

# Reasoning with Epistemic Rights and Duties: Automating a Dynamic Logic of the Right to Know in LogiKEY

Lara Lawniczak<sup>a</sup>, Luca Pasetto<sup>b,\*</sup>, Christoph Benzmüller<sup>a,c</sup>, Xu Li<sup>b</sup> and Réka Markovich<sup>b</sup>

<sup>a</sup>Otto-Friedrich-Universität Bamberg

<sup>b</sup>University of Luxembourg

<sup>c</sup>Freie Universität Berlin

**Abstract.** It is not straightforward to reason about specific legal concepts such as *epistemic rights and duties*, which are crucial in AI systems that have to make autonomous decisions based on who knows what, who is entitled to know, and under what conditions information should be shared or withheld. Such issues are central to responsible AI, data governance, and regulatory compliance. A concrete application arises in the context of the GDPR, where a data subject has a *right to know* whether and for what purpose her personal data is being processed, creating a *duty to tell* for the controller when asked. On the other hand, if the software used for the processing is proprietary, the data subject does not have the right to know its exact mechanisms, so her asking to know them does not create a corresponding duty for the data controller.

In this paper, a shallow semantical embedding (SSE) of the Dynamic Logic of the Right to Know (LRK) in Higher-Order Logic is presented. The embedding is proven faithful, and it is encoded and experimented with in the *Isabelle/HOL* proof assistant. The SSE is then used to reason with the GDPR example encoded in LRK.

The embedding of LRK differs from existing ones in how it represents the dynamic updating of the model: instead of performing changes on the domain of possible worlds, the provided SSE maintains the accessibility and neighborhood relations within the context of a formula. Updates are then handled by updating the relations, while the domain of possible worlds stays the same.

The work presented in this paper contributes to the LogiKEY knowledge engineering methodology and framework, which enables experimentation with logics and logic combinations, with general and domain knowledge, and with concrete use cases.

## 1 Introduction

*Epistemic rights*, such as the *right to know*, concern an individual's ability to access, receive, and control information, playing a crucial role in our information-centric society [62]. Such rights raise critical questions about transparency, accountability, and compliance with legal frameworks and are increasingly relevant to AI systems that handle personal data, make autonomous decisions, or mediate access to information [59]. For example, the General Data Protection Regulation (GDPR) [26] establishes that the data subject has a right to (or right of) access to some specific information, such as information about whether her personal data is being processed, and if so, for what purpose. As Article 15 puts it, the data subject has a *right to*

*obtain* this information by asking for it. However, we also know from Recital 63 [26] that this right should not adversely affect the rights or freedoms of others, including trade secrets or intellectual property and in particular the copyright protecting the software, meaning that the data subject is not allowed to know this information. Reasoning with such an epistemic right is not straightforward, as it requires formally representing dynamic knowledge changes, contextual dependencies, and the interplay between obligations, permissions, and prohibitions.

In this example, the right to know is a *power* as described in the theory of normative positions based on the work of Hohfeld [34], where power is characterized as the potential of the agent to execute an action resulting in a change in the counterparty's normative positions, for instance by creating a duty [46]. The Dynamic Logic of the Right to Know (LRK) [42] is a logic that explicitly deals with the right to know as a power to know whether something is the case. LRK allows to represent and reason with scenarios where agents are entitled to access certain information under specific conditions, an issue highly relevant to the GDPR and to AI regulation. These regulations concern hundreds of thousand of entities for the procedures of whom automation and transparent, explainable computational solutions are highly valuable, and hence subject to very active research within the field of AI and Law.

The objective is to enable *machines* to reason correctly and verifiably on right to know issues, and to develop practical means to support this. However, neither LRK nor any other logic with a similar purpose has been automated or applied yet. In this paper, we provide a Shallow Semantic Embedding (SSE) of LRK and demonstrate that, and how, this logic can be seen and elegantly handled as a fragment of Classical Higher-Order Logic (HOL) [21, 8]. We show that our embedding is faithful, that is, sound and complete, and we use it to successfully encode and reason with the above GDPR example.

This work constitutes an important addition to the LogiKEY [15] logic-pluralistic knowledge representation and reasoning methodology represented in Fig. 1. LogiKEY's unifying formal framework is based on SSEs of 'object' logics (and their combinations) in HOL, enabling the provision of powerful tool support: off-the-shelf theorem provers and model finders for HOL [9], as provided, e.g., in the proof assistant system *Isabelle/HOL* [50], are assisting the LogiKEY knowledge engineer to flexibly experiment with underlying logics and their combinations, with general and domain knowledge, and with concrete use cases—all at the same time. The approach is capable of supporting model checking, model finding and theorem proving,

\* Corresponding Author. Email: Luca.Pasetto@uni.lu.

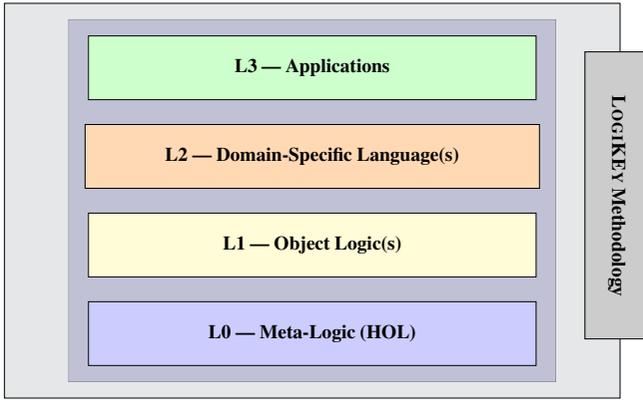


Figure 1. LOGiKEY KR&R methodology

while continuous improvements of off-the-shelf provers and model finders boost the reasoning performance in LogiKEY without further ado. One particular focus of the framework is on ethico-legal applications [17] and normative reasoning [52]. At layer L0 of LogiKEY, HOL is serving as the underlying metalogic to encode (combinations of) the object logics of layer L1, allowing for the expression of domain theories at layer L2, which enables experimentation with concrete applications and examples at layer L3. While the LogiKEY methodology has been applied successfully to automate a number of logics over the years, in this work we correctly encode a non-trivial logic like LRK and we automate it for the first time.

The structure of the paper is as follows: Section 2 briefly recaps HOL, and Sect. 3 sketches LRK. In Sect. 4 a shallow semantical embedding of LRK in HOL is presented and discussed, and faithfulness of the embedding is proved. Section 5 tests the embedding and uses it for automating the above GDPR example. Section 6 discusses related work and concludes the paper.

## 2 Classical Higher-Order Logic

Classical Higher-Order Logic (HOL), a polymorphic version of Church's *type theory* [21, 8], is sketched below (adapted from [5]).

**Syntax of HOL.** The *syntax* of HOL is defined by the grammar

$$s, t := P_\alpha \mid x_\alpha \mid (\lambda x_\alpha s_\beta)_{\alpha \Rightarrow \beta} \mid (s_{\alpha \Rightarrow \beta} t_\alpha)_\beta$$

where  $\alpha, \beta, o \in \mathcal{T}$  and with  $\mathcal{T}$  being a set of *simple types* defined by  $\alpha, \beta := o \mid i \mid (\alpha \Rightarrow \beta)$ . Type  $o$  denotes truth values,  $i$  individuals, and  $\Rightarrow$  is the function type constructor. The  $P_\alpha$  are typed constants symbols in the signature of HOL, and the  $x_\alpha$  are typed variable symbols (distinct from  $P_\alpha$ ). Complex HOL terms are constructed from given HOL terms via  $\lambda$ -abstraction  $(\lambda x_\alpha s_\beta)_{\alpha \Rightarrow \beta}$  and function application  $(s_{\alpha \Rightarrow \beta} t_\alpha)_\beta$ , which both involve type constraints. HOL is thus a logic of terms defined on top of the simply typed  $\lambda$ -calculus, and terms of type  $o$  are called *formulas*. The type of each term is given as a subscript and may be omitted if obvious in context. As *primitive logical connectives* we choose  $\neg_{o \Rightarrow o}, \vee_{o \Rightarrow o \Rightarrow o}, =_{\alpha \Rightarrow \alpha \Rightarrow o}$  (short:  $=^\alpha$ ) and  $\Pi_{(\alpha \Rightarrow o) \Rightarrow o}$  (short:  $\Pi^\alpha$ ). Other logical connectives can be introduced as abbreviations resp. shorthand notations; e.g.  $\forall x_\alpha \varphi_o = \Pi^\alpha (\lambda x_\alpha \varphi)$  and  $\rightarrow_{o \Rightarrow o \Rightarrow o} = \lambda x_o \lambda y_o (\neg x \vee y)$ .<sup>1</sup>

<sup>1</sup> Dot-notation is used in the remainder, where the scope opened by  $\cdot$  is reaching as far to the right as consistent with the formula structure. Moreover, types may be omitted if obvious in context; we may thus write  $\lambda x y. \neg x \vee y$  instead of  $\lambda x_o \lambda y_o (\neg x_o \vee y_o)$ .

