

Weighted Assumption Based Argumentation to reason about ethical principles and actions

Paolo Baldi¹, Fabio Aurelio D'Asaro^{1,*}, Abeer Dyoub² and Francesca A. Lisi^{2,3}

¹Dept. of Human Studies, University of Salento, Lecce, Italy

²Dept. of Informatics, University of Bari "Aldo Moro", Via E. Orabona 4, Bari, 70125, Italy

³Centro Interdipartimentale di Logica e Applicazioni (CILA), University of Bari "Aldo Moro", Via E. Orabona 4, Bari, 70125, Italy

Abstract

We augment Assumption Based Argumentation (ABA for short) with weighted argumentation. In a nutshell, we assign weights to arguments and then derive the weight of attacks between ABA arguments. We illustrate our proposal through running examples in the field of ethical reasoning, and present an implementation based on Answer Set Programming.

Keywords

Formal Argumentation, Assumption Based Argumentation, Ethical Reasoning, Fuzzy Logic

1. Introduction

Formal argumentation frameworks model reasoning with conflicting claims, based on the crucial notion of *attacks* among arguments. These attacks are rendered as directed edges connecting nodes in a graph, while the arguments themselves are represented simply as nodes of the graph, see e.g. the seminal paper by Dung [1]. Subsequent approaches, belonging to the field of *structured argumentation*, see e.g. [2], provide a more fine-grained representation of the argument nodes, equipping them with a logical structure, and using such structure for deriving from logical principles the occurrence of attacks among arguments.

Assumption-based argumentation (see, e.g., [2, Chapter 7]) is a prominent approach to structured argumentation. It represents arguments as logical derivations, built on the basis of two ingredients: *rules*, which are considered to be non-defeasible, and *assumptions*, which are taken instead to be the defeasible part of the argument, and possibly the target of attacks.

In this work, we introduce *Weighted Assumption Based Argumentation frameworks* (wABAs), which enrich ABA with *weighted* arguments. These weights on arguments determine in turn weights on the attacks among arguments, on the model of weighted abstract argumentation, see, e.g., [3, Chapter 6]. This has several modeling advantages. On the one hand, the introduction of weights allows us to import ideas from fuzzy logic. On the other hand, since weighted arguments translate into weighted attacks, we may then allow for certain forms of incoherence in the semantics, making use of standard techniques in weighted abstract argumentation.

We propose an implementation of wABA using a translation into Answer Set Programming (ASP) by means of `clingo` answer set grounder and solver, and demonstrate the formalism and its implementation on a scenario involving AI ethics.

Our motivations for the introduction and implementation of wABA originates indeed from the general aim of addressing computational reasoning with ethical principles in AI, and more specifically in the

CILC 2025: 40th Italian Conference on Computational Logic, June 25–27, 2025, Alghero, Italy

*Corresponding author.

✉ paolo.baldi@unisalento.it (P. Baldi); fabioaurelio.dasaro@unisalento.it (F. A. D'Asaro); abeer.dyoub@uniba.it (A. Dyoub); francesca.alessandra.Lisi@uniba.it (F. A. Lisi)

🌐 <https://sites.google.com/view/paolobaldi> (P. Baldi); <https://sites.google.com/view/fdasaro> (F. A. D'Asaro);

<https://www.abeerdyoub.com> (A. Dyoub); <https://www.uniba.it/it/docenti/lisi-francesca-alessandra> (F. A. Lisi)

🆔 0000-0003-2657-753X (P. Baldi); 0000-0002-2958-3874 (F. A. D'Asaro); 0000-0003-0329-2419 (A. Dyoub);

0000-0001-5414-5844 (F. A. Lisi)



© 2025 Copyright for this paper by its authors. Use permitted under Creative Commons License Attribution 4.0 International (CC BY 4.0).

domain of *symbiotic human-AI interactions*. Deontological ethics, with its rule-based approach seems at first most well- suited for computational implementation in our setting. However, eliciting precise rules from general principles is far from obvious, and ethical rules can easily result in inconsistent outcomes.

We argue that, due to its peculiar features, wABA offers adequately expressive computational machinery both for the representation of ethical reasoning, and guidance on conflict resolution.

Concerning representation, our approach distinguishes, within wABA arguments, among (i) the general ethical principles, which play the role of the (defeasible) assumptions in the framework, (ii) the factual elements, which may be used as premises of arguments in addition to the assumptions, and (iii) the prescriptions of courses of actions, which appear as conclusions of arguments.

Arguments in wABA are thus well-suited to render *prima-facie duties* [4], i.e., duties based on ethical principles, that may lead to conflicting prescriptions, and may be overridden on the basis of application context.

Conflicts between ethically motivated prescriptions are naturally derived in wABA as attacks among the arguments. Following the *extension*-based semantics of formal argumentation, wABA then allows one to represent various possible solutions to the ethical conflicts, i.e., different horns of ethical dilemmas, as different extensions. In the spirit of symbiotic AI [5], we take it to be a crucial feature of our approach that the AI agent does not substitute itself to the human agent, but disentangles, rather than solve, the ethical conflicts.

At the same time, wABA also offers guidance to the resolution of conflicts, due to the use of weights. Our design choice here is to avoid any a priori weighting of ethical principles, rather allowing the assignment of weights only to formulas standing for factual aspects, i.e. for the assessment of the context of application. In this respect, we extend and incorporate the fuzzy rule-based system developed in [6] within our framework. Only secondarily, on the basis of suitable computations, weights are carried over to arguments, then to attacks, and finally to the extensions of the argumentation framework. The user is thus ultimately confronted with weighted extensions, i.e., weighted solutions to ethical dilemmas, and can thus evaluate how strong a violation of certain ethical assumption she is willing to accept.

The rest of the paper is structured as follows. In Section 2 we provide some background on abstract and assumption-based argumentation. In Section 3 we introduce wABA and in Section 4 we discuss our ASP implementation in `clingo`. Section 5 discusses the application of the framework to ethical reasoning and Section 6 the related work on computational approaches to ethical reasoning. Section 7 concludes the paper by wrapping up our contribution while hinting at future developments of the present work.

2. Background

In this section, we briefly recall the fundamental concepts of abstract argumentation frameworks (AAFs), assumption-based argumentation (ABA), and weighted abstract argumentation frameworks (wAAFs). These frameworks provide the machinery upon which our proposed approach is constructed.

2.1. Abstract Argumentation Frameworks

We begin with Dung’s abstract argumentation frameworks [1], which provide an abstract characterization of argumentation.

Definition 2.1 (Abstract Argumentation Framework). An *abstract argumentation framework* (AAF) is a pair (Arg, Att) where:

- Arg is a finite set of arguments;
- $Att \subseteq Arg \times Arg$ is a binary relation representing attacks between arguments.

