

# On the Semantics of Belief Revision Systems<sup>\*†</sup>

Gösta Grahne<sup>‡</sup>  
 Dept. of Comp. Sci.  
 Univ. of Helsinki  
 Teollisuuskatu 23  
 SF-00510 Helsinki

Alberto O. Mendelzon  
 Dept. of Comp. Sci.  
 Univ. of Toronto  
 Toronto, Canada  
 M5S 1A4

Ray Reiter<sup>§</sup>  
 Dept of Comp. Sci.  
 Univ. of Toronto  
 Toronto, Canada  
 M5S 1A4

## Abstract

We give semantics to belief revision operators that satisfy the Alchourrón-Gärdenfors-Makinson postulates by presenting an epistemic logic such that, for any such revision operator, the result of a revision can be described by a sentence in this logic.

In our logic, the fact that the agent's set of beliefs is  $\phi$  is represented by the sentence  $0\phi$ , where  $0$  is Levesque's 'only know' operator. Intuitively,  $0\phi$  is read as ' $\phi$  is *all* that is believed.' The fact that the agent believes  $\psi$  is represented by the sentence  $B\psi$ , read in the usual way as ' $\psi$  is believed'. The connective  $\diamond$  represents *update* as defined by Katsuno and Mendelzon. The revised beliefs are represented by the sentence  $0\phi \diamond B\psi$ . We show that for every revision operator that satisfies the *AGM* postulates, there is a model for our epistemic logic such that the beliefs implied by the sentence  $0\phi \diamond B\psi$  in this model correspond exactly to the sentences implied by the theory that results from revising  $\phi$  by  $\psi$ . This means that reasoning about changes in the agent's beliefs reduces to model checking of certain epistemic sentences.

The negative result in the paper is that this type of formal account of revision cannot be extended to the situation where the agent is able to reason about its beliefs. A fully introspective agent cannot use our construction to reason about the results of its own revisions, on pain of triviality.

## 1 Introduction

Consider an agent that holds a set of beliefs and is sometimes required to change these beliefs according to its perceptions or other inputs. The problem of *belief revision* is: what are rational (or desirable, or typical, or economical) ways of revising a set of beliefs given some input? These questions have been studied in philosophy [7], in statistics [15], in Artificial Intelligence [3], and in the theory of databases [5].

---

<sup>\*</sup>This work was supported by the Natural Sciences and Engineering Research Council of Canada and by the Institute for Robotics and Intelligent Systems.

<sup>†</sup>Authors e-mail addresses: [grahne@cs.helsinki.fi](mailto:grahne@cs.helsinki.fi), [mendel@db.toronto.edu](mailto:mendel@db.toronto.edu), [reiter@cs.toronto.edu](mailto:reiter@cs.toronto.edu)

<sup>‡</sup>Work performed while visiting University of Toronto, Department of Computer Science.

<sup>§</sup>Fellow of the Canadian Institute of Advanced Research.

For concreteness, suppose, as in [16], that the agent's beliefs are represented by a propositional theory  $\phi$  and that the new input is a propositional sentence  $\psi$  that must be incorporated into the theory. A growing body of work [2, 21, 22] takes as a departure point the *rationality postulates* proposed of Alchourrón, Gärdenfors and Makinson (the *AGM* postulates, [1]). These are rules that every adequate revision operator should be expected to satisfy. For example: the new fact  $\psi$  must be a consequence of the revised theory.

Katsuno and Mendelzon [17] reason that no such set of postulates will be adequate for every application. In particular, they make a distinction between two kinds of modifications to a theory. The first one, *update*, consists of bringing the theory up to date when the world described by it changes. For example, most database updates are of this variety, e.g. 'increase Joe's salary by 5%'. Another example is the incorporation into the theory of changes in the world caused by the actions of an agent such as a robot [9, 23, 24]. The second type of modification, *revision*, is used when the agent obtains new information about a static world. For example, the agent may be trying to diagnose a faulty circuit and wanting to incorporate into the knowledgebase the results of successive tests, where newer results may contradict old ones. Katsuno and Mendelzon conclude that the *AGM* postulates describe only revision and they give a modified set of postulates for update (the *KM* postulates).

As an aside, we note that distinctions of a similar nature have been suggested in other contexts. Gibbard [8] compares indicative and subjunctive conditionals of revision and update. In the context of Bayesian models of belief revision, Lewis [20] and Gärdenfors [7] propose *imaging* as an alternative to *conditionalization*. As suggested to us in private communications by H. Arló-Costa and P. Gärdenfors, imaging appears to be a probabilistic counterpart of update.

