# Attention and evidence

**Abstract**

It is well known that agents may prefer to avoid cost-free evidence if they are uncertain whether they will conditionalize on that evidence. This tendency has historically been regarded as a type of irrationality. In this paper, I build on recent studies of Bayesian persuasion and rational inattention to show how agents with limited attentional capacities may prefer to receive partial rather than full information about the world, even if they are certain that they will optimally attend to incoming information and update rationally in response to the information attended to. I argue that some cases of this type are plausibly rational, then discuss implications for duties to gather evidence, bounded rationality, and internalist theories of evidence and epistemic justification. I also discuss a novel and moderate pattern of cost-free evidence aversion that arises in these cases: the agents studied may strictly prefer partial information to full information, but cannot strictly prefer no information to partial or full information.

## 1  Introduction

When and why should rational agents gather evidence? They may, of course, decline costly evidence if the evidence is not worth the cost. But may rational agents ever prefer not to receive cost-free evidence?

The classic answer is that they may not. I.J. Good (1966) and David Blackwell (1953) independently proved that under various conditions, rational agents cannot expect cost-free evidence to decrease the quality of their decisions. Patrick Maher (1990b) proved that agents also cannot expect cost-free evidence to reduce the accuracy of their beliefs. This suggests that if agents seek to maximize expected accuracy or utility, they cannot rationally prefer to avoid cost-free evidence.

Recent years have revealed various conditions under which rational agents might arguably prefer to avoid cost-free evidence. In particular, it has been suggested that rational agents may prefer to avoid cost-free evidence if they are risk averse (Buchak 2010; Campbell-Moore and Salow 2020), have priors that are imprecise (Bradley and Steele 2016) or not countably additive (Kadane et al. 2008), have imperfect control over their future selves (Maher 1990a), or face act-dependent world states (Adams and Rosenkrantz 1980), as well as on an externalist conception of evidence (Das 2023).

It is important to continue adding to this list of exceptions for two reasons. First, none of these exceptions apply to all agents and decision problems, so it is important to ask whether some further reasons for declining cost-free evidence might apply when known exceptions give out. Second, there is considerable debate over the rationality of many exceptions. Some might hold that rational agents should have countably additive priors, have full control over their future selves, or be risk-neutral, and externalism about evidence is controversial.

It has long been known that agents may prefer to avoid cost-free evidence if they are unsure whether they will conditionalize on that evidence. Evidence certainly can harm agents if they draw the wrong conclusions from it. Historically, this exception has been dismissed as a type of irrational aversion to evidence gathering. However, recent work by Sven Neth (forthcoming) revives the thought by suggesting that agents may be rationally modest, uncertain whether they will be able to follow the demanding dictates of Bayesian conditionalization, and on these grounds they may rationally prefer to avoid cost-free evidence.

In this paper, I present Good's Theorem and Neth's generalization (Section 2). I show how Neth's presentation depends on three limiting assumptions: that agents lack control over their future selves; that agents are uncertain they will update by conditionalization; and that agents are uncertain their future selves will behave rationally (Section 3). Then I present a related result based on the idea that rational agents may have limited capacities to attend to incoming information (Section 4). I show how rational agents who are certain

that they will optimally attend to incoming information and conditionalize on the results of attending to incoming information may nonetheless rationally prefer to receive partial rather than full information about the world.

I show how this result removes at least two, and arguably all of the limitations of Neth's result (Section 5), strengthening the robustness of the result that agents may rationally decline evidence if they suspect they will not conditionalize on the total evidence gathered. Section 6 concludes by discussing four important issues raised by my result: the relationship between attention and costly evidence-gathering (Section 6.1); a moderate feature of the rationalized pattern of cost-free evidence aversion (Section 6.2); the relationship between my result and the framework of bounded rationality (Section 6.3); and a surprising avenue of support for internalist theories of evidence (Section 6.4). Section 7 concludes.

## 2   Updating and Good's theorem

In this section, I present a restricted version of Good's Theorem (Section 2.1) and Neth's result (Section 2.2), with notation and assumptions selected for continuity with the discussion in Section 4.

### 2.1   Good's theorem

Suppose that an agent $S$ has credences $c$ on the algebra generated by a state space $\Omega$. She must choose among a set $\mathcal{A}$ of acts according to her utility function $u$. If forced to act now (Figure 1a), she faces the decision problem $\Gamma = (c, u, \Omega, \mathcal{A})$. The value of facing $\Gamma$ is the expected utility of the best act in $\Gamma$. That is, $V(\Gamma) = \max_a E_c[u(a)]$.

However, suppose $S$ is offered the opportunity to gather evidence before acting. Specifically, she may choose to receive an item $E$ from some evidence partition $\mathcal{E}$ of $\Omega$. Since she does not know which evidence she will receive, for each item of evidence $E$, $S$ assigns credence $c(E)$ to the prospect that gathering evidence from $\mathcal{E}$ will yield evidence $E$.

3

**(1a) No information**    $c \xrightarrow{\quad\quad\quad\quad E_c[u(a)] \quad\quad\quad\quad} a_c^*$

**(1b) Good**    $c \xrightarrow{\quad\quad c(*|E) \quad\quad} c(*|E) \xrightarrow{\quad\quad E_{c(*|E)}[u(a)] \quad\quad} a_{c(*|E)}^*$

**(1c) Neth**    $c \xrightarrow{\quad\quad ? \quad\quad} c_E \xrightarrow{\quad\quad E_{c_E}[u(a)] \quad\quad} a_{c_E}^*$
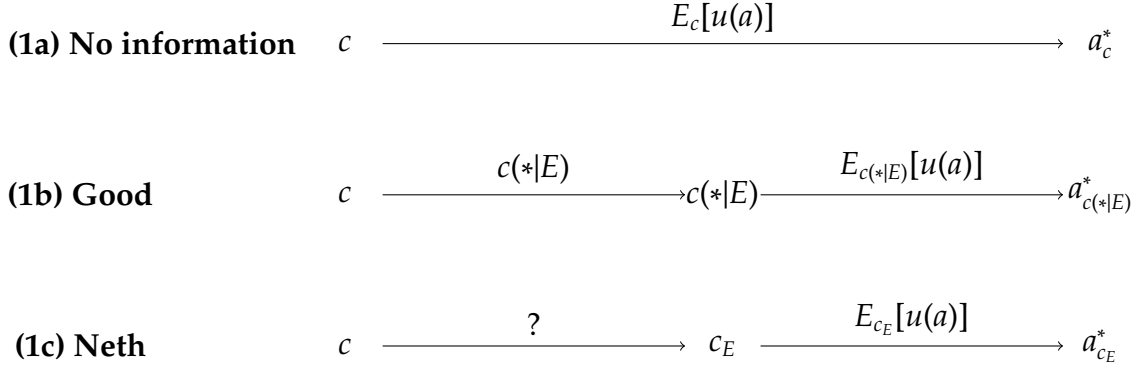
Figure 1: (1a) Agents who do not gather evidence maximize $E_c[u]$. (1b) In Good's setup, agents update by conditionalization to $c(*|E)$ and maximize $E_{c(*|E)}[u]$. (1c) In Neth's setup, agnets update by an unknown process to $c_E$ and maximize $E_{c_E}[u]$.

Suppose that $S$ gathers some evidence (Figure 1b) and learns some $E \in \mathcal{E}$. She updates her credences to $c(*|E)$ so that she now faces the problem $\Gamma_E = (c(*|E), u, \Omega, \mathcal{A})$. The value of facing $\Gamma_E$ is the expected utility of the best act in $\Gamma_E$, that is, $V(\Gamma_E) = \max_a E_{c(*|E)}[u(a)]$.

