

University of Groningen

## Expressiveness of SETAFs and Support-Free ADFs under 3-valued Semantics

Keshavarzi Zafarghandi, Atefeh; Woltran, Stefan ; Dvorak, Wolfgang

*Published in:*  
 IOSS Press

**IMPORTANT NOTE: You are advised to consult the publisher's version (publisher's PDF) if you wish to cite from it. Please check the document version below.**

*Document Version*  
 Early version, also known as pre-print

*Publication date:*  
 2020

[Link to publication in University of Groningen/UMCG research database](#)

*Citation for published version (APA):*

Keshavarzi Zafarghandi, A., Woltran, S., & Dvorak, W. (2020). Expressiveness of SETAFs and Support-Free ADFs under 3-valued Semantics. Manuscript submitted for publication. In *IOSS Press arXiv*.

### Copyright

Other than for strictly personal use, it is not permitted to download or to forward/distribute the text or part of it without the consent of the author(s) and/or copyright holder(s), unless the work is under an open content license (like Creative Commons).

### Take-down policy

If you believe that this document breaches copyright please contact us providing details, and we will remove access to the work immediately and investigate your claim.

Downloaded from the University of Groningen/UMCG research database (Pure): <http://www.rug.nl/research/portal>. For technical reasons the number of authors shown on this cover page is limited to 10 maximum.

# Expressiveness of SETAFs and Support-Free ADFs under 3-valued Semantics

Wolfgang Dvořák <sup>a</sup>      Atefeh Keshavarzi Zafarghandi <sup>b</sup>  
Stefan Woltran <sup>a</sup>

<sup>a</sup> Institute of Logic and Computation, TU Wien, Austria

<sup>b</sup> Department of Artificial Intelligence, Bernoulli Institute,  
University of Groningen, The Netherlands

July 8, 2020

## Abstract

Generalizing the attack structure in argumentation frameworks (AFs) has been studied in different ways. Most prominently, the binary attack relation of Dung frameworks has been extended to the notion of collective attacks. The resulting formalism is often termed SETAFs. Another approach is provided via abstract dialectical frameworks (ADFs), where acceptance conditions specify the relation between arguments; restricting these conditions naturally allows for so-called support-free ADFs. The aim of the paper is to shed light on the relation between these two different approaches. To this end, we investigate and compare the expressiveness of SETAFs and support-free ADFs under the lens of 3-valued semantics. Our results show that it is only the presence of unsatisfiable acceptance conditions in support-free ADFs that discriminate the two approaches.

## 1 Introduction

Abstract argumentation frameworks (AFs) as introduced by Dung [1] are a core formalism in formal argumentation. A popular line of research investigates extensions of Dung AFs that allow for a richer syntax (see, e.g. [2]). In this work we investigate two generalisations of Dung AFs that allow for a more flexible attack structure (but do not consider support between arguments).

The first formalism we consider are SETAFs as introduced by Nielsen and Parsons [3]. SETAFs extend Dung AFs by allowing for collective attacks such that a set of arguments  $B$  attacks another argument  $a$  but no proper subset of  $B$  attacks  $a$ . Argumentation frameworks with collective attacks have received increasing interest in the last years. For instance, semi-stable, stage, ideal, and eager semantics have been adapted to SETAFs in [4, 5]; translations between SETAFs and other abstract argumentation

formalisms are studied in [6]; [7] observed that for particular instantiations, SETAFs provide a more convenient target formalism than Dung AFs. The expressiveness of SETAFs with two-valued semantics has been investigated in [4] in terms of signatures. Signatures have been introduced in [8] for AFs. In general terms, a signature for a formalism and a semantics captures all possible outcomes that can be obtained by the instances of the formalism under the considered semantics. Besides that, signatures are recognized as crucial for operators in dynamics of argumentation (cf. [9]).

The second formalism we consider are support-free abstract dialectical frameworks (SFADFs), a subclass of abstract dialectical frameworks (ADFs) [10] which are known as an advanced abstract formalism for argumentation, that is able to cover several generalizations of AFs [2, 6]. This is accomplished by acceptance conditions which specify, for each argument, its relation to its neighbour arguments via propositional formulas. These conditions determine the links between the arguments which can be, in particular, attacking or supporting. SFADFs are ADFs where each link between arguments is attacking; they have been introduced in a recent study on different sub-classes of ADFs [11].

For comparison of the two formalisms, we need to focus on 3-valued (labelling) semantics [12, 13], which are integral for ADF semantics [10]. In terms of SETAFs, we can rely on the recently introduced labelling semantics in [5]. We first define a new class of ADFs (SETADFs) where the acceptance conditions strictly follow the nature of collective attacks in SETAFs and show that SETAFs and SETADFs coincide for the main semantics, i.e. the  $\sigma$ -labellings of a SETAF are equal to the  $\sigma$ -interpretations of the corresponding SETADF. We then provide exact characterisations of the 3-valued signatures for SETAFs (and thus for SETADFs) for most of the semantics under consideration. While SETADFs are a syntactically defined subclass of ADFs, the second formalism we study can be understood as semantical subclass of ADFs. In fact, for SFADFs it is not the syntactic structure of acceptance conditions that is restricted but their semantic behavior, in the sense that all links need to be attacking. The second main contribution of the paper is to determine the exact difference in expressiveness between SETADFs and SFADFs.

We briefly discuss related work. The expressiveness of SETAFs has first been investigated in [14] where different sub-classes of ADFs, i.e. AFs, SETAFs and Bipolar ADFs, are related w.r.t. their signatures of 3-valued semantics. Moreover, they provide an algorithm to decide realizability in one of the formalisms under different semantics. However, no explicit characterisations of the signatures are given. Recently, Pührer [15] presented explicit characterisations of the signatures of general ADFs (but not for the sub-classes discussed above). In contrast, [4] provides explicit characterisations of the two-valued signatures of SETAFs and shows that SETAFs are more expressive than AFs. In both works all arguments are relevant for the signature, while in [5] it is shown that when allowing to add extra arguments to an AF which are not relevant for the signature, i.e. the extensions/labellings are projected on common arguments, then SETAFs and AFs are of equivalent expressiveness. Other recent work [16] already implicitly showed that SFADFs with satisfiable acceptance conditions can be equivalently represented as SETAFs. This provides a sufficient condition for rewriting an ADF as SETAF and raises the question whether it is also a necessary condition. In fact, we will show that a SFADF has an equivalent SETAF if and only if all accep-

tance conditions are satisfiable. Different sub-classes of ADFs (including SFADFs) have been compared in [11], but no exact characterisations of signatures as we provide here are given in that work.

To summarize, the main contributions of our paper are as follows:

- We embed SETAFs under 3-valued labeling based semantics [5] in the more general framework of ADFs. That is, we show 3-valued labeling based SETAF semantics to be equivalent to the corresponding ADF semantics. As a side result, this also shows the equivalence of the 3-valued SETAF semantics in [14] and [5].
- We investigate the expressiveness of SETAFs under 3-valued semantics by providing exact characterizations of the signatures for preferred, stable, grounded and conflict-free semantics, thus complementing the investigations on expressiveness of SETAFs [4] in terms of extension-based semantics.
- We study the relations between SETAFs and support-free ADFs (SFADFs). In particular we give the exact difference in expressiveness between SETAFs and SFADFs under conflict-free, admissible, preferred, grounded, complete, stable and two-valued model semantics.

Some technical details had to be omitted but are available in an appendix.

## 2 Background

In this section we briefly recall the necessary definitions for SETAFs and ADFs.

**Definition 1.** A set argumentation framework (SETAF) is an ordered pair  $F = (A, R)$ , where  $A$  is a finite set of arguments and  $R \subseteq (2^A \setminus \{\emptyset\}) \times A$  is the attack relation.

The semantics of SETAFs are usually defined similarly to AFs, i.e., based on extensions. However, in this work we focus on 3-valued labelling based semantics, cf. [5].

**Definition 2.** A (3-valued) labelling of a SETAF  $F = (A, R)$  is a total function  $\lambda : A \mapsto \{\text{in}, \text{out}, \text{undec}\}$ . For  $x \in \{\text{in}, \text{out}, \text{undec}\}$  we write  $\lambda_x$  to denote the sets of arguments  $a \in A$  with  $\lambda(a) = x$ . We sometimes denote labellings  $\lambda$  as triples  $(\lambda_{\text{in}}, \lambda_{\text{out}}, \lambda_{\text{undec}})$ .

**Definition 3.** Let  $F = (A, R)$  be a SETAF. A labelling is called conflict-free in  $F$  if (i) for all  $(S, a) \in R$  either  $\lambda(a) \neq \text{in}$  or there is a  $b \in S$  with  $\lambda(b) \neq \text{in}$ , and (ii) for all  $a \in A$ , if  $\lambda(a) = \text{out}$  then there is an attack  $(S, a) \in R$  such that  $\lambda(b) = \text{in}$  for all  $b \in S$ . A labelling  $\lambda$  which is conflict-free in  $F$  is

- *admissible* in  $F$  iff for all  $a \in A$  if  $\lambda(a) = \text{in}$  then for all  $(S, a) \in R$  there is a  $b \in S$  such that  $\lambda(b) = \text{out}$ ;
- *complete* in  $F$  iff for all  $a \in A$  (i)  $\lambda(a) = \text{in}$  iff for all  $(S, a) \in R$  there is a  $b \in S$  such that  $\lambda(b) = \text{out}$ , and (ii)  $\lambda(a) = \text{out}$  iff there is an attack  $(S, a) \in R$  such that  $\lambda(b) = \text{in}$  for all  $b \in S$ ;
- *grounded* in  $F$  iff it is complete and there is no  $\lambda'$  with  $\lambda'_{\text{in}} \subset \lambda_{\text{in}}$  complete in  $F$ ;

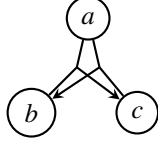


Figure 1: The SETAF of Example 1.

- *preferred* in  $F$  iff it is complete and there is no  $\lambda'$  with  $\lambda'_{\text{in}} \supset \lambda_{\text{in}}$  complete in  $F$ ;
- *stable* in  $F$  iff  $\lambda_{\text{undec}} = \emptyset$ .

The set of all  $\sigma$  labellings for a SETAF  $F$  is denoted by  $\sigma_{\mathcal{L}}(F)$ , where  $\sigma \in \{cf, adm, com, grd, prf, stb\}$  abbreviates the different semantics in the obvious manner.

**Example 1.** The SETAF  $F = (\{a, b, c\}, \{(\{a, b\}, c), (\{a, c\}, b)\})$  is depicted in Figure 1. For instance,  $(\{a, b\}, c) \in R$  says that there is a joint attack from  $a$  and  $b$  to  $c$ . This represents that neither  $a$  nor  $b$  is strong enough to attack  $c$  by themselves. Further,  $\{a \mapsto \text{in}, b \mapsto \text{undec}, c \mapsto \text{in}\}$  is an instance of a conflict-free labelling, that is not an admissible labelling (since  $c$  is mapped to  $\text{in}$  but neither  $a$  nor  $b$  is mapped to  $\text{out}$ ). The labelling that maps all argument to  $\text{undec}$  is not a complete labelling, however, it is an admissible labelling. Further,  $\{a \mapsto \text{in}, b \mapsto \text{undec}, c \mapsto \text{undec}\}$  is an admissible, the unique grounded and a complete labelling, which is not a preferred labelling because  $\lambda_{\text{in}} = \{a\}$  is not  $\subseteq$ -maximal among all complete labellings. Moreover,  $prf_{\mathcal{L}}(F) = stb_{\mathcal{L}}(F) = \{\{a \mapsto \text{in}, b \mapsto \text{out}, c \mapsto \text{in}\}, \{a \mapsto \text{in}, b \mapsto \text{in}, c \mapsto \text{out}\}\}$ .

We next turn to abstract dialectical frameworks [17].

