Thinking Is Believing

Abstract

The idea that people can entertain propositions without believing them is widespread, intuitive, and most probably false. The main goal of this essay is to argue against the claim that people can entertain a proposition without believing it. Evidence is presented demonstrating that we cannot withhold assent from any proposition we happen to consider. A model of belief fixation is then sketched and used to explain hitherto disparate, recalcitrant, and somewhat mysterious psychological phenomena. The proposed model is one where beliefs are the automatic output of a computationally null belief acquisition reflex. In short, the model holds that the mere tokening of a mentally represented truth apt proposition leads to immediately believing it. The essay concludes by considering some consequences that the proposed model of belief acquisition has for our concept of rationality.

Thinking Is Believing

"Keep Moving and Faith Will Come"-Jean le Rond D'alembert

1. Belief Fixation, Doxastic Deliberation, and Rationality

Suppose that you have just stubbed your toe on a rock. If you are like some people, you will, at least momentarily, be angry *at the rock*. Even though you might know that the rock is not an appropriate recipient of your reactive attitude, often enough you can't help but be angry at it.

Cases like the one just described are common: we frequently feel emotions that are, even by our own lights, rationally groundless. But we tend to assume that this is not equally true of our beliefs. If I ask you to please not believe what I'm about to say (because, e.g., I'm merely parroting someone else's falsehood), it seems plausible that you will be able to not believe what I'm about to say. If I tell you that I'm about to read a list of sentences, all of which are false, and then I read the sentences, it seems plausible that you would not automatically believe these sentences in the way that you may, for example, automatically get excited when hearing of a rare and tantalizing opportunity.

However, in what follows I will argue that this plausible assumption is false: just as we get angry with the rock while knowing full well that it's not an appropriate object of our anger, so too we believe what people say even when we know that what they are saying is false.¹ That is, just as emotions are insensitive to our background beliefs, so too is belief formation initially insensitive to our background beliefs.² More specifically, I will argue for the claim that whenever we entertain a proposition we automatically believe that proposition.

The idea that we can contemplate a proposition without believing it has been accepted in philosophy since at least the time of the Stoics (Long and Sedley 1987, p. 438–61), and remains widespread in contemporary debates concerning everything from modularity theory to epistemology. To take one representative example, Jerry Fodor says,

To a first approximation, we can assume that the mechanisms that affect [the fixation of perceptual belief] work like this: they [central systems] look simultaneously at the representations delivered by the various input systems and at the information currently in memory and they arrive at a best (i.e., best available) hypothesis about how the world must be, given these various sorts of data. (Fodor 1983, p. 102)

Note that this story assumes that our central systems examine how different propositions are analyzed in light of our background beliefs. Fodor assumes that background beliefs interact with propositions we entertain because he thinks that belief fixation is a conservative, gradual process that (ideally) takes into account all the relevant data in one's information store before assenting to any proposition.³ Here Fodor's view is quite indicative of the field at large. Belief fixation is hypothesized to be a slow, conservative process, in part, to allow for the idea that we have the ability

¹ Or a different formulation for those who think that you can't believe that p and know that not-p: we will believe someone's testimony even while knowing that the testifier claims to be lying.

² It's plausible that the initial process of belief formation is even more encapsulated than the elicitation of emotions. I will argue that belief formation is completely informationally encapsulated, so much so that it can be fruitfully seen as completely reflexive.

³ Hence Fodor writes things like "the fixation of perceptual belief is the evaluation of such hypotheses in light of the *totality of background theory*" (emphasis added; Fodor 1990, p. 248).

to contemplate the truth of a proposition before assenting to that proposition. This intuitive view is at odds with a theory in which any proposition that is entertained is simply automatically and reflexively believed. So, if it were true that belief fixation is reflexive (such that every thought that is contemplated is believed) and interacts with no background information, then it would be a very interesting and surprising fact about the mind.⁴

The consequences of such a radical departure from the standard view would extend far beyond the topic of belief fixation. The ability to withhold assent from propositions that we entertain is a crucial part of our picture of an important variety of impartial doxastic deliberation: the ability to impartially consider propositions while suspending judgment. When first encountering a proposition, we take ourselves to be able to consider it while remaining neutral as to its truth. Furthermore, impartial doxastic deliberation is integral to our conception of what it is to be a person because we take people to be paradigmatically rational creatures, and impartial doxastic deliberation appears to be a necessary condition on rationality. If we found creatures that couldn't help but believe any idea that they entertained, we would be inclined to regard them as massively irrational. Sadly, we seem to be such creatures.

A critique of rationality stemming from our inability to impartially deliberate differs from the contemporary "rationality wars" criticisms (Samuels et al. 2002). Recent decades have brought heated debates over how rational people are, but these debates cluster around whether people tend to answer some particular problem correctly. One needn't look hard to find claims that people are irrational because they, for example, fall prey to cognitive illusions, use fast and frugal heuristics, let emotions dictate their moral reasoning, and so forth. Throughout these debates, a cornerstone of our rationality has remained beyond critique: our ability to entertain propositions without believing them.⁵ This ability has received scant attention and has endured very few serious critiques. Yet when one looks closely at our actual doxastic capacities, the picture that arises is surprising and quite epistemologically troubling.

If the theory I propose is correct, then we will have to reconsider the nature of doxastic deliberation and whether we are able to engage in it. This is because if the proposed theory is correct, then truly impartial doxastic deliberation is impossible. Consequently, the theory of belief fixation defended here is somewhat radical and unintuitive. My goal is not to establish the truth of the theory beyond a doubt, rather my aim is more modest: to convince you that it is a plausible model of our cognitive architecture that demands further investigation.

But before we get there, let's first be clear about the notion of belief with which we'll be working. The notion of belief that is operative throughout this paper will be the quotidian one that is operative in the cognitive sciences, with belief understood as a relational, gradable, functional state. This notion of belief, being gradable, allows that one can believe things to stronger or weaker degrees. For current purposes, belief will not be understood as merely a binary relation where one either does or does not believe that P.⁶ Rather, belief will be understood similarly to the way one understands credences.⁷

⁴ A note about the scope of the claim: the arguments for the vast majority of the essay cover initial belief fixation, but remain neutral as to what the capacities underlie cases where one is reconsidering a proposition already believed, though this topic will briefly arise in the final section.

⁵ For example, when arguing over whether the use of heuristics is ecologically rational, all parties assume that the information that heuristics process can initially be rejected. It is only after the acceptance of information that the question of the efficacy of our information processing techniques arises.

⁶ I say 'merely' because the gradable notion still allows for some binary notion of belief.

⁷ Thus, one can interpret my theory as stating that whenever you entertain a proposition, you raise your credence in that proposition. How high is credence raised? Is it to a high degree or just to a non-zero degree? To a first approximation,

Additionally, the essay will assume a certain amount of 'realism' about belief. This project is of a piece with the search for a state in cognitive science that shares a certain "spiritual similarity" (Fodor 1987b) with the folk psychological notion of beliefs. As such, the essay works toward building a psychofunctional theory of beliefs. But one needn't buy into any psychofunctional theory in order to assess the arguments in the paper. For the present purposes, all that one needs to assume is that beliefs are causally efficacious states—ones that, for example, are caused by perception, serve as the premise of inferences, and interact with desire to cause behavior. Not much else is assumed about belief, at least nothing that will affect the arguments that follow.

In what follows, I compare two theories of belief fixation. Ultimately, I argue that one of these theories is false and that the other theory can unify and explain a plethora of seemingly disparate phenomena and so should be taken quite seriously. But for now, let's peruse a ubiquitous and influential theory of cognitive architecture: the Cartesian theory of belief fixation.

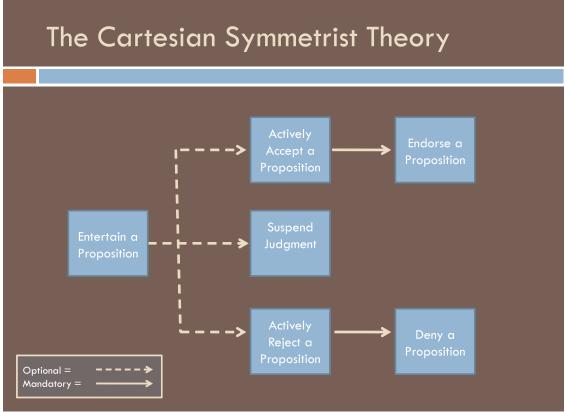
2. Theories of Belief Fixation

The methodical withholding of assent is part of a venerable epistemological tradition: if surety is what one desires, then one should be skeptical of what one thinks, waiting for the ideas that pass through one's mind to be 'clear and distinct,' or at least well justified. Surety was Descartes stated goal in the *Meditations* (1641/1988). But it's worth asking: when Descartes was sitting beside the fire contemplating which propositions to believe in, what was he actually trying to do? He was attempting to first *entertain* an idea, then *contemplate* its truth, and finally *decide* what to assent to and what to *withhold* belief from. Descartes' attempt presupposed a serial model of belief fixation, according to which one first entertains a proposition.⁸ This type of serial model presupposes that a)

the credence is raised to a level that would generally produce behavior (in combination with the appropriate desires). Presumably a belief with a credence of .0001 won't produce any behavior; on the other hand, a belief needn't have a credence of .9 in order for the belief to have behavioral consequences. I take it as an open empirical question how high one's credences have to be for a belief to regularly eventuate in behavior. The operative claim in the text is that entertaining propositions causes one's credences to go at least that high. This said, there will be little talk of credences in what will follow, for my preferred analysis of what credences are is cashed out in terms of resistance to disconfirming evidence and not something akin to betting procedures. For a more fleshed out explanation of this claim see Mandelbaum 2010. For more on the specific properties applicable to belief, see section 4.4. ⁸ Although this scenario admittedly paints Descartes with a broad brush, some relevant literature has interpreted Descartes as attempting the project I sketch out (e.g., Gilbert 1991, Huebner 2009). Nevertheless, there are some reasons to believe that Descartes actually wasn't a Cartesian in this sense. Some historians like Alan Nelson (personal communication) interpret Descartes' epistemic methodology as such: assume Descartes wants to assess the truth of the proposition that Santa Claus exists. Call this proposition S. Descartes' first step in assessing S is to token the thought WITHHOLD ASSENT FROM S (actually Nelson's take on this seems to be that the first step is to token the thought: THINK WITHHOLD ASSENT FROM S; I'll ignore this element, which strikes me as regress prone.) The next step is to think of situations that would entail the falsity of S-for example, imagining an empty North Pole. The reason we think of an empty North Pole as opposed to thinking NOT S is that Descartes doesn't believe one can just think of negation as such. Nelson's Descartes holds a variation on the view that I'm promoting; he believes that people believe everything they think because they do not have the ability to withhold assent. Rather, what people can do instead is constantly have a belief swamped by a contrary belief. In essence, this reading of Descartes interprets the withholding of assent as a type of thought suppression: your belief that S is weak if it immediately leads to a different belief, and it is super-weak if it leads to a different belief that would entail the falsity of S (ironically, if this reading is right then my analysis of credences is very similar to Descartes). A strong belief is a belief that doesn't automatically lead to a second

belief, which destroys our consciousness of the first belief. So perhaps Descartes wasn't a Cartesian in the sense expressed in the main text. That doesn't really matter because an overwhelming majority of contemporary philosophers and cognitive scientists are. If one would like they can substitute Pollock or Fodor (or anyone else who has the modular/central systems distinction) in for Descartes (see Pollock 1986 and Fodor 1975, 1983, 1998).

the faculty of entertaining a proposition is a separate faculty from the faculty of believing a proposition; and that b) the workings of the former faculty are prior to the workings of the latter (see Figure 1). These assumptions are at the heart of the serial model of belief fixation, which I (following Gilbert) will term 'the Cartesian theory of belief fixation.'



(Figure 1.) (NB: The dotted lines represent optional links, and solid lines necessary links)

What I am here calling the "Cartesian Theory" consists of the following claims:

- 1) People have the ability to contemplate propositions that arise in the mind, whether through perception or imagination, before believing those propositions.
- 2) Accepting and rejecting a proposition use the same mental processes, and consequently, should be affected by performance constraints in similar ways.⁹ I will sometimes refer to the Cartesian position as a 'symmetrist' position, because it treats accepting and rejecting symmetrically.

⁹ I'll use the phrases 'accepting a proposition' and 'believing a proposition' interchangeably; likewise for 'rejecting a proposition' and 'disbelieving a proposition.' No doubt, in the present climate doing so is controversial (see, e.g., Velleman 2000; Stalnaker 1984; Cohen 1992; Bratman 1992; Tuomela 2000; van Fraasen 1980), but I haven't the space to argue for such usage. The idea that the states discussed are indeed beliefs and not some form of 'mere acceptances' (and in fact that the theory proposed here can serve as a model for a full-fledged psychofunctional theory of belief) is discussed in section 4.4.

