

Evaluating Networks of Arguments: A Case Study in Mīmāṃsā Dialectics*

Kees van Berkel¹, Agata Ciabattoni¹, Elisa Freschi², and Sanjay Modgil³

¹ TU Wien, Vienna, Austria {kees,agata}@logic.at

² ÖAW and UniWien, Vienna, Austria elisa.freschi@gmail.com

³ King’s College London, United Kingdom sanjay.modgil@kcl.ac.uk

Abstract. We formalize networks of authored arguments. These networks are then mapped to *ASPIC*⁺ theories that subsequently instantiate Extended Argumentation Frameworks. Evaluation of arguments in the latter determines the status of the arguments in the source networks. The methodology is illustrated through a collaboration between scholars of South Asian philosophy, logicians and formal argumentation theorists, analyzing excerpts of Sanskrit texts concerning a controversial normative debate within the philosophical school of Mīmāṃsā.

Keywords: Instantiated Arguments · Extended Argumentation Frameworks · *ASPIC*⁺ · Argument Networks · Dialectics · Mīmāṃsā philosophy.

1 Introduction

Dung’s seminal theory of argumentation [7] provides foundations for dynamic and distributed nonmonotonic reasoning [3]. Given a set of logical formulae, one defines arguments (sets of logical formulae and their inferred conclusions) and a binary (attack) relation amongst them, encoding that one argument is a counter-argument to another. The status of arguments in the resulting argumentation framework (*AF*) is then evaluated, and the claims of the winning arguments identify the non-monotonic inferences from the ‘instantiating’ set of logical formulae. However, whereas the above procedure is often static, argumentation in practice is typically dynamic and dialectic, where arguments are authored incrementally rather than being defined by a given, fixed set of formulae. Moreover, in practice not only attacks, but also collective attacks [15], support relations [1, 16], attacks on attacks [10] etc. are specified as holding between arguments, thus defining networks of authored arguments (see [11] for more details).

In this paper, we formalise a methodology that was informally proposed in [11], and that accounts for the above described dynamic authoring of argument networks. Here, structured arguments—i.e., arguments whose internal logical structure is specified—are first related by attacks and supports, and can express preferences over arguments. These arguments are then mapped to their constituent formulae and rules so as to define an *ASPIC*⁺ theory [14] (a promising

* Work partially funded by the projects WWTF MA16-028 and FWF W1255-N23.

formal approach to structured-argumentation [9]). The *ASPIC*⁺ theory, subsequently, instantiates an Extended Argumentation Framework (*EAF*) [10] (an extension of Dung’s *AF* that accommodates arguments expressing preferences over other arguments through attacking attack relations). Finally, the evaluation of these arguments in the resulting *EAF* then determines the status of arguments in the source network, and consequently the inferences defined by the constituent formulae and rules in the original authored arguments. As argued in [11], this methodology is proposed as a more rigorous approach to evaluating arguments—and the defined inferences—in authored networks, as compared with directly evaluating the arguments in the source network (as typically done by scholars defining networks relating arguments by relations other than attacks).

The usefulness of the proposed formalization, corresponding to the aforementioned methodology, will be exemplified through a case study which resulted from a collaboration between scholars in South Asian philosophy, argumentation theorists and logicians. It involves a formal analysis of some excerpts of (Sanskrit) commentaries by philosophers of the school of *Mīmāṃsā*, and their application to South Asian jurisprudence. This school originated in ancient India more than two millennia ago and was devoted to the analysis of normative statements in the Vedas, the sacred texts of the so-called Hinduism. The dialectic nature of *Mīmāṃsā* argumentation, its structured analyses and its use of abstract logical principles, makes it particularly suitable for exhibiting the formal extensions introduced in this paper. In particular, we analyzed a portion of the debate on the immolation of widows on their husbands’ funeral pyre, i.e., the so-called *satī* ritual. This debate has had deep socio-political implications in South Asia since the 9th c. until today (e.g., see [4, 18]) and has been broadly dealt with by South Asian jurists and philosophers, primarily of the *Mīmāṃsā* school.

Plan of the paper: We assume familiarity with abstract argumentation and briefly recap *EAF*s and *ASPIC*⁺ instantiations of *EAF*s in Sec.2 (for a review see [2, 17]). In Sec.3, we define *ASPIC*⁺ argument networks which represent dynamically authored arguments and their relations, as specified by domain experts using some putative authoring—i.e. argument diagramming—tool. We also define the mapping of these networks to *ASPIC*⁺ theories. We use the *Mīmāṃsā* debate on *satī* as a case study to exemplify the formalised methodology in Sec.4.

2 Background: *EAF*s and *ASPIC*⁺

Extended Argumentation Frameworks. Along with the usual binary attack relation (\mathcal{C}) over arguments, *Extended Argumentation Frameworks* (*EAF*s) [10] extend Argumentation Frameworks (*AF*s) [7] to also include a *pref-attack* relation \mathcal{D} : i.e. an argument expresses that Y is *preferred* to X attacks the binary attack from X to Y , so that the latter attack does not succeed as a defeat.

Definition 1. *An **EAF** is a tuple $(\mathcal{A}, \mathcal{C}, \mathcal{D})$, \mathcal{A} is a set of arguments, $\mathcal{C} \subseteq \mathcal{A} \times \mathcal{A}$, $\mathcal{D} \subseteq \mathcal{A} \times \mathcal{C}$, and if $(Z, (X, Y))$, $(Z', (Y, X)) \in \mathcal{D}$ then (Z, Z') , $(Z', Z) \in \mathcal{C}$.*

Notice that the constraint on the relation \mathcal{D} ensures that if arguments Z and Z' respectively *pref-attack* (X, Y) and (Y, X) , then Z and Z' express contradictory

preferences—i.e., Y is preferred to X , respectively X is preferred to Y —and so themselves symmetrically (i.e. mutually) attack each other. Henceforth, we focus on *bounded hierarchical EAFs*, stratified so that attacks at a level i are only pref-attacked by arguments at the next level (preserving rationality [13]).

Definition 2. $\Delta=(\mathcal{A},\mathcal{C},\mathcal{D})$ is a bounded hierarchical EAF (**bh-EAF**) iff there exists a partition $\Delta_H = (((\mathcal{A}_1,\mathcal{C}_1),\mathcal{D}_1), \dots, ((\mathcal{A}_n,\mathcal{C}_n),\mathcal{D}_n))$ s.t. $\mathcal{D}_n = \emptyset$, and:

- $\mathcal{A}=\bigcup_{i=1}^n \mathcal{A}_i$, $\mathcal{C}=\bigcup_{i=1}^n \mathcal{C}_i$, $\mathcal{D}=\bigcup_{i=1}^n \mathcal{D}_i$, and for $1\leq i\leq n$, $(\mathcal{A}_i,\mathcal{C}_i)$ is an AF
- $(C, (A,B)) \in \mathcal{D}_i$ implies $(A,B) \in \mathcal{C}_i$, $C \in \mathcal{A}_{i+1}$

The notion of a successful attack (i.e. *defeat*) is then parameterised with respect to the preferences specified by some given set S of arguments: i.e., Y *defeats_S* X (denoted $Y \rightarrow^S X$) iff $(Y,X) \in \mathcal{C}$ and $\neg\exists Z \in S$ s.t. $(Z,(Y,X)) \in \mathcal{D}$.

