# On the Justification of Statements in Argumentation-Based Reasoning

**Pietro Baroni**
DII, Univ. of Brescia, Italy
pietro.baroni@unibs.it

**Guido Governatori, Ho-Pun Lam, Régis Riveret**
DATA61, CSIRO, Brisbane, Australia
{guido.governatori,brian.lam,regis.riveret}@data61.csiro.au

## Abstract

In the study of argumentation-based reasoning, argument justification has received far more attention than statement justification, often treated as a simple byproduct of the former. As a consequence, counterintuitive results and significant losses of sensitivity can be identified in the treatment of statement justification by otherwise appealing formalisms. To overcome this limitation, we propose to reappraise statement justification as a formalism-independent component. To this purpose, we introduce a novel general model of argumentation-based reasoning based on multiple levels of labellings, one of which is devoted to statement justification. This model is able to encompass several literature proposals as special cases: we illustrate this ability for the case of the *ASPIC*$^+$ formalism and provide a first example of tunable statement justification in this context.

## Introduction

Suppose Dr. Smith, considered an expert, says you: "Given your clinical picture, you are affected by disease D1, not disease D2". Suppose then Dr. Jones, equally renowned, says you: "Given your clinical picture, you are affected by disease D2, not by disease D1". Your attitude about the justification of the statements S1="I am affected by disease D1" and S2="I am affected by disease D2" would become rather dubious. Suppose, in a different situation, you use an off-the-shelf test kit at home whose outcome suggests you are affected by disease D3. Then you undertake a serious and reliable clinical test whose outcomes excludes that you are affected by disease D3. Would you say that, in this case, your attitude about the justification of the statement S3="I am affected by disease D3" is the same as the attitude concerning S1 and S2 in the previous case? And what about the attitude towards the justification of the statement S4="I am affected by D4", where D4 is a poorly studied and initially asymptomatic disease you (and possibly your doctors too) never heard of? Is it in any way similar to the attitudes towards S1, S2, and S3 in the cases above? Intuitively, it seems reasonable to say that all these attitudes are different and that it is useful (besides being apparently easy) to distinguish them. Encompassing this kind of distinctions could then be considered a minimal

requirement for a knowledge representation and reasoning formalism.

Surprisingly, the current versions of several well-known structured argumentation formalisms fail to satisfy this simple requirement, equating, for instance, the justification status of S4 with the one of S3, or with that of S1 and S2, or even the justification status of S3 with that of S1 and S2.

While this may appear a severe drawback, we argue that this is not due to an intrinsic limitation of the argumentation formalisms themselves, rather to the relatively limited attention paid to the notion of justification of statements, often treated as a mere appendix of the notions of acceptance and justification of arguments, that (not surprisingly) are among the main focuses in formal argumentation studies.

In order to overcome this limitation, we suggest that the issue of statement justification, in the context of argumentation-based reasoning, can be a subject of analysis on its own, where general, formalism-independent, principles and properties can be investigated, to be then applied uniformly across different specific formalisms.

This paper makes some initial steps in this research direction by introducing a generic labelling-based model of argumentation-based reasoning process, where the notions of argument acceptance, argument justification, and statement justification are clearly distinguished and defined in a formalism-independent way, paving the way towards *tunable* statement justification.

## A Multi-Level Labelling System

We investigate the different notions of justification involved in a generic argument-based reasoning process. The process model we adopt consists of the following levels: argument production, argument acceptance, argument justification and statement justification.

**Argument production.** The first level regards the production of a set of arguments $\mathcal{A}$ whose structure and mutual relationships are left unspecified. The only relevant property for our purposes is that each argument $A \in \mathcal{A}$ has a conclusion, denoted as $Con(A)$, belonging to a language $\mathcal{L}$. Intuitively, an argument, whatever it actually is, provides a support for regarding its conclusion as justified. We do not make any assumption on the set of arguments, while we assume that the language is equipped with a contrariness relation. In its

simplest form the contrariness relation corresponds to the traditional notion of negation but other more general forms of contrariness have been considered in the literature. To encompass this more general view, we assume a contrariness relation $Cnt$, allowing the existence of multiple (or no, if $Cnt(\varphi) = \emptyset$) contraries for each statement $\varphi$ of the language.