**Semantics of HOL.** A *frame*  $\mathcal{D}$  for HOL is a collection  $\{\mathcal{D}_\alpha\}_{\alpha \in \mathcal{T}}$  of nonempty sets  $\mathcal{D}_\alpha$ , such that  $\mathcal{D}_o = \{T, F\}$  (for true and false).  $\mathcal{D}_i$  is chosen freely and  $\mathcal{D}_{\alpha \Rightarrow \beta}$  are collections of functions mapping  $\mathcal{D}_\alpha$  into  $\mathcal{D}_\beta$ . A *model* for HOL is a tuple  $\mathcal{M} = (\mathcal{D}, I)$ , where  $\mathcal{D}$  is a frame, and  $I$  is a family of typed interpretation functions mapping constant symbols  $P_\alpha$  to appropriate elements of  $\mathcal{D}_\alpha$ , called the *denotation* of  $P_\alpha$ . The logical connectives  $\neg, \vee, \Pi$  and  $=$  are always given their expected standard denotations:  $I(\neg) = \text{not} \in \mathcal{D}_{o \Rightarrow o}$  s.t.  $\text{not}(T) = F$  and  $\text{not}(F) = T$ ;  $I(\vee) = \text{or} \in \mathcal{D}_{o \Rightarrow o \Rightarrow o}$  s.t.  $\text{or}(a, b) = T$  iff  $(a = T \text{ or } b = T)$ ;  $I(=^\alpha) = \text{id} \in \mathcal{D}_{\alpha \Rightarrow \alpha \Rightarrow o}$  s.t.  $\forall a, b \in \mathcal{D}_\alpha, \text{id}(a, b) = T$  iff  $a$  is identical to  $b$ ;  $I(\Pi^\alpha) = \text{all} \in \mathcal{D}_{(\alpha \Rightarrow o) \Rightarrow o}$  s.t.  $\forall s \in \mathcal{D}_{\alpha \Rightarrow o}, \text{all}(s) = T$  iff  $s(a) = T$  for all  $a \in \mathcal{D}_\alpha$ .

A *variable assignment*  $g$  maps variables  $x_\alpha$  to elements in  $\mathcal{D}_\alpha$ .  $g[d/x]$  denotes the  $g'$ , that is identical to  $g$ , except for variable  $x$ , which is now mapped to  $d$ . The *denotation*  $\llbracket s_\alpha \rrbracket^{\mathcal{M}, g}$  of an HOL term  $s_\alpha$  on a model  $\mathcal{M} = \langle \mathcal{D}, I \rangle$  under assignment  $g$  is an element  $d \in \mathcal{D}_\alpha$  defined in the following way:

$$\begin{aligned} \llbracket P_\alpha \rrbracket^{\mathcal{M}, g} &= I(P_\alpha) \\ \llbracket x_\alpha \rrbracket^{\mathcal{M}, g} &= g(x_\alpha) \\ \llbracket (s_{\alpha \Rightarrow \beta} t_\alpha)_\beta \rrbracket^{\mathcal{M}, g} &= \llbracket s_{\alpha \Rightarrow \beta} \rrbracket^{\mathcal{M}, g} (\llbracket t_\alpha \rrbracket^{\mathcal{M}, g}) \\ \llbracket (\lambda x_\alpha s_\beta)_{\alpha \Rightarrow \beta} \rrbracket^{\mathcal{M}, g} &= \text{the } f: \mathcal{D}_\alpha \mapsto \mathcal{D}_\beta \text{ s.t. for all } d \in \mathcal{D}_\alpha \\ & f(d) = \llbracket s_\beta \rrbracket^{\mathcal{M}, g[d/x_\alpha]} \end{aligned}$$

It follows:  $\llbracket \forall x_\alpha \varphi_o \rrbracket^{\mathcal{M}, g} = T$  iff  $\llbracket \varphi_o \rrbracket^{\mathcal{M}, g[d/x_\alpha]} = T$  for all  $d \in \mathcal{D}_\alpha$ .

In a *standard model* a domain  $\mathcal{D}_{\alpha \Rightarrow \beta}$  is defined as the set of all total functions from  $\mathcal{D}_\alpha$  to  $\mathcal{D}_\beta$ :  $\mathcal{D}_{\alpha \Rightarrow \beta} = \{f \mid f: \mathcal{D}_\alpha \mapsto \mathcal{D}_\beta\}$ . In a *Henkin model* (or general model) [33] function spaces are not necessarily required to be the full:  $\mathcal{D}_{\alpha \Rightarrow \beta} \subseteq \{f \mid f: \mathcal{D}_\alpha \mapsto \mathcal{D}_\beta\}$ . However, it is required that every term still denotes.

Term  $s_o$  is *valid in  $\mathcal{M}$  under assignment  $g$* , denoted as  $\mathcal{M}, g \models^{\text{HOL}} s_o$ , iff  $\llbracket s_o \rrbracket^{\mathcal{M}, g} = T$ .  $s_o$  is *valid in  $\mathcal{M}$* , denoted as  $\mathcal{M} \models^{\text{HOL}} s_o$ , iff  $\mathcal{M}, g \models^{\text{HOL}} s_o$  for all assignments  $g$ , and  $s_o$  is *valid*, denoted as  $\models^{\text{HOL}} s_o$ , iff  $s_o$  is valid in all Henkin models  $\mathcal{M}$ .

Due to Gödel [31] a sound and complete mechanization of HOL with standard semantics cannot be achieved. For HOL with Henkin semantics sound and complete calculi exist, cf. [9, 12] and the references therein. Each standard model is obviously also a Henkin model. Consequently, when a HOL formula is Henkin-valid, it is also valid in all standard models.

## 3 A Dynamic Logic of the Right to Know

The Dynamic Logic of the Right to Know (LRK), introduced in [42], can be used to represent scenarios involving communication between two agents, where information is communicated directly from the sender to the receiver (indicated by  $s$  and  $r$ , respectively). Announcements can only be made by the sender, and the receiver may ask questions. The behavior of the sender is subject to additional restrictions, for instance coming from information security policies. The main objective of LRK is to characterize and reason about the receiver's so-called *power* [46] type of right to know whether  $\varphi$ , which means that the sender is obliged to announce the answer *if* the receiver asks the question whether  $\varphi$  is the case—the question itself creating the duty which did not exist before. The main notions of LRK are introduced as required for the rest of the paper; for details see [42].

**Syntax of LRK.** The *syntax* of LRK is defined by the grammar (where  $p \in \text{PROP}$  are atomic propositional symbols):

$$\varphi, \psi := p \mid \neg \varphi \mid \varphi \rightarrow \psi \mid U\varphi \mid Q\varphi \mid K_r \varphi \mid \mathbb{O}_s \varphi \mid [\varphi?] \psi \mid [\varphi!] \psi$$

Other Boolean connectives can be defined as usual, and the following abbreviations are introduced:  $R_r \varphi := U(Q\varphi \vee Q\neg\varphi)$ ,  $[r: \varphi?] \psi := (\neg R_r \varphi \wedge \psi) \vee (R_r \varphi \wedge [\varphi?] \psi)$ , and  $[s: \varphi!] \psi := (\varphi \rightarrow [\varphi!] \psi)$ .

In LRK,  $U\varphi$  is the familiar universal modality expressing that “ $\varphi$  is true in all worlds”, while  $Q\varphi$  is a technical modality that first appeared in [60]. The fundamental idea that “The receiver has the power to know the answer to the question  $\varphi$ ?” is captured by the formula  $R_r\varphi := U(Q\varphi \vee Q\neg\varphi)$ . The expression  $K_r\varphi$  means “The receiver knows  $\varphi$ ”, while  $\mathbb{O}_s\varphi$  stands for “The sender is obliged to announce  $\varphi$ ”. The technical operators  $[\varphi?]\varphi$  and  $[\varphi!]\varphi$  are used to define more complex operators. The concept that “After the receiver has asked the question  $\varphi?$ ,  $\psi$  holds” is represented by  $[r : \varphi?]\psi := (\neg R_r\varphi \wedge \psi) \vee (R_r\varphi \wedge [\varphi?]\psi)$ . Finally, the notion of truthful announcements is expressed by  $[s : \varphi!]\psi := (\varphi \rightarrow [\varphi!]\psi)$ , which reads as “After the sender truthfully announces  $\varphi$ , it follows that  $\psi$  holds”.

**Semantics of LRK.** An LRK model is a tuple  $M = (W, \sim, \approx, N, V)$ , where  $W$  is a non-empty set of possible worlds,  $\sim$  and  $\approx$  are two equivalence relations on  $W$ ,  $N : W \rightarrow \wp(W)$  is such that  $w \in X$  for all  $w \in W$  and  $X \in N(w)$ , and  $V : \text{PROP} \rightarrow \wp(W)$  is a valuation function that assigns a set of worlds to each atomic proposition. Vice versa, each world can be identified with the set of propositions that are validated in it. A pointed model is a pair  $M, w$  such that  $w$  is a state of  $M$ . For every state  $w \in W$ ,  $\sim(w)$  denotes the set  $\{v \in W \mid w \sim v\}$ , and analogously for  $\approx(w)$ . The relation  $\sim$  is the epistemic indistinguishability relation of the receiver. In what follows, the partition generated by the equivalence classes of  $\approx$  will also be denoted by  $\approx$ , and this partition  $\approx$  encodes the set of questions to which the receiver has the power to know the answers, as prescribed for instance by some given information security policies.  $N$  is a neighborhood function and each subset  $X \in N(w)$  is an *ideal epistemic state* for the receiver at  $w$ , i.e., the epistemic state  $X$  is compliant with the given policies specified at  $w$ .