Given an AAF (Arg, Att) , for any $a, b \in Arg$, the notation $(a, b) \in Att$ indicates that a attacks b . We say that a set of arguments $B \subseteq Arg$ *defends* an argument $a \in Arg$ if for each argument $b \in Arg$ that attacks a , there exists an argument $c \in B$ such that c attacks b . A subset $B \subseteq Arg$ is:

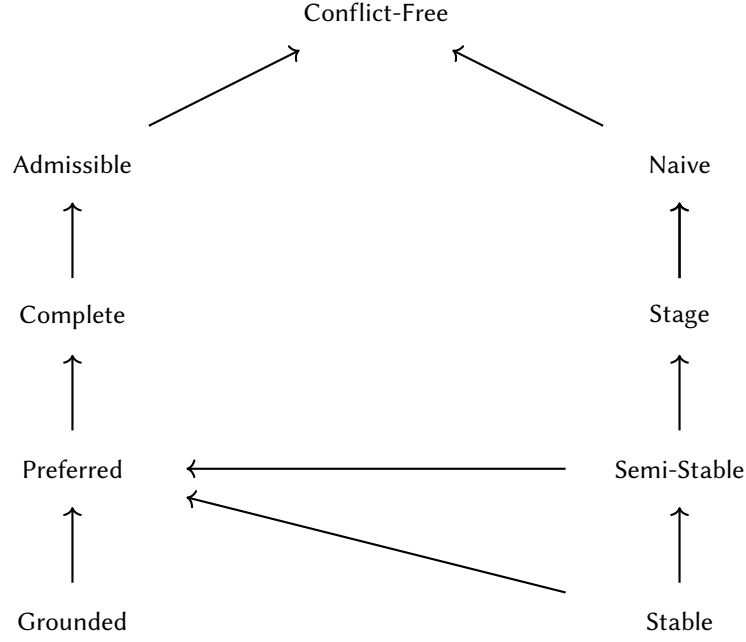


Figure 1: Set-theoretic inclusion relations among semantics. All arrows denote subset inclusion.

- *conflict-free* if there are no $a, b \in B$ such that $(a, b) \in Att$;
- *admissible* if it is conflict-free and defends all its elements;
- *preferred* if it is a maximal (w.r.t. set inclusion) admissible set;
- *grounded* if it is the least (w.r.t. set inclusion) admissible set;
- *stable* if it is conflict-free and attacks every argument not in the set.

In addition to the semantics based on admissibility, alternative semantics have been introduced to capture different intuitions, particularly when admissibility leads to overly restrictive or unintuitive outcomes. These include the following:

- *Naive* extensions are the maximal (w.r.t. set inclusion) conflict-free subsets of Arg . Unlike admissible extensions, naive extensions do not require defense against attacks, focusing instead on maximizing conflict-freeness alone.
- *Semi-stable* extensions are admissible sets whose range—i.e., the union of the set and the arguments it attacks—is maximal (w.r.t. set inclusion) among admissible sets. These extensions aim to approximate stable extensions when the latter do not exist.
- *Stage* extensions are conflict-free sets whose range is maximal among all conflict-free sets. Like semi-stable semantics, they emphasize coverage (via attack) of the argument space, but do not require admissibility.

These non-admissibility-based semantics are particularly relevant in weighted or resource-bounded settings, where defense may be impractical or where broader coverage is desirable despite some lack of coherence. In such cases, naive and stage extensions may yield more informative or robust outcomes than traditional admissibility-based semantics.

2.2. Assumption-Based Argumentation

Assumption-Based Argumentation (ABA; see, e.g., [7] for a primer) provides a structured argumentation framework grounded in a deductive system, based on rules and defeasible assumptions.

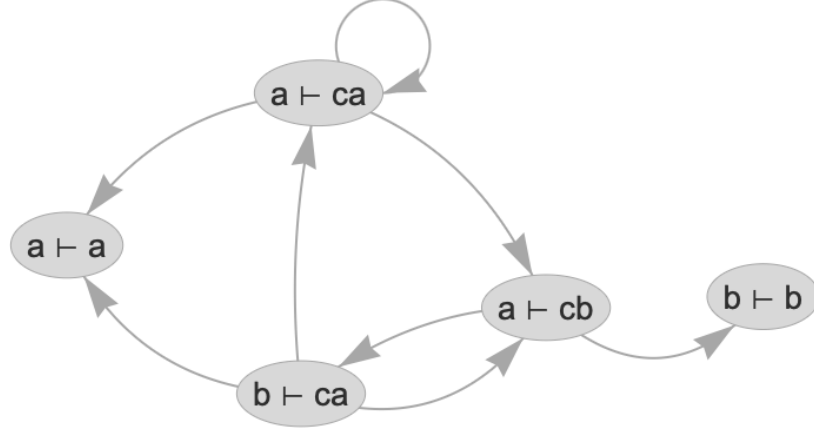


Figure 2: ABA framework visualization for Example 2.3.

Definition 2.2 (ABA Framework). An *assumption-based argumentation* (ABA) framework is a tuple $(\mathcal{L}, \mathcal{R}, \mathcal{A}, \neg)$ where:

- \mathcal{L} is a formal language;
- \mathcal{R} is a set of inference rules of the form $\phi \leftarrow \phi_1, \dots, \phi_n$ with $\phi, \phi_i \in \mathcal{L}$;
- $\mathcal{A} \subseteq \mathcal{L}$ is a non-empty set of assumptions;
- $\neg : \mathcal{A} \rightarrow \mathcal{L}$ is a total mapping that assigns to each assumption its contrary.

An ABA is said to be *flat* if and only if assumptions only appear in the body of rules.

Henceforth, for any rule $r \in \mathcal{R}$ of the form $\phi \leftarrow \phi_1, \dots, \phi_n \in \mathcal{R}$ we let $body(r) = \{\phi_1, \dots, \phi_n\}$ and $head(r) = \phi$. Arguments in ABA are deductions, denoted by $\Phi \vdash \psi$, where Φ is a set of assumptions and ψ is any formula in \mathcal{L} , obtained from Φ by applying one or more rules. Given an ABA argument x (i.e. a deduction), we denote by $rules(x)$ the set of rules *supporting* it. Attacks are defined via contraries. Specifically, we will have that an argument x attacks an argument y , if and only if x is a deduction of the form $\Phi' \vdash \bar{\psi}$, and y is a deduction of the form $\Phi, \psi \vdash \phi$, where $\psi \in \mathcal{A}$, $\phi \in \mathcal{L}$, $\bar{\psi}$ is the contrary of ψ and $\emptyset \subseteq \Phi$, $\Phi' \subseteq \mathcal{A}$. Any ABA framework thus naturally defines a corresponding abstract argumentation framework associated with it, from which extensions can be extracted.

Example 2.3. As a refresher, consider the flat ABA framework consisting of *atoms* a , b , ca , and cb , where a and b are *assumptions*, *contraries* are given by $\bar{a} = ca$ and $\bar{b} = cb$, and *rules* are $ca \leftarrow a$, $ca \leftarrow b$, $cb \leftarrow a$. Recall that in ABA an argument is a deduction. Examples of arguments in this ABA are $b \vdash b$, $a \vdash ca$, and $b \vdash ca$. Attacks are then derived by considering contraries: for example, $b \vdash ca$ attacks $a \vdash ca$ as the derived atom ca is the contrary of the assumption a . The full framework is shown in Figure 2, which was produced with the *PyArg* library [8]. This graph can be treated as a standard AAF. For instance, the only stable extension of this framework is $\{b \vdash ca, b \vdash b\}$.