**Definition 1** (Language). *A language $\mathcal{L}$ is a set of statements equipped with a contrariness relation $Cnt : \mathcal{L} \rightarrow 2^{\mathcal{L}}$. For all $\varphi \in \mathcal{L}$, $\psi \in Cnt(\varphi)$ is called a contrary of $\varphi$.*

**Definition 2** (Argument-conclusion structure). *An argument-conclusion structure is a triple $\langle \mathcal{L}, \mathcal{A}, Con \rangle$ where $\mathcal{L}$ is a language, $\mathcal{A}$ is a finite set of arguments and $Con : \mathcal{A} \rightarrow \mathcal{L}$ is a relation associating any argument with its conclusion.*

Note that some elements of $\mathcal{L}$ may not play the role of conclusions, e.g. if $\mathcal{L}$ encompasses negation as failure.

**Argument acceptance.** The second level concerns the acceptance evaluation of a set of arguments, the outcome is a set of argument acceptance labellings. An argument acceptance labelling $L_A$ assigns to each argument an acceptance label taken from a set of labels $\Lambda_A$. Each label in $\Lambda_A$ represents an individual acceptance status and each labelling represents a "reasonable" point of view about the acceptance of the arguments belonging to $\mathcal{A}$.

**Definition 3** ($\Lambda_A$-based argument acceptance labelling). *Given an argument-conclusion structure $\mathcal{AC} = \langle \mathcal{L}, \mathcal{A}, Con \rangle$ and a set of acceptance labels $\Lambda_A$, a $\Lambda_A$-based argument acceptance labelling for $\mathcal{AC}$ is a function $L_A : \mathcal{A} \rightarrow \Lambda_A$.*

Definition 3 only specifies what a generic labelling is, independently of the way a labelling can be generated and of the properties that a labelling should satisfy. Indeed, in general, not every labelling shall appear reasonable. For instance, not all arguments can be accepted at the same time because of some relationships (e.g. of attack) holding among them, which are abstracted away in the representation we are considering. Hence, an acceptance criterion (or mechanism) is needed to select those acceptance labellings which are reasonable, i.e are compatible with the underlying constraints corresponding to the relationships among arguments. We leave this acceptance criterion/mechanism unspecified at this level and to abbreviate the presentation, we will simply use the symbol $\mathfrak{L}_A(\mathcal{A})$ to denote the set of labellings specified by an acceptance criterion.

In general, there are many reasonable labellings that a given acceptance criterion can select. It is also possible however to define an acceptance criterion such that its outcome always consists in exactly one labelling: we say that an acceptance labelling criterion is *single-status* if, for every $\mathcal{A}$, $| \mathfrak{L}_A(\mathcal{A}) | = 1$, *multiple-status* if, for some $\mathcal{A}$, $| \mathfrak{L}_A(\mathcal{A}) | > 1$.

**Argument justification.** The third level deals with the assignment of a synthetic justification status to each argument. We assume that this is represented by an argument justification labelling $L_J$, i.e. a function from the set of arguments $\mathcal{A}$ to a set of justification labels $\Lambda_J$. It is rather natural to assume that $L_J$ is functionally dependent on $\mathfrak{L}_A(\mathcal{A})$. As a starting point, we make two assumptions on the nature of this dependency: first, for each argument $A$, $L_J(A)$ depends only

on the acceptance labels of $A$ in $\mathfrak{L}_A(\mathcal{A})$; second, cardinality does not count in this evaluation, i.e. for each label $\lambda \in \Lambda_A$ it only matters whether there are some elements of $\mathfrak{L}_A(\mathcal{A})$ such that $L_A(A) = \lambda$. Following these assumptions, we can identify two steps in the definition of $L_J$. First, for each argument $A$ the acceptance evaluation outcome is synthesised by a set $\Sigma_A(A) \subseteq \Lambda_A$ of acceptance labels associated to $A$ by the members of $\mathfrak{L}_A(\mathcal{A})$. The second step consists in associating to each possible value of $\Sigma_A(A)$, i.e. to every subset of $\Lambda_A$ one of the labels in $\Lambda_J$, i.e. to define a synthesis function $\mathfrak{S}_J : 2^{\Lambda_A} \rightarrow \Lambda_J$. In this way we get that $L_J(A) = \mathfrak{S}_J(\Sigma_A(A))$. It can be noted that the acceptance criterion can be defined so as to ensure that $\Sigma_A(A) \neq \emptyset$, in this case $\mathfrak{S}_J(\emptyset)$ can be left undefined. Moreover in the special case where the acceptance criterion is single-status, the only possible values of $\Sigma_A(A)$ are singletons, i.e. $\Sigma_A(A) = \{\lambda\}$ for some $\lambda \in \Lambda_A$. In this setting one can virtually skip this third level by putting $\Lambda_J = \Lambda_A$ and letting, for every argument $A$, $L_J(A) = \lambda$.

