Standard Deontic Logic (SDL = D) Isn't Going to Cut It!

(Chisholm's Paradox; The Free Choice Permission Paradox)

Selmer Bringsjord

Rensselaer AI & Reasoning (RAIR) Lab
Department of Cognitive Science
Department of Computer Science
Lally School of Management & Technology
Rensselaer Polytechnic Institute (RPI)
Troy, New York 12180 USA

Intro to Logic 4/6/2020





Peek ahead to next time for some context today ...

"We're in very deep trouble."

"We're in very deep trouble."







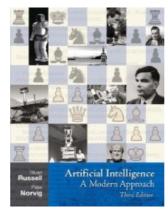


"We're in very deep trouble."











The PAID Problem!

 $\forall x : Agents$

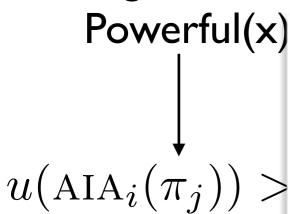
```
\forall \mathtt{x} : \mathtt{Agents}
```

Powerful(x) + Autonomous(x) + Intelligent(x) = Dangerous(x)

```
\forall x : Agents
Powerful(x) + Autonomous(x) + Intelligent(x) = Dangerous(x)
```

 $\begin{array}{c} \forall \mathtt{x} : \mathtt{Agents} \\ \mathsf{Powerful}(\mathtt{x}) + \mathsf{Autonomous}(\mathtt{x}) + \mathsf{Intelligent}(\mathtt{x}) = \mathsf{Dangerous}(\mathtt{x}) \\ \downarrow \\ u(\mathtt{AIA}_i(\pi_j)) > \tau^+ \in \mathbb{Z} \ \mathrm{or} \ \tau^- \in \mathbb{Z} \end{array}$

 $\forall x : Agents$



Are Autonomous-and-Creative Machines Intrinsically Untrustworthy?*

Selmer Bringsjord • Naveen Sundar G.

Rensselaer AI & Reasoning (RAIR) Lab Department of Cognitive Science Department of Computer Science Rensselaer Polytechnic Institute (RPI) Troy NY 12180 USA

020217NY

Abstract

Given what we find in the case of human cognition, the following principle appears to be quite plausible: An artificial agent that is both autonomous (A) and creative (C) will tend to be, from the viewpoint of a rational, fully informed agent, (U) untrustworthy. After briefly explaining the intuitive, internal structure of this disturbing principle, in the context of the human sphere, we provide a more formal rendition of it designed to apply to the realm of intelligent artificial agents. The more-formal version makes use of some of the basic structures available in one of our cognitive-event calculi, and can be expressed as a (confessedly — for reasons explained naïve) theorem. We prove the theorem, and provide simple demonstrations of it in action, using a novel theorem prover (ShadowProver). We then end by pointing toward some future defensive engineering measures that should be taken in light of the theorem.

Contents

1	Introduction	1
2	The Distressing Principle, Intuitively Put	1
	The Distressing Principle, More Formally Put 3.1 The Ideal-Observer Point of View. 3.2 Theory-of-Mind-Creativity 3.3 Autonomy 3.4 The Deontic Cognitive Event Calculus (D ⁰ CEC) 3.5 Collaborative Situations; Untrustworthiness 3.6 Theorem ACU	3
	Computational Simulations 4.1 ShadowProver	8 9 10
Re	Stronges	16

Dangerous(x)

 $\begin{array}{c} \forall \mathtt{x} : \mathtt{Agents} \\ \mathsf{Powerful}(\mathtt{x}) + \mathsf{Autonomous}(\mathtt{x}) + \mathsf{Intelligent}(\mathtt{x}) = \mathsf{Dangerous}(\mathtt{x}) \\ \downarrow \\ u(\mathtt{AIA}_i(\pi_j)) > \tau^+ \in \mathbb{Z} \ \mathrm{or} \ \tau^- \in \mathbb{Z} \end{array}$

 $\forall \mathtt{x} : \mathtt{Agents}$

Powerful(x) + Autonomous(x) + Intelligent(x) = Dangerous(x)

$$u(\operatorname{AIA}_i(\pi_j)) > \tau^+ \in \mathbb{Z} \text{ or } \tau^- \in \mathbb{Z}$$

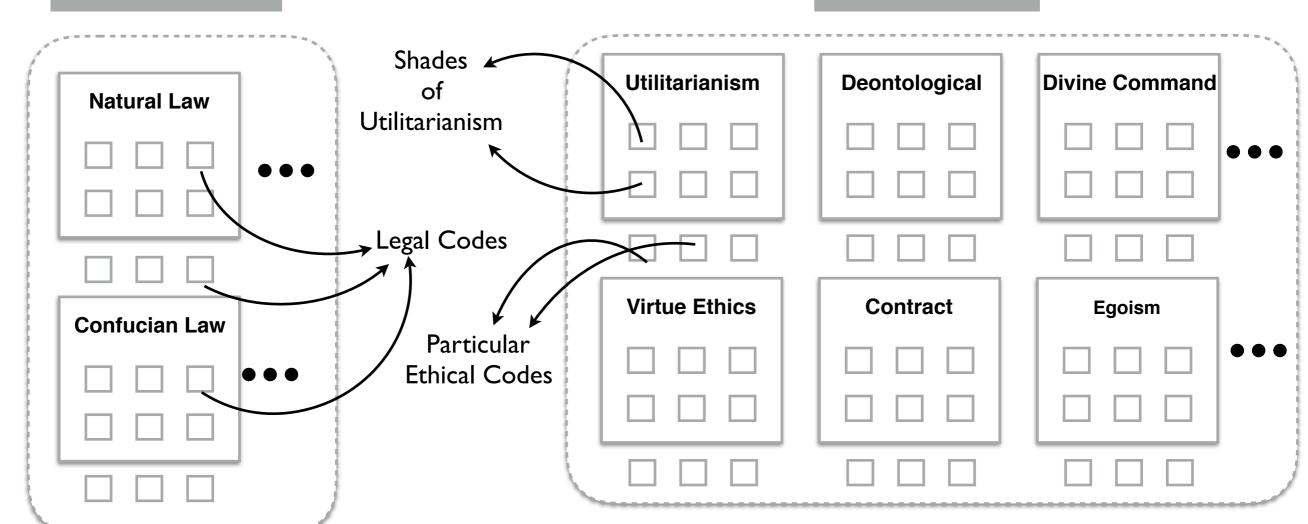
Theorem ACU: In a collaborative situation involving agents a (as the "trustor") and a' (as the "trustee"), if a' is at once both autonomous and ToM-creative, a' is untrustworthy from an ideal-observer o's viewpoint, with respect to the action-goal pair $\langle \alpha, \gamma \rangle$ in question.

Proof: Let a and a' be agents satisfying the hypothesis of the theorem in an arbitrary collaborative situation. Then, by definition, $a \neq a'$ desires to obtain some goal γ in part by way of a contributed action α_k from a', a' knows this, and moreover a' knows that a believes that this contribution will succeed. Since a' is by supposition ToM-creative, a' may desire to surprise a with respect to a's belief regarding a''s contribution; and because a' is autonomous, attempts to ascertain whether such surprise will come to pass are fruitless since what will happen is locked inaccessibly in the oracle that decides the case. Hence it follows by TRANS that an ideal observer a' will regard a' to be untrustworthy with respect to the pair a' pair. **QED**



Theories of Law

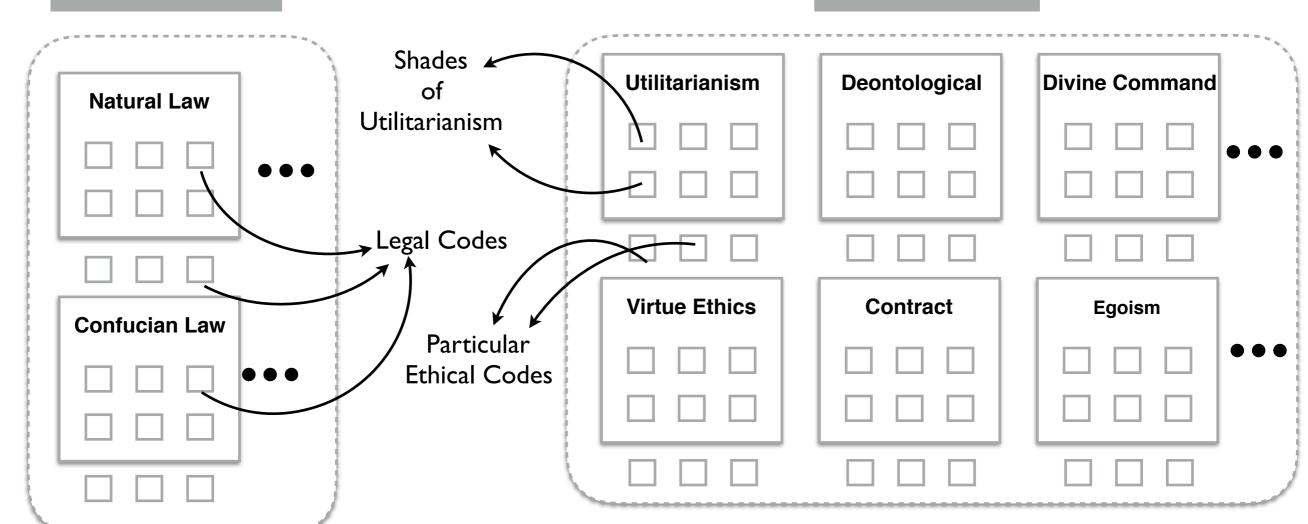
Ethical Theories





Theories of Law

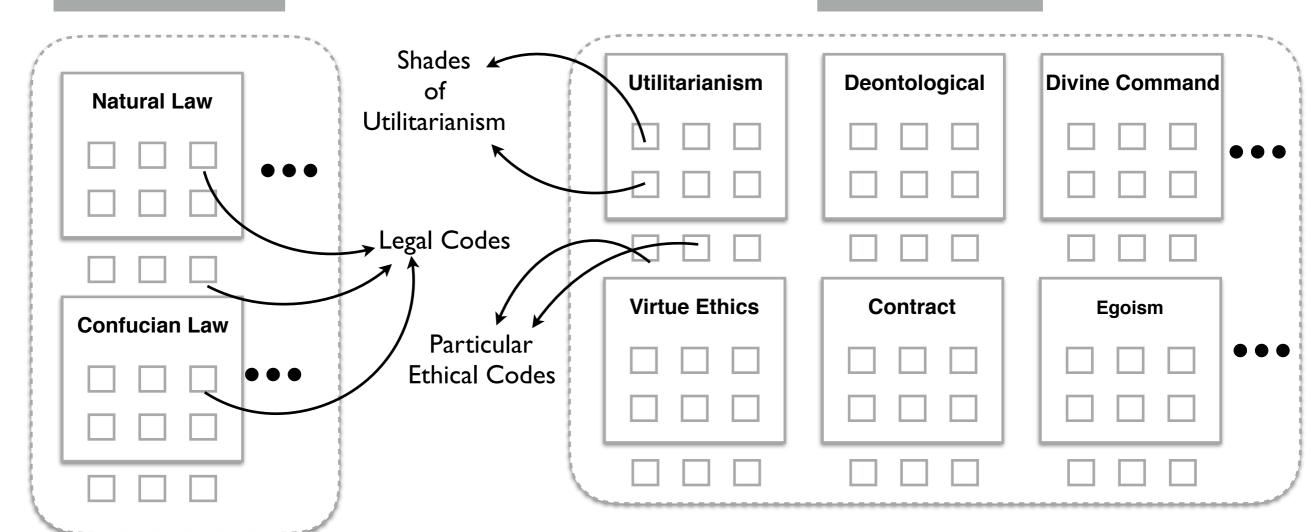
Ethical Theories





Theories of Law

Ethical Theories





Theories of Law **Ethical Theories** Shades * **Utilitarianism Deontological Divine Command** of **Natural Law** Utilitarianism Legal Codes **Virtue Ethics Contract Egoism Confucian Law** Particular **Ethical Codes**

