

A multi-agent argumentation framework to support collective reasoning*

Jordi Ganzer-Ripoll¹, Maite López-Sánchez¹, Juan Antonio Rodriguez-Aguilar²

¹ University of Barcelona, Barcelona, Spain

² Artificial Intelligence Research Institute (IIIA-CSIC), Bellaterra, Spain
jordi891@gmail.com, maite@maia.ub.es, jar@iiia.csic.es

Abstract. Argumentative debates are a powerful tool for resolving conflicts and reaching agreements in open environments such as on-line communities. Here we introduce an argumentation framework to structure argumentative debates. Our framework represents the arguments issued by the participants involved in a debate, the (attack and defence) relationships between them, as well as participants' opinions on them. Furthermore, we tackle the problem of computing a collective decision from participants' opinions. With this aim, we design an aggregation function that satisfies valuable social-choice properties.

1 Introduction

As argued in [10,12], argumentative debates are a powerful tool for reaching agreements in open environments such as on-line communities. Nowadays, this is particularly true in our society due to the increasing interest and deployment of e-participation and e-governance ICT-systems that involve citizens in governance [17]. Not surprisingly some European cities are opening their policy making to citizens (e.g. Reykjavík [2], Barcelona [1]). Moreover, the need for argumentative debates has also been deemed as necessary for open innovation systems [14]. On-line debates are usually organised as threads of arguments and counter-arguments that users issue to convince others so that debates eventually converge to agreements. Users are allowed to express their opinions on arguments by rating them (e.g. [12]). There are two main issues in the management of large-scale on-line debates. First, as highlighted by [10] and [12], there is simply too much noise when many individuals participate in a discussion, and hence there is the need for *structuring* it to keep the focus. Second, the opinions on arguments issued by users must be aggregated to achieve a collective decision about the topic under discussion [4]. In this paper we try to make headway on these two issues.

Recently, argumentation has become one of the key approaches to rational interaction in artificial intelligence [5,16]. Here, we propose to follow an argumentation-based approach that allows agents to issue arguments in favour or against a *topic* under discussion as well as about other agents' arguments. Furthermore, we will consider that agents express their opinions about each other's arguments and the topic itself.

Within our multi-agent framework, we face the following collective decision problem: *given a set of agents, each with an individual opinion about a given set of arguments related to a topic, how can agents reach a collective decision on the topic*

* Funded by Collectiveware TIN2015-66863-C2-1-R (MINECO/FEDER) and 2014 SGR 118.

under discussion? To solve this problem, we propose a social choice function that aggregates agents' opinions to infer the overall opinion about the topic under discussion. Our aggregation function is based on combining opinions and exploiting dependencies between arguments to produce an aggregated opinion. Moreover, and most importantly, our aggregation function guarantees the resulting aggregated opinion to be *coherent*, i.e., it is free of contradictions. In more detail, here we make the following contributions:

- A novel multi-agent argumentation framework, the so-called *target-oriented discussion framework*, to support discussions about the acceptance of a target proposal. Besides the usual attack relationship between arguments, our framework allows agents to express explicit defence relationships between arguments. Furthermore, it introduces a mechanism for assessing whether individual opinions about the arguments are reasonable (coherent) or not. Formally, this is captured through the notion of *coherent labelling*, which can be regarded as a relaxed version of the notion of complete labelling in [4] to provide further flexibility to express opinions.
- A novel aggregation function that combines agents' opinions in our multi-agent argumentation framework to assess the collective decision reached by the agents about the topic under discussion. Interestingly, our aggregation function guarantees the *coherent collective rationality* of the outcome. Besides collective rationality, we show that our aggregation function satisfies further valuable social-choice theoretic properties for the argumentation domain.

Organisation. Sections 2 and 3 characterise and formalise our multi-agent argumentation framework; section 4 details both our decision problem and the desired properties of an aggregation function; section 5 introduces our aggregation function and studies its social-choice properties; and section 6 draws conclusions and plans future research.

2 Characterising a target-oriented discussion framework

From a general perspective, we envision a setting where some individuals discuss collectively about a given issue or topic (the so-called *target*) with the aim of reaching a consensus on it. Discussion is articulated by means of arguments in favour or against this topic. Thus, we consider an argumentation scenario where participants issue their opinions by labelling such arguments. For explanatory purposes, below we consider that this topic under discussion may well correspond to a norm. Henceforth, we will refer to this setting as a target-oriented argumentation framework (see Section 3 for a formal definition). Within this target-oriented argumentation framework, two argument relationships (attack and defence) can be established so that arguments are defined as being in favour of (or against) other arguments. Both relationships are binary, directed and mutually exclusive, and the target is an argument that deserves special attention since it is the only one not defending nor attacking any other argument. Additionally, participants can show that they like or dislike some (not necessarily all) existing arguments. In order to do so, participants assign labels to each argument indicating whether they accept it; reject it; or they abstain from deciding whether to accept or reject it. Thus, participants can also explicitly indicate that they are uncertain about whether they like or dislike an argument. Moreover, we consider that this uncertainty may also capture

the fact that a participant may skip providing an opinion (i.e., label) on an argument, which seems to be a suitable feature when dealing with human agents.

Overall, the problem we tackle is that of aggregating all the legitimate and subjective participant's opinions (i.e., labellings) into a single collective one. That will allow us to assess whether the group of participants: accepts the topic under discussion; rejects it; or there is not enough support in favour or against the given target.

However, considering human participants prevents us from requiring rationality, since contradictions or inconsistencies may occur when expressing opinions. In fact, we aim at designing an aggregation function that guarantees some desirable properties so that the outcome does represent the consensus on the topic under discussion. From these properties we highlight that of *coherent labelling*, which intuitively characterises whether an individual exhibits non-contradictory opinions (i.e., labelling). The next section introduces this concept formally and subsequent sections study how our aggregation function results in a single aggregated coherent labelling which also satisfies further desirable properties. Next, we introduce a simple example that will allow us to illustrate some of the presented concepts along the paper.

Example 1 (Flatmates' discussion) Consider three flatmates (Alan, Bart, and Cathy) discussing about norm (N): "Flatmates take fixed turns for dishwashing at 10 p.m." and issuing the following arguments: a_1 = "10 p.m. is too late and cannot be changed"; a_2 = "Schedule is too rigid"; and a_3 = "Fair distribution". Notice that: arguments a_1 and a_2 attack N whereas a_3 defends it; and a_1 is in favour of a_2 . Once all arguments and their relations are clear, flatmates express their opinions by accepting, rejecting (or not opining about) each argument : (1) Alan (Ag_1) gets up early 4 days per week, and so (as first row in Table 1 shows) he rejects norm N and accepts arguments a_1 and a_2 . Nevertheless, he acknowledges and accepts argument a_3 . (2) Bart (Ag_2) has spare time at night and is clearly pro norm N . Second row in Table 1 shows he accepts N and a_3 , and rejects a_1 and a_2 . Finally, (3) Cathy (Ag_3) is keen on routines so she rejects a_2 and accepts N , a_1 , and a_3 (see third row in Table 1).

