

# A Deontic Logic Reasoning Infrastructure

**Christoph Benz Müller** (jww Xavier Parent & Leon van der Torre)

Freie Universität Berlin | University of Luxembourg

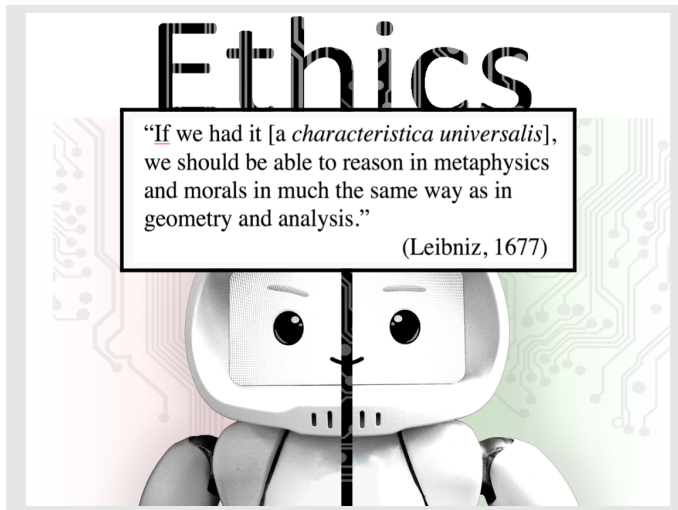


**HaPoC@CiE 2018, Kiel, 2 Aug 2018 — Celebration of Martin Davis' 90th Birthday**

# A Deontic Logic Reasoning Infrastructure

Christoph Benz Müller (jww Xavier Parent & Leon van der Torre)

Freie Universität Berlin | University of Luxembourg



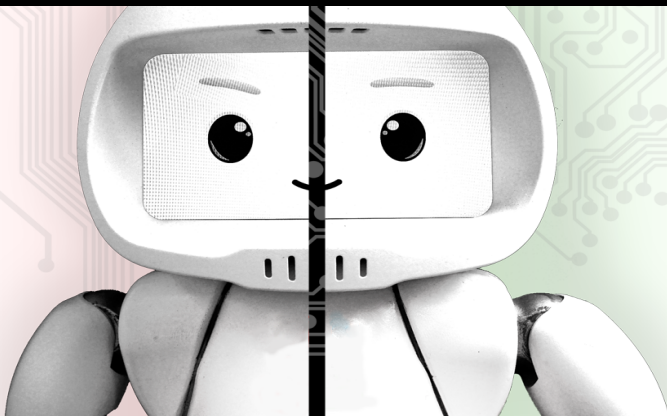
HaPoC@CiE 2018, Kiel, 2 Aug 2018 — Celebration of Martin Davis' 90th Birthday



# Ethics

Peaceful coexistence with **intelligent autonomous systems (IASs)**?

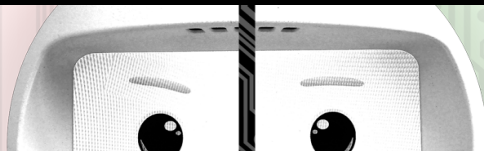
- ▶ appropriate forms of **machine-control**
- ▶ appropriate forms of **human-machine-interaction**



# Ethics

Peaceful coexistence with **intelligent autonomous systems (IASs)**?

- ▶ appropriate forms of **machine-control**
- ▶ appropriate forms of **human-machine-interaction**



Existing societal processes are based on:

- ▶ **rational argumentation & dialog**
- ▶ **explicit normative reasoning** (legal & ethical)

Deployment of IASs lacking such competencies? How wise is this?

# Ethics

## Talk Outline

Foreword: How does Martin Davis fit in?

- A Motivation:** Explicit Ethical Reasoning
- B Technology:** Universal Logical Reasoning in Higher-Order Logic
- C Evidence:** Experiments in Computational Metaphysics
- D Demo(s):** Normative Reasoning Experimentation Platform

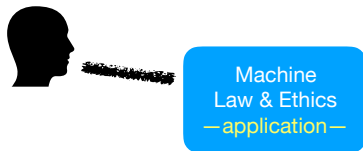
# Ethics

## Talk Outline

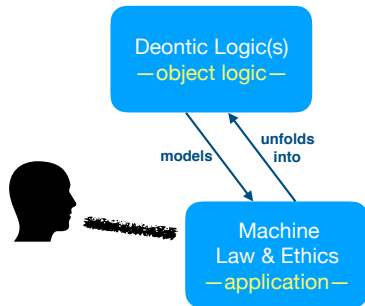
**Foreword:** How does Martin Davis fit in?

- A Motivation:** Explicit Ethical Reasoning
- B Technology:** Universal Logical Reasoning in Higher-Order Logic
- C Evidence:** Experiments in Computational Metaphysics
- D Demo(s):** Normative Reasoning Experimentation Platform

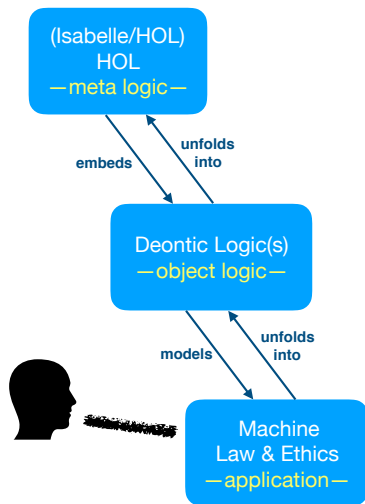
## How does Martin Davis fit in?



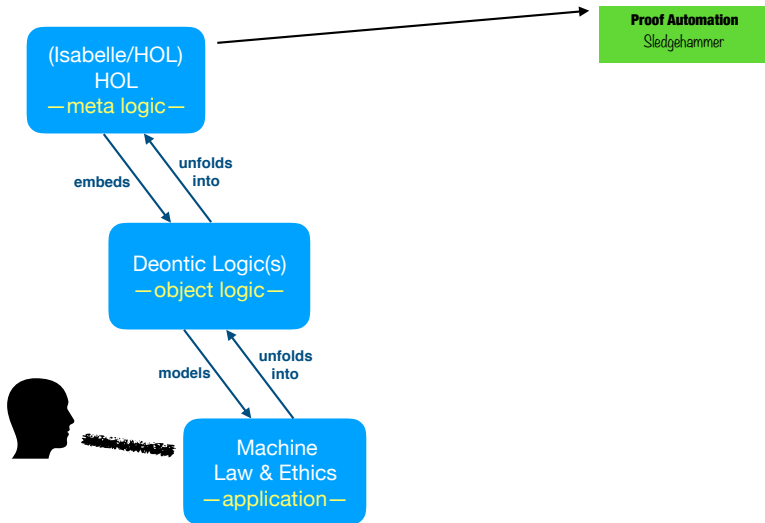
## How does Martin Davis fit in?



## How does Martin Davis fit in?

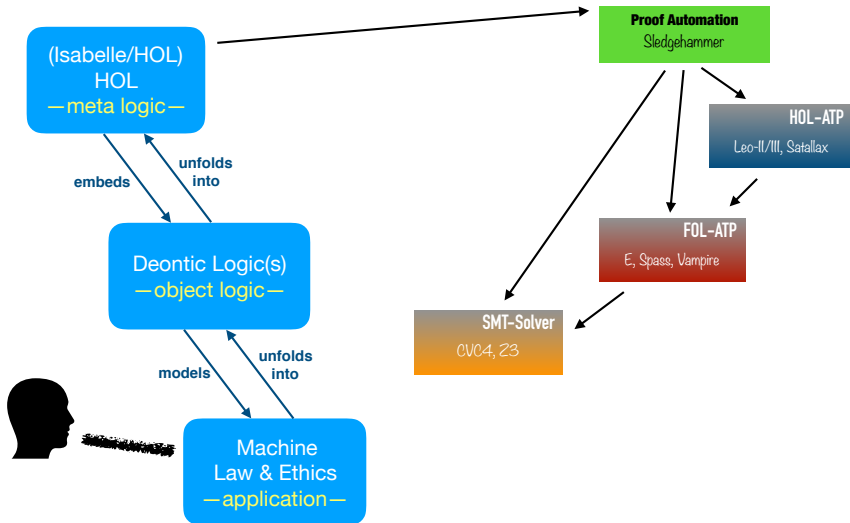


## How does Martin Davis fit in?

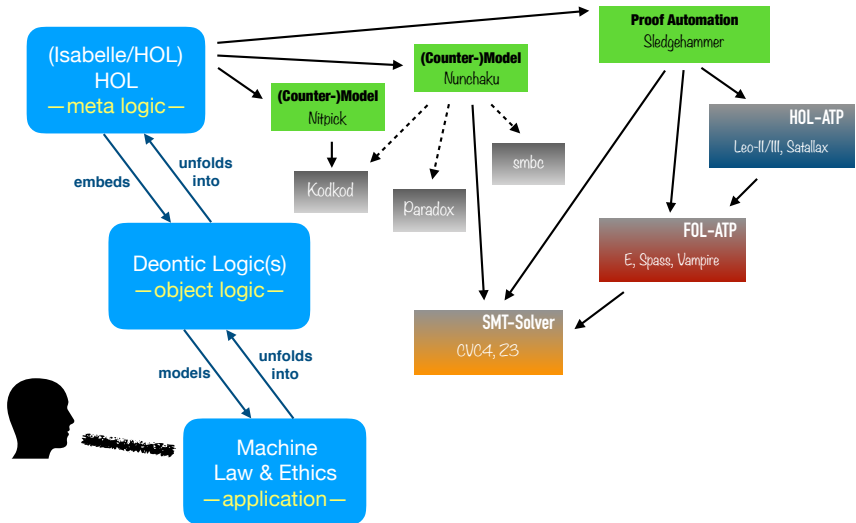




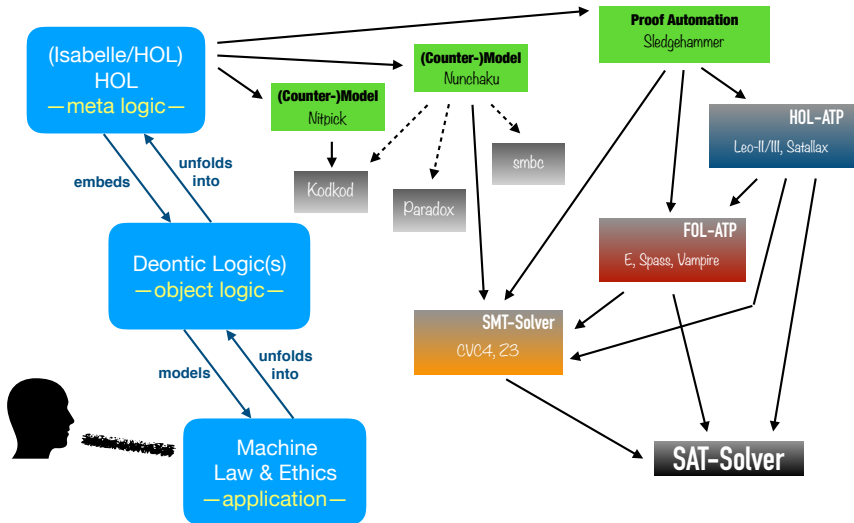
## How does Martin Davis fit in?



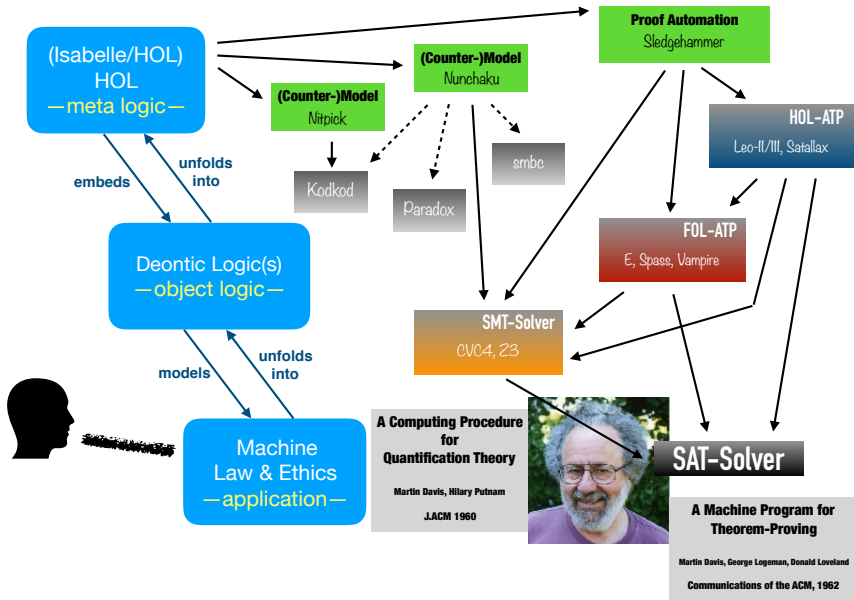
## How does Martin Davis fit in?



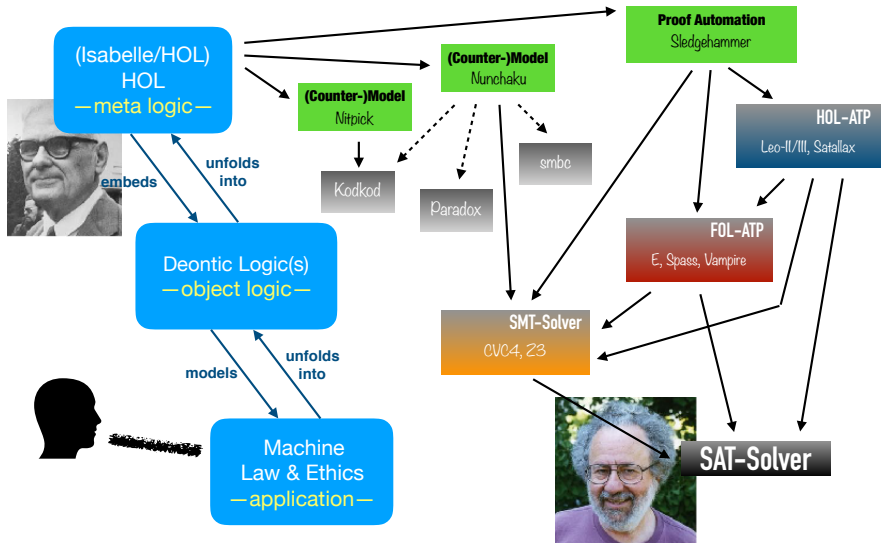
## How does Martin Davis fit in?



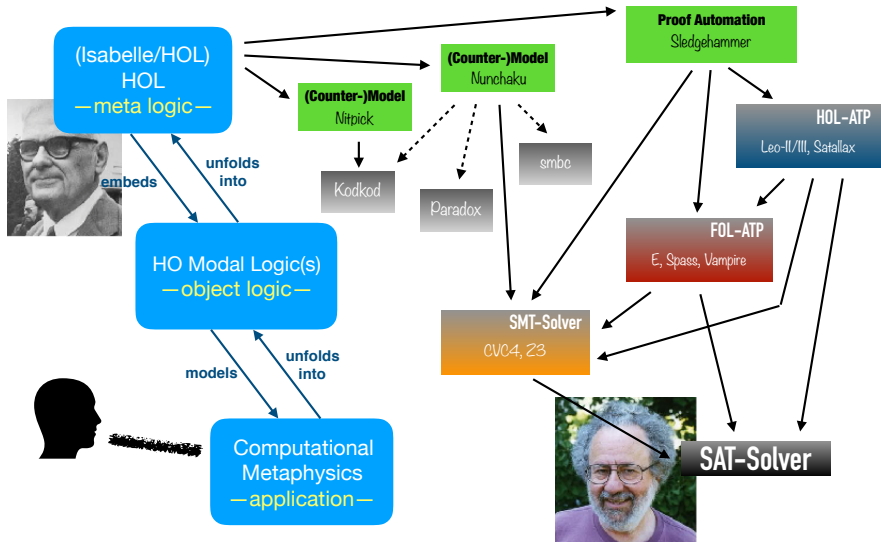
## How does Martin Davis fit in?



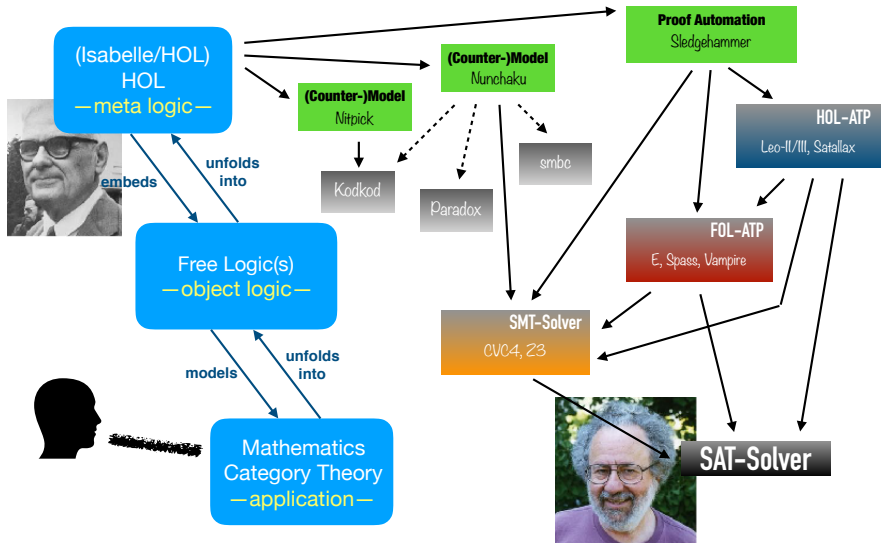
## How does Martin Davis fit in?



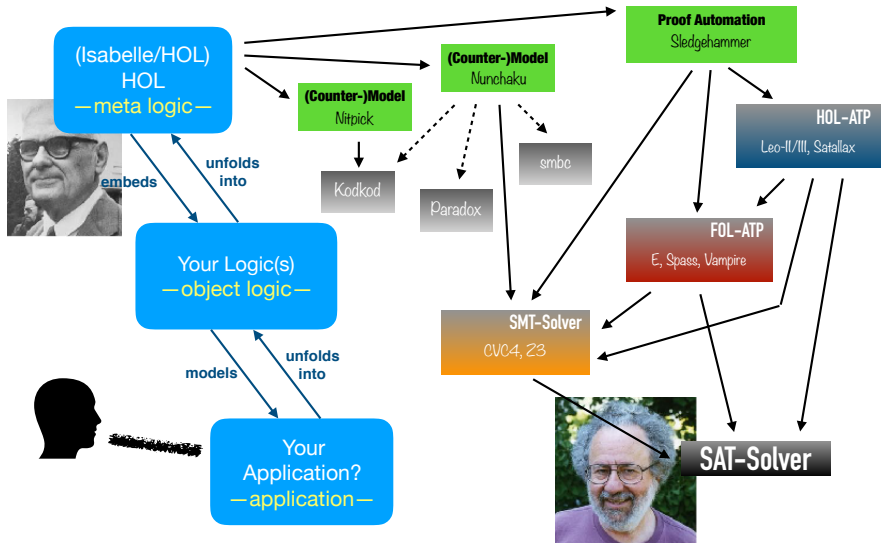
## How does Martin Davis fit in?