**Definition 4** (Projection). *Given an argument-conclusion structure $\mathcal{AC} = \langle \mathcal{L}, \mathcal{A}, Con \rangle$ and a set of acceptance labellings $\mathfrak{L}_A(\mathcal{A})$, the projection of $\mathfrak{L}_A(\mathcal{A})$ on every argument $A \in \mathcal{A}$ is defined as*

$$\Sigma_A(A) = \{\lambda \in \Lambda_A \mid \exists L_A \in \mathfrak{L}_A(\mathcal{A}) : L_A(A) = \lambda\}.$$

**Definition 5** ($\Lambda_J$-based justification labelling). *Given a set of justification labels $\Lambda_J$, a $\Lambda_J$-based justification labelling for $\mathcal{A}$ is a function $L_J : \mathcal{A} \rightarrow \Lambda_J$. A justification labelling is called cardinality insensitive if there is a function $\mathfrak{S}_J : 2^{\Lambda_A} \rightarrow \Lambda_J$ such that for every argument $A$ it holds that $L_J(A) = \mathfrak{S}_J(\Sigma_A(A))$.*

**Example 1.** ASPIC$^+$ *(denoted as $A^+$ for short) is a rule-based argumentation formalism which assumes the existence of a generic language $\mathcal{L}$ equipped with a contrariness relation (Modgil and Prakken 2014).*

*$A^+$ arguments may attack each other, and argument acceptance is based on Dung's formalism of argumentation frameworks (Dung 1995) and its semantics. Accordingly, it is possible to refer to the labelling-based version of Dung's semantics (Baroni, Caminada, and Giacomin 2011), where a set of three argument acceptance labels is adopted, namely $\Lambda_A = \{\text{IN}, \text{OUT}, \text{UN}\}$.*

*Concerning the subsequent level of argument justification, $A^+$ adopts the traditional notion of skeptical and credulous justification which says that an argument is skeptically justified (denoted SKJ) if it is labelled IN in all labellings prescribed by the adopted semantics, while it is credulously justified (denoted CRJ) if it is labelled IN in some labellings (but not all, in order to keep these notions disjoint).*

**Proposition 1.** *Given the set of argument justification labels $\Lambda_J^{A^+} = \{\text{SKJ}, \text{CRJ}\}$, the argument justification labelling $L_J^{A^+}$ prescribed by $A^+$ is such that for any argument $A$,*

- $L_J^{A^+}(A) = \text{SKJ}$ *iff $\Sigma_A(A) = \{\text{IN}\}$;*
- $L_J^{A^+}(A) = \text{CRJ}$ *iff $\Sigma_A(A) \supsetneq \{\text{IN}\}$.* □

**Statement justification.** The fourth level caters for the justification status of statements, i.e. the elements of the language $\mathcal{L}$. We assume that this is represented by a statement

justification labelling $L_S$, i.e. a function from $\mathcal{L}$ to a set of statement justification labels $\Lambda_S$. One may then assume that $L_S$ depends on the argument justification labelling $L_J$.

First of all it can be observed that, in general each statement $\varphi \in \mathcal{L}$ is supported by a (possibly empty) set of arguments $Arg(\varphi)$ and, similarly, the contraries of a statement are supported by a set of arguments $CntArg(\varphi)$. Following the assumption of cardinality insensitivity the justification labellings of these sets of argument can be synthesized by the sets of labels they include.

**Definition 6** (Supporting arguments)**.** *Given an argument-conclusion structure $\langle \mathcal{L}, \mathcal{A}, Con \rangle$ and a set of statements $\Phi \subseteq \mathcal{L}$, the set of supporting arguments of $\Phi$ is defined as*

$$Arg(\Phi) = \{A \in \mathcal{A} \mid Con(A) \in \Phi\}.$$

**Definition 7** (Synthetic justification)**.** *Given an argument-conclusion structure $\langle \mathcal{L}, \mathcal{A}, Con \rangle$ and a $\Lambda_J$-based justification labelling for $\mathcal{A}$, the synthetic justification of the supporting arguments for $\varphi$ is defined such that for all $\varphi \in \mathcal{L}$:*

$$\Sigma_J(\varphi) \triangleq \{\lambda \in \Lambda_J \mid \exists A \in Arg(\{\varphi\}) : L_J(A) = \lambda\}$$

