

THE LOGIC OF PERMISSION AND OBLIGATION IN THE
FRAMEWORK OF ALX.3:
HOW TO AVOID THE PARADOXES OF DEONTIC LOGICS.

Zhisheng HUANG and Michael MASUCH

Abstract

Standard deontic logic features fairly serious so-called “paradoxes” (technically: counterintuitive validities). Much energy in deontic logic has been spent on avoiding these “paradoxes”. We suggest a reformulation of deontic logic in terms of a multi-agent logic, ALX.3, where a “super-agent” (think of a legislature) lays down the law of the land and other agents have to follow its rules. In particular, obligations are reformulated in terms of “preferences” of the superagent. We can show that our approach avoids the classical “paradoxes” of deontic logic, thanks to the properties of the preference operator of ALX.3.

1. *Introduction*

Deontic logic is a branch of modal logic for reasoning about social norms by means of modal operators denoting states of *obligation* (written as **O**), *permission* (written as **P**), and *prohibition* (written as **F**, from “forbidden”). Deontic logic has many potential applications in areas such as law, computer science, and sociology, but the standard versions of deontic logic suffer from serious “paradoxes” — paradoxes not in a technical sense but in the sense of counterintuitive validities — that have kept deontic logic from living up to its full potential. Here are some examples of such validities:

- Ross' Paradox: $O\phi \rightarrow O(\phi \vee \psi)$.
This validity would justify propositions such as: “if one is obliged to mail the letter, one is obliged to either mail the letter or burn it”.
- Penitent's Paradox: $F\phi \rightarrow F(\phi \wedge \psi)$.
This validity would justify propositions such as: “if it is forbidden to commit a crime, then it is also forbidden to commit a crime and do penitence for it”.

- Good Samaritan Paradox: $\phi \rightarrow \psi \Rightarrow \mathbf{O}\phi \rightarrow \mathbf{O}\psi$

This validity implies, for example, that the fact that a good Samaritan helps a victim if the victim has been robbed, together with the obligation to help victims after a robbery, implies the obligation to rob the victim in the first place.

The underlying reason for these paradoxes is technical: at least one of the modal operators is introduced as a primitive modality in the context of a normal modal logic (roughly speaking: a logic with Kripke-style semantics and no absurd worlds). The Good Samaritan paradox is a consequence of the monotonicity of normal modal logics, while Ross's paradox is just a special case of this monotonicity with the appropriate substitutions made. Penitent's paradox, is, in fact, the dual version of Ross's paradox with the obligation operator replaced by the prohibition operator.

Not only deontic logics suffer from the side effects of normality; similar problems arise in epistemic logics, where the epistemic operators are also acting as normal modalities (logical omniscience is a typical side-effect of normality in epistemic logics).

Various attempts have been made to circumvent the paradoxes. Anderson, for example [1], defines the prohibition operator with the help of a propositional constant V as follows:

$$\mathbf{F}\phi \stackrel{\text{def}}{\Leftrightarrow} \Box(\phi \rightarrow V)$$

where the meaning of the constant V is "liable to sanction or punishment", and the box-operator assumes the standard meaning of alethic modal logic: a state ϕ is forbidden if and only if the state ϕ necessarily implies sanctioning the agent. As it turned out however, Anderson's approach missed its goal [10]. The Good Samaritan Paradox does not go away, and neither do Ross' and Penitent's paradoxes—being special cases of the former one.

Inspired by Anderson's reduction to alethic modal logic, J.-J. Meyer proposes another solution in [9]. There, Meyer uses propositional dynamic logic, still employing Anderson's special violation atom V . One of the consequences of the use of dynamic logic is the distinction between propositions and actions. In Meyer's approach, the deontic operators are defined through dynamic expressions as follows: $\mathbf{F}\alpha \stackrel{\text{def}}{\Leftrightarrow} [a]V$, $\mathbf{P}\alpha \stackrel{\text{def}}{\Leftrightarrow} \neg \mathbf{F}\alpha$ and $\mathbf{O}\alpha \stackrel{\text{def}}{\Leftrightarrow} \mathbf{F}(\neg\alpha)$. As pointed out in [10], Ross' paradox remains in Meyer's reformulation.

Inspired by Anderson's idea of reformulating deontic logics, we propose an approach that uses the preference operator instead of the logical implica-

tion or dynamic actions. In other words, we try to define the **F**-operator in terms of the preference operator in a multi-agent logic. We assume that there exists a superagent, next to other agents, who lays down the law of the land. Informally, a state is forbidden for an agent i if and only if the superagent assumes a preference against that state. We find that the preference operator in ALX is suitable to fulfill such a task [5, 7].

There are three versions of ALX logics. The first version is a propositional action logic for agents with bounded rationality. The second version adds a first-order description language, while the third version introduces multiple agents.

The paper is organized as follows: in section 2, we briefly review the preference operator in ALX and its semantics, discuss the formal properties of preferences, and define goodness and badness operators in terms of the preference operator. Then, in the section 3, we reformulate the deontic logic in terms of the new operators, discuss the formal properties of this new logic and show how those “paradoxes” can be avoided. Section 4 has concluding remarks and discusses future directions.

2. ALX logics and its Preference Operator

2.1. ALX Logics

In [5, 7], we propose a modal action logic that combines ideas from H.A. Simon’s *bounded rationality*, S. Kripke’s *possible world semantics*, G. H. von Wright’s *preference logic*, Pratt’s *dynamic logic*, Stalnaker’s *minimal change*, and more recent approaches to *update semantics*. ALX (the x ’s action logic) is sound, complete, and decidable, making it the first complete logic for two-place preference operators. ALX avoids important drawbacks of other action logics, especially the counterintuitive necessitation rule for goals (every theorem must be a goal) and the equally counterintuitive closure of goals under logical implication.

In this paper, we use ALX’s preference operator to reformulate deontic logic. In particular, we use the preference operator of ALX.3. ALX.3 is discussed in detail in [8]. A short overview is given in the appendix. In the following, we employ a simplified version of ALX.3.

2.2. Preferences

Preferences provide the basis for rational action in ALX. Following von Wright [12], a preference statement is understood as a statement about situations. For example, the statements that “I prefer oranges to apples” is interpreted as the fact that “I prefer the states in which I have an orange to the

states in which I have an apple.” Following von Wright again, we assume that an agent who claims to prefer oranges to apples should prefer a situation where he has an orange but *no* apple to a situation where he has an apple but *no* orange, call it *the conjunction expansion principle*. Preferences are expressed via two-place modal operators; if the agent prefers the proposition ϕ to the proposition ψ , we write $\phi \mathbf{P}_i \psi$.