## How does Martin Davis fit in?

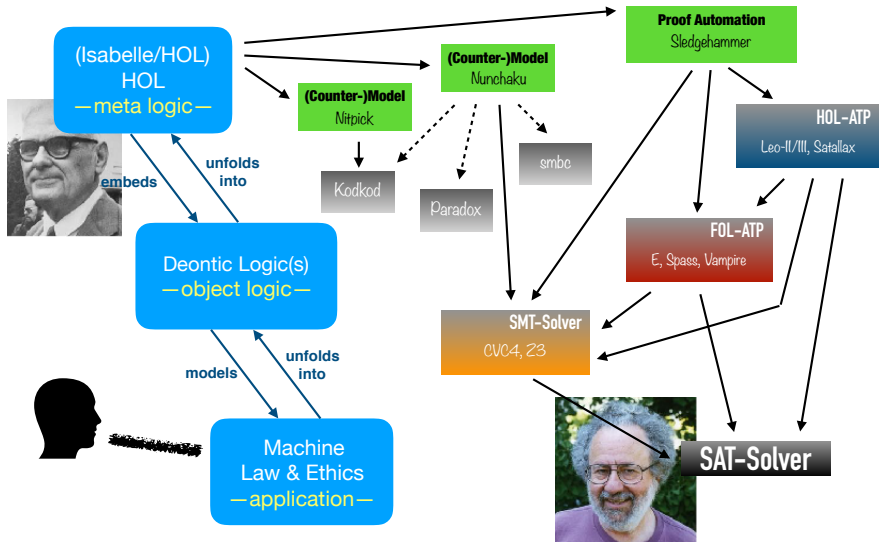


## How does Martin Davis fit in?





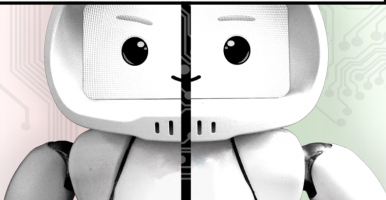
## How does Martin Davis fit in?



# Ethics

“If we had it [a *characteristica universalis*], we should be able to reason in metaphysics and morals in much the same way as in geometry and analysis.”

(Leibniz, 1677)



## Part A — Motivation: Explicit Ethical Reasoning

### Long-term: Emerging Superintelligence

Really? Anyhow ...

- ▶ How to prevent Superintelligence from turning against humanity?

### Medium-term: Development of pseudo-ethical skills in IASs

- ▶ Which norms? Which reasoning principles?
- ▶ What architectural design? What functionalities?
- ▶ How to implement, deploy and verify?

### Different kinds of systems and approaches:

- ▶ [Moor, 2009]:
  - ethical impact agents (ethical consequences to actions)
  - implicit ethical agents (ethical reactions to given situations)
  - explicit ethical agents (reasoning with ethical theories/rules)
  - full ethical agents (conscious, intentional, free will)
- ▶ bottom-up vs. top-down
- ▶ [DoranEtAl., 2017]:  
opaque — comprehensible — interpretable — explainable AI

### Long-term: Emerging Superintelligence

Really? Anyhow ...

- ▶ How to prevent Superintelligence from turning against humanity?

### Medium-term: Development of pseudo-ethical skills in IASs

- ▶ Which norms? Which reasoning principles?
- ▶ What architectural design? What functionalities?
- ▶ How to implement, deploy and verify?

### Different kinds of systems and approaches:

- ▶ [Moor, 2009]:
  - ethical impact agents (ethical consequences to actions)
  - implicit ethical agents (ethical reactions to given situations)
  - explicit ethical agents (reasoning with ethical theories/rules)
  - full ethical agents (conscious, intentional, free will)
- ▶ bottom-up vs. top-down
- ▶ [DoranEtAl., 2017]:  
opaque — comprehensible — interpretable — explainable AI

### Long-term: Emerging Superintelligence

Really? Anyhow ...

- ▶ How to prevent Superintelligence from turning against humanity?

### Medium-term: Development of pseudo-ethical skills in IASs

- ▶ Which norms? Which reasoning principles?
- ▶ What architectural design? What functionalities?
- ▶ How to implement, deploy and verify?

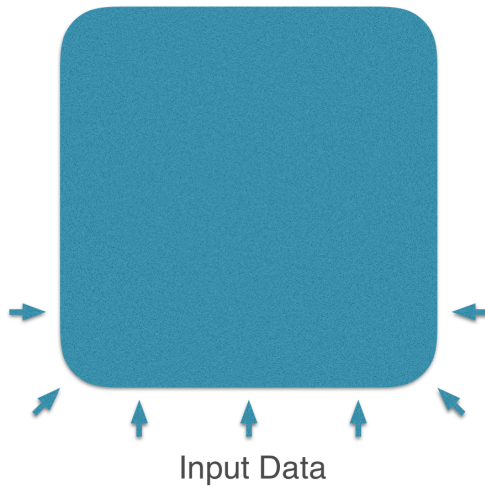
### Different kinds of systems and approaches:

- ▶ [Moor, 2009]:
  - ethical impact agents (ethical consequences to actions)
  - implicit ethical agents (ethical reactions to given situations)
  - **explicit ethical agents** (reasoning with ethical theories/rules)
  - full ethical agents (conscious, intentional, free will)
- ▶ bottom-up vs. **top-down**
- ▶ [DoranEtAl., 2017]:  
opaque — comprehensible — interpretable — **explainable** AI

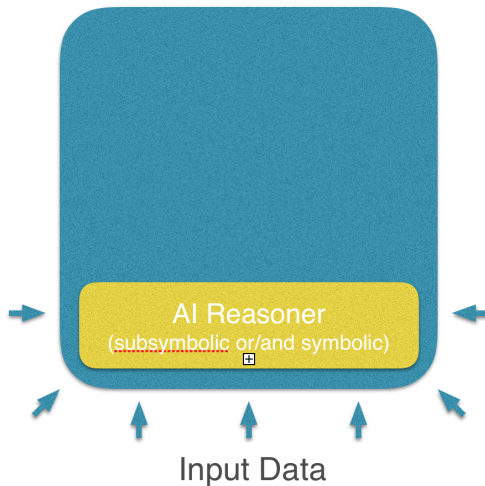


IAS

## Pseudo-Ethical IAS (medium-term)

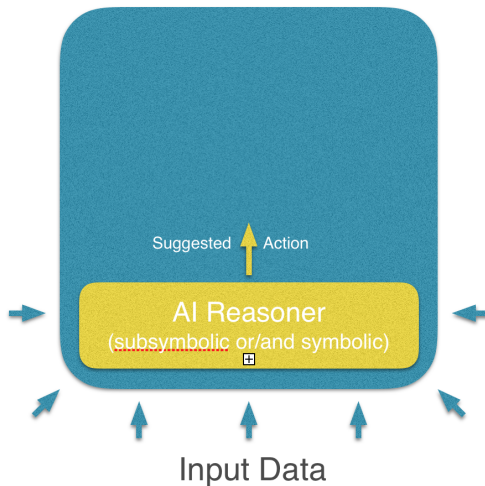


## Pseudo-Ethical IAS (medium-term)

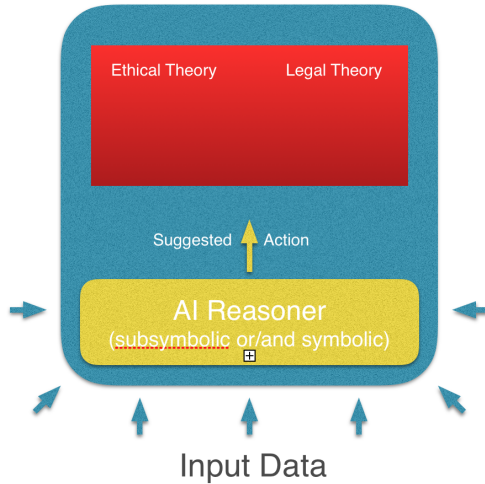




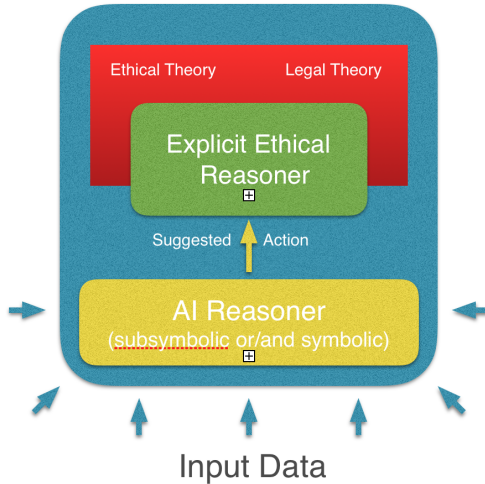
## Pseudo-Ethical IAS (medium-term)



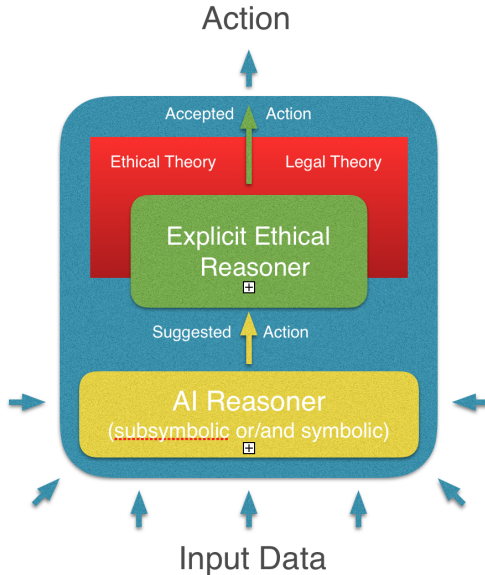
## Pseudo-Ethical IAS (medium-term)



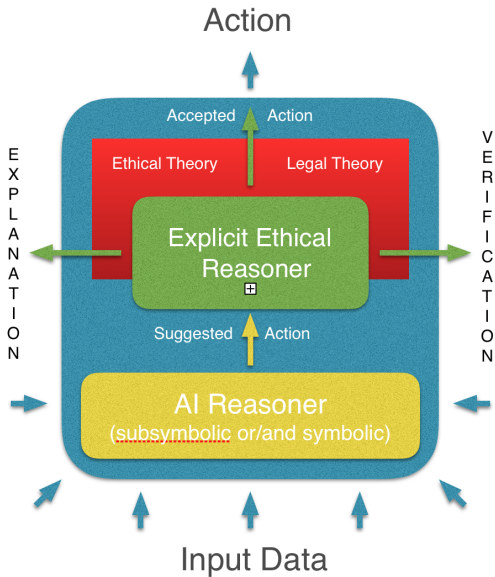
## Pseudo-Ethical IAS (medium-term)



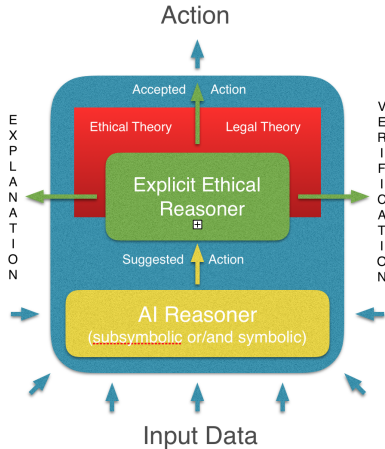
## Pseudo-Ethical IAS (medium-term)



## Pseudo-Ethical IAS (medium-term)



## Pseudo-Ethical IAS (medium-term)



### Related Work

- ▶ Artificial Moral Agents
  - ▶ [Wallach&Allen, 2008]
- ▶ Ethical Governors
  - ▶ [ArkinEtAl., 2009, 2012]
  - ▶ [Dennis&Fisher, 2017]
- ▶ Ethical Deliberation in ART
  - ▶ [Dignum, 2017]
- ▶ Programming Machine Ethics
  - ▶ [Pereira&Saptawijaya, 2016]
- ▶ ...

## Which Reasoning Formalisms?

“If we had it [a *characteristica universalis*], we should be able to reason in metaphysics and morals in much the same way as in geometry and analysis.”

(Leibniz, 1677)

### Challenges for Explicit Ethical Reasoning Engines: Which Logic(s)?

- ▶ Dilemmas, conflicting theories, etc.
- ▶ Appropriate handling of notion of **obligation**
  - ▶ Contrary-to-duty (CTD) scenarios

## Which Reasoning Formalisms?

“If we had it [a *characteristica universalis*], we should be able to reason in metaphysics and morals in much the same way as in geometry and analysis.”

(Leibniz, 1677)

### Challenges for Explicit Ethical Reasoning Engines: Which Logic(s)?

- ▶ Dilemmas, conflicting theories, etc.
- ▶ Appropriate handling of notion of **obligation**
  - ▶ Contrary-to-duty (CTD) scenarios

#### Standard CTD structure (Chisholm)

1. obligatory '*a*'
  2. obligatory 'if *a* then not *b*'
  3. if 'not *a*' then obligatory '*b*'
  4. 'not *a*'
- (in a given situation)

**Danger:** Paradox/inconsistency — ex falso quodlibet!



## Which Reasoning Formalisms?

“If we had it [a *characteristica universalis*], we should be able to reason in metaphysics and morals in much the same way as in geometry and analysis.”

(Leibniz, 1677)

### Challenges for Explicit Ethical Reasoning Engines: Which Logic(s)?

- ▶ Dilemmas, conflicting theories, etc.
- ▶ Appropriate handling of notion of **obligation**
  - ▶ Contrary-to-duty (CTD) scenarios

#### CTD example (X. Parent): **EU General Data Protection Regulation (GDPR)**

1. Personal data shall be processed lawfully. (Art. 5)  
E.g., the data subject must have given consent to the processing. (Art. 6/1.a)
2. **Implicit:** The data shall be kept, for the agreed purposes, if processed lawfully.
3. If personal data has been processed unlawfully, the controller has the obligation to erase the personal data in question without delay. (Art. 17.d, right to be forgotten)
4. **Given situation:** Some personal data has been processed unlawfully.

**Danger:** Paradox/inconsistency — ex falso quodlibet!

## Which Reasoning Formalisms?

“If we had it [a *characteristica universalis*], we should be able to reason in metaphysics and morals in much the same way as in geometry and analysis.”

(Leibniz, 1677)

### Challenges for Explicit Ethical Reasoning Engines: Which Logic(s)?

- ▶ Dilemmas, conflicting theories, etc.
- ▶ Appropriate handling of notion of **obligation**
  - ▶ Contrary-to-duty (CTD) scenarios

#### Deontic Logic

- ▶ Reasoning about obligations and permissions
- ▶ Two groups of approaches:
  - Possible worlds
    - ▶ standard deontic logic CTD: **no**
    - ▶ dyadic deontic logic CTD: **yes**
  - Norm-based semantics
    - ▶ input/output logic CTD: **yes**



L. van der Torre



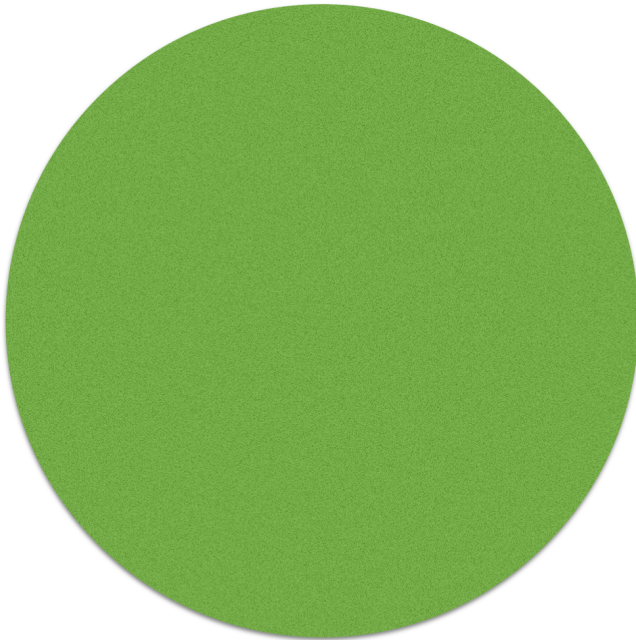
X. Parent



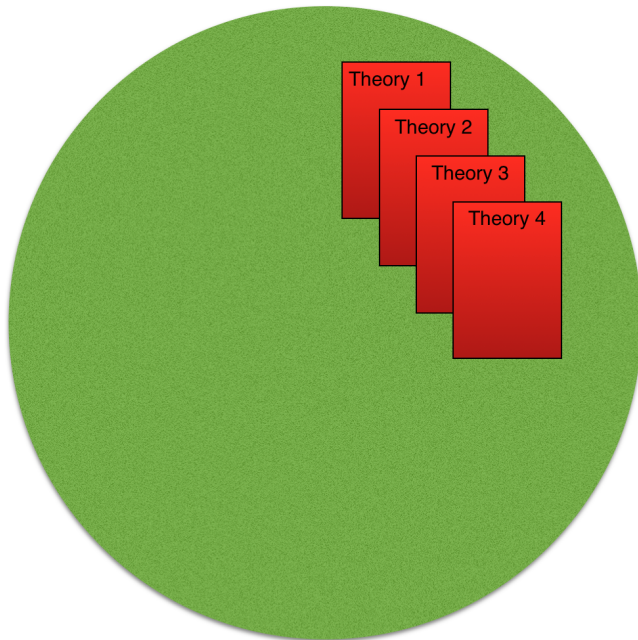
A. Farjami

### Further interests and challenges

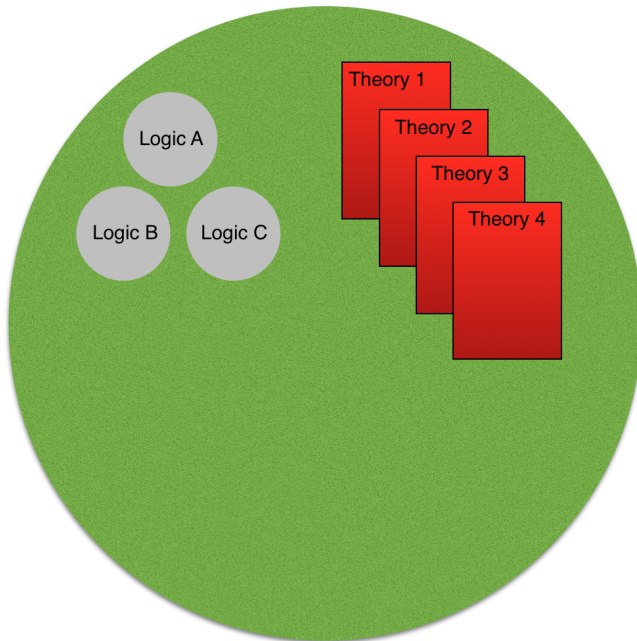
- ▶ Combination with other logics (other modalities)
- ▶ Propositional deontic logic(s) will hardly be sufficient in practice