**Definition 4.** An abstract dialectical framework (ADF) is a tuple  $D = (S, L, C)$  where:

- $S$  is a finite set of arguments (statements, positions);
- $L \subseteq S \times S$  is a set of links among arguments;
- $C = \{\varphi_s\}_{s \in S}$  is a collection of propositional formulas over arguments, called acceptance conditions.

An ADF can be represented by a graph in which nodes indicate arguments and links show the relation among arguments. Each argument  $s$  in an ADF is attached by a propositional formula, called acceptance condition,  $\varphi_s$  over  $par(s)$  such that,  $par(s) = \{b \mid (b, s) \in L\}$ . Since in ADFs an argument appears in the acceptance condition of an argument  $s$  if and only if it belongs to the set  $par(s)$ , the set of links  $L$  of an ADF is given implicitly via the acceptance conditions. The acceptance condition of each argument clarifies under which condition the argument can be accepted and determines the type of links (see Definition 6 below). An *interpretation*  $v$  (for  $F$ ) is a function  $v : S \mapsto \{\mathbf{t}, \mathbf{f}, \mathbf{u}\}$ , that maps arguments to one of the three truth values true ( $\mathbf{t}$ ), false ( $\mathbf{f}$ ), or undecided ( $\mathbf{u}$ ). Truth values can be ordered via information ordering relation  $<_i$  given by  $\mathbf{u} <_i \mathbf{t}$  and  $\mathbf{u} <_i \mathbf{f}$  and no other pair of truth values are related by  $<_i$ . Relation  $\leq_i$  is the reflexive and transitive closure of  $<_i$ . An interpretation  $v$  is *two-valued* if it maps each argument to either  $\mathbf{t}$  or  $\mathbf{f}$ . Let  $\mathcal{V}$  be the set of all interpretations for an ADF

$D$ . Then, we call a subset of all interpretations of the ADF,  $\mathbb{V} \subseteq \mathcal{V}$ , an *interpretation-set*. Interpretations can be ordered via  $\leq_i$  with respect to their information content, i.e.  $w \leq_i v$  if  $w(s) \leq_i v(s)$  for each  $s \in S$ . Further, we denote the update of an interpretation  $v$  with a truth value  $x \in \{\mathbf{t}, \mathbf{f}, \mathbf{u}\}$  for an argument  $b$  by  $v|_x^b$ , i.e.  $v|_x^b(b) = x$  and  $v|_x^b(a) = v(a)$  for  $a \neq b$ . Finally, the partial valuation of acceptance condition  $\varphi_s$  by  $v$ , is given by  $\varphi_s^v = v(\varphi_s) = \varphi_s[p/\top : v(p) = \mathbf{t}][p/\perp : v(p) = \mathbf{f}]$ , for  $p \in \text{par}(s)$ .

Semantics for ADFs can be defined via a *characteristic operator*  $\Gamma_D$  for an ADF  $D$ . Given an interpretation  $v$  (for  $D$ ), the characteristic operator  $\Gamma_D$  for  $D$  is defined as

$$\Gamma_D(v) = v' \text{ such that } v'(s) = \begin{cases} \mathbf{t} & \text{if } \varphi_s^v \text{ is irrefutable (i.e., a tautology),} \\ \mathbf{f} & \text{if } \varphi_s^v \text{ is unsatisfiable,} \\ \mathbf{u} & \text{otherwise.} \end{cases}$$

**Definition 5.** Given an ADF  $D = (S, L, C)$ , an interpretation  $v$  is

- *conflict-free* in  $D$  iff  $v(s) = \mathbf{t}$  implies  $\varphi_s^v$  is satisfiable and  $v(s) = \mathbf{f}$  implies  $\varphi_s^v$  is unsatisfiable;
- *admissible* in  $D$  iff  $v \leq_i \Gamma_D(v)$ ;
- *complete* in  $D$  iff  $v = \Gamma_D(v)$ ;
- *grounded* in  $D$  iff  $v$  is the least fixed-point of  $\Gamma_D$ ;
- *preferred* in  $D$  iff  $v$  is  $\leq_i$ -maximal admissible in  $D$ ;
- a *(two-valued) model* of  $D$  iff  $v$  is two-valued and for all  $s \in S$ , it holds that  $v(s) = v(\varphi_s)$ ;
- a *stable model* of  $D$  if  $v$  is a model of  $D$  and  $v^t = w^t$ , where  $w$  is the grounded interpretation of the *stb*-reduct  $D^v = (S^v, L^v, C^v)$ , where  $S^v = v^t$ ,  $L^v = L \cap (S^v \times S^v)$ , and  $\varphi_s[p/\perp : v(p) = \mathbf{f}]$  for each  $s \in S^v$ .

The set of all  $\sigma$  interpretations for an ADF  $D$  is denoted by  $\sigma(D)$ , where  $\sigma \in \{cf, adm, com, grd, prf, mod, stb\}$  abbreviates the different semantics in the obvious manner.

**Example 2.** An example of an ADF  $D = (S, L, C)$  is shown in Figure 2. To each argument a propositional formula is associated, the acceptance condition of the argument. For instance, the acceptance condition of  $c$ , namely  $\varphi_c : \neg a \vee \neg b$ , states that  $c$  can be accepted in an interpretation where either  $a$  or  $b$  (or both) are rejected.

In  $D$  the interpretation  $v = \{a \mapsto \mathbf{u}, b \mapsto \mathbf{u}, c \mapsto \mathbf{t}\}$  is conflict-free. However,  $v$  is not an admissible interpretation, because  $\Gamma_D(v) = \{a \mapsto \mathbf{u}, b \mapsto \mathbf{u}, c \mapsto \mathbf{u}\}$ , that is,  $v \not\leq_i \Gamma_D(v)$ . The interpretation  $v_1 = \{a \mapsto \mathbf{f}, b \mapsto \mathbf{t}, c \mapsto \mathbf{u}\}$  on the other hand is an admissible interpretation. Since  $\Gamma_D(v_1) = \{a \mapsto \mathbf{f}, b \mapsto \mathbf{t}, c \mapsto \mathbf{t}\}$  and  $v_1 \leq_i \Gamma_D(v_1)$ . Further,  $\text{prf}(D) = \text{mod}(D) = \{\{a \mapsto \mathbf{t}, b \mapsto \mathbf{f}, c \mapsto \mathbf{t}\}, \{a \mapsto \mathbf{f}, b \mapsto \mathbf{t}, c \mapsto \mathbf{t}\}\}$ , but only the first interpretation in this set is a stable model. This is because for  $v = \{a \mapsto \mathbf{t}, b \mapsto \mathbf{f}, c \mapsto \mathbf{t}\}$  the unique grounded interpretation  $w$  of  $D^v$  is  $\{a \mapsto \mathbf{t}, c \mapsto \mathbf{t}\}$  and  $v^t = w^t$ . The interpretation  $v' = \{a \mapsto \mathbf{f}, b \mapsto \mathbf{t}, c \mapsto \mathbf{t}\}$  is not a stable model, since the unique grounded interpretation  $w'$  of  $D^{v'}$  is  $\{b \mapsto \mathbf{u}, c \mapsto \mathbf{t}\}$  and  $v'^t \neq w'^t$ . Actually,  $v'$  is not

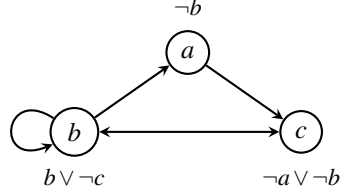


Figure 2: The ADF of Example 2.

a stable model because the truth value of  $b$  in  $v'$  is since of self-support. Moreover, the unique grounded interpretation of  $D$  is  $v = \{a \mapsto \mathbf{u}, b \mapsto \mathbf{u}, c \mapsto \mathbf{u}\}$ . In addition, we have  $com(D) = prf(D) \cup grd(D)$ .

In ADFs links between arguments can be classified into four types, reflecting the relationship of attack and/or support that exists among the arguments. In Definition 6 we consider two-valued interpretations that are only defined over the parents of  $a$ , that is, only give values to  $par(a)$ .

**Definition 6.** Let  $D = (S, L, C)$  be an ADF. A link  $(b, a) \in L$  is called

- *supporting* (in  $D$ ) if for every two-valued interpretation  $v$  of  $par(a)$ ,  $v(\varphi_a) = \mathbf{t}$  implies  $v|_{\mathbf{t}}^b(\varphi_a) = \mathbf{t}$ ;
- *attacking* (in  $D$ ) if for every two-valued interpretation  $v$  of  $par(a)$ ,  $v(\varphi_a) = \mathbf{f}$  implies  $v|_{\mathbf{t}}^b(\varphi_a) = \mathbf{f}$ ;
- *redundant* (in  $D$ ) if it is both attacking and supporting;
- *dependent* (in  $D$ ) if it is neither attacking nor supporting.

The classification of the types of the links of ADFs is also relevant for classifying ADFs themselves. One particularly important subclass of ADFs is that of *bipolar* ADFs or BADFs for short. In such an ADF each link is either attacking or supporting (or both; thus, the links can also be redundant). Another subclass of ADFs, having only attacking links, is defined in [18], called *support free ADFs* (SFADFs) in the current work, defined formally as follows.

**Definition 7.** An ADF is called support-free if it has only attacking links.

For SFADFs, it turns out that the intention of stable semantics, i.e. to avoid cyclic support among arguments, becomes immaterial, thus  $mod(D) = stb(D)$  for any ADF  $D$ ; the property is called weakly coherent in [18].

**Proposition 1.** For every SFADF  $D$  it holds that  $mod(D) = stb(D)$ .

*Proof.* The result follows from the following observation: Let  $D = (S, L, C)$  be an ADF, let  $v$  be a model of  $D$  and let  $s \in S$  be an argument such that all parents of  $s$  are attackers. Thus,  $\varphi_s^v$  is irrefutable if and only if  $\varphi_s[p/\perp : v(p) = \mathbf{f}]$  is irrefutable.  $\square$

### 3 Embedding SETAFs in ADFs

As observed by Polberg [19] and Linsbichler et.al [14], the notion of collective attacks can also be represented in ADFs by using the right acceptance conditions. We next introduce the class SETADFs of ADFs for this purpose.

**Definition 8.** An ADF  $D = (S, L, C)$  is called SETAF-like (SETADF) if each of the acceptance conditions in  $C$  is given by a formula (with  $\mathcal{C}$  a set of non-empty clauses)

$$\bigwedge_{cl \in \mathcal{C}} \bigvee_{a \in cl} \neg a.$$

That is, in a SETADF each acceptance condition is either  $\top$  (if  $\mathcal{C}$  is empty) or a proper CNF formula over negative literals. SETADFs and SETAFs can be embedded in each other as follows.

**Definition 9.** Let  $F = (A, R)$  be a SETAF. The ADF associated to  $F$  is a tuple  $D_F = (S, L, C)$  in which  $S = A$ ,  $L = \{(a, b) \mid (B, b) \in R, a \in B\}$  and  $C = \{\varphi_a\}_{a \in S}$  is the collection of acceptance conditions defined, for each  $a \in S$ , as

$$\varphi_a = \bigwedge_{(B, a) \in R} \bigvee_{a' \in B} \neg a'.$$

Let  $D = (S, L, C)$  be a SETADF. We construct the SETAF  $F_D = (A, R)$  in which,  $A = S$ , and  $R$  is constructed as follows. For each argument  $s \in S$  with acceptance formula  $\bigwedge_{cl \in \mathcal{C}} \bigvee_{a \in cl} \neg a$  we add the attacks  $\{(cl, s) \mid cl \in \mathcal{C}\}$  to  $R$ .

Clearly the ADF  $D_F$  associated to a SETAF  $F$  is a SETADF and  $D$  is the ADF associated to the constructed SETAF  $F_D$ . We next deal with the fact that SETAF semantics are defined as three-valued labellings while semantics for ADFs are defined as three valued interpretations. In order to compare these semantics we associate the *in* label with  $t$ , the *out* label with  $f$ , and the *undec* label with  $u$ .

