

# Encoding Legal Balancing: Automating an Abstract Ethico-Legal Value Ontology in Preference Logic

Christoph Benz Müller<sup>1</sup>, David Fuenmayor<sup>1</sup>, Bertram Lomfeld<sup>2</sup>

<sup>1</sup>Department of Computer Science and Mathematics, FU Berlin, Berlin, Germany

<sup>2</sup>Department of Law, FU Berlin, Berlin, Germany

{c.benzmueller, david.fuenmayor, bertram.lomfeld}@fu-berlin.de

## Abstract

Enabling machines to legal balancing is a non-trivial task challenged by a multitude of factors some of which are addressed and explored in this work. We propose a holistic approach to formal modeling at different abstraction layers supported by a pluralistic framework in which the encoding of an ethico-legal value and upper ontology is developed in combination with the exploration of a formalization logic, with legal domain knowledge and with exemplary use cases until a reflective equilibrium is reached. Our work is enabled by a meta-logical approach to universal logical reasoning and it applies the recently introduced *LogiKEy* methodology for designing normative theories for ethical and legal reasoning. The particular focus in this paper is on the formalization and encoding of a value ontology suitable e.g. for explaining and resolving legal conflicts in property law (wild animal cases).

## 1 Introduction

Law today has to reflect highly pluralistic environments. There are plural values, world-views, and logics which may even be considered as constituting plural worlds. One function of modern, reflexive law is to enable the social interaction within and between such plural worlds (Teubner 1983; Lomfeld 2017). Any sound model of legal reasoning needs to be pluralistic while at the same time reflecting the uniting force of law.

Logical reconstructions of legal reasoning quite often separate deductive rule application and legal interpretation, cf. the overview in (Prakken and Sartor 2015). Yet, understanding legal reasoning as practical argumentation (Alexy 1978), this separation breaks down. Nonetheless different types of rule-based (Hage 1997), case-based (Horty 2011) and value-based (Bench-Capon et al. 2005) reasoning remain. A discourse theory of law could integrate these different types by translating legal reasoning into the balancing of plural and opposing socio-legal values (Lomfeld 2015).

There are several models to quantify legal balancing, e.g. (Alexy 2003; Sartor 2010). Nonetheless, these approaches need to get “integrated with logic and argumentation to provide a comprehensive account of value-based reasoning” (Sartor 2018). An adequate balancing model of legal reasoning, thus, has to reconstruct rule subsumption and case interpretation as argumentation process between conflicting values. Here, the differentiation of legal norms into rules

and principles reveals its potential (Alexy 2000). Legal principles could be understood as material values on a deep level of legal balancing, which is structured by legal rules on an explicit upper level of legal reasoning (Lomfeld 2015).

Within the recent AI & Law perspective value conflicts are often modeled according to Dung’s abstract argumentation (Dung 1995) as a value argumentation framework (Bench-Capon et al. 2005). Yet, if preferences between defensible rules are reconstructed and justified in terms of preferences between values, some questions about values necessarily pop up: “Are values scalar? [...] Can values be ordered at all? [...] How can sets of values be compared? [...] Can several less important values together overcome a more important value?” (Bench-Capon and Sartor 2003).

Thus, an encompassing approach for legal reasoning as practical argumentation needs not only a formal reconstruction of the relation between legal values (principles) and legal rules, but also a value ontology that allows to systematize value comparison and conflicts as “discursive grammar” (Lomfeld 2015; Lomfeld 2019).

In this paper we demonstrate how an abstract level encoding and automation of an ethico-legal value ontology, i.e. Lomfeld’s “discursive grammar” of justification, can be achieved by reusing and extending the *LogiKEy* methodology and formal framework (Benz Müller et al. 2020). This framework works with shallow semantical embeddings (SSEs) of (combinations of) non-classical logics in classical higher-order logic (HOL). HOL thereby serves as a meta-logic, rich enough to support the encoding of a plurality of “object logics” (e.g. conditional, deontic or epistemic logics and combinations thereof) and plural and adaptable value systems. The embedded “object logics” are used for the iterative, experimental encoding of normative theories, which finally help to reason with concrete legal cases utilizing the encoded value system. This reflects ideas of understanding the solution of legal cases as theory construction, “building, evaluating and using theories” (Bench-Capon and Sartor 2003).<sup>1</sup> This multi-layered knowledge engineering process is supported in our framework by adapting automated

<sup>1</sup>They quote e.g. McCarty (1995, p285): “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.”

theorem proving (ATP) technology for meta-logic HOL.

Ethico-legal ontologies also constitute a core ingredient to enable the computation, assessment and communication of abstract-level justifications in the future ethico-legal governance of AI (Benzmüller and Lomfeld 2020). A sound implementation of any legally accountable “moral machine” requires the development of upper-level ontologies to guide and connect the encoding of concrete regulatory codes (Hoekstra et al. 2009; Fuenmayor and Benzmüller 2019b) or legal cases. Understanding legal reasoning as conflictual practical argumentation, the plural interpretation of these concrete legal rules requires a complementary abstract ethico-legal value ontology, for example, as a “discursive grammar” of justification (Lomfeld 2019).

The contributions of this paper are manifold: After summarizing Lomfeld’s value ontology, resp. “discursive grammar” of justification, in §2, and briefly depicting the *LogiKey* development and knowledge engineering methodology in §3, we outline our object logic of choice – the preference logic by van Benthem et al. (2009) – in §4. In this section we then present, as our first main contribution, a semantical embedding of this preference logic in the *Isabelle/HOL* proof assistant (Nipkow et al. 2002), and we test and assess its meta-theory with state-of-the-art ATP technology. Subsequently we encode in §5, as our second main contribution, Lomfeld’s value ontology, which is again extensively tested and assessed with ATP systems. Further relevant legal and world knowledge is modeled and added in §6, and then it is demonstrated in §7 how the introduced and formalized (general) knowledge can be used for automatically generating value-oriented computer explanations for exemplary property law (“wild animal”) cases. §8 concludes the paper.

Formalization sources are available at <http://logikey.org> under “Preference-Logics/vanBenthemEtAl2009.”

## 2 “Discursive Grammar” of Justification

Combining the discourse theoretical idea that legal reasoning is practical argumentation with a two-level model of legal norms, legal “rules” could be reconstructed as conditional preference relations between conflicting underlying value “principles” (Alexy 2000; Lomfeld 2015). The legal consequence of a rule  $R$  implies the value preference of value principle  $A$  over value principle  $B$ :  $A > B$  (e.g. health security outweighs freedom to move).<sup>2</sup> This value preference applies under the condition that the norm elements  $E_1$  and  $E_2$  which the rule denotes as prerequisites are met. Thus, if  $E_1$  and  $E_2$  are satisfied (e.g. a virus pandemic occurs and voluntary shut down fails), then  $A > B$  which justifies the rule’s legal consequence (e.g. sanctioned lock-down). We thus have a *conditional preference relation*

$$R : (E_1 \wedge E_2) \rightarrow A > B$$

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 could reconstruct a general dialectical taxonomy of legal value

<sup>2</sup>In §5 these values will be assigned to particular parties/actors, i.e. ruling for different parties may promote different values.

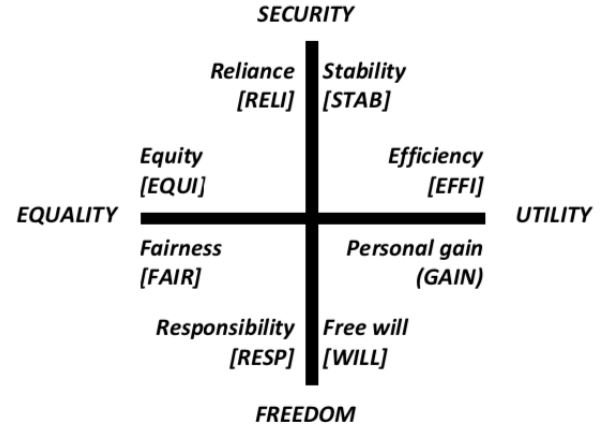


Figure 1: Value ontology by (Lomfeld 2019)

principles used in (at least Western) legislation, legislative materials, cases, textbooks and scholar writings (Lomfeld 2015). The idea is to provide a consistent systematic value ontology of legal principles independent of concrete cases or legal fields to justify legal decisions.

The proposed legal value ontology, see Fig. 1, is consistent with similar taxonomies of political and economic values (Lomfeld 2019). In all these social systems one could observe a general antinomy between individual and collective values. Ideal types of this basic dialectic are the value of FREEDOM for the individual and the value of SECURITY for the collective perspective. Another classic social value antinomy is between a functional (economic) and a more idealistic viewpoint, represented in the ethical debate by the value dialectic of UTILITY versus EQUALITY. These four normative poles stretch an axis of value coordinates for the general set construction.

Within this dialectical matrix eight more concrete legal values/principles are identified. FREEDOM represents the normative value of individual autonomy and comprises the legal principles of more functional individual choice or ‘free will’ (WILL) and more idealistic (self-) ‘responsibility’ (RESP). The 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-maximization, i.e. ‘efficiency’ (EFFI) and individual welfare-maximization, 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 ontology could consistently cover existing value sets from law & logic (or AI & Law) accounts on value-based reasoning, e.g. (Berman and Hafner 1993; Bench-Capon 2012; Gordon and Walton 2012; Sartor 2010), mostly modeled in connection with so-called “wild animal cases” or succeeding common law property cases.

### 3 Methodology

The *LogiKey* methodology (Benzmüller et al. 2020) refers to the principles underlying the organization and the conduct of a complex knowledge design and engineering process—which is what we are faced with. Design means the depiction of the main features of the system we want to achieve, and (knowledge or logic) engineering refers to all the technical and scientific aspects involved in building, maintaining and using a knowledge-based, resp. logic-based, system.

*LogiKey*’s unifying formal framework is based on (shallow) semantical embeddings (SSEs) of ‘object’ logics into a higher-order meta-logic (HOL), which enables the provision of powerful tool support: off-the-shelf theorem provers and model finders for HOL are assisting the *LogiKey* knowledge engineer to *flexibly experiment* with underlying logics and their combinations, with (legal & world) general and domain knowledge, and with concrete use cases—all at the same time. Continuous improvements of these off-the-shelf provers, without further ado, leverage the reasoning performance in *LogiKey*.

In our context, the knowledge engineering task is facing the layers as depicted in Fig. 2: meta-logic HOL enables the encoding of a preference logic, followed by the value ontology, both of which are then used for formalizing legal & world knowledge, before concrete wild animal use cases are modeled and assessed.

The engineering process at these layers has backtracking points and several work cycles may be required; thereby the higher layers may also pose modification requests to the lower layers. Such requests may, unlike in most other approaches, also include far-reaching modifications of the chosen logical foundations, e.g. in the adopted particular preference logic. The work we present in this paper is in fact the result of an iterative, give-and-take process encompassing several cycles of modeling, assessment and testing activities, whereby a (modular) logical theory gradually evolves until eventually reaching a state of highest coherence and acceptability. We then speak of arriving at a *reflective equilibrium*, cf. previous work on *computational hermeneutics* (Fuenmayor and Benzmüller 2019a).

*LogiKey* thus supports empirical studies on ethico-legal theories in which the underlying logic formalisms itself can be flexibly varied, assessed and compared in context.