If $S$ were certain that she would receive some item $E$ of evidence, the expected value of gathering information would be the value difference between her new decision problem $\Gamma_E$ and her old decision problem $\Gamma$, that is, $V(\Gamma_E) - V(\Gamma)$. This quantity might well be negative. However, $S$ is uncertain which item of evidence she might gather, so the value of information depends on $S$'s prior beliefs $c(E)$ about the likelihood of gathering each item of evidence $E$, as

$$VOI_G = \sum_E c(E)V(\Gamma_E) - V(\Gamma). \tag{1}$$

Good (1966) proved that $VOI_G$ must be greater than or equal to zero. In this sense, an expected-utility maximizer cannot strictly prefer to avoid gathering cost-free evidence.

## 2.2  Neth's extension

Neth (forthcoming) observes that this result assumes the agent will update her beliefs by conditionalization. If we relax this result (Figure 1c), then we cannot assume that the agent's beliefs $c_E$ after receiving evidence $E$ correspond to the conditional probability function $c(*|E)$. This means that after receiving some item $E$ of evidence, the agent faces

4

decision problem $\Gamma'_E = (c_E, u, \Omega, \mathcal{A})$. The value of $\Gamma'_E$ is determined as before, as $V(\Gamma'_E) = \max_a E_{c_E}[u(a)]$. This induces a small, but important change in the value of information:

$$VOI_N = \sum_E c(E)V(\Gamma'_E) - V(\Gamma). \tag{2}$$

Although Good's value of information $VOI_G$ cannot be negative, Neth shows that this generalized value of information $VOI_N$ can be negative.

Borrowing an example from Neth, suppose I offer Ann a bad bet. More specifically, I flip a fair coin twice, showing neither flip to Ann. Then, I offer her a bet $B$ worth 1 util if the second flip landed heads, and $-2$ utils if the second flip landed tails. Since the coin is fair, Ann assigns expected utility $-1/2$ to $B$ and tells me to get lost, achieving total utility 0.

Now, suppose I offer Ann the chance to learn the outcome of the first coinflip. If Ann is certain that she will conditionalize, then she knows that whatever she learns, she will still assign expected utility $-1/2$ to $B$, so that gathering information can do no harm, in accordance with Good's Theorem. However, suppose Ann assigns small probability $\epsilon$ to the claim that she will commit a strong form of the gambler's fallacy, such that upon observing heads she will be highly confident that the next flip will be tails, and upon observing tails she will be highly confident that the next flip will be heads. If that happens, Ann will take $B$, a prospect to which she currently assigns expected utility $-1/2$. Otherwise, she will decline the bet and achieve utility 0. In this case, Ann assigns utility $-\epsilon/2$ to gathering evidence, so she strictly prefers to avoid cost-free disclosure of the value of the first coinflip.

More generally, call an agent *modest* if they are uncertain whether they will conditionalize.[1] Learning from Adams and Rosenkrantz (1980), assume that the agent's beliefs are not act-dependent, at least in the sense that conditional on the learned event $E$, your updated credences $c_E$ are independent of what action is best. Finally, make the richness assumption that for any $x$ in the unit interval $[0, 1]$, there is an outcome $o$ with $u(o) = x$.

---

[1] That is, for some $E \subseteq \Omega$, we have $c(c_E = c(*|E)) < 1$.

Neth proves that under the modesty and richness assumptions, for any modest agent $(c, u)$ we can always find an evidence partition $\mathcal{E}$ and a choice set $\mathcal{A}$ such that $VOI_N$ is negative on $\mathcal{E}$ and $\mathcal{A}$. In this sense, modest agents will always be vulnerable to violations of Good's Theorem.

Neth argues that modesty is rationally permissible. After all, we all know agents who have failed to conditionalize, and we have no good grounds to suspect that we will always do better. If this is right, then it may be rationally permissible for many of us to prefer to avoid receiving cost-free evidence in some decision problem.

## 3  Limitations

Neth's result provides an important step forward in our understanding of the interaction between updating policies and evidence gathering. Like all results, this result is limited in several ways which it would be productive for future research to remove.

First, Neth's result requires agents to now have imperfect control over their future behavior. In particular, agents must not be able to bind themselves to follow a given updating rule in the future. If Ann could now, before observing the first coinflip, choose an updating policy, she would choose to update by conditionalization on her current beliefs, on exactly the same grounds generally given for the rationality of updating by conditionalization (Briggs and Pettigrew 2020; Greaves and Wallace 2006). If this is right, then insofar as Ann suspects she might not update by conditionalization, she must not now have full control over her future updating policies.

It is becoming increasingly common to study agents with limited control over their future selves (Bermúdez 2018; Elster 2009; Gauthier 1997). However, previous results already suggest that agents who lack control over their future selves may rationally prefer to avoid cost-free evidence (Maher 1990a). It is natural to ask whether a similar result can be delivered, even for agents who have full control over their future selves, including an ability to bind themselves to follow update rules of their choice.

Second and relatedly, Neth's exposition discusses agents who are unsure whether they will conditionalize. They might, Neth suggests, commit the gambler's fallacy or weight items of evidence more heavily than their priors license. This is a natural restriction, since most agents have good reason to suspect that they will sometimes fail to conditionalize. However, many philosophers take conditionalization to be a rational requirement (Briggs and Pettigrew 2020; Greaves and Wallace 2006), and it is common to study agents who are certain that they will update by conditionalization. Could a similar result be proved for agents who are certain to update by conditionalization?

Third, Neth's result requires agents to be immodest, in the sense that they are uncertain whether they will behave rationally. While there are arguments for the permissibility of immodesty (Christensen 2007), immodesty is also known to cause problems (Titelbaum 2015) as a result of which some authors may reject the rational permissibility of immodesty. More to the point, the standard proof of Good's Theorem assumes a rationality model in which agents are certain to behave rationally. It is perhaps not overly surprising that Good's Theorem, like many results proven within rationality models, may fail once the assumption of rationality is relaxed. Indeed, it is already known that Good's Theorem may fail if agents suspect they may act irrationally after updating. Could a similar result be proved for agents who are certain that they will behave rationally?

In Section 4, I explore the case of an attentionally limited agent. I show how this agent may rationally prefer to receive less rather than more evidence. Then in Section 5, I argue that this case removes all of the limitations of Neth's result. The agent in question has full control over her future behavior, is certain that she will conditionalize and certain that she will respond rationally to current and future decision problems, but nonetheless is led to prefer less evidence to more. Then in Section 6, I discuss objections and draw out philosophical consequences of this discussion.

# 4 Attention

## 4.1 Rational inattention

In an information-rich economy, attention is becoming a scarce and precious resource (Castro and Pham 2020; Davenport and Beck 2001; Simon 1971). By one estimate, the average American is exposed to 34 gigabytes of information every day (Bohn and Short 2009). Agents can attend only to a small fraction of incoming information, so the choice of what to attend to exerts a substantial effect on the beliefs that agents go on to form and the actions they take on the basis of their beliefs.

Increasingly, philosophers (Archer 2021; Siegel 2017; Irving and Glasser 2020), economists (Maćkowiak et al. 2023; Sims 2003) and cognitive scientists (Lieder and Griffiths 2020; van den Berg and Ma 2018) have suggested that rational agents must make appropriate use of scarce attentional resources. On this view, wasting attention is just as irrational as wasting any other scarce resources, such as time, money, or computing power.

