

Citation for published version:

Shams, Z, De Vos, M, Oren, N & Padget, J 2016, Normative practical reasoning via argumentation and dialogue. in S Kambhampati (ed.), *Proceedings of the 25th International Joint Conference on Artificial Intelligence (IJCAI), 2016*. Proceedings (IJCAI International Joint Conference on Artificial Intelligence), AAAI Press, Palo Alto, U. S. A., pp. 1244-1250, 25th International Joint Conference on Artificial Intelligence (IJCAI), 2016, New York, USA United States, 9/07/16. <https://doi.org/10.5555/3060621.3060794>

DOI:

[10.5555/3060621.3060794](https://doi.org/10.5555/3060621.3060794)

Publication date:

2016

Document Version

Peer reviewed version

[Link to publication](#)

(C) ACM 2018. This is the author's version of the work. It is posted here by permission of ACM for your personal use. Not for redistribution. The definitive version was published in IJCAI '16: Proceedings of the Twenty-Fifth International Joint Conference on Artificial Intelligence. Available via: <https://dl.acm.org/citation.cfm?id=3060794>

University of Bath

Alternative formats

If you require this document in an alternative format, please contact:
openaccess@bath.ac.uk

General rights

Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

Take down policy

If you believe that this document breaches copyright please contact us providing details, and we will remove access to the work immediately and investigate your claim.

Normative Practical Reasoning via Argumentation and Dialogue

Zohreh Shams^a, Marina De Vos^a, Nir Oren^b and Julian Padget^a

^a Department of Computer Science, University of Bath, UK

{z.shams, m.d.vos, j.a.padget}@bath.ac.uk

^b Department of Computing Science, University of Aberdeen, UK

n.oren@abdn.ac.uk

Abstract

In a normative environment an agent’s actions are not only directed by its goals but also by norms. Here, potential conflicts among the agent’s goals and norms makes decision-making challenging. We therefore seek to answer the following questions: (i) how should an agent act in a normative environment? and (ii) how can the agent explain why it acted in a certain way? We propose a solution in which a normative planning problem serves as the basis for a practical reasoning approach based on argumentation. The properties of the best plan(s) w.r.t. goal achievement and norm compliance are mapped to arguments that are used to explain why a plan is justified, using a dialogue game.

1 Introduction

Agents in normative systems must be able to reason about actions in pursuit of their goals, but must also consider the *regulative* norms imposed on them. Such norms define obligations and prohibitions on their behavior, and to avoid punishment, agents must comply with norms while pursuing their goals. However, if norm compliance hinders a more important goal or norm, an agent should consider violating it. To decide how to act, an agent thus needs to generate all plans and weigh up the importance of goal achievement and norm compliance against the cost of goal failure and norm violation in different plans. Although some reasoning frameworks do this [Broersen *et al.*, 2001; Kollingbaum and Norman, 2003], little attention has been paid to explaining the agents’ decision making in such frameworks. Such explanation is important in contexts including human-agent teams and agent debugging, and to provide explanation, we propose utilising formal argumentation.

Argumentation has been applied to inconsistency handling and decision-making [Dung, 1995; Amgoud and Prade, 2009], and its dialogical interpretation makes it an appropriate tool to generate explanations for decisions [Fan and Toni, 2015; Caminada *et al.*, 2014b]. Although argumentation has been extensively used in practical reasoning (e.g., [Atkinson and Bench-Capon, 2007]), integrating the reasoning and dialogical aspect of argumentation for decision-making and its explanation has not been addressed by existing approaches.

In this paper we propose an argumentation-based approach to normative practical reasoning using a dialogue game to provide an intuitive overview of agent’s reasoning. In achieving this aim, the following contributions are made: (i) we formalise a set of argument schemes and critical questions [Walton, 1996] aimed at checking plan *justifiability* with respect to goal satisfaction and norm compliance/violation; (ii) we offer a novel decision criterion that identifies the best plan(s) both in the presence and absence of preferences over goals and norms; and (iii) we investigate the properties of the best plan(s). These properties, together with Caminada’s Socratic dialogu game [Caminada *et al.*, 2014a], are used to generate an explanation for the justifiability of the best plan(s).

2 Model

This section introduces a model for normative practical reasoning based on STRIPS planning [Fikes and Nilsson, 1971].

Definition 1 (Normative Planning Problem). A normative planning problem is a tuple $P = \langle FL, \Delta, A, G, N \rangle$ where FL is a set of fluents; $\Delta \subseteq FL$ is the initial state; A is a finite, non-empty set of durative actions; G is the set of agent goals; and N is a set of action-based norms imposed on the agent.

Fluents FL is a set of domain fluents. A literal l is a fluent or its negation. For a set of literals L , we define $L^+ = \{fl \text{ s.t. } fl \in L\}$ and $L^- = \{fl \text{ s.t. } \neg fl \in L\}$. L is *well-defined* if $L^+ \cap L^- = \emptyset$. For a state $s \subseteq FL$, s^+ are fluents considered *true*, and $s^- = FL \setminus s^+$. A state s satisfies literal fl , denoted as $s \models fl$, if $fl \in s$, and satisfies literal $\neg fl$, denoted $s \models \neg fl$, if $fl \notin s$.

Actions An action $a = \langle pr, ps, d \rangle$ is composed of well-defined sets of literals pr, ps that represent a ’s pre- and postconditions respectively, and a number $d \in \mathbb{N}$ representing the action’s duration. Given an action $a = \langle pr, ps, d \rangle$, we write $pr(a)$, $ps(a)$ and $d(a)$ for pr, ps , and d . Postconditions are divided into *add* $(ps(a)^+)$ and *delete* $(ps(a)^-)$ postcondition sets. An action a can be executed in state s iff the state satisfies its preconditions. The postconditions of a durative action are applied in the state s at which the action ends, by adding the positive postconditions belonging to $ps(a)^+$ and deleting the negative postconditions belonging to $ps(a)^-$.

Goals *Achievement* goals instantaneously achieve a certain state of affairs. Each $g \in G$ is a well-defined set of literals

$g = \{r_1, \dots, r_n\}$, known as *goal requirements* (denoted as r_i), that should be satisfied in the state to satisfy the goal.

