

# A Probabilistic Deontic Argumentation Framework

Régis Riveret

*Commonwealth Scientific and Industrial Research Organisation, Australia*

Nir Oren

*University of Aberdeen, United Kingdom*

Giovanni Sartor

*University of Bologna, Italy*

---

## Abstract

What does it mean that something is probably obligatory? And how does it relate to the probability that it is permitted or prohibited? In this paper, we provide a possible answer by merging deontic argumentation and probabilistic argumentation into a probabilistic deontic argumentation framework. This framework allows us to specify a semantics for the probability of deontic statuses. The deontic argumentation part builds on standard concepts from the study of computational models of argument: rule-based arguments, argumentation graphs, argument labelling semantics and statement labelling semantics. We then encapsulate this deontic composition with the approach of probabilistic labellings to probabilistic argumentation, in order to associate deontic statements with probability values. The framework is illustrated with a scenario featuring a violation and a contrary-to-duty obligation.

*Keywords:* Probabilistic argumentation, deontic argumentation.

---

## 1. Introduction

Research in the area of argumentation investigates how conclusions can be drawn from a set of conflicting arguments and how such arguments can be constructed from an inconsistent and incomplete knowledge base. Argumentation has become a well-established approach to defeasible reasoning by virtue of its perceived simplicity and explanatory power. Because of its logic-based underpinnings, argumentation theory has traditionally approached uncertainty from a qualitative perspective. However, recent work has begun integrating probabilistic quantitative aspects into the argumentation process, see e.g. [1, 2, 3]. Typically, arguments are associated with probability values, and quantitative relationships amongst these arguments are studied to add a probabilistic dimension to defeasible reasoning as addressed by argumentation.

It has long been recognised that normative and deontic reasoning is defeasible [4] or non-monotonic [5]. Capturing such reasoning within an argumentative context could – potentially – provide an enhanced understanding of it and foster useful applications, and consequently,

normative or deontic argumentation frameworks have attracted increasing attention, see e.g. [6, 7, 8, 9].

Integrating probabilistic uncertainty into deontic argument-based reasoning may be useful for various applications. For example, a probabilistic deontic argumentation framework may provide semantics for legal expert systems associating normative effects with probability measures. It may also be useful to norm-governed cyber-physical systems where sensors provide uncertain readings as to the state of the environment, and probability measures are attached to system's normative states. Multiple applications can be conceived.

In general, probabilistic reasoning upon deontic concepts is important for systems where norm-related statements have an associated probabilistic uncertainty. As a running example, we will consider the following simple scenario (a follow-up from [9], inspired by the Hart-Fuller debate [10, 11]).

**Example 1.** *Vehicles are prohibited from entering a park unless there is an emergency. At some point, sensors suggest that a vehicle has entered the park, while other sensors indicate that an emergency situation may be occurring. Given this situation, we then seek to compute the likelihood that the prohibition is violated; doing so would allow us — for example — to decide whether to dispatch a law enforcement agent to the park or not.*

In the example above, probabilities are associated with the facts matching the antecedents of conditional norms, a prohibition and an exception to it, and one is interested in the probability that the prohibition is complied with or violated. In turn, such measures can then be used to obtain descriptive or predictive insights into the system. If no proper probabilistic framework is adopted, then outcomes and uncertainty measures may be inconsistent or incoherent.

In this paper we take a first step in combining probability theory and deontic argumentation, investigating how probabilistic argumentation can be coupled with normative and deontic reasoning. Our approach builds on the notion of probabilistic labellings [3] and we show how such labellings can be applied to the types of deontic argumentation frameworks described in [9] where pieces of doctrine are expressed within defeasible theories.

The combination of deontic and probabilistic notions is certainly not new, and there exist various works where probabilistic concepts are used to capture different aspects of normative notions. In the utilitarian tradition, for example, an action is considered obligatory if it has the highest expected aggregate utility when compared to the available alternatives, and this view has influenced fields as diverse as ethics, law, and economics. Our focus here is different: we concentrate on normative conclusions, i.e. on the normative statuses (unconditional obligations, permissions, or other normative properties) that are established by a given set of norms. We seek to determine how likely it is that such normative statuses are triggered by facts matching the antecedents of conditional norms. With regard to obligations or prohibitions we also examine how likely it is that they are violated or rather complied with.

In this work we ignore dynamic and practical aspects of normative reasoning, considering only how uncertain knowledge impacts on normative concepts. In other words, we consider only epistemic issues, ignoring, for example, how to select specific actions to pursue by utilising probabilistic argumentation [12]. We believe that determining the likelihood of normative statuses is an important subject on its own, and intend to investigate the links with — for example — practical decision making as part of future work.

Our work builds on probabilistic argumentation, which has been applied to several aspects of legal and normative domains. For example, probability theory has been combined with

argument and scenario approaches to evidential reasoning [1]. In another line of work, early probabilistic argumentation frameworks (see e.g. [13]) were used together with game theory to identify optimal strategies in dialogue games [14]. More recent research has coupled probabilistic argumentation with reinforcement learning to determine what actions to perform in environments featuring basic conditional obligations [12]. However, to the best of our knowledge, probabilistic argumentation has never been integrated with any proper deontic argumentation framework, and it is precisely this gap we seek to address in the work reported here.

The remainder of the paper is structured as follows. In the next section we outline a labelling-based approach to probabilistic argumentation and describe its semantics. In Section 3 we describe how probabilistic labellings can be used to reason about deontic concepts. Section 4 provides a detailed example of our approach, after which we conclude.

## 2. Probabilistic Argumentation Framework

In this section we present a simple ASPIC-like argumentation system (cf. e.g. [15, 16]) and its probabilistic development [3]. We begin by defining the language for expressing rules, and then describe how arguments, attack and support relations can be constructed. Given a set of arguments, we identify which are accepted, and in turn which statements are acceptable. Acceptance is specified through labellings, which are then integrated into a probabilistic setting.

### 2.1. Language

Our language is propositional, with literals being constructed from propositions, negations of propositions, and modal operators (which are used later to capture deontic modalities).

**Definition 2.1.** A *literal statement* is either a plain literal statement or a modal literal statement, where

- a *plain literal statement* is either an atomic proposition  $p$  or the negation of an atomic proposition (i.e.  $\neg p$ ) and
- a *modal literal statement* is a statement of the form  $\Box\gamma$  or  $\neg\Box\gamma$ , such that  $\Box$  is a placeholder for any modal operator and  $\gamma$  is a plain literal statement.

**Notation 2.1.** For notational convenience we write the complement of a literal statement  $\gamma$  as  $\bar{\gamma}$ . That is, if  $\gamma$  is of the form  $p$ , then  $\bar{\gamma}$  is  $\neg p$ , while if  $\gamma$  is  $\neg p$ , then  $\bar{\gamma}$  is  $p$ .

Defeasible rules constructed using literals represent conditionals of the form ‘if... then ... unless...’. While many argumentation formalisms distinguish between strict and defeasible rules, for simplicity we consider only defeasible rules in this work<sup>1</sup>.

**Definition 2.2.** A *defeasible rule* over a set of literal statements  $\mathcal{S}$  is a construct of the form:

$$r : \varphi_1, \dots, \varphi_n, \sim \varphi'_1, \dots, \sim \varphi'_m \Rightarrow \varphi$$

---

<sup>1</sup>The translation of our results to strict rules can be achieved relatively simply through techniques similar to those proposed by Li and Parsons [17].

where  $0 \leq n$  and  $0 \leq m$ ,  $r$  is a unique identifier for the rule, while for any  $0 \leq i \leq n$ ,  $0 \leq j \leq m$ ,  $\varphi_i, \varphi'_j, \varphi \in \mathcal{S}$  are all literal statements.

A rule  $r : \varphi_1, \dots, \varphi_n, \sim \varphi'_1, \dots, \sim \varphi'_m \Rightarrow \varphi$  can be read as follows: “if  $\varphi_1$  and ... and  $\varphi_n$  are supported, and all of  $\varphi'_1, \dots, \varphi'_m$  are not supported, then  $\varphi$  is defeasibly supported. While premises prefixed with  $\sim$  can be interpreted as a form of negation as failure, we view such premises as exceptions to the (default) application of the rule.

**Notation 2.2.** Let  $r$  be a defeasible rule  $r$  as in Definition 2.2, and  $Rules$  a set of rules.

- $Body(r)$  denotes the body of  $r$ , i.e.  $Body(r) = \{\varphi_1, \dots, \varphi_n, \sim \varphi'_1, \dots, \sim \varphi'_m\}$ .
- $Head(r)$  denotes the head of  $r$ , i.e.  $Head(r) = \{\varphi\}$ .
- $Prop(r)$  denotes the set of propositions of  $r$ , i.e.  $Prop(r) = \{p \mid p, \neg p, \sim p, \sim \neg p, \Box p, \neg \Box p, \Box \neg p, \sim \Box p, \sim \neg \Box p, \sim \Box \neg p, \sim \neg \Box \neg p \in Body(r) \cup Head(r)\}$ .
- $Prop(Rules)$  denotes the set of propositions of  $Rules$ , i.e.  $Prop(Rules) = \bigcup_{r \in Rules} Prop(r)$ .

Certain sets of literals are mutually inconsistent. In other words, they conflict, or are incompatible with each other. Perhaps the simplest form of such a conflict is between a proposition and its negation, but more complex conflicts can also exist. As done in many other systems [18, 15], we introduce an abstract conflict relation to encode incompatibilities between literal statements.

**Definition 2.3.** A **conflict relation** ‘Conflicts’ over a set of literal statements  $\mathcal{S}$  is a binary relation over  $\mathcal{S}$ , i.e.  $Conflicts \subseteq \mathcal{S} \times \mathcal{S}$ .

**Notation 2.3.** We write  $Prop(Conflicts)$  to denote the propositions found within a conflict relation, i.e.  $Prop(Conflicts) = \{p \mid (\varphi, \varphi') \in Conflicts : \varphi = p, \neg p, \Box p, \neg \Box p, \Box \neg p, \neg \Box \neg p, \text{ or } \varphi' = p, \neg p, \Box p, \neg \Box p, \Box \neg p, \neg \Box \neg p\}$ .

The conflict relation should adhere to certain principles, in which case it is said to be *well-formed*. For example, for a purely propositional system, it should be the case that  $Conflicts(\gamma, \bar{\gamma})$ . For a deontic setting, we will specify the requirements for the conflict relation in Section 3. Note that the conflict relation can be asymmetric or symmetric to capture, for example, contrary or contradictory relationships [15], but we do not consider such aspects here.

Preferences over rules are commonly used to resolve conflicts. We capture such preferences via a *superiority* relation  $\succ$  over rules. Informally,  $s \succ r$  means that rule  $s$  prevails over rule  $r$ .

**Definition 2.4.** A **superiority relation**  $\succ$  over a set of rules  $Rules$  is an asymmetric binary relation over  $Rules$ , i.e.  $\succ \subseteq Rules \times Rules$ .

As the superiority relation is asymmetric, for any rule  $r$  it is not the case that  $r \succ r$ , and for two distinct rules  $r$  and  $r'$  we cannot have both  $r \succ r'$  and  $r' \succ r$ .

## 2.2. Defeasible theories and argumentation graphs

Having introduced rules, conflicts between literal statements and superiority relationships between rules, we now turn our attention to how rules can be combined to build arguments. Given that these arguments may conflict, we describe later how sets of accepted arguments can be identified.

A defeasible theory lists a set of rules, a conflict relation and a superiority relation. Such a defeasible theory forms the basic structure over which inferences can take place.

**Definition 2.5.** A *defeasible theory* is a tuple  $\langle Rules, Conflicts, \succ \rangle$  where

- *Rules* is a set of rules, and
- *Conflicts* is a conflict relation, and
- $\succ$  is a superiority relation over *Rules*.

**Notation 2.4.** Given a defeasible theory  $T = \langle Rules, Conflicts, \succ \rangle$ ,

- $Rules(T)$ ,  $Conflicts(T)$ , and  $\succ(T)$  denote the set of rules *Rules*, the relation *Conflicts*, and the relation  $\succ$  of theory  $T$  respectively,
- $Prop(T)$  denotes the set of propositions of  $T$ , i.e.  $Prop(T) = Prop(Rules) \cup Prop(Conflicts)$ .

Inferences over a defeasible theory are performed by chaining its defeasible rules to form arguments, which we define below. Our definition of arguments is inspired by, and similar to, those used in other rule-based argumentation frameworks, as described for example in [15, 16, 19].

**Definition 2.6.** An *argument*  $A$  constructed from a defeasible theory  $\langle Rules, Conflicts, \succ \rangle$  is a finite construct of the form:

$$A : A_1, \dots, A_n, \sim \varphi_1, \dots, \sim \varphi_m \Rightarrow_r \varphi$$

with  $0 \leq n$  and  $0 \leq m$ , and where

- $A$  is the unique identifier of the argument;
- $\varphi$  is the conclusion of the argument, denoted  $con(A)$ ;
- $A_1, \dots, A_n$  are arguments constructed from the defeasible theory  $\langle Rules, Conflicts, \succ \rangle$ ;
- $r \in Rules$  is the top rule of the argument, and it is of the form  $r : con(A_1), \dots, con(A_n), \sim \varphi_1, \dots, \sim \varphi_m \Rightarrow \varphi$ .

**Notation 2.5.** Given an argument  $A$  as in Definition 2.6, we use the following notations.

- $Sub(A)$  denotes the set of subarguments of  $A$ , i.e.  $Sub(A) = Sub(A_1) \cup \dots \cup Sub(A_n) \cup \{A\}$ .
- $DirectSub(A)$  denotes the direct subarguments of  $A$ , i.e.  $DirectSub(A) = \{A_1, \dots, A_n\}$ .
- $TopRule(A)$  denotes the top rule of  $A$ , i.e.  $TopRule(A) = (r : con(A_1), \dots, con(A_n), \sim \varphi_1, \dots, \sim \varphi_m \Rightarrow \varphi)$ .
- $Rules(A)$  denotes the rules used within  $A$ , i.e.  $Rules(A) = Rules(A_1) \cup \dots \cup Rules(A_n) \cup \{TopRule(A)\}$ .

We can remark that Definition 2.6 assumes that arguments are finite, i.e. any argument has a finite set of subarguments. Accordingly, any argument has exactly one conclusion and ‘bottoms out’ in arguments with an empty set of subarguments (taking the form  $A : \sim \varphi_1, \dots, \sim \varphi_m \Rightarrow \varphi$  with  $0 \leq m$ ). Nevertheless, for a given defeasible theory, we may have an infinite number of finite arguments.