Normally, the meaning of a preference statement is context dependent, even when this is not made explicit. An agent may claim to prefer an apple to an orange — and actually mean it — but he may prefer an orange to an apple later — perhaps because then he already had an apple. To capture this context dependency, we borrow the notion of minimal change from Stalnaker’s approach to conditionals [11]. The idea is to apply the conjunction expansion principle only to situations that are minimally different from the agent’s present situation — just as different as they really need to be in order to make the propositions true about which preferences are expressed. We introduce a binary function, cw , to the semantics that determines a set of “closest” states relative to a given state, such that the new states fulfill some specified conditions, but resemble the old state as much as possible in all other respects. For situations (sets of states), we apply cw to each element of the situation separately.

Let W be the set of all possible worlds in a semantic model M . Semantically, a closest world function cw is a function $W \times \mathcal{P}(^cW) \rightarrow \mathcal{P}(^cW)$, which assigns a set of possible worlds to each world. In other words, $cw(w, [\phi]_M^v) = [\psi]_M^v$ means that $[\psi]_M^v$ is the set of the closest- ϕ -world to the world w , where $[\phi]_M^v$ as usual is a set of world in which ϕ holds, i.e. $[\phi]_M^v = \{w \in W : M, w, v \models \phi\}$.

The semantic component of the preference in the model is a function $\succ : AGENT \rightarrow \mathcal{P}(\mathcal{P}(^cW) \times \mathcal{P}(^cW))$, which assigns a comparison relation for preferences to each agent.

Moreover, in the models, \succ must satisfy the following conditions

(*NORM*):

$(\emptyset \not\succ_i X), (X \not\succ_i \emptyset)$, where $\succ_i = \succ(i)$ for each agent $i \in AGENT$

(*TRAN*):

$cw(w, X \cap \bar{Y}) \succ_i cw(w, Y \cap \bar{X})$ and $cw(w, Y \cap \bar{Z}) \succ_i cw(w, Z \cap \bar{Y})$
 $\Rightarrow cw(w, X \cap \bar{Z}) \succ_i cw(w, Z \cap \bar{X})$, where $\bar{X} = W - X$

(*NORM*) and (*TRAN*) constrain the semantic preference relation. (*NORM*) is required in support of the logical axiom (*N*) (normality), which, in turn, protects the preference logic against counterintuitive consequences. (*TRAN*) guarantees the soundness of the logical axiom (*TR*) which, in turn, assures transitivity for preferences.

The meaning function of the preference relation is:

$$M, w, v \models \phi \mathbf{P}_i \psi \text{ iff } cw(w, [\phi \wedge \neg \psi]_M^v) \succ_i cw(w, [\psi \wedge \neg \phi]_M^v).$$

The interpretation of $\phi \mathbf{P}_i \psi$ assures the conjunction expansion principle.

2.3. Formal Properties of Preferences

The preference operator has the following axioms and inference rules:

Axioms

- (CEP) $\phi \mathbf{P}_i \psi \leftrightarrow (\phi \wedge \neg \psi) \mathbf{P}_i (\neg \phi \wedge \psi)$
 (N) $\neg(\perp \mathbf{P}_i \phi), \neg(\phi \mathbf{P}_i \perp)$
 (TR) $(\phi \mathbf{P}_i \psi) \wedge (\psi \mathbf{P}_i \chi) \rightarrow (\phi \mathbf{P}_i \chi)$

Inference Rules

- (MP) $\vdash \phi \ \& \ \vdash \phi \rightarrow \psi \Rightarrow \vdash \psi$
 (SUBP) $\vdash(\phi \leftrightarrow \phi') \ \& \ \vdash(\psi \leftrightarrow \psi') \Rightarrow \vdash(\phi \mathbf{P}_i \psi) \leftrightarrow (\phi' \mathbf{P}_i \psi')$

(CEP) states the conjunction expansion principle. (N) establishes “normality” and (TR) transitivity. As noted before, (TR) would go if its semantic equivalent, (TRAN), goes, so we could have non-transitive preferences. (CEP) and (N) together imply the irreflexivity and contraposition (CP) of the \mathbf{P} operator[7]. We have the modus ponens (MP) for obvious reasons. Furthermore, logically equivalent propositions are substitutable in preference formulae (SUBP). Note that we do *not* have monotonicity for preferences. Because of this, we are able to avoid the counterintuitive deductive closure of goals that mars other action logics.

Furthermore, preferences in this semantics have pleasant logical properties. In particular, the preference operator can *avoid* the following counterintuitive properties:

- Necessitation rule for preferences:

$$\models \phi \Rightarrow \models \phi \mathbf{P}_i \psi \text{ and } \models \phi \Rightarrow \models \psi \mathbf{P}_i \phi$$

The first half of this property is exemplified by the statement “if it is necessary that the sun rises in the morning, then the state of the sun’s rising in the morning is always preferred to any other state”. This is definitely counterintuitive. It is also easy to find a counterexample for the second half of this property.

- Closure for preference:

$$\models (\phi \rightarrow \psi) \Rightarrow \models (\phi \mathbf{P}_i \phi' \rightarrow \psi \mathbf{P}_i \phi') \text{ and}$$

$$\models (\phi \rightarrow \psi) \Rightarrow \models (\phi' \mathbf{P}_i \phi \rightarrow \phi' \mathbf{P}_i \psi)$$

This property means that if I prefer tea to coffee, then I prefer tea or one million dollars to coffee, since having tea always implies having tea or one million dollars.

- Conjunction extension:

$$\models \phi \mathbf{P}_i \psi \rightarrow (\phi \wedge \phi') \mathbf{P}_i \psi \text{ and } \models \psi \mathbf{P}_i \phi \rightarrow \psi \mathbf{P}_i (\phi \wedge \phi')$$

This property has a consequence that if I prefer tea to coffee, then I prefer tea and poison to coffee.

Disjunction extension

$$\models \phi \mathbf{P}_i \psi \rightarrow (\phi \vee \phi') \mathbf{P}_i \psi \text{ and } \models \psi \mathbf{P}_i \phi \rightarrow \psi \mathbf{P}_i (\phi \vee \phi')$$

This property is a simplified case of the closure for preferences. We can use the same counter-example for this property.

