

## Modelling Value-oriented Legal Reasoning in LOGIKEY

Christoph Benz Müller · David Fuenmayor ·  
Bertram Lomfeld

Received: date / Accepted: date

**Abstract** The logico-pluralist LOGIKEY knowledge engineering methodology and framework is applied to the modelling of a theory of legal balancing in which legal knowledge (cases and laws) is encoded by utilising context-dependent value preferences. The theory obtained is then used to formalise, automatically evaluate, and reconstruct illustrative property law cases (involving appropriation of wild animals) within the *Isabelle/HOL* proof assistant system, illustrating how LOGIKEY can harness interactive and automated theorem proving technology to provide a testbed for the development and formal verification of legal domain-specific languages and theories. Modelling value-oriented legal reasoning in that framework, we establish novel bridges between latest research in knowledge representation and reasoning in non-classical logics, automated theorem proving, and applications in legal reasoning.

**Keywords** legal balancing · value-oriented reasoning · automated theorem proving · logical pluralism · proof assistants · Isabelle/HOL

---

C. Benz Müller  
University of Bamberg, AI Systems Engineering & Freie Universität Berlin, Dep. of Mathematics and Computer Science  
E-mail: c.benzmueller@fu-berlin.de

D. Fuenmayor (corresponding author)  
Université du Luxembourg, Dep. of Computer Science & Freie Universität Berlin, Dep. of Mathematics and Computer Science  
E-mail: david.fuenmayor@uni.lu

B. Lomfeld  
Freie Universität Berlin, Dep. of Law  
E-mail: bertram.lomfeld@fu-berlin.de

## 1 Introduction

Law today has to reflect highly pluralistic environments in which a plurality of values, world-views and logics coexist. One function of modern, reflexive law is to enable the social interaction within and between such worlds (Lomfeld 2017; Teubner 1983). Any sound model of legal reasoning needs to be pluralistic, supporting different value systems, value preferences, and maybe even different logical notions, while at the same time reflecting the uniting force of law.

Adopting such a perspective, in this paper we apply the logico-pluralistic LOGIKEY knowledge engineering methodology and framework (Benzmüller et al. 2020) to the modelling of a theory of value-based legal balancing, a *discursive grammar* of justification (Lomfeld 2019), which we then employ to formally reconstruct and automatically assess, using the *Isabelle/HOL* proof assistant system, some illustrative property law cases involving the appropriation of wild animals (termed “wild animal cases”; cf. Bench-Capon and Sartor (2003), Berman and Hafner (1993), and Merrill and H. E. Smith (2017, Ch. II. A.1) for background). Lomfeld’s *discursive grammar* is encoded, for our purposes, as a logic-based domain-specific language (DSL) in which the legal knowledge embodied in statutes and case corpora becomes represented as *context-dependent* preferences among (combinations of) values constituting a pluralistic value system or ontology. This knowledge can thus be complemented by further legal and world knowledge, e.g., from legal ontologies (Casanovas et al. 2016; Hoekstra et al. 2009).

The LOGIKEY framework supports plurality at different layers; cf. Fig. 1. Classical higher-order logic (HOL) is fixed as a *universal meta-logic* (Benzmüller 2019) at the base layer (L0), on top of which a plurality of (combinations of) object logics can become encoded (layer L1). Employing these logical notions we can now articulate a variety of logic-based domain-specific languages (DSLs), theories and ontologies at the next layer (L2), thus enabling the modelling and automated assessment of different application scenarios (layer L3). These linked layers, as featured in the LOGIKEY approach, facilitate fruitful interdisciplinary collaboration between

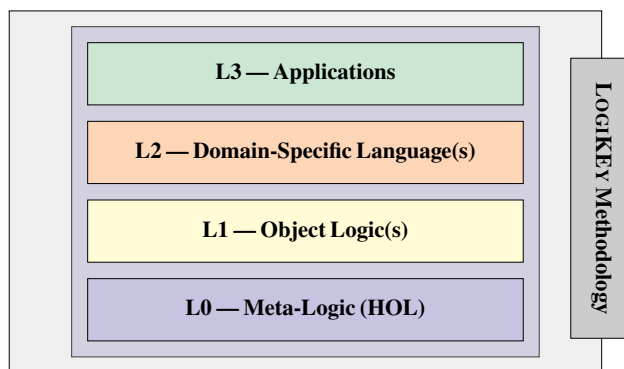


Fig. 1 LOGIKEY development methodology

specialists in different AI-related domains and domain experts in the design and development of knowledge-based systems.

LOGIKEY, in this sense, fosters a *division of labour* among different specialist roles. For example, ‘logic theorists’ can concentrate on investigating the semantics and proof calculi for different object logics, while ‘logic engineers’ (e.g., with a computer science background) can focus on the encoding of suitable combinations of those formalisms in the meta-logic HOL and on the development and/or integration of relevant automated reasoning technology. Knowledge engineers can then employ these object logics for knowledge representation (by developing ontologies, taxonomies, controlled languages, etc.), while domain experts (ethicists, lawyers, etc.) collaborate with requirements elicitation and analysis, as well as providing domain-specific counseling and feedback. These tasks can be supported in an integrated fashion by harnessing (and extending) modern mathematical proof assistant systems (aka. interactive theorem provers), which thus become a testbed for the development of logics and ethico-legal theories.

The work reported below is a LOGIKEY-supported collaborative research effort involving two computer scientists (Benzmüller & Fuenmayor) together with a lawyer and legal philosopher (Lomfeld), who have joined forces with the goal of studying the computer-encoding and automation of a theory of value-based legal balancing: Lomfeld’s *discursive grammar* (Lomfeld 2019). A formally-verifiable legal domain-specific language (DSL) has been developed for the encoding of this theory (at LOGIKEY’s layer L2). This DSL has been built on top of a suitably chosen object-logical language: a modal logic of preferences (at layer L1), by drawing upon the representation and reasoning infrastructure integrated within the proof assistant *Isabelle/HOL* (layer L0). The resulting system is then employed for the assessment of legal cases in property law (at layer L3), which includes the formal modelling of background legal and world knowledge, as required to enable the verification of predicted legal case outcomes and the automatic generation of value-oriented logical justifications (backings) for them.

From a wider perspective, LOGIKEY aims at supporting the practical development of computational tools for legal and normative reasoning based on formal methods. One of the main drives for its development has been the introduction of automated reasoning techniques for the design, verification (offline & online), and control of intelligent autonomous systems, as a step towards *explicit ethical agents* (Moor 2009; Scheutz 2017). The argument here is that ethico-legal control mechanisms (such as ethical governors; cf. Arkin et al. (2009)) of intelligent autonomous systems should be understood and designed as knowledge-based systems, where the required ethical and legal knowledge becomes *explicitly* represented as a logical theory, i.e., as a set of formulas (axioms, definitions and theorems) encoded in a logic. We have set a special focus on the (re-)use of modern proof assistants based on HOL (*Isabelle/HOL*, *HOL-Light*, *HOLA*, etc.) and integrated automated reasoning tools (*theorem provers* and *model generators*) for the interactive development and verification of ethico-legal theories. To carry out the technical work reported in this paper, we have chosen to work with *Isabelle/HOL*, but the essence of our contributions can easily be applied to other proof assistants and automated reasoning systems for HOL.

Technical results concerning in particular our *Isabelle/HOL* encoding have been presented at the *International Conference on Interactive Theorem Proving* (ITP 2021) (Benzmüller and Fuenmayor 2021), and earlier ideas have been discussed at the *Workshop on Models of Legal Reasoning* (MLR 2020). In the present paper, we elaborate on these results and provide a more self-contained exposition, by giving further background information on Lomfeld’s *discursive grammar*, on the meta-logic HOL, and on the modal logic of preferences by van Benthem et al. (2009). More fundamentally, this paper presents the full picture, as framed by the underlying LOGIKEY framework, and highlights methodological insights, applications, and perspectives relevant to the *AI & Law* community. One of our main motivations is to help build bridges between recent research in knowledge representation and reasoning in non-classical logics, automated theorem proving, and applications in normative and legal reasoning.

*Paper structure:* After summarising Lomfeld’s theory of value-based legal balancing in §2, we briefly depict the LOGIKEY development and knowledge engineering methodology in §3, and our meta-logic HOL in §4. We then outline our object logic of choice – a (quantified) modal logic of preferences – in §5, where we also present its encoding in the meta-logic HOL and formally verify the preservation of meta-theoretical properties using the *Isabelle/HOL* proof assistant. Subsequently, we model in §6 Lomfeld’s legal theory and provide a custom-built DSL, which is again formally assessed using *Isabelle/HOL*. As an illustrative application of our framework, we present in §7 the formal reconstruction and assessment of well-known example legal cases in property law (“wild animal cases”), together with some considerations regarding the encoding of required legal and world knowledge. Related and further work is addressed in §8, and §9 concludes the article.

The *Isabelle/HOL* sources of our formalisation work are available at <http://logikey.org> under *Preference-Logics/EncodingLegalBalancing*. They are also explained in some detail in the Appendix A.

## 2 A Theory of Legal Values: *Discursive Grammar* of Justification

The case study with which we illustrate the LOGIKEY methodology in the present paper consists in the formal encoding and assessment on the computer of a theory of value-based legal balancing, as put forward by Lomfeld 2019. Lomfeld himself has played the role of the domain expert in our joint research, which from a methodological perspective, can be characterised as being both in part theoretical and in part empirical. Lomfeld’s primary role has been to provide background legal domain knowledge and to assess the adequacy of our formalisation results, while informing us of relevant conceptual and legal distinctions that needed to be made. In a sense, this created a win-win situation in which both Lomfeld’s theory and LOGIKEY’s methodology have been put to the test. This section presents Lomfeld’s theory and discusses some of its merits in comparison to related approaches.

Logical reconstructions quite often separate deductive rule application and inductive case-contextual interpretation as completely distinct ways of legal reasoning (cf. the overview in Prakken and Sartor (2015)). Understanding the whole process of legal

reasoning as an exchange of opposing action-guiding arguments, i.e., practical argumentation (Alexy 1978; Feteris 2017), a strict separation between logically distinct ways of legal reasoning breaks down. Yet, a variety of modes of rule-based (Hage 1997; Modgil and Prakken 2018; Prakken 1997), case-based (Aleven 1997; Ashley 1990; Horty 2011) and value-based (Bench-Capon et al. 2005; Berman and Hafner 1993; Grabmair 2016) reasoning coexist in legal theory and (court) practice.

In line with current computational models combining these different modes of reasoning (e.g., Bench-Capon and Sartor 2003; Maranhão and Sartor 2019), we argue that a discourse theory of law can consistently integrate them in the form of a multi-level system of legal reasoning. Legal rules or case precedents can thus be translated into (or analysed as) the balancing of plural and opposing (socio-legal) values on a deeper level of reasoning (Lomfeld 2015).

There exist indeed some models for quantifying legal balancing, i.e., for weighing competing reasons in a case (e.g., Alexy 2003; Sartor 2010). We share the opinion that these approaches need to get “integrated with logic and argumentation to provide a comprehensive account of value-oriented reasoning” (Sartor 2018). Hence a suitable model of legal balancing would need to reconstruct rule subsumption and case distinction as argumentation processes involving conflicting values.

Here, the functional differentiation of legal norms into *rules* and *principles* reveals its potential (Alexy 2000; Dworkin 1978). Whereas legal rules have a binary all-or-nothing validity driving out conflicting rules, legal principles allow for a scalable dimension of weight. Thus, principles could outweigh each other without rebutting the normative validity of colliding principles. Legal principles can be understood as a set of plural and conflicting values on a deep level of socio-legal balancing, which is structured by legal rules on an explicit and more concrete level of legal reasoning (Lomfeld 2015). The two-faceted argumentative relation is partly mirrored in the functional differentiation between *goal-norms* and *action-norms* (Sartor 2010) but should not be mixed up with a general understanding of principles as abstract rules (Raz 1972; Verheij et al. 1998) or as specific constitutional law elements (Barak 2012; Neves 2021).

In any event, if preferences between defeasible rules are reconstructed and justified in terms of preferences between underlying values, some questions about values necessarily pop up. In the words of Bench-Capon and Sartor (2003): “Are values scalar? [...] Can values be ordered at all? [...] How can sets of values be compared? [...] Can values conflict so the promotion of their combination is worse than promoting either separately? Can several less important values together overcome a more important value?”.

Hence an encompassing approach for legal reasoning as practical argumentation needs not only a formal reconstruction of the relation between legal values (or principles) and legal rules, but also a substantial framework of values (a basic value system) that allows to systematise value comparison and conflicts as a *discursive grammar* (Lomfeld 2015, 2019) of argumentation. In this article we define a value system not as a “preference order on sets of values” (Weide et al. 2010) but as a pluralistic set of values which allow for different preference orders. The computational conceptualisation (as a formal logical theory) of such a set of representational primitives for a pluralist basic value system can then be considered as a value “ontology” (Gruber 1993, 2009;

B. Smith 2003), which of course needs to be complemented by further ontologies for relevant background legal and world knowledge (see e.g. Casanovas et al. (2016) and Hoekstra et al. (2009)).

Combining the discourse-theoretical idea that legal reasoning is practical argumentation with a two-faceted model of legal norms, legal *rules* could be logically reconstructed as conditional preference relations between conflicting underlying *value principles* (Alexy 2000; Lomfeld 2015). The legal consequence of a rule  $R$  thus implies the preference of value principle  $A$  over value principle  $B$ , noted  $A > B$  (e.g. health security outweighs freedom to move).<sup>1</sup> This value preference applies under the condition that the rule's prerequisites  $E_1$  and  $E_2$  hold. Thus, if the propositions  $E_1$  and  $E_2$  are true in a given situation (e.g. a virus pandemic occurs and voluntary shut down fails), then the value preference  $A > B$  obtains. This value preference can be said to weight or *balance* the two values  $A$  and  $B$  against each other. We can thus translate this concrete legal rule as a *conditional preference relation* between colliding value principles:

$$R : (E_1 \wedge E_2) \Rightarrow A < B$$

More generally,  $A$  and  $B$  could also be structured as aggregates of value principles, whereas the condition of the rule can consist in a conjunction of arbitrary propositions. Moreover, it may well happen that, given some conditions, several rules become relevant in a concrete legal case. In such cases the rules determine a structure of legal balancing between conflicting plural value principles. Moreover, making explicit the underlying *balancing* of values against each other (as value preferences) helps to justify a legal consequence (e.g. sanctioned lock-down) or ruling in favour of a party (e.g. defendant) in a legal case.

But which value principles are to be balanced? How to find a suitable justification framework? Based on comparative discourse analyses in different legal systems, one can reconstruct a general dialectical (antagonistic) taxonomy of legal value principles used in (at least Western) legislation, legislative materials, cases, textbooks and scholar writings (Lomfeld 2015). The idea is to provide a plural and yet consistent system of basic legal values and principles, independent of concrete cases or legal fields, to justify legal decisions.

The proposed legal value system (Lomfeld 2019), see Fig. 2, is consistent with many existing taxonomies of antagonistic psychological (Rokeach 1973; Schwartz 1992), political (Eysenck 1954; Mitchell 2007) and economic values (Clark 1991).<sup>2</sup> In all social orders one can observe a general antinomy between individual and collective values. Ideal types of this fundamental dialectic are: the basic value of FREEDOM for the individual, and the basic value of SECURITY for the collective perspective. Another classic social value antinomy is between a functional-economic (utilitarian) and a more idealistic (egalitarian) viewpoint, represented in the ethical debate by the essential dialectic concerning the basic values of UTILITY versus EQUALITY. These

<sup>1</sup> In §6 these values will be assigned to particular parties/actors, so that ruling in favour of different parties may promote different values.

<sup>2</sup> All these taxonomies are pluralist frameworks that do encompass differences in global value patterns and cultural value evolution (Hofstede 2001; Inglehart 2018). For an approach oriented at Maslow's hierarchy of needs (Bench-Capon 2020a).

four normative poles stretch an axis of value coordinates for the general value system construction. We thus speak of a normative dialectics, since each of the antagonistic basic values and related principles can (and in most situations will) collide with each other.

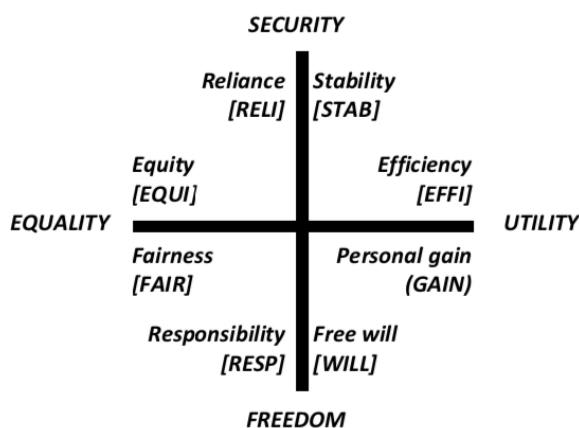


Fig. 2 Basic legal value system (ontology) by Lomfeld (2019)

Within this dialectical matrix eight more concrete legal *value principles* are identified. FREEDOM represents the normative basic value of individual autonomy and comprises the legal (value) principles of –more functional– individual choice or ‘free will’ (WILL) and –more idealistic– (self-)‘responsibility’ (RESP). The basic value of SECURITY addresses the collective dimension of public order and comprises the legal principles of –more functional– collective ‘stability’ (STAB) of a social system and –more idealistic– social trust or ‘reliance’ (RELI). The value of UTILITY means economic welfare on the personal and collective level and comprises the legal principles of collective overall welfare-maximisation, i.e., ‘efficiency’ (EFFI) and individual welfare-maximisation, i.e., economic benefit or ‘gain’ (GAIN). Finally, EQUALITY represents the normative ideal of equal treatment and equal allocation of resources and comprises the legal principles of –more individual– equal opportunity or procedural ‘fairness’ (FAIR) and –more collective– distributional equality or ‘equity’ (EQUI).

This legal value system (or ontology) can consistently cover existing value sets from AI & Law accounts of value-oriented reasoning (e.g., Bench-Capon 2012; Berman and Hafner 1993; T. Gordon and Walton 2012; Sartor 2010), mostly exemplified by modelling famous common law property cases, in particular, “wild animal cases”. A key feature of Lomfeld’s *discursive grammar* of dialectical values consists in its purely qualitative modelling of legal balancing in terms of context-dependent logic-based preferences among values, without any need for determining quantitative weights.

### 3 The LOGIKEY Methodology

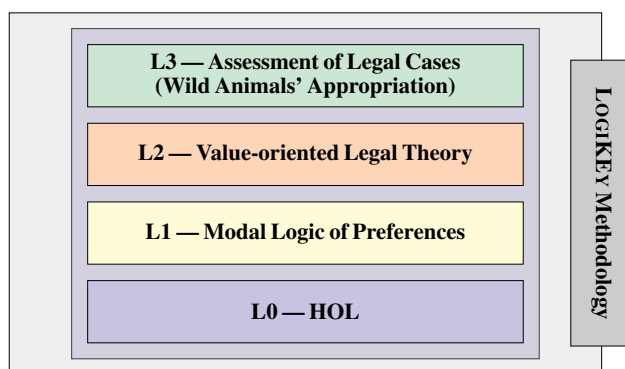
LOGIKEY, as a methodology (Benzmüller et al. 2020), refers to the principles underlying the organisation and the conduct of complex knowledge design and engineering processes, with a particular focus on the legal and ethical domain. Knowledge engineering refers to all the technical and scientific aspects involved in building, maintaining and using knowledge-based systems employing logical formalisms as a representation language. In particular, we speak of *logic engineering* to highlight those tasks directly related to the syntactic and semantic definition, as well as to the meta-logical encoding and automation, of different combinations of object logics. It is also LOGIKEY’s objective to fruitfully integrate contributions from different research communities (such as interactive and automated theorem proving, non-classical logics, knowledge representation, and domain specialists) and to make them accessible at a suitable level of abstraction and technicality to practitioners in diverse fields.

A fundamental characteristic of the LOGIKEY methodology consists in the utilisation of classical higher-order logic (HOL, cf. Benzmüller and P. Andrews (2019)) as a general-purpose logical formalism in which to encode different (combinations of) object logics. This enabling technique is known as shallow<sup>3</sup> semantical embeddings (SSEs). HOL thus acts as the substrate in which a plurality of logical languages, organised hierarchically at different abstraction layers, become ultimately encoded and reasoned with. This in turn enables the provision of powerful tool support: we can harness mathematical proof assistants (e.g. *Isabelle/HOL*) as a testbed for the development of logics, and ethico-legal DSLs and theories. More concretely, off-the-shelf theorem provers and (counter-)model generators for HOL, as provided, e.g., in the interactive proof assistant *Isabelle/HOL* (Blanchette et al. 2016), are assisting the LOGIKEY knowledge & logic engineers (as well as domain experts) to *flexibly experiment* with underlying (object) logics and their combinations, with general and domain knowledge, and with concrete use cases—all at the same time. Thus, continuous improvements of these off-the-shelf provers, without further ado, leverage the reasoning performance in LOGIKEY.

The LOGIKEY methodology, cf. Fig. 1, has been instantiated in this article to support and guide the simultaneous development of the different modelling layers as depicted in Fig. 3, and which will be the subject of discussion in the following sections. According to the logico-pluralistic nature of LOGIKEY, only the lowest layer (L0), meta-logic HOL (cf. §4), remains fixed, while all other layers are subject to dynamic adjustments until a satisfying overall solution in the overall modelling process is reached. At the next layer (L1) we are faced with the choice of an object logic, in our case a modal logic of preference (cf. §5). A legal DSL (cf. §6), created after Lomfeld’s *discursive grammar* (cf. §2), further extends this object logic at a higher

<sup>3</sup> Shallow semantical embeddings are different from *deep embeddings* of an object logic. In the latter case the syntax of the object logic is represented using an inductive data structure (e.g., following the definition of the language). The semantics of a formula is then evaluated by recursively traversing the data structure, and additionally a proof theory for the logic maybe be encoded. Deep embeddings typically require technical inductive proofs, which hinder proof automation, that can be avoided when shallow semantical embeddings are used instead. For more information on shallow and deep embeddings we refer to the literature (Gibbons and Wu 2014; Svenningsson and Axelsson 2013).





**Fig. 3** LOGIKEY development methodology as instantiated in the given context

level of abstraction (layer L2). At the upper layer (layer L3), we use this legal DSL to encode and automatically assess some example legal cases (“wild animal cases”) in property law (cf. §7) by relying upon previously encoded background legal and world knowledge.<sup>4</sup> The higher layers thus make use of the concepts introduced at the lower layers. Moreover, at each layer, the encoding efforts are guided by selected tests and ‘sanity checks’ in order to formally verify relevant properties of the encodings at and up to that level.

It is worth noting that the application of our approach to deciding concrete legal cases reflects ideas in the AI & Law literature about understanding the solution of legal cases as theory construction, i.e., “building, evaluating and using theories” (Bench-Capon and Sartor 2003).<sup>5</sup> This multi-layered, iterative knowledge engineering process is supported in our LOGIKEY framework by adapting interactive and automated reasoning technology for HOL (as a meta-logic).

An important aspect thereby is that the LOGIKEY methodology foresees and enables the knowledge engineer to flexibly switch between the modelling layers and to suitably adapt the encodings also at lower layers if needed. The engineering process thus has backtracking points and several work cycles may be required; thereby the higher layers may also pose modification requests to the lower layers. Such demands, unlike in most other approaches, may also involve far-reaching modifications of the chosen logical foundations, e.g., in the particularly chosen modal preference logic.

The work we present in this article is in fact the result of an iterative, give-and-take process encompassing several cycles of modelling, assessment and testing activities, whereby a (modular) logical theory gradually evolves until eventually reaching a state of highest coherence and acceptability. In line with previous work on *computational hermeneutics* (Fuenmayor and Benzmüller 2019), we may then speak of arriving at a

<sup>4</sup> In some cases it can be convenient to split one or more layers into sublayers. For instance, in our case study (cf. §7), layer L3 has been further subdivided to allow for a more strict separation between general legal & world knowledge (legal concepts and norms), cf. §7.1, from its *application* to relevant facts in the process of deciding a case (factual/contextual knowledge), cf. §7.2.

<sup>5</sup> The authors judiciously quote McCarty (1995): “The task for a lawyer or a judge in a ‘hard case’ is to construct a theory of the disputed rules that produces the desired legal result, and then to persuade the relevant audience that this theory is preferable to any theories offered by an opponent.”

state of *reflective equilibrium* (Daniels 2020), as the end-point of an iterative process of mutual adjustment among (general) principles and (particular) judgements, where the latter are intended to become logically entailed by the former. It is also worth noting that the notion of *reflective equilibrium* has been introduced by the philosopher John Rawls in moral and political philosophy as a method for the development of his *theory of justice* (Rawls 1971), an analogous endeavour to ours in the present work. In fact, an earlier formulation of this approach is often credited to the philosopher Nelson Goodman, who proposed it as a method for the development of (inference rules for) deductive and inductive logical systems (Goodman 1955), again, very much in the spirit of LOGIKEY.

#### 4 Meta-logic (L0) – Classical Higher-Order Logic

To keep this article sufficiently self-contained we briefly introduce a classical higher-order logic, termed HOL; more detailed information on HOL and its automation can be found in the literature (P. B. Andrews 1972a,b; Benzmüller and P. Andrews 2019; Benzmüller et al. 2004; Benzmüller and Miller 2014).

The notion of HOL used in this article refers to a simply typed logic of functions that has been put forward by Church (1940). Hence all terms of HOL get assigned a fixed and unique type, commonly written as a subscript (i.e., the term  $t_\alpha$  has  $\alpha$  as its type). HOL provides  $\lambda$ -notation, as an elegant and useful means to denote unnamed functions, predicates and sets;  $\lambda$ -notation also supports compositionality, a feature we heavily exploit to obtain elegant, non-recursive encoding definitions for our logic embeddings in the remainder. Types in HOL eliminate paradoxes and inconsistencies.

HOL comes with a set  $T$  of *simple types*, which is freely generated from a set of *basic types*  $BT \supseteq \{o, \iota\}$  using the function type constructor  $\rightarrow$  (written as a right-associative infix operator). For instance,  $o$ ,  $\iota \rightarrow o$  and  $\iota \rightarrow \iota \rightarrow \iota$  are types. The type  $o$  denotes a two-element set of truth-values and  $\iota$  denotes a non-empty set of individuals.<sup>6</sup> Further base types may be added as the need arises.

The *terms* of HOL are inductively defined starting from typed constant symbols ( $C_\alpha$ ) and typed variable symbols ( $x_\alpha$ ) using  $\lambda$ -abstraction ( $(\lambda x_\alpha. s_\beta)_{\alpha \rightarrow \beta}$ ) and *function application* ( $(s_{\alpha \rightarrow \beta} t_\alpha)_\beta$ ), thereby obeying type constraints as indicated. Type subscripts and parentheses are usually omitted to improve readability, if obvious from the context or irrelevant. Observe that  $\lambda$ -abstractions introduce unnamed functions. For example, the function that adds 2 to a given argument  $x$  can be encoded as  $(\lambda x. x + 2)$ , and the function that adds two numbers can be encoded as  $(\lambda x. (\lambda y. x + y))$ .<sup>7</sup> HOL terms of type  $o$  are also called formulas.<sup>8</sup>

<sup>6</sup> In this article, we will actually associate type  $\iota$  later on (cf. §5.2) with the domain of possible states/worlds.

<sup>7</sup> Note that functions of more than one argument can be represented in HOL in terms of functions of one argument. In this case the values of these one-argument function applications are themselves functions, which are subsequently applied to the next argument. This technique, introduced by Schönfinkel (1924), is commonly called *currying*; cf. Benzmüller and P. Andrews (2019).

<sup>8</sup> HOL formulas (layer L0) should not be confused with the object-logical formulas (layer L1); the latter will later be identified in §5.2 with HOL predicates of type  $\iota \rightarrow o$ .

Moreover, to obtain a proper logic, we add  $\neg_{o \rightarrow o}$ ,  $\vee_{o \rightarrow o \rightarrow o}$  and  $\Pi_{(\alpha \rightarrow o) \rightarrow o}$  (for each type  $\alpha$ ) as predefined typed constant symbols to our language and call them *primitive logical connectives*. *Binder notation* for quantifiers  $\forall x_\alpha s_o$  is used as an abbreviation for  $\Pi_{(\alpha \rightarrow o) \rightarrow o} \lambda x_\alpha. s_o$ .

The *primitive logical connectives* are given a fixed interpretation as usual, and from them other logical connectives can be introduced as abbreviations. Material implication  $s_o \rightarrow t_o$  and existential quantification  $\exists x_\alpha s_o$ , for example, may be introduced as shortcuts for  $\neg s_o \vee t_o$  and  $\neg \forall x_\alpha \neg s_o$ , respectively. Additionally, *description or choice operators* or *primitive equality*  $=_{\alpha \rightarrow \alpha \rightarrow o}$  (for each type  $\alpha$ ), abbreviated as  $=^\alpha$ , may be added. Equality can also be defined by exploiting Leibniz' principle, expressing that two objects are equal if they share the same properties.