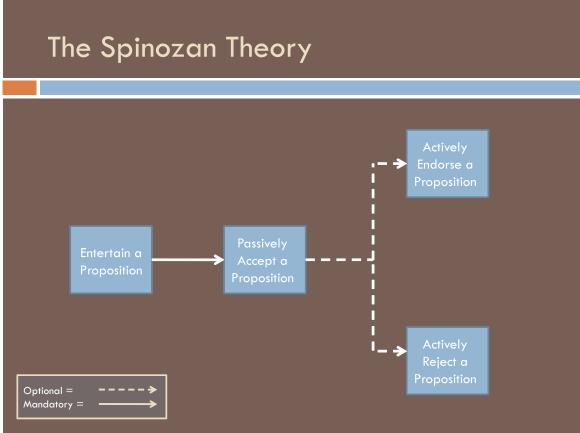
3) Forming a belief is an active endeavor. Since accepting a proposition and rejecting a proposition are underwritten by the same mental processes, rejecting a proposition is also an active endeavor.¹⁰

The Cartesian theory is intuitive, widely accepted, and rarely, if ever, argued for. It is assumed throughout many areas of both philosophy and psychology (e.g., Quine, 1960; Milgram 1974; Fodor 1975, 1983, 1998; Dennett 1987; Pylyshyn 1989; Ford and Pylyshyn 1996; and Cooper 2007). Moreover, the serial view of belief fixation that the theory presupposes underwrites our conception of impartial doxastic deliberation. However, there is reason to suppose that the Cartesian theory is more venerable myth than hard fact, and so we can be thankful that the Cartesian view isn't the only available theory of belief formation. Spinoza (1677/1991) had a competing view of belief formation, according to which contemplating a proposition's truth coincides with assenting to a proposition. In lieu of the Cartesian view, I propose a version of a Spinozan theory of belief fixation, one in which tokening¹¹ an idea is sufficient for believing that idea.¹² In the Spinozan theory, one automatically and passively accepts whatever ideas one tokens, and only after the initial acceptance can one effortfully reject one's newly acquired belief (see Figure 2).

¹⁰ Suspending one's judgment can be either active (as when one decides that there is not enough information to decide one way or the other) or passive (as when one's head becomes momentarily attached to a fast moving brick thus making the decision process moot). According to my view, even a fast moving brick cannot derail one's passive assent.

¹¹ I use 'tokening' because it strikes me as the most neutral and general verb for covering the category of heterogeneous mental acts addressed by my theory. These acts include understanding, entertaining, contemplating, and related activities. If you are having trouble envisioning the thesis, assume that there is a language of thought (LOT). My thesis is that every time a truth-apt sentence is tokened in one's LOT, one believes that sentence.

¹² Having no metaphysical axes to grind, I don't particularly care whether we believe propositions or whether we believe ideas. I will thus use these descriptions interchangeably. The difference between the two does not affect my main points, but if you prefer you can substitute one for the other throughout.



(Figure 2) (NB: The dotted lines represent optional links, and solid lines necessary links)

What I am here calling the "Spinozan Theory" consists of the following claims:

- 1) People do not have the ability to contemplate propositions that arise in the mind, whether through perception or imagination, before believing them; that is, because of our mental architecture, it is (nomologically) impossible for one to withhold assent from propositions that one tokens.¹³
- 2) Accepting a proposition is accomplished by a different system than rejecting a proposition. Because different systems are at play, the processes of accepting and rejecting should be affected by performance constraints in different ways. I will sometimes refer to the Spinozan position as an 'asymmetrist' position, because it treats accepting and rejecting asymmetrically.¹⁴

¹³ The impossibility claim is there to rule out that one has a heuristic that makes people tend to believe what they perceive. This is not to suppose that people actually perceive propositions. Rather, a phrase like 'believing what you perceive' is shorthand for 'believing what normally comes to mind when you perceive X.' The idea behind this is quite tame: many perceptual situations lead to the corresponding automatic tokening of thoughts.

¹⁴ The reader may see a certain affinity between the current claim and certain so-called dual process theories of reasoning (e.g., Frankish 2004). However the superficial resemblances are misleading. I will argue that the unconscious states (putatively 'system 1' states) count as beliefs because they are the ones that act in law-like manners and cause behavior, whereas the conscious states (putatively 'system 2' states) do not seem to behave in any law-like fashion and so are unfit for psychofunctional theories. That said, no doubt this quick note gives short shrift to an important issue, just one that is beyond the scope of the current endeavor.

3) Forming a belief is a passive endeavor. However, rejecting a proposition is an active and effortful mental action, which can only happen after a belief has been acquired. Consequently, one can effortlessly form new beliefs while being mentally taxed, but rejecting an already held belief will become more difficult the more mentally taxed one is. For the Spinozan, every proposition that is entertained is necessarily accepted, but every proposition that is accepted is not necessarily endorsed.¹⁵

My version of the Spinozan theory takes on an extra commitment on an issue about which the Cartesian theory is agnostic: the relation between rejection and negation. Because the Spinozan theory dictates that accepting and rejecting are subserved by different mental processes, it's natural for such a theory to give some idea of what rejecting is. As opposed to analyzing rejection in terms of negation, I follow Price (1990) in analyzing negation in terms of rejection.

4) To negate a thought is, in part, to reject it.

The Spinozan view is consistent with the idea that, phylogenetically speaking, cognition grew out of perception. Since our perceptual faculties were by and large veridical, the cognitive faculties that first evolved just took the deliverances of perception at face value. The ability to reject information was a later evolutionary development. The Spinozan view sees this phylogenetic story mirrored in our current cognitive processing: just as the ability to accept information arose before the ability to reject information, so too do we automatically accept information before being able to reject it.

Now for a few non-obvious consequences of the Spinozan view. The Spinozan sees acceptance and rejection as different propositional attitudes. However, the logical relations between these attitudes can differ based on one's tastes. For example, a Spinozan who denied property 4 could hold that accepting not-p does not entail rejecting p, though a Spinozan of my variety has to allow the entailment. However, no Spinozan can allow that one can reject p without also accepting p. Consequently, any Spinozan will predict that people (de facto) believe many contradictions.¹⁶

As per property 3, exercising the faculty of rejection is effortful. However, the Spinozan does not predict that rejection is effortful merely because it is the second step in the system; rather, rejection is effortful because the connection between acceptance and rejection is not mandatory. For our current use, all mandatory processing connections should be thought of as effortless and all non-mandatory processing connections as effortful. This is because all mandatory connections are automatic, like a reflex.

As per property 4, negating involves rejecting. Since negating involves rejecting, and since rejecting is effortful, negating is effortful too. Thus, negative sentences/thoughts should be more difficult to process (e.g., take longer and be more error prone) than affirmative sentences/thoughts. Furthermore, because negation involves rejection and because one can only reject complete propositions, the Spinozan theory predicts that negations can only be processed after a complete

¹⁵ For current purposes, endorsing a proposition is something that happens at the person level. One consciously chooses what to endorse, whereas accepting needn't be conscious nor volitional. In the Spinozan ontology, denying is the negative complement to endorsing (also a person-level phenomenon), whereas rejecting complements accepting (and both are sub-personal phenomena).

¹⁶ Of course, this does not entail that people will *assert* contradictions. What one asserts is tied to what one endorses and on this picture endorsements are a species of judgment, not belief (for more on the relations between endorsing/denying and believing/rejecting, see the end of section 4.3). For the Spinozan judgments are a person-level phenomenon whereas beliefs are subpersonal.

(affirmative) proposition has been formed. As a consequence, when processing negative clauses negations will be processed last.¹⁷

As a consequence of properties 2 and 4, the Spinozan view not only treats acceptance and rejection asymmetrically, but also treats negation and affirmation asymmetrically. The Cartesian position officially makes no predictions about negations, but it's quite natural for the Cartesian to be a symmetrist about negation as well as belief.

The big picture: on the Spinozan view, any propositional thought¹⁸ one tokens one thereby believes. Only after a belief is acquired can decision procedures be brought to bear on the belief. If one tokens a dubious proposition, one can effortfully attend to the proposition and reject it. Further contemplation can toggle the strengths of these beliefs, reducing the strength of the affirmative belief and raising the strength of the negated counterpart.

The Cartesian and Spinozan theories create quite different predictions. For example, if the Cartesian view is right then we should be able to dismantle the belief-fixating process after the understanding has happened but before the believing (or disbelieving) has occurred. In such a case the Cartesian view predicts that the system will be agnostic about the truth of the proposition. Consequently, since cognitive load is a disabling performance constraint, the Cartesian theory predicts that deciding about the truth of a proposition should not normally occur under cognitive load. Additionally, because the Cartesian theory treats assenting and rejecting identically, it predicts that cognitive load will affect both processes identically.

In contrast, if the Spinozan view is right, then the belief-fixating process can be dismantled by invoking some performance constraints prior to rejecting a proposition, but never before accepting a proposition. Because the Spinozan theory posits that believing is reflexive, believing should occur even when one is under cognitive load. That is to say, since the Spinozan view treats accepting and rejecting differently, with rejection being effortful, it predicts that load should only affect rejecting a proposition, not assenting to it.

We will return to these predictions throughout the paper. For now, let's turn our attention to some evidence that should make us quite wary of the Cartesian view.

3. The Case against the Cartesian

One needn't look very hard to find evidence that cast doubt on the Cartesian view.¹⁹ Space constraints prevent me from recounting every datum, so the following reviews are meant to be more illustrative than comprehensive. Each piece of evidence below is evidence that the Cartesian theory cannot account for and the Spinozan theory can. Some are more suggestive than others, but in combination they make for a daunting challenge to the Cartesian theory.

¹⁷ Importantly, the claims in the text regarding negation do not pertain to syntax; rather, they pertain to understanding negation. Additionally, the claims about negation apply to propositions, not necessarily sentences. So, for example, the theory handles embedded negations, like the one in 'John believes that Jesse is not a communist,' by stating that the negation is processed after the clause sans negation (i.e., 'Jesse is a communist') is processed, rather than after the entire sentence sans negation ('John believe that Jesse is a communist') is processed. The theory is supposed to hold over both sentential ('I do not regret going to your party') and constituent negation, but the negation processed before the rest of the outlying sentence.

¹⁸ Again, 'propositional thought' means 'thought that is truth-apt.'

¹⁹ For example, see Festinger and Maccoby (1964), Gilbert et al. (1993), Gilbert (2002), Anderson et al. (1980), and Kruger (1999). The data against the Cartesian view have remained invisible yet well circulated. For a much more comprehensive review of the data, see chapter 2 of Mandelbaum (2010).

3.1 Memory Asymmetries between Truths and Falsehoods

The quintessential anti-Cartesian experimental paradigm is one that exploits asymmetries in people's memory of truths and falsehoods. In a typical experiment, participants are asked to partake in a learning task while intermittently under cognitive load and are then tested about what they learned. (E.g., they would learn about fictional criminal acts and then decide on an appropriate sentence or they would learn about random people and then assess their mental states.) In one (more or less arbitrarily chosen, though typical) experiment participants were asked to learn nonsense word meanings. They watched a computer screen where sentences of the form "An X is a Y" appeared, in which the "X" was a nonsense word and the "Y" was a word in English (e.g., "A suffa is a cloud," from Gilbert et al. 1990). Right after participants read a sentence the screen flashed either the word 'true' or the word 'false,' indicating whether the previous statement was accurate or not. Participants were also told to be on guard for a tone that would occur; the tone would occasionally bellow and when it did the participants were to push a button as soon as possible. The tone was introduced in order to induce cognitive load. During the critical trials, participants read six true and six false claims. While reading four of these claims (two true, two false), the participants were interrupted by the tone. At the end of the trials the sentences were then turned into questions (e.g., "Is a suffa a cloud?") which the participants then answered.

The Cartesian view predicts that the tone task should affect both true and false statements equally since although contemplation has occurred, the participants haven't yet had the time to integrate the information properly because of the cognitive burden brought on by the tone task. The Spinozan view predicts that during interrupted trials participants should mistake false claims as true, but not true claims as false, the reason being that the belief-fixation system's processing gets shut down by the cognitive load *after* comprehension but *before* rejection. The Cartesian view predicts incorrectly: the added cognitive load made participants reliably encode true statements as true, but consistently incorrectly encode false statements as true.²⁰

This type of asymmetry can be seen throughout the literature: a person put under cognitive load is apt to remember statements that they are told to be false as true but not statements they are told to be true as false. This robust asymmetry helps to confirm the second and third properties of the Spinozan theory. The experiment above displays that accepting a proposition (i.e., remembering the proposition as true) comes much easier than rejecting a proposition (i.e., remembering the proposition as false). Accepting is easier because it is a passive process, whereas rejecting is an active one. The added cognitive load helps to shortcut the active rejection, but does not interfere with passive acceptance because the passive process is automatic and load does not affect a reflex. Compare how counting backwards from one hundred by increments of five would affect *seeing* a crossword puzzle vs. *completing* the puzzle. The former will not be affected while the latter will be greatly affected. Rejecting a proposition is more like thinking than seeing, while accepting is more like seeing than thinking.

