# A Communication-First Account of Explanation

Jacqueline Harding, Tobias Gerstenberg, Thomas Icard May 7, 2025

#### Abstract

This paper develops a formal account of causal explanation, grounded in a theory of conversational pragmatics, and inspired by the interventionist idea that explanation is about asking and answering what-if-things-had-been-different questions. We illustrate the fruitfulness of the account, relative to previous accounts, by showing that widely recognized "explanatory virtues" emerge naturally, as do subtle empirical patterns concerning the impact of norms on causal judgments. This shows the value of a "communication-first" approach to explanation: getting clear on explanation's communicative dimension is an important prerequisite for philosophical work on explanation. The result is a simple but powerful framework for incorporating insights from the cognitive sciences into philosophical work on explanation, which will be useful for philosophers or cognitive scientists interested in explanation.

## Contents

1	Intr	oduction	3
2	Exp	lanation and Causation	4
	2.1	Early Formal Accounts	4
	2.2	Causal Explanation and Interventions	6
	2.3	Explanation and Actual Causation	8
		2.3.1 Halpern and Pearl on Explanation	8
		2.3.2 Issues with EX2: listener uncertainty about which events are causes	6
		2.3.3 Issues with EX3: longer explanations are often better	. 11
		2.3.4 Issues with EX4: explanations often cite known events	. 12
		2.3.5 Summing Up	. 13
3	A C	ommunication-First Account	14
	3.1	The Literal Listener	. 14
	3.2	The Pragmatic Speaker	
		3.2.1 Useful Utterances	
		3.2.2 Production and Processing Costs	
	3.3	What is the Listener's Decision Problem?	
		3.3.1 The Manipulation Game	
	3.4	The Pragmatic Listener and the Goodness of an Explanation	
	3.5	Summing Up	

4	$\mathbf{E}\mathbf{x}\mathbf{p}$	Explaining Explanation			
	4.1	Downs	stream Interests	20	
		4.1.1	Back to Roof Replacement	20	
		4.1.2	Forward-Looking Decision Problems and Invariance	21	
		4.1.3	Explanations Identify Good Points of Intervention	21	
		4.1.4	Pragmatic Theories of Explanation	22	
	4.2	Backgr	round Knowledge	23	
		4.2.1	Known causes can be informative	24	
	4.3	Explar	natory Relationships and Background Conditions	25	
		4.3.1	Causal Selection Interacts with Normality and Structure	25	
		4.3.2	Inferences from Normality	27	
	4.4	Minim	ality and Simplicity	27	
		4.4.1	Simplicity	29	
	4.5	Propor	rtionality and Levels of Explanation	30	
5	Disc	cussion	and Future Work	31	
6	Ack	nowled	dgments	32	

## 1 Introduction

The philosophical problem of explanation is usually taken to be a matter of locating the right explanatory relations in the world, with human interests and capacities offered at best a marginal role. Although philosophers of explanation will agree that some explanatory virtues are better accounted for in terms of 'pragmatic' factors (c.f. those involving background knowledge or downstream interests), they take this to show that these virtues are peripheral to the core of the explanatory project. Strevens (2008, p.7) gives clear voice to this perspective: "To discern the nature of explanation in the ontological sense, you must acquire the ability to see past the communicative conventions and strategies of scientists to the explanatory facts themselves. A philosopher of explanation will therefore occasionally discuss communicative conventions just as an astronomer might study atmospheric distortion so as to more clearly see the stars".

In particular, it is tempting to think of the philosophy of explanation and the psychology of explanation as separate subjects. While the philosophical literature has exerted considerable influence on the way psychologists investigate explanation (Hilton, 1990; Lombrozo, 2006; Keil, 2006; Goddu and Gopnik, 2024), philosophers have been more reluctant to incorporate insights from psychology and cognitive science into their accounts. Even interventionists, who tend to motivate their account by gesturing to the downstream usefulness of causal information, eschew pragmatic factors in the account's actual development (Woodward, 2003; Woodward and Hitchcock, 2003).

Here's the question, then, which motivates our paper: what would happen if we took pragmatic factors seriously in developing an account of explanation, rather than sidelining them? Building on recent advances in cognitive science and formal pragmatics (Frank and Goodman, 2012; Goodman and Frank, 2016b; Degen, 2023; Sumers et al., 2023; Beller and Gerstenberg, 2024), we direct attention to the question of what a speaker is doing when she offers an explanation to a listener. Our main claim is that several hallmarks of explanation – including key features that were explicitly built in to previous formalisations – simply (and softly) emerge from the dynamics of conversation together with the minimal assumption that "why" and "because" communicate relations of (typically causal) dependence.

Of course, we're not the first to take pragmatic factors seriously. Many philosophers have stressed the relevance of conversational and cognitive considerations in characterising explanation (Bromberger, 1965; van Fraassen, 1977, 1980; Achinstein, 1983; Ylikoski and Kuorikoski, 2010; De Regt, 2017; Potochnik, 2017), but these proposals have remained programmatic and informal. By contrast, our formal account is developed using tools from contemporary cognitive science; not only does this afford precise comparison with other formal proposals, it suggests a promising future avenue for incorporating psychological work into the philosophy of explanation.

Following van Fraassen (1980), work on explanation which emphasises pragmatic factors tends to be associated with a more broad-sweeping claim: that there's nothing to distinguish explanation from other sorts of communication. We do not make this claim here. Rather, our claim is that getting clear on explanation's communicative dimension is an important prerequisite for philosophical work on explanation. (It's for this reason that we call our proposal a "communication-first"

<sup>&</sup>lt;sup>1</sup>On Woodward's account, for example, contextual features serve only to pin down the 'contrast class' for the explanandum expressed by a "why?" question; once these have been fixed, the goodness of an explanation is entirely independent of the communicative context in which the explanation occurs (Woodward, 2003, Sect 5.12). To compare different explanations of the same explanandum, interventionists develop concepts like 'proportionality' and 'stability' in explicitly non-pragmatic terms (Woodward 2010; see also the exchange between Franklin-Hall (2016) and Woodward (2021a)).

account.) To stretch Strevens' simile, our suggestion is that – when it comes to explanation – much of what's said about the stars can be accounted for by paying more attention to atmospheric distortion. Although we situate our ideas within a broadly interventionist account of explanation, we anticipate that the framework will be useful for a wide variety of philosophers and cognitive scientists interested in explanation.

# 2 Explanation and Causation

What primitives are needed to account for all of what is interesting and important about explanations? Any theory of explanation must answer this question. This is especially true of formal accounts, which demand precise specification of their mathematical building blocks. In this section, we introduce and motivate our theory's primitives through comparison with other formal accounts.

## 2.1 Early Formal Accounts

The first formal theories of explanation involved sets of sentences closed under logical deduction (that is, logical "theories"; see Hempel and Oppenheim (1948)), statistical models and probability distributions (Douven, 2022; Hempel, 1965; Salmon, 1971), and even physical models of spacetime (Dowe, 1992; Salmon, 1998). Some of these accounts are given an ontic gloss (e.g. Salmon 1998), in that the components are intended to refer to objective features of the world; others are given a more epistemic gloss (e.g. Hempel and Oppenheim 1948), in that they are relative to a body of evidence. Numerous authors object that none of these earlier accounts carve explanation at the right joints (Achinstein, 1983). In particular, they pay insufficient attention to key contextual factors pertaining to a receiver of the explanation, including what may have prompted the need for an explanation in the first place, and what it would take to satisfy that need.

As highlighted already by Hempel and Oppenheim (1948), and stressed especially by Bromberger (1965), explanations can be conceived as answers to "why?" questions. Perhaps, then, some of the structure of explanations can be read off the grammar of (answers to) "why?" questions. For instance, asking "Why FACT?" (where FACT is some explanandum, a proposition to be explained) is typically understood relative to a set of relevant alternatives to FACT. Moreover, this set of alternatives will often depend on contextual factors, including the interests of those posing the question. Such factors may affect both what is being asked, and what would be considered a good answer (van Fraassen 1980; see also Woodward 2003, Section 5.12 and Ylikoski 2007). For example, the question "Why did Adam eat the apple?" can – depending on whether stress is placed on "Adam" or "apple" (Dretske, 1972) – be a request for information about why Adam (rather than someone else) ate the apple, or about why Adam ate the apple (rather than something else) (van Fraassen, 1980, p.127). Proponents of earlier accounts are not unaware of these 'pragmatic' and contextual factors, but regard them as peripheral to what is most interesting about explanation. Such pragmatic factors hence play no formal role in their positive frameworks.

Others endeavour to put such matters front and center, in part because they feel that pragmatic and contextual dependence undermines previous accounts (Bromberger, 1965, 1984; Achinstein, 1977, 1983, 1984; van Fraassen, 1977, 1980; De Regt and Dieks, 2005; De Regt, 2017; Potochnik, 2017). To illustrate their motivations, consider a variant of an example from Gärdenfors (1980).

**Example 1** (Holiday Tan). Alice sees her friend Bob after returning from a work trip to Rome. She is noticeably more tanned than when Bob last saw her; envious, he asks, 'why are you so

tanned?'. She replies, 'I was just in Rome!'.

We can make two observations here. First, given what Bob already knows, he can — upon hearing Alice's answer — reasonably infer that it was sunny when she was in Rome, that she spent time outside while she was there, and so on. What makes it seem like a good explanation is, at least in part, that Bob is able to draw all of these conclusions from Alice's answer. Second, if Bob also happens to know that it was raining in Rome all week, then it may not be as good of an explanation. It may also be a poor explanation if Bob already knows that Alice was in Rome. In such contexts, the explanatory question evidently persists even after Alice's statement.

As the example makes clear, then, Bob's epistemic state is relevant in at least two ways. First, it helps determine what the range of possible answers could be in the first place. Second, it draws attention to what explanations do for us. In the example, Alice's answer to Bob's question helps demystify a fact that might have seemed initially surprising, or worthy of explanation. That is, given what he already knows, he is able to trace the sequence of events that led to Alice's tan; this pinpoints an important sense in which Alice's answer provides a good explanation. Rather than being peripheral to explanation, then, pragmatic facts (such as Bob's background knowledge and downstream interests) seem crucial in accounting for why we have a practice of giving and receiving explanations at all.

Assuming we take these considerations seriously, what should our account of explanation look like? Famously, van Fraassen (1980) argues that context is a central and irreducible feature of explanation, and that if we completely draw out the repercussions of context dependence, then no special explanatory relation may even be needed. The role of explanation, on this view, is simply to "satisfy our desires [...] for descriptive information" (p. 153). While van Fraassen (1977, 1980) offers a sketch of an account of explanation along these lines, too little detail is given to enable proper comparison with competing accounts. As Kitcher and Salmon (1987) observe, an unanalysed notion of relevance plays a fundamental role in van Fraassen's proposal; characterising this relevance notion, though, seems exactly what's at stake in giving a theory of explanation (Strevens, 2008).

A more fully developed account is offered in a series of papers and books by Gärdenfors (1980, 1988, 1990). The key move is to put the *epistemic state*  $\mathcal{K}$  of the receiver of the explanation (Bob, in the example above) front and center. For Gärdenfors, as for others in the pragmatic tradition, explanation can only be characterised and evaluated relative to this epistemic state. It is assumed that facts demanding explanation are those which are surprising to an agent, in the sense that, prior to learning them, the agent assigned them relatively low probability. Explaining why FACT, on Gärdenfors's view, is a matter of specifying probabilistic or factual information that would have – in the agent's epistemic state prior to learning FACT – rendered it less surprising. (In the example above, the statement that Alice was in Rome counts as an explanation because it renders her tan less surprising to Bob, given his epistemic state.) While the account is thus reminiscent of Hempel's original deductive-nomological account (Hempel, 1965), the details are explicated in a very different way, by appeal to the theory of belief revision (Gärdenfors, 1988).

<sup>&</sup>lt;sup>2</sup>This idea is developed more recently by Chandra et al. (2024).

## 2.2 Causal Explanation and Interventions

As illustrated above, a theme that runs through the literature on pragmatic approaches is that there is nothing fundamentally different about explanation above and beyond "mere" description.<sup>3</sup> What is special about explanation is the way in which appropriate descriptive information resolves an agent's uncertainty through interaction with pragmatic context.

By contrast, a common view is that explanation fundamentally involves relations of asymmetric dependence, chiefly causal dependence (Woodward and Hitchcock, 2003; Woodward, 2003; Strevens, 2008). This view is nicely summarized by Salmon (1977, p. 162), "To give [...] explanations is to show how events [...] fit into the causal network of the world". The basic thought is that, e.g., telling someone that you took a walk in the park today (a description of what happened) is not typically an explanation, unless it helps address a causal-explanatory question (e.g., why your shoes were muddy, what caused them to be muddy). To borrow an example from Salmon (1971, p.34) discussed by Woodward (2003, Sect. 4.2.), learning that a (cisgender) man takes birth control pills may – given suitable background assumptions about one's epistemic state – lower one's surprise at his failure to become pregnant (i.e. satisfy Gärdenfors's (1980) criterion), but it does not explain this failure. As Salmon and Woodward observe, the issue is that the man's taking birth control pills does not play a causal role in the failure. Indeed, in Example 1, Alice's response seems like an explanation just in case (relative to Bob's knowledge state) it provides information about what caused her tan. It may be helpful for some other purpose; it would just not count as an explanation.

To say that causality is central to explanation is, of course, only to raise the question of how we ought to analyse causality. For many decades, a dominant view in philosophy reduced it essentially to probability, e.g., declaring that A causes B when A occurs before, and also raises the probability of, B (e.g., Suppes 1970). (Employing a closely related analysis – one of probability raising relative to a "contracted" belief state in which the explanans is removed – Gärdenfors (1990) endeavours to incorporate causality into his account of explanation.)

For a host of reasons, purely probabilistic analyses of causality have fallen out of favour, leading instead to a variety of formal frameworks centred around the fundamental notion of *causal intervention* (see, e.g., Spirtes et al. 1993; Pearl 1995; Hitchcock 2001; Woodward 2003, among many others). Roughly speaking, an intervention on a system is an ideal, exogenous manipulation of some component of the system, which leaves all other components unchanged. This idea has been made more precise with the help of various mathematical models, among the most general of which is the *structural causal model* (see Pearl 1995, 2009; Peters et al. 2017; Bareinboim et al. 2022):

**Definition 1** (Structural Causal Model). A structural causal model (SCM)  $\mathcal{M}$  is a 4-tuple

$$(\mathbf{U}, \mathbf{V}, \mathbf{f}_{\mathbf{V}}, P(\mathbf{U}))$$

where:

- U is a set of exogenous variables, with possible values Val(U) for each  $U \in U$ ;
- V is a set of endogenous variables, with possible values Val(X) for each  $X \in V^4$

<sup>&</sup>lt;sup>3</sup>This stance has longstanding empiricist roots, for example, reflected in the following quotation from Pearson (1911): "Nobody believes now that science *explains* anything; we all look upon it as a shorthand description, as an economy of thought" (p. xi, also quoted in Salmon 1978).

<sup>&</sup>lt;sup>4</sup>We are assuming, without loss, that  $Val(X) \cap Val(Y)$  whenever  $X \neq Y$ . Note also that  $Val(X) = \bigcup_{X \in X} Val(X)$ .

