

Reasoning about other agents' beliefs under bounded resources

Natasha Alechina, Brian Logan, Nguyen Hoang Nga, and Abdur Rakib*

School of Computer Science
University of Nottingham
Nottingham NG8 1BB, UK
{nza,bsl,hnn,rza}@cs.nott.ac.uk

Abstract. There exists a considerable body of work on epistemic logics for bounded reasoners where the bound can be time, memory, or the amount of information the reasoners can exchange. In much of this work the epistemic logic is used as a meta-logic to reason about beliefs of the bounded reasoners from an external perspective. In this paper, we present a formal model of a system of bounded reasoners which reason about each other's beliefs, and propose a sound and complete logic in which such reasoning can be expressed. Our formalisation highlights a problem of incorrect belief ascription in resource-bounded reasoning about beliefs, and we propose a possible solution to this problem, namely adding reasoning strategies to the logic.

1 Introduction

The purpose of this paper is to investigate a multi-agent epistemic logic which results from taking seriously the idea that agents have bounded time, memory and communication resources, and are reasoning about each other's beliefs. The main contribution of the paper is to generalise several existing epistemic logics for resource-bounded reasoners by adding an ability for reasoners to reason about each other's beliefs. We show that a problem of incorrect belief ascription arises as a result, and propose a possible solution to this problem.

To give the reader an idea where the current proposal fits into the existing body of research on epistemic logics for bounded reasoners, we include a brief survey of existing approaches, concentrating mostly on the approaches which have influenced the work presented here.

In standard epistemic logic (see e.g. [1, 2] for a survey) an agent's (implicit) knowledge is modelled as closed under logical consequence. This can clearly pose a problem when using an epistemic logic to model resource-bounded reasoners, whose set of beliefs is not generally closed with respect to their reasoning rules. Various proposals to modify possible worlds semantics in order to solve this problem of logical omniscience (e.g., introducing impossible worlds as in [3, 4], or non-classical assignment as in [5]) result in agent's beliefs still being logically closed, but with respect to a weaker logic.

* This work was supported by the UK Engineering and Physical Sciences Research Council [grant number EP/E031226].

Our work builds on another approach to solving this problem, namely treating beliefs as syntactic objects rather than propositions (sets of possible worlds). In [6], Fagin and Halpern proposed a model of limited reasoning using the notion of awareness: an agent explicitly believes only the formulas which are in a syntactically defined awareness set (as well as in the set of its implicit beliefs). Implicit beliefs are still closed under consequence, but explicit beliefs are not, since a consequence of explicit beliefs is not guaranteed to belong to the awareness set. However, the awareness model does not give any insight into the connection between the agent's awareness set and the agent's resource limitations, which is what we try to do in this paper.¹ Konolige [7] proposed a different model of non-omniscient reasoners, the deduction model of belief. Reasoners were parameterised with sets of rules which could, for example, be incomplete. However, the deduction model of belief still models beliefs of a reasoner as closed with respect to reasoner's deduction rules; it does not take into account the time it takes to produce this closure, or any limitations on the agent's memory. Step logic, introduced in [8], gives a syntactic account of beliefs as theories indexed by time points; each application of inference rules takes a unit of time. No fixed bound on memory was considered, but the issue of bounded memory was taken into account. An account of epistemic logic called algorithmic knowledge, which treats explicit knowledge as something which has to be computed by an agent, was introduced in [9], and further developed in e.g. [1, 10]. In the algorithmic knowledge approach, agents are assumed to possess a procedure which they use to produce knowledge. In later work [10] this procedure is assumed to be given as a set of rewrite rules which are applied to the agent's knowledge to produce a closed set, so, like Konolige's approach, algorithmic knowledge is concerned with the result rather than the process of producing knowledge. In [11, 12] Duc proposed logics for non-omniscient epistemic reasoners which will believe all consequences of their beliefs *eventually*, after some interval of time. It was shown in [13] that Duc's system is complete with respect to semantics in which the set of agent's beliefs is always finite. Duc's system did not model the agents' reasoning about each others' beliefs. Other relevant approaches where epistemic logics were given a temporal dimension and each reasoning step took a unit of time are, for example, [14], where each inference step is modelled as an action in the style of dynamic logic, and [15] which proposes a logic for verification of response-time properties of a system of communicating rule-based agents (each rule firing or communication takes a unit of time). In a somewhat different direction, [16] proposed a logic where agents reason about each others beliefs, but have no explicit time or memory limit; however there is a restriction on the depth of belief nestings (context switching by the agents). Epistemic logics for bounded-memory agents were investigated in, for example, [17–20], and the interplay between bounded recall and bounded memory (ability to store strategies of only bounded size) was studied in [21].

An epistemic logic BMCL for communicating agents with communication limits on the number of exchanged messages (and connections to space complexity of proofs and communication complexity) was investigated in [20]. In this paper we expand BMCL by adding rules for reasoning about other agents' beliefs, demonstrate that epistemic

¹ We also completely dispense with the notion of implicit beliefs.

reasoning done in resource-bounded fashion has an inherent problem of incorrect belief ascription, and propose the use of reasoning strategies as a solution to this problem.

2 Model of reasoning agents

The logic BMCL presented in [20] formalises reasoning about the beliefs of a system of reasoners who reason using propositional resolution and can exchange information to solve a problem together. The set up is similar to, for example, [22]. BMCL models each inference rule application as taking a single time step, introduces an explicit bound on the set of beliefs of each reasoner, and a bound on the number of messages the reasoners can exchange. In this paper, we generalise this approach by assuming that agents can also reason about each other's beliefs. Namely, they assume that other agents use a certain set of inference rules, and they reason about what another agent may believe at the next step. For example, if agent A believes that agent B believes two clauses c_1 and c_2 and these two clauses are resolvable to a clause c , and agent A assumes that agent B reasons using resolution, then it is reasonable for agent A to believe that agent B may believe c at the next step.

We assume a set of n agents. Each agent i has a set of inference rules, a set of premises KB_i , and a *working memory*. To infer from the premises in KB_i , the relevant formulas must first be read into working memory. We assume that each agent's working memory is bounded by n_M , which is the maximal number of formulas an agent can believe at the same time. We also set a limit on the possible size of a formula, or rather on the depth of nesting of belief operators, n_B , and a limit, n_C , on the maximal number of communications an agent can make. For simplicity, we assume that these bounds are the same for all agents, but this can be easily relaxed by introducing functions $n_M(i)$, $n_B(i)$ and $n_C(i)$ which assign a different limit to each agent i .