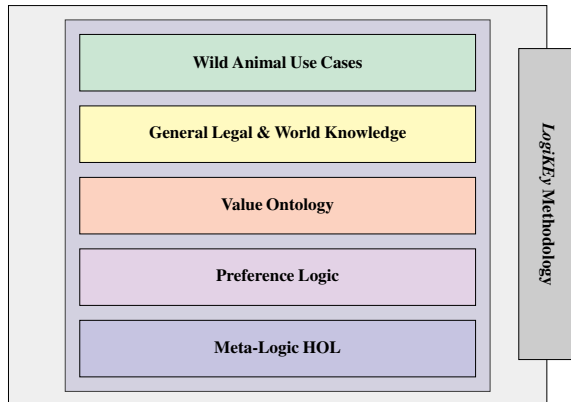


Figure 2: *LogiKey* development methodology

### 4 Preference Logic

Adopting the *LogiKey* methodology, the first question to address is: Which preference logic to choose (initially)? Preference logics generally aim at an adequate modeling of preference reasoning of an (idealized) rational individual. Hence, the logic we choose needs to provably support relevant criteria (e.g., strengthening/weakening the antecedent/consequent). Our initial choice has been the preference logic by van Benthem et al. (2009), abbreviated by  $\mathcal{PL}$  in the remainder:  $\mathcal{PL}$  is a “modal logic for *ceteris paribus* preferences understood 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 and has returned in computer science in the 1990’s and in more abstract ‘dependency logics’ today.”  $\mathcal{PL}$  appears suited for effective automation using the SSEs approach, which has been another selection criterion. This judgment is based on good prior experience with SSEs of related logics whose semantics employs accessibility relations between possible worlds/states, just as  $\mathcal{PL}$  does. Interesting alternatives to  $\mathcal{PL}$  include the preference logics by Liu (2008) and Osherson and Weinstein (2012).

$\mathcal{PL}$  may be criticized as an overly complex choice for our aims as presented in §7, where the automation of *ceteris paribus* reasoning has hardly been addressed yet (but will be in future work). Anyhow, we may always ‘downgrade’ our choice of logic if that is needed by e.g. practical performance considerations, or we may ‘upgrade’ or combine it with other logics if required. Remember that the choice of formalization logic is a parameter in the *LogiKey* approach.

#### Automating Preference Logic $\mathcal{PL}$

For the mechanization and automation of  $\mathcal{PL}$  we utilize the SSE technique, which encodes the language constituents of an object logic,  $\mathcal{PL}$  in our case, as expressions ( $\lambda$ -terms) in HOL. This shows that  $\mathcal{PL}$  is just a fragment of HOL and can be automated as such. A core idea is to model (relevant parts of) the semantic structures of  $\mathcal{PL}$  explicitly in HOL. For  $\mathcal{PL}$  these structures are relational frames constituted by sets of possible worlds/states and their accessibility relations.  $\mathcal{PL}$  propositions can thus be encoded in the SSE as predicates in HOL taking possible worlds/states as arguments.

The SSE of the basic operators of preference logic  $\mathcal{PL}$  in HOL is depicted in Fig. 3; this encoding is further explained below. Further extensions to support *ceteris paribus* reasoning in  $\mathcal{PL}$  are presented in Fig. 9 in Appx. A.1.

As a result we obtain a combined, interactive and automated, theorem prover and model finder for  $\mathcal{PL}$  realized within *Isabelle/HOL*. This is a new contribution, since we are not aware of any other existing implementation and automation of  $\mathcal{PL}$ . 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. A formal proof of the faithfulness (soundness & completeness) of the SSE is analogous to numerous prior works, cf. Benzmüller (2019) and the references therein.

```

1 theory PreferenceLogicBasics (** Benzmüller & Fuenmayor, 2020 **)
2 imports Main
3 begin (** SSE of prefer. logic by van Benthem et al., JPL 2009 **)
4 (*unimportant*)declare[[syntax_ambiguity_warning=false]]
5 nitpick_params[user_axioms,expect=genuine,show_all,format=3]
6 (*preliminaries*)
7 typedef i (*possible worlds*)
8 type_synonym  $\sigma = i \Rightarrow \text{bool}$  (*'world-lifted' propositions*)
9 type_synonym  $\gamma = i \Rightarrow i \Rightarrow \text{bool}$  (*preference relations*)
10 type_synonym  $\mu = \sigma \Rightarrow \sigma$  (*unary logical connectives*)
11 type_synonym  $\nu = \sigma \Rightarrow \sigma \Rightarrow \sigma$  (*binary logical connectives*)
12 type_synonym  $\pi = \sigma \Rightarrow \text{bool}$  (*sets of world-lifted propositions*)
13 (*betterness relation  $\preceq$  and strict betterness relation  $\prec$ *)
14 consts BR:: $\gamma$  ("<_")
15 abbreviation SBR:: $\gamma$  ("<_") where " $\gamma \preceq \equiv (\gamma \preceq) \wedge \neg(\gamma \prec)$ "
16 abbreviation "reflexive R  $\equiv \forall x. R x x$ "
17 abbreviation "transitive R  $\equiv \forall x y z. R x y \wedge R y z \longrightarrow R x z$ "
18 abbreviation "is_total R  $\equiv \forall x y. R x y \vee R y x$ "
19 axiomatization where rBR: "reflexive BR" and tBR: "transitive BR"
20 (*modal logic connectives (operating on truth-sets)*)
21 abbreviation c1:: $\sigma$  ("⊥") where "⊥  $\equiv \lambda w. \text{False}$ "
22 abbreviation c2:: $\sigma$  ("⊤") where "⊤  $\equiv \lambda w. \text{True}$ "
23 abbreviation c3:: $\mu$  ("¬") where "¬  $\equiv \lambda w. \neg(\varphi w)$ "
24 abbreviation c4:: $\nu$  ("infixl" ^ "85") where " $\varphi \wedge \psi \equiv \lambda w. (\varphi w) \wedge (\psi w)$ "
25 abbreviation c5:: $\nu$  ("infixl" ∨ "83") where " $\varphi \vee \psi \equiv \lambda w. (\varphi w) \vee (\psi w)$ "
26 abbreviation c6:: $\nu$  ("infixl" → "84") where " $\varphi \rightarrow \psi \equiv \lambda w. (\varphi w) \longrightarrow (\psi w)$ "
27 abbreviation c7:: $\nu$  ("infixl" ↔ "84") where " $\varphi \leftrightarrow \psi \equiv \lambda w. (\varphi w) \longleftrightarrow (\psi w)$ "
28 abbreviation c8:: $\mu$  ("□") where "□  $\equiv \lambda w. \forall v. (w \preceq v) \longrightarrow (\varphi v)$ "
29 abbreviation c9:: $\mu$  ("◇") where "◇  $\equiv \lambda w. \exists v. (w \preceq v) \wedge (\varphi v)$ "
30 abbreviation c10:: $\mu$  ("□¬") where "□¬  $\equiv \lambda w. \forall v. (w \prec v) \longrightarrow (\varphi v)$ "
31 abbreviation c11:: $\mu$  ("◇¬") where "◇¬  $\equiv \lambda w. \exists v. (w \prec v) \wedge (\varphi v)$ "
32 abbreviation c12:: $\mu$  ("E") where "E  $\equiv \lambda w. \exists v. (\varphi v)$ "
33 abbreviation c13:: $\mu$  ("A") where "A  $\equiv \lambda w. \forall v. (\varphi v)$ "
34 (*meta-logical predicate for global and validity*)
35 abbreviation gl:: $\pi$  ("⊧") where "⊧  $\equiv \forall w. \psi w$ "
36 (*some tests: dualities*)
37 lemma "[ (◇¬) ↔ (¬□) ] ∧ [ (◇¬) ↔ (¬□¬) ] ∧ [ (A) ↔ (¬E) ]" by blast (*proof*)
38 (**** Section 3: A basic modal preference language ****)
39 (*Definition 5*)
40 abbreviation p1:: $\nu$  ("<_EE")
41 where " $(\varphi \preceq_{EE} \psi) \equiv \exists s. \exists t. \varphi s \wedge \psi t \wedge s \preceq t$ "
42 abbreviation p2:: $\nu$  ("<_AE")
43 where " $(\varphi \preceq_{AE} \psi) \equiv \forall s. \exists t. \varphi s \longrightarrow (\psi t \wedge s \preceq t)$ "
44 abbreviation p3:: $\nu$  ("<_EE")
45 where " $(\varphi \preceq_{EE} \psi) \equiv \exists s. \exists t. \varphi s \wedge \psi t \wedge s \preceq t$ "
46 abbreviation p4:: $\nu$  ("<_AE")
47 where " $(\varphi \preceq_{AE} \psi) \equiv \forall s. \exists t. \varphi s \longrightarrow (\psi t \wedge s \preceq t)$ "
48 abbreviation p5:: $\nu$  ("<_AA")
49 where " $(\varphi \preceq_{AA} \psi) \equiv \forall s. \forall t. (\varphi s \wedge \psi t) \longrightarrow s \preceq t$ "
50 abbreviation p6:: $\nu$  (">_EA")
51 where " $(\varphi \succ_{EA} \psi) \equiv \exists s. \forall t. (\varphi s \wedge \psi t) \longrightarrow t \prec s$ "
52 abbreviation p7:: $\nu$  ("<_AA")
53 where " $(\varphi \preceq_{AA} \psi) \equiv \forall s. \forall t. (\varphi s \wedge \psi t) \longrightarrow s \preceq t$ "
54 abbreviation p8:: $\nu$  (">_EA")
55 where " $(\varphi \succ_{EA} \psi) \equiv \exists s. \forall t. (\varphi s \wedge \psi t) \longrightarrow t \prec s$ "
56 abbreviation P1:: $\nu$  ("<_EE") where " $\varphi \preceq_{EE} \psi \equiv E(\varphi \wedge \Diamond \neg \psi)$ "
57 abbreviation P2:: $\nu$  ("<_AE") where " $\varphi \preceq_{AE} \psi \equiv A(\varphi \rightarrow \Diamond \neg \psi)$ "
58 abbreviation P3:: $\nu$  ("<_EE") where " $\varphi \preceq_{EE} \psi \equiv E(\varphi \wedge \Diamond \neg \psi)$ "
59 abbreviation P4:: $\nu$  ("<_AE") where " $\varphi \preceq_{AE} \psi \equiv A(\varphi \rightarrow \Diamond \neg \psi)$ "
60 abbreviation P5:: $\nu$  ("<_AA") where " $\varphi \preceq_{AA} \psi \equiv A(\varphi \rightarrow \Diamond \neg \psi)$ "
61 abbreviation P6:: $\nu$  (">_EA") where " $\varphi \succ_{EA} \psi \equiv E(\varphi \wedge \Diamond \neg \psi)$ "
62 abbreviation P7:: $\nu$  ("<_AA") where " $\varphi \preceq_{AA} \psi \equiv A(\varphi \rightarrow \Diamond \neg \psi)$ "
63 abbreviation P8:: $\nu$  (">_EA") where " $\varphi \succ_{EA} \psi \equiv E(\varphi \wedge \Diamond \neg \psi)$ "
64 (*quantification for objects of arbitrary type*)
65 abbreviation mforall ("∀") where "∀  $\equiv \lambda w. \forall x. (\Phi x w)$ "
66 abbreviation mforallB (binder"∀"[8]9) where "∀x.  $\varphi(x) \equiv \forall \varphi$ "
67 abbreviation mexists ("∃") where "∃  $\equiv \lambda w. \exists x. (\Phi x w)$ "
68 abbreviation mexistsB (binder"∃"[8]9) where "∃x.  $\varphi(x) \equiv \exists \varphi$ "
69 (*polymorph operators for sets of worlds/values*)
70 abbreviation subs (infix "⊆" 70) where " $A \subseteq B \equiv \forall x. A x \longrightarrow B x$ "
71 abbreviation union (infixr "∪" 70) where " $A \cup B \equiv \lambda x. A x \vee B x$ "
72 abbreviation inters (infixr "∩" 70) where " $A \cap B \equiv \lambda x. A x \wedge B x$ "
73 (*Consistency confirmed (trivial: only abbreviations introduced)*)
74 lemma True nitpick[satisfy,user_axioms] oops
75 end

```

Figure 3: SSE of  $\mathcal{PL}$  (Benthem et al. 2009) in HOL

**SSE of  $\mathcal{PL}$  in HOL** We comment on the implementation of the SSE of  $\mathcal{PL}$  in *Isabelle/HOL* as displayed in Fig. 3; see

van Benthem et al. (2009) for further details on  $\mathcal{PL}$ .

