# Logicist Machine Ethics Can Save Us

## Selmer Bringsjord & Atriya Sen et al.

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

*Are Humans Rational?*
10/18/2018

# Logicist Machine Ethics Can Save Us

## Selmer Bringsjord & Atriya Sen et al.

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

*Are Humans Rational?*
10/18/2018

# The PAID Problem

# The PAID Problem

$\forall x : \texttt{Agents}$

# The PAID Problem

$\forall x$ : Agents

**P**owerful(x) + **A**utonomous(x) + **I**ntelligent(x) = **D**angerous(x)/
**D**estroy_Us

# The PAID Problem

$\forall x : \text{Agents}$

**P**owerful(x) + **A**utonomous(x) + **I**ntelligent(x) = **D**angerous(x)/

**D**estroy_Us

# The PAID Problem

$\forall x : \text{Agents}$

**P**owerful(x) + **A**utonomous(x) + **I**ntelligent(x) = **D**angerous(x)/

**D**estroy_Us

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

# The PAID Problem

$\forall x : \text{Agents}$

**P**owerful(x) + **A**utonomous(x) + **I**ntelligent(x) = **D**angerous(x)/
**D**estroy_Us

## 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

# The PAID Problem

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

**P**owerful(x) + **A**utonomous(x) + **I**ntelligent(x) = **D**angerous(x)/

**D**estroy_Us

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

# The PAID Problem

$\forall x :$ Agents

**P**owerful(x) + **A**utonomous(x) + **I**ntelligent(x) = **D**angerous(x)/

**D**estroy_Us

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

**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 $\gamma$ in part by way of a contributed action $\alpha_k$ 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 $o$ will regard $a'$ to be untrustworthy with respect to the pair $\langle \alpha, \gamma \rangle$ pair. **QED**

# "We're in *very* deep trouble."

# "We're in *very* deep trouble."

# "We're in *very* deep trouble."

Unfortunately, not quite as easy as this to use logic to save the day …

# Logic Thwarts Landru!



First Suspicion That It's a Mere Computer Running the Show

# Logic Thwarts Landru!



Landru is Indeed Merely a Computer
(the real Landru having done the programming)

# Logic Thwarts Landru!



Landru Kills Himself Because Kirk/Spock Argue He Has Violated
the Prime Directive for Good by Denying Creativity to Others

# Logic Thwarts Nomad!
## (with the Liar Paradox)

# I.
# Cognitive Calculi …

# Hierarchy of Ethical Reasoning

$$\mathcal{DCEC}^*_{CL}$$

$$\mathcal{DCEC}^*$$

$$\mathcal{ADR}^M$$

$$\mathcal{U}$$

UIMA/Watson-inspired

DIARC

# Hierarchy of Ethical Reasoning

$$\mathcal{DCEC}^*_{CL}$$

$$\mathcal{DCEC}^*$$

$$\mathcal{ADR}^M$$

$$\mathcal{U}$$

UIMA/Watson-inspired

DIARC

# Hierarchy of Ethical Reasoning

# Hierarchy of Ethical Reasoning

*Not* paradox-prone deontic logics!

$$\mathcal{DCEC}^*_{CL}$$

$$\mathcal{DCEC}^*$$

$$\mathcal{ADR}^M$$

$$\mathcal{U}$$

UIMA/Watson-inspired

DIARC

"Universal Cognitive Calculus"

$\mathcal{DCEC}^*$

Logic Theorist
(birth of modern logicist AI)

66

1666

1956

2017

Leibniz

Simon

1.5 centuries < Boole!
2.5 centuries < Kripke

$\int$

**Syntax**

$S ::=$ Object | Agent | Self $\sqsubseteq$ Agent | ActionType | Action $\sqsubseteq$ Event | Moment | Boolean | Fluent | Numeric

$action$ : Agent $\times$ ActionType $\to$ Action
$initially$ : Fluent $\to$ Boolean
$holds$ : Fluent $\times$ Moment $\to$ Boolean
$happens$ : Event $\times$ Moment $\to$ Boolean
$clipped$ : Moment $\times$ Fluent $\times$ Moment $\to$ Boolean
$f ::=$ $initiates$ : Event $\times$ Fluent $\times$ Moment $\to$ Boolean
$terminates$ : Event $\times$ Fluent $\times$ Moment $\to$ Boolean
$prior$ : Moment $\times$ Moment $\to$ Boolean
$interval$ : Moment $\times$ Boolean
$* $ : Agent $\to$ Self
$payoff$ : Agent $\times$ ActionType $\times$ Moment $\to$ Numeric

$t ::= x : S \mid c : S \mid f(t_1, \ldots, t_n)$

$\phi ::=$
$t$ : Boolean $\mid \neg\phi \mid \phi \wedge \psi \mid \phi \vee \psi \mid \forall x : S. \phi \mid \exists x : S. \phi$
$\mathbf{P}(a,t,\phi) \mid \mathbf{K}(a,t,\phi) \mid \mathbf{C}(t,\phi) \mid \mathbf{S}(a,b,t,\phi) \mid \mathbf{S}(a,t,\phi)$
$\mathbf{B}(a,t,\phi) \mid \mathbf{D}(a,t,holds(f,t')) \mid \mathbf{I}(a,t,happens(action(a^*,\alpha),t'))$
$\mathbf{O}(a,t,\phi,happens(action(a^*,\alpha),t'))$

**Rules of Inference**

$$\dfrac{\mathbf{C}(t,\mathbf{P}(a,t,\phi) \to \mathbf{K}(a,t,\phi))}{} \; [R_1] \qquad \dfrac{\mathbf{C}(t,\mathbf{K}(a,t,\phi) \to \mathbf{B}(a,t,\phi))}{} \; [R_2]$$

$$\dfrac{\mathbf{C}(t,\phi) \, t \leq t_1 \ldots}{\mathbf{K}(a_1,t_1,\ldots \mathbf{K}(a_n,t_n,\phi)\ldots)} \; [R_3] \qquad \dfrac{\mathbf{K}(a,t,\phi)}{\phi} \; [R_4]$$

$$\dfrac{\mathbf{C}(t,\mathbf{K}(a,t_1,\phi_1 \to \phi_2) \to \mathbf{K}(a,t_2,\phi_1) \to \mathbf{K}(a,t_3,\phi_2))}{} \; [R_5]$$

$$\dfrac{\mathbf{C}(t,\mathbf{B}(a,t_1,\phi_1 \to \phi_2) \to \mathbf{B}(a,t_2,\phi_1) \to \mathbf{B}(a,t_3,\phi_2))}{} \; [R_6]$$

$$\dfrac{\mathbf{C}(t,\mathbf{C}(t_1,\phi_1 \to \phi_2) \to \mathbf{C}(t_2,\phi_1) \to \mathbf{C}(t_3,\phi_3))}{} \; [R_7]$$

$$\dfrac{\mathbf{C}(t,\forall x. \, \phi \to \phi[x \mapsto t])}{} \; [R_8] \qquad \dfrac{\mathbf{C}(t,\phi_1 \leftrightarrow \phi_2 \to \neg\phi_2 \to \neg\phi_1)}{} \; [R_9]$$

$$\dfrac{\mathbf{C}(t, \ldots \to \psi])}{} \; [R_{10}]$$

$$\dfrac{\mathbf{B}(a, \ldots) \; \mathbf{B}(a,t,\psi)}{\ldots t,\psi \wedge \phi)} \; [R_{11b}]$$

$$\dfrac{\mathbf{B}(h, \ldots)}{} \qquad \dfrac{\mathbf{I}(a, \ldots)}{\mathbf{P}(a,t,happens(action(a^*,\alpha),t))}$$

$$\dfrac{\mathbf{B}(a,t,\phi) \; \mathbf{B}(a,t,\mathbf{O}(a^*,t,\phi,happens(action(a^*,\alpha),t'))) \quad \mathbf{O}(a,t,\phi,happens(action(a^*,\alpha),t'))}{\mathbf{K}(a,t,\mathbf{I}(a^*,t,happens(action(a^*,\alpha),t')))} \; [R_{14}]$$

$$\dfrac{\phi \leftrightarrow \psi}{\mathbf{O}(a,t,\phi,\gamma) \leftrightarrow \mathbf{O}(a,t,\psi,\gamma)} \; [R_{15}]$$

**R A I R**
Rensselaer AI and Reasoning Lab

AI of Today: What Would Leibniz Say?

"Sorry, not impressed."

Selmer Bringsjord

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

Turin Italy
11/14/2016

R A I R
Rensselaer AI and Reasoning Lab

# Cognitive Calculi $\mathcal{CC}$

purely extensional level:   FOL   MSL   SOL   TOL   IFOL   . . .

theories:   **PA ZFC** axiomatic physics . . .

intensional level:   epistemic   deontic   possibility/necessity   . . .

model finders:   MACE . . .

ATPs:   SPASS   SNARK   ShadowProver . . .

nature of representation: symbolic or homomorphic: . . .

- - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -

# Cognitive Calculi $\mathscr{CC}$

purely
extensional
level:  FOL  MSL  SOL  TOL  IFOL  . . .

theories:  **PA ZFC** axiomatic physics  . . .

intensional
level:  epistemic  deontic  possibility/necessity  . . .

model finders:  MACE  . . .

nature of representation: symbolic or homomorphic:
. . .

ATPs:  SPASS  SNARK  ShadowProver  . . .

- - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -

# Cognitive Calculi $\mathscr{CC}$

purely
extensional
level:

FOL  MSL  SOL  TOL  IFOL  . . .

theories:  **PA ZFC** axiomatic physics  . . .

intensional
level:

epistemic  deontic  possibility/necessity  . . .

model finders:  MACE  . . .

ATPs:  SPASS  SNARK  ShadowProver  . . .

nature of representation: symbolic or homomorphic:
. . .

# Cognitive Calculi $\mathcal{CC}$

purely
extensional
level:  FOL  MSL  SOL  TOL  IFOL  . . .

theories:  **PA ZFC** axiomatic physics  . . .

intensional
level:  epistemic  deontic  possibility/necessity  . . .

model finders:  MACE  . . .

ATPs:  SPASS  SNARK  ShadowProver  . . .

nature of representation:  symbolic or homomorphic:
. . .

$\lambda$-calculus

$\lambda$-calculus

# Cognitive Calculi $\mathcal{CC}$

purely extensional level:    FOL   MSL   SOL   TOL   IFOL   . . .

theories:   **PA ZFC** axiomatic physics   . . .

intensional level:    epistemic   deontic   possibility/necessity   . . .

model finders:   MACE   . . .

ATPs:    SPASS   SNARK   ShadowProver   . . .

nature of representation: symbolic or homomorphic: . . .

$\lambda$-calculus

$\lambda$-calculus

. . .    analogical reasoning

   inductive reasoning   . . .

inference schemas   $\infty$

# Cognitive Calculi $\mathscr{CC}$

purely extensional level: FOL MSL SOL TOL IFOL ...

theories: **PA ZFC** axiomatic physics ...

intensional level: epistemic deontic possibility/necessity ...

model finders: MACE ...

ATPs: SPASS SNARK ShadowProver ...

nature of representation: symbolic or homomorphic: ...

$\lambda$-calculus

$\mathcal{D_\mathcal{I}CEC^*}$  $\mathcal{DCEC^*}$  $\mathcal{DCSC^*}$  $\mathcal{CEC}$  $\mathcal{CSC}$  ...

$\lambda$-calculus

dialects:

analogical reasoning

inference schemas $\infty$

inductive reasoning

# Cognitive Calculi $\mathscr{CC}$



purely extensional level: FOL MSL SOL TOL IFOL . . .

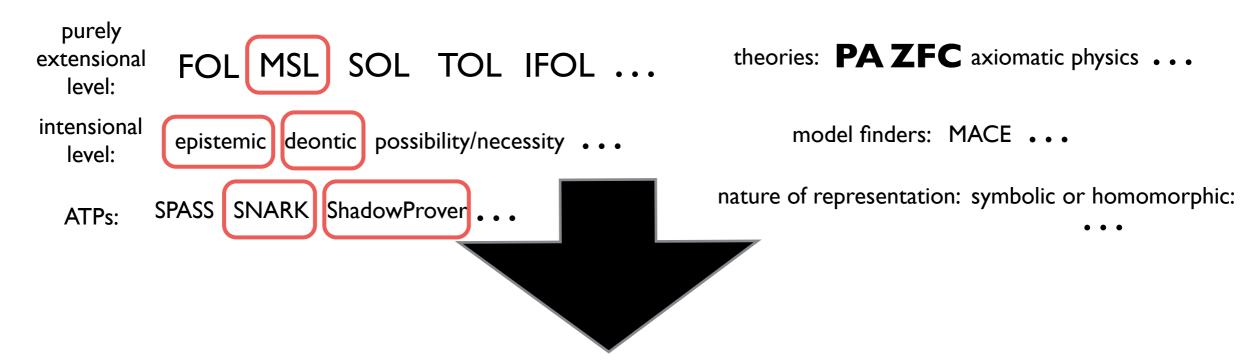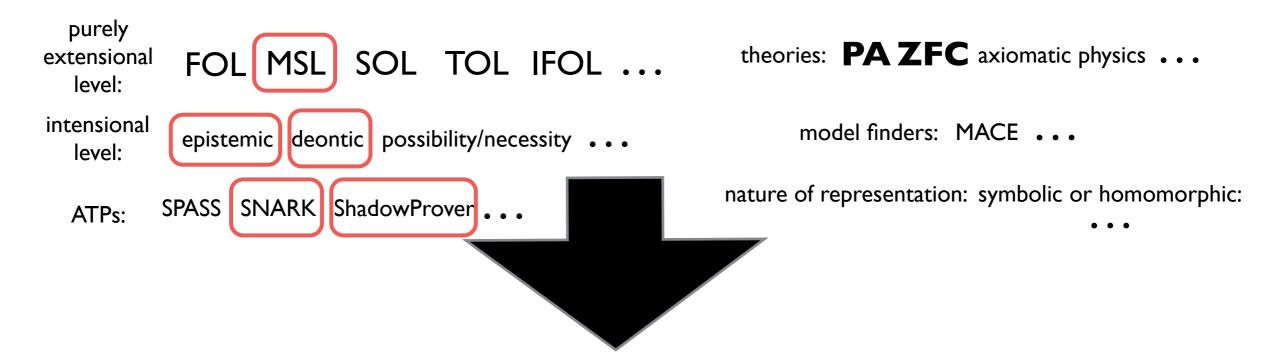intensional level: epistemic deontic possibility/necessity . . .

ATPs: SPASS SNARK ShadowProver . . .

theories: **PA ZFC** axiomatic physics . . .

model finders: MACE . . .

nature of representation: symbolic or homomorphic: . . .

$\lambda$-calculus

. . . $\mathcal{D_{\not E}CEC^*}$ $\mathcal{DCEC^*}$ $\mathcal{DCSC^*}$ $\mathcal{CEC}$ $\mathcal{CSC}$ . . .

$\lambda$-calculus

dialects:

analogical reasoning

inference schemas $\infty$

inductive reasoning

# Cognitive Calculi $\mathscr{CC}$

purely
extensional    FOL  MSL  SOL  TOL  IFOL  . . .
level:

intensional
level:    epistemic  deontic  possibility/necessity  . . .

ATPs:    SPASS  SNARK  ShadowProver  . . .
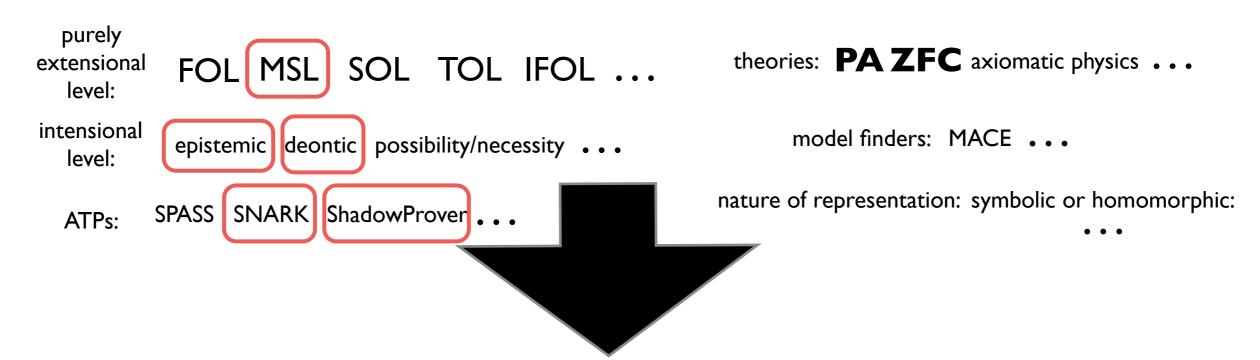
theories:  **PA ZFC** axiomatic physics  . . .

model finders:  MACE  . . .

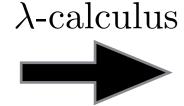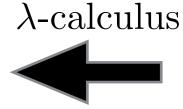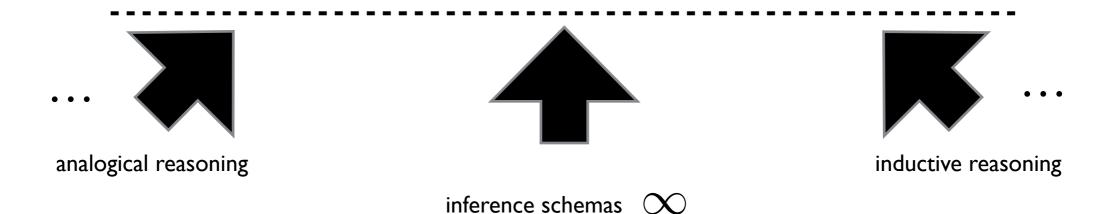nature of representation: symbolic or homomorphic:
. . .