**Theorem 2.** For  $\sigma \in \{cf, adm, com, prf, grd, stb\}$ , a SETAF  $F$  and its associated SETADF  $D$ , we have that  $\sigma_{\mathcal{L}}(F)$  and  $\sigma(D)$  are in one-to-one correspondence with each labelling  $\mathbb{L} \in \sigma_{\mathcal{L}}(F)$  corresponding to an interpretation  $v \in \sigma(D)$  such that  $v(s) = \mathbf{t}$  iff  $\lambda(s) = \mathbf{in}$ ,  $v(s) = \mathbf{f}$  iff  $\lambda(s) = \mathbf{out}$ , and  $v(s) = \mathbf{u}$  iff  $\lambda(s) = \mathbf{undec}$ .

Notice that by the above theorem we have that the 3-valued SETAF semantics introduced in [14] coincide with the 3-valued labelling based SETAF semantics of [5] and the model semantics of [14] corresponds to the stable semantics of [5].

### 4 3-valued Signatures of SETAFs

We adapt the concept of signatures [8] towards our needs first.

**Definition 10.** The signature of SETAFs under a labelling-based semantics  $\sigma_{\mathcal{L}}$  is defined as  $\Sigma_{SETAF}^{\sigma_{\mathcal{L}}} = \{\sigma_{\mathcal{L}}(F) \mid F \in SETAF\}$ . The signature of an ADF-subclass  $\mathcal{C}$  under a semantics  $\sigma$  is defined as  $\Sigma_{\mathcal{C}}^{\sigma} = \{\sigma(D) \mid D \in \mathcal{C}\}$ .



By Theorem 2 we can use labellings of SETAFs and interpretations of the SETADF class of ADFs interchangeably, yielding that  $\Sigma_{SETAF}^{\sigma_{\mathcal{L}}} \equiv \Sigma_{SETADF}^{\sigma}$ , i.e. the 3-valued signatures of SETAFs and SETADFs only differ in the naming of the labels. For convenience, we will use the SETAF terminology in this section.

**Proposition 3.** *The signature  $\Sigma_{SETAF}^{stb_{\mathcal{L}}}$  is given by all sets  $\mathbb{L}$  of labellings such that*

1. *all  $\lambda \in \mathbb{L}$  have the same domain  $\text{ARGS}_{\mathbb{L}}$ ;  $\lambda(s) \neq \text{undec}$  for all  $\lambda \in \mathbb{L}$ ,  $s \in \text{ARGS}_{\mathbb{L}}$ .*
2. *If  $\lambda \in \mathbb{L}$  assigns one argument to out then it also assigns an argument to in.*
3. *For arbitrary  $\lambda_1, \lambda_2 \in \mathbb{L}$  with  $\lambda_1 \neq \lambda_2$  there is an argument  $a$  such that  $\lambda_1(a) = \text{in}$  and  $\lambda_2(a) = \text{out}$ .*

*Proof.* We first show that for each SETAF  $F$  the set  $stb_{\mathcal{L}}(F)$  satisfies the conditions of the proposition. First clearly all  $\lambda \in stb_{\mathcal{L}}(F)$  have the same domain and by the definition of stable semantics do not assign undec to any argument. That is the first condition is satisfied. For Condition (2), towards a contradiction assume that the domain is non-empty and  $\lambda \in stb_{\mathcal{L}}(F)$  assigns all arguments to out. Consider an arbitrary argument  $a$ . By definition of stable semantics  $a$  is only labeled out if there is an attack  $(B, a)$  such that all arguments in  $B$  are labeled in, a contradiction. Thus we obtain that there is at least one argument  $a$  with  $\lambda(a) = \text{in}$ . For Condition (3), towards a contradiction assume that for all arguments  $a$  with  $\lambda_1(a) = \text{in}$  also  $\lambda_2(a) = \text{in}$  holds. As  $\lambda_1 \neq \lambda_2$  there is an  $a$  with  $\lambda_2(a) = \text{in}$  and  $\lambda_1(a) = \text{out}$ . That is, there is an attack  $(B, a)$  such that  $\lambda_1(b) = \text{in}$  for all  $b \in B$ . But then also  $\lambda_2(b) = \text{in}$  for all  $b \in B$  and by  $\lambda_2(a) = \text{in}$  we obtain that  $\lambda_2 \notin cf_{\mathcal{L}}(F)$ , a contradiction.

Now assume that  $\mathbb{L}$  satisfies all the conditions. We give a SETAF  $F_{\mathbb{L}} = (A_{\mathbb{L}}, R_{\mathbb{L}})$  with  $A_{\mathbb{L}} = \text{ARGS}_{\mathbb{L}}$  and  $R_{\mathbb{L}} = \{(\lambda_{\text{in}}, a) \mid \lambda \in \mathbb{L}, \lambda(a) = \text{out}\}$ . We show that  $stb_{\mathcal{L}}(F_{\mathbb{L}}) = \mathbb{L}$ .

To this end we first show  $stb_{\mathcal{L}}(F_{\mathbb{L}}) \supseteq \mathbb{L}$ . Consider an arbitrary  $\lambda \in \mathbb{L}$ : By Condition (1) there is no  $a \in \text{ARGS}_{\mathbb{L}}$  with  $\lambda(a) = \text{undec}$  and it only remains to show  $\lambda \in cf_{\mathcal{L}}(F_{\mathbb{L}})$ . First, if  $\lambda(a) = \text{out}$  for some argument  $a$  then by construction of  $R_{\mathbb{L}}$  and Condition (2) we have an attack  $(\lambda_{\text{in}}, a)$  and thus  $a$  is legally labeled out. Now towards a contradiction assume there is a conflict  $(B, a)$  such that  $B \cup \{a\} \subseteq \lambda_{\text{in}}$ . Then, by construction of  $R_{\mathbb{L}}$  there is a  $\lambda' \in \mathbb{L}$  with  $\lambda'_{\text{in}} = B$  and  $\lambda_{\text{in}} \neq B$  (as  $a \in \lambda_{\text{in}}$ ). That is,  $\lambda'_{\text{in}} \subset \lambda_{\text{in}}$ , a contradiction to Condition (3). Thus,  $\lambda \in cf_{\mathcal{L}}(F_{\mathbb{L}})$  and therefore  $\lambda \in stb_{\mathcal{L}}(F_{\mathbb{L}})$ .

To show  $stb_{\mathcal{L}}(F_{\mathbb{L}}) \subseteq \mathbb{L}$ , consider  $\lambda \in stb_{\mathcal{L}}(F_{\mathbb{L}})$ . If  $\lambda$  maps all arguments to in then there is no attack in  $R_{\mathbb{L}}$  which means that  $\mathbb{L}$  contains only the labelling  $\lambda$ . Thus, we assume that there is  $a$  with  $\lambda(a) = \text{out}$  and there is  $(B, a) \in R_{\mathbb{L}}$  with  $B \subseteq \lambda_{\text{in}}$ . By construction there is  $\lambda' \in \mathbb{L}$  such that  $\lambda'_{\text{in}} = B$ . Then by construction we have  $(B, c) \in R_{\mathbb{L}}$  for all  $c \notin B$  and thus  $\lambda'_{\text{in}} = B = \lambda_{\text{in}}$  and moreover  $\lambda'_{\text{out}} = \lambda_{\text{out}}$  and thus  $\lambda = \lambda'$ .  $\square$

We now turn to the signature for preferred semantics. Compared to the conditions for stable semantics, labelling may now assign undec to arguments. Note that stable is the only semantics allowing for an empty labelling set.

**Proposition 4.** The signature  $\Sigma_{SETAF}^{prf_{\mathcal{L}}}$  is given by all non-empty sets  $\mathbb{L}$  of labellings s.t.

1. all labellings  $\lambda \in \mathbb{L}$  have the same domain  $\text{ARGS}_{\mathbb{L}}$ .
2. If  $\lambda \in \mathbb{L}$  assigns one argument to *out* then it also assigns an argument to *in*.
3. For arbitrary  $\lambda_1, \lambda_2 \in \mathbb{L}$  with  $\lambda_1 \neq \lambda_2$  there is an argument  $a$  such  $\lambda_1(a) = \text{in}$  and  $\lambda_2(a) = \text{out}$ .

*Proof sketch.* We first show that for each SETAF  $F$  the set  $prf_{\mathcal{L}}(F)$  satisfies the conditions of the proposition. The first condition is satisfied as all  $\lambda \in prf_{\mathcal{L}}(F)$  have the same domain. The second condition is satisfied by the definition of conflict-free labellings. Condition (3) is by the  $\subseteq$ -maximality of  $\lambda_{\text{in}}$  which implies that there is a conflict between each two preferred extensions.

Now assume that  $\mathbb{L}$  satisfies all the conditions. We give a SETAF  $F_{\mathbb{L}} = (A_{\mathbb{L}}, R_{\mathbb{L}})$  with  $A_{\mathbb{L}} = \text{ARGS}_{\mathbb{L}}$  and  $R_{\mathbb{L}} = \{(\lambda_{\text{in}}, a) \mid \lambda \in \mathbb{L}, \lambda(a) = \text{out}\} \cup \{(\lambda_{\text{in}} \cup \{a\}, a) \mid \lambda \in \mathbb{L}, \lambda(a) = \text{undec}\}$ . It remains to show that  $prf_{\mathcal{L}}(F_{\mathbb{L}}) = \mathbb{L}$ . To show  $prf_{\mathcal{L}}(F_{\mathbb{L}}) \supseteq \mathbb{L}$ , consider an arbitrary  $\lambda \in \mathbb{L}$ .  $\lambda \in cf_{\mathcal{L}}(F_{\mathbb{L}})$  can be seen by construction, and  $\lambda \in adm_{\mathcal{L}}(F_{\mathbb{L}})$  since argument labelled *out* is attacked by  $\lambda$ ; finally  $\lambda \in prf_{\mathcal{L}}(F_{\mathbb{L}})$  is guaranteed since the arguments  $a$  with  $\lambda(a) = \text{undec}$  are involved in self-attacks. To show  $prf_{\mathcal{L}}(F_{\mathbb{L}}) \subseteq \mathbb{L}$  consider  $\lambda \in prf_{\mathcal{L}}(F_{\mathbb{L}})$ . It can be checked that  $\lambda$  satisfies all the conditions of the proposition.  $\square$

**Proposition 5.** The signature  $\Sigma_{SETAF}^{cf_{\mathcal{L}}}$  is given by all non-empty sets  $\mathbb{L}$  of labellings s.t.

1. all  $\lambda \in \mathbb{L}$  have the same domain  $\text{ARGS}_{\mathbb{L}}$ .
2. If  $\lambda \in \mathbb{L}$  assigns one argument to *out* then it also assigns an argument to *in*.
3. For  $\lambda \in \mathbb{L}$  and  $C \subseteq \lambda_{\text{in}}$  also  $(C, \emptyset, \text{ARGS}_{\mathbb{L}} \setminus C) \in \mathbb{L}$ .
4. For  $\lambda \in \mathbb{L}$  and  $C \subseteq \lambda_{\text{out}}$  also  $(\lambda_{\text{in}}, \lambda_{\text{out}} \setminus C, \lambda_{\text{undec}} \cup C) \in \mathbb{L}$ .
5. For  $\lambda, \lambda' \in \mathbb{L}$  with  $\lambda_{\text{in}} \subseteq \lambda'_{\text{in}}$  also  $(\lambda'_{\text{in}}, \lambda_{\text{out}} \cup \lambda'_{\text{out}}, \lambda_{\text{undec}} \cap \lambda'_{\text{undec}}) \in \mathbb{L}$ .
6. For  $\lambda, \lambda' \in \mathbb{L}$  and  $C \subseteq \lambda_{\text{out}}$  (s.t.  $C \neq \emptyset$ ) we have  $\lambda_{\text{in}} \cup C \not\subseteq \lambda'_{\text{in}}$ .