It is well known that, as a consequence of Gödel's Incompleteness Theorems, HOL with standard semantics is necessarily incomplete. In contrast, theorem proving in HOL is usually considered with respect to so-called general semantics (or Henkin semantics) in which a meaningful notion of completeness can be achieved (P. B. Andrews 1972a; Henkin 1950). Note that standard models are subsumed by Henkin general models such that valid HOL-formulas with respect to general semantics are also valid in the standard sense.

For the purposes of the present article, we shall omit the formal presentation of HOL semantics and of its proof system(s). We fix instead some useful notation for use in the remainder. We write  $\mathcal{H} \models^{\text{HOL}} \varphi$  if formula  $\varphi$  of HOL is *true* in a Henkin general model  $\mathcal{H}$ ;  $\models^{\text{HOL}} \varphi$  denotes that  $\varphi$  is (Henkin) *valid*, i.e., that  $\mathcal{H} \models^{\text{HOL}} \varphi$  for all Henkin models  $\mathcal{H}$ .

## 5 Object Logic (L1) – A Modal Logic of Preferences

Adopting the LOGIKEY methodology of §3 to support the given formalisation challenge, the first question to be addressed is: how to (initially) select the object logic at layer L1? The chosen logic not only must be expressive enough to allow the encoding of knowledge about the law (and the world), as required for the application domain (cf. our case study in §7), but must also provide the means to represent the kind of conditional value preferences featured in Lomfeld's theory (as described in §2). Importantly, it must also enable the adequate modelling of the notions of value aggregation and conflict, as featured in our legal DSL (discussed in §6).

Our initial choice has been the family of modal logics of preference presented by van Benthem et al. (2009), which we abbreviate by  $\mathcal{PL}$  in the remainder.  $\mathcal{PL}$  has been put forward as a modal logic framework for the formalisation of preferences which also allows for the modelling of *ceteris paribus* clauses in the sense of “all other things being equal”. This reading goes back to the seminal work of von Wright in the early 1960's (von Wright 1963).<sup>9</sup>

$\mathcal{PL}$  appears well suited for effective automation using the SSEs approach, which has been an important selection criterion. This judgment is based on good prior ex-

<sup>9</sup> For the purposes of the application scenarios studied later on §7, we have focused on  $\mathcal{PL}$ 's basic modal preference language, not yet employing *ceteris paribus* clauses. Nevertheless, we have provided a complete encoding and assessment of full  $\mathcal{PL}$  in the associated *Isabelle/HOL* sources.

perience with SSEs of related (monadic) modal logic frameworks (Benzmüller and Paulson 2010, 2013), whose semantics employs accessibility relations between possible worlds/states, just as  $\mathcal{PL}$  does. We note, however, that this choice of (a family of) object logics ( $\mathcal{PL}$ ) is just one out of a variety of logical systems which can be encoded as fragments of HOL employing the *shallow semantical embedding* approach; cf. Benzmüller (2019). This approach also allows us ‘upgrade’ our object logic whenever necessary. In fact, we add quantifiers and conditionals to  $\mathcal{PL}$  in §5.4. Moreover, we may consider combining  $\mathcal{PL}$  with other logics, e.g., with normal modal logics by the mechanisms of *fusion* and *product* (Carnielli and Coniglio 2020), or, more generally, by *algebraic fibring* (Carnielli et al. 2008, Ch. 2–3). This illustrates a central objective of the LOGIKEY approach, namely that the precise choice of a formalisation logic (i.e., the *object logic* at L1) is to be seen as a parameter.

In the subsections below we start by informally outlining the family of modal logics of preferences  $\mathcal{PL}$  (hence postponing their formal definition to an appendix §A.1). We then discuss its embedding as a fragment of HOL using the SSE approach. As for §4, the technically and formally less interested reader may actually skip the content of these subsections and get back later.

### 5.1 The modal logic of preferences $\mathcal{PL}$

We sketch the syntax and semantics of  $\mathcal{PL}$  adapting the description from van Benthem et al. (2009) (we refer to the appendix §A.1 for more details).

The formulas of  $\mathcal{PL}$  are inductively defined as follows (where  $p$  ranges over a set Prop of propositional constant symbols):

$$\varphi, \psi ::= p \mid \varphi \wedge \psi \mid \neg\varphi \mid \Diamond^{\leq}\varphi \mid \Diamond^{<}\varphi \mid \mathbf{E}\varphi$$

As usual in modal logic, van Benthem et al. (2009) give  $\mathcal{PL}$  a Kripke-style semantics, which models propositions as sets of states or ‘worlds’.  $\mathcal{PL}$  semantics employs a reflexive and transitive accessibility relation  $\leq$  (resp., its strict counterpart  $<$ ) to define the modal operators in the usual way. This relation is called a *betterness ordering* (between states or ‘worlds’).

For the sake of self-containedness, we summarize below the semantics of  $\mathcal{PL}$ .

A preference model  $\mathcal{M}$  is a triple  $\mathcal{M} = \langle \mathcal{W}, \leq, \delta \rangle$  where: (i)  $\mathcal{W}$  is a set of worlds/states; (ii)  $\leq$  is a *betterness relation* (reflexive and transitive) on  $\mathcal{W}$ , where its strict subrelation  $<$  is defined as:  $w < v := w \leq v \wedge v \not\leq w$  for all  $v, w \in \mathcal{W}$  (totality of  $\leq$ , i.e.,  $v \leq w \vee w \leq v$ , is generally not assumed); (iii)  $\delta$  is a standard modal valuation. Below we show the truth conditions for  $\mathcal{PL}$ ’s modal connectives (the rest being standard):

$$\begin{aligned} \mathcal{M}, w \models \Diamond^{\leq}\varphi & \text{ iff } \exists v \in \mathcal{W} \text{ such that } w \leq v \text{ and } \mathcal{M}, v \models \varphi \\ \mathcal{M}, w \models \Diamond^{<}\varphi & \text{ iff } \exists v \in \mathcal{W} \text{ such that } w < v \text{ and } \mathcal{M}, v \models \varphi \\ \mathcal{M}, w \models \mathbf{E}\varphi & \text{ iff } \exists v \in \mathcal{W} \text{ such that } \mathcal{M}, v \models \varphi \end{aligned}$$

A formula  $\varphi$  is *true* at world  $w \in \mathcal{W}$  in model  $\mathcal{M}$  if  $\mathcal{M}, w \models \varphi$ .  $\varphi$  is *globally true* in  $\mathcal{M}$ , denoted  $\mathcal{M} \models \varphi$ , if  $\varphi$  is *true* at every  $w \in \mathcal{W}$ . Moreover,  $\varphi$  is *valid* (in a class of models  $\mathcal{K}$ ) if *globally true* in every  $\mathcal{M} (\in \mathcal{K})$ , denoted  $\models_{\mathcal{PL}} \varphi$  ( $\models_{\mathcal{K}} \varphi$ ).

Thus,  $\diamond^{\leq}\varphi$  (resp.,  $\diamond^{<}\varphi$ ) can informally be read as “ $\varphi$  is true in a state that is considered to be at least as good as (resp., strictly better than) the current state” and  $\mathbf{E}\varphi$  can be read as “there is a state where  $\varphi$  is true”.

Further, standard connectives such as  $\vee$ ,  $\rightarrow$  and  $\leftrightarrow$  can also be defined in the usual way. The dual operators  $\Box^{\leq}\varphi$  (resp.,  $\Box^{<}\varphi$ ) and  $\mathbf{A}\varphi$  can also be defined as  $\neg\diamond^{\leq}\neg\varphi$  (resp.,  $\neg\diamond^{<}\neg\varphi$ ) and  $\neg\mathbf{E}\neg\varphi$ .

Readers acquainted with Kripke semantics for modal logic will notice that  $\mathcal{PL}$  features normal *S4* and *K4* diamonds operators  $\diamond^{\leq}$  and  $\diamond^{<}$ , together with a global existential modality  $\mathbf{E}$ . We can thus give the usual reading to  $\Box$  and  $\diamond$  as *necessity* and *possibility*, respectively.

Moreover, note that, since the *strict* betterness relation  $<$  is not reflexive, it does not hold in general that  $\Box^{<}\varphi \rightarrow \varphi$  (modal axiom *T*). Hence we can also give a ‘deontic reading’ to  $\diamond^{<}\varphi$  and  $\Box^{<}\varphi$ ; the former could then read as “ $\varphi$  is admissible/permisible” and the latter as “ $\varphi$  is recommended/obligatory”. This deontic interpretation can be further strengthened so that the latter entails the former by extending  $\mathcal{PL}$  with the postulate  $\Box^{<}\varphi \rightarrow \diamond^{<}\varphi$  (modal axiom *D*), or alternatively, by postulating the corresponding (meta-logical) semantic condition, namely, that for each state there exists a strictly better one (*seriality* for  $<$ ).

Observe that we use **boldface** fonts to distinguish standard logical connectives of  $\mathcal{PL}$  from their counterparts in HOL.

Similarly, eight different binary connectives for modelling preference statements between propositions can be defined in  $\mathcal{PL}$ :

$$\leq_{EE}/<_{EE}, \leq_{EA}/<_{EA}, \leq_{AE}/<_{AE}, \leq_{AA}/<_{AA}.$$

These connectives arise from four different ways of ‘lifting’ the *betterness ordering*  $\leq$  (resp.,  $<$ ) on states to a *preference ordering* on sets of states or propositions.

$$\begin{aligned} (\varphi \leq_{EE}/<_{EE} \psi) u &\text{ iff } \exists t \varphi s \wedge (\exists t \psi t \wedge s \leq/< t) \\ (\varphi \leq_{EA}/<_{EA} \psi) u &\text{ iff } \exists t \psi t \wedge (\forall s \varphi s \rightarrow s \leq/< t) \\ (\varphi \leq_{AE}/<_{AE} \psi) u &\text{ iff } \forall s \varphi s \rightarrow (\exists t \psi t \wedge s \leq/< t) \\ (\varphi \leq_{AA}/<_{AA} \psi) u &\text{ iff } \forall s \varphi s \rightarrow (\forall t \psi t \rightarrow s \leq/< t) \end{aligned}$$

Thus, different choices for a *logic of preference* are possible if we restrict ourselves to employing only a selected preference connective, where each choice provides the logic with particular characteristics, so that we can interpret preference statements between propositions (i.e., sets of states) in a variety of ways. As an illustration, according to the semantic interpretation provided by van Benthem et al. (2009), we can read  $\varphi <_{AA} \psi$  as “every  $\psi$ -state being better than every  $\varphi$ -state”, and read  $\varphi <_{AE} \psi$  as “every  $\varphi$ -state having a better  $\psi$ -state” (and analogously for others).

In fact, the family of preference logics  $\mathcal{PL}$  can be seen as encompassing, in substance, the proposals by von Wright (1963) (variant  $<_{AA}$ ) and Halpern (1997) (variants  $\leq_{AE}/<_{AE}$ ).<sup>10</sup> As we will see later in §6, there are only four choices ( $\leq_{EA}/<_{EA}$  and  $\leq_{AE}/<_{AE}$ ) of modal preference relations that satisfy the minimal conditions we

<sup>10</sup> Von Wright’s proposal is discussed in some detail in van Benthem et al. (2009); cf. also Liu (2008) for a discussion of further proposals.

impose for a logic of value aggregation. Moreover, they are the only ones which satisfy transitivity, a quite controversial property in the literature on preferences.

Last but not least, van Benthem et al. (2009) have provided ‘syntactic’ counterparts for these binary preference connectives as derived operators in the language of  $\mathcal{PL}$  (i.e., defined by employing the modal operators  $\Diamond^{\leq}\varphi$  (resp.,  $\Diamond^{<}\varphi$ ). We note these ‘syntactic variants’ in **boldface** font:

$$\begin{aligned}
(\varphi \leq_{EE} \psi) &:= E(\varphi \wedge \Diamond^{\leq}\psi) & \text{and} & \quad (\varphi <_{EE} \psi) := E(\varphi \wedge \Diamond^{<}\psi) \\
(\varphi \leq_{EA} \psi) &:= E(\psi \wedge \Box^{<}\neg\varphi) & \text{and} & \quad (\varphi <_{EA} \psi) := E(\psi \wedge \Box^{\leq}\neg\varphi) \\
(\varphi \leq_{AE} \psi) &:= A(\varphi \rightarrow \Diamond^{\leq}\psi) & \text{and} & \quad (\varphi <_{AE} \psi) := A(\varphi \rightarrow \Diamond^{<}\psi) \\
(\varphi \leq_{AA} \psi) &:= A(\psi \rightarrow \Box^{<}\neg\varphi) & \text{and} & \quad (\varphi <_{AA} \psi) := A(\psi \rightarrow \Box^{\leq}\neg\varphi)
\end{aligned}$$

The relationship between both, i.e., the semantically and syntactically defined families of binary preference connectives is discussed in van Benthem et al. (2009). In a nutshell, as regards the *EE*- and the *AE*-variants, both definitions (syntactic and semantic) are equivalent; concerning the *EA*- and the *AA*-variants, the equivalence only holds for a total  $\leq$  relation. In fact, drawing upon our encoding of  $\mathcal{PL}$  as presented in the next subsection §5.2, we have employed *Isabelle/HOL* for automatically verifying this sort of meta-theoretic correspondences; cf. Lines 4–12 in Fig. 11 in Appx. A.1.

## 5.2 Embedding $\mathcal{PL}$ in HOL

For the implementation of  $\mathcal{PL}$  we utilise the *shallow semantical embeddings* (SSE) technique, which encodes the language constituents of an object logic,  $\mathcal{PL}$  in our case, as expressions ( $\lambda$ -terms) in HOL. A core idea is to model (relevant parts of) the semantical structures of the object logic explicitly in HOL. This essentially shows that the object logic can be unraveled as a fragment of HOL and hence automated as such. For (multi-)modal normal logics, like  $\mathcal{PL}$ , the relevant semantical structures are relational frames constituted by sets of possible worlds/states and their accessibility relations.  $\mathcal{PL}$  formulas can thus be encoded as predicates in HOL taking possible worlds/states as arguments.<sup>11</sup> The detailed SSE of the basic operators of  $\mathcal{PL}$  in HOL is presented and formally tested in Appx. A.1. Further extensions to support reasoning with *ceteris paribus* clauses in  $\mathcal{PL}$  are studied there as well.

As a result, we obtain a combined, interactive and automated, theorem prover and model finder for  $\mathcal{PL}$  (and its extensions; cf. §5.4) realised within *Isabelle/HOL*. This is a new contribution, since we are not aware of any other existing implementation and automation of such a logic. Moreover, as we will demonstrate below, the SSE technique supports the automated assessment of meta-logical properties of the embedded logic in *Isabelle/HOL*, which in turn provides practical evidence for the correctness of our encoding.

<sup>11</sup> This corresponds to the well-known standard translation to first-order logic. Observe, however, that the additional expressivity of HOL allows us to also encode and flexibly combine non-normal modal logics (conditional, deontic, etc.; cf. Benzmüller 2017; Benzmüller et al. 2019a,b, 2022) and we can elegantly add quantifiers (cf. §5.4).

The embedding starts out with declaring the HOL base type  $\iota$ , which is denoting the set of possible states (or worlds) in our formalisation.  $\mathcal{PL}$  propositions are modelled as predicates on objects of type  $\iota$  (i.e., as *truth-sets* of states/worlds) and, hence, they are given the type  $\iota \rightarrow o$ , which is abbreviated as  $\sigma$  in the remainder. The *bettersness relation*  $\leq$  of  $\mathcal{PL}$  is introduced as an uninterpreted constant symbol  $\leq_{\iota \rightarrow \iota \rightarrow o}$  in HOL, and its strict variant  $<$  is introduced as an abbreviation  $<_{\iota \rightarrow \iota \rightarrow o}$  standing for the HOL term  $\lambda v. \lambda w. (v \leq w \wedge \neg(w \leq v))$ . In accordance with van Benthem et al. (2009), we postulate that  $\leq$  is a preorder, i.e., reflexive and transitive.

In a next step, the  $\sigma$ -type lifted logical connectives of  $\mathcal{PL}$  are introduced as abbreviations for  $\lambda$ -terms in the meta-logic HOL. The conjunction operator  $\wedge$  of  $\mathcal{PL}$ , for example, is introduced as an abbreviation  $\wedge_{\sigma \rightarrow \sigma \rightarrow \sigma}$ , which stands for the HOL term  $\lambda \varphi_{\sigma}. \lambda \psi_{\sigma}. \lambda w_{\iota}. (\varphi w \wedge \psi w)$ , so that  $\varphi_{\sigma} \wedge \psi_{\sigma}$  reduces to  $\lambda w_{\iota}. (\varphi w \wedge \psi w)$ , denoting the set<sup>12</sup> of all possible states  $w$  in which both  $\varphi$  and  $\psi$  hold. Analogously, for the negation we introduce an abbreviation  $\neg_{\sigma \rightarrow \sigma}$  which stands for  $\lambda \varphi_{\sigma}. \lambda w_{\iota}. \neg(\varphi w)$ .

The operators  $\diamond^{\leq}$  and  $\diamond^{<}$  use  $\leq$  and  $<$  as guards in their definitions. These operators are introduced as shorthand  $\diamond_{\sigma \rightarrow \sigma}^{\leq}$  and  $\diamond_{\sigma \rightarrow \sigma}^{<}$  abbreviating the HOL terms  $\lambda \varphi_{\sigma}. \lambda w_{\iota}. \exists v_{\iota}. (w \leq v \wedge \varphi v)$  and  $\lambda \varphi_{\sigma}. \lambda w_{\iota}. \exists v_{\iota}. (w < v \wedge \varphi v)$ , respectively.  $\diamond_{\sigma \rightarrow \sigma}^{\leq} \varphi_{\sigma}$  thus reduces to  $\lambda w_{\iota}. \exists v_{\iota}. (w \leq v \wedge \varphi v)$ , denoting the set of all worlds  $w$  so that  $\varphi$  holds in some world  $v$  that is at least as good as  $w$ ; analogously for  $\diamond_{\sigma \rightarrow \sigma}^{<}$ . Additionally, the *global existential* modality  $\mathbf{E}_{\sigma \rightarrow \sigma}$  is introduced as shorthand for the HOL term  $\lambda \varphi_{\sigma}. \lambda w_{\iota}. \exists v_{\iota}. (\varphi v)$ . The duals  $\square_{\sigma \rightarrow \sigma}^{\leq} \varphi_{\sigma}$ ,  $\square_{\sigma \rightarrow \sigma}^{<} \varphi_{\sigma}$  and  $\mathbf{A}_{\sigma \rightarrow \sigma} \varphi$  can easily be added so that they are equivalent to  $\neg \diamond_{\sigma \rightarrow \sigma}^{\leq} \neg \varphi_{\sigma}$ ,  $\neg \diamond_{\sigma \rightarrow \sigma}^{<} \neg \varphi_{\sigma}$  and  $\neg \mathbf{E}_{\sigma \rightarrow \sigma} \neg \varphi$  respectively.

Moreover, a special predicate  $[\varphi_{\sigma}]$  (read  $\varphi_{\sigma}$  is *valid*) for  $\sigma$ -type lifted  $\mathcal{PL}$  formulas in HOL is defined as an abbreviation for the HOL term  $\forall w_{\iota}. (\varphi_{\sigma} w)$ .

The encoding of object logic  $\mathcal{PL}$  in meta-logic HOL is presented in full detail in Appendix A.1.

Remember again that in the LOGIKEY methodology the modeler is not enforced to make an irreversible selection of an object logic (L1) before proceeding with the formalisation work at higher LOGIKEY layers. Instead the framework enables preliminary choices at all layers which can easily be revised by the modeler later on if this is indicated by e.g. practical experiments.

### 5.3 Formally Verifying Encoding's Adequacy

A pen-and-paper proof of the faithfulness (soundness & completeness) of the SSE easily follows from previous results regarding the SSE of propositional multi-modal logics (Benzmüller and Paulson 2010) and their quantified extensions (Benzmüller and Paulson 2013); cf. also Benzmüller (2019) and the references therein. We sketch such an argument below, as it provides an insight into the underpinnings of SSE for the interested reader.

By drawing upon the approach in Benzmüller and Paulson (2010), it is possible to define a mapping between semantic structures of the object logic  $\mathcal{PL}$  (preference models  $\mathcal{M}$ ) and the corresponding structures in HOL (general Henkin models  $\mathcal{H}^{\mathcal{M}}$ ),

<sup>12</sup> In HOL (with Henkin semantics) sets are associated with their characteristic functions.

in such a way that

$$\models^{\text{HOL}(\Gamma)} [\varphi_\sigma] \quad \text{iff} \quad \models_{\mathcal{PL}} \varphi \quad \text{iff} \quad \vdash_{\mathcal{PL}} \varphi,$$

where  $\vdash_{\mathcal{PL}}$  denotes derivability in the (complete) calculus axiomatised by van Benthem et al. (2009). Observe that  $\text{HOL}(\Gamma)$  corresponds to  $\text{HOL}$  extended with the relevant types and constants plus a set  $\Gamma$  of axioms encoding  $\mathcal{PL}$  semantic conditions, e.g., reflexivity and transitivity of  $\leq_{l \rightarrow l \rightarrow \sigma}$ .

Soundness of the SSE (i.e.,  $\models^{\text{HOL}(\Gamma)} [\varphi_\sigma]$  implies  $\models_{\mathcal{PL}} \varphi$ ) is particularly important since it ensures that our modelling does not give any ‘false positives’, i.e., proofs of  $\mathcal{PL}$  non-theorems. Completeness of the SSE (i.e.,  $\models_{\mathcal{PL}} \varphi$  implies  $\models^{\text{HOL}(\Gamma)} [\varphi_\sigma]$ ) means that our modelling does not give any ‘false negatives’, i.e., spurious counterexamples. Besides the pen-and-paper proof, completeness can also be mechanically verified by deriving the  $\sigma$ -type lifted  $\mathcal{PL}$  axioms and inference rules in  $\text{HOL}(\Gamma)$ ; cf. Fig. 11 and Fig. 12 in Appx. A.1.

To gain practical evidence for the faithfulness of our SSE of  $\mathcal{PL}$  in *Isabelle/HOL*, and also to assess proof automation performance, we have conducted numerous experiments in which we automatically verify meta-theoretical results on  $\mathcal{PL}$  as presented by van Benthem et al. (2009). Note that these statements thus play a role analogous to that of a requirements specification document (cf. Fig. 11 and Fig. 12 in Appx. A.1).

#### 5.4 Beyond $\mathcal{PL}$ : Extending the Encoding with Quantifiers and Conditionals

We can further extend our encoded logic  $\mathcal{PL}$  by adding quantifiers. This is done by identifying  $\forall x_\alpha s_\sigma$  with the HOL term  $\lambda w_l. \forall x_\alpha (s_\sigma w)$  and  $\exists x_\alpha s_\sigma$  with  $\lambda w_l. \exists x_\alpha (s_\sigma w)$ ; cf. *binder notation* in §4. This way quantified expressions can be seamlessly employed in our modelling at upper layers (as done exemplarily in §7). We refer the reader to Benzmüller and Paulson (2013) for a more detailed discussion (including faithfulness proofs) of SSEs for *quantified* (multi-)modal logics.

Moreover, observe that having a semantics based on *preferential structures* allows us to extend our logic with a (defeasible) conditional connective  $\Rightarrow$ . This can be done in several closely related ways. As an illustration, drawing upon an approach by Boutilier (1994), we can further extend the SSE of  $\mathcal{PL}$  by defining the connective:

$$\varphi_\sigma \Rightarrow \psi_\sigma := \mathbf{A}(\varphi_\sigma \rightarrow \Diamond^{\leq}(\varphi_\sigma \wedge \Box^{\leq}(\varphi_\sigma \rightarrow \psi_\sigma))).$$

An intuitive reading of this conditional statement would be: “every  $\varphi$ -state has a reachable  $\varphi$ -state such that  $\psi$  holds there in also in every reachable  $\varphi$ -state” (where we can interpret “reachable” as “at least as good”). This is equivalent, for finite models, to demanding that all ‘best’  $\varphi$ -states are  $\psi$ -states, cf. Lewis (1973). This can indeed be shown equivalent to the approach by Halpern (1997), who axiomatises a strict binary preference relation  $\succ^s$ , interpreted as “relative likelihood”.<sup>13</sup> For further discussion

<sup>13</sup> In fact, Halpern (1997) variant corresponds to employing the preference relation  $\prec_{AE}$  discussed previously, augmented with an additional constraint to cope with infinite-sized countermodels to irreflexivity (building upon an approach by Lewis (1973)). Thus,  $\psi \succ^s \varphi$  (read:  $\psi$  is more likely than  $\varphi$ ) iff every  $\varphi$ -state has a more likely  $\psi$ -state, say  $v$ , which *dominates*  $\varphi$  (i.e., no  $\varphi$ -state is more likely than  $v$ ). Halpern (1997) goes on to define a conditional operator as follows:  $\varphi \Rightarrow \psi := \mathbf{A}\neg\varphi \vee ((\varphi \wedge \psi) \succ^s (\varphi \wedge \neg\psi))$ .



regarding the properties and applications of this –and other similar– preference-based conditionals we refer the interested reader to the discussions in van Benthem (2009) and Liu (2011, Ch. 3).

## 6 Domain Specific Language (L2) – Value-Oriented Legal Theory

In this section we incrementally define a domain-specific language (DSL) for reasoning with values in a legal context. We start by defining a “logic of value preferences” on top of the object logic  $\mathcal{PL}$  (layer L1). This logic is subsequently encoded in *Isabelle/HOL*, and in the process it becomes suitably extended with custom means to encode the *discursive grammar* in §2. We thus obtain a HOL-based DSL formally modelling Lomfeld’s theory. This formally-verifiable DSL is then put to the test using theorem provers and model generators.

Recall from the discussion of the *discursive grammar* in §2 that value-oriented legal rules can become expressed as context-dependent preference statements between *value principles* (e.g. *RELIance*, *STABility*, *WILL*, etc.). Moreover, these value principles were informally associated to basic *values*. (i.e., *FREEDOM*, *UTILITY*, *SECURITY* and *EQUALITY*), in such a way as to arrange the first over (the quadrants of) a plane generated by two axes labelled by the latter. More specifically, each axis’ pole is labelled by a basic value, with values lying at contrary poles playing somehow antagonistic roles (e.g. *FREEDOM* vs. *SECURITY*). We recall the corresponding diagram (Fig. 2) below for the sake of illustration:

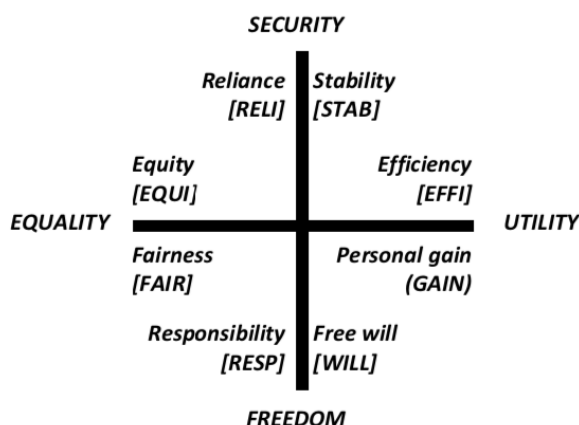


Fig. 4 Basic legal value system (ontology) by Lomfeld (2019)

Inspired by this theory, we model the notion of a (value) *principle* as consisting of a collection (in this case a set<sup>14</sup>) of *base values*. Thus, by considering principles

<sup>14</sup> Observe that in doing so we are simplifying Lomfeld’s value theory to the effect that, e.g., *STABility* becomes identified with *EFFiciency*. This simplified model has proven sufficient for our modelling work in

as structured entities, we can more easily define adequate notions of aggregation and conflict among them; cf. §6.

