

# Reasoning With and About Norms in Logical Argumentation

Kees VAN BERKEL <sup>a,1</sup> and Christian STRASSER <sup>b</sup>

<sup>a</sup>*Institute of Logic and Computation, TU Wien, Austria*

<sup>b</sup>*Institute for Philosophy II, Ruhr University Bochum, Germany*

This is a preprint of the article published in the proceedings of the 9th International Conference on Computational Models of Argument (COMMA 2022).

**Abstract.** Normative reasoning is inherently defeasible. Formal argumentation has proven to be a unifying framework for representing nonmonotonic logics. In this work, we provide an argumentative characterization of a large class of Input/Output logics, a prominent defeasible formalism for normative reasoning. In many normative reasoning contexts, one is not merely interested in knowing whether a specific obligation holds, but also in why it holds despite other norms to the contrary. We propose sequent-style argumentation systems called Deontic Argument Calculi (DAC), which serve transparency and bring meta-reasoning about the inapplicability of norms to the object language level. We prove soundness and completeness between DAC-instantiated argumentation frameworks and constrained Input/Output logics. We illustrate our approach in view of two deontic paradoxes.

**Keywords.** Nonmonotonic logic, Argumentation, Deontic logic, Normative reasoning

## 1. Introduction

Obligations and norms fulfil a crucial role in a variety of fields, including law, ethics, AI, and everyday life [1]. The logical study of normative reasoning investigates reasoning with such concepts in formal systems of logic, e.g., deontic logics. Its importance increases with the development of intelligent autonomous systems. Complex normative systems often require reasoning with normative conflicts, exceptions, preferences and priorities [1]. A central challenge is to provide transparent formal models of the underlying reasoning processes, e.g., by means of nonmonotonic logics.

Over the past decades, abstract argumentation has proven to be a unifying framework for the representation of large classes of nonmonotonic logics [2]. Formal argumentation provides both a natural and a transparent model of conflicts and their resolution in terms of conflicting arguments. In this way, it provides a promising basis for tackling the challenging requirements of normative reasoning. The logical analysis of normative reasoning is well-established [1] with the Input/Output framework (I/O) being one of the central approaches [3]. Nonmonotonicity is captured in constrained I/O logics through

---

<sup>1</sup>Corresponding Authors: Kees van Berkel (e-mail: kees@logic.at) and Christian Straßer (e-mail: christian.strasser@rub.de). Acknowledgments: This work is supported by WWTF and FWF projects MA16-028 and W1255-N23. We thank Leon van der Torre for discussions that significantly shaped this work.

considering maximal consistent families of norms. In recent years, also argumentative representations of deontic logics have attracted increasing interest [4,5,6,7,8,9].

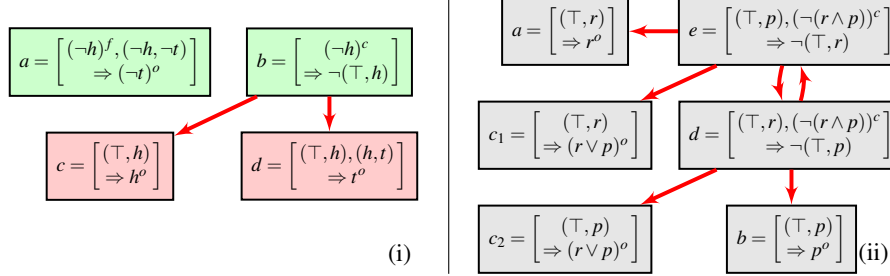
This paper is the first to provide argumentative characterizations for a significant class of I/O logics, including all original logics from [3]. In this way, we are able to combine the advantages of I/O logic with those of formal argumentation. On the one hand, I/O is a highly expressive and robust framework with two decades of developments, including many applications (e.g., priorities, constitutive norms, cognitive modeling, causal reasoning [1,10]). On the other hand, it does not provide the level of transparency that comes with the explicit representation of conflicts in formal argumentation.

In particular, I/O leaves some central challenges of normative reasoning unaddressed. When answering the question as to *why* an obligation holds, one must state *reasons*. Moreover, often it does not suffice to know why a specific obligation holds, one must know why other obligations to the contrary *do not hold*. E.g., in order to understand why “I am permitted to overtake on the left, *despite* having to drive on the right” one must know how the first norm relates to the second. In this case, the first is an exception that renders the latter *inapplicable* in the context of “overtaking another vehicle”. Common approaches to I/O logics—as well as deontic logics—do not provide means for making explicit the reasons why certain obligations are not derivable. Despite their central role in ethics and explanation [11], a general lack of explicit modeling of reasons in formal systems has been recently identified [12] (with some exceptions, e.g., [13]). Support and defeat relations are central in the context of reasons as well as in formal argumentation, which makes the latter an ideal framework to reason with and about reasons.

We address these problems by introducing a class of rule-based proof systems called *Deontic Argument Calculi* (DAC) for normative reasoning by means of argumentation.

Our conceptual contribution is twofold: First, we use labels on formulae to make the presentation transparent on the object level, i.e., we can syntactically distinguish between facts, obligations, and constraints without “burdening” the logics with modalities [10]. Second, we internalize some of the meta-reasoning in the I/O formalism by referring to the inapplicability of norms on the object language level. Consequently, our calculi generate both arguments that provide *explicit reasons* for obligations and arguments that defeat other arguments by giving explicit reasons for why certain norms are inapplicable. The second type of arguments concerns the nonmonotonicity of normative reasoning. The possibility to reason about the inapplicability of norms on the object language level distinguishes our work from other systems such as [14,15]. We illustrate the utility of our approach using the notion of *related admissibility* [16] to explain why some obligation holds *despite* certain norms to the contrary.

The technical contribution of this work contains two types of completeness results: First, we show adequacy between DAC and a significant class of monotonic I/O consequence relations. Second, we prove that formal argumentation frameworks instantiated with DAC arguments characterize a large class of nonmonotonic I/O logics. This makes our work the first to argumentatively characterize I/O logic. Moreover, DAC enjoys a modularity particularly suitable for expansions and our calculi are modular with respect to a large class of base logics. Last, our work contributes to previous representation results in formal argumentation concerning systems based on maximal consistent sets [2].



**Figure 1.** Defeasible normative reasoning examples: (i) The Chisholm scenario (Example 1). Arrows denote defeat relations between arguments, relative to  $\mathcal{C}' = \{-h^c\}$  (Example 2). (ii) A deontic conflict (Example 3). Argument  $e$  defends  $\{b, c_2, e\}$ , whereas argument  $d$  defends  $\{a, c_1, d\}$ .

## 2. Basic Terminology and Benchmark Examples

We introduce basic terminology by considering two examples. Developments in deontic logic are driven by challenging examples [17]. Here, we focus on *contrary-to-duty* reasoning and *deontic conflicts*. Both can be effectively addressed using nonmonotonic reasoning [10] (for alternative approaches see [1]). The language  $\mathcal{L}$  is defined as follows:

$$\varphi ::= p \mid \top \mid \perp \mid \neg\varphi \mid \varphi \vee \varphi \mid \varphi \wedge \varphi \mid \varphi \rightarrow \varphi$$

with  $p \in \text{Atoms}$ . All connectives are primitive in order to be modular with respect to the base logic (Section 3). We use  $p, q, r, \dots$  for atoms, and reserve  $\varphi, \psi, \theta, \dots$  for arbitrary formulae of  $\mathcal{L}$ . In order to increase transparency we *label* formulae of  $\mathcal{L}$ , i.e.,  $\mathcal{L}^i = \{\varphi^i \mid \varphi \in \mathcal{L}\}$  for  $i \in \{f, o, c\}$ . We have formulae expressing *facts*  $\mathcal{L}^f$ , *obligations*  $\mathcal{L}^o$ , and *constraints*  $\mathcal{L}^c$ . Moreover, we employ pairs of formulae  $\mathcal{L}^n = \{(\varphi, \psi) \mid \varphi, \psi \in \mathcal{L}\}$ . A pair  $(\varphi, \psi)$  represents a *norm*: i.e., “given fact  $\varphi$ , it is obligatory that  $\psi$ ” [3].

We work with knowledge bases of the type  $\langle \mathcal{F}, \mathcal{N}, \mathcal{C} \rangle$ , where  $\mathcal{F} \subseteq \mathcal{L}^f$  constitutes the factual context,  $\mathcal{N} \subseteq \mathcal{L}^n$  denotes a system of norms, and  $\mathcal{C} \subseteq \mathcal{L}^c$  represents the constraints with which output must be consistent. The basic idea is that facts (input) trigger norms, from which obligations are detached (output). Moreover, constraints control the output to ensure consistency. The above is in the spirit of constrained I/O logic [3].