One reason for taking the allocation of attention to be a normative question is that we already make normative claims about attentional allocation. For example, we criticize agents for insufficient and harmful patterns of attentional allocation towards racial and ethnic minorities (Siegel 2017), and there are increasing calls to restrict the demands that technology can put on users' attention (Watzl 2021). Another motivation comes from the tradition of bounded rationality (Section 6.3), which has long stressed the need for agents to make appropriate use of scarce cognitive resources in order to improve the quality of the beliefs they form and the actions they take on the basis of those beliefs (Simon 1971; Thorstad forthcoming). Philosophers (Watzl 2021) and scientists (Lieder and Griffiths 2020; Sims 2003) have increasingly offered accounts of rational attention. This paper draws on a popular *rational inattention* framework due to Christopher Sims (2003), although it should be possible to develop similar arguments within competing frameworks.

There is good evidence that agents avoid cost-free disclosures of excessive information,

and that they often form more accurate beliefs and make better decisions as a result (Ben-Shahar and Schneider 2016; Hertwig and Engel 2021; Golman et al. 2017). For example, exclusion of evidence during legal fact-finding may increase the accuracy of fact-finding (Lester et al. 2012), and disclosure of complex financial information about loans may lead agents to pick inferior loans (Lacko and Pappalardo 2004).

One natural explanation for these phenomena is attentional: agents avoid excessive intake of cost-free evidence as a means of allocating scarce attention towards the most important and decision-relevant evidence.[2] As a result, they are able to form more accurate beliefs and make better decisions because their attention is focused on what matters most to improving the accuracy of their beliefs and the quality of their decisions. If that is right, then we have good reason to expect that theories of rational inattention should permit agents to turn down cost-free evidence as a means to improve the quality of their decisions and the accuracy of their beliefs.

It turns out that making good on this insight is not so easy. It is easy to see how cost-free evidence may fail to help agents without the attention to absorb it, but how can cost-free evidence harm agents? They could, it may seem, simply ignore evidence they do not wish to attend to. And in fact, some very idealized theories of rational inattention have the consequence that cost-free evidence cannot harm agents.[3] However, with a bit of care, we can construct a plausible story about why rational attentionally limited agents may turn down cost-free evidence in order to increase the quality of their beliefs and decisions. To keep the parallel with Good's Theorem, I focus on decision quality over belief accuracy.

In the rest of this section, I adapt a model of information transmission from the literature on Bayesian persuasion (Kamenica and Gentzkow 2011). I show that a benevolent communicator may withhold information from an agent in order to allow the receiving agent to produce better decision outcomes by appropriately allocating attention to the most decision-relevant information. Then I show how the same model can be used to

---

[2]To be clear, I do not suggest that this attention is the only contributory factor. Analyzing other factors should reveal further reasons for rational refusal of cost-free evidence.

[3]This includes the classic models of Sims (2003) and Caplin and Dean (2015).
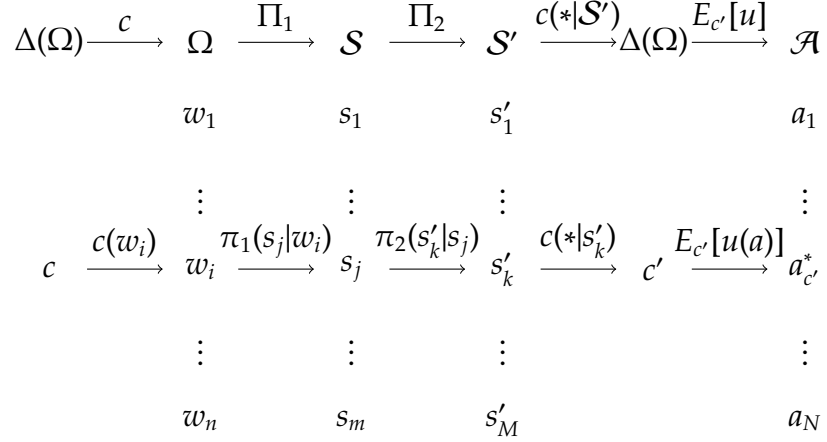
$$\Delta(\Omega) \xrightarrow{\phantom{xx}c\phantom{xx}} \Omega \xrightarrow{\phantom{x}\Pi_1\phantom{x}} \mathcal{S} \xrightarrow{\phantom{x}\Pi_2\phantom{x}} \mathcal{S}' \xrightarrow{c(*|\mathcal{S}')} \Delta(\Omega) \xrightarrow{E_{c'}[u]} \mathcal{A}$$

$$w_1 \qquad\qquad s_1 \qquad\qquad s_1' \qquad\qquad\qquad a_1$$

$$\vdots \qquad\quad \vdots \qquad\quad \vdots \qquad\qquad\qquad \vdots$$

$$c \xrightarrow{c(w_i)} w_i \xrightarrow{\pi_1(s_j|w_i)} s_j \xrightarrow{\pi_2(s_k'|s_j)} s_k' \xrightarrow{c(*|s_k')} c' \xrightarrow{E_{c'}[u(a)]} a_{c'}^*$$

$$\vdots \qquad\quad \vdots \qquad\quad \vdots \qquad\qquad\qquad \vdots$$

$$w_n \qquad\qquad s_m \qquad\qquad s_M' \qquad\qquad\qquad a_N$$

Figure 2: Communication by a benevolent sender. The sender transmits an informational signal $s_j$ by policy $\Pi_1$. Attentionally limited processing by $\Pi_2$ provides the receiver with signal $s_k'$. Receiver updates their priors $c$ to $c'$ by conditionalization on the signal $s_k'$ and maximizes expected utility given $c'$.

ground an argument for the rationality of a single agent opting to ignore cost-free evidence in order to attend to the most important evidence.

## 4.2 Benevolent informants

Suppose that agent $R$, a receiver, must choose an act $a$ from set $\mathcal{A}$. Perhaps $R$ is sick and must choose whether to undergo an experimental surgery, a standard surgery, or no surgery. She distinguishes among a finite set $\Omega$ of world-states and has credences $c$ on the algebra generated by $\Omega$. $R$ assigns utilities to each possible outcome $a(w)$ of performing an act $a \in \mathcal{A}$ in world $w \in \Omega$. If $R$ is forced to act now, she will maximize expected utility given her current, rather uninformed credences, choosing the act $a_c^*$ maximizing $E_c[u(a)]$.

However, another agent $S$, the sender, wants to help. Perhaps $S$ is $R$'s doctor. $S$ has the benevolent goal of providing $R$ with an informational signal about the world that will help $R$ to do as well as possible by $R$'s own lights. Perhaps $S$ must choose whether to send an extensive, detailed and complex disclosure of medical information, or a simpler, less detailed but more intuitive explanation of key facts.

Concretely, $S$ chooses a *signaling strategy* $\Pi_1 = (\mathcal{S}, \pi_1)$ allowing her to send different

signals in different states of the world (Figure 2). Choosing a signaling strategy $\Pi_1$ amounts to choosing two things. First, $S$ chooses a finite set $\mathcal{S}$ of possible signals to send. For example, $\mathcal{S}$ may consist of the detailed and intuitive disclosure documents. Next, $S$ chooses the probabilities $\pi_1(s_j|w_i)$ that each signal $s_j \in \mathcal{S}$ will be sent in each world-state $w_i \in \Omega$. There are no limits on the contents of $\mathcal{S}$ and $\pi_1$ beyond the requirements that $\pi_1$ be a probability function, so that $\sum_j \pi_1(s_j|w_i) = 1$ for all $i$.