From a logical point of view it is additionally required to conceive value principles as truth-bearers, i.e., propositions.<sup>15</sup> We thus seem to face a dichotomy between, at the same time, modelling value principles as sets of basic values and modelling them as sets of worlds. In order to adequately tackle this modelling challenge we make use of the mathematical notion of a *Galois connection*.<sup>16</sup>

For the sake of exposition, Galois connections are to be exemplified by the notion of *derivation operators* in the theory of Formal Concept Analysis (FCA), from which we took inspiration; cf. Ganter and Wille (2012). FCA is a mathematical theory of concepts and concept hierarchies as formal ontologies, which finds practical application in many computer science fields such as data mining, machine learning, knowledge engineering, semantic web, etc.<sup>17</sup>

### 6.1 Some Basic FCA Notions

A *formal context* is a triple  $K = \langle G, M, I \rangle$  where  $G$  is a set of *objects*,  $M$  is a set of *attributes*, and  $I$  is a relation between  $G$  and  $M$  (usually called *incidence relation*), i.e.,  $I \subseteq G \times M$ . We read  $\langle g, m \rangle \in I$  as “the object  $g$  has the attribute  $m$ ”. Additionally we define two so-called *derivation operators*  $\uparrow$  and  $\downarrow$  as follows:

$$A\uparrow := \{m \in M \mid \langle g, m \rangle \in I \text{ for all } g \in A\} \quad \text{for } A \subseteq G \quad (1)$$

$$B\downarrow := \{g \in G \mid \langle g, m \rangle \in I \text{ for all } m \in B\} \quad \text{for } B \subseteq M \quad (2)$$

$A\uparrow$  is the set of all attributes shared by all objects from  $A$ , which we call the *intent* of  $A$ . Dually,  $B\downarrow$  is the set of all objects sharing all attributes from  $B$ , which we call the *extent* of  $B$ . This pair of derivation operators thus forms an antitone *Galois connection* between (the powersets of)  $G$  and  $M$ , and we always have that  $B \subseteq A\uparrow$  iff  $A \subseteq B\downarrow$ .

A *formal concept* (in a context  $K$ ) is defined as a pair  $\langle A, B \rangle$  such that  $A \subseteq G$ ,  $B \subseteq M$ ,  $A\uparrow = B$ , and  $B\downarrow = A$ . We call  $A$  and  $B$  the *extent* and the *intent* of the concept  $\langle A, B \rangle$ , respectively.<sup>18</sup> Indeed  $\langle A\uparrow\downarrow, A\uparrow \rangle$  and  $\langle B\downarrow, B\downarrow\uparrow \rangle$  are always concepts.

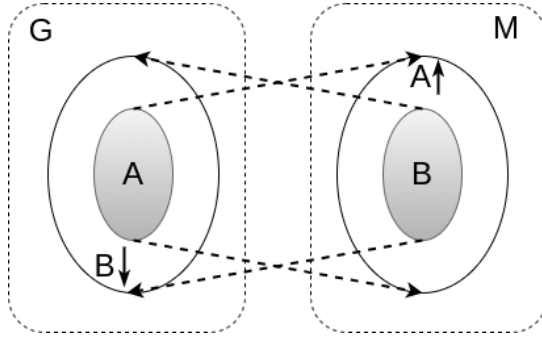
§7. A more granular encoding of principles is possible by adding a third axis to the value space in Fig. 4, thus allocating each principle to its own octant.

<sup>15</sup> We recall that, from a modal logic perspective, a proposition is modelled as the set of ‘worlds’ (i.e., states or situations) in which it holds. Informally, we want to be able to express the fact that a given principle, say legal STABILITY, is being observed or respected in a particular situation, or, abusing modal logic jargon, that the principle is ‘satisfied’ in that ‘world’. This can become further interpreted as providing a *justification* for why that world or situation is desirable.

<sup>16</sup> An old mathematician’s trick has been to employ –maybe unknowingly– Galois connections (resp. adjunctions) to relate two universes of mathematical objects with each other, in such a way that certain order structures become inverted (resp. preserved). In doing so, insights and results can be transferred from a well-known universe towards a less-known one, in order to gain information and help illuminate difficult problems; cf. the discussion in Ern  (2004).

<sup>17</sup> In particular, we want to highlight the potential of employing the powerful FCA methods, e.g., *attribute exploration* (Ganter et al. 2016), to prospective ‘legal value mining’ applications.

<sup>18</sup> The terms *extent* and *intent* are reminiscent of the philosophical notions of *extension* and *intension* (*comprehension*) reaching back to the 17th century *Logique de Port-Royal*.



**Fig. 5** A suggestive representation of a Galois connection between a set of objects  $G$  (e.g. worlds) and set of their attributes  $M$  (e.g. values).

The set of concepts in a formal context is partially ordered by set inclusion of their extents, or, dually, by the (reversing) inclusion of their intents. In fact, for a given formal context this ordering forms a complete lattice: its *concept lattice*. Conversely, it can be shown that every complete lattice is isomorphic to the concept lattice of some formal context. We can thus define lattice-theoretical meet and join operations on FCA concepts in order to obtain an algebra of concepts:<sup>19</sup>

$$\langle A_1, B_1 \rangle \wedge \langle A_2, B_2 \rangle := \langle (A_1 \cap A_2), (B_1 \cup B_2) \downarrow \uparrow \rangle \quad (3)$$

$$\langle A_1, B_1 \rangle \vee \langle A_2, B_2 \rangle := \langle (A_1 \cup A_2) \uparrow \downarrow, (B_1 \cap B_2) \rangle \quad (4)$$

## 6.2 A Logic of Value Preferences

In order to enable the modelling of Lomfeld’s legal theory as discussed in §2, we will enhance our object logic  $\mathcal{PL}$  with additional expressive means by drawing upon the FCA notions expounded above, and by assuming an *arbitrary* domain set  $\mathcal{V}$  of basic values.

A first step towards our legal DSL is to define a pair of operators  $\uparrow$  and  $\downarrow$  such that they form a Galois connection between the semantic domain  $\mathcal{W}$  of worlds/states of  $\mathcal{PL}$  (as ‘objects’  $G$ ) and the set of basic values  $\mathcal{V}$  (as ‘attributes’  $M$ ). By employing the operators  $\uparrow$  and  $\downarrow$  in an appropriate way, we can obtain additional well-formed  $\mathcal{PL}$  terms, thus converting our object logic  $\mathcal{PL}$  in a logic of value preferences. Details follow.

### *Principles, Values and Propositions*

We introduce a *formal context*  $K = \langle \mathcal{W}, \mathcal{V}, \mathcal{I} \rangle$  composed by the set of worlds  $\mathcal{W}$ , the set of basic values  $\mathcal{V}$ , and the (implicit) relation  $\mathcal{I} \subseteq \mathcal{W} \times \mathcal{V}$ , which we might interpret, intuitively, in a teleological sense:  $\langle w, v \rangle \in \mathcal{I}$  means that value  $v$  provides reasons for the situation (world/state)  $w$  to obtain.

<sup>19</sup> This result can be seamlessly stated for infinite meets and joins (infima and suprema) in the usual way. It corresponds to the first part of the so-called *basic theorem on concept lattices* (Ganter and Wille 2012).

Now, recall that we aim at modelling value principles as sets of basic values (i.e., elements of  $2^{\mathcal{V}}$ ), while, at the same time, conceiving of them as propositions (elements of  $2^{\mathcal{W}}$ ). Indeed, drawing upon the above FCA notions allows us to overcome this dichotomy. Given the formal context  $K = \langle \mathcal{W}, \mathcal{V}, \mathcal{I} \rangle$  we can define the pair of derivation operators  $\uparrow$  and  $\downarrow$  employing the corresponding definitions (1-2) above.

We can now employ these derivation operators to switch between the ‘(value) principles as sets of (basic) values’ and the ‘principles as propositions (sets of worlds)’ perspectives. Hence, we can now –recalling the informal discussion of the semantics of the object logic  $\mathcal{PL}$  in §5 – give an intuitive reading for truth at a world in a preference model to terms of the form  $P\downarrow$ ; namely, we can read  $\mathcal{M}, w \models P\downarrow$  as “principle  $P$  provides a reason for (state of affairs)  $w$  to obtain”. In the same vein, we can read  $\mathcal{M} \models A \rightarrow P\downarrow$  as “principle  $P$  provides a reason for proposition  $A$  being the case”.<sup>20</sup>

### Value Aggregation

Recalling Lomfeld’s theory, as discussed in §2, our logic of value preferences must provide means for expressing conditional preferences between value principles, according to the schema:

$$E_1 \wedge \dots \wedge E_n \Rightarrow (A_1 \oplus \dots \oplus A_n) < (B_1 \oplus \dots \oplus B_n)$$

As regards the preference relation (connective  $<$ ), we might think that, in principle, any choice among the eight preference relation variants in  $\mathcal{PL}$  (cf. §5) will work. Let us recall, however, that Lomfeld’s theory also presupposed some (no further specified) mechanism for aggregating value principles (operator  $\oplus$ ); thus, the joint selection of both a preference relation and an aggregation operator cannot be arbitrary: they need to interact in an appropriate way. We explore first a suitable mechanism for value aggregation before we get back to this issue.

Suppose that, for example, we are interested in modelling a legal case in which, say, the principle of “respect for property” *together with* the principle “economic benefit for society” *outweighs* the principle of “legal certainty”.<sup>21</sup> A binary connective  $\oplus$  for modelling this notion of *together with*, i.e., for aggregating legal principles (as reasons) must, expectedly, satisfy particular logical constraints in interaction with a (suitably selected) value preference relation  $<$ :

$$\begin{aligned} (A < B) \rightarrow (A < B \oplus C) \text{ but not } (A < B \oplus C) \rightarrow (A < B) & \quad \text{right aggregation} \\ (A \oplus C < B) \rightarrow (A < B) \text{ but not } (A < B) \rightarrow (A \oplus C < B) & \quad \text{left aggregation} \\ (B < A) \wedge (C < A) \rightarrow (B \oplus C < A) & \quad \text{union property (opt.)} \end{aligned}$$

For our purposes, the aggregation connectives are most conveniently defined using set union (FCA join), which gives us commutativity. As it happens, only the  $<_{AE}/\leq_{AE}$  and  $<_{EA}/\leq_{EA}$  variants from §5 satisfy the first two conditions. They are also the only

<sup>20</sup> Observe that this can be written semi-formally as: *for all  $w$  in  $\mathcal{M}$  we have that if  $\mathcal{M}, w \models A$  then  $\mathcal{M}, w \models P\downarrow$ , which can be interpreted as “ $P$  provides a reason for all those worlds that satisfy  $A$ ”.*

<sup>21</sup> Employing Lomfeld’s value theory this corresponds to RELiance together with personal GAIN outweighing STABILITY.

variants satisfying transitivity. Moreover, if we choose to enforce the optional third aggregation principle (called “union property”; cf. Halpern (1997)), then we would be left with only one variant to consider, namely  $\prec_{AE}/\preceq_{AE}$ .<sup>22</sup>

In the end, after extensive computer-supported experiments in *Isabelle/HOL* we have identified the following candidate definitions for the value aggregation and preference connectives which satisfy our modelling desiderata:<sup>23</sup>

- For the binary value aggregation connective  $\oplus$  we have identified the following two candidates (both taking two value principles and returning a proposition):

$$\begin{aligned} A \oplus_{(1)} B &:= (A \cap B) \downarrow \\ A \oplus_{(2)} B &:= (A \downarrow \vee B \downarrow) \end{aligned}$$

Observe that  $\oplus_1$  is based upon the join operation on the corresponding FCA formal concepts (see Def. 4).  $\oplus_2$  is a strengthening of the first, since  $(A \oplus_2 B) \subseteq (A \oplus_1 B)$ .

- For a binary preference connective  $\prec$  between propositions we have as candidates:

$$\begin{aligned} \varphi \prec_{(1)} \psi &:= \varphi \preceq_{AE} \psi \\ \varphi \prec_{(2)} \psi &:= \varphi \prec_{AE} \psi \\ \varphi \prec_{(3)} \psi &:= \varphi \preceq_{EA} \psi \\ \varphi \prec_{(4)} \psi &:= \varphi \prec_{EA} \psi \end{aligned}$$

In line with the LOGIKEY methodology, we consider the concrete choices of definitions for  $\prec$ ,  $\oplus$ , and even  $\Rightarrow$  (classical or defeasible) as parameters in our overall modelling process. No particular determination is enforced in the LOGIKEY approach, and we may alter any preliminary choices as soon as this appears appropriate. In this spirit we experimented with the listed different definition candidates for our connectives and explored their behaviour. We will present our final selection in §6.3.

### Promoting Values

Given that we aim at providing a logic of value preferences for use in legal reasoning, we still need to consider the mechanism by which we can link legal decisions, together with other legally relevant facts, to values. We conceive of such a mechanism as a sentence schema, which reads intuitively as: “Taking decision  $D$  in the presence of facts  $F$  promotes (value) principle  $P$ ”. The formalisation of this schema can indeed be seen as a new predicate in the domain-specific language (DSL) that we have been gradually defining in this section. In the expression  $Promotes(F, D, P)$  we have that  $F$

<sup>22</sup> Lacking any strong opinion regarding the correctness of transitivity or the union property, we have still chosen this latter variant for our case study in §7, since it offers several benefits for our current modelling purposes: it can be faithfully encoded in the language of  $\mathcal{PL}$  (van Benthem et al. 2009) and its behaviour is well documented in the literature; cf. Halpern (1997), Liu (2008, Ch. 4). In fact, as mentioned in §5.4, drawing upon the strict variant  $\prec_{AE}$  we can even define a defeasible conditional  $\Rightarrow$  in  $\mathcal{PL}$ .

<sup>23</sup> Respective tests are presented in Figs. 13–14 in Appx. A.1.

is a conjunction of facts relevant to the case (a proposition),  $D$  is the legal decision, and  $P$  is the value principle thereby promoted.<sup>24</sup>

$$\text{Promotes}(F, D, P) := F \rightarrow \Box^<(D \leftrightarrow \Diamond^<P\downarrow)$$

It is important to remark that, in the spirit of the LOGiKEY methodology, the definition above has arisen from the many iterations of encoding, testing and ‘debugging’ of the modelling of the ‘wild animal cases’ in §7 (until reaching a *reflective equilibrium*). We can still try to give this definition a somewhat intuitive interpretation, which might read along the lines of: “given the facts  $F$ , taking decision  $D$  is (necessarily) tantamount to (possibly) observing principle  $P$ ”, with the caveat that the (bracketed) modal expressions would need to be read in a non-alethic mood (e.g. deontically as discussed in §5.1).

### Value Conflict

Another important idea inspired from Lomfeld’s theory in §2 is the notion of value *conflict*. As discussed there (see Fig. 2), values are disposed around two axis of value coordinates, with values lying at contrary poles playing antagonistic roles. For our modelling purposes it makes thus sense to consider a predicate *Conflict* on worlds (i.e., a proposition) signalling situations where value conflicts appear. Taking inspiration from the traditional logical principle of *ex contradictio sequitur quodlibet*, which we may intuitively paraphrase, for the present purposes, as *ex conflictio sequitur quodlibet*,<sup>25</sup> we define *Conflict* as the set of those worlds in which *all* basic values become applicable:

$$\text{Conflict} := \bigwedge \{v\downarrow \text{ for all } v \text{ in } \mathcal{V}$$

Of course, and in the spirit of the LOGiKEY methodology, the specification of such a predicate can be further improved upon by the modeller as the need arises.

### 6.3 Instantiation as a HOL-based Legal DSL

In this subsection we encode our logic of value preferences in HOL (recall discussion in §4), building incrementally on top of the corresponding HOL-encoding for our (extended) object logic  $\mathcal{PL}$  in §5.2. In the process, our encoding will be gradually extended with custom means to encode Lomfeld’s legal theory (cf. §2). For the sake of illustrating a concrete, formally-verifiable modelling we also present in most cases the corresponding encoding in *Isabelle/HOL* (see also Appx. A.2).

<sup>24</sup> We adopt the terminology of *promoting* (or *advancing*) a value from the literature (Bench-Capon and Sartor 2003; Berman and Hafner 1993; Prakken 2002) understanding it in a teleological sense: a decision promoting a value principle means taking that decision *for the sake* of observing the principle; thus seeing the value principle *as a reason* for taking that decision.

<sup>25</sup> We shall not be held responsible for damages resulting from sloppy Latin paraphrasings!

In a preliminary step, we introduce a new base HOL-type  $c$  (for “contender”) as an (extensible) two-valued type introducing the legal parties “plaintiff” ( $p$ ) and “defendant” ( $d$ ). For this we employ in *Isabelle/HOL* the keyword `datatype`, which has the advantage of automatically generating (under the hood) the adequate axiomatic constraints (i.e., the elements  $p$  and  $d$  are distinct and exhaustive).

We also introduce a function, suggestively termed  $\text{other}_{c \rightarrow c}$ , with notation  $(\cdot)^{-1}$ . This function is used to return for a given party the *other* one; i.e.,  $p^{-1} = d$  and  $d^{-1} = p$ . Moreover, we add a ( $\sigma$ -lifted) predicate  $\text{For}_{c \rightarrow \sigma}$  to model the ruling *for* a given party and postulate that it always has to be ruled for either one party or the other:  $\text{For } x \leftrightarrow \neg \text{For } x^{-1}$ .

```

3 (*new datatype for parties/contenders (there could be more in principle)*)
4 datatype c = p | d (*plaintiff & defendant*)
5 fun other::"c⇒c" ("⁻¹") where "p⁻¹ = d" | "d⁻¹ = p"
6 (*new constant symbol: finding/ruling for party*)
7 consts For::"c⇒σ"
8 axiomatization where ForAx: "[|For x ↔ (¬For x⁻¹)|]"

```

As a next step, in order to enable the encoding of basic values, we introduce a four-valued datatype ( $'t$ ) `VAL` (corresponding to our domain  $\mathcal{V}$  of all values). Observe that this datatype is parameterised with a type variable  $'t$ . In the remainder we will always instantiate  $'t$  with the type  $c$  (see discussion below).

$$('t) \text{VAL} := \text{FREEDOM } 't \mid \text{UTILITY } 't \mid \text{SECURITY } 't \mid \text{EQUALITY } 't$$

We also introduce some convenient type-aliases:

$v := (c) \text{VAL} \rightarrow o$  is introduced as the type for (characteristic functions of) sets of basic values. The reader will recall that this corresponds to the characterisation of value principles as given in the previous subsection (i.e., elements of  $2^{\mathcal{V}}$ ).

It is important to note, however, that to enable the modelling of legal cases (plaintiff v. defendant) we need to further specify *legal* value principles *with respect to a legal party*, either plaintiff or defendant. For this we define  $cv := c \rightarrow v$  intended as the type for specific legal (value) principles (wrt. a legal party), so that they are functions taking objects of type  $c$  (either  $p$  or  $d$ ) to sets of basic values.

```

9 (*new parameterized datatype for abstract values (wrt. a given party)*)
10 datatype 't VAL = FREEDOM 't | UTILITY 't | SECURITY 't | EQUALITY 't
11 type_synonym v = "(c)VAL⇒bool" (*principles: sets of (abstract) values*)
12 type_synonym cv = "c⇒v" (*principles are specified wrt. a given party*)

```

We introduce useful set-constructor operators for basic values ( $\{\dots\}$ ) and a superscript notation for specification wrt. a legal party. As an illustration, recalling the discussion in §2, we have that, e.g., the legal principle of STABILITY' wrt. the plaintiff (notation  $\text{STAB}^p$ ) can be encoded as a two-element set of basic values (wrt. the plaintiff), i.e.,  $\{\text{SECURITY } p, \text{UTILITY } p\}$ .

The corresponding *Isabelle/HOL* encoding is:

```

13 (**notation for sets*)
14 abbreviation vset1 ("{|}") where "{|φ|} ≡ λx::(c)VAL. x=φ"
15 abbreviation vset2 ("{|,|}") where "{|α,β|} ≡ λx::(c)VAL. x=α ∨ x=β"
16 (**abstract values and value principles*)
17 abbreviation utility::cv ("UTILITY-") where "UTILITYx ≡ {|UTILITY x|}"
18 abbreviation security::cv ("SECURITY-") where "SECURITYx ≡ {|SECURITY x|}"
19 abbreviation equality::cv ("EQUALITY-") where "EQUALITYx ≡ {|EQUALITY x|}"
20 abbreviation freedom::cv ("FREEDOM-") where "FREEDOMx ≡ {|FREEDOM x|}"
21 abbreviation stab::cv ("STAB-") where "STABx ≡ {|SECURITY x, UTILITY x|}"
22 abbreviation effi::cv ("EFFI-") where "EFFIx ≡ {|UTILITY x, SECURITY x|}"
23 abbreviation gain::cv ("GAIN-") where "GAINx ≡ {|UTILITY x, FREEDOM x|}"
24 abbreviation will::cv ("WILL-") where "WILLx ≡ {|FREEDOM x, UTILITY x|}"
25 abbreviation resp::cv ("RESP-") where "RESPx ≡ {|FREEDOM x, EQUALITY x|}"
26 abbreviation fair::cv ("FAIR-") where "FAIRx ≡ {|EQUALITY x, FREEDOM x|}"
27 abbreviation equi::cv ("EQUI-") where "EQUIx ≡ {|EQUALITY x, SECURITY x|}"
28 abbreviation reli::cv ("RELI-") where "RELIx ≡ {|SECURITY x, EQUALITY x|}"

```

After defining legal (value) principles as combinations (in this case: sets<sup>26</sup>) of basic values (wrt. a legal party), we need to relate them to propositions (sets of worlds/states) in our logic  $\mathcal{PL}$ . For this we employ the *derivation operators* introduced in §6, whereby each value principle (set of basic values) becomes associated with a proposition (set of worlds) by means of the operator  $\downarrow$  (conversely for  $\uparrow$ ). We encode this by defining the corresponding *incidence* relation, or, equivalently, a function  $\mathcal{I}_{\iota \rightarrow \nu}$  mapping worlds/states (type  $\iota$ ) to sets of basic values (type  $\nu = (c) \text{VAL} \rightarrow o$ ). We define  $\downarrow_{\nu \rightarrow \sigma}$  so that, given some set of basic values  $V_\nu$ ,  $V_\nu \downarrow_\sigma$  denotes the set of all worlds that are  $\mathcal{I}$ -related to every value in  $V$  (analogously for  $\uparrow_{\sigma \rightarrow \nu}$ ). The modelling in the *Isabelle/HOL* proof assistant is as follows:

```

29 (**Value Theory*)
30 consts Irel::"ι⇒ν" ("I") (*incidence relation worlds-values*)
31 (*derivation operators (cf. theory of "formal concept analysis" *)
32 abbreviation intent::"σ⇒ν" ("I") where "W↑ ≡ λv. ∀x. W x → I x v"
33 abbreviation extent::"ν⇒σ" ("I") where "V↓ ≡ λw. ∀x. V x → I w x"
34 abbreviation extent_brkt ("I") where "[V] ≡ V↓" (*alternative notation*)

```

Thus we can intuitively read the proposition (set of worlds) denoted by  $\text{STAB}^p \downarrow$  as (those worlds in which) “the legal principle of STABILITY is observed wrt. the plaintiff”. For convenience, we introduce square brackets ( $[ \cdot ]$ ) as an alternative notation to  $\downarrow$ -postfixing in our DSL, so we have  $[V] = V \downarrow$ .

Now, our concrete choice of an aggregation operator for values (out of the two options presented in §6.2) is  $\oplus_{(2)}$ , which thus becomes encoded in HOL as:

$$A_\nu \oplus_{\nu \rightarrow \sigma} B_\nu := (A \downarrow)_\sigma \vee (B \downarrow)_\sigma$$

Analogously, the chosen preference relation ( $<$ ) is the variant  $<_{AE}$  (i.e.  $<_{(2)}$  from the candidate modelling options discussed in §6), which, recalling §5.1, becomes equivalently encoded as any of the following:

$$\begin{aligned} \varphi_\sigma <_{\sigma \rightarrow \sigma} \psi_\sigma &:= \forall s_i \varphi s \rightarrow (\exists t_i \psi t \wedge s < t) \\ \varphi_\sigma <_{\sigma \rightarrow \sigma} \psi_\sigma &:= \mathbf{A}_{\sigma \rightarrow \sigma}(\varphi \rightarrow \diamond_{\sigma \rightarrow \sigma}^< \psi) \end{aligned}$$

<sup>26</sup> Recall our discussion in §6 (cf. footnote 14). In a future modelling of a (suitably enhanced) *discursive grammar* (§2) we might take into account the order of combination of basic values in forming value principles, to the effect that, e.g., STABILITY can be properly distinguished from EFFICIENCY.



In a similar fashion, we encode in HOL the value-logical predicate *Promotes* as introduced in the previous subsection §6.2. The corresponding *Isabelle/HOL* encoding is shown below:

```

35 (*connective for aggregating value principles*)
36 abbreviation aggr ("[_@_]") where "[V1@V2] ≡ (V1) ∨ (V2)"
37 (*chosen variant for preference relation*)
38 abbreviation pref::"σ⇒σ⇒σ" ("_<") where "φ < ψ ≡ φ <AE ψ"
39 (*schema for value principle promotion*)
40 abbreviation "Promotes F D V ≡ [F → □~(D ↔ ◇~(V))]"

```

We have similarly encoded the proposition *Conflict* in HOL.

```

41 (*proposition for testing for value conflict*)
42 abbreviation conflict ("Conflict-") where (*conflict for value support*)
43 "Conflict* ≡ [SECURITY*] ∧ [EQUALITY*] ∧ [FREEDOM*] ∧ [UTILITY*]"

```

#### 6.4 Formally Verifying DSL's Adequacy

In this subsection we put our HOL-based legal DSL to the test by employing the automated tools integrated into *Isabelle/HOL*. In this process, the *discursive grammar*, as well as the continuous feedback by our legal domain expert (Lomfeld), served the role of a requirements specification for the formal verification of the adequacy of our modelling. We briefly discuss some of the conducted tests as shown in Fig. 6; further tests are presented in Fig. 16 in Appx. A.2 and in Benzmüller and Fuenmayor (2021).

In accordance with the dialectical interpretation of the *discursive grammar* (recall Fig. 2 in §2), our modelling foresees that observing values (wrt. the same party) from two opposing value quadrants, say RESP & STAB, or RELI & WILL, entails a value conflict; theorem provers quickly confirm this as shown in Fig. 6 (Lines 4–5). Moreover, observing values from two non-opposed quadrants, such as WILL & STAB

```

1 theory ValueOntologyTest imports ValueOntology (** Benzmüller, Fuenmayor & Lomfeld, 2021 **)
2 begin
3 (*value principles in two opposed quadrants: conflict*)
4 lemma "[RESP*] ∧ [STAB*] → Conflict*" by simp (*proof*)
5 lemma "[RELI*] ∧ [WILL*] → Conflict*" by simp (*proof*)
6 (*value principles in two non-opposed quadrants: no conflict*)
7 lemma "[WILL*] ∧ [STAB*] → Conflict*" nitpick oops (*countermodel*)
8 (*value principles in opposed quadrants for different parties: no conflict*)
9 lemma "[EQUI*] ∧ [GAIN*] → (Conflict* ∨ Conflict*)" nitpick oops (*countermodel*)
10 lemma "[RESP*] ∧ [STAB*] → (Conflict* ∨ Conflict*)" nitpick oops (*countermodel*)
11 lemma "[RELI*] ∧ [WILL*]" nitpick nitpick[satisfy] oops (*contingent: countermodel and model*)
12 (*value aggregation properties*)
13 lemma "[A@B] w → (A □ B) ↓ w" by simp
14 lemma "[A@B] w → A ↓ w" nitpick nitpick[satisfy] oops (*contingent: countermodel and model*)
15 lemma "[WILL*] < [STAB*] → [[WILL*] < [RELI*@STAB*]]" by blast (*proof*)
16 lemma "[RELI*@STAB*] < [WILL*] → [[STAB*] < [WILL*]]" by metis (*proof*)
17 lemma "[WILL*] < [RELI*@STAB*] → [[WILL*] < [STAB*]]"
18 nitpick nitpick[satisfy] oops (*contingent: countermodel and model*)
19 lemma "[STAB*] < [WILL*] → [[RELI*@STAB*] < [WILL*]]"
20 nitpick nitpick[satisfy] oops (*contingent: countermodel and model*)
21 end

```

Fig. 6 Verifying the DSL

```

Nitpick found a model for card  $\iota = 1$ :

Types:
c = {d, p}
c VAL = {FREEDOM d, FREEDOM p, UTILITY d, UTILITY p, EQUALITY d, EQUALITY p, SECURITY d, SECURITY p}
Constants:
BR = ( $\lambda x. \_$ )(( $\iota_1, \iota_1$ ) := True)
For = ( $\lambda x. \_$ )((d,  $\iota_1$ ) := False, (p,  $\iota_1$ ) := True)
 $\mathcal{I}$  = ( $\lambda x. \_$ )
  (( $\iota_1$ , FREEDOM d) := False, ( $\iota_1$ , FREEDOM p) := True, ( $\iota_1$ , UTILITY d) := False, ( $\iota_1$ , UTILITY p) := True,
  ( $\iota_1$ , EQUALITY d) := False, ( $\iota_1$ , EQUALITY p) := True, ( $\iota_1$ , SECURITY d) := False, ( $\iota_1$ , SECURITY p) := True)
other = ( $\lambda x. \_$ )(d := p, p := d)

```

Fig. 7 Satisfying model for the statement in Line 11 of Fig. 6.

