

# Keep Changing Your Beliefs, Aiming for the Truth

Alexandru Baltag · Sonja Smets

Received: 6 May 2010 / Accepted: 7 January 2011 / Published online: 21 September 2011  
© The Author(s) 2011. This article is published with open access at Springerlink.com

**Abstract** We investigate the process of *truth-seeking by iterated belief revision with higher-level doxastic information*. We elaborate further on the main results in Baltag and Smets (Proceedings of TARK, 2009a, Proceedings of WOLLIC’09 LNAI 5514, 2009b), applying them to the issue of *convergence to truth*. We study the conditions under which the belief revision induced by a series of truthful iterated upgrades eventually stabilizes on true beliefs. We give two different conditions ensuring that beliefs converge to “full” (complete) truth, as well as a condition ensuring only that they converge to true (but not necessarily complete) beliefs.

## 1 Introduction

What happens in the long term with an agent’s “epistemic state” (comprising her beliefs, her knowledge and her conditional beliefs), when receiving a *sequence of truthful but uncertain* pieces of information, formulated in a way that *may refer (directly or indirectly) to that agent’s own epistemic state*? Do the agent’s beliefs (or knowledge, or conditional beliefs) reach a *fixed point*? Or do they exhibit instead a *cyclic behavior*, oscillating forever? Finally, in case that beliefs do stabilize on a fixed point, what conditions would ensure that they stabilize on the *truth*?

---

A. Baltag (✉)

Institute for Logic, Language and Information, University of Amsterdam, Amsterdam, The Netherlands

e-mail: thealexandrumbaltag@gmail.com

S. Smets

Rijksuniversiteit Groningen, Groningen, The Netherlands

e-mail: S.J.L.Smets@rug.nl

S. Smets

IEG, Oxford University, Oxford, UK

In this paper we present a setting to investigate and provide some answers to the above questions. Our work lies at the interface between Belief Revision, Dynamic Epistemic Logic, Formal Learning Theory and the logical research area known as Truth Approximation. As such, it is related to the work by Kelly (1998a, b), and Hendricks et al. (1997), on the learning power of iterated belief revision. These authors adopt the paradigm of *function learning*: the problem of *learning (the future of) an infinite stream of incoming data*, based on its past (i.e. on the data already observed). This is a problem of *prediction*, which connects to one of our results (Corollary 6). But, in Learning Theory terms, our focus is rather on the so-called *set learning* (or “*language learning*”): the problem of *learning the set of data that are true* (in the “real world”, or in a given language), *based on the finite sequence of data that were already observed*. In this setting, it is assumed that *the data are observed in more or less random manner*, so that *predicting the entire future sequence in all its details is not feasible, or even relevant*. Another difference from the above authors, as well as from the AGM-based setting for Belief Revision (Alchourrón et al. 1985), is that we also consider learning “data” that embody higher-level doxastic information. In this sense, our work lies within the paradigm of Dynamic Epistemic Logic, and can be seen as an extension and continuation of the work of Dégremont and Gierasimczuk (2009).

We start by giving an overview of the most important results in Baltag and Smets (2009a, b) and investigate further *under what conditions does iterated belief revision reach a fixed point, and when do beliefs stabilize on the “truth”*. In particular the new results in this paper show what type of belief revision method and what infinite streams of information ensure that *the agent’s beliefs will converge to the “full truth” in finitely many steps*. By full truth we mean that, from some point onwards, all beliefs will come to be true and all true sentences (from a given language) will be believed. So we give sufficient conditions ensuring that the agent’s beliefs are *guaranteed to eventually fully match the truth (from some point onwards), but without the agent necessarily ever knowing for sure whether or not she has already reached the truth*. We also indicate weaker conditions, in which an agent’s beliefs are *guaranteed to converge to true (but not necessarily complete) beliefs*.

Among our conditions we require the *finiteness of the initial set of possibilities*. So, in this paper, we restrict ourselves to the case of finite plausibility frames. This is a simplifying condition, which is neither necessary nor sufficient for convergence. We adopt it here mainly for simplicity, to avoid the problems connected to induction, but also because this condition induces a *natural “expectation of convergence”*: indeed, for *finite* models it seems almost “obvious” that any process of learning (only) true information should eventually come to an end. But in fact, our counterexamples show that this natural expectation is totally unjustified! Hence, *further (and rather intricate) conditions are needed for convergence*.

These results relate to the concept of “*identifiability in the limit*” from Formal Learning Theory, which captures the above property of eventually reaching a *stable (complete) true belief, without necessarily ever knowing when this happens*. In contrast, the concept of “*finite identifiability*” from Learning Theory corresponds to reaching *knowledge* of full truth. Dégremont and Gierasimczuk (2009) first

introduced this distinction in the context of belief revision, and investigated the connection between repeated updates and finite identifiability. Of course, the interested case in Formal Learning Theory is the *infinite* one, corresponding to e.g. learning a natural language and connected to the problem of *induction* in experimental sciences. In current, unpublished work (Baltag et al. 2010), we fully engage with Learning Theory, by investigating the *infinite-case analogues* of the work presented in this paper.

While the classical literature on iterated belief revision (Booth and Meyer 2006; Darwiche and Pearl 1997) deals only with *propositional* information, we are interested in the case where the new information may itself refer to *the agent's own beliefs, knowledge or belief-revision plans*. Hence, our framework allows us to reason explicitly about (revision with) *higher-order beliefs* (beliefs about beliefs etc.). In particular, the issue of *convergence* of beliefs during *iterated belief revision with the same sentence* (starting on a finite model) is highly *non-trivial* in our setting, while one can easily see that it becomes *trivial* when we restrict the incoming information to purely ontic, non-doxastic sentences.

The toolbox we use in our investigation includes the belief-revision-friendly version of Dynamic Epistemic Logic, as developed in Aucher (2003), Baltag and Smets (2006a, b, 2008a, b), van Benthem (2007), van Ditmarsch (2005). We work with a specific type of Kripke structures, namely finite “plausibility” structures, to represent the beliefs, knowledge and epistemic state of an agent. These structures consist of a finite set of possible worlds equipped with a total preorder to represent the worlds’ *relative plausibility*. In Belief Revision Theory, plausibility structures are known as special cases of Halpern’s “preferential models” (Halpern 2003), Grove’s “sphere” models (Grove 1988), Spohn’s ordinal-ranked models (Spohn 1988) and Board’s “belief-revision structures” (Board 2002). As models we use *pointed plausibility models*, which are plausibility structures together with a designated state, representing *the real world*, and a valuation map capturing the ontic “facts” in each possible world.