Norms An action-based norm is defined as a tuple $n = \langle d_{\text{o}}, a_{\text{con}}, a_{\text{sub}}, dl \rangle$, where $d_{\text{o}} \in \{o, f\}$ is the deontic operator denoting obligation or prohibition; $a_{\text{con}} \in A$ is the action that activates the norm; $a_{\text{sub}} \in A$ is the action that is the subject of the obligation or prohibition; and $dl \in \mathbb{N}$ is the norm deadline relative to the completion of the execution of the action a_{con} , the activation condition of the norm.

2.1 Semantics

Let $P = \langle FL, \Delta, A, G, N \rangle$ be a normative planning problem. Also let $\pi = \langle (a_0, 0), \dots, (a_n, t_{a_n}) \rangle$ be a sequence of actions such that $\nexists (a_i, t_{a_i}), (a_j, t_{a_j}) \in \pi$ s.t. $t_{a_i} \leq t_{a_j} < t_{a_i} + d(a_i), (a_i, a_j) \in cf_{\text{action}}$, where cf_{action} is defined below.

Definition 2 (Conflicting Actions). Actions a_i and a_j have a *concurrency conflict* iff the preconditions or postconditions of a_i contradict the preconditions or postconditions of a_j .

$$cf_{\text{action}} = \{(a_i, a_j) \text{ s.t. } \exists r \in pr(a_i) \cup ps(a_i), \\ \neg r \in pr(a_j) \cup ps(a_j)\}$$

The duration of a sequence of actions π is calculated as $Makespan(\pi) = \max(t_{a_i} + d(a_i))$. The execution of π from a starting state s_0 brings about a sequence of states $S(\pi) = \langle s_0, \dots, s_m \rangle$ for every discrete time interval from 0 to $m = Makespan(\pi)$. The transition relation between two states is as follows. Let A_k be the set of action, time pairs such that the actions end at state s_k . State s_k results from removing all delete postconditions and adding all add postconditions of actions in A_k to state s_{k-1} . I.e., $\forall 0 < k \leq m$:

$$s_k = \begin{cases} (s_{k-1} \setminus \bigcup_{a \in A_k} ps(a)^-) \cup \bigcup_{a \in A_k} ps(a_i)^+ & A_k \neq \emptyset \\ s_{k-1} & A_k = \emptyset \end{cases}$$

π **satisfies a goal** if there is a state that satisfies the goal: $\pi \models g$ iff $\exists s_k \in S(\pi)$ s.t. $s_k \models g$. The set of satisfied goals by π is denoted G_π .

π **complies with an obligation** if the action that is the subject of the obligation, a_{sub} , occurs during the compliance period (i.e., between when the condition holds and when the deadline expires):

$$\pi \models n \text{ iff } (a_{\text{con}}, t_{a_{\text{con}}}), (a_{\text{sub}}, t_{a_{\text{sub}}}) \in \pi \text{ s.t.} \\ t_{a_{\text{sub}}} \in [t_{a_{\text{con}}} + d(a_{\text{con}}), dl + t_{a_{\text{con}}} + d(a_{\text{con}}))$$

If a_{sub} does not occur during the compliance period, the obligation is violated: $\pi \not\models n$.

π **complies with a prohibition** if the prohibition's subject action a_{sub} does not occur during the compliance period:

$$\pi \models n \text{ iff } (a_{\text{con}}, t_{a_{\text{con}}}) \in \pi, \nexists (a_{\text{sub}}, t_{a_{\text{sub}}}) \in \pi \text{ s.t.} \\ t_{a_{\text{sub}}} \in [t_{a_{\text{con}}} + d(a_{\text{con}}), dl + t_{a_{\text{con}}} + d(a_{\text{con}}))$$

If a_{sub} occurs during the compliance period, the prohibition is violated: $\pi \not\models n$.

We assume that all norm deadlines end before $Makespan(\pi)$. Therefore, all activated norms in π (denoted as N_π) are either complied with (denoted $N_{\text{cmp}}(\pi)$) or

violated (denoted $N_{\text{vol}}(\pi)$) by time m . Another assumption made is deontic detachment [Andrighetto *et al.*, 2013], meaning that norm instances are unique, even if a norm is invoked several times.

2.2 Conflict

In this section we consider several types of conflict, defining which sequences of action within a plan are determined to be in conflict. We consider a running example where an agent has the goals of going on strike; submitting a report; and getting a certificate of some sort. However, if the agent goes on maternity leave, it cannot go to the office and submit the report. Moreover, if it goes on strike, it cannot go to office or attend meetings.

Definition 3 (Conflicting Goals). Goal g_i and g_j are in conflict if satisfying them requires bringing about conflicting state of affairs.

$$cf_{\text{goal}} = \{(g_i, g_j) \text{ s.t. } \exists r \in g_i, \neg r \in g_j\}$$

Example 1. The goal *strike*, made up of fluents $\{union_member, \neg at_office, \neg meeting_attended\}$ and goal *submission*, with fluents $\{at_office, report_finalised\}$ are in conflict.

Definition 4 (Conflicting Obligations and Goals). Norm $n = \langle o, a_{\text{con}}, a_{\text{sub}}, dl \rangle$ and goal g are in conflict if executing action a_{sub} — the subject of the obligation — brings about postconditions that are in conflict with the requirements of g .

$$cf_{\text{goalobl}} = \{(g, n) \text{ s.t. } \exists r \in g, \neg r \in ps(a_{\text{sub}})\}$$

Example 2. Goal *strike* and norm $n_1 = \langle o, get_company_funding, attend_meeting, 2 \rangle$, obliging meeting attendance if company funding is used are in conflict, since postcondition *meeting_attended* of *attend_meeting* is incompatible with the fluents of *strike*.

Definition 5 (Conflicting Prohibitions and Goals). A prohibition norm $n = \langle f, a_{\text{con}}, a_{\text{sub}}, dl \rangle$ and a goal g are in conflict, if the postconditions of a_{sub} contribute to satisfying g via r (and r cannot be brought about by any other action.), but executing action a_{sub} is prohibited by norm n .

$$cf_{\text{goalpro}} = \{(g, n) \text{ s.t. } \exists r \in g, r \in ps(a_{\text{sub}})\}$$

Example 3. *submission* = $\{at_office, report_finalised\}$ and $n_2 = \langle f, take_maternity_leave, go_to_office, 6 \rangle$ are in conflict since taking maternity leave prevents the agent from going to the office and hence prevents fulfilling the goal of *submission*: $(submission, n_2) \in cf_{\text{goalpro}}$.