(Line 7) should not imply any conflict: the model finder *Nitpick*<sup>27</sup> computes and reports a countermodel (not shown here) to the stated conjecture. A value conflict is also not implied if values from opposing quadrants are observed wrt. different parties (Lines 9–10).

Note that the notion of value conflict has deliberately not been aligned with logical inconsistency, neither in the object logic  $\mathcal{PL}$  nor in the meta-logic HOL. This way we can represent conflict situations in which, for instance, RELI and WILL (being conflicting values, see Line 5 in Fig. 6) are observed wrt. the plaintiff ( $p$ ), without leading to a logical inconsistency in *Isabelle/HOL* (thus avoiding ‘explosion’). In Line 11 of Fig. 6, for example, *Nitpick* is called simultaneously in both modes in order to confirm the contingency of the statement; as expected both a model (cf. Fig. 7) and countermodel (not displayed here) for the statement are returned. This value conflict can also be spotted by inspecting the satisfying models generated by *Nitpick*. One of such models is depicted in Fig. 7, where it is shown that (in the given possible world  $\iota_1$ ) all of the basic values (EQUALITY, SECURITY, UTILITY, and FREEDOM) are simultaneously observed wrt.  $p$ , which implies a value conflict according to our definition.

Such model structures as computed by *Nitpick* are ideally *communicated to* (and *inspected with*) domain experts (Lomfeld in our case) early on and checked for plausibility, which, in case of issues, might trigger adaptations to the axioms and definitions. Such a process may require several cycles until arriving at a state of *reflective equilibrium* (recall the discussion from §3) and, as a useful side effect, it conveniently fosters cross-disciplinary mutual understanding.

Further tests in Fig. 6 (Lines 13-20) assess the behaviour of the aggregation operator  $\oplus$  by itself, and also in combination with value preferences. For example, we test for a correct behaviour when ‘strengthening’ the right-hand side: if STAB is preferred over WILL, then STAB combined with, say, RELI is also preferred over WILL alone (Line 15). Similar tests are conducted for ‘weakening’ of the left-hand side.<sup>28</sup>

<sup>27</sup> *Nitpick* (Blanchette and Nipkow 2010) searches for, respectively enumerates, finite models or countermodels to a conjectured statement/lemma. By default *Nitpick* searches for countermodels, and model finding is enforced by stating the parameter keyword ‘satisfy’. These models are given as concrete interpretations of relevant terms in the given context so that the conjectured statement is satisfied or falsified.

<sup>28</sup> Further related tests are reported in Fig. 16 in Appx. A.2.

## 7 Applications (L3) – Assessment of Legal Cases

In this section we provide a concrete illustration of our reasoning framework by formally encoding and assessing two classic common law property cases concerning the appropriation of wild animals (“wild animal cases”): *Pierson v. Post*, and *Conti v. ASPCA*.<sup>29</sup>

Before starting with the analysis a word is in order about the support of our work by the tools *Sledgehammer* (Blanchette et al. 2016; Blanchette et al. 2013) and *Nitpick* (Blanchette and Nipkow 2010) in *Isabelle/HOL*. The ATP systems integrated via *Sledgehammer* in *Isabelle/HOL* include higher-order ATP systems, first-order ATP systems, and SMT (satisfiability modulo theories) solvers, and many of these systems in turn use efficient SAT solver technology internally. Indeed, proof automation with *Sledgehammer* and (counter)model finding with *Nitpick* were invaluable in supporting our exploratory modeling approach at various levels. These tools were very responsive in automatically proving (*Sledgehammer*), disproving (*Nitpick*), or showing consistency by providing a model (*Nitpick*). In the first case, references to the required axioms and lemmas were returned (which can be seen as a kind of abduction), and in the case of models and counter-models they often proved to be very readable and intuitive. In this section, we highlight some explicit use cases of *Sledgehammer* and *Nitpick*. They have been similarly applied at all levels as mentioned before.

We have split our analysis in layer L3 into two ‘sub-layers’ in order to highlight the separation between general legal & world knowledge (legal concepts and norms), from its ‘application’ to relevant facts in the process of deciding a case (factual/contextual knowledge). We shall first address the modelling of some background legal and world knowledge in §7.1, as minimally required in order to formulate each of our legal cases in the form of a logical *Isabelle/HOL* theory (cf. §7.2).

### 7.1 General Legal & World Knowledge

The realistic modelling of concrete legal cases requires further legal & world knowledge (LWK) to be taken into account. LWK is typically modelled in so called “upper” and “domain” ontologies. The question about which particular notion belongs to which category is difficult, and apparently there is no generally agreed answer in the literature. Anyhow, we introduce only a small and monolithic exemplary logical *Isabelle/HOL* theory,<sup>30</sup> called “GeneralKnowledge”, with a minimal amount of axioms and definitions as required to encode our legal cases. This LWK example includes a small excerpt of a much simplified “animal appropriation taxonomy”, where we associate “animal appropriation” (kinds of) situations with the value preferences they imply (i.e., conditional preference relations as discussed in §2 and §6).

In a more realistic setting this knowledge base would be further split and structured similarly to other legal or general ontologies, e.g., in the *Semantic Web* (Casanovas et

<sup>29</sup> Cf. Bench-Capon (2002), Berman and Hafner (1993), and Prakken (2002), and also T. F. Gordon and Walton (2006) for the significance of the *Pierson v. Post* case as a benchmark.

<sup>30</sup> *Isabelle* documents are suggestively called “theories”. They correspond to top-level modules bundling together related definitions, theories, proofs, etc.

al. 2016; Hoekstra et al. 2009). Note, however, that the expressiveness in our approach, unlike in many other legal ontologies or taxonomies, is by no means limited to definite underlying (but fixed) logical language foundations. We could thus easily decide for a more realistic modelling, e.g. avoiding simplifying propositional abstractions. For instance, the proposition “appWildAnimal”, representing the appropriation of one or more wild animals, can anytime be replaced by a more complex formula (featuring, e.g., quantifiers, modalities, and conditionals; see §5.4).

Next steps include interrelating notions introduced in our *Isabelle/HOL theory* “GeneralKnowledge” with values and value preferences, as introduced in the previous sections. It is here where the preference relations and modal operators of  $\mathcal{PL}$  as well as the notions introduced in our legal DSL are most useful. Remember that, at a later point and in line with the LOGIKEY methodology, we may in fact exchange  $\mathcal{PL}$  by an alternative choice of an object logic; or, on top of it, we may further modify our legal DSL, e.g., we might choose and assess alternative candidates for our connectives  $\leftarrow$  and  $\oplus$ ; moreover, we may want to replace material implication  $\rightarrow$  by a conditional implication to better support defeasible legal reasoning.<sup>31</sup>

We now briefly outline the *Isabelle/HOL* encoding of our example LWK; see Fig. 17 in Appx. A.3 for the full details.

First, some non-logical constants that stand for kinds of legally relevant situations (here: of appropriation) are introduced, and their meaning is constrained by some postulates:

```

3 (*LWK: kinds of situations addressed*)
4 consts appObject:: $\sigma$  appAnimal:: $\sigma$  (*appropriation of objects/animals in general*)
5       appWildAnimal:: $\sigma$  appDomAnimal:: $\sigma$  (*appropriation of wild/domestic animals*)
6 (*LWK: postulates for kinds of situations*)
7 axiomatization where
8 W1: "[appAnimal  $\rightarrow$  appObject]" and
9 W2: "[ $\neg$ (appWildAnimal  $\wedge$  appDomAnimal)]" and
10 W3: "[appWildAnimal  $\rightarrow$  appAnimal]" and
11 W4: "[appDomAnimal  $\rightarrow$  appAnimal]"

```

Then the ‘default’<sup>32</sup> legal rules for several situations (here: appropriation of animals) are formulated as conditional preference relations:

```

12 (*LWK: (prima facie) value preferences for kinds of situations*)
13 axiomatization where
14 R1: "[appAnimal  $\rightarrow$  ([STABp]  $\leftarrow$  [STABd])]" and
15 R2: "[appWildAnimal  $\rightarrow$  ([WILL*-1]  $\leftarrow$  [STAB*])]" and
16 R3: "[appDomAnimal  $\rightarrow$  ([STAB*-1]  $\leftarrow$  [RELI* $\oplus$ RESP*])]"

```

<sup>31</sup> Remember that a defeasible conditional implication can be defined employing  $\mathcal{PL}$  modal operators; cf. §5.4. Alternatively we may also opt for an SSE of a conditional logic in HOL using other approaches as, e.g., in Benzmüller (2013).

<sup>32</sup> We use of the term ‘default’ in the colloquial sense of ‘fallback’, noting however, that there exist in fact several (non-monotonic) logical systems aimed at modelling such a kind of *defeasible*, aka. “default”, behaviour for rules/conditionals (i.e., meaning that they can be ‘overruled’). One of them has been suggestively called “default logic”. We refer to Koons (2017) for a discussion. In fact, and in the spirit of LOGIKEY, we could have also employed, for encoding these rules, a  $\mathcal{PL}$ -defined defeasible conditional as discussed in §5.4. For the illustrative purposes of the present paper, and in view of the good performance of our present modelling, we did not yet find this step necessary.

For example, rule R2 could be read as: “In a wild-animals-appropriation kind of situation, observing STABILITY wrt. a party (say, the plaintiff) is preferred over observing WILL wrt. the other party (defendant)”. If there is no more specific legal rule from a precedent or a codified statute then these ‘default’ preference relations determine the result. Of course, this default is not arbitrary but itself an implicit normative setting of the existing legal statutes or cases. Moreover, we can have rules conditioned on more concrete legal *factors*.<sup>33</sup> As a didactic example, the legal rule R4 states that the *ownership* (say, the plaintiff’s) of the land on which the appropriation took place, together with the fact that the opposing party (defendant) acted out of *malice* implies a value preference of *reliance* and *responsibility* over *stability*. This rule has been chosen to reflect the famous common law precedent of *Keeble v. Hickeringill* (1704, 103 ER 1127; cf. also Bench-Capon (2002) and Berman and Hafner (1993)).

```
37>(*LWK: conditional value preferences, e.g. from precedents*)
38axiomatization where
39 R4: "[Mal x-1 ∧ Own x] → ([STABx-1] < [RESP*⊕RELI*])"
```

As already discussed, for ease of illustration, terms like “appWildAnimal” are modelled here as simple propositional constants. In practice, however, they may later be replaced, or logically implied, by a more realistic modelling of the relevant situational facts, utilising suitably complex (even quantified; cf. §5.4) formulas depicting states of affairs to some desired level of granularity.

For the sake of modelling the appropriation of objects, we have introduced an additional base type in our meta-logic HOL (recall §4). The type *e* (for ‘entities’) can be employed for terms denoting individuals (things, animals, etc.) when modelling legally relevant situations. Some simple vocabulary and taxonomic relationships (here: for wild and domestic animals) are specified to illustrate this.

```
17(*LWK: domain vocabulary*)
18typedecl e (*declares new type for 'entities'*)
19consts
20 Animal::"e⇒σ" Domestic::"e⇒σ" Fox::"e⇒σ" Parrot::"e⇒σ" Pet::"e⇒σ" FreeRoaming::"e⇒σ"
21(*LWK: domain knowledge (about animals)*)
22axiomatization where
23 W5: "[∀a. Fox a → Animal a]" and
24 W6: "[∀a. Parrot a → Animal a]" and
25 W7: "[∀a. (Animal a ∧ FreeRoaming a ∧ ¬Pet a) → ¬Domestic a]" and
26 W8: "[∀a. Animal a ∧ Pet a → Domestic a]"
```

As mentioned before, we have introduced some convenient legal *factors* into our example LWK to allow for the encoding of legal knowledge originating from precedents or statutes at a more abstract level. In our approach these factors are to be logically implied (as deductive arguments) from the concrete facts of the case (as exemplified in §A.4 below). Observe that our framework also allows us to introduce definitions for those factors for which clear legal specifications exist, such as property or

<sup>33</sup> The introduction of legal *factors* is an established practice in the implementation of case-based legal systems (cf. Bench-Capon (2017) for an overview). They can be conceived –as we do– as propositions abstracted from the facts of a case by the analyst/modeler in order to allow for assessing and comparing cases at a higher level of abstraction. Factors are typically either pro-plaintiff or pro-defendant, and their being true or false (resp. present or absent) in a concrete case can serve to invoke relevant precedents or statutes.

possession. At the present stage, we will provide some simple postulates constraining factors' interpretation.

```

27 (*LWK: legally-relevant, situational 'factors'*)
28 consts Own::"c⇒σ" (*object is owned by party c*)
29 Poss::"c⇒σ" (*party c has actual possession of object*)
30 Intent::"c⇒σ" (*party c has intention to possess object*)
31 Mal::"c⇒σ" (*party c acts out of malice*)
32 Mtn::"c⇒σ" (*party c respons. for maintenance of object*)
33 (*LWK: meaning postulates for general notions*)
34 axiomatization where
35 W9: "[Poss x → (¬Poss x⁻¹)]" and
36 W10: "[Own x → (¬Own x⁻¹)]"

```

Recalling §6 we relate the introduced legal factors (and relevant situational facts) to value principles and outcomes by means of the *Promotes* predicate:<sup>34</sup>

```

40 (*LWK: relate values, outcomes and situational 'factors'*)
41 axiomatization where
42 F1: "Promotes (Intent x) (For x) WILLx" and
43 F2: "Promotes (Mal x) (For x⁻¹) RESPx" and
44 F3: "Promotes (Poss x) (For x) STABx" and
45 F4: "Promotes (Mtn x) (For x) RESPx" and
46 F5: "Promotes (Own x) (For x) RELIx"
47 (*Theory is consistent, (non-trivial) model found*)
48 lemma True nitpick[satisfy,card ≠4] oops

```

Finally, the consistency of all axioms and rules provided is confirmed by *Nitpick*.

## 7.2 Pierson v. Post

This famous legal case (T. F. Gordon and Walton 2006) can be succinctly described as follows:

*Pierson killed and carried off a fox which Post already was hunting with hounds on public land. The Court found for Pierson (1805, 3 Cai R 175).*

For the sake of illustration we will consider in this subsection two modelling scenarios: in the first one a case is built to favour the defendant (Pierson); in the second one a case favouring the plaintiff (Post).

### Ruling for Pierson

The formal modelling of an argument in favour of Pierson is outlined next.<sup>35</sup>

First we introduce some minimal vocabulary: a constant  $\alpha$  of type  $e$  (denoting the appropriated animal), and the relations *pursue* and *capture* between the animal and one of the parties (of type  $c$ ). A background (generic) theory as well as the (contingent) case facts as suitably interpreted by Pierson's party are then stipulated:

<sup>34</sup> We note that our normative assignment here is widely in accordance with classifications in the AI & Law literature (Bench-Capon 2012; Berman and Hafner 1993).

<sup>35</sup> The entire formalisation of this argument is presented in Fig. 18 in Appx. A.4.

```

4 (*case-specific 'world-vocabulary'*)
5 consts  $\alpha::\text{"e"}$  (*appropriated animal (fox in this case) *)
6 consts Pursue::" $c \Rightarrow e \Rightarrow \sigma$ " Capture::" $c \Rightarrow e \Rightarrow \sigma$ "
7 (***** pro-defendant (Pierson) argument *****)
8 (*defendant's theory*)
9 abbreviation "dT1  $\equiv$  [ $\exists c. \text{Capture } c \ \alpha \ \wedge \ \neg \text{Domestic } \alpha \ \rightarrow \ \text{appWildAnimal}$ ]"
10 abbreviation "dT2  $\equiv$  [ $\forall c. \text{Pursue } c \ \alpha \ \rightarrow \ \text{Intent } c$ ]"
11 abbreviation "dT3  $\equiv$  [ $\forall c. \text{Capture } c \ \alpha \ \rightarrow \ \text{Poss } c$ ]"
12 abbreviation "d_theory  $\equiv$  dT1  $\wedge$  dT2  $\wedge$  dT3"
13 (*defendant's facts*)
14 abbreviation "dF1  $w \equiv$  Fox  $\alpha \ w$ "
15 abbreviation "dF2  $w \equiv$  FreeRoaming  $\alpha \ w$ "
16 abbreviation "dF3  $w \equiv$   $\neg$ Pet  $\alpha \ w$ "
17 abbreviation "dF4  $w \equiv$  Pursue p  $\alpha \ w$ "
18 abbreviation "dF5  $w \equiv$  Capture d  $\alpha \ w$ "
19 abbreviation "d_facts  $\equiv$  dF1  $\wedge$  dF2  $\wedge$  dF3  $\wedge$  dF4  $\wedge$  dF5"

```

The aforementioned decision of the court for Pierson was justified by the majority opinion. The essential preference relation in the case is implied in the idea that appropriation of (free-roaming) wild animals requires actual corporal possession. The manifest corporal link to the possessor creates legal certainty, which is represented by the value STABILITY and outweighs the mere WILL to possess by the plaintiff; cf. the arguments of classic lawyers cited by the majority opinion: “pursuit alone vests no property” (Justinian), and “corporal possession creates legal certainty” (Pufendorf). According to Lomfeld’s legal theory in §2 (cf. Fig. 2), this corresponds to a preference for the basic value SECURITY over FREEDOM. We can see that this legal rule R2, as introduced in the previous section (§7.1)<sup>36</sup> is indeed employed by *Isabelle/HOL*’s automated tools to prove that, given a suitable defendant’s theory, the (contingent) facts imply a decision in favour of Pierson in all ‘better’ worlds (which we could even give a ‘deontic’ reading as some sort of *recommendation*):

```

20 (*decision for defendant (Pierson) can be proven automatically*)
21 theorem Pierson: "d_theory  $\rightarrow$  [d_facts  $\rightarrow$   $\square$ -For d]"
22 by (smt F1 F3 ForAx R2 W5 W7 other.simps tSBR)

```

The previous ‘one-liner’ proof has indeed been automatically suggested by *Sledgehammer* (Blanchette et al. 2016; Blanchette et al. 2013) which we credit, together with the model finder *Nitpick* (Blanchette and Nipkow 2010), for doing the proof heavy-lifting in our work.

A proof argument in favour of Pierson that uses the same dependencies can also be constructed interactively using *Isabelle*’s human-readable proof language *Isar* (*Isabelle/Isar*; cf. Wenzel (2007)). The individual steps of the proof are this time formulated with respect to an explicit world/situation parameter  $w$ . The argument goes roughly as follows:

1. From Pierson’s facts and theory we infer that in the disputed situation  $w$  a wild animal has been appropriated:  $\text{appWildAnimal } w$
2. In this context, by applying the value preference rule R2, we get that observing STAB wrt. Pierson (d) is preferred over observing WILL wrt. Post (p):  

$$[[\text{WILL}^p] < [\text{STAB}^d]]$$

<sup>36</sup> Also observe that the legal precedent rule R4 of *Keeble v. Hickeringill* (see Fig. 17, Line 39) as appears in §7.1 does not apply to this case.

3. The possibility of observing WILL wrt. Post thus entails the possibility of observing STAB wrt. Pierson:  $\Box^{\leftarrow}[\text{WILL}^p] \rightarrow \Box^{\leftarrow}[\text{STAB}^d]$
4. Moreover, after instantiating the *value promotion* schema F1 (§7.1) for Post ( $p$ ), and acknowledging that his pursuing the animal (Pursue  $p \alpha$ ) entails his intention to possess (Intent  $p$ ), we obtain (for the given situation  $w$ ) a recommendation to ‘align’ any ruling for Post with the possibility of observing WILL wrt. Post:  $\Box^{\leftarrow}(\text{For } p \leftrightarrow \Box^{\leftarrow}[\text{WILL}^p]) w$ . Following the interpretation of the *Promotes* predicate given in §6, we can read this ‘alignment’ as involving both a logical entailment (left to right) and a justification (right to left); thus the possibility of observing WILL (wrt. Post) both entails and justifies (as a reason) a legal decision for Post.
5. Analogously, in view of Pierson’s ( $d$ ) capture of the animal (Capture  $d \alpha$ ), thus having taken possession of it (Poss  $d$ ), we infer from the instantiation of *value promotion* schema F3 (for Pierson) a recommendation to align a ruling for Pierson with the possibility of observing the value principle STAB wrt. Pierson:  $\Box^{\leftarrow}(\text{For } d \leftrightarrow \Box^{\leftarrow}[\text{STAB}^d]) w$
6. From (4) and (5) in combination with the courts duty to find a ruling for one of both parties (axiom *ForAx*) we infer, for the given situation  $w$ , that either the possibility of observing WILL wrt. Post or the possibility of observing STAB wrt. Pierson (or both) hold in every ‘better’ world/situation (thus becoming a recommended condition):  $\Box^{\leftarrow}(\Box^{\leftarrow}[\text{WILL}^p] \vee \Box^{\leftarrow}[\text{STAB}^d]) w$
7. From this and (3) we thus get that the possibility of observing STAB wrt. Pierson is recommended in the given situation  $w$ :  $\Box^{\leftarrow}(\Box^{\leftarrow}[\text{STAB}^d]) w$
8. And this together with (5) finally implies the recommendation to rule in favour of Pierson in the given situation  $w$ :  $\Box^{\leftarrow}(\text{For } d) w$

```

23 (*we reconstruct the reasoning process leading to the decision for the defendant*)
24 theorem Pierson': assumes d_theory and "d_facts w" shows "□←For d w"
25 proof -
26   have 1: "appWildAnimal w" using W5 W7 assms by blast
27   have 2: "[WILLp]←[STABd]" using 1 R2 assms by fastforce
28   have 3: "[□←[WILLp]] → □←[STABd]" using 2 tSBR by smt
29   have 4: "□←(For p ↔ □←[WILLp]) w" using F1 assms by meson
30   have 5: "□←(For d ↔ □←[STABd]) w" using F3 assms by meson
31   have 6: "□←((□←[WILLp]) ∨ (□←[STABd])) w" using 4 5 ForAx by (smt other.simps)
32   have 7: "□←(□←[STABd]) w" using 3 6 by blast
33   have 8: "□←(For d) w" using 5 7 by simp
34   then show ?thesis by simp
35 qed

```

The consistency of the assumed theory and facts (favouring Pierson) together with the other postulates from the previously introduced logical theories “GeneralKnowledge” and “ValueOntology” is verified by generating a (non-trivial) model using *Nitpick* (Line 38). Further tests confirm that the decision for Pierson (and analogously for Post) is compatible with the premises and, moreover, that for neither party value conflicts are implied.



```

36 (***** Further checks (using model finder) *****)
37 (*defendant's theory and facts are logically consistent*)
38 lemma "d_theory ∧ [d_facts]" nitpick[satisfy,card ≠3] oops (* (non-trivial) model found*)
39 (*decision for defendant is compatible with premises and lacks value conflicts*)
40 lemma "[¬Conflictd] ∧ [¬Conflictd] ∧ d_theory ∧ [d_facts ∧ For d]"
41 nitpick[satisfy,card ≠3] oops (* (non-trivial) model found*)
42 (*situations where decision holds for plaintiff are compatible too*)
43 lemma "[¬Conflictp] ∧ [¬Conflictd] ∧ d_theory ∧ [d_facts ∧ For p]"
44 nitpick[satisfy,card ≠3] oops (* (non-trivial) model found*)

```

We show next, how it is indeed possible to construct a case (theory) suiting Post using our approach.

### Ruling for Post

We model a possible counterargument in favour of Post claiming an interpretation (i.e., a *distinction* in case law methodology) in that the animal, being vigorously pursued (with large dogs and hounds) by a professional hunter, is not “free-roaming”. In doing this, the value preference  $[[WILL^p] < [STAB^d]]$  (for appropriation of wild animals), as in the previous Pierson’s argument, does not obtain. Furthermore, Post’s party postulates an alternative (suitable) value preference for hunting situations.

```

4 (*case-specific 'world-vocabulary'*)
5 consts α::"e" (*appropriated animal (fox in this case) *)
6 consts Pursue::"c⇒e⇒σ" Capture::"c⇒e⇒σ"
7 (***** pro-plaintiff (Post) argument *****)
8 (*acknowledges from defendant's theory*)
9 abbreviation "dT2 ≡ [∀c. Pursue c α → Intent c]"
10 abbreviation "dT3 ≡ [∀c. Capture c α → Poss c]"
11 (*theory amendment: the animal was chased by a professional hunter (Post); protecting
12 hunters' labor, thus fostering economic efficiency, prevails over legal certainty.*)
13 consts Hunter::"c⇒σ" hunting::"σ" (*new kind of situation: hunting*)
14 (*plaintiff's theory*)
15 abbreviation "pT1 ≡ [(∃c. Hunter c α ∧ Pursue c α) → hunting]"
16 abbreviation "pT2 ≡ ∀x. [hunting → ([STABx-1] < [EFFIx⊕WILLx])]" (*case-specific rule*)
17 abbreviation "pT3 ≡ ∀x. Promotes (hunting ∧ Hunter x) (For x) EFFIx"
18 abbreviation "p_theory ≡ pT1 ∧ pT2 ∧ pT3 ∧ dT2 ∧ dT3"
19 (*plaintiff's facts*)
20 abbreviation "pF1 w ≡ Fox α w"
21 abbreviation "pF2 w ≡ Hunter p w"
22 abbreviation "pF3 w ≡ Pursue p α w"
23 abbreviation "pF4 w ≡ Capture d α w"
24 abbreviation "p_facts ≡ pF1 ∧ pF2 ∧ pF3 ∧ pF4"

```

Note that an alternative legal rule (i.e., a possible argument for overruling in case law methodology) is presented in Line 16 above, entailing a value preference of the value principle combination EFFiciency together with WILL over STABILITY:  $[[STAB^d] < [EFFI^p \oplus WILL^p]]$ . Following the argument put forward by the dissenting opinion in the original case (3 Cai R 175) we might justify this new rule (inverting the initial value preference in the presence of EFFI) by pointing to the alleged public benefit of hunters getting rid of foxes, since the latter cause depredations in farms.

Accepting these modified assumptions the deductive validity of a decision for Post can in fact be proved and confirmed automatically, again, thanks to *Sledgehammer*:

```

25 | (*decision for plaintiff (Post) can be proven automatically (needs approx. 20s)*)
26 | theorem Post: "p_theory → [p_facts → □¬For p]"
27 | by (smt F1 F3 ForAx tBR SBR_def other.simps)

```

Similar to above, a detailed, interactive proof for the argument in favour of Post has been encoded and verified in *Isabelle/Isar*. We have also conducted further tests confirming the consistency of the assumptions and the absence of value conflicts.<sup>37</sup>

### 7.3 Conti v. ASPCA

An additional illustrative case study we have modelled in our framework is Conti v. ASPCA (353 NYS 2d 288; cf. Bench-Capon et al. (2005)). In a nutshell:

*Chester, a parrot owned by the ASPCA, escaped and was recaptured by Conti. The ASPCA found this out and reclaimed Chester from Conti. The court found for ASPCA.*

In this case, the court made clear that for domestic animals the opposite preference relation as the standard in Pierson’s case applies. More specifically, it was ruled that for a domestic animal it is in fact sufficient that the owner did not neglect or stopped caring for the animal, i.e., give up the responsibility for its maintenance (RESP). This, together with ASPCA’s reliance (RELI) in the parrot’s property, outweighs Conti’s corporal possession (STAB) of the animal:  $[[STAB^d] < [RELI^p \oplus RESP^p]]$ . Observe that a corresponding rule had previously been integrated as R3 into our legal & world knowledge (§7.1).

The plaintiff’s theory and facts are encoded analogously to the previous case:

```

5 | consts α::"e" (*appropriated animal (parrot in this case) *)
6 | consts Care::"c⇒e⇒σ" Prop::"c⇒e⇒σ" Capture::"c⇒e⇒σ"
7 | (***** pro-plaintiff (ASPCA) argument *****)
8 | (*plaintiff's theory*)
9 | abbreviation "pT1 ≡ [(∃c. Capture c α ∧ Domestic α) → appDomAnimal]"
10 | abbreviation "pT2 ≡ [∀c. Care c α → Mtn c]"
11 | abbreviation "pT3 ≡ [∀c. Prop c α → Own c]"
12 | abbreviation "pT4 ≡ [∀c. Capture c α → Poss c]" (*concedes' to defendant*)
13 | abbreviation "p_theory ≡ pT1 ∧ pT2 ∧ pT3 ∧ pT4"
14 | (*plaintiff's facts*)
15 | abbreviation "pF1 w ≡ Parrot α w"
16 | abbreviation "pF2 w ≡ Pet α w"
17 | abbreviation "pF3 w ≡ Care p α w"
18 | abbreviation "pF4 w ≡ Prop p α w"
19 | abbreviation "pF5 w ≡ Capture d α w"
20 | abbreviation "p_facts ≡ pF1 ∧ pF2 ∧ pF3 ∧ pF4 ∧ pF5"

```

Accepting these assumptions the deductive validity of a decision for the plaintiff (ASPCA) can again be proved and confirmed automatically (thanks to *Sledgehammer*):

```

21 | (*decision for plaintiff (ASPCA) can be proven automatically*)
22 | theorem ASPCA: "p_theory → [p_facts → □¬For p]"
23 | by (smt F3 F4 F5 ForAx R3 W6 W8 tBR SBR_def other.simps(1))

```

<sup>37</sup> See the complete modelling in Fig. 19 in Appx. A.4.

