

# Modeling Contrastiveness in Argumentation

AnneMarie Borg<sup>1</sup>, Floris Bex<sup>1,2</sup>

<sup>1</sup>Department of Information and Computing Sciences, Utrecht University, Utrecht, The Netherlands

<sup>2</sup>Tilburg Institute for Law, Technology, and Society, Tilburg University, Tilburg, The Netherlands

## Abstract

Modeling contrastive explanations for the use in artificial intelligence (AI) applications is an important research branch within the field of explainable AI (XAI). However, most of the existing contrastive XAI approaches are not based on the findings in the literature from the social sciences on contrastiveness in human reasoning and human explanations. In this work we collect the various types of contrastiveness proposed in the literature and model these with formal argumentation. The result is a variety of argumentation-based methods for contrastive explanations, based on the available literature and applicable in a wide variety of AI-applications.

## 1. Introduction

Explainable Artificial Intelligence (XAI) is an important and fast growing research area, contributing to closing the gap between AI-application and its human user. To this end it is essential that XAI approaches incorporate findings from the humanities and social sciences on how humans request, generate, interpret and evaluate explanations (e.g., research on explanations from philosophy [1] or law [2], see [3] for an extensive survey).

Argumentation plays an important role in XAI, because it is central to all human reasoning [4] including explanation [5, 6]. In addition to existing argumentation-based systems being inherently interpretable (cf. [7]), argumentation has also been used to explain the output of other less interpretable AI models such as machine learning classification models (e.g., [8]), deep reinforcement learning models [9] and probabilistic Bayesian Networks [10] (see [11] for a recent overview). Such explanations in terms of arguments can come in many forms, taking into account the above-mentioned literature from the humanities and social sciences. For example, argument-based explanations can be minimal [12], in the form of a dialogue [13], or include only necessary and sufficient reasons [14].

One aspect of everyday explanations that has so far received relatively little attention in the literature on argumentative explanations is that of *contrastiveness*: when asking *why P?* we often expect the explanation to answer *why P rather than Q?* Contrastive explanations can be used to compare a surprising outcome with the expected outcome [1, 15, 16]. Additionally, by using contrastiveness, explanations can be interactively customized and personalized [17], and become easier to derive and cognitively less demanding [1, 18].

There are a lot of studies on contrastive explanations in XAI, but very few of these base their


---

Workshop on Computational Models of Natural Argument (CMNA'22), Cardiff, UK

✉ a.borg@uu.nl (A. Borg); f.j.bex@uu.nl (F. Bex)



© 2022 Copyright for this paper by its authors. Use permitted under Creative Commons License Attribution 4.0 International (CC BY 4.0).

 CEUR Workshop Proceedings (CEUR-WS.org)

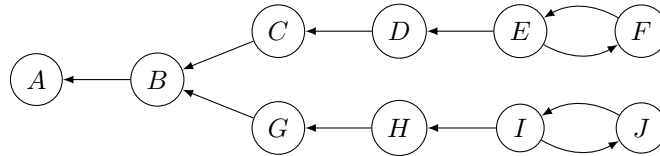
explanation definitions on the theoretical foundations from the humanities and social sciences literature on contrastiveness in human reasoning and human explanation [19]. In this work we study the concept of contrastiveness as can be found in the literature on explanations and formalize our findings in terms of formal argumentation. After a brief introduction to formal argumentation, we discuss the various types of contrastiveness as proposed in the literature as well as the requirements placed on them [1, 19, 16, 18, 20, 21]. We then show how these types of contrastiveness can be modeled in abstract argumentation [22] (Section 3.2). Due to space restrictions we only study contrastiveness for the (non-)acceptance of arguments in abstract argumentation. However, the proposed explanations could be applied to a structured setting (e.g., one of those in [23]) in a similar way as was done in [24].

With this paper we start a discussion on contrastiveness for argumentative XAI. We provide the basis for this discussion by giving a literature overview and showing how notions from the literature can be implemented in abstract argumentation.

## 2. Preliminaries: argumentation and explanation

### 2.1. Abstract Argumentation

An *abstract argumentation framework* (AF) [22] is a pair  $\mathcal{AF} = \langle \text{Args}, \text{Att} \rangle$ , where  $\text{Args}$  is a set of *arguments* and  $\text{Att} \subseteq \text{Args} \times \text{Args}$  is an *attack relation* on these arguments. An AF can be viewed as a directed graph, see Figure 1 for an example, in which the nodes represent arguments and the arrow represent attacks between arguments.



**Figure 1:** Graphical representation of the argumentation framework  $\mathcal{AF}_1$ .

Given an argumentation framework  $\mathcal{AF}$ , Dung-style semantics [22] can be applied to it, to determine what combinations of arguments (called *extensions*) can collectively be accepted.<sup>1</sup>

**Definition 1.** Let  $\mathcal{AF} = \langle \text{Args}, \text{Att} \rangle$  be an AF,  $S \subseteq \text{Args}$  be a set of arguments and let  $A \in \text{Args}$ . Then:  $S$  attacks  $A$  if there is an  $A' \in S$  such that  $(A', A) \in \text{Att}$ ;  $S$  defends  $A$  if  $S$  attacks every attacker of  $A$ ;  $S$  is conflict-free if there are no  $A_1, A_2 \in S$  such that  $(A_1, A_2) \in \text{Att}$ ;  $S$  is admissible if it is conflict-free and it defends all of its elements.

An admissible extension that contains all the arguments that it defends is a complete extension (Cmp). The grounded extension (Grd) is the minimal (w.r.t.  $\subseteq$ ) complete extension and a preferred extension (Prf) is a maximal (w.r.t.  $\subseteq$ ) complete extension.