The purpose of this paper is to give semantics for revision by presenting an epistemic logic such that, for any revision operator that satisfies the *AGM* postulates, the result of a revision can be described by a sentence in this logic.

In our logic, the fact that the agent's set of beliefs is  $\phi$  is represented by the sentence  $\mathbf{O}\phi$ , where  $\mathbf{O}$  is Levesque's 'only know' operator [19]. Intuitively,  $\mathbf{O}\phi$  is read as ' $\phi$  is *all* that is believed.' The fact that the agent believes  $\psi$  is represented by the sentence  $\mathbf{B}\psi$ , read in the usual way as ' $\psi$  is believed'. The connective  $\diamond$  represents *update* as defined in [17]. The revised beliefs are represented by the sentence  $\mathbf{O}\phi \diamond \mathbf{B}\psi$ .<sup>1</sup> We show that for every revision operator that satisfies the *AGM* postulates, there is a model for our epistemic logic such that the beliefs implied by the sentence  $\mathbf{O}\phi \diamond \mathbf{B}\psi$  in this model correspond exactly to the sentences implied by the theory that results from revising  $\phi$  by  $\psi$ . This means that reasoning about changes in the agent's beliefs reduces to model checking [13] of certain epistemic sentences.

Note that a revision in the agent's universe of discourse is translated into an update at the meta level. Intuitively, the reason is that update, and not revision, is the operator that serves to accommodate changes in the world. At the meta level, the world includes the agent's set of beliefs, since we can assert that the agent does or does not believe certain sentences; hence update is appropriate to incorporate the change in the world that occurs

<sup>1</sup>We shall leave out parenthesis symbols from sentences, with the convenience that unary operators bind stronger than binary, and that among the binary connectives, the update  $\diamond$  has the strongest binding. Thus for instance  $((\mathbf{O}\phi) \diamond (\mathbf{B}\psi)) \rightarrow (\mathbf{B}\chi)$  is written as  $\mathbf{O}\phi \diamond \mathbf{B}\psi \rightarrow \mathbf{B}\chi$ .

when the agent changes its mind about something.

For example, in a distributed system an agent (a node in a network) may revise its beliefs about the state of a computation. An account of this revision is necessary to reason formally about properties of such systems. An interesting extension would be to include multiple agents and a notion of ‘common knowledge’, as in Fagin, Halpern, Moses and Vardi [4], since a formal analysis of their ‘muddy children puzzle’ requires the concept of belief revision to describe how the children’s beliefs change as a consequence of the father’s statements.

It might be illuminating at this point to consider Levesque’s distinction between *subjective* and *objective* uses of logic: If the agent used logic to reason about its epistemic states, it would be using logic subjectively. The metalogic is used here objectively to *describe* the epistemic states of the agent.

## 2 Belief Revision Systems

In this section we define the language and the apparatus that the agent uses to represent, reason about, and change its beliefs. The language and reasoning is classical propositional logic, and the change mechanism is a revision operator.

Let  $\mathcal{L}$  be a propositional language with a countably infinite set  $P$  of propositional letters. The set  $S$  of all (well formed) sentences (in  $\mathcal{L}$ ) is defined in the usual way. Sentences will be denoted by  $\phi, \psi, \chi, \dots$ . A *theory* (in  $\mathcal{L}$ ) is a subset  $\Gamma$  of  $S$ , such that  $Cn(\Gamma) = \Gamma$ , where  $Cn$  denotes the propositional consequence operator. The set of all theories (in  $\mathcal{L}$ ) is denoted  $T$ . A sentence  $\phi$  is said to be a *theorem*, if  $\phi \in Cn(\emptyset)$ , and  $\phi$  is said to be *consistent*, if  $\neg\phi \notin Cn(\emptyset)$ . A theory  $\Gamma$  is said to be consistent, if  $\Gamma \neq S$ . A theory  $\Gamma$  is said to be *finitely axiomatizable* (in  $\mathcal{L}$ ), if there is a  $\phi \in S$ , such that  $Cn(\{\phi\}) = \Gamma$ .

