

Inferentialist Expressivism for Explanatory Vocabulary

April 9, 2017

^{JM}I'm using the following convention: 'Theorem' is reserved for the main results of the paper; 'Lemma' is used for claims that are needed in proofs of Theorems; and 'Proposition' is reserved for claims that are important but not the main results.

Outline Thesis: By treating explanation as a form of inference, we give an inferentialist-expressivist account of explanatory vocabulary that specifies the conditions under which IBE is a reliable form of ampliative inference.

1 Introduction

By treating explanation as a form of inference, we give an inferentialist-expressivist account of explanatory vocabulary that specifies the conditions under which IBE is a reliable form of ampliative inference.

2 Motivation

Depending on your philosophical interests, an inferentialist-expressivist approach to explanation may have different attractions. To that end, we offer customized advertisements to inferentialists and to philosophers of science. Ideally, we hope that you are motivated by both sets of considerations!

2.1 For the Inferentialists

For those with broad sympathies to inferentialist-expressivism, there are three main attractions to treating explanatory vocabulary in a manner akin to traditional semantic vocabulary. First, despite the fact that a sizable chunk of our inferential lives traffic in the inductive, surprisingly little work has been done to develop an inferentialist-expressivist approach to inductive vocabulary. Rather, the paradigmatic texts in inferentialist-expressivism focus on deductive vocabulary. There is a potential goldmine that is largely unexplored: probabilistic, statistical, and (of course) explanatory vocabularies might be glittering with unclaimed inferentialist-expressivist riches.

While few inferentialists have reflected upon the nature of inductive and explanatory vocabulary, those that have offer some tantalizing insights. The grandfather of inferentialism, Wilfrid Sellars¹⁹⁵⁴; Sellars¹⁹⁶³, saw the scientific activities of describing and explaining the world as *the* meaning-conferring practices of empirical vocabulary. In other words, it is their role in the game of giving and asking for *explanations* that determines the content of our empirical concepts. Moreover, Sellars thought that it was the business of lawlike statements or *nomologicals*, as he called them, to make explicit the material inferences that underwrite this practice.

In his most recent work, **Brandom2015** has spelled out Sellars' view in great detail, but his earlier work also offered reflections on induction. **Brandom1994; Brandom2008** claims that material inferences come in two flavors—those that preserve commitment and which are analogs to deduction, and those that preserve entitlement and are analogs to induction. Motivating this distinction is the idea that entitlement-preserving inferences are akin to inductive inferences insofar as they permit or license rather than oblige an agent to believe or assert the conclusion. Our work here represents the first attempt to develop this view of material-inductive inferences in a systematic fashion.

Finally, as a species of expressivism, such an approach would highlight—and settle—certain metaphysical debts that have been overlooked in the explanation literature. The dominant philosophical accounts of scientific explanation hold that explanations represent various kinds of “dependency” relationships—paradigmatically causal ones (**Lewis1986; Salmon1984; Strevens2008; Woodward2003**). More provocatively, we might say that these dependency relationships are truthmakers for their respective explanations. Yet, virtually all of these dependency relationships are modally-loaded concepts. This, in turn, begets standard “placement problems,” especially about how modal objects and properties fit within a naturalistically respectable ontology (**Price2011**). Anti-descriptivism—shared by expressivists of all stripes—would appear to be one way of tidying up the ontological backdrop of the scientific explanation literature. Closely related, by “inferentializing” IBE, inferential-expressivists rob those who partake in a broader array of metaphysical extravagances of one of their most powerful gateway drugs. Various kinds of realism (mathematical, modal, etc.) are frequently justified in virtue of their ability to explain otherwise “miraculous” facts. If the view we are suggesting is correct, then explanatory claims can fit within a larger inferentialist-expressivist framework, which tends towards deflationism. Furthermore, there would be a sense in which explanation is just not the right tool for justifying substantive metaphysical claims: why would we ever think that certain things exist (at least in any non-deflationary way) on the basis of what we can do with an expressive device? ^{KK}This might be a bit too strong; rein me in here. ^{JM}I think this is ok, but then again, I'm a bad judge of über-deflationary claims.

- Explanatory role of reasons (Sellars @ Explanatory Coherence). Space of reasons includes explanatory arguments. Hat-tip to IBE.
- General expressivist motivations.

2.2 For the Philosophers of Science

For philosophers of science, an inferentialist-expressivist account of explanation promises three kinds of attractions. First, many of the previous section's selling points will be tantalizing to those philosophers of science with metaphysical proclivities that overlap with expressivists'—i.e. that tend towards the austere.

Second, going back as far as Hempel, many theorists of explanation have endorsed a view that we call *Explanation As Inference* (EAI). On this view, explanations are a kind of inference. ¹ Hempel famously analyzed certain explanation statements in terms of sound deductive arguments whose premises contain at least one law of nature. Likewise, Kitcher's unificationist account treats explanations as arguments belonging to a class of formal derivations circumscribed by semantic constraints on substitutions—so-called ‘filling instructions.’ However, EAI has faced many searching counterexamples (**Salmon1989**). Elsewhere, we identify the source of these problems: earlier proponents of EAI have marshaled the wrong

consequence relations. The current paper offers a precise account of these consequence relations (**Khalifaforthcoming**). By itself, this would simply be a technical modification (albeit significant) to an otherwise-correct account of explanation. However, the current paper also suggests a second (and deeper) way in which inferentialist-expressivism can further the cause of EAI: its earlier proponents of EAI saw explanations as *reducible* to or *explicated* by a class of arguments, yet, if our approach is correct, then there are *prima facie* reasons to think that the proper relationship is that of *expression*. In other words, the tight connection that has already been established between explanations and inferences by earlier EAI'ers may be readily exploited and profitably recast by expressivists.

A final attraction of approaching explanation as we do is the prospect of identifying the conditions wherein IBE is reliable. The spectrum of positions on this particular topic is overwhelming, with some theorists voicing skepticism about the probity of IBE (**Laudan1981; Stanford2006; Fraassen1989**), others claiming it is reliable because it tracks with a more fundamental form of inference (**Day1994; Rappaport1996; Khalifaforthcoming**) (add **Nortonms,Roche201?**), and still others claiming that IBE is both reliable and fundamental (**Harman1965; Harman1986; Lycan1988; Poston2014**). These disagreements have always had an air of intractability because IBE has not been stated with adequate precision. By clarifying the relationship between explanation and inference, as is demanded by an inferentialist semantics, we thereby shed some light on when it is licit to infer the best explanation.

- Heir to Hempelian explanation (see other papers)
- Ascertaining when IBE is reliable
- IBE as inference-first and explanation-second.

3 Inferentialism, Expressivism, and Inferentialist Expressivism

If there are two positions that unite the work of Wilfrid Sellars, Robert Brandom, and Jaroslav Peregrin, then they are undoubtedly *inferentialism* and *expressivism*. ^{JM}This was just a line to jump-start my thinking. Help me think of a better way to start the paper off. I can see the paper starting off with a programmatic statement of the thesis followed by this discussion of expressivism. ^{KK}we'll get back to this once we have something more complete. I think that will help with programmatic statements.

Roughly put, *inferentialism* is the view that the meaning of a linguistic expression is determined by its inferential role, that is, by what follows from it and/or by what it follows from.² In this sense, inferentialism stands in opposition to the traditional approaches to semantic analysis which take the meaning of linguistic expressions to be functions of their representational purport and success. For instance, inferentialism holds that the meaning of logical constants is determined first and foremost by the rules governing their introduction and/or their elimination, rather than by the function of constituent truth-values. More ambitiously, inferentialism may be applied to the analysis of ordinary empirical vocabulary, in which case, the notion of *inference* is expanded to include not only proper transitions among assertions but also those to and from non-verbal behavior.

Indeed, in order to successfully explicate the meaning of any vocabulary inferentialism needs to provide some account of what qualifies as an inference. The account favored by

¹We introduce EAI to distinguish it from IBE, and from the broader sense in which we offer an inferentialist semantics of "best explains."

²We formulate the position of inferentialism using the explicitly inclusive 'and/or' in order to represent as approaches that identify meaning-constitutive rules with those for introduction (**Prawitz2006**), elimination (**Dummett1991; Prawitz2007**), or both (**Brandom1994; Brandom2008**).

Sellars, Brandom, and Peregrin treats inferences (as well as assertions) as theoretical posits sufficient for describing the behavior of agents playing a ‘game of giving and asking for reasons.’ Such a practice involves the acquisition, maintenance, and discharge of normative statuses. These statuses comprise various kinds of *commitments* and *entitlements*. Thus, for example, in making an assertion, a speaker entitles others to reassert her claim and commits herself, *inter alia*, to providing reasons for it when it is challenged. *Inference* is the theoretical concept that enables us to describe this practice—e.g. when acknowledging one assertional commitment commits a speaker to undertake another, we may say that the second claim *follows* from the first.

Unlike inferentialism, *expressivism*, in the sense that describes the positions of these thinkers, is not a thesis about linguistic meaning. Rather, it is a claim about the role or function of some vocabulary, about what this vocabulary permits speakers to *do*. Broadly construed, expressivism about some vocabulary *X* holds that the function of *X* is not to describe or represent features of the world, but rather to express something or to make something explicit that would otherwise remain implicit, tacit, or latent. On this broad construal, deflationists about truth count as expressivists, for they treat the function of the truth predicate as enabling us to express *generalizations* that would otherwise require knowledge we might lack or repetition we might rather do without (Horwich1999). Likewise, Thomasson2013 has argued that the role of metaphysical modal vocabulary is to make explicit constitutive semantic rules and to do so in a prescriptive manner while employing the indicative mood. On a narrower conception of expressivism, neither generalizations nor rules are expressed. Rather, our *attitudes* are expressed. Thus, moral expressivists treat the function of ethical discourse as expressing *evaluative attitudes* rather than as tracking moral facts or properties (Blackburn1993; Gibbard2003). Similarly, Brandom’s logical expressivism holds that the function of logical vocabulary is to make explicit our *endorsement* of various patterns of inference rather than to represent features of a ‘logical realm’ (Brandom1994; Brandom2000; Brandom2008). Of course, even the “non-attitudinal” expressivists can be recast so that their target vocabulary expresses an attitude—an *endorsement* of a generalization, a *prescription* of how to use certain words.

Whether narrow or broad, an expressivist treatment of a vocabulary’s function need not be wedded to an inferentialist account of its meaning. Nonetheless, the two approaches have a natural affinity. Both eschew traditional programs in the study of language that, by taking the word-world relation as fundamental, privilege representational concepts and descriptive functions. More importantly, certain benefits accrue to the theorist who holds both positions for the same target. An inferentialist account of the meaning of vocabulary *X* may, indeed, *should*, go some way toward explaining why *X* has the peculiar expressive function that it does. For instance, the introduction and elimination rules for the conditional in classical logic permit the derivation of the Deduction Theorem ($A \vdash B \Leftrightarrow \vdash A \rightarrow B$) which can be thought of as articulating the expressive function of the conditional—the theorem says that an object-language operator gives expression to or makes explicit what would otherwise be given as a rule of inference in the meta-language. Finally, the example of the conditional also illustrates the way in which inferentialism supports a particular brand of expressivism, one distinguished by the claim that the expressive role of the target vocabulary is to make explicit speakers’ *commitments to inferences*.