One of the important ingredients in our setting are the so-called *belief upgrades*. A belief upgrade is a type of *model transformer*, i.e. a function taking plausibility models as input and returning “upgraded” models as output. Examples of upgrades have been considered in the literature on Belief Revision e.g. by Boutilier (1996) and Rott (1989), in the literature on dynamic semantics for natural language by Veltman (1996), and in Dynamic Epistemic Logic by van Benthem (2007). In this paper we consider *three types of upgrades* (“update”, “radical” upgrade and “conservative” upgrade), corresponding to *three different levels of trust in the reliability of the new information*. Update (known as “conditioning” in the Belief Revision literature) corresponds to having an absolute guarantee that the newly received information is truthful; radical upgrade (known also as “lexicographic revision”) corresponds to having a high degree of trust (a “strong belief”) that the new information is truthful; while conservative upgrade (known also as “minimal revision”) corresponds to having a simple (minimal, or “bare”) belief in the reliability of the source. We show that *each of these belief-revision methods has a different behavior with respect to convergence and truth-approximation*.

## 2 Beliefs and Upgrades in Plausibility Models

In this section, we review some notions and examples from Baltag and Smets (2009). The basic semantic concept is that of a (*finite, single-agent*) “*plausibility*” frame. This is a structure  $\mathbf{S} = (S, \leq)$ , consisting of a finite set  $S$  of “states” (or possible worlds) and a *total preorder*  $\leq \subseteq S \times S$ .<sup>1</sup> We use the notation  $s < t$  for the corresponding *strict order* (i.e.  $s < t$  iff  $s \leq t$  but  $t \not\leq s$ ), and  $s \cong t$  for the corresponding *equi-plausibility* relation ( $s \cong t$  iff both  $s \leq t$  and  $t \leq s$ ).

**(Pointed) Plausibility Models** A (*finite, single-agent, pointed*) *plausibility model* is a structure  $\mathbf{S} = (S, \leq, \|\cdot\|, s_0)$ , consisting of a plausibility frame  $(S, \leq)$  together with a *valuation map*  $\|\cdot\| : \Phi \rightarrow \mathcal{P}(S)$  and a designated state  $s_0 \in S$ , called the “actual state”. The valuation maps every element  $p$  of some given set  $\Phi$  of “atomic sentences” into a set of states  $\|p\| \subseteq S$ .

**Knowledge and (Conditional) Belief** Given a plausibility model  $\mathbf{S}$  and sets  $P, Q \subseteq S$  of states, we define:  $bestP = Min_{\leq} P := \{s \in P : s \leq s' \text{ for all } s' \in P\}$ ,  $best := bestS$ ,  $KP := \{s \in S : P = S\}$ ,  $BP := \{s \in S : best \subseteq P\}$ ,  $B^Q P := \{s \in S : bestQ \subseteq P\}$ .

Note that the operators  $K, B, B^Q$  are in fact *Kripke modalities* for the following accessibility relations:  $R_K = S \times S$ ,  $R_B = S \times best$ ,  $R^{B^Q} = S \times bestQ$ . This shows that the plausibility-model semantics for knowledge and (conditional) belief is a special case of the Kripke semantics. Moreover, note that, for all  $P, Q \subseteq S$ , we have that  $KP, BP, B^Q P \in \{S, \emptyset\}$ . In other words, *the truth value of these propositions does not depend on the choice of the actual world  $s_0$* . This is natural: “knowledge” (in the fully introspective sense given by the  $K$  operator), belief and conditional belief are attitudes that are “*internal*” to the subject, and so the truth value of such statements should only depend on the agent’s doxastic structure, *not* on the real world.

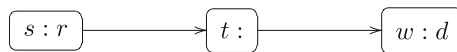
**Interpretation** The elements of  $S$  represent the *possible states*, or “possible worlds”. The “actual state”  $s_0$  (as given by the pointed plausibility model) represents the *correct* description of the real world and might well be unknown to the agent. The atomic sentences  $p \in \Phi$  represent “*ontic*” (*non-doxastic*) facts, that might hold or not in a given state. The valuation tells us which facts hold at which worlds.  $KP$  is read as “the agent knows  $P$ ”. The plausibility relation  $\leq$  is the agent’s *plausibility order* between her “epistemically possible” states: we read  $s \leq t$  as “the agent considers  $s$  at least as plausible as  $t$ ”. This is meant to capture the agent’s (*conditional*) beliefs about the state of the system. Here  $B^Q P$  is read as “the agent believes  $P$  conditional on  $Q$ ” and means that, if the agent would receive some further (certain) information  $Q$  (to be added to what she already knows) then she would believe that  $P$  was the case. The above definition says that  $B^Q P$  holds iff  $P$  holds in all the “best” (i.e. the most plausible)  $Q$ -states. In particular, a *simple* (*non-conditional*) belief  $BP$  holds iff  $P$  holds in all the best states that are epistemically possible.

<sup>1</sup> A “total preorder”  $\leq$  on a state space  $S$  is a reflexive and transitive binary relation on  $S$  such that any two states are comparable: for all  $s, t \in S$ , we have either  $s \leq t$  or  $t \leq s$ .

**Doxastic Propositions** We introduce a notion of proposition that can be interpreted in *any* model: A *doxastic proposition* is a map  $\mathbf{P}$  assigning to each model  $\mathbf{S}$  some set  $\mathbf{P}_S \subseteq S$  of states in  $S$ . We write  $s \models_S \mathbf{P}$ , and say that the proposition  $\mathbf{P}$  is *true at state*  $s \in S$  *in the model*  $\mathbf{S}$ , iff  $s \in \mathbf{P}_S$ . We say that a doxastic proposition  $\mathbf{P}$  is *true at (the pointed model)*  $\mathbf{S}$ , and write  $\mathbf{S} \models \mathbf{P}$  if it is true at the “actual state”  $s_0$  in the model  $\mathbf{S}$ , i.e. if  $s_0 \models_S \mathbf{P}$ . We have the “always true”  $\top$  and “always false”  $\perp$  propositions  $\perp_S := \emptyset, \top_S := S$ . And for each atomic sentence  $p$ , the corresponding doxastic proposition  $\mathbf{p}$  is given by  $\mathbf{p}_S = \|p\|_S$ . All the operations on sets can be “lifted” *pointwise* to propositions: negation  $(\neg\mathbf{P})_S := S \setminus \mathbf{P}_S$ , conjunction  $(\mathbf{P} \wedge \mathbf{R})_S := \mathbf{P}_S \cap \mathbf{R}_S$  etc, the “best” operator  $(\text{best } \mathbf{P})_S := \text{best } \mathbf{P}_S$ , all the *modalities*  $(\mathbf{K}\mathbf{P})_S := \mathbf{K}\mathbf{P}_S, (\mathbf{B}\mathbf{P})_S := \mathbf{B}\mathbf{P}_S, (\mathbf{B}^Q\mathbf{P})_S := \mathbf{B}^Q\mathbf{P}_S$  etc.

**Questions and Answers** A (binary) *question*  $\mathcal{Q} = \{\mathbf{P}, \neg\mathbf{P}\}$  is a pair consisting of any proposition  $\mathbf{P}$  and its negation. A binary question  $\mathcal{Q}$  induces a *partition*  $\{\mathbf{P}_S, \neg\mathbf{P}_S\}$  on any model  $\mathbf{S}$ . Any of the propositions  $\mathbf{A} \in \{\mathbf{P}, \neg\mathbf{P}\}$  is a possible *answer* to question  $\mathcal{Q}$ . The *correct answer* to  $\mathcal{Q}$  in a pointed model  $\mathbf{S}$  is the unique answer  $\mathbf{A} \in \mathcal{Q}$  such that  $\mathbf{S} \models \mathbf{A}$ .