Different types of inconsistencies can appear between arguments, causing them to *attack* each other. We consider two types of attacks between arguments: *rebuttals* (where two argument’s conclusions are incompatible), and *undercuts*<sup>2</sup> (where exceptions prevent an argument’s conclusion from being drawn). In the ASPIC family of argumentation frameworks, attack is differentiated from *defeat*, with the latter taking preferences between arguments into account. For simplicity, we make no such distinction in this work, instead integrating such

<sup>2</sup>Note that this term is overloaded in the argumentation literature, and is used with different meanings in different contexts [15].

preferences into our definition of attack. While diverse approaches to lifting preferences over rules to preferences over arguments have been described [6], in this work we utilise a simple last-link ordering [15] to compute preferences over arguments. More specifically, argument  $A$  is preferred to argument  $B$  (written  $A \succ B$ ) iff  $\text{TopRule}(A) \succ \text{TopRule}(B)$ .

**Definition 2.7.**

- An **attack relation**  $\rightsquigarrow$  over a set of arguments  $\mathcal{A}$  is a binary relation over  $\mathcal{A}$ , i.e.  $\rightsquigarrow \subseteq \mathcal{A} \times \mathcal{A}$ .
- Let  $A$  and  $B$  be two arguments constructed from a defeasible theory  $\langle \text{Rules}, \text{Conflicts}, \succ \rangle$ , argument  $A$  attacks  $B$ , i.e.  $(A, B) \in \rightsquigarrow$ , iff  $A$  rebuts or undercuts  $B$ , where
  - $A$  rebuts  $B$  (on  $B'$ ) iff  $\exists B' \in \text{Sub}(B)$  such that  $\text{Conflicts}(\text{con}(A), \text{con}(B'))$ , and  $B' \neq A$ ;
  - $A$  undercuts  $B$  (on  $B'$ ) iff  $\exists B' \in \text{Sub}(B)$  such that  $(\sim \text{con}(A)) \in \text{Body}(\text{TopRule}(B'))$ .

In what follows, it will be useful to identify which arguments are the immediate subarguments of other arguments, leading us to the following definition of *direct subargument relations*.

**Definition 2.8.**

- A **direct subargument relation**  $\Rightarrow$  over a set of arguments  $\mathcal{A}$  is a binary relation over  $\mathcal{A}$ , i.e.  $\Rightarrow \subseteq \mathcal{A} \times \mathcal{A}$ .
- Let  $A$  and  $B$  be two arguments constructed from a defeasible theory, argument  $B$  is a direct subargument of  $A$ , written  $B \Rightarrow A$ , iff  $B$  belongs to the set of direct subarguments of  $A$ , i.e.  $B \in \text{DirectSub}(A)$ .

Since an argument is not a direct subargument of itself and cannot be a subargument of its direct subarguments, the direct subargument relation over arguments constructed from a defeasible theory is antireflexive and acyclic. In the remainder of this work, for the sake of simplicity, we say that an argument  $A$  supports an argument  $B$  if  $A$  is a direct subargument of  $B$ . We note that other notions of support may be considered, but leave them to future work.

Given the arguments and attacks obtained from a defeasible theory, and by making use of the subargument relation, we can create an argumentation graph as defined below, cf. [20, 21, 22, 23, 24].

**Definition 2.9.** An **argumentation graph** constructed from a defeasible theory  $T$  is a tuple  $\langle \mathcal{A}, \rightsquigarrow, \Rightarrow \rangle$  where  $\mathcal{A}$  is the set of all arguments constructed from  $T$ ,  $\rightsquigarrow$  is an attack relation over  $\mathcal{A}$ , and  $\Rightarrow$  is a direct subargument relation over  $\mathcal{A}$ .

In the rest of the paper, we assume that all argumentation graphs are constructed from a defeasible theory, which may be left unspecified.

As to the terminology, an argumentation graph may be called a *bipolar* argumentation graph/framework [21], as long as the support relation is understood as a (direct) subargument relation, every argument has a conclusion which is a statement, and such bipolar graphs enjoy all the (probabilistic) characteristics and properties discussed in the remainder of this work.

**Notation 2.6.** Given an argumentation graph  $G = \langle \mathcal{A}, \rightsquigarrow, \Rightarrow \rangle$ , we may write  $\mathcal{A}_G, \rightsquigarrow_G$  and  $\Rightarrow_G$  to denote the graph's arguments  $\mathcal{A}$ , attacks  $\rightsquigarrow$  and subargument relation  $\Rightarrow$  respectively.

**Example 2.** An argumentation graph is illustrated in Figure 1, with solid arrows denoting attacks and hollow arrows representing supports. From our definitions, arguments  $A$  and  $B$  are direct subarguments of  $C$  and therefore provide conjunctive support to  $C$ .

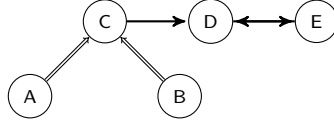


Figure 1: An argumentation graph. Arguments A and B support argument C. Argument C attacks D. Arguments D and E attack each other.

□

We note that given an argumentation graph  $G$  constructed from a defeasible theory  $T$ , since  $\mathcal{A}_G$  is the set of all arguments constructed from  $T$  then the relations  $\rightsquigarrow_G$  and  $\Rightarrow_G$  of the graph are defined over  $\mathcal{A}_G$ . Hence  $\rightsquigarrow_G$  is obtained according to Definition 2.7, and  $\Rightarrow_G$  complies with specification of Definition 2.8 on direct subargument relations.

Given an argumentation graph, we may wish to determine what conclusions can be drawn if some arguments are omitted. We may therefore wish to deal with subgraphs of an argumentation graph.

**Definition 2.10.** Let  $G$  denote an argumentation graph. The **subgraph**  $H$  of  $G$  induced by a set of arguments  $\mathcal{A}_H \subseteq \mathcal{A}_G$  is an argumentation graph such that  $H = \langle \mathcal{A}_H, \rightsquigarrow_G \cap (\mathcal{A}_H \times \mathcal{A}_H), \Rightarrow_G \cap (\mathcal{A}_H \times \mathcal{A}_H) \rangle$ .

**Notation 2.7.** Given an argumentation graph  $G$ , we denote all of its subgraphs as  $\text{Sub}(G)$ , i.e.  $\text{Sub}(G) = \{ \langle \mathcal{A}_H, \rightsquigarrow_G \cap (\mathcal{A}_H \times \mathcal{A}_H), \Rightarrow_G \cap (\mathcal{A}_H \times \mathcal{A}_H) \rangle \mid \mathcal{A}_H \subseteq \mathcal{A}_G \}$ .

Clearly, if an argument  $A$  appears within an argumentation graph, then any subgraph should contain not only  $A$ , but also  $A$ 's subarguments and any attacks related to these subarguments. Graphs obeying this property are said to be *subargument-complete*.

**Definition 2.11.** A subgraph  $H$  of an argumentation graph  $G$  induced by a set of arguments  $\mathcal{A}_H$  is **subargument-complete** iff for every argument  $A \in \mathcal{A}_H$ , if  $B \Rightarrow_G A$  then  $B \in \mathcal{A}_H$ .

**Example 2** (continued). The argumentation graphs in Figures 2 (a) and (b) are subgraphs of the argumentation graph in Figure 1, and they are induced by the set of arguments  $\{A, B, C, D\}$  and  $\{A, C, D\}$  respectively. The graph of Figure 2 (a) is a subargument-complete subgraph, whereas the subgraph in Figure 2 (b) is not.



Figure 2: Two subgraphs of the argumentation graph from Figure 1.

□

Through the definitions above, given a defeasible theory we are able to build arguments and identify attacks and subarguments across arguments, from which we can construct an argumentation graph. Next, we consider how the acceptance status of arguments and statements can be determined through the use of labellings.

### 2.3. Argument labelling

Given an argumentation graph, the selection of arguments considered acceptable or justified is performed on the basis of a formal specification traditionally called argumentation semantics. This evaluation can be carried out in terms of sets of arguments called extensions [25], or in terms of labellings [26]. For our purposes, we adopt the labelling approach.

Argument labelling approaches usually assign one of three labels (IN, OUT and UND) to arguments, reflecting their acceptance status with respect to a specific argumentation semantics. We extend such labellings to consider whether an argument is, or is not included when evaluating acceptance by introducing an additional OFF label. The addition of this label will allow us to distinguish between the probability of using an argument as part of the construction of an argumentation graph, and the probability of its acceptance.

Argument labellings are commonly formalised as follows.

**Definition 2.12.** Let  $G$  be an argumentation graph, and  $\text{ArgLab}$  a set of labels for arguments. An  $\text{ArgLab}$ -labelling of a set of arguments  $\mathcal{A} \subseteq \mathcal{A}_G$  is a total function  $L : \mathcal{A} \rightarrow \text{ArgLab}$ .

While argument labellings map sets of arguments to a label, we will abuse notation and may speak of the labelling of an argument. Labellings involving all the arguments of an argumentation graph play a special role and deserve a specific terminology and notation.

**Definition 2.13.** Let  $\text{ArgLab}$  be a set of labels for arguments. An  $\text{ArgLab}$ -labelling of an argument graph  $G$  is a total function  $L : \mathcal{A}_G \rightarrow \text{ArgLab}$ .

#### Notation 2.8.

- We write  $\mathcal{L}_{\text{ArgLab}}(\mathcal{A})$  to denote the universe of all possible  $\text{ArgLab}$ -labelling assignments to a set of arguments  $\mathcal{A}$ , and we write  $\mathcal{L}_{\text{ArgLab}}(G)$  for the set of all possible  $\text{ArgLab}$ -labellings of an argumentation graph  $G$ .
- If a labelling  $L$  assigns the label  $l$  to a set of arguments  $A$ , we may write  $l(L)$ . In other words,  $l(L) = \{A \mid L(A) = l\}$ . Thus for example, if IN is a label, then  $\text{IN}(L) = \{A \mid L(A) = \text{IN}\}$ .

As mentioned above, standard argument labellings make use of  $\{\text{IN}, \text{OUT}, \text{UND}\}$ -labellings. Here, each argument is associated with a label representing its status with regards to an argumentation semantics [26]. Intuitively, an IN labelling means that the argument is accepted, while an OUT labelling indicates its rejection. An UND labelling states that the status of an argument is undecided, i.e. it is neither accepted nor rejected. Different labelling functions may label arguments in different ways, and a single labelling function may label the same graph in multiple ways. In this work, we focus on complete labellings and from these, adopt grounded labellings.

**Definition 2.14.** A **complete**  $\{\text{IN}, \text{OUT}, \text{UND}\}$ -labelling of an argumentation graph  $G$  is a  $\{\text{IN}, \text{OUT}, \text{UND}\}$ -labelling such that for every argument  $A$  in  $\mathcal{A}_G$ :

- $A$  is labelled IN if, and only if, all attackers of  $A$  are labelled OUT, and
- $A$  is labelled OUT if, and only if,  $A$  has an attacker labelled IN.

**Definition 2.15.** A **grounded**  $\{\text{IN}, \text{OUT}, \text{UND}\}$ -labelling  $L$  of an argumentation graph  $G$  is a complete  $\{\text{IN}, \text{OUT}, \text{UND}\}$ -labelling of  $G$  such that  $\text{IN}(L)$  is minimal (with respect to set inclusion) among all complete  $\{\text{IN}, \text{OUT}, \text{UND}\}$ -labellings of  $G$ .



**Example 2** (continued). The grounded  $\{\text{IN}, \text{OUT}, \text{UND}\}$ -labelling of the argumentation graph is illustrated in Figure 3.

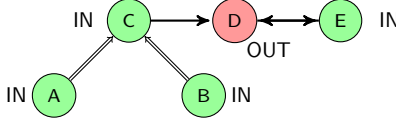


Figure 3: Argumentation graph and its grounded  $\{\text{IN}, \text{OUT}, \text{UND}\}$ -labelling.

□

For any argumentation graph, the grounded labelling of the graph exists and is unique. Moreover, for any finite argumentation graph, it can be computed in polynomial time using for example Algorithm 1 [27]. The algorithm begins by labelling IN all arguments not being attacked or whose attackers are OUT (line 4), and then it iteratively labels OUT any argument attacked by an argument labelled IN (line 5). The iteration continues until no more arguments can be labelled IN or OUT (line 6); and it terminates by labelling UND all unlabelled arguments (line 7).

---

**Algorithm 1** Computation of a grounded  $\{\text{IN}, \text{OUT}, \text{UND}\}$ -labelling.

---

- 1: **input** A finite argumentation graph  $G$ ,
  - 2:  $L_0 = \langle \emptyset, \emptyset, \emptyset \rangle$ ,
  - 3: **repeat**
  - 4:  $\text{IN}(L_{i+1}) \leftarrow \text{IN}(L_i) \cup \{A \mid A \in \mathcal{A}_G \text{ is not labelled in } L_i, \text{ and } \forall B \in \mathcal{A}_G : \text{if } B \text{ attacks } A \text{ then } B \in \text{OUT}(L_i)\}$
  - 5:  $\text{OUT}(L_{i+1}) \leftarrow \text{OUT}(L_i) \cup \{A \mid A \in \mathcal{A}_G \text{ is not labelled in } L_i, \text{ and } \exists B \in \mathcal{A}_G : B \text{ attacks } A \text{ and } B \in \text{IN}(L_{i+1})\}$
  - 6: **until**  $L_i = L_{i+1}$
  - 7: **return**  $\langle \text{IN}(L_i), \text{OUT}(L_i), \mathcal{A}_G \setminus (\text{IN}(L_i) \cup \text{OUT}(L_i)) \rangle$
- 

We can extend the grounded  $\{\text{IN}, \text{OUT}, \text{UND}\}$ -labelling semantics through the introduction of an OFF label. Intuitively, OFF labelled arguments are not evaluated when computing argument acceptance.

**Definition 2.16.** Let  $H$  be a subargument-complete subgraph of an argumentation graph  $G$ . A **grounded**  $\{\text{IN}, \text{OUT}, \text{UND}, \text{OFF}\}$ -labelling of  $G$  with respect to  $H$  is a  $\{\text{IN}, \text{OUT}, \text{UND}, \text{OFF}\}$ -labelling of  $G$  such that:

- every argument in  $\mathcal{A}_H$  is labelled according to the grounded  $\{\text{IN}, \text{OUT}, \text{UND}\}$ -labelling of  $H$ , and
- every argument in  $\mathcal{A}_G \setminus \mathcal{A}_H$  is labelled OFF.

An argumentation graph  $G$  has a unique grounded  $\{\text{IN}, \text{OUT}, \text{UND}\}$ -labelling, but it has as many grounded  $\{\text{IN}, \text{OUT}, \text{UND}, \text{OFF}\}$ -labellings as there are subargument-complete subgraphs of  $G$ .

**Example 2** (continued). A grounded  $\{\text{IN}, \text{OUT}, \text{UND}, \text{OFF}\}$ -labelling of the argumentation graph is illustrated in Figure 4.