Then, a set S is *EAF conflict free* when it does not admit arguments that symmetrically attack, but S can contain some Y and X such that Y *asymmetrically* attacks X , given a $Z \in S$ that pref-attacks the attack from Y to X .

Furthermore, since attacks can themselves be attacked, these attacks need to be reinstated (defended) by attacking arguments that pref-attack. That is, the acceptability of an argument X w.r.t. a set S requires that there is a *reinstatement set* for any reinstating defeat:

Definition 3. Let $S \subseteq \mathcal{A}$ in $(\mathcal{A},\mathcal{C},\mathcal{D})$. Let $R_S=\{X_1 \rightarrow^S Y_1, \dots, X_n \rightarrow^S Y_n\}$ s.t. for $1\leq i\leq n$, $X_i \in S$. We call R_S a *reinstatement set* for $A \rightarrow^S B$, iff $A \rightarrow^S B \in R_S$, and $\forall X \rightarrow^S Y \in R_S$, $\forall Y'$ s.t. $(Y', (X,Y)) \in \mathcal{D}$, $\exists X' \rightarrow^S Y' \in R_S$.

Furthermore, X is *acceptable* w.r.t. $S \subseteq \mathcal{A}$ iff for all $Y \in \mathcal{A}$ s.t. $Y \rightarrow^S X$, there is a $Z \in S$ s.t. $Z \rightarrow^S Y$, and there is a reinstatement set for $Z \rightarrow^S Y$.

For *bh-EAFs*, the semantic extensions are defined as for *AFs*. That is, let S be a conflict free set: S is an *admissible* extension iff all arguments in S are acceptable w.r.t. S ; S is *complete* iff it is admissible and all arguments acceptable w.r.t. S are in S ; S is *preferred* iff it is a set inclusion maximal complete extension; S is the (unique) *grounded* extension iff it is the set inclusion minimal complete extension; S is *stable* iff $\forall Y \notin S$, $\exists X \in S$ s.t. $X \rightarrow^S Y$. Lastly, for $e \in \{\text{complete, preferred, grounded, stable}\}$, $X \in \mathcal{A}$ is *credulously* (*sceptically*) justified under the e semantics, if X belongs to at least one (all) e extension(s).

ASPIC⁺ Instantiations of EAFs. *ASPIC⁺* [14] is a general framework in which one is free to choose a logical language \mathcal{L} . One is also free to specify defeasible and strict inference rules, as well as ‘axiom’ and ‘ordinary’ premises for construction of arguments. Furthermore, it facilitates preference relations over arguments, used to determine when attacks succeed as defeats. Defeasible rules are typically domain specific, while strict rules may either encode domain specific infallible inferences or inference rules of some deductive logic. In this system, only the fallible ordinary premises and fallible consequents of defeasible rules can be attacked. Axiom premises are infallible and conclusions of strict rules cannot be

attacked. A partial function assigns names (wff in \mathcal{L}) to defeasible rules, so that applications of defeasible rules can be invalidated by arguments that claim the negation of the rule name. Finally, $ASPIC^+$ allows one to specify a contrary function specifying when formulae in \mathcal{L} are said to be in conflict. In this paper, we assume that such conflicts are symmetric.

$ASPIC^+$ poses constraints on the above choices to ensure that the outcomes of evaluating the Dung frameworks instantiated by $ASPIC^+$ arguments and defeats, are rational [5]. In this work, the following review of $ASPIC^+$ [14] suffices:

Definition 4. An argumentation theory is a tuple $AT = (\mathcal{L}, -, \mathcal{R}, n, \mathcal{K})$ where \mathcal{L} is a logical language, and:

- $\mathcal{R} = \mathcal{R}_s \cup \mathcal{R}_d$ is a set of strict (\mathcal{R}_s) and defeasible (\mathcal{R}_d) inference rules, respectively of the form $\varphi_1, \dots, \varphi_n \rightarrow \varphi$ and $\varphi_1, \dots, \varphi_n \Rightarrow \varphi$ (where φ_i and φ are metavariables ranging over wff in \mathcal{L});
- $n : \mathcal{R}_d \mapsto \mathcal{L}$ is a partial naming function;
- $\mathcal{K} = \mathcal{K}_n \cup \mathcal{K}_p$ where $\mathcal{K} \subseteq \mathcal{L}$, \mathcal{K}_n is a set of axiom premises, \mathcal{K}_p is a set of ordinary premises, and $\mathcal{K}_n \cap \mathcal{K}_p = \emptyset$.
- for all wff ϕ in \mathcal{L} , if $\varphi \in \overline{\psi}$ then $\psi \in \overline{\varphi}$. (In this case, we say that ψ and φ are contradictories, which is denoted by $\varphi = -\psi$.)

Henceforth, for convenience we write ‘ $\delta : \varphi_1, \dots, \varphi_n \Rightarrow \varphi$ ’ instead of explicitly declaring that n assigns the wff δ to the defeasible rule $\varphi_1, \dots, \varphi_n \Rightarrow \varphi$.

Definition 5. An $ASPIC^+$ argument A on the basis of an $AT (\mathcal{L}, -, \mathcal{R}, n, \mathcal{K})$ is:

1. φ if $\varphi \in \mathcal{K}$ with: $\text{Prem}(A) = \{\varphi\}$; $\text{Conc}(A) = \varphi$; $\text{Sub}(A) = \{\varphi\}$; $\text{Rules}(A) = \emptyset$; $\text{DefRules}(A) = \emptyset$; $\text{TopRule}(A) = \text{undefined}$.
2. $A_1, \dots, A_n \rightarrow / \Rightarrow \psi$ if A_1, \dots, A_n are arguments such that there exists a strict/defeasible rule $\text{Conc}(A_1), \dots, \text{Conc}(A_n) \rightarrow / \Rightarrow \psi$ in $\mathcal{R}_s / \mathcal{R}_d$, with:
 - $\text{Prem}(A) = \text{Prem}(A_1) \cup \dots \cup \text{Prem}(A_n)$; $\text{Conc}(A) = \psi$;
 - $\text{Sub}(A) = \text{Sub}(A_1) \cup \dots \cup \text{Sub}(A_n) \cup \{A\}$;
 - $\text{Rules}(A) = \bigcup_{i=1}^n \text{Rules}(A_i) \cup \{\text{Conc}(A_1), \dots, \text{Conc}(A_n) \rightarrow / \Rightarrow \psi\}$;
 - $\text{TopRule}(A) = \text{Conc}(A_1), \dots, \text{Conc}(A_n) \rightarrow / \Rightarrow \psi$;
 - $\text{DefRules}(A) = \bigcup_{i=1}^n \text{DefRules}(A_i) \cup \{A_1, \dots, A_n \Rightarrow \psi\}$ if
 - $\text{TopRule}(A) = A_1, \dots, A_n \Rightarrow \psi$ and $\bigcup_{i=1}^n \text{DefRules}(A_i)$ otherwise.

