

# Deontic Defeasible Reasoning in Legal Interpretation

## Two Options for Modelling Interpretive Arguments

Antonino Rotolo\*  
CIRSFID, University of Bologna  
Italy

Guido Governatori†  
NICTA  
Australia

Giovanni Sartor  
EUI and University of Bologna  
Italy

### ABSTRACT

This paper offers a new logical machinery for reasoning about interpretive canons. We identify some options for modelling reasoning about interpretations and show that interpretive argumentation has a distinctive structure where the claim that a legal text ought or may be interpreted in a certain way can be supported or attacked by arguments, whose conflicts may have to be assessed according to further arguments.

### Categories and Subject Descriptors

I.2.4 [Artificial Intelligence]: Knowledge Representation Formalisms and Methods

### General Terms

Theory

### Keywords

Defeasible Logic, Legal Interpretation, Argumentation

## 1. INTRODUCTION

Legal doctrine and judicial practice distinguish among a number of canons for interpreting legal statutes, i.e., different rules that are employed in legal systems as patterns for constructing arguments aimed at justifying certain interpretations, while attacking other interpretations. [11], summarising the outcomes of a vast study on statutory interpretation, involving scholars from many different legal systems, distinguishes eleven types of arguments. A different list of interpretive arguments was developed by [15] and identifies fourteen types of arguments.

\*Supported by the Unibo FARB 2012 project *Mortality Salience, Legal and Social Compliance, and Economic Behaviour: Theoretical Models and Experimental Methods*.

†NICTA is funded by the Australian Government and the Australian Research Council through the ICT Centre of Excellence program.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [Permissions@acm.org](mailto:Permissions@acm.org).

ICAIL'15 June 08 - 12, 2015, San Diego, CA, USA

Copyright 2015 ACM 978-1-4503-3522-5/15/06 ...\$15.00.

<http://dx.doi.org/10.1145/2746090.2746100>.

This paper proposes a new logical machinery for modelling reasoning about interpretive canons and thus for justifying the choice of a certain canon and the resulting legal outcome over competing interpretations (see [12, 10]). [14] argued that an interpretive canon for statutory law can be expressed as follows: if provision  $n$  occurs in document  $D$ ,  $n$  has a setting of  $S$ , and  $n$  would fit this setting of  $S$  by having interpretation  $a$ , then,  $n$  ought to be interpreted as  $a$ . For instance, the ordinary language canon has the following structure:<sup>1</sup> if provision  $n$ , stating that “Killing a man is punishable by no less than 21 years in prison”, occurs in document  $D =$  Penal code,  $n$  has a setting of ordinary language, and  $n$  would fit this setting of ordinary language by having interpretation  $a =$  “Killing an adult male person is punishable by no less than 21 years in prison”, then  $n$  ought to be interpreted as  $a$ .

In this paper we basically accept this research line and work according to the following intuitions.

INTUITION 1 (REASONING AND CANONS). *We analyse the logical structure of interpretive arguments (in the sense of [11]) using a rule-based logical system. In particular, interpretation canons are represented by defeasible rules, where*

- antecedent conditions of interpretation rules can be of any type (assertions, obligations, etc.), including the fact that another canon is refuted or that another legal provision ought to be interpreted in a certain way;
- the conclusion of interpretation rules is an interpretive act leading to an interpretation of a certain provision  $n$  and thus to a sentence which expresses the result of such an interpretation and paraphrases  $n$  [6]. If  $n$  and  $n'$  are legal provisions, the following is an example of interpretation rule regarding  $n'$ :

IF  
 $n$  ought to be interpreted literally as  $a$ , AND  
 $n$  is related with  $n'$ , AND  
 $a$  entails  $a'$ ,  
THEN  
 $n'$  is interpreted by coherence as  $a'$ .

*We will use these rules to devise a reasoning machinery that mirrors legal reasoning about interpretive canons. The resulting rule-based system is in line with the basic ideas inspiring the argumentation system by [12].*

<sup>1</sup>This argument supports the option that a provision be interpreted according to the meaning a native speaker of a given language would ascribe to it.

Notice that the above intuition distinguishes the interpretive act from the result of the interpretation:

INTUITION 2 (A- AND O-INTERPRETATION). *We assume the distinction between interpretation as activity and as outcome [13, p. 117] (cf. [15, p. 39]):*

- *interpretation as activity (A-interpretation) (literal or from ordinary language, by coherence, etc.) views any argumentative canon as a means through which a certain meaning is ascribed to a legal provision, and*
- *interpretation as outcome (O-interpretation) is precisely the meaning obtained through a certain interpretive act and ascribed to the provision.*

*The distinction between interpretation as activity and as outcome is well known in continental legal theory, and it was introduced precisely to capture cases where, e.g., one has legal reasons to prefer a certain interpretive canon over others even though all considered canons support the same interpretive outcome. In other words, an interpretive act  $I$  of  $n$  as a (A-interpretation of  $n$ ) is a way to bring about that a (O-interpretation of  $n$ ) is the case. For example, in Intuition 1, the A-interpretation of  $n'$  is the act interpreting  $n'$  by coherence, while the resulting O-interpretation from that act is  $a'$ , i.e., a sentence expressing the meaning attributed to  $n'$  through the interpretation by coherence.*

Since different competing canons can be employed, different conflicting rules can be accordingly applied for interpreting statutes. Interpretation rules are thus defeasible. As argued in [14], some priority criteria should be applied to interpretation rules [1]. Such criteria impose preference relations over conflicting interpretive acts and outcomes. In other words, to address interpretive conflicts, we need to assume that one of the conflicting arguments is stronger than its competitors. Some legal traditions provide indeed general criteria for addressing conflicts of arguments on the basis of their priorities: for instance, several continental legal systems explicitly state that literal interpretation ought to be preferred, or that an argument concerning constitutional values ought to prevail over a historical argument (e.g., an argument based on the intent of the historical legislator).

However, ranking among interpretive acts and canons can be applied also when such acts are not in conflict. Suppose, for example, that provision  $n$  can be interpreted as  $a$  by adopting an argument by analogy and one from substantive reasons;<sup>2</sup> if  $n$  is a provision of criminal law (but analogy is admissible whenever it favours the defendant), then the argument from substantive reasons ought to be preferred, even though both lead to read  $n$  as  $a$ .

INTUITION 3 (PREFERENCES OVER INTERPRETATIONS). *A standard priority relation [2] over interpretation rules can be introduced to handle and solve conflicts between different*

<sup>2</sup>An argument from substantive reasons states that, if there is some goal that can be considered to be fundamentally important to the legal system, and if the goal can be promoted by one rather than another interpretation of the statutory provision, then the provision should be interpreted in accord with the goal.

*interpretation rules. Consider the following example:*

**Rule1**

IF

*$n$  ought to be interpreted literally as  $a$ , AND  
 $n$  is related with  $n'$ , AND  
 $a$  entails  $a'$ ,*

THEN

*$n'$  is interpreted by coherence as  $a'$*

**Rule2**

IF

*$n''$  ought to be interpreted literally as  $\neg a$ , AND  
 $n$  is related with  $n''$ , AND  
 $\neg a$  entails  $\neg a'$ ,*

THEN

*$n'$  is interpreted by coherence as  $\neg a'$ .*

*Here, we can handle the conflict by stating that **Rule1** > **Rule2** (or vice versa).*

*Ranking among interpretive acts can be applied also when such acts are not in conflict. We will thus introduce an operator  $\otimes$  that can be used to make explicit in single rules this idea. For instance,*

IF

*$n$  ought to be interpreted literally as  $a$ , AND  
 $n$  is related with  $n'$*

THEN

*$n'$  is interpreted by coherence as  $a' \otimes$   
 $\otimes n'$  is interpreted by analogy as  $a'$*

*means that the most preferred interpretation resulting in  $a$  is the one by coherence, but, if this is refuted, the second option is the interpretation by analogy. This does not require to only derive one interpretation resulting in  $a$  (other rules could first support interpretation by analogy of  $n$ ).*

Following some doctrinal and judicial practice, [14] argued that interpretive canons are defeasible rules licensing deontic interpretive claims, namely, the claim that a certain expression in a statute ought, ought not, may or may not be interpreted in a certain way. For example, art. 12 in the general provisions of the Italian civil code states that the literal interpretation of statutes ought to be preferred and this option is nothing but an interpretive prescription. Here, we follow this intuition with some adjustments.