The entire set of conflicting goals and norms is defined as:

$$cf_{\text{goalnorm}} = cf_{\text{goalobl}} \cup cf_{\text{goalpro}}$$

Definition 6 (Conflicting Obligations). $n_1 = \langle o, a_{\text{con}}, a_{\text{sub}}, dl \rangle$ and $n_2 = \langle o, b_{\text{con}}, b_{\text{sub}}, dl' \rangle$ are in conflict in the context of π if the obliged actions in n_1 , i.e., a_{sub} , and n_2 , i.e., b_{sub} have a concurrency conflict; and action a_{sub} is in progress during the entire period over which the agent is obliged to execute action b_{sub} .

$$cf_{\text{oblobl}}^\pi = \{(n_1, n_2) \text{ s.t. } (a_{\text{con}}, t_{a_{\text{con}}}), (b_{\text{con}}, t_{b_{\text{con}}}) \in \pi; \\ (a_{\text{con}}, b_{\text{sub}}) \in cf_{\text{action}}; t_{a_{\text{sub}}} \in \\ [t_{a_{\text{con}}} + d(a_{\text{con}}), t_{a_{\text{con}}} + d(a_{\text{con}}) + dl]; [t_{b_{\text{con}}} + \\ d(b_{\text{con}}), t_{b_{\text{con}}} + d(b_{\text{con}}) + dl'] \subseteq [t_{a_{\text{sub}}}, t_{a_{\text{sub}}} + d(a_{\text{sub}})]\}$$

Example 4. Due to the concurrency conflict between actions *attend_meeting* and *attend_interview*, in $n_1 = \langle o, get_company_funding, attend_meeting, 2 \rangle$ and $n_4 = \langle o, take_theory_test, attend_interview, 2 \rangle$ and depending on the way actions are sequenced in a plan, it is possible that in some π : $(n_1, n_4) \in cf_{obl}^\pi$.

Definition 7 (Conflicting Obligations and Prohibitions). An obligation $n_1 = \langle o, a_{con}, a_{sub}, dl \rangle$ and a prohibition $n_2 = \langle f, b_{con}, a_{sub}, dl' \rangle$ are in conflict in the context of π if n_2 forbids the agent to execute action a_{sub} during the entire period over which obligation n_1 obliges the agent to take a_{sub} .

$$cf_{oblpro}^\pi = \{(n_1, n_2) \text{ s.t. } (a_{con}, t_{a_{con}}), (b_{con}, t_{b_{con}}) \in \pi; \\ [t_{a_{con}} + d(a_{con}), t_{a_{con}} + d(a_{con}) + dl] \subseteq \\ [t_{b_{con}} + d(b_{con}), t_{b_{con}} + d(b_{con}) + dl']\}$$

Example 5. The obligation and prohibition $n_1 = \langle o, get_company_funding, attend_meeting, 2 \rangle$, and $n_3 = \langle f, take_maternity_leave, attend_meeting, 6 \rangle$ can be in conflict in some π as they require and forbid attending a meeting. Thus, for some π , $(n_1, n_3) \in cf_{oblpro}^\pi$.

The two sets cf_{obl}^π and cf_{oblpro}^π constitute the set of conflicting norms: $cf_{norm}^\pi = cf_{obl}^\pi \cup cf_{oblpro}^\pi$.

Definition 8 (Plan). A sequence of actions $\pi = \langle (a_0, 0), \dots, (a_n, t_{a_n}) \rangle$ s.t. $\nexists (a_i, t_{a_i}), (a_j, t_{a_j}) \in \pi$ s.t. $t_{a_i} \leq t_{a_j} < t_{a_i} + d(a_i), (a_i, a_j) \in cf_{action}$ is a plan for the normative planning problem $P = (FL, \Delta, A, G, N)$ iff:

- Only the fluents in Δ hold in the initial state: $s_0 = \Delta$
- the preconditions of action a_i holds at time t_{a_i} and throughout the execution of a_i :

$$\forall k \in [t_{a_i}, t_{a_i} + d(a_i)], s_k \models pr(a_i)$$

- the set of goals satisfied by plan π is a non-empty ($G_\pi \neq \emptyset$) consistent subset of goals:

$$G_\pi \subseteq G \text{ and } \nexists g_i, g_j \in G_\pi \text{ s.t. } (g_i, g_j) \in cf_{goal}$$

- there is no conflict between the goals satisfied and norms complied with:

$$\nexists g \in G_\pi \text{ and } n \in N_{cmp(\pi)} \text{ s.t. } (g, n) \in cf_{goalnorm}$$

Note that since norms are action-based and there is no possibility of executing conflicting actions, there will be no conflict between the norms complied with in a plan.

We consider a set of plans Π , and in the next section deal with the problem of choosing the best plan from this set.

3 Identifying the Best Plan

The conflict between an agent's goals and norms often makes it impossible for the agent to satisfy all its goals while complying with all norms triggered in a plan. We begin by considering how to treat each plan as a proposal of actions and how to use argumentation schemes to check the justifiability of a plan proposal with respect to conflicts and preferences. In Section 3.2, we then identify the best plan from the justified subset.

3.1 Generating Arguments

An argumentation framework (AF) consists of a set of arguments and attacks between them [Dung, 1995]: $AF = \langle Arg, Att \rangle, Att \subseteq Arg \times Arg$. In scheme-based approaches [Walton, 1996] arguments are expressed in natural language and a set of critical questions is associated with each scheme, identifying how the scheme can be attacked. Below, we introduce a set of argument schemes and critical questions to reason about a plan proposal with respect to the goals it satisfies and norms it complies with or violates.