The notation in Def.5 is generalised to sets of arguments in the usual way: e.g., letting $E = \{A_1, \dots, A_n\}$, then $\text{DefRules}(E) = \bigcup_{i=1}^n \text{DefRules}(A_i)$.

Definition 6. Let A, B and B' be $ASPIC^+$ arguments.

- A undercuts argument B (on B') iff $\text{Conc}(A) \in \overline{n(r)}$ for some $B' \in \text{Sub}(B)$ such that $\text{TopRule}(B') = r$.
- A rebuts argument B on (B') iff $\text{Conc}(A) = -\varphi$ for some $B' \in \text{Sub}(B)$ of the form $B_1'', \dots, B_n'' \Rightarrow \varphi$.
- A undermines B (on $B' = \varphi$) iff $\text{Conc}(A) = -\varphi$ for some $\varphi \in \text{Prem}(B) \setminus \mathcal{K}_n$.

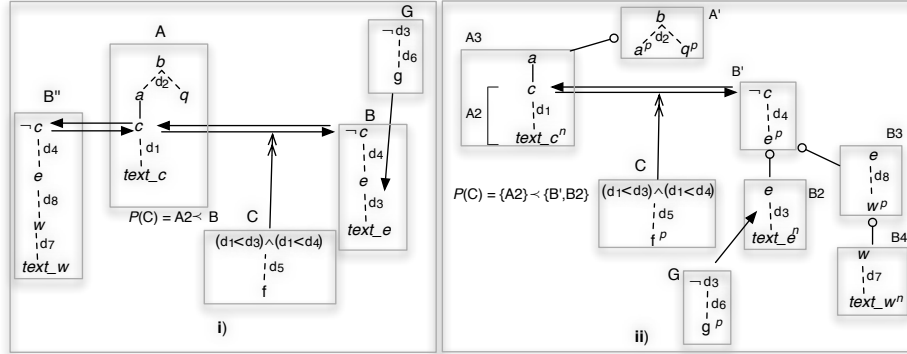


Fig. 1. Figure i) depicts the bh -EAF instantiated by the $ASPIC^+$ theory of Example 1. Figure ii) represents an $ASPIC^+$ EANS mapped to the $ASPIC^+$ theory of Example 1. Solid and dashed lines denote the application of strict and defeasible rules (resp.). We illustrate $(X, Y) \in \mathcal{C}$ with $X \rightarrow Y$, and $(Z, (X, Y)) \in \mathcal{D}$ with $Z \rightarrow (X \rightarrow Y)$.

When $ASPIC^+$ arguments instantiate an AF , a preference relation over the arguments is used to decide whether rebut or undermine attacks succeed as defeats; i.e., an attack from A to B succeeds only if $A \not\prec B'$. Undercuts succeed as defeats independently of preferences. Following [13], we will instantiate bh -EAFs in such a way that $ASPIC^+$ arguments may themselves conclude preferences over arguments (rather than assuming a given strict ordering \prec over arguments). Pref-attacks originating from these arguments may then target binary attacks, denying the success of the latter as defeats. As in [12], we assume a function \mathcal{P} that maps the conclusion of an individual argument to strict preferences over other arguments; e.g., given A and B with respective defeasible rules $\{r_1\}$ and $\{r_2, r_3\}$, if argument C concludes $(r_1 < r_2) \wedge (r_1 < r_3)$, then $\mathcal{P}(\text{Conc}(C)) = A \prec B$ (under the *Elitist* set ordering of [14]):

Definition 7. Let \mathcal{A} be a set of $ASPIC^+$ arguments, $A, B \in \mathcal{A}$ and $\mathcal{C} \subseteq \mathcal{A} \times \mathcal{A}$, s.t. $(A, B) \in \mathcal{C}$ iff A rebuts, undermines or undercuts B . Let $\mathcal{P} : \mathcal{L} \mapsto \prec$, where $\prec \subseteq \mathcal{A} \times \mathcal{A}$ is a strict partial ordering over \mathcal{A} . Then an $ASPIC^+$ instantiated EAF is a tuple $(\mathcal{A}, \mathcal{C}, \mathcal{D})$ defined as in Def.1, where $(C, (A, B)) \in \mathcal{D}$ iff A rebuts or undermines B on B' , and $A \prec B' \in \mathcal{P}(\text{Conc}(C))$.

Example 1. To illustrate the above, suppose the *argumentation theory* $AT = (\mathcal{L}, \neg, \mathcal{R}, n, \mathcal{K})$, with propositional language \mathcal{L} , a strict priority relation $<$ and:

- $\mathcal{K}_n = \{\text{text}_c, \text{text}_e, \text{text}_w\}$; $\mathcal{K}_p = \{f, g, q\}$; $\mathcal{R}_s = \{c \rightarrow a\}$;
- $\mathcal{R}_d = \{d_1: \text{text}_c \Rightarrow c; d_2: a, q \Rightarrow b; d_3: \text{text}_e \Rightarrow e; d_4: e \Rightarrow \neg c; d_5: f \Rightarrow (d_1 < d_3) \wedge (d_1 < d_4); d_6: g \Rightarrow \neg d_3\}$; $d_7: \text{text}_w \Rightarrow w$; $d_8: w \Rightarrow e$.
- $\varphi = \neg\psi$ just in case $\varphi = \neg\psi$ or $\psi = \neg\varphi$

We obtain the instantiated bh -EAF (see Fig.1-i), consisting of the following:

$$\mathcal{A} = \left\{ \begin{array}{l} A1 = [\textit{text}_c], A2 = [A1 \Rightarrow c], A3 = [A2 \rightarrow a], A4 = [q], A = [A3, A4 \Rightarrow b] \\ B1 = [\textit{text}_e], B2 = [B1 \Rightarrow e], B = [B2 \Rightarrow \neg c], \\ C1 = [f], C = [C1 \Rightarrow (d_1 < d_3) \wedge (d_1 < d_4)], G1 = [g], G = [g \Rightarrow \neg d_3], \\ B5 = [\textit{text}_w], B4 = [B5 \Rightarrow w], B3 = [B4 \Rightarrow e], B'' = [B3 \Rightarrow \neg c] \end{array} \right\}$$

$$\mathcal{C} = \{(B, A)^*, (B, A2), (A2, B), (G, B), (G, B2), (B'', A)^\dagger, (B'', A2), (A2, B'')\}$$

$$\mathcal{D} = \{(C, (A2, B))\} \quad (\text{NB. } (*) B \text{ attacks } A \text{ on } A2 \text{ and } (\dagger) B'' \text{ attacks } A \text{ on } A2.)$$

The single *grounded* extension of this EAF is the set $E = \{G1, G, C1, C, B1, A1, A4, B5, B4, B3\}$. The two *preferred/stable* extensions are $EU\{B''\}$ and $EU\{A, A2, A3\}$.