The theory that an agent holds changes as new inputs are assimilated. To abstract the change mechanism we use a revision operator that maps a theory and a sentence to a new theory. Following Gärdenfors [7], a ( $\mathcal{L}$ -) *belief revision system* is a triple

$$\langle Cn, T, \oplus \rangle,$$

where  $\oplus$  is a function from  $T \times S$  to  $T$ , such that the following conditions hold:

- (R1)  $\phi \in \Gamma \oplus \phi$ .
- (R2) If  $\Gamma \cup \{\phi\}$  is consistent, then  $\Gamma \oplus \phi = Cn(\Gamma \cup \{\phi\})$ ,
- (R3) If  $\phi$  is consistent, then  $\Gamma \oplus \phi$  is consistent.
- (R4) If  $\phi \leftrightarrow \psi$  is a theorem, then  $\Gamma \oplus \phi = \Gamma \oplus \psi$ .
- (R5)  $\Gamma \oplus (\phi \wedge \psi) \subseteq Cn(\Gamma \oplus \phi \cup \{\psi\})$ .
- (R6) If  $\Gamma \oplus \phi \cup \{\psi\}$  is consistent, then  $Cn(\Gamma \oplus \phi \cup \{\psi\}) \subseteq \Gamma \oplus (\phi \wedge \psi)$ .

For an intuitive motivation of the conditions (the ‘revision postulates’, or ‘rationality criteria’), see [1, 7].

In any practical application, the agent's beliefs  $\Gamma$  will be of the form  $Cn(\phi)$ , for some  $\phi \in S$ . This practical aspect will however be ruined, if the same is no longer true after a revision. Consider therefore the following condition.

**(R0)**  $Cn(\{\phi\}) \oplus \psi$  is finitely axiomatizable.

If a belief revision  $\langle Cn, T, \oplus \rangle$  also satisfies condition (R0), we say that that  $\langle Cn, T, \oplus \rangle$  is a ( $\mathcal{L}$ -) *preservative* belief revision system.

### 3 The Update Language and its Interpretation

We now present the epistemic language that we use to describe an agent's beliefs and how they change under revision. This language  $\mathcal{L}_{\mathcal{L}}$  is constructed from the set  $P$  of propositional letters, the Boolean operators, and the following intensional connectives: a binary connective  $\diamond$  (update), and the unary connectives **B** and **N** (believe at least and believe at most). Intuitively, **B** $\phi$  will mean that the agent believes sentence  $\phi$  (and perhaps other things too), while **N** $\phi$  will mean that the agent does not believe anything that is not a consequence of  $\phi$ . **N** $\phi$  is only introduced for technical reasons; we will actually use  $\mathbb{O}\phi$ , which is defined in terms of **B** $\phi$  and **N** $\phi$  and means the agent believes precisely  $\phi$ .

In the set of well formed sentences of  $\mathcal{L}_{\mathcal{L}}$  we shall not include for instance sentences of the form **BB** $\phi$ , which intuitively would mean 'the agent believes that it believes  $\phi$ '. The reason for this will become clear in Section 5, where we consider the scenario where the agent is fully introspective. The set of all well formed sentences of  $\mathcal{L}_{\mathcal{L}}$  is thus defined as follows: Any propositional sentence (i.e. any sentence in  $S$ ) is a  $\mathcal{L}_{\mathcal{L}}$ -sentence, and if  $\phi$  is in  $S$ , then **B** $\phi$  and **N** $\phi$  are  $\mathcal{L}_{\mathcal{L}}$ -sentences. Furthermore, if  $\phi$  and  $\psi$  are  $\mathcal{L}_{\mathcal{L}}$ -sentences, then so is  $\phi \diamond \psi$ .

Now let us define formally the interpretation of  $\mathcal{L}_{\mathcal{L}}$ -sentences. An *update model* is a quadruple

$$\mathcal{M} = \langle W, R, \leq, \llbracket \cdot \rrbracket \rangle,$$

where  $W$  is a non-empty set of worlds  $w, v, u, \dots$ , where  $R$  is a binary relation over  $W$ , where  $\leq$  is a function that assigns a total pre-order  $\leq_w$  on  $W$  to each member  $w$  of  $W$ , and where  $\llbracket \cdot \rrbracket$  is a function that assigns a subset  $\llbracket \phi \rrbracket$  of  $W$  to each  $\mathcal{L}_{\mathcal{L}}$ -sentence  $\phi$ . The relation  $R$  is an *accessibility relation*. We use it to model belief as in the standard semantics of modal logic, that is,  $\phi$  will be believed at some world  $w$  iff  $\phi$  is true at all possible worlds accessible from  $w$ . For any world  $w \in W$  we shall use the notation  $R_w$  for the set of worlds accessible from  $w$ , that is,  $\{v \in W : wRv\}$ . The relation  $\leq_w$  is a *comparative similarity ordering of worlds w.r.t. world  $w$* . This relation will be used to define the semantics of update as in [10, 17]. The sentence  $\phi \diamond \psi$  will be true at a world  $v$  if  $v$  satisfies  $\psi$  and there exists a world  $w$  that satisfies  $\phi$  such that  $v$  is 'as close' to  $w$  as any other world that satisfies  $\psi$ , w.r.t. the comparative similarity ordering  $\leq_w$ .

