



Online Handbook of Argumentation for AI

Volume 3

Edited by

Federico Castagna
Jack Mumford
Ştefan Sarkadi
Andreas Xydis

December 2022

Preface

This volume contains revised versions of the papers selected for the third volume of the Online Handbook of Argumentation for AI (OHAAI). Previously, formal theories of argument and argument interaction have been proposed and studied, and this has led to the more recent study of computational models of argument. Argumentation, as a field within artificial intelligence (AI), is highly relevant for researchers interested in symbolic representations of knowledge and defeasible reasoning. The purpose of this handbook is to provide an open access and curated anthology for the argumentation research community. OHAAI is designed to serve as a research hub to keep track of the latest and upcoming PhD-driven research on the theory and application of argumentation in all areas related to AI. The handbook's goals are to:

1. Encourage collaboration and knowledge discovery between members of the argumentation community.
2. Provide a platform for PhD students to have their work published in a citable peer-reviewed venue.
3. Present an indication of the scope and quality of PhD research involving argumentation for AI.

The papers in this volume are those selected for inclusion in OHAAI Vol.3 following a back-and-forth peer-review process. The volume thus presents a strong representation of the contemporary state of the art research of argumentation in AI that has been strictly undertaken during PhD studies. Papers in this volume are listed alphabetically by author. We hope that you will enjoy reading this handbook.

Editors

Federico Castagna
Jack Mumford
Ştefan Sarkadi
Andreas Xydis

December 2022

Acknowledgements

We thank the senior researchers in the area of Argumentation and Artificial Intelligence for their efforts in spreading the word about the OHAAI project with early-career researchers.

We are especially thankful to Sylwia Polberg and the COMMA 2022 organisers for collaborating with us at the conference summer school SSA 2022, during which OHAAI ran a student programme that provided the attending students the opportunity to present and discuss their PhD research amongst peers and academics in a friendly forum.

We are also grateful to ArXiv for their willingness to publish this handbook.

Our sincere gratitude to Costanza Hardouin for her fantastic work in designing the OHAAI logo.

We owe many thanks to Sanjay Modgil for helping to form the motivation for the handbook, and to Elizabeth Black and Simon Parsons for their advice and guidance that enabled the OHAAI project to come to fruition.

The success of the OHAAI project depends upon the quality feedback provided by our reviewers. We have been fortunate in securing a diligent and thoughtful program committee that produced reviews of a reliably high standard. Our thanks go to our PC:

Andreas Brännström, Théo Duchatelle, Timotheus Kampik, Isabelle Kuhlmann, Lars Malmqvist, Mariela Morveli-Espinoza, Stipe Pandžić, Christos Rodosthenous, Robin Schaefer, Kenneth Skiba, Luke Thorburn, Antonio Yuste-Ginel, and Heng Zheng.

Special thanks must go to the contributing authors:

Lars Bengel, Elfia Bezou-Vrakatseli, Lydia Blümel, Federico Castagna, Giulia D'Agostino, Daphne Odekerken, Minal Suresh Patil, Jordan Robinson, Hao Wu, and Andreas Xydis. Thank you for making the world of argumentation greater!

Contents

On Serialisability for Argumentative Explanations <i>Lars Bengel</i>	2
Debating Ethics: Using Argumentation to Support Dialogue <i>Elfia Bezou-Vrakatseli</i>	7
Characterization of Unresolvable Conflicts in Abstract Argumentation <i>Lydia Blümel</i>	12
Towards a Fully-fledged Formal Protocol for the Explanation-Question-Response Dialogue <i>Federico Castagna</i>	17
Argumentation without Opposition? <i>Giulia D'Agostino</i>	22
Justification, Stability and Relevance for Transparent and Efficient Human-in-the-Loop Decision Support <i>Daphne Odekerken</i>	28
Towards Preserving Semantic Structure in Argumentative Multi-Agent via Abstract Interpretation <i>Minal Suresh Patil</i>	33
Distributed Hypothesis Generation and Evaluation <i>Jordan Robinson</i>	38
Exploring Internal Structures of an Argumentation System and Improving Reasoning Efficiency with Backward Searching Framework <i>Hao Wu</i>	43
Discussing Soundness and Completeness for Dialogues that Account for Enthymemes <i>Andreas Xydias</i>	47

On Serialisability for Argumentative Explanations

Lars Bengel

Artificial Intelligence Group, University of Hagen, Germany

Abstract

We consider the recently proposed notion of *serialisability* of semantics for abstract argumentation frameworks. This notion describes a method for the serialised non-deterministic construction of extensions through iterative addition of non-empty minimal admissible sets. This paper provides an overview on serialisability and proposes some promising applications of this concept in the context of *explainable AI*. In particular, we outline how serialisability could be employed to provide more structured explanations. We also discuss how serialisability could be utilised in discussion games.

1 Introduction

Since the introduction of *abstract argumentation* by Dung in his seminal paper [Dung, 1995], the research field has received great attention in the literature. The goal of abstract argumentation is modelling real-world exchange of arguments and using this model for reasoning. An important concept in argumentation is *admissibility* which characterises sets of arguments (extensions) that are conflict-free and defend itself against all attackers.

In recent years, the need for human understandable explanations for AI models has grown stronger and as a result of that the field of *eXplainable AI* (XAI) [Adadi and Berrada, 2018] has seen many new proposals. One very natural approach to explanations is in fact formal argumentation [Antaki and Leudar, 1992] and there already exist many approaches to XAI that utilise argumentation [Čyras et al., 2021, Ulbricht and Wallner, 2021].

In [Xu and Cayrol, 2016] the notion of *initial sets*, i. e., non-empty minimal admissible sets, has been introduced. The intuition behind initial sets is that they each solve an atomic conflict of the argumentation framework. Built on that idea, the concept of *serialisability* [Thimm, 2022] has emerged as a novel means for constructing extensions, with initial sets as the building blocks.

In this work, we take a closer look at serialisability. We will recall some work that we have already done on this topic [Bengel and Thimm, 2022] and also provide an outlook on some applications of serialisability, in particular in the domain of explainability. We consider how serialisability could be used to provide structured explanations for the acceptance of an extension. The serialised construction of an extension essentially gives us a decomposition of the extension into a series of initial sets of the respective reducts. This helps by providing a clearer view on why the extension is acceptable and allows for the natural construction of an explanation for its acceptance.

In interesting and related field are the *discussion games* [Caminada, 2015, Caminada, 2018]. In these games, two agents try to win a discussion by taking turns providing arguments that refute those of the opponent. The idea is that the presence of a winning strategy for a semantics σ guarantees the existence of a σ -extension. Caminada also shows that there are different types of discussion games, corresponding to different semantics. The discussions induced by these games can be seen as a kind of explanation for the acceptance of arguments. In the following, we will also explore how serialisability could be applied in the context of the discussion games.

The remainder of this work is structured as follows. In Section 2, we introduce the necessary background on argumentation and serialisability. Following that, in Section 3 we recall some of the work that we have already done on the topic and discuss in more detail the above mentioned potential applications of serialisability. Section 4 concludes the paper.

2 Method

We consider the *abstract argumentation frameworks* as introduced by [Dung, 1995]. Let \mathfrak{A} denote a universal set of arguments. An argumentation framework (AF) is represented as a graph whose nodes are arguments and the directed edges are conflicts between them, i. e., an edge between two arguments a and b means that a attacks b .

Definition 1. An *abstract argumentation framework* F is a tuple $F = (A, R)$ where $A \subseteq \mathfrak{A}$ is a finite set of arguments and R is a relation $R \subseteq A \times A$.

With $\mathfrak{A}\mathfrak{F}$ we denote the set of all argumentation frameworks. For a set $X \subseteq A$, we denote by $F|_X = (X, R \cap (X \times X))$ the projection of F on X . For a set $S \subseteq A$ we define $S^+ = \{a \in A \mid \exists b \in S : bRa\}$ and $S^- = \{a \in A \mid \exists b \in S : aRb\}$.

We say that a set $S \subseteq A$ is *conflict-free* if for all $a, b \in S$ we do not have that aRb . A set S *defends* an argument $b \in A$ if for all a with aRb there is an argument $c \in S$ such that cRa . A conflict-free set S is called *admissible* if S defends all $a \in S$. Let $\text{adm}(F)$ denote the set of admissible sets of F .

We define different semantics by imposing constraints on admissible sets [Baroni et al., 2018].

Definition 2. Let $F = (A, R)$ be an argumentation framework. An admissible set E

- is a *complete* (co) extension iff for all $a \in A$, if E defends a then $a \in E$,
- is a *grounded* (gr) extension iff E is complete and minimally so,
- is a *preferred* (pr) extension iff E is maximal.

All statements on minimality/maximality are with respect to set inclusion. For $\sigma \in \{\text{co}, \text{gr}, \text{st}, \text{pr}\}$ let $\sigma(F)$ denote the set of σ -extensions of F .

In [Xu and Cayrol, 2016], the authors introduce the notion of *initial sets*. They are defined as the non-empty, minimal admissible sets of an argumentation framework.

Definition 3. For $F = (A, R)$, a set $S \subseteq A$ with $S \neq \emptyset$ is called an *initial set* if S is admissible and there is no admissible $S' \subsetneq S$ with $S' \neq \emptyset$. Let $\text{IS}(F)$ denote the set of initial sets of F .

The intuition behind initial sets is that they each solve an atomic conflict of the argumentation framework. Due to the minimality, an initial set S contains only those arguments that actually contribute to the defense of S . In other words, all elements of S are relevant to the conflict that S solves. The work [Thimm, 2022] introduces a useful distinction between three different types of initial sets.

Definition 4. For $F = (A, R)$, $S \in \text{IS}(F)$, we say that

1. S is *unattacked* iff $S^- = \emptyset$,
2. S is *unchallenged* iff $S^- \neq \emptyset$ and there is no $S' \in \text{IS}(F)$ with $S'RS$,
3. S is *challenged* iff there is $S' \in \text{IS}(F)$ with $S'RS$.

In the following, we will denote with $\text{IS}^\neq(F)$, $\text{IS}^{\neq}(F)$ and $\text{IS}^{\neq\neq}(F)$ the set of unattacked, unchallenged, and challenged initial sets, respectively.

Based on the initial sets, [Thimm, 2022] introduced the notion of *serialisability*. This is a new approach for iteratively constructing the extensions of an admissible-based semantics via initial sets. For that, we first recall the notion of the *reduct* from [Baumann et al., 2020b].

Definition 5. For $F = (A, R)$ and $S \subseteq A$, the S -reduct F^S is defined via $F^S = F|_{A \setminus (S \cup S^+)}$.

Intuitively, the construction process works as follows: First, we solve an atomic conflict in F by selecting some initial set S . Afterwards, we move to the reduct F^S which may reveal new conflicts and therefore different initial sets. We continue this process until some termination criterion is satisfied.

Formally, this method for constructing extensions is represented as a transition system. For the step

of selecting an initial set (for the transition) we need a *selection function* α . Additionally, we also require a criterion β for determining if the construction of an extension is finished. The following concepts have been defined for this purpose.

Definition 6. A *state* T is a tuple $T = (F, S)$ with $F = (A, R)$ and $S \subseteq A$.

Definition 7. A *selection function* α is any function $\alpha : 2^{2^{\mathfrak{A}}} \times 2^{2^{\mathfrak{A}}} \times 2^{2^{\mathfrak{A}}} \rightarrow 2^{2^{\mathfrak{A}}}$ with $\alpha(X, Y, Z) \subseteq X \cup Y \cup Z$ for all $X, Y, Z \subseteq 2^{\mathfrak{A}}$.

The selection function will be applied as $\alpha(\text{IS}^{\neq}(F), \text{IS}^{\neq}(F), \text{IS}^{\leftrightarrow}(F))$ for some argumentation framework F . So α selects a subset of the initial sets that are eligible to be selected in the construction.

Definition 8. A *termination function* β is any function $\beta : \mathfrak{AF} \times 2^{\mathfrak{A}} \rightarrow \{0, 1\}$.

The termination function β is used to indicate when the construction of an extension is finished (this will be the case if $\beta(F, S) = 1$).

The transition rule, for some selection function α , is defined as follows:

$$(F, S) \xrightarrow{S' \in \alpha(\text{IS}^{\neq}(F), \text{IS}^{\neq}(F), \text{IS}^{\leftrightarrow}(F))} (F^{S'}, S \cup S').$$

If (F', S') can be reached from (F, S) via a finite number of steps (including no steps at all) with the above rule and also satisfies some termination criterion of β we write $(F, S) \rightsquigarrow^{\alpha, \beta} (F', S')$. Given a concrete instance of α and β , let $\mathcal{E}^{\alpha, \beta}(F)$ be the set of all S with $(F, \emptyset) \rightsquigarrow^{\alpha, \beta} (F', S)$ (for some F').

Definition 9. A semantics σ is *serialisable* if there exists a selection function α_σ and a termination function β_σ with $\sigma(F) = \mathcal{E}^{\alpha_\sigma, \beta_\sigma}(F)$ for all F .

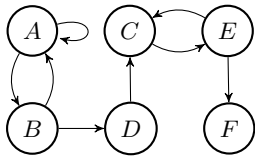


Figure 1: An AF with the initial sets: $\{B\}$ and $\{E\}$.

Example 1. Consider the preferred semantics, serialisable via $\alpha_{ad}(X, Y, Z) = X \cup Y \cup Z$ and

$$\beta_{pr}(F, S) = \begin{cases} 1 & \text{if IS}(F) = \emptyset \\ 0 & \text{otherwise} \end{cases}$$

We compute the preferred extensions of the AF F in Figure 1. We start the transition system in the state (F, \emptyset) and are allowed to select both $\{B\}$ and $\{E\}$. Assume we select $\{B\}$, thus we move to the state $(F^{\{B\}}, \{B\})$. Now, in this state α_{ad} returns $\{\{C\}$ and $\{E\}\}$. Selecting $\{C\}$ then leads to the state $(F^{\{B, C\}}, \{B, C\})$. Finally, only $\{F\}$ can be selected and we obtain the preferred extension $\{B, C, F\}$ of F . Alternatively, selecting $\{E\}$ first leads to the only other preferred extension $\{B, E\}$.

Most admissible-based semantics are serialisable, namely admissible, strong admissible, complete, grounded, preferred and stable semantics [Thimm, 2022].

3 Discussion

We start with recalling some of the results already achieved in this area. In [Bengel and Thimm, 2022], we investigated the serialisability principle in more detail. In particular, we considered its relationship to other principles from the literature. Due to space limitations, we refer to [Baroni et al., 2005, Baroni and Giacomin, 2007, Baumann et al., 2020a] for the definitions of the mentioned principles. In our investigation, we found that every serialisable semantics also satisfies conflict-freeness, admissibility and modularization. On the other hand, serialisability does not imply directionality or SCC-recursiveness and vice versa.

We also introduced the property of $\alpha\beta$ -closure for serialisable semantics, which is satisfied if every path of the corresponding transition system terminates, i. e., leads to some extension. Interestingly, this leads to the following connection.

Theorem 1. *If a semantics σ is serialisable via α_σ and β_σ and is $\alpha_\sigma\beta_\sigma$ -closed, then σ satisfies directionality.*

Furthermore, in [Bengel and Thimm, 2022] we also took a closer look at the unchallenged semantics from [Thimm, 2022] defined via the selection function $\alpha_{uc}(X, Y, Z) = X \cup Y$ and the termination function

$$\beta_{uc}(F, S) = \begin{cases} 1 & \text{if } IS^{\neq}(F) \cup IS^{\neq}(S) = \emptyset \\ 0 & \text{otherwise} \end{cases}$$

Example 2. Consider the AF F in Figure 1. Both $\{B\}$ and $\{E\}$ are unchallenged. If we select $\{B\}$ in the first step, then in the reduct $F^{\{B\}}$ we have that $\{C\}$ and $\{E\}$ are now challenged initial sets. Thus $\beta_{uc}(F^{\{B\}}, \{B\}) = 1$ and $\{B\}$ is an unchallenged extension of F . However, if we select $\{E\}$ first, then $\{B\}$ can still be selected in the state $(F^{\{E\}}, \{E\})$ and we reach the only other unchallenged extension $\{B, E\}$.

Essentially, this semantics amounts to exhaustively adding unattacked and unchallenged initial sets. The unchallenged semantics is α_{uc}, β_{uc} -closed and thus satisfies directionality. It also satisfies reinstatement, but it does for example not satisfy SCC-recursiveness and I-maximality. We have also analysed the computational complexity of the relevant reasoning tasks.

We now turn to some applications of serialisability that we plan to explore in the future.

Using an argumentative approach to provide explanations is very natural [Antaki and Leudar, 1992]. An extension E contains all relevant arguments to justify its own acceptance. However, a larger AF may consist of many different conflicts and for each only a subset of E is of relevance. This is where serialisable semantics can be very useful. With the serialised construction of an extension E we additionally obtain a decomposition of E into a series of initial sets S_1, \dots, S_n of the respective reducts. Each of these initial sets solves an atomic conflict of the AF and only contains the arguments relevant to that conflict. Not only does this decompose an extension into initial sets, it also gives us information about their order and how some initial sets are only revealed once other conflicts are addressed. This information could, for example, be utilised to generate better, more structured explanations for an extension. In particular, in my PhD thesis, I want to investigate the relation to other recent approaches for acceptance explanations from the literature, such as related admis-

sibility [Fan and Toni, 2015] or explanation schemes [Baumann and Ulbricht, 2021].

Another interesting application of serialisability could be the discussion games [Caminada, 2018]. In the course of such a game, each agent essentially iteratively constructs his own extension by adding one argument each round. This is somewhat similar to how the serialised construction works, where we select initial sets instead of individual arguments. Thus, it would be interesting to explore the relation between both concepts. There exist different discussion games for the different semantics. Notably, the notion of serialisability allows us to define completely new semantics by defining only a selection and a termination function. For example, confronted with some argument by our opponent, we may want to reply with the strongest refutation instead of simply providing any counterargument. For this purpose, we can define a selection function with some heuristic that evaluates the strength of initial sets (or arguments) and only allows us to select the strongest. This way, we could construct a strong, discussion-like explanation for the acceptance of our starting argument.

4 Conclusion

In this work we discussed the recently introduced notion of serialisability for argumentation semantics that provides a non-deterministic procedure for constructing extensions. We looked at some results relating serialisability to other principles from the literature. We also considered one instantiation, the unchallenged semantics, and its properties. For my PhD thesis, we outlined that the decomposition of an extension into initial sets could be used to provide more structured explanations. Furthermore, we also discussed the planned investigation of the similarities to discussion games [Caminada, 2015] and utilising serialisability for them.

Acknowledgements The research reported here was partially supported by the Deutsche Forschungsgemeinschaft (grant 375588274).

References

- [Adadi and Berrada, 2018] Adadi, A. and Berrada, M. (2018). Peeking inside the black-box: a survey on explainable artificial intelligence (XAI). *IEEE access*, 6:52138–52160.
- [Antaki and Leudar, 1992] Antaki, C. and Leudar, I. (1992). Explaining in conversation: Towards an argument model. *European Journal of Social Psychology*, 22(2):181–194.
- [Baroni et al., 2018] Baroni, P., Caminada, M., and Giacomin, M. (2018). Abstract argumentation frameworks and their semantics. In Baroni, P., Gabbay, D., Giacomin, M., and van der Torre, L., editors, *Handbook of Formal Argumentation*, pages 159–236. College Publications.
- [Baroni and Giacomin, 2007] Baroni, P. and Giacomin, M. (2007). On principle-based evaluation of extension-based argumentation semantics. In *Artificial Intelligence*, volume 171, pages 675–700. Elsevier.
- [Baroni et al., 2005] Baroni, P., Giacomin, M., and Guida, G. (2005). SCC-recursiveness: a general schema for argumentation semantics. *Artificial Intelligence*, 168(1–2):162–210.
- [Baumann et al., 2020a] Baumann, R., Brewka, G., and Ulbricht, M. (2020a). Comparing weak admissibility semantics to their Dung-style counterparts—reduct, modularization, and strong equivalence in abstract argumentation. In *Proceedings of the International Conference on Principles of Knowledge Representation and Reasoning*, pages 79–88.
- [Baumann et al., 2020b] Baumann, R., Brewka, G., and Ulbricht, M. (2020b). Revisiting the foundations of abstract argumentation—semantics based on weak admissibility and weak defense. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 2742–2749.
- [Baumann and Ulbricht, 2021] Baumann, R. and Ulbricht, M. (2021). Choices and their consequences—explaining acceptable sets in abstract argumentation frameworks. In *KR*, pages 110–119.
- [Bengel and Thimm, 2022] Bengel, L. and Thimm, M. (2022). Serialisable semantics for abstract argumentation. In *Proceedings of COMMA 2022*. to appear.
- [Caminada, 2015] Caminada, M. (2015). A discussion game for grounded semantics. In *International Workshop on Theory and Applications of Formal Argumentation*, pages 59–73. Springer.
- [Caminada, 2018] Caminada, M. (2018). Argumentation semantics as formal discussion. In Baroni, P., Gabbay, D., Giacomin, M., and van der Torre, L., editors, *Handbook of formal argumentation*, pages 487–518. College Publications.
- [Čyras et al., 2021] Čyras, K., Rago, A., Albin, E., Baroni, P., and Toni, F. (2021). Argumentative xai: a survey. *arXiv preprint arXiv:2105.11266*.
- [Dung, 1995] Dung, P. M. (1995). On the Acceptability of Arguments and its Fundamental Role in Non-monotonic Reasoning, Logic Programming and n-Person Games. *Artificial Intelligence*, 77(2):321–358.
- [Fan and Toni, 2015] Fan, X. and Toni, F. (2015). On computing explanations in argumentation. In *Twenty-Ninth AAAI Conference on Artificial Intelligence*.
- [Thimm, 2022] Thimm, M. (2022). Revisiting initial sets in abstract argumentation. *Argument & Computation*.
- [Ulbricht and Wallner, 2021] Ulbricht, M. and Wallner, J. P. (2021). Strong explanations in abstract argumentation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 6496–6504.
- [Xu and Cayrol, 2016] Xu, Y. and Cayrol, C. (2016). Initial sets in abstract argumentation frameworks. In *Proceedings of the 1st Chinese Conference on Logic and Argumentation (CLAR’16)*, volume 1811, pages 72–85.