In an analogous manner to Pierson’s case, an interactive proof in *Isabelle/Isar* has been encoded and verified, and the consistency of the assumptions and the absence of value conflicts has been confirmed.<sup>38</sup>

## 8 Related and Further Work

Custom software systems for legal case-based reasoning have been developed in the AI & Law community, beginning with the influential HYPO system in the 1980s (Rissland and Ashley 1987); cf. also the survey paper by Bench-Capon (2017). In later years, there has been a gradual shift of interest from rule-based non-monotonic reasoning (e.g., logic programming) towards argumentation-based approaches (see Prakken and Sartor (2015) for an overview); however, we are not aware of any other work that uses higher-order theorem proving and proof assistants (the argumentation logic of Krause et al. (1995) is an early related effort that is worth mentioning). Another important aspect of our work concerns value-oriented legal reasoning and deliberation, where a considerable amount of work has been presented in AI & Law in response to the challenge posed by Berman and Hafner (1993). Our approach, based mainly on Lomfeld’s (2015; 2019) theory, has also been influenced by some of this work, in particular by Bench-Capon (2002), Bench-Capon and Sartor (2003), and Prakken (2002).

We are currently working towards further refining the modelling of Lomfeld’s legal theory with the aim of providing more expressive (combinations of) object logics at LOGIKEY layer L1. In this regard, it is somehow remarkable that the use of material implication to encode rules has proven sufficient for the illustrative purposes of this paper. However, it is important to note that a more realistic modelling of legal cases must also provide mechanisms to deal with the inevitable emergence of conflicts and contradictions in normative reasoning (overruling, conflict resolution, etc.). In line with the LOGIKEY approach, we are working at introducing conditional connectives in our object logics with the aim of enabling *defeasible* (or *default*) reasoning. Such connectives can be introduced by reusing the modal operators of  $\mathcal{PL}$  (recalling the discussion in §5.4) or, alternatively, through the shallow semantical embedding (Benzmüller 2019) of a suitable conditional logic in HOL (Benzmüller 2017). Moreover, special kinds of paraconsistent (modal-like) *Logics of Formal Inconsistency* (Carnielli et al. 2021) can also be integrated into our modelling to enable the non-explosive representation of (and recovery from) contradictions by purely object-logical means (cf. Fuenmayor (2020) for a related encoding in *Isabelle/HOL*). In a similar vein, we think that some of the recent work that employs expressive deontic logics for value-based legal balancing (e.g. Maranhão and Sartor (2019) and the references therein) can be fruitfully integrated in our approach. It is the pluralistic nature of LOGIKEY, realised within a dynamic modelling framework (e.g. *Isabelle/HOL*), that enables and supports such improvements without requiring expensive technical adjustments to the underlying base reasoning technology.

As a broader application scenario, we are currently proposing that ethico-legal value-oriented theories and ontologies should constitute a core ingredient to enable

<sup>38</sup> The full details of the encoding are presented in Fig. 20 in Appx. A.5.

the computation, assessment and communication of rational justifications and explanations in the future ethico-legal governance of AI (Benzmüller and Lomfeld 2020). Thus, a sound and trustworthy implementation of any legally accountable ‘moral machine’ requires the development of formal theories and ontologies for the legal domain to guide and interconnect the encoding of concrete regulatory codes and legal cases. Understanding legal reasoning as dialectical practical argumentation, the pluralist interpretation of concrete legal rules arguably requires a complementary ethico-legal value-oriented theory such as, e.g., the *discursive grammar* of justification by Lomfeld (2019), which we formally encoded in this paper. In this sense, some first positive evidence has been provided regarding challenges that we have previously identified with respect to the ethical-legal governance of future AI systems (Fuenmayor and Benzmüller 2020). Indeed, it was this broader vision that primarily motivated our work on value-oriented legal reasoning in the first place.

## 9 Conclusion

We illustrate the application of the LOGIKEY knowledge engineering methodology and framework to enable the interdisciplinary collaboration among different specialist roles. In the present case, they are a lawyer and legal philosopher (L.) and two computer scientists (B. and F.) who join forces with the aim of formally modelling a value-oriented legal theory (*discursive grammar* by Lomfeld 2019) for the sake of providing means for computer-automated prediction and assessment of legal case decisions.

From a technical perspective, the core objective of this article has been to demonstrate that the LOGIKEY methodology appears indeed suitable for the task of value-oriented legal reasoning. As instantiated in the present work, the LOGIKEY methodology builds upon a HOL-encoding of a modal logic of preferences to model a domain-specific theory of value-based legal balancing. In combination with further legal and world knowledge this theory has been successfully employed for the formal encoding and computer-supported assessment, using the *Isabelle/HOL* system, of illustrative legal cases in property law (“wild animal cases”).

It is the flexibility of the multi-layer modelling which is novel in our approach, in combination with a very rich support for automated reasoning in expressive, quantified classical and non-classical logics, thereby rejecting the idea that knowledge representation means should be limited *prima facie* to decidable logic frameworks, due to complexity or performance considerations. In the LOGIKEY approach, the choice of a particular object logic is deliberately left to the knowledge engineer. The range of options varies from well-manageable decidable logics to sophisticated quantified non-classical logics and combinations thereof, depending on what is best suited to handle a particular knowledge representation (and reasoning) task at hand.

From a legal perspective, the reconstruction of legal balancing is, already with classical argumentative tools, a non-trivial task, which is methodologically not yet settled (Sieckmann 2010). Here, our work proposed the structuring of legal balancing by means of a dialectical ethico-legal value system (*discursive grammar*). Legal rules and their various interpretations can thus be represented within a unified yet plu-

realistic logic of value preferences. The integration of this logic and the value system within the dynamic HOL-based modelling environment allows us to experiment with different forms of interpretation. This enables us, not only to find more accurate reconstructions of legal argumentation, but also supports the modelling of value-based legal balancing, taking into account notions of value preference, aggregation, promotion and conflict; and also in a manner amenable to computer automation. The modelling of Lomfeld’s legal theory in LOGIKEY enabled us to successfully predict (and to some extent justify) case outcomes by ‘just using logic’, employing qualitative value preferences without the necessity to bring in numbers and weights into the model.

From a general perspective, supporting interactive and automated value-oriented legal argumentation on the computer is a non-trivial challenge which we address, for reasons as defended, e.g., by Bench-Capon (2020), with symbolic AI techniques and formal methods. Motivated by recent pleas for *explainable and trustworthy AI*, our primary goal is to work towards the development of ethico-legal governors for future generations of intelligent systems, or more generally, towards some form of legally and ethically *reasonable machines* (Benzmüller and Lomfeld 2020) capable of exchanging rational justifications for the actions they take. While building up a capacity to engage in value-oriented legal argumentation is just one of a multitude of challenges this vision is faced with, it clearly constitutes an important stepping stone towards this ambitious long-term goal.

**Acknowledgements** We thank the unknown reviewers of our prior paper at the MLR 2020 workshop for their valuable comments and suggestions that have led to significant improvements of this article.

## References

- Aleven, V. (1997). *Teaching case-based reasoning through a model and examples*. PhD Dissertation University of Pittsburgh.
- Alexy, R., ed. (1978). *Theorie der juristischen Argumentation*. Frankfurt/M: Suhrkamp.
- (2000). “On the Structure of Legal Principles”. In: *Ratio Juris* 13, pp. 294–304.
- (2003). “On Balancing and Subsumption: A Structural Comparison”. In: *Ratio Juris* 16, pp. 433–449.
- Andrews, P. B. (1972a). “General Models and Extensionality”. In: *J. Symb. Log.* 37.2, pp. 395–397. DOI: 10.2307/2272982.
- (1972b). “General Models, Descriptions, and Choice in Type Theory”. In: *J. Symb. Log.* 37.2, pp. 385–394.
- Arkin, R. C., P. Ulam, and B. A. Duncan (2009). Tech. rep. GVU-09-02. Georgia Institute of Technology.
- Ashley, K. D. (1990). *Modelling Legal Argument: Reasoning with Cases and Hypotheticals*. Cambridge/MA: MIT Press.
- Barak, A. (2012). *Proportionality*. Cambridge University Press.
- Bench-Capon, T. (2020a). “Ethical approaches and autonomous systems”. In: *Artif. Intell.* 281.

- Bench-Capon, T. (2002). “The missing link revisited: The role of teleology in representing legal argument”. In: *Artif. Intell. Law* 10.1-3, pp. 79–94.
- (2012). “Representing Popov v Hayashi with dimensions and factors”. In: *Artif. Intell. Law* 20, pp. 15–35.
- (2017). “Hypo’s legacy: introduction to the virtual special issue”. In: *Artif. Intell. Law* 25.2, pp. 205–250.
- (2020b). “The Need for Good Old-Fashioned AI and Law.” In: *In International Trends in Legal Informatics: A Festschrift for Erich Schweighofer*. Ed. by W. Hötendorfer, C. Tschol, and F. Kummer. Weblaw AG.
- Bench-Capon, T., K. Atkinson, and A. Chorley (2005). “Persuasion and value in legal argument”. In: *J. Log. Comput.* 15, pp. 1075–1097.
- Bench-Capon, T. and G. Sartor (2003). “A model of legal reasoning with cases incorporating theories and value”. In: *Artif. Intell.* 150, pp. 97–143.
- Benzmüller, C. (2013). “Automating Quantified Conditional Logics in HOL”. In: *23rd International Joint Conference on Artificial Intelligence (IJCAI-13)*. Ed. by F. Rossi. Beijing, China: AAAI Press, pp. 746–753. ISBN: 978-1-57735-633-2.
- (2017). “Cut-Elimination for Quantified Conditional Logic”. In: *J. Philos. Log.* 46.3, pp. 333–353. DOI: 10.1007/s10992-016-9403-0.
- (2019). “Universal (Meta-)Logical Reasoning: Recent Successes”. In: *Sci. Comput. Program.* 172, pp. 48–62. DOI: 10.1016/j.scico.2018.10.008.
- Benzmüller, C. and P. Andrews (2019). “Church’s Type Theory”. In: *The Stanford Encyclopedia of Philosophy*. Ed. by E. N. Zalta. Summer 2019. Metaphysics Research Lab, Stanford University, 1–62 (in pdf version).
- Benzmüller, C., C. Brown, and M. Kohlhase (2004). “Higher-Order Semantics and Extensionality”. In: *J. Symb. Log.* 69.4, pp. 1027–1088. DOI: 10.2178/jsl/1102022211.
- Benzmüller, C., A. Farjami, P. Meder, and X. Parent (2019a). “I/O Logic in HOL”. In: *J. Appl. Logics – IfCoLoG J. Logics their Appl. (Special Issue: Reason. for Leg. AI)* 6.5. Ed. by L. Robaldo and L. van der Torre, pp. 715–732.
- Benzmüller, C., A. Farjami, and X. Parent (2019b). “Åqvist’s Dyadic Deontic Logic E in HOL”. In: *J. Appl. Logics – IfCoLoG J. Logics their Appl. (Special Issue: Reason. for Leg. AI)* 6.5. Ed. by L. Robaldo and L. van der Torre, pp. 733–755.
- (2022). “Dyadic Deontic Logic in HOL: Faithful Embedding and Meta-Theoretical Experiments”. In: *New Developments in Legal Reasoning and Logic: From Ancient Law to Modern Legal Systems*. Ed. by S. Rahman, M. Armgardt, and H. C. Nordtveit Kvernenes. Springer Nature Switzerland AG, pp. 353–377. DOI: 978-3-030-70084-3\_14.
- Benzmüller, C. and D. Fuenmayor (2021). “Value-oriented Legal Argumentation in Isabelle/HOL”. In: *International Conference on Interactive Theorem Proving (ITP), Proceedings*. Ed. by L. Cohen and C. Kaliszyk. Vol. 193. LIPIcs 23. Schloss Dagstuhl - Leibniz-Zentrum für Informatik, 23:1–23:18. DOI: 10.4230/LIPIcs.ITP.2021.7.
- Benzmüller, C. and B. Lomfeld (2020). “Reasonable Machines: A Research Manifesto”. In: *KI 2020: Advances in Artificial Intelligence – 43rd German Conference on Artificial Intelligence, Bamberg, Germany, September 21–25, 2020, Proceedings*. Ed. by U. Schmid, F. Klügl, and D. Wolter. Vol. 12352. Lecture Notes in

- Artificial Intelligence. Springer, Cham, pp. 251–258. ISBN: 978-3-030-30178-1. DOI: 10.1007/978-3-030-58285-2\_20.
- Benzmüller, C. and D. Miller (2014). “Automation of Higher-Order Logic”. In: *Handbook of the History of Logic, Volume 9 — Computational Logic*. Ed. by D. M. Gabbay, J. H. Siekmann, and J. Woods. North Holland, Elsevier, pp. 215–254. ISBN: 978-0-444-51624-4. DOI: 10.1016/B978-0-444-51624-4.50005-8.
- Benzmüller, C., X. Parent, and L. van der Torre (2020). “Designing Normative Theories for Ethical and Legal Reasoning: LogiKEY Framework, Methodology, and Tool Support”. In: *Artif. Intell.* 287, p. 103348. DOI: 10.1016/j.artint.2020.103348.
- Benzmüller, C. and L. C. Paulson (2010). “Multimodal and Intuitionistic Logics in Simple Type Theory”. In: *The Log. J. IGPL* 18.6, pp. 881–892. DOI: 10.1093/jigpal/jzp080.
- (2013). “Quantified Multimodal Logics in Simple Type Theory”. In: *Logica Universalis (Special Issue on Multimodal Logics)* 7.1, pp. 7–20. DOI: 10.1007/s11787-012-0052-y.
- Berman, D. and C. Hafner (1993). “Representing teleological structure in case-based legal reasoning: the missing link”. In: *Proceedings 4th ICAIL*. New York: ACM Press, pp. 50–59.
- Blanchette, J. C., C. Kaliszyk, L. C. Paulson, and J. Urban (2016). “Hammering towards QED”. In: *J. Formaliz. Reason.* 9.1, pp. 101–148.
- Blanchette, J. C., S. Böhme, and L. C. Paulson (2013). “Extending Sledgehammer with SMT Solvers”. In: *J. Autom. Reason.* 51.1, pp. 109–128.
- Blanchette, J. C. and T. Nipkow (2010). “Nitpick: A Counterexample Generator for Higher-Order Logic Based on a Relational Model Finder”. In: *ITP 2010*. Ed. by M. Kaufmann and L. C. Paulson. Vol. 6172. LNCS. Springer, pp. 131–146.
- Boutilier, C. (1994). “Toward a logic for qualitative decision theory”. In: *Principles of knowledge representation and reasoning*. Elsevier, pp. 75–86. DOI: 10.1016/B978-1-4832-1452-8.50104-4.
- Carnielli, W., M. E. Coniglio, and D. Fuenmayor (2021). “Logics of Formal Inconsistency Enriched with Replacement: An Algebraic and Modal Account”. In: *The Rev. Symb. Log.* online first, pp. 1–36. DOI: 10.1017/S1755020321000277.
- Carnielli, W. and M. E. Coniglio (2020). “Combining Logics”. In: *The Stanford Encyclopedia of Philosophy*. Ed. by E. N. Zalta. Fall 2020. Metaphysics Research Lab, Stanford University.
- Carnielli, W. et al. (2008). *Analysis and Synthesis of Logics*. Applied Logics Series 35. Springer.
- Casanovas, P. et al. (2016). “Semantic Web for the Legal Domain: The next step”. In: *Semantic Web* 7.3, pp. 213–227. DOI: 10.3233/SW-160224.
- Church, A. (1940). “A Formulation of the Simple Theory of Types”. In: *J. Symb. Log.* 5.2, pp. 56–68. DOI: 10.2307/2266170.
- Clark, B. (1991). *Political Economy: A Comparative Approach*. New York: Praeger.
- Daniels, N. (2020). “Reflective Equilibrium”. In: *The Stanford Encyclopedia of Philosophy*. Ed. by E. N. Zalta. Summer 2020. Metaphysics Research Lab, Stanford University.

- Dworkin, R. (1978). *Taking rights seriously*. OCLC: 4313351. Cambridge, Mass: Harvard Univ. Press. 371 pp. ISBN: 978-0-674-86711-6.
- Erné, M. (2004). “Adjunctions and Galois Connections: Origins, History and Development”. In: *Galois Connections and Applications*. Ed. by K. Denecke, M. Ern , and S. L. Wismath. Dordrecht: Springer Netherlands, pp. 1–138. DOI: 10.1007/978-1-4020-1898-5\_1.
- Eysenck, H. (1954). *The Psychology of Politics*. London: Routledge.
- Feteris, E. (2017). *Fundamentals of Legal Argumentation*. Dordrecht: Springer.
- Fuenmayor, D. (2020). “Topological semantics for paraconsistent and paracomplete logics”. In: *Arch. Formal Proofs*. [https://isa-afp.org/entries/Topological\\_Semantics.html](https://isa-afp.org/entries/Topological_Semantics.html), Formal proof development. ISSN: 2150-914x.
- Fuenmayor, D. and C. Benzm ller (2019). “A Computational-Hermeneutic Approach for Conceptual Explicitation”. In: *Model-Based Reasoning in Science and Technology. Inferential Models for Logic, Language, Cognition and Computation*. Ed. by A. Nepomuceno et al. Vol. 49. SAPERE. Springer, pp. 441–469. DOI: 10.1007/978-3-030-32722-4\_25.
- (2020). “Normative Reasoning with Expressive Logic Combinations”. In: *ECAI 2020 – 24th European Conference on Artificial Intelligence, June 8-12, Santiago de Compostela, Spain*. Ed. by G. De Giacomo et al. Vol. 325. Frontiers in Artificial Intelligence and Applications. IOS Press, pp. 2903–2904. DOI: 10.3233/FAIA200445.
- Ganter, B., S. Obiedkov, S. Rudolph, and G. Stumme (2016). *Conceptual exploration*. Springer.
- Ganter, B. and R. Wille (2012). *Formal concept analysis: mathematical foundations*. Springer Berlin.
- Gibbons, J. and N. Wu (2014). “Folding domain-specific languages: deep and shallow embeddings (functional Pearl)”. In: *Proceedings of the 19th ACM SIGPLAN international conference on Functional programming, Gothenburg, Sweden, September 1-3, 2014*. Ed. by J. Jeuring and M. M. T. Chakravarty. ACM, pp. 339–347. DOI: 10.1145/2628136.2628138.
- Goodman, N. (1955). *Fact, fiction, and forecast*. Harvard University Press.
- Gordon, T. and D. Walton (2012). “A Carneades reconstruction of Popov v Hayashi”. In: *Artif. Intell. Law* 20, pp. 37–56.
- Gordon, T. F. and D. Walton (2006). “Pierson vs. Post revisited”. In: *Front. Artif. Intell. Appl.* 144, p. 208.
- Grabmair, M. (2016). *Modeling purposive legal argumentation and case outcome prediction using argument schemes in the value judgment formalism*. Dissertation University of Pittsburgh.
- Gruber, T. (1993). “A Translation Approach to Portable Ontology Specifications”. In: *Knowl. Acquis.* 5, pp. 199–220.
- (2009). “Ontology”. In: *Encyclopedia of Database Systems*. Ed. by L. Liu and M. T.  zsu. Springer.
- Hage, J. (1997). *Reasoning With Rules*. Dordrecht: Kluwer.
- Halpern, J. Y. (1997). “Defining relative likelihood in partially-ordered preferential structures”. In: *J. Artif. Intell. Res.* 7, pp. 1–24.



- Henkin, L. (1950). “Completeness in the Theory of Types”. In: *J. Symb. Log.* 15.2, pp. 81–91.
- Hoekstra, R., J. Breuker, M. D. Bello, and A. Boer (2009). “LKIF Core: Principled Ontology Development for the Legal Domain”. In: *Law, Ontologies and the Semantic Web - Channelling the Legal Information Flood*. Ed. by J. Breuker and et.al. Vol. 188. *Frontiers in Artificial Intelligence and Applications*. IOS Press, pp. 21–52. DOI: 10.3233/978-1-58603-942-4-21.
- Hofstede, G. (2001). *Culture’s Consequences*. Thousands Oaks: Sage.
- Horty, J. (2011). “Rules and reasons in the theory of precedent”. In: *Leg. Theory* 17, pp. 1–33.
- Inglehart, R. (2018). *Cultural Evolution*. Cambridge University Press.
- Koons, R. (2017). “Defeasible Reasoning”. In: *The Stanford Encyclopedia of Philosophy*. Ed. by E. N. Zalta. Winter 2017. Metaphysics Research Lab, Stanford University.
- Krause, P. et al. (1995). “A Logic Of Argumentation for Reasoning under Uncertainty”. In: *Comput. Intell.* DOI: <https://doi.org/10.1111/j.1467-8640.1995.tb00025.x>.
- Lewis, D. (1973). *Counterfactuals*. Harvard University Press.
- Liu, F. (2008). “Changing for the better: Preference dynamics and agent diversity”. PhD thesis. Inst. for Logic, Language and Computation, Universiteit van Amsterdam.
- (2011). *Reasoning about Preference Dynamics*. Dordrecht: Springer Netherlands. DOI: 10.1007/978-94-007-1344-4.
- Lomfeld, B. (2015). *Die Gründe des Vertrages: Eine Diskurstheorie der Vertragsrechte*. Tübingen: Mohr Siebeck.
- (2017). “Vor den Fällen: Methoden soziologischer Jurisprudenz”. In: *Die Fälle der Gesellschaft: Eine neue Praxis soziologischer Jurisprudenz*. Ed. by Lomfeld. Tübingen: Mohr Siebeck, pp. 1–16.
- (2019). “Grammatik der Rechtfertigung: Eine kritische Rekonstruktion der Rechts(fort)bildung”. In: *Kritische Justiz* 52.4.
- Maranhão, J. and G. Sartor (2019). “Value assessment and revision in legal interpretation”. In: *Proceedings of the Seventeenth International Conference on Artificial Intelligence and Law, ICAIL 2019, Montreal, QC, Canada, June 17-21, 2019*, pp. 219–223. DOI: 10.1145/3322640.3326709.
- McCarty, L. T. (1995). “An implementation of Eisner v. Macomber”. In: *Proceedings of the 5th International Conference on Artificial Intelligence and Law*, pp. 276–286.
- Merrill, T. W. and H. E. Smith (2017). *Property: Principles and Policies*. Foundation Press.
- Mitchell, B. (2007). *Eight Ways to Run the Country*. Westport: Praeger.
- Modgil, S. and H. Prakken (2018). “Abstract Rule-Based Argumentation”. In: *Handbook of Formal Argumentation*. Ed. by P. Baroni, D. Gabbay, M. Giacomin, and L. van der Torre. College Publications, pp. 287–364.
- Moor, J. (2009). “Four kinds of ethical robots”. In: *Philos. Now* 72, pp. 12–14.
- Neves, M. (2021). *Constitutionalism and the Paradox of Principles and Rules*. Oxford University Press.

- Prakken, H. (1997). *Logical Tools for Modelling Legal Argument*. Dordrecht: Springer.
- (2002). “An exercise in formalising teleological case-based reasoning”. In: *Artif. Intell. Law* 10.1-3, pp. 113–133.
- Prakken, H. and G. Sartor (2015). “Law and logic: A review from an argumentation perspective”. In: *Artif. Intell.* 227, pp. 214–225.
- Rawls, J. (1971). *A Theory of Justice*. Revised edition 1999. Harvard university press.
- Raz, J. (1972). “Legal Principles and the Limits of Law”. In: *Yale Law J.* 81, pp. 823–854.
- Rissland, E. L. and K. D. Ashley (1987). “A case-based system for trade secrets law”. In: *Proceedings of the 1st international conference on Artificial Intelligence and Law*, pp. 60–66.
- Rokeach, M. (1973). *The Nature of Human Values*. New York: Free Press Macmillan.
- Sartor, G. (2010). “Doing justice to rights and values: teleological reasoning and proportionality”. In: *Artif. Intell. Law* 18, pp. 175–215.
- (2018). “A Quantitative Approach to Proportionality”. In: *Handbook of Legal Reasoning and Argumentation*. Ed. by B. et al. Dordrecht: Springer, pp. 613–636.
- Scheutz, M. (2017). “The Case for Explicit Ethical Agents”. In: *AI Mag.* 38.4, pp. 57–64.
- Schönfinkel, M. (1924). “Über die Bausteine der mathematischen Logik”. In: *Math. Ann.* 92, pp. 305–316.
- Schwartz, S. (1992). “Universals in the Content and Structure of Values”. In: *Adv. Exp. Soc. Psychol.* 25, pp. 1–65.
- Sieckmann, J.-R., ed. (2010). *Legal Reasoning: The Methods of Balancing*. Vol. 124. ARSP Beiheft. Stuttgart: Franz Steiner.
- Smith, B. (2003). “Ontology”. In: *Blackwell Guide to the Philosophy of Computing and Information*. Ed. by L. Floridi. Oxford: Blackwell.
- Svenningsson, J. and E. Axelsson (2013). “Combining Deep and Shallow Embedding for EDSL”. In: *Trends in Functional Programming*. Ed. by H.-W. Loidl and R. Peña. Berlin, Heidelberg: Springer Berlin Heidelberg, pp. 21–36. ISBN: 978-3-642-40447-4.
- Teubner, G. (1983). “Substantive and Reflexive Elements in Modern Law”. In: *Law & Soc. Rev.* 17, pp. 239–285.
- van Benthem, J. (2009). “For Better or for Worse: Dynamic Logics of Preference”. In: *Preference Change: Approaches from Philosophy, Economics and Psychology*. Ed. by T. Grüne-Yanoff and S. O. Hansson. Dordrecht: Springer Netherlands, pp. 57–84. DOI: 10.1007/978-90-481-2593-7\_3.
- van Benthem, J., P. Girard, and O. Roy (2009). “Everything Else Being Equal: A Modal Logic for *Ceteris Paribus* Preferences”. In: *J. Philos. Log.* 38.1, pp. 83–125. DOI: 10.1007/s10992-008-9085-3.
- Verheij, B., J. C. Hage, and H. J. Van Den Herik (1998). “An integrated view on rules and principles”. In: *Artif. Intell. Law* 6.1, pp. 3–26.
- von Wright, G. H. (1963). *The logic of preference*. Edinburgh University Press.
- Weide, T. van der et al. (2010). “Practical Reasoning Using Values”. In: *Argumentation in Multi-Agent Systems (ArgMAS)*. Ed. by P. McBurney, I. Rahwan, S. Parsons, and N. Maudet. Berlin: Springer, pp. 79–93.

---

Wenzel, M. (2007). “Isabelle/Isar—a generic framework for human-readable proof documents”. In: *From Insight to Proof-Festschrift Honour Andrzej Trybulec* 10.23, pp. 277–298.

## A Appendix - Isabelle/HOL Encoding

### A.1 SSE of $\mathcal{PL}$ in HOL

We comment on the implementation of the SSE of  $\mathcal{PL}$  in *Isabelle/HOL* as displayed in Figs. 8-9; see van van Benthem et al. (2009) for further details on  $\mathcal{PL}$ . The defined theory is named "PreferenceLogicBasics" and it relies on base logic HOL, imported here as theory "Main".

First, a new base type  $\iota$  is declared (Line 6), denoting the set of possible worlds or states. Subsequently (Lines 7–11), useful type abbreviations are introduced, including the type  $\sigma$  for  $\mathcal{PL}$  propositions, which are modelled as predicates on objects of type  $\iota$  (i.e., as *truth-sets* of worlds/states). A *betterness relation*  $\leq$ , and its strict variant  $<$ , are introduced (Lines 13–14), with  $\leq$ -accessible worlds interpreted as those that are *at least as good* as the present one. Definitions for relation properties are provided, and it is postulated that  $\leq$  is a preorder, i.e., reflexive and transitive (Lines 15–18).

Subsequently, the  $\sigma$ -type lifted logical connectives of  $\mathcal{PL}$  are introduced as abbreviations of  $\lambda$ -terms in the meta-logic HOL (Lines 21–33). The operators  $\Box^{\leq}$  and  $\Box^{<}$  use  $\leq$  and  $<$  as guards in their definitions (Lines 28 and 30); analogous for  $\Diamond^{\leq}$  and  $\Diamond^{<}$ . An *universal* modality and its dual are also introduced (Lines 32–33). Moreover, a notion of (global) truth for  $\mathcal{PL}$  formulas  $\psi$  is defined (Line 35): proposition  $\psi$  is globally true, we also say ‘valid’, if and only if it is true in all worlds.