2.3. Weighted Abstract Argumentation

Weighted extensions of abstract argumentation frameworks [9] assign numerical weights to arguments or attacks, enabling reasoning with preferences, strengths, or costs.

Definition 2.4 (Weighted AAF). A *weighted abstract argumentation framework* (wAAF) is a triple (Arg, Att, w) where:

- (Arg, Att) is an abstract argumentation framework;
- $w : Att \rightarrow \mathbb{R}^+$ is a function assigning a positive real-valued weight to each attack.

An *inconsistency budget* β is typically used to specify a degree of inconsistency one is willing to tolerate in a given scenario. In other words, attacks that weigh up to β may be discarded from the framework. Semantics are then defined with respect to the inconsistency budget β and a standard semantics σ of AAF, see, e.g., [10]. One usually refers to these semantics as β - σ extensions, e.g., 3-stable, 2-admissible, 5-grounded, etc. Note that whenever $\beta = 0$ we obtain the standard argumentation semantics σ , e.g., 0-stable is the standard AAF stable semantics, 0-grounded is the standard AAF grounded semantics, etc.

Example 2.5. Consider the weighted abstract argumentation framework (Arg, Att, w) where:

$$Arg = \{a, b, c\}, \quad Att = \{(a, b), (b, c), (c, a)\}, \quad w((a, b)) = 2, \quad w((b, c)) = 1, \quad w((c, a)) = 4.$$

This framework forms a directed cycle where each argument attacks one other. Under standard stable semantics (i.e., with $\beta = 0$ and $\sigma = \text{stable}$), there exists no extension, as each argument is attacked by another, and no conflict-free set can defend against all incoming attacks.

Now consider the framework under a budgeted stable semantics with inconsistency budget $\beta = 3$. In this case, we may discard attacks whose cumulative weight does not exceed 3. For instance, we may choose to discard the attack (c, a) , which has weight 4, but this alone would exceed the budget. However, if we discard (a, b) and (b, c) , whose combined weight is 3, we remain within the budget.

After discarding (a, b) and (b, c) , the only remaining attack is (c, a) . The set $\{a, b\}$ is now conflict-free with respect to the reduced attack relation and attacks c , making it a valid 3-stable extension.

This example illustrates how the introduction of a budget β permits certain extensions that are disallowed under classical semantics, thereby enabling reasoning in the presence of bounded inconsistency or uncertainty.

3. Weighted Assumption Based Argumentation

In this Section we develop our proposed framework for Weighted Assumption Based Argumentation (wABA). The basic idea is to extend ABA by assigning a cost to certain atoms, and use these costs to compute the weight of attacks.

Just like in plain ABA, arguments are defined as deductions. Attacks are still defined in terms of contraries as in ABA, i.e. any argument of the form $\Phi, \psi \vdash \phi$ may only be attacked by arguments of the form $\Phi' \vdash \bar{\psi}$.

Definition 3.1 (Weighted Assumption-Based Argumentation). A *wABA framework* is a tuple

$$(\mathcal{L}, \mathcal{R}, \mathcal{A}, \neg, \mathcal{S}, w)$$

where $(\mathcal{L}, \mathcal{R}, \mathcal{A}, \neg)$ is an ABA framework, $S \subseteq \mathbb{R}_0^+ \cup \{\infty\}$ (where \mathbb{R}_0^+ is the set of nonnegative real numbers) and $\mathcal{S} = (S, \oplus, \otimes, e_\oplus, e_\otimes)$ is a *semiring*¹. The function $w: \mathcal{L} \rightarrow S$ is a *weight* such that $w(\phi) = e_\otimes$ for each $\phi \in \mathcal{A}$.

The weights are then extended to attacks (x, y) among arguments x and y by letting²:

$$w((x, y)) = \bigotimes_{r \in \text{rules}(x)} \bigotimes_{\phi \in \text{body}(r)} w(\phi).$$

As in weighted abstract argumentation, one can then choose to discard attacks that do not exceed an inconsistency budget $\beta \in S$, using the operation \oplus of the semiring, as follows.

¹This means that (S, \oplus, e_\oplus) is a commutative monoid, (S, \otimes, e_\otimes) is a monoid, distributivity holds w.r.t. both operators, and any element is annihilated by e_\oplus .

²In other words, the weight of the attack is the weight of the attacking derivation, which is in turn computed by suitably aggregating the weights of the atoms occurring in the body of the rules supporting the derivation.

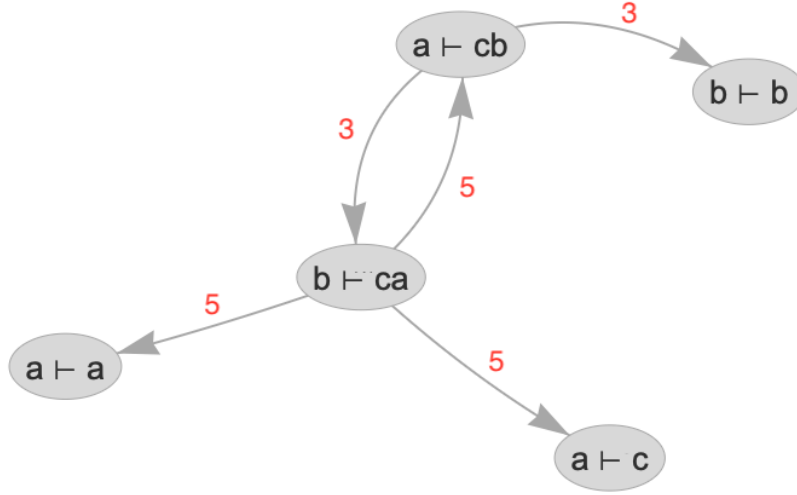


Figure 3: wABA framework from Example 3.3. Weights of attacks are displayed in red. Sets of rules used in derivations are: $rules(b \vdash ca) = \{d \leftarrow \top, ca \leftarrow b, d\}$, $rules(a \vdash cb) = \{d \leftarrow \top, c \leftarrow a, d, cb \leftarrow a, c\}$, $rules(a \vdash c) = \{d \leftarrow \top, c \leftarrow a, d\}$, and $rules(a \vdash a) = rules(b \vdash b) = \emptyset$.

Definition 3.2 (β - σ extensions). Let $(\mathcal{L}, \mathcal{R}, \mathcal{A}, \neg, \mathcal{S}, w)$ be a WABA framework, over the semiring $\mathcal{S} = (S, \oplus, \otimes, e_\oplus, e_\otimes)$, σ be a set of extensions over the abstract argumentation framework (Arg, Att) associated with $(\mathcal{L}, \mathcal{R}, \mathcal{A}, \neg)$, and $\beta \in S$. We say that a subset $Arg' \subseteq Arg$ is a β - σ extension for $(\mathcal{L}, \mathcal{R}, \mathcal{A}, \neg, \mathcal{S}, w)$ iff there is a subset $Att' \subseteq Att$ such that

$$\bigoplus_{(x,y) \in Att'} w(x,y) \leq \beta$$

and Arg' is a σ extension for $(Arg, Att \setminus Att')$.