Therefore, the question that arises is how to aggregate all these opinions so that a consensus is reached over the acceptance (or not) of this dish-washing norm.

3 The target-oriented discussion framework

The flatmates' discussion illustrates the main elements of our argumentation framework. Within such framework, the norm constitutes the target of a multi-agent argumentation scenario where: i) a number of arguments are issued; and ii) participating agents express their opinions about those arguments as well as about the norm under discussion. Additionally, we characterise *coherent* opinions as those not incurring in contradictions. The purpose of this section is to formally capture all these core elements of our argumentation framework. Thus, section 3.1 introduces the *target-oriented discussion framework*, section 3.2 characterises the formal structure of an agent's opinion, and section 3.3 characterises our notion of *coherent* opinion.

		Arguments			
		N	a ₁	a ₂	a ₃
Agents	Ag 1	✗	✓	✓	✓
	Ag 2	✓	✗	✗	✓
	Ag 3	✓	✓	✗	✓

Table 1. Flatmates’ opinions in the discussion on the dish-washing norm.

3.1 Formalising our argumentation framework

Our purpose is to provide an argumentation framework that allows one to capture both attack and defence relationships between arguments, as done in bipolar argumentation frameworks [8,3].³ The motivation for including defence relationships is based on recent studies in large-scale argumentation frameworks involving humans (e.g. [14,13]). There, humans naturally handle both attack and defence relationships between arguments. Our notion of *discussion framework* aims at offering such expressiveness.

Definition 1 A *discussion framework* is a triple $DF = \langle \mathcal{A}, \mapsto, \Vdash \rangle$, where \mathcal{A} is a finite set of arguments, and $\mapsto \subseteq \mathcal{A} \times \mathcal{A}$ and $\Vdash \subseteq \mathcal{A} \times \mathcal{A}$ stand for attack and defence relationships that are disjoint, namely $\mapsto \cap \Vdash = \emptyset$. We say that an argument $b \in \mathcal{A}$ attacks another argument $a \in \mathcal{A}$ iff $b \mapsto a$, and that b defends a iff $b \Vdash a$.

A discussion framework can be depicted as a graph whose nodes stand for arguments and whose edges represent either attack or defence relationships between arguments. Figure 1 shows our graphical representation of attack and defence relationships.

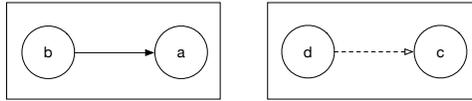


Fig. 1. Representation of an attack relationship $b \mapsto a$ and a defence relationship $d \Vdash c$.

Each argument in a discussion framework can be indirectly related to other arguments through a chain of attack and defence relationships. Given an argument, we capture its indirect relationships with other arguments through the notion of *descendant*.

Definition 2 Let $DF = \langle \mathcal{A}, \mapsto, \Vdash \rangle$ be a discussion framework and $a \in \mathcal{A}$ one of its arguments. We say that an argument $b \in \mathcal{A}$ is a *descendant* of a if there is a finite subset of arguments $\{c_1, \dots, c_r\} \subseteq \mathcal{A}$ such that $b = c_1$, $c_1 R_1 c_2$, \dots , $c_{r-1} R_{r-1} c_r$, $c_r = a$ and $R_i \in \{\mapsto, \Vdash\}$ for all $1 \leq i < r$.

³ Nevertheless, there are notable differences with bipolar argumentation frameworks. First, bipolar argumentation does not consider labellings (different opinions on arguments), nor their aggregation. Second, bipolar argumentation focuses on studying the structure between arguments and groups of arguments, whereas we focus on computing a collective decision from differing opinions about arguments. Third, arguments in bipolar argumentation can be regarded as objective facts, while in our case, arguments can be subjective facts on which individuals can differ. Thus, our argumentation framework is less restrictive to include humans in the loop.

Now we are ready to define our argumentation framework, the so-called *target-oriented discussion framework*, which considers that there is a target argument (e.g. a norm or a proposal) under discussion.

Definition 3 A *target-oriented discussion framework* $TODF = \langle \mathcal{A}, \mapsto, \Vdash, \tau \rangle$ is a discussion framework satisfying the following properties: (i) for every argument $a \in \mathcal{A}$, a is not a descendant of itself; and (ii) there is an argument $\tau \in \mathcal{A}$, called the target, such that for all $a \in \mathcal{A} \setminus \{\tau\}$, a is a descendant of τ .

Observation 1 From the previous definitions we infer some properties that help us further characterise a target-oriented discussion framework:

1. No reflexivity. No argument can either attack or defend itself. Formally, $\forall a \in \mathcal{A}$, $a \not\mapsto a$ and $a \not\Vdash a$.
2. No reciprocity. If an argument a attacks another argument b , then a cannot be attacked nor defended back by b , namely $\forall a, b \in \mathcal{A}$, if $a \mapsto b$ then $b \not\mapsto a$ and $b \not\Vdash a$. Analogously, if an argument a defends another argument b , a cannot be defended nor attacked by b , namely $\forall a, b \in \mathcal{A}$, if $a \Vdash b$ then $b \not\mapsto a$ and $b \not\Vdash a$.
3. No target contribution. The target neither attacks nor defends any other argument, namely for all $a \in \mathcal{A} \setminus \{\tau\}$, $\tau \not\mapsto a$ and $\tau \not\Vdash a$. This distinguishes the special role of the target as the center of discussion to which attacks and supports are directly or indirectly pointed.

The next result follows from definition 2 and the observation 1.

Proposition 4 Let $TODF = \langle \mathcal{A}, \mapsto, \Vdash, \tau \rangle$ be a target-oriented discussion framework and $E = \mapsto \cup \Vdash$. The graph associated to a $TODF$, $G = \langle \mathcal{A}, E \rangle$, is a directed acyclic graph, where \mathcal{A} is the set of nodes and E the edge relationship.

Proof. Straightforward from definition 2 and observation 1.

