# A Formalisation of Abstract Argumentation in Higher-Order Logic

Alexander Steen   David Fuenmayor

{*alexander.steen, david.fuenmayor*}*@uni.lu*
*University of Luxembourg, Department of Computer Science*
*6, avenue de la Fonte*
*L-4364 Esch-sur-Alzette, Luxembourg*

**Abstract**

We present an approach for representing abstract argumentation frameworks based on an encoding into classical higher-order logic. This provides a uniform framework for computer-assisted assessment of abstract argumentation frameworks using interactive and automated reasoning tools. This enables the formal analysis and verification of meta-theoretical properties as well as the flexible generation of extensions and labellings with respect to well-known argumentation semantics.

*Keywords:* Abstract Argumentation, Higher-Order Logic, Automated Reasoning, Proof Assistants, Isabelle/HOL.

## 1  Introduction

Argumentation theory is a relevant and active field of research in artificial intelligence. Argumentation frameworks [37] constitute the central concept in abstract argumentation. An argumentation framework essentially is a directed graph in which the nodes of the graph represent abstract arguments (carrying no further structural properties) and the edges of the graph represent attacks between arguments. The exact interpretation of the arguments depends on the application context: Argumentation frameworks have many topical applications in, among others, non-monotonic reasoning, logic programming and multi-agent systems [5]. As an example, in non-monotonic reasoning the abstract arguments may be regarded as defeasible reasons (or proofs) for certain claims; and the attack relation then formalises which of them act as counter-arguments against others.

Since the original formulation of Dung in the 1990s, a lot of research has been conducted concerning algorithmic procedures, complexity aspects, as well as various extended and related formalisms, cf. [6] and references therein. In this paper, we propose to investigate argumentation frameworks from the perspective of extensional type theory (ExTT), also commonly simply referred to as higher-order logic [17]. To that end, we present a novel encoding of argumentation frameworks and their semantics into higher-order logic, enabling

the employment of off-the-shelf automated theorem provers. We argue that this constitutes a uniform approach to assess argumentation frameworks under different argumentation semantics, while at the same time enabling computer-assisted exploration and verification of meta-theoretical properties within the same logical formalism. We furthermore argue that our approach is flexible in the sense that it allows the instantiation of the abstract arguments with arbitrary structures, and to compute acceptable subsets of arguments satisfying complex (higher-order) properties. Up to the authors' knowledge, there does not exist any other formalisation of abstract argumentation frameworks within higher-order logic and, in particular, existing proof assistants. Although there exists related encodings into less expressive logical formalisms, these cannot allow for both object-level and meta-level reasoning within the same framework (cf. §6 for a more thorough discussion of related work).

In particular, we exemplarily demonstrate how our approach can be used to (a) flexibly synthesise extensions and labellings for argumentation frameworks and (b) to conduct (explorative) assessment of meta-theoretical properties of argumentation frameworks. The experiments presented in this article were conducted using the well-established proof assistant Isabelle/HOL [49]. The corresponding Isabelle/HOL source files for this work are freely available at GitHub [44].

This article is an extended version of an earlier version of this work [43]. It contains the following novel contributions:

 (i) The encoding is generalised, compared to [43], to allow for interactive and automated reasoning with instantiated argumentation frameworks. This is done by defining a *relativised* notion of argumentation semantics (cf. §3).

 (ii) Important meta-theoretical properties of the presented encoding are now fully verified in Isabelle/HOL for both extension-based semantics and labellings-based semantics (cf. §4).

(iii) The application of interactive proof assistants for meta-theoretical reasoning is exemplified by formally assessing the adequacy of the formalised notions within the Isabelle/HOL proof assistant, and by applying its integrated automated tools for theory exploration (cf. §4).

(iv) The flexibility of our approach is demonstrated on interactive experiments, including the generation of extensions and labellings for abstract argumentation frameworks and the assessment on the existence of such extensions (cf. §5).

The remainder of this article is structured as follows: Technical preliminaries about higher-order logic and abstract argumentation are introduced in §2. The encoding of abstract argumentation frameworks and their semantics is presented in §3. Subsequently, §4 and §5 present the meta-theoretical assessment of the presented encoding and its applications, respectively. Finally, §6 concludes and discusses related and further work.

## 2 Preliminaries

We start this section with a brief exposition of a higher-order logic (HOL) loosely adopted from earlier work of the first author [53]. Subsequently, we introduce the notion of abstract argumentation frameworks. In the present work, the former formalism is employed as an expressive logical language in order to encode the latter.

### 2.1 Higher-Order Logic

The term *higher-order logic* refers in general to expressive logical formalisms that allow for quantification over predicate and function variables. In the context of automated reasoning, higher-order logic commonly refers to systems based on a simply typed $\lambda$-calculus, as originally introduced in the works of Church, Henkin and several others [34,47]. In the present work, higher-order logic (abbreviated as HOL) is used interchangeably with Henkin's Extensional Type Theory, cf. [53, §2], which constitutes the basis of most contemporary higher-order automated reasoning systems. HOL provides $\lambda$-notation as an expressive binding mechanism to denote unnamed functions, predicates and sets (by their characteristic functions), and it comes with built-in principles of Boolean and functional extensionality as well as type-restricted comprehension (cf. further below).

**Syntax and Semantics.** HOL is a typed logic; and all terms of HOL get assigned a fixed and unique type. The set $\mathcal{T}$ of types is freely generated from a set of base types $\mathcal{BT}$ and the function type constructor $\Rightarrow$ (written as a right-associative infix operator).[1] Traditionally, the generating set $\mathcal{BT}$ is taken to include at least two base types, $\mathcal{BT} \subseteq \{\iota, o\}$, where $\iota$ is interpreted as the type of individuals and $o$ as the type of (bivalent) Boolean truth values.

HOL terms of are given by the following abstract syntax ($\tau, \nu \in \mathcal{T}$):

$$s, t ::= c_\tau \in \Sigma \mid X_\tau \in \mathcal{V} \mid (\lambda X_\tau . \, s_\nu)_{\tau \Rightarrow \nu} \mid (s_{\tau \Rightarrow \nu} \, t_\tau)_\nu$$

where $\Sigma$ is a set of constant symbols and $\mathcal{V}$ a set of variable symbols. The different forms of terms above are called *constants*, *variables*, *abstractions* and *applications*, respectively. We assume that $\Sigma$ contains equality predicate symbols $=^\tau_{\tau \Rightarrow \tau \Rightarrow o}$ for each $\tau \in \mathcal{T}$. All remaining logical connectives (including conjunction $\wedge_{o \Rightarrow o \Rightarrow o}$, disjunction $\vee_{o \Rightarrow o \Rightarrow o}$, material implication $\longrightarrow_{o \Rightarrow o \Rightarrow o}$, negation $\neg_{o \Rightarrow o}$, universal quantification for predicates over type $\tau$ denoted $\Pi^\tau_{(\tau \Rightarrow o) \Rightarrow o}$) can be defined as abbreviations using equality and the other syntactical structures [53, §2.1].[2]

---

[1] In order to minimise syntactical differences with respect to the formalisation in Isabelle/HOL, we use to symbol $\Rightarrow$ to denote the function type constructor (despite the fact that other different yet similar symbols are often used in the literature). For the same reason, $\longrightarrow$ will denote material implication throughout the article. This will help to avoid confusion between the different (meta-logical) arrow-like symbols used in Isabelle/HOL.

[2] It is worth noting that in HOL there is no strict differentiation between formulas and terms, as in first-order formalisms. Terms of type $o$ are customarily referred to as "formulas".

For simplicity, the binary logical connectives may be written in infix notation, e.g., the term/formula $p_o \vee q_o$ formally represents the application $(\vee_{o \Rightarrow o \Rightarrow o} \ p_o \ q_o)$. Also, so-called *binder notation* [17] is used for universal and existential quantification: The term $\forall X_\tau. s_o$ is used as a short-hand for $\Pi^\tau_{(\tau \Rightarrow o) \Rightarrow o} (\lambda X_\tau. s_o)$ and analogously for existential quantification $\exists X_\tau. s_o$. To improve readability, type-subscripts and parentheses are usually omitted if there is no risk of confusion. Note that, by construction, HOL syntax only admits functions that take one parameter; $n$-ary function applications are represented using *currying* [17], e.g., a first-order-like term such as $f(a, b)$ involving a binary function $f$ and two constants $a$ and $b$ is represented in HOL by consecutive applications of the individual constants, as in $((f_{\iota \Rightarrow \iota \Rightarrow \iota} \ a_\iota) \ b_\iota)$, or simply $f \ a \ b$ if omitting parentheses and type subscripts. Here, the term $(f \ a)$ itself represents a function that is subsequently applied to the argument $b$. Also, functional terms may be only *partially applied*, e.g., occurring in terms like $(g_{(\iota \Rightarrow \iota) \Rightarrow \iota} \ (f \ a))$, where $f$ is the "binary" function from above and $g_{(\iota \Rightarrow \iota) \Rightarrow \iota}$ is a higher-order function symbol taking a functional expression of type $\iota \Rightarrow \iota$ as argument.

HOL automation is usually investigated with respect to so-called *general semantics*, due to Henkin [47], for which complete proof calculi can be achieved. Note that standard models for HOL are subsumed by general models such that every valid formula with respect to general semantics is also valid in the standard sense. We omit the formal exposition to HOL semantics at this point and instead refer to the literature (cf., e.g., [53,17] and the references therein). For the remainder of this article, HOL with general semantics is assumed.

**HOL automation.** Automated theorem proving (ATP) systems are computer programs that, given a set of assumptions and a conjecture, try to prove that the conjecture is a logical consequence of the assumptions. This is done completely autonomously, i.e., without any interaction from the outside by the user. In contrast, proof assistants – also called interactive theorem provers – allow for the computer-assisted creation and assessment of verified formal proofs, and also facilitate interactive (and possibly incremental) experiments with such formal representations. In the interactive scenario, it is the user that will manually construct, formalise and enter the proofs into the system. These proofs are then assessed for correctness by the system. Proof assistants are usually based on (extensions of) higher-order formalisms. Isabelle/HOL [49] is a well-established HOL-based proof assistant that is employed in a wide range of applications, including this work. Further well-known proof assistants include, e.g., Coq, Lean, HOL4 and HOL-Light.

One of the most relevant practical features of Isabelle/HOL is the Sledgehammer system [21] that bridges between the proof assistant and external ATP systems, such as the first-order ATP system E [52] or the higher-order ATP system Leo-III [55], and SMT solvers such as Z3 [35] and CVC4 [7]. The idea

---

Analogously, terms of type $\tau \Rightarrow o$ (for $\tau \in \mathcal{T}$) are suggestively called "predicates" (over type $\tau$).

is to use these systems to automatically resolve open proof obligations and to import the resulting proofs into the verified context of Isabelle/HOL. The employment of Sledgehammer is of great practical importance and usually a large amount of laborious proof engineering work can be solved by the ATP systems. In fact, most of the formal proofs presented in the remainder of this article were automatically constructed using Sledgehammer. Additionally, Isabelle/HOL integrates so-called *model finders* such as Nitpick [23] that can generate (counter-)models to given formulas. Also, most non-theorems presented in this work were automatically refuted by Nitpick.

## 2.2 Abstract Argumentation

The subsequent brief introduction of abstract argumentation frameworks largely follows the survey paper by Baroni, Caminada and Giacomin [3] with occasional references to Dung's seminal paper [37]. In the present treatment we will not, however, presuppose finiteness for classes of arguments. Any required cardinality assumptions for argumentation frameworks will be stated explicitly when necessary. We refer the interested reader to [10] (and references therein) for further details on infinite argumentation frameworks.

In the theory of abstract argumentation of Dung [37], arguments are represented as abstract objects and constitute the nodes of a directed graph.

**Definition 2.1** An *argumentation framework* $AF$ is a pair $AF = (\mathcal{A}, \rightsquigarrow)$, where $\mathcal{A}$ is a set (finite or infinite) and $\rightsquigarrow \subseteq A \times A$ is a binary relation on $\mathcal{A}$. The elements of $\mathcal{A}$ are called *arguments*, and $\rightsquigarrow$ is called the *attack relation*.

An argumentation framework formally captures how arguments interact (i.e., conflict with each other). Given an argumentation framework $AF$, the primary problem is to determine subsets of arguments that can reasonably be accepted together; those sets are called *extensions*. Restrictions on this selection are imposed by so-called argumentation semantics. The set of designators (names) for the different argumentation semantics is denoted $\mathcal{S}em$ in the following; the restrictions they impose on the set of extensions are then assigned by an interpretation function.

**Definition 2.2** An *extension-based semantics interpretation* $\mathcal{E}$ associates with each argumentation framework $AF = (\mathcal{A}, \rightsquigarrow)$ and an argumentation semantics $S \in \mathcal{S}em$ a set of *extensions*, denoted $\mathcal{E}_S(AF)$, where $\mathcal{E}_S(AF) \subseteq 2^{\mathcal{A}}$.

Roughly speaking, each $E \in \mathcal{E}_S(AF)$ is a subset of arguments that can be accepted (under the criterion specified by $\mathcal{E}_S$), while all arguments in $\mathcal{A} \setminus E$ are rejected. In fact, we will show in §3 how to encode the criteria imposed by $\mathcal{E}_S$ as HOL predicates for several well-known semantics in the literature [3].

Alternatively, argumentation semantics can be specified in terms of labelling functions (this approach can be traced back to [25]). We loosely follow [3] below:

**Definition 2.3** Let $AF = (\mathcal{A}, \rightsquigarrow)$ be an argumentation framework. A *labelling* of $AF$ is a total function $\mathcal{L}ab : \mathcal{A} \Rightarrow \{\texttt{In}, \texttt{Out}, \texttt{Undec}\}$; the set of all labellings of $AF$ is denoted $\mathfrak{L}(AF)$. A *labelling-based semantics interpretation* $\mathcal{L}$ then asso-

ciates with each $AF$ and argumentation semantics $S \in \mathcal{S}em$ a set of *labellings*, denoted $\mathcal{L}_S(AF)$, where $\mathcal{L}_S(AF) \subseteq \mathfrak{L}(AF)$.