We denote by  $\text{Sem}(\mathcal{AF})$  the set of all extensions of  $\mathcal{AF}$  under the semantics  $\text{Sem} \in \{\text{Grd}, \text{Cmp}, \text{Prf}\}$ . Based on the semantics, an argument can be accepted and/or not accepted.

<sup>1</sup>Other semantics can be found in, e.g., [25]. We do not discuss these, since the specific semantics is not relevant for our discussion.

**Definition 2.** Let  $\mathcal{AF} = \langle \text{Args}, \text{Att} \rangle$  be an AF,  $\text{Sem} \in \{\text{Grd}, \text{Cmp}, \text{Prf}\}$  a semantics and let  $A \in \text{Args}$  be an argument. Then:  $A$  is accepted if there is some  $\mathcal{E} \in \text{Sem}(\mathcal{AF})$  such that  $A \in \mathcal{E}$ ; and  $A$  is not accepted (or non-accepted) if there is some  $\mathcal{E} \in \text{Sem}(\mathcal{AF})$  such that  $A \notin \mathcal{E}$ .

Note that, since the grounded extension is unique, if  $\text{Sem} = \text{Grd}$ , an argument is either accepted or not accepted. For the other semantics an argument can be accepted, not accepted or both (i.e., there might be extensions with the argument and extensions without the argument).

**Example 1.** In  $\mathcal{AF}_1$  from Figure 1 the grounded extension is empty:  $\text{Grd}(\mathcal{AF}_1) = \{\emptyset\}$  and there are four preferred extensions  $\text{Prf}(\mathcal{AF}_1) = \{\mathcal{E}_1, \mathcal{E}_2, \mathcal{E}_3, \mathcal{E}_4\}$  where  $\mathcal{E}_1 = \{A, C, E, G, I\}$ ,  $\mathcal{E}_2 = \{A, C, E, H, J\}$ ,  $\mathcal{E}_3 = \{A, D, F, G, I\}$  and  $\mathcal{E}_4 = \{B, D, F, H, J\}$ . It follows that all arguments of  $\mathcal{AF}_1$  are both accepted and non-accepted for  $\text{Sem} \in \{\text{Cmp}, \text{Prf}\}$ .

## 2.2. Explanations in Argumentation

We now recall the main definitions for the explanation framework from [24], in which basic explanations for both the acceptance and non-acceptance of some argument  $A$  under various semantics were defined. We will apply a general notation for acceptance and non-acceptance explanations, without any restrictions:

**Notation 1.** Let  $\mathcal{AF} = \langle \text{Args}, \text{Att} \rangle$  be an AF,  $A \in \text{Args}$ ,  $\text{Sem} \in \{\text{Grd}, \text{Cmp}, \text{Prf}\}$  be a semantics and  $\mathcal{E} \in \text{Sem}(\mathcal{AF})$  an extension. Then:

- $\text{SemAcc}(A, \mathcal{E})$  denotes an acceptance explanation for  $A$  in the specific extension  $\mathcal{E}$ ; and
- $\text{SemNotAcc}(A, \mathcal{E})$  denotes a non-acceptance explanation for  $A$  in the specific extension  $\mathcal{E}$ .

By using this general notation, other argumentative explanation methods can be used as well, e.g., [12, 26, 27, 28]. Note that these explanations differ slightly from those introduced in [24], in that they also require an extension  $\mathcal{E}$  in addition to an argument  $A$  – the reason for this will become apparent in Section 3.2. Because explanations in abstract argumentation are defined in terms of attack and defense between arguments, we also need the following definition.

**Definition 3.** Given  $\mathcal{AF} = \langle \text{Args}, \text{Att} \rangle$ , we say that  $A \in \text{Args}$  directly defends  $B$  if there is some  $C$  such that  $(C, B) \in \text{Att}$  and  $(A, C) \in \text{Att}$ , and  $A$  indirectly defends  $B$  if  $A$  defends  $C$  and  $C$  defends  $B$ . Similarly,  $A$  directly attacks  $B \in \text{Args}$  if  $(A, B) \in \text{Att}$ ,  $A$  indirectly attacks  $B$  if  $A$  attacks some  $C \in \text{Args}$  and  $C$  defends  $B$ . We will often say that  $A$  defends [resp. attacks]  $B$  when  $A$  directly or indirectly defends [resp. attacks]  $B$ .

Now we can define basic explanations for (non-)accepted arguments as follows.

**Definition 4.** Let  $\mathcal{AF} = \langle \text{Args}, \text{Att} \rangle$  be an AF and let  $\text{Sem} \in \{\text{Grd}, \text{Cmp}, \text{Prf}\}$ . First, let  $A \in \text{Args}$  be accepted and let  $\mathcal{E} \in \text{Sem}(\mathcal{AF})$  such that  $A \in \mathcal{E}$ , then:

$$\text{SemAcc}(A, \mathcal{E}) = \mathbb{D}^{\text{acc}}(A, \mathcal{E}).$$

Now, let  $B \in \text{Args}$  be non-accepted and let  $\mathcal{E} \in \text{Sem}(\mathcal{AF})$  such that  $B \notin \mathcal{E}$ , then:

$$\text{SemNotAcc}(B, \mathcal{E}) = \mathbb{D}^{\text{nacc}}(B, \mathcal{E}).$$

The explanations use a generic function  $\mathbb{D}$  that returns a set of arguments as an explanation given a (non-)accepted argument and an extension. This function allows us to vary which attacking or defending arguments are in an explanation for an acceptance explanation ( $\mathbb{D}^{\text{acc}}$ ) or a non-acceptance explanation ( $\mathbb{D}^{\text{nacc}}$ ). The two common instantiations of  $\mathbb{D}^{\text{acc}}$  and  $\mathbb{D}^{\text{nacc}}$  that we will use throughout this paper are as follows.<sup>2</sup>

- $\text{Defending}(A, \mathcal{E}) = \{B \in \mathcal{E} \mid B \text{ (in)directly defends } A\}$ ;
- $\text{NoDefAgainst}(A, \mathcal{E}) = \{B \in \mathcal{E} \mid B \text{ attacks } A \text{ and } \mathcal{E} \text{ does not attack } B\}$ .

