

# A Basic Framework for Explanations in Argumentation

AnneMarie Borg and Floris Bex

**Abstract**—We discuss explanations for formal (abstract and structured) argumentation – the question whether and why a certain argument or claim can be accepted (or not) under various extension-based semantics. We introduce a flexible framework, which can act as the basis for many different types of explanations. For example, we can have simple or comprehensive explanations in terms of arguments for or against a claim, arguments that (indirectly) defend a claim, the evidence (knowledge base) that supports or is incompatible with a claim, and so on. We show how different types of explanations can be captured in our basic framework, discuss a real-life application and formally compare our framework to existing work.

**Index Terms**—Artificial intelligence, knowledge representation formalisms and methods, nonmonotonic reasoning and belief revision



## 1 INTRODUCTION

RECENTLY, *explainable AI* (XAI) has received much attention, mostly directed at new techniques for explaining decisions of (subsymbolic) machine learning algorithms [1]. However, explanations also play an important role in (symbolic) knowledge-based systems [2]. Argumentation is one research area in symbolic AI that is frequently mentioned in relation to XAI. For example, arguments can be used to provide reasons for or against decisions [2], [3]. The focus can also be on the argumentation itself, where it is explained whether and why a certain argument or claim can be accepted under certain semantics for computational argumentation [4], [5], [6], [7]. It is the latter type of explanations we are interested in.

Two central concepts in argumentation are *abstract argumentation frameworks* [8] – sets of arguments and the attack relations between them – and *structured or logical argumentation frameworks* (e.g., [9]) – where arguments are constructed from a knowledge base and a set of rules and the attack relation is based on the individual elements in the arguments. For both abstract and structured argumentation frameworks we can determine extensions, sets of arguments that can collectively be considered as acceptable, under different semantics [8]. In XAI terms [10], this is a *global explanation* – what can we conclude from the model as a whole? However, as argumentation is being applied in real-life AI systems with lay-users, we would rather have simpler, more compact explanations for the acceptability of individual arguments – a *local explanation* for a particular decision or conclusion. We noticed the need for such explanations when deploying an argumentation system at the Dutch National Police, which assists citizens in filing online reports and complaints [11], [12].

We propose a basic framework for explanations in struc-

tured and abstract argumentation, with which explanations for (non-)accepted arguments and (sub-)conclusions can be generated. Though some work on explanations for argumentation-based conclusions exists in the literature ([4], [5], [6], [7], Section 5), our framework is *generic* in that the underlying argumentation framework does not have to be adjusted and the definitions are semantics-independent – for example, the explanations based on the new semantics of Fan and Toni [4] are a special case of our framework. The framework is also *flexible*, as the contents of explanations can be varied. For example, rather than returning all defending or attacking arguments, we can return only those that can defend themselves, or the ones that directly attack an argument. Furthermore, we are the first to use the structure of arguments for explanations: not just arguments for a conclusion, but also elements of these arguments (e.g., premises or rules) can be returned as an explanation.

## 2 PRELIMINARIES

An *abstract argumentation framework* (AF) [8] 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, in which the nodes represent arguments and the arrows represent attacks between arguments.

*Example 1.* Consider the AF  $\mathcal{AF}_1 = \langle \text{Args}_1, \text{Att}_1 \rangle$  where  $\text{Args}_1 = \{A_1, A_2, A_3, A_4\}$  and  $\text{Att}_1 = \{(A_2, A_1), (A_3, A_2), (A_3, A_4), (A_4, A_3)\}$ .

Given an AF  $\mathcal{AF}$ , Dung-style semantics [8] can be applied to it, to determine what combinations of arguments (called *extensions*) can collectively be accepted.

**Definition 1.** Let  $\mathcal{AF} = \langle \text{Args}, \text{Att} \rangle$  be an AF,  $S \subseteq \text{Args}$  a set of arguments and let  $A \in \text{Args}$ . Then:

- $S$  *attacks*  $A$  if there is an  $A' \in S$  such that  $(A', A) \in \text{Att}$ ,  $S^+$  denotes the set of all arguments attacked by  $S$ ;
- $S$  *defends*  $A$  if  $S$  attacks every attacker of  $A$ ;

- Both authors are with the Department of Information and Computing Science, Utrecht University, The Netherlands.  
E-mail: {a.borg, f.j.bex}@uu.nl
- Floris Bex is also with the Department of Law, Technology, Markets, and Society, Tilburg University, The Netherlands.

Manuscript received August 5, 2020.

- $S$  is *conflict-free* if there are no  $A_1, A_2 \in S$  such that  $(A_1, A_2) \in \text{Att}$ ; and
- $S$  is *admissible* if it is conflict-free and it defends all of its elements.

An admissible set that contains all the arguments that it defends is a *complete extension* (cmp).

- The *grounded extension* (grd) is the minimal (with respect to  $\subseteq$ ) complete extension;
- A *preferred extension* (prf) is a maximal (with respect to  $\subseteq$ ) complete extension; and
- A *semi-stable extension* (sstb)  $S$  is a complete extension where  $S \cup S^+$  is maximal.

$\text{Ext}_{\text{sem}}(\mathcal{AF})$  denotes the set of all the extensions of  $\mathcal{AF}$  under the semantics  $\text{sem} \in \{\text{cmp}, \text{grd}, \text{prf}, \text{sstb}\}$ .

Where  $\mathcal{AF} = \langle \text{Args}, \text{Att} \rangle$  is an AF,  $\text{sem}$  a semantics and  $\text{Ext}_{\text{sem}}(\mathcal{AF}) \neq \emptyset$ , it is said that  $A \in \text{Args}$  is *skeptically* [resp. *credulously*] *accepted* if  $A \in \bigcap \text{Ext}_{\text{sem}}(\mathcal{AF})$  [resp.  $A \in \bigcup \text{Ext}_{\text{sem}}(\mathcal{AF})$ ]. These acceptability strategies are denoted by  $\bigcap$  [resp.  $\bigcup$ ].  $A$  is said to be *skeptically* [resp. *credulously*] *non-accepted* in  $\mathcal{AF}$  if for some [resp. all]  $\mathcal{E} \in \text{Ext}_{\text{sem}}(\mathcal{AF})$ ,  $A \notin \mathcal{E}$ . When these are arbitrary, result in the same or are clear from the context, we will refer to accepted respectively non-accepted arguments.

The notions of attack and defense can also be defined between arguments:

**Definition 2.** Let  $\mathcal{AF} = \langle \text{Args}, \text{Att} \rangle$  be an AF,  $A, B \in \text{Args}$  and  $\mathcal{E} \in \text{Ext}_{\text{sem}}(\mathcal{AF})$  for some sem.  $A$  can defend  $B$  directly or indirectly:  $A$  *directly defends*  $B$  if there is some  $C \in \text{Args}$  such that  $(C, B) \in \text{Att}$  and  $(A, C) \in \text{Att}$ , and  $A$  *indirectly defends*  $B$  if  $A$  defends  $C \in \text{Args}$  and  $C$  defends  $B$ . It is said that  $A$  *defends*  $B$  in  $\mathcal{E}$  if  $A$  defends  $B$  and  $A \in \mathcal{E}$ .

Similarly,  $A$  can attack  $B$  directly or indirectly:  $A$  *directly attacks*  $B$  if  $(A, B) \in \text{Att}$  and  $A$  *indirectly attacks*  $B$  if  $A$  attacks some  $C \in \text{Args}$  and  $C$  defends  $B$ .

Next we define two notions that will be used in the basic definitions of explanations. The first, used for acceptance explanations, denotes the set of arguments that defend the argument  $A$ , while the last, used for non-acceptance explanations, denotes the set of arguments that attack  $A$  and for which there is no defense in the given extension.

**Definition 3.** Let  $\mathcal{AF} = \langle \text{Args}, \text{Att} \rangle$  be an AF,  $A \in \text{Args}$  and  $\mathcal{E} \in \text{Ext}_{\text{sem}}(\mathcal{AF})$  an extension for some semantics sem.

- $\text{DefBy}(A) = \{B \in \text{Args} \mid B \text{ defends } A\}$ ;
- $\text{DefBy}(A, \mathcal{E}) = \text{DefBy}(A) \cap \mathcal{E}$  denotes the set of arguments that defend  $A$  in  $\mathcal{E}$ ;
- $\text{NotDef}(A, \mathcal{E}) = \{B \in \text{Args} \mid B \text{ attacks } A \text{ and } \mathcal{E} \text{ does not attack } B\}$ , denotes the set of all attackers of  $A$  for which no defense exists from  $\mathcal{E}$ .

*Example 2.* In  $\mathcal{AF}_1$  (recall Example 1), example conflict-free sets are  $\{A_1, A_3\}$  and  $\{A_2, A_4\}$ .  $\text{Ext}_{\text{cmp}}(\mathcal{AF}_1) = \{\emptyset, \{A_1, A_3\}, \{A_2, A_4\}\}$ , while  $\text{Ext}_{\text{prf}}(\mathcal{AF}_1) = \text{Ext}_{\text{sstb}}(\mathcal{AF}_1) = \{\{A_1, A_3\}, \{A_2, A_4\}\}$  and  $\text{Ext}_{\text{grd}}(\mathcal{AF}_1) = \{\emptyset\}$ . None of the arguments in  $\text{Args}_1$  is skeptically accepted, while all of them are credulously accepted for  $\text{sem} \in \{\text{cmp}, \text{prf}, \text{sstb}\}$ .