It is this variant of expressivism, which we will call *inferentialist expressivism*, that Peregrin and Brandom apply to traditional logical operators and that the latter extends to semantic, modal, and representational vocabulary. Indeed, this same approach can be detected in Sellars’ assertion that statements of natural law express commitments to patterns of inference

among empirical claims. The inferentialist-expressivist approach to these various vocabularies stands as both an invitation to others to expand its application and as a paradigm of how to proceed in such an endeavor. In this paper, we pursue just such an expansion, aiming to provide an inferentialist-expressivist treatment of *explanatory vocabulary*.

Our first task is to develop a base logic that formalizes the inferences expressed by explanatory vocabulary. For convenience, we refer to these inferences as *explanatory arguments*. To make good on the expressivist promise, the logic of explanatory arguments must exhibit the salient properties of explanatory vocabulary. As we shall see, this is a tall order and consequently much of our argumentative energy is spent demonstrating that the inferences in our base logic have just those properties they would need if they were explicitly endorsed by the use of explanatory claims.

Our next task is to show that the object-language of this logic can be extended to include a connective that makes (commitment to) explanatory arguments explicit. To do so, we must provide rules for the this connectives' introduction and elimination. Since one of our goals is to show that the meaning of explanatory vocabulary is partly constituted by the role it plays in IBE, one of the rules for this connective must be a recognizable form of IBE. From these rules, it should be possible to establish a deduction theorem. However, since the logic of explanatory arguments contains two distinct classes of consequence relations, we will not be able to establish the traditional version of this theorem. Instead, we will show that the use of the explanatory connective in one class of consequence relations codifies a rule in the other. Since this captures the expressivist purport of the classical deduction theorem, we dub this a *quasi-deduction theorem*.³

Finally, we must defend our explanatory connective against charges, such as that leveled by **Prior1960**, that it trivializes a consequence relation and thus fails to be meaningful. In response to this worry, most inferentialists follow **Belnap1962** and say that the rules by which we introduce a new bit of logical vocabulary must extend the consequence relation we start with in a conservative manner. This means that an inference that does not contain the new connective holds in the extended consequence relation if and only if it already held in the unextended one. Since there is a cut-elimination theorem for our base logic, the easiest way to establish the conservativity of our extension is to show that all proofs that use cut in the extension can be reduced to proofs which do not. In other words, we will demonstrate that the cut-elimination theorem holds in our extended logic.

4 Explanatory Vocabulary

To start, let us specify our target. We aim to show that the role that a particular class of expressions plays in a language is that of enabling speakers to express their commitment to certain inferences, rather than to describe features of the world. The class we have in mind includes those expressions which are given in English as “That... (best, possibly, actually) explains why...,” and its nominalization “That... is the (best, possible, actual) explanation of why...,” where the ellipses are filled by declarative sentences. As with most complex locutions in natural language, the analysis of these expressions faces challenges ranging from issues of polysemy to the “non-intersective” behavior of the adjectival predicates. While we certainly wish to see the insights of inferentialist-expressivism brought to bear in the analysis of natural languages, we do not intend to pursue such an ambitious undertaking. Instead, we aim to capture at least part of the notion of *best explanation*. Unfortunately, even here, there is ambiguity. For instance, “A explains B” is typically factive. On a simple reading, then “A

³Our quasi-deduction theorem can be thought of as an instance of what XXXX calls a “Generalized Deduction Theorem”.

best explains B ” is also factive. After all, it is natural to think of our best explanations as a proper subset of the totality of actual explanations. However, our chief aim is to capture the meaning of “ A best explains B ” as it functions as a major premise in the pattern of inference dubbed “Inference to the Best Explanation” (IBE). In its simplest rendition, IBE takes the following form:

$$\begin{array}{c} A \text{ best explains why } B \\ \hline B \\ \therefore A \end{array}$$

As IBE’s foremost defender, Peter Lipton argues that to treat “ A best explains B ” as factive or “actual” in this case is “like a dessert recipe that says start with a soufflé.” (Lipton2004). For this reason, Lipton suggests that the best explanation ought to be construed as the best *potential* explanation.⁴ Following a common convention in the literature on scientific explanation, let us stipulate that “ A potentially explains B ” requires neither “ A ” nor “ B ” to be true. Then our claim is that “ A best explains B ” is not factive, but must still provide good inductive grounds for detaching the explanans, A , when the explanandum, B , obtains.

In addition to being non-factive, the sense of *best explanation* that we aim to represent is *immediate* and *exhaustive*. We say that A is an immediate explanation of B so long as it does not explain B merely by explaining something else, C , which in turn explains B . In other words, the explanations we have in mind are not transitive. By “exhaustive,” we mean that nothing needs to be added to A in order for it to explain B . The target of our account is thus expressed by the locution “... is the loveliest potential, immediate, exhaustive explanation of why ...”. To avoid having to repeat this phrase, we will henceforth speak of *best explanations* or simply *explanations*.

^{JM}I just built non-transitivity into our target. This seems cleaner to me. Now it’s no longer a property we have to justify.

5 Properties of Explanatory Arguments

If, *ex hypothesi*, explanatory expressions *do* serve to make explicit certain inferential commitments, the first question that arises is, what kinds of inferences are being made explicit? Let us call these target phenomena *explanatory arguments*. Our question then becomes: what are the properties of explanatory arguments that distinguish them from other inferences? We can certainly reconstruct a Hempelian answer to this question. For Hempel, (certain) explanatory arguments are just inferences with true premises that include at least one lawlike statement. However, given the notorious difficulties that plagued Hempel’s deductive-nomological account, we should resist cleaving too closely to this model. Instead, we propose to work backwards from scientific practice concerning the appropriate use of explanatory expressions and the general inferential relationships they exhibit. In other words, we proceed from the assumption that the properties of explanatory arguments may be read off of the behavior of explanatory practices.

Explanatory arguments exhibit six distinctive properties. First, they are *defeasible*—an argument can cease to be explanatory when new information is added to its premises. For instance, while the claim that “The liquid’s acidity explains why the blue litmus paper turned

⁴Lipton also argues that IBE should construe the best explanation as the “loveliest” explanation, and then vaguely describes “loveliness” as a combination of various “theoretical virtues,” such as simplicity, scope, fit with background belief, etc. While our view does appeal to some superlative—what we call “sturdiness”—it appears to be something quite different than this. A detailed comparison of loveliness and sturdiness exceeds the scope of the current paper.

red,” may constitute a good explanation, strengthening the explanans can easily produce a bad one: “The liquid’s acidity and the presence of chlorine gas explains why the blue litmus paper turned red.” In this case, the additional information is incompatible with the explanandum. In other cases, however, what is added to an explanans undercuts the explanatory relation itself. For instance, that the match was struck may explain why the match lit, but if it is discovered that match was damp, then the striking no longer, by itself, explains why it lit (perhaps the match was dried before it was struck). In both sorts of cases, we say that explanatory arguments in question are defeated. Since classical deductive inferences are monotonic—i.e. the premises of a valid inference can be arbitrarily expanded and the resulting inference is also valid—capturing the defeasibility of explanatory arguments cannot be done within a strictly deductive system.

Second, explanatory arguments are *minimal* in the sense that all of their premises are needed to infer their conclusion. More precisely, minimality insists that (M1) no proper subset of the premises of an explanatory argument are explanatory (of the same explanandum), (M2) nor any set comprised of conjunctions of members of a proper subset of the premises, (M3) nor any set comprised of disjunctions of two or more members of a subset of the premises. For example, if the set of sentences $\{A, B, C\}$ is a minimal explanatory argument for D , then none of the following are explanatory arguments for D : $\{A, B\}$, $\{A \wedge B\}$, $\{B \wedge C\}$, $\{A \vee B\}$, $\{A \vee (B \vee C)\}$. In general, minimality prohibits the addition of irrelevant information to explanatory arguments. For example, if the liquid’s acidity explains why the litmus paper turns red, then it does not follow that the liquid’s acidity *and its potability* explains why the litmus paper turned red, even if the liquid is in fact potable. The minimality of explanatory arguments is thus another reason why they cannot be modeled by deductive, i.e. monotonic inferences.

Third, since contradictions do not explain, explanatory arguments must have consistent premises. Classically valid arguments adhere to the principle of *ex contradictione quodlibet* (ECQ)—i.e. $A, \neg A \vdash B$ for any arbitrary B —meaning that any argument with inconsistent premises is classically valid. Logical systems that renounce ECQ are called either paraconsistent or premise-consistent. Paraconsistent systems are able to represent reasoning that is *tolerant* of inconsistency and that may, in some cases, even validate contradictions. By contrast, a *premise-consistent* logic aims, in the opposite direction, to represent reasoning that is *intolerant* of inconsistency. Since contradictions never explain, explanatory arguments are premise-consistent, not paraconsistent. Excluding theorems with contradictory premises, however, comes at the cost of logical strength. In our system, for example, disjunctive syllogism is not a theorem, so it can only be retained as a rule. Finally, note that premise consistency describes yet another form of nonmonotonic behavior.

Fourth, explanatory arguments are *irreflexive*.⁵ For instance, “the litmus paper turned red because it turned it red” is not an explanation. More generally, no explanatory arguments should be of the form $A \vdash A$. Even partial self-explanations seem unacceptable; $A \wedge B \vdash A$ and its ilk should also be prohibited from qualifying as explanatory arguments. Like ECQ, reflexivity partly defines the consequence relation of classical logic. However, while many logics abandon ECQ, few logics abandon the principle of reflexivity, let alone try to preserve the property of irreflexivity.⁶ One particular difficulty arises from the fact that irreflexive logics cannot easily recover Modus Ponens (MP), yet many explanations require Modus Ponens. Our irreflexive system preserves Modus Ponens.

⁵While a *non-reflexive* system would be one in which the principle of reflexivity fails, perhaps only on occasion, an *irreflexive* system is one in which no instance of reflexivity holds.

⁶While there have been some logics in which reflexivity fails, irreflexivity holds in very few. Notable exceptions include input-output logics, logics of grounding, and quantum logics.

⁰Logicians of inferentialist and proof-theoretic persuasion have already explored systems in which transitivity fails (Ripley2011;

Fifth, explanations are *stable*, which philosophers of science have analyzed in different ways (Hempel1965; Lange2009; Mitchell2003; Skyrms1980; Woodward2003). In its most general form, X is said to be stable if X remains unchanged as other conditions C change. For instance, suppose that a patient’s rash is explained by a particular bacterial infection, though an alternative potential explanation is that the rash is caused by an allergic reaction. The explanation is stable insofar as she would have a rash regardless of whether she had had an allergic reaction. Typically, the fundamental bearers of stability are taken to be laws or generalizations. By contrast, as inferentialists, we take explanatory arguments to be the fundamental bearers of stability. In what follows, we develop a distinctively inferentialist brand of stability that we call “sturdiness.”

We baptize as *non-trivial* those defeasible inferences that are premise-consistent and ir-reflexive. Sturdiness is then a comparative property among non-trivial inferences. A non-trivial inference is sturdy just in case it succeeds when all other non-trivial inferences that share its conclusion fail, where by ‘failure’ we mean *defeat* and by ‘success’, the absence of failure. More precisely, the failure we have in mind is that which obtains when the premises of an inference are false. Since premise-consistent inferences do not hold when their premises are false, a sturdy inference is one that remains undefeated when the premises of its competitors—i.e. those non-trivial inferences that share its conclusion—are false. It should be clear from what has been said thus far that representing sturdiness requires us to treat premise-*inconsistency* as a form of defeat. As we shall see, our formalism enables us to do just this.