3 Towards Formalizing Networks of Authored Arguments

Many extensions of Dung *AFs* are motivated by natural language examples in which arguments and their relations are *dynamically* specified, rather than being instantiated by a given *static* set of formulae. Following this observation, [11] argues that *networks* of arguments related by attacks, supports, collective attacks, recursive attacks on attacks etc., are thus more properly motivated in argument authoring contexts in which (human) domain experts specify and relate arguments incrementally, and hence, dynamically. A principled way to then evaluate these networks is to map their contents to an *ASPIC⁺* theory that, subsequently, instantiates an *AF* or *EAF*.

This section formally realises the above informal proposal in [11]. We define networks of *ASPIC⁺* arguments authored by domain experts who specify the contents of these arguments—that is, axiom and ordinary premises, strict and defeasible inference rules—as well as support, attack and pref-attack relations. An argument Y is used to support X only if Y supplies the rationale (argument) for an ordinary (i.e. fallible) premise in X ; axiom premises, which typically encode empirically validated information and so cannot be challenged, need not be supported. Hence, when authoring arguments one must distinguish between ordinary and axiom premises (with the respective superscripts p and n). We first define networks related by attacks and supports (Def.8), and then hierarchies of such networks that include pref-attacks (Def.10):

Definition 8. An *ASPIC⁺ANS* (*Argument Network with Support*) is $\langle \mathcal{A}, \mathcal{C}, \mathcal{S} \rangle$:

- \mathcal{A} are *ASPIC⁺* arguments such that for all $X \in \mathcal{A}$, and for all $\alpha \in \text{prem}(X)$, α is labelled by p or n ;
- $(X, Y) \in \mathcal{C}$ implies (X, Y) is an *ASPIC⁺* attack as defined in Def.6, and;
- $(X, Y) \in \mathcal{S}$ implies $\exists \alpha^p \in \text{Prem}(Y)$ such that $\text{Conc}(X) = \alpha$, in which case we say that X supports Y on α . We may write $X \multimap Y$ to denote $(X, Y) \in \mathcal{S}$.

Since we assume authoring of *ANSs* by humans, we account for the possibility that not all valid attacks may be explicitly authored. Thus, Def.8 accommodates

that although for some $X, Y \in \mathcal{A}$, X attacks Y according to Def.6, this attack might not be diagrammed as such, and so $(X, Y) \notin \mathcal{C}$.

Next we define a chain of supporting arguments, and the ‘collapsing’ of a chain into a single argument, with each supported premise replaced by its supporting argument.

Definition 9. Let $\Delta = \langle \mathcal{A}, \mathcal{C}, \mathcal{S} \rangle$ be an ASPIC⁺ANS. Then \mathcal{S}_{chain} is a set of ordered sets defined as follows:

$$\mathcal{S}_{chain}(\Delta) = \{ \{A_1, \dots, A_n\} \mid \bigcup_{i=1}^n A_i \subseteq \mathcal{A}, \neg \exists X, \neg \exists Y \in \mathcal{A} \text{ s.t. } (A_1, X), (Y, A_n) \in \mathcal{S}, \text{ and for } i = 1 \dots n-1, (A_{i+1}, A_i) \in \mathcal{S} \}$$

The function $coll$ takes as input a chain of supporting arguments Γ , and returns a single argument if $|\Gamma| = 2$, else it returns a chain of supporting arguments Γ' in the case that $|\Gamma| > 2$:

- $coll(\{A_1, A_2\}) = A$, where A_2 supports A_1 on α , and A is the argument A_1 with A_2 replacing premise α in A_1 ;
- $coll(\{A_1, \dots, A_n\}) = coll(\{A_1, \dots, coll(A_{n-1}, A_n)\})$ if $n > 2$.

We now define bounded hierarchies of networks of attacking and supporting arguments, in which *pref*-attacks are directed at attacks in the next level down the hierarchy. Since arguments may be ‘backward extended’ by supporting arguments, so as to define chains, we propose a definition of *pref*-attacks originating from arguments whose conclusion is specified as mapping (via a function \mathcal{P}_{set}) to a preference over *chains of supporting* arguments; i.e., a preference ordering over sets rather than single arguments (cf. Def.7).

Definition 10. An ASPIC⁺EANS (Extended Argument Network with Support) is a tuple $\Delta = \langle \mathcal{A}, \mathcal{C}, \mathcal{S}, \mathcal{D} \rangle$ iff there exists a partition $\Delta_H = \langle ((\mathcal{A}_1, \mathcal{C}_1, \mathcal{S}_1), \mathcal{D}_1), \dots, ((\mathcal{A}_n, \mathcal{C}_n, \mathcal{S}_n), \mathcal{D}_n) \rangle$ such that $\mathcal{D}_n = \emptyset$, and:

- $\mathcal{A} = \bigcup_{i=1}^n \mathcal{A}_i$, $\mathcal{C} = \bigcup_{i=1}^n \mathcal{C}_i$, $\mathcal{S} = \bigcup_{i=1}^n \mathcal{S}_i$, $\mathcal{D} = \bigcup_{i=1}^n \mathcal{D}_i$, and for $i = 1 \dots n$, $\langle \mathcal{A}_i, \mathcal{C}_i, \mathcal{S}_i \rangle$ is an ASPIC⁺ANS.
- $(C, (A, B)) \in \mathcal{D}_i$ iff $C \in \mathcal{A}_{i+1}$, $(A, B) \in \mathcal{C}_i$, where A undermine or rebut attacks B on B' , and $\exists \{B', \dots, B_m\}, \exists \{A, \dots, A_n\} \in \mathcal{S}_{chain}((\mathcal{A}_i, \mathcal{C}_i, \mathcal{S}_i))$ s.t. $(\{A, \dots, A_n\} \prec \{B', \dots, B_m\}) \in \mathcal{P}_{set}(\text{conc}(C))$, where $\mathcal{P}_{set} : \mathcal{L} \mapsto \prec_s$, and $\prec_s \subseteq 2^{\mathcal{A}} \times 2^{\mathcal{A}}$ is a strict partial ordering over sets of arguments.