Argument  $A_3$  directly attacks  $A_4$ , and attacks  $A_2$  both directly and indirectly.  $A_3$  defends  $A_1$  directly against  $A_2$

and indirectly against  $A_4$ . Moreover,  $\text{DefBy}(A_1) = \{A_3\}$ ,  $\text{DefBy}(A_1, \{A_1, A_3\}) = \{A_3\}$  and  $\text{NotDef}(A_3, \{A_2, A_4\}) = \{A_4\}$ .

## 2.1 ASPIC<sup>+</sup>

We investigate explanations for a well-known approach to structured argumentation: ASPIC<sup>+</sup> [9], which allows for two types of premises – *axioms* that cannot be questioned and *ordinary premises* that can be questioned – and two types of rules – *strict* rules that cannot be questioned and *defeasible* rules. We choose ASPIC<sup>+</sup> as the structured argumentation approach in this paper since it allows to vary the form of the explanations in many ways (see Section 4). The definitions in this section are based on [9].

**Definition 4.** An *argumentation system* is a tuple  $\text{AS} = \langle \mathcal{L}, \mathcal{R}, n \rangle$ , where:

- $\mathcal{L}$  is a propositional language closed under classical negation ( $\neg$ ), we denote  $\psi = \neg\phi$  if  $\psi = \neg\phi$  or  $\phi = \neg\psi$ .
- $\mathcal{R} = \mathcal{R}_s \cup \mathcal{R}_d$  is a set of strict ( $\mathcal{R}_s$ ) and defeasible ( $\mathcal{R}_d$ ) inference rules of the form  $\phi_1, \dots, \phi_n \rightarrow \phi$  resp.  $\phi_1, \dots, \phi_n \Rightarrow \phi$ , such that  $\{\phi_1, \dots, \phi_n, \phi\} \subseteq \mathcal{L}$  and  $\mathcal{R}_s \cap \mathcal{R}_d = \emptyset$ . Where  $r \in \mathcal{R}$ ,  $\text{Ant}(r) = \{\phi_1, \dots, \phi_n\}$  are the *antecedents* of the rule and  $\text{Cons}(r) = \phi$  is the *consequent* of the rule. Moreover,  $\text{Rules}(\mathcal{R}, \phi) = \{r \in \mathcal{R} \mid \text{Cons}(r) = \phi\}$ .
- $n : \mathcal{R}_d \rightarrow \mathcal{L}$  is a naming convention for defeasible rules.

A *knowledge base* in an argumentation system  $\langle \mathcal{L}, \mathcal{R}, n \rangle$  is a set of formulas  $\mathcal{K} \subseteq \mathcal{L}$  which contains two disjoint subsets:  $\mathcal{K} = \mathcal{K}_p \cup \mathcal{K}_n$ , the set of *axioms*  $\mathcal{K}_n$  and the set of *ordinary premises*  $\mathcal{K}_p$ .

Arguments in ASPIC<sup>+</sup> are constructed in an argumentation system from a knowledge base.

**Definition 5.** An *argument*  $A$  on the basis of a knowledge base  $\mathcal{K}$  in an argumentation system  $\langle \mathcal{L}, \mathcal{R}, n \rangle$  is:

- 1)  $\phi$  if  $\phi \in \mathcal{K}$ , where  $\text{Prem}(A) = \text{Sub}(A) = \{\phi\}$ ,  $\text{Conc}(A) = \phi$ ,  $\text{Rules}(A) = \emptyset$  and  $\text{TopRule}(A) = \text{undefined}$ ;
- 2)  $A_1, \dots, A_n \rightsquigarrow \psi$ , where  $\rightsquigarrow \in \{\rightarrow, \Rightarrow\}$ , if  $A_1, \dots, A_n$  are arguments such that there exists a rule  $\text{Conc}(A_1), \dots, \text{Conc}(A_n) \rightsquigarrow \psi$  in  $\mathcal{R}_s$  if  $\rightsquigarrow = \rightarrow$  and in  $\mathcal{R}_d$  if  $\rightsquigarrow = \Rightarrow$ .  
 $\text{Prem}(A) = \text{Prem}(A_1) \cup \dots \cup \text{Prem}(A_n)$ ;  $\text{Conc}(A) = \psi$ ;  $\text{Sub}(A) = \text{Sub}(A_1) \cup \dots \cup \text{Sub}(A_n) \cup \{A\}$ ;  
 $\text{Rules}(A) = \text{Rules}(A_1) \cup \dots \cup \text{Rules}(A_n) \cup \{\text{Conc}(A_1), \dots, \text{Conc}(A_n) \rightsquigarrow \psi\}$ ;  
 $\text{DefRules}(A) = \{r \in \mathcal{R}_d \mid r \in \text{Rules}(A)\}$ ;  
 $\text{TopRule}(A) = \text{Conc}(A_1), \dots, \text{Conc}(A_n) \rightsquigarrow \psi$ .

The above notation can be generalized to sets. For example, where  $S$  is a set of arguments  $\text{Prem}(S) = \bigcup\{\text{Prem}(A) \mid A \in S\}$ ,  $\text{Conc}(S) = \{\text{Conc}(A) \mid A \in S\}$  and  $\text{DefRules}(S) = \bigcup\{\text{DefRules}(A) \mid A \in S\}$ .

*Example 3.*  $\text{AS}_2 = \langle \mathcal{L}_2, \mathcal{R}_2, n \rangle$  is an argumentation system where  $\mathcal{R}_2 = \mathcal{R}_s^2 \cup \mathcal{R}_d^2$  such that  $\mathcal{R}_s^2 = \emptyset$ ,  $\mathcal{R}_d^2 = \{d_1, \dots, d_5\}$

(the application of these rules is shown in the arguments below), let  $\mathcal{K}_2 = \mathcal{K}_n^2 \cup \mathcal{K}_p^2$  where  $\mathcal{K}_n^2 = \{t\}$  and  $\mathcal{K}_p^2 = \{r\}$ . The following arguments can be constructed:

$$\begin{array}{ll} A_1 : t & B_1 : r \\ A_2 : A_1 \xrightarrow{d_3} \neg r & B_2 : B_1 \xrightarrow{d_2} p \\ A_3 : A_1, A_2 \xrightarrow{d_4} q & B_3 : B_1 \xrightarrow{d_5} \neg q \\ A_4 : A_3 \xrightarrow{d_1} p & \end{array}$$

We denote the set of arguments constructed from  $\text{AS}_2$  and  $\mathcal{K}_2$  by  $\text{Args}_2$ . For  $A_4$  we have that  $\text{Prem}(A_4) = \{t\}$ ,  $\text{Conc}(A_4) = p$ ,  $\text{Sub}(A_4) = \{A_1, A_2, A_3, A_4\}$  and  $\text{Rules}(A_4) = \{d_1, d_3, d_4\}$ . Furthermore,  $\text{Rules}(\mathcal{R}_2, p) = \{d_1, d_2\}$ .

Attacks on an argument are based on the rules and premises applied in the construction of that argument.

**Definition 6.** An argument  $A$  attacks an argument  $B$  iff  $A$  *undercuts*, *rebuts* or *undermines*  $B$ , where:

- $A$  *undercuts*  $B$  (on  $B'$ ) iff  $\text{Conc}(A) = \neg n(r)$  for some  $B' \in \text{Sub}(B)$  such that  $B'$ 's top rule  $r$  is defeasible, it denies a rule;
- $A$  *rebuts*  $B$  (on  $B'$ ) iff  $\text{Conc}(A) = \neg \phi$  for some  $B' \in \text{Sub}(B)$  of the form  $B'_1, \dots, B'_n \Rightarrow \phi$ , it denies a conclusion;
- $A$  *undermines*  $B$  (on  $\phi$ ) iff  $\text{Conc}(A) = \neg \phi$  for some  $\phi \in \text{Prem}(B) \setminus \mathcal{K}_n$ , it denies a premise.

Argumentation theories and their corresponding Dung-style argumentation frameworks can now be defined.

**Definition 7.** An *argumentation theory* is a pair  $\text{AT} = \langle \text{AS}, \mathcal{K} \rangle$ , where  $\text{AS}$  is an argumentation system and  $\mathcal{K}$  is a knowledge base.

A *structured argumentation framework* (SAF) defined by an argumentation theory  $\text{AT}$  is a pair  $\mathcal{AF}(\text{AT}) = \langle \text{Args}, \text{Att} \rangle$ , where  $\text{Args}$  is the set of all arguments constructed from  $\text{AT}$  and  $(A, B) \in \text{Att}$  iff  $A$  attacks  $B$  according to Definition 6.

Dung-style semantics, as in Definition 1, can be applied to SAFs in the same way as they are applied in AFs.