When  $\mathbb{D}^{\text{acc}} = \text{Defending}$ , an acceptance explanation for argument  $A$  in  $\mathcal{E}$  is the set of arguments from  $\mathcal{E}$  that defends  $A$  against all its attackers, and for  $\mathbb{D}^{\text{nacc}} = \text{NoDefAgainst}$ , a non-acceptance explanation for argument  $A \notin \mathcal{E}$  is the set of attackers of  $A$  against which  $\mathcal{E}$  provides no defence.

**Example 2.** Recall that for  $\mathcal{AF}_1$  we have the preferred extensions  $\mathcal{E}_1, \mathcal{E}_2, \mathcal{E}_3$  and  $\mathcal{E}_4$ , then:

- $\text{PrfAcc}(A, \mathcal{E}_1) = \{C, E, G, I\}$ ,  $\text{PrfAcc}(A, \mathcal{E}_2) = \{C, E\}$ ,  $\text{PrfAcc}(A, \mathcal{E}_3) = \{G, I\}$  and  $\text{PrfAcc}(B, \mathcal{E}_4) = \{D, F, H, J\}$ ;
- $\text{PrfNotAcc}(A, \mathcal{E}_4) = \{B, D, F, H, J\}$  and  $\text{PrfNotAcc}(B, \mathcal{E}_1) = \{C, E, G, I\}$ ,  $\text{PrfNotAcc}(B, \mathcal{E}_2) = \{C, E\}$ ,  $\text{PrfNotAcc}(B, \mathcal{E}_3) = \{G, I\}$ .

### 3. Contrastiveness in argumentation

The explanations in Section 2.2 are simple, non-contrastive explanations, answers to a *Why (not) P?* explanation-seeking question, where  $P$  is understood as the (non-)acceptance of an argument. In its most basic form a contrastive explanation-seeking question is then of the form *Why P rather than Q?*, where  $P$  is called the *fact* and  $Q$  is called the *foil* [1]. However, there is no agreement on exactly what is or constitutes a contrastive explanation [19]. Like Rudin [7], we assume that more than one definition is possible, and below we will discuss the different ideas on contrastiveness and how they can be interpreted for argumentative explanations.

The specific contrastive question may vary. First, there is the *negation* question, where the fact is compared with the situation in which it does not occur: *Why P rather than not-P?* [16]. Second, it is possible to compare the *properties* of a certain object, resulting in questions like: *Why P rather than the default value for P?* [16] or *Why does object a have property P, rather than property Q?* [15]. Finally, different *objects* can be compared as well: *Why P rather than Q?* [16] and *Why does object a have property P, while object b has property Q?* [15].

- I. Considering explanations in abstract argumentation, the fact  $P$  concerns the (non-)acceptance of an argument, which is also the only property of the argument (i.e., whether it is accepted or not). Thus, both negation-type and property-type questions are captured by comparing why an argument is accepted as opposed to not accepted. Object-type questions then compare different arguments (objects) to each other.

<sup>2</sup>Other instantiations of  $\mathbb{D}$  are, e.g.,  $\text{DirDefending}$  and  $\text{DirAttacking}$ , which take only the direct defenders or attackers of  $A$ , see [14, 24].

The idea of contrastiveness usually assumes that fact and foil are incompatible [18, 20]. Even for a seemingly compatible fact and foil, what we are really interested in is some sort of (hypothetical) contrast. Take, for example, the question ‘Why does John have disease  $X$  and Harry does not?’ In this case, it would be perfectly possible for Harry to also have disease  $X$ , so fact and foil are not inherently incompatible. However, as Ylikoski [18] notes, we are not really interested in Harry (why he does not have  $X$ ), but rather treat him as a surrogate for John, for whom we want to know why he has  $X$  while someone similar to him (Harry) does not.

- II. We will assume that fact and foil are not always compatible. This means that fact and foil are not always part of the same extensions, i.e., there is at least one extension in which fact or foil is accepted and the other is not. This can be further restricted by, e.g., requiring that fact and foil are never part of the same extension (i.e., always incompatible).

Note that it is assumed that fact and foil should be similar or closely related. For example, they should have common presuppositions or a common cause, or share similar causal histories [1, 18, 20]. For example, it makes little sense to ask ‘Why does John have disease  $X$  while my cat does not?’. We assume that John and Harry are related in that they are comparable persons (age, gender, background, etc.) who share a similar history, i.e., the same (type of) events that have caused John to contract disease  $X$  could have also affected Harry.

- III. We will assume that fact and foil are *relevant* for each other, which means that they are at least connected via the attack relation.

Sokol and Flach [17] differentiate between *Why  $P$  despite  $Q$ ?* and a *Why  $P$  given  $Q$ ?* The despite-question requires an explanation that is not conditioned on  $Q$ , that is, the explanation for the fact  $P$  does not take into account the explanation for the foil  $Q$ . The given-question requests an explanation in which the explanation for fact  $P$  is conditioned on the explanation for foil  $Q$ : what do the explanations for fact and foil have in common? Take, for example, the question ‘Why did John get a loan and Harry did not?’. A despite-explanation would then be ‘because John has an income of more than 100k’ (John’s income is not directly relevant for Harry), and a given-explanation would be ‘because John was first and only one person can get a loan’ (John being first is directly relevant for Harry).