At one extreme, $S$ may fully disclose the world-state by choosing $\mathcal{S}$ with $|\mathcal{S}| = |\Omega|$ and $\pi(s_j|w_i) = 1$ if $i = j$ and 0 otherwise. An agent receiving this signal would, in principle, be able to fully determine the state of the world: signal $s_j$ indicates state $w_j$. Alternatively, $S$ may partially disclose the world-state, for example by letting $\pi(s_j|w_i)$ lie strictly between 0 and 1 for some $i, j$. An agent receiving $s_j$ would then be unsure whether the world-state was $w_i$ or some other state.

If $R$ had unlimited capacity to attend to medical information, full disclosure would always be optimal by reasoning similar to Good's Theorem. However, we saw above that rational agents often must make do with limited attentional capacities. Let's be very generous and assume that $R$ learns the precise nature of $S$'s signaling strategy, so that $\Pi_1 = (\mathcal{S}, \pi_1)$ is known to $R$. This assumption is common in models of rational inattention, and might be justified by the assumption that $R$ has learned to adapt her attention well to similar situations in the past. However, models of behavioral inattention waive this optimality assumption, and matters only get worse for cost-free evidence gathering from there (Gabaix 2014).

Receiver $R$ must pick an *attention policy* $\Pi_2 = (\mathcal{S}', \pi_2)$ after the signaling strategy $\Pi_1$ is revealed to her. $\Pi_2$ introduces a garbling of the original signal, reflecting the fact that $R$ may pay incomplete attention to the incoming signal. The set $\mathcal{S}'$ is the set of possible signals that $R$ can receive after attentionally limited processing of the incoming signal. The conditional probabilities $\pi_2(s'_k|s_j)$ determine the likelihood that each signal $s'_k \in \mathcal{S}'$ will result from attention to each possible incoming signal $s_k \in \mathcal{S}$. As every doctor knows, a complete and thorough explanation $s$ of how things stand is not usually received as $s$ by

their patient (Falagas et al. 2009).

At one extreme, $R$ may pay full attention to $\Pi_1$ so that $|\mathcal{S}'| = |\mathcal{S}|$ and $\pi_2(s'_k|s_j) = 1$ if $k = j$ and 0 otherwise. Alternatively, incomplete attention may introduce normally distributed noise $\epsilon$ into the signal, so that $|\mathcal{S}'| = |\mathcal{S}|$ but $\pi_2 = \pi_1 + \epsilon$. Or, as we shall shortly see, the agent may prefer to preferentially direct her attention towards aspects of the signal that are most relevant to the decision at hand.

After receiving some signal $s'_k$ from $\Pi_2$, $R$ updates her credences by conditionalization to $c' = c(*|s'_k)$, using the known distributions $\pi_1, \pi_2$, the known signal $s'_k$ and her priors $c$.[4] Now instead of an uninformed action, she takes the act $a^*_{c'}$ maximizing expected utility $E_{c'}[u(a)]$ on her new credences. If the doctor has done her job, $R$ now has a better understanding of the virtues of each treatment and expects a better outcome from her choice.

Because attention is a scarce resource for $R$, she incurs a cost for her attentional policy $\Pi_2$. If we understand the function of attention as the extraction of information from the incoming signal $s$, then the cost of attentional policies should increase in the amount of information that they are expected to extract. If we had a measure $I(c)$ of the information contained by a credence function $c$, then we would treat attentional costs as increasing in $E[I(c')] - I(c)$, the difference between the known information content of the prior $c$ and the unknown information content of the posterior $c'$.

Most authors add a scalar $\kappa$ to capture the varying cost of attention across problems.[5] The cost of attentional policy $\Pi_2$ should then be

$$C(\Pi_2) = \kappa(E[I(c')] - I(c)) \tag{3}$$

for some $\kappa > 0$ measuring the relative cost of attention for agent $R$ in her current situation.

---

[4]That is, $c(w_i|s'_k) = \frac{c(s'_k|w_i)c(w_i)}{c(s'_k)}$ with $c(s'_k|w_i) = \sum_j \pi_1(s_j|w_i)\pi_2(s'_k|s_j)$ and $c(s'_k) = \sum_i c(s'_k|w_i)c(w_i)$.

[5]This can be justified on two grounds. First, attentional costs are often opportunity costs: attention paid to a medical decision takes attention away from other problems, leading to downstream utility losses. Opportunity costs vary according to the agent's current suite of decisionmaking problems. Second, agents differ in their capacities to attend to information, so attentional costs should be indexed to agents.

How should the information content $I(c)$ of opinions $c$ be understood? Many choices are possible here, but most models of rational inattention follow information theorists in measuring information content using Shannon entropy (Shannon 1948).[6]

Receiver $R$ must choose an attentional policy $\Pi_2$ based on her priors $c$ and the known signaling strategy $\Pi_1$ of the sender $S$. In our example, she must choose how much and what kinds of attention she will pay to the extensive and intuitive medical disclosures. Because $R$ is rewarded for the value of her act $a \in \mathcal{A}$, less the cost $C(\Pi_2)$ of her attentional strategy, she chooses $\Pi_2^*$ to maximize net reward, as

$$\Pi_2^* = \operatorname*{argmax}_{\Pi_2}[E_{c'}[u(a_{c'}^*)] - C(\Pi_2)] \tag{4}$$

where the optimal non-attentive act $a_{c'}^*$ is chosen to maximize $E_{c(*|s_k')}$ with uncertain signal $s_k'$ depending on the world state and the chosen policy $\Pi_2$.

## 4.3   Optimal disclosure

Let's think again about our benevolent sender $S$. Suppose $S$ knows a good deal about her patient $R$: $S$ knows $R$'s utilities $u$, priors $c$, and attentional cost $\kappa$. $S$ also knows that her patient $R$ will choose according to (4). Should $S$ fully disclose everything to $R$? Or should $S$ send a more modest signal that helps $R$ to focus on what is important?

In the case $|\Omega| = 1$ of a single world state, there is no uncertainty and the problem is trivial. The smallest nontrivial case occurs when $|\Omega| = 2$, so that the patient is uncertain among two possible world-states. There is this much to be said in favor of full disclosure: in the highly unlikely event that $R$ is uncertain over a binary state space, full disclosure is optimal. That is, we can show that the signaling policy $\Pi_1^F = (\mathcal{S}, \pi_1)$ with $|\mathcal{S}| = |\Omega|$ and $\pi_1(s_j|w_i) = 1$ if $i = j$ and 0 otherwise allows $R$ to select a higher-value attentional policy by the lights of (4) than she would be able to select if $S$ chose another signaling policy (Lipnowski et al. 2020).

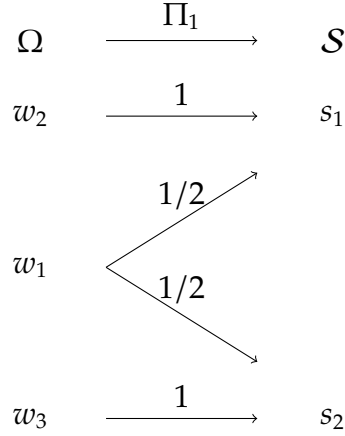---

[6]That is, $I(c) = E_c[-log(c)]$.

Figure 3: Partial disclosure policy.