Finally, we define a mapping from an ASPIC⁺EANS to an ASPIC⁺ theory and the corresponding instantiation of a *bh-EAF*, which allows us to calculate the theory’s extensions. Notice that, if an argument X (not of the form $[\alpha]$) is available to support a premise α^p , then α is not included as a premise in the ASPIC⁺theory (given that a rationale has been provided for why α holds). Also observe that, in line with our remark on attacks following Def.8, X may not have been explicitly moved to support α .

Definition 11. Let $\Delta = \langle \mathcal{A}, \mathcal{C}, \mathcal{S}, \mathcal{D} \rangle$ be an ASPIC⁺EANS, \mathcal{P}_{set} a user specified function s.t. $\mathcal{P}_{set} : \mathcal{L} \mapsto \prec_s$, with strict partial ordering $\prec_s \subseteq 2^{\mathcal{A}} \times 2^{\mathcal{A}}$, and \mathcal{L} and $-$ a given language and contrary function, respectively. Then $AT_{\Delta} = \langle \mathcal{L}, -, \mathcal{R}, n, \mathcal{K} = \mathcal{K}_n \cup \mathcal{K}_p \rangle$ is defined as follows:

1. $\mathcal{R} = \text{Rules}(\mathcal{A})$;
2. $\mathcal{K}_n = \{\alpha \mid \alpha^n \in \text{prem}(\mathcal{A})\}$;
3. $\mathcal{K}_p = \{\alpha \mid \alpha^p \in \text{prem}(\mathcal{A}), \neg \exists X \in \mathcal{A} \text{ s.t. } \text{conc}(X) = \alpha \text{ and } X \text{ is not of the form } [\alpha]\}$;
4. $\forall r \in \text{DefRules}(\mathcal{A}), n(r) = \alpha$, where α does not appear in $\mathcal{K}_n \cup \mathcal{K}_p$, and α does not appear in the antecedent or consequent of a rule in \mathcal{R} .

Let \mathcal{A}' be a set of ASPIC⁺ arguments defined by AT_Δ . Then:

- $\forall X, A, B \in \mathcal{A}', \mathcal{P}(\text{conc}(X)) = A \prec B$ iff $\exists Y \in \mathcal{A} \text{ s.t. } \text{conc}(X) = \text{conc}(Y)$, and $\mathcal{P}_{\text{set}}(\text{conc}(Y)) = \Gamma_A \prec \Gamma_B, A = \text{coll}(\Gamma_A)$ and $B = \text{coll}(\Gamma_B)$.

Let \mathcal{C}' be the attack relation defined over \mathcal{A}' , such that $\forall (A, B) \in \mathcal{C}, (A, B) \in \mathcal{C}'$ iff A rebuts, undermines or undercuts B . Then $(\mathcal{A}, \mathcal{C}, \mathcal{D})$ is defined as in Definition 1, where $(C, (A, B)) \in \mathcal{D}$ iff A rebut or undermines B on B' and $A \prec B' \in \mathcal{P}(\text{Conc}(C))$.

Example 2. Consider the network of arguments in Fig.1-ii (mapped to the ASPIC⁺ theory of Ex.1) as authored by one or more users in the consecutive order $A', A3, B', B2, C, G, B3, B4$. Note that $B2$ supplies the rationale (i.e. argument) for the premise e^p in B' , and so supports B' . Hence e is not included as an ordinary premise in the ASPIC⁺ theory of Ex.1. In Fig.1-i the *bh-EAF* instantiated by the theory is presented and we obtain the credulously justified arguments $A, A2, A3$ and B'' (under preferred and stable semantics), and so the conclusions c, a, b and $\neg c$ are credulously supported.

4 Case Study: The *satī* Ritual

We now apply the methodology formalized in Sec. 3, and analyze (part of) the controversy surrounding widows immolating themselves on their husbands' funeral pyre (the *satī* ritual). Despite the numerous arguments available, for space reasons, we limit our analysis to a single Mīmāṃsā author, namely, Medhātithi (9th–10th c. Kashmir). The analysis captures the arguments (in the form of an ASPIC⁺EANS) as they are successively elucidated (and augmented with contextual information in the form of basic reasoning principles). As will be seen, Medhātithi argues that *satī* should not be performed.

Basic Mīmāṃsā principles. Over the last two millennia, philosophers of the Mīmāṃsā school have thoroughly analyzed prescriptive statements in the Vedas. They distinguish between three classes of normative statements (see, e.g., [8]): obligations, recommendations, and prohibitions. Prohibitions lead to no result if respected but to a sanction if not observed; recommendations, which are driven by a desire, lead to a result if fulfilled and to no sanction otherwise; obligations lead to a result if fulfilled and to a sanction if disregarded. Hence, for instance, if something is obligatory, it is not recommended. For our formalization of the *satī* debate, we can rely on some basic reasoning principles, which are either explicitly formulated or implicitly endorsed by all Mīmāṃsā authors, and are strict or defeasible. The list of principles is presented in Def.12 below.

Last, when dealing with Mīmāṃsā we distinguish between two levels of normative statements: the ones that can be directly found in the Vedas or in *smṛti* texts based on the Vedas, and those obtained from applying metarules identified by Mīmāṃsā authors. We will refer to the former as *prima facie* norms and to the latter as *derived* normative statements.

Definition 12. *The following list of principles are are Mīmāṃsā metarules:*

1. STRICT CONTEXTUAL PRINCIPLES:
 - D1 *Prima facie prohibitions and prima facie obligations are mutually exclusive.*^{††}
 - D2 *Prima facie recommendations and prima facie obligations are mutually exclusive.*[†]
 - D3 *Prohibitions and obligations are mutually exclusive.*^{††}
2. DEFAULT (DEFEASIBLE) CONTEXTUAL PRINCIPLES:
 - D4 *An obligation/prohibition/recommendation on the prima facie level, is also an obligation/prohibition/recommendation on the derived level.*
 - D5 *If an obligatory/prohibited action necessarily presupposes some (other) action, then that action is also obligatory/prohibited.*[†]
 - D6 *An argument supported by a rationale (i.e. a justification) is the preferred argument in case of a conflict between equipollent claims.*
 - D7 *If the Vedas/smṛtis prescribes an obligation/prohibition/recommendation, then we take the obligation/prohibition/recommendation to hold prima facie.*^{††}
 - D8 *a) If two actions cause identical effects and have equal normative status, they are analogous. b) Conclusions drawn for one case apply to analogue cases.*^{††}
 - D9 *If the Vedas/smṛtis explicitly mention a reward for a prescribed action, then a) the action brings about that result and b) it is prima facie recommended.*[†]
 - D10 *If an action causes some effect, which subsequently implies another effect, then the action causes the second effect as well.*^{††}

The symbol † indicates that the principle is explicitly stated by Mīmāṃsā authors, those with †† are not stated as rules, yet explicitly applied in Mīmāṃsā reasoning. The remaining rules are implicit assumptions that other metarules presuppose.