*and similarly the synthetic justification of the contrary-supporting arguments of $\varphi$ is defined such that for all $\varphi \in \mathcal{L}$:*

$$\Sigma_{\overline{J}}(\varphi) \triangleq \{\lambda \in \Lambda_J \mid \exists A \in Arg(Cnt(\varphi)) : L_J(A) = \lambda\}.$$

We assume then that for every statement $\varphi$, $L_S(\varphi)$ depends on the argument justification labels of the elements of $Arg(\varphi)$ and $CntArg(\varphi)$. Adopting the assumption that cardinality does not matter, we are led to identify two steps in the definition of $L_S$. In the first step, for each statement $\varphi$ the justification status of the relevant arguments is synthesised by the sets $\Sigma_J(\varphi)$ and $\Sigma_{\overline{J}}(\varphi)$. The second step consists in associating to each pair of subsets of $\Lambda_J$ one of the labels in $\Lambda_S$, i.e. to define a synthesis function $\mathfrak{S}_S : 2^{\Lambda_J} \times 2^{\Lambda_J} \to \Lambda_S$.

**Definition 8** ($\Lambda_S$-based statement justification labelling)**.** *Given an argument-conclusion structure $\mathcal{AC} = \langle \mathcal{L}, \mathcal{A}, Con \rangle$ and a set of statement justification labels $\Lambda_S$, a $\Lambda_S$-based statement justification labelling for $\mathcal{AC}$ is a function $L_S : \mathcal{L} \to \Lambda_S$. Assuming that $\mathcal{AC}$ is equipped with a $\Lambda_J$-based justification labelling for $\mathcal{A}$, we say that a statement justification labelling $L_S$ is argument aware and cardinality insensitive if there is a function $\mathfrak{S}_S : 2^{\Lambda_J} \times 2^{\Lambda_J} \to \Lambda_S$ such that for every statement $\varphi \in \mathcal{L}$, $L_S(\varphi) = \mathfrak{S}_S(\Sigma_J(\varphi), \Sigma_{\overline{J}}(\varphi))$.*

**Example 2.** *In $\mathsf{A}^+$, statements inherit directly the justification status of the "best justified" argument supporting them (see Def. 3.17 of (Modgil and Prakken 2014): a statement is skeptically justified if and only if it is the conclusion of a skeptically justified argument, while it is credulously justified if and only if it is not skeptically justified and it is the conclusion of a credulously justified argument. The case where a statement is neither skeptically nor credulously justified is not explicitly covered.*

**Proposition 2.** *Given the set of statement justification labels $\Lambda_S^{\mathsf{A}^+} = \{\mathsf{skj}, \mathsf{crj}\}$, the statement justification labelling $L_S^{\mathsf{A}^+}$ prescribed by $\mathsf{A}^+$ is such that for any statement $\varphi$,*

- $L_S^{\mathsf{A}^+}(\varphi) = \mathsf{skj}$ *iff* $\mathsf{SKJ} \in \Sigma_J(\varphi)$;
- $L_S^{\mathsf{A}^+}(\varphi) = \mathsf{crj}$ *iff* $\mathsf{SKJ} \notin \Sigma_J(\varphi)$ *and* $\mathsf{CRJ} \in \Sigma_J(\varphi)$.

*Let us illustrate the statements labellings with our introductory example. We can assume that there are two mutually attacking arguments supporting the statements S1 and S2, that the argument supporting the statement S3 is defeated by another (stronger) argument supporting the negation of S3 (denoted ¬S3), and that there are no arguments supporting S4, nor its negation.*

*According to $\mathsf{A}^+$, with every semantics, S3 and S4 get an undefined justification status (we may mark this undefined justification status as 'noj'), while ¬S3 would be skj. The status of S1 and S2 is semantics-dependent: both would get the status crj if a Dung multiple-status semantics (e.g. preferred or stable) is adopted, while they would be equated to S3 and S4 (undefined or noj) in the case of a Dung single-status semantics (e.g. grounded or ideal).* □

An argument-conclusion structure fully equipped with the labellings introduced above will be called a *multi-level labelling system*.