In addition the following conditions have to be fulfilled.

**(S1)** For each sentence  $\phi$ , the set  $\min_{\leq_w}(\llbracket \phi \rrbracket)$  is non-empty, whenever  $\llbracket \phi \rrbracket$  is non-empty.<sup>2</sup>

<sup>2</sup>For  $V \subseteq W$  and  $w \in W$ ,  $\min_{\leq_w}(V)$  denotes the set  $\{v \in V : \text{if } u \in V \text{ and } u \leq_w v, \text{ then } v \leq_w u\}$ .

$$(S2) \min_{\leq_w}(W) = \{w\}.$$

$$(VC) \llbracket \neg\phi \rrbracket = W \setminus \llbracket \phi \rrbracket.$$

$$(VA) \llbracket \phi \wedge \psi \rrbracket = \llbracket \phi \rrbracket \cap \llbracket \psi \rrbracket.$$

$$(VU) \llbracket \phi \diamond \psi \rrbracket = \bigcup_{w \in \llbracket \phi \rrbracket} \min_{\leq_w}(\llbracket \psi \rrbracket).$$

$$(VB) \llbracket B\phi \rrbracket = \{w \in W : R_w \subseteq \llbracket \phi \rrbracket\}$$

$$(VN) \llbracket N\phi \rrbracket = \{w \in W : R_w \supseteq \llbracket \phi \rrbracket\}$$

Using (VU) and a representation theorem in [17], the following property can be shown.

**Theorem 3.1** *The  $\diamond$ -operator satisfies the KM postulates for update.*

In the sequel we shall use the following abbreviation:

$$\mathbb{O}\phi =_{df} B\phi \wedge N\phi.$$

It is obvious that

$$\llbracket \mathbb{O}\phi \rrbracket = \{w \in W : R_w = \llbracket \phi \rrbracket\}.$$

A sentence  $\phi$  is said to be *valid in a model*  $\mathcal{M} = \langle W, R, \leq, \llbracket \cdot \rrbracket \rangle$ , if  $\llbracket \phi \rrbracket = W$ . This fact is denoted  $\mathcal{M} \models \phi$ .

## 4 Interpreting Belief Revision Systems

In this section we show that the semantics of  $\mathcal{L}$ -preservative belief revision systems can be captured by update models to the extent that the question  $\chi \in Cn(\{\phi\}) \oplus \psi$  can be reduced to deciding whether the  $\mathcal{L}_{\mathcal{L}}$ -sentence  $\mathbb{O}\phi \diamond B\psi \rightarrow B\chi$  is valid in a particular  $\mathcal{L}_{\mathcal{L}}$  update model  $\mathcal{M}^{\oplus}$ . Formally, we have the following.

**Theorem 4.1** *For every  $\mathcal{L}$ -preservative belief revision system  $\langle Cn, T, \oplus \rangle$ , there is an  $\mathcal{L}_{\mathcal{L}}$  update model  $\mathcal{M}^{\oplus}$ , such that for all sentences  $\phi$  and  $\psi$  in  $S$ , we have*

$$(RU) Cn(\{\phi\}) \oplus \psi = \{\chi \in S : \mathcal{M}^{\oplus} \models \mathbb{O}\phi \diamond B\psi \rightarrow B\chi\}.$$

**Proof.** (*Construction*). Let  $\mathcal{M}^{\oplus} = \langle W, R, \leq, \llbracket \cdot \rrbracket \rangle$ , where  $W = \{w_{\lfloor \phi \rfloor} : \phi \in S\}$ . By  $\lfloor \phi \rfloor$  we mean the equivalence class  $\{\psi \in S : Cn(\{\psi\}) = Cn(\{\phi\})\}$ . In order not to complicate the notation we shall write  $w_{\lfloor \phi \rfloor}$  simply as  $w_{\phi}$ .

Then choose some (arbitrary) values  $\llbracket p \rrbracket$ , for each  $p \in P$ , and extend  $\llbracket \cdot \rrbracket$  to all  $\phi \in S$  through equations (VC) and (VA). After this, put  $w_{\phi} R w_{\psi}$  if and only if  $w_{\psi} \in \llbracket \phi \rrbracket$ . Thus  $R_{w_{\phi}} = \llbracket \phi \rrbracket$ .