*Example 4.* (Example 3 continued) Consider the argumentation theory  $\text{AT}_2 = \langle \text{AS}_2, \mathcal{K}_2 \rangle$ . Figure 1 contains the graphical representation of  $\mathcal{AF}(\text{AT}_2) = \langle \text{Args}_2, \text{Att}_2 \rangle$ . In this framework there are no undercuts, all the attacks from  $A_2$  are underminers and all the other attacks are rebuts.

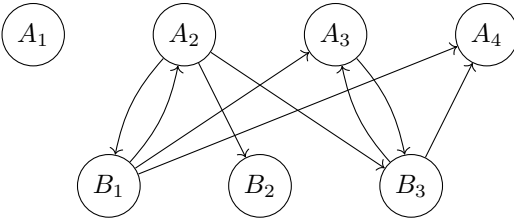


Fig. 1. Graphical representation of  $\mathcal{AF}(\text{AT}_2)$ .

Then:  $\text{Ext}_{\text{grd}}(\mathcal{AF}(\text{AT}_2)) = \{A_1\}$ ; and  $\text{Ext}_{\text{sem}}(\mathcal{AF}(\text{AT}_2)) = \{\{A_1, A_2, A_3, A_4\}, \{A_1, B_1, B_2, B_3\}\}$ , for  $\text{sem} \in \{\text{prf}, \text{sstb}\}$ .

Entailment relations, induced by the structured argumentation framework and a semantics, are defined by:

**Definition 8.** Let  $\mathcal{AF}(\text{AT}) = \langle \text{Args}, \text{Att} \rangle$  for a semantics  $\text{sem}$ ,  $\text{Ext}_{\text{sem}}(\mathcal{AF}) \neq \emptyset$  and let some  $\phi \in \mathcal{L}$ . We define:

- *Credulous entailment*:  $\text{AT} \vdash_{\text{sem}}^{\cup} \phi$  iff for some  $\mathcal{E} \in \text{Ext}_{\text{sem}}(\mathcal{AF})$  there is an argument  $A \in \mathcal{E}$  with  $\text{Conc}(A) = \phi$ , it is said that  $\phi$  is *credulously accepted*;
- *Skeptical entailment*:  $\text{AT} \vdash_{\text{sem}}^{\cap} \phi$  iff for each  $\mathcal{E} \in \text{Ext}_{\text{sem}}(\mathcal{AF})$  there is some  $A \in \mathcal{E}$  such that  $\text{Conc}(A) = \phi$ , it is said that  $\phi$  is *skeptically accepted*.

When arbitrary or clear from the context, the superscript will be omitted (e.g.,  $\vdash_{\text{grd}}$  as  $\vdash_{\text{grd}}^{\cup}$  and  $\vdash_{\text{grd}}^{\cap}$  coincide).

*Example 5* (Example 4 continued). For  $\mathcal{AF}(\text{AT}_2) = \langle \text{Args}_2, \text{Att}_2 \rangle$  we have that:

- 1)  $\text{AT}_2 \not\vdash_{\text{grd}} \phi$  and  $\text{AT}_2 \not\vdash_{\text{sem}}^{\cap} \phi$  for  $\phi \in \{q, \neg q, r, \neg r\}$ , and  $\text{sem} \in \{\text{cmp}, \text{prf}, \text{sstb}\}$ ; while
- 2)  $\text{AT}_2 \vdash_{\text{sem}}^{\cup} \phi$  for any  $\phi \in \{p, q, \neg q, r, \neg r, t\}$  and  $\text{sem} \in \{\text{cmp}, \text{prf}, \text{sstb}\}$ ;
- 3)  $\text{AT}_2 \vdash_{\text{grd}} t$  and  $\text{AT}_2 \vdash_{\text{sem}}^{\cap} t$  for  $\text{sem} \in \{\text{cmp}, \text{prf}, \text{sstb}\}$ ; and
- 4)  $\text{AT}_2 \vdash_{\text{sem}}^{\cap} p$  for  $\text{sem} \in \{\text{prf}, \text{sstb}\}$  but  $\text{AT}_2 \not\vdash_{\text{grd}} p$ .

This follows since each argument from  $\text{Args}_2$  is part of at least one extension, but only  $A_1$  is part of every extension. The last item follows since each  $\text{sem}$ -extension of  $\mathcal{AF}(\text{AT}_2)$  contains either  $A_4$  or  $B_2$  for  $\text{sem} \in \{\text{prf}, \text{sstb}\}$ .

## 2.2 Necessary Notation

This notation is meant to keep the definitions of explanations in Section 3 general and short.

*Notation 1.* Let  $\mathcal{AF} = \langle \text{Args}, \text{Att} \rangle$  be an AF,  $A \in \text{Args}$  and  $S \subseteq \text{Args}$ . Then, for some  $\text{sem} \in \{\text{grd}, \text{cmp}, \text{prf}, \text{sstb}\}$ :

- $\mathcal{E}_A^{\text{sem}} = \{\mathcal{E} \in \text{Ext}_{\text{sem}}(\mathcal{AF}) \mid A \in \mathcal{E}\}$  denotes the set of  $\text{sem}$ -extensions of  $\mathcal{AF}$  which contain  $A$ ;
- $\mathcal{E}_X^{\text{sem}} = \{\mathcal{E} \in \text{Ext}_{\text{sem}}(\mathcal{AF}) \mid A \notin \mathcal{E}\}$  denotes the set of  $\text{sem}$ -extensions of  $\mathcal{AF}$  which do not contain  $A$ .

The set of arguments that can be used to explain the acceptance of a formula differs depending on the acceptance strategy. For this the following notation will be applied.

*Notation 2.* Let  $\mathcal{AF}(\text{AT}) = \langle \text{Args}, \text{Att} \rangle$  be an SAF,  $\phi \in \mathcal{L}$  and let  $\text{sem} \in \{\text{grd}, \text{cmp}, \text{prf}, \text{sstb}\}$ . Then:

- $\text{Args}_{\phi} = \{A \in \text{Args} \mid \text{Conc}(A) = \phi\}$  denotes the set of all arguments of  $\mathcal{AF}(\text{AT})$  with conclusion  $\phi$ ;
- $\text{Args}_{\phi}^{\text{sem}, \cup} = \{A \in \cup \text{Ext}_{\text{sem}}(\mathcal{AF}(\text{AT})) \mid \text{Conc}(A) = \phi\}$  denotes the set of all arguments of  $\mathcal{AF}(\text{AT})$  with conclusion  $\phi$  that are part of at least one  $\text{sem}$ -extension (i.e., that are credulously accepted);
- $\text{Args}_{\phi}^{\text{sem}, \cap} = \begin{cases} \emptyset & \text{if } \text{AT} \not\vdash_{\text{sem}}^{\cap} \phi \\ \text{Args}_{\phi}^{\text{sem}, \cup} & \text{otherwise} \end{cases}$   
is the same as  $\text{Args}_{\phi}^{\text{sem}, \cup}$  if  $\phi$  is skeptically accepted and  $\emptyset$  if it is not skeptically accepted.

*Example 6.* (Example 4 continued) Whenever  $\text{Args}_p^{\text{sem}, \cap} \neq \emptyset$ , there is no difference between  $\cup$  and  $\cap$ . But  $\text{Args}_q = \text{Args}_q^{\text{sem}, \cup} = \{A_3\}$  while  $\text{Args}_q^{\text{sem}, \cap} = \emptyset$  for  $\text{sem} \in \{\text{cmp}, \text{prf}, \text{sstb}\}$ .

Next it is defined what it means for two formulas to be connected in an argumentation system.

**Definition 9.** Let  $AS = \langle \mathcal{L}, \mathcal{R}, n \rangle$  be an argumentation system. Then,  $\phi$  is *connected* to  $\psi$  if  $\phi = \psi$ , or:

- there is some  $r \in \mathcal{R}$  with  $\text{Cons}(r) = \psi$  and  $\phi \in \text{Ant}(r)$ ;
- there is some  $\gamma \in \mathcal{L}$  such that  $\phi$  is connected to  $\gamma$  and  $\gamma$  is connected to  $\psi$ .

The set of all connected formulas of  $\psi$  is denoted by:

- $\text{Connected}(\psi) = \{\phi \in \mathcal{L} \mid \phi \text{ is connected to } \psi\}$ .

In explanations for formulas for which no argument exists the following notation will be used:

*Notation 3.* Let  $\mathcal{AF}(AT) = \langle \text{Args}, \text{Att} \rangle$  be an SAF and let  $\phi \in \mathcal{L}$  be such that there is no argument for it in  $\text{Args}$ . Then:

- $\text{NoArgAnt}(\phi) = \{\psi \mid \psi \in \bigcup \{\text{Ant}(r) \mid r \in \text{Rules}(\mathcal{R}, \phi)\} \text{ and } \nexists A \in \text{Args s.t. } \text{Conc}(A) = \psi\}$  denotes the set of formulas in antecedents of rules for  $\phi$  for which no argument exists.
- $\text{NoArgPrem}(\phi) = \{\psi \in \text{Connected}(\phi) \mid \text{Rules}(\mathcal{R}, \psi) = \emptyset \text{ and } \psi \notin \mathcal{K}\}$  denotes the set of formulas that are connected to  $\phi$  but that are not part of  $\mathcal{K}$  and for which no rules exist.