In sum, the observed asymmetry can be predicted if we assume what the Spinozan view asks us to: that when propositions are initially processed they are encoded as true and can only subsequently be marked as false. Not only does the Cartesian view miss the asymmetry between acceptance and rejection, but it misses that acceptance is automatic. The Cartesian view predicts that

²⁰ Participants answered correctly on the true statements 55 percent of the time when uninterrupted and 58 percent of the time when interrupted, but participants answered correctly on the false statements 55 percent of the time when uninterrupted but only 35 percent of the time when interrupted.

load should shut down acceptance because it assumes that acceptance is active. Contra the Cartesian view, load seems to increase, not decrease, people's disposition to accept propositions.

3.2 Belief Perseverance

Another telling set of experiments comes from the literature on belief perseverance in the face of experimental debriefing. For example, in a typical experiment, an experimenter asks participants to read a collection of suicide notes and to sort the real ones from the fakes. In Ross et al. (1975), participants encountered twenty-five pairs of notes and were told that one note from each pair was a real note, the other a fake. After seeing each pair participants would judge which note was real and which fake and were then given feedback on their performance. After receiving the feedback the participants were (partially) debriefed. During the debriefing the participants were told that all the feedback they received was fictitious, it being arbitrarily determined beforehand regardless of the participants' responses. After the debriefing the participants were asked to estimate both how many times they actually answered correctly and how many correct answers an average person would give. Sadly, the information in the debriefing session did not affect participants' opinions about their ability: if the participant originally received positive false feedback (e.g., twenty-four out of twenty-five correct), they believed that they were better than average at the task, and if they received negative false feedback (e.g., seven out of twenty-five correct), they believed that an average at picking out real suicide notes from fake ones.

The aforementioned experiment is generally not taken to illuminate anything about belief acquisition per se. It seems that the participants formed their beliefs in a reasonable way, based on the experimental feedback. Once they are told that the feedback was non-veridical they may just have trouble updating their beliefs. Perhaps beliefs are 'sticky,' in that once one has a belief, that belief is hard to relinquish. If so, then the debriefing effect wouldn't tell us anything about belief acquisition per se, but rather belief perseverance.

But what happens if the people are briefed before they take part in the study and receive false feedback (call such a technique 'prebriefing')? What if before sorting the notes they are told that the feedback they are about to receive is bogus? The Cartesian view predicts that if we tell people beforehand that what they are about to read is false, and they have no reason to distrust what we tell them, then, ceteris paribus, they will approach the stimuli skeptically, withholding the formation of any beliefs about their ability if those beliefs are based on the bogus data. In contrast, the Spinozan view predicts that since people believe everything they token, they'll be stuck believing propositions that they encounter even if they know beforehand that they are false.

As predicted by the Spinozan view but not the Cartesian view, prebriefing the participants does not impact the participants' judgments about their ability. Wegner et al. (1985) replicated the Ross study except the participants were told *prior* to the task that the feedback would be dubious. Yet, even after the explicit prebriefing the participants continued to behave as if the feedback was veridical. They were unable to reject the feedback they received, even though they knew the feedback was bogus. These perseverance effects are easily explicable on the Spinozan view: the knowledge of the feedback persists because the participants automatically believe the feedback when they hear it, even though they know the feedback is false. Since they are engaged in a relatively fast-paced experiment, the participants lack the mental energy to override the false beliefs. The prebriefing effect helps to verify the first and third properties of the Spinozan theory. Equally important, these persistence effects are anomalous on the Cartesian theory, particularly casting doubt on the first property of the view.

3.3 Personality Metrics

Evidence for the Spinozan theory comes from a wide array of sources. The next example comes from an unlikely place: personality psychology. When studying personality psychology, researchers often present participants with a list of personality attributes and ask them to evaluate how much each attribute describes their personality. Consider a personality survey where participants are given twenty statements and are asked to answer, for each statement, whether the statement applies to them or not. The participants answer 'yes' when the statement applies to them and 'no' when it doesn't. For ten of the questions an answer of "yes" corresponds with being an introvert and for the other ten an answer of "yes" corresponds with being an extrovert. On such a scale a 'perfect' introvert would be one who answered "yes" to the ten introversion questions and "no" to the ten extroversion questions, while a perfect extrovert would reverse their answers. When using such methods researchers have found that their data are sometimes compromised by 'yea-sayers'; that is, people who are apt to respond affirmatively to whatever question they are asked. For example, a perfect yea-sayer would respond to the aforementioned study by answering "yes" to all twenty questions, thus confounding the personality metric. The perfect 'nay-sayer' would reverse the pattern of the perfect yea-sayer.

If negations are processed subsequent to affirmations, as the Spinozan view would have it, then we should expect nay-saying to take more energy, and thus more time, than yea-saying. This is because, for the Spinozan, the first stage of encoding/accepting is passive and effortless whereas the second stage of rejecting is active and effortful. Thus, the Spinozan nay-sayer would have to first encode the property as applying to them and would then have to go back and reject the property, whereas the acquiescing yea-sayer would just need to passively encode the property. Additionally, the Spinozan view predicts that if people are put under cognitive load while answering one of these personality metrics, then yea-saying should increase relative to an administered personality metric that lacks any load-inducing element. This is because the load makes the participant more cognitively enervated and therefore less able to summon the energy to reject the proposition. In contrast, the Cartesian symmetrist position predicts that because accepting and rejecting are products of the same underlying process, yea-sayers should take the same amount of time as nay-sayers to complete the survey, and both should be equally affected by cognitive load.

Both Spinozan predictions were borne out in Knowles and Condon (1999). In their study, participants received a counterbalanced one-hundred-item personality questionnaire and had their reaction times measured. Yea-sayers were operationalized as those who answered affirmatively on fifty-three or more of the items, and nay-sayers as those who answered affirmatively on forty-seven or fewer of the items. The middle group counted as appropriate responders. The response times for yea-sayers were significantly quicker than the response times for either of the other two groups. In fact, when we look closer we can see that the response patterns perfectly conforms to the Spinozan hypothesis, with yea-sayers taking longer than appropriate responders, who took longer than nay-sayers. This response pattern is directly at odds with the first Cartesian prediction.

Cognitive load also affects yea-saying in the way predicted by the Spinozan, but not Cartesian, hypothesis. In a related study participants were split into two groups, both of which were asked to answer twenty counterbalanced personality questions. Intermittent music was playing in the background for both sets of participants. One set of participants was put under cognitive load by being asked to listen to the music and distinguish notes that came from the piano from those that came from other instruments. The non-loaded group heard the same sounds but wasn't asked to attend to them. The group under load was significantly more apt to answer affirmatively to the questionnaire, thus confirming the second Spinozan prediction and disconfirming the second Cartesian prediction. A theory that sees acceptance as passive and automatic but rejection as active and effortful, as the Spinozan theory does, predicts that load affects nay-saying differently than yeasaying because only the former is active and effortful, thus only the former is a viable candidate for being affected by load. A theory that sees acceptance and rejection as part of the same underlying active mental process, as the Cartesian theory does, cannot explain such findings.²¹

In sum, there are some strong reasons to be skeptical of the Cartesian theory of belief acquisition. The evidence that we've encountered so far doesn't just tell against the Cartesian theory, it also provides support for the Spinozan model. But there are still looming objections to the Spinozan model. In the next section I will present the two most popular objections.

4. Objections and Replies

4.1 The Gullibility Heuristic

Instead of using the presented evidence to support an architectural-processing story (viz. the Spinozan theory), one may be inclined to see it as support for an explanation that appeals to a pretty simple heuristic. I will call such a heuristic "the gullibility heuristic." The gullibility heuristic is a (putative) rule that states that one should accept whatever one perceives as true.²²

Heuristics are posited as cognitive short-cuts. Roughly, the idea behind heuristics is that the tougher the computational task, the more apt one is to use a heuristic. If the problem one is dealing with is too computationally demanding (e.g., making a probability judgment), then one typically doesn't engage in the demanding processing and instead uses a rule of thumb (like trading in representative categories for probabilistic distributions; see Kahneman and Tversky 1981).

Figuring out what to believe is a very computationally demanding problem. It's difficult enough that it sometimes goes by a proper name: The Frame Problem. One version of The Frame Problem is the problem of figuring out which beliefs to acquire (or update) and which to ignore based on one's current evidence, stock of beliefs, and recent actions (Dennett 1998). Some have taken the problem to be so intractable that they see the scientific study of belief acquisition and updating (and central cognition in general) as a fruitless venture (e.g. Fodor 1983, 1987a, 2000).

The gullibility heuristic could be used to solve the frame problem. Perhaps what people do is initially believe everything most of the time, and then later toggle belief strengths in different ways.²³ In short, the problem of belief acquisition is exactly the type of problem that is ripe for a heuristic solution, so perhaps we should pursue that line of inquiry, and not look for an architectural solution. Furthermore, since many of the studies I've presented in support of the Spinozan theory depend on getting someone to believe some stimulus that is presented exogenously (e.g., the memory asymmetry studies), perhaps their results could be explained by merely positing the gullibility heuristic.²⁴

²¹ More grist for the Spinozan mill: nay-sayers are apt to have high scores on the 'Need for Cognition' scale (a scale that ranks how much cognitive effort one is apt to engage in), and yea-sayers are apt to have low scores, just as the Spinozan predicts (Cacioppo and Petty 1982). Those who acquiesce more often do so because they do not want to expend more mental energy, so they end up believing whatever they token and then never reconsider these beliefs. Those who nay-say do so because they want more mental exercise and thus are willing to expend more mental energy, making them more apt to reject their extant beliefs.

²² I thank both Jerry Fodor and Bill Lycan for (separately) raising this objection.

²³ Of course, this way of 'solving' the problem just pushes it one step back: now the problem will arise for updating beliefs as opposed to acquiring beliefs.

²⁴ The main modern proponent of Spinozan theories is Dan Gilbert, who proposed a forerunner to the view described here (e.g., Gilbert 1991). If Spinoza deserves to be the namesake of this view, then Gilbert should at least be considered the modern intellectual progenitor of it. However, strictly speaking, Gilbert is just committed to the second and third

Yet there are some strong reasons for doubting the gullibility heuristic hypothesis. For one thing, it can't explain the belief perseverance effects. The participants in the belief perseverance studies have all the time they'd like to form their own thoughts about their abilities. Furthermore, most heuristics can be 'turned off' or overcome in certain situations, especially when the participants are told that the heuristic they are using is inapplicable to the situation at hand (see, e.g., Nisbett and Ross 1980; Chapman and Johnson 2002). So, when participants are told that the feedback they are about to receive is false, they should override the putative gullibility heuristic and withhold from forming beliefs based on the dubious feedback (after all, they are explicitly being told that the situation they're about to encounter is one where the heuristic does not apply). However, as we've seen, this is not how people behave. Thus, it seems like the belief perseverance effects are incompatible with the gullibility heuristic hypothesis.²⁵

4.2 The Informativeness Objection

Some theorists have proposed that people do have the ability to contemplate propositions without believing them, but only when the propositions are informative when they are false and not when they are uninformative when they are false (Hasson et al. 2005). If this were the case then all of the memory asymmetry evidence (3.1) used against the Cartesian view would be undercut.

To support their hypothesis Hasson et al. set up an experiment where participants were given statements that were paired with faces. Upon seeing a statement, participants were told whether that particular statement was true or false. The experiment was designed so that some of the statements were informative when true but not when false (e.g. 'this person walks barefoot to work'), some were informative when false but not when true (e.g., 'this person owns a television'), some were informative when either true or false (e.g., 'this person is a liberal'), and some were uninformative when both true and false (e.g., 'this person drinks tea for breakfast'). During the learning phase participants were instructed to memorize the statement/face pairs for later testing. Additionally, for some face/statement pairs participants were put under cognitive load.²⁶ In the testing phase participants revisited the faces and were asked to determine whether the accompanying statements were true or false of the person whose face they viewed.