However, this is the most that can be said for full disclosure. In the common case that there are at least three world-states, we can always find an agent and decision problem for which full disclosure is suboptimal. That is:

**Theorem 1 (Lipnowski et al. 2020):** The following are equivalent.

i Full disclosure ($\Pi_1^F$) is optimal for every $u, c, \kappa, \mathcal{A}$.

ii $|\Omega| < 3$.

To illustrate the intuition behind Theorem 1, it may help to borrow an example from (Lipnowski et al. 2020). Suppose there are three acts $a_1, a_2, a_3$: for example, no surgery, standard surgery, or experimental surgery. There are three world-states $w_1, w_2, w_3$ with each $a_i$ performing better than its competitors in the corresponding state $w_i$. However, suppose that it is most decision-relevant for $R$ to learn about the values of her surgery options $a_2, a_3$, for example because she regards $w_1$ as unlikely or because the experimental surgery $a_3$ represents a costly mistake in the world $w_2$. Then $R$ would like her doctor to inform her about the most decision-relevant states $w_2, w_3$, allowing her to fully attend to this information without filtering out low-value information about $w_1$.

The doctor does this by fully disclosing the states $w_2$ and $w_3$, but giving no information about $w_1$. That is, she chooses $\mathcal{S}$ with $|\mathcal{S}| = |\Omega|$, and $\pi(s_i|w_i) = 1$ for $i = 2, 3$, randomly

14

mixing signals in $w_1$ as $\pi(s_2|w_1) = \pi(s_3|w_1) = 1/2$ to convey no information about $w_1$ (Figure 3). We can show that for suitable specifications of $u, c$ and reasonable attentional costs $\kappa$ that this partial disclosure policy outperforms full disclosure (Lipnowski et al. 2020).

## 4.4 Cutting out the middleman

So far, we have seen that a benevolent sender may choose to withhold information from a receiver. The sender avoids full disclosure in order to direct the receiver's attention to the most important and decision-relevant information, allowing the receiver to realize more value by her own lights. We saw that this result is far from isolated: beyond the near-trivial case of a binary or unary state space, we can always find a decision problem and an agent whose credences, utilities, and attention costs make it the case that she benefits from partial rather than full disclosure.

At this point, readers may ask what all of this has to do with evidence-gathering. It is all well and good to say that benevolent informants might provide agents with incomplete information. But how does this bear on the question of whether agents themselves may rationally seek out incomplete information in a context where further information is available and cost-free?

To see the relevance, note what we did not need to assume in our discussion above. We did not need to assume that the receiver $R$ pays for the cost of the informational signal that $S$ prepares for her. As far as $R$ is concerned, the evidence gathered by $S$ is cost-free. Nor did we need to assume that $R$ or $S$ are ignorant of the setup. In fact, we assumed that $R$ knows with certainty the signaling process $\Pi_1$ generating informational signals, and we even assumed that those signals were chosen by an agent who knew exactly how $R$ would react to them.

All of this suggests that we should cut out the middleman and consider a reformulated model in which $R$ seeks out informational signals that are available to her in her environment (Figure 4). Suppose that $R$ wants to learn about her surgery options. $R$ can gather evidence by attending to cost-free explanations available in her environment.
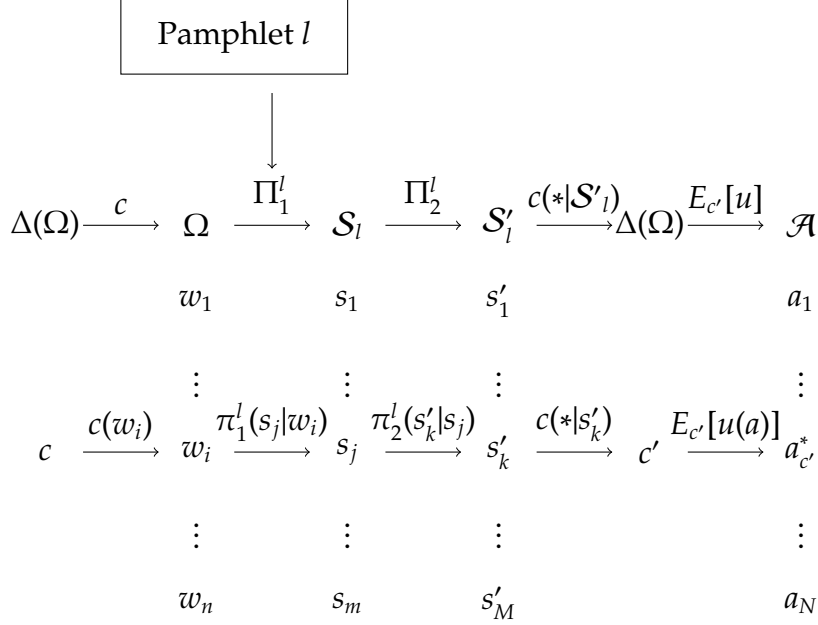
$$\boxed{\text{Pamphlet } l}$$

$$\Delta(\Omega) \xrightarrow{\;c\;} \Omega \xrightarrow{\;\Pi_1^l\;} \mathcal{S}_l \xrightarrow{\;\Pi_2^l\;} \mathcal{S}'_l \xrightarrow{c(*|\mathcal{S}'_l)} \Delta(\Omega) \xrightarrow{E_{c'}[u]} \mathcal{A}$$

$$w_1 \qquad s_1 \qquad s'_1 \qquad\qquad a_1$$

$$\vdots \qquad\quad \vdots \qquad\quad \vdots \qquad\qquad \vdots$$

$$c \xrightarrow{c(w_i)} w_i \xrightarrow{\pi_1^l(s_j|w_i)} s_j \xrightarrow{\pi_2^l(s'_k|s_j)} s'_k \xrightarrow{c(*|s'_k)} c' \xrightarrow{E_{c'}[u(a)]} a^*_{c'}$$

$$\vdots \qquad\quad \vdots \qquad\quad \vdots \qquad\qquad \vdots$$

$$w_n \qquad s_m \qquad s'_M \qquad\qquad a_N$$

Figure 4: Choice of medical pamphlets. Pamphlet $l$ transmits an informational signal $s_j^l$ by policy $\Pi_1^l$. Attentionally limited processing by $\Pi_2^l$ provides the receiver with signal $s'_k$. Receiver updates their priors $c$ to $c'_l$ by conditionalization on the signal $s'_k$ and chooses $a^*_{c'_l}$ maximizing expected utility given $c'_l$. The pamphlet maximizing the expected utility of the resulting action is chosen.

For example, perhaps the doctor's office offers a range of pamphlets, from a simple infographic to a complete compendium of all relevant scientific literature. Even under the strict assumption that $R$ knows the processes $\Pi_1^1, \ldots, \Pi_1^N$ used to generate the pamphlets, and even if some pamphlet is generated in a way that fully discloses the relevant world state, she may rationally choose some pamphlet generated by an incomplete disclosure process $\Pi_1^l$ over the pamphlet generated by a complete disclosure process. She may, that is, turn down a cost-free opportunity to gain perfect information for an opportunity to gain imperfect information about her environment.

Theorem 1 guarantees that outside of unary and binary state-spaces, we can always find an agent and a decision problem for which it is optimal to turn down cost-free full disclosure in favor of partial disclosure. In fact, such agents will be willing to pay a *cost* to avoid full disclosure. And the examples cited in Section 4.1 suggest that this phenomenon may be relatively common: attentionally limited agents seek to avoid costless disclosures

of excessive information, and do better for it. This need not occur because they are irrational. Rather, they may be limiting their information intake as a strategy for rationing scarce attention in the direction of the most decision-relevant information.