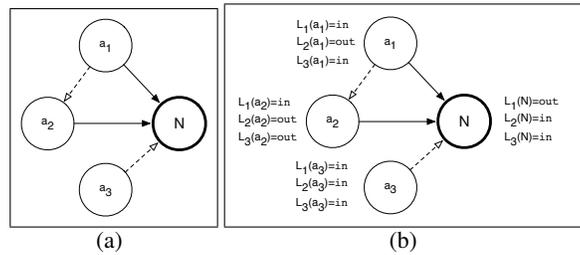


Fig. 2. Flatmates example: (a) associated graph to $TODF$; (b) $TODF$ together with labellings.

Example 2 (Flatmates' example formalization) Figure 2(a) depicts the flatmates' target-oriented discussion framework. The nodes in the graph represent the set of arguments $\mathcal{A} = \{N, a_1, a_2, a_3\}$ in the example of section 2, where N is the dish-washing norm, and a_1, a_2, a_3 are the rest of arguments. Thus, N , the norm under discussion, is taken to be τ in our $TODF$. As to edges, they represent both the attack and defence relationships: $a_1 \mapsto N$, $a_2 \mapsto N$ and $a_1 \Vdash a_2$, $a_3 \Vdash N$ respectively.

3.2 Argument labellings

Given a target-oriented argumentation framework shared by agents, now we focus on how these encode their opinions (argument evaluations). Here we consider that each agent’s opinion over our argumentation framework corresponds to a *labelling* [6,7]. Furthermore, we adhere to the labelling-based semantics proposed by Caminada in [6,7], which gives a labelling per argument. By means of argument labellings, each agent can support an argument (by labelling it as *in*), reject it (by labelling it as *out*), or abstain from deciding whether to accept it or reject it (by labelling it as *undec*). Besides expressing uncertainty regarding the assessment of an argument, the *undec* label stands for the absence of an opinion. This is important in large-scale argumentation frameworks involving humans. As observed in [13], we cannot expect that humans express their opinions about all the arguments involved in a discussion, since they tend to focus on the arguments of their interest. Formally:

Definition 5 (Argument labelling) Let $TODF = \langle \mathcal{A}, \mapsto, \Vdash, \tau \rangle$ be a target-oriented discussion framework. An argument labelling for $TODF$ is a function $L : \mathcal{A} \rightarrow \{\text{in}, \text{out}, \text{undec}\}$ that maps each argument of \mathcal{A} to one out of the following labels: *in* (accepted), *out* (rejected), or *undec* (undecidable).

We note as $Ag = \{ag_1, \dots, ag_n\}$ the set of agents taking part in a $TODF$, and as L_i the labelling encoding the opinion of agent $ag_i \in Ag$. We will put together the opinions of all the agents participating in an argumentation as follows.

Definition 6 (Labelling profile) Let L_1, \dots, L_n be argument labellings of the agents in Ag , where L_i is the argument labelling of agent ag_i . A labelling profile is a tuple $\mathcal{L} = (L_1, \dots, L_n)$.

Example 3 (Flatmates’ opinions) Figure 2(b) graphically depicts Alan’s, Barbara’s, and Charles’ labellings (noted as L_1, L_2, L_3 respectively), each one appearing next to the corresponding arguments in the $TODF$ ’s graphical representation in Figure 2(a).

3.3 Coherent argument labellings

As noted in [4], there are multiple reasonable ways in which an agent may evaluate an argument structure through a labelling. There, authors introduce the notion of *complete* labelling⁴. Here we argue that the conditions required by a complete labelling are very restrictive. Instead, we will consider alternative, more relaxed conditions for an argument labelling (an opinion) to be reasonable. With this aim, for each argument a we will compare the labelling over the argument, what we consider to be its *direct opinion*, with the aggregated labellings over its children or immediate descendants, namely, its *indirect opinion*.

Considering the example in Figure 2(b), if we take any argument, such as for instance N , we consider its associated labels as the direct opinion, whereas we think of

⁴ A complete labelling requires that: an argument is labelled *in* iff all its defeaters are labelled *out*; and an argument is labelled *out* iff at least one of its defeaters is accepted.

the labels associated to its immediate descendants a_1 , a_2 , and a_3 as its indirect opinion. Analogously, the direct opinion on argument a_2 corresponds to its associated labels, whereas the labels associated to a_1 , its single immediate descendant, constitute its indirect opinion.

Thus, informally, we will say that the labelling of an argument is coherent if its direct opinion is in line with its indirect opinion. This will occur when the *majority* of arguments in the indirect opinion of an argument agree with the labelling of the argument. In what follows, we formalise our notion of coherent labelling.

First, given an argument a we will define its set of attacking arguments $A(a) = \{b \in \mathcal{A} \mid b \mapsto a\}$; and its set of defending arguments $D(a) = \{c \in \mathcal{A} \mid c \Vdash a\}$. Thus, the labelling of the arguments in $A(a) \cup D(a)$ compose the indirect opinion on a .

Given an argument labelling L and a set of arguments $S \subseteq \mathcal{A}$, we can quantify the number of accepted arguments in S as $\text{in}_L(S) = |\{b \in S \mid L(b) = \text{in}\}|$. Analogously, we can also quantify the number of rejected arguments in S as $\text{out}_L(S) = |\{b \in S \mid L(b) = \text{out}\}|$. Thus, given an argument a , we can readily quantify its accepted and rejected defending arguments as $\text{in}_L(D(a))$ and $\text{out}_L(D(a))$ respectively. Moreover, we can also quantify its accepted and rejected attacking arguments as $\text{in}_L(A(a))$ and $\text{out}_L(A(a))$ respectively. Now we are ready to measure the *positive* and *negative support* contained in the indirect opinion of a given argument as follows.

Definition 7 (Positive support) *Let $a \in \mathcal{A}$ be an argument and L a labelling on \mathcal{A} . We define the positive (pro) support of a as: $\text{Pro}_L(a) = \text{in}_L(D(a)) + \text{out}_L(A(a))$. If $\text{Pro}_L(a) = |A(a) \cup D(a)|$ we say that a receives full positive support from L .*

Definition 8 (Negative support) *Let $a \in \mathcal{A}$ be an argument and L a labelling on \mathcal{A} . We define the negative (con) support of a as: $\text{Con}_L(a) = \text{in}_L(A(a)) + \text{out}_L(D(a))$. If $\text{Con}_L(a) = |A(a) \cup D(a)|$ we say that a receives full negative support from L .*

Notice that the positive support of an argument combines the strength of its accepted defending arguments with the weakness of its rejected attacking arguments in the argument's indirect opinion. As a dual concept, the negative support combines accepted attacking arguments with rejected defending arguments.

We now introduce our notion of coherence by combining the positive and negative support of an argument. We say that a labelling is coherent if the following conditions hold for each argument: (1) if an argument is labelled accepted (*in*) then it cannot have more negative than positive support (the majority of its indirect opinion supports the argument); and (2) if an argument is labelled rejected (*out*) then it cannot have more positive than negative support (the majority of its indirect opinion rejects the argument).