Debating Ethics: Using Argumentation to Support Dialogue

Elfia Bezou-Vrakatseli

King's College London

Abstract

This project proposes the exploration and analysis of natural language texts, specifically of moral debates, via argument schemes and critical questions. It aims to develop new natural language datasets and AI algorithms through defining semi-automated approaches for the identification and extraction of argument schemes. These new datasets will be in the form of philosophical debates on society's ethical issues and will serve the purpose of creating a new corpus that will enable the automatization of argument mining from texts. Subsequently, the obtained argument schemes will be used to support dialogical exchanges between humans and AI systems with respect to transparent and rational reasoning. This paper describes progress towards said goals, starting from identifying argument schemes and scheme taxonomies specialised in ethical reasoning.

1 Introduction

Artificial Intelligence (AI) is becoming more and more powerful, efficient, and autonomous. Therefore, the need for safe and trustworthy systems increases as well; systems for which there is assurance about the correctness of their behaviour and that inspire confidence for their users about their decision-making process.

Interacting and communicating with AI systems is, thus, of vital importance and at the core of the current research in the field of AI. Dialogue is a powerful tool of interaction and communication and multiple applications of AI make direct use of it (e.g., conversational agents). Dialogical exchanges enhance the

understanding of the involved parties not only by conveying information but also as a means of exploring the reasoning of the interlocutors. Therefore, dialogues also facilitate joint reasoning.

This project proposes the use of argumentation tools to support the dialogue between humans and humans and AI systems, by generalising argumentation-based characterisations of non-monotonic (*nm*) inference relations to situations where agents exchange locutions. Argumentation, at the core of human reasoning, is indispensable to the manner people understand causality and explain facts; “the majority of what might look like causal attributions turn out to look like argumentative claim-backings” [Antaki and Leudar, 1992]. Indeed argumentation has already been used for various kinds of explanation in AI; from explaining the decision making of recommender systems [Rago et al., 2020] to enabling AI to explain outputs determined by humans or by machines [Čyras et al., 2019].

In order to support dialogue that can enable transparent and rational joint reasoning, we draw from recent advances in natural language processing (NLP), and in particular argument mining (AM) [Lawrence and Reed, 2020], an area of research which focuses on extracting arguments and their relations from text. To this end, natural language (NL) texts are analysed with the use of argument schemes and critical questions. Argument schemes represent stereotypical patterns of reasoning and consist of a set of premises and a conclusion [Walton, 1996]. In particular, this project focuses on NL debates in matters of ethics and the specific argumentative mechanisms that underline the dialogues involved in said debates.

As previously stated, representing dialogue with argumentation tools can facilitate joint human and AI reasoning. Defeasible schemes enable agents to arrive at presumptive conclusions on their course of action in a situation where the prolongation of collecting evidence may be costly [Macagno et al., 2017]. The focus on moral debates can result in ameliorating the decision-making process on ethical issues (e.g., in the case of autonomous cars [Matthias, 2004]). Besides enabling the decision of what the best choice is in given circumstances (practical cognition), schemes are also useful for goals of epistemic cognition (getting to the truth of a matter) [Macagno et al., 2017]. Ultimately, the enabling of joint human and AI reasoning can help solve the *value alignment problem* [Modgil, 2017].

Despite AM being a very promising field with rapid growth, there are still various aspects that need improvement (e.g., the level of detail of its techniques, the lack of consistently annotated argument data, etc. [Lawrence and Reed, 2020]). Our research enhances AM by going beyond its standard techniques (i.e., to determine the relation between premise and conclusion and identify support/attack relations between arguments) by making use of informal logic. Moreover, to achieve a higher level of annotation consistency, we purport to tackle one of the biggest issues of argument schemes: the existing vast diversity of them and of the taxonomies they fall under, which renders them hard to use in AI research [Katzav and Reed, 2004][Hornikx, 2013]. We aim to tackle these issues by addressing three main questions:

1. How can one reconcile the existing various argument classifications so as to devise a theoretically well-founded, as well as practically useful, hybrid classification, which can be specialised for application in ethical debates?
2. Given such a classification system, what is an effective way to develop a semi-automated way to map NL arguments, used in ethical debates, to argument schemes?
3. How can the outputs of the previous techniques be exploited to support the dialogue both be-

tween humans and between humans and AI systems, with emphasis on matters of ethics?

2 Method

The initial step of this research consists of annotating arguments from moral debates taken from the user-generated platform Kialo¹, using the two most used argument scheme classifications: Walton’s argument schemes [Walton et al., 2008] and Wagemans’ Periodic Table of Argument (PTA) [Wagemans, 2016]. The former is more comprehensive, while the latter is more practically useful and can be seen as an intermediate between the semantic detail of the former and the relation between premise-conclusion used in AM. We purport to reconcile the two and develop a new classification that can be applied to ethical debates. An example of the annotating process can be found below.

Classification of an argument taken from the debate *Pro-life vs Pro-choice: Should abortion be legal?*

- Thesis statement *T*: Pregnant people should have the right to choose abortion.
- First support claim *A*: Access to legal abortion improves the health and safety of women.

The first argument to be categorised in this debate is the claim stated in the second bullet point above. The initial step is examining whether the argument in question is an enthymeme; argument *A* is an enthymeme. The next step consists of assuming the implicit premise/conclusion; in this example, the conclusion of the argument was assumed to be the thesis statement (*T*). Thus, the argument under analysis becomes:

A': Access to legal abortion improves the health and safety of women, so pregnant people should have the right to choose abortion.

For Walton’s taxonomy, the main guideline was the Argument Scheme Key (ASK); a dichotomous identification key proposed by [Visser et al., 2021], which

¹<https://www.kialo.com/>

consists of a series of disjunctive choices based on the distinctive features of argument schemes. In Figure 1, the 25 initial steps of ASK can be found, along with the path that corresponds to the steps taken in classifying argument A' . The end of the path leads to the classification of the argument. In this case, argument A' is an *argument from positive consequences*.

Walton’s representation of the *argument from positive consequences* scheme is defined by a Premise: “If A is brought about, good consequences will (plausibly) occur” and Conclusion: “Therefore, A should be brought about”, with the following critical questions (CQs): CQ1: “How strong is the likelihood that the cited consequences will (may, must) occur?”; CQ2: “What evidence supports the claim that the cited consequences will occur and is it sufficient to support the strength of the claim adequately?”; “CQ3: Are there opposite consequences (bad as opposed to good) that should be taken into account?”. Walton’s taxonomy is the most widely used taxonomy, containing 60 argument schemes. The most prominent argument schemes from this taxonomy that our method identified in the analyzed debates are: *argument from example*, *argument from positive/negative consequences*, *argument from values*, *argument from cause to effect*, *argument from analogy*, *argument from alternatives*, *argument from expert opinion*.

Critical questions represent a key aspect of argument schemes, as they can be pointers to counter-arguments or supporting arguments, since the argument answering a critical question attacks/supports the original argument by virtue of the critical question. Since A' is an *argument from positive consequences*, one of the corresponding CQS is: *Are there other opposite consequences that should be taken into account?*. Argument B : *Abortion has harmful mental and physical consequences for the woman involved*. serves as an answer to this question and is a counter-argument to A' . Moreover, argument B is an *argument from negative consequences*.

Contrary to Walton’s taxonomy, Wagemans’ Periodic Table of Arguments is based on an *a priori* constrained set of possible combinations between different characterisations of argument. To annotate with Wagemans’ taxonomy, the three characterisations of argument described in [Wagemans, 2016] are

1. Argument relies on a source’s opinion or character	2
- Argument does not depend on a source’s opinion or character	17
2. Argument is about the source’s character	3
- Argument is about the source’s opinion	9
3. Argument establishes the source’s character	
- Argument refers to the source’s existing character	4
4. Argument relies on the source’s good character	<i>Ethotic argument</i>
- Argument relies on bad character	5
5. Source is biased	6
- Argument is not related to bias	7
6. Source does not take both sides into account	<i>Argument from bias</i>
- Source’s opinion is not acceptable	<i>Bias ad hominem</i>
7 (5). Source is of bad overall character	<i>Generic ad hominem</i>
- The source’s actions are not compatible with their commitments	8
8. Source’s actions contradict the advocated position	<i>Pragmatic inconsistency</i>
- Source is not credible due to inconsistent commitments	
	<i>Circumstantial ad hominem</i>
9 (2). Argument establishes a source’s opinion	10
- Argument is based on an existing opinion	11
10. Commitment at issue is consistent with existing commitments	
- Commitment at issue is not consistent with existing commitments	
	<i>Argument from inconsistent commitment</i>
11 (9). Source is a general group of people	<i>Argument from popular opinion</i>
- Source is a specific individual	12
12. Source is an expert in the subject domain	<i>Argument from expert opinion</i>
- Source’s credibility is not based on domain knowledge	13
13. Source is a witness	<i>Argument from witness testimony</i>
- Source is not a witness	14
14. Argument is based on the source’s memories	<i>Argument from memory</i>
- Argument does not explicitly refer to memories	15
15. Argument is based on the source’s visual perception	<i>Argument from perception</i>
- Argument does not explicitly refer to perception	16
16. Conclusion is about a course of action	<i>Two-person practical reasoning</i>
- Argument is not action-oriented	<i>Argument from position to know</i>
17 (1). Conclusion is about a course of action	18
- Conclusion is not specifically action-oriented	32
18. Argument focuses on the outcome of an action	22
- Argument hinges on another motivation for the action	19
19. Course of action follows an established practice	20
- Course of action is compared to a similar or alternative action	21
20. Course of action is explicitly regulated	<i>Argument from rules</i>
- Course of action follows general practices	<i>Argument from popular practice</i>
21 (19). Action is best alternative on the basis of prior commitments	
	<i>Argument from sunk costs</i>
- Action is directly compared to another	<i>Practical reasoning from analogy</i>
22 (18). Conclusion promotes a positive outcome	23
- Conclusion prevents a negative outcome	26
23. Course of action assists someone else	24
- Course of action does not offer help	25
24. Course of action relieves suffering	<i>Argument from distress</i>
- Argument does not mention suffering	<i>Argument from need for help</i>
25 (23). Course of action promotes a goal	<i>Argument from (positive) consequences</i>
- Course of action is not related to an explicit goal	<i>Practical reasoning</i>

Figure 1: An example of classifying an argument via ASK. The path that corresponds to the steps taken in classifying argument A' is the following, with the answer in brackets: Argument relies on a source’s opinion or character (No); Conclusion is about a course of action (Yes); Argument focuses on the outcome of the action (Yes); Conclusion promotes a positive outcome (Yes); Course of action assists someone else (No); Course of action promotes a goal (Yes); The argument is an *argument from positive consequences*.

considered (i.e., subject vs predicate arguments, first vs second order arguments, and argument substance). In particular, A' is a first-order (as it is not source-based), predicate argument as the premise and the conclusion share the same subject (i.e. “access to legal abortion”). Note that A' needs to be rephrased slightly in order to match the definition of a predicate argument. Lastly, the conclusion of A' is a proposition of *policy* (“pregnant people should have the right . . .”), while the premise is a *fact* (“legal access improves the health and safety of pregnant people”), thus the argument is of type *1-pre-PF*, equivalent to a *pragmatic argument*.

3 Discussion

This project constitutes, thus, a study of NL texts with the tools of argumentation to the end of supporting transparent reasoning and rendering systems safer and more trustworthy. So far, the focus of the project has been on building an appropriate taxonomy and new argument schemes, both specialised in ethical reasoning. The reconciliation of the two predominant taxonomies in a hybrid one leverages the strengths of both, while identifying the schemes particular to ethical reasoning.

This research goes beyond standard argument mining techniques by making use of informal logic. In particular, the use of argument schemes and critical questions offers a *semantically* richer approach to inner- and inter- argument classification, as premises support a conclusion *by virtue* of instantiating a scheme and support/attack relations are instigated *in response* to critical questions.

The aforementioned example offers insights into the approach adopted towards this hybrid taxonomy. Firstly, observing the co-occurrences of argument schemes in both taxonomies allows us to detect correspondences between the two in order to reconcile them. For example, Walton’s *argument from positive consequences* is often classified as *1-pre-PF* using Wagemans’ PTA (something also independently observed in [Visser et al., 2021]). Secondly, comparing and contrasting the annotation guidelines of each taxonomy helps reflect on them and their usefulness.

For instance, deciding if an argument is second or first order (source-based or not) is a criterion in both taxonomies, which serves as an indicator of the necessity of this distinction/characteristic, rendering it an indispensable criterion of our hybrid taxonomy as well.

The immediate next step of this process consists of completing the development of said hybrid taxonomy, which will also take into account the corresponding critical questions of each scheme and which will incorporate: (1) the existing argument schemes that are predominantly used in ethical reasoning; (2) groupings of existing schemes (with the help of the clustering nature of the ASK algorithm along with the criteria of the PTA); (3) new schemes, either object-level, specialised in ethical reasoning, or meta-schemes (i.e., schemes of another level that enable commentary on object-level reasoning).

Lastly, the annotation process explained in this paper has also resulted in the creation of a large dataset of annotated ethical arguments. The following steps are aimed at developing a semi-automated way to map NL argument to argument schemes. The final step of this research consists in combining the outputs of all the aforementioned techniques in order to support the dialogue both between humans and between humans and AI systems in matters of ethics.

4 Conclusion

This paper describes an initial step towards realising the long-term research goal of supporting dialogue between humans and between humans and AI systems. The first step consists of developing a new taxonomy, as well as new argument schemes, specialised in ethical reasoning. To this end, arguments from ethical debates were annotated using Walton’s and Wagemans’ taxonomies; an important step both for identifying schemes and taxonomies specialised in ethical reasoning and for the development of a dataset. The use of argument schemes and critical questions goes beyond standard annotation approaches for argument mining and proposes the use of argumentation to achieve a semantically richer approach to argument annotation.

Acknowledgements

I would like to thank Oana Cocarascu and Sanjay Modgil for their guidance, valuable insights, and suggestions in this project.

References

- [Antaki and Leudar, 1992] Antaki, C. and Leudar, I. (1992). Explaining in conversation: Towards an argument model. *European Journal of Social Psychology*, 22(2):181–194.
- [Čyras et al., 2019] Čyras, K., Letsios, D., Misener, R., and Toni, F. (2019). Argumentation for explainable scheduling. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 2752–2759.
- [Hornikx, 2013] Hornikx, J. (2013). A bayesian perspective on argument quality—the ad populum argument under the microscope. *Language Mastery Journal*, 35(2):128–143.
- [Katzav and Reed, 2004] Katzav, J. and Reed, C. A. (2004). On argumentation schemes and the natural classification of arguments. *Argumentation*, 18(2):239–259.
- [Lawrence and Reed, 2020] Lawrence, J. and Reed, C. (2020). Argument mining: A survey. *Computational Linguistics*, 45(4):765–818.
- [Macagno et al., 2017] Macagno, F., Walton, D., and Reed, C. (2017). Argumentation schemes. history, classifications, and computational applications. *History, Classifications, and Computational Applications (December 23, 2017)*. Macagno, F., Walton, D. & Reed, C, pages 2493–2556.
- [Matthias, 2004] Matthias, A. (2004). The responsibility gap: Ascribing responsibility for the actions of learning automata. *Ethics and information technology*, 6(3):175–183.
- [Modgil, 2017] Modgil, S. (2017). Dialogical scaffolding for human and artificial agent reasoning. In *AIC*, pages 58–71.
- [Rago et al., 2020] Rago, A., Cocarascu, O., Bechlivanidis, C., and Toni, F. (2020). Argumentation as a framework for interactive explanations for recommendations. In *Proceedings of the International Conference on Principles of Knowledge Representation and Reasoning*, volume 17, pages 805–815.
- [Visser et al., 2021] Visser, J., Lawrence, J., Reed, C., Wagemans, J., and Walton, D. (2021). Annotating argument schemes. *Argumentation*, 35(1):101–139.
- [Wagemans, 2016] Wagemans, J. (2016). Constructing a periodic table of arguments. In *Argumentation, objectivity, and bias: Proceedings of the 11th international conference of the Ontario Society for the Study of Argumentation (OSSA)*, Windsor, ON: OSSA, pages 1–12.
- [Walton, 1996] Walton, D. (1996). Argumentation schemes for presumptive reasoning lawrence erlbaum associates mahwah. *New Jersey*.
- [Walton et al., 2008] Walton, D., Reed, C., and Macagno, F. (2008). *Argumentation schemes*. Cambridge University Press.

Characterization of Unresolvable Conflicts in Abstract Argumentation

Lydia Blümel

Artificial Intelligence Group, University of Hagen

Abstract

We investigate Abstract Argumentation Frameworks (AFs) for which no nonempty admissible extension exists, meaning no argument is accepted. Although a few observations regarding such unresolvable AFs have been stated in the literature, notably that they must include at least one odd cycle, little is known about the structure of larger AFs with this property. We discuss the question how unresolvable conflicts in AFs can be broken down in order to better explain their behaviour and how they can be characterized in terms of syntactic criteria in general.

1 Introduction

Why? Humans, and especially researchers, are on a constant search for explanations for the inner workings of the systems surrounding them, be they natural or artificial, motivated by the desire for good decision making, adaptation to, and improvement of the world they live in. Towards this goal artificial argumentation has proven to be a fruitful research area with numerous applications and a fast growing number of publications in the last decades [Atkinson et al., 2017]. On a more recent note, the First International Workshop on Argumentation for eXplainable AI was held in September of this year, demonstrating both the potential of argumentation and the challenges it has to face as a theory for explanations. One such challenge is to provide an explanation of the underlying basic argumentation framework introduced by Dung itself [Dung, 1995]. With the number of applications in XAI growing this task has recently received a lot of attention [Fan and Toni, 2015a, Ulbricht and Baumann, 2019,

Saribatur et al., 2020].

As it turns out, non-acceptability of arguments is harder to explain than acceptability. Since non-acceptance is coNP-complete, providing an explanation is not a matter of existence but of exhaustion [Dvorák and Dunne, 2018]. The problem we want to discuss here constitutes a worst-case scenario of this - argumentation frameworks in which no argument can be accepted under the admissible semantics. We will call such argumentation frameworks *unresolvable*. Take a look at the AF in Fig.1. Although it consists of only seven arguments, even an experienced reader will have trouble judging on first glance whether it contains a nonempty admissible extension (it does not). The reason why this decision problem is so hard for a framework like the one depicted in Fig.1 is that it consists of a single strongly connected component (SCC) so the SCC-recursive scheme [Baroni et al., 2005] normally used for speeding up decisions on credulous acceptance [Liao, 2013] does not apply here. Since AFs can get very large for some applications [Sakai et al., 2018] this poses a serious problem in practice which is also elaborated in [Lafages, 2021].

Apart from the computational aspects, there is the question of which form an explanation for unresolvability can or should take? Existing approaches [Fan and Toni, 2015b, Ulbricht and Baumann, 2019, Saribatur et al., 2020] give explanations in the form of sets of problematic arguments or subframeworks. The high number of involved arguments in case of large SCCs poses a challenge here as well. We want to take a different approach and propose an investigation through syntactic criteria instead, to give a general account of what makes a constellation of arguments an unresolvable

conflict.

The remainder of this paper is as follows. The next section contains the necessary background definitions to formally introduce unresolvable conflicts and we make some preliminary observations. In the discussion part we look at related existing work and outline two main future work directions. We conclude with a short summary in the last section.

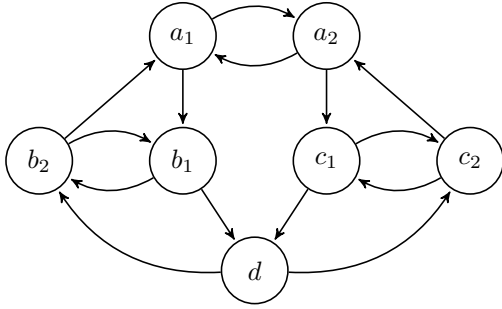


Figure 1: A nontrivial SCC with no admissible extension

2 Method

2.1 Outlining the Problem

Let us begin by introducing the syntax of abstract argumentation frameworks and the admissible semantics which were first defined in [Dung, 1995]. An *Abstract Argumentation Framework (AF)* is a tuple $AF = (A, R)$ where A is a finite set of arguments and $R \subseteq A \times A$ an attack relation representing conflicts between arguments. We say an argument a attacks an argument b iff $(a, b) \in R$. For some $E \subseteq A$ we define

$$E^- = \{a \in A \mid \exists e \in E : (a, e) \in R\}$$

$$E^+ = \{a \in A \mid \exists e \in E : (e, a) \in R\}$$

to be the set of all attackers of E and the set of all arguments attacked by E , respectively, and say E is unattacked if $E^- \cap (A \setminus E) = \emptyset$.