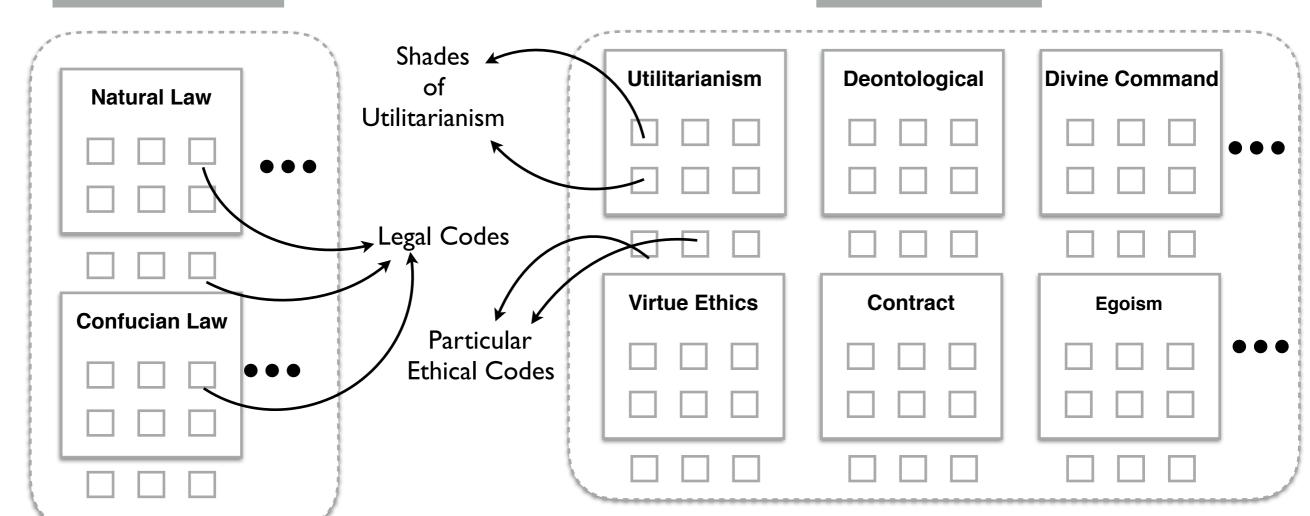
Step I

- I. Pick a theory
- 2. Pick a code
- 3. Run through EH.



Theories of Law

Ethical Theories

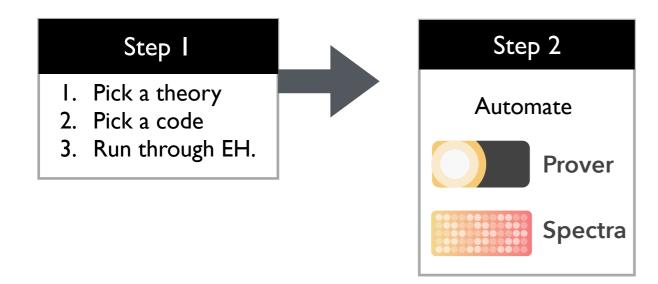


Step I

- I. Pick a theory
- 2. Pick a code
- 3. Run through EH.

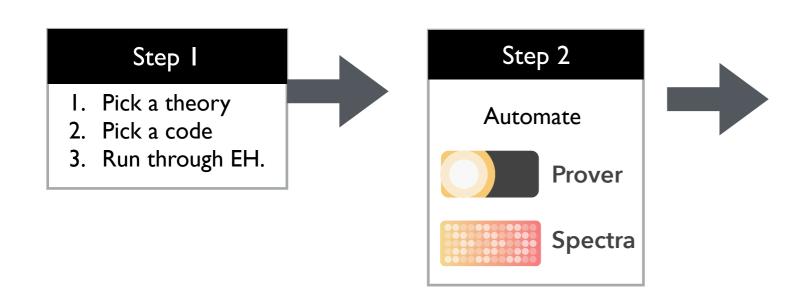


Theories of Law **Ethical Theories** Shades * **Utilitarianism Deontological Divine Command** of **Natural Law** Utilitarianism Legal Codes **Virtue Ethics** Contract **Egoism Confucian Law Particular Ethical Codes**

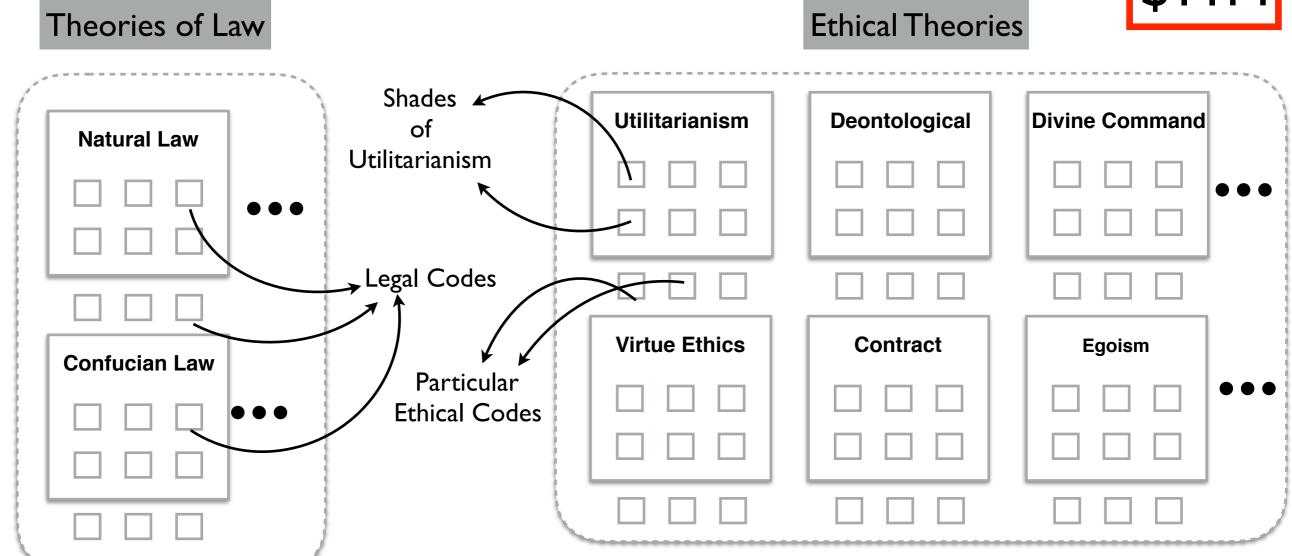


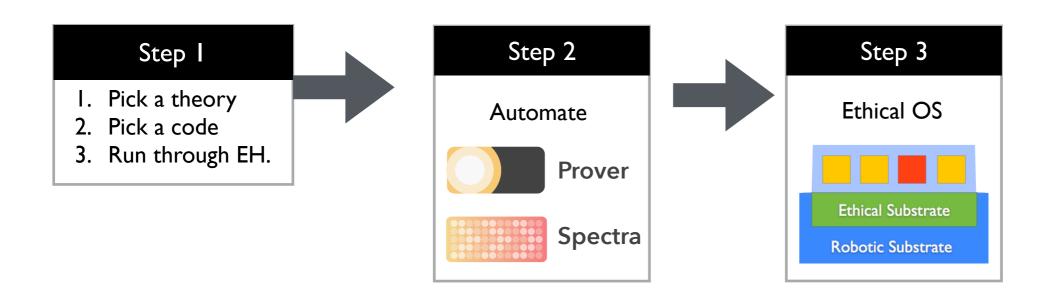


Theories of Law **Ethical Theories** Shades * **Utilitarianism Deontological Divine Command** of **Natural Law** Utilitarianism Legal Codes **Virtue Ethics** Contract **Egoism Confucian Law Particular Ethical Codes**