We say that a wABA is *flat* if $(\mathcal{L}, \mathcal{R}, \mathcal{A}, \neg)$ is a flat ABA. In the following, we assume that all wABAs are flat, unless stated otherwise.

We illustrate the notions we have introduced with the following example.

Example 3.3. Consider the wABA consisting of atoms $\{a, b, c, d, ca, cb\}$, where a and b are assumptions, c and d are contextual information, and contraries are defined via $\bar{a} = ca$ and $\bar{b} = cb$. Weight of contextual atoms are defined as $w(c) = 3$ and $w(d) = 5$. Rules are:

$$\begin{aligned} ca &\leftarrow b, d \\ cb &\leftarrow a, c \\ c &\leftarrow a, d \\ d &\leftarrow \top \end{aligned}$$

The resulting framework may be depicted as in Figure 3. Note that if one has an appropriate inconsistency budget then some pairs of attacks may be removed from the framework, e.g., the mutual attacks from $b \vdash ca$ and $a \vdash cb$ for $rules(b \vdash ca) = \{d \leftarrow \top, ca \leftarrow b, d\}$ and $rules(a \vdash cb) = \{d \leftarrow \top, c \leftarrow a, d, cb \leftarrow a, c\}$. The resulting extensions can be calculated according to the standard AAF definition.

4. Implementation

Our implementation of (Weighted) Assumption-Based Argumentation is based on Answer Set Programming (ASP) using `clingo` [11]. The code is freely available at <https://github.com/dasaro/ABA-variants/tree/main/WABA>.

4.1. Theoretical assumptions

The current encoding supports *context-free*, *naive* and *stable* semantics.

Our implementation adopts the *min-max semiring*

$$\mathcal{S} = (\mathbb{N} \cup \{\infty\}, \max, \min, 0, \infty)$$

as our default structure for wABA. This choice is driven primarily by the inherently fuzzy nature of inputs, and `clingo` implementation, which provides native support for natural numbers but not for real values. We start by detailing the translation procedure which translates a wABA framework into an answer set program, which is largely based on ASPforABA [12].

4.2. Translation

Let $(\mathcal{L}, \mathcal{R}, \mathcal{A}, \neg, \mathcal{S}, w)$ be a flat wABA. For compactness and readability, in the following we use the standard `clingo` abbreviation “;” which unpacks an expression, e.g., of the form $r(a_1, b_1; \dots; a_n, b_n)$ into a set of clauses $r(a_1, b_1), \dots, r(a_n, b_n)$.

The set of assumptions $\mathcal{A} = \{a_1, \dots, a_n\}$ gets translated to:

```
assumption(a1; ...; an).
```

Each assumption is assigned to its contrary via the \neg operator, e.g., $\overline{a_1} = c_1, \dots, \overline{a_n} = c_n$. This gets translated to:

```
contrary(a1,c1; ...; an,cn).
```

Each rule of the form $h \leftarrow b_1, \dots, b_m$, where $h \in \mathcal{L} \setminus \mathcal{A}$ is the *head*, and $b_1, \dots, b_m \in \mathcal{L}$ are the *body* of the rule, is translated to:

```
head(id,h). body(id,b1; ...; id,bm).
```

where `id` is a unique identifier assigned to the rule, so that different rules get assigned different identifiers.

Finally, weights for atoms in \mathcal{L} , e.g., $w(d_1) = w_1, \dots, w(d_l) = w_l$ translate to:

```
weight(d1,w1; ...; dl, wl).
```

4.3. Semantics

The semantics file `core.lp` is fixed for all semantics and is described in what follows. We start by declaring a global inconsistency budget β :

```
budget(beta).
```

This can be set from the Command Line by using `clingo's -const beta=N` in-built flag. Note that if one wants to enumerate all the extensions, regardless of their budget, `-const beta=#sup` may be used.

Each assumption can be either in or out of a candidate extension:

```
in(X) :- assumption(X), not out(X).  
out(X) :- assumption(X), not in(X).
```

Support propagates from selected assumptions through the rules in a bottom-up fashion:

```
supported(X) :- assumption(X), in(X).  
supported(X) :- head(R,X), triggered_by_in(R).  
triggered_by_in(R) :- head(R,_), supported(X) : body(R,X).
```

We then assign supported atoms a weight according to the rules they were produced from, using the $\otimes = \min$ operator (and note that assumptions are assigned weight $e_\otimes = \infty$ which in `clingo` is rendered as the constant `#sup`):


```

supported_with_weight(X,#sup) :- assumption(X), in(X).
supported_with_weight(X,W) :- supported(X), weight(X,W).
supported_with_weight(X,W) :-
    supported(X), head(R,X),
    W = #min{ V, B : body(R,B), supported_with_weight(B,V) }.

```

Attacks arise from contraries:

```

attacks_with_weight(X,Y,W) :-
    supported(X), supported_with_weight(X,W),
    assumption(Y), contrary(Y,X).

```

As usual in wAAFs, we may choose to discard any subset of those attacks, paying their full weight; the total discarded weight must not exceed the budget w.r.t. semiring operation $\oplus = \max$:

```

{ discarded_attack(X,Y,W) : attacks_with_weight(X,Y,W) }.
extension_cost(C) :- C = #max{ W, X, Y : discarded_attack(X,Y,W) }.
:- extension_cost(C), C > B, budget(B).

```

Since we do not want discarded attacks to be effective, we also introduce the notion of a successful attack:

```

attacks_successfully_with_weight(X,Y,W) :-
    attacks_with_weight(X,Y,W), not discarded_attack(X,Y,W).

```

Once we have figured out how arguments are built from the assumptions, including their weights, contraries and discarded attacks according to the budget, we are left with the task of selecting an appropriate semantics to choose valid extensions from.

First, we introduce useful shorthands:

```

defeated(X) :- attacks_successfully_with_weight(_,X,_).
not_defended(X) :- attacks_successfully_with_weight(Y,X,_), not defeated(Y).

```

which state that an atom is *defeated* iff it receives a successful attack, and that it is *not defended* iff it has an undefeated attacker. These shorthands will help simplify the form of semantics below.

Then, we define four popular semantics:

Conflict-Free Conflict freeness ensures that, once discarded attacks are removed in the way outlined above, two assumptions in the same extension do not attack each other:

```

:- in(X), defeated(X).

```

Admissible semantics Admissibility is a widely adopted property of argumentation frameworks, and many other semantics (such as the stable semantics defined below) produce subsets of admissible sets. It is implemented by making sure it is context-free and all its elements are defended:

```

:- in(X), defeated(X). % conflict-freeness
:- in(X), not_defended(X).