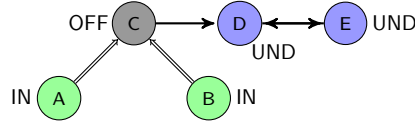


Figure 4: A grounded {IN,OUT,UND,OFF}-labelling.

□

Argument labellings which are equivalent to an argumentation semantics are also called *argument acceptance labellings*. Given a set of labels from a set of argument acceptance labellings, one can compute justification statuses for arguments to obtain an *argument justification labelling*. From an argument justification labelling, individual statements can then be labelled, yielding a *statement justification labelling*. However, one can also directly label statements from argument acceptance labellings, and it is this latter approach the one we use in this paper.

#### 2.4. Statement labelling

Given a set of statements, a *statement labelling* of this set is a (preferably total) function associating any statement with a label. Different specifications for statement labellings are possible [28, 29], but for our purposes, we consider what is perhaps the simplest meaningful labelling, namely the bivalent labelling semantics. Under this labelling, a statement is either accepted or not accepted. If accepted, then the statement is labelled ‘in’, otherwise it is labelled ‘ni’.

As statements are labelled via argument acceptance labellings, the resultant statement labellings are acceptance bivalent {in, ni}-labellings, but we may simply call them bivalent {in, ni}-labellings.

**Definition 2.17.** Let  $\mathcal{L}$  be a set of {IN, OUT, UND, OFF}-labellings,  $\mathcal{S}$  a set of literal statements. A **bivalent** {in, ni}-labelling of  $\mathcal{S}$  from  $\mathcal{L}$  is a total function  $K : \mathcal{L}, \mathcal{S} \rightarrow \{\text{in}, \text{ni}\}$  such that for any argument labelling  $L \in \mathcal{L}$  and any statement  $\varphi \in \mathcal{S}$ :

- $K(L, \varphi) = \text{in}$  iff  $\text{IN} \in \{L(A) \mid \text{con}(A) = \varphi\}$ , and
- $K(L, \varphi) = \text{ni}$  otherwise.

Bivalent labellings cannot differentiate between a statement being deemed unjustified and undecidable (the latter by being the conclusion of UND labelled arguments). Trivalent labellings [28] are a simple extension of bivalent labellings which include undecidability.

**Definition 2.18.** Let  $\mathcal{L}$  be a set of {IN, OUT, UND, OFF}-labellings,  $\mathcal{S}$  a set of literal statements. A **trivalent** {in, und, niund}-labelling of  $\mathcal{S}$  and  $\mathcal{L}$  is a total function  $K : \mathcal{L}, \mathcal{S} \rightarrow \{\text{in}, \text{und}, \text{niund}\}$  such that for any argument labelling  $L \in \mathcal{L}$  and any statement  $\varphi \in \mathcal{S}$ :

- $K(L, \varphi) = \text{in}$  iff  $\text{IN} \in \{L(A) \mid \text{con}(A) = \varphi\}$ , and
- $K(L, \varphi) = \text{und}$  iff  $\text{UND} \in \{L(A) \mid \text{con}(A) = \varphi\}$  and  $\text{IN} \notin \{L(A) \mid \text{con}(A) = \varphi\}$ , and
- $K(L, \varphi) = \text{niund}$  otherwise.

**Notation 2.9.** A bivalent {in, ni}-labelling or trivalent {in, und, niund}-labelling  $K$  may be represented as a tuple  $\langle \text{in}(K), \text{ni}(K) \rangle$  or  $\langle \text{in}(K), \text{und}(K), \text{niund}(K) \rangle$  respectively, with the obvious meaning.

The distinction between bivalent and trivalent labellings is later exploited in our deontic setting. For now we note the following proposition.

**Proposition 2.1.** *Let  $L$  be a grounded  $\{\text{IN}, \text{OUT}, \text{UND}, \text{OFF}\}$ -labelling of an argumentation graph constructed from a defeasible theory  $T$ ,  $\mathcal{S}$  a set of literal statements.*

- *Let  $K$  be a bivalent  $\{\text{in}, \text{ni}\}$ -labelling of  $\mathcal{S}$  and from  $\{L\}$ . For any statement  $\varphi_1, \varphi_2 \in \mathcal{S}$ , such that  $(\varphi_1, \varphi_2) \in \text{Conflicts}(T)$ , if  $K(L, \varphi_1) = \text{in}$  then  $K(L, \varphi_2) = \text{ni}$ .*
- *Let  $K$  be a trivalent  $\{\text{in}, \text{und}, \text{niund}\}$ -labelling of  $\mathcal{S}$  and from  $\{L\}$ . For any statement  $\varphi_1, \varphi_2 \in \mathcal{S}$ , such that  $(\varphi_1, \varphi_2) \in \text{Conflicts}(T)$ , if  $K(L, \varphi_1) = \text{in}$  then  $K(L, \varphi_2) = \text{niund}$ .*

This proposition is used to show later some results in our probabilistic deontic investigation.

### 2.5. Probabilistic labellings

To extend our work to the probabilistic setting we take the approach of probabilistic labellings according to which, given an argumentation graph, specific sets of argument labellings are sample spaces [3]. For our purposes, we will focus on grounded labellings and thus we will work with so-called *grounded probabilistic labelling frames*.

**Definition 2.19.** *A grounded probabilistic labelling frame (or grounded PLF) based on an argumentation graph  $G$  is a tuple  $\langle G, \langle \Omega, \mathcal{F}, P \rangle \rangle$  where  $\langle \Omega, \mathcal{F}, P \rangle$  is a probability space such that:*

- *the sample space  $\Omega$  is the set of grounded  $\{\text{IN}, \text{OUT}, \text{UND}, \text{OFF}\}$ -labellings of  $G$ ;*
- *the  $\sigma$ -algebra  $\mathcal{F}$  is the power set of  $\Omega$ , i.e.  $\mathcal{F} = 2^\Omega$ ;*
- *the function  $P$  from  $\mathcal{F}$  to  $[0, 1]$  is a probability measure (or probability distribution) satisfying Kolmogorov axioms.*

**Example 2** (continued). *Let us consider the grounded PLF based on the argumentation graph given in Figure 1 along with a uniform probability distribution. The sample space and the distribution are illustrated in Figure 5.* □

A grounded PLF defines a probability space, and thus we can work with random variables (written using upper case letters such as  $X, Y$  or  $Z$ ) from  $\Omega$  to another set of elements. Therefore, we introduce a categorical random variable, which we refer to as a *random labelling*. The random labelling for an argument  $A$ , denoted  $L_A$  maps from  $\Omega$  to the set of labels  $\{\text{IN}, \text{OUT}, \text{UND}, \text{OFF}\}$ . We may write, for example,  $L_A = \text{IN}$  as shorthand for the outcomes  $\{L \in \Omega \mid L(A) = \text{IN}\}$  (i.e. capturing those outcomes where  $A$  is labelled IN). We can follow a similar approach to the labelling of statements. That is, for any statement  $\varphi$ , we can introduce a categorical random variable  $K_\varphi$  which maps  $\Omega$  to the set of labels  $\{\text{in}, \text{ni}\}$  or  $\{\text{in}, \text{und}, \text{niund}\}$  (depending on whether we consider bivalent or trivalent labellings). Hence,  $P(K_\varphi = l) = 1$ , for example, is interpreted as ‘the probability that statement  $\varphi$  is labelled  $l$  equals one’, while  $P(K_{\varphi'} = l' \mid K_\varphi = l) = 1$  is interpreted as ‘the probability that statement  $\varphi'$  is labelled  $l'$ , given that statement  $\varphi$  is labelled  $l$ , equals one’. In the remainder of this paper, we assume that all specified conditional probabilities are not undefined (a conditional probability  $P(A \mid B)$  is undefined if  $P(B) = 0$ ).

**Example 2** (continued). *Referring to the probability distribution in Figure 5, we can easily compute that, for example,  $P(L_A = \text{IN}) = 12/20$ ,  $P(L_A = \text{OFF}) = 8/20$  or  $P(L_E = \text{IN} \mid L_D = \text{OUT}) = 1/2$ .* □

A	B	C	D	E	$P(\cdot)$
IN	IN	IN	OUT	IN	1/20
IN	IN	IN	OUT	OFF	1/20
IN	IN	IN	OFF	IN	1/20
IN	IN	IN	OFF	OFF	1/20
IN	IN	OFF	UND	UND	1/20
IN	IN	OFF	IN	OFF	1/20
IN	IN	OFF	OFF	IN	1/20
IN	IN	OFF	OFF	OFF	1/20
IN	OFF	OFF	UND	UND	1/20
IN	OFF	OFF	IN	OFF	1/20
IN	OFF	OFF	OFF	IN	1/20
IN	OFF	OFF	OFF	OFF	1/20
OFF	IN	OFF	UND	UND	1/20
OFF	IN	OFF	IN	OFF	1/20
OFF	IN	OFF	OFF	IN	1/20
OFF	IN	OFF	OFF	OFF	1/20
OFF	OFF	OFF	UND	UND	1/20
OFF	OFF	OFF	IN	OFF	1/20
OFF	OFF	OFF	OFF	IN	1/20
OFF	OFF	OFF	OFF	OFF	1/20

Figure 5: Sample space and probability distribution.

This probabilistic argumentation framework has multiple properties, see [3] for some of them. For our purposes, let us consider the following proposition and corollary.

**Proposition 2.2.** *Let  $\langle G, \langle \Omega, \mathcal{F}, P \rangle \rangle$  be a grounded PLF. Let  $L$  denote any grounded  $\{\text{IN}, \text{OUT}, \text{UND}, \text{OFF}\}$ -labelling such that  $L \in \Omega$ ,  $\mathcal{S}$  a set of literal statements, and  $K$  a bivalent  $\{\text{in}, \text{ni}\}$ -labelling or trivalent  $\{\text{in}, \text{und}, \text{niund}\}$ -labelling of  $\mathcal{S}$  and from  $\{L\}$ . For any  $\varphi, \varphi' \in \mathcal{S}$  such that if  $K(L, \varphi) = l$  then  $K(L, \varphi') = l'$ ,*

$$P(K_{\varphi'} = l' \mid K_{\varphi} = l) = 1.$$

**Corollary 2.1.** *Let  $\langle G, \langle \Omega, \mathcal{F}, P \rangle \rangle$  be a grounded PLF. Let  $L$  denote any grounded  $\{\text{IN}, \text{OUT}, \text{UND}, \text{OFF}\}$ -labelling such that  $L \in \Omega$ ,  $\mathcal{S}$  a set of literal statements, and  $K$  a bivalent  $\{\text{in}, \text{ni}\}$ -labelling or trivalent  $\{\text{in}, \text{und}, \text{niund}\}$ -labelling of  $\mathcal{S}$  and from  $\{L\}$ . For any  $\varphi, \varphi' \in \mathcal{S}$  such that if  $K(L, \varphi) = l$  then  $K(L, \varphi') = l'$ ,*

$$P(K_{\varphi} = l) \leq P(K_{\varphi'} = l').$$

We will also use the proposition and corollary below [3] when later studying the probabilistic relationships amongst deontic statements.

**Proposition 2.3.** *Let  $\langle G, \langle \Omega, \mathcal{F}, P \rangle \rangle$  be a grounded PLF where  $G$  is constructed from a defeasible theory  $T = \langle \text{Rules}, \text{Conflicts}, \succ \rangle$ ,  $\mathcal{S}$  a set of literal statements. For any  $\varphi_1, \varphi_2 \in \mathcal{S}$  such that  $\text{Conflicts}(\varphi_1, \varphi_2)$ ,*

$$P(K_{\varphi_1} = \text{in}) + P(K_{\varphi_2} = \text{in}) \leq 1.$$

**Corollary 2.2.** *Let  $\langle G, \langle \Omega, \mathcal{F}, P \rangle \rangle$  be a grounded PLF where  $G$  is constructed from a defeasible theory  $T = \langle \text{Rules}, \text{Conflicts}, \succ \rangle$ ,  $\mathcal{S}$  a set of literal statements. For any  $\varphi_1, \varphi_2 \in \mathcal{S}$  such that  $\text{Conflicts}(\varphi_1, \varphi_2)$ ,*

$$P(K_{\phi_1} = \text{in}) \leq P(K_{\phi_2} \neq \text{in})$$

where

- $P(K_{\phi_2} \neq \text{in}) = P(K_{\phi_2} = \text{ni})$  in the case of bivalent  $\{\text{in}, \text{ni}\}$ -labellings;
- $P(K_{\phi_2} \neq \text{in}) = P(K_{\phi_2} = \text{und}) + P(K_{\phi_2} = \text{niund})$  in the case of trivalent  $\{\text{in}, \text{und}, \text{niund}\}$ -labellings.

These properties of the framework are used later to obtain probabilistic results in our deontic investigation.

The ‘source’ and interpretation of probability values can be diverse. In a classical approach for example, one may simply assume that all the possible outcomes are equally possible and thus that they should share the same probability value (as illustrated in Example 2 above). Alternatively, in a frequentist interpretation, the probability of any event is the frequency of occurrences of the event in a collection of outcomes. Clearly, there are various ways to associate an event with probability values and to interpret these probabilities. In this paper, we assume that probability distributions are simply given *a priori*.

This probabilistic argumentation framework subsumes, or can be related to, a number of existing approaches to probabilistic argumentation [3], and consequently it can be applied to various investigations in the legal domain, such as game-theoretical investigations of optimal strategies in dialogue games [14]. It can also be used in machine learning endeavours, such as structure learning [30, 31] or combinations of probabilistic argumentation with reinforcement learning to determine actions to pursue in environments featuring basic conditional obligations [12, 32]. In the remainder of the paper, we consider how deontic concepts can be represented and reasoned about by exploiting probabilistic labellings.

To recap, we described a propositional language with modal operators, which — in the next section — will be instantiated into deontic modalities. We have also described the syntax of defeasible rules, and how such rules can be combined to form arguments. We introduced a labelling approach to argument semantics, from which the acceptance status of statements can be computed. Finally, we described how uncertainty and probability values can be associated with arguments and statements within our system.

### 3. Probabilistic Deontic Argumentation

Having laid out a simple probabilistic rule-based argumentation system, we can now specify a deontic version of it. To do so, we adopt the deontic argumentation framework of [9] and develop it within our probabilistic setting.

#### 3.1. Deontic language

In law and moral reasoning multiple normative concepts occur. We do not propose to consider all these concepts, but rather focus on three basic deontic notions, namely obligations, prohibitions, and permissions.

We assume that such deontic operators are interdefinable: the prohibition of something is the obligation of its opposite, and the permission of something is the negation of its prohibition. Therefore, we focus on a single deontic obligation operator  $O$ , and accordingly, we assume a language  $\mathcal{L}_D$  whose literal statements are defined as follows, cf. Definition 2.1.