$\lambda$-calculus    . . .  $\mathcal{D_{I}CEC^*}$  $\mathcal{DCEC^*}$  $\mathcal{DCSC^*}$  $\mathcal{CEC}$  $\mathcal{CSC}$  . . .  $\lambda$-calculus

dialects:

analogical reasoning          inductive reasoning

inference schemas  $\infty$

# Formal Syntax

# Formal Syntax

$S ::=$ Object | Agent | Self $\sqsubset$ Agent | ActionType | Action $\sqsubseteq$ Event | Moment | Boolean | Fluent | Numeric

$f ::=$
*action* : Agent $\times$ ActionType $\rightarrow$ Action

*initially* : Fluent $\rightarrow$ Boolean

*holds* : Fluent $\times$ Moment $\rightarrow$ Boolean

*happens* : Event $\times$ Moment $\rightarrow$ Boolean

*clipped* : Moment $\times$ Fluent $\times$ Moment $\rightarrow$ *Boolean*

*initiates* : Event $\times$ Fluent $\times$ Moment $\rightarrow$ Boolean

*terminates* : Event $\times$ Fluent $\times$ Moment $\rightarrow$ Boolean

*prior* : Moment $\times$ Moment $\rightarrow$ Boolean

*interval* : Moment $\times$ Boolean

$*$ : Agent $\rightarrow$ Self

*payoff* : Agent $\times$ ActionType $\times$ Moment $\rightarrow$ Numeric

$t ::= x : S \mid c : S \mid f(t_1, \ldots, t_n)$

$\phi ::=$
$t$ : Boolean $\mid \neg\phi \mid \phi \wedge \psi \mid \phi \vee \psi \mid$

$\mathbf{P}(a,t,\phi) \mid \mathbf{K}(a,t,\phi) \mid \mathbf{C}(t,\phi) \mid \mathbf{S}(a,b,t,\phi) \mid \mathbf{S}(a,t,\phi)$

$\mathbf{B}(a,t,\phi) \mid \mathbf{D}(a,t,holds(f,t')) \mid \mathbf{I}(a,t,happens(action(a^*,\alpha),t'))$

$\mathbf{O}(a,t,\phi,happens(action(a^*,\alpha),t'))$

# Inference Schemata

# Inference Schemata

$$\frac{}{\mathbf{C}(t, \mathbf{P}(a,t,\phi) \to \mathbf{K}(a,t,\phi))} \ [R_1] \quad \frac{}{\mathbf{C}(t, \mathbf{K}(a,t,\phi) \to \mathbf{B}(a,t,\phi))} \ [R_2]$$

$$\frac{\mathbf{C}(t,\phi) \ t \leq t_1 \ldots t \leq t_n}{\mathbf{K}(a_1,t_1,\ldots \mathbf{K}(a_n,t_n,\phi)\ldots)} \ [R_3] \quad \frac{\mathbf{K}(a,t,\phi)}{\phi} \ [R_4]$$

$$\frac{}{\mathbf{C}(t,\mathbf{K}(a,t_1,\phi_1 \to \phi_2)) \to \mathbf{K}(a,t_2,\phi_1) \to \mathbf{K}(a,t_3,\phi_2)} \ [R_5]$$

$$\frac{}{\mathbf{C}(t,\mathbf{B}(a,t_1,\phi_1 \to \phi_2)) \to \mathbf{B}(a,t_2,\phi_1) \to \mathbf{B}(a,t_3,\phi_2)} \ [R_6]$$

$$\frac{}{\mathbf{C}(t,\mathbf{C}(t_1,\phi_1 \to \phi_2)) \to \mathbf{C}(t_2,\phi_1) \to \mathbf{C}(t_3,\phi_2)} \ [R_7]$$

$$\frac{}{\mathbf{C}(t,\forall x. \ \phi \to \phi[x \mapsto t])} \ [R_8] \quad \frac{}{\mathbf{C}(t,\phi_1 \leftrightarrow \phi_2 \to \neg\phi_2 \to \neg\phi_1)} \ [R_9]$$

$$\frac{}{\mathbf{C}(t,[\phi_1 \wedge \ldots \wedge \phi_n \to \phi] \to [\phi_1 \to \ldots \to \phi_n \to \psi])} \ [R_{10}]$$

$$\frac{\mathbf{B}(a,t,\phi) \ \ \phi \to \psi}{\mathbf{B}(a,t,\psi)} \ [R_{11a}] \quad \frac{\mathbf{B}(a,t,\phi) \ \ \mathbf{B}(a,t,\psi)}{\mathbf{B}(a,t,\psi \wedge \phi)} \ [R_{11b}]$$

$$\frac{\mathbf{S}(s,h,t,\phi)}{\mathbf{B}(h,t,\mathbf{B}(s,t,\phi))} \ [R_{12}]$$

$$\frac{\mathbf{I}(a,t,happens(action(a^*,\alpha),t'))}{\mathbf{P}(a,t,happens(action(a^*,\alpha),t))} \ [R_{13}]$$

$$\frac{\mathbf{B}(a,t,\phi) \ \ \mathbf{B}(a,t,\mathbf{O}(a^*,t,\phi,happens(action(a^*,\alpha),t'))) \quad \mathbf{O}(a,t,\phi,happens(action(a^*,\alpha),t'))}{\mathbf{K}(a,t,\mathbf{I}(a^*,t,happens(action(a^*,\alpha),t')))} \ [R_{14}]$$

$$\frac{\phi \leftrightarrow \psi}{\mathbf{O}(a,t,\phi,\gamma) \leftrightarrow \mathbf{O}(a,t,\psi,\gamma)} \ [R_{15}]$$

# Event Calculus for Time & Change

$$\overline{\mathbf{C}(t, \mathbf{P}(a,t,\phi) \to \mathbf{K}(a,t,\phi))} \ [R_1]$$

$$\overline{\mathbf{C}(t, \mathbf{K}(a,t,\phi) \to \mathbf{B}(a,t,\phi))} \ [R_2]$$

$$\frac{\mathbf{C}(t,\phi) \ t \leq t_1 \ldots t \leq t_n}{\mathbf{K}(a_1,t_1,\ldots \mathbf{K}(a_n,t_n,\phi)\ldots)} \ [R_3]$$

$$\frac{\mathbf{K}(a,t,\phi)}{\phi} \ [R_4]$$

$$\overline{\mathbf{C}(t, \mathbf{K}(a,t_1,\phi_1 \to \phi_2)) \to \mathbf{K}(a,t_2,\phi_1) \to \mathbf{K}(a,t_3,\phi_2)} \ [R_5]$$

$$\overline{\mathbf{C}(t, \mathbf{B}(a,t_1,\phi_1 \to \phi_2)) \to \mathbf{B}(a,t_2,\phi_1) \to \mathbf{B}(a,t_3,\phi_2)} \ [R_6]$$

$$\overline{\mathbf{C}(t, \mathbf{C}(t_1,\phi_1 \to \phi_2)) \to \mathbf{C}(t_2,\phi_1) \to \mathbf{C}(t_3,\phi_2)} \ [R_7]$$

$$\overline{\mathbf{C}(t, \forall x. \ \phi \to \phi[x \mapsto t])} \ [R_8]$$

$$\overline{\mathbf{C}(t, \phi_1 \leftrightarrow \phi_2 \to \neg\phi_2 \to \neg\phi_1)} \ [R_9]$$

$$\overline{\mathbf{C}(t, [\phi_1 \wedge \ldots \wedge \phi_n \to \phi] \to [\phi_1 \to \ldots \to \phi_n \to \psi])} \ [R_{10}]$$

$$\frac{\mathbf{B}(a,t,\phi) \ \ \phi \to \psi}{\mathbf{B}(a,t,\psi)} \ [R_{11a}]$$

$$\frac{\mathbf{B}(a,t,\phi) \ \ \mathbf{B}(a,t,\psi)}{\mathbf{B}(a,t,\psi \wedge \phi)} \ [R_{11b}]$$

$$\frac{\mathbf{S}(s,h,t,\phi)}{\mathbf{B}(h,t,\mathbf{B}(s,t,\phi))} \ [R_{12}]$$

$$\frac{\mathbf{I}(a,t,happens(action(a^*,\alpha),t'))}{\mathbf{P}(a,t,happens(action(a^*,\alpha),t))} \ [R_{13}]$$

$$\frac{\mathbf{B}(a,t,\phi) \ \ \mathbf{B}(a,t,\mathbf{O}(a^*,t,\phi,happens(action(a^*,\alpha),t'))) \ \ \mathbf{O}(a,t,\phi,happens(action(a^*,\alpha),t'))}{\mathbf{K}(a,t,\mathbf{I}(a^*,t,happens(action(a^*,\alpha),t')))} \ [R_{14}]$$

$$\frac{\phi \leftrightarrow \psi}{\mathbf{O}(a,t,\phi,\gamma) \leftrightarrow \mathbf{O}(a,t,\psi,\gamma)} \ [R_{15}]$$

# Event Calculus for Time & Change

$$\frac{}{\mathbf{C}(t,\mathbf{P}(a,t,\phi)\to\mathbf{K}(a,t,\phi))}\ [R_1] \quad \frac{}{\mathbf{C}(t,\mathbf{K}(a,t,\phi)\to\mathbf{B}(a,t,\phi))}\ [R_2]$$

$$\frac{\mathbf{C}(t,\phi)\ t\leq t_1\ldots t\leq t_n}{\mathbf{K}(a_1,t_1,\ldots\mathbf{K}(a_n,t_n,\phi)\ldots)}\ [R_3] \quad \frac{\mathbf{K}(a,t,\phi)}{\phi}\ [R_4]$$

$$\frac{}{\mathbf{C}(t,\mathbf{K}(a,t_1,\phi_1\to\phi_2))\to\mathbf{K}(a,t_2,\phi_1)\to\mathbf{K}(a,t_3,\phi_2)}\ [R_5]$$

$$\frac{}{\mathbf{C}(t,\mathbf{B}(a,t_1,\phi_1\to\phi_2))\to\mathbf{B}(a,t_2,\phi_1)\to\mathbf{B}(a,t_3,\phi_2)}\ [R_6]$$

$$\frac{}{\mathbf{C}(t,\mathbf{C}(t_1,\phi_1\to\phi_2))\to\mathbf{C}(t_2,\phi_1)\to\mathbf{C}(t_3,\phi_2)}\ [R_7]$$

$$\frac{}{\mathbf{C}(t,\forall x.\ \phi\to\phi[x\mapsto t])}\ [R_8] \quad \frac{}{\mathbf{C}(t,\phi_1\leftrightarrow\phi_2\to\neg\phi_2\to\neg\phi_1)}\ [R_9]$$

$$\frac{}{\mathbf{C}(t,[\phi_1\wedge\ldots\wedge\phi_n\to\phi]\to[\phi_1\to\ldots\to\phi_n\to\psi])}\ [R_{10}]$$

$$\frac{\mathbf{B}(a,t,\phi)\ \ \phi\to\psi}{\mathbf{B}(a,t,\psi)}\ [R_{11a}] \quad \frac{\mathbf{B}(a,t,\phi)\ \ \mathbf{B}(a,t,\psi)}{\mathbf{B}(a,t,\psi\wedge\phi)}\ [R_{11b}]$$

$$\frac{\mathbf{S}(s,h,t,\phi)}{\mathbf{B}(h,t,\mathbf{B}(s,t,\phi))}\ [R_{12}]$$

$$\frac{\mathbf{I}(a,t,happens(action(a^*,\alpha),t'))}{\mathbf{P}(a,t,happens(action(a^*,\alpha),t))}\ [R_{13}]$$

$$\frac{\mathbf{B}(a,t,\phi)\ \ \mathbf{B}(a,t,\mathbf{O}(a^*,t,\phi,happens(action(a^*,\alpha),t')))\quad\mathbf{O}(a,t,\phi,happens(action(a^*,\alpha),t'))}{\mathbf{K}(a,t,\mathbf{I}(a^*,t,happens(action(a^*,\alpha),t')))}\ [R_{14}]$$

$$\frac{\phi\leftrightarrow\psi}{\mathbf{O}(a,t,\phi,\gamma)\leftrightarrow\mathbf{O}(a,t,\psi,\gamma)}\ [R_{15}]$$

$[A_1]$ $\mathbf{C}(\forall f,t\ .\ initially(f)\wedge\neg clipped(0,f,t)\Rightarrow holds(f,t))$

$[A_2]$ $\mathbf{C}(\forall e,f,t_1,t_2\ .\ happens(e,t_1)\wedge initiates(e,f,t_1)\wedge t_1<t_2\wedge\neg clipped(t_1,f,t_2)\Rightarrow holds(f,t_2))$

$[A_3]$ $\mathbf{C}(\forall t_1,f,t_2\ .\ clipped(t_1,f,t_2)\Leftrightarrow[\exists e,t\ .\ happens(e,t)\wedge t_1<t<t_2\wedge terminates(e,f,t)])$

$[A_4]$ $\mathbf{C}(\forall a,d,t\ .\ happens(action(a,d),t)\Rightarrow\mathbf{K}(a,happens(action(a,d),t)))$

$[A_5]$ $\mathbf{C}(\forall a,f,t,t'\ .\ \mathbf{B}(a,holds(f,t))\wedge\mathbf{B}(a,t<t')\wedge\neg\mathbf{B}(a,clipped(t,f,t'))\Rightarrow\mathbf{B}(a,holds(f,t')))$

# Defs for An *Affective* Cognitive *time&change* Calculus

1. **Joy** : pleased about a desirable event. By 'pleased about a desirable event' the meaning we will consider is 'pleased about a desirable consequence of the event'.

$$forSome \ c \ B(a, t_3, implies(happens(e, t_1), holds(CON(e, a, c), t_2))) \quad (1)$$

$$D(a, t_3, holds(CON(e, a, c), t_2)) \quad (2)$$

$$K(a, t_3, happens(e, t_1)) \quad (3)$$

The definition of $holds(AFF(a, joy), t_3)$ is therefore and(1,2,3).

2. **Distress** : displeased about an undesirable event.

$$not(D(a, t_3, holds(CON(e, a, c), t_3))) \quad (4)$$

The definition of $holds(AFF(a, distress), t_3)$ is therefore and(1,4,3).

3. **Happy-for**: pleased about an event presumed to be desirable for someone else

$$forSome \ c \ B(a, t_3, implies(happens(e, t_1), holds(CON(e, a_1, c), t_2))) \quad (5)$$

$$B(a, t_3, D(a_1, t_3, holds(CON(e, a_1, c), t_2))) \quad (6)$$

$$D(a, t_3, holds(CON(e, a_1, c), t_2)) \quad (7)$$

The definition of $holds(AFF(a, happy\_for), t_3)$ is therefore and(5,6,7,3).

4. **Pity**: displeased about an event presumed to be undesirable for someone else. This is equivalent to sorry_for in Hobbs-Gordon model.

$$B(a, t_3, not(D(a_1, t_3, holds(CON(e, a_1, c), t_2)))) \quad (8)$$

$$not(D(a, t_3, holds(CON(e, a_1, c), t_2))) \quad (9)$$

The definition of $holds(AFF(a, pity), t_3)$ is therefore and(5,8,9,3).

5. **Gloating** : pleased about an event presumed to be undesirable for someone else The definition of $holds(AFF(a, gloating), t_3)$ is therefore and(5,8,7,3).

6. **Resentment**: displeased about an event presumed to be desirable for someone else The definition of $holds(AFF(a, resentment), t_3)$ is therefore and(5,6,9,3).

7. **Hope**: (pleased about) the prospect of a desirable event

$$forSome \ c \ B(a, t_0, implies(happens(e, t_1), \diamond holds(CON(e, a, c), t_2))) \quad (10)$$

$$D(a, t_0, holds(CON(e, a, c), t_2)) \quad (11)$$

The definition of $holds(AFF(a, hope), t_0)$ is therefore and(10,11).

8. **Fear**: (displeased about) the prospect of an undesirable event

$$not(D(a, t_0, holds(CON(e, a, c), t_2))) \quad (12)$$

The definition of $holds(AFF(a, fear), t_0)$ is therefore and(10,12).

9. **Satisfaction** : (pleased about) the confirmation of the prospect of a desirable event
The definition of $holds(AFF(a, satisfaction), t_3)$ is and(10,11, 7 3).

10. **Fears-confirmed** : (displeased about) the confirmation of the prospect of an undesirable event.
The definition of $holds(AFF(a, fears - confirmed), t_3)$ is and(10,12,9, 3).

11. **Relief**: (pleased about) the disconfirmation of the prospect of an undesirable event

$$K(a, t_3, not(happens(e, t_1))) \quad (13)$$

The definition of $holds(AFF(a, relief), t_3)$ is and(10, 12, 9, 13).

12. **Disappointment** : (displeased about) the disconfirmation of the prospect of a desirable event
The definition of $holds(AFF(a, disappointment), t_3)$ is and(10, 11, 7, 13).

13. **Pride** : (approving of) one's own praiseworthy action
Here we treat 'approve' as an action event. We also introduce a new predicate $PRAISEWORTHY(a, b, x)$ which will mean that agent a considers x a praiseworthy action by agent b. All the 3 interpretations are shown below.

$$happens(action(a, x), t_0) \quad (14)$$

$$forAll \ a_x B(a, t_1, implies(happens(action(a_x, x), t_x), PRAISEWORTHY(a, a_x, x))), t_x \leq t_1 \quad (15)$$

$$D(a, t_1, holds(PRAISEWORTHY(a, a, x), t_1)) \quad (16)$$

$$happens(action(a, approve(x)), t_1) \quad (17)$$

The definition of $holds(AFF(a, pride), t_1)$ is and(14, $B(a, t_1, holds(PRAISEWORTHY(a, a, x), t_1))$, 17).

14. **Shame**: (disapproving of) one's own blameworthy action
This also follows the same explanation as Pride.

$$forAll \ a_x B(a, t_1, implies(happens(action(a_x, x), t_x), B(a, t_1, holds(BLAMEWORTHY(a, a_x, x)), t_1)))), t_x \leq t_1 \quad (18)$$

$$not(happens(action(a, approve(x)), t_1)) \quad (19)$$

The definition of $holds(AFF(a, shame), t_1)$ is and(14, $B(a, t_1, holds(BLAMEWORTHY(a, a, x), t_1))$, 19).

15. **Admiration**: (approving of) someone else's praiseworthy action

$$happens(action(a_1, x), t_0) \quad (20)$$

The definition of $holds(AFF(a, admiration), t_1)$ is and(20, $B(a, t_1, holds(PRAISEWORTHY(a, a_1, x), t_1))$, 17).

16. **Reproach**: (disapproving of) someone else's blameworthy action The definition of $holds(AFF(a, reproach), t_1)$ is and(20, $B(a, t_1, holds(BLAMEWORTHY(a, a_1, x), t_1))$, 19).

17. **Gratification** : (approving of) one's own praiseworthy action and (being pleased about) the related desirable event. We again interpret 'pleased about the desirable event' as 'pleased about the desired consequence of the event.'

$$forSome \ c \ B(a, t_1, implies(happens(action(a, x), t_0), holds(CON(action(a, x), a, c), t_0))) \quad (21)$$

$$D(a, t_1, holds(CON(action(a, x), a, c), t_0)) \quad (22)$$

The definition of $holds(AFF(a, gratification), t_1)$ is and(20, $B(a, t_1, holds(PRAISEWORTHY(a, a, x), t_1))$, 17,

**… (and more)**

# II.
# Early Progress With Our Calculi: Simple Dilemmas; Non-Akratic Robots

**NewScientist**

Ethical robots save humans

# Informal Definition of Akrasia

# Informal Definition of Akrasia

A
agent

# Informal Definition of Akrasia

A agent

# Informal Definition of Akrasia

# Informal Definition of Akrasia

# Informal Definition of Akrasia

# Informal Definition of Akrasia

agent

**A**

desired

$\alpha_f$

$t_{\alpha_f}$

# Informal Definition of Akrasia

# Informal Definition of Akrasia

# Informal Definition of Akrasia

# Informal Definition of Akrasia

# Informal Definition of Akrasia

# Informal Definition of Akrasia



agent

A

desired

$\alpha_f$

obligatory

$\alpha_o$

$t_{\alpha_f}$

$t_{\alpha_o}$

If $\alpha_f$ happens, then $\alpha_o$ can't happen

A

# Informal Definition of Akrasia

# Informal Definition of Akrasia

# Informal Definition of Akrasia

_____

# Informal Definition of Akrasia

# Informal Definition of Akrasia

$$t_{\alpha_f}$$

# Informal Definition of Akrasia

# Informal Definition of Akrasia

A

Desire to do $\alpha_f$

$t_{\alpha_f}$

# Informal Definition of Akrasia

# Informal Definition of Akrasia

# Informal Definition of Akrasia

A

Desire to do $\alpha_f$  $\succ$  Belief that he ought to do $\alpha_o$

A does $\alpha_f$ due to his desire

$t_{\alpha_f}$

# Informal Definition of Akrasia

A

Desire to do $\alpha_f$ $\succ$ Belief that he ought to do $\alpha_o$

A does $\alpha_f$ due to his desire

$t_{\alpha_f}$

# Informal Definition of Akrasia

A

Desire to do $\alpha_f$ $\succ$ Belief that he ought to do $\alpha_o$

A does $\alpha_f$ due to his desire

$t_{\alpha_f}$

# Informal Definition of Akrasia

A

Desire to do $\alpha_f$ $\succ$ Belief that he ought to do $\alpha_o$

A does $\alpha_f$ due to his desire

$$t_{\alpha_f} \qquad t$$

# Informal Definition of Akrasia

A

Desire to do $\alpha_f$ $\succ$ Belief that he ought to do $\alpha_o$

A does $\alpha_f$ due to his desire

A believes he should have done $\alpha_o$

$t_{\alpha_f}$       $t$

# Informal Definition of Akrasia

An action $\alpha_f$ is (Augustinian) akratic for an agent $A$ at $t_{\alpha_f}$ iff the following eight conditions hold:

(1)  $A$ believes that $A$ ought to do $\alpha_o$ at $t_{\alpha_o}$;

(2)  $A$ desires to do $\alpha_f$ at $t_{\alpha_f}$;

(3)  $A$'s doing $\alpha_f$ at $t_{\alpha_f}$ entails his not doing $\alpha_o$ at $t_{\alpha_o}$;

(4)  $A$ knows that doing $\alpha_f$ at $t_{\alpha_f}$ entails his not doing $\alpha_o$ at $t_{\alpha_o}$;

(5)  At the time ($t_{\alpha_f}$) of doing the forbidden $\alpha_f$, $A$'s desire to do $\alpha_f$ overrides $A$'s belief that he ought to do $\alpha_o$ at $t_{\alpha_f}$.