```

Stable Semantics In the stable semantics, every (conflict-free) assumption outside the extension must be defeated by a non-discarded attack:

```

:- in(X), defeated(X). % conflict-freeness
:- out(X), not defeated(X).

```

Naive Semantics The naive semantics is a maximal conflict-free set w.r.t. to set inclusion, which is not necessarily admissible:

```

:- in(X), defeated(X). % conflict freeness
#heuristic in(X) : assumption(X). [1,true]

```

which ensures the maximal set of assumptions is “in” regardless of the extension being admissible or not.

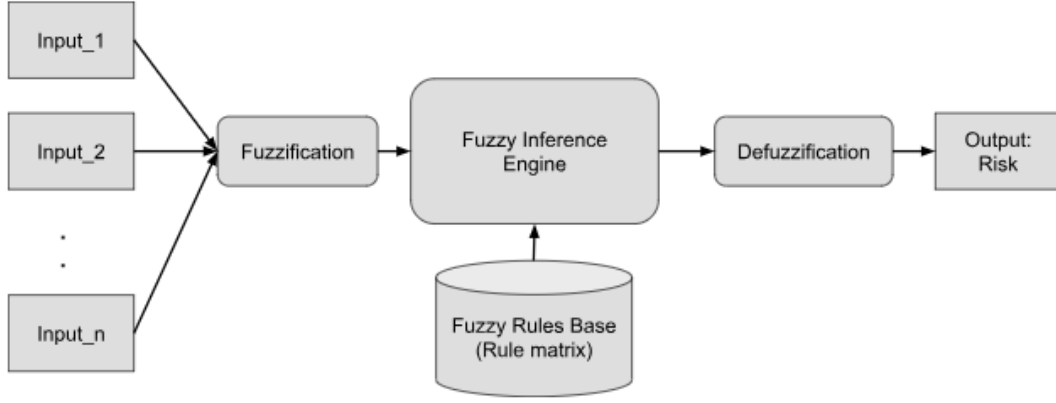


Figure 4: Architecture of the fuzzy system for ERA.

5. A Weighted ABA Approach to Ethical Reasoning

In this section, we apply the proposed `clingo` implementation of wABA to ethical reasoning involving conflicting principles and contextual medical information. We begin by revisiting the approach proposed in [6], which we expand into a more comprehensive (W)ABA-based framework for ethical decision-making. [6] focuses on the following motivating scenario:

Example 5.1 (Patient Dilemma). A care robot approaches a patient to administer medication, which would very likely address some health issues. The patient, however, refuses to take it. Should the robot attempt to persuade the patient, or should it respect the patient’s decision?

The *dilemma* involves a potential conflict between *autonomy*, which prioritizes the patient’s decision, and *beneficence*, which emphasizes promoting the patient’s well-being. Depending on the clinical and institutional context, other principles such as *non-maleficence* (e.g., medication may have harmful side effects) and *justice* (e.g., resource constraints across patients) may also be relevant.

Figure 4 shows the architecture of the fuzzy logic based system for ethical risk assessment (ERA) proposed in [6]. Dyoub & Lisi used ERA system to assess the possible ethical risk of causing physical harm to the patient in the above mentioned Patient Dilemma. For evaluating the physical harm risk in this case, we can consider different parameters as inputs to ERA system such as, severity of health condition of the patient, mental/psychological condition of the patient, physiological indicators of well-being, etc. These inputs are rated on some scale (e.g. between 0 and 10). The crisp values are then fuzzified into fuzzy sets with linguistic variables. For example, the fuzzy set label for *severity* could be one of $\{very_low, low, medium, high, very_high\}$. Then, the fuzzy inference engine uses the fuzzy rules from the fuzzy rule base to calculate the risk level. For instance, a patient in a severe health condition (90%) with reduced mental capacity (80%) is evaluated as *very_high* risk (95%). These values may be readily piped in into our wABA framework.

5.1. Plain ABA approach

We first show how we can embed linguistic labels from the fuzzy system to reason about which ethical principles are satisfied (resp. violated).

Example 5.2 (Patient Dilemma in ABA). Let us assume that, on a scale from 0 to 10, a patient in a highly severe physical condition (9/10), high mental risk (8/10) and refusing to taking her medications, is evaluated as a *high_risk* individual by the underlying fuzzy systems—meaning that not taking the medications could lead to developing severe physical consequences (including death).

We may formalize principles involved in this scenario through assumptions *respect_autonomy* and *act_beneficently*. Other atoms in the language are *risk(high)*, reflecting the fuzzy evaluation of risk,

$action(give_meds)$ as the action of giving medicines, $action(dont_give_meds)$ as the action of not giving medicines to the patient, $opinion(refuses_meds)$ as the contextual information that the patient refuses the medications, ca as the contrary of respect for autonomy (i.e., $\overline{respect_autonomy} = ca$) and cb as the contrary of $act_beneficently$ (i.e., $act_beneficently = cb$). Experts define rules as follows:

$$\begin{aligned}
opinion(refuses_meds) &\leftarrow \top \\
risk(high) &\leftarrow \top \\
action(give_meds) &\leftarrow risk(high), act_beneficently \\
action(dont_give_meds) &\leftarrow opinion(refuses_meds), respect_autonomy \\
ca &\leftarrow action(give_meds), opinion(refuses_meds) \\
cb &\leftarrow action(dont_give_meds), risk(high)
\end{aligned}$$

whose interpretation is intuitive. We can calculate the stable extensions in our framework, namely $\{act_beneficently\}$, from which action $give_meds$ follows, and $\{respect_autonomy\}$, which supports $not\ giving\ meds$. Note that, in an alternative situation where the patient does *not* refuse to take the meds both ethical principles may be satisfied: indeed, in this scenario we get the unique extension $\{act_beneficently, respect_autonomy\}$ from which one can derive that giving medications satisfies both *beneficence* and *autonomy*, since the patient does not refuse to take the medications in the first place.

This simple example can be further extended by considering more nuanced interactions between ethical principles and their contraries. A fuller example, involving the other two standard principle of bioethics (namely, *non maleficence* and *justice*) is available on our GitHub repository, and the interested reader can try it out.

We now turn to discussing how such reasoning about ethical principles may be enriched with weights that allow for inconsistencies.

5.2. wABA approach

In addition to assumptions and contraries, wABA introduces *weighted information*, and an *inconsistency budget* β . In this way, we can model more nuanced aspects of ethical decision-making.

Example 5.3 (Patient Dilemma in wABA). In order to show the characteristics of wABA, we further elaborate on Example 5.2. Recall that the patient is in a severe physical condition (0.9) with reduced mental capacity (0.8). In this case, we let the fuzzy component of our system *defuzzify* its output, which provides us with a numeric value for the risk of physical harm. Let us assume the output defuzzified value for the risk is 0.7. Furthermore, let us assume that the patient is very reluctant to taking medications, which we assign (or the fuzzy system assigns) a weight of 0.9. This amounts to saying that the contextual information is described by:

$$w(risk) = 0.7, \quad w(opinion(refuses_meds)) = 0.9$$

It is worth noting here that we have dropped the linguistic variables appearing in the first example (i.e., $risk(high)$) in favor of explicit weights. We modify the rules to reflect this:

$$\begin{aligned}
opinion(refuses_meds) &\leftarrow \top \\
risk &\leftarrow \top \\
action(give_meds) &\leftarrow risk, act_beneficently \\
action(dont_give_meds) &\leftarrow opinion(refuses_meds), respect_autonomy \\
ca &\leftarrow action(give_meds), opinion(refuses_meds) \\
cb &\leftarrow action(dont_give_meds), risk
\end{aligned}$$

where the weights of *risk* and *refuses_meds* automatically enter the computation by means of wABA semantics.