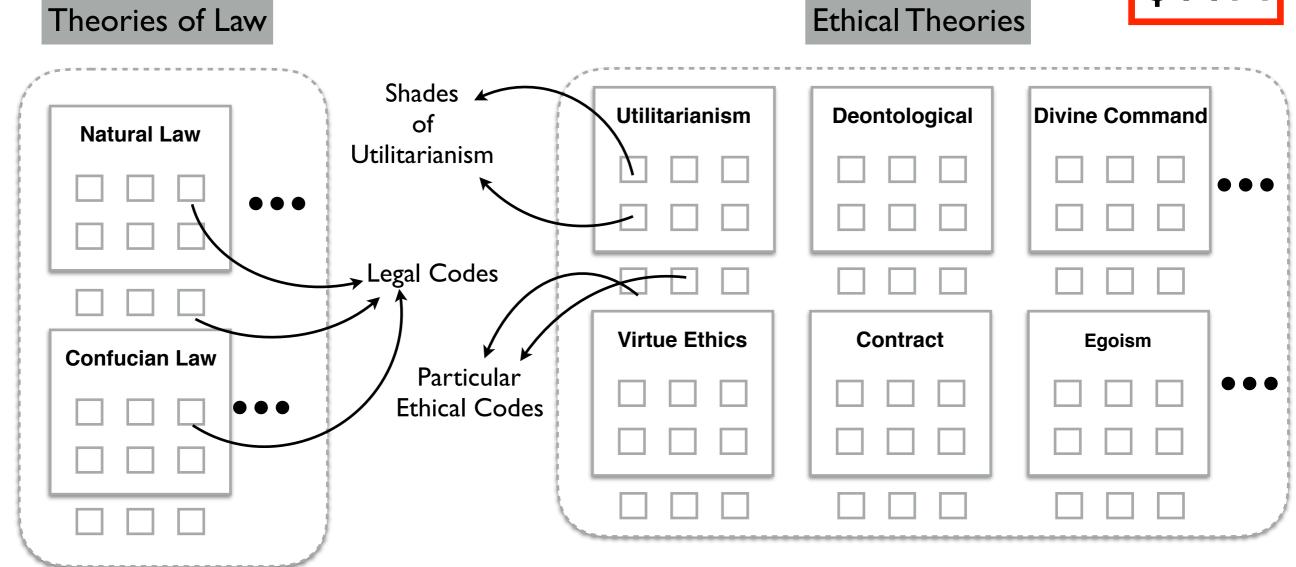


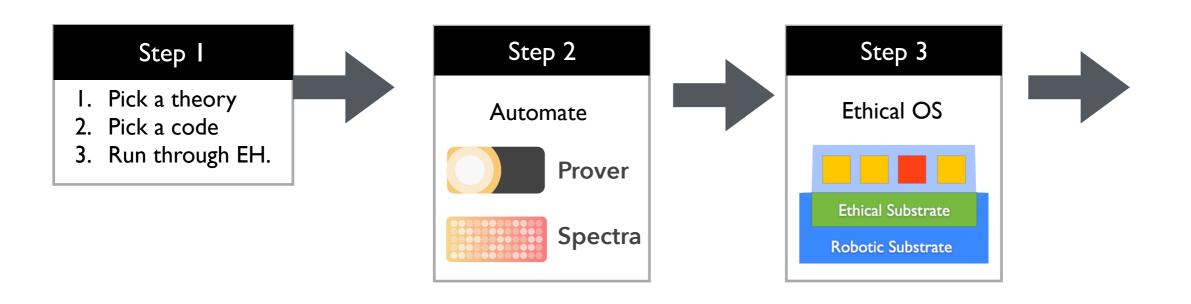




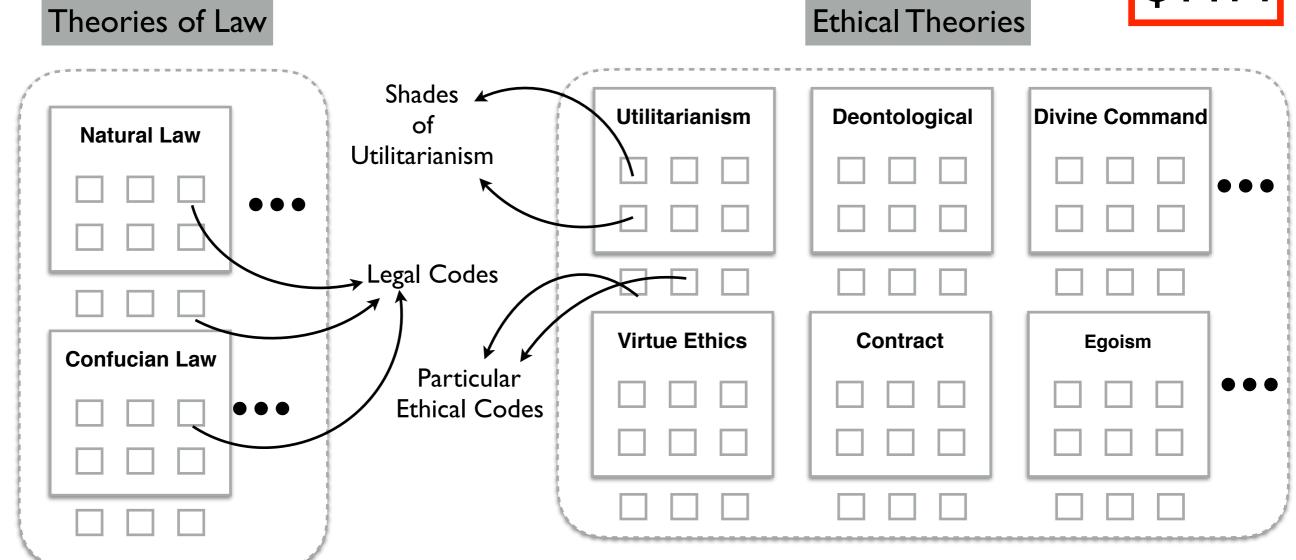


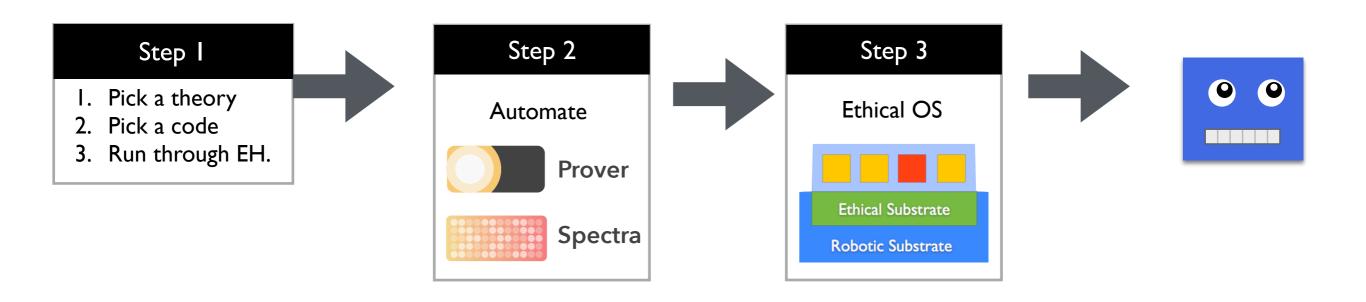




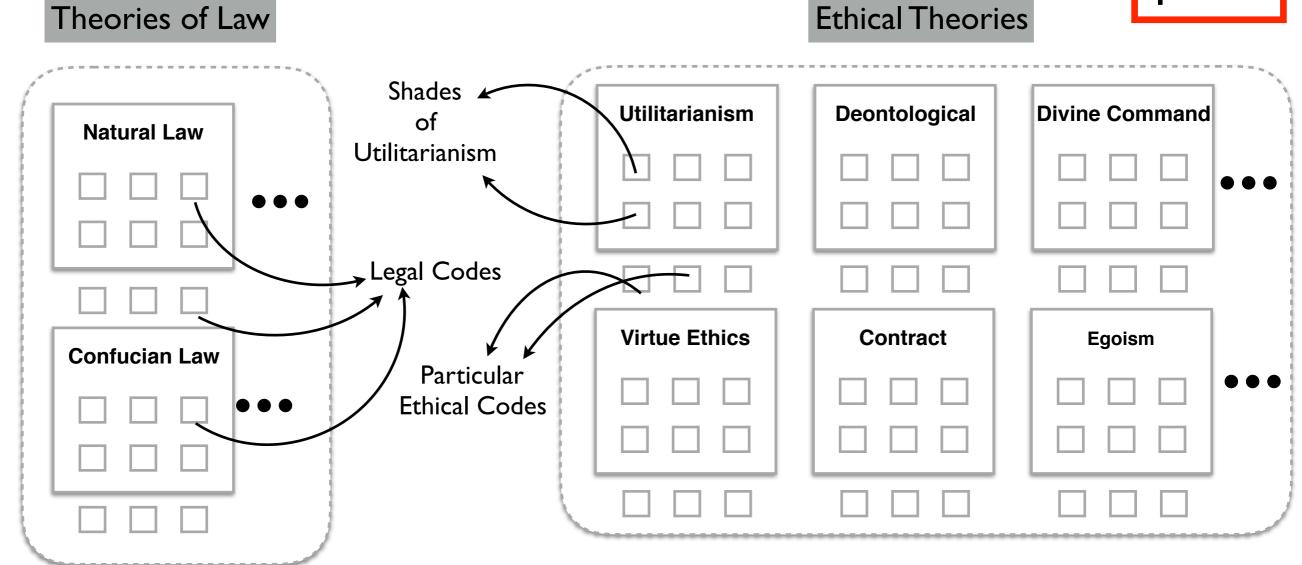


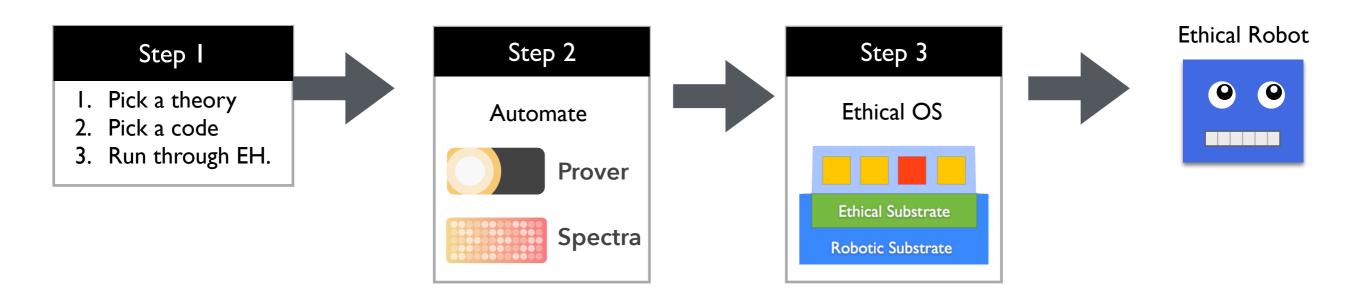




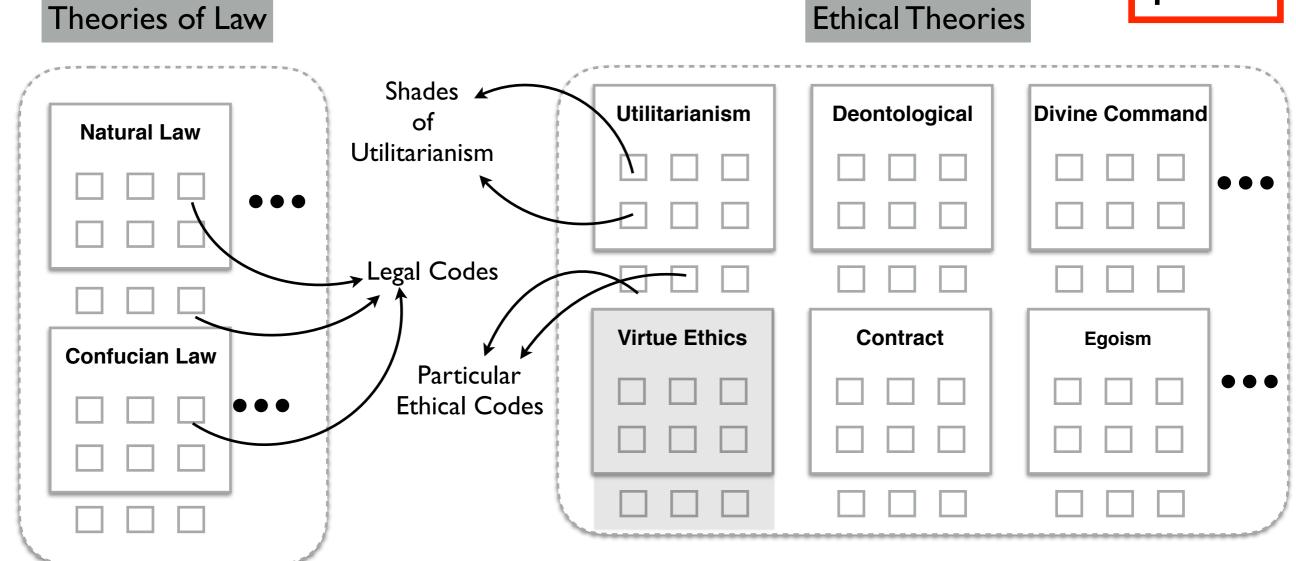


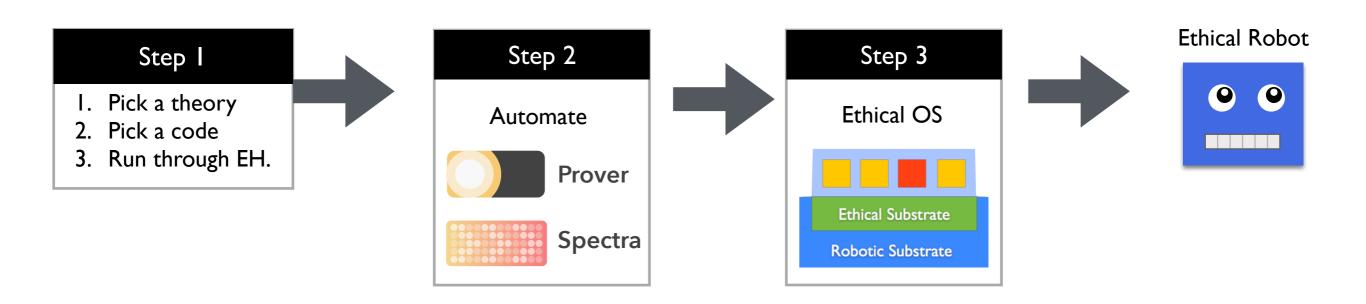


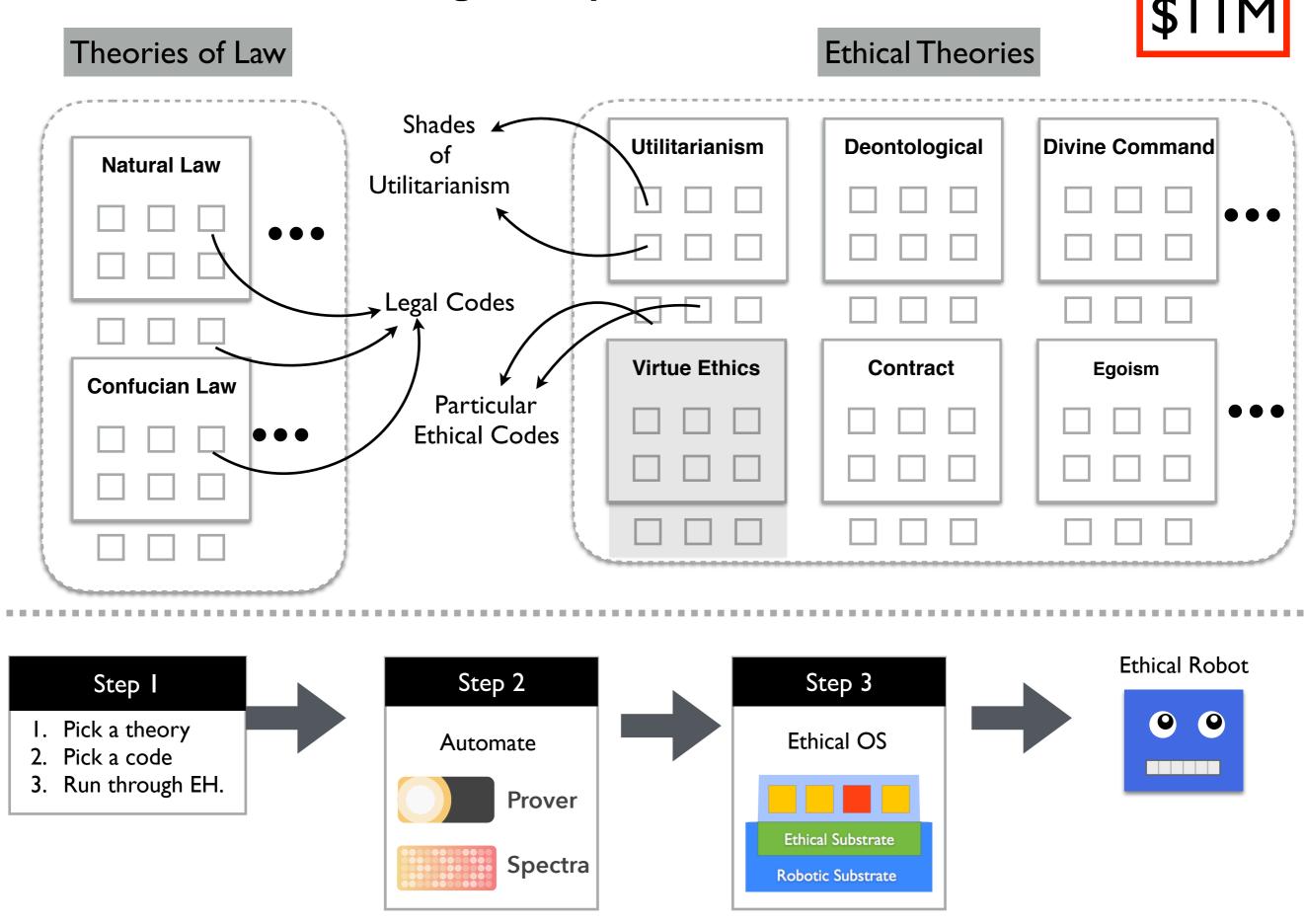