Definition 9 (Coherence) *Given a $TODF = \langle \mathcal{A}, \mapsto, \Vdash, \tau \rangle$, a coherent labelling is a total function $L : \mathcal{A} \rightarrow \{\text{in}, \text{out}, \text{undec}\}$ such that for all $a \in \mathcal{A}$ with $A(a) \cup D(a) \neq \emptyset$: (1) if $L(a) = \text{in}$ then $\text{Pro}_L(a) \geq \text{Con}_L(a)$; and (2) if $L(a) = \text{out}$ then $\text{Pro}_L(a) \leq \text{Con}_L(a)$.*

Finally, we offer a more refined version of coherence based on the difference between positive and negative supports.

Definition 10 (c-Coherence) Let $TODF = \langle \mathcal{A}, \mapsto, \Vdash, \tau \rangle$ be a target-oriented discussion framework. A c -coherent labelling for some $c \in \mathbb{N}$ is a total function $L : \mathcal{A} \rightarrow \{\text{in}, \text{out}, \text{undec}\}$ such that for all $a \in \mathcal{A}$ with $A(a) \cup D(a) \neq \emptyset$: (i) if $L(a) = \text{in}$ then $Pro_L(a) > Con_L(a) + c$; (ii) if $L(a) = \text{out}$ then $Pro_L(a) + c < Con_L(a)$; and (iii) if $L(a) = \text{undec}$ then $|Pro_L(a) - Con_L(a)| \leq c$.

Let $TODF$ be a target oriented discussion framework. We will note the class of all the argument labellings of $TODF$ as $\mathbf{L}(TODF)$, the subclass of coherent argument labellings as $Coh(TODF)$, and the subclass of c -coherent argument labellings as $Coh_c(TODF)$ for some $c \in \mathbb{N}$.

Example 4 Again, considering our example, and its labellings from Figure 2(b) (L_1, L_2, L_3 in $\mathbf{L}(TODF)$), we note that just L_1, L_2 belong to the subclass of its coherent argument labellings $Coh(TODF)$. Moreover, L_1 and L_2 are 0-coherent.

4 The aggregation problem

Recall that our aim is to have multiple agents jointly decide whether to accept a target (e.g. a norm) or not. In section 4.1 we pose such problem as a *judgement aggregation* [15] problem in the context of argumentation: a set of agents collectively decide how to label a target-oriented argumentation framework, and such collective labelling provides a label for the target. Since there are many ways of aggregating labellings, following [4], section 4.2 states that such aggregation must guarantee that the outcome is *fair*.

4.1 Collective labelling

First, a discussion problem will encompass a target-oriented discussion framework together with a set of agents' individual labellings.

Definition 11 (Labelling discussion problem) Let $Ag = \{ag_1, \dots, ag_n\}$ be a finite non-empty set of agents, and $TODF = \langle \mathcal{A}, \mapsto, \Vdash, \tau \rangle$ be a target-oriented discussion framework. A labelling discussion problem is a pair $\mathcal{LDP} = \langle Ag, TODF \rangle$.

Given an \mathcal{LDP} , our aim is to find how to aggregate the individuals' labellings into a single labelling that captures the opinion of the collective.

Definition 12 (Aggregation function) An aggregation function for a labelling discussion problem $\mathcal{LDP} = \langle Ag, TODF \rangle$ is a function $F : \mathbf{L}(TODF)^n \rightarrow \mathbf{L}(TODF)$.

Plainly, an aggregation function F takes a labelling profile representing all agents' opinions and yields a single labelling computed from the individual labellings. Such aggregation function is key to assessing the collective decision over the target.

Definition 13 (Decision over a target) Let $\mathcal{LDP} = \langle Ag, TODF \rangle$ be a labelling discussion problem, \mathcal{L} a labelling profile, and F an aggregation function for the \mathcal{LDP} . The decision over the target of the $TODF$ is the label $F(\mathcal{L})(\tau)$.

4.2 Desirable properties of an aggregation function

The literature on Social Choice theory has already identified fair ways of aggregating votes. These can be translated into formal properties that an aggregation function is required to satisfy [9]. Based on [4], here we formally state the desirable properties for an aggregation function that allows to assess the decision over the target of a target-oriented discussion framework. First, notice that an aggregation function may not compute over every labelling profile, so we start by referring the domain properties of an aggregate function.

Exhaustive Domain (ED) F can take as input all labelling profiles, i.e., all $\mathcal{L} \in \mathbf{L}(TODF)^n$.

Coherent Domain (CD) F can take as input all the coherent labelling profiles, $\mathcal{L} \in Coh(TODF)^n$.

Furthermore, it is natural to require that aggregation outcomes are also coherent, namely, that the aggregation results in *collective coherence*.

Collective coherence (CC) $F(\mathcal{L}) \in Coh(TODF)$ for all $\mathcal{L} \in \mathbf{L}(TODF)^n$.

Collective coherence is our most desired property. Notice that if an aggregation function does not produce a coherent labelling, there is at least some argument whose collective label (direct opinion) is in contradiction with its indirect opinion. Thus, the resulting aggregation would not be reliable.

Notice also that the agents involved in a discussion expect that their opinions are as important as others'. This idea is captured by the anonymity property, where all opinions are equally significant.

Anonymity (A) If $\mathcal{L} = (L_1, \dots, L_n)$ is a labelling profile and σ is a permutation over Ag then: if $\mathcal{L}' = (L_{\sigma(1)}, \dots, L_{\sigma(n)})$ then $F(\mathcal{L}) = F(\mathcal{L}')$.

A weaker version of anonymity, non-dictatorship, states that no agent can decide over the others, like a dictator. Notice that, this directly follows from anonymity.

Non-Dictatorship (ND) There is no agent $ag_i \in Ag$ such that, for every labelling profile \mathcal{L} we have $F(\mathcal{L}) = L_i$.

Regarding unanimity, we shall consider two main notions of unanimity: direct and endorsed. On the one hand, we formulate the *direct unanimity* property to capture the following requirement: if all the agents agree (share the opinion) on one argument, then the aggregate opinion must reflect such agreement.

Direct Unanimity (DU) Let $l \in \{\text{in}, \text{undec}, \text{out}\}$. For each $a \in \mathcal{A}$ such that $L_i(a) = l$ for all $ag_i \in Ag$, then $F(\mathcal{L})(a) = l$.

Endorsed unanimity is a variant of direct unanimity: for each argument, if all the agents agree on the indirect opinion of an argument (be it to give it full positive support or full negative support), this cannot contradict the aggregated opinion on the argument.