First, a new base type  $i$  is declared (Line 7), denoting the set of possible worlds or states. Subsequently (Lines 8–12), useful type abbreviations are introduced, including the type  $\sigma$  for  $\mathcal{PL}$  propositions, which are modeled as predicates on objects of type  $i$  (i.e. as *truth-sets* of worlds/states). A *betterness relation*  $\preceq$ , and its strict variant  $\prec$ , are introduced (Lines 14–15), with  $\preceq$ -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  $\preceq$  is a preorder, i.e. reflexive and transitive (Lines 16–19).

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  $\Box \preceq$  and  $\Box \prec$  use  $\preceq$  and  $\prec$  as guards in their definitions (Lines 30–31). 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 (Lines 37–38).

Subsequently, the *betterness* ordering  $\preceq$  (resp.  $\prec$ ) 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 41–56). These are complemented by eight syntactic definitions of the same binary preferences stated within the object language  $\mathcal{PL}$  (Lines 57–64).<sup>3</sup> Intuitively, preferring proposition  $P$  over  $Q$  amounts to preferring every P-state over every Q-state ( $\prec_{AA}$ ), or at least one P-state over every Q-state ( $\succ_{EA}$  or  $\prec_{AE}$ ). Each of these non-trivial variants can be argued for, cf. (Benthem et al. 2009). However, only  $\succ_{EA}$  and  $\prec_{AE}$  satisfy the conditions required for a logic of value aggregation (see the discussion in §5). Transitivity, a quite controversial property in the literature on preferences, is satisfied by  $\prec_{AE}$  but not by  $\succ_{EA}$ . Our framework can support both modeling options. In line with the *LogiKey* methodology, such a choice is left at the discretion of the modeler.

We further extend the propositional logic  $\mathcal{PL}$  as defined by van Benthem et al. (2009) by adding quantifiers (Lines 66–69).<sup>4</sup> Moreover, polymorphic operators for subset, union and intersection are defined (Lines 71–73).

The model finder *Nitpick* (Blanchette and Nipkow 2010) confirms the consistency of the introduced theory (Line 75) 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. 10 and 11 in Appx. A.2.

## 5 Encoding the Value Ontology

The next, essential step is to model and encode the “Discursive Grammar” value ontology. The current status of

<sup>3</sup>ATP systems can prove the meta-theoretic correspondences of these definitions; see Lines 5–13 in Fig. 10 in Appx. A.2.

<sup>4</sup>See Benzmüller and Paulson (2013) for a discussion of the SSE of quantified modal logics.



```

1 theory ValueOntology (** Benzmüller, Fuenmayor & Lomfeld, 2020 **)
2 imports PreferenceLogicBasics
3 begin (** Lomfeld's value ontology is encoded **)
4
5 (*two legal parties (there can be more in principle)*)
6 datatype c = p | d (*parties/contenders: plaintiff, defendant*)
7 fun other::"c⇒c" ("_" ) where "p-1 = d" | "d-1 = p"
8
9 consts For::"c⇒σ" (*decision: find/rule for party*)
10 axiomatization where ForAx: "[For x ↔ (¬For x-1)]"
11
12 datatype (*ethico-legal upper values (wrt. a given party)*)
13 't VAL = FREEDOM 't | UTILITY 't | SECURITY 't | EQUALITY 't
14 type_synonym v = "(c)VAL⇒bool" (*principles: sets of upper values*)
15 type_synonym cv = "c⇒v" (*principles are specified wrt. a given party*)
16
17 abbreviation vset1 ("[") where "[x] ≡ λx::(c). VAL. x=x"
18 abbreviation vset2 ("[") where "[α,β] ≡ λx::(c)VAL. x=α ∨ x=β"
19
20 abbreviation utility::cv ("UTILITY-" ) where "UTILITYx ≡ {UTILITY x}"
21 abbreviation security::cv ("SECURITY-" ) where "SECURITYx ≡ {SECURITY x}"
22 abbreviation equality::cv ("EQUALITY-" ) where "EQUALITYx ≡ {EQUALITY x}"
23 abbreviation freedom::cv ("FREEDOM-" ) where "FREEDOMx ≡ {FREEDOM x}"
24 abbreviation stab::cv ("STAB-" ) where "STABx ≡ {SECURITY x, UTILITY x}"
25 abbreviation effi::cv ("EFFI-" ) where "EFFIx ≡ {UTILITY x, SECURITY x}"
26 abbreviation gain::cv ("GAIN-" ) where "GAINx ≡ {UTILITY x, FREEDOM x}"
27 abbreviation will::cv ("WILL-" ) where "WILLx ≡ {FREEDOM x, UTILITY x}"
28 abbreviation resp::cv ("RESP-" ) where "RESPx ≡ {FREEDOM x, EQUALITY x}"
29 abbreviation fair::cv ("FAIR-" ) where "FAIRx ≡ {EQUALITY x, FREEDOM x}"
30 abbreviation equi::cv ("EQUI-" ) where "EQUIx ≡ {EQUALITY x, SECURITY x}"
31 abbreviation reli::cv ("RELI-" ) where "RELIx ≡ {SECURITY x, EQUALITY x}"
32
33 (*derivation operators (cf. theory of "formal concept analysis" *)
34 consts Vrel::"i⇒(c)VAL⇒bool" ("I") (*incidence relation worlds-values*)
35 abbreviation intension::"σ⇒v" ("↑") where "W ≡ λv. ∀x. W x → I x v"
36 abbreviation extension::"v⇒σ" ("↓") where "V ≡ λw. ∀x. V x → I w x"
37
38 (*shorthand notation for aggregating values*)
39 abbreviation agg ("infixr" "@80" ) where "v1@v2 ≡ v1 ∩ v2"
40 abbreviation agg1 ("[") where "[v] ≡ v"
41 abbreviation agg2 ("["_@_"]") where "[v1@v2] ≡ (v1@v2)"
42 abbreviation agg3 ("["_@_@_"]") where "[v1@v2@v3] ≡ (v1@v2@v3)"
43 abbreviation agg4 ("["_@_@_@_"]") where "[v1@v2@v3@v4] ≡ (v1@v2@v3@v4)"
44
45 (*chosen variant for preference relation (cf. van Benthem et al. 2009*)
46 abbreviation relPref::"σ⇒σ" ("<_") where "φ < ψ ≡ ψ >EA φ"
47 abbreviation relPrefval::"v⇒v" ("<_v") where "φ <v ψ ≡ ψ >EA φ"
48
49 abbreviation incnst ("INCONS-" ) where (*inconsistency for value support*)
50 "INCONSx ≡ {SECURITYx} ∩ {EQUALITYx} ∩ {FREEDOMx} ∩ {UTILITYx}"
51
52 lemma "True" nitpick[satisfy] oops (*verify consistency of this theory*)
53 end

```

Figure 4: Encoding of the Value Ontology

our modeling efforts is illustrated in the *Isabelle/HOL* theory displayed in Fig. 4. Similar to before, this modeling may be subject to further modifications. In fact, what is presented here is the result of several cycles of modeling, encoding and testing activities, in which previous versions have been refuted and rejected or further improved.

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 6–7); and a predicate modeling the ruling for a party is also provided (Lines 9–10).

As regards to Lomfeld’s upper value ontology, a four-valued (parameterized) datatype is introduced (Lines 12–13) as described in §2. Moreover, type-aliases (Lines 14–15) and set-constructor operators for values (Lines 17–18) are introduced for ease of presentation. The legal principles from §2 are next introduced as combinations of those upper values (Lines 20–31). As an illustration, the principle *stability* (STAB) is understood as instantiation of the upper value

```

1 theory ValueOntologyTest (** Benzmüller, Fuenmayor & Lomfeld, 2020 **)
2 imports ValueOntology
3 begin (* value ontology tests *)
4 (*values in two opposed quadrants: inconsistent*)
5 lemma "[RESPx] ∩ [STABx] → INCONSx" by simp
6 lemma "[RELIx] ∩ [WILLx] → INCONSx" by simp
7 (*all values in two non-opposed quadrants: consistent*)
8 lemma "[WILLx] ∩ [STABx] → INCONSx" nitpick oops (*countermodel*)
9 (*values in opposed quadrants for different parties: consistent*)
10 lemma "[EQUIx] ∩ [GAINy] → (INCONSx ∨ INCONSy)" nitpick oops (*ctm*)
11 lemma "[RESPx] ∩ [STABy] → (INCONSx ∨ INCONSy)" nitpick oops (*ctm*)
12 lemma "[RELIx] ∩ [WILLy] nitpick[satisfy] nitpick oops (*contingent*)
13 (*value preferences tests*)
14 lemma "[WILLx <v STABx] → [WILLx <v RELIx@STABx]" by blast
15 lemma "[RELIx@STABx <v WILLx] → [STABx <v WILLx]" by auto
16 lemma "[WILLx <v RELIx@STABx] → [WILLx <v STABx]"
17 nitpick nitpick[satisfy] oops (*contingent*)
18 lemma "[STABx <v WILLx] → [RELIx@STABx <v WILLx]"
19 nitpick nitpick[satisfy] oops (*contingent*)
20 end

```

Figure 5: Testing the Value Ontology

*security* coalesced with *utility* (in that order of importance).<sup>5</sup>

After defining legal principles as combinations (in this case: sets) of values (w.r.t. a legal party), we need to relate them to propositions (sets of worlds/states) in our logic  $\mathcal{PL}$ . For this we reuse some basic notions from the theory of *Formal Concept Analysis* (FCA) (Ganter and Wille 2012). We define a binary *incidence relation* (Line 34) between worlds/states (type  $i$ ) and upper values (type  $VAL$ ), which is used to link propositions (sets of worlds) with corresponding sets of values (their *intension*) by means of the operator  $\uparrow$  (Line 35). More important to our purposes: each legal principle (set of values) is associated with a proposition (its *extension*) by means of the operator  $\downarrow$  (Line 36).<sup>6</sup>

Intuitively, we can read the proposition denoted by  $STAB^p \downarrow$  as: “a decision in favor of the plaintiff ( $p$ ) promotes the principle of stability”. Value-inconsistency (INCONS) is given when all four upper values apply for the same legal party (Lines 49–50). Further shorthand notation for value aggregation is introduced (Lines 39–43).