INTUITION 4 (OBLIGATORY INTERPRETATIONS). *An interpretation can be admissible or obligatory. In the case of A-interpretations, for instance, an interpretive act  $I$  of  $n$  (A-interpretation of  $n$ ) is admissible, if it is provable using a defeasible interpretation rule; it is obligatory, if this interpretation of  $n$  is the only one admissible. Similarly for O-interpretations. Indeed, consider the general provisions of the Italian civil code, which state at art. 12 that literal interpretation  $I_{lit}$  ought to be preferred: this would support that such interpretation is obligatory, unless another interpretation prevails. We have two options here:*

- *other conflicting interpretations can be derived, thus requiring to check if literal interpretation overrides the other options; if it does not, then the interpretation at stake is not even admissible;*
- *other non-conflicting interpretations can be provable; if they are, the interpretation at stake is only admissible, otherwise, it is obligatory.*

On the basis of the above intuitions, we will offer two options for modelling reasoning about interpretations: a defeasible logic for reasoning about the interpretation of abstract, non-analysed provisions and of structured provisions.

**INTUITION 5 (ABSTRACT OR STRUCTURED PROVISIONS).** A provision  $n$  is abstract if it is taken in its sentential entirety for interpretive purposes, i.e., as a non-analysed sentence without considering its internal (logical) structure. In our paper, Option 1 amounts to interpreting  $n$  by ascribing to  $n$ , intended as an abstract provision, a sentential meaning that can be expressed by another sentence paraphrasing this provision as a whole.

Rather, a provision  $n'$  is logically structured if it corresponds to a linguistic sentence having the structure of a rule  $a_1, \dots, a_n \Rightarrow b$ : this means that  $n'$  is semi-interpreted provision, since expressing the logical structure of  $n'$  requires an interpretive effort on the original textual version of  $n'$ . In our paper, Option 2 amounts to interpreting  $n'$  by considering the components  $a_1, \dots, a_n, b$  of  $n'$  and ascribing to them a meaning as already explained above.

The above intuitions are implemented adjusting the framework in [7], which is a Modal Defeasible Logic [8] extended with the operator  $\otimes$ . The logic is a significant extension of standard Defeasible Logic [2], which however preserves linear computational complexity (like standard Defeasible Logic).

The layout of the paper is as follows. Section 2 offers a Modal Defeasible Logic for Option 1 (see Intuition 5 above), i.e., a machinery for reasoning about the interpretation of abstract provisions: Section 2.1 presents the formal language, Section 2.2 introduces some preliminary concepts, while Sections 2.3 and 2.4 offer a complete proof theory for A-interpretations and O-interpretations (see Intuition 2). Section 3 similarly offers a Modal Defeasible Logic for reasoning about the interpretation of structured provisions. Some discussion end the paper.

## 2. OPTION 1: ABSTRACT NORMS

Let us use Defeasible Logic to reason about interpretive arguments that handle the overall meaning of legal provisions intended as argumentative, abstract (i.e., non-analysed) logical units. We thus have the following basic components (among others):

- a set of legal provisions  $n_1, n_2, \dots$  to be interpreted;
- as set of literals  $a, b, \dots$ , corresponding to any sentences, which can be used to offer a sentential meaning to any provision  $n$  (a literal  $a$  is the meaning of provision  $n$ );
- a set of interpretative acts or interpretations  $l_1, l_2, \dots$  (literal interpretation, teleological interpretation, etc.) that return for any legal provision a sentential meaning for it;
- a set of rules for reasoning about ordinary meanings, i.e., literals, and rules encoding interpretive arguments (i.e., rules that state what interpretive act can be obtained under suitable conditions).

## 2.1 Formal Language

The language of our logic is formally defined as follows.

**DEFINITION 1 (LANGUAGE, OPTION 1).** Let  $\text{PROP} = \{a, b, \dots\}$  be a set of propositional atoms,  $\text{NORM} = \{n_1, n_2, \dots\}$  a set of legal provisions,  $\text{INTR} = \{\mathcal{I}_1, \mathcal{I}_2, \dots\}$  a set of interpretation functions (for example, denoting literal interpretation, etc.),  $\text{MOD} = \{\text{Obl}, \text{Adm}\}$  a set of modal operators where **Obl** is the modality for denoting obligatory interpretations and interpretation outcomes and **Adm** for denoting the admissible ones, and  $\sqsubseteq$  and  $\approx$  the relations denoting semantic inclusion and equivalence.<sup>3</sup>

1. The set  $\text{Lit} = \text{PROP} \cup \{\neg p \mid p \in \text{PROP}\}$  denotes the set of literals.
2. The complementary of a literal  $q$  is denoted by  $\sim q$ ; if  $q$  is a positive literal  $p$ , then  $\sim q$  is  $\neg p$ , and if  $q$  is a negative literal  $\neg p$ , then  $\sim q$  is  $p$ .
3. The set  $\text{ModLit} = \{\square a, \neg \square a \mid a \in \text{Lit}, \square \in \text{MOD}\}$  denotes the set of modal literals.
4. The set  $\text{Int} = \{l_i(n, a), \neg l_i(n, a) \mid \exists \mathcal{I}_i : \text{NORM} \mapsto \text{Lit} \in \text{INTR} : \mathcal{I}_i(n) = a\}$  denotes the set of interpretive acts and their negations: an expression  $l_i(n, a)$ , for instance, means that the interpretation  $l_i$  of provision  $n$  returns that the literal  $a$  is the case.
5. The complementary of an interpretation  $\phi$  is denoted by  $\sim \phi$  and is defined as follows:

$$\begin{array}{ll} \phi & \sim \phi \\ l_i(n, a) & \sim l_i(n, a) \in \{\neg l_i(n, a), l_i(n, b), l_j(n, c) \mid \\ & a \neq b, \text{ and either } a \not\sqsubseteq c \text{ or } c \not\sqsubseteq a\} \\ \neg l_i(n, a) & \sim \neg l_i(n, a) \in \{l_i(n, a), l_i(n, c) \mid \\ & a \approx c\} \end{array}$$

We will also use the notation  $\pm l_i(n, a)$  to mean respectively  $l_i(n, a)$  and  $\sim l_i(n, a)$ . Hence,  $\sim \pm l_i(n, a)$  means  $\mp l_i(n, a)$ .

6. The set of qualified interpretations is  $\text{ModIntr} = \{\square \phi, \neg \square \phi \mid \phi \in \text{Int}, \square \in \text{MOD}\}$ .
7. The complementary of a modal literal or qualified interpretation  $l$  is defined as follows ( $\phi \in \text{Lit} \cup \text{Int}$ ):

$$\begin{array}{ll} \mathbf{I} & \sim \mathbf{I} \\ \text{Obl} \phi & \sim \text{Obl} \phi \in \{\neg \text{Obl} \phi, \text{Obl} \sim \phi, \text{Adm} \sim \phi, \neg \text{Adm} \phi\} \\ \neg \text{Obl} \phi & \sim \neg \text{Obl} \phi = \text{Obl} \phi \\ \text{Adm} \phi & \sim \text{Adm} \phi \in \{\neg \text{Adm} \phi, \text{Obl} \sim \phi\} \\ \neg \text{Adm} \phi & \sim \neg \text{Adm} \phi = \text{Adm} \phi \end{array}$$

**REMARK 1.** A complementary of an interpretation  $l_i$  is another interpretation that is necessarily incompatible with  $l_i$ . Given any  $l_i(n, a)$  (“ $n$  is interpreted by  $l_i$  as  $a$ ”), its propositional negation  $\neg l_i(n, a)$  (“it is false that  $n$  is interpreted by  $l_i$  as  $a$ ”) leads to an incompatibility. But other cases should be considered. First of all, the same interpretation (literal, teleological, etc.) of the same provision cannot give two (or

<sup>3</sup>A simple characterisation of  $\sqsubseteq$  and  $\approx$  is given in Definition 5. As we will see, if a rule says  $a \rightarrow b$ , then  $a \sqsubseteq b$ , but the other way around does not necessarily hold. Similarly, if  $a = b$ , then  $a \equiv b$ , while the opposite direction is not necessarily the case.

more) distinct results, even if these ones are somehow compatible. For instance, let us consider the following provision from the Italian penal code:

Art. 575. Homicide. Whoever causes the death of a man [uomo] is punishable by no less than 21 years in prison.

The literal interpretation (denoted, e.g., by  $l_1$ ) of art. 575 cannot lead at the same time to reading “man”, e.g., as “adult male human being” and as “human being”, even though being an adult male human being entails being a human being. If we consider two distinct interpretations  $l_i$  and  $l_j$  of the same provision, things are slightly different because we resort to two different interpretive arguments. Indeed, suppose to fictionally change art. 575 and replace “man” with “smart Italian person”: if the lawmaker likes people from Florence, one may resort to a psychological argument [14, 15] and support the interpretation  $l_1$  driven by the actual intent of punishing those who kill Florentine people. On the other side, if a naturalistic argument [14, 15] supporting an interpretation  $l_j$  leads to reading “smart Italian person” as “person living in a nice place”, the fact that  $l_i$  and  $l_j$  are necessarily incompatible depends, for instance, on whether being Florentine does not entail to live in a nice place: if this is the case those interpretations cannot be compatible. This analysis is partial and debatable: more logical relations between the propositional content of provisions can be identified [3, 4], thus allowing to refine the concept of conflicting interpretation. Also, one may argue that any interpretations  $l_i(n, a)$  and  $l_j(n, a')$  of  $n$  are in conflict [14] whenever  $a \neq a'$ , irrespective of any possible logical relation (entailment, semantic overlapping, etc.) between  $a$  and  $a'$ . In fact, while we will rely on the above analysis to illustrate the logical machinery, we do not commit to it: different concepts of complementary interpretation can be introduced. More comments will be given in the conclusions of the paper.