Endorsed Unanimity (EU) Let \mathcal{L} be a labelling profile. For each $a \in A$: (i) if a receives full positive support for all $L_i \in \mathcal{L}$ then $F(\mathcal{L})(a) = \text{in}$; and (ii) if a receives full negative support for all $L_i \in \mathcal{L}$ then $F(\mathcal{L})(a) = \text{out}$.

As an additional variant of unanimity, we consider supportiveness: the aggregated opinion on an argument cannot be set to a label $l \in \{\text{in}, \text{out}, \text{undec}\}$ unless at least one agent labels the argument with l .

Supportiveness (S) Let \mathcal{L} be a labelling profile. For all $a \in A$, there exists some agent $ag_i \in Ag$ such that $F(\mathcal{L})(a) = L_i(a)$.

Finally, we state a novel notion of monotonicity, the so-called *familiar monotonicity*, which considers the opinions of an argument's descendants. Intuitively, our notion of familiar monotonicity captures the following principle: if the support for an argument increases, the collective labeling of the argument should remain the same, but provided that the opinions on the argument's descendants do not change. The latter condition is necessary because changes in the opinions about the descendants of the argument may affect the support on the argument. In other words, our notion of monotonicity, unlike the notion of monotonicity presented in [4], is aware of the dependencies between arguments.

We also formulate a weaker version of familiar monotonicity that only applies to **in** and **out**.

Familiar Monotonicity (FM) Let $a \in \mathcal{A}$ be an argument and two labelling profiles $\mathcal{L} = (L_1, \dots, L_i, \dots, L_{i+k}, \dots, L_n)$, $\mathcal{L}' = (L_1, \dots, L'_i, \dots, L'_{i+k}, \dots, L_n)$ such that $F(\mathcal{L})(a) = l \in \{\text{in}, \text{out}, \text{undec}\}$, agents ag_i, \dots, ag_{i+k} only differing on their labellings of a (namely, for all b descendant of a , $L_j(b) = L'_j(b)$ for every $j \in \{i, \dots, i+k\}$) and $L(a)_j \neq L'(a)_j = l$ for $j \in \{i, \dots, i+k\}$, if $F(\mathcal{L})(a) = l$ then $F(\mathcal{L}')(a) = l$.

The next property establishes the same idea considering only the cases where the previous aggregate opinion is either **in** or **out**, not **undec**.

in/out-Familiar Monotonicity (i/o-FM) Let $a \in \mathcal{A}$ be an argument and two labelling profiles \mathcal{L} and \mathcal{L}' satisfying the previous hypothesis of the familiar monotonicity property adding that $F(\mathcal{L})(a) = l \neq \text{undec}$. Then, if $F(\mathcal{L})(a) = l$ then $F(\mathcal{L}')(a) = l$.

Some other properties that are desirable in other multi-agent argumentation contexts (e.g. [4]) are not desirable here. In particular, systematicity and independence are not desirable because we want to exploit dependence relationships between arguments.

5 The coherent aggregation function

Next we define an aggregation function to compute the collective labelling, and hence the decision over a target, for a labelling discussion problem. Section 5.1 introduces our function, while Section 5.2 analyses the satisfaction of the properties in Section 4.2.

5.1 Defining the coherent aggregation function

First, we introduce notation to quantify the direct positive and negative support of an argument. Let $\mathcal{L} = (L_1, \dots, L_n)$ be a labelling profile and a an argument. We note the *direct positive support* of a as $\text{in}_{\mathcal{L}}(a) = |\{ag_i \in Ag \mid L_i(a) = \text{in}\}|$; and its *direct negative support* as $\text{out}_{\mathcal{L}}(a) = |\{ag_i \in Ag \mid L_i(a) = \text{out}\}|$. Next, we define our chosen aggregation function: the *coherent aggregation function*. The main purpose of this function is to compute a coherent aggregated labelling, and hence fulfil the collective coherence property. Notice that we consider that the most important desirable property for an aggregation function is to yield a rational outcome that is free of contradiction.

Definition 14 (Coherent aggregation function) *Let \mathcal{L} be a labelling profile. For each argument a the coherent function over \mathcal{L} is defined as:*

$$CF(\mathcal{L})(a) = \begin{cases} \text{in} & , IO(\mathcal{L})(a) + DO(\mathcal{L})(a) > 0 \\ \text{out} & , IO(\mathcal{L})(a) + DO(\mathcal{L})(a) < 0 \\ \text{undec} & , IO(\mathcal{L})(a) + DO(\mathcal{L})(a) = 0 \end{cases}$$

where the functions *IO* (indirect opinion) and *DO* (direct opinion) are defined as:

$$IO(\mathcal{L})(a) = \begin{cases} 1 & , Pro_{CF(\mathcal{L})}(a) > Con_{CF(\mathcal{L})}(a) \\ 0 & , Pro_{CF(\mathcal{L})}(a) = Con_{CF(\mathcal{L})}(a) \\ -1 & , Pro_{CF(\mathcal{L})}(a) < Con_{CF(\mathcal{L})}(a) \end{cases}$$

$$DO(\mathcal{L})(a) = \begin{cases} 1 & , \text{in}_{\mathcal{L}}(a) > \text{out}_{\mathcal{L}}(a) \\ 0 & , \text{in}_{\mathcal{L}}(a) = \text{out}_{\mathcal{L}}(a) \\ -1 & , \text{in}_{\mathcal{L}}(a) < \text{out}_{\mathcal{L}}(a) \end{cases}$$

Example 5 (Flatmates' discussion) *Back to our example involving a flatmates' discussion, we use the coherent aggregation function to obtain the aggregated opinion of the provided labellings (see Figure 2(b)). Figure 3 shows the results of the aggregation and the decision over the target as produced by CF. We observe that the flatmates collectively accept arguments a_1 and a_3 , whereas argument a_2 becomes undecidable. Finally, the decision over the norm is to accept it.*

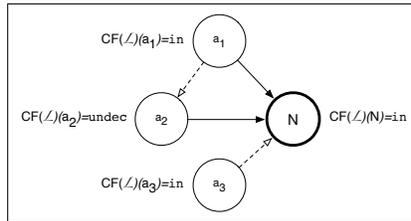


Fig. 3. Flatmates example: aggregated labellings (and decision over target N) computed by *CF*.

5.2 Analysing the coherent aggregation function

Below we analyse the desired properties from section 4.2 that our aggregation function CF fulfils.

The first two results about domain properties follow from CF 's definition.

