

A Brief Outline of ALX3, a Multi-Agent Action Logic

Zhisheng Huang and Michael Masuch

Center for Computer Science in Organization and Management (CCSOM)

University of Amsterdam

Oude Turfmarkt 151, 1012 GC Amsterdam, The Netherlands

email: {huang,michael}@ccsom.uva.nl

Abstract

ALX3 is a multi-agent version of ALX with a first-order description language. ALX3 is sound and complete, and is already proving its practical use in the formal representation of modern organization theory.

1 Introduction

Action logics are usually developed for the (hypothetical) use by intelligent robots [4; 7; 17] or as a description language of program behavior [8]. Our effort is motivated by a different concern. We want to develop a formal language for social science theories, especially for theories of organizations. The difference in motivation leads to a new approach to action logic. It combines ideas from various strands of thought, notably H.A. Simon's notion of *bounded rationality*, Kripke's *possible world semantics*, V.R. Pratt's *dynamic logic*, Stalnaker's notion of *minimal change*, G. H. von Wright's approach to *preferences*, and J. Hintikka's approach to *knowledge and belief*. ALX3, the action logic presented in this paper, has a first order description language with multiple agents, and four modal operator types.¹ We outline the language, and discuss some of its important properties. A fuller presentation is given in [11] and [12]. In a companion paper, ALX3's potential for knowledge representation is extensively demonstrated in the formalization of an important organization theory, J.D. Thompson's *Organizations in Action* [14].

2 ALX's Background

Most social science theories are expressed in natural language, but natural language does not provide a formal scaffold for checking a theory's logical properties. As a

consequence, the social sciences have acquired a reputation for "softness" — a soft way of saying that the logical properties of their theories are often dubious. Reformulating a social theory in a formal language with known logical properties would facilitate the tasks of consistency checking, disambiguation, or the examination of other important logical properties, such as contingency (whether or not the theory is falsifiable).

We focus on action logic as a formal language, because actions are key to the understanding of social phenomena. In fact, most social scientist agree that action theory provides the underlying framework for the social sciences in general [2; 6; 8; 13; 16; 18]. Yet actions involve attitudes and engender change, and both phenomena are notoriously hard to grasp in the extensional context of first order languages [5]. This explains our attempt to develop a new logic, rather than taking First Order Logic off the shelf.

Herbert A. Simon's conceptualization of *bounded rationality* [20] serves as a point of departure. His approach is intuitively appealing, and had great impact on the postwar social sciences. Simon wanted to overcome the omniscience claims of the traditional conceptualizations of rational action. He assumed (1) an agent with (2) a set of behavior alternatives, (3) a set of future states of affairs (each such state being the outcome of a choice among the behavior alternatives), and (4) a preference order over future states of affairs. The omniscient agent, endowed with "perfect rationality", would know all behavior alternatives and the exact outcome of each alternative; the agent would also have a complete preference ordering for those outcomes. An agent with bounded rationality, in contrast, might not know all alternatives, nor need it know the exact outcome of each; also, the agent might lack a complete preference ordering for those outcomes.

Kripke's *possible world semantics* provides a natural setting for Simon's conceptualization. We assume a set of possible worlds with various relations defined over this set (we may also call those possible worlds *states*). One can see a behavior alternative as a mapping from states to states, so each behavior alternative constitutes an ac-

¹ALX stands for the *x*'s Action Logic. ALX1, the first version, had a propositional description language and a (backward-looking) update operator instead of the conditional; it was a single-agent language [9]. ALX2, the intermediate version, is not multi-agent.[10]

cessibility relation. An accessibility relation, in turn, can be interpreted as an opportunity for action, that is, as an opportunity for changing the world by moving from a given state to another state. Accessibility relations are expressed by indexed one-place modal operators, as in dynamic logic [8]. For example, the formula $\langle a_i \rangle$ expresses the fact that the agent has an action a at its disposal such that effecting a in the present situation would result in the situation denoted by ϕ .

