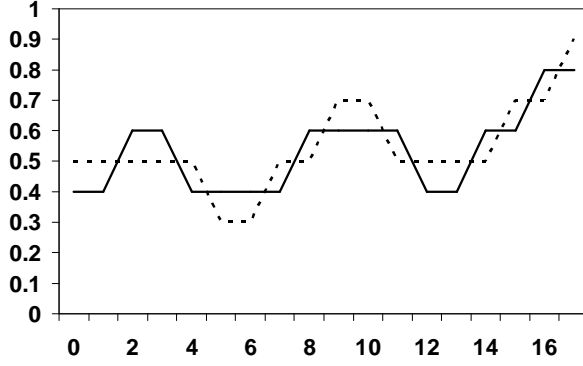[4]The same phenomenon occurs in the "Zeno effect" of quantum mechanics .

Figure 2: Alice's expectation $E_{A,t}$ (solid line), and Bob's expectation $E_{B,t}$ (dashed line), as a function of $t$

**Proof.** Let

$$n = \frac{1}{64\varepsilon^2 \ln \frac{6}{(1-\delta)\varepsilon^2}}$$

(throughout we omit floor and ceiling signs for convenience). The state space $\Omega$ consists of all pairs $(x, y)$, where $x = x_1 \ldots x_n$ and $y = y_1 \ldots y_n$ belong to $\{-1, 1\}^n$. The prior distribution $\mathcal{D}$ is uniform over $\Omega$. Let

$$F(x, y) = \frac{1}{2} + 2\varepsilon \sum_{i=1}^{n} (y_{i-1} x_i + x_i y_i)$$

where $y_0 = 1$. Then the function that interests the agents is

$$f(x, y) = \begin{cases} F(x, y) & \text{if } F(x, y) \in [0, 1] \\ 0 & \text{if } F(x, y) < 0 \\ 1 & \text{if } F(x, y) > 1 \end{cases}.$$

For simplicity, we first consider $F$ (which need not be bounded in $[0, 1]$), and later analyze the "edge effects" that arise in switching to $f$. We claim that, if the agents use the continuous standard protocol to evaluate $F$, then $|E_{A,t} - E_{B,t}| = 2\varepsilon$ at all steps $t < 2n$, where $E_{A,t}$ and $E_{B,t}$ are Alice's and Bob's expectations of $F$ respectively after $t$ messages have been exchanged. For initially $E_{A,0} = 1/2 + 2\varepsilon x_1$ and $E_{B,0} = 1/2$. Most of the terms in the sum defining $F(x, y)$ simply average to 0 for both agents, since Alice does not know the $y_i$'s and Bob does not know the $x_i$'s. In the first step, however, the expectation that Alice sends to Bob reveals $x_1$ to him. This causes $E_{B,1}$ to become $1/2 + 2\varepsilon x_1 + 2\varepsilon x_1 y_1$, which differs from $E_{A,0} = 1/2 + 2\varepsilon x_1$ by $2\varepsilon$. Then in the second step, the expectation that Bob sends to Alice reveals $y_1$ to her, thereby "unlocking" the terms $x_1 y_1$ and $y_1 x_2$ in her expectation, and so on. It follows that until all $2n$ bits $x_1 \ldots x_n$ and $y_1 \ldots y_n$ have been exchanged, the agents disagree by $2\varepsilon$ with certainty (see Figure 2).

In switching from $F$ to $f$, the key observation is that Alice's expectation $E_{A,t}(f)$ of $f$ is a function of her expectation

$$E_{A,t} = \frac{1}{2} + 2\varepsilon \left( x_1 + x_1 y_1 + y_1 x_2 + \cdots + x_{(t-1)/2} y_{(t-1)/2} + y_{(t-1)/2} x_{(t+1)/2} \right)$$

of $F$. For from Alice's point of view, the later terms $x_{(t+1)/2} y_{(t+1)/2}$, $y_{(t+1)/2} x_{(t+3)/2}$, and so on are steps in an unbiased random walk with starting point $E_{A,t}$, step size $2\varepsilon$, and "snapping barriers" at 0 and 1. (A snapping barrier is neither absorbing nor reflecting: it allows a particle through, but if the particle is found on the wrong side of the barrier after the walk ends, then the particle is moved back to the barrier.) Let $E_{A,t}^*$ be the ending point of this walk; then $E_{A,t}(f) = \text{EX}\left[E_{A,t}^*\right]$ is a function of $E_{A,t}$. Likewise, $E_{B,t}(f) = \text{EX}\left[E_{B,t}^*\right]$ is the expected ending point of an unbiased walk with starting point $E_{B,t} = E_{A,t} + 2\varepsilon x_{(t+1)/2} y_{(t+1)/2}$, step size $2\varepsilon$, and snapping barriers at 0 and 1.

10

The lower bound for the standard protocol now follows from two claims: first, that $E_{A,t} \in [1/4, 3/4]$ and $E_{B,t} \in [1/4, 3/4]$ for all $t \in \{0, \ldots, 2n\}$ with probability at least $\delta$. Second, that whenever $E_{A,t}$ and $E_{B,t}$ belong to $[1/4, 3/4]$, we have $|E_{A,t}(f) - E_{B,t}(f)| > \varepsilon$. For the first claim, choose $z_1, \ldots, z_{2n}$ uniformly and independently from $\{0, 1\}$; then Theorem 2 says that

$$\Pr\left[|z_1 + \cdots + z_{2n} - n| > \alpha(2n)\right] \leq 2e^{-4\alpha^2 n}.$$

Setting $\alpha = \frac{1/4}{2\varepsilon(2n)}$, this implies that for any fixed $t$,

$$\Pr\left[|E_{A,t} - 1/4| > 1/4\right] \leq 2e^{-1/(64\varepsilon^2 n)} \leq \frac{1-\delta}{2n}$$

and similarly for $E_{B,t}$. The claim now follows from the union bound. For the second claim, a bound similar to the above implies that

$$\Pr\left[|E_{A,t}^* - E_{A,t}| > 1/4\right] \leq 2e^{-1/(64\varepsilon^2 n)} \leq \frac{\varepsilon}{3}$$

and similarly for $E_{B,t}^*$. This in turn implies that $|E_{A,t}(f) - E_{A,t}| \leq \varepsilon/3$ and $|E_{B,t}(f) - E_{B,t}| \leq \varepsilon/3$, from whence it follows that $|E_{A,t}(f) - E_{B,t}(f)| > \varepsilon$ by the triangle inequality.

We now give the $O(\log 1/\delta)$ upper bound. It suffices to give a protocol for $F$, since it is not hard to see that switching from $F$ to $f$ can only decrease $|E_{A,t} - E_{B,t}|$. Let $k = 8 \ln 2/\delta$. For each $i \in \{1, \ldots, k\}$, Alice sends Bob a bit that is uniformly random with probability $i/k$ and $x_i$ otherwise. Likewise, Bob sends Alice a bit that is uniformly random with probability $i/k$ and $y_i$ otherwise. Then Alice's final expectation is

$$E_{A,2} = \frac{1}{2} + 2\varepsilon \sum_{i=1}^{k} \left(\frac{i-1}{k} y_{i-1} x_i + \frac{i}{k} x_i y_i\right)$$

while Bob's is

$$E_{B,2} = \frac{1}{2} + 2\varepsilon \sum_{i=1}^{k} \left(\frac{i}{k} y_{i-1} x_i + \frac{i}{k} x_i y_i\right).$$

So

$$E_{B,2} - E_{A,2} = \frac{2\varepsilon}{k} \sum_{i=1}^{k} y_{i-1} x_i,$$

and hence

$$\Pr\left[|E_{A,2} - E_{B,2}| > \varepsilon\right] = \Pr\left[|z_1 + \cdots + z_k - k/2| > k/4\right]$$

where $z_i = (y_{i-1} x_i + 1)/2$. Since the $z_i$'s are uniform, independent samples from $\{0, 1\}$, the above probability is at most $2e^{-2(1/4)^2 k} = \delta$ by Theorem 2. ∎

The main defect of Theorem 7 is that the function $f$ had to be tailored to a particular $\varepsilon$. The next theorem fixes this defect, although the advantage of the attenuated protocol over the standard one is not quite as dramatic as in Theorem 7. For simplicity, in stating the theorem we fix $\delta = 1/2$.

**Theorem 8** *For all $\gamma \in (0, 1)$, there exist $f, \mathcal{D}$ such that for all $\varepsilon \geq 1/n^{1/(2-\gamma)}$:*

  *(i) Using the standard protocol, Alice and Bob need to exchange $\Omega\left(1/\varepsilon^{2-\gamma}\right)$ messages before they $(\varepsilon, 1/2)$-agree.*

  *(ii) Using the attenuated protocol, they need only exchange $2$ messages, both consisting of $O(1/\varepsilon)$ bits.*

  **Sketch.** Again we let $\mathcal{D}$ be uniform over $x = x_1 \ldots x_n$ and $y = y_1 \ldots y_n$ in $\{-1, 1\}^n$. We then let

$$F(x, y) = \frac{1}{2} + \frac{\sqrt{\gamma}}{10} \sum_{i=1}^{n} \frac{y_{i-1} x_i + x_i y_i}{i^{1/(2-\gamma)}}$$

and

$$f(x, y) = \begin{cases} F(x, y) & \text{if } F(x, y) \in [0, 1] \\ 0 & \text{if } F(x, y) < 0 \\ 1 & \text{if } F(x, y) > 1 \end{cases}.$$

The rest of the proof is almost identical to that of Theorem 7, so we omit it here. ∎
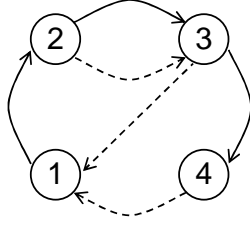
Figure 3: For a sample graph $G$, spanning tree $\mathcal{T}_1$ is shown in solid lines, and $\mathcal{T}_2$ in dashed lines.