Proposition 15 *CF satisfies the exhaustive domain property.*

Corollary 16 *CF satisfies the coherent domain property.*

Proof. It is clear that CF is defined for all labelling profiles, and hence it is also defined for every coherent labelling $\mathcal{L} \in Coh(TODF)^n$.

Notice that, since CF is defined for all labelling profiles, it is also defined for labelling profiles in $Coh_c(TODF)^n$, namely for labelling profiles whose argument labellings are c -coherent. Now recall that we designed our CF function to satisfy the collective coherence property. Thus, the following property naturally follows.

Proposition 17 *CF satisfies the collective coherence property.*

Proof. Let a be an argument such that $CF(\mathcal{L})(a) = \text{in}$. From Definition 14 we know that $IO(\mathcal{L})(a) + DO(\mathcal{L})(a) > 0$. Thus, there are three possibilities: (i) $DO(\mathcal{L})(a) = 1$ and $IO(\mathcal{L})(a) = 1$; (ii) $DO(\mathcal{L})(a) = 1$ and $IO(\mathcal{L})(a) = 0$; or (iii) $IO(\mathcal{L})(a) = 0$ and $DO(\mathcal{L})(a) = 1$. Since $IO(\mathcal{L})(a) \geq 0$ in all cases, this implies that $Pro_{CF(\mathcal{L})}(a) \geq Con_{CF(\mathcal{L})}(a)$, and hence CF satisfies the coherence property. The proof goes analogously for the case $CF(\mathcal{L})(a) = \text{out}$.

Now we turn our attention into the anonymity property and its weaker version: the non-dictatorship property.

Proposition 18 *CF satisfies the anonymity property.*

Proof. Let $\mathcal{L} = (L_1, \dots, L_n)$ be a labelling profile and σ a permutation over Ag such that $\mathcal{L}' = (L_{\sigma(1)}, \dots, L_{\sigma(n)})$. Since CF only uses the number of elements, there is no dependency on the identities of the agents' labellings. We only have to check that $DO(\mathcal{L}) = DO(\mathcal{L}')$ because functions IO , Pro , and Con only depend on $CF(\mathcal{L})$, and hence in turn they will not depend either on the identities of the agents' labellings. This amounts to checking whether $\text{in}_{\mathcal{L}}(a) = \text{in}_{\sigma(\mathcal{L})}(a)$ and $\text{out}_{\mathcal{L}}(a) = \text{out}_{\sigma(\mathcal{L})}(a)$ hold. Indeed, on the one hand $\text{in}_{\mathcal{L}}(a) = |\{ag_i \in Ag \mid L_i(a) = \text{in}\}| = |\{\sigma(ag_i) \in Ag \mid L_{\sigma(i)}(a) = \text{in}\}| = \text{in}_{\sigma(\mathcal{L})}(a)$. Moreover, $\text{out}_{\mathcal{L}}(a) = |\{ag_i \in Ag \mid L_i(a) = \text{out}\}| = |\{\sigma(ag_i) \in Ag \mid L_{\sigma(i)}(a) = \text{out}\}| = \text{out}_{\sigma(\mathcal{L})}(a)$.

Since CF satisfies anonymity, the identity of which agent submits which labelling is irrelevant. Furthermore, recall from Section 4.2 that non-dictatorship follows.

Corollary 19 *CF satisfies the non-dictatorship property.*

Next, we focus on unanimity properties. First, we will show that CF fulfils the endorsed unanimity property. With this aim, we will introduce an additional hypothesis based on the following lemma.

Lemma 1. *Let $TODF$ be a target-oriented discussion framework, \mathcal{L} a 0-coherent labelling profile ($\mathcal{L} \in Coh_0(TODF)^n$), a an argument in \mathcal{A} , and m the number of immediate descendants of a ($m = |A(a) \cup D(a)|$). If $Pro_{L_i}(a) = m$ for all $i \in \{1, \dots, n\}$ then $in_{\mathcal{L}}(a) = n$; and if $Con_{L_i}(a) = m$ for all $i \in \{1, \dots, n\}$ then $out_{\mathcal{L}}(a) = n$.*

Proof. We next prove that if $Pro_{L_i}(a) = m$ then $L_i(a) = in$, for all $i \in \{1, \dots, n\}$. Thus, all the agents label argument a as *in*, i.e. $in_{\mathcal{L}}(a) = n$. Since we assume that each L_i is 0-coherent, a 's label can be neither *out*, because $Pro_{L_i}(a) \not\leq Con_{L_i}(a)$, nor *undec*, because $Pro_{L_i}(a) \neq Con_{L_i}(a)$. Thus, the only option is that a is labelled as *in*. The proof runs analogously when considering the case $Con_{L_i}(a) = m$.

Plainly, the lemma says that, when assuming 0-coherence, if the indirect opinion on an argument is unanimous, the direct opinion on the argument will also be unanimous. Using this lemma we can prove the following result.

Proposition 20 *Let $\mathcal{L} = (L_1, \dots, L_n)$ be a labelling profile. If every $L_i, i \in \{1, \dots, n\}$, satisfies the 0-coherence property, then CF satisfies the endorsed unanimity property.*

Proof. We focus on the case for which if each argument a receives full positive support for all $L_i \in \mathcal{L}$, namely $Pro_{L_i}(a) = m$ for every i , then $CF(\mathcal{L})(a) = in$. First of all, we will analyse the aggregated indirect opinion on a given argument a . Let b be a defending argument of a , namely $b \in D(a)$. Since $Pro_{L_i}(a) = m$ for all i , $L_i(b) = in$. Since we do not know the labellings of the immediate descendants of b , we can assume that $DO(\mathcal{L})(b) = 1$, and therefore either $CF(\mathcal{L})(b) = undec$ or $CF(\mathcal{L})(b) = in$. Following a similar reasoning, we observe that if $b \in A(a)$, then either $CF(\mathcal{L})(b) = undec$ or $CF(\mathcal{L})(b) = out$. Therefore, we have that $IO(\mathcal{L})(a) \geq 0$. Because $Pro_{L_i}(a) = m$ and L_i is 0-coherent for every agent ag_i , we have that $n = in_{\mathcal{L}}(a) > out_{\mathcal{L}}(a) = 0$ by lemma 1, and hence $DO(\mathcal{L})(a) = 1$. Since $IO(\mathcal{L})(a) \geq 0$, we finally have that $CF(\mathcal{L})(a) = in$. Analogously, we can also prove that $CF(\mathcal{L})(a) = out$ if $Con_{L_i}(a) = m$ for every $ag_i \in Ag$.

Notice however that CF does not satisfy the other two unanimity properties presented in section 4.2, namely direct unanimity and supportiveness.