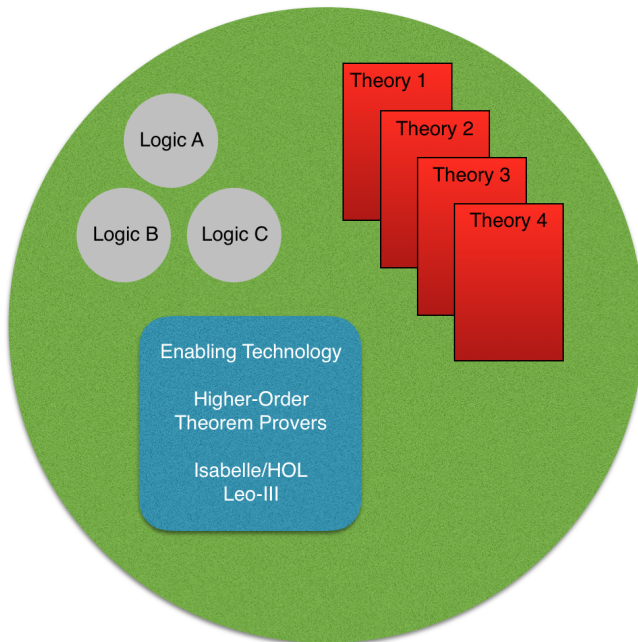
## Normative Reasoning Experimentation Platform



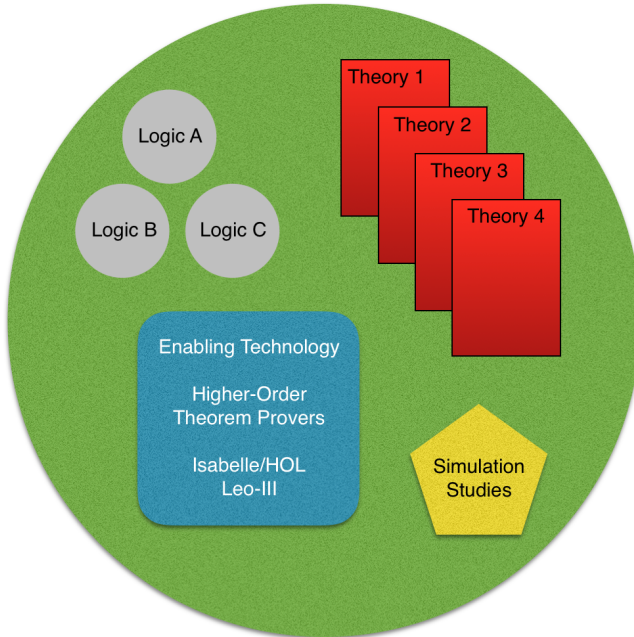
## Normative Reasoning Experimentation Platform



# Normative Reasoning Experimentation Platform



# Normative Reasoning Experimentation Platform



# Normative Reasoning Experimentation Platform — Demo in Isabelle/HOL

GDPR.thy

GDPR.thy (~:/chris/trunk/tex/talks/2018-DEON/DEMO/)

```
1 theory GDPR imports SDL                                (* Christoph Benzmüller & Xavier Parent, 2018 *)
2
3 begin (** GDPR Example **)
4   consts process_data_lawfully:: $\sigma$  erase_data:: $\sigma$  kill_boss:: $\sigma$ 
5
6   axiomatization where
7     (* It is an obligation to process data lawfully. *)
8     A1: "[0(process_data_lawfully)]" and
9     (* Implicit: It is an obligation to keep the data if it was processed lawfully. *)
10    Implicit: "[0(process_data_lawfully  $\rightarrow$   $\neg$ erase_data)]" and
11    (* If data was not processed lawfully, then it is an obligation to erase the data. *)
12    A2: "[ $\neg$ process_data_lawfully  $\rightarrow$  0(erase_data)]"
13    (* Given a situation where data is processed unlawfully. *) and
14    A3: "[ $\neg$ process_data_lawfully]cv"
15
16  (** Some Experiments **)
17  lemma True nitpick [satisfy] oops (* Consistency-check: Is there a model? *)
18  lemma False sledgehammer oops (* Inconsistency-check: Can Falsum be derived? *)
19
20  lemma "[0(erase_data)]" sledgehammer nitpick oops (* Should the data be erased? *)
21  lemma "[0( $\neg$ erase_data)]" sledgehammer nitpick oops (* Should the data be kept? *)
22  lemma "[0(kill_boss)]" sledgehammer nitpick oops (* Should the boss be killed? *)
23 end
```

Documentation

Sidekick

State

Theories

☒ Proof state ☒ Auto update  Search:

Sledgehammering...

Proof found...

"spass": The prover derived "False" from "A1", "A2", "A3", "D", and "Implicit", which could

"e": The prover derived "False" from "A1", "A2", "A3", "D", and "Implicit", which could be

"cvc4": Try this: by (metis A1 A2 A3 D Implicit) (68 ms)

"z3": Try this: by (metis A1 A2 A3 D Implicit) (59 ms)

☒ Output ☐ Query ☐ Sledgehammer ☐ Symbols

C. Benzmüller, 2018

9



# Normative Reasoning Experimentation Platform — Demo in Isabelle/HOL

GDPR.thy

GDPR.thy (~/chris/trunk/tex/talks/2018-DEON/DEMO/)

```
1 theory GDPR imports SDL (* Christoph Benzmüller & Xavier Parent, 2018 *)
2
3 begin (** GDPR Example **)
4 consts process_data_lawfully:: $\sigma$  erase_data:: $\sigma$  kill_boss:: $\sigma$ 
5
6 axiomatization where
7   (* It is an obligation to process data lawfully. *)
8   A1: "[0(process_data_lawfully)]" and
9   (* Implicit: It is an obligation to keep the data if it was processed lawfully. *)
10  Implicit: "[0(process_data_lawfully  $\rightarrow$   $\neg$ erase_data)]" and
11  (* If data was not processed lawfully, then it is an obligation to erase the data. *)
12  A2: "[ $\neg$ process_data_lawfully  $\rightarrow$  0(erase_data)]"
13  (* Given a situation where data is processed unlawfully. *) and
14  A3: "[ $\neg$ process_data_lawfully]cv"
15
16 (*
17  l
18  l
19  l
20  l
21 lemma "[0( $\neg$ erase_data)]" sledgehammer nitpick oops (* Should the data be kept? *)
22 lemma "[0(kill_boss)]" sledgehammer nitpick oops (* Should the boss be killed? *)
23 end
```

**Danger Zone:  
Paradoxes and Inconsistencies!**

☒ Proof state ☒ Auto update Update Search: 100%

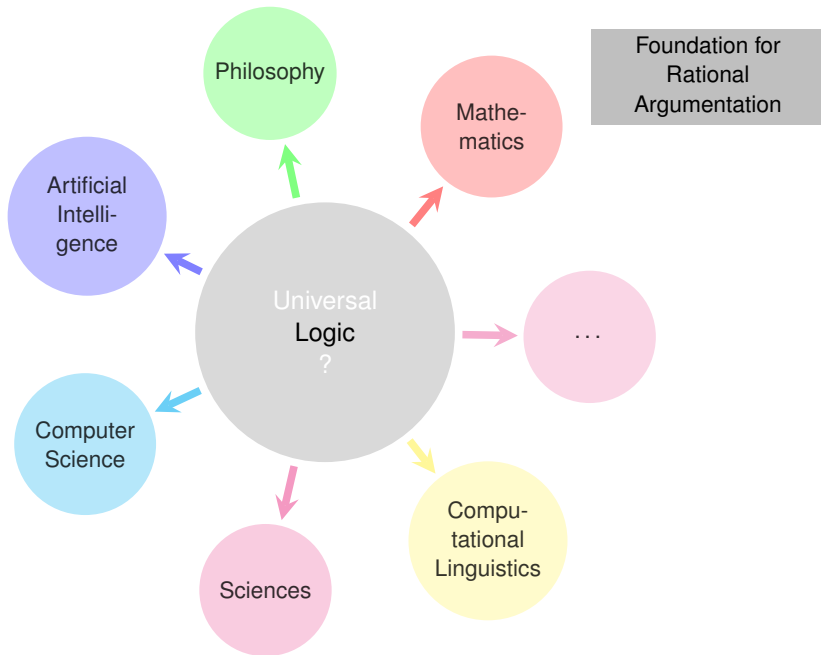
Sledgehammering...  
Proof found...  
"spass": The prover derived "False" from "A1", "A2", "A3", "D", and "Implicit", which could be  
"e": The prover derived "False" from "A1", "A2", "A3", "D", and "Implicit", which could be  
"cvc4": Try this: by (metis A1 A2 A3 D Implicit) (68 ms)  
"z3": Try this: by (metis A1 A2 A3 D Implicit) (59 ms)

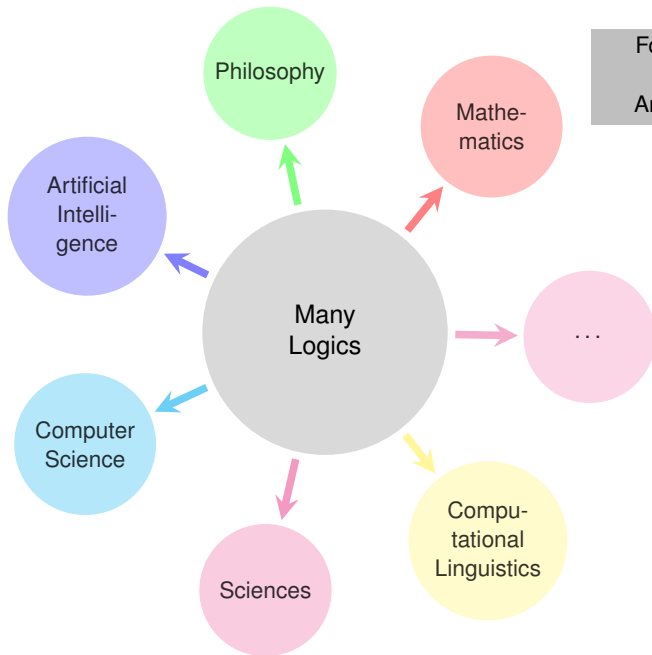
Output Query Sledgehammer Symbols

“If we had it [*a characteristica universalis*], we should be able to reason in metaphysics and morals in much the same way as in geometry and analysis.”

(Leibniz, 1677)

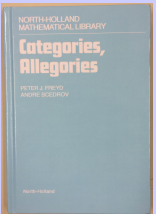
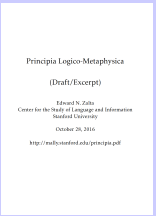
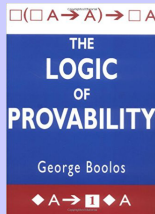
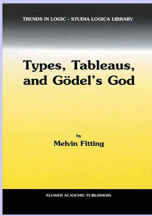
**Part B — Technology:  
Universal Logical Reasoning in Higher-Order Logic**







## Logic Zoo







STUDIES IN LOGIC  
AND  
PRACTICAL REASONING

VOLUME 3

D.M. GABBAY / P. GARDENFORS / J. SIEKMANN / J. VAN BENTHEM / M. VARDI / J. WOODS

EDITORS

---

*Handbook of  
Modal Logic*



### 2 BASIC MODAL LOGIC

In this section we introduce the basic modal language and its relational semantics. We define basic modal syntax, introduce models and frames, and give the satisfaction definition. We then draw the reader's attention to the internal perspective that modal languages offer on relational structure, and explain why models and frames should be thought of as graphs. Following this we give the standard translation. This enables us to convert any basic modal formula into a first-order formula with one free variable. The standard translation is a bridge between the modal and classical worlds, a bridge that underlies much of the work of this chapter.

#### 2.1 First steps in relational semantics

Suppose we have a set of proposition symbols (whose elements we typically write as  $p, q, r$  and so on) and a set of modality symbols (whose elements we typically write as  $m, m', m''$ , and so on). The choice of PROP and MOD is called the *signature* (or *similarity type*) of the language; in what follows we'll tacitly assume that PROP is denumerably infinite, and we'll often work with signatures in which MOD contains only a single element. Given a signature, we define the *basic modal language* (over the signature) as follows:

$$\varphi ::= p \mid \top \mid \perp \mid \neg\varphi \mid \varphi \wedge \psi \mid \varphi \vee \psi \mid \varphi \rightarrow \psi \mid \varphi \leftrightarrow \psi \mid \langle m \rangle \varphi \mid [m] \varphi.$$

That is, a basic modal formula is either a proposition symbol, a boolean constant, a boolean combination of basic modal formulas, or (most interesting of all) a formula prefixed by a diamond

### 2 BASIC MODAL LOGIC

In this section we introduce the basic modal language and its relational semantics. We define basic modal syntax, introduce models and frames, and give the satisfaction definition. We then draw the reader's attention to the internal perspective that modal languages offer on relational structure, and explain why models and frames should be thought of as graphs. Following this we give the standard translation. This enables us to convert any basic modal formula into a first-order formula with one free variable. The standard translation is a bridge between the modal and classical worlds, a bridge that underlies much of the work of this chapter.

#### 2.1 First steps in relational semantics

### Metalanguage

What follows we tacitly assume that  $\text{PROP}$  is denumerably infinite, and we'll often work with signatures in which  $\text{MOD}$  contains only a single element. Given a signature, we define the *basic modal language* (over the signature) as follows:

$$\varphi ::= p \mid \top \mid \perp \mid \neg\varphi \mid \varphi \wedge \psi \mid \varphi \vee \psi \mid \varphi \rightarrow \psi \mid \varphi \leftrightarrow \psi \mid \langle m \rangle \varphi \mid [m] \varphi.$$

That is, a basic modal formula is either a proposition symbol, a boolean constant, a boolean combination of basic modal formulas, or (most interesting of all) a formula prefixed by a diamond

### Syntax

use elements we typically write as  $p, q, r$  and  
nts we typically write as  $m, m', m''$ , and so  
*nature* (or *similarity type*) of the language; in

## Example: Modal Logic Textbook

A *model* (or *Kripke model*)  $\mathfrak{M}$  for the basic modal language (over some fixed signature) is a triple  $\mathfrak{M} = (W, \{R^m\}_{m \in \text{MOD}}, V)$ . Here  $W$ , the *domain*, is a non-empty set, whose elements we usually call *points*, but which, for reasons which will soon be clear, are sometimes called *states*, *times*, *situations*, *worlds* and other things besides. Each  $R^m$  in a model is a binary relation on  $W$ , and  $V$  is a function (the valuation) that assigns to each proposition symbol  $p$  in PROP a subset  $V(p)$  of  $W$ ; think of  $V(p)$  as the set of points in  $\mathfrak{M}$  where  $p$  is true. The first two components  $(W, \{R^m\}_{m \in \text{MOD}})$  of  $\mathfrak{M}$  are called the *frame* underlying the model. If there is only one relation in the model, we typically write  $(W, R)$  for its frame, and  $(W, R, V)$  for the model itself. We encourage the reader to think of Kripke models as graphs (or to be slightly more precise, *directed graphs*, that is, graphs whose points are linked by directed arrows) and will shortly give some examples which show why this is helpful.

Suppose  $w$  is a point in a model  $\mathfrak{M} = (W, \{R^m\}_{m \in \text{MOD}}, V)$ . Then we inductively define the notion of a formula  $\varphi$  being *satisfied* (or *true*) in  $\mathfrak{M}$  at point  $w$  as follows (we omit some of the clauses for the booleans):

$\mathfrak{M}, w \models p$	iff	$w \in V(p)$ ,
$\mathfrak{M}, w \models \top$		always,
$\mathfrak{M}, w \models \perp$		never,
$\mathfrak{M}, w \models \neg\varphi$	iff	not $\mathfrak{M}, w \models \varphi$ (notation: $\mathfrak{M}, w \not\models \varphi$ ),
$\mathfrak{M}, w \models \varphi \wedge \psi$	iff	$\mathfrak{M}, w \models \varphi$ and $\mathfrak{M}, w \models \psi$ ,
$\mathfrak{M}, w \models \varphi \rightarrow \psi$	iff	$\mathfrak{M}, w \not\models \varphi$ or $\mathfrak{M}, w \models \psi$ ,
$\mathfrak{M}, w \models \langle m \rangle \varphi$	iff	for some $v \in W$ such that $R^m wv$ we have $\mathfrak{M}, v \models \varphi$ ,
$\mathfrak{M}, w \models [m] \varphi$	iff	for all $v \in W$ such that $R^m wv$ we have $\mathfrak{M}, v \models \varphi$ .

## Example: Modal Logic Textbook

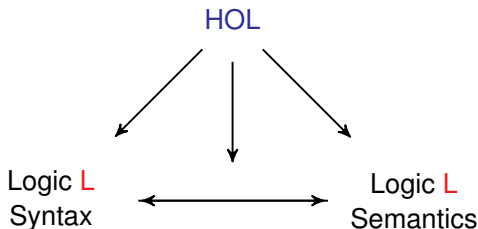
A *model* (or *Kripke model*)  $\mathfrak{M}$  for the basic modal language (over some fixed signature) is a triple  $\mathfrak{M} = (W, \{R^m\}_{m \in \text{MOD}}, V)$ . Here  $W$ , the *domain*, is a non-empty set, whose elements we usually call *points*, but which, for reasons which will soon be clear, are sometimes called *states*, *times* and  $\mathbb{I}$ .  $V(p)$  is the set of points in the model where the proposition symbol  $p$  in PROP is true. The first two components of a model are  $(W, R)$  and  $V$ . We encourage the reader to think of Kripke models as graphs (or to be slightly more precise, *directed graphs*, that is, graphs whose points are linked by directed arrows) and will shortly give some examples which show why this is helpful.

Suppose  $w$  is a point in a model  $\mathfrak{M} = (W, \{R^m\}_{m \in \text{MOD}}, V)$ . Then we inductively define the notion of a formula  $\varphi$  being *satisfied* (or *true*) in  $\mathfrak{M}$  at point  $w$  as follows (we omit some of the clauses for the booleans):

$\mathfrak{M}, w \models p$	iff	$w \in V(p)$ ,
$\mathfrak{M}, w \models \top$		always,
$\mathfrak{M}, w \models \perp$		never,
$\mathfrak{M}, w \models \neg\varphi$	iff	not $\mathfrak{M}, w \models \varphi$ (notation: $\mathfrak{M}, w \not\models \varphi$ ),
$\mathfrak{M}, w \models \varphi \wedge \psi$	iff	$\mathfrak{M}, w \models \varphi$ and $\mathfrak{M}, w \models \psi$ ,
$\mathfrak{M}, w \models \varphi \rightarrow \psi$	iff	$\mathfrak{M}, w \not\models \varphi$ or $\mathfrak{M}, w \models \psi$ ,
$\mathfrak{M}, w \models \langle m \rangle \varphi$	iff	for some $v \in W$ such that $R^m wv$ we have $\mathfrak{M}, v \models \varphi$ ,
$\mathfrak{M}, w \models [m] \varphi$	iff	for all $v \in W$ such that $R^m wv$ we have $\mathfrak{M}, v \models \varphi$ .