# 5   Limitations removed

In Section 3, we saw that Neth's original result establishes the rationality of ignoring cost-free evidence given three limitations. First, agents have incomplete control over their future acts: in particular, they cannot now bind themselves to later follow a given update policy. Second, agents are uncertain whether they will update by conditionalization. Third, agents are uncertain whether they will behave rationally in the future. In this section, I argue that the attentionally limited model from Section 4 removes all three limitations of the Neth model.

First, the agent in this model has full control over her future update policy. Her update rule is fully determined by her attentional policy $\Pi_2$, and she chooses $\Pi_2$ prior to receiving the incoming signal $\mathcal{S}$. If we like, we might generalize the problem and see her as choosing $\Pi_2$ prior even to knowing the structure of incoming information. For example, she might choose a map from possible signaling strategies $\Pi_1$ to attentional policies $\Pi_2$. The agent is certain that she will follow the chosen policy $\Pi_2$ and hence has no need to bind herself, but if desired we can add an arbitrarily strong binding mechanism without changing rational behavior.

Second, the agent in this model is certain that she will update by conditionalization. She knows that upon receiving attentionally processed signal $s'_k$ she will update from prior beliefs $c$ to posterior beliefs $c(*|s'_k)$ through the ordinary process of Bayesian conditionalization. What is uncertain is not *whether* she will conditionalize, but *what* she will conditionalize on. Because incoming information is filtered through attentionally limited processes, the agent conditionalizes on the attentionally filtered signal $\mathcal{S}'$ which may, but need not coincide with the signal $\mathcal{S}$. One way to understand this result would be to say

17

that although the agent is certain she will conditionalize, she is not certain that she will conditionalize on her total evidence: the signal $S'$ may not convey her total evidence $S$. However, we will see (Section 6.4) that this gloss is controversial: there are some understandings on which the agent not only conditionalizes, but also conditionalizes on her total evidence.

Finally, the agent in this model need not suspect herself of irrationality. For one thing, we saw in Section 4 that the agent's behavior is rationalized by a leading model of rational attention, and is likely to be rationalized by many competing models. For another, the agent is certain that she does and will always do all of the following: have probabilistic credences, update those credences by Bayesian conditionalization, and take the available acts which maximize expected utility.

There are, of course, important philosophical issues raised by the agent's bounded attentional capacities, and a full discussion of the rational importance of these capacities will have to wait for a fuller discussion of bounded rationality in Section 6.3. However, at the very least, we should pause now to be impressed by the difference between this agent and agents such as Ann. Unlike Ann, the agent in this problem does not suspect herself of paradigmatic forms of irrationality such as overconfidence and the gambler's fallacy. She has detailed and well-informed beliefs about her own capacities and about the structure of the information she receives, and takes both mental acts (updating) and non-mental acts (options from $\mathcal{A}$) to maximize expected utility given the structure of her decision problem. Many traditional arguments for immodesty, directed at agents such as Ann, will do nothing to impugn the rationality of this agent, so if there is something irrational about the agent described above, it will need to be spelled out with some care.

Summing up, the agent in Section 4 seems to lack all three limitations of the Neth result. She has full control over her future behavior, including updating behavior. She is certain that she will conditionalize. And she need not, on many views, be or suspect herself to be irrational. If that is right, then it tends to strengthen the case for the rationality of declining cost-free evidence by removing limitations on which that case might have

18

been thought to depend.

# 6  Discussion

This paper developed a model in which attentionally limited agents allocate attention optimally to incoming information on the basis of full knowledge of the statistical distribution from which that information is drawn. We saw that agents in this model may strictly prefer partial rather than full disclosure of act-relevant world-states. And we saw that this result does not depend on three limitations of the Neth (forthcoming) result. In this section, I discuss four important implications raised by the results of this paper.

## 6.1  Cost-free evidence

Some readers might ask whether the argument in this paper makes the familiar point that while agents might not prefer to avoid cost-free evidence, they can certainly prefer to avoid costly evidence. After all, they may not want to pay the cost. It might be asked whether attention assumes the familiar role of evidential costs in Section 4. Certainly, this model does not assume a that evidence *gathering* is costly: the cost of evidence gathering is born by the sender $S$. However, the model might assume that evidence *processing* is costly, insofar as it consumes limited attention.

Here it is important to stress at least three points. First, we could change the model to treat attention as not costly but limited. We would assume that agents have the capacity to extract no more than a given amount of information from incoming signals. Formally, this would be done by replacing the cost function with one on which $C(\Pi_2)$ is zero for expected entropy differences $E[I(c')] - I(c)$ below an ability-determined threshold $T$, and infinite for $E[I(c')] - I(c) \geq T$. Sims (2003) originally studied models of this sort. In such models, it is hard to argue that attention is being treated as a costly parameter: attention has the function of limiting the options (attention policies) that it is possible for agents to take, rather than imposing costs on those options. If that is still not satisfactory, we could

19

do away with $C$ altogether and simplify specify the space of available attention policies as those meeting the entropy constraint set by the agent's abilities. This is no more than the standard ability limitation imposed on choice sets by the constraint that options must be things which an agent can do (Hedden 2012), and need not amount to reintroducing a cost for gathering or processing evidence.

Second, it is a near-universal fact that agents must process incoming information through attentionally limited processes. This means that if the notion of cost-free evidence is meant to apply only to agents who are not attentionally limited, there may be few if any actual instances of cost-free evidence gathering in the world today. This is important because recent work in bounded rationality and non-ideal theory has stressed that rationality results derived by abstracting away from paradigmatic cognitive bounds may provide misleading guidance to agents once their bounds are reintroduced (Lipsey and Lancaster 1956; Thorstad 2024; Wiens 2020). This puts some pressure on the demand to restrict study to agents for whom not only evidence-gathering, but also evidence-processing is cost-free, insofar as constraints on evidence-processing are a ubiquitous and important fact of agents' lives, and insofar as rational evidence-gathering under processing constraints may be most helpfully studied by incorporating those constraints into normative models directly.

Third, even if the model of Section 4 is read as making a point about costly evidence, it is not identical to the two most familiar points. One familiar point is that when evidence *gathering* is costly, agents may prefer to avoid gathering evidence. But that is not the point made in this paper: the costs of evidence gathering are borne by the sender $S$.

A second point is that when evidence *processing* is costly, agents may choose not to process some items of incoming evidence because processing is not worth the expense (Lieder and Griffiths 2020). Typically, this point is made by assuming that items of information are fed to the agent in succession, and she chooses whether or not to process each incoming item of information. Information is processed if the expected of processing exceeds the expected cost, and otherwise the information is discarded. Such models

explain why agents may be indifferent to receiving full information, rather than partial disclosure. However, they do not explain why agents may strictly prefer to avoid full disclosure. Below (Section 6.2) I show that the model of Section 4 has this feature.

## 6.2 A moderate pattern

Consider two questions. First, may agents rationally prefer receiving no information to receiving partial information? Second, may agents rationally prefer receiving partial information to receiving full information? Typically, these questions are answered in the same way. Standard philosophical appeals to risk-aversion (Buchak 2010; Campbell-Moore and Salow 2020), conditionalization failures (Neth forthcoming), act-dependent states (Adams and Rosenkrantz 1980) and other factors rationalize a preference for no information over partial information, and also rationalize a preference for partial information over full information.