*Proof sketch.* Let  $F$  be an arbitrary SETAF we show that  $cf_{\mathcal{L}}(F)$  satisfies the conditions of the proposition. The first two conditions are clearly satisfied by the definition of conflict-free labelling. For Condition (3), towards a contradiction assume that  $(C, \emptyset, \text{ARGS}_{\mathbb{L}} \setminus C)$  is not conflict-free. Then there is an attack  $(B, a)$  such that  $B \cup \{a\} \subseteq C \subseteq \lambda_{\text{in}}$ , and thus  $\lambda \notin cf_{\mathcal{L}}(F)$ , a contradiction. Condition (4) is satisfied as in the definition of conflict-free labellings there are no conditions for labeling an argument *undec*. Further, the conditions that allow to label an argument *out* solely depend on the *in* labeled arguments. For Condition (5), consider  $\lambda, \lambda' \in cf_{\mathcal{L}}(F)$  with  $\lambda_{\text{in}} \subseteq \lambda'_{\text{in}}$  and  $\lambda^* = (\lambda'_{\text{in}}, \lambda_{\text{out}} \cup \lambda'_{\text{out}}, \lambda_{\text{undec}} \cap \lambda'_{\text{undec}})$ . Since  $\lambda, \lambda' \in \mathbb{L}$ , it is easy to check that  $\lambda^*$  is a well-founded labelling and  $\lambda^* \in cf_{\mathcal{L}}(F)$ . For Condition (6), consider  $\lambda, \lambda' \in cf_{\mathcal{L}}(F)$  and a set  $C \subseteq \lambda_{\text{out}}$  containing an argument  $a$  such that  $\lambda(a) = \text{out}$ . That is, there is an attack  $(B, a)$  with  $B \subseteq \lambda_{\text{in}}$  and thus  $\lambda_{\text{in}} \cup C \not\subseteq \lambda'_{\text{in}}$ . That is, Condition (6) is satisfied.

Now assume that  $\mathbb{L}$  satisfies all the conditions. We give a SETAF  $F_{\mathbb{L}} = (A_{\mathbb{L}}, R_{\mathbb{L}})$  with  $A_{\mathbb{L}} = \text{ARGS}_{\mathbb{L}}$  and  $R_{\mathbb{L}} = \{(\lambda_{\text{in}}, a) \mid \lambda \in \mathbb{L}, \lambda(a) = \text{out}\} \cup \{(B, b) \mid b \in B, \nexists \lambda \in \mathbb{L} : \lambda_{\text{in}} = B\}$ . To complete the proof it remains to show that  $cf_{\mathcal{L}}(F_{\mathbb{L}}) = \mathbb{L}$ .  $\square$

Finally, we give an exact characterisation of the signature of grounded semantics.

**Proposition 6.** *The signature  $\Sigma_{\text{SETAF}}^{\text{grd}_{\mathcal{L}}}$  is given by sets  $\mathbb{L}$  of labellings such that  $|\mathbb{L}| = 1$ , and if  $\lambda \in \mathbb{L}$  assigns one argument to  $\text{out}$  then  $\lambda_{\text{in}} \neq \emptyset$ .*

Notice that Proposition 6 basically exploits that grounded semantics is a unique status semantics based on admissibility. The result thus immediately extends to other semantics satisfying these two properties, e.g. to ideal or eager semantics [5].

So far, we have provided characterisations for the signatures  $\Sigma_{\text{SETAF}}^{\text{stb}_{\mathcal{L}}}$ ,  $\Sigma_{\text{SETAF}}^{\text{prf}_{\mathcal{L}}}$ ,  $\Sigma_{\text{SETAF}}^{\text{cf}_{\mathcal{L}}}$ ,  $\Sigma_{\text{SETAF}}^{\text{grd}_{\mathcal{L}}}$ . By Theorem 2 we get analogous characterizations of  $\Sigma_{\text{SETADF}}^{\sigma}$  for the corresponding ADF semantics.

We have not yet touched admissible and complete semantics. Here, the exact characterisations seem to be more cumbersome and are left for future work. However, for admissible semantics the following proposition provides necessary conditions for an labelling-set to be *adm*-realizable, but it remains open whether they are also sufficient.

**Proposition 7.** *For each  $\mathbb{L} \in \Sigma_{\text{SETAF}}^{\text{adm}_{\mathcal{L}}}$  we have:*

1. *all  $\lambda \in \mathbb{L}$  have the same domain  $\text{ARGS}_{\mathbb{L}}$ .*
2. *If  $\lambda \in \mathbb{L}$  assigns one argument to  $\text{out}$  then it also assigns an argument to  $\text{in}$ .*
3. *For  $\lambda, \lambda' \in \mathbb{L}$  and  $C \subseteq \lambda_{\text{out}}$  (s.t.  $C \neq \emptyset$ ) we have  $\lambda_{\text{in}} \cup C \not\subseteq \lambda'_{\text{in}}$ .*
4. *For arbitrary  $\lambda, \lambda' \in \mathbb{L}$  either (a)  $(\lambda_{\text{in}} \cup \lambda'_{\text{in}}, \lambda_{\text{out}} \cup \lambda'_{\text{out}}, \lambda_{\text{undec}} \cap \lambda'_{\text{undec}}) \in \mathbb{L}$  or (b) there is an argument  $a$  such  $\lambda(a) = \text{in}$  and  $\lambda'(a) = \text{out}$ .*
5. *For  $\lambda, \lambda' \in \mathbb{L}$  with  $\lambda_{\text{out}} \subseteq \lambda'_{\text{out}}$ , and  $C \subseteq \lambda_{\text{in}} \setminus \bigcup_{\lambda^* \in \mathbb{L} : \lambda^*_{\text{in}} = \lambda'_{\text{in}}} \lambda^*_{\text{out}}$  we have  $(\lambda'_{\text{in}} \cup C, \lambda'_{\text{out}}, \lambda'_{\text{undec}} \setminus C) \in \mathbb{L}$ .*
6. *For  $\lambda, \lambda' \in \mathbb{L}$  with  $\lambda_{\text{in}} \subseteq \lambda'_{\text{in}}$ , and  $C \subseteq \lambda_{\text{out}}$  we have  $(\lambda'_{\text{in}}, \lambda'_{\text{out}} \cup C, \lambda'_{\text{undec}} \setminus C) \in \mathbb{L}$ .*
7. *For  $\lambda, \lambda' \in \mathbb{L}$  with  $\lambda_{\text{in}} \subseteq \lambda'_{\text{in}}$  and  $\lambda_{\text{out}} \supseteq \lambda'_{\text{out}}$  we have  $(\lambda_{\text{in}}, \lambda'_{\text{out}}, \text{ARGS}_{\mathbb{L}} \setminus (\lambda_{\text{in}} \cup \lambda'_{\text{out}})) \in \mathbb{L}$ .*
8.  *$(\emptyset, \emptyset, \text{ARGS}_{\mathbb{L}}) \in \mathbb{L}$ .*

*Proof.* We show that for each SETAF  $F$  the set  $\text{adm}_{\mathcal{L}}(F)$  satisfies the conditions of the proposition. Conditions (1)–(3) are by the fact that  $\text{adm}_{\mathcal{L}}(F) \subseteq \text{cf}_{\mathcal{L}}(F)$ . For Condition (4), let  $\lambda, \lambda' \in \text{adm}_{\mathcal{L}}(F)$  with  $\lambda_{\text{in}} \cap \lambda'_{\text{out}} = \{\}$  (since each admissible labelling defends itself,  $\lambda'_{\text{in}} \cap \lambda_{\text{out}} = \{\}$ ). Thus,  $\lambda^* = (\lambda_{\text{in}} \cup \lambda'_{\text{in}}, \lambda_{\text{out}} \cup \lambda'_{\text{out}}, \lambda_{\text{undec}} \cap \lambda'_{\text{undec}})$  is a well-defined labelling. Further, since  $\lambda, \lambda' \in \text{adm}_{\mathcal{L}}(F)$  it is easy to check that  $\lambda^* \in \text{adm}_{\mathcal{L}}(F)$ . For Condition (5), let  $\lambda^* = (\lambda'_{\text{in}} \cup C, \lambda'_{\text{out}}, \lambda'_{\text{undec}} \setminus C)$ . First,  $\lambda^*$  is a well-defined labelling. Notice that the set  $C$  contains arguments defended by  $\lambda$  and not attacked by  $\lambda'_{\text{in}}$ . Now, it is easy to check that  $\lambda^*$  meets the condition for being

an admissible labelling. For Condition (6), let  $\lambda^* = (\lambda'_{\text{in}}, \lambda'_{\text{out}} \cup C, \lambda'_{\text{undec}} \setminus C)$ . Notice that the set  $C$  contains only arguments attacked by  $\lambda_{\text{in}}$  and thus are also attacked by  $\lambda'_{\text{in}}$ . Thus, starting from the admissible labelling  $\lambda'$  we can relabel arguments in  $C$  to out and obtain that  $\lambda^*$  is also an admissible labelling. For Condition (7), let  $\lambda^* = (\lambda_{\text{in}}, \lambda'_{\text{out}}, \text{ARGS}_{\perp} \setminus (\lambda_{\text{in}} \cup \lambda'_{\text{out}}))$ . First,  $\lambda^*$  is a well-defined labelling. We have that setting  $\lambda'_{\text{out}}$  to out is sufficient to make all the in labels for arguments in  $\lambda'_{\text{in}}$  valid and thus are also sufficient to make the in labels for arguments  $\lambda_{\text{in}} \subseteq \lambda'_{\text{in}}$  valid. Moreover, as  $\lambda_{\text{out}} \supseteq \lambda'_{\text{out}}$  also labelling arguments  $\lambda_{\text{in}}$  with in is sufficient to make the out labels for  $\lambda'_{\text{out}}$  valid. Hence,  $\lambda^*$  is admissible. For Condition (8), the conditions of admissible labelling for arguments labelled in or out in  $(\emptyset, \emptyset, \text{ARGS}_{\perp})$  are clearly met, since there are no such arguments.  $\square$

## 5 On the Relation between SETAFs and Support-Free ADFs

In order to compare SETAFs with SFADFs, we can rely on SETADFs (recall Theorem 2). In particular, we will compare the signatures  $\Sigma_{\text{SETADF}}^{\sigma}$  and  $\Sigma_{\text{SFADF}}^{\sigma}$ , cf. Definition 10. We start with the observation that each SETADF can be rewritten as an equivalent SETADF that is also a SFADF.<sup>1</sup>

**Lemma 8.** *For each SETADF  $D = (S, L, C)$  there is an equivalent SETADF  $D' = (S, L', C')$  that is also a SFADF, i.e. for each  $s \in S$ ,  $\varphi_s \in C$ ,  $\varphi'_s \in C'$  we have  $\varphi_s \equiv \varphi'_s$ .*

*Proof.* Given a SETADF  $D$ , by Definition 8, each acceptance condition is a CNF over negative literals and thus does not have any support link which is not redundant. We can thus obtain  $L'$  by removing the redundant links from  $L$  and  $C'$  by, in each acceptance condition, deleting the clauses that are super-sets of other clauses.  $\square$

By the above we have that  $\Sigma_{\text{SETADF}}^{\sigma} \subseteq \Sigma_{\text{SFADF}}^{\sigma}$ . Now consider the interpretation  $\nu = \{a \mapsto \mathbf{f}\}$ . We have that for all considered semantics  $\sigma$ ,  $\nu$  is a  $\sigma$ -interpretation of the SFADF  $D = (\{a\}, \{\varphi_a = \perp\})$  but there is no SETADF with  $\nu$  being a  $\sigma$ -interpretation. We thus obtain  $\Sigma_{\text{SETADF}}^{\sigma} \subsetneq \Sigma_{\text{SFADF}}^{\sigma}$ .

**Theorem 9.**  $\Sigma_{\text{SETADF}}^{\sigma} \subsetneq \Sigma_{\text{SFADF}}^{\sigma}$ , for  $\sigma \in \{cf, adm, stb, mod, com, prf, grd\}$ .

In the remainder of this section we aim to characterise the difference between  $\Sigma_{\text{SETADF}}^{\sigma}$  and  $\Sigma_{\text{SFADF}}^{\sigma}$ . To this end we first recall a characterisation of the acceptance conditions of SFADF that can be rewritten as collective attacks.