For statements that were uninformative-when-false, the Spinozan prediction held: the cognitive load (interruption) had no effect on the recollection of true statements, but it did increase participants' tendency to report false statements as true. On the contrary, interruption had no effect on statements that were informative when false. For these statements, interruption affected true and

²⁶ The load was induced by another tone task. The participants would hear a tone and they were instructed to detect whether the tone was high pitched or low pitched and then push a button corresponding to the pitch.

properties of the Spinozan theory as I've described it here. Though he takes no specific stance on the architectural vs. heuristic question (and thus no stance on the question of the nomological *impossibility* of contemplating without believing, property 1), it is most natural to read Gilbert as espousing this type of heuristic view.

²⁵ Even if we put aside the perseverance effects, there are still other insurmountable hurdles for the heuristic proposal. For example, people seem to believe everything they think, even when the ideas are *self-generated* and the participants aren't under load. In a series of studies (Epley 2004; Epley et al. 2004; Epley and Gilovich 2001, 2006) researchers tweaked the traditional anchoring and adjustment paradigm (a subject to which I will return in section 5.3) and showed that there are self-generated anchoring and adjustment effects. Since the gullibility heuristic says to believe what you *perceive*, it should only range over exogenously given stimuli, but the self-generated effects show that people believe endogenously created stimuli even when they know their creations are fatuous. Additionally, a heuristic explanation would predict that neurological damage shouldn't cause a dissociation between acceptance and rejection, yet it seems that people with Capgras syndrome show such a dissociation. For a fuller treatment of the heuristic hypothesis see chapter 2 of Mandelbaum 2010.

false statements equally. That is, for informative-when-false statements, participants remembered true and false statements equally well regardless of cognitive load.

This appears to be a decidedly anti-Spinozan datum, for it seems to show that people do have the ability to withhold assent from propositions when those propositions are informative when false. The experimenters write, "These results support the idea that the effect of resource depletion on the encoding of falsity ultimately depends on whether or not the proposition's false version is informative" (ibid, p. 568). If their hypothesis were correct, then at best the Spinozan hypothesis's scope would be severely restricted. However, there is good reason to resist their conclusion.

First, it is important to note how odd the consequences of the informativeness hypothesis are. If the hypothesis were true, then people wouldn't be able to entertain a proposition without believing that proposition when the proposition is uninformative-when-false. People could only entertain without believing when they are thinking about propositions that are informative-when-false. Prima facie, this situation is theoretically untenable. Suppose that you encounter a proposition, P. If not-P is informative, then you will be able to contemplate P without believing it. However, in order for you to determine whether not-P is informative, you must first parse and to some extent consider P (such considerations needn't be conscious). But what happens when you consider P? When you're considering P, do you believe P or not? In other words, what relation do you bear to P before you have figured out whether not-P is informative? When first considering P you either believe it or you don't, but the informativeness hypothesis can't seem to account for this fact. According to the informativeness hypothesis, if the proposition you are about to consider is uninformative-when-false you will believe it upon first hearing it, and if it's informative-when-false you won't, in which case you as perceiver of propositions have to be able to somehow see into the future to determine your disposition toward a proposition. But of course no one (pace Daryl Bem) wants to affirm a theory that entails parapsychological powers. There is a way out of positing such powers and keeping the informativeness hypothesis, but it too is not the most promising route.

The reader may be thinking that Hasson et al. can just assume that we hold no relation when first encountering a proposition; that is, that we initially withhold assent like the Cartesian view supposes. But consider the following: it seems overwhelmingly plausible that before a person can determine how informative a proposition is, the person must first entertain the proposition (though, again, a person needn't consciously consider it). Consequently, it appears that the informativeness hypothesis must entail that people can withhold assent regardless of the (subjective) informational content of a proposition. This would in turn imply that after one has withheld assent, one goes and marks propositions as true *only when they are uninformative-when-false*. This is quite an odd situation. The informativeness hypothesis dictates that people believe propositions after they've considered a proposition they've been told is false when that proposition is uninformative-when-false. How could such a situation come about? How would the mind possibly evolve such an odd processing system? If we have the ability to withhold assent then why wouldn't we use this ability in situations where we are told a statement is false; why would we only use it when we are told a statement is false and that statement is informative-when false?

Showing that the explanation underwriting Hasson et al. is, to say the least, unclear still doesn't explain why they got the data they did. It would be nice if there were an explanation for what exactly caused these data to be generated. I'll now attempt to give you one that will happily be quite consistent with the Spinozan worldview.

I suspect that Hasson et al.'s ascertained their results because their study was flawed. Consider being a subject who has just seen two sentences, both of which you were told were false, one that is uninformative-when-false ('this person drinks tea for breakfast') and one informative-

when-false ('this person owns a television'). Why would we be more apt to remember that the latter was false? Perhaps we would be because the latter is more shocking and vivid. When we encounter abnormal situations we are more apt to think longer and harder about the abnormal situation (in this example one might think: who doesn't own a television?).²⁷ Finding out that someone doesn't drink tea for breakfast doesn't really get one's mental juices flowing but finding out that someone doesn't own a television immediately raises some questions (e.g., is this person a Humanities professor? A Communist? Is she poor? Does she live in the woods?). Unsurprisingly, the more you think about something, the more you are apt to remember it (Petty and Cacioppo 1986). The subjects in this study were probably thinking about the statements that were informative-when-false for a much longer amount of time than they were the statements that were uninformative-when-false. Participants were probably considering the situation where one doesn't own a television for longer than they would consider the situation where one doesn't drink tea for breakfast (certainly the former would startle undergraduates, the study's participants, more than the latter). Accordingly, the subjects would perseverate on the thought DOES NOT OWN A TELEVISION, more than they would meditate on DOES NOT DRINK TEA FOR BREAKFAST, thus they would be more apt to remember the former than the latter. Seen in this light, Hasson et al.'s data tells us nothing about the processing of belief per se.

One last reason to think that my above explanation is correct: the informativeness criterion coincides with the ease of imagining a situation. When one considers someone who doesn't drink tea for breakfast what comes to mind? There is no concrete mental image that occurs. However, when one considers someone who doesn't own a television, then many mental images pop up (try this on yourself). In fact, one can see the difference in these statements as on par with the difference between the abstract and concrete innuendo effects. In studies of the perseverance of innuendos (e.g., Wegner 1984) we find that innuendos make a deeper impression when they are concrete rather than when they are abstract. People can more easily ignore innuendos that are abstract (e.g., 'Audrey is not sour') than they can for innuendos that are concrete (e.g., 'Audrey did not rob Toys R Us'). Presumably this is because 'not sour' can be immediately translated into 'sweet.' Moreover, we know that people will flip negative statements into the equivalent positive statement whenever possible (see Wason and Johnson-Laird 1972). One can easily paraphrase and flip the abstract statements, but how could one do the same for the concrete statements? What comes to mind when I tell you that Audrey didn't rob Toys R Us? Was she at home sleeping? Did she attempt to rob it but was foiled by the Pinkertons? In sum, the concrete statements stick because it's hard to envision a particular situation that holds when the statement is false. The difficulty of envisioning does not occur in the abstract statements because they have a quick negative counterpart.

Similarly, the informative statements in Hasson et al.'s studies can be easily envisioned when negated. When considering that this person doesn't own a television you may immediately think of a person living in a log cabin in the woods (or perhaps you envision someone reading, or a big old-timey radio, or a television that's been turned into a diorama). However, when I tell you that this person doesn't drink tea for breakfast what is the first thing that comes to mind? Do you envision a person sitting at a table with no drinks? Do you envision a coffee cup? The uninformative situations

²⁷ There is some evidence that deals with this line of thought (see Brigard et al. 2009; Mandelbaum and Ripley 2012; Uttich and Lomborozo 2010). For example, the main thesis of Mandelbaum and Ripley is that people have a belief that when a norm is broken, an agent must have broken the norm. The idea is that one gleans more (sometimes false) information about a person's mental states when they break norms than when they follow norms. To repurpose an example from Uttich and Lombrozo, if you see me on a tuxedo at a fancy wedding, then you don't learn nearly as much about my mental states as if you see me in a tuxedo at the beach.

are hard to visualize when false. Thus, it is no wonder that people have a harder time remembering the veracity of uninformative statements than the veracity of informative statements. People will think about the latter more often and will thus be able to answer more correctly. Seen in this light, Hasson et al.'s results tell us nothing about the relation between contemplation and belief per se, and do not cast doubt on the Spinozan hypothesis. In fact, in order to explain their results one needs the Spinozan hypothesis in order to explain why participants represent uninformative statements as true even when they are told they are false. As opposed to attacking the Spinozan hypothesis, the data collected in Hasson et al. helps to support Spinozan view.

4.3 Argument from Introspection

Now let's examine an objection I frequently encounter. The objection starts by someone attempting to consider some fantastically odd proposition, like *dogs are made out of paper*. The objector then proclaims that after some quick consideration she is sure that she does not believe that dogs are made out of paper. Since the Spinozan theory says otherwise, she concludes that the Spinozan theory is false. An incredulous stare is often added for good measure.

The intuition behind this type of argument is robust. In general, people think that they know what they believe and they know it straightforwardly through introspection. This intuition presupposes that beliefs are the types of things that are consciously accessible through introspection. However, I, following many self-respecting philosophers and psychologists, do not think that beliefs are in general accessible through introspection (e.g., Bem 1970; Gopnik and Meltzoff 1994; Lycan 1986, 2008; Dretske 1995, 2004; Williamson 2000; Carruthers 2009, 2010). For example, when discussing Daryl Bem's work on belief, Joel Cooper writes,

We do not always have insight into our own attitudes and beliefs, especially when they are not very strong or salient...When asked about our opinion toward most political issues or attitude objects we engage the very same process to infer our attitudes as we use to infer the attitudes of others. We look at our behavior, analyze the environmental stimuli, and make a logical inference about our attitudes. (Cooper 2007, p. 37)²⁸

If I ask you whether you like pinto beans, you may immediately know the answer (perhaps you have a pinto-based diet), but more likely you first recall your history with pintos. Perhaps you ordered them last week, so you infer you must like them or else you wouldn't have ordered them. We generally infer what we believe by examining our past behavior (even if such an examination is reflexive, unconscious, and instantaneous). Of course, in the paradigm instances of belief, the belief has been made salient to us so often that we needn't engage in any elaborate inferential process: if I ask you whether you love your spouse, you generally know what the response is (or at least should be); if I ask you whether you believe that 2+2=4, you can quickly respond because it is a question you have frequently answered.

The intuition that we have the ability to introspect our beliefs is a cognitive illusion caused by the paradigmatic cases of belief. When we are asked what we believe about a topic that is strongly

²⁸ Lycan takes an even stronger line, asserting that beliefs are never conscious. Lycan writes, "It is an interesting question whether we can ever introspect beliefs. On both phenomenological and theoretical grounds I doubt that too; what we introspect, in the way of cognitive items, are judgments, and we infer our knowledge of our beliefs from these" (1986, p. 64). Though I am sympathetic to Lycan's claim, the Spinozan hypothesis can rebuff the introspective objection with the weaker thesis in the text.

affectively valenced the answer arises instantaneously. Yet, the vast majority of the beliefs we hold are not strongly valenced like our belief that we love our spouse, or that genocide is abhorrent. Rather, the vast majority of our beliefs are more like our belief in the tastiness of pinto beans. It is the salience of paradigm cases that lead us to infer that all cases of belief are like our cases of strong belief. Once one spots how the salient cases differ from the majority of cases, the intuition pushing against the Spinozan view should be tempered.

Moreover, the last thirty years of psychology have shown how opaque our minds are to us. Philosophers have overplayed how much we can introspect because they often parochially focus on the contents of thought instead of on mental processes. Mental contents are sometimes available for report (of course, sometimes they aren't too); however, our thought processes are almost never available to report (see, e.g., Nisbett and Wilson 1977). If they were available, then the cottage industry of priming studies wouldn't be thriving. Not only are our mental processes unavailable for report, but we are even apt to misreport our emotional states (Dutton and Aron 1974). It is reasonable to think that propositional attitudes are more like emotional states and mental processes than mental contents (and therefore just as difficult to report).²⁹

One reason for thinking beliefs are unlike contents is because, although they have contents as proper parts, they are more than just contents: a belief is a content with a certain functional role, and to have a functional role is to *play a part* in our mental economy. Functional roles also should strike us as reasonably similar to mental processes (both are operations on contents, not contents themselves) and as such shouldn't be introspectable.

A slightly different way of arguing for the same point is to note that beliefs, like propositional attitudes in general, are *relations* to contents. Being able to introspect the content of beliefs does not imply the ability to introspect beliefs themselves.³⁰ To do that we'd need to be able to introspect a certain relation to a content, and there is little reason to believe that we have introspective access to this relation. After all, these relations are to be spelled out in terms of the nomological connections of the belief states, and to figure out these connections we generally need to engage in empirical psychology.