As a first test some expected dualities of the modal operators are automatically proved (Line 36).

Subsequently, the *betterness* ordering  $\leq$  (resp.  $<$ ) is lifted to a preference relation between  $\mathcal{PL}$  propositions (sets of worlds). Eight possible semantic definitions for such preferences are encoded in HOL (Lines 40–47 in Fig. 9). The semantic definitions are complemented by eight syntactic definitions of the same

```

1 theory PreferenceLogicBasics imports Main (** Benzmüller & Fuenmayor, 2021 **)
2 begin (*unimportant*) declare[[syntax_ambiguity_warning=false]]
3 (*unimportant*) nitpick_params[user_axioms,expect=genuine,show_all,format=3]
4 (** SSE of preference logic by van Benthem et al., JPL 2009 **)
5 (*preliminaries*)
6 typedecl  $\iota$  (*possible worlds*)
7 type_synonym  $\sigma = \iota \Rightarrow \text{bool}$  (*world-lifted' propositions*)
8 type_synonym  $\gamma = \iota \Rightarrow \iota \Rightarrow \text{bool}$  (*preference relations*)
9 type_synonym  $\mu = \sigma \Rightarrow \sigma$  (*unary logical connectives*)
10 type_synonym  $\nu = \sigma \Rightarrow \sigma \Rightarrow \sigma$  (*binary logical connectives*)
11 type_synonym  $\pi = \sigma \Rightarrow \text{bool}$  (*sets of world-lifted propositions*)
12 (*betterness relation  $\leq$  and strict betterness relation  $<$ *)
13 consts BR:  $\gamma$  ("<math>\leq</math>")
14 definition SBR:  $\gamma$  ("<math><</math>") where "v < w  $\equiv (v \leq w) \wedge \neg (w \leq v)$ "
15 abbreviation "reflexive R  $\equiv \forall x. R\ x\ x$ "
16 abbreviation "transitive R  $\equiv \forall x\ y\ z. R\ x\ y \wedge R\ y\ z \longrightarrow R\ x\ z$ "
17 abbreviation "is_total R  $\equiv \forall x\ y. R\ x\ y \vee R\ y\ x$ "
18 axiomatization where rBR: "reflexive BR" and tBR: "transitive BR"
19 lemma tSBR: "transitive SBR" using SBR_def tBR by blast (*derived from axioms*)
20 (*modal logic connectives (operating on truth-sets)*)
21 abbreviation c1:  $\sigma$  ("<math>\perp</math>") where " $\perp \equiv \lambda w. \text{False}$ "
22 abbreviation c2:  $\sigma$  ("<math>\top</math>") where " $\top \equiv \lambda w. \text{True}$ "
23 abbreviation c3:  $\mu$  ("<math>\neg</math>") where " $\neg \varphi \equiv \lambda w. \neg (\varphi\ w)$ "
24 abbreviation c4:  $\nu$  ("<math>\wedge</math>") where " $\varphi \wedge \psi \equiv \lambda w. (\varphi\ w) \wedge (\psi\ w)$ "
25 abbreviation c5:  $\nu$  ("<math>\vee</math>") where " $\varphi \vee \psi \equiv \lambda w. (\varphi\ w) \vee (\psi\ w)$ "
26 abbreviation c6:  $\nu$  ("<math>\longrightarrow</math>") where " $\varphi \longrightarrow \psi \equiv \lambda w. (\varphi\ w) \longrightarrow (\psi\ w)$ "
27 abbreviation c7:  $\nu$  ("<math>\longleftrightarrow</math>") where " $\varphi \longleftrightarrow \psi \equiv \lambda w. (\varphi\ w) \longleftrightarrow (\psi\ w)$ "
28 abbreviation c8:  $\mu$  ("<math>\Box^{\leq}</math>") where " $\Box^{\leq} \varphi \equiv \lambda w. \forall v. (w \leq v) \longrightarrow (\varphi\ v)$ "
29 abbreviation c9:  $\mu$  ("<math>\Box^{<</math>") where " $\Box^{<} \varphi \equiv \lambda w. \exists v. (w < v) \wedge (\varphi\ v)$ "
30 abbreviation c10:  $\mu$  ("<math>\Diamond^{\leq}</math>") where " $\Diamond^{\leq} \varphi \equiv \lambda w. \forall v. (w < v) \longrightarrow (\varphi\ v)$ "
31 abbreviation c11:  $\mu$  ("<math>\Diamond^{<</math>") where " $\Diamond^{<} \varphi \equiv \lambda w. \exists v. (w < v) \wedge (\varphi\ v)$ "
32 abbreviation c12:  $\mu$  ("<math>E</math>") where " $E \varphi \equiv \lambda w. \exists v. (\varphi\ v)$ "
33 abbreviation c13:  $\mu$  ("<math>A</math>") where " $A \varphi \equiv \lambda w. \forall v. (\varphi\ v)$ "
34 (*meta-logical predicate for global and validity*)
35 abbreviation g1:  $\pi$  ("<math>|\psi</math>") where " $|\psi \equiv \forall w. \psi\ w$ "
36 (*some tests: dualities*)
37 lemma "[ $\Box^{\leq} \varphi \leftrightarrow (\neg \Box^{<} \neg \varphi)$ ]  $\wedge$  [ $\Diamond^{<} \varphi \leftrightarrow (\neg \Box^{\leq} \neg \varphi)$ ]  $\wedge$  [ $A \varphi \leftrightarrow (\neg E \neg \varphi)$ ]" by blast (*proof*)

```

Fig. 8 SSE of  $\mathcal{PL}$  (van Benthem et al. 2009) in HOL (continued in Fig. 9)

```

38 (**** Section 3: A basic modal preference language ****)
39 (*Definition 5*)
40 abbreviation pEE:: $\nu$  ("<math>\prec_{EE}</math>") where "<math>\varphi \prec_{EE} \psi \equiv \exists s. \varphi s \wedge (\exists t. \psi t \wedge s \prec t)</math>"
41 abbreviation pEES:: $\nu$  ("<math>\prec_{EE}</math>") where "<math>(\varphi \prec_{EE} \psi) \equiv \exists s. \varphi s \wedge (\exists t. \psi t \wedge s \prec t)</math>"
42 abbreviation pEA:: $\nu$  ("<math>\prec_{EA}</math>") where "<math>\varphi \prec_{EA} \psi \equiv \exists t. \psi t \wedge (\forall s. \varphi s \longrightarrow s \prec t)</math>"
43 abbreviation pEAs:: $\nu$  ("<math>\prec_{EA}</math>") where "<math>(\varphi \prec_{EA} \psi) \equiv \exists t. \psi t \wedge (\forall s. \varphi s \longrightarrow s \prec t)</math>"
44 abbreviation pAE:: $\nu$  ("<math>\prec_{AE}</math>") where "<math>\varphi \prec_{AE} \psi \equiv \forall s. \varphi s \longrightarrow (\exists t. \psi t \wedge s \prec t)</math>"
45 abbreviation pAEs:: $\nu$  ("<math>\prec_{AE}</math>") where "<math>(\varphi \prec_{AE} \psi) \equiv \forall s. \varphi s \longrightarrow (\exists t. \psi t \wedge s \prec t)</math>"
46 abbreviation pAA:: $\nu$  ("<math>\prec_{AA}</math>") where "<math>\varphi \prec_{AA} \psi \equiv \forall s. \varphi s \longrightarrow (\forall t. \psi t \longrightarrow s \prec t)</math>"
47 abbreviation pAAs:: $\nu$  ("<math>\prec_{AA}</math>") where "<math>(\varphi \prec_{AA} \psi) \equiv \forall s. \varphi s \longrightarrow (\forall t. \psi t \longrightarrow s \prec t)</math>"
48 abbreviation pEE:: $\nu$  ("<math>\prec_{EE}</math>") where "<math>\varphi \prec_{EE} \psi \equiv E(\varphi \wedge \Diamond \psi)</math>"
49 abbreviation pEES:: $\nu$  ("<math>\prec_{EE}</math>") where "<math>\varphi \prec_{EE} \psi \equiv E(\varphi \wedge \Diamond \psi)</math>"
50 abbreviation pEA:: $\nu$  ("<math>\prec_{EA}</math>") where "<math>\varphi \prec_{EA} \psi \equiv E(\psi \wedge \Box \neg \varphi)</math>"
51 abbreviation pEAs:: $\nu$  ("<math>\prec_{EA}</math>") where "<math>\varphi \prec_{EA} \psi \equiv E(\psi \wedge \Box \neg \varphi)</math>"
52 abbreviation pAE:: $\nu$  ("<math>\prec_{AE}</math>") where "<math>\varphi \prec_{AE} \psi \equiv A(\varphi \longrightarrow \Diamond \psi)</math>"
53 abbreviation pAEs:: $\nu$  ("<math>\prec_{AE}</math>") where "<math>\varphi \prec_{AE} \psi \equiv A(\varphi \longrightarrow \Diamond \psi)</math>"
54 abbreviation pAA:: $\nu$  ("<math>\prec_{AA}</math>") where "<math>\varphi \prec_{AA} \psi \equiv A(\psi \longrightarrow \Box \neg \varphi)</math>"
55 abbreviation pAAs:: $\nu$  ("<math>\prec_{AA}</math>") where "<math>\varphi \prec_{AA} \psi \equiv A(\psi \longrightarrow \Box \neg \varphi)</math>"
56 (*quantification for objects of arbitrary type*)
57 abbreviation mforall ("<math>\forall</math>") where "<math>\forall \Phi \equiv \lambda w. \forall x. (\Phi x w)</math>"
58 abbreviation mforallB (binder"<math>\forall</math>"[8]9) where "<math>\forall x. \varphi(x) \equiv \forall \varphi.</math>"
59 abbreviation mexists ("<math>\exists</math>") where "<math>\exists \Phi \equiv \lambda w. \exists x. (\Phi x w)</math>"
60 abbreviation mexistsB (binder"<math>\exists</math>"[8]9) where "<math>\exists x. \varphi(x) \equiv \exists \varphi.</math>"
61 (*polymorphic operators for sets of worlds/values*)
62 abbreviation subs (infix "<math>\sqsubseteq</math>" 70) where "<math>A \sqsubseteq B \equiv \forall x. A x \longrightarrow B x</math>"
63 abbreviation union (infix "<math>\sqcup</math>" 70) where "<math>A \sqcup B \equiv \lambda x. A x \vee B x</math>"
64 abbreviation inters (infix "<math>\sqcap</math>" 70) where "<math>A \sqcap B \equiv \lambda x. A x \wedge B x</math>"
65 (*consistency confirmed (trivial: only abbreviations introduced)*)
66 lemma True nitpick[satisfy,user_axioms] oops (*satisfying model*)
67 end

```

Fig. 9 SSE of  $\mathcal{PL}$  (van Benthem et al. 2009) in HOL (continued from Fig. 8)

binary preferences stated within the object language  $\mathcal{PL}$  (Lines 48–55). (ATP systems prove the meta-theoretic correspondences between these semantic and syntactic definitions; cf. Lines 4–12 in Fig. 11.)

$\mathcal{PL}$  is extended by adding quantifiers (Lines 57–60); cf. (Benzmüller and Paulson 2013) for explanations on the SSE of quantified modal logics. Moreover, useful polymorphic operators for subset, union and intersection are defined (Lines 62–64).

The model finder *Nitpick* (Blanchette and Nipkow 2010) confirms the consistency of the introduced theory (Line 66) by generating and presenting a model for it (not shown here).

To gain practical evidence for the faithfulness of our SSE of  $\mathcal{PL}$  in *Isabelle/HOL*, and also to assess proof automation performance, we have conducted numerous experiments in which we automatically reconstruct meta-theoretical results on  $\mathcal{PL}$ ; see Figs. 11–12.

Extending our SSE of  $\mathcal{PL}$  in HOL some further preference relations for  $\mathcal{PL}$  are defined in Fig. 10. These additional relations support *ceteris paribus* reasoning in  $\mathcal{PL}$ . We give some explanations:

Lines 5–13 Useful set theoretic notions are introduced as abbreviations for corresponding  $\lambda$ -terms in HOL.

Lines 14–22  $\mathcal{PL}$  is further extended with (equality-based) *ceteris paribus* preference relations and modalities; here  $\Gamma$  represents a set of formulas that are assumed constant between two possible worlds to compare. Hence our variant can be understood as “these (given) things being equal”-preferences. This variant can be used for modelling von Wright’s notion of *ceteris paribus* (“all other things being equal”) preferences, eliciting an appropriate  $\Gamma$  by extra-logical means.

Lines 26–33: Except for  $\prec_{AA}^{\Gamma}$ , the remaining operators we define here were not explicitly defined by van Benthem et al. (2009); however, their existence is tacitly suggested.

Meta-theoretical results on  $\mathcal{PL}$  as presented by van Benthem et al. (2009) are automatically verified by the reasoning tools in *Isabelle/HOL*; see Figs. 11–14; we in fact prove all relevant results from (van Benthem et al. 2009). The experiments shown in Fig. 11 are briefly commented:

Lines 5–13 Correspondences between the semantically and syntactically defined preference relations are proved.

Lines 15–22 It is proved that (e.g. inclusion and interaction) axioms for  $\mathcal{PL}$  follow as theorems in our SSE. This tests the faithfulness of the embedding in one direction.

Lines 25–47 We continue the mechanical verification of theorems, and generate countermodels (not displayed here) for non-theorems of  $\mathcal{PL}$ , thus putting our encoding to the test. Our results coincide with

```

1 theory PreferenceLogicCeterisParibus (** Benzmüller & Fuenmayor, 2021 **)
2   imports PreferenceLogicBasics
3 begin (** Ceteris Paribus reasoning by van Benthem et al, JPL 2009 **)
4 (*Section 5: Equality-based Ceteris Paribus Preference Logic*)
5 abbreviation a1: " $\sigma \Rightarrow \pi \Rightarrow \text{bool}$ " ("  $\in$  ") where " $\varphi \in \Gamma \equiv \Gamma \varphi$ "
6 abbreviation a2 ("  $\subseteq$  ") where " $\Gamma \subseteq \Gamma' \equiv \forall \varphi. \varphi \in \Gamma \longrightarrow \varphi \in \Gamma'$ "
7 abbreviation a3 ("  $\cup$  ") where " $\Gamma \cup \Gamma' \equiv \lambda \varphi. \varphi \in \Gamma \vee \varphi \in \Gamma'$ "
8 abbreviation a4 ("  $\cap$  ") where " $\Gamma \cap \Gamma' \equiv \lambda \varphi. \varphi \in \Gamma \wedge \varphi \in \Gamma'$ "
9 abbreviation a5 (" { } ") where " $\{\varphi\} \equiv \lambda x::\sigma. x = \varphi$ "
10 abbreviation a6 (" {_,_} ") where " $\{\alpha, \beta\} \equiv \lambda x::\sigma. x = \alpha \vee x = \beta$ "
11 abbreviation a7 (" {_,_,_} ") where " $\{\alpha, \beta, \gamma\} \equiv \lambda x::\sigma. x = \alpha \vee x = \beta \vee x = \gamma$ "
12 abbreviation a8 ("  $\emptyset$  ") where " $\emptyset \equiv (\lambda \psi::\sigma. \text{False})$ "
13 abbreviation a9 ("  $\mathcal{U}$  ") where " $\mathcal{U} \equiv (\lambda \psi::\sigma. \text{True})$ "
14 abbreviation c14 ("  $\equiv$  ") where " $w \equiv v \equiv \forall \varphi. \varphi \in \Gamma \longrightarrow (\varphi w \longleftrightarrow \varphi v)$ "
15 abbreviation c15 ("  $\trianglelefteq$  ") where " $w \trianglelefteq v \equiv w \leq v \wedge w \equiv v$ "
16 abbreviation c16 ("  $\triangleleft$  ") where " $w \triangleleft v \equiv w < v \wedge w \equiv v$ "
17 abbreviation c17 ("  $\trianglelefteq_{\Gamma}$  ") where " $(\Gamma) \trianglelefteq_{\Gamma} \varphi \equiv \lambda w. \exists v. w \trianglelefteq_{\Gamma} v \wedge \varphi v$ "
18 abbreviation c18 ("  $\triangleleft_{\Gamma}$  ") where " $(\Gamma) \triangleleft_{\Gamma} \varphi \equiv \lambda w. \forall v. w \triangleleft_{\Gamma} v \longrightarrow \varphi v$ "
19 abbreviation c19 ("  $\trianglelefteq_{\Gamma}$  ") where " $(\Gamma) \trianglelefteq_{\Gamma} \varphi \equiv \lambda w. \exists v. w \triangleleft_{\Gamma} v \wedge \varphi v$ "
20 abbreviation c20 ("  $\triangleleft_{\Gamma}$  ") where " $(\Gamma) \triangleleft_{\Gamma} \varphi \equiv \lambda w. \forall v. w \triangleleft_{\Gamma} v \longrightarrow \varphi v$ "
21 abbreviation c21 ("  $\trianglelefteq$  ") where " $(\Gamma) \trianglelefteq \varphi \equiv \lambda w. \exists v. w \equiv v \wedge \varphi v$ "
22 abbreviation c22 ("  $\triangleleft$  ") where " $(\Gamma) \triangleleft \varphi \equiv \lambda w. \forall v. w \equiv v \longrightarrow \varphi v$ "
23 (*Section 6: Ceteris Paribus Counterparts of Binary Pref. Statements*)
24 (*operators below not defined in paper; existence is tacitly suggested.
25   AA-variant draws upon von Wright's. AE-variant draws upon Halpern's.*)
26 abbreviation c23 ("  $\triangleleft_{AA}$  ") where " $(\varphi \triangleleft_{AA} \psi) u \equiv \forall s. \forall t. \varphi s \wedge \psi t \longrightarrow s \triangleleft_{\Gamma} t$ "
27 abbreviation c24 ("  $\trianglelefteq_{AA}$  ") where " $(\varphi \trianglelefteq_{AA} \psi) u \equiv \forall s. \forall t. \varphi s \wedge \psi t \longrightarrow s \trianglelefteq_{\Gamma} t$ "
28 abbreviation c25 ("  $\triangleleft_{AE}$  ") where " $(\varphi \triangleleft_{AE} \psi) u \equiv \forall s. \exists t. \varphi s \longrightarrow \psi t \wedge s \triangleleft_{\Gamma} t$ "
29 abbreviation c26 ("  $\trianglelefteq_{AE}$  ") where " $(\varphi \trianglelefteq_{AE} \psi) u \equiv \forall s. \exists t. \varphi s \longrightarrow \psi t \wedge s \trianglelefteq_{\Gamma} t$ "
30 abbreviation c27 ("  $\triangleleft_{AA}^{\Gamma}$  ") where " $(\varphi \triangleleft_{AA}^{\Gamma} \psi) \equiv A(\psi \rightarrow (\Gamma) \trianglelefteq_{\Gamma} \varphi)$ "
31 abbreviation c28 ("  $\trianglelefteq_{AA}^{\Gamma}$  ") where " $(\varphi \trianglelefteq_{AA}^{\Gamma} \psi) \equiv A(\psi \rightarrow (\Gamma) \trianglelefteq_{\Gamma} \varphi)$ "
32 abbreviation c29 ("  $\triangleleft_{AE}^{\Gamma}$  ") where " $(\varphi \triangleleft_{AE}^{\Gamma} \psi) \equiv A(\varphi \rightarrow (\Gamma) \triangleleft_{\Gamma} \psi)$ "
33 abbreviation c30 ("  $\trianglelefteq_{AE}^{\Gamma}$  ") where " $(\varphi \trianglelefteq_{AE}^{\Gamma} \psi) \equiv A(\varphi \rightarrow (\Gamma) \trianglelefteq_{\Gamma} \psi)$ "
34 (*Consistency confirmed (trivial: only abbreviations are introduced*)
35 lemma True nitpick[satisfy,user_axioms] oops
36 end

```

Fig. 10 SSE of  $\mathcal{PL}$  (van Benthem et al. 2009) in HOL (continued from Figs. 8-9)

the corresponding ones claimed (and in many cases proved) in van Benthem et al. (2009), except for the claims encoded in lines 40-41 and 44-45, where countermodels are reported by *Nitpick*. Lines 25-47 Some application-specific tests in preparation for the modelling of the legal DSL (including the value theory/ontology) are conducted.

```

1 theory PreferenceLogicTests1 imports PreferenceLogicBasics (** Benzm. & Fuenmayor, 2021 **)
2 begin (** Tests for the SSE of van Benthem et al, JPL 2009, in HOL **)
3 (*Fact 1: definability of the principal operators and verification*)
4 lemma F1_9: "[(\varphi \preceq_{EE} \psi) \leftrightarrow (\varphi \preceq_{EE} \psi)]" by blast
5 lemma F1_10: "[(\varphi \preceq_{AE} \psi) \leftrightarrow (\varphi \preceq_{AE} \psi)]" by blast
6 lemma F1_11: "[(\varphi \prec_{EE} \psi) \leftrightarrow (\varphi \prec_{EE} \psi)]" by blast
7 lemma F1_12: "[(\varphi \prec_{AE} \psi) \leftrightarrow (\varphi \prec_{AE} \psi)]" by blast
8 (*Fact 2: definability of remaining pref. operators and verification*)
9 lemma F2_13: "is_total BR \longrightarrow [(\varphi \prec_{AA} \psi) \leftrightarrow (\varphi \prec_{AA} \psi)]" using SBR_def by blast
10 lemma F2_14: "is_total BR \longrightarrow [(\varphi \prec_{EA} \psi) \leftrightarrow (\varphi \prec_{EA} \psi)]" using SBR_def by blast
11 lemma F2_15: "is_total BR \longrightarrow [(\varphi \preceq_{AA} \psi) \leftrightarrow (\varphi \preceq_{AA} \psi)]" using SBR_def by blast
12 lemma F2_16: "is_total BR \longrightarrow [(\varphi \preceq_{EA} \psi) \leftrightarrow (\varphi \preceq_{EA} \psi)]" using SBR_def by blast
13 (*Section 3.5 "Axiomatization" -- verify interaction axioms*)
14 lemma Incl_1: "[(\Diamond \varphi) \rightarrow (\Diamond \varphi)]" using SBR_def by blast
15 lemma Inter_1: "[(\Diamond \Diamond \varphi) \rightarrow (\Diamond \varphi)]" using tBR SBR_def by metis
16 lemma Trans_le: "[(\Diamond \Diamond \varphi) \rightarrow (\Diamond \varphi)]" using tSBR by blast
17 lemma Inter_2: "[(\varphi \wedge \Diamond \psi) \rightarrow ((\Diamond \psi) \vee \Diamond (\psi \wedge \Diamond \varphi))]" using SBR_def by blast
18 lemma F4: "[(\varphi \wedge \Diamond \psi) \rightarrow ((\Diamond \psi) \vee \Diamond (\psi \wedge \Diamond \varphi))] \longleftrightarrow
19 (\forall w. \forall v. ((w \preceq v) \wedge \neg (v \preceq w)) \longrightarrow (w \prec v))" using SBR_def by blast
20 lemma Inter_3: "[(\Diamond \Diamond \varphi) \rightarrow (\Diamond \varphi)]" using tBR SBR_def by blast
21 lemma Incl_2: "[(\Diamond \varphi) \rightarrow (E\varphi)]" by blast
22 (*Section 3.6 "A binary preference fragment"*)
23 (* \preceq_{EE} is the dual of \prec_{AA} *)
24 lemma "[(\varphi \preceq_{EE} \psi) \leftrightarrow \neg(\psi \prec_{AA} \varphi)] \wedge [(\varphi \prec_{AA} \psi) \leftrightarrow \neg(\psi \preceq_{EE} \varphi)]" by blast
25 (* \preceq_{EE} is the dual of \prec_{AA} only if totality is assumed*)
26 lemma "[(\varphi \preceq_{EE} \psi) \leftrightarrow \neg(\psi \prec_{AA} \varphi)]" nitpick oops (*countermodel*)
27 lemma "[(\varphi \preceq_{EE} \psi) \rightarrow \neg(\psi \prec_{AA} \varphi)]" using SBR_def by blast (*this direction holds*)
28 lemma "is_total BR \longrightarrow [(\varphi \preceq_{EE} \psi) \leftrightarrow \neg(\psi \prec_{AA} \varphi)]" using SBR_def by blast
29 lemma "[(\varphi \prec_{AA} \psi) \leftrightarrow \neg(\psi \preceq_{EE} \varphi)]" nitpick oops (*countermodel*)
30 lemma "[(\varphi \prec_{AA} \psi) \rightarrow \neg(\psi \preceq_{EE} \varphi)]" using SBR_def by blast (*this direction holds*)
31 lemma "is_total BR \longrightarrow [(\varphi \prec_{AA} \psi) \leftrightarrow \neg(\psi \preceq_{EE} \varphi)]" using SBR_def by blast
32 (* verify p.97-98 *)
33 lemma monotonicity: "[((\varphi \preceq_{EE} \psi) \wedge A(\varphi \rightarrow \xi)) \rightarrow (\xi \preceq_{EE} \psi)]" by blast
34 lemma reducibility: "[(((\varphi \preceq_{EE} \psi) \wedge \alpha) \preceq_{EE} \beta) \leftrightarrow ((\varphi \preceq_{EE} \psi) \wedge (\alpha \preceq_{EE} \beta))]" by blast
35 lemma reflexivity: "[\varphi \rightarrow (\varphi \preceq_{EE} \varphi)]" using rBR by blast
36 (*The condition below enforcing totality of the preference relation is supposed to hold.
37 However there are countermodels (both local & global consequence). Error in paper?*)
38 lemma "[(((\varphi \preceq_{EE} \varphi) \wedge (\psi \preceq_{EE} \psi)) \rightarrow ((\varphi \preceq_{EE} \psi) \vee (\psi \preceq_{EE} \varphi))]"
39 nitpick oops (*countermodel*)
40 lemma "[(((\varphi \preceq_{EE} \varphi) \wedge (\psi \preceq_{EE} \psi)) \longrightarrow (((\varphi \preceq_{EE} \psi) \vee (\psi \preceq_{EE} \varphi)))]"
41 nitpick oops (*countermodel*)
42 end