## 3.3 $N$ Agents

We have seen that two Bayesian agents can reach rapid agreement, provided they communicate directly with each other. An obvious followup question is, what if there are three or more agents, each of which talks only to its 'neighbors'? Will the agents still reach agreement, and if so, after how long?

Formally, let $G$ be a directed graph with vertices $1, \dots, N$, each representing an agent. Suppose messages can only be sent from agent $i$ to agent $j$ if $(i, j)$ is an edge in $G$. We need to assume $G$ is strongly connected, since otherwise reaching agreement could be impossible for trivial reasons. In this setting, a *standard protocol* consists of a sequence of edges $(i_1, j_1), \dots, (i_t, j_t), \dots$ of $G$. At the $t^{th}$ step, agent $i_t$ sends its current expectation $E_{i_t, t-1}$ of $f$ to agent $j_t$, whereupon $j_t$ updates its expectation accordingly. Call the protocol *fair* if every edge occurs infinitely often in the sequence. Parikh and Krasucki [15] proved the following important theorem.

**Theorem 9 ([15])** *For all $f, \mathcal{D}$, any fair protocol will cause all the agents' expectations to agree after a finite number of messages.*

Indeed, the agents will reach *consensus* after finitely many messages, meaning it will be common knowledge among them that $E_{1,t} = \cdots = E_{N,t}$. Here, though, we care only about the weaker condition of agreement.

Our goal is to cause every pair of agents to $(\varepsilon, \delta)$-agree,[5] after a number of steps polynomial in $N$, $1/\delta$, and $1/\varepsilon$. We can achieve this via the following "spanning-tree protocol." Let $\mathcal{T}_1$ and $\mathcal{T}_2$ be two spanning trees of $G$ of minimum diameter, both rooted at agent 1. As illustrated in Figure 3, $\mathcal{T}_1$ points outward from 1 to the other $N-1$ agents; $\mathcal{T}_2$ points inward back to 1. Let $\mathcal{O}_1$ be an ordering of the edges of $\mathcal{T}_1$, in which every edge originating at $i$ is preceded by an edge terminating at $i$, unless $i = 1$. Likewise let $\mathcal{O}_2$ be an ordering of the edges of $\mathcal{T}_2$, in which every edge originating at $i$ is preceded by an edge terminating at $i$, unless $i$ is a leaf of $\mathcal{T}_2$. Then the protocol is simply for agents to send their current expectations along edges of $G$ in the order $\mathcal{O}_1, \mathcal{O}_2, \mathcal{O}_1, \mathcal{O}_2, \dots$.

**Theorem 10** *For all $f, \mathcal{D}$, the spanning-tree protocol causes every pair of agents to $(\varepsilon, \delta)$-agree after $O\left(\frac{Nd^2}{\delta \varepsilon^2}\right)$ messages, where $d$ is the diameter of $G$.*

**Proof.** We will track $\eta_t = \min_i \|E_{i,t}\|_2^2$. Observe that, if the $t^{th}$ message is from agent $i$ to agent $j$, then the partition $\{\Omega_{j,t+1}\}_{\omega \in \Omega}$ refines both $\{\Theta_{i,t}\}_{\omega \in \Omega}$ and $\{\Omega_{j,t}\}_{\omega \in \Omega}$, and therefore

$$\|E_{j,t+1}\|_2^2 \geq \max\left\{ \|E_{i,t}\|_2^2, \|E_{j,t}\|_2^2 \right\}$$

by Proposition 4. Also observe that, in any window of $4N$ messages, the spanning-tree protocol "sends information" from every agent to every other. Together these observations imply that $\eta_{t+4N} \geq \max_i \|E_{i,t}\|_2^2$. So as long as there exists an $i$ such that $\|E_{i,t}\|_2^2 \gg \eta_t$, the protocol makes significant progress.

It may happen, though, that $\|E_{i,t}\|_2^2$ is nearly constant as we range over $i$. Assume $\Pr\left[|E_{i,t} - E_{j,t}| > \varepsilon\right] \geq \delta$ for two agents $i, j$. Consider a path from $i$ to $j$ in $G$, obtained by first following $i$ to 1 in $\mathcal{T}_2$ and then

---

[5]If we want every pair of agents to agree within $\varepsilon$ with *global* probability $1 - \delta$, then we want every pair to $(\varepsilon, \delta/N^2)$-agree.

following 1 to $j$ in $\mathcal{T}_1$. This path has at most $2d$ edges. So by the triangle inequality, there exist consecutive agents $A, B$ along the path such that

$$\|E_{A,t} - E_{B,t}\|_2 \geq \frac{1}{2d} \|E_{i,t} - E_{j,t}\|_2 > \frac{\sqrt{\delta\varepsilon^2}}{2d}.$$

Imagine that the $t^{th}$ message is from $A$ to $B$. Then since $\{\Omega_{B,t+1}\}_{\omega\in\Omega}$ refines both $\{\Theta_{A,t}\}_{\omega\in\Omega}$ and $\{\Omega_{B,t}\}_{\omega\in\Omega}$, Proposition 4 yields

$$\|E_{B,t+1} - E_{A,t}\|_2^2 = \|E_{B,t+1}\|_2^2 - \|E_{A,t}\|_2^2,$$
$$\|E_{B,t+1} - E_{B,t}\|_2^2 = \|E_{B,t+1}\|_2^2 - \|E_{B,t}\|_2^2.$$

Also, by the triangle inequality either

$$\|E_{B,t+1} - E_{A,t}\|_2^2 \geq \frac{1}{4} \|E_{A,t} - E_{B,t}\|_2^2$$

or

$$\|E_{B,t+1} - E_{B,t}\|_2^2 \geq \frac{1}{4} \|E_{A,t} - E_{B,t}\|_2^2.$$

Therefore

$$\|E_{B,t+1}\|_2^2 > \min\left\{\|E_{A,t}\|_2^2, \|E_{B,t}\|_2^2\right\} + \frac{1}{4}\left(\frac{\delta\varepsilon^2}{4d^2}\right) \geq \eta_t + \frac{\delta\varepsilon^2}{16d^2}.$$

It remains only to show why the above result is not spoiled if $A$ or $B$ receive other messages before $A$ sends its message to $B$. Let $u$ be the first time step after $t$ in which $A$ sends a message to $B$, and suppose the steps between $t$ and $u$ somehow reduce the distance between $E_A$ and $E_B$:

$$\|E_{A,u} - E_{B,u}\|_2^2 \leq \frac{\delta\varepsilon^2}{16d^2}.$$

Then by the triangle inequality (again!):

$$\|E_{A,u} - E_{A,t}\|_2 + \|E_{B,u} - E_{B,t}\|_2 \geq \|E_{B,t} - E_{A,t}\|_2 - \|E_{B,u} - E_{A,u}\|_2 > \sqrt{\frac{\delta\varepsilon^2}{4d^2}} - \sqrt{\frac{\delta\varepsilon^2}{16d^2}}$$

so either

$$\|E_{A,u} - E_{A,t}\|_2^2 > \frac{\delta\varepsilon^2}{64d^2}$$

or

$$\|E_{B,u} - E_{B,t}\|_2^2 > \frac{\delta\varepsilon^2}{64d^2}.$$

Suppose the former without loss of generality. Then since $\{\Omega_{A,u}\}_{\omega\in\Omega}$ refines $\{\Omega_{A,t}\}_{\omega\in\Omega}$,

$$\|E_{A,u}\|_2^2 = \|E_{A,t}\|_2^2 + \|E_{A,u} - E_{A,t}\|_2^2 > \eta_t + \frac{\delta\varepsilon^2}{64d^2}.$$

We have shown that $\max_i \|E_{i,t+2N}\|_2^2 = \eta_t + \Omega\left(\delta\varepsilon^2/d^2\right)$, from which it follows that $\eta_{t+6N} = \eta_t + \Omega\left(\delta\varepsilon^2/d^2\right)$. Hence the constraint $\eta_t \leq 1$ yields an upper bound of $O\left(Nd^2/\left(\delta\varepsilon^2\right)\right)$ on the number of messages. ∎

Let us make three remarks about Theorem 10. First, naturally one can combine Theorems 10 and 6, to obtain an $N$-agent protocol in which the messages are discrete. We omit the details here. Second, all we really need about the *order* of messages is that information gets propagated from any agent in $G$ to any other in a reasonable number of steps. Our spanning-tree construction was designed to guarantee this, but sending messages in a random order (for example) would also work. Third, it seems fair to assume that many agents send messages in parallel; if so, our complexity bound can almost certainly be improved.

# 4 Computational Complexity

The previous sections have weakened the idea that communication cost is a fundamental barrier to agreement. However, we have glossed over the issue of *computational* cost entirely. A protocol that requires only $O\left(1/\left(\delta\varepsilon^2\right)\right)$ messages has little real-world relevance if it would take Alice and Bob billions of years to calculate the messages! Moreover, all protocols discussed above seem to have that problem, since the number of possible states $|\Omega|$ could be exponential in the length $n$ of the agents' inputs.

Recognizing this issue, Hanson [13] introduced the notion of a "Bayesian wannabe": a computationally-bounded agent that can still make sense of what its expectations would be if it had enough computational power to be a Bayesian. He then showed that under certain assumptions, if two Bayesian wannabes agree to disagree about the expectation of a function $f(\omega)$, then they must also disagree about some variable that is independent of the state of the world $\omega \in \Omega$. However, Hanson's result does not suggest a *protocol* by which two Bayesian wannabes who agree about all state-independent variables could come to agree about $f$ as well.

Admittedly, if the two wannabes have *very* limited abilities, it might be trivial to get them to agree. For example, if Alice and Bob both ignore all their evidence and estimate $f(\omega) = 1/3$, then they agree before exchanging even a single message. But this example seems contrived: after all, if one the agents (with equal justification) estimated $f(\omega) = 2/3$, then no sequence of messages would ever cause them to agree within $\varepsilon < 1/3$. So informally, what we really want to know is whether two wannabes will always agree, having put in a "good-faith effort" to emulate Bayesian rationality.