**Example 1** Consider a pollster (Charles), holding the following beliefs about how a voter (Mary) will vote:



In this single-agent model (with Charles as the only real ‘agent’, since Mary’s voting is considered here as part of ‘Nature’), the arrows represent the converse plausibility relations (but note that for simplicity of our drawing we skipped the loops and the arrows that can be obtained by transitivity). There are only three possible worlds  $s$  (in which Mary votes Republican),  $w$  (in which she votes Democrat) and  $t$  (in which she doesn’t vote). The atomic sentences are  $r$  (for “voting Republican”) and  $d$  (for “voting Democrat”). The valuation is given:  $\|r\| = \{s\}, \|d\| = \{w\}$ . We assume the real world is  $s$ , so in reality Mary will vote Republican ( $r$ ). This is unknown to Charles who believes that she will vote Democrat ( $d$ ) - because  $d$  is true in the most plausible world  $w$  - ; and in case this turns out wrong, he’d rather believe that she won’t vote ( $\neg d \wedge \neg r$ ) than thinking that she votes Republican. If we ask Charles the question “will Mary vote Democrat?”, this is represented by  $\mathcal{Q} = \{d, \neg d\}$ . The correct answer would be  $\neg d$ .

**Learning New Information: Doxastic Upgrades** We move on now to *dynamics*: what happens when some proposition  $\mathbf{P}$  is *learnt*? According to Dynamic Epistemic Logic, this induces a *change of model*: a “model transformer”. However, the specific change depends on *the agent’s attitude to the plausibility* of the announcement: *how certain is the new information*? Three main possibilities have been discussed in the literature: (1) the announcement  $\mathbf{P}$  is *certainly* true: the agent learns with absolute certainty the correct answer  $\mathbf{P}$  to a question  $\mathcal{Q} = \{\mathbf{P}, \neg\mathbf{P}\}$ ; (2) the announcement is *strongly believed* to be true but the agent is not absolutely certain. The agent learns that an answer  $\mathbf{P}$  is more plausible than its negation; (3) the announcement is (*simply*) *believed*: it is *believed (in a simple, “weak” sense)* that

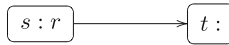
the speaker tells the truth. In this case the agent considers a different question, namely  $\{best\mathbf{P}, \neg best\mathbf{P}\}$  and accepts the first answer to be more plausible. These three alternatives correspond to *three forms of “learning”*  $\mathbf{P}$ , forms discussed in van Benthem (2007, 2009) in a Dynamic Epistemic Logic context: “update”<sup>2</sup>  $!\mathbf{P}$ , “radical upgrade”  $\uparrow \mathbf{P}$  and “conservative upgrade”  $\uparrow \mathbf{P}$ .

We will use “upgrades” as a general term for all these three model transformers, and denote them in general by  $\dagger \mathbf{P}$ , where  $\dagger \in \{!, \uparrow, \uparrow\}$ . Formally, each of our upgrades is a (possibly partial) function taking as inputs pointed models  $\mathbf{S} = (S, \leq, \|\|, s_0)$  and returning new (“upgraded”) pointed models  $\dagger \mathbf{P}(\mathbf{S}) = (S', \leq', \|\|', s'_0)$ , with  $S' \subseteq S$ . Since upgrades are purely doxastic, they *won’t affect the real world or the “ontic facts” of each world*: i.e. they all satisfy  $s'_0 = s_0$  and  $\|p\|' = \|p\| \cap S'$ , for atomic  $p$ . So, in order to completely describe a given upgrade, we only have to specify (a) *its possible inputs*  $\mathbf{S}$ , (b) *the new set of states*  $S'$ ; (c) *the new relations*  $\leq'$ .

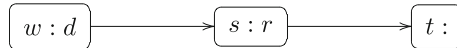
- (1) **Learning Certain information: “Update”**. The update  $!\mathbf{P}$  is an operation on pointed models which is *executable* (on a pointed model  $\mathbf{S}$ ) *iff*  $\mathbf{P}$  is true (at  $\mathbf{S}$ ) and which *deletes all the non- $\mathbf{P}$ -worlds from the pointed model, leaving everything else the same*. Formally, an update  $!\mathbf{P}$  is an upgrade such that: (a) it takes as inputs only pointed models  $\mathbf{S}$ , such that  $\mathbf{S} \models \mathbf{P}$ ; (b) the new set of states is  $S' = \mathbf{P}_S$ ; (c)  $s \leq' t$  iff  $s \leq t$  and  $s, t \in S'$ .
- (2) **Learning from a Strongly Trusted Source: “Radical” Upgrade**. The “radical” (or “lexicographic”) upgrade  $\uparrow \mathbf{P}$ , as an operation on pointed models, can be described as “promoting” all the  $\mathbf{P}$ -worlds so that they become “better” (more plausible) than all  $\neg \mathbf{P}$ -worlds, while keeping everything else the same: the valuation, the actual world and the relative ordering between worlds within either of the two zones ( $\mathbf{P}$  and  $\neg \mathbf{P}$ ) stay the same. Formally, a radical upgrade  $\uparrow \mathbf{P}$  is (a) a *total* upgrade (taking as input *any* model  $\mathbf{S}$ ), such that (b)  $S' = S$ , and (c)  $s \leq' t$  holds iff we have: *either*  $t \notin \mathbf{P}_S$  and  $s \in \mathbf{P}_S$ , or  $s, t \in \mathbf{P}_S$  and  $s \leq t$ , or *else*  $s, t \notin \mathbf{P}_S$  and  $s \leq t$ .
- (3) **“Barely believing” what you hear: “Conservative” Upgrade**. The so-called “conservative upgrade”  $\uparrow \mathbf{P}$  [called “minimal conditional revision” in Boutilier (1996)] performs in a sense the minimal modification of a model that is forced by believing the new information  $\mathbf{P}$ . As an operation on pointed models, it can be described as “promoting” only the “best” (most plausible)  $\mathbf{P}$ -worlds, so that they become the most plausible, while keeping everything else the same. Formally,  $\uparrow \mathbf{P}$  is (a) a *total* upgrade, such that (b)  $S' = S$ , and (c)  $s \leq' t$  holds iff: *either*  $s \in best\mathbf{P}_S$ , or *else*  $s \notin best\mathbf{P}_S$  but  $s \leq t$ .

**Examples:** Consider first the case in which Charles learns from an *absolutely infallible* authority (a truth-telling “Oracle”) that *Mary will definitely not vote Democrat*. This is an *update*  $!(\neg \mathbf{d})$ , resulting in the updated model:

<sup>2</sup> Note that in Belief Revision, the term “belief update” is used for a totally different operation [the Katzuno-Mendelzon update (Katsuno and Mendelzon 1992)], while what we call “update” is known as “conditioning”. We choose to follow here the terminology used in Dynamic Epistemic Logic.