Given a model  $M = (W, \sim, \approx, N, V)$ , for all  $w \in W$  and  $\varphi \in \mathcal{L}$ , the satisfaction relation  $M, w \models \varphi$  is inductively defined as:

$$\begin{aligned} M, w \models p & \quad \text{iff } w \in V(p) \\ M, w \models \neg\varphi & \quad \text{iff } M, w \not\models \varphi \\ M, w \models \varphi \rightarrow \psi & \quad \text{iff } M, w \not\models \varphi \text{ or } M, w \models \psi \\ M, w \models U\varphi & \quad \text{iff } \forall v \in W, M, v \models \varphi \\ M, w \models Q\varphi & \quad \text{iff } \forall v \in W, w \approx v \text{ implies } M, v \models \varphi \\ M, w \models K_r\varphi & \quad \text{iff } \forall v \in W, w \sim v \text{ implies } M, v \models \varphi \\ M, w \models \mathbb{O}_s\varphi & \quad \text{iff } \forall X \in N(w), X \subseteq \sim(w) \text{ implies } X \subseteq \llbracket \varphi \rrbracket_M \\ M, w \models [\varphi?]\psi & \quad \text{iff } M_{\varphi?}, w \models \psi \\ M, w \models [\varphi!]\psi & \quad \text{iff } M_{\varphi!}, w \models \psi \end{aligned}$$

where  $\llbracket \varphi \rrbracket_M = \{x \in W \mid M, x \models \varphi\}$  and  $M_{\varphi?}$  and  $M_{\varphi!}$  are defined as follows:

$$\begin{aligned} M_{\varphi?} & := (W, \sim, \approx, N_{\varphi?}, V) \text{ where for all } x \in W, N_{\varphi?}(x) = \\ & \{X \in N(x) \mid X \subseteq \llbracket \varphi \rrbracket_M \text{ or } X \subseteq \llbracket \neg\varphi \rrbracket_M\}; \text{ and} \\ M_{\varphi!} & := (W, \sim_{\varphi!}, \approx, N, V) \text{ where } \sim_{\varphi!} = \{(u, v) \in \sim \mid M, u \models \varphi \text{ iff } M, v \models \varphi\}. \end{aligned}$$

The notion of validity, written as  $\models^{LRK} \varphi$ , is defined as usual (cf. Sect. 2). In addition, if the above conditions on  $\sim$  and  $\approx$  to be equivalence relations are dropped and the extra condition on the neighborhood function  $N$  is removed,  $M = (W, \sim, \approx, N, V)$  is called an LRK<sup>-</sup> model; all other definitions apply analogously. An LRK<sup>-</sup> model thus constitutes a basic model structure, analogous to basic logic  $K$  in the well-known modal logic cube, which leaves the accessibility relation between worlds unrestricted with regard to further conditions such as reflexivity, symmetry, or transitivity.

The convention adopted by the authors of [42] is that an operator without an index serves only as a technical modality, such as  $U$  or  $Q$ . The LRK semantics of  $R_r\varphi$  is non-standard but straightforward: the receiver has the power to know the answer to  $\varphi?$  if and only if the partition  $\approx$  “settles” the question, meaning that each cell in  $\approx$  lies entirely within either the truth set of  $\varphi$  or  $\neg\varphi$ . Since  $Q\varphi$

behaves as a normal modality for  $\approx$ , this idea is precisely captured by  $U(Q\varphi \vee Q\neg\varphi)$ . The formula  $[r : \varphi?]\psi$  expresses that “after the receiver asks  $\varphi?$ ,  $\psi$  holds.” If the receiver does not have the power to know the answer, nothing changes after  $[r : \varphi?]$ . Otherwise, the sender is obliged to answer, and model updating ensures that non-answering epistemic states cease to be ideal. The dynamic operator  $[\varphi?]\psi$  captures what holds after this update, leading to the formalization of  $[r : \varphi?]\psi$  as  $(\neg R_r\varphi \wedge \psi) \vee (R_r\varphi \wedge [\varphi?]\psi)$ . The semantics of  $K_r\varphi$  and  $[s : \varphi!]\psi$  is standard, except that in  $M_{\varphi!}$  links are deleted between  $\varphi$  and  $\neg\varphi$ -states rather than removing  $\neg\varphi$ -states from the model. Note that  $N$  assigns, to each possible world  $w$ , a set of ideal epistemic states  $N(w)$ : for  $\mathbb{O}_s\varphi$ , the sender is obliged to announce  $\varphi$  if it is “known” in all ideal epistemic states achievable through further announcements, as  $X \subseteq \sim(w)$  says that the ideal epistemic state  $X$  must be achievable from the current epistemic state  $\sim(w)$ .

## 4 Modeling LRK as a Fragment of HOL

Before providing an embedding of LRK into HOL, a few words about the SSE approach employed here are in order. An SSE of a target logic into HOL provides a translation between the two logics in such a way that the former is identified and characterized as a proper fragment of the latter.<sup>2</sup> Once such an SSE is obtained, all that is needed to prove (or refute) conjectures in the target logic is to provide the SSE, encoded in an input file, to the HOL prover or model finder in addition to the encoded conjecture. The HOL tool can then be used as-is to solve problems in the target logic, i.e. without making any changes to its source code.

### 4.1 Shallow Semantical Embedding of LRK in HOL

In this subsection an SSE for target logic LRK in HOL is presented and faithfulness of the embedding is proved. To define our SSE, the type of LRK propositions is mapped to HOL type  $\tau := \gamma \Rightarrow \gamma \Rightarrow v \Rightarrow \sigma$ , where the following (additional) type abbreviations are used:  $\gamma := i \Rightarrow i \Rightarrow o$ ,  $v := i \Rightarrow ((i \Rightarrow o) \Rightarrow o)$  and  $\sigma := i \Rightarrow o$ . Note how HOL type  $\tau$  captures the dependencies of LRK formulas on two accessibility relations (of type  $\gamma$ ), a neighborhood function (of type  $v$ ) and on possible worlds (of type  $i$ ). These dependencies are visible in the definition of their semantical evaluation in Sect. 3, and in principle we could avoid the dependency of the second type  $\gamma$ , since only one of the accessibility relations is updated during semantical evaluation while the other one remains fixed. In general, this idea to explicitly capture and maintain relevant dependencies is analogous to previous work [6, 7, 11], but a bit more complicated here, because more dependencies need to be maintained and dynamically updated.

For each propositional symbol  $p^k$  of LRK, the corresponding HOL signature is assumed to contain a corresponding predicate constant symbol  $p_{i \Rightarrow o}^k$ , which is (rigidly) denoting the set of all those worlds in which  $p^k$  holds. The mapping  $\llbracket \cdot \rrbracket$  translates a LRK formula  $\varphi$  into a HOL term  $\llbracket \varphi \rrbracket$  of type  $\tau$ . As first and second arguments such a mapped HOL term  $\llbracket \varphi \rrbracket$  accepts indistinguishability relation terms  $t_\gamma$  and  $q_\gamma$ , and as third argument a neighborhood function term  $n_v$ . Finally, it takes a current world term  $w_i$  with respect to which the evaluation is performed. These dependencies are passed recursively through the recursive evaluation structure captured in the definitions below, and can be modified/updated on the fly. Such updates occur when either the current world in the evaluation changes (this happens for the modal connectives  $U$ ,  $Q$ ,  $K_r$ , and  $\mathbb{O}_s$ ), or when either

<sup>2</sup> The SSE technique is not to be confused with higher-order abstract syntax [54]. For differences between shallow and deep embeddings see [29, 19].

the neighborhood function (see  $[?\varphi]\psi$ ) or the indistinguishability relation  $\sim$  (see  $[\varphi]\psi$ ) is modified. Since only one of the two indistinguishability relations is modified recursively, while the other is always kept fixed, the captured dependencies are not minimal and could be reduced to just one indistinguishability relation of type  $\tau$ , but we decided for the sake of clarity to represent both of them explicitly. As mentioned in Sect. 3, the connective  $Q$  is a technical (normal) modality that is needed to define  $R_\tau$ . We keep this approach for the SSE: since the semantics for  $R_\tau\varphi$  is highly non-standard, the embedding and proof would be much harder if  $R_\tau\varphi$  were primitive in the language.

The mapping  $[\cdot]$  is defined recursively as follows:

$$\begin{aligned} [p^k] &= A_{\sigma \Rightarrow \tau}(p^k) \text{ with } A_{\sigma \Rightarrow \tau} = \lambda p_\sigma t_\gamma q_\gamma n_v w_i. p w \\ [\neg\varphi] &= \neg_{\tau \Rightarrow \tau} [\varphi] \text{ with } \neg_{\tau \Rightarrow \tau} = \lambda \varphi_\tau t_\gamma q_\gamma n_v w_i. \neg(\varphi t q n w) \\ [\varphi \rightarrow \psi] &= \Rightarrow_{\tau \Rightarrow \tau} [\varphi] [\psi] \text{ with } \Rightarrow_{\tau \Rightarrow \tau} = \\ &\quad \lambda \varphi_\tau \psi_\tau t_\gamma q_\gamma n_v w_i. \neg \varphi t q n w \vee \psi t q n w \\ [U\varphi] &= U_{\tau \Rightarrow \tau} [\varphi] \text{ with } U_{\tau \Rightarrow \tau} = \lambda \varphi_\tau t_\gamma q_\gamma n_v w_i. \forall v_i. \varphi t q n w \\ [Q\varphi] &= Q_{\tau \Rightarrow \tau} [\varphi] \text{ with } Q_{\tau \Rightarrow \tau} = \\ &\quad \lambda \varphi_\tau t_\gamma q_\gamma n_v w_i. \forall v_i. q w v \longrightarrow \varphi t q n w \\ [Kr\varphi] &= Kr_{\tau \Rightarrow \tau} [\varphi] \text{ with } Kr_{\tau \Rightarrow \tau} = \\ &\quad \lambda \varphi_\tau t_\gamma q_\gamma n_v w_i. \forall v_i. t w v \longrightarrow \varphi t q n w \\ [\mathbb{O}_s\varphi] &= \mathbb{O}_{s\tau \Rightarrow \tau} [\varphi] \text{ with } \mathbb{O}_{s\tau \Rightarrow \tau} = \lambda \varphi_\tau t_\gamma q_\gamma n_v w_i. \forall H_\sigma. \\ &\quad n w H \longrightarrow (\forall v_i. H v \longrightarrow t w v) \longrightarrow (\forall u_i. H u \longrightarrow \varphi t s n u) \\ [[\varphi?]\psi] &= ?_{\tau \Rightarrow \tau} [\varphi] [\psi] \text{ with } ?_{\tau \Rightarrow \tau} = \\ &\quad \lambda \varphi_\tau \psi_\tau t_\gamma q_\gamma n_v w_i. \psi t q (\text{UPDATE}N\varphi t q n) w \\ [[\varphi!]\psi] &= !_{\tau \Rightarrow \tau} [\varphi] [\psi] \text{ with } !_{\tau \Rightarrow \tau} = \\ &\quad \lambda \varphi_\tau \psi_\tau t_\gamma q_\gamma n_v w_i. \psi (\text{UPDATE}T\varphi t q n) q n w \end{aligned}$$

where  $\text{UPDATE}N\varphi t q n = \lambda w_i X_\sigma. n w X \wedge ((\forall v_i. X v \longrightarrow \varphi t q n w) \vee (\forall v_i. X v \longrightarrow \neg \varphi t q n w))$  and  $\text{UPDATE}T\varphi t q n = \lambda u_i v_i. t w v \wedge (\varphi t q n u \longleftrightarrow \varphi t q n v)$ .