To repeat: ALX *avoids* these properties.

2.4. Good and Bad States

Following von Wright, we define a “good” state ϕ as a state that agent i prefers to its negation, and conversely for a bad state:

$$Good_i(\phi) \stackrel{\text{def}}{\Leftrightarrow} (\phi \mathbf{P}_i \neg \phi) \quad Bad_i(\phi) \stackrel{\text{def}}{\Leftrightarrow} (\neg \phi \mathbf{P}_i \phi)$$

Proposition 1 (More properties of goodness and badness)

- (a) $\phi \mathbf{P}_i \psi \wedge Good_i \psi \rightarrow Good_i \phi$.
- (b) $\phi \mathbf{P}_i \psi \wedge Bad_i \phi \rightarrow Bad_i \psi$.
- (c) $Good_i \phi \leftrightarrow Bad_i \neg \psi$.
- (d) $Good_i \phi \rightarrow \neg Bad_i \phi$.
- (e) $Bad_i \phi \rightarrow Bad_i \neg \phi$.

Proof:

(a)

$$\begin{aligned}
 & \vdash \phi \mathbf{P}_i \psi \wedge \text{Good}_i \psi \\
 \Rightarrow & \vdash \phi \mathbf{P}_i \psi \wedge \psi \mathbf{P}_i \neg \psi \\
 \Rightarrow & \vdash \phi \mathbf{P}_i \psi \wedge \phi \mathbf{P}_i \neg \psi \quad (TR) \\
 \Rightarrow & \vdash \phi \mathbf{P}_i \psi \wedge \psi \mathbf{P}_i \neg \phi \quad (CP) \\
 \Rightarrow & \vdash \phi \mathbf{P}_i \neg \phi \quad (TR) \\
 \Rightarrow & \vdash \text{Good}_i \phi
 \end{aligned}$$

The proof for (b) is similar to the proof of (a). (c)-(e) are straightforward from the definitions. \square

The notion of goodness and badness are crucial notions used to define the deontic operators in this paper. Since these two operators are defined in terms of the preference operator, they also avoid the counterintuitive properties of preferences.

3. Deontic logic in the framework of ALX.3

3.1. Defining the operations of prohibition, obligation, and permission

ALX.3 is a multi-agent action logic. We can assume that there exists a super-agent, written *sg*, next to other agents. This super-agent need not be a dictator. It could be something like legislature. It lays down the law of land according to its preferences. Therefore, we can define the prohibition **F** operator (“forbidding”) as follows:

$$\mathbf{F}_i \phi \stackrel{\text{def}}{\Leftrightarrow} \text{Bad}_{sg} \text{Good}_i \phi$$

Furthermore, we define the obligation and permission operators in the standard way, namely,

$$\begin{aligned}
 \mathbf{O}_i \phi & \stackrel{\text{def}}{\Leftrightarrow} \mathbf{F}_i (\neg \phi), \\
 \mathbf{P}_i \phi & \stackrel{\text{def}}{\Leftrightarrow} \neg \mathbf{O}_i (\neg \phi).
 \end{aligned}$$

Adding the above two definitions to ALX.3, we obtain deontic ALX, called DALX. Since ALX.3 is sound and complete [7], adding more definitions actually does not change any formal properties of the original logic. As a

consequence, DALX is sound and complete as well. Furthermore, DALX keeps a lot of nice validities of ordinary deontic logics.

The following are theorems of the deontic ALX:

- (a) Consistency of prohibition

$$\mathbf{F}_i\phi \rightarrow \neg\mathbf{F}_i\neg\phi.$$

- (b) Consistency of obligation

$$\mathbf{O}_i\phi \rightarrow \neg\mathbf{O}_i\neg\phi.$$

- (c) Connection between prohibited and permitted states

$$\mathbf{P}_i\phi \rightarrow \neg\mathbf{F}_i\phi.$$

- (d) Obligation implies permission

$$\mathbf{O}_i\phi \rightarrow \mathbf{P}_i\phi.$$

$$\textit{Proof: } \mathbf{O}_i\phi \Rightarrow \mathbf{F}_i\neg\phi \Rightarrow \neg\mathbf{F}_i\phi \Rightarrow \mathbf{P}_i\phi \quad \square$$

- (e) Prohibited states are not permitted

$$\mathbf{F}_i\phi \leftrightarrow \neg\mathbf{P}_i\phi.$$

- (f) No contradictory obligation

$$\neg\mathbf{O}_i(\phi \wedge \neg\phi)$$

- (g) Permission property

$$\neg\mathbf{P}_i\phi \rightarrow \mathbf{P}_i\neg\phi.$$

$$\textit{Proof: } \neg\mathbf{P}_i\phi \Rightarrow \mathbf{F}_i\phi \Rightarrow \neg\mathbf{F}_i\neg\phi \Rightarrow \mathbf{P}_i\neg\phi \quad \square$$

- (h) Substitution rule for prohibited states

$$\models(\phi \leftrightarrow \psi) \Rightarrow \models(\mathbf{F}_i\phi \leftrightarrow \mathbf{F}_i\psi).$$

- (i) Substitution rule for obligation

$$\models(\phi \leftrightarrow \psi) \Rightarrow \models(\mathbf{O}_i\phi \leftrightarrow \mathbf{O}_i\psi).$$

- (j) Substitution rule for permission

$$\models(\phi \leftrightarrow \psi) \Rightarrow \models(\mathbf{P}_i\phi \leftrightarrow \mathbf{P}_i\psi).$$

3.2. Avoiding the Paradoxes

This new deontic logic does not only preserve a lot of nice validities of ordinary deontic logics. More importantly, we can avoid the paradoxes. Just observe the connection between those paradoxes and the counterintuitive properties:

- Ross's paradox is an example of the disjunction extension property.
- Penitent's paradox is an example of the conjunction extension property.
- The Good Samaritan paradox is an example of the closure under logical implication.

Since all three deontic operations in DALX are defined in terms of preferences, those paradoxes are avoided.

Proposition 2 The deontic logic ALX logic DALX can avoid (i) Ross's paradox, (ii) Penitent's paradox, and (iii) Good Samaritan paradox.

Proof: For more technical details of the proof, see the Appendix 2. \square

How does DALX do with respect to the other paradoxes of deontic logic? Actually, one can see that most of them can be reduced the above three typical paradoxes.