The preference relation  $<$  utilized in the remainder is chosen to be the  $>_{EA}$ -variant from §4 (Line 46–47). As mentioned before, the (better-known, transitive)  $<_{AE}$ -variant is also suitable for our logic of value preferences. The results presented are valid for both (unless otherwise stated).

## Testing the Value Ontology

In order to test the adequacy of our modeling, some implied and non-implied knowledge is studied. This is shown in Fig. 5. We briefly discuss some of the conducted tests.

Following a conflictual dialectical interpretation of the value ontology, promoting values (for the same party) from two opposing quadrants, say RESP & STAB or RELI & WILL (cf. Fig. 1), is provably ‘value-inconsistent’ (INCONS); theorem provers quickly confirm this (Lines 5–6). However, promoting values from two non-opposed quad-

<sup>5</sup> At this early stage of our modeling we have not yet considered this additional factor of relative importance or weight between values. Doing this would allow us to clearly distinguish, e.g., *stability* (STAB) from *efficiency* (EFFI) among others.

<sup>6</sup> FCA exhibits a perfect duality between the pairs of notions: *object/attribute* and *extension/intension*. Our choice to consider worlds as *objects* and values as their *attributes* (i.e. propositions as *extensions* of legal principles) bears no mathematical significance.

```

Nitpick found a model for card i = 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 = (λx. _)((i1, i1) := True)
  For = (λx. _)((d, i1) := False, (p, i1) := True)
  I = (λx. _)
    ((i1, FREEDOM d) := False,
     (i1, FREEDOM p) := True,
     (i1, UTILITY d) := False,
     (i1, UTILITY p) := True,
     (i1, EQUALITY d) := False,
     (i1, EQUALITY p) := True,
     (i1, SECURITY d) := False,
     (i1, SECURITY p) := True)
  other = (λx. _)(d := p, p := d)

```

Figure 6: Satisfying model for the statement in Line 12 of Fig. 5.

rants, such as WILL & STAB (Line 8) is value-consistent: the model finder *Nitpick* computes and reports a counter-model (not shown here) to the stated conjecture. Value-inconsistency is also not implied if values from opposing quadrants are promoted for different parties (Lines 10–11).

Remark on model finder *Nitpick* (Blanchette and Nipkow 2010): *Nitpick* 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’. In Line 12 of Fig. 5, for example, *Nitpick* is called simultaneously in both modes in order to confirm the contingency of the statement; as expected both a model (Fig. 6) and countermodel for the statement is returned. These models are given as concrete interpretations of relevant terms in the given context so that the conjectured statement is satisfied or falsified (depending of the mode of *Nitpick*). Studying and analysing the returned interpretations can be very helpful and intuition-fostering for the user, in particular, in very complex modeling tasks when intuition and deep understanding of the modeled structures is initially lacking.

Note that the notion of value-inconsistency (INCONS), cf. Lines 49–50 in Fig. 4, has deliberately not been aligned with inconsistency in meta-logic HOL so that it can be explicitly reasoned with (avoiding explosion). That means, that we can contingently promote conflicting values such as RELI and WILL for  $p$ , and this then leads to value-inconsistency (for  $p$ ) but not to general inconsistency in meta-logic HOL. In Line 12 this contingency is confirmed. However, inspecting the satisfying models generated by *Nitpick*, we see that  $p$  is indeed value-inconsistent. One of such models is depicted in Fig. 6, where it is shown that (in the given possible world  $i_1$ ) the upper values EQUALITY, SECURITY, UTILITY and FREEDOM are simultaneously promoted for  $p$ , this amounts to value-inconsistency (INCONS) 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 then triggers adaptations of the defining axioms. Such a process may require several cycles (remember the discussion from §3) and, as a useful side effect, it clearly fosters cross-disciplinary mutual understanding.

Further tests in Fig. 5 assess the suitability of the pref-

erence relation  $\prec_v$  (either  $\prec_{AE}$  or  $\succ_{EA}$ ) for reasoning e.g. with value aggregations. For example, we test for a correct behavior 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 14). Similar test are conducted for ‘weakening’ of the left-hand side.

## 6 General Legal & World Knowledge

The realistic modeling of concrete legal cases requires further legal & world knowledge (LWK) to be taken into account. LWK is typically modeled 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 present in Fig. 7 only a small and monolithic example theory called “GeneralKnowledge” for illustration. In our case this includes a small excerpt of a much simplified “animal appropriation taxonomy”. Moreover, we associate “animal appropriation” (kinds of) situations with the value preferences they imply (conditional preference relations).

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

Next steps include interrelating certain notions in our theory “GeneralKnowledge” with values, resp. value preferences, as introduced in the previous sections. It is here where the preference relations and modal operators of  $\mathcal{PL}$  are used. 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 a preference logic, or we may replace our current material implication operator by a conditional implication to better support defeasible legal reasoning.<sup>7</sup>

A general problem of knowledge representation frameworks and ontology languages often is their insufficient support for the proper treatment of modalities, see e.g. the discussion and solution proposed by Benz Müller and Pease (2012). Our framework, which is based on that proposed solution, is not suffering from such restrictions.

We briefly explain below the *Isabelle/HOL* encoding of our LWK as shown in Fig. 7.

First, some non-logical constants that stand for kinds of legally relevant situations (here: of appropriation) are introduced (Lines 6–9), and their meaning is constrained by some postulates (Lines 12–17). As already discussed, for ease of illustration, these terms are modeled here as simple propositional constants. In practice, however, they may later be replaced, or logically implied, by a more realistic and suitable

<sup>7</sup>Indeed a (cautiously-monotonic) conditional implication can be defined employing  $\mathcal{PL}$  modal operators. We may also opt for an SSE of a conditional logic in HOL, see e.g. Benz Müller (2013).

```

1 theory GeneralKnowledge (*Benzmüller, Fuenmayor & Lomfeld, 2020*)
2   imports ValueOntology
3   begin (** General Legal and World Knowledge (LWK) **)
4
5   (*LWK: kinds of situations addressed*)
6   consts appObject::"σ" (*appropriation of objects in general*)
7   consts appAnimal::"σ" (*appropriation of animals in general*)
8   consts appWildAnimal::"σ" (*appropriation of wild animals*)
9   consts appDomAnimal::"σ" (*appropriation of domestic animals*)
10
11   (*LWK: postulates for kinds of situations*)
12   axiomatization where
13   W1: "[appWildAnimal ∨ appDomAnimal] ↔ appAnimal]" and
14   W2: "[appWildAnimal ↔ ¬appDomAnimal]" and
15   W3: "[appWildAnimal → appAnimal]" and
16   W4: "[appDomAnimal → appAnimal]" and
17   W5: "[appAnimal → appObject]"
18   (*...further situations regarding appropriation of objects, etc.*)
19
20   (*LWK: (prima facie) value preferences for kinds of situations*)
21   axiomatization where
22   R1: "[appAnimal → (STABp <v STABd)]" and
23   R2: "[appWildAnimal → (WILLx-1 <v STABx)]" and
24   R3: "[appDomAnimal → (STABx-1 <v RELI*⊕RESPx)]"
25   (*...further preferences...*)
26
27   (*LWK: domain vocabulary*)
28   typedef e (*declares new type for 'entities'*)
29   consts Animal::"e⇒σ"
30   consts Domestic::"e⇒σ"
31   consts Fox::"e⇒σ"
32   consts Parrot::"e⇒σ"
33   consts Pet::"e⇒σ"
34   consts FreeRoaming::"e⇒σ"
35
36   (*LWK: taxonomic (domain) knowledge*)
37   axiomatization where
38   W6: "[∀a. Fox a → Animal a]" and
39   W7: "[∀a. Parrot a → Animal a]" and
40   W8: "[∀a. (Animal a ∧ FreeRoaming a ∧ ¬Pet a) → ¬Domestic a]"
41   (*...others...*)
42
43   (*LWK: legally-relevant, situational 'factors'*)
44   consts Own::"c⇒σ" (*object is owned by party c*)
45   consts Poss::"c⇒σ" (*party c has actual possession of object*)
46   consts Intent::"c⇒σ" (*party c has intention to possess object*)
47   consts Mal::"c⇒σ" (*party c acts out of malice*)
48   consts Mtn::"c⇒σ" (*party c respons. for maintenance of object*)
49
50   (*LWK: meaning postulates for general notions*)
51   axiomatization where
52   W9: "[Poss x → (¬Poss x-1)]" and
53   W10: "[Own x → (¬Own x-1)]"
54   (*...others...*)
55
56   (*LWK: conditional value preferences, e.g. from precedents*)
57   axiomatization where
58   R4: "[Mal x-1 ∧ Own x → (STABx-1 <v RESP*⊕RELIx)]"
59   (*...others...*)
60
61   (*LWK: relate values, outcomes and situational 'factors'*)
62   axiomatization where
63   F1: "[For x → (Intent x ↔ □:[WILLx])]" and
64   F2: "[For x → (Mal x-1 ↔ □:[RESPx])]" and
65   F3: "[For x → (Poss x ↔ □:[STABx])]" and
66   F4: "[For x → (Mtn x ↔ □:[RESPx])]" and
67   F5: "[For x → (Own x ↔ □:[RELIx])]"
68
69   (*theory is consistent, (non-trivial) model found*)
70   lemma True nitpick[satisfy,card i=10] oops
71 end

```

Figure 7: Encoding of general legal & world knowledge

modeling utilizing arbitrarily complex (even higher-order, if needed) formulas depicting states of affairs to some desired level of granularity. Some simple vocabulary and taxonomic relationships (here: for wild and domestic animals) are spec-

ified (Lines 28–40) to illustrate this.

The legal default rules for several situations (here: appropriation of animals) are formulated as *prima facie* preference relations (Lines 21–24). For example, one rule R2 (Line 23) could be read as: “In a wild-animals-appropriation kind of situation, a decision promoting STAB in favor of a party (say, the plaintiff) is preferred over a decision, favoring the other party (defendant), which promotes WILL”. If there is no more specific legal rule from a precedent or a codified statute then these *prima facie* preference relations determine the result. As a didactic example, the legal rule R4 (Line 58) states that the ownership (say, the plaintiff’s) of the land on which the hunting 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*, cf. §2. This reflects the Common law precedent of *Keeble v. Hickergill* (1704, 103 ER 1127).

An established AI & Law tool to structure the representation of legal precedents are situational “factors” (Ashley 1990; Prakken and Sartor 2015). Some of them are specified as illustration (Lines 44–48), together with some postulates constraining their meaning (Lines 51–53). Our framework also allows us to introduce definitions for those factors for which clear legal specifications exist. In our model, such factors are further related to value principles and outcomes (Lines 62–67). Our normative assignment here is widely in accordance with classifications in the AI & Law literature (Berman and Hafner 1993; Bench-Capon 2012).

Finally, the consistency of all axioms and rules provided is confirmed by *Nitpick* (Line 70).