We first prove faithfulness of the embedding for the more general  $\text{LRK}^-$  models and then extend the result to  $\text{LRK}$  models. This two step approach is not only technically simpler, it also shows that faithfulness results can be obtained for a wider class of models, independent of the particular constraints assumed for  $\sim$ ,  $\approx$  and  $N$ .

**Lemma 1** (Faithfulness:  $\text{LRK}^-$  models). *Consider an  $\text{LRK}^-$  model  $M = (W, \sim, \approx, N, V)$  and a HOL model  $\mathcal{M} = (\mathcal{D}, I)$  such that*

1. *the set of worlds  $W$  in  $M$  is identified with the domain  $\mathcal{D}_i$  in  $\mathcal{M}$ ,*
2. *for each  $p^k \in \text{PROP}$  of  $\text{LRK}$ , there is a constant symbol  $p_\sigma^k$  in HOL s.t. for all  $w \in W = \mathcal{D}_i$ :  $I(p_\sigma^k)(w) = T$  iff  $w \in V(p^k)$ ,*
3. *the HOL signature contains  $\sim_\gamma, \approx_\gamma$  and  $N_v$  s.t.  $I(\sim_\gamma) = \sim$ ,  $I(\approx_\gamma) = \approx$  and  $I(N_v) = N$ ;  $\sim$  and  $\approx$  are the mentioned  $\text{LRK}$  equivalence relations, and  $N$  is s.t.  $w \in X$  for all  $w \in W$  and  $X \in N(w)$ , resp.  $I(N_v)$  satisfies  $\forall w_i. \forall X_\sigma. (N w X \longrightarrow X w)$ .*

*Then, for all  $w \in W$ , and  $X_i$  not occurring free in  $[\varphi]$ , we have:*

$$M, w \models^{\text{LRK}} \varphi \text{ iff } [[[\varphi] \sim_\gamma \approx_\gamma N_v X_i]]^{\mathcal{M}, g[w/X_i]} = T \text{ in HOL.}$$

*Proof.* By induction over  $\text{LRK}$  formula  $\varphi$ ; we only sketch it here.

In the base case, where  $\varphi$  is  $p^k \in \text{PROP}$ , the statement follows from semantical evaluation and assumption 2.

The cases for  $\neg\varphi$  and  $\varphi \vee \psi$  follow from the induction hypothesis.

In the case for  $U\varphi$  we need to show  $M, w \models U\varphi$  in  $\text{LRK}$  iff  $[[[U\varphi] \sim_\gamma \approx_\gamma N_v X_i]]^{\mathcal{M}, g[w/X_i]} = T$  in HOL. Since  $[U\varphi]$   $\lambda$ -converts into  $\lambda t_\gamma q_\gamma n_v w_i. \forall v_i. [\varphi] t q n v$ , the latter is equivalent to  $[[\forall v_i. [\varphi] \sim_\gamma \approx_\gamma N_v v_i]]^{\mathcal{M}, g[w/X_i]} = T$ , respectively to  $[[\forall v_i. [\varphi] \sim_\gamma \approx_\gamma N_v v_i]]^{\mathcal{M}, g} = T$ , since  $X_i$  is not occurring in  $[\varphi]$  by assumption. This is equivalent to  $[[[\varphi] \sim_\gamma \approx_\gamma N_v X_i]]^{\mathcal{M}, g[w/X_i]} = T$  for all  $w \in D_i$  by the semantics of HOL, which by induction hypothesis, and since  $D_i = W$  by assumption 1, is equivalent to  $M, w \models \varphi$  for all  $w \in W$  in  $\text{LRK}$ . By the semantics of  $\text{LRK}$  it

follows  $M, w \models U\varphi$  as intended.

The cases for  $Q\varphi$ ,  $Kr\varphi$ , and  $\mathbb{O}_s\varphi$  follow analogously by semantic evaluation, induction and the application of the assumptions.

In the case for  $[\varphi!]\psi$  we need to show  $M, w \models [\varphi!]\psi$  in  $\text{LRK}$  iff  $[[[[\varphi!]\psi] \sim_\gamma \approx_\gamma N_v X_i]]^{\mathcal{M}, g[w/X_i]} = T$  in HOL. This is equivalent to  $[[[\psi] (\text{UPDATE}T[\varphi] \sim_\gamma \approx_\gamma N_v) \approx_\gamma N_v X_i]]^{\mathcal{M}, g[w/X_i]} = T$ , since  $[[[\varphi!]\psi]$   $\lambda$ -converts into  $\lambda t_\gamma q_\gamma n_v w_i. [\psi] (\text{UPDATE}T[\varphi] t q n) q n w$ . By definition,  $\text{UPDATE}T[\varphi] \sim_\gamma \approx_\gamma N_v$   $\lambda$ -converts into  $\lambda u_i v_i. \sim_\gamma u v \wedge ([\varphi] \sim_\gamma \approx_\gamma N_v u \longleftrightarrow [\varphi] \sim_\gamma \approx_\gamma N_v v)$ , which by the induction hypothesis on  $\varphi$  and the assumptions corresponds to the updated accessibility relation  $\sim_{\varphi!} = \{(u, v) \in \sim \mid M, u \models \varphi \text{ iff } M, v \models \varphi\}$ . By this, the assumptions and the induction hypothesis on  $\psi$  we have  $[[[\psi] (\text{UPDATE}T[\varphi] \sim_\gamma \approx_\gamma N_v) \approx_\gamma N_v X_i]]^{\mathcal{M}, g[w/X_i]} = T$  equivalent to  $M_{\varphi!}, w \models \psi$ .

The case for  $[\varphi?]\psi$  follows analogously.  $\square$

**Theorem 2** (Faithfulness:  $\text{LRK}$  models). *Suppose that the requirements given for Lemma 1 are satisfied for all considered  $\text{LRK}$  models  $M$  and HOL models  $\mathcal{M}$ . Furthermore, let  $A_X$  be a set of HOL axioms postulating the relation symbols  $\sim_\gamma$  and  $\approx_\gamma$  (the HOL counterparts of  $\sim$  and  $\approx$ ) to denote equivalence relations, and  $N_v$  (the HOL counterpart of  $N$ ) to satisfy the condition  $\forall w_i. \forall U_\sigma. (n w U \longrightarrow U w)$ . Then, we have:  $\models^{\text{LRK}} \varphi$  iff  $A_X \models^{\text{HOL}} [\varphi]$*

*Proof.* Corollary of Lemma 1; the additional  $\text{LRK}$ -conditions on  $\sim$ ,  $\approx$  and  $N$  are enforced by the axioms  $A_X$  postulated in HOL.  $\square$

## 4.2 Encoding into Isabelle/HOL

The presented SSE of  $\text{LRK}$  in HOL has been encoded in the proof assistant *Isabelle/HOL* to enable automated reasoning with  $\text{LRK}$  for the first time.<sup>3</sup> All necessary types can be modeled in a straightforward way. We declare type  $i$  (using **typedef** in Isabelle) to denote possible worlds and then introduce type aliases (using **type synonym**) for  $\sigma, \gamma, v$  and  $\tau$ , as mentioned in Sect.4.1. Type **bool** represents (the bivalent set of) truth values. Type  $\sigma$  is the type for formulas independent from a context, and type  $\tau$  relates formulas to their context, which consists of an accessibility and a neighborhood relation.  $\gamma$  is the type of the functions that are used to update the accessibility and neighborhood relations.

The relations  $\approx$ ,  $\sim$ , and  $N$  are declared as infix constant symbols.  $\approx$  defines the set of questions the receiver is allowed to know the answers to, whereas  $\sim$  is the usual indistinguishability relation for the receiver. The ideal epistemic state of the receiver is captured via the relation  $N$ .

In order for  $\approx$  and  $\sim$  to be equivalence relations, we constrain them to be reflexive, transitive, and symmetric:

**axiomatization** where **til\_ref**: " $\forall w. w \sim w$ " and **til\_sym**: " $\forall w v. w \sim v \longrightarrow v \sim w$ " and **til\_trans**: " $\forall w v u. w \sim v \wedge v \sim u \longrightarrow w \sim u$ " and **dtil\_ref**: " $\forall w. w \approx w$ " and **dtil\_sym**: " $\forall w v. w \approx v \longrightarrow v \approx w$ " and **dtil\_trans**: " $\forall w v u. w \approx v \wedge v \approx u \longrightarrow w \approx u$ " and **Nax**: " $\forall w. \forall H. (N w H \longrightarrow H w)$ "

The axiom **Nax** constrains the neighborhood relation  $N$  by stating that all sets of worlds that are in the set of sets of worlds returned by  $(N w)$  contain  $w$  itself. This implies that none of the sets is empty.

The operator for atomic propositions  $A(\cdot)$  of type  $\sigma \Rightarrow \tau$  checks whether an atom  $p$  holds in a given world  $w$ :

**abbreviation** atom: " $\sigma \Rightarrow \tau$ " (" $A\_$ ") where  $A_p \equiv \lambda t q n w. p w$ "

<sup>3</sup> The full sources of our encoding can be found at <http://logikey.org> in subfolder **LRK** (<https://github.com/cbenzmueller/LogiKEy/tree/master/LRK>).