Definition 9 (Plan Argument Scheme Arg_π). A plan argument claims that a proposed sequence of actions should be executed because it satisfies a set of goals, and complies with a set of norms while violating some other norms:

- In the initial state Δ
- The agent should execute sequence of actions π
- Which will satisfy set of goals G_π and complies with set of norms $N_{cmp(\pi)}$ and violates set of norms $N_{vol(\pi)}$

Definition 10 (Goal Argument Scheme Arg_g). A goal argument claims that a feasible goal should be satisfied:

- Goal g is *feasible*¹ for the agent
- Therefore, satisfying g is required.

The set of goal argument is denoted as Arg_G .

Definition 11 (Norm Argument Scheme Arg_n). A norm argument claims that an activated norm should be complied with:

- n is an activated norm imposed on the agent in plan π
- Therefore, complying with n is required in π .

The set of norm argument for a plan is denoted as Arg_{N_π} .

Critical Questions for the Plan Argument Scheme

CQ1: Is there a goal argument which attacks Arg_π ? This CQ results in an undercut attack (asymmetric by nature) from a goal argument to a plan argument, when the goal is not satisfied in the plan:

$$\forall Arg_g \in Arg_G \text{ if } \pi \not\models g \text{ then } (Arg_g, Arg_\pi) \in Att$$

CQ2: Is there a norm argument which attacks Arg_π ? This CQ results in an undercut from a norm argument to a plan argument, when the norm is violated in the plan:

$$\forall Arg_n \in Arg_{N_\pi} \text{ if } \pi \not\models n \text{ then } (Arg_n, Arg_\pi) \in Att$$

Critical Questions for the Goal Argument Scheme

CQ3: What goal arguments might attack Arg_g ? This CQ results in a rebut attack (symmetric by definition) between arguments for conflicting goals:

$$\forall Arg_g, Arg_{g'} \in Arg_G$$

$$\text{if } (g, g') \in cf_{goal} \text{ then } (Arg_g, Arg_{g'}) \in Att$$

CQ4: What norm arguments might attack Arg_g ? This CQ results in a rebut attack between arguments for a goal and a norm that are in conflict:

$$\forall Arg_g \in Arg_G, Arg_n \in Arg_{N_\pi}$$

$$\text{if } (g, n) \in cf_{goalnorm} \text{ then } (Arg_g, Arg_n) \in Att$$

¹A goal is feasible if there is at least one plan that satisfies it.

Critical Questions for the Norm Argument Scheme

CQ4: What goal arguments might attack the norm presented by Arg_n ? The previous critical question is associated with argument schemes for norms as well as goals, hence the repetition of the critical question.

$$\begin{aligned} &\forall Arg_g \in Arg_G, Arg_n \in Arg_{N_\pi} \\ &\text{if } (n, g) \in cf_{goalnorm} \text{ then } (Arg_n, Arg_g) \in Att \end{aligned}$$

CQ5: What norm arguments might attack the norm presented by Arg_n ? Conflict between two norms is defined as a contextual conflict that depends upon the context of the plan in which the norms are activated.

$$\begin{aligned} &\forall Arg_n, Arg_{n'} \in Arg_{N_\pi} \\ &\text{if } (n, n') \in cf_{norm}^\pi \text{ then } (Arg_n, Arg_{n'}) \in Att \end{aligned}$$

Preferences between arguments distinguish an attack from a *defeat* (i.e., a successful attack [Amgoud and Cayrol, 2002]). The attack from one argument to another is a defeat if the latter argument is not preferred over the former. However, as discussed in [Prakken, 2012], rebuttal attacks are preference-dependent, whereas undercuts are preference-independent. Thus, attacks due to CQ3, CQ4 and CQ5 need preferences to be resolved, while attacks caused by CQ1 and CQ2 are preference-independent, always resulting in defeat.

We define \succeq^{gn} as a partial preorder on $G \cup N$. \succ^{gn} denotes the strict relation corresponding to \succeq^{gn} . Also, $(\alpha, \beta) \in \sim^{gn}$ iff $(\alpha, \beta) \in \succeq^{gn}$ and $(\beta, \alpha) \in \succeq^{gn}$. The preferences between the goal and norm arguments result from the preference between these entities: $(Arg_\alpha, Arg_\beta) \in \succeq$ iff $(\alpha, \beta) \in \succeq^{gn}$.

An AF for a plan proposal consists of the argument for the plan itself, a set of arguments for goals and arguments for norms that are activated in that plan. Although the set of goal arguments in AFs for plan proposals remain the same across the AFs, the set of norm arguments differs between AFs depending on the norms that are activated by the plan proposal in each AF.

Definition 12 (Plan Proposal AF). The AF for plan proposal π is $AF_\pi = \langle Arg, Def \rangle$, where $Arg = Arg_\pi \cup Arg_G \cup Arg_{N_\pi}$ and Def is defined as: $\forall Arg_\alpha, Arg_\beta \in Arg, (Arg_\alpha, Arg_\beta) \in Def$ iff $(Arg_\alpha, Arg_\beta) \in Att_{CQ1-5}$ and $(Arg_\beta, Arg_\alpha) \notin \succ$.

The next section explains how an AF for a plan proposal is evaluated and used to identify the best plan(s).

3.2 Evaluating the Argumentation Framework

Argumentation semantics [Dung, 1995] are a means for evaluating arguments in an AF. Among the proposed semantics are the *credulous preferred* semantics which several authors have suggested [Caminada, 2006; Prakken, 2006; Oren, 2013] are appropriate for reasoning about actions. Caminada [2006] provides an intuitive way to identify the status of arguments w.r.t. various semantics through labellings. An argument is, respectively, labelled *in*, *out* and *undec*, if it is acceptable, rejected and undecided under a certain semantics. In a complete labelling (\mathcal{L}_{cmp}), an argument is labelled *in* iff all its attackers are labelled *out*, and is labelled *out* iff there exists an attacker for it that is labelled *in*. A complete

labelling in which the set of arguments labelled *in* are maximal (w.r.t. set inclusion) is a preferred labelling (\mathcal{L}_{pr}). An argument is credulously accepted under preferred semantics if it is labelled *in* by at least one preferred labelling.

Definition 13 (Justified Plans). Plan π is *justified* if Arg_π is labelled *in* by at least one preferred labelling for AF_π : $\exists \mathcal{L}_{pr}$ s.t. $Arg_\pi \in in(\mathcal{L})$.