The set of reasoning actions is as follows:

- Read KB:** an agent can retrieve information from its KB and put it into its working memory using the *Read* action. Since an agent has a fixed size memory, adding a formula to its memory may require erasing some belief already in memory (if the limit n_M would otherwise be exceeded). The same applies to other reasoning actions which add a new formula, in that adding a new formula may involve overwriting a formula currently in working memory.
- Resolution:** an agent can derive a new clause if it has two resolvable clauses in its memory.
- Copy:** an agent can communicate with another agent to request a clause from the memory of the other agent. We assume that communication is always successful if the other agent has the requested clause. If agent A has clause c in memory, then a copy by B will result in agent B believing that A believes c . *Copy* is only enabled if the agent has performed fewer than n_C copy actions in the past and the prefix of the resulting belief has nesting of at most n_B .
- Idle:** an agent may idle (do nothing) at any time step. This means that at the next time point of the system, the agent does not change its state of memory.
- Erase:** an agent may remove a formula from its working memory. This action is introduced for technical reasons to simplify the proofs.

In addition to the actions listed above, we introduce actions that enable agents to reason about other agents' beliefs, essentially epistemic axioms K (ascribing propositional reasoning to the other agent) and 4 (positive introspection about the agent's own beliefs, and ascribing positive introspection to other agents). The reasons we do not adopt for example KD45 are as follows. If the agent's knowledge base is inconsistent, we want it to be able to derive $B\perp$ (or $B\Box$ where \Box is the empty clause). Negative introspection is also problematic in a resource-bounded setting, in that the agent may derive $\neg B\alpha$ if α is not in its current set of beliefs, and then derive α from its other beliefs, ending up with an inconsistent set of beliefs ($\neg B\alpha$ and $B\alpha$ by positive introspection from α), even if its knowledge base is consistent. We could have adopted a restricted version of negative introspection (see, e.g., [12]) but in this paper we omit it for simplicity.

In addition to the reasoning actions listed above, we therefore add the following actions:

Other's Resolution: an agent A can perform this action if it believes that another agent B believes two resolvable clauses c_1 and c_2 . Then A can conclude that B will believe in the resolvent clause c of c_1 and c_2 in the next time point. As a general case, we can extend the chain *agent-believes ... agent-believes*. For example, if agent A believes that agent B believes that agent C believes two resolvable clauses c_1 and c_2 , then it is possible in the next time point that agent A believes that agent B believes that agent C believes c which is the resolvent of c_1 and c_2 .

Positive Introspection: if an agent A believes a clause c , it can perform this action to reach a state where it believes that it believes c .

Other's Positive Introspection: if an agent A believes that another agent B believes a clause c , it can perform this action to reach a state where it believes that B believes that B believes c .

The reasoning actions *Positive Introspection* and *Other's Positive Introspection* are only enabled if the derived formula has a depth of nesting of at most n_B .

Note that the assumption that the agents reason using resolution and positive introspection is not essential for the main argument of this paper. This particular set of inference rules has been chosen to make the logic concrete; we could have, for example, assumed that the agents reason using modus ponens and conjunction introduction instead of resolution. In what follows, we give a formal definition of an epistemic logic for communicating agents which reason in a step-wise, memory-bounded fashion using some well-defined set of inference rules.

3 Syntax and semantics of *ERBL*

In this section, we give the syntax and semantics of the logic *ERBL* which formalises the ideas sketched in the previous section. *ERBL* (Epistemic Resource Bounded Logic) is an epistemic and temporal meta-language in which we can talk about beliefs expressed in the agents' internal language.

Let the set of agents be $A = \{1, 2, \dots, n_A\}$. We assume that all agents agree on a finite set *PROP* of propositional variables, and that all *belief formulas* of the internal

language of the agents are in the form of *clauses* or clauses preceded by a prefix of belief operators of fixed length.

From the set of propositional variables, we have the definition of all literals as follows:

$$LPROP = \{p, \neg p \mid p \in PROP\}$$

Then, the set of all clauses is $\Omega = \wp(LPROP)$. Finally, the set of all belief formulas is defined as follows:

$$B\Omega ::= \{B_{i_1} \dots B_{i_k} c \mid c \in \Omega, 0 \leq k \leq n_B\},$$

where $i_j \in A$. Note that we only include in the set of belief formulas those whose belief operator nesting is limited by n_B . Therefore, $B\Omega$ is finite.

Each agent $i \in A$ is assumed to have a knowledge base $KB_i \subseteq B\Omega$.

For convenience, the negation of a literal L is defined as $\neg L$, where:

$$\neg L = \begin{cases} \neg p & \text{if } L = p \text{ for some } p \in PROP \\ p & \text{if } L = \neg p \text{ for some } p \in PROP \end{cases}$$

The form of resolution rule which will be used in formal definitions below is as follows: given two clauses c_1 and $c_2 \in \Omega$ such that one contains a literal L and the other has its negation $\neg L$, we can derive a new clause which is the union $c_1 \setminus \{L\} \cup c_2 \setminus \{\neg L\}$.

The syntax of *ERBL* is then defined inductively as follows.

- \top is a well-formed formula (wff) of *ERBL*.
- *start* is a wff of *ERBL*; it is a marker for the start state.
- $cp_i^{\bar{n}}$ (the number of communication actions performed by agent i) is a wff of *ERBL* for all $n = 0, \dots, n_C$, and $i \in A$; it is used as a communication counter in the language.
- If $\alpha \in B\Omega$, then $B_i\alpha$ (agent i believes α) is a wff of *ERBL*, $i \in A$.
- If φ and ψ are wffs, then so are $\neg\varphi$, $\varphi \wedge \psi$.
- If φ and ψ are wffs, then so are $X\varphi$ (φ holds in the next moment of time), $\varphi U \psi$ (φ holds until ψ), and $A\varphi$ (φ holds on all paths).

Classical abbreviations for \vee , \rightarrow , \leftrightarrow are defined as usual. We also have $\perp \equiv \neg\top$, $F\varphi \equiv \top U \varphi$ (φ holds some time in the future), $E\varphi \equiv \neg A \neg \varphi$ (φ holds on some path). For convenience, let $CP_i = \{cp_i^{\bar{n}} \mid n = \{0, \dots, n_C\}\}$ and $CP = \bigcup_{i \in A} CP_i$.