Intuitively, the labels `In` and `Out` represent the status of accepting and rejecting a given argument, respectively. Arguments labelled `Undec` are left undecided, either because one explicitly refrains from accepting resp. rejecting it, or because it cannot be labelled otherwise. Given a labelling $\mathcal{L}ab$ we write ($in\ \mathcal{L}ab$), ($out\ \mathcal{L}ab$), and ($undec\ \mathcal{L}ab$) (read as *in-set*, *out-set*, *undec-set*, respectively) for the subset of arguments labelled by $\mathcal{L}ab$ as `In`, `Out` and `Undec`, respectively.

In this paper we will work mainly with the different extension-based semantics introduced by Dung, together with their labelling-based counterparts, following the exposition by [3]. Accepted sets of arguments under each of these semantics are termed: *conflict-free*, *admissible*, *complete*, *grounded*, *preferred* and *stable* [37] extensions.[3] We will also consider some further semantics, namely, *semi-stable* [28], *ideal* [38] and *stage* [58] semantics. For each of these semantics there exists an equivalent labelling-based formulation. In fact, there is a one-to-one correspondence between the extension-based semantics and their labelling counterparts for the semantics listed below, so that for each extension a corresponding labelling can be found and vice versa [29,3]. This is also why the names of the argumentation semantics stand for both extension-based and labelling-based approaches, i.e., in the remainder of this article we consider the set of argumentation semantics $\mathcal{S}$ to be defined as follows:

$$\mathcal{S}em = \{\texttt{conflictfree}, \texttt{admissible}, \texttt{complete}, \texttt{grounded},$$
$$\texttt{preferred}, \texttt{stable}, \texttt{semistable}, \texttt{ideal}, \texttt{stage}\}$$

We omit the formal definitions for the different extension- and labelling-based semantics at this point, as they will be subject of the discussion in §3.

## 3    Encoding of Abstract Argumentation in HOL

In this section we discuss the encoding of argumentation semantics in HOL (as introduced in §2.1). For the sake of the formal assessment of the encoding and the verification of our results, the proof assistant Isabelle/HOL [49] is employed. A few remarks are in order.

In the discussion below, we will often employ type variables that stand for fixed but arbitrary (base or functional) types. Following Isabelle/HOL's notation, type variables will be represented using letters preceded by a single quote, e.g., '$a$ is a type variable. Throughout this work we make generous use of definitions as understood in the context of Isabelle/HOL. A *definition* defines a new symbol that can, for the purposes of this paper, be regarded as an abbreviation for the respective terms. We write $c := s$ to denote the

---

[3] For reasons of uniformity, we refer to all types of 'acceptable' argument sets (according to some given criteria) as *extensions*. In particular, this includes conflict-free sets and admissible sets as well, to which we will also refer as *extensions* in the following.
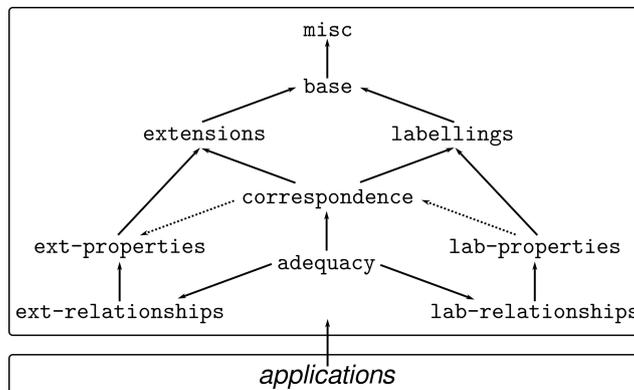
Fig. 1. Structure of the encoding as implemented in Isabelle/HOL. The individual nodes of the graph represent topically self-contained parts of the overall encoding. The solid arrows indicate a dependency relation in which the respective part of the encoding reuses notions and definitions of the underlying part (e.g., both the definitions of the different labellings and extensions use low-level definitions collected in the `base` theory). A dotted arrow indicates auxiliary usage where certain lemmas are used as parts of larger proofs.

introduction of a new symbol $c$, with definition $s$, where $s$ is some HOL term. A *type synonym* is similar to a term definition but rather introduces a new type symbol that abbreviates a (complex) type expression.

We will often mention several results concerning abstract argumentation as having been proven using Isabelle/HOL. By this we mean formal and internally verified proofs that have been automatically generated by different theorem proving systems integrated into the Isabelle proof assistant. [4]

The Isabelle/HOL sources for the presented encoding into HOL [44] have been hierarchically organised into several different files (also referred to as *theory* files). The general layout is depicted in Fig. 1. Theory `misc` collects general purpose notions such as set-theoretic definitions and notions related to orderings (cf. §3.1). Building on this, the theory file `base` contains general definitions related to abstract argumentation frameworks (cf. §3.2). The different argumentation semantics are then defined in `extensions` for extension-based semantics and in `labellings` for labelling-based semantics (cf. §3.3 and §3.4, respectively). The remaining theory files contain meta-theoretical results of the presented encodings, including formal proofs of correspondences between extension-based and labelling-based semantics, the relationship between different argumentation semantics, and further fundamental properties. The meta-theoretical assessment results are presented separately in §4 and are collected

---

[4] Isabelle allows for the manual formalisation of proofs using the general purpose proof language *Isar* [60]. Isabelle also supports the invocation of external *state-of-the-art* theorem provers via *Sledgehammer* [22]. These provers can be run on a local installation or remotely via *System on TPTP* (`http://www.tptp.org/cgi-bin/SystemOnTPTP`).

in theory file `adequacy`; the remainder of this section focuses on the encoding in HOL itself.

### 3.1   Basic Notions for Sets and Orderings

We start by introducing useful type synonyms for the types of sets and relations, which will be represented by predicates on objects of some type $`a$. We thus define $`a$ `Set` and $`a$ `Rel` as type synonyms for the functional (predicate) types $`a \Rightarrow o$ and $`a \Rightarrow `a \Rightarrow o$, respectively. Set-theoretic operations can be defined by anonymous functions making use of the respective underlying HOL logical connectives:

$$\cap := \lambda A. \lambda B. \lambda x. (A\ x) \wedge (B\ x) \qquad \cup := \lambda A. \lambda B. \lambda x. (A\ x) \vee (B\ x)$$
$$- := \lambda A. \lambda x. \neg (A\ x)$$

where $\cap$ and $\cup$ are both terms of type $`a$ `Set` $\Rightarrow `a$ `Set` $\Rightarrow `a$ `Set` and $-$ is of type $`a$ `Set` $\Rightarrow `a$ `Set`. By convention we write binary operations defined this way as infix operators in the remainder, e.g., we write $A \cup B$ instead of $(\cup\ A\ B)$.

Set equality and the subset relation, and also a few other notions defined further below, will often be used in a *relativised* fashion in the remainder of this article, i.e., restricted to a certain subset $D$ of elements of interest. This is mainly for technical reasons: In HOL, a type $`a$ intuitively represents a set of objects that inhabit the type. Given some type $`a$ representing the arguments in our encoding of abstract argumentation, we need to be able to represent that only a subset of all possible objects of type $`a$ are considered as the domain of arguments (denoted by $\mathcal{A}$; cf. §2.2) in the argumentation framework. Intuitively, the extra parameter $D$ (of type $`a$ `Set`) restricts the domain of set-theoretic and logical operations to those elements (usually those contained in $\mathcal{A}$). In a dependently typed logical formalism as employed, e.g., by the proof assistant Coq, it is possible to encode this information as part of the type of the operations. In the simply-typed discipline of HOL, however, this needs to be encoded as part of the term language. This design decision also allows for the instantiation of the abstract argumentation network with arbitrary (complex) objects, cf. §3.2 for a more thorough discussion on this topic. [5]

Relativised set equality $\approx^D$ and the relativised subset relation $\subseteq^D$, both of type $`a$ `Set` $\Rightarrow `a$ `Set` $\Rightarrow `a$ `Set` $\Rightarrow o$, are defined by restricting the domain of quantification over the elements in $D$ only, in the following way:

$$\approx^D := \lambda A. \lambda B. \forall x. (D\ x) \longrightarrow \big((A\ x) \longleftrightarrow (B\ x)\big)$$
$$\subseteq^D := \lambda A. \lambda B. \forall x. (D\ x) \longrightarrow \big((A\ x) \longrightarrow (B\ x)\big)$$

---

[5]  In the Isabelle/HOL sources, we also provide a simplified, non-relativised variant of these operations. They can be employed when *all objects of the given type* are considered to be the arguments under consideration. This comes handy in several applications in abstract argumentation, but not in general (e.g. when working with instantiated, structured arguments). For reasons of conciseness and legibility, we omit the presentation of these variants in the article, as they can be obtained seamlessly by simply dropping the superscripts ($D$, resp. $\mathcal{A}$).

In a similar spirit, we can define domain-restricted versions of the logical quantifiers, as featured in *free logics*.[6] Thus, we define the restricted (or relativised) universal quantification predicate $\Pi^D$ as follows:

$$\Pi^D := \lambda P. \forall x. (D\ x) \longrightarrow (P\ x)$$

We again employ binder notation in the relativised case and write $\forall^D x.\ P$ instead of $\Pi^D (\lambda X.\ P)$ and let $\exists^D x.\ P := \neg(\forall^D x.\ \neg(P\ x))$.

For the sake of illustration, we note that the above definitions for relativised equality and inclusion can be, alternatively, encoded in a more succinct fashion by employing restricted quantification (we will continue doing so in the sequel).

$$\approx^D\ =\ \lambda A.\ \lambda B.\ \forall^D x.\ (A\ x) \longleftrightarrow (B\ x)$$
$$\subseteq^D\ =\ \lambda A.\ \lambda B.\ \forall^D x.\ (A\ x) \longrightarrow (B\ x)$$

Because of their importance in various argumentation semantics, we additionally provide generic notions for representing minimal and maximal (resp. least and greatest) sets, with respect to set inclusion: Let $Obj$ be a term of some type 'a, and $Prop$ a predicate of type 'a $\Rightarrow$ o. We formalise the statement that the set $\varphi(Obj)$ induced by $Obj$ is minimal/maximal/least/greatest among all objects $O$ satisfying property $Prop$ wrt. a domain of quantification $D$ as follows:

$$\mathtt{minimal}^D := \lambda Prop.\ \lambda Obj.\ \lambda \varphi.\ (Prop\ Obj)\ \wedge$$
$$\left(\forall O.\ \big((Prop\ O) \wedge (\varphi\ O) \subseteq^D (\varphi\ Obj)\big) \longrightarrow (\varphi\ O) \approx^D (\varphi\ Obj)\right)$$
$$\mathtt{maximal}^D := \lambda Prop.\ \lambda Obj.\ \lambda \varphi.\ (Prop\ Obj)\ \wedge$$
$$\left(\forall O.\ \big((Prop\ O) \wedge (\varphi\ Obj) \subseteq^D (\varphi\ O)\big) \longrightarrow (\varphi\ O) \approx^D (\varphi\ Obj)\right)$$
$$\mathtt{least}^D := \lambda Prop.\ \lambda Obj.\ \lambda \varphi.\ (Prop\ Obj)\ \wedge$$
$$\big(\forall O.\ (Prop\ O) \longrightarrow (\varphi\ Obj) \subseteq^D (\varphi\ O)\big)$$
$$\mathtt{greatest}^D := \lambda Prop.\ \lambda Obj.\ \lambda \varphi.\ (Prop\ Obj)\ \wedge$$
$$\big(\forall O.\ (Prop\ O) \longrightarrow (\varphi\ O) \subseteq^D (\varphi\ Obj)\big)$$

We formally verified in Isabelle/HOL several well-known properties of least (greatest) and minimal (maximal) sets while successfully obtaining countermodels for non-theorems using model finder Nitpick. As an example, it has been formally verified that least and greatest elements are unique and that minimal (maximal) elements collapse to the least (greatest) one when the latter exist [44, misc]:

---

[6] Free logics [50] are quantified, non-classical logics in which terms do not necessarily denote objects in the domain of quantification (sometimes suggestively called "existing" objects). Common approaches introduce domain-restricted versions of universal/existential quantifiers, leaving out one or more objects. Free logics have previously been employed to model different notions of partiality in the encoding of axiomatic category theory in Isabelle/HOL [16].