For  $w_{\phi} \in W$ , define  $w_{\psi} \leq_{w_{\phi}} w_{\chi}$  if and only if either  $Cn(\{\psi\}) = Cn(\{\chi\})$ , or  $Cn(\{\phi\}) \oplus (\psi \vee \chi) = Cn(\{\psi\})$ .

From the definition it immediately follows that  $\leq_{w_\phi}$  is reflexive. The transitivity of  $\leq_{w_\phi}$  follows from the revision postulates, as do (S1) and (S2).

Then, extend  $\llbracket \cdot \rrbracket$  to all  $\mathcal{L}_{\mathcal{L}}$ -sentences through the equations (VU), (VB), and (VN).

Thus  $\mathcal{M}^\oplus$  is indeed an  $\mathcal{L}_{\mathcal{L}}$  update model.<sup>3</sup> Now the equality (RU) can be verified.

**Corollary 4.1** *Let  $Cn(\{\phi\}) \oplus \psi = Cn(\{\chi\})$ . Then  $\mathcal{M}^\oplus \models \Box\phi \diamond B\psi \leftrightarrow \Box\chi$ .*

**Corollary 4.2**  $Cn(\emptyset) = \{\phi \in S : \mathcal{M}^\oplus \models B\phi\}$ .

## 5 On the Limitations of Introspection

There is considerable evidence to the fact that revision is a meta level concept, not expressible in the object language whose theories are being revised. The most striking piece of evidence is the *triviality theorem* of Gärdenfors [7]. To arrive at his theorem Gärdenfors invokes the so called *Ramsey rule* that relates conditional propositions and theory change. The Ramsey rule is summarized by Gärdenfors as follows:

Accept a proposition of the form ‘If  $A$  then  $B$ ’ in a state of belief  $K$  if and only if the minimal change of  $K$  needed to accept  $A$  also requires accepting  $B$ .

In other words, if  $\psi > \chi$  represents the conditional ‘If  $\psi$  were true, then  $\chi$  would also be true’, then Gärdenfors’ interpretation of the Ramsey rule can be expressed as

**(RR’)**  $\psi > \chi \in Cn(\{\phi\})$  iff  $\chi \in Cn(\{\phi\}) \oplus \psi$ .

Gärdenfors has shown that this rule is incompatible with revision postulates (R1)–(R3).<sup>4</sup> The incompatibility means that if a belief revision system satisfies (RR’), and (R1)–(R3), then the system is trivial.<sup>5</sup>

The fundamental conflict seems to be between (RR’) and (R2). Much discussion has centered around the question whether to therefore abandon (RR’) or (R2) (for an overview, see [7]). However, Grahne [10] has constructed a non-trivial logic that has both a conditional connective  $>$  and an update operator  $\diamond$  in the object language, such that  $\diamond$  satisfies the *KM*-postulates for update, and  $>$  has the interpretation

$$\llbracket \psi > \chi \rrbracket = \{w \in W : \min_{\leq_w}(\psi) \subseteq \llbracket \chi \rrbracket\},$$

while the two operators are connected by the bidirectional derivation rule

**(RR)**  $\psi > \chi \in Cn(\{\phi\})$  iff  $\chi \in Cn(\{\phi \diamond \psi\})$ ,

<sup>3</sup>Note that if  $\mathcal{L}$  is finitary, then  $\mathcal{M}^\oplus$  is finite.

<sup>4</sup>As a matter of fact, Gärdenfors obtains the result already for certain weakened versions of (R1)–(R3).

<sup>5</sup>A belief revision system  $\langle T, \oplus \rangle$  is said to be *non-trivial*, if there are sentences  $\phi$ ,  $\psi$ , and  $\chi$  in  $S$ , and a theory  $\Gamma$  in  $T$ , such that the sentences  $\phi$ ,  $\psi$ , and  $\chi$  are consistent, the sentences  $\phi \wedge \psi$ ,  $\psi \wedge \chi$ , and  $\phi \wedge \chi$  are inconsistent, and the theories  $\Gamma$ ,  $\Gamma \oplus \phi$ ,  $\Gamma \oplus \psi$ , and  $\Gamma \oplus \chi$  are consistent. Otherwise, the system is *trivial* [7].

which is nothing less than the Ramsey rule. In particular, our current update language  $\mathcal{L}_C$  could be extended with a  $>$ -operator and the following semantic version of the Ramsey rule:  $\mathcal{M} \models \phi \rightarrow \psi > \chi$  if and only if  $\mathcal{M} \models \phi \diamond \psi \rightarrow \chi$ , for all update models  $\mathcal{M}$ . Since update and conditional thus belong to the same language level, and revision becomes an update on the meta level, it should not be too surprising that confusing the two levels by adopting (RR') and the revision postulate (R2) leads to triviality.