The semantics of *ERBL* is defined by *ERBL* transition systems which are based on ω -tree structures (standard CTL* models as defined in [23]).

Let (T, R) be a pair where T is a set and R is a binary relation on T . Let the relation $<$ be the irreflexive and transitive closure of R , namely the set of pairs of states $\{(s, t) \in T \times T \mid \exists n \geq 0, t_0 = s, \dots, t_n = t \in T \text{ such that } t_i R t_{i+1} \text{ for all } i = 0, \dots, n-1\}$. (T, R) is a ω -tree frame iff the following conditions are satisfied.

1. T is a non-empty set.
2. R is total, i.e., for all $t \in T$, there exists $s \in T$ such that tRs .
3. For all $t \in T$, the past $\{s \in T \mid s < t\}$ is linearly ordered by $<$.
4. There is a smallest element called the root, denoted by t_0 .

5. Each maximal linearly $<$ - ordered subset of T is order-isomorphic to the natural numbers.

A branch of (T, R) is an ω -sequence (t_0, t_1, \dots) such that t_0 is the root and $t_i R t_{i+1}$ for all $i \geq 0$. We denote by $B(T, R)$ the set of all branches of (T, R) .

A *ERBL* transition system M is defined as a triple (T, R, V) where:

- (T, R) is a ω -tree frame,
- $V : T \times A \rightarrow \wp(B\Omega \cup CP)$ such that for all $s \in T$ and $i \in A$: $V(s, i) = Q \cup \{cp_i^{\leq n}\}$ for some $Q \subseteq B\Omega$ and $0 \leq n \leq n_C$. We denote by $V^*(s, i)$ the set $V(s, i) \setminus \{cp_i^{\leq n} | 0 \leq n\}$.

For a branch $\sigma \in B(T, R)$, σ_i denotes the element t_i of σ and $\sigma_{\leq i}$ is the prefix (t_0, t_1, \dots, t_i) of σ .

The truth of a *ERBL* formula at a point n of a path $\sigma \in B(T, R)$ is defined inductively as follows:

- $M, \sigma, n \models \top$,
- $M, \sigma, n \models B_i \alpha$ iff $\alpha \in V(s, i)$,
- $M, \sigma, n \models \text{start}$ iff $n = 0$,
- $M, \sigma, n \models cp_i^{\leq m}$ iff $cp_i^{\leq m} \in V(s, i)$,
- $M, \sigma, n \models \neg \varphi$ iff $M, \sigma, n \not\models \varphi$,
- $M, \sigma, n \models \varphi \wedge \psi$ iff $M, \sigma, n \models \varphi$ and $M, \sigma, n \models \psi$,
- $M, \sigma, n \models X \varphi$ iff $M, \sigma, n+1 \models \varphi$,
- $M, \sigma, n \models \varphi U \psi$ iff $\exists m \geq n$ such that $\forall k \in [n, m)$ $M, \sigma, k \models \varphi$ and $M, \sigma, m \models \psi$,
- $M, \sigma, n \models A \varphi$ iff $\forall \sigma' \in BR$ such that $\sigma'_{\leq n} = \sigma_{\leq n}$, $M, \sigma', n \models \varphi$.

The set of possible transitions in a model is defined as follows. Definition 1 below describes possible outcomes of various actions. For example, performing a resolution results in adding the resolvent to the set of beliefs. Definition 2 describes when an action is possible or enabled. For example, resolution is enabled if the agent has two resolvable clauses in memory.

Definition 1. Let (T, R, V) be a tree model. The set of effective transitions R_a for an action a is defined as a subset of R and satisfies the following conditions, for all $(s, t) \in R$:

1. $(s, t) \in R_{\text{Read}_{i, \alpha, \beta}}$ iff $\alpha \in KB_i$, $\alpha \notin V(s, i)$ and $V(t, i) = V(s, i) \setminus \{\beta\} \cup \{\alpha\}$. This condition says that s and t are connected by agent i 's *Read* transition if the following is true: α is in i 's knowledge base but not in $V(s, i)$, α is added to the set of i 's beliefs at t , and $\beta \in B\Omega$ is removed from the agent's set of beliefs. The argument β stands for a formula which is overwritten in the transition. If $\beta \in V(s, i)$ then the agent actually loses a belief in the transition, if $\beta \notin V(s, i)$ then the transition only involves adding a formula α without removing any beliefs.
2. $(s, t) \in R_{\text{Res}_{i, \alpha_1, \alpha_2, L, \beta}}$ where $\alpha_1 = B_{i_1} \dots B_{i_{k-1}} B_{i_k} c_1$ and $\alpha_2 = B_{i_1} \dots B_{i_{k-1}} B_{i_k} c_2$ iff $\alpha_1 \in V(s, i)$, $\alpha_2 \in V(s, i)$, $L \in c_1$, $\neg L \in c_2$, $\alpha = B_{i_1} \dots B_{i_{k-1}} B_{i_k} c$ and $V(t, i) = V(s, i) \setminus \{\beta\} \cup \{\alpha\}$ where $c = c_1 \setminus \{L\} \cup c_2 \setminus \{\neg L\}$. This condition says that s and t are connected by agent i 's *Res* transition if in s agent i believes

two resolvable clauses α_1 and α_2 but not α , possibly preceded by the same sequence of belief operators, and in t agent i believes their resolvent, preceded by the same prefix. Again, $\beta \in B\Omega$ is overwritten if it is in the set of agent's beliefs in s .