We are thus led to the following question. Is there an agreement protocol that

(i) would cause two computationally-unbounded Bayesians to $(\varepsilon, \delta)$-agree after a small number of messages, and

(ii) can be simulated using a small amount of computation?

We will say shortly what we mean by a "small amount of computation." By "simulate," we mean that a computationally-unbounded referee, given the state $\omega \in \Omega$ together with a transcript $M = (m_1, \ldots, m_R)$ of all messages exchanged during the protocol, should be unable to decide (with non-negligible bias) whether Alice and Bob were Bayesians following the protocol exactly, or Bayesian wannabes merely simulating it. More formally, let $\mathcal{B}(\omega)$ be the probability distribution over message transcripts, assuming Alice and Bob are Bayesians and the state of the world is $\omega$. Likewise, let $\mathcal{W}(\omega)$ be the distribution assuming Alice and Bob are wannabes. Then we require that for all Boolean functions $\Phi(\omega, M)$,

$$\left| \Pr_{\omega \in \mathcal{D}, M \in \mathcal{B}(\omega)} [\Phi(\omega, M) = 1] - \Pr_{\omega \in \mathcal{D}, M \in \mathcal{W}(\omega)} [\Phi(\omega, M) = 1] \right| \leq \zeta \tag{*}$$

where $\zeta$ is a parameter that can be made as small as we like (say 0.00001).

A consequence of the requirement (*) is that even if Alice is computationally unbounded, she cannot decide with bias greater than $\zeta$ whether Bob is also unbounded, judging only from the messages he sends to her. For if Alice could decide, then so could our hypothetical referee, who learns at least as much about Bob as Alice does. Though a little harder to see, another consequence is that if Alice is unbounded, but knows Bob to be bounded and *takes his algorithm into account* when computing her expectations, her messages will still be statistically indistinguishable from what they would have been had she believed that Bob was unbounded. Indeed, no beliefs, beliefs about beliefs, etc., about whether either agent is bounded or not can significantly affect the sequence of messages, since the truth or falsehood of those beliefs is almost irrelevant to predicting the agents' future messages. Also, if Alice is unbounded for some steps of the protocol but bounded for others, then Bob will never notice these changes, and would hardly behave any differently were he told of them.

Because of these considerations, we claim that, while simulating a Bayesian agreement protocol might not be the *only* way for two Bayesian wannabes to reach an "honest" agreement, it is certainly a *sufficient* way. Therefore, if we can show how to meet even the stringent requirement (*), this will provide strong evidence that computation time is not a fundamental barrier to agreement.

But what do we mean by computation time? We assume the state space $\Omega$ is a subset of $\{0,1\}^n \times \{0,1\}^n$, so that Alice's initial knowledge is an $n$-bit string $x$, and Bob's is an $n$-bit string $y$. Given the prior

distribution $\mathcal{D}$ over $(x, y)$ pairs, let $\mathcal{D}_{A,x}$ be Alice's posterior distribution over $y$ conditioned on $x$, and let $\mathcal{D}_{B,y}$ be Bob's posterior distribution over $x$ conditioned on $y$. The following two computational assumptions are the only ones that we make:

(1) Alice and Bob can both evaluate $f(\omega)$ for any $\omega \in \Omega$.

(2) Alice and Bob can both sample from $\mathcal{D}_{A,x}$ for any $x \in \{0,1\}^n$, and from $\mathcal{D}_{B,y}$ for any $y \in \{0,1\}^n$.

Our simulation procedure will *not* have access to descriptions of $f$ or $\mathcal{D}$; it can learn about them only by calling subroutines for (1) and (2) respectively. The complexity of the procedure will then be expressed in terms of the number of subroutine calls, other computations adding a negligible amount of time. Thus, we might stipulate that both subroutines should run in time polynomial in $n$. On the other hand, $n$ could be extremely large—otherwise the agents would simply exchange their entire inputs and be done! So we probably want to be even stricter, and stipulate that the subroutines should use time (say) *logarithmic* in $n$, albeit with many parallel processors. The latter seems like a better model for the human brain; after all, to reach an opinion based on our current knowledge, we do not contemplate every fact we know in sequential order, but instead zero in quickly on the relevant facts. In any case, the simulation procedure will treat the subroutines purely as "black boxes," so decisions about their implementation will not affect our results.

The justification for assumptions (1) and (2) is that without them, it is hard to see how the agents could estimate their expectations even before they started talking to each other. In other words, we have to assume the agents enter the conversation with minimal tools for reasoning about their universe of discourse. We do *not* assume that those tools extend to reasoning about each other's expectations, expectations of expectations, etc., conditioned on a sequence of messages exchanged. That the tools do extend in this way is what we intend to prove.

The one assumption that seems debatable to us is that Alice can sample from Bob's distribution $\mathcal{D}_{B,y}$, and Bob can sample from Alice's distribution $\mathcal{D}_{A,y}$. How can an agent possibly be expected to possess "someone else's" sampling subroutine? On further reflection, though, this question is simply a variant of an earlier question: why can we assume that Alice knows Bob's set of possible states $\Omega_B(\omega)$, and that Bob knows $\Omega_A(\omega)$? For if Alice knows $\Omega_B(\omega)$ as well as Bob does, then there is no particular reason why she should not be able to sample from it as well as he can. Again, the reason the agents know each other's partitions is that the state of the world $\omega \in \Omega$ includes both agents' mental states as part of it. None of this seems *too* out of line with everyday experience—for whenever we use what we know to try and figure out what someone else might be thinking, a Bayesian would say we are sampling an $\omega$ from our set of possible states, then sampling from what the other person's set of possible states would be if the state of the world were $\omega$.

Finally, let us note that assumptions (1) and (2) can both be relaxed. In particular, it is enough to approximate $f(\omega)$ to within an additive factor $\eta$ with probability at least $1 - \eta$, in time that increases polynomially in $1/\eta$. It is also enough to sample from a distribution whose variation distance from $\mathcal{D}_{A,x}$ or $\mathcal{D}_{B,y}$ is at most $\eta$, in time polynomial in $1/\eta$. Indeed, since the probabilities and $f$-values are real numbers, we will generally *need* to approximate in order to represent them with finite precision. For ease of presentation, though, we assume exact algorithms in what follows.

## 4.1 Smoothed Standard Protocol

Naïvely, requirement (*) seems impossible to satisfy. All of the agreement protocols discussed earlier in this paper—for example, that of Theorem 6—are easy to distinguish from any efficient simulation of them. For consider Alice's first message to Bob. If Alice's expectation $E_{A,0}$ is below some threshold $c$, she sends one message, whereas if $E_{A,0} \geq c$, she sends a different message. Even if we fix $f$, and limit probabilities and $f$-values to (say) $n$ bits of precision, we can arrange things so that $E_{A,0}(\omega)$ is exponentially close to $c$, sometimes greater and sometimes less, with high probability over $\omega$. Then to decide which message to send, Alice needs to evaluate $f$ exponentially many times.

We resolve this issue by having the agents add random noise to their messages ("smoothing" them), even if they are unbounded Bayesians. This noise does not prevent the agents from reaching $(\varepsilon, \delta)$-agreement. On the other hand, it makes their messages easier to simulate. For unlike real numbers $a \neq b$, which are perfectly distinguishable no matter how close they are, two probability distributions with close means may be hard to distinguish, like wavepackets in quantum mechanics.
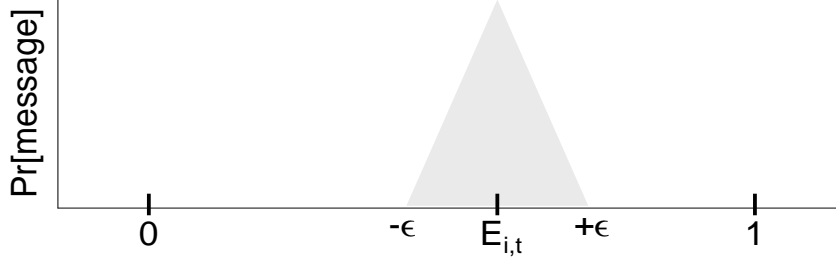
Figure 4: Agent $i$ "smoothes" its expectation $E_{i,t}$ with triangular noise before sending it.

In the *smoothed standard protocol*, Alice generates her messages to Bob as follows. Let $b \geq \log_2(200/\varepsilon)$ be a positive integer to be specified later. Then let $\epsilon$ be an integer multiple of $2^{-b}$ between $\varepsilon/50$ and $\varepsilon/40$, and let $L = 2^b \epsilon$. First Alice rounds her current expectation $E_{A,t}$ of $f$ to the nearest multiple of $2^{-b}$. Denote the result by $\text{round}(E_{A,t})$. She then draws an integer $r \in \{-L, \ldots, L\}$, according to a *triangular distribution* in which $r = j$ with probability $(L - |j|)/L^2$ (see Figure 4). The message she sends Bob is $m_{t+1} = \text{round}(E_{A,t}) + 2^{-b}r$. Observe that since $m_{t+1} \in [-\epsilon, 1 + \epsilon]$, there are at most $2^b(1 + 2\epsilon) + 1$ possible values of $m_{t+1}$—meaning Alice's message takes only $b + 1$ bits to specify. After receiving the message, Bob updates his expectation of $f$ using Bayes' rule, then draws an integer $r \in \{-L, \ldots, L\}$ according to the same triangular distribution and sends Alice $m_{t+2} = \text{round}(E_{B,t+1}) + 2^{-b}r$. The two agents continue to send messages in this way.