- IV. For a given-question, the explanation of the fact is conditioned on the foil (the explanation contains the part of the explanations that are the same for both fact and foil). For a despite-question, the explanation of the fact is not conditioned on the foil (the explanation contains the part of the explanation for the fact without the foil).

The final important difference made in the literature is the one between alternative and actual foils, which is made explicit in [18]. On the alternative side there are two incompatible and alternative outcomes of the same process, one the fact and one the foil. On the actual side there are two actual (possibly compatible) outcomes, fact and foil, from two actual and different processes. From these two types of contrasts, two questions are formulated and modeled in [29]. The *alternative* question asks *Why  $P$  rather than  $Q$ ?*, where fact  $P$  occurred and foil  $Q$  did not, while the *actual* question asks *Why  $P$  but  $Q$ ?*, where  $P$  occurred in the current situation and  $Q$  occurred in some other situation.

- V. We will consider both types of foils. To this end we will take a fixed argumentation framework (i.e., for now we do not allow for dynamic settings) and take extensions within that framework to represent different situations. The alternative question is then modeled given a specific extension (i.e., within an extension), while the actual question is modeled between extensions.

### 3.1. Relevance and Compatibility

In view of considerations II. and III. above we introduce the notions of compatibility and relevance. Compatibility concerns the mutual acceptance of arguments in the same extensions.

**Definition 5.** Let  $\mathcal{AF} = \langle \text{Args}, \text{Att} \rangle$  be an AF and  $A, B \in \text{Args}$  be two arguments. Then, for  $\text{Sem} \in \{\text{Cmp}, \text{Grd}, \text{Prf}\}$ :  $A$  and  $B$  are compatible w.r.t.  $\text{Sem}$  if for some  $\mathcal{E} \in \text{Sem}(\mathcal{AF})$  it holds that  $A, B \in \mathcal{E}$ ;  $A$  and  $B$  are incompatible w.r.t.  $\text{Sem}$  if there is no  $\mathcal{E} \in \text{Sem}(\mathcal{AF})$  such that  $A, B \in \mathcal{E}$ ; and  $A$  and  $B$  are not always compatible w.r.t.  $\text{Sem}$  if there is some  $\mathcal{E} \in \text{Sem}(\mathcal{AF})$  such that  $A \in \mathcal{E}$  and  $B \notin \mathcal{E}$ .

The above can be generalized to sets of arguments. Given two sets of arguments  $S_1, S_2 \subseteq \text{Args}$ :  $S_1$  and  $S_2$  are compatible w.r.t.  $\text{Sem}$  if for each pair  $(A, B) \in S_1 \times S_2$ , there is some  $\mathcal{E} \in \text{Sem}(\mathcal{AF})$  such that  $A, B \in \mathcal{E}$ ;  $S_1$  and  $S_2$  are incompatible w.r.t.  $\text{Sem}$  if there is no  $\mathcal{E} \in \text{Sem}(\mathcal{AF})$  such that  $S_1 \cap \mathcal{E} \neq \emptyset$  and  $S_2 \cap \mathcal{E} \neq \emptyset$ ;  $S_1$  and  $S_2$  are not always compatible w.r.t.  $\text{Sem}$  if there is some  $\mathcal{E} \in \text{Sem}(\mathcal{AF})$  such that  $S_1 \cap \mathcal{E} \neq \emptyset$  and  $S_2 \cap \mathcal{E} = \emptyset$ .

Arguments or sets of arguments are compatible if these can be part of the same extension, incompatible if there is no extension that contains both arguments or arguments from both sets and not always compatible if there is an extension that only contains one argument or arguments from one set and not the other. Given II. we will assume that fact and foil are not always compatible. Given III. we will assume that fact and foil are relevant for each other:

**Definition 6.** Let  $\mathcal{AF} = \langle \text{Args}, \text{Att} \rangle$  be an AF and  $A, B \in \text{Args}$  be two arguments. Then  $A$  is relevant for  $B$  if:  $A$  (in)directly attacks or defends  $B$ ; or there is some  $C \in \text{Args}$  such that  $C$  is relevant for both  $A$  and  $B$ . We say that  $A$  is conflict-relevant for  $B$  if it (in)directly attacks  $B$  and  $A$  is defending-relevant for  $B$  if it (in)directly defends  $B$ .

Arguments are relevant for each other when they are connected by the attack relation, or if there is an argument that is connected by the attack relation to both of the considered arguments.

**Example 3.** For the argumentation framework  $\mathcal{AF}_1$  we have that  $E$  and  $H$  are relevant for  $B$  and  $A$  but not for each other. Also, for  $\text{Sem} = \text{Prf}$ ,  $E$  and  $H$  are compatible (there is the extension  $\{A, C, E, H, J\}$ ), but not always (there is the extension  $\{A, C, E, G, I\}$ ) and although  $H$  is compatible with  $B$  (there is the extension  $\{B, D, F, H, J\}$ ),  $E$  is not compatible with  $B$ .

### 3.2. Contrastive Explanations

We can now formalize the notions of contrastiveness discussed earlier. Recall that we consider (I.) negation-type and object-type questions; (II.), (III.) relevant, but not always compatible fact and foil; (IV.) given- and despite-questions; and (V.) alternative (within an extension) and actual

	Neg + despite	Obj + despite	Neg + given	Obj + given
Within	Note 1	Def. 7 & Ex. 4	Note 1	Def. 8 & Ex. 5
Between	Def. 9 & Ex. 7	Def. 9 & Ex. 7	Def. 10 & Ex. 8	Def. 10 & Ex. 8

**Table 1**

Overview of the types of contrastiveness for considerations I. to V.

(between extensions) questions. The resulting eight types of contrastive explanations are shown in Table 1.