3. $(s, t) \in R_{Copy_{i,\alpha,\beta}}$ iff $\alpha \in V(s, j)$ for some $j \in A$ and $j \neq i$, for any $cp_i^{\bar{n}} \in V(s, i)$ such that $n < n_C$, $B_j\alpha \notin V(s, i)$ and $V(t, i) = V(s, i) \setminus \{cp_i^{\bar{n}} | cp_i^{\bar{n}} \in V(s, i)\} \cup \{cp_i^{\bar{n}+1} | cp_i^{\bar{n}} \in V(s, i)\} \setminus \{\beta\} \cup \{B_j\alpha\}$. s and t are connected by a *Copy* transition of agent i if in t , i adds to its beliefs a formula $B_j\alpha$ where α is an agent j 's belief in s , and i has previously copied fewer than n_C formulas. Again some $\beta \in B\Omega$ is possibly overwritten.
4. $(s, t) \in R_{Idle_i}$ iff $V(t, i) = V(s, i)$. The *Idle* transition does not change the state.
5. $(s, t) \in R_{Erase_{i,\beta}}$ iff $V(t, i) = V(s, i) \setminus \{\beta\}$. *Erase* removes one of the agent's beliefs.
6. $(s, t) \in R_{PI_{i,\alpha,\beta}}$ iff $\alpha \in V(s, i)$, $B_i\alpha \notin V(s, i)$ and $V(t, i) = V(s, i) \setminus \{\beta\} \cup \{B_i\alpha\}$. *PI* is i 's positive introspection: s and t are connected by i 's *PI* transition if in s it believes α but not $B_i\alpha$ and in t it believes $B_i\alpha$.
7. $(s, t) \in R_{OPI_{i,B_{i_1} \dots B_{i_{k-1}}, B_{i_k}\alpha,\beta}}$ iff $B_{i_1} \dots B_{i_{k-1}} B_{i_k}\alpha \in V(s, i)$ but not $B_{i_1} \dots B_{i_{k-1}} B_{i_k} B_{i_k}\alpha$, $V(t, i) = V(s, i) \setminus \{\beta\} \cup \{B_{i_1} \dots B_{i_{k-1}} B_{i_k} B_{i_k}\alpha\}$. This corresponds to ascribing positive introspection to agent i_k .

This specifies the effects of actions. Below, we specify when an action is possible. Note that we only enable deriving a formula if this formula is not already in the set of the agent's beliefs.

Definition 2. Let (T, R, V) be a tree model. The set $Act_{s,i}$ of possible actions that an agent i can perform at a state $s \in T$ is defined as follows:

1. $Read_{i,\alpha,\beta} \in Act_{s,i}$ iff $\alpha \notin V(s, i)$, $\alpha \in KB_i$ and $\beta \in V(s, i)$ if $|V^*(s, i)| \geq n_M$.
2. $Res_{i,\alpha_1,\alpha_2,L,\beta} \in Act_{s,i}$ iff $c = (c_1 \setminus L) \cup (c_2 \setminus \neg L) \notin V(s, i)$, $\alpha_1 = B_{i_1} \dots B_{i_{k-1}} B_{i_k} c_1$, $\alpha_2 = B_{i_1} \dots B_{i_{k-1}} B_{i_k} c_2$, $L \in c_1$, $\neg L \in c_2$, $\alpha_1, \alpha_2 \in V(s, i)$, and $\beta \in V(s, i)$ if $|V^*(s, i)| \geq n_M$.
3. $Copy_{i,\alpha,\beta} \in Act_{s,i}$ iff $B_j\alpha \notin V(s, i)$, $\alpha \in V(s, j)$ for some $j \in A$ and $j \neq i$, $n < n_C$ for any $cp_i^{\bar{n}} \in V(s, i)$ and $\beta \in V(s, i)$ if $|V^*(s, i)| \geq n_M$.
4. It is always the case that $Idle_i \in Act_{s,i}$.
5. $PI_{i,\alpha,\beta} \in Act_{s,i}$ iff $i\alpha \notin V(s, i)$, $\alpha \in V(s, i)$ and $\beta \in V(s, i)$ if $|V^*(s, i)| \geq n_M$.
6. $OPI_{i,B_{i_1} \dots B_{i_{k-1}}, B_{i_k}\alpha,\beta} \in Act_{s,i}$ iff $B_{i_1} \dots B_{i_{k-1}} B_{i_k} B_{i_k}\alpha \notin V(s, i)$, $B_{i_1} \dots B_{i_{k-1}} B_{i_k}\alpha \in V(s, i)$ and $\beta \in V(s, i)$ if $|V^*(s, i)| \geq n_M$.

There are no specified conditions for enabling $Erase_{i,\beta}$. This action is introduced for technical reasons, to simplify the proofs.

Finally, the definition of the set of models corresponding to a system of reasoners is given below:

Definition 3. $M(KB_1, \dots, KB_{n_A}, n_B, n_M, n_C)$ is the set of models (T, R, V) which satisfies the following conditions:

1. $|V^*(s, i)| \leq n_M$ for all $s \in T$ and $i \in A$.
2. $cp_i^{\bar{0}} \in V(t_0, i)$ where t_0 is the root of (T, R) for all $i \in A$.
3. $R = \bigcup_{a \in A} R_a$.
4. For all $s \in T$, $a_i \in Act_{s,i}$, there exists $t \in T$ such that $(s, t) \in R_{a_i}$ for all $i \in A$.

4 Axiomatisation

In this section, we introduce an axiom system which is sound and complete with respect to the set of models defined in the previous section.

Below are some abbreviations which will be used in the axiomatisation:

- $ByRead_i(\alpha, n) = \neg B_i \alpha \wedge cp_i^{\overline{n}}$. This formula describes the state before the agent comes to believe formula α by the *Read* transition. n is the value of i 's communication counter.
- $ByRes_i(\alpha, n) = \perp$ if $\alpha = B_{i_1} \dots B_{i_{k-1}} \neg B_{i_k} c$ for some $c \in \Omega$ and $1 \leq k \leq n_B$; otherwise $ByRes_i(\alpha, n) = \neg B_i \alpha \wedge \bigvee_{(\alpha_1, \alpha_2) \in Res^{-1}(\alpha)} (B_i \alpha_1 \wedge B_i \alpha_2)$ where $Res^{-1}(B_{i_1} \dots B_{i_{k-1}} B_{i_k} c) = \{(B_{i_1} \dots B_{i_{k-1}} B_{i_k} c_1, B_{i_1} \dots B_{i_{k-1}} B_{i_k} c_2) \mid \exists L \in LPROP \text{ such that } c = c_1 \setminus \{L\} \cup c_2 \setminus \{\neg L\}\}$. This formula describes the state of the system before i derives α by resolution. Note that it may not be possible to derive an arbitrary formula α by resolution; in that case, the state is described by falsum \perp .
- $ByCopy_i(\alpha, n) = \perp$ if $\alpha \neq B_j \alpha'$ for some $j \neq i$ or $n \leq 0$; otherwise $ByCopy_i(B_j \alpha', n) = \neg B_i B_j \alpha' \wedge B_j \alpha' \wedge cp_i^{\overline{n-1}}$.
- $ByPI_i(\alpha, n) = \perp$ if $\alpha \neq B_i \alpha'$; otherwise $ByPI_i(\alpha, n) = \neg B_i B_i \alpha' \wedge B_i \alpha' \wedge cp_i^{\overline{n}}$.
- $ByOPI_i(\alpha, n) = \perp$ if $\alpha \neq B_{i_1} \dots B_{i_{k-1}} B_{i_k} B_{i_k} \alpha'$; otherwise $ByOPI_i(\alpha, n) = \neg B_i B_{i_1} \dots B_{i_{k-1}} B_{i_k} B_{i_k} \alpha' \wedge B_i B_{i_1} \dots B_{i_{k-1}} B_{i_k} \alpha' \wedge cp_i^{\overline{n}}$.