Sixth, explanations are *detachable*. If the bacterial infection is the best explanation of the rash, then we may conclude that the patient has a bacterial infection. By detaching the explanans, we may make further predictions, design interventions, and construct new models. This is the animating idea behind Inference to the Best Explanation. A consequence relation that captures explanation must thereby underwrite abduction.

^{JM}Transition

6 The Logic of Explanatory Arguments and IBE

In this section, we develop sequent calculi for explanatory arguments and the vocabulary that makes them explicit. We begin by introducing our base logic LEA, defined over the standard propositional language, \mathcal{L} . LEA is composed of two parts. The first part, LK^\ominus , is based on a variant of Gentzen’s sequent calculus for classical logic that was proposed by Piazza2015 with the aim of representing nonmonotonicity in terms of an inference’s context-sensitivity. Although we have altered most of the definitions, terminology, and rules with which Piazza2015 introduced their calculus, we retain the key technical features of LK^S and inherit their proof of the cut-elimination theorem. Our modifications are intended, on one hand, to represent a concept of *defeat* more appropriate to the behavior of explanations, and, on the other, to extend the axioms of the system to non-logical, material inferences. The second part of the base system, LE^\blacktriangleright introduces an additional class of sequents and rules that govern their interaction with those of LK^\ominus , including a rule for abduction or IBE. Finally, we describe a system LEA^+ defined over $\mathcal{L}_\blacktriangleright$ that extends \mathcal{L} to include an object-language connective, \blacktriangleright , that makes commitments to explanatory arguments explicit.

The need for a base logic that includes two consequence relations stems from the following observation: If the function of locutions like *best explains why* is to express commitment to explanatory arguments, and if IBE is a legitimate rule of inference, then IBE cannot itself

Tennant2014; Hlobil2016).

be an explanatory argument. Here is the support for this claim. Suppose that the locution *A best explains why B* does make explicit a commitment to an explanatory argument. Now suppose (for *reductio*) that IBE is also an explanatory argument. From these suppositions it follows that the *best-explains-why*-connective expresses an explanatory argument from *B* to *A*, since this is the only way for there to be a good explanatory argument with the form of IBE: *A best explains why B*, *B* / \therefore *A*. But it also follows that such an inference can itself be made explicit by the *best-explains-why*-connective, thereby yielding the sentence $[(A \text{ best explains why } B) \wedge B] \text{ best explains why } A$. In so far as sentences of this form are even intelligible, it is far from obvious that they *claim* what an application of IBE *shows*. Thus, we should not conclude that IBE is an explanatory argument. Rather, IBE is one sort of inference and that made explicit by explanatory vocabulary is another.⁷ This means that any logical system designed to represent both explanatory arguments and IBE must appeal to two distinct consequence relations.

In recognition of this point, LEA contains two consequence relations, or, more precisely, two *classes* of consequence relations, for as we shall see the system contains infinitely many consequence relations. The class whose rules are given by LK^Θ is intended to represent those sets of inferences to which both IBE and *candidate* explanatory arguments belong. The class of consequence relations introduced by $\text{LE}^\blacktriangleright$ and denoted by $\vdash^\blacktriangleright$ is supposed to capture the behavior of explanatory arguments themselves. The chief results of this section are (1) that the theorems of LEA constructed with $\vdash^\blacktriangleright$ exhibit all of the properties associated with explanatory arguments (Theorem 2), (2) that LEA can be extended (LEA^+) to include an object-language expression for making explicit commitments to explanatory arguments (Theorem 3), and (3) that this extension is conservative (Corollary 4.1).

In what follows, we assume a propositional language, \mathcal{L} , for classical logic, that consists of a countably infinite set of atomic sentences $At = \{p_1, \dots, p_n, q_1, \dots, q_n\}$, the binary connectives \wedge and \vee , and the unary connective \neg . Let A, B, C, D range over formulas; let $\Gamma, \Delta, \Sigma, \Theta$ range over sets of formulas; let S, T, U range over sets of sets of formulas, and let X, Y, Z range over sets of atoms.

We begin by employing the standard sequent notations—e.g. $\Gamma, A \vdash B, \Delta$. Formulas on the left side of the turnstile are called the *antecedent*; on the right side they are called the *succedent*. Commas in the antecedent are read ‘conjunctively’ and those on the right are read ‘disjunctively’. The formula with the connective in a rule is the *principal* formula of that rule, and its components in the premises are the *active* formulas.

The sequents in our calculi depart from the standard form in two respects: just below our turnstile we add a set of formulas, Θ , called a *defeater set* and to the far left of the turnstile we add another, Σ , called a *background set*.

$$\Sigma \mid \Gamma \vdash_{\Theta} \Delta$$

Defeater sets contain information whose addition to the premises would defeat the inference represented by the sequent. Roughly put, a sequent is defeated whenever its antecedent or background set contains a formula that is logically equivalent to a subset of the defeater set. Thus, as the name suggests, defeater sets are sets of inference-defeaters.

As noted, we adopt and modify the sequents and rules of the calculus LK^S developed by **Piazza2015**. The general idea captured by this calculus is that any application of the rules along a derivation ought to preserve not only the validity, but also the defeat-status of sequents. In our system, LEA, this means that for any derivation tree, π , constructed by

⁷Treatments of abduction in the computer science literature have long recognized this point, if only implicitly, insofar as they have viewed abduction as a form of (restricted) deduction in reverse.

recursive applications of the rules (excluding *cut*), the following holds: if a defeated sequent occurs in a branch of π , then all sequents below it are also defeated. By tracking the defeat of sequents, our proof theory can identify if and when a line of defeasible reasoning goes astray.

We interpret background sets as consisting of information that is available to a reasoner when she draws an inference, but which does not serve as a premise and from which the conclusion is not said to follow. Having such a device in our formalism enables us to capture an important aspect of defeasibility, namely, that the introduction of new information may jeopardize prior inferential commitments even when that information does not serve as fodder for new inferences in its own right. These sets also serve a technical role in our system since they often expand in the course of a derivation, picking up traces of those formulas shifted by the rules from the left to the right side of the turnstile. (See $\neg \vdash$ in Figure 1). It is this latter feature which permits the implementation of a Gentzen-style normalization procedure and proof of cut-elimination. Indeed, by keeping a *record*, so to speak, of formulas that have moved from the antecedent in a premise to the succedent of a conclusion, background sets ensure that all the sequents in a cut-free derivation are undefeated just in case its end-sequent is undefeated.⁸

Because the provability of any sequent in LEA depends, in part, upon the contents of its defeater set, there is not one or two but $\mathcal{P}(\mathcal{L})$ -many consequence relations represented by the calculi. Some are classical, i.e. $\Theta = \emptyset$; many are non-monotonic, i.e. $\Theta \neq \emptyset$; others defy even the most ubiquitous structural properties, such as reflexivity, i.e. $\Delta \subseteq \Theta$. By specifying the contents of defeater sets, the rules of our calculi are able to exploit this panoply so as to home in on the class of consequence relations that bears precisely those properties we associate with explanatory arguments.

Before delving into the details of the system, we offer an informal gloss on our defeasible sequents. As noted, we follow **Piazza2015** to the extent that proofs in LEA preserve both validity and undefeatedness. However, since the rules in our calculi, and those of LE^\triangleright in particular, are not deductive, it cannot be deductive ‘validity’ that is preserved. Instead, we follow **Brandom2008** and treat the sequents that belong to a proof as inferences that preserve *entitlement*. In keeping with this view, we offer the following reading of the defeasible sequent above: *Anyone entitled to (every member of) Γ is entitled to (at least one member of) Δ , given the background of Σ .*

^{JM}Issues with normative pragmatics interpretations of multiple succedent sequents.

^{JM}transition

Definition 1 (Defeater sets, Background sets, Defeasible Sequents). Defeater sets are sets of formulas that defeat an inference (see below). Background sets are sets of formulas that represent the background knowledge of the inference. A defeasible sequent is a standard sequent with a *background set*, Σ , and a *defeater set*, Θ , attached:⁹

$$\Sigma \mid \Gamma \mid_{\Theta} \Delta$$

When no background sets have been specified (i.e. $\Sigma = \emptyset$) we write:

$$\cdot \mid \Gamma \mid_{\Theta} \Delta$$

In order to state the conditions under which a defeasible sequent is defeated, we must introduce some preliminary concepts.

⁸In **Piazza2015** what we call background sets are simply referred to as *repositories* and while they play the same technical role, they are not provided with an substantive interpretation.

⁹In general, the respective roles played by *control sets*, *compatibility*, and *soundness* in Piazza and Pulcini’s LK^S are played by *defeater sets*, (our notion of) *compatibility*, and *undefeatedness* in our system. The crucial difference between the two approaches is that while in **Piazza2015**, the occurrence of a disjunctive formula in the antecedent would render the sequent defeated (resp. unsound) if either of the

Definition 2 (\mathcal{D} -Rules). Let \mathcal{D} be the following set of rules:

$$\mathcal{D}_\wedge: A \in \mathcal{L} \Rightarrow A \wedge B \in \mathcal{L}$$

$$\mathcal{D}_\vee: A \in \mathcal{L} \text{ and } B \in \mathcal{L} \Leftrightarrow A \vee B \in \mathcal{L}$$

$$\mathcal{D}_\neg: A \in \mathcal{L} \Leftrightarrow \neg\neg A \in \mathcal{L}$$

$$\mathcal{D}_{\neg\vee\wedge}: \neg A \vee \neg B \in \mathcal{L} \Leftrightarrow \neg(A \wedge B) \in \mathcal{L}$$

$$\mathcal{D}_{\neg\wedge\vee}: \neg A \wedge \neg B \in \mathcal{L} \Leftrightarrow \neg(A \vee B) \in \mathcal{L}$$

These rules ought to be rather familiar. \mathcal{D}_\neg is the principle of Double Negation and $\mathcal{D}_{\neg\vee\wedge}/\mathcal{D}_{\neg\wedge\vee}$ are *de Morgan's Laws*. Read left to right, \mathcal{D}_\vee is Conjunction Introduction with \wedge substituted for \vee and similarly, \mathcal{D}_\wedge is just Disjunction Introduction with the converse substitution. The motivation for these substitutions as well as the peculiar lack of a biconditional in \mathcal{D}_\wedge will become clear in a moment.

Definition 3 ($\mathcal{D}(\Theta)$). Let $\mathcal{D}(\Theta)$ be the closure of Θ under the \mathcal{D} -rules as well as the commutativity, associativity, and distributivity of conjunction and disjunction, respectively.

Definition 4 (Compatibility). A set of formulas, Γ , is said to be compatible with a defeater set, Θ , just in case the conjunction of the members of Γ is not included in $\mathcal{D}(\Theta)$. We use ' \succsim ' to symbolize compatibility.

$$\Gamma \succsim \Theta \text{ iff } \bigwedge \Gamma \notin \mathcal{D}(\Theta)$$

Example 1. $\{p_1 \vee q_1, \neg p_2, p_3 \wedge q_2\} \succsim \{p_1, p_2, \neg p_3 \vee \neg q_1\}$

Example 2. $\{\neg(p_1 \vee q_1), \neg\neg\neg p_2, \neg p_3 \wedge q_2\} \succsim \{p_1 \wedge p_2, p_3, \neg q_2\}$

Example 3. $\{\neg p_1 \vee q_1, \neg p_2, \neg p_3 \wedge q_2\} \not\succsim \{q_2\}$

Example 4. $\{\neg(p_1 \vee q_1), \neg\neg\neg p_2, \neg p_3 \wedge q_2\} \not\succsim \{\neg p_1, \neg p_2\}$

Remark 1. For any set of formulas Γ , $\Gamma \succsim \emptyset$.