**Definition 3.1.** *A literal statement of a language  $\mathcal{L}_D$  is either a plain literal statement or a deontic literal statement, where:*

- a **plain literal statement** is either an atomic proposition  $p$  or the negation of an atomic proposition, i.e.  $\neg p$ , and
- a **deontic literal statement** is a statement of the form  $O\gamma$  or  $\neg O\gamma$  such that  $\gamma$  is a plain literal statement.

As indicated above, prohibitions and permissions are captured by assuming that a prohibition on something ( $F\gamma$ ) is equivalently expressed by the obligation of its opposite ( $O\bar{\gamma}$ ), and a permission for something ( $P\gamma$ ) is equivalent to the negation of the obligation of the opposite ( $\neg O\bar{\gamma}$ ).

**Notation 3.1.** As syntactic sugar, we may write  $O\bar{\gamma}$  as  $F\gamma$ , and  $\neg O\bar{\gamma}$  as  $P\gamma$ .

A normative system may not be complete, i.e. it may not regulate every possible state of affairs. This means that the system does not specify, for every state of affairs, whether that state is obligatory, prohibited or permitted: there may be gaps, namely states of affairs for which no deontic position is specified. In our model, we only consider those states of affairs that are described by plain literals. Thus, by *normative completeness*, we refer to the quality of a normative system having deontic positions specified for all plain literals. We focus on normative systems that are meant to be complete, and whose gaps are filled according to the ‘principle of prohibition’ (see [33] p.125).

The principle of prohibition can be formulated as follows: ‘everything that is not prohibited is permitted’. Although the principle is rather clear at first sight, it may not be given a unique interpretation. We can adopt a mere tautological reading, based on the above notational convention according to which being permitted for something means not being obligatory of the opposite, i.e. according to which  $P\gamma$  is defined as  $\neg O\bar{\gamma}$ . Following this reading, we can say that a normative system specifies that a state of affairs is permitted when it specifies that there is no prohibition of it: the normative system states that  $\gamma$  is permitted, when it entails that  $\neg O\bar{\gamma}$ . This reading does not add anything to the content of the concerned normative system. In the absence of any norms specifying that  $\neg O\bar{\gamma}$  or that  $O\gamma$  (in which case  $\gamma$  would be permitted, assuming that obligation entails permission) the deontic status of  $\gamma$  remains undetermined.

Here we adopt a different reading of the principle of prohibition, namely the reading according to which everything is permitted unless it is prohibited. More precisely, according to this reading, a normative system based on the principle of permission grants to a state of affairs the status of being permitted whenever it does not grant to the same state of affairs the status of being prohibited. Following this interpretation, the principle of prohibition is no longer tautological. Rather, it becomes a normative principle included in a normative system, which is used for its completion: every gap in the system based on the principle of permission is closed by generating permissions.

We use the terms ‘strong permission’ and ‘weak permission’ to distinguish the permission that is derived from a specific norm (whose consequent denies that a state of affair is prohibited, or equivalently, permits it) from the permission that is obtained according to the general principle of prohibition, or also, as we shall see, from the principle that whatever is obligatory is permitted.

Note that our concept of a weak permission differs from the concept of a weak permission, as characterised by G. H. von Wright [34], or C.E. Alchourrón and E. Bulygin [33]. For these authors a weak permission does not derive from a normative specification (for example from a permissive norm): it is an assertion about a normative system, namely, the assertion that the

normative system does not entail a prohibition. On the contrary, we view a weak permission as an additional normative specification, which is added to a normative system, according to the principle of prohibition, whenever the normative system does not entail a corresponding prohibition. Finally note that the distinction between strong and weak permissions is not directly represented in our deontic language; this distinction depends in the way in which the permission at issue has been derived, from an explicit permissive norm, or rather according to the general principle of prohibition.

A defeasible rule can specify varied relationships amongst (deontic) literal statements of a given language  $\mathcal{L}_D$ . Such rules are called normative defeasible rules.

**Definition 3.2.** *Given a language  $\mathcal{L}_D$ , a **normative defeasible rule** is a defeasible rule over a set of literal statements in  $\mathcal{L}_D$ .*

For the sake of simplicity, norms potentially captured by ‘modalised rules’, e.g. rules of the form  $O(r : \varphi_1, \dots, \varphi_n, \sim \varphi'_1, \dots, \sim \varphi'_m \Rightarrow \varphi)$ , are not accounted for in this paper. Such constructs and their meanings are left for future work.

Normative rules are partitioned into *foreground rules* and *background rules*. Foreground rules provide substantive normative regulations for particular normative domains, while background rules express general deontic principles underlying the normative system being considered.

Foreground rules are domain-dependent, and represent primary norms such as *constitutive rules* or *regulative rules*. The effect of a constitutive rule is to define a term as understood in a given situation or to ‘create’ an institutional entity from a set of brute or institutional facts. A regulative rule, on the other hand, determines the conditions (premises) when a ‘deontic’ effect (obligation, prohibition, permission) is in force. While constitutive and regulative norms have been formally approached in various (and sometimes sophisticated) ways in the literature [35], the distinction is simply addressed in the present system: the consequent of the rule is a plain literal for constitutive rules, and a deontic literal for regulative rules. A regulative rule whose head is a (strong) permission would typically be used to specify an exception to a prohibition (as discussed by A. Ross [36]), but such a rule can also be used to stress a permission and clarify its conditions without being an exception to any existing prohibitions.

Background rules are domain-independent and apply to a whole legal system, or to broad sections of it (e.g. it has often been argued that the principle of prohibition applies to criminal law). They express general deontic assumptions. These background rules can be viewed as defeasible rule schemata which are isomorphic to some pieces of very basic legal doctrines or some axioms of a deontic system. Various background defeasible rules may be proposed, we provide some examples of such rules below.

d. $\gamma$ :  $O\gamma \Rightarrow P\gamma$  An obligation  $O\gamma$  implies a permission  $P\gamma$  (cf. Axiom ‘D’ in deontic logics),

p. $\gamma$ :  $\Rightarrow P\gamma$  Anything is permitted prima facie.

k. $\gamma$ :  $\sim O\bar{\gamma} \Rightarrow P\gamma$  Anything that is not prohibited is permitted.

These background rules can be employed to build arguments supporting permissions. These arguments are successful as long as no corresponding prohibitions can be accepted. We

recognise that such rules do not capture all possible deontic schemata, but note they are relevant in many contexts.

Different sets of background rules result in systems with different behaviour. In this work, we focus on two sets of background rules:  $\{d.\gamma, p.\gamma\}$  and  $\{d.\gamma, k.\gamma\}$ .

**Definition 3.3.** *A set of background defeasible rule schemata  $B$  is*

- *a prima facie permissive set of background defeasible rule schemata iff  $B = \{d.\gamma, p.\gamma\}$ ;*
- *a Kelsenian permissive set of background defeasible rule schemata iff  $B = \{d.\gamma, k.\gamma\}$ ;*
- *a permissive set of background defeasible rule schemata iff  $B = \{d.\gamma, p.\gamma\}$  or  $B = \{d.\gamma, k.\gamma\}$ .*

‘Prima facie permissive’ and ‘Kelsenian permissive’ sets of rules both indicate that anything is defeasibly permitted. However, the Kelsenian permissive set may better reflect the principle of prohibition as exposed by the legal theorist H. Kelsen (thus its name). For both sets, we will see that such background rules are not enough to obtain normative completeness when using bivalent labellings.

Whatever the set of background defeasible rules, we will ground the rules over a set of propositions. For our purposes, we do so over the propositions of an input (domain-dependent) defeasible theory to yield a set of background rules. Then background rules are combined with the rules of the theory to obtain backgrounded rules.

**Definition 3.4.** *A set of rules is a set of background rules with respect to a defeasible theory  $T$  and a set of background defeasible rule schemata  $B$ , denoted  $\text{BackRules}(T, B)$ , iff*

- $\text{BackRules}(T, B) = \{d.\gamma, d.\bar{\gamma}, p.\gamma, p.\bar{\gamma} \mid \gamma \in \text{Prop}(T)\}$  if  $B$  is a prima facie permissive set of background defeasible rule schemata;
- $\text{BackRules}(T, B) = \{d.\gamma, d.\bar{\gamma}, k.\gamma, k.\bar{\gamma} \mid \gamma \in \text{Prop}(T)\}$  if  $B$  is a Kelsenian permissive set of background defeasible rule schemata.

**Definition 3.5.** *A set of rules  $\text{Rules}$  is a backgrounded set of rules with respect to a defeasible theory  $T$  and a set of background defeasible rule schemata  $B$  iff  $\text{Rules} = \text{Rules}(T) \cup \text{BackRules}(T, B)$ .*

**Example 3.** *Returning to our park example, let us assume the presence of a policy stating that vehicles are prohibited from entering the park, unless there is an emergency. Let us consider the following atoms and informal meanings.*

vehi: a vehicle stays at the entrance of the park.  
emer: there is an emergency.  
enter: the vehicle enters into the park.

Accordingly,  $F_{\text{enter}}$  means that the vehicle is forbidden to enter into the park, and  $P_{\text{enter}}$  means that the vehicle is permitted to enter.  $F_{\text{vehi}}$  means that it is forbidden that a vehicle stays at the entrance of the park,  $P_{\text{vehi}}$  means that it is permitted that the vehicle stays at the park entrance, and so on.

The policy may be easily formalised by the foreground defeasible theory  $\langle \{r\}, \emptyset, \emptyset \rangle$  where

$r : \text{vehi}, \sim \text{emer} \Rightarrow F_{\text{enter}}$

The Kelsenian permissive set of background rules with respect to theory  $\langle \{r\}, \emptyset, \emptyset \rangle$  includes all the following rules.



d.vehi :	Ovehi	$\Rightarrow$	Pvehi		d. $\neg$ vehi :	O $\neg$ vehi	$\Rightarrow$	P $\neg$ vehi
k.vehi :	$\sim$ Fvehi	$\Rightarrow$	Pvehi		k. $\neg$ vehi :	$\sim$ F $\neg$ vehi	$\Rightarrow$	P $\neg$ vehi
d.emer :	Oemer	$\Rightarrow$	Pemer		d. $\neg$ emer :	O $\neg$ emer	$\Rightarrow$	P $\neg$ emer
k.emer :	$\sim$ Femer	$\Rightarrow$	Pemer		k. $\neg$ emer :	$\sim$ F $\neg$ emer	$\Rightarrow$	P $\neg$ emer
d.enter :	Oenter	$\Rightarrow$	Penter		d. $\neg$ enter :	O $\neg$ enter	$\Rightarrow$	P $\neg$ enter
k.enter :	$\sim$ Fenter	$\Rightarrow$	Penter		k. $\neg$ enter :	$\sim$ F $\neg$ enter	$\Rightarrow$	P $\neg$ enter

□

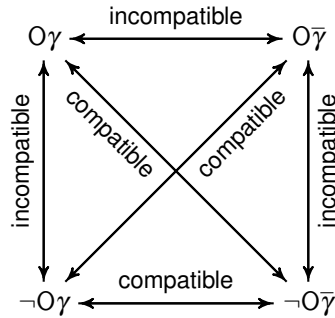


Figure 6: Deontic square of compatibility relation.

Concerning conflicts, we have normative conflicts, and any normative conflict can be a foreground conflict or a background conflict, cf. [9]. A normative conflict has the form  $(\gamma, \bar{\gamma})$  or  $(O\gamma, O\bar{\gamma})$  or  $(\neg O\gamma, O\gamma)$  or  $(O\gamma, \neg O\gamma)$ . The latter three types of conflict are deontic conflicts, and they can be visualised in the deontic square drawn in Figure 6.

**Definition 3.6.** A *normative conflict* is a conflict of the form  $(\gamma, \bar{\gamma})$  or  $(O\gamma, O\bar{\gamma})$  or  $(\neg O\gamma, O\gamma)$  or  $(O\gamma, \neg O\gamma)$ .

For example, for any atomic proposition  $p$ , a conflict  $(Op, \neg Op)$  is captured by the conflict of the form  $(O\gamma, \neg O\gamma)$  where  $\gamma = p$ . Similarly, a conflict  $(O\neg p, \neg O\neg p)$  is captured by  $(O\gamma, \neg O\gamma)$  where  $\gamma = \neg p$ .

Foreground conflicts are normative conflicts. They are meant to be conflicts which are specified in foreground theories (see Definition 3.11). However, such foreground conflicts may be incomplete — some (intuitively desirable) conflicts may not be included in the foreground relation. To ensure completeness of conflicts, we utilise background conflicts. Background conflicts are domain-independent in normative reasoning, and a conflict relation is backgrounded by such conflicts if, and only if, they are included in the conflict relation.

**Definition 3.7.** A set of conflicts is a *set of background conflicts* with respect to a defeasible theory  $T$ , denoted  $\text{BackConflicts}(T)$ , iff  $\text{BackConflicts}(T) = \{(\gamma, \bar{\gamma}), (\bar{\gamma}, \gamma), (O\gamma, O\bar{\gamma}), (O\bar{\gamma}, O\gamma), (O\gamma, \neg O\gamma), (\neg O\gamma, O\gamma), (O\bar{\gamma}, \neg O\bar{\gamma}), (\neg O\bar{\gamma}, O\bar{\gamma}) \mid \gamma \in \text{Prop}(T)\}$ .

**Definition 3.8.** A conflict relation  $\text{Conflicts}$  is a *backgrounded conflict relation* with respect to a defeasible theory  $T$  iff  $\text{Conflicts} = \text{Conflicts}(T) \cup \text{BackConflicts}(T)$ .