A set E is *conflictfree* iff $E \cap E^+ = \emptyset$, i.e. no two arguments in E are attacking each other. E *defends* itself iff $E^- \subseteq E^+$, meaning every argument attacking E is

in turn attacked by some argument in E . An argument set which both defends itself and is conflictfree is called *admissible*. The admissible semantics assigns to each AF the set of its admissible extensions

$$adm(AF) = \{E \subseteq A \mid E \cap E^+ = \emptyset \wedge E^- \subseteq E^+\}$$

We say an argument a is *credulously accepted* under the admissible semantics if an admissible extension E with $a \in E$ exists.

For example, in Fig. 1 the argument set $\{a_1, d\}$ is conflictfree, but has an attacker, c_1 , which it does not attack back. Therefore it is not admissible. If we add the argument a_2 to our set, say $E = \{a_1, a_2, d\}$, we get defense, since the set of attackers $E^- = \{b_1, c_1, b_2, c_2, a_1, a_2\}$ is the same as the set of attacked arguments E^+ . But a_1 and a_2 attack each other, so E is not conflictfree and again not admissible. Repeating these observations for other combinations allows us to reject all arguments of this AF resulting in the empty set being the only admissible extension and no argument being credulously accepted. From this point on, any AF with this property will be called an *unresolvable AF*.

Definition 2.1. Let $AF = (A, R)$ be an AF. We say AF is unresolvable iff $adm(AF) = \{\emptyset\}$.

Our aim is to find a syntactic explanation for the unresolvability of an AF. For this we will first examine in which cases unresolvability can be reduced to the unresolvability of subframeworks. Following the notation from [Baroni et al., 2005], by $AF \downarrow_S = (S, R \cap (S \times S))$ we denote the restriction of AF on some $S \subseteq A$.

Definition 2.2. Let $AF = (A, R)$ be an AF. A set of arguments $C \subseteq A$ is an *unresolvable conflict* iff $AF \downarrow_C$ is unresolvable.

At the end of this section let us introduce some useful terminology for discussing properties of unresolvable AFs. A path $[a, b]$ from an argument a to another argument b is a finite sequence (a_0, \dots, a_n) of pairwise distinct arguments such that $a = a_0$, $b = a_n$ and $(a_i, a_{i+1}) \in R$ for $0 \leq i < n$. We call n the length of $[a, b]$. If n is even, $[a, b]$ is called an *even/support path*, else an *odd/attack path*. We will allow $a = b$ and call such a path $[a, a]$ a cycle. We say an $S \subseteq A$ is *strongly connected* iff for any $a, b \in S$ either $a = b$ or some paths $[a, b]$ and $[b, a]$ in

$AF|_S$ exist. S is a *strongly connected component (SCC)* iff S is \subseteq -maximal among the strongly connected subsets of A [Baroni et al., 2005].

2.2 Preliminary Investigations

Some simple observations regarding unresolvability can be made, e.g. an unresolvable AF cannot contain any unattacked argument. On a more general note, the fact that the admissible semantics satisfies directionality narrows the problem down to the strongly connected components (SCCs) of an AF. The principle of directionality [Baroni et al., 2005], which holds for the admissible semantics, states that given an unattacked set of arguments S a set $E \subseteq A$ is admissible in AF if and only if $E \cap S$ is admissible in $AF \downarrow_S$. A direct consequence of this property for the case of unresolvable AFs is the following statement.

Proposition 2.3. *If an AF is unresolvable then all its unattacked SCCs are unresolvable conflicts.*

As a result, we can focus on the study of syntactic peculiarities (if they exist) of unresolvable AFs consisting of a single strongly connected component. Strongly connected components are the current atoms for computing argumentation semantics. When confronted with a single large SCC, no systematic approach for exploiting its inner structure for explanations has been proposed so far.¹ We will now give an example where an SCC decomposition can be easily spotted and conducted.

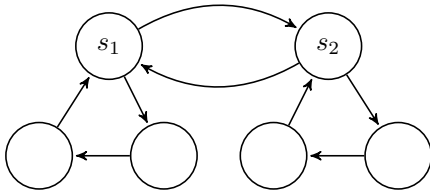


Figure 2: Mutual attack between two 3-cycles

The AF in Fig.2 consists of two cycles of length 3 which are connected by a single pair of arguments s_1, s_2 with mutual attacks $(s_1, s_2), (s_2, s_1) \in R$ on each other, thus forming a single SCC S . We can now argue that S is

¹[Lafages, 2021] makes use of sparseness to improve computation

unresolvable by observing that an isolated 3-cycle is unresolvable and that the connection via a single mutual attack makes it impossible for an argument from one such cycle to help an argument of the other to become admissible. This observation can be generalized as a decomposability criterion for unresolvable SCCs which to the best of our knowledge has not been stated in the literature, yet.

Proposition 2.4. *Let S be an AF consisting of a single SCC, $S_1, S_2 \subset A$, $S_1 \cap S_2 = \emptyset$, $S_1 \cup S_2 = A$ and $s_1 \in S_1, s_2 \in S_2$ such that $(s_1, s_2), (s_2, s_1) \in R$ and every directed path $[a, b]$ from S_1 to S_2 goes via (s_1, s_2) and vice versa.² Then S is unresolvable iff S_1 and S_2 are both unresolvable.*

Proof. (\Rightarrow) Suppose $E \subseteq S_1$ is admissible in $AF \downarrow_{S_1}$. Then either $s_1 \in E$ in which case adding S_2 preserves admissibility since s_1 attacks s_2 back. Or $s_1 \notin E$, then no attack on E is added, so E continues to be defended and conflictfree. So $E \in adm(S)$.

(\Leftarrow) Suppose now both S_1 and S_2 are unresolvable and w.l.o.g. let $E \subseteq S_1$ be conflictfree. Then E has an attacker a against which it cannot defend itself. Since s_2 is the only argument of S_2 attacking S_1 , E can only add s_2 to reach admissibility in S as a whole. But for any conflictfree $E' \subseteq S_2$ containing s_2 we know an attacker a' in S_2 exists against which E' cannot defend itself. E cannot defend E' either, since the only attack of S_1 on S_2 is on s_2 itself. Therefore $E \cup E'$ cannot be admissible in AF. \square

Prop. 2.4 describes a simple method to decompose SCCs with an arguably limited range of cases where it applies. However, it can serve as a motivation to conduct a systematic search for preprocessing techniques to exploit syntactic properties for solving large SCCs. At the same time it also constitutes a clear case of how two small unresolvable conflicts can explain the unresolvability of their union given a certain type of connection. For future work we plan to expand on this result to find more classes of decomposable SCCs and to deepen our understanding of multiple interconnected unresolvable conflicts.

²To be precise we ask that for any $[a, b] = (a, a_1, \dots, a_{n-1}, b)$ if $a \in S_1, b \in S_2$ then for some $0 \leq i \leq n$ we have $a_i = s_1, a_{i+1} = s_2$ and vice versa a path $[b, a]$ with $a \in S_1, b \in S_2$ has to contain $a_i = s_2, a_{i+1} = s_1$.

3 Discussion

In this section, we discuss different approaches towards characterizing unresolvable argumentation frameworks in the context of existing work.

As mentioned in the introduction a lot of work has already been done on explaining the non-acceptability of single arguments [Fan and Toni, 2015b, Ulbricht and Baumann, 2019, Saribatur et al., 2020]. In this context, unresolvability appears as a special case of credulous non-acceptance which has been systematically studied in [Ulbricht and Baumann, 2019] among other cases. The viewpoint taken by the authors there is that unresolvability is an undesirable state and a diagnosis is needed in order to find the problematic parts of the AF and to "repair" it. Towards this end their method is to find minimal sets of arguments resp. attacks that, if excluded from the AF, result in a resolvable AF. They call such a set a diagnosis of the AF. They investigate the computational complexity of their approach and whether certain structural properties allow one to draw conclusions on the form of the diagnosis, e.g. symmetrical AFs³ are always resolvable and therefore have an empty credulous diagnosis for the admissible semantics. The influence of directionality we stated in Prop.2.3 is exploited by making use of so called Splitting, which provides a structural account of diagnoses. On a single SCC structural insights are difficult to come by with this form of an explanation which is more like an ideal solution. As such, it serves to outline the problem while the way to find such a solution in an explainable manner needs further investigations on the structure of unresolvable AFs along the lines of Prop. 2.4. A similar problem arises with the strongly rejecting subframeworks proposed in [Saribatur et al., 2020] to explain non-acceptability. Given some AF = (A, R) a set of arguments E strongly rejects an argument a if a is not credulously accepted in $AF_{\downarrow E'}$ for any $E \subseteq E' \subseteq A$. Again this gives us no insight on how to recognize an unresolvable SCC. An interesting requirement of an explanation is introduced, though, namely that a has to be unaccepted in any of the steps inbetween $AF_{\downarrow E}$ and AF. Requesting a certain degree of stability from an explanation is desirable in case of unresolvable AFs as well when trying to decompose large unresolvable conflicts

into smaller ones. It rules out generic approaches for SCC-decomposition like [Baroni and Giacomin, 2006], and by doing so emphasizes the idea of explaining. Rather than requesting stability for any step though we consider the investigation of types of inbetween steps one can take without violating unresolvability to be promising for future research.

A conceptually different approach worth pursuing can be found in some side results of other works. Apparently, the presence of odd cycles plays a big role for the unresolvability of an AF. To this day, a complete description of what exactly that role is remains an open question, but two insightful criteria regarding unresolvability in the literature are both related to odd cycles. In [Dung, 1995], it is observed that an AF has to contain at least one odd cycle to be unresolvable. Interesting is the observation in [Dvorák, 2012] that an argumentation framework with only odd cycles has only one complete extension (the grounded). Applying the second result to a single non-trivial SCC yields unresolvability for SCCs which contain only odd cycles, since the grounded extension of nontrivial SCCs is always empty. It is remarkable that both of these criteria, the necessary by Dung and the sufficient by Dvorak are linked to the presence of odd cycles. For a purely syntactical characterization of unresolvable conflicts we consider this as another viable and interesting route to take in the future.

4 Conclusion

We formally introduce the problem of unresolvability as a special case of non-acceptance and highlight its importance for ongoing research on computation and explanation in Abstract Argumentation. We then show how the problem can be reduced to the case of nontrivial strongly connected components and provide a first result towards a further structural decomposition. Related results from the literature are presented and discussed with the overall objective of a syntactical characterization of unresolvability in mind. Two directions for future work are identified - decomposing SCCs by investigating the composition of large unresolvable conflicts out of smaller unresolvable conflicts and the investigation of odd cycles for general criteria.

³AFs where all attacks are mutual

Acknowledgements The research reported here was partially supported by the Deutsche Forschungsgemeinschaft (grant 375588274).

References

- [Atkinson et al., 2017] Atkinson, K., Baroni, P., Giacomin, M., Hunter, A., Prakken, H., Reed, C., Simari, G. R., Thimm, M., and Villata, S. (2017). Towards artificial argumentation. *AI Mag.*, 38(3):25–36.
- [Baroni and Giacomin, 2006] Baroni, P. and Giacomin, M. (2006). Refining scc decomposition in argumentation semantics: a first investigation.
- [Baroni et al., 2005] Baroni, P., Giacomin, M., and Guida, G. (2005). Scc-recursiveness: a general schema for argumentation semantics. *Artif. Intell.*, 168(1-2):162–210.
- [Dung, 1995] Dung, P. M. (1995). On the acceptability of arguments and its fundamental role in nonmonotonic reasoning, logic programming and n-person games. *Artif. Intell.*, 77(2):321–358.
- [Dvorák, 2012] Dvorák, W. (2012). *Computational Aspects of Abstract Argumentation*. PhD thesis, Technische Universität Wien.
- [Dvorák and Dunne, 2018] Dvorák, W. and Dunne, P. E. (2018). Computational problems in formal argumentation and their complexity. *Handbook of formal argumentation*, 4:631–688.
- [Fan and Toni, 2015a] Fan, X. and Toni, F. (2015a). On computing explanations in argumentation. In Bonet, B. and Koenig, S., editors, *Proceedings of the Twenty-Ninth AAAI Conference on Artificial Intelligence, January 25-30, 2015, Austin, Texas, USA*, pages 1496–1502. AAAI Press.
- [Fan and Toni, 2015b] Fan, X. and Toni, F. (2015b). On explanations for non-acceptable arguments. In Black, E., Modgil, S., and Oren, N., editors, *Theory and Applications of Formal Argumentation - Third International Workshop, TFA 2015, Buenos Aires, Argentina, July 25-26, 2015, Revised Selected Papers*, volume 9524 of *Lecture Notes in Computer Science*, pages 112–127. Springer.
- [Lafages, 2021] Lafages, M. (2021). *Algorithms for enriched abstract argumentation frameworks for large-scale cases*. PhD thesis, Université Paul Sabatier-Toulouse III.
- [Liao, 2013] Liao, B. (2013). Toward incremental computation of argumentation semantics: A decomposition-based approach. *Ann. Math. Artif. Intell.*, 67(3-4):319–358.
- [Sakai et al., 2018] Sakai, K., Inago, A., Higashinaka, R., Yoshikawa, Y., Ishiguro, H., and Tomita, J. (2018). Creating large-scale argumentation structures for dialogue systems. In Calzolari, N., Choukri, K., Cieri, C., Declerck, T., Goggi, S., Hasida, K., Isahara, H., Maegaard, B., Mariani, J., Mazo, H., Moreno, A., Odijk, J., Piperidis, S., and Tokunaga, T., editors, *Proceedings of the Eleventh International Conference on Language Resources and Evaluation, LREC 2018, Miyazaki, Japan, May 7-12, 2018*. European Language Resources Association (ELRA).
- [Saribatur et al., 2020] Saribatur, Z. G., Wallner, J. P., and Woltran, S. (2020). Explaining non-acceptability in abstract argumentation. In Giacomo, G. D., Catalá, A., Dilkina, B., Milano, M., Barro, S., Bugarín, A., and Lang, J., editors, *ECAI 2020 - 24th European Conference on Artificial Intelligence, 29 August-8 September 2020, Santiago de Compostela, Spain, August 29 - September 8, 2020 - Including 10th Conference on Prestigious Applications of Artificial Intelligence (PAIS 2020)*, volume 325 of *Frontiers in Artificial Intelligence and Applications*, pages 881–888. IOS Press.
- [Ulbricht and Baumann, 2019] Ulbricht, M. and Baumann, R. (2019). If nothing is accepted - repairing argumentation frameworks. *J. Artif. Intell. Res.*, 66:1099–1145.

Towards a Fully-fledged Formal Protocol for the Explanation-Question-Response Dialogue

Federico Castagna

School of Computer Science, University of Lincoln, UK

Abstract

The increasing interest in eXplainable Artificial Intelligence (XAI) lies within the opacity and convoluted AI-based systems output generation, which is more than obscure for laypeople. A dialectical interaction with such systems may enhance the users' understanding and build a more robust trust towards AI. Commonly employed as specific formalisms for modelling intra-agents communications, *argumentation-based dialogue games* prove to be useful tools to rely upon when dealing with user's explanation needs. Among the literature, there exist some dialogical protocols that expressly handle explanations and their delivery. This research outlines the novel Explanation-Question-Response (EQR) dialogue and its properties, whose main feature is to provide satisfactory information whilst ensuring a simplified protocol (in comparison with other existing approaches) for both humans and artificial agents.

1 Introduction

The paradigm shift that described computation as multi-agent distributed cognition and interactions required the design of means of communication between such intelligent agents (i.e., software entities capable of flexible autonomous action in dynamic and unpredictable domains [Luck et al., 2005]). The challenge was to identify a widely applicable normative language that would not concede an indefinite number of replies. The choice fell then upon formal di-

alogue models: equipped with powerful expressivity, although still subject to specific regulations.

The urge to overcome ethical issues involving AI-based systems (e.g., consequences resulting from self-driving cars or recommendation systems), along with distrust from their users due to *black-box* algorithms, constitutes the reason for the recent interest in the field of eXplainable AI (XAI). The idea is that the trustworthiness of AIs can be improved by building more transparent and interpretable tools capable of explaining salient information about systems operations [Bellotti and Edwards, 2001]. Noticeably, Article 22 of the General Data Protection Regulation (GDPR)¹ introduces the right to obtain an explanation of the inferences produced by automated decision-making models. The necessity to abide by this new regulation contributes to making explainability a current hot topic in the AI research landscape. Although the most effective way to structure explanations is still an open problem among scholars, it is not uncommon to employ dialogical patterns. These patterns are characterised by a set of questions (e.g., how, why and what) and the answers presented by an explainer to an explainee who may acknowledge the responses or challenge them with counterfactual examples [Vilone and Longo, 2021]. The combination of compelling explainability demands with such dialogical frameworks engendered different protocols. Among these dialectical structures, the Explanation-Question-Response (EQR)² is a novel type of dia-

¹Regulation (EU) 2016/679

²First sketched in [McBurney and Parsons, 2021].

logue seeking to supply information to agents eschewing formalisms that would unnecessarily complicate the protocol. This paper provides a brief outline of the fully-fledged EQR dialogue and its properties, whose main features are: (1) the embeddings of multiple protocol types whilst (2) ensuring a simplified structure (both for humans and artificial agents) and (3) conveying satisfactory information.

2 Method

This research is underpinned by the theoretical model of *argumentation-based dialogue games*. Dialogue games are *rule-governed interactions* among players (i.e., agents with their own beliefs, goals, desires and only a small, possibly none, amount of information regarding the other players) that take turns in making utterances [McBurney and Parsons, 2009]. Dialogue games are commonly categorized according to elements such as information possessed by the participants at the commencement of the interaction, their individual goals, and the knowledge and goals they share with other agents [Walton and Krabbe, 1995]. Table 1 describes a (non-exhaustive) example of different dialogue types. Among these, of particular interest is the *query* type: a formal model that envisages when the questioner wants to hear and understand not just a claim itself but the arguments for the claim (i.e., a sort of ‘naive’ explanation). The selection and transitions between different dialogues can instead be rendered via a *Control Layer* [McBurney and Parsons, 2002], defined in terms of *atomic dialogue-types* and *control dialogues*. The latter are meta-structures that have as their topics other dialogues and contribute to the management of the protocols combinations and their transitions. Overall, a formal dialogue main components can be identified as:

Syntax The syntax of a dialogue game specifies the utterances available to the participants and the rules that govern the interactions among such utterances. In addition, it is standard to consider these utterances as composed of two layers: the innermost comprising the topics of

Dialogue Types	Description
<i>Information seeking/giving</i> [Hulstijn, 2000]	one participant seeks the answer to some question(s) from another agent (the information ‘giver’), who is believed by the first to know the answer(s)
<i>Persuasion</i> [Prakken, 2006]	one participant seeks to persuade another to accept a proposition it does not currently endorse
<i>Query</i> [Cogan et al., 2005]	one participant, X, keeps challenging the answer about p from another participant, Y. X’s interest lies more on Y’s argument for p rather than if Y believes p or not.

Table 1: Example of different dialogue protocols.

discussion and the outermost comprising the locutions.

Semantics Semantics differ according to the specific focus and final deployment of the dialogue:

- *Axiomatic*: defines each locution in terms of its pre and (possibly) post-conditions;
- *Operational*: considers each locution as a computational instruction that operates successively on the states of some abstract machine;
- *Denotational*: assigns, for each element of the language syntax, a relationship to an abstract mathematical entity, its denotation.

Pragmatics Pragmatics deals with those aspects of the language that do not involve considerations about truth and falsity, such as the illocutionary force of the utterances along with

speech acts, i.e., non-propositional utterances intended to or perceived to change the state of the world [McBurney and Parsons, 2007, McBurney et al., 2013].

3 Discussion

Argumentative models of explanation dialogues have already been introduced in the literature. For instance, [Bex and Walton, 2016] presented a protocol on such a topic with a complete list of locutions. However, to evaluate the provided explanation, the explainee needs to resort to a different dialogue (denoted *examination*). Similarly, [Madumal et al., 2019] devised a study for modelling explanation dialogues by following a data-driven approach. The resulting formalisation embeds (possibly several) argumentation-based dialogues nested in the outer layer of the explanation protocol. Finally, the dialogical structure proposed by [Sassoon et al., 2019] in the context of explanations for wellness consultation also exploits multiple dialogue types (i.e., persuasion, deliberation and information seeking) and their respective protocols whilst mostly focusing on the course of action to undertake. This is different from the anticipated EQR dialogue (sketched in [McBurney and Parsons, 2021] as Explanation-Question-Response), whose protocol is halfway between *persuasion*, *information-giving/seeking* and *query* and more comprehensively incorporates locutions for handling each of these tasks without the need for adopting a *Control Layer* or switching between dialogues. This allows for a simpler structure since it provides a model more suited to deliver explanations eschewing formalisms that would unnecessarily complicate the protocol.

3.1 Outline of the EQR Structure

The dialogue involves an explainer (the proponent) and an explainee (the opponent). The goal of the first is to deliver all the requested information to its interlocutor and persuade it into acknowledging the explanation. The aim of the latter is instead to challenge every argument posited by the proponent until it ob-

tains a satisfactory explanation. The proponent will be the first to start the EQR dialogue by stating the requested explanation after which the opponent will make its move. Henceforth, the two participants will alternate in turns. Figure 1 and Figure 2 identify the ordered sequences of locutions describing the turns of each player. The dashed arrows denote moves that must be performed during the first turn only (e.g., in all the subsequent turns, the player will start from the locution *turn start* rather than *enter dialogue*). Notice also that the opponent will always prefer to utter a *deny initial 'something'* rather than a *deny 'something'* or *ask 'something'*. This preference is emphasized in Figure 2 by the different (higher) positions of the locutions.