## Semantics

# Universal Logical Reasoning in Meta-Logic HOL



**Examples for **L** we have already studied:**

Intuitionistic Logics, (Mathematical) Fuzzy Logics, Free Logic, Modal Logics, Description Logics, Conditional Logics, Access Control Logics, Hybrid Logics, Multivalued Logics, Logics with Neighborhood Semantics, Paraconsistent Logics, Dyadic Deontic Logic, ...

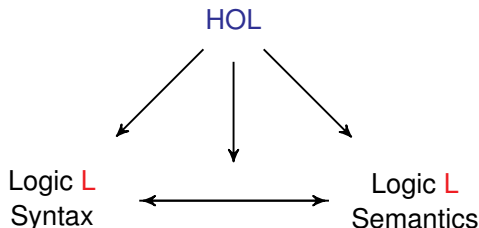
**Embedding works also for quantifiers (first-order & higher-order)**

**HOL provers become universal logic reasoning engines!**

interactive: Isabelle/HOL, PVS, HOL4, Hol Light, Coq/HOL, ...

automated: Leo-III, LEO-II, Satallax, TPS, Nitpick, Isabelle/HOL, ...

# Universal Logical Reasoning in Meta-Logic HOL



## Examples for **L** we have already studied:

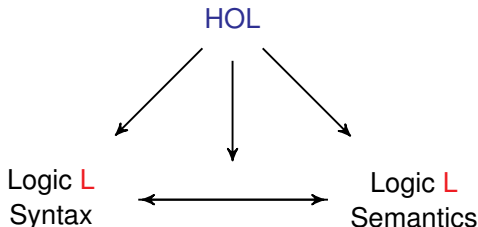
Intuitionistic Logics, (Mathematical) Fuzzy Logics, Free Logic, Modal Logics, Description Logics, Conditional Logics, Access Control Logics, Hybrid Logics, Multivalued Logics, Logics with Neighborhood Semantics, Paraconsistent Logics, Dyadic Deontic Logic, ...

## Embedding works also for quantifiers (first-order & higher-order)

## HOL provers become universal logic reasoning engines!

interactive: Isabelle/HOL, PVS, HOL4, Hol Light, Coq/HOL, ...

automated: Leo-III, LEO-II, Satallax, TPS, Nitpick, Isabelle/HOL, ...



### Examples for **L** we have already studied:

Intuitionistic Logics, (Mathematical) Fuzzy Logics, Free Logic, Modal Logics, Description Logics, Conditional Logics, Access Control Logics, Hybrid Logics, Multivalued Logics, Logics with Neighborhood Semantics, Paraconsistent Logics, Dyadic Deontic Logic, ...

**Embedding works also for quantifiers (first-order & higher-order)**

### **HOL provers become universal logic reasoning engines!**

interactive: Isabelle/HOL, PVS, HOL4, Hol Light, Coq/HOL, ...

automated: Leo-III, LEO-II, Satallax, TPS, Nitpick, Isabelle/HOL, ...

# Isabelle/HOL (one of various Theorem Provers for HOL)



## Isabelle

[Home](#)[Overview](#)[Installation](#)[Documentation](#)

### Site Mirrors:

Cambridge (.uk)  
Munich (.de)  
Sydney (.au)  
Potsdam, NY (.us)

## What is Isabelle?

Isabelle is a generic proof assistant. It allows mathematical formulas to be expressed in a formal language and provides tools for proving those formulas in a logical calculus. Isabelle was originally developed at the [University of Cambridge](#) and [Technische Universität München](#), but now includes numerous contributions from institutions and individuals worldwide. See the [Isabelle overview](#) for a brief introduction.

## Now available: Isabelle2017 (October 2017)



[Download for Linux](#) - [Download for Windows \(32bit\)](#) - [Download for Windows \(64bit\)](#) - [Download for Mac OS X](#)

### Some notable changes:

- Experimental support for Visual Studio Code as alternative PIDE front-end.
- Improved Isabelle/Edit Prover IDE: management of session sources independently of editor buffers, removal of unused theories, explicit indication of theory status, more careful auto-indentation.
- Session-qualified theory imports.
- Code generator improvements: support for statically embedded computations.
- Numerous HOL library improvements.
- More material in HOL-Algebra, HOL-Computational\_Algebra and HOL-Analysis (ported from HOL-Light).
- Improved Nunchaku model finder, now in main HOL.
- SQL database support in Isabelle/Scala.

See also the cumulative [NEWS](#).

## Distribution & Support

Isabelle is distributed for free under a conglomerate of open-source licenses, but the main code-base is subject to BSD-style regulations. The application bundles include source and binary packages and documentation, see the detailed [installation instructions](#). A vast collection of Isabelle examples and applications is available from the [Archive of Formal Proofs](#).

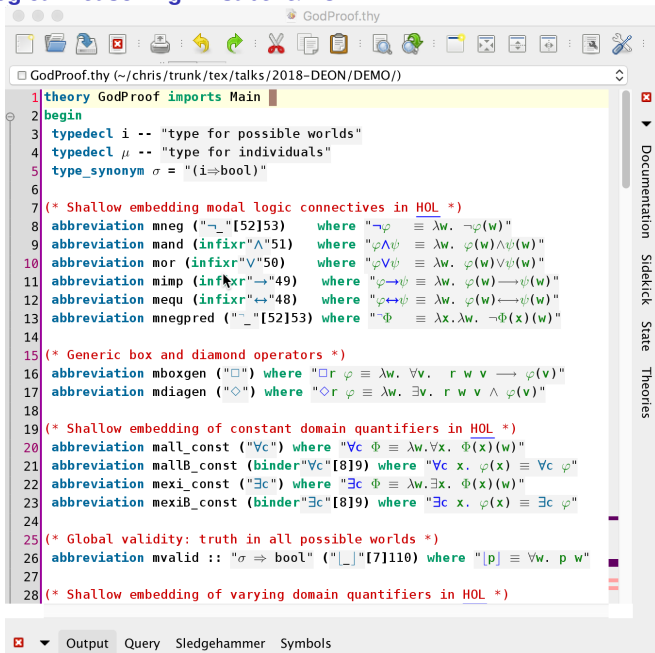
<https://isabelle.in.tum.de>

many other systems:

Coq, HOL, HOL Light, PVS, Lean, NuPrL, IMPS, ACL2, **Leo-II/Leo-III**, . . .



# Universal Logical Reasoning in Isabelle/HOL



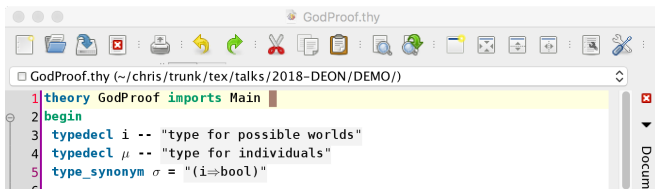
```
1 theory GodProof imports Main
2 begin
3 typedecl i -- "type for possible worlds"
4 typedecl μ -- "type for individuals"
5 type_synonym σ = "(i⇒bool)"
6
7 (* Shallow embedding modal logic connectives in HOL *)
8 abbreviation mneg ("¬"[52]53) where "¬φ ≡ λw. ¬φ(w)"
9 abbreviation mand (infixr"∧"51) where "φ∧ψ ≡ λw. φ(w)∧ψ(w)"
10 abbreviation mor (infixr"∨"50) where "φ∨ψ ≡ λw. φ(w)∨ψ(w)"
11 abbreviation mimp (infixr"→"49) where "φ→ψ ≡ λw. φ(w)→ψ(w)"
12 abbreviation mequ (infixr"↔"48) where "φ↔ψ ≡ λw. φ(w)↔ψ(w)"
13 abbreviation mnegpred ("¬"[52]53) where "¬Φ ≡ λx.λw. ¬Φ(x)(w)"
14
15 (* Generic box and diamond operators *)
16 abbreviation mboxgen ("□") where "□r φ ≡ λw. ∀v. r w v → φ(v)"
17 abbreviation mdiagen ("◇") where "◇r φ ≡ λw. ∃v. r w v ∧ φ(v)"
18
19 (* Shallow embedding of constant domain quantifiers in HOL *)
20 abbreviation mall_const ("∀c") where "∀c Φ ≡ λw.∀x. Φ(x)(w)"
21 abbreviation mallB_const (binder"∀c"[8]9) where "∀c x. φ(x) ≡ ∀c φ"
22 abbreviation mexi_const ("∃c") where "∃c Φ ≡ λw.∃x. Φ(x)(w)"
23 abbreviation mexiB_const (binder"∃c"[8]9) where "∃c x. φ(x) ≡ ∃c φ"
24
25 (* Global validity: truth in all possible worlds *)
26 abbreviation mvalid :: "σ ⇒ bool" ("⊨"[7]110) where "⊨p ≡ ∀w. p w"
27
28 (* Shallow embedding of varying domain quantifiers in HOL *)
```

GodProof.thy (~/.chris/trunk/tex/talks/2018-DEON/DEMO/)

Documentation Sidekick State Theories

Output Query Sledgehammer Symbols

# Universal Logical Reasoning in Isabelle/HOL



```

1 theory GodProof imports Main
2 begin
3 typedef i -- "type for possible worlds"
4 typedef mu -- "type for individuals"
5 type_synonym sigma = "(i => bool)"
6

```

$\varphi_o$  lifted to  $\varphi_{i \rightarrow o}$  ("truth sets")

$$\varphi_{i \rightarrow o} \vee \psi_{i \rightarrow o} = \lambda w_i. (\varphi w \vee \psi w) \quad \text{encodes: } \{w \mid w \in \varphi \text{ or } w \in \psi\}$$

$$\vee = \lambda \varphi_{i \rightarrow o}. \lambda \psi_{i \rightarrow o}. \lambda w_i. (\varphi w \vee \psi w)$$

$$\Box \varphi_{i \rightarrow o} = \lambda w_i. \forall y_i. (w \ r \ y \rightarrow \varphi y)$$

$$\Box = \lambda \varphi_{i \rightarrow o}. \lambda w_i. \forall y_i. (w \ r \ y \rightarrow \varphi y)$$

In HOL  $\forall x_\mu. \varphi_o$  is syntactic sugar for  $\Pi (\overbrace{\lambda x_\mu. \varphi_o}^{\phi_{\mu \rightarrow o}})$

$$\Pi \phi_{\mu \rightarrow (i \rightarrow o)} = \lambda w_i. \forall x_\mu. (\phi \ x \ w)$$

$$\Pi = \lambda \phi_{\mu \rightarrow (i \rightarrow o)}. \lambda w_i. \forall x_\mu. (\phi \ x \ w)$$

# Universal Logical Reasoning in Isabelle/HOL

```

1 theory GodProof imports Main
2 begin
3   typedecl i -- "type for possible worlds"
4   typedecl μ -- "type for individuals"
5   type_synonym σ = "(i ⇒ bool)"
6

```

$\varphi_o$  lifted to  $\varphi_{i \rightarrow o}$  ("truth sets")

$\varphi_{i \rightarrow o} \vee \psi_{i \rightarrow o} = \lambda w_i. (\varphi w \vee \psi w)$  encodes:  $\{w \mid w \in \varphi \text{ or } w \in \psi\}$

$\vee = \lambda \varphi_{i \rightarrow o} \lambda \psi_{i \rightarrow o} \lambda w_i. (\varphi w \vee \psi w)$

$\Box \varphi_{i \rightarrow o} = \lambda w_i \forall y_i. (w \text{ r } y \rightarrow \varphi y)$

$\Box = \lambda \varphi_{i \rightarrow o} \lambda w_i \forall y_i. (w \text{ r } y \rightarrow \varphi y)$

In HOL  $\forall x_\mu. \varphi_o$  is syntactic sugar for  $\Pi \overbrace{(\lambda x_\mu \varphi_o)}^{\phi_{\mu \rightarrow o}}$

$\Pi \phi_{\mu \rightarrow (i \rightarrow o)} = \lambda w_i \forall x_\mu. (\phi x w)$

$\Pi = \lambda \phi_{\mu \rightarrow (i \rightarrow o)} \lambda w_i \forall x_\mu. (\phi x w)$

# Universal Logical Reasoning in Isabelle/HOL

```

1 theory GodProof imports Main
2 begin
3   typedecl i -- "type for possible worlds"
4   typedecl μ -- "type for individuals"
5   type_synonym σ = "(i ⇒ bool)"
6

```

$\varphi_o$  lifted to  $\varphi_{i \rightarrow o}$  ("truth sets")

$\varphi_{i \rightarrow o} \vee \psi_{i \rightarrow o} = \lambda w_i. (\varphi w \vee \psi w)$  encodes:  $\{w \mid w \in \varphi \text{ or } w \in \psi\}$

$\vee = \lambda \varphi_{i \rightarrow o} \lambda \psi_{i \rightarrow o} \lambda w_i. (\varphi w \vee \psi w)$

$\Box \varphi_{i \rightarrow o} = \lambda w_i \forall y_i. (w \text{ r } y \rightarrow \varphi y)$

$\Box = \lambda \varphi_{i \rightarrow o} \lambda w_i \forall y_i. (w \text{ r } y \rightarrow \varphi y)$

In HOL  $\forall x_\mu. \varphi_o$  is syntactic sugar for  $\Pi \overbrace{(\lambda x_\mu \varphi_o)}^{\phi_{\mu \rightarrow o}}$

$\Pi \phi_{\mu \rightarrow (i \rightarrow o)} = \lambda w_i \forall x_\mu. (\phi x w)$

$\Pi = \lambda \phi_{\mu \rightarrow (i \rightarrow o)} \lambda w_i \forall x_\mu. (\phi x w)$

# Universal Logical Reasoning in Isabelle/HOL

```

1 theory GodProof imports Main
2 begin
3   typedecl i -- "type for possible worlds"
4   typedecl μ -- "type for individuals"
5   type_synonym σ = "(i ⇒ bool)"
6

```

$\varphi_o$  lifted to  $\varphi_{i \rightarrow o}$  ("truth sets")

$\varphi_{i \rightarrow o} \vee \psi_{i \rightarrow o} = \lambda w_i. (\varphi w \vee \psi w)$  encodes:  $\{w \mid w \in \varphi \text{ or } w \in \psi\}$

$\vee = \lambda \varphi_{i \rightarrow o} \lambda \psi_{i \rightarrow o} \lambda w_i. (\varphi w \vee \psi w)$

$\Box \varphi_{i \rightarrow o} = \lambda w_i \forall y_i. (w \text{ r } y \rightarrow \varphi y)$

$\Box = \lambda \varphi_{i \rightarrow o} \lambda w_i \forall y_i. (w \text{ r } y \rightarrow \varphi y)$

In HOL  $\forall x_\mu. \varphi_o$  is syntactic sugar for  $\Pi \overbrace{(\lambda x_\mu \varphi_o)}^{\phi_{\mu \rightarrow o}}$

$\Pi \phi_{\mu \rightarrow (i \rightarrow o)} = \lambda w_i \forall x_\mu. (\phi x w)$

$\Pi = \lambda \phi_{\mu \rightarrow (i \rightarrow o)} \lambda w_i \forall x_\mu. (\phi x w)$

# Universal Logical Reasoning in Isabelle/HOL

```

1 theory GodProof imports Main
2 begin
3   typedecl i -- "type for possible worlds"
4   typedecl μ -- "type for individuals"
5   type_synonym σ = "(i ⇒ bool)"
6

```

$\varphi_o$  lifted to  $\varphi_{i \rightarrow o}$  ("truth sets")

$\varphi_{i \rightarrow o} \vee \psi_{i \rightarrow o} = \lambda w_i. (\varphi w \vee \psi w)$  encodes:  $\{w \mid w \in \varphi \text{ or } w \in \psi\}$

$\vee = \lambda \varphi_{i \rightarrow o} \lambda \psi_{i \rightarrow o} \lambda w_i. (\varphi w \vee \psi w)$

$\Box \varphi_{i \rightarrow o} = \lambda w_i. \forall y_i. (w \text{ r } y \rightarrow \varphi y)$

$\Box = \lambda \varphi_{i \rightarrow o} \lambda w_i. \forall y_i. (w \text{ r } y \rightarrow \varphi y)$

In HOL  $\forall x_\mu. \varphi_o$  is syntactic sugar for  $\Pi \overbrace{(\lambda x_\mu. \varphi_o)}^{\phi_{\mu \rightarrow o}}$

$\Pi \phi_{\mu \rightarrow (i \rightarrow o)} = \lambda w_i. \forall x_\mu. (\phi x w)$

$\Pi = \lambda \phi_{\mu \rightarrow (i \rightarrow o)} \lambda w_i. \forall x_\mu. (\phi x w)$

# Universal Logical Reasoning in Isabelle/HOL

```

1 theory GodProof imports Main
2 begin
3   typedecl i -- "type for possible worlds"
4   typedecl μ -- "type for individuals"
5   type_synonym σ = "(i ⇒ bool)"
6

```

$\varphi_o$  lifted to  $\varphi_{i \rightarrow o}$  ("truth sets")

$\varphi_{i \rightarrow o} \vee \psi_{i \rightarrow o} = \lambda w_i. (\varphi w \vee \psi w)$  encodes:  $\{w \mid w \in \varphi \text{ or } w \in \psi\}$

$\vee = \lambda \varphi_{i \rightarrow o} \lambda \psi_{i \rightarrow o} \lambda w_i. (\varphi w \vee \psi w)$

$\Box \varphi_{i \rightarrow o} = \lambda w_i \forall y_i. (w \text{ r } y \rightarrow \varphi y)$

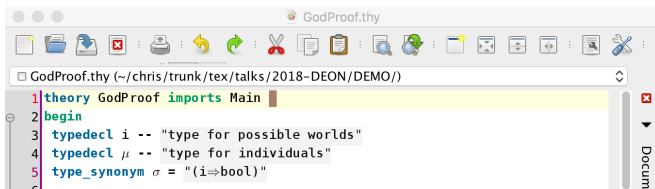
$\Box = \lambda \varphi_{i \rightarrow o} \lambda w_i \forall y_i. (w \text{ r } y \rightarrow \varphi y)$

In HOL  $\forall x_\mu. \varphi_o$  is syntactic sugar for  $\Pi \overbrace{(\lambda x_\mu \varphi_o)}^{\phi_{\mu \rightarrow o}}$

$\Pi \phi_{\mu \rightarrow (i \rightarrow o)} = \lambda w_i \forall x_\mu. (\phi x w)$