Definition 5 (Defeat). A defeasible sequent, $\Sigma \mid \Gamma \overline{\Delta}$, is said to be *undefeated* whenever $\Sigma \cup \Gamma \succsim \Theta$ and defeated otherwise.

The \mathcal{D} -Rules are designed to generate closed sets of formulas whose occurrence in the antecedent defeats the sequent to which the defeater set is attached. Since a formula of the form $A \wedge B$ defeats an inference just in case one or more of its constituents belongs to the defeater set, \mathcal{D}_\wedge is constructed so that a set closed under it will contain all those conjunctions for which at least one member of the original set, Θ , is a conjunct. Conversely, if the defeater set contains a conjunction but neither of its conjuncts, then we know that the two formulas *together* defeat the sequent but not whether either formula by itself would be sufficient for defeat. (The need for defeat to occur when both conjuncts appear on their own in the antecedent is handled by the conjunctive formulation of *compatibility*.) Thus, unlike the other

disjuncts occurred in the sequent's defeater set (resp. control set), in our system, the sequent would only be defeated if both disjuncts were present in the defeater set. We believe that the latter property captures a more intuitive, less cautious conception of defeat. Unfortunately, this conception of defeat could only be purchased at the cost of attaching a *proviso* to $\vee \vdash$ that restricts the rule's application to sequents whose active formulas are compatible with their respective defeater sets. On balance, we were willing to trade the elegance of the rules for the more accurate portrayal of defeat. In order to realize this conception in our definitions, we found it necessary to deploy substantially different operations and have re-named the resultant concepts so as to avoid confusion. With this said, Definition 6 and Lemma 1 are taken over from [Piazza2015](#) with very little modification.

\mathcal{D} -rules, \mathcal{D}_\wedge only permits the construction of more complex formulas—hence the absence of a biconditional in its formulation.

In contrast, a disjunctive formula defeats an inference only when both of the disjuncts occur in its defeater set. The \mathcal{D}_\vee -rule captures this intuition—read left-to-right—by constructing disjunctions when both disjuncts are present in the defeater set. Conversely, a disjunction in a defeater set means that the presence of either disjunct in the antecedent will defeat the sequent. Thus, the \mathcal{D}_\vee -rule is formulated—from right-to-left—so that a disjunction in the \mathcal{D} -closure of a defeater set will contain both disjuncts. The result of comparing the antecedent (and background set) with the closure of the defeater set under the \mathcal{D} -Rules is a definition of compatibility that underwrites an intuitive conception of defeat. We illustrate the nature of compatibility with the following lemma.

Lemma 1. 1. If $\Gamma \cup \Delta \succsim \Theta$ and $\Lambda \subset \Theta$, then $\Delta \succsim \Lambda$.

2. $\Gamma \cup \{A \wedge B\} \succsim \Theta$ iff $\Gamma \cup \{A, B\} \succsim \Theta$
3. $\Gamma \cup \{A \vee B\} \succsim \Theta$ iff $\Gamma \cup \{A\} \succsim \Theta$ or $\Gamma \cup \{B\} \succsim \Theta$
4. $\Gamma \cup \{\neg\neg A\} \succsim \Theta$ iff $\Gamma \cup \{A\} \succsim \Theta$

Proof. For the first sub-theorem, suppose for *reductio* that $\Delta \not\succsim \Lambda$. It would then follow that $\bigwedge \Delta \in \mathcal{D}(\Lambda)$. But then, by hypothesis, $\bigwedge \Delta \in \mathcal{D}(\Theta)$ and thus, against our assumption, $\Gamma \cup \Delta \not\succsim \Theta$. The remaining sub-theorems follow directly from Definitions 2, 3, and 4. ■

6.1 LK^Θ

Our notions of compatibility and defeat enable us to capture premise-consistency and defeasibility. To demonstrate this, we invite the reader to consider the rules for LK^Θ in Figure 1. These rules are designed to generate trees that preserve validity *downward* and undefeatedness *upward*. The most natural way to read the rules is bottom-up as follows: “If [the conclusion sequent] is undefeated, then so is/are [the premise sequent/s].” Alternatively, the rules may be read top-down as permitting the conclusion, given the premises, so long as the former is not defeated. On either reading, the rules are formulated to ensure that a cut-free derivation whose end-sequent is undefeated will contain only undefeated sequents throughout. This property is important both for the establishment of cut-eliminability and, more primitively, for the fact that proofs in a nonmonotonic system should not contain defeated sequents.

Definition 6 (Proof, Paraproof). For a rooted, finitely branching tree π whose nodes are sequents of LK^Θ (LEA), and which is recursively built up from axioms by means of the rules of LK^Θ (LEA), if each sequent in π is undefeated, then π is said to be a proof of LK^Θ (LEA), otherwise π is called a paraproof.

The axioms of LK^Θ come in two varieties. *Logical axioms* are the familiar ‘initial sequents’ of LK to which defeater sets are attached. *Proper axioms*, on the other hand, are sequents composed of nonempty, non-overlapping sets of atoms on the left and right of the turnstile.¹⁰ These axioms are intended to represent the non-logical, material inferences that figure in various types of scientific reasoning. Since such inferences are often the products of concrete empirical inquiries, we insist that they be introduced with non-empty *background sets* of (possibly complex) formulas that reflect the epistemic context of their use.

¹⁰It is well-known that the addition of non-tautological axioms to the system of classical logic will lead to inconsistency if those axioms are taken to be closed under universal substitution (US). Even if closure under US is abandoned for proper axioms, their addition to

Logical Axioms

$$\frac{}{\cdot \mid p \mid_{\Theta} p} \text{ log. ax.}$$

Proper Axioms

$$\frac{}{\Sigma \mid X \mid_{\Theta} Y} \quad \Sigma, X, Y \text{ are nonempty; } X \cap Y = \emptyset \quad \text{prop. ax.}$$

Cut Rule

$$\frac{\Sigma \mid \Gamma \mid_{\Theta} A, \Delta \quad \Sigma' \mid \Gamma', A \mid_{\Psi} \Delta'}{\Sigma', \Sigma \mid \Gamma', \Gamma \mid_{\Theta \cup \Psi} \Delta, \Delta'} \text{ cut}$$

Structural Rules

$$\frac{\Sigma \mid \Gamma \mid_{\Theta} \Delta}{\Sigma \mid \Gamma, A \mid_{\Theta} \Delta} \text{ LW}$$

$$\frac{\Sigma \mid \Gamma \mid_{\Theta} \Delta}{\Sigma \mid \Gamma \mid_{\Theta} \Delta, A} \text{ RW}$$

$$\frac{\Sigma \mid \Gamma \mid_{\Theta} \Delta}{\Sigma \mid \Gamma \mid_{\Theta \cup \Psi} \Delta} \text{ DE}$$

$$\frac{\Sigma \mid \Gamma \mid_{\Theta} \Delta}{\Sigma, A \mid \Gamma \mid_{\Theta} \Delta} \text{ BE}$$

Logical Rules

$$\frac{\Sigma \mid \Gamma, A, B \mid_{\Theta} \Delta}{\Sigma \mid \Gamma, A \wedge B \mid_{\Theta} \Delta} \wedge \vdash$$

$$\frac{\Sigma \mid \Gamma \mid_{\Theta} A, \Delta \quad \Sigma' \mid \Gamma' \mid_{\Psi} B, \Delta'}{\Sigma', \Sigma \mid \Gamma', \Gamma \mid_{\Theta \cup \Psi} A \wedge B, \Delta, \Delta'} \vdash \wedge$$

$$\frac{\Sigma \mid \Gamma, A \mid_{\Theta} \Delta \quad \Sigma' \mid \Gamma', B \mid_{\Psi} \Delta'}{\Sigma', \Sigma \mid \Gamma', \Gamma, A \vee B \mid_{\Theta \cup \Psi} \Delta, \Delta'} \vee \vdash^{\dagger}$$

$$\frac{\Sigma \mid \Gamma \mid_{\Theta} A, B, \Delta}{\Sigma \mid \Gamma \mid_{\Theta} A \vee B, \Delta} \vdash \vee$$

$$\frac{\Sigma \mid \Gamma \mid_{\Theta} A, \Delta}{\Sigma \mid \Gamma, \neg A \mid_{\Theta} \Delta} \neg \vdash$$

$$\frac{\Sigma \mid \Gamma, A \mid_{\Theta} \Delta}{\Sigma, A \mid \Gamma \mid_{\Theta} \neg A, \Delta} \vdash \neg$$

\dagger Provided that $\{A\} \lesssim \Theta$ and $\{B\} \lesssim \Psi$.

Figure 1: Rules for LK^{Θ}

In order to avoid having defeated axioms—i.e. non-starters—we must place certain constraints on the defeater sets of initial sequents.

Definition 7 (Constraints on Defeater Sets for Axioms). If $\cdot \mid p \mid_{\Theta} p$ is a logical axiom and $\Sigma \mid X \mid_{\Psi} Y$ is a proper axiom, then

- (i) Θ and Ψ are sets of literals.
- (ii) $p \notin \Theta$ and $\forall p \in X (p \notin \Psi)$
- (iii) $\Sigma \cap \Psi = \emptyset$

The first constraint restricts the defeater sets of axioms to a low level of formula complexity in order to limit the lacuna that occurs when a conjunction but neither of its conjuncts belongs to a defeater set. The second ensures that initial sequents are not defeated by their antecedents and, in the case of logical axioms, that equivalence among atoms is preserved. While this move protects reflexivity in LK^{Θ} , the rules of LE^{\blacktriangleright} will prevent this property from being transferred to properly explanatory arguments. Finally, the third constraint prevents proper axioms from being defeated by their initial background sets.

Before demonstrating that these constraints allow us to produce a nonmonotonic system with premise-consistent theorems, we pause to explain some of the more exotic features of LK^{Θ} . First, note that, with the exception of DE, single-premise rules have no effect on defeater sets, whereas, all the two-premise rules yield defeater sets in the conclusion that are the union of those in the premises. This fact guarantees that no information about potential defeaters is lost along a derivation.

Proposition 1. A. If π is a tree whose root is $\Sigma \mid \Gamma \mid_{\Theta'} \Delta$ and $\cdot \mid p \mid_{\Theta} p$ occurs in the leaves of π and $A \in \Theta$, then $A \in \Theta'$.

B. If π is a tree whose root is $\Sigma \mid \Gamma \mid_{\Theta'} \Delta$ and $\Sigma \mid X \mid_{\Theta} Y$ occurs in the leaves of π and $A \in \Theta$, then $A \in \Theta'$.

Proof. Follows directly from the rules of LK^{Θ} . ■

Exempting DE from this requirement is justified by the idea that reasoners ought to be able to add new, extra-logical information about defeaters to their arguments as it becomes available. The DE rule (which stands for *Defeater Expansion*) allows one to do so as long as it does not defeat the sequent in question.

Second, all of the rules either transfer or combine background sets from premises to conclusions, except for BE and $\vdash \neg$. The former (whose name abbreviates *Background Expansion*) permits the addition of arbitrary formulas to a sequent's background set. At one level, this rule simply captures the way new contextual information is added in the course of scientific reasoning. But at a deeper level, it attempts to represent the way reasoners might *probe* the defeasibility of an inference by discharging it in different contexts. In this sense, BE codifies certain patterns of *experimental* reasoning.