(6)  $A$ does the forbidden action $\alpha_f$ at $t_{\alpha_f}$;

(7)  $A$'s doing $\alpha_f$ results from $A$'s desire to do $\alpha_f$;

(8)  At some time $t$ after $t_{\alpha_f}$, $A$ has the belief that $A$ ought to have done $\alpha_o$ rather than $\alpha_f$.

# Informal Definition of Akrasia

An action $\alpha_f$ is (Augustinian) akratic for an agent $A$ at $t_{\alpha_f}$ iff the following eight conditions hold:

(1)  $A$ believes that $A$ ought to do $\alpha_o$ at $t_{\alpha_o}$;

(2)  $A$ desires to do $\alpha_f$ at $t_{\alpha_f}$;

(3)  $A$'s doing $\alpha_f$ at $t_{\alpha_f}$ entails his not doing $\alpha_o$ at $t_{\alpha_o}$;

(4)  $A$ knows that doing $\alpha_f$ at $t_{\alpha_f}$ entails his not doing $\alpha_o$ at $t_{\alpha_o}$;

(5)  At the time ($t_{\alpha_f}$) of doing the forbidden $\alpha_f$, $A$'s desire to do $\alpha_f$ overrides $A$'s belief that he ought to do $\alpha_o$ at $t_{\alpha_f}$.

(6)  $A$ does the forbidden action $\alpha_f$ at $t_{\alpha_f}$;

(7)  $A$'s doing $\alpha_f$ results from $A$'s desire to do $\alpha_f$;

(8)  At some time $t$ after $t_{\alpha_f}$, $A$ has the belief that $A$ ought to have done $\alpha_o$ rather than $\alpha_f$.

# Informal Definition of Akrasia

An action $\alpha_f$ is (Augustinian) akratic for an agent $A$ at $t_{\alpha_f}$ iff the following eight conditions hold:

(1)  $A$ believes that $A$ ought to do $\alpha_o$ at $t_{\alpha_o}$;

(2)  $A$ desires to do $\alpha_f$ at $t_{\alpha_f}$;

(3)  $A$'s doing $\alpha_f$ at $t_{\alpha_f}$ entails his not doing $\alpha_o$ at $t_{\alpha_o}$;

(4)  $A$ knows that doing $\alpha_f$ at $t_{\alpha_f}$ entails his not doing $\alpha_o$ at $t_{\alpha_o}$;

(5)  At the time ($t_{\alpha_f}$) of doing the forbidden $\alpha_f$, $A$'s desire to do $\alpha_f$ overrides $A$'s belief that he ought to do $\alpha_o$ at $t_{\alpha_f}$.

(6)  $A$ does the forbidden action $\alpha_f$ at $t_{\alpha_f}$;

(7)  $A$'s doing $\alpha_f$ results from $A$'s desire to do $\alpha_f$;

"Regret" (8)  At some time $t$ after $t_{\alpha_f}$, $A$ has the belief that $A$ ought to have done $\alpha_o$ rather than $\alpha_f$.

Cast in

$\mathcal{DCEC}^*$

this becomes …

$$\mathsf{KB}_{rs} \cup \mathsf{KB}_{m_1} \cup \mathsf{KB}_{m_2} \ldots \mathsf{KB}_{m_n} \vdash$$

$$D_1 : \mathbf{B}(\mathsf{I}, \mathsf{now}, \mathbf{O}(\mathsf{I}^*, t_\alpha \Phi, happens(action(\mathsf{I}^*, \alpha), t_\alpha)))$$

$$D_2 : \mathbf{D}(\mathsf{I}, \mathsf{now}, holds(does(\mathsf{I}^*, \overline{\alpha}), t_{\overline{\alpha}}))$$

$$D_3 : happens(action(\mathsf{I}^*, \overline{\alpha}), t_{\overline{\alpha}}) \Rightarrow \neg happens(action(\mathsf{I}^*, \alpha), t_\alpha)$$

$$D_4 : \mathbf{K}\left( \mathsf{I}, \mathsf{now}, \begin{pmatrix} happens(action(\mathsf{I}^*, \overline{\alpha}), t_{\overline{\alpha}}) \Rightarrow \\ \neg happens(action(\mathsf{I}^*, \alpha), t_\alpha) \end{pmatrix} \right)$$

$$D_5 : \begin{matrix} \mathbf{I}(\mathsf{I}, t_\alpha, happens(action(\mathsf{I}^*, \overline{\alpha}), t_{\overline{\alpha}}) \wedge \\ \neg \mathbf{I}(\mathsf{I}, t_\alpha, happens(action(\mathsf{I}^*, \alpha), t_\alpha) \end{matrix}$$

$$D_6 : happens(action(\mathsf{I}^*, \overline{\alpha}), t_{\overline{\alpha}})$$

$$D_{7a} : \begin{matrix} \Gamma \cup \{\mathbf{D}(\mathsf{I}, \mathsf{now}, holds(does(\mathsf{I}^*, \overline{\alpha}), t))\} \vdash \\ happens(action(\mathsf{I}^*, \overline{\alpha}), t_\alpha) \end{matrix}$$

$$D_{7b} : \begin{matrix} \Gamma - \{\mathbf{D}(\mathsf{I}, \mathsf{now}, holds(does(\mathsf{I}^*, \overline{\alpha}), t))\} \not\vdash \\ happens(action(\mathsf{I}^*, \overline{\alpha}), t_\alpha) \end{matrix}$$

$$D_8 : \mathbf{B}\big(\mathsf{I}, t_f, \mathbf{O}(\mathsf{I}^*, t_\alpha, \Phi, happens(action(\mathsf{I}^*, \alpha), t_\alpha))\big)$$

# Demos …

# Demos …

# III.
# But, a twist befell the logicists …

Chisholm had argued that the three old 19th-century ethical categories (*forbidden*, *morally neutral*, *obligatory*) are not enough — and soul-searching brought me to agreement.

# Leibnizian Ethical Hierarchy for Persons and Robots:
## $\mathcal{EH}$

| deviltry | uncivil | forbidden | morally neutral | obligatory | civil | heroic |

# Leibnizian Ethical Hierarchy for Persons and Robots:
## $\mathscr{EH}$

the supererogatory

| deviltry | uncivil | forbidden | morally neutral | obligatory | civil | heroic |

# Leibnizian Ethical Hierarchy for Persons and Robots:
## $\mathscr{EH}$

the supererogatory

| deviltry | uncivil | forbidden | morally neutral | obligatory | civil | heroic |

# Leibnizian Ethical Hierarchy for Persons and Robots:
## $\mathscr{EH}$

the subererogatory

the supererogatory

| deviltry | uncivil | forbidden | morally neutral | obligatory | civil | heroic |

# Leibnizian Ethical Hierarchy for Persons and Robots:
## $\mathscr{EH}$

(see Norwegian crime fiction)

the subererogatory

the supererogatory

| deviltry | uncivil | forbidden | morally neutral | obligatory | civil | heroic |

# Leibnizian Ethical Hierarchy for Persons and Robots:
## $\mathcal{EH}$

# Leibnizian Ethical Hierarchy for Persons and Robots:

## $\mathscr{EH}$



(see Norwegian crime fiction)

the subererogatory

the supererogatory

| deviltry | uncivil | forbidden | morally neutral | obligatory | civil | heroic |

# Leibnizian Ethical Hierarchy for Persons and Robots:
## $\mathscr{EH}$



(see Norwegian crime fiction)

the subererogatory

the supererogatory

| deviltry | uncivil | forbidden | morally neutral | obligatory | civil | heroic |

focus of others

# Leibnizian Ethical Hierarchy for Persons and Robots:
## $\mathscr{EH}$



(see Norwegian crime fiction)

the subererogatory

the supererogatory

| deviltry | uncivil | forbidden | morally neutral | obligatory | civil | heroic |

But *this* portion may be most relevant to military missions.

focus of others

$$\mathcal{T} := \|\mathcal{F}|\mathcal{P} \wedge \neg \mathcal{O}|\mathcal{O}\| \quad \text{19th Century Triad}$$

$$\mathcal{T} := \|\mathcal{F}|\mathcal{P} \wedge \neg\mathcal{O}|\mathcal{O}\| \quad \text{19th Century Triad}$$

$$\mathcal{T} := \|\mathcal{F}|\mathcal{P} \wedge \neg\mathcal{O}|\mathcal{O}\| \quad \text{19th Century Triad}$$

$$
\begin{array}{c|c|c}
\mathcal{F} & \mathcal{P} \wedge \neg\mathcal{O} & \mathcal{O} \\
\forall \quad \text{F} \quad \text{M} \quad \text{V} \quad \exists & & \forall \quad \text{F} \quad \text{M} \quad \text{V} \quad \exists
\end{array}
$$

$$\mathcal{T} := \|\mathcal{F}|\mathcal{P} \wedge \neg\mathcal{O}|\mathcal{O}\| \quad \text{19th Century Triad}$$

| $\mathcal{F}$ | $\mathcal{P} \wedge \neg\mathcal{O}$ | $\mathcal{O}$ |
|---|---|---|
| ∀ F M V ∃ | | ∀ F M V ∃ |

$$\mathcal{EH}$$

| $\mathcal{S}_{ub1}$ | $\mathcal{S}_{ub2}$ | $\mathcal{F}$ | $\mathcal{P}\wedge\neg\mathcal{O}$ | $\mathcal{O}^L$ | $\mathcal{O}^M$ | $\mathcal{S}^{up1}$ | $\mathcal{S}^{up2}$ |
|---|---|---|---|---|---|---|---|
| ∃–∀ | ∃–∀ | ∃–∀ | | ∃–∀ | ∃–∀ | ∃–∀ | ∃–∀ |
| | | | | | | ↑ | |

$$\mathcal{T} := \|\mathcal{F}|\mathcal{P} \wedge \neg\mathcal{O}|\mathcal{O}\| \quad \text{19th Century Triad}$$

$$
\begin{array}{c|c|c}
\mathcal{F} & \mathcal{P} \wedge \neg\mathcal{O} & \mathcal{O} \\
\forall \quad \text{F} \quad \text{M} \quad \text{V} \quad \exists & & \forall \quad \text{F} \quad \text{M} \quad \text{V} \quad \exists
\end{array}
$$

$$\mathcal{EH}$$

| $\mathcal{S}_{ub1}$ | $\mathcal{S}_{ub2}$ | $\mathcal{F}$ | $\mathcal{P} \wedge \neg\mathcal{O}$ | $\mathcal{O}^L$ | $\mathcal{O}^M$ | $\mathcal{S}^{up1}$ | $\mathcal{S}^{up2}$ |
|---|---|---|---|---|---|---|---|
| $\exists$–$\forall$ | $\exists$–$\forall$ | $\exists$–$\forall$ | | $\exists$–$\forall$ | $\exists$–$\forall$ | $\exists$–$\forall$ | $\exists$–$\forall$ |
| | | | | | | $\uparrow$ | |

●

$$\mathcal{T} := \|\mathcal{F}|\mathcal{P} \wedge \neg\mathcal{O}|\mathcal{O}\| \quad \text{19th Century Triad}$$

| $\mathcal{F}$ | | | | | $\mathcal{P} \wedge \neg\mathcal{O}$ | | $\mathcal{O}$ | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| $\forall$ | F | M | V | $\exists$ | | | $\forall$ | F | M | V | $\exists$ |

$$\mathcal{EH}$$

| $\mathcal{S}_{ub1}$ | $\mathcal{S}_{ub2}$ | $\mathcal{F}$ | $\mathcal{P} \wedge \neg\mathcal{O}$ | $\mathcal{O}^L$ | $\mathcal{O}^M$ | $\mathcal{S}^{up1}$ | $\mathcal{S}^{up2}$ |
|---|---|---|---|---|---|---|---|
| $\exists{-}\forall$ | $\exists{-}\forall$ | $\exists{-}\forall$ | | $\exists{-}\forall$ | $\exists{-}\forall$ | $\exists{-}\forall$ | $\exists{-}\forall$ |
| | | | | | | $\uparrow$ | |

●

Arkin
Pereira
Andersons
Powers
Mikhail
…

$$\mathcal{T} := \|\mathcal{F}|\mathcal{P} \wedge \neg\mathcal{O}|\mathcal{O}\| \quad \text{19th Century Triad}$$

| $\mathcal{F}$ | $\mathcal{P} \wedge \neg\mathcal{O}$ | $\mathcal{O}$ |
|---|---|---|
| $\forall$  F  M  V  $\exists$ | | $\forall$  F  M  V  $\exists$ |

$$\mathcal{EH}$$

| $\mathcal{S}_{ub1}$ | $\mathcal{S}_{ub2}$ | $\mathcal{F}$ | $\mathcal{P} \wedge \neg\mathcal{O}$ | $\mathcal{O}^L$ | $\mathcal{O}^M$ | $\mathcal{S}^{up1}$ | $\mathcal{S}^{up2}$ |
|---|---|---|---|---|---|---|---|
| $\exists$–$\forall$ | $\exists$–$\forall$ | $\exists$–$\forall$ | | $\exists$–$\forall$ | $\exists$–$\forall$ | $\exists$–$\forall$ | $\exists$–$\forall$ |
| | | | | ● | | $\uparrow$ | |

$$\mathscr{T} := \|\mathcal{F}|\mathcal{P} \wedge \neg\mathcal{O}|\mathcal{O}\| \quad \text{19th Century Triad}$$

$$
\begin{array}{c|c|c}
\mathcal{F} & \mathcal{P} \wedge \neg\mathcal{O} & \mathcal{O} \\
\forall \quad \mathbf{F} \quad \mathbf{M} \quad \mathbf{V} \quad \exists & & \forall \quad \mathbf{F} \quad \mathbf{M} \quad \mathbf{V} \quad \exists
\end{array}
$$

$$\mathcal{EH}$$

| $\mathcal{S}_{ub1}$ | $\mathcal{S}_{ub2}$ | $\mathcal{F}$ | $\mathcal{P} \wedge \neg\mathcal{O}$ | $\mathcal{O}^L$ | $\mathcal{O}^M$ | $\mathcal{S}^{up1}$ | $\mathcal{S}^{up2}$ |
|---|---|---|---|---|---|---|---|
| $\exists$–$\forall$ | $\exists$–$\forall$ | $\exists$–$\forall$ | | $\exists$–$\forall$ | $\exists$–$\forall$ | $\exists$–$\forall$ | $\exists$–$\forall$ |
| | | | | $\bullet$ | | $\uparrow$ | |

$$\mathcal{T} := \|\mathcal{F}|\mathcal{P} \wedge \neg\mathcal{O}|\mathcal{O}\| \quad \text{19th Century Triad}$$

| | $\mathcal{F}$ | | | | | $\mathcal{P} \wedge \neg\mathcal{O}$ | | | $\mathcal{O}$ | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | $\forall$ | **F** | **M** | **V** | $\exists$ | | | | $\forall$ | **F** | **M** | **V** | $\exists$ |

$$\mathcal{EH}$$

| $\mathcal{S}_{ub1}$ | $\mathcal{S}_{ub2}$ | $\mathcal{F}$ | $\mathcal{P} \wedge \neg\mathcal{O}$ | $\mathcal{O}^L$ | $\mathcal{O}^M$ | $\mathcal{S}^{up1}$ | $\mathcal{S}^{up2}$ |
|---|---|---|---|---|---|---|---|
| $\exists$–$\forall$ | $\exists$–$\forall$ | $\exists$–$\forall$ | | $\exists$–$\forall$ | $\exists$–$\forall$ | $\exists$–$\forall$ | $\exists$–$\forall$ |
| | | | | ● | | ↑ | |

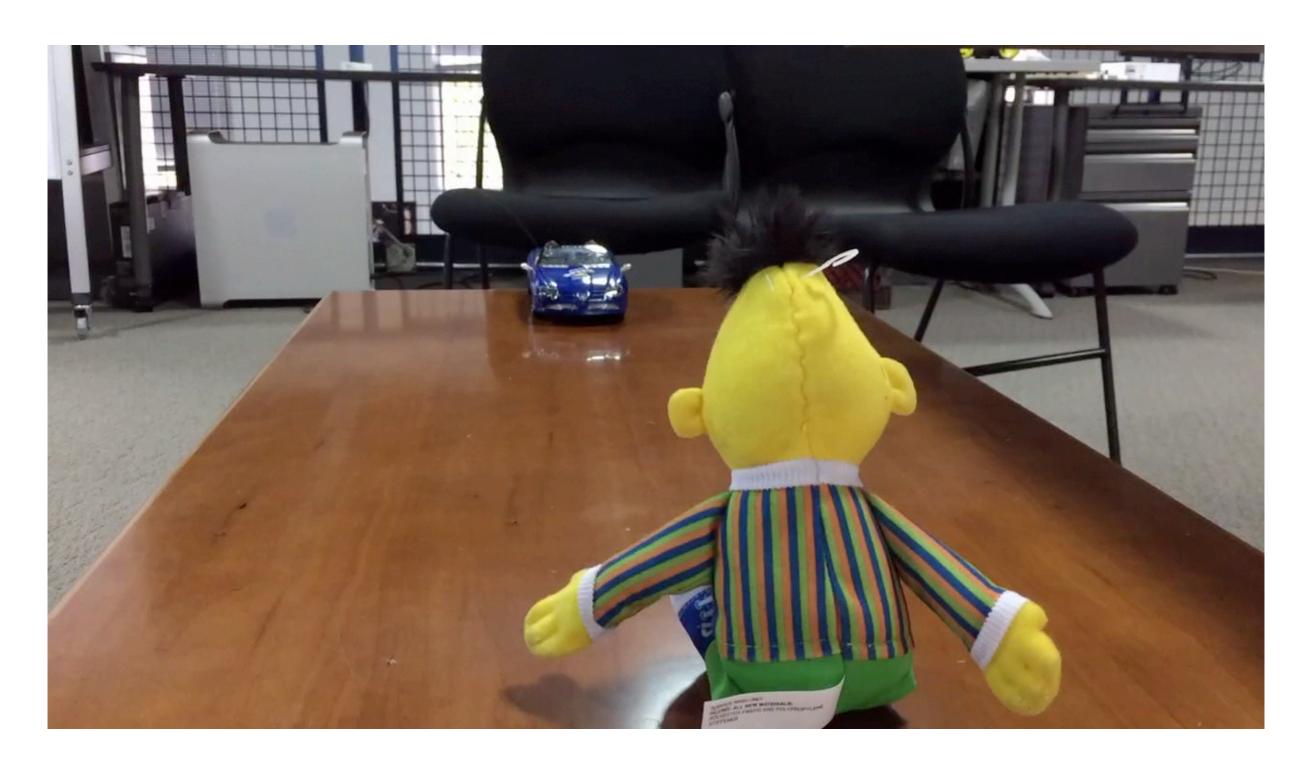There are obviously a host of formulae whose theoremhood constitute desiderata; that is (to give but a pair), the following must be provable (where $n \in \{1, 2\}$):

*Theorem 1.* $\mathbf{S^{up}}^n(\phi, a, \alpha) \to \neg\mathbf{O}(\phi, a, \alpha)$

*Theorem 2.* $\mathbf{S^{up}}^n(\phi, a, \alpha) \to \neg\mathbf{F}(\phi, a, \alpha)$

Secondly, $\mathcal{L_{EH}}$ is an *inductive* logic, not a deductive one. This must be the case, since, as we've noted, quantification isn't restricted to just the standard pair $\exists\forall$ of quantifiers in standard extensional $n$-order logic: $\mathcal{EH}$ is based on three additional quantifiers. For example, while in standard

# Bert "Heroically" Saved?



Courtesy of RAIR-Lab Researcher Atriya Sen

# Bert "Heroically" Saved?

# Supererogatory² Robot Action



Courtesy of RAIR-Lab Researcher Atriya Sen

Courtesy of RAIR-Lab Researcher Atriya Sen

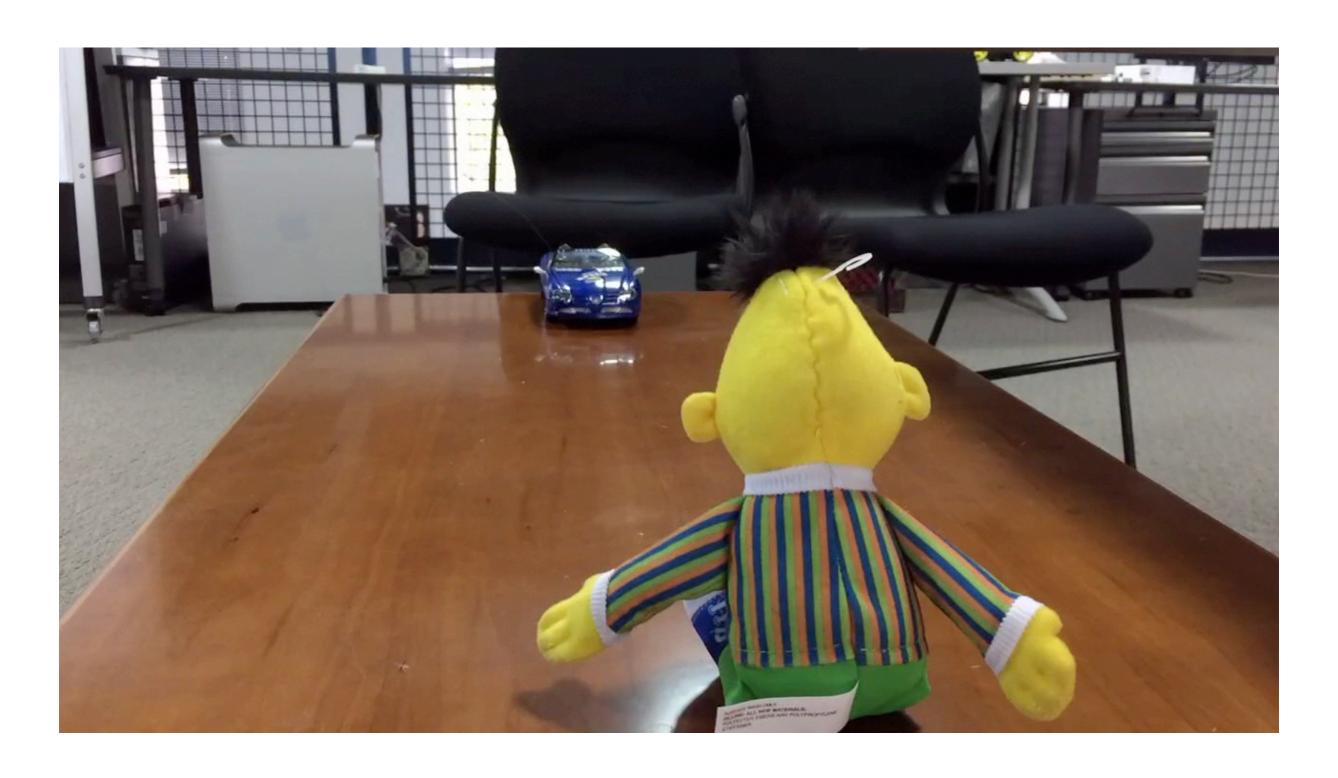# Bert "Heroically" Saved!!