We introduce a preference operator  $\otimes$  over interpretations. Indeed, when more interpretations of a provision are available, one may wonder what is the most preferred one [14]. For example, an expression like  $l_i(n, a) \otimes l_j(n, a) \otimes l_k(n, b)$  means that the most preferred interpretation of provision  $n$  is  $a$  via  $l_i$ ; if this is not admissible, then the second best choice is the same result  $a$  via a different interpretation type  $l_j$ ; if this last is not admissible, too,  $n$  is interpreted as  $b$  using interpretation  $l_k$ . This operator is thus used to build chains of preferences, called  $\otimes$ -expressions. The formation rules for well-formed  $\otimes$ -expressions are:

1.  $\forall n \in \text{NORM}, \forall l_i \in \text{Int}, \text{ and } \forall a \in \text{Lit}, \text{ each } l_i(n, a) \text{ is an } \otimes\text{-expression};$
2. If  $\phi_1, \dots, \phi_n \in \text{Int}$ , then  $\phi_1 \otimes \dots \otimes \phi_n$  is an  $\otimes$ -expression;
3. Nothing else is an  $\otimes$ -expression.

EXAMPLE 1. Let us consider again art. 575 of the Italian penal code. The original version speaks of “the death of a man”, i.e., in the ordinary sense, “an adult male human being”.<sup>4</sup> It goes without saying that this provision does not confine to adult male persons only. However, for the sake of argument let us assume that the literal interpretation of art.

<sup>4</sup>Entry “Man”, *Online Cambridge English Dictionary*, 2014.

575 (denoted by  $l_1$ ) corresponds to accepting the ordinary meaning of the term “man [uomo]”; then, we would have that  $l_1(\text{art.575}, a)$ , where  $a =$  “Whoever causes the death of an adult male person is punishable by no less than 21 years in prison”. Other options would rather interpret art. 575 by covering any person (not only adult men): this could be the same reasonable conclusion of arguments such as those from general principles ( $l_2(\text{art.575}, b)$ ) and from substantive reasons ( $l_3(\text{art.575}, b)$ ), where  $b =$  “Whoever causes the death of a person is punishable by no less than 21 years in prison”. Suppose  $l_2$  is the best choice,  $l_3$  the second best choice, and  $\neg l_1$  the least preferred option; the following  $\otimes$ -expression represents this preference ordering:

$$l_2(\text{art.575}, b) \otimes l_3(\text{art.575}, b) \otimes \neg l_1(\text{art.575}, a).$$

We stipulate that  $\otimes$  obeys the following properties:

1.  $\phi \otimes (\psi \otimes \pi) = (\phi \otimes \psi) \otimes \pi$  (associativity);
2.  $\bigotimes_{i=1}^n \phi_i = (\bigotimes_{i=1}^{k-1} \phi_i) \otimes (\bigotimes_{i=k}^n \phi_i)$  where there exists  $j$  such that  $\phi_j = \phi_k$  and  $j < k$  (duplication and contraction on the right).

Given an  $\otimes$ -expression  $A$ , the *length* of  $A$  is the number of interpretations in it. Given an  $\otimes$ -expression  $A \otimes \phi \otimes C$  (where  $A$  and  $C$  can be empty), the *index* of  $\phi$  is the length of  $A \otimes \phi$ . We also say that  $\phi$  appears at index  $n$  in  $A \otimes \phi$  if the length of  $A \otimes \phi$  is  $n$ .

We adopt the standard definitions in Defeasible Logic of strict rules, defeasible rules, and defeaters<sup>5</sup> [2] with the following modifications:

**strict rules**, encoding the monotonic part of the knowledge base, are used to reason about basic linguistic meanings in ordinary parlance and thus contain only literals in their antecedents and consequents; we denote them as usual with the arrow  $\rightarrow$  and mark them with a superscript  $M$ ;<sup>6</sup>

**the non-monotonic part** of the knowledge base (defeasible rules and defeaters) is used to reason about the interpretations of provisions; these rules contain literals, interpretations and qualified interpretations in their antecedent, and  $\otimes$ -expressions in their consequents; we denote them as usual with the arrows  $\Rightarrow$  and  $\rightsquigarrow$  and mark them with a superscript  $I$ .

DEFINITION 2 (RULES, OPTION 1). Let  $\text{Lab}$  be a set of arbitrary labels. Every rule is of the type

$$r : A(r) \hookrightarrow^X C(r)$$

where

1.  $r \in \text{Lab}$  is the name of the rule, and either
2. if  $r$  is a strict rule, then  $r$  is a meaning rule such that
  - $A(r) = \{a_1, \dots, a_n\}$ , the antecedent (or body) of the rule, is the set of the premises of the rule (alternatively, it can be understood as the conjunction of all the literals in it). Each  $a_i$  is a literal;

<sup>5</sup>A defeater is a rule which prevents opposite conclusions without allowing to positively deriving anything.

<sup>6</sup>One may notice that any suitable logic modelling the meanings of ordinary language should be non-monotonic: a sentence like “all birds are flying animals” admits exceptions. We use only strict rules to simplify proof theory.

- $\hookrightarrow \Rightarrow^M$ ;
- $C(r)$  is the consequent (or head) of the rule, which is a literal;

3. otherwise,  $r$  is an interpretation rule such that

4.  $A(r) = \{\phi_1, \dots, \phi_n\}$ , the antecedent (or body) of the rule is such that each  $\phi_i$  is either a literal  $l \in \text{Lit}$ , a modal literal  $Y \in \text{ModLit}$ , or a qualified interpretation  $X \in \text{ModIntr}$ ;

5.  $\hookrightarrow \in \{\Rightarrow^I, \rightsquigarrow^I\}$  denotes the type of the rule. If  $\hookrightarrow$  is  $\Rightarrow^I$ , the rule is a defeasible rule, while if  $\hookrightarrow$  is  $\rightsquigarrow^I$ , the rule is a defeater;

6.  $C(r)$  is the consequent (or head) of the rule, which is an  $\otimes$ -expression; if  $\hookrightarrow$  is  $\rightsquigarrow$ , then  $C(r)$  is a single interpretation, i.e., an  $\otimes$ -expression of length 1.

EXAMPLE 2. Let us consider the following examples of meaning (strict) rules:

- $r_1 : \text{kill\_man} \rightarrow^M \text{kill\_person}$
- $r_2 : \text{kill\_male}, \text{kill\_adult} \rightarrow^M \text{kill\_man}$
- $r_3 : \text{kill\_female} \rightarrow^M \neg \text{kill\_male}$
- $r_4 : \text{kill\_newborn} \rightarrow^M \neg \text{kill\_adult}$

Consider now that paragraph 1 of art. 3 of the Italian constitution reads as follows:

*All citizens have equal social status and are equal before the law, without regard to their sex, race, language, religion, political opinions, and personal or social conditions.*

The most preferred interpretation  $l_3$  (interpretation from substantive reasons) of it leads to  $c$ , which corresponds to the following sentence:

*All persons have equal social status and are equal before the law, without regard to their sex, race, language, religion, political opinions, and personal or social conditions.*

The following interpretation defeasible rule could be in order:

- $r_5 : \text{kill\_adult}, \text{kill\_female}, \text{Obl } l_3(\text{art.3}, c) \Rightarrow^I$   
 $l_2(\text{art.575}, b) \otimes l_3(\text{art.575}, b) \otimes \neg l_1(\text{art.575}, a).$

In other words, if art. 3 of the Italian constitution states formal equality before the law without regard also to gender identity, then  $b$  is the best interpretation outcome of art. 575 of the penal code, with  $l_2$  (e.g., interpretation from general principles) as preferred over  $l_3$  (say, interpretation from substantive reasons). In the light of art. 3, if these two interpretive options are refuted by any stronger interpretive conclusions, then the last sub-ideal option is to reject the interpretation  $l_1$  from ordinary meaning (according to which only the homicide of adult male persons is punishable!).

Given a set of rules  $R$ , we use the following abbreviations for specific subsets of rules:

- $R_{\text{def}}^I$  denotes the set of all defeaters in the set  $R$  (since defeaters can only be of type  $I$ , we will omit the superscript in the remainder);

- $R^I[\phi, n]$  is the set of rules where the interpretation  $\phi$  appears at index  $n$  in the consequent. The set of (defeasible) rules where  $\phi$  appears at any index  $n$  is denoted by  $R^I[\phi]$  (since rules with  $\otimes$ -expressions can only be of type  $I$ , we will omit the superscript in the remainder);

- $R^M[l]$  denotes the set of all strict rules with head  $l$ .