Intuitively,  $\text{NoArgAnt}$  determines the formulas for which arguments are missing in order for an argument for  $\phi$  to be available, while  $\text{NoArgPrem}$  determines the formulas that are not derivable from  $\mathcal{AF}(AT)$  (neither from  $\mathcal{K}$  nor as a conclusion of some rule) and which could be part of the derivation of an argument for  $\phi$ .

*Example 7.* Consider  $AS_2$  from Example 3, but let  $\mathcal{K}'_2 = \mathcal{K}_p^2$  (i.e.,  $\mathcal{K}'_n = \emptyset$ ). It follows that the arguments  $A_1, A_2, A_3$  and  $A_4$  no longer exist. Thus there is no argument for  $\neg r$  nor for  $q$  (though there is still an argument for  $p$ :  $B_2$ ). We have that:  $\text{NoArgAnt}(q) = \{t, \neg r\}$ ,  $\text{Connected}(q) = \{t, \neg r\}$  and  $\text{NoArgPrem}(q) = \{t\}$ .

### 3 BASIC EXPLANATIONS

We now define basic explanations in terms of two functions.  $\mathbb{D}$  determines the *depth* of the explanation, how “far away” we should look when considering attacking and defending arguments as explanations.  $\mathbb{F}$  determines the *form* of the explanation, whether we want, for example, an argument as an explanation or only its premises. A formal definition of these functions is not provided since domain ( $\mathbb{F}$ ) and codomain ( $\mathbb{D}$  and  $\mathbb{F}$ ) are not fixed. We will sometimes use the superscripts  $\text{acc}$  and  $\text{na}$  to denote the function used in the context of acceptance [resp. non-acceptance] explanations.

See the online appendix for an algorithm that computes the basic explanations.

#### 3.1 Basic Explanations for Acceptance

We define two types of acceptance explanations, where  $\cap$ -explanations provide all the reasons why an argument or formula can be accepted by a skeptical reasoner, while  $\cup$ -explanations provide one reason why an argument or formula can be accepted by a credulous reasoner. For the purpose of this section let  $\mathbb{D}^{\text{acc}}(A, S) = \text{DefBy}(A, S)$  and  $\mathbb{F}^{\text{acc}}(T) = \text{id}(T) = T$  (i.e.,  $\text{id}(S) = S$  for any set  $S$ ).

##### 3.1.1 Explanations for Accepted Arguments

An argument explanation for an accepted argument  $A$  consists of the arguments that defend it, depending on the extensions considered according to the acceptability strategy.

**Definition 10** (Argument explanation). Let  $\mathcal{AF} = \langle \text{Args}, \text{Att} \rangle$  be an AF and let  $A \in \text{Args}$  be an accepted argument, given some  $\text{sem} \in \{\text{cmp}, \text{grd}, \text{prf}, \text{sstb}\}$  and an acceptance strategy ( $\cap$  or  $\cup$ ). Then:

$$\text{Acc}_{\text{sem}}^{\cap}(A) = \bigcup_{\mathcal{E} \in \text{Ext}_{\text{sem}}(\mathcal{AF})} \mathbb{D}^{\text{acc}}(A, \mathcal{E});$$

$$\text{Acc}_{\text{sem}}^{\cup}(A) \in \{\mathbb{D}^{\text{acc}}(A, \mathcal{E}) \mid \mathcal{E} \in \mathfrak{E}_A^{\text{sem}}\}.$$

$\text{Acc}_{\text{sem}}^{\cap}(A)$  provides for each  $\text{sem}$ -extension  $\mathcal{E}$  the arguments that defend  $A$  in  $\mathcal{E}$ , and  $\text{Acc}_{\text{sem}}^{\cup}(A)$  the arguments that defend  $A$  in one of the  $\text{sem}$ -extensions.

*Example 8* (Example 2 continued). Recall  $\mathcal{AF}_1 = \langle \text{Args}_1, \text{Att}_1 \rangle$ . We have that:

- $\text{Acc}_{\text{prf}}^{\cup}(A_2) = \{A_4\}$ ;
- $\text{Acc}_{\text{prf}}^{\cup}(A_3) = \{A_3\}$ .

##### 3.1.2 Explanations for Accepted Formulas

In structured argumentation explanations for the acceptance of a formula  $\phi$  can be requested, in addition to argument explanations. For  $\phi$  to be accepted, at least one argument for  $\phi$  must exist. Therefore, the existence of such an argument is part of the explanation as well.

**Definition 11** (Formula explanation). Let  $\mathcal{AF}(AT) = \langle \text{Args}, \text{Att} \rangle$  be an SAF and let  $\phi \in \mathcal{L}$  be such that  $AT \sim_{\text{sem}}^{\star} \phi$ , for  $\text{sem} \in \{\text{cmp}, \text{grd}, \text{prf}, \text{sstb}\}$  and  $\star \in \{\cap, \cup\}$ . Here  $S = \text{Args}_{\phi}^{\text{sem}, \cap}$ ,  $A \in \text{Args}_{\phi}^{\text{sem}, \cup}$  and  $S_A \in \{\mathbb{D}^{\text{acc}}(A, \mathcal{E}) \mid \mathcal{E} \in \mathfrak{E}_A^{\text{sem}}\}$ :

$$\text{Acc}_{\text{sem}}^{\cap}(\phi) = \left\langle \mathbb{F}^{\text{acc}}(S), \mathbb{F}^{\text{acc}} \left( \bigcup_{B \in S} \bigcup_{\mathcal{E} \in \mathfrak{E}_B^{\text{sem}}} \mathbb{D}^{\text{acc}}(B, \mathcal{E}) \right) \right\rangle;$$

$$\text{Acc}_{\text{sem}}^{\cup}(\phi) = \left\langle \mathbb{F}^{\text{acc}}(A), \mathbb{F}^{\text{acc}}(S_A) \right\rangle.$$

The first part of the explanation denotes arguments for  $\phi$  (recall Notation 2) – all arguments in the case of  $\text{Acc}_{\text{sem}}^{\cap}(\phi)$  and one argument in the case of  $\text{Acc}_{\text{sem}}^{\cup}(\phi)$ . The second part of the explanation is similar to the set of arguments in an argument explanation, although now the function  $\mathbb{F}$  is applied to it. This makes it possible to change the form of the explanation (e.g., premises instead of arguments). The main difference with argument explanations is that more than one argument for  $\phi$  may be considered in the  $\cap$ -explanation. The (skeptical)  $\cap$ -explanation again takes all extensions in  $\mathfrak{E}_B^{\text{sem}}$  into account to determine the arguments that defend  $B$ , while for the (credulous)  $\cup$ -explanation again the defending arguments for  $A$  from just one extension in  $\mathfrak{E}_A^{\text{sem}}$  are taken.

*Example 9.* (Example 5 continued) Consider the SAF  $\mathcal{AF}(AT_2)$  for  $AT_2 = \langle AS_2, \mathcal{K}_2 \rangle$ . Recall that  $AT_2 \sim_{\text{prf}}^{\cap} p$ , hence:

- $\text{Acc}_{\text{prf}}^{\cap}(p) = \langle \{A_4, B_2\}, \{A_2, A_3, B_1\} \rangle$ .

For other formulas the  $\text{Acc}_{\text{sem}}^{\cap}$ -explanation does not apply, since none of these are skeptically accepted. However:

- $\text{Acc}_{\text{prf}}^{\cup}(q) = \langle \{A_3\}, \{A_2, A_3\} \rangle$ ;
- $\text{Acc}_{\text{prf}}^{\cup}(\neg q) = \langle \{B_3\}, \{B_1, B_3\} \rangle$ .

### 3.2 Basic Explanations for Non-Acceptance

Similar to acceptance explanations, there are two types of non-acceptance explanations:  $\cap$ -explanations for why an argument or formula is not accepted in some extensions (i.e., is not skeptically accepted), and  $\cup$ -explanations for why an argument or formula is not accepted in all extensions (i.e., is not credulously accepted). For this let  $\mathbb{D}^{\text{na}}(A, S) = \text{NotDef}(A, S)$  and  $\mathbb{F}^{\text{na}}(T) = \text{id}(T) = T$ .

#### 3.2.1 Explanations for Non-Accepted Arguments

In any Dung-style semantics based on the complete semantics, an argument is not accepted if it is attacked and it is not defended by an accepted argument. Hence, intuitively, the explanation for the non-acceptance of an argument is the set of arguments for which no defense exists.