For some readers, this may be a desirable consequence, and those readers are welcome to supplement the present account with their favorite of the rationalizing factors surveyed above. However, some readers may balk at the idea of strictly preferring no information to partial disclosure. Surely, they might object, agents should be at most indifferent between these prospects, since it is always an option to simply ignore incoming information. Moreover, many authors have noticed that the predicted situations in no information is strictly preferred to partial disclosure often look like rational defects (Campbell-Moore and Salow 2020).

The model in this paper grounds a more moderate pattern of preference. The agent modeled may strictly prefer partial disclosure to full disclosure, on the grounds that disclosure of unimportant information overwhelms her ability to attend to what is important. But she cannot strictly prefer no information to partial disclosure. After all, it is always possible for her to pay no attention to incoming information. She does this by choosing an attentional policy $\Pi_2$ which provides no information about the world state. For example, she might set $|\mathcal{S}'| = 1$ and $\pi_2(s'_1|s_j) = 1$ for all $j$. She is then certain that her priors will be

unchanged by conditionalization on $\mathcal{S}'$, yielding no improvement in the expected value of her resulting action, but also no cost of attending to $\mathcal{S}'$.

An advantage of this result is that it captures one core intuition behind Good's Theorem: that an agent can never strictly prefer receiving no information to some information. However, this result shows that the intuition does not ground the further claim that agents cannot strictly prefer receiving partial information to full information. It may have been thought that this further claim would follow for exactly the same reasons as the first. However, the model in this paper shows that it is not merely possible, but perhaps also normatively necessary to distinguish these two claims.

## 6.3   Bounded rationality

We saw in Section 4 that the model in this paper is best understood through the lens of bounded rationality. Theories of bounded rationality stress that agents have limited cognitive abilities, and incur costs for using those abilities in given problem environments (Lieder and Griffiths 2020; Gigerenzer and Selten 2001; Simon 1959). Limited abilities should be incorporated in rational models because ought implies can, and costs should be incorporated because they worsen the outcome of agents' acts (Thorstad 2024).

Bounded rationality is a theory of rationality, not irrationality. It says that once an agent's decision problem is fully described, behaviors that may previously have seemed irrational are in fact fully rational, and behaviors that may have seemed rational are in fact irrational. For example, once we recognize that agents are attentionally limited, we see that seeking full over partial disclosure may sometimes be irrational due to an agent's inability to process all disclosed information, or to the costs of processing fully disclosed information.

Some readers may wish to distinguish theories of bounded rationality from theories of a distinctive kind of unbounded or ideal rationality, which does not take into account agents' cognitive limitations or the costs of exercising them (Carr 2022). While this is not my view, readers are welcome to read the discussion in this paper with such a view in

mind. In this case, the right conclusion would be that it may be boundedly, though not unboundedly rational to strictly prefer partial to full disclosure.

Readers are also invited to apply a distinction between the rationality of an agent's inquisitive *acts*, such as gathering and attending to evidence, and the rationality of the resulting belief *states* (Thorstad 2024). The model in this paper suggests that the action of gathering partial rather than full evidence may sometimes be rational, and likewise that the action of fully attending to all incoming evidence may be irrational. However, it does not yet pronounce on the rationality of the resulting belief state. For example, it might be maintained that although the act of gathering and processing full evidence would be wastefully irrational, the belief that results would be rational so long as it results from conditionalization on an agent's total evidence. In this way, we may begin to see a split between conditionalization as a diachronic norm on belief states and as a norm on the processing of incoming information. The discussion in this paper suggests that it is not always rational, or even possible for agents to take the action of attending to all incoming information during processing, and in this sense there may be no rational process whose outputs coincide with Bayesian conditionalization. However, it does not yet cast doubt on the idea that after learning some evidence $E$, the credal state $c(*|E)$ would be rational, perhaps even uniquely rational, for the agent to occupy.[7]

## 6.4  Conditionalization, learning and total evidence

I said in Section 5 that while our agent is certain that she will conditionalize, she is not certain that she will conditionalize on her total evidence. Roughly the idea was that an agent's total evidence expands after a learning experience by adding the totality of what is learned. In this case, $R$ learns the contents $s$ of the incoming signal $\mathcal{S}$, so conditionalizing on her total evidence requires her to update to $c(*|s)$. However, $R$ updates instead on an attentionally distorted signal $s'$, and hence fails to update on her total evidence.

---

[7]This distinction may be helpful in responding to proceduralist attacks on conditionalization based on rational delay (Podgorski 2017) and other bounded rationality phenomena.

Some internalist readers may disagree with this verdict.[8] For example, mentalists (Conee and Feldman 2001, 2004) hold that only mental states can justify belief. Insofar as evidence justifies belief, it seems plausible that incoming information which is not attended to cannot be evidence for a mentalist, since it is not incorporated into the agent's mental states. Likewise, access internalists (BonJour 1985; Chisholm 1966) might hold that evidence is subject to strong accessibility requirements which are unlikely to be met by information which is not attended to and is thereby lost to the agent forever.

Familiar debates ensue, and it is not my intention to re-open these debates here. For my part, I am not ready to accept anything like the accounts described above. I would rather accept that it is sometimes rational to prefer partial disclosure to full disclosure than to accept internalist accounts of evidence and justification. But I would be remiss if I did not acknowledge that the present discussion reveals an important advantage of some internalist accounts. If the right conclusion is that internalist accounts of evidence are correct, then that is an important lesson too.

# 7  Conclusion

This paper explored the possibility that agents may rationally decline cost-free evidence because they are uncertain whether they will conditionalize on the evidence gathered. Section 2 explored one foundation for this possibility based on rational modesty. Section 3 noted that Neth's presentation of this possibility has three limitations: it requires agents to lack control over their future actions, to suspect that they may fail to conditionalize, and that their future selves may be irrational.

Section 4 adapted a model from the literature on Bayesian persuasion to show how attentionally limited agents may prefer partial over full disclosure of their current situation, even if they are certain that they will optimally attend to incoming information based on full knowledge of the statistical distribution through which it is generated. Section 5

---

[8]That is not to say that we need anything approaching the stronger externalist conditions of Das (2023) to escape internalist reinterpretations of the cases in this paper.

showed how this model removes the limitations of the modesty-based argument. This has the effect of strengthening the case for thinking that agents may rationally decline cost-free evidence due to uncertainty about whether they will conditionalize on the evidence gathered. Section 6 concluded by discussing the relationship between attention and costly evidence; a moderate feature of the pattern of evidence-aversion exhibited by attentionally limited agents; connections between attentional limitations and bounded rationality; and a surprising avenue of support for internalist theories of evidence.

It may be productive for future work to extend this result and the results of Neth (forthcoming) to explore further ways in which agents may rationally decline cost-free evidence because they are uncertain whether they will conditionalize on the evidence gathered. For example, future work might investigate the relationship between upstream failures to gather evidence due to attentional limitations and downstream failures to gather evidence because of challenges associated with storing and retrieving evidence from memory once it is gathered (Harman 1986; Schooler and Hertwig 2005). It may also be productive to explore the robustness of this paper's results to alternative accounts of rational attention. Finally, the discussion in this paper is conditional on a successful defense of the relevant theory of bounded rationality, so the results of this paper give renewed importance to elaborating the nature, scope and justification for theories of bounded rationality.

# References

Adams, Ernest and Rosenkrantz, Roger. 1980. "Applying the Jeffrey decision model to rational betting and information acquisition." *Theory and Decision* 12:1–20.