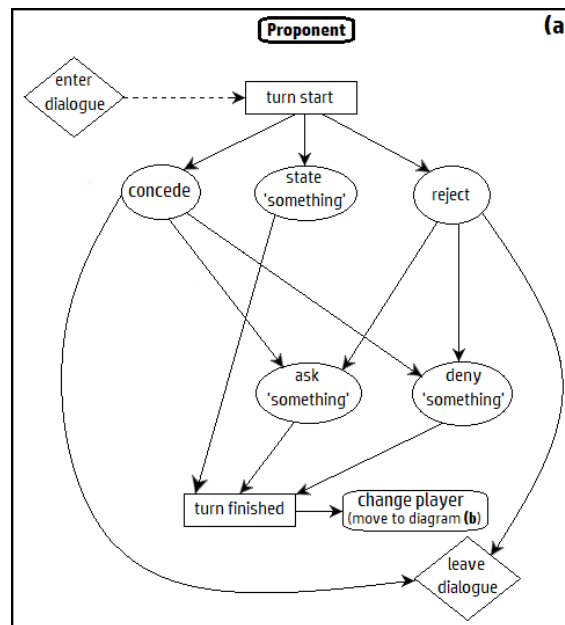


Figure 1: Overview of PRO's available moves.

Intuitively, we can informally define the main locutions as follows:

- **Ask** aims at demanding the ground on which it is the case that 'something' (say a) and not otherwise. As such, it can be seen as an argument attacking another argument on a.

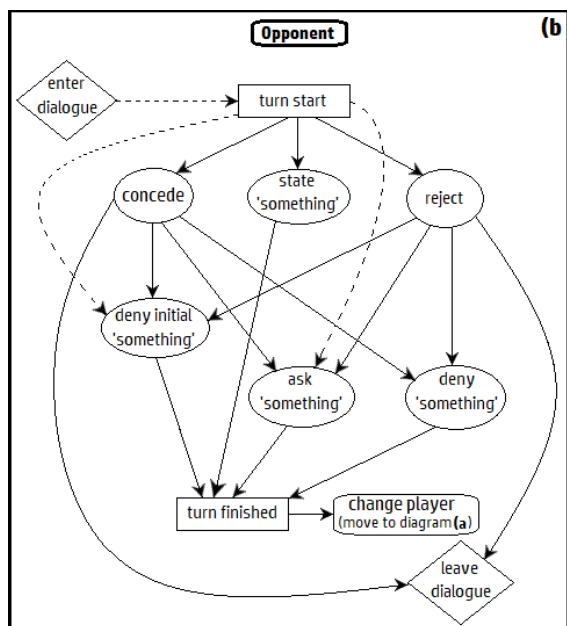


Figure 2: Overview of OPP’s available moves.

- **State** answers the query moved by the ask locution³. This will establish the reason why the questioned ‘something’ is the case. State can then be seen as identifying such rationale via an argument attacking the argument that posited the previous question.
- **Deny** denotes a refutation against the ‘something’ (say a) it is addressed. As such, it can be straightforwardly seen as an argument attacking another argument on a. The same reasoning can be applied to deny initial (which is merely a deny locution that directly targets the initial explanation).
- The role of locutions as **enter/leave dialogue**, **turn start/finished**, **concede**, **reject** and **change player** is, instead, to administer the protocol in a more structural way by identifying specific phases of each turn and different stages of the overall dialogue.

³Except during the proponent’s first turn where *state* corresponds just with the utterance of the initial explanation.

Finally, the EQR dialogue embeds an axiomatic semantics that presents the pre-conditions necessary for the legal utterance of each locution, and any post-conditions arising from their legal utterance. Such conditions influence the *commitment store* of each agent participating in the dialogue. These commitment stores are public statements that the agents have to defend in the dialogue, but they might not correspond to the agent’s real beliefs or intentions.

4 Conclusion

Stemming as a model for addressing XAI concerns, in this paper, we have outlined the fundamental characteristics of the Explanation-Question-Response dialogue, a novel approach to argumentation-based dialectical explanations. These features include (1) the embeddings of multiple protocol types whilst (2) ensuring a simplified structure for the involved agents (both artificial and human) and (3) conveying satisfactory information. Future works will shift from a fully-fledged formalisation of the dialogue and its locutions (objective already achieved, although space constraints prevented us from appreciating the overall research output) to concrete applications of the established protocol and their evaluations via user studies. Indeed, with the exponential growth of interest towards AI decision-support systems and smart personal assistants [Knote et al., 2019], the argumentative interactions between virtual agents and laypeople are increasing alongside. This may also entail additional users’ explanation demands, which provide a suited setting for EQR dialogue implementations. Similarly, chatbots could benefit from the EQR protocol when handling their interlocutors’ information requests. For example, the conversational bot presented in [Castagna et al., 2022] (that also introduced a new argument scheme, the EQR scheme, which draws upon the essential properties of the homonymous dialogue) would improve its communication delivery by deploying the Explanation-Question-Response protocol.

Acknowledgements

This work was supported by a PhD scholarship from King's College London. I would also like to thank Peter McBurney for his keen insights and valuable guidance in the field of dialogue games.

References

- [Bellotti and Edwards, 2001] Bellotti, V. and Edwards, K. (2001). Intelligibility and accountability: human considerations in context-aware systems. *Human-Computer Interaction*, 16(2-4):193–212.
- [Bex and Walton, 2016] Bex, F. and Walton, D. (2016). Combining explanation and argumentation in dialogue. *Argument & Computation*, 7(1):55–68.
- [Castagna et al., 2022] Castagna, F., Parsons, S., Sassoon, I., and Sklar, E. I. (2022). Providing explanations via the EQR argument scheme. In *Computational Models of Argument: Proceedings of COMMA 2022*.
- [Cogan et al., 2005] Cogan, E., Parsons, S., and McBurney, P. (2005). New types of inter-agent dialogues. In *International Workshop on Argumentation in Multi-Agent Systems*, pages 154–168. Springer.
- [Hulstijn, 2000] Hulstijn, J. (2000). *Dialogue models for inquiry and transaction*. Citeseer.
- [Knote et al., 2019] Knote, R., Janson, A., Söllner, M., and Leimeister, J. M. (2019). Classifying smart personal assistants: An empirical cluster analysis.
- [Luck et al., 2005] Luck, M., McBurney, P., Shehory, O., and Willmott, S. (2005). Agent technology: computing as interaction (a roadmap for agent based computing).
- [Madumal et al., 2019] Madumal, P., Miller, T., Sonnenberg, L., and Vetere, F. (2019). A grounded interaction protocol for explainable artificial intelligence. *arXiv preprint arXiv:1903.02409*.
- [McBurney and Parsons, 2002] McBurney, P. and Parsons, S. (2002). Games that agents play: A formal framework for dialogues between autonomous agents. *Journal of logic, language and information*, 11(3):315–334.
- [McBurney and Parsons, 2007] McBurney, P. and Parsons, S. (2007). Retraction and revocation in agent deliberation dialogs. *Argumentation*, 21(3):269–289.
- [McBurney and Parsons, 2009] McBurney, P. and Parsons, S. (2009). Dialogue games for agent argumentation. In *Argumentation in artificial intelligence*, pages 261–280. Springer.
- [McBurney and Parsons, 2021] McBurney, P. and Parsons, S. (2021). Argument schemes and dialogue protocols: Doug walton's legacy in artificial intelligence. *Journal of Applied Logics*, 8(1):263–286.
- [McBurney et al., 2013] McBurney, P., Parsons, S., Atkinson, K., Prakken, H., and Wyner, A. (2013). Talking about doing.
- [Prakken, 2006] Prakken, H. (2006). Formal systems for persuasion dialogue. *The knowledge engineering review*, 21(2):163–188.
- [Sassoon et al., 2019] Sassoon, I., Kökciyan, N., Sklar, E., and Parsons, S. (2019). Explainable argumentation for wellness consultation. In *International Workshop on Explainable, Transparent Autonomous Agents and Multi-Agent Systems*, pages 186–202. Springer.
- [Vilone and Longo, 2021] Vilone, G. and Longo, L. (2021). Notions of explainability and evaluation approaches for explainable artificial intelligence. *Information Fusion*, 76:89–106.
- [Walton and Krabbe, 1995] Walton, D. and Krabbe, E. C. (1995). *Commitment in dialogue: Basic concepts of interpersonal reasoning*. SUNY press.

Argumentation without Opposition?

Giulia D'Agostino

Institute of Argumentation, Linguistics and Semiotics, Università della Svizzera Italiana,
Switzerland

Abstract

Argumentation is a practice that originates from the challenging of a standpoint in a dialectical environment, one would expect to find as the departing standpoint in a text about argumentation in context. What we challenge here, however, is the assumption that all instances of argumentation need to arise from an opposition. In such a context, traditional argumentation mining techniques might not be effective and in need for redirection and, therefore, we propose the application of a newly identified pipeline for automatic extraction of non-contrastive arguments. The proof of concept will be left to an unassuming category yet initiator of argumentative patterns – namely, requests of elaboration.

1 Introduction

Many theories of argumentation focus on the conflict that drives arguments forward.

From a computational perspective, taking into account Dung's abstract argumentation theory [Dung, 1995] as primary reference, we introduce such fundamental notion by recalling that an Argumentation Framework is a directed graph in which arguments (nodes) are related by binary attacks (arcs) [Modgil, 2014]. Further and subsequent extensions of the original framework too (see for instance the bipolar (BAF, [Amgoud et al., 2008]) and weighted (WAF, [Dunne et al., 2011]) ones, not to mention the framework including recursive attacks (AFRA, [Cayrol et al., 2018])) always foresee the attack rela-

tion as the underlying driving force of argumentative occurrences.

On the other hand, in a more analogue perspective, the argumentation theory domain – though not always agreeing on some aspects such as the nature, number, and characterization of argument schemes – consistently illustrate that each instance of argumentation (a) is dialectical in nature (if not dialogical) and (b) starts with a confrontation between parties, i.e., the explicit or implicit challenging of a standpoint on which the parties disagree. This is already present in Toulmin's work [Toulmin, 1958], but it's probably Pragma-Dialectics that fully displays the potential and the implications of such view [van Eemeren and Grootendorst, 2004], with the explicit listing of stages and functional aspects that at the same time characterize and define an argumentative exchange. As a result, we should be able to acknowledge that the purpose of adopting an argumentative exchange is to settle some disputed issue between parties [Walton, 2012].

Deriving from the general definitions above, the application of a model belonging to either perspective foresee the inescapable instantiation of attacks (or conflicts), which formally represent practical cases either of disagreement or challenge. This is however not always the case: in some instances – the validity of which may be defined by the activity type they pertain to – argumentation does not derive from a conflict between parties or the manifestation of contrasting standpoints. Thus, such occurrences are critical on two levels: first and foremost, they are not trivial to trace and annotate even manually; besides,

it is challenging to justify procedurally or generalize inductively the course of action that led to their identification. As a consequence of the latter, they add a further layer of complexity to the argumentation mining task – which then needs to retrieve hidden features of covert contrast.

Therefore, the present paper poses two preliminary objectives: (1) shedding light on a theoretical issue which is common among similar activity types and (2) proposing a pipeline for annotation and subsequent automated retrieval in such domain – aware of the fact that some traditional assumptions might not work, or not be appropriate.

2 Argumentation in context

2.1 The structure of a Q&A session

The context for which the theoretical perspective outlined above does not seem to apply are "Question and Answer" (Q&A) types of interaction. In Q&A situations, only a challenge would overtly indicate a difference of opinion between actors and across moves: the speaker calls out the interlocutor on the justification of a stance of which they are known or expected to be an advocate. However, challenges seem to be almost nonexistent in corpora that gather such exchanges in certain, rather formulaic, domains: a clear instance is provided by the QT30 corpus [Hautli-Janisz et al., 2022], where challenges represent 0.3% of all moves; but also our corpus of ongoing development (see further) shows the same tendency, since requests of justification constitute the absolute minority: see Table 1 for reference. Therefore, even in the case there *is* an opposition, such a condition is not in any way formally traceable – even though a human interpreter would in most cases recognize there is some sort of conflict. Hence, the issue at stake is (a) finding out the nature of the human impression of disagreement, and ultimately (b) translate it into machine-readable instructions.

On the basis of qualitative observations, we claim that any instance of Q&A session – regardless of the specific context that generated it – shares certain characteristics:

Request type	count	percentage
clarification	83	9.47
commitment	15	1.71
confirmation	64	7.31
data	59	6.74
elaboration	373	42.58
explanation	91	10.39
justification	5	0.57
opinion	186	21.23
<i>Total</i>	<i>876</i>	<i>100.00</i>

Table 1: Distribution of request types, ECCs corpus.

- Three alternative formulations: they can either be "open", request the "polar" verification of a proposed assertion (yes/no), propose a set of "alternative" answers, among which the respondent is expected to pick one and only one.
- Not all utterances that occur formally as questions are real questions (i.e., they don't contemplate the elicitation of an answer) and, conversely, not all questions have an interrogative form.

Additionally, the vast majority of questions display a preface (in our corpus: more than 71%; for our intended meaning of such functional unit and its role in the activity type see [Lucchini et al., 2022]). This means that there is a prevalence of questions that explicitly justify either the appropriateness of the questioning act, or the validity of the content that is introduced as a premise to the question, or both. In other words, such questions result to be supported by the argumentative sequence accompanying them.

2.2 Corpus and case study

The corpus on which we performed a first cycle of manual annotation consists of 24 Earnings Conference Calls (ECCs) of 6 listed companies for the financial year 2021, for a total of 1,002,388 words. ECCs play a pivotal role in the financial communication domain – being a standardized activity where any deviation from the norm necessarily hints at some deeper imbalance, the effects of which are retrievable in the

fluctuations of stock prices in the financial market. Moreover, neighboring domains might as well take advantage of the clear-cut results drawn from ECCs and apply them to akin exchanges. Within the dialogical exchanges performed in ECCs, our project aims at the retrieval of so-called argumentative patterns, first defined by [van Eemeren, 2016] and revised by our research group as a working concept [Rocci et al., 2022].

Given such context, in this paper we will provide insight into the peculiar case represented by patterns originating from the performing of a request of elaboration. Requests of elaboration are indeed particularly relevant in the activity type and for the issue at stake because they don't show any sign of opposition whatsoever – in fact, not even human raters would argue they represent a challenging move, aside from the baseline concept that, being a question, it prompts the interlocutor to somewhat broaden the common ground. Nonetheless, they correspond to the majority of request types.

3 Pipeline

3.1 A two-stage manual annotation

The first stage of our study consisted in a multi-layered manual annotation on a general-purpose platform (INCEpTION, [Klie et al., 2018]), followed by a detailed argumentative reconstruction in OVA [Janier et al., 2014], thus implementing a (preliminary) pipeline that outputs *strong labeled data* [Shnarch et al., 2018]. For a further description of the procedure see [Lucchini and D'Agostino, 2022].

With respect to the context and for the specific type of requests described above, the aim is to progress from a basic "manual annotation and extraction" process to a more refined way to describe the pattern and automate its recognition – or at least, some features of it. The reason why we need to define, test, and evaluate a new technique for the mining of our patterns is the realization that even if we applied traditional means of retrieval of argumentation, they would probably fall short of a (correct) extraction for two main reasons: (1) the argumen-

tative structure is not juxtaposed but split between the question and the answer, where the latter might be quite distant from the former. We shall therefore refine and optimize for the specific context some ML techniques already proven to be effective in text segmentation and Q&A recognition, such as those described in [Devlin et al., 2019] and (2) there's no traceable conflict between the two components, as just stated in Section 1. Moreover, there's usually no overt (cross)reference between standpoints or arguments of question and answer: each turn is, apparently, formally enclosed.

Therefore, the finalized pipeline we define below will make use of argumentation mining techniques with the primary goal of retrieving patterns initiated by requests of elaboration; to do so, three sub-tasks are identified, namely to recognize and compute the relationship between:

- A question and its related answer
- Each move and the argumentation provided within the same turn
- Argumentation(s) in the two moves

3.2 Argument mining addendum: a proposal

Recent advancements in AM undoubtedly showed the benefits offered by argumentation mining approaches (see [Lawrence and Reed, 2019] for a comprehensive review). The starting point for a computational perspective on argument(ation) identification and evaluation is the acknowledgement that it requires three main stages: (1) the identification, segmentation, and classification of argumentative discourse units (ADUs), (2) the identification and classification of the relations between ADUs [Peldszus and Stede, 2013], and (3) the identification of argument schemes, namely the implicit and explicit inferential relations within and across ADUs [Macagno and Walton, 2014]. It should be noted that the process applies both for manual and automatic annotation and, therefore, may characterize both approaches in compatible ways – hence naturally leading to an operative chain or, introducing a

feedback function for dynamic improvement, a cycle. Having already applied it to a fair extent to a first full cycle of manual annotation and analysis, times appear to be mature enough to draw a first proposal of how we intend to arrange the incorporation of a computational module at the bottom of what we found to be a functional, analog pipeline.

The segmentation of ADUs is somewhat critical for a start, but we might be provisionally satisfied with proposition segmentation by means of two Naïve Bayes classifiers, as proposed in [Lawrence et al., 2014]. Provided the segmentation task could be taken for granted and carried out automatically with satisfactory performance, the most meaningful contribution would be an informed detection and analysis of argumentative patterns in the context defined above. On this aspect, our proposal is twofold:

- Development of AI techniques (models) with particular consideration to the type and peculiar characteristics of the object under scrutiny for retrieval - both in terms of contextual information and defining features of the pattern itself
- Creation of an optimal and representative dataset of manually annotated examples with which to feed a ML supervised learning cycle, later to be (again manually) verified for further refinement

4 Results and discussion

Preliminary results show that argumentation is not absent in the answers to requests of elaboration; moreover, there seems to be a positive correlation between argumentation in the question and argumentation in the answer, as shown in [D’Agostino, 2022] for an exploratory sample: a question the validity of which has been argued for has high(er) chances to receive an argumentative response (or semi-argumentative, including here explanations; see again [D’Agostino, 2022] for an account on the soundness of such a choice). The relationship between the two turns is therefore to be accounted as of (implicit)

opposition, since it prompts formally unrequited justification – though likely not of attack, since there’s no contrasting standpoint.

The meaning of such a spontaneous instance of argumentation is unclear. We propose two alternative interpretations:

- It is a preemptive justification based on a combination of presumptions’ exploitation, assumptions on the questioner’s standpoint on the matter, and on the anticipation of (plausible) attacks to the answerer’s standpoint (i.e., *procatalepsis*)
- Is it a tentative construction of what we could call argued undefeasibility; to effortlessly (try to) include the – in fact defeasible – standpoint to the *knowledge base* of the interaction

Either way, to further investigate the phenomenon and possibly determine which (if any) interpretation is correct we first need more data for both qualitative and quantitative analysis. Hence, argumentation mining will play a crucial role in the retrieval of potential patterns, both as an intermediate stage and as a final goal.

5 Conclusion

In this contribution, an alternative view on argumentation – not relying on the concept of opposition/attack – has been proposed. An empirical domain constituting an *activity type* was presented as a context of application of such a claim, having already been the object of an exploratory qualitative study based on manual annotation. Moreover, an AI module for argumentation mining in such context was presented and discussed as an integration and expansion of the pipeline previously applied to the study. Further discussed was a plan to collect a significant number of occurrences of the pattern both to enhance a more comprehensive understanding of this instance of non-contrastive argumentation and to develop a well-established ML model for its retrieval.

Acknowledgements

The author thanks Andrea Rocci and Chris Reed for their advice and support.