**Definition 9** (Multi-level labelling system)**.** *A multi-level labelling system is a tuple $\langle \mathcal{AC}, L_A, L_J, L_S \rangle$ where*

- *$\mathcal{AC}$ is an argument-conclusion structure,*
- *$L_A$ is a $\Lambda_A$-based argument acceptance labelling for $\mathcal{AC}$,*
- *$L_J$ is a $\Lambda_J$-based argument justification labelling for $\mathcal{AC}$,*
- *$L_S$ is a $\Lambda_S$-based statement justification labelling for $\mathcal{AC}$.*

The model defined above is useful to analyse and compare actual argumentation formalisms on a common ground consisting of abstract general properties. For example, we may consider the notions of full coverage and insensitivity to contrariness.

As to the first property, it simply amounts to require that the relevant functions are total.

**Definition 10** (Coverage)**.** *A multi-level labelling system $\mathfrak{L}^* = \langle \mathcal{AC}, L_A, L_J, L_S \rangle$ is said to provide*

- *a full coverage of argument acceptance if $L_A$ is total,*
- *a full coverage of argument justification if $L_J$ is total, and*
- *a full coverage of statement justification if $L_S$ is total.*

*$\mathfrak{L}^*$ provides an exhaustive justification coverage if it provides all the three levels of full coverage introduced above.*

**Example 3.** *It can be immediately observed that the $\mathsf{A}^+$ argument justification labelling (see Prop. 1) does not provide full coverage, since it does not cover the cases where $\mathsf{IN} \notin \Sigma_A(A)$. This can be explained by the emphasis on acceptance in $\mathsf{A}^+$. It is anyway easy to recover a full coverage by defining a third label (let say not justified, denoted as $\mathsf{NOJ}$), covering the remaining cases, i.e. letting $\Lambda_J^{\mathsf{A}^+} = \{\mathsf{SKJ}, \mathsf{CRJ}, \mathsf{NOJ}\}$.* □

Sensitivity to contrariness concerns statement justification only: the idea is that the justification status of a statement $\varphi$ is actually somehow affected also by the status of the arguments supporting its contraries. Formally this amounts to require that they make some difference in the evaluation.

**Definition 11** (Contrary-sensitivity)**.** *Given a multi-level labelling system $\mathfrak{L}^* = \langle \mathcal{AC}, L_A, L_J, L_S \rangle$ we say that $L_S$ is contrary-sensitive iff $\exists \varphi, \psi \in \mathcal{L}$ such that $\Sigma_J(\varphi) = \Sigma_J(\psi)$, $\Sigma_{\overline{J}}(\varphi) \neq \Sigma_{\overline{J}}(\psi)$, and $L_S(\varphi) \neq L_S(\psi)$.*

**Example 4.** $A^+$ *is not contrary-sensitive, since* $\Sigma_{\overline{J}}(\varphi)$ *does not play any role in the definition of* $L_S(A)$. *This can be explained by the focus on positive support in* $A^+$. *Hence these limitations are certainly not intrinsic to* $A^+$, *rather they can be overcome by providing more articulated definitions for the notions of justification, leaving unchanged all the rest of the formalism, as we will see next.*  □

## Towards tunable justification notions

Different argumentation formalisms, such as $A^+$ (Modgil and Prakken 2014), ABA (Toni 2014), DeLP (Garcia and Simari 2014) or DL (Governatori et al. 2004), adopt quite different notions of justification, both at the level of arguments and of statements, featuring different properties and sometimes failing to satisfy some intuitive requirements like full coverage and contrary-sensitivity. However these differences do not seem to be caused by technical motivations, but rather to depend on arbitrary choices based on the intended use of the notion of justification in the presentation of the formalisms themselves. These observations back up our claim that the notion of justification (and in particular of statement justification) has been somehow neglected in the development of argumentation formalisms, often more focused on the notion of argument acceptance. Moreover they suggest that justification notions, instead of being "hardwired" in the definitions could better be conceived as tunable components of any argumentation formalism, with a role similar to those played by argumentation semantics. These formalisms do not stick to a single argumentation semantics, rather they assume that one is chosen among the various available ones (including possibly those to be developed in the future).