The axiomatisation is as follows.

- A1.** All axioms and inference rules of CTL^* [24].
- A2.** $\bigwedge_{\gamma \in Q} B_i \gamma \wedge cp_i^{\overline{n}} \wedge \neg B_i \alpha \rightarrow EX(\bigwedge_{\gamma \in Q} B_i \gamma \wedge cp_i^{\overline{n}} \wedge B_i \alpha)$ for all $\alpha \in KB_i$, and $Q \subseteq B\Omega$ such that $|Q| < n_M$.
Intuitively, this axiom says that it is always possible to make a transition to a state where agent i believes a formula from its knowledge base KB_i . In addition, the communication counter of the agent does not increase, and a set of beliefs Q of cardinality less than n_M can also be carried over to the same state.
- Axioms **A3** - **A6** similarly describe transitions made by resolution (given that resolvable clauses are available), copy (with communication counter increased), and positive introspection (applied by agent i or ascribed by i to another agent).
- A3.** $\bigwedge_{\gamma \in Q} B_i \gamma \wedge B_i B_{i_1} \dots B_{i_{k-1}} B_{i_k} c_1 \wedge B_i B_{i_1} \dots B_{i_{k-1}} B_{i_k} c_2 \wedge cp_i^{\overline{n}} \wedge \neg B_i B_{i_1} \dots B_{i_{k-1}} B_{i_k} c \rightarrow EX(\bigwedge_{\gamma \in Q} B_i \gamma \wedge cp_i^{\overline{n}} \wedge B_i B_{i_1} \dots B_{i_{k-1}} B_{i_k} c)$ for all $c_1, c_2 \in \Omega$ such that $L \in c_1, \neg L \in c_2$ and $c = c_1 \setminus \{L\} \cup c_2 \setminus \{\neg L\}$, $k \geq 0$, and $Q \subseteq B\Omega$ such that $|Q| < n_M$.
- A4.** $\bigwedge_{\gamma \in Q} B_i \gamma \wedge B_j \alpha \wedge cp_i^{\overline{n}} \wedge \neg B_i B_j \alpha \rightarrow EX(\bigwedge_{\gamma \in Q} B_i \gamma \wedge B_i B_j \alpha \wedge cp_i^{\overline{n+1}})$ for any $\alpha \in B\Omega$, $j \in A$, $j \neq i$, $n < n_C$, and $Q \subseteq B\Omega$ such that $|Q| < n_M$.
- A5.** $\bigwedge_{\gamma \in Q} B_i \gamma \wedge B_i \alpha \wedge cp_i^{\overline{n}} \wedge \neg B_i B_i \alpha \rightarrow EX(\bigwedge_{\gamma \in Q} B_i \gamma \wedge B_i B_i \alpha \wedge cp_i^{\overline{n}})$ for any $\alpha \in B\Omega$ and $Q \subseteq B\Omega$ such that $|Q| < n_M$.
- A6.** $\bigwedge_{\gamma \in Q} B_i \gamma \wedge B_i B_{i_1} \dots B_{i_{k-1}} B_{i_k} \alpha \wedge cp_i^{\overline{n}} \wedge \neg B_i B_{i_1} \dots B_{i_{k-1}} B_{i_k} B_{i_k} \alpha \rightarrow EX(\bigwedge_{\gamma \in Q} B_i \gamma \wedge B_i B_{i_1} \dots B_{i_{k-1}} B_{i_k} B_{i_k} \alpha \wedge cp_i^{\overline{n}})$ for any $\alpha \in B\Omega$, $k \geq 0$ and $Q \subseteq B\Omega$ such that $|Q| < n_M$.

A7. $EX(B_i\alpha \wedge B_i\beta) \rightarrow B_i\alpha \vee B_i\beta$.

This axiom says that at most one new belief is added in the next state.

A8. $EX(\neg B_i\alpha \wedge \neg B_i\beta) \rightarrow \neg B_i\alpha \vee \neg B_i\beta$.

This axiom says that at most one belief is deleted in the next state.

A9. $EX(B_i\alpha \wedge cp_i^{-n}) \rightarrow B_i\alpha \vee ByRead_i(\alpha, n) \vee ByRes_i(\alpha, n) \vee ByCopy_i(\alpha, n) \vee ByPI_i(\alpha, n) \vee ByOPI_i(\alpha, n)$ for $\alpha \in KB_i$.

This axiom says that a new belief which is an element of the agent's knowledge base can only be added by one of the valid reasoning actions.

A10. $EX(B_i\alpha \wedge cp_i^{-n}) \rightarrow B_i\alpha \vee ByRes_i(\alpha, n) \vee ByCopy_i(\alpha, n) \vee ByPI_i(\alpha, n) \vee ByOPI_i(\alpha, n)$ for $\alpha \notin KB_i$.

This axiom describes possible ways in which a new belief which is not in the agent's knowledge base can be added.

A11. $B_i\alpha_1 \wedge \dots \wedge B_i\alpha_{n_M} \rightarrow \neg B_i\alpha_{n_M+1}$ for all $i \in A$, $\alpha_j \in B\Omega$ where $j = 1, \dots, n_M + 1$ and all α_j are pairwise different.

This axiom states that an agent cannot have more than n_M different beliefs.

A12a $start \rightarrow cp_i^{-0}$ for all $i \in A$.