- Derived Obligation 1: $O\phi \rightarrow O(\psi \rightarrow \phi)$.

The derived obligation 1 is logically equivalent to the expression $O\phi \rightarrow O(\neg\psi \vee \phi)$. Since the formula ψ in this expression is arbitrary, the derived obligation amounts to Ross' Paradox.

- Derived Obligation 2: $O\neg\phi \rightarrow O(\phi \rightarrow \psi)$

Similarly, this paradox can be reduced to the expression $O\neg\phi \rightarrow O(\neg\phi \vee \psi)$. Again, this is actually Ross' Paradox.

- Chisholm's Paradox: $O\phi \wedge O(\phi \rightarrow \psi) \wedge (\neg\phi \rightarrow O\neg\psi) \wedge \neg\phi \equiv \text{false}$

Here are Chisholm's contrary-to-duty imperatives: (i) it ought to be that someone goes to the assistance of his neighbours, (ii) it ought to be that if he does go he tell them he is coming, (iii) if he does not go then he ought not to tell them he is coming, and (iv) he does not go. However, in ordinary deontic logics, the above four statements together imply a contradiction. The reasoning is as follows: From $(\neg\phi \rightarrow O\neg\psi) \wedge \neg\phi$, we have $O\neg\psi$, which is valid in any logic. Then, by K -axiom $O\phi \wedge O(\phi \rightarrow \psi)$, we have $O\psi$. (the K -axiom is valid in any standard Kripke semantics). Fortunately, there is no K -axiom for the preference in ALX.3. Therefore, DALX can avoid Chisholm's Paradox.

3.3. Comparison

Anderson's reformulation cannot avoid most paradoxes, since this approach uses the logical implication to define the deontic operators. The monotonicity rule is valid in his logic. However, the monotonicity rule is exactly the Good Samaritan Paradox. Ross' Paradox and Penitent Paradox are just special cases of the Good Samaritan Paradox. Furthermore, Derived Obligations are just special cases of Ross Paradox. Therefore, Anderson's approach cannot avoid Derived Obligations. Since the K -axiom is valid in Anderson's deontic logic, this approach cannot avoid Chisholm's Paradox as well.

Meyer's approach uses dynamic actions instead of the logical implication to define the deontic operators. This approach can get rid of most of the nasty paradoxes, including Chisholm's Paradox. Furthermore, some of these paradoxes are not even expressible in his language any more, such as the Derived Obligation $O\neg\phi \rightarrow O(\phi \rightarrow \psi)$. However, Ross' Paradox remains in Meyer's approach, since his approach uses dynamic actions, and in any standard dynamic logic the axiom $[\alpha]\phi \rightarrow [\alpha \cup \beta]\phi$ is always valid. Furthermore, although $F\alpha \rightarrow F(\alpha \wedge \beta)$ is not expressible any more, $F\alpha \rightarrow F(\alpha \& \beta)$ is still one of the theorems of this logic, where "&" means a parallel composition. Therefore, Penitent Paradox remains, too.

4. Concluding Remarks

We have proposed a deontic logic, DALX, in terms of ALX without introducing any new primitive operator. The main idea of DALX is that we assume a super-agent and define the deontic operators for an ordinary agent i in terms of the superagent's and the agent's preferences. Thanks to the properties of the preference operator in ALX, DALX avoids the paradoxes of ordinary deontic logics.

Although the approach depends on the existence of super-agent, we need not make any metaphysical commitments. Also the super-agent need not be a dictator. As a matter of fact, the assumption is purely conventional. The super-agent may have different interpretations, super-agent may be legislature, or anything else may have the right to establish norms for other agents.

5. Appendix 1: Formal Syntax and Semantics of ALX.3

5.1. Formal Syntax

ALX.3 has a sorted first-order description language. There are predicate letters, regular variables, variables reserved for agent and actions respectively, plus the corresponding constant letters:

- (1) For each natural number $n(\geq 1)$, a countable set of n -place predicate letters, PRE_n , written as p_i, p_j, \dots
- (2.1) A countable set of regular variables, $RVAR$, written as x, x_1, y, z, \dots
- (2.2) A countable set of action variables, $AVAR$, written as a, a_1, b, \dots
- (2.3) A countable set of agent variables, $AGVAR$, written as i, i_1, j, \dots
- (3.1) A countable set of regular constants, $RCON$, written as c, c_1, c_2, \dots
- (3.2) A countable set of actions constants, $ACON$, written as ac, ac_1, ac_2, \dots
- (3.3) A countable set of agent constants, $AGCON$, written as ag, ag_1, ag_2, \dots

Furthermore, we have the usual booleans, an existential quantifier, a unary operator for beliefs, binary operators for preferences and causality, a dynamic operator type for actions, and operators that establish a sequence of, or indeterminate choice between, actions. Finally, there are comma and brackets:

- (4) The symbols \neg (negation), \wedge (conjunction), **B**(belief), \exists (existential quantifier), **P**(preference), \leadsto (conditional), $;$ (sequence), \cup (choice), \langle, \rangle , (, and) .

Definition 1 (Variable) Define the set of variables VAR as follows:
 $VAR = RVAR \cup AVAR \cup AGVAR$.

Definition 2 (Constant) Define the set of constants CON as follows:
 $CON = RCON \cup ACON \cup AGCON$.

Definition 3 (Term) Define the set of terms $TERM$ as follows:
 $TERM = VAR \cup CON$.

Definition 4 (Action Term) Define the set of action terms $ATERM$ as follows:
 $ATERM = AVAR \cup ACON$.

Definition 5 (Agent Term) Define the set of agent terms $AGTERM$ as follows:
 $AGTERM = AGVAR \cup AGCON$.

We use t, t_1, \dots , to denote terms, a, a_1, \dots , to denote action terms, i, j, \dots , to denote agent terms, if that does not cause any ambiguity.

An atomic first-order formula is defined as usual:

Definition 6 (ATOM) Define the set of atomic formulas $ATOM$ as follows:
 $ATOM =_{df} \{p(t_1, t_2, \dots, t_n) : p \in PRE_n, t_1, t_2, \dots, t_n \in TERM\}$

Primitive actions carry an agent index. Compound actions need not carry an agent index; this allows for the sequencing of actions carried out by different agents:

Definition 7 (ACTION) Define the set of action expressions $ACTION$ recursively as follows:

- $a \in ATERM, i \in AGTERM \Rightarrow a_i \in ACTION$.
- $a, b \in ACTION \Rightarrow (a; b), (a \cup b) \in ACTION$.