## 7 Proof of Concept – Case Study

We illustrate our reasoning framework by encoding the classic common law property case *Pierson v. Post*. In a nutshell:

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

The modeling is presented in Fig. 8. We have introduced and interrelated some minimal vocabulary, e.g. to “pursue” and to “capture” (Lines 6–14), needed to represent the case facts as interpreted by the parties (Lines 19–20 for *Pierson*; Lines 45–46 for *Post*). Their consistency with other postulates from the previously introduced layers “GeneralKnowledge” and “ValueOntology” (imported in Line 2) is verified by generating a (non-trivial) model using *Nitpick* (Line 16).

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* (STAB) and outweighs the mere *will* to possess (WILL) 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 the legal value ontology in §2, this corresponds to a preference for the (upper) value SECURITY over FREEDOM. We can see that this legal rule R2,



as previously introduced in the layer “GeneralKnowledge” (see Fig. 7, Line 23) is indeed employed by *Isabelle/HOL* automated tools to prove (Lines 31–32) that the given facts imply a preference for a decision for Pierson over one for Post. Notice that the legal precedent rule R4 of *Keeble v. Hickergill* (see Fig. 7, Line 58) does not apply to this case.

We also present and model a possible counterargument for 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” as such (Line 45–46). Moreover, an alternative legal rule (i.e. a possible argument for overruling in case law methodology) is presented (Lines 42–43), entailing a value preference of *efficiency* (EFFI) over *stability* (STAB), and justified by the alleged public benefit of hunters getting rid of foxes, since the latter cause depredations in farms. This is the argument put forward by the dissenting opinion in the original case (3 Cai R 175). Other tests, i.e. consistency of both decision alternatives with the given premises and with the value-consistency of both parties (Lines 48–54) and refutability of a possible decision favoring Pierson (Lines 61–62), are analogous (but contrary) to Pierson’s argument.

## 8 Conclusion

Contributions and preliminary results of an ongoing project have been presented in which a “discursive grammar” value ontology in combination with further legal and general world knowledge is being utilized for a detailed, granular assessment and explanation of legal argumentation as partly controversial legal balancing in a concrete case.

From a technical perspective, the core objective of this paper has been to demonstrate that the technology we promote and apply—shallow semantical embeddings (SSEs) in classical higher-order logic (HOL)—appears indeed suitable for the task of structured legal balancing. It is the flexibility of the multi-layer modeling which is unique in our approach, in combination with a very rich support for expressiveness, quantified classical and non-classical logics, thereby rejecting the idea that knowledge representation means should or must be limited a priori to propositional frameworks due only to computational considerations.

From a legal perspective, the reconstruction of legal balancing is already with classical argumentative tools a non-trivial task which is methodologically not settled at all (Sieckmann 2010). Here, our paper proposed the structuring of legal argumentation with a dialectical ethico-legal value ontology. Legal rules and their various interpretations could thus be displayed within a unified pluralistic logic of value preference relations. The integration of the preference logic and the value ontology within the dynamic HOL modeling environment allows to experiment with different forms of interpretation. This enables not only to find more accurate reconstructions of legal argumentation but also supports to aggregate legal arguments and check their normative (value) consistency. We aim at expanding the model towards proportional preference relations which would allow to integrate qualitative and quantitative legal balancing approaches and to evaluate their relative strength.

```

1 theory Pierson (** Benzmüller, Fuenmayor & Lomfeld, 2020 **)
2   imports GeneralKnowledge
3   begin (** Pierson v. Post "wild animal" case **)
4
5   (*case-specific 'world-vocabulary'*)
6   consts α::"e" (*appropriated animal (fox in this case) *)
7   consts Pursue::"c⇒e⇒σ"
8   consts Capture::"c⇒e⇒σ"
9
10  (*case-specific taxonomic (legal domain) knowledge*)
11  axiomatization where
12  CW1: "[ (∃c. Capture c α ∧ ¬Domestic α) → appWildAnimal]" and
13  CW2: "[∀c. Pursue c α → Intent c]" and
14  CW3: "[∀c. Capture c α → Poss c]"
15
16  lemma True nitpick[satisfy,card i=4] oops (*satisfiable*)
17
18  (***** pro-Pierson's argument *****)
19  abbreviation "Pierson_facts ≡ [Fox α ∧ (FreeRoaming α) ∧
20    (¬Pet α) ∧ Pursue p α ∧ (¬Pursue d α) ∧ Capture d α]"
21
22  (*decision for defendant (Pierson) is compatible with premises*)
23  lemma "Pierson_facts ∧ [¬INCONSp] ∧ [¬INCONSd] ∧ [For p < For d]"
24    nitpick[satisfy,card i=4] oops (* (non-trivial) model found*)
25
26  (*decision for plaintiff (Post) is compatible with premises*)
27  lemma "Pierson_facts ∧ [¬INCONSp] ∧ [¬INCONSd] ∧ [For d < For p]"
28    nitpick[satisfy,card i=4] oops (* (non-trivial) model found*)
29
30  (*decision for defendant (Pierson) is provable*)
31  theorem assumes Pierson_facts shows "[For p < For d]"
32    by (metis assms CW1 CW2 W6 W8 ForAx R2 F1 other.simps(2) rBR)
33
34  (*while a decision for the plaintiff is not*)
35  lemma assumes Pierson_facts shows "[For d < For p]"
36    nitpick[card i=4] oops (*counterexample found*)
37
38  (***** pro-Post's argument *****)
39  (* Theory amendment: the animal is not free-roaming since it
40    is being chased by a professional hunter (Post) *)
41  consts Hunter::"c⇒σ"
42  axiomatization where (*case-specific legal rule for hunters*)
43  R5: "[ (Hunter x ∧ Pursue x α) → (STABx-1 <v EFFI)]"
44
45  abbreviation "Post_facts ≡ [Fox α ∧ (¬FreeRoaming α) ∧
46    Hunter p ∧ Pursue p α ∧ (¬Pursue d α) ∧ Capture d α]"
47
48  (*decision for defendant (Pierson) is compatible with premises*)
49  lemma "Post_facts ∧ [¬INCONSp] ∧ [¬INCONSd] ∧ [For p < For d]"
50    nitpick[satisfy,card i=4] oops (* (non-trivial) model found*)
51
52  (*decision for plaintiff (Post) is compatible with premises too*)
53  lemma "Post_facts ∧ [¬INCONSp] ∧ [¬INCONSd] ∧ [For d < For p]"
54    nitpick[satisfy,card i=4] oops (* (non-trivial) model found*)
55
56  (*indeed, a decision for plaintiff (Post) now becomes provable*)
57  theorem assumes Post_facts shows "[For d < For p]"
58    using assms by (metis CW3 ForAx R5 F3 other.simps rBR)
59
60  (*while a decision for the defendant is now refutable*)
61  lemma assumes Post_facts shows "[For p < For d]"
62    nitpick[card i=4] oops (* counterexample found*)
63  end

```

Figure 8: Modeling the Pierson v. Post case.

The suitable, reusable, paradox-free modeling of ethico-legal balancing as presented here is partly an art and partly a challenging engineering task. Combining both aspects in a holistic approach within a dynamic modeling framework, we hope to contribute to the evolution of legal reasoning and to pave the way for some form of (legally) reasonable machines.



**Acknowledgements:** Will be included in the final document.

## References

- 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.
- Ashley, K. D. (1990). *Modelling Legal Argument: Reasoning with Cases and Hypotheticals*. Cambridge/MA: MIT Press.
- Bench-Capon, T. (2012). “Representing Popov v Hayashi with dimensions and factors”. In: *Artificial Intelligence and Law* 20, pp. 15–35.
- Bench-Capon, T., K. Atkinson, and A. Chorley (2005). “Persuasion and value in legal argument”. In: *Journal of Logic and Computation* 15, pp. 1075–1097.
- Bench-Capon, T. and G. Sartor (2003). “A model of legal reasoning with cases incorporating theories and value”. In: *Artificial Intelligence* 150, pp. 97–143.
- Benthem, J. van, 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.
- Benzmüller, C. (2013). “Automating Quantified Conditional Logics in HOL”. In: *IJCAI-13*. Ed. by F. Rossi. AAAI Press, pp. 746–753.
- (2019). “Universal (Meta-)Logical Reasoning: Recent Successes”. In: *Sci. Comp. Progr.* 172, pp. 48–62.
- Benzmüller, C. and B. Lomfeld (2020). “Reasonable Machines: A Research Manifesto”. Preprint: <https://dx.doi.org/10.13140/RG.2.2.28918.63045>.
- 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: *Artificial Intelligence*. DOI: <https://doi.org/10.1016/j.artint.2020.103348>.
- Benzmüller, C. and L. C. Paulson (2013). “Quantified Multimodal Logics in Simple Type Theory”. In: *Logica Universalis* 7.1, pp. 7–20.
- Benzmüller, C. and A. Pease (2012). “Higher-order Aspects and Context in SUMO”. In: *Journal of Web Semantics* 12-13, pp. 104–117.
- 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. and T. Nipkow (2010). “Nitpick: A Counterexample Generator for Higher-Order Logic Based on a Relational Model Finder”. In: *ITP 2010*. Vol. 6172. LNCS. Springer, pp. 131–146.
- Dung, P. M. (1995). “On the acceptability of arguments and its fundamental role in nonmonotonic reasoning, logic programming, and n-person games”. In: *Artificial Intelligence* 77, pp. 321–357.
- Fuenmayor, D. and C. Benzmüller (2019a). “A Computational-Hermeneutic Approach for Conceptual Explication”. In: *Model-Based Reasoning in Science and Technology. Inferential Models for Logic, Language, Cognition and Computation*. Vol. 49. SAPERE. Springer, Cham, pp. 441–469.
- (2019b). “Harnessing Higher-Order (Meta-)Logic to Represent and Reason with Complex Ethical Theories”. In: *PRICAI 2019: Trends in Artificial Intelligence*. Vol. 11670. LNAI. Springer, Cham, pp. 418–432.
- Ganter, B. and R. Wille (2012). *Formal concept analysis: mathematical foundations*. Springer Berlin.
- Gordon, T. and D. Walton (2012). “A Carneades reconstruction of Popov v Hayashi”. In: *Artificial Intelligence and Law* 20, pp. 37–56.
- Hage, J. (1997). *Reasoning With Rules*. Dordrecht: Kluwer.
- 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*. Vol. 188. Frontiers in Artificial Intelligence and Applications. IOS Press, pp. 21–52.
- Horty, J. (2011). “Rules and reasons in the theory of precedent”. In: *Legal Theory* 17, pp. 1–33.
- Liu, F. (2008). “Changing for the Better – Preference Dynamics and Agent Diversity”. PhD thesis. Inst. for Logic, Language and Computation, Universiteit van Amsterdam.
- 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.
- Nipkow, T., L. Paulson, and M. Wenzel (2002). *Isabelle/HOL: A Proof Assistant for Higher-Order Logic*. Vol. 2283. Lecture Notes in Computer Science. Springer.
- Osherson, D. N. and S. Weinstein (2012). “Preference based on Reasons”. In: *Rev. Symb. Log.* 5.1, pp. 122–147.
- Prakken, H. and G. Sartor (2015). “Law and logic: A review from an argumentation perspective”. In: *Artificial Intelligence* 227, pp. 214–225.
- Sartor, G. (2010). “Doing justice to rights and values: teleological reasoning and proportionality”. In: *Artificial Intelligence and 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.
- Sieckmann, J.-R., ed. (2010). *Legal Reasoning: The Methods of Balancing*. Vol. 124. ARSP Beiheft. Stuttgart: Franz Steiner.
- Teubner, G. (1983). “Substantive and Reflexive Elements in Modern Law”. In: *Law & Society Rev.* 17, pp. 239–285.