The formulas  $\top$  and  $\perp$  are defined represent truth and falsity (independent of the given world):

**abbreviation**  $\text{lrkTop}::\tau$  (" $\top$ ") **where** " $\top \equiv \lambda t q n w. \text{True}$ "

**abbreviation**  $\text{lrkBot}::\tau$  (" $\perp$ ") **where** " $\perp \equiv \lambda t q n w. \text{False}$ "

We then define the Boolean connectives that are primitive in LRK. To distinguish between HOL connectives (e.g.,  $\neg$ ) and the lifted LRK connectives (e.g.,  $\neg\tau \Rightarrow \tau$ ), we use boldface fonts:

**abbreviation**  $\text{lrkNot}::\tau \Rightarrow \tau$  (" $\neg$ ")

**where** " $\neg\varphi \equiv \lambda t q n w. \neg(\varphi t q n w)$ "

**abbreviation**  $\text{lrkImp}::\tau \Rightarrow \tau \Rightarrow \tau$  (" $\rightarrow$ ")

**where** " $\varphi \rightarrow \psi \equiv \lambda t q n w. (\varphi t q n w) \rightarrow (\psi t q n w)$ "

For convenience, we also define the remaining Boolean connectives as abbreviations.

**abbreviation**  $\text{lrkAnd}::\tau \Rightarrow \tau \Rightarrow \tau$  (" $\wedge$ ")

**where** " $\varphi \wedge \psi \equiv \lambda t q n w. (\varphi t q n w) \wedge (\psi t q n w)$ "

**abbreviation**  $\text{lrkOr}::\tau \Rightarrow \tau \Rightarrow \tau$  (" $\vee$ ")

**where** " $\varphi \vee \psi \equiv \lambda t q n w. (\varphi t q n w) \vee (\psi t q n w)$ "

**abbreviation**  $\text{lrkEquiv}::\tau \Rightarrow \tau \Rightarrow \tau$  (" $\leftrightarrow$ ")

**where** " $\varphi \leftrightarrow \psi \equiv \lambda t q n w. (\varphi t q n w) \leftrightarrow (\psi t q n w)$ "

To build formulas in LRK, embeddings of the specific connectives of the logic are needed (see Sect. 3). Here,  $U$  is the universal modality,  $Rr$  represents the right to know of the receiver,  $Kr$  the knowledge of the receiver,  $Os$  the obligation of the sender to announce something truthfully,  $Q$  is a technical modality.

**abbreviation**  $\text{lrkU}::\tau \Rightarrow \tau$  (" $U$ ")

**where** " $U\varphi \equiv \lambda t q n w. \forall v. \varphi t q n v$ "

**abbreviation**  $\text{lrkQ}::\tau \Rightarrow \tau$  (" $Q$ ")

**where** " $Q\varphi \equiv \lambda t q n w. \forall v. (q w v) \rightarrow (\varphi t q n v)$ "

**abbreviation**  $\text{lrkRr}::\tau \Rightarrow \tau$  (" $Rr$ ")

**where** " $Rr\varphi \equiv \lambda t q n w. \forall va. \forall v. (q va v) \rightarrow$

$(\varphi t q n v) \vee (\forall v. q va v \rightarrow (\neg\varphi t q n v))$ "

**abbreviation**  $\text{lrkKr}::\tau \Rightarrow \tau$  (" $Kr$ ")

**where** " $Kr\varphi \equiv \lambda t q n w. \forall v. (t w v) \rightarrow (\varphi t q n v)$ "

**abbreviation**  $\text{lrkOs}::\tau \Rightarrow \tau$  (" $Os$ ") **where**

" $Os\varphi \equiv \lambda t q n w. \forall H. (n w H) \rightarrow$

$(\forall v. H v \rightarrow (t w v)) \rightarrow (\forall i. H i \rightarrow (\varphi t q n i))$ "

Formulas  $[\varphi?]\varphi$  and  $[\varphi!]\varphi$  are defined using the dynamic updates mentioned above, and they are then used to encode  $[r : \varphi?]\psi$  and  $[s : \varphi!]\psi$ :

**abbreviation**  $\text{lrkQuestion}::\tau \Rightarrow \tau \Rightarrow \tau$  (" $[?]\_$ ")

**where** " $[\varphi?]\psi \equiv \lambda t q n w. (\psi t q (\text{update\_N } \varphi n t q w))$ "

**abbreviation**  $\text{lrkExclamation}::\tau \Rightarrow \tau \Rightarrow \tau$  (" $[!]\_$ ")

**where** " $[\varphi!]\psi \equiv \lambda t q n w. (\psi (\text{update\_t } \varphi n t q) q n w)$ "

**abbreviation**  $\text{lrkAfterQuestion}::\tau \Rightarrow \tau \Rightarrow \tau$  (" $[r : ?]\_$ ")

**where** " $[r : \varphi?]\psi \equiv (((\neg Rr\varphi) \wedge \psi) \vee ((Rr\varphi) \wedge [\varphi?]\psi))$ "

**abbreviation**  $\text{lrkAfterExclamation}::\tau \Rightarrow \tau \Rightarrow \tau$  (" $[s : !]\_$ ")

**where** " $[s : \varphi!]\psi \equiv \varphi \rightarrow [\varphi!]\psi$ "

Update of the accessibility relation and neighborhood function is:

**abbreviation**  $\text{update\_t}::\tau \Rightarrow v \Rightarrow \gamma \Rightarrow \gamma \Rightarrow \gamma$

**where** " $\text{update\_t } \varphi n t q \equiv \lambda u v. t u v \wedge (\varphi t q n u \leftrightarrow \varphi t q n v)$ "

**abbreviation**  $\text{update\_N}::\tau \Rightarrow v \Rightarrow \gamma \Rightarrow \gamma \Rightarrow v$

**where** " $\text{update\_N } \varphi n t q \equiv \lambda w. \lambda X. n w X \wedge$

$((\forall v. (X v \rightarrow (\varphi t q n v))) \vee (\forall v. (X v \rightarrow ((\neg\varphi) t q n v))))$ "

Finally, the notion of validity is embedded, and the initial accessibility relations  $\approx$ ,  $\sim$  and neighborhood function  $N$  are used for this. Their value is then updated when the operators  $[\varphi?]\varphi$  and  $[\varphi!]\varphi$  are used in the evaluated formula.

**abbreviation**  $\text{lrkValidLocal}::\tau \Rightarrow i \Rightarrow \text{bool}$  (" $[?]\_$ ")

**where** " $[?]\_w \equiv \varphi \text{ til } \text{dtil } n w$ "

**abbreviation**  $\text{lrkValidGlobal}::\tau \Rightarrow \text{bool}$  (" $[?]\_$ ")

**where** " $[?]\_ \equiv \forall w. \varphi \text{ til } \text{dtil } n w$ "

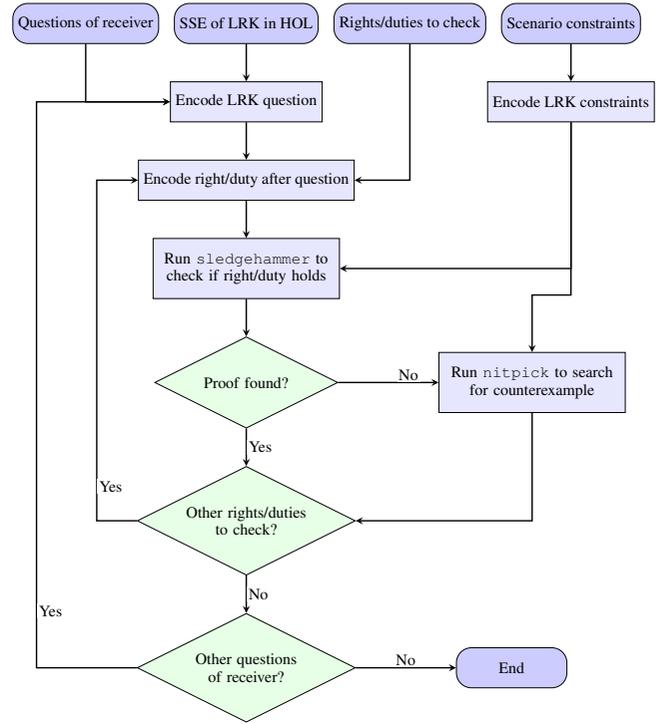


Figure 2. Workflow diagram

## 5 Case Study: GDPR Scenario

In a first experiment, the SSE of LRK in HOL has been tested by checking that relevant propositions, semantic results and the axioms postulated for LRK in [42] are actually implied in the embedded logic as expected; see the online repository mentioned in Footnote 3.

A second experiment demonstrates the use of the embedding on the GDPR example introduced in Sect. 1. This example has been represented in LRK and encoded in *Isabelle/HOL* using the shallowly embedded logic. Recall that according to the GDPR, a data subject has the right to access some specific information: for instance, whether her personal data is being processed, and if so, with what purpose [26]. However, this right should not adversely affect the rights or freedoms of others, such as the copyright protecting the software used to process data.

Consider the following situation: a data controller processes the personal data of a data subject with the purpose of optimizing advertisement placement. In doing so, the controller uses proprietary software using mechanism  $X$ . When the data subject asks the controller whether her personal data is processed and if so, whether this happens in order to optimize what advertisement she sees on the website, the questions create the controller's duty to answer them. However, due to the software being proprietary, the controller has a duty to not tell whether it applies mechanism  $X$ . Hence when the data subject asks about it, according to the GDPR, this question does not establish the controller's duty to tell.