The Argument Against *satī*. A synopsis of Medhātithi’s argument against *satī*, as found in the Sanskrit source is presented in Fig.2. The arguments presented here are translated and interpreted by Sanskritists.⁴ We will elaborate on the separate steps of the argument, identifying the involved rules and premises, as well as the individual arguments and their relations. We process the above as it consecutively appears in the source, thus capturing the pivotal dialectic aspect of Mīmāṃsā argumentation in an *ASPIC*⁺ network.

The formal language used in our case study consists of unary predicates $O(X)$ to express ‘ X is obligatory’, and similarly predicates F and R expressing prohibitions and recommendations, respectively. We reserve $*$ as a superscript for *prima facie* norms (e.g., $O^*(satī)$), whereas the absence of $*$ indicates a *derived*

⁴ Different interpretations of these arguments might be implemented in *ASPIC*⁺, and compared and evaluated on their logical consequences.

-
- (A) [OPENING] The performance of *sati* causes a widow to take her life. The latter is, as an act of violence, prohibited for women as it is for men.
- (B) [OBJECTION] Performance of *sati* is obligatory, because this prescription is derived from an explicit occurrence in a *smṛti* text.
- (C) [FIRST REPLY] The referred *smṛti* text prescribing *sati* mentions a result, namely heaven, and therefore *sati* is recommended, not obligatory.
- (D) [SECOND REPLY] The ritual of *sati* is similar to the *śyena* sacrifice; that is, (i) both are performed due to the desire for their respective results and (ii) the performance of each transgresses a prohibition, namely, that of committing violence. By analogy, since *śyena* is prohibited due to a prohibition being violated, *sati* is prohibited too.
- (E) [ADDITIONAL ARGUMENT] The claim that the mentioned *smṛti* prescribes the performance of *sati*, expressed in (B), is based on a misinterpretation of the text. Hence, it does not follow from the *smṛti* that *sati* is obligatory.
-

Fig. 2. Summary of Medhātithi's Argument Against *sati*

norm. Furthermore, we interpret $cs(X, Y)$ as ‘ X causes Y ’; $eff(X, Y)$ as ‘ X has Y as an effect’; and we read $txtO^*(X)$ as ‘the authoritative texts state that X is obligatory’. Also, $sim(X, Y)$ expresses that ‘ X and Y are similar’, and $mis(X)$ express that ‘ X has been misinterpreted’. We chronologically label arguments with A, B, C, \dots etc. The usage of the other terms will be clear from the context: e.g., we use *sati* for *sati* and *hvn* for ‘heaven’. Note also that predicate names for defeasible inference rules, will take as arguments the variables and constants that appear in the rule named. Last, the *contrary function* is defined so that $\phi = -\psi$ iff $\phi = \neg\psi$ or $\psi = \neg\phi$, and $F(X) = -O(X)$, $F^*(X) = -O^*(X)$, $R^*(X) = -O^*(X)$ (the latter three correspond to D1–D3 of Def.12). Recall that $-$ determines *contraries* (Def.4). Hence, for example $F(X) = -O(X)$ denotes that obligations and prohibitions are mutually exclusive: i.e., ‘if X is obligatory, then X is not forbidden and vice versa’ (observe that $-$ is not to be confused with logical negation \neg). Let us proceed to the first sub-argument, put forward by Medhātithi.

Argument (A) (shown in Fig. 3) claims that *sati* is prohibited because *sati* is a form of taking one’s life, which equates with self-violence: $cs(sati, s_vio)$. Furthermore, self-violence is an instance of violence in general: $eff(s_vio, vio)$. Any performance of violence, however, is *prima facie* prohibited: $F^*(vio)$. Hence, it is concluded, *sati* must be prohibited too: $F(sati)$. In the above, concrete instances of the following generic rules were applied:

$$\mathcal{R}_d(A) = \left\{ \begin{array}{l} d_{10}(Act, Eff_1, Eff_2): cs(Act, Eff_1), eff(Eff_1, Eff_2) \Rightarrow cs(Act, Eff_2); \\ d_5(Act, Eff): cs(Act, Eff), F^*(Eff) \Rightarrow F(Act); \end{array} \right\}$$

In the corresponding formal argument A , the variables Act, Eff_1 and Eff_2 are respectively substituted by *sati*, *s_vio* and *vio*. Note that the labelling of the rules in $\mathcal{R}_d(A)$ corresponds to the list of Mīmāṃsā principles presented in Def.12.

Subsequently, in argument (B), an opponent objects to (A) by asserting that *sati* is instead obligatory— $O(sati)$ —since the obligation is *prima facie*: $O^*(sati)$.

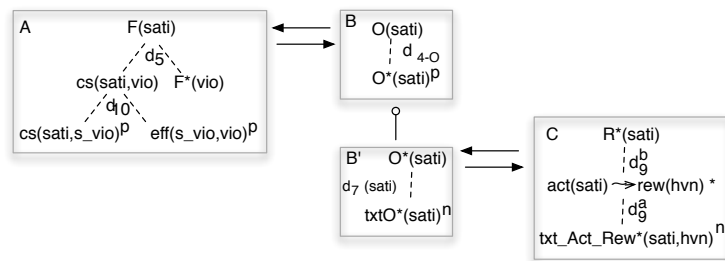


Fig. 3. Arguments A, B, B', C . Rules d_i are shown without instantiated variables.

Claim $O^*(sati)$ is itself supported by argument **(B')** referencing the passage with the prescription ‘the widow should die after her husband’: $txtO^*(sati)$. The above reasoning uses instantiated applications of the following rules:

$$\mathcal{R}_d(B) = \{d_{4-O}(X): O^*(X) \Rightarrow O(X); \quad d_7(X): txtO^*(X) \Rightarrow O^*(X).\}$$

In reply to **(B)**, argument **(C)** asserts that $sati$ is instead a *prima facie* recommendation: $R^*(sati)$. The claim is substantiated by the observation that, (i) the *smṛti* passage mentioning $sati$ explicitly relates the performance of $sati$ to a specific reward, namely the reward of heaven: $txt_Act_Rew^*(sati, hvn)$. (ii) Explicit mention of a reward identifies a norm as a *prima facie* recommendation:

$$\mathcal{R}_d(C) = \left\{ \begin{array}{l} d_9^a(Act, Rew): txt_Act_Rew^*(Act, Rew) \Rightarrow act(Act) \rightsquigarrow rew(Rew)^*; \\ d_9^b(Act, Rew): act(Act) \rightsquigarrow rew(Rew)^* \Rightarrow R^*(Act). \end{array} \right\}$$

The contrary function \sim implies a symmetric attack between arguments C and B' . Fig.3 shows argument A, B, B' and C , where B supports B' and A and B symmetrically attack each other on the basis of the defined contrary function.