The definition of formulas is standard:

Definition 8 (FORMULA) Define the set of formulae FML recursively as follows:

- $ATOM \subseteq FML$.
- $\phi \in FML \Rightarrow \neg\phi \in FML$.
- $\phi, \psi \in FML \Rightarrow (\phi \wedge \psi) \in FML$.
- $\phi \in FML, x \in VAR \Rightarrow (\exists x\phi) \in FML$.
- $\phi \in FML, a \in ACTION \Rightarrow (\langle a \rangle \phi) \in FML$.

- $\phi, \psi \in FML \Rightarrow (\phi \leadsto \psi) \in FML$.
- $\phi, \psi \in FML, i \in AGTERM \Rightarrow (\phi \mathbf{P}_i \psi) \in FML$.
- $\phi \in FML, i \in AGTERM \Rightarrow \mathbf{B}_i \phi \in FML$.

5.2. Semantics

Definition 9 (ALX.3 Model) Call

$$M = \langle O, PA, AGENT, W, cw, \succ, \mathfrak{R}, \mathfrak{B}, I \rangle$$

an ALX.3 model, if

- O is a set of objects,
- PA is a set of primitive actions,
- $AGENT$ is a set of agents,
- W is a set of possible worlds,
- $cw: W \times \mathcal{P}(W) \rightarrow \mathcal{P}(W)$ is a closest world function,
- $\succ: AGENT \rightarrow \mathcal{P}(\mathcal{P}(W) \times \mathcal{P}(W))$ is a function that assigns a comparison relation for preferences to each agent,
- $\mathfrak{R}: AGENT \times PA \rightarrow \mathcal{P}(W \times W)$ is a function that assigns an accessibility relation to each agent and each primitive action,
- $\mathfrak{B}: AGENT \rightarrow \mathcal{P}(W \times W)$ is a function that assigns an accessibility relation for the belief operation to each agent,
- I is a pair $\langle I_P, I_C \rangle$,
 where I_P is a predicate interpretation function that assigns to each n -place predicate letter $p \in PRE_n$ and each world $w \in W$ a set of n tuples $\langle u_1, \dots, u_n \rangle$, where each of the u_i, \dots, u_n is in $D = O \cup PA \cup AGENT$, called a domain, and I_C is a constant interpretation function that assigns to each regular constants $c \in RCON$ an object $d \in O$, assigns to each action constant $ac \in ACON$ a primitive action $a_p \in PA$, and assigns to each agent constant $g \in AGCON$ an agent $a_g \in AGENT$.

and if cw, \succ, \mathfrak{B} satisfy the following conditions respectively:

- (CS1): $cw(w, X) \subseteq X$.
- (CS2): $w \in X \Rightarrow cw(w, X) = \{w\}$.
- (CS3): $cw(w, X) = \emptyset \Rightarrow cw(w, Y) \cap X = \emptyset$.
- (CS4): $cw(w, X) \subseteq Y$ and $cw(w, Y) \subseteq X \Rightarrow cw(w, X) = cw(w, Y)$.
- (CS5): $cw(w, X) \cap Y \neq \emptyset \Rightarrow cw(w, X \cap Y) \subseteq cw(w, X)$.

For each agent $i \in AGENT$

(*NORM*): $(\emptyset \not\prec_i X), (X \not\prec_i \emptyset)$, where $\succ_i = \succ(i)$

(*TRAN*): $cw(w, X \cap \bar{Y}) \succ_i cw(w, Y \cap \bar{X})$ and $cw(w, Y \cap \bar{Z}) \succ_i cw(w, Z \cap \bar{Y}) \Rightarrow cw(w, X \cap \bar{Z}) \succ_i cw(w, Z \cap \bar{X})$,
where $\bar{X} = W - X$

(*SEB*): $\forall w \exists w' (\langle w, w' \rangle \in \mathfrak{B}_i)$, where $\mathfrak{B}_i = \mathfrak{B}(i)$

(*TRB*): $\langle w, w' \rangle \in \mathfrak{B}_i$ and $\langle w', w'' \rangle \in \mathfrak{B}_i \Rightarrow \langle w, w'' \rangle \in \mathfrak{B}_i$

(*CS#*) constrain the closest world function. Note that we do not require uniqueness for closest world. (*NORM*) and (*TRAN*) constrain the semantic preference relation. (*NORM*) is required to support the logical axiom (*N*) (normality), which, in turn, protects the preference logic against counterintuitive consequences. Its direct effect is to rule out the occurrence of \perp in preference statements. By conjunction expansion principle, (*NORM*) actually implies (*IRE*): $\neg(\phi \mathbf{P}_i \phi)$,¹ that requires irreflexivity, since we are working with a “strong” preference. (*TRAN*) guarantees the soundness of the logical axiom (*TR*) which, in turn, assures transitivity for preferences. (*SEB*) establishes the seriality of the beliefs and prevents agents from believing \perp , while (*TRB*) assures positive introspection. If (*TRB*) would go, the logical axiom (*4B*) would go as well, so we could have a version of *ALX* without positive introspection. (*SEB*) and (*TRB*) make the relation \mathfrak{B} serial and transitive. They are standard requirements for the semantics of beliefs.

Definition 10 (Valuation of Variables) A valuation of variables v in the domain D of an *ALX.3* model M is a mapping that assigns to each variable $x \in VAR$ $v(x) \in OBJECT$, $v(a) \in PA$, and $v(i) \in AGENT$ for any $x \in RVAR$, $a \in AVAR$, and $i \in AGVAR$.

Definition 11 (Valuation of terms) For an *ALX.3* model $M = \langle O, PA, AGENT, W, cw, \succ, \mathfrak{R}, \mathfrak{B}, I \rangle$ and a valuation of variables v , a valuation of terms v_I is a function that assigns to each term $t \in TERM$ an element in the domain D , which is defined as follows:

$$t \in CON \Rightarrow v_I(t) = I_C(t)$$

$$t \in VAR \Rightarrow v_I(t) = v(t)$$

¹ since $\phi \mathbf{P}_i \phi \Rightarrow (\phi \wedge \neg \phi) \mathbf{P}_i (\phi \wedge \neg \phi) \Rightarrow \perp \mathbf{P}_i \perp \Rightarrow \text{False}$.