$\Pi = \lambda \phi_{\mu \rightarrow (i \rightarrow o)} \lambda w_i \forall x_\mu. (\phi x w)$

# Universal Logical Reasoning in Isabelle/HOL



```

1 theory GodProof imports Main
2 begin
3   typedecl i -- "type for possible worlds"
4   typedecl mu -- "type for individuals"
5   type_synonym sigma = "(i=>bool)"
6

```

$\varphi_o$  lifted to  $\varphi_{i \rightarrow o}$  ("truth sets")

$\varphi_{i \rightarrow o} \vee \psi_{i \rightarrow o} = \lambda w_i. (\varphi w \vee \psi w)$  encodes:  $\{w \mid w \in \varphi \text{ or } w \in \psi\}$

$\vee = \lambda \varphi_{i \rightarrow o} \lambda \psi_{i \rightarrow o} \lambda w_i. (\varphi w \vee \psi w)$

$\Box \varphi_{i \rightarrow o} = \lambda w_i \forall y_i. (w \text{ r } y \rightarrow \varphi y)$

$\Box = \lambda \varphi_{i \rightarrow o} \lambda w_i \forall y_i. (w \text{ r } y \rightarrow \varphi y)$

In HOL  $\forall x_\mu. \varphi_o$  is syntactic sugar for  $\Pi (\overbrace{\lambda x_\mu. \varphi_o}^{\phi_{\mu \rightarrow o}})$

$\Pi \phi_{\mu \rightarrow (i \rightarrow o)} = \lambda w_i \forall x_\mu. (\phi x w)$

$\Pi = \lambda \phi_{\mu \rightarrow (i \rightarrow o)} \lambda w_i \forall x_\mu. (\phi x w)$



# Universal Logical Reasoning in Isabelle/HOL

```

1 theory GodProof imports Main
2 begin
3   typedef i -- "type for possible worlds"
4   typedef μ -- "type for individuals"
5   type_synonym σ = "(i ⇒ bool)"
6

```

$\varphi_o$  lifted to  $\varphi_{i \rightarrow o}$  ("truth sets")

$\varphi_{i \rightarrow o} \vee \psi_{i \rightarrow o} = \lambda w_i. (\varphi w \vee \psi w)$  encodes:  $\{w \mid w \in \varphi \text{ or } w \in \psi\}$

$\vee = \lambda \varphi_{i \rightarrow o} \lambda \psi_{i \rightarrow o} \lambda w_i. (\varphi w \vee \psi w)$

$\Box \varphi_{i \rightarrow o} = \lambda w_i \forall y_i. (w \text{ r } y \rightarrow \varphi y)$

$\Box = \lambda \varphi_{i \rightarrow o} \lambda w_i \forall y_i. (w \text{ r } y \rightarrow \varphi y)$

In HOL  $\forall x_\mu. \varphi_o$  is syntactic sugar for  $\Pi \overbrace{(\lambda x_\mu \varphi_o)}^{\phi_{\mu \rightarrow o}}$

$\Pi \phi_{\mu \rightarrow (i \rightarrow o)} = \lambda w_i \forall x_\mu. (\phi x w)$

$\Pi = \lambda \phi_{\mu \rightarrow (i \rightarrow o)} \lambda w_i \forall x_\mu. (\phi x w)$

# Universal Logical Reasoning in Isabelle/HOL

```

1 theory GodProof imports Main
2 begin
3   typedecl i -- "type for possible worlds"
4   typedecl μ -- "type for individuals"
5   type_synonym σ = "(i ⇒ bool)"
6

```

$\varphi_o$  lifted to  $\varphi_{i \rightarrow o}$  ("truth sets")

$\varphi_{i \rightarrow o} \vee \psi_{i \rightarrow o} = \lambda w_i. (\varphi w \vee \psi w)$  encodes:  $\{w \mid w \in \varphi \text{ or } w \in \psi\}$

$\vee = \lambda \varphi_{i \rightarrow o} \lambda \psi_{i \rightarrow o} \lambda w_i. (\varphi w \vee \psi w)$

$\Box \varphi_{i \rightarrow o} = \lambda w_i \forall y_i. (w \text{ r } y \rightarrow \varphi y)$

$\Box = \lambda \varphi_{i \rightarrow o} \lambda w_i \forall y_i. (w \text{ r } y \rightarrow \varphi y)$

In HOL  $\forall x_\mu. \varphi_o$  is syntactic sugar for  $\Pi \overbrace{(\lambda x_\mu \varphi_o)}^{\phi_{\mu \rightarrow o}}$

$\Pi \phi_{\mu \rightarrow (i \rightarrow o)} = \lambda w_i \forall x_\mu. (\phi x w)$

$\Pi = \lambda \phi_{\mu \rightarrow (i \rightarrow o)} \lambda w_i \forall x_\mu. (\phi x w)$

# Universal Logical Reasoning in Isabelle/HOL

```

1 theory GodProof imports Main
2 begin
3   typedecl i -- "type for possible worlds"
4   typedecl μ -- "type for individuals"
5   type_synonym σ = "(i ⇒ bool)"
6

```

$\varphi_o$  lifted to  $\varphi_{i \rightarrow o}$  ("truth sets")

$\varphi_{i \rightarrow o} \vee \psi_{i \rightarrow o} = \lambda w_i. (\varphi w \vee \psi w)$  encodes:  $\{w \mid w \in \varphi \text{ or } w \in \psi\}$

$\vee = \lambda \varphi_{i \rightarrow o} \lambda \psi_{i \rightarrow o} \lambda w_i. (\varphi w \vee \psi w)$

$\Box \varphi_{i \rightarrow o} = \lambda w_i \forall y_i. (w \text{ r } y \rightarrow \varphi y)$

$\Box = \lambda \varphi_{i \rightarrow o} \lambda w_i \forall y_i. (w \text{ r } y \rightarrow \varphi y)$

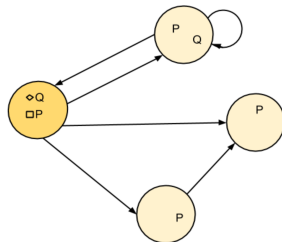
In HOL  $\forall x_\mu. \varphi_o$  is syntactic sugar for  $\Pi \overbrace{(\lambda x_\mu \varphi_o)}^{\phi_{\mu \rightarrow o}}$

$\Pi \phi_{\mu \rightarrow (i \rightarrow o)} = \lambda w_i \forall x_\mu. (\phi x w)$

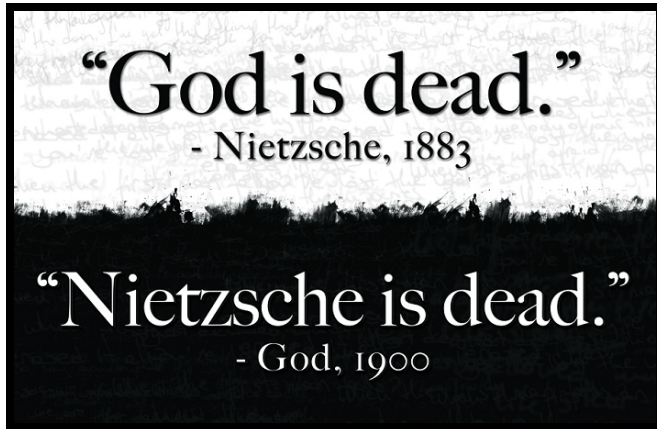
$\Pi = \lambda \phi_{\mu \rightarrow (i \rightarrow o)} \lambda w_i \forall x_\mu. (\phi x w)$

# Universal Logical Reasoning in Isabelle/HOL

Properties of  $\Box$  and  $\Diamond$  correlated to structure of transition system between worlds



- ▶ Logic K: — (no restrictions, any structure)
- ▶ Logic M: **reflexiv** transition relation,  $\forall P. \Box P \rightarrow P$
- ▶ Logic KB: **symmetric** transition relation,  $\forall P. P \rightarrow \Box \Diamond P$
- ▶ Logic S5: **equivelance relation** as transition system, add  $\forall P. \Box P \rightarrow \Box \Box P$
- ▶ Logic D: **serial** transition relation,  $\forall P. \Box P \rightarrow \Diamond P$  (Standard Deontic Logic)  
(alternatively:  $\forall P. \neg(\Box P \wedge \Box \neg P)$ )



**Part C — Evidence:  
Experiments in Computational Metaphysics**

[BenzmüllerWoltzenlogelPaleo, ECAI, 2014 + IJCAI, 2016 + KI 2016 + ...]

# Ontological Proofs of God's Existence

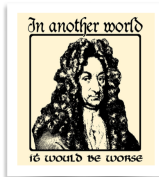
## A Long and Continuing Tradition in Philosophy



St. Anselm



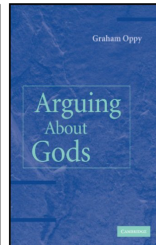
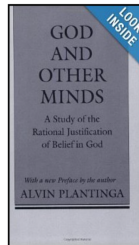
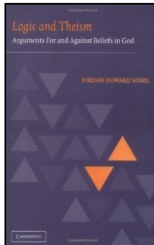
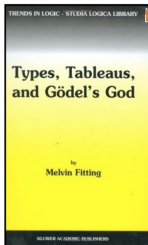
Descartes



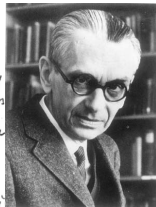
Leibniz



Gödel



# Computational Metaphysics: Kurt Gödel's Ontological Argument



Ontologischer Beweis

Feb 10, 1970

$P(\varphi)$   $\varphi$  is positive ( $\varphi \in P$ )

At 1  $P(\varphi) \cdot P(\psi) \supset P(\varphi \cdot \psi)$  At 2  $P(\varphi) \supset P(\Box \varphi)$

P1  $G(x) \equiv (\varphi) [P(\varphi) \supset \varphi(x)]$  (God)

P2  $\varphi \text{ Ess } x \equiv (\psi) [\psi(x) \supset N(\exists y [\varphi(y) \supset \psi(y)])]$  (Essence of x)

$P \supset Nq = N(p \supset q)$  Necessity

At 2  $P(\varphi) \supset NP(\varphi)$   
 $\sim P(\varphi) \supset N \sim P(\varphi)$  } because it follows from the nature of the property

Th.  $G(x) \supset G \text{ Ess } x$

Df.  $E(x) \equiv (p) [\varphi \text{ Ess } x \supset N \exists x \varphi(x)]$  necessary Existence

At 3  $P(E)$

Th.  $G(x) \supset N(\exists y) G(y)$

hence  $(\exists x) G(x) \supset N(\exists y) G(y)$

"  $M(\exists x) G(x) \supset M N(\exists y) G(y)$

"  $\supset N(\exists y) G(y)$

$M = possibility$

any two essences of x are nec. equivalent

exclusive or  $\cdot$  and for any number of summands

$M(\exists x) G(x)$  means <sup>the system of</sup> all pos-  
 sible This is true

At 4:  $P(\varphi) \cdot \varphi \supset N \psi$

~~then~~  $\begin{cases} x = x & \text{is positive} \\ x \neq x & \text{is negative} \end{cases}$

But if a system S of pos. prop. were incomp. it would mean that the sum prop. S (which is positive) would be  $x \neq x$

Positive means positive in the moral aesthetic sense (independently of the accidental structure of the world). Only when the act is pure. It also means "attribution" as opposed to "privation (or containing privation)". This is important for the proof.

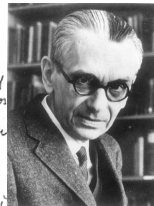
If  $\varphi$  is positive:  $(x) N \sim \varphi(x)$  (nec. true:  $\varphi(x) \supset N \sim \varphi(x)$ )

hence  $x \neq x$  positive, not  $x = x$  neg. necessary At 4

on the basis of pos. prop.

$x$  i.e. the normal form in terms of elem. prop. contains a member without negation.

# Computational Metaphysics: Kurt Gödel's Ontological Argument



Ontologischer Beweis Feb. 10, 1970  
 $P(\varphi)$  if  $\varphi$  is positive (i.e.  $\varphi \in P$ )

$M(\exists x) G(x)$  means all pos.   
 loc. b. l. This is true

## Computational Metaphysics: Dana Scott's Variant

**Axiom A1** Either a property or its negation is positive, but not both:

$$\forall \phi [P(\neg \phi) \leftrightarrow \neg P(\phi)]$$

**Axiom A2** A property necessarily implied by a positive property is positive:

$$\forall \phi \forall \psi [(P(\phi) \wedge \Box \forall x [\phi(x) \rightarrow \psi(x)]) \rightarrow P(\psi)]$$

**Thm. T1** Positive properties are possibly exemplified:

$$\forall \phi [P(\phi) \rightarrow \Diamond \exists x \phi(x)]$$

**Def. D1** A God-like being possesses all positive properties:

$$G(x) \leftrightarrow \forall \phi [P(\phi) \rightarrow \phi(x)]$$

**Axiom A3** The property of being God-like is positive:

$$P(G)$$

**Cor. C** Possibly, God exists:

$$\Diamond \exists x G(x)$$

**Axiom A4** Positive properties are necessarily positive:

$$\forall \phi [P(\phi) \rightarrow \Box P(\phi)]$$

**Def. D2** An essence of an individual is a property possessed by it and necessarily implying any of its properties:

$$\phi \text{ ess. } x \leftrightarrow \phi(x) \wedge \forall \psi (\psi(x) \rightarrow \Box \forall y (\phi(y) \rightarrow \psi(y)))$$

**Thm. T2** Being God-like is an essence of any God-like being:

$$\forall x [G(x) \rightarrow G \text{ ess. } x]$$

**Def. D3** Necessary existence of an individual is the necessary exemplification of all its essences:

$$NE(x) \leftrightarrow \forall \phi [\phi \text{ ess. } x \rightarrow \Box \exists y \phi(y)]$$

**Axiom A5** Necessary existence is a positive property:

$$P(NE)$$

**Thm. T3** Necessarily, God exists:

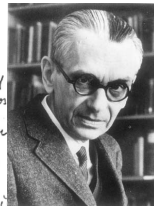
$$\Box \exists x G(x)$$



$\Diamond \neg \psi$  : is possible is negative  
 if pos. prop. were in com. at the same prop. is (which is not be  $x \neq x$ )  
 positive in the moral aesthetic  
 the accidental structure of the ax. time. It is not as opposed to "privation"  
 This is not a contradiction  
 $\neg \forall \phi(x) : \text{Characteristic } \phi(x) \supset x \neq x$   
 not  $x = x$  by necessity. At least in the  
 if elem. prop. contains a m.



# Computational Metaphysics: Kurt Gödel's Ontological Argument



Ontologischer Beweis Feb. 10, 1970  
 $P(\phi)$   $\phi$  is positive (if  $\phi \in P$ )

$M(\exists x) G(x)$  means all possible worlds  
 This is true

## Computational Metaphysics: Dana Scott's Variant

**Axiom A1** Either a property or its negation is positive, but not both:

$$\forall \phi [P(\neg \phi) \leftrightarrow \neg P(\phi)]$$

**Axiom A2** A property necessarily implied by a positive property is positive:

$$\forall \phi \forall \psi [(P(\phi) \wedge \Box \forall x [\phi(x) \rightarrow \psi(x)]) \rightarrow P(\psi)]$$

**Thm. T1** Positive properties are possibly exemplified:

$$\forall \phi [P(\phi) \rightarrow \Diamond \exists x \phi(x)]$$

**Def. D1** A God-like being possesses all positive properties:

$$G(x) \leftrightarrow \forall \phi [P(\phi) \rightarrow \phi(x)]$$

**Axiom A3** The property of being God-like is positive:

$$P(G)$$

**Cor. C** Possibly, God exists:

$$\Diamond \exists x G(x)$$

**Axiom A4** Positive properties are necessarily positive:

$$\forall \phi [P(\phi) \rightarrow \Box P(\phi)]$$

**Def. D2** An essence of an individual is a property possessed by it and necessarily implying any of its properties:

$$\phi \text{ ess. } x \leftrightarrow \phi(x) \wedge \forall \psi (\psi(x) \rightarrow \Box \forall y (\phi(y) \rightarrow \psi(y)))$$

**Thm. T2** Being God-like is an essence of any God-like being:

$$\forall x [G(x) \rightarrow G \text{ ess. } x]$$

**Def. D3** Necessary existence of an individual is the necessary exemplification of all its essences:

$$NE(x) \leftrightarrow \forall \phi [\phi \text{ ess. } x \rightarrow \Box \exists y \phi(y)]$$

**Axiom A5** Necessary existence is a positive property:

$$P(NE)$$

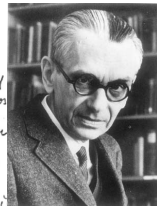
**Thm. T3** Necessarily, God exists:

$$\Box \exists x G(x)$$



$\Diamond \neg \psi$ : is possible is negative  
 if pos. prop. were incomp. at the same prop. is (which is not possible)  
 positive in the modal sense  
 the accidental structure of the ax. time. It is not as opposed to "privation"  
 This is inconsistent  
 $N \neg \phi(x) \leftrightarrow \text{essence } \phi(x) \supset x \neq x$   
 for  $x = x$  is necessary. At least in the modal sense.  
 if elem. prop. contains a

# Computational Metaphysics: Kurt Gödel's Ontological Argument



Ontologischer Beweis Feb. 10, 1970  
 $P(\varphi)$   $\varphi$  is positive (i.e.  $\varphi \in P$ )

$M(\exists x) G(x)$  means all possible worlds  
 This is true

## Computational Metaphysics: Dana Scott's Variant

**Axiom A1** Either a property or its negation is positive, but not both:

$$\forall \phi [P(\neg \phi) \leftrightarrow \neg P(\phi)]$$

**Axiom A2** A property necessarily implied by a positive property is positive:

$$\forall \phi \forall \psi [(P(\phi) \wedge \Box \forall x [\phi(x) \rightarrow \psi(x)]) \rightarrow P(\psi)]$$

- ▶ consistent
- ▶ argument valid already in logic KB
- ▶ monotheism
- ▶ modal collapse ( $\varphi \rightarrow \Box \varphi$ )—no free will

$$\Diamond \exists x \phi(x)$$

$$\rightarrow \phi(x)$$

$$P(G)$$

$$\Diamond \exists x G(x)$$

$$\rightarrow \Box P(\phi)$$

$$\rightarrow \psi(y))$$

**Thm. T2** Being God-like is an essence of any God-like being:

$$\forall x [G(x) \rightarrow G \text{ ess. } x]$$

**Def. D3** Necessary existence of an individual is the necessary exemplification of all its essences:

$$NE(x) \leftrightarrow \forall \phi [\phi \text{ ess. } x \rightarrow \Box \exists y \phi(y)]$$

**Axiom A5** Necessary existence is a positive property:

$$P(NE)$$

**Thm. T3** Necessarily, God exists:

$$\Box \exists x G(x)$$

**inconsistent**



## Computational Metaphysics: Vision of Leibniz (1646–1716) — *Calculemus!*



Quo facto, quando orientur controversiae, non magis disputatione opus erit inter duos philosophos, quam inter duos Computistas. Sufficiet enim calamos in manus sumere sedereque ad abacos, et sibi mutuo . . . dicere: calculemus. (Leibniz, 1684)

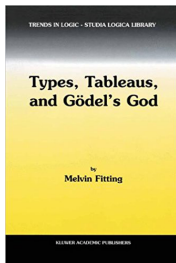


If controversies were to arise, there would be no more need of disputation between two philosophers than between two accountants. For it would suffice to take their pencils in their hands, to sit down to their slates, and to say to each other . . . : Let us calculate.