References

- [Amgoud et al., 2008] Amgoud, L., Cayrol, C., Lagasque-Schiex, M. C., and Livet, P. (2008). On bipolarity in argumentation frameworks. *International Journal of Intelligent Systems*, 23(10):1062–1093.
- [Cayrol et al., 2018] Cayrol, C., Fandinno, J., Fariñas del Cerro, L., and Lagasque-Schiex, M.-C. (2018). Argumentation frameworks with recursive attacks and evidence-based supports. In Ferrarotti, F. and Woltran, S., editors, *10th International Symposium on Foundations of Information and Knowledge Systems (FoIKS 2018)*, volume 10833 of *Lecture Notes in Computer Science book series (LNCS)*, pages 150–169. Springer.
- [D’Agostino, 2022] D’Agostino, G. (2022). ”(so long, and) thanks for all the color”. requests of elaboration and answers they trigger in earnings conference calls. In *Proceedings of the 22nd Edition of the Workshop on Computational Models of Natural Argument (CMNA 22)*, pages 80–85. CEUR.
- [Devlin et al., 2019] Devlin, J., Chang, M.-W., Lee, K., and Toutanova, K. (2019). BERT: Pre-training of deep bidirectional transformers for language understanding. In *NAACL HLT 2019 - 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies - Proceedings of the Conference*, pages 4171–4186.
- [Dung, 1995] Dung, P. M. (1995). On the acceptability of arguments and its fundamental role in nonmonotonic reasoning, logic programming and n-person games. *Artificial Intelligence*, 77(2):321–357.
- [Dunne et al., 2011] Dunne, P. E., Hunter, A., McBurney, P., Parsons, S., and Wooldridge, M. (2011). Weighted argument systems: Basic definitions, algorithms, and complexity results. *Artificial Intelligence*, 175(2):457–486.
- [Hautli-Janisz et al., 2022] Hautli-Janisz, A., Kikiteva, Z., Siskou, W., Gorska, K., Becker, R., and Reed, C. (2022). QT30: A corpus of argument and conflict in broadcast debate. In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pages 3291–3300. European Language Resources Association.
- [Janier et al., 2014] Janier, M., Lawrence, J., and Reed, C. (2014). OVA+: an argument analysis interface. *Computational Models of Argument*, pages 463–464.
- [Klie et al., 2018] Klie, J.-C., Bugert, M., Boullosa, B., Eckart de Castilho, R., and Gurevych, I. (2018). The INCEpTION platform: Machine-assisted and knowledge-oriented interactive annotation. In *Proceedings of the 27th International Conference on Computational Linguistics: System Demonstrations*, pages 5–9. Association for Computational Linguistics.
- [Lawrence and Reed, 2019] Lawrence, J. and Reed, C. (2019). Argument mining: A survey. *Computational Linguistics*, 45(4):765–818.
- [Lawrence et al., 2014] Lawrence, J., Reed, C., Allen, C., McAlister, S., and Ravenscroft, A. (2014). Mining arguments from 19th century philosophical texts using topic based modelling. In *Proceedings of the First Workshop on Argumentation Mining*, pages 79–87. Association for Computational Linguistics.
- [Lucchini and D’Agostino, 2022] Lucchini, C. and D’Agostino, G. (2022). Good answers, better questions. building an annotation scheme for financial dialogues. Unpublished manuscript.
- [Lucchini et al., 2022] Lucchini, C., Rocci, A., and D’Agostino, G. (2022). Annotating argumentation within questions. prefaced questions as genre

- specific argumentative pattern in earnings conference calls. In *Proceedings of the 22nd Edition of the Workshop on Computational Models of Natural Argument (CMNA 22)*, pages 61–66. CEUR.
- [Macagno and Walton, 2014] Macagno, F. and Walton, D. (2014). Argumentation schemes and topical relations. In Gobber, G. and Rocci, A., editors, *Language, reason and education*, pages 185–216. Peter Lang.
- [Modgil, 2014] Modgil, S. (2014). Revisiting Abstract Argumentation Frameworks. In Black, E., Modgil, S., and Oren, N., editors, *Theory and Applications of Formal Argumentation*, Lecture Notes in Computer Science, pages 1–15, Berlin, Heidelberg. Springer.
- [Peldszus and Stede, 2013] Peldszus, A. and Stede, M. (2013). From argument diagrams to argumentation mining in texts. *International Journal of Cognitive Informatics and Natural Intelligence*, 7(1):1–31.
- [Rocci et al., 2022] Rocci, A., Yaskorska-Shah, O., D’Agostino, G., and Lucchini, C. (2022). Argumentative patterns initiated by closed-list questions in accountability dialogues. a corpus study of financial conference calls. In *Proceedings of the 4th European Conference on Argumentation - ECA 2022*.
- [Shnarch et al., 2018] Shnarch, E., Alzate, C., Dankin, L., Gleize, M., Hou, Y., Choshen, L., Aharonov, R., and Slonim, N. (2018). Will it blend? blending weak and strong labeled data in a neural network for argumentation mining. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 599–605. Association for Computational Linguistics.
- [Toulmin, 1958] Toulmin, S. (1958). *The Uses of Argument*. Cambridge University Press.
- [van Eemeren, 2016] van Eemeren, F. H. (2016). Identifying argumentative patterns: A vital step in the development of pragma-dialectics. *Argumentation*, 30(1):1–23.
- [van Eemeren and Grootendorst, 2004] van Eemeren, F. H. and Grootendorst, R. (2004). *A Systematic Theory of Argumentation: The Pragma-dialectical Approach*. Cambridge University Press.
- [Walton, 2012] Walton, D. (2012). Using argumentation schemes for argument extraction: A bottom-up method. *International Journal of Cognitive Informatics and Natural Intelligence*, 6(3):33–61.

Justification, Stability and Relevance for Transparent and Efficient Human-in-the-Loop Decision Support

Daphne Odekerken^{1,2}

¹Department of Information and Computing Sciences, Utrecht University

²National Police Lab AI, Netherlands Police

Abstract

One of the promises of artificial intelligence is improving efficiency in various processes, including decision-making. For specific decisions it is vital that human experts understand and are able to influence machine-made advice. In my dissertation research, I design and study argumentation-based systems for transparent human-in-the-loop decision support. Based on a domain-specific argumentation setting, these systems are able to construct an initial advice on some decision (*justification*); investigate the possibility that additional, yet uncertain, information can change the conclusion (*stability*) and if so, which information is worth investigating (*relevance*). The systems' requirements of detecting justification, stability and relevance correspond to theoretical problems in computational argumentation, most of which are in high complexity classes. In order to achieve reasonable estimations for these problems in polynomial time, I develop and investigate not only exact algorithms but also approximations.

1 Introduction

Artificial intelligence (AI) is developing fast, promising quality and efficiency benefits in various processes where (repetitive) tasks are outsourced from human analysts to machines. At the Netherlands Police, two examples of tasks in which AI could be helpful are the intake of complaints and the classification of suspect web shops. Especially for complex or high-risk

decision-making tasks it is essential that domain experts understand machine-made decisions or advice and have the possibility to correct possible mistakes [European Commission, 2021]. In addition, there are decisions that require input that cannot be obtained by a machine; instead, an analyst has to be consulted. We therefore need AI systems that are able to explain their decisions and to keep the human in the loop. Since computational argumentation is a research topic concerning reasoning with incomplete or inconsistent information in a way similar to human reasoning [Atkinson et al., 2017], it seems promising to use argumentation-based techniques for developing the required systems. In my dissertation research, I propose a general decision-making approach centered around three theoretical problems in computational argumentation: justification, stability and relevance.

Informally, the problem (or task) of determining *justification* can be seen as giving an initial advice regarding a specific topic, thereby only considering information that is currently available. The topic satisfies *stability* if the corresponding advice will not change, regardless of currently unavailable information that could still be added and/or currently available information that could be removed. In situations where the topic is not stable, one should identify information that is *relevant*, in the sense that investigation into its presence possibly leads to a stable topic. I will define these three problems on three settings: structured and abstract argumentation frameworks, as well as precedent-based reasoning.

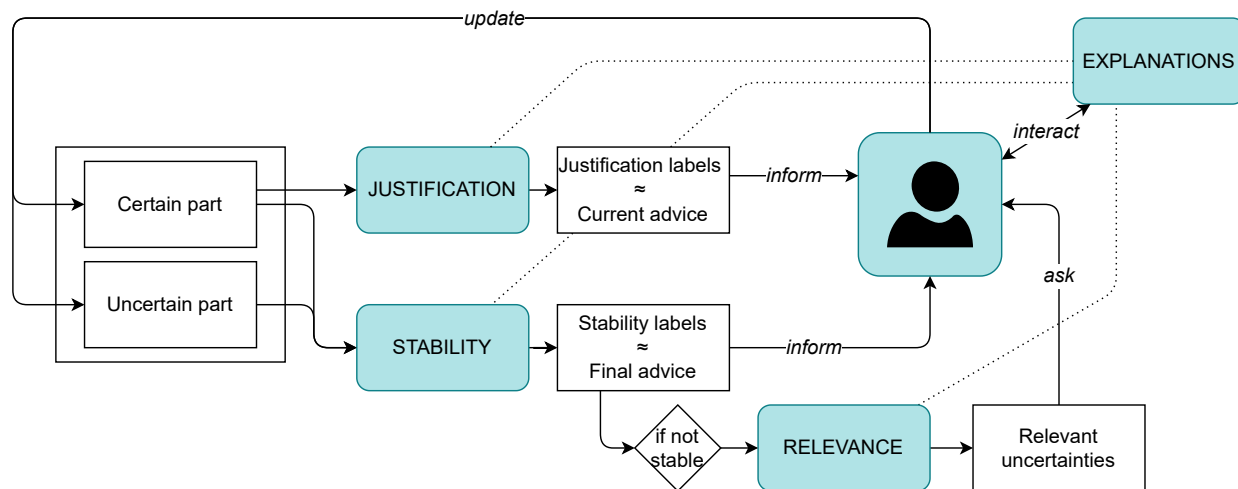


Figure 1: High-level overview of the proposed human-in-the-loop decision-making process, involving algorithms for justification, stability and relevance identification.

2 Method

Before formally defining the problems of justification, stability and relevance, I first outline the proposed human-in-the-loop decision-making procedure, illustrated in Figure 1, which relies on practical solutions for these theoretical problems. First, an algorithm for justification is used to obtain an initial decision or advice, given only the current information. Additionally, an algorithm for stability detection is applied on both the current and yet uncertain information; if this algorithm identifies a stable situation, the advice can be considered final. Otherwise, the algorithm for relevance identifies those uncertainties that should be investigated by the analyst. After investigation, the analyst can return findings to the system. This process is repeated until a final advice is found (or the analyst decides not to investigate any further).

In the remainder of this section, I define the problems of justification, stability and relevance on both structured and abstract argumentation frameworks, as well as for *a fortiori reasoning* based on precedents. Due to limited space, I only give formal definitions of these problems for incomplete (abstract) argumentation frameworks and provide an informal description for the other two settings.

2.1 Definitions for IAFs

An argumentation framework (AF) $\langle \mathcal{A}, \mathcal{R} \rangle$ consists of a set \mathcal{A} of arguments and attack relation $\mathcal{R} \subseteq \mathcal{A} \times \mathcal{A}$, where $(A, B) \in \mathcal{R}$ indicates that argument A attacks argument B [Dung, 1995]. Given a specific semantics (see e.g. [Baroni et al., 2011]), one can determine the justification status of arguments in an AF, based on their presence in an extension (i.e. set of arguments):

Definition 1 (Argument justification status). *Let $AF = \langle \mathcal{A}, \mathcal{R} \rangle$ be an argumentation framework and σ some semantics. Let A be some argument in \mathcal{A} .*

- A is σ -scep-IN (resp. σ -cred-IN) iff A belongs to each (resp. some) σ -extension of AF ;
- A is σ -scep-OUT (resp. σ -cred-OUT) iff for each (resp. some) σ -extension S of AF , A is attacked by some argument in S ;
- A is σ -scep-UNDEC (resp. σ -cred-UNDEC) iff for each (resp. some) σ -extension of AF , A is not in S , nor attacked by any argument in S .

Incomplete argumentation frameworks (IAFs) are an extension to AFs that encode qualitative uncertainty regarding the presence of arguments and attacks [Baumeister et al., 2021]. An IAF is a tuple

$\langle \mathcal{A}, \mathcal{A}^?, \mathcal{R}, \mathcal{R}^? \rangle$, where $\mathcal{A} \cap \mathcal{A}^? = \emptyset$, $\mathcal{R} \cap \mathcal{R}^? = \emptyset$; \mathcal{A} is the set of certain arguments; $\mathcal{A}^?$ is the set of uncertain arguments; $\mathcal{R} \subseteq (\mathcal{A} \cup \mathcal{A}^?) \times (\mathcal{A} \cup \mathcal{A}^?)$ is the certain attack relation and $\mathcal{R}^? \subseteq (\mathcal{A} \cup \mathcal{A}^?) \times (\mathcal{A} \cup \mathcal{A}^?)$ is the uncertain attack relation. An IAF can be *specified* by investigating the presence of uncertain elements: a specification is an IAF $\langle \mathcal{A}', \mathcal{A}'^?, \mathcal{R}', \mathcal{R}'^? \rangle$, where $\mathcal{A} \subseteq \mathcal{A}' \subseteq \mathcal{A} \cup \mathcal{A}^?$; $\mathcal{R} \subseteq \mathcal{R}' \subseteq \mathcal{R} \cup \mathcal{R}^?$; $\mathcal{A}'^? \subseteq \mathcal{A}^?$ and $\mathcal{R}'^? \subseteq \mathcal{R}^?$. Given some IAF \mathcal{I} , we denote all possible specifications for \mathcal{I} by $F(\mathcal{I})$.

In [Odekerken et al., 2022b], we introduced the notion of stability for IAFs, where an argument is stable if and only if its justification status remains the same, regardless of the way the uncertain arguments and attacks would turn out to be present or absent.

Definition 2 (Stability). *Given an IAF $\mathcal{I} = \langle \mathcal{A}, \mathcal{A}^?, \mathcal{R}, \mathcal{R}^? \rangle$ with $A \in \mathcal{A}$ and justification status j , A is stable- j w.r.t. \mathcal{I} iff for each $\langle \mathcal{A}', \mathcal{A}'^?, \mathcal{R}', \mathcal{R}'^? \rangle$ in $F(\mathcal{I})$, A is j w.r.t. $\langle \mathcal{A}', \mathcal{R}' \cap (\mathcal{A}' \times \mathcal{A}'^?) \rangle$.*

In situations where the topic argument is stable, one can give a final advice; in all other situations, further investigation could lead to a different advice. In order to decide which arguments or attacks are worth investigating, we define relevance on IAFs in Definition 4. The notion of relevance is defined based on minimal stable specifications, which we introduced in [Odekerken et al., 2022b] and recall next.

Definition 3 (Minimal stable specification). *Given an IAF $\mathcal{I} = \langle \mathcal{A}, \mathcal{A}^?, \mathcal{R}, \mathcal{R}^? \rangle$, a certain argument $A \in \mathcal{A}$ and a justification status j , a minimal stable- j specification for A w.r.t. \mathcal{I} is a specification $\mathcal{I}' = \langle \mathcal{A}', \mathcal{A}'^?, \mathcal{R}', \mathcal{R}'^? \rangle$ in $F(\mathcal{I})$ such that A is stable- j in \mathcal{I}' and there is no $\mathcal{I}'' = \langle \mathcal{A}'', \mathcal{A}''^?, \mathcal{R}'', \mathcal{R}''^? \rangle$ in $F(\mathcal{I})$ such that A is stable- j in \mathcal{I}'' , $\mathcal{I}'' \neq \mathcal{I}'$ and $\mathcal{I}' \in F(\mathcal{I}'')$.*

Definition 4 (Relevance). *Given an IAF $\mathcal{I} = \langle \mathcal{A}, \mathcal{A}^?, \mathcal{R}, \mathcal{R}^? \rangle$, certain argument $A \in \mathcal{A}$, uncertain element $U \in \mathcal{A}^? \cup \mathcal{R}^?$ and justification status j ,*

- *Addition of U is j -relevant for A w.r.t. \mathcal{I} iff there is some minimal stable- j specification $\langle \mathcal{A}', \mathcal{A}'^?, \mathcal{R}', \mathcal{R}'^? \rangle$ for A w.r.t. \mathcal{I} such that $U \in \mathcal{A}' \cup \mathcal{R}'$;*
- *Removal of U is j -relevant for A w.r.t. \mathcal{I} iff there is some minimal stable- j specification*

$\langle \mathcal{A}', \mathcal{A}'^?, \mathcal{R}', \mathcal{R}'^? \rangle$ for A w.r.t. \mathcal{I} such that $U \notin \mathcal{A}' \cup \mathcal{A}'^? \cup \mathcal{R}' \cup \mathcal{R}'^?$.

Note that this argumentation-based approach offers ample opportunity for explanation. For example, the justification status IN of some argument A can be explained by returning one or more extensions containing A (see e.g. [Borg and Bex, 2021]). A stability status can be explained by showing all specifications and appropriate extensions, but, alternatively, also by returning a minimal stable specification containing A . One way of explaining relevance of U for A would be to give a specification where U is certainly present and A has a given justification status j , whereas A would not be j in the variation in which U is certainly absent.

2.2 Structured argumentation

Similar to (abstract) AFs, the notions of justification status, stability and relevance can be defined on (a dynamic version of) structured argumentation frameworks, such as ASPIC⁺. In [Odekerken et al., 2020], we extended ASPIC⁺ argumentation theories with a set of queryables \mathcal{Q} , containing those literals for which it is not yet known if they will be added to the knowledge base. Using the definition of justification for ASPIC⁺ [Modgil and Prakken, 2018], we defined stability on this extension in [Odekerken et al., 2020]. In future work, I plan to define and study the notion of relevance for this ASPIC⁺ extension as well.

2.3 Precedent-based reasoning

The idea of precedent-based reasoning is that decisions on cases generalise to, and thereby restrict possible outcomes of new, yet undecided, cases [Horty, 2011]. Normally, it is assumed that the factors of a new case are certain and this case is compared to a case base with earlier cases and the corresponding decisions. We view the task of decision-making given only certain information as the justification problem. However, in reality it is not always known which factors are present in a case: sometimes, this can only be determined by additional investigation. Just like for abstract and structured argu-

mentation frameworks, the framework for precedent-based reasoning can be extended with an uncertain component, on which the problems of stability and relevance can be defined. Specifically, we assume that the factors of cases in the case base are all certain, while a new case may have a combination of certain and uncertain factors. Whereas the justification task is to make a decision based on the *certain* factors of the case and the case base, the stability task is to decide if the addition or removal of *uncertain* information could still change this outcome; if so, identifying this information corresponds to the relevance task.

3 Computational Issues

The decision-making procedure proposed in the previous section requires efficient algorithms for detecting justification status, stability and relevance. However, formal complexity analysis reveals that most of the problems are in high complexity classes. For example, in [Odekerken et al., 2020] we proved that the problem of detecting stability is CoNP-complete under grounded semantics, given a simple implementation of ASPIC⁺. In [Odekerken et al., 2022b] we studied stability and relevance problems for IAFs and observed that these are highly complex as well.

It is therefore far from trivial to find a fast solution for each instantiation of the problem. This issue can be handled in various ways, as shown in the argumentation literature for (other) problems in high complexity classes. One possible solution are exact algorithms based on SAT-solvers (see e.g. [Baumeister et al., 2021]) or Answer Set Programming (e.g. [Lehtonen et al., 2020]). These algorithms may be relatively fast in practice, but fast computation is not guaranteed as they are exponential. A second option are approximation algorithms that are learned from data [Craandijk and Bex, 2020]: once trained, such algorithms are fast and quite accurate, but a theoretical accuracy analysis is not possible. I am therefore particularly interested in a third alternative: developing approximation algorithms that can be evaluated not only empirically, but also theoretically. In [Odekerken et al., 2020], we describe an approximation algorithm for estimating stability and

show that it is polynomial and sound, but not complete. In [Odekerken et al., 2022a], we compare this algorithm to an exact algorithm. In future work, I plan to develop and study approximation algorithms for the relevance problem as well. To assess the performance of these algorithms, I plan to conduct a theoretical analysis of time complexity and identify when the algorithm gives an exact solution, similar to our stability study in [Odekerken et al., 2022a].

At this moment, the general human-in-the-loop decision-making procedure is already applied for two specific applications at the police: an ASPIC⁺-based inquiry system that helps citizens to decide if they should submit a complaint on online trade fraud [Schraagen et al., 2018, Testerink et al., 2019, Odekerken et al., 2020], as well as a human-in-the-loop classifier of suspect web shops, based on a combination of structured argumentation and precedent-based reasoning [Odekerken and Bex, 2020]. In a future user study, I will investigate the analysts' experience and performance using the latter system, studying e.g. how they use the suggestions for relevance.

4 Conclusion

The application of argumentation-based techniques is a promising approach towards transparent human-in-the-loop support for complex or high-risk decisions. This abstract summarized my proposal for a general approach for decision support, centered around three theoretical problems: justification, stability and relevance. These problems are defined for various settings in computational argumentation. Given that most of these problems are situated in high complexity classes, I develop algorithms that obtain an estimation in polynomial time and evaluate them empirically and theoretically. The algorithms are applied in decision support systems at the Netherlands Police.

Acknowledgements

I would like to thank Floris Bex, AnneMarie Borg and Henry Prakken for their support throughout my PhD project.

References

- [Atkinson et al., 2017] Atkinson, K., Baroni, P., Giacomini, M., Hunter, A., Prakken, H., Reed, C., Simari, G., Thimm, M., and Villata, S. (2017). Towards artificial argumentation. *AI magazine*, 38(3):25–36.
- [Baroni et al., 2011] Baroni, P., Caminada, M., and Giacomin, M. (2011). An introduction to argumentation semantics. *The Knowledge Engineering Review*, 26(4):365–410.
- [Baumeister et al., 2021] Baumeister, D., Järvisalo, M., Neugebauer, D., Niskanen, A., and Rothe, J. (2021). Acceptance in incomplete argumentation frameworks. *Artificial Intelligence*, 295:103470.
- [Borg and Bex, 2021] Borg, A. and Bex, F. (2021). A basic framework for explanations in argumentation. *IEEE Intelligent Systems*, 36(2):25–35.
- [Craandijk and Bex, 2020] Craandijk, D. and Bex, F. (2020). Deep learning for abstract argumentation semantics. *Proceedings of the Twenty-Ninth International Joint Conference on Artificial Intelligence*, pages 1667–1673.
- [Dung, 1995] Dung, P. M. (1995). On the acceptability of arguments and its fundamental role in nonmonotonic reasoning, logic programming and n-person games. *Artificial Intelligence*, 77:321–357.
- [European Commission, 2021] European Commission (2021). Proposal for a regulation of the European Parliament and of the council laying down harmonised rules on artificial intelligence (Artificial Intelligence Act) and amending certain union legislative acts. <https://eur-lex.europa.eu/legal-content/EN/TXT/?uri=CELEX:52021PC0206> [Online; accessed 17 June 2022].
- [Horty, 2011] Horty, J. F. (2011). Rules and reasons in the theory of precedent. *Legal Theory*, 17:1–33.
- [Lehtonen et al., 2020] Lehtonen, T., Wallner, J. P., and Järvisalo, M. (2020). An answer set programming approach to argumentative reasoning in the ASPIC+ framework. In *Proceedings of the International Conference on Principles of Knowledge Representation and Reasoning*, volume 17, pages 636–646.
- [Modgil and Prakken, 2018] Modgil, S. and Prakken, H. (2018). Abstract rule-based argumentation. In Baroni, P., Gabbay, D., Giacomin, M., and van der Torre, L., editors, *Handbook of Formal Argumentation*, volume 1, pages 286–361. College Publications.
- [Odekerken and Bex, 2020] Odekerken, D. and Bex, F. (2020). Towards transparent human-in-the-loop classification of fraudulent web shops. In *Legal Knowledge and Information Systems*, pages 239–242.
- [Odekerken et al., 2022a] Odekerken, D., Bex, F., Borg, A., and Testerink, B. (2022a). Approximating stability for applied argument-based inquiry. *Intelligent Systems with Applications*, 16:200110.
- [Odekerken et al., 2020] Odekerken, D., Borg, A., and Bex, F. (2020). Estimating stability for efficient argument-based inquiry. In *Computational Models of Argument. Proceedings of COMMA 2020*, pages 307–318.
- [Odekerken et al., 2022b] Odekerken, D., Borg, A., and Bex, F. (2022b). Stability and relevance in incomplete argumentation frameworks. In *Computational Models of Argument. Proceedings of COMMA 2022*, pages 272–283.
- [Schraagen et al., 2018] Schraagen, M., Bex, F., Odekerken, D., and Testerink, B. (2018). Argumentation-driven information extraction for online crime reports. In *Proceedings of International Workshop on Legal Data Analytics and Mining*, pages 20–25.
- [Testerink et al., 2019] Testerink, B., Odekerken, D., and Bex, F. (2019). A method for efficient argument-based inquiry. In *Proceedings of the 13th International Conference on Flexible Query Answering Systems*, pages 114–125. Springer International Publishing.