Suppose we have a single fact  $\mathcal{F} = \{p^f\}$ , a norm system  $\mathcal{N} = \{(p, q)\}$ , and no constraints, then an argument concluding that  $q$  is obligatory is of the following form,

$$p^f, (p, q) \Rightarrow q^o$$

The left-hand side (lhs) gives reasons for the conclusion on the right-hand side (rhs).

**Example 1** (Chisholm scenario [1], Figure 1-i). *Jones is under the obligation to go and help her neighbors  $(\top, h)$ .<sup>2</sup> Furthermore, Jones knows if she goes to help, she must tell them she goes  $(h, t)$ . Now, if Jones does not go, she ought not to tell them she goes  $(\neg h, \neg t)$ . It turns out that Jones does not go to help  $\neg h^f$ . Clearly, Jones has violated her primary obligation to go and help. Let the knowledge base be  $\mathcal{F} = \{\neg h^f\}$  and  $\mathcal{N} = \{(\top, h), (h, t), (\neg h, \neg t)\}$ . Figure 1-i presents arguments  $a, c$ , and  $d$  that can be constructed from the knowledge base (we explain the meaning of  $b$  and the arrows in Exam-*

<sup>2</sup>A norm  $(\top, \varphi)$  with a precondition  $\top$  is triggered by default, that is, even by an empty factual context.

ple 2); e.g., in argument  $a$ , the reasons for not telling  $\neg t^o$  are the fact  $\neg h^f$  and the norm  $(\neg h, \neg t)$ . What must Jones do in this contrary-to-duty scenario? The desired answer is that she ought not to tell the neighbors she goes  $\neg t^o$ . Formalizations of this scenario cause problems for (monadic) deontic logics, e.g., both  $t$  and  $\neg t$  become obligatory.

Arguments do not only provide reasons in support of an obligation, but also defend them from potential *defeaters*. A rebuttal defeat opposes the conclusion of an argument without pinpointing the reason as to why. In contrast, attacks on reasons—i.e., *undercuts*—are arguments that express *which reasons are inapplicable* given some context. We adopt undercuts since they are more transparent about attacks. Recall that constraints are consistency requirements and suppose  $\mathcal{C} = \{\neg q^c\}$ .<sup>3</sup> Then, a defeating argument

$$p^f, \neg q^c \Rightarrow \neg(p, q)$$

expresses that, if output is to be consistent with the constraint  $\neg q^c$ , in context  $p^f$  the norm  $(p, q)$  cannot be consistently asserted as a reason (since it would detach  $q^o$ ). Hence,  $\neg(p, q)$  expresses that this norm is *inapplicable* given  $\mathcal{F}$  and  $\mathcal{C}$ . An argumentation framework is then simply a set of arguments with defeat relations holding between them.

**Example 2** (Example 1 cont.). *We want to know what Jones must do in the light of her violation  $\neg h^f$ . Thus, we impose the constraint that the output must be consistent with the fact that Jones does not help  $\mathcal{C} = \{\neg h^c\}$  (i.e.,  $\mathcal{C} = \mathcal{F}$  modulo relabelling). This constraint gives us the argument  $b : \neg h^c \Rightarrow \neg(\top, h)$  expressing that given consistency requirement  $\neg h$ , the norm  $(\top, h)$  may not be asserted as a reason (it would output the inconsistent  $h^o$ ). This argument serves as a defeater of any argument which appeals to  $(\top, h)$  in its reasons, in this case both  $c$  and  $d$ ; see the defeat arrows in Figure 1-i. So, that Jones ought not to tell  $\neg t^o$  is explained by argument  $a$  together with the fact that arguments  $c$  and  $d$  concluding helping  $h^o$ , respectively telling  $t^o$ , cannot be defended in view of  $b$ . Namely,  $c$  and  $d$  both employ reasons that are inapplicable given  $\mathcal{C}$ .*

What makes this approach more transparent is the use of labels in arguments to indicate different types of information (factual, obligations, constraints), the internalized meta-reasoning about inapplicability of norms, and the argumentation framework revealing the contrastive dimension of defeasible reasoning. In Figure 1-i, the question “why shouldn’t Jones help, *despite* argument  $c$ ?” is answered by “since argument  $b$  attacks  $c$  and  $b$  is not attacked.” These notions will be made precise in subsequent sections.

**Example 3** (Deontic Conflict [1], Figure 1-ii). *Suppose Smith has an obligation to return a borrowed weapon to a colleague  $(\top, r)$ . Smith knows the colleague is planning to commit a crime with this weapon and Smith is under the obligation to prevent crime  $(\top, p)$ . Furthermore, the constraint is that Smith cannot secure both  $r$  and  $p$ . What should Smith do? This is a deontic conflict. The knowledge base is  $\mathcal{F} = \emptyset$ ,  $\mathcal{N} = \{(\top, r), (\top, p)\}$ , and  $\mathcal{C} = \{\neg(r \wedge p)^c\}$ . Suppose we reason classically, e.g.,  $p$  entails  $p \vee r$ . The arguments that can be constructed are presented in Figure 1-ii. The two defeating arguments,  $d$  and  $e$ , express that given the constraints one of either two norms cannot be asserted.*

Intuitively, the defensible set  $\{a, c_1, d\}$  justifies the obligation that Smith ought to return the weapon in Example 3, whereas  $\{b, c_2, e\}$  does this for the prevention of crime.

<sup>3</sup>Since we do not allow for formulae with mixed labels, we can safely omit brackets w.r.t. using labels.

$$\begin{array}{cccc}
\frac{}{(\top, \top)} \text{T} & \frac{}{(\varphi, \varphi)} \text{ID} & \frac{(\varphi, \psi) \quad \psi \vdash \gamma}{(\varphi, \gamma)} \text{WO} & \frac{(\varphi, \psi) \quad (\varphi, \gamma)}{(\varphi, \psi \wedge \gamma)} \text{AND} \\
\frac{(\varphi, \psi) \quad \gamma \vdash \varphi}{(\gamma, \psi)} \text{SI} & \frac{(\varphi, \psi) \quad (\varphi \wedge \psi, \gamma)}{(\varphi, \gamma)} \text{CT} & \frac{(\varphi, \psi) \quad (\gamma, \psi)}{(\varphi \vee \gamma, \psi)} \text{OR} & 
\end{array}$$

**Figure 2.** Rules for constructing  $\text{deriv}_{\mathcal{R}, \mathbb{L}}$  proof-systems. The minimal set of deriv-rules is  $\{\text{WO}, \text{AND}, \text{SI}\}$ .

Likewise, one can justify the floating conclusion  $(r \vee p)^o$  in Figure 1-ii, by arguing that in every defensible stance *either*  $c_1$  *or*  $c_2$  is selected (cf. disjunctive response [13]). However, following a more skeptical reasoning style one can argue why  $r \vee p$  is not obligatory since there is no single argument concluding  $(r \vee p)^o$  that is selected in *every* defensible stance. Defeasible reasoning by means of argumentation gives rise to various reasoning styles, including the aforementioned. We will discuss these in Section 5.

### 3. Constrained Input/Output Logic

We briefly recall the basics of *Constrained Input/Output logic*, the systems for which we provide argumentative characterizations. The formalism was developed by Makinson and van der Torre [3] and is particularly suitable for normative reasoning [10]. Its central feature is the employment of syntactic objects of the form  $(\varphi, \psi)$ , called *norms*.