Thus, in order to be able to truly introspect beliefs we would need to be able to introspect the functional roles and relations of beliefs. But it is difficult to see how one could introspect these dispositional and relational facts when even our empirical psychology has trouble ascertaining such facts (for a similar argument put to a different use, see Goldman 2006).

However, these philosophical arguments only go so far. The idea that beliefs can be merely introspected is deeply held, so it would be nice if there were some direct empirical evidence against the view that we can reliably introspect our own beliefs. Happily there is some recent evidence that speaks directly to the question of our access to our beliefs. Gweon et al. (2011) designed an experiment where subjects formed beliefs about certain pictures (e.g., a partially occluded picture was actually a picture of a snake) and experimenters recorded these judgments. The subjects were then told what a different set of people believed about the pictures (e.g., that the partially occluded picture was a fish). Some of these beliefs turned out to be true, others false, and the subjects were aware of whether they had formed true or false beliefs. Fifty minutes later the subjects were asked to recall what they believed and two very interesting results were found: first, subjects were at remembering whether they had true or false beliefs than they were at remembering whether others

²⁹ Even our metaphors for the attitudes (e.g., the 'belief *box*') show that beliefs are unlike contents—after all, the 'belief box' is supposed to be the 'place' you put certain contents.

³⁰ Compare: being able to introspect a part of visual processing (e.g., the output) does not imply being able to introspect visual processing.

had true or false beliefs. Moreover, there is no evidence that the subjects felt as if they were misremembering—to them they were just telling the experimenters what they previously believed. Second, the same neural network that was activated when subjects were thinking about their own beliefs was activated when subjects were thinking about others beliefs. In particular the neural networks that were activated when thinking about one's own beliefs (or others) were the same that are activated in traditional 'theory of mind' experiments (the right and left temporo-parietal junction and the dorsomedial prefrontal cortex) Since it's clear that we don't have introspective access to other's beliefs, it's reasonable to suppose that this data is a strong indicator that we don't have introspective access to our own beliefs. But even if one wants to ignore the neural data, the behavioral data should speak for itself.

As a consequence of the previous data, I think it wise to remain unmoved by an objection that crucially relies on introspection. Rather, like Lycan, Bem, Dretske, and others, I suggest that we find out what we believe by simulating what others would do in our position, by watching our own behavior, by inferring from past instances, by inferring from what is reasonable to believe, or through some other investigative methods. Moreover, our lack of introspective access to our beliefs is central to solving multiple psychological puzzles: it's why there are so many implicit racists who make sincere avowals of their egalitarian beliefs; it's why people fall in love when traveling on other continents, mistaking fear for lust; it's why writing a counterattitudinal essay will sway what we report our beliefs to be. Many beliefs we think we don't harbor we do, only we can't figure that out merely through introspection—that's why we have clever psychologists.³¹

Returning to our original proposition, a ridiculous belief such as *dogs are made out of paper* isn't a belief that's going to eventuate in much behavior, certainly not in the millisecond after contemplating it and before reporting you don't harbor the belief. Since this belief has such a low chance of causing any behavior, you couldn't come to find out that you harbor this belief even if you were excellent at reading your beliefs off of your behavior.³² If you considered a more sensible though still outlandish proposition, such as *all dogs carry deadly viruses*, you would probably also claim to not believe it after consideration. But for all that, you would probably show subtle signs of harboring the belief. For instance, if after considering that proposition you were presented with dogs, you'd probably start lightly sweating, the galvanic skin response being an effect of having considered the proposition. And note, that even if galvanic skin response isn't the paradigm of intentional behavior, it is nonetheless behavior that needs to be explained and the proffered hypothesis can do so.

When you consider a ridiculous proposition, you generally attempt to falsify the correspondingly acquired belief immediately. Assuming you're not under cognitive load, you can normally do this quickly. What then becomes available for introspection is your *judgment* that dogs aren't made out of paper. From this you can rightly infer that you believe that dogs aren't made out of paper (after all, you can't judge that X without believing that X).³³ Thus, in many of these cases people will both believe that dogs are made out of paper and believe that dogs aren't made out of paper, but they'll only think they have the latter belief because they have access to the judgment that

³¹ None of this is meant to entail that we don't have a type of special first-person access to our beliefs; rather, it just implies that if we do have such special access it is not gained through a Lockean style of introspection.

³² One may object by saying that the belief could show up in behavior at a later time. For example, if you believed that dogs were made out of paper then why would you ever give your dog a bath? However, the Spinozan can respond that you probably also have a much stronger belief that dogs aren't made out of paper, and we'd expect stronger beliefs to win out (in most contexts) over weaker ones.

³³ This follows on the tame assumption that conscious thoughts also involve tokenings in (e.g.,) a language of thought. If so, then the inference from judgment to belief is secure even for the Spinozan.

accords with that belief.³⁴ However, the consideration process just serves to change the relative strengths of these beliefs. A well-informed deliberator will raise the strength of the negated belief, but will still have formed the affirmative belief.

Perhaps an example will make you feel more comfortable with the idea that we don't have introspective access to our beliefs. Let me introduce Fred, a man in his early thirties who has always been pretty skinny and has imbibed a substantial amount of beer in his day. After turning thirty, Fred became worried that he had a beer gut, although he didn't actually have one. He would walk shirtless to the shower, and looking down would see a slight bulge in his belly, from which he inferred that he had a beer gut. After several weeks of this routine, Fred made a self-deprecating joke about his beer gut to his friends, who acted astonished at the suggestion. Fred then asked his friends if he has a beer gut and his friends said that he didn't. Fred trusts his friends and believes that they are giving him a sincere response. His friends' adamant denials of the beer gut serve as the best evidence he has; he now happily reports that he is, in fact, not a skinny guy with a beer gut.

However, every time Fred looks down at his stomach he *sees* a beer gut. Because he trusts his friends' words, Fred tries to discount these perceptions. For example, if Fred is asked if he has a beer gut, he asserts that he doesn't have a beer gut. One might thus reasonably suspect that Fred doesn't believe that he has a beer gut. Yet if you want to predict the majority of Fred's behavior, your best bet is to believe that Fred believes that he has a beer gut. When Fred walks by a mirror, he's apt to turn sideways so see if he has a bulge; when Fred walks to the shower, he still looks down and gets a spike of anxiety; when Fred approaches the buffet table, he thinks twice about the fried chicken; when Fred sees a beer commercial, he winces; when Fred goes clothes shopping, he opts for baggy shirts instead of more form-fitting ones. Yet Fred sincerely reports that he believes he doesn't have a beer gut. So, what's going on with Fred?³⁵

One important datum in explaining this situation is realizing that Fred looks down and sees a beer gut much more frequently than he hears that he doesn't have a beer gut. Fred looks down and sees the beer gut every day, whereas Fred's friends' interventions happen quite infrequently. Although Fred discounts his beer-gut perceptions as optical illusions and he trusts his friends' reports that he's beer gut–less, he still acts as if he has a beer gut. The Spinozan theory proposes that Fred acts as if he believes he has a beer gut, he's continually tokening the thought that he has a beer gut, and these tokenings are sufficient for believing that he has a beer gut. Moreover, the relative strengths of beliefs (between say, believing one does have a beer gut versus believing one doesn't) are in part a function of how often each belief is tokened.³⁶ Since Fred tokens the belief that he does have a beer gut more often than the belief that he doesn't, he believes that he does more strongly

³⁴ Of course, I now commit myself to the existence of contradictory beliefs, but I do so happily, for the uncovering of contradictory beliefs is legion in psychological inquiry. For some recent examples of contradictory moral beliefs see Cushman and Greene 2012, Strickland et al. 2011, Ripley 2011).

³⁵ Although I am a proponent of dissonance theory and think that its explanatory powers are often overlooked, as I have argued elsewhere (Mandelbaum and Ripley 2012), dissonance will be of little use here because dissonance theory posits that people abhor inconsistencies. Dissonance theory would predict that Fred's behaviors would align with his assertions (or vice versa), in which case the person who keeps asserting he has no beer gut should start acting as if he didn't have a beer gut.

³⁶ This is why the therapeutic advice of self-affirmation theory (saying what you want to believe over and over again) actually works (Steele 1988). It is strange to think that just saying over and over again, T'm a good, smart, likeable person' actually makes a difference to one's beliefs about one's goodness, intelligence, and likeability. That it does is a stark datum that is wholly explicable on the view I'm offering (see Sherman and Cohen 2006 for a fairly comprehensive overview of this phenomenon).

(and hence you see it in his behavior more clearly).³⁷ Nevertheless, he judges that he does not believe that he has a beer gut, because when he's discussed the issue in the past he has come to the sensible conclusion that he doesn't have one. However, since he can't introspect his beliefs, he only reports the belief that seems most reasonable, which is his judgment that he does not have a beer gut.

Fred's case is by no means unique. The moral to be taken from such cases is that we need to make a distinction between belief reports and beliefs. What we can introspect are the former, and not the latter. Beliefs are unconscious propositional attitudes. In contrast, the Spinozan views belief *reports* as a species of judgments (a person-level phenomenon) that can be affected by all sorts of pragmatic factors. The beliefs that we report having are beliefs that on reflection we catalogue as normatively respectable. In essence, the beliefs that we report are the beliefs that we *endorse*, and we generally are wont to endorse only normatively justifiable propositions. Consequently, what we endorse is affected by a slew of heterogeneous factors, such as social pressure, anxiety, face-saving techniques, etc. We endorse propositions that seem reasonable to us, and when we are 'introspecting our beliefs' we are generally just reasoning about what seems rational to believe, not searching our actual stock of beliefs (Evans 1982). What we end up sincerely reporting as beliefs may have little in common with what we actually believe. What we endorse is a social matter, but what we believe is a brute architectural matter; we believe what we think, even if we think many things that we would never want to publicly endorse.

One might object: "Beliefs are the types of things that play a role in practical reason. How could beliefs play these roles if they are never conscious? Either your 'beliefs' don't play these roles and so aren't beliefs, or they do play this role and so are available for conscious introspection." However, this line of thought is misguided. In effect, the Spinozan theory accuses the folk view of behavior of making too few distinctions. The Spinozan sees something akin to practical reasoning occurring on two levels: one at the conscious level and one at the unconscious level. At the conscious level, judgments—not beliefs—play a critical role; at the unconscious level, beliefs—not judgments—take center stage. Thus, the Spinozan can allow that beliefs still play the same role that they always did, they are just not accessible in ways we might have pretheoretically thought they were.

4.4 But Why Are These States Beliefs?

The last objections I will consider is why we must conclude that the states I've been discussing are *beliefs* and not some other mental state. One might take accept all the arguments I've given so far, and yet still not want to identify the states under discussion with belief. Perhaps one would rather identify them as aliefs (Gendler 2008a, 2008b) or as some yet to be named mental state. The motivation behind such a move is understandable: many of the interesting properties of belief, from an epistemological standpoint are missing from these states (these states are accepted arationally, you can hold contradictory ones, etc). At this point, I could capitulate and accept that the view so far glossed is interesting enough without bestowing the term 'belief' onto these states. However, I do not think that calling these states beliefs is just to bestow some honorific upon them. Rather, I think that if we are to have an empirically adequate theory of beliefs, then any of the conditions on what a

³⁷ Additionally, every time Fred tokens the negated thought *I don't have a beer gut*, he tokens *I have a beer gut*, which raises the strength of the affirmative belief. This situation raises the question of how negative thoughts can ever become stronger than positive thoughts, which is a topic discussed in chapter 5 of Mandelbaum 2010. In short, the answer is partly based around the unsurprising observation that strength of belief is not just a function of tokenings.

belief must be will be met by the states under discussion. Some the issues here are deep and involved; for space requirements I'll only be able to scratch the surface of some important debates.³⁸

I take it that whatever beliefs are, they must have at least the following properties: they must be semantically evaluable (have conditions of satisfaction), be able to be acquired by perception, give rise to Moore's paradox, interact with desires and other motivational states to cause behavior (so beliefs must have actual causal powers, in which case, one might want them to be mental particulars, pace Dennett), and perhaps most importantly, be able to serve as the premises in inferences (henceforth called 'inferential promiscuity'; Stich 1978). Although it would be nice if we could find an empirically adequate account of a mental state that met all of these conditions,³⁹ these last two conditions are non-negotiable. If we are to identify a given operationalized state as a belief, that state better be able to interact with motivational states to cause behavior⁴⁰ and, in particular, it better be the type of state that can serve as the premise in inferences. Beliefs can interact with other beliefs in order to generate new beliefs/knowledge. This last condition is not just a philosophical doctrine; it is this condition that has been used to separate full-blown beliefs from belief-like intramodular propositional states (such as our visual system's information that there is only one overhead light source).