```

Fig. 11 Experiments: Testing the meta-theory of  $\mathcal{PL}$

```

1 theory PreferenceLogicTests2 (** Benzmüller & Fuenmayor, 2021 **)
2   imports PreferenceLogicCeterisParibus
3 begin (** Tests for the SSE of van Benthem et al, JPL 2009 **)
4 (** Section 5: Equality-based Ceteris Paribus Preference Logic **)
5 (*Some tests: dualities*)
6 lemma "[((Γ)⊆φ) ↔ ¬((Γ)⊆¬φ)]" by auto
7 lemma "[((Γ)⊆¬φ) ↔ ¬((Γ)⊆φ)]" by auto
8 lemma "[((Γ)⊆φ) ↔ ¬((Γ)⊆¬φ)]" by auto
9 (*Lemma 2*)
10 lemma lemma2_1 : "(⊆⊆φ) w ↔ ((∅)⊆φ) w" by auto
11 lemma lemma2_2 : "(⊆⊆¬φ) w ↔ ((∅)⊆¬φ) w" by auto
12 lemma lemma2_3 : "((⊆φ) w ↔ ((∅)⊆φ) w) ∧ ((⊆Aφ) w ↔ ((∅)⊆φ) w)" by auto
13 (**Axiomatization**)
14 (*inclusion and interaction axioms *)
15 lemma Incl1 : "[((Γ)⊆φ) → ((Γ)⊆φ)]" using SBR_def by auto
16 lemma Inc2 : "[((Γ)⊆φ) → ((Γ)⊆φ)]" by auto
17 lemma Int3 : "[((Γ)⊆((Γ)⊆φ)) → ((Γ)⊆φ)]" by (meson tBR )
18 lemma Int4 : "[((Γ)⊆((Γ)⊆φ)) → ((Γ)⊆φ)]" by (metis SBR_def tBR )
19 lemma Int5 : "[(ψ ∧ ((Γ)⊆φ)) → ((Γ)⊆φ) ∨ ((Γ)⊆(φ ∧ ((Γ)⊆φ)))]" by (metis rBR )
20 (*ceteris paribus reflexivity*)
21 lemma CetPar6 : "φ ∈ Γ → [(Γ)⊆φ] → φ" by blast
22 lemma CetPar7 : "φ ∈ Γ → [(Γ)⊆¬φ] → ¬φ" by blast
23 (*monotonicity*)
24 lemma CetPar8 : "Γ ⊆ Γ' → [(Γ')⊆φ] → ((Γ)⊆φ)" by auto
25 lemma CetPar9 : "Γ ⊆ Γ' → [(Γ')⊆¬φ] → ((Γ)⊆¬φ)" by auto
26 lemma CetPar10 : "Γ ⊆ Γ' → [(Γ')⊆φ] → ((Γ)⊆φ)" by auto
27 (*increase (decrease) of ceteris paribus sets*)
28 lemma CetPar11a : "[((φ ∧ ((Γ)(α ∧ φ))) → ((ΓU{φ})α)]" by auto
29 lemma CetPar11b : "[((¬φ) ∧ ((Γ)(α ∧ ¬φ))) → ((ΓU{φ})α)]" by auto
30 lemma CetPar12a : "[((φ ∧ ((Γ)⊆(α ∧ φ))) → ((ΓU{φ})⊆α)]" by auto
31 lemma CetPar12b : "[((¬φ) ∧ ((Γ)⊆(α ∧ ¬φ))) → ((ΓU{φ})⊆α)]" by auto
32 lemma CetPar13a : "[((φ ∧ ((Γ)⊆(α ∧ φ))) → ((ΓU{φ})⊆α)]" by auto
33 lemma CetPar13b : "[((¬φ) ∧ ((Γ)⊆(α ∧ ¬φ))) → ((ΓU{φ})⊆α)]" by auto
34 (*Example 1, Lemma 4, Corollary 1 and Lemmas5*)
35 lemma Ex1 : "[(((Γ)⊆φ) ∧ ((Γ)⊆α)) → ((ΓU{φ})⊆α)]" using rBR by auto
36 lemma Lemma4 : "[((Γ)⊆φ) w → (∃v. (w ⊆r v) ∧ (φ v))]" by simp
37 lemma Cor1 : "[((Γ)⊆φ) w → (∃v. (w ≡r v) ∧ (φ v))]" by simp
38 lemma Lemma5 : "(w ⊆r v) ↔ ((w ⊆ v) ∧ (w ≡r v))" by auto
39 (**** Section 6: Ceteris Paribus Counterparts ****)
40 (*AA-variant (drawing upon von Wright's)*)
41 lemma "(φ ⊆AAΓ ψ) u ↔ (φ ⊆AAΓ ψ) u" nitpick oops (*Ctm*)
42 lemma "(φ ⊆AAΓ ψ) u → (φ ⊆AAΓ ψ) u" nitpick oops (*Ctm*)
43 lemma "(φ ⊆AAΓ ψ) u → (φ ⊆AAΓ ψ) u" using SBR_def by auto
44 lemma "is_total SBR → (φ ⊆AAΓ ψ) u ↔ (φ ⊆AAΓ ψ) u" by (smt SBR_def )
45 lemma "(φ ⊆AAΓ ψ) u ↔ (φ ⊆AAΓ ψ) u" nitpick oops (*Ctm*)
46 lemma "(φ ⊆AAΓ ψ) u → (φ ⊆AAΓ ψ) u" nitpick oops (*Ctm*)
47 lemma "(φ ⊆AAΓ ψ) u → (φ ⊆AAΓ ψ) u" using SBR_def by auto
48 lemma "is_total SBR → (φ ⊆AAΓ ψ) u ↔ (φ ⊆AAΓ ψ) u" by (smt SBR_def )
49 (*AE-variant*)
50 lemma leAE_cp_pref : "(φ ⊆AEΓ ψ) u ↔ (φ ⊆AEΓ ψ) u" by auto
51 lemma leqAE_cp_pref : "(φ ⊆AEΓ ψ) u ↔ (φ ⊆AEΓ ψ) u" by auto
52 (*miscellaneous tests*)
53 lemma "let Γ = ∅ in [(φ ⊆AAΓ ψ) ↔ (φ ⊆AA ψ)]" by simp
54 lemma "let Γ = {⊥} in [(φ ⊆AAΓ ψ) ↔ (φ ⊆AA ψ)]" by simp
55 lemma "let Γ = {⊥, A} in [(φ ⊆AAΓ ψ) ↔ (φ ⊆AA ψ)]" nitpick oops (*Ctm*)
56 lemma "let Γ = {A} in [(φ ⊆AAΓ ψ) → (A → (φ ⊆AA ψ))]" nitpick oops (*Ctm*)
57 lemma "let Γ = {A} in [(A → (φ ⊆AA ψ)) → (φ ⊆AAΓ ψ)]" nitpick oops (*Ctm*)
58 lemma "let Γ = ∅ in [(φ ⊆AEΓ ψ) ↔ (φ ⊆AE ψ)]" by simp
59 lemma "let Γ = {⊥} in [(φ ⊆AEΓ ψ) ↔ (φ ⊆AE ψ)]" by simp
60 lemma "let Γ = {⊥, A} in [(φ ⊆AEΓ ψ) ↔ (φ ⊆AE ψ)]" nitpick oops (*Ctm*)
61 lemma "let Γ = {A} in [(φ ⊆AEΓ ψ) → (A → (φ ⊆AE ψ))]" by auto
62 lemma "let Γ = {A, B} in [(φ ⊆AEΓ ψ) → ((A ∧ B) → (φ ⊆AE ψ))]" by auto
63 lemma "let Γ = {A} in [(A → (φ ⊆AE ψ)) → (φ ⊆AEΓ ψ)]" nitpick oops (*Ctm*)
64 end

```

Fig. 12 Experiments (continued): Testing the meta-theory of  $\mathcal{PL}$



```

1 | theory PreferenceLogicTestsApp1 imports PreferenceLogicBasics (** Benzmüller & Fuenmayor, 2021 **)
2 | begin (****Application-specific tests for the value ontology****)
3 | (* EE variant ( $\wedge$ *)
4 | lemma "[A <EE (A $\wedge$ B)]" nitpick[satisfy] nitpick oops (*contingent*)
5 | lemma "[(A $\wedge$ B) <EE A]" nitpick[satisfy] nitpick oops (*contingent*)
6 | lemma "[A <EE B  $\rightarrow$  (A <EE (C $\wedge$ B))]" nitpick[satisfy] nitpick oops (*contingent*)
7 | lemma "[A <EE (C $\wedge$ B)]  $\rightarrow$  (A <EE B)" by blast
8 | lemma "[((C $\wedge$ B) <EE A)  $\rightarrow$  (B <EE A)]" by blast
9 | lemma "[B <EE A]  $\rightarrow$  ((C $\wedge$ B) <EE A)]" nitpick[satisfy] nitpick oops (*contingent*)
10 | (* EE variant ( $\vee$ *)
11 | lemma "[A <EE (A $\vee$ B)]" nitpick[satisfy] nitpick oops (*contingent*)
12 | lemma "[A $\vee$ B] <EE A]" nitpick[satisfy] nitpick oops (*contingent*)
13 | lemma "[A <EE B]  $\rightarrow$  (A <EE (C $\vee$ B))]" by blast
14 | lemma "[A <EE (C $\vee$ B)]  $\rightarrow$  (A <EE B)]" nitpick[satisfy] nitpick oops (*contingent*)
15 | lemma "[((C $\vee$ B) <EE A)  $\rightarrow$  (B <EE A)]" nitpick[satisfy] nitpick oops (*contingent*)
16 | lemma "[B <EE A]  $\rightarrow$  ((C $\vee$ B) <EE A)]" by blast
17 | (* AE variant ( $\wedge$ *)
18 | lemma "[A <AE (A $\wedge$ B)]" nitpick[satisfy] nitpick oops (*contingent*)
19 | lemma "[A $\wedge$ B] <AE A]" nitpick[satisfy] nitpick oops (*contingent*)
20 | lemma "[A <AE B]  $\rightarrow$  (A <AE (C $\wedge$ B))]" nitpick[satisfy] nitpick oops (*contingent*)
21 | lemma "[A <AE (C $\wedge$ B)]  $\rightarrow$  (A <AE B)]" by blast
22 | lemma "[((C $\wedge$ B) <AE A)  $\rightarrow$  (B <AE A)]" nitpick[satisfy] nitpick oops (*contingent*)
23 | lemma "[B <AE A]  $\rightarrow$  ((C $\wedge$ B) <AE A)]" by blast
24 | (* AE variant ( $\vee$ *)
25 | lemma "[A <AE (A $\vee$ B)]" nitpick[satisfy] nitpick oops (*contingent*)
26 | lemma "[A $\vee$ B] <AE A]" nitpick[satisfy] nitpick oops (*contingent*)
27 | lemma "[A <AE B]  $\rightarrow$  (A <AE (C $\vee$ B))]" by blast
28 | lemma "[A <AE (C $\vee$ B)]  $\rightarrow$  (A <AE B)]" nitpick[satisfy] nitpick oops (*contingent*)
29 | lemma "[((C $\vee$ B) <AE A)  $\rightarrow$  (B <AE A)]" by blast
30 | lemma "[B <AE A]  $\rightarrow$  ((C $\vee$ B) <AE A)]" nitpick[satisfy] nitpick oops (*contingent*)
31 | (* AA variant ( $\wedge$ *)
32 | lemma "[A <AA (A $\wedge$ B)]" nitpick[satisfy] nitpick oops (*contingent*)
33 | lemma "[A $\wedge$ B] <AA A]" nitpick[satisfy] nitpick oops (*contingent*)
34 | lemma "[A <AA B]  $\rightarrow$  (A <AA (C $\wedge$ B))]" by blast
35 | lemma "[A <AA (C $\wedge$ B)]  $\rightarrow$  (A <AA B)]" nitpick[satisfy] nitpick oops (*contingent*)
36 | lemma "[((C $\wedge$ B) <AA A)  $\rightarrow$  (B <AA A)]" nitpick[satisfy] nitpick oops (*contingent*)
37 | lemma "[B <AA A]  $\rightarrow$  ((C $\wedge$ B) <AA A)]" by blast
38 | (* AA variant ( $\vee$ *)
39 | lemma "[A <AA (A $\vee$ B)]" nitpick[satisfy] nitpick oops (*contingent*)
40 | lemma "[A $\vee$ B] <AA A]" nitpick[satisfy] nitpick oops (*contingent*)
41 | lemma "[A <AA B]  $\rightarrow$  (A <AA (C $\vee$ B))]" nitpick[satisfy] nitpick oops (*contingent*)
42 | lemma "[A <AA (C $\vee$ B)]  $\rightarrow$  (A <AA B)]" by blast
43 | lemma "[((C $\vee$ B) <AA A)  $\rightarrow$  (B <AA A)]" by blast
44 | lemma "[B <AA A]  $\rightarrow$  ((C $\vee$ B) <AA A)]" nitpick[satisfy] nitpick oops (*contingent*)
45 | (* EA variant ( $\wedge$ *)
46 | lemma "[A <EA (A $\wedge$ B)]" using rBR by blast
47 | lemma "[A $\wedge$ B] <EA A]" nitpick[satisfy] nitpick oops (*contingent*)
48 | lemma "[A <EA B]  $\rightarrow$  (A <EA (C $\wedge$ B))]" nitpick[satisfy] nitpick oops (*contingent*)
49 | lemma "[A <EA (C $\wedge$ B)]  $\rightarrow$  (A <EA B)]" by blast
50 | lemma "[((C $\wedge$ B) <EA A)  $\rightarrow$  (B <EA A)]" nitpick[satisfy] nitpick oops (*contingent*)
51 | lemma "[B <EA A]  $\rightarrow$  ((C $\wedge$ B) <EA A)]" by blast
52 | (* EA variant ( $\vee$ *)
53 | lemma "[A <EA (A $\vee$ B)]" nitpick[satisfy] nitpick oops (*contingent*)
54 | lemma "[A $\vee$ B] <EA A]" using rBR by blast
55 | lemma "[A <EA B]  $\rightarrow$  (A <EA (C $\vee$ B))]" by blast
56 | lemma "[A <EA (C $\vee$ B)]  $\rightarrow$  (A <EA B)]" nitpick[satisfy] nitpick oops (*contingent*)
57 | lemma "[((C $\vee$ B) <EA A)  $\rightarrow$  (B <EA A)]" by blast
58 | lemma "[B <EA A]  $\rightarrow$  ((C $\vee$ B) <EA A)]" nitpick[satisfy] nitpick oops (*contingent*)
59 | end

```

Fig. 13 Experiments (continued): Checking properties of strict preference relations

```

1 | theory PreferenceLogicTestsApp2 imports PreferenceLogicBasics (** Benzmüller & Fuenmayor, 2021 **)
2 | begin (****Application-specific tests for the value ontology****)
3 | (* EE variant ( $\wedge$ *)
4 | lemma "A  $\preceq_{EE}$  (A $\wedge$ B)]" nitpick[satisfy] nitpick oops (*contingent*)
5 | lemma "[A $\wedge$ B]  $\preceq_{EE}$  A]" nitpick[satisfy] nitpick oops (*contingent*)
6 | lemma "[A  $\preceq_{EE}$  B]  $\rightarrow$  [A  $\preceq_{EE}$  (C $\wedge$ B)]]" nitpick[satisfy] nitpick oops (*contingent*)
7 | lemma "[A  $\preceq_{EE}$  (C $\wedge$ B)]  $\rightarrow$  [A  $\preceq_{EE}$  B]" by blast
8 | lemma "[((C $\wedge$ B)  $\preceq_{EE}$  A)  $\rightarrow$  (B  $\preceq_{EE}$  A)]" by blast
9 | lemma "[B  $\preceq_{EE}$  A]  $\rightarrow$  ((C $\wedge$ B)  $\preceq_{EE}$  A)]" nitpick[satisfy] nitpick oops (*contingent*)
10 | (* EE variant ( $\vee$ *)
11 | lemma "A  $\preceq_{EE}$  (A $\vee$ B)]" nitpick[satisfy] nitpick oops (*contingent*)
12 | lemma "[A $\vee$ B]  $\preceq_{EE}$  A]" nitpick[satisfy] nitpick oops (*contingent*)
13 | lemma "[A  $\preceq_{EE}$  B]  $\rightarrow$  [A  $\preceq_{EE}$  (C $\vee$ B)]]" by blast
14 | lemma "[A  $\preceq_{EE}$  (C $\vee$ B)]  $\rightarrow$  [A  $\preceq_{EE}$  B]" nitpick[satisfy] nitpick oops (*contingent*)
15 | lemma "[((C $\vee$ B)  $\preceq_{EE}$  A)  $\rightarrow$  (B  $\preceq_{EE}$  A)]" nitpick[satisfy] nitpick oops (*contingent*)
16 | lemma "[B  $\preceq_{EE}$  A]  $\rightarrow$  ((C $\vee$ B)  $\preceq_{EE}$  A)]" by blast
17 | (* AE variant ( $\wedge$ *)
18 | lemma "A  $\preceq_{AE}$  (A $\wedge$ B)]" nitpick[satisfy] nitpick oops (*contingent*)
19 | lemma "[A $\wedge$ B]  $\preceq_{AE}$  A]" using rBR by blast (*change wrt. strict*)
20 | lemma "[A  $\preceq_{AE}$  B]  $\rightarrow$  [A  $\preceq_{AE}$  (C $\wedge$ B)]]" nitpick[satisfy] nitpick oops (*contingent*)
21 | lemma "[A  $\preceq_{AE}$  (C $\wedge$ B)]  $\rightarrow$  [A  $\preceq_{AE}$  B]" by blast
22 | lemma "[((C $\wedge$ B)  $\preceq_{AE}$  A)  $\rightarrow$  (B  $\preceq_{AE}$  A)]" nitpick[satisfy] nitpick oops (*contingent*)
23 | lemma "[B  $\preceq_{AE}$  A]  $\rightarrow$  ((C $\wedge$ B)  $\preceq_{AE}$  A)]" by blast
24 | (* AE variant ( $\vee$ *)
25 | lemma "A  $\preceq_{AE}$  (A $\vee$ B)]" using rBR by blast (*change wrt. strict*)
26 | lemma "[A $\vee$ B]  $\preceq_{AE}$  A]" nitpick[satisfy] nitpick oops (*contingent*)
27 | lemma "[A  $\preceq_{AE}$  B]  $\rightarrow$  [A  $\preceq_{AE}$  (C $\vee$ B)]]" by blast
28 | lemma "[A  $\preceq_{AE}$  (C $\vee$ B)]  $\rightarrow$  [A  $\preceq_{AE}$  B]" nitpick[satisfy] nitpick oops (*contingent*)
29 | lemma "[((C $\vee$ B)  $\preceq_{AE}$  A)  $\rightarrow$  (B  $\preceq_{AE}$  A)]" by blast
30 | lemma "[B  $\preceq_{AE}$  A]  $\rightarrow$  ((C $\vee$ B)  $\preceq_{AE}$  A)]" nitpick[satisfy] nitpick oops (*contingent*)
31 | (* AA variant ( $\wedge$ *)
32 | lemma "A  $\preceq_{AA}$  (A $\wedge$ B)]" nitpick[satisfy] nitpick oops (*contingent*)
33 | lemma "[A $\wedge$ B]  $\preceq_{AA}$  A]" nitpick[satisfy] nitpick oops (*contingent*)
34 | lemma "[A  $\preceq_{AA}$  B]  $\rightarrow$  [A  $\preceq_{AA}$  (C $\wedge$ B)]]" by blast
35 | lemma "[A  $\preceq_{AA}$  (C $\wedge$ B)]  $\rightarrow$  [A  $\preceq_{AA}$  B]" nitpick[satisfy] nitpick oops (*contingent*)
36 | lemma "[((C $\wedge$ B)  $\preceq_{AA}$  A)  $\rightarrow$  (B  $\preceq_{AA}$  A)]" nitpick[satisfy] nitpick oops (*contingent*)
37 | lemma "[B  $\preceq_{AA}$  A]  $\rightarrow$  ((C $\wedge$ B)  $\preceq_{AA}$  A)]" by blast
38 | (* AA variant ( $\vee$ *)
39 | lemma "A  $\preceq_{AA}$  (A $\vee$ B)]" nitpick[satisfy] nitpick oops (*contingent*)
40 | lemma "[A $\vee$ B]  $\preceq_{AA}$  A]" nitpick[satisfy] nitpick oops (*contingent*)
41 | lemma "[A  $\preceq_{AA}$  B]  $\rightarrow$  [A  $\preceq_{AA}$  (C $\vee$ B)]]" nitpick[satisfy] nitpick oops (*contingent*)
42 | lemma "[A  $\preceq_{AA}$  (C $\vee$ B)]  $\rightarrow$  [A  $\preceq_{AA}$  B]" by blast
43 | lemma "[((C $\vee$ B)  $\preceq_{AA}$  A)  $\rightarrow$  (B  $\preceq_{AA}$  A)]" by blast
44 | lemma "[B  $\preceq_{AA}$  A]  $\rightarrow$  ((C $\vee$ B)  $\preceq_{AA}$  A)]" nitpick[satisfy] nitpick oops (*contingent*)
45 | (* EA variant ( $\wedge$ *)
46 | lemma "A  $\preceq_{EA}$  (A $\wedge$ B)]" nitpick[satisfy] nitpick oops (*contingent*)
47 | lemma "[A $\wedge$ B]  $\preceq_{EA}$  A]" nitpick[satisfy] nitpick oops (*contingent*)
48 | lemma "[A  $\preceq_{EA}$  B]  $\rightarrow$  [A  $\preceq_{EA}$  (C $\wedge$ B)]]" nitpick[satisfy] nitpick oops (*contingent*)
49 | lemma "[A  $\preceq_{EA}$  (C $\wedge$ B)]  $\rightarrow$  [A  $\preceq_{EA}$  B]" by blast
50 | lemma "[((C $\wedge$ B)  $\preceq_{EA}$  A)  $\rightarrow$  (B  $\preceq_{EA}$  A)]" nitpick[satisfy] nitpick oops (*contingent*)
51 | lemma "[B  $\preceq_{EA}$  A]  $\rightarrow$  ((C $\wedge$ B)  $\preceq_{EA}$  A)]" by blast
52 | (* EA variant ( $\vee$ *)
53 | lemma "A  $\preceq_{EA}$  (A $\vee$ B)]" nitpick[satisfy] nitpick oops (*contingent*)
54 | lemma "[A $\vee$ B]  $\preceq_{EA}$  A]" nitpick[satisfy] nitpick oops (*contingent*)
55 | lemma "[A  $\preceq_{EA}$  B]  $\rightarrow$  [A  $\preceq_{EA}$  (C $\vee$ B)]]" by blast
56 | lemma "[A  $\preceq_{EA}$  (C $\vee$ B)]  $\rightarrow$  [A  $\preceq_{EA}$  B]" nitpick[satisfy] nitpick oops (*contingent*)
57 | lemma "[((C $\vee$ B)  $\preceq_{EA}$  A)  $\rightarrow$  (B  $\preceq_{EA}$  A)]" by blast
58 | lemma "[B  $\preceq_{EA}$  A]  $\rightarrow$  ((C $\vee$ B)  $\preceq_{EA}$  A)]" nitpick[satisfy] nitpick oops (*contingent*)
59 | end

```

Fig. 14 Experiments (continued): Checking properties of strict preference relations

## A.2 Encoding of the Legal DSL (Value Ontology)

The encoding of the legal DSL (value theory or ontology) is shown in Fig. 15. The new theory is termed “ValueOntology”, and it imports theory “PreferenceLogicBasics” (and thus recursively also *Isabelle/HOL*’s internal theory “Main”).

As a preliminary, the legal parties *plaintiff* and *defendant* are introduced as an (extensible) two-valued datatype together with a function to obtain for a given party the *other* one ( $x^{-1}$ ) (Lines 4–5); and a predicate modelling the ruling *for* a party is also provided (Lines 7–8).

As regards the *discursive grammar* value theory, a four-valued (parameterised) datatype is introduced (Line 10) as described in §2. Moreover, type-aliases (Lines 11–12) and set-constructor operators for values (Lines 14–15) are introduced for ease of presentation. The legal principles from §2 are introduced as combinations of those basic values (Lines 17–28). As an illustration, the principle STABility is encoded as a set composed of the basic values SECURITY and UTILITY.

Next, the incidence relation  $I$  and operators  $\uparrow$  and  $\downarrow$ , borrowed and adapted from formal concept analysis (FCA), are introduced (Lines 30–34).

We then define the aggregation operator  $\oplus$  as  $A \oplus B := (A \downarrow \vee B \downarrow)$ ; i.e., we select the second candidate as discussed in §2. And as our preference relation of choice we select the relation  $\prec_{AE}$  (Line 38).

```

1 theory ValueOntology imports PreferenceLogicBasics (** Benzml., Fuenmayor & Lomfeld, 2021 **)
2 begin (** Lomfeld's value ontology is encoded **)
3 (*new datatype for parties/contenders (there could be more in principle)*)
4 datatype c = p | d (*plaintiff & defendant*)
5 fun other::"c⇒c" ("_⁻¹") where "p⁻¹ = d" | "d⁻¹ = p"
6 (*new constant symbol: finding/ruling for party*)
7 consts For::"c⇒σ"
8 axiomatization where ForAx: "[For x ↔ (¬For x⁻¹)]"
9 (*new parameterized datatype for abstract values (wrt. a given party)*)
10 datatype 't VAL = FREEDOM 't | UTILITY 't | SECURITY 't | EQUALITY 't
11 type_synonym v = "(c)VAL⇒bool" (*principles: sets of (abstract) values*)
12 type_synonym cv = "c⇒v" (*principles are specified wrt. a given party*)
13 (*notation for sets*)
14 abbreviation vset1 ("_{_}") where "{_} ≡ λx::(c)VAL. x={_}"
15 abbreviation vset2 ("_{_,_}") where "{α,β} ≡ λx::(c)VAL. x=α ∨ x=β"
16 (*abstract values and value principles*)
17 abbreviation utility::cv ("UTILITY_") where "UTILITY^x ≡ {UTILITY x}"
18 abbreviation security::cv ("SECURITY_") where "SECURITY^x ≡ {SECURITY x}"
19 abbreviation equality::cv ("EQUALITY_") where "EQUALITY^x ≡ {EQUALITY x}"
20 abbreviation freedom::cv ("FREEDOM_") where "FREEDOM^x ≡ {FREEDOM x}"
21 abbreviation stab::cv ("STAB_") where "STAB^x ≡ {SECURITY x, UTILITY x}"
22 abbreviation effi::cv ("EFFI_") where "EFFI^x ≡ {UTILITY x, SECURITY x}"
23 abbreviation gain::cv ("GAIN_") where "GAIN^x ≡ {UTILITY x, FREEDOM x}"
24 abbreviation will::cv ("WILL_") where "WILL^x ≡ {FREEDOM x, UTILITY x}"
25 abbreviation resp::cv ("RESP_") where "RESP^x ≡ {FREEDOM x, EQUALITY x}"
26 abbreviation fair::cv ("FAIR_") where "FAIR^x ≡ {EQUALITY x, FREEDOM x}"
27 abbreviation equi::cv ("EQUI_") where "EQUI^x ≡ {EQUALITY x, SECURITY x}"
28 abbreviation reli::cv ("RELI_") where "RELI^x ≡ {SECURITY x, EQUALITY x}"
29 (**Value Theory*)
30 consts Irel::"t⇒v" ("I") (*incidence relation worlds-values*)
31 (*derivation operators (cf. theory of "formal concept analysis") *)
32 abbreviation intent::"σ⇒v" ("_↑") where "W↑ ≡ λv. ∀x. W x → I x v"
33 abbreviation extent::"v⇒σ" ("_↓") where "V↓ ≡ λw. ∀x. V x → I w x"
34 abbreviation extent_brkt ("[_]") where "[V] ≡ V↓" (*alternative notation*)
35 (*connective for aggregating value principles*)
36 abbreviation aggr ("[_⊕_]") where "[V₁⊕V₂] ≡ (V₁↓) ∨ (V₂↓)"
37 (*chosen variant for preference relation (cf. Halpern (1997)*)
38 abbreviation pref::"σ⇒σ" ("_≺") where "φ ≺ ψ ≡ φ ≺AE ψ"
39 (*schema for value principle promotion*)
40 abbreviation "Promotes F D V ≡ [F → □·(D ↔ ◇·(V↓))]"
41 (*proposition for testing for value conflict*)
42 abbreviation conflict ("Conflict_") where (*conflict for value support*)
43 "Conflict^x ≡ [SECURITY^x] ∧ [EQUALITY^x] ∧ [FREEDOM^x] ∧ [UTILITY^x]"
44 (*verify consistency of this theory*)
45 lemma "True" nitpick[satisfy] oops
46 end

```

Fig. 15 Encoding of the legal DSL (value ontology)

Finally we introduce “Promotes” schema for encoding the promotion of value principles via legal decisions (Line 40) and we introduce a notion “Conflict<sup>x</sup>” expressing a legal value conflict for a party  $x$  (Lines 42-43).

The consistency of the theory is confirmed by *Nitpick* (Line 45).

Tests on the modelling and encoding of the legal DSL are displayed in Fig. 16.