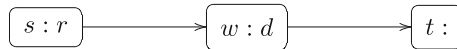


Next, consider the same scenario, except that instead of an infallible oracle we now have only a *highly trusted but fallible source*. This is a *radical upgrade*  $\uparrow(\neg\mathbf{d})$  of the original model from Example 1, resulting in:



Now, Charles believes that Mary will not vote; and if he later would learn this was not the case then he would still keep his newly acquired belief that she doesn't vote Democrat (so in this eventuality he'd conclude, correctly, that she votes Republican).

Finally, consider the case in which Charles *“barely trusts”* the source; e.g. he just hears a *rumor* that Mary will not vote Democrat. He believes the rumor, but *not strongly*: in case he later is forced to revise his beliefs, he would immediately give up the belief in the veracity of the rumor. We can interpret the learning event as a *conservative upgrade*  $\uparrow(\neg\mathbf{d})$  of the original model from Example 1, resulting in:



As in the previous case, Charles believes now that Mary will not vote. But if he would later learn this was not the case, he would immediately revert back to his older belief that Mary votes Democrat.

**Truthful Upgrades and Truthful Answers** An upgrade  $\uparrow\mathbf{P}$  is *truthful* in a pointed model  $\mathbf{S}$  if  $\mathbf{P}$  is true at  $\mathbf{S}$ . For any question  $Q = \{\mathbf{P}, \neg\mathbf{P}\}$  and any upgrade type  $\uparrow \in \{!, \uparrow, \uparrow\}$ , we can consider *the action of truthfully answering the question via upgrades of type  $\uparrow$* . This action, denoted by  $\uparrow Q$ , is a model transformer that behaves like the truthful upgrade  $\uparrow\mathbf{A}$ , whenever it is applied to a model  $\mathbf{S}$  in which  $\mathbf{A} \in Q$  is the correct answer to  $Q$ .

### 3 Iterated Belief Upgrades

As we focus in this paper on *long-term learning processes via iterated belief revision*, we need to look at *sequences (or so-called “streams”) of upgrades*. We are in particular interested in the case where an agent repeatedly faces true information, and so in iterating *truthful upgrades*. In the remainder of this section we give an overview of our recent results from Baltag and Smets (2009) concerning the issue of *belief convergence during iterated revision with true information*.

**Redundancy, Informativity, Fixed Points** An upgrade  $\uparrow\mathbf{P}$  is *redundant on a pointed model  $\mathbf{S}$*  if the two pointed models  $\uparrow\mathbf{P}(\mathbf{S})$  and  $\mathbf{S}$  are bisimilar (i.e.  $\uparrow\mathbf{P}(\mathbf{S}) \simeq \mathbf{S}$ ) in the usual sense of modal logic (Blackburn et al. 2001). Essentially, this means that  $\uparrow\mathbf{P}$  doesn't change anything (when applied to model  $\mathbf{S}$ ): all simple beliefs, conditional beliefs and knowledge *stay the same* after the upgrade. A question  $Q$  is *redundant for a type  $\uparrow \in \{!, \uparrow, \uparrow\}$  on a pointed model  $\mathbf{S}$*  if the action

$\dagger Q$  of truthfully answering the question (via upgrade type  $\dagger$ ) is a redundant upgrade. An upgrade, or question, is *informative (on S)* if it is not redundant. A model  $S$  is a *fixed point* of  $\dagger P$  if  $S \simeq \dagger P(S)$ , i.e. if  $\dagger P$  is redundant on  $S$ .

**Upgrade Streams** An *upgrade stream*  $\vec{\dagger P} = (\dagger P_n)_{n \in \mathbb{N}}$  is an infinite sequence of upgrades  $\dagger P_n$  of the same type  $\dagger \in \{!, \uparrow, \uparrow\}$ . An upgrade stream is *definable in a logic L* if all  $P_n$  are of the form  $P_n = \|\varphi_n\|$  for some  $\varphi_n \in L$ .

Any upgrade stream  $\vec{\dagger P}$  induces a function mapping every pointed model  $S$  into an infinite sequence  $\vec{\dagger P}(S) = (S_n)_{n \in \mathbb{N}}$  of pointed models, defined inductively by:

$$S_0 = S, \text{ and } S_{n+1} = \dagger P_n(S_n).$$

The upgrade stream  $\vec{\dagger P}$  is *truthful* if every  $\dagger P_n$  is truthful with respect to  $S_n$  (i.e.  $S_n \models P_n$ ). A *repeated truthful upgrade* is a truthful upgrade stream of the form  $(\dagger P_n)_{n \in \mathbb{N}}$ , where  $P_n \in \{P, \neg P\}$  for some proposition  $P$ . In other words, it consists in repeatedly learning the answer to the “same” question  $P$ ?

We say that a stream  $\vec{\dagger P}$  *stabilizes a (pointed) model S* if there exists some  $n \in \mathbb{N}$  such that  $S^n \simeq S^m$  for all  $m > n$ . Obviously, a *repeated upgrade* stabilizes  $S$  if it reaches a fixed point of  $\dagger P$  or of  $\dagger(\neg P)$ .

We say that  $\vec{\dagger P}$  *stabilizes all (simple, non-conditional) beliefs on the model S* if the process of belief-changing induced by  $\vec{\dagger P}$  on  $S$  reaches a fixed point; i.e. if there exists some  $n \in \mathbb{N}$  such that  $S_n \models BP$  iff  $S_m \models BP$ , for all  $m > n$  and all doxastic propositions  $P$ . Equivalently, iff there exists some  $n \in \mathbb{N}$  such that  $best_{S_n} = best_{S_m}$  for all  $m > n$ .

Similarly, we say that  $\vec{\dagger P}$  *stabilizes all conditional beliefs on the model S* if the process of conditional-belief-changing induced by  $\vec{\dagger P}$  on  $S$  reaches a fixed point; i.e. if there exists  $n \in \mathbb{N}$  such that  $S_n \models B^R P$  iff  $S_m \models B^R P$ , for all  $m > n$  and all doxastic propositions  $P, R$ . Equivalently, iff there exists  $n \in \mathbb{N}$  such that  $(bestR)_{S_n} = (bestR)_{S_m}$  for all  $m > n$  and all  $R$ .

**Stabilization Results for Conditional Beliefs and Knowledge** In Baltag and Smets (2009b) we proved the following results:

**Proposition 1** (cf (Baltag and Smets 2009b)) If an upgrade stream  $\vec{\dagger P}$  stabilizes a pointed model  $S$ , then  $\vec{\dagger P}$  stabilizes all conditional beliefs on  $S$  and vice versa. Also, if this is the case then  $\vec{\dagger P}$  also stabilizes all knowledge and all (simple) beliefs. Every upgrade stream stabilizes all *knowledge*. Every *update* stream stabilizes every model on which it is executable.<sup>3</sup>

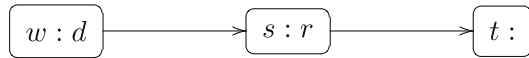
Note however that the analogue of the last result is *not true* for arbitrary upgrade streams, *not even for truthful upgrade streams*. It is not even true for repeated truthful upgrades:

**Counterexample 2** Suppose we are in the situation of Example 1 and that Charles learns from *strongly trusted (but fallible) source* the following true statement  $P$ : “If

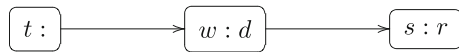
<sup>3</sup> An update stream  $(!P_n)_{n \in \mathbb{N}}$  is *executable* on  $S$  if each  $!P_n$  is executable at its turn, i.e. if  $S_n \models P_n$  for all  $n$ .