**Example 3** (continued). *The background conflict pairs are as follows.*

(vehi, $\neg$ vehi)	(emer, $\neg$ emer)	(enter, $\neg$ enter)
( $\neg$ vehi, vehi)	( $\neg$ emer, emer)	( $\neg$ enter, enter)
(Ovehi, $O\neg$ vehi)	(Oemer, $O\neg$ emer)	(Oenter, $O\neg$ enter)
( $O\neg$ vehi, Ovehi)	( $O\neg$ emer, Oemer)	( $O\neg$ enter, Oenter)
(Ovehi, $\neg$ Ovehi)	(Oemer, $\neg$ Oemer)	(Oenter, $\neg$ Oenter)
( $\neg$ Ovehi, Ovehi)	( $\neg$ Oemer, Oemer)	( $\neg$ Oenter, Oenter)
( $O\neg$ vehi, $\neg$ $O\neg$ vehi)	( $O\neg$ emer, $\neg$ $O\neg$ emer)	( $O\neg$ enter, $\neg$ $O\neg$ enter)
( $\neg$ $O\neg$ vehi, $O\neg$ vehi)	( $\neg$ $O\neg$ emer, $O\neg$ emer)	( $\neg$ $O\neg$ enter, $O\neg$ enter)

□

Given a defeasible theory  $T$  where any conflict in the relation  $\text{Conflicts}(T)$  is a normative conflict, we have that  $\text{Conflicts}(T) \subseteq \text{BackConflicts}(T)$ . Consequently, if  $\text{Conflicts}(T)$  is specified as a set of foreground conflicts which is backgrounded, then every conflict in  $\text{Conflicts}(T)$  is a foreground conflict which is also a background conflict in  $\text{BackConflicts}(T)$ , while in some cases a background conflict may not be a foreground conflict. Thus, if one works with backgrounded conflicts relations of a foreground defeasible theory, as we will do, foreground conflicts are not necessary within the domain specification. Foreground conflicts are nevertheless necessary when specifying conflicts of foreground theories (as per Definition 3.11 below), especially in any setting where such conflicts would not be backgrounded

It is worth recapping our setting on conflicts. We have defined normative conflicts which are conflicts capturing propositional conflicts, as well as deontic conflicts. A normative conflict can be a foreground conflict or a background conflict. Foreground conflicts will be those conflicts which are made explicit when specifying a foreground theory. A backgrounded conflict relation is one for which all possible normative conflicts are included.

Foreground and background deontic rules may have conflicting heads, and to ensure correct reasoning patterns, background superiorities can be proposed. We assume that background rules are inferior to any foreground rules.

**Definition 3.9.** *A superiority relation is a **background superiority relation** with respect to a defeasible theory  $T$  and a set of background defeasible rule schemata  $B$ , denoted  $\text{BackSup}(T, B)$ , iff  $\text{BackSup}(T, B) = \{(s, r) \mid s \in \text{Rules}(T), r \in \text{BackRules}(T, B)\}$ .*

**Definition 3.10.** *A superiority relation  $\succ$  is a **backgrounded superiority relation** with respect to a defeasible theory  $T$  and a set of background defeasible rule schemata  $B$  iff  $\succ = \succ(T) \cup \text{BackSup}(T, B)$ .*

**Example 3** (continued). *The background pairs in the superiority relation are as follows.*

(r, d.vehi)	(r, d.emer)	(r, d.enter)
(r, k.vehi)	(r, k.emer)	(r, k.enter)
(r, d. $\neg$ vehi)	(r, d. $\neg$ emer)	(r, d. $\neg$ enter)
(r, k. $\neg$ vehi)	(r, k. $\neg$ emer)	(r, k. $\neg$ enter)

□

### 3.2. Deontic defeasible theory and argumentation graphs

We now propose to ‘background’ defeasible theories where rules, conflicts and superiority relationships are backgrounded with respect to any foreground defeasible theory. A foreground

defeasible theory is a defeasible theory where rules are not background rules, i.e. rules whose identifiers are not identifiers of any background rules.

**Definition 3.11.** *A defeasible theory  $\langle Rules, Conflicts, \succ \rangle$  is a **foreground defeasible theory** iff*

- *every defeasible rule in Rules is a (foreground) normative defeasible rule which is not a background defeasible rule, and*
- *every conflict in Conflicts is a (foreground) normative conflict.*

**Definition 3.12.** *A defeasible theory  $\langle Rules, Conflicts, \succ \rangle$  is a **backgrounded defeasible theory** of a foreground defeasible theory  $T$  with a set of background defeasible rule schemata  $B$  iff*

- *Rules is a backgrounded set of rules with respect to  $T$  and  $B$ , and*
- *Conflicts is a backgrounded conflict relation with respect to  $T$ , and*
- *$\succ$  is a backgrounded superiority relation with respect to  $T$  and  $B$ .*

In practice, we will first construct a foreground defeasible theory which can be backgrounded to obtain a backgrounded defeasible theory. In the remainder of the paper, we assume that any defeasible theory is backgrounded with a permissive set of background defeasible rule schemata to obtain a permissive defeasible theory.

**Definition 3.13.** *A defeasible theory is a **permissive defeasible theory** iff it is a backgrounded defeasible theory with a permissive set of background defeasible rule schemata.*

We can build arguments from a backgrounded defeasible theory. When building arguments, chaining rules implicitly use the detachment of the consequent of rules. In that regard, deontic studies usually distinguish between factual detachment and deontic detachment. In this work, we restrict ourselves to factual detachment, leaving (defeasible) deontic detachments (which are considered somewhat controversial in the literature [37]) to future work. Once arguments are built we can form an argumentation graph, and then label arguments and (deontic) statements to determine their statuses, as discussed next.

### 3.3. Probabilistic labellings

Given a backgrounded defeasible theory and an argumentation graph built from it, we can now label (deontic) literal statements and associate their acceptance statuses with probability values.

To label arguments and (deontic) statements, we first consider labelling semantics as described in Section 2. Hence, given the argumentation graph built from a backgrounded defeasible theory, arguments are labelled according to the grounded  $\{IN, OUT, UND, OFF\}$ -labelling semantics. Then, literal statements are labelled with acceptance statement labelling semantics. Such labellings are thus a straightforward application of standard labelling semantics corresponding to some common modes of reasoning.

However, bare  $\{IN, OUT, UND, OFF\}$ -labellings are cumbersome in that they allow any ‘doctrinal’ arguments, i.e. arguments based on background principles, to be excluded, i.e. labelled OFF, whereas such doctrinal arguments should always be included as long as the corresponding background principles are endorsed in the legal system being considered. For this reason, we may introduce ‘legitimate’ (legit) grounded  $\{IN, OUT, UND, OFF\}$ -labellings and legit grounded probabilistic labelling frames.

**Definition 3.14.** Let  $G$  be an argumentation graph constructed from a permissive defeasible theory. A grounded  $\{\text{IN}, \text{OUT}, \text{UND}, \text{OFF}\}$ -labelling  $L$  of  $G$  is a **legit grounded**  $\{\text{IN}, \text{OUT}, \text{UND}, \text{OFF}\}$ -labelling of  $G$  iff for any argument  $A \in \mathcal{A}_G$ :

- if  $\text{TopRule}(A) = \text{d.}\gamma$  and  $\text{DirectSub}(A) = \{B\}$  and  $L(B) \neq \text{OFF}$  then  $L(A) \neq \text{OFF}$ , and
- if  $\text{TopRule}(A) = \text{p.}\gamma$  or  $\text{TopRule}(A) = \text{k.}\gamma$  then  $L(A) \neq \text{OFF}$ .

**Example 4.** Let us consider the following arguments, along with the associated argumentation graph  $G$  drawn in Figure 7. All the legit grounded  $\{\text{IN}, \text{OUT}, \text{UND}, \text{OFF}\}$ -labellings of  $G$  are given in Table 1.

O1 :	$\Rightarrow_r$	Oa	O2 :	$\Rightarrow_{r'}$	O¬a				
P1 :	O1	$\Rightarrow_{\text{d.a}}$	Pa	P2 :	O2	$\Rightarrow_{\text{d.¬a}}$	P¬a		
W1 :	$\sim$	Fa	$\Rightarrow_{\text{k.a}}$	Pa	W2 :	$\sim$	F¬a	$\Rightarrow_{\text{k.¬a}}$	P¬a

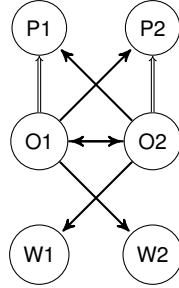


Figure 7: Argumentation graph  $G$ .

O1	O2	P1	P2	W1	W2
OFF	OFF	OFF	OFF	IN	IN
IN	OFF	IN	OFF	IN	OUT
OFF	IN	OFF	IN	OUT	IN
UND	UND	UND	UND	UND	UND

Table 1: Legit grounded  $\{\text{IN}, \text{OUT}, \text{UND}, \text{OFF}\}$ -labellings of argumentation graph  $G$ .

□

**Definition 3.15.** Let  $G$  be an argumentation graph constructed from a permissive defeasible theory. A grounded PLF  $\langle G, \langle \Omega, \mathcal{F}, P \rangle \rangle$  is a **legit grounded probabilistic labelling frame** iff for any  $\{\text{IN}, \text{OUT}, \text{UND}, \text{OFF}\}$ -labelling  $L \in \Omega$  such that  $L$  is not a legit grounded  $\{\text{IN}, \text{OUT}, \text{UND}, \text{OFF}\}$ -labelling,  $P(L) = 0$ .

Moving to statement labellings, (legit) grounded probabilistic labelling frames have multiple properties. Let us mention a few of them. First of all, bivalent  $\{\text{in}, \text{ni}\}$ -labelling and trivalent  $\{\text{in}, \text{und}, \text{niund}\}$ -labelling semantics imply that two conflicting deontic statements cannot be simultaneously labelled in: if a deontic statement is labelled in then any conflicting statement is labelled ni (in the case of  $\{\text{in}, \text{ni}\}$ -labellings) or niund (in the case of  $\{\text{in}, \text{und}, \text{niund}\}$ -labellings), see Proposition 2.1.

**Proposition 3.1.** *Let  $L$  be a grounded  $\{\text{IN}, \text{OUT}, \text{UND}, \text{OFF}\}$ -labelling of an argumentation graph constructed from a permissive defeasible theory  $T$ ,  $\mathcal{S}$  a set of literal statements such that  $\mathcal{S} = \{p, \neg p \mid p \in \text{Prop}(T)\}$ .*

- *Let  $K$  be a bivalent  $\{\text{in}, \text{ni}\}$ -labelling of  $\mathcal{S}$  and from  $\{L\}$ . For any  $\gamma \in \mathcal{S}$ :*
  - *if  $K(L, O\gamma) = \text{in}$  then  $K(L, O\bar{\gamma}) = \text{ni}$ ;*
  - *if  $K(L, \neg O\gamma) = \text{in}$  then  $K(L, O\gamma) = \text{ni}$ ;*
  - *if  $K(L, O\gamma) = \text{in}$  then  $K(L, \neg O\gamma) = \text{ni}$ .*
- *Let  $K$  be a trivalent  $\{\text{in}, \text{und}, \text{niund}\}$ -labelling of  $\mathcal{S}$  and from  $\{L\}$ . For any  $\gamma \in \mathcal{S}$ :*
  - *if  $K(L, O\gamma) = \text{in}$  then  $K(L, O\bar{\gamma}) = \text{niund}$ ;*
  - *if  $K(L, \neg O\gamma) = \text{in}$  then  $K(L, O\gamma) = \text{niund}$ ;*
  - *if  $K(L, O\gamma) = \text{in}$  then  $K(L, \neg O\gamma) = \text{niund}$ .*

By Proposition 2.2 and Corollary 2.1, results in Proposition 3.1 can be understood in our probabilistic development as follows.

**Proposition 3.2.** *Let  $\langle G, \langle \Omega, \mathcal{F}, P \rangle \rangle$  be a grounded PLF where  $G$  is constructed from a permissive defeasible theory  $T$ , and  $\mathcal{S}$  is a set of literal statements such that  $\mathcal{S} = \{p, \neg p \mid p \in \text{Prop}(T)\}$ .*

- *In case of bivalent  $\{\text{in}, \text{ni}\}$ -labellings, for any  $\gamma \in \mathcal{S}$ :*
  - $P(K_{O\bar{\gamma}} = \text{ni} \mid K_{O\gamma} = \text{in}) = 1$ ;
  - $P(K_{O\gamma} = \text{ni} \mid K_{\neg O\gamma} = \text{in}) = 1$ ;
  - $P(K_{\neg O\gamma} = \text{ni} \mid K_{O\gamma} = \text{in}) = 1$ .
- *In case of trivalent  $\{\text{in}, \text{und}, \text{niund}\}$ -labellings, for any  $\gamma \in \mathcal{S}$ :*
  - $P(K_{O\bar{\gamma}} = \text{niund} \mid K_{O\gamma} = \text{in}) = 1$ ;
  - $P(K_{O\gamma} = \text{niund} \mid K_{\neg O\gamma} = \text{in}) = 1$ ;
  - $P(K_{\neg O\gamma} = \text{niund} \mid K_{O\gamma} = \text{in}) = 1$ .

*Proof.* Let us provide the proof for the first item only in the case of bivalent  $\{\text{in}, \text{ni}\}$ -labellings (proofs for other items follow the same structure). From Proposition 3.1, if  $K(L, O\gamma) = \text{in}$  then  $K(L, O\bar{\gamma}) = \text{ni}$ . Therefore, by Proposition 2.2,  $P(K_{O\bar{\gamma}} = \text{ni} \mid K_{O\gamma} = \text{in}) = 1$ .  $\square$

**Corollary 3.1.** *Let  $\langle G, \langle \Omega, \mathcal{F}, P \rangle \rangle$  be a grounded PLF where  $G$  is constructed from a permissive defeasible theory  $T$ , and  $\mathcal{S}$  is a set of literal statements such that  $\mathcal{S} = \{p, \neg p \mid p \in \text{Prop}(T)\}$ .*

- *In case of bivalent  $\{\text{in}, \text{ni}\}$ -labellings, for any  $\gamma \in \mathcal{S}$ :*
  - $P(K_{O\gamma} = \text{in}) \leq P(K_{O\bar{\gamma}} = \text{ni})$ ;
  - $P(K_{\neg O\gamma} = \text{in}) \leq P(K_{O\gamma} = \text{ni})$ ;
  - $P(K_{O\gamma} = \text{in}) \leq P(K_{\neg O\gamma} = \text{ni})$ .
- *In case of trivalent  $\{\text{in}, \text{und}, \text{niund}\}$ -labellings, for any  $\gamma \in \mathcal{S}$ :*
  - $P(K_{O\gamma} = \text{in}) \leq P(K_{O\bar{\gamma}} = \text{niund})$ ;
  - $P(K_{\neg O\gamma} = \text{in}) \leq P(K_{O\gamma} = \text{niund})$ ;
  - $P(K_{O\gamma} = \text{in}) \leq P(K_{\neg O\gamma} = \text{niund})$ .

*Proof.* Let us provide the proof for the first item only in the case of bivalent  $\{\text{in}, \text{ni}\}$ -labellings (proofs for other items follow the same structure). From Proposition 3.1, if  $K(L, O\gamma) = \text{in}$  then  $K(L, O\bar{\gamma}) = \text{ni}$ . Therefore, by Corollary 2.1,  $P(K_{O\gamma} = \text{in}) \leq P(K_{O\bar{\gamma}} = \text{ni})$ .  $\square$

On the same note, Proposition 2.3 can be instantiated with deontic statements which are

inherently conflicting as specified in Definition 3.6. For instance

$$P(K_{F\gamma} = \text{in}) + P(K_{O\gamma} = \text{in}) \leq 1 \quad (1)$$

$$P(K_{F\gamma} = \text{in}) + P(K_{P\gamma} = \text{in}) \leq 1 \quad (2)$$

and thus by Corollary 2.2 we obtain

$$P(K_{O\gamma} = \text{in}) \leq P(K_{F\gamma} = \text{ni}) \quad (3)$$

$$P(K_{O\gamma} = \text{in}) \leq P(K_{F\gamma} = \text{niund}) + P(K_{F\gamma} = \text{und}). \quad (4)$$

$$P(K_{P\gamma} = \text{in}) \leq P(K_{F\gamma} = \text{ni}) \quad (5)$$

$$P(K_{P\gamma} = \text{in}) \leq P(K_{F\gamma} = \text{niund}) + P(K_{F\gamma} = \text{und}). \quad (6)$$

These results are more specific than results given in Corollary 3.1 in the case of trivalent  $\{\text{in}, \text{und}, \text{niund}\}$ -labellings.