**Definition 12** (Non-acceptance argument explanation). Let  $\mathcal{AF} = \langle \text{Args}, \text{Att} \rangle$  be an AF and let  $A \in \text{Args}$  be an argument that is not accepted, given some  $\text{sem} \in \{\text{cmp}, \text{grd}, \text{prf}, \text{sstb}\}$  and some  $\star \in \{\cap, \cup\}$ . Then:

$$\begin{aligned} \text{NotAcc}_{\text{sem}}^{\cap}(A) &= \bigcup_{\mathcal{E} \in \mathcal{E}_{\text{sem}}^{\cap}} \mathbb{D}^{\text{na}}(A, \mathcal{E}); \\ \text{NotAcc}_{\text{sem}}^{\cup}(A) &= \bigcup_{\mathcal{E} \in \text{Ext}_{\text{sem}}(\mathcal{AF})} \mathbb{D}^{\text{na}}(A, \mathcal{E}). \end{aligned}$$

So the non-acceptance argument explanation contains all the arguments in  $\text{Args}$  that attack  $A$  and for which no defense exists in: some sem-extensions (for  $\cap$ ) of which  $A$  is not a member; all sem-extensions (for  $\cup$ ). That for  $\cap$  only some extensions have to be considered follows since  $A$  is skeptically non-accepted as soon as  $\mathcal{E}_{\text{sem}}^{\cap} \neq \emptyset$ , while  $A$  is credulously non-accepted when  $\mathcal{E}_{\text{sem}}^{\cup} = \text{Ext}_{\text{sem}}(\mathcal{AF})$ .

*Example 10.* (Example 5 continued) Recall  $\mathcal{AF}(\text{AT}_2)$ . Then:

- $\text{NotAcc}_{\text{grd}}^{\cup}(A_3) = \{B_1, B_3\}$ ;
- $\text{NotAcc}_{\text{prf}}^{\cup}(B_3) = \{A_2, A_3\}$ .

#### 3.2.2 Explanations for Non-Accepted Formulas

The non-acceptance of a formula  $\phi$  can have two causes: either there is no argument for  $\phi$  at all (i.e., it is not derivable) or all arguments for  $\phi$  are attacked. In the first case  $\phi$  is not part of the knowledge base  $\mathcal{K}$ . Moreover, if there are rules with  $\phi$  as consequent, for each rule there is at least one antecedent for which no argument exists.

**Definition 13** (Non-derivability explanation). Let  $\mathcal{AF}(\text{AT})$  be an SAF and let  $\phi$  be some non-derivable formula. Then:

$$\begin{aligned} \text{NotDer}(\phi) &= \langle \text{Rules}(\mathcal{R}, \phi), \\ &\quad \text{NoArgAnt}(\phi), \text{NoArgPrem}(\phi) \rangle. \end{aligned}$$

The idea is that the explanation points out the gaps in the argumentation theory: the missing knowledge base elements and/or missing rules. If there are rules for  $\phi$  these are collected in the first part of the explanation, the second part contains the missing antecedents of these rules (if there would be arguments for all antecedents, there would be an argument for  $\phi$ ) and the third part contains the formulas that are connected to  $\phi$  but for which no rule exists (i.e., formulas which are neither part of the knowledge base nor the consequent of a rule).

*Example 11.* (Example 7 continued) Consider again  $\text{AS}_2$  from Example 3, with the knowledge base  $\mathcal{K}'_2$  from Example 7 (i.e.,  $\mathcal{K}'_2 = \mathcal{K}_2 \setminus \{t\}$ ). There are no arguments for  $\neg r$  and  $q$ :

- $\text{NotDer}(\neg r) = \langle \{d_3\}, \{t\}, \{t\} \rangle$ ;
- $\text{NotDer}(q) = \langle \{d_4\}, \{t, \neg r\}, \{t\} \rangle$ .

This follows since, although there is a rule for  $q$  (i.e.,  $d_4 \in \mathcal{R}_d^2$ ) [resp. for  $\neg r$  (i.e.,  $d_3 \in \mathcal{R}_d^2$ )], there is some  $\psi \in \text{Ant}(d_4)$  [resp.  $\psi \in \text{Ant}(d_3)$ ] (i.e.,  $\psi = t$  [resp.  $\psi = \neg r$ ]) such that there is no argument for  $t$  [resp.  $\neg r$ ] in  $\mathcal{AF}(\text{AT}_2)$  and when looking at the missing premises to derive  $q$  [resp.  $\neg r$ ] the formula  $t$ , necessary for  $d_3$  is found.

Like for non-acceptance argument explanations, if an argument for  $\phi$  exists but it is not accepted, there has to be an attacker for which there is no defense.

**Definition 14** (Non-acceptance formula explanation). Let  $\mathcal{AF}(\text{AT}) = \langle \text{Args}, \text{Att} \rangle$  be an SAF and let  $\phi \in \mathcal{L}$  be such that  $\text{AT} \not\vdash_{\text{sem}}^{\star} \phi$ , given some  $\text{sem} \in \{\text{cmp}, \text{grd}, \text{prf}, \text{sstb}\}$  and  $\star \in \{\cap, \cup\}$ . Here,  $S_{\phi} = \text{Args}_{\phi}$ .

$$\begin{aligned} \text{NotAcc}_{\text{sem}}^{\cap}(\phi) &= \left\langle \mathbb{F}^{\text{na}}(S_{\phi}), \mathbb{F}^{\text{na}} \left( \bigcup_{A \in S_{\phi}} \bigcup_{\mathcal{E} \in \mathcal{E}_{\text{sem}}^{\cap}} \mathbb{D}^{\text{na}}(A, \mathcal{E}) \right) \right\rangle \\ \text{NotAcc}_{\text{sem}}^{\cup}(\phi) &= \left\langle \mathbb{F}^{\text{na}}(S_{\phi}), \mathbb{F}^{\text{na}} \left( \bigcup_{A \in S_{\phi}} \bigcup_{\mathcal{E} \in \text{Ext}_{\text{sem}}(\mathcal{AF})} \mathbb{D}^{\text{na}}(A, \mathcal{E}) \right) \right\rangle. \end{aligned}$$

These explanations consist of the existing arguments for  $\phi$  and the arguments for which no defense exists from  $\mathcal{E}$  under  $\mathbb{D}^{\text{na}}$ . Similar to non-acceptance argument explanations, for  $\cap$  only the extensions without any argument for  $\phi$  have to be considered, while for  $\cup$  all extensions have to be accounted for. By assumption  $S_{\phi} \neq \emptyset$ , since otherwise the explanation for the non-acceptance of  $\phi$  would be its non-derivability.

*Example 12.* (Example 9 continued) Consider again  $\mathcal{AF}(\text{AT}_2)$ . Recall that all arguments are credulously accepted, we do however have:

- $\text{NotAcc}_{\text{prf}}^{\cup}(q) = \langle \{A_3\}, \{B_1, B_3\} \rangle$ ;
- $\text{NotAcc}_{\text{prf}}^{\cup}(\neg q) = \langle \{B_3\}, \{A_2, A_3\} \rangle$ .

## 4 VARYING $\mathbb{D}$ AND $\mathbb{F}$

This section proposes several variations for  $\mathbb{D}$  and  $\mathbb{F}$ , the main purpose of which is to show the flexibility of the basic framework. We focus on notions of defense, which are suitable for the completeness-based semantics in this paper. For, for example, naive semantics, one might want to base  $\mathbb{D}$  on conflicts instead. In Section 4.4 these variations are discussed in the context of a real-life application.

### 4.1 Notions of Defense

We start by only considering the arguments that defend themselves against all attacks.

**Definition 15.** Let  $\mathcal{AF} = \langle \text{Args}, \text{Att} \rangle$  be an AF,  $A, B \in \text{Args}$  and let  $\mathcal{E} \in \text{Ext}_{\text{sem}}(\mathcal{AF})$  for some semantics sem. Then:

$$\begin{aligned} \text{FinalDef}(A, \mathcal{E}) = \{ & B \in \text{DefBy}(A, \mathcal{E}) \mid \forall C \in \text{Args s.t. } (C, B) \\ & \in \text{Att}, (B, C) \in \text{Att} \} \cup \bigcup \{ \text{DefBy}(B, \mathcal{E}) \mid B \in \text{DefBy}(A, \mathcal{E}), \\ & \forall C \in \text{DefBy}(B, \mathcal{E}), \text{DefBy}(C, \mathcal{E}) = \text{DefBy}(B, \mathcal{E}) \text{ and } \nexists D \\ & \in \text{DefBy}(B, \mathcal{E}) \text{ s.t. } \forall E \in \text{Args s.t. } (E, D) \in \text{Att}, (D, E) \in \text{Att} \} \end{aligned}$$

denotes the set of arguments that defend  $A$  in  $\mathcal{E}$  and that are not attacked at all, defend themselves against any attacker or are part of an even cycle that is not attacked.

Intuitively this means that these arguments that defend  $A$  do not need other arguments to be defended and, given  $\mathcal{E}$ , can be considered as safe to be accepted. To see why even cycles should be regarded, take a look at the following example:

*Example 13.* (Figure 2(a)) Note that  $\text{Ext}_{\text{grd}}(\mathcal{AF}_3) = \emptyset$ , while  $\text{Ext}_{\text{sem}}(\mathcal{AF}_3) = \{ \{A, D, F, H\}, \{A, D, F, I\}, \{B, C, E, H\}, \{B, C, E, I\} \}$  for  $\text{sem} \in \{\text{prf}, \text{sstb}\}$ . Let  $\mathcal{E} = \{A, D, F, H\}$ . Then  $\text{FinalDef}(F, \mathcal{E}) = \{A, D, H\}$ . This follows since  $H$  defends itself against the attack from  $I$  and  $\{A, D\}$  is part of an even cycle that is not attacked. If even cycles would not be covered by  $\text{FinalDef}$ , the defense of the attack  $(E, F)$  would not be accounted for.