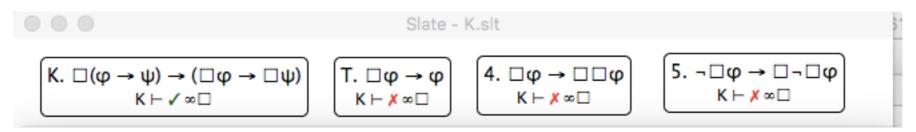


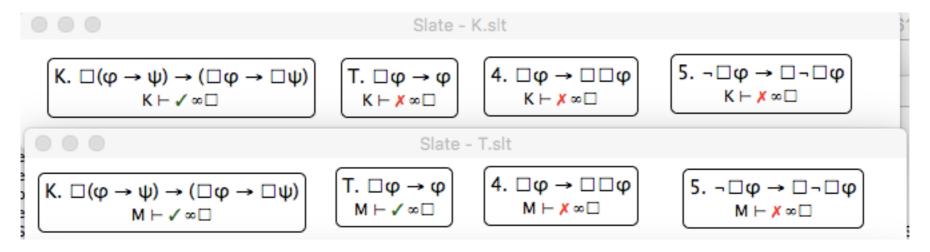
"Toward the Engineering of Virtuous Robots" Naveen, Selmer et al.

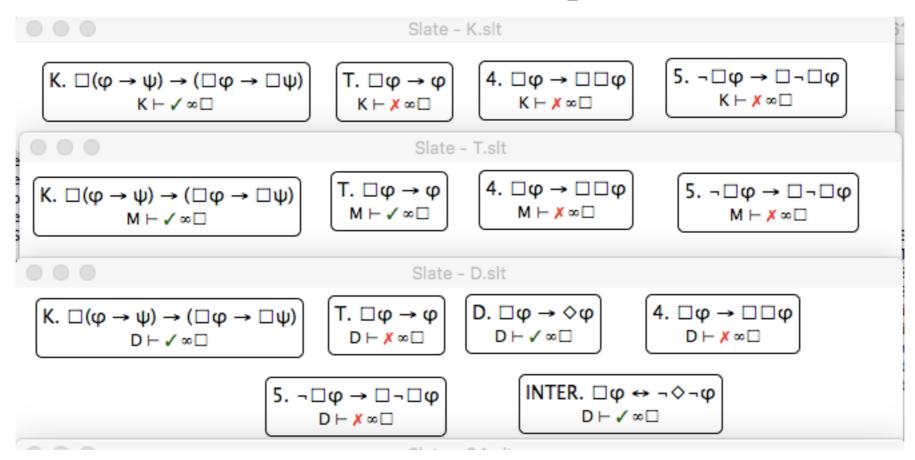
Well, maybe, but at any rate, what logic??

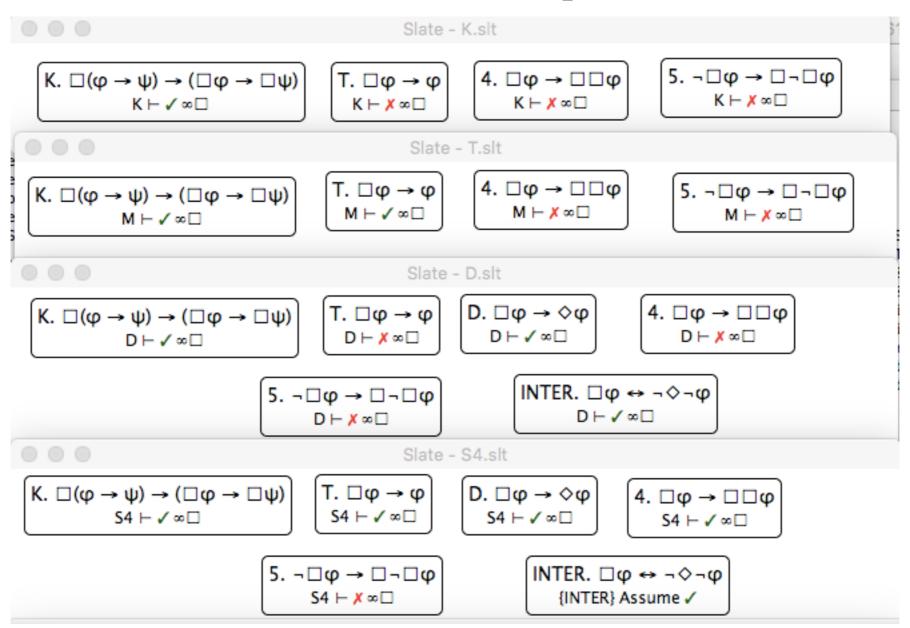
Well, maybe, but at any rate, what logic??

Perhaps **D** = **SDL**? ...