Preferences – not goals – provide the basic rationale for rational action in ALX3. Following von Wright [21], 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 says that she prefers oranges to apples should prefer a situation where she has an orange but *no* apple to a situation where she has an apple but *no* orange. We call this principle *conjunction expansion principle* and restrict attention to preference statements that obey it. Preferences are expressed via two-place modal operators; if the agent prefers the proposition ϕ to the proposition ψ , we write $\phi P_i \psi$.

Normally, the meaning of a preference statement is context dependent, even if this is not made explicit. An agent may say to prefer an apple to an orange – and actually mean it – but she may prefer an orange to an apple later – perhaps because then she already had an apple. To capture this context dependency, we borrow the notion of minimal change from Stalnaker's approach to conditionals [19]. 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 the set of "closest" states relative to a given state, such that the new states fulfill some specified conditions, but resembles the old state as much as possible in all other respects.

The syntactic equivalent of the closest world function is the wiggled "causal arrow". It appears in expressions such as $\phi \rightsquigarrow \psi$ where it denotes: in all closest worlds where ϕ holds, ψ also holds. The causal arrow expresses the conditional notion of a causal relation between ϕ and ψ : if ϕ were the case, ψ would also be the case.

The last primitive operator of ALX3 is the indexed belief operator. In a world of bounded rationality, an agent's beliefs do not necessarily coincide with reality, and in order to make this distinction, we must be able to distinguish between belief and reality; $B_i(\phi)$ will denote the fact that agent i believes ϕ . As the logical axioms characterizing the belief operator show, B represents a sense of "subjective knowledge", not metaphysical attachment, or epistemic uncertainty.

3 Syntax and Semantics

3.1 Formal Syntax

ALX3 has the following primitive symbols:

- (1) For each natural number $n(\geq 1)$, a countable set of n -place predicate letters, p_i, p_j, \dots
- (2.1) A countable set of regular variables, x, x_1, y, z, \dots
- (2.2) A countable set of action variables, a, a_1, b, \dots
- (2.3) A countable set of agent variables, i, i_1, j, \dots
- (3.1) A countable set of regular constants, c, c_1, c_2, \dots
- (3.2) A countable set of actions constants, ac, ac_1, ac_2, \dots
- (3.3) A countable set of agent constants, ag, ag_1, ag_2, \dots
- (4) The symbols \neg (negation), \wedge (conjunction), B (belief), \exists (existential quantifier), P (preference), \rightsquigarrow (conditional), $\langle \rangle$ (sequence), \cup (choice), $\{, \}$, $(,)$, and

Furthermore, ALX3 has the following syntax rules:

(Variable)	::= (Regular variable) (Action variable) (Agent variable)
(Constant)	::= (Regular constant) (Action constant) (Agent constant)
(Term)	::= (Variable) (Constant)
(Action term)	::= (Action variable) (Action constant)
(Agent term)	::= (Agent variable) (Agent constant)
(Atom)	::= (Predicate)((Term), ..., (Term))
(Action)	::= (Action term)(Agent term) (Action); (Action) (Action) \cup (Action)
(Formula)	::= (Atom) \neg (Formula) (Formula) \wedge (Formula) \exists (Variable)(Formula) (Action)(Formula) (Formula) \rightsquigarrow (Formula) (Formula) $P_{(Agent\ term)}$ (Formula) $B_{(Agent\ term)}$ (Formula)

3.2 Semantics

Definition 1 (ALX3 Model)

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

an ALX3 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,

- $\mathcal{R} : AGENT \times PA \rightarrow \mathcal{P}(W \times W)$ is a function that assigns an accessibility relation to each agent and each primitive action,
- $B : AGENT \rightarrow \mathcal{P}(W \times W)$ 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_1, \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 $a_g \in AGENT$.