Although all justified plans are internally consistent, they can still be disagreed with externally. That is, there might be further criteria to take into account when identifying the best plan among justified plans. We define the criteria for the best plan(s) using an established set ordering principle in argumentation, the *Democratic* principle: $(S_i, S_j) \in \succeq$ iff $\forall \beta \in S_j \setminus S_i, \exists \alpha \in S_i \setminus S_j$ s.t. $(\alpha, \beta) \in \succ$. Since preferences over goals and norms are partial, comparing two plans based on the set of goals and norms is not always possible. Therefore, absent such preference information, the best plan(s) satisfies the most goals while violating the fewest norms. We start by defining the *goal-dominant* and *norm-dominant* plans, based on which a *better than* relation between plans is defined.

Let $G_\Pi = \{G_{\pi_1}, G_{\pi_2}, \dots, G_{\pi_n}\}$, where G_{π_i} is the set of goals satisfied in plan π_i .

Definition 14 (Goal-dominance). Plan π_i goal-dominates π_j denoted as $(\pi_i, \pi_j) \in \succeq_G$ if

1. \succeq_G is a total preorder on G_Π and $(G_{\pi_i}, G_{\pi_j}) \in \succeq_G$; or
 2. $|G_{\pi_i}| \geq |G_{\pi_j}|$ (i.e., if \succeq_G is not a total preorder on G_Π).
- \succ_G and \sim_G are strictly and equally dominant version of \succeq_G .

Let $N_{vol}(\Pi) = \{N_{vol}(\pi_1), N_{vol}(\pi_2), \dots, N_{vol}(\pi_n)\}$, where $N_{vol}(\pi_i)$ is the set of norms violated in plan π_i .

Definition 15 (Norm-dominance). Plan π_i norm-dominates π_j denoted as $(\pi_i, \pi_j) \in \succeq_N$ if

1. \succeq_N is a total preorder on $N_{vol}(\Pi)$ and $(N_{vol}(\pi_i), N_{vol}(\pi_j)) \in \succeq_N$; or
2. $|N_{vol}(\pi_i)| \geq |N_{vol}(\pi_j)|$ (i.e., if \succeq_N is not a total preorder on $N_{vol}(\Pi)$).

\succ_N and \sim_N are strictly and equally dominant version of \succeq_N .

Definition 16 (Plan Comparison). Plan π_i is *better than* π_j , denoted $(\pi_i, \pi_j) \in \succ_\pi$, iff:

1. π_i is justified and π_j is not; or
2. π_i and π_j are both justified and $(\pi_i, \pi_j) \in \succ_G$; or
3. π_i and π_j are both justified and $(\pi_i, \pi_j) \in \sim_G$ but $(\pi_j, \pi_i) \in \succ_N$.

Plan π_i is *as good as* π_j , denoted $(\pi_i, \pi_j) \in \sim_\pi$, iff $(\pi_i, \pi_j) \notin \succ_\pi$ and $(\pi_j, \pi_i) \notin \succ_\pi$.

The relation \succ_π is irreflexive, asymmetric and transitive, while \sim_π is an equivalence relation on Π .

Definition 17 (Equivalence Classes). Given $\pi \in \Pi$, let $[\pi_i]$ denote the equivalence class to which π_i belongs. $([\pi_i], [\pi_j]) \in \succeq$ iff $(\pi_i, \pi_j) \in \succ_\pi$ or $(\pi_i, \pi_j) \in \sim_\pi$.

Definition 18 (Best Plan(s)). Plan π_i is (one of) the best plan(s) for the agent to execute iff

- π_i is justified, and
- $\nexists \pi_j$ such that $([\pi_j], [\pi_i]) \in \succeq$.

Example 6. Assume an agent with three goals *strike*, *submission* (c.f., Example 1), and *certificate* =

$\{course_fee_paid, theory_test_done, interviewed\}$ and four norms n_1, n_2, n_3 , and n_4 (c.f., Examples 2, 3, 4, and 5). Assume that the agent prefers satisfying goal *submission* to complying with norm n_2 : $submission \succeq n_2$, and prefers complying with n_4 over n_1 : $n_4 \succeq n_1$. Let $\pi_1, \pi_2, \pi_3, \pi_4 \in \Pi$:

- $\pi_1 \models submission, N_{\pi_1} = \{n_1\}, N_{cmp(\pi_1)} = \{n_1\}$
- $\pi_2 \models submission, certificate, N_{\pi_2} = \{n_1, n_4\}, N_{cmp(\pi_2)} = \{n_4\}, N_{vol(\pi_2)} = \{n_1\}$
- $\pi_3 \models submission, certificate, N_{\pi_3} = \{n_1, n_2, n_3, n_4\}, N_{cmp(\pi_3)} = \{n_3, n_4\}, N_{vol(\pi_3)} = \{n_1, n_2\}$
- $\pi_4 \models strike, certificate, N_{\pi_4} = \{n_1, n_2, n_3, n_4\}, N_{cmp(\pi_4)} = \{n_2, n_3, n_4\}, N_{vol(\pi_4)} = \{n_1\}$

Figure 1 displays the argumentation graph associated with each of these plans². Plan π_1 is not justified, whereas π_2, π_3 and π_4 all are. Thus, the first condition in Definition 18 holds for the last three plans. Since the preferences provided over goals and norms is minimal, in this example the number of goals satisfied and norms violated determines the best plans as follows: although $|G_{\pi_2}| = |G_{\pi_3}| = |G_{\pi_4}|, |N_{vol(\pi_2)}| = |N_{vol(\pi_4)}| < |N_{vol(\pi_3)}|$. Therefore, $\pi_2 \succ \pi_3, \pi_4 \succ \pi_3$, and $\pi_2 \sim \pi_4$, which makes π_2 and π_4 the best plans.

3.3 Properties

We now consider the properties of our system with regards to the rationality postulates [Caminada and Amgoud, 2007], and investigate the properties of the best plan(s) and the preferred extensions that include it.