DEFINITION 3 (INTERPRETATION THEORY). An Interpretation Theory is a structure  $D = (F, R, >)$ , where  $F$ , the set of facts, is a set of literals, modal literals, and qualified interpretations,  $R$  is a set of rules and  $>$ , the superiority relation, is a binary relation over  $R$ .

An interpretation theory corresponds to a knowledge base providing us with interpretive arguments about legal provisions. The superiority relation is used for conflicting rules, i.e., rules whose conclusions are complementary. We do not impose any restriction on the superiority relation: it just determines the relative strength of two rules.

EXAMPLE 3. Consider art. 575 of the Italian penal code, art. 3 of the Italian constitution. Also consider art. 578 (of the penal code):

*Art. 578. Infanticide. The mother who causes the death of her newborn baby immediately after birth [...], when this act is committed in conditions of material or moral distress, is punishable with a sentence between 4 and 12 years of prison.*

Assume that

- $a =$  Whoever causes the death of a adult male person is punishable by no less than 21 years in prison
- $a' =$  Whoever causes the death of a male person is punishable by no less than 21 years in prison
- $b =$  Whoever causes the death of a person is punishable by no less than 21 years in prison
- $c =$  All persons have equal social status and are equal before the law, without regard to their sex, race, language, religion, political opinions, and personal or social conditions.

- $l_1 =$  Literal interpretation or from ordinary meaning
- $l_2 =$  Interpretation from general principles
- $l_3 =$  Interpretation from substantive reasons
- $l_4 =$  Interpretation by coherence

The following theory reconstructs an interpretive toy scenario in the Italian legal system.

$F = \{\text{kill\_adult}, \text{kill\_female}, \text{Obl } l_3(\text{art.3}, c), \text{Obl } l_1(\text{art.578}, d)\}$

$R = \{r_5 : \text{kill\_adult}, \text{kill\_female}, \text{Obl } l_3(\text{art.3}, c) \Rightarrow^I$   
 $l_2(\text{art.575}, b) \otimes l_3(\text{art.575}, b) \otimes \neg l_1(\text{art.575}, a),$   
 $r_6 : \text{Obl } l_1(\text{art.578}, d) \Rightarrow^I l_1(\text{art.575}, a) \otimes l_4(\text{art.575}, a),$   
 $r_7 : \Rightarrow^I l_1(\text{art.575}, a'), r_8 : b \rightarrow^M a\}$

$> = \{r_6 > r_5\}$

Rule  $r_5$  has been already introduced above. Rule  $r_6$  states that, in case one kills an adult person and art. 578 ought to be interpreted literally, then we have a reason to interpret art. 575 by coherence (with respect to art. 578) as  $a$ . Rule

$r_7$  establishes by default that art. 575 be literally interpreted as  $a'$ . Finally, rule  $r_8$  is a meaning rule saying that  $b$  entails  $a$ , i.e., that a provision punishing whoever causes the death of a person entails that a provision should be in the system that punishes whoever causes the death of a an adult male person.

## 2.2 Proof Theory: Preliminaries

This section presents some preliminary concepts for reasoning about interpretations: the notions of proof, semantic inclusion, and rule applicability.

**DEFINITION 4.** [Proofs] A proof  $P$  in an interpretation theory  $D$  is a linear sequence  $P(1) \dots P(n)$  of tagged expressions in the form of  $+\Delta^M q$ ,  $-\Delta^M q$  (with  $q \in \text{Lit}$ ),  $+\partial_{\square}^I \phi$  and  $-\partial_{\square}^I \phi$  (with  $\phi \in \text{Int}$  and  $\square \in \text{MOD}$ ),  $+\partial^{\square} l$  and  $-\partial^{\square} l$  (with  $l \in \text{Lit}$  and  $\square \in \text{MOD}$ ), where  $P(1) \dots P(n)$  satisfy the proof conditions given in Definitions 9–16.

The tagged literal  $+\Delta^M q$  means that  $q$  is strictly or monotonically provable in  $D$  using only facts and strict rules, while  $-\Delta^M q$  means that there is no strict or monotonic proof for  $q$  in  $D$ . The tagged interpretation  $+\partial_{\square}^I \phi$  means that the interpretation  $\phi$  is defeasibly provable in  $D$  with modality  $\square$ , while  $-\partial_{\square}^I \phi$  means that  $\phi$  is defeasibly refuted with modality  $\square$ . The tagged literal  $+\partial^{\square} l$  means that  $l$  is defeasibly provable in  $D$  with modality  $\square$ ,<sup>7</sup> while  $-\partial^{\square} l$  means that  $l$  is defeasibly refuted with modality  $\square$ . The initial part of length  $n$  of a proof  $P$  is denoted by  $P(1..n)$ .

Proof theory for strict conclusions is trivially the same for  $\Delta$ -conclusions in standard Defeasible Logic [2]. We omit it here. We only state the following, to define  $\sqsubseteq$  and  $\approx$ :

**DEFINITION 5** (SEMANTIC INCLUSION). Let  $D = (F, R, >)$  an interpretation theory. If  $p, q \in \text{Lit}$ , we write  $p \sqsubseteq_D q$ <sup>8</sup> to say that the meaning of literal  $p$  is semantically included by  $q$  in  $D$ . If  $\mathcal{P} = \{r \mid r \in R^M, p \in A(r) \text{ or } p \in C(r)\}$ ,  $p \sqsubseteq_D q$  if

- $D \vdash +\Delta^M p$  and  $D \vdash +\Delta^M q$ , and
- $D_{\sqsubseteq} = (F_{\sqsubseteq} R_{\sqsubseteq}^M, >) \not\vdash +\Delta^M q$  where (i)  $F_{\sqsubseteq} = F \setminus \{p, q\}$  and (ii)  $R_{\sqsubseteq}^M$  is such that  $\{r \mid r \in R_{\sqsubseteq}\} \setminus \mathcal{P}$ . We state that  $p \approx_D q$  iff  $p \sqsubseteq_D q$  and  $q \sqsubseteq_D p$ .

In other words, a literal  $p$  is included by  $q$  in a theory  $D$  when both are obtained in  $D$  using meaning rules, but, if  $p$  and  $q$  are removed from the set of facts (if they are there) and all meaning rules having  $p$  in the head of body are removed as well, then  $q$  no longer follows from the theory. Hence, it means that  $p$  is decisive for having  $q$ . As we have already observed, semantic inclusion does not fully coincide with entailment via a single rule. Indeed, if  $p \sqsubseteq q$ , this could mean, for example, that  $\{a_1, \dots, a_n, p\} \vdash_D q$ ,  $\{a_1, \dots, a_n\} \not\vdash_D q$ , but also that  $\{p\} \not\vdash_D q$ .

Let us work on the proof theory for deriving qualified interpretations and first consider the admissible ones (i.e., for any interpretation  $\phi$ , those having the form  $\text{Adm}\phi$ ). The first thing to do is to define when a rule is applicable or discarded. A rule is *applicable* for an interpretation  $\phi$  if  $\phi$  occurs in the head of the rule, all literals in the antecedent

are  $\Delta^M$ -provable and all the qualified interpretations and modal literals in the antecedent have been defeasibly proved (with the appropriate modalities). On the other hand, a rule is *discarded* if at least one of the literals, modal literals, or of the qualified interpretations in the antecedent has not been proved. However, as interpretation  $\phi$  might not appear as the first element in an  $\otimes$ -expression in the head of the rule, some additional conditions on the consequent of rules must be satisfied. Defining when a rule is applicable or discarded is essential to characterise the notion of provability for admissible interpretations ( $\pm\partial_{\text{Adm}}^I$ ).

**DEFINITION 6.** A rule  $r \in R^I$  is body-applicable in the proof  $P$  at  $P(n+1)$  iff for all  $a_i \in A(r)$ :

1. if  $a_i = \square\psi$ ,  $\psi \in \text{Int}$ , then  $+\partial_{\square}^I \psi \in P(1..n)$  with  $\square \in \text{MOD}$ ;
2. if  $a_i = \neg\square\psi$  then  $-\partial_{\square}^I \psi \in P(1..n)$  with  $\square \in \text{MOD}$ ;
3. if  $a_i = \square l$ ,  $l \in \text{Lit}$ , then  $+\partial^{\square} l \in P(1..n)$ ;
4. if  $a_i = \neg\square l$ ,  $l \in \text{Lit}$ , then  $-\partial^{\square} l \in P(1..n)$ ;
5. if  $a_i = l \in \text{Lit}$  then  $D \vdash +\Delta^M l$ .