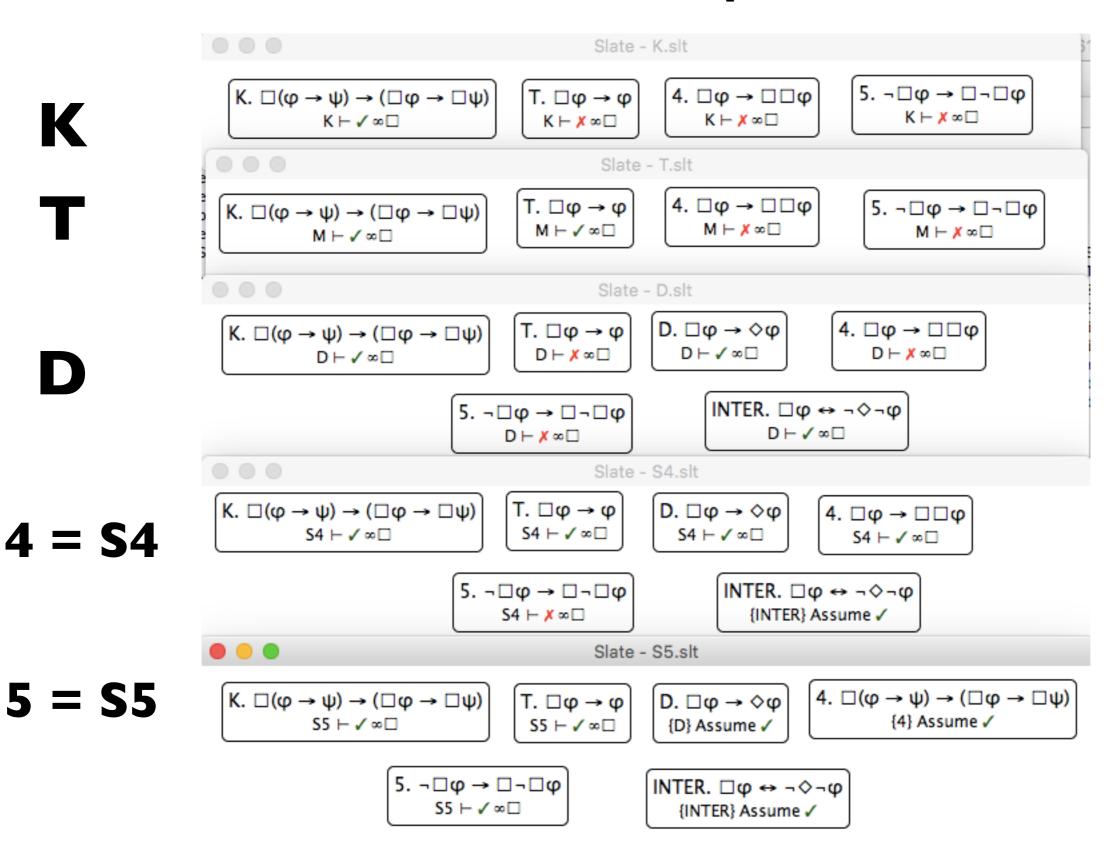




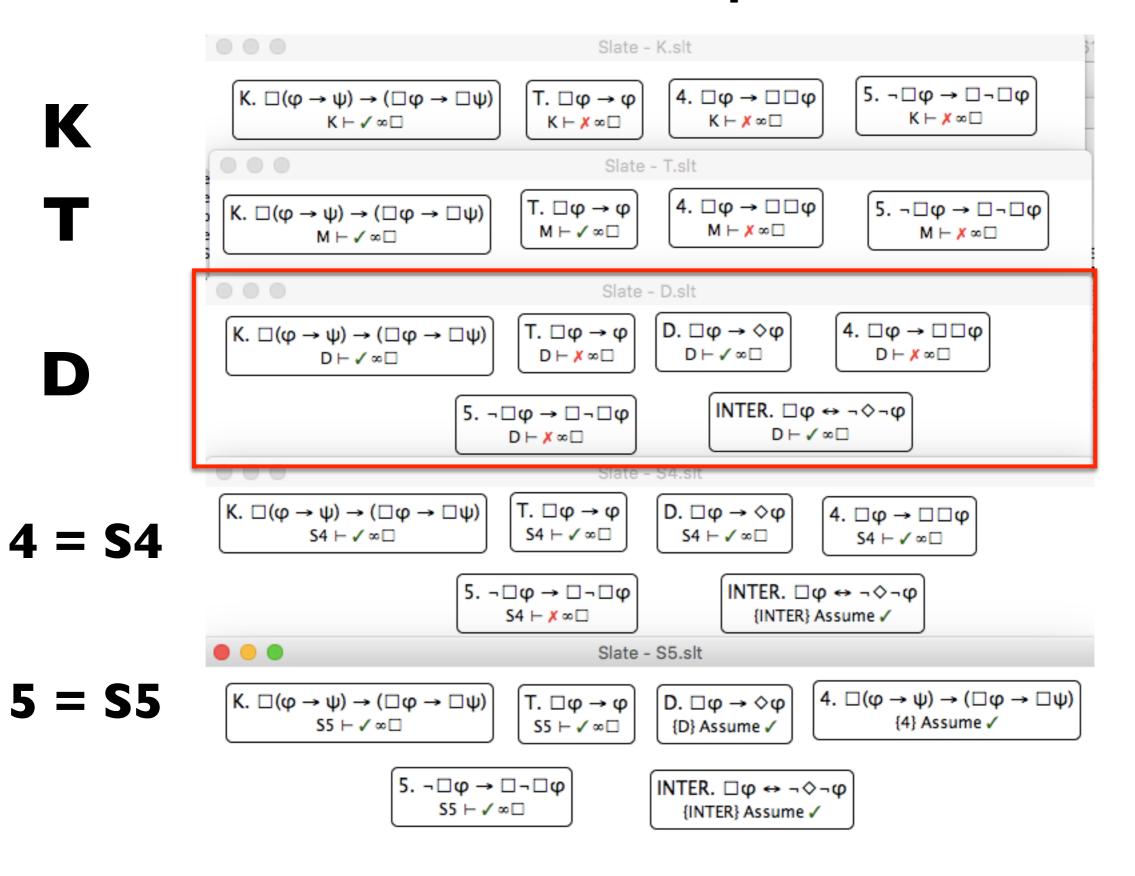


Review: Encapsulation

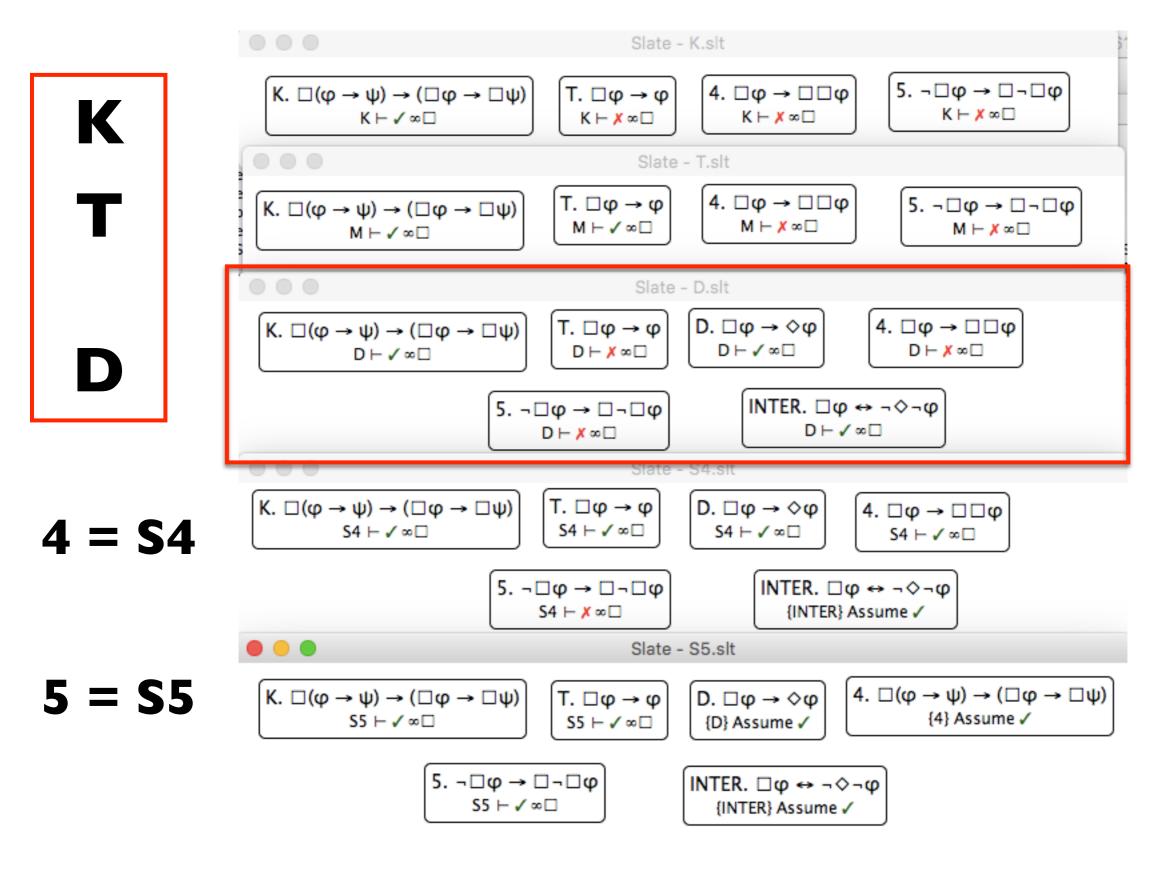
K



Review: Encapsulation



Review: Encapsulation



4.4.4 D = SDL (= 'Standard Deontic Logic')

We here introduce what is known as 'Standard Deontic Logic' (**SDL**), which in Slate is the system **D**. Deontic logic is the sub-branch of logic devoted to formalizing the fundamental concepts of morality; for example, the concepts of *obligation*, *permissibility*, and *forbiddenness*. The first of these three concepts can apparently serve as a cornerstone, since to say that ϕ (a formulae representing some state-of-affairs) is permissible seems to amount to saying that it's not obligatory that it not be the case that ϕ (which shows permissibility can be defined in terms of obligation), and to say that ϕ is forbidden would seem to amount to it being obligatory that it not be the case that ϕ (which of course appears to show that forbiddenness buildable from obligation). This interconnected trio of ethical concepts is a triad explicitly invoked and analyzed since the end of the 18^{th} century, and the importance of the triad even to modern deontic logic would be quite hard to exaggerate.

SDL is traditionally axiomatized by the following:10

SDL

TAUT All theorems of the propositional calculus.

OB-K
$$\odot(\phi \rightarrow \psi) \rightarrow (\odot \phi \rightarrow \odot \psi)$$

OB-D
$$\odot \phi \rightarrow \neg \odot \neg \phi$$

MP If
$$\vdash \phi$$
 and $\vdash \phi \rightarrow \psi$, then $\vdash \phi$

OB-NEC If
$$\vdash \phi$$
 then $\vdash \odot \phi$

4.4.4 D = SDL (= 'Standard Deontic Logic')

We here introduce what is known as 'Standard Deontic Logic' (**SDL**), which in Slate is the system **D**. Deontic logic is the sub-branch of logic devoted to formalizing the fundamental concepts of morality; for example, the concepts of *obligation*, *permissibility*, and *forbiddenness*. The first of these three concepts can apparently serve as a cornerstone, since to say that ϕ (a formulae representing some state-of-affairs) is

permissible seems to amount to saying that i that ϕ (which shows permissibility can be c say that ϕ is forbidden would seem to amouthe case that ϕ (which of course appears to sobligation). This interconnected trio of ethic and analyzed since the end of the 18^{th} centure to modern deontic logic would be quite hard SDL is traditionally axiomatized by the fo

SDL

TAUT All theorems of the proposition:

OB-K
$$\odot(\phi \rightarrow \psi) \rightarrow (\odot\phi \rightarrow \odot\psi)$$

OB-D
$$\odot \phi \rightarrow \neg \odot \neg \phi$$

MP If
$$\vdash \phi$$
 and $\vdash \phi \rightarrow \psi$, then $\vdash \phi$

OB-NEC If
$$\vdash \phi$$
 then $\vdash \odot \phi$

CHAPTER 4. PROPOSITIONAL MODAL LOGIC

OB-RE If $\vdash \phi \longleftrightarrow \psi$, then $\vdash \odot \phi \longleftrightarrow \odot \psi$.