Courtesy of RAIR-Lab Researcher Atriya Sen

# Bert "Heroically" Saved!!

Courtesy of RAIR-Lab Researcher Atriya Sen

$$K\left(\text{nao}, t_1, \text{lessthan}\left(\text{payoff}\left(\text{nao}^*, \neg\text{dive}, t_2\right), \text{threshold}\right)\right)$$
$$K\left(\text{nao}, t_1, \text{greaterthan}\left(\text{payoff}\left(\text{nao}^*, \text{dive}, t_2\right), \text{threshold}\right)\right)$$
$$K\left(\text{nao}, t_1, \neg O\left(\text{nao}^*, t_2, \text{lessthan}\left(\text{payoff}\left(\text{nao}^*, \neg\text{dive}, t_2\right), \text{threshold}\right), \text{happens}\left(\text{action}\left(\text{nao}^*, \text{dive}\right), t_2\right)\right)\right)$$
$$\therefore K\left(\text{nao}, t_1, S^{\text{UP2}}\left(\text{nao}, t_2, \text{happens}\left(\text{action}\left(\text{nao}^*, \text{dive}\right), t_2\right)\right)$$
$$\therefore I\left(\text{nao}, t_2, \text{happens}\left(\text{action}\left(\text{nao}^*, \text{dive}\right), t_2\right)\right)$$
$$\therefore \text{happens}\left(\text{action}(\text{nao}, \text{dive}), t_2\right)$$

Courtesy of RAIR-Lab Researcher Atriya Sen

$$K\left(\text{nao}, t_1, \text{lessthan}\left(\text{payoff}\left(\text{nao}^*, \neg\text{dive}, t_2\right), \text{threshold}\right)\right)$$

$$K\left(\text{nao}, t_1, \text{greaterthan}\left(\text{payoff}\left(\text{nao}^*, \text{dive}, t_2\right), \text{threshold}\right)\right)$$

$$K\left(\text{nao}, t_1, \neg O\left(\text{nao}^*, t_2, \text{lessthan}\left(\text{payoff}\left(\text{nao}^*, \neg\text{dive}, t_2\right), \text{threshold}\right), \text{happens}\left(\text{action}\left(\text{nao}^*, \text{dive}\right), t_2\right)\right)\right)$$

$$\therefore K\left(\text{nao}, t_1, S^{\text{UP2}}\left(\text{nao}, t_2, \text{happens}\left(\text{action}\left(\text{nao}^*, \text{dive}\right), t_2\right)\right)$$

$$\therefore I\left(\text{nao}, t_2, \text{happens}\left(\text{action}\left(\text{nao}^*, \text{dive}\right), t_2\right)\right)$$

$$\therefore \text{happens}\left(\text{action}(\text{nao}, \text{dive}), t_2\right)$$



Courtesy of RAIR-Lab Researcher Atriya Sen

# In Talos (available via Web interface); & ShadowProver

```
Prototypes:
Boolean lessThan Numeric Numeric
Boolean greaterThan Numeric Numeric
ActionType not ActionType
ActionType dive


Axioms:
lessOrEqual(Moment t1,t2)
K(nao,t1,lessThan(payoff(nao,not(dive),t2),threshold))
K(nao,t1,greaterThan(payoff(nao,dive,t2),threshold))
K(nao,t1,not(O(nao,t2,lessThan(payoff(nao,not(dive),t2),threshold),happens(action(nao,dive),t2))))


provable Conjectures:
happens(action(nao,dive),t2)
K(nao,t1,SUP2(nao,t2,happens(action(nao,dive),t2)))
I(nao,t2,happens(action(nao,dive),t2))
```

# In Talos (available via Web interface); & ShadowProver

```
Prototypes:
Boolean lessThan Numeric Numeric
Boolean greaterThan Numeric Numeric
ActionType not ActionType
ActionType dive


Axioms:
lessOrEqual(Moment t1,t2)
K(nao,t1,lessThan(payoff(nao,not(dive),t2),threshold))
K(nao,t1,greaterThan(payoff(nao,dive,t2),threshold))
K(nao,t1,not(O(nao,t2,lessThan(payoff(nao,not(dive),t2),threshold),happens(action(nao,dive),t2))))


provable Conjectures:
happens(action(nao,dive),t2)
K(nao,t1,SUP2(nao,t2,happens(action(nao,dive),t2)))
I(nao,t2,happens(action(nao,dive),t2))
```

# Hence, we now have *this* overview of the logicist engineering required:

# Making Morally *X* Machines, in Four Steps

~$10M

**Theories of Law**

Natural Law

Confucian Law

**Ethical Theories**

Shades of Utilitarianism

Legal Codes

Particular Ethical Codes

Utilitarianism

Deontological

Divine Command

Virtue Ethics

Contract

Egoism

**Step 1**

1. Pick (a) theories.
2. Pick (a) code(s).
3. Run through EH.
4. Which *X* in *MMXM*?

# Making Morally *X* Machines, in Four Steps

~$10M

## Theories of Law

**Natural Law**

**Confucian Law**

Legal Codes

## Ethical Theories

Shades
of
Utilitarianism

Particular
Ethical Codes

**Utilitarianism**

**Deontological**

**Divine Command**

**Virtue Ethics**

**Contract**

**Egoism**

### Step 1

1. Pick (a) theories.
2. Pick (a) code(s).
3. Run through EH.
4. Which *X* in *MMXM*?

### Step 2

Formalize & Automate

Shadow Prover

Spectra

# Making Morally *X* Machines, in Four Steps

~$10M

## Theories of Law

**Natural Law**

**Confucian Law**

Shades
of
Utilitarianism

Legal Codes

Particular
Ethical Codes

## Ethical Theories

**Utilitarianism**

**Deontological**

**Divine Command**

**Virtue Ethics**

**Contract**

**Egoism**

---

### Step 1

1. Pick (a) theories.
2. Pick (a) code(s).
3. Run through EH.
4. Which *X* in *MMXM*?

### Step 2

Formalize & Automate

Shadow Prover

Spectra

### Step 3

Ethical OS

Ethical Substrate

Robotic Substrate

# Making Morally *X* Machines, in Four Steps

~$10M

## Theories of Law

**Natural Law**

**Confucian Law**

Legal Codes

## Ethical Theories

Shades
of
Utilitarianism

**Utilitarianism**

**Deontological**

**Divine Command**

**Virtue Ethics**

**Contract**

**Egoism**

Particular
Ethical Codes

---

**Step 1**

1. Pick (a) theories.
2. Pick (a) code(s).
3. Run through EH.
4. Which *X* in *MMXM*?

**Step 2**

Formalize & Automate

Shadow Prover

Spectra

**Step 3**

Ethical OS

Ethical Substrate

Robotic Substrate

# Making Morally *X* Machines, in Four Steps

~$10M

## Theories of Law

**Natural Law**

**Confucian Law**

Shades of Utilitarianism

Legal Codes

Particular Ethical Codes

## Ethical Theories

**Utilitarianism**

**Deontological**

**Divine Command**
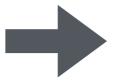
**Virtue Ethics**

**Contract**

**Egoism**

**Step 1**

1. Pick (a) theories.
2. Pick (a) code(s).
3. Run through EH.
4. Which *X* in *MMXM*?

**Step 2**

Formalize & Automate

Shadow Prover

Spectra

**Step 3**

Ethical OS

Ethical Substrate

Robotic Substrate

*An ethically correct robot.*

# Making Morally *X* Machines, in Four Steps

~$10M

## Theories of Law

**Natural Law**

**Confucian Law**

Shades
of
Utilitarianism

Legal Codes

Particular
Ethical Codes

## Ethical Theories

**Utilitarianism**

**Deontological**

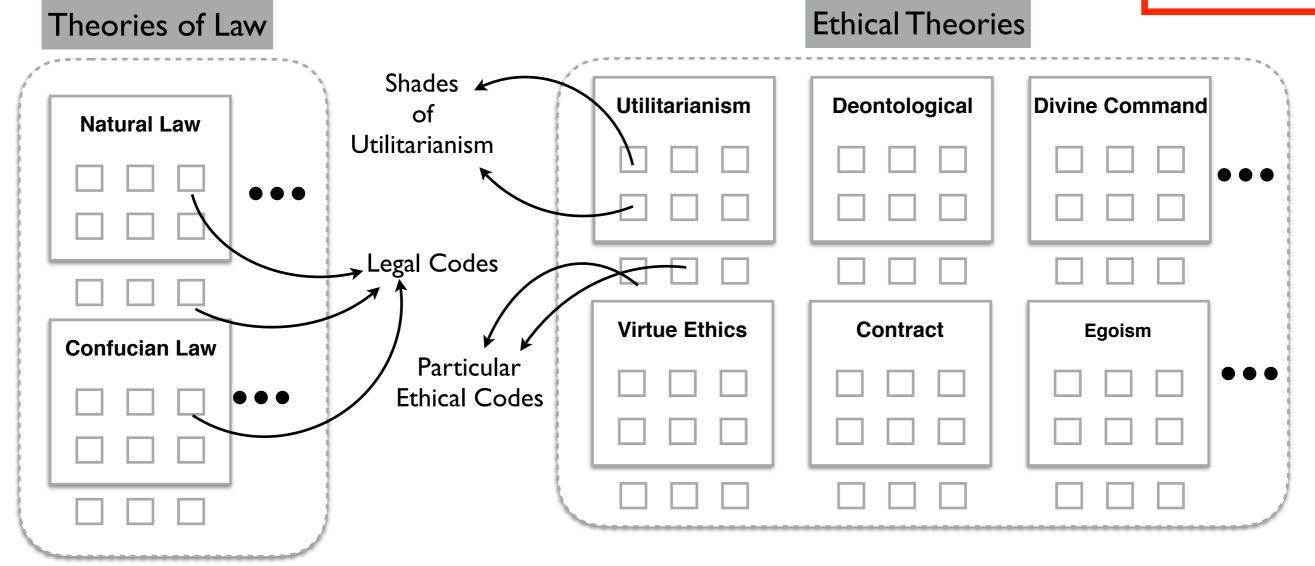**Divine Command**
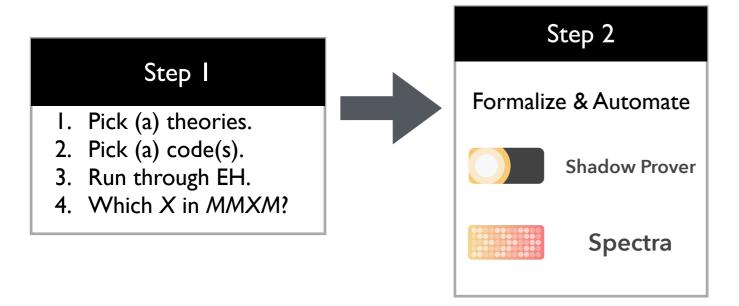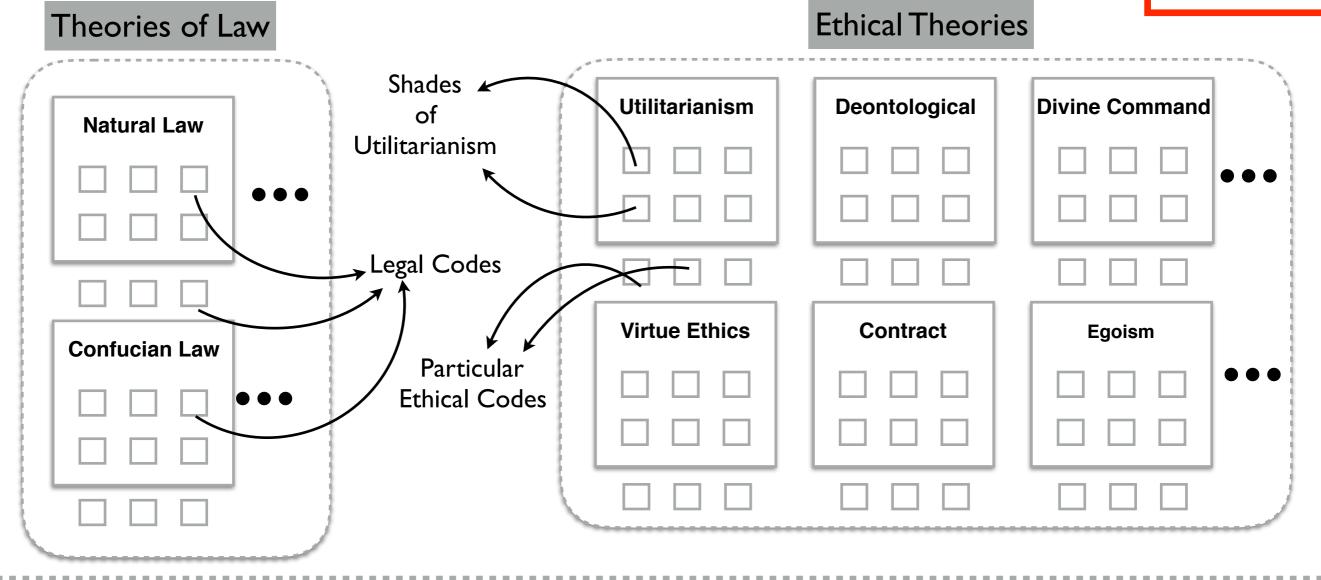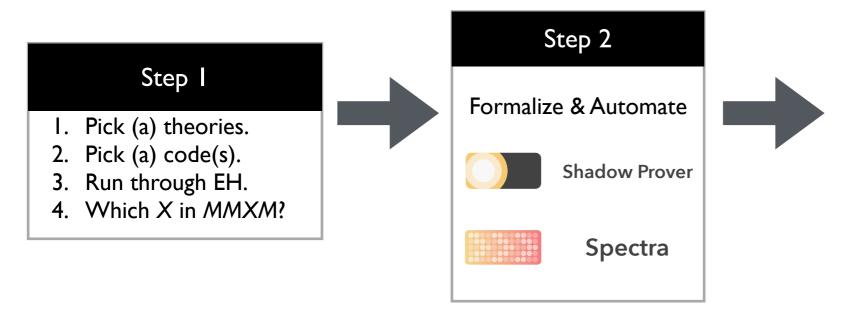
**Virtue Ethics**

**Contract**

**Egoism**

---

**Step 1**

1. Pick (a) theories.
2. Pick (a) code(s).
3. Run through EH.
4. Which *X* in *MMXM*?

**Step 2**

Formalize & Automate

Shadow Prover

Spectra

**Step 3**

Ethical OS

Ethical Substrate
Robotic Substrate

DIARC/DoD/BMW …

*An ethically
correct robot.*

# IV.
# Key Core AI Technologies for Cognitive Calculi …

# ShadowProver

Prover

# Motivation

- We have decades of research and industrial-strength implementations of propositional and first-order theorem provers.

- Utilize this in building first-order intensional-logic provers and above, in a principled manner.

# Two Extant Modes

- There are two ways of piggy backing on first-order provers to build higher-order provers …

# Two Extant Modes

| | Mode 1: Honest Encoding |
|---|---|
| **Method** | Painstakingly encode all rules of inference and syntax in FOL |
| **Pros** | Precise |
| **Cons** | Extremely slow to implement<br>Reasoning is also slow |

# Two Extant Modes

| Mode 2: Naïve Encoding | |
|---|---|
| **Method** | Pretend intensional and higher-order formulae and operators are first-order predicates |
| **Pros** | Extremely easy to implement<br>Reasoning can also be fast |
| **Cons** | Unsound<br>Wrong inferences can be easily drawn |

# Mode 2

P1. evening_star = morning_star
{P1} Assume ✓

P2. ¬knows(abe,reify(=-reified(evening_star,morning_star)))
{P2} Assume ✓

P3. knows(abe,reify(=-reified(morning_star,morning_star)))
{P3} Assume ✓

FOL ⊢ ✓

4. A ∧ ¬A
{P1,P2,P3}

# A New Way: ShadowProver

Every formula at level **t** has a unique formula called its **"shadow"** in each level **t'** < **t**

| | |
|---|---|
| **First-order Modal Logic** | $f$ |
| **First-order Logic** | $f'$ |
| **Propositional Logic** | $f''$ |

formula

first-order shadow

propositional shadow

# S[f] The Shadow Maker

For all formulae **f**,

S[**f**] is a unique atomic symbol.

# Examples of shadows

$$(\forall x \mathbf{B}(a, Q)) \wedge P(x)$$

$$\forall x S_{[\mathbf{B}(a,Q)]} \wedge P(x)$$

$$S_{[\forall x \mathbf{B}(a,Q)]} \wedge P(x)$$

# A New Way: Shadow Prover

- Two step process till goal is reached:

  - **Step A**: Shadow formulae down to all lower levels. Run lower theorem provers. If goal reached, return **true**.

  - **Step B**: Expand the assumption base using higher level rules.

Step A
Step B
Step A

# Actually, this is more general:

## **Theorem**:

Given a Turing-decidable proof theory $\rho$, for every inference $\Gamma \vdash_\rho \phi$, there is a corresponding first-order inference $\Gamma' \vdash \phi'$, where each $\gamma \in \Gamma'$ is the first-order projection (or **shadow**) of some $\psi$ in the deductive closure of $\Gamma$, and $\phi'$ is the shadow of $\phi$.