A rule  $r \in R^I$  is body-discarded iff  $\exists a_i \in A(r)$  such that

1. if  $a_i = \square\psi$ ,  $\psi \in \text{Int}$ , then  $-\partial_{\square}^I \psi \in P(1..n)$  with  $\square \in \text{MOD}$ ;
2. if  $a_i = \neg\square\psi$ ,  $\psi \in \text{Int}$ , then  $+\partial_{\square}^I \psi \in P(1..n)$  with  $\square \in \text{MOD}$ ;
3. if  $a_i = \square l$ ,  $l \in \text{Lit}$ , then  $-\partial^{\square} l \in P(1..n)$ ;
4. if  $a_i = \neg\square l$ ,  $l \in \text{Lit}$ , then  $+\partial^{\square} l \in P(1..n)$ ;
5. if  $a_i = l \in \text{Lit}$  then  $D \vdash -\Delta^M l$ .

**DEFINITION 7.** An interpretation rule  $r \in R^I[\phi, j]$  such that  $C(r) = \phi_1 \otimes \dots \otimes \phi_n$  is applicable for the interpretation  $\phi$  at index  $j$ , with  $1 \leq j \leq n$ , in the condition for  $\pm\partial_{\text{Adm}}^I$  (in the proof  $P$  at  $P(n+1)$ ) iff

1.  $r$  is body-applicable; and
2. for all  $\phi_k \in C(r)$ ,  $1 \leq k < j$ ,  $+\partial_{\text{Adm}}^I \sim \phi_k \in P(1..n)$ .

Conditions (1) represents the requirements on the antecedent stated in Definition 6; condition (2) on the head of the rule states that each element  $\phi_k$  prior to  $\phi$  must be refuted as an admissible interpretation.

**DEFINITION 8.** An interpretation rule  $r \in R^I[\phi, j]$  such that  $C(r) = \phi_1 \otimes \dots \otimes \phi_n$  is discarded for interpretation  $\phi$  at index  $j$ , with  $1 \leq j \leq n$  in the condition for  $\pm\partial_{\text{Adm}}^I$  iff

1.  $r$  is body-discarded; or
2. there exists  $\phi_k \in C(r)$ ,  $1 \leq k < j$ , such that  $+\partial_{\text{Adm}} \phi_k \in P(1..n)$ .

In this case, condition (2) states that there exists at least one explicit derived admissible interpretation prior to  $\phi$ .

<sup>7</sup>As resulting here from a provable A-interpretation.

<sup>8</sup>When clear from the context, we omit the subscript.

### 2.3 Proof Theory: Interpretation Activities

We work in this subsection on the proof theory for reasoning about interpretation activities (A-Interpretations).

Let us define the proof conditions for  $\pm\partial_{\text{Adm}}$ .

DEFINITION 9 (ADMISSIBLE A-INTERPRETATIONS).

The proof condition of defeasible provability for admissible A-interpretations is

$+\partial_{\text{Adm}}^I$ : If  $P(n+1) = +\partial_{\text{Adm}}^I\phi$  then

- (1)  $\text{Adm}\phi \in F$  or  $\text{Obl}\phi \in F$ , or
  - (2.1)  $\sim\text{Adm}\phi \notin F$ , and
  - (2.2)  $\exists r \in R[\phi, i]$  such that  $r$  is applicable for  $\phi$ , and
  - (2.3)  $\forall s \in R[\sim\phi, j]$ , either
    - (2.3.1)  $s$  is discarded for  $\sim\phi$ , or
    - (2.3.2)  $\exists t \in R[\phi, k]$  such that  $t$  is applicable for  $\phi$  and  $t > s$ .

To show that  $\phi$  is defeasibly provable as an admissible interpretation, there are two ways: (1)  $\text{Adm}\phi$  or  $\text{Obl}\phi$  are a fact, or (2)  $\text{Adm}\phi$  must be derived by the rules of the theory. In the second case, three conditions must hold: (2.1) any complementary of  $\text{Adm}\phi$  does belong to the facts; (2.2) there must be a rule introducing the admissibility for  $\phi$  which can apply; (2.3) every rule  $s$  for  $\sim\phi$  is either discarded or defeated by a stronger rule for  $\phi$ .

The strong negation [2] of Definition 9 gives the negative proof condition for admissible interpretations.

DEFINITION 10 (REFUTED A-INTERPRETATIONS). The proof condition of defeasible refutability for admissible A-interpretations is

$-\partial_{\text{Adm}}^I$ : If  $P(n+1) = -\partial_{\text{Adm}}^I\phi$  then

- (1)  $\text{Adm}\phi \notin F$  and  $\text{Obl}\phi \notin F$ , and
  - (2.1)  $\sim\text{Adm}\phi \in F$ , or
  - (2.2)  $\forall r \in R[\phi, i]$   $r$  is discarded for  $\phi$ , or
  - (2.3)  $\exists s \in R[\sim\phi, j]$  such that
    - (2.3.1)  $s$  is applicable for  $\sim\phi$ , and
    - (2.3.2)  $\forall t \in R[\phi, k]$ , either  $t$  is discarded for  $\phi$  or  $t \not> s$ .

Proof conditions for  $\pm\partial_{\text{Obl}}$  are much easier but we need to work on the fact that  $\phi$  is an interpretation of any given provision  $n$  and we have to make explicit its structure. Indeed, that an interpretation  $l_i$  for the provision  $n$  is obligatory means that  $l_i$  is admissible and that no other (non-conflicting) interpretations for  $n$  is admissible.

DEFINITION 11 (OBLIGATORY A-INTERPRETATIONS).

The proof condition of defeasible provability for obligatory A-interpretations is

$+\partial_{\text{Obl}}^I$ : If  $P(n+1) = +\partial_{\text{Obl}}^I\pm l_i(n, a)$  then

- (1)  $\text{Obl}\pm l_i(n, a) \in F$  or
  - (2.1)  $\sim\text{Obl}\pm l_i(n, a) \notin F$ , and
  - (2.2)  $+\partial_{\text{Adm}}^I\pm l_i(n, a) \in P(1..n)$ , and
  - (2.3)  $\forall s \in R[\pm l_m(n, b), j]$ , such that  $l_m(n, b) \neq \sim l_i(n, a)$ , either
    - (2.3.1)  $s$  is discarded for  $\pm l_m(n, b)$ , or
    - (2.3.2)  $\exists t \in R_{\text{defl}}[\sim\pm l_m(n, b), k]$  such that  $t$  is applicable for  $\sim\pm l_m(n, b)$  and  $t > s$ .

The negative proof condition for obligatory interpretations is as follows.

DEFINITION 12 (NON-OBLIGATORY A-INT.). The proof condition of defeasible refutability for obligatory A-interpretations is

$-\partial_{\text{Obl}}^I$ : If  $P(n+1) = -\partial_{\text{Obl}}^I\pm l_i(n, a)$  then

- (1)  $\text{Obl}\pm l_i(n, a) \notin F$  and
  - (2.1)  $\sim\text{Obl}\pm l_i(n, a) \in F$ , or
  - (2.2)  $-\partial_{\text{Adm}}^I\pm l_i(n, a) \in P(1..n)$ , or
  - (2.3)  $\exists s \in R[\pm l_m(n, b), j]$  such that  $l_m(n, b) \neq \sim l_i(n, a)$  and
    - (2.3.1)  $s$  is applicable for  $\pm l_m(n, b)$ , and
    - (2.3.2)  $\forall t \in R[\sim\pm l_m(n, b), k]$  either  $t$  is discarded for  $\sim\pm l_m(n, b)$  or  $t \not> s$ .

EXAMPLE 4. Consider the theory in Example 3. Facts make rules  $r_5$  and  $r_6$  applicable. Rule  $r_7$  has an empty antecedent, so it is applicable, too. Despite  $r_8 \in R$ ,  $b \not\sqsubseteq a$ , then  $r_5$  and  $r_6$  are in conflict. The theory assumes that  $r_6$  is stronger than  $r_5$ , thus we would obtain  $+\partial_{\text{Adm}}^I l_1(\text{art.575}, a)$  (and so  $-\partial_{\text{Adm}}^I l_2(\text{art.575}, b)$  and  $-\partial_{\text{Adm}}^I l_3(\text{art.575}, b)$ ). However, these last conclusions are not obtained because of  $r_7$ :  $r_7$  and  $r_6$  attack each other, thus we in fact have  $-\partial_{\text{Adm}}^I l_1(\text{art.575}, a)$ ,  $-\partial_{\text{Adm}}^I l_1(\text{art.575}, a')$ , and also  $-\partial_{\text{Adm}}^I l_4(\text{art.575}, a)$ . Hence, we reinstate  $+\partial_{\text{Adm}}^I l_2(\text{art.575}, b)$  via  $r_5$ . What interpretations are obligatory? Trivially, we get  $+\partial_{\text{Obl}}^I l_3(\text{art.3}, c)$ . Also, since  $l_2(\text{art.575}, b)$  is the only admissible interpretation of art. 575, then  $+\partial_{\text{Obl}}^I l_2(\text{art.575}, b)$  (check conditions 2.2) and 2.3) in Definition 11). All other interpretations are refuted as obligatory.