The above results are based on the fact that if a deontic statement is labelled in then any conflicting statement is labelled ni or niund. Let us now inspect probabilistic relationships between obligations and permissions. We can first note that if an obligation  $O\gamma$  is labelled in then the implied permission  $\neg O\bar{\gamma}$  (i.e.  $P\gamma$ ) is also labelled in.

**Proposition 3.3.** *Let  $L$  be a legit grounded  $\{\text{IN}, \text{OUT}, \text{UND}, \text{OFF}\}$ -labelling of an argumentation graph constructed from a permissive defeasible theory  $T$ ,  $\mathcal{S}$  a set of literal statements such that  $\mathcal{S} = \{p, \neg p \mid p \in \text{Prop}(T)\}$ , and  $K$  a bivalent  $\{\text{in}, \text{ni}\}$ -labelling or trivalent  $\{\text{in}, \text{und}, \text{niund}\}$ -labelling of  $\mathcal{S}$  and from  $\{L\}$ . For any  $\gamma \in \mathcal{S}$ , if  $K(L, O\gamma) = \text{in}$  then  $K(L, \neg O\bar{\gamma}) = \text{in}$ .*

*Proof.* Given an argumentation graph  $G$  constructed from a permissive defeasible theory, if  $K(L, O\gamma) = \text{in}$ , then there exist an argument  $A \in \mathcal{A}_G$  and  $B \in \mathcal{A}_G$  such that  $\text{con}(A) = O\gamma$  and  $B : A \Rightarrow_{d,\gamma} P\gamma$ . Let  $\mathcal{A}^{\rightsquigarrow} \subseteq \mathcal{A}_G$  be the set of attackers of  $A$ ,  $\mathcal{A}^{\rightsquigarrow\rightsquigarrow} \subseteq \mathcal{A}_G$  the set of arguments attacked by  $A$ , and  $\mathcal{B}^{\rightsquigarrow} \subseteq \mathcal{A}_G$  the set of attackers of  $B$ . We have that  $\mathcal{B}^{\rightsquigarrow} \subseteq \mathcal{A}^{\rightsquigarrow} \cup \mathcal{A}^{\rightsquigarrow\rightsquigarrow}$ . By Definition 2.16, if  $L(A) = \text{IN}$  then for any argument  $C \in \mathcal{A}^{\rightsquigarrow} \cup \mathcal{A}^{\rightsquigarrow\rightsquigarrow}$   $L(C) = \text{OUT}$  or  $L(C) = \text{OFF}$ . Since  $\mathcal{B}^{\rightsquigarrow} \subseteq \mathcal{A}^{\rightsquigarrow} \cup \mathcal{A}^{\rightsquigarrow\rightsquigarrow}$ , for any argument  $C \in \mathcal{B}^{\rightsquigarrow}$   $L(C) = \text{OUT}$  or  $L(C) = \text{OFF}$ , and thus  $L(B) = \text{IN}$ . Therefore, if  $L(A) = \text{IN}$  then  $L(B) = \text{IN}$ , and thus if  $K(L, O\gamma) = \text{in}$  then  $K(L, P\gamma) = \text{in}$ , i.e. if  $K(L, O\gamma) = \text{in}$  then  $K(L, \neg O\bar{\gamma}) = \text{in}$ .  $\square$

In our probabilistic framework, Proposition 3.3 tells us that for any outcome where  $O\gamma$  is labelled in, the implied permission  $P\gamma$  is also labelled in. Consequently, the probability that something is permitted given that it is obligatory is one, and the probability that something is permitted is necessarily greater or equal than the probability that it is obligatory.

**Proposition 3.4.** *Let  $\langle G, \langle \Omega, \mathcal{F}, P \rangle \rangle$  be a legit grounded PLF where  $G$  is constructed from a permissive defeasible theory  $T$ , and  $\mathcal{S}$  is a set of literal statements such that  $\mathcal{S} = \{p, \neg p \mid p \in \text{Prop}(T)\}$ . For any  $\gamma \in \mathcal{S}$ :*

$$P(K_{P\gamma} = \text{in} \mid K_{O\gamma} = \text{in}) = 1.$$

*Proof.* From Proposition 3.3, if  $K(L, O\gamma) = \text{in}$  then  $K(L, \neg O\bar{\gamma}) = \text{in}$ . By Proposition 2.2  $P(K_{P\gamma} = \text{in} \mid K_{O\gamma} = \text{in}) = 1$ .  $\square$

**Corollary 3.2.** *Let  $\langle G, \langle \Omega, \mathcal{F}, P \rangle \rangle$  be a legit grounded PLF where  $G$  is constructed from a permissive defeasible theory  $T$ , and  $\mathcal{S}$  is a set of literal statements such that  $\mathcal{S} = \{p, \neg p \mid p \in \text{Prop}(T)\}$ . For any  $\gamma \in \mathcal{S}$ :*

$$P(K_{O\gamma} = \text{in}) \leq P(K_{P\gamma} = \text{in}).$$

*Proof.* From Proposition 3.3, if  $K(L, O\gamma) = \text{in}$  then  $K(L, \neg O\bar{\gamma}) = \text{in}$ . By Corollary 2.1,  $P(K_{O\gamma} = \text{in}) \leq P(K_{P\gamma} = \text{in})$ .  $\square$

We can also observe that any backgrounded defeasible theory along with a trivalent labelling semantics leads to a third interpretation of the principle of prohibition in terms of labelling: if something is not prohibited (the prohibition is labelled *niund*) then it is permitted (the permission is labelled *in*).

**Proposition 3.5.** *Let  $L$  be a legit grounded  $\{\text{IN}, \text{OUT}, \text{UND}, \text{OFF}\}$ -labelling of an argumentation graph constructed from a permissive defeasible theory  $T$ ,  $\mathcal{S}$  a set of literal statements such that  $\mathcal{S} = \{p, \neg p \mid p \in \text{Prop}(T)\}$ . Let  $K$  be a trivalent  $\{\text{in}, \text{und}, \text{niund}\}$ -labelling of  $\mathcal{S}$  and from  $\{L\}$ . For any  $\gamma \in \mathcal{S}$ : if  $K(L, F\gamma) = \text{niund}$  then  $K(L, P\gamma) = \text{in}$ .*

*Proof.* Given an argumentation graph constructed from a permissive defeasible theory  $T$ , for any  $\gamma \in \mathcal{S} = \{p, \neg p \mid p \in \text{Prop}(T)\}$ , there exists a unique argument  $W : \sim F\gamma \Rightarrow_{k,\gamma} P\gamma$ . All attackers of  $W$  are arguments whose conclusion is  $F\gamma$  (i.e.  $O\bar{\gamma}$ ). If  $K(L, O\bar{\gamma}) = \text{niund}$  then all attackers of  $W$  are *OUT* or *OFF*, and thus  $W$  is labelled *IN*, and  $P\gamma$  is labelled *in*. i.e.  $K(L, \neg O\bar{\gamma}) = \text{in}$ . Therefore, if  $K(L, F\gamma) = \text{niund}$  then  $K(L, P\gamma) = \text{in}$ .  $\square$

**Corollary 3.3.** *Let  $L$  be a legit grounded  $\{\text{IN}, \text{OUT}, \text{UND}, \text{OFF}\}$ -labelling of an argumentation graph constructed from a permissive defeasible theory  $T$ , and  $\mathcal{S}$  a set of literal statements such that  $\mathcal{S} = \{p, \neg p \mid p \in \text{Prop}(T)\}$ . Let  $K$  be a trivalent  $\{\text{in}, \text{und}, \text{niund}\}$ -labelling of  $\mathcal{S}$  and from  $\{L\}$ . For any  $\gamma \in \mathcal{S}$ ,  $K(L, F\gamma) = \text{niund}$  iff  $K(L, P\gamma) = \text{in}$ .*

*Proof.* From Proposition 3.1, if  $K(L, \neg O\gamma) = \text{in}$  then  $K(L, O\gamma) = \text{niund}$ , that is, if  $K(L, P\gamma) = \text{in}$  then  $K(L, F\gamma) = \text{niund}$ . In addition, from Proposition 3.5, if  $K(L, F\gamma) = \text{niund}$  then  $K(L, P\gamma) = \text{in}$ . Hence,  $K(L, F\gamma) = \text{niund}$  iff  $K(L, P\gamma) = \text{in}$ .  $\square$

In our probabilistic framework, Corollary 3.3 implies that for any outcome, something is permitted if, and only if, it is not prohibited. Consequently, the probability that something is permitted and the probability that it is not prohibited are equal.

**Proposition 3.6.** *Let  $\langle G, \langle \Omega, \mathcal{F}, P \rangle \rangle$  be a legit grounded PLF where  $G$  is constructed from a permissive defeasible theory  $T$ ,  $\mathcal{S}$  a set of literal statements such that  $\mathcal{S} = \{p, \neg p \mid p \in \text{Prop}(T)\}$ . For any  $\gamma \in \mathcal{S}$ :*

$$P(K_{F\gamma} = \text{niund}) = P(K_{P\gamma} = \text{in}).$$

Now, on the basis of legit grounded  $\{IN, OUT, UND, OFF\}$ -labellings, jurists may argue that bivalent  $\{in, ni\}$ -labellings are not satisfactory. Given an argumentation graph from any backgrounded defeasible theory (which can contain arguments supporting implicit permissions due to background rules), it may be the case that all arguments are labelled UND. Consequently, using a bivalent  $\{in, ni\}$ -labelling, deontic statements  $O\gamma$ ,  $F\gamma$  and  $P\gamma$  may be labelled ni. This is illustrated in Example 4 below, and suggests that bivalent labellings are insufficient to capture normative completeness. For this reason, we discard bivalent  $\{in, ni\}$ -labelling semantics, and instead propose trivalent  $\{in, und, niund\}$ -labelling semantics to cater for deontic reasoning to address normative completeness.

**Example 4** (continued). *Let us consider the legit grounded  $\{IN, OUT, UND, OFF\}$ -labelling drawn in Figure 8.*

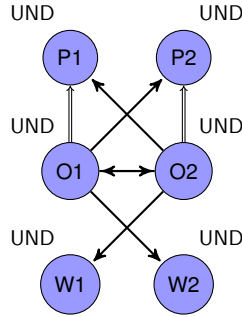


Figure 8: A legit grounded  $\{IN, OUT, UND\}$ -labelling.

First, naive common sense (bivalent) reasoning can be modelled via the following acceptance bivalent  $\{in, ni\}$ -labelling:

$$\langle \emptyset, \{Oa, Pa, O\neg a, \neg Pa\} \rangle.$$

However, this bivalent labelling is problematic from a legal stance because statement  $a$  is here not obligated, permitted or prohibited. To address this gap, we can employ a trivalent  $\{in, und, niund\}$ -labelling:

$$\langle \emptyset, \{Oa, Pa, O\neg a, \neg Pa\}, \emptyset \rangle$$

according to which the deontic status of  $a$  is undecided. □

More formally, the definition of normative gaps — which represent statements that are not permitted, obliged or prohibited — as we may conceive them in terms of statement labellings, depends on whether bivalent  $\{in, ni\}$ -labellings or trivalent  $\{in, und, niund\}$ -labellings are employed.

**Definition 3.16.** *Let  $L$  be a grounded  $\{IN, OUT, UND, OFF\}$ -labelling of an argumentation graph constructed from a defeasible theory  $T$ , and  $\mathcal{S}$  the set of literal statements such that  $\mathcal{S} = \{p, \neg p \mid p \in \text{Prop}(T)\}$ .*

- *Let  $K$  denote a bivalent  $\{in, ni\}$ -labelling of  $\mathcal{S}$  and from  $\{L\}$ . There is a  $\{in, ni\}$ -labelling **normative gap** iff there exists a literal statement  $\gamma \in \mathcal{S}$  such that*



$$K(L, O\gamma) = \text{ni} \text{ and } K(L, O\bar{\gamma}) = \text{ni} \text{ and } K(L, \neg O\bar{\gamma}) = \text{ni}.$$

- Let  $K'$  denote a trivalent  $\{\text{in}, \text{und}, \text{niund}\}$ -labelling of  $\mathcal{S}$  and from  $\{L\}$ . There is a  $\{\text{in}, \text{und}, \text{niund}\}$ -**labelling normative gap** iff there exists a literal statement  $\gamma \in \mathcal{S}$  such that

$$K'(L, O\gamma) = \text{niund} \text{ and } K'(L, O\bar{\gamma}) = \text{niund} \text{ and } K'(L, \neg O\bar{\gamma}) = \text{niund}.$$

In the case of a legit grounded  $\{\text{IN}, \text{OUT}, \text{UND}, \text{OFF}\}$ -labelling of an argumentation graph constructed from a permissive defeasible theory, the definition of gaps may be simplified. We can remark, for example, that the conjunction in the first item can be simplified into ' $K(L, O\bar{\gamma}) = \text{ni}$  and  $K(L, \neg O\bar{\gamma}) = \text{ni}$ ', because  $K(L, \neg O\bar{\gamma}) = \text{ni}$  implies  $K(L, O\gamma) = \text{ni}$  (by Proposition 3.3). In some other deontic frameworks, however, these two conjuncts may hold while  $K(L, O\gamma) = \text{in}$ , suggesting a gap where in fact there is none. Hence, a conjunction specifying the labelling of the three deontic statuses may appear as a stronger definition of gaps, and eventually it may be easier to grasp.

In the case of a legit grounded  $\{\text{IN}, \text{OUT}, \text{UND}, \text{OFF}\}$ -labelling of an argumentation graph constructed from a permissive defeasible theory, and as illustrated in Example 4, bivalent  $\{\text{in}, \text{ni}\}$ -labellings may lead to  $\{\text{in}, \text{ni}\}$ -labelling normative gaps, whereas, trivalent  $\{\text{in}, \text{und}, \text{niund}\}$ -labellings have no normative gaps, i.e. they provide normative completeness.