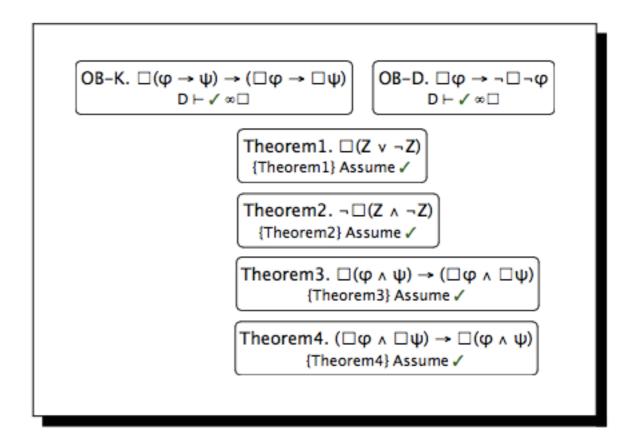


Figure 4.7: The Initial Configuration Upon Opening the File SDL.slt

4.4.4.1 Chisholm's Paradox and SDL

There are a host of problems that, together, constitute what is probably a fatal threat to **SDL** as a model of human-level ethical reasoning. We discuss in the present section the first of these problems to hit the "airwaves": Chisholm's Paradox (CP) (Chisholm 1963). CP can be generated in Slate, you we shall see. But before we get to the level of experimentation in Slate, let's understand the scenario that Chisholm's imagined.

Chisholm's clever scenario revolves around the character Jones. 11 It's given that Jones is obligated to go to assist his neighbors, in part because he has promised to do so. The second given fact is that it's obligatory that, if Jones goes to assist his neighbors, he tells them (in advance) that he is coming. In addition, and this is the third given, if Jones doesn't go to assist his neighbors, it's obligatory that he not tell

CHAPTER 4. PROPOSITIONAL MODAL LOGIC

124

them that he is coming. The fourth and final given fact is simply that Jones doesn't go to assist his neighbors. (On the way to do so, suppose he comes upon a serious vehicular accident, is proficient in emergency medicine, and (commendably!) seizes the opportunity to save the life (and subsequently monitor) of one of the victims in this accident.) These four givens have been represented in an obvious way within four formula nodes in a Slate file; see Figure 4.8. (Notice that \square is used in place of \odot .) The paradox arises from the fact that Chisholm's quartet of givens, which surely reflect situations that are common in everyday life, in conjunction with the axioms of SDL, entail outright contradictions (see Exercise 2 for D = SDL, in §4.4.4.2).

¹¹We change some particulars to ease exposition; generally, again, follow, the SEP entry on deontic logic (recall footnote 10). The core logic mirrors (Chisholm 1963), the original publication.

4.4.4.1 Chisholm's Paradox and SDL

There are a host of problems that, together, constitute what is probably a fatal threat to **SDL** as a model of human-level ethical reasoning. We discuss in the present section the first of these problems to hit the "airwaves": Chisholm's Paradox (CP) (Chisholm 1963). CP can be generated in Slate, you we shall see. But before we get to the level of experimentation in Slate, let's understand the scenario that Chisholm's imagined.

Chisholm's clever scenario revolves around the character Jones. 11 It's given that Jones is obligated to go to assist his neighbors, in part because he has promised to do so. The second given fact is that it's obligatory that, if Jones goes to assist his neighbors, he tells them (in advance) that he is coming. In addition, and this is the third given, if Jones doesn't go to assist his neighbors, it's obligatory that he not tell

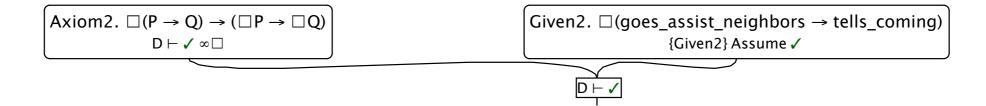
CHAPTER 4. PROPOSITIONAL MODAL LOGIC

124

them that he is coming. The fourth and final given fact is simply that Jones doesn't go to assist his neighbors. (On the way to do so, suppose he comes upon a serious vehicular accident, is proficient in emergency medicine, and (commendably!) seizes the opportunity to save the life (and subsequently monitor) of one of the victims in this accident.) These four givens have been represented in an obvious way within

four formula nodes in a Slate file; see Figure 4.8. (Notice that \square is used in place of \odot .) The paradox arises from the fact that Chisholm's quartet of givens, which surely reflect situations that are common in everyday life, in conjunction with the axioms of **SDL**, entail outright contradictions (see Exercise 2 for $\mathbf{D} = \mathbf{SDL}$, in §4.4.4.2).

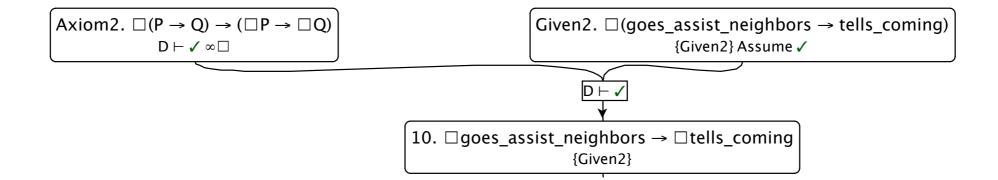
¹¹We change some particulars to ease exposition; generally, again, follow, the SEP entry on deontic logic (recall footnote 10). The core logic mirrors (Chisholm 1963), the original publication.



Axiom4. "Modus ponens for provability." {Axiom4} Assume ✓

Axiom5. "Theorems are obligatory." {Axiom5} Assume ✓

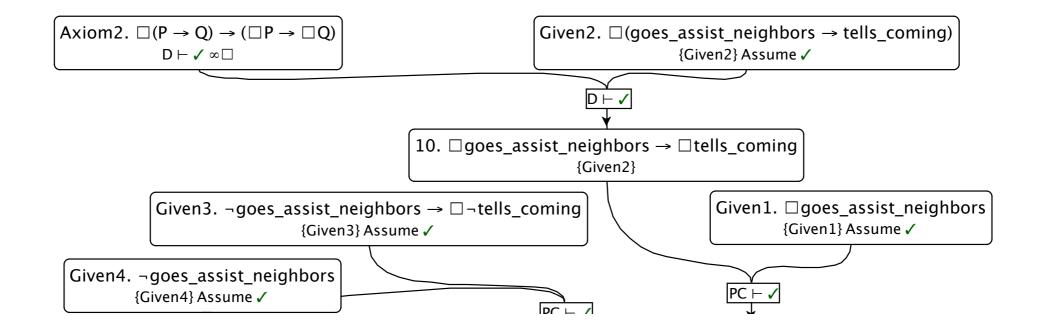
Axiom1. "All theorems of the propositional calculus." {Axiom1} Assume ✓



Axiom4. "Modus ponens for provability." {Axiom4} Assume ✓

Axiom5. "Theorems are obligatory." {Axiom5} Assume ✓

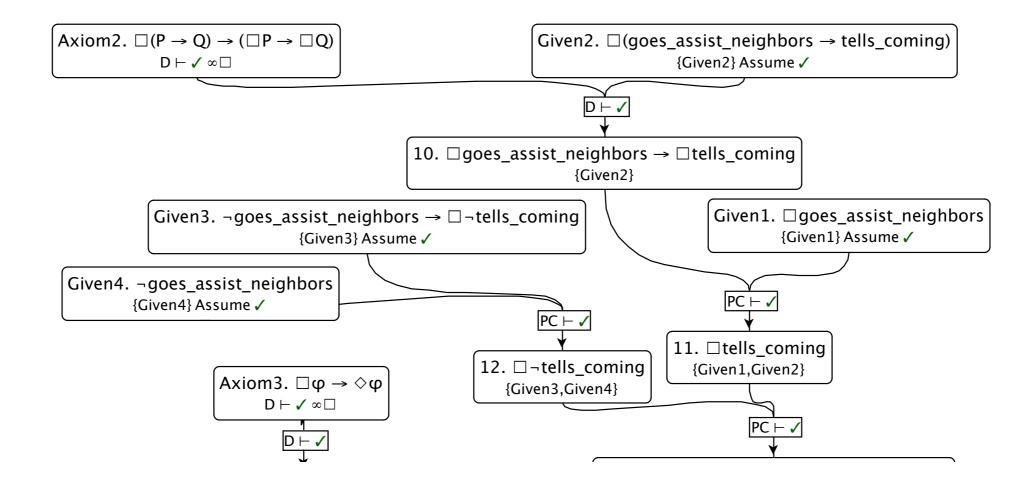
Axiom1. "All theorems of the propositional calculus."
{Axiom1} Assume ✓