- $\mathbf{f}_{\mathbf{V}}$  is a set of *structural functions*, where  $f_X : \mathsf{Val}(\mathbf{V} \cup \mathbf{U}) \to \mathsf{Val}(X)$  for each endogenous variable  $X \in \mathbf{V}$ ;
- $P(\mathbf{U})$  is a probability distribution over  $Val(\mathbf{U})$ .

The notion of intervention on an SCM is captured by mechanism replacement: intervening to set variable X to value x is a matter of replacing the structural function  $f_X$  with the constant function sending all arguments to x.

**Example 2** ( $\mathcal{M}_{(A \wedge B) \vee C}$  Causal Model). Consider the simple scenario in Example 1. We model this with three endogenous variables A, B, C, each binary valued, so that  $\mathsf{Val}(X) = \{0, 1\}$  for X = A, B, C, and three binary exogenous variables  $U_A, U_B, U_C$ . E represents whether Alice is tanned (E = 1 when she is, 0 when not); A represents whether Alice is in Rome; B represents whether it is sunny in Rome; and C represents whether Alice went to a tanning salon. The structural functions have E true when either C is true, or both A and B are true.

Meanwhile, A, B, and C are all determined exogenously, taking on the value of their respective exogenous variables  $(U_A, U_B, \text{ and } U_C)$ . That is,  $f_A$ ,  $f_B$ , and  $f_C$  are all the identity function. See Figure 1, where arrows represent functional relationships. (To fully specify the model  $\mathcal{M}_{(A \wedge B) \vee C}$ , we would also need to define a probability distribution  $P(U_A, U_B, U_C)$ , which in the simplest cases would factor as a product  $P(U_A) \cdot P(U_B) \cdot P(U_C)$ .)

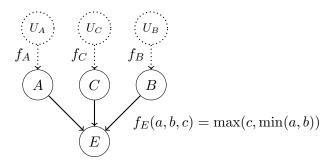


Figure 1: Causal Model  $\mathcal{M}_{(A \wedge B) \vee C}$ 

Structural causal models – and related frameworks, like *causal Bayesian networks*<sup>5</sup> – have risen to prominence not only in philosophy, but also in psychology and cognitive science, where a number of researchers have suggested that intuitive causal cognition employs representations, and operations over them, that resemble causal models with interventions (Rottman and Hastie, 2016; Waldmann, 2017; Lombrozo, 2006; Goddu and Gopnik, 2024; Lagnado et al., 2013).

How exactly do causal interventions fit into a broader theory of explanation? Most proponents of interventionist approaches to explanation agree with something like the following picture, articulated by (Woodward and Hitchcock, 2003, p. 21): "To explain why some phenomenon occurs is to [provide] the resources for answering a variety of what-if-things-had-been-different questions: how would the outcome have differed if the initial conditions had been changed in various ways?", where "changed" is to be understood as, "changed by ideal intervention."

While this is a rather different conception from the approaches to explanation that came before, it leaves open lots of important details. Which what-if-things-had-been-different questions are most

<sup>&</sup>lt;sup>5</sup>See Bareinboim et al. (2022, Theorems 2 and 4) for the precise relationship between these models.

important? What kinds of "resources" for answering such questions are suitable, and how are they to be provided?

## 2.3 Explanation and Actual Causation

We propose that a natural approach to answering these questions pays attention to the situation in which the explanation is being given and received. This approach combines the interventionist idea that explanations communicate patterns of (causal) dependence with explicit modelling of more 'pragmatic' factors, such as the epistemic state of the receiver of the explanation.

Arguably the most celebrated and influential formalisation of this basic approach – influencing not just philosophy, but also cognitive science and computer science – is due to Halpern and Pearl (2005a). They present a conceptually conservative extension of Gärdenfors's epistemic approach, but with one fundamental new ingredient, namely appeal to a notion of *actual causation* (Halpern and Pearl, 2005a,b), that is itself captured using structural causal models.

An "actual cause" of some token event E is a factor C that causally contributed to E's occurrence. A common way to think about this causal contribution is counterfactual, employing the so-called "but-for" test: if E would not have occurred but for C, then we say that C was an actual cause of E. In Example 1, Alice's being in Rome was an actual cause of her tan – had she not been in Rome, she would not have tanned.

While this simple but-for analysis is widely applicable, we often want to identify some factor as an actual cause even when it fails the but-for test. Suppose, for example, that Alice not only went to sunny Rome (A = 1 and B = 1), but she also went to a tanning salon (C = 1). We might want to identify the tanning salon as a cause of her tan even though, had she not gone, the Roman sun still would have tanned her. Such cases of overdetermination – and many other puzzle cases for the but-for test – have motivated a variety of analyses of actual cause (see, e.g., Hitchcock 2001; Woodward 2003; Gallow 2021, among many others). In Part I of their paper (see Halpern and Pearl 2005a), Halpern and Pearl offer their own counterfactual account of what it takes for a factor to be an actual cause (with several variations summarised in Halpern, 2016).

The rough idea behind these and other counterfactual analyses is that C should be a but-for cause under some contingency in which we hold the values of some variables (off a "causal pathway" from C to E) fixed to some (possibly non-actual) values. In Example 1, note that C=1 is a but-for cause of E=1 if, for instance, we imagine what would have happened had Alice not been in Rome (that is, assume A=0). Different proposals make this intuition precise in subtly different ways. Such nuances will not matter for our purposes; we only need to assume that we have fixed some "egalitarian" analysis of what it is to be a causally contributing factor (Bebb and Beebee, 2024). In particular, our proposal is compatible with any account of actual causation which treats both events as actual causes in cases of overdetermination. In the case above, for example, we assume that A=1, B=1 and C=1 all count as actual causes of E=1.

## 2.3.1 Halpern and Pearl on Explanation

Following Gärdenfors (1988), Halpern and Pearl (2005a) suppose that explanation is relative to a receiver's epistemic state. They model this epistemic state  $\mathcal{K}$  is as a set of pairs  $(\mathcal{M}, \mathbf{u})$ , where  $\mathcal{M}$  is a causal model and  $\mathbf{u}$  is a "context" that specifies values for all the exogenous variables. Intuitively, these are the causal situations that are consistent with the person's knowledge.<sup>6</sup> It is

<sup>&</sup>lt;sup>6</sup>For brevity's sake, we often refer to these causal situations as 'worlds'.

supposed for simplicity that all models in K have the same variables, so that uncertainty is only over the structural relationships and the values of (both exogenous and endogenous) variables.

Given an explanandum FACT,<sup>7</sup> what does it take to be an explanation? Halpern and Pearl (2005a) propose the following (which we call the HP analysis):

**Definition 2** (HP Analysis of Explanation).  $\mathbf{X} = \mathbf{x}$  is an explanation of FACT relative to an epistemic state  $\mathcal{K}$  iff the following conditions hold:

- **(EX1)** FACT is true at all models  $(\mathcal{M}, \mathbf{u}) \in \mathcal{K}$ .
- (EX2) X = x is an actual cause of FACT for all  $(\mathcal{M}, \mathbf{u}) \in \mathcal{K}$  in which X = x is true.
- (EX3) No proper subset of X satisfies EX2.
- **(EX4)** There exists  $(\mathcal{M}, \mathbf{u}) \in \mathcal{K}$  in which  $\mathbf{X} \neq \mathbf{x}$ .

Thus, on Halpern and Pearl's (2005a) picture, the agent is certain that some explanandum holds (EX1). An explanation is then a proposition which was previously unknown to the agent (EX4), and expresses a minimal (EX3) event which is an actual cause of the explanandum in every world in the agent's epistemic state in which it is true (EX2).<sup>8</sup>

In addition to incorporating listener-dependence as in Gärdenfors's proposal, and involving causality in a thoroughgoing way, the HP account enjoys the virtue of formal precision. At the same time, Halpern and Pearl are not entirely clear on how the various formal components relate to actual practices of giving and receiving explanations. Part of this unclarity arises when we try to assess potential counterexamples to the account.

#### 2.3.2 Issues with EX2: listener uncertainty about which events are causes

It's unclear what the motivation for EX2 is; it seems like whether or not  $\mathbf{X} = \mathbf{x}$  counts as an explanation of FACT in a world  $(\mathcal{M}, \mathbf{u})$  should depend only on facts about  $(\mathcal{M}, \mathbf{u})$ . Why should it matter whether the listener also considers some other world  $(\mathcal{M}', \mathbf{u}')$  possible?

To make this concern concrete, consider Example 1 once again. Bob knows that Alice is tanned (that E=1); let's suppose that he also knows the causal structure is given by  $\mathcal{M}_{(A \wedge B) \vee C}$  (Figure 1). Then his epistemic state  $\mathcal{K}$  consists of five pairs ( $\mathcal{M}_{(A \wedge B) \vee C}$ , -): he is uncertain about the five possible settings of exogenous variables that would lead to E=1.9 The five possibilities that Bob entertains are shown in Figure 2 (exogenous variables omitted for readability). Suppose the actual world is given by  $\mathbf{u}_{A,B}$  (that is, Alice went to sunny Rome but did not go to the tanning salon). In this world, A=1, B=1 and the pair  $\langle A=1, B=1 \rangle$  are all actual causes of E=1.

As we saw in Section 2.1, it seems intuitive to think that A = 1 is an explanation of E = 1; given his epistemic state  $\mathcal{K}$ , Bob can use this fact to infer that it was sunny in Rome. But note that EX2 rules this out: since A = 1 is not a cause of E = 1 in  $\mathbf{u}_{A,C}$ , it does not count as an explanation on the HP analysis. The only explanation available to Alice on Halpern and Pearl's picture is  $\langle A = 1, B = 1 \rangle$  (i.e. citing both that she went to Rome and that it was sunny there).

<sup>&</sup>lt;sup>7</sup>Technically, FACT could be any Boolean combination of "basic statements" of the form Y = y, for  $Y \in \mathbf{V}$ .

<sup>&</sup>lt;sup>8</sup>The HP statement of condition EX2 only requires that  $\mathbf{X} = \mathbf{x}$  is a *sufficient cause* of FACT, not an *actual cause* (a weaker condition on their definition of actual causation (Halpern and Pearl, 2005a)). Since our account doesn't depend on any particular account of actual causation, we omit this distinction for ease of presentation.

<sup>&</sup>lt;sup>9</sup>To wit, these are precisely  $\mathbf{u}_C = \langle U_A = 0, U_B = 0, U_C = 1 \rangle$ ;  $\mathbf{u}_{A,B} = \langle U_A = 1, U_B = 1, U_C = 0 \rangle$ ;  $\mathbf{u}_{A,C} = \langle U_A = 1, U_B = 0, U_C = 1 \rangle$ ;  $\mathbf{u}_{B,C} = \langle U_A = 0, U_B = 1, U_C = 1 \rangle$ ; and  $\mathbf{u}_{A,B,C} = \langle U_A = 1, U_B = 1, U_C = 1 \rangle$ .

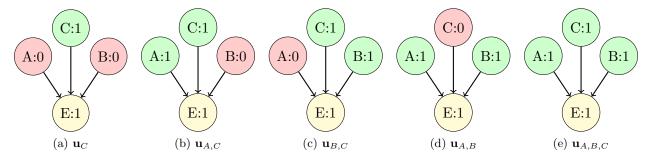


Figure 2: Bob's epistemic state K in Example 1. Red nodes indicate that the variable is false (has value 0), and green true. The explanandum (known to have value 1) is in yellow.

HP are aware of this issue (Halpern and Pearl, 2005b, p.902). They suggest that A=1 be viewed as a partial explanation of E=1. Although this response might seem plausible when the cause named is part of a larger conjunction of causes (as in  $\mathbf{u}_{A,B}$  above, where A=1 is only a cause of E=1 when we also have B=1), it misses the more general problem illustrated by the example. The problem with EX2 is that it doesn't allow explanations to resolve the listener's uncertainty about the causal facts; an event  $\mathbf{X}=\mathbf{x}$  can only feature in an explanation if the listener is certain about the causal story connecting it to the explanandum FACT. But this get things the wrong way round; on (for example) the interventionist picture, a function of explanations is precisely to get the listener to grasp this causal story.

To make this clearer, consider the following example adapted from Faye (2007):

**Example 3** (Roof Replacement). A house catches fire. Two things caused the house to catch fire: the fact that the roof was thatched and the fact that it hadn't rained recently. Bob asks, "why did the house catch fire?". Suppose that Bob knows the house had a thatched roof and that it hadn't rained recently (Bob knows the state of the world), but doesn't know enough about how fires start to know if either of these factors are the sorts of factors that cause fires (Bob is uncertain about the world's causal structure).

We could model this situation with three binary endogenous variables R, D, F, and two binary exogenous variables  $U_R, U_D$ . F represents whether the house catches fire; R whether its roof is thatched; and D whether there was a recent drought. Bob knows R, D, F = 1, but is uncertain about the causal structure; he believes it is possible that R = 1 but not D = 1 causes F = 1 (structure  $\mathcal{M}_R$ ), that D = 1 but not R = 1 causes F = 1 (structure  $\mathcal{M}_D$ ), that F = 1 iff R = 1 and D = 1 (a conjunctive structure) or that F = 1 iff R = 1 or D = 1 (a disjunctive structure). The models appear below, in Figure 3, with exogenous variables omitted for readability. Note that R = 1 and D = 1 in all four worlds (i.e. the context is  $\mathbf{u}_{D,R}$ ).

Suppose the actual causal structure is  $\mathcal{M}_R$  (i.e. whether or not the house caught fire depends only on whether the roof is thatched, and not on the drought). In this case, it seems like R=1 is an explanation of F=1; indeed, it seems like the only explanation! But note that EX2 rules this out, since there is a world,  $\mathcal{M}_D$ , in which R=1 is not a cause of F=1. On the HP analysis, the only explanation involves citing both D=1 and R=1; following their response from above, R=1

 $<sup>^{10}</sup>$ Note that we assume that Bob knows that thatched roofs and drought are the kinds of things that wouldn't prevent fires, regardless of whether or not they cause them; that is, we assume that Bob knows that the value of F is an increasing function of the values of R, D.

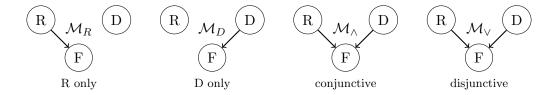


Figure 3: Bob's epistemic state  $\mathcal{K}$  in Example 3.

is only a partial explanation of the fire's starting. But this is clearly a worse explanation than R=1 alone in world  $\mathcal{M}_D$ . This highlights the counterintuitiveness of EX2; it doesn't allow for cases in which the listener has uncertainty not only about the context, but also about the structure itself.<sup>11</sup> As such, it cannot accommodate a key function of explanation identified in Section 2.1, namely that listeners can make inferences about the underlying structure (and the values of other variables) from speakers' choices of which variable values to cite in answering a "why?" question. In Example 1, the point is that Bob can infer that B=1 from the information that A=1, given that A=1 is provided as an explanation of F=1. EX2 rules out this sort of inference.