**Lemma 10.** [16] *Let  $D = (S, L, C)$  be a SFADF. If  $s \in S$  has at least one incoming link then the acceptance condition  $\varphi_s$  can be written in CNF containing only negative literals.*

It remains to consider those arguments in an SFADF with no incoming links. Such arguments allow for only two acceptance conditions  $\top$  and  $\perp$ . While condition  $\top$  is

<sup>1</sup> As discussed in [6], in general, SETAFs translate to bipolar ADFs that contain attacking and redundant links. However, when we first remove redundant attacks from the SETAF we obtain a SFADF.

unproblematic (it refers to an initial argument in a SETAF), an argument with unsatisfiable acceptance condition cannot be modeled in a SETADF. In fact, the different expressiveness of SETADFs and SFADFs is solely rooted in the capability of SFADFs to set an argument to  $\mathbf{f}$  via a  $\perp$  acceptance condition.

We next give a generic characterisations of the difference between  $\Sigma_{\text{SETADF}}^\sigma$  and  $\Sigma_{\text{SFADF}}^\sigma$ .

**Theorem 11.** *For  $\sigma \in \{cf, adm, stb, mod, com, prf, grd\}$ , we have  $\Delta_\sigma = \Sigma_{\text{SFADF}}^\sigma \setminus \Sigma_{\text{SETADF}}^\sigma$  with*

$$\Delta_\sigma = \{\mathbb{V} \in \Sigma_{\text{SFADF}}^\sigma \mid \exists v \in \mathbb{V} \text{ s.t. } \forall a : v(a) \in \{\mathbf{f}, \mathbf{u}\} \wedge \exists a : v(a) = \mathbf{f}\}.$$

*Proof sketch.* First for  $\mathbb{V} \in \Delta_\sigma$  the interpretation  $v$  cannot be realized in a SETADF as we cannot have  $v(a) \in \mathbf{f}$  without  $v(b) \in \mathbf{t}$  for some other argument  $b$ . On the other hand one can show that when  $\mathbb{V} \in \Sigma_{\text{SFADF}}^\sigma$  is such that each  $v \in \mathbb{V}$  assigns some argument to  $\mathbf{t}$  one can construct a SETADF  $D$  with  $\sigma(D) = \mathbb{V}$ . This is by the fact that we can rewrite acceptance conditions via Lemma 10 and replace  $\perp$  acceptance conditions by collective attacks, i.e. for each interpretation we add collective attacks from the arguments set to  $\mathbf{t}$  to all argument with  $\perp$  acceptance condition.  $\square$

Next, we provide stronger characterisations of  $\Delta_\sigma$  for preferred and stable semantics.

**Proposition 12.** *For  $\mathbb{V} \in \Delta_\sigma$  and  $\sigma \in \{stb, mod, prf\}$  we have  $|\mathbb{V}| = 1$ . For  $\sigma \in \{stb, mod\}$  the unique  $v \in \mathbb{V}$  assigns all arguments to  $\mathbf{f}$ .*

*Proof sketch.* If a SFADF has a  $\sigma$ -interpretation  $v$  that assigns some arguments to  $\mathbf{f}$  without assigning an argument to  $\mathbf{t}$  then we have that the arguments assigned to  $\mathbf{f}$  are exactly the arguments with acceptance condition  $\perp$ . For *stb* and *mod* semantics this means all arguments have acceptance condition  $\perp$  and the result follows. Each preferred interpretation assigns arguments with acceptance condition  $\perp$  to  $\mathbf{f}$  and thus the existence of another preferred interpretation would violate the  $\leq_I$ -maximality of  $v$ .  $\square$

In other words each interpretation-set which is  $\sigma$ -realizable in SFADFs and contains at least two interpretations can be realized in SETADFs, for  $\sigma \in \{stb, prf, mod\}$ . We close this section with an example illustrating that the above characterisation thus not hold for *cf*, *adm*, and *com*.

**Example 3.** Let  $D = (\{a, b, c\}, \{\varphi_a = \perp, \varphi_b = \neg c, \varphi_c = \neg b\})$ . We have  $com(D) = \{\{a \mapsto \mathbf{f}, b \mapsto \mathbf{u}, c \mapsto \mathbf{u}\}, \{a \mapsto \mathbf{f}, b \mapsto \mathbf{t}, c \mapsto \mathbf{f}\}, \{a \mapsto \mathbf{f}, b \mapsto \mathbf{f}, c \mapsto \mathbf{t}\}\}$ . By Theorem 11,  $com(D)$  cannot be realized as SETADF. Moreover, as  $com(D) \subseteq adm(D) \subseteq cf(D)$  for every ADF  $D$ , we have that, despite all three contain more than one interpretation, none of them can be realized via a SETADF.

## 6 Discussion

In this paper, we have characterised the expressiveness of SETAFs under 3-valued signatures. The more fine-grained notion of 3-valued signatures reveals subtle differences of the expressiveness of stable and preferred semantics which are not present in the

2-valued setting [4] and enabled us to compare the expressive power of SETAFs and SFADFs, a subclass of ADFs that allows only for attacking links. In particular, we have exactly characterized the difference for conflict-free, admissible, complete, stable, preferred, and grounded semantics; this difference is rooted in the capability of SFADFs to set an initial argument to false. Together with our exact characterisations on signatures of SETAFs for stable, preferred, grounded, and conflict-free semantics, this also yields the corresponding results for SFADFs. Exact characterisations for admissible and complete semantics are subject of future work. Another aspect to be investigated is to which extent our insights on labelling-based semantics for SETAFs and SFADFs can help to improve the performance of reasoning systems.

**Acknowledgments** This research has been supported by FWF through projects I2854, P30168. The second researcher is currently embedded in the Center of Data Science & Systems Complexity (DSSC) Doctoral Programme, at the University of Groningen.

## References

- [1] Phan Minh Dung. On the acceptability of arguments and its fundamental role in nonmonotonic reasoning, logic programming and n-person games. *Artif. Intell.*, 77(2):321–357, 1995.
- [2] Gerhard Brewka, Sylwia Polberg, and Stefan Woltran. Generalizations of Dung frameworks and their role in formal argumentation. *IEEE Intelligent Systems*, 29(1):30–38, 2014.
- [3] Søren Holbech Nielsen and Simon Parsons. A generalization of Dung’s abstract framework for argumentation: Arguing with sets of attacking arguments. In *Proc. ArgMAS*, LNCS 4766, pages 54–73, 2006.
- [4] Wolfgang Dvořák, Jorge Fandinno, and Stefan Woltran. On the expressive power of collective attacks. *Argument & Computation*, 10(2):191–230, 2019.
- [5] Giorgos Flouris and Antonis Bikakis. A comprehensive study of argumentation frameworks with sets of attacking arguments. *Int. J. Approx. Reason.*, 109:55–86, 2019.
- [6] Sylwia Polberg. *Developing the abstract dialectical framework*. PhD thesis, TU Wien, Institute of Information Systems, 2017.
- [7] Bruno Yun, Srdjan Vesic, and Madalina Croitoru. Toward a more efficient generation of structured argumentation graphs. In *Proc. COMMA*, pages 205–212. IOS Press, 2018.
- [8] Paul E. Dunne, Wolfgang Dvořák, Thomas Linsbichler, and Stefan Woltran. Characteristics of multiple viewpoints in abstract argumentation. *Artif. Intell.*, 228:153–178, 2015.

- [9] Ringo Baumann and Gerhard Brewka. Extension removal in abstract argumentation - an axiomatic approach. In *Proc. AAAI*, pages 2670–2677. AAAI Press, 2019.
- [10] Gerhard Brewka, Stefan Ellmauthaler, Hannes Strass, Johannes P. Wallner, and Stefan Woltran. Abstract Dialectical Frameworks: An Overview. In *Handbook of Formal Argumentation*, chapter 5. College Publications, February 2018.
- [11] Martin Diller, Atefeh Keshavarzi Zafarghandi, Thomas Linsbichler, and Stefan Woltran. Investigating subclasses of abstract dialectical frameworks. *Argument & Computation*, 11(1), 2020.
- [12] Bart Verheij. Two approaches to dialectical argumentation: admissible sets and argumentation stages. *Proc. NAIC*, 96:357–368, 1996.
- [13] Martin W. A. Caminada and Dov M. Gabbay. A logical account of formal argumentation. *Studia Logica*, 93(2-3):109–145, 2009.
- [14] Thomas Linsbichler, Jörg Pührer, and Hannes Strass. A uniform account of realizability in abstract argumentation. In *Proc. ECAI*, pages 252–260. IOS Press, 2016.
- [15] Jörg Pührer. Realizability of three-valued semantics for abstract dialectical frameworks. *Artif. Intell.*, 278, 2020.
- [16] Johannes Peter Wallner. Structural constraints for dynamic operators in abstract argumentation. *Argument & Computation*, 11(1-2): 151-190, 2020.
- [17] Gerhard Brewka, Stefan Ellmauthaler, Hannes Strass, Johannes P. Wallner, and Stefan Woltran. Abstract dialectical frameworks revisited. In *Proc. IJCAI*, pages 803–809, 2013.
- [18] Atefeh Keshavarzi Zafarghandi. Investigating subclasses of abstract dialectical frameworks. Master’s thesis, TU Wien, 2017.
- [19] Sylwia Polberg. Understanding the abstract dialectical framework. In *Proc. JELIA*, LNCS 10021, pages 430–446, 2016.

## A Full Proofs

### Proof of Proposition 1

We first show the following result.

**Lemma 13.** *Let  $D = (S, L, C)$  be an ADF, let  $v$  be a model of  $D$  and let  $s \in S$  be an argument such that all parents of  $s$  are attackers. Thus,  $\varphi_s^v$  is irrefutable if and only if  $\varphi_s[p/\perp : v(p) = \mathbf{f}]$  is irrefutable.*

*Proof.* Assume that  $D = (S, L, C)$  is an ADF and  $v$  is a model of  $D$ . Further, assume  $s \in S$  such that  $\forall p \in \text{par}(s)$ ,  $(p, s)$  is an attacking link in  $D$ . Clearly if  $\varphi_s[p/\perp : v(p) = \mathbf{f}]$  is irrefutable then also  $\varphi_s^v = \varphi_s[p/\top : v(p) = \mathbf{t}][p/\perp : v(p) = \mathbf{f}]$  is irrefutable. It remains to show that if  $\varphi_s^v$  is irrefutable then also  $\varphi_s[p/\perp : v(p) = \mathbf{f}]$  is irrefutable. Let  $\varphi_s' = \varphi_s[p/\perp : v(p) = \mathbf{f}]$ . Towards a contradiction, assume that  $\varphi_s^v$  is irrefutable and  $\varphi_s'$  is not irrefutable. That is, either  $\varphi_s'$  is unsatisfiable or it is undecided. In both cases,  $\varphi_s'[p/\top : v(p) = \mathbf{t}]$  is unsatisfiable (as all the links are attacking). Thus,  $\varphi_s^v = \varphi_s'[p/\top : v(p) = \mathbf{t}]$  is unsatisfiable as well. This is a contradiction with the assumption that  $\varphi_s^v$  is irrefutable.  $\square$

*Proof of Proposition 1.* Let  $D = (S, L, C)$  be a SFADF. Since  $\text{stb}(D) \subseteq \text{mod}(D)$  for each ADF  $D$ , it remains to show that each model of  $D$  is also a stable model of  $D$ . Towards a contradiction assume that  $\text{mod}(D) \not\subseteq \text{stb}(D)$ . Thus, there exists a model  $v$  of  $D$  which is not a stable model. Let  $D^v$  be a *stb*-reduct of  $D$  and let  $w$  be the unique grounded interpretation of  $D^v$ . Since it is assumed that  $v$  is not a stable model,  $v^{\mathbf{t}} \neq w^{\mathbf{t}}$ . That is, there exists  $s \in S$  such that  $v(s) = \mathbf{t}$  and  $w(s) \neq \mathbf{t}$ . Thus,  $\varphi_s[p/\perp : v(p) = \mathbf{f}]$  is not irrefutable. Since,  $D$  is a SFADF, all parents of  $s$  are attackers. Hence, By Lemma 13,  $\varphi_s^v$  is not irrefutable, that is,  $v(s) \neq \mathbf{t}$ . This is a contradiction by the assumption that  $v(s) = \mathbf{t}$ . Thus, the assumption that  $D$  consists of a model which is not a stable model is incorrect.  $\square$

### Proof of Theorem 2

We first introduce some notation.