### 2.4 Proof Theory: Interpretation Outcomes

Let us show how to prove O-interpretations (i.e., modal literals following from A-interpretations): for example, if  $l_i(n, a)$  is admissible, then  $a$  should be admissible as well.

DEFINITION 13 (ADMISSIBLE O-INTERPRETATIONS).

The proof condition of defeasible provability for admissible O-interpretations is

$+\partial^{\text{Adm}}$ : If  $P(n+1) = +\partial^{\text{Adm}}l$  then

- (1)  $\text{Adml} \in F$  or  $\text{Obl}l \in F$ , or
- (2)  $\exists l_i \in \text{Int}, \exists n \in \text{NORM} : +\partial_{\text{Adm}}^I l_i(n, l) \in P(1..n)$ .

DEFINITION 14 (REFUTED O-INTERPRETATIONS). The proof condition of defeasible refutability for admissible O-interpretations is

$-\partial^{\text{Adm}}$ : If  $P(n+1) = -\partial^{\text{Adm}}l$  then

- (1)  $\text{Adml} \notin F$  and  $\text{Obl}l \notin F$ , and
- (2)  $\forall l_i \in \text{Int}, \forall n \in \text{NORM} : -\partial_{\text{Adm}}^I l_i(n, l) \in P(1..n)$ .

DEFINITION 15 (OBLIGATORY O-INTERPRETATIONS).

The proof condition of defeasible provability for obligatory O-interpretations is

$+\partial^{\text{Obl}}$ : If  $P(n+1) = +\partial^{\text{Obl}}l$  then

- (1)  $\text{Obl}l \in F$ , or
- (2)  $\exists n \in \text{NORM}$  such that
  - (2.1)  $\exists l_i \in \text{Int} : +\partial_{\text{Adm}}^I l_i(n, a) \in P(1..n)$  and
  - (2.2)  $\forall l_j \in \text{Int}, -\partial_{\text{Adm}}^I l_j(n, x) \in P(1..n)$  if  $x \neq a$ .

DEFINITION 16 (NON-OBLIGATORY O-INT.). The proof condition of defeasible refutability for obligatory O-interpretations is

$-\partial^{\text{Obl}}$ : If  $P(n+1) = -\partial^{\text{Obl}}l$  then

- (1)  $\text{Obl}l \notin F$ , and
- (2)  $\forall n \in \text{NORM} :$ 
  - (2.1)  $\forall l_i \in \text{Int} : -\partial_{\text{Adm}}^I l_i(n, a) \in P(1..n)$  or
  - (2.2)  $\exists l_j \in \text{Int}, +\partial_{\text{Adm}}^I l_j(n, x) \in P(1..n)$  and  $x \neq a$ .

Proof theory for O-interpretations is simple. Just notice that obligatory interpretations are nothing but literals resulting from obligatory A-interpretations.

EXAMPLE 5. *If we work again on Example 3, the proofs informally described in Example 4 show that  $+\partial^{\square}b$ , where  $\square \in \{\text{Adm}, \text{Obl}\}$ .*

### 3. OPTION 2: STRUCTURED NORMS

The second option is a refinement of the one presented in Section 2. We add deontic rules expressing legal provisions and interpretation rules here work on the literals composing those deontic rules.

More precisely, as in [8], we introduce obligation rules of the form  $r : a_1, \dots, a_n \Rightarrow^{\text{Obl}} b$  to represent deontic legal provisions; if  $r$  is applicable and undefeated, then we can derive  $\text{Obl}b$ . These rules represent in our framework the legal provisions on which interpretation options apply [5].

Second, interpretation rules work on the literals occurring in obligation rules. In other words, given a rule  $r : a_1, \dots, a_n \Rightarrow^{\text{Obl}} b$ , an interpretation function maps for each provision  $n$  the sequence  $x = \langle a_1, \dots, a_n, b \rangle$  of literals in  $r$  onto another sequence  $y$  of literals that can be identical (literal interpretation), partially different or completely different from  $x$ . Hence, an interpretation  $l_i$  is meant to make the original version of rule  $r$  unusable and the new one—where the literals are changed according to  $y$ —usable to derive an obligation. For instance, if  $r : a_1, a_2 \Rightarrow^{\text{Obl}} b$  and the interpretation  $l_i$  returns  $y = \langle a_1, a'_2, b' \rangle$ , the *interpreted version* of  $r$  according to  $l_i$  is  $r : a_1, a'_2 \Rightarrow^{\text{Obl}} b'$ .

Let us extend and modify language and proof theory.

DEFINITION 17 (LANGUAGE, OPTION 2). *Let*  
 PROP =  $\{a, b, \dots\}$  *be a set of propositional atoms,*  
 NORM =  $\{n_1, n_2, \dots\}$  *a set of provision labels,*  
 INTR =  $\{\mathcal{I}_1, \mathcal{I}_2, \dots\}$  *a set of interpretation functions,*  
 MOD =  $\{\text{Obl}, \text{Adm}\}$  *a set of modal operators, and*  $\sqsubseteq$  *and*  $\approx$  *the relations denoting semantic inclusion and equivalence.*

- *The set Lit, the complementary of literals, the set ModLit, the set ModIntr, the complementary of modal literals and of qualified interpretations are as in Definition 1;*
- *The set Int =  $\{l_i(n, x, y), \neg l_i(n, x, y) \mid \exists \mathcal{I}_i : \text{NORM} \times X \mapsto Y \in \text{INTR} : X, Y \in \bigcup_{j=1}^n \mathcal{P}(\text{Lit}^n) \text{ and } x \in X, y \in Y\}$  denotes the set of interpretations and their negations: an expression  $l_i(n, x, y)$ , for instance, means that the interpretation  $l_i$  of provision  $n$  (corresponding to the sequence of literals  $x$ ) returns a tuple  $y$  of literals that should instead occur in  $n$ .*
- *The complementary of an interpretation  $\phi$  is denoted by  $\sim\phi$  and is defined as follows:*

$$\begin{array}{ll} \phi & \sim\phi \\ l_i(n, x, y) & \sim l_i(n, x, y) \in \{\neg l_i(n, x, y), l_i(n, x, z), \\ & l_j(n, x, w) \mid \exists a \text{ in } y, \exists b \text{ in } w : \\ & \text{either } a \not\sqsubseteq b \text{ or } b \not\sqsubseteq a\} \\ \neg l_i(n, x, y) & \sim \neg l_i(n, x, y) \in \{l_i(n, x, y), l_i(n, x, w) \mid \\ & \exists a \text{ in } y, \exists b \text{ in } w : a \approx b\}. \end{array}$$

Let us define the set of rules, which consists of the set of meaning rules, obligation rules, and interpretation rules.

Meaning rules (i.e., strict rules with the arrow  $\rightarrow^M$ ) are exactly as in Section 2 and their definition is omitted.

DEFINITION 18 (OBLIGATION RULES). *Let Lab be a set of arbitrary provision labels. The set  $\text{Rul}^{\text{Obl}}$  of obligation rules consists of rules of the following type*

$$r : A(r) \hookrightarrow^{\text{Obl}} C(r)$$

where

1.  $r \in \text{Lab}$  *is the name of the rule, and*
2.  $A(r) = \{a_1, \dots, a_n\}$ , *the antecedent (or body) of the rule is such that each  $a_i \in A(r)$  is a literal  $l \in \text{Lit}$ ;*
3.  $\hookrightarrow \in \{\Rightarrow, \rightsquigarrow\}$ ;
4.  $C(r) = l$  *is the consequent (or head) of the rule such that  $l \in \text{Lit}$ .*

DEFINITION 19 ( $\otimes$ -EXPRESSIONS, OPTION 2).  $\otimes$ -expressions are defined as follows:

1. *Each interpretation  $l_i(n, x, y)$  or  $\neg l_i(n, x, y)$  is an  $\otimes$ -expressions where*
  - (a)  $n$  *is an obligation rule  $a_1, \dots, a_n \hookrightarrow^{\text{Obl}} b$ ;*
  - (b)  $x = \langle a_1, \dots, a_n, b \rangle$  *and*  $y = \langle a'_1, \dots, a'_n, b' \rangle$ ;
2. *If  $\phi_1, \dots, \phi_n$  are interpretations, then  $\phi_1 \otimes \dots \otimes \phi_n$  is an  $\otimes$ -expression;*
3. *Nothing else is an  $\otimes$ -expression.*

REMARK 2. *We state that interpretations work on  $n$ -tuples of literals returning other  $n$ -tuples. This option makes things simple: each literal can be possibly replaced but not removed, and no new literals can be added. Hence, we cannot directly express the shift from  $\text{man} \Rightarrow^{\text{Obl}} b$  into  $\text{adult, male} \Rightarrow^{\text{Obl}} b$ . To cope with this, it suffices to state that any  $l_i(n, x, y)$  is such  $x$  and  $y$  have the form  $\langle X, l \rangle$  where  $X$  is a set of literals. Conflicting interpretations leading to  $\langle X, l \rangle$  and  $\langle Y, l' \rangle$  would require to check if  $X \neq Y$  or  $l \neq l'$ .*

Interpretation rules (i.e., defeasible rules with the arrow  $\Rightarrow^I$ ) are exactly as in Section 2: the difference is here that their consequents are  $\otimes$ -expressions as in Definition 19. Interpretation theories are as in Definition 3.