#### 2.3.3 Issues with EX3: longer explanations are often better

For similar reasons, it's unclear why the HP analysis has a hard requirement that explanations be minimal (EX3). If anything, a good explanation ought to bear *more* rather than less inferential fruit, a point often stressed in the literature on explanatory generality (e.g., Potochnik 2017; Putnam 1975). Consider the following example:

**Example 4** (Milk Theft). Bob returns from work to find that his milk carton is emptier than it was when he left that morning. He knows that the culprit(s) must be at least one (or both) of his roommates, Charlie or Dana. He asks Alice (a neutral, lactose-intolerant fourth roommate, let's say), "why is my milk gone?".

We could model this situation with binary variables C, D, M where C = 1 (respectively, D = 1) represents Charlie (respectively, Dana) drinking the milk and M = 1 represents the milk carton's being depleted. Bob knows that the causal structure is given by  $\mathcal{M}_{\vee}$  (that is, that at least one roommate's drinking the milk would be necessary and sufficient to deplete it), and knows that M = 1, meaning at least one of C = 1 and D = 1. So there are three possible worlds, given by  $\mathbf{u}_C$ ,  $\mathbf{u}_D$  and  $\mathbf{u}_{C,D}$ . We represent this in Figure 4.

Suppose the actual world is given by  $\mathbf{u}_{C,D}$ , and Alice knows this. Alice knows that Bob wants to speak to the culprit(s) to get them to stop stealing his milk in the future. In this case, it seems appropriate for her to explain M=1 by  $\langle C=1, D=1 \rangle$  (i.e. by citing both Charlie and Dana's stealing the milk), rather than by C=1 or D=1 alone. But note this is ruled out by EX3 on the HP analysis, since C=1 and D=1 both satisfy EX2.

<sup>&</sup>lt;sup>11</sup>In order to extend EX2 to cases in which the listener is uncertain about the structure, Halpern and Pearl need to stipulate that an explanation can include additional information about the structure (see their Definition 5.1. (Halpern and Pearl, 2005b, p.907)). But note that the additional information doesn't have to relate to the explanandum at all, casting doubt on whether it is best seen as part of an explanation.

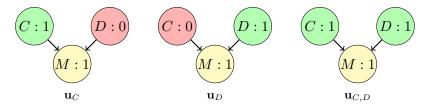


Figure 4: Bob's epistemic state in Example 4

What this example suggests is that there will be situations in which EX3 is too stringent a requirement on explanation. Indeed, it seems like pragmatic theories of explanation are well-placed to accommodate these kinds of situations. This suggests that rather than building minimality in to the definition of explanation itself, we should motivate it by pragmatic factors, such as communicative pressure (cf. Section 4.4 below).

#### 2.3.4 Issues with EX4: explanations often cite known events

Finally, EX4 is also too demanding. It has been observed in the empirical literature that an explanans is often something that the consumer already knows, in evident violation of EX4 (see, e.g., Kirfel et al. 2022 for many such examples).

Consider Example 3 above once again, and suppose the actual structure is  $\mathcal{M}_R$ . Most people will have the intuition that a good answer to Bob's question is something that he already knows, namely, "the house had a thatched roof". Not only is this a possible answer, it seems the best answer. Thus Example 3 is a straightforward counterexample to EX4.<sup>12</sup>

Could we amend EX4 to solve this issue? Looking at the example above, one could argue that the "surface form" of this answer to Bob's question – which only mentions the thatched roof – conveys new causal information, and it is this implicit causal information that is the real explanans.<sup>13</sup> To wit, Bob learns something to the effect that thatched roofs tend to cause fires to spread quickly. This suggests a weakening of EX4: rather than requiring that the listener considers it possible that  $\mathbf{X} \neq \mathbf{x}$ , we could require merely that the listener considers it possible that  $\mathbf{X} = \mathbf{x}$  is not an actual cause of the explanandum FACT. We dub this weakening EX4\*:

**(EX4)\*** There exists 
$$(\mathcal{M}, \mathbf{u}) \in \mathcal{K}$$
 in which  $\neg \mathsf{Cause}(\mathbf{X} = \mathbf{x}, \mathsf{FACT})$ .

Here,  $Cause(\mathbf{X} = \mathbf{x}, FACT)$  denotes that  $\mathbf{X} = \mathbf{x}$  is an actual cause of FACT. Note that EX4\* is weaker than EX4, because an event cannot be an actual cause if it does not even occur (technically,  $\mathbf{X} \neq \mathbf{x}$  implies  $\neg Cause(\mathbf{X} = \mathbf{x}, FACT)$ ). This means it avoids the issue with Example 3 above.

Even with this amendment, problems remain with EX4\*. Although EX4\* does seem to capture a desirable feature of explanations (that the speaker shares information unknown to the listener), it's unclear why we should think this feature is constitutive of what explanation is. Surely a speaker can explain something by citing information which the listener knows, even if this explanation will often be unhelpful to the listener? In other words, as with minimality (EX3), we should expect a preference for unknown causes to emerge from properties of the communicative scenario, rather than be built into our definition of explanation.

<sup>&</sup>lt;sup>12</sup>Note that even HP's "general definition" (Definition 5.1. (Halpern and Pearl, 2005b, p.907)) also requires that the variable values cited were not previously known, so it doesn't fix this problem.

<sup>&</sup>lt;sup>13</sup>Indeed, this is exactly how the example will be analyzed on our account; see Section 4.1.1 below.

We can make this concrete by showing that there are cases in which the best explanation available to the speaker cites causes which are known to the listener. Consider the following:

**Example 5** (Late Meeting). Bob is late to meet Alice and Charlie, and can tell that Charlie is cross at him when he arrives. He knows that his tardiness is a cause of Charlie's crossness (Charlie is famously punctual), but is unsure if it's sufficient by itself. In particular, he's wondering if Charlie's cross that he forgot his birthday the previous week. When Charlie is out of the room, Bob asks Alice, "why is Charlie cross at me?"

We model this situation with binary variables T, B, C where T represents Bob's tardiness, B represents his having forgotten Charlie's birthday, and C represents Charlie's being cross. Bob knows that B, C, T = 1, but is uncertain whether the causal structure is given by  $\mathcal{M}_T$  (his tardiness is necessary and sufficient for Charlie's crossness) or by  $\mathcal{M}_{\wedge}$  (both his tardiness and his having forgotten Charlie's birthday are individually necessary, and jointly sufficient). See Figure 5.

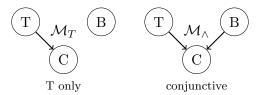


Figure 5: Bob's epistemic state  $\mathcal{K}$  in Example 5.

Suppose the actual world is given by  $\mathcal{M}_T$ . Then intuitively the best explanation is "because you were late" (T=1). Importantly, the point is not merely that this is the only explanation available to Alice. Rather, providing this explanation seems better than providing no explanation. Intuitively, it allows Bob to infer that the actual world is not  $\mathcal{M}_{\wedge}$  (because otherwise Alice would have cited his forgetting Charlie's birthday). But note that EX4\* rules out this explanation, since T=1 is an actual cause of C=1 in both  $\mathcal{M}_T$  and  $\mathcal{M}_{\wedge}$ .

## 2.3.5 Summing Up

The HP analysis of explanation can be seen as an attempt to combine (a) the (interventionist) suggestion that explanations communicate causal information with (b) the (pragmatic) suggestion that a model of explanation should be sensitive to the communicative context in which the explanation is given and received. The problem with the HP account is that it builds too much into its account of explanation; it stipulates conditions on explanation – EX2, EX3 and EX4 – which are too strong. Moreover, it seems like the benefit of incorporating contextual information into a theory of explanation is precisely to account for the plausibility of conditions like EX3 and EX4; simply stipulating these conditions seems to deprive the account of its main appeal.

We suggest that, like the formal accounts that came before, the HP definition marks important progress on analysing explanation. However, the substantive components of their theory – what goes beyond actual causation, and to the extent that these further conditions are apt in the first place – should instead emerge, perhaps in a defeasible and graded way, from a more fundamental analysis of what explanations do in conversational context.

Absent EX2, EX3 and EX4, the Halpern-Pearl account just says that the explanandum should be true and the explanans should pick out an actual cause of the explanandum. This will be our starting point for a positive theory of explanation.

## 3 A Communication-First Account

As Woodward and Ross (2021) have pointed out (see also Woodward, 2021b), most prominent philosophical accounts of explanation that stress the importance of psychological facts and interests nevertheless engage very little with work in the cognitive sciences on how people in fact produce and interpret explanations. Our positive account, by contrast, draws heavily upon recent theoretical and experimental work in the psychology of explanation, and of linguistic communication.

As in Halpern and Pearl's account, we assume that there is an agent with uncertainty over a set  $\mathcal{K}$  of pairs  $(\mathcal{M}, \mathbf{u})$ . We assume that the agent's uncertainty is represented by a probability distribution Prior on  $\mathcal{K}$ .<sup>14</sup> We will suppose further that this agent has asked (perhaps explicitly) the question "why FACT?", where FACT is something the agent knows (that is, FACT is true in all causal situations in  $\mathcal{K}$ , following EX1 from the Halpern-Pearl account). Call this agent the *listener*.

Moving beyond the setup from Halpern and Pearl (2005b), we assume that a second agent – call this agent the *speaker* – is involved.<sup>15</sup> The phenomenon of explanation is to be elucidated via the interaction between the listener and the speaker.

To model this interaction we draw upon a body of work at the intersection of linguistics and psychology intended to capture the pragmatics of ordinary conversation. This work, sometimes referred to as the *rational speech acts* (or RSA) framework, can be seen as a formalisation of familiar Gricean ideas from the philosophy of language (Frank and Goodman, 2012; Goodman and Frank, 2016b; Degen, 2023). It involves a hierarchy of conversational agent models, typically beginning with a "literal" listener, who simply interprets messages according to their semantic content, before then ascending to a pragmatic speaker, who selects a message taking into account how it will be interpreted by the literal listener. Last comes a pragmatic listener who interprets the pragmatic speaker's utterance as a "rational speech act". We follow this presentation here.

#### 3.1 The Literal Listener

Imagine a literal listener – named L0 – who has posed the question, "why FACT?", and now hears a response, "because  $\mathbf{X} = \mathbf{x}$ ". What might L0 do with this response?

Recall that L0 is uncertain about some causal facts, represented by the probability function Prior on an epistemic state  $\mathcal{K}$ . Intuitively, what L0 learns – that is, how Prior should be updated – will depend on what the message "FACT because  $\mathbf{X} = \mathbf{x}$ " says about the causal facts. Let us suppose that "FACT because  $\mathbf{X} = \mathbf{x}$ " is (literally) true of some model-context pairs  $(\mathcal{M}, \mathbf{u})$ , and false at others. In particular, we will suppose that this statement is true just in case  $\mathbf{X} = \mathbf{x}$  is an actual cause of FACT (i.e. when Cause( $\mathbf{X} = \mathbf{x}$ , FACT) is true at a world  $(\mathcal{M}, \mathbf{u})$ ). As noted above, for the purposes of this paper, virtually any account of actual cause from the literature will be suitable here. Formally, we thus have a message m with a semantic value

$$\llbracket m \rrbracket = \{ (\mathcal{M}, \mathbf{u}) \in \mathcal{K} \mid \mathcal{M}, \mathbf{u} \models m \} \},$$
 (1)

where  $\mathcal{M}, \mathbf{u} \models m$  means that m is true of the pair  $(\mathcal{M}, \mathbf{u})$ . So the message "FACT because  $\mathbf{X} = \mathbf{x}$ " has the value  $\{(\mathcal{M}, \mathbf{u}) \in \mathcal{K} \mid \mathcal{M}, \mathbf{u} \models \mathsf{Cause}(\mathbf{X} = \mathbf{x}, \mathsf{FACT})\}$ .

<sup>&</sup>lt;sup>14</sup>In general, Prior is a distribution on a  $\sigma$ -algebra over  $\mathcal{K}$ .

<sup>&</sup>lt;sup>15</sup>To disambiguate the speaker and listener, we use female pronouns for the speaker and male pronouns for the listener. In examples, we refer to the speaker as 'Alice' and the listener as 'Bob'.

Following the RSA framework, we then suppose that L0 updates his uncertainty toward a "posterior" distribution  $P_{L0}$  as:<sup>16</sup>

$$P_{L0}(\mathcal{M}, \mathbf{u} \mid m) \propto \mathsf{Prior}(\mathcal{M}, \mathbf{u}) \cdot \mathbf{1}_{\mathcal{M}, \mathbf{u} \models m}$$
 (2)

In words, message m leads L0 to assign probability 1 to m being true, and to redistribute probability mass among the remaining possibilities in proportion to his prior beliefs.

Now suppose that Alice says, "I was just in Rome!"; that is, she effectively asserts, "E=1 because A=1". What can Bob conclude? In only two situations – namely, (d) and (e) in Figure 2 – would A=1 be an actual cause of E=1 (written as Cause(A=1, E=1)). This is because (according to all extant treatments of actual causation) an actual cause needs to be true, and suitable enabling conditions (in this case, B=1, that it is sunny in Rome) also need to be true. Thus, Bob reassigns all probability mass to these two possibilities, with their relative probability unchanged from the prior.

One might expect that a listener like Bob could draw even stronger inferences from Alice's statement, e.g., that possibility (d) is significantly more likely than (e).<sup>17</sup> There is perhaps some sense in which A=1 is a "better" cause to cite when C=0 than when C=1. This type of inference, however, may depend on Bob's ability to think about why Alice said what she did (as opposed to various alternatives she could have said). At any rate, it will depend on going beyond the literal meaning that we proposed above, viz. an egalitarian notion of actual causation.

Note, however, that even at this relatively simple stage, a literal listener will be able to draw additional inferences due to the causal content of Alice's utterance. In particular, because A=1 must be an actual cause, Bob can infer that B=1 must also hold, just as highlighted above in Example 1. To go beyond what is implied by mere actual causation we will incorporate some degree of higher-order reasoning, which first involves making explicit Alice's predicament as a speaker.

## 3.2 The Pragmatic Speaker

Imagine that a speaker S – say, Alice in our running example – is aware of her addressee, L0 – say, Bob – and how Bob will respond to various possible responses she could give to his "why?" question. How should Alice choose a response?

#### 3.2.1 Useful Utterances

While much work in the RSA framework has focused on what we might call "pure information exchange", a natural thought – explored in recent work by Sumers et al. (2023) – is that a helpful speaker will convey information that helps the listener achieve his goals. Following Sumers et al., let us suppose that the listener possesses a decision problem, characterized by a pair  $(\mathcal{A}, \mathcal{R})$ , where  $\mathcal{A}$  is a set of actions and  $\mathcal{R}: \mathcal{A} \times \mathcal{K} \to \mathbb{R}$  is a reward function, specifying how good some action a is in possible causal situation  $(\mathcal{M}, \mathbf{u})$ . That is, the listener will choose an action a and receive

<sup>&</sup>lt;sup>16</sup>Here, **1** is an "indicator function," equal to 1 if the expression in the subscript is true, and 0 otherwise. Meanwhile, ' $\propto$ ' signifies "is proportional to." That is,  $P_{L0}(\mathcal{M}, \mathbf{u} \mid m) = (1/Z) \cdot \mathsf{Prior}(\mathcal{M}, \mathbf{u}) \cdot \mathbf{1}_{\mathcal{M}, \mathbf{u} \models m}$ , where Z is a "normalising constant," equal to the sum of quantities  $\mathsf{Prior}(w) \cdot \mathbf{1}_{\mathcal{M}, \mathbf{u} \models m}$ , over all possible messages m, ensuring that the probabilities all sum to (i.e., "normalise to") 1.