Property 1. *Closure: The conclusions of any extension (in labelled arguments) are closed under strict rules.*

Proof. Since all arguments are built on defeasible rules, the property follows immediately. \square

Property 2. *Direct Consistency: The conclusions of any extension are consistent.*

Proof. Suppose the conclusions of extension E are inconsistent, i.e., there are arguments $Arg_\alpha, Arg_\beta \in E$ such that:

- Arg_α 's conclusion requires executing plan π and Arg_β 's conclusion requires satisfying goal g /complying with norm n , while g is not satisfied/ n is violated in π . Thus, Arg_β defeats Arg_α ; E is not conflict-free and cannot be an extension.
- Arg_α 's conclusion requires satisfying goal g /complying with norm n and Arg_β 's conclusion requires satisfying goal g /complying with norm n' , while g/n and g/n' are inconsistent. Thus, Arg_α attacks Arg_β and vice versa. Due to the preferences, at least one of these attacks is identified as defeat and therefore E is not conflict-free and not an extension. \square

Property 3. *Indirect Consistency: The closure under strict rules of the conclusions of any extension is consistent.*

Proof. Follows immediately from lack of strict rules. \square

Property 4. *If a plan argument is labelled in by preferred labelling \mathcal{L} , the arguments representing all the goals that it does not satisfy and norms it violates are labelled out by \mathcal{L} and vice versa:*

$$Arg_\pi \in in(\mathcal{L}) \Leftrightarrow Arg_{g \in G \setminus G_\pi} \cup Arg_{n \in N_{vol(\pi)}} \subseteq out(\mathcal{L}).$$

²st, sub and cer stand for *strike, submission* and *certificate*.

Proof. Every preferred labelling is a complete labelling. An argument is labelled *in* by a complete labelling iff all its attackers are labelled *out*. Therefore, a plan argument is labelled *in* by a preferred labelling iff all its attackers, namely the arguments for goals that it does not satisfy and norms that it violates, are labelled *out* by that labelling. \square

Property 5. *If a plan argument is labelled in by preferred labelling \mathcal{L} , the arguments representing all the goals that it satisfies and norms it complies with are also labelled in:*

$$Arg_\pi \in in(\mathcal{L}) \Rightarrow Arg_{g \in G_\pi} \cup Arg_{N_{cmp(\pi)}} \subseteq in(\mathcal{L}).$$

Proof. Since $Arg_\pi \in in(\mathcal{L})$, from Property 4 we know that $Arg_{g \in G \setminus G_\pi} \cup Arg_{n \in N_{vol(\pi)}} \subseteq out(\mathcal{L})$. We also know from the definition of a plan that $Arg_{g \in G_\pi} \cup Arg_{n \in N_{cmp(\pi)}}$ is conflict free. Since all possible attackers of $Arg_{g \in G_\pi} \cup Arg_{n \in N_{cmp(\pi)}}$ belong to $Arg_{g \in G \setminus G_\pi} \cup Arg_{n \in N_{vol(\pi)}}$ and $Arg_{g \in G \setminus G_\pi} \cup Arg_{n \in N_{vol(\pi)}}$ are all labelled *out*, we conclude that $Arg_{g \in G_\pi} \cup Arg_{n \in N_{cmp(\pi)}} \subseteq in(\mathcal{L})$. \square

Note that from $Arg_{g \in G_\pi} \cup Arg_{n \in N_{cmp(\pi)}} \subseteq in(\mathcal{L})$ one cannot conclude that $Arg_\pi \in in(\mathcal{L})$, as there might be justified goals or norms not satisfied or complied with in the plan.

Property 6. *There is no more than one preferred labelling in which $Arg_\pi \in in(\mathcal{L})$.*

Proof. From Properties 4 and 5 we know that if $Arg_\pi \in in(\mathcal{L})$ then $Arg_{g \in G \setminus G_\pi} \cup Arg_{n \in N_{vol(\pi)}} \subseteq out(\mathcal{L})$ and $Arg_{g \in G_\pi} \cup Arg_{n \in N_{cmp(\pi)}} \subseteq in(\mathcal{L})$. Since every preferred labelling is a complete labelling and the following property holds for complete labellings: if $out(\mathcal{L}_{cmp1}) = out(\mathcal{L}_{cmp2})$ then $\mathcal{L}_{cmp1} = \mathcal{L}_{cmp2}$; we conclude that there is no more than one preferred labelling in which $Arg_\pi \in in(\mathcal{L})$. \square

Property 7. *If $Arg_\pi \in in(\mathcal{L})$, \mathcal{L} is a stable labelling.*

Proof. In Property 4 we showed that if $Arg_\pi \in in(\mathcal{L})$ then $Arg_{g \in G \setminus G_\pi} \cup Arg_{n \in N_{vol(\pi)}} \subseteq out(\mathcal{L})$ and $Arg_{g \in G_\pi} \cup Arg_{n \in N_{cmp(\pi)}} \subseteq in(\mathcal{L})$, which makes the $undec(\mathcal{L}) = \emptyset$. A preferred labelling with $undec(\mathcal{L}) = \emptyset$ is a stable labelling. Therefore, \mathcal{L} is a stable labelling. \square

Property 8. *Let \succeq^{gn} be a total preorder on $G \cup N$ and therefore \succeq be a total preorder on goal and norm arguments. If $Arg_\pi \in in(\mathcal{L})$, and the set of arguments for the most preferred goals and norms, $Pref(Arg)$, is conflict free, all arguments belong to $Pref(Arg)$ are labelled in by \mathcal{L} .*

Proof. Elements of set $Pref(Arg)$ cannot be defeated, since the set is conflict-free and the remaining arguments belong to $Arg \setminus Pref(Arg)$. The latter cannot defeat elements of $Pref(Arg)$, because this would imply an attack from a less preferred argument to a more preferred one has resulted in a defeat, which is contrary to assumption. Assume that $\exists Arg_\alpha \in Pref(Arg)$ such that $Arg_\alpha \notin in(\mathcal{L})$. If $\nexists Arg_\beta \in in(\mathcal{L})$ s.t. $(Arg_\alpha, Arg_\beta) \in Def$ then Arg_α should have been labelled *in* by \mathcal{L} otherwise it is contrary to the assumption of maximality of preferred labellings. If $\exists Arg_\beta \in in(\mathcal{L})$ s.t. $(Arg_\alpha, Arg_\beta) \in Def$ then $\exists Arg_\gamma \in in(\mathcal{L})$ s.t. $(Arg_\gamma, Arg_\alpha) \in Def$, which is contradictory to