The reader might be wondering why we chose triangular noise, and whether other types of noise would work equally well. The answer is that we want the message distribution to have three basic properties. First, it should be concentrated about a mean of $E_{i,t}$ with variance at most $\tilde{}\epsilon^2$. Second, shifting the mean by $\eta \leq \epsilon$ should shift the distribution by at most $\tilde{}\eta/\epsilon$ in variation distance. And third, the derivative of the probably density function should never exceed $\tilde{}\eta/\epsilon^2$ in absolute value. Thus, Gaussian noise would also work, though it is somewhat harder to analyze than triangular noise. However, noise that is uniform over $[-\epsilon, \epsilon]$ would *not* work (so far as we could tell), since it violates the third property.

Before we analyze the protocol, we need to develop some notation. Let $M_t = (m_1, \ldots, m_t)$ consist of the first $t$ messages that Alice and Bob exchange. Since messages are now probabilistic, the agents' expectations of $f$ at step $t$ depend not only on the initial state of the world $\omega$, but also on $M_t$. When we want to emphasize this, we denote the agents' expectations by $E_{A,t}(\omega, M_t)$ and $E_{B,t}(\omega, M_t)$ respectively. Another important consequence of messages being probabilistic is that after an agent has received a message, its posterior distribution over $\omega$ is no longer obtainable by restricting the prior distribution $\mathcal{D}$ to a subset of possible states. Thus, we let $\Omega_i(\omega) = \Omega_{i,0}(\omega)$, since we will never refer to $\Omega_{i,t}(\omega)$ for $t > 0$.

Say the agents $(\varepsilon, \delta)$-agree after the $t^{th}$ message if

$$\Pr_{\omega \in \Omega, M_t}[|E_{A,t}(\omega, M_t) - E_{B,t}(\omega, M_t)| > \varepsilon] \leq \delta.$$

Also, let

$$\|E_{i,t}\|_2^2 = \underset{\omega \in \Omega, M_t}{\text{EX}}\left[E_{i,t}(\omega, M_t)^2\right].$$

**Theorem 11** *For all $f, \mathcal{D}$, the smoothed standard protocol causes Alice and Bob to $(\varepsilon, \delta)$-agree after at most $2/(\delta\varepsilon^2)$ messages.*

**Proof.** Similarly to Theorem 6, we let $E_{C,t}$ be the expectation of a third party Charlie who sees all messages between Alice and Bob, but who knows neither their inputs nor the random bits that they use to produce their messages. We then track $\|E_{C,t}\|_2^2$.

Assume that $\Pr[|E_{A,t} - E_{B,t}| > \varepsilon] \geq \delta$ and that Alice sends the $t^{th}$ message $m_t$. Notice that $m_t$ cannot deviate from Alice's expectation $E_{A,t} = E_{A,t-1}$ by more than $2\epsilon$, since $|\text{round}(E_{A,t}) - E_{A,t}| \leq \epsilon$ and $|m_t - \text{round}(E_{A,t})| \leq \epsilon$. So keeping $M_t$ fixed,

$$|E_{A,t}(\omega, M_t) - E_{A,t}(\omega', M_t)| \leq 4\epsilon$$

16

for all $\omega, \omega'$. Now Charlie's expectation $E_{C,t}(\omega, M_t)$ is just an average of $E_{A,t}(\omega', M_t)$'s, so it follows that

$$|E_{C,t}(\omega, M_t) - E_{A,t}(\omega, M_t)| \leq 4\epsilon$$

as well. Similarly, after Bob sends the $(t+1)^{st}$ message,

$$|E_{C,t+1}(\omega, M_{t+1}) - E_{B,t+1}(\omega, M_{t+1})| \leq 4\epsilon.$$

Therefore

$$\Pr[|E_{C,t+1} - E_{C,t}| > \varepsilon - 8\epsilon] \geq \Pr[|E_{A,t} - E_{B,t}| > \varepsilon] \geq \delta,$$

using the triangle inequality and the fact that $E_{B,t+1} = E_{B,t}$. The final observation is that Charlie's partition of $\Omega \times M_{t+1}$ at step $t+1$ refines his partition at step $t$, so by Proposition 4,

$$\|E_{C,t+1}\|_2^2 - \|E_{C,t}\|_2^2 = \|E_{C,t+1} - E_{C,t}\|_2^2 > \delta (\varepsilon - 8\epsilon)^2.$$

Since $\|E_{C,t}\|_2^2 \leq 1$, this yields an upper bound of $1/\left(\delta(\varepsilon - 8\epsilon)^2\right) < 2/\left(\delta\varepsilon^2\right)$ on the number of messages. ∎

## 4.2 Simulating the Smoothed Protocol

Having proved that the smoothed standard protocol works, in this section we explain how Alice and Bob can simulate the protocol. In the ideal case—where the agents have unlimited computational power—they use the following recursive formulas. Let

$$\Delta(m_t, E_{i,t-1}) = \begin{cases} 1 - |m_t - \text{round}(E_{i,t-1})|/\epsilon & \text{if } |m_t - \text{round}(E_{i,t-1})| \leq \epsilon \\ 0 & \text{otherwise} \end{cases}$$

be proportional to the probability that agent $i$ sends message $m_t$, given that its expectation is $E_{i,t-1}$. Also, let $q_t(\omega, M_t)$ be proportional to the joint probability of messages $m_1, \ldots, m_t$ assuming the true state of the world is $\omega$. Then assuming $t$ is even and suppressing dependencies on $M_t$, for all $X, Y$ we have

$$q_t(Y) = q_{t-2}(Y) \Delta(m_t, E_{B,t-1}(Y)),$$
$$q_{t-1}(X) = q_{t-3}(X) \Delta(m_{t-1}, E_{A,t-2}(X)),$$
$$E_{A,t}(X) = \frac{\text{EX}_{Y \in \Omega_A(X)}[q_t(Y) f(Y)]}{\text{EX}_{Y \in \Omega_A(X)}[q_t(Y)]},$$
$$E_{B,t-1}(Y) = \frac{\text{EX}_{X \in \Omega_B(Y)}[q_{t-1}(X) f(X)]}{\text{EX}_{X \in \Omega_B(Y)}[q_{t-1}(X)]}$$

with the base cases $q_0(Y) = q_{-1}(X) = 1$ for all $X, Y$. The correctness of these formulas follows from simple Bayesian manipulations. Having computed $E_{i,t}(\omega)$ by the formulas above (note that this does not require knowledge of $\omega$), all agent $i$ needs to do is draw $r \in \{-L, \ldots, L\}$ from the triangular distribution, then send the message

$$m_{t+1} = \text{round}(E_{i,t}(\omega)) + 2^{-b}r.$$

In the real case, the agents are computationally bounded, and can no longer afford the luxury of taking expectations over the exponentially large sets $\Omega_i$. A natural idea is to compensate by somehow *sampling* those sets. But since we never assumed the ability to sample $\Omega_i$ conditioned on messages $m_1, \ldots, m_t$, it is not obvious how that make that idea work. Our solution will consist of two phases: the construction of "sampling-trees," which involves no communication, followed by a message-by-message simulation of the ideal protocol. Let us describe these phases in turn.

   **(I) Sampling-Tree Construction.** Alice creates a tree $\mathcal{T}_A$ with height $R$ and branching factor $K$. Here $R < 2/\left(\delta\varepsilon^2\right)$ is the number of messages, and $K$ is a parameter to be specified later. Let $\text{root}_A$ be the root node of $\mathcal{T}_A$, and let $S(v)$ be the set of children of node $v$. Then Alice labels each of the $K$ nodes $w \in S(\text{root}_A)$ by a sample $Y_w \in \Omega_A(\omega)$, drawn independently from her posterior distribution $\mathcal{D}_{A,x}$. Next, for each $w \in S(\text{root}_A)$, she labels each of the $K$ nodes $v \in S(w)$ by a sample $X_v \in \Omega_B(Y_w)$, drawn independently from Bob's distribution $\mathcal{D}_{B,y}$ where $Y_w = (x, y)$. She continues recursively in this manner,

labeling each $v$ an even distance from the root with a sample $X_v \in \Omega_B (Y_w)$ where $w$ is the parent of $v$, and each $w$ an odd distance from the root with a sample $Y_w \in \Omega_A (X_v)$ where $v$ is the parent of $w$. Thus her total number of samples is

$$K + K^2 + \cdots + K^R = \frac{K^{R+1} - 1}{K - 1} - 1.$$

Similarly, Bob creates a tree $\mathcal{T}_B$ with height $R$ and branching factor $K$. Let $\text{root}_B$ be the root of $\mathcal{T}_B$; then Bob labels each $v \in S (\text{root}_B)$ by a sample $X_v \in \Omega_B (\omega)$, each child $w \in S (v)$ of each $v \in S (\text{root}_B)$ by a sample $Y_w \in \Omega_A (X_v)$, and so on, alternating between $\Omega_B$ and $\Omega_A$ at successive levels. As a side remark, if the agents share a random string, then there is no reason for them not to use the same set of samples. However, we cannot assume that such a string is available.