Showing that the states under discussion do have these properties would be the first step to showing that these states are belief-like. Equally importantly, it would be one of the only serious attempts I know of to identify beliefs with states that appear in empirical cognitive science.⁴¹ To start, notice that the states under discussion have clearly been caused by perception and are semantically evaluable; in all of the experimental designs the beliefs are acquired perceptually and generally given a truth signal, so their semantic evaluability shouldn't be in doubt. Oddly enough, even though this account divorces belief from assertion in many ways, the account can still give rise to Moore's paradox: one cannot assert "p, but I don't believe p" for whatever proposition one entertains one believes and one has to entertain any state that one asserts (tokening by linguistic

³⁸ In particular I won't be able to discuss issues about the metaphysics of belief, comparing analytic, teleological, and psychofunctionalism. For a short discussion of how these issues intersect with the current topic at hand see Mandelbaum (2010). Suffice to say, that I think it's hard to square any analytic functionalism (or Dennett-style Interpretationalism) with a desire to have an empirically compatible, to say nothing of rigorous, theory of belief.
³⁹ And it would be very nice if an account of such states could a) explain how we can (more or less) believe anything we can assert, b) analyze beliefs as two part relations (between a believer and something else, such as a proposition); c) entail that beliefs can have the same content as other propositional attitudes (you can believe that there is a table in front of you, or desire it to be true); d) explain the opacity of belief; and e) mesh well with an empirical theory of mental processes (Fodor 1981). The first four conditions can be met if we assume that beliefs are relations to syntactically structured mental representations; I take it that this paper's argument is, in part, an attempt to meet this last condition.
⁴⁰ Of course, how this interaction works (and is implemented) is an empirical question. I see no reason why one must be committed to a picture where beliefs can only cause behavior through decision-theoretic means as opposed to heuristic (or other) means.

⁴¹ Perhaps the most lauded recent attempt to do this is to be found in Schwitzgebel (2002). However, to my eyes Schwitzgebel's account runs afoul of two very central desiderata: it cannot make sense of any of the acquisition data canvassed here and it is a deeply anti-realist account. It's Rylean dispositionalism is cleverly defended but by making beliefs dispositional it is unclear how beliefs are supposed to do any real causal work. Since it's unclear how dispositions can be proximal causes, it's unclear how beliefs can be understood to actually cause behavior. Since I take it that it's nonnegotiable that beliefs must cause behavior, it's reasonable to suppose that Schwitzgebel's account has an important looming lacuna. (Of course Schwitzgebel's account deserves more exposition and discussion that it can possibly receive here).

production mechanisms is just another way of entertaining).⁴² But the really interesting work to be done is to show that these states cause behavior and are inferentially promiscuous. To do so, let's revisit the experimental paradigm introduced in 3.1.

Lest one think that the asymmetry between remember truths and falsehoods holds just over 'mere memory' (whatever exactly that is supposed to be), perhaps one more example would help to show how this acquired information is used just like beliefs. In Gilbert et al. (1993) participants were asked to watch a video screen with two crawling lines of text on it, one on top of the other. The top scroll contained text reports of two unrelated crime incidents. Participants were told that they would read both true and false details about the incidents, true statements appearing in black, false statements appearing in red. The bottom crawl did not contain any text, but instead had digits that slowly moved across the screen. Half the participants were told to ignore these digits (the unburdened participants) whereas half were told to peruse the digit crawl and to push a button anytime the number 5 appeared (the burdened participants).

At the conclusion of the video, participants were asked to recommend a prison sentence for the offenses, ranging from zero to twenty years and were also asked to assess the criminal's personality (in particular, how much they liked them, how dangerous they were, and how much counseling would help them). The false statements the participants read during the first phase of the experiment either exacerbated or mitigated the severity of the crime. The participants in the burdened condition were significantly more likely to be persuaded by the false information. The participants in the unburdened condition recommended a sentence of six years when the false information was extenuating and seven when it was exacerbating (not a significant difference), whereas their burdened counterparts recommended five years in jail in the extenuating condition and eleven years in jail in the exacerbating (a statistically significant difference). Significant differences were also found across the board when looking at the defendant's likeability, ability to benefit from counseling, and overall dangerousness. Thus, it appears that the burdened participants believed the lies they read more than the unburdened. Moreover, the falsehoods did not just get into the participants 'belief box' and then were parroted out as responses; rather, the falsehoods became integrated with the participants' beliefs and affected a robust range of their responses. The propositions that the participants encountered while under load rippled through their cognitive system. In the first part of the study the participants not only processed the lies fed to them, but they made (presumably unconscious) inferences from those states which then informed their judgments concerning the duration of the sentence and the character's likeability. The attitudes the participants formed infiltrated and interacted with (presumably some subset) of their web of belief in order to produce the behavior the experiment detected.

The experiment just discussed is one of many (overlooked) studies that can help execute some real philosophical work. They serve as evidence that the states I've been discussing do indeed serve as premises in inferences and they show that the states are honest-to-god causally active in producing behavior. The inferential promiscuity is a very strong requirement—not any run-of-the-mill mental states can achieve it.⁴³

⁴² That said, these states do allow a kind of inverse paradox to arise: one can assert 'not-p but I believe p.' Of course, if you sever the relation between belief and assertion, like I'm inclined to do, the Moorean paradoxes are mere window dressing.

⁴³ In particular, aliefs do not appear to be capable of inferential promiscuity. The content of an alief is not a truth apt proposition, in general it's just a single mental representation and, it is, to put it lightly, unclear how a single mental representation (such as DOG) can serve as a premise in an inference. There is much more to be said on the fecund topic of alief; for a discussion of what contents aliefs can take and aliefs' ability to pick out psychological kinds see Mandelbaum (2012).

In sum, it appears that the states we've covered do fit with our criteria for beliefs. Now that I have a backbone of an account of beliefs, it's time to put the states to work and see what type of explanatory power these states can have.

5. The Case for the Spinozan: Explanatory Power

As we've seen, the Spinozan theory can account for important data that the Cartesian theory cannot. Additionally, the Spinozan theory can unify and explain many disparate, hitherto mysterious phenomena. In this section, I will briefly examine a subset of previously recalcitrant phenomena that the theory can elucidate.⁴⁴ The unificatory and explanatory power of the Spinozan theory gives us strong inductive reasons for believing the theory. One should read this section as an abductive argument. The Spinozan theory of belief formation is the only theory on offer that can give a unified explanation of what has previously appeared to be a slew of disconnected and problematic phenomena.

5.1 The 'Mere Possibilities' Version of the Confirmation Bias

The 'confirmation bias' refers to people's tendency to search for confirmatory, but not disconfirmatory, evidence for the hypotheses they believe (Lord et al. 1979, Klayman and Ha 1987). The bias is explicable on a basic dissonance theory (see, e.g., Festinger 1957). A potted explanation goes something like this: if we believe in X and we find evidence that speaks against X, ascertaining such evidence will put us into a dissonant state. Since dissonant states feel bad (Zanna and Cooper 1974), they act as negative reinforcers, and through classical conditioning they reinforce us to not search for disconfirming evidence. This type of explanation articulates why the confirmation bias arises in the case of a previously held belief. However, the 'mere possibilities' version of the confirmation bias also arises in cases where people are merely considering a proposition. For example, if people are asked to consider if they are happy with their social life, they generally respond that they are, but when people are instead asked if they are unhappy with their social life, they also respond that they are! (Kunda et al. 1993). In these cases people search their memory for information that would confirm the question and then stop their search once they have reached such information. Dissonance theories have trouble explaining the mere possibilities formulation because they suppose that people aren't yet invested in thoughts that they merely entertain. Thus, the mere possibilities formulation of the confirmation bias is a standing mystery. The Spinozan theory can explain it by positing that propositions that are merely entertained are thereby automatically believed. Since one believes merely contemplated hypotheses, the dissonance explanation can get a foothold and start doing its explanatory work.

5.2 Anchoring and Adjustment

In the paradigmatic anchoring and adjustment paradigm (see, e.g., Kahneman and Tversky 1974) experimenters ask participants to give numerical values in answer to some arbitrary questions, like "How old was Gandhi when he died?," and "What is the freezing point of vodka?"⁴⁵ Before participants are allowed to answer the target question, the experimenter arbitrarily selects a number

⁴⁴ Some of the topics that the Spinozan view can help explain but which are not discussed here because of space limitations are the Fundamental Attribution Error (Jones 1979), source monitoring errors (Sherman and Bessenoff 1999); self-affirmation theory (Sherman and Cohen 2006); the 'fearing fictions' phenomenon (Walton 1978); the appearance of 'aliefs' (Gendler 2008); the efficacy of certain sorts of propaganda (Skinner 2002); the 'abstract/concrete' paradoxes (Sinnott-Armstrong 2008; Brigard et al. 2008); the ubiquity of implicit racism (Nosek et al. 2007); and the recovered memories phenomenon (Schacter et al. 1997). These topics are discussed in chapter 3 of Mandelbaum 2010. ⁴⁵ Gandhi died at seventy-eight; vodka freezes at approximately -16.51 degrees Fahrenheit (for 80-proof vodka).

(e.g., by spinning a wheel, or by using a participant's social security number, or by a randomly chosen card), which serves as an 'anchor.' Participants are then asked whether the answer to the target question is higher or lower than the arbitrarily picked number. After answering this question, participants are allowed to answer the original question. The randomly generated anchors make a significant impact on the subjects' answers.⁴⁶ For example, people will guess that Gandhi died at 50 if they first have to decide whether he died before or after he was 9, and they'll guess he died at 67 if they receive 140 as the anchor (Strack and Mussweiler 1997).

Explanations of the anchoring and adjustment effect are scant at best. For example, the traditional 'explanation' of the effect is that people anchor onto a value and then adjust up or down from that value (Kahneman and Tversky 1974). This explanation is just a restatement of the phenomenon. A more recent explanation of the effect is that it is produced by "increased accessibility of anchor-consistent information" (Epley and Gilovich 2001; Mussweiler and Strack 1999, 2000). Although this seems like an explanation, this explanation itself is just an instance of a broader trend, the bias toward searching for confirmatory evidence; the confirmation bias. Hence, the confirmation bias is supposed to explain the anchoring and adjustment effect. But as we've just seen, the confirmation bias itself presupposes the Spinozan theory.⁴⁷

As previously discussed, the confirmation bias is only supposed to be in play when one is searching for evidence to confirm an *already held* belief, so by accepting the confirmation bias explanation the non-Spinozan theorist just doubles her mysteries, for she also needs to explain why merely contemplated hypotheses are believed. But the Spinozan theory can eliminate these mysteries. Anchoring and adjustment effects arise because participants believe that the anchor they're given is actually the answer to the question they have been posed. Participants believe that the anchor sare the correct answer merely because they entertained that possibility, and entertainment causes belief.⁴⁸

5.3 Negation

A Spinozan theory that accepts property 4—which states that to negate a thought is, in part, to reject it—makes many predictions regarding negation. Two in particular are germane to this short discussion: the prediction that negations are difficult to process and the prediction that negations are held back in the initial processing of a sentence. The former prediction is well-known, so I will keep my discussion of it short (for the full suite of evidence and arguments see chapter 3 of Mandelbaum 2010).

5.3.1 The Difficulty of Negation

On the Spinozan theory rejection can occur only after acceptance. But it's not just the greater number of steps involved that makes rejection difficult; rather, it's that since starting the rejection

⁴⁶ The participants given an anchoring number generally choose a number halfway closer to the anchors than the numbers chosen by participants who don't encounter an anchor (Jacowitz and Kahneman 1995, p. 1163).

⁴⁷ One may be inclined to claim that anchor-consistent information becomes available through mere semantic priming. Maybe one's 'accumulator' (see Gallistel and Gelman 1992) is active when the number 140 comes up, and this primes other closer numbers. If this is so, then we'd expect that what participants do when they adjust is to continually slide along the number line until they reach a limit, one presumably dictated by the extent of the priming effect. However, Epley and Gilovich (2001) present data that speaks against the sliding hypothesis and propose instead that the adjustment phase is a series of jumps. Such jumps are inexplicable on the priming hypothesis.

⁴⁸ Here, as elsewhere, lies a tacit ceteris paribus clause. When participants are asked a question and then given the anchor, they must form a thought that turns the interrogative into a declarative. Fascinating evidence that participants do perform such a transformation can be found in Chapman and Johnson (2002).

process is optional, one has to expend effort every time one rejects a proposition. The effort needed to reject a proposition is thus greater than the effort needed to accept a proposition. Since negations are a subset of rejections, applying a negation should also be an effortful, and thus difficult, task. This is a theoretical coup for the Spinozan because practically anywhere one looks, one can find data showing that negations are hard to process.