**Lemma 3.1** *For every predicate $P$ of type $\text{'}a\ \mathtt{Set} \Rightarrow o$, elements $O$, $O'$ of type $\text{'}a$ and transformation function $\varphi$ of type $\text{'}a \Rightarrow \text{'}b\ \mathtt{Set}$ it holds that*

(i) $(\mathtt{least}^D\ P\ O\ \varphi) \wedge (\mathtt{least}^D\ P\ O'\ \varphi) \longrightarrow (\varphi\ O) \approx^D (\varphi\ O')$

(ii) $(\mathtt{greatest}^D\ P\ O\ \varphi) \wedge (\mathtt{greatest}^D\ P\ O'\ \varphi) \longrightarrow (\varphi\ O) \approx^D (\varphi\ O')$

(iii) $(\mathtt{least}^D\ P\ O\ \varphi) \wedge (\mathtt{minimal}^D\ P\ O'\ \varphi) \longrightarrow \mathtt{least}^D\ P\ O'\ \varphi$

(iv) $(\mathtt{greatest}^D\ P\ O\ \varphi) \wedge (\mathtt{maximal}^D\ P\ O'\ \varphi) \longrightarrow \mathtt{least}^D\ P\ O'\ \varphi$

<div align="right">□</div>

Monotonicity of functions over sets is a property that plays an important role in argumentation. This notion is, again, relativised and encoded as follows:

$$\mathtt{MONO}^D := \lambda F.\, \forall A.\, \forall B.\, (A \subseteq^D B) \longrightarrow (F\ A) \subseteq^D (F\ B)$$

We formalised a *fixed point* notion; namely, we speak of sets of arguments being fixed points of operations on sets (e.g. the so-called *characteristic function* of argumentation frameworks [37]). For a given operation $\varphi$ (of type $\text{'}a\ \mathtt{Set} \Rightarrow \text{'}a\ \mathtt{Set}$) we define in the usual way a fixed-point predicate on sets:

$$\mathtt{fixpoint}^D := \lambda\varphi.\, \lambda X.\, (\varphi\ X) \approx^D X$$

We formally verified a weak formulation of the Knaster-Tarski theorem, whereby any monotone function on a powerset lattice has a least (greatest) fixed point [44, $\mathtt{misc}$]:

**Lemma 3.2** *For every function $f$ of type $\text{'}a\ \mathtt{Set} \Rightarrow \text{'}a\ \mathtt{Set}$ it holds that*

(i) $(\mathtt{MONO}^D\ f) \longrightarrow \exists S.\, \mathtt{least}^D\ (\mathtt{fixpoint}^D\ f)\ S\ \mathtt{id}$

(ii) $(\mathtt{MONO}^D\ f) \longrightarrow \exists S.\, \mathtt{greatest}^D\ (\mathtt{fixpoint}^D\ f)\ S\ \mathtt{id}$

<div align="right">□</div>

## 3.2   Basic Notions for Abstract Argumentation

We encode definitions involving some given argumentation framework $AF = (\mathcal{A}, \rightsquigarrow)$ as HOL terms (i.e., functions) that take as parameters, among others, the encoded set of arguments $\mathcal{A}$ of type $\text{'}a\ \mathtt{Set}$ and the encoded attack relation $att$ of type $\text{'}a\ \mathtt{Rel}$. Note that, for reasons of legibility, the attack relation $att$ will often be referred to as $\rightsquigarrow$ in infix notation, i.e., $a \rightsquigarrow b$ formally stands for $(att\ a\ b)$. Sets of arguments are represented by HOL terms of type $\text{'}a\ \mathtt{Set}$. It is worth noting that an abstract argumentation framework $AF$ can, in principle, be completely characterised in HOL by simply encoding its underlying attack relation $att$; this way, the underlying set of arguments $\mathcal{A}$ is given implicitly as the collection of objects of type $\text{'}a$ (i.e., the carrier of $att$) [43]. In order to generalise from this specific case, we do not consider this simplification in the following. This subject is motivated and discussed at the end of this section.

For a given set of arguments $S \subseteq \mathcal{A}$, we encode its set of attacked ($S^+$) and attacking ($S^-$) arguments in a relativised fashion as follows:

$$[\mathcal{A}|att|S]^+ := \lambda b.\, \exists^{\mathcal{A}} a.\, (S\ a) \wedge a \rightsquigarrow b \qquad [\mathcal{A}|att|S]^- := \lambda b.\, \exists^{\mathcal{A}} a.\, (S\ a) \wedge b \rightsquigarrow a$$

We now encode the fundamental notion of *defense* of arguments (called *acceptability* by Dung [37]): We say that a set of arguments $S$ *defends* an argument $a$ iff each argument $b \in \mathcal{A}$ attacking $a$ is itself attacked by at least one argument $z$ in $S$. This is formalised as a HOL predicate `defends` of type '$a$ `Set` $\Rightarrow$ '$a$ `Rel` $\Rightarrow$ '$a$ `Set` $\Rightarrow$ '$a$ $\Rightarrow o$ by:

$$\mathtt{defends}^{\mathcal{A}} = \lambda att. \, \lambda S. \, \lambda a. \, \forall^{\mathcal{A}} b. \, b \rightsquigarrow a \longrightarrow (\exists^{\mathcal{A}} z. \, (S \, z) \wedge z \rightsquigarrow b)$$

In fact, it can be verified automatically in Isabelle/HOL that the condition imposed by $\mathtt{defends}^{\mathcal{A}} \, att \, S \, a$ can equivalently expressed by $S^{+}$ and $S^{-}$ as follows [44, `base`, ll. 53-54]:

**Lemma 3.3** *For every argumentation framework encoded by* $(\mathcal{A}, att)$, *subset of arguments $S$ and argument $a \in \mathcal{A}$ it holds that*

$$\mathtt{defends}^{\mathcal{A}} \, att \, S \, a \longleftrightarrow [\mathcal{A}|att|\{a\}]^{-} \subseteq^{\mathcal{A}} [\mathcal{A}|att|S]^{+}$$

$\square$

Due to the fact that sets are represented in HOL by predicates (i.e., terms with functional types returning an element of type $o$), the notion of *characteristic function* $\mathcal{F}^{\mathcal{A}}$ of an argumentation framework [37] is actually an extensionally equivalent alias for the function $\mathtt{defends}^{\mathcal{A}}$, yielding:

$$\mathcal{F}^{\mathcal{A}} := \lambda att. \, \lambda S. \, \lambda a. \, \mathtt{defends}^{\mathcal{A}} \, att \, S \, a$$

It is also formally verified that $\mathcal{F}^{\mathcal{A}}$ (i.e., $\mathtt{defends}^{\mathcal{A}}$) is indeed a monotone function and that it has both a least and a greatest fixed point, drawing upon the previously formalised Knaster-Tarski theorem.

**Lemma 3.4** *For every argumentation framework encoded by* $(\mathcal{A}, att)$ *it holds that*

(i) $\mathtt{MONO}^{\mathcal{A}} \, (\mathcal{F}^{\mathcal{A}} \, att)$

(ii) $\exists S. \, \mathtt{least}^{\mathcal{A}} \, \big(\mathtt{fixpoint}^{\mathcal{A}} \, (\mathcal{F}^{\mathcal{A}} \, att)\big) \, S \, id$

(iii) $\exists S. \, \mathtt{greatest}^{\mathcal{A}} \, \big(\mathtt{fixpoint}^{\mathcal{A}} \, (\mathcal{F}^{\mathcal{A}} \, att)\big) \, S \, id$

*where* $id := \lambda x. \, x$ *is the identity function.* $\square$

Recall that argument sets (i.e., potential argument *extensions*) are encoded as functions mapping objects of an arbitrary type '$a$ (i.e., arguments) to the bivalent Boolean type $o$. Generalising this, we can now define argument *labellings* as functions into some arbitrary but finite co-domain of *labels*. Following the usual approach in the literature [3], we assume a fixed set of three labels $\{\mathtt{In}, \mathtt{Out}, \mathtt{Undec}\}$. This is encoded in HOL as the three-valued type `Label`:[7]

$$\mathtt{Label} := \mathtt{In} \mid \mathtt{Out} \mid \mathtt{Undec},$$

---

[7] For convenience, in our formalisation work [44] we have encoded `Label` as an Isabelle/HOL datatype, noting that any datatypes can be, in turn, encoded into plain HOL.

together with the type synonym $‘a\ \texttt{Labelling} := ‘a \Rightarrow \texttt{Label}$ as type of labellings.

It is convenient to encode the basic notions of in-set, out-set and undec-set of an labelling $\mathcal{L}ab$:

$$\texttt{in} := \lambda \mathcal{L}ab.\,\lambda a.\,(\mathcal{L}ab\ a) = \texttt{In}$$
$$\texttt{out} := \lambda \mathcal{L}ab.\,\lambda a.\,(\mathcal{L}ab\ a) = \texttt{Out}$$
$$\texttt{undec} := \lambda \mathcal{L}ab.\,\lambda a.\,(\mathcal{L}ab\ a) = \texttt{Undec}$$

Using these definitions above we can represent the *as-is-state* of a given argument $a$ wrt. a given labelling; for instance, $(\texttt{in}\ \mathcal{L}ab\ a)$ means that argument $a$ is labelled $\texttt{In}$ by $\mathcal{L}ab$.

Moreover, we also want to represent a *target-state* in which an argument $a$ is said to be adequately or *legally* labelled. For our particular purposes, we slightly generalise the usual definitions [3] and say of an argument $a$ that it is *legally in* if all of its attackers are labelled $\texttt{Out}$. Similarly, $a$ is said to be *legally out* if it has at least one attacker that is labelled $\texttt{In}$. Finally, $a$ is said to be *legally undecided* if it is neither *legally in* nor *legally out*. These notions are encoded in HOL as follows:

$$\texttt{legallyIn}^{\mathcal{A}} := \lambda att.\,\lambda \mathcal{L}ab.\,\lambda a.\,\forall^{\mathcal{A}} b.\,(b \rightsquigarrow a) \longrightarrow (\texttt{out}\ \mathcal{L}ab)\ b$$
$$\texttt{legallyOut}^{\mathcal{A}} := \lambda att.\,\lambda \mathcal{L}ab.\,\lambda a.\,\exists^{\mathcal{A}} b.\,(b \rightsquigarrow a) \wedge (\texttt{in}\ \mathcal{L}ab)\ b$$
$$\texttt{legallyUndec}^{\mathcal{A}} := \lambda att.\,\lambda \mathcal{L}ab.\,\lambda a.\,\neg(\texttt{legallyIn}^{\mathcal{A}}\ att\ \mathcal{L}ab\ a) \wedge \neg(\texttt{legallyOut}^{\mathcal{A}}\ att\ \mathcal{L}ab\ a)$$

**Remark on Relativisation.** At this point, it might not be apparent what are the benefits of using relativised encodings, as opposed to an earlier version of this work in which simpler definitions were presented [43]. Recall that HOL is a typed formalism in which each term of the language is associated with a unique and fixed type. Assuming that abstract arguments are presented by terms of type $‘a$ in HOL, then the attack relation $att$ is a term of type $‘a\ \texttt{Rel}$ which abbreviates the type $‘a \Rightarrow ‘a \Rightarrow o$. In particular, in [43], argumentation frameworks were completely characterised by their attack relation; the set of underlying arguments was implicitly assumed to be the carrier set of the attack relation, i.e., *all* objects of type $‘a$.

While this seems appealing for assessing properties of argumentation semantics from an abstract perspective, the simplified approach is too rigid when concrete arguments are being studied, e.g., when $‘a$ is instantiated with a type representing formulas of some logical language. For the sake of the argument, let us assume that $‘a$ is instantiated with a type $PROP$ of classical propositional logic formulae. [8] It is clear that the type $PROP$ is generally inhabited by infinitely many objects since infinitely many syntactically different propositional

---

[8] Propositional logic can easily be encoded into HOL via so-called *deep embeddings* in which a new type $PROP$ is postulated, and axiomatised inductively to contain the respective syntactical elements of propositional logic formulae. See, e.g., the work by Michaelis and Nipkow on encoding propositional logic in Isabelle/HOL [48].

logic formulae can be constructed (assuming a non-empty set of propositional variables). This, in turn, implies that every argumentation framework instantiated with $PROP$ and encoded by the simplified representation from [43] will be of infinite size, as there is no possibility to control which objects of type $PROP$ are contained in the set of arguments $\mathcal{A}$ and which are not.

Clearly this is an undesired effect of implicitly representing the set of arguments as the carrier of the attack relation. This is mitigated by the here presented relativisation in which only a subset of arguments of a certain type are assumed to be members of the argumentation framework (i.e., those included in the set $\mathcal{A}$). While relativisation slightly complicates the encoding itself, it allows for a more fine-grained control of participating arguments and automatically allows for instantiating the abstract arguments with arbitrary objects, enabling the assessment of *instantiated* argumentation frameworks using the very same encoding. While instantiation is not the primary aim of this work, we focus on the relativised encoding to allow for future extensions with instantiated arguments.

### 3.3 Extension-based semantics

The well-known extension-based semantics by Dung [37] have been encoded drawing upon the notions introduced in the previous section. For each of the discussed semantics in this section, we first give an informal definition, leaving the underlying argumentation framework $AF = (\mathcal{A}, att)$ implicit. We then complement those informal definitions with their corresponding formalisation in HOL, which we encode in their most general form (i.e., relativised wrt. the underlying domain or universe of arguments $\mathcal{A}$).

**Conflict-free and admissible extensions**

**Definition 3.5** A set of arguments $S$ is termed *conflict-free* if it does not contain two arguments that attack each other.

**Definition 3.6** A set of arguments $S$ is termed *admissible* if it is conflict-free and defends each of its arguments.

**Formalisation.** Employing the notions encoded in §3.2, we can formalise the interpretations $\mathcal{E}_{\texttt{conflictfree}}$ and $\mathcal{E}_{\texttt{admissible}}$ as HOL predicates `conflictfreeExt` and `admissibleExt`, respectively, of type `'a Set ⇒ 'a Rel ⇒ 'a Set ⇒ o`:

$$\texttt{conflictfreeExt}^{\mathcal{A}} := \lambda att. \lambda S. \forall^{\mathcal{A}} a. \forall^{\mathcal{A}} b. (S\ a) \wedge (S\ b) \longrightarrow \neg(a \rightsquigarrow b)$$

$$\texttt{admissibleExt}^{\mathcal{A}} := \lambda att. \lambda S. (\texttt{conflictfreeExt}^{\mathcal{A}}\ att\ S) \wedge$$
$$\forall^{\mathcal{A}} a. (S\ a) \longrightarrow (\texttt{defends}^{\mathcal{A}}\ att\ S\ a).$$

**Complete extensions**

**Definition 3.7** An set of arguments $S$ is called a *complete extension* if it is admissible and contains each argument defended by it.

**Formalisation.** The above definition is analogously encoded in HOL:

$$\texttt{completeExt}^{\mathcal{A}} := \ \lambda att.\, \lambda S.\, (\texttt{admissibleExt}^{\mathcal{A}}\ att\ S) \land$$
$$\forall^{\mathcal{A}} a.\, (\texttt{defends}^{\mathcal{A}}\ att\ S\ a) \longrightarrow (S\ a).$$

### Preferred and grounded extensions

We now discuss complete extensions which are maximal or minimal wrt. set inclusion. They are termed *preferred* and *grounded* extensions respectively.