Argument **(D)** is a reply to **(B)**, claiming that $sati$ is in fact prohibited: $F(sati)$. This claim follows from the assertions that (i) $sati$ is similar to the *śyena* sacrifice—i.e., $sim(sye, sati)$ —and (ii) the performance of *śyena* is prohibited: $F(sye)$. Note, the *śyena* sacrifice is a controversial Vedic ritual which results in the death of one’s enemy; e.g. see [6].) By *analogy*, since a performance of $sati$ violates the prohibition of violence too, we conclude that $sati$ must also be prohibited. Clearly, B and D symmetrically attack each other.

A successive argument **(D')** then supports the premise $sim(sye, sati)$ of **(D)**: Both *śyena* and $sati$ are recommendations due to fact that they depend on a desired result: $R(sye)$ and $R(sati)$. In particular, $R(sati)$ is justified given the earlier argument **(C)**, whose claim $R^*(sati)$ is again included as a premise (i.e., lemma) in **(D')**, and is used to infer $R(sati)$ via the principle (D4) (i.e., rule d_{4-R}). Subsequently, the performance of *śyena* implies violence—i.e., $cs(sye, vio)$ —as does the performance of $sati$. Note that $sati$ causing violence was justified earlier in **(A)**; hence, in **(D')** this fact is included as a premise rather than repeated as an argument (cf. the use of lemmas). Therefore, *śyena* and $sati$ are similar.

In support of the premise $F(sye)$ of **(D)** the argument **(D'')** is added, explaining that the *śyena* sacrifice is prohibited because performing *śyena* implies

violence – $cs(sye, vio)$ – and violence is *prima facie* forbidden: $F^*(vio)$. The rules applied in the corresponding formal arguments D, D' and D'' are as follows:

$$\mathcal{R}_d(D) = \left\{ \begin{array}{l} d_5(Act, Eff): cs(Act, Eff), F^*(Eff) \Rightarrow F(Act); \\ d_8^b(Act_1, Act_2): sim(Act_1, Act_2), F(Act_1) \Rightarrow F(Act_2); \\ d_8^a(Act_1, Act_2, Eff): cs(Act_1, Eff), cs(Act_2, Eff), R(Act_1), R(Act_2) \\ \quad \Rightarrow sim(Act_1, Act_2); \\ d_{4-R}(X): R^*(X) \Rightarrow R(X). \end{array} \right\}$$

Lastly, Medhātithi argues in **(E)** that the interpretation of the *smṛti* prescribing *sati*, as purported in **(B')**, is based on a misinterpretation of the word ‘after’: *mis(after)*. The correct interpretation of the *smṛti* prescription is not ‘dying *immediately* after’, but rather ‘dying *sometime* after’; the latter interpretation is in harmony with the Vedic prescription ‘one should not depart before one’s natural lifespan’. Hence, the *smṛti* does not prescribe *sati*. Therefore, E attacks B' via an *undercut* on the *instance* of the d_7 -rule; that is, the interpretative inference step encoded in $d_7(sati)$ is invalidated.

Additionally, the Mīmāṃsā principle (D6) gives rise to an additional argument **(F)**: namely, **(C)** uses a rule that encodes a rationale justifying why *sati* is a *prima facie* recommendation, in contrast to **(B')**’s rule which merely claims that *sati* is *prima facie* obligatory. Hence, prioritising the former rule over the latter licenses $\mathcal{P}_{set}(\text{conc}()F) = \{B'\} \prec \{C\}$. The rules applied in the formal correspondents E and F are, respectively:

$$\begin{aligned} \mathcal{R}_d(E) &= \{ d_E(\text{after}, \text{sati}): \text{mis}(\text{after}) \Rightarrow \neg d_7(\text{sati}). \} \\ \mathcal{R}_d(F) &= \{ d_F(\text{sati}, \text{hvn}): \text{True} \Rightarrow d_7(\text{sati}) < d_9^b. \} \end{aligned}$$

The resulting formal theory of Medhātithi’s argument is defined accordingly:

Definition 13. *The following presents the ASPIC⁺ Argumentation Theory of Medhātithi’s argument against *sati*:*

1. $\mathcal{R}_s = \emptyset$
2. $\mathcal{R}_d = \left\{ \begin{array}{l} d_{10}(Act, Eff_1, Eff_2): cs(Act, Eff_1), eff(Eff_1, Eff_2) \Rightarrow cs(Act, Eff_2); \\ d_5(Act, Eff): cs(Act, Eff), F^*(Eff) \Rightarrow F(Act); \\ d_{4-O}(X): O^*(X) \Rightarrow O(X); \\ d_{4-R}(X): R^*(X) \Rightarrow R(X); \\ d_{4-F}(X): F^*(X) \Rightarrow F(X); \\ d_7(X): \text{txt}O^*(X) \Rightarrow O^*(X); \\ d_9^a(Act, Rew): \text{txt_Act_Rew}^*(Act, Rew) \Rightarrow \text{act}(Act) \rightsquigarrow \text{rew}(Rew)^*; \\ d_9^b(Act, Rew): \text{act}(Act) \rightsquigarrow \text{rew}(Rew)^* \Rightarrow R^*(Act); \\ d_E(\text{after}, \text{sati}): \text{mis}(\text{after}) \Rightarrow \neg d_7(\text{sati}); \\ d_F(\text{sati}, \text{hvn}): \text{True} \Rightarrow d_7(\text{sati}) < d_{10}^b; \\ d_8^b(Act_1, Act_2): sim(Act_1, Act_2), F(Act_1) \Rightarrow F(Act_2); \\ d_8^a(Act_1, Act_2, Eff): cs(Act_1, Eff), cs(Act_2, Eff), R(Act_1), R(Act_2) \\ \quad \Rightarrow sim(Act_1, Act_2). \end{array} \right\}$
3. $\mathcal{K}_p = \{ cs(\text{sati}, s.vio); cs(sye, vio); eff(s.vio, vio); F^*(vio); R(sye); \text{mis}(\text{after}) \}$
4. $\mathcal{K}_n = \{ \text{txt}O^*(\text{sati}); \text{txt_Act_Rew}^*(\text{sati}, \text{hvn}); \text{True} \}$
5. $F(X) = -O(X), \quad F^*(X) = -O^*(X), \quad R^*(X) = -O^*(X)$
6. $\mathcal{P}(\text{conc}(F)) = B' \prec C$