Definition 12 (Accessibility Relations for Actions) We define an accessibility relation R^a in a model $M = \langle O, PA, AGENT, W, cw, \succ, \mathfrak{R}, \mathfrak{B}, I \rangle$ and a valuation v for each action $a \in ACTION$ as follows.

- $a \in ATERM, i \in AGTERM \Rightarrow R^a = \mathfrak{R}(v_I(a), v_I(i))$,
- $a, b \in ACTION \Rightarrow R^{(a;b)} = R^a \circ R^b = \{ \langle w, w' \rangle \in W \times W : (\exists w_1 \in W) (R^a w w_1 \text{ and } R^b w_1 w') \}$,
- $a, b \in ACTION \Rightarrow R^{(a \cup b)} = R^a \cup R^b$.

Definition 13 (Meaning function) Let FML be as above and let $M = \langle O, PA, AGENT, W, cw, \succ, \mathfrak{R}, \mathfrak{B}, I \rangle$ be an ALX.3 model. Let furthermore v be a valuation of variables in the domain D . Then the meaning function $\llbracket \cdot \rrbracket_M^v : FML \rightarrow \mathcal{P}(W)$ is defined as follows:

$$\begin{aligned}
 \llbracket p(t_1, \dots, t_n) \rrbracket_M^v &= \{ w \in W : \langle v_I(t_1), v_I(t_2), \dots, v_I(t_n) \rangle \in Ip(p, w) \} \text{ where } p \in PRE_n. \\
 \llbracket \neg \phi \rrbracket_M^v &= W \setminus \llbracket \phi \rrbracket_M^v. \\
 \llbracket \phi \wedge \psi \rrbracket_M^v &= \llbracket \phi \rrbracket_M^v \cap \llbracket \psi \rrbracket_M^v. \\
 \llbracket \exists x \phi \rrbracket_M^v &= \{ w \in W : (\exists d \in D)(w \in \llbracket \phi \rrbracket_M^{v(d/x)} \} \}. \\
 \llbracket \langle a \rangle \phi \rrbracket_M^v &= \{ w \in W : (\exists w' \in W)(R^a w w' \text{ and } w' \in \llbracket \phi \rrbracket_M^v) \}. \\
 \llbracket \phi \leadsto \psi \rrbracket_M^v &= \{ w \in W : cw(w, \llbracket \phi \rrbracket_M^v) \subseteq \llbracket \psi \rrbracket_M^v \}. \\
 \llbracket \phi P_i \psi \rrbracket_M^v &= \{ w \in W : cw(w, \llbracket \phi \wedge \neg \psi \rrbracket_M^v) \succ_{v_I(i)} cw(w, \llbracket \psi \wedge \neg \phi \rrbracket_M^v) \}. \\
 \llbracket B_i \phi \rrbracket_M^v &= \{ w \in W : (\forall w') (\langle w, w' \rangle \in \mathfrak{B}_{v_I(i)} \Rightarrow w' \in \llbracket \phi \rrbracket_M^v) \}.
 \end{aligned}$$

Definition 14 (The logic ALX.3) Let FML be as above, let Mod be the class of all ALX.3 models, and let $\llbracket \cdot \rrbracket_M^v$ be as above, defined for every model $M \in Mod$. We call the logic $\langle FML, Mod, \llbracket \cdot \rrbracket_M^v \rangle$ ALX.3 logic.

\models is defined as usual:

Let $M = \langle O, PA, AGENT, W, cw, \succ, \mathfrak{R}, \mathfrak{B}, I \rangle$.

$$M \models \phi \stackrel{def}{\Leftrightarrow} (\forall v \in V_D)(\forall w \in W)(M, w, v \vdash \phi).$$

$$M \models \Gamma \stackrel{def}{\Leftrightarrow} (\forall \gamma \in \Gamma)(M \vdash \gamma).$$

$$Mod(\Gamma) \stackrel{def}{\Leftrightarrow} \{ M \in Mod : M \vdash \Gamma \}.$$

$$\Gamma \models \phi \stackrel{def}{\Leftrightarrow} Mod(\Gamma) \subseteq Mod(\{\phi\}).$$

Definition 15 (ALX.3 inference system) Let ALX3S be the following set of axioms and rules of inference.