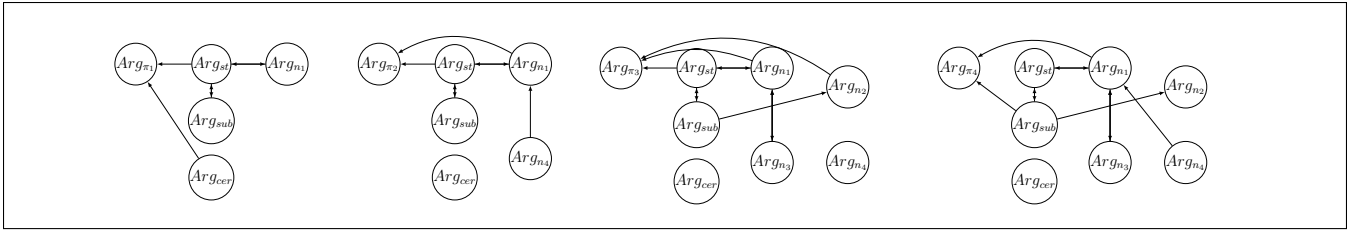


Figure 1: Argumentation Graph for plans $\pi_1, \pi_2, \pi_3, \pi_4$

the fact that Arg_α cannot be defeated. Therefore, all elements of $Pref(Arg)$ are labelled in by $in(\mathcal{L})$. \square

4 Explaining the Justifiability of the Best Plan

In this section we exploit an existing dialogue for preferred semantics known as *Socratic Discussion* [Caminada *et al.*, 2014a] to provide an explanation for the justifiability of the best plan(s). Deciding if an argument is in at least one preferred extension amounts to deciding if it is in at least in one admissible extension (i.e. it is labelled *in* by at least one admissible labelling). In an admissible labelling if an argument is labelled *in*, all its attackers are labelled *out*, and if an argument is labelled *out*, it has an attacker that is labelled *in*.

Definition 19 (Socratic Discussion [Caminada *et al.*, 2014a]). Let $AF = \langle Arg, Def \rangle$. The sequence of moves $[\Delta_1, \Delta_2, \dots, \Delta_n]$ ($n \geq 1$) is a Socratic discussion iff: (i) each odd move (M-move) is an argument labelled *in*; (ii) each even move (S-move) is an argument labelled *out*; (iii) each argument moved by *S* attacks an argument moved by *M* earlier in the dialogue; (iv) each argument moved by *M* attacks an argument moved by *S* in the previous step; (v) S-moves cannot be repeated. Player *S* wins the discussion if there is an M-move and an S-move containing the same argument. Otherwise, the winner is the player that makes the last move.

Given that the agent's best plans(s) π is labelled *in* by at least one preferred labelling, player *M* is guaranteed a winning strategy in a Socratic discussion with $\Delta_1 = in(Arg_\pi)$. The even moves in the rest of dialogue are arguments labelled *out*, which according to Property 4 are goals not satisfied or norms violated in π . On the other hand, the rest of odd moves in the dialogue are arguments labelled *in*, which according to Property 5 are goals satisfied or norms complied with in π . Since each odd move attacks the even move in the previous step, during a dialogue the agent is able to dialectically explain why it did not satisfy a goal or violate a norm, which are the two causes of attacks on plan proposals.

Example 7. This example shows a Socratic discussion $\Delta = [in(Arg_{\pi_4}), out(Arg_{sub}), in(Arg_{st}), out(Arg_{n_1}), in(Arg_{n_4})]$ for plan π_4 .

- M: Plan π_4 is (one of) the best plan(s) and is justifiable.
- S: Why does the plan not satisfy goal *submission*?
- M: As it satisfies goal *strike*, attacking goal *submission*.
- S: Why does the plan violate norm n_1 ?
- M: Because the plan satisfies norm n_4 that attacks norm n_1 .

5 Related Work

One of the most well-known scheme-based approach in practical reasoning is [Atkinson and Bench-Capon, 2007]. Recently, [Oren, 2013] has proposed a similar scheme-based approach for normative practical reasoning where arguments are constructed for a sequence of actions. Similar to the latter approach, we construct arguments for plans rather than actions. [Oren, 2013] assumes that conflicts within and between goals and norms are inferred from paths, rather than being obtained from the model. Thus, although it is possible to explain why one path is preferred over another, it is not possible to understand why a path does not satisfy a goal or violate a norm. In contrast, we explicitly consider why an agent does not satisfy a goal, or violate a norm. In addition, in this work the explanation of justifiability of why a plan is (one of) the best plan(s) for the agent to execute is formulated using a dialogue game for preferred semantics.

There are several applications where dialogue games are used for explanation purposes [Zhong *et al.*, 2014; Fan and Toni, 2015; Caminada *et al.*, 2014b]. In [Zhong *et al.*, 2014] and [Fan and Toni, 2015] admissible dispute trees developed for Assumption-based Argumentation [Dung *et al.*, 2009] are used to provide explanation for why a certain decision is better than another. In [Caminada *et al.*, 2014b] a dialogical proof based on the grounded semantics [Caminada and Podlaskowski, 2012] is created to justify the actions executed in a plan. Despite the popularity of the preferred semantics, they have not been used for explanation in such contexts, and our work is the first to do so in the practical reasoning domain.

6 Conclusion and Future Work

In contrast to existing argument based practical reasoning approaches, we propose a framework that integrates both the reasoning and dialogical aspects of argumentation to perform normative practical reasoning. In doing so, we answer the question of how an agent should act in a normative environment under conflicting goals and norms. Moreover, our approach can generate explanation for agent behaviour using an argumentation-based dialogue.

In future work we will investigate temporal solutions to addressing goal-goal and goal-norm conflict, similar to how conflicts between norms are handled. We also intend to empirically evaluate the effectiveness of our explanations, determining how likely a human is to accept the recommendation of a system regarding the best plan(s).