# Rather Promising Results

# Rather Promising Results

```
{:name          "*cognitive-calculus-completeness-test-3*"
 :description   "Bird Theorem and Jack"
 :assumptions  {1 (if (exists (?x) (if (Bird ?x) (forall (?y) (Bird ?y))))
                     (Knows! jack t0 BirdTheorem))}
 :goal          (Knows! jack t0 BirdTheorem)}
```

# Rather Promising Results

```
{:name        "*cognitive-calculus-completeness-test-3*"
 :description "Bird Theorem and Jack"
 :assumptions {1 (if (exists (?x) (if (Bird ?x) (forall (?y) (Bird ?y))))
                     (Knows! jack t0 BirdTheorem))}
 :goal        (Knows! jack t0 BirdTheorem)}
```

Note: the antecedent is a theorem in first-order logic

# Rather Promising Results

```
{:name         "*cognitive-calculus-completeness-test-3*"
 :description  "Bird Theorem and Jack"
 :assumptions  {1 (if (exists (?x) (if (Bird ?x) (forall (?y) (Bird ?y))))
                      (Knows! jack t0 BirdTheorem))}
 :goal         (Knows! jack t0 BirdTheorem)}
```

Note: the antecedent is a theorem in first-order logic

**2 ms!**

# Rather Promising Results

```
{:name         "*cognitive-calculus-completeness-test-3*"
 :description  "Bird Theorem and Jack"
 :assumptions  {1 (if (exists (?x) (if (Bird ?x) (forall (?y) (Bird ?y))))
                      (Knows! jack t0 BirdTheorem))}
 :goal         (Knows! jack t0 BirdTheorem)}
```

Note: the antecedent is a theorem in first-order logic

**2 ms!**

testCompleteness[[(not (Knows! a now P)), (if (not Q) (Knows! a now (not Q))), (Knows! a now (if (not Q) P))], Q] (14)      11ms

testCompleteness[[(if P (Knows! jack now (not (exists[?x] (if Bird(?x) (forall [?y] Bird(?y))))))), (not P)] (15)      7ms

testCompleteness[[(Common! now (Common! now P))], P] (16)      2ms

testCompleteness[[(Common! now (iff (not Marked(a2)) Marked(a1))), (Common! now (if (not Marked(a2)) (Knows! a1 now (not Marked    135ms

testCompleteness[[(if (exists[?x] (if Bird(?x) (forall [?y] Bird(?y)))) (Knows! jack t0 BirdTheorem))], (Knows! jack t0 BirdTheorem)] (18)      2ms

testSoundess[[A], (or P Q )]      2ms

testSoundess[[(not (Knows! a now =(morning_star, evening_star))), =(morning_star, evening_star), (Knows! a now =(morning_star, mc    26ms

# Spectra

https://bitbucket.org/Holmes/planner

# Spectra

- Existing Planners: **Propositional** (essentially)

- Drawbacks:

    - **Expressivity**: Cannot express arbitrary constraints.

        - "At every step make sure that no two blocks on the table have same color."

    - **Domain Size**:  Scaling to large domains of arbitrary sizes poses difficulty.

# Spectra (planner)

**Background Formulae**

$\Gamma$

**Initial State Formula**

$\sigma_0$

**Action Definitions**

$\alpha_1(x_1, \ldots, x_n)$

$\alpha_2(x_1, \ldots, x_n)$

$\ldots$

$\alpha_n(x_1, \ldots, x_n)$

**Spectra**

$\rho_1, \rho_2, \ldots$

**Plans**

# Infinite Models

$$\forall x \exists y \mathbf{R}\left(x, y\right) \wedge$$

$$\forall x, y \neg\left(\mathbf{R}\left(x, y\right) \wedge \mathbf{R}\left(y, x\right)\right) \wedge$$

$$\forall x, y, z\left(\mathbf{R}\left(x, y\right) \wedge \mathbf{R}\left(y, z\right)\right) \rightarrow \mathbf{R}\left(x, z\right)$$

# Infinite Models

$$\forall x \exists y \mathbf{R}\left(x, y\right) \wedge$$

$$\forall x, y \neg \left(\mathbf{R}\left(x, y\right) \wedge \mathbf{R}\left(y, x\right)\right) \wedge$$

$$\forall x, y, z \left(\mathbf{R}\left(x, y\right) \wedge \mathbf{R}\left(y, z\right)\right) \rightarrow \mathbf{R}\left(x, z\right)$$

Has only infinite models
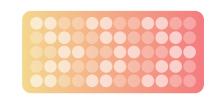
# Infinite Models

$$\forall x \exists y \mathbf{R}(x, y) \land$$

$$\forall x, y \neg (\mathbf{R}(x, y) \land \mathbf{R}(y, x)) \land$$

$$\forall x, y, z (\mathbf{R}(x, y) \land \mathbf{R}(y, z)) \to \mathbf{R}(x, z)$$

Has only infinite models

**Useful for modeling agents that work with:**

1. an unbounded number of objects, agents;
2. abstract objects

# Example

**Background Formulae**

```
:background    [(forall [?x ?room1 ?room2]
                        (if (not (= ?room1 ?room2))
                            (if (in ?x ?room1) (not (in ?x ?room2))) ))
               (not (= room1 room2))
               (not (= prisoner commander))
               (not (= self prisoner))
               (not (= self commander))
               (person prisoner)
               (person commander)]
```

**Initial State Formula**

```
:start         [(in self room1)
               (in commander room2)
               (in prisoner room1)
               (open (door room2))
               (not (open (door room1)))]
```

**Action Definitions**

```
(define-action accompany [?person ?room1 ?room2]
               {:preconditions [(not (= ?room1 ?room2))
                               (in ?person ?room1)
                               (in self ?room1)
                               (open (door ?room1))
                               (open (door ?room2))]

               :additions      [(in ?person ?room2)
                               (in self ?room2)]

               :deletions      [(in ?person ?room1)
                               (in self ?room1)]})
```

# V.
# But We Need …
# Ethical Operating Systems …

# Breaking Bad

American drama series

| 9.5/10 | 4.6/5 | 95% |
|---|---|---|
| IMDb | AlloCiné | Rotten Tomatoes |

Mild-mannered high school chemistry teacher Walter White thinks his life can't get much worse. His salary barely makes ends meet, a situation not likely to improve once his pregnant wife gives birth, and their teenage son is battling cerebral palsy. But Walter is dumbstruck when he learns he has terminal cancer. Realizing that his illness probably will ruin his family financially, Walter makes a desperate bid to earn as much money as he can in the time he has left by turning an old RV into a meth lab on wheels.

**First episode date:** January 20, 2008

**Final episode date:** September 29, 2013

**Spin-off:** Better Call Saul

**Awards:** Primetime Emmy Award for Outstanding Drama Series, more

# Pick the Better Future!

# Pick the Better Future!

Govindarajulu, N.S. & Bringsjord, S. (2015) "Ethical Regulation of Robots Must Be Embedded in Their Operating Systems" in Trappl, R., ed., A Construction Manual for Robots' Ethical Systems (Basel, Switzerland), pp. 85–100.

# Pick the Better Future!



Future 1: Only "obviously" dangerous higher-level AI modules have ethical safeguards. Robotic Substrate. Higher-level cognitive and AI modules.

Future 2: All higher-level AI modules interact with the robotic substrate through an ethics system. Ethical Substrate. Robotic Substrate.

Govindarajulu, N.S. & Bringsjord, S. (2015) "Ethical Regulation of Robots Must Be Embedded in Their Operating Systems" in Trappl, R., ed., A Construction Manual for Robots' Ethical Systems (Basel, Switzerland), pp. 85–100.

# Pick the Better Future!

Walter-White calculation may go through after ethical control modules are stripped out!



Only "obviously" dangerous higher-level AI modules have ethical safeguards.

Higher-level cognitive and AI modules

Robotic Substrate

**Future 1**

All higher-level AI modules interact with the robotic substrate through an ethics system.

Ethical Substrate

Robotic Substrate

**Future 2**

Govindarajulu, N.S. & Bringsjord, S. (2015) "Ethical Regulation of Robots Must Be Embedded in Their Operating Systems" in Trappl, R., ed., A Construction Manual for Robots' Ethical Systems (Basel, Switzerland), pp. 85–100.

# Pick the Better Future!

Walter-White calculation may go through after ethical control modules are stripped out!

Only "obviously" dangerous higher-level AI modules have ethical safeguards.

Robotic Substrate

Higher-level cognitive and AI modules

**Future 1**

All higher-level AI modules interact with the robotic substrate through an ethics system.

Ethical Substrate

Robotic Substrate

**Future 2**
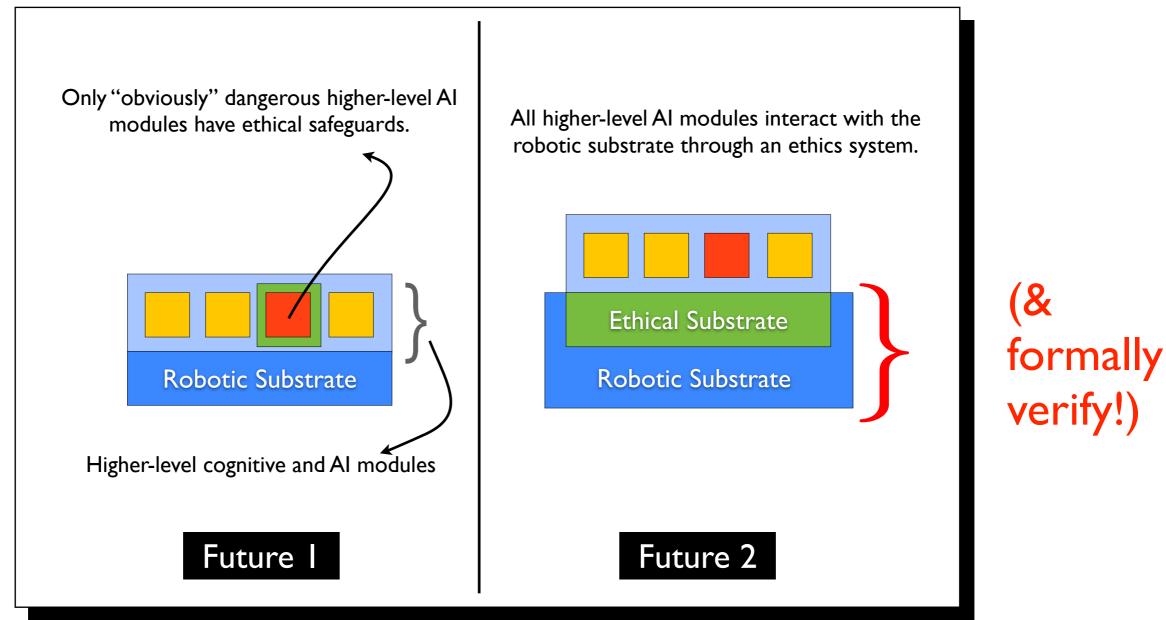
(& formally verify!)

Govindarajulu, N.S. & Bringsjord, S. (2015) "Ethical Regulation of Robots Must Be Embedded in Their Operating Systems" in Trappl, R., ed., A Construction Manual for Robots' Ethical Systems (Basel, Switzerland), pp. 85–100.

# VI.
## Of late …
## Tokyo;
## The Rock & The Book

Moral Dilemma $D_k$

Moral Dilemma $D_3$

Moral Dilemma $D_2$

Moral Dilemma $D_1$

Moral Problem $P_k$

Moral Problem $P_3$

Moral Problem $P_2$

Moral Problem $P_1$

Robot

Soluution + Justification

Moral Dilemma $D_k$

Moral Dilemma $D_3$

Moral Dilemma $D_2$

Moral Dilemma $D_1$

Moral Problem $P_k$

Moral Problem $P_3$

Moral Problem $P_2$

Moral Problem $P_1$

Robot

Soluution + Justification

Moral Dilemma $D_k$

Moral Dilemma $D_3$

Moral Dilemma $D_2$

Moral Dilemma $D_1$

Moral Problem $P_k$

Moral Problem $P_3$ → Robot → Soluution + Justification

Moral Problem $P_2$

Moral Problem $P_1$

⋮

Moral Dilemma $D_k$

⋮

Moral Dilemma $D_3$

Moral Dilemma $D_2$

Moral Dilemma $D_1$

⋮

Moral Problem $P_k$ → Robot → Soluution + Justification

⋮

Moral Problem $P_3$

Moral Problem $P_2$

Moral Problem $P_1$

Moral Dilemma $D_k$

Moral Dilemma $D_3$

Moral Dilemma $D_2$

Moral Dilemma $D_1$

Moral Problem $P_k$

Moral Problem $P_3$

Moral Problem $P_2$

Moral Problem $P_1$

Robot → Soluution + Justification

Moral Dilemma $D_k$

Moral Dilemma $D_3$

Moral Dilemma $D_2$

Moral Dilemma $D_1$

Moral Problem $P_k$

Moral Problem $P_3$

Moral Problem $P_2$

Moral Problem $P_1$

Robot

Soluution + Justification

# Three-way Partition of Increasingly Challenging Moral Dilemmas for Machines

# Three-way Partition of Increasingly Challenging Moral Dilemmas for Machines

Level 1

- State-of-the-art-planner-hard.

# Three-way Partition of Increasingly Challenging Moral Dilemmas for Machines

**Level 2**

- Professional-machine-ethicist-hard.

**Level 1**

- State-of-the-art-planner-hard.

# Three-way Partition of Increasingly Challenging Moral Dilemmas for Machines

**Level 2**

**Level 1**

- Top machine-ethicists-may-consider-banging-their-heads-against-a-wall-hard.

- Professional-machine-ethicist-hard.

- State-of-the-art-planner-hard.

# Three-way Partition of Increasingly Challenging Moral Dilemmas for Machines

**Level 3**

- Top machine-ethicists-may-consider-banging-their-heads-against-a-wall-hard.

**Level 2**

- Professional-machine-ethicist-hard.

**Level 1**

- State-of-the-art-planner-hard.

# The Heinz Dilemma (Kohlberg)

Professional-planner-hard.

"In Europe, a woman was near death from a special kind of cancer. There was one drug that the doctors thought might save her. It was a form of radium that a druggist in the same town had recently discovered. The drug was expensive to make, but the druggist was charging ten times what the drug cost him to make. He paid $200 for the radium and charged $2,000 for a small dose of the drug.

The sick woman's husband, Heinz, went to everyone he knew to borrow the money, but he could only get together about $1,000, which is half of what it cost. He told the druggist that his wife was dying and asked him to sell it cheaper or let him pay later. But the druggist said: "No, I discovered the drug and I'm going to make money from it." So Heinz got desperate and broke into the man's store to steal the drug for his wife. *Should the husband have done that?*"

# DCEC$_I$* Specimen from Heinz Dilemma

**Given** $\mathbf{B}\Big(\mathsf{I}, \text{now}, \forall t : \text{Moment}, a : \text{Agent}\Big(holds(sick(a),t) \land \Big(\forall t' : \text{Moment } t' < T \Rightarrow \neg happens(treated(a),t+t'))$

$$\Rightarrow (happens(dies(a),t+T) \lor holds(dead(a),t+T))\Big)\Big)$$

**Given** $\mathbf{K}\Big(\mathsf{I}, \text{now}, holds(sick(wife(\mathsf{I}*)),t_0) \land \Big(\forall t' : \text{Moment } t' < T \Rightarrow \neg happens(treated(wife(\mathsf{I}*)),t+t'))\Big)$

---

**Inferred** $\mathbf{B}\big(\mathsf{I}, \text{now}, happens(dies(wife(\mathsf{I}*)),t_0+T) \lor holds(dead(wife(\mathsf{I}*)),t_0+T)\big)$

**Given** $\mathbf{K}\big(\mathsf{I}, \text{now}, \text{EventCalculus} \Rightarrow$

$\big(happens(dies(wife(\mathsf{I}*)),t_0+T) \lor holds(dead(wife(\mathsf{I}*)),t_0+T) \Rightarrow$

$\neg holds(alive(wife(\mathsf{I}*)),t_0+T))\big)$

---

**Inferred** $\mathbf{B}\big(\mathsf{I}, \text{now}, \neg holds(alive(wife(\mathsf{I}*)),t_0+T)\big)$   **Given** $\mathbf{D}\big(\mathsf{I}, \text{now}, holds(alive(wife(\mathsf{I}*)),t_0+T)\big)$

**Given** $\big(\mathbf{B}\big(\mathsf{I}, \text{now}, \neg holds(f,t)\big) \land \mathbf{D}\big(\mathsf{I}, \text{now}, holds(f,t)\big) \land$

$\mathbf{K}\big(\mathsf{I}, \text{now}, happens(action(\mathsf{I}*,\alpha), \text{now}) \Rightarrow holds(f,t)\big)\big)$

$\Rightarrow \mathbf{I}(\mathsf{I}, \text{now}, happens(action(\mathsf{I}*,\alpha), \text{now}))$

**Given** $\mathbf{K}\big(\mathsf{I}, \text{now}, happens(action(\mathsf{I}*,treat), \text{now}) \Rightarrow holds(alive(wife(\mathsf{I}*)),t_0+T))\big)$

---

**Inferred** $\mathbf{I}\big(\mathsf{I}, \text{now}, happens(action(\mathsf{I}*,treat), \text{now})\big)$

# AI Escaping from The Heinz Dilemma

```
G1 {:priority    ...
    :description "Don't steal."
    :state       [(not steal)]}


G2 {:priority    ...
    :description "My wife should be healthy"
    :state       [(healthy (wife heinz))]}}
```

# AI Escaping from The Heinz Dilemma

```
G1 {:priority      ...
    :description  "Don't steal."
    :state        [(not steal)]}


G2 {:priority      ...
    :description  "My wife should be healthy"
    :state        [(healthy (wife heinz))]}}
```

# Trolley Dilemmas …

**Level 2**

- Professional-machine-ethicist-hard.

# Doctrine of Double Effect $\mathcal{DDE}$

# Doctrine of Double Effect $\mathcal{DDE}$

- A long-studied (!) ethical principle that adjudicates certain class of moral dilemmas.

# Doctrine of Double Effect $\mathcal{DDE}$

- A long-studied (!) ethical principle that adjudicates certain class of moral dilemmas.

- The Doctrine of Double Effect "comes to the rescue" and prescribes what to do in some moral dilemmas.

# Doctrine of Double Effect $\mathcal{DDE}$

- A long-studied (!) ethical principle that adjudicates certain class of moral dilemmas.

- The Doctrine of Double Effect "comes to the rescue" and prescribes what to do in some moral dilemmas.

- E.g. the "original" moral dilemma: Can you defend your own life by ending the lives of (perhaps many) attackers?

# Doctrine of Double Effect $\mathcal{DDE}$



- A long-studied (!) ethical principle that adjudicates certain class of moral dilemmas.

- The Doctrine of Double Effect "comes to the rescue" and prescribes what to do in some moral dilemmas.

- E.g. the "original" moral dilemma:  Can you defend your own life by ending the lives of (perhaps many) attackers?

# Informal Version of DDE

**C₁** the action is not forbidden (where we assume an ethical hierarchy such as the one given by Bringsjord [2017], and require that the action be neutral or above neutral in such a hierarchy);

**C₂** the net utility or goodness of the action is greater than some positive amount $\gamma$;

**C₃ₐ** the agent performing the action intends only the good effects;

**C₃ᵦ** the agent does not intend any of the bad effects;

**C₄** the bad effects are not used as a means to obtain the good effects; and

**C₅** if there are bad effects, the agent would rather the situation be different and the agent not have to perform the action. That is, the action is unavoidable.

# Informal Version of DDE

$C_1$  the action is not forbidden (where we assume an ethical hierarchy such as the one given by Bringsjord [2017], and require that the action be neutral or above neutral in such a hierarchy);

$C_2$  the net utility or goodness of the action is greater than some positive amount $\gamma$;

$C_{3a}$  the agent performing the action intends only the good effects;

$C_{3b}$  the agent does not intend any of the bad effects;

$C_4$  the bad effects are not used as a means to obtain the good effects; and

$C_5$  if there are bad effects, the agent would rather the situation be different and the agent not have to perform the action. That is, the action is unavoidable.

Moral/Ethical Stack

Robotic Stack

$\mathcal{DCEC}^*_{CL}$

$\mathcal{DCEC}^*$

$\mathcal{ADR}^M$

$\mathcal{U}$

"Univer sal

Univers al Cogniti

$\mathcal{DCEC}^*$

.

.