**Definition 3.8** A set of arguments $S$ is termed a *preferred extension* if it is a maximal (wrt. set inclusion) complete extension.

**Definition 3.9** A set of arguments $S$ is termed a *grounded extension* if it is a minimal (wrt. set inclusion) complete extension.

**Formalisation.** The above definitions are encoded as HOL predicates in an analogous fashion:

$$\texttt{preferredExt}^{\mathcal{A}} := \ \lambda att.\, \lambda S.\, (\texttt{maximal}^{\mathcal{A}}\ (\texttt{completeExt}^{\mathcal{A}}\ att)\ S\ id)$$
$$\texttt{groundedExt}^{\mathcal{A}} := \ \lambda att.\, \lambda S.\, (\texttt{minimal}^{\mathcal{A}}\ (\texttt{completeExt}^{\mathcal{A}}\ att)\ S\ id)$$

Recalling the definitions of maximality (minimality) from §3.1, these unfold to:

$$\texttt{preferredExt}^{\mathcal{A}}\ att\ S\ =\ \texttt{completeExt}^{\mathcal{A}}\ att\ S\ \land$$
$$(\forall E.\ \texttt{completeExt}^{\mathcal{A}}\ att\ E \land S \subseteq^{\mathcal{A}} E \longrightarrow E \approx^{\mathcal{A}} S)$$
$$\texttt{groundedExt}^{\mathcal{A}}\ att\ S\ =\ \texttt{completeExt}^{\mathcal{A}}\ att\ S\ \land$$
$$(\forall E.\ \texttt{completeExt}^{\mathcal{A}}\ att\ E \land E \subseteq^{\mathcal{A}} S \longrightarrow E \approx^{\mathcal{A}} S).$$

### Ideal extensions

We will be concerned with admissible sets of arguments that are contained in every preferred extension, which we call *ideal sets*. In any $AF$ the family of ideal sets has indeed an unique maximal (hence greatest) element, which is termed the *ideal extension* [38].

**Definition 3.10** An *ideal set* is an admissible set of arguments that is a subset of every preferred extension. The (unique) maximal/greatest *ideal set* is called the *ideal extension*.

**Formalisation.** The above are encoded, analogously, as HOL predicates:

$$\texttt{idealSet}^{\mathcal{A}} := \ \lambda att.\, \lambda S.\, (\texttt{admissibleExt}^{\mathcal{A}}\ att\ S) \land$$
$$\forall E.\, (\texttt{preferredExt}^{\mathcal{A}}\ att\ E) \longrightarrow S \subseteq^{\mathcal{A}} E$$
$$\texttt{idealExt}^{\mathcal{A}} := \ \lambda att.\, \lambda S.\, (\texttt{greatest}^{\mathcal{A}}\ (\texttt{idealSet}^{\mathcal{A}}\ att)\ S\ id)$$

### Stable, semi-stable and stage extensions

For convenience of exposition we define the *range* of a set of arguments $S$ as the union of $S$ with the set $S^{+}$ of its attacked arguments. We now turn to sets of arguments whose range satisfies particular maximality requirements.

**Definition 3.11** A set of arguments $S$ is termed a *stable extension* if it is conflict-free and its range is the whole universe (i.e., every possible argument belongs either to $S$ or to $S^+$).

Stable extensions are in fact complete. However, in contrast to the previous extensions, they do not always exist [37]. Semi-stable extensions were introduced independently by Verheij [57] and Caminada [26] as an approximate, existence-entailing notion.

**Definition 3.12** A set of arguments $S$ is termed a *semi-stable extension* if it is a complete extension and its range is maximal among all complete extensions.

While the notion of semi-stable extensions aims at maximising the *range* under the condition of admissibility, stage extensions do so under the (weaker) condition of conflict-freeness.

**Definition 3.13** A set of arguments $S$ is termed a *stage extension* if it is conflict-free and its range is maximal among all conflict-free sets of arguments.

**Formalisation.** The definition of *range* is encoded in HOL as a function of type '$a$ Set $\Rightarrow$ '$a$ Rel $\Rightarrow$ '$a$ Set $\Rightarrow$ '$a$ Set, while the extension predicates are encoded analogously as before (terms of type '$a$ Set $\Rightarrow$ '$a$ Rel $\Rightarrow$ '$a$ Set $\Rightarrow o$):

$$\texttt{range}^{\mathcal{A}} := \lambda att. \lambda S.\ S \cup [\mathcal{A}|att|S]^+$$
$$\texttt{stableExt}^{\mathcal{A}} := \lambda att. \lambda S.\ (\texttt{conflictfreeExt}^{\mathcal{A}}\ att\ S) \wedge \mathcal{A} \subseteq (\texttt{range}^{\mathcal{A}}\ att\ S)$$
$$\texttt{semistableExt}^{\mathcal{A}} := \lambda att. \lambda S.\ \texttt{maximal}^{\mathcal{A}}\ (\texttt{completeExt}^{\mathcal{A}}\ att)\ S\ (\texttt{range}^{\mathcal{A}}\ att)$$
$$\texttt{stageExt}^{\mathcal{A}} := \lambda att. \lambda S.\ \texttt{maximal}^{\mathcal{A}}\ (\texttt{conflictfreeExt}^{\mathcal{A}}\ att)\ S\ (\texttt{range}^{\mathcal{A}}\ att)$$

### 3.4 Labelling-based semantics

Analogous to the previous section, we provide informal definitions followed by their corresponding formalisation in HOL, relativised wrt. the underlying domain or universe $\mathcal{A}$.

**Conflict-free and admissible labellings**

**Definition 3.14** A labelling $\mathcal{L}ab$ is termed *conflict-free* if (i) every In-labelled argument is not *legally out*; and (ii) every Out-labelled argument is *legally out*.

**Definition 3.15** A labelling $\mathcal{L}ab$ is termed *admissible* if (i) every In-labelled argument is *legally in*; and (ii) every Out-labelled argument is *legally out*.

**Formalisation.** Employing the notions encoded in §3.2 for argument labellings, we can formalise the interpretations $\mathcal{L}_{\texttt{conflictfree}}$ and $\mathcal{L}_{\texttt{admissible}}$ as HOL predicates $\texttt{conflictfreeLab}$ and $\texttt{admissibleLab}$, respectively, of type '$a$ Set $\Rightarrow$ '$a$ Rel $\Rightarrow$ '$a$ Labelling $\Rightarrow o$):

$$\texttt{conflictfreeLab}^{\mathcal{A}} := \lambda att. \lambda \mathcal{L}ab. \forall^{\mathcal{A}} x.\ (\texttt{in}\ \mathcal{L}ab\ x \longrightarrow \neg\texttt{legallyOut}^{\mathcal{A}}\ att\ \mathcal{L}ab\ x)$$
$$\wedge\ (\texttt{out}\ \mathcal{L}ab\ x \longrightarrow \texttt{legallyOut}^{\mathcal{A}}\ att\ \mathcal{L}ab\ x)$$
$$\texttt{admissibleLab}^{\mathcal{A}} := \lambda att. \lambda \mathcal{L}ab. \forall^{\mathcal{A}} x.\ (\texttt{in}\ \mathcal{L}ab\ x \longrightarrow \texttt{legallyIn}^{\mathcal{A}}\ att\ \mathcal{L}ab\ x)$$
$$\wedge\ (\texttt{out}\ \mathcal{L}ab\ x \longrightarrow \texttt{legallyOut}^{\mathcal{A}}\ att\ \mathcal{L}ab\ x)$$

We have, in fact, employed Isabelle to verify automatically that admissible labellings always exist (e.g., consider a labelling where all arguments are `Undec`) and also that admissible labellings are indeed conflict-free.

Moreover, it has been proven automatically that, for admissible labellings, if an argument is *legally undec* then it is labelled `Undec`, but not the other way round (counter-models provided by Nitpick). Interestingly, one can also verify, again by generating counter-models with Nitpick, that for admissible labellings, a *legally in* (resp. *legally out*) argument is not generally labelled `In` (resp. `Out`). This situation changes, however, when we start considering complete labellings.

### Complete labellings

**Definition 3.16** A labelling $\mathcal{L}ab$ is termed *complete* if (i) it is admissible; and (ii) every `Undec`-labelled argument is *legally undec*.

**Formalisation.** The above definition is analogously encoded in HOL:

$$\texttt{completeLab}^{\mathcal{A}} := \ \lambda att. \ \lambda \mathcal{L}ab. \ (\texttt{admissibleLab}^{\mathcal{A}} \ att \ \mathcal{L}ab) \ \wedge$$
$$\forall^{\mathcal{A}}x. \ (\texttt{undec} \ \mathcal{L}ab \ x) \ \longrightarrow \ (\texttt{legallyUndec}^{\mathcal{A}} \ att \ \mathcal{L}ab \ x)$$

Using the Sledgehammer tool integrated into Isabelle/HOL it can be proven automatically that for complete labellings, *legally in* (resp. *legally out*) arguments are indeed labelled `In` (resp. `Out`). In fact, the following alternative characterisation for complete labellings has been verified as a theorem:

**Lemma 3.17**

$$\texttt{completeLab}^{\mathcal{A}} \ att \ \mathcal{L}ab = \forall^{\mathcal{A}}x. \ (\texttt{in} \ \mathcal{L}ab \ x \longleftrightarrow \texttt{legallyIn} \ att \ \mathcal{L}ab \ x)$$
$$\wedge \ (\texttt{out} \ \mathcal{L}ab \ x \longleftrightarrow \texttt{legallyOut} \ att \ \mathcal{L}ab \ x)$$

$\square$

Additionally, it is verified that for complete labellings, we have that *in/out-sets* completely determine the labelling, i.e., the following two statements hold:

**Lemma 3.18**

$(i) \ (\texttt{completeLab}^{\mathcal{A}} \ att \ L_1) \wedge (\texttt{completeLab}^{\mathcal{A}} \ att \ L_2) \longrightarrow$
$$\big((\texttt{in} \ L_1) \approx^{\mathcal{A}} (\texttt{in} \ L_2) \longrightarrow \forall^{\mathcal{A}}x. \ (L_1 \ x) = (L_2 \ x)\big)$$
$(ii) \ (\texttt{completeLab}^{\mathcal{A}} \ att \ L_1) \wedge (\texttt{completeLab}^{\mathcal{A}} \ att \ L_2) \longrightarrow$
$$\big((\texttt{out} \ L_1) \approx^{\mathcal{A}} (\texttt{out} \ L_2) \longrightarrow \forall^{\mathcal{A}}x. \ (L_1 \ x) = (L_2 \ x)\big).$$

$\square$

By generating counter-examples with Nitpick it is verified that, in contrast, *undec-sets* do not completely determine the (complete) labellings.

**Preferred and grounded labellings**

We now turn to the notions of minimality and maximality for complete labellings, drawing upon the definitions provided in §3.1. With these notions we can now discuss complete labellings where *in-sets* are maximal or minimal. They correspond to the so-called *preferred* and *grounded* labellings, respectively.

**Definition 3.19** A labelling $\mathcal{L}ab$ is termed *preferred* if it is a complete labelling whose *in-set* is maximal (wrt. set inclusion) among all the complete labellings.

**Definition 3.20** A labelling $\mathcal{L}ab$ is termed *grounded* if it is a (in fact: *the*) complete labelling whose *in-set* is minimal (wrt. set inclusion) among all the complete labellings.

**Formalisation.** The above definitions are encoded as HOL predicates in an analogous fashion:

$$\texttt{preferredLab}^{\mathcal{A}} := \lambda att.\ \lambda \mathcal{L}ab.\ \texttt{maximal}^{\mathcal{A}}\ (\texttt{completeLab}^{\mathcal{A}}\ att)\ \mathcal{L}ab\ \texttt{in}$$
$$\texttt{groundedLab}^{\mathcal{A}} := \lambda att.\ \lambda \mathcal{L}ab.\ \texttt{minimal}^{\mathcal{A}}\ (\texttt{completeLab}^{\mathcal{A}}\ att)\ \mathcal{L}ab\ \texttt{in}$$

Recalling the definitions of maximal (resp. minimal) in §3.1, they unfold into:

$$\texttt{preferredLab}^{\mathcal{A}}\ att\ \mathcal{L}ab\ =\ \texttt{completeLab}^{\mathcal{A}}\ att\ \mathcal{L}ab\ \wedge$$
$$\forall L.\ (\texttt{completeLab}^{\mathcal{A}}\ att\ L) \wedge (\texttt{in}\ \mathcal{L}ab) \subseteq^{\mathcal{A}} (\texttt{in}\ L) \longrightarrow (\texttt{in}\ L) \approx^{\mathcal{A}} (\texttt{in}\ \mathcal{L}ab)$$
$$\texttt{groundedLab}^{\mathcal{A}}\ att\ \mathcal{L}ab\ =\ \texttt{completeLab}^{\mathcal{A}}\ att\ \mathcal{L}ab\ \wedge$$
$$\forall L.\ (\texttt{completeLab}^{\mathcal{A}}\ att\ L) \wedge (\texttt{in}\ L) \subseteq^{\mathcal{A}} (\texttt{in}\ \mathcal{L}ab) \longrightarrow (\texttt{in}\ L) \approx^{\mathcal{A}} (\texttt{in}\ \mathcal{L}ab).$$

**Ideal labellings**

The notion of *ideal sets* and their maximal/greatest element (*ideal extension*) from §3.3 can analogously be lifted to labellings (cf. [27]). In order to do this, an ordering relation on labellings needs to be introduced first:

**Definition 3.21** Let $L_1$ and $L_2$ be two labellings. We say that $L_1$ is less or equally committed than $L_2$ if both the *in-set* resp. *out-set* of $L_1$ are contained in the *in-set* resp. *out-set* of $L_2$. We use the notation $L_1 \sqsubseteq L_2$.