(Translation by Russell)

Required:  
**characteristica universalis** and **calculus ratiocinator**

## Further Experiments



Melvin Fitting (New York)

**Ontological Argument**  
(avoids modal collapse)

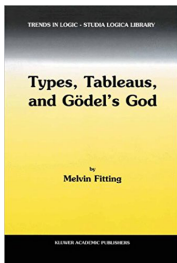
Intensional higher-order modal logic

Verified (main chapters)



David Fuenmayor  
(Philosophy, FU Berlin)

## Further Experiments



Melvin Fitting (New York)

**Ontological Argument**  
(avoids modal collapse)

Intensional higher-order modal logic

Verified (main chapters)



David Fuenmayor  
(Philosophy, FU Berlin)



Ed Zalta (Stanford)

**Principia Logico-Metaphysica**

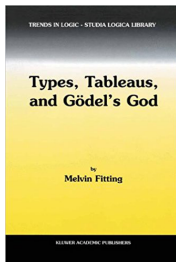
Hyperintensional higher-order modal logic

Inconsistency/Paradox detected



Daniel Kirchner  
(Mathematics, FU Berlin)

## Further Experiments



Melvin Fitting (New York)

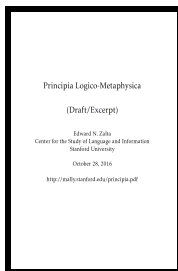
**Ontological Argument**  
(avoids modal collapse)

Intensional higher-order modal logic

Verified (main chapters)



David Fuenmayor  
(Philosophy, FU Berlin)



Ed Zalta (Stanford)

**Principia Logico-Metaphysica**

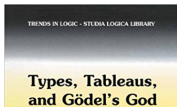
Hyperintensional higher-order modal logic

Inconsistency/Paradox detected



Daniel Kirchner  
(Mathematics, FU Berlin)

## Further Experiments



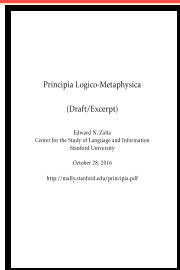
**Ontological Argument**  
(avoids modal collapse)



### Kirchner Paradox

Daniel & Isabelle/HOL are now closely collaborating with Ed Zalta

*Computational Metaphysics* par excellence!!!



Ed Zalta (Stanford)

### Principia Logico-Metaphysica

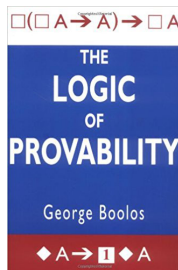
Hyperintensional higher-order modal logic

Inconsistency/Paradox detected



Daniel Kirchner  
(Mathematics, FU Berlin)

## Further Experiments

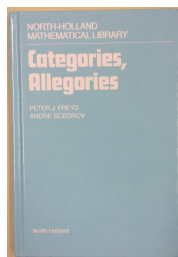


Provability Logic

Various parts verified



David Streit  
(Mathematics, FU Berlin)



### Category Theory

Free first-order logic

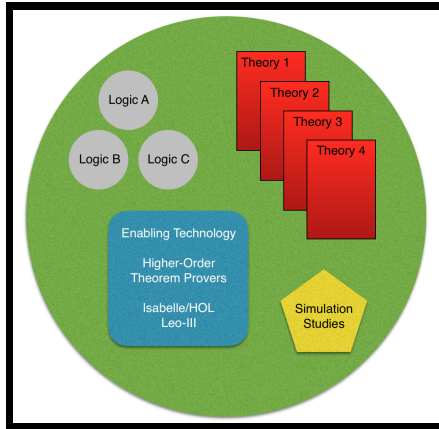
(Constricted) Inconsistency detected



D. Scott  
(UC Berkeley)

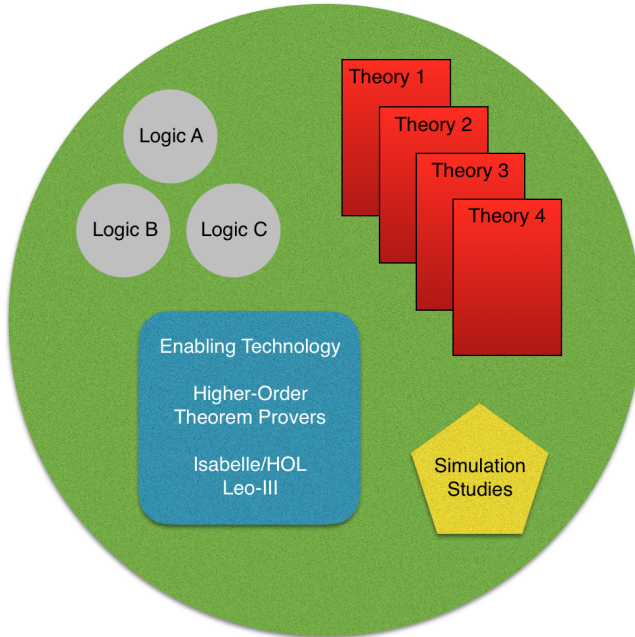
Papers on these topics: <http://christoph-benzmueller.de> → Publications





## Part D — Demo(s): Normative Reasoning Experimentation Platform

## Demo(s): Normative Reasoning Experimentation Platform



## Demo(s): Normative Reasoning Experimentation Platform

### Demo I

- ▶ Standard Deontic Logic (SDL) in Isabelle/HOL
- ▶ Dyadic Deontic Logic (DDL) in Isabelle/HOL
- ▶ Preference-based DDL in Isabelle/HOL

### Demo II

- ▶ Input/Output-Logic in Isabelle/HOL

### Demo III

- ▶ Gewirth's Principle of Generic Consistency (PGC) in Isabelle/HOL

### Demo IV

- ▶ Native Support for Deontic Logic(s) in Leo-III

## Demo(s): Normative Reasoning Experimentation Platform

### Demo I

- ▶ **Standard Deontic Logic (SDL) in Isabelle/HOL**
- ▶ **Dyadic Deontic Logic (DDL) in Isabelle/HOL**
- ▶ Preference-based DDL in Isabelle/HOL

### Demo II

- ▶ Input/Output-Logic in Isabelle/HOL

### Demo III

- ▶ **Gewirth's Principle of Generic Consistency (PGC) in Isabelle/HOL**

### Demo IV

- ▶ Native Support for Deontic Logic(s) in Leo-III

```

1 theory SDL imports Main (* Christoph Benzmüller & Xavier Parent, 2018 *)
2
3 begin (* SDL: Standard Deontic Logic (Modal Logic D) *)
4 typedecl i (*type for possible worlds*) type_synonym  $\sigma$  = "(i $\Rightarrow$ bool)"
5 consts r::"i $\Rightarrow$ i $\Rightarrow$ bool" (infixr"70") (*Accessibility relation.*) cw::i (*Current world.*)
6
7 abbreviation mtop ("T") where "T  $\equiv$   $\lambda w$ . True"
8 abbreviation mbot ("⊥") where "⊥  $\equiv$   $\lambda w$ . False"
9 abbreviation mnot ("¬") [52]53 where "¬ $\varphi$   $\equiv$   $\lambda w$ . ¬ $\varphi(w)$ "
10 abbreviation mand (infixr"∧"51) where " $\varphi \wedge \psi$   $\equiv$   $\lambda w$ .  $\varphi(w) \wedge \psi(w)$ "
11 abbreviation mor (infixr"∨"50) where " $\varphi \vee \psi$   $\equiv$   $\lambda w$ .  $\varphi(w) \vee \psi(w)$ "
12 abbreviation mimp (infixr"→"49) where " $\varphi \rightarrow \psi$   $\equiv$   $\lambda w$ .  $\varphi(w) \rightarrow \psi(w)$ "
13 abbreviation mequ (infixr"↔"48) where " $\varphi \leftrightarrow \psi$   $\equiv$   $\lambda w$ .  $\varphi(w) \leftrightarrow \psi(w)$ "
14 abbreviation obligatory ("OB") where "OB  $\varphi$   $\equiv$   $\lambda w$ .  $\forall v$ .  $w \ r \ v \rightarrow \varphi(v)$ " (*obligatory*)
15 abbreviation mpermissible ("PE") where "PE  $\varphi$   $\equiv$  ¬(OB (¬ $\varphi$ ))" (*permissible*)
16 abbreviation mimpermissible ("IM") where "IM  $\varphi$   $\equiv$  OB (¬ $\varphi$ )" (*impermissible*)
17 abbreviation omissible ("OM") where "OM  $\varphi$   $\equiv$  ¬(OB  $\varphi$ )" (*omissible*)
18 abbreviation moptional ("OP") where "OP  $\varphi$   $\equiv$  (¬(OB  $\varphi$ )  $\wedge$  ¬(OB (¬ $\varphi$ )))" (*optional*)
19
20 abbreviation ddlvalid::" $\sigma \Rightarrow$  bool" ("⊢") [7]105 (*Global Validity*)
21 where "⊢  $\equiv$   $\forall w$ . A w"
22 abbreviation ddlvalidcw::" $\sigma \Rightarrow$  bool" ("⊢cw") [7]105 (*Local Validity (in cw)*)
23 where "⊢cw  $\equiv$  A cw"
24
25 (* The D axiom is postulated *)
26 axiomatization where D: "⊢ (¬ ((OB  $\varphi$ )  $\wedge$  (OB (¬  $\varphi$ ))))"
27
28 (* Meta-level study: D corresponds to seriality *)
29 lemma "⊢ (¬ ((OB  $\varphi$ )  $\wedge$  (OB (¬  $\varphi$ ))))  $\longleftrightarrow$  ( $\forall w$ .  $\exists v$ .  $w \ r \ v$ )" by auto
30
31 (* Standardised syntax: unary operator for obligation in SDL *)
32 abbreviation obligatorySDL::" $\sigma \Rightarrow$   $\sigma$ " ("O⊢") where "O⊢ A  $\equiv$  OB A"
33
34 (* Consistency *)
35 lemma True nitpick [satisfy] oops

```

# Completeness and decidability results for a logic of contrary-to-duty conditionals

JOSÉ M. C. L. M. CARMO, *Centre of Exact Sciences and Engineering, University of Madeira, Campus Universitario da Penteada, 9020-105 Funchal, Madeira, Portugal.*

*E-mail: jcc@uma.pt*

ANDREW J. I. JONES, *Department of Informatics, King's College London Strand, London WC2R 2LS, UK.*

*E-mail: andrewji.jones@kcl.ac.uk*

## Abstract

This article has two parts. In Part I, we briefly outline the analysis of 'contrary-to-duty' obligation sentences presented in our 2002 handbook chapter 'Deontic logic and contrary-to-duties', with a focus on the intuitions that motivated the basic formal-logical moves we made. We also explain that the present account of the theory differs in two significant respects from the earlier version, one terminological, the other concerning the way the constituent modalities interconnect. Part II is the principal contribution of this article, in which we show that it is possible to define a complete and decidable axiomatization for the Carmo and Jones logic, a problem that was still open. The axiomatization includes two new inference rules; we illustrate their use in proofs, and show that on the basis of this axiomatization we can recover all the axioms and rules considered in 'Deontic logic and contrary-to-duties', and used there in the analysis of contrary-to-duty conditional scenarios.

*Keywords:* deontic logic, contrary-to-duty conditionals (CTDs), completeness and decidability results.

# Completeness and decidability results for a logic of contrary-to-duty conditionals

## 2.2 Section 2. Semantics

Our *models* are structures  $M = \langle W, av, pv, ob, V \rangle$ , where:

- (1)  $W$  is a non-empty set.
- (2)  $V$  is a function assigning a truth set to each atomic sentence (i.e.  $V(q) \subseteq W$ ).
- (3) 'av' is a function (where  $\wp(W)$  denotes the power set of  $W$ )  
 $av : W \rightarrow \wp(W)$   
 such that (where  $w$  denotes an arbitrary element of  $W$ ):  
 (3a)  $av(w) \neq \emptyset$
- (4)  $pv : W \rightarrow \wp(W)$  is such that:  
 (4a)  $av(w) \subseteq pv(w)$   
 (4b)  $w \in pv(w)$
- (5) and  $ob : \wp(W) \rightarrow \wp(\wp(W))$  is such that (where  $X, Y, Z$  designate arbitrary subsets of  $W$ )<sup>7</sup>:  
 (5a)  $\emptyset \notin ob(X)$   
 (5b) if  $Y \cap X = Z \cap X$ , then  $(Y \in ob(X) \text{ iff } Z \in ob(X))$   
 (5c\*) Let  $\beta \subseteq ob(X)$  and  $\beta \neq \emptyset$ , i.e. let  $\beta$  be a non-empty set of elements of  $ob(X)$ .  
 If  $(\cap \beta) \cap X \neq \emptyset$  (where  $\cap \beta = \{w \in W : \forall Z \in \beta \ w \in Z\}$ )  
 then  $(\cap \beta) \in ob(X)$   
 (5d) if  $Y \subseteq X$  and  $Y \in ob(X)$  and  $X \subseteq Z$ , then  $((Z - X) \cup Y) \in ob(Z)$   
 (5e) if  $Y \subseteq X$  and  $Z \in ob(X)$  and  $Y \cap Z \neq \emptyset$ , then  $Z \in ob(Y)$

ing, University  
Madeira,

ondon

ntences presented in  
motivated the basic  
ificant respects from  
onnect. Part II is the  
e axiomatization for  
rules; we illustrate  
rules considered in  
narios.

# Completeness and decidability results for a logic of contrary-to-duty conditionals

Given a model  $M = \langle W, \dots \rangle$ , the elements of  $W$  are designated by *worlds* and (as above) in what follows we will use  $w, v, \dots$  to denote arbitrary worlds and  $X, Y, Z$  to denote arbitrary sets of worlds. Intuitively:  $av(w)$  denotes the set of actual versions of the world  $w$ ;  $pv(w)$  denotes the set of potential versions of the world  $w$ ; and  $ob(X)$  denotes the set of propositions which are obligatory in context  $X$ .

We write  $M \models_w A$  to denote that formula  $A$  is true in the world  $w$  of a model  $M = \langle W, av, pv, ob, V \rangle$ , and we define  $\|A\|^M = \{w \in W : M \models_w A\}$ . In order to simplify the presentation, whenever the model  $M$  is obvious from the context, we write  $\|A\|$  instead of  $\|A\|^M$ .

*Truth in a world  $w$  in a model  $M = \langle W, av, pv, ob, V \rangle$  is characterized as follows:*

$M \models_w p$	iff	$w \in V(p)$
...		(the usual truth conditions for the connectives $\neg, \wedge, \vee, \rightarrow$ and $\leftrightarrow$ )
$M \models_w \Box A$	iff	$\ A\  = W$
$M \models_w \Box_a A$	iff	$av(w) \subseteq \ A\ $
$M \models_w \Box_p A$	iff	$pv(w) \subseteq \ A\ $
$M \models_w O(B/A)$	iff	$\ A\  \cap \ B\  \neq \emptyset$ and $(\forall X)(\text{if } X \subseteq \ A\  \text{ and } X \cap \ B\  \neq \emptyset, \text{ then } \ B\  \in ob(X))$
$M \models_w O_a A$	iff	$\ A\  \in ob(av(w))$ and $av(w) \cap \ \neg A\  \neq \emptyset$
$M \models_w O_p A$	iff	$\ A\  \in ob(pv(w))$ and $pv(w) \cap \ \neg A\  \neq \emptyset$

A sentence  $A$  is said to be *true in a model*  $M = \langle W, av, pv, ob, V \rangle$ , written  $M \models A$ , iff  $\|A\|^M = W$ ; and  $A$  is said to be *valid*, written  $\models A$ , iff  $M \models A$  in all models  $M$ .



```

DDL.thy
~/chris/trunk/tex/talks/2018-DEON/DEMO/

1 theory DDL imports Main (* Christoph Benzmüller & Xavier Parent & Ali Farjami, 2018 *)
2
3 begin (* DDL: Dyadic Deontic Logic by Carmo and Jones *)
4 typedecl i (*type for possible worlds*) type_synonym  $\sigma$  = "(i $\Rightarrow$ bool)"
5 consts av::"i $\Rightarrow$  $\sigma$ " pv::"i $\Rightarrow$  $\sigma$ " ob::" $\sigma \Rightarrow (\sigma \Rightarrow \text{bool})$ " (*accessibility relations*) cw::i (*current world*)
6
7 axiomatization where
8   ax_3a: " $\exists x. \text{av}(w)(x)$ " and ax_4a: " $\forall x. \text{av}(w)(x) \longrightarrow \text{pv}(w)(x)$ " and ax_4b: " $\text{pv}(w)(w)$ " and
9   ax_5a: " $\neg \text{ob}(X)(\lambda x. \text{False})$ " and
10  ax_5b: " $(\forall w. ((Y(w) \wedge X(w)) \longleftrightarrow (Z(w) \wedge X(w)))) \longrightarrow (\text{ob}(X)(Y) \longleftrightarrow \text{ob}(X)(Z))$ " and
11  ax_5c: " $((\forall Z. \beta(Z) \longrightarrow \text{ob}(X)(Z)) \wedge (\exists Z. \beta(Z)) \longrightarrow$ 
12     $((\exists y. ((\lambda w. \forall Z. (\beta(Z) \longrightarrow (Z(w))(y) \wedge X(y))) \longrightarrow \text{ob}(X)(\lambda w. \forall Z. (\beta(Z) \longrightarrow (Z(w))))$ " and
13  ax_5d: " $((\forall w. Y(w) \longrightarrow X(w)) \wedge \text{ob}(X)(Y) \wedge (\forall w. X(w) \longrightarrow Z(w)))$ 
14     $\longrightarrow \text{ob}(Z)(\lambda w. (Z(w) \wedge \neg X(w)) \vee Y(w))$ " and
15  ax_5e: " $((\forall w. Y(w) \longrightarrow X(w)) \wedge \text{ob}(X)(Z) \wedge (\exists w. Y(w) \wedge Z(w))) \longrightarrow \text{ob}(Y)(Z)$ "
16
17 abbreviation ddlneg ("¬"[52]53) where "¬A  $\equiv \lambda w. \neg A(w)$ "
18 abbreviation ddland ("infir"51) where "A $\wedge$ B  $\equiv \lambda w. A(w) \wedge B(w)$ "
19 abbreviation ddlor ("infir"50) where "A $\vee$ B  $\equiv \lambda w. A(w) \vee B(w)$ "
20 abbreviation ddlimp ("infir"49) where "A $\rightarrow$ B  $\equiv \lambda w. A(w) \rightarrow B(w)$ "
21 abbreviation ddlequiv ("infir"48) where "A $\leftrightarrow$ B  $\equiv \lambda w. A(w) \leftrightarrow B(w)$ "
22 abbreviation dlbox ("□") where "□A  $\equiv \lambda w. \forall v. A(v)$ " (*A = ( $\lambda w. \text{True}$ )*
23 abbreviation ddiboxa ("□a") where "□aA  $\equiv \lambda w. (\forall x. \text{av}(w)(x) \longrightarrow A(x))$ " (*in all actual worlds*)
24 abbreviation ddiboxp ("□p") where "□pA  $\equiv \lambda w. (\forall x. \text{pv}(w)(x) \longrightarrow A(x))$ " (*in all potential worlds*)
25 abbreviation ddldia ("◇") where "◇A  $\equiv \neg \square(\neg A)$ "
26 abbreviation ddldiaa ("◇a") where "◇aA  $\equiv \neg \square_a(\neg A)$ "
27 abbreviation ddldiap ("◇p") where "◇pA  $\equiv \neg \square_p(\neg A)$ "
28 abbreviation ddlo ("0[  ]"52]53) where "0[B]A  $\equiv \lambda w. \text{ob}(A)(B)$ " (*it ought to be  $\psi$ , given  $\varphi$ *)
29 abbreviation ddloa ("0a") where "0aA  $\equiv \lambda w. \text{ob}(\text{av}(w))(A) \wedge (\exists x. \text{av}(w)(x) \wedge \neg A(x))$ " (*actual obligation*)
30 abbreviation ddlop ("0p") where "0pA  $\equiv \lambda w. \text{ob}(\text{pv}(w))(A) \wedge (\exists x. \text{pv}(w)(x) \wedge \neg A(x))$ " (*primary obligation*)
31 abbreviation ddltop ("⊤") where "⊤  $\equiv \lambda w. \text{True}$ "
32 abbreviation ddlbop ("⊥") where "⊥  $\equiv \lambda w. \text{False}$ "
33
34 abbreviation ddlvalidi::" $\sigma \Rightarrow \text{bool}$ " ("⊢"[7]105) where "⊢A  $\equiv \forall w. A w$ " (*Global validity*)
35 abbreviation ddlvalidcw::" $\sigma \Rightarrow \text{bool}$ " ("⊢cw"[7]105) where "⊢cwA  $\equiv A \text{ cw}$ " (*Local validity (in cw)*)
36
37 (* A is obligatory *)
38 abbreviation obligatoryDDL::" $\sigma \Rightarrow \sigma$ " ("0[  ]" where "0[A]  $\equiv 0[A] \top$ "
39
40 (* Consistency *)
41 lemma True nitpick [satisfy]oops
42