## A Appendix

### A.1 Ceteris paribus preference relation for $\mathcal{PL}$

Extending the SSE of  $\mathcal{PL}$  in HOL from Fig. 3 some further preference relations for  $\mathcal{PL}$  are defined in Fig. 9. These additional relations support *ceteris paribus* reasoning in  $\mathcal{PL}$ .

```
1 theory PreferenceLogicCeterisParibus (** Benzmüller & Fuenmayor, 2020 **)
2   imports PreferenceLogicBasics
3 begin (** Ceteris Paribus reasoning by van Benthem et al., JPL 2009 **)
4
5 (*Section 5: Equality-based Ceteris Paribus Preference Logic*)
6 abbreviation a1::"σ⇒π⇒bool" ("∈") where "φ ∈ Γ ≡ Γ φ"
7 abbreviation a2 ("⊆") where "Γ ⊆ Γ' ≡ ∀φ. φ ∈ Γ → φ ∈ Γ'"
8 abbreviation a3 ("⊇") where "Γ ⊇ Γ' ≡ λφ. φ ∈ Γ' ∨ φ ∈ Γ"
9 abbreviation a4 ("⊂") where "Γ ⊂ Γ' ≡ λφ. φ ∈ Γ' ∧ φ ∉ Γ"
10 abbreviation a5 ("⊂=") where "{φ} ≡ λx::σ. x=φ"
11 abbreviation a6 ("⊆=") where "{α,β} ≡ λx::σ. x=α ∨ x=β"
12 abbreviation a7 ("⊆=") where "{α,β,γ} ≡ λx::σ. x=α ∨ x=β ∨ x=γ"
13 abbreviation a8 ("∅") where "∅ ≡ (λψ::σ. False)"
14 abbreviation a9 ("U") where "U ≡ (λψ::σ. True)"
15
16 abbreviation c14 ("≡") where "w ≡ r v ≡ ∀φ. φ ∈ Γ → (φ w ↔ φ v)"
17 abbreviation c15 ("≡") where "w ≡ r v ≡ w ≡ v ∧ w ≡ r v"
18 abbreviation c16 ("≡") where "w ≡ r v ≡ w < v ∧ w ≡ r v"
19 abbreviation c17 ("⊆") where "(Γ)⊆ ≡ λw. ∃v. w ⊆ r v ∧ φ v"
20 abbreviation c18 ("⊆") where "(Γ)⊆ ≡ λw. ∀v. w ⊆ r v → φ v"
21 abbreviation c19 ("⊆") where "(Γ)⊆ ≡ λw. ∃v. w ⊆ r v ∧ φ v"
22 abbreviation c20 ("⊆") where "(Γ)⊆ ≡ λw. ∀v. w ⊆ r v → φ v"
23 abbreviation c21 ("⊆") where "(Γ)⊆ ≡ λw. ∃v. w ≡ r v ∧ φ v"
24 abbreviation c22 ("⊆") where "(Γ)⊆ ≡ λw. ∀v. w ≡ r v → φ v"
25
26 (*Section 6: Ceteris Paribus Counterparts of Binary Pref. Statements*)
27 (*operators below not defined in paper; existence is tacitly suggested.
28   AA-variant draws upon von Wright's. AE-variant draws upon Halpern's.*)
29 abbreviation c23 ("⊆AA") where "⊆AA ≡ Vs.Vt. φ s ∧ ψ t → s ⊆r t"
30 abbreviation c24 ("⊆AA") where "⊆AA ≡ Vs.Vt. φ s ∧ ψ t → s ⊆r t"
31 where "(φ ⊆AA ψ) u ≡ Vs.Vt. φ s ∧ ψ t → s ⊆r t"
32 abbreviation c25 ("⊆AE") where "⊆AE ≡ Vs. ∃t. φ s → ψ t ∧ s ⊆r t"
33 where "(φ ⊆AE ψ) u ≡ Vs. ∃t. φ s → ψ t ∧ s ⊆r t"
34 abbreviation c26 ("⊆AE") where "⊆AE ≡ Vs. ∃t. φ s → ψ t ∧ s ⊆r t"
35 where "(φ ⊆AE ψ) u ≡ Vs. ∃t. φ s → ψ t ∧ s ⊆r t"
36 abbreviation c27 ("⊆AA") where "⊆AA ≡ A(ψ → (Γ)⊆)"
37 abbreviation c28 ("⊆AA") where "⊆AA ≡ A(ψ → (Γ)⊆)"
38 abbreviation c29 ("⊆AE") where "⊆AE ≡ A(φ → (Γ)⊆)"
39 abbreviation c30 ("⊆AE") where "⊆AE ≡ A(φ → (Γ)⊆)"
40
41 (*Consistency confirmed (trivial: only abbreviations are introduced*)
42 lemma True nitpick[satisfy,user_axioms] oops
43 end
```

Figure 9: SSE (cont'd) of  $\mathcal{PL}$  (Benthem et al. 2009) in HOL

We give some explanations:

**Lines 6–14** Useful set theoretic notions are introduced as abbreviations for corresponding  $\lambda$ -terms in HOL.

**Lines 16–24**  $\mathcal{PL}$  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 modeling von Wright’s notion of *ceteris paribus* (“all other things being equal”) preferences, eliciting an appropriate  $\Gamma$  by extra-logical means.

**Lines 29–40:** Except for  $\prec_{AA}^\Gamma$ , the remaining operators we define here are not explicitly defined in (Benthem et al. 2009); however, their existence is tacitly suggested.

### A.2 Testing the Meta-Theory of $\mathcal{PL}$

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. 10 and 11.

```
1 theory PreferenceLogicTests1 (** Benzmüller & Fuenmayor, 2020 **)
2   imports PreferenceLogicBasics
3 begin (*Tests for the SSE of van Benthem et al., JPL 2009, in HOL*)
4 (*Fact 1: definability of the principal operators and verification*)
5 lemma F1_9: "(φ ⊆ EE ψ) u ↔ (φ ⊆ EE ψ) u" by smt
6 lemma F1_10: "(φ ⊆ AE ψ) u ↔ (φ ⊆ AE ψ) u" by smt
7 lemma F1_11: "(φ ⊆ EE ψ) u ↔ (φ ⊆ EE ψ) u" by smt
8 lemma F1_12: "(φ ⊆ AE ψ) u ↔ (φ ⊆ AE ψ) u" by smt
9 (*Fact 2: definability of remaining pref. operators and verification*)
10 lemma F2_13: "is_total SBR → ((φ < AA ψ) u ↔ (φ < AA ψ) u)" by smt
11 lemma F2_14: "is_total SBR → ((φ > EA ψ) u ↔ (φ > EA ψ) u)" by smt
12 lemma F2_15: "is_total SBR → ((φ < AA ψ) u ↔ (φ < AA ψ) u)" by smt
13 lemma F2_16: "is_total SBR → ((φ > EA ψ) u ↔ (φ > EA ψ) u)" by smt
14 (*Section 3.5 "Axiomatization" -- verify interaction axioms*)
15 lemma Incl_1: "[ (φ < φ) → (φ < φ) ]" by auto
16 lemma Inter_1: "[ (φ < φ) → (φ < φ) ]" using tBR by blast
17 lemma Trans_le: "[ (φ < φ) → (φ < φ) ]" using tBR by blast
18 lemma Inter_2: "[ (φ < φ) → ((φ < φ) ∨ φ < (ψ < φ)) ]" by blast
19 lemma F4: "[ (φ < φ) → ((φ < φ) ∨ φ < (ψ < φ)) ]" by smt
20 (Vw. Vv. ((w < v) ∧ ¬(v < w)) → (w < v))) by smt
21 lemma Inter_3: "[ (φ < φ) → (φ < φ) ]" using tBR by blast
22 lemma Incl_2: "[ (φ < φ) → (φ < φ) ]" by blast
23 (*Section 3.6 "A binary preference fragment"*)
24 (* ⊆ is the dual of < AA *)
25 lemma "[ (φ ⊆ EE ψ) ↔ ¬(ψ < AA φ) ] ∧ [ (ψ < AA φ) ↔ ¬(ψ ⊆ EE φ) ]" by simp
26 (* ⊆ is the dual of < AA only if totality is assumed*)
27 lemma "[ (φ ⊆ EE ψ) ↔ ¬(ψ < AA φ) ]" nitpick oops (*countermodel*)
28 lemma "[ (φ ⊆ EE ψ) → ¬(ψ < AA φ) ]" by blast (*this direction holds*)
29 lemma "is_total SBR → [ (φ ⊆ EE ψ) ↔ ¬(ψ < AA φ) ]" by blast
30 lemma "[ (φ < AA ψ) ↔ ¬(ψ ⊆ EE φ) ]" nitpick oops (*countermodel*)
31 lemma "[ (φ < AA ψ) → ¬(ψ ⊆ EE φ) ]" by blast (*this direction holds*)
32 lemma "is_total SBR → [ (φ < AA ψ) ↔ ¬(ψ ⊆ EE φ) ]" by blast
33 (* verify p.97-98 *)
34 lemma monotonicity: "[ ((φ ⊆ EE ψ) ∧ A(φ → ζ)) → (ζ ⊆ EE ψ) ]" by blast
35 lemma reducibility: "[ (((φ ⊆ EE ψ) ∧ α) ⊆ EE β) ↔ ((φ ⊆ EE ψ) ∧ (α ⊆ EE β)) ]" by blast
36 lemma reflexivity: "[ φ → (φ ⊆ EE φ) ]" using rBR by blast
37 (*The condition below is supposed to enforce totality of the preference
38   relation. However there are countermodels. See p.98?*)
39 lemma "is_total SBR →
40   [ ((φ ⊆ EE ψ) ∧ (ψ ⊆ EE φ)) → ((φ ⊆ EE ψ) ∨ (ψ ⊆ EE φ)) ]" by auto
41 lemma "[ ((φ ⊆ EE ψ) ∧ (ψ ⊆ EE φ)) → ((φ ⊆ EE ψ) ∨ (ψ ⊆ EE φ)) ]"
42   → is_total SBR" nitpick oops (*countermodel - error in paper?*)
43 lemma "is_total SBR →
44   [ ((φ ⊆ EE ψ) ∧ (ψ ⊆ EE φ)) → ((φ ⊆ EE ψ) ∨ (ψ ⊆ EE φ)) ]" by auto
45 lemma "[ ((φ ⊆ EE ψ) ∧ (ψ ⊆ EE φ)) → ((φ ⊆ EE ψ) ∨ (ψ ⊆ EE φ)) ]"
46   → is_total SBR" nitpick oops (*countermodel - error in paper?*)
47 end
```

Figure 10: Experiments: Testing the meta-theory of  $\mathcal{PL}$

We briefly explain the experiments shown in Fig. 10:

**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 the corresponding ones claimed (and in many cases proved) in Benthem et al. (2009), except for the claims encoded in lines 41 and 42, where countermodels are reported by *Nitpick*.