**(II) Simulation.** We now explain how the agents can use the samples from (I) to simulate the smoothed standard protocol. First Alice estimates her expectation $E_{A,0}$ by the quantity

$$\langle E_{A,0} (\text{root}_A) \rangle_A = \underset{w \in S(\text{root}_A)}{\text{EX}} [f (Y_w)] = \frac{1}{K} \sum_{w \in S(\text{root}_A)} f (Y_w).$$

She then chooses a random $r \in \{-L, \ldots, L\}$ and sends Bob

$$m_1 = \text{round} \left( \langle E_{A,0} (\text{root}_A) \rangle_A \right) + 2^{-b} r.$$

On receiving the message, for each $v \in S (\text{root}_B)$ Bob computes

$$\langle E_{A,0} (v) \rangle_B = \frac{1}{K} \sum_{w \in S(v)} f (Y_w),$$

his estimate of $E_{A,0} (X_v)$ assuming $\omega = X_v$. He then defines

$$\langle q_0 (v) \rangle_B = \Delta \left( m_1, \langle E_{A,0} (v) \rangle_B \right)$$

and estimates his own expectation $E_{B,1} (\omega)$ by

$$\langle E_{B,1} (\text{root}_B) \rangle_B = \frac{\sum_{v \in S(\text{root}_B)} \langle q_0 (v) \rangle_B f (X_v)}{\sum_{v \in S(\text{root}_B)} \langle q_0 (v) \rangle_B}.$$

Finally, he chooses a random $r \in \{-L, \ldots, L\}$ and sends Alice

$$m_2 = \text{round} \left( \langle E_{B,1} (\text{root}_B) \rangle_B \right) + 2^{-b} r.$$

In general, if $t$ is even then the recursive formulas for agent $i$ are

$$\langle q_t (w) \rangle_i = \langle q_{t-2} (w) \rangle_i \Delta \left( m_t, \langle E_{B,t-1} (w) \rangle_i \right),$$
$$\langle q_{t-1} (v) \rangle_i = \langle q_{t-3} (v) \rangle_i \Delta \left( m_{t-1}, \langle E_{A,t-2} (v) \rangle_i \right),$$
$$\langle E_{A,t} (v) \rangle_i = \frac{\sum_{w \in S(v)} \langle q_t (w) \rangle_i f (Y_w)}{\sum_{w \in S(v)} \langle q_t (w) \rangle_i},$$
$$\langle E_{B,t-1} (w) \rangle_i = \frac{\sum_{v \in S(w)} \langle q_{t-1} (v) \rangle_i f (X_v)}{\sum_{v \in S(w)} \langle q_{t-1} (v) \rangle_i}$$

with the base cases $\langle q_0 (w) \rangle_i = \langle q_{-1} (v) \rangle_i = 1$ for all $w, v$. Agent $i$ computes a message $m_t$ in the obvious way, from its expectation at the root of $\mathcal{T}_i$:

$$m_t = \text{round} \left( \langle E_{i,t-1} (\text{root}_i) \rangle_i \right) + 2^{-b} r.$$

That completes the description of the simulation procedure. Its complexity is easily determined: let $T_1$ be the number of computational steps needed to sample from $\mathcal{D}_{A,x}$ or $\mathcal{D}_{B,y}$, and let $T_2$ be number of steps needed to evaluate $f$. Then both agents use $O \left( K^R (T_1 + T_2) \right)$ steps, where we have summed over all $R$ communication rounds. Thus, the complexity is exponential in $R \approx 2/ (\delta \varepsilon^2)$; on the other hand, it has no dependence on $n$.

## 4.3 Analysis

Our goal is to show that the message sequence in the simulated protocol is statistically indistinguishable from the sequence in the ideal protocol, for some reasonable sample size $K$. Here 'reasonable', unfortunately, is still quite huge: of order $(11/\epsilon)^{R^2}/\zeta^2$, where $\zeta$ is the maximum bias with which a referee can distinguish the conversations. So assuming $\epsilon \geq \varepsilon/50$ and $R \leq 2/\left(\delta\varepsilon^2\right)$, the total number of computational steps is of order

$$\left(\frac{(11/\epsilon)^{R^2}}{R^2}\right)^R (T_1 + T_2) = \exp\left(\frac{8\ln\left(550/\varepsilon\right)}{\delta^3\varepsilon^6} + \frac{4\ln\left(1/\zeta\right)}{\delta\varepsilon^2}\right)(T_1 + T_2).$$

The reader might complain that this bound is not at all reasonable: for example, if $\varepsilon = \delta = 1/2$, then it translates into more than $2^{36864}$ subroutine calls! Let us make two points in response. First, we do show that the number of subroutine calls needed is independent of $n$, and that it grows "only" exponentially in a polynomial in $1/\delta$ and $1/\varepsilon$. Theoretical computer scientists often see cases in which the first polynomial-time algorithm for a problem has a completely impractical complexity, say $n^{40}$. However, once the problem is known to be in polynomial time, it is usually possible to reduce the exponent to obtain a truly practical algorithm. In our case, we conjecture that the factor of $1/\left(\delta^3\varepsilon^6\right)$ in the exponent could be reduced to $1/\left(\delta^2\varepsilon^4\right)$ or even $1/\left(\delta\varepsilon^2\right)$; certainly the constants in the exponent can be reduced. The second point is that the complexity is so large only because we never assumed the agents can sample from their sets of possible states *conditioned* on messages exchanged. So the best they can do is to sample a huge number of states from their original sets $\Omega_A$ and $\Omega_B$, then retain the few that are compatible with the messages. However, it seems likely that agents would have at least some ability to sample conditioned on messages. After all, we assumed that they enter the conversation with the ability to sample, and presumably they have had other conversations in the past! In practice, then, the complexity will probably be better than the worst-case estimate above.

How do we prove the simulation theorem? In one sense, the proof is 'merely' an exercise in error analysis and large deviation bounds. However, the details are extremely subtle and difficult to get right. The problem is that if a message has probability $q$ from its recipient's point of view, then order $1/q$ samples are needed to find even a single input that could have caused the sender to produce that message. Fortunately, low-probability messages are unlikely to be sent, for almost tautological reasons that we spell out in Lemma 14. However, because the sample trees $\mathcal{T}_i$ are so large, with overwhelming probability they contain *some* nodes $v$ with miniscule values of $\langle q_t(v)\rangle_i$. We need to argue that the errors introduced by these "bad nodes" are washed out by the good nodes before they can propagate to the root.

The proof will repeatedly use the Chernoff-Hoeffding bound (Theorem 2). As shown by the following corollary, Theorem 2 sometimes lets us estimate the mean of a random variable, even if we cannot sample that variable directly.

**Corollary 12** *Let $p_1, \ldots, p_n$ and $x_1, \ldots, x_n$ belong to $[0, 1]$, and let $P = p_1 + \cdots + p_n$ and $x = p_1 x_1 + \cdots + p_n x_n$. If we choose $K$ indices $i(1), \ldots, i(K)$ uniformly at random from $\{1, \ldots, n\}$, then*

$$\Pr\left[\left|\frac{p_{i(1)}x_{i(1)} + \cdots + p_{i(K)}x_{i(K)}}{p_{i(1)} + \cdots + p_{i(K)}} - \frac{x}{P}\right| > \alpha\right] \leq 4e^{-\alpha^2(P/n)^2 K/2}.$$

**Proof.** Let

$$\widetilde{P} = \frac{n}{K}\left(p_{i(1)} + \cdots + p_{i(K)}\right),$$

$$\widetilde{X} = \frac{n}{K}\left(p_{i(1)}x_{i(1)} + \cdots + p_{i(K)}x_{i(K)}\right).$$

Then since $\widetilde{X} \leq \widetilde{P}$,

$$\left|\frac{\widetilde{X}}{\widetilde{P}} - \frac{x}{P}\right| = \frac{\left|\widetilde{X}\left(P - \widetilde{P}\right) - \widetilde{P}\left(x - \widetilde{X}\right)\right|}{\widetilde{P}P} \leq \frac{\left|\widetilde{P} - P\right|}{P} + \frac{\left|\widetilde{X} - X\right|}{P}.$$

So

$$\Pr\left[\left|\frac{\widetilde{X}}{\widetilde{P}} - \frac{X}{P}\right| > \alpha\right] \le \Pr\left[\left|\widetilde{P} - P\right| > \frac{\alpha P}{2}\right] + \Pr\left[\left|\widetilde{X} - X\right| > \frac{\alpha P}{2}\right].$$

By Theorem 2,

$$\Pr\left[\frac{K}{n}\left|\widetilde{P} - P\right| > \frac{\alpha P}{2n}K\right] \le 2e^{-\alpha^2 (P/n)^2 K/2}$$

and similarly for $\left|\widetilde{X} - X\right|$. ∎

We will also need a bound for a sum of exponentially distributed variables, which can be found in [6] for example.

**Theorem 13** *Let* $x_1, \ldots, x_K \in [0, \infty)$ *be independent and exponentially distributed with mean 1 (that is,* $\Pr[x_i \ge x] = e^{-\omega}$). *Then*

$$\Pr\left[x_1 + \cdots + x_K \ge (1 + \alpha) K\right] \le \left(\frac{e^\alpha}{1 + \alpha}\right)^{-K}.$$

For convenience, we will state our results in terms of Alice's tree $\mathcal{T}_A$, with the understanding that they apply equally well to $\mathcal{T}_B$. Throughout, we assume that $t$ is even and that the $t^{th}$ message $m_t$ is sent from Bob to Alice. Let $Q_t = \sum_{Y \in \Omega_A(\omega)} q_t(Y)$ measure the "likelihood" of Alice's situation at step $t$. Then $Q_t/Q_{t-2}$ measures the likelihood of the $t^{th}$ message, conditioned on Alice's situation just before she receives it. The following lemma says essentially that "unlikely messages are unlikely."

**Lemma 14** *For all inputs* $x$ *of Alice, message sequences* $M_{t-1}$, *and constants* $\gamma > 0$,

$$\Pr_{m_t}\left[\frac{Q_t}{Q_{t-2}} \le \frac{\gamma\epsilon}{2}\right] < \gamma.$$

**Proof.** For all $m \in [-\epsilon, 1 + \epsilon]$,

$$\Pr[m_t = m] = \sum_{j \in \{-L, \ldots, L\}} \left(\Pr_Y\left[\text{round}\left(E_{B,t-1}(Y)\right) = m + 2^{-b}j\right] \cdot \frac{\Delta\left(m, m + 2^{-b}j\right)}{L}\right)$$

$$= \frac{1}{L}\frac{\sum_{Y \in \Omega_A(\omega)} q_{t-2}(Y) \Delta\left(m, E_{B,t-1}(Y)\right)}{\sum_{Y \in \Omega_A(\omega)} q_{t-2}(Y)} = \frac{1}{L}\frac{Q_t}{Q_{t-2}}$$

from Alice's point of view. So it suffices to observe that

$$\Pr_m\left[\Pr_{m_t}[m_t = m] \le \frac{\gamma\epsilon}{2L}\right] \le \frac{\gamma\epsilon}{2L}\frac{L(1 + 2\epsilon) + \epsilon}{\epsilon} < \gamma.$$

Here the first inequality follows from elementary probability theory, together with the fact that there are at most $(1 + 2\epsilon)/2^{-b} + 1$ possible messages $m$, and hence the mean of $\Pr_{m_t}[m_t = m]$ over $m$ chosen uniformly at random is at least

$$\frac{1}{(1 + 2\epsilon)/2^{-b} + 1} = \frac{\epsilon}{L(1 + 2\epsilon) + \epsilon}.$$

The second inequality follows since $\epsilon < 1/4$. ∎

A consequence of Lemma 14 is that unlikely *sequences* of messages are unlikely. For the remainder of this section, let $g = \frac{4e}{\epsilon}\ln K$.