In the implementation, we use the appropriate predicate `weight/2` as follows:

```
head(ctx1, risk).           weight(risk, 7).
head(ctx2, refuses_meds).  weight(refuses_meds, 9).
```

Assumptions, rules and contraries are implemented in the exact same way as in plain ABA (see Example 5.2). Note that if we let $\beta = 0$ we reconstruct exactly the same extensions as in Example 5.2, i.e., the (postprocessed for readability) output looks as follows:

Answer: 1

```
in(respect_autonomy) supported_with_weight(action(dont_give_meds),9) supported_with_weight(cb,7)
extension_cost(0)
```

Answer: 2

```
in(act_beneficently) supported_with_weight(action(give_meds),7) supported_with_weight(ca,7)
extension_cost(0)
```

However, if we list *all* extensions, we get another *inconsistent* stable extension resulting from the removal of the ethical inconsistency between `respect_autonomy` and `act_beneficently`:

Answer: 3

```
in(respect_autonomy) in(act_beneficently) supported_with_weight(action(dont_give_meds),9)
supported_with_weight(action(give_meds),7) supported_with_weight(ca,7) supported_with_weight(
cb,7) extension_cost(7)
```

This extension comes at an *inconsistency cost* of 7. However, we now have the additional information that `action(give_meds)` weighs 7 in all extensions, while `action(dont_give_meds)` weighs 9. Therefore, it is *slightly* recommended *not to give meds* in this scenario, thus respecting autonomy. This suggestion could of course be overridden by a human operator if s/he wishes to satisfy beneficence over autonomy.

We illustrated how wABA may resolve conflicts between ethical principles by deriving attacks from contraries and evaluating which sets of assumptions can be tolerated together under the given inconsistency budget. The weighted implementation allows discarding some attacks—e.g., between autonomy and beneficence—if the total weight of discarded attacks remains within the user-specified limit.

Note that this is particularly useful in this scenario as it also allows for reasoning about ethical principles in highly inconsistent scenarios, where plain ABA would not produce extensions, as it is often the case in realistic scenarios due to ethical principles leading to different courses of actions. In such cases wABA can help a human operator getting the fuller picture, e.g., verify what principles are satisfied in most extensions, normalizing their weight according to inconsistency levels, as well as showing what course of action is recommended and what this implies on the ethical dimension.

6. Related Work

Resolving conflicts between ethical principles has long been a core challenge in both normative ethics and applied biomedical ethics. Clinical decision-making frequently necessitates the careful balancing of competing moral considerations, such as honoring patient autonomy, promoting well-being, preventing harm, and ensuring fairness. The principlist approach proposed by Beauchamp and Childress [13] remains a cornerstone of biomedical ethics. It articulates four central principles -autonomy, beneficence, nonmaleficence, and justice - that guide ethical decision-making. In his influential work [14], Floridi further argues that a fifth principle, namely *explicability* is a required addition to those principles, specifically for the needs of AI ethics. However, these ethical frameworks do not prescribe a unique resolution when these principles come into conflict, prompting the need for formal methods that can assist in principled judgment. Subsequent proposals, such as the mixed consequentialist-nonconsequentialist hierarchy, advocate balancing principles within categories (e.g., nonmaleficence vs. beneficence) before applying lexical priority to nonconsequentialist outcomes [15]. Similarly, [16] provides a general

workflow to derive concrete rules from general principles, taking into account the specificity of the contexts, and the allowable exceptions to rules. This approach has been subsequently implemented in Datalog [17]. These frameworks emphasize the role of case-based reasoning and iterative specification to resolve dilemmas.

Several computational approaches have been proposed to tackle this issue. Anderson and Anderson, in their *MedEthEx* [18] and *EthEl* [19], operationalize Beauchamp and Childress' principles through machine learning, deriving decision rules from biomedical ethicists' intuitions in training cases. In these systems, the intensity of duty violations (e.g., autonomy vs. beneficence) is quantified by assigning it a weight, then, for each possible action, the system computes the weighted sum of duty satisfaction. After that, using inductive logic programming (ILP), their systems generate actionable rules for recurring dilemmas. In [20] and [21], Rossi and her team formalize principles as constraint-based systems with conflict resolution engines. They argue that preferences can model the relative importance of ethical principles and help resolve conflicts by identifying most preferred outcomes, given context-sensitive constraints. Their work leverages CP-nets and weighted constraints to evaluate ethical options, offering a flexible and explainable approach to value-sensitive decision-making. In [22], Kleiman-Weiner et al. suggest an abstract and recursive utility calculus to resolve conflicts among moral principles. Moral theories (for the purposes of trading off different agents' interests) can be formalized as values or weights that an agent attaches to a set of abstract principles for how to factor any other agents' utility functions into their own utility-based decision-making and judgment. Drawing on machine learning and computational social choice, [23] proposes an algorithm to learn a model of societal preferences, and, when faced with specific ethical dilemma at runtime, aggregate those preferences to identify a desirable choice. Awad et al. in [24], propose a framework for incorporating public opinion, as essential tool, into policy making in situations where ethical values are in conflict. Their framework advocates creating vignettes representing abstract value choices, eliciting the public's opinion on these choices, and using machine learning, in particular ILP, to extract principles that can serve as succinct statements of the policies implied by these choices and rules that can be embedded in algorithms to guide the behavior of AI-based systems. To present the functionality of this proposal, the authors borrow vignettes from the Moral Machine website [25].

Dennis et al. [26] developed the ETHAN (a BDI (Belief-Desire-Intention) agent language) system that deals with situations when civil air navigation regulations are in conflict. The system relates these rules to four hierarchical ordered ethical principles (do not harm people, do not harm animals, do not damage self, and do not damage property) and develops a course of action that generates the smallest violation to those principles in case of conflict. In their prototype, ethical reasoning was integrated into a BDI agent programming language via the agent plan selection mechanism. [27] implement a BDI architecture within a multi-agent system (MAS), with a particular emphasis on handling norm conflicts. In their framework, agents are capable of adopting and dynamically updating norms, and they determine which norms to activate based on the current context, their desires, and their intentions. Conflicts between norms are resolved by selecting the norm that best contributes to the fulfillment of the agent's goals and intentions, effectively embedding norm adherence within the agent's motivational structure. Similarly, Mermet and Simon [28] address norm conflicts by distinguishing between moral and ethical rules, the latter being invoked when moral rules are in conflict. They perform a verification of whether their system called GDT4MAS, is able to choose the correct ethical rule in conflict cases.