I/O logics are construed over the *non*-labelled propositional language  $\mathcal{L}$  (Section 2) and a base logic  $\mathbb{L}$ . We use capital Greek letters  $\Delta, \Gamma, \dots$  for finite sets of  $\mathcal{L}$ -formulae and write  $\wedge \Delta$  to denote the conjunction of elements of  $\Delta$ . Let  $\vdash$  denote the consequence relation of the base logic  $\mathbb{L}$ . We assume that  $\vdash$  is reflexive ( $\Gamma, \varphi \vdash \varphi$ ), transitive ( $\Gamma \vdash \psi$  and  $\Gamma', \psi \vdash \varphi$  implies  $\Gamma, \Gamma' \vdash \varphi$ ) and monotonic ( $\Gamma \vdash \varphi$  implies  $\Gamma, \Gamma' \vdash \varphi$ ). We also assume the presence of a conjunction  $\wedge$ , for which  $\Gamma \vdash \psi \wedge \varphi$  iff  $\Gamma \vdash \psi$  and  $\Gamma \vdash \varphi$ , a negation  $\neg$ , for which  $\Gamma \vdash \varphi$  iff  $\Gamma, \neg \varphi \vdash \perp$ , a disjunction  $\vee$ , for which  $\Gamma, \varphi_1 \vdash \psi$  and  $\Gamma, \varphi_2 \vdash \psi$  iff  $\Gamma, \varphi_1 \vee \varphi_2 \vdash \psi$ , and a falsum constant  $\perp$  for which  $\perp \vdash \varphi$  and  $\varphi, \neg \varphi \vdash \perp$ . We assume  $\mathbb{L}$  has an adequate *sequent calculus* LC, i.e.,  $\Delta \vdash \varphi$  iff the sequent  $\Delta \Rightarrow \varphi$  is LC-derivable.

Constrained I/O logics work with knowledge bases of the type  $\langle \mathcal{F}, \mathcal{N}, \mathcal{C} \rangle$ , where  $\mathcal{F} \subseteq \mathcal{L}$  is the *factual* input,  $\mathcal{N} \subseteq \mathcal{L} \times \mathcal{L}$  a *normative system*, and  $\mathcal{C} \subseteq \mathcal{L}$  a set of *constraints* containing the formulae with which output must be consistent. We assume  $\mathcal{F}$  and  $\mathcal{C}$  to be consistent, i.e.,  $\mathcal{F} \not\vdash \perp$  and  $\mathcal{C} \not\vdash \perp$ . The traditional I/O proof systems are only available for a class of *monotonic* I/O logics [10]. The system is referred to as “deriv” and contains inference rules that derive I/O pairs from other I/O pairs (Figure 2).

**Definition 1.** Let  $\text{deriv}_{\mathcal{R}, \mathbb{L}}$  be a proof-system, with  $\mathcal{R}$  a set of rules from Table 2. Let  $\mathbb{L}$  be the base logic, and let  $\mathcal{N} \subseteq \mathcal{L}^n$ . A derivation of  $(\varphi, \psi) \in \text{deriv}_{\mathcal{R}, \mathbb{L}}(\mathcal{N})$  is a tree of rule-applications of  $\mathcal{R}$  where the leaves are either members of  $\mathcal{N}$  or instances of T and ID (if  $T, ID \in \mathcal{R}$ ), all members of  $\mathcal{N}$  are among the leaves, and the root is  $(\varphi, \psi)$ .

We say  $\psi$  is obligatory (detached) under  $\mathcal{N}$  and  $\mathcal{F}$  if  $(\varphi, \psi) \in \text{deriv}_{\mathcal{R}, \mathbb{L}}(\mathcal{N}')$  with  $\mathcal{F} \vdash \varphi$  and  $\mathcal{N}' \subseteq \mathcal{N}$ . We write  $\psi \in \text{deriv}_{\mathcal{R}, \mathbb{L}}(\Delta, \mathcal{N})$  if  $(\wedge \Delta, \varphi) \in \text{deriv}_{\mathcal{R}, \mathbb{L}}(\mathcal{N})$ .

Paradigmatic I/O logics are characterized by the sets of rules  $\mathcal{R}_1 = \{\text{T}, \text{WO}, \text{SI}, \text{AND}\}$ ,  $\mathcal{R}_2 = \{\text{OR}\} \cup \mathcal{R}_1$ ,  $\mathcal{R}_3 = \mathcal{R}_1 \cup \{\text{CT}\}$ , and  $\mathcal{R}_4 = \mathcal{R}_2 \cup \mathcal{R}_3$ . The system  $\mathcal{R}_1$  represents a *single deontic detachment* procedure which allows for weakening of the output (WO), combining output (AND), and strengthening of the input (SI). All propositional tautolo-

gies are among the output (T). System  $\mathcal{R}_2$  extends  $\mathcal{R}_1$  with *reasoning by cases* (OR), i.e., if both  $\varphi$  and  $\gamma$  generate output  $\psi$ , then  $\varphi \vee \gamma$  generates  $\psi$  too. System  $\mathcal{R}_3$  extends  $\mathcal{R}_1$  with reusability (CT) allowing for iterations of *successive deontic detachment* (cf. chaining reasons in Example 5). Last,  $\mathcal{R}_4$  combines  $\mathcal{R}_2$  and  $\mathcal{R}_3$ . The above systems may be closed under *throughput* (ID), i.e., input is ‘put through’ as output. We write  $\mathcal{R}_i^+ = \mathcal{R}_i \cup \{\text{ID}\}$  for  $i \in \{1, 2, 3, 4\}$ . The resulting eight systems are sound and complete with respect to their semantic characterizations [10]. We omit the semantics here.

The above systems are still monotonic. As Example 1 and 3 demonstrate, we require defeasible detachment. Constrained I/O logics enable this [3]. Constrained I/O logics work with maximal families of norms  $\mathcal{N}' \subseteq \mathcal{N}$  under which the output remains consistent with the constraints  $\mathcal{C}$ . If the output is required to be consistent *per se*, we let  $\mathcal{C} = \emptyset$ . If the output is to be consistent with the input, we take  $\mathcal{F} \subseteq \mathcal{C}$  (e.g., Example 2).

**Definition 2.** Let  $\text{deriv}_{\mathcal{R},L}$  be a system from Figure 2 and let  $\mathcal{K} = \langle \mathcal{F}, \mathcal{N}, \mathcal{C} \rangle$  be a knowledge base. The set of maximal consistent families of  $\mathcal{N}$  (maxfam) is defined as:

- $\text{maxfam}_{\mathcal{R},L}(\mathcal{K})$  is the set of max-elements of  $\{\mathcal{N}' \subseteq \mathcal{N} \mid \text{for all } (\varphi, \psi) \in \text{deriv}_{\mathcal{R},L}(\mathcal{N}'), \text{ if } \mathcal{F} \vdash \varphi, \text{ then } \mathcal{C}, \psi \not\vdash \perp\}$ .

We define *sceptic nonmonotonic inference*  $\vdash^s$  for constrained I/O logic as follows:

- $\mathcal{K} \vdash_{\mathcal{R},L}^s \varphi$  iff  $\forall \mathcal{N}' \in \text{maxfam}_{\mathcal{R},L}(\mathcal{K}), \exists (\psi, \varphi) \in \text{deriv}_{\mathcal{R},L}(\mathcal{N}') \text{ with } \mathcal{F} \vdash \psi$ .

**Example 4** (Example 1 cont.). Consider  $\mathcal{R}_3$  with  $L$  a classical logic,  $\mathcal{F} = \{-h\}$ , and  $\mathcal{N} = \{(\top, h), (h, t), (-h, \neg t)\}$ . For  $\mathcal{C} = \emptyset$ , we have  $\text{maxfam}_{\mathcal{R}_3,L}(\mathcal{F}, \mathcal{N}, \mathcal{C}) = \{\{(\top, h), (h, t)\}, \{(-h, \neg t), (\top, h)\}, \{(-h, \neg t), (h, t)\}\}$ . We derive  $(\top, h \wedge t) \in \text{deriv}_{\mathcal{R}_3,L}(\{(\top, h), (h, t)\})$  as follows:

$$\frac{\frac{\frac{(\top, h)}{(\top, h)} \quad \frac{\frac{(h, t)}{(\top \wedge h, t)} \quad \top \wedge h \vdash h}{(\top \wedge h, t)} \text{SI}}{(\top, t)} \text{CT}}{(\top, h \wedge t)} \text{AND}$$

with  $\mathcal{F} \vdash \top$  and  $\mathcal{C}, h \wedge t \not\vdash \perp$ . However, once we set the constraints to Jones’ violation, i.e.,  $\mathcal{C}' = \mathcal{F}$ , we obtain a single maxfam member  $\mathcal{N}' = \{(-h, \neg t), (h, t)\}$  since now  $\mathcal{C}', h \vdash \perp$  whereas  $\mathcal{C}', \neg t \not\vdash \perp$  (note that  $(h, t) \in \mathcal{N}'$  cannot be triggered by  $\mathcal{F}$ ). Given  $\mathcal{C}'$ , Jones is obliged to not tell, i.e.,  $\mathcal{K} \vdash_{\mathcal{R},L}^s \neg t$ , and is not obliged to help, i.e.,  $\mathcal{K} \not\vdash_{\mathcal{R},L}^s h$ .