References

- [Amgoud and Cayrol, 2002] Leila Amgoud and Claudette Cayrol. A reasoning model based on the production of acceptable arguments. *Annals of Mathematics Artificial Intelligence*, 34(1-3):197–215, 2002.
- [Amgoud and Prade, 2009] Leila Amgoud and Henri Prade. Using arguments for making and explaining decisions. *Artificial Intelligence*, 173(3-4):413–436, 2009.
- [Andrighetto *et al.*, 2013] Giulia Andrighetto, Guido Governatori, Pablo Noriega, and Leendert van der Torre, editors. *Normative Multi-Agent Systems*, volume 4 of *Dagstuhl Follow-Ups*. Schloss Dagstuhl - Leibniz-Zentrum fuer Informatik, 2013.
- [Atkinson and Bench-Capon, 2007] Katie Atkinson and Trevor J. M. Bench-Capon. Practical reasoning as presumptive argumentation using action based alternating transition systems. *Artificial Intelligence*, 171(10-15):855–874, 2007.
- [Broersen *et al.*, 2001] Jan Broersen, Mehdi Dastani, Joris Hulstijn, Zisheng Huang, and Leendert van der Torre. The boid architecture: Conflicts between beliefs, obligations, intentions and desires. In *Proceedings of the Fifth International Conference on Autonomous Agents*, AGENTS '01, pages 9–16, New York, NY, USA, 2001. ACM.
- [Caminada and Amgoud, 2007] Martin Caminada and Leila Amgoud. On the evaluation of argumentation formalisms. *Artificial Intelligence*, 171(5-6):286–310, 2007.
- [Caminada and Podlaszewski, 2012] Martin Caminada and Mikolaj Podlaszewski. Grounded semantics as persuasion dialogue. In Bart Verheij, Stefan Szeider, and Stefan Woltran, editors, *Computational Models of Argument - Proceedings of COMMA 2012, Vienna, Austria, September 10-12, 2012*, volume 245 of *Frontiers in Artificial Intelligence and Applications*, pages 478–485. IOS Press, 2012.
- [Caminada *et al.*, 2014a] Martin Caminada, Wolfgang Dvork, and Srdjan Vesic. Preferred semantics as socratic discussion. *Journal of Logic and Computation (JLC)*, 2014.
- [Caminada *et al.*, 2014b] Martin Caminada, Roman Kutlak, Nir Oren, and Wamberto Weber Vasconcelos. Scrutable plan enactment via argumentation and natural language generation. In Ana L. C. Bazzan, Michael N. Huhns, Alessio Lomuscio, and Paul Scerri, editors, *International conference on Autonomous Agents and Multi-Agent Systems, AAMAS '14, Paris, France, May 5-9, 2014*, pages 1625–1626. IFAAMAS/ACM, 2014.
- [Caminada, 2006] Martin Caminada. On the issue of reinstatement in argumentation. In Michael Fisher, Wiebe van der Hoek, Boris Konev, and Alexei Lisitsa, editors, *Logics in Artificial Intelligence, 10th European Conference, JELIA 2006, Liverpool, UK, September 13-15, 2006, Proceedings*, volume 4160 of *Lecture Notes in Computer Science*, pages 111–123. Springer, 2006.
- [Dung *et al.*, 2009] Phan Minh Dung, Robert Kowalski, and Francesca Toni. Assumption-based argumentation. In *Argumentation in Artificial Intelligence*. Springer, 2009.
- [Dung, 1995] Phan Minh Dung. On the acceptability of arguments and its fundamental role in nonmonotonic reasoning, logic programming and n-person games. *Artificial Intelligence*, 77(2):321–358, 1995.
- [Fan and Toni, 2015] Xiuyi Fan and Francesca Toni. On computing explanations in argumentation. In Blai Bonet and Sven Koenig, editors, *Proceedings of the Twenty-Ninth AAAI Conference on Artificial Intelligence, January 25-30, 2015, Austin, Texas, USA.*, pages 1496–1502. AAAI Press, 2015.
- [Fikes and Nilsson, 1971] Richard E. Fikes and Nils J. Nilsson. Strips: A new approach to the application of theorem proving to problem solving. In *Proceedings of the 2Nd International Joint Conference on Artificial Intelligence, IJCAI'71*, pages 608–620, San Francisco, CA, USA, 1971. Morgan Kaufmann Publishers Inc.
- [Kollingbaum and Norman, 2003] Martin J. Kollingbaum and Timothy J. Norman. Noa - A normative agent architecture. In Georg Gottlob and Toby Walsh, editors, *IJCAI-03, Proceedings of the Eighteenth International Joint Conference on Artificial Intelligence, Acapulco, Mexico, August 9-15, 2003*, pages 1465–1466. Morgan Kaufmann, 2003.
- [Oren, 2013] Nir Oren. Argument schemes for normative practical reasoning. In Elizabeth Black, Sanjay Modgil, and Nir Oren, editors, *TFAA*, volume 8306 of *Lecture Notes in Computer Science*, pages 63–78. Springer, 2013.
- [Prakken, 2006] Henry Prakken. Combining sceptical epistemic reasoning with credulous practical reasoning. In Paul E. Dunne and Trevor J. M. Bench-Capon, editors, *Computational Models of Argument: Proceedings of COMMA 2006, September 11-12, 2006, Liverpool, UK*, volume 144 of *Frontiers in Artificial Intelligence and Applications*, pages 311–322. IOS Press, 2006.
- [Prakken, 2012] Henry Prakken. Some reflections on two current trends in formal argumentation. In Alexander Artikis, Robert Craven, Nihan Kesim Cicekli, Babak Sadighi, and Kostas Stathis, editors, *Logic Programs, Norms and Action - Essays in Honor of Marek J. Sergot on the Occasion of His 60th Birthday*, volume 7360 of *Lecture Notes in Computer Science*, pages 249–272. Springer, 2012.
- [Walton, 1996] Douglas N. Walton. *Argumentation Schemes for Presumptive Reasoning*. L. Erlbaum Associates, 1996.
- [Zhong *et al.*, 2014] Qiaoting Zhong, Xiuyi Fan, Francesca Toni, and Xudong Luo. Explaining best decisions via argumentation. In Andreas Herzig and Emiliano Lorini, editors, *Proceedings of the European Conference on Social Intelligence (ECSI-2014), Barcelona, Spain, November 3-5, 2014.*, volume 1283 of *CEUR Workshop Proceedings*, pages 224–237. CEUR-WS.org, 2014.