**Note 1.** *An argument cannot be part of an extension and, at the same time, not be part of that extension. Given II. (not always compatible) and the fact that within explanations look at just one extension, a contrastive explanation for within will always be an object-type explanation (i.e., between two different arguments).*

We start with the contrastive explanation for *despite* and *within*: *Why A despite B in the current situation?* To ensure a form of contrast, since the fact is accepted, the foil is non-accepted if it is defending-relevant for the fact and the foil is accepted if it is conflict-relevant for the fact. We do not limit the foil to be a single argument, rather, we define the explanations between one argument (the fact) and a set of arguments (the foil).

**Definition 7.** *Given  $\mathcal{AF} = \langle \text{Args}, \text{Att} \rangle$ ,  $A \in \text{Args}$ ,  $S \subseteq \text{Args}$  and  $\text{Sem} \in \{\text{Grd}, \text{Cmp}, \text{Prf}\}$ . Let  $\text{Expl} \in \{\text{Acc}, \text{NotAcc}\}$  determine the type of explanation and let  $\mathcal{E} \in \text{Sem}(\mathcal{AF})$ , then:*

$$\begin{aligned} \text{SemCont}((A, \mathcal{E}), (S, \mathcal{E})) &= \text{SemAcc}(A, \mathcal{E}) \setminus \left( \bigcup_{B \in S} \text{SemExpl}(B, \mathcal{E}) \right) \\ \text{SemContN}((A, \mathcal{E}), (S, \mathcal{E})) &= \text{SemNotAcc}(A, \mathcal{E}) \setminus \left( \bigcup_{B \in S} \text{SemExpl}(B, \mathcal{E}) \right). \end{aligned}$$

The explanation contains the reasons for the acceptance of the fact that are not part of the acceptance (if  $\text{Expl} = \text{Acc}$ ) or non-acceptance (if  $\text{Expl} = \text{NotAcc}$ ) explanation of the foil and similarly for a non-accepted fact.

**Example 4.** *In the running example with  $\mathcal{AF}_1$ , there are three preferred extensions in which A is accepted and one in which it is not accepted. Now, when A is accepted it has to be defended against B, which C can do. However, A might still be accepted without C (i.e., in  $\mathcal{E}_3$ ).*

- Why A despite not-C? Recall that  $\text{PrfAcc}(A, \mathcal{E}_3) = \{G, I\}$  and that  $\text{PrfNotAcc}(C, \mathcal{E}_3) = \{D, F\}$ . We therefore have, for  $\text{Expl} = \text{NotAcc}$  (i.e., comparing the acceptance of the fact with the non-acceptance of the foil):  $\text{PrfCont}((A, \mathcal{E}_3), (\{C\}, \mathcal{E}_3)) = \{G, I\}$ .
- Why not-B despite not-C? This question is very similar to the above one, where we will now apply  $\text{PrfNotAcc}(B, \mathcal{E}_3) = \{G, I\}$ . We then have, again for  $\text{Expl} = \text{NotAcc}$  now comparing two non-accepted arguments:  $\text{PrfContN}((B, \mathcal{E}_3), (\{C\}, \mathcal{E}_3)) = \{G, I\}$ .
- Why not-B despite not-G? The explanation for the non-acceptance of B was:  $\text{PrfNotAcc}(B, \mathcal{E}_2) = \{C, E\}$ . Similarly to the item above, with  $\text{Expl} = \text{NotAcc}$ :  $\text{PrfContN}((B, \mathcal{E}_2), (\{G\}, \mathcal{E}_2)) = \{C, E\}$ .



For the contrastive explanation *given* and *within* we answer *Why A given its alternative B in the current situation?* Similar to the case of *despite* and *within* contrasts, independent of whether the fact is accepted or non-accepted, the foil can be accepted or non-accepted. However, for accepted [resp. non-accepted] fact, if the foil is conflict-relevant for the fact it should be non-accepted [resp. accepted] and if it is defending-relevant for the fact it should be accepted [resp. non-accepted], this to ensure a contrast between fact and foil.

**Definition 8.** Given  $\mathcal{AF} = \langle \text{Args}, \text{Att} \rangle$ ,  $A \in \text{Args}$ ,  $S \subseteq \text{Args}$  and  $\text{Sem} \in \{\text{Grd}, \text{Cmp}, \text{Prf}\}$ . Let  $\text{Expl} \in \{\text{Acc}, \text{NotAcc}\}$  determine the type of explanation and let  $\mathcal{E} \in \text{Sem}(\mathcal{AF})$ , then:

$$\begin{aligned} \text{SemCont}((A, \mathcal{E}), (S, \mathcal{E})) &= \text{SemAcc}(A, \mathcal{E}) \cap \left( \bigcup_{B \in S} \text{SemExpl}(B, \mathcal{E}) \right) \\ \text{SemContN}((A, \mathcal{E}), (S, \mathcal{E})) &= \text{SemNotAcc}(A, \mathcal{E}) \cap \left( \bigcup_{B \in S} \text{SemExpl}(B, \mathcal{E}) \right). \end{aligned}$$

While in Definition 7 the explanation contained the reasons for the fact without the reasons for the foil, now the intersection of the explanations for fact and foil are taken. This follows since the question assumes that the foil is *given* and should be part of the explanation.

**Example 5.** For the running example with  $\mathcal{AF}_1$ :

- Why *A* given not-*B*? Since there are several non-acceptance explanations for *B*, we have that  $\text{PrfCont}((A, \mathcal{E}_2), (\{B\}, \mathcal{E}_2)) = \{C, E\}$  and  $\text{PrfCont}((A, \mathcal{E}_3), (\{B\}, \mathcal{E}_3)) = \{G, I\}$ .
- Why *B* given not-*G*? Recall that there is only one extension with *B*. Therefore:  $\text{PrfCont}((B, \mathcal{E}_4), (\{G\}, \mathcal{E}_4)) = \{H, J\}$ .