First, maxfam sets (of arbitrary size) do not provide formal ways of pinpointing the reasons why some norms are inapplicable, e.g., why  $(\top, h)$  in Example 4 is inapplicable given  $\mathcal{C} = \mathcal{F}$ . Second, deriv is unsuitable for generating transparent arguments, e.g., as a certificate the derivation in Example 4 may justify that  $(\top, h \wedge t)$  is derivable, its conclusion does not explain why  $h \wedge t$  is obligatory. In fact, although in general a derivation is a justification, it is not necessarily an explanation. Our calculi address both challenges.

#### 4. Deontic Argument Calculi (DAC)

In order to generate more transparent I/O arguments, we label propositional formulae as facts  $\mathcal{L}^f$ , obligations  $\mathcal{L}^o$ , and constraints  $\mathcal{L}^c$  (Section 2). What is more, we allow for Boolean operations over the more complex meta-logical objects  $(\varphi, \psi)$  denoting norms. Operations over these higher-order syntactic objects enable undercuts that explain why

$$\begin{array}{c}
\mathbf{Ax} \vdash_{\text{LC}} \Delta^i \Rightarrow \Gamma^i, \text{ for } i \in \{f, o, c\} \quad \mathbf{Taut} \Rightarrow (\top, \top) \quad \mathbf{Detach} \quad \varphi^f, (\varphi, \psi) \Rightarrow \psi^o \quad \mathbf{TP} \quad \varphi^f \Rightarrow \varphi^o \\
\hline
\mathbf{R-C} \frac{\Delta \Rightarrow \varphi^o}{\Delta, (\neg\varphi)^c \Rightarrow} \quad \mathbf{R-N} \frac{\Delta, (\varphi, \psi) \Rightarrow}{\Delta \Rightarrow \neg(\varphi, \psi)} \quad \mathbf{L-CT}^a \frac{\varphi^f, \Delta \Rightarrow \Theta}{\varphi^o, \Delta \Rightarrow \Theta} \\
\mathbf{L-OR}^b \frac{\Delta, \varphi^f \Rightarrow \Theta \quad \Delta', \psi^f \Rightarrow \Theta}{\Delta, \Delta', (\varphi \vee \psi)^f \Rightarrow \Theta} \quad \mathbf{Cut}^c \frac{\Delta \Rightarrow \varphi \quad \varphi, \Delta' \Rightarrow \Theta}{\Delta, \Delta' \Rightarrow \Theta}
\end{array}$$

**Figure 3.** Rules for building  $\text{DAC}_{\mathcal{S}}$  (Definition 3). The upper level contains initial sequents and the lower level logical and structural rules. Side-condition (a) denotes  $\Delta \cap \mathcal{L}^n \neq \emptyset$ ; (b) denotes that if  $\mathbf{TP} \notin \mathcal{S}$ , then  $\Delta \cap \mathcal{L}^n \neq \emptyset$  and  $\Delta' \cap \mathcal{L}^n \neq \emptyset$ ; and (c) that  $\varphi \in \mathcal{L}^{io}$ .

certain norms should (not) be applied. For the present work, it suffices to consider negation only. Let  $\overline{\mathcal{L}^n} = \{\neg(\varphi, \psi) \mid (\varphi, \psi) \in \mathcal{L}^n\}$ . The *language of norms* is defined as  $\mathcal{L}^n \cup \overline{\mathcal{L}^n}$ . Furthermore, let  $\mathcal{L}^{io} = \mathcal{L}^f \cup \mathcal{L}^o \cup \mathcal{L}^c \cup \mathcal{L}^n \cup \overline{\mathcal{L}^n}$  be the full labelled I/O language. In  $\mathcal{L}^{io}$ , norms are integrated into the object-level language. We write  $\varphi$  for an arbitrary formula of  $\mathcal{L}^{io}$  and write  $\Delta^i$  to denote that  $\Delta^i \subseteq \mathcal{L}^i$  for  $i \in \{f, o, c, n\}$ .

We introduce *Deontic Argument Calculi* (DAC) for I/O logic. These calculi are *sequent-style* calculi, which are rule-based proof systems employing syntactic objects of the form  $\Delta \Rightarrow \Gamma$ , with  $\Delta, \Gamma \subseteq \mathcal{L}^{io}$  and ‘ $\Rightarrow$ ’ as a sequent arrow. We call  $\Delta \Rightarrow \Gamma$  a sequent or an argument, where  $\Delta$  denotes the reasons for  $\Gamma$  (Section 2). Furthermore,  $\Delta$  is interpreted as a regular finite set and  $\Gamma$  is restricted to at most one formula. The use of regular sets instead of multi-sets is more modular w.r.t. the base logic L. Let LC be an adequate sequent calculus for the base logic L, then, intuitively, DAC takes *labelled* versions of any LC-derivable  $\Delta \Rightarrow \Gamma$  as an initial sequent (i.e.,  $\Delta^i \Rightarrow \Gamma^i$  for each  $i \in \{f, o, c\}$ ) and contains logical- and structural rules for transforming labelled formulae of  $\mathcal{L}^{io}$  (see Figure 3).

**Definition 3.** Let DAC be the base system with the underlying logic L, containing the rules **Ax**, **Detach**, **R-C**, **R-N**, and **Cut** from Figure 3. The calculus  $\text{DAC}_{\mathcal{S}}$  extends DAC with the set of rules  $\mathcal{S} \subseteq \{\mathbf{Taut}, \mathbf{TP}, \mathbf{L-OR}, \mathbf{L-CT}\}$ , leading to 16 DAC-axiomatizations.

A  $\text{DAC}_{\mathcal{S}}$ -derivation of  $\Delta \Rightarrow \Gamma$  is a tree whose leaves are initial sequents of  $\text{DAC}_{\mathcal{S}}$ , whose root is  $\Delta \Rightarrow \Gamma$ , and whose rule-applications are instances of the rules of  $\text{DAC}_{\mathcal{S}}$ . We write  $\vdash_{\mathcal{S}} \Delta \Rightarrow \Gamma$  (resp.  $\vdash_{\mathcal{S}}^n \Delta \Rightarrow \Gamma$ ) if  $\Delta \Rightarrow \Gamma$  is  $\text{DAC}_{\mathcal{S}}$ -derivable (in at most  $n$  steps).

Since  $\text{DAC}_{\mathcal{S}}$  takes labelled LC-derivable sequents as initial sequents, the rules of LC are not part of  $\text{DAC}_{\mathcal{S}}$ . Still, LC rules can be straightforwardly shown admissible in DAC due to the presence of **Cut**. The rule **Taut** ensures that all propositional tautologies are considered as output. The rule **Detach** is an initial explanatory argument stating that the fact  $\varphi$  and the norm  $(\varphi, \psi)$  are reasons for the obligation  $\psi$ . Instead of deriving pairs from other pairs (as in deriv), we keep norms as primitive reasons from a given normative code  $\mathcal{N}$  and only modify facts, obligations, and constraints. This gives us some explanatory advantages (see **R-C** and **R-N** below). The rule **TP** corresponds to throughput. The rule **L-CT** corresponds to successive detachment, expressing that a norm may likewise be triggered by the output of some other norm (cf. Example 5). **L-OR** reflects reasoning by cases over input. The side-condition on **L-OR** is dropped for **TP**  $\in \mathcal{S}$  due to reasoning by cases with throughput. **Cut** suffices as the only structural rule.

More interesting are the rules **R-C** and **R-N**. Concerning **R-C**, think of a sequent with an empty right-hand side as an argument expressing inconsistent reasons. For instance, an argument  $\varphi^f, (\varphi, \psi), (\neg\psi)^c \Rightarrow$  explains that the fact  $\varphi$  and the norm  $(\varphi, \psi)$

(which are reasons for  $\psi$ ) are inconsistent whenever the output must be consistent with  $\neg\psi$ . What is more, whenever such an argument expresses inconsistent reasons, we know at least one of its norms is inapplicable. The rule **R-N** expresses this: from  $\varphi^f, (\varphi, \psi), (\neg\psi)^c \Rightarrow$  we obtain the defeating argument  $\varphi^f, (\neg\psi)^c \Rightarrow \neg(\varphi, \psi)$ . Hence,  $\varphi^f$  and  $(\neg\psi)^c$  are reasons for the *inapplicability* of the norm  $(\varphi, \psi)$ .  $\text{DAC}_{\mathcal{S}}$  sequents will be the building blocks for the desired argumentative characterizations (Section 5).