**Lemma 15** *For all* $\gamma > 0$ *and all* $x$,

$$\Pr_{y, M_t}[Q_t \le \gamma] < g^{t/2}\max\left\{\gamma, \frac{1}{K}\right\}.$$

20

**Proof.** For all $u \in \{2, 4, \ldots, t\}$, let $x_u = \ln(\epsilon Q_{u-2}/2Q_u)$. Then

$$Q_t = \frac{2Q_t}{\epsilon Q_{t-2}} \frac{2Q_{t-2}}{\epsilon Q_{t-4}} \cdots \frac{2Q_2}{\epsilon Q_0} \left(\frac{\epsilon}{2}\right)^{t/2} = e^{-x_2 - x_4 - \cdots - x_t} \left(\frac{\epsilon}{2}\right)^{t/2}$$

since $Q_0 = 1$. Furthermore, Lemma 14 implies that for each $u$,

$$\Pr_{m_u} [x_u \geq x] = \Pr_{m_u} \left[\frac{Q_u}{Q_{u-2}} \leq \frac{e^{-x}\epsilon}{2}\right] < e^{-x},$$

even conditioned on $x_2, \ldots, x_{u-2}$. Therefore $x_2 + \cdots + x_t$ is stochastically dominated by a sum of $t/2$ independent exponential variables each with mean 1. So by Theorem 13,

$$\Pr\left[x_2 + \cdots + x_t \geq (1+\alpha)\frac{t}{2}\right] < \left(\frac{e^\alpha}{1+\alpha}\right)^{-t/2}.$$

Setting $\gamma = e^{-(1+\alpha)t/2} (\epsilon/2)^{t/2}$ and solving to obtain $\alpha = (2/t)\ln\left((\epsilon/2)^{t/2}/\gamma\right) - 1$, it follows that

$$\Pr_{y, M_t} [Q_t \leq \gamma] < \left(\frac{e^{(2/t)\ln\left((\epsilon/2)^{t/2}/\gamma\right) - 1}}{(2/t)\ln\left((\epsilon/2)^{t/2}/\gamma\right)}\right)^{-t/2} < \left(\frac{4e}{\epsilon}\ln\frac{1}{\gamma}\right)^{t/2} \gamma \leq g^{t/2} \max\left\{\gamma, \frac{1}{K}\right\}.$$

∎

In the next four results, we fix a particular node $v \in \mathcal{T}_A$, then study how the error at $v$ depends on the errors at its children $w \in S(v)$. For simplicity, we assume $v$ is an even distance from the root, but our results will apply equally to nodes an odd distance from the root. We need to upper-bound the expected difference between Alice's actual expectation $\langle E_{A,t}(v)\rangle_A$, and her ideal expectation $E_{A,t}(X_v)$. To this end, it will be helpful to define the following "hybrid" between $\langle E_{A,t}(v)\rangle_A$ and $E_{A,t}(X_v)$:

$$E_{A,t}^*(v) = \frac{\sum_{w \in S(v)} q_t(Y_w) f(Y_w)}{\sum_{w \in S(v)} q_t(Y_w)}.$$

To compute $E_{A,t}^*$, we use the ideal weights $q_t(Y_w)$, but we average over Alice's $K$ samples $\{Y_w\}_{w \in S(v)}$ only, not over all of $\Omega_A(X_v)$. By the triangle inequality, to upper-bound $\left|\langle E_{A,t}(v)\rangle_A - E_{A,t}(X_v)\right|$ it suffices to upper-bound $\left|\langle E_{A,t}(v)\rangle_A - E_{A,t}^*(v)\right|$ and $\left|E_{A,t}^*(v) - E_{A,t}(X_v)\right|$. We start with the latter.

**Lemma 16**

$$\underset{y, M_t, S(v)}{\mathrm{EX}} \left[\left|E_{A,t}^*(v) - E_{A,t}(X_v)\right|\right] \leq \frac{7g^{t/2+1}}{\sqrt{K}}.$$

**Proof.** Assuming $Q_t = Q$,

$$\Pr\left[\left|E_{A,t}^*(v) - E_{A,t}(X_v)\right| \geq \omega\right] \leq 4e^{-\omega^2 Q^2 K/2}$$

by Corollary 12. Furthermore, since $E_{A,t}^*(v)$ and $E_{A,t}(X_v)$ are in $[0, 1]$, we have the trivial but important

bound $\left|E_{A,t}^*(v) - E_{A,t}(X_v)\right| \leq 1$. Therefore

$$\mathrm{EX}\left[\left|E_{A,t}^*(v) - E_{A,t}(X_v)\right|\right] = \int_0^1 \Pr\left[\left|E_{A,t}^*(v) - E_{A,t}(X_v)\right| \geq x\right] dx$$

$$\leq 4\int_0^1 \mathop{\mathrm{EX}}_{Q_t}\left[e^{-x^2 Q_t^2 K/2}\right] dx$$

$$= 4\int_0^1 \int_0^1 \Pr\left[e^{-x^2 Q_t^2 K/2} \geq x\right] dx dx$$

$$= 4\int_0^1 \int_0^1 \Pr\left[Q_t \leq \frac{1}{x}\sqrt{\frac{2}{K}\ln\frac{1}{x}}\right] dx dx$$

$$\leq 4\int_0^1 \int_0^1 \min\left\{1, \max\left\{\frac{1}{x}\sqrt{\frac{2}{K}\ln\frac{1}{x}}, \frac{1}{K}\right\} g^{t/2}\right\} dx dx$$

$$\leq 4g^{t/2}\left(\frac{1}{K} + \int_0^1 \int_0^1 \min\left\{g^{-t/2}, \frac{1}{x}\sqrt{\frac{2}{K}\ln\frac{1}{x}}\right\} dx dx\right)$$

$$= 4g^{t/2}\left(\frac{1}{K} + \int_{x=0}^1 \sqrt{\frac{2}{K}\ln\frac{1}{x}}\left(x_{\min}(x) + \int_{x=x_{\min}(x)}^1 \frac{1}{x}dx\right) dx\right)$$

Here the fifth line uses Lemma 15, and

$$x_{\min}(x) = g^{t/2}\sqrt{\frac{2}{K}\ln\frac{1}{x}}.$$

By straightforward integral approximations, the last expression is at most $7g^{t/2+1}/\sqrt{K}$ for sufficiently large $K$. ∎

For each child $w \in S(v)$, let

$$\eta_t(w) = \sum_{u \in \{1,3,\ldots,t-1\}} \left|\langle E_{B,u}(w)\rangle_A - E_{B,u}(Y_w)\right|$$

measure the total error in Alice's estimates of $E_{B,u}(Y_w)$, summed over all time steps $u \leq t$. The following proposition shows that to upper-bound the error in $\langle q_t(w)\rangle_A$, it suffices to upper-bound $\eta_t(w)$. For this proposition to hold, we need the function $\Delta$ to have bounded derivative. That is why we chose triangular instead of uniform noise when defining the protocol.

**Proposition 17**

$$\left|\langle q_t(w)\rangle_A - q_t(Y_w)\right| \leq \frac{\eta_t(w)}{\epsilon}.$$

**Proof.** From the definition of $\Delta$,

$$\left|\Delta\left(m_{u+1}, \langle E_{B,u}(w)\rangle_A\right) - \Delta\left(m_{u+1}, E_{B,u}(Y_w)\right)\right| \leq \frac{1}{\epsilon}\left|\langle E_{B,u}(w)\rangle_A - E_{B,u}(Y_w)\right|.$$

Furthermore, $\Delta\left(m_{u+1}, \langle E_{B,u}(w)\rangle_A\right)$ and $\Delta\left(m_{u+1}, E_{B,u}(Y_w)\right)$ are both bounded in $[0,1]$. It follows that

$$\left|\langle q_t(w)\rangle_A - q_t(Y_w)\right| = \left|\prod_u \Delta\left(m_{u+1}, \langle E_{B,u}(w)\rangle_A\right) - \prod_u \Delta\left(m_{u+1}, E_{B,u}(Y_w)\right)\right|$$

$$\leq \sum_u \frac{1}{\epsilon}\left|\langle E_{B,u}(w)\rangle_A - E_{B,u}(Y_w)\right| = \frac{\eta_t(w)}{\epsilon}$$

where $u$ ranges over $\{1,3,\ldots,t-1\}$. ∎

Now let

$$H = \sum_{w \in S(v)} q_t \left( Y_w \right),$$

$$F = \sum_{w \in S(v)} q_t \left( Y_w \right) f \left( Y_w \right),$$

$$\langle H \rangle_A = \sum_{w \in S(v)} \langle q_t \left( w \right) \rangle_A,$$

$$\langle F \rangle_A = \sum_{w \in S(v)} \langle q_t \left( w \right) \rangle_A f \left( Y_w \right),$$

so that $E_{A,t}^* \left( v \right) = F/H$ and $\langle E_{A,t} \left( v \right) \rangle_A = \langle F \rangle_A / \langle H \rangle_A$. Using Lemma 15, we can upper-bound the probability that $H$ is too much smaller than its mean value.