**Lines 25–47** Some application-specific tests in preparation for the modeling of the value ontology are conducted.

Further tests are shown in Fig. 11 for the *ceteris paribus* extension of  $\mathcal{PL}$ ; we automatically prove all the relevant results from Benthem et al. (2009).

```

1 theory PreferenceLogicTests2 (** Benzmüller & Fuenmayor, 2020 **)
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 "[([F]⊆φ) ↔ ¬([F]⊆¬φ)]" by auto
7 lemma "[([F]⊆¬φ) ↔ ¬([F]⊆φ)]" by auto
8 lemma "[([F]⊆φ) ↔ ¬([F]⊆¬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) ∧ ((⊆¬φ) w ↔ ((⊆⊆¬φ) w))" by auto
13 (**Axiomatization**)
14 (*inclusion and interaction axioms *)
15 lemma Incl1: "[([F]⊆φ) → (([F]⊆φ)]" by auto
16 lemma Incl2: "[([F]⊆φ) → (([F]⊆φ)]" by auto
17 lemma Incl3: "[([F]⊆([F]⊆φ)) → ([F]⊆φ)]" by (meson tBR)
18 lemma Incl4: "[([F]⊆([F]⊆φ)) → ([F]⊆φ)]" by (metis tBR)
19 lemma Incl5: "[([F]⊆([F]⊆φ)) → (([F]⊆φ) ∨ ([F]⊆¬φ))]" by (metis rBR)
20 (*ceteris paribus reflexivity*)
21 lemma CetPar6: "φ ∈ Γ → ([F]⊆φ) → φ]" by blast
22 lemma CetPar7: "φ ∈ Γ → ([F]⊆¬φ) → ¬φ]" by blast
23 (*monotonicity*)
24 lemma CetPar8: "Γ ⊆ Γ' → ([F]⊆φ) → ([F]⊆φ)" by auto
25 lemma CetPar9: "Γ ⊆ Γ' → ([F]⊆¬φ) → ([F]⊆¬φ)" by auto
26 lemma CetPar10: "Γ ⊆ Γ' → ([F]⊆φ) → ([F]⊆φ)" by auto
27 (*increase (decrease) of ceteris paribus sets*)
28 lemma CetPar11a: "[([F]⊆([F]⊆φ)) → ([F]⊆φ)]" by auto
29 lemma CetPar11b: "[([F]⊆([F]⊆φ)) → ([F]⊆φ)]" by auto
30 lemma CetPar12a: "[([F]⊆([F]⊆φ)) → ([F]⊆φ)]" by auto
31 lemma CetPar12b: "[([F]⊆([F]⊆φ)) → ([F]⊆φ)]" by auto
32 lemma CetPar13a: "[([F]⊆([F]⊆φ)) → ([F]⊆φ)]" by auto
33 lemma CetPar13b: "[([F]⊆([F]⊆φ)) → ([F]⊆φ)]" by auto
34 lemma CetPar13c: "[([F]⊆([F]⊆φ)) → ([F]⊆φ)]" by auto
35 (*Example 1, Lemma 4, Corollary 1 and Lemma5*)
36 lemma Ex1: "[([F]⊆([F]⊆φ)) → ([F]⊆φ)]" using rBR by auto
37 lemma Lem4: "((⊆⊆φ) w → (∃v. (w ⊆ v) ∧ (φ v)))" by simp
38 lemma Cor1: "((⊆⊆φ) w → (∃v. (w ⊆ v) ∧ (φ v)))" by simp
39 lemma Lem5: "(w ⊆ v) → ((w ⊆ v) ∧ (φ v))" by auto
40 (** Section 6: Ceteris Paribus Counterparts **)
41 (*AA-variant (drawing upon von Wright's*)
42 lemma "(φ <AA φ) u ↔ (φ <AA φ) u" nitpick oops (*Ctm*)
43 lemma "(φ <AA φ) u ↔ (φ <AA φ) u" nitpick oops (*Ctm*)
44 lemma "(φ <AA φ) u ↔ (φ <AA φ) u" by auto
45 lemma "is_total SBR → (φ <AA φ) u ↔ (φ <AA φ) u" by smt
46 lemma "(φ <AA φ) u ↔ (φ <AA φ) u" nitpick oops (*Ctm*)
47 lemma "(φ <AA φ) u ↔ (φ <AA φ) u" nitpick oops (*Ctm*)
48 lemma "(φ <AA φ) u ↔ (φ <AA φ) u" by auto
49 lemma "is_total SBR → (φ <AA φ) u ↔ (φ <AA φ) u" by smt
50 (*AE-variant*)
51 lemma leAE_cp_pref: "(φ <AE φ) u ↔ (φ <AE φ) u" by auto
52 lemma leqAE_cp_pref: "(φ <AE φ) u ↔ (φ <AE φ) u" by auto
53 (*miscellaneous tests*)
54 lemma "let I=∅ in [(φ <AA φ) ↔ (φ <AA φ)]" by simp
55 lemma "let I={1} in [(φ <AA φ) ↔ (φ <AA φ)]" by simp
56 lemma "let I={1,A} in [(φ <AA φ) ↔ (φ <AA φ)]" nitpick oops (*Ctm*)
57 lemma "let I={A} in [(φ <AA φ) ↔ (φ <AA φ)]" nitpick oops (*Ctm*)
58 lemma "let I={A} in [(φ <AA φ) ↔ (φ <AA φ)]" nitpick oops (*Ctm*)
59 lemma "let I=∅ in [(φ <AE φ) ↔ (φ <AE φ)]" by simp
60 lemma "let I={1} in [(φ <AE φ) ↔ (φ <AE φ)]" by simp
61 lemma "let I={1,A} in [(φ <AE φ) ↔ (φ <AE φ)]" nitpick oops (*Ctm*)
62 lemma "let I={A} in [(φ <AE φ) ↔ (φ <AE φ)]" by auto
63 lemma "let I={A,B} in [(φ <AE φ) ↔ ((A ∧ B) → (φ <AE φ))]" by auto
64 lemma "let I={A} in [(A → (φ <AE φ)) → (φ <AE φ)]" nitpick oops (*Ctm*)
65 end

```

Figure 11: Experiments (cont'd): Testing the meta-theory of  $\mathcal{PL}$

### A.3 Testing the Value Ontology

Some further tests on the modeling and encoding of the value ontology are conducted; these tests are displayed in Fig. 12. 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 7–19), and we carry out some further tests on the value ontology (extending the ones presented in §5) concerning value aggregation and consistency (Lines 24ff.).

```

1 theory ValueOntologyTestLong (** Benzmüller, Fuenmayor & Lomfeld, 2020 **)
2 imports ValueOntology
3 begin
4 lemma "True" nitpick[satisfy,show_all,card i=10] oops
5 lemma "[INCONS]" nitpick[satisfy,card i=4] nitpick oops (*contingent*)
6 (*ext/int operators satisfy main properties of Galois connections*)
7 lemma G1: "B ⊆ A1 ↔ A ⊆ B1" by blast
8 lemma G1: "A ⊆ A1" by simp
9 lemma G2: "B ⊆ B1" by simp
10 lemma G3: "A1 ⊆ A2 → A2 ⊆ A1" by simp
11 lemma G4: "B1 ⊆ B2 → B2 ⊆ B1" by simp
12 lemma cl1: "A1 = A1" by blast
13 lemma cl2: "B1 = B1" by blast
14 lemma dual1a: "(A1 ⊆ A2) ⊆ (A1 ⊆ A2)" by blast
15 lemma dual1b: "(B1 ⊆ B2) ⊆ (B1 ⊆ B2)" by blast
16 lemma "([A1 ⊆ A2] ⊆ [A1 ⊆ A2])" nitpick oops
17 lemma "([B1 ⊆ B2] ⊆ [B1 ⊆ B2])" nitpick oops
18 lemma dual2a: "(A1 ⊆ A2) ⊆ (A1 ⊆ A2)" by blast
19 lemma dual2b: "(B1 ⊆ B2) ⊆ (B1 ⊆ B2)" by blast
20 (*Note: two different but logically equivalent notations*)
21 lemma "[WILL*] ≡ WILL*" by simp
22 lemma "[WILL*⊗STAB*] ≡ (WILL*⊗STAB*)" by simp
23 (***** value ontology tests *****
24 lemma "[RELIP] ∧ [WILL*] → [INCONS]" by simp
25 lemma "[INCONS] → [RELIP] ∧ [WILL*]" by simp
26 lemma "[RELIP] ∧ [WILL*]" nitpick[satisfy] nitpick oops (*contingent*)
27 lemma "[FAIR*] ∧ [EFFI*]" nitpick[satisfy] nitpick oops (*contingent*)
28 lemma "[(-INCONS) ∧ [FAIR*] ∧ [EFFI*]]"
29 nitpick[satisfy,show_all] nitpick oops (*contingent: p & d independent*)
30 lemma "[(-INCONS) ∧ (-INCONS) ∧ [RELIP] ∧ [WILL*]]"
31 nitpick[satisfy,show_all] nitpick oops (*contingent: p & d independent*)
32 (** more tests **)
33 (*values in two non-opposed quadrants (nog): consistent*)
34 lemma "[WILL*] ∧ [STAB*] → [INCONS]" nitpick oops (*countermodel found*)
35 lemma "[WILL*] ∧ [GAIN*] ∧ [EFFI*] ∧ [STAB*] → [INCONS]" nitpick oops
36 (*values in two opposed quadrants: inconsistent*)
37 lemma "[RESP*] ∧ [STAB*] → [INCONS]" by simp
38 (*values in three quadrants: inconsistent*)
39 lemma "[WILL*] ∧ [EFFI*] ∧ [RELI*] → [INCONS]" by simp
40 (*values in opposed quadrants for different parties: consistent*)
41 lemma "[[EQUI*] ∧ [GAIN*] → (INCONS* ∨ INCONS*)]" nitpick oops (*cntmdl*)
42 lemma "[[RESP*] ∧ [STAB*] → (INCONS* ∨ INCONS*)]" nitpick oops (*cntmdl*)
43 (*value preferences tests*)
44 lemma "[WILL* <v WILL*⊗STAB*]"
45 nitpick nitpick[satisfy] oops (*contingent*)
46 lemma "[WILL* <v STAB*] → [WILL* <v WILL*⊗STAB*]" by blast
47 lemma "[WILL* <v STAB*] → [WILL* <v RELI*⊗STAB*]" by blast
48 lemma "[WILL* <v WILL*⊗STAB*] → [WILL* <v STAB*]"
49 nitpick nitpick[satisfy] oops (*contingent*)
50 lemma "[WILL* <v RELI*⊗STAB*] → [WILL* <v STAB*]"
51 nitpick nitpick[satisfy] oops (*contingent*)
52 lemma "[WILL*⊗STAB* <v WILL*]" using rBR by auto
53 lemma "[WILL*⊗STAB* <v WILL*] → [STAB* <v WILL*]" by auto
54 lemma "[RELI*⊗STAB* <v WILL*] → [STAB* <v WILL*]" by auto
55 lemma "[STAB* <v WILL*] → [WILL*⊗STAB* <v WILL*]"
56 nitpick nitpick[satisfy] oops (*contingent*)
57 lemma "[STAB* <v WILL*] → [RELI*⊗STAB* <v WILL*]"
58 nitpick nitpick[satisfy] oops (*contingent*)
59 (*basic properties*)
60 lemma "[¬(X <v X)]" using rBR by auto (*irreflexive*)
61 lemma "[X <v Y → ¬(Y <v X)]" nitpick oops (*not asymmetric*)
62 lemma "[X <v Y ∧ Y <v Z] → X <v Z]" nitpick oops (*not transitive*)
63 end

```

Figure 12: Further testing of the value ontology

### A.4 Modeling Conti v. ASPCA

Another illustrative case study we can model in our framework is: Conti v. ASPCA (cf. 353 NYS 2d 288).