**Definition 11.** The function  $\text{Lab2Int}(\cdot)$  maps three-valued labellings to three-valued interpretations such that

- (a)  $\text{Lab2Int}(\lambda)(s) = \mathbf{t}$  iff  $\lambda(s) = \text{in}$ ,
- (b)  $\text{Lab2Int}(\lambda)(s) = \mathbf{f}$  iff  $\lambda(s) = \text{out}$ , and
- (c)  $\text{Lab2Int}(\lambda)(s) = \mathbf{u}$  iff  $\lambda(s) = \text{undec}$ .

For a labelling  $\lambda$  and an interpretation  $I$  we write  $\lambda \equiv I$  iff  $\text{Lab2Int}(\lambda) = I$ . For a set  $\mathcal{L}$  of labellings and a set  $\mathbb{V}$  of interpretations we write  $\mathcal{L} \equiv \mathbb{V}$  iff  $\{\text{Lab2Int}(\lambda) \mid \lambda \in \mathcal{L}\} = \mathbb{V}$ .

With the above notation we can restate Theorem 2 as follows: For a SETAF  $F$  and its associated SETADF  $D$  we have  $\sigma_{\mathcal{L}}(F) \equiv \sigma(D)$  for  $\sigma \in \{\text{cf}, \text{adm}, \text{com}, \text{prf}, \text{grd}, \text{stb}\}$ .



*Proof of Theorem 2.* Let  $F = (A, R)$  be a SETAF and  $D = (S, L, C)$  be its corresponding SETADF. We show that  $\{Lab2Int(\lambda) \mid \lambda \in \sigma_{\mathcal{L}}(F)\} = \sigma(D)$ . Let  $\lambda$  be an arbitrary three-valued labelling and let  $v = Lab2Int(\lambda)$ . We investigate that  $\lambda \in \sigma_{\mathcal{L}}(F)$  if and only if  $v \in \sigma(D)$ .

- Let  $\sigma = adm$ . We first assume that  $\lambda \in adm_{\mathcal{L}}(F)$  and show that  $v \in adm(D)$ . Consider  $s \in S$  and the acceptance condition  $\varphi_s = \bigwedge_{(B,s) \in R} \bigvee_{a \in B} \neg a$ . If  $v(s) = \mathbf{t}$  we have that  $\lambda(s) = \mathbf{in}$  and thus that for all  $(B, s) \in R$  there exists  $b \in B$  s.t.  $\lambda(b) = \mathbf{out}$ . The latter holds iff for all  $(B, s) \in R$  there exists  $b \in B$  s.t.  $v(b) = \mathbf{f}$  iff partial evaluation of  $\varphi_s$  under  $v$  is irrefutable iff  $\Gamma_D(v)(s) = \mathbf{t}$ . If  $v(s) = \mathbf{f}$  we have that  $\lambda(s) = \mathbf{out}$  and thus that there exists  $(B, s) \in R$  s.t. for all  $b \in B$ :  $\lambda(b) = \mathbf{in}$ . The latter holds iff there exists  $(B, s) \in R$  s.t. for all  $b \in B$ :  $v(b) = \mathbf{t}$  iff  $\varphi_s^v$  is unsatisfiable iff  $\Gamma_D(v)(s) = \mathbf{f}$ . We thus obtain that  $v \leq_i \Gamma_D(v)$  and therefore  $v \in adm(D)$ .

Now we assume  $v \in adm(D)$  and show that  $\lambda \in adm_{\mathcal{L}}(F)$ . That is for each  $s$  with  $\lambda(s) = \mathbf{in}$  we have  $\Gamma_D(v)(s) = \mathbf{t}$  and, as argued above, that for all  $(B, s) \in R$  there exists  $b \in B$  s.t.  $\lambda(b) = \mathbf{out}$ . Moreover for each  $s$  with  $\lambda(s) = \mathbf{out}$  we have  $\Gamma_D(v)(s) = \mathbf{f}$  and, as argued above, that there exists  $(B, s) \in R$  s.t. for all  $b \in B$ :  $\lambda(b) = \mathbf{in}$ . We obtain  $\lambda \in adm_{\mathcal{L}}(F)$ .

- Let  $\sigma \in \{com, prf, grd\}$ . Let  $\lambda \in com_{\mathcal{L}}(F)$  and let  $\varphi_s = \bigwedge_{(B,s) \in R} \bigvee_{a \in B} \neg a$  be the acceptance condition of  $s \in S$  in  $D$ . For complete semantics it is enough to show that  $\lambda(s) = \mathbf{in}$  iff  $\Gamma_D(v)(s) = \mathbf{t}$  and  $\lambda(s) = \mathbf{out}$  iff  $\Gamma_D(v)(s) = \mathbf{f}$ .
  - It holds that  $\lambda(s) = \mathbf{in}$  (i.e.  $v(s) = \mathbf{t}$ ) iff for all  $(B, s) \in R$  there exists  $b \in B$  s.t.  $\lambda(b) = \mathbf{out}$  iff for all  $(B, s) \in R$  there exists  $b \in B$  s.t.  $v(b) = \mathbf{f}$  iff partial evaluation of  $\varphi_s$  under  $v$  is irrefutable iff  $\Gamma_D(v)(s) = \mathbf{t}$ .
  - On the other hand,  $\lambda(s) = \mathbf{out}$  (i.e.  $v(s) = \mathbf{f}$ ) iff there exists  $(B, s) \in R$  s.t. for all  $b \in B$ :  $\lambda(b) = \mathbf{in}$  iff there exists  $(B, s) \in R$  s.t. for all  $b \in B$ :  $v(b) = \mathbf{t}$  iff  $\varphi_s^v$  is unsatisfiable iff  $\Gamma_D(v)(s) = \mathbf{f}$ .

Now as complete semantics coincide it is easy to verify that also the maximal, i.e. the preferred, extensions and the minimal, i.e. the grounded, extension coincide.

- Let  $\sigma = stb$ . Recall that, by Proposition 1, on SETADFs we have that stable and models semantics coincide. We will show that  $\lambda \in stb_{\mathcal{L}}(F)$  iff  $v \in mod(D)$ . That is we show that for each  $s \in S$  we have (i)  $\lambda(s) = \mathbf{in}$  iff  $v(\varphi_s) = \mathbf{t}$  and (ii)  $\lambda(s) = \mathbf{out}$  iff  $v(\varphi_s) = \mathbf{f}$ . To this end let  $\varphi_s = \bigwedge_{(B,s) \in R} \bigvee_{a \in B} \neg a$  be the acceptance condition of  $s$ .
  - It holds that  $\lambda(s) = \mathbf{in}$  (i.e.  $v(s) = \mathbf{t}$ ) iff for all  $(B, s) \in R$  there exists  $b \in B$  s.t.  $\lambda(b) = \mathbf{out}$  iff for all  $(B, s) \in R$  there exists  $b \in B$  s.t.  $v(b) = \mathbf{f}$  iff  $v(\varphi_s) = \mathbf{t}$ .
  - On the other hand,  $\lambda(s) = \mathbf{out}$  (i.e.  $v(s) = \mathbf{f}$ ) iff there exists  $(B, s) \in R$  s.t. for all  $b \in B$ :  $\lambda(b) = \mathbf{in}$  iff there exists  $(B, s) \in R$  s.t. for all  $b \in B$ :  $v(b) = \mathbf{t}$  iff  $v(\varphi_s) = \mathbf{f}$ .

- Finally let  $\sigma = cf$ . We first assume that  $\lambda \in cf_{\mathcal{L}}(F)$  and show that  $v \in cf(D)$ . Consider  $s \in S$  and the acceptance condition  $\varphi_s = \bigwedge_{(B,s) \in R} \bigvee_{a \in B} \neg a$ . If  $v(s) = \mathbf{t}$  we have that  $\lambda(s) = \mathbf{in}$  and thus that for all  $(B,s) \in R$  there exists  $b \in B$  s.t.  $\lambda(b) \neq \mathbf{in}$ . The latter holds iff for all  $(B,s) \in R$  there exists  $b \in B$  s.t.  $v(b) \neq \mathbf{t}$  iff  $\varphi_s^v$  is satisfiable. If  $v(s) = \mathbf{f}$  we have that  $\lambda(s) = \mathbf{out}$  and thus that there exists  $(B,s) \in R$  s.t. for all  $b \in B$ :  $\lambda(b) = \mathbf{in}$ . The latter holds iff there exists  $(B,s) \in R$  s.t. for all  $b \in B$ :  $v(b) = \mathbf{t}$  iff  $\varphi_s^v$  is unsatisfiable. We thus obtain that  $v \in cf(D)$ .

Now we assume  $v \in cf(D)$  and show that  $\lambda \in cf_{\mathcal{L}}(F)$ . That is for each  $s$  with  $\lambda(s) = \mathbf{in}$  we have  $\varphi_s^v$  is satisfiable and, as argued above, that for all  $(B,s) \in R$  there exists  $b \in B$  s.t.  $\lambda(b) \neq \mathbf{in}$ . Moreover for each  $s$  with  $\lambda(s) = \mathbf{out}$  we have  $\varphi_s^v$  is unsatisfiable and, as argued above, that there exists  $(B,s) \in R$  s.t. for all  $b \in B$ :  $\lambda(b) = \mathbf{in}$ . We obtain  $\lambda \in cf_{\mathcal{L}}(F)$ .  $\square$

#### Proof of Proposition 4

We first show that for each SETAF  $F$  the set  $prf_{\mathcal{L}}(F)$  satisfies the conditions of the proposition. The first condition is satisfied as clearly all  $\lambda \in prf_{\mathcal{L}}(F)$  have the same domain. Now, assume that  $\lambda \in prf_{\mathcal{L}}(F)$  assigns an argument  $a$  to  $\mathbf{out}$ . By the definition of conflict-free labellings there is an attack  $(B,a)$  such that all arguments  $b \in B$  are labeled  $\mathbf{in}$ . Thus Condition (2) is satisfied. For Condition (3), consider  $\lambda, \lambda' \in prf_{\mathcal{L}}(F)$ . Notice that there must be a conflict  $(S,a)$  with  $S \cup \{a\} \subseteq \lambda_{\mathbf{in}} \cup \lambda'_{\mathbf{in}}$  as otherwise  $(\lambda_{\mathbf{in}} \cup \lambda'_{\mathbf{in}}, \lambda_{\mathbf{out}} \cup \lambda'_{\mathbf{out}}, \lambda_{\mathbf{undec}} \cap \lambda'_{\mathbf{undec}})$  would be a larger admissible labelling. If  $a \in \lambda'_{\mathbf{in}}$  then, by the definition of admissible labellings, there is an attack  $(B,b)$  with  $B \subseteq \lambda'_{\mathbf{in}}$  and  $b \in S \cap \lambda_{\mathbf{in}}$ . Thus  $b$  is an argument with  $\lambda(b) = \mathbf{in}$  and  $\lambda'(b) = \mathbf{out}$ . Otherwise if  $a \in \lambda_{\mathbf{in}}$  then, by the definition of admissible labellings, there is an attack  $(B,b)$  with  $B \subseteq \lambda_{\mathbf{in}}$  and  $b \in S \cap \lambda'_{\mathbf{in}}$ . Then, again by the definition of admissible labellings, there is an attack  $(C,c)$  with  $C \subseteq \lambda'_{\mathbf{in}}$  and  $c \in B \subseteq \lambda_{\mathbf{in}}$ . Thus  $c$  is an argument with  $\lambda(c) = \mathbf{in}$  and  $\lambda'(c) = \mathbf{out}$ .