**Theorem 3.1.** *Let  $L$  be a legit grounded  $\{\text{IN}, \text{OUT}, \text{UND}, \text{OFF}\}$ -labelling of an argumentation graph constructed from a permissive defeasible theory  $T$ ,  $\mathcal{S}$  a set of literal statements such that  $\mathcal{S} = \{p, \neg p \mid p \in \text{Prop}(T)\}$ , and  $K$  a trivalent  $\{\text{in}, \text{und}, \text{niund}\}$ -labelling of  $\mathcal{S}$  and from  $\{L\}$ . For any  $\gamma \in \mathcal{S}$ :*

$$K(L, O\gamma) \neq \text{niund} \text{ or } K(L, O\bar{\gamma}) \neq \text{niund} \text{ or } K(L, \neg O\bar{\gamma}) \neq \text{niund}.$$

*Proof.* Any trivalent  $\{\text{in}, \text{und}, \text{niund}\}$ -labelling  $K$  is a total function. Thus, for any  $\gamma \in \mathcal{S}$ , there are three cases which must be considered: 1.  $K(L, O\bar{\gamma}) = \text{in}$ , 2.  $K(L, O\bar{\gamma}) = \text{und}$ , and 3.  $K(L, O\bar{\gamma}) = \text{niund}$ . If  $K(L, O\bar{\gamma}) = \text{niund}$  then  $K(L, P\gamma) = \text{in}$  (Proposition 3.5), and thus in the last case,  $K(L, P\gamma) = \text{in}$ , i.e.  $K(L, \neg O\bar{\gamma}) = \text{in}$ . Therefore in any case,  $K(L, O\gamma) \neq \text{niund}$  or  $K(L, O\bar{\gamma}) \neq \text{niund}$  or  $K(L, \neg O\bar{\gamma}) \neq \text{niund}$ .  $\square$

Hence, trivalent  $\{\text{in}, \text{und}, \text{niund}\}$ -labellings address normative completeness by means of the status 'undecided' for deontic statements. Eventually, such undecidedness can be disentangled in various ways, typically by a competent authority such as a judge. The trivalent labellings make such undecidedness explicit.

In our probabilistic framework, Theorem 3.1 implies that for any outcome there is no normative gap, and therefore the probability of a normative gap is zero.

**Corollary 3.4.** *Let  $\langle G, \langle \Omega, \mathcal{F}, P \rangle \rangle$  be a legit grounded PLF where  $G$  is constructed from a permissive defeasible theory  $T$ ,  $\mathcal{S}$  a set of literal statements such that  $\mathcal{S} = \{p, \neg p \mid p \in \text{Prop}(T)\}$ . For any  $\gamma \in \mathcal{S}$ :*

$$P(K_{O\gamma} = \text{niund}, K_{F\gamma} = \text{niund}, K_{P\gamma} = \text{niund}) = 0.$$

The results presented above hold for any permissive theory. Consequently, they hold for any backgrounded defeasible theory of a (foreground) defeasible theory with a prima facie or Kelsenian permissive set of background defeasible rule schemata. In general, it turns out that both sets yield the same trivalent  $\{\text{in}, \text{und}, \text{niund}\}$ -labelling. To show this result, let us first introduce some definitions for ease of presentation.

**Definition 3.17.** Two arguments  $A$  and  $B$  are ***k-equivalent arguments*** iff

- if  $A$  is of the form  $A : \Rightarrow_{p,\gamma} P\gamma$  then  $B$  is of the form  $B : \sim F\gamma \Rightarrow_{k,\gamma} P\gamma$ , and
- if  $A$  is of the form  $A : A_1, \dots, A_n, \sim \varphi_1, \dots, \sim \varphi_m \Rightarrow_r \varphi$  ( $r \neq p,\gamma$ ) then  $B$  is of the form  $B : B_1, \dots, B_n, \sim \varphi_1, \dots, \sim \varphi_m \Rightarrow_r \varphi$ , where  $A_i$  and  $B_i$  ( $1 \leq i \leq n$ ) are *k-equivalent*.

**Definition 3.18.** Argumentation graphs  $G_1$  and  $G_2$  are ***k-equivalent argumentation graphs*** iff there exists a bijection  $f : \mathcal{A}_{G_1} \rightarrow \mathcal{A}_{G_2}$  such that for any argument  $A \in \mathcal{A}_{G_1}$ :

- $f(A) = B$  iff  $A$  and  $B$  are *k-equivalent*, and
- $(A, A') \in \rightsquigarrow_{G_1}$  iff  $(f(A), f(A')) \in \rightsquigarrow_{G_2}$ .

**Definition 3.19.** Let

- $G_1$  and  $G_2$  be two argumentation graphs;
- $H_1$  ( $H_2$  resp.) be a subargument-complete subgraph of  $G_1$  ( $G_2$  resp.);
- $L_1$  ( $L_2$  resp.) be a legit grounded  $\{\text{IN}, \text{OUT}, \text{UND}, \text{OFF}\}$ -labelling of the argumentation graph  $G_1$  ( $G_2$  resp.) with respect to subgraph  $H_1$  ( $H_2$  resp.);

$L_1$  and  $L_2$  are ***k-equivalent labellings*** iff  $G_1$  and  $G_2$  are *k-equivalent*, and  $H_1$  and  $H_2$  are *k-equivalent*.

The work in [9] showed that both types of permissive theories yield the same trivalent  $\{\text{in}, \text{und}, \text{niund}\}$ -labelling for  $\{\text{IN}, \text{OUT}, \text{UND}\}$ -labellings. Theorem 3.2 is a generalisation of this result to grounded  $\{\text{IN}, \text{OUT}, \text{UND}, \text{OFF}\}$ -labellings.

**Theorem 3.2.** Let

- $T$  be a (foreground) defeasible theory;
- $U$  be the backgrounded defeasible theory of  $T$  with a prima facie permissive set of background defeasible rule schemata;
- $V$  be the backgrounded defeasible theory of  $T$  with a Kelsenian permissive set of background defeasible rule schemata;
- $L_U$  ( $L_V$  resp.) be a legit grounded  $\{\text{IN}, \text{OUT}, \text{UND}, \text{OFF}\}$ -labelling of the argumentation graph constructed from  $U$  ( $V$  resp.);
- $\mathcal{S}$  be a set of literal statements such that  $\mathcal{S} = \{p, \neg p \mid p \in \text{Prop}(T)\}$ , and  $K$  a trivalent  $\{\text{in}, \text{und}, \text{niund}\}$ -labelling of  $\mathcal{S}$  and from  $\{L\}$ .

For any  $\gamma \in \mathcal{S}$ , if  $L_U$  and  $L_V$  are *k-equivalent* then

$$K(L_U, \gamma) = K(L_V, \gamma) \text{ and } K(L_U, O\gamma) = K(L_V, O\gamma) \text{ and } K(L_U, \neg O\gamma) = K(L_V, \neg O\gamma).$$

*Proof.* Let  $G_U$  ( $G_V$  resp.) denote the argumentation graph constructed from  $U$  ( $V$  resp.), and  $H_U$  ( $H_V$  resp.) a subgraph of  $G_U$  ( $G_V$  resp.) such that  $L_U$  ( $L_V$  resp.) the legit grounded  $\{\text{IN}, \text{OUT}, \text{UND}, \text{OFF}\}$ -labelling with respect to  $H_U$  ( $H_V$  resp.). If  $L_U$  and  $L_V$  are *k-equivalent* then  $G_U$  and  $G_V$  are *k-equivalent*, and  $H_U$  and  $H_V$  are *k-equivalent*. By Theorem 3.2 in [9], for any  $\gamma \in \mathcal{S}$ :  $K(L_U, \gamma) = K(L_V, \gamma)$  and  $K(L_U, O\gamma) = K(L_V, O\gamma)$  and  $K(L_U, \neg O\gamma) = K(L_V, \neg O\gamma)$ .  $\square$

Theorem 3.2 extends representation results to our probabilistic setting. The principle of prohibition ‘anything that is not prohibited is permitted’ as a schema  $k,\gamma : \sim O\bar{\gamma} \Rightarrow P\gamma$  or a principle such as ‘anything is permitted prima facie’ as a schema  $p,\gamma : \Rightarrow P\gamma$  are two alternatives to cater for normative completeness, and both alternatives lead to the same results in terms of statement labellings in the case of grounded  $\{\text{IN}, \text{OUT}, \text{UND}\}$ -labellings [9]. Theorem 3.2 extends the result to legit grounded  $\{\text{IN}, \text{OUT}, \text{UND}, \text{OFF}\}$ -labellings and thus to our probabilistic setting.

## 4. Worked Example

To illustrate our system, let us reappraise the policy stating that vehicles are forbidden to enter a park unless there is an emergency. This policy and its associated assumptions may be formalised in various ways. We illustrate our system with one option which is developed in the remainder of this section.

### 4.1. Backgrounded defeasible theory

We assume that there is a vehicle at the entrance of the park, and that there may *possibly* be an emergency. Let us capture this with the foreground defeasible theory  $\langle \{rv, re, r\bar{e}, r\}, \emptyset, \emptyset \rangle$  where

$$\begin{aligned} rv &\Rightarrow \text{vehi} \\ re &\Rightarrow \text{emer} \\ r\bar{e} &\Rightarrow \neg\text{emer} \\ r &: \text{vehi}, \sim \text{emer} \Rightarrow \text{Fenter} \end{aligned}$$

The foreground theory can be then backgrounded to yield a backgrounded theory featuring, amongst others, background rules as exposed in Example 3.

### 4.2. Argument and argumentation graph construction

Let us consider a Kelsenian permissive set of background defeasible rule schemata. Accordingly, we can build the following arguments from background rules:

$$\begin{aligned} W1: \sim F\text{vehi} &\Rightarrow_{k,\text{vehi}} P\text{vehi} & W4: \sim F\neg\text{vehi} &\Rightarrow_{k,\neg\text{vehi}} P\neg\text{vehi} \\ W2: \sim F\text{emer} &\Rightarrow_{k,\text{emer}} P\text{emer} & W5: \sim F\neg\text{emer} &\Rightarrow_{k,\neg\text{emer}} P\neg\text{emer} \\ W3: \sim F\text{enter} &\Rightarrow_{k,\text{enter}} P\text{enter} & W6: \sim F\neg\text{enter} &\Rightarrow_{k,\neg\text{enter}} P\neg\text{enter} \end{aligned}$$

In addition, we can build the following arguments:

$$\begin{aligned} A1: &\Rightarrow_{rv} \text{vehi} & B1: &\Rightarrow_{re} \text{emer} \\ A2: A1, \sim \text{emer} &\Rightarrow_r \text{Fenter} & C1: &\Rightarrow_{r\bar{e}} \neg\text{emer} \\ A3: A2 &\Rightarrow_{d,\neg\text{enter}} P\neg\text{enter} \end{aligned}$$

Consequently, we can build the argumentation graph  $G$  such that:  $\mathcal{A}_G = \{A1, A2, A3, B1, C1, W1, W2, W3, W4, W5, W6\}$ , and  $\rightsquigarrow_G = \{(B1, C1), (C1, B1), (B1, A2), (B1, A3), (A2, W3)\}$ , and  $\models_G = \{(A1, A2), (A2, A3)\}$  (see Figure 9).

We note that we have built arguments to support weak/doctrinal permissions, thus we can argue and present full-fledged arguments about such permissions (here arguments  $W1 \dots W6$ ).

### 4.3. Argument labelling

The grounded  $\{\text{IN}, \text{OUT}, \text{UND}\}$ -labelling of argumentation graph  $G$  is illustrated in Figure 9, and all legit grounded  $\{\text{IN}, \text{OUT}, \text{UND}, \text{OFF}\}$ -labelling outcomes which are deemed possible are exposed in Table 2.

We assume that the cyber-physical system is such that arguments  $B1$  and  $C1$  cannot both be labelled  $\text{OFF}$ , so that it is always indicated whether it is the case that there is an emergency, or whether the case is undecided. We also assume that rule  $r$  is always applied when its antecedents hold. Finally we adopt the principle of indifference on

all possible legit grounded  $\{IN, OUT, UND, OFF\}$ -labelling outcomes, so that they are equally distributed as in Table 2.

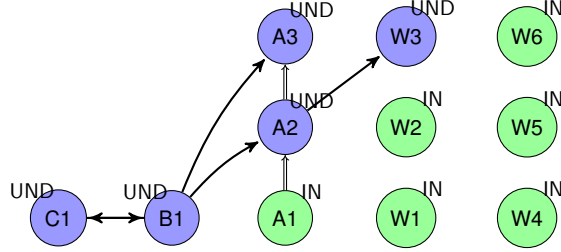


Figure 9: Grounded  $\{IN, OUT, UND\}$ -labelling of argumentation graph G.

A1	A2	A3	B1	C1	W1	W2	W3	W4	W5	W6	$P(\cdot)$
OFF	OFF	OFF	IN	OFF	IN	IN	IN	IN	IN	IN	1/6
OFF	OFF	OFF	OFF	IN	IN	IN	IN	IN	IN	IN	1/6
OFF	OFF	OFF	UND	UND	IN	IN	IN	IN	IN	IN	1/6
IN	OUT	OUT	IN	OFF	IN	IN	IN	IN	IN	IN	1/6
IN	IN	IN	OFF	IN	IN	IN	OUT	IN	IN	IN	1/6
IN	UND	UND	UND	UND	IN	IN	UND	IN	IN	IN	1/6

Table 2: Legit grounded  $\{IN, OUT, UND, OFF\}$ -labellings of argumentation graph G with non-zero probabilities.

Given probabilities in Table 2, we can easily see that for example

$$P(L_{A3} = IN) = 1/6 \quad (7)$$

$$P(L_{A3} = IN \mid L_{B1} = IN) = 0. \quad (8)$$

In words, the marginal probability that argument A3 is labelled IN is 1/6; and the probability that argument A3 is labelled IN, given that B1 is labelled IN, is zero.

#### 4.4. Statement labelling

Let L denote the grounded  $\{IN, OUT, UND\}$ -labelling of argumentation graph G, i.e. the labelling in the last row in Table 2. From this, we can obtain the bivalent  $\{in, ni\}$ -labelling and trivalent  $\{in, und, niund\}$ -labelling as exposed in Table 3.

	vehi	¬vehi	emer	¬emer	enter	¬enter
K(L,·)	in	ni	ni	ni	ni	ni
K(L,·)	in	niund	und	und	niund	niund
	Ovehi	O¬vehi	Oemer	O¬emer	Oenter	O¬enter
K(L,·)	ni	ni	ni	ni	ni	ni
K(L,·)	niund	niund	niund	niund	niund	und
	Pvehi	P¬vehi	Pemer	P¬emer	Penter	P¬enter
K(L,·)	in	in	in	in	ni	in
K(L,·)	in	in	in	in	und	in

Table 3: Bivalent  $\{in, ni\}$ -labelling and trivalent  $\{in, und, niund\}$ -labelling.