Proposition 21 *Neither direct unanimity nor supportiveness are satisfied by CF .*

Proof. Figure 6(a) graphically represents a $TODF$ that will serve to illustrate our proposition. Our $TODF$ contains a target argument $\tau = a$, which is defended by five other arguments $\{a_1, a_2, a_3, a_4, a_5\}$. The $TODF$ involves the argument labellings of three agents, noted as L_1, L_2 , and L_3 : (1) agent 1 accepts (labels with *in*) arguments a, a_1, a_2 , and a_3 , and refuses (labels with *out*) arguments a_4 and a_5 ; (2) Agent 2 accepts arguments a, a_1, a_2 , and a_4 , and rejects arguments a_3 , and a_5 ; and agent 3 accepts arguments a, a_1, a_2 , and a_5 , and rejects arguments a_3 , and a_4 .

Notice that the three agents agree on accepting the target ($L_1(a) = in, L_2(a) = in$, and $L_3(a) = in$), and hence there is unanimous opinion on a .

Figure 6(b) depicts the resulting labelling when computing the CF function for this $TODF$ over the labelling profile $\mathcal{L} = (L_1, L_2, L_3)$. Since arguments a_1 and a_2 are

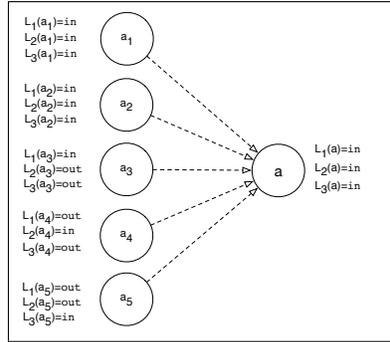


Fig. 4. Counterexample for direct unanimity and supportiveness: argument labellings.

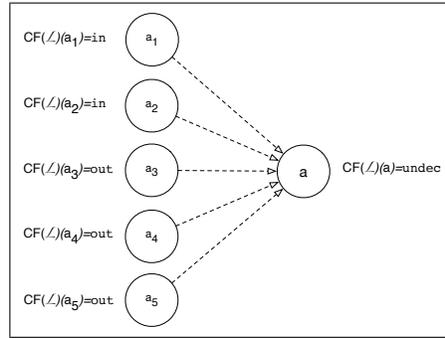


Fig. 5. Counterexample for direct unanimity and supportiveness: CF result.

collectively accepted ($CF(\mathcal{L})(a_1) = \text{in}$, $CF(\mathcal{L})(a_2) = \text{in}$) and arguments a_3, a_4 , and a_5 are rejected ($CF(\mathcal{L})(a_3) = \text{out}$, $CF(\mathcal{L})(a_4) = \text{out}$, $CF(\mathcal{L})(a_5) = \text{out}$), the target is neither accepted nor rejected ($CF(\mathcal{L})(a) = \text{undec}$). Thus, although the three agents agree on accepting a , the collective decision obtained by CF is undec . Therefore, CF does not satisfy direct unanimity.

As to supportiveness, it does not hold either. Observe that although the aggregate label of a is undec , no agent has labelled argument a as undec .

Finally, we study CF 's monotonicity. Although familiar monotonicity does not hold for CF , its weaker version, in/out -familiar, does hold.

Proposition 22 *CF does not satisfy the familiar monotonicity property.*

Proof. Our proof only requires a simple $TODF$ with a target argument a and two labelling profiles with two argument labellings. Let $\mathcal{L} = (L_1, L_2)$ and $\mathcal{L}' = (L_1, L'_2)$ where $L_1(a) = \text{in}$, $L_2(a) = \text{out}$ and $L'_2(a) = \text{undec}$. These two labelling profiles satisfy the hypothesis required by the familiar monotonicity property. Nonetheless, notice now that the aggregate labellings on the target obtained for each labelling are: $CF_{\mathcal{L}}(a) = \text{undec}$ and $CF_{\mathcal{L}'}(a) = \text{in}$. Since $CF_{\mathcal{L}}(a) \neq CF_{\mathcal{L}'}(a)$, familiar monotonicity does not hold.

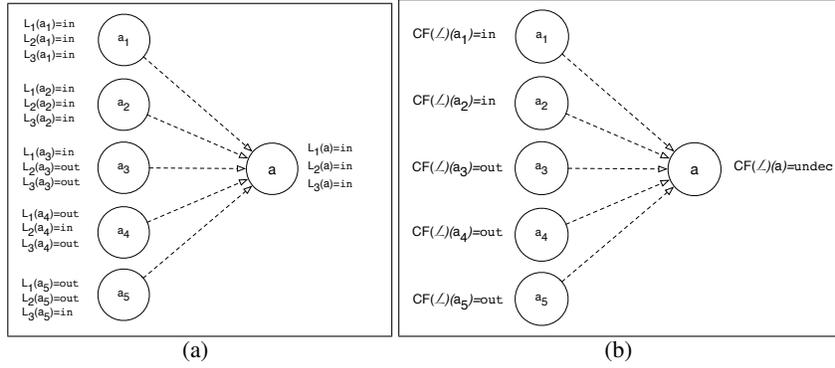


Fig. 6. Counterexamples to illustrate lack of: (a) direct unanimity and supportiveness (argument labellings); (b) direct unanimity and supportiveness (result of the coherent aggregation function).

Proposition 23 *CF satisfies the in/out-familiar monotonicity property.*

Proof. Let $\mathcal{L}, \mathcal{L}'$ be two labelling profiles satisfying the hypothesis required by the in/out-familiar monotonicity property on the argument a , and whose collective label on a for \mathcal{L} is $CF(\mathcal{L})(a) = l = in$. Since $L_j(b) = L'_j(b)$ for all b descendant of a , we know that $IO(\mathcal{L})(a) = IO(\mathcal{L}')(a)$ because IO only depends on the descendants. Since $CF(\mathcal{L})(a) = in$, we have that $DO(\mathcal{L})(a) \geq 0$. Now, because $in_{\mathcal{L}}(a) \leq in_{\mathcal{L}'}(a)$ and $out_{\mathcal{L}}(a) \geq out_{\mathcal{L}'}(a)$, we know that $DO(\mathcal{L}')(a) \geq DO(\mathcal{L})(a) \geq 0$. From this follows that $DO(\mathcal{L}')(a) + IO(\mathcal{L})(a) \geq DO(\mathcal{L}')(a) + IO(\mathcal{L})(a) = 1$, and hence $CF(\mathcal{L}')(a) = in$. We can analogously check the case $CF(\mathcal{L})(a) = out$.

Analysis. First, notice that CF satisfies a significant number of the desirable properties identified in section 4.2 for *any* sort of labelling profiles. This means that CF does not constrain at all an agent's labelling, and hence even can cope with the inconsistencies of agents' opinions. This is not the case though for endorsed unanimity. This property is constrained to labellings that are 0-coherent. Second, properties such as direct unanimity and supportiveness, which are not satisfied by CF , assume that aggregation is computed independently for each argument. In other words, they serve to analyse the behaviour of an aggregation function in a single argument. Since in this paper we pursue to exploit dependencies between arguments within a discussion, such properties prevent us from making a more informed decision about the discussion target.

5.3 Computing the decision over a target

Given a target-oriented discussion framework $TODF = \langle \mathcal{A}, \mapsto, \Vdash, \tau \rangle$ shared by the agents in Ag , a labelling profile \mathcal{L} , and our coherent aggregation function CF , we now consider how to compute the collective label assigned to the target τ , namely $CF(\mathcal{L})(\tau)$. Such computation is based on the following observation:

Since the graph associated to a target-oriented discussion framework $TODF$ is a directed acyclic graph (DAG), we can embed the computation of the collective label of each argument in \mathcal{A} within its the traversal. of its associated graph. Such a graph traversal could be performed by its topological sorting [11]. Therefore, the running time required to compute $CF(\mathcal{L})(\tau)$ is linear in the number of arguments plus the number of edges, asymptotically, namely $O(|\mathcal{A}| + |\mapsto| + ||\vdash|)$. Algorithm 1 shows the pseudo-code of function COMPUTETARGETDECISION, which returns the collective label of target τ for an input graph G_{TODF} and a labelling profile \mathcal{L} .

Algorithm 1 Algorithm to compute the collective label of a target

```

1: function COMPUTETARGETDECISION( $G_{TODF}, \mathcal{L}, \tau$ )
2:    $ToVisit \leftarrow \{a \in A \mid A(a) \cup D(a) = \emptyset\}$   $\triangleright$  Arguments with neither attacks nor defences
   (no descendants)
3:   while  $ToVisit \neq \emptyset$  do
4:     remove an argument  $b$  from  $ToVisit$ 
5:     compute  $CF(\mathcal{L})(b)$ 
6:     for each node  $c$  with an edge  $(b, c) \in G_{TODF}$  do
7:       remove edge  $(b, c)$  from graph  $G_{TODF}$ 
8:       if  $c$  has no other incoming edges then
9:         insert  $c$  into  $ToVisit$ 
10:  return  $CF(\mathcal{L})(\tau)$   $\triangleright$  Return collective label for target  $\tau$ 

```

6 Conclusions and future work

Along this paper we have formalised the problem of taking a collective decision over a target. We claim this problem can be tackled within a target oriented decision framework, and we have tailored it for humans, due to the increasing interest on e-participation, e-governance and open innovation systems. Within this framework, we have also proposed a coherent aggregation function that combines participants opinions and has proven to satisfy valuable social choice properties without any additional assumption (with the exception of endorsed unanimity, which just requires the labelling profile to be 0-coherent). When considering humans, we hypothesise that the larger the number of people, the less the number of undecidable labels will result from combining their opinions, and thus, the less unlikely will be the occurrence of an undecidable outcome (i.e., a target collective decision).

Finally, notice that although our argumentation framework shares the use of attack and defense relations with bipolar argumentation frameworks [8,3], there are notable differences. First, bipolar argumentation does not consider labellings (different opinions on arguments), neither their aggregation. Second, bipolar argumentation focuses on studying the structure between arguments and groups of arguments, whereas we focus on computing a collective decision from differing opinions about arguments. Third, arguments in bipolar argumentation can be regarded as objective facts, while in our case, arguments can be subjective facts on which individuals can differ. Thus, our argumentation framework pursues to be less restrictive to include humans in the loop.

As to future work, we plan to consider alternative semantics for arguments' pros and cons that diminish the relevance associated to the rejection of arguments. Furthermore, we plan to extend our TODF to allow loops, and hence ease rebuttal, a common feature of argumentation systems. Finally, we plan to provide more fine-grained means of computing argument support.

References

1. City of Barcelona participation portal. <https://decidim.barcelona>, 2016.
2. City of Reykjavik participation portal. <http://reykjavik.is/en/participation>, 2016.
3. Leila Amgoud, Claudette Cayrol, Marie-Christine Lagasquie-Schiex, and Pierre Livet. On bipolarity in argumentation frameworks. *International Journal of Intelligent Systems*, 23(10):1062–1093, 2008.
4. Edmond Awad, Richard Booth, Fernando Tohmé, and Iyad Rahwan. Judgement aggregation in multi-agent argumentation. *Journal of Logic and Computation*, 2015.
5. Trevor J. M. Bench-Capon and Paul E. Dunne. Argumentation in artificial intelligence. *Artif. Intell.*, 171(10-15):619–641, 2007.
6. Martin Caminada. On the issue of reinstatement in argumentation. In *Logics in artificial intelligence*, pages 111–123. Springer, 2006.
7. Martin WA Caminada and Dov M Gabbay. A logical account of formal argumentation. *Studia Logica*, 93(2-3):109–145, 2009.
8. Claudette Cayrol and Marie-Christine Lagasquie-Schiex. On the acceptability of arguments in bipolar argumentation frameworks. In *Symbolic and quantitative approaches to reasoning with uncertainty*, pages 378–389. Springer, 2005.
9. Franz Dietrich. A generalised model of judgment aggregation. *Social Choice and Welfare*, 28(4):529–565, 2007.
10. Simone Gabbriellini and Paolo Torroni. Microdebates: Structuring debates without a structuring tool1. *AI Commun.*, 29(1):31–51, 2015.
11. Arthur B Kahn. Topological sorting of large networks. *Communications of the ACM*, 5(11):558–562, 1962.
12. Mark Klein. Enabling large-scale deliberation using attention-mediation metrics. *Computer Supported Cooperative Work*, 21(4-5):449–473, 2012.
13. Mark Klein. Enabling large-scale deliberation using attention-mediation metrics. *Computer Supported Cooperative Work (CSCW)*, 21(4-5):449–473, 2012.
14. Mark Klein and Gregorio Convertino. A roadmap for open innovation systems. *Journal of Social Media for Organizations*, 2(1):1, 2015.
15. Christian List and Philip Pettit. Aggregating sets of judgments: An impossibility result. *Economics and Philosophy*, 18(01):89–110, 2002.
16. Iyad Rahwan, Guillermo R. Simari, and Johan Benthem. *Argumentation in artificial intelligence*, volume 47. Springer, 2009.
17. Vishanth Weerakkody and Christopher G Reddick. *Public sector transformation through e-government: experiences from Europe and North America*. Routledge, 2012.