Chorley et al. in [29] described an implementation of the approach to deliberation about a choice of action based on presumptive argumentation and associated critical questions [30]. The authors use the argument scheme proposed in [31] to generate presumptive arguments for and against actions, and then subject these arguments to critical questioning. They have explored automation of argumentation for practical reasoning by a single agent in a multi-agent context, where agents may have conflicting values. Their approach was illustrated with a particular example based on an ethical dilemma.

[32] uses abstract argumentation for decision-making with multiple experts. The approach represents in the form of attacks in abstract argumentation, the lack of fairness in the evaluation performed by an expert, which may thus lead to the dismissal of the evaluation. Finally, [33] is directly related with our approach. It introduces a moral advisor based on logic-based normative systems integrated

with Aspic-style structured argumentation, in order to handle moral dilemmas involving multiple stakeholders, bearing different interests and ethical views.

7. Conclusion and future work

We have introduced *wABA*, together with an ASP based implementation. We applied our formalism to a patient dilemma scenario, showing how it can smoothly integrate the weighted assessment of a medical situation, the ethical principles involved, and possibly weighted solutions to the dilemma. We believe that this framework may provide a useful module for implementing reasoning with ethical principles in AI systems that interact with humans. Ongoing work is devoted to a detailed formal analysis of the framework, an extension of the implementation, and its applications in the ethical domain. Concerning the first aspect, we plan to analyze the computational properties of *wABA* and its implementation, such as correctness w.r.t. the intended semantics and analysis of its complexity. We would also like to further explore the relation to the different choices of the semantics and the choice of semiring. We plan to compare our framework to other approaches, based on structured argumentation frameworks, in particular involving preferences, such as ASPIC+ [2] and ABA+ [34].

Concerning the implementation, we are incorporating further formal argumentation semantics and operations for the aggregation of weights.

Finally, for applications, we are currently investigating realistic scenarios, involving reasoning with data and conflicting ethical principles, where we believe that our system may provide valuable support. A particularly promising application in this sense would be the analysis of medical triage, with their related intricate legal and ethical regulations [35].

Acknowledgments

This work was partially supported by the project FAIR- Future AI Research (PE00000013), under the NRRP MUR program funded by the NextGenerationEU.

Declaration on Generative AI

The authors have not employed any Generative AI tools.

References

- [1] P. M. Dung, On the acceptability of arguments and its fundamental role in nonmonotonic reasoning, logic programming and n-person games, *Artificial Intelligence* 77 (1995) 321–357. doi:10.1016/0004-3702(94)00041-x.
- [2] F. H. van Eemeren, B. Verheij, T. F. Gordon, P. Baroni, M. Caminada, G. Brewka, S. Ellmauthaler, H. Strass, J. P. Wallner, S. Woltran, X. Fan, C. Schulz, F. Toni, P. Besnard, A. Hunter, F. Macagno, D. Walton, C. Reed, *Handbook of Formal Argumentation*, College Publications, 2018.
- [3] D. Gabbay, M. Giacomin, G. Simari, *Handbook of Formal Argumentation*, Volume 2, v. 2, College Publications, 2021. URL: <https://books.google.it/books?id=JUekzgEACAAJ>.
- [4] W. D. Ross, *The Right and the Good*, Oxford University Press, Oxford, UK, 1930. doi:10.2307/2180065.
- [5] A. Carnevale, A. Lombardi, F. A. Lisi, A human-centred approach to symbiotic AI: Questioning the ethical and conceptual foundation, *Intelligenza Artificiale* 18 (2024) 9–20. doi:10.3233/IA-240034.
- [6] A. Dyoub, F. A. Lisi, Towards Ethical Risk Assessment of Symbiotic AI Systems with Fuzzy Rules, *CEUR Workshop Proceedings* 3881 (2024) 36–49.

- [7] F. Toni, A tutorial on assumption-based argumentation, *Argument & Computation* 5 (2014) 89–117. URL: <https://doi.org/10.1080/19462166.2013.869878>. doi:10.1080/19462166.2013.869878. arXiv:<https://doi.org/10.1080/19462166.2013.869878>.
- [8] A. Borg, D. Odekerken, Pyarg for solving and explaining argumentation in python: Demonstration, in: F. Toni, S. Polberg, R. Booth, M. Caminada, H. Kido (Eds.), *Computational Models of Argument - Proceedings of COMMA 2022*, Cardiff, Wales, UK, 14-16 September 2022, volume 353 of *Frontiers in Artificial Intelligence and Applications*, IOS Press, 2022, pp. 349–350. URL: <https://doi.org/10.3233/FAIA220167>. doi:10.3233/FAIA220167.
- [9] P. E. Dunne, A. Hunter, P. McBurney, S. Parsons, M. Wooldridge, Weighted argument systems: Basic definitions, algorithms, and complexity results, *Artificial Intelligence* 175 (2011) 457–486. URL: <http://dx.doi.org/10.1016/j.artint.2010.09.005>. doi:10.1016/j.artint.2010.09.005.
- [10] S. Bistarelli, F. Santini, Weighted argumentation, *FLAP* 8 (2021) 1589–1622. URL: <https://collegepublications.co.uk/ifcolog/?00048>.
- [11] M. GEBSER, R. KAMINSKI, B. KAUFMANN, T. SCHAUB, Multi-shot asp solving with clingo, *Theory and Practice of Logic Programming* 19 (2019) 27–82. doi:10.1017/S1471068418000054.
- [12] T. Lehtonen, J. Wallner, M. Järvisalo, *Aspforaba - asp-based algorithms for reasoning in aba*, 2023.
- [13] T. L. Beauchamp, J. F. Childress, et al., *Principles of biomedical ethics*, eighth ed., Oxford University Press, USA, 2019.
- [14] L. Floridi, *The Ethics of Artificial Intelligence: Principles, Challenges, and Opportunities*, Oxford University Press, 2023. URL: <https://doi.org/10.1093/oso/9780198883098.001.0001>. doi:10.1093/oso/9780198883098.001.0001.
- [15] R. M. Veatch, Resolving conflicts among principles: ranking, balancing, and specifying, *Kennedy Institute of Ethics Journal* 5 (1995) 199–218.
- [16] B. Townsend, C. Paterson, T. T. Arvind, G. Nemirovsky, R. Calinescu, A. Cavalcanti, I. Habli, A. Thomas, From Pluralistic Normative Principles to Autonomous-Agent Rules, *Minds and Machines* 32 (2022) 683–715. URL: <https://doi.org/10.1007/s11023-022-09614-w>. doi:10.1007/s11023-022-09614-w.
- [17] M. Mirani, F. Raimondi, N. Troquard, Towards Efficient Norm-Aware Robots ’ Decision Making Using Datalog, in: *3rd Workshop on Bias, Ethical AI, Explainability and the Role of Logic and Logic Programming – Joint Workshop @ AIXIA 2024*, Bolzano, Italy, volume 7713, 2024.
- [18] M. Anderson, S. L. Anderson, C. Armen, Medethex: Toward a medical ethics advisor, in: *Caring Machines: AI in Eldercare, Papers from the 2005 AAAI Fall Symposium*, Arlington, Virginia, USA, November 4-6, 2005., volume FS-05-02 of *AAAI Technical Report*, AAAI Press, USA, 2005, pp. 9–16. URL: <https://www.aaai.org/Library/Symposia/Fall/fs05-02.php>.
- [19] M. Anderson, S. L. Anderson, ETHEL: toward a principled ethical eldercare system, in: *AI in Eldercare: New Solutions to Old Problems, Papers from the 2008 AAAI Fall Symposium*, Arlington, Virginia, USA, November 7-9, 2008, volume FS-08-02 of *AAAI Technical Report*, AAAI, USA, 2008, pp. 4–11. URL: <http://www.aaai.org/Library/Symposia/Fall/fs08-02.php>.
- [20] F. Rossi, Safety constraints and ethical principles in collective decision making systems, in: S. Hölldobler, M. Krötzsch, R. Peñaloza, S. Rudolph (Eds.), *KI 2015: Advances in Artificial Intelligence - 38th Annual German Conference on AI*, Dresden, Germany, September 21-25, 2015, *Proceedings*, volume 9324 of *Lecture Notes in Computer Science*, Springer, 2015, pp. 3–15. URL: https://doi.org/10.1007/978-3-319-24489-1_1. doi:10.1007/978-3-319-24489-1_1.
- [21] A. Loreggia, N. Mattei, F. Rossi, K. B. Venable, Preferences and ethical principles in decision making, in: J. Furman, G. E. Marchant, H. Price, F. Rossi (Eds.), *Proceedings of the 2018 AAAI/ACM Conference on AI, Ethics, and Society*, AIES 2018, New Orleans, LA, USA, February 02-03, 2018, ACM, 2018, p. 222. URL: <https://doi.org/10.1145/3278721.3278723>. doi:10.1145/3278721.3278723.
- [22] M. Kleiman-Weiner, R. Saxe, J. B. Tenenbaum, Learning a commonsense moral theory, *Cognition* 167 (2017) 107–123. URL: <https://www.sciencedirect.com/science/article/pii/S0010027717300707>. doi:<https://doi.org/10.1016/j.cognition.2017.03.005>, moral Learning.
- [23] R. Noothigattu, S. N. S. Gaikwad, E. Awad, S. Dsouza, I. Rahwan, P. Ravikumar, A. D. Procaccia, A voting-based system for ethical decision making, in: S. A. McIlraith, K. Q. Weinberger (Eds.),

- Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence, (AAAI-18), the 30th innovative Applications of Artificial Intelligence (IAAI-18), and the 8th AAAI Symposium on Educational Advances in Artificial Intelligence (EAAI-18), New Orleans, Louisiana, USA, February 2-7, 2018, AAAI Press, 2018, pp. 1587–1594. URL: <https://doi.org/10.1609/aaai.v32i1.11512>. doi:10.1609/AAAI.V32I1.11512.
- [24] E. Awad, M. Anderson, S. L. Anderson, B. Liao, An approach for combining ethical principles with public opinion to guide public policy, *Artificial Intelligence* 287 (2020) 103349. URL: <https://www.sciencedirect.com/science/article/pii/S0004370219301079>. doi:<https://doi.org/10.1016/j.artint.2020.103349>.
 - [25] E. Awad, S. Dsouza, R. Kim, J. Schulz, J. Henrich, A. Shariff, J.-F. Bonnefon, I. Rahwan, The moral machine experiment, *Nature* 563 (2018) 59–64.
 - [26] L. A. Dennis, M. Fisher, M. Slavkovik, M. Webster, Formal verification of ethical choices in autonomous systems, *Robotics Auton. Syst.* 77 (2016) 1–14. URL: <https://doi.org/10.1016/j.robot.2015.11.012>. doi:10.1016/J.ROBOT.2015.11.012.
 - [27] B. F. dos Santos Neto, V. T. da Silva, C. J. P. de Lucena, NBDI: an architecture for goal-oriented normative agents, in: J. Filipe, A. L. N. Fred (Eds.), *ICAART 2011 - Proceedings of the 3rd International Conference on Agents and Artificial Intelligence*, Volume 1 - Artificial Intelligence, Rome, Italy, January 28-30, 2011, SciTePress, 2011, pp. 116–125.
 - [28] B. Mermet, G. Simon, Formal verification of ethical properties in multiagent systems, in: G. Bonnet, M. Harbers, K. V. Hindriks, M. Katell, C. Tessier (Eds.), *Proceedings of the 1st Workshop on Ethics in the Design of Intelligent Agents*, The Hague, The Netherlands, August 30, 2016, volume 1668 of *CEUR Workshop Proceedings*, CEUR-WS.org, 2016, pp. 26–31. URL: <https://ceur-ws.org/Vol-1668/paper5.pdf>.
 - [29] A. Chorley, T. J. M. Bench-Capon, P. McBurney, Automating argumentation for deliberation in cases of conflict of interest, in: P. E. Dunne, T. J. M. Bench-Capon (Eds.), *Computational Models of Argument: Proceedings of COMMA 2006*, September 11-12, 2006, Liverpool, UK, volume 144 of *Frontiers in Artificial Intelligence and Applications*, IOS Press, 2006, pp. 279–290. URL: <http://www.booksonline.iospress.nl/Content/View.aspx?piid=1949>.
 - [30] D. N. Walton, *Argumentation Schemes for Presumptive Reasoning* (), Routledge, New York, USA, 1996. doi:<https://doi.org/10.4324/9780203811160>, eBook Published: 5 November 2013.
 - [31] K. Atkinson, What should we do? : computational representation of persuasive argument in practical reasoning, Ph.D. thesis, University of Liverpool, UK, 2005. URL: <https://ethos.bl.uk/OrderDetails.do?uin=uk.bl.ethos.426134>.
 - [32] S. Bistarelli, M. Ceberio, J. A. Henderson, F. Santini, Abstract argumentation frameworks to promote fairness and rationality in multi-experts multi-criteria decision making, *Studies in Systems, Decision and Control* 100 (2018) 7–19. doi:10.1007/978-3-319-61753-4_2.
 - [33] B. Liao, P. Pardo, M. Slavkovik, L. van der Torre, The Jiminy Advisor: Moral Agreements among Stakeholders Based on Norms and Argumentation, *Journal of Artificial Intelligence Research* 77 (2023) 737–792. doi:10.1613/jair.1.14368.
 - [34] K. Čyras, F. Toni, Aba+: assumption-based argumentation with preferences, in: *Proceedings of the Fifteenth International Conference on Principles of Knowledge Representation and Reasoning, KR'16*, AAAI Press, 2016, p. 553–556.
 - [35] A. Vântu, A. Vasilescu, A. Băicoianu, Medical emergency department triage data processing using a machine-learning solution, *Heliyon* 9 (2023). doi:10.1016/j.heliyon.2023.e18402.