We can remark that the  $\{in, ni\}$ -bivalent labelling results in a normative gap (enter is neither obligatory nor prohibited nor permitted), whereas the trivalent  $\{in, und, niund\}$ -labelling fills this gap by labelling the permission to enter as undecided.

More generally, for every legit grounded  $\{IN, OUT, UND, OFF\}$ -labelling, we can draw a statement labelling, and from all the statement labellings we can determine probability values on the acceptance statuses of statements. For instance, from the argument labelling probabilities as given in Table 2, we can compute that:

$$P(K_{vehi} = in) = 1/2 \quad (9)$$

$$P(K_{Fenter} = in) = 1/6 \quad (10)$$

$$P(K_{Fenter} = in \mid K_{vehi} = in) = 1/3. \quad (11)$$

In words, the marginal probability that there is a vehicle is  $1/2$ , and the marginal probability that there is a detached prohibition from entering the park is  $1/6$ . If it is accepted that there is a vehicle (a vehicle is perceived), then the probability that the vehicle should not enter is  $1/3$  (there may be an emergency, or the case is undecided). Such likelihoods may be useful, for example, in determining whether to monitor specific portions of the system.

#### 4.5. Violation and contrary-to-duty obligation

Let us extend the example with the formalisation of a violation and a contrary-to-duty obligation. Such obligations can be a pitfall for deontic formalisms which have a more sophisticated conception of deontic modalities [38, 39], and we would like to illustrate how such obligations can be handled in our probabilistic deontic argumentation framework.

Let us suppose that the park policy also states that a violation of the prohibition would be sanctioned by a fine (the amount does not matter for our purposes). To capture such a policy, we can add the following rules.

$$\begin{array}{ll} v : & \text{Fenter, enter} \Rightarrow \text{violation} & \bar{v} : & \Rightarrow \neg \text{violation} \\ f : & \text{violation} \Rightarrow \text{fine} & \bar{f} : & \Rightarrow \neg \text{fine} \end{array}$$

such that  $v \succ \bar{v}$  and  $f \succ \bar{f}$ . Rules  $\bar{v}$  and  $\bar{f}$  specify that, by default, we can derive that there is neither violation nor fine, unless the contrary is shown.

Furthermore, assume that the park's management adds the following contrary-to-duty obligation: if the prohibition is violated then one should stop driving (in the park). We can thus add the following rule.

s: violation  $\Rightarrow$  Ostop

A sequence of compensatory obligations can be added along similar lines. We also assume that a vehicle enters in the park:

e:  $\Rightarrow$  enter

In addition to the previous arguments, we can thus build the following arguments (amongst others).

E1:  $\Rightarrow_e$  enter  
V1: E1, A2  $\Rightarrow_r$  violation  
F1: V1  $\Rightarrow_f$  fine  
V2:  $\Rightarrow_{\bar{v}}$   $\neg$ violation  
F2:  $\Rightarrow_{\bar{f}}$   $\neg$ fine

A1	A2	A3	B1	C1	W1	W2	W3	W4	W5	W6	E1	V1	$P(\cdot)$
OFF	OFF	OFF	IN	OFF	IN	IN	IN	IN	IN	IN	OFF	OFF	1/9
OFF	OFF	OFF	OFF	IN	IN	IN	IN	IN	IN	IN	OFF	OFF	1/9
OFF	OFF	OFF	UND	UND	IN	IN	IN	IN	IN	IN	OFF	OFF	1/9
IN	OUT	OUT	IN	OFF	IN	IN	IN	IN	IN	IN	OFF	OFF	1/9
IN	IN	IN	OFF	IN	IN	IN	OUT	IN	IN	IN	OFF	OFF	1/9
IN	UND	UND	UND	UND	IN	IN	UND	IN	IN	IN	OFF	OFF	1/9
IN	OUT	OUT	IN	OFF	IN	IN	IN	IN	IN	IN	IN	OUT	1/9
IN	IN	IN	OFF	IN	IN	IN	OUT	IN	IN	IN	IN	IN	1/9
IN	UND	UND	UND	UND	IN	IN	UND	IN	IN	IN	IN	UND	1/9

Table 4: Legit grounded {IN,OUT,UND,OFF}-labellings with non-zero probabilities (not all possible arguments are labelled due to the lack of space).

Again, for every legit grounded {IN,OUT,UND,OFF}-labelling of the resulting argumentation graph, we can draw a statement labelling, and from all the statement labellings we can determine probability values on the labellings of statements. For instance, from the argument labelling probabilities given in Table 4, and assuming that rules v, f,  $\bar{v}$ ,  $\bar{f}$  and s apply in all outcomes where their antecedents hold, and assuming that there is no doubt that a vehicle has entered in the park, we can compute the following probabilities.

$$\begin{aligned}
P(K_{\text{violation}} = \text{in} \mid K_{\text{vehi}} = \text{in}) &= 1/6 & P(K_{\text{violation}} = \text{in} \mid K_{\text{vehi}} = \text{in}, K_{\text{enter}} = \text{in}) &= 1/3 \\
P(K_{\text{fine}} = \text{in} \mid K_{\text{vehi}} = \text{in}) &= 1/6 & P(K_{\text{fine}} = \text{in} \mid K_{\text{vehi}} = \text{in}, K_{\text{enter}} = \text{in}) &= 1/3 \\
P(K_{\text{Ostop}} = \text{in} \mid K_{\text{vehi}} = \text{in}) &= 1/6 & P(K_{\text{Ostop}} = \text{in} \mid K_{\text{vehi}} = \text{in}, K_{\text{enter}} = \text{in}) &= 1/3
\end{aligned}$$

In words, if a vehicle is perceived, then a violation, its sanction, and the associated contrary-to-duty obligation will occur with a probability value 1/6. If, in addition the vehicle has entered the park, then the probability values increase to 1/3.

## 5. Conclusion

We have introduced a probabilistic deontic argumentation framework to reason about probability values attached to deontic statements and arguments supporting these statements. Given a set of norms and a state of affairs modelled as a defeasible theory, the framework allows one to associate acceptance statuses of deontic statements and related violations with probability values, in a principled way.

To do so, we have combined a deontic argumentation approach where deontic principles are reified [9] and a probabilistic argumentation approach where argument labellings are associated with probabilities values [3]. The deontic argumentation composition relies on common concepts taken from computational models of argument: rule-based arguments, argumentation graphs, argument labelling semantics and statement labelling semantics. In this framework, normative completeness is addressed by reifying the principle of prohibition. This deontic composition is then combined with probabilistic labellings used in probabilistic argumentation, enabling us to associate statement acceptance statuses with probability values.

We have learnt that the construction of a probabilistic deontic argumentation framework can be achieved by using standard constructs from computational models of argument and probabilistic labellings. Then, we could have appreciated that argumentation properties of the deontic framework can be coupled with probabilistic counterparts, i.e. relationships amongst acceptance statuses of deontic statements can be understood in meaningful relationships of probabilities of acceptance statuses of these statements.

Throughout the paper we have identified several potential directions for future research. In addition, we note that no algorithms to learn or compute probability statuses of arguments and statements have been considered, and exploring efficient algorithms for this purpose would be interesting. Finally, throughout the paper, we have avoided relating the probabilistic framework with possible probability interpretations (such as classical, frequentist, or Bayesian views on these matters), and thus it may be interesting to investigate the impact of such interpretations on probabilistic deontic issues.

## References

- [1] B. Verheij, F. Bex, S. T. Timmer, C. S. Vlek, J.-J. C. Meyer, S. Renooij, H. Prakken, Arguments, scenarios and probabilities: connections between three normative frameworks for evidential reasoning, *Law, Probability and Risk* 15 (1) (2015) 35–70.
- [2] A. Hunter, M. Thimm, Probabilistic reasoning with abstract argumentation frameworks, *Journal of Artificial Intelligence Research* 59 (2017) 565–611.
- [3] R. Riveret, P. Baroni, Y. Gao, G. Governatori, A. Rotolo, G. Sartor, A labelling framework for probabilistic argumentation, *Annals of Mathematics and Artificial Intelligence* 83 (1) (2018) 21–71.
- [4] H. Hart, *The Concept of Law*, Oxford University Press, 1961.
- [5] J. Horty, Nonmonotonic foundations for deontic logic, in: *Defeasible Deontic Logic*, Kluwer Academic Publishers, 1997, pp. 17–44.
- [6] B. Liao, N. Oren, L. Van Der Torre, S. Villata, Prioritized norms in formal argumentation, *Journal of Logic and Computation* 29 (2) (2018) 215–240.
- [7] M. Beirlaen, J. Heyninck, C. Straßer, Structured argumentation with prioritized conditional obligations and permissions, *Journal of Logic and Computation* 29 (2) (2019) 187–214.
- [8] G. Pigozzi, L. van der Torre, Arguing about constitutive and regulative norms, *Journal of Applied Non-Classical Logics* 28 (2-3) (2018) 189–217.

- [9] R. Riveret, A. Rotolo, G. Sartor, A deontic argumentation framework towards doctrine reification, *IfCoLog Journal of Logics and their Applications* 6 (5) (2019) 903–939.
- [10] L. L. Fuller, Positivism and fidelity to law: A reply to Professor Hart, *Harvard Law Review* 71 (4) (1958) 630–672.
- [11] H. L. A. Hart, Positivism and the separation of law and morals, *Harvard Law Review* 71 (4) (1958) 593–629.
- [12] R. Riveret, Y. Gao, G. Governatori, A. Rotolo, J. Pitt, G. Sartor, A probabilistic argumentation framework for reinforcement learning agents - towards a mentalistic approach to agent profiles, *Autonomous Agents and Multi-Agent Systems* 33 (1-2) (2019) 216–274.
- [13] R. Riveret, A. Rotolo, G. Sartor, H. Prakken, B. Roth, Success chances in argument games: A probabilistic approach to legal disputes, in: *Proc. of the 20th Conference on Legal Knowledge and Information Systems*, IOS Press, 2007, pp. 99–108.
- [14] R. Riveret, H. Prakken, A. Rotolo, G. Sartor, Heuristics in argumentation: A game theory investigation, in: *Proc. of the 2nd Conference on Computational Models of Argument*, IOS Press, 2008, pp. 324–335.
- [15] S. Modgil, H. Prakken, The ASPIC<sup>+</sup> framework for structured argumentation: a tutorial, *Argument & Computation* 5 (1) (2014) 31–62.
- [16] M. Caminada, L. Amgoud, On the evaluation of argumentation formalisms, *Artificial Intelligence* 171 (5-6) (2007) 286–310.
- [17] Z. Li, S. Parsons, On argumentation with purely defeasible rules, in: *Proc. of the 9th International Conference Scalable Uncertainty Management*, Vol. 9310, Springer, 2015, pp. 330–343.
- [18] A. Bondarenko, P. M. Dung, R. A. Kowalski, F. Toni, An abstract, argumentation-theoretic approach to default reasoning, *Artificial Intelligence* 93 (1-2) (1997) 63–101.
- [19] G. A. Vreeswijk, Abstract argumentation systems, *Artificial Intelligence* 90 (1) (1997) 225 – 279.
- [20] H. Prakken, On support relations in abstract argumentation as abstractions of inferential relations, in: *Proc. of the 21st European Conference on Artificial Intelligence*, IOS Press, 2014, pp. 735–740.
- [21] C. Cayrol, M.-C. Lagasque-Schiex, Bipolarity in argumentation graphs: Towards a better understanding, *International Journal of Approximate Reasoning* 54 (7) (2013) 876 – 899.
- [22] A. Cohen, S. Gottifredi, A. J. Garca, G. R. Simari, A survey of different approaches to support in argumentation systems, *The Knowledge Engineering Review* 29 (5) (2014) 513550.
- [23] S. Polberg, N. Oren, Revisiting support in abstract argumentation systems, in: *Proc. of the 5th International Conference on Computational Models of Argument*, 2014, pp. 369–376.



- [24] A. Cohen, S. Parsons, E. I. Sklar, P. McBurney, A characterization of types of support between structured arguments and their relationship with support in abstract argumentation, *International Journal of Approximate Reasoning* 94 (2018) 76 – 104.
- [25] P. M. Dung, On the acceptability of arguments and its fundamental role in nonmonotonic reasoning, logic programming and n-person games, *Artificial Intelligence* 77 (2) (1995) 321–358.
- [26] P. Baroni, M. Caminada, M. Giacomin, An introduction to argumentation semantics, *Knowledge Engineering Review* 26 (4) (2011) 365–410.
- [27] S. Modgil, M. Caminada, Proof theories and algorithms for abstract argumentation frameworks, in: *Argumentation in Artificial Intelligence*, Springer, 2009, pp. 105–129.
- [28] P. Baroni, G. Governatori, R. Riveret, On labelling statements in multi-labelling argumentation, in: *Proc. of the 22nd European Conference on Artificial Intelligence*, IOS Press, 2016, pp. 489–497.
- [29] P. Baroni, R. Riveret, Enhancing statement evaluation in argumentation via multi-labelling systems, *Journal of Artificial Intelligence Research* 66 (2019) 793–860.
- [30] R. Riveret, G. Governatori, On learning attacks in probabilistic abstract argumentation, in: *Proc. of the 15th International Conference on Autonomous Agents & Multiagent Systems*, ACM, 2016, pp. 653–661.
- [31] R. Riveret, On searching explanatory argumentation graphs, *Journal of Applied Non-Classical Logics* 30 (2) (2020) 123–192.
- [32] R. Riveret, A. Artikis, J. Pitt, E. G. Nepomuceno, Self-governance by transfiguration: From learning to prescription changes, in: *Proc. of the 8th International Conference on Self-Adaptive and Self-Organizing Systems*, 2014, pp. 70–79.
- [33] C. E. Alchourrón, E. Bulygin, *Normative systems*, Springer-Verlag, 1971.
- [34] G. H. von Wright, *Norm and Action: A Logical Enquiry*, Routledge and Kegan Paul, 1963.
- [35] D. Gabbay, J. Horty, X. Parent, R. van der Meyden, L. van der Torre (Eds.), *Handbook of Deontic Logic and Normative Systems*, College Publications, 2013.
- [36] A. Ross, *Directives and Norms*, Humanities Press, 1967.
- [37] R. Mullins, Legal positivism and deontic detachment, *Ratio Juris* 31 (1) (2018) 4–8.
- [38] J. Carmo, A. J. I. Jones, *Deontic Logic and Contrary-to-Duties*, Springer, 2002, pp. 265–343.
- [39] P. McNamara, Deontic logic, in: D. Gabbay, J. Woods (Eds.), *The Handbook of the History of Logic, vol. 7: Logic and the Modalities in the Twentieth Century*, Elsevier Press, 2006, pp. 197–288.