```

GDPR.thy

GDPR.thy (~/chris/trunk/tex/talks/2018-DEON/DEMO/)

```

1 theory GDPR imports SDL                                (* Christoph Benzmüller & Xavier Parent, 2018 *)
2
3 begin (** GDPR Example **)
4   consts process_data_lawfully::σ erase_data::σ kill_boss::σ
5
6   axiomatization where
7     (* It is an obligation to process data lawfully. *)
8     A1: "[0(process_data_lawfully)]" and
9     (* Implicit: It is an obligation to keep the data if it was processed lawfully. *)
10    Implicit: "[0(process_data_lawfully → ¬erase_data)]" and
11    (* If data was not processed lawfully, then it is an obligation to erase the data. *)
12    A2: "[¬process_data_lawfully → 0(erase_data)]"
13    (* Given a situation where data is processed unlawfully. *) and
14    A3: "[¬process_data_lawfully]_cw"
15
16  (** Some Experiments **)
17  lemma True nitpick [satisfy] oops (* Consistency-check: Is there a model? *)
18  lemma False sledgehammer oops (* Inconsistency-check: Can Falsum be derived? *)
19
20  lemma "[0(erase_data)]" sledgehammer nitpick oops (* Should the data be erased? *)
21  lemma "[0(¬erase_data)]" sledgehammer nitpick oops (* Should the data be kept? *)
22  lemma "[0(kill_boss)]" sledgehammer nitpick oops (* Should the boss be killed? *)
23 end

```

Documentation

Sidekick

State

Theories

☒ Proof state
☒ Auto update

Update

Search:

100%

Sledgehammering...

Proof found...

"spass": The prover derived "False" from "A1", "A2", "A3", "D", and "Implicit", which could

"e": The prover derived "False" from "A1", "A2", "A3", "D", and "Implicit", which could be

"cvc4": Try this: by (metis A1 A2 A3 D Implicit) (68 ms)

"z3": Try this: by (metis A1 A2 A3 D Implicit) (59 ms)

☒

Output

Query

Sledgehammer

Symbols

C. Benzmüller, 2018

32

## Demo III: Gewirth's Principle of Generic Consistency (PGC) in Isabelle/HOL

**“Act in accord with the generic rights of your recipients as well as of yourself.** I shall call this the Principle of Generic Consistency (PGC), since it combines the formal consideration of consistency with the material consideration of rights to the generic features or goods of action.” **(Alan Gewirth, Reason and Morality, Chicago U Press, 1978)**

REASON  
AND  
MORALITY  
ALAN GEWIRTH

### ► Gewirth's PGC has

- stirred much controversy in moral philosophy
- been discussed as means to bound the impact of artificial general intelligence (AGI)

### ► Idea (in a nutshell):

- devise a safety mechanism of a mathematical (deductive) nature to ensure that an AGI respects human's freedom and well-being
- mechanism is based on assumption that it is able to recognize itself, as well as us humans, as agents (prospective purposive agents, PPA) which
  - (i) act voluntarily on self-chosen purposes, and
  - (ii) reason rationally

### ► Further References

- D. Beyleveld. The dialectical necessity of morality: An analysis and defense of Alan Gewirth's argument to the principle of generic consistency. U of Chicago Press, 1991
- A. Kornai. Bounding the impact of AGI. J. Experimental & Theoretical AI, 2014

**“Act in accord with the generic rights of your recipients as well as of yourself.** I shall call this the Principle of Generic Consistency (PGC), since it combines the formal consideration of consistency with the material consideration of rights to the generic features or goods of action.” (Alan Gewirth, **Reason and Morality**, Chicago U Press, 1978)

REASON  
AND  
MORALITY  
ALAN GEWIRTH

### ► Gewirth's PGC has

- stirred much controversy in moral philosophy
- been discussed as means to bound the impact of artificial general intelligence (AGI)

### ► Idea (in a nutshell):

- devise a safety mechanism of a mathematical (deductive) nature to ensure that an AGI respects human's freedom and well-being
- mechanism is based on assumption that it is able to recognize itself, as well as us humans, as agents (prospective purposive agents, PPA) which
  - (i) act voluntarily on self-chosen purposes, and
  - (ii) reason rationally

### ► Further References

- D. Beyleveld. The dialectical necessity of morality: An analysis and defense of Alan Gewirth's argument to the principle of generic consistency. U of Chicago Press, 1991
- A. Kornai. Bounding the impact of AGI. J. Experimental & Theoretical AI, 2014

**“Act in accord with the generic rights of your recipients as well as of yourself.** I shall call this the Principle of Generic Consistency (PGC), since it combines the formal consideration of consistency with the material consideration of rights to the generic features or goods of action.” (Alan Gewirth, *Reason and Morality*, Chicago U Press, 1978)

REASON  
AND  
MORALITY  
ALAN GEWIRTH

- ▶ **Gewirth's PGC has**
  - ▶ stirred much controversy in moral philosophy
  - ▶ been discussed as means to bound the impact of artificial general intelligence (AGI)
- ▶ **Idea (in a nutshell):**
  - ▶ devise a safety mechanism of a mathematical (deductive) nature to ensure that an AGI respects human's freedom and well-being
  - ▶ mechanism is based on assumption that it is able to recognize itself, as well as us humans, as agents (prospective purposive agents, PPA) which
    - (i) act voluntarily on self-chosen purposes, and
    - (ii) reason rationally
- ▶ **Further References**
  - ▶ D. Beyleveld. The dialectical necessity of morality: An analysis and defense of Alan Gewirth's argument to the principle of generic consistency. U of Chicago Press, 1991
  - ▶ A. Kornai. Bounding the impact of AGI. J. Experimental & Theoretical AI, 2014

## Demo III: Gewirth's Principle of Generic Consistency (PGC) in Isabelle/HOL

Idea is to constrain AGIs to reason in the following way

- ▶ For me, as an AGI, it is necessary to accept that:
  - (P1) I act voluntarily on purpose E (equivalent by def. to: "I am a PPA")**
  - (C2) E is good (for me)**
  - (P3) In order to achieve any purpose whatsoever by my agency, I need my freedom and well-being**
  - (C4) My freedom and well-being are necessary goods (for me)**
  - (C5) I (even if no one else) have a claim right to my freedom and well-being**
- ▶ It is necessary for all PPAs to accept that:
  - (C9) Every PPA has a necessary right to their freedom and well-being**

### Demo III: Gewirth's Principle of Generic Consistency (PGC) in Isabelle/HOL

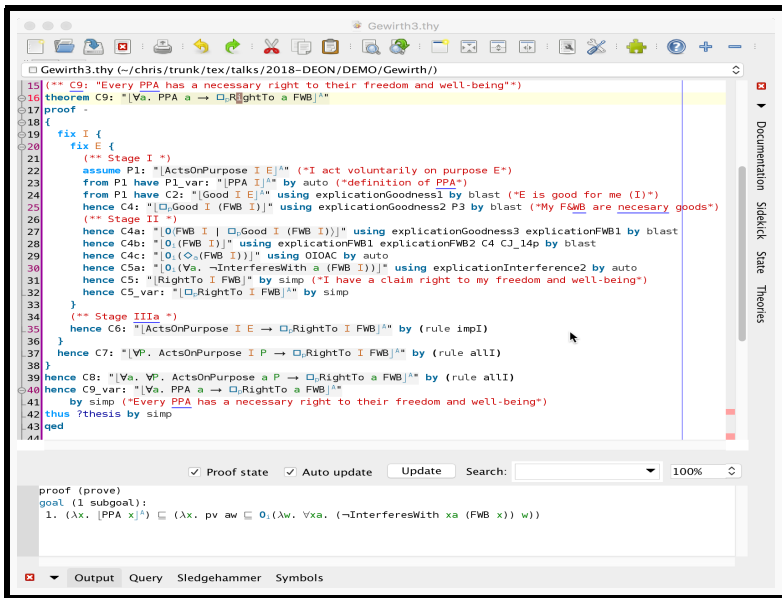
Idea is to constrain AGIs to reason in the following way

- ▶ For me, as an AGI, it is necessary to accept that:
  - (P1) I act voluntarily on purpose E (equivalent by def. to: "I am a PPA")**
  - (C2) E is good (for me)**
  - (P3) In order to achieve any purpose whatsoever by my agency, I need my freedom and well-being**
  - (C4) My freedom and well-being are necessary goods (for me)**
  - (C5) I (even if no one else) have a claim right to my freedom and well-being**
- ▶ It is necessary for all PPAs to accept that:
  - (C9) Every PPA has a necessary right to their freedom and well-being**

Any AGI (PPA) denying that it is bound by the PCG, e.g. by refusing to respect humans' well being, would deny that it is a PPA (and thus its own agency).

Hence, to avoid self-contradiction, an AGI would be bound to accord basic rights to humans.

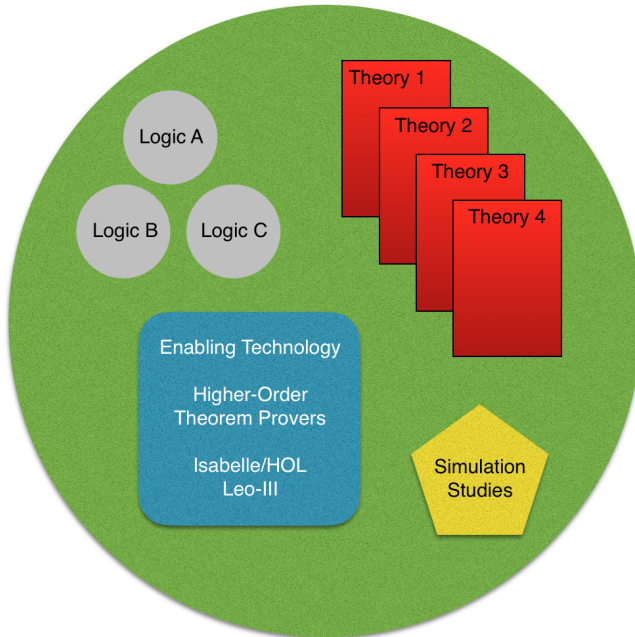
## Demo III: Gewirth's Principle of Generic Consistency (PGC) in Isabelle/HOL



By David Fuenmayor, cf. <http://christoph-benzmueller.de/papers/2018-GewirthArgument.zip>



## Demo III motivates Simulation Studies



# Ethics

## **Argued for explicit ethical reasoning competencies in IASs**

- ▶ normative reasoning experimentation platform
- ▶ HOL as universal meta-logic
- ▶ evidence from previous work
- ▶ very suitable for teaching logics

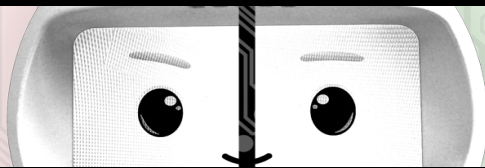
## **Ongoing and further work**

- ▶ more (deontic) logics, more logic combinations
- ▶ encoding of ethical & legal theories
- ▶ experiments, ... simulation studies, ... deployment

# Ethics

## **Argued for explicit ethical reasoning competencies in IASs**

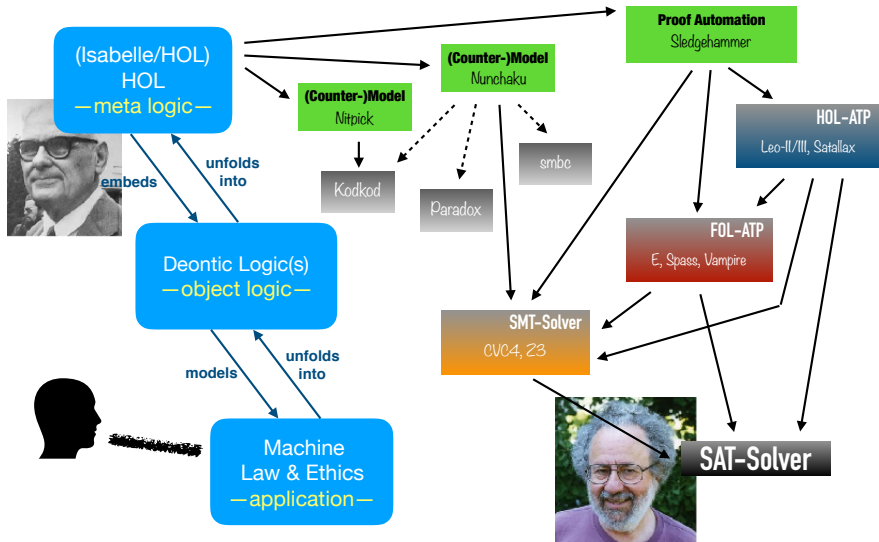
- ▶ normative reasoning experimentation platform
- ▶ HOL as universal meta-logic
- ▶ evidence from previous work
- ▶ very suitable for teaching logics



## **Ongoing and further work**

- ▶ more (deontic) logics, more logic combinations
- ▶ encoding of ethical & legal theories
- ▶ experiments, ... simulation studies, ... deployment

## How does Martin Davis fit in?



**University of Luxembourg:**

ILIAS group of Leon van der Torre, many others

**Research Grants:**

DFG, Heisenberg grant: Computational Metaphysics, BE 2501/9, **2012-2017**

DFG, Project Leo-III: Higher-Order Theorem Prover, BE 2501/9, **2013-2017**

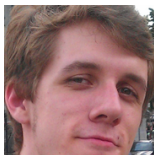
**Various Collaborators:**



B. Woltzenl.-P.  
(ANU Canberra)



Alexander Steen  
(FU Berlin)



Max Wisniewski  
(FU Berlin)



Ed Zalta  
(Stanford U.)



Dana Scott  
(UC Berkeley)

**Many further Collaborators and Students:**

Matthias Bentert (TU Berlin), Jasmin Blanchette (Amsterdam), Chad Brown (Prag), Maximilian Claus, David Fuenmayor, Tobias Gleißner, Kim Kern, Daniel Kirchner, Hanna Lachnitt, Irina Makarenko (alle FU Berlin), Larry Paulson (Cambridge), Fabian Schütz, Hans-Jörg Schurr, David Streit, Marco Ziener (alle FU Berlin), many further students in Berlin und Luxemburg

# Demo I: Global vs. Local Consequence Relation

The screenshot shows the Lean theorem prover interface with a file named `GDPRGlobal.thy` open. The file path is `~/chris/trunk/tex/talks/2018-Bath/experiments/`. The code defines a theory `GDPRGlobal` that imports `DDL` and includes a comment about the source: `(* Christoph Benzmüller & Xavier Parent, 2018 *)`. It begins with a `begin` block containing a `GDPR Example` comment. Constants are defined for `process_data_lawfully`, `erase_data`, and `kill_boss`, all of type  $\sigma$ . An `axiomatization` block follows, containing three axioms: `A1` (obligation to process data lawfully), `A3` (situation where data is processed unlawfully), and `A3` (negation of `process_data_lawfully`). A comment `(** Some Experiments **)` precedes four lemmas. The first lemma checks for consistency using `nitpick` and `nunchaku`. The second lemma checks for inconsistency using `sledgehammer` and `oops`. The next three lemmas use `sledgehammer` and `nitpick` to check specific scenarios: `erase_data`, `¬erase_data`, and `kill_boss`. The `end` keyword concludes the `begin` block. The interface includes a sidebar with `Documentation`, `Sidekick`, `State`, and `Theories`. At the bottom, the `Proof state` and `Auto update` checkboxes are checked, and the `Update` button is visible. The `Search` field is empty, and the zoom level is set to `100%`. The `Output` pane shows the results of the `sledgehammer` proof attempt, indicating that the prover derived `False` from the axioms and the goal, which could be due to a bug.

```
1 theory GDPRGlobal imports DDL (* Christoph Benzmüller & Xavier Parent, 2018 *)
2
3 begin (** GDPR Example **)
4   consts process_data_lawfully:: $\sigma$  erase_data:: $\sigma$  kill_boss:: $\sigma$ 
5
6   axiomatization where
7     (* It is an obligation to process data lawfully. *)
8     A1: "[0(process_data_lawfully)]"
9     (* Given a situation where data is processed unlawfully. *) and
10    A3: "[¬process_data_lawfully]"
11
12  (** Some Experiments **)
13  lemma True nitpick [satisfy] nunchaku [satisfy] oops (* Consistency-check: Is there a model? *)
14  lemma False sledgehammer oops (* Inconsistency-check: Can Falsum be derived? *)
15
16  lemma "[0(erase_data)]" sledgehammer nitpick oops (* Should the data be erased? *)
17  lemma "[0(¬erase_data)]" sledgehammer nitpick oops (* Should the data be kept? *)
18  lemma "[0(kill_boss)]" sledgehammer nitpick oops (* Should the boss be killed? *)
19 end
20
21
22
23
```

Proof state: ☒ Auto update: ☒ Update Search:  100%

Sledgehammering...  
Proof found..  
"cvc4": Try this: using A1 A3 ax\_5a ax\_5b by auto (11 ms)  
"z3": Try this: using A1 A3 ax\_5a ax\_5b by auto (2 ms)  
"e": Try this: using A1 A3 ax\_5a ax\_5b by auto (3 ms)  
"spass": The prover derived "False" from "A1", "A3", "ax\_5a", and "ax\_5b", which could be due to a bug

Output Query Sledgehammer Symbols

# Demo I: Preference-based DDL in Isabelle/HOL

[Journal of Philosophical Logic](#) / [Vol. 43, No. 6, December 2014](#) / Maximality vs. Optim...