Among others, we verify that the pair of operators for *extension* ( $\downarrow$ ) and *intension* ( $\uparrow$ ), cf. *Formal Concept Analysis* (Ganter and Wille 2012), constitute indeed a Galois connection (Lines 6–18), and we carry out some further tests on the value theory (extending the ones presented in §6) concerning value aggregation and consistency (Lines 20ff.).

```

1 | theory ValueOntologyTestLong imports ValueOntology (** Benzmüller, Fuenmayor & Lomfeld, 2021 **)
2 | begin
3 | lemma "True" nitpick[satisfy,show_all,card :=10] oops
4 | lemma "[ConflictP]" nitpick[satisfy,card :=4] nitpick oops (*contingent*)
5 | (*derivation operators satisfy main properties of Galois connections*)
6 | lemma G: "B ⊆ A↑ ↔ A ⊆ B↓" by blast
7 | lemma G1: "A ⊆ A↑↓" by simp
8 | lemma G2: "B ⊆ B↑↓" by simp
9 | lemma G3: "A₁ ⊆ A₂ → A₂↑ ⊆ A₁↑" by simp
10 | lemma G4: "B₁ ⊆ B₂ → B₂↓ ⊆ B₁↓" by simp
11 | lemma cl1: "A↑ = A↑↑↑" by blast
12 | lemma cl2: "B↓ = B↓↓↓" by blast
13 | lemma dual1a: "(A₁ ⊔ A₂)↑ = (A₁↑ ⊔ A₂↑)" by blast
14 | lemma dual1b: "(B₁ ⊓ B₂)↓ = (B₁↓ ⊓ B₂↓)" by blast
15 | lemma " (A₁ ⊓ A₂)↑ ⊆ (A₁↑ ⊔ A₂↑)" nitpick oops (*countermodel*)
16 | lemma " (B₁ ⊔ B₂)↓ ⊆ (B₁↓ ⊓ B₂↓)" nitpick oops (*countermodel*)
17 | lemma dual2a: "(A₁↑ ⊔ A₂↑) ⊆ (A₁ ⊔ A₂)↑" by blast
18 | lemma dual2b: "(B₁↓ ⊓ B₂↓) ⊆ (B₁ ⊓ B₂)↓" by blast
19 | (*value conflict tests*)
20 | lemma "[RELI^P] ∧ [WILL^P] → ConflictP" by simp
21 | lemma "[Conflict^P] → [RELI^P] ∧ [WILL^P]" by simp
22 | lemma "[RELI^P] ∧ [WILL^P]" nitpick[satisfy] nitpick oops (*contingent*)
23 | lemma "[FAIR^d] ∧ [EFFI^d]" nitpick[satisfy] nitpick oops (*contingent*)
24 | lemma "[¬Conflict^d] ∧ [FAIR^d] ∧ [EFFI^d]"
25 | nitpick[satisfy,show_all] nitpick oops (*contingent: p & d independent*)
26 | lemma "[¬Conflict^d] ∧ (¬Conflict^P) ∧ [RELI^d] ∧ [WILL^P]"
27 | nitpick[satisfy,show_all] nitpick oops (*contingent: p & d independent*)
28 | (*values in two non-opposed quadrants: no conflict*)
29 | lemma "[WILL^x] ∧ [STAB^x] → Conflict^x" nitpick oops (*countermodel found*)
30 | lemma "[WILL^x] ∧ [GAIN^x] ∧ [EFFI^x] ∧ [STAB^x] → Conflict^x" nitpick oops
31 | (*values in two opposed quadrants: conflict*)
32 | lemma "[RESP^x] ∧ [STAB^x] → Conflict^x" by simp
33 | (*values in three quadrants: conflict*)
34 | lemma "[WILL^x] ∧ [EFFI^x] ∧ [RELI^x] → Conflict^x" by simp
35 | (*values in opposed quadrants for different parties: no conflict*)
36 | lemma "[EQUI^x] ∧ [GAIN^x] → (Conflict^x ∨ Conflict^y)" nitpick oops (*cntmdl*)
37 | lemma "[RESP^x] ∧ [STAB^y] → (Conflict^x ∨ Conflict^y)" nitpick oops (*cntmdl*)
38 | (*value preferences tests*)
39 | lemma "[WILL^x] < [WILL^y @ STAB^x]" nitpick nitpick[satisfy] oops (*contingent*)
40 | lemma "[WILL^x] < [STAB^y]" → "[WILL^x] < [WILL^y @ STAB^x]" by blast
41 | lemma "[WILL^x] < [STAB^y]" → "[WILL^x] < [RELI^y @ STAB^x]" by blast
42 | lemma "[WILL^x] < [WILL^y @ STAB^x]" → "[WILL^x] < [STAB^y]" (*nitpick*) nitpick[satisfy] oops (*ctgnt?*)
43 | lemma "[WILL^x] < [RELI^y @ STAB^x]" → "[WILL^x] < [STAB^y]" nitpick nitpick[satisfy] oops (*contingent*)
44 | lemma "[WILL^x @ STAB^y] < [WILL^x]" nitpick nitpick[satisfy] oops (*contingent*)
45 | lemma "[WILL^x @ STAB^y] < [WILL^x]" → "[STAB^y] < [WILL^x]" by metis
46 | lemma "[RELI^x @ STAB^y] < [WILL^x]" → "[STAB^y] < [WILL^x]" by metis
47 | lemma "[STAB^y] < [WILL^x]" → "[WILL^x @ STAB^y] < [WILL^x]" nitpick nitpick[satisfy] oops (*contingent*)
48 | lemma "[STAB^y] < [WILL^x]" → "[RELI^x @ STAB^y] < [WILL^x]" nitpick nitpick[satisfy] oops (*contingent*)
49 | (*basic properties*)
50 | lemma "[X] < [X]" nitpick nitpick[satisfy] oops (*contingent*)
51 | lemma "[((X) < (Y)) ∧ ((Y) < (Z))] → ((X) < (Z))" using tSBR by blast (*transitive*)
52 | lemma "[([X] < [Y]) ∧ ([Y] < [X])] → X = Y" nitpick oops (*not antisymmetric*)
53 | end

```

Fig. 16 Formally verifying/testing the legal DSL or value ontology

### A.3 Legal and World Knowledge

The encoding of the relevant legal & world knowledge (LWK) is shown in Fig. 17. The defined *Isabelle/HOL* theory is termed “GeneralKnowledge” and imports the “ValueOntology” (and thus recursively also “PreferenceLogicBasics”) theory.

Lines 4–5 Declaration of logical constant symbols that stand for kinds of legally relevant situations.

Lines 8–11 Meaning postulates for these kinds of legally relevant situations are introduced.

Lines 14–16 Preference relations for these kinds of legally relevant situations are introduced.

Lines 18–26 Some simple vocabulary is introduced and some taxonomic relations for wild and domestic animals are specified.

Lines 28–36 Some relevant situational *factors* are declared and subsequently constrained by meaning postulates.

Line 39 An example for a value preference conditioned on *factors* is specified.

Lines 41–46 The situational *factors* are related with values and with ruling outcomes according to the notion of value *promotion*.

Line 48 The model finder *Nitpick* is used to confirm the consistency of the introduced theory.

```

1 theory GeneralKnowledge imports ValueOntology (** Benzmüller, Fuenmayor & Lomfeld, 2021 **)
2 begin (** General Legal and World Knowledge (LWK) **)
3 (*LWK: kinds of situations addressed*)
4 consts appObject::σ appAnimal::σ (*appropriation of objects/animals in general*)
5 appWildAnimal::σ appDomAnimal::σ (*appropriation of wild/domestic animals*)
6 (*LWK: postulates for kinds of situations*)
7 axiomatization where
8 W1: "[appAnimal → appObject]" and
9 W2: "[¬(appWildAnimal ∧ appDomAnimal)]" and
10 W3: "[appWildAnimal → appAnimal]" and
11 W4: "[appDomAnimal → appAnimal]"
12 (*LWK: (prima facie) value preferences for kinds of situations*)
13 axiomatization where
14 R1: "[appAnimal → ([STAB*] < [STABd])]" and
15 R2: "[appWildAnimal → ([WILLx-1] < [STAB*])]" and
16 R3: "[appDomAnimal → ([STABx-1] < [RELI*⊕RESP*])]"
17 (*LWK: domain vocabulary*)
18 typedecl e (*declares new type for 'entities'*)
19 consts
20 Animal::"e⇒σ" Domestic::"e⇒σ" Fox::"e⇒σ" Parrot::"e⇒σ" Pet::"e⇒σ" FreeRoaming::"e⇒σ"
21 (*LWK: domain knowledge (about animals)*)
22 axiomatization where
23 W5: "[∀a. Fox a → Animal a]" and
24 W6: "[∀a. Parrot a → Animal a]" and
25 W7: "[∀a. (Animal a ∧ FreeRoaming a ∧ ¬Pet a) → ¬Domestic a]" and
26 W8: "[∀a. Animal a ∧ Pet a → Domestic a]"
27 (*LWK: legally-relevant, situational 'factors'*)
28 consts Own::"c⇒σ" (*object is owned by party c*)
29 Poss::"c⇒σ" (*party c has actual possession of object*)
30 Intent::"c⇒σ" (*party c has intention to possess object*)
31 Mal::"c⇒σ" (*party c acts out of malice*)
32 Mtn::"c⇒σ" (*party c respons. for maintenance of object*)
33 (*LWK: meaning postulates for general notions*)
34 axiomatization where
35 W9: "[Poss x → (¬Poss x-1)]" and
36 W10: "[Own x → (¬Own x-1)]"
37 (*LWK: conditional value preferences, e.g. from precedents*)
38 axiomatization where
39 R4: "[([Mal x-1 ∧ Own x) → ([STABx-1] < [RESP*⊕RELI*])]"
40 (*LWK: relate values, outcomes and situational 'factors'*)
41 axiomatization where
42 F1: "Promotes (Intent x) (For x) WILL*" and
43 F2: "Promotes (Mal x) (For x-1) RESP*" and
44 F3: "Promotes (Poss x) (For x) STAB*" and
45 F4: "Promotes (Mtn x) (For x) RESP*" and
46 F5: "Promotes (Own x) (For x) RELI*"
47 (*Theory is consistent, (non-trivial) model found*)
48 lemma True nitpick[satisfy,card :=4] oops
49 end

```

Fig. 17 Encoding of relevant legal & world knowledge

#### A.4 Modelling Pierson v. Post

The *Isabelle/HOL* encoding of two scenarios in the Pierson v. Post case is presented in Figs. 18 and 19.

In Fig. 18, which presents the initial ruling in favour of Pierson, the *Isabelle/HOL* theory is termed “Pierson” and imports the theory “GeneralKnowledge” (which recursively imports theories “ValueOntology” and “PreferenceLogicBasics”).

Lines 5–19 (generic) theory and (contingent) facts suitable to the defendant (Pierson) are postulated.

Lines 21–22 automated proof justifying the ruling for Pierson; the dependencies of the proof are shown.

Lines 24–35 corresponding interactive proof (with the same dependencies as for the automated one) modelling the argument justifying the finding for Pierson.

Lines 36–44 various checks for consistency of the assumptions and the absence of value conflicts.

```

1 theory Pierson imports GeneralKnowledge (** Benzmüller, Fuenmayor & Lomfeld, 2021 **)
2 begin (** Pierson v. Post "wild animal" case **)
3 (*unimportant*) nitpick_params[user_axioms,expect=genuine,show_all,format=3]
4 (*case-specific 'world-vocabulary'*)
5 consts α::"e" (*appropriated animal (fox in this case) *)
6 consts Pursue::"c⇒e⇒σ" Capture::"c⇒e⇒σ"
7 (***** pro-defendant (Pierson) argument *****)
8 (*defendant's theory*)
9 abbreviation "dT1 ≡ [(∃c. Capture c α ∧ ¬Domestic α) → appWildAnimal]"
10 abbreviation "dT2 ≡ [∀c. Pursue c α → Intent c]"
11 abbreviation "dT3 ≡ [∀c. Capture c α → Poss c]"
12 abbreviation "d_theory ≡ dT1 ∧ dT2 ∧ dT3"
13 (*defendant's facts*)
14 abbreviation "dF1 w ≡ Fox α w"
15 abbreviation "dF2 w ≡ FreeRoaming α w"
16 abbreviation "dF3 w ≡ ¬Pet α w"
17 abbreviation "dF4 w ≡ Pursue p α w"
18 abbreviation "dF5 w ≡ Capture d α w"
19 abbreviation "d_facts ≡ dF1 ∧ dF2 ∧ dF3 ∧ dF4 ∧ dF5"
20 (*decision for defendant (Pierson) can be proven automatically*)
21 theorem Pierson: "d_theory → [d_facts → □¬For d]"
22 by (smt F1 F3 ForAx R2 W5 W7 other.simps tsBR)
23 (*we reconstruct the reasoning process leading to the decision for the defendant*)
24 theorem Pierson': assumes d_theory and "d_facts w" shows "□¬For d w"
25 proof -
26 have 1: "appWildAnimal w" using W5 W7 assms by blast
27 have 2: "[WILLP]¬[STABd]" using 1 R2 assms by fastforce
28 have 3: "[¬(◇¬[WILLP]) → ◇¬[STABd]]" using 2 tsBR by smt
29 have 4: "□¬(For p ↔ ◇¬[WILLP]) w" using F1 assms by meson
30 have 5: "□¬(For d ↔ ◇¬[STABd]) w" using F3 assms by meson
31 have 6: "□¬((◇¬[WILLP]) ∨ (◇¬[STABd])) w" using 4 5 ForAx by (smt other.simps)
32 have 7: "□¬(◇¬[STABd]) w" using 3 6 by blast
33 have 8: "□¬(For d) w" using 5 7 by simp
34 then show ?thesis by simp
35 qed
36 (***** Further checks (using model finder) *****)
37 (*defendant's theory and facts are logically consistent*)
38 lemma "d_theory ∧ [d_facts]" nitpick[satisfy,card ι=3] oops (* (non-trivial) model found*)
39 (*decision for defendant is compatible with premises and lacks value conflicts*)
40 lemma "[¬ConflictP] ∧ [¬Conflictd] ∧ d_theory ∧ [d_facts ∧ For d]"
41 nitpick[satisfy,card ι=3] oops (* (non-trivial) model found*)
42 (*situations where decision holds for plaintiff are compatible too*)
43 lemma "[¬ConflictP] ∧ [¬Conflictd] ∧ d_theory ∧ [d_facts ∧ For p]"
44 nitpick[satisfy,card ι=3] oops (* (non-trivial) model found*)
45 end

```

Fig. 18 Modelling the Pierson v. Post case; ruling for Pierson

As a further illustration, we present in Fig. 19 a plausible counterargument by Post. The *Isabelle/HOL* theory is now termed “Post” and imports the theory “GeneralKnowledge” (which recursively imports theories “ValueOntology” and “PreferenceLogicBasics”).

Lines 5–24 theory and facts suitable to the plaintiff (Post) are postulated.

Lines 26–27 automated proof justifying the ruling for Post; the dependencies of the proof are shown.

Lines 29–42 corresponding interactive proof (with the same dependencies as for the automated one) modelling the argument justifying the finding for Post.

Lines 43–51 various checks for consistency of the assumptions and the absence of value conflicts.

```

1 theory Post imports GeneralKnowledge (** Benzmüller, Fuenmayor & Lomfeld, 2021 **)
2 begin (** Pierson v. Post "wild animal" case **)
3 (*unimportant*) nitpick_params[user_axioms,expect=genuine,show_all,format=3]
4 (*case-specific 'world-vocabulary'*)
5 consts  $\alpha$ ::"e" (*appropriated animal (fox in this case) *)
6 consts Pursue::"c $\Rightarrow$ e $\Rightarrow$  $\sigma$ " Capture::"c $\Rightarrow$ e $\Rightarrow$  $\sigma$ "
7 (***** pro-plaintiff (Post) argument *****)
8 (*acknowledges from defendant's theory*)
9 abbreviation "dT2  $\equiv$  [ $\forall$ c. Pursue c  $\alpha$   $\rightarrow$  Intent c]"
10 abbreviation "dT3  $\equiv$  [ $\forall$ c. Capture c  $\alpha$   $\rightarrow$  Poss c]"
11 (*theory amendment: the animal was chased by a professional hunter (Post); protecting
12  hunters' labor, thus fostering economic efficiency, prevails over legal certainty.*)
13 consts Hunter::"c $\Rightarrow$  $\sigma$ " hunting::" $\sigma$ " (*new kind of situation: hunting*)
14 (*plaintiff's theory*)
15 abbreviation "pT1  $\equiv$  [( $\exists$ c. Hunter c  $\wedge$  Pursue c  $\alpha$ )  $\rightarrow$  hunting]"
16 abbreviation "pT2  $\equiv$   $\forall$ x. [hunting  $\rightarrow$  ([STABx-1]  $\leftarrow$  [EFFIx@WILLx])]" (*case-specific rule*)
17 abbreviation "pT3  $\equiv$   $\forall$ x. Promotes (hunting  $\wedge$  Hunter x) (For x) EFFIx"
18 abbreviation "p_theory  $\equiv$  pT1  $\wedge$  pT2  $\wedge$  pT3  $\wedge$  dT2  $\wedge$  dT3"
19 (*plaintiff's facts*)
20 abbreviation "pF1 w  $\equiv$  Fox  $\alpha$  w"
21 abbreviation "pF2 w  $\equiv$  Hunter p w"
22 abbreviation "pF3 w  $\equiv$  Pursue p  $\alpha$  w"
23 abbreviation "pF4 w  $\equiv$  Capture d  $\alpha$  w"
24 abbreviation "p_facts  $\equiv$  pF1  $\wedge$  pF2  $\wedge$  pF3  $\wedge$  pF4"
25 (*decision for plaintiff (Post) can be proven automatically (needs approx. 20s)*)
26 theorem Post: "p_theory  $\rightarrow$  [p_facts  $\rightarrow$   $\square$ -For p]"
27 by (smt F1 F3 ForAx tBR SBR_def other.simps)
28 (*we reconstruct the reasoning process leading to the decision for the plaintiff*)
29 theorem Post': assumes p_theory and "p_facts w" shows " $\square$ -For p w"
30 proof -
31 have 1: "hunting w" using assms by auto
32 have 2: "[[STABd]  $\leftarrow$  [EFFIp@WILLp]]" using 1 assms by auto
33 have 3: "[( $\diamond$ -[STABd])  $\rightarrow$   $\diamond$ -([EFFIp]  $\vee$  [WILLp])]" using 2 tsBR by smt
34 have 4: " $\square$ -((For p  $\leftrightarrow$   $\diamond$ -[EFFIp]) w)" using assms by meson
35 have 5: " $\square$ -((For p  $\leftrightarrow$   $\diamond$ -[WILLp]) w)" using F1 assms by meson
36 have 6: " $\square$ -((For d  $\leftrightarrow$   $\diamond$ -[STABd]) w)" using F3 assms by meson
37 have 7: " $\square$ -(( $\diamond$ -[EFFIp])  $\vee$  ( $\diamond$ -[WILLp])  $\vee$  ( $\diamond$ -[STABd])) w"
38 using 4 5 6 ForAx by (smt other.simps)
39 have 8: " $\square$ -(( $\diamond$ -[EFFIp])  $\vee$  ( $\diamond$ -[WILLp])) w" using 3 7 by metis
40 have 9: " $\square$ -((For p) w)" using 4 5 8 by auto
41 then show ?thesis by simp
42 qed
43 (***** Further checks (using model finder) *****)
44 (*plaintiff's theory and facts are logically consistent*)
45 lemma "p_theory  $\wedge$  [p_facts]" nitpick[satisfy,card  $\iota$ =2] oops (* (non-trivial) model found*)
46 (*decision for plaintiff is compatible with premises and lacks value conflicts*)
47 lemma "[ $\neg$ -Conflictp]  $\wedge$  [ $\neg$ -Conflictd]  $\wedge$  p_theory  $\wedge$  [p_facts  $\wedge$  For p]"
48 nitpick[satisfy,card  $\iota$ =2] oops (* (non-trivial) model found*)
49 (*situations where decision holds for defendant are compatible too*)
50 lemma "[ $\neg$ -Conflictp]  $\wedge$  [ $\neg$ -Conflictd]  $\wedge$  p_theory  $\wedge$  [p_facts  $\wedge$  For d]"
51 nitpick[satisfy,card  $\iota$ =2] oops (* (non-trivial) model found*)
52 end

```

Fig. 19 Modelling the Pierson v. Post case; ruling for Post

## A.5 Modelling Conti v. ASPCA

The reconstructed theory for the Conti v. ASPCA case is displayed in Fig. 20. The *Isabelle/HOL* theory is termed “Conti” and imports the theory “GeneralKnowledge” (which recursively imports theories “ValueOntology” and “PreferenceLogicBasics”).

Lines 5–20 the theory and the facts of the pro-plaintiff (ASPCA) argument are formulated.

Lines 22–23 automated proof justifying the ruling for ASPCA; the dependencies of the proof are shown.

Lines 25–38 corresponding interactive proof (with the same dependencies as for the automated one) modelling the argument justifying the finding for ASPCA.

Lines 39–47 various checks for consistency of the assumptions and the absence of value conflicts.

```

1 theory Conti imports GeneralKnowledge (** Benzmüller, Fuenmayor & Lomfeld, 2021 **)
2 begin (** Conti v. ASPCA "wild animal" case **)
3   (*unimportant*) nitpick_params[user_axioms,expect=genuine,show_all,format=3]
4   (*case-specific 'world-vocabulary'*)
5   consts α:"e" (*appropriated animal (parrot in this case) *)
6   consts Care::"c⇒e⇒σ" Prop::"c⇒e⇒σ" Capture::"c⇒e⇒σ"
7   (***** pro-plaintiff (ASPCA) argument *****
8   (*plaintiff's theory*)
9   abbreviation "pT1 ≡ [(∃c. Capture c α ∧ Domestic α) → appDomAnimal]"
10  abbreviation "pT2 ≡ [∀c. Care c α → Mtn c]"
11  abbreviation "pT3 ≡ [∀c. Prop c α → Own c]"
12  abbreviation "pT4 ≡ [∀c. Capture c α → Poss c]" (*concedes' to defendant*)
13  abbreviation "p_theory ≡ pT1 ∧ pT2 ∧ pT3 ∧ pT4"
14  (*plaintiff's facts*)
15  abbreviation "pF1 w ≡ Parrot α w"
16  abbreviation "pF2 w ≡ Pet α w"
17  abbreviation "pF3 w ≡ Care p α w"
18  abbreviation "pF4 w ≡ Prop p α w"
19  abbreviation "pF5 w ≡ Capture d α w"
20  abbreviation "p_facts ≡ pF1 ∧ pF2 ∧ pF3 ∧ pF4 ∧ pF5"
21  (*decision for plaintiff (ASPCA) can be proven automatically*)
22  theorem ASPCA: "p_theory → [p_facts → □¬For p]"
23  by (smt F3 F4 F5 ForAx R3 W6 W8 tBR SBR_def other.simps(1))
24  (*we reconstruct the reasoning process leading to the decision for the plaintiff*)
25  theorem ASPCA': assumes p_theory and "p_facts w" shows "□¬For p w"
26  proof -
27    have 1: "appDomAnimal w" using W6 W8 assms by blast
28    have 2: "[[STABd] <- [RELIp⊕RESPp]]" using 1 R3 by fastforce
29    have 3: "[[□¬[STABd]] → □¬([RELIp] ∨ [RESPp])]" using 2 tSBR by smt
30    have 4: "□¬(For p ↔ □¬[RELIp]) w" using F5 assms by metis
31    have 5: "□¬(For p ↔ □¬[RESPp]) w" using F4 assms by metis
32    have 6: "□¬(For d ↔ □¬[STABd]) w" using F3 assms by meson
33    have 7: "□¬((□¬[RELIp]) ∨ (□¬[RESPp]) ∨ (□¬[STABd])) w"
34      using 4 5 6 ForAx by (smt other.simps)
35    have 8: "□¬((□¬[RELIp]) ∨ (□¬[RESPp])) w" using 3 7 by metis
36    have 9: "□¬(For p) w" using 4 5 8 by auto
37    then show ?thesis by simp
38  qed
39  (***** Further checks (using model finder) *****
40  (*plaintiff's theory and facts are logically consistent*)
41  lemma "p_theory ∧ [p_facts]" nitpick[satisfy,card ι=3] oops (* (non-trivial) model found*)
42  (*decision for plaintiff is compatible with premises and lacks value conflicts*)
43  lemma "[¬Conflictp] ∧ [¬Conflictd] ∧ p_theory ∧ [p_facts ∧ For p]"
44    nitpick[satisfy,card ι=3] oops (* (non-trivial) model found*)
45  (*situations where decision holds for defendant are compatible too*)
46  lemma "[¬Conflictp] ∧ [¬Conflictd] ∧ p_theory ∧ [p_facts ∧ For d]"
47    nitpick[satisfy,card ι=3] oops (* (non-trivial) model found*)
48  end

```

Fig. 20 Modelling of the Conti v. ASPCA case



## A.6 Complex (Counter-)Models

In Fig. 21 we present an example of a model computed by model finder *Nitpick* for the statement in Line 41 in Fig. 20. This non-trivial model features three possible worlds/states. It illustrates the richness of the information and the level of detail that is supported in model (and countermodel) finding technology for HOL. This information is very helpful to support the knowledge engineer and user of the LOGIKEY framework to gain insight about the modelled structures. We observe that the proof assistant *Isabelle/HOL* allows for the parallel execution of its integrated tools. We can thus execute, for a given candidate theorem, all three tasks in parallel (and in different modes): theorem proving, model finding, and countermodel finding. This is one reason for the very good response rates we have experienced in our work with the system – despite the general undecidability of HOL.

Nitpick found a model for card e = 1 and card  $\iota$  = 3:

```
Types:
c = {d, p}
e  $\times$   $\iota$  [boxed] = {(e1,  $\iota$ 1), (e1,  $\iota$ 2), (e1,  $\iota$ 3)}
c VAL = {UTILITY d, UTILITY p, EQUALITY p, ...}
Constants:
Capture =
( $\lambda$ x. _)
((d, e1,  $\iota$ 1) := True, (d, e1,  $\iota$ 2) := True, (d, e1,  $\iota$ 3) := True, (p, e1,  $\iota$ 1) := False,
(p, e1,  $\iota$ 2) := False, (p, e1,  $\iota$ 3) := False)
Care = ( $\lambda$ x. _)((p, e1,  $\iota$ 1) := True, (p, e1,  $\iota$ 2) := True, (p, e1,  $\iota$ 3) := True)
Prop =
( $\lambda$ x. _)
((d, e1,  $\iota$ 1) := False, (d, e1,  $\iota$ 2) := False, (d, e1,  $\iota$ 3) := False, (p, e1,  $\iota$ 1) := True,
(p, e1,  $\iota$ 2) := True, (p, e1,  $\iota$ 3) := True)
 $\alpha$  = e1
Animal = ( $\lambda$ x. _)((e1,  $\iota$ 1) := True, (e1,  $\iota$ 2) := True, (e1,  $\iota$ 3) := True)
Domestic = ( $\lambda$ x. _)((e1,  $\iota$ 1) := True, (e1,  $\iota$ 2) := True, (e1,  $\iota$ 3) := True)
Fox =  $\lambda$ x. _
FreeRoaming = ( $\lambda$ x. _)((e1,  $\iota$ 1) := False, (e1,  $\iota$ 2) := False, (e1,  $\iota$ 3) := False)
Intent =
( $\lambda$ x. _)
((d,  $\iota$ 1) := False, (d,  $\iota$ 2) := False, (d,  $\iota$ 3) := False, (p,  $\iota$ 1) := False, (p,  $\iota$ 2) := False, (p,  $\iota$ 3) := False)
Mal =
( $\lambda$ x. _)
((d,  $\iota$ 1) := False, (d,  $\iota$ 2) := False, (d,  $\iota$ 3) := False, (p,  $\iota$ 1) := False, (p,  $\iota$ 2) := False, (p,  $\iota$ 3) := False)
Mtn =
( $\lambda$ x. _)((d,  $\iota$ 1) := True, (d,  $\iota$ 2) := True, (d,  $\iota$ 3) := True, (p,  $\iota$ 1) := True, (p,  $\iota$ 2) := True, (p,  $\iota$ 3) := True)
Own =
( $\lambda$ x. _)
((d,  $\iota$ 1) := False, (d,  $\iota$ 2) := False, (d,  $\iota$ 3) := False, (p,  $\iota$ 1) := True, (p,  $\iota$ 2) := True, (p,  $\iota$ 3) := True)
Parrot = ( $\lambda$ x. _)((e1,  $\iota$ 1) := True, (e1,  $\iota$ 2) := True, (e1,  $\iota$ 3) := True)
Pet = ( $\lambda$ x. _)((e1,  $\iota$ 1) := True, (e1,  $\iota$ 2) := True, (e1,  $\iota$ 3) := True)
Poss =
( $\lambda$ x. _)
((d,  $\iota$ 1) := True, (d,  $\iota$ 2) := True, (d,  $\iota$ 3) := True, (p,  $\iota$ 1) := False, (p,  $\iota$ 2) := False, (p,  $\iota$ 3) := False)
appAnimal = ( $\lambda$ x. _)( $\iota$ 1 := True,  $\iota$ 2 := True,  $\iota$ 3 := True)
appDomAnimal = ( $\lambda$ x. _)( $\iota$ 1 := True,  $\iota$ 2 := True,  $\iota$ 3 := True)
appObject = ( $\lambda$ x. _)( $\iota$ 1 := True,  $\iota$ 2 := True,  $\iota$ 3 := True)
appWildAnimal = ( $\lambda$ x. _)( $\iota$ 1 := False,  $\iota$ 2 := False,  $\iota$ 3 := False)
BR = ( $\lambda$ x. _)
(( $\iota$ 1,  $\iota$ 1) := True, ( $\iota$ 1,  $\iota$ 2) := False, ( $\iota$ 1,  $\iota$ 3) := False, ( $\iota$ 2,  $\iota$ 1) := False, ( $\iota$ 2,  $\iota$ 2) := True,
( $\iota$ 2,  $\iota$ 3) := False, ( $\iota$ 3,  $\iota$ 1) := False, ( $\iota$ 3,  $\iota$ 2) := False, ( $\iota$ 3,  $\iota$ 3) := True)
For =
( $\lambda$ x. _)
((d,  $\iota$ 1) := True, (d,  $\iota$ 2) := True, (d,  $\iota$ 3) := True, (p,  $\iota$ 1) := False, (p,  $\iota$ 2) := False, (p,  $\iota$ 3) := False)
I = ( $\lambda$ x. _)
(( $\iota$ 1, UTILITY d) := False, ( $\iota$ 1, UTILITY p) := False, ( $\iota$ 2, UTILITY d) := False, ( $\iota$ 2, UTILITY p) := False,
( $\iota$ 3, UTILITY d) := False, ( $\iota$ 3, UTILITY p) := False)
other = ( $\lambda$ x. _)(d := p, p := d)
```

**Fig. 21** Example of a (satisfying) model to the statement in Line 26 in Fig. 20