In the start state, the agent has not performed any *Copy* actions.

A12b $\neg EX start$ (*start* only holds at the root of the tree).

A13. $\bigvee_{n=0..n_C} cp_i^{-n}$ for all $i \in A$.

There is always a number n between 0 and n_C corresponding to the number of *Copy* actions agent i has performed.

A14. $cp_i^{-n} \rightarrow \neg cp_i^{-n'}$ for all $i \in A$ and $n' \neq n$.

The number of previous *Copy* actions by i in each state is unique.

A15. $\varphi \rightarrow EX\varphi$ where φ does not contain *start*.

It is always possible to make a transition to a state where all agents have the same beliefs and communication counter values as in the current state (essentially an *Idle* transition by all agents)

A16. $\bigwedge_{i \in A} EX(\bigwedge_{\gamma \in Q_i} B_i\gamma \wedge cp_i^{-n_i}) \rightarrow EX \bigwedge_{i \in A} (\bigwedge_{\gamma \in Q_i} B_i\gamma \wedge cp_i^{-n_i})$ for any $Q_i \subseteq B\Omega$ such that $|Q_i| \leq n_M$.

If each agent i can separately reach a state where it believes formulas in Q_i , then all agents together can reach a state where for each i , agent i believes formulas in Q_i .

Notice that since the depth of the nesting of belief operators is restricted by n_B , for any subformula $B_i\alpha$ appearing in any above axiom, $\alpha \in B\Omega$.

Definition 4. $L(KB_1, \dots, KB_{n_A}, n_B, n_M, n_C)$ is the logic defined by the axiomatisation **A1–A16**.

We have the following result.

Theorem 1. $L(KB_1, \dots, KB_{n_A}, n_B, n_M, n_C)$ is sound and complete with respect to $M(KB_1, \dots, KB_{n_A}, n_B, n_M, n_C)$.

The proof is omitted due to lack of space; it is based on the proof technique used in [24].

5 Discussion

Systems of step-wise reasoners with bounded memory and a communication limit are faithful models of systems of distributed resource-limited reasoners, and various resource requirements of such systems can be effectively verified, e.g. by model-checking, as in for example [20]. However, adding reasoning about beliefs poses a significant challenge, both in the complexity of the system and in the way this reasoning is modelled. The branching factor of the models is much larger when reasoning about beliefs is considered, making model-checking less feasible. The main problem however has to do with the correctness of an agent's belief ascription. We describe this problem below and propose a tentative solution.

In the system proposed in this paper, agents correctly ascribe reasoning mechanisms to each other, and in the limit, their predictions concerning other agents' beliefs are correct: if agent j believes that eventually agent i will believe α , then eventually agent i will believe α , and vice versa. More precisely, for every model M and state s ,

$$\{\alpha : M, s \models EFB_j B_i \alpha\} = \{\alpha : M, s \models EFB_i \alpha\}$$

However, in spite of this, the agents are almost bound to make wrong predictions when trying to second-guess what other reasoners will believe in the next state. More precisely,

$$\{\alpha : M, s \models B_j B_i \alpha\} \not\subseteq \{\alpha : M, s \models B_i \alpha\}$$

i.e. agent j may believe that i believes some α when i does not believe α .

Consider the following example. Suppose there are two agents, 1 and 2, each with a memory limit of two formulas, communication limit of one formula, belief nesting limit of two, and knowledge bases $KB_1 = \{p\}$ and $KB_2 = \{q\}$. A possible run of the system is shown in Figure 1.

| State | Agent 1 | Agent 2 |
|-------------|----------------|----------------|
| t_0 | { } | { } |
| transition: | Read | Read |
| t_1 | { p } | { q } |
| transition: | Copy | Copy |
| t_2 | { $p, B_2 q$ } | { $q, B_1 p$ } |

Fig. 1. A possible run of the system

Note that this is only one possible run, and other transitions are possible. For example, in t_0 , one or both agents can idle. In t_1 , one or both agents can idle, or make a positive introspection transition. In state t_2 , the agents' beliefs about each other's beliefs are correct. However, in most successor states of t_2 , agent 1 will have incorrect beliefs about agent 2's beliefs, and vice versa. Indeed, the options of agent 1 in t_2 are: read p , idle, erase p , erase $B_2 q$, apply positive introspection to derive $B_1 p$ or $B_1 B_2 q$, ascribe introspection to agent 2 to derive $B_2 B_2 q$. Agent 2 has similar choices. In only

two of these cases do the agents make non-trivial (that is, new compared to the ones already existing in t_2) correct belief ascriptions, namely if agent 1 derives B_1p and agent 2 derives B_1B_1p , and vice versa when agent 2 derives B_2q and agent 1 derives B_2B_2q (see Figure 2).

| State | Agent 1 | Agent 2 |
|-------------|----------------------|---------------------|
| t_2 | $\{p, B_2q\}$ | $\{q, B_1p\}$ |
| transition: | PI, overwrite B_2q | OPI, overwrite q |
| t_3 | $\{p, B_1p\}$ | $\{B_1p, B_1B_1p\}$ |

Fig. 2. Continuing from t_2 : a correct ascription

Figure 3 shows one of many possible incorrect ascriptions. Note that agent 1's ascription is now incorrect because agent 2 has forgotten q , and agent 2's ascription is incorrect because it assumed agent 1 will use positive introspection to derive B_1p , which it did not.

| State | Agent 1 | Agent 2 |
|-------------|---------------|---------------------|
| t_2 | $\{p, B_2q\}$ | $\{q, B_1p\}$ |
| transition: | Idle | OPI, overwrite q |
| t_4 | $\{p, B_2q\}$ | $\{B_1p, B_1B_1p\}$ |

Fig. 3. Continuing from t_2 : an incorrect ascription

This suggests an inherent problem with modelling agents reasoning about each other's beliefs in a step-wise, memory-bounded fashion. Note that this problem is essentially one of belief ascription, i.e., of correctly predicting what another agent will believe given limited information about what it currently believes (of deriving correct conclusions from correct premises), rather than a problem of belief revision [25], i.e., what an agent should do if it discovers the beliefs it has ascribed to another agent are incorrect. It is also distinct from the problem of determining the consequences of information updates as studied in dynamic epistemic logic (e.g. [26]). Adding new true beliefs in a syntactic approach such as ours is straightforward compared to belief update in dynamic epistemic logic, which interprets beliefs as sets of possible worlds. Essentially, in dynamic epistemic logic an agent acquires a new logically closed set of beliefs at the next 'step' after an announcement is made, while we model the gradual process of deriving consequences from a new piece of information (and the agent's previous beliefs).