**Example 5** (Example 1 cont.). *The DAC-argument  $d$  (Figure 1-i), concluding that Jones should tell her neighbors she is coming to help, is derived through chaining  $(\top, h)$  and  $(h, t)$ . The following  $\text{DAC}_{\mathcal{S}}$ -derivation (left) shows this, where **L-CT**  $\in \mathcal{S}$ :*

$$\frac{\frac{\frac{}{\top^f, (\top, h) \Rightarrow h^o} \text{Detach}}{\top^f, (\top, h), (h, t) \Rightarrow t^o} \text{Detach}}{\top^f, (\top, h), (h, t) \Rightarrow t^o} \text{Cut}}{\frac{\frac{\frac{}{h^f, (h, t) \Rightarrow t^o} \text{Detach}}{h^o, (h, t) \Rightarrow t^o} \text{L-CT}}{\top^f, (\top, h), (h, t) \Rightarrow t^o} \text{Cut}}{\frac{\frac{\frac{}{\top^f, (\top, h) \Rightarrow h^o} \text{Detach}}{\top^f, (\neg h)^c, (\top, h) \Rightarrow} \text{R-C}}{\top^f, (\neg h)^c \Rightarrow \neg(\top, h)} \text{R-N}}{\top^f, (\neg h)^c \Rightarrow \neg(\top, h)} \text{R-N}}$$

Given  $\mathcal{C}' = \{\neg h^c\}$ , “why should Jones not tell, despite argument  $d$ ?” is answered by the (right) derivable argument  $b$  (Figure 1-i). The fact  $\top^f$  is omitted by a **Cut** with  $\Rightarrow \top^f$ .

**Example 6** (Example 3 cont.). *In the dilemma, Smith cannot both return the weapon and prevent the crime. So, we find  $(\top, r)$  applicable if and only if  $(\top, p)$  is inapplicable. This is expressed by arguments  $e$  and  $f$ . The  $\text{DAC}_{\mathcal{S}}$ -derivations of  $e$  and  $f$  from Figure 1-ii are obtained similarly to argument  $b$  in Example 5, using **Detach** twice, the DAC-admissible rule from LC for right conjunction introduction, **R-C**, and **R-N** consecutively.*

## 5. Argumentation and DAC-Instantiations

DAC arguments are of two types: they either give reasons for obligations, or they give reasons for why certain norms are inapplicable, i.e., defeated. The latter arguments capture the defeasibility of normative reasoning and define the interaction among arguments. We define DAC-induced *argumentation frameworks* (AFs) to model this interaction.

**Definition 4.** *Let  $\text{DAC}_{\mathcal{S}}$  be a calculus and  $\mathcal{K} = \langle \mathcal{F}, \mathcal{N}, \mathcal{C} \rangle$  a labelled knowledge base (i.e.,  $\mathcal{F} \subseteq \mathcal{L}^f, \mathcal{N} \subseteq \mathcal{L}^n$ , and  $\mathcal{C} \subseteq \mathcal{L}^c$ ). We define  $\text{AF}_{\mathcal{S}}(\mathcal{K}) = \langle \text{Arg}, \text{Att} \rangle$  as follows:*

- $\Delta \Rightarrow \Gamma \in \text{Arg}$  iff  $\Delta \Rightarrow \Gamma$  is  $\text{DAC}_{\mathcal{S}}$ -derivable,  $\Delta \subseteq \mathcal{F} \cup \mathcal{N} \cup \mathcal{C}$ , and  $\Gamma \subseteq \mathcal{L}^{io}$ ;
- $a$  defeats  $b$ , i.e.,  $(a, b) \in \text{Att}$  iff  $a = \Delta \Rightarrow \neg(\varphi, \psi)$  and  $b = \Gamma, (\varphi, \psi) \Rightarrow \Theta$ .

We write  $\text{Arg}(\Sigma)$  to denote the set of  $\text{DAC}_{\mathcal{S}}$ -arguments  $\Delta \Rightarrow \Gamma$  for which  $\Delta \subseteq \Sigma \subseteq \mathcal{L}^{io}$ .

For an  $\text{AF}_{\mathcal{S}}(\mathcal{K})$  it suffices to only consider arguments relevant to  $\mathcal{K}$ , i.e.,  $\text{Arg}(\mathcal{F} \cup \mathcal{N} \cup \mathcal{C})$ . We are interested in what combinations of arguments (*extensions*) can be collectively accepted given an AF. For our purpose, stable extensions suffice.

**Definition 5.** *Let  $\langle \text{Arg}, \text{Att} \rangle$  be an AF and let  $\mathcal{E} \subseteq \text{Arg}$ :*

- $\mathcal{E}$  defeats an argument  $a \in \text{Arg}$  if there is a  $b \in \mathcal{E}$  that defeats  $a$ , i.e.,  $(b, a) \in \text{Att}$ ;
- $\mathcal{E}$  is conflict-free if it does not defeat any of its own elements;
- $\mathcal{E}$  is stable if it is conflict-free and defeats all  $b \in \text{Arg} \setminus \mathcal{E}$ .



**Table 1.** Lemmas for  $\vdash_{\mathcal{S}}$ . Let  $\Delta^\downarrow$  and  $\varphi^\downarrow$  be the set of formulae in  $\Delta$ , resp.  $\varphi$  stripped from any labels.

Lemma :	if	then
1	$\vdash_{\mathcal{S}} \Delta, \Gamma_1^c \Rightarrow \Sigma$ and $\mathcal{C} \vdash \wedge \Gamma_1$	$\exists \Gamma_2 \subseteq \mathcal{C} : \vdash_{\mathcal{S}} \Delta, \Gamma_2^c \Rightarrow \Sigma$ and $\Gamma_2 \vdash \wedge \Gamma_1$
2	$\vdash_{\mathcal{S}}^n \Delta \Rightarrow \neg(\varphi, \psi)$	$\vdash_{\mathcal{S}}^n \Delta, (\varphi, \psi) \Rightarrow$
3		$\Delta^\downarrow \vdash \gamma^\downarrow$ , where $\Delta \subseteq \mathcal{L}^f \cup \mathcal{L}^o \cup \{(\top, \top)\}$ , $\gamma \in \mathcal{L}^f \cup \mathcal{L}^o$
4		$\vdash \gamma^\downarrow$ , where $\mathbf{TP} \notin \mathcal{S}$ , $\Delta \subseteq \mathcal{L}^f \cup \{(\top, \top)\}$ , $\gamma \in \mathcal{L}^o$
5	$\vdash_{\mathcal{S}}^n \Delta \Rightarrow$	$\vdash_{\mathcal{S}}^n \Delta \setminus \mathcal{L}^c \Rightarrow \varphi^o$ s.t. $\varphi \vdash \neg \wedge (\Delta \cap \mathcal{L}^c)^\downarrow$ , where $\neg \wedge \emptyset =_{\text{df}} \perp$

Let **Stable** be the set of stable extensions of AF. We define sceptic ( $s$ ), sceptic\* ( $s^*$ ), and credulous ( $c$ ) nonmonotonic inference as follows:

- $\text{AF} \sim_{\text{stable}}^s \varphi$  iff for each  $\mathcal{E} \in \text{Stable}$ , there is an  $a \in \mathcal{E}$  concluding  $\varphi$ ;
- $\text{AF} \sim_{\text{stable}}^{s^*} \varphi$  iff there is an  $a \in \bigcap \text{Stable}$  concluding  $\varphi$ ;
- $\text{AF} \sim_{\text{stable}}^c \varphi$  iff there is a  $\mathcal{E} \in \text{Stable}$  s.t. there is an  $a \in \mathcal{E}$  concluding  $\varphi$ .

The use of DAC-arguments introduces nuances in sceptic inference: e.g., the distinction between  $s$  and  $s^*$  corresponds to the discussion of floating conclusions in Section 2.