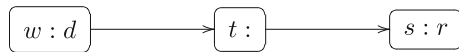
you would truthfully learn (from some infallible source) that Mary won't vote Republican, then your resulting belief about whether or not she votes Democrat would be wrong". The statement  $\mathbf{P}$  can be formally written as  $[\!-\mathbf{r}](\mathbf{B}\mathbf{d} \Leftrightarrow \neg\mathbf{d})$  (using dynamic update modalities to capture the possibility of truthfully learning from an infallible source), and in its turn this last expression can be seen to be equivalent to  $\mathbf{r} \vee (\mathbf{B}^{-\mathbf{r}}\mathbf{d} \Leftrightarrow \neg\mathbf{d})$ . The act of learning this sentence  $\mathbf{P}$  in this context is a *truthful radical upgrade*  $\uparrow \mathbf{P}$  (since the proposition  $\mathbf{P}$  is true in the actual world  $s$ , as well as in  $t$ , though the source is not known to be infallible, but is only strongly trusted). Applying this transformation we obtain the model.



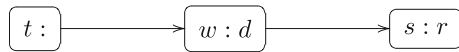
The same sentence is again true in (the real world)  $s$  and in  $w$ , so  $\uparrow \mathbf{P}$  is again truthful, resulting in:



Another truthful upgrade with the same proposition produces



then another truthful upgrade  $\uparrow \mathbf{P}$  gets us back to the previous model:



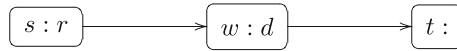
From now on the last two models will keep reappearing in a cycle. The conditional beliefs of Charles never stabilize in this example. But his simple beliefs are the same in these last two models, since  $s$  is the most plausible world in both: so the set of simple (non-conditional) beliefs stays the same from now on. This is a symptom of a more general converge phenomenon:

**Proposition 3** (Belief Convergence Theorem, cf (Baltag and Smets 2009b)) Every truthful radical upgrade stream  $(\uparrow \mathbf{P}_n)_{n \in \mathbb{N}}$  stabilizes all (simple, non-conditional) beliefs (even if it doesn't stabilize the model).

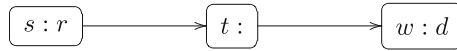
**Corollary 4** (cf (Baltag and Smets 2009b)) Every repeated truthful radical upgrade definable in doxastic-epistemic logic (i.e. in the language of simple belief and knowledge operators, without any conditional beliefs) stabilizes every model (with respect to which it is correct), and thus stabilizes all conditional beliefs.

The analogue of this is not true for conservative upgrades:

**Counterexample 5** (cf (Baltag and Smets 2009b)) Suppose that, in the situation described in Example 1, Charles hears a ("barely believable") rumor  $\mathbf{P}'$  saying that "Either Mary will vote Republican or else your beliefs about whether or not she votes Democrat are wrong". The statement  $\mathbf{P}'$  can be written as:  $\mathbf{r} \vee (\mathbf{B}\mathbf{d} \Leftrightarrow \neg\mathbf{d})$ . This is a *truthful conservative upgrade*  $\uparrow \mathbf{P}'$ , since  $\mathbf{P}'$  is true in  $s$  (as well as in  $t$ ). Its result is:



Now the sentence  $\mathbf{P}'$  is again true in the actual state  $s$  (as well as in  $w$ ), so  $\uparrow \mathbf{P}'$  can again be applied, resulting in:



So these two models (supporting *opposite beliefs!*) will keep reappearing, in an endless cycle. Surprisingly enough, this counterexample shows that we may get into an infinite belief-upgrade cycle, even if the dynamic revision is “directed” towards the real world, i.e. even if all upgrades are truthful!

### 4 Stabilizing on the Truth

As we’ve seen, simple beliefs will eventually stabilize if we consider an infinite series of truthful radical upgrades. Our question in this section is: under what conditions these beliefs stabilize on the *truth*?

**Strongly Informative Upgrades, Questions and Streams** Call an upgrade  $\uparrow \mathbf{P}$  “strongly informative” on a pointed model  $\mathbf{S}$  iff  $\mathbf{P}$  is not already believed at  $\mathbf{S}$ , i.e.  $\mathbf{S} \models \neg \mathbf{B}\mathbf{P}$ . A question  $\mathcal{Q}$  is “strongly informative” for an upgrade type  $\uparrow$  on a pointed model  $\mathbf{S}$  if the action  $\uparrow \mathbf{A}$  of truthfully answering it is strongly informative. An upgrade stream  $(\uparrow \mathbf{P}_n)_{n \in \mathbb{N}}$  is called “strongly informative” if each of the upgrades is strongly informative at the time when it is announced; i.e. if  $\mathbf{S}_n$  satisfies  $\neg \mathbf{B}\mathbf{P}_n$ , for all  $n$ . Finally, a question stream  $\vec{\mathcal{Q}} = (\mathcal{Q}_n)_{n \in \mathbb{N}}$  is strongly informative for an upgrade type  $\uparrow$  if the corresponding truthful upgrade stream  $(\uparrow \mathbf{A}_n)_{n \in \mathbb{N}}$  is strongly informative.

**Belief Correcting Upgrades, Questions and Streams** Call an upgrade  $\uparrow \mathbf{P}$  “belief-correcting” on  $\mathbf{S}$  iff  $\mathbf{P}$  is believed to be false at  $\mathbf{S}$ , i.e.  $\mathbf{S} \models \mathbf{B}\neg \mathbf{P}$ . An upgrade stream is “belief-correcting” if all upgrades are belief-correcting at the time when they are announced; i.e. if  $\mathbf{S}_n \models \mathbf{B}\neg \mathbf{P}_n$ , for all  $n$ . “Belief correcting” implies “strongly informative”, but the converse fails.

**Predicting the True Answers to Future Questions** A consequence of the results in the previous section is that, if the answers to an infinite sequence of questions are successively given by using *updates or radical upgrades*, then *there is a finite stage after which the agent can “predict” the correct answers to all the future questions*. More precisely:

**Corollary 6** For every infinite sequence of questions, the agent’s beliefs converge, after some finitely many updates or radical upgrades with the correct answers, to a belief that correctly answers all future questions in the given sequence. In other words: *for every infinite stream of updates, or truthful radical upgrades, there is a finite stage after which no subsequent upgrades are strongly informative*. Hence, *there are no (infinite) strongly informative update, or radical upgrade, streams; and there are no (infinite) belief-correcting update, or radical upgrade, streams*.