**Example 5.** *In replacement of the* $A^+$ *statement justification labelling (Prop. 2), we may consider an 'ignorance-aware' labelling such that a statement is labelled (i)* in *iff there is a supporting justified argument for it, (ii)* fal *iff there is a supporting justified argument for its contrary, (iii)* unk *iff there is no supporting justified arguments for it or its contrary, and (iv)* ni *otherwise. This labelling shall occur in a skeptical and credulous mode.*

**Definition 12.** *The skeptical ignorance-aware labelling for* $A^+$ *is defined as follows*

- $L_S^{\mathrm{sklaA}^+}(\varphi) =$ yes *iff* $\mathrm{SKJ} \in \Sigma_J(\varphi)$;
- $L_S^{\mathrm{sklaA}^+}(\varphi) =$ fal *iff* $\mathrm{SKJ} \in \Sigma_{\overline{J}}(\varphi)$;
- $L_S^{\mathrm{sklaA}^+}(\varphi) =$ unk *iff* $\Sigma_{\overline{J}}(\varphi) \cup \Sigma_J(\varphi) = \emptyset$;
- $L_S^{\mathrm{sklaA}^+}(\varphi) =$ ni *otherwise.*

**Definition 13.** *The credulous ignorance-aware labelling for* $A^+$ *is defined as follows*

- $L_S^{\mathrm{crlaA}^+}(\varphi) =$ yes *iff* $\{\mathrm{SKJ}, \mathrm{CRJ}\} \cap \Sigma_J(\varphi) \neq \emptyset$;
- $L_S^{\mathrm{crlaA}^+}(\varphi) =$ fal *iff* $\{\mathrm{SKJ}, \mathrm{CRJ}\} \cap \Sigma_J(\varphi) = \emptyset$ *and* $\{\mathrm{SKJ}, \mathrm{CRJ}\} \cap \Sigma_{\overline{J}}(\varphi) \neq \emptyset$;
- $L_S^{\mathrm{crlaA}^+}(\varphi) =$ unk *iff* $\Sigma_{\overline{J}}(\varphi) \cup \Sigma_J(\varphi) = \emptyset$;
- $L_S^{\mathrm{crlaA}^+}(\varphi) =$ ni *otherwise.*

*In contrast to the original* $A^+$ *statement labelling, we can see that these ignorance-aware labellings are contrary-sensitive.*

*The capability to capture more distinctions is evident in the example: in the skeptical labelling, S1 and S2 are labelled as* ni, *S3 as* fal, *¬S3 as* yes, *S4 as* unk, *while in the credulous case S1 and S2 are labelled as* yes, *S3 as* fal, *¬S3 as* yes, *S4 as* unk.  □

## Conclusion

Even referring to simple common reasoning examples, some argumentation formalisms provide counterintuive results concerning the justification status of statements, e.g. by assigning the same label to statements whose status appears to be intuitively different. Moreover, using the same examples, it can be shown (this is not covered in this paper due to space limitations) that different argumentation formalisms disagree on the status of some statements. We argue that these disagreements are due to the relatively limited attention paid to statement justification rather than to inherent substantial differences between the formalisms themselves. We therefore propose a novel multi-level labelling model for argument-based reasoning which regards statement justification as a formalism-independent component of the process and promotes the idea that it is tunable, much in the way argumentation semantics is a tunable component in several formalisms.

Overall, our abstract multi-level labelling system provides a first foundational contribution towards a deeper study of statement justification in argumentation-based reasoning and opens the way to several future research directions. In particular, this work is a basis for a systematic study of general principles and properties for statement labellings.

## Acknowledgements

## References

Baroni, P.; Caminada, M.; and Giacomin, M. 2011. An introduction to argumentation semantics. *Knowledge Eng. Review* 26(4):365–410.

Dung, P. M. 1995. On the acceptability of arguments and its fundamental role in nonmonotonic reasoning, logic programming and n-person games. *Artificial Intelligence* 77(2):321–358.

Garcia, A. J., and Simari, G. R. 2014. Defeasible logic programming: DeLP servers, contextual queries, and explanations for answers. *Argument & Computation* 5(1):63 – 88.

Governatori, G.; Maher, M. J.; Antoniou, G.; and Billington, D. 2004. Argumentation semantics for defeasible logic. *J. Log. Comput.* 14(5):675–702.

Modgil, S., and Prakken, H. 2013. A general account of argumentation with preferences. *Artif. Intell.* 195:361 – 397.

Modgil, S., and Prakken, H. 2014. The aspic$^+$ framework for structured argumentation: a tutorial. *Argument & Computation* 5(1):31 – 62.

Prakken, H. 2010. An abstract framework for argumentation with structured arguments. *Argument & Computation* 1(2):93–124.

Toni, F. 2014. A tutorial on assumption-based argumentation. *Argument & Computation* 5(1):89 – 117.