For the fact-foil pairs from Example 4, where the foil is now accepted:

- Why *A* given *C*? or Why not-*B* given *C*? From the acceptance explanation for *C* (i.e.,  $\text{PrfAcc}(C, \mathcal{E}_2) = \{E\}$ ), for  $\text{Expl} = \text{Acc}$ :  $\text{PrfCont}((A, \mathcal{E}_2), (\{C\}, \mathcal{E}_2)) = \{E\}$  and  $\text{PrfContN}((B, \mathcal{E}_2), (\{C\}, \mathcal{E}_2)) = \{E\}$ .
- Why not-*B* given *G*? The explanation is very similar as the explanations in the above item, now based on  $\text{PrfAcc}(G, \mathcal{E}_3) = \{I\}$ . Then, for  $\text{Expl} = \text{Acc}$ :  $\text{PrfContN}((B, \mathcal{E}_3), (\{G\}, \mathcal{E}_3)) = \{I\}$ .

We now turn to the comparison *between* extensions. These differ from the contrastive explanations within an extension, in that we do not take the difference or intersection with (a set of) explanations for a foil from the same extension, but rather with the second extension as containing the foil. In other words, we compare the actual situation (extension) of the fact with the actual situation (extension) of the foil.

**Example 6.** Suppose we want to compare the acceptance of *B* in one extension with the acceptance of *A* in another. For example:  $\text{PrfAcc}(B, \mathcal{E}_4) = \{D, F, H, J\}$  and  $\text{PrfAcc}(A, \mathcal{E}_2) = \{C, E\}$ . Taking *A* as the specific foil does not change the explanation, since the explanations do not intersect. If we instead look at  $\mathcal{E}_2$  as a whole (recall  $\mathcal{E}_2 = \{A, C, E, H, J\}$ ) we can take the difference and obtain the explanation  $\{D, F\}$ : *B* is accepted in  $\mathcal{E}_4$  despite *A* being accepted in  $\mathcal{E}_2$ , since in  $\mathcal{E}_4$  *D* and *F* are accepted as well.



We start with the contrastive explanation for *despite*: *Why A despite B in another situation?* Fact and foil can both be (non-)accepted, since we compare different extensions.

**Definition 9.** Given  $\mathcal{AF} = \langle \text{Args}, \text{Att} \rangle$ ,  $A \in \text{Args}$ ,  $S \subseteq \text{Args}$  and  $\text{Sem} \in \{\text{Cmp}, \text{Prf}\}$ . Let  $\text{Expl} \in \{\text{Acc}, \text{NotAcc}\}$  determine the type of explanation and let  $\mathcal{E}, \mathcal{E}' \in \text{Sem}(\mathcal{AF})$ , where  $\mathcal{E} \neq \mathcal{E}'$ . Moreover,  $\text{SemExpl}(B, \mathcal{E}')$  is an explanation according to Notation 1 for all  $B \in S$ . Then:

$$\begin{aligned} \text{SemCont}((A, \mathcal{E}), (S, \mathcal{E}')) &= \text{SemAcc}(A, \mathcal{E}) \setminus \mathcal{E}' \\ \text{SemContN}((A, \mathcal{E}), (S, \mathcal{E}')) &= \text{SemNotAcc}(A, \mathcal{E}) \setminus \mathcal{E}'. \end{aligned}$$

These explanations present the difference of the explanation for the fact in one extension and the extension of the foil. It is assumed that the foil argument(s) are (non-)accepted in the extension  $\mathcal{E}'$  and that proper explanations exist. This is necessary to represent the fact-foil construction of contrastive explanations.

Because fact and foil are now in different extensions, we can also ask negation-type questions, where the acceptance of an argument is compared with its non-acceptance.

**Example 7.** For the running example with  $\mathcal{AF}_1$ , suppose that we are interested in the extensions  $\mathcal{E}_2$  and  $\mathcal{E}_4$  for  $\text{Sem} = \text{Prf}$ . We start by comparing the acceptance and non-acceptance of  $B$  for these settings:

- Why  $B$  in  $\mathcal{E}_4$  despite not- $B$  in  $\mathcal{E}_2$ ? Then, for  $\text{Expl} = \text{NotAcc}$ :  $\text{PrfCont}((B, \mathcal{E}_4), (\{B\}, \mathcal{E}_2)) = \{D, F\}$ .
- Why not- $B$  in  $\mathcal{E}_2$  despite  $B$  in  $\mathcal{E}_4$ ? This results, for  $\text{Expl} = \text{NotAcc}$ , in:  $\text{PrfContN}((B, \mathcal{E}_2), (\{B\}, \mathcal{E}_4)) = \{C, E\}$ .

We can also compare the arguments from  $A$ ,  $B$ ,  $C$  and  $H$ :

- Why  $A$  in  $\mathcal{E}_2$  despite  $B$  in  $\mathcal{E}_4$ ? For  $\text{Expl} = \text{Acc}$ :  $\text{PrfCont}((A, \mathcal{E}_2), (\{B\}, \mathcal{E}_4)) = \{C, E\}$ .
- Why  $B$  in  $\mathcal{E}_4$  despite  $A$  in  $\mathcal{E}_2$ ? The explanation is similar to the explanation above, for  $\text{Expl} = \text{Acc}$ :  $\text{PrfCont}((B, \mathcal{E}_4), (\{A\}, \mathcal{E}_2)) = \{D, F\}$ .
- Why  $A$  in  $\mathcal{E}_2$  despite not- $C$  in  $\mathcal{E}_3$ ? For  $\text{Expl} = \text{NotAcc}$ :  $\text{PrfCont}((A, \mathcal{E}_2), (\{C\}, \mathcal{E}_3)) = \{C, E\}$ .
- Why not- $A$  in  $\mathcal{E}_4$  despite  $H$  in  $\mathcal{E}_2$ ? For  $\text{Expl} = \text{Acc}$ :  $\text{PrfContN}((A, \mathcal{E}_4), (\{H\}, \mathcal{E}_2)) = \{B, D, F\}$ .