**Example 7** (Example 3 cont.). *Smith is in a dilemma of conflicting duties. The AF of Figure 1-ii represents this conflict, where  $\text{Arg} = \{a, b, c_1, c_2, d, e\}$  and  $\text{Att} = \{(e, a), (e, c_1), (e, d), (d, e), (d, c_2), (d, b)\}$ . It has two stable extensions  $\{a, c_1, d\}$  and  $\{b, c_2, e\}$ , defending the views that Smith ought to return the weapon, resp. prevent the crime. Hence,  $\text{AF} \sim_{\text{stable}}^c r^o, p^o$ , whereas  $\text{AF} \not\sim_{\text{stable}}^c (r \wedge p)^o$ ,  $\text{AF} \not\sim_{\text{stable}}^s r^o$ , and  $\text{AF} \not\sim_{\text{stable}}^s p^o$ . For the floating conclusion  $(r \vee p)^o$  we have  $\text{AF} \sim_{\text{stable}}^s (r \vee p)^o$  but  $\text{AF} \not\sim_{\text{stable}}^{s^*} (r \vee p)^o$ . (The AF of Example 2 in Figure 1-i has one stable extension  $\{a, b\}$ , and so  $\text{AF} \sim_{\text{stable}}^{s, s^*, c} (\neg t)^o$ .)*

To illustrate the utility of our approach, we consider the notion of related admissibility [16]. An extension  $\mathcal{E}$  is *admissible* if it is conflict-free and  $\mathcal{E}$  defeats all arguments defeating some  $a \in \mathcal{E}$ . An argument  $a$  *defends*  $b$  iff  $a = b$ , or there is a  $c$  s.t.  $a$  defeats  $c$  and  $c$  defeats  $b$ , or there is a  $c$  s.t.  $a$  defends  $c$  and  $c$  defends  $b$ . A set  $\mathcal{E}_a \subseteq \text{Arg}$  is *related admissible with topic*  $a$  iff  $a \in \mathcal{E}_a$ , for all  $b \in \mathcal{E}_a$ ,  $b$  defends  $a$ , and  $\mathcal{E}_a$  is admissible. Thus, a related admissible set  $\mathcal{E}_a$  identifies the relevant arguments that justify the acceptability of  $a$ . Let  $\mathcal{E}^+ = \{a \in \text{Arg} \mid \mathcal{E} \text{ defeats } a\}$  and  $\mathcal{E}^- = \{a \in \text{Arg} \mid a \text{ defeats some } b \in \mathcal{E}\}$ . In Example 3, the answer to “why is Smith obliged to prevent crime ( $b$ )?” is given by the related admissible set  $\mathcal{E}_b = \{b, e\}$  where  $\mathcal{E}_b^- = \{d\}$  and  $\{d\}^- \cap \mathcal{E}_b = \{e\}$  explain that the only counterargument to  $b$  is  $d$  which is defeated by  $e$  expressing that the norm  $(\top, r)$  used in  $d$  is inapplicable given the reasons  $(\top, p)$  and  $\neg(r \wedge p)^c$  offered in  $e$ . Hence, using only undercuts enables a more refined analysis of the *relevant* norms explaining the (non-)acceptability of certain arguments and obligations. The DAC approach is therefore more precise compared to using maximal consistent families of norms in traditional I/O.

## 6. Metatheory: Soundness and Completeness

We demonstrate two soundness and completeness results: First, we prove adequacy between I/O proof systems and DAC (Theorem 1). Second, we prove adequacy between constrained I/O logics and DAC-based argumentation frameworks (Theorem 2). We pro-

**Table 2.** Correspondence between  $\text{deriv}_{\mathcal{R},L}$  rules and  $\text{DAC}_{\mathcal{S}}$  rules with the underlying logic  $L$ . For instance,  $\{\text{ID}, \text{OR}\} \subseteq \mathcal{R}$  iff  $\{\text{TP}, \text{L-OR}\} \subseteq \mathcal{S}$ . The first column represents the minimal sets the systems must contain.

Rules of $\text{deriv}_{\mathcal{R},L}$	$\{\text{WO}, \text{AND}, \text{SI}\}$	$\text{T}$	$\text{ID}$	$\text{CT}$	$\text{OR}$
Rules of $\text{DAC}_{\mathcal{S}}$	$\{\text{Ax}, \text{Detach}, \text{R-C}, \text{R-N}, \text{Cut}\}$	$\text{Taut}$	$\text{TP}$	$\text{L-CT}$	$\text{L-OR}$

vide explicit proofs of the main results. Table 1 lists several technical lemmas whose proofs are omitted: Lemma 1 follows by the compactness of  $L$ , while Lemmas 2 to 5 are proven by a straightforward induction on the length of the derivation.

We first show adequacy between  $\text{deriv}$  and  $\text{DAC}$ . Both systems are modular and correspondence between the rules of these systems is defined in Table 2. In referring to  $\text{deriv}_{\mathcal{R},L}$  and  $\text{DAC}_{\mathcal{S}}$  we assume this correspondence. We state the two directions of Theorem 1 separately, we prove Lemma 7, and omit the similar proof of Lemma 6.

**Lemma 6.** *Let  $\Theta \subseteq \mathcal{L}^n$ , If  $\vdash_{\mathcal{S}} \Delta^f, \Theta \Rightarrow \varphi^o$ , then  $\varphi \in \text{deriv}_{\mathcal{R},L}(\Delta, \Theta)$ .*

**Lemma 7.** *If  $(\varphi, \psi) \in \text{deriv}_{\mathcal{R},L}(\Theta)$ , then  $\vdash_{\mathcal{S}} \varphi^f, \Theta \Rightarrow \psi^o$ .*

*Proof.* By induction on the length of the  $\text{deriv}_{\mathcal{R},L}$ -derivation of  $(\varphi, \psi)$ . *Base case.* Case  $\{(\varphi, \psi)\} = \Theta$ . By **Detach**,  $\vdash_{\mathcal{S}} \varphi^f, (\varphi, \psi) \Rightarrow \psi^o$ . Case  $(\top, \top)$  is derived by **T** with  $\Theta = \emptyset$ . By **Detach**,  $\vdash_{\mathcal{S}} \top^f, (\top, \top) \Rightarrow \top^o$  and by **Taut**,  $\vdash_{\mathcal{S}} \Rightarrow (\top, \top)$ . By **Cut**,  $\vdash_{\mathcal{S}} \top^f \Rightarrow \top^o$ . Case  $(\varphi, \varphi)$  is derived by **ID** with  $\Theta = \emptyset$ . By **TP**,  $\varphi^f \Rightarrow \varphi^o$ .

*Inductive step.* To illustrate, we consider the case of **CT**. The other cases are similar or straightforward. Suppose that  $(\varphi, \psi)$  is derived from  $(\varphi, \sigma) \in \text{deriv}_{\mathcal{R},L}(\Theta_1)$  and  $(\varphi \wedge \sigma, \psi) \in \text{deriv}_{\mathcal{R},L}(\Theta_2)$  by **CT**, where  $\Theta = \Theta_1 \cup \Theta_2$ . By the IH,  $\vdash_{\mathcal{S}} \varphi^f, \Theta_1 \Rightarrow \sigma^o$  and  $\vdash_{\mathcal{S}} (\varphi \wedge \sigma)^f, \Theta_2 \Rightarrow \psi^o$ . By **R $\wedge$ 2**,  $\varphi, \sigma \vdash \varphi \wedge \sigma$ . By **Ax**,  $\vdash_{\mathcal{S}} \varphi^f, \sigma^f \Rightarrow (\varphi \wedge \sigma)^f$  and by **Cut**,  $\varphi^f, \sigma^f, \Theta_2 \Rightarrow \psi^o$ . Then, if  $\emptyset \neq \Theta_2$ , by **L-CT**,  $\vdash_{\mathcal{S}} \varphi^f, \sigma^o, \Theta_2 \Rightarrow \psi^o$  and by **Cut**,  $\vdash_{\mathcal{S}} \varphi^f, \Theta \Rightarrow \psi^o$ . Else,  $\Theta_2 = \emptyset$  (and hence  $\Theta = \Theta_1$ ). We consider: (i) **TP**  $\in \mathcal{S}$  and (ii) **TP**  $\notin \mathcal{S}$ . Ad (i). By Lemma 3.1,  $\varphi, \sigma \vdash \psi$  and by **Ax**,  $\vdash_{\mathcal{S}} \varphi^o, \sigma^o \Rightarrow \psi^o$ . By **TP**,  $\vdash_{\mathcal{S}} \varphi^f \Rightarrow \varphi^o$  and by twice **Cut**,  $\vdash_{\mathcal{S}} \varphi^f, \Theta \Rightarrow \psi^o$ . Ad (ii). By Lemma 3.2,  $\vdash \psi$  and so  $\sigma \vdash \psi$ . By **Ax**,  $\vdash_{\mathcal{S}} \sigma^o \Rightarrow \psi^o$ . By **Cut**,  $\vdash_{\mathcal{S}} \varphi^f, \Theta \Rightarrow \psi^o$ .  $\square$

**Theorem 1.** *Let  $\Delta \subseteq \mathcal{L}$ ,  $\psi \in \mathcal{L}$ , and  $\Theta \subseteq \mathcal{L}^n$ . Then,  $\vdash_{\mathcal{S}} \Delta^f, \Theta \Rightarrow \psi^o$  iff  $\psi \in \text{deriv}_{\mathcal{R},L}(\Delta, \Theta)$ .*

*Proof.*  $(\Rightarrow)$  This is Lemma 6.  $(\Leftarrow)$  Suppose  $\psi \in \text{deriv}_{\mathcal{R},L}(\Delta, \Theta)$ . So,  $(\wedge \Delta, \psi) \in \text{deriv}_{\mathcal{R},L}(\Theta)$ . By Lemma 7,  $\vdash_{\mathcal{S}} (\wedge \Delta)^f, \Theta \Rightarrow \psi^o$ . Since  $\Delta \vdash \wedge \Delta$ ,  $\vdash_{\mathcal{S}} \Delta^f \Rightarrow (\wedge \Delta)^f$  by **Ax**. By **Cut**,  $\vdash_{\mathcal{S}} \Delta^f, \Theta \Rightarrow \psi^o$ .  $\square$

We now prove our second adequacy result concerning constrained I/O logics and  $\text{DAC}$ -instantiated argumentation frameworks.