From a semantic viewpoint, the hierarchical approach of Theorem 4.1 is well justified. In particular, the equation (RU) in Theorem 4.1 could be justifiable as a postulate to require for any belief revision system. From a proof theoretic viewpoint, Levi [18] arrives at a similar meta level approach. In his approach, Levi distinguishes between the agent's beliefs (called the *corpus*), and a set of statements about the agent's beliefs (called the *metacorpus*). That is, if the agent's beliefs are represented by a theory  $\Gamma$ , then the metacorpus, called  $Poss(\Gamma)$ , contains  $B\phi$  iff  $\phi \in \Gamma$ , and contains  $\neg B\phi$  iff  $\phi \notin \Gamma$ . Furthermore,  $\Gamma$  is to be a subset of  $Poss(\Gamma)$ . In particular, Levi points out that the relationship between  $Poss(\Gamma)$  and  $Poss(\Gamma \oplus \psi)$  is *not* that of a revision (but he does not say exactly what that relationship is), and that  $Poss(\Gamma)$  is formulated in a language that properly extends the language in which  $\Gamma$  is formulated.

Our results in the previous section verify that observation: If we assume that the agent's corpus is represented by a sentence  $\phi$ , then Levi's metacorpus would correspond to the sentence  $0\phi$  interpreted in the model  $\mathcal{M}^\oplus$ . As we have seen, the relationship between the 'metacorpus' for  $Cn(\{\phi\})$  and for  $Cn(\{\phi\}) \oplus \psi$  is that of an *update*, and not a revision. We could thus, in analogy with Levi, relate revision and conditional through the following reasoning:

$$\chi \in Cn(\{\phi\}) \oplus \psi \text{ iff } \mathcal{M}^\oplus \models 0\phi \diamond B\psi \rightarrow B\chi \text{ iff } \mathcal{M}^\oplus \models 0\phi \rightarrow B\psi > B\chi.$$

We could thus also account for *iterated* conditionals, that is propositions of the form 'If A then (if B then C)'. Let  $\varphi$  be a sentence in  $S$  such that  $Cn(\{\phi\}) \oplus \psi = Cn(\{\varphi\})$ . Then we would simply have  $\mathcal{M}^\oplus \models 0\phi \rightarrow B\psi > (B\chi > B\xi)$  iff  $\mathcal{M}^\oplus \models 0\phi \diamond B\psi \rightarrow B\chi > B\xi$  iff  $\mathcal{M}^\oplus \models (0\phi \diamond B\psi) \diamond B\chi \rightarrow B\xi$  iff  $\mathcal{M}^\oplus \models (0\varphi) \diamond B\chi \rightarrow B\xi$  iff  $\xi \in Cn(\{\varphi\}) \oplus \chi$  iff  $\xi \in (Cn(\{\phi\}) \oplus \psi) \oplus \chi$ .

Returning to the triviality issue, we saw that adopting (RR') and (R2) gives the agent some abilities to reason about its revision strategies. We shall now consider the scenario where the agent has such abilities as a consequence of *introspection*: Suppose that instead of propositional logic, the agent is reasoning in some *epistemic* logic, for example in the logic *K45* of belief. This modal logic is propositional logic augmented with a modal operator  $\Box$ . The axioms and derivation rules of *K45* are those of modal logic *K* plus the axiom of *positive introspection*

$$\Box\phi \rightarrow \Box\Box\phi,$$

and the axiom of *negative introspection*

$$\neg\Box\phi \rightarrow \Box\neg\Box\phi.$$

If the operator  $\Box$  is taken to mean belief, then the positive introspection axiom 'says' that if the agent believes  $\phi$  then it believes that it believes  $\phi$ . The axiom for negative introspection says that if the agent does not believe  $\phi$  then it believes that it does not believe  $\phi$ . In case the agent's logic is both positively and negatively introspective, the agent is said to be *fully introspective*. For a more detailed account of *K45* see e.g. [12, 14].