We now employ the definition above to lift the corresponding definition for ideal extensions from §3.3: (i) *ideal sets* become *quasi-ideal labellings*, and (ii) *ideal extensions* (greatest ideal sets wrt. $\subseteq$) become *ideal labellings* (greatest quasi-ideal labellings wrt. $\sqsubseteq$). Let us now make this definition official:

**Definition 3.22** A labelling is termed *quasi-ideal* if it is admissible and is less or equally committed than every preferred labelling. The 'most committed' among all quasi-ideal labellings (i.e., greatest wrt. $\sqsubseteq$) is called the *ideal labelling*.

**Formalisation.** The above definitions are encoded, in an analogous manner, as HOL predicates (of type '$a$ `Set` $\Rightarrow$ '$a$ `Rel` $\Rightarrow$ '$a$ `Labelling` $\Rightarrow o$):

$$\sqsubseteq^{\mathcal{A}} := \quad \lambda L_1.\ \lambda L_2.\ (\text{in } L_1 \subseteq^{\mathcal{A}} \text{in } L_2)\ \wedge\ (\text{out } L_1 \subseteq^{\mathcal{A}} \text{out } L_2)$$

$$\texttt{qidealLab}^{\mathcal{A}} := \lambda att.\ \lambda \mathcal{L}ab.\ (\texttt{admissibleLab}^{\mathcal{A}}\ att\ \mathcal{L}ab)\ \wedge$$
$$\forall L.\ (\texttt{preferredLab}^{\mathcal{A}}\ att\ L) \longrightarrow \mathcal{L}ab \sqsubseteq^{\mathcal{A}} L$$

$$\texttt{idealLab}^{\mathcal{A}} := \lambda att.\ \lambda \mathcal{L}ab.\ (\texttt{qidealLab}^{\mathcal{A}}\ att\ \mathcal{L}ab)\ \wedge$$
$$\forall L.\ (\texttt{qidealLab}^{\mathcal{A}}\ att\ L) \longrightarrow L \sqsubseteq^{\mathcal{A}} \mathcal{L}ab$$

### Stable, semi-stable and stage labellings

We now turn to those labellings whose *undec-sets* must satisfy some particular minimality requirements. Observe that, for their *in-sets*, this works analogously to the corresponding conditions in §3.3 involving maximality of their *range*.

**Definition 3.23** A labelling $\mathcal{L}ab$ is termed *stable* if it is a complete labelling whose *undec-set* is empty, i.e., no argument is labelled `Undec`.

**Definition 3.24** A labelling $\mathcal{L}ab$ is termed *semi-stable* if it is a complete labelling whose *undec-set* is minimal (wrt. set inclusion) among all complete labellings.

**Definition 3.25** A labelling $\mathcal{L}ab$ is termed *stage* if it is a conflict-free labelling whose *undec-set* is minimal (wrt. set inclusion) among all conflict-free labellings.

**Formalisation.** The above definitions are encoded in HOL analogously:

$$\texttt{stableLab}^{\mathcal{A}} := \lambda att.\ \lambda \mathcal{L}ab.\ (\texttt{completeLab}^{\mathcal{A}}\ att\ \mathcal{L}ab) \wedge \forall x.\ (\mathcal{L}ab\ x) \neq \texttt{Undec}$$

$$\texttt{semistableLab}^{\mathcal{A}} := \lambda att.\ \lambda \mathcal{L}ab.\ \texttt{minimal}^{\mathcal{A}}\ (\texttt{completeLab}^{\mathcal{A}}\ att)\ \mathcal{L}ab\ \texttt{undec}$$

$$\texttt{stageLab}^{\mathcal{A}} := \lambda att.\ \lambda \mathcal{L}ab.\ \texttt{minimal}^{\mathcal{A}}\ (\texttt{conflictfreeLab}^{\mathcal{A}}\ att)\ \mathcal{L}ab\ \texttt{undec}$$

## 4 Assessment of Meta-Theoretical Properties

In this section, meta-theoretic properties of the presented encoding are analysed. We proceed by first briefly exemplifying the concept of interactive theory exploration in §4.1 as a powerful application of computer-assisted (meta-theoretical) reasoning. Here, the utilisation of proof assistants for exploring and synthesising meta-theoretical concepts and insights in a dialogue-like setting is discussed. Secondly, in §4.2, the formal verification of the argumentation semantics' relationships as well as selected further fundamental properties is presented. By doing so, we implicitly present adequacy claims of the encoding of abstract argumentation into HOL.

### 4.1 Interactive Theory Exploration

Consider the following lemma, called *Fundamental lemma* by Dung [37]:

**Lemma 4.1 (Fundamental Lemma)** *Let $AF = (\mathcal{A}, \rightsquigarrow)$ and $S \subseteq \mathcal{A}$ be an admissible set of arguments. For any $a \in \mathcal{A}$ it holds that if $S$ defends $a$, then $S \cup \{a\}$ is admissible.* [9]                                                    □

Suppose that we want to formulate a corresponding fundamental lemma for labelling-based semantics. It is well-known that a labelling $\mathcal{L}ab$ can be translated into a corresponding extension by taking its set of In-labelled arguments $in(\mathcal{L}ab)$. Following this intuition, a naive (and, indeed, wrong) adaption of Lemma 4.1 for labellings could amount to the following statement:

Let $AF = (\mathcal{A}, \rightsquigarrow)$ *and $\mathcal{L}ab$ be an admissible labelling for $AF$. For any $a \in \mathcal{A}$ it holds that if $in(\mathcal{L}ab)$ defends $a$, then $\mathcal{L}ab'$ is an admissible labelling for $AF$, where $\mathcal{L}ab'(a) = $* In *and $\mathcal{L}ab'(x) = \mathcal{L}ab(x)$ for every $x \neq a$.*

Intuitively, in this approach the labelling $\mathcal{L}ab$ is extended to a new labelling $\mathcal{L}ab'$ that is identical to $\mathcal{L}ab$ except that $a$ is now labelled In by $\mathcal{L}ab'$. Based on the encoding from §3, this is formalised in HOL as follows:

$$\texttt{admissibleLab}^{\mathcal{A}} \, att \, \mathcal{L}ab \wedge \texttt{defends} \, att \, (\texttt{in} \, \mathcal{L}ab) \, a \longrightarrow$$
$$\texttt{admissibleLab}^{\mathcal{A}} \, att \, (\lambda x.\texttt{if} \ x = a \ \texttt{then In else} \ (\mathcal{L}ab \, x))$$

The new labelling $\mathcal{L}ab'$ is thereby given by the anonymous function

$$\big(\lambda x.\texttt{if} \ x = a \ \texttt{then In else} \ (\mathcal{L}ab \, x)\big)$$

that corresponds to the following function in a more conventional mathematical notation:

$$\mathcal{L}ab' : x \mapsto \begin{cases} \texttt{In} & \text{if } x = a, \\ \mathcal{L}ab(x) & \text{otherwise} \end{cases}$$

Note that the if-then-else statement is merely a syntactic abbreviation (so-called *syntactic sugar*) and can adequately be represented itself in HOL. But since Isabelle/HOL supports this elegant and concise representation as well, we adopt it in the following.

When attempting to prove this statement in Isabelle/HOL using the built-in automated reasoning tools, an invocation of the counter-model generator **nitpick** quickly (in less than one second) yields a counter-example to the proposed lemma. The formulation in Isabelle/HOL and the original output of **nitpick** is displayed in Fig. 2a. The false lemma, named `NaiveDungFundLemma1Lab`, is displayed in the upper part of the window; the lower half contains the output of the counter-model generator.

The counter-example autonomously found by **nitpick**, cf. Fig. 2a, specifies an argumentation framework $AF = (\mathcal{A}, \rightsquigarrow)$ with $\mathcal{A} = \{A, B, C\}$, $\rightsquigarrow = \{(A, B), (B, C)\}$, and an offending initial labelling $\mathcal{L}ab$ given by $(\{A\}, \emptyset, \{B, C\})$. While this indeed provides a feasible counter-example, it

---

[9] Note that in the original formulation of Dung [37] the notion of *defence* was rather referred to as "acceptability".

```
36 lemma NaiveDungFundLemma1Lab: ‹admissibleLabᴬ att Lab ∧ defendsᴬ att (in Lab) a ⟶
37                                  admissibleLabᴬ att ( λx. if (x = a) then In else (Lab x) )›
38   nitpick
39
```

| ✓ Proof state  ✓ Auto update   Update   Search: _____ ▾  100% |

```
Nitpick found a counterexample for card 'a = 3:
  Free variables:
    Lab = (λx. _)(A := In, B := Undec, C := Undec)
    𝒜 = (λx. _)(A := True, B := True, C := True)
    a = C
    att =
      (λx. _)
      ((A, A) := False, (A, B) := True, (A, C) := False, (B, A) := False, (B, B) := False,
       (B, C) := True, (C, A) := False, (C, B) := False, (C, C) := False)
```

(a) Formalisation of the incorrect fundamental lemma adaption in Isabelle/HOL. The (counter-)model generator **nitpick** is invoked by writing its name after the lemma statement. The generated counter-example is printed in the lower part of the window; it encodes a concrete argumentation framework and a labelling as human readable (but technical) output. Free variables of the statement (in dark blue, left-hand side) are assigned to concrete interpretations (in black, right-hand side), e.g., the free variable A is assigned to a predicate that maps the synthesised arguments A, B and C to true, i.e., representing the set {A, B, C} of arguments (by its characteristic function). Similarly, Lab is assigned to a function that maps A to In, and B and C to Undec, representing a concrete labelling function $\mathcal{L}ab$.

```
68 lemma NaiveDungFundLemma1Lab: ‹admissibleLabᴬ att Lab ∧ defendsᴬ att (in Lab) a ⟶
69                                  admissibleLabᴬ att ( λx. if (x = a) then In else (Lab x) )›
70   nitpick [ eval = "legallyInᴬ att Lab a"]
```

| ✓ Proof state  ✓ Auto update   Update   Search: _____ ▾  100% |

```
Evaluated term:
  legallyInᴬ att Lab a = False
```
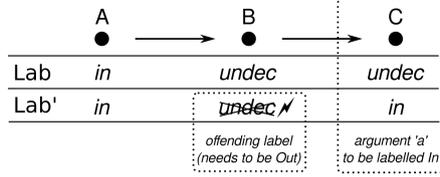
(b) An extended query to **nitpick** requesting to evaluate ("eval") a concrete expression within the provided counter-model. Here, we ask whether the argument assigned to the free variable a is *legally in*; the result in the lower part of the windows yields the result that this is not the case (the expression is assigned to the value False representing falsehood).

Fig. 2. Interactive assessment of an exemplary meta-theoretical statement on admissible labellings.

may not be completely apparent why this is the case. This is why **nitpick** will also output the offending assignment of any relevant free variable contained in the statement, here the free variable a: In this counter-example, a is assigned to $C \in \mathcal{A}$; indicating that the statement is not valid when assuming the above framework $AF$, (admissible) labelling $\mathcal{L}ab$ and considering to extend $\mathcal{L}ab$ by additionally labelling $C$ with In. If it is still not clear why labelling $C$ with In would result in a labelling that is not admissible, we can ask **nitpick** to evaluate further statements in the context of the found counter-model, e.g., whether $C$ would be *legally in* in the new labelling $\mathcal{L}ab'$. This is visualised in Fig. 2b. Of course, any statement can be assessed in this fashion if further information are required by the user, e.g., in more complex scenarios.

The counter-example, together with the other information provided by the proof assistant, is summarised in Fig 3a. Of course, argument $C$ cannot be

(a) Visualisation of the synthesised counter-example from Fig. 2a. Argument $C$ is chosen as the argument to be accepted by the updated labelling, i.e, assigned to the free variable $a$ in the lemma formulation. The updated labelling $\mathcal{L}ab'$ assigns In to argument $C$ which, in turn, requires argument $B$ to be labelled Out.



(b) Encoding and verification of the updated adaption of the fundamental lemma for admissible labellings. The tool **sledgehammer** automatically invokes external ATP systems and suggests proofs that are subsequently formally verified by Isabelle/HOL. Consequently, the one-line proof at line 91 ("**by** ...") is a fully computer-verified proof of Lemma 4.2.

Fig. 3. Analysis of the computer-generated counter-example and formal proof of the updated fundamental lemma for labellings.

labelled *legally in* since, by definition, any attacker of it would need to be labelled out (but is, in fact, labelled Undec in the counter-example). Using this insights, we may propose an updated adapted fundamental lemma for labellings:

**Lemma 4.2 (Fundamental Lemma for labellings)** *Let $AF = (\mathcal{A}, \rightsquigarrow)$ and $\mathcal{L}ab$ be an admissible labelling for $AF$. For any $a \in \mathcal{A}$ it holds that if $in(\mathcal{L}ab)$ defends $a$, then $\mathcal{L}ab'$ is an admissible labelling for $AF$, where*

  (i) $\mathcal{L}ab'(a) = \texttt{In}$,

 (ii) $\mathcal{L}ab'(x) = \texttt{Out}$ *if* $x \rightsquigarrow a$, *and*

(iii) $\mathcal{L}ab'(x) = \mathcal{L}ab(x)$ *otherwise.*

$\square$

This lemma is translated into HOL syntax as follows:

$$\texttt{admissibleLab}^{\mathcal{A}}\ att\ \mathcal{L}ab \wedge \texttt{defends}\ att\ (\texttt{in}\ \mathcal{L}ab)\ a \longrightarrow$$
$$\texttt{admissibleLab}^{\mathcal{A}}\ att\ \big(\lambda x.\texttt{if}\ x = a\ \texttt{then In else}$$
$$(\texttt{if}\ att\ x\ a\ \texttt{then Out else}\ (\mathcal{L}ab\ x))\big)$$

The updated labelling function $\mathcal{L}ab'$ is again represented by an anonymous

function $(\lambda x. \ldots)$ that contains a nested `if-then-else` statement following the definition of $\mathcal{L}ab'$ from Lemma 4.2.