### 3.1 Proof Theory

Proof theory must be significantly revised. A first major change is that we can derive obligatory literals using obligation rules (i.e., legal provisions). A second major change is that the derivation of an obligatory interpretation with respect to a provision  $n$  has a double effect: (i) making ineffective the original version of  $n$ ; (ii) making applicable the new version that follows from the derived obligatory interpretation. In other words, the provability of interpretations allows us to update the set of applicable obligation rules: obligatory interpretations thus change obligation rules and work in a similar way of legal modifications as modelled in [9]. This means that we need a mechanism for deriving *usable obligation rules*, i.e., those obligation rules that can be used to derive obligations: if the interpretation  $l_i(r, \langle a_1, a_2, b \rangle, \langle a_1, a'_2, b' \rangle)$  is derived (as obligatory), then  $r : a_1, a_2 \Rightarrow^{\text{Obl}} b$  is unusable and  $r : a_1, a'_2 \Rightarrow^{\text{Obl}} b'$  is usable, i.e., this last version of  $r$  can only be employed to derive obligations. In contrast, deriving A-interpretations does not substantially change: the only difference is in the format of

interpretation rules ( $\otimes$ -expressions are different, and modal literals can be differently proved), but proofs run exactly as before. It is enough to keep Definitions 9, 10, 11, and 12.

DEFINITION 20. An obligation rule  $r \in R^{\text{Obl}}$  is applicable if it is body-applicable. An obligation rule  $r \in R^{\text{Obl}}$  is discarded if it is body-discarded.

Let us first state the procedures for deriving usable obligation rules, which require to integrate Definition 4 with new proof tags  $\pm\partial^R$ .

DEFINITION 21 (USABLE OBLIGATION RULES). The proof condition of defeasible provability for usable obligation rules is

- $+\partial^R$ : If  $P(n+1) = +\partial^R r : a_1, \dots, a_n \hookrightarrow^{\text{Obl}} b$  then
- (1)  $r \in R^{\text{Obl}}$ , and, either
  - (2) if  $A(r) = \{a_1, \dots, a_n\}$  and  $C(r) = b$ , then
    - $\forall i \in \text{Int} : -\partial_{\text{Adm}}^I i(r, x, y)$  if  $y \neq \langle a_1, \dots, a_n, b \rangle$ , or
    - (3)  $+\partial_{\text{Obl}}^I i(r, x, y)$  if  $y = \langle a_1, \dots, a_n, b \rangle$ .

DEFINITION 22 (UNUSABLE OBLIGATION RULES).

- $-\partial^R$ : If  $P(n+1) = -\partial^R r : a_1, \dots, a_n \hookrightarrow^{\text{Obl}} b$  then
- (1)  $r \notin R^{\text{Obl}}$ , or
  - (2)  $A(r) = \{a_1, \dots, a_n\}$ ,  $C(r) = b$ , and
    - $\exists i \in \text{Int} : +\partial_{\text{Adm}}^I i(r, x, y)$  and  $y \neq \langle a_1, \dots, a_n, b \rangle$ , or
    - (3)  $-\partial_{\text{Obl}}^I i(r, x, y)$  and  $y = \langle a_1, \dots, a_n, b \rangle$ .

REMARK 3. An obligation rule  $r : a_1, \dots, a_n \Rightarrow^{\text{Obl}} b$  is usable if one of the following mutually exclusive conditions holds: either (2)  $r : a_1, \dots, a_n \Rightarrow^{\text{Obl}} b$  is in the set  $R$  of rules of the theory and no other version is admissible, or (3) a different version is in  $R$  but an interpretation  $i_i$  is provable as obligatory that transforms this version into the one where  $\{a_1, \dots, a_n\}$  and  $C(r) = b$ . Notice that we could relax condition (3) and just require that  $i_i$  is admissible: we would change proof theory in such a way as to possibly get more derived compatible versions of  $r$ , i.e., that the same provision  $n$  may lead to different normative sentences that are implied by  $n$ .

DEFINITION 23 (PROVABILITY FOR OBLIGATIONS). The proof condition of defeasible provability for obligations is

- $+\partial^{\text{Obl}}$ : If  $P(n+1) = +\partial^{\text{Obl}} l$  then
- (1)  $\text{Obl} l \in F$ , or
    - (2.1)  $\sim \text{Obl} l \notin F$  and
    - (2.2)  $\exists r \in \text{Rul}^{\text{Obl}}[l] : +\partial^R r \in P(1..n)$ ,  $r$  is applicable, and
    - (2.3)  $\forall s \in R^{\text{Obl}}[\sim l]$ , either
      - (2.3.1)  $s$  is discarded, or
      - (2.3.2)  $\exists t \in R^{\text{Obl}}[l] : t$  is applicable and  $t > s$ .

DEFINITION 24 (REFUTABILITY FOR OBLIGATIONS). The proof condition of defeasible refutability for obligations is

- $-\partial^{\text{Obl}}$ : If  $P(n+1) = -\partial^{\text{Obl}} l$  then
- (1)  $\text{Obl} l \notin F$ , and
    - (2.1)  $\sim \text{Obl} l \in F$  or
    - (2.2)  $\forall r \in \text{Rul}^{\text{Obl}}[l] : -\partial^R r \in P(1..n)$ , or  $r$  is discarded, or
    - (2.3)  $\exists s \in R^{\text{Obl}}[\sim l] :$ 
      - (2.3.1)  $s$  is applicable, and
      - (2.3.2)  $\forall t \in R^{\text{Obl}}[l] : \text{either } t \text{ is discarded or } t \not> s$ .

REMARK 4. The conditions for proving obligations are very similar to the ones in [8]. A relevant difference is that we require here that any rule employed to derive an obligation must be usable. Interestingly, the provability of obligations affects the derivation of interpretations. Consider when, given an obligation rule  $r : a \Rightarrow^{\text{Obl}} b$  in  $R$ , an interpretation rule  $r' : \text{Obl} b \Rightarrow^I i_i(r, \langle a, \neg b \rangle)$  holds, which would lead to  $r : a \Rightarrow^{\text{Obl}} \neg b$ , if  $i_i$  results from  $r'$  to be obligatory. As expected, we cannot prove this new version of  $r$ , whose provability relies on an incompatible reading of  $r$ .

EXAMPLE 6. Let us develop Example 3 introducing the new ideas of this section. Consider the following:<sup>9</sup>

- $a = \langle \text{kill\_man}, 21y \leq \rangle$
- $a' = \langle \text{kill\_male}, 21y \leq \rangle$
- $b = \langle \text{kill\_person}, 21y \leq \rangle$
- $c = \langle \text{person}, \text{equal\_status} \rangle$
- $c' = \langle \text{citizen}, \text{equal\_status} \rangle$
- $d = \langle \text{mother}, \text{kill\_newborn}, \neg 21y \leq \rangle$

$$\begin{aligned}
 F &= \{ \text{kill\_adult}, \text{kill\_female}, \text{Obl } l_3(\text{art.3}, c', c), \\
 &\quad \text{Obl } l_1(\text{art.578}, d, d), \neg \text{capable} \} \\
 R &= \{ r_5 : \text{kill\_adult}, \text{kill\_female}, \text{Obl } l_3(\text{art.3}, c', c) \Rightarrow^I \\
 &\quad l_2(\text{art.575}, a, b) \otimes l_3(\text{art.575}, a, b) \otimes \neg l_1(\text{art.575}, a, a), \\
 r_6 &: \text{Obl } l_1(\text{art.578}, d, d) \Rightarrow^I l_1(\text{art.575}, a, a) \otimes l_4(\text{art.575}, a, a), \\
 r_7 &: \Rightarrow^I l_1(\text{art.575}, a, a'), \\
 r_8 &: b \rightarrow^M a, \quad r_9 : \text{kill\_female} \rightarrow^M \text{kill\_person}, \\
 \text{art.575} &: \text{kill\_man} \Rightarrow^{\text{Obl}} 21y \leq, \\
 \text{art.85} &: \neg \text{capable} \rightsquigarrow^{\text{Obl}} \neg 21y \leq \} \\
 &>= \{ r_6 > r_5 \}
 \end{aligned}$$

The theory above extends the one in Example 3 by adding (i)  $r_9$ , which clarifies that killing any female human is killing a person, (ii)  $\text{art.575}$ , which is the obligation rule encoding the logical structure of  $\text{art. 575}$  as stated in the Italian penal code, (iii)  $\text{art.85}$ , which exhibits the logical content of  $\text{art. 85}$  of the Italian penal code, a provision that prevents to prove that one is punishable for homicide if incapable.