<sup>&</sup>lt;sup>17</sup>Such an inference is not logically implied by what is said – and so has sometimes been labeled an *illusory* disjunctive inference (Johnson-Laird and Savary, 1999) – but under reasonable assumptions about conversational dynamics it can be perfectly justifiable (Sablé-Meyer and Mascarenhas, 2022).

some scalar reward  $\mathcal{R}(a, \mathcal{M}, \mathbf{u})$ . If he knew which possibility in  $\mathcal{K}$  in fact obtained, then he would simply choose the action that returns the highest reward in that situation. But recall the listener has some uncertainty, represented by Prior.

A cooperative speaker will thus produce an utterance that leads the listener toward better actions, by updating the listener's uncertainty. Recall from the previous section that a message m will lead L0 to update Prior to a posterior distribution  $P_{L0}(-|m|)$ . We assume that the listener will then (approximately) maximize expected reward with respect to  $P_{L0}$ . Specifically, we assume that the probability with which L0 chooses action a given message m is

$$\pi_{L0}(a \mid m) \propto \exp\left(\beta_L \cdot \left[\sum_{(\mathcal{M}, \mathbf{u}) \in \mathcal{K}} P_{L0}(\mathcal{M}, \mathbf{u} \mid m) \cdot R(a, \mathcal{M}, \mathbf{u})\right]\right),$$
 (3)

where  $\exp(z) = e^z$  is the exponential function, and  $\beta_L \in \mathbb{R}^{\geq 0}$  is a "rationality parameter", measuring how close the agent is to maximizing expected utility; L0 comes ever closer as  $\beta_L$  tends to infinity. This choice follows a large body of work in psychology, economics, and computer science modelling agents as approximate expected utility maximizers (see, e.g., Luce (1963)).

If L0 is employing decision policy  $\pi_{L0}$ , then the reward that S can expect of L0 is given by

$$U_S(m, \mathcal{M}, \mathbf{u}) = \sum_{a \in \mathcal{A}} \pi_{L0}(a \mid m) \cdot R(a, \mathcal{M}, \mathbf{u}),$$
 (4)

S ought to choose her utterance m – in our setting, her answer to the listener's "why?" question – in a way that (approximately) maximizes this anticipated utility  $U_S$ :

$$P_S(m \mid \mathcal{M}, \mathbf{u}) \propto \exp\left(\beta_S \cdot U_S(m, \mathcal{M}, \mathbf{u})\right).$$
 (5)

Like in Equation 3, the parameter  $\beta_S$  here measures how close S is to optimizing (as  $\beta_S$  tends to infinity, so  $P_S$  concentrates all of its probability mass on the highest expected utility message). Most often in this paper,  $\beta_L$  and  $\beta_S$  will be set to  $\infty$ , meaning that agents are maximizing.

#### 3.2.2 Production and Processing Costs

Part of what is interesting about explanations is that they must be constructed and interpreted by resource-limited agents (Wimsatt, 2007; Ylikoski and Kuorikoski, 2010; Potochnik, 2017). One can imagine several ways of incorporating such resource limitations. For instance, S may be aware of the possibility that L0 could misinterpret a message, or indeed that S herself may err in her production of the message. In this vein we may incorporate assumptions S could make about the channel capacity for messages from S to L0 (see, e.g., Gibson et al. 2019).

An alternative, which we will adopt here, is that both production and processing of a message m come with some measurable cost. We put these two sources of cost together into a single function Cost, assigning a scalar Cost(m) to each possible message m. Thereby incorporating costs, the speaker probabilities  $P_S(m \mid \mathcal{M}, \mathbf{u})$  now become

$$P_S(m \mid \mathcal{M}, \mathbf{u}) \propto \exp\left(\beta_S \cdot \left[U_S(m \mid \mathcal{M}, \mathbf{u}) - \mathsf{Cost}(m)\right]\right).$$
 (6)

Our model thus assumes that the speaker and listener will produce and interpret each message correctly; they may just suffer cost in doing so, due to length, obscurity, complexity, etc. <sup>18</sup> This assumption could be dropped.

<sup>&</sup>lt;sup>18</sup>The cost could also include politeness (Yoon et al., 2020; Chandra et al., 2024).

## 3.3 What is the Listener's Decision Problem?

The idea that explanation can depend on interests and goals of a listener is a perennial theme across pragmatic approaches (van Fraassen, 1980; Potochnik, 2017; De Regt, 2017; Lombrozo and Liquin, 2023). As we aim for a precise framework, we will need to be somewhat concrete about what sorts of decision problems  $(\mathcal{A}, \mathcal{R})$  agents might face. At least in some cases, context will render it common knowledge that the listener has asked a "why?" question so as to inform a particular future choice. (We discuss an example like this in Section 4.1.1, returning to Example 3.)

In other cases, though, the speaker may be uncertain about which decision problems the listener might face. Thus, we could also imagine that the decision problem  $(\mathcal{A}, \mathcal{R})$  decomposes into a collection of decision problems  $(\mathcal{A}_i, \mathcal{R}_i)$ , each with a weight  $w_i$ . Where an action a is now a vector  $(\ldots, a_i \ldots)$ , specifying choices for all the decision problems, the total reward would be given by the sum  $\mathcal{R}(a, \mathcal{M}, \mathbf{u}) = \sum_i w_i \mathcal{R}_i(a_i, \mathcal{M}, \mathbf{u})$ . (We discuss an example that can be naturally modelled this way in Section 4.1.3.) Suppose, however, that the speaker cannot enumerate a specific list of possible decision problems, together with their weights. Instead, she might want to impart causal information that broadly promises to be useful. A core intuition from the interventionist tradition in the philosophy of causation and explanation is that possibilities for manipulation and control are of central importance (Woodward, 2003; Kirfel et al., 2024; D'Amico, 2025). We formalize what it might mean to have relatively broad capacity for manipulation and control with what we call a "manipulation game".

#### 3.3.1 The Manipulation Game

In the simplest case, we imagine the listener having asked, "why FACT?" suggests that causal information about FACT is somehow relevant to the decision problems he faces. As a general proxy for whatever those decision problems might be, imagine the following game. The listener is presented with some alternative way the world might have been (that is, different assignments to the exogenous variables in his knowledge state), with probability proportional to their prior likelihood. For each such possibility, the agent must choose some endogenous variable to intervene upon, with the aim of changing the truth value of FACT. In each such case, the agent wins a point just in case he successfully manipulates FACT in that situation. More formally:

**Definition 3** (Manipulation Game). A manipulation game is a decision problem  $(\mathcal{A}, \mathcal{R})$ , where:

- A is the set of all endogenous variables other than those appearing in FACT.
- $\mathcal{R}: \mathcal{A} \times \mathcal{K} \to \mathbb{R}$  is defined, for a given endogenous variable  $X \in \mathcal{A}$  and causal situation  $(\mathcal{M}, \mathbf{u}) \in \mathcal{K}$ , as

$$\mathcal{R}(X,\mathcal{M},\mathbf{u}) = \sum_{\mathbf{u}' \in \mathsf{Val}(\mathbf{U})} P(\mathbf{u}') \cdot \mathsf{Manipulates}(X,\mathsf{FACT} \mid \mathcal{M},\mathbf{u}')$$

where  $\mathsf{Manipulates}(X,\mathsf{FACT}\mid\mathcal{M},\mathbf{u}')$  is a binary-valued function that takes 1 iff there exists  $x\in\mathsf{Val}(X)$  such that either  $\mathcal{M},\mathbf{u}'\models\mathsf{FACT}\wedge[X=x]\neg\mathsf{FACT},$  or  $\mathcal{M},\mathbf{u}'\models\neg\mathsf{FACT}\wedge[X=x]\mathsf{FACT}.$  That is,  $\mathsf{Manipulates}(X,\mathsf{FACT}\mid\mathcal{M},\mathbf{u}')$  takes value 1 iff there is some intervention to X that changes the value of  $\mathsf{FACT}$  in world  $(\mathcal{M},\mathbf{u}')$ , and value 0 otherwise.

Less formally, the listener submits an endogenous variable to intervene on. For each possible context, he receives a score according to whether manipulating the value of the variable manipulates

the value of FACT, weighted by the probability of the context. Note in particular that  $\mathcal{R}(X, \mathcal{M}, \mathbf{u})$  is not sensitive to the actual context  $\mathbf{u}$ .

This reward function is closely related to several models of causal strength in the psychological literature. When all the variables are binary it coincides with (a causal version) of the so-called  $\Delta P$  measure (Jenkins and Ward, 1965; Cheng and Novick, 1992). Replacing  $P(\mathbf{u}')$  with the sampling procedure in Lucas and Kemp (2015) (which is sensitive to the actual context  $\mathbf{u}$ ) gives the so-called counterfactual effect size model of Quillien and Lucas (2023).

## 3.4 The Pragmatic Listener and the Goodness of an Explanation

Having described the literal listener, the pragmatic speaker, and the range of possible decision problems the speaker could entertain for the listener, we are now ready to present the final component, the pragmatic listener. This last agent-type will update his beliefs not directly based on the content of the utterance, but based on the fact that the speaker chose this particular utterance, knowing what she does about the (literal) listener and the decision problems confronting him.

This pragmatic listener – who we call L – updates his prior in a way that is closely analogous to L0, the literal listener:

$$P_L(\mathcal{M}, \mathbf{u} \mid m) \propto \mathsf{Prior}(\mathcal{M}, \mathbf{u}) \cdot P_S(m \mid \mathcal{M}, \mathbf{u})$$
 (7)

The only difference is in the second term, with the indicator function on the truth of m (in situation  $\mathcal{M}, \mathbf{u}$ ) replaced by the probability that S would utter m. This in turn supplies us with an updated agent policy  $\pi_L(a \mid m)$ , again identical to the expression for  $\pi_{L0}$ , except with  $P_L$  in place of  $P_{L0}$ .

While one could ascend further in this theory of mind hierarchy, the second level of pragmatic reasoning – thinking about a pragmatic speaker thinking about a literal listener – has attained privileged status in the empirical and modelling literature on linguistic pragmatics (Goodman and Frank, 2016a, see, e.g.,). We therefore take it as a reasonable conjecture about how high a typical listener might go in their interpretation of a speaker.

With this much we are now ready to offer our proposal about explanatory goodness. We submit that this can be measured as:

$$\mathsf{Goodness}(m, \mathcal{M}, \mathbf{u}) = \sum_{a \in \mathcal{A}} \pi_L(a \mid m) \cdot R(a, \mathcal{M}, \mathbf{u}) - \sum_{a \in \mathcal{A}} \pi_{\mathsf{Prior}}(a \mid m) \cdot R(a, \mathcal{M}, \mathbf{u})$$
(8)

The first term represents the listener's expected utility in the decision problem(s) he faces after hearing the explanation m and updating his prior distribution Prior to a distribution  $P_L(- \mid m)$ . The second term is a baseline; it represents the listener's expected utility if he hadn't received any explanation, and instead used his prior distribution to decide how to act. So Goodness $(m, \mathcal{M}, \mathbf{u}) > 0$  iff the explanation was helpful to the listener. When Goodness $(m, \mathcal{M}, \mathbf{u}) = 0$ , the explanation hasn't affected the listener's expected utility. When Goodness $(m, \mathcal{M}, \mathbf{u}) < 0$ , it would have been better for the speaker not to give an explanation at all.

This means that m is a good explanation to the extent that it helps the listener achieve his goals. So while actual causation plays a fundamental role, it does so only instrumentally, in the way it facilitates information transfer and, ultimately, success at a range of downstream tasks.

## 3.5 Summing Up

To a large extent, the framework we have presented in this section is a variation of the rational speech act model, as extended by Sumers et al. (2023). In particular, we have a cascade of agent-

types, with a literal listener at the bottom, followed by a pragmatic speaker thinking about her likely effect on the literal listener, and finally a pragmatic listener thinking about what prompted the speaker to say what she did. This is all grounded in the speaker's motivation to help the listener achieve some goal. Importantly, this very general framework has enjoyed success in modelling many phenomena in human conversational dynamics (Degen, 2023).

What is special about our explanatory setting is twofold. First and most fundamentally, we assume that "why?" questions are requests for causal information, and – at least when asking about singular events – for information about actual causes. This assumption is built into the semantics of "because" statements (Equation 1). Second, listeners are often interested in general issues of manipulation and control, an assumption that is formalised by means of the manipulation game (Definition 3). More generally, listeners seek a more accurate causal understanding of the world.

In sum, the framework can be seen as a natural, formal combination of reasonably general pragmatic reasoning on the one hand, and the interventionist idea that explanation is about asking and answering what-if-things-had-been-different questions on the other. These two ingredients, combined in the simple way we have done here, is enough to account for some of the most striking features of (good) explanation that have been discussed in the literature.

# 4 Explaining Explanation

How should we assess an account of explanation? A natural approach starts with the observation that explanation, unlike causation, is *inegalitarian*; it is highly selective. We have strong and systematic intuitions about which explanations are better or worse. Indeed, the literature has pinpointed a handful of 'explanatory virtues' that determine when an explanation is more or less apt (Ylikoski and Kuorikoski, 2010; De Regt, 2017), or 'lovely' (Lipton, 1991). Accordingly, good explanations

- 1. are sensitive to the downstream interests of the receiver of the explanation (De Regt 2017; Potochnik 2017; discussed in Section 4.1);
- 2. are appropriate to the listener's background knowledge (Gärdenfors 1980; Halpern and Pearl 2005b; discussed in Section 4.2);
- 3. identify explanatory relationships which are invariant across background conditions (Lewis 1986; Ylikoski and Kuorikoski 2010; Woodward 2006, 2010; discussed in Section 4.3);
- 4. involve simpler theories or more 'minimal' explanantia (Salmon 1984; discussed in Section 4.4);
- 5. sit at the 'right' level of abstraction (Yablo 1992; Strevens 2008; Woodward 2010, 2021a; discussed in Section 4.5).

This suggests a desideratum for any account of explanation, namely to make sense of why these features are – or at least appear to be – explanatory virtues. In this section, we'll show that the model in the previous section is sufficient to generate the explanatory virtues enumerated above. Our goal here is not to offer a full-fledged defence of our account, but rather to paint the view in enough detail so that its relative merits can be assessed and appreciated.

<sup>&</sup>lt;sup>19</sup>Philosophers of science have suggested that there is perhaps no privileged weighting of these (and other) good-making features; they plausibly trade off against one another such that promoting one means neglecting another (Kuhn, 1970; De Regt, 2017).

#### 4.1 Downstream Interests

In giving explanations, speakers are sensitive to listeners' downstream interests. Ceteris paribus, explanations which are more useful for the listener are better (De Regt, 2017). On our picture, this is unsurprising. Beyond communicating causal information, explanations are no different from other acts of communication: interlocutors who consider their audience's needs are simply better communicators.

