

All They Know: A Study in Multi-Agent Autoepistemic Reasoning

—PRELIMINARY REPORT—

Gerhard Lakemeyer
Institute of Computer Science III
University of Bonn
Romerstr. 164
5300 Bonn 1, Germany
gerhard@cs.uni-bonn.de

Abstract

With few exceptions the study of nonmonotonic reasoning has been confined to the single-agent case. However, it has been recognized that intelligent agents often need to reason about other agents and their ability to reason non-monotonically. In this paper we present a formalization of multi-agent autoepistemic reasoning, which naturally extends earlier work by Levesque. In particular, we propose an n -agent modal belief logic, which allows us to express that a formula (or finite set of them) is *all* an agent knows, which may include beliefs about what other agents believe. The paper presents a formal semantics of the logic in the possible-world framework. We provide an axiomatization, which is complete for a large fragment of the logic and sufficient to characterize interesting forms of multi-agent autoepistemic reasoning. We also extend the stable set and stable expansion ideas of single-agent autoepistemic logic to the multi-agent case.

1 Introduction

While the study of nonmonotonic reasoning formalisms has been at the forefront of foundational research in knowledge representation for quite some time, work in this area has concentrated on the single-agent case with only few exceptions.

This focus on single agents is somewhat surprising, since there is little doubt that agents, who have been invested with nonmonotonic reasoning mechanisms, should be able to reason about other agents and their ability to reason nonmonotonically as well. For example, if we assume the common default that birds normally fly and if Jill tells Jack that she has just bought a bird, then Jill should be able to infer that Jack thinks that her bird flies. Multi-agent nonmonotonic reasoning is also crucial when agents need to coordinate their activity. For example, assume I promised a friend (who is always on time) to meet him at a restaurant at 7PM. If I leave my house knowing that I will not make it there by 7PM, I will probably not change my plans and still go to the restaurant. After all, I know that my friend has

no reason to believe that I am not on my way to meet him and that he will therefore wait for me. Note that I reason about my friends default assumption to wait in this case. Other examples from areas like planning and temporal projection can be found in [Mor90].

One of the main formalisms of nonmonotonic reasoning is *autoepistemic logic* (e.g. [Moo85]). The basic idea is that the beliefs of agents are closed under *perfect introspection*, that is, they know¹ what they know and do not know. Nonmonotonic reasoning comes about in this framework in that agents can draw inferences on the basis of their own ignorance. The following example by Moore illustrates this feature: I can reasonably conclude that I have no older brother simply because I do not *know* of any older brother of mine. A particular formalization of autoepistemic reasoning is due to Levesque [Lev90], who proposes a logic of *only-knowing* (*OL*), which is a classical modal logic extended by a new modality to express that a formula is *all* an agent believes. An advantage of this approach is that, rather than having to appeal to non-standard inference rules as in Moore's original formulation, *OL* captures autoepistemic reasoning using only the classical notions of logical consequence and theoremhood.

In this paper, we propose a propositional multi-agent extension of *OL*. We provide a formal semantics within the possible-world framework and a proof theory, which is complete for a large fragment of the logic. The new logic also leads us to natural extensions of notions like *stable sets* and *stable expansions*, which were originally developed for the single-agent case.

General multi-agent nonmonotonic reasoning formalisms have received very little attention until recently.² A notable exception is work by Morgenstern and Guerreiro [Mor90, MG92], who consider both multi-agent autoepistemic reasoning and multi-agent circumscription theories. On the autoepistemic side, they propose multi-agent versions of stable sets, which al-

¹While we are concerned with belief and, in particular, allow agents to have false beliefs, we nevertheless use the terms knowledge and belief interchangeably.

²There has also been work applying nonmonotonic theories to special multi-agent settings such as speech acts (e.g. [Per87, AK88]). Unlike our work, these approaches are not concerned with general purpose multi-agent nonmonotonic reasoning.

low agents to reason about other agents' nonmonotonic inferences. In contrast to our work, however, these stable sets are not justified by an independent semantic account. Recently, and independently of our work, Halpern [Hal93] also extended Levesque's logic OL to the multi-agent case. While Halpern's logic and ours share many (but not all) properties, the respective model theories are quite different. In particular, while our approach remains within classical possible-world semantics, Halpern uses concepts very much related to the so-called *knowledge structures* of [FHV91]. Using the same technique, Halpern also extends the notion of only-knowing proposed in [HM84] to the multi-agent case. There, however, agents are not capable of reasoning about other agent's nonmonotonic inferences because only-knowing is used only as a meta-logical concept.

The rest of the paper is organized as follows. Section 2 defines the logic OL_n , extending Levesque's logic of only-knowing to many agents. Besides a formal semantics we also provide a proof theory, which is complete for a large fragment of OL_n . Furthermore, we look at the properties of formulas in OL_n which uniquely determine the beliefs of an agent and are thus of particular interest to knowledge representation. Section 3 considers examples of multi-agent autoepistemic reasoning as modeled by OL_n . Section 4 shows how OL_n yields natural multi-agent versions of *stable sets* and *stable expansions*. Finally, we summarize the results of the paper and point to some future work in Section 5.

2 The Logic OL_n

After introducing the syntax of the logic, we define the semantics in two stages. First we describe that part of the semantics that does not deal with only-knowing. In fact, this is just an ordinary possible-world semantics for n agents with perfect introspection. Then we introduce the necessary extensions that give us the semantics of only-knowing. Finally, we present a proof theoretic account, which is complete for a large fragment of the logic, and discuss properties of the logic which are important in the context of knowledge representation.

2.1 Syntax

Definition 1 The Language OL_n

The primitives of OL_n consist of a countably infinite set of atomic propositions (or atoms), the connectives \vee , \neg , and the modal operators L_i and O_i for $1 \leq i \leq n$. (Agents are referred to as $1, 2, \dots, n$.) Formulas are formed in the usual way from these primitives.³ $L_i\alpha$ should be read as "the agent i believes α " and $O_i\alpha$ as " α is all agent i believes." A formula α is called *basic* iff there are no occurrences of O_i ($1 \leq i \leq n$) in α .

Definition 2 A modal operator occurs at depth n of a formula α iff it occurs within the scope of exactly n modal operators.

³We will freely use other connectives like \wedge , \supset and \equiv , which should be understood as syntactic abbreviations of the usual kind.

For example, given $\alpha = p \wedge L_1 L_2 (L_3 q \vee \neg O_2 r)$, L_1 occurs at depth 0, L_2 at depth 1, and L_3 and O_2 occur both at depth 2.

Definition 3 A formula α is called *i-objective* (for $i = 1, \dots, n$) iff every modal operator at depth 0 is of the form O_j or L_j with $i \neq j$.

In other words, *i-objective* formulas talk about the external world from agent i 's point of view, which includes beliefs of other agents but not his own. For example, $(p \vee L_2 q) \wedge \neg O_3 L_1 p$ is *1-objective*, but $(p \vee L_2 q) \wedge \neg L_1 O_3 p$ is not.

Definition 4 A formula α is called *i-subjective* iff every modality at depth 0 is of the form L_i or O_i and every atom occurs within the scope of a modal operator.

That is, *i-subjective* formulas talk only about what agent i believes. For example, $L_1 p \wedge \neg O_1 p$ is *1-subjective*, $L_1 p \wedge q$ and $L_1 p \wedge L_2 q$ are not.

2.2 The Semantics of Basic Formulas

Basic formulas are given a standard possible-world semantics [Kri63, Hin62, Hin71], which the reader is assumed to be familiar with.⁴ Roughly, a *possible-world model* consists of *worlds*, which determine the truth of atomic propositions, and binary *accessibility* relations between worlds. An agent's beliefs at a given world w are determined by what is true in all those worlds that are accessible to the agent from w . Since we are concerned with agents who possess *perfect introspection*, we restrict the accessibility relations in the usual way, that is, they are both transitive and Euclidean.⁵ The resulting logic is called $K45_n$ or *weak S5_n*.⁶

Definition 5 A $K45_n$ -Model

$M = \langle W, \pi, R_1, \dots, R_n \rangle$ is called a $K45_n$ -model (or simply model) iff

1. W is a set (of worlds).
2. π is a mapping from the set of atoms into 2^W .
3. $R_i \subseteq W \times W$ for $1 \leq i \leq n$.
4. R_i is transitive and Euclidean for $1 \leq i \leq n$.

Given a model $M = \langle W, \pi, R_1, \dots, R_n \rangle$ and a world $w \in W$, the possible-world semantics of basic formulas is defined as follows: Let p be an atom and α and β arbitrary basic formulas.

$$\begin{aligned} w \models p & \iff w \in \pi(p) \\ w \models \neg \alpha & \iff w \not\models \alpha \\ w \models \alpha \vee \beta & \iff w \models \alpha \text{ or } w \models \beta \\ w \models L_i \alpha & \iff \text{for all } w', \text{ if } w R_i w' \text{ then } w' \models \alpha \end{aligned}$$

A set of basic formulas Γ is *satisfiable* iff there is a model $M = \langle W, \pi, R_1, \dots, R_n \rangle$ and $w \in W$ such that $w \models \gamma$ for all $\gamma \in \Gamma$.

In the single-agent case, Levesque made use of the fact that the semantics of $K45_1$ has a much simpler formulation. There we can assume a fixed set of worlds, which is the set of all truth assignments. A model then consists

⁴See [HC84, HM92] for an introduction.

⁵ R_i is Euclidean iff $\forall w, w', w'', \text{ if } w R_i w' \text{ and } w R_i w'', \text{ then } w' R_i w''$.

⁶The subscript n indicates that we are concerned with the n -agent case.

of a world w , which corresponds intuitively to the real world, and a set of worlds W' , which determine the beliefs of the agent. There is no need for an explicit accessibility relation, since the worlds in M are globally accessible from every world and a sentence is believed just in case it is true at all worlds in W . Unfortunately, such a simple model does not extend to the multi-agent case and we are forced to a more complicated semantics with explicit accessibility relations as defined above.⁷ For this reason, the extension of Levesque's logic OL to many agents turns out to be a non-trivial exercise.

2.3 The Canonical Model

It is well known that, as far as basic formulas are concerned, it suffices to look at just one, the so-called canonical model [HC84, HM92]. This canonical model will be used later on to define the semantics of only-knowing.

The central idea behind canonical models are maximally consistent sets.

Definition 6 Maximally consistent sets

Given any proof theory of $K45_n$ and the usual notion of theoremhood and consistency, a set of basic formulas T is called maximally consistent iff T is consistent and for every basic α , either α or $\neg\alpha$ is contained in T .

The canonical $K45_n$ -model M_c has as worlds precisely all the maximally consistent sets and a world w' is inaccessible from w just in case all of i 's beliefs at w are included in w' .

Definition 7 The Canonical $K45_n$ -Model M_c

The canonical model $M_c = (W_c, \pi, R_1, \dots, R_n)$ is a Kripke structure such that

1. $W_c = \{w \mid w \text{ is a maximally consistent set}\}$.
2. For all atoms p and $w \in W_c$, $w \in \pi(p)$ iff $p \in w$.
3. wR_iw' iff for all formulas $L_i\alpha$, if $L_i\alpha \in w$ then $\alpha \in w'$.

The following (well known) theorem tells us that nothing is lost from a logical point of view if we confine our attention to the canonical model.

Theorem 1 M_c is a $K45_n$ -model and for every set of basic formulas T , T is satisfiable iff it is satisfiable in M_c .

2.4 The Semantics of All They Know

Given this classical possible-world framework, what does it mean for an agent i to only-know, say, an atom p at some world w in a model M ? Certainly, i should believe p , that is, all worlds that are i -accessible from w should make p true. Furthermore, i should believe as little else as possible apart from p . For example, i should neither believe q nor believe that j believes p etc. Minimizing knowledge using possible worlds simply means maximizing the number of accessible worlds. Thus, in our example, there should be an accessible world where q is false and another one where j does not believe p and so on. It should be clear that in order for w to satisfy

⁷In essence, if we have more than 1 agent and a global set of worlds for each agent, the agents would be mutually introspective, which is not what we want.

only-knowing α this way, the model M must have a huge supply of worlds that are accessible from w . While not essential for the definition of only-knowing, it turns out to be very convenient to simply restrict our attention to models that are guaranteed to contain a sufficient supply of worlds.⁸ In fact, we will consider just one, namely the canonical model of $K45_n$. Let us call the set of all formulas that are true at some world w in some model of $K45_n$ a world state. The canonical model has the nice property that it contains precisely one world for every possible world state, since world states are just maximally consistent sets.

With that agent i is said to only-know a formula α at some world w (in the canonical model) just in case α is believed and any world w' which satisfies α and from which the same worlds are i -accessible as from w is itself i -accessible from w . We now turn to the formal definitions.

Definition 8 Given a model $M = (W, \pi, R_1, \dots, R_n)$ and worlds w and w' in W , we say that w and w' are i -equivalent ($w \approx_i w'$) iff for all worlds $w'' \in W$, wR_iw'' iff $w'R_iw''$.

Given an arbitrary formula α of OL_n , a world w in a model M , let

$$w \models O_i\alpha \iff \text{for all } w' \text{ s.t. } w \approx_i w', wR_iw' \text{ iff } w' \models \alpha.$$

A formula α of OL_n is a logical consequence of a set of formulas Γ iff for all worlds w in the canonical model M_c , if $w \models \gamma$ for all $\gamma \in \Gamma$, then $w \models \alpha$. As usual, we say that α is valid ($\models \alpha$) iff $\{\} \models \alpha$. A formula α is satisfiable iff $\neg\alpha$ is not valid.

Note that Theorem 1 guarantees that OL_n restricted to basic formulas is still simply $K45_n$.

A superficial comparison of the rules for O_i and L_i suggests that they differ at two places. For one the O_i -rule quantifies over i -equivalent worlds only. For another, the "if ... then" of the L_i -rule is replaced by an "iff." While the latter is significant, it can be easily shown that we can restrict the L_i -rule to quantify over i -equivalent worlds as well.

Lemma 2.1

For any $K45_n$ -model M and world w in M ,

$$w \models L_i\alpha \iff \text{for all } w' \text{ s.t. } w \approx_i w', \text{ if } wR_iw' \text{ then } w' \models \alpha$$

2.5 A Proof Theory

In order to obtain a proof theory for a large fragment of OL_n , we apply the same idea Levesque used to axiomatize his logic OL . The idea is to add a dual operator N_i for every L_i to the language OL_n , which we refer to as ONL_n . While $L_i\alpha$ can be read as "agent i knows at least α ," $N_i\alpha$ should be read as "agent i knows at most that α is false." $O_i\alpha$ now becomes simply an abbreviation for $L_i\alpha \wedge N_i\neg\alpha$, that is, agent i believes only α iff he believes at least and at most α . We extend the notion of an i -subjective or i -objective formula to cover the N_i -operators as well. For example, $\neg L_1p \wedge N_1\neg L_2p$ is both 1-subjective and 2-objective.

⁸In Levesque's OL , this is automatically given since worlds are drawn from the fixed set of all truth assignments.

The semantics of N_i is as follows. $N_i\alpha$ is true at a world w iff α is true at all i -equivalent worlds of w which are *not* accessible from w . In the formalization it is convenient to add new accessibility relations \bar{R}_i for just these non-accessible worlds to the canonical $K45_n$ -model. This way, the N_i behave just like ordinary $K45$ belief operators.

Definition 9 *The Extended Canonical Model*
Given the canonical model $M_c = (W_c, \pi, R_1, \dots, R_n)$, let the extended canonical model be

$$M_{ec} = (W_c, \pi, R_1, \dots, R_n, \bar{R}_1, \dots, \bar{R}_n),$$

where for all $w, w' \in W_c$, $w\bar{R}_i w'$ iff $w \approx_i w'$ and $wR_i w'$.

Lemma 2.2 *The \bar{R}_i are transitive and Euclidean.*

Given the extended canonical model M_{ec} and a world $w \in W_c$, the semantic rule for N_i is simply:

$$w \models N_i\alpha \iff \text{for all } w', \text{ if } w\bar{R}_i w' \text{ then } w' \models \alpha.$$

Notions like logical consequence and validity for this extended logic are defined as for OL_n with M_c replaced by M_{ec} .

Given Lemma 2.2, it is obvious that the N_i have all the properties of a $K45_n$ -operator. Moreover, if we view L_i and N_i as two different agents, then the L_i -agent knows exactly which worlds the N_i -agent can see and vice versa. In other words, the L_i -agent and N_i -agent are mutually introspective. For example $\neg L_i p \supset N_i \neg L_i p$ is valid. Finally, note that for every equivalence class of worlds W_i under \approx_i , R_i and \bar{R}_i cover all worlds in W_i exhaustively, that is, taken together, the two relations form a complete subgraph over W_i . This feature is reflected in valid formulas of the form $N_i\alpha \supset \neg L_i\alpha$ for falsifiable α .

Axioms:

- A1 Axioms of propositional logic
- A2 $L_i(\alpha \supset \beta) \supset (L_i\alpha \supset L_i\beta)$
- A3 $N_i(\alpha \supset \beta) \supset (N_i\alpha \supset N_i\beta)$
- A4 $\sigma \supset L_i\sigma \wedge N_i\sigma$ for all i -subjunctive σ
- A5 $N_i\alpha \supset \neg L_i\alpha$ for all basic i -obj. α falsif. in $K45_n$
- A6 $O_i\alpha \equiv (L_i\alpha \wedge N_i\neg\alpha)$ for all α

Inference Rules:

- MP From α and $\alpha \supset \beta$ infer β .
- Nec From α infer $L_i\alpha$ and $N_i\alpha$.

Note that A2–A4 imply that the L_i and N_i have all the properties of regular $K45_n$ -operators. Also, A4 not only gives us the regular introspection axioms of $K45_n$ but also cross axioms (mutual introspection) such as $\neg L_i\alpha \supset N_i\neg L_i\alpha$. A5 provides the crucial link between L_i and N_i capturing the relationship between “knowing at least” and “knowing at most.”⁹

Theorem 2 Soundness

Given the usual definition of theoremhood (\vdash) with respect to the above proof theory, for all α in ONL_n , if $\vdash\alpha$ then $\models\alpha$.

⁹Levesque, in his axiomatization of OL , uses a special case of this axiom. In his case, α ranges merely over the falsifiable formulas of classical propositional logic (no modalities).

Notice that axiom A5 assumes that α ranges only over basic i -objective formulas. We need this restriction in order to appeal to falsifiability in the existing logic $K45_n$.¹⁰ For the axiomatization to be complete, such a restriction essentially requires that arbitrary formulas $N_i\beta$ or $L_i\beta$ are reducible to equivalent forms $N_i\beta^*$ and $L_i\beta^*$, where β^* is a basic formula, that is a formula without any N_j 's. That, however, is not true in general. For example, $L_iN_j p$ and $N_iL_jN_i p$ are not reducible for distinct i and j . However, if we rule out such cases, the proof theory is indeed complete for the restricted language, which we call ONL_n^- .

Definition 10 *The Language ONL_n^-*
 α is a formula of ONL_n^- iff α is a formula of ONL_n and, after replacing every occurrence of $O_i\beta$ within α by its definition $L_i\beta \wedge N_i\neg\beta$, no N_j may occur within the scope of an N_i or L_i for $i \neq j$.

For example, while $N_iL_i\neg N_i p$ and $N_i(L_j p \vee N_i\neg p)$ are in ONL_n^- , $N_iN_j p$ and $N_iL_jN_i p$ are not for distinct i and j .

Theorem 3 Completeness For ONL_n^-

For all $\alpha \in ONL_n^-$, if $\models\alpha$ then $\vdash\alpha$.

Proof: The proof is an adaptation of Levesque's completeness proof for OL [Lev90]. ■

(Halpern recently proved that the above proof theory is indeed *not* complete for all of OL_n .)

Given Theorem 2 and 3, it is not hard to prove that OL_1 is equivalent to the propositional version of Levesque's logic OL because Levesque's proof theory is a special case of ours.

Theorem 4 OL_1 and OL are equivalent

A formula α in OL_1 is a theorem of OL_1 iff α is a theorem of OL .

We introduced the operator N_i mainly for the purpose of obtaining a proof theory for OL_n . Except for a formal derivation using the axioms in Section 3, we will no longer be concerned with the N_i -operators in the rest of the paper and go back to the original OL_n and its language OL_n .

2.6 i -Determinate Sentences

We now turn to a class of formulas which are of particular interest to knowledge representation, since they determine, when only-known, precisely what an agent believes and does not believe.

Definition 11 i -Determinate Formulas¹¹

A formula $\alpha \in OL_n$ is i -determinate iff

1. there is a world w in M_c such that $w \models O_i\alpha$ and
2. for all worlds w' , if $w' \models O_i\alpha$ then $w \approx_i w'$.

The following theorem demonstrates that i -determinate formulas indeed deserve their name.

Theorem 5 For all $\alpha \in OL_n$, α is i -determinate iff for all β , exactly one of $O_i\alpha \supset L_i\beta$ and $O_i\alpha \supset \neg L_i\beta$ is valid.

¹⁰Note that this peculiar axiom schema is recursive since falsifiability in propositional $K45_n$ is decidable.

¹¹In the single-agent case, the analogous concept of *determinate* formulas was defined in [Lev90].

What are examples of determinate formulas? In the single-agent case, it has been shown that all objective formulas (no modalities at all) are determinate. Not surprisingly, objective formulas are also i -determinate. In fact, this result can be generalized to include all basic i -objective formulas.

Theorem 6

All basic i -objective formulas are i -determinate.

Other examples of i -determinate formulas, which are not i -objective, include $\neg L_i \neg p \supset p$ and $\neg L_i L_j p \supset \neg L_j p$, which allow agents to reason nonmonotonically and are discussed in more detail in Section 3.

So far we have only considered basic i -determinate sentences. Does the above result extend to non-basic i -objective formulas as well? The answer, surprisingly, is: sometimes but not always! For example, it is not hard to show that the formula $O_j p$ (where p is an atom) is also i -determinate, that is, the beliefs of agent i are uniquely determined if all i knows is that all j knows is p . However, the formula $\neg O_j p$ is not i -determinate because $\models \neg O_i \neg O_j p$.¹² In other words, it is impossible for i to only-know that j does not only-know p . Intuitively, for i to know that j does not only-know p , i needs to have some evidence in terms of a basic belief or non-belief of j . It is because of properties like $\models \neg O_i \neg O_j p$ that our axiomatization is not complete for all of ONL_n .

3 Multi-agent Nonmonotonic Reasoning

While our axiomatization is complete only for a subset of OL_n , it is nevertheless strong enough to model interesting cases of multi-agent nonmonotonic reasoning. Here are two examples:

1. Let p be agent i 's secret and suppose i makes the following assumption: unless i know that j knows my secret assume that j does not know it. We can prove in OL_n that if this assumption is all i believes then he indeed believes that j does not know his secret. Formally $\vdash O_i(\neg L_i L_j p \supset \neg L_j p) \supset L_i \neg L_j p$.¹³

A formal derivation of this theorem of OL_n can be obtained as follows. Let $\alpha = \neg L_i L_j p \supset \neg L_j p$. The justifications in the following derivation indicate which axioms or previous derivations have been used to derive the current line. PL or $K45_n$ indicate that reasoning in either standard propositional logic or $K45_n$ is used is without further analysis.

1	$O_i \alpha \supset L_i \alpha$	A6; PL
2	$O_i \alpha \supset N_i \neg \alpha$	A6; PL
3	$(L_i \alpha \wedge \neg L_i L_j p) \supset L_i \neg L_j p$	$K45_n$
4	$N_i \neg \alpha \supset (N_i \neg L_i L_j p \wedge N_i L_j p)$	$K45_n$
5	$N_i L_j p \supset \neg L_i L_j p$	A5
6	$O_i \alpha \supset \neg L_i L_j p$	2;4;5; PL
7	$O_i \alpha \supset L_i \neg L_j p$	1;3;6; PL

To see that i 's beliefs may evolve nonmonotonically given

¹²In contrast, $O_i \neg O_j p$ is satisfiable in Halpern's logic.

¹³Note that if we replace $L_j p$ by p we obtain regular single-agent autoepistemic reasoning.

that i only-knows α , assume that i finds out that j has found out about the secret. Then i 's belief that j does not believe the secret will be retracted. Formally $\vdash O_i(L_j p \wedge (\neg L_i L_j p \supset \neg L_j p)) \supset (\neg L_i \neg L_j p \wedge L_i L_j p)$. Notice that, while OL_n is itself a regular monotonic logic, the nonmonotonicity of agent i 's beliefs is hidden within the O_i -operator.

Finally, $\neg L_i L_j p \supset \neg L_j p$ is also i -determinate.

In particular, $O_i(\neg L_i L_j p \supset \neg L_j p) \equiv O_i \neg L_j p$ is valid.

2. Now let p stand for the old "Tweety flies." As expected, we obtain $\vdash O_j(\neg L_j \neg p \supset p) \supset L_j p$. (Not surprisingly, $\neg L_j \neg p \supset p$ is j -determinate.) By the rule of necessitation and the distribution axiom for belief we immediately get $\vdash L_i O_j(\neg L_j \neg p \supset p) \supset L_i L_j p$. In other words, i is able to reason about j 's ability to reason nonmonotonically, essentially by simulating j 's reasoning.

It should be pointed out that an assumption like i knows that all j knows is α is certainly unrealistic in general and is related to what Morgenstern calls *arrogance* [Mor90]. It would be much more reasonable if we could say that i knows that α is all j knows about *some relevant subject*, say Tweety. In fact, we have proposed such a notion of *only-knowing-about* as an extension of OL_n in [Lak93]. There we also show that only-knowing is often a good approximation of *only-knowing-about* in the sense that, if we restrict ourselves to beliefs that are only about the subject matter in question,¹⁴ then the beliefs that follow from only-knowing-about are the same as those that follow from only-knowing.

4 i -Stable Sets and i -Stable Expansions

Single-agent autoepistemic logic was originally developed using the concepts of *stable sets* [Sta80] and *stable expansions* [Moo85]. Here we define natural n -agent extensions of these notions and show how they relate to OL_n . For the purposes of this paper, we confine ourselves to *basic* formulas only.

Definition 12 i -Epistemic State

A set of basic formulas Γ is called an i -epistemic state iff there is a world w in M_c such that for all basic γ , $w \models L_i \gamma$ iff $\gamma \in \Gamma$.

Definition 13 i -Stable Sets

Let Γ be a set of basic formulas. Γ is called i -stable iff

1. Γ contains all valid formulas of $K45_n$.
2. If $\alpha \in \Gamma$ and $\alpha \supset \beta \in \Gamma$ then $\beta \in \Gamma$.
3. If $\alpha \in \Gamma$ then $L_i \alpha \in \Gamma$.
4. If $\alpha \notin \Gamma$ then $\neg L_i \alpha \in \Gamma$.

Note that the only difference between i -stable sets and the original definition of *stable sets* lies in condition 1. While *stable sets* are only required to contain the tautologies of regular propositional logic, i -stable sets must contain all $K45_n$ -valid formulas.

The next theorem proves that i -stable sets and i -epistemic states are equivalent notions. In other words, the definition of i -stability falls out naturally given the semantics of OL_n .

¹⁴The subject matter is defined as a set of atomic propositions π . The restriction says that the atoms occurring within a belief must be contained in π .

Theorem 7

Let Γ be a set of basic formulas. Γ is *i-stable* iff Γ is an *i-epistemic state*.

Halpern arrived at the same notion of *i-stability* independently [Hal93]. He also considered definitions based on logics other than $K45_n$.

Given a set of basic formulas Γ , let $\bar{\Gamma} = \{\gamma \mid \gamma \text{ is basic and } \gamma \notin \Gamma\}$, $L_i\Gamma = \{L_i\gamma \mid \text{for all } \gamma \in \Gamma\}$, and $\neg L_i\bar{\Gamma} = \{\neg L_i\gamma \mid \text{for all } \gamma \in \bar{\Gamma}\}$.

Definition 14 *i-Stable Expansion*

Let A be a set of basic formulas and let \models_{K45} denote logical consequence in $K45_n$. Γ is called an *i-stable expansion* of A iff $\Gamma = \{\text{basic } \gamma \mid A \cup L_i\Gamma \cup \neg L_i\bar{\Gamma} \models_{K45} \gamma\}$.

The definition of *i-stable expansions* looks exactly like Moore's definition of *stable expansions* except that we use logical consequence in $K45_n$ instead of logical consequence in propositional logic. Using the stronger $K45_n$ is necessary in the multi-agent case, since an agent knows that other agents are also perfectly introspective. For example, if agent i believes $\neg L_j p$ for a different agent j then he also believes $L_j \neg L_j p$.

Finally, the following theorem demonstrates that the *i-stable expansions* of a formula α correspond precisely to the different *i-epistemic states* of agent i who only-knows α .

Theorem 8 *Only-Knowing and i-Stable Expansions*

Let α be basic and $w \in W_c$ with $\Gamma = \{\text{basic } \alpha \mid w \models L_i \alpha\}$. Then $w \models O_i \alpha$ iff Γ is an *i-stable expansion* of $\{\alpha\}$.

An analogous result was obtained by Halpern [Hal93].

Corollary 4.1 *Basic i-objective formulas have a unique i-stable expansion.*

5 Conclusion

We proposed a multi-agent logic of only-knowing that extends earlier work by Levesque regarding the single-agent case. Our logic gives a semantic and proof theoretic account of autoepistemic reasoning for many knowers. Notions like *stable set* and *stable expansion* fall out as natural extensions of single-agent autoepistemic logic.

As for future work, it would be interesting to obtain a complete axiomatization for all of OL. Also, as noted earlier, there are subtle differences between OL_n and Halpern's logic. Halpern showed that every valid sentence in his logic is also valid in OL_n and that the valid sentences of both logics coincide when restricted to ONL_n^- . However, there are sentences such as $\neg O_i \neg O_j p$ that are valid in OL_n but not in Halpern's case. For a better comparison of the two approaches it would be interesting to see which modifications are necessary to obtain identical logics. Finally, a more expressive first-order language should be considered to make OL_n more applicable in real world domains. We conjecture that an approach as in [Lev90] could be adapted for this purpose without great difficulty.

Acknowledgements

I would like to thank Joe Halpern and Hector Levesque for fruitful discussions on this subject.

References

- [AK88] Appelt, D. and Konolige, K., A Practical Non-monotonic Theory of Reasoning about Speech Acts, in Proc. of the 26th Conf of the ACL, 1988.
- [FHV91] A Model-Theoretic Analysis of Knowledge, Journal of the ACM 91(2), 1991, pp. 382-428.
- [Hal93] Halpern, J. Y., Reasoning about only knowing with many agents, in Proc. of the 11th National Conference on Artificial Intelligence (AAAI-93).
- [HM84] Halpern, J. Y. and Moses, Y. O., Towards a Theory of Knowledge and Ignorance: Preliminary Report, in Proceedings of The Non-Monotonic Workshop, New Paltz, NY, 1984, pp.125-143.
- [HM92] Halpern, J. Y. and Moses, Y. O., A Guide to Completeness and Complexity for Modal Logics of Knowledge and Belief, Artificial Intelligence 54, 1992, pp. 319-379.
- [HC84] Hughes, G. E. and Cresswell, M. J., A Companion to Modal Logic, Methuen & Co., London, 1984.
- [Hin62] Hintikka, J., Knowledge and Belief: An Introduction to the Logic of the Two Notions, Cornell University Press, 1962.
- [Hin71] Hintikka, J., Semantics for Propositional Attitudes, in L. Linsky (ed.), Reference and Modality, Oxford University Press, Oxford, 1971.
- [Kri63] Kripke, S. A., Semantical Considerations on Modal Logic, Acta Philosophica Fennica 16, 1963, pp. 83-94.
- [Lak93] All They Know About, to appear in: Proc. of the 11th National Conference on Artificial Intelligence (AAAI-93), Washington DC, 1993.
- [Lev90] Levesque, H. J., All I Know: A Study in Autoepistemic Logic, Artificial Intelligence, North Holland, 42, 1990, pp. 263-309.
- [Moo85] Moore, R., Semantical Considerations on Non-monotonic Logic, Artificial Intelligence 25, 1985, pp. 75-94.
- [Mor90] Morgenstern, L., A Theory of Multiple Agent Nonmonotonic Reasoning, in Proc. of AAAI-90, 1990, pp. 538-544.
- [MG92] Morgenstern, L. and Guerreiro, R., Epistemic Logics for Multiple Agent Nonmonotonic Reasoning!, Symposium on Formal Reasoning about Beliefs, Intentions, and Actions, Austin, TX, 1992.
- [Per87] Perrault, R., An Application of Default Logic to Speech Act Theory, in Proc. of the Symposium on Intentions and Plans in Communication and Discourse, Monterey, 1987.
- [Sta80] Stalnaker, R. C., A Note on Nonmonotonic Modal Logic, Department of Philosophy, Cornell University, 1980.