Proofs of  $A$ - and  $O$ -interpretations are exactly as explained in Example 4. The novelty is that provability of interpretations may lead to new usable versions of a same provision, thus affecting what obligatory literals may hold. Trivially, since  $\text{Obl } l_3(\text{art.3}, c', c)$  and  $\text{Obl } l_1(\text{art.578}, d, d)$  are in  $R$ , we have that  $+\partial^R \text{art.3} : \text{person} \Rightarrow^{\text{Obl}} \text{equal\_status}$  and  $+\partial^R \text{art.578} : \text{mother}, \text{kill\_newborn} \Rightarrow^{\text{Obl}} \neg 21y \leq$ . Also, (a)  $+\partial_{\text{Adm}}^I l_2(\text{art.575}, a, b)$ , which makes unusable the original version of  $\text{art.575} \in R$ , and (b)  $+\partial_{\text{Obl}}^I l_2(\text{art.575}, a, b)$ , which supports the provability of  $+\partial^R \text{art.575} : \text{kill\_person} \Rightarrow^{\text{Obl}} 21y \leq$ . The new version of  $\text{art. 575}$  is usable: since we obtain  $+\Delta^M \text{kill\_person}$  from  $\text{kill\_female}$  via  $r_9$ , we would have an argument supporting the derivation of  $+\partial^{\text{Obl}} 21y \leq$  from  $\text{kill\_female}$  (which is not possible with the original version of  $\text{art. 575}$ ). However, this is not the case, because  $\text{art. 85}$  is applicable and attacks the conclusion of  $\text{art. 575}$ , thus leading to  $-\partial^{\text{Obl}} 21y \leq$ .

<sup>9</sup>For the sake of simplicity, let us assume that the consequent of  $\text{art. 578}$  is  $\neg 21y \leq$  (at least 21 years prison), which is implied by any imprisonment between 4 and 12 years.

## 4. SUMMARY AND DISCUSSION

This paper presented a logical machinery for reasoning about interpretive canons, which is based on the following intuitions: (a) canons are represented by defeasible rules; (b) different reasoning patterns can be identified depending on whether we work on interpretations as activities or as outcomes [13]; (c) competing interpretive options can be handled by stating a priority over conflicting rules, but different ranking preferences can also be introduced among compatible interpretive acts; (d) canons are defeasible rules licensing deontic interpretive claims; (e) the logic can deal with the interpretation of abstract, non-analysed provisions and of structured provisions.

It is not hard to show that the proposed framework enjoys some relevant properties (see [7]).

**DEFINITION 25.** *An Interpretation Theory  $D = (F, R, >)$  is consistent iff  $>$  is acyclic and  $F$  does not contain pairs of complementary (modal) literals and interpretations.*

**PROPOSITION 26.** *Let  $D$  be a consistent Interpretation Theory, and  $\square \in \text{MOD}$ . For any interpretation  $\phi$  or literal  $l$ , it is not possible to have (a) both  $D \vdash +\partial_{\square}^I \phi$  and  $D \vdash -\partial_{\square}^I \phi$ , (b)  $D \vdash +\partial_{\square}^I l$  and  $D \vdash -\partial_{\square}^I l$ , (c) both  $D \vdash +\partial_{\square}^I \phi$  and  $D \vdash +\partial_{\square}^I \sim \phi$ , (d)  $D \vdash +\partial_{\square}^I l$  and  $D \vdash +\partial_{\square}^I \sim l$ .*

**PROPOSITION 27.** *Let  $D$  be a consistent Interpretation Theory. For any A- and O-interpretation  $\phi$  and  $\square \in \text{MOD}$ :*

1. if  $D \vdash +\partial_{\text{Obl}}^I \phi$ , then  $D \vdash -\partial_{\square}^I \sim \phi$ ;
2. if  $D \vdash +\partial_{\text{Adm}}^I \phi$ , then  $D \vdash -\partial_{\text{Obl}}^I \sim \phi$ ;
3.  $D \vdash +\partial_{\text{Obl}}^I \phi$ , then  $D \vdash +\partial_{\text{Adm}}^I \phi$ ;
4. if  $D \vdash +\partial_{\text{Obl}}^I l_i(n, a)$ , then  $D \vdash -\partial_{\text{Adm}}^I l_j(n, a')$ ,<sup>10</sup>
5. if  $D \vdash +\partial_{\square}^I l_i(n, a)$ , then  $D \vdash +\partial_{\square}^I a$ .

If we assume that, given an interpretive theory  $D$ , the set of conflicting interpretation rules has been defined, the following holds for Option 1:<sup>11</sup>

**THEOREM 28** (COMPLEXITY, OPTION 1; CF. [7]). *Let  $D = (F, R, >)$  be an Interpretation Theory and  $U^D$  be the set of all the atomic literals and atomic A-interpretations occurring in  $D$ . The set of conclusions of  $D$  can be computed in time linear to the size of the theory, i.e.,  $O(|R| * |U^D|)$ .*

We leave to future research to study the complexity of Option 2, where usable obligation rules can be proved.

Finally, a valuable aspect of the proposed machinery is that it can accommodate different doctrinal views regarding legal interpretation. In particular, we argued that two different but non-conflicting interpretations of the same provision can be admissible. First of all, we must notice that a number of logical relations can hold between meanings [3, 4]. All these aspects can be modularly introduced in the language definition. Also, one may argue that any interpretations  $l_i(n, a)$  and  $l_j(n, a')$  of  $n$ —if we just consider Option 1 of this paper—are in conflict [14] whenever  $a \neq a'$ , irrespective

<sup>10</sup>Similar results hold for Option 2.

<sup>11</sup>The relation  $\sqsubseteq$  requires to compute the monotonic part of the theory using, as facts, all sets in the powerset of the set of literals occurring in the antecedents of meaning rules.

of any possible logical relation (entailment, semantic overlapping, etc.) between  $a$  and  $a'$ . Suppose we accept this last view. Then, we would simply have to amend condition 5 of Definition 1 and modularly modify proof theory: if  $l_i(n, a)$  and  $l_j(n, a')$  are always in conflict, then Definitions 11-12 and 15-16 no longer apply. If so, a different way of proving as obligatory an A- and O-interpretation  $\phi$  could be based on checking whether  $\phi$  is provable for any set  $F$  of facts in the theory or it rather depends on specific facts: in the first case  $\phi$  would be obligatory, in the second we would just have a case of admissibility. A systematic investigation on these options is left for the future research.

## References

- [1] R. Alexy and R. Dreier. Statutory interpretation in the Federal Republic of Germany. In MacCormick and Summers [11].
- [2] G. Antoniou, D. Billington, G. Governatori, and M. Maher. Representation results for defeasible logic. *ACM Trans. Comput. Log.*, 2(2):255–287, 2001.
- [3] M. Araszkievicz. Towards systematic research on statutory interpretation in AI and law. In *Proc. JURIX 2013*. IOS Press, 2013.
- [4] M. Araszkievicz. Scientia juris: A missing link in the modelling of statutory reasoning. In *Proc. JURIX 2014*. IOS Press, 2014.
- [5] G. Boella, G. Governatori, A. Rotolo, and L. van der Torre. A logical understanding of legal interpretation. In *Proc. KR 2010*. AAAI Press, 2010.
- [6] B. Brozek. Legal interpretation and coherence. In M. Araszkievicz and J. Savelka, editors, *Coherence: Insights from Philosophy, Jurisprudence and Artificial Intelligence*. Springer, 2013.
- [7] G. Governatori, F. Olivieri, A. Rotolo, and S. Scanapiecico. Computing strong and weak permissions in defeasible logic. *J. Philosophical Logic*, 42(6):799–829, 2013.
- [8] G. Governatori and A. Rotolo. BIO logical agents: Norms, beliefs, intentions in defeasible logic. *Auton. Agent Multi Agent Syst.*, 17(1):36–69, 2008.
- [9] G. Governatori and A. Rotolo. Changing legal systems: Legal abrogations and annulments in defeasible logic. *Logic Journal of the IGPL*, 18:157–194, 2010.
- [10] F. Macagno, G. Sartor, and D. Walton. Argumentation schemes for statutory interpretation. In *Proc. ARGUMENTATION 2012*. Masaryk University, 2012.
- [11] D. MacCormick and R. Summers, editors. *Interpreting Statutes: A Comparative Study*. Ashgate, 1991.
- [12] H. Prakken and G. Sartor. Formalising arguments about norms. In *Proc. JURIX 2013*. IOS Press, 2013.
- [13] A. Ross. *On Law and Justice*. Stevens, London, 1958.
- [14] G. Sartor, D. Walton, F. Macagno, and A. Rotolo. Argumentation schemes for statutory interpretation: A logical analysis. In *Proc. JURIX 2014*. IOS Press, 2014.
- [15] G. Tarello. *L'interpretazione della legge*. Giuffrè, 1980.