On the other hand $\vdash \neg$ adds the (active) antecedent of the premise to the background set of the conclusion. This behavior is of a piece with the explanation given for background sets above—they act as a kind of *record* of those formulas that have been shifted from the left to

a sequent system such as LK, threatens the cut-elimination theorem. Fortunately, **Piazza2016** have shown how to generate non-logical axiomatic extensions of classical propositional logic that admit cut elimination. While such extensions are obviously not complete—they are *post-compet*—axioms can be formulated so as to preserve consistency. The trick to doing so is to ensure that the empty sequent does belong to the set of proper axioms (See Theorem 3.7 in **Piazza2016**). Our proper axioms have been formulated in conformity with this constraint.

the right of a turnstile.¹¹ An informal interpretation of the rule (read upwards) can be given as follows: if one is entitled to $\neg A$ while A is in one's background set of entitlements, then one is entitled to whatever follows from A once it has been removed from that background set and entitlement to its negation has been renounced. The formulation of $\vdash \neg$ in this manner is critical to the upwards preservation of undefeatedness in cut-free proofs.

Proposition 2. Any cut-free paraproof in LK^\ominus is a proof if and only if its end-sequent is undefeated, i.e. undefeatedness is preserved upwards in cut-free proofs.

Proof. For all of the rules except $\vee \vdash$, the undefeatedness of the premises follows directly from that of the conclusion by way of Lemma 1.1, 1.2 or 1.4. In the case of $\vee \vdash$, it follows from Lemma 1.3 that it would be possible for the conclusion to be undefeated while exactly one of the premises is defeated. This possibility, however, is blocked by the proviso (\dagger) that restricts the rule's application to sequents whose active formulas are compatible with their respective defeater sets. ■

The preservation of undefeatedness in cut-free proofs is, in turn, critical to cut-elimination because by permuting cut upwards in derivations, the normalization procedure occasionally turns proofs into paraproof. Proposition 2, however, ensures that for any such paraproof, there is a cut-free proof of its end-sequent.

Lemma 2. Any sequent that is provable in LK^\ominus has a cut-free proof.

Proof. As noted above, the similarity between our LK^\ominus and Piazza et al.'s LK^S enables us to inherit the latter's cut-elimination theorem. We refrain from presenting the proof of the theorem here, but the interested reader is invited to see **Piazza2015** Theorem 19. ■

6.2 LE^\blacktriangleright and LEA

With the rules for LK^\ominus now in place, we can turn to the other part of our base system, LE^\blacktriangleright , where we introduce a class of consequence relations, denoted by $\vdash_{\ominus}^\blacktriangleright$, that represents explanatory arguments. We will call sequents constructed with this turnstile \blacktriangleright -sequents. Before discussing the rules for this system, we need a few preliminary concepts.

Definition 8 ($\neg(\Gamma)$). Let $\neg(\Gamma)$ stand for the following operation:

$$\neg(\Gamma) =_{df} \{\neg A : A \in \Gamma\}$$

The operation easily extends to sets of sets:

$$\neg(\mathbf{S}) =_{df} \{\neg(\Gamma) : \Gamma \in \mathbf{S}\}$$

Definition 9 ($\bar{\Gamma}$). Let $\bar{\Gamma}$ stand for the following set function:

$$\bar{\Gamma} =_{df} \neg(\Gamma) \cup \{A \wedge \neg A : A \in \mathcal{L}\}$$

By definition, $\mathcal{D}(\bar{\Gamma})$ contains all of the contradictory formulas in \mathcal{L} . This set function will thus be instrumental to securing premise-consistency for \blacktriangleright -sequents.

¹¹By this same reasoning, the *cut* rule also ought to add the active formula in its premises to the background set of the conclusion. However, the *cut* rule is exempted on the grounds that such a rule is just a statement about the conditions under which information may be *removed* from a proof.

¹²See **Piazza2015** for an example.

Definition 10 (Competitor Set, S^Δ). Let S^Δ stand for the following set:

$$S^\Delta =_{df} \{ \Omega_i : \Sigma_i \mid \Omega_i \xrightarrow[\Psi_i \cup \overline{\Omega}_i \cup \Delta]{LK^\Theta} \Delta \}$$

where $i \in \mathbb{N}$ and LK^Θ above a turnstile indicates that the sequent is a theorem of LK^Θ .

Despite its complex appearance, S^Δ is nothing more than the set containing the antecedents of all those nontrivial (i.e. premise-consistent and irreflexive) theorems of LK^Θ that have Δ as its conclusion. For ease of reference, the definition indexes background sets and defeater sets by the antecedents of the relevant sequent. This definition is crucial to securing sturdiness.

Definition 11 (Disjunction Deletion, $\lfloor S \rfloor$).

$$\lfloor S \rfloor =_{df} \{ \Delta \setminus \{A \vee B\} : \Delta \in S \text{ and } A \in \bigcup S \}$$

The set denoted by $\lfloor S \rfloor$ is the result of deleting from the members of S any disjunction one of whose disjuncts belongs to a member of S .

Definition 12 ($\Delta \parallel \Gamma$).

$$\Delta \parallel \Gamma =_{df} \begin{cases} \Delta & \text{if } \Delta \subseteq \Gamma \\ \Delta \setminus \Gamma & \text{otherwise} \end{cases}$$

We extend the definition to sets of sets as follows:

$$S \parallel \Gamma =_{df} \{ \Delta \parallel \Gamma : \Delta \in S \}$$

The set denoted by $S \parallel \Gamma$ consists of those members of S that are subsets of Γ and the set-complements relative to Γ of those members that are not subsets of Γ . As we shall see, Definitions 11 and 12 are needed to secure minimality.

We will now consider the rules for LE^\blacktriangleright , presented in Figure 2. Since the structural rules for LE^\blacktriangleright are just the DE and BE rules, *modulo* the \blacktriangleright turnstile, we begin with what we have called the *Mixed Rules*. These are both the most unusual—as they contain two different types consequence relations—and the most important for the purposes of representing explanatory arguments.

As mentioned above, *sturdiness* is our proposal for how to understand the property of *stability* associated with explanations. In slogan form, inferences are sturdy just in case they succeed where all others fail. To say that one inference succeeds when another fails, we can imagine the following procedure:

Step 1: Line up all of the nontrivial inferences that have the explanandum, B , as their conclusion. For each of these inferences, all other nontrivial inferences leading to B are its “competitors.”

Step 2: For each A that has B as a nontrivial consequence, suppose that all of A ’s competitors’ premises are false.

Step 3: If the falsehood of any of these competitors defeats the inference from A to B , then the latter is not sturdy; otherwise, it is sturdy.

Our *sturdy* rule aims to formalize this procedure. The first step is represented by the fact that the succedent of the premise is Δ and its background set includes a set obtained from S^Δ .

Structural Rules

$$\frac{\Sigma \mid \Gamma \mid_{\Theta}^{\blacktriangleright} \Delta}{\Sigma \mid \Gamma \mid_{\Theta \cup \Psi}^{\blacktriangleright} \Delta} \text{DE}^{\blacktriangleright}$$

$$\frac{\Sigma \mid \Gamma \mid_{\Theta}^{\blacktriangleright} \Delta}{\Sigma, A \mid \Gamma \mid_{\Theta}^{\blacktriangleright} \Delta} \text{BE}^{\blacktriangleright}$$

Mixed Rules

$$\frac{\bigcup \neg([\mathbf{S}^{\Delta}] \parallel \Gamma), \Sigma \mid \Gamma \mid_{\Theta \cup \bar{\Gamma} \cup \Delta}^{\blacktriangleright} \Delta}{\Sigma \mid \Gamma \mid_{\Theta \cup \bar{\Gamma} \cup \Delta}^{\blacktriangleright} \Delta} \text{sturdy}^{\ddagger}$$

$$\frac{\Sigma \mid \Gamma, A \mid_{\Theta}^{\blacktriangleright} \Delta \quad \Sigma' \mid \Gamma' \mid_{\Psi \cup \bar{\Gamma}' \cup \Delta}^{\blacktriangleright} \Delta}{\Sigma', \Sigma, A \mid \Gamma', \Gamma \mid_{\Theta \cup \Psi \cup \bar{\Gamma}' \cup \Delta}^{\blacktriangleright} A} \text{abduct.}$$

Logical Rules

$$\frac{\Sigma \mid \Gamma, A, B \mid_{\Theta}^{\blacktriangleright} \Delta}{\Sigma \mid \Gamma, A \wedge B \mid_{\Theta}^{\blacktriangleright} \Delta} \wedge \vdash^{\blacktriangleright}$$

$$\frac{\Sigma \mid \Gamma \mid_{\Theta}^{\blacktriangleright} A, \Delta \quad \Sigma' \mid \Gamma \mid_{\Psi}^{\blacktriangleright} B, \Delta}{\Sigma', \Sigma \mid \Gamma \mid_{\Theta \cup \Psi}^{\blacktriangleright} A \wedge B, \Delta} \vdash^{\blacktriangleright} \wedge$$

\ddagger Provided that (i) all sets are nonempty, (ii) $\Sigma \subseteq \Sigma_i$, (iii) $\Gamma \cap \Sigma_i = \emptyset$, (iv) $\Omega_i \cap \Sigma = \emptyset$.

\star Provided that $\Gamma \cap \Gamma' \cap \{A\} = \emptyset$.

Figure 2: Rules for $\text{LE}^{\blacktriangleright}$

The second step is captured by the fact that $\bigcup \neg([\mathbf{S}^{\Delta}] \parallel \Gamma)$ appears in the background set of the premise. Roughly put, this set contains the negations of all the antecedents of nontrivial inferences whose succedent is Δ , with the caveat that the members of Γ are removed from any antecedent that is a superset of Γ . Finally, the third step is reached when the premise sequent carries down into the conclusion where it appears with the explanatory-argument-denoting turnstile, i.e. $\mid_{\Theta}^{\blacktriangleright}$.

By restricting premise sequents to those whose defeater sets contain $\bar{\Gamma}$ —where Γ is the antecedent—the *sturdy* rule ensures that the antecedents of \blacktriangleright -sequents are consistent and thereby enables these sequents to represent premise-consistent inferences.

Lemma 3. The sequent $\Sigma \mid \Gamma, A, \neg A \mid_{\Theta}^{\blacktriangleright} \Delta$ is not a theorem of LEA.

Proof. Suppose for *reductio* that $\Sigma \mid \Gamma, A, \neg A \mid_{\Theta}^{\blacktriangleright} \Delta$ is a theorem and hence is undefeated. It must be derived via an application of *sturdy* whose premise, if we omit the background set, is $\Gamma, A, \neg A \mid_{\Theta \cup \{A, \neg A\} \cup \Delta}^{\blacktriangleright} \Delta$. By Definition 9, $\bigwedge \{A, \neg A\} \in \mathcal{D}(\overline{\{A, \neg A\}})$, since $A \wedge \neg A \in \{B \wedge \neg B : B \in \mathcal{L}\}$. Thus, *contra* our supposition, $\Sigma \mid \Gamma, A, \neg A \mid_{\Theta}^{\blacktriangleright} \Delta$ is defeated. ■

Lemma 4. The sequent $\Sigma \mid \Gamma, A \wedge \neg A \mid_{\Theta}^{\blacktriangleright} \Delta$ is not a theorem of LEA.