**Note 2.** With these explanation it becomes apparent that the choice of the right foil is essential when formulating the explanation. While the above explanations make sense when an argument is compared with its negation (see the first part of Example 7), this might be less so when two arguments are compared. This has to be accounted for when choosing the formalization, extensions, fact and foil. For example (last two bullets in Example 7), why  $A$  in  $\mathcal{E}_2$  despite not- $C$  in  $\mathcal{E}_3$  might not be an informative explanation, while why not- $A$  in  $\mathcal{E}_4$  despite  $H$  in  $\mathcal{E}_2$  might be informative, since  $H$  is a reason for the non-acceptance of  $A$  in  $\mathcal{E}_4$  but they are both accepted in  $\mathcal{E}_2$ . What conditions are placed on the foil depends on the applications and situation at hand and is left for future work.

Finally, we consider the contrastive explanation for *given* and *between*: *Why A given B in another situation?* Again, rather than taking the intersection with a set of explanations, we take the intersection with the second extension, to ensure a comparison between the fact and its extension with the situation of the foil.

**Definition 10.** Given  $\mathcal{AF} = \langle \text{Args}, \text{Att} \rangle$ ,  $A \in \text{Args}$ ,  $S \subseteq \text{Args}$  and  $\text{Sem} \in \{\text{Cmp}, \text{Prf}\}$ . Let  $\text{Expl} \in \{\text{Acc}, \text{NotAcc}\}$  determine the type of explanation and let  $\mathcal{E}, \mathcal{E}' \in \text{Sem}(\mathcal{AF})$ , where  $\mathcal{E} \neq \mathcal{E}'$ . Moreover,  $\text{SemExpl}(B, \mathcal{E}')$  is an explanation according to Notation 1 for all  $B \in S$ . Then:

$$\begin{aligned} \text{SemCont}((A, \mathcal{E}), (S, \mathcal{E}')) &= \text{SemAcc}(A, \mathcal{E}) \cap \mathcal{E}' \\ \text{SemContN}((A, \mathcal{E}), (S, \mathcal{E}')) &= \text{SemNotAcc}(A, \mathcal{E}) \cap \mathcal{E}'. \end{aligned}$$

In words, these explanations contain the common reasons for the fact in one extension and the arguments in another extension. The other extension is such that it provides proper explanations for the argument(s) in the foil.

**Example 8.** We are again interested in comparing  $\mathcal{E}_2$  and  $\mathcal{E}_4$ . However, rather than taking the difference between the two, we will now take the intersection.

- Why  $B$  in  $\mathcal{E}_4$  given not- $B$  in  $\mathcal{E}_2$ ? The acceptance of  $B$  in  $\mathcal{E}_4$  has several reasons (i.e.,  $B$  has to be defended against the attack from both  $C$  and  $G$ ), by specifying the extension with which the explanation has to be compared, we can select the specific reasons. In particular, for  $\text{Expl} = \text{NotAcc}$ :  $\text{PrfCont}((B, \mathcal{E}_4), (\{B\}, \mathcal{E}_2)) = \{H, J\}$ .
- Why  $A$  in  $\mathcal{E}_2$  given  $B$  in  $\mathcal{E}_4$ ? We take the intersection of the explanations for  $A$  and  $B$  in their respective extensions, where  $\text{Expl} = \text{Acc}$ :  $\text{PrfCont}((A, \mathcal{E}_2), (\{B\}, \mathcal{E}_4)) = \emptyset$ .

## 4. Conclusion

We have proposed a variety of contrastive explanations in the context of abstract argumentation, based on the existing humanities and social sciences literature on contrastiveness. Contrastive explanations have been studied extensively in the XAI literature [19] and some argumentative formalizations have been proposed (see, e.g., [30, 31, 32]). What is new, and different, about this paper, is that we start from the literature on contrastiveness and propose several ways in which one might turn (argumentative) explanations contrastive. To the best of our knowledge, this is the first work in which the variety of notions on contrastiveness have been formalized in one setting, for use in XAI. It is therefore a starting point for a discussion on contrastiveness in (argumentative) XAI.

With this study we pave the way for other researchers to work with contrastive explanations based on specific notions in the literature. By avoiding restrictions on the (non-)acceptance explanations (recall Notation 1), researchers outside the argumentation community can benefit from this work as well.

Now that we have general notions of contrastiveness, we and other researcher working on contrastive explanations can study when to apply which definition and how to apply them outside abstract argumentation. Moreover, we have now restricted our study to arguments, we can extend this to consider extensions, (dynamic) argumentation frameworks or elements of arguments in the structured setting.