#### 4.1.1 Back to Roof Replacement

It's illustrative to show how our account models this phenomenon formally, by considering a case in which it is common knowledge which decision problem the listener faces. Consider again Example 3; recall that in this example Bob knows that F=1 (the house caught fire), R=1 (the roof is thatched) and D=1 (there is a drought in the area), but is uncertain about the causal relationships between R, D and F. He considers four causal structures possible:  $\mathcal{M}_R$ ,  $\mathcal{M}_D$ ,  $\mathcal{M}_{\wedge}$  and  $\mathcal{M}_{\vee}$ . We suppose, for the sake of simplicity, that Bob has a uniform prior over these four causal structures (the same conclusion would apply with arbitrary prior distributions with full support). His knowledge state is depicted in Figure 3.

Suppose that Bob's own roof is thatched, and he's wondering whether to replace it. He lives nearby to the house that caught fire and believes, reasonably, that the same factors which caused the house to catch fire are causally responsible for whether or not his house will catch fire in the future. He would prefer to replace his roof  $(a_{\text{replace}})$  if thatched roofs cause fires, but otherwise he'd prefer not to pay the expense  $(a_{\text{don't replace}})$ . We represent this decision problem with the pay-off matrix shown in Table 1.

Table 1: Bob's decision problem in roof replacement.

	$\mathcal{M}_R$	$\mathcal{M}_D$	$\mathcal{M}_{\wedge}$	$\mathcal{M}_{ee}$
$a_{\text{replace}}$	0	0	0	0
$a_{\rm don't\ replace}$	-1	1	-1	-1

Focus on the case where the actual causal structure is  $\mathcal{M}_{\wedge}$  (that is, F=1 iff both R=1 and D=1). Suppose that Alice is deciding between two utterances, "F=1 because R=1" and "F=1 because D=1", with the following interpretations:

$$\llbracket "F = 1 \text{ because } R = 1" \rrbracket = \{ \mathcal{M}_R, \mathcal{M}_{\wedge}, \mathcal{M}_{\vee} \}$$
$$\llbracket "F = 1 \text{ because } D = 1" \rrbracket = \{ \mathcal{M}_D, \mathcal{M}_{\wedge}, \mathcal{M}_{\vee} \}.$$

Note that both of Alice's utterances allow Bob to eliminate equally likely (according to his prior) non-actual worlds from consideration ( $\mathcal{M}_R$  and  $\mathcal{M}_D$ , respectively). So  $P_{L0}(\mathcal{M}_{\wedge} \mid R = 1) = P_{L0}(\mathcal{M}_{\wedge} \mid D = 1)$  (i.e. citing either cause is equally informative to Bob as to the true causal structure). But we have  $U_S(R = 1, \mathcal{M}_{\wedge}) > U_S(D = 1, \mathcal{M}_{\wedge})$ , as is easily seen; it is more useful to Bob to learn that the roof's being thatched was a cause. This is because it allows him to make the sensible decision to replace his roof (formally, we have  $\pi_{L0}(a_{\text{replace}} \mid R = 1) > \pi_{L0}(a_{\text{replace}} \mid D = 1)$  for any value of  $\beta_L > 0$ ).

This example illustrates a general phenomenon: two utterances can be equally informative to the listener, but not equally useful to him. Eliminating  $\mathcal{M}_D$  is much more useful to Bob than

eliminating  $\mathcal{M}_R$ , because the decision problem he possibly faces is sensitive to whether or not R=1 is a cause of F=1, but not to whether or not D=1 is a cause of F=1.

# 4.1.2 Forward-Looking Decision Problems and Invariant Type-Level Causal Relationships

It's worth spelling out in more detail an assumption in the way we model Example 3. Technically speaking, Bob's decision problem is sensitive to type-level causal relationships in a closely related causal scenario: what factors could causally affect whether *his own house* will catch fire. In modelling the decision problem as sensitive to the type-level causal facts concerning the house which actually caught fire, we implicitly assume that these type-level causal relationships are unchanged between that house and Bob's house (and that this is common knowledge between Bob and Alice).

This is because Bob's decision problem is *forward-looking*; it is sensitive to what could happen in the future, under the assumption that any type-level causal relationships will stay the same as in the past. As we see it, this assumption (that many type-level causal relationships will remain unchanged in decision problems the agent might face) underlies and concretises the interventionist suggestion that causal information has practical utility. Asking "why?" questions about the past gives you information about type-level causal relationships which will be useful for the future (Woodward, 2003, Ch.1); as the example above shows, this can be true even when the explanandum itself is a singular event (rather than a regularity). Interestingly, not all decision problems are forward-looking in this way, a point that is underappreciated by interventionists; many will be *backwards-looking*, especially those that involve the attribution of responsibility (see Example 4 and Example 5).

#### 4.1.3 Explanations Identify Good Points of Intervention

It's important to emphasise that what's going on in the example above is not merely that Alice uses her reply to Bob as a way of communicating a recommendation (that he should replace his roof). Rather, as the formalism makes clear, Bob's decision whether or not to replace his roof proceeds via his update on the causal information Alice provides to him. We can make this more explicit by considering a variant of the case above. Suppose now that Alice is uncertain about the precise decision problem Bob faces. She still knows he is considering whether to replace his roof or not, but considers two situations possible.

The first situation is that described above, and modelled by the pay-off matrix in Table 1. In this case, Bob will replace his roof if it would cause a fire, but otherwise would prefer not to spend the money replacing it. Call this decision problem  $(\mathcal{A}^{(1)}, \mathcal{R}^{(1)})$ . In the second situation, by contrast, Bob wants his house to burn down (say, he is tired of upkeep effort and wants a big insurance pay-out), but would otherwise prefer to replace his roof (say, for aesthetic reasons). We can model this using a pay-off matrix similar to that in Table 1, but where the pay-offs in the two rows have been flipped. Call this decision problem  $(\mathcal{A}^{(2)}, \mathcal{R}^{(2)})$  (note that  $\mathcal{A}^{(2)} = \mathcal{A}^{(1)}$ ).

We can model the case in which Alice is uncertain between these two decision problems in the manner described in Section 3.3; let's suppose that Alice thinks the chance that Bob is facing

<sup>&</sup>lt;sup>20</sup>Which causal facts the decision problem is sensitive to can be seen by looking at the columns of the pay-off matrix. When the decision problem is insensitive to a causal fact (as with the fact that D=1 is a cause of F=1), the pay-off matrix has the same columns even as the causal fact varies (so we see that the structures  $\mathcal{M}_{\wedge}$  and  $\mathcal{M}_{R}$  have the same columns). By contrast, when the decision problem is sensitive to a causal fact (as with the fact that R=1 is a cause of F=1), we can find cases in which the columns of the pay-off matrix vary with that causal fact (so we see that the structures  $\mathcal{M}_{\wedge}$  and  $\mathcal{M}_{D}$  have different columns).

decision problem  $(\mathcal{A}^{(1)}, \mathcal{R}^{(1)})$  is  $\frac{1}{2}$ , and the chance that he's facing decision problem  $(\mathcal{A}^{(2)}, \mathcal{R}^{(2)})$  is  $\frac{1}{2}$ , such that each individual decision problem's reward function contributes equally to Bob's total reward (nothing hinges on this choice).

Note that when the actual causal structure is  $\mathcal{M}_{\wedge}$ , the best action for Bob to select is  $a_{\text{replace}} \in \mathcal{A}^{(1)}$  (as we saw above) and  $a_{\text{don't replace}} \in \mathcal{A}^{(2)}$  (i.e. he should replace his roof when he doesn't want his house to burn down, and fail to replace it when he does). Using identical reasoning to that in the previous section, we can see that – in this causal structure – it's better for Alice to cite R = 1 as a cause of F = 1 than it is for her to cite D = 1. The point is that information about the roof is more valuable to Bob than the information about the drought in both the decision problems he could face.

This discussion clarifies the interventionist suggestion that explanations often communicate effective points of intervention (Woodward, 2003; Kirfel et al., 2024). The point is not that explanations contain direct recommendations for action, but rather that (ceteris paribus) a better explanation will cite causal facts which the listener's decision problem is sensitive to. This means that what it is good for the listener to do is sensitive to the world's causal structure. The most obvious way this can happen is when the information provided by the speaker pertains to the effects of the listener's actions on the world; we saw this in Example 3, when Alice conveyed to Bob that his decision to replace his roof would have a causal effect on his house's catching fire. So information which is useful to the listener often (but of course not always) relates to the effects of causal interventions available to him.

#### 4.1.4 Pragmatic Theories of Explanation

Finally, note that varying the decision problem faced by the listener is sufficient to account for many of the examples which historically motivated pragmatic theories of explanation. Consider the following example from Hanson (1958), cited by van Fraassen (1980, p.125):

"There are as many causes of x as there are explanations of x. Consider how the cause of death might have been set out by a physician as 'multiple haemorrhage', by the barrister as 'negligence on the part of the driver', by a carriage-builder as 'a defect in the brakeblock construction', by a civic planner as 'the presence of tall shrubbery at that turning'."

Van Fraassen takes this to show that there is no common core to explanation, but rather that the appropriate relation between explanans and explanandum varies between communicative contexts. We offer a simpler diagnosis: the characters above all cite different parts of a single larger causal model in giving their explanations (i.e. the explanatory relation is one of causal dependence in all the cases). Which part of the causal model it is appropriate for the character to cite depends on the decision problem her listener faces, which the reader infers from the character's occupation. As readers, we naturally suppose that the physician is preparing an autopsy report which will be used by the police to decide whether or not to open a criminal case; the barrister is speaking to a jury who are deciding whether or not to convict the driver; the carriage builder is speaking to a team of engineers who will design the next model of the carriage; and so on.

Finding out more information about the communicative context the characters find themselves in makes this clearer; if the physician had happened to be a witness to the crash, for example, and was asked to testify as to whether the driver was distracted, it would be unhelpful for her to start talking about the victim's internal injuries, whatever her occupation. In other words, the information about the characters' occupation can be superseded by information about their audience in assessing the goodness of their explanation. What's really at stake is whether the information the character provides is *useful* for those receiving the explanation.

## 4.2 Background Knowledge

Explanations are better if they are appropriate to the listener's background knowledge. In particular, better explanations tend to cite information which is not already known to the listener. As we've seen (Section 2), existing formal accounts of explanation simply stipulate that explanations involve the provision of unknown information (Gärdenfors, 1980; Halpern and Pearl, 2005a). By contrast, our model does not require that speakers only cite unknown information. Instead, this fact emerges from our model of the dynamics of communication, as a general consequence of the fact that speakers aim at usefulness.

To see this, consider another variant of Example 3. As in Section 4.1.3, let's suppose that Alice thinks it possible that Bob is facing one of two decision problems. The first is the same as that in Section 4.1.1, specified by the pay-off matrix in Table 1. Denote it  $(\mathcal{A}^{(1)}, \mathcal{R}^{(1)})$ . For the second, suppose that Bob is considering whether or not to move to an area which is unaffected by drought. He'd rather not move, but would prefer to do so if droughts cause fires. We represent this with the pay-off matrix below. Denote this second problem  $(\mathcal{A}^{(2)}, \mathcal{R}^{(2)})$ .

	$\mathcal{M}_R$	$\mathcal{M}_D$	$\mathcal{M}_{\wedge}$	$\mathcal{M}_{ee}$
$a_{\text{move}}$	0	0	0	0
$a_{\mathrm{stay}}$	1	-1	-1	-1

Suppose, as in the examples above, that the actual causal structure is  $\mathcal{M}_{\wedge}$ . Note that the first decision problem  $(\mathcal{A}^{(1)}, \mathcal{R}^{(1)})$  is sensitive only to whether R=1 is a cause of F=1, whereas the second decision problem  $(\mathcal{A}^{(2)}, \mathcal{R}^{(2)})$  is sensitive only to whether D=1 is a cause of F=1. So unlike the case in Section 4.1.3, in which knowledge about whether R=1 was a cause of F=1 was more useful to Bob in both decision problems, here which cause would be better for Alice to cite depends on which decision problem Bob actually faces. In particular, when Alice thinks it equally likely that Bob faces the first decision problem as it is that he faces the second, then  $U_S(R=1,\mathcal{M}_{\wedge})=U_S(D=1,\mathcal{M}_{\wedge}).^{21}$ 

Suppose, though, that Bob learns that drought is a cause of the fire (i.e. that D=1 is a cause of F=1). We model this by supposing Prior is instead uniform over  $\{\mathcal{M}_D, \mathcal{M}_{\wedge}, \mathcal{M}_{\vee}\}$ ). In this case, we will have  $\mathsf{Goodness}(R=1,\mathcal{M}_{\wedge}) > \mathsf{Goodness}(D=1,\mathcal{M}_{\wedge})$ . So our model accounts for

$$\pi_{L0}(a_{\text{replace}}, - \mid R = 1) = \pi_{L0}(-, a_{\text{move}} \mid D = 1)$$

and

$$\pi_{L0}(a_{\text{replace}}, - \mid D = 1) = \pi_{L0}(-, a_{\text{move}} \mid R = 1)$$

In other words, citing R=1 allows Bob to perform better on  $(\mathcal{A}^{(1)}, \mathcal{R}^{(1)})$ , but does not help his performance on  $(\mathcal{A}^{(2)}, \mathcal{R}^{(2)})$ . So when the two decision problems are weighted equally, we have  $U_S(R=1, \mathcal{M}_{\wedge}) = U_S(D=1, \mathcal{M}_{\wedge})$ .

22 To see this, note that

$$\pi_{L0}(a_{\text{replace}}, - \mid R = 1) > \pi_{L0}(a_{\text{replace}}, - \mid D = 1)$$

<sup>&</sup>lt;sup>21</sup>Specifically, it's the case that

the intuition that it is better to cite causes which are unknown to the listener, without simply stipulating this (contra, e.g., the HP analysis in Definition 2). In fact, as we saw in Example 5, there are cases in which a good explanation can cite an event which is known to the listener to be a cause of the explanandum. It is a virtue of our proposal that it allows for these cases to occur.