*Proof* Let  $(\dagger\mathbf{P}_n)_{n \in \mathbb{N}}$  be a stream of truthful updates, or truthful radical upgrades. Let  $\mathbf{S}$  be a given finite pointed model having  $s$  as the real world, and let  $(\mathbf{S}_n)_{n \in \mathbb{N}}$  be the infinite sequence of pointed models obtained by performing the successive upgrades according to the above upgrade stream. By Propositions 1 and 3, the set of most plausible states stabilizes at some finite stage: there exists  $m$  such that  $best_{\mathbf{S}_m} = best_{\mathbf{S}_n}$  for all  $n > m$ . We show that none of the upgrades  $\dagger\mathbf{P}_n$ , with  $n > N$ , is strongly informative. Suppose, towards a contradiction, that one of these upgrades would be strongly informative, i.e.  $\mathbf{S}_n \models B-\mathbf{P}_n$ . This means that  $best_{\mathbf{S}_n} \cap (\mathbf{P}_n)_{\mathbf{S}_n} \neq \emptyset$ . But on the other hand all our upgrades have the AGM property that, *after an upgrade* with a truthful proposition, the proposition is *believed to have been true before the upgrade*: i.e. in the new model  $\mathbf{S}_{n+1}$  after the upgrade, we have that  $best_{\mathbf{S}_{n+1}} \subseteq (\mathbf{P}_n)_{\mathbf{S}_n}$ . This shows that we must have  $best_{\mathbf{S}_{n+1}} \neq best_{\mathbf{S}_n}$ , which contradicts the stabilization at  $m$  of the set of most plausible states (i.e. the fact that  $best_{\mathbf{S}_m} = best_{\mathbf{S}_n} = best_{\mathbf{S}_{n+1}}$ , for all  $n \geq m$ ).  $\square$

This result is connected to issues in Formal Learning Theory, and especially to the area of *function learning*: the issue of identifying a function (infinite data stream), in the sense of being able to *predict* the entire future data stream after seeing only an initial segment. See (Hendricks et al. 1997; Kelly 1998a, b) for an analysis of the function-learning power of iterated belief revision.

**Convergence to Full Truth via Complete Question Streams** As a consequence, we can show that, if a question stream is “complete”, in the sense of being enough to completely characterize the real world, then after finitely many (updates or) radical upgrades with truthful answers, the agent’s beliefs stabilize on full truth.

A (truthful) upgrade stream  $(\dagger\mathbf{P}_n)_{n \in \mathbb{N}}$  is *complete* for a language  $\mathcal{L}$  on a pointed model  $\mathbf{S}$ , if it is enough to (truthfully) settle every issue definable in  $\mathcal{L}$ . In other words: for every world  $t \in S$ , if the stream  $(\dagger\mathbf{P}_n)_{n \in \mathbb{N}}$  is truthful in the world  $t$  (of the same plausibility frame as the one of  $\mathbf{S}$ ), then  $t$  satisfies the same sentences in  $\mathcal{L}$  as the real world  $s$  of the (pointed) model  $\mathbf{S}$ . Similarly, a question stream  $(\mathcal{Q}_n)_{n \in \mathbb{N}}$  is *complete with respect to an upgrade type  $\dagger$*  if the corresponding truthful upgrading stream  $(\dagger\mathcal{Q}_n)_{n \in \mathbb{N}}$  (that successively upgrades with the correct answers) is complete.

It is obvious that, only for a complete upgrade/question stream, one can expect the agent’s beliefs to stabilize on the “full truth” after finitely many answers: how else can they do it, when not even the information contained in the whole infinite stream of answers is enough to determine the full truth?

And indeed, in this case we do have a positive result:

**Corollary 7** *Every complete (and truthful) upgrade stream  $(\dagger\mathbf{P}_n)_{n \in \mathbb{N}}$  of types  $\dagger \in \{!, \uparrow\}$ , starting on a given finite model  $\mathbf{S}$ , stabilizes the beliefs on the “full truth” in finitely many steps: eventually, all  $\mathcal{L}$ -expressible beliefs are true and all true  $\mathcal{L}$ -sentences are believed. In other words, from some stage onwards, a sentence in  $\mathcal{L}$  is believed if and only if it is true.*

*Proof* Let  $(\dagger\mathbf{P}_n)_{n \in \mathbb{N}}$ ,  $\mathbf{S}$ ,  $(\mathbf{S}_n)_{n \in \mathbb{N}}$  and  $m$  be as in the proof of Corollary 6. By Corollary 6, the set  $best_{\mathbf{S}_n}$  stabilizes at  $m$  onto a set of states  $best_{\mathbf{S}_m}$  for which no subsequent upgrade in the stream is strongly informative. This means that, for all

$n \geq m$ , we have  $best_{S_m} = best_{S_n} \subseteq (\mathbf{P}_n)_{S_n}$ . It is also so (by the same AGM property of upgrades used above in the proof of Corollary 6) that we have  $best_{S_m} \subseteq (\mathbf{P}_k)_{S_k}$  for all  $k < m$  as well. Putting these together, it follows that  $best_{S_m} \subseteq (\mathbf{P}_k)_{S_k}$  for all  $k \in N$ . But then the completeness of the upgrade stream implies that every state  $t \in best_{S_m} = best_{S_n}$  (for any  $n > m$ ) satisfies the same  $\mathcal{L}$ -sentences as the real state  $s$ . This means that the agent's beliefs stabilize on "full truth": from stage  $m$  onwards, every  $\mathcal{L}$ -sentence is believed iff it is true (in the real world  $s$ ).  $\square$

While in the case of updates, the above revision process comes to an end, at which stage the agent actually *knows* the full truth, this is *not necessarily the case for radical upgrades*. Indeed, for iterated radical upgrades, it is possible that the agent will never know the full truth at any finite stage of revision, even though she will come to *believe* the full truth. Moreover, *even when after the stage when her beliefs match the full truth, she won't necessarily know that!* Indeed, for all she knows, her beliefs might still be revised in the future: *she never gets to know that her beliefs have stabilized*, and consequently *she never knows that she has reached the truth* (even if we assume she knows that the stream is complete and truthful!).

This result is connected to "language learning" (or "set learning") in Formal Learning Theory. In our terms, this can be stated as the problem of learning *the set of 'all true 'data'' (true  $\mathcal{L}$ -sentences), based only on the finite sequence of data that were already observed*. In Learning Theory terms, our result says that *every finite set of worlds is "identifiable in the limit from positive and negative data"*, by a specific "*belief-revision-based*" learning method: namely, start with any preorder on the finite set as the "plausibility relation", and do belief-revision via *either updates or radical upgrades*. The concept of identifiability in the limit is related to *reaching a stable belief in full truth*; while *reaching knowledge of full truth* (as in the case of repeated updates) corresponds to "*finite identifiability*" (in the terminology used in Learning Theory). Dégremont and Gierasimczuk (2009) first introduced this distinction in the context of belief revision, and investigated the connection between repeated updates and finite identifiability.