## References

- [1] P. Lipton, Contrastive explanation, *Royal Institute of Philosophy Supplement* 27 (1990) 247–266.
- [2] K. Atkinson, T. Bench-Capon, D. Bollegala, Explanation in AI and law: Past, present and future, *Artificial Intelligence* 289 (2020) 103387.
- [3] T. Miller, Explanation in artificial intelligence: Insights from the social sciences, *Artificial Intelligence* 267 (2019) 1–38.
- [4] H. Mercier, D. Sperber, Why do humans reason? Arguments for an argumentative theory, *Behavioral and Brain Sciences* 34 (2011) 57–74.
- [5] C. Antaki, I. Leudar, Explaining in conversation: Towards an argument model, *European Journal of Social Psychology* 22 (1992) 181–194.
- [6] F. Bex, D. Walton, Combining explanation and argumentation in dialogue, *Argument & Computation* 7 (2016) 55–68.
- [7] C. Rudin, Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead, *Nature Machine Intelligence* 1 (2019) 206–215.
- [8] H. Prakken, R. Ratsma, A top-level model of case-based argumentation for explanation: formalisation and experiments, *Argument & Computation* (2021) 1–36.
- [9] D. Kazhdan, Z. Shams, P. Liò, Marleme: A multi-agent reinforcement learning model extraction library, in: *2020 International Joint Conference on Neural Networks (IJCNN)*, IEEE, 2020, pp. 1–8.
- [10] S. Timmer, J.-J. Meyer, H. Prakken, S. Renooij, B. Verheij, A two-phase method for extracting explanatory arguments from bayesian networks, *International Journal of Approximate Reasoning* 80 (2017) 475–494.
- [11] K. Čyras, A. Rago, E. Albini, P. Baroni, F. Toni, Argumentative XAI: A survey, in: Z. Zhou (Ed.), *Proceedings of IJCAI'21*, ijcai.org, 2021, pp. 4392–4399.
- [12] X. Fan, F. Toni, On computing explanations in argumentation, in: B. Bonet, S. Koenig (Eds.), *Proceedings of AAAI'15*, AAAI Press, 2015, pp. 1496–1502.
- [13] A. Arioua, M. Croitoru, Formalizing explanatory dialogues, in: C. Beierle, A. Dekhtyar (Eds.), *Scalable Uncertainty Management*, Springer, 2015, pp. 282–297.
- [14] A. Borg, F. Bex, Necessary and sufficient explanations for argumentation-based conclusions, in: *Proceedings of ECSQARU'21*, Springer, 2021, pp. 45–58.
- [15] J. Van Bouwel, E. Weber, Remote causes, bad explanations?, *Journal for the Theory of Social Behaviour* 32 (2002) 437–449.
- [16] D. Hilton, Conversational processes and causal explanation, *Psychological Bulletin* 107 (1990) 65–81.
- [17] K. Sokol, P. Flach, One explanation does not fit all: The promise of interactive explanations for machine learning transparency, *Künstliche Intelligenz* 34 (2020) 235–250.
- [18] P. Ylikoski, The idea of contrastive explanandum, in: J. Persson, P. Ylikoski (Eds.), *Rethinking Explanation*, Springer, 2007, pp. 27–42.
- [19] I. Stepin, J. Alonso, A. Catala, M. Pereira-Fariña, A survey of contrastive and counterfactual explanation generation methods for explainable artificial intelligence, *IEEE Access* 9 (2021) 11974–12001.
- [20] E. Barnes, Why P rather than Q? The curiosities of fact and foil, *Philosophical Studies* 73

- (1994) 35–53.
- [21] K. Sokol, P. Flach, Counterfactual explanations of machine learning predictions: Opportunities and challenges for AI safety, in: 2019 AAAI Workshop on Artificial Intelligence Safety (SafeAI'19), volume 2301, CEUR Workshop Proceedings, 2019.
  - [22] P. M. Dung, On the acceptability of arguments and its fundamental role in nonmonotonic reasoning, logic programming and n-person games, *Artificial Intelligence* 77 (1995) 321–357.
  - [23] P. Besnard, A. García, A. Hunter, S. Modgil, H. Prakken, G. Simari, F. Toni, Introduction to structured argumentation, *Argument & Computation* 5 (2014) 1–4.
  - [24] A. Borg, F. Bex, A basic framework for explanations in argumentation, *IEEE Intelligent Systems* 36 (2021) 25–35.
  - [25] P. Baroni, M. Caminada, M. Giacomin, Abstract argumentation frameworks and their semantics, in: P. Baroni, D. Gabay, M. Giacomin, L. van der Torre (Eds.), *Handbook of Formal Argumentation*, College Publications, 2018, pp. 159–236.
  - [26] X. Fan, F. Toni, On explanations for non-acceptable arguments, in: E. Black, S. Modgil, N. Oren (Eds.), *Proceedings of TAFE'15, LNCS 9524*, Springer, 2015, pp. 112–127.
  - [27] Z. Saribatur, J. Wallner, S. Woltran, Explaining non-acceptability in abstract argumentation, in: *Proceedings of the 24th European Conference on Artificial Intelligence (ECAI'20)*, volume 325 of *Frontiers in Artificial Intelligence and Applications*, IOS Press, 2020, pp. 881–888.
  - [28] M. Ulbricht, J. P. Wallner, Strong explanations in abstract argumentation, in: *Proceedings of the 35th AAAI Conference on Artificial Intelligence (AAAI'21)*, volume 35, 2021, pp. 6496–6504.
  - [29] T. Miller, Contrastive explanation: a structural-model approach, *The Knowledge Engineering Review* 36 (2021) e14.
  - [30] L. Amgoud, Explaining black-box classification models with arguments, in: *Proceedings of ICTAI'21*, IEEE, 2021, pp. 791–795.
  - [31] P. Besnard, S. Doutre, T. Duchatelle, M.-C. Lagasque-Schiex, Question-Based Explainability in Abstract Argumentation, Research Report IRIT/RR–2022–01–FR, IRIT : Institut de Recherche en Informatique de Toulouse, France, 2022. URL: <https://ut3-toulouseinp.hal.science/hal-03647896>.
  - [32] A. Borg, F. Bex, Contrastive explanations for argumentation-based conclusions, in: *Proceedings of AAMAS'22*, 2022, pp. 1551–1553.