Adding negations to a sentence exponentially increases the difficulty in understanding the sentence with each additional negation. One doesn't need much data to see the point: it's easier to understand 'Jane kicked the ball' than it is to understand 'Jane didn't kick the ball,' which is much easier still than 'Jane didn't not kick the ball,' and so on.

Negations also cause trouble when they are used as a search criterion. For example, people sort much faster and more accurately when they're asked to use a criterion that is positively formulated rather than negatively formulated (Wason 1972). Thus, people are much quicker at sorting when they're asked to sort out the spades and hearts from a pack of cards than when asked to sort out the non-clubs and non-diamonds (Fodor 1975). We would expect both faster performance and fewer errors when using a criterion that involves less mental energy, and the Spinozan theory states that the processing of affirmations uses less energy than the processing of their negative counterparts.

5.3.2 Psycholinguistic Processing of Negation

The second main prediction of the Spinozan theory regarding negation is more tendentious, though there is evidence that suggests that the prediction is accurate. The Spinozan predicts that negations, as a subspecies of rejections, can only be added to whole propositions, and this addition can be completed only after the proposition is formed. That is, the Spinozan theory predicts that in sentence comprehension people should process negative statements initially as affirmatives, processing the negation secondarily. This prediction is supported by Hasson and Glucksberg (2006). In their study, participants received affirmative and negative assertions and were then asked to perform a lexical decision task. For example, participants were asked to read sentences like "The kindergarten is/isn't a zoo' and 'Lawyers are/aren't sharks.' All of the statements participants read were metaphors, as to not allow for regular semantic priming effects to affect their data.

After reading the statements the participants were shown a string of letters on a screen and they were told to assess whether the letter string spelled an English word or not (i.e., they were given a lexical decision task). The experimenters varied the delay intervals between the metaphors and the lexical decision task and then looked at the participants' response times. Responses to affirmative-related targets were significantly faster than negative-related targets. Furthermore, the response latencies showed that *both* affirmative and negative sentences facilitated affirmative-related primes. However, the negative-related primes were *not* facilitated in the affirmative sentence conditions. For example, the negative sentence 'Surgeons aren't butchers' equally primed the affirmative-related prime 'clumsy,' as it did the negative-related 'precise,' whereas the affirmative sentence 'Surgeons are butchers' primed 'clumsy' but did not prime 'precise.' The negative-related prime 'precise' only arose in the negative context, whereas the positive-related prime arose in both contexts. This evidence shows the type of asymmetry the Spinozan hypothesis predicts and lends strong evidence to the view that negations are processed by first processing the corresponding affirmation.

The preceding evidence shows that people process affirmatives quicker than, and prior to, their negative counterparts. When processing a sentence, the negation is held back from the initial processing and appears online only after the initial processing happens; negations are not initially integrated in the construction of sentence meaning. Hasson and Glucksberg's study gives us a

glimpse of the actual time it takes negations to be processed. They conclude that negation doesn't take hold in processing until between five hundred and one thousand milliseconds after the negative sentence has been read, which is an enormous amount of time in linguistic processing. To illustrate, Hasson and Glucksberg non-metaphorically assert, "We found that terms related to the affirmative meaning of the metaphor were accessible immediately after reading the affirmative metaphors, indicating that the affirmative meaning was arrived at immediately" (p. 1027; for other work showing that affirmatives are processed immediately see Blasko and Connine 1993). The Spinozan view (but not the Cartesian) predicts this startling psycholinguistic data.

-***

In sum, the Spinozan hypothesis can help to explain quite disparate phenomena. The anchoring and adjustment effect and the confirmation bias are two of the most robust and mysterious psychological effects, and both can be elucidated using the Spinozan hypothesis, while the Cartesian model sheds no light on them. Additionally, the Spinozan theory can illuminate results about negation from many different branches of psychology: social psychology, cognitive psychology, and psycholinguistics, while the Cartesian theory stays silent on the matter of negation. Moreover, the Spinozan theory can explain one of the most well-supported findings in all of psychology: the fact that processing negation is hard. In contrast, the Cartesian model cannot even account for basic belief acquisition data. The Cartesian model continually misses the asymmetry between acceptance and rejection. For example, the Cartesian model makes all of the participants in the belief perseverance experiments look like exceptions even though their behavior is the rule. Surely, theories that can explain the relevant data are preferable to theories that claim that each datum is an exception. Thus, we should at least take the possibility of discarding the Cartesian hypothesis in favor of the Spinozan theory quite seriously. By doing so, we can explain and unify a number of mysterious phenomena in a theoretically respectable way. Even if there were no data against the Cartesian theory, this type of consilience should make one genuinely consider the merits of the Spinozan research program.

6. Rationality Imperiled

As mentioned at the outset, our conception of rationality is directly impacted by a Spinozan theory of belief acquisition. The Spinozan theory creates the following dilemma: either the ability to impartially doxastically deliberate is not a precondition on rationality, or people are necessarily irrational. Neither option is particularly appealing. Part of our concept of rationality is the ability to be a judicious cognizer; as philosophers we particularly pride ourselves on being able to justify our beliefs, and we have the expectation that these justifications aren't just post-hoc rationalizations. However, if the Spinozan theory is right then we don't have the ability to deliberate about a proposition before believing it.

That's just the start of the trouble for impartial deliberation. If the Spinozan theory is correct, not only would we be unable to withhold assent from propositions, but we would also be unable to impartially evaluate the beliefs that we do hold. Because of the confirmation bias (and the principles of dissonance in general) we will have a biased deliberation strategy, one where we tend to search for confirming information while ignoring disconfirming information. But this brings us to a depressing moral: at no point in our doxastic lives will we be able to consider propositions in a nonbiased way. Yet it seems that our normative standards demand that a rational cognizer at least be able to impartially consider propositions.⁴⁹ So, the first horn of the dilemma is quite unappealing.

The second horn is also unpalatable. For years research has been mounting that shows that people tend to be irrational in all sorts of domains; for example, we ignore base rates, we're Dutch bookable, we have trouble working out probabilities, etc. However, all these cognitive illusions are set against a background presumption of rationality. We consider ourselves irrational in these ventures as compared to our normal rational conception of ourselves. The rational conception of ourselves is central to many theories of intentional ascription (e.g., Davidson 2001; Dennett 1987). If we were to give up our conception of ourselves as rational creatures, then it is unclear what the paradigm of a rational creature would be.

I raise this dilemma not to solve it, but only to point out that our concept of rationality is imperiled in a new way. If the Spinozan theory is correct, we will have to reconsider either our standards of rationality or our conceptions of ourselves. Perhaps a cherished metaphor will help drive home the Spinozan challenge to rationality. The Spinozan theory gives us another way to understand the metaphor of Neurath's boat: we are always reconstructing our boat at sea because we never have any fixed point from which to adjudicate our beliefs. We're stuck with our beliefs, and even when we reject some, we are constantly drifting in the direction of the beliefs we hold, even if that direction is not particularly justifiable. We drift because our beliefs guide our searches towards confirming what we already believe, which is in part a function of whatever we happen to have entertained. We act as the epigram from D'Alembert tells us to: we just start moving forward with whatever propositions we happened to encounter and our psychology makes it so that, often enough, the faith in those propositions comes too. And of course, the propositions we happen to encounter are often a hodge-podge. Sometimes a thought pops in one's head, not because of some reasonable inferential process, but instead because of one's dinner choice. And presumably we wouldn't have wanted our epistemology held hostage to our gustation.

7. Conclusion

I've argued that there is a slew of evidence against the intuitive and ubiquitous Cartesian theory of belief fixation. As an alternative, I have offered a Spinozan theory of belief fixation. My goal has not been to argue that the theory is necessarily true, rather my aim has been the milder end of establishing that the theory is a respectable hypothesis about belief acquisition. And respectable hypotheses are what we need, for we have an overwhelming dearth of plausible theories of belief acquisition. The Spinozan theory is a fecund program, one with wide-ranging applications and one that can unify and explain quite disparate findings in psychology while diffusing some philosophical paradoxes. If you don't find the theory plausible upon first read, I recommend rereading the paper, preferably while in a bustling café.

References

Anderson, C., M. Lepper, and L. Ross. 1980. "Perseverance of Social Theories: The Role of Explanation in the Persistence of Discredited Information." *Journal of Personality and Social Psychology* 39 (6): 1037–49.

Bem, D. J. 1970. Beliefs, Attitudes, and Human Affairs. Belmont, Calif.: Brooks/Cole.

Blasko, D., and C. Connine. 1993. "Effects of Familiarity and Aptness on Metaphor Processing." *Journal of Experimental Psychology: Learning, Memory, and Cognition* 19 (2): 295–308.

Bratman, M. 1992. "Practical Reasoning and Acceptance in a Context." Mind, 101 (401): 1-16.

⁴⁹ If one wanted to deny that ought implies can (in epistemology), then our normative standards wouldn't necessarily demand that we have the capacity to live up to them. This isn't an unreasonable position, though it's not a way out of the dilemma either—it's tantamount to impaling oneself on the second horn.

Brigard, F., E. Mandelbaum, and D. Ripley. 2008. "Responsibility and the Brain Sciences." *Ethical Theory and Moral Practice* 12 (5): 511–24.

Cacioppo, J. T., and R. E. Petty. 1982. "The Need for Cognition." Journal of Personality and Social Psychology 42 (1): 116-31.

Carruthers, P. 2009. "How We Know Our Own Minds: The Relationship between Mindreading and Metacognition," Behavioral and Brain Sciences, 32: 121–82.

. 2010. "Introspection: Divided and Partly Eliminated," Philosophy and Phenomenological Research, 80: 76-111.

- Chapman, G and E. Johnson. 2002. "Incorporating the Irrelevant: Anchors in Judgments of Belief and Value." In T. Gilovich, D. Kahneman, and A. Tversky, eds. *Heuristics and Biases: The Psychology of Intuitive Judgment*. New York: Cambridge University Press.
- Cohen, L. 1992. An Essay on Belief and Acceptance. NY: Oxford University Press.
- Cooper, J. 2007. Cognitive Dissonance: 50 Years of a Classic Theory. London: Sage Publications Ltd.
- Cushman, F., and Greene, J. 2012. "The Philosopher in the Theater." In M. Mikulincer, P.R. Shaver, eds., *The Social Psychology of Morality: Exploring the Causes of Good and Evil.* Washington D.C: APA Press.
- Davidson, D. 2001. Inquiries into Truth and Interpretation. Oxford: Clarendon Press.
- Descartes, R. 1988. The Philosophical Writings of Descartes. Cambridge: Cambridge University Press.
- Dennett, D. 1987. The Intentional Stance. Cambridge, Mass.: MIT Press.
- . 1998. Brainchildren. Cambridge, Mass.: MIT Press.
- Dretske, F. 1995. Naturalizing the Mind. Cambridge: MIT Press.

———. 2004."Knowing What You Think vs. Knowing that You Think It." In Richard Schantz, ed., The Externalist Challenge. Berlin: Walter de Gruyter, 389–99.