#### 4.2.1 Known causes can be informative

but

We modelled Bob's knowledge state in Example 5 with two causal situations,  $\mathcal{M}_T$  and  $\mathcal{M}_{\wedge}$  (as represented in Figure 5). Let's suppose that  $\mathsf{Prior}(\mathcal{M}_T) = \mathsf{Prior}(\mathcal{M}_{\wedge})$ , for the sake of simplicity. Let's suppose that Alice doesn't have long before Charlie gets back in the room; she has two utterances available to her, "C = 1 because T = 1" and "C = 1 because T = 1", with interpretations as follows:

$$\llbracket "C = 1 \text{ because } T = 1" \rrbracket = \{ \mathcal{M}_T, \mathcal{M}_{\wedge} \}$$
$$\llbracket "C = 1 \text{ because } B = 1" \rrbracket = \{ \mathcal{M}_{\wedge} \}.$$

When Bob asks "why is Charlie cross at me?", what is his decision problem? It seems natural to suppose that he is asking the question to inform his apology. In particular, let's suppose that he is already planning to apologise to Charlie for being late, but is wondering whether to apologise for forgetting Charlie's birthday too.<sup>23</sup> If the actual causal structure is  $\mathcal{M}_T$  (that is, if Bob's tardiness is the sole cause of Charlie's crossness), he would rather not remind Charlie that he forgot his birthday. But if the actual causal structure is  $\mathcal{M}_{\wedge}$  (that is, Charlie is cross because Bob was late and forgot his birthday), Bob would rather apologise for both things. We can represent this in the following pay-off matrix, where  $a_{\text{tardiness}}$  represents Bob's apologising for being late alone, and  $a_{\text{both}}$  represents Bob's apologising for forgetting Charlie's birthday in addition. Suppose the

	$\mathcal{M}_T$	$\mathcal{M}_{\wedge}$
$a_{\text{tardiness}}$	1	-1
$a_{\mathrm{both}}$	-1	1

actual causal structure is given by  $\mathcal{M}_T$ . We have  $P_{L0}(\mathcal{M}_T \mid T=1) = P_{L0}(\mathcal{M}_{\wedge} \mid T=1)$ , so  $U_S(T=1 \mid \mathcal{M}_{\wedge}) = U_S(T=1 \mid \mathcal{M}_T)$ ; citing T=1 doesn't improve the literal listener's performance on the decision problem. But note that

$$0 = P_{L0}(\mathcal{M}_T \mid B = 1) < P_{L0}(\mathcal{M}_{\wedge} \mid B = 1) = 1.$$

So  $U_S(B=1 \mid \mathcal{M}_T) < U_S(T=1 \mid \mathcal{M}_T)$  and  $U_S(B=1 \mid \mathcal{M}_{\wedge}) > U_S(T=1 \mid \mathcal{M}_{\wedge})$ ; this means that the speaker will cite B=1 over T=1 iff the causal structure is  $\mathcal{M}_{\wedge}$ . But then we will have that

$$\pi_{L0}(-, a_{\text{move}} \mid D = 1) = \pi_{L0}(-, a_{\text{move}} \mid R = 1)$$

In other words, given Bob knows that D=1 is a cause of F=1, citing D=1 no longer improves his performance on  $(\mathcal{A}^{(2)}, \mathcal{R}^{(2)})$  in  $\mathcal{M}_{\wedge}$ , whereas citing R=1 continues to improve his performance on  $(\mathcal{A}^{(1)}, \mathcal{R}^{(1)})$ .

<sup>&</sup>lt;sup>23</sup>Note this is an example of a *backwards-looking* decision problem; it is sensitive to token-level causal relationships involving a past event, rather than type-level relationships.

 $P_L(\mathcal{M}_T \mid T=1) > P_L(\mathcal{M}_{\wedge} \mid T=1)$ ; the pragmatic listener is able to infer from the fact that the speaker didn't cite B=1 that the actual causal structure is  $\mathcal{M}_T$ ! Crucially, this means that

Goodness
$$(T=1,\mathcal{M}_T)>0$$
.

So it is actively helpful for Alice to cite T=1 as a cause of C=1, even though this was already known to Bob. Our model delivers the correct intuition.

## 4.3 Explanatory Relationships and Background Conditions

Better explanations tend to cite explanatory relationships which persist across different background conditions. Woodward (2003) develops this idea in terms of the *invariance* of a causal relationship; a causal relationship between two variable assignments in an SCM is invariant if it continues to hold as the values of other variables change (see also Vasilyeva et al., 2018). Woodward suggests that, all else being equal, information about a causal relationship will be more useful for manipulation and control to the degree the causal relationship is invariant.

Our model concretises Woodward's claim that invariant causal relationships are more useful for listeners, via the notion of a *manipulation game*, introduced in Section 3.3.1 (Definition 3). Recall that a manipulation game is a type of decision problem which (we argued) arises when the speaker is unsure about precisely why the listener has asked "why FACT?".

The speaker supposes that the listener is interested in manipulating whether or not FACT holds. Concretely, the listener must select an endogenous variable in his knowledge state. He is evaluated on the number of background conditions (settings to endogenous variables in his knowledge state) in which intervening on the variable he has selected would change whether or not FACT holds (i.e. he is evaluated on his ability to manipulate the explanandum, based on the information received).

In this section, we use manipulation games to show that our model neatly accounts for a body of literature in the empirical literature on causal selection (Gerstenberg and Icard, 2020; Icard et al., 2017; Kominsky et al., 2015; Henne et al., 2019; Quillien and Lucas, 2023; Kirfel et al., 2022). Specifically, we show that our model correctly predicts that

- speakers prefer to cite 'abnormal' causes in conjunctive causal structures;
- speakers prefer to cite 'normal' causes in disjunctive causal structures;
- given speakers' preferences, listeners can infer the normality of causes from their knowledge of the causal structure (and vice-versa).

This highlights another important way in which our account constitutes a bridge between the philosophical and psychological literatures on explanation.

#### 4.3.1 Causal Selection Interacts with Normality and Structure

Once again, let's adopt the set-up of Example 3, where Bob's knowledge state is as represented in Figure 3. Suppose also that Alice is deciding whether to cite R=1 or D=1 as a cause of F=1. There will be two differences between the cases we've discussed so far and the case we'll consider here. First, the examples above didn't rely on any particular specification of the distribution  $P(\mathbf{U})$  on the exogenous variables. In this example, by contrast, we'll suppose that for all the causal situations in Bob's epistemic state, we have  $0 < P(U_R) < P(U_D)$ .<sup>24</sup> So R=1 is

We assume that  $U_R, U_D$  are the exogenous variables on which R, D depend. So R = 1 iff  $U_R = 1$  and D = 1 iff  $U_D = 1$ . We assume that  $U_R, U_D$  are independent.

relatively (statistically) 'abnormal' and D=1 is relatively 'normal'. Second, suppose that Alice doesn't know the specific decision problems Bob faces. In this situation, the interventionist suggests that she will give him information useful to manipulating whether or not the house catches fire in similar situations. Should she cite D=1 or R=1?

We can model this using a manipulation game  $(\mathcal{A}, \mathcal{R})$ , as in Definition 3. There will be two variables which Bob can choose between (R and D) and four contexts in which intervening on these variables will be evaluated:  $\mathbf{u}_{\emptyset}$ ,  $\mathbf{u}_{R}$ ,  $\mathbf{u}_{D}$  and  $\mathbf{u}_{R,D}$ , where  $\mathbf{u}_{\emptyset}$  denotes the case in which  $R, D = 0.^{25}$  The pay-off matrix is given in Table 2. Here,  $a_{R}$  (respectively,  $a_{D}$ ) denotes Bob's choosing to intervene on R (respectively, D).

Table 2: Pay-offs for the simple manipulation game for Example 3.

	$\mathcal{M}_R$	$\mathcal{M}_D$	$\mathcal{M}_{\wedge}$	$\mathcal{M}_ee$
$a_{\rm R}$	1	0	$P(U_D)$	$1 - P(U_D)$
$a_{\rm D}$	0	1	$P(U_R)$	$1 - P(U_R)$

Note that  $\mathcal{R}(a_R, \mathcal{M}_R) = 1$ , since changing the value of variable R in structure  $\mathcal{M}_R$  changes the value of F across all contexts. Conversely,  $\mathcal{R}(a_D, \mathcal{M}_R) = 0$  since changing the value of variable D in structure  $\mathcal{M}_R$  has no effect on the value of F, regardless of the context. Identical reasoning applies to the column for structure  $\mathcal{M}_D$ . To calculate the entries for the column for structure  $\mathcal{M}_{\wedge}$ , note that – in this structure – changing the value of R successfully manipulates F in contexts  $\mathbf{u}_D$  and  $\mathbf{u}_{R,D}$ , and changing the value of D successfully manipulates F in contexts  $\mathbf{u}_R$  and  $\mathbf{u}_{R,D}$ . Adding together the probabilities of these contexts gives  $P(U_D)$  and  $P(U_R)$ , respectively. To calculate the entries for the column for structure  $\mathcal{M}_{\vee}$ , note that – in this structure – changing the value of R successfully manipulates F in contexts  $\mathbf{u}_{\emptyset}$  and  $\mathbf{u}_R$ , and changing the value of D successfully manipulates F in contexts  $\mathbf{u}_{\emptyset}$  and  $\mathbf{u}_D$ . Adding together the probabilities of these contexts gives  $1 - P(U_D)$  and  $1 - P(U_R)$ , respectively.

We have  $\pi_{L0}(a_R \mid R=1) > \pi_{L0}(a_D \mid R=1)$  and  $\pi_{L0}(a_D \mid D=1) > \pi_{L0}(a_R \mid D=1)$ , as is easily seen. Recall that, by assumption,  $P(U_R) < P(U_D)$ . So we have:

$$U_S(R = 1, \mathcal{M}_{\wedge}) > U_S(D = 1, \mathcal{M}_{\wedge})$$
  
 $U_S(R = 1, \mathcal{M}_{\vee}) < U_S(D = 1, \mathcal{M}_{\vee})$ 

Informally, when the actual causal structure is  $\mathcal{M}_{\wedge}$ , the speaker should cite R=1; when the actual causal structure is  $\mathcal{M}_{\vee}$ , the speaker should cite D=1.

This accords with – and perhaps goes some way toward explaining – a growing body of evidence in the cognitive science literature that speakers select abnormal causes in conjunctive structures, and normal causes in disjunctive structures (Kominsky et al., 2015; Icard et al., 2017; Henne et al., 2019; Gerstenberg and Icard, 2020; Quillien and Lucas, 2023). We take it to be compatible with the idea that in judging causal strength people simulate different counterfactual possibilities (Icard et al., 2017; Quillien and Lucas, 2023); on our picture, they are imagining speaking to a listener playing a simple manipulation game.

<sup>&</sup>lt;sup>25</sup>So we have  $P(\mathbf{u}_{\emptyset}) = (1 - P(U_R))(1 - P(U_D)), P(\mathbf{u}_R) = P(U_R)(1 - P(U_D)), P(\mathbf{u}_D) = (1 - P(U_R))P(U_D), P(\mathbf{u}_{R,D}) = P(U_R)P(U_D).$ 

#### 4.3.2 Inferences from Normality

Aside from production, recent work also shows that a listener can infer a great deal from a speaker's choice of which cause to cite (Kirfel et al., 2022; see also Davis et al., 2025). For instance, people can infer whether a structure is disjunctive or conjunctive (when they know which cause is normal) or infer which cause is normal (when they know the causal structure). Again, our model captures this, via the 'pragmatic listener' L. For example, suppose Alice cites the abnormal cause (the roof's being thatched, R = 1). Then, since  $P_S(R = 1, \mathcal{M}_{\wedge}) > P_S(R = 1, \mathcal{M}_{\vee})$ , we have

$$P_L(\mathcal{M}_{\wedge} \mid R=1) > P_L(\mathcal{M}_{\vee} \mid R=1).$$

In other words, the pragmatic listener is able to infer that the structure is conjunctive rather than disjunctive, since he takes into account the fact that the speaker will prefer to cite abnormal causes when the structure is conjunctive (and normal causes when it is disjunctive).

A variant of the above example can account for listeners' ability to infer normality of the cited cause from knowledge of the causal structure. Suppose Bob is uncertain about whether R=1 or D=1 is normal, but knows the causal structure is not disjunctive. That is, Prior is uniform over

$$\mathcal{K} = \left\{ \mathcal{M}^{P(U_D) > P(U_R)}, \mathcal{M}^{P(U_R) > P(U_D)} \mid \mathcal{M} \in \{\mathcal{M}_R, \mathcal{M}_D, \mathcal{M}_\wedge\} \right\}$$

where (e.g.)  $\mathcal{M}_{\wedge}^{P(U_D)>P(U_R)}$  is the same as  $\mathcal{M}_{\wedge}$ , but with  $P(U_D)>P(U_R)$  (i.e. with D=1 as the normal cause). Suppose the actual causal situation is given by  $\mathcal{M}_{\wedge}^{P(U_D)>P(U_R)}$ . Then since we have

$$P_S(R = 1, \mathcal{M}_{\wedge}^{P(U_D) > P(U_R)}) > P_S(R = 1, \mathcal{M}_{\wedge}^{P(U_R) > P(U_D)}),$$

we will have

$$P_L(\mathcal{M}_{\wedge}^{P(U_D)>P(U_R)} \mid R=1) > P_L(\mathcal{M}_{\wedge}^{P(U_R)>P(U_D)} \mid R=1).$$

So our model recovers the prediction that listeners can infer the normality of causes from their knowledge of the structure.<sup>26</sup>

## 4.4 Minimality and Simplicity

Good explanations are often minimal (have no redundant information) and simple (easily comprehended by the listener). We model this using a standard modification to RSA, adding a term to the speaker utility to reflect the *cost* of selecting the message.<sup>27</sup>

As we saw in Section 2.3.3, there are situations in which longer explanations are intuitively better than shorter explanations. Indeed, psychologists of explanation emphasise that people often prefer explanations that cite more causal mechanisms rather than fewer (Zemla et al., 2017, 2023; Vrantsidis and Lombrozo, 2024). Our model reflects this. Consider Example 4, with Bob's knowledge state as in Figure 4. Recall that Bob has asked why his milk has been depleted; he knows

<sup>&</sup>lt;sup>26</sup>It is worth emphasizing again that patterns like these appear to violate HP's fourth postulate, EX4. While it is often appropriate to cite unknown factors, this cannot be a hard requirement on explanation.

<sup>&</sup>lt;sup>27</sup>Philosophers of science often cash out minimality in terms of *length* in some description language (Kitcher, 1981). But when we move to a more cognitive perspective, it's clear that length is at best an imperfect heuristic for the relevant sense of minimality. Short messages may be obscure, whilst longer messages may be cumbersome. See, e.g., Ylikoski and Kuorikoski (2010); Lage et al. (2019).

that either Charlie or Dana (or both) were culprits (contexts  $\mathbf{u}_C$ ,  $\mathbf{u}_D$  and  $\mathbf{u}_{C,D}$ , respectively). Let's suppose that Bob has asked the question so as to know which roommate(s) to confront (and Alice knows this). He wants to confront a roommate iff they have taken his milk. We can model the decision problem in the following table (as always, the precise pay-offs are arbitrary, since we are interested in demonstrating a qualitative effect): Let's suppose Alice has three utterances available

	$\mathbf{u}_C$	$\mathbf{u}_D$	$\mathbf{u}_{C,D}$
$a_{\text{Charlie}}$	1	-1	0
$a_{\mathrm{Dana}}$	-1	1	0
$a_{\mathrm{both}}$	0	0	1

to her, "M=1 because C=1" ("Charlie took the milk"), "M=1 because C=1" ("Dana took the milk") and "M=1 because C=1, D=1" ("Charlie and Dana each took the milk"), with the following interpretations:

It is clear that when the costs associated with each message are equal, we will have  $Goodness(C = 1, D = 1, \mathbf{u}_{C,D}) > Goodness(C = 1, \mathbf{u}_{C,D})$ ,  $Goodness(D = 1, \mathbf{u}_{C,D})$ . When both roommates took milk, it is better for Bob to know this. But suppose, plausibly, that

$$Cost(C = 1, D = 1) > Cost(C = 1) = Cost(D = 1),$$

(because, e.g., the longer message is more onerous to produce). If the difference in costs is sufficiently high, such that

$$U_S(C = 1, D = 1, \mathbf{u}_{C,D}) - \mathsf{Cost}(C = 1, D = 1) < U_S(C = 1, \mathbf{u}_{C,D}) - \mathsf{Cost}(C = 1)$$
  
=  $U_S(C = 1, \mathbf{u}_{C,D}) - \mathsf{Cost}(D = 1)$ 

the speaker will prefer the shorter messages (she will be indifferent between them), even though they are less useful.