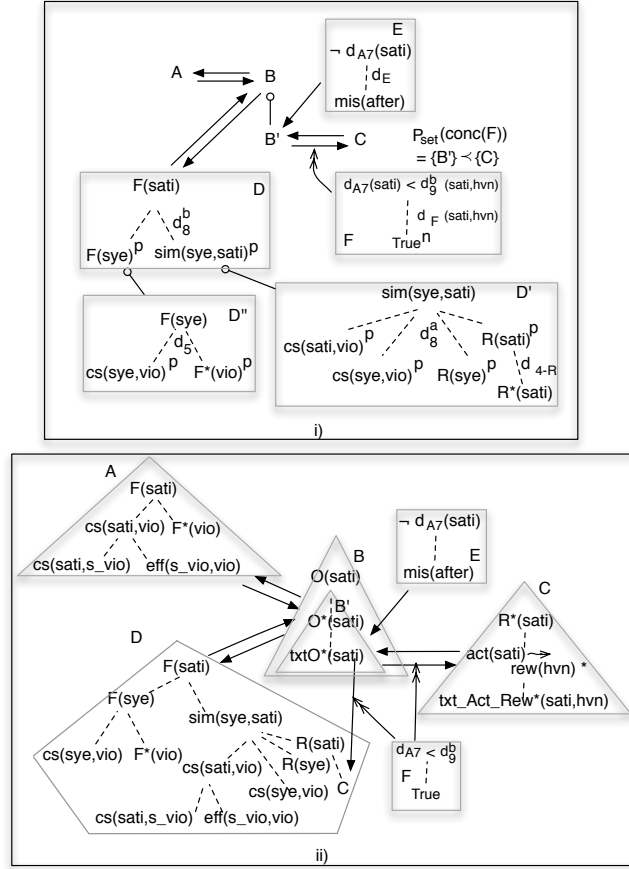


Fig. 4. i) Authored *ASPIC*⁺ network generated from Medhātithi’s analysis of *sati* ii) Some of the arguments and attacks in *EAF* constructed from mapping of i) to an *ASPIC*⁺ theory.

The final argument network is presented in Fig.4-i and is mapped to the argumentation theory of Def.13, subsequently instantiating a *bh-EAF* in Fig.4-ii (some of whose arguments are shown). In the network, both $R^*(sati)$ and $cs(sye, vio)$, are effectively incorporated in D' as lemmas, since arguments justifying these claims are included elsewhere in the network. Consequently, in the *bh-EAF* in Figure 4-ii, the argument C concluding $R^*(sati)$ is also shown as a sub-argument of D , and B also attacks D on C , where this attack is itself pref-attacked by F . Evaluating the *bh-EAF*, we obtain a single grounded, preferred and stable extension containing $A, D, E, C, F, [txtO^*(sati)]$, and their sub-arguments. In line with Medhātithi’s conclusion, we thus obtain justified arguments in favour of prohibiting *sati*, while also keeping it as a recommendation.

Concluding remark. The above case study highlights the advantages of providing formal argumentative support for scholars: helping to reveal and clarify

the structure of the dialectical commentaries being studied, as well as disclosing implicitly used assumptions and rendering these explicit for further analysis (including assumptions as to why some arguments are preferred to others). It also testifies to the utility, and hence promising future developments, of computational tools enabling the authoring of networks, their mapping to *ASPIC*⁺ theories, and evaluation of instantiated *EAF*s. While in this work the authoring, mapping and evaluation was done by hand, our future aim is to provide automated support for each step.

References

1. Amgoud, L., Cayrol, C., Lagasquie-Schiex, M., Livet, P.: On bipolarity in argumentation frameworks. *International Journal of Intelligent Systems*, vol.23(10), pp.1062–1093 (2008)
2. Baroni P., Caminada M., Giacomin M., An introduction to argumentation semantics. *The Knowledge Engineering Review*, 26 (4), pp.365–410 (2011)
3. Bondarenko, A., Dung, P.M., Kowalski, R.A., Toni, F.: An abstract, argumentation-theoretic approach to default reasoning. *Artificial Intelligence* 93, pp.63–101 (1997)
4. Brick, D.: The Dharmasāstric Debate on Widow-Burning. *Journal of the American Oriental Society* vol.130 (2), pp.203–223 (2010)
5. Caminada, M., Amgoud, L.: On the evaluation of argumentation formalisms. *Artificial Intelligence*, 171(5-6), pp.286–310 (2007)
6. Ciabattoni, A., Freschi, E., Genco, F.A., Lellmann, B.: Mīmāṃsā deontic logic: proof theory and applications. In: *TABLEAUX 2015*, vol.9323. Springer, pp.323–338 (2015)
7. Dung, P.M.: On the acceptability of arguments and its fundamental role in non-monotonic reasoning, logic programming and n-person games. *Artificial Intelligence*, vol.77(2), pp.321–358 (1995)
8. Freschi, E., Ollett, A., Pascucci, M.: Duty and Sacrifice. A Logical Analysis of the Mīmāṃsā Theory of Vedic Injunctions. *History and Philosophy of Logic*, (forthcoming) (2019)
9. Hunter, A., Modgil, S., Prakken, H., Simari, G., Besnard, P., Garcia, A., Toni, F.: Tutorials on structured argumentation. *Argument and Computation*, vol.5(1) (2014)
10. Modgil, S: Reasoning about preferences in argumentation frameworks. *Artificial Intelligence*, 173(9-10), pp.901–934 (2009)
11. Modgil, S: Revisiting abstract argumentation frameworks. In: *TAFa 2013*. Springer Berlin , pp.1–15 (2014)
12. Modgil, S: Towards a General Framework for Dialogues That Accommodate Reasoning About Preferences. In: *TAFa 2017*. Springer, pp.175–191 (2018)
13. Modgil, S., Prakken, H.: Reasoning about preferences in structured extended argumentation frameworks. In *Proc. COMMA 2010*, pp.347–358 (2010)
14. Modgil, S., Prakken, H.: A general account of argumentation with preferences. *Artificial Intelligence*, 195(0), pp.361–397 (2013)
15. Nielsen, S.H., Parsons, S.: A generalization of dung’s abstract framework for argumentation: Arguing with sets of attacking arguments. In: *ArgMAS 2006*. Springer, vol.4766, pp.54–73 (2007)
16. Oren, N., Norman, T.J.: Semantics for evidence-based argumentation. In: *Computational Models of Argument*, (COMMA 2008). IOS press, pp.276–284 (2008)
17. Prakken, H.: Historical overview of formal argumentation. In: *Handbook of Formal Argumentation*. College Publications, pp.75–144 (2018)
18. Sakuntala, N.: *Sati, widow burning in India*. Doubleday, New Delhi: Viking (1992)