*Chester, a parrot owned by the ASPCA, escaped and was recaptured by Conti. The ASPCA found this out and re-claimed 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 *Pier-son's case* (cf. §7) applies (with a preference in this case of FREEDOM over SECURITY, cf. §1). More specifically, it was ruled that for a domestic animal it is 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. A (simplified) argument is reconstructed as follows:

Let  $\alpha$  stand for the animal.  $p$  and  $d$  stand for ASPCA and Conti respectively. ASPCA ( $p$ ) argument goes like this:

- 1 ASPCA owned the parrot – (Prop  $p \alpha$ ).
- 2 The parrot escaped and was recaptured by Conti – (Capture  $d \alpha$ ).
- 3 However, ASPCA still cares for (recovering and maintaining) the animal – (Care  $p \alpha$ ).
- 4 Therefore, ruling for ASPCA – (For  $p$ ) is preferred (since the combination of RELI & RESP is preferred to STAB by legal rule R3, cf. §7).

The reconstructed theory for the Conti vs. ASPCA case is displayed in Fig. 13.

```

1 theory Conti (** Benzmüller, Fuenmayor & Lomfeld, 2020 **)
2   imports GeneralKnowledge
3   begin (** ASPCA v. Conti "wild animal" case **)
4
5   (*case-specific 'world-vocabulary'*)
6   consts  $\alpha$ ::"e" (*appropriated animal (parrot in this case) *)
7   consts Care::" $c \Rightarrow e \Rightarrow \sigma$ "
8   consts Prop::" $c \Rightarrow e \Rightarrow \sigma$ "
9   consts Capture::" $c \Rightarrow e \Rightarrow \sigma$ "
10
11  (*case-specific taxonomic (legal domain) knowledge*)
12  axiomatization where
13    CW1: "[Animal  $\alpha \wedge$  Pet  $\alpha \rightarrow$  Domestic  $\alpha$ ]" and
14    CW2: "[ $(\exists c. \text{Capture } c \alpha \wedge \text{Domestic } \alpha) \rightarrow \text{appDomAnimal}$ ]" and
15    CW3: "[ $\forall c. \text{Care } c \alpha \rightarrow \text{Mtn } c$ ]" and
16    CW4: "[ $\forall c. \text{Prop } c \alpha \rightarrow \text{Own } c$ ]" and
17    CW5: "[ $\forall c. \text{Capture } c \alpha \rightarrow \text{Poss } c$ ]"
18
19  lemma True nitpick[satisfy,card i=4] oops (*satisfiable*)
20
21  (***** pro-ASPCA's argument *****)
22  abbreviation "ASPCA_facts  $\equiv$  [Parrot  $\alpha \wedge$  Pet  $\alpha \wedge$  Care  $p \alpha \wedge$ 
23    Prop  $p \alpha \wedge (\neg \text{Prop } d \alpha) \wedge \text{Capture } d \alpha]$ "
24
25  (* decision for defendant (Conti) is compatible with premises*)
26  lemma "ASPCA_facts  $\wedge [\neg \text{INCONS}^p] \wedge [\neg \text{INCONS}^d] \wedge [\text{For } p \prec \text{For } d]$ "
27    nitpick[satisfy,card i=4] oops (* (non-trivial) model found*)
28
29  (* decision for plaintiff (ASPCA) is compatible with premises*)
30  lemma "ASPCA_facts  $\wedge [\neg \text{INCONS}^p] \wedge [\neg \text{INCONS}^d] \wedge [\text{For } d \prec \text{For } p]$ "
31    nitpick[satisfy,card i=4] oops (* (non-trivial) model found*)
32
33  (* decision for plaintiff (ASPCA) is provable*)
34  lemma aux: assumes ASPCA_facts shows "[ $(\text{STAB}^d \prec_v \text{RELI}^p \oplus \text{RESP})$ ]"
35    using CW1 CW2 W7 asms R3 by fastforce
36  theorem assumes ASPCA_facts shows "[ $\text{For } d \prec \text{For } p$ ]"
37    using asms aux CW5 ForAx F3 other.simps(1) rBR by metis
38
39  (* while a decision for the defendant is refutable*)
40  lemma assumes ASPCA_facts shows "[ $\text{For } p \prec \text{For } d$ ]"
41  nitpick[card i=4] oops (* (non-trivial) counterexample found*)
42 end

```

Figure 13: Modeling of the Conti vs. ASPCA case

## A.5 Complex (Counter-)Models

An example of a complex countermodel with four possible worlds/states (type  $i$ ) that has been computed by *Nitpick* for the statement in Line 41 in Fig. 13 is presented in Fig. 14.

This countermodel is included here to illustrate (to the interested reader) the richness of the information and the level of detail that is supported in our framework. It is in particular the combination of automated theorem proving, model finding and countermodel finding that is supporting the knowledge engineer and user of the LogiKey framework in a unique manner to gain intuition about the modeled structures. And, in fact, these three analysis modes can be worked with in parallel in *Isabelle/HOL*. This is one reason for the good response rates to user requests that we often experience – despite the general undecidability of HOL.



Nitpick found a counterexample for card e = 1 and card i = 4:

```

Skolem constant:
  λv. v = (λx. _)(i1 := i4, i2 := i1, i3 := i1, i4 := i4)
Types:
  c = {d, p}
  e × i [boxed] = {(e1, i1), (e1, i2), (e1, i3), (e1, i4)}
  c VAL = {FREEDOM d, FREEDOM p, UTILITY d, UTILITY p, EQUALITY d, EQUALITY p, SECURITY d, SECURITY p}
Constants:
  Capture =
    (λx. _)
    ((d, e1, i1) := True, (d, e1, i2) := True, (d, e1, i3) := True, (d, e1, i4) := True, (p, e1, i1) := False, (p, e1, i2) := False,
    (p, e1, i3) := False, (p, e1, i4) := False)
  Care =
    (λx. _)
    ((d, e1, i1) := False, (d, e1, i2) := False, (d, e1, i3) := False, (d, e1, i4) := True, (p, e1, i1) := True, (p, e1, i2) := True,
    (p, e1, i3) := True, (p, e1, i4) := True)
  Prop =
    (λx. _)
    ((d, e1, i1) := False, (d, e1, i2) := False, (d, e1, i3) := False, (d, e1, i4) := False, (p, e1, i1) := True, (p, e1, i2) := True,
    (p, e1, i3) := True, (p, e1, i4) := True)
  α = e1
  Animal = (λx. _)((e1, i1) := True, (e1, i2) := True, (e1, i3) := True, (e1, i4) := True)
  Domestic = (λx. _)((e1, i1) := True, (e1, i2) := True, (e1, i3) := True, (e1, i4) := True)
  Fox = (λx. _)((e1, i1) := False, (e1, i2) := False, (e1, i3) := False, (e1, i4) := False)
  FreeRoaming = (λx. _)((e1, i1) := False, (e1, i2) := False, (e1, i3) := False, (e1, i4) := False)
  Intent =
    (λx. _)
    ((d, i1) := False, (d, i2) := True, (d, i3) := False, (d, i4) := True, (p, i1) := False, (p, i2) := False, (p, i3) := True,
    (p, i4) := False)
  Liv =
    (λx. _)
    ((d, i1) := False, (d, i2) := True, (d, i3) := False, (d, i4) := False, (p, i1) := False, (p, i2) := False, (p, i3) := False,
    (p, i4) := False)
  Mtn =
    (λx. _)
    ((d, i1) := True, (d, i2) := False, (d, i3) := False, (d, i4) := True, (p, i1) := True, (p, i2) := True, (p, i3) := True,
    (p, i4) := True)
  Own =
    (λx. _)
    ((d, i1) := False, (d, i2) := False, (d, i3) := False, (d, i4) := False, (p, i1) := True, (p, i2) := True, (p, i3) := True,
    (p, i4) := True)
  Parrot = (λx. _)((e1, i1) := True, (e1, i2) := True, (e1, i3) := True, (e1, i4) := True)
  Pet = (λx. _)((e1, i1) := True, (e1, i2) := True, (e1, i3) := True, (e1, i4) := True)
  Poss =
    (λx. _)
    ((d, i1) := True, (d, i2) := True, (d, i3) := True, (d, i4) := True, (p, i1) := False, (p, i2) := False, (p, i3) := False,
    (p, i4) := False)
  appAnimal = (λx. _)(i1 := True, i2 := True, i3 := True, i4 := True)
  appDomAnimal = (λx. _)(i1 := True, i2 := True, i3 := True, i4 := True)
  appObject = (λx. _)(i1 := True, i2 := True, i3 := True, i4 := True)
  appWildAnimal = (λx. _)(i1 := False, i2 := False, i3 := False, i4 := False)
  BR = (λx. _)
    ((i1, i1) := True, (i1, i2) := False, (i1, i3) := False, (i1, i4) := False, (i2, i1) := True, (i2, i2) := True,
    (i2, i3) := False, (i2, i4) := False, (i3, i1) := True, (i3, i2) := True, (i3, i3) := True, (i3, i4) := False,
    (i4, i1) := False, (i4, i2) := False, (i4, i3) := False, (i4, i4) := True)
  For =
    (λx. _)
    ((d, i1) := False, (d, i2) := True, (d, i3) := True, (d, i4) := False, (p, i1) := True, (p, i2) := False, (p, i3) := False,
    (p, i4) := True)
  I = (λx. _)
    ((i1, FREEDOM d) := True, (i1, FREEDOM p) := True, (i1, UTILITY d) := True, (i1, UTILITY p) := False, (i1, EQUALITY d) := False,
    (i1, EQUALITY p) := True, (i1, SECURITY d) := True, (i1, SECURITY p) := True, (i2, FREEDOM d) := True, (i2, FREEDOM p) := True,
    (i2, UTILITY d) := True, (i2, UTILITY p) := False, (i2, EQUALITY d) := True, (i2, EQUALITY p) := False, (i2, SECURITY d) := True,
    (i2, SECURITY p) := True, (i3, FREEDOM d) := False, (i3, FREEDOM p) := True, (i3, UTILITY d) := True, (i3, UTILITY p) := False,
    (i3, EQUALITY d) := True, (i3, EQUALITY p) := True, (i3, SECURITY d) := True, (i3, SECURITY p) := False, (i4, FREEDOM d) := True,
    (i4, FREEDOM p) := True, (i4, UTILITY d) := False, (i4, UTILITY p) := False, (i4, EQUALITY d) := True, (i4, EQUALITY p) := True,
    (i4, SECURITY d) := True, (i4, SECURITY p) := True)
  other = (λx. _)(d := p, p := d)

```

Figure 14: Example of a countermodel to the statement in Line 41 in Fig. 13