The result is *not true for infinite models*<sup>4</sup>, *nor for other belief-revision methods* (such as conservative upgrade). Indeed, Counterexample 5 in the previous section shows that the analogue of Corollary 7 fails for conservative upgrades. It is easy to see that *the infinite sequence of truthful upgrades in that example is complete* (for any language) on the given pointed model: it completely determines the real world. (In fact, *any two successive upgrades on that sequence are already complete*: only the real world  $s$  makes both of them truthful!) Nevertheless, the beliefs never stabilize.

**Special Case: Maximally Strongly informative streams** An upgrade stream or sequence is a "*maximal*" *strongly-informative* (or *maximal belief-correcting*), truthful stream/sequence if: (1) it is strongly-informative (or belief-correcting) and truthful, and (2) it cannot properly be extended to any stream/sequence having property (1).

<sup>4</sup> In on-going work with N. Gierasimczuk, we are investigating in detail the infinite case, proving various limited analogues of this result.

So a strongly informative truthful stream is “maximal” if either (a) it is infinite, or else (b) it is finite (say, of length  $n$ ) but there exists no upgrade  $\uparrow\mathbf{P}_{n+1}$  which would be truthful and strongly informative on *the last model*  $\mathbf{S}_n$ . The next result says that the first alternative (a) is *impossible*, and moreover *maximal strong informativity ensures convergence to full truth*:

**Proposition 8** For  $\uparrow \in \{!, \uparrow\}$ , every maximal strongly-informative truthful upgrade stream/sequence  $(\uparrow\mathbf{P}_k)$  starting on a given finite model  $\mathbf{S}$  is *finite and complete*. Hence, it reaches a fixed point (final model)  $\mathbf{S}_n$ , in which all beliefs are true and all true sentences are believed. In other words, in the final model, a sentence is believed if and only if it is true.

*Proof* By Corollary 6, for  $\uparrow \in \{!, \uparrow\}$ , there are no infinite strongly informative  $\uparrow$ -streams. Hence,  $(\uparrow\mathbf{P}_k)$  must be *finite*. Let  $\mathbf{S}_n$  be its last model. Suppose there exists a truth that is not believed at the last stage, i.e. there exists a  $\mathcal{L}$ -expressible proposition  $\mathbf{P}$  such that  $\mathbf{S}_n$  satisfies  $\mathbf{P}$  and  $\neg\mathbf{BP}$  at the same time. Then we can add to the stream a new upgrade, namely  $\uparrow\mathbf{P}$ . This is still truthful and strongly informative, so the stream was not maximal and we have reached a contradiction. Hence, in  $\mathbf{S}_n$ , every true  $\mathcal{L}$ -sentence is believed. But then the converse must also be the case: every  $\mathcal{L}$ -expressible belief must be true! Indeed, otherwise, let  $\mathbf{Q}$  be an  $\mathcal{L}$ -expressible proposition such that  $\mathbf{S}_n$  satisfies  $\mathbf{BQ}$  and  $\neg\mathbf{Q}$  in the same time. Then  $\mathbf{P} := \neg\mathbf{Q}$  is also an  $\mathcal{L}$ -expressible proposition, such that both  $\mathbf{P}$  and  $\mathbf{B}\neg\mathbf{P}$  are true. By the first part, since  $\mathbf{P}$  is true, it must be believed, so  $\mathbf{BP}$  is also true. But then in  $\mathbf{S}_n$ , we have both  $\mathbf{B}\neg\mathbf{P}$  and  $\mathbf{BP}$ , which contradicts the principle of Consistency of Beliefs (valid in all plausibility models).  $\square$

Note that *the last result is not necessarily equivalent to saying that the set of most plausible worlds coincides in the end with only the real world!* The reason is that the language may not be expressive enough to distinguish the real world from some of other ones.

**Convergence to Partial Truth** If we only want to ensure that beliefs will stabilize on *true (but not necessarily complete!) beliefs*, then we can weaken our conditions:

**Proposition 9** For  $\uparrow \in \{!, \uparrow\}$ , every maximal belief-correcting truthful upgrade stream/sequence  $(\uparrow\mathbf{P}_k)$  (starting on a finite model  $\mathbf{S}$ ) is *finite and converges to true beliefs*; i.e. in its last model  $\mathbf{S}_n$ , all the beliefs are true.

*Proof* By Corollary 6, there are no infinite belief-revising  $\uparrow$ -streams, for  $\uparrow \in \{!, \uparrow\}$ . So  $(\uparrow\mathbf{P}_k)$  must be finite. Let  $\mathbf{S}_n$  be its last model. Suppose there exists a false belief at this last stage, i.e. there exists a proposition  $\mathbf{P}$  such that  $\mathbf{S}_n$  satisfies  $\mathbf{BP}$  and  $\neg\mathbf{P}$  at the same time. Then we can add to the stream a new upgrade, namely  $\uparrow\neg\mathbf{P}$ . This is still truthful and belief-correcting, so the stream was not maximal: contradiction!  $\square$

## 5 Conclusions

We focused on *the long-term behavior of doxastic structures under iterated revision with higher-level doxastic information*. In contexts when the new information is

purely propositional, this problem has been studied in Learning Theory. But as far as we are aware, there are no results in the literature on convergence to truth via iterated revision with higher-level doxastic information. So our investigation here is a first step in tackling this question.

The truth-convergence results of this paper (presented in Sect. 4) are easy consequences of our results in Baltag and Smets (2009b; reproduced in Sect. 3). Although their proofs are straightforward, given the work in Baltag and Smets (2009b), we think that they have an independent conceptual value. Our investigation shows that *iterated revision with doxastic information is highly non-trivial, even in the single-agent case. Surprisingly, the revision process is not always guaranteed to reach a fixed point, even when indefinitely repeating the same correct upgrade (repeatedly answering the same question!). However, both knowledge and beliefs are eventually stabilized by updates, or radical upgrades, with the true answers to a sequence of questions; moreover the stabilized beliefs are truthful and complete with respect to the answers to all the future questions (in the sequence): from some moment onwards, the agent can correctly predict all future answers (without necessarily knowing that she can do that).*

We further show that *answering complete question streams via iterated truthful updates, or radical upgrades, makes the agent's beliefs converge to full truth in finitely many steps. Again, the exact stage when this happens may remain unknown to the agent, in the case of radical upgrades: her beliefs will be right from some moment on, without her ever knowing that!*

In particular, convergence of beliefs to full truth happens for *maximal strongly-informative question streams*; while in the case of *maximal belief-correcting streams*, the eventually-stabilized beliefs will only be guaranteed to be true (but not complete). The analogues of these results for infinite models, or for other belief-revision methods (e.g. conservative upgrades), are false.

**Connections with Verisimilitude Theory** The results of this paper do not assume any specific theory of Truth Approximation. Our total preorders have no intrinsic relation with any notion of “closeness to the truth”: they capture a purely subjective notion of plausibility, encoding the agent’s contingency plans for belief revision. In particular, our use of possible world semantics does *not* automatically fall prey to Miller’s hot, rainy and windy argument (Miller 1974): that argument assumes a specific measure of verisimilitude, given by the (cardinal of the) set of atomic sentences whose believed truth values differ from their “real” truth values (in the actual world).