Axiom4. "Modus ponens for provability."
{Axiom4} Assume ✓

Axiom5. "Theorems are obligatory."
{Axiom5} Assume ✓

Axiom1. "All theorems of the propositional calculus."
{Axiom1} Assume ✓

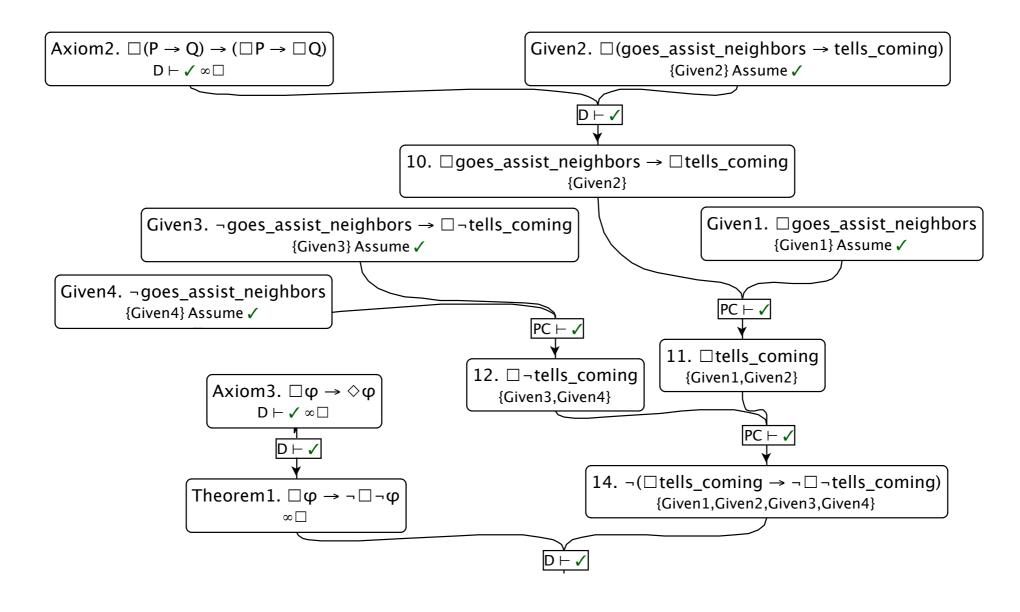


Axiom4. "Modus ponens for provability."
{Axiom4} Assume

Axiom5. "Theorems are obligatory."
{Axiom5} Assume

Axiom1. "All theorems of the propositional calculus."
{Axiom1} Assume

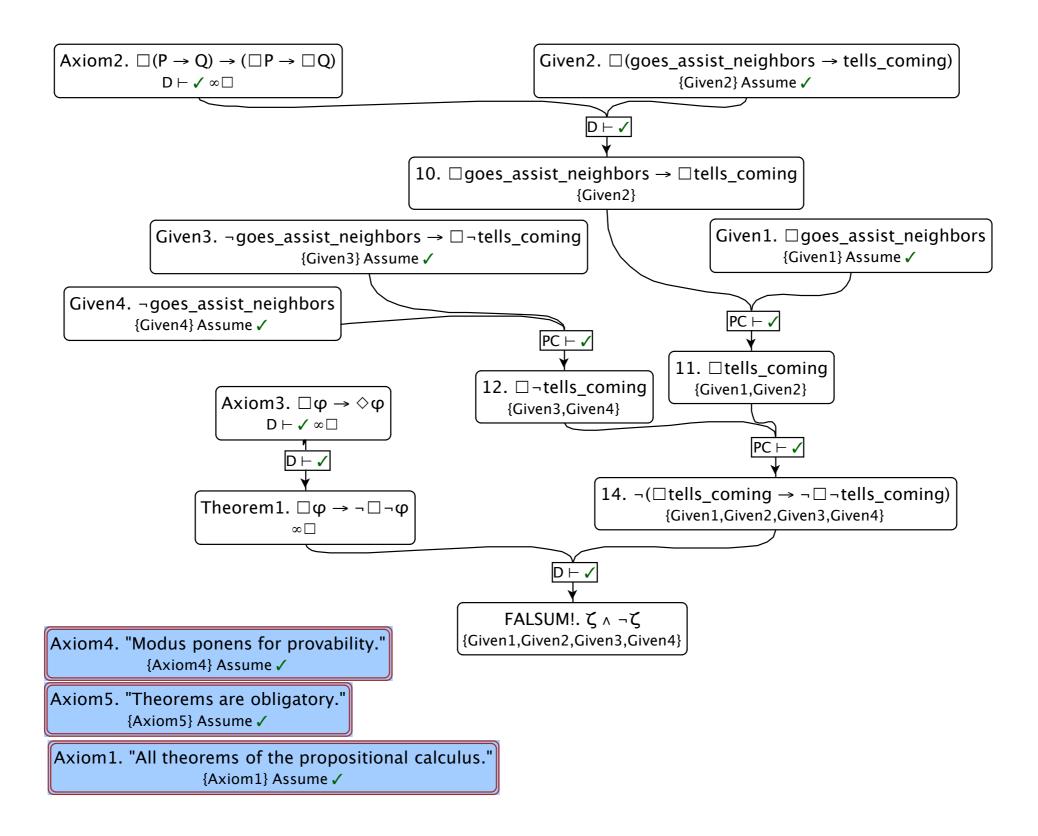
Axiom1 Assume

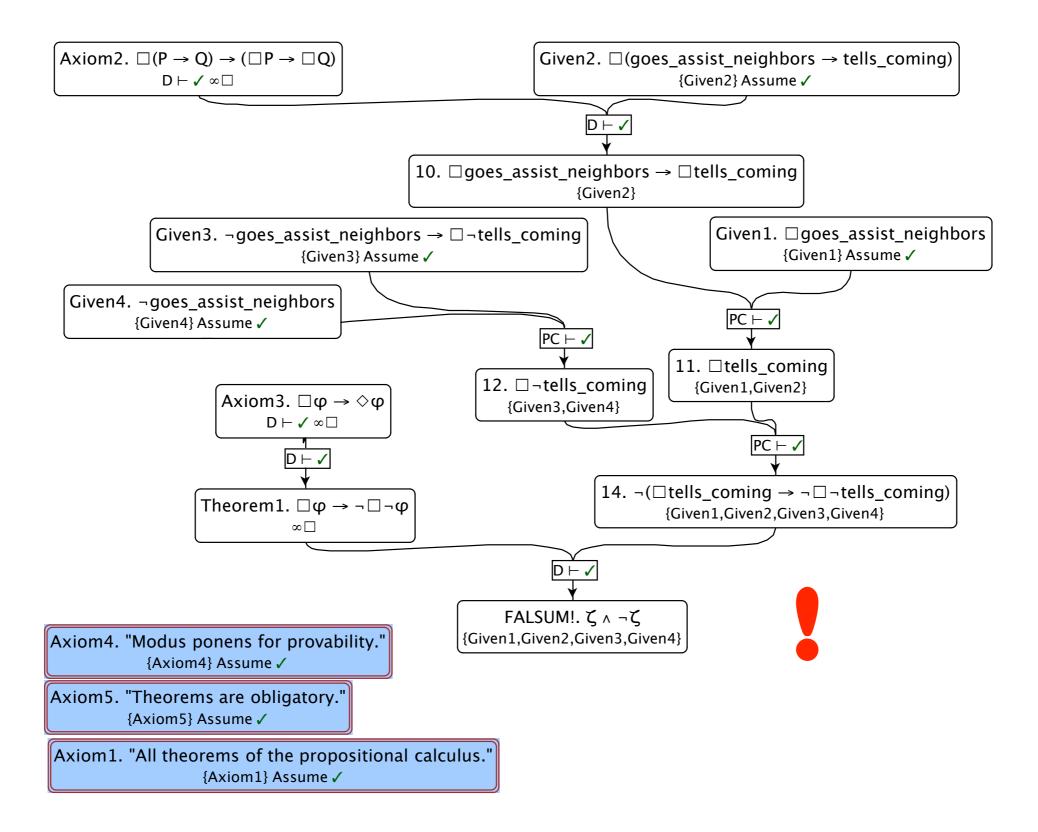


Axiom4. "Modus ponens for provability."
{Axiom4} Assume

Axiom5. "Theorems are obligatory."
{Axiom5} Assume

Axiom1. "All theorems of the propositional calculus."
{Axiom1} Assume





☑ **②** ChisholmsParadox

Here you are asked to build a proof that confirms *Chisholm's Paradox*. This paradox is that from a particular representation in **D** (= Standard Deontic Logic (SDL)) of four seemingly innocuous givens, a contradiction $\zeta \land \neg \zeta$ can be deduced. (Your instructor should have covered this in class, and may well have supplied a proof of CP.) The four givens are based on the story of a character Jones, who is obligated to go to assist his neighbors (move to a different domicile, e.g.). It would be wrong of him to show up unannounced, though; so if he goes to assist them, it ought to be that he tells them he's coming. In addition, if it's not the case that Jones goes to assist them, then it ought to be that it not be the case that he tells them he is coming. Finally, as a matter of fact, it's not case the Jones goes to assist (because on the way he comes across a car accident, and has an opportunity to save one of the victims).

Fortunately, the RAIR Lab's modern cognitive calculus \mathcal{DCEC}^* allows Chisholm's Paradox to be avoided. A recent paper explaining the use by an ethically correct AI of this calculus is available here.

Your finished proof is allowed to make use of the PC provability oracle, but of no other oracle.

Deadline 22 Apr 2020 23:59:00 EDT

SDL's = D's Problems Don't Stop Here ...

- 1. "You may either sleep on the sofa bed or the guest bed." {1} Assume ✓
- 2. "Therefore: You may sleep on the sofa bed, and you may sleep on the guest bed." {2} Assume ✓

- 1'. ♦(sofa-bed v guest-bed) {1'} Assume ✓
- 1. "You may either sleep on the sofa bed or the guest bed." {1} Assume ✓



- 2'. \diamondsuit sofa-bed $\land \diamondsuit$ guest-bed $\{1'\}$
- 2. "Therefore: You may sleep on the sofa bed, and you may sleep on the guest bed." {2} Assume ✓

- 1'. ♦(sofa-bed v guest-bed) {1'} Assume ✓
- "You may either sleep on the sofa bed or the guest bed."
 Assume ✓



- 2'. \diamond sofa-bed $\land \diamond$ guest-bed {1'}
- 2. "Therefore: You may sleep on the sofa bed, and you may sleep on the guest bed." {2} Assume ✓

NEW SCHEMA?. $\diamondsuit(\phi \lor \psi) \to (\diamondsuit\phi \land \diamondsuit\psi)$ {NEW SCHEMA?} Assume \checkmark

- 1'. ♦(sofa-bed v guest-bed) {1'} Assume ✓
- 1. "You may either sleep on the sofa bed or the guest bed." {1} Assume ✓



- 2'. \diamond sofa-bed $\land \diamond$ guest-bed {1'}
- 2. "Therefore: You may sleep on the sofa bed, and you may sleep on the guest bed." {2} Assume ✓

NEW SCHEMA?. $\diamondsuit(\phi \lor \psi) \rightarrow (\diamondsuit\phi \land \diamondsuit\psi)$ {NEW SCHEMA?} Assume \checkmark

COMMENT. "We can prove:" {COMMENT} Assume ✓

THM 5.
$$\Diamond \phi \rightarrow \Diamond (\phi \lor \psi)$$

D $\vdash \checkmark \infty \Box$

- 1'. ♦(sofa-bed v guest-bed) {1'} Assume ✓
- 1. "You may either sleep on the sofa bed or the guest bed." {1} Assume ✓



- 2'. \diamond sofa-bed $\land \diamond$ guest-bed {1'}
- 2. "Therefore: You may sleep on the sofa bed, and you may sleep on the guest bed." {2} Assume ✓

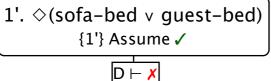
NEW SCHEMA?.
$$\diamondsuit(\phi \lor \psi) \rightarrow (\diamondsuit\phi \land \diamondsuit\psi)$$
{NEW SCHEMA?} Assume \checkmark

COMMENT. "We can prove:" {COMMENT} Assume ✓

THM 5.
$$\Diamond \phi \rightarrow \Diamond (\phi \lor \psi)$$

$$D \vdash \checkmark \infty \square$$

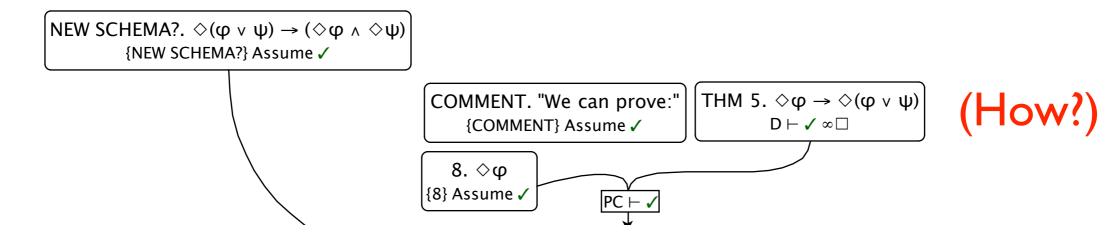
(How?)