The verification of Lemma 4.2 is displayed at Fig. 3b, where its contents are formalised as lemma `DungFundLemma1Lab` in Isabelle/HOL. As depicted in Fig. 3b, the tool **sledgehammer** may be used to attempt automated proof search. In this case a proof is found after a few seconds, provided by the external reasoning system CVC4 [7], an SMT solver [8], which is then automatically reconstructed by and verified in Isabelle/HOL. Hence, we have found a new meta-theoretical property – a labelling-based fundamental lemma – using computer-assisted reasoning in an interactive way.

Of course, such an approach may as well be applied to more complex statements in the context of interactive and semi-automated theory exploration. Note that the encoding of abstract argumentation into HOL provides a uniform framework for both reasoning *about* abstract argumentation (i.e., meta-theory) as well as *within* abstract argumentation (i.e., computation of extensions and labellings) using the same tools. The former case is illustrated in more detail next, the latter case is addressed in §5.

## 4.2 Adequacy of the Formalisation

One of the key advantages of the presented encoding of abstract argumentation is that the surrounding logical framework provided by the HOL formalism allows for the computer-assisted verification of the encoding's adequacy – while providing computational means for generating extensions/labellings of argumentation frameworks at the same time. Adequacy here means that the semantics of the encoded structures coincide with the intended ones, e.g., that any extension or labelling produced via the encoding is sound with respect to the given argumentation semantics. For this purpose, we establish the central properties of argumentation frameworks and their semantics for the presented encoding. The properties are thereby either taken from Dung's original work [37] or the survey by Baroni et al. [3]. Each of these properties have been formally verified within Isabelle/HOL – usually proven automatically by external ATP systems in a few seconds using **sledgehammer**. This is a strong practical argument for the feasibility of the presented approach and proof assistants in general. We present the following properties in their formally encoded variant within HOL, exactly as given to the proof assistant; free variables are implicitly universally quantified. We start by providing general properties of admissible and conflict-free sets represented by the encoding:

There are no self-attacking arguments in conflict-free extensions [3, p. 8]:

**Lemma 4.3** `conflictfreeExt`$^{\mathcal{A}}$ $att$ $E \longrightarrow \neg\big(\exists^{\mathcal{A}}a.\,(E\,a) \wedge a \rightsquigarrow a\big).$      □

The characteristic function preserves conflict-freeness (monotonicity was already shown in §3.2):

**Lemma 4.4**
`conflictfreeExt`$^{\mathcal{A}}$ $att$ $E \longrightarrow$ `conflictfreeExt`$^{\mathcal{A}}$ $att$ $(\mathcal{F}^{\mathcal{A}}\,att\,E).$      □

A conflict-free set $E$ is admissible iff $E \subseteq \mathcal{F}^{\mathcal{A}}(E)$ [37, Lemma 18]:

**Lemma 4.5**
$\texttt{conflictfreeExt}^{\mathcal{A}} \; att \; E \longrightarrow \left( \texttt{admissibleExt}^{\mathcal{A}} \; att \; E \longleftrightarrow E \subseteq^{\mathcal{A}} \mathcal{F}^{\mathcal{A}} \; att \; E \right)$.
$\square$

Admissible sets can be extended with defended arguments (Dung's fundamental lemma) [37, Lemma 10]:

**Lemma 4.6**

(i) $\left( \texttt{admissibleExt}^{\mathcal{A}} \; att \; E \wedge \texttt{defends} \; att \; E \; a \right)$
$\qquad\qquad\qquad\qquad \longrightarrow \texttt{admissibleExt}^{\mathcal{A}} \; att \; (E \cup a)$

(ii) $\left( \texttt{admissibleExt}^{\mathcal{A}} \; att \; E \wedge \texttt{defends} \; att \; E \; a \wedge \texttt{defends} \; att \; E \; a' \right)$
$\qquad\qquad\qquad\qquad \longrightarrow \texttt{defends}^{\mathcal{A}} \; att \; (E \cup a) \; a' \qquad \square$

Admissible sets form a $(\omega-)$complete partial order (CPO) with respect to set inclusion [37, Theorem 11]. In fact, a stronger statement could be formally verified: The collection of admissible sets form a directed CPO (dCPO):

**Lemma 4.7**

(i) $\omega\texttt{-cpo}^{\mathcal{A}} \; (\texttt{admissibleExt}^{\mathcal{A}} \; att)$

(ii) $\texttt{dcpo}^{\mathcal{A}} \; (\texttt{admissibleExt}^{\mathcal{A}} \; att)$

*where* $\omega\texttt{-cpo}^{\mathcal{A}}$ *and* $\texttt{dcpo}^{\mathcal{A}}$ *encode the notions of (directed) complete partial orders in HOL [44, misc].* $\square$

From this, it can be verified that for each admissible set $S$ there exists a preferred extensions extending $S$ [37, Theorem 11]:

**Lemma 4.8**
$\texttt{admissibleExt}^{\mathcal{A}} \; att \; S \longrightarrow \exists E. \, S \subseteq^{\mathcal{A}} E \wedge \texttt{preferredExt}^{\mathcal{A}} \; att \; E \qquad \square$

We proceed by highlighting a few properties of selected classes of extensions resp. labellings. Conflict-free sets are complete iff they are fixed-points of $\mathcal{F}$ [37, Lemma 24]:

**Lemma 4.9**
$\texttt{conflictfreeExt}^{\mathcal{A}} \; att \; S \longrightarrow \left( \texttt{completeExt}^{\mathcal{A}} \; att \; S \longleftrightarrow S \approx^{\mathcal{A}} \mathcal{F}^{\mathcal{A}} \; att \; S \right).$ $\square$

Complete, preferred and grounded extensions always exist [37]:

**Lemma 4.10**

(i) $\exists E. \, \texttt{completeExt}^{\mathcal{A}} \; att \; E$

(ii) $\exists E. \, \texttt{preferredExt}^{\mathcal{A}} \; att \; E$

(iii) $\exists E. \, \texttt{groundedExt}^{\mathcal{A}} \; att \; E$ $\qquad\qquad \square$

Grounded labellings can equivalently be characterised by minimal in-sets and minimal out-sets [3, Prop. 5]:

**Lemma 4.11**
$\texttt{minimal}^{\mathcal{A}} \; (\texttt{complete}^{\mathcal{A}} \; att) \; \mathcal{L}ab \; \texttt{in} \longleftrightarrow \texttt{minimal}^{\mathcal{A}} \; (\texttt{complete}^{\mathcal{A}} \; att) \; \mathcal{L}ab \; \texttt{out}$
$\square$

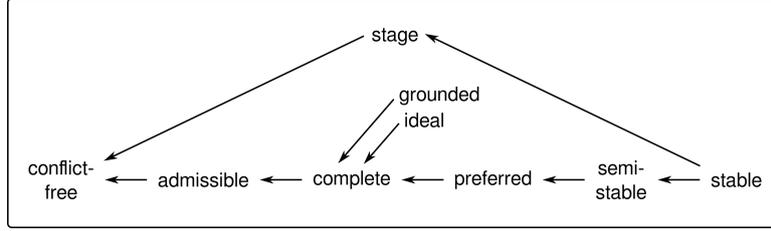Furthermore, grounded extensions resp. labellings are unique [3, Prop. 4]:

Fig. 4. Relationship between each of the different argumentation semantics considered in this article. Each entry represents both the respective extension-based and labelling-based semantics. An arrow from one semantics to another can be read as *"is a"*, it symbolises a generalisation relation, i.e., the latter is more general than the former. Transitive arrows are omitted.

**Lemma 4.12**

(i) $\texttt{groundedExt}^{\mathcal{A}}\ att\ S \longleftrightarrow \texttt{least}^{\mathcal{A}}\ (\texttt{completeExt}^{\mathcal{A}}\ att)\ S\ id$

(ii) $\texttt{groundedLab}^{\mathcal{A}}\ att\ \mathcal{L}ab \longleftrightarrow \texttt{least}^{\mathcal{A}}\ (\texttt{completeLab}^{\mathcal{A}}\ att)\ \mathcal{L}ab\ \texttt{in}$          □

Analogously, preferred labellings can equivalently be characterised by maximal in-sets and minimal out-sets [3, Prop. 8]:

**Lemma 4.13**
$\texttt{maximal}^{\mathcal{A}}\ (\texttt{complete}^{\mathcal{A}}\ att)\ \mathcal{L}ab\ \texttt{in} \longleftrightarrow \texttt{maximal}^{\mathcal{A}}\ (\texttt{complete}^{\mathcal{A}}\ att)\ \mathcal{L}ab\ \texttt{out}$

□

Additionally to the properties of the individual argumentation semantics, there exist well-known relationships between them: As an example, every complete labelling is also an admissible labelling; and every stable extension is also a conflict-free one [3]. Fig. 4 shows these relationships for a subset of argumentation semantics considered in this article. As another case study for the assessment of meta-theoretical properties of argumentation semantics using proof assistants, we have verified all of the displayed links in Isabelle/HOL [44].

Even more, the non-inclusion of the inverse directions of the arrows from Fig. 4 has been established, for each case, using the counter-model generator **nitpick**. Here, small counter-examples are generated automatically that illustrate the non-validity of the inverse statement. As an example, **nitpick** is able to refute the statement that every stage labelling is also a stable labelling, encoded in HOL as

$$\texttt{stage}^{\mathcal{A}}\ att\ \mathcal{L}ab \longrightarrow \texttt{stable}^{\mathcal{A}}\ att\ \mathcal{L}ab,$$

in less than one second. It is easy to check that the labelling returned by **nitpick** is indeed a stage labelling but not a stable one.

Furthermore we formally verified the different correspondence results between extensions and labellings [44] for the semantics listed below.

$$\mathcal{S}em = \{\texttt{conflictfree}, \texttt{admissible}, \texttt{complete}, \texttt{grounded},$$
$$\texttt{preferred}, \texttt{stable}, \texttt{semistable}, \texttt{stage}\}$$

For this sake we have formalised the well known translation mappings between extensions and labellings [3] as the HOL functions: `Lab2Ext` and `Ext2Lab` of types '$a$ `Labelling` $\Rightarrow$ '$a$ `Set` and '$a$ `Set` $\Rightarrow$ '$a$ `Rel` $\Rightarrow$ '$a$ `Set` $\Rightarrow$ '$a$ `Labelling` respectively. Observe that the function `Lab2Ext`, in contrast to `Ext2Lab`, can be defined independently of the underlying argumentation framework.

$$\begin{aligned} \texttt{Lab2Ext} \;&:=\; \lambda \mathcal{L}ab.\; \texttt{in } \mathcal{L}ab \\ \texttt{Ext2Lab}^{\mathcal{A}} \;&:=\; \lambda att.\; \lambda E.\; \lambda a.\; \texttt{if } (E\ a) \texttt{ then In else} \\ &\qquad\qquad \texttt{if } ([\mathcal{A}|att|E]^{+}\ a) \texttt{ then Out else Undec} \end{aligned}$$

As an example, the following statements relate preferred extensions and preferred labellings (analogous results hold for the other semantics listed above):

**Lemma 4.14**

(i) $\texttt{preferredLab}^{\mathcal{A}}\ att\ \mathcal{L}ab \longrightarrow \texttt{preferredExt}^{\mathcal{A}}\ att\ (\texttt{Lab2Ext}\ \mathcal{L}ab)$

(ii) $\texttt{preferredExt}^{\mathcal{A}}\ att\ S \longrightarrow \texttt{preferredLab}^{\mathcal{A}}\ att\ (\texttt{Ext2Lab}\ S)$.

(iii) $\texttt{preferredLab}^{\mathcal{A}}\ att\ (\texttt{Ext2Lab}\ S) \longrightarrow \texttt{preferredExt}^{\mathcal{A}}\ att\ S$. $\qquad\qquad\square$

It is worth noting that the model generator *nitpick* has, in fact, found counter-models to the converse implication to item (i) above, as expected [3].

In summary, the presented HOL encoding provides not only the definition of abstract arguments and its different semantics, but also provides a formal and computer-assisted verification of different meta-theoretical properties including relationships between argumentation semantics, correspondences, and equivalent alternative characterisations.

# 5 Flexible Generation of Extensions and Labellings

The encoding of the different argumentation semantics presented above captures the structure and the logical behaviour of abstract argumentation frameworks within HOL. Building on top of that, we can make use of *model finders*, e.g., Nitpick [23] and Nunchaku [51], for generating concrete extensions and labellings for a given argumentation frameworks. To this end we here employ the model finder Nitpick that is readily integrated into the Isabelle/HOL proof assistant. [10]

## 5.1 Generating Standard Extensions and Labellings

Figure 5 displays a few representative examples of argumentation frameworks, taken from [3], that serve as use cases to illustrate the generation of extensions and labellings employing the model finder Nitpick. We will discuss only a few results. The rest can be consulted in the corresponding Isabelle/HOL sources that can be found in [44, `model-generation`].

---

[10] Of course, other automated theorem provers and model generators for higher-order logics, as provided in other proof assistants (e.g. HOL, Coq, Lean) can also be employed for our purposes. To the best of our knowledge, the level of proof automation featured in the Isabelle/HOL ecosystem is currently unmatched.
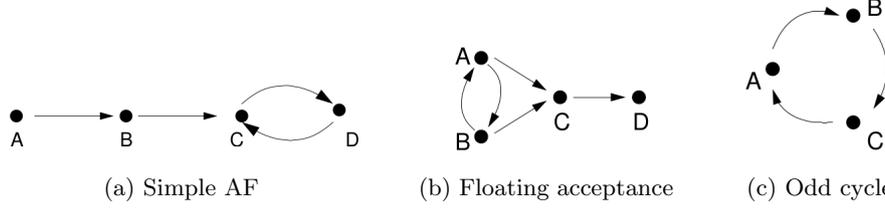
(a) Simple AF         (b) Floating acceptance        (c) Odd cycle

Fig. 5. Some representative examples of argumentation frameworks by Baroni, Cam-inada and Giacomin [3, Fig. 4-6].