- Dutton, D., and A. Aron. 1974. "Some Evidence for Heightened Sexual Attraction under Conditions of High Anxiety." Journal of Personality and Social Psychology 30 (4): 510–17.
- Epley, N. 2004. "A Tale of Tuned Decks? Anchoring as Accessibility and Anchoring as Adjustment." In D. J. Koehler and N. Harvey, eds., *The Blackwell Handbook of Judgment and Decision Making*. Oxford: Blackwell Publishers.
- Epley, N., and T. Gilovich. 2001. "Putting Adjustment Back in the Anchoring and Adjustment Heuristic: Differential Processing of Self-Generated and Experimenter-Provided Anchors." *Psychological Science* 12 (5): 391–96.
 2006. "The Anchoring-and-Adjustment Heuristic: Why the Adjustments Are Insufficient." *Psychological Science* 17
 - (4): 311–18.
- Epley, N., B. Keysar, L. Van Boven, and T. Gilovich. 2004. "Perspective Taking as Egocentric Anchoring and Adjustment." *Journal of Personality and Social Psychology* 87 (3): 327–39.
- Evans, G. (1982). Varieties of Reference. Oxford: Oxford University Press.
- Festinger, L. 1957. Theory of Cognitive Dissonance. Palo Alto: Stanford University Press.
- Festinger, L., and N. Maccoby. 1964. "On Resistance to Persuasive Communication." Journal of Abnormal and Social Psychology 68 (4): 359–66.
- Fodor, J. 1968. *Psychological Explanation: An Introduction to the Philosophy of Psychology*. New York: Random House. ______. 1975. *The Language of Thought*. New York: Thomas Y. Crowell.
- -------. 1981. "Propositional Attitudes." In *RePresentations: Philosophical Essays on the Foundations of Cognitive Science*. Cambridge, Mass.: MIT Press.
- . 1983. Modularity of Mind. Cambridge, Mass.: MIT Press.
- ———. 1987a. "Modules, Frames, Fridgeons, Sleeping Dogs, and the Music of the Spheres." In J. Garfield, ed., Modularity in Knowledge Representation and Natural-Language Understanding. Cambridge, Mass.: MIT Press.
- . 1987b. Psychosemantics: The Problem of Meaning in Philosophy of Mind. Cambridge, MA: MIT Press.
- . 1998 Concepts: Where Cognitive Science Went Wrong. New York: Oxford University Press.
- . 2000. The Mind Doesn't Work That Way. Cambridge, Mass.: MIT Press.
- Ford, K., and Z. Pylyshyn, eds. 1996. The Robot's Dilemma Revisited: The Frame Problem in Artificial Intelligence. Norwood, N.J.: Greenwood Publishing Group.
- Frankish, K. 2004. Mind and Supermind. Cambridge: Cambridge University Press.
- Gallistel, C., and R. Gelman. 1992. "Preverbal and Verbal Counting and Computation." Cognition 44 (1-2): 43-74.
- Gendler, T. 2000. "The Puzzle of Imaginative Resistance." Journal of Philosophy 97 (2): 55-81.
- _____. 2008. "Alief and Belief." Journal of Philosophy 105 (10): 634–63.
- . 2008. "Alief in Action (and Reaction)." Mind and Language 23 (5): 552-85.
- Gilbert, D. 1991. "How Mental Systems Believe." American Psychologist 46 (2): 107-19.
- ———. 1993. The Ascent of Man: Mental Representation and the Control of Belief." In D. Wegner and J. Pennebaker, eds., *The Handbook of Mental Control*. Englewood Cliffs, N.J.: Prentice-Hall.
- ———. 2002. "Inferential Correction." In T. Gilovich, D. Kahneman, and A. Tversky, eds., Heuristics and Biases: The Psychology of Intuitive Judgment. New York: Cambridge University Press.

- Gilbert, D., D. Krull, and M. Malone. 1990. "Unbelieving the Unbelievable: Some Problems in the Rejection of False Information." *Journal of Personality and Social Psychology* 59 (4): 601–13.
- Gilbert, D., R. Tafarodi, and P. Malone. 1993. "You Can't Not Believe Everything You Read." *Journal of Personality and Social Psychology* 65 (2): 221–33.

Goldman, A. 2006. Simulating Minds. Oxford: Oxford University Press.

- Gopnik, A., and A. Meltzoff. 1994. "Minds, Bodies and Persons: Young Children's Understanding of the Self and Others as Reflected in Imitation and "Theory of Mind' Research." In Sue Taylor Parker, Robert W. Mitchell, and Maria L. Boccia eds., *Self-Awareness in Animals and Humans*. New York: Cambridge University Press, 166– 186.
- Gweon, H., Young, L., and Saxe, R. 2011. "Theory of Mind For You, and For Me: Behavioral and Neural Similarities and Differences in Thinking about Beliefs of the Self and Other." *Proceedings of the 33rd Annual Meeting of the Cognitive Science Society.*
- Hasson, U., Simmons, J. and Todorov, A. 2005. "Believe It or Not: On the Possibility of Suspending Belief." *Psychological Science* 16 (7): 566-571.
- Hasson, U., and S. Glucksberg. 2006. "Does Negation Entail Affirmation? The Case of Negated Metaphors." *Journal of Pragmatics* 38: 1015–32.
- Huebner, B. 2009. "Troubles with Stereotypes for Our Spinozan Psychology." Philosophy of Social Science 39 (1): 63-92.
- Jacowitz, K., and D. Kahneman. 1995. "Measures of Anchoring in Estimation Tasks." *Personality and Social Psychology Bulletin* 21 (11): 1161–66.
- Jones, E. 1979. "The Rocky Road from Acts to Dispositions." American Psychologist 34 (2): 107-17.
- Kahneman, D., and A. Tversky. 1974. "Judgment under Uncertainty: Heuristics and Biases." *Science* 185 (4157): 1124–31. ______. 1981. "The Framing of Decisions and the Psychology of Choice." *Science* 211 (4481): 453–58.
- Klayman, J., and Y. Ha. 1987. "Confirmation, Disconfirmation, and Information in Hypothesis Testing." *Psychological Review* 94 (2): 211–28.
- Knowles, E., and C. Condon. 1999. "Why People Say 'Yes': A Dual Process Theory of Acquiescence." Journal of Personality and Social Psychology 77 (2): 379–86.
- Kruger, J. 1999. "Lake Wobegon Be Gone!: The 'Below Average Effect' and the Egocentric Nature of Comparative Ability Judgments." *Journal of Personality and Social Psychology* 77 (2): 1121–34.
- Kunda, Z., G. Fong, R. Sanitoso, and E. Reber. 1993. "Directional Questions Direct Self-Conceptions." Journal of Experimental Social Psychology 29 (1): 63–86.
- Long, A., and D., Sedley, eds. 1987. The Hellenistic Philosophers: Translations of the Principal Sources with Philosophical Commentary, Vol. 1. Cambridge: Cambridge University Press.
- Lord, C., Ross, L., and Lepper, M. (1979). "Biased Assimilation and Attitude Polarization: The Effect of Prior Theories on Subsequently Considered Evidence." *Journal of Personality and Social Psychology*, 37 (11): 2098-2109.
- Lycan, W. 1986. "Tacit Beliefs." In R. Bogdan, ed., *Belief.* Oxford: Oxford University Press. ———. 2008. "Phenomenal Intentionalities." *American Philosophical Quarterly* 45 (3): 233–52.
- Mandelbaum, E. 2010. The Architecture of Belief: An Essay on the Unbearable Automaticity of Believing. Doctoral Dissertation, UNC Chapel Hill.
- Mandelbaum, E. (2012). "Against Alief." Philosophical Studies. DOI: 10.1007/s11098-012-9930-7.
- Mandelbaum, E., and Ripley, D. 2012. "Explaining the Abstract/Concrete Paradoxes in Moral Psychology: The NBAR Hypothesis." *Review of Philosophy and Psychology* 3 (3): 351-68.
- Milgram, S. 1974. Obedience to Authority. New York: Harper Perennial.
- Mitchell, D. 2006. "Non-Conscious Priming after 17 Years: Invulnerable Implicit Memory?" *Psychological Science* 17 (11): 925–29.
- Mussweiler, T., and F. Strack. 1999. "Hypothesis-Consistent Testing and Semantic Priming in the Anchoring Paradigm: A Selective Accessibility Model." *Journal of Experimental Social Psychology* 35 (2): 136–64.
- ———. 2000. "Numeric Judgment under Uncertainty: The Role of Knowledge in Anchoring." Journal of Experimental Social Psychology 36 (5): 495–518.
- Nisbett, R., and L. Ross. 1980. Human Inference: Strategies and Shortcomings of Human Inference. Englewood Cliffs, N.J.: Prentice Hall.
- Nisbett, R., and T. Wilson. 1977. "Telling More than We Can Know: Verbal Reports on Mental Processes." *Psychological Review* 84 (3): 231–59.
- Nosek, B., Smyth, F., Hansen, J., Devos, T., Lindner, N., Ratliff, K., Smith, C., Olson, K., Chugh, D., Greenwald, A., and Banaji, M. 2007. "Pervasiveness and Correlates of Implicit Attitudes and Stereotypes." *European Review of Social Psychology*, 18, 36-88.

- Petty, R., and Cacioppo, J. 1986. Communication and Persuasion: Central and Peripheral Routes to Attitude Change. New York: Springer-Verlag.
- Pollock, J. 1986. Contemporary Theories of Knowledge. Totowa, NJ: Rowan and Littlefield.
- Price, H. 1990. "Why 'Not." Mind 99 (394): 221-38.
- Prinz, J. 2004. Furnishing the Mind: Concepts and Their Perceptual Basis. New York: Oxford University Press.
- Pryor, J. 2000. "The Skeptic and the Dogmatist." Nous 44 (4): 517-49.
- Pylyshyn, Z. 1989. "Computing in Cognitive Science." In M. I. Posner, ed., Foundations of Cognitive Science. Cambridge, Mass.: MIT Press.
- Quine, W. V. 1960. Word and Object. Cambridge, Mass.: MIT Press.
- Ripley, D. 2011. "Contradictions at the Borders." In Rick Nouwen, Robert van Rooij, Hans-Christian Schmitz, and Uli Sauerland, eds., *Vagueness in Communication*. Heidelberg: Springer LNCS.
- Ross, L., M. Lepper, and M. Hubbard. 1975. "Perseverance in Self-Perception and Social Perception: Biased Attributional Processes in the Debriefing Paradigm." *Journal of Personality and Social Psychology* 32 (5): 880–92.
- Samuels, R., S. Stich, and M. Bishop. 2002. "Ending the Rationality Wars: How to Make Disputes about Human Rationality Disappear." In R. Elio, ed., *Common Sense, Reasoning and Rationality*. New York: Oxford University Press.
- Schacter, D., K. Norman, and W. Koutstaal. 1997. "The Recovered Memories Debate: A Cognitive Neuroscience Perspective." In C. Martin, ed., Recovered Memories and False Memories: Debates in Psychology. New York: Oxford University Press.
- Schwitzgebel, E. 2002. "A Phenomenal Dispositional Account of Belief." Nous 36: 249-75.
- Sherman, J., and G. Bessenoff. 1999. "Stereotypes as Source-Monitoring Cues: On the Interaction between Episodic and Semantic Memory." *Psychological Science* 10 (2): 106–10.
- Sherman, D. K., and G. L. Cohen. 2006. "The Psychology of Self-Defense: Self-Affirmation Theory." In M. P. Zanna, ed., *Advances in Experimental Social Psychology* 38. San Diego: Academic Press.
- Sinnott-Armstrong, W. 2008. "Abstract + Concrete = Paradox." In J. Knobe and S. Nichols, eds. *Experimental Philosophy*. New York: Oxford University Press.
- Skinner, B. 2002. Beyond Freedom and Dignity. Indianapolis: Hackett Publishing Company.
- Spinoza, B. 1677/1991. Ethics. Indianapolis: Hackett.
- Stalnaker, R. 1984. Inquiry. Cambridge, MA: MIT Press.
- Strack, F., and T. Mussweiler. 1997. "Explaining the Enigmatic Anchoring Effect: Mechanisms of Selective Accessibility." *Journal of Personality and Social Psychology* 73 (3): 437–46.
- Steele, C. M. 1988. "The Psychology of Self-Affirmation: Sustaining the Integrity of the Self." In L. Berkowitz, ed., Advances in Experimental Social Psychology 21. San Diego: Academic Press.
- Stich, S. 1978. "Beliefs and Subdoxastic States." Philosophy of Science 45: 499-518.
- Strickland, B., Fisher, M., Peyroux, E., and Keil, F. 2011. "Syntactic Biases in Intentionality Judgments." XXXIII Proceedings of the Cognitive Science Society.
- Tuomela, R. 2000. Belief Vs. Acceptance. Philosophical Explorations. 3 (2): 122-37.
- Uttich, K., and Lombrozo, T. 2010. "Norms Inform Mental State Ascriptions: A Rational Explanation for The Side-Effect." *Cognition* 116: 87–100.
- Van Fraassen, B. 1980. The Scientific Image. Oxford: Oxford University Press.
- Velleman, D. 2000. "On The Aim of Belief" in The Possibility of Practical Reason. Oxford: Oxford University Press.
- Walton, K. 1978. "Fearing Fictions." Journal of Philosophy 75 (1): 5-27.
- Wason, P., and P. Johnson-Laird. 1972. Psychology of Reasoning: Structure and Content. Cambridge, Mass.: Harvard University Press.
- Wegner, D. 1984. "Innuendo and Damage to Reputation." Advances in Consumer Research 11: 694-96.
- Wegner, D., G. Coulton, and R. Wenzloff. 1985. "The Transparency of Denial: Briefing in the Debriefing Paradigm." Journal of Personality and Social Psychology 49 (2): 338–46.
- Williamson, T. 2000. Knowledge and Its Limits. Oxford: Oxford University Press.
- Zanna, M. P., and J. Cooper. 1974. "Dissonance and the Pill: An Attribution Approach to Studying the Arousal Properties of Dissonance." *Journal of Personality and Social Psychology* 29 (5): 703–9.