By invoking a representation theorem of Grove [11], it is possible to show that Theorem 4.1 also holds for *K45*. Note however that the operators  $B$  and  $\Box$  will be given different meanings by this construction. If the agent believes the sentence  $\Box\phi$ , this will be translated into the meta level sentence  $B\Box\phi$ , meaning intuitively that ‘the agent believes  $\Box\phi$ ’, which is not the same as ‘the agent believes  $B\phi$ ’, or ‘the agent believes “the agent believes  $\phi$ ”’. Suppose we try to eliminate this difference by identifying the object level  $\Box$  with the meta level  $B$ , allowing the agent to hold beliefs about its own beliefs. We show below that this attempt is doomed to end in triviality, as also happens when (RR’) and (R2) are adopted. A result with a similar flavour but somewhat different assumptions has been obtained by Fuhrmann [6]. To state our result formally, we shall need the following definition: A  $\mathcal{L}$ -belief revision system is said to be *strongly trivial* if all sentences are complete.<sup>6</sup>

We can now show that if the difference between  $B$  and  $\Box$  is erased, the construction of Theorem 4.1 is only applicable to strongly trivial belief revision systems.

**Theorem 5.1** *Let  $\langle Cn, T, \oplus \rangle$  be a  $\mathcal{L}$ -belief revision system satisfying at least (R0)–(R3), where  $\mathcal{L}$  and  $Cn$  extend propositional logic, the consequence operator is finitistic<sup>7</sup>, and where  $B\phi \rightarrow BB\phi \in Cn(\emptyset)$ , and  $\neg B\phi \rightarrow B\neg B\phi \in Cn(\emptyset)$ . Suppose furthermore that the equation (RU) of Theorem 4.1 holds. Then the belief revision system is strongly trivial.*

**Proof.** Suppose to the contrary that the belief revision system is not strongly trivial, and let  $\phi$  and  $\psi$  be sentences in  $\mathcal{L}$  such that  $\psi \notin Cn(\{\phi\})$ , and  $\neg\psi \notin Cn(\{\phi\})$ . Thus, by (R2),  $\psi \notin Cn(\{\phi\}) \oplus \phi$ . Then we have by (RU)

$$\mathcal{M}^\oplus \not\models \Box\phi \diamond B\phi \rightarrow B\psi,$$

and from the definition of  $\diamond$  we get

$$\mathcal{M}^\oplus \not\models \Box\phi \rightarrow B\psi.$$

It now follows from the definition of the  $\Box$ -operator that

$$\mathcal{M}^\oplus \models \Box\phi \rightarrow \neg B\psi,$$

and by negative introspection we get

$$\mathcal{M}^\oplus \models \Box\phi \rightarrow B\neg B\psi,$$

Applying (RU) again we get

$$\neg B\psi \in Cn(\{\phi\}).$$

From (R2) it now follows that

$$\neg B\psi \in Cn(\{\phi\}) \oplus \psi.$$

On the other hand, from the definition of  $\diamond$  it follows that

$$\mathcal{M}^\oplus \models \Box\phi \diamond B\psi \rightarrow B\psi,$$

<sup>6</sup>A sentence  $\phi$  is complete if for any other sentence  $\psi$ , either  $\psi \in Cn(\{\phi\})$ , or  $\neg\psi \in Cn(\{\phi\})$ .

<sup>7</sup>A consequence operator  $Cn$  is said to be finitistic, if for all sets of sentences  $X$  we have  $Cn(X) = \bigcup\{Cn(Y) : Y \subseteq X \text{ and } Y \text{ is finite}\}$

and applying positive introspection we get

$$\mathcal{M}^\oplus \models \mathbf{0}\phi \diamond \mathbf{B}\psi \rightarrow \mathbf{B}\mathbf{B}\psi.$$

Then (RU) gives

$$\mathbf{B}\psi \in \mathit{Cn}(\{\phi\}) \oplus \psi.$$

This means that  $\mathit{Cn}(\{\phi\}) \oplus \psi$  is inconsistent. From (R3) it then follows that  $\psi$  must be inconsistent. Therefore  $\neg\psi \in \mathit{Cn}(\emptyset)$ , and consequently  $\neg\psi \in \mathit{Cn}(\phi)$ , a contradiction to the assumption that the belief revision system is not strongly trivial.

We see that the agent cannot be fully introspective with respect to the operator  $\mathbf{B}$ . Whatever introspection means, it cannot thus mean introspection with respect to the way the agent revises its beliefs.

## 6 Conclusion

We have given a metalanguage in which it is possible to reason formally about how an agent changes its beliefs. The agent is performing revisions, that is, it is adjusting its beliefs to accommodate new, possibly conflicting information about a static world. At the meta level, the world—which includes the agent’s state of beliefs—is changing as the agent changes its mind. These changes are modeled by updates. Levesque’s ‘only know’ modal operator is used at the meta level to describe the totality of beliefs held by an agent.

The negative result in the paper is that this type of formal account of revision cannot be extended to the situation where the agent is able to reason about its beliefs. A fully introspective agent cannot use our construction to reason about the results of its own revisions, on pain of triviality.