Towards Preserving Semantic Structure in Argumentative Multi-Agent via Abstract Interpretation

Minal Suresh Patil

Umeå Universitet

Abstract

Over the recent twenty years, argumentation has received considerable attention in the fields of knowledge representation, reasoning, and multi-agent systems. However, argumentation in dynamic multi-agent systems encounters the problem of significant arguments generated by agents, which comes at the expense of representational complexity and computational cost. In this work, we aim to investigate the notion of *abstraction* from the model-checking perspective, where several arguments are trying to defend the same position from various points of view, thereby reducing the size of the argumentation framework whilst preserving the semantic flow structure in the system.

particular advantages on argumentation-based techniques, including transparent decision-making and the capability to provide a defensible rationale for outcomes. In our recent work, we propose the use of explanations in autonomous pedagogical scenarios [Patil, 2022] i.e. how explanations should be tailored in multi-agent systems (MAS) (teacher-learner interaction) as shown in Figure 1. It is rational to assume that autonomous agents in open, dynamic, and distributed systems will conform to a linguistic system for expressing their knowledge in terms of one or more ontologies that reflect the salient domain. Agents consequently must agree on the semantics (e.g. privacy) of the terms they use to organize the information, contextualise the environment, and represent different entities in order to engage or cooperate jointly. Abstract argumentation frameworks (AFs) [Dung, 1995, Bench-Capon and Dunne, 2007] are naturally employed for modelling and effectively resolving such types of challenges. In both multi-agent [Maudet et al., 2006] and single-agent [Amgoud and Prade, 2009] decision-making situations, AFs have been extensively utilised to describe behaviours since they can innately represent and reason with opposing information. Moreover, argumentative models have been presented due to the dialectic nature of AFs so that agents can cooperatively resolve issues or arrive at decisions by communicating implicitly [Dung et al., 2009].

Present studies of AFs, however, may not be immediately applicable to multi-agent scenarios where agents could come across certain unexpected circumstances in their environment. AFs are naturally

1 Introduction

Humans must possess beliefs in order to engage with their surrounding environments successfully, coordinate their activities, and be capable of communicating. Humans sometimes use arguments to influence others to act or realise a particular approach, to reach a reasonable agreement, and to collaborate together to seek the optimal possible solution to a particular problem. In light of this, it is not unexpected that many recent efforts to represent artificially intelligent agents have incorporated arguments and beliefs of their environment. Argumentation-based decision-making approaches are anticipated to be more in line with how people reason, consider possibilities and achieve objectives. This confers

used for modelling dynamic systems since, in actuality, the argumentation process is inherently dynamic in nature [Falappa et al., 2011, Booth et al., 2013] and this comes with high computational complexity [Dunne and Wooldridge, 2009, Dunne, 2009].

To give a practical example, autonomous Intent-Based Networking (IBN) [Campanella, 2019] captures and translates business intent into network policies that can be automated and applied consistently across the network. The goal is for the network to continuously monitor and adjust its performance to assure the desired business outcome. Intent allows the agent to understand the global utility and the value of its actions. Consequently, the autonomous agents can evaluate situations and potential action strategies rather than being limited to following instructions that human developers have specified in policies. In these cases, agents may adjust their model of the environment as well as their strategy according to information provided by the environment. There are several other circumstances in which the agent may not be able to guarantee a specific status of specific arguments and would necessitate assistance from other agents. Agents may not always know the optimal strategy until they form a coalition. In such circumstances, agents cannot merely compute semantics/conclusions from the ground up since it is not feasible. “Abstracting” AFs from the original (concrete domain) AF to via Abstract Interpretation can help compute semantics on a much smaller AF. Abstraction inherently are necessary if specific properties or specifications of the AF that are abstracted away from them are maintained.

The main contribution of this work is to investigate the semantic properties of the “abstract” AF from the “concrete” AF during the multi-agent interactions. The term “abstraction” in this work pertains to the notion of abstraction from model checking. Abstraction of the state space may reduce the AF to a manageable size by clustering similar concrete states into abstract states, which can further facilitate verifying these abstract states. We summarise the primary research question as follows: *Given a MAS in an uncertain environment, each with a specific subjective evaluation of a given set of conflicting arguments, how can agents reach a consensus whilst*

preserving specific semantic properties?

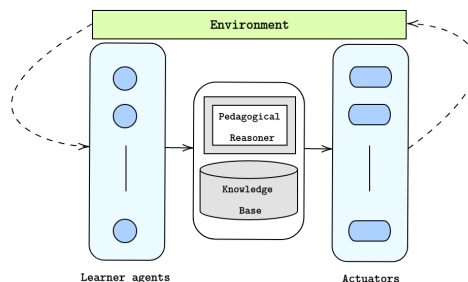


Figure 1: Pedagogical Multi-Agent Reasoning

2 Method

Model checking [Clarke et al., 2000] is widely accepted as a powerful automatic verification technique for the verification of finite-state systems. Halpern and Vardi proposed the use of model checking as an alternative to the deduction for logics of knowledge [Halpern and Vardi, 1991]. Since then, model checking has been extended to multi-agent systems [Hoek and Wooldridge, 2002]. The state explosion issue is the main impediment to the tractability of model checking. Nevertheless, significant research has been done on this well-known issue, and a variety of approaches have been proposed to circumvent the model checking limitation, including such symbolic methods with binary decision diagrams [Burch et al., 1992], SAT solvers [Biere et al., 1999], partial order reduction [Peled and Pnueli, 1994] and abstraction [Clarke et al., 2000].

In this work, we focus on abstract interpretation for computing dynamic semantics in MAS. The main point of abstract interpretation [Cousot and Cousot, 1977] is to replace the formal semantics of a system with an abstract semantics computed over a domain of abstract objects, which describe the properties of the system we are interested in. It formalises formal methods and allows to discuss the guarantees they provide such as soundness

(the conclusions about programs are always correct under suitable explicitly stated hypotheses), completeness (all true facts are provable), or incompleteness (showing the limits of applicability of the formal method). Abstract interpretation is mainly applied to design semantics, proof methods, and static analysis of programs. The semantics of programs formally defines all their possible executions at various levels of abstraction. Proof methods can be used to prove (manually or using theorem provers) that the semantics of a program satisfy some specification, that is a property of executions defining what programs are supposed to do. Now, we provide a brief technical primer on key concepts in abstraction interpretation.

Posets: A partially ordered set (poset) $\langle \mathcal{D}, \sqsubseteq \rangle$ is a set \mathcal{D} equipped with a partial order \sqsubseteq that is (1) reflexive: $\forall x \in \mathcal{D} \cdot x \sqsubseteq x$; (2) antisymmetric: $\forall x, y \in \mathcal{D} \cdot ((x \sqsubseteq y) \wedge (y \sqsubseteq x)) \Rightarrow (x = y)$; and (3) transitive: $\forall x, y, z \in \mathcal{D} \cdot ((x \sqsubseteq y) \wedge (y \sqsubseteq z)) \Rightarrow (x \sqsubseteq z)$. Let $S \in \wp(\mathcal{D})$ be a subset of the poset $\langle \mathcal{D}, \sqsubseteq \rangle$, then the least upper bound (*lub/join*) of S (if any) is denoted as $\sqcup S$ such that $\forall x \in S \cdot x \sqsubseteq \sqcup S$ and $\forall u \in S \cdot (\forall x \in S \cdot x \sqsubseteq u) \Rightarrow \sqcup S \sqsubseteq u$, and the greatest lower bound (*glb/meet*) of S (if any) is denoted as $\sqcap S$ such that $\forall x \in S \cdot \sqcap S \sqsubseteq x$ and $\forall l \in S \cdot (\forall x \in S \cdot l \sqsubseteq x) \Rightarrow l \sqsubseteq \sqcap S$. The poset \mathcal{D} has a supremum (or top) T if and only if $T = \sqcup \mathcal{D} \in \mathcal{D}$, and has an infimum (or bottom) \perp iff $\perp = \sqcap \mathcal{D} \in \mathcal{D}$.

Lattice and Complete Partial Order (CPO): A CPO is a poset $\langle \mathcal{D}, \sqsubseteq, \perp, \sqcup \rangle$ with infimum \perp such that any denumerable ascending chain $\{x_i \in \mathcal{D} \mid i \in \mathbb{N}\}$ has a least upper bound $\sqcup_{i \in \mathbb{N}} x_i \in \mathcal{D}$. A lattice is a poset $\langle \mathcal{D}, \sqsubseteq, \sqcup, \sqcap \rangle$ such that every pair of elements x, y has a *lub* $x \sqcup y$ and a *glb* $x \sqcap y$ in \mathcal{D} , thus every finite subset of \mathcal{D} has a *lub* and *glb*. A complete lattice $\langle \mathcal{D}, \sqsubseteq, \perp, \top, \sqcup, \sqcap \rangle$ is lattice with arbitrary subset $S \in \wp(\mathcal{D})$ has a *lub* $\sqcup S$, hence a complete lattice has a supremum $\top = \sqcup \mathcal{D}$ and an infimum $\perp = \sqcap \emptyset$.

Abstraction and Galois connection: In the framework for abstract interpretation, Galois connections are used to formalise the correspondence between concrete properties (like sets of traces) and abstract properties (like sets of reachable states), in case there is always a most precise abstract property over-approximating any concrete property. Given two posets $\langle \mathcal{D}, \sqsubseteq \rangle$ (*concrete domain*) and $\langle \mathcal{D}^\#, \sqsubseteq^\# \rangle$ (*ab-*

stract domain), the pair $\langle \alpha, \gamma \rangle$ of functions $\alpha \in \mathcal{D} \mapsto \mathcal{D}^\#$ (known as *abstraction function*) and $\gamma \in \mathcal{D}^\# \mapsto \mathcal{D}$ (known as *concretisation function*) forms a *Galois connection* iff $\forall x \in \mathcal{D} \cdot \forall y^\# \in \mathcal{D}^\# \cdot \alpha(x) \sqsubseteq^\# y^\# \Leftrightarrow x \sqsubseteq \gamma(y^\#)$ which is mathematically represented as $\langle \mathcal{D}, \sqsubseteq \rangle \xleftrightarrow[\alpha]{\gamma} \langle \mathcal{D}^\#, \sqsubseteq^\# \rangle$ such that (1) α and γ are monotonic; (2) $\gamma \circ \alpha$ is extensive (i.e. $\forall x \in \mathcal{D} \cdot x \sqsubseteq \gamma(\alpha(x))$); (3) $\alpha \circ \gamma$ is reductive (i.e. $\forall y^\# \in \mathcal{D}^\# \cdot y^\# \sqsubseteq \alpha(\gamma(y^\#))$). The rationale underpinning Galois connections is that the concrete properties in \mathcal{D} are approximated by abstract properties in $\mathcal{D}^\#$: $\alpha(x)$ is the most precise sound over-approximation of x in the abstract domain \mathcal{D} and $\gamma(y^\#)$ is the least precise element of \mathcal{D} that can be over-approximated by $y^\#$. The abstraction of a concrete property $x \in \mathcal{D}$ is said to be *exact* whenever $\gamma(\alpha(x)) = x$, in other words, abstraction $\alpha(x)$ of property x loses no information at all. Furthermore, we can say $y^\# \in \mathcal{D}^\#$ is a *sound approximation* of $x \in \mathcal{D}$ iff $x \sqsubseteq \gamma(y^\#)$.

3 Discussion

In this section, we illustrate the nature of *abstraction* in AF and leverage the accrual of arguments whilst preserving the semantic information between them.

Example 3.1. Consider the example provided in [Nielsen and Parsons, 2006], consisting of the following abstract arguments:

- *A1: Joe does not like Jack;*
- *A2: There is a nail in Jack's antique coffee table;*
- *A3: Joe hammered a nail into Jack's antique coffee table;*
- *A4: Joe plays golf, so Joe has full use of his arms;*
- *A5: Joe has no arms, so Joe cannot use a hammer, so Joe did not hammer a nail into Jack's antique coffee table.*

As we can see in Figure 2, that the argument A5 attacks the argument A3, whereas arguments A3 and A4 directly attacks and defeats the argument A5.

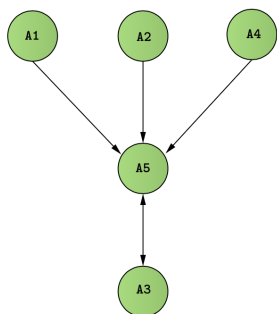


Figure 2: Original AF of Jack and Joe situation



Figure 3: Abstraction of Jack and Joe situation

Given a set of arguments A in the concrete domain AF, a mapping G_m is a Galois surjective mapping $G_m : A \mapsto \hat{A}$ for a set \hat{A} in the abstract domain AF. An abstract domain AF is, formally speaking, an AF that corresponds to a mapping G_m , which identifies arguments that are mapped onto the abstract domain AF. An attack (a, b) in R is mapped like arguments to the abstract domain AF into $(G_m(a), G_m(b))$; in other words, many attacks could map to the same attack. The mapping G_m from concrete domain AF to abstract domain AF can be considered a semantic preserving function where arguments and attacks preserve their semantic relationship/similarity in the abstract domain AF.

In our work, employing the abstraction interpretation technique, the semantic relationship between arguments $A1, A2$ and $A4$ can be strengthened as AX , as shown in Figure 3. Through this simple example, we can reduce the representational complexity of

large AFs which further reduces computational cost. This abstraction in multi-agent dynamic AFs can be extended to many realms of argumentation, where auxiliary information (apart from simply winning or losing the argument) come into consideration. One such consideration involves hiding certain information from an opponent e.g. agents abstracting away sensitive and confidential information. Future work will investigate and formally define the properties of the abstract domain AF, such as conflict-freeness, admissibility and stable extension and further investigate the complexity results of computing the abstract domain AF from the concrete domain AF.

4 Conclusion

In this work, we introduced the notion of reducing the complexity of an abstract argumentation framework in a multi-agent setting using abstraction principles from model checking to reduce representational as well as computational cost, which is usually caused due to increased number of arguments in the framework. Furthermore, due to the abstraction of the AF, it would be possible to develop succinct explanations for humans or other agents in the system.

Acknowledgements

The author thanks Timotheus Kampik for guidance and valuable insights in this project and the anonymous reviewers for their suggestions and feedback. This work was partially funded by the Knut and Alice Wallenberg Foundation.

References

- [Amgoud and Prade, 2009] Amgoud, L. and Prade, H. (2009). Using arguments for making and explaining decisions. *Artificial Intelligence*, 173(3-4):413–436.
- [Bench-Capon and Dunne, 2007] Bench-Capon, T. J. and Dunne, P. E. (2007). *Argumentation*

- in artificial intelligence. *Artificial intelligence*, 171(10-15):619–641.
- [Biere et al., 1999] Biere, A., Cimatti, A., Clarke, E., and Zhu, Y. (1999). Symbolic model checking without bdds. In *International conference on tools and algorithms for the construction and analysis of systems*, pages 193–207. Springer.
- [Booth et al., 2013] Booth, R., Kaci, S., Rienstra, T., and Torre, L. v. d. (2013). A logical theory about dynamics in abstract argumentation. In *International Conference on Scalable Uncertainty Management*, pages 148–161. Springer.
- [Burch et al., 1992] Burch, J. R., Clarke, E. M., McMillan, K. L., Dill, D. L., and Hwang, L.-J. (1992). Symbolic model checking: 1020 states and beyond. *Information and computation*, 98(2):142–170.
- [Campanella, 2019] Campanella, A. (2019). Intent based network operations. In *2019 Optical Fiber Communications Conference and Exhibition (OFC)*, pages 1–3. IEEE.
- [Clarke et al., 2000] Clarke, E., Grumberg, O., and Peled, D. (2000). Model checking cambridge.
- [Cousot and Cousot, 1977] Cousot, P. and Cousot, R. (1977). Abstract interpretation: a unified lattice model for static analysis of programs by construction or approximation of fixpoints. In *Proceedings of the 4th ACM SIGACT-SIGPLAN symposium on Principles of programming languages*, pages 238–252.
- [Dung, 1995] Dung, P. M. (1995). On the acceptability of arguments and its fundamental role in nonmonotonic reasoning, logic programming and n-person games. *Artificial intelligence*, 77(2):321–357.
- [Dung et al., 2009] Dung, P. M., Kowalski, R. A., and Toni, F. (2009). Assumption-based argumentation. In *Argumentation in artificial intelligence*, pages 199–218. Springer.
- [Dunne, 2009] Dunne, P. E. (2009). The computational complexity of ideal semantics. *Artificial Intelligence*, 173(18):1559–1591.
- [Dunne and Wooldridge, 2009] Dunne, P. E. and Wooldridge, M. (2009). Complexity of abstract argumentation. In *Argumentation in artificial intelligence*, pages 85–104. Springer.
- [Falappa et al., 2011] Falappa, M. A., Garcia, A. J., Kern-Isberner, G., and Simari, G. R. (2011). On the evolving relation between belief revision and argumentation. *The Knowledge Engineering Review*, 26(1):35–43.
- [Halpern and Vardi, 1991] Halpern, J. Y. and Vardi, M. Y. (1991). Model checking vs. theorem proving: a manifesto. *Artificial intelligence and mathematical theory of computation*, 212:151–176.
- [Hoek and Wooldridge, 2002] Hoek, W. v. d. and Wooldridge, M. (2002). Model checking knowledge and time. In *International SPIN Workshop on Model Checking of Software*, pages 95–111. Springer.
- [Maudet et al., 2006] Maudet, N., Parsons, S., and Rahwan, I. (2006). Argumentation in multi-agent systems: Context and recent developments. In *International workshop on argumentation in multi-agent systems*, pages 1–16. Springer.
- [Nielsen and Parsons, 2006] Nielsen, S. H. and Parsons, S. (2006). A generalization of dung’s abstract framework for argumentation: Arguing with sets of attacking arguments. In *International Workshop on Argumentation in Multi-Agent Systems*, pages 54–73. Springer.
- [Patil, 2022] Patil, M. S. (2022). Explainability in autonomous pedagogically structured scenarios. In *36th AAAI 2022 Workshop on Explainable Agency in Artificial Intelligence*.
- [Peled and Pnueli, 1994] Peled, D. and Pnueli, A. (1994). Proving partial order properties. *Theoretical Computer Science*, 126(2):143–182.

Distributed Hypothesis Generation and Evaluation

Jordan Robinson

Department of Electrical Engineering and Electronics, University of Liverpool, UK

Abstract

Intelligence analysis is currently conducted by distributed teams of expert human agents who use their domain knowledge, combined with a variety of structured analytical techniques, to generate and evaluate sets of conflicting hypotheses to inform potential high-stake decision making. Analysis can be tedious but it requires the full attention of human agents as the context can be such that what would otherwise be a minor detail has a significant impact on the likelihood of a hypothesis, and so on the downstream decision making. This project aims to enhance the speed and scale of intelligence analyses through the development of decision-support tools which combine explainable artificial intelligence algorithms with human expertise, in the form of human-machine teams, to aid intelligence analysts in evaluating complex and competing hypotheses. The tools created will combine techniques found within the natural language processing and computational argumentation literature and should enable analysts to focus their attention where it's needed most by assisting them throughout the analytical pipeline.

1 Introduction

Threats to countries are forever present on a domestic and an international scale, such as the attacks on the Twin Towers on 11th September 2001 or the current (2022) war in Ukraine, and their realisation can cause catastrophic and irreversible damage, if overlooked. Much of the information available to assess these risks is fragmented, incomplete, conflicting, dy-

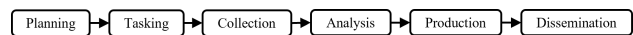


Figure 1: A typical intelligence analysis pipeline.

namic, obscured, and potentially deceptive, which makes monitoring situations increasingly complex.

The stages within a typical intelligence cycle are planning, tasking, collection, analysis, production, and dissemination of a finished intelligence product to decision makers, stakeholders and policymakers (Fig. 1) [U.S., 2000]. All information is funnelled through this process which aims to produce succinct and accurate intelligence. After *planning*, an all-source analyst raises a Request For Information (RFI), in response to which a single-source analyst is *tasked* to *collect*, *analyse* and *produce* a summary report which is *disseminated* back to the all-source analyst. Single-source analysts gather their information using the five disciplines of intelligence collection, which are defined in [Lo Clwenthal, 2015] as: Human Intelligence, Signals Intelligence, Geo-spatial Intelligence, Measurement and Signature Intelligence, and Open Source Intelligence. Once the RFI has been satisfied, the all-source analyst's job is to produce a finished intelligence product using all the provided data and a series of structured analytical techniques [Heuer Jr. and Pherson, 2014].