The disparity between agent i 's beliefs and the beliefs agent j ascribes to i at each step is due both to the fact that at most one formula is derived by each agent at any given step (and agent j may guess incorrectly which inference rule agent i is going to use) and to memory limitations which cause agents to forget formulas. An obvious alternative is to do tentative ascription of beliefs to other agents, namely conclude that

the other agent will be in *one of several* possible belief sets in the next state, e.g.

$$B_2B_1p \rightarrow EX(B_2((B_1p \wedge B_1B_1p) \vee (B_1p \wedge \neg B_1B_1p) \vee \dots))$$

However, this implies that one of the agents (agent 2 in this case) has a much larger (exponentially larger!) memory and a more expressive internal language to reason about the other agent's beliefs.

It is clearly not sufficient for correct belief prediction for the reasoners to ascribe to other agents just a set of inferences rules or a logic such as KD45. They need to be able to ascribe to other agents a *reasoning strategy*, or a preference order on the set reasoning actions used by the other agents which constrains the possible transitions of each reasoner, and directs each agent's reasoning about the beliefs of other agents. As a simple example, suppose agent 2 believes that agent 1's strategy is to apply positive introspection to formula p in preference to all other actions. Then in state t_2 agent 2 will derive B_1B_1p from B_1p . If agent 2's ascription of strategy to agent 1 is correct, agent 1 will indeed derive B_1p from p in the next state, making agent 2's belief prediction correct.

6 ERBL with strategies

In this section, we modify the semantics of *ERBL* to introduce reasoning strategies.

First we need to define strategies formally. A *reasoning strategy for agent i* , \prec_i , is a total order on the set Act_i of all reasoning actions of i and their arguments:

$$Act_i = \{Read_{i,\alpha,\beta}, Res_{i,\alpha_1,\alpha_2,L,\beta}, Copy_{i,\alpha,\beta}, \\ Erase_{i,\beta}, Idle_i, PI_{i,\alpha,\beta}, OPI_{i,B_{i_1} \dots B_{i_{k-1}}, B_{i_k} \alpha,\beta} \mid \alpha, \beta, \alpha_1, \alpha_2 \in B\Omega\}$$

A simple example of a reasoning strategy for i would be a lexicographic order on Act_i which uses two total orders: an order on the set of transitions, e.g. $Res < PI < OPI < Copy < Read < Idle$, and an order on $B\Omega$.

Recall that in Definition 2 we specified which actions are enabled in state s , $Act_{s,i} \subseteq Act_i$. We required in Definition 3 that for each enabled action, there is indeed a transition by that action out of s . The simple change that we make to Definition 3 is that for every agent i we only enable *one* action, namely the element of $Act_{s,i}$ which is minimal in \prec_i .

Definition 5. *The set of reasoning strategy models $M^{strat}(KB_1, \dots, KB_{n_A}, n_B, n_M, n_C)$ is the set of models (T, R, V) which satisfies conditions 1-3 from Definition 3 and the following condition:*

- 4'.** *For all $s \in T$, there exists a unique state t such that $(s, t) \in R_{a_i}$ for all $i \in A$, where a_i is the minimal element with respect to \prec_i in $Act_{s,i}$.*

Observe that in the reasoning strategy models, the transition relation is a linear order.

Finally, we give one possible definition of a correct ascription of a reasoning strategy which allows an agent j to have a correct and complete representation of the beliefs

of another agent i , namely ensuring that $B_i\alpha \leftrightarrow B_jB_i\alpha$ at each step. Such perfect matching of i 's beliefs by j is possible if

$$KB_j = \{B_i\alpha : \alpha \in KB_i\}$$

and agent i does not use the *Copy* action (intuitively, because in order to match *Copy* by i , agent j has to add two modalities in one step: when agent i derives $B_l\alpha$ from α being in agent l 's belief set, agent j has to derive $B_iB_l\alpha$). Below, we also assume that j is allowed one extra nesting of belief modalities ($n_B(j) = n_B(i) + 1$).

Definition 6. *Agent j has a strategy which matches the strategy of agent i if for every natural number k , the following correspondence holds between the k th element of \prec_j and the k th element of \prec_i :*

- if the k th element of \prec_i is $Read_{i,\alpha,\beta}$, then the k th element of \prec_j is $Read_{j,B_i\alpha,B_i\beta}$
- if the k th element of \prec_i is $Res_{i,\alpha_1,\alpha_2,L,\beta}$, then the k th element of \prec_j is $Res_{j,B_i\alpha_1,B_i\alpha_2,L,B_i\beta}$
- if the k th element of \prec_i is $PI_{i,\alpha,\beta}$, then the k th element of \prec_j is $OPI_{j,B_i\alpha,B_i\beta}$
- if the k th element of \prec_i is $OPI_{i,B_l\alpha,\beta}$, then the k th element of \prec_j is $OPI_{j,B_iB_l\alpha,B_i\beta}$
- if the k th element of \prec_i is $Erase_{i,\beta}$, then the k th element of \prec_j is $Erase_{j,B_i\beta}$
- if the k th element of \prec_i is $Idle_i$, then the k th element of \prec_j is $Idle_j$.

Theorem 2. *If agent j 's strategy matches agent i 's strategy and agent j has complete and correct beliefs about agent i 's beliefs in state s : $M, s \models B_i\alpha \leftrightarrow B_jB_i\alpha$, then agent j will always have correct beliefs about agent i 's beliefs: $M, s \models AG(B_i\alpha \leftrightarrow B_jB_i\alpha)$.*

Other more realistic matching strategies, for example, those which allow the agent to have a less than complete representation of other agent's beliefs, are possible, and their formal investigation is a subject of future work.