**Corollary 18** *For all $\gamma > 0$,*

$$\Pr_{y, M_t, S(v)} \left[ H \le \gamma K \right] < 3g^{t/2} \max \left\{ \gamma, \frac{4 \ln K}{K} \right\}.$$

**Proof.** By the principle of deferred decisions, we can think of each $q_t \left( Y_w \right)$ as an independent sample of a $[0,1]$ random variable with mean $Q_t$. Then $H$ is a sum of $K$ such samples. Setting $\Gamma = \max \left\{ 2\gamma, 8 \left( \ln K \right) / K \right\}$, by Lemma 15 we have

$$\Pr_{y, M_t} \left[ Q_t \le \Gamma \right] < 2g^{t/2} \max \left\{ \gamma, \frac{4 \ln K}{K} \right\}.$$

Furthermore, assuming $Q_t > \Gamma$, Theorem 2 yields

$$\Pr_{S(v)} \left[ H \le \gamma K \right] \le \exp \left( -Q_t \left( 1 - \frac{\gamma}{\Gamma} \right)^2 \frac{K}{2} \right) \le e^{-\Gamma K/8} \le \frac{1}{K}.$$

The corollary now follows by the union bound. ∎

The last piece of the puzzle is to upper-bound the difference between $\langle E_{A,t} \left( v \right) \rangle_A$ and $E_{A,t}^* \left( v \right)$, using techniques similar to those of Lemma 16. Let $\eta = \mathrm{EX}_{w \in S(v)} \left[ \eta_t \left( w \right) \right]$ and $\widehat{\eta} = \mathrm{EX}_{y, M_t, \mathcal{T}_A} \left[ \eta \right]$.

**Lemma 19** *Assuming $\eta \ge 1/K$ for all $y, M_t, \mathcal{T}_A$,*

$$\mathrm{EX}_{y, M_t, \mathcal{T}_A} \left[ \left| \langle E_{A,t} \left( v \right) \rangle_A - E_{A,t}^* \left( v \right) \right| \right] \le 18g^{t/2+1} \widehat{\eta}.$$

**Proof.** Using the fact that $\langle F \rangle_A \le \langle H \rangle_A$,

$$\left| \langle E_{A,t} \left( v \right) \rangle_A - E_{A,t}^* \left( v \right) \right| = \left| \frac{\langle F \rangle_A}{\langle H \rangle_A} - \frac{F}{H} \right| \le \frac{\left| \langle H \rangle_A - H \right|}{H} + \frac{\left| \langle F \rangle_A - F \right|}{H}$$

by the same trick as in Corollary 12. Furthermore, it follows from Proposition 17 together with the triangle inequality that $\left| \langle H \rangle_A - H \right| \le \eta K / \epsilon$ and $\left| \langle F \rangle_A - F \right| \le \eta K / \epsilon$. So we can upper-bound $\left| \langle E_{A,t} \left( v \right) \rangle_A - E_{A,t}^* \left( v \right) \right|$ by $2\eta K / \left( \epsilon H \right)$, as well as (of course) by 1. Fix $\eta$; then

$$\mathrm{EX}_H \left[ \min \left\{ 1, \frac{2\eta K}{\epsilon H} \right\} \right] = \int_0^1 \Pr_H \left[ \frac{2\eta K}{\epsilon H} \ge x \right] dx$$

$$\le \int_0^1 \min \left\{ 1, 3 \max \left\{ \frac{2\eta}{\epsilon x}, \frac{4 \ln K}{K} \right\} g^{t/2} \right\} dx$$

$$\le 3g^{t/2} \left( \frac{4 \ln K}{K} + \int_0^1 \min \left\{ \frac{1}{3g^{t/2}}, \frac{2\eta}{\epsilon x} \right\} dx \right)$$

$$= 3g^{t/2} \left( \frac{4 \ln K}{K} + \frac{x_{\min}}{3g^{t/2}} + \int_{x_{\min}}^1 \frac{2\eta}{\epsilon x} dx \right)$$

23

where the second line uses Corollary 18 and $x_{\min} = (6\eta/\epsilon)\,g^{t/2}$. This in turn is at most

$$3g^{t/2}\left(\frac{4\ln K}{K} + \frac{2\eta}{\epsilon}\ln\frac{1}{\eta}\right).$$

Assuming $\eta \geq 1/K$ always, the expectation of the above quantity over $\eta$ is at most $18g^{t/2+1}\widehat{\eta}$. ∎

We are finally ready to put everything together, and show that the referee can distinguish the real and ideal conversations with bias at most $\zeta$.

**Theorem 20** *By setting $b = \lceil \log_2 R/(\zeta\epsilon)\rceil + 2$ and $K = O\left((11/\epsilon)^{R^2}/\zeta^2\right)$, it is possible to achieve*

$$\left|\Pr_{\omega\in\mathcal{D},M\in\mathcal{W}(\omega)}[\Phi(\omega,M_R)=1] - \Pr_{\omega\in\mathcal{D},M\in\mathcal{B}(\omega)}[\Phi(\omega,M_R)=1]\right| \leq \zeta$$

*for all Boolean functions $\Phi$.*

**Proof.** Combining Lemmas 16 and 19,

$$\underset{y,M_t,\mathcal{T}_A}{\mathrm{EX}}\left[\left|\langle E_{A,t}(v)\rangle_A - E_{A,t}(X_v)\right|\right] \leq g^{t/2+1}\left(\frac{7}{\sqrt{K}} + 18\widehat{\eta}\right).$$

Let $\mathcal{L}_j$ be the set of nodes at the $j^{th}$ level of Alice's tree $\mathcal{T}_A$. Then if $j$ is even, let

$$\lambda_j = \underset{v\in\mathcal{L}_j}{\mathrm{EX}}\left[\sum_{t\in\{j,j+2,\ldots,R\}}\underset{y,M_t,\mathcal{T}_A}{\mathrm{EX}}\left[\left|\langle E_{A,t}(v)\rangle_A - E_{A,t}(X_v)\right|\right]\right],$$

$$\lambda_{j+1} = \underset{w\in\mathcal{L}_{j+1}}{\mathrm{EX}}\left[\sum_{t\in\{j+1,j+3,\ldots,R-1\}}\underset{y,M_t,\mathcal{T}_A}{\mathrm{EX}}\left[\left|\langle E_{B,t}(w)\rangle_A - E_{B,t}(Y_w)\right|\right]\right].$$

By linearity of expectation,

$$\lambda_j \leq \left(\frac{R}{2}+1\right)g^{R/2+1}\left(\frac{7}{\sqrt{K}} + 18\lambda_{j+1}\right).$$

Solving this recurrence relation, we find that at the root node,

$$\lambda_0 \leq (9R+18)^R\,g^{R^2/2+R}\frac{7}{\sqrt{K}},$$

and similarly for the root of Bob's tree $\mathcal{T}_B$. So in particular, $\mathrm{EX}_{\omega,M_R,\mathcal{T}_i}[\partial_t] \leq \lambda_0 + 2^{-b+1}$ for all $i,t$, where

$$\partial_t = \left|\mathrm{round}\left(\langle E_{i,t}(\mathrm{root}_i)\rangle_i\right) - \mathrm{round}\left(E_{i,t}(\omega)\right)\right|.$$

Now observe that, if we let $\mathcal{W}_{t+1}$ be the distribution over message $m_{t+1}$ in the wannabe case, and let $\mathcal{B}_{t+1}$ be the distribution in the unbounded Bayesian case, then

$$\|\mathcal{W}_{t+1} - \mathcal{B}_{t+1}\|_1 = \frac{1}{2}\sum_{r=1}^{\partial_t/2^{-b}}\frac{2(L-r+1)}{L^2} \leq \frac{\partial_t/2^{-b}}{L} = \frac{\partial_t}{\epsilon}$$

where $\|\ \|_1$ denotes variation distance. So the referee can distinguish the whole conversations with bias at most

$$\frac{1}{\epsilon}\mathrm{EX}[\partial_0 + \cdots + \partial_{R-1}] \leq \frac{1}{\epsilon}\left(\lambda_0 + 2^{-b+1}\right)R$$

since variation distance satisfies the triangle inequality. Therefore, we can achieve the goal of simulation by taking $\lambda_0 \leq \zeta\epsilon/R - 2^{-b+1} \leq \zeta\epsilon/2R$, or equivalently

$$K = \frac{196R^2}{\zeta^2\epsilon^2}(9R+18)^{2R}\left(\frac{4e}{\epsilon}\right)^{R^2+2R}\left(\ln\left(\frac{196R^2}{\zeta^2\epsilon^2}(9R+18)^{2R}\left(\frac{4e}{\epsilon}\right)^{R^2+2R}\right)\right)^{2R} = O\left(\frac{1}{\zeta^2}\left(\frac{11}{\epsilon}\right)^{R^2}\right).$$

∎

# 5  Discussion

> "We publish this observation with some diffidence, since once one has the appropriate framework,
> it is mathematically trivial.  Intuitively, though, it is not quite obvious..." —Aumann [2], on
> his original agreement result

This paper has studied agreement protocols from the quantitative perspective of theoretical computer science.  If nothing else, we hope to have shown that adopting that perspective leads to rich mathematical questions.  Here are a few of the more interesting open problems raised by our results:

- How tight is our $O\left(1/\left(\delta\varepsilon^2\right)\right)$ upper bound?  Can we improve Theorem 6 to show that the discretized standard protocol uses only $O\left(1/\varepsilon^2\right)$ messages, independently of $\delta$?  More importantly, is there a scenario where Alice and Bob must exchange $\Omega\left(1/\varepsilon\right)$ or $\Omega\left(1/\varepsilon^2\right)$ bits to $(\varepsilon, 1/2)$-agree, regardless of what protocol they use?  Recall that the best lower bound we currently know is $\Omega\left(\log 1/\varepsilon\right)$, from Proposition 3.