**Theorem 2.** *Let  $\mathcal{K} = \langle \mathcal{F}, \mathcal{N}, \mathcal{C} \rangle$  be a knowledge base. Let  $\mathcal{R}$  be a set of deriv-rules and  $\mathcal{S}$  the set of corresponding  $\text{DAC}$ -rules (Table 2). Let  $\text{AF} = \text{AF}_{\mathcal{S}}(\mathcal{K}) = \langle \text{Arg}, \text{Att} \rangle$ .*

1. *If  $\mathcal{N}' \in \text{maxfam}_{\mathcal{R},L}(\mathcal{K})$  then  $\text{Arg}(\mathcal{F}^f \cup \mathcal{N}' \cup \mathcal{C}^c)$  is stable in  $\text{AF}$ .*
2. *If  $\mathcal{A}$  is stable in  $\text{AF}$  then there is a  $\mathcal{N}' \subseteq \mathcal{N}$  such that  $\mathcal{N}' \in \text{maxfam}_{\mathcal{R},L}(\mathcal{K})$  for which  $\mathcal{A} = \text{Arg}(\mathcal{F}^f \cup \mathcal{N}' \cup \mathcal{C}^c)$ .*

*Proof.* (1) Let  $\mathcal{N}' \in \text{maxfam}_{\mathcal{R},L}(\mathcal{K})$  and  $\mathcal{A} = \text{Arg}(\mathcal{F}^f \cup \mathcal{N}' \cup \mathcal{C}^c)$ . For conflict-freeness assume towards a contradiction that there are  $a = \Delta^f, \Theta, \Gamma^c \Rightarrow \neg(\varphi, \psi) \in \mathcal{A}$

(where  $\Theta \subseteq \mathcal{N}'$ ) and  $b = \Omega, (\varphi, \psi) \Rightarrow \Sigma \in \mathcal{A}$  such that  $a$  attacks  $b$ . By Lemma 2 and since  $(\varphi, \psi) \in \mathcal{N}'$ , we have,  $\Delta^f, \Theta, \Gamma^c, (\varphi, \psi) \Rightarrow \in \mathcal{A}$ . By Lemma 5,  $\Delta^f, \Theta, (\varphi, \psi) \Rightarrow \sigma^o \in \mathcal{A}$  for some  $\sigma$  for which  $\sigma \vdash \neg \wedge \Gamma$ . By Theorem 1,  $\sigma \in \text{deriv}_{\mathcal{R}, \mathcal{L}}(\Delta, \Theta \cup \{(\varphi, \psi)\})$ , which contradicts the  $\mathcal{C}$ -consistency of  $\mathcal{N}'$ .

For  $\mathcal{A}$  defeats all  $b \in \text{Arg} \setminus \mathcal{A}$  let  $a = \Delta_1^f, \Theta_1, \Gamma_1^c \Rightarrow \Sigma \in \text{Arg} \setminus \mathcal{A}$ , where  $\Theta_1 \subseteq \mathcal{L}^n$ . So, there is a  $(\varphi, \psi) \in \Theta_1 \setminus \mathcal{N}'$ . By the maximal consistency of  $\mathcal{N}'$ ,  $\mathcal{N}' \cup \{(\varphi, \psi)\}$  is inconsistent with  $\mathcal{C}$ . So, there is a  $\theta \in \text{deriv}_{\mathcal{R}, \mathcal{L}}(\Delta_2, \Theta_2)$  for some  $\Delta_2 \subseteq \mathcal{F}$  and  $\Theta_2 \subseteq \mathcal{N}' \cup \{(\varphi, \psi)\}$  such that  $\mathcal{C} \vdash \neg \theta$ . By Theorem 1,  $\Delta_2, \Theta_2 \Rightarrow \theta^o \in \text{Arg}$ . Note that  $(\varphi, \psi) \in \Theta_2$  since otherwise  $\Theta_2 \subseteq \mathcal{N}'$  in contradiction to the consistency of  $\mathcal{N}'$ . By **R-C** and **R-N**,  $\vdash_{\mathcal{S}} \Delta_2, \Theta_2 \setminus \{(\varphi, \psi)\}, (\neg \theta)^c \Rightarrow \neg(\varphi, \psi)$ . By Lemma 1,  $b = \Delta_2, \Theta_2 \setminus \{(\varphi, \psi)\}, \Gamma_2^c \Rightarrow \neg(\varphi, \psi) \in \text{Arg}$  for some  $\Gamma_2 \subseteq \mathcal{C}$  for which  $\Gamma_2 \vdash \neg \theta$ . Note,  $b \in \mathcal{A}$  and  $b$  attacks  $a$ .

(2) Let  $\mathcal{A}$  be a stable extension of  $AF(\mathcal{K})$ . Let  $\mathcal{N}' = \{(\varphi, \psi) \in \mathcal{N} \mid \neg \exists a = \Delta \Rightarrow \neg(\varphi, \psi) \in \mathcal{A}\}$ . We first show that  $\mathcal{A} = \text{Arg}(\mathcal{F}^f \cup \mathcal{N}' \cup \mathcal{C}^c)$ :

“( $\supseteq$ )” Let  $a \in \text{Arg}(\mathcal{F}^f \cup \mathcal{N}' \cup \mathcal{C}^c)$ . By the definition of  $\mathcal{N}'$  there is no  $b \in \mathcal{A}$  that attacks  $a$  and since  $a \in \text{Arg}$  and by the stability of  $\mathcal{A}$ ,  $a \in \mathcal{A}$ . “( $\subseteq$ )” Let  $a \in \text{Arg} \setminus \text{Arg}(\mathcal{F}^f \cup \mathcal{N}' \cup \mathcal{C}^c)$  with  $a = \Delta \Rightarrow \Gamma$ . So, there is a  $(\varphi, \psi) \in \Delta$  for which there is a  $b \in \mathcal{A}$  with  $b = \Theta \Rightarrow \neg(\varphi, \psi)$ . So  $b$  attacks  $a$  and by the stability of  $\mathcal{A}$ ,  $a \notin \mathcal{A}$ .

We now show that  $\mathcal{N}' \in \text{maxfam}_{\mathcal{R}, \mathcal{L}}(\mathcal{K})$ . Assume for a contradiction that  $\mathcal{N}'$  is inconsistent with  $\mathcal{C}$ . So, there is a  $\theta \in \text{deriv}_{\mathcal{R}, \mathcal{L}}(\Delta, \Theta)$  for some  $\Delta \subseteq \mathcal{F}$  and  $\Theta \subseteq \mathcal{N}'$  for which  $\mathcal{C} \vdash \neg \theta$ . By Theorem 1,  $a = \Delta^f, \Theta \Rightarrow \theta^o \in \mathcal{A}$ . Assume first that  $\Theta = \emptyset$ .

If **TP**  $\notin \mathcal{S}$ , by Lemma 3.3,  $\vdash \theta$  and thus  $\mathcal{C}$  is inconsistent which is a contradiction. Thus,  $\Theta \neq \emptyset$ . If **TP**  $\in \mathcal{S}$  then, by Lemma 3.2,  $\Delta \vdash \theta$ . But then  $\mathcal{F} \cup \mathcal{C}$  is inconsistent, a contradiction and so  $\Theta \neq \emptyset$ . So, in both cases  $\Theta \neq \emptyset$ .