Another option is to consider only the arguments that directly defend the considered argument.

**Definition 16.** Let  $\mathcal{AF} = \langle \text{Args}, \text{Att} \rangle$  be an AF,  $A, B \in \text{Args}$  and let  $\mathcal{E} \in \text{Ext}_{\text{sem}}(\mathcal{AF})$  for some semantics sem. Then:  $\text{DirDef}(A, \mathcal{E}) = \{B \in \mathcal{E} \mid B \text{ directly defends } A\}$ , denotes the set of arguments in  $\mathcal{E}$  that *directly* defend  $A$ .

One reason for looking at direct conflicts might be that direct conflicts are often more clear from the context than indirect conflicts.

*Example 14.* (Figure 2(c)). Here  $\text{Ext}_{\text{sem}}(\mathcal{AF}_4) = \{ \{A_1, A_3, A_5\} \}$  for any  $\text{sem} \in \{\text{grd}, \text{cmp}, \text{prf}, \text{sstb}\}$ . Moreover:

- $\text{Acc}(A_1) = \{A_3, A_5\}$  for  $\mathbb{D} = \text{DefBy}$ ;
- $\text{Acc}(A_1) = \{A_5\}$  for  $\mathbb{D} = \text{FinalDef}$ ; and
- $\text{Acc}(A_1) = \{A_3\}$  for  $\mathbb{D} = \text{DirDef}$ .

This minimal example can be seen as a discussion in the form of a sequence of arguments attacking and defending the topic  $A_1$ . When at the end an explanation for the acceptance of  $A_1$  is requested:  $\text{DefBy}$  returns all arguments that defend  $A_1$ ;  $\text{FinalDef}$  returns the last argument that was put forward, which is uncontested; and  $\text{DirDef}$  returns the argument against the direct attacker of the topic.

*Example 15.* (Example 9 continued) Consider  $\mathcal{AF}(\text{AT}_2)$ . Then, for  $\mathbb{F}^{\text{acc}} = \text{id}$ :

- $\text{Acc}_{\text{prf}}^{\square}(p) = \langle \{A_4, B_2\}, \{A_2, A_3, B_1\} \rangle$ , for  $\mathbb{D}^{\text{acc}} = \text{DirDef}$ ;
- $\text{Acc}_{\text{prf}}^{\square}(p) = \langle \{A_4, B_2\}, \{A_2, B_1\} \rangle$ , for  $\mathbb{D}^{\text{acc}} = \text{FinalDef}$ .

In the case of non-acceptance explanations,  $\mathbb{D}$  was defined as the set of all attacking arguments against which no defense exists. The next definition considers only those attackers that  $A$  does not (in)directly attack itself.

**Definition 17.** Let  $\mathcal{AF} = \langle \text{Args}, \text{Att} \rangle$  be an AF,  $A, B \in \text{Args}$  and let  $\mathcal{E} \in \text{Ext}_{\text{sem}}(\mathcal{AF})$  be an extension for some semantics sem. Then:  $\text{NoDir}(A, \mathcal{E}) = \{B \in \text{NotDef}(A, \mathcal{E}) \mid A \text{ does not (in)directly attack } B\}$  denotes the set of arguments that attack  $A$  for which no defense exists in  $\mathcal{E}$  and which are not attacked by  $A$  itself.

Intuitively, the members of  $\text{NoDir}(A, \mathcal{E})$  attack  $A$  but in order to defend  $A$  against the attack another argument than  $A$  itself is necessary.

*Example 16.* Let  $\mathcal{AF}_5 = \langle \{A, B\}, \{(A, B), (B, A)\} \rangle$ . Here  $\text{Ext}_{\text{prf}}(\mathcal{AF}_5) = \{ \{A\}, \{B\} \}$ ,  $\text{NotAcc}^{\square}(A) = \{B\}$  for  $\mathbb{D} = \text{NotDef}$  but  $\text{NotAcc}_{\text{prf}}^{\square}(A) = \emptyset$  for  $\mathbb{D} = \text{NoDir}$  since by accepting  $A$ ,  $A$  can indeed be concluded. Now let  $\mathcal{AF}'_5$  as in Figure 2(d). Then  $\text{Ext}_{\text{prf}}(\mathcal{AF}'_5) = \{ \{A, D\}, \{B, C\}, \{B, D\} \}$ ,  $\text{NotAcc}_{\text{prf}}^{\square}(A) = \{B, C\}$  for  $\mathbb{D} = \text{NotDef}$  and  $\text{NotAcc}_{\text{prf}}^{\square}(A) = \{C\}$  for  $\mathbb{D} = \text{NoDir}$ , since in order to defend  $A$ , just accepting  $A$  is not enough,  $D$  is needed to defend against the attack from  $C$ .

*Example 17.* (Example 12 continued) Consider  $\mathcal{AF}(\text{AT}_2)$  from Example 3. Then, for  $\mathbb{F}^{\text{acc}} = \text{id}$  and  $\mathbb{D}^{\text{na}} = \text{NoDir}$ :

- $\text{NotAcc}_{\text{prf}}^{\square}(q) = \langle \{A_3\}, \{B_1\} \rangle$ ;
- $\text{NotAcc}_{\text{prf}}^{\square}(\neg q) = \langle \{B_3\}, \{A_2\} \rangle$ .

## 4.2 Element Explanations

In structured argumentation, one can provide full arguments as the explanation (e.g.,  $\mathbb{F} = \text{id}$ ), but the structure of the arguments provides other possibilities as well.

**Definition 18.** Let  $\mathcal{AF}(\text{AT}) = \langle \text{Args}, \text{Att} \rangle$  be an SAF and  $S \subseteq \text{Args}$  a set of formulas. Then  $\text{AntTop}(S) = \{ \text{Ant}(\text{TopRule}(A)) \mid A \in S \}$  denotes the set of antecedents of the top rule of all arguments in  $S$ .

The above definition, combined with the introduced notation in Definition 5, provides some ideas of how  $\mathbb{F}$  can be defined. For example, explanations in terms of premises explain the conclusion in terms of knowledge base items. The notion  $\text{AntTop}$  provides explanations in terms of closely related information and the rule with which the conclusion is derived from that information.

*Example 18.* (Examples 9 and 12 continued) Consider  $\mathcal{AF}(\text{AT}_2)$  from Example 3. Then, for  $\mathbb{D}^{\text{acc}} = \text{DefBy}$  and  $\mathbb{D}^{\text{na}} = \text{NotDef}$ :

- $\text{Acc}_{\text{prf}}^{\square}(p) = \langle \{t, r\}, \{t, r\} \rangle$  for  $\mathbb{F}^{\text{acc}} = \text{Prem}$ ;
- $\text{Acc}_{\text{prf}}^{\square}(p) = \langle \{q, r\}, \{t, \neg r\} \rangle$  for  $\mathbb{F}^{\text{acc}} = \text{AntTop}$ ;
- $\text{NotAcc}_{\text{prf}}^{\cup}(q) = \langle \{t\}, \{r\} \rangle$  for  $\mathbb{F}^{\text{na}} = \text{Prem}$ ;
- $\text{NotAcc}_{\text{prf}}^{\cup}(q) = \langle \{\neg r, t\}, \{r\} \rangle$  for  $\mathbb{F}^{\text{na}} = \text{AntTop}$ .

## 4.3 Comparing the Size of Explanations

When choosing a definition for  $\mathbb{D}$  and  $\mathbb{F}$  the size of the resulting explanation might be one of the considerations. While for  $\mathbb{F}$  this depends on the AF (e.g., an argument might have many premises or the top rule might have only one antecedent), for  $\mathbb{D}$  the size of the different definitions can be compared. We will apply  $\leq$  to the size of the sets, i.e.,  $S_1 \leq S_2$  denotes  $|S_1| \leq |S_2|$ .

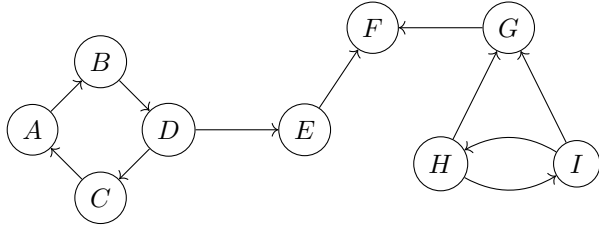
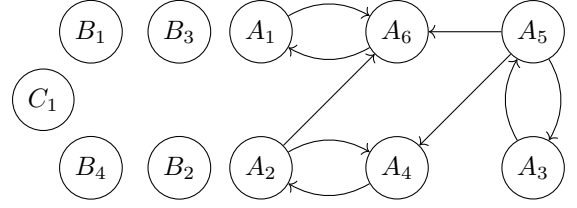
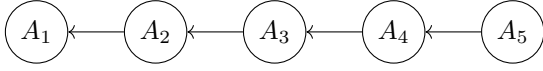
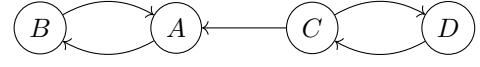
(a)  $\mathcal{AF}_3$ , Example 13.(b)  $\mathcal{AF}(\text{AT}_6)$ , Example 19.(c)  $\mathcal{AF}_4$ , Example 14.(d)  $\mathcal{AF}'_5$ , Example 16.