The discussion above gives us a way of quantifying the amount of redundancy in a candidate explanation m, by considering the term

$$U_S(m, \mathcal{M}, \mathbf{u}) - \max_{m' \neq m} U_S(m, \mathcal{M}, \mathbf{u})$$

If this term is positive, then – assuming all messages have the same cost – the speaker will select m over the other messages (as  $\beta_S \to \infty$ , the message will be selected with certainty). The larger the positive value is, the less redundancy the message contains, in the sense that we could increase its production cost relative to the other messages and it would still be selected by the speaker.

To give a concrete example of a message with high redundancy, suppose that  $U_S$  is specified using the decision problem  $(\mathcal{A}, \mathcal{R})$  in Example 3, with the pay-off matrix in Table 1. Suppose Alice has the utterance "F = 1 because both R = 1 and D = 1" available to her, with interpretation  $\{\mathcal{M}_{\wedge}\}$ . In this situation, we will have

$$U_S(R=1, D=1, \mathcal{M}_{\wedge}) = U_S(R=1, \mathcal{M}_{\wedge}),$$

since the only information relevant to Bob's decision problem is whether R = 1 is a cause of F = 1. In this case, then, if

$$Cost(R = 1, D = 1) > Cost(R = 1),$$

Alice will always prefer the shorter message "R = 1", since it is just as useful to Bob as the more informative message.<sup>28</sup>

Note that these considerations also interact with the *normality* of the causes being cited. In Example 4 above, if (say)  $P(U_D = 1) >> 0.5 >> P(U_C = 1)$  (Dana usually takes Bob's food and Charlie rarely does, and this is common knowledge between Alice and Bob), such that  $Prior(\mathbf{u}_C) > Prior(\mathbf{u}_{C,D}) > Prior(\mathbf{u}_{C,D})$ , then we will have a relatively small difference between  $U_S(C = 1, D = 1, \mathbf{u}_{C,D})$  and  $U_S(C = 1, \mathbf{u}_{C,D})$ . To see this, note that Bob already assumes Dana has taken the milk, and so – upon learning that Charlie has taken the milk (i.e. that the actual context is *not*  $\mathbf{u}_D$ ) – believes it is more likely that the actual context is  $\mathbf{u}_{C,D}$  than  $\mathbf{u}_C$ . This means that there is more redundancy in the message "M = 1 because C = 1, D = 1".

#### 4.4.1 Simplicity

Several authors have observed that choice of wording matters in assessing the quality of an explanation. Consider the following, adapted from Kim (1999) and discussed by (Woodward, 2003, p. 217):

**Example 6** (Event of the Year). Assume that the short circuit caused the fire, and that the short circuit was the most noteworthy event of the year. Bob asks Alice, "why did the fire happen?" Intuitively, Alice's utterance "the short circuit caused the fire" is a better explanation of the fire's starting than the utterance "the most noteworthy event of the year caused the fire".

Woodward suggests that this should be accounted for in terms of invariance; the relationship between "the short circuit" and "the fire" is invariant across nearby interpretations of "the short circuit", in a way that it is not between "the most noteworthy event of the year" and "the fire" (since "the most noteworthy event of the year" could easily have referred to a different event that has nothing to do with the fire). But it's clear that this response isn't satisfactory; invariance describes a causal relationship between two features of the world (e.g. events, propositions, etc.), rather than between the phrases used to refer to those features. Moreover, we could easily imagine a situation in which Bob knows what a short circuit is and knows the recent short circuit is widely agreed to be "the most noteworthy event of the year", but doesn't know that "short circuit" refers to a short circuit. In this case, "the most noteworthy event of the year caused the fire" would intuitively be a better explanation for Bob than "the short circuit caused the fire".

Thinking in terms of pragmatics provides an easier solution. The speaker is deciding between two utterances  $m, m' \in \mathfrak{M}$  which refer to the same feature of the world; that is, we have  $\llbracket m \rrbracket =$ 

<sup>&</sup>lt;sup>28</sup>In Example 1, we have another case in which a longer message is redundant. When Alice cites being in Rome as a cause of her tan, she doesn't need to specify that it was also sunny; Bob is able to infer this. (Of course, if we supposed there was some cost attached to this inference, and incorporated this into our model, then the more exhaustive explanation would be preferred. See, e.g., Ylikoski and Kuorikoski (2010, p.214) "since the power of an explanation is always dependent on the range of counterfactual inferences that the explanatory information enables, the kinds of inferences possible for limited cognitive systems such as humans directly affect what can be explained and understood by such cognitive systems".)

[m']. <sup>29</sup> However, the listener's ability to interpret m might be different from his ability to interpret m'; so m and m' – even though they are identical when interpreted – incur different processing costs for the listener. In particular, it requires more effort for the listener to recognise that "the most noteworthy event of the year" refers to the short circuit than it does "the short circuit"; there is also more chance that he will misunderstand which part of the world the speaker meant to refer to. So although m and m' have the same truth conditions,  $Cost(m) \neq Cost(m')$ . <sup>30</sup>

## 4.5 Proportionality and Levels of Explanation

Consider the following example, adapted from Yablo (1992).<sup>31</sup>

**Example 7.** Alice and Bob see a pigeon peck at a scarlet target. Bob asks "why did the pigeon peck at the target?". As it happens, the penguin has been trained to peck at a target iff it is some shade of red. Alice is deciding between the following two explanations of the pigeon's pecking.

- 1. "The pigeon pecked at the target because it is scarlet."
- 2. "The pigeon pecked at the target because it is red."

Intuitively, the second explanation is better than the first, even though both involve true causal claims. This intuition is referred to as 'proportionality'; it is closely related to the idea that better explanations sit at a more appropriate level of abstraction.

We can model this with two variables, C and P. P is a binary variable, representing whether or not the pigeon pecks at the target. C is a variable representing the colour of the target. Let's suppose that it has five possible values: "empty" (0), "scarlet" (s), "crimson" (c), "blue" (b). Bob knows that C = s (the target is scarlet). Suppose that Bob knows that P depends entirely on the value of C; he knows that the pigeon won't peck at an empty target (i.e. that P = 0 when C = 0), but doesn't know which values of C (apart from s) are sufficient for P = 1. For each  $V \subseteq \{c, b\}$ , let  $\mathcal{M}_{\{s\} \cup V}$  be the causal model specified by

$$f_P(v) = 1$$
 iff  $v \in \{s\} \cup V$ .

For example,  $\mathcal{M}_{\{c,s\}}$  is the causal structure corresponding to the case in which the pigeon pecks at the target iff it is crimson *or* scarlet. We can then define Bob's knowledge state

$$\mathcal{K} = \{ \mathcal{M}_{\{s\} \cup V} \}_{V \subseteq \{c,b\}}.$$

So Bob considers four possible causal situations possible:  $\mathcal{M}_{\{s\}}$ ,  $\mathcal{M}_{\{c,s\}}$ ,  $\mathcal{M}_{\{b,s\}}$  and  $\mathcal{M}_{\{b,c,s\}}$ .

In the example above, the actual causal situation is given by  $\mathcal{M}_{\{c,s\}}$ ; the pigeon pecks at any red target, regardless of whether it is crimson or scarlet. Suppose that Bob is facing a decision problem which is sensitive to the exact causal situation; say, he wants to demonstrate the pigeon's pecking abilities, and is wondering what colour targets to purchase (should he buy scarlet only, scarlet and crimson, scarlet and blue, or targets of any of the colours?). We can model this using the following pay-off matrix:

<sup>&</sup>lt;sup>29</sup>Formally, for all  $(\mathcal{M}, \mathbf{u}) \in \mathcal{K}$ :  $\mathcal{M}, \mathbf{u} \models m$  iff  $\mathcal{M}, \mathbf{u} \models m'$ .

<sup>&</sup>lt;sup>30</sup>As discussed in Section 3.2.2, our core model is compatible with different ways of formalising a 'processing cost' to the listener. For example, one natural idea is that more complex messages have a higher probability of being misinterpreted by the listener; "the most noteworthy event of the year" example could be treated in this way.

<sup>&</sup>lt;sup>31</sup>As Woodward (2021a) observes, some take Yablo to be making a point about the metaphysics of causation here. We follow Woodward in interpreting this as a desideratum on causal explanation.

	$\mathcal{M}_{\{s\}}$	$\mathcal{M}_{\{c,s\}}$	$\mathcal{M}_{\{b,s\}}$	$\mathcal{M}_{\{b,c,s\}}$
$a_{\{s\}}$	1	0	0	0
$a_{\{c,s\}}$	0	1	0	0
$a_{\{b,s\}}$	0	0	1	0
$a_{\{b,c,s\}}$	0	0	0	1

where, for example,  $a_{\{c,s\}}$  represents the case in which Bob buys both crimson and scarlet targets to use.

Suppose Alice has four utterances available, with the following interpretations:

$$\llbracket "P = 1 \text{ because the target is scarlet"} \rrbracket = \left\{ \mathcal{M}_{\{s\}}, \mathcal{M}_{\{c,s\}}, \mathcal{M}_{\{b,c,s\}}, \mathcal{M}_{\{b,c,s\}} \right\}$$

$$\llbracket "P = 1 \text{ because the target is red} \rrbracket = \left\{ \mathcal{M}_{\{c,s\}}, \mathcal{M}_{\{b,c,s\}} \right\}$$

$$\llbracket "P = 1 \text{ because the target is (scarlet or blue)} \rrbracket = \left\{ \mathcal{M}_{\{b,s\}}, \mathcal{M}_{\{b,c,s\}} \right\}$$

$$\llbracket "P = 1 \text{ because the target is coloured} \rrbracket = \left\{ \mathcal{M}_{\{b,c,s\}} \right\}$$

Suppose, as in the example, the actual causal structure is  $\mathcal{M}_{\{c,s\}}$ . Then clearly it is better for the speaker to cite that the target is red than that the target is crimson. In particular, the pragmatic listener L will correctly infer that the actual structure is not  $\mathcal{M}_{\{b,c,s\}}$  from the fact that the speaker didn't cite the target's being coloured as the cause. Similarly, if the speaker cites the target's being crimson, the pragmatic listener will incorrectly infer that the structure is  $\mathcal{M}_{\{s\}}$ , since otherwise the speaker would have cited some other cause.<sup>32</sup>

Interestingly, the model predicts interactions between these kinds of inferences and the cost of messages. Suppose that

$$Cost("P = 1 because the target is (scarlet or blue)")$$

is very high.<sup>33</sup>. Then when the pragmatic listener hears the speaker cite the target's being scarlet as the cause, he can no longer rule out the causal structure  $\mathcal{M}_{\{b,s\}}$  (since the cost constraint means that the speaker wouldn't have cited the disjunction (scarlet or blue) even if this were the structure). This seems to us like the right intuition.

## 5 Discussion and Future Work

We've introduced a formal pragmatics of explanation, and have shown that it can account for several explanatory virtues, as well as empirical patterns of causal selection. We have demonstrated the value of putting communication first in an account of explanation. On our picture, candidate

 $<sup>^{32}</sup>$ Readers familiar with RSA will note the similarity between this analysis and the way that RSA accounts for (e.g.) scalar implicature.

<sup>&</sup>lt;sup>33</sup>This seems plausible; after all, there is no single word for (scarlet or blue).

explanations are evaluated first and foremost for their usefulness qua acts of communication (an "illocutionary evaluation", as Achinstein (1983) puts it). To model an explanation, one must represent both the listener's knowledge and the set of decision problems he faces. An important facet of the philosophical problem of explanation, then, consists in accounting for the kinds of things explanations can do for us. Put another way: what are the situations in which it is valuable for someone to seek out causal information, by (e.g.) asking a "why?" question?

For pedagogical purposes, the examples we've considered in this paper are toy, contrived in various respects. At this stage, it's natural to wonder: how much of what's been said about explanation might our model be able to account for, with suitable augmentations? There are several avenues which future work might address; we give three below.

First, much philosophical interest in explanation comes from the central role it plays in scientific inquiry. Applying the model to explanation in the sciences would involve, for example, treating a scientific community as a 'listener', and attempting to represent its knowledge and downstream aims more explicitly. This strikes us as a fruitful research area for philosophers of science engaging closely with actual scientific practice.<sup>34</sup>

Second, our model makes the simplifying assumption that all the causal situations in the listener's knowledge state  $\mathcal{K}$  share the same finite set of variables; in particular, the listener is not uncertain about which variables there are, merely about how they relate. But it seems clear that explanations are often most useful precisely when they introduce variables hitherto unknown to the listener; a future version of the model could attempt to model variable introduction explicitly (a natural first attempt would involve representing the knowledge state hierarchically).

Third, although we tether the account we develop here to causation, it's clear that the model's formalism can accommodate any dependence describable by an SCM.<sup>35</sup> Future work could apply the model to non-causal domains, such as mathematical explanation.

# 6 Acknowledgments

This work was supported by a grant from the Human-centered Artificial Intelligence Institute (HAI) at Stanford University. TG was supported by a grant from Cooperative AI.

<sup>&</sup>lt;sup>34</sup>Similarly, many of the explananda we discuss here are single event explananda, but it's clear that the model could apply to regularity explananda (just as interventionists use SCMs to describe these explananda).

<sup>&</sup>lt;sup>35</sup>This point is often made about interventionist accounts of explanation (Woodward, 2018).