Proof. Same as for Lemma 3. ■

The *sturdy* rule also prevents \blacktriangleright -sequents from being reflexive. It does so by restricting the defeater sets of premise sequents to those which contain the succedent. We can now show that both partial and complete self-explanations will be prohibited from $\text{LE}^{\blacktriangleright}$.

Lemma 5. The sequent $\Sigma \mid \Gamma, A \mid_{\Theta}^{\blacktriangleright} A, \Delta$ is not a theorem of LEA.

Proof. Suppose for *reductio* that $\Sigma \mid \Gamma, A \mid_{\Theta}^{\triangleright} A, \Delta$ is a theorem and hence is undefeated. It must be derived via an application of *sturdy* whose premise, if we omit the background set, is $\Gamma, A \mid_{\Theta \cup \{A\} \cup \Delta}^{\triangleright} A, \Delta$. But, since (by Definition 3) $\bigwedge(\Gamma \cup \{A\}) \in \mathcal{D}(A)$, this sequent is defeated, contradicting our supposition. ■

Lemma 6. The sequent $\Sigma \mid \Gamma, A \wedge B \mid_{\Theta}^{\triangleright} A, \Delta$ is not a theorem of LEA.

Proof. Same as for Lemma 5. ■

Furthermore, neither premise-consistency, nor irreflexivity prevent instances of **MP** from standing as candidates for explanatory arguments, i.e. as premises of *sturdy*. For the remaining proofs, we omit arbitrary background sets (i.e. Σ) whenever possible.

Proposition 3. If the sequent $\Sigma \mid A, \neg A \vee B \mid_{\Theta}^{\triangleright} B$ is a theorem of LEA, then so is $\Sigma \mid A, \neg A \vee B \mid_{\Theta \cup \{A, \neg A \vee B\} \cup \{B\}}^{\triangleright} B$.

Proof. Suppose for *reductio* that there is no proof π of $A, \neg A \vee B \mid_{\Theta \cup \{A, \neg A \vee B\} \cup \{B\}}^{\triangleright} B$ from $A, \neg A \vee B \mid_{\Theta}^{\triangleright} B$. It follows that either π is not a paraproof or $A, \neg A \vee B \mid_{\Theta \cup \{A, \neg A \vee B\} \cup \{B\}}^{\triangleright} B$ is defeated. But $A, \neg A \vee B \mid_{\Theta \cup \{A, \neg A \vee B\} \cup \{B\}}^{\triangleright} B$ follows from $A, \neg A \vee B \mid_{\Theta}^{\triangleright} B$ by a single application of DE, and thus π is at least a paraproof. If $A, \neg A \vee B \mid_{\Theta \cup \{A, \neg A \vee B\} \cup \{B\}}^{\triangleright} B$ is defeated, then $\bigwedge\{A, \neg A \vee B\} \in \mathcal{D}(\Theta \cup \overline{\{A, \neg A \vee B\}} \cup \{B\})$. But $A \wedge (\neg A \vee B) \notin \mathcal{D}(\{B\})$, and by Definition 3 and 9 it follows that $A \wedge (\neg A \vee B) \notin \mathcal{D}(\overline{\{A, \neg A \vee B\}})$. Therefore, it must be the case that $A \wedge (\neg A \vee B) \in \mathcal{D}(\Theta)$. But by hypothesis $A, \neg A \vee B \mid_{\Theta}^{\triangleright} B$ is undefeated. Thus, *contra* our supposition, there is a proof of $A, \neg A \vee B \mid_{\Theta \cup \{A, \neg A \vee B\} \cup \{B\}}^{\triangleright} B$. ■

Less obvious than the achievement of irreflexivity is the fact that *sturdy* also ensures that \triangleright -sequents are minimal in the sense articulated above by conditions (M1)–(M3).

Lemma 7. If the sequent $\Sigma \mid \Gamma, A \mid_{\Theta}^{\triangleright} \Delta$ is a theorem of LEA, then $\Sigma \mid \Gamma \mid_{\Theta}^{\triangleright} \Delta$ is not.

Proof. We proceed by proving the contrapositive. Suppose for conditional proof that $\Gamma \mid_{\Theta}^{\triangleright} \Delta$ is a theorem of LEA. It follows from *sturdy* that $\Gamma \mid_{\Theta}^{\triangleright} \Delta$ is a theorem and from Definition 10 that $\Gamma \in \lfloor \mathbf{S}^{\Delta} \rfloor$. Next suppose for *reductio* that $\Gamma, A \mid_{\Theta}^{\triangleright} \Delta$ is also a theorem. From *sturdy* it follows that $\bigcup \neg(\lfloor \mathbf{S}^{\Delta} \rfloor \parallel \Gamma \cup \{A\}) \mid \Gamma, A \mid_{\Theta \cup \Gamma \cup \{A\} \cup \Delta}^{\triangleright} \Delta$ is a theorem. Since $\Gamma \subseteq \Gamma \cup \{A\}$, it follows from Definition 10 and 12 that $\neg(\Gamma) \subset \bigcup \neg(\lfloor \mathbf{S}^{\Delta} \rfloor \parallel \Gamma \cup \{A\})$. But from Definition 3 it follows that $\bigwedge \neg(\Gamma) \in \mathcal{D}(\overline{\Gamma \cup \{A\}})$, so, according to Definition 4 $\bigcup \neg(\lfloor \mathbf{S}^{\Delta} \rfloor \parallel \Gamma \cup \{A\}) \mid \Gamma, A \mid_{\Theta \cup \Gamma \cup \{A\} \cup \Delta}^{\triangleright} \Delta$ is defeated, contradicting our supposition. Thus, if $\Gamma \mid_{\Theta}^{\triangleright} \Delta$ is a theorem, then $\Gamma, A \mid_{\Theta}^{\triangleright} \Delta$ is not. It follows by contraposition and double negation that if the sequent $\Gamma, A \mid_{\Theta}^{\triangleright} \Delta$ is a theorem of LEA, then $\Gamma \mid_{\Theta}^{\triangleright} \Delta$ is not. ■

Lemma 8. If the sequent $\Sigma \mid A, B, C \mid_{\Theta}^{\triangleright} \Delta$ is a theorem of LEA, then $\Sigma \mid B \wedge C \mid_{\Theta}^{\triangleright} \Delta$ is not.

Proof. It suffices to show that $\bigcup \neg([\mathbf{S}^\Delta] \parallel B \wedge C), \Sigma \mid B \wedge C \frac{\triangleright}{\Theta} \Delta$ is defeated when $\{A, B, C\} \in [\mathbf{S}^\Delta]$. First note that $\bigwedge \neg(\{A, B, C\}) = \neg A \wedge \neg B \wedge \neg C$ and that $\neg B \vee \neg C \in \mathcal{D}(\overline{B \wedge C})$ and thus, by the \mathcal{D}_\vee -rule, $\{\neg B, \neg C\} \subset \mathcal{D}(\overline{B \wedge C})$. It follows by the \mathcal{D}_\wedge -rule that $\neg A \wedge \neg B \wedge \neg C \in \mathcal{D}(\overline{B \wedge C})$, and thus $\bigwedge \neg(\{A, B, C\}) \in \mathcal{D}(\overline{B \wedge C})$. The rest of the proof mirrors that of Lemma 7 and is left as an exercise for the reader. ■

Lemma 9. If the sequents $\Sigma' \mid A \frac{}{\Theta'} \Delta$ and $\Sigma'' \mid B \frac{}{\Theta''} \Delta$ are theorems of LEA then $\Sigma \mid A \vee B \frac{\triangleright}{\Theta} \Delta$ is not.

Proof. Suppose for *reductio* that $\Sigma, \mid A \vee B \frac{\triangleright}{\Theta} \Delta$ is a theorem. It follows by hypothesis and *sturdy* that $\Sigma, \bigcup \neg(\{\{A, B\}\} \parallel \{A \vee B\}) \mid A \vee B \frac{}{\Theta \cup \overline{A \vee B} \cup \Delta} \Delta$ is a theorem. $\{\{A, B\}\} \parallel \{A \vee B\} = \{\{A, B\}\}$ and by Definition 8, $\bigcup \neg(\{\{A, B\}\}) = \{\neg A, \neg B\}$. Since $\bigwedge \{\neg A, \neg B\} = \neg A \wedge \neg B$, we test for compatibility by determining whether $\neg A \wedge \neg B \in \mathcal{D}(\overline{A \vee B})$. The latter is indeed obtained by Definitions 3 and 9. It follows that $\Sigma, \bigcup \neg(\{\{A, B\}\} \parallel \{A \vee B\}) \mid A \vee B \frac{}{\Theta \cup \overline{A \vee B} \cup \Delta} \Delta$ is defeated and thus that $\Sigma \mid A \vee B \frac{\triangleright}{\Theta} \Delta$ is not a theorem. ■

Remark 2. The consequent of Lemma 9 also holds under the hypothesis that $\Sigma' \mid A, B \frac{}{\Theta'} \Delta$ is a theorem.

Lemma 9 holds that a disjunctive candidate is never sturdy when it must compete against both of its disjuncts. Unfortunately, without a special constraint, disjunctive competitors would also block their disjuncts from obtaining sturdiness.

Fact 1. $\bigcup \neg(\{\{A \vee B\}\} \parallel \{A\}) \in \mathcal{D}(\neg A)$.

Proof. It suffices to show that $\bigcup \neg(\{\{A \vee B\}\} \parallel \{A\}) = \{\neg A \wedge \neg B\}$ and $\neg A \wedge \neg B \in \mathcal{D}(\neg A)$. ■

In order to prevent this unwanted result, we deploy Disjunction Deletion (Definition 11) in the formulation of *sturdy*. Thus, as the following theorem states, a candidate explanans is never forced to compete against a set that contains its disjunction with an arbitrary formula.

Proposition 4. If the sequent $\bigcup \neg([\mathbf{S}^\Delta] \parallel \Gamma \cup \{A\}), \Sigma \mid \Gamma, A \frac{}{\Theta \cup \overline{\Gamma \cup \{A\}} \cup \Delta} \Delta$ is the premise in an application of *sturdy*, then $A \vee B \notin \bigcup \{[\mathbf{S}^\Delta] \parallel \Gamma \cup \{A\}\}$ for any formula B .

Proof. It suffices to note that according to Definition 11, if $A \in \bigcup \mathbf{S}^\Delta$, then $A \vee B \notin \bigcup [\mathbf{S}^\Delta]$ for any arbitrary formula B . ■

The provisos on *sturdy* are intended to prevent applications of the rule in cases where there is no genuine comparison between a candidate explanans and its competitors—e.g. if a candidate explanans were to be smuggled into the background set of a competitor ($\Gamma \subseteq \Sigma_i$) or vice versa ($\Omega_i \subseteq \Sigma$). Similarly, the requirement that the background set of the candidate explanans form a subset of those of its competitors ($\Sigma \subseteq \Sigma_i$) provides a common set of assumptions against which comparisons can be made. Thus, these provisos provide a level playing field on which candidate explanans may compete for sturdiness.

In combination with Definition 10, this last constraint enables the set of competitors to be culled. For instance, one can prevent an antecedent, Ω_i , from belonging to \mathbf{S}^Δ by constructing a background set for the candidate that includes information that defeats the inference from Ω_i to Δ . This procedure of culling the competitor set describes how a reasoner goes about holding certain pieces of information (e.g. actual causes) “fixed”.