However, our results do have some obvious consequences for the field of Truth Approximation<sup>5</sup>. First, the Learning-Theoretic perspective adopted in the last section can be used to define a “dynamic” *measure of verisimilitude* of a given (finite, pointed) plausibility model: simply define it by comparing the *lengths of the longest maximally strongly informative question stream*, and declaring the model with the shortest such length to be closer to the truth. Then the above results show

<sup>5</sup> For a good survey of the field, see e.g. (Niiniluoto 1998).

that *performing a strongly informative (update or) radical upgrade on a (finite, pointed) model always leads to a model that is closer to the truth.*

Our negative results have a more general relevance to *any* reasonable theory of verisimilitude. Namely, they point to an *inherent tension* between the “rational” demand of “*conservatism*” of belief revision (underlying the AGM paradigm, and asking for always keeping as much as possible of the “old” belief structure when accommodating the new information) and the equally “rational” demand of *increasing the verisimilitude*. *The more “conservative” the belief revision method, the “harder” is to get closer to the truth*, even when receiving only true, and non-redundant, new information! *Updates* (which put absolute trust in the new information, and ruthlessly dismiss all the old possibilities that are incompatible with it) *always increase verisimilitude*. *Radical upgrades* with new true information (which put strong, but not absolute, trust in it, and hence induce a sweeping revision of the old belief structure, downgrading all the old possibilities that are incompatible with the new piece of data) *only increase verisimilitude when they are strongly informative*. Finally, the most conservative revision method, the so-called *conservative upgrade*, *fails in general to increase verisimilitude*, occasionally leading to infinite belief-revision loops.

**Acknowledgments** Sonja Smets thanks Theo Kuipers for organizing an inspiring workshop on Belief Revision Aiming at Truth Approximation during the EPSA conference in Amsterdam. Smets’ contribution to this paper was made possible by the VIDI research programme number 639.072.904, financed by the Netherlands Organization for Scientific Research. The authors thank Johan van Benthem, Nina Gierasimczuk, Vincent Hendricks and David Makinson for insightful discussions of the issues investigated in this paper. Finally, we thank the anonymous referees and the editors of this volume for their very useful comments and feedback.

**Open Access** This article is distributed under the terms of the Creative Commons Attribution Non-commercial License which permits any noncommercial use, distribution, and reproduction in any medium, provided the original author(s) and source are credited.

## References

- Alchourrón, C. E., Gärdenfors, P., & Makinson, D. (1985). On the logic of theory change: partial meet contraction and revision functions. *The Journal of Symbolic Logic*, 50, 510–530.
- Aucher, G. (2003). *A combined system for update logic and belief revision*. Master’s thesis, ILLC. Amsterdam. The Netherlands: University of Amsterdam.
- Baltag, A., Gierasimczuk, N., & Smets, S. (2010). *Truth-tracking by belief-revising* (Unpublished Manuscript).
- Baltag, A., & Smets, S. (2006a). Conditional doxastic models: A qualitative approach to dynamic belief revision. *Electronic Notes in Theoretical Computer Science*, 165, 5–21.
- Baltag, A., & Smets, S. (2006b). Dynamic belief revision over multi-agent plausibility models. In *Proceedings of LOFT’06*, pp. 11–24.
- Baltag, A., & Smets, S. (2008a). The logic of conditional doxastic actions. *Texts in Logic and Games*, 4, 9–32.
- Baltag, A., & Smets, S. (2008b). A qualitative theory of dynamic interactive belief revision. *Texts in Logic and Games*, 3, 9–58.
- Baltag, A., & Smets, S. (2009a). Group belief dynamics under iterated revision: Fixed points and cycles of joint upgrades. In *Proceedings of TARK* (Vol. 12, pp. 41–50).

- Baltag, A., & Smets, S. (2009b). Learning by questions and answers: from belief-revision cycles to doxastic fixed points. In R. de Queiroz, H. Ono, & M. Kanazawa (Eds.), *Proceedings of WOLLIC'09, LNAI 5514* (Vol. 5514, pp. 124–139).
- Blackburn, P., de Rijke, M., & Venema, Y. (2001). *Modal Logic*. Cambridge: Cambridge University Press.
- Board, O. (2002). Dynamic interactive epistemology. *Games and Economic Behaviour*, 49, 49–80.
- Booth, R., & Meyer, T. (2006). Admissible and restrained revision. *Journal of Artificial Intelligence Research*, 26, 127–151.
- Boutilier, C. (1996). Iterated revision and minimal change of conditional beliefs. *Journal of Philosophical Logic*, 25(3), 262–305.
- Darwiche, A., & Pearl, A. (1997). On the logic of iterated belief revision. *Artificial Intelligence*, 89, 1–29.
- Dégremont, C., & Gierasimczuk, N. (2009). Can doxastic agents learn? On the temporal structure of learning. *Lectures Notes in Artificial Intelligence*, 5834, 90–104.
- Grove, A. (1988). Two modellings for theory change. *Journal of Philosophical Logic*, 17, 157–170.
- Halpern, J. Y. (2003). *Reasoning about uncertainty*. Cambridge MA: MIT Press.
- Hendricks, V., Kelly, K., & Schulte, O. (1997). Reliable belief revision. *Logic and scientific methods*.
- Katsuno, H., & Mendelzon, A. (1992). On the difference between updating a knowledge base and revising it. In *Cambridge tracts in theoretical computer science* (pp. 183–203).
- Kelly, K. (1998a). Iterated belief revision, reliability, and inductive amnesia. *Erkenntnis*, 50, 11–58.
- Kelly, K. (1998b). The learning power of iterated belief revisions. In I. Gilboa (Ed.), *Proceedings of the seventh TARK conference* (pp. 111–125).
- Miller, D. (1974). Popper's qualitative theory of verisimilitude. *The British Journal for the Philosophy of Science*, 25(2), 166–177.
- Niiniluoto, I. (1998). Verisimilitude: The third period. *The British Journal for the Philosophy of Science*, 49, 1–29.
- Rott, H. (1989). Conditionals and theory change: Revisions, expansions, and additions. *Synthese*, 81, 91–113.
- Spohn, W. (1988). Ordinal conditional functions: A dynamic theory of epistemic states. *Causation in Decision, Belief Change, and Statistics, II*, 105–134.
- van Benthem, J. F. A. K. (2007). Dynamic logic of belief revision. *Journal of Applied Non-Classical Logics*, 17(2), 129–155.
- van Benthem, J. F. A. K. (2009). Logical dynamics of information and interaction. In *Manuscript* (To appear).
- van Ditmarsch, H. P. (2005). Prolegomena to dynamic logic for belief revision. *Synthese*, 147, 229–275.
- Veltman, F. (1996). Defaults in update semantics. *Journal of Philosophical Logic*, 25, 221–261.