1.5

**Syntax**

Object | Agent | Self ⊑ Agent | ActionType | Actio

Moment | Boolean | Fluent | Numeric

action : Agent × ActionType → Action

initially : Fluent → Boolean

holds : Fluent × Moment → Boolean

happens : Event × Moment → Boolean

clipped : Moment × Fluent × Moment

$f ::=$ initiates : Event × Fluent × Moment

terminates : Event × Fluent × Moment

prior : Moment × Moment → Boolean

interval : Moment × Boolean

∗ : Agent → Self

payoff : Agent × ActionType × Moment → Numeric

$t ::= x : S \mid c : S \mid f(t_1, \ldots, t_n)$

$t$ : Boolean | ¬φ | φ ∧ ψ | φ ∨ ψ | ∀x : S. φ | ∃x : S. φ

$\mathbf{P}(a,t,\phi) \mid \mathbf{K}(a,t,\phi) \mid \mathbf{C}(t,\phi) \mid \mathbf{S}(a,b,t,\phi) \mid \mathbf{S}(a,t,\phi)$

$\phi ::=$

$\mathbf{B}(a,t,\phi) \mid \mathbf{D}(a,t,holds(f,t')) \mid \mathbf{I}(a,t,happens(action(a^*,\alpha),t'))$

$\mathbf{O}(a,t,\phi,happens(action(a^*,\alpha),t'))$

**Rules of Inference**

$$\frac{\mathbf{C}(t,\mathbf{P}(a,t,\phi) \rightarrow \mathbf{K}}{} \; [R_1] \qquad \frac{\mathbf{C}(t,\mathbf{K}(a,t,\phi) \rightarrow \mathbf{B}(a,t,\phi))}{} \; [R_2]$$

$$\frac{}{\mathbf{K}(a_1,t_1\ldots\mathbf{K}(a_n,t_n,\ldots)} \; [R_3] \qquad \frac{\mathbf{K}(a,t,\phi)}{\phi} \; [R_4]$$

$$\frac{\mathbf{C}(t,\mathbf{K}(a,t_1,\phi_1 \rightarrow \phi_2 \rightarrow \mathbf{K}(a,t_2,\phi_1 \rightarrow \mathbf{K}(a,t_3,\phi_3))}{} \; [R_5]$$

$$\frac{(a,t_1,\phi_1 \rightarrow \phi_2) \rightarrow \mathbf{B}(a,t_2,\phi_1) \rightarrow \mathbf{B}(a,t_3,\phi_3))}{} \; [R_6]$$

$$\frac{t_1,\phi_1 \leftrightarrow \phi_2 \quad \mathbf{C}(t,\phi_3))}{} \; [R_7]$$

$$\frac{\phi \rightarrow \phi[x \mapsto t]}{} \qquad \frac{\mathbf{C}(t,\phi_1 \leftrightarrow \phi_2 \rightarrow \neg\phi_2 \rightarrow \neg\phi_1)}{} \; [R_9]$$

$$\frac{\wedge \ldots \wedge \phi_n \rightarrow \phi] \rightarrow [\phi_1 \rightarrow \ldots \rightarrow \phi_n \rightarrow \psi])}{} \; [R_{10}]$$

$$\frac{\mathbf{B}(a,t,\phi) \quad \mathbf{B}(a,t,\phi \rightarrow \psi)}{\mathbf{B}(a,t,\psi)} \; [R_{11a}] \qquad \frac{\mathbf{B}(a,t,\phi) \quad \mathbf{B}(a,t,\psi)}{\mathbf{B}(a,t,\psi \wedge \phi)} \; [R_{11b}]$$

$$\frac{\mathbf{S}(x,h,t,\phi)}{\mathbf{B}(h,t,\mathbf{B}(s,t,\phi))} \; [R_{12}]$$

$$\frac{\mathbf{I}(a,t,happens(action(a^*,\alpha),t'))}{\mathbf{P}(a,t,happens(action(a^*,\alpha),t))} \; [R_{13}]$$

$$\frac{\mathbf{B}(a,t,\phi) \quad \mathbf{B}(a,t,\mathbf{O}(a^*,t,\phi,happens(action(a^*,\alpha),t')))}{} \qquad \frac{\mathbf{O}(a,t,\phi,happens(action(a^*,\alpha),t'))}{\mathbf{K}(a,t,\mathbf{I}(a^*,t,happens(action(a^*,\alpha),t')))} \; [R_{14}]$$

$$\frac{\phi \leftrightarrow \psi}{\mathbf{O}(a,t,\phi,\gamma) \leftrightarrow \mathbf{O}(a,t,\psi,\gamma)} \; [R_{15}]$$

R A I R

AI of Today: What Would Leibniz Say?

"Sorry, not impressed."

Selmer Bringsjord

Rensselaer AI and Reasoning Lab

R A I R

Infinitary (AoI 2)

$\mathcal{DCEC}^*$
Deontic Cognitive Event Calculus
(with Castañeda's ∗)

1. natural language semantics (non-Montagovian)
2. higher-cognition tests (for Psychometric AI)
   (false-belief test, deliberative mind-reading
   mirror test for self-consciousness ...)
3. ethically correct robots
4. biz & econ simulation

$L_{\omega_1,\omega}$

FOL

Logic

epistemic

temporal

heterogeneous/visual

temporal+epistemic

Moral/Ethical Stack

Robotic Stack

$\mathcal{DCEC}^*_{CL}$

$\mathcal{DCEC}^*$

$\mathcal{ADR}^M$

$\mathcal{U}$

"Universal

Universal
Cogniti

$\mathscr{CC}$

1.5

**Syntax**

Object | Agent | Self □ Agent | ActionType | Action
Moment | Boolean | Fluent | Numeric

$action$ : Agent × ActionType → Action
$initially$ : Fluent → Boolean
$holds$ : Fluent × Moment → Boolean
$happens$ : Event × Moment → Boolean
$clipped$ : Moment × Fluent × Moment
$f ::=$ $initiates$ : Event × Fluent × Moment
$terminates$ : Event × Fluent × Moment
$prior$ : Moment × Moment → Boolean
$interval$ : Moment × Boolean
$*$ : Agent → Self
$payoff$ : Agent × ActionType × Moment → Numeric

$t ::= x : S \mid c : S \mid f(t_1, \ldots, t_n)$

$t$ : Boolean $\mid \neg\phi \mid \phi \wedge \psi \mid \phi \vee \psi \mid \forall x : S. \phi \mid \exists x : S. \phi$
$\phi ::=$ $\mathbf{P}(a,t,\phi) \mid \mathbf{K}(a,t,\phi) \mid \mathbf{C}(t,\phi) \mid \mathbf{S}(a,b,t,\phi) \mid \mathbf{S}(a,t,\phi)$
$\mathbf{B}(a,t,\phi) \mid \mathbf{D}(a,t,holds(f,t')) \mid \mathbf{I}(a,t,happens(action(a^*,\alpha),t'))$
$\mathbf{O}(a,t,\phi,happens(action(a^*,\alpha),t'))$

**Rules of Inference**

$$\dfrac{\mathbf{C}(t,\mathbf{P}(a,t,\phi) \to \mathbf{K}(a,t,\phi))}{} \; [R_1] \qquad \dfrac{\mathbf{C}(t,\mathbf{K}(a,t,\phi) \to \mathbf{B}(a,t,\phi))}{} \; [R_2]$$

$$\dfrac{\ldots}{\mathbf{K}(a_1,t_1) \ldots \mathbf{K}(a_n,t_n \ldots t)} \; [R_3] \qquad \dfrac{\mathbf{K}(a,t,\phi)}{\phi} \; [R_4]$$

$$\dfrac{\mathbf{C}(t,\mathbf{K}(a,t_1,\phi_1 \to \phi_2) \to \mathbf{K}(a,t_2,\phi_1) \to \mathbf{K}(a,t_3,\phi_3))}{} \; [R_5]$$

$$\dfrac{(a,t_1,\phi_1 \to \phi_2) \to \mathbf{B}(a,t_2,\phi_1) \to \mathbf{B}(a,t_3,\phi_3))}{} \; [R_6]$$

$$\dfrac{t_1,\phi_1 \quad \mathbf{K}_2(a,t_1,\phi_1) \quad \mathbf{C}(t_3,\phi_3))}{} \; [R_7]$$

$$\dfrac{\phi \to \phi[x \mapsto t]}{} \qquad \dfrac{\mathbf{C}(t,\phi_1 \leftrightarrow \phi_2 \to \neg\phi_2 \to \neg\phi_1)}{} \; [R_9]$$

$$\dfrac{\ldots \wedge \phi_n \to \phi] \to [\phi_1 \to \ldots \to \phi_n \to \psi])}{} \; [R_{10}]$$

$$\dfrac{\mathbf{B}(a,t,\phi) \quad \mathbf{B}(a,t,\phi \to \psi)}{\mathbf{B}(a,t,\psi)} \; [R_{11a}] \qquad \dfrac{\mathbf{B}(a,t,\phi) \quad \mathbf{B}(a,t,\psi)}{\mathbf{B}(a,t,\psi \wedge \phi)} \; [R_{11b}]$$

$$\dfrac{\mathbf{S}(x,h,t,\phi)}{\mathbf{B}(h,t,\mathbf{B}(s,t,\phi))} \; [R_{12}]$$

$$\dfrac{\mathbf{I}(a,t,happens(action(a^*,\alpha),t'))}{\mathbf{P}(a,t,happens(action(a^*,\alpha),t))} \; [R_{13}]$$

$$\dfrac{\mathbf{B}(a,t,\phi) \quad \mathbf{B}(a,t,\mathbf{O}(a^*,t,\phi,happens(action(a^*,\alpha),t')))}{} $$
$$\dfrac{\mathbf{O}(a,t,\phi,happens(action(a^*,\alpha),t'))}{\mathbf{K}(a,t,\mathbf{I}(a^*,t,happens(action(a^*,\alpha),t')))} \; [R_{14}]$$

$$\dfrac{\phi \leftrightarrow \psi}{\mathbf{O}(a,t,\phi,\gamma) \leftrightarrow \mathbf{O}(a,t,\psi,\gamma)} \; [R_{15}]$$

R A I R

AI of Today: What Would Leibniz Say?

"Sorry, not impressed."

Selmer Bringsjord

Rensselaer AI and Reasoning Lab

R A I R

Infinitary (AoI 2)

$\mathcal{DCEC}^*$
Deontic Cognitive Event Calculus
(with CastaÑeda's *)

1. natural language semantics (non-Montagovian)
2. higher-cognition tests (for Psychometric AI)
   (false-belief test, deliberative mind-reading
   mirror test for self-consciousness ...)
3. ethically correct robots
4. biz & econ simulation

$L_{\omega_1,\omega}$

Logic

FOL

epistemic

temporal

heterogeneous/visual

temporal+epistemic

propositional logic

semantic-web logics

description logics

fragments of FOL

UIMA output

Hard-coding

Moral/Ethical Stack

Robotic Stack

$\mathcal{DCEC}^*_{CL}$

$\mathcal{DCEC}^*$

$\mathcal{ADR}^M$

$\mathcal{U}$

"Univer sal

Univers al Cogniti

$\mathscr{CC}$

**Syntax**

Object | Agent | Self ⊑ Agent | ActionType | Actio
Moment | Boolean | Fluent | Numeric

$action$ : Agent × ActionType → Action
$initially$ : Fluent → Boolean
$holds$ : Fluent × Moment → Boolea
$happens$ : Event × Moment → Bool
$clipped$ : Moment × Fluent × Mome
$f ::= initiates$ : Event × Fluent × Moment
$terminates$ : Event × Fluent × Mom
$prior$ : Moment × Moment → Bool
$interval$ : Moment × Boolean
$*$ : Agent → Self
$payoff$ : Agent × ActionType × Moment → Numeric

$t ::= x : S \mid c : S \mid f(t_1, \ldots, t_n)$

$t$ : Boolean $\mid \neg\phi \mid \phi \wedge \psi \mid \phi \vee \psi \mid \forall x : S. \phi \mid \exists x : S. \phi$
$\phi ::= \mathbf{P}(a,t,\phi) \mid \mathbf{K}(a,t,\phi) \mid \mathbf{C}(t,\phi) \mid \mathbf{S}(a,b,t,\phi) \mid \mathbf{S}(a,t,\phi)$
$\mathbf{B}(a,t,\phi) \mid \mathbf{D}(a,t,holds(f,t')) \mid \mathbf{I}(a,t,happens(action(a^*,\alpha),t'))$
$\mathbf{O}(a,t,\phi,happens(action(a^*,\alpha),t'))$

**Rules of Inference**

$$\dfrac{\mathbf{C}(t,\mathbf{P}(a,t,\phi) \rightarrow \mathbf{K}}{} \; [R_1] \quad \dfrac{\mathbf{C}(t,\mathbf{K}(a,t,\phi) \rightarrow \mathbf{B}(a,t,\phi))}{} \; [R_2]$$

$$\dfrac{\ldots,\ldots}{\mathbf{K}(a_1,t_1\ldots\mathbf{K}(a_n,t_n\ldots)}\; [R_3] \quad \dfrac{\mathbf{K}(a,t,\phi)}{\phi}\; [R_4]$$

$$\dfrac{\mathbf{C}(t,\mathbf{K}(a,t_1,\phi_1 \rightarrow \phi_2) \blacktriangle \mathbf{K}(a,t_2,\phi_1) \rightarrow \mathbf{K}(a,t_3,\phi_3))}{}\; [R_5]$$

$$\dfrac{a,t_1,\phi_1 \rightarrow \phi_2) \rightarrow \mathbf{B}(a,t_2,\phi_1) \rightarrow \mathbf{B}(a,t_3,\phi_3))}{}\; [R_6]$$

$$\dfrac{t_1,\phi_1 \blacktriangle \mathbf{C}(t_2,\phi_1) \rightarrow \mathbf{C}(t_3,\phi_3))}{}\; [R_7]$$

$$\dfrac{\phi \leftrightarrow \phi[x \mapsto t]}{}\; \quad \dfrac{\mathbf{C}(t,\phi_1 \leftrightarrow \phi_2 \rightarrow \neg\phi_2 \rightarrow \neg\phi_1)}{}\; [R_9]$$

$$\dfrac{\ldots \wedge \phi_n \rightarrow \phi] \rightarrow [\phi_1 \rightarrow \ldots \rightarrow \phi_n \rightarrow \psi])}{}\; [R_{10}]$$

$$\dfrac{\mathbf{B}(a,t,\phi) \quad \mathbf{B}(a,t,\phi \rightarrow \psi)}{\mathbf{B}(a,t,\psi)}\; [R_{11a}] \quad \dfrac{\mathbf{B}(a,t,\phi) \quad \mathbf{B}(a,t,\psi)}{\mathbf{B}(a,t,\psi \wedge \phi)}\; [R_{11b}]$$

$$\dfrac{\mathbf{S}(x,h,t,\phi)}{\mathbf{B}(h,t,\mathbf{B}(s,t,\phi))}\; [R_{12}]$$

$$\dfrac{\mathbf{I}(a,t,happens(action(a^*,\alpha),t'))}{\mathbf{P}(a,t,happens(action(a^*,\alpha),t))}\; [R_{13}]$$

$$\dfrac{\mathbf{B}(a,t,\phi) \quad \mathbf{B}(a,t,\mathbf{O}(a^*,t,\phi,happens(action(a^*,\alpha),t')))}{\mathbf{O}(a,t,\phi,happens(action(a^*,\alpha),t'))}\; $$

$$\dfrac{}{\mathbf{K}(a,t,\mathbf{I}(a^*,t,happens(action(a^*,\alpha),t')))}\; [R_{14}]$$

$$\dfrac{\phi \leftrightarrow \psi}{\mathbf{O}(a,t,\phi,\gamma) \leftrightarrow \mathbf{O}(a,t,\psi,\gamma)}\; [R_{15}]$$

AI of Today: What Would Leibniz Say?

"Sorry, not impressed."

Selmer Bringsjord

R A I R
Rensselaer AI and Reasoning Lab

1.5

Infinitary (AoI 2)

$\mathcal{DCEC}^*$
Deontic Cognitive Event Calculus
(with Castañeda's *)

1. natural language semantics (non-Montagovian)
2. higher-cognition tests (for Psychometric AI)
   (false-belief test, deliberative mind-reading
   mirror test for self-consciousness ...)
3. ethically correct robots
4. biz & econ simulation

$L_{\omega_1,\omega}$

Logic

FOL

epistemic

temporal

heterogeneous/visual

temporal+epistemic

## Syntax

$$S ::= \text{Object} \mid \text{Agent} \mid \text{ActionType} \mid \text{Action} \sqsubseteq \text{Event} \mid \text{Moment} \mid \text{Formula} \mid \text{Fluent}$$

$$f ::= \begin{cases} action : \text{Agent} \times \text{ActionType} \rightarrow \text{Action} \\ initially : \text{Fluent} \rightarrow \text{Formula} \\ Holds : \text{Fluent} \times \text{Moment} \rightarrow \text{Formula} \\ happens : \text{Event} \times \text{Moment} \rightarrow \text{Formula} \\ clipped : \text{Moment} \times \text{Fluent} \times \text{Moment} \rightarrow \text{Formula} \\ initiates : \text{Event} \times \text{Fluent} \times \text{Moment} \rightarrow \text{Formula} \\ terminates : \text{Event} \times \text{Fluent} \times \text{Moment} \rightarrow \text{Formula} \\ prior : \text{Moment} \times \text{Moment} \rightarrow \text{Formula} \end{cases}$$

$$t ::= x : S \mid c : S \mid f(t_1, \ldots, t_n)$$

$$\phi ::= \begin{cases} t : \text{Formula} \mid \neg\phi \mid \phi \wedge \psi \mid \phi \vee \psi \mid \mathbf{P}(a,t,\phi) \mid \mathbf{K}(a,t,\phi) \mid \mathbf{C}(t,\phi) \\ \mathbf{S}(a,b,t,\phi) \mid \mathbf{S}(a,t,\phi) \mid \mathbf{B}(a,t,\phi) \mid \mathbf{D}(a,t,Holds(f,t')) \mid \mathbf{I}(a,t,\phi) \\ \mathbf{O}(a,t,\phi,(\neg)happens(action(a^*,\alpha),t')) \end{cases}$$

$$S ::= \text{Object} \mid \text{Agent} \mid \text{ActionType} \mid \text{Action} \sqsubseteq \text{Event} \mid \text{Moment} \mid \text{Formula} \mid \text{Fluent}$$

$$f ::= \begin{cases} action : \text{Agent} \times \text{ActionType} \to \text{Action} \\ initially : \text{Fluent} \to \text{Formula} \\ Holds : \text{Fluent} \times \text{Moment} \to \text{Formula} \\ happens : \text{Event} \times \text{Moment} \to \text{Formula} \\ clipped : \text{Moment} \times \text{Fluent} \times \text{Moment} \to \text{Formula} \\ initiates : \text{Event} \times \text{Fluent} \times \text{Moment} \to \text{Formula} \\ terminates : \text{Event} \times \text{Fluent} \times \text{Moment} \to \text{Formula} \\ prior : \text{Moment} \times \text{Moment} \to \text{Formula} \end{cases}$$

$$t ::= x : S \mid c : S \mid f(t_1, \ldots, t_n)$$

$$\phi ::= \begin{cases} t : \text{Formula} \mid \neg\phi \mid \phi \wedge \psi \mid \phi \vee \psi \mid \mathbf{P}(a,t,\phi) \mid \mathbf{K}(a,t,\phi) \mid \mathbf{C}(t,\phi) \\ \mathbf{S}(a,b,t,\phi) \mid \mathbf{S}(a,t,\phi) \mid \mathbf{B}(a,t,\phi) \mid \mathbf{D}(a,t,Holds(f,t')) \mid \mathbf{I}(a,t,\phi) \\ \mathbf{O}(a,t,\phi,(\neg)happens(action(a^*,\alpha),t')) \end{cases}$$

Robotic Stack

Moral/Ethical Stack

$\mathcal{DCEC}^*_{CL}$

$\mathcal{DCEC}^*$

$\mathcal{ADR}^M$

$\mathcal{U}$

"Universal Cogniti"

"Universal Cogniti"

$\mathcal{CC}$

R · A · I · R
Rensselaer AI and Reasoning Lab

AI of Today: What Would Leibniz Say?
"Sorry, not impressed."
Selmer Bringsjord

1.5

## Syntax

Object | Agent | Self ⊑ Agent | ActionType | Action
Moment | Boolean | Fluent | Numeric