Archer, Sophie (ed.). 2021. *Salience: A philosophical inquiry*. Routledge.

Ben-Shahar, Omri and Schneider, Carl. 2016. *More than you wanted to know: The failure of mandated disclosure*. Princeton University Press.

Bermúdez, José Luis (ed.). 2018. *Self-control, decision theory and rationality: New essays*. Cambridge University Press.

Blackwell, David. 1953. "Equivalent comparisons of experiments." *The Annals of Mathematical Statistics* 24:265–72.

Bohn, Roger and Short, James. 2009. "How much information? 2009 Report on American Consumers." Global Information Industry Center, http://hmi.ucsd.edu/howmuchinfo.php.

BonJour, Lawrence. 1985. *The structure of empirical knowledge*. Harvard University Press.

Bradley, Seamus and Steele, Katie. 2016. "Can free evidence be bad? Value of information for the imprecise probabilist." *Philosophy of Science* 83:1–28.

Briggs, R.A. and Pettigrew, Richard. 2020. "An accuracy-dominance argument for conditionalization." *Noûs* 54:162–81.

Buchak, Lara. 2010. "Instrumental rationality, epistemic rationality, and evidence-gathering." *Philosophical Perspectives* 24:85–120.

Campbell-Moore, Catrin and Salow, Bernhard. 2020. "Avoiding risk and avoiding evidence." *Australasian Journal of Philosophy* 98:495–515.

Caplin, Andrew and Dean, Mark. 2015. "Revealed preference, rational inattention, and costly information acquisition." *American Economic Review* 105:2183–2203.

Carr, Jennifer. 2022. "Why ideal epistemology?" *Mind* 131:1131–62.

Castro, Clinton and Pham, Adam. 2020. "Is the attention economy noxious?" *Philosophers' Imprint* 20:1–13.

Chisholm, Roderick. 1966. *Theory of knowledge*. Prentice-Hall.

Christensen, David. 2007. "Epistemology of disagreement: The good news." *Philosophical Review* 116:187–217.

Conee, Earl and Feldman, Richard. 2001. "Internalism defended." *American Philosophical Quarterly* 38:1–18.

—. 2004. *Evidentialism*. Oxford University Press.

Das, Nilanjan. 2023. "The value of biased information." *British Journal for the Philosophy of Science* 74:25–55.

Davenport, Thomas and Beck, John. 2001. *The attention economy: Understanding the new economy of business*. Harvard Business School Press.

Elster, Jon. 2009. *Ulysses unbound: Studies in rationality, precommitment, and constraints*. Cambridge University Press.

Falagas, Matthew, Korbila, Joanna, Giannopoulou, Konstantina, Kondilis, Barbara, and Peppas, George. 2009. "Informed consent: How much and what do patients understand?" *American Journal of Surgery* 198:420–35.

Gabaix, Xavier. 2014. "A sparsity-based model of bounded rationality." *Journal of Economic Literature* 1661–1710.

Gauthier, David. 1997. "Resolute choice and rational deliberation: A critique and a defense." *Noûs* 31:1–25.

Gigerenzer, Gerd and Selten, Reinhard (eds.). 2001. *Bounded rationality: The adaptive toolbox*. MIT Press.

Golman, Russell, Hagmann, David, and Lowenstein, George. 2017. "Information avoidance." *Journal of Economic Literature* 55:96–135.

Good, I.J. 1966. "On the principle of total evidence." *British Journal for the Philosophy of Science* 17:319–321.

Greaves, Hilary and Wallace, David. 2006. "Justifying conditionalization: Conditionalization maximizes expected epistemic utility." *Mind* 115:607–32.

Harman, Gilbert. 1986. *Change in view*. MIT Press.

Hedden, Brian. 2012. "Options and the subjective ought." *Philosophical Studies* 158:343–360.

Hertwig, Ralph and Engel, Christoph (eds.). 2021. *Deliberate ignorance*. MIT Press.

Irving, Zachary and Glasser, Aaron. 2020. "Mind wandering: A philosophical guide." *Philosophy Compass* 15:e12644.

Kadane, Joseph, Schervish, Mark, and Seidenfeld, Teddy. 2008. "Is ignorance bliss?" *Journal of Philosophy* 105:5–36.

Kamenica, Emir and Gentzkow, Matthew. 2011. "Bayesian persuasion." *American Economic Review* 101:2590–2615.

Lacko, James and Pappalardo, Janis. 2004. "The effect of mortgage broker compensation disclosures on consumers and competition: A controlled experiment." Federal Trade Commission Bureau of Economics Staff Report, https://www.ftc.gov/reports/effect-mortgage-broker-compensation-disclosures-consumers-competition-controlled-experiment.

Lester, Benjamin, Persico, Nicola, and Visschers, Ludo. 2012. "Information acquisition and the exclusion of evidence in trials." *Journal of Law, Economics and Organization* 28:163–82.

Lieder, Falk and Griffiths, Thomas. 2020. "Resource-rational analysis: Understanding human cognition as the optimal use of limited computational resources." *Behavioral and Brain Sciences* 43:E1.

Lipnowski, Elliott, Mathevet, Lauren, and Wei, Dong. 2020. "Attention management." *American Economic Review: Insights* 2:17–32.

Lipsey, Richard and Lancaster, Kelvin. 1956. "The general theory of second best." *Review of Economic Studies* 24:11–32.

Maćkowiak, Bartosz, Matŏjka, Filip, and Wiederholt, Mirko. 2023. "Rational inattention: A review." *Journal of Economic Literature* 61:226–73.

Maher, Patrick. 1990a. "Symptomatic acts and the value of evidence in causal decision theory." *Philosophy of Science* 57:479–98.

—. 1990b. "Why scientists gather evidence." *British Journal for the Philosophy of Science* 41:103–119.

Neth, Sven. forthcoming. "Rational aversion to information." *British Journal for the Philosophy of Science* forthcoming.

Podgorski, Abelard. 2017. "Rational delay." *Philosophers' Imprint* 17:1–19.

Schooler, Lael J. and Hertwig, Ralph. 2005. "How forgetting aids heuristic inference." *Psychological Review* 112:610–28.

Shannon, Claude. 1948. "A mathematical theory of communication." *The Bell System Technical Journal* 27:379–423.

Siegel, Susanna. 2017. *The rationality of perception*. Oxford University Press.

Simon, Herbert. 1959. *Models of man*. Wiley.

—. 1971. "Designing organizations for an information-rich world." In Martin Greenberger (ed.), *Computers, communications, and the public interest*, 37–72. Johns Hopkins Press.

Sims, Christopher. 2003. "Implications of rational inattention." *Journal of Monetary Economics* 50:665–90.

Thorstad, David. 2024. *Inquiry under bounds*. Oxford University Press.

—. forthcoming. "Why bounded rationality (in epistemology)?" *Philosophy and Phenomenological Research* forthcoming.

Titelbaum, Michael. 2015. "Rationality's fixed point." *Oxford Studies in Epistemology* 5:253–94.

van den Berg, Ronald and Ma, Wei Ji. 2018. "A resource-rational theory of set size effects in human visual working memory." *eLife* 7:e34963.

Watzl, Sebastian. 2021. "The ethics of attention: An argument and a framework." In Sophie Archer (ed.), *Salience: A philosophical inquiry*, 89–112. Routledge.

Wiens, David. 2020. "The general theory of the second best is more general than you think." *Philosophers' Imprint* 5:1–26.