Definition 2 (Meaning function) Let FML be as above and let

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

be an ALX3 model. Let furthermore v be a valuation of variables in the domain D . Then the meaning function $[\]_M^v$ is defined as follows:

$$\begin{aligned} [p(t_1, \dots, t_n)]_M^v &= \{w : \langle v_I(t_1), \dots, v_I(t_n) \rangle \in I_P(p, w)\} \\ [\neg\phi]_M^v &= W \setminus [\phi]_M^v \\ [\phi \wedge \psi]_M^v &= [\phi]_M^v \cap [\psi]_M^v \\ [\exists x\phi]_M^v &= \{w : (\exists d \in D)(w \in [\phi]_M^{v(d/x)})\} \\ [\langle a \rangle \phi]_M^v &= \{w : (\exists w')(R^a w w' \text{ and } w' \in [\phi]_M^v)\} \\ [\phi \leadsto \psi]_M^v &= \{w : cw(w, [\phi]_M^v) \subseteq [\psi]_M^v\} \\ [\phi P_i \psi]_M^v &= \{w : cw(w, [\phi \wedge \neg\psi]_M^v) \succ_{v_I(i)} cw(w, [\psi \wedge \neg\phi]_M^v)\} \\ [B_i \phi]_M^v &= \{w : (\forall w')(\langle w, w' \rangle \in B_{v_I(i)} \Rightarrow w' \in [\phi]_M^v)\} \end{aligned}$$

The interpretation of the atomic formulas, the boolean connectives and the existential quantifier is straightforward. The interpretation of $\langle a \rangle \phi$ yields the set of worlds from where the agent can access at least one ϕ -world via action a . The interpretation of $\phi \leadsto \psi$ yields the set of worlds in which the closest ϕ -worlds are also ψ -worlds. So, $\phi \leadsto \psi$ is true at a world if ϕ would *ceteris paribus* entail ψ . This is the standard counter-factual conditional used to express a causal relation between ϕ and ψ . Note that our wiggled arrow does not require actual counter-factuality, so ϕ may be true in the actual world. The interpretation of $\phi P_i \psi$ gives a set of worlds such that the agent will prefer at each of those worlds the closest ϕ -and-not- ψ -worlds to the closest ψ -and-not- ϕ worlds. This interpretation assures the “conjunction expansion” principle established by von Wright.

Definition 3 (ALX3 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) : $[a] \forall x \phi \rightarrow \forall x [a] \phi$
- (ID) : $\psi \leadsto \psi$
(MPC) : $(\psi \leadsto \phi) \rightarrow (\psi \rightarrow \phi)$
(CC) : $(\psi \leadsto \phi) \wedge (\psi \leadsto \phi') \rightarrow (\psi \leadsto \phi \wedge \phi')$
(MOD) : $(\neg \psi \leadsto \psi) \rightarrow (\phi \leadsto \psi)$
(CSO) : $[(\psi \leadsto \phi) \wedge (\phi \leadsto \psi)] \rightarrow [(\psi \leadsto \chi) \leftrightarrow (\phi \leadsto \chi)]$
(CV) : $[(\psi \leadsto \phi) \wedge \neg(\psi \leadsto \neg \chi)] \rightarrow [(\psi \wedge \chi) \leadsto \phi]$
(CS) : $(\psi \wedge \phi) \rightarrow (\psi \leadsto \phi)$
- (CEP) : $\phi P_i \psi \leftrightarrow (\phi \wedge \neg \psi) P_i (\neg \phi \wedge \psi)$
(N) : $\neg(\perp P_i \phi), \neg(\phi P_i \perp)$
(TR) : $(\phi P_i \psi) \wedge (\psi P_i \chi) \rightarrow (\phi P_i \chi)$
- (PC) : $(\phi P_i \psi) \rightarrow \neg((\phi \wedge \neg \psi) \leadsto \neg(\phi \wedge \neg \psi)) \wedge \neg((\psi \wedge \neg \phi) \leadsto \neg(\psi \wedge \neg \phi))$
- (KB) : $B_i \phi \wedge B_i (\phi \rightarrow \psi) \rightarrow B_i \psi$
(DB) : $\neg B_i \perp$
(4B) : $B_i \phi \rightarrow B_i B_i \phi$
- (BFB) : $\forall x B_i \phi \rightarrow B_i \forall x \phi$
- (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) : $\vdash \phi \Rightarrow \vdash B_i \phi$
(MONA) : $\vdash \langle a \rangle \phi \ \& \ \vdash \phi \rightarrow \psi \Rightarrow \vdash \langle a \rangle \psi$
(MONC) : $\vdash \phi \leadsto \psi \ \& \ \vdash \psi \rightarrow \psi' \Rightarrow \vdash \phi \leadsto \psi'$
(SUBA) : $\vdash (\phi \leftrightarrow \phi') \Rightarrow \vdash ((\langle a \rangle \phi) \leftrightarrow (\langle a \rangle \phi'))$
(SUBC) : $\vdash (\phi \leftrightarrow \phi') \ \& \ \vdash (\psi \leftrightarrow \psi') \Rightarrow \vdash (\phi \leadsto \psi) \leftrightarrow (\phi' \leadsto \psi')$
(SUBP) : $\vdash (\phi \leftrightarrow \phi') \ \& \ \vdash (\psi \leftrightarrow \psi') \Rightarrow \vdash (\phi P_i \psi) \leftrightarrow (\phi' P_i \psi')$