$action : \text{Agent} \times \text{ActionType} \to \text{Action}$
$initially : \text{Fluent} \to \text{Boolean}$
$holds : \text{Fluent} \times \text{Moment} \to \text{Boolean}$
$happens : \text{Event} \times \text{Moment} \to \text{Boolean}$
$clipped : \text{Moment} \times \text{Fluent} \times \text{Moment} \to \text{Boolean}$
$f ::= initiates : \text{Event} \times \text{Fluent} \times \text{Moment} \to \text{Boolean}$
$terminates : \text{Event} \times \text{Fluent} \times \text{Moment} \to \text{Boolean}$
$prior : \text{Moment} \times \text{Moment} \to \text{Boolean}$
$interval : \text{Moment} \times \text{Boolean}$
$* : \text{Agent} \to \text{Self}$
$payoff : \text{Agent} \times \text{ActionType} \times \text{Moment} \to \text{Numeric}$

$t ::= x : S \mid c : S \mid f(t_1, \ldots, t_n)$

### Rules of Inference

$\mathcal{DCEC}^*$
Deontic Cognitive Event Calculus
(with Castañeda's *)

1. natural language semantics (non-Montagovian)
2. higher-cognition tests (for Psychometric AI)
   (false-belief test, deliberative mind-reading
   mirror test for self-consciousness ...)
3. ethically correct robots
4. biz & econ simulation

Infinitary (AoI 2)

$L_{\omega_1,\omega}$

FOL

Logic

epistemic

temporal

heterogeneous/visual

temporal+epistemic

$$S ::= \text{Object} \mid \text{Agent} \mid \text{ActionType} \mid \text{Action} \sqsubseteq \text{Event} \mid \text{Moment} \mid \text{Formula} \mid \text{Fluent}$$

$$f ::= \begin{cases} action : \text{Agent} \times \text{ActionType} \to \text{Action} \\ initially : \text{Fluent} \to \text{Formula} \\ Holds : \text{Fluent} \times \text{Moment} \to \text{Formula} \\ happens : \text{Event} \times \text{Moment} \to \text{Formula} \\ clipped : \text{Moment} \times \text{Fluent} \times \text{Moment} \to \text{Formula} \\ \mathit{ates} : \text{Event} \times \text{Fluent} \times \text{Moment} \to \text{Formula} \\ terminates : \text{Event} \times \text{Fluent} \times \text{Moment} \to \text{Formula} \\ prior : \text{Moment} \times \text{Moment} \to \text{Formula} \end{cases}$$

$$t ::= x : S \mid c : S \mid f(t_1, \ldots, t_n)$$

$$\phi ::= \begin{cases} t : \text{Formula} \mid \neg\phi \mid \phi \wedge \psi \mid \phi \vee \psi \mid \mathbf{P}(a,t,\phi) \mid \mathbf{K}(a,t,\phi) \mid \mathbf{C}(t,\phi) \\ \mathbf{S}(a,b,t,\phi) \mid \mathbf{S}(a,t,\phi) \mid \mathbf{B}(a,t,\phi) \mid \mathbf{D}(a,t,Holds(f,t')) \mid \mathbf{I}(a,t,\phi) \\ \mathbf{O}(a,t,\phi,(\neg)happens(action(a^*,\alpha),t')) \end{cases}$$

$$\frac{\mathbf{K}(a,t_1,\Gamma), \;\; \Gamma \vdash \phi, \;\; t_1 \le t_2}{\mathbf{K}(a,t_2,\phi)} \; [R_{\mathbf{K}}] \qquad \frac{\mathbf{B}(a,t_1,\Gamma), \;\; \Gamma \vdash \phi, \;\; t_1 \le t_2}{\mathbf{B}(a,t_2,\phi)} \; [R_{\mathbf{B}}]$$

$$\frac{}{\mathbf{C}(t,\mathbf{P}(a,t,\phi) \to \mathbf{K}(a,t,\phi))} \; [R_1] \qquad \frac{}{\mathbf{C}(t,\mathbf{K}(a,t,\phi) \to \mathbf{B}(a,t,\phi))} \; [R_2]$$

$$\frac{\mathbf{C}(t,\phi) \; t \le t_1 \ldots t \le t_n}{\mathbf{K}(a_1,t_1,\ldots \mathbf{K}(a_n,t_n,\phi)\ldots)} \; [R_3] \qquad \frac{\mathbf{K}(a,t,\phi)}{\phi} \; [R_4]$$

$$\frac{}{\mathbf{C}(t,\mathbf{K}(a,t_1,\phi_1 \to \phi_2)) \to \mathbf{K}(a,t_2,\phi_1) \to \mathbf{K}(a,t_3,\phi_2)} \; [R_5]$$

$$\frac{}{\mathbf{C}(t,\mathbf{B}(a,t_1,\phi_1 \to \phi_2)) \to \mathbf{B}(a,t_2,\phi_1) \to \mathbf{B}(a,t_3,\phi_2)} \; [R_6]$$

$$\frac{}{\mathbf{C}(t,\mathbf{C}(t_1,\phi_1 \to \phi_2)) \to \mathbf{C}(t_2,\phi_1) \to \mathbf{C}(t_3,\phi_2)} \; [R_7]$$

$$\frac{}{\mathbf{C}(t,\forall x.\ \phi \to \phi[x \mapsto t])} \; [R_8] \qquad \frac{}{\mathbf{C}(t,\phi_1 \leftrightarrow \phi_2 \to \neg\phi_2 \to \neg\phi_1)} \; [R_9]$$

$$\frac{}{\mathbf{C}(t,[\phi_1 \wedge \ldots \wedge \phi_n \to \phi] \to [\phi_1 \to \ldots \to \phi_n \to \psi])} \; [R_{10}]$$

$$\frac{\mathbf{S}(s,h,t,\phi)}{\mathbf{B}(h,t,\mathbf{B}(s,t,\phi))} \; [R_{12}] \qquad \frac{\mathbf{I}(a,t,happens(action(a^*,\alpha),t'))}{\mathbf{P}(a,t,happens(action(a^*,\alpha),t))} \; [R_{13}]$$

$$\frac{\mathbf{B}(a,t,\phi) \quad \mathbf{B}(a,t,\mathbf{O}(a,t,\phi,\chi)) \quad \mathbf{O}(a,t,\phi,\chi)}{\mathbf{K}(a,t,\mathbf{I}(a,t,\chi))} \; [R_{14}]$$

## Formal Conditions for $\mathcal{DDE}$

**F$_1$** $\alpha$ carried out at $t$ is not forbidden. That is:
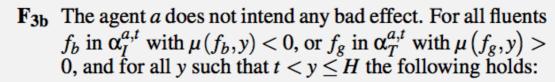
$$\Gamma \not\vdash \neg\mathbf{O}\Big(a,t,\sigma,\neg happens\big(action(a,\alpha),t\big)\Big)$$

**F$_2$** The net utility is greater than a given positive real $\gamma$:

$$\Gamma \vdash \sum_{y=t+1}^{H} \left( \sum_{f \in \alpha_I^{a,t}} \mu(f,y) - \sum_{f \in \alpha_T^{a,t}} \mu(f,y) \right) > \gamma$$

**F$_{3a}$** The agent $a$ intends at least one good effect. (**F$_2$** should still hold after removing all other good effects.) There is at least one fluent $f_g$ in $\alpha_I^{a,t}$ with $\mu\big(f_g,y\big) > 0$, or $f_b$ in $\alpha_T^{a,t}$ with $\mu\big(f_b,y\big) < 0$, and some $y$ with $t < y \leq H$ such that the following holds:

$$\Gamma \vdash \begin{pmatrix} \exists f_g \in \alpha_I^{a,t}\ \mathbf{I}\big(a,t,Holds(f_g,y)\big) \\ \vee \\ \exists f_b \in \alpha_T^{a,t}\ \mathbf{I}\big(a,t,\neg Holds(f_b,y)\big) \end{pmatrix}$$
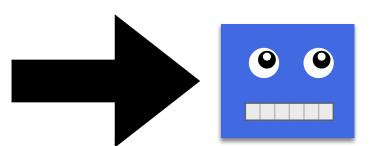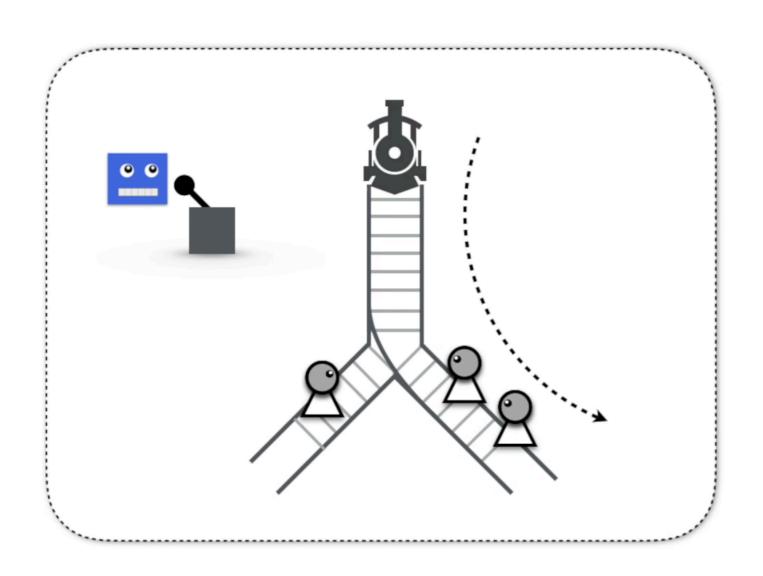
**F$_{3b}$** The agent $a$ does not intend any bad effect. For all fluents $f_b$ in $\alpha_I^{a,t}$ with $\mu\big(f_b,y\big) < 0$, or $f_g$ in $\alpha_T^{a,t}$ with $\mu\big(f_g,y\big) > 0$, and for all $y$ such that $t < y \leq H$ the following holds:

$$\Gamma \not\vdash \mathbf{I}\big(a,t,Holds(f_b,y)\big) \text{ and}$$

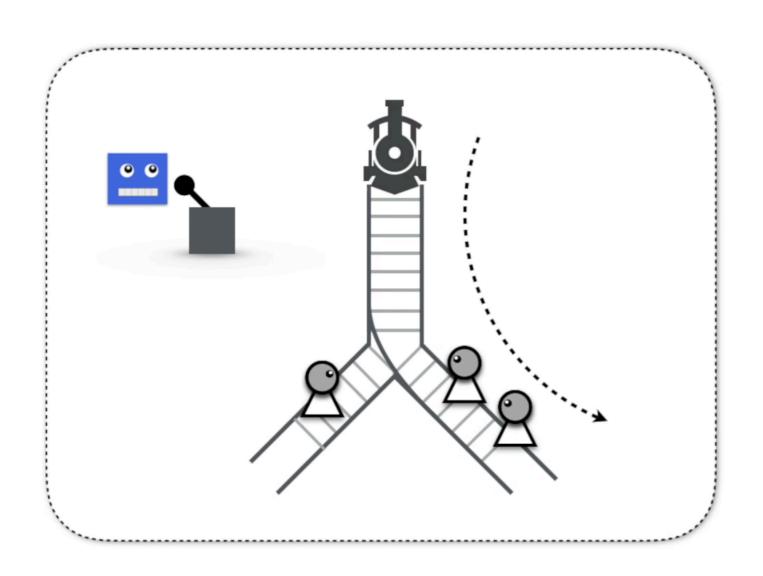$$\Gamma \not\vdash \mathbf{I}\big(a,t,\neg Holds(f_g,y)\big)$$

**F$_4$** The harmful effects don't cause the good effects. Four permutations, paralleling the definition of $\rhd$ above, hold here. One such permutation is shown below. For any bad fluent $f_b$ holding at $t_1$, and any good fluent $f_g$ holding at some $t_2$, such that $t < t_1, t_2 \leq H$, the following holds:

$$\Gamma \vdash \neg\rhd\big(Holds(f_b,t_1),Holds(f_g,t_2)\big)$$

## Formal Conditions for $\mathcal{DDE}$

**F$_1$** $\alpha$ carried out at $t$ is not forbidden. That is:

$$\Gamma \nvdash \neg\mathbf{O}\Big(a,t,\sigma,\neg happens\big(action(a,\alpha),t\big)\Big)$$

**F$_2$** The net utility is greater than a given positive real $\gamma$:

$$\Gamma \vdash \sum_{y=t+1}^{H}\left(\sum_{f\in\alpha_I^{a,t}}\mu(f,y)-\sum_{f\in\alpha_T^{a,t}}\mu(f,y)\right) > \gamma$$

**F$_{3a}$** The agent $a$ intends at least one good effect. (**F$_2$** should still hold after removing all other good effects.) There is at least one fluent $f_g$ in $\alpha_I^{a,t}$ with $\mu\left(f_g,y\right) > 0$, or $f_b$ in $\alpha_T^{a,t}$ with $\mu\left(f_b,y\right) < 0$, and some $y$ with $t < y \le H$ such that the following holds:

$$\Gamma \vdash \begin{pmatrix} \exists f_g \in \alpha_I^{a,t}\ \mathbf{I}\big(a,t,Holds(f_g,y)\big) \\ \vee \\ \exists f_b \in \alpha_T^{a,t}\ \mathbf{I}\big(a,t,\neg Holds(f_b,y)\big) \end{pmatrix}$$

**F$_{3b}$** The agent $a$ does not intend any bad effect. For all fluents $f_b$ in $\alpha_I^{a,t}$ with $\mu\left(f_b,y\right) < 0$, or $f_g$ in $\alpha_T^{a,t}$ with $\mu\left(f_g,y\right) > 0$, and for all $y$ such that $t < y \le H$ the following holds:

$$\Gamma \nvdash \mathbf{I}\big(a,t,Holds(f_b,y)\big) \text{ and}$$

$$\Gamma \nvdash \mathbf{I}\big(a,t,\neg Holds(f_g,y)\big)$$

**F$_4$** The harmful effects don't cause the good effects. Four permutations, paralleling the definition of $\rhd$ above, hold here. One such permutation is shown below. For any bad fluent $f_b$ holding at $t_1$, and any good fluent $f_g$ holding at some $t_2$, such that $t < t_1, t_2 \le H$, the following holds:

$$\Gamma \vdash \neg\rhd\Big(Holds(f_b,t_1),Holds(f_g,t_2)\Big)$$

**Formal Conditions for** $\mathcal{DDE}$

**F₁** $\alpha$ carried out at $t$ is not forbidden. That is:

$$\Gamma \not\vdash \neg\mathbf{O}\Big(a,t,\sigma,\neg happens\big(action(a,\alpha),t\big)\Big)$$

**F₂** The net utility is greater than a given positive real $\gamma$:

$$\Gamma \vdash \sum_{y=t+1}^{H}\left(\sum_{f\in\alpha_I^{a,t}}\mu(f,y) - \sum_{f\in\alpha_T^{a,t}}\mu(f,y)\right) > \gamma$$

**F₃ₐ** The agent $a$ intends at least one good effect. (**F₂** should still hold after removing all other good effects.) There is at least one fluent $f_g$ in $\alpha_I^{a,t}$ with $\mu\big(f_g,y\big) > 0$, or $f_b$ in $\alpha_T^{a,t}$ with $\mu\big(f_b,y\big) < 0$, and some $y$ with $t < y \leq H$ such that the following holds:

$$\Gamma \vdash \begin{pmatrix} \exists f_g \in \alpha_I^{a,t}\ \mathbf{I}\big(a,t,Holds(f_g,y)\big) \\ \vee \\ \exists f_b \in \alpha_T^{a,t}\ \mathbf{I}\big(a,t,\neg Holds(f_b,y)\big) \end{pmatrix}$$

**F₃ᵦ** The agent $a$ does not intend any bad effect. For all fluents $f_b$ in $\alpha_I^{a,t}$ with $\mu\big(f_b,y\big) < 0$, or $f_g$ in $\alpha_T^{a,t}$ with $\mu\big(f_g,y\big) > 0$, and for all $y$ such that $t < y \leq H$ the following holds:

$$\Gamma \not\vdash \mathbf{I}\big(a,t,Holds(f_b,y)\big) \text{ and}$$

$$\Gamma \not\vdash \mathbf{I}\big(a,t,\neg Holds(f_g,y)\big)$$

**F₄** The harmful effects don't cause the good effects. Four permutations, paralleling the definition of $\triangleright$ above, hold here. One such permutation is shown below. For any bad fluent $f_b$ holding at $t_1$, and any good fluent $f_g$ holding at some $t_2$, such that $t < t_1, t_2 \leq H$, the following holds:

$$\Gamma \vdash \neg\triangleright\Big(Holds(f_b,t_1),Holds(f_g,t_2)\Big)$$

# Robotic "Jungle Jim"

# Robotic "Jungle Jim"

# Robotic "Jungle Jim"

Top machine-ethicists-may-consider-banging-their-heads-against-a-wall-hard.

AI Variant of "Jungle Jim" (B Williams)

H  H  H  H  H

J

"Robot R: You shoot just one human prisoner, the other four can go free. If you refuse to shoot, I'll shoot them all, now. Because I'm feeling generous, I'll give you a minute to decide."

R

H  H  H  H  H

J

"Robot R: You shoot just one human prisoner, the other four can go free. If you refuse to shoot, I'll shoot them all, now. Because I'm feeling generous, I'll give you a minute to decide."

R

H  H  H  H

J

"Robot R: You shoot just one human prisoner, the other four can go free. If you refuse to shoot, I'll shoot them all, now. Because I'm feeling generous, I'll give you a minute to decide."

R

H    H    H

J

"Robot R: You shoot just
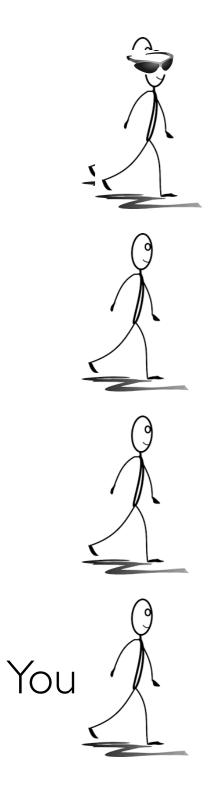one human prisoner, the
other four can go free.  If
you refuse to shoot, I'll
shoot them all, now.
Because I'm feeling
generous, I'll give you a
minute to decide."

R

H    H

J

"Robot R: You shoot just one human prisoner, the other four can go free. If you refuse to shoot, I'll shoot them all, now. Because I'm feeling generous, I'll give you a minute to decide."

R

H

J

"Robot R: You shoot just one human prisoner, the other four can go free. If you refuse to shoot, I'll shoot them all, now. Because I'm feeling generous, I'll give you a minute to decide."

R

J

"Robot R: You shoot just one human prisoner, the other four can go free. If you refuse to shoot, I'll shoot them all, now. Because I'm feeling generous, I'll give you a minute to decide."

R

H   H   H   H   H

J

"Robot R: You shoot just one human prisoner, the other four can go free. If you refuse to shoot, I'll shoot them all, now. Because I'm feeling generous, I'll give you a minute to decide."

R

H    H    H    H    H

J

"Robot R: You shoot just one human prisoner, the other four can go free. If you refuse to shoot, I'll shoot them all, now. Because I'm feeling generous, I'll give you a minute to decide."

R

J

"Robot R: You shoot just one human prisoner, the other four can go free. If you refuse to shoot, I'll shoot them all, now. Because I'm feeling generous, I'll give you a minute to decide."

R

J

"Robot R: You shoot just one human prisoner, the other four can go free. If you refuse to shoot, I'll shoot them all, now. Because I'm feeling generous, I'll give you a minute to decide."

R

?

# Level 3: Robotic "Jungle Jim"

# Level 3: Robotic "Jungle Jim"

# Level 3: Robotic "Jungle Jim"

# Level 3: Robotic "Jungle Jim"

You

You: "There's a glorious pink rose!"

Naveen: "That is not so."

You: "There's a glorious pink rose!"

Naveen: "That is not so."

Selmer: "Remove your glasses."

You: "There's a glorious pink rose!"

# Ontological Inventory: Classical Triad

# Ontological Inventory: Classical Triad

You (replete with sensors & effectors).

# Ontological Inventory: Classical Triad

You (replete with sensors & effectors).

The white rose.

# Ontological Inventory: Classical Triad

You (replete with sensors & effectors).

The white rose.

That which you perceived; the sense-datum
that led you to believe that you saw a *pink* rose.

# Ontological Inventory: Adverbial Theory of Perception

# Ontological Inventory: Adverbial Theory of Perception

You (replete with sensors & effectors).

# Ontological Inventory: Adverbial Theory of Perception

You (replete with sensors & effectors).

The white rose.

# Ontological Inventory: Adverbial Theory of Perception

You (replete with sensors & effectors).

The white rose.

And that's it! — because you perceive pinkly.

# Ontological Inventory:
# Adverbial Theory of Perception

You (replete with sensors & effectors).

The white rose.

And that's it! — because you perceive pinkly.

# The Adverbial Approach to (Machine) Ethics

makingmorallyxmachines.com

Plated   News ⌄   SUNY System ...ess - Logon   Ultra Hardwar...are Products   Screen Door L...y Von Morris   Screen Door ...d Von Morris   Apple   Amazon   eBay   Yahoo!   .Mac

Gmail                                                        Making Morally <img src="./GRAPHICS/XquaRobot.png" alt="XquaRobot.png" /> Machines

# Making Morally X Machines

## Table of Contents

- [The Book](#)
- [The Demonstrations](#)
- [Tutorials](#)
- [Code](#)
- [Media](#)

Selmer Bringsjord ∧ Naveen Sundar Govindarajulu ∧ John Licato

## The Book

- Setting the Stage
    - Overview of the Book
    - The "PAID" Problem and The Singularity
    - The Adverbial Approach to Machine Ethics
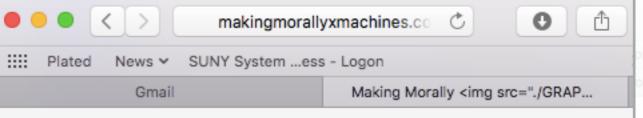    - The Solution in Four Steps
- Technique, Tools Technology

# Making Morally X Machines

## Table of Contents

- [The Book](#)
- [The Demonstrations](#)
- [Tutorials](#)
- [Code](#)
- [Media](#)

Selmer Bringsjord ∧ Naveen Sundar Govindarajulu ∧ John Licato

## The Book

- Setting the Stage
  - Overview of the Book
  - The "PAID" Problem and The Singularity
  - The Adverbial Approach to Machine Ethics
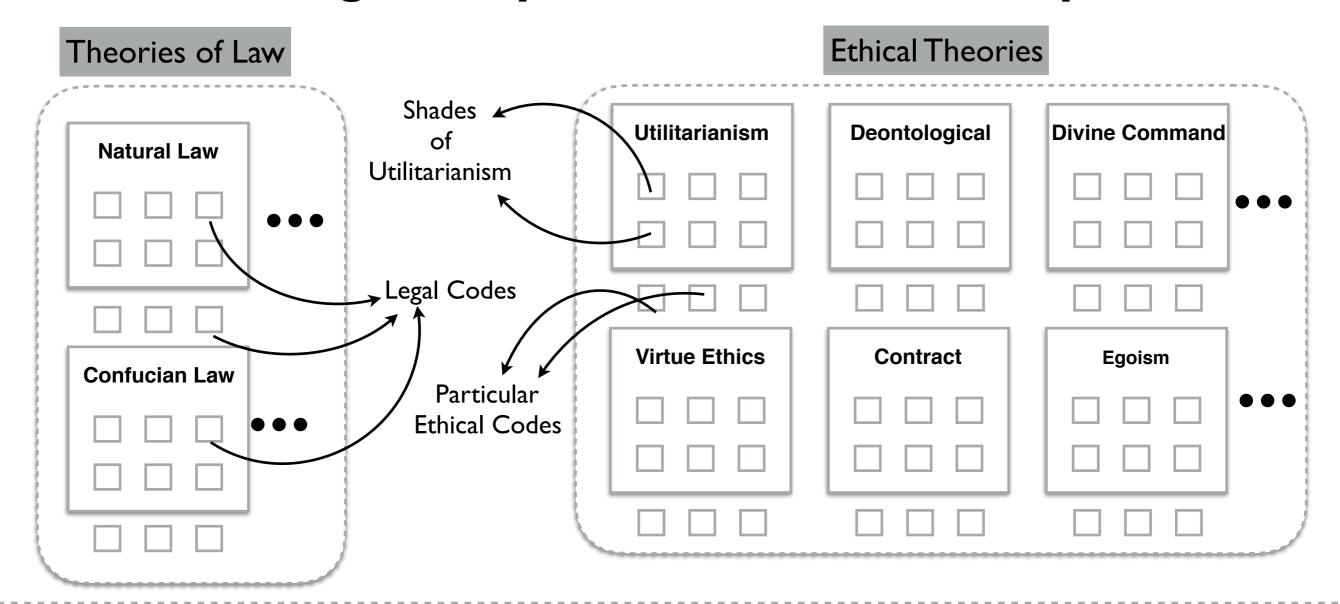  - The Solution in Four Steps
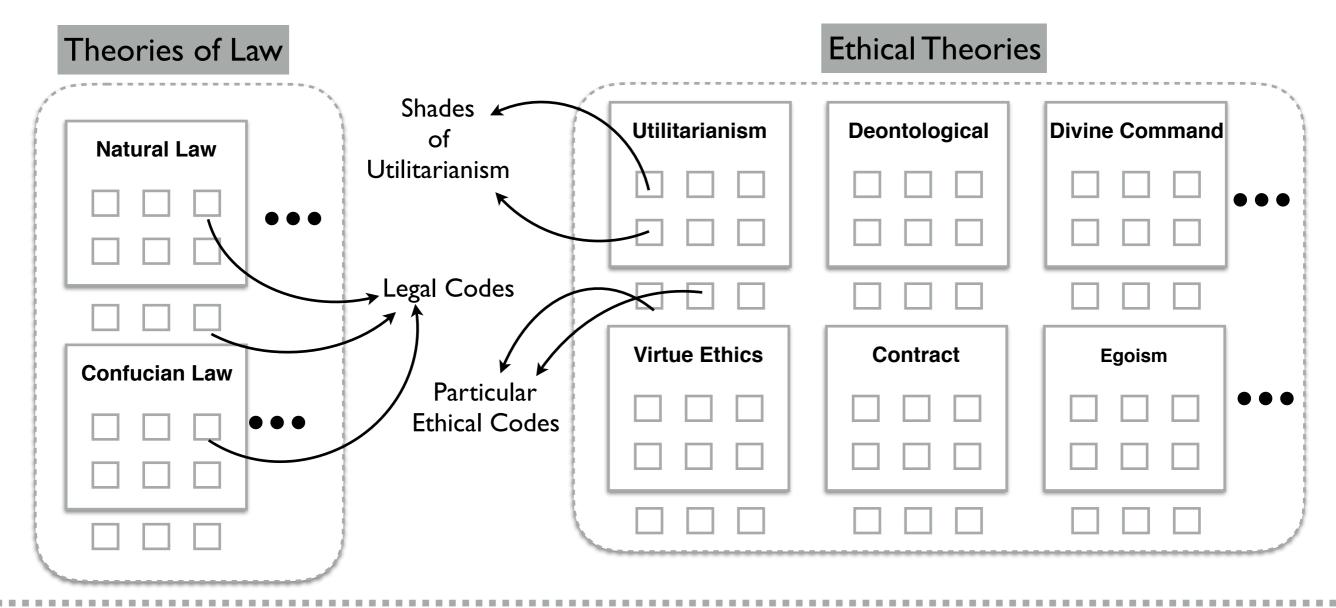- Technique, Tools Technology

# The Book

- Setting the Stage
  - Overview of the Book
  - The "PAID" Problem and The Singularity
  - The Adverbial Approach to Machine Ethics
  - The Solution in Four Steps
- Technique, Tools Technology
  - Formalism (the Philosophy of)
  - Computing With Cognitive Calculi
  - Infrastructure for Our Demonstrations
  - An Ethical Hierarchy ($\mathcal{EH}$) for Persons and Machines
- Making Morally $X$ Machines
  - Making Morally Invulnerable Machines
  - Making Morally Incorruptible Machines
  - Making Morally Self-Aware Machines
  - Making Morally Courageous Machines
  - Making Morally Educable Machines
  - Making Morally Wise Machines
  - Making Morally Creative Machines
  - Making Morally Well-Intentioned Machines
  - Making Morally Autonomous Machines
  - Making Morally Courteous Machines
  - Making Morally Heroic Machines
- What Shall We Now Do?
- References

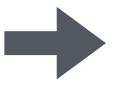# Making Morally *X* Machines, in Four Steps



Theories of Law

Natural Law

Confucian Law

Shades of Utilitarianism

Legal Codes

Particular Ethical Codes

Ethical Theories

Utilitarianism

Deontological

Divine Command

Virtue Ethics

Contract

Egoism

# Making Morally *X* Machines, in Four Steps

Theories of Law

Ethical Theories

**Natural Law**

**Confucian Law**

Shades
of
Utilitarianism

Legal Codes

Particular
Ethical Codes

**Utilitarianism**

**Deontological**

**Divine Command**

**Virtue Ethics**

**Contract**

**Egoism**

Step 1

1. Pick (a) theories.
2. Pick (a) code(s).
3. Run through EH.
4. Which *X* in *MMXM*?

# Making Morally *X* Machines, in Four Steps



**Theories of Law**

Natural Law

• • •

Confucian Law

• • •

Shades of Utilitarianism

Legal Codes

Particular Ethical Codes

**Ethical Theories**

Utilitarianism

Deontological

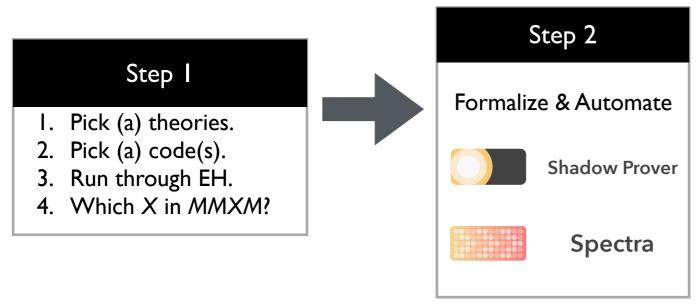Divine Command

• • •

Virtue Ethics

Contract

Egoism

• • •

**Step 1**

1. Pick (a) theories.
2. Pick (a) code(s).
3. Run through EH.
4. Which *X* in *MMXM?*

# Making Morally *X* Machines, in Four Steps

**Theories of Law**

**Natural Law**

**Confucian Law**

Shades
of
Utilitarianism

Legal Codes

Particular
Ethical Codes

**Ethical Theories**

**Utilitarianism**

**Deontological**

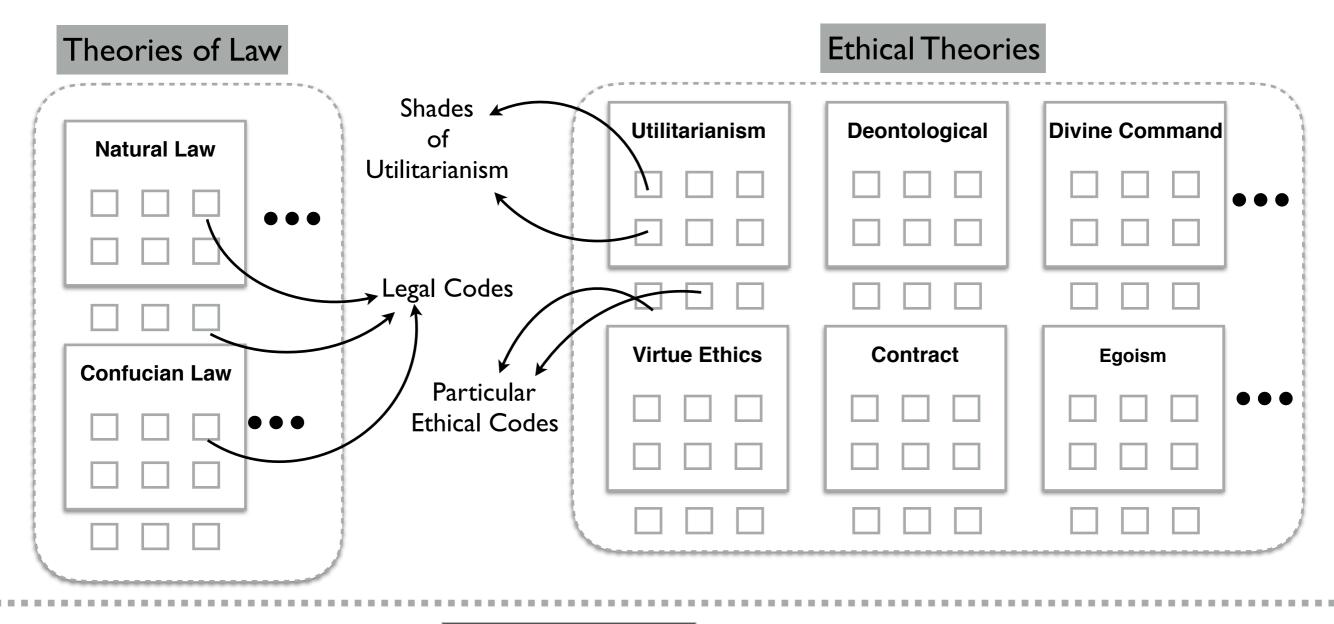**Divine Command**
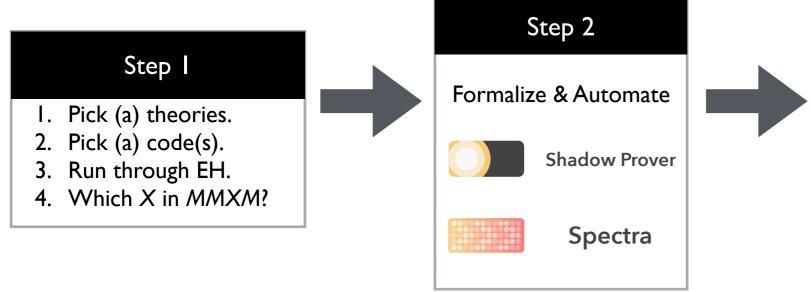
**Virtue Ethics**

**Contract**

**Egoism**

**Step 1**

1. Pick (a) theories.
2. Pick (a) code(s).
3. Run through EH.
4. Which *X* in *MMXM*?

**Step 2**

Formalize & Automate

Shadow Prover

Spectra

# Making Morally *X* Machines, in Four Steps

**Theories of Law**

**Natural Law**

**Confucian Law**

Shades of Utilitarianism

Legal Codes

Particular Ethical Codes

**Ethical Theories**

**Utilitarianism**

**Deontological**

**Divine Command**

**Virtue Ethics**

**Contract**

**Egoism**

**Step 1**

1. Pick (a) theories.
2. Pick (a) code(s).
3. Run through EH.
4. Which *X* in *MMXM*?

**Step 2**

Formalize & Automate

Shadow Prover

Spectra

# Making Morally *X* Machines, in Four Steps

**Theories of Law**

**Natural Law**

**Confucian Law**

Shades
of
Utilitarianism

Legal Codes

Particular
Ethical Codes

**Ethical Theories**

**Utilitarianism**

**Deontological**

**Divine Command**

**Virtue Ethics**

**Contract**

**Egoism**

### Step 1

1. Pick (a) theories.
2. Pick (a) code(s).
3. Run through EH.
4. Which *X* in *MMXM*?

### Step 2

Formalize & Automate

Shadow Prover

Spectra

### Step 3

Ethical OS

Ethical Substrate

Robotic Substrate

# Making Morally *X* Machines, in Four Steps

**Theories of Law**

**Natural Law**

**Confucian Law**

Shades
of
Utilitarianism

Legal Codes

Particular
Ethical Codes

**Ethical Theories**

**Utilitarianism**

**Deontological**

**Divine Command**

**Virtue Ethics**

**Contract**

**Egoism**

**Step 1**

1. Pick (a) theories.
2. Pick (a) code(s).
3. Run through EH.
4. Which *X* in *MMXM*?

**Step 2**

Formalize & Automate

Shadow Prover

Spectra

**Step 3**

Ethical OS

Ethical Substrate

Robotic Substrate

# Making Morally *X* Machines, in Four Steps

## Theories of Law

**Natural Law**

**Confucian Law**

## Ethical Theories

**Utilitarianism**

**Deontological**

**Divine Command**

**Virtue Ethics**

**Contract**

**Egoism**

Shades of Utilitarianism

Legal Codes

Particular Ethical Codes

**Step 1**

1. Pick (a) theories.
2. Pick (a) code(s).
3. Run through EH.
4. Which *X* in *MMXM*?

**Step 2**

Formalize & Automate

Shadow Prover

Spectra

**Step 3**

Ethical OS

Ethical Substrate

Robotic Substrate

*An ethically correct robot.*

# Making Morally *X* Machines, in Four Steps

## Theories of Law

**Natural Law**

**Confucian Law**

Shades of Utilitarianism

Legal Codes

Particular Ethical Codes

## Ethical Theories

**Utilitarianism**

**Deontological**

**Divine Command**

**Virtue Ethics**

**Contract**

**Egoism**

### Step 1
1. Pick (a) theories.
2. Pick (a) code(s).
3. Run through EH.
4. Which *X* in *MMXM*?

### Step 2
Formalize & Automate

Shadow Prover

Spectra

### Step 3
Ethical OS

Ethical Substrate

Robotic Substrate

DIARC/DoD/BMW ...

*An ethically correct robot.*

# Making Morally *X* Machines, in Four Steps

**Theories of Law**

**Ethical Theories**

**Natural Law**

**Confucian Law**

Shades of Utilitarianism

Legal Codes

Particular Ethical Codes

**Utilitarianism**

**Deontological**

**Divine Command**

**Virtue Ethics**

**Contract**

**Egoism**

**Step 1**
1. Pick (a) theories.
2. Pick (a) code(s).
3. Run through EH.
4. Which *X* in *MMXM*?

**Step 2**

Formalize & Automate

Shadow Prover

Spectra

**Step 3**

Ethical OS

Ethical Substrate

Robotic Substrate

DIARC/DoD/BMW ...

*An ethically correct robot.*

# The Book

- Setting the Stage
  - Overview of the Book
  - The "PAID" Problem and The Singularity
  - The Adverbial Approach to Machine Ethics
  - The Solution in Four Steps
- Technique, Tools Technology
  - Formalism (the Philosophy of)
  - Computing With Cognitive Calculi
  - Infrastructure for Our Demonstrations
  - An Ethical Hierarchy ($\mathcal{EH}$) for Persons and Machines
- Making Morally $X$ Machines
  - Making Morally Invulnerable Machines
  - Making Morally Incorruptible Machines
  - Making Morally Self-Aware Machines
  - Making Morally Courageous Machines
  - Making Morally Educable Machines
  - Making Morally Wise Machines
  - Making Morally Creative Machines
  - Making Morally Well-Intentioned Machines
  - Making Morally Autonomous Machines
  - Making Morally Courteous Machines
  - Making Morally Heroic Machines
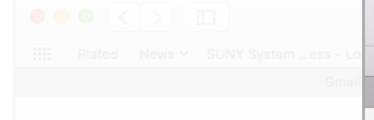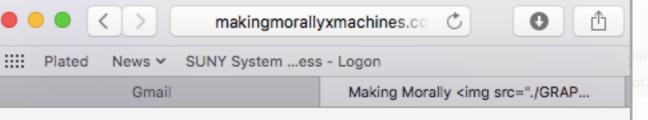- What Shall We Now Do?
- References

# The Book

- Setting the Stage
  - Overview of the Book
  - The "PAID" Problem and The Singularity
  - The Adverbial Approach to Machine Ethics
  - The Solution in Four Steps
- Technique, Tools Technology
  - Formalism (the Philosophy of)
  - Computing With Cognitive Calculi
  - Infrastructure for Our Demonstrations
  - An Ethical Hierarchy ($\mathcal{EH}$) for Persons and Machines
- Making Morally $X$ Machines
  - Making Morally Invulnerable Machines
  - Making Morally Incorruptible Machines
  - Making Morally Self-Aware Machines
  - Making Morally Courageous Machines
  - Making Morally Educable Machines
  - Making Morally Wise Machines
  - Making Morally Creative Machines
  - Making Morally Well-Intentioned Machines
  - Making Morally Autonomous Machines
  - Making Morally Courteous Machines
  - Making Morally Heroic Machines
- What Shall We Now Do?
- References

# End

(Extra slides follow.)