We turn now to the *abduct.* rule. As the name suggests, this rule is intended to capture the *detachability* of the premises of explanatory arguments. Roughly put, *abduct.* says that if Γ, A is an explanatory argument for Δ , and Δ is the non-trivial consequence of Γ' , then together, Γ and Γ' license the inference to A . Since A only forms part of the explanatory argument for Δ , we ought to read *abduct.* as licensing the detachment of a *partial* explanation of Δ . The formulation of the rule thus makes detachability a manifest property of \blacktriangleright -sequents.

There are, however, some peculiarities to the rule that deserve discussion. First, the explanandum (Δ) disappears from the conclusion, leaving only the partial explanans. This feature accords with our desire to present IBE in the strongest form possible. If IBE only licensed inferences to best explanations *or* their explananda, its legitimacy would hardly have roused debate—though its utility might have. Unfortunately, the absence of the explananda in the succedent of *abduct.*'s conclusion means that information is lost in any derivation that contains an application of the rule. The effect, like that of proofs that employ cut, is the failure of analyticity—i.e. some derivations will contain formulas that are not subformulas of those in the end-sequent. The loss of the subformula property is not all that surprising given the standard characterization of IBE as a form of ampliative inference. We are reassured by the fact that despite this loss, the cut-elimination theorem holds for LEA (Theorem 1).

Second, the defeater set attached to the second premise indicates that the inference that entitles us to the explanandum must be non-trivial. This restriction is justified on the grounds that a tautology should never count as evidence for an explanandum's obtaining. Third, the proviso on *abduct.* prevents the antecedents of the premises from overlapping. This constraint follows from the idea that abductive inferences are only licensed when one has evidence for the explanandum that is independent of the explanans. Note that this means that the second premise of *abduct.* will contain a sequent whose antecedent may have 'competed' with the explanans for sturdiness. This is as it should be, since the competitors include not only potential explanations, but also non-explanatory, evidential inferences. ^{JM}[Should I elaborate on this?](#) Finally, in addition to appearing in the succedent of the conclusion, the (partial) explanans also appears in the background set. This feature is consistent with our understanding of these sets as keeping track of formulas that have shifted from the left to the right of the turnstile.

The logical rules of $\text{LE}^\blacktriangleright$ are just the left and right rules for conjunction in LK^\ominus , *modulo* \blacktriangleright -sequents. We have chosen to restrict the logical operations permitted on \blacktriangleright -sequents to these out of an abundance of caution. To our ears, inferences from conjunctive explanantia and to conjunctive explananda sound far more natural than those involving disjunctive or negated explanantia/explananda. ^{JM}[Should we say more?](#)

Now that we have in place its constituent systems, let us reflect on the global properties of LEA. Perhaps most striking is the absence of classical structural rules for \blacktriangleright -sequents. Neither weakening nor cut is permitted. If the antecedents of \blacktriangleright -sequents could be arbitrarily weakened, then the minimality condition, which *sturdy* secures, would be compromised. Right weakening, on the other hand, does not directly conflict with the properties of explanatory arguments; rather, we deny it because the weakening of explananda by arbitrary disjuncts appears to us as both unnatural and unreasonable. Finally, the absence of cut follows from the fact that our target vocabulary is that which expresses the notion of *immediate* explanation.

Proposition 5. It is possible that $\Sigma \mid \Gamma \frac{\blacktriangleright}{\ominus} A, \Delta$ and $\Sigma' \mid \Gamma', A \frac{\blacktriangleright}{\Psi} \Delta'$ are theorems of LEA, but $\Sigma', \Sigma \mid \Gamma', \Gamma \frac{\blacktriangleright}{\ominus \cup \Psi} \Delta, \Delta'$ is not.

Proof. Follows from the fact that there is no cut rule for \blacktriangleright -sequents. ■

While there is no $\text{cut}_{\blacktriangleright}$ rule, proofs in LEA contain an application of cut that was not available in LK^\ominus , namely, one where the cut formula appears in the succedent that follows the application of abduct . Fortunately, the cut-elimination theorem can be extended to cover these cases.

Theorem 1. Any sequent which is provable in LEA has a cut-free proof.

Proof. Lemma 2 gives us cut-elimination for the proofs in LK^\ominus . The following reduction covers the one application of cut that appears in proofs of LEA that does not appear in proofs of LK^\ominus .

$$\begin{array}{c}
\frac{\frac{\bigcup \neg([\mathbf{S}^\Delta] \parallel \Gamma), \Sigma \mid \Gamma, A \mid_{\Theta \bar{\Gamma} \cup \Delta} \Delta}{\Sigma \mid \Gamma, A \mid_{\Theta \bar{\Gamma} \cup \Delta} \Delta} \text{ sturdy} \quad \frac{\pi \vdots \Sigma'' \mid \Gamma'' \mid_{\Psi \bar{\Gamma}'' \cup \Delta} \Delta}{\Sigma'', \Sigma, A \mid \Gamma'', \Gamma \mid_{\Theta \bar{\Gamma} \cup \Delta \cup \Psi \bar{\Gamma}''} A} \text{ abduct.}}{\frac{\Sigma \mid \Gamma, A \mid_{\Theta} \Delta}{\Sigma'', \Sigma, A \mid \Gamma'', \Gamma \mid_{\Theta \bar{\Gamma} \cup \Delta \cup \Psi \bar{\Gamma}''} \Delta} \text{ cut}} \\
\\
\downarrow \\
\begin{array}{c}
\pi \vdots \\
\frac{\Sigma \mid \Gamma'' \mid_{\Psi} \Delta}{\Sigma \mid \Gamma'', \Gamma \mid_{\Psi} \Delta} \text{ LW} \\
\frac{\Sigma \mid \Gamma'', \Gamma \mid_{\Psi} \Delta}{\Sigma, A \mid \Gamma'', \Gamma \mid_{\Psi} \Delta} \text{ BE} \\
\frac{\Sigma, A \mid \Gamma'', \Gamma \mid_{\Psi} \Delta}{\Sigma'', \Sigma, A \mid \Gamma'', \Gamma \mid_{\Psi} \Delta} \text{ BE} \\
\frac{\Sigma'', \Sigma, A \mid \Gamma'', \Gamma \mid_{\Psi} \Delta}{\Sigma'', \Sigma, A \mid \Gamma'', \Gamma \mid_{\Delta \cup \Psi} \Delta} \text{ DE} \\
\frac{\Sigma'', \Sigma, A \mid \Gamma'', \Gamma \mid_{\Delta \cup \Psi} \Delta}{\Sigma'', \Sigma, A \mid \Gamma'', \Gamma \mid_{\bar{\Gamma} \cup \Delta \cup \Psi} \Delta} \text{ DE} \\
\frac{\Sigma'', \Sigma, A \mid \Gamma'', \Gamma \mid_{\bar{\Gamma} \cup \Delta \cup \Psi} \Delta}{\Sigma'', \Sigma, A \mid \Gamma'', \Gamma \mid_{\Theta \bar{\Gamma} \cup \Delta \cup \Psi} \Delta} \text{ DE} \\
\frac{\Sigma'', \Sigma, A \mid \Gamma'', \Gamma \mid_{\Theta \bar{\Gamma} \cup \Delta \cup \Psi} \Delta}{\Sigma'', \Sigma, A \mid \Gamma'', \Gamma \mid_{\Theta \bar{\Gamma} \cup \Delta \cup \Psi \bar{\Gamma}''} \Delta} \text{ DE}
\end{array}
\end{array}$$

■

There are two facts about the proof of Theorem 1 that are of particular significance. First, in contrast with standard normalization procedures for LK , the normalized proof above does not permute the application of cut upwards. Indeed, there is no application of cut whatsoever. Thus, unlike the normalization that secures cut-elimination in LK^\ominus , there is no need here to guarantee cut-free proofs of the end-sequents of cut-laden paraproofs. The relaxation of this demand is a welcome result, not least because the *sturdy*-rule fails to preserve undefeatedness upwards. (We discuss this point in the conclusion).

Second, there are no \blacktriangleright -sequents in the normalized proof. While applications of cut in cut-free calculi always involve a ‘detour’ through unnecessary steps, in this case the application of cut to the conclusion of abduct . renders the prior application of the *sturdy*-rule, and hence the deployment of the \blacktriangleright -subsystem, particularly gratuitous. Since the normalized proof is a

non-branching tree whose leaf is a non-trivial inference—indeed the very inference whose role in *abduct.* is to provide non-explanatory evidence that the ‘explanandum’ obtains—we have in the reduced proof an instance of reasoning that proceeds through explanatory arguments with no epistemic gain. ^{JM}Are there broader implications that could drawn out here?

Finally, we can now see that the theorems of LEA constructed with $\vdash_{\blacktriangleright}$ exhibit all of the properties associated with explanatory arguments.

Theorem 2. The \blacktriangleright -sequents that are theorems of LEA are (1) defeasible, (2) minimal, (3) premise-consistent, (4) irreflexive, (5) sturdy, and (6) detachable.

Proof. (1) Defeasibility follows from Definition 5.

(2) Minimality follows from Lemmas 7, 8, and 9.

(3) Premise-consistency follows from Lemma 3 and Corollary 4.

(4) Irreflexivity follows from Lemma 5 and Corollary 6.

(5) Sturdiness follows from the *sturdy* rule.

(6) Detachability follows from the *abduct.* rule. ■

6.3 The extension LEA^+

We shall now demonstrate how the system LEA and the language \mathcal{L} over which it is defined may be extended to include an object-language expression for *best explains why*. We begin with the syntax of the extended language $\mathcal{L}_{\blacktriangleright}$.

Definition 13 (Syntax of $\mathcal{L}_{\blacktriangleright}$).

- (1) If $A \in \mathcal{L}$ then $A \in \mathcal{L}_{\blacktriangleright}$.
- (2) If $A, B \in \mathcal{L}$ then $A \blacktriangleright B \in \mathcal{L}_{\blacktriangleright}$.

The expressions ‘ $A \blacktriangleright B$ ’ is intended to be read as ‘ A best explains why B .’ Note that the syntactic definition for \blacktriangleright is not recursive with respect to $\mathcal{L}_{\blacktriangleright}$. Consequently, the operator \blacktriangleright is non-iterative. We impose this syntactic constraint on the grounds that the “... best explains why...” locution does not appear to iterate in natural languages—at least not in English.

The rules in Figure 3 define the calculus LEA^+ over $\mathcal{L}_{\blacktriangleright}$. While the rules of LK^{Θ} apply to all the formulas in $\mathcal{L}_{\blacktriangleright}$, the rules for $\text{LE}^{\blacktriangleright}$ are restricted to the fragment $\mathcal{L} \cap \mathcal{L}_{\blacktriangleright}$. This restriction is intended to prevent the generation of ill-formed formulas along a derivation, e.g. $A \blacktriangleright (B \blacktriangleright C)$.

As promised, the extension LEA^+ provides introduction (right) and elimination (left) rules for the *best-explains-why* operator, i.e. \blacktriangleright . The $\blacktriangleright \vdash$ rule ought to be familiar—it is essentially the *abduct.* rule with the abducible formula joined to a member of the succedent of the second premise (i.e. B) by the \blacktriangleright -connective and appended to the antecedent of the conclusion. In fact, the $\blacktriangleright \vdash$ rule is derivable from *abduct.* and LW.

Proposition 6. $\blacktriangleright \vdash$ is a derivable rule in LEA^+ .