Now assume that  $\mathbb{L}$  satisfies all the conditions. We give a SETAF  $F_{\mathbb{L}} = (A_{\mathbb{L}}, R_{\mathbb{L}})$  with  $prf_{\mathcal{L}}(F_{\mathbb{L}}) = \mathbb{L}$ . We use

$$A_{\mathbb{L}} = \text{ARGS}_{\mathbb{L}}$$

$$R_{\mathbb{L}} = \{(\lambda_{\mathbf{in}}, a) \mid \lambda \in \mathbb{L}, \lambda(a) = \mathbf{out}\} \cup \{(\lambda_{\mathbf{in}} \cup \{a\}, a) \mid \lambda \in \mathbb{L}, \lambda(a) = \mathbf{undec}\}$$

We first show  $prf_{\mathcal{L}}(F_{\mathbb{L}}) \supseteq \mathbb{L}$ : Consider an arbitrary  $\lambda \in \mathbb{L}$ : We first show  $\lambda \in cf_{\mathcal{L}}(F_{\mathbb{L}})$ . We first consider  $\mathbf{out}$  labeled arguments. First, if  $\lambda(a) = \mathbf{out}$  for some argument  $a$  then by construction and Condition (2) we have an attack  $(\lambda_{\mathbf{in}}, a)$  and thus  $a$  is legally labeled  $\mathbf{out}$ . Now towards a contradiction assume there is a conflict  $(B,a)$  such that  $B \cup \{a\} \subseteq \lambda_{\mathbf{in}}$ .

If  $|\mathbb{L}| = 1$ , by the construction of  $F_{\mathbb{L}}$  there is no  $(B,a) \in R_{\mathbb{L}}$  such that  $a \in \lambda_{\mathbf{in}}$ . That is,  $a$  is legally labeled  $\mathbf{in}$ . If  $|\mathbb{L}| > 1$ , by construction there is a  $\lambda' \in \mathbb{L}$  with  $\lambda'_{\mathbf{in}} = B \setminus \{a\}$ , a contradiction to Condition (3). Thus,  $\lambda \in cf_{\mathcal{L}}(F_{\mathbb{L}})$ . Next we show that  $\lambda \in adm_{\mathcal{L}}(F_{\mathbb{L}})$ . Consider an argument  $a$  with  $\lambda(a) = \mathbf{in}$  and an attack  $(B,a)$ . Then, by construction there is a  $\lambda' \in \mathbb{L}$  with  $\lambda'_{\mathbf{in}} = B \setminus \{a\}$  and, by Condition (3), an argument  $b \in B$  such that  $\lambda(b) = \mathbf{out}$ . Thus,  $\lambda \in adm_{\mathcal{L}}(F_{\mathbb{L}})$ . Finally we show that  $\lambda \in prf_{\mathcal{L}}(F_{\mathbb{L}})$ . Towards a contradiction assume that there is a  $\lambda' \in adm_{\mathcal{L}}(F_{\mathbb{L}})$  with

$\lambda_{\text{in}} \subset \lambda'_{\text{in}}$ . Let  $a$  be an argument such that  $\lambda'(a) = \text{in}$  and  $\lambda(a) \in \{\text{out}, \text{undec}\}$ . By construction there is either an attack  $(\lambda_{\text{in}}, a)$  or an attack  $(\lambda_{\text{in}} \cup \{a\}, a)$ . In both cases  $\lambda' \notin \text{adm}_{\mathcal{F}}(F_{\mathbb{L}})$ , a contradiction. Hence,  $\lambda \in \text{prf}_{\mathcal{F}}(F_{\mathbb{L}})$ .

We complete the proof by showing  $\text{prf}_{\mathcal{F}}(F_{\mathbb{L}}) \subseteq \mathbb{L}$ : Consider  $\lambda \in \text{prf}_{\mathcal{F}}(F_{\mathbb{L}})$ : If  $\lambda$  maps all arguments to  $\text{in}$  then there is no attack in  $R_{\mathbb{L}}$  which means that  $\mathbb{L}$  contains only the labelling  $\lambda$ . Thus we can assume that  $\lambda(a) = \text{out}$  for some argument  $a$  and there is  $(B, a) \in R_{\mathbb{L}}$  with  $\lambda(b) = \text{in}$  for all  $b \in B$ . By construction there is  $\lambda' \in \mathbb{L}$  such that  $\lambda'_{\text{in}} = B$ . Then by construction we have  $(B, c) \in R_{\mathbb{L}}$  for all  $c$  with  $\lambda'(c) = \text{out}$  and  $(B \cup \{c\}, c) \in R_{\mathbb{L}}$  for all  $c$  with  $\lambda'(c) = \text{undec}$ . We obtain that  $\lambda'_{\text{in}} = B = \lambda_{\text{in}}$  and thus  $\lambda = \lambda'$ .

## Proof of Proposition 5

We first show that for each SETAF  $F$  the set  $\text{cf}_{\mathcal{F}}(F)$  satisfies the conditions of the proposition. The first condition is satisfied as clearly all  $\lambda \in \text{cf}_{\mathcal{F}}(F)$  have the same domain. Now, assume that  $\lambda \in \text{cf}_{\mathcal{F}}(F)$  assigns an argument  $a$  to  $\text{out}$ . By the definition of conflict-free labellings there is an attack  $(B, a)$  such that all arguments  $b \in B$  are labeled  $\text{in}$ . Thus Condition (2) is satisfied. For Condition (3), towards a contradiction assume that  $(C, \emptyset, \text{ARGS}_{\mathbb{L}} \setminus C)$  is not conflict-free. Then there is an attack  $(B, a)$  such that  $B \cup \{a\} \subseteq C$ . But then also  $B \cup \{a\} \subseteq \lambda_{\text{in}}$  and thus  $\lambda \notin \text{cf}_{\mathcal{F}}(F)$ , a contradiction. Condition (4) is satisfied as in the definition of conflict-free labellings there are no conditions for label an argument  $\text{undec}$ . Further, the conditions that allow to label an argument  $\text{out}$  solely depend on the  $\text{in}$  labeled arguments. Since  $\lambda_{\text{out}} \setminus C \subseteq \lambda_{\text{out}}$ , the condition for arguments labeled  $\text{out}$  is satisfied. For Condition (5) consider  $\lambda, \lambda' \in \text{cf}_{\mathcal{F}}(F)$  with  $\lambda_{\text{in}} \subseteq \lambda'_{\text{in}}$  and  $\lambda^* = (\lambda'_{\text{in}}, \lambda_{\text{out}} \cup \lambda'_{\text{out}}, \lambda_{\text{undec}} \cap \lambda'_{\text{undec}})$ . First there cannot be an attack  $(B, a)$  such that  $B \cup \{a\} \subseteq \lambda^*_{\text{in}}$  as  $\lambda' \in \text{cf}_{\mathcal{F}}(F)$ . Hence,  $\lambda'_{\text{in}} \cap \lambda_{\text{out}} = \emptyset$  and thus  $\lambda^*$  is a well-defined labelling. Moreover, for each  $a$  with  $\lambda^*(a) = \text{out}$  there is an attack  $(B, a)$  with  $B \subseteq \lambda^*_{\text{in}}$  as either  $\lambda(a) = \text{out}$  or  $\lambda'(a) = \text{out}$ . Thus,  $\lambda^* \in \text{cf}_{\mathcal{F}}(F)$  and therefore the condition holds. For Condition (6) consider  $\lambda, \lambda' \in \text{cf}_{\mathcal{F}}(F)$  and a set  $C \subseteq \lambda_{\text{out}}$  containing an argument  $a$  such that  $\lambda(a) = \text{out}$ . That is, there is an attack  $(B, a)$  with  $B \subseteq \lambda_{\text{in}}$  and thus  $\lambda_{\text{in}} \cup C \not\subseteq \lambda'$ . That is, Condition (6) is satisfied.

Now assume that  $\mathbb{L}$  satisfies all the conditions. We give a SETAF  $F_{\mathbb{L}} = (A_{\mathbb{L}}, R_{\mathbb{L}})$  satisfying  $\text{cf}_{\mathcal{F}}(F_{\mathbb{L}}) = \mathbb{L}$ , where

$$\begin{aligned} A_{\mathbb{L}} &= \text{ARGS}_{\mathbb{L}} \\ R_{\mathbb{L}} &= \{(\lambda_{\text{in}}, a) \mid \lambda \in \mathbb{L}, \lambda(a) = \text{out}\} \cup \{(B, b) \mid b \in B, \nexists \lambda \in \mathbb{L} : \lambda_{\text{in}} = B\} \end{aligned}$$

We first show  $\text{cf}_{\mathcal{F}}(F_{\mathbb{L}}) \supseteq \mathbb{L}$ : Consider an arbitrary  $\lambda \in \mathbb{L}$ : First, if  $\lambda(a) = \text{out}$  for some argument  $a$  then by construction and Condition (2) we have an attack  $(\lambda_{\text{in}}, a)$  and thus  $a$  is legally labeled  $\text{out}$ . Now towards a contradiction assume there is a conflict  $(B, a)$  such that  $B \cup \{a\} \subseteq \lambda_{\text{in}}$ . By Condition (3) it cannot be the case that  $a \in B$ . Thus, by construction there is a  $\lambda' \in \mathbb{L}$  with  $\lambda'_{\text{in}} = B$ , a contradiction to Condition (6). Thus,  $\lambda \in \text{cf}_{\mathcal{F}}(F_{\mathbb{L}})$ .

We complete the proof by showing  $\text{cf}_{\mathcal{F}}(F_{\mathbb{L}}) \subseteq \mathbb{L}$ : Consider  $\lambda \in \text{cf}_{\mathcal{F}}(F_{\mathbb{L}})$ : If  $\lambda$  maps all arguments to  $\text{in}$  then there is no attack in  $R_{\mathbb{L}}$  which means that  $\mathbb{L}$  contains only the labelling  $\lambda$ . Thus we can assume that  $\lambda(a) \in \{\text{out}, \text{undec}\}$  for some argument  $a$ . If

$\lambda_{\text{in}} \neq \lambda'_{\text{in}}$  for all  $\lambda' \in \mathbb{L}$  then by construction of the second part of  $R_{\mathbb{L}}$  there would be attacks  $(\lambda_{\text{in}}, b)$  for all  $b \in \lambda_{\text{in}}$ , which is in contradiction to  $\lambda \in \text{cf}_{\mathcal{G}}(F_{\mathbb{L}})$ . Thus, there is  $\lambda' \in \mathbb{L}$  such that  $\lambda'_{\text{in}} = \lambda_{\text{in}}$ . For arguments  $a$  with  $\lambda(a) = \text{out}$  there is an attack  $(B, a)$  with  $B \subseteq \lambda_{\text{in}}$  and, by construction, a  $\lambda^* \in \mathbb{L}$  such that  $\lambda^*_{\text{in}} = B$  and  $\lambda^*(a) = \text{out}$ . By the existence of  $\lambda' \in \mathbb{L}$  and Condition (5) we have that there exists  $\lambda'' \in \mathbb{L}$  such that  $\lambda_{\text{in}} = \lambda''_{\text{in}}$ ,  $\lambda'_{\text{out}} \subseteq \lambda''_{\text{out}}$  and  $a \in \lambda''_{\text{out}}$ . By iteratively applying this argument for each argument  $a$  with  $\lambda(a) = \text{out}$  we obtain that there is a labelling  $\hat{\lambda} \in \mathbb{L}$  such that  $\lambda_{\text{in}} = \hat{\lambda}_{\text{in}}$  and  $\lambda_{\text{out}} \subseteq \hat{\lambda}_{\text{out}}$ . By Condition (4) we obtain that  $\lambda \in \mathbb{L}$ .