Fig. 2. Graphical representations of the AFs in Section 4.

**Proposition 1.** Let  $\mathcal{AF} = \langle \text{Args}, \text{Att} \rangle$  be an AF, let  $A \in \text{Args}$  and let  $\mathcal{E} \in \text{Ext}_{\text{sem}}(\mathcal{AF})$  be an extension for it. Where  $\preceq \in \{\leq, \subseteq\}$ :

- 1)  $\text{DirDef}(A, \mathcal{E}) \preceq \text{DefBy}(A, \mathcal{E})$ ;
- 2)  $\text{FinalDef}(A, \mathcal{E}) \preceq \text{DefBy}(A, \mathcal{E})$ ;
- 3)  $\text{NoDir}(A, \mathcal{E}) \preceq \text{NotDef}(A, \mathcal{E})$ .

This follows since  $\text{DirDef}(A, \mathcal{E})$  and  $\text{FinalDef}(A, \mathcal{E})$  are always subsets of  $\text{DefBy}(A, \mathcal{E})$  and  $\text{NoDir}(A, \mathcal{E})$  is always a subset of  $\text{NotDef}(A, \mathcal{E})$ . Indeed,  $\text{Acc}_{\text{prf}}^{\cap}(p)$  is both  $\leq$ - and  $\subseteq$ -smaller for  $\mathbb{D}^{\text{acc}} = \text{DirDef}$  than for  $\mathbb{D}^{\text{acc}} = \text{DefBy}$  (see Example 15). Similarly,  $\text{NotAcc}_{\text{prf}}^{\cap}(q)$  is  $\leq$ - and  $\subseteq$ -smaller for  $\mathbb{D}^{\text{na}} = \text{NoDir}$  than for  $\mathbb{D}^{\text{na}} = \text{NotDef}$  (see Example 17).

#### 4.4 Applying the Basic Framework

One of the inspirations for this paper is an argumentation-based system in use by the Dutch National Police, which assists citizens who might have been the victim of internet trade fraud (e.g., malicious web shops or traders) in filing a criminal report [11], [12]. From this report basic observations such as ‘money was paid by the complainant to the counterparty’ or ‘no package was delivered to the complainant’ are collected, and these observations are used as premises in legal arguments to infer whether or not the report concerns a possible case of fraud. This conclusion is then provided to the complainant who filed the report. The system is based on ASPIC<sup>+</sup> [9], with axioms (the observations) and defeasible rules (based on Dutch law concerning fraud), and all attacks are rebuts. The next example illustrates such an argumentation framework.

*Example 19.* Let  $\text{AS}_6 = \langle \mathcal{L}_6, \mathcal{R}_6, n \rangle$  be an argumentation system, where  $\mathcal{L}_6$  contains the propositions  $p$  (the complainant paid),  $w$  (the wrong package arrived),  $fk$  (the product is fake),  $su$  (the product looks suspicious),  $re$  (counterparty states that the product is real),  $cd$  (the complainant delivered),  $cpd$  (the counterparty delivered) and  $f$  (it is fraud) and their negations

and where  $\mathcal{R}_6$  is such that the following arguments can be derived from  $\mathcal{K}_6 = \mathcal{K}_n^6 = \{p, w, su, re\}$ :

$$\begin{array}{lll} B_1 : p & C_1 : B_1 \Rightarrow cd & \\ B_2 : w & A_1 : B_2 \Rightarrow \neg f & A_4 : A_3 \Rightarrow \neg cpd \\ B_3 : su & A_2 : B_2 \Rightarrow cpd & A_5 : B_4 \Rightarrow \neg fk \\ B_4 : re & A_3 : B_3 \Rightarrow fk & A_6 : C_1, A_4 \Rightarrow f \end{array}$$

Figure 2(b) shows the corresponding SAF  $\mathcal{AF}(\text{AT}_6)$ . The preferred extensions of  $\mathcal{AF}(\text{AT}_6)$ , only mentioning the  $A$  arguments, are  $\{A_1, A_2, A_3\}$ ,  $\{A_1, A_2, A_5\}$ ,  $\{A_1, A_3, A_4\}$  and  $\{A_3, A_4, A_6\}$ . None of  $A_1, \dots, A_6$  is skeptically accepted and all are credulously accepted. Take conclusion  $f$ , where  $\mathcal{E} = \{A_3, A_4, A_6, B_1, B_2, B_3, B_4, C_1\}$ . Then:

- $\text{Acc}_{\text{prf}}^{\cup}(f) = \langle \{A_6\}, \{A_3, A_4, A_6\} \rangle$  for  $\mathbb{F}^{\text{acc}} = \text{id}$  and  $\mathbb{D}^{\text{acc}} \in \{\text{DefBy}, \text{DirDef}\}$ ;
- $\text{Acc}_{\text{prf}}^{\cup}(f) = \langle \{p, su\}, \{p, su\} \rangle$  for  $\mathbb{F}^{\text{acc}} = \text{Prem}$  and  $\mathbb{D}^{\text{acc}} \in \{\text{DefBy}, \text{DirDef}\}$ ;
- $\text{Acc}_{\text{prf}}^{\cup}(f) = \langle \{cd, \neg cpd\}, \{su\} \rangle$  for  $\mathbb{F}^{\text{acc}} = \text{AntTop}$  and  $\mathbb{D}^{\text{acc}} = \text{FinalDef}$ ;
- $\text{NotAcc}_{\text{prf}}^{\cup}(\neg f) = \langle \{A_1\}, \{A_3, A_4, A_6\} \rangle$  for  $\mathbb{F}^{\text{na}} = \text{id}$  and  $\mathbb{D}^{\text{acc}} = \text{NotDef}$ ;
- $\text{NotAcc}_{\text{prf}}^{\cup}(\neg f) = \langle \{A_1\}, \{A_3, A_4\} \rangle$  for  $\mathbb{F}^{\text{na}} = \text{id}$  and  $\mathbb{D}^{\text{acc}} = \text{NoDir}$ .

Looking at the different possibilities for  $\mathbb{F}$ , we see that instead of the full arguments we can also return just the premises (observations) of the supporting arguments, so ‘ $f$  because  $p$  and  $su$ ’. This is what the police system currently does. The reasoning behind this is that citizens understand these more factual observations better than more legal concepts such as delivering under a contract. On the other hand, for the public prosecutor involved in the processing of complaints, an explanation in legal terms – ‘ $f$  because  $cd$  and  $\neg cpd$ ’ (based on AntTop) – might make more sense.

For  $\mathbb{D}$  there are also different options. For example, FinalDef returns arguments that do not need other arguments to defend them. That  $A_3$  is such an argument w.r.t.  $A_6$  means that this argument  $A_3$  for  $fk$  is the ‘main reason’ we accept  $f$ , that is, without  $A_3$  the conclusion  $f$  will never be accepted. With NoDir, no directly conflicting arguments

are given (e.g.,  $A_6$  which directly conflicts with  $A_1$ ). This avoids explanations such as ‘(the argument for)  $\neg f$  is not accepted because (there is an argument for)  $f$ ’.

#### 4.5 Overview

In this section we have considered variations for the functions  $\mathbb{D}$  and  $\mathbb{F}$ . Acceptance explanations can be given in terms of all the defending arguments ( $\mathbb{D} = \text{DefBy}$ ), the arguments that need no further defense ( $\mathbb{D} = \text{FinalDef}$ ), and arguments that defend against direct conflicts ( $\mathbb{D} = \text{DirDef}$ ). Non-acceptance explanations can be given in terms of all the attackers for which no defense exists ( $\mathbb{D} = \text{NotDef}$ ) and those arguments that need to be defended by another argument ( $\mathbb{D} = \text{NoDir}$ ). In a structured setting (e.g., in  $\text{ASPIC}^+$ ), the form of these explanations can be varied. We discussed sets of arguments ( $\mathbb{F} = \text{id}$ ), sets of premises/observations ( $\mathbb{F} = \text{Prem}$ ) and sets of antecedents of the last applied rule ( $\mathbb{F} = \text{AntTop}$ ).

## 5 RELATED WORK

Fan and Toni [4] define relevant explanations for a single topic argument in the form of a new *related admissibility* semantics, and show how explanations can be derived from related admissible sets for abstract argumentation and ABA. A set of arguments is called related admissible if it is admissible and each argument in it defends the topic. An explanation for an argument  $A$  (called here *RA-explanation* to avoid confusion) is then defined as a related admissible set of arguments with topic  $A$ . In the next proposition we show how RA-explanations can be expressed in our framework.