The workflow is portrayed in Fig. 2: for each question of the data subject  $ds$  to the data controller  $dc$ , we check which rights or duties hold using the SSE of LRK in HOL on the encoded scenario. Results of proven and disproven rights and duties after each new question are reported in Table 1. The proposition "the data is processed by  $dc$ " is represented by propositional variable  $pa$ ; "the data is processed with the purpose of optimizing advertisement placement" is represented by  $pb$ ; and the proposition "the software used by  $dc$  uses mechanism

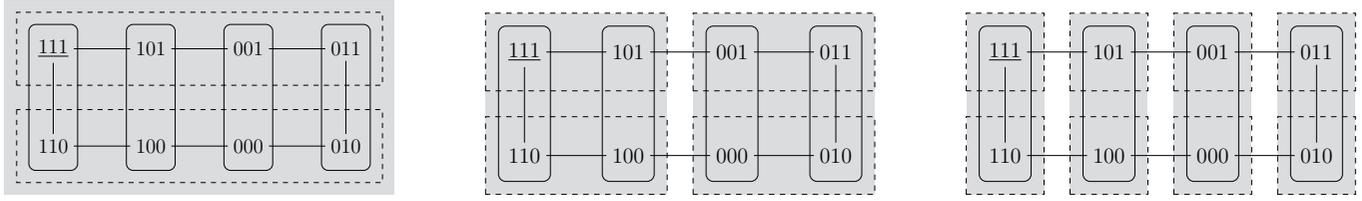


Figure 3. Initial model (left) —  $pa?$  Updated Model (middle) —  $pa?pb?$  Updated Model (right)

$X$  is represented by  $pc$ . The situation is displayed in Fig. 3 (left): states are identified by 3-bits sequences, corresponding to the value of  $pa$ ,  $pb$ , and  $pc$ , respectively. We are in state 111, where all three propositions are true. The indistinguishability relation  $\sim$  is indicated by the straight line, whereas the equivalence relation  $\approx$  (the right of  $ds$  to know whether  $pa$ , and whether  $pb$ ) is pictured by the rectangles with the rounded corners. Finally, for every state  $s$ , neighborhood function  $Ns$  contains every subset  $H$  such that: (1)  $H$  contains  $s$  itself; (2)  $H$  is contained in a shaded area (corresponding to what the  $dc$  is obligated to announce); and (3)  $H$  is not contained in one of the dashed rectangles ( $ds$  is prohibited to know  $pc$ ).

We start to encode the scenario by declaring constants of type  $\sigma$  for the three propositions and constants of type  $i$  for the eight different worlds, each representing one combination of values of  $pa$ ,  $pb$ , and  $pc$ . The worlds are named with letters from  $a$  to  $h$ , where  $a$  corresponds to assignment  $\{pa = 0, pb = 0, pc = 0\}$  and  $h$  corresponds to  $\{pa = 1, pb = 1, pc = 1\}$ .

**consts**  $pa:\sigma$   $pb:\sigma$   $pc:\sigma$   $a:i$   $b:i$   $c:i$   $d:i$   $e:i$   $f:i$   $g:i$   $h:i$

Next, this initial model is represented in the embedded logic. This includes ensuring the distinctness of the eight worlds, as well as declaring the truth values of each of the propositions in each possible world (axioms  $aa$ – $ah$ ). Moreover, we define the elements that are included in the accessibility relations  $\sim$  and  $\approx$  (axioms  $til1$ – $til8$ ,  $dtil1$ – $dtil10$ ) and neighborhood function (axiom Neigh). The  $\sim$  relation initially contains all possible pairs of worlds. Since  $\sim$  is reflexive, transitive, and symmetric, not all pairs of worlds need to be listed explicitly. To define  $N$  and  $\approx$ , we rely on what  $ds$  is allowed to know, or forbidden to know.

**axiomatization where**

$a0: "(\forall x. (x=a \vee x=b \vee x=c \vee x=d \vee x=e \vee x=f \vee x=g \vee x=h))"$  **and**  
 $a00: "(\forall x. (x=a \rightarrow \neg(x=b \vee x=c \vee x=d \vee x=e \vee x=f \vee x=g \vee x=h))"$  **and**  
 $a01: "(\forall x. (x=b \rightarrow \neg(x=a \vee x=c \vee x=d \vee x=e \vee x=f \vee x=g \vee x=h))"$  **and**  
 $a02: "(\forall x. (x=c \rightarrow \neg(x=a \vee x=b \vee x=d \vee x=e \vee x=f \vee x=g \vee x=h))"$  **and**  
 $a03: "(\forall x. (x=d \rightarrow \neg(x=a \vee x=b \vee x=c \vee x=e \vee x=f \vee x=g \vee x=h))"$  **and**  
 $a04: "(\forall x. (x=e \rightarrow \neg(x=a \vee x=b \vee x=c \vee x=d \vee x=f \vee x=g \vee x=h))"$  **and**  
 $a05: "(\forall x. (x=f \rightarrow \neg(x=a \vee x=b \vee x=c \vee x=d \vee x=e \vee x=g \vee x=h))"$  **and**  
 $a06: "(\forall x. (x=g \rightarrow \neg(x=a \vee x=b \vee x=c \vee x=d \vee x=e \vee x=f \vee x=h))"$  **and**  
 $a07: "(\forall x. (x=h \rightarrow \neg(x=a \vee x=b \vee x=c \vee x=d \vee x=e \vee x=f \vee x=g))"$  **and**  
 $aa: "\neg[Apa]_a \wedge \neg[Apb]_a \wedge \neg[Apc]_a"$  **and**  
 $ab: "\neg[Apa]_b \wedge \neg[Apb]_b \wedge [Apc]_b"$  **and**  
 $ac: "\neg[Apa]_c \wedge [Apb]_c \wedge \neg[Apc]_c"$  **and**  
 $ad: "\neg[Apa]_d \wedge [Apb]_d \wedge [Apc]_d"$  **and**  
 $ae: "[Apa]_e \wedge \neg[Apb]_e \wedge \neg[Apc]_e"$  **and**  
 $af: "[Apa]_f \wedge \neg[Apb]_f \wedge [Apc]_f"$  **and**  
 $ag: "[Apa]_g \wedge [Apb]_g \wedge \neg[Apc]_g"$  **and**  
 $ah: "[Apa]_h \wedge [Apb]_h \wedge [Apc]_h"$  **and**  
 $til1: "a \sim b"$  **and**  $til2: "b \sim c"$  **and**  $til3: "c \sim d"$  **and**  
 $til4: "d \sim e"$  **and**  $til5: "e \sim f"$  **and**  $til6: "f \sim g"$  **and**  
 $til7: "g \sim h"$  **and**  $til8: "h \sim a"$  **and**  $dtil1: "a \approx b"$  **and**  
 $dtil2: "c \approx d"$  **and**  $dtil3: "e \approx f"$  **and**  $dtil4: "g \approx h"$  **and**

Questions	Proven rights/duties	Disproven rights/duties
–	$R_r pa$ $R_r pb$	$R_r pc$ $\mathcal{O}_s pa$ $\mathcal{O}_s pb$ $\mathcal{O}_s pc$
$pa?$	$[r:pa?] R_r pa$ $[r:pa?] R_r pb$ $[r:pa?] \mathcal{O}_s pa$	$[r:pa?] R_r pc$ $[r:pa?] \mathcal{O}_s pb$ $[r:pa?] \mathcal{O}_s pc$
$pa?pb?$	$[r:pb?]( [r:pa?] R_r pa )$ $[r:pb?]( [r:pa?] R_r pb )$ $[r:pb?]( [r:pa?] \mathcal{O}_s pa )$ $[r:pb?]( [r:pa?] \mathcal{O}_s pb )$	$[r:pb?]( [r:pa?] R_r pc )$ $[r:pb?]( [r:pa?] \mathcal{O}_s pc )$
$pa?pc?$	$[r:pc?]( [r:pa?] R_r pa )$ $[r:pc?]( [r:pa?] R_r pb )$ $[r:pc?]( [r:pa?] \mathcal{O}_s pa )$	$[r:pc?]( [r:pa?] R_r pc )$ $[r:pc?]( [r:pa?] \mathcal{O}_s pb )$ $[r:pc?]( [r:pa?] \mathcal{O}_s pc )$

Table 1. Proven and disproven rights/duties at each new question

$dtil5: "\neg(a \approx c)"$  **and**  $dtil6: "\neg(a \approx e)"$  **and**  $dtil7: "\neg(a \approx g)"$  **and**  
 $dtil8: "\neg(c \approx e)"$  **and**  $dtil9: "\neg(c \approx g)"$  **and**  $dtil10: "\neg(e \approx g)"$  **and**  
 Neigh: " $\forall w. \forall H. ((N w) H \leftrightarrow (H w \wedge \neg(\text{check\_subset } H (\lambda x. (x=a \vee x=c \vee x=e \vee x=g)))) \wedge (\neg(\text{check\_subset } H (\lambda x. (x=b \vee x=d \vee x=f \vee x=h))))))"$

In the initial model (Fig. 3 (left); Table 1 (first row)), the following LRK literals for rights hold in world  $h$ :  $R_r pa$  ( $ds$  has the right to know whether  $pa$ ) and  $R_r pb$  ( $ds$  has the right to know whether  $pb$ ), while  $R_r pc$  does not hold ( $ds$  does not have the right to know  $pc$ ). In the encoded SSE, we prove that  $[R_r^A pa]_h$  and  $[R_r^A pb]_h$  hold using sledgehammer, and using nitpick we find counterexamples for  $[R_r^A pc]_h$ ,  $[O_s^A pa]_h$ ,  $[O_s^A pb]_h$  and  $[O_s^A pc]_h$  ( $dc$  has no obligation yet to announce anything).

Next,  $ds$  asks whether  $pa$  holds and we show that the updated model changes dynamically as expected; see Fig. 3 (middle) and Table 1 (second row). In the  $pa?$  updated model,  $ds$  (still) has the right to know whether  $pa$  holds and now it is obligatory for  $dc$  to announce it. Also,  $ds$  (still) has the right to know whether  $pb$  holds, but not whether  $pc$  holds. In the encoded SSE, we use sledgehammer to prove that  $[r:pa?] R_r^A pa$ ,  $[r:pa?] R_r^A pb$ , and  $[r:pa?] O_s^A pa$  hold.