Intelligence analysts synthesise both contradicting and supporting information from a multitude of sources where the resolution of those differences can be thought of as an argumentative process. Much of the analysis undertaken by both single-source and all-source analysts is performed manually. The wealth of intelligence within summary reports is a great asset,

in that all-source analysts could use it to inform their analysis by verifying claims with evidence. However, incorporating it into analysis is a time-consuming task which can be cognitively demanding and could lead to information overload. In the future, there is a necessity for the intelligence community to employ more data from single-source analysts to assess hypotheses. Artificial Intelligence (AI) offers the potential for computers to enhance a (human) analyst's work-flow by performing some of the argumentative reasoning involved in their duties. More specifically, argument mining – an emerging topic within the field of AI – could help an analyst to assimilate how multiple pieces of intelligence may be both contradicting and supporting each another.

As humans, we communicate using an arrangement of propositions which are amalgamated to form arguments in an attempt to transfer our understanding to one another. An argument consists of a conclusion which is supported by one or more premises [Walton et al., 2008]. Argumentation is inherently *defeasible* which is why a plethora of propositions can be asserted in a debate, if a person believes that another's argument is strong enough to defeat theirs [Eemeren et al., 2014]. Reasoning is *defeasible* when the premises supporting a conclusion are tentatively presumed to be true until additional information becomes available which might invalidate it.

Argumentation schemes capture stereotypical patterns of inference in commonly used arguments, and their critical questions are a way of casting doubt on an argument's support for its conclusion¹. The schemes were combined with structured analytical techniques for intelligence analysis where they improved an analyst's critical thinking [Murukannaiah et al., 2015]. They have also been applied within a legal domain to model hypothetical reasoning, a common type of reasoning within the intelligence domain [Bench-Capon and Prakken, 2010, Grabmair and Ashley, 2010].

Computational Argumentation (CA) is gaining more traction within the AI community due to its effectiveness in modelling logical and defeasible

reasoning, making it a good candidate when modelling debates, decision-making and investigations [Dung et al., 2007, Bench-Capon and Dunne, 2007]. Argumentation mining, first proposed in [Mochales and Moens, 2011], is a new subset of CA which could enable analysts to synthesise more data due to its ability to automatically extract structured argument data from unstructured natural language, which can then be reasoned with using a computational model of argument. Since the field's inception, there have been many advancements and improvements to each part of the argument mining pipeline – refer to [Lippi and Torroni, 2016], [Cabrio and Villata, 2018] and [Lawrence and Reed, 2020] for a comprehensive review on the state of the art.

While argumentation schemes can capture the structure of unstructured natural language corpora, they have not been used to aid analysts in the reasoning about those arguments nor can they fully model the changes to a hypothesis as new information becomes available. Abstract argumentation, first proposed in [Dung, 1995], is an interesting approach to this reasoning problem as it can model the acceptability of arguments, where arguments and attack relations can be respectively modelled as nodes and edges in a directed graph. There have been many accounts which extend Dung's seminal work to broaden the concepts captured in an argumentation framework, however they are not discussed in this letter. To name two examples, in [Robinson, 2021a], they quantify the value of additional arguments within Dung-style frameworks and then test their approach, in [Robinson, 2021b], to assess the net benefit of collecting certain data within an intelligence setting.

The evaluation of argumentation frameworks results in sets of conflict-free and acceptable arguments which fall into two categories, extension- [Baroni and Giacomin, 2009] or labelling-based [Wu and Caminada, 2010] approaches. The extensions or labellings can be thought of as a potential conclusion (or hypothesis) which depend on a graph's topology. Labelling-based semantics also enable the quantification of uncertainty about the inclusion and acceptance status of arguments in frameworks, first proposed in [Riveret et al., 2017].

¹See [Walton et al., 2008] for a comprehensive overview on argumentation schemes and their critical questions.

2 Research Aim & Method

The aim of this research is to investigate, adapt, design, and develop the techniques found within the natural language processing and CA literature to aid human agents during hypothesis generation and evaluation within intelligence settings. The application of CA in an intelligence domain is novel as it has not been well-studied within the current literature.

To date, this research has modelled the argumentative mechanics between an all-source analyst and their single-source counterparts. The approach captured this using Dung-style argument frameworks, labelling-based semantics, and an uncertainty quantification about the inclusion and acceptance status of arguments within frameworks. The model assumed that each single-source agent had been *tasked* and had *collected* the information required to satisfy an all-source agent’s RFI (Fig. 1). The single-source agent then evaluated their evidence locally and disseminated their conclusions with an all-source agent who used those semantic evaluations to arrive at a global conclusion, which was assumed to be synonymous with a finished intelligence product.

To initialise the simulation, a random Dung-style argument system \mathcal{G} was instantiated using a user-defined number of arguments n_A , relations n_R and symmetric attacks n_{sym} . The number of symmetric attacks are important as increasing them increases the potential number of complete labellings within an evaluated argumentation framework, and thus the uncertainty between the sets of acceptable labellings. The framework was evaluated to discover the complete labellings, or the ground truth, of the global graph. The resulting argument labels are used to compute the probability of each argument’s acceptance status $P(arg = label | \mathcal{L}_C)^2$. The global framework \mathcal{G} was then partitioned into several local frameworks, which correspond to the number of single-source analysts n_{ssa} chosen to work on the problem. It should be noted that no arguments or attacks were removed when partitioning the global framework into a set of n_{ssa} local ones – i.e., a union of the local argumentation frameworks results in the undivided,

²where \mathcal{L}_C refers to the set of complete labellings.

global framework. Sampling relations, instead of arguments, was the method chosen to dissect the global frameworks as it was believed that this better represented how single-source analysts obtain data during collection. This ensured that the argument labels for each local framework were diagnostic – i.e., the argument labels for each local framework included the labels “in”, “out”, and “undecided”, and not just “in” if the sampled arguments happened to be unrelated.

Once the local perspectives had been instantiated, each single-source agent evaluated their graph to discover their sets of complete labellings. The labellings were then divulged to the all-source agent³ who counted the number of times an argument was labelled *in*, *out*, or *undecided* in each of the single-source agent’s frameworks. The counts were normalised using the total number of times that an argument featured in each of the single-source agent’s labellings, which enabled the all-source agent to compute the probability of an argument’s acceptance status given each local agent’s evaluation of their graph $P(arg = label | \mathcal{L}_{C,ssa})^4$. The error \mathbb{E} (or difference) between the all-source agent’s prediction and the ground truth was computed for each argument to assess the algorithm’s accuracy.

3 Discussion

A parametric study is currently underway on 10 nodes on Barkla, the University of Liverpool’s High Performance Computing cluster. This experiment aims to assess whether increasing the number of single-source agents n_{ssa} (Fig. 2) or the amount of local knowledge an agent possesses (i.e., increasing the number of local relations n_r sampled) (Fig. 3), increases the all-source agent’s prediction accuracy. Some of the preliminary results are presented for one global framework which had 25 arguments and 30 total relations with no symmetric attacks and one complete labelling only, for 100 test runs.

For the two initial results presented in this letter,

³It is important for the reader to note that there were no dialogue protocols employed in this algorithm.

⁴where $\mathcal{L}_{C,ssa}$ refers to the set of complete labellings shared by the single-source agent with the all-source agent.

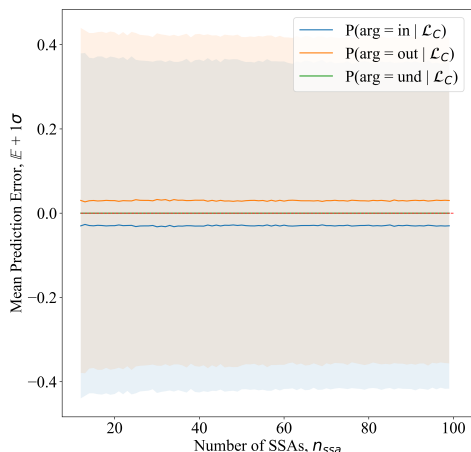


Figure 2: Measures the change in the mean prediction error when increasing the number of single-source agents n_{ssa} from 10 to 100, where $n_r = 5$.

the all-source agent’s prediction error for arguments labelled *in* and *out* was symmetrical, and the probabilities of arguments labelled *undecided* was zero (Fig. 2 and 3). This is most likely due to a lack of symmetric attacks within the framework, but this will be reassessed after the tests conclude.

Increasing the number of agents evaluating the problem did not decrease the uncertainty, or variance, in the all-source agent’s prediction error (Fig. 2). However, increasing the number of local relations available to each single-source agent dramatically reduced the all-source agent’s prediction error (Fig. 3). This is quite intuitive and seems to affirm an old saying, *quality over quantity*.

4 Conclusion

To conclude, the preliminary results discussed have provided this research with an initial understanding of the application of argumentative mechanics within intelligence settings. As the discussed system is an abstraction, the next steps are to apply this technique to a realistic scenario to ensure that the gen-

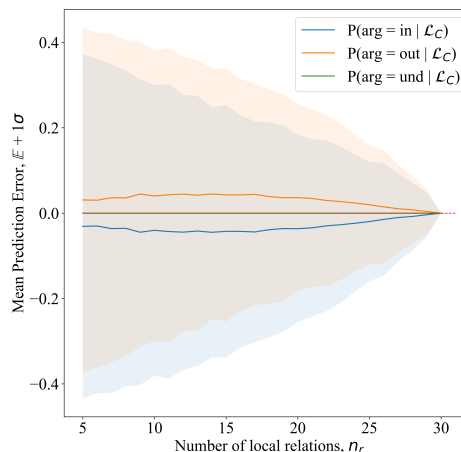


Figure 3: Examines the change in the mean prediction error when increasing the number of sampled local relations n_r from 5 to 30, where $n_{ssa} = 10$.

erated frameworks are a correct mapping of this interaction. Moreover, the incorporation of arguments from real-world scenarios could provide an opportunity to extend the theoretical work on merging Incomplete Argumentation Frameworks, proposed in [Baumeister et al., 2018], by combining single-source agents’ local argumentation frameworks to create a global one for semantic evaluation by the all-source agent. The envisaged system could benefit the intelligence enterprise by automatically extracting and evaluating sets of conflicting arguments found within single-source analysts’ natural language summary reports, and thus aiding an all-source analyst’s analysis of competing hypotheses in a human-machine multi-agent argument mining team.

Acknowledgements

I would like to thank Prof. Katie Atkinson, Prof. Simon Maskell, Prof. Chris Reed, and Todd Robinson for their contributions to this research thus far, and the Defence, Science and Technology Laboratory and the UK Engineering and Physical Sciences Research

Council (Grant number: EP/S023445/1) for financially supporting this project.

References

- [Baroni and Giacomin, 2009] Baroni, P. and Giacomin, M. (2009). *Semantics of Abstract Argument Systems*, pages 25–44. Springer.
- [Baumeister et al., 2018] Baumeister, D., Neugebauer, D., Rothe, J., and Schadrack, H. (2018). Verification in incomplete argumentation frameworks. *Artificial Intelligence*, 264:1–26.
- [Bench-Capon and Dunne, 2007] Bench-Capon, T. and Dunne, P. E. (2007). Argumentation in artificial intelligence. *Artificial Intelligence*, 171(10):619–641.
- [Bench-Capon and Prakken, 2010] Bench-Capon, T. J. M. and Prakken, H. (2010). Using argument schemes for hypothetical reasoning in law. *Artificial Intelligence and Law*, 18:153–174.
- [Cabrio and Villata, 2018] Cabrio, E. and Villata, S. (2018). Five years of argument mining: a data-driven analysis. In *Proc. of IJCAI 2018*.
- [Dung et al., 2007] Dung, P., Mancarella, P., and Toni, F. (2007). Computing ideal sceptical argumentation. *Artificial Intelligence*, 171(10):642–674.
- [Dung, 1995] Dung, P. M. (1995). On the acceptability of arguments and its fundamental role in nonmonotonic reasoning, logic programming and n-person games. *Artificial Intelligence*, 77(2):321–357.
- [Eemeren et al., 2014] Eemeren, F. H. v., Garssen, B., Verheij, B., Krabbe, E. C. W., Snoeck Henkemans, A. F., and Wagemans, J. H. M. (2014). *Handbook of Argumentation Theory*. Springer.
- [Grabmair and Ashley, 2010] Grabmair, M. and Ashley, K. (2010). Argumentation with value judgments - an example of hypothetical reasoning. volume 223, pages 67–76.
- [Heuer Jr. and Pherson, 2014] Heuer Jr., R. and Pherson, R. (2014). *Structured Analytic Techniques for Intelligence Analysis*. CQ Press.
- [Lawrence and Reed, 2020] Lawrence, J. and Reed, C. (2020). Argument mining: A survey. *Comput. Linguist.*, 45(4):765–818.
- [Lippi and Torroni, 2016] Lippi, M. and Torroni, P. (2016). Argumentation mining: State of the art and emerging trends. *ACM Trans. Internet Technol.*, 16(2).
- [Lo Clwenthal, 2015] Lo Clwenthal, M.M. andark, R. (2015). *The Five Disciplines of Intelligence Collection*. SAGE Publications.
- [Mochales and Moens, 2011] Mochales, R. and Moens, M.-F. (2011). Argumentation mining. *Artificial Intelligence and Law*, 19:1–22.
- [Murukannaiah et al., 2015] Murukannaiah, P. K., Kalia, A. K., Telangy, P. R., and Singh, M. P. (2015). Resolving goal conflicts via argumentation-based analysis of competing hypotheses. *IEEE 23rd International Requirements Engineering Conference*.
- [Riveret et al., 2017] Riveret, R., Baroni, P., Gao, Y., Gao, Y., Gao, Y., Governatori, G., Rotolo, A., and Sartor, G. (2017). A labelling framework for probabilistic argumentation. *Artificial Intelligence*.
- [Robinson, 2021a] Robinson, T. (2021a). Value of information for argumentation based intelligence analysis. *arXiv: Artificial Intelligence*.
- [Robinson, 2021b] Robinson, T. (2021b). Value of information for argumentation based intelligence analysis. *arXiv: Artificial Intelligence*.
- [U.S., 2000] U.S., I. O. S. S. S. (2000). Intelligence threat handbook.
- [Walton et al., 2008] Walton, D., Reed, C., and Macagno, F. (2008). *Argumentation Schemes*. Cambridge University Press.
- [Wu and Caminada, 2010] Wu, Y. and Caminada, M. (2010). A labelling based justification status of arguments. *Studies in Logic*, 3:12–29.

Exploring Internal Structures of an Argumentation System and Improving Reasoning Efficiency with Backward Searching Framework

Hao Wu

University of Aberdeen, Scotland

Abstract

The first step to determine an argument's acceptability in an argumentation system normally takes place by computing extensions based on the work of (Dung, 1995). Doing so — for many semantics — is computationally expensive, particularly when dealing with large or rapidly updated knowledge bases. My PhD work focuses on identifying the relationship between the dynamic process of constructing an argumentation system from a given knowledge base and determining the acceptability of arguments. Characterising this relationship could aid in investigating the internal structure of an argumentation system, highlighting critical arguments related to the acceptability of a given argument, and providing heuristic information for developing efficient query-answering algorithms by avoiding unnecessary computation.

searching along attack relations backwards would potentially avoid the need to consider any arguments that are unrelated to that argument whose status is being determined. Similarly, in a structured context (Wu et al., 2022), starting with the conclusion of an argument and building the argument backwards would provide us with a chance to stop argument construction when sufficient information exists (e.g., due to preferences between elements of the knowledge base), potentially avoiding the need to build complete arguments. The combination of these two strategies would improve computational efficiency when dealing with large argumentation systems built upon a real-world knowledge base. To perform this backwards search, we must be able to determine how much we learn about the acceptability of a specific argument via a step in the backwards search.

1 Introduction

In his seminal paper, (Dung, 1995) defines various semantics to characterise different notions of acceptability in an argumentation system, with an argument's acceptability defined by whether it is present in a set of arguments, or *extension*, computed according to the semantics. When many arguments are involved, the underlying knowledge base is constantly updated, or some semantics are involved, computing such extensions can be computationally expensive. My work builds on the notion of backwards search; intuitively, starting with an argument and

2 Background

We begin by describing standard notions from argumentation theory. An abstract argumentation framework F is a pair (\mathbb{A}, \mathbb{R}) , where \mathbb{A} is a set of arguments and $\mathbb{R} \subseteq \mathbb{A} \times \mathbb{A}$ is a set of attack relations. For two arguments $A, B \in \mathbb{A}$, if B attacks A , we write $(B, A) \in \mathbb{R}$. A set of arguments $S \subseteq \mathbb{A}$ is said to be conflict-free if $\forall A, B \in S, (A, B) \notin \mathbb{R}$. A conflict-free set of arguments S is said to be admissible if for each $A \in S$ we know that $\forall B \in \mathbb{A} \wedge (B, A) \in \mathbb{R}$ there exists $C \in S$ such that $(C, B) \in \mathbb{R}$. An argument $A \in \mathbb{A}$ is said to be defended by a set of arguments $S \subseteq \mathbb{A}$ if for all $B \in \mathbb{A}$ such that $(B, A) \in \mathbb{R}$ there is an argument $C \in S$ such

that $(C, B) \in \mathbb{R}$. We say an admissible set S is a complete extension if for all $A \in S$ and A is defended by S ; we say a complete extension S is a preferred extension if S is maximal w.r.t. set inclusion (Dung, 1995). We note in passing that other semantics have been described in the literature (e.g., stable (Dung, 1995) and ideal (Cerutti et al., 2017)).

3 Research Method and Goals

In our research, we propose to use a tree structure to describe the backward search process and its outputs. At the abstract level, we define a set of basic structures and operators to compute argument acceptability by composing different types of sub-structures of trees. Adopting these tree structures will assist us in identifying and focusing on critical structures in an argumentation system. These tree structures also connect the dynamic searching process to the acceptability of a given argument. Ultimately, our work aims to improve the efficiency of (argument-based) reasoning over a given knowledge base and offer a set of higher-level definitions for characterising an argumentation framework.

4 Result and Application

We define different roles of paths and arguments in such a tree structure in terms of their effects on the acceptability of the root argument. By composing these basic definitions with properly defined operations, it is possible to determine an argument's acceptability and simultaneously avoid traversing arguments that do not affect its status. At first glance, some definitions in this framework seem to resemble the concepts of dispute trees (Dung et al., 2007), dialogue trees (Amgoud et al., 2000), and proof dialogues (Modgil and Caminada, 2009). However, this framework is a formalisation of argumentation theory which provides a metalanguage aiming to describe and explore the properties of a dynamic argumentation framework. Primitive structures and their compositions in this framework help us identify and classify critical computational factors and explain the role of an argument in an argumentation space. Due to the page limitation, only two definitions and related examples are introduced next.

Definition 1 (*o-primitive*). An *o-primitive* is denoted as o' , and is defined as either :

- A path (B, A) where $arg(B)$ attacks $arg(A)$. B is known as an *o-node*, A is known as an *e-node*, arg is a surjective function which maps a node to an argument; or
- Given an *o-primitive* o'_1 , the result of
 1. attaching an o' to an *o-node* of o'_1 .
 2. attaching an e' to an *e-node* of o'_1 .
 3. attaching an o' to an *e-node* of o'_1 .

Definition 2 (*e-primitive*). An *e-primitive* is denoted as e' and is defined as either:

- A path (C, B, A) where $arg(c)$ attacks $arg(C)$, $arg(B)$ attacks $arg(A)$ ¹. C is an *e-node*, B is an *o-node* and A is an *e-node*.
- Given an *e-primitive* e'_1 , the result of
 1. attaching an o' to an *e-node* of e'_1 .
 2. attaching an e' to an *o-node* of e'_1 .
 3. attaching an o' to an *o-node* of e'_1 .

The generalization of the second bullet point of above definitions are (1) $o'::t$, (2) $e'::\bar{t}$ (3) $o'::\bar{t}$ w.r.t. the given t -primitive, where $t \in \{o', e'\}$, $\bar{t} = \{o', e'\} \setminus t$, $::t$ means attaching to a t -node of the given t -primitive. The corresponding semantics are: (1) to introduce a conjunctively effective component, (2) to introduce a disjunctively effective component, (3) to introduce an ineffective component. Intuitively, an *e-primitive* indicates all its effective *e-nodes* are admissible, an *o-primitive* indicates all its effective *o-nodes* are admissible, effective nodes are nodes on effective components, effective components are sub structures that contribute to the admissibility of effective nodes and block the influence of ineffective components in the meantime.

Example 1. Starting from (step 0) an *o-primitive*, (step 1) we can introduce two conjunctively effective components: focus on the *o-primitive* introduced in step 0, according

¹A single node A is also an *e-primitive*, we use an approximate definition here for simplicity.