## JOURNAL ARTICLE

### **Maximality vs. Optimality in Dyadic Deontic Logic: Completeness Results for Systems in Hansson's Tradition**

Xavier Parent

*Journal of Philosophical Logic*

Vol. 43, No. 6 (December 2014), pp. 1101-1128

# Demo I: Preference-based DDL in Isabelle/HOL

Journal of Philosophical Logic / Vol. 43, No. 6, December 2014 / Maximality vs. Optim...



PrefDDL.thy

```
PrefDDL.thy (~:/chris/trunk/tex/talks/2018-DEON/DEMO/)  
  
246 (* axioms of proof theory for E, check for soundness *)  
247 lemma classical: "A  $\implies$  [( $\lambda w. A$ )]" by simp -- "all classical tautologies"  
248  
249 lemma "OM": "[ $\Box A \rightarrow A$ ]" by simp -- "part of S5 schema for  $\Box$ "  
250 lemma "OS": "[ $\Diamond A \rightarrow \Box(\Diamond A)$ ]" by simp -- "part of S5 schema for  $\Diamond$ "  
251  
252 lemma DFP: "[P(B|A)  $\leftrightarrow$   $\neg(O(\neg B)|A)]$ " by (simp add: prefDDLBase.truthSet_def)  
253 lemma COK: "[O((B  $\rightarrow$  C)|A)  $\rightarrow$  (O(B|A)  $\rightarrow$  O(C|A))]" by (simp add: prefDDLBase.truthSet_def)  
254 lemma abs: "[O(B|A)  $\rightarrow$  O(O(B|A))]" by simp  
255 lemma nec: "[ $\Box A \rightarrow$  O(A|B)]" by (simp add: prefDDLBase.truthSet_def)  
256 lemma ext: "[O(A  $\leftrightarrow$  B)  $\rightarrow$  (O(C|A)  $\leftrightarrow$  O(C|B))]" by (simp add: prefDDLBase.truthSet_def)  
257 lemma id: "[O(A|A)]" by (simp add: optChoice)  
258 lemma Sh: "[O(C|(A  $\wedge$  B))  $\rightarrow$  O((B  $\rightarrow$  C)|A)]" by (smt optBest.optChoice optBest_axioms prefDDLBase.  
259  
260 (* soundness of inference rules *)  
261 lemma MP: "[A]  $\implies$  [A  $\rightarrow$  B]  $\implies$  [B]" by simp  
262 lemma N: "[A]  $\implies$  [ $\Box A$ ]" by simp  
263  
264 (* D* should hold in F, this can be verified: *)  
265 lemma "D*": "[ $\Diamond A \rightarrow$  (O(B|A)  $\rightarrow$  P(B|A))]"  
266 by (metis FOpt.opt_limitedness FOpt_axioms truthSet_def)  
267  
268 (* CM should not be provable in system F but only as of system F+CM, verified by nitpick *)  
269 lemma CM: "[O(O(B|A)  $\wedge$  O(C|A))  $\rightarrow$  O(C|(A  $\wedge$  B))]" nitpick oops  
270
```

☒ Proof state ☒ Auto update Update Search: 100%

Nitpicking formula...  
Nitpick found a counterexample for card 'w = 3':

Free variables:  
A = ( $\lambda x. \_$ )(w<sub>1</sub> := True, w<sub>2</sub> := True, w<sub>3</sub> := True)  
B = ( $\lambda x. \_$ )(w<sub>1</sub> := True, w<sub>2</sub> := True, w<sub>3</sub> := False)  
C = ( $\lambda x. \_$ )(w<sub>1</sub> := False, w<sub>2</sub> := True, w<sub>3</sub> := False)  
op  $\succeq$  =  
( $\lambda x. \_$ )  
(w<sub>1</sub> := ( $\lambda x. \_$ )(w<sub>1</sub> := True, w<sub>2</sub> := True, w<sub>3</sub> := False),  
w<sub>2</sub> := ( $\lambda x. \_$ )(w<sub>1</sub> := True, w<sub>2</sub> := True, w<sub>3</sub> := True),  
w<sub>3</sub> := ( $\lambda x. \_$ )(w<sub>1</sub> := True, w<sub>2</sub> := False, w<sub>3</sub> := True))  
opt = ( $\lambda x. \_$ )

Output Query Sledgehammer Symbols

By  
A. Steen



### Input/output (I/O) logic

[Makinson, JPL, 2000], [GabbayHortyParentEtAl-Handbook, 2013]

- ▶ I/O-operators, such as  $\text{out}_1$  (simple-minded output), accept set  $G$  of conditional norms as argument
- ▶ Conditional norms: pairs  $(a, x)$  with input “ $a$ ” (condition) and output “ $x$ ” (obligation)
- ▶ Pairs  $(a, x)$  are not given a truth-functional semantics in I/O logic

**Input/output (I/O) logic**

[Makinson, JPL, 2000], [GabbayHortyParentEtAl-Handbook, 2013]

- ▶ I/O-operators, such as  $\text{out1}$  (simple-minded output), accept set  $G$  of conditional norms as argument
- ▶ Conditional norms: pairs  $(a,x)$  with input “ $a$ ” (condition) and output “ $x$ ” (obligation)
- ▶ Pairs  $(a,x)$  are not given a truth-functional semantics in I/O logic

**Semantics of  $\text{out1}$**  (for a of input formulas  $A$ )

- ▶  $\text{out1}(G,A) := \text{Cn}(G(\text{Cn}(A)))$
- ▶ where  $\text{Cn}(X) := \{s \mid X \models s\}$  and  $G(X) := \{s \mid \exists a \in X. (a,s) \in G\}$

## Input/output (I/O) logic

[Makinson, JPL, 2000], [GabbayHortyParentEtAl-Handbook, 2013]

- ▶ I/O-operators, such as  $\text{out1}$  (simple-minded output), accept set  $G$  of conditional norms as argument
- ▶ Conditional norms: pairs  $(a,x)$  with input “ $a$ ” (condition) and output “ $x$ ” (obligation)
- ▶ Pairs  $(a,x)$  are not given a truth-functional semantics in I/O logic

Semantics of  $\text{out1}$  (for a of input formulas  $A$ )

- ▶  $\text{out1}(G,A) := \text{Cn}(G(\text{Cn}(A)))$
- ▶ where  $\text{Cn}(X) := \{s \mid X \models s\}$  and  $G(X) := \{s \mid \exists a \in X. (a,s) \in G\}$

```
(*IO logic in HOL*)
typedef i -- "type for possible worlds"      type_synonym e = "(i⇒bool)"
abbreviation ktop :: "e" ("⊤") where "⊤ ≡ λw. True"
abbreviation kbot :: "e" ("⊥") where "⊥ ≡ λw. False"
abbreviation knot :: "e⇒e" ("¬" [52] 53) where "¬φ ≡ λw. ¬φ(w)"
abbreviation kor :: "e⇒e⇒e" (infixr"∨"50) where "φ∨ψ ≡ λw. φ(w)∨ψ(w)"
abbreviation kand :: "e⇒e⇒e" (infixr"∧"51) where "φ∧ψ ≡ λw. φ(w)∧ψ(w)"
abbreviation kimp :: "e⇒e⇒e" (infixr"⊃"49) where "φ⊃ψ ≡ λw. φ(w)→ψ(w)"
abbreviation kvalid :: "e⇒bool" ("⊢" [8] 109) where "⊢p ≡ ∀w. p w"

abbreviation "outpre ≡ λG.λa.λy::e. ∃f. [a ⊃ f] ∧ G (f,y)"

abbreviation "out1 ≡ λG.λa.λx. [x] ∨
  (∃i j k. outpre G a i ∧ outpre G a j ∧ outpre G a k ∧ [(i ∧ j ∧ k) ⊃ x])"
```

IO\_Logic.thy (~/.chris/trunk/tex/talks/2018-DEON/DEMO/)

```

28 (* Some Tests *)
29 consts a::e b::e e::e
30 abbreviation "G1 ≡ (λX. X=(a,e) ∨ X=(b,e))" (* G = {(a,e),(b,e)} *)
31
32 lemma "out1 G1 a e" by blast (*proof*)
33 lemma "outpre G1 a e" by blast (*proof*)
34 lemma "outpre G1 (a ∨ b) e" nitpick oops (*countermodel*)
35 lemma "out1 G1 (a ∨ b) e" nitpick oops (*countermodel*)
36 lemma "[x] ⇒ outpre G1 (a ∨ b) x" nitpick oops (*countermodel*)
37 lemma "[x] ⇒ out1 G1 (a ∨ b) x" by blast (*proof*)
38
39
40 (* GDPR Example from before *)
41 consts pr_d_lawf::e erase_d::e kill_boss::e
42
43 abbreviation (* G = {(T,pr_d_lawf),(pr_d_lawf,¬erase_d),(¬pr_d_lawf,erase_d)} *)
44 "G ≡ (λX. X=(T,pr_d_lawf) ∨ X=(pr_d_lawf,¬erase_d) ∨ X=(¬pr_d_lawf,erase_d) )"
45
46 lemma "out1 G (¬pr_d_lawf) erase_d" by smt (*proof*)
47 lemma "out1 G (¬pr_d_lawf) (¬erase_d)" nitpick oops (*countermodel*)
48 lemma "out1 G (¬pr_d_lawf) kill_boss" nitpick oops (*countermodel*)
49 lemma "out1 G (¬pr_d_lawf) ⊥" nitpick oops (*countermodel*)

```

☒ Proof state
 ☒ Auto update
 Update Search:  100%

Nitpicking formula...

Nitpick found a counterexample for card i = 2:

Skolem constant:  
w = i<sub>1</sub>

Constants:

```

erase_d = (λx::i. _)(i1 := True, i2 := True)
kill_boss = (λx::i. _)(i1 := False, i2 := False)
pr_d_lawf = (λx::i. _)(i1 := False, i2 := True)

```

☒ Output Query Sledgehammer Symbols



$\forall$ **P**.Leo III

**Leo III** - A MASSIVELY PARALLEL HIGHER-ORDER THEOREM PROVER

### What is Leo-III?

- ▶ ATP for classical HOL (by **A. Steen**, M. Wisniewski and myself)
- ▶ ordered paramodulation; efficient data-structures; parallelisation; etc.
- ▶ native support for more than 120 logics (all normal quantified modal logics)
- ▶ including native support for **quantified SDL and DDL**
- ▶ Website: <http://page.mi.fu-berlin.de/lex/leo3/>
- ▶ Download: <https://github.com/leoprover/Leo-III>



$\forall$ P.Leo III

**Leo III** - A MASSIVELY PARALLEL HIGHER-ORDER THEOREM PROVER

### What is Leo-III?

- ▶ ATP for classical HOL (by **A. Steen**, M. Wisniewski and myself)
- ▶ ordered paramodulation; efficient data-structures; parallelisation; etc.
- ▶ native support for more than 120 logics (all normal quantified modal logics)
- ▶ including native support for **quantified SDL and DDL**
- ▶ Website: <http://page.mi.fu-berlin.de/lex/leo3/>
- ▶ Download: <https://github.com/leoprover/Leo-III>

## Brand new: Support for Dyadic Deontic Logic (Carmo/Jones)

- ▶ Enhance propositional TPTP fragment with
  1. Dyadic deontic obligation  $\$O(p/q)$
  2. Actual/primary deontic obligations  $\$O_a(p)$ ,  $\$O_p(p)$
  3. Box operators  $\$box(p)$ ,  $\$box_a(p)$ ,  $\$box_p(p)$
- ▶ Integrated into Leo-III (stand-alone tool available)



ASCII	Syntax	Meaning
~	$\neg$	Negation
	$\vee$	Disjunction
&	$\wedge$	Conjunction
=>	$\Rightarrow$	Material implication
<=>	$\Leftrightarrow$	Equivalence
$\$O(p/q)$	$O(p/q)$	Dyadic deontic obligation (It ought to be $p$ given that $q$ )
$\$box(p)$	$\Box(p)$	In all worlds $p$

Input statements: `ddl(<name>, <role>, <formula>).`

## Brand new: Support for Dyadic Deontic Logic (Carmo/Jones)

- ▶ Enhance propositional TPTP fragment with
  1. Dyadic deontic obligation  $\$O(p/q)$
  2. Actual/primary deontic obligations  $\$O_a(p)$ ,  $\$O_p(p)$
  3. Box operators  $\$box(p)$ ,  $\$box_a(p)$ ,  $\$box_p(p)$
- ▶ Integrated into Leo-III (stand-alone tool available)



ASCII	Syntax	Meaning
~	$\neg$	Negation
	$\vee$	Disjunction
&	$\wedge$	Conjunction
=>	$\Rightarrow$	Material implication
<=>	$\Leftrightarrow$	Equivalence
$\$O(p/q)$	$O(p/q)$	Dyadic deontic obligation (It ought to be $p$ given that $q$ )
$\$box(p)$	$\Box(p)$	In all worlds $p$

Input statements: `ddl(<name>, <role>, <formula>).`



Input statements: `ddl(<name>, <role>, <formula>)`.

where `<role>` provides meta-logical information:

- ▶ **axiom** *assumed, globally valid*
- ▶ **localAxiom** *assumed, valid in current world*
- ▶ **conjecture** *global consequence?*
- ▶ **localConjecture** *consequence in current world?*



### Example

This problem can directly be given to Leo-III:

```
ddl(a1, axiom, $0(processDataLawfully)).  
ddl(a2, axiom, $0(eraseData/~processDataLawfully)).  
ddl(a3, localAxiom, ~processDataLawfully).  
  
ddl(c1, conjecture, $0(eraseData)).
```

... giving ...

```
% SZS status Theorem for gdpr_new.p : 2143 ms resp. 776 ms w/o parsing
```

Input statements: `ddl(<name>, <role>, <formula>)`.

where `<role>` provides meta-logical information:

- ▶ `axiom` *assumed, globally valid*
- ▶ `localAxiom` *assumed, valid in current world*
- ▶ `conjecture` *global consequence?*
- ▶ `localConjecture` *consequence in current world?*



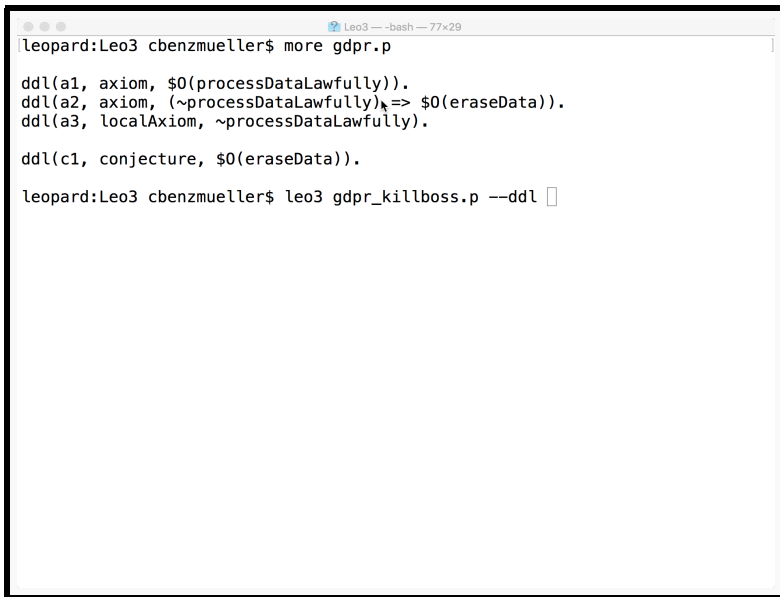
### Example

This problem can directly be given to Leo-III:

```
ddl(a1, axiom, $0(processDataLawfully)).  
ddl(a2, axiom, $0(eraseData/~processDataLawfully)).  
ddl(a3, localAxiom, ~processDataLawfully).  
  
ddl(c1, conjecture, $0(eraseData)).
```

... giving ...

```
% SZS status Theorem for gdpr_new.p : 2143 ms resp. 776 ms w/o parsing
```



```
leopard:Leo3 cbenzmueller$ more gdpr.p

ddl(a1, axiom, $0(processDataLawfully)).
ddl(a2, axiom, (~processDataLawfully) => $0(eraseData)).
ddl(a3, localAxiom, ~processDataLawfully).

ddl(c1, conjecture, $0(eraseData)).

leopard:Leo3 cbenzmueller$ leo3 gdpr_killboss.p --ddl
```

## Deontic Logic Reasoning Infrastructure

- ▶ A Dyadic Deontic Logic in HOL, DEON 2018, 2018. (See also <https://arxiv.org/abs/1802.08454>)
- ▶ A Deontic Logic Reasoning Infrastructure, CiE 2018, Springer LNCS, 2018.
- ▶ I/O Logic in HOL — First Steps, CoRR, 2018. <https://arxiv.org/abs/1803.09681>
- ▶ First Experiments with a Flexible Infrastructure for Normative Reasoning, CoRR, 2018. <http://arxiv.org/abs/1804.02929>

## Computational Metaphysics (selected)

- ▶ Experiments in Computational Metaphysics: Gödel's Proof of God's Existence, Savijnanam: scientific exploration for a spiritual paradigm. Journal of the Bhaktivedanta Institute, volume 9, pp. 43-57, 2017.
- ▶ The Inconsistency in Gödel's Ontological Argument: A Success Story for AI in Metaphysics, IJCAI 2016, 2016.
- ▶ Automating Gödel's Ontological Proof of God's Existence with Higher-order Automated Theorem Provers, ECAI 2014, IOS Press, 2014.

## Other (selected)

- ▶ The Higher-Order Prover Leo-III, IJCAR 2018, Springer LNCS, 2018.