## Acknowledgement

We would like to thank David Israel for pointing out several relevant references.

## References

- [1] C. E. Alchourrón, P. Gärdenfors & D. Makinson. On the logic of theory change: partial meet contraction and revision functions. *Journal of Symbolic Logic*, (50) 510–530, 1985.
- [2] M. Dalal. Investigations into a theory of knowledge base revision: Preliminary Report. In *Proceedings of the Seventh National Conference on Artificial Intelligence*, pages 475–479, 1988.
- [3] J. Doyle. Rational belief revision (preliminary report). In *Proceedings of the Second International Conference on Principles of Knowledge Representation and Reasoning*, pages 163–174, 1991.

- [4] R. Fagin, J. Y. Halpern, Y. Moses, & M. Y. Vardi. *Reasoning about Knowledge*. To appear in 1992.
- [5] R. Fagin, J. D. Ullman & M. Y. Vardi. On the semantics of updates in databases. In *Proceedings of the Second ACM SIGACT-SIGMOD-SIGART Symposium on Principles of Database Systems*, pages 352–365, 1983.
- [6] A. Fuhrmann. Reflective modalities and theory change. *Synthèse*, (81) 115–1344, 1989.
- [7] P. Gärdenfors. *Knowledge in Flux: Modeling the Dynamics of Epistemic States*. Bradford Books, MIT Press, Cambridge, MA, 1988.
- [8] A. Gibbard. Two recent theories of conditionals. In: *Ifs: Conditionals, Belief, Decision, Chance and Time*, W. L. Harper, R. Stalnaker & G Pearce (eds.). D. Reidel Publ. Co., Dordrecht, 1981, pp. 211–248.
- [9] M. L. Ginsberg. Counterfactuals. *Artificial Intelligence*, (30) 35–79, 1986.
- [10] G. Grahne. Updates and Counterfactuals. In: *Proceedings of the Second International Conference on Principles of Knowledge Representation and Reasoning*, pages 269–276, 1991.
- [11] A. Grove. Two modelings for theory change. *Journal of Philosophical Logic*, (17) 157–70, 1988.
- [12] J. Y. Halpern & Y. Moses. A guide to the modal logics of knowledge and belief. In: *Proceedings of the Ninth International Joint Conference on Artificial Intelligence*, pages 480–490, 1985.
- [13] J. Y. Halpern & M. Y. Vardi. Model checking vs. theorem proving: a manifesto. In: *Proceedings of the Second International Conference on Principles of Knowledge Representation and Reasoning*, pages 325–334, 1991.
- [14] G. E. Hughes & M. J. Cresswell. *A Companion to Modal Logic*. Methuen, London, 1984.
- [15] W.L. Harper & C.H. Hooker (eds). *Foundations of Probability Theory, Statistical Inference, and Statistical Theories of Science*. D. Reidel Publ. Co., Dordrecht, 1976.
- [16] H. Katsuno & A. O. Mendelzon. Propositional knowledgebase revision and minimal change. *Artificial Intelligence*, (52) 263–294, 1991.  
A unified view of propositional knowledge base updates. In: *Proceedings of the 11th International Joint Conference on Artificial Intelligence*, pages 1413–1419, 1989.
- [17] H. Katsuno & A. O. Mendelzon. On the difference between updating a knowledge base and revising it. In: *Proceedings of the Second International Conference on Principles of Knowledge Representation and Reasoning*, pages 387–394, 1991.
- [18] I. Levi. Iteration of Conditionals and the Ramsey test. *Synthèse*, (76) 49–81, 1988.

- [19] H. J. Levesque. All I Know: A study in autoepistemic logic. *Artificial Intelligence*, (42) 263–309, 1990.
- [20] D. K. Lewis. Counterfactuals and comparative possibility. *Journal of Philosophical Logic* (7) 418–446, 1973.
- [21] B. Nebel. A knowledge level analysis of belief revision. In: *Proceedings of the First International Conference on Principles of Knowledge Representation and Reasoning*, pages 301–311, 1989.
- [22] A. S. Rao & N. Y. Foo. Minimal change and maximal coherence: A basis for belief revision and reasoning about action. In *Proceedings of the Eleventh International Joint Conference on Artificial Intelligence*, pages 966–971, 1989.
- [23] M. Winslett. Reasoning about action using a possible models approach. In: *Proceedings of the Seventh National Conference on Artificial Intelligence*, pages 89–93, 1988.
- [24] M. Winslett. *Updating Logical Databases*. Cambridge University Press, 1990.