To encode the above argumentation graphs in Isabelle/HOL we employ the simplified encoding approach as discussed in §3. In this approach the universe of arguments is implicitly given as the carrier of the *attack* relation. Thus, for the type of arguments we define a new datatype: `Arg`, consisting only of the distinct terms `A`, `B`, `C` (and `D` when required). Next, we encode the attack relation `att` as binary predicate such that `att` $X$ $Y$ if and only if $X$ attacks $Y$ according to the corresponding graph in Fig. 5. This is displayed in the corresponding Isabelle/HOL setup in Figure 6.

```
datatype Arg = A | B | C | D          datatype Arg = A | B | C | D          datatype Arg = A | B | C
fun att :: ‹Arg Rel› where              fun att :: ‹Arg Rel› where            fun att::‹Arg Rel› where
  "att A B = True" |                      "att A B = True" |                    "att A B = True" |
  "att B C = True" |                      "att B A = True" |                    "att B C = True" |
  "att C D = True" |                      "att A C = True" |                    "att C A = True" |
  "att D C = True" |                      "att B C = True" |                    "att _ _  = False"
  "att _ _  = False"                      "att C D = True" |
                                          "att _ _  = False"
```

(a) Simple AF      (b) Floating acceptance      (c) Odd cycle
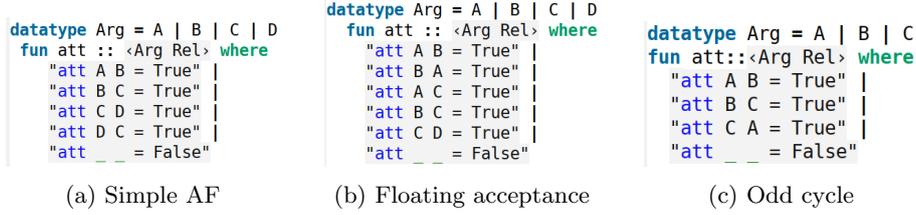
Fig. 6. Encoding the argumentation frameworks in Fig. 5.

We can now ask the model finder Nitpick to generate, say, all preferred labellings for the *AF* in Fig. 5a. Indeed, Nitpick produces the following two labellings (cf. the original output of Nitpick displayed in Fig. 7):

$$\mathcal{L}ab_1 : x \mapsto \begin{cases} \text{In} & \text{if } x = A \\ \text{Out} & \text{if } x = B \\ \text{In} & \text{if } x = C \\ \text{Out} & \text{if } x = D \end{cases} \quad \text{and} \quad \mathcal{L}ab_2 : x \mapsto \begin{cases} \text{In} & \text{if } x = A \\ \text{Out} & \text{if } x = B \\ \text{Out} & \text{if } x = C \\ \text{In} & \text{if } x = D \end{cases}$$

which represent the labelling $\mathcal{L}ab_1$ and $\mathcal{L}ab_2$ such that $in(\mathcal{L}ab_1) = \{A, C\}$, $out(\mathcal{L}ab_1) = \{B, D\}$ and $undec(\mathcal{L}ab_1) = \emptyset$, and $in(\mathcal{L}ab_2) = \{A, D\}$, $out(\mathcal{L}ab_2) = \{B, C\}$ and $undec(\mathcal{L}ab_2) = \emptyset$.

In the example above, we employed a specially engineered function `findFor`, given by

$$\texttt{findFor} := \lambda AF. \, \lambda Prop. \, \lambda S. \, \forall Lab. \, (S \, Lab) \longleftrightarrow (Prop \, AF \, Lab),$$

in such a way that Nitpick tries to satisfy the statement

$$\texttt{findFor} \; att \; \texttt{preferredLab} \; Labs$$

```
65 (* This gives us exactly the 2 sets predicted in [BCG2011] p. 12: {A, C}, {A, D} *)
66 lemma ‹findFor' att preferredExt Exts› nitpick[satisfy, eval = "card Exts"]
67 (* This gives us exactly the two labellings predicted in [BCG2011] p. 12
68  {(A := In, B := Out, C := In,  D := Out),
69   (A := In, B := Out, C := Out, D := In)} *)
70 lemma ‹findFor' att preferredLab Labs› nitpick[satisfy,box=false, eval = "card Labs"]
71
```

☑ Proof state  ☑ Auto update   Update   Search: [                    ] ▼   100%

```
Nitpick found a model:
  Free variable:
    Labs =
      {(λx. _)(A := In, B := Out, C := In, D := Out), (λx. _)
       (A := In, B := Out, C := Out, D := In)}
  Evaluated term:
    card Labs = 2
```

Fig. 7. Nitpick output enumerating all preferred labellings of the argumentation framework from Fig. 5a.

```
30 (* Admissible labelling where A is In *)
31 lemma ‹admissibleLab att Lab ∧ (Lab a) = In› nitpick[satisfy]
32
33 (* Admissible labelling where Lab is surjective *)
34 lemma ‹admissibleLab att Lab ∧ (surjective Lab)› nitpick[satisfy]
35
36 (* Admissible labelling where there are more than two arguments labelled In *)
37 lemma ‹admissibleLab att Lab ∧ (card({x. in(Lab) x}) > 2)› nitpick[satisfy]
```

Fig. 8. Asking for a specific labelling that satisfies additional properties in Isabelle/HOL.

by generating a model (and, hence, enumerating all the labellings). Nitpick provides all preferred labellings by finding the value given to the free variable *Labs* above. This is equivalent to finding the value of *Labs* such that

$$\forall Lab. \; (Labs \; Lab) \; \longleftrightarrow \; (\texttt{preferredLab} \; att \; Lab)$$

holds, where `preferredLab` is the predicate as defined in §3.4. We observe that the reported results are in fact as described in [3]. The same holds for the remaining argumentation semantics and examples from Fig. 5.

### 5.2 Generating Flexibly-Constrained Extensions and Labellings

In addition to the above – quite standard – applications, we can now make use of the expressive surrounding logical framework to ask for *specific* labellings, flexibly constrained by means of an arbitrary (higher-order) predicate. Consider the example displayed in Fig. 8 extending the example from Fig. 5b:

In the three lemmas, we ask Nitpick to generate admissible labellings where, additionally, (1) argument *A* is labelled `In`, (2) *Lab* is a surjective function, and (3) there are more than two arguments labelled `In`, respectively. In the first two cases, suitable labellings are provided, in the third case no such labelling can be found (visualised by the red background colour indicating an error). Indeed, no such labelling exists.

Similarly, we can prove in Isabelle/HOL that for Fig. 5c no admissible

labelling other than the trivial one exists. This is expressed by the formula

$$(\texttt{admissibleLab}\ att\ \mathcal{L}ab) \longrightarrow \forall x.\,(\mathcal{L}ab\ x) = \texttt{Undec}$$

which is proven automatically by Sledgehammer within a few seconds. If not interested in a specific labelling or extension, it is also possible to merely prove the (non-)existence of one using automated provers for HOL. Additionally notions of skeptical and credulous argument justification [3, Defs. 56 and 57] have been encoded in HOL [44, base], although this has not been a subject of focus here.

## 6   Conclusion

In this article, an encoding of abstract argumentation frameworks as well as various argumentation semantics into higher-order logic (HOL) was presented. To that end, sets are identified with their characteristic function and represented by typed predicates within the HOL formalism. Similarly, the attack relation of argumentation frameworks is encoded as a binary predicate of appropriate type. Finally, argumentation semantics are represented by higher-order predicates on extensions (labellings) that hold, by construction, if and only if the given extension (labelling) indeed satisfies the constraints imposed by the respective argumentation semantics.

The presented encoding was exemplarily implemented within the well-known proof assistant Isabelle/HOL, enabling the employment of various interactive and automated deduction tools including, in particular, the model finder Nitpick and the automated meta theorem prover sledgehammer. The resulting source files of the encoding are readily and freely available at GitHub [44] for further usage. It is important to note that the encoding presented in this article is not fixed to any one specific reasoning system; we merely chose to use Isabelle/HOL for demonstration purposes because of its flexibility and user-friendliness, including, e.g., a graphical user interface for interactive experimentation and its portfolio of integrated automated reasoning tools.

Due to the expressiveness of the higher-order formalism, the encoding of abstract argumentation in HOL allows for both meta-level reasoning (i.e., reasoning *about* notions of abstract argumentation) as well as object-level reasoning (i.e., reasoning *with* argumentation networks) using the same portfolio of first-order and higher-order automated reasoning tools. Both aspects were highlighted in the article in the context of Isabelle/HOL applications: Meta-level reasoning was exemplified, firstly, by utilising Isabelle/HOL for interactively exploring the meta-theory of abstract argumentation; secondly, the adequacy of the formalisation itself was verified by formally encoding and proving the central properties, relationships and correspondences from abstract argumentation literature, while obtaining counter-models for well-known non-theorems. Subsequently, we demonstrated how to use the encoding for object-level reasoning, i.e., for generating extensions and labellings for given argumentation frameworks. Here, quite specific extensions and labellings can be generated that should additionally satisfy arbitrarily complex (higher-order) properties.

Since the computation is based on a computer-verified encoding, we automatically know that the results of the concrete outputs are correct as well. Up to the author's knowledge, there does not exist any other approach capable of all the above aspects. We hence argue that an encoding of argumentation into HOL provides a uniform framework for assessing abstract argumentation from the perspective of automated reasoning.

It has to be pointed out that the presented approach is not meant to provide an alternative to well-established means for efficiently computing extensions (or labellings) for large-scale argumentation frameworks. Of course, special-purpose procedures or SAT-based approaches do not make use of formalisms as expressive as HOL, and hence admit decidable and, in some sense, efficient routines. In contrast, HOL automated theorem provers are semi-decidable only. Nevertheless our aim is quite orthogonal and rather aims at providing generic means for interactively and (partly) automatically assess abstract argumentation within a rich ecosystem of reasoning tools. We hence provide a bridge between the landscape of abstract argumentation one the one side and automated deduction on the other. In fact, this is in line with the motivation put forward by the LogiKEy framework [14] that employs generic higher-order reasoning for assessing normative theories for ethical and legal reasoning. As a side contribution of this work we thus extend the LogiKEy framework with generic means of abstract argumentation, e.g., allowing experiments in legal argumentation [13] to be based on more principled notions of argumentation.

**Related work.** Besides the well-known reduction to logic programs [56], the encoding of constraints enforced by argumentation semantics into other formalisms has become a standard technique for implementing abstract argumentation [32]. Early works on logical encodings into propositional logic, as proposed by Dunne and Bench-Capon [39], and Besnard and Doutre [18], reduce the problem of finding admissible-based extensions as a logical satisfiability problem. This work has paved the way for later work on harnessing SAT-solvers [20] for this task [59,32]. This technique form the basis of various tools for argumentation such as, e.g., Cegartix [40], LabSATSolver [11] and jArgSem-SAT [33]. These approaches mostly focus on generating adequate extensions (labellings) and do not allow for introspection, i.e., support for meta-theoretical reasoning about abstract argumentation.

In a similar vein, other approaches rely on encoding abstract argumentation in more expressive logical formalisms than propositional logic. They make use of this increased expressivity for capturing many different argumentation semantics under a purely logical umbrella. An early approach towards the encoding of abstract argumentation in quantified propositional logic (QBF) has been proposed by Egly and Woltran [41]. An extended QBF-based approach has been introduced by Arieli and Caminada [2] to represent in an uniform way a wide range of extension-based semantics. Their approach allows for automated verification of some semantical properties, e.g., the existence of stable extensions and some inclusion and equivalence relations between extensions. Several (restricted) first-order logical formalisms have also

been proposed. Dupin de Saint-Cyr et al. [36] introduce a first-order language (YALLA) for the encoding of argumentation semantics and study dynamic aspects of abstract argumentation. More recently, Cayrol and Lagasquie-Schiex have proposed a first-order logical encoding of extended abstract argumentation frameworks featuring higher-order attacks and support relations [31]. We refer the reader to Gabbay [45] for further survey and discussion about logical encodings, including modal and second-order logic-based approaches. We note that, while those approaches are much in the same spirit as the one presented in this article, they are less expressive and generic, since many of the meta-theoretic analyses presented in §4 cannot be carried out within them as they make essential use of higher-order constructs. This sort of expressivity limitations for existent approaches has been, of course, a conscious design choice, given the well-known expressivity vs. scalability trade-off. In this respect, our HOL-based approach complements rather than competes with them; and it has the added value of enabling the utilisation in interactive proof assistants.

The presented encoding is closely related to the so-called *shallow semantical embeddings* [12]. Such embeddings allow for the representation of domain-specific expressions by merely considering the defined concepts as abbreviations of a host meta-language (in our case, of HOL) that can be unfolded exhaustively, yielding an ordinary (but complex) formula in the meta-logic. Such shallow embeddings were already studied for encoding non-classical logics into HOL, e.g., modal logics [15] and many-valued logics [54]. Via such embeddings, any higher-order reasoning system can be turned into a reasoner for the respective non-classical logic [12,46]. This approach has previously been utilised for encoding networks of structured arguments in [42] and also in [13] in the context of legal reasoning.

**Further work.** We plan to further extend the collection of encoded argumentation semantics towards other recent proposals in the literature (e.g., eager, CF2, and stage2 semantics) as well as to conduct extended meta-theoretical studies based on them.