Most axioms are straightforward. As usual, we have the tautologies (BA). Since ALX3 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 syntactically specify the meaning of “ceteris paribus” in ALX3. 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 to a system 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 “wiggled arrow” 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

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

yields Stalnaker’s original system, 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 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 give a characterization of 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 to draw correct logical inferences when necessary, so not having (KB) may 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 inference rules (MP) and (NECB), they turn the belief operation into a weak S4 system. As shown in [11], we could weaken the belief operator considerably, but these weaker alternatives have their own problems that would overload ALX. (One radical alternative would be an empty belief operator.)

The remaining expressions characterize ALX3’s inference rules. We have the 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 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.

4 Formal Properties of ALX.3

Proposition 1 *ALX3S is sound and complete, i.e., for an arbitrary set of formulas Σ and an arbitrary formula ϕ ,*

$$\Sigma \vdash \phi \Leftrightarrow \Sigma \models \phi$$

PROOF: See [11]. □

5 Conclusions and Summary

ALX is the first action logic modeled on the decision cycle of potentially rational agents. Perhaps its most important feature is its preference operator. The preference operator has a closest-world semantics in combination with the conjunction expansion principle, so agents prefer ϕ to ψ if they prefer the closest ϕ -and-not- ψ worlds to the closest ψ - and-not- ϕ worlds. This facilitates the representation of situation-dependent preferences, but does not exclude the representation of stable preferences. The preference operator is “normal” in the sense that agents cannot have preferences with respect to absurd worlds. This normality protects the conjunction expansion principle against counterintuitive utilizations [3]. On the other hand, the preference operator is not “tarskian”, i.e., it does not distribute over disjunction [15; 1] and this protects all intensional operators built on the preference operator against the necessitation rule and against undesired closure properties. For example, goal operators can be defined as preferred states subject to additional qualifications (e.g., the best accessible state, the best state not believed to be inaccessible). The preference operator is transitive, but a non-transitive version is easily generated by removing the constraint (TRAN) on the semantic preference relation. The conditional operator is adapted from Stalnaker’s system but allows for non-unique closest worlds. It allows for an easy representation of the notion of “ceteris paribus”, and hence for the standard notion of causality. This, in turn, greatly facilitates the representation of causal effects, side effects, and similar non-monotonic relations that would have otherwise to be represented by the (monotonic) material conditional.

An important property of ALX is the virtual absence of interaction between modal operators. This property may raise eyebrows in philosophical circles, but we designed ALX as a flexible knowledge-representation tool, and such a tool should, in our view, not preempt the structure of domains to be represented.

As a flexible tool, ALX3 allows for the definition of various additional intensional operators, such as the alethic modalities, goal operators, intention-operators,

and for the characterization of other action-related notions, such as ability, effect, side effect, intended effect, etc [12].

As demonstrated in a companion paper [14], ALX3 serves already as a versatile tool of knowledge representation. However, there are some desiderata left to be satisfied. ALX3 has no explicit notion of time in its semantics, and this complicates a direct representation of events. For example, we cannot define a “do”-operator, and hence no actions that are not deliberate. Second, we think that the notion of closest worlds needs closer inspection. The constraints on the closest world function are relatively weak; stronger constraints may be required, in particular if one wants to combine the notion of action with the notion of the closest world. Third, and related, one may want to strengthen ALX so that it allows for a *calculation* of causal outcomes. Such a calculation would be an answer to the frame problem, but it requires a more specific notion of possible worlds.

References

- [1] van Benthem, J., *Essays in Logical Semantics*, (D. Reidel Publishing Company, 1986).
- [2] Blumer, H., *Symbolic Interactionism: Perspective and Methods*, (Englewood Cliffs, NJ, Prentice-Hall, 1969).
- [3] Chisholm, R., and Sosa, E., Intrinsic preferability and the problem of supererogation, *Synthese* 16 (1966), 321-331.
- [4] Cohen, P. R. and Levesque, H. J., Intention is choice with commitment. *Artificial Intelligence* 42 (1990) 213-261.
- [5] Gamut, L.T.F., *Logic, Language, and Meaning*, (The University of Chicago Press, 1991).
- [6] Giddens, A., *Central Problems in Social Theory: Action, Structures, and Contradiction in Social Analysis*, (Berkeley, CA, University of California Press, 1979).
- [7] Ginsberg, M. L., and Smith, D. E., Reasoning about action I: a possible worlds approach, in: M. Ginsberg, ed., *Readings in Non-monotonic Reasoning*, (Morgan Kaufman, Los Altos, 1987).
- [8] Harel, D., Dynamic logic, in: D. Gabbay and F. Guenther, eds., *Handbook of Philosophical Logic*, Vol.II, (D. Reidel Publishing Company, 1984) 497-604.
- [9] Huang, Z., Masuch, M., and Pólos, L., ALX, an action logic for agents with bounded rationality, *Artificial Intelligence* (forthcoming).
- [10] Huang, Z., Masuch, M., and Pólos, L., ALX2, a quantifier ALX logic, CCSOM Working Paper 93-99.
- [11] Huang, Z., *Logics for Agents with Bounded Rationality*, ILLC Dissertation series 1994-10, University of Amsterdam, (1994).
- [12] Huang, Z., and Masuch, M., ALX3, a Multi-agent Action Logic, CCSOM Technical Report 94-102.
- [13] Luhmann, N., *The Differentiation of Society*, (New York, Columbia University Press, 1982).
- [14] Masuch, M., and Huang, Z., A Logical Deconstruction: Formalizing J.D. Thompson's *Organizations in Action* in a Multi-agent Action Logic, CCSOM Working Paper 94-120.
- [15] Marx, M., *Algebraic Relativization and Arrow Logic*, ILLC Dissertation series 1995-3, University of Amsterdam, 1995.
- [16] Parsons, T., *The Structure of Social Action*, (Glencoe, IL, Free Press, 1937).
- [17] Rao, A. S. and Georgeff, M. P., Modeling rational agents within a BDI- architecture, in: J. Allen, R. Fikes, and E. Sandewall, eds., *Proceedings of the Second International Conference on Principles of Knowledge Representation and Reasoning*, Morgan Kaufmann Publishers, San Mateo, CA, (1991) 473-484.
- [18] Schutz, A., *The Phenomenology of the Social World*, (Evanston, IL, Northwestern University Press, 1967).
- [19] Stalnaker, R., A theory of conditionals, in: *Studies in Logical Theory*, *American Philosophical Quarterly* 2 (1968) 98-122
- [20] Simon, H. A., A behavioral model of rational choice, *Quarterly Journal of Economics* 69 (1955) 99-118.
- [21] von Wright, G. H., *The Logic of Preference*, (Edinburgh, 1963).