(BA):	all tautologies of the first order logic	
(A1):	$\langle a \rangle \perp$	$\leftrightarrow \perp$.
(A2):	$\langle a \rangle (\phi \vee \psi)$	$\leftrightarrow \langle a \rangle \phi \vee \langle a \rangle \psi$.
(A3):	$\langle a; b \rangle \phi$	$\leftrightarrow \langle a \rangle \langle b \rangle \phi$.
(A4):	$\langle a \cup b \rangle \phi$	$\leftrightarrow \langle a \rangle \phi \vee \langle b \rangle \phi$.
(AU):	$[a] \forall x \phi$	$\leftrightarrow \forall x [a] \phi$.
(ID):	$\psi \rightsquigarrow \psi$.	
(MPC):	$(\psi \rightsquigarrow \phi)$	$\rightarrow (\psi \rightarrow \phi)$
(CC):	$(\psi \rightsquigarrow \phi) \wedge (\psi \rightsquigarrow \phi')$	$\rightarrow (\psi \rightsquigarrow \phi \wedge \phi')$.
(MOD):	$(\neg \psi \rightarrow \psi)$	$\rightarrow (\phi \rightsquigarrow \psi)$
(CSO):	$[(\psi \rightsquigarrow \phi) \wedge (\phi \rightsquigarrow \psi)]$	$\rightarrow [(\psi \rightsquigarrow \chi) \leftrightarrow (\phi \rightsquigarrow \chi)]$.
(CV):	$[(\psi \rightsquigarrow \phi) \wedge \neg(\psi \rightsquigarrow \neg \chi)]$	$\rightarrow [(\psi \wedge \chi) \rightsquigarrow \phi]$
(CS):	$(\psi \wedge \phi)$	$\rightarrow (\psi \rightsquigarrow \phi)$.
(CEP):	$\phi \mathbf{P}_i \psi$	$\leftrightarrow (\phi \wedge \neg \psi) \mathbf{P}_i (\neg \phi \wedge \psi)$.
(N):	$\neg(\perp \mathbf{P}_i \phi), \neg(\phi \mathbf{P}_i \perp)$.	
(TR):	$(\phi \mathbf{P}_i \psi) \wedge (\psi \mathbf{P}_i \chi)$	$\rightarrow (\phi \mathbf{P}_i \chi)$.
(PC):	$(\phi \mathbf{P}_i \psi)$	$\rightarrow \neg((\phi \wedge \neg \psi) \rightsquigarrow \neg(\phi \wedge \neg \psi)) \wedge$ $\neg((\psi \wedge \neg \phi) \rightsquigarrow \neg(\psi \wedge \neg \phi))$.
(KB):	$\mathbf{B}_i \phi \wedge \mathbf{B}_i (\phi \rightarrow \psi)$	$\rightarrow \mathbf{B}_i \psi$.
(DB):	$\neg \mathbf{B}_i \perp$.	
(4B):	$\mathbf{B}_i \phi$	$\rightarrow \mathbf{B}_i \mathbf{B}_i \phi$.
(BFB):	$\forall x \mathbf{B}_i \phi$	$\leftrightarrow \mathbf{B}_i \forall x \phi$.
(MP):	$\vdash \phi \ \& \vdash \phi \rightarrow \psi$	$\Rightarrow \vdash \psi$.
(G)	$\vdash \phi$	$\Rightarrow \vdash \forall x \phi$.
(NECA):	$\vdash \phi$	$\Rightarrow \vdash [a] \phi$.
(NECB):	$\vdash \phi$	$\Rightarrow \vdash \mathbf{B}_i \phi$.
(MONA)	$\vdash \langle a \rangle \phi \ \& \vdash \phi \rightarrow \psi$	$\Rightarrow \vdash \langle a \rangle \psi$.
(MONC):	$\vdash \phi \rightsquigarrow \psi \ \& \vdash \psi \rightarrow \psi'$	$\Rightarrow \vdash \phi \rightsquigarrow \psi'$.
(SUBA):	$\vdash (\phi \rightarrow \phi')$	$\Rightarrow \vdash (\langle a \rangle \phi) \leftrightarrow (\langle a \rangle \phi')$.
(SUBC):	$\vdash (\phi \rightarrow \phi') \ \& \vdash (\psi \rightarrow \psi')$	$\Rightarrow \vdash (\phi \rightsquigarrow \psi) \leftrightarrow (\phi' \rightsquigarrow \psi')$.
(SUBP):	$\vdash (\phi \rightarrow \phi') \ \& \vdash (\psi \rightarrow \psi')$	$\Rightarrow \vdash (\phi \mathbf{P}_i \psi) \leftrightarrow (\phi' \mathbf{P}_i \psi')$.

Most axioms are straightforward. As usual, we have the tautologies (BA). Since ALX.3 is a normal modal logic, the absurdum is not true anywhere, so it is not accessible (A1). The action modalities behave as usual, so they distribute over disjunction both ways (A2) (they also distribute over conjunction in one direction, but the corresponding axiom is redundant). (A3) characterizes the sequencing operator ‘;’ and (A4) does the same for the indeterminate choice of actions. (AU) establishes the Barcan formula for universal action modalities. We have the Barcan formula because the underlying domain D is the same in all possible worlds.

The next seven axioms characterize the intensional conditional. Informally speaking, they specify syntactically the meaning of “ceteris paribus” in ALX.3. They are fairly standard, and, with the exception of (CC), they already provide a characterization of Lewis’ system VC, which, in turn, is an adaptation of Stalnaker’s conditional logic for non-unique closest worlds. (ID) establishes the triviality that ψ is true in all closest ψ worlds; (MPC) relates the intensional and the material conditional in the obvious way: so if ϕ would hold given ψ then, if ψ actually does hold, ϕ must also hold. Conjunction distributes over the causality operator in one way (CC). (MOD) rules out the eventuality of closest absurd worlds; (CSO) gives an identity condition for closest worlds, (CV) establishes a cautious monotony for the intensional conditional, and (CS) relates the conjunction to the intensional conditional. Replacing (CS) by

$$(\phi \rightarrow \psi) \vee (\phi \rightarrow \neg \psi)$$

would return Stalnaker’s original systems, as the new axiom would require the uniqueness of the closest possible world.

The next four axioms characterize the preference relation. (CEP) states the conjunction expansion principle. (IRE) confirms the irreflexivity of the P_i operator. (N) establishes “normality” and (TR) transitivity. As noted before, (TR) would go if its semantic equivalent, (TRAN), goes, so we could have non-transitive preferences. The axiom (PC) says that if an agent i prefers ϕ to ψ , then both $\phi \wedge \neg \psi$ and $\psi \wedge \neg \phi$ are possible.

The last four axioms characterize the belief operator. As pointed out above, our belief operator is designed to represent subjective knowledge. (KB) is standard in epistemic logic, but it is often criticised, since it requires logical omniscience with respect to the material conditional. On the other hand, one would expect rational agents to draw correct logical inferences, when necessary, so not having (KB) might be worse. (DB) rules out the belief in absurdities, (4B) establishes positive self-introspection for beliefs, and (BFB) is the Barcan formula for beliefs. These four axioms give a standard characterization of subjective knowledge. Together with the infer-

ence rules (MP) and (NECB), they turn the belief operation into a weak S4 system.

The remaining expressions characterize ALX.3's inference rules.

We have modus ponens and generalization for obvious reasons. By the same token, we have the necessitation rule for the universal action modality: if indeed, ϕ is true in all worlds, then all activities will lead to ϕ -worlds; by the same token, we have the necessitation rule for beliefs. (MONA) connects the meaning of the action modality with the meaning of the material conditional. We have right monotonicity for the intensional conditional but *not* left monotonicity. Furthermore, logically equivalent propositions are substitutable in action- conditional- and preference formulae (SUBA), (SUBC), (SUBP). Note that we are *not* having monotonicity for preferences. Because of this, we are able to avoid the counterintuitive deductive closure of goals that mars other action logics.

6 Appendix 2: Some Proofs

Proposition 3 The deontic ALX.3 logic can avoid (i) Ross's paradox, (ii) Penitent's paradox, and (iii) Good Samaritan paradox.