## Proof of Proposition 6

We first show that for each SETAF  $F$  the set  $\text{grd}_{\mathcal{G}}(F)$  satisfies the conditions of the proposition. Towards a contradiction assume that there are  $\lambda, \lambda' \in \text{grd}_{\mathcal{G}}$  with  $\lambda \neq \lambda'$ . By the definition of grounded labelling  $\lambda_{\text{in}}, \lambda'_{\text{in}}$  are  $\subseteq$ -minimal among all complete labellings, thus,  $\lambda_{\text{in}} = \lambda'_{\text{in}}$ . Assume that  $\lambda_{\text{out}} \subset \lambda'_{\text{out}}$ . Since each grounded labelling is conflict-free, for each  $a$  with  $a \in \lambda'_{\text{out}}$  there is  $(B, a)$  such that  $B \subseteq \lambda'_{\text{in}}$ . Since  $\lambda_{\text{in}} = \lambda'_{\text{in}}$ ,  $a \in \lambda_{\text{out}}$ . Therefore,  $\lambda = \lambda'$ . Now, assume that  $\lambda \in \text{grd}_{\mathcal{G}}(F)$  assigns an argument  $a$  to out. By the definition of conflict-free labeling there is an attack  $(B, a)$  such that  $B \subseteq \lambda_{\text{in}}$ .

Now assume that  $\mathbb{L}$  satisfies all the conditions. We give a SETAF  $F_{\mathbb{L}} = (A_{\mathbb{L}}, R_{\mathbb{L}})$  with  $\text{grd}_{\mathcal{G}}(F_{\mathbb{L}}) = \mathbb{L}$ . We set

$$\begin{aligned} A_{\mathbb{L}} &= \text{ARGS}_{\mathbb{L}} \\ R_{\mathbb{L}} &= \{(\lambda_{\text{in}}, a) \mid \lambda \in \mathbb{L}, \lambda(a) = \text{out}\} \cup \{(\lambda_{\text{in}} \cup \{a\}, a) \mid \lambda \in \mathbb{L}, \lambda(a) = \text{undec}\} \end{aligned}$$

Consider the unique  $\lambda \in \mathbb{L}$  and the unique  $\lambda^G \in \text{grd}_{\mathcal{G}}(F_{\mathbb{L}})$ . For each argument  $a \in \lambda_{\text{in}}$  we have that  $a$  is not attacked in  $F_{\mathbb{L}}$  and thus  $a \in \lambda^G_{\text{in}}$ . For each argument  $a \in \lambda_{\text{out}}$  there is an attack  $(\lambda_{\text{in}}, a)$  in  $F_{\mathbb{L}}$  and as  $\lambda_{\text{in}} \subseteq \lambda^G_{\text{in}}$  by the definition of complete labellings we have  $a \in \lambda^G_{\text{out}}$ . Finally for each argument  $a \in \lambda_{\text{undec}}$  the attack  $(\lambda_{\text{in}} \cup \{a\}, a)$  is the only attack towards  $a$  in  $F_{\mathbb{L}}$ . Thus, by the definition of complete labellings, we have that  $a$  is neither labelled in nor out in  $F_{\mathbb{L}}$  and therefore  $a \in \lambda^G_{\text{undec}}$ . We obtain that  $\lambda^G = \lambda$  and thus  $\text{grd}_{\mathcal{G}}(F_{\mathbb{L}}) = \mathbb{L}$ .

## Proof of Theorem 9

$\Sigma_{\text{SETADF}}^{\sigma} \subseteq \Sigma_{\text{SFADF}}^{\sigma}$  follows from Lemma 8. For showing  $\Sigma_{\text{SETADF}}^{\text{adm}} \subsetneq \Sigma_{\text{SFADF}}^{\text{adm}}$ , let  $\mathbb{V} = \{\{a \mapsto \mathbf{u}, b \mapsto \mathbf{u}\}, \{a \mapsto \mathbf{u}, b \mapsto \mathbf{f}\}, \{a \mapsto \mathbf{t}, b \mapsto \mathbf{f}\}\}$  be an interpretation-set. A witness of adm-realizability of  $\mathbb{V}$  in SFADFs is  $D = (\{a, b\}, \{\varphi_a = \neg a \vee \neg b, \varphi_b = \perp\})$ . However,  $\mathbb{V}$  is not realizable by any SETADF for admissible interpretations (cf. Proposition 7). To show  $\Sigma_{\text{SFADF}}^{\sigma} \not\subseteq \Sigma_{\text{SETADF}}^{\sigma}$ , for  $\sigma \in \{\text{stb}, \text{mod}, \text{com}, \text{prf}, \text{grd}\}$ , let  $\mathbb{V} = \{\{a \mapsto \mathbf{f}\}\}$ . The interpretation  $\mathbb{V}$  is  $\sigma$ -realizable in SFADFs for  $\sigma \in \{\text{stb}, \text{mod}, \text{com}, \text{prf}, \text{grd}\}$ , and a witness of  $\sigma$ -realizability of  $\mathbb{V}$  in SFADFs is  $D = (\{a\}, \{\varphi_a = \perp\})$ . However,  $\mathbb{V}$  cannot be realized by any SETADF for semantics  $\sigma \in \{\text{adm}, \text{stb}, \text{prf}, \text{grd}\}$  (cf. Propositions 3–6). The result for  $\sigma = \text{mod}$  follows from Proposition 1 and for  $\sigma = \text{com}$  by  $|\mathbb{V}| = 1$  (i.e. complete and grounded semantics have to coincide). Further,  $\text{cf}(D)$  is not cf-realizable with any SETADF.

**Lemma 14.** *Given an interpretation-set  $\mathbb{V} \in \Delta_\sigma$ , for  $\sigma \in \{adm, stb, mod, com, prf, grd\}$ . Let  $v \in \mathbb{V}$  be a non-trivial interpretation in which  $v(a) = \mathbf{f}/\mathbf{u}$ , for each argument  $a$ . In all SFADFs that realize  $\mathbb{V}$  under  $\sigma$ , the acceptance conditions of all arguments assigned to  $\mathbf{f}$  by  $v$  are equal to  $\perp$ .*

*Proof.* Let  $D$  be a SFADF that realizes  $\mathbb{V}$  under  $\sigma$ , for  $\sigma \in \{adm, stb, mod, com, prf, grd\}$ . Let  $v \in \mathbb{V}$  be an non-trivial interpretation that assigns all arguments either to  $\mathbf{f}$  or  $\mathbf{u}$ . Towards a contradiction, assume that there exists an argument  $a$  which is assigned to  $\mathbf{f}$  by  $v$ , and  $\varphi_a \neq \perp$  in  $D$ . First we show that  $\mathbb{V}$  cannot be *adm*-realizable in SFADFs. Since  $a$  is assigned to  $\mathbf{f}$  in  $v$  the acceptance condition of  $a$  cannot be equal to  $\top$ . By Lemma 10, the acceptance condition of  $a$  is in CNF and having only negative literals. Since all  $b \in par(a)$  are either assigned to  $\mathbf{f}$  or  $\mathbf{u}$  by  $v$ ,  $\varphi_a^v$  cannot be unsatisfiable. That is,  $v(a) \not\leq_i \Gamma_D(v)(a)$ . Therefore,  $v$  is not an admissible interpretation of  $D$ . Thus, any  $\mathbb{V}$  that contains  $v$  is not *adm*-realizable in SFADF. To complete the proof it remains to see that for each of the remaining semantics, each  $\sigma$ -interpretation is also admissible.  $\square$

### Proof of Theorem 11

To show that  $\Delta_\sigma = \{\mathbb{V} \in \Sigma_{SFADF}^\sigma \mid \exists v \in \mathbb{V} \text{ s.t. } \forall a : v(a) \in \{\mathbf{f}, \mathbf{u}\} \wedge \exists a : v(a) = \mathbf{f}\}$ , let  $\mathbb{V}$  be an arbitrary interpretation-set of  $\Delta_\sigma$ . By the definition of  $\Delta_\sigma$ ,  $\mathbb{V} \in \Sigma_{SFADF}^\sigma$  and  $\mathbb{V} \notin \Sigma_{SETADF}^\sigma$ . It remains to show that there exists  $v \in \mathbb{V}$  that assigns at least an argument to  $\mathbf{f}$  but none of the arguments to  $\mathbf{t}$ . Towards a contradiction, assume that there exists no such interpretation and let  $D = (S, L, C)$  be an arbitrary SFADF with  $\sigma(SFADF) = \mathbb{V}$ . Notice that by Lemma 10 all acceptance conditions of  $D$  that are not equal to  $\perp$  can be transformed to be in SETADF form. Thus we can focus on the arguments with acceptance condition  $\perp$ . As, under the above assumption, each  $v \in \mathbb{V}$  that assigns an argument to  $\mathbf{f}$  also assigns an argument  $b$  to  $\mathbf{t}$  it is easy to verify that we can replace  $\perp$  acceptance conditions by  $\bigwedge_{s \in S} \neg s$  without changing the semantics. That is, we can transform  $D$  to an equivalent SETADF and thus  $\mathbb{V} \in \Sigma_{SETADF}^\sigma$ . This is a contradiction by the definition of  $\Delta_\sigma$  and we obtain that there exists  $v \in \mathbb{V}$  that assigns all arguments to either  $\mathbf{f}$  or  $\mathbf{u}$ .

On the other hand, let  $\mathbb{V}$  be an interpretation-set that is  $\sigma$ -realizable in SFADF such that there exists  $v \in \mathbb{V}$  that assigns at least one argument to  $\mathbf{f}$  and none of the arguments to  $\mathbf{t}$ . We show that  $\mathbb{V} \notin \Sigma_{SETADF}^\sigma$ . By Lemma 14, in any SFADF with  $\sigma(SFADF) = \mathbb{V}$  the acceptance conditions of all arguments assigned to  $\mathbf{f}$  by  $v$  are equal to  $\perp$ . Therefore,  $D$  is not  $\sigma$ -realizable in any SETADF. That is,  $\mathbb{V} \in \Delta_\sigma$ .

### Proof of Proposition 12

Consider  $\mathbb{V} \in \Delta_\sigma$ , for  $\sigma \in \{stb, mod, prf\}$  and let  $v \in \mathbb{V}$  be an interpretation that assigns all arguments to either  $\mathbf{f}$  or  $\mathbf{u}$  (since  $\mathbb{V} \in \Delta_\sigma$ , such a  $v$  exists). By Lemma 14, the acceptance condition of all arguments that are assigned to  $\mathbf{f}$  by  $v$  is equal to  $\perp$  in all SFADFs that realize  $\mathbb{V}$  under  $\sigma \in \{stb, mod, prf\}$ . Let  $D = (S, L, C)$  be a witness of  $\sigma$ -realizability of  $\mathbb{V}$  in SFADFs, under  $\sigma \in \{stb, mod, prf\}$ .

First, if all arguments are assigned to  $\mathbf{f}$  in  $v$ , the acceptance conditions of all arguments are  $\perp$  in SFADF  $D$  and  $|\sigma(D)| = 1$ . Now assume that  $v$  assigns some arguments

to  $\mathbf{u}$ . Thus,  $V$  cannot be *mod* or *stb*-realized in any ADF. It remains to consider *prf* semantics. Let  $B = \{s \in S \mid v(b) = \mathbf{u}\}$ . For each  $s \in S \setminus B$ , by Lemma 14,  $\varphi_s = \perp$  in  $D$ . Therefore, in all  $v' \in \mathbb{V}$ ,  $v'(s) = \mathbf{f}$  for  $s \in S \setminus B$ . For each  $v' \neq v$  in  $\mathbb{V}$  there exists at least  $b \in B$  such that  $v'(b) \neq \mathbf{u}$ , therefore,  $v < v'$ . By the definition of preferred interpretations  $v$  cannot be a preferred interpretation. Thus,  $|\text{prf}(D)| = 1$  and therefore, the assumption  $|\mathbb{V}| = 1$ . Summarizing the two cases we have that interpretation set  $\mathbb{V} \in \Delta_\sigma$ , for  $\sigma \in \{\text{stb}, \text{mod}, \text{prf}\}$  consist of only one interpretation.