

Discursive Input/Output Logic: Deontic Modals, and Computation

Ali Farjami

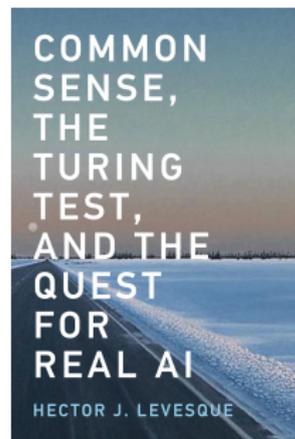
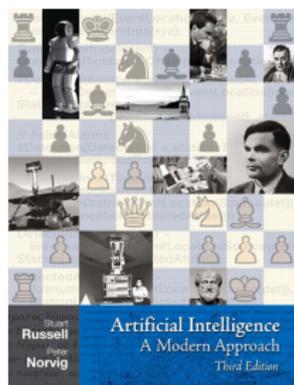
University of Luxembourg



PhD Defense

October 1, 2020

Logic-based AI



- ▶ Logic in computer science (1980)
 - ▶ Relational database: formulas defines queries
 - ▶ Boolean satisfiability (SAT)
 - ▶ Non-monotonic logic
 - ▶ Knowledge representation and reasoning (KR) (1990)
 - ▶ Machine learning + Symbolic logic (2020)
 - ▶ Trustworthy and responsible AI
 - ▶ Deontic Logic
- Codd 1981 (Turing award)
SAT solvers
Common sense reasoning*
- Semantic Web*
- Neuro-Symbolic AI*
- Contrary-to-duty*

The Miners Example

	<i>in_A</i>	<i>in_B</i>
<i>block_A</i>	All live	All die
<i>block_B</i>	All die	All live
$\neg(\textit{block_A} \vee \textit{block_B})$	Nine live	Nine live



Deontic Modals

Modal Logic Approach

Danielsson (1968), Hansson (1969), Føllesdal and Hilpinen (1970), van Fraassen (1973), and Lewis (1974), Kratzer (1977), ...

- ▶ A set of accessible worlds



Norm-based Approach

Makinson (1998), Makinson, van der Torre (2000,2001), Horty (2012), Hansen (2008), ...

- ▶ A set of norms
- ▶ Input/output logic
- ▶ Inference patterns



Research Objectives

Norms + Informational Modalities

- ▶ How can we use an algebraic setting such as Boolean algebras instead of a logical setting for building input/output logic on top of it?
- ▶ How can we introduce two groups of I/O operations similar to syntactical characterization of box and diamond in modal logic?
- ▶ How can we integrate conversational background informations, from Kratzerian framework, into input/output logic framework to build a more fruitful unified semantics for deontic modals?

Norms + Preferences

- ▶ How can we integrate input/output logic with Hansson and Lewis's conditional theory for building a new compositional theory about conditional deontic modals?
 - ▶ Resolving contrary-to-duty problems
 - ▶ Non-monotonic defeat mechanism within Hansson and Lewis's conditionals

Normative Reasoning in Computer Systems

- ▶ Providing a (faithful) embedding of some well-known deontic logics in HOL
- ▶ Encoding the logical embeddings in Isabelle/HOL

Research Objectives

Norms + Informational Modalities

- ▶ How can we use an algebraic setting such as Boolean algebras instead of a logical setting for building input/output logic on top of it?
- ▶ How can we introduce two groups of I/O operations similar to syntactical characterization of box and diamond in modal logic?
- ▶ How can we integrate conversational background informations, from Kratzerian framework, into input/output logic framework to build a more fruitful unified semantics for deontic modals?

Norms + Preferences

- ▶ How can we integrate input/output logic with Hansson and Lewis's conditional theory for building a new compositional theory about conditional deontic modals?
 - ▶ Resolving contrary-to-duty problems
 - ▶ Non-monotonic defeat mechanism within Hansson and Lewis's conditionals

Normative Reasoning in Computer Systems

- ▶ Providing a (faithful) embedding of some well-known deontic logics in HOL
- ▶ Encoding the logical embeddings in Isabelle/HOL

Research Objectives

Norms + Informational Modalities

- ▶ How can we use an algebraic setting such as Boolean algebras instead of a logical setting for building input/output logic on top of it?
- ▶ How can we introduce two groups of I/O operations similar to syntactical characterization of box and diamond in modal logic?
- ▶ How can we integrate conversational background informations, from Kratzerian framework, into input/output logic framework to build a more fruitful unified semantics for deontic modals?

Norms + Preferences

- ▶ How can we integrate input/output logic with Hansson and Lewis's conditional theory for building a new compositional theory about conditional deontic modals?
 - ▶ Resolving contrary-to-duty problems
 - ▶ Non-monotonic defeat mechanism within Hansson and Lewis's conditionals

Normative Reasoning in Computer Systems

- ▶ Providing a (faithful) embedding of some well-known deontic logics in HOL
- ▶ Encoding the logical embeddings in Isabelle/HOL

Research Objectives

Norms + Informational Modalities

- ▶ How can we use an algebraic setting such as Boolean algebras instead of a logical setting for building input/output logic on top of it?
- ▶ How can we introduce two groups of I/O operations similar to syntactical characterization of box and diamond in modal logic?
- ▶ How can we integrate conversational background informations, from Kratzerian framework, into input/output logic framework to build a more fruitful unified semantics for deontic modals?

Norms + Preferences

- ▶ How can we integrate input/output logic with Hansson and Lewis's conditional theory for building a new compositional theory about conditional deontic modals?
 - ▶ Resolving contrary-to-duty problems
 - ▶ Non-monotonic defeat mechanism within Hansson and Lewis's conditionals

Normative Reasoning in Computer Systems

- ▶ Providing a (faithful) embedding of some well-known deontic logics in HOL
- ▶ Encoding the logical embeddings in Isabelle/HOL

Research Objectives

Norms + Informational Modalities

- ▶ How can we use an algebraic setting such as Boolean algebras instead of a logical setting for building input/output logic on top of it?
- ▶ How can we introduce two groups of I/O operations similar to syntactical characterization of box and diamond in modal logic?
- ▶ How can we integrate conversational background informations, from Kratzerian framework, into input/output logic framework to build a more fruitful unified semantics for deontic modals?

Norms + Preferences

- ▶ How can we integrate input/output logic with Hansson and Lewis's conditional theory for building a new compositional theory about conditional deontic modals?
 - ▶ Resolving contrary-to-duty problems
 - ▶ Non-monotonic defeat mechanism within Hansson and Lewis's conditionals

Normative Reasoning in Computer Systems

- ▶ Providing a (faithful) embedding of some well-known deontic logics in HOL
- ▶ Encoding the logical embeddings in Isabelle/HOL

Methodology

Normative Reasoning

- ▶ Algebraic approach to input/output logic *Connection to modal logic*
 - ▶ Gabbay, Parent, and van der Torre: a geometrical view of I/O logic
 - ▶ *Upward-closed set of the infimum of A* instead of $C_n(A)$ *Compactness (?),...*
 - ▶ We use *upward-closed set* of A *Removing AND*
 - ▶ Reversibility of inference rules *Adding AND*
- ▶ Non-adjunctive input/output operations

$$\{\varphi_1, \dots, \varphi_n\} \vdash \psi \implies \begin{cases} \varphi_1 \wedge \dots \wedge \varphi_n \vdash \psi \\ \varphi_i \vdash \psi \quad \varphi_i \in \{\varphi_1, \dots, \varphi_n\} \end{cases}$$

- ▶ Semantical unification: Detachment + Conversational backgrounds
 - ▶ Syntactical unification *Adaptive logic*
- ▶ A semantical characterization of constrained input/output logic *Preferences*
 - ▶ Syntactical characterization *Adaptive logic characterizations of I/O logic*
 - ▶ No need to AND, SI and EQ required for syntactical translation

Logic Engineering

- ▶ *Shallow semantical embedding*
- ▶ Translating into Higher-order logic (HOL)
- ▶ Benchmark examples

CTD

Conversational Backgrounds

Examples: knowledge, beliefs, relevant facts, desires, plans,...

Functions from evaluation worlds to sets of propositions

- ▶ *Modal base* determines the set of accessible worlds ($f(w)$)
- ▶ *Ordering source* induces the ordering on worlds ($g(w)$)



Quantification

$$[[\text{be-allowed-to}]]^{w,f,g} = \lambda x (Best_{g(w)}(\bigcap f(w)) \cap x \neq \emptyset)$$

$$[[\text{have-to}]]^{w,f,g} = \lambda x (Best_{g(w)}(\bigcap f(w)) \subseteq x)$$

where $Best_{g(w)}(\bigcap f(w))$ is given as follows:

$$\{w' \in \bigcap f(w) : \neg \exists w'' \in \bigcap f(w) \text{ such that } \exists y \in g(w) : w'' \in y \text{ and } w' \notin y\}$$

Quantification \implies Detachment

Conversational Backgrounds

Examples: knowledge, beliefs, relevant facts, desires, plans,...

Functions from evaluation worlds to sets of propositions

- ▶ *Modal base* determines the set of accessible worlds ($f(w)$)
- ▶ *Ordering source* induces the ordering on worlds ($g(w)$)



Quantification

Compatibility $[[\text{be-allowed-to}]]^{w,f,g} = \lambda x (Best_{g(w)}(\bigcap f(w)) \cap x \neq \emptyset)$

Entailment $[[\text{have-to}]]^{w,f,g} = \lambda x (Best_{g(w)}(\bigcap f(w)) \subseteq x)$

where $Best_{g(w)}(\bigcap f(w))$ is given as follows:

$$\{w' \in \bigcap f(w) : \neg \exists w'' \in \bigcap f(w) \text{ such that } \exists y \in g(w) : w'' \in y \text{ and } w' \notin y\}$$

Quantification \implies Detachment

Chisholm's Paradox

$$\text{DD} \frac{\bigcirc(t/g) \quad \bigcirc(g)}{\bigcirc(t)}$$

$$\text{FD} \frac{\bigcirc(\neg t/\neg g) \quad \neg g}{\bigcirc(\neg t)}$$

Norm-based Semantics: Input/output logic

- “ x is obligatory if a ” \rightsquigarrow “ x can be detached in context a ”
- **Output operation:** $x \in \text{out}(N^O, A)$ Normative system N^O
- *Detachment vs Quantification*

Detachment in Discursive Context : Out(N, Discursive Context)

- **Context** in a discourse **Modal base or ordering source**
- Modal bases **Factual**
- Ordering sources Possible inconsistency

Out(N, modal base/ordering source)

Chisholm's Paradox

$$\text{DD} \frac{\bigcirc(t/g) \quad \bigcirc(g)}{\bigcirc(t)}$$

$$\text{FD} \frac{\bigcirc(\neg t/\neg g) \quad \neg g}{\bigcirc(\neg t)}$$

Norm-based Semantics: Input/output logic

- “ x is obligatory if a ” \rightsquigarrow “ x can be detached in context a ”
- **Output operation:** $x \in \text{out}(N^O, A)$ Normative system N^O
- *Detachment vs Quantification*

Detachment in Discursive Context : Out(N, Discursive Context)

- **Context** in a discourse **Modal base or ordering source**
- Modal bases **Factual**
- Ordering sources **Possible inconsistency**

Out(N, modal base/ordering source)

On a Fundamental Problem of Deontic Logic

David Makinson

Les Etangs B2, La Ronce, 92410 Ville d'Avray, France

Email: d.makinson@unesco.org

DAVID MAKINSON and LEENDERT VAN DER TORRE

INPUT/OUTPUT LOGICS

Normative Systems

(Received on 16 November 1999; final version received on 13 March 2000)

- ▶ T: infer every (\top, \top)
- ▶ SI: from (a, x) and $\vdash b \rightarrow a$, infer (b, x)
- ▶ WO: from (a, x) and $\vdash x \rightarrow y$, infer (a, y)
- ▶ AND: from (a, x) and (a, y) , infer $(a, x \wedge y)$
- ▶ OR: from (a, x) and (b, x) , infer $(a \vee b, x)$
- ▶ CT: from (a, x) and $(a \wedge x, y)$, infer (a, y)

- **Unconstrained input/output logic**
- **Constrained input/output logic**

AGM theory/Contrary-to-duty

Input/Output Logic: Output operations

- ▶ Simple-Minded Output:

$$out_1(N, A) = Cn(N(Cn(A)))$$

- ▶ Basic Output:

$$out_2(N, A) = \bigcap \{Cn(N(V)) \mid A \subseteq V, V \text{ complete}\}$$

- ▶ Simple-Minded Reusable Output:

$$out_3(N, A) = \bigcap \{Cn(N(B)) \mid A \subseteq B = Cn(B) \supseteq N(V)\}$$

- ▶ Basic Reusable Output:

$$out_4(N, A) = \bigcap \{Cn(N(V)) \mid A \subseteq V \supseteq N(V), V \text{ complete}\}$$

$deriv_i(N)$	Rules
$deriv_1(N)$	{T, SI, WO, AND}
$deriv_2(N)$	{T, SI, WO, OR, AND}
$deriv_3(N)$	{T, SI, WO, CT, AND}
$deriv_4(N)$	{T, SI, WO, OR, CT, AND}

Input/Output Logic: Output operations

- ▶ Simple-Minded Output:

$$out_1(N, A) = \cancel{Cn(N(\cancel{Cn(A))}} \quad Up(N(Up(A)))$$

- ▶ Basic Output:

$$out_2(N, A) = \bigcap \{\cancel{Cn(N(V))} \mid A \subseteq V, V \text{ complete}\} \quad \dots$$

- ▶ Simple-Minded Reusable Output:

$$out_3(N, A) = \bigcap \{\cancel{Cn(N(B))} \mid A \subseteq B = Cn(B) \supseteq N(V)\} \quad \dots$$

$deriv_i(N)$	Rules
$deriv_1(N)$	{T, SI, WO, A , D }
$deriv_2(N)$	{T, SI, WO, OR, A , D }
$deriv_3(N)$	{T, SI, WO, T , A , D }

- I/O operations over Boolean algebras

- Stone's representation theorem

Possible world semantics

- $Up(X) = \{x \in B \mid \exists y \in X, y \leq x\}$

$$a \wedge b \notin Up(a, b)$$

Non-adjunctive Logical Systems

Deriving the conjunctive formula $\varphi \wedge \psi$ from the set $\{\varphi, \psi\}$ fails



Discursive Systems

*"[...] the **joining** of a thesis to a discursive system has a different intuitive meaning than has assertion in an ordinary system."* jaskowski1969

$$A \rightarrow_d B$$

$$\diamond A \rightarrow B$$

$$Up(A \cup B) = Up(A) \cup Up(B); \quad out_i(N, A) = \bigcup_{a \in A} out_i(N, a)$$

- ▶ Simple-Minded Output :

$$out_1^{\mathcal{B}}(N, A) = Up(N(Up(A)))$$

- ▶ Basic Boolean I/O operation:

$$out_2^{\mathcal{B}}(N, A) = \bigcap \{Up(N(V)), A \subseteq V, V \text{ is saturated}\}$$

- ▶ Reusable Boolean I/O operation:

$$out_3^{\mathcal{B}}(N, A) = \bigcap \{Up(N(V)), A \subseteq V = Up(V) \supseteq N(V)\}$$

Discursive Input/Output Logic: Output operations

- ▶ Zero Boolean I/O operation: ($out_R(N, A) = Eq(N(A))$,
 $out_L(N, A) = N(Eq(A))$)

$$out_0^B(N, A) = Eq(N(Eq(A)))$$

- ▶ Simple-I Boolean I/O operation:

$$out_I^B(N, A) = Eq(N(Up(A)))$$

- ▶ Simple-II Boolean I/O operation:

$$out_{II}^B(N, A) = Up(N(Eq(A)))$$

- ▶ Simple-Minded Output :

$$out_1^B(N, A) = Up(N(Up(A)))$$

- ▶ Basic Boolean I/O operation:

$$out_2^B(N, A) = \bigcap \{Up(N(V)), A \subseteq V, V \text{ is saturated}\}$$

- ▶ Reusable Boolean I/O operation:

$$out_3^B(N, A) = \bigcap \{Up(N(V)), A \subseteq V = Up(V) \supseteq N(V)\}$$

Discursive Input/Output Logic: Proof system

$(A, x) \in \text{deriv}_i(N)$ if $(a, x) \in \text{deriv}(N)$ for some $a \in A$

deriv_i^B	Rules				
deriv_R^B	{EQO}	EQO	$\frac{(a, x) \quad x = y}{(a, y)}$	WO	$\frac{(a, x) \quad x \leq y}{(a, y)}$
deriv_L^B	{EQI}				
deriv_0^B	{EQI, EQO}	EQI	$\frac{(a, x) \quad a = b}{(b, x)}$	OR	$\frac{(a, x) \quad (b, x)}{(a \vee b, x)}$
deriv_I^B	{SI, EQO}				
deriv_{II}^B	{WO, EQI}				
deriv_1^B	{SI, WO}				
deriv_2^B	{SI, WO, OR}	SI	$\frac{(a, x) \quad b \leq a}{(b, x)}$	T	$\frac{(a, x) \quad (x, y)}{(a, y)}$
deriv_3^B	{SI, WO, T}				

◇ VS □

$(\Box A \wedge \Box B) \rightarrow \Box(A \wedge B)$

Discursive Input/Output Logic: Proof system

$(A, x) \in \text{deriv}_i(N)$ if $(a, x) \in \text{deriv}(N)$ for some $a \in A$

deriv_i^B	Rules				
deriv_R^B	{EQO}	EQO	$\frac{(a, x) \quad x = y}{(a, y)}$	WO	$\frac{(a, x) \quad x \leq y}{(a, y)}$
deriv_L^B	{EQI}				
deriv_0^B	{EQI, EQO}	EQI	$\frac{(a, x) \quad a = b}{(b, x)}$	OR	$\frac{(a, x) \quad (b, x)}{(a \vee b, x)}$
deriv_I^B	{SI, EQO}				
deriv_{II}^B	{WO, EQI}				
deriv_1^B	{SI, WO}				
deriv_2^B	{SI, WO, OR}	SI	$\frac{(a, x) \quad b \leq a}{(b, x)}$	T	$\frac{(a, x) \quad (x, y)}{(a, y)}$
deriv_3^B	{SI, WO, T}				

◇ VS □

$(\Box A \wedge \Box B) \rightarrow \Box(A \wedge B)$

Adding Other Rules

$$\text{AND} \frac{(p, q) \quad (p, r)}{(p, q \wedge r)}$$

$$\text{CT} \frac{(p, q) \quad (p \wedge q, r)}{(p, r)}$$

Reversibility of Inference Rules

	SI	WO	CT		AND	OR
SI		✓	none?	none?	none?	✓
WO	✓		SI, CT	✓	✓	none?
CT	✓	✓			SI, AND, CT	none?
AND	✓	✓	SI, CT	✓		WO, OR, AND
OR	✓	✓	SI, CT, OR	SI, CT, OR	SI, AND, OR	

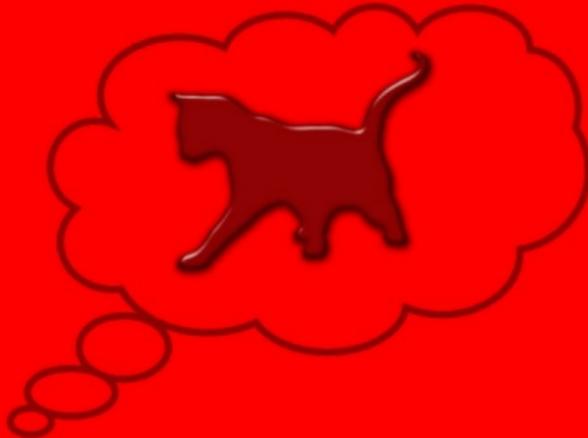
(Makinson and van der Torre 2000)

Adding AND: Output operations + Iteration of AND

$$\begin{aligned} out_i^{AND^0}(N, A) &= out_i^B(N, A) \\ out_i^{AND^{n+1}}(N, A) &= out_i^{AND^n}(N, A) \cup \\ &\quad \{y \wedge z : y, z \in out_i^{AND^n}(N, \{a\}), a \in A\} \\ out_i^{AND}(N, A) &= \bigcup_{n \in \mathbb{N}} out_i^{AND^n}(N, A) \end{aligned}$$

$deriv_i^X$	Rules
$deriv_{II}^{AND}$	{WO, EQI, AND}
$deriv_1^{AND}$	{SI, WO, AND}
$deriv_2^{AND}$	{SI, WO, OR, AND}

Semantics



Unification: Consistent premise sets

Suppose $\bigcap f(w) \neq \emptyset$

$$[[\text{be-allowed-to}]]^{wf} = \lambda N^P \lambda x (x \in \text{out}(N^P, \{\bigcap f(w)\}))$$

$$[[\text{have-to}]]^{wf} = \lambda N^O \lambda x (x \in \text{out}(N^O, \{\bigcap f(w)\}))$$

$$M, w \models \Box a_1 \wedge \dots \wedge \Box a_n = \Box(a_1 \wedge \dots \wedge a_n)$$

$$a_i \in f(w)$$

Unification: Inconsistent premise sets

Suppose $\bigcap g(w) = \emptyset$, and

$\text{Maxfamily}^{\cap}(g(w)) = \{\bigcap A \mid A \subseteq g(w) \text{ and } A \text{ is consistent and maximal}\}$

$$[[\text{be-allowed-to}]]^{w,g} = \lambda N^P \lambda x (x \in \text{out}(N^P, \text{Maxfamily}^{\cap}(g(w))))$$

$$[[\text{have-to}]]^{w,g} = \lambda N^O \lambda x (x \in \text{out}(N^O, \text{Maxfamily}^{\cap}(g(w))))$$

$$M, w \models \diamond \overbrace{(a_1 \wedge \dots \wedge a_n)}^{\text{Maxfamily}} \wedge \dots \wedge \diamond \overbrace{(b_1 \wedge \dots \wedge b_m)}^{\text{Maxfamily}} \quad a_i, b_j \in g(w)$$

The Miners Example

1- Either the miners are in shaft A or in shaft B.

2- If the miners are in shaft A, we should block shaft A.

3- If the miners are in shaft B, we should block shaft B.

4- We should block neither shaft.

- ▶ Syntactical analysis (Not satisfactory)
- ▶ Not allowed by the baseline algorithm of Kratzerian framework (Cariani 2020)
- ▶ Kolodny and MacFarlane 2010 (Modus ponens is invalid)

$$N = \{(ShA, blA), (ShB, blB), (\top, \neg blA \wedge \neg blB)\}$$

▶ $M, w \models \Box(shA \vee shB)$

$$f(w) = \{shA \vee shB\}$$

a set of **factual informations**

$$\neg blA \wedge \neg blB \in out(N^O, \{shA \vee shB\})$$

▶ $M, w \models \Diamond shA \wedge \Diamond shB$

$$g(w) = \{shA, shB\}$$

a set of possible **inconsistent informations**

$$blA, blB, \neg blA \wedge \neg blB \in out(N^O, \{shA, shB\})$$

▶ $C = \{blA\}$

$$f(w) = \{shA\}$$

constrained I/O logic

$$blA \in out_c(N^O, \{shA\})$$

If P, then Q

antecedent consequent

$$\mathbf{Fm}(X) = \langle Fm(X), \wedge^{\mathbf{Fm}(X)}, \vee^{\mathbf{Fm}(X)}, \neg^{\mathbf{Fm}(X)}, \top^{\mathbf{Fm}(X)}, \perp^{\mathbf{Fm}(X)} \rangle$$

$\Gamma \vDash_{\mathbf{BA}} \varphi$ if and only if $\Gamma \vdash_C \varphi$

$$(p, q) \in \mathit{derive}_i^{\mathbf{Fm}(X)}(N)$$

if and only if

$V(q) \in \mathit{out}_i^{\mathcal{B}}(N^V, \{V(p)\})$ for every $\mathcal{B} \in \mathbf{BA}$, for every valuation V on \mathcal{B}

Neighborhood Characterization of I/O Operations

$$f : P(W) \rightarrow P(P(W))$$

Input/output logic + Constraints (preferences)

$$\varphi > \bigcirc\psi \in \text{derive}_i^O(N^O)$$

if and only if

$$(\varphi, \psi) \in \text{derive}_i^{\mathbf{Fm}(X)}(N^O) \text{ and}$$

For every preference Boolean algebra $M = \langle \mathcal{B}, \mathcal{V}, \succeq_f \rangle$,
for every valuation $V_i \in \text{opt}_{\succeq_f}(\varphi)$ we
have $V_i(\psi) = 1_{\mathcal{B}}$

- ▶ \mathcal{B} is a Boolean algebra,
- ▶ $\mathcal{V} = \{V_i\}_{i \in I}$ is the set of valuations from $\mathbf{Fm}(X)$ on \mathcal{B} ,
- ▶ $\succeq_f \subseteq \mathcal{V} \times \mathcal{V}$: \succeq_f is a betterness or comparative goodness relation over valuations from $\mathbf{Fm}(X)$ to \mathcal{B} such that $V_i \succeq_f V_j$ iff $(\{\varphi|V_i(\varphi) = 1_{\mathcal{B}}\}, \{\psi|V_j(\psi) = 1_{\mathcal{B}}\}) \in f$.

$$\varphi > P\psi \in \text{derive}_i^P(N^P)$$

if and only if

$$(\varphi, \psi) \in \text{derive}_i^{\mathbf{Fm}(X)}(N^P) \text{ and}$$

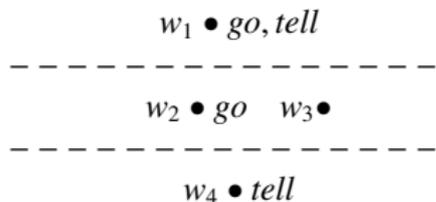
For every preference Boolean algebra $M = \langle \mathcal{B}, \mathcal{V}, \succeq_f \rangle$,
there is a valuation $V_i \in \text{opt}_{\succeq_f}(\varphi)$
such that $V_i(\psi) = 1_{\mathcal{B}}$

Chisholm's Paradox

It ought to be that a certain man go to help his neighbors.

It ought to be that if he goes he tell them he is coming.

If he does not go, he ought not to tell them he is coming.



$$N^O = \{(\top, g), (g, t), (\neg g, \neg t)\}$$

- ▶ $\top > \bigcirc g \in \text{derive}_i^O(N^O)$
- ▶ $g > \bigcirc t \in \text{derive}_i^O(N^O)$
- ▶ $\neg g > \bigcirc(\neg t) \in \text{derive}_i^O(N^O)$

KR Tools: Higher-order logic theorem provers



- ▶ Church's simple theory of types
- ▶ HOL provers
 - ▶ interactive:
 - ▶ automated:
- ▶ Isabelle/HOL
 - ▶ Bridges to external theorem provers
 - ▶ Model finders
 - ▶ Sophisticated user interaction

λ -calculus/Henkin models

Isabelle/HOL, HOL4, Hol Light, Coq/HOL
LEO-II, Satallax, Nitpick, Isabelle/HOL

Sledgehammer tool
Nitpick

Semantical Embedding

Aligning Henkin models $\langle D, I \rangle$ with Kripke models $\langle S, R, V \rangle$

Possible worlds $s \in S$	Set of individuals $s_i \in D_i$
Acceptability relation R sRu	Binary predicates $r_{i \rightarrow i \rightarrow o}$ $I r_{i \rightarrow i \rightarrow o}(s_i, u_i) = \top$
Propositional letters p^j Valuation function $s \in V(p^j)$	Unary predicates $p_{i \rightarrow o}^j$ Interpretation function $I p_{i \rightarrow o}^j(s_i) = \top$

$$\begin{aligned}
 \llbracket p^j \rrbracket &= p_{\tau}^j \\
 \llbracket \neg \varphi \rrbracket &= \neg_{\tau \rightarrow \tau} \llbracket \varphi \rrbracket \\
 \llbracket \varphi \vee \psi \rrbracket &= \vee_{\tau \rightarrow \tau \rightarrow \tau} \llbracket \varphi \rrbracket \llbracket \psi \rrbracket \\
 \llbracket \Box \varphi \rrbracket &= \Box_{\tau \rightarrow \tau} \llbracket \varphi \rrbracket
 \end{aligned}$$

$$\begin{aligned}
 \neg_{\tau \rightarrow \tau} &= \lambda A_{\tau} \lambda X_i \neg(A X) \\
 \vee_{\tau \rightarrow \tau \rightarrow \tau} &= \lambda A_{\tau} \lambda B_{\tau} \lambda X_i (A X \vee B X) \\
 \Box_{\tau \rightarrow \tau} &= \lambda A_{\tau} \lambda X_i \forall Y_i (\neg(r_{i \rightarrow i \rightarrow o} X Y) \vee A Y)
 \end{aligned}$$

- ▶ Modal translation of I/O operations in HOL □
- ▶ Åqvist dyadic deontic logic **E** in HOL $\bigcirc(/)$
- ▶ Dyadic deontic logic by Carmo and Jones in HOL $\Box_p, \Box_a, \bigcirc(/), O_p, O_a$

Isabelle/HOL: An infrastructure for deontic reasoning

```

theory IOBoolean
  imports Main

begin

typedcl i (* type for boolean elements *)
type_synonym  $\tau$  = "(i $\Rightarrow$ bool)"
consts N :: "i $\Rightarrow$ i $\Rightarrow$ bool" ("N") (* Nor
consts dis :: "i $\Rightarrow$ i $\Rightarrow$ i" (infixr" $\vee$ "50)
consts con :: "i $\Rightarrow$ i $\Rightarrow$ i" (infixr" $\wedge$ "60)
consts neg :: "i $\Rightarrow$ i" ("¬"[52]53)
consts top :: i ("1")
consts bot :: i ("0")

axiomatization where
COMdis : " $\forall X. \forall Y. (X \vee Y) = (Y \vee X)$ " and
COMcon : " $\forall X. \forall Y. (X \wedge Y) = (Y \wedge X)$ " and
ASSdis : " $\forall X. \forall Y. \forall Z. (X \vee (Y \vee Z)) = (X \vee (Y \vee Z))$ " and
ASScon : " $\forall X. \forall Y. \forall Z. (X \wedge (Y \wedge Z)) = (X \wedge (Y \wedge Z))$ " and
IDEdis : " $\forall X. (X \vee 0) = X$ " and
IDEcon : " $\forall X. (X \wedge 1) = X$ " and
COMPdis : " $\forall X. (X \vee \neg X) = 1$ " and
COMPcon : " $\forall X. (X \wedge (\neg X)) = 0$ " and
Ddiscon : " $\forall X. \forall Y. \forall Z. (X \vee (Y \wedge Z)) = ((X \vee Y) \wedge (X \vee Z))$ " and
Dcondis : " $\forall X. \forall Y. \forall Z. (X \wedge (Y \vee Z)) = ((X \wedge Y) \vee (X \wedge Z))$ "

```

$$\begin{aligned}
 [p^j] &= p_i^j & p^j \in X \\
 [\top] &= \top_i \\
 [\perp] &= \perp_i \\
 [\neg\varphi] &= \neg_{i \rightarrow i}([\varphi]) \\
 [\varphi \vee \psi] &= \vee_{i \rightarrow i}([\varphi][\psi]) \\
 [\varphi \wedge \psi] &= \wedge_{i \rightarrow i}([\varphi][\psi]) \\
 [d_i(N)(\varphi, \psi)] &= (\odot_i(N)_{\tau \rightarrow \tau}\{[\varphi]\})[\psi]
 \end{aligned}$$

$(\varphi, \psi) \in \text{derive}_i^{\text{Fm}(X)}(N)$ iff $V(\psi) \in \text{out}_i^{\mathcal{B}}(N^V, \{V(\varphi)\})$ in all Boolean normative models
 $(N = \langle \mathcal{B}, V, N^V \rangle)$

Faithful embedding

```

definition ordeIOB :: "i ⇒ τ" (infixr "≤" 80) where "X ≤ Y ≡ ((X ∧ Y) = X)"
definition satuIOB :: "τ ⇒ bool" ("Saturated") where
  "Saturated V ≡ ∀X. ∀Y. (((V (X ∨ Y)) → (V X ∨ V Y)) ∧ ((V X ∧ (X ≤ Y)) → V Y))"
definition UpwardIOB :: "τ ⇒ τ" ("Up") where "Up V ≡ λX. (∃Z. (V Z ∧ Z ≤ X))"

definition outI :: "α ⇒ τ ⇒ τ" ("O1<_>")
  where "O1<M>A ≡ λX. ∃U. (∃Y. (∃Z. (A Z ∧ (Z=Y) ∧ M Y U ∧ (U ≤ X) ) ) )"

definition outII :: "α ⇒ τ ⇒ τ" ("OII<_>")
  where "OII<M>A ≡ λX. ∃U. (∃Y. (∃Z. (A Z ∧ (Z≤Y) ∧ M Y U ∧ (U = X) ) ) )"

definition out1 :: "α ⇒ τ ⇒ τ" ("O1<_>")
  where "O1<M>A ≡ λX. ∃U. (∃Y. (∃Z. (A Z ∧ (Z≤Y) ∧ M Y U ∧ (U ≤ X) ) ) )"

definition out2 :: "α ⇒ τ ⇒ τ" ("O2<_>")
  where "O2<M>A ≡ λX. (∀V. ( (Saturated V) ∧ (∀U. (A U → V U))
    → (∃Y. (∃Z. ( (V Y) ∧ (M Y Z) ∧ (Z≤X) ) ) ) )"

definition out3 :: "α ⇒ τ ⇒ τ" ("O3<_>")
  where "O3<M>A ≡ λX. (∀V. ( ((V = Up V) ∧ (∀U. (A U → V U)) ∧ (∀W. (∃Y. (V Y ∧ (M Y W)) → V W))
    → (∃Y. (∃Z. ((Z≤X) ∧ N Y Z ∧ V Y) ) ) )"
  
```

$$\left\{ \begin{array}{l} \text{Sub-rel } R Q \equiv \forall uv. Ruv \rightarrow Quv \\ \text{Close-AND } Q \equiv \forall uvw. (Quv \wedge Quw \rightarrow Qu(v \wedge w)) \\ \text{TCAND } R \equiv \lambda XY. \forall Q. \text{Close-AND } Q \rightarrow (\text{Sub-rel } R Q \rightarrow QXY) \end{array} \right.$$

```

theory outoperation imports IOBoolean
begin

definition Rout :: " $\alpha \Rightarrow \tau \Rightarrow i \Rightarrow i \Rightarrow \text{bool}$ " ("Rout<_ ; >")
  where "Rout<M;A>  $\equiv \lambda Z. \lambda X. \exists U. (\exists Y. (A Z \wedge (Z \leq Y) \wedge M Y U \wedge (U \leq X)))$ "
definition Sub_rel :: " $\alpha \Rightarrow \alpha \Rightarrow \text{bool}$ " where "Sub_rel R Q  $\equiv \forall u v. R u v \rightarrow Q u v$ "

(* OUT1 original *)
definition Close_AND :: " $\alpha \Rightarrow \text{bool}$ " where "Close_AND Q  $\equiv \forall u v w. (Q u v \wedge Q u w \rightarrow (Q u (v \wedge w)))$ "
definition TCAND :: " $\alpha \Rightarrow \alpha$ " where "TCAND R  $\equiv \lambda X Y. \forall Q. \text{Close\_AND } Q \rightarrow (\text{Sub\_rel } R Q \rightarrow Q X Y)$ "
definition outAND :: " $\alpha \Rightarrow \tau \Rightarrow \tau$ " ("OAND<_ ; >") where "OAND<M;A>  $\equiv \lambda X. \exists Y. \text{TCAND } (\text{Rout}<M;A>) Y X$ "
(* OUT2 original *)
definition Close_OR :: " $\alpha \Rightarrow \text{bool}$ " where "Close_OR Q  $\equiv \forall u v w. (Q v u \wedge Q w u \rightarrow (Q (v \vee w) u))$ "
definition TCOR :: " $\alpha \Rightarrow \alpha$ " where "TCOR R  $\equiv \lambda X Y. \forall Q. \text{Close\_OR } Q \rightarrow (\text{Sub\_rel } R Q \rightarrow Q X Y)$ "
definition outOR :: " $\alpha \Rightarrow \tau \Rightarrow \tau$ " ("OR<_ ; >") where "OR<M;A>  $\equiv \lambda X. \exists Y. \text{TCOR } (\text{Rout}<M;A>) Y X$ "
definition outORAND :: " $\alpha \Rightarrow \tau \Rightarrow \tau$ " ("ORAND<_ ; >")
  where "ORAND<M;A>  $\equiv \lambda X. \exists Y. \text{TCAND } (\text{TCOR } (\text{Rout}<M;A>)) Y X$ "

```

Isabelle/HOL: I/O proof systems in HOL

```
(*Derive2-0b*)
definition derSIWOORAND :: " $\alpha \Rightarrow \alpha$ " ("derSIWOORAND<_>")
  where "derSIWOORAND<M>  $\equiv$  TCAND (TCOR (TCWO (TCSI (M))))"

(*Derive3-0b*)
definition derSIWOCT :: " $\alpha \Rightarrow \alpha$ " ("derSIWOCT<_>")
  where "derSIWOCT<M>  $\equiv$  TCCT (TCWO (TCSI (M)))"
definition derSIWOCTAND :: " $\alpha \Rightarrow \alpha$ " ("derSIWOCTAND<_>")
  where "derSIWOCTAND<M>  $\equiv$  TCAND (TCCT (TCWO (TCSI (M))))"

(*Derive4-0b*)
definition derSIWOCTORAND :: " $\alpha \Rightarrow \alpha$ " ("derSIWOCTORAND<_>")
  where "derSIWOCTORAND<M>  $\equiv$  TCAND (TCOR (TCCT (TCWO (TCSI (M)))))"
```

```
lemma "Close_AND (TCAND N)" unfolding Defst TCAND_def
  by metis

lemma "(M a b  $\vee$  ( $\exists y. M y b \wedge (a \leq y)$ ))  $\longrightarrow$  derSI<M> a b"
using Sub_rel_def Close_SI_def TCSI_def
  unfolding Defst and Defs derSI_def
  by metis

(*OUT1 completeness*)
lemma "( $\bigcirc_1 < N; ((\lambda X. X = a)) > y$ )  $\longrightarrow$  derSIW0<N> a y"
using Sub_rel_def Close_SI_def Close_W0_def TCSI_def TCW0_def
  unfolding Defst and Defs derSI_def Sub_rel_def TCW0_def TCSI_def
  by metis
```

Normative Reasoning

- Discursive input/output logic: **Detachment** vs Quantification
 - Input/output logic for permission: Removing AND rule
 - Input/output logic for obligation: Adding AND rule
 - **Semantic unification**: Integrating input/output logic into Kratzerian framework
- Normative reasoning + **Preferences** /Normality
 - A compositional theory of conditional obligation and permission
- Algebraic method: I/O framework on top of any **abstract logic**
 - Input/output methodology: Secretarial assistant ($\mathcal{A} = \langle \mathcal{L}, C \rangle$)

Logic Engineering

- A dataset for normative reasoning: **LogiKEy** methodology
 - Faithful embedding of some deontic logics in HOL
 - Isabelle/HOL: An infrastructure for deontic reasoning

Normative Reasoning

- Discursive input/output logic: **Detachment** vs Quantification
 - Input/output logic for permission: Removing AND rule
 - Input/output logic for obligation: Adding AND rule
 - **Semantic unification**: Integrating input/output logic into Kratzerian framework
- Normative reasoning + **Preferences** /Normality
 - A compositional theory of conditional obligation and permission
- Algebraic method: I/O framework on top of any **abstract logic**
 - Input/output methodology: Secretarial assistant ($\mathcal{A} = \langle \mathcal{L}, C \rangle$)

Logic Engineering

- A dataset for normative reasoning: **LogiKEy** methodology
 - Faithful embedding of some deontic logics in HOL
 - Isabelle/HOL: An infrastructure for deontic reasoning

Normative Reasoning

- Discursive input/output logic: **Detachment** vs Quantification
 - Input/output logic for permission: Removing AND rule
 - Input/output logic for obligation: Adding AND rule
 - **Semantic unification**: Integrating input/output logic into Kratzerian framework
- Normative reasoning + **Preferences** /Normality
 - A compositional theory of conditional obligation and permission
- Algebraic method: I/O framework on top of any **abstract logic**
 - Input/output methodology: Secretarial assistant ($\mathcal{A} = \langle \mathcal{L}, C \rangle$)

Logic Engineering

- A dataset for normative reasoning: **LogiKEy** methodology
 - Faithful embedding of some deontic logics in HOL
 - Isabelle/HOL: An infrastructure for deontic reasoning

Normative Reasoning

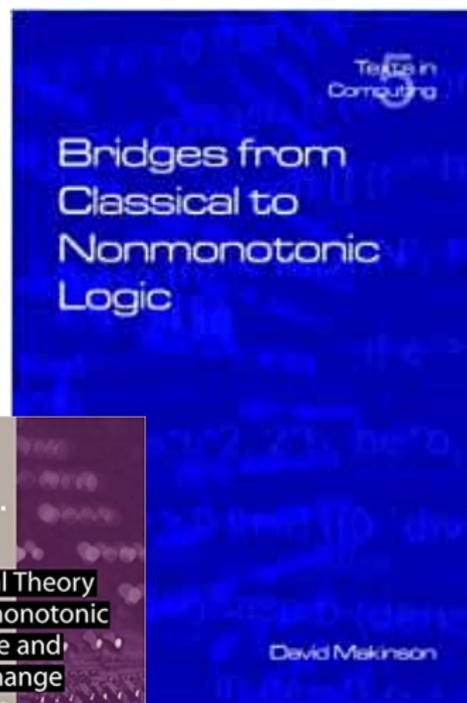
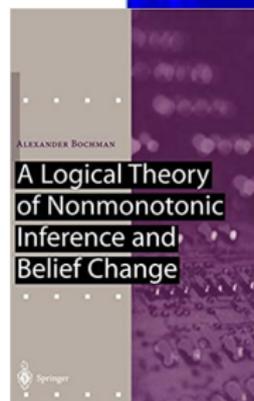
- Discursive input/output logic: **Detachment** vs Quantification
 - Input/output logic for permission: Removing AND rule
 - Input/output logic for obligation: Adding AND rule
 - **Semantic unification**: Integrating input/output logic into Kratzerian framework
- Normative reasoning + **Preferences** /Normality
 - A compositional theory of conditional obligation and permission
- Algebraic method: I/O framework on top of any **abstract logic**
 - Input/output methodology: Secretarial assistant ($\mathcal{A} = \langle \mathcal{L}, C \rangle$)

Logic Engineering

- A dataset for normative reasoning: **LogiKEy** methodology
 - Faithful embedding of some deontic logics in HOL
 - Isabelle/HOL: An infrastructure for deontic reasoning

Future Work: Normative reasoning

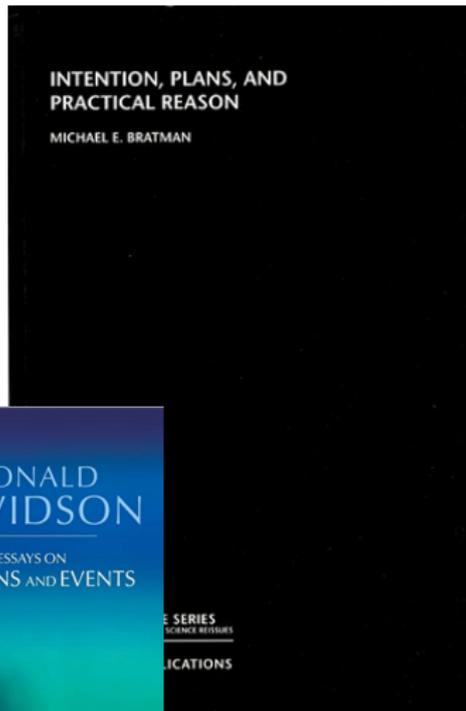
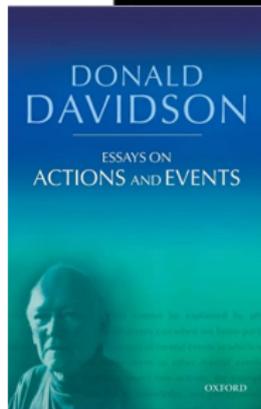
- ▶ Adding rules to logic:
A full characterization
- ▶ Logic without reflexivity:
Application in other domains such
as causality
- ▶ Credulous (brave) inference:
Belief change



Future Work: Practical reasoning

How do norms interact with informational modalities such as beliefs and knowledge, and motivational modalities such as intentions and desires? (Ten Problems of Deontic Logic and Normative Reasoning in Computer Science)

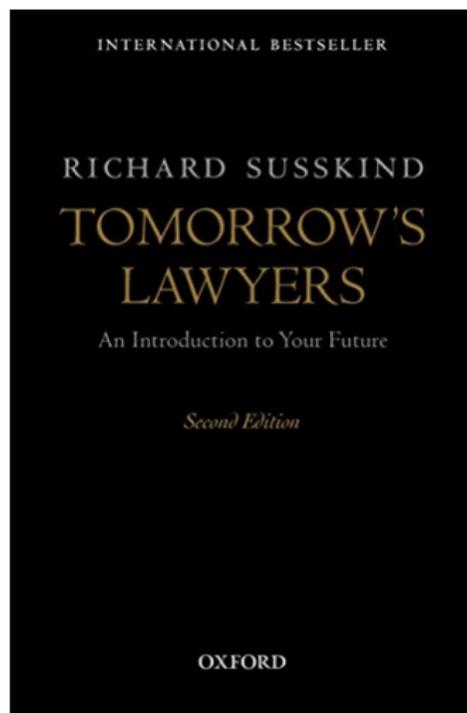
- ▶ Anankastic conditionals (means-end reasoning)
- ▶ Rational architecture: BOID
- ▶ Human-Computer Interaction



Future Work: Online legal guidance systems

“what if we, as lawyers, could make our knowledge and expertise available through a wide range of online legal services, whether for the drafting of documents or for the resolution of disputes?” (Susskind)

- ▶ Improving normative expressivity of the implemented logics in HOL
- ▶ A domain for individuals
- ▶ Logic and ontology
- ▶ Higher-order deontic logic



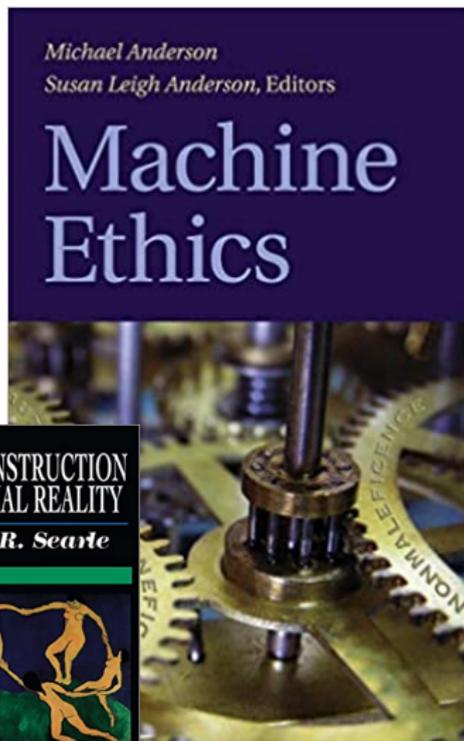
Future Work: Autonomous vehicles

Can ethical frameworks and rules derived for human behavior be implemented as control algorithms in automated vehicles? (Implementable Ethics for Autonomous Vehicles)

- ▶ Cost Functions and consequentialism
- ▶ Constraints and deontological Ethics
- ▶ Norms and preferences
- ▶ How does the deontological approach fare with uncertainty?

(Normative) multi-agent systems

- ▶ Constitutive norms, Regulative norms
- ▶ Privacy policies and Knowledge management



Many Thanks!

#سیاسگزارم