Proof: (i) In order to prove that the deontic logic in ALX.3 can avoid Ross's paradoxes, we have to show that $O_i \wedge \neg O_i(\phi \vee \psi)$ is satisfiable.

$$\begin{aligned}
 & O_i \phi \wedge \neg O_i(\phi \vee \psi) \text{ is satisfiable} \\
 \Leftrightarrow & F_i \neg \phi \wedge \neg F_i \neg(\phi \vee \psi) \text{ is satisfiable} && (\text{Definition of } O_i) \\
 \Leftrightarrow & F_i \neg \phi \wedge \neg F_i(\neg \phi \wedge \neg \psi) \text{ is satisfiable} && (\text{Meta reasoning}) \\
 \Leftrightarrow & \mathbf{Bad}_{sg} \mathbf{Good}_i \neg \phi \wedge \mathbf{Bad}_{sg} \mathbf{Good}_i(\neg \phi \wedge \neg \psi) && (\text{Definition of } F_i) \\
 & \text{is satisfiable} \\
 \Leftrightarrow & \neg(\mathbf{Good}_i \neg \phi) P_{sg}(\mathbf{Good}_i \neg \phi) \wedge && \\
 & \neg((\neg \mathbf{Good}_i(\neg \phi \wedge \neg \psi)) P_{sg} \mathbf{Good}_i(\neg \phi \wedge \neg \psi)) && \\
 & \text{is satisfiable} && (\text{Definition of } \mathbf{BAD})
 \end{aligned}$$

In the following, we will construct a model M such that there exists a world which

$$\begin{aligned}
 M, w, v \models & \neg(\mathbf{Good}_i \neg p) P_{sg}(\mathbf{Good}_i \neg p) \wedge \\
 & \neg((\mathbf{Good}_i(\neg p \wedge \neg q)) P_{sg} \mathbf{Good}_i(\neg p \wedge \neg q)).
 \end{aligned}$$

The model is constructed as follows:

$$M = \langle O, PA \text{ AGENT}, W, cw, \succ, \mathfrak{R}, \mathfrak{B}, I \rangle,$$

where

$$O = \emptyset$$

$$PA = \emptyset$$

$$AGENT = \{sg, ag\},$$

$$W = \{w_0, w_1, w_2, w_3\},$$

cw is defined by the minimal difference between two worlds,

$$\succ(sg) = \{\langle\{w_1, w_2\}, \{w_0\}\rangle\};$$

$$\succ(ag) = \{\langle\{w_2\}, \{w_0\}\rangle\},$$

$$\mathfrak{R} = \emptyset,$$

$$\mathfrak{B} = \emptyset,$$

$$I_p(p, w_0) = \{\langle ag \rangle\},$$

$$I_p(p, w_1) = \{\langle ag \rangle\},$$

$$I_p(q, w_0) = \{\langle ag \rangle\},$$

$$I_p(q, w_2) = \{\langle ag \rangle\}.$$

Let v be a variable evaluation $v(i) = ag$.

So, we have $\llbracket \neg p(i) \mathbf{P}_i p(i) \rrbracket_M^v = \llbracket \neg(\neg p(i) \mathbf{P}_i p(i)) \rrbracket_M^v = \{w_1, w_2, w_3\}$. Furthermore, $\llbracket (\neg p(i) \wedge \neg q(i)) \mathbf{P}_i \neg(\neg p(i) \wedge \neg q(i)) \rrbracket_M^v = \emptyset = \llbracket \perp \rrbracket_M^v$.

Therefore, $cw(w_0, \llbracket \neg(\neg p(i) \mathbf{P}_i p(i)) \rrbracket_M^v) = \{w_1, w_2\}$. By the truth condition, we have

$$M, w_0, v \models \neg(\mathbf{Good}_i \neg p(i)) \mathbf{P}_{sg} (\mathbf{Good}_i \neg p(i)).$$

Furthermore, by the normality, we have

$$M, w_0, v \models \neg((\neg \mathbf{Good}_i (\neg p(i) \wedge \neg q(i))) \mathbf{P}_{sg} \mathbf{Good}_i (\neg p(i) \wedge \neg q(i))).$$

So,

$$\neg(\mathbf{Good}_i \neg p(i)) \mathbf{P}_{sg} (\mathbf{Good}_i \neg p(i)) \wedge \neg((\neg \mathbf{Good}_i (\neg p(i) \wedge \neg q(i))) \mathbf{P}_{sg} \mathbf{Good}_i (\neg p(i) \wedge \neg q(i)))$$

is satisfiable. That completes the proof for (i).

For (ii) and (iii), we follow the same type of the proof. We do not go into the details here. \square

REFERENCES

- [1] Anderson, A. R., A reduction of deontic logic to alethic modal logic, *Mind* 67, 1958, 100-103.
- [2] Chisholm, R., and Sosa, E., On the logic of "intrinsically better", *American Philosophical Quarterly* 3 (1966) 244-249.
- [3] Danielsson, S., *Preference and obligation*, (Filosofiska föreningen, Uppsala, 1968).
- [4] Harel, D. Dynamic logic, in: D. Gabbay and F. Guentner, eds., *Handbook of Philosophical Logic*, Vol.II, (D. Reidel Publishing Company, 1984) 497-604.
- [5] Huang, Z., Masuch, M., and Pólos, L., ALX, an action logic for agents with bounded rationality, *Artificial Intelligence* 82 (1996) 100-153.
- [6] Huang, Z., Masuch, M., and Pólos, L., ALX2, a quantifier ALX logic, CCSOM Research Report 93-99, (1993).
- [7] Huang, Z., *Logics for Agents with Bounded Rationality*, ILLC Dissertation Series 1994-10.
- [8] Huang, Z., and Masuch, M., ALX.3: a multi-agent action logic, *the Proceedings of AAAI95 Spring Symposium*, (1995).
- [9] Meyer, J-J., A different approach to deontic logic: deontic logic viewed as a variant of dynamic logic, *Notre Dame J. of Formal Logic* 29 (1988) 109-136.
- [10] Meyer, J-J., and Wieringa, R. J., Deontic logic: a concise overview, in Meyer, J.-J., and Wieringa, R. J., (eds.) *Deontic Logic in Computer Science*, (John Wiley & Sons, 1993).
- [11] Stalnaker, R., A theory of conditionals, in: *Studies in Logical Theory*, *American Philosophical Quarterly* 2 (1968) 98-122.
- [12] von Wright, G. H., *The Logic of Preference*, (Edinburgh, 1963).
- [13] von Wright, G. H., The logic of preference reconsidered, *Theory and Decision* 3 (1972) 140-169.