**Proposition 2.** *Let  $\mathcal{AF} = \langle \text{Args}, \text{Att} \rangle$  be an AF and let  $A \in \text{Args}$ . Then  $\{\text{DefBy}(A, \mathcal{E}) \mid \mathcal{E} \in \mathfrak{C}_A^{\text{adm}}\}$  is the set of all RA-explanations for  $A$ .*

*Proof.* Let  $\mathcal{AF} = \langle \text{Args}, \text{Att} \rangle$  be an AF and let  $A \in \text{Args}$ . Suppose that  $\mathfrak{C}_A^{\text{adm}} \neq \emptyset$ . Let  $S \in \{\text{DefBy}(A, \mathcal{E}) \mid \mathcal{E} \in \mathfrak{C}_A^{\text{adm}}\}$ , we first show that  $S$  is related admissible:

$S$  **defends**  $A$ . This follows by the definition of  $S = \text{DefBy}(A, \mathcal{E})$ .

$S$  **is admissible**. Note that  $S \subseteq \mathcal{E}$  for some  $\mathcal{E} \in \mathfrak{C}_A^{\text{adm}}$ , therefore  $S$  is conflict-free. Suppose that there is some  $B \in S$  such that  $B$  is not defended against an attack from  $C \in \text{Args}$ . By definition of  $\text{DefBy}$ ,  $C$  (in)directly attacks  $A$ . Since  $A, B \in \mathcal{E}$ , there is some  $D \in \mathcal{E}$  such that  $D$  defends  $A$  and  $B$  against  $C$ . By assumption,  $D \notin S$ . A contradiction with the definition of  $\text{DefBy}$ . Therefore  $S$  defends all of its arguments and is thus admissible.

Now suppose that there is some  $S'$  which is an RA-explanation for  $A$  but  $S' \notin \{\text{DefBy}(A, \mathcal{E}) \mid \mathcal{E} \in \mathfrak{C}_A^{\text{adm}}\}$ . By definition of related admissible sets  $A \in S'$ ,  $S' \in \mathfrak{C}_A^{\text{adm}}$  and for each  $B \in S'$ ,  $B = A$  or  $B$  defends  $A$  in  $S'$ , thus  $B \in \text{DefBy}(A, \mathcal{E})$ , a contradiction. Hence any RA-explanation for  $A$  is in  $\{\text{DefBy}(A, \mathcal{E}) \mid \mathcal{E} \in \mathfrak{C}_A^{\text{adm}}\}$ .  $\square$

This shows that any  $\text{Acc}_{\text{adm}}^{\cup}$ -explanation is an RA-explanation and that therefore our framework is a more general version of [4].

García et al. [6] study explanations for abstract argumentation and DELP. Explanations for a claim are defined as triples of dialectical trees that provide a warrant for

the claim, dialectical trees that provide a warrant for the contrary of the claim, and dialectical trees for the claim and its contrary that provide no warrant. This means, on the one hand, that explanations might contain many arguments and, on the other hand, that the receiver of the explanation is expected to understand argumentation and dialectical trees. With real-life applications in mind, we believe that explanations that rely less on the underlying AF and that can be adjusted to the application are more useful. Therefore, in our framework an explanation consists of a set of (parts of) arguments, that could be embedded in a natural language sentence to be presented to a user, as suggested in Section 4.4.

Explanations for non-accepted arguments in abstract argumentation are studied in [5], [7], both of which focus on the structure of the AF and credulous non-acceptance under admissible semantics. Note that we consider skeptical and credulous non-acceptance for several Dung-style semantics. In [5] an explanation consists of either a set of arguments or a set of attacks, the removal of which would make the argument admissible. In structured argumentation it is not always possible to remove exactly one argument (or attack). In the AF of Figure 1,  $A_3$  would become skeptically acceptable for any semantics, if  $B_1$  would be removed. However, when looking at the underlying argumentation theory (recall Example 3), when  $B_1$  is removed, the arguments  $B_2$  and  $B_3$  do no longer exist and thus  $\neg q$  is no longer a credulous conclusion. Therefore, in this paper the basic definition for non-accepted arguments is defined in terms of the arguments for which no defense exist and no suggestion is made how to change the AF in order to get the considered argument accepted. In [7], explanations are sub-frameworks, such that the considered argument is credulously non-accepted in that sub-framework and any of its super-frameworks. Though a note was added on the applicability of such explanations in a structured setting, this is not formally investigated in that paper.

Summarizing, our basic framework is (formally) shown to be more general, more flexible and specifically adjustable to the receiver of the explanation. Furthermore, none of the above-mentioned works consider the structure of the arguments when providing explanations.

## 6 CONCLUSIONS AND FUTURE WORK

We have introduced a generic, flexible basic framework for explanations in structured and abstract argumentation. With this framework, specialized local explanations for the (non-)acceptance of arguments can be given, taking into account credulous and skeptical reasoners.

In future work, we plan to extend our framework with preferences – although showing preferences is sometimes considered less effective when providing explanations [3], the (non-)acceptance of arguments very often depends directly on them, making a preference the direct reason for (not) accepting an argument.

Given our basic framework, we will further study how our explanations formally relate to acceptance strategies and different semantics, and investigate the necessity and sufficiency of arguments and how to implement this in explanations.



Aside from formal investigations, we also want to look at how findings from the social sciences on what good explanations are (see e.g., [1], [3]) can be integrated, and how different types of explanations are evaluated by human users. Important in this respect is that explanations are *contrastive*: while people may ask *why A?*, they often mean *why A rather than B?*, where *A* is called the *fact* and *B* is called the *foil*. The goal is then to explain as much of the differences between fact and foil as possible. One of the challenges for an AI system is that the foil is not always explicit. We plan to study contrastive explanations within our framework by combining acceptance and non-acceptance and the knowledge of conflicting arguments and contraries in the case of an implicit foil.

## ACKNOWLEDGMENTS

This research has been partly funded by the Dutch Ministry of Justice and the Dutch National Police.

## REFERENCES

- [1] W. Samek, T. Wiegand, and K.-R. Müller, “Explainable artificial intelligence: Understanding, visualizing and interpreting deep learning models,” *arXiv preprint arXiv:1708.08296*, 2017.
- [2] C. Lacave and F. J. Diez, “A review of explanation methods for heuristic expert systems,” *The Knowledge Engineering Review*, vol. 19, no. 2, pp. 133–146, 2004.
- [3] T. Miller, “Explanation in artificial intelligence: Insights from the social sciences,” *Artificial Intelligence*, vol. 267, pp. 1 – 38, 2019.
- [4] X. Fan and F. Toni, “On computing explanations in argumentation,” in *Proceedings of the 29th AAAI Conference on Artificial Intelligence (AAAI’15)*, B. Bonet and S. Koenig, Eds. AAAI Press, 2015, pp. 1496–1502.
- [5] —, “On explanations for non-acceptable arguments,” in *Proceedings of the 3rd International Workshop on Theory and Applications of Formal Argumentation (TFAFA’15)*, ser. LNCS 9524, E. Black, S. Modgil, and N. Oren, Eds. Springer, 2015, pp. 112–127.
- [6] A. García, C. Chesñevar, N. Rotstein, and G. Simari, “Formalizing dialectical explanation support for argument-based reasoning in knowledge-based systems,” *Expert Systems with Applications*, vol. 40, no. 8, pp. 3233 – 3247, 2013.
- [7] Z. Saribatur, J. Wallner, and S. Woltran, “Explaining non-acceptability in abstract argumentation,” in *Proceedings of the 24th European Conference on Artificial Intelligence (ECAI’20)*, ser. Frontiers in Artificial Intelligence and Applications, vol. 325. IOS Press, 2020, pp. 881–888.
- [8] P. M. Dung, “On the acceptability of arguments and its fundamental role in nonmonotonic reasoning, logic programming and n-person games,” *Artificial Intelligence*, vol. 77, no. 2, pp. 321–357, 1995.
- [9] H. Prakken, “An abstract framework for argumentation with structured arguments,” *Argument & Computation*, vol. 1, no. 2, pp. 93–124, 2010.
- [10] L. Edwards and M. Veale, “Slave to the algorithm: Why a ‘right to an explanation’ is probably not the remedy you are looking for,” *Duke Law & Technology Review*, vol. 16, no. 1, pp. 18–84, 2017.
- [11] F. Bex, B. Testerink, and J. Peters, “AI for online criminal complaints: From natural dialogues to structured scenarios,” in *Workshop proceedings of Artificial Intelligence for Justice at ECAI 2016*, 2016, pp. 22–29.
- [12] D. Odekerken, A. Borg, and F. Bex, “Estimating stability for efficient argument-based inquiry,” in *Proceedings of the 8th International Conference on Computational Models of Argument (COMMA’20)*, vol. 326. IOS Press, 2020, pp. 307–318.

**AnneMarie Borg** is a postdoctoral researcher in the Police-lab AI at Utrecht University. Her research interests include formal argumentation and logic. She received her PhD from Ruhr University Bochum in 2019. Contact her at Politielab AI, Departement Informatica, Princetonplein 5, 3508 TB Utrecht, The Netherlands, a.borg@uu.nl.

**Floris Bex** is scientific director of the Police-lab AI at Utrecht University, and professor of Data Science and the judiciary at Tilburg University. His interests include argumentation and AI & Law. He received his PhD from the University of Groningen in 2009. Contact him at Politielab AI, Departement Informatica, Princetonplein 5, 3508 TB Utrecht, The Netherlands, f.j.bex@uu.nl.