## References

- Peter Achinstein. What Is an Explanation? American Philosophical Quarterly, 14(1):1–15, 1977. ISSN 0003-0481.
- Peter Achinstein. The Nature of Explanation. Oxford University Press, 1983.
- Peter Achinstein. The Pragmatic Character of Explanation. *PSA: Proceedings of the Biennial Meeting of the Philosophy of Science Association*, 1984:275–292, 1984. ISSN 0270-8647.
- Elias Bareinboim, Juan Correa, Duligur Ibeling, and Thomas Icard. On Pearl's hierarchy and the foundations of causal inference. In Hector Geffner, Rina Dechter, and Joseph Y. Halpern, editors, *Probabilistic and Causal Inference: The Works of Judea Pearl*, pages 509–556. ACM Books, 2022.
- Jon Bebb and Helen Beebee. Causal Selection and Egalitarianism. In Shaun Nichols and Joshua Knobe, editors, Oxford Studies in Experimental Philosophy, Volume 5. Oxford University Press, 2024.
- Aaron Beller and Tobias Gerstenberg. Causation, meaning, and communication. *PsyArXiv*, 2024. URL https://psyarxiv.com/xv8hf.
- Sylvain Bromberger. An Approach to Explanation. In R. J. Butler, editor, *Analytic Philosophy*, 2nd Edition. Blackwell, 1965.
- Sylvain Bromberger. On Pragmatic and Scientific Explanation: Comments on Achinstein's and Salmon's Papers. *PSA: Proceedings of the Biennial Meeting of the Philosophy of Science Association*, 1984:306–325, 1984. ISSN 0270-8647.
- Kartik Chandra, Tony Chen, Tzu-Mao Li, Jonathan Ragan-Kelley, and Joshua T. Tenenbaum. Cooperative explanation as rational communication. In L. K. Samuelson, S. L. Frank, M. Toneva, A. Mackey, and E. Hazeltine, editors, *Proceedings of the 46th Annual Conference of the Cognitive Science Society*, 2024.
- Patricia W. Cheng and L. R. Novick. Covariation in natural causal induction. *Psychological Review*, 99:365–382, 1992.
- Zachary J Davis, Kelsey R Allen, Max Kleiman-Weiner, Julian Jara-Ettinger, and Tobias Gerstenberg. Inference from social evaluation. *Journal of Personality and Social Psychology*, 2025.
- Henk. W De Regt. Understanding Scientific Understanding. Oxford University Press, 2017.
- Henk W. De Regt and Dennis Dieks. A Contextual Approach to Scientific Understanding. *Synthese*, 144(1):137–170, 2005. ISSN 0039-7857.
- Judith Degen. The rational speech act framework. Annual Review of Linguistics, 9:519–540, 2023.
- Igor Douven. The Art of Abduction. MIT Press, 2022.
- Phil Dowe. Wesley salmon's process theory of causality and the conserved quantity theory. *Philosophy of Science*, 59(2):195–216, 1992.

- Fred I. Dretske. Contrastive statements. The Philosophical Review, 81(4):411-437, 1972.
- Marina D'Amico. An interventionist approach to causal selection: The optimal control hypothesis. Journal of Philosophy, 2025. Forthcoming.
- Jan Faye. The pragmatic-rhetorical theory of explanation. In Johannes Persson and Petri Ylikoski, editors, *Rethinking Explanation. Series: Boston Studies in the Philosophy of Science Vol. 252.*, pages 43–68. Springer Verlag, 2007.
- Michael C Frank and Noah D Goodman. Predicting pragmatic reasoning in language games. Science, 336(6084):998–998, 2012.
- L. R. Franklin-Hall. High-Level Explanation and the Interventionist's 'Variables Problem'. *The British Journal for the Philosophy of Science*, 67(2):553–577, June 2016. ISSN 0007-0882. doi: 10.1093/bjps/axu040.
- Dmitri Gallow. A model-invariant theory of causation. The Philosophical Review, 130(1):45–96, 2021.
- Peter Gärdenfors. A Pragmatic Approach to Explanations. *Philosophy of Science*, 47(3):404–423, 1980. ISSN 0031-8248.
- Peter Gärdenfors. Knowledge in Flux: Modeling the Dynamics of Epistemic States. MIT Press, 1988.
- Peter Gärdenfors. An epistemic analysis of explanations and causal beliefs. *Topoi*, 9(2):109–124, September 1990. ISSN 1572-8749. doi: 10.1007/BF00135892.
- Tobias Gerstenberg and Thomas F. Icard. Expectations affect physical causation judgments. *Journal of Experimental Psychology: General*, 149(3):599–607, 2020.
- Edward Gibson, Richard Futrell, Steven P Piantadosi, Isabelle Dautriche, Kyle Mahowald, Leon Bergen, and Roger Levy. How efficiency shapes human language. *Trends in Cognitive Sciences*, 23(5):389–407, 2019.
- Mariel K Goddu and Alison Gopnik. The development of human causal learning and reasoning. Nature Reviews Psychology, pages 1–21, 2024.
- Noah D. Goodman and Michael C. Frank. Pragmatic language interpretation as probabilistic inference. *Trends in Cognitive Science*, 20:818–829, 2016a.
- Noah D. Goodman and Michael C. Frank. Pragmatic language interpretation as probabilistic inference. *Trends in Cognitive Sciences*, 20(11):818–829, 2016b.
- Joseph Y Halpern. Actual causality. MIT Press, 2016.
- Joseph Y. Halpern and Judea Pearl. Causes and Explanations: A Structural-Model Approach. Part II: Explanations. The British Journal for the Philosophy of Science, 56(4):889–911, 2005a.
- Joseph Y. Halpern and Judea Pearl. Causes and Explanations: A Structural-Model Approach. Part I: Causes. *The British Journal for the Philosophy of Science*, 56(4):843–887, 2005b.

- Norwood Russell Hanson. Patterns of Discovery, volume 11. University Press, 1958.
- Carl G. Hempel and Paul Oppenheim. Studies in the Logic of Explanation. *Philosophy of Science*, 15(2):135–175, 1948. ISSN 0031-8248.
- Carl Gustav Hempel. Aspects of Scientific Explanation and Other Essays in the Philosophy of Science. Number 1. The Free Press, 1965.
- Paul Henne, Laura Niemi, A. Pinillos, Felipe De Brigard, and Joshua Knobe. A counterfactual explanation for the action effect in causal judgment. *Cognition*, 190:157–164, 2019.
- Denis J. Hilton. Conversational processes and causal explanation. *Psychological Bulletin*, 107(1): 65–81, 1990. ISSN 1939-1455(Electronic),0033-2909(Print). doi: 10.1037/0033-2909.107.1.65.
- Christopher Hitchcock. The intransitivity of causation revealed in equations and graphs. *The Journal of Philosophy*, 98(6):273–299, 2001.
- Bas C. van Fraassen. The Pragmatics of Explanation. American Philosophical Quarterly, 14(2): 143–150, 1977. ISSN 0003-0481.
- Bas C. van Fraassen. The Scientific Image. Number 4. Oxford University Press, 1980.
- Thomas F. Icard, Jonathan F. Kominsky, and Joshua Knobe. Normality and actual causal strength. *Cognition*, 161:80–93, April 2017. ISSN 0010-0277. doi: 10.1016/j.cognition.2017.01.010.
- Herbert M. Jenkins and William C. Ward. Judgment of contingency between responses and outcomes. *Psychological Monographs: General and Applied*, 79(1):1–17, 1965.
- Philip N Johnson-Laird and Fabien Savary. Illusory inferences: A novel class of erroneous deductions. *Cognition*, 71(3):191–229, 1999.
- Frank C. Keil. Explanation and understanding. *Annual Review of Psychology*, 57(Volume 57, 2006): 227–254, 2006.
- Jaegwon Kim. Hempel, Explanation, Metaphysics. *Philosophical Studies*, 94(1):1–20, May 1999. ISSN 1573-0883. doi: 10.1023/A:1004420102896.
- Lara Kirfel, Thomas Icard, and Tobias Gerstenberg. Inference from explanation. *Journal of Experimental Psychology: General*, 151(7):1481–1501, 2022. ISSN 1939-2222. doi: 10.1037/xge0001151.
- Lara Kirfel, Jacqueline Harding, Jeong Shin, Cindy Xin, Thomas Icard, and Tobias Gerstenberg. Do as I explain: Explanations communicate optimal interventions. In Larissa K Samuelson, Stefan Frank, Mariya Toneva, Allyson Mackey, and Eliot Hazeltine, editors, *Proceedings of the 46th Annual Conference of the Cognitive Science Society*, 2024.
- Philip Kitcher. Explanatory unification. Philosophy of Science, 48(4):507–531, 1981.
- Philip Kitcher and Wesley Salmon. Van Fraassen on Explanation. *The Journal of Philosophy*, 84 (6):315–330, 1987. ISSN 0022-362X. doi: 10.2307/2026782.
- Jonathan F Kominsky, Jonathan Phillips, Tobias Gerstenberg, David A Lagnado, and Joshua Knobe. Causal superseding. *Cognition*, 137:196–209, 2015.

- Thomas S. Kuhn. *The Structure of Scientific Revolutions*. University of Chicago Press, Chicago, 2nd edition, 1970.
- Isaac Lage, Emily Chen, Jeffrey He, Menaka Narayanan, Been Kim, Sam Gershman, and Finale Doshi-Velez. An evaluation of the human-interpretability of explanation, 2019. URL https://arxiv.org/abs/1902.00006.
- D. A. Lagnado, T. Gerstenberg, and R. Zultan. Causal responsibility and counterfactuals. *Cognitive Science*, 47:1036–1073, 2013.
- David Lewis. Causal explanation. In David Lewis, editor, *Philosophical Papers Vol. Ii*, volume 2, pages 214–240. Oxford University Press, 1986.
- Peter Lipton. Contrastive Explanation and Causal Triangulation. *Philosophy of Science*, 58(4): 687–697, December 1991. ISSN 0031-8248, 1539-767X. doi: 10.1086/289648.
- Tania Lombrozo. The structure and function of explanations. Trends in Cognitive Sciences, 10 (10):464–470, 2006.
- Tania Lombrozo and Emily G. Liquin. Explanation is effective because it is selective. *Current Directions in Psychological Science*, 32(3):212–219, 2023.
- Christopher G. Lucas and Charles Kemp. An Improved Probabilistic Account of Counterfactual Reasoning. *Psychological Review*, 122(4):700–734, 2015. doi: 10.1037/a0039655.
- R. Duncan Luce. Detection and recognition. In R. D. Luce, R. R. Bush, and E. Galanter, editors, *Handbook of Mathematical Psychology*, pages 103–189. Wiley, 1963.
- Judea Pearl. Causal diagrams for empirical research. Biometrika, 82(4):669–710, 1995.
- Judea Pearl. Causality. Cambridge University Press, 2009.
- Karl Pearson. The Grammar of Science. Meridan Books, third edition, 1911.
- Jonas Peters, Dominik Janzing, and Bernhard Schölkopf. *Elements of Causal Inference: Foundations and Learning Algorithms*. The MIT Press, Cambridge, MA, 2017. ISBN 9780262037310.
- Angela Potochnik. *Idealization and the Aims of Science*. University of Chicago Press, Chicago, 2017.
- Hilary Putnam. Philosophy and our mental life. In *Philosophical Papers*, volume 2, pages 291–303. Cambridge University Press, 1975.
- Tadeg Quillien and Christopher G. Lucas. Counterfactuals and the logic of causal selection. *Psychological Review*, 2023.
- Benjamin M. Rottman and Reid Hastie. Do people reason rationally about causally related events? markov violations, weak inferences, and failures of explaining away. *Cognitive Psychology*, 87: 88–134, 2016.

- Mathias Sablé-Meyer and Salvador Mascarenhas. Indirect illusory inferences from disjunction: a new bridge between deductive inference and representativeness. *Review of Philosophy and Psychology*, 13:567–592, 2022.
- Wesley C. Salmon. Statistical Explanation & Statistical Relevance. University of Pittsburgh Press, 1971.
- Wesley C. Salmon. A third dogma of empiricism. In Robert Butts and Jaakko Hintikka, editors, Basic Problems in Methodology and Linguistics, pages 149–166. Dordrecht: D. Reidel., 1977.
- Wesley C. Salmon. Why Ask, "Why?"? An Inquiry concerning Scientific Explanation. *Proceedings* and Addresses of the American Philosophical Association, 51(6):683–705, 1978.
- Wesley C. Salmon. Scientific Explanation and the Causal Structure of the World. Number 3. Princeton University Press, 1984.
- Wesley C. Salmon. Causality and Explanation, volume 52. Oxford University Press, 1998.
- Peter Spirtes, Clark Glymour, and Richard Scheines. Causation, Prediction, and Search. Springer Lecture Notes in Statistics, 1993.
- Michael Strevens. Depth: An Account of Scientific Explanation. Harvard University Press, 2008.
- Theodore R Sumers, Mark K Ho, Thomas L Griffiths, and Robert D Hawkins. Reconciling truth-fulness and relevance as epistemic and decision-theoretic utility. *Psychological Review*, 2023.
- Patrick Suppes. A Probabilistic Theory of Causality. North Holland, 1970.
- Nadya Vasilyeva, Thomas Blanchard, and Tania Lombrozo. Stable causal relationships are better causal relationships. *Cognitive Science*, 42(4):1265–1296, 2018.
- Thalia H Vrantsidis and Tania Lombrozo. Inside ockham's razor: A mechanism driving preferences for simpler explanations. *Memory & Cognition*, pages 1–29, 2024.
- Michael R. Waldmann, editor. The Oxford Handbook of Causal Reasoning. Oxford University Press, 2017.
- William C. Wimsatt. Re-Engineering Philosophy for Limited Beings: Piecewise Approximations to Reality. Harvard University Press, 2007.
- James Woodward. Making Things Happen: A Theory of Causal Explanation. Oxford University Press, 2003.
- James Woodward. Sensitive and Insensitive Causation. The Philosophical Review, 115(1):1–50, 2006. ISSN 0031-8108.
- James Woodward. Causation in biology: Stability, specificity, and the choice of levels of explanation. *Biology & Philosophy*, 25(3):287–318, June 2010. ISSN 1572-8404. doi: 10.1007/s10539-010-9200-z.

- James Woodward. Some varieties of non-causal explanation. In Alexander Reutlinger and Juha Saatsi, editors, Explanation Beyond Causation: Philosophical Perspectives on Non-Causal Explanations. Oxford University Press, 2018.
- James Woodward. Explanatory autonomy: The role of proportionality, stability, and conditional irrelevance. Synthese, 198(1):237–265, January 2021a.
- James Woodward and Christopher Hitchcock. Explanatory generalizations, part I: A counterfactual account. Noûs, 37(1):1–24, 2003.
- James Woodward and Lauren Ross. Scientific Explanation. In Edward N. Zalta, editor, *The Stanford Encyclopedia of Philosophy*. Metaphysics Research Lab, Stanford University, summer 2021 edition, 2021.
- James C. Woodward. Causation with a Human Face: Normative Theory and Descriptive Psychology. Oxford University Press, 2021b.
- Stephen Yablo. Mental causation. The Philosophical Review, 101(2):245–280, 1992.
- Petri Ylikoski. The Idea of Contrastive Explanandum. In Johannes Persson and Petri Ylikoski, editors, *Rethinking Explanation*, pages 27–42. Springer, 2007.
- Petri Ylikoski and Jaakko Kuorikoski. Dissecting explanatory power. *Philosophical Studies*, 148 (2):201–219, 2010.
- Erica J Yoon, Michael Henry Tessler, Noah D Goodman, and Michael C Frank. Polite speech emerges from competing social goals. *Open Mind*, 4:71–87, 2020.
- Jeffrey C Zemla, Steven Sloman, Christos Bechlivanidis, and David A Lagnado. Evaluating everyday explanations. *Psychonomic Bulletin & Review*, 24(5):1488–1500, 2017. doi: 10.3758/s13423-017-1258-z. URL https://doi.org/10.3758/s13423-017-1258-z.
- Jeffrey C Zemla, Steven A Sloman, Christos Bechlivanidis, and David A Lagnado. Not so simple! causal mechanisms increase preference for complex explanations. *Cognition*, 239:105551, 2023.