- Can Alice and Bob $(\varepsilon, \delta)$-agree after a small number of steps, even if the "true" distribution over $\omega$ differs from their shared prior distribution $\mathcal{D}$?  Or is there a scenario where regardless of what protocol they use, there exists a state $\omega$ for which they must exchange $\Omega(n)$ bits to agree within $\varepsilon$ on $\omega$?  (It is easy to construct a scenario where the discretized standard protocol needs $\Omega(n)$ bits for some $\omega$.)

- Can the simulation procedure of Section 4.2 be made practical?  That is, can we reduce the number of subroutine calls to (say) $c^{1/\left(\delta\varepsilon^2\right)}$, or even to a polynomial in $1/\delta$ and $1/\varepsilon$?  Alternatively, can we prove a lower bound showing that such reductions are impossible?

- Can we obtain a better simulation procedure if $\mathcal{D}$ is represented in a compact form, for example a graphical model?

Stepping back, have the results of this paper taught us anything about the origins of disagreement?  As mentioned in Section 1, it is easy to list plausible reasons why people might disagree, Aumann's theorem notwithstanding: indifference to truth, misconstrual, vagueness, dishonesty, self-deceit, mistrust, stupidity, systematic cognitive biases, no priors, different priors, different indexicality assumptions, diagonalization (as discussed in Section 3), communication cost, and computation cost, among others.  But which of these reasons, if any, are fundamental?  In other words, were we forced to identify a single point at which the assumptions of Aumann's theorem diverge from reality, what would it be?

Before we undertook the research described in this paper, we would have said *either* that

(1) imposing reasonable communication and computation bounds is likely to change everything, or

(2) at least one party to any persistent disagreement must be dishonest, irrational, or indifferent to truth.[6]

Today, however, we would make an argument less technical than (1) and less misanthropic than (2): that even in idealized models, *we should not treat agents as initially-identical Bayesian "containers" that later get filled with different experiences.*  In particular, the Common Prior Assumption (CPA) is fundamentally misguided.

Presumably no one would claim that the CPA is empirically true for human beings.  It seems obvious that, when five-year-olds go to Sunday school, they are not updating a shared prior over possible religions conditioned on what their teacher tells them.  Rather, their priors are being "initialized" to some extent.  Furthermore, the existence of a common prior would be astonishing from the perspectives of physics, evolutionary biology, and neuroscience, since nothing in those fields predicts or requires one.  However, as Aumann [3] rightly emphasizes, the question is not whether the CPA is "true" but whether it is a useful idealization.  What we suggest is that, when trying to understand the origins of disagreement, the CPA is *not* a useful idealization.  There are two main reasons for this.

First, the CPA presents difficulties with transtemporal identity.  Are you really the "same" person as you were when you were two months old?  If not, then why must your posterior be obtained by updating

---

[6] Here "indifference to truth" means choosing opinions according to their novelty, social acceptability, value in attracting sexual partners, etc. rather than evidence.

the two-month-old's prior? The difficulties become even more severe if we adopt the many-worlds view of quantum mechanics. For then there are millions of basis states containing beings very much like you. Suppose we fix which one of those beings is "really" you at time $t$; then which one is you at times $t-1$ or $t+1$? Quantum mechanics does not fix an answer; more than that, it does not even fix the probabilities of *possible* answers.[7] That is why Bohmian mechanics and its many variants can all be compatible with quantum mechanics, despite having different equations of motion.

Second, the CPA begs the question of what determines the common prior. Some might argue that human beings' shared genetic heritage causes them (or rather, should cause them) to share a prior. But if your prior is to fix your initial opinions about *everything*, then it must assign a probability to your future experiences being consistent with those of (say) a five-legged extraterrestrial. Presumably that probability decreases dramatically once you condition on the indexical fact of your humanity. But it ought to start nonzero and stay nonzero, for instance because of quantum fluctuations. This raises a question: why shouldn't your prior *equal* the extraterrestrial's? After all, the extraterrestrial has to assign a probability to *its* future experiences being consistent with yours—and at a hypothetical time before either of you knows who "you" will become, why should the two of you reason differently? We can similarly imagine beings governed by different laws of physics; and these, too, should share our prior. It follows that "the" common prior, if it exists, is not determined by anything in our genetic makeup or even the physical world.

This leaves the possibility that mathematics or logic could determine the common prior. Along these lines, Schmidhuber [17] has advocated a prior in which the probability of any sequence of experiences $x$ is proportional to $2^{-K(x)}$, where $K$ is the Kolmogorov complexity of $x$—that is, the length of the shortest computer program that outputs $x$. This idea has several problems, though. First, our actual experiences seem to have gratuitously high Kolmogorov complexity. Believers in the Kolmogorov prior are forced to say, without evidence, that this is an illusion. Second, why should we use Kolmogorov complexity, rather than (say) time-bounded Kolmogorov complexity, or perhaps the length of the shortest program that outputs $x$ given an oracle for the halting problem? Third, whenever we wish to compare the probabilities of a few "equally complex" events, the probabilities will depend less on the events themselves than on our choice of programming language, so we face another arbitrary choice.

So it seems that a common prior would be independent of the physical world and even of mathematics, yet would somehow be readily available to and unquestioningly accepted by every rational agent. Agents equipped with this prior would live a 'preprogrammed' existence, meaning that they would never change, only conditionalize. We have argued that this picture of the world presents serious intrinsic problems, even setting aside its naked implausibility. So perhaps the common prior should be jettisoned with the ether.

But is there any principled basis for prior differences, then? Consider Shakespeare's Julius Caesar, debating whether to venture outside on the Ides of March. From his dismissals of omens, we know that Caesar bases his final decision on a *belief* that he will not be in particular danger, rather than just a *preference* for risky actions. Yet the process of reaching the belief seems to have nothing to do with conditioning on evidence—or rather, it starts after the conditioning is already done. Our proposal is to view the process as that of Caesar *choosing his prior, and thereby choosing what sort of person he is*. In other words, Caesar assigns a low prior probability to his getting killed, for the sole reason that had he assigned a high one, he would no longer be Caesar but someone else.[8] On this view, not only can Alice and Bob have different priors because they are different people, but the fact that they have different priors is a large part of what *makes* them different people, rather than the same person filling two pairs of shoes.

In saying this, we are not taking the relativist stance that any prior is "rational" for the sort of person who would hold that prior. If no priors are objectively more rational than others, then the word "rational" is meaningless, since there exists a prior to justify essentially any belief. But the question remains: is the number of rational priors exactly one? We have already seen an argument of Hanson [12] that it should be, based on the concept of a "pre-prior" (that is, a prior over all possible priors). Why should Alice give her own prior any more weight than Bob's? Our response is simply to point out that there is a tremendous

---

[7]What quantum mechanics does fix are the probabilities of possible outcomes of a measurement. But those probabilities will only be meaningful to you if you are not part of the system being measured.

[8][D]anger knows full well
That Caesar is more dangerous than he:
We are two lions litter'd in one day,
And I the elder and more terrible...
—*Julius Caesar*, Act 2, Scene 2

gap between empathizing with someone else's perspective and adopting it, or between calculating what your expectation would be under someone else's prior and willing that expectation to be yours. No matter how long she talks to Bob, in the end Alice must confront the irreducible fact of her individuality. As Clarence Darrow famously put it, "I don't like spinach, and I'm glad I don't, because if I liked it I'd eat it, and I just hate it."

# 6 Acknowledgments

# References

[1] M. Allais and O. Hagen (eds). *Expected Utility Hypotheses and the Allais Paradox*, Dordrecht, 1979.

[2] R. J. Aumann. Agreeing to disagree, *Annals of Statistics* 4(6):1236–1239, 1976.

[3] R. J. Aumann. Reply to Gul, *Econometrica* 66(4):929–938, 1998.

[4] T. Cowen and R. Hanson. Are disagreements honest?, submitted, 2003.

[5] P. Cross. Not *can* but *will* college teaching be improved?, *New Directions for Higher Education*, 17:1–15, 1977. Cited by Gilovich [9].

[6] D. P. Dubhashi and A. Panconesi. Concentration of measure for computer scientists, draft at http://www.cs.unibo.it/~pancones/master.ps.

[7] J. Earman. Old evidence, new theories: two unresolved problems in Bayesian confirmation theory, *Pacific Philosophical Quarterly* 70:323–340, 1989.

[8] J. D. Geanakoplos and H. M. Polemarchakis. We can't disagree forever, *J. Economic Theory* 28:192–200, 1982.

[9] T. Gilovich. *How We Know What Isn't So*, Free Press, 1993.

[10] F. Gul. A comment on Aumann's Bayesian view, *Econometrica* 66(4):923–927, 1998.

[11] R. Hanson. Disagreement is unpredictable, *Economics Letters* 77(3):365–369, November 2002.

[12] R. Hanson. Uncommon priors require origin disputes, submitted, 2004.

[13] R. Hanson. For savvy Bayesian wannabes, are disagreements not about information?, *Theory and Decision* 54(2):105–123, 2003.

[14] E. Kushilevitz and N. Nisan. *Communication Complexity*, Cambridge, 1996.

[15] R. Parikh and P. Krasucki. Communication, consensus, and knowledge, *J. Economic Theory* 52:178–189, 1990.

[16] S. Russell and P. Norvig. *Artificial Intelligence: A Modern Approach* (2nd edition), Prentice Hall, 2002.

[17] J. Schmidhuber. A computer scientist's view of life, the universe, and everything, in *Foundations of Computer Science: Potential - Theory - Cognition* (C. Freksa, ed.), pp. 201–208, Springer, 1997.

[18] A. M. Turing. Computing machinery and intelligence, *Mind* 59:433–460, 1950.

[19] A. Tversky and D. Kahneman. Judgment under uncertainty: heuristics and biases, *Science* 185:1124–1131, 1974. See also D. Kahneman, P. Slovic, and A. Tversky (eds.), *Judgment under Uncertainty*, Cambridge, 1982.