Let  $(\varphi, \psi) \in \Theta$ . By **R-N** and **R-C**,  $\vdash_{\mathcal{S}} \Delta, \Theta \setminus \{(\varphi, \psi)\}, (\neg \theta)^c \Rightarrow \neg(\varphi, \psi)$ . By Lemma 1, there is a  $\Gamma \subseteq \mathcal{C}$  for which  $b = \Delta, \Theta \setminus \{(\varphi, \psi)\}, \Gamma^c \Rightarrow \neg(\varphi, \psi) \in \mathcal{A}$ . Since  $b$  attacks  $a$ , this contradicts conflict-freeness of  $\mathcal{A}$ , which shows  $\mathcal{N}'$  is consistent with  $\mathcal{C}$ .

Assume for a contradiction that there is a  $(\varphi, \psi) \in \mathcal{N} \setminus \mathcal{N}'$  such that  $\mathcal{N}' \cup \{(\varphi, \psi)\}$  is consistent with  $\mathcal{C}$  (i.e.,  $\mathcal{N}'$  is not maximal). By the definition of  $\mathcal{N}'$ , there is a  $b = \Delta^f, \Theta, \Gamma^c \Rightarrow \neg(\varphi, \psi) \in \mathcal{A}$ . By Lemma 2,  $\vdash_{\mathcal{S}} \Delta^f, \Theta, (\varphi, \psi), \Gamma^c \Rightarrow$ . By Lemma 5,  $\vdash_{\mathcal{S}} \Delta^f, \Theta, (\varphi, \psi) \Rightarrow \sigma^o$  such that  $\sigma \vdash \neg \wedge \Gamma$ . By Theorem 1,  $\sigma \in \text{deriv}_{\mathcal{R}, \mathcal{L}}(\Delta, \Theta \cup \{(\varphi, \psi)\})$  which shows that  $\mathcal{N}' \cup \{(\varphi, \psi)\}$  is inconsistent with  $\mathcal{C}$  (note that  $\Gamma \subseteq \mathcal{C}$ ). This completes our proof since it shows that  $\mathcal{N}' \in \text{maxfam}_{\mathcal{R}, \mathcal{L}}(\mathcal{K})$ .  $\square$

**Corollary 1.** Let  $\mathcal{K}$  be a knowledge base,  $\mathcal{R}$  a set of deriv-rules, and  $\mathcal{S}$  a set of corresponding DXC-rules (Table 2). For  $i \in \{s, c\}$ ,  $AF_{\mathcal{S}}(\mathcal{K}) \sim_{stable}^i \varphi$  iff  $\mathcal{K} \sim_{\mathcal{R}, \mathcal{L}}^i \varphi$ .

## 7. Related Work and Conclusion

In [14], a sequent-style system for monotonic I/O logics without constraints is presented. It utilizes a correspondence between I/O and conditional logics. In [15] proof systems for constrained I/O logic are developed, where modalities for ‘input’ and ‘output’ allow for meta-reasoning in the object language. DAC uses labels instead of modalities and additionally allows for meta-reasoning about the (in)applicability of norms. In [9], sequent argumentation is used for defeasible reasoning with deontic logic. Norms are modelled with material implications which allows for less fine-tuning of norms than in DAC.

In [5,6,18] argumentative characterizations of normative systems employing priority orderings are studied. Their language is restricted to literals only, whereas our approach adopts a full propositional language. In [6,18] arguments consist only of (sets of) norms. In future work, we aim to incorporate priority and preference reasoning in the more transparent context of DAC. Moreover, the I/O formalism has other applications including reasoning with consistency checks, permissions, and constitutive norms [8,17]. In particular, we aim to exploit the internalization of meta-reasoning in DAC to characterize various types of permission [17], for instance, negative permissions as defined in terms of the absence of applicable norms to the contrary.

An alternative approach to model reasoning with norms is to instantiate ASPIC<sup>+</sup> [2] with conditionals representing norms and a defeasible modus ponens rule. This approach leads to a “greedier” style of reasoning than our approach. Consider  $\mathcal{F} = \emptyset$  and  $\mathcal{N} = \{(\top, p), (p, q), (\top, \neg q)\}$ . An ASPIC<sup>+</sup>-based approach yields the obligation to  $p$  with stable semantics since the argument for  $p$  from  $(\top, p)$  is unchallenged. In contrast, our approach generates the argument  $(\top, \neg q), (p, q) \Rightarrow \neg(\top, p)$  concluding the inapplicability of  $(\top, p)$ . The latter is in line with the I/O approach to normative reasoning.

We illustrated our approach with the notion of related admissibility [16]. For future work, we will investigate other argumentative approaches to explanation and how these can be used in the context of DAC, e.g., explicit reasoning about the inapplicability of norms in DAC can be harnessed to explain the non-acceptability of arguments [19].

Last, explanations typically occur in dialogues, through an exchange of reasons, questions, and arguments [20]. Consequently, explanations are often tailored to the background of the explainee. We will adopt our approach to dialogue models in future work.

In conclusion, in normative reasoning contexts one is not just interested in whether a specific obligation holds, but also in why it holds despite other norms to the contrary. To address this challenge, we developed Deontic Argument Calculi (DAC) which are rule-based proof systems that use labels to facilitate transparency and incorporate meta-normative reasoning with norms into the object language.

## References

- [1] Gabbay D, Horty J, Parent X, van der Meyden R, van der Torre L. Handbook of Deontic Logic and Normative Systems. vol. 1. United Kingdom: College Publications; 2013.
- [2] Gabbay D, Giacomini M, Simari GR, Thimm M. Handbook of Formal Argumentation. vol. 2. United Kingdom: College Publications; 2021.
- [3] Makinson D, van der Torre L. Constraints for Input/Output Logics. *Journal of Philosophical Logic*. 2001;30(2):155-85.
- [4] Beirlaen M, Straßer C, Heynink J. Structured argumentation with prioritized conditional obligations and permissions. *Journal of Logic and Computation*. 2018;29(2):187-214.
- [5] Governatori G, Rotolo A, Riveret R. A deontic argumentation framework based on deontic defeasible logic. In: *Proceedings of PRIMA 2018*. Springer; 2018. p. 484-92.
- [6] Liao B, Oren N, van der Torre L, Villata S. Prioritized norms in formal argumentation. *Journal of Logic and Computation*. 2018;29(2):215-40.
- [7] Peirera C, Tettamanzi AG, Villata S, Liao B, Malerba A, Rotolo A, et al. Handling norms in multi-agent system by means of formal argumentation. *IfCoLog*. 2017;4(9):1-35.
- [8] Pigozzi G, van der Torre L. Arguing about constitutive and regulative norms. *Journal of Applied Non-Classical Logics*. 2018;28(2-3):189-217.
- [9] Straßer C, Arieli O. Normative reasoning by sequent-based argumentation. *Journal of Logic and Computation*. 2015 07;29(3):387-415.

- [10] Parent X, van der Torre L. Introduction to deontic logic and normative systems. College Publications; 2018.
- [11] Brunero J. Reasons, Evidence, and Explanations. Oxford Handbooks Online. 2018 Jul.
- [12] Nair S, Horty J. The Logic of Reasons. Oxford Handbooks Online. 2018 Jul.
- [13] Horty JF. Reasons as defaults. Oxford University Press, USA; 2012.
- [14] Lellmann B. From Input/Output Logics to Conditional Logics via Sequents – with Provers. In: Automated Reasoning with Analytic Tableaux and Related Methods. Cham: Springer; 2021. p. 147-64.
- [15] Straßer C, Beirlaen M, Van De Putte F. Adaptive logic characterizations of input/output logic. *Studia Logica*. 2016;104(5):869-916.
- [16] Fan X, Toni F. On computing explanations in argumentation. In: AAAI; 2015. p. 1496-502.
- [17] Tosatto SC, Boella G, van der Torre L, Villata S. Abstract normative systems: Semantics and proof theory. In: Proceedings of KR12; 2012. .
- [18] Straßer C, Pardo P. Prioritized Defaults and Formal Argumentation. In: Proceedings of 14th International Conference of Deontic Logic and Normative Systems. College Publications; 2021. p. 427-46.
- [19] Borg A, Bex F. A Basic Framework for Explanations in Argumentation. *IEEE Int Systems*. 2021:25-35.
- [20] Walton D. A dialogue system specification for explanation. *Synthese*. 2010 Apr;182(3):349–374.