Then, if now  $ds$  asks whether  $pb$  holds, the model is again updated; see Fig. 3 (right) and Table 1 (third row). In the  $pb?$  (after  $pa?$ ) updated model,  $ds$  still has the right to know  $pa$  and  $pb$ , and no right to know  $pc$ . The obligation for  $dc$  to announce  $pa$  remains, and now it is also obligatory for  $dc$  to announce  $pb$ . We find proofs for these formulas to hold in the encoded SSE:  $[r:pa?] ([r:pa?] R_r^A pa)_h$ ,  $[r:pa?] ([r:pa?] R_r^A pb)_h$ ,  $[r:pa?] ([r:pa?] O_s^A pa)_h$ ,  $[r:pb?] ([r:pa?] O_s^A pb)_h$ .

However, if after asking whether  $pa$  holds,  $ds$  demands to know also  $pc$ , no obligation for  $dc$  to announce  $pc$  arises, since  $ds$  does not have the right to know  $pc$ . Hence, after this question the model is still as in Fig. 3 (middle); see also Table 1 (fourth row). We

show that asking whether  $pc$  holds does not result in an obligation for  $dc$  to announce it, while the right of  $ds$  to know  $pa$  and  $pb$ , but not  $pc$ , as well as the obligation for  $dc$  to announce  $pa$  remain. As expected, the LRK formula  $[r : pc?](\{r : pa?\} \circ_s pc)$  stating that after such questions  $dc$  is obligated to announce  $pc$  is disproven: in the encoded SSE `nitpick` provides a counterexample to  $\llbracket [r : pc?](\{r : pa?\} \circ_s pc) \rrbracket_h$ , where  $Nh = \{g, h\}$ . To understand this, notice that in LRK the formula unfolds to  $(\neg R_r pc) \wedge ([r : pa?](\circ_s pc) \vee ((R_r pc) \wedge ([pc?](\{r : pa?\} \circ_s pc))))$ . Since  $ds$  does not have the right to know  $pc$ ,  $\neg R_r pc$  holds and we consider the first part of the disjunction  $(\neg R_r pc) \wedge ([r : pa?](\circ_s pc))$ . Formula  $([r : pa?](\circ_s pc))$  corresponds to  $((\neg R_r pa) \wedge (\circ_s pc)) \vee (R_r pa \wedge [pa?](\circ_s pc))$ .  $R_r pa$  holds, so we only need to check that  $[pa?](\circ_s pc)$  is falsified by the provided counterexample with  $Nh = \{g, h\}$ . Notice that  $\circ_s pc$  would imply that since  $g \sim h$ , then  $pc$  holds for all worlds in  $\{g, h\}$ . However,  $pc$  does not hold in world  $g$ , which corresponds to the evaluation  $pa = 1$ ,  $pb = 1$ ,  $pc = 0$ . Consequently, the original LRK formula  $[r : pc?](\{r : pa?\} \circ_s pc)$  does not hold in the current world.

## 6 Related Work and Conclusion

The right to know is a form of legal right called an *epistemic right* [62]. Implementing logics for epistemic rights is increasingly important for AI systems, as they now frequently operate in domains where access to information is governed by legal, ethical, or institutional norms [59]: an intelligent agent must not only manage and reason about knowledge but also respect who is entitled to know what and when, for example when handling health data [28]. Studies on the logical analysis of legal rights include [37, 44, 45, 36], with more recent works exploring power-type rights in [46, 58, 25]. These works focus on general rights, while a specific logic for epistemic rights is LRK, the Dynamic Logic of the Right to Know [42], that we mechanized in this paper. One reason such a specific logic is needed is that certain reasoning patterns are valid for epistemic rights but not for general rights.<sup>4</sup> As noted by [25], there are two main approaches to formalizing legal power: the first [37, 44] defines power in terms of obligations, permissions, and actions; while in the second approach [45, 36] power cannot be reduced to static normative positions, aligning with Hohfeld’s original distinction (see [46]). In LRK, the second approach is followed, as the power to know is treated as a notion distinct from obligatory announcements.

Dynamic Epistemic Logic (DEL) is a family of modal logics of model change, and the type of public announcements used in LRK, removing links and not states, originates from [27] and can be represented as arrow updates [38]. LRK is also related to logics of questions [32, 60] and inquisitive semantics [22]. Other works that use deontic logic for epistemic actions include [4, 2, 43].

The work presented in this paper follows the LogiKEy methodology outlined in [15, 14]. This approach has focused on using SSEs in HOL [6] to enable knowledge engineers to explore and experiment with different object logics and their combinations, integrating both general and domain-specific knowledge. Other works based on LogiKEy are for example [17, 53], and work in [52, 51, 16, 18, 13] specifically focuses on SSEs of deontic logics. Work in [11] presents an SSE of Public Announcement Logic (PAL), where the embedded formulas depend only on evaluation domains, while additional parameters that are more complex to maintain and update are needed

for LRK. Another contribution of this paper has been to show that the LogiKEy approach scales well to support the embedding of a non-trivial logic like LRK. Of course, the automation of non-trivial object logics such as LRK with HOL argumentation tools will eventually reach its limits, while specialized LRK provers may ultimately still be responsive. However, it should be noted that developing powerful native LRK provers and model finders is a very ambitious task. And if such tools eventually exist, they can also be integrated into our LogiKEy framework as additional external tools.

The SSE approach has been specifically employed for the automation of quantified modal logics [10]. A variety of other methods are available for automated reasoning in modal logics. For propositional modal logic it is possible to use encodings to SAT [57] or to SMT solvers [1]. Authors have provided resolution calculi for modal logics [47, 49] and corresponding model construction algorithms [35] that are implemented in the resolution-based multimodal logic prover KSP [48]. Algorithms for unification in modal logics exist [3], as well as syntactic abstraction methods [24] to reason about first-order modal logics. Authors in [61] studied how DEL models can be faithfully represented as knowledge structures for symbolic model checking with an approach based on Binary Decision Diagram reasoning.

There are also a number of tools that can be used for legal reasoning and compliance checking, for instance implementations of Answer Set Programming (ASP) [40], its extension implemented in DLV [41], implementations of defeasible deontic logic [39, 23], or PROLEG (PROlog based LEGal reasoning support system) [56]. A survey of automated reasoners for compliance checking can be found in [55]. However, none of these can explicitly represent epistemic rights such as the right to know.

Compared to other approaches for automating a logic, the SSE approach offers several advantages: (1) it handles all conversions directly within HOL, making the encoding concise and readable; (2) it naturally extends from propositional to first-order and higher-order logic, enabling more expressive reasoning; (3) it supports deep and shallow embeddings [19] and mechanized faithfulness proofs between them; and (4) it benefits from state-of-the-art theorem provers and model finders integrated with *Isabelle/HOL* [20]. In general, this approach may have a performance loss compared to domain-specific solutions, but not always [30].

To the best of our knowledge, this work presents the first automation of a logic capable of reasoning about the right to know. The SSE approach enables us to represent and experiment with the right to know by effectively handling model updates, as shown in the discussed GDPR use case. Directions for future research include applying the embedded logic to represent other use cases involving epistemic rights, and extending the embedding to multi-agent scenarios, with multiple sender and receiver agents, possibly with overlapping roles. Another important direction that is often required in practical applications is to combine LRK with other logics, such as doxastic or temporal logics, for more expressive reasoning: the meta-logical nature of the LogiKEy framework supports such combinations.

## Acknowledgements

This work was supported by the Luxembourg National Research Fund (FNR) through the project Logical methods for Deontic Explanations (INTER/DFG/23/17415164/LODEX) and the project Deontic Logic for Epistemic Rights (OPEN O20/14776480). We also thank the anonymous reviewers for their fruitful feedback.

<sup>4</sup> As observed in [42], consider, e.g.,  $R_r p \rightarrow R_r \neg p$ , which is valid in LRK, but would be counterintuitive for operator  $R_r p$  interpreted as a general power-right.