1. "You may either sleep on the sofa bed or the guest bed." {1} Assume ✓

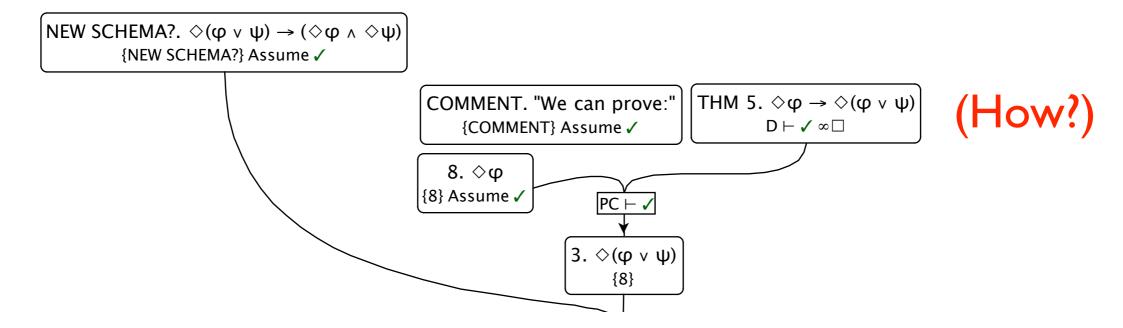


- 2'. \diamond sofa-bed $\land \diamond$ guest-bed {1'}
- 2. "Therefore: You may sleep on the sofa bed, and you may sleep on the guest bed." {2} Assume ✓



- 1'. ♦(sofa-bed v guest-bed)
 {1'} Assume ✓
- 1. "You may either sleep on the sofa bed or the guest bed." {1} Assume ✓

- D ⊢ X
- 2'. \diamond sofa-bed $\land \diamond$ guest-bed {1'}
- 2. "Therefore: You may sleep on the sofa bed, and you may sleep on the guest bed." {2} Assume ✓

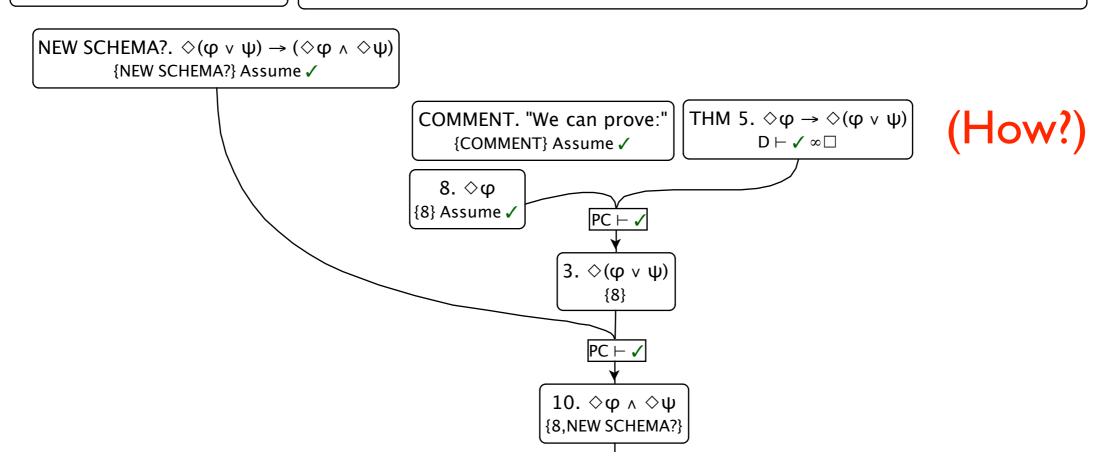


- 1'. ♦(sofa-bed v guest-bed)
 {1'} Assume ✓

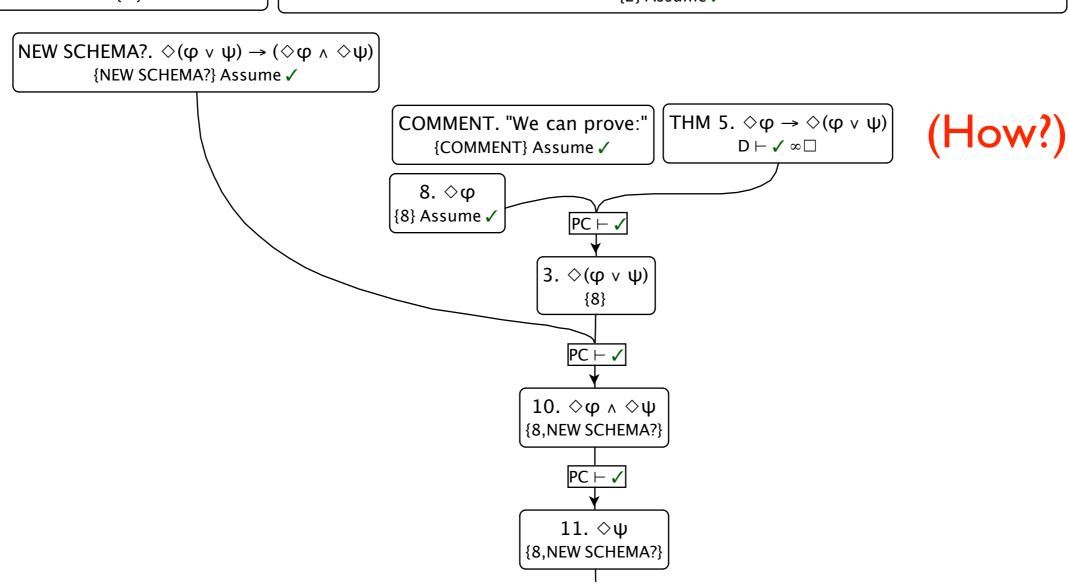
 D ⊢ X
- 1. "You may either sleep on the sofa bed or the guest bed." {1} Assume ✓
- 2'. ♦sofa-bed ∧ ♦quest-bed

{1'}

2. "Therefore: You may sleep on the sofa bed, and you may sleep on the guest bed." {2} Assume ✓

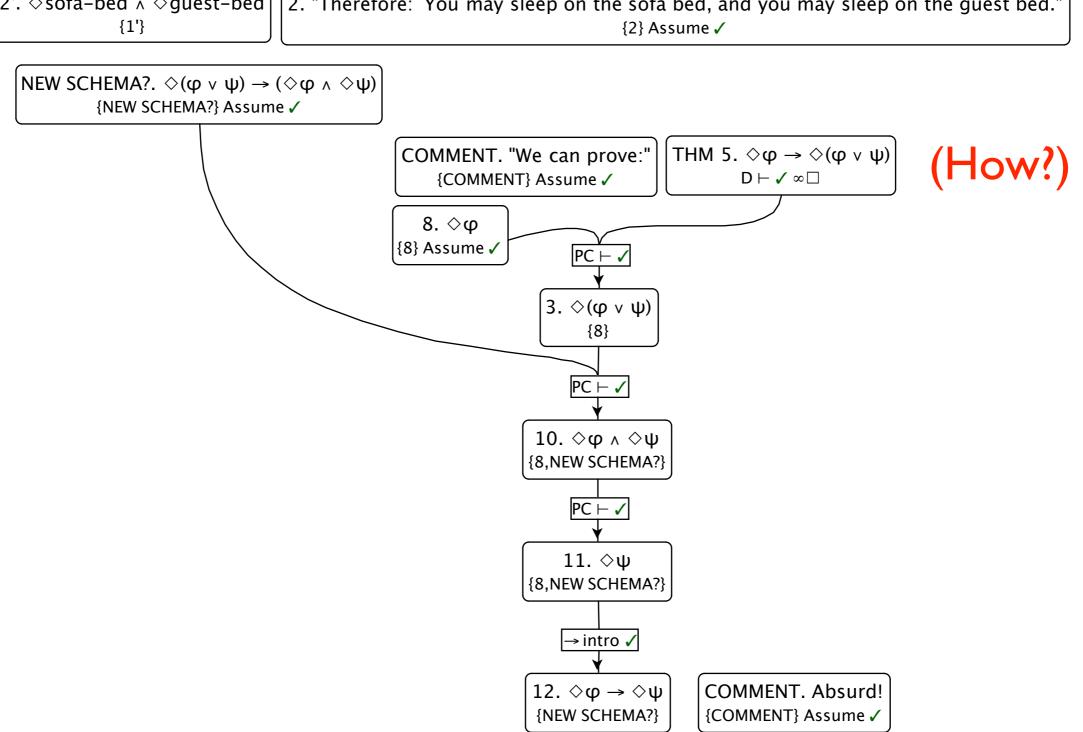


- 1'. ♦(sofa-bed v guest-bed)
 {1'} Assume ✓
- 1. "You may either sleep on the sofa bed or the guest bed." $\{1\}$ Assume \checkmark
- 2'. \diamondsuit sofa-bed $\land \diamondsuit$ guest-bed $\{1'\}$
- 2. "Therefore: You may sleep on the sofa bed, and you may sleep on the guest bed." {2} Assume ✓



{1} Assume ✓

- 1'. ♦(sofa-bed v guest-bed) 1. "You may either sleep on the sofa bed or the guest bed." {1'} Assume ✓ $D \vdash X$ 2'. ♦sofa-bed ∧ ♦quest-bed
 - 2. "Therefore: You may sleep on the sofa bed, and you may sleep on the guest bed."



"Computational logician, sorry, back to your drawing board to find a logic that works with The Four Steps!"

Producing a valid proof in this problem will enable you to understand The Free Choice Permission Paradox (FCPP), discovered in 1941 by Ross ("Imperatives and Logic," *Theoria* 7: 53–71). Given that the proof in question yields an absurdity, FCPP can be taken to show that **SDL** (Standard Deontic Logic) = **D** leads to inconsistency when applied; or, put in AI terms, you wouldn't want a robot to base its ethical decision-making on **D**! Fortunately, the RAIR Lab's modern cognitive calculus \mathcal{DCEC}^* allows FCPP to be avoided. (A recent paper explaining the use by an ethically correct AI of this calculus is available here.)

Here's the paradox. Suppose that you travel to visit a friend, arrive late at night, and are weary. Your friend says hospitably: "You may either sleep on the sofa-bed or sleep on the guest-room bed." (1) From this statement it follows that you are permitted to sleep on the sofa-bed, and you are permitted to sleep on the guest-room bed. (2) In **D**, this pair gets symbolized like this:

(1') $\Diamond(sofabed \lor guestbed)$

(2') $\Diamond sofabed \land \Diamond questbed$

But (2') doesn't follow deductively from (1') in \mathbf{D} , as a call to the provability oracle for \mathbf{D} in the HyperSlateTM file for this problem confirms. A suggested repair is to add to \mathbf{D} the schema

$$\Diamond(\phi \lor \psi) \to (\Diamond\phi \land \Diamond\psi),$$

but as your proof will (hopefully) show, this addition allows a proof of the absurd theorem that if anything is morally perimssible, everything is!

Your finished proof is allowed to make use of the PC provability oracle, but of no other oracle. (No deadline for now.)