Moreover, the expressivity of HOL also allows us to extend the scope of our work towards other extensions of abstract argumentation frameworks beyond Dung's approach, including joint and higher-order attacks [24,9,4], as well as bipolar argumentation that adds support relations between arguments [30,1].

Given the trade-off between expressivity and scalability, we are currently exploring the limits of our approach for larger inputs. This analysis is quite non-trivial, as it involves substantial engineering work regarding the effective orchestration of the different automated tools in the portfolio (automated theorem provers, SAT/SMT-solvers, model generators, etc.) for this particular family of applications.

Our approach also allows, in fact, for the instantiation of abstract arguments by complex objects, such as sets or tuples of formulas in a (non-)classical logic. This suggests a seamless extension of our application to instantiated argumentation frameworks [19]. Some preliminary experiments involve the instantiation of arguments as tuples composed of a pair of formulas (`support`, `claim`) in a

formal (recursively defined) logical language. We can then employ the shallow semantical embedding approach [12] to give a semantics to these formulas, including an appropriate definition for a logical consequence relation. Subsequently, we can then employ the latter to instantiate the corresponding attack relation between arguments in several different ways, namely, as rebutting, undermining and undercutting, as given by the semantics of the embedded logic. A detailed exploration of this is, however, further work.

## Acknowledgements

## References

[1] Amgoud, L., C. Cayrol, M. Lagasquie-Schiex and P. Livet, *On bipolarity in argumentation frameworks*, Int. J. Intell. Syst. **23** (2008), pp. 1062–1093.

[2] Arieli, O. and M. W. A. Caminada, *A QBF-based formalization of abstract argumentation semantics*, J. Appl. Log. **11** (2013), pp. 229–252.

[3] Baroni, P., M. Caminada and M. Giacomin, *An introduction to argumentation semantics*, Knowl. Eng. Rev. **26** (2011), pp. 365–410.

[4] Baroni, P., F. Cerutti, M. Giacomin and G. Guida, *Afra: Argumentation framework with recursive attacks*, International Journal of Approximate Reasoning **52** (2011), pp. 19–37.

[5] Baroni, P., D. M. Gabbay, M. Giacomin and L. van der Torre, "Handbook of formal argumentation," College Publications, 2018.

[6] Baroni, P., F. Toni and B. Verheij, *On the acceptability of arguments and its fundamental role in nonmonotonic reasoning, logic programming and n-person games: 25 years later*, Argument Comput. **11** (2020), pp. 1–14.

[7] Barrett, C. W., C. L. Conway, M. Deters, L. Hadarean, D. Jovanovic, T. King, A. Reynolds and C. Tinelli, *CVC4*, in: *CAV*, Lecture Notes in Computer Science **6806** (2011), pp. 171–177.

[8] Barrett, C. W., R. Sebastiani, S. A. Seshia and C. Tinelli, *Satisfiability modulo theories*, in: *Handbook of Satisfiability*, Frontiers in Artificial Intelligence and Applications **185**, IOS Press, 2009 pp. 825–885.

[9] Barringer, H., D. M. Gabbay and J. Woods, *Temporal dynamics of support and attack networks: From argumentation to zoology*, in: *Mechanizing Mathematical Reasoning*, Lecture Notes in Computer Science **2605** (2005), pp. 59–98.

[10] Baumann, R. and C. Spanring, *A Study of Unrestricted Abstract Argumentation Frameworks*, in: *IJCAI* (2017), pp. 807–813.

[11] Beierle, C., F. Brons and N. Potyka, *A software system using a SAT solver for reasoning under complete, stable, preferred, and grounded argumentation semantics*, in: *KI*, Lecture Notes in Computer Science **9324** (2015), pp. 241–248.

[12] Benzmüller, C., *Universal (meta-)logical reasoning: Recent successes*, Science of Computer Programming **172** (2019), pp. 48–62.

[13] Benzmüller, C. and D. Fuenmayor, *Value-Oriented Legal Argumentation in Isabelle/HOL*, in: *ITP*, LIPIcs **193** (2021), pp. 7:1–7:20.

[14] Benzmüller, C., X. Parent and L. W. N. van der Torre, *Designing normative theories for ethical and legal reasoning:* LogiKEy *framework, methodology, and tool support*, Artif. Intell. **287** (2020), p. 103348.

[15] Benzmüller, C. and L. C. Paulson, *Quantified multimodal logics in simple type theory*, Logica Universalis **7** (2013), pp. 7–20.

[16] Benzmüller, C. and D. S. Scott, *Automating free logic in HOL, with an experimental application in category theory*, Journal of Automated Reasoning **64** (2020), pp. 53–72.

[17] Benzmüller, C. and P. Andrews, *Church's Type Theory*, in: E. N. Zalta, editor, *The Stanford Encyclopedia of Philosophy*, Metaphysics Research Lab, Stanford University, 2019, summer 2019 edition .

[18] Besnard, P. and S. Doutre, *Checking the acceptability of a set of arguments*, in: *NMR*, 2004, pp. 59–64.

[19] Besnard, P. and A. Hunter, *A review of argumentation based on deductive arguments*, in: P. Baroni, D. M. Gabbay, M. Giacomin and L. van der Torre, editors, *Handbook of formal argumentation*, College Publications, 2018 pp. 437–484.

[20] Biere, A., M. Heule, H. van Maaren and T. Walsh, editors, "Handbook of Satisfiability," Frontiers in Artificial Intelligence and Applications **185**, IOS Press, 2009.

[21] Blanchette, J. C., S. Böhme and L. C. Paulson, *Extending sledgehammer with SMT solvers*, J. Autom. Reason. **51** (2013), pp. 109–128.

[22] Blanchette, J. C., C. Kaliszyk, L. C. Paulson and J. Urban, *Hammering towards QED*, Journal of Formalized Reasoning **9** (2016), pp. 101–148.

[23] Blanchette, J. C. and T. Nipkow, *Nitpick: A counterexample generator for higher-order logic based on a relational model finder*, in: M. Kaufmann and L. C. Paulson, editors, *ITP 2010*, LNCS **6172** (2010), pp. 131–146.

[24] Brewka, G., S. Ellmauthaler, H. Strass, J. P. Wallner and S. Woltran, *Abstract dialectical frameworks*, in: P. Baroni, D. M. Gabbay, M. Giacomin and L. van der Torre, editors, *Handbook of formal argumentation*, College Publications, 2018 pp. 237–285.

[25] Caminada, M., *On the issue of reinstatement in argumentation*, in: *European Workshop on Logics in Artificial Intelligence*, Springer, 2006, pp. 111–123.

[26] Caminada, M., *Semi-stable semantics*, in: *COMMA*, Frontiers in Artificial Intelligence and Applications **144** (2006), pp. 121–130.

[27] Caminada, M. and G. Pigozzi, *On judgment aggregation in abstract argumentation*, Autonomous Agents and Multi-Agent Systems **22** (2011), pp. 64–102.

[28] Caminada, M. W. A., W. A. Carnielli and P. E. Dunne, *Semi-stable semantics*, J. Log. Comput. **22** (2012), pp. 1207–1254.

[29] Caminada, M. W. A. and D. M. Gabbay, *A logical account of formal argumentation*, Stud Logica **93** (2009), pp. 109–145.

[30] Cayrol, C. and M. Lagasquie-Schiex, *On the acceptability of arguments in bipolar argumentation frameworks*, in: *ECSQARU*, Lecture Notes in Computer Science **3571** (2005), pp. 378–389.

[31] Cayrol, C. and M. Lagasquie-Schiex, *Logical encoding of argumentation frameworks with higher-order attacks and evidential supports*, Int. J. Artif. Intell. Tools **29** (2020), pp. 2060003:1–2060003:50.

[32] Cerutti, F., S. A. Gaggl, M. Thimm and J. P. Wallner, *Foundations of implementations for formal argumentation*, in: P. Baroni, D. M. Gabbay, M. Giacomin and L. van der Torre, editors, *Handbook of formal argumentation*, College Publications, 2018 pp. 689–768.

[33] Cerutti, F., M. Vallati and M. Giacomin, *An Efficient Java-Based Solver for Abstract Argumentation Frameworks: jArgSemSAT*, Int. J. Artif. Intell. Tools **26** (2017), pp. 1750002:1–1750002:26.

[34] Church, A., *A formulation of the simple theory of types*, J. Symb. Log. **5** (1940), pp. 56–68.

[35] de Moura, L. M. and N. Bjørner, *Z3: an efficient SMT solver*, in: C. R. Ramakrishnan and J. Rehof, editors, *Tools and Algorithms for the Construction and Analysis of Systems, 14th International Conference, TACAS 2008. Proceedings*, Lecture Notes in Computer Science **4963** (2008), pp. 337–340.

[36] de Saint-Cyr, F. D., P. Bisquert, C. Cayrol and M. Lagasquie-Schiex, *Argumentation update in YALLA (yet another logic language for argumentation)*, Int. J. Approx. Reason. **75** (2016), pp. 57–92.

[37] Dung, P. M., *On the acceptability of arguments and its fundamental role in nonmonotonic reasoning, logic programming and n-person games*, Artif. Intell. **77** (1995), pp. 321–358.

[38] Dung, P. M., P. Mancarella and F. Toni, *Computing ideal sceptical argumentation*, Artif. Intell. **171** (2007), pp. 642–674.

[39] Dunne, P. E. and T. J. M. Bench-Capon, *Coherence in finite argument systems*, Artif. Intell. **141** (2002), pp. 187–203.

[40] Dvorák, W., M. Järvisalo, J. P. Wallner and S. Woltran, *Complexity-sensitive decision procedures for abstract argumentation*, Artif. Intell. **206** (2014), pp. 53–78.

[41] Egly, U. and S. Woltran, *Reasoning in argumentation frameworks using quantified boolean formulas*, in: P. E. Dunne and T. J. M. Bench-Capon, editors, *Computational Models of Argument: Proceedings of COMMA 2006, September 11-12, 2006, Liverpool, UK*, Frontiers in Artificial Intelligence and Applications **144** (2006), pp. 133–144.

[42] Fuenmayor, D. and C. Benzmüller, *Computer-supported analysis of arguments in climate engineering*, in: *CLAR*, Lecture Notes in Computer Science **12061** (2020), pp. 104–115.

[43] Fuenmayor, D. and A. Steen, *A Flexible Approach to Argumentation Framework Analysis using Theorem Proving*, in: B. Liao, J. Luo and L. van der Torre, editors, *Logics for New-Generation AI 2021* (2021), pp. 18–32.

[44] Fuenmayor, D. and A. Steen, *Isabelle/HOL sources associated with this paper*, Online available at Github: `https://github.com/aureleeNet/formalizations` (2021).

[45] Gabbay, D. M., "Meta-logical Investigations in Argumentation Networks," Studies in Logic – Mathematical Logic and Foundations **44**, College Publications, 2013.

[46] Gleißner, T., A. Steen and C. Benzmüller, *Theorem provers for every normal modal logic*, in: T. Eiter and D. Sands, editors, *LPAR-21, 21st International Conference on Logic for Programming, Artificial Intelligence and Reasoning, Maun, Botswana, May 7-12, 2017*, EPiC Series in Computing **46** (2017), pp. 14–30.

[47] Henkin, L., *Completeness in the theory of types*, J. Symb. Log. **15** (1950), pp. 81–91.

[48] Michaelis, J. and T. Nipkow, *Propositional proof systems*, Archive of Formal Proofs (2017), `https://isa-afp.org/entries/Propositional_Proof_Systems.html`, Formal proof development.

[49] Nipkow, T., L. C. Paulson and M. Wenzel, "Isabelle/HOL - A Proof Assistant for Higher-Order Logic," Lecture Notes in Computer Science **2283**, Springer, 2002.

[50] Nolt, J., *Free Logic*, in: E. N. Zalta, editor, *The Stanford Encyclopedia of Philosophy*, Metaphysics Research Lab, Stanford University, 2021, Fall 2021 edition .

[51] Reynolds, A., J. C. Blanchette, S. Cruanes and C. Tinelli, *Model finding for recursive functions in SMT*, in: *IJCAR*, Lecture Notes in Computer Science **9706** (2016), pp. 133–151.

[52] Schulz, S., *E – a brainiac theorem prover*, AI Commun. **15** (2002), pp. 111–126.

[53] Steen, A., "Extensional Paramodulation for Higher-Order Logic and its Effective Implementation Leo-III," DISKI – Dissertations in Artificial Intelligence **345**, Akademische Verlagsgesellschaft AKA GmbH, Berlin, 2018.

[54] Steen, A. and C. Benzmüller, *Sweet SIXTEEN: Automation via embedding into classical higher-order logic*, Logic and Logical Philosophy **25** (2016), pp. 535–554.

[55] Steen, A. and C. Benzmüller, *Extensional Higher-Order Paramodulation in Leo-III*, J. Autom. Reason. **65** (2021), pp. 775–807.

[56] Toni, F. and M. Sergot, *Argumentation and answer set programming*, in: M. Balduccini and T. C. Son, editors, *Logic Programming, Knowledge Representation, and Nonmonotonic Reasoning: Essays in Honor of Michael Gelfond*, Lecture Notes in Computer Science **6565** (2011), pp. 164–180.

[57] Verheij, B., *Two approaches to dialectical argumentation: admissible sets and argumentation stages*, Proc. NAIC **96** (1996), pp. 357–368.

[58] Verheij, B., *Deflog: on the logical interpretation of prima facie justified assumptions*, Journal of Logic and Computation **13** (2003), pp. 319–346.

[59] Wallner, J. P., G. Weissenbacher and S. Woltran, *Advanced SAT techniques for abstract argumentation*, in: *CLIMA*, Lecture Notes in Computer Science **8143** (2013), pp. 138–154.

[60] Wenzel, M., *Isabelle/Isar—a generic framework for human-readable proof documents*, From Insight to Proof-Festschrift in Honour of Andrzej Trybulec **10** (2007), pp. 277–298.