## References

- [1] C. Areces, P. Fontaine, and S. Merz. Modal satisfiability via SMT solving. In R. Hennicker and R. de Nicola, editors, *Software, Services, and Systems*, volume 8950 of *LNCS*, pages 30–45. Springer, 2015.
- [2] G. Aucher, G. Boella, and L. van der Torre. A dynamic logic for privacy compliance. *Artificial Intelligence and Law*, 19(2):187, 2011.
- [3] F. Baader and S. Ghilardi. Unification in modal and description logics. *Logic Journal of the IGPL*, 19(6):705–730, 2010.
- [4] P. Balbiani and P. Seban. Reasoning about permitted announcements. *Journal of Philosophical Logic*, 40(4):445–472, 2011.
- [5] C. Benzmüller. Cut-elimination for quantified conditional logic. *Journal of Philosophical Logic*, 46(3):333–353, 2017.
- [6] C. Benzmüller. Universal (meta-)logical reasoning: Recent successes. *Science of Computer Programming*, 172:48–62, 2019.
- [7] C. Benzmüller. Universal (meta-)logical reasoning: The wise men puzzle (Isabelle/HOL Dataset). *Data in Brief*, 24(103823):1–5, 2019.
- [8] C. Benzmüller and P. Andrews. Church’s Type Theory. In E. N. Zalta and U. Nodelman, editors, *The Stanford Encyclopedia of Philosophy*. Stanford University, Spring 2024 edition, 2024.
- [9] C. Benzmüller and D. Miller. Automation of higher-order logic. In D. M. Gabbay et al., editors, *Handbook of the History of Logic. Volume 9 — Computational Logic*, pages 215–254. North Holland, 2014.
- [10] C. Benzmüller and L. C. Paulson. Quantified multimodal logics in simple type theory. *Logica Universalis*, 7(1):7–20, 2013.
- [11] C. Benzmüller and S. Reiche. Automating public announcement logic with relativized common knowledge as a fragment of HOL in LogiKEy. *Journal of Logic and Computation*, 33(6):1243–1269, 2023.
- [12] C. Benzmüller, C. Brown, and M. Kohlase. Higher-order semantics and extensionality. *Journal of Symbolic Logic*, 69(4):1027–1088, 2004.
- [13] C. Benzmüller, A. Farjami, and X. Parent. Åqvist’s dyadic deontic logic E in HOL. *FLAP*, 6(5):733–755, 2019.
- [14] C. Benzmüller, A. Farjami, D. Fuenmayor, et al. LogiKEy workbench: Deontic logics, logic combinations and expressive ethical and legal reasoning (Isabelle/HOL dataset). *Data in Brief*, 33(106409):1–10, 2020.
- [15] C. Benzmüller, X. Parent, and L. van der Torre. Designing normative theories for ethical and legal reasoning: LogiKEy framework, methodology, and tool support. *Artificial Intelligence*, 287:103348, 2020.
- [16] C. Benzmüller, A. Farjami, and X. Parent. Dyadic deontic logic in HOL: Faithful embedding and meta-theoretical experiments. In S. Rahman et al., editors, *New Developments in Legal Reasoning and Logic: From Ancient Law to Modern Legal Systems*, volume 23 of *Logic, Argumentation & Reasoning*. Springer Nature, 2022.
- [17] C. Benzmüller, D. Fuenmayor, and B. Lomfeld. Modelling value-oriented legal reasoning in LogiKEy. *Logics*, 2(1):31–78, 2024.
- [18] C. Benzmüller et al. I/O logic in HOL. *FLAP*, 6(5):715–732, 2019.
- [19] C. Benzmüller. Faithful logic embeddings in HOL — deep and shallow. In *CADE-30*, volume 15943 of *LNCS*, pages 280–302. Springer, 2025.
- [20] J. Blanchette, S. Böhme, and L. Paulson. Extending Sledgehammer with SMT solvers. *Journal of Automated Reasoning*, 51:116–130, 2011.
- [21] A. Church. A formulation of the simple theory of types. *Journal of Symbolic Logic*, 5(2):56–68, 1940.
- [22] I. A. Ciardelli and F. Roelofsen. Inquisitive dynamic epistemic logic. *Synthese*, 192:1643–1687, 2015.
- [23] M. Cristani et al. The architecture of a reasoning system for defeasible deontic logic. *Procedia Comput. Sci.*, 225:4214–4224, 2023.
- [24] D. Doligez, J. Kriener, L. Lammport, et al. Coalescing for reasoning in first-order modal logics. In C. Benzmüller et al., editors, *ARQNL 2014*, volume 33 of *EPiC Series in Computing*, pages 1–16. EasyChair, 2014.
- [25] H. Dong and O. Roy. Dynamic logic of legal competences. *J. Log. Lang. Inf.*, 30(4):701–724, 2021.
- [26] European Parliament and Council of the European Union. Regulation (EU) 2016/679 (GDPR), 2016.
- [27] J. Gerbrandy and W. Groeneveld. Reasoning about information change. *Journal of Logic, Language and Information*, 6:147–169, 1997.
- [28] S. Gerke, T. Minssen, and G. Cohen. Ethical and legal challenges of artificial intelligence-driven healthcare. In A. Bohr et al., editors, *Artificial Intelligence in Healthcare*, pages 295–336. Academic Press, 2020.
- [29] J. Gibbons and N. Wu. Folding domain-specific languages: deep and shallow embeddings (functional pearl). In J. Jeuring and M. M. T. Chakravarty, editors, *ICFP 2014*, pages 339–347. ACM, 2014.
- [30] T. Gleißner, A. Steen, and C. Benzmüller. Theorem provers for every normal modal logic. In T. Eiter et al., editors, *Proc. of LPAR-21*, volume 46 of *EPiC Series in Computing*, pages 14–30. EasyChair, 2017.
- [31] K. Gödel. Über formal unentscheidbare Sätze der Principia Mathematica und verwandter Systeme I. *Monatshefte für Mathematik und Physik*, 38(1):173–198, 1931.
- [32] J. Groenendijk and M. Stokhof. Chapter 19 - questions. In J. van Ben-  
them and A. ter Meulen, editors, *Handbook of Logic and Language*, pages 1055–1124. North-Holland, Amsterdam, 1997.
- [33] L. Henkin. Completeness in the theory of types. *The Journal of Symbolic Logic*, 15(2):81–91, 1950.
- [34] W. N. Hohfeld. Fundamental legal conceptions applied in judicial reasoning. In W. W. Cook, editor, *Fundamental Legal Conceptions Applied in Judicial Reasoning and Other Legal Essays*, pages 23–64. New Haven: Yale University Press, 1923.
- [35] U. Hustadt et al. Model construction for modal clauses. In *Automated Reasoning*, volume 14740 of *LNCS*, pages 3–23. Springer, 2024.
- [36] A. J. I. Jones and M. Sergot. A formal characterisation of institutionalised power. *Log. J. IGPL*, 4(3):427–443, 1996.
- [37] S. Kanger and H. Kanger. Rights and parliamentarism. *Theoria*, 32(2):85–115, 1966.
- [38] B. Kooi and B. Renne. Arrow update logic. *The Review of Symbolic Logic*, 4(4):536–559, 2011.
- [39] H.-P. Lam and G. Governatori. The making of SPINdle. In G. Governatori et al., editors, *Rule Representation, Interchange and Reasoning on the Web*, number 5858 of *LNCS*, pages 315–322. Springer, 2009.
- [40] N. Leone and F. Ricca. Answer set programming: A tour from the basics to advanced development tools and industrial applications. In *Reasoning Web*, volume 9203 of *LNCS*, pages 308–326. Springer, 2015.
- [41] N. Leone, G. Pfeifer, et al. The DLV system for knowledge representation and reasoning. *ACM Trans. Comput. Log.*, 7(3):499–562, 2006.
- [42] X. Li and R. Markovich. A dynamic logic of the right to know. *FLAP*, 12(2):221–250, 2025.
- [43] X. Li, D. Gabbay, and R. Markovich. Dynamic Deontic Logic for Permitted Announcements. In *Proc. of the 19th Intl. Conf. on Principles of Knowledge Representation and Reasoning*, pages 226–235, 2022.
- [44] L. Lindahl. *Position and change: A study in law and logic*. Synthese Library. Springer Dordrecht, 1 edition, 1977.
- [45] D. Makinson. On the formal representation of rights relations. *Journal of Philosophical Logic*, 15(4):403–425, 1986.
- [46] R. Markovich. Understanding Hohfeld and Formalizing Legal Rights: the Hohfeldian Conceptions and Their Conditional Consequences. *Studia Logica*, 108, 2020.
- [47] C. Nalon and C. Dixon. Clausal resolution for normal modal logics. *Journal of Algorithms*, 62(3):117–134, 2007.
- [48] C. Nalon, U. Hustadt, and C. Dixon. A resolution-based theorem prover for  $k_n, kn$ : Architecture, refinements, strategies and experiments. *Journal of Automated Reasoning*, 64, 2018.
- [49] C. Nalon, C. Dixon, and U. Hustadt. Modal resolution: Proofs, layers, and refinements. *ACM Trans. Comput. Log.*, 20(4), 2019.
- [50] T. Nipkow, L. C. Paulson, and M. Wenzel. *Isabelle/HOL: A Proof Assistant for Higher-Order Logic*. Springer, Berlin, Heidelberg, 2002.
- [51] X. Parent and C. Benzmüller. Automated verification of deontic correspondences in Isabelle/HOL – first results. In *Proc. of ARQNL 2022*, volume 3326, pages 92–108. CEUR-WS.org, 2023.
- [52] X. Parent and C. Benzmüller. Conditional normative reasoning as a fragment of HOL. *J. Appl. Non-Class. Log.*, 34:561–592, 2024.
- [53] L. Pasetto and C. Benzmüller. Implementing the Fatio protocol for multi-agent argumentation in LogiKEy. In *Proc. of ARQNL 2024*, volume CEUR-3875, pages 38–47. CEUR-WS.org, 2024.
- [54] F. Pfennig and C. Elliott. Higher-order abstract syntax. In R. L. Wexelblat, editor, *Proc. of the ACM SIGPLAN’88 Conf. on Progr. Language Design and Implementation (PLDI)*, pages 199–208. ACM, 1988.
- [55] L. Robaldo et al. Compliance checking on first-order knowledge with conflicting and compensatory norms: a comparison among currently available technologies. *Artif. Intell. Law*, 32(2):505–555, 2023.
- [56] K. Satoh, K. Asai, T. Kogawa, et al. Proleg: An implementation of the presupposed ultimate fact theory of Japanese civil code by PROLOG technology. In T. Onada et al., editors, *New Frontiers in Artificial Intelligence*, pages 153–164. Springer Berlin Heidelberg, 2011.
- [57] R. Sebastiani and M. Vescovi. Automated reasoning in modal and description logics via sat encoding: the case study of  $k(m)/alc$ -satisfiability. *J. Artif. Intell. Res.*, 35:343–389, 2009.
- [58] G. Sileno and M. Pascucci. Disentangling deontic positions and abilities: a modal analysis. In F. Calimeri et al., editors, *Proc. of the 35th Edition of the Italian Conf. on Comp. Logic (CILC)*, pages 36–50, 2020.
- [59] UNESCO. Recommendation on the ethics of ai. <https://unesdoc.unesco.org/ark:/48223/pf0000381137>, 2022. Accessed: 2025-05-06.
- [60] J. van Benthem and Ş. Minićă. Toward a dynamic logic of questions. *Journal of Philosophical Logic*, 41(4):633–669, 2012.
- [61] J. van Benthem et al. Symbolic model checking for dynamic epistemic logic—S5 and beyond. *J. Log. Comp.*, 28(2):367–402, 2018.
- [62] L. Watson. *The Right to Know: Epistemic Rights and Why We Need Them*. Routledge Focus on Philosophy. Routledge, 1 edition, 2021.