Proof.

$$\frac{\frac{\Sigma \mid \Gamma, A \mid_{\Theta}^{\triangleright} \Delta \quad \Sigma' \mid \Gamma' \mid_{\Psi \cup \Gamma' \cup \Delta}^{\overline{}} \Delta}{\Sigma', \Sigma, A \mid \Gamma', \Gamma \mid_{\Theta \cup \Psi \cup \Gamma' \cup \Delta}^{\overline{}} A} \text{abduct.}}{\Sigma', \Sigma, A \mid \Gamma', \Gamma, A \triangleright B \mid_{\Theta \cup \Psi \cup \Gamma' \cup \Delta}^{\overline{}} A} \text{LW}$$

■

^{JM}Is this proposition worth having?

The $\vdash \triangleright$ rule, on the other hand, represents something quite novel. While it resembles the right rule for \rightarrow in LK, it is distinguished by two features. First, the active formula in the antecedent of the premise occurs in the background set of the conclusion. This peculiarity is justified by the need to preserve undefeatedness upwards, much in the way that $\neg \vdash$ does. Second, while the premise is a \triangleright -sequent, the conclusion is not. This feature captures the sense in which formulas whose main operator is \triangleright are *explicitly* explanatory claims. As such, these claims can enter into reasoning patterns that do not consist in the making of explanatory arguments, and thus they belong to the class of sequents whose turnstile is unadorned by \triangleright .

We are now in a position to make good on our promise to provide an expressivist treatment of explanatory vocabulary. Since the deduction theorem serves as the model for logical expressivist theses, it is incumbent upon us to show that a similar theorem holds in LEA^+ . In order to do so, we must prove the invertibility of $\vdash \triangleright$. As we noted above, the fact that latter invokes two distinct classes of consequence relations means that, upon proof of invertibility, the resultant theorem will say something different. Namely, that an expression used in one logic (that of the unadorned turnstile, $\mid_{\Theta}^{}$) encodes the rules of another logic ($\mid_{\Theta}^{\triangleright}$). To distinguish this claim from the standard deduction theorem, we refer to it as a *quasi-deduction theorem*.

Theorem 3 (Quasi-Deduction Theorem). $\Sigma \mid \Gamma, A \mid_{\Theta}^{\triangleright} B, \Delta$ is provable in LEA^+ if and only if $\Sigma, A \mid \Gamma \mid_{\Theta}^{} A \triangleright B, \Delta$ is.

Proof. (\Rightarrow) Follows from $\vdash \triangleright$.

(\Leftarrow) By induction on proof-height. **Base Case:** If $\Sigma, A \mid \Gamma \mid_{\Theta}^{} A \triangleright B, \Delta$ is an axiom, then, since $A \triangleright B$ is not atomic, $\Sigma \mid \Gamma, A \mid_{\Theta}^{\triangleright} B, \Delta$ would have to be an axiom. However, there are no axioms with \triangleright -sequents. Rather, all such sequents are derived via *sturdy*. Thus, $\Sigma \mid \Gamma, A \mid_{\Theta}^{\triangleright} B, \Delta$ follows by *sturdy*.

Inductive Step: Assume inversion up to height n and let $\Sigma, A \mid \Gamma \mid_{\Theta}^{} A \triangleright B, \Delta$ be the root of a proof of height $n+1$. There are two cases:

Case 1: If $A \triangleright B$ is not principal in the last rule, it has one or two premises, $\Sigma', [A] \mid \Gamma' \mid_{\Theta}^{} A \triangleright B, \Delta'$ and $\Sigma'', [A] \mid \Gamma'' \mid_{\Theta}^{} A \triangleright B, \Delta''$ of derivation height $\leq n$, where $[A]$ indicates that A is in at least one of the premises (if there are two). By inductive hypothesis $\Sigma' \mid \Gamma', A \mid_{\Theta}^{\triangleright} B, \Delta'$ and $\Sigma'' \mid \Gamma'', A \mid_{\Theta}^{\triangleright} B, \Delta''$ are obtained by a proof of height n . Now apply the last rule to obtain $\Sigma \mid \Gamma, A \mid_{\Theta}^{\triangleright} B, \Delta$ with a proof of height $\leq n+1$.

Case 2: If $A \triangleright B$ is principal in the last rule, the premise $\Sigma \mid \Gamma, A \mid_{\Theta}^{\triangleright} B, \Delta$ has a proof of height $\leq n$. ■

^{JM}Discuss quasi-deduction theorem and show how it captures the spirit of logical expressivism.

We can now show that LEA^+ is a conservative extension of LEA. First, we prove that the cut elimination theorem holds for LEA^+ .

The rules of LK^Θ apply to all the formulas in $\mathcal{L}_\blacktriangleright$.

The rules of LE^\blacktriangleright apply to all those formulas in the fragment $\mathcal{L} \cap \mathcal{L}_\blacktriangleright$.

Rules for \blacktriangleright

$$\frac{\Sigma \mid \Gamma, A \mid_{\Theta} B, \Delta \quad \Sigma' \mid \Gamma' \mid_{\Psi \cup \Gamma' \cup \Delta} B, \Delta}{\Sigma', \Sigma, A \mid \Gamma', \Gamma, A \blacktriangleright B \mid_{\Theta \cup \Psi \cup \Gamma' \cup \Delta} A} \blacktriangleright \vdash^* \quad \frac{\Sigma \mid \Gamma, A \mid_{\Theta} B, \Delta}{\Sigma, A \mid \Gamma \mid_{\Theta} A \blacktriangleright B, \Delta} \vdash \blacktriangleright$$

* Provided that $\Gamma \cap \Gamma' \cap \{A\} = \emptyset$.

Figure 3: Rules of LEA^+ .

Theorem 4. Any sequent provable in LEA^+ has a cut-free proof.

Proof. Since Theorem 2 establishes cut elimination for LEA, we need to show that every application of cut that occurs in proofs of LEA^+ but not in proofs of LEA is eliminable. There is only one such application, namely, that which cuts the formula $A \blacktriangleright B$ from the conclusions of $\blacktriangleright \vdash$ and $\vdash \blacktriangleright$. The following reduction covers this application.

$$\begin{array}{c} \begin{array}{ccc} \pi_1 & & \pi_2 \\ \vdots & & \vdots \\ \Sigma \mid \Gamma, A \mid_{\Theta} B, \Delta & \quad \quad & \Sigma' \mid \Gamma' \mid_{\Psi \cup \Gamma' \cup \Delta} B, \Delta \\ \hline \Sigma', \Sigma, A \mid \Gamma', \Gamma, A \blacktriangleright B \mid_{\Theta \cup \Psi \cup \Gamma' \cup \Delta} A & \blacktriangleright \vdash & \frac{\Sigma \mid \Gamma, A \mid_{\Theta} B, \Delta}{\Sigma, A \mid \Gamma \mid_{\Theta} A \blacktriangleright B, \Delta} \vdash \blacktriangleright \\ \hline \Sigma', \Sigma, A \mid \Gamma', \Gamma \mid_{\Theta \cup \Psi \cup \Gamma' \cup \Delta} A, \Delta & & \text{cut} \end{array} \\ \downarrow \\ \begin{array}{ccc} \pi_1 & & \pi_2 \\ \vdots & & \vdots \\ \Sigma \mid \Gamma, A \mid_{\Theta} B, \Delta & \quad \quad & \Sigma' \mid \Gamma' \mid_{\Psi \cup \Gamma' \cup \Delta} B, \Delta \\ \hline \Sigma', \Sigma, A \mid \Gamma', \Gamma \mid_{\Theta \cup \Psi \cup \Gamma' \cup \Delta} A & \text{abduct.} & \\ \hline \Sigma', \Sigma, A \mid \Gamma', \Gamma \mid_{\Theta \cup \Psi \cup \Gamma' \cup \Delta} A, \Delta & \text{RW} & \end{array} \end{array}$$

■

The conservativity of LEA^+ follows immediately from cut's eliminability.

Corollary 4.1. Every theorem in LEA^+ that only contains formulas from \mathcal{L} is a theorem of LEA.

Proof. Since $\blacktriangleright \vdash$ and $\vdash \blacktriangleright$ are the only rules in LEA^+ that are not in LEA, the only source of new theorems formulated in \mathcal{L} are those derived via an application of cut to the conclusions of these two rules. It follows from Theorem 4 that this application of cut is eliminable. ■

7 Conclusion

In this paper, we have argued that there is a viable inferentialist-expressivist treatment of explanatory vocabulary. More specifically, we have shown how explanatory arguments and IBE can be represented by a sequent calculus and how that calculus can be conservatively extended to a language that contains explicitly explanatory expressions. We conclude by discussing some of the peculiarities and limitations of the current approach as well as the prospects for future work.

Since our calculi are constructed on the basis of LK, our defeasible sequents represent relations among sets of formulas, while our connectives are operations on formulas. In LK this discrepancy between the relata of consequence relations and the relata of connectives is completely natural. However, because our \blacktriangleright -sequents are intended to capture an exhaustive relation—namely, the relation of loveliest potential *exhaustive* immediate explanation—the fact that the \blacktriangleright connective only holds among formulas might threaten a change in meaning.

Fortunately, this concern is easily allayed. What $\vdash \blacktriangleright$ says is that if Γ, A is the best *exhaustive* immediate explanation of B , then A is the best *exhaustive* immediate explanation of B given (all of) Γ . This reading of the rule corresponds nicely to the pragmatics of explanation. We are rarely in search of *complete* explanations. Rather, our inquiries aim at *the* (best) explanation of some phenomena. The latter may be thought of as a *selection* from the former. To proffer a *selected* explanation is to put forth one claim as an explanans while *assuming* that the remaining components of an exhaustive explanation hold. For instance, to explain why George ordered the chocolate cake, we might appeal to a stable preference—e.g. “Chocolate cake is George’s favorite dessert.” But the sufficiency of this explanation rests on a number of assumptions—for instance, that a preference for chocolate cake would lead someone to choose chocolate over carrot cake. Our formalism does justice to this fact: someone entitled to an exhaustive explanatory argument is entitled to claim that some element of the premises (A) is the best explanation of the conclusion (B), so long as one is entitled to the remaining premises (Γ). In fact, our insistence that defeasible sequents be paired with background sets (Σ) supports the view that even the propriety of exhaustive explanatory arguments depends upon what background assumptions are in play.

Of course, what makes one selected explanation acceptable in practice as opposed to another is no doubt in part a function of speaker interests. Consequently, a satisfactory theory of best explanations must account for the role that the practical commitments of speakers (i.e. plans, projects, preferences, interests) play in determining *which* components of an exhaustive explanation can serve as the best explanation. To do so, the theory must provide a means for conceptualizing the pragmatics of explanation.

Wish List

1. Our calculi provide opportunities for the formalization of different target vocabularies. For instance, while our aim was to offer an expressivist treatment of immediate explanation, there is no reason to think that a similar treatment of *mediate* explanation is unavailable. Such an account would need to permit a version of cut in LEA. Adding the cut rule from LK^Θ would mean that transitivity can fail when one of the sequents with the cut formula is defeated. Is this just the sort of non-transitivity we want?
2. Unlike many nonmonotonic logics, cautious monotonicity does not hold in LEA. How would we include this rule? Would we want to? Think about Ulf’s line on this: CM says that you can always add implicit content *explicitly* to the premises.
3. What additional logical rules should be included in LE^\blacktriangleright ?

4. How to represent probabilistic inferences? Treat defeater sets as things that would bring the inference below a threshold probability.