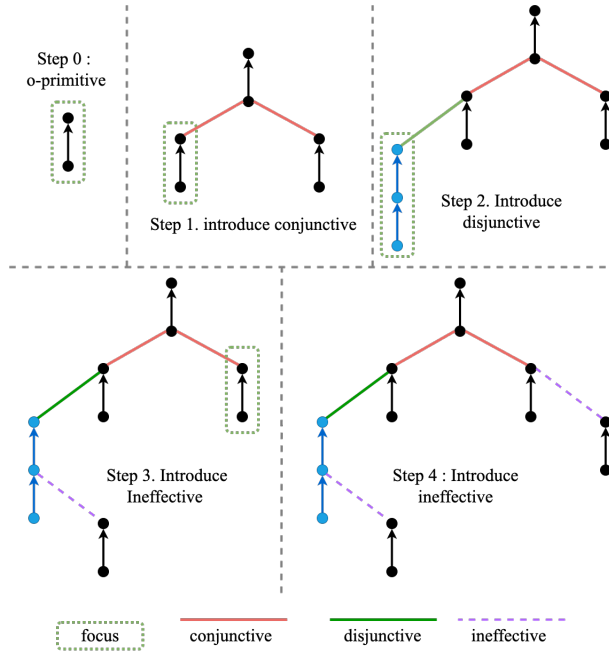


Figure 1: The composition of primitive components

to (1), we attach two o' to an o -node. (step 2) We can introduce disjunctively effective components: focus on one o -primitive introduced in step 1, according to (2), we attach e' to an e -node. (step 3) We can introduce an ineffective component: focus on the e -primitive introduced in step 2, according to (3), we attach o' to an e -node. (step 4) We can introduce another ineffective component: focus on the other o -primitive introduced in step 1, according to (3), we attach o' to an e -node.

Each tree as shown in Figure 1 captures the core structures of attack relation of a family of argumentation frameworks, An instance of the tree at step 4 is shown in Figure 2. There are two levels of effective components, 0 and 1. Each and all elements in level 0, each element in the Cartesian product of level 0 and level 1, and the union of level 0 and level 1 correspond to all complete extensions. Arguments of consecutive levels play different roles but can also belong to the same complete extension in Dung's semantics.

Example 1 demonstrates that an argumentation framework can be represented as a composition of primitive

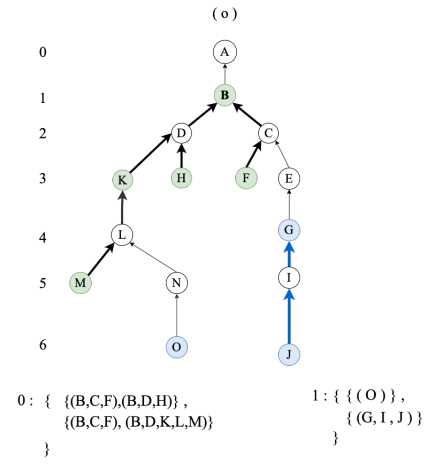


Figure 2: An instance of step 4

structures, and there are corresponding relations between different semantics and compositions of these primitive structures². An argumentation framework might need to be represented by multiple such tree structures. We can work with range-based semantics by considering the overall effects of multiple trees such as the stage(Lagniez et al., 2020) or preferred semantics. Moreover, by monitoring these local structures, their effects and the range of the effects can be computed instantly. Many important concepts in this framework are not covered in this paper such as primitives that introduce important structures corresponding to the existence of different preferred extensions, and trees that share the common sub-tree, etc.

The works introduced in the previous section allow us to handle dynamic argumentation frameworks at the abstract level. Additionally, in our previous work (Wu et al., 2022), we present a framework that constructs an argument backwards from its conclusion and eliminates unnecessary computation when there is enough information to determine related properties.

Example 2. Given a defeasible theory $\mathcal{D} = \{b \rightarrow a, \{c, d\} \rightarrow b, e \Rightarrow c, g \Rightarrow e, \emptyset \Rightarrow g, f \Rightarrow d, \emptyset \Rightarrow f\}$ ³.

²There are another two types of primitives, a -primitive and n -primitive. They represent different kinds of circular attack relations. Due to space limitations, they are not introduced here.

³We use \rightarrow for a strict rule, \Rightarrow for a defeasible rule, and strict (respectively defeasible) rules with empty bodies are axioms (respectively

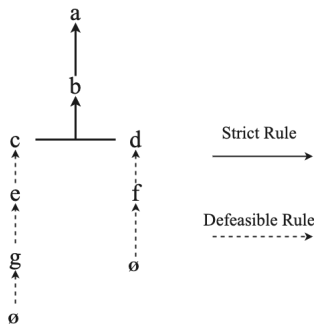


Figure 3: Building a partial argument from its conclusion

In our work (Wu et al., 2022), the building process of an argument concludes ‘a’ could be stopped at different phases. In the scenario of computing preference orderings, all rules in \mathcal{D} are required when adopting the weakest-link, and rules $\{b \rightarrow a, \{c, d\} \rightarrow b, e \Rightarrow c, f \Rightarrow d\}$ are enough when adopting last-link as defined in (Modgil and Prakken, 2014). So the building process of an argument will avoid involving computations that are not necessary.

5 Conclusion and Future Works

Based on the idea of searching along attack relation backwards and adopting a tree structure, we determine the acceptability of an argument in a given argumentation framework. More specifically, compositions of primitive structures provide us with a new viewpoint of an argumentation framework and Dung’s extensions. There is still a significant amount of work that must be done to explore the potential of this framework, guarantee the correctness, and bridge the gap between abstract-level representation and structured-level instantiation (Toni, 2013; Wu et al., 2022). Next, we intend to develop algorithms based on our work to provide an effective solution for tasks such as the ICCMA structured track 2022 (Lagniez et al., 2020). Furthermore, the expressive power provided by this framework could also aid in customising benchmarks of various properties. As a result, researchers in this field will have more solid evaluation criteria.

ordinary premises). When the body of the rule has only one element, such as $\{a\} \rightarrow b$, it will be denoted as $a \rightarrow b$ for simplicity.

Acknowledgements

This research would not be possible without the encouragement and great support from my supervisors, Professor Nir Oren and Dr Bruno Yun.

References

- Amgoud, L., Maudet, N., and Parsons, S. (2000). Modelling dialogues using argumentation. In *Proceedings Fourth International Conference on MultiAgent Systems*, pages 31–38. IEEE.
- Cerutti, F., Gaggl, S. A., Thimm, M., and Wallner, J. (2017). Foundations of implementations for formal argumentation. *IfCoLog Journal of Logics and their Applications*, 4(8):2623–2705.
- Dung, P. M. (1995). On the acceptability of arguments and its fundamental role in nonmonotonic reasoning, logic programming and n-person games. *Artificial Intelligence*, 77(2):321–357.
- Dung, P. M., Mancarella, P., and Toni, F. (2007). Computing ideal sceptical argumentation. *Artificial Intelligence*, 171(10-15):642–674.
- Lagniez, J.-M., Lonca, E., Mailly, J.-G., and Rossit, J. (2020). The fourth international competition on computational models of argumentation.
- Modgil, S. and Caminada, M. (2009). Proof theories and algorithms for abstract argumentation frameworks. In *Argumentation in artificial intelligence*, pages 105–129. Springer.
- Modgil, S. and Prakken, H. (2014). The aspic+ framework for structured argumentation: a tutorial. *Argument & Computation*, 5(1):31–62.
- Toni, F. (2013). A generalised framework for dispute derivations in assumption-based argumentation. *Artificial Intelligence*, 195:1–43.
- Wu, H., Yun, B., and Oren, N. (2022). Improving reasoning efficiency in aspic+ with backwards chaining and partial arguments. In *The Fourth International Workshop on Systems and Algorithms for Formal Argumentation 2022*, pages 86–94.

Discussing Soundness and Completeness for Dialogues that Account for Enthymemes

Andreas Xydis

Department of Informatics, King's College London, UK

Abstract

In previous volumes of OHAAI I discussed how certain locutions can be used to handle enthymemes (i.e., arguments with incomplete logical structure) in dialogical settings, and how a dialogue framework (i.e., a graph which represents the moves made in the dialogue as nodes, and the various relationships amongst them as edges) can be exploited to evaluate enthymemes exchanged during the dialogue. In this paper, I continue from these points by arguing that, under certain conditions, the status of moves made during a dialogue conforming to my dialogue system, corresponds with the status of arguments in the Dung argument framework instantiated by the contents of the moves made at that stage in the dialogue.

missing claim of the intended argument A from which E was constructed including or excluding the claim of A , and misunderstandings that may occur between them due to the enthymemes' lack of structure (see Table 1). Moreover, in [Xydis et al., 2020] the authors provided a dialogue system that integrates these locutions, whereas in [Xydis, 2021] I proposed how to instantiate a *dialogue framework* which can be used to describe the relationships between the moves of the dialogue and, later on, exploit it to assess the acceptability status of enthymemes moved during a dialogue conforming to our system. Notice that although there are some works that focus on dialogue systems which accommodate enthymemes, such as [Black and Hunter, 2007], [Dupin de Saint-Cyr, 2011], [Hosseini et al., 2017] and [Prakken, 2005], none of these deal with misunderstandings that may occur during the dialogue because of the use of enthymemes, nor with the evaluation of enthymemes exchanged during the dialogue (notable exception is [Prakken, 2005] which only addresses backward extension).

The goal of my PhD is to develop a sound and complete dialogue system, which accommodates arguments, backward and forward extension of enthymemes, and misunderstandings that may occur between the participants of a dialogue, with regards to the acceptability of these arguments and enthymemes (moved in the dialogue) in the argument framework instantiated by the contents of the moves made in the dialogue, under the complete, preferred, grounded and stable semantics. After showing soundness and completeness for a dialogue system that deals with arguments and forward extension of enthymemes under the complete semantics in [Xydis et al., 2021], and proposing a dialogue system in

1 Introduction

Formal systems modelling arguments and dialogues usually assume that arguments are complete, i.e. that they are fully logically structured. Practically, however, humans tend to omit parts of the arguments they intend to get across to their discussants. Hence, if we are to provide normative support for human-human debate and enable AIs and humans to jointly reason, we need to investigate how to process enthymemes during dialogues.

In [Xydis, 2020] I provided locutions (witnessed in real world dialogues) which allow the participants of a dialogue to handle *backward extension* of enthymemes, i.e. how to question and provide a justification for a premise ϕ of an enthymeme, *forward extension* of enthymemes, i.e. how to request and provide the missing components between the internal elements of an enthymeme E and the

[Xydis et al., 2022a] which handles arguments and misunderstandings rising from the use of enthymemes during the dialogue and is sound and complete under the complete, preferred, grounded and stable semantics (where the technical details are given in [Xydis et al., 2022b]), I managed to achieve my PhD’s objective under the assumption that the participants of the dialogue are *sensible* and *honest*, and that the dialogue is *exhaustive*. In Section 2 I describe the techniques used for my PhD, and in Section 3 I present the conditions of sensibleness, honesty and exhaustiveness, as well as the theorem proved in my PhD. In Section 4 I provide a conclusion for my paper, and potential future work.

Locution	Meaning
assert	Assert an enthymeme.
why	Question a particular element of a previous enthymeme, which is a request for the other participant to provide a backward extension on that element.
because	Provide a backward extension on a questioned element.
and-so	Request a forward extension of a previous enthymeme.
hence	Provide a forward extension of a previous enthymeme.
w.d.y.t.i.m.b.	Check the other participant’s understanding of an enthymeme by asking “ <i>what did you think I meant by ...</i> ”.
assumed	Provide their own interpretation of an enthymeme.
meant	Correct the other participant’s interpretation of an enthymeme.
agree	Confirm the other participant’s interpretation of an enthymeme.

Table 1: Table of locutions.

2 Method

For my research, I decided to use the $ASPIC^+$ framework [Modgil and Prakken, 2013] to formalise arguments as it allows for evaluating them under Dung’s classic semantics [Dung, 1995], accessing their internal structure,

exploring the nature of attacks between them and accommodating other existing argumentation formalisms, e.g. deductive argumentation [Besnard and Hunter, 2008] and ABA [Bondarenko et al., 1997]. Additionally, it allows one to define their own way of constructing arguments into a given logic, and so it enables the use of argument-trees [Hosseini et al., 2014] which preserve the principles of $ASPIC^+$ and based on which the structure of enthymemes is defined (see [Xydis et al., 2020] for details on how we structure enthymemes).

To develop a dialogue framework DF, which represents the various relations between the moves made during a dialogue, I explored how a bipolar argument framework is constructed [Cayrol and Lagasque-Schiex, 2005] as well as how extended argument frameworks behave [Modgil and Prakken, 2010]. The reason was to capture the attacking and the supporting relationship that two enthymemes can have in the same way as arguments do in a bipolar argument framework. In other words, the backward expansion E' of an enthymeme E and the forward expansion E'' of E in a dialogue are perceived as E' supporting E and E supporting E'' , respectively. Additionally, I wanted to explore how attacks on attacks work and how this influences the acceptability of arguments in the framework. Finally, I defined a labelling function (inspired by [Caminada, 2006]) on DF to determine the dialogical status of moves made in a dialogue. In this way I was able to compare the acceptability of arguments and enthymemes moved in the dialogue to their acceptability in the argument framework instantiated by the contents of the moves made in the dialogue.

3 Discussion

When an agent moves an enthymeme, it has in mind a complete argument and so I call this the *intended argument* of the move (denoted as $\text{IntArg}_{s(m_i)}(m_i)$, where $s(m_i)$ is the sender of the move m_i). Each move has a content which can be an enthymeme, an argument or \emptyset . However, it is possible that we deal with nefarious agents so I do not insist that an agent intends a complete argument, nor that the intended argument does in fact extend the enthymeme moved (meaning that the intended argument of a move can be \emptyset , and the content of a move might not be part of the intended argument

of the move). In order to prove that the dialectical status of the moves in the DF (determined by a σ labelling) is sound and complete in relation to the dialectical status of their intended arguments in the argument framework AF instantiated by the contents of the moves made in the dialogue (determined by a σ labelling), where $\sigma \in \{complete, grounded, preferred, stable\}$, the participants of the dialogue need to be sensible and honest, and the dialogue needs to be exhaustive. Below I explain these conditions:

- **Sensibleness:** Agents are sensible only when they challenge their counterpart's enthymemes and extend their own enthymemes;
- **Honesty:** Agents are honest only if when they assert an enthymeme, this is constructed from the intended argument of their move and such that it does defeat the intended targeted argument (i.e., what the sender assumes to be the target move's intended argument) according to the *ASPIC*⁺ definition of defeat. It also means that whenever they reveal their understanding of the argument from which their counterpart's enthymeme E is constructed, the argument does indeed extend E , and whenever they reveal the intended argument of their own move, this is indeed the argument they intended;
- **Exhaustiveness:** A dialogue is exhaustive only when each agent has made all their available moves. In other words, if a participant can move an argument (constructed by the contents of the moves made in the dialogue) as a defeat against the content of another move made in the dialogue, or question a defeat relation or reply to a move, they indeed do so.

Based on the above I proved the following theorem:

Theorem 1. *Let $d = [m_0, \dots, m_l]$ be an exhaustive dialogue between participants $Part = \{Prop, Op\}$ who are sensible and honest with respect to d . Let the dialogue framework of d and the argument framework instantiated by AT_d be $DF = \langle M \cup \{YES, NO\}, T, Rep, Sup, Exp \rangle$ and $AF = \langle A_{AT_d}, Dfs \rangle$, respectively. It follows that:*

- (i) *For every σ labelling function L_{AF} on AF, there exists a σ labelling function L_{DF} on DF such that for*

every $m_i \in M$ such that $IntArg_{\mathcal{S}(m_i)}(m_i) \neq \emptyset$ we have that:

$$L_{DF}(m_i) = L_{AF}(IntArg_{\mathcal{S}(m_i)}(m_i));$$

- (ii) *For every σ labelling function L_{DF} on DF, there exists a σ labelling function L_{AF} on AF such that for every $A \in A_{AT_d}$, if there is some $m_i \in M$ such that $A = IntArg_{\mathcal{S}(m_i)}(m_i)$ then:*

$$L_{AF}(A) = L_{DF}(m_i);$$

where:

- M is a subset of the set of moves made in d , and YES and NO are two auxiliary elements used for confirming and rejecting the other participant's interpretation of an enthymeme, respectively;
- T is a binary defeat relationship that is determined by the target relationship between moves;
- Rep is a binary reply relationship that is determined by the reply relationship between moves;
- Sup is a binary support relationship that is also determined by the reply relationship between moves (such that m_i supports m_j if and only if the contents of m_i and m_j are enthymemes and there is some m_k such that m_i replies to m_k and m_k replies to m_j);
- Exp is a new binary expansion relationship that is determined by the target and reply relationships between moves (such that m_i expands m_j if and only if the content of m_j is an enthymeme of the content of m_i);
- AT_d is the argumentation theory instantiated by the elements revealed in d ;
- A_{AT_d} is the set of arguments instantiated in d ;
- Dfs is a set of defeat relationships between arguments in A_{AT_d} ;
- $\sigma \in \{complete, grounded, preferred, stable\}$.

Due to lack of space I cannot describe here all the technical details needed to support our claim (e.g., the protocol of the dialogue, the characteristics of a move conforming to our dialogue system, the formalisation of the conditions given in this Section, the proof of our theorem etc.), but for a better understanding the reader can refer to [Xydis et al., 2020, Xydis et al., 2021, Xydis et al., 2022a, Xydis et al., 2022b].

4 Conclusion

In this paper I presented my main PhD goal, my work towards achieving it, and the final result of my research, together with the conditions for reaching it.

Specifically, during my PhD I managed to develop a sound and complete dialogue system which captures both arguments and enthymemes, formalised within the *ASPIC*⁺ framework, and accounts for the backward and forward extension of enthymemes, and the misunderstandings that may occur between the participants of a dialogue due to the use of enthymemes. To prove this soundness and completeness the dialogue needs to be exhaustive, and its participants need to be sensible and honest. This is important since it ensures that the use of enthymemes does not prevent the agents from reaching the appropriate conclusion according to the information they have shared.

To the best of my knowledge, dialogues conforming to our system are the first to deal with backward and forward extension of enthymemes as well as the misunderstandings that may arise between the participants of such a dialogue, while at the same time being shown to be sound and complete. Essentially, with my work I try to close the gap between formal logic-based models of dialogue and the kinds of dialogue studied by the informal logic community, which focus on more human oriented models of dialogue, since enthymemes are a common feature of real-world dialogues.

Future work includes investigating how enthymemes can be used for strategic purposes during dialogues as well as implementing our dialogue system in the form of a dialogue manager that supports and guides humans interlocutors as to how they should rationally engage in jointly inquiry dialogues or debates.

Acknowledgements

The research described in this paper would not be possible without the help of E. Black, C. Hampson and S. Modgil.

References

- [Besnard and Hunter, 2008] Besnard, P. and Hunter, A. (2008). *Elements of argumentation*, volume 47. MIT press Cambridge.
- [Black and Hunter, 2007] Black, E. and Hunter, A. (2007). A generative inquiry dialogue system. In *Proceedings of the 6th International Joint Conference on Autonomous Agents and Multiagent systems*, pages 1–8.
- [Bondarenko et al., 1997] Bondarenko, A., Dung, P. M., Kowalski, R. A., and Toni, F. (1997). An abstract, argumentation-theoretic approach to default reasoning. *Artificial intelligence*, 93(1-2):63–101.
- [Caminada, 2006] Caminada, M. (2006). On the issue of reinstatement in argumentation. In *European Workshop on Logics in Artificial Intelligence*, pages 111–123.
- [Cayrol and Lagasquie-Schiex, 2005] Cayrol, C. and Lagasquie-Schiex, M. (2005). On the acceptability of arguments in bipolar argumentation frameworks. In *European Conference on Symbolic and Quantitative Approaches to Reasoning and Uncertainty*, pages 378–389.
- [Dung, 1995] Dung, P. M. (1995). On the acceptability of arguments and its fundamental role in nonmonotonic reasoning, logic programming and n-person games. *Artificial intelligence*, 77(2):321–357.
- [Dupin de Saint-Cyr, 2011] Dupin de Saint-Cyr, F. (2011). Handling enthymemes in time-limited persuasion dialogs. *Proceeding of International Conference on Scalable Uncertainty Management*, pages 149–162.
- [Hosseini et al., 2014] Hosseini, S. A., Modgil, S., and Rodrigues, O. (2014). Enthymeme construction in dialogues using shared knowledge. In *Proceedings of Computational Models of Argument*, pages 325–332.

- [Hosseini et al., 2017] Hosseini, S. A., Modgil, S., and Rodrigues, O. (2017). Dialogues incorporating enthymemes and modelling of other agents' beliefs.
- [Modgil and Prakken, 2010] Modgil, S. and Prakken, H. (2010). Reasoning about preferences in structured extended argumentation frameworks. In *Proceedings of Computational Models of Argument*, pages 347–358.
- [Modgil and Prakken, 2013] Modgil, S. and Prakken, H. (2013). A general account of argumentation with preferences. *Artificial Intelligence*, 195:361–397.
- [Prakken, 2005] Prakken, H. (2005). Coherence and flexibility in dialogue games for argumentation. *Journal of logic and computation*, 15(6):1009–1040.
- [Xydis, 2020] Xydis, A. (2020). Speech acts and enthymemes in argumentation-based dialogues. *Online Handbook of Argumentation for AI (OHAAI)*, 1:53–57.
- [Xydis, 2021] Xydis, A. (2021). A dialogue framework for enthymemes. *Online Handbook of Argumentation for AI*, pages 62–66.
- [Xydis et al., 2020] Xydis, A., Hampson, C., Modgil, S., and Black, E. (2020). Enthymemes in dialogues. *Computational Models of Argument: Proceedings of COMMA 2020*, 326:395–402.
- [Xydis et al., 2021] Xydis, A., Hampson, C., Modgil, S., and Black, E. (2021). Towards a sound and complete dialogue system for handling enthymemes. In *International Conference on Logic and Argumentation*, pages 437–456.
- [Xydis et al., 2022a] Xydis, A., Hampson, C., Modgil, S., and Black, E. (2022a). A sound and complete dialogue system for handling misunderstandings.
- [Xydis et al., 2022b] Xydis, A., Hampson, C., Modgil, S., and Black, E. (2022b). Technical report: A sound and complete dialogue system for handling misunderstandings.