7 Conclusion

We presented a formal model of resource-bounded reasoners reasoning about each other's beliefs, and a sound and complete logic, *ERBL*, for reasoning about such systems. Our formalisation highlighted a problem of incorrect belief ascription, and we showed that this problem can be overcome by extending the framework with reasoning strategies. In future work we plan to extend the framework in a number of ways, including producing correct belief ascription under less strict matching between agents' strategies, and introducing reasoning about other agent's resource limitations. At the moment the agents have no way of forming beliefs about another agent's memory limit n_M or belief nesting bound n_B (note that we can also easily make those limits different for different agents). If they could represent those limitations, then one agent could infer that another agent does not believe some formula on the grounds that the latter agent's memory is bounded.

References

1. Fagin, R., Halpern, J.Y., Moses, Y., Vardi, M.Y.: Reasoning about Knowledge. MIT Press, Cambridge, Mass. (1995)
2. Meyer, J.J., van der Hoek, W.: Epistemic Logic for Computer Science and Artificial Intelligence. Cambridge University Press (1995)
3. Hintikka, J.: Knowledge and belief. Cornell University Press, Ithaca, NY (1962)
4. Rantala, V.: Impossible worlds semantics and logical omniscience. *Acta Philosophica Fennica* **35** (1982) 106–115
5. Fagin, R., Halpern, J.Y., Vardi, M.Y.: A non-standard approach to the logical omniscience problem. In Parikh, R., ed.: *Theoretical Aspects of Reasoning about Knowledge: Proceedings of the Third Conference*, Morgan Kaufmann (1990) 41–55
6. Fagin, R., Halpern, J.Y.: Belief, awareness and limited reasoning: Preliminary report. In: *Proceedings of the 9th International Joint Conference on Artificial Intelligence*. (1985) 491–501
7. Konolige, K.: *A Deduction Model of Belief*. Morgan Kaufmann, San Francisco, Calif. (1986)
8. Elgot-Drapkin, J.J., Perlis, D.: Reasoning situated in time I: Basic concepts. *Journal of Experimental and Theoretical Artificial Intelligence* **2** (1990) 75–98
9. Halpern, J.Y., Moses, Y., Vardi, M.Y.: Algorithmic knowledge. In Fagin, R., ed.: *Theoretical Aspects of Reasoning about Knowledge: Proceedings of the Fifth Conference (TARK 1994)*. Morgan Kaufmann, San Francisco (1994) 255–266
10. Pucella, R.: Deductive algorithmic knowledge. *J. Log. Comput.* **16**(2) (2006) 287–309
11. Duc, H.N.: Logical omniscience vs. logical ignorance on a dilemma of epistemic logic. In Pinto-Ferreira, C.A., Mamede, N.J., eds.: *Progress in Artificial Intelligence, 7th Portuguese Conference on Artificial Intelligence, EPIA '95, Funchal, Madeira Island, Portugal, October 3-6, 1995, Proceedings*. Volume 990 of *Lecture Notes in Computer Science.*, Springer (1995) 237–248
12. Duc, H.N.: Reasoning about rational, but not logically omniscient, agents. *Journal of Logic and Computation* **7**(5) (1997) 633–648
13. Ågotnes, T., Alechina, N.: The dynamics of syntactic knowledge. *Journal of Logic and Computation* **17**(1) (2007) 83–116
14. Sierra, C., Godo, L., de Mántaras, R.L., Manzano, M.: Descriptive dynamic logic and its application to reflective architecture. *Future Gener. Comput. Syst.* **12**(2-3) (1996) 157–171
15. Alechina, N., Jago, M., Logan, B.: Modal logics for communicating rule-based agents. In Brewka, G., Coradeschi, S., Perini, A., Traverso, P., eds.: *Proceedings of the 17th European Conference on Artificial Intelligence (ECAI 2006)*, IOS Press (2006) 322–326
16. Fisher, M., Ghidini, C.: Programming resource-bounded deliberative agents. In: *Proceedings of the Sixteenth International Joint Conference on Artificial Intelligence (IJCAI 1999)*, Morgan-Kaufmann (1999) 200–205
17. Ågotnes, T.: *A Logic of Finite Syntactic Epistemic States*. Ph.D. thesis, Department of Informatics, University of Bergen, Norway (2004)
18. Ågotnes, T., Alechina, N.: Knowing minimum/maximum n formulae. In Brewka, G., Coradeschi, S., Perini, A., Traverso, P., eds.: *Proceedings of the 17th European Conference on Artificial Intelligence (ECAI 2006)*, IOS Press (2006) 317–321
19. Albore, A., Alechina, N., Bertoli, P., Ghidini, C., Logan, B., Serafini, L.: Model-checking memory requirements of resource-bounded reasoners. In: *Proceedings of the Twenty-First National Conference on Artificial Intelligence (AAAI 2006)*, AAAI Press (2006) 213–218
20. Alechina, N., Logan, B., Nga, N.H., Rakib, A.: Verifying time, memory and communication bounds in systems of reasoning agents. In Padgham, L., Parkes, D., Müller, J., Parsons,

- S., eds.: Proceedings of the Seventh International Conference on Autonomous Agents and Multiagent Systems (AAMAS 2008). Volume 2., Estoril, Portugal, IFAAMAS, IFAAMAS (May 2008) 736–743
21. Ågotnes, T., Walther, D.: Towards a logic of strategic ability under bounded memory. In: Proceedings of the Workshop on Logics for Resource-Bounded Agents. (2007)
 22. Adjiman, P., Chatalic, P., Goasdoué, F., Rousset, M.C., Simon, L.: Scalability study of peer-to-peer consequence finding. In Kaelbling, L.P., Saffiotti, A., eds.: Proceedings of the Nineteenth International Joint Conference on Artificial Intelligence (IJCAI-05), Edinburgh, Scotland, Professional Book Center (2005) 351–356
 23. Emerson, E.A.: Temporal and modal logic. In: Handbook of Theoretical Computer Science, Volume B: Formal Models and Semantics (B). Elsevier and MIT Press (1990) 995–1072
 24. Reynolds, M.: An axiomatization of full computation tree logic. *Journal of Symbolic Logic* **66**(3) (2001) 1011–1057
 25. Alchourrón, C.E., Gärdenfors, P., Makinson, D.: On the logic of theory change: Partial meet functions for contraction and revision. *Journal of Symbolic Logic* **50** (1985) 510–530
 26. Baltag, A., Moss, L.S., Solecki, S.: The logic of public announcements, common knowledge, and private suspicions. In: Proceedings of the 7th conference on Theoretical aspects of rationality and knowledge (TARK'98). (1998)