# Experimenting with Carmo and Jones' DDL

Lavanya Singh

March 31, 2021

## Contents

**theory** *carmojones-DDL*
  **imports**
    *Main*

**begin**

Referencing Benzmuller and Parent's implementation [1]

This theory contains the axiomatization of the system and some useful abbreviations.

# 1 System Definition

## 1.1 Definitions

This section contains definitions and constants necessary to construct a DDL model.

**typedecl** $i$ — i is the type for a set of possible worlds."

**type-synonym** $t = (i \Rightarrow bool)$
— t represents a set of DDL formulas.
— this set is defined by its truth function, mapping the set of worlds to the formula set's truth value.

— accessibility relations map a set of worlds to:
**consts** $av{::}i \Rightarrow t$ — actual versions of that world set
  — these worlds represent what is "open to the agent"
  — for example, the agent eating pizza or pasta for dinner might constitute two different actual worlds

**consts** $pv{::}i \Rightarrow t$ — possible versions of that world set
  — these worlds represent was was "potentially open to the agent"
  — for example, what someone across the world eats for dinner might constitute a possible world, — since the agent has no control over this

**consts** $ob{::}t \Rightarrow (t \Rightarrow bool)$ — set of propositions obligatory in this "context"
  — ob(context)(term) is True if t is obligatory in the context

**consts** $cw{::}i$ — current world

## 1.2 Axiomatization

This subsection contains axioms. Because the embedding is semantic, these are just constraints on models.

This axiomatization comes from [2] p6 and the HOL embedding defined in Benzmuller and Parent

**axiomatization where**
*ax-3a*: $\forall\, w.\exists\, x.\ av(w)(x)$
 — every world has some actual version

**and** *ax-4a*: $\forall\, w\ x.\ av(w)(x) \longrightarrow pv(w)(x)$
— all actual versions of a world are also possible versions of it

**and** *ax-4b*: $\forall\, w.\ pv(w)(w)$
— every world is a possible version of itself

**and** *ax-5a*: $\forall\, X.\neg ob(X)(\lambda w.\ \textit{False})$
— in any arbitrary context X, something will be obligatory

**and** *ax-5b*: $\forall\, X\ Y\ Z.\ (\forall\, w.\ ((X(w) \wedge Y(w)) \longleftrightarrow (X(w) \wedge Z(w)))) \longrightarrow (ob(X)(Y) \longleftrightarrow ob(X)(Z))$ — note that X(w) denotes w is a member of X
— X, Y, and Z are sets of formulas
— If X ∩ Y = X ∩ Z then the context X obliges Y iff it obliges Z

— ob(X)(λ w. Fw) can be read as F ∈ ob(X)

**and** *ax-5c2*: $\forall\, X\ Y\ Z.\ (((\exists\, w.\ (X(w) \wedge Y(w) \wedge Z(w))) \wedge ob(X)(Y) \wedge ob(X)(Z))) \longrightarrow ob(X)(\lambda w.\ Y(w) \wedge Z(w))$

**and** *ax-5d*: $\forall\, X\ Y\ Z.\ ((\forall\, w.\ Y(w) \longrightarrow X(w)) \wedge ob(X)(Y) \wedge (\forall\, w.\ X(w) \longrightarrow Z(w)))$

 $\longrightarrow ob(Z)(\lambda w.(Z(w) \wedge \neg X(w)) \vee Y(w))$
— If some subset Y of X is in ob(X) then in a larger context Z, any obligatory proposition must either be in Y or in Z-X

**and** *ax-5e*: $\forall\, X\ Y\ Z.\ ((\forall\, w.\ Y(w) \longrightarrow X(w)) \wedge ob(X)(Z) \wedge (\exists\, w.\ Y(w) \wedge Z(w))) \longrightarrow ob(Y)(Z)$
— If Z is obligatory in context X, then Z is obligatory in a subset of X called Y, if Z shares some elements with Y

## 1.3   Abbreviations

These abbreviations are defined in @citeBenzmullerParent p9

These are all syntactic sugar for HOL expressions, so evaluating these symbols will be light-weight

— propositional logic symbols
**abbreviation** *ddlneg*::$t \Rightarrow t$ ($\neg$)
  **where** $\neg A \equiv \lambda w.\ \neg A(w)$
**abbreviation** *ddlor*::$t \Rightarrow t \Rightarrow t$ (-∨-)
  **where** $A \vee B \equiv \lambda w.\ (A(w) \vee B(w))$
**abbreviation** *ddland*::$t \Rightarrow t \Rightarrow t$ (-∧-)
  **where** $A \wedge B \equiv \lambda w.\ (A(w) \wedge B(w))$
**abbreviation** *ddlif*::$t \Rightarrow t \Rightarrow t$ (-→-)
  **where** $A \rightarrow B \equiv (\lambda w.\ A(w) \longrightarrow B(w))$

**abbreviation** *ddlequiv::t⇒t⇒t* (-≡-)
  **where** $(A \equiv B) \equiv ((A \rightarrow B) \land (B \rightarrow A))$

— modal operators
**abbreviation** *ddlbox::t⇒t* ($\square$)
  **where** $\square\ A \equiv \lambda w. \forall y.\ A(y)$
**abbreviation** *ddldiamond::t ⇒ t* ($\lozenge$)
  **where** $\lozenge A \equiv \neg(\square(\neg A))$

— O{B|A} can be read as "B is obligatory in the context A"
**abbreviation** *ddlob::t⇒t⇒t* (O{-|-})
  **where** $O\{B|A\} \equiv \lambda\ w.\ ob(A)(B)$

— modal symbols over the actual and possible worlds relations
**abbreviation** *ddlboxa::t⇒t* ($\square_a$)
  **where** $\square_a A \equiv \lambda x. \forall y.\ (\neg\ av(x)(y) \lor A(y))$
**abbreviation** *ddldiamonda::t⇒t* ($\lozenge_a$)
  **where** $\lozenge_a A \equiv \neg(\square_a(\neg A))$
**abbreviation** *ddlboxp::t⇒t* ($\square_p$)
  **where** $\square_p A \equiv \lambda x. \forall y.\ (\neg\ pv(x)(y) \lor A(y))$
**abbreviation** *ddldiamondp::t⇒t* ($\lozenge_p$)
  **where** $\lozenge_p A \equiv \neg(\square_a(\neg A))$

— obligation symbols over the actual and possible worlds
**abbreviation** *ddloba::t⇒t* ($O_a$)
  **where** $O_a\ A \equiv \lambda x.\ ob(av(x))(A) \land (\exists y.(av(x)(y) \land \neg A(y)))$
**abbreviation** *ddlobp::t⇒t* ($O_p$)
  **where** $O_p\ A \equiv \lambda x.\ ob(pv(x))(A) \land (\exists y.(pv(x)(y) \land \neg A(y)))$

— syntactic sugar for a "monadic" obligation operator
**abbreviation** *ddltrue::t* ($\top$)
  **where** $\top \equiv \lambda w.\ True$
**abbreviation** *ddlob-normal::t⇒t* (O {-})
  **where** $(O\ \{A\}) \equiv (O\{A|\top\})$

— validity
**abbreviation** *ddlvalid::t⇒bool* ($\models$-)
  **where** $\models A \equiv \forall w.\ A\ w$
**abbreviation** *ddlvalidcw::t⇒bool* ($\models_c$-)
  **where** $\models_c A \equiv A\ cw$

## 1.4 Consistency

Consistency is so easy to show in Isabelle!

**lemma** *True* **nitpick** [*satisfy,user-axioms,show-all,format=2*] ⟨*proof*⟩

**end**

**theory** *carmojones-DDL-completeness* **imports** *carmojones-DDL*

**begin**

This theory shows completeness for this logic with respect to the models presented in carmojonesDDl.thy.

# 2 Inference Rules

## 2.1 Basic Inference Rules

These inference rules are common to most modal and propostional logics

**lemma** *modus-ponens*: **assumes** $\models A$ **assumes** $\models (A \rightarrow B)$
  **shows** $\models B$
  $\langle proof \rangle$

**lemma** *nec*: **assumes** $\models A$ **shows** $\models (\Box A)$
  $\langle proof \rangle$

**lemma** *nec-a*: **assumes** $\models A$ **shows** $\models (\Box_a A)$
  $\langle proof \rangle$
**lemma** *nec-p*: **assumes** $\models A$ **shows** $\models (\Box_p A)$
  $\langle proof \rangle$

## 2.2 Fancier Inference Rules

These are new rules that Carmo and Jones introduced for this logic.

**lemma** *Oa-boxaO*:
  **assumes** $\models (B \rightarrow ((\neg(\Box((O_a\ A) \rightarrow ((\Box_a w) \wedge O\{A|w\}))))))$
  **shows** $\models (B \rightarrow (\neg(\Diamond(O_a\ A))))$
  $\langle proof \rangle$
**lemma** *Oa-boxpO*:
  **assumes** $\models (B \rightarrow ((\neg(\Box((O_p\ A) \rightarrow ((\Box_p w) \wedge O\{A|w\}))))))$
  **shows** $\models (B \rightarrow (\neg(\Diamond(O_p\ A))))$
  $\langle proof \rangle$

B and A must not contain w. not sure how to encode that requirement. one option is to define a new free variables predicate and use that, but that requires a deeper embedding than I have. If Benzmuller and Parent can survive without these inference rules, so can I

# 3 Axioms

## 3.1 Box

— $\Box$ is an S5 modal operator, which is where these axioms come from.

**lemma** *K*:
  **shows** $\models ((\Box(A \rightarrow B)) \rightarrow ((\Box A) \rightarrow (\Box B)))$
  $\langle proof \rangle$

**lemma** *T*:
  **shows** $\models ((\Box A) \rightarrow A)$
  $\langle proof \rangle$

**lemma** *5*:
  **shows** $\models ((\Diamond A) \rightarrow (\Box(\Diamond A)))$
  $\langle proof \rangle$

## 3.2   O

This characterization of O comes from Carmo and Jones p 593

**lemma** *O-diamond*:
  **shows** $\models (O\{A|B\} \rightarrow (\Diamond(B \wedge A)))$
  $\langle proof \rangle$

**lemma** *O-C*:
  **shows** $\models (((\Diamond(A \wedge (B \wedge C))) \wedge (O\{B|A\} \wedge O\{C|A\})) \rightarrow (O\{B \wedge C|A\}) )$
  $\langle proof \rangle$

**lemma** *O-SA*:
  **shows** $\models (((\Box(A \rightarrow B)) \wedge ((\Diamond(A \wedge C)) \wedge O\{C|B\})) \rightarrow (O\{C|A\}))$
  $\langle proof \rangle$

**lemma** *O-REA*:
  **shows** $\models ((\Box(A \equiv B)) \rightarrow (O\{C|A\} \equiv O\{C|B\}))$
  $\langle proof \rangle$

**lemma** *O-contextual-REA*:
  **shows** $\models ((\Box(C \rightarrow (A \equiv B))) \rightarrow (O\{A|C\} \equiv O\{B|C\}))$
  $\langle proof \rangle$

**lemma** *O-nec*:
  **shows** $\models (O\{B|A\} \rightarrow (\Box O\{B|A\}))$
  $\langle proof \rangle$

**lemma** *ax-5b″*:
  **shows** $ob\ X\ Y \longleftrightarrow ob\ X\ (\lambda z.\ (Y\ z) \wedge (X\ z))$
  $\langle proof \rangle$

**lemma** *O-to-O*:
  **shows** $\models (O\{B|A\} \rightarrow O\{(A \rightarrow B)|\top\})$
$\langle proof \rangle$

## 3.3 Possible Box

— $\Box_p$ is a KT modal operator.

**lemma** *K-boxp*:
  **shows** $\models((\Box_p(A \to B)) \to ((\Box_p A) \to (\Box_p B)))$
  $\langle proof \rangle$
**lemma** *T-boxp*:
  **shows** $\models((\Box_p A) \to A)$
  $\langle proof \rangle$

## 3.4 Actual Box

— $\Box_a$ is a KD modal operator.

**lemma** *K-boxa*:
  **shows** $\models((\Box_a(A \to B)) \to ((\Box_a A) \to (\Box_a B)))$
  $\langle proof \rangle$
**lemma** *D-boxa*:
  **shows** $\models((\Box_a A) \to (\Diamond_a A))$
  $\langle proof \rangle$

## 3.5 Relations Between the Modal Operators

— Relation between $\Box$, $\Box_a$, and $\Box_p$.

**lemma** *box-boxp*:
  **shows** $\models((\Box A) \to (\Box_p A))$
  $\langle proof \rangle$
**lemma** *boxp-boxa*:
  **shows** $\models((\Box_p A) \to (\Box_a A))$
  $\langle proof \rangle$
**lemma** *not-Oa*:
  **shows** $\models((\Box_a A) \to ((\neg(O_a\ A)) \land (\neg(O_a\ (\neg A)))))$
  $\langle proof \rangle$
**lemma** *not-Op*:
**shows** $\models((\Box_p A) \to ((\neg(O_p\ A)) \land (\neg(O_p\ (\neg A)))))$
  $\langle proof \rangle$
**lemma** *equiv-Oa*:
  **shows** $\models((\Box_a(A \equiv B)) \to ((O_a\ A) \equiv (O_a\ B)\ ))$
  $\langle proof \rangle$
**lemma** *equiv-Op*:
  **shows** $\models((\Box_p(A \equiv B)) \to ((O_p\ A) \equiv (O_p\ B)\ ))$
  $\langle proof \rangle$
**lemma** *factual-detach-a*:
  **shows** $\models(((O\{B|A\} \land (\Box_a A)) \land ((\Diamond_a B) \land (\Diamond_a(\neg B)))) \to (O_a\ B))$
  $\langle proof \rangle$
**lemma** *factual-detach-p*:
  **shows** $\models(((O\{B|A\} \land (\Box_p A)) \land ((\Diamond_p B) \land (\Diamond_p(\neg B)))) \to (O_p\ B))$
  $\langle proof \rangle$

**end**

**theory** *categorical-imperative-1* **imports** *carmojones-DDL-completeness*

**begin**

# 4 The Categorical Imperative

## 4.1 Simple Formulation of the Kingdom of Ends

This is my first attempt at formalizing the concept of the Kingdom of Ends

NOTE: this attempt revealed a bug in my embedding. I've included it as an artifact, but none of these theorems hold anymore (hence the oops).

**abbreviation** *ddlpermissable*::$t{\Rightarrow}t$ (*P-*)
  **where** $(P\ A) \equiv (\neg(O\ \{\neg A\}))$
— This operator represents permissibility
— Will be useful when discussing the categorical imperative
— Something is permissible if it is not prohibited
— Something is prohibited if its negation is obligatory


**lemma** *kingdom-of-ends-1*:
  **shows** $\models ((O\ \{A\}) \to (\Box\ (P\ A)))$
  ⟨*proof*⟩


**lemma** *kingdom-of-ends-2*:
  **shows** $\models ((\Box\ (P\ A)) \to (O\ \{A\}))$
  ⟨*proof*⟩


**lemma** *permissible-to-ob*:
  **shows** $\models ((P\ A) \to (O\ \{A\}))$
  ⟨*proof*⟩

**lemma** *weaker-permissible-to-ob*:
  **shows** $\models ((\Diamond\ (P\ A)) \to O\ \{A\})$
   ⟨*proof*⟩

**lemma** *contradictory-obligations*:
  **shows** $\models (\neg\ ((O\ \{A\}) \land (O\ \{\neg\ A\})))$
  ⟨*proof*⟩

Sidebar: the above theorem is really intuitive - it seems like we wouldn't want contradictory things to be obligatory in any logic. But for some reason, not only is it not a theorem of Carmo and Jones' logic, it actually implies some strange conclusions, including that everything is either permissible or obligatory. It's not clear to me from a semantic perspective why this would be the case. In fact this theorem seems like a desirable property. Potential

avenue for exploration

Did some debugging. What was the problem? A misplaced parentheses in the definition of ax5b that led to a term being on the wrong side of an implication. Computer Science :(

After the debugging, all of this is no longer true! On to the next attempt :)

**end**

**theory** *categorical-imperative-naive* **imports** *carmojones-DDL-completeness*

**begin**

# 5 The Categorical Imperative

## 5.1 Simple Formulation of the Formula of Universal Law

This is my second attempt at formalizing the Formula of Universal Law

**abbreviation** *ddlpermissable*::$t \Rightarrow t$ (*P-*)
  **where** $(P\ A) \equiv (\neg(O\ \{\neg A\}))$
— This operator represents permissibility
— Will be useful when discussing the categorical imperative
— Something is permissible if it is not prohibited
— Something is prohibited if its negation is obligatory

Let's consider a naive reading of the Formula of Universal Law (FUL). From the Groundwork, 'act only in accordance with that maxim through which you can at the same time will that it become a universal law'. What does this mean in DDL? One interpretation is if A is not necessarily permissible, then its negation is obligated.

**axiomatization where**
*FUL-1*: $\models ((\neg(\Box\ (P\ A))) \rightarrow (O\ \{(\neg A)\}))$

## 5.2 Basic Tests

**lemma** *True* **nitpick** [*satisfy*,*user-axioms*,*format=2*] ⟨*proof*⟩

**lemma** *something-is-obligatory*:
  **shows** $\forall\ w.\ \exists\ A.\ O\ \{A\}\ w$
  **nitpick** [*user-axioms*]
  ⟨*proof*⟩

**lemma** *something-is-obligatory-2*:
  **shows** $\forall\ w.\ \exists\ A.\ O\ \{A\}\ w$
  **nitpick** [*user-axioms*, *falsify=false*]
  ⟨*proof*⟩

**lemma** *something-is-obligatory-relaxed*:
  **shows** $\exists\ A\ w.\ O\ \{A\}\ w$
  **nitpick** [*user-axioms*]
  ⟨*proof*⟩

**lemma** *something-is-obligatory-relaxed-2*:
  **shows** $\exists\ A\ w.\ O\ \{A\}\ w$
  **nitpick** [*user-axioms*, *falsify=false*]
  ⟨*proof*⟩

## 5.3 Specifying the Model

Let's specify the model. What if we add something impermissible?

**consts** *M*::*t*
**abbreviation** *murder-wrong*::*bool* **where** *murder-wrong* ≡ ⊨($O$ {¬ *M*})

**lemma** *something-is-obligatory-2*:
  **assumes** *murder-wrong*
  **shows** ∀ *w*. ∃ *A*. *O* {*A*} *w*
  ⟨*proof*⟩

**abbreviation** *poss-murder-wrong*::*bool* **where** *poss-murder-wrong* ≡ ⊨(◇ (*O* {¬ *M*}))

**lemma** *wrong-if-posibly-wrong*:
  **assumes** *poss-murder-wrong*
  **shows** *murder-wrong*
  ⟨*proof*⟩

Let's try an even weaker assumption: Not everyone can lie.

**typedecl** *person*
**consts** *lies*::*person*⇒*t*
**consts** *me*::*person*

**lemma** *breaking-promises*:
  **assumes** ¬ (∀ *x*. *lie*(*x*) *cw*) ∧ (*lie*(*me*) *cw*)
  **shows** (*O* {¬ (*lie*(*me*))}) *cw*
  **nitpick** [*user-axioms*]
  ⟨*proof*⟩

**lemma** *universalizability*:
  **assumes** ⊨ *O* {(*lie*(*me*))}
  **shows** ∀ *x*. ⊨ (*O* {(*lie*(*x*))})
  **nitpick** [*user-axioms*] ⟨*proof*⟩

## 5.4 Consistent Sentences

The above section tested validity. We might also be interested in some weaker properties

Let's test whether certain sentences are consistent - can we find a model that makes them true?

**lemma** *permissible*:
  **fixes** *A*
  **shows** ((¬ (*O* {*A*})) ∧ (¬ (*O* {¬ *A*}))) *w*
  **nitpick** [*user-axioms*, *falsify=false*] ⟨*proof*⟩

**lemma** *conflicting-obligations*:

11

**fixes** *A*
**shows** (*O* {*A*} ∧ *O* {¬ *A*}) *w*
**nitpick** [*user-axioms*, *falsify=false*] ⟨*proof*⟩

## 5.5 Metaethical Tests

**lemma** *FUL-alternate*:
  **shows** ⊨ ((◊ (*O* {¬ *A*})) → (*O* {¬ *A*}))
  ⟨*proof*⟩

**lemma** *arbitrary-obligations*:
  **fixes** *A*::*t*
  **shows** *O* {*A*} *w*
  **nitpick** [*user-axioms=true*] ⟨*proof*⟩

**lemma** *removing-conflicting-obligations*:
  **assumes** ∀ *A*. ⊨ (¬ (*O* {*A*} ∧ *O* {¬ *A*}))
  **shows** *True*
  **nitpick** [*satisfy*,*user-axioms*,*format=2*] ⟨*proof*⟩

**lemma** *implied-contradiction*:
  **fixes** *A*::*t*
  **fixes** *B*::*t*
  **assumes** ⊨(¬ (*A* ∧ *B*))
  **shows** ⊨(¬ (*O* {*A*} ∧ *O* {*B*}))
  **nitpick** [*user-axioms*]
⟨*proof*⟩

**lemma** *distribute-obligations-if*:
  **assumes** ⊨ *O* {*A* ∧ *B*}
  **shows** ⊨ (*O* {*A*} ∧ *O* {*B*})
  **nitpick** [*user-axioms*, *falsify=true*, *verbose*]
  ⟨*proof*⟩

**lemma** *distribute-boxes*:
  **assumes** ⊨( □(*A* ∧ *B*))
  **shows** ⊨ ((□*A*) ∧ (□*B*))
  ⟨*proof*⟩

**lemma** *distribute-obligations-onlyif*:
  **assumes**  ⊨ (*O* {*A*} ∧ *O* {*B*})
  **shows** ⊨ *O* {*A* ∧ *B*}
  **nitpick** [*user-axioms*] ⟨*proof*⟩

**lemma** *ought-implies-can*:
  **shows** ∀ *A*. ⊨ (*O* {*A*} → (◊*A*))
  ⟨*proof*⟩

**end**

# References

[1] C. Benzmüller, A. Farjami, and X. Parent. Faithful semantical embedding of a dyadic deontic logic in HOL. *CoRR*, abs/1802.08454, 2018.

[2] J. Carmo and A. Jones. Completeness and decidability results for a logic of contrary-to-duty